

Pre-Training and Prompting for Few-Shot Node Classification on Text-Attributed Graphs

Huanjing Zhao[†]

Tsinghua University
Beijing, China

hwangyeong@mail.tsinghua.edu.cn

Beining Yang^{†‡}

University of Edinburgh
Edinburgh, UK

B.Yang-32@sms.ed.ac.uk

Yukuo Cen

Zhipu AI
Beijing, China

yukuo.cen@zhipuai.cn

Junyu Ren

Tsinghua University
Beijing, China

renjy19@mails.tsinghua.edu.cn

Chenhui Zhang

Zhipu AI
Beijing, China
chenhui.zhang@zhipuai.cn

Yuxiao Dong

Tsinghua University
Beijing, China
yuxiaod@tsinghua.edu.cn

Evgeny Kharlamov

Bosch Center for Artificial
Intelligence

Renningen, Germany
evgeny.kharlamov@de.bosch.com

Shu Zhao

Anhui University
Anhui, China
zhaoshuzs@ahu.edu.cn

Jie Tang^{*}

Tsinghua University
Beijing, China
jietang@tsinghua.edu.cn

ABSTRACT

The text-attributed graph (TAG) is one kind of important real-world graph-structured data with each node associated with raw texts. For TAGs, traditional few-shot node classification methods directly conduct training on the pre-processed node features and do not consider the raw texts. The performance is highly dependent on the choice of the feature pre-processing method. In this paper, we propose P2TAG¹, a framework designed for few-shot node classification on TAGs with graph pre-training and prompting. P2TAG first pre-trains the language model (LM) and graph neural network (GNN) on TAGs with self-supervised loss. To fully utilize the ability of language models, we adapt the masked language modeling objective for our framework. The pre-trained model is then used for the few-shot node classification with a mixed prompt method, which simultaneously considers both text and graph information. We conduct experiments on six real-world TAGs, including paper citation networks and product co-purchasing networks. Experimental results demonstrate that our proposed framework outperforms existing graph few-shot learning methods on these datasets with +18.98% ~ +35.98% improvements.

CCS CONCEPTS

• **Mathematics of computing** → **Graph algorithms**; • **Computing methodologies** → **Learning latent representations**.

[†] Both authors contributed equally to this research.

[‡] This work was done when the author was interned at Zhipu AI.

^{*} Corresponding author: JT.

¹ Our code is available at <https://github.com/THUDM/P2TAG>.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

KDD '24, August 25–29, 2024, Barcelona, Spain

© 2024 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-0490-1/24/08

<https://doi.org/10.1145/3637528.3671952>

KEYWORDS

graph self-supervised learning, text-attributed graphs, few-shot node classification, graph neural networks

ACM Reference Format:

Huanjing Zhao, Beining Yang, Yukuo Cen, Junyu Ren, Chenhui Zhang, Yuxiao Dong, Evgeny Kharlamov, Shu Zhao, Jie Tang. 2024. Pre-Training and Prompting for Few-Shot Node Classification on Text-Attributed Graphs. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '24)*, August 25–29, 2024, Barcelona, Spain. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3637528.3671952>

1 INTRODUCTION

The few-shot node classification task involves identifying the classes of nodes in a given graph structure using only a limited number of labeled examples. This task has practical applications in areas such as social network analysis, recommendation systems, and more. Inspired by the successful experiences in the field of Computer Vision (CV), several works, such as G-Meta and TENT [7, 16, 34, 43], apply meta-learning to graphs to address the few-shot node classification problem. These methods learn transferable knowledge from meta-tasks, enabling rapid adaptation to new, unseen labels. Unlike images, graphs represent a form of structured data. Through graph pre-training, models can also accumulate a substantial amount of domain-specific knowledge.

Recently, a series of graph pre-training methods [30, 40, 50, 54, 57] emerged with self-supervised learning (SSL) to yield generalized node representations in the absence of labels. These methods mainly include contrastive and generative ones. Contrastive SSL methods utilize data augmentation to generate multiple views of data for contrastive loss. Generative SSL methods such as GraphMAE [15] aim to reconstruct the (masked) attributes of the graph. These self-supervised learning methods have greatly contributed to many aspects, such as graph augmentation and graph-based pretext tasks. They only consider part of the self-supervised learning on TAGs, making the training of GNNs independent of LMs text encoding.

Table 1: Related works of few-shot learning on TAGs.

Method	Few-shot	Text	Prompt Type
G-Meta [16]	✓	-	-
TENT [34]	✓	-	-
Prog [29]	✓	-	Graph Prompt
TAPE [13]	-	✓	Hard Text Prompt
ENG [47]	✓	✓	Hard Text Prompt
GPrompt [17]	✓	✓	Hard Text Prompt
G2P2 [36]	✓	✓	Soft Text Prompt
P2TAG (Ours)	✓	✓	Mixed Prompt

The nodes in the graph usually contain rich textual information [18] such as titles and abstracts in citation networks. Some recent works [3, 41, 51] attempt to integrate the training of LMs and GNNs. These methods directly deal with the original texts and the topological structure in the graph, and achieve better performance on the downstream tasks. GIANT [3] aims to obtain stronger textual representations by the language model with neighborhood prediction in TAGs. Although GIANT utilizes graph information to fine-tune the language models, the text encoding and graph propagation steps are still independent. GraphFormer [41] proposes a nested architecture of LMs and GNNs and jointly trains them with a link-prediction-based objective. GLEM [51] alternatively optimizes LMs and GNNs using label information from the downstream tasks.

In few-shot node classification, prompting helps guide the pre-trained model to generalize from a few examples by providing specific input cues. Prog [29] constructs prompts at the graph level, which is not directly applicable to TAGs. With the leapfrog development of large language models (LLMs), some works attempt to construct prompts by processing raw texts in the TAG, the name we refer to as hard text prompt, exemplified by TAPE [13] and ENG [47]. GPrompt [17] leverages the GNN model as a downstream adapter, using hard text prompts to improve adaptation efficiency. These works rely on the quality of constructed prompts. It is highly challenging since there are both text (LM side) and structure information (GNN side) on TAGs. G2P2 [36] tries to solve this problem solely from the LM side, named soft text prompt. They reweight the pre-trained GNN embedding by the neighboring nodes' texts and label texts. However, the potential capability of GNN is neglected in this way and may lead to sub-optimal results.

Present work: P2TAG. In this paper, we propose the **Pre-training and Prompting** for few-shot node classification on TAGs. In the pre-training phase, different from previous works, our framework integrates the LM and GNN effectively with joint training in a self-supervised way. Some solutions are proposed to address the challenges of integration of LMs and GNNs. We incorporate the encoder of GNNs as a supplement to LMs and keep the original self-supervised loss of the LM model, i.e., the masked language modeling objective, which is a very hard pretext task and can alleviate the over-fitting issue. To train the GNNs and the LMs jointly, our framework first samples mini-batch subgraphs through a random-walk-based sampler. Each mini-batch will be fed to the LMs for text encoding, and the output of the classification token (i.e., [CLS] token) will be considered as the representation of the node (or the text sequence). The output embeddings of [CLS] tokens will be fed

to the GNNs to aggregate the information from neighbors through graph propagation. For the masked language modeling, we incorporate the aggregation of neighbor information into the outputs of [MASK] tokens. Finally, the cross-entropy loss of the outputs and the original tokens is computed.

In the prompting phase, to bridge the gap between pre-train tasks and the downstream ones, we process with a mixed approach of LM and GNN, which differs from previous works as shown in Table 1. Specifically, we jointly concatenate the *label text initialized graph prompt* and a simple *LM output initialized text prompt* to mimic the training paradigm. As illustrated in Figure 2, by leveraging the label text embedding, we can easily align the node to the downstream label space without discrete and intuitive hand-craft prompts. The empirical results presented in Section 4.3 attest to the superior efficacy of our prompt design, which achieves a substantial improvement of 11.1% ~ 25.2% over the existing soft text prompted G2P2 method [36].

The key contributions of our framework are summarized in the following.

- We propose a unified framework to jointly train language models and graph neural networks via the masked language modeling objective.
- We propose a new graph-text mixed prompt learning method with text and structure information to mitigate the gap between pre-training and few-shot downstream tasks.
- Our proposed P2TAG achieves an improvement with +18.98% ~ +35.98% over meta-learning methods across six TAG datasets.

2 PRELIMINARIES

In this section, we introduce the background of our paper including text-attributed graph and few-shot node classification.

Notations. Denote a graph $G = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} is a set of n nodes and \mathcal{E} is a set of edges between nodes. $A \in \mathbb{R}^{n \times n}$ is the adjacency matrix where its entry $A(i, j) \geq 0$, if nonzero, denotes there is an edge between node i and j with edge weight $A(i, j)$. In practice, the network could be either directed or undirected. If G is directed, we have $A(i, j) \neq A(j, i)$; if G is undirected, we have $A(i, j) \equiv A(j, i)$.

2.1 Text-Attributed Graph

DEFINITION 1 (TEXT-ATTRIBUTED GRAPH). A *text-attributed graph (TAG)* is a graph $G = (\mathcal{V}, \mathcal{E}, S)$, where each node $v_i \in \mathcal{V}$ is associated with a text sequence $s_i \in S$ and \mathcal{E} represents the set of edges between nodes.

Given the above definition, we can formally define our problem for self-supervised learning on TAG.

PROBLEM 1 (SELF-SUPERVISED LEARNING ON TAGs). Given a TAG $G = (\mathcal{V}, \mathcal{E}, S)$, the problem of *self-supervised learning (SSL) on TAG* is to learn a unified low-dimensional space representation of each text-attributed node $v_i \in \mathcal{V}$ with text sequence s_i . The goal is to find a function $f : \mathcal{V} \rightarrow \mathbb{R}^d$, where $d \ll |\mathcal{V}|$.

Previous methods implement self-supervised learning on TAG through separated text encoders and GNNs. Built upon these works, we explore self-supervised learning by integrating LMs and GNNs.

The joint training for the problem of SSL on TAGs can be formulated as follows:

$$\theta_1^*, \theta_2^*, \theta_3^* = \arg \min_{\theta_1, \theta_2, \theta_3} \mathcal{L}_{SSL} \left(f_{\theta_1}^{GNN}, f_{\theta_2}^{LM}, f_{\theta_3}^{Head}, G \right),$$

where $f_{\theta_1}^{GNN}$, $f_{\theta_2}^{LM}$, and $f_{\theta_3}^{Head}$ are the GNN encoder, LM encoder, and the prediction head for the SSL. Then we can obtain the final embeddings via the following:

$$Z = f_{\theta_1}^{GNN}(G, X), \text{ where } X = f_{\theta_2}^{LM}(S).$$

2.2 Few-shot Node Classification

Few-shot node classification are tasks that involve training a machine learning model with a limited number of labeled nodes, and subsequently predicting the classes of other nodes based on these initial labels. The N -way K -shot few-shot node classification tasks \mathcal{T} on a graph is described as follows with the node label set C . Each task $\mathcal{T}_t \in \mathcal{T}$ selects N labels from C , forming a subset $C_t = \{C_1, C_2, \dots, C_N\} \in C$. The \mathcal{T}_t consists of two parts as follows:

$$\mathcal{T}_t = \{S_t, Q_t\},$$

where S_t means the support set and Q_t means the query set. The machine learning model is trained on S_t and subsequently evaluated on Q_t . Both the support set and the query set consist of nodes and their corresponding labels, which shown as follows:

$$\begin{aligned} S_t &= \{(v_1, c_1), (v_2, c_2), \dots, (v_{N \times K}, c_{N \times K})\}, \\ Q_t &= \{(v'_1, c'_1), (v'_2, c'_2), \dots, (v'_{N \times Q}, c'_{N \times Q})\}, \end{aligned}$$

where v and v' represent nodes, their corresponding labels are denoted as c and c' . In the support and query sets, each label $C_i \in C_t$ associated with K nodes and Q nodes, respectively. For the support set, there exists $c_{(i-1) \times K + j} \in C_i$, where $i \in \{1, 2, \dots, N\}$ and $j \in \{1, 2, \dots, K\}$. Similarly, there exists $c'_{(i-1) \times Q + j} \in C_i$, where $i \in \{1, 2, \dots, N\}$ and $j \in \{1, 2, \dots, Q\}$ for the query set.

3 METHODOLOGY

We tackle few-shot node classification by adhering to the pre-training and prompting paradigm, as depicted in Figure 1. To better leverage the unique characteristics of the TAGs, we propose a new self-supervised framework for pre-training.

3.1 Pre-training Framework

In this part, we introduce our proposed pre-training framework in detail. For self-supervised learning on TAGs, previous works usually take two separate steps. The first step is to encode the raw texts to node features using bag-of-words, word2vec [24], or pre-trained language models [4, 6, 8, 12, 19, 22, 42]. Most graph self-supervised methods [15, 30, 50, 54, 57] use the node features pre-processed in the first step to construct a self-supervised objective. The two-step process can bring non-negligible information loss. To address this issue, we propose an end-to-end self-supervised framework to train directly on the TAGs. Inspired by the recent prosperity of pre-trained language models (LM), we choose an effective LM to encode the raw texts. To model the relations between nodes (texts), we can utilize powerful graph neural networks (GNNs). As illustrated in Figure 1, our framework consists of two core modules, including a pre-trained language model and a GNN encoder. The GNN encoder

is to help make better node representations. In our framework, we mainly use the DeBERTa-base [12] with 100M parameters as the LM of our framework. DeBERTa utilizes two techniques based on BERT [6] and RoBERTa [22] and significantly improves the performance on the natural language understanding and generation tasks. The choice of the LMs is flexible, and we also explore other LMs in our experiments.

The pre-training of LMs and GNNs is much more challenging because we need to 1) choose an appropriate self-supervised training objective to avoid over-fitting and 2) sample small mini-batches to address the high computational and space costs of large LMs.

Self-supervised training objective. The training objective plays an important role in self-supervised learning. Different from conventional graph self-supervised learning, the design of self-supervised objective for TAGs are more challenging. Our model architecture contains two different modules with different scales of model parameters (large LM v.s. small GNN), making the training much more difficult. Simple self-supervised objectives will probably make the model overfitted. To make a harder self-supervised objective, we adopt the masked language modeling (MLM) introduced in [6] as our objective. Given a node v_i , we first feed the text sequence $s_i = [T_1, T_2, \dots, T_{n_i}]$ into a language model to get the hidden representations associated with v_i where T_t is t -th token of s_i and n_i is the length of s_i . In the training stage, we randomly mask a pre-defined portion of tokens by replacing them with a special token named [MASK]. We use $s'_i = [T'_1, T'_2, \dots, T'_{n_i}]$ to denote the text sequence after masking and each token T'_t is a random variable with the following distribution:

$$\Pr(T'_t = [\text{MASK}]) = p \text{ and } \Pr(T'_t = T_t) = 1 - p. \quad (1)$$

Here, $p \in (0, 1)$ is a hyper-parameter representing the mask rate. When working with BERT-like language models, we usually add a starting token (e.g., [CLS]) and an ending token (e.g., [SEP]) to the sequence. Therefore, each node is represented as a sequence of hidden vectors:

$$[\mathbf{o}_0, \mathbf{o}_1, \dots, \mathbf{o}_{n_i}, \mathbf{o}_{n_i+1}] = f^{LM} \left([[\text{CLS}], T'_1, T'_2, \dots, T'_{n_i}, [\text{SEP}]] \right), \quad (2)$$

where $\mathbf{o}_t \in \mathbb{R}^d$ is the hidden representations of t -th token of s'_i . Notice that the first hidden vector \mathbf{o}_0 corresponds to the special token [CLS] and can be treated as the “summary” representation of the whole text sequence s_i . To capture the correlations among nodes, we then use a graph neural network to propagate the hidden representations of nodes where the input of node v_i , denoted as x_i , is the first hidden vector \mathbf{o}_0 of s_i . The node representations after passing a GNN encoder are denoted as $\mathbf{H} = f^{GNN}(G, X)$. The propagated node presentations are exploited to construct the self-supervised training objective.

For each node v_i , we concatenate the hidden representation \mathbf{h}_i with the output vector \mathbf{o}_t of each token and then feed the concatenated vector to an MLP as:

$$\mathbf{z}_t = \text{MLP}([\mathbf{h}_i^\top, \mathbf{o}_t^\top]), \text{ where } t = 0, \dots, (n_i + 1). \quad (3)$$

The objective of the masked language model (MLM) is to predict the masked tokens in the original text sequence s_i , which can be

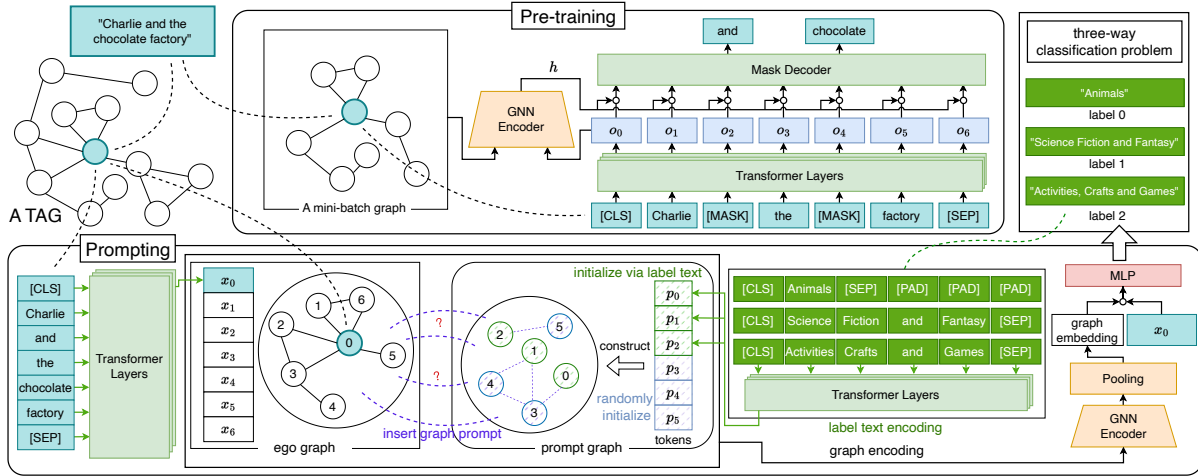


Figure 1: Our proposed framework P2TAG. A toy example illustrating the 3-way classification of children’s books. (1) In the pre-training phase, we jointly train the LM and GNN using a masked language modeling objective. For the GNN module, a subgraph-based sampler employing random walks generates mini-batch subgraphs for training. (2) In the prompting phase, we construct the ego graph for each target node and generate node features with the pre-trained LM. Graph tokens are utilized to learn the graph structure and create a prompt graph. These tokens are initialized either by encoding label text or through a random initialization process.

expressed as:

$$\mathcal{L}_i = \sum_{t=1}^{n_i} \mathbb{1}(T'_t = [MASK]) \log \Pr(T_t | \mathbf{z}_t). \quad (4)$$

The indicator $\mathbb{1}(T'_t = [\text{MASK}])$ ensures that the loss is only applied on masked tokens. We can control the difficulty of the SSL task via a flexible hyperparameter, i.e., mask rate.

Algorithm 1: The pre-training stage of our framework.

- ```

1 Input TAG $G = (\mathcal{V}, \mathcal{E}, \mathcal{S})$, training steps M , mask rate p ,
 GraphSAINT sampler Sample ;
2 Output Model parameters θ of both the LM and GNN;
3 Pre-processing: compute normalization coefficient α and λ
 according to GraphSAINT sampler [48];
4 for $\text{step} \leftarrow 1$ to M do
5 $G_s = (\mathcal{V}_s, \mathcal{E}_s) \leftarrow \text{Sample}(G)$;
6 $S \leftarrow \{s_i : v_i \in \mathcal{V}_s\}$;
7 $S' \leftarrow \{\text{MaskToken}(s_i, p) : s_i \in S\}$;
8 $O \leftarrow f^{LM}(S')$;
9 $H \leftarrow f^{GNN}(G_s, O[:, 0])$;
10 $O \leftarrow \text{Concat}(O, H)$;
11 $Z \leftarrow \text{MLP}(O)$;
12 Compute the loss via Equation 4;
13 Update θ based on the gradient of the loss;
14 return Model parameters θ ;

```

### 3.2 Mini-batch Training

Due to the heavy language models, we need to adopt a mini-batch training strategy even if we train our model on small graph datasets.

There are several genres of mini-batch training techniques for GNNs. To trade off the efficiency and flexibility, we choose a subgraph-based method, GraphSAINT [48], as our training strategy. At every training step, GraphSAINT constructs a mini-batch by sampling a subgraph from the original graph and generates node representations according to the sampled subgraph. In this work, we adopt the random walk sampler to preserve the connectivity of the whole graph. The sampling process starts by selecting  $r$  root nodes from the whole node set  $\mathcal{V}$  uniformly at random. Starting from each root node, a random walk of length  $l$  is sampled from the original graph structure. We then have a sampled node set  $\mathcal{V}_s$  by adding all the nodes that occurred in the random walks and the subgraph  $G_s$  induced by  $\mathcal{V}_s$  is used in the GNN encoder to generate node representations in the minibatch.

### 3.3 Graph-Text Mixed Prompt Learning

To prevent tuning the whole pre-trained model in the few-shot downstream tasks, we adopt the prompt learning paradigm, *i.e.*, using a few trainable parameters to be the substitution of full-scale model parameters. **The goal of prompt design is to mitigate the gap between the pre-train objective and the downstream tasks.** It is highly challenging since there are both text (LM side) and structure information (GNN side) on TAG graphs. We try to investigate the joint prompting method, namely the *graph prompt* and the *text prompt*, the detailed design as follows. To simplify the discussion, we start with a target node  $v$  (e.g., the 0-th node in Figure 1). We further introduce the concept of the ego graph node set  $\mathcal{V}_{ego}$ . We select up to 100 first-order neighbors of the given node, along with the node itself, to form the node set of the ego graph. Then we define the induced ego graph from  $\mathcal{V}_{ego}$  as  $G_{ego}$ .

**3.3.1 Graph prompt design.** Inspired by the recently proposed graph prompt literature [28, 29], we aim to bridge gap through



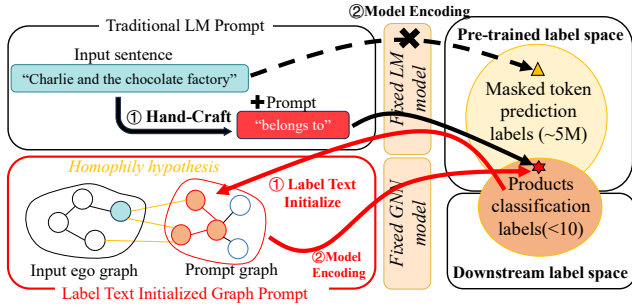
a small yet informative prompt graph. We first craft a prompt graph  $G_p = (\mathcal{V}_p, \mathcal{E}_p)$  serving as an ancillary component connected to the target ego node set  $\mathcal{V}_{ego}$ .  $G_p$  comprises  $\|G_p\|$  nodes (i.e., tokens), the internal edges between tokens can be defined as:

$$A_{G_p}(i, j) = \begin{cases} 1 & \text{if } \text{Sim}(\mathbf{p}_i, \mathbf{p}_j) \geq \sigma_{inner}, \\ 0 & \text{otherwise,} \end{cases} \quad (5)$$

where the  $\text{Sim}(\cdot)$  denotes the similarity function, and the  $\mathbf{p}_i$  denotes the feature of  $i$ -th tokens; in this paper, we use dot product as the implicit operation. Then, let the  $\mathbf{x}_k$  denotes on the  $k$ -th node feature of the  $G_{ego}$ , the inter-connectivity between  $G_p$  and  $G_{ego}$  is given by:

$$A_{G_{ego}, G_p}(k, j) = \begin{cases} 1 & \text{if } \text{Sim}(\mathbf{x}_k, \mathbf{p}_j) \geq \sigma_{inter}, \\ 0 & \text{otherwise.} \end{cases} \quad (6)$$

Here, the  $\sigma_{inner}$  and  $\sigma_{inter}$  are two hyper-parameters indicating the pruning value, i.e., the entries on the matrix higher than  $\sigma_{inner}$  and  $\sigma_{outer}$  are formed as edge respectively.



**Figure 2: Initialize tokens via the label text.** Crafting the prompt graph via downstream label texts under the **homophily hypothesis**, thereby making the target node embedding more aligned with the downstream label space.

**Initialize tokens via the label text.** While the prompt graph serves as an efficient substitute for full-scale model parameters, we empirically find that the naive  $G_p$  performs sub-optimally. We attribute this to the inconsistency in the pre-trained task objectives between our masked token prediction task and previous graph-based ones [28, 29] (e.g., the link prediction task).

Referencing Figure 2 as an example, during the pre-training phase where token prediction tasks are undertaken, the model is equipped to align input with the pre-trained label space. Consequently, the prompt should slightly modify the input distribution to approximate the downstream label space more closely. Hence, unlike traditional cumbersome and hand-crafted text prompt templates, we discovered that employing label text embeddings for prompt graph initialization is more direct and efficient. **Under the homophily hypothesis, the target node’s embedding is immediately positioned more closer to the desired label space.**

Specifically, we utilize the pre-trained  $f^{LM}$  as our label text encoder, each label text (e.g., “Animals”, “Science Fiction and Fantasy”, and “Activities, Crafts and Games”) can be represented as  $s_{y_i}$ . Then we can formalize the initialization of  $G_p$  as follows:

$$\mathbf{p}_i^{(0)} = \begin{cases} f^{LM}(s_{y_i}) & \text{if } i \leq N, \\ \text{Random Initialization} & \text{otherwise.} \end{cases} \quad (7)$$

Here, the  $\mathbf{p}_i^{(0)}$  denotes the initial embedding of  $i$ -th prompt node, and for the additional prompt nodes, we adopt the random initialization. Then the inner structure of each nodes will be constructed through the Eq. 5. The links among the prompt nodes and the target ego graph  $G_{ego}$  are built by the Eq. 6.

**3.3.2 Text prompt design.** As shown in Equation 3, we use the joint representation of the LM model output (text information) and the GNN model output (structure information). Therefore, to align the pre-train and downstream tasks, we use a similar strategy to concat both output features in the downstream scenario.

Then, instead of making intricate text-based prompts, we use a rather simple way: letting the second concat feature be a trainable parameter, namely  $\mathbf{w}_t \in \mathbb{R}^{1 \times d}$ , i.e., treating the LM model output itself as a tunable parameter. The initialization process can be defined similarly as:

$$\mathbf{w}_t^0 = f^{LM}(s_v), \quad (8)$$

the  $s_v$  shows the text sequence with the target node  $v$ . Our ablation study in Section 4.3 shows our text prompt’s effectiveness.

**3.3.3 Final prompt learning forward function.** Assuming we have a node ego graph  $G_v$  and its text  $s_v$ , the prompt learning forward function is:

$$\tilde{z}_v = \text{MLP}(\text{READOUT}(f^{GNN}(G_v; G_p)), \mathbf{w}_t). \quad (9)$$

The parameters of  $f^{LM}$  and  $f^{GNN}$  are fixed during prompt learning, the parameters of  $G_p$ , the  $\mathbf{w}_t$  and the task head (e.g., the MLP layers) are all trainable.

## 3.4 Model Inference

Finally, we incorporate the pretrained models  $f^{LM}$ ,  $f^{GNN}$ , and the prompt graph  $G_p$  in the model inference stage. We outline three inference methods. The first utilizes the trained LM, referred to as **P2TAG (LM)**. It computes the output of the [CLS] token of the LM for each node in the graph via a mini-batch manner. The second method involves further feeding the node [CLS] outputs from P2TAG (LM) into a GNN encoder, referred to as **P2TAG (GNN)**. In the training stage, the node representation of  $v_i$  is estimated according to the sampled neighborhood in the subgraph, and thus, varies with the sampling process. The randomness introduced by the sampler is unpleasant during the test stage. To get the exact node representation, we use the full neighbor sampling strategy when performing inference. That is, The node representation of  $v_i$  at the  $\ell$ -th layer is obtained by aggregating the representations of its full neighborhood at the  $(\ell - 1)$ -th layer. The computation process is performed layer by layer. Taking the node representations at the  $(\ell - 1)$ -th layer as the input, we can compute a batch of node representations at the  $\ell$ -th by first loading their full neighborhoods and then aggregating the representations of those neighbor nodes at the  $(\ell - 1)$ -th layer. The third method involves conducting few-shot node classification through prompting, as described in Section 3.3. The **P2TAG** inference loop of our framework is listed in Algorithm 2.

**Algorithm 2:** The inference of P2TAG

---

```

1 Input TAG $G = (\mathcal{V}, \mathcal{E}, \mathcal{S})$, batch size b , the number of GNN
 layers L ;
2 Output Node embeddings;
/* Obtain [CLS] embeddings from LMs for all
 nodes in a mini-batch manner. */
3 $X \leftarrow 0, i \leftarrow 0, j \leftarrow \min(b, n)$;
4 while $i < j$ do
5 $O = f^{LM}(S[i : j])$;
6 $X[i : j] = O[:, 0]$;
7 $i \leftarrow j, j \leftarrow \min(j + b, n)$
/* initialize the graph and text prompt. */
8 $p_i^{(0)} = \begin{cases} f^{LM}(s_{y_i}) & \text{if } i \leq N, \\ \text{Random Initialization} & \text{otherwise.} \end{cases}$
 $w_t^{(0)} \leftarrow f^{LM}(s_v)$;
/* Propagate node embeddings via a GNN encoder
 layer by layer. */
9 $H^{(0)} \leftarrow [X; G_p]$
10 for $\ell \leftarrow 1$ to L do
11 $H^{(\ell)} \leftarrow 0, i \leftarrow 0, j \leftarrow \min(b, n)$;
12 while $i < j$ do
13 $H^{(\ell)}[i : j] = \text{ReLU}(\hat{A}[i : j] H^{(\ell-1)} W^{(\ell)})$;
14 $i \leftarrow j, j \leftarrow \min(j + b, n)$
15 $H^{(L)} \leftarrow [H^{(L)}; w_t]$
16 $Z \leftarrow \text{MLP}(H^{(L)})$
17 return Probabilities of the classes Z ;

```

---

**Table 2: Dataset statistics of the used TAGs.**

| Dataset       | #Nodes    | #Edges     | Avg. Degree |
|---------------|-----------|------------|-------------|
| ogbn-arxiv    | 169,343   | 1,166,243  | 13.7        |
| ogbn-products | 2,449,029 | 61,859,140 | 50.5        |
| Children      | 76,875    | 1,554,578  | 40.4        |
| History       | 41,551    | 358,574    | 17.2        |
| Computers     | 87,229    | 721,081    | 16.5        |
| Photo         | 48,362    | 500,928    | 20.7        |

## 4 EXPERIMENT

In this section, we conduct comprehensive experiments to validate the effectiveness of our framework, particularly emphasizing the enhanced performance for few-shot node classification tasks.

### 4.1 Experimental Setup

**Datasets.** We conduct experiments across two categories of datasets. All these datasets are publicly available and widely used in TAG-related research. The statistics of the datasets are given in Table 2.

- **Open Graph Benchmark (OGB).** *ogbn-arxiv* is a directed citation graph between all computer science (CS) ArXiv papers indexed by Microsoft Academic Graph (MAG) [32]. We convert the abbreviated labels into their full forms as specified on

the arXiv website, such as transforming "NA" into "Numerical Analysis". *ogbn-products* is an undirected and unweighted graph representing an Amazon product co-purchasing network [1].

- **Amazon Review.** Amazon Review includes graphs composed of products, which are constructed based on co-purchased and co-viewed patterns. We use *Children*, *History*, *Computers* and *Photo* datasets compiled by [39]. For the four datasets, we extract bag-of-words features and utilize principal component analysis (PCA) to reduce dimensions, generating a 100-dimensional feature for each node.

In terms of the class split, ogbn-arxiv adopts the same partitioning strategy as TENT [34]. Meanwhile, for other datasets, we employ a random division approach. Notice that our P2TAG differs from meta-learning methods, as it does not require the establishment of specific training and validation tasks. Instead, we ensure fairness using the same test tasks as other baselines.

**Compared methods.** We choose three types of methods for comparison, all of which can address the problem of few-shot node classification. These include methods based on meta-learning on graphs, self-supervised pre-training methods on text attribute graphs, and methods following the paradigm of pre-training and prompting. Specifically, meta-learning methods such as GPN [7], G-Meta [16] and TENT [34]; GIANT [3], proposed P2TAG (LM) and P2TAG (GNN) as self-supervised pre-training methods; G2P2 [36] follows the same pre-training and prompting paradigm as proposed P2TAG. In the pre-training phase, we concurrently train both the LM and GNN, considering two scenarios during inference. P2TAG (LM) derives vector representations from the original text of nodes. Meanwhile, P2TAG (GNN) leverages these node representations as inputs to the GNN, resulting in the output obtained.

**Evaluation.** GPN, G-Meta, and TENT are methods based on meta-learning, where they are learned on train tasks and are subsequently evaluated on test tasks. In contrast, the proposed P2TAG along with G2P2 and GIANT, are evaluated exclusively on test tasks. Specifically, P2TAG (LM), P2TAG (GNN), and GIANT infer the classes of nodes in the query set by fitting a logistic regression classifier on the support set. Furthermore, P2TAG and G2P2 extend their approach by constructing prompts through the support set. For a detailed comparison, we conduct experiments in four different settings for OGB datasets: 5-way 3-shot, 5-way 5-shot, 10-way 3-shot, and 10-way 5-shot. Considering the number of labels, we use three settings for Amazon Review datasets: 3-way 3-shot, 3-way 5-shot, and 3-way 10-shot.

**Parameter configuration.** For our framework, the selection of LMs and GNNs is flexible. In our experiment, we choose a representative LM – DeBERTa-base and a powerful GAT model for the main experiments. The DeBERTa-base is a 100M-parameter pre-trained language model with a hidden size of 768. We keep the same hidden size of the GAT model with DeBERTa-base. We also explore other LMs in the ablation studies. We use AdamW optimizer [23] with learning rate  $\text{lr} = 1e-5$  for model optimization. We run 3 epochs for all datasets. We construct 5 groups of test tasks for each  $N$ -way  $K$ -shot setting, with each group consisting of 50 tasks specifically formulated from the test set label. In each task, the support set length is  $K$ , and the query set length  $Q$  is set to 10. We calculate

**Table 3: The accuracy of few-shot node classification on two OGB datasets.**

| Dataset      | ogbn-arxiv        |                   |                   |                   | ogbn-products     |                   |                   |                   |
|--------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
|              | 5-way 3-shot      | 5-way 5-shot      | 10-way 3-shot     | 10-way 5-shot     | 5-way 3-shot      | 5-way 5-shot      | 10-way 3-shot     | 10-way 5-shot     |
| OGB features | 53.24±0.77        | 59.31±0.71        | 38.62±0.85        | 44.59±0.61        | 44.15±1.30        | 50.57±1.83        | 32.79±0.88        | 38.75±0.61        |
| GPN          | 61.65±0.67        | 66.34±0.78        | 47.36±0.25        | 53.60±0.86        | 64.04±1.92        | 69.86±1.90        | 52.68±0.66        | 58.36±0.45        |
| G-Meta       | 60.06±1.56        | 63.77±2.53        | 48.36±0.94        | 52.78±1.04        | 64.06±1.18        | 66.02±0.71        | 51.17±1.48        | 57.81±1.40        |
| TENT         | 62.24±0.35        | 65.80±0.37        | 45.08±0.54        | 51.84±0.86        | 62.69±2.78        | 66.63±1.48        | 47.63±0.59        | 51.53±0.90        |
| GIANT        | 73.50±1.59        | 78.18±1.25        | 60.05±1.20        | 64.96±0.69        | 69.44±1.40        | 75.46±1.53        | 59.23±0.64        | 65.67±0.38        |
| G2P2         | 70.93±1.29        | 73.07±1.22        | 56.73±0.68        | 59.47±0.69        | > 30 days         | > 30 days         | > 30 days         | > 30 days         |
| P2TAG (LM)   | <u>78.14±1.18</u> | <u>81.33±1.34</u> | 65.21±1.26        | 69.90±0.94        | 75.38±1.27        | 78.66±1.55        | 61.87±0.91        | 67.06±0.49        |
| P2TAG (GNN)  | 77.73±1.76        | 81.31±0.84        | <u>66.21±0.82</u> | <u>70.42±0.62</u> | <b>80.95±1.61</b> | <u>83.65±1.24</u> | <u>70.04±0.33</u> | <u>73.67±0.35</u> |
| P2TAG        | <b>80.18±0.27</b> | <b>84.76±0.89</b> | <b>67.81±0.97</b> | <b>70.88±0.96</b> | <u>80.54±1.59</u> | <b>83.77±1.45</b> | <b>71.14±0.26</b> | <b>75.99±0.41</b> |

**Table 4: The accuracy of few-shot node classification on four Amazon Review datasets.**

| Dataset       | Children          |                   |                   | History           |                   |                   |
|---------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
|               | 3-way 3-shot      | 3-way 5-shot      | 3-way 10-shot     | 3-way 3-shot      | 3-way 5-shot      | 3-way 10-shot     |
| node features | 39.68±0.61        | 43.33±1.78        | 49.00±0.77        | 37.05±1.63        | 38.16±1.33        | 40.82±1.07        |
| GPN           | 54.77±0.42        | 60.51±1.16        | 63.89±0.93        | 38.48±0.74        | 40.63±0.82        | 43.88±1.01        |
| G-Meta        | 54.00±2.38        | 57.76±1.43        | 61.62±1.52        | 40.41±1.32        | 41.11±0.76        | 42.50±0.86        |
| TENT          | 53.23±0.76        | 60.52±1.73        | 64.32±0.81        | 38.30±0.72        | 37.73±0.43        | 41.47±2.05        |
| G2P2          | 51.33±1.84        | 55.88±1.18        | 58.85±0.84        | 47.21±0.37        | 49.89±0.97        | 55.63±1.11        |
| P2TAG (LM)    | <u>73.64±1.13</u> | <u>77.43±1.46</u> | 80.53±1.27        | <u>68.88±2.10</u> | <u>72.36±1.28</u> | <u>76.81±0.47</u> |
| P2TAG (GNN)   | 73.20±1.34        | 76.71±0.82        | 80.67±1.54        | 63.33±1.00        | 68.15±1.41        | 72.68±0.98        |
| P2TAG         | <b>77.35±1.09</b> | <b>81.43±1.24</b> | <b>84.67±1.90</b> | <b>72.35±1.85</b> | <b>76.19±1.25</b> | <b>80.96±0.74</b> |

| Dataset       | Computers         |                   |                   | Photo             |                   |                   |
|---------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
|               | 3-way 3-shot      | 3-way 5-shot      | 3-way 10-shot     | 3-way 3-shot      | 3-way 5-shot      | 3-way 10-shot     |
| node features | 37.19±1.36        | 40.11±1.72        | 44.55±1.11        | 40.43±0.64        | 44.41±1.68        | 50.57±1.38        |
| GPN           | 71.73±1.67        | 70.47±1.72        | 71.56±1.91        | 74.73±1.75        | 76.01±1.12        | 70.77±0.97        |
| G-Meta        | 71.58±2.19        | 71.14±2.52        | 72.36±0.63        | 69.68±1.08        | 72.41±1.09        | 73.54±1.22        |
| TENT          | 60.68±1.94        | 62.13±1.17        | 65.96±1.92        | 68.24±1.20        | 71.20±1.30        | 72.19±2.03        |
| G2P2          | 61.36±1.20        | 65.13±1.90        | 68.56±1.26        | 70.97±1.51        | 72.49±1.23        | 76.52±2.37        |
| P2TAG (LM)    | 66.73±2.16        | 71.17±0.54        | 74.77±1.06        | 73.49±1.77        | 75.68±0.99        | 77.89±0.94        |
| P2TAG (GNN)   | 86.95±1.71        | <u>88.67±0.51</u> | 92.15±0.88        | 86.47±1.45        | 88.20±0.49        | <b>90.72±0.58</b> |
| P2TAG         | <b>87.24±1.38</b> | <b>89.55±0.54</b> | <b>93.76±0.63</b> | <b>86.73±2.14</b> | <b>88.73±0.98</b> | <u>90.68±0.76</u> |

the average performance within each group and report the overall mean and standard deviation across these 5 groups of tasks.

## 4.2 Performance Analysis

Our main results are summarized in Table 3 and Table 4. The proposed P2TAG (LM) outperforms the meta-learning methods, increasing the average accuracy on six datasets from +2.27% ~ +35.51%; the P2TAG (GNN) achieves an average improvement with +16.38% ~ +27.55%; the P2TAG performs best on most datasets, with an average improvement of +18.98% ~ +35.98%. Compared with the pre-training method that also utilize raw text such as GIANT and G2P2, our P2TAG still has better performance, which demonstrates its effectiveness. On the ogbn-arxiv dataset, although G2P2 constructs prompts on the basis of pre-training, its performance is lower than that of GIANT, P2TAG (LM) and P2TAG (GNN). This underscores the importance of a well pre-trained model for

few-shot node classification. Pre-training on TAGs often comes with increased time expenditure, such as G2P2 taking over 30 days on the ogbn-products dataset. The proposed P2TAG employs a more general approach to jointly train LM and GNN, consuming less than one day on this dataset. The P2TAG, employing mixed prompts, further enhances performance on the pre-trained model P2TAG (LM) and P2TAG (GNN). It secures the best outcomes on most datasets, with an average enhancement of 2.89% compared to P2TAG (GNN), demonstrating the effectiveness of prompts. On the History dataset, P2TAG (LM) achieves the second-best results; however, after passing through the GNN encoder, P2TAG (GNN) experiences an average of 4.63% decrease. This might be attributed to the quality of the topological structure within the data. This hypothesis is further supported by the minimal improvement observed when comparing meta-learning methods that do not utilize raw texts (GPN, G-Meta, TENT) to the node features method. The

**Table 5: Ablation study of language models on ogbn-arxiv dataset.** We choose various LMs, then report the classification accuracy achieved with P2TAG (LM).

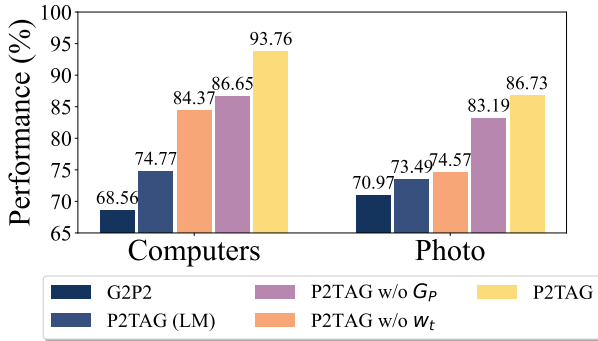
| Setting       | 5-way             |                   | 10-way            |                   |
|---------------|-------------------|-------------------|-------------------|-------------------|
|               | 3-shot            | 5-shot            | 3-shot            | 5-shot            |
| DeBERTa-base  | 78.14±1.18        | 81.33±1.34        | 65.21±1.26        | 69.90±0.94        |
| DeBERTa-large | <b>78.64±1.24</b> | <b>83.55±1.05</b> | <b>66.00±1.12</b> | <b>71.96±1.06</b> |
| e5-v2-base    | 78.62±1.00        | 82.44±1.22        | 65.79±0.64        | 71.02±0.60        |
| e5-v2-large   | 77.31±1.68        | 81.87±1.09        | 64.60±0.72        | 70.08±0.60        |

proposed graph prompt improvement strategy, which initializes through label text, mitigates this issue. On the History dataset, P2TAG further enhances performance compared to P2TAG (LM), achieving an average improvement of 3.8% across three settings.

### 4.3 Ablation Studies

In the previous sections, we demonstrate the powerful performance of the P2TAG (LM), P2TAG (GNN), and P2TAG. This part analyzes the impact of the different types of prompting and LMs.

**Effect of LMs.** To better analyze the impact of LMs, we explore other LMs such as e5-v2-base with 110M parameters [33]. We also try larger LMs such as DeBERTa-large with 350M parameters and e5-v2-large with 330M parameters. The results are reported in Table 5. Generally, the results of LMs are quite similar, with differences within 1.5%. The reason that e5-large does not achieve better results may be attributed to insufficient training iterations. This paper selects DeBERTa-base, intending to address the joint learning problem of LMs and GNNs in a more general manner. There remains room for further exploration in the specific choice of LMs.

**Figure 3: Effect of different prompt types.** We omit label initialization on the prompt graph (w/o  $P_G$ ) and node text embedding (w/o  $w_t$ ) and report the classification results on two datasets.

**Effect of different prompt types.** In Section 3.3, we discussed two significant enhancements. Here, we further investigate these improvements by dissecting them and presenting the findings in Figure 3. We use the P2TAG w/o  $G_p$  when the label text embedding is not utilized to initialize our prompt graph, and the P2TAG w/o  $w_t$  when node text embedding is not employed to bolster the prompted graph embedding (*i.e.*, the graph embedding is directly inputted into

**Table 6: The accuracy of classification by P2TAG with different hyperparameters on ogbn-arxiv dataset.** We report the results with P2TAG (GNN) and P2TAG, where the underscored annotations indicate the parameters used in the main result tables.

| Setting                         |             | 5-way 5-shot      | 10-way 5-shot     |
|---------------------------------|-------------|-------------------|-------------------|
| P2TAG (GNN):<br>sequence length | 32          | 79.74±1.32        | 67.81±0.48        |
|                                 | 64          | 81.40±1.26        | 69.61±0.52        |
|                                 | <u>128</u>  | 81.31±0.84        | <b>70.42±0.62</b> |
| P2TAG (GNN):<br>mask rate       | 256         | <b>81.61±0.95</b> | 70.40±0.70        |
|                                 | 0.15        | 43.93±0.41        | 27.02±0.48        |
|                                 | 0.30        | 60.66±1.00        | 46.34±0.48        |
|                                 | 0.50        | 80.48±1.10        | 69.96±0.69        |
| P2TAG (GNN):<br>walk length     | <u>0.75</u> | <b>81.31±0.84</b> | <b>70.42±0.62</b> |
|                                 | 5           | 81.28±1.07        | 70.18±0.54        |
|                                 | <u>10</u>   | 81.31±0.84        | <b>70.42±0.62</b> |
|                                 | 20          | 81.28±1.43        | 70.20±0.59        |
| P2TAG:<br>token num             | 50          | <b>81.70±0.61</b> | 70.11±0.74        |
|                                 | 2           | 84.31±0.95        | 70.04±1.42        |
|                                 | <u>5</u>    | <b>84.76±0.89</b> | 70.16±0.78        |
|                                 | <u>10</u>   | 84.35±1.07        | <b>70.88±0.96</b> |
|                                 | 20          | 84.06±1.12        | 69.71±0.69        |

MLP layers). We report the results on Computers and Photo datasets. The results reveal that our prompting approach outperforms both the P2TAG (LM) and baseline G2P2. Furthermore, the absence of  $P_G$  has a marginally greater impact than omitting  $w_t$  across both datasets, underscoring the importance of structural effectiveness in the TAG context.

### 4.4 Hyperparameter Analysis

In this part, we conduct a comprehensive analysis of hyperparameters to elucidate their impact on the performance of our framework. We analyze the four key hyperparameters in our experiments. The "sequence length" means the reserved length of raw texts. The "mask rate" represents the proportion of masked tokens in training. The "walk length" means the length of random walks used in the GraphSAINT sampler for mini-batch training. The "token num" denotes the number of tokens in the graph prompt. The results of hyperparameter experiments are shown in Table 6.

**Preprocessing: length of truncated sequences.** Since each batch will be fed to both the language model and the GNN encoder in training, we need to increase the batch size to better leverage the ability of GNNs. However, the batch size is limited due to the large memory cost of the LM. Therefore, considering both training efficiency and information loss, we explore using a smaller truncated length in our experiments to utilize a larger batch size.

**Pre-training: mask rate.** A smaller mask rate represents an easier self-supervised task and might make the model overfitted. To make the self-supervised task harder, we need to choose a large enough mask rate to let the model learn better. From the results, we find that our framework achieves the best performance with a 75% mask rate on the ogbn-arxiv dataset.

**Sampling: lengths of random walks.** For each batch, the sampler conducts multiple random walks from different root nodes with the same length. The GNN encoder of our framework relies on the topology of the sampled mini-batch graph to propagate and aggregate the information between nodes. Therefore, the length of random walks used to construct the mini-batch will influence the model performance.

**Prompting: number of tokens.** Incorporating the prompt graph as a trainable parameter highlights the importance of token quantity for downstream tasks. Intuitively, more tokens might create a richer internal structure, offering enhanced information. However, managing more parameters often complicates convergence in few-shot scenarios. To accurately incorporate the textual information of labels into prompting, we typically set the number of tokens to match the number of classes in the task. Surprisingly, 2 tokens also yield results, implying a minimal internal structure. This suggests the effectiveness of token features in improving performance.

## 5 RELATED WORK

In this section, we introduce the related work, including graph representation learning and few-shot node classification.

### 5.1 Graph Representation Learning

Graph neural networks (GNNs) provide the foundation for applying deep learning on graphs and yielding good results on several downstream tasks. The earlier works [5, 31, 35, 38] perform convolution on small-scale graphs using all topological relations in a semi-supervised way. The subsequent works focus on sampling strategies [2, 10, 48] and model architecture [26, 37] to enhance the scalability of GNNs and apply them to large-scale graphs. GraphSAGE [10] first proposed the idea of neighborhood sampling, and later it was applied in a real-world recommendation system by PinSAGE [45]. GraphSAINT [48], first samples subgraphs [20] and runs full-batch GNNs on sampled subgraphs. Additionally, there are works [52, 53] focusing on the hierarchy within the graph to obtain better representations. The self-supervised learning methods on GNNs are developed via contrastive and generative ways. The contrastive methods [11, 25, 30, 46, 49, 57] adopt data augmentation with or without negative sampling to construct samples to optimize the contrastive loss under different augmentations. GRACE [57] aims to maintain node uniformity across views. BGRL [30] designs two encoders for two views with data augmentation. Some studies focus on graph generative learning. GraphMAE [15] and GraphMAE2 [14] are proposed to reconstruct masked attributes of corrupted nodes for representation learning.

The semantics and graph topology of TAGs can express real-world relationships between entities or objects. Most GNNs do not consider the text processing in the TAG but directly use the numerical features, which are generated through text encoding, as attributes of nodes. Recent studies [3, 44] utilize graph topology to enhance the representation during text pre-training. Despite the promising results, the language model is still independent of the GNNs [21, 56]. GraphFormer [41] designs a nested architecture of GNNs and Transformers. GLEM [51] implements the fusion of LMs and GNNs on large-scale text-attributed graphs with the variational expectation-maximization framework. Although these

methods progress in integrating LMs and GNNs, there is still a lack of discussions on self-supervised learning on large-scale text-attributed graphs. Our proposed pre-train framework enhances LMs utilizing GNNs and achieves joint training with the objective of self-supervision. Yan et al. [39] release the benchmark on TAGs, and provide multiple clean TAG datasets. These datasets also lay a solid foundation for our experiments.

### 5.2 Few-shot Node Classification

Few-shot node classification on graphs to categorize nodes within a graph with limited labeled nodes. Drawing inspiration from the success of meta-learning in few-shot classification tasks within computer vision [9, 27], several studies apply meta-learning techniques to graph-based tasks. GFL [43] and GPN [7] utilize prototype networks to learn the distance from nodes to classes. Meta-GNN [55] optimizes the graph neural network with model-agnostic meta-learning. G-Meta [16] addresses the meta-learning problem on both single and multiple graphs by extracting local subgraphs. TENT [34] enhances the model's generalization capabilities by adapting at three levels: nodes, edges, and tasks.

With the increasing attention towards LLMs, several works attempt to follow the paradigm of pre-training and prompt tuning to address the problem of few-shot node classification. Prog [29] incorporates both node-level and edge-level tasks by constructing prompts at the graph level, but lacks analysis of the text attribute. ENG [47] employs LLM with designed prompts to refine node texts and reconstructs the adjacency matrix semantically, which is then used as input for the GNN. G2P2 [36] enhances text information through graph structure, aligning text representations in three forms during the pre-training phase. In the tuning phase, it uses the neighborhood text of the target node and label text to generate the initial parameters of prompts while freezing the parameters of the LM and GNN during the tuning process.

## 6 CONCLUSION

Our paper focuses on few-shot node classification on the TAG. We address this problem by employing a graph pre-training and prompting approach. The proposed framework P2TAG utilizes a masked language modeling objective for the joint training of the language model and GNN model. We also propose a new prompting method that mixes graph and text information, enabling the pre-trained model on TAG to better adapt to downstream few-shot node classification tasks. We conduct experiments on six real-world TAG datasets and our P2TAG framework achieves state-of-the-art results on the six datasets with +18.98% ~ +35.98% improvement.

## 7 ACKNOWLEDGEMENT

This work is supported by National Key R&D Program of China 2021ZD0113304, Natural Science Foundation of China (NSFC) 62276148 and 62425601, CCF-Zhipu AI Large Model Fund (Grant 202213), Zhipu AI - Anhui University Joint Research Center on Foundation Model and the University Synergy Innovation Program of Anhui Province (GXXT-2023-050), the New Cornerstone Science Foundation through the XPLOER PRIZE, and Tsinghua-Bosch Joint ML Center.

## REFERENCES

- [1] K. Bhatia, K. Dahiya, H. Jain, P. Kar, A. Mittal, Y. Prabhu, and M. Varma. 2016. The extreme classification repository: Multi-label datasets and code. <http://manikvarma.org/downloads/XC/XMLRepository.html>
- [2] Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. 2019. Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks. In *KDD'19*. 257–266.
- [3] Eli Chien, Wei-Cheng Chang, Cho-Jui Hsieh, Yu Hsiang-Fu, Jiong Zhang, Olga Milenkovic, and Inderjit S. Dhillon. 2022. Node Feature Extraction by Self-Supervised Multi-scale Neighborhood Prediction. In *ICLR'22*.
- [4] Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. 2020. Electra: Pre-training text encoders as discriminators rather than generators. In *ICLR'20*.
- [5] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. 2016. Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering. In *NeurIPS'16*. 3837–3845.
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL-HLT'19 (1)*. 4171–4186.
- [7] Kaize Ding, Jianling Wang, Jundong Li, Kai Shu, Chenghao Liu, and Huan Liu. 2020. Graph prototypical networks for few-shot learning on attributed networks. In *CIKM'20*.
- [8] Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. GLM: General language model pretraining with autoregressive blank infilling. In *ACL'22*. 320–335.
- [9] Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML'17*.
- [10] William L. Hamilton, Zitao Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In *NeurIPS'17*.
- [11] Kaveh Hassani and Amir Hosein Khasahmadi. 2020. Contrastive multi-view representation learning on graphs. In *ICML'20*. 4116–4126.
- [12] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. DeBERTa: Decoding-enhanced BERT with Disentangled Attention. In *ICLR'21*.
- [13] Xiaoxin He, Xavier Bresson, Thomas Laurent, Adam Perold, Yann LeCun, and Bryan Hooi. 2023. Harnessing explanations: LLM-to-LM interpreter for enhanced text-attributed graph representation learning. *arXiv preprint arXiv:2305.2305* (2023).
- [14] Zhenyu Hou, Yufei He, Yukuo Cen, Xiao Liu, Yuxiao Dong, Evgeny Kharlamov, and Jie Tang. 2023. GraphMAE2: A Decoding-Enhanced Masked Self-Supervised Graph Learner. In *WWW'23*. 737–746.
- [15] Zhenyu Hou, Xiao Liu, Yukuo Cen, Yuxiao Dong, Hongxia Yang, Chunjie Wang, and Jie Tang. 2022. GraphMAE: Self-Supervised Masked Graph Autoencoders. In *KDD'22*. 594–604.
- [16] Kexin Huang and Marinka Zitnik. 2020. Graph meta learning via local subgraphs. *NeurIPS'20* (2020).
- [17] Xuanwen Huang, Kaiqiao Han, Dezheng Bao, Quanjin Tao, Zhisheng Zhang, Yang Yang, and Qi Zhu. 2023. Prompt-based node feature extractor for few-shot learning on text-attributed graphs. *arXiv preprint arXiv:2309.02848* (2023).
- [18] Shaoxiong Ji, Shirui Pan, Erik Cambria, Pekka Marttinen, and S Yu Philip. 2021. A survey on knowledge graphs: Representation, acquisition, and applications. *IEEE transactions on neural networks and learning systems* 33, 2 (2021), 494–514.
- [19] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. Albert: A lite bert for self-supervised learning of language representations. In *ICLR'20*.
- [20] Jure Leskovec and Christos Faloutsos. 2006. Sampling from large graphs. In *KDD'06*.
- [21] Chaozhao Li, Bochen Pang, Yuming Liu, Hao Sun, Zheng Liu, Xing Xie, Tianqi Yang, Yanling Cui, Liangjie Zhang, and Qi Zhang. 2021. Adsgnn: Behavior-graph augmented relevance modeling in sponsored search. In *SIGIR'21*. 223–232.
- [22] Yinhan Liu, MyLe Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [23] Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101* (2017).
- [24] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013).
- [25] Jiezhong Qiu, Qibin Chen, Yuxiao Dong, Jing Zhang, Hongxia Yang, Ming Ding, Kuansan Wang, and Jie Tang. 2020. GCC: Graph contrastive coding for graph neural network pre-training. In *KDD'20*. 1150–1160.
- [26] Emanuele Rossi, Fabrizio Frasca, Ben Chamberlain, Davide Eynard, Michael Bronstein, and Federico Monti. 2020. Sign: Scalable inception graph neural networks. *arXiv preprint arXiv:2004.11198* 7 (2020), 15.
- [27] Jake Snell, Kevin Swersky, and Richard S. Zemel. 2017. Prototypical networks for few-shot learning. In *NeurIPS'17*.
- [28] Mingchen Sun, Kaixiong Zhou, Xin He, Ying Wang, and Xin Wang. 2022. Gpnt: Graph pre-training and prompt tuning to generalize graph neural networks. In *KDD'22*. 1717–1727.
- [29] Xiangguo Sun, Hong Cheng, Jia Li, Bo Liu, and Jihong Guan. 2023. All in One: Multi-Task prompting for graph neural networks. In *KDD'23*.
- [30] Shantanu Thakoor, Corentin Tallec, Mohammad Gheshlaghi Azar, Mehdi Azabou, Eva L Dyer, Remi Munos, Petar Veličković, and Michal Valko. 2022. Large-scale representation learning on graphs via bootstrapping. In *ICLR'22*.
- [31] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. In *ICLR'18*.
- [32] Kuansan Wang, Zhihong Shen, Chiyuan Huang, Chieh-Han Wu, Yuxiao Dong, and Anshul Kanakia. 2020. Microsoft academic graph: When experts are not enough. *Quantitative Science Studies* 1, 1 (2020), 396–413.
- [33] Liang Wang, Nan Yang, Xiaolong Huang, Bingxin Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533* (2022).
- [34] Song Wang, Kaize Ding, Chuxu Zhang, Chen Chen, and Jundong Li. 2022. Task-adaptive few-shot node classification. In *KDD'22*.
- [35] Max Welling and Thomas N Kipf. 2017. Semi-supervised classification with graph convolutional networks. In *ICLR'17*.
- [36] Zhihao Wen and Yuan Fang. 2023. Augmenting low-Resource text classification with graph-grounded pre-training and prompting. In *SIGIR'23*.
- [37] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. 2019. Simplifying graph convolutional networks. In *ICML'19*.
- [38] Bingbing Xu, Huawei Shen, Qi Cao, Yunqi Qiu, and Xueqi Cheng. 2019. Graph wavelet neural network. In *ICLR'19*.
- [39] Hao Yan, Chaozhao Li, Ruosong Long, Chao Yan, Jianan Zhao, Wenwen Zhuang, Jun Yin, Peiyan Zhang, Weihao Han, Hao Sun, et al. 2023. A comprehensive study on text-attributed graphs: Benchmarking and rethinking. In *NeurIPS'23*.
- [40] Beining Yang, Kai Wang, Qingyun Sun, Cheng Ji, Xingcheng Fu, Hao Tang, Yang You, and Jianxin Li. 2023. Does Graph Distillation See Like Vision Dataset Counterpart? In *NeurIPS'23*.
- [41] Junhan Yang, Zheng Liu, Shitao Xiao, Chaozhao Li, Defu Lian, Sanjay Agrawal, Amit Singh, Guangzhong Sun, and Xing Xie. 2021. GraphFormers: GNN-nested Transformers for Representation Learning on Textual Graph. In *NeurIPS'21*.
- [42] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. 2019. XLNet: Generalized Autoregressive Pretraining for Language Understanding. In *NeurIPS'19*. 5754–5764.
- [43] Huaxiu Yao, Chuxu Zhang, Ying Wei, Meng Jiang, Suhang Wang, Junzhou Huang, Nitesh Chawla, and Zhenhui Li. 2020. Graph few-shot learning via knowledge transfer. In *AAAI'20*.
- [44] Michihiro Yasunaga, Jure Leskovec, and Percy Liang. 2022. LinkBERT: Pretraining Language Models with Document Links. In *ACL'22*.
- [45] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L. Hamilton, and Jure Leskovec. 2018. Graph convolutional neural networks for web-scale recommender systems. In *KDD'18*.
- [46] Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. 2020. Graph Contrastive Learning with Augmentations. In *NeurIPS'20*.
- [47] Jianxiang Yu, Yuxiang Ren, Chenghua Gong, Jiaqi Tan, Xiang Li, and Xuecang Zhang. 2023. Empower text-attributed graphs learning with large language models (llms). *arXiv preprint arXiv:2310.09872* (2023).
- [48] Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor K. Prasanna. 2020. GraphSAINT: Graph Sampling Based Inductive Learning Method. In *ICLR'20*.
- [49] Jiaqi Zeng and Pengtao Xie. 2021. Contrastive self-supervised learning for graph classification. In *AAAI'21*, Vol. 35. 10824–10832.
- [50] Hengrui Zhang, Qitian Wu, Junchi Yan, David Wipf, and Philip S. Yu. 2021. From Canonical Correlation Analysis to Self-supervised Graph Neural Networks. In *NeurIPS'21*. 76–89.
- [51] Jianan Zhao, Meng Qu, Chaozhao Li, Hao Yan, Qian Liu, Rui Li, Xing Xie, and Jian Tang. 2023. Learning on Large-scale Text-attributed Graphs via Variational Inference. In *ICLR'23*.
- [52] Shu Zhao, Jialin Chen, Jie Chen, Yanping Zhang, and Jie Tang. 2022. Hierarchical label with imbalance and attributed network structure fusion for network embedding. *AI Open* 3 (2022), 91–100.
- [53] Shu Zhao, Ziwei Du, Jie Chen, Yanping Zhang, Jie Tang, and S Yu Philip. 2023. Hierarchical representation learning for attributed networks. *IEEE Transactions on Knowledge and Data Engineering* 35, 3 (2023), 2641–2656.
- [54] Yizhen Zheng, Shirui Pan, Vincent Cs Lee, Yu Zheng, and Philip S Yu. 2022. Rethinking and Scaling Up Graph Contrastive Learning: An Extremely Efficient Approach with Group Discrimination. In *NeurIPS'22*.
- [55] Fan Zhou, Chengtai Cao, Kunpeng Zhang, Goce Trajcevski, Ting Zhong, and Ji Geng. 2019. Meta-gnn: On few-shot node classification in graph meta-learning. In *CIKM'2019*.
- [56] Jason Zhu, Yanling Cui, Yuming Liu, Hao Sun, Xue Li, Markus Pelger, Tianqi Yang, Liangjie Zhang, Ruofei Zhang, and Huasha Zhao. 2021. Textgnn: Improving text encoder via graph neural network in sponsored search. In *WWW'21*. 2848–2857.
- [57] Yanqiao Zhu, Yichen Xu, Feng Yu, Qiang Liu, Shu Wu, and Liang Wang. 2020. Deep graph contrastive representation learning. *arXiv preprint arXiv:2006.04131* (2020).

## A APPENDIX

### A.1 Implementation Notes

**Running environment.** The experiments are conducted on a Linux machine with AMD EPYC 7642 48-Core Processor, 1T RAM, and 8 NVIDIA A100 (80G). Our code is implemented with PyTorch 1.10 and Python 3.8.

**Model Configuration.** Our framework involves the joint training of LM and GNN, with some of the key parameters presented in Table 7. For more detailed configurations, please find in our code.

Table 7: Key parameters of the P2TAG.

| Parameters                   | Value                         |
|------------------------------|-------------------------------|
| language model               | microsoft/deberta-base        |
| sequence length              | 128                           |
| mask rate                    | 0.75                          |
| GNN encoder                  | graph attention network       |
| GNN hidden size              | 768                           |
| walk length                  | 10                            |
| number of root               | 10                            |
| length of query set          | 10                            |
| number of test groups        | 5                             |
| number of task in each group | 50                            |
| number of tokens             | the number of classes of task |
| prompt epochs                | 50                            |

### A.2 Comparing P2TAG with GPrompt

Though P2TAG shares a similar task and intuition with GPrompt [17], there are significant differences in the prompt design. Since the prompting parts are crucial for the pre-training and prompting paradigm, solely considering the task and intuition may not be fair. Therefore, we further elaborate our design here:

**Difference in Fundamental Settings of Prompting.** GPrompt uses the entire GNN model as the prompt, which requires fine-tuning the GNN for each downstream task. In contrast, we freeze the GNN parameters and instead use an auxiliary prompt graph to alter the input. This approach, we believe, introduces a significant change to the community. Our method not only preserves the knowledge encapsulated within the GNN model but also reduces the number of training parameters to a minimal set of node embeddings, thereby enhancing computational efficiency.

**Novel Label Text Prompt Initialization Strategy.** Furthermore, we introduce a novel strategy for initializing the prompt graph by encoding the label texts to embeddings, which is a unique and effective element in the TAG scenario. By incorporating the label text embedding, we can easily initialize the prompt graph and align the label space more closely with the downstream task. For example, a straightforward text prompt design might be: “Given labels *student*, *teacher*, *worker*, the person with attributes [feature] belongs to [cls]”. In the graph domain, we initialize the prompt graph with these label text embeddings. As a result, the target node embedding naturally aligns closely with the label text embedding via message

passing, replicating the prompt initialization process. Moreover, the label text embeddings within the prompt graph can adaptively influence the target node, potentially offering more effectiveness than static text prompts with uniform contributions.

### A.3 Baselines

We conduct five baselines, among which GPN, G-Meta, and TENT are based on meta-learning and do not utilize raw texts. GIANT employs raw texts for self-supervised learning on graphs. G2P2 also leverages raw texts and follows the paradigm of pre-training and prompting. All of them address the few-shot node classification problem. For each  $N$ -way  $K$ -shot few-shot setting, we sample 5 groups, each containing 50 tasks for testing. To ensure a fair comparison, we save the constructed tasks, and all baselines, including our proposed model, are tested on these tasks. G-Meta is an exception due to its more complex sampling process, but we ensure the same number of testing tasks.

- GPN<sup>1</sup> applies the prototype network from meta-learning to graphs. We set the number of episodes for training to 2000.
- G-Meta<sup>2</sup> implements the meta-learning method for graphs at three levels using local subgraphs. To match our number of test tasks, the batch size was set to 50 during the generation of test tasks, and this process is repeated 5 times.
- TENT<sup>3</sup> tunes the model across multiple tasks. We modify this method by Deep Graph Library (DGL) to enable its application on larger-scale graphs, such as the ogbn-products.
- GIANT<sup>4</sup> performs extreme multi-label classification on the raw texts, resulting in stronger node features. We directly download their pre-trained node features for evaluation on the ogbn-arxiv and ogbn-products datasets.
- G2P2<sup>5</sup> enhances text representation with graph structure, aligning at three levels, and subsequently performs few-shot classification through prompting. We set the hidden layer dimension to 128 and the number of epochs to 3.

### A.4 Analysis of P2TAG Inference Cost

We analyze the inference cost of P2TAG and baselines with respect to the text sequence length  $K$ , embedding dimension  $d$ , number of nodes  $N$ , number of nodes  $b$  in the ego-graph, and the number of graph tokens  $t$ .

The cost of P2TAG (LM) stems from the language model (LM), with a complexity of  $O(N \times (Kd^2 + K^2d))$ . P2TAG (GNN) additionally utilizes a GNN encoder with negligible computation overhead. The complexity is  $O(N \times (Kd^2 + K^2d))$ . P2TAG includes two stages: pre-training and prompting. The complexity is  $O((N + b + t) \times (Kd^2 + K^2d))$ . The complexity of GIANT is  $O(N \times (Kd^2 + K^2d) + rd \log(N))$ , where  $r$  is a hyperparameter. The complexity of G2P2 is  $O((N + c + 1) \times (Kd^2 + K^2d))$ , where  $c$  represents the number of neighbors.

Overall, the proposed P2TAG exhibits computational costs similar to other baselines, yet achieves improved performance in few-shot classification.

<sup>1</sup>[https://github.com/kaize0409/GPN\\_Graph-Few-shot](https://github.com/kaize0409/GPN_Graph-Few-shot)

<sup>2</sup><https://github.com/mims-harvard/G-Meta>

<sup>3</sup><https://github.com/SongW-SW/TENT>

<sup>4</sup><https://github.com/amzn/pecos/tree/mainline/examples/giant-xrt>

<sup>5</sup><https://github.com/WenZhihao666/G2P2>



### A.5 Ablation Study about the GNN Backbone

We use GAT due to its ability to adaptively transfer information from our prompt graph to the target node, matching our graph prompt design. Our ablation study on GCN and GraphSAGE within P2TAG on the ogbn-arxiv dataset shows:

**Table 8: Ablation study results on the ogbn-arxiv dataset.**

| GNN Encoder | 5-way 3-shot (%) | 5-way 5-shot (%) | 10-way 3-shot (%) | 10-way 5-shot (%) |
|-------------|------------------|------------------|-------------------|-------------------|
| GCN         | 75.41 $\pm$ 1.21 | 78.90 $\pm$ 0.61 | 60.40 $\pm$ 1.26  | 64.54 $\pm$ 0.69  |
| GraphSAGE   | 74.02 $\pm$ 1.43 | 77.58 $\pm$ 0.73 | 59.60 $\pm$ 1.40  | 64.45 $\pm$ 0.95  |
| GAT (Ours)  | 80.18 $\pm$ 0.27 | 84.76 $\pm$ 0.89 | 67.81 $\pm$ 0.97  | 70.88 $\pm$ 0.96  |

GAT outperforms the other encoders, as expected. The performance gap can be attributed to the lack of adaptive selection when introducing information (including noise) from the prompt graph into the target node embedding.

### A.6 Comparing More Classes Settings in Few-Shot Experiments

We chose the 3-way, 5-way, and 10-way settings across all datasets for two reasons. First, most works on few-shot node classification [7, 16, 34, 43] retain these settings, so we chose to follow them. Second, the limited number of classes in the datasets, such as the Amazon dataset with only slightly more than 10 classes, restricts us to these settings. Some meta-learning methods require splitting the classes

into training, validation, and test sets, which limits us to the 3-way setting. We further conducted new experiments on the ogbn-arxiv dataset under the 15-way setting for nine models (including variants). The results are as follows:

**Table 9: Comparison under the 15-way setting on the ogbn-arxiv dataset.**

| Model        | 15-way 3-shot (%)                   | 15-way 5-shot (%)                   |
|--------------|-------------------------------------|-------------------------------------|
| Raw Features | 32.24 $\pm$ 0.26                    | 37.37 $\pm$ 0.82                    |
| GPN          | 38.38 $\pm$ 0.39                    | 44.46 $\pm$ 0.68                    |
| G-Meta       | 39.36 $\pm$ 0.28                    | 44.85 $\pm$ 0.48                    |
| TENT         | 37.61 $\pm$ 0.32                    | 40.20 $\pm$ 0.37                    |
| GIANT        | 52.77 $\pm$ 0.50                    | 57.74 $\pm$ 0.86                    |
| G2P2         | 48.48 $\pm$ 0.97                    | 52.24 $\pm$ 0.85                    |
| P2TAG (LM)   | 57.31 $\pm$ 0.57                    | 62.71 $\pm$ 0.66                    |
| P2TAG (GNN)  | 58.49 $\pm$ 0.53                    | 63.60 $\pm$ 0.68                    |
| P2TAG        | 59.42 $\pm$ 0.42 (6.65 $\uparrow$ ) | 64.05 $\pm$ 0.56 (6.31 $\uparrow$ ) |

We can find that increasing the number of classes indeed makes the few-shot node classification task more challenging for all models. For example, the raw feature method achieves 59.31% in the 5-way setting, but only 44.59% in the 10-way setting. Nevertheless, P2TAG consistently outperforms the baselines, improving over the next best method by 6.65% and 6.31%, respectively.