

Node Classification on Graphs with Few-Shot Novel Labels via Meta Transformed Network Embedding

Lin Lan¹, Pinghui Wang¹, Xuefeng Du¹, Kaikai Song², Jing Tao¹, Xiaohong Guan¹

¹ Xi'an Jiaotong University, China ² Huawei Noah's Ark Lab

llan@sei.xjtu.edu.cn, {phwang, jtao, xhguan}@mail.xjtu.edu.cn

xuefengdu1@gmail.com, caesar.song@huawei.com

Abstract

We study the problem of node classification on graphs with few-shot novel labels, which has two distinctive properties: (1) There are novel labels to emerge in the graph; (2) The novel labels have only a few representative nodes for training a classifier. The study of this problem is instructive and corresponds to many applications such as recommendations for newly formed groups with only a few users in online social networks. To cope with this problem, we propose a novel Meta Transformed Network Embedding framework (MetaTNE), which consists of three modules: (1) A *structural module* provides each node a latent representation according to the graph structure. (2) A *meta-learning module* captures the relationships between the graph structure and the node labels as prior knowledge in a meta-learning manner. Additionally, we introduce an *embedding transformation function* that remedies the deficiency of the straightforward use of meta-learning. Inherently, the meta-learned prior knowledge can be used to facilitate the learning of few-shot novel labels. (3) An *optimization module* employs a simple yet effective scheduling strategy to train the above two modules with a balance between graph structure learning and meta-learning. Experiments on four real-world datasets show that MetaTNE brings a huge improvement over the state-of-the-art methods.

1 Introduction

Graphs are ubiquitously used to represent data in a wide range of fields, including social network analysis, bioinformatics, recommender systems, and computer network security. Accordingly, graph analysis tasks, such as node classification, link prediction, and community detection, have a significant impact on our lives in reality. In this paper, we focus on the task of node classification. Particularly, we consider the classification of *few-shot novel labels*, which means there are some novel labels to emerge in the graph of interest and the novel labels usually have only a few representative nodes including the positive and the negative (i.e., holding and not holding the novel labels, respectively)¹. The study of *Node Classification on graphs with Few-shot Novel Labels* (NCFNL) is instructive for many practical applications. Let us consider the following scenarios.

Motivating Examples. (1) Some organizations in online social networks, such as Facebook, Twitter, and Flickr, may distribute advertisements about whether users are interested in their new features or are willing to join their new social media groups. Through NCFNL, these organizations can predict other users' preferences based on positive and negative responses of a few users and provide better services or recommendations without too much bother for users. (2) For biological protein-protein networks, some researchers may discover a new biological function of certain proteins. Given a few proteins with and without a specific function, the study of NCFNL could predict whether other proteins have the function, which helps recommend new directions for wet laboratory experimentation.

¹ Hereafter, we refer to the *available positive and negative nodes* of a label as the *support nodes* of that label.

Some straightforward ways could be derived from existing unsupervised or semi-supervised network embedding methods while suffer from low performance, and please refer to § 2 for detailed discussions. To tackle this problem, we argue that different labels in a graph share some intrinsic evolution patterns (e.g., the way a label propagates along the graph structure according to the proximities between nodes). Assuming that there are a set of labels that have sufficient support nodes (e.g., interest groups that have existed and evolved for a long time in online social networks and protein functions that biologists are already familiar with), we desire to extract the common patterns from the graph structure and these labels and then utilize the found patterns to help recognize few-shot novel labels. However, the relationships between the graph structure and node labels are complex and there could be various propagation patterns between nodes. It remains challenging to design a model to capture all the patterns, and how to apply them to novel labels still needs to be further studied.

Overview of Our Approach. Inspired by recent advances in few-shot learning through meta-learning [22, 6], we cast the problem of NCFNL as a meta-learning problem and propose a novel Meta Transformed Network Embedding framework, namely MetaTNE, which allows us to exploit the common patterns. As shown in Fig. 1, our proposed framework consists of three modules: *the structural module*, *the meta-learning module*, and *the optimization module*. Given a graph and a set of labels (called known labels) with sufficient support nodes, the structural module first learns a latent representation for each node according to the graph structure. Then, considering that we ultimately expect to recognize few-shot novel labels, we propose the meta-learning module to simulate the few-shot scenario during the training phase instead of directly performing optimization over all known labels. Moreover, most existing meta-learning works [6, 25] focus on image- and text-related tasks, while the graph structure is more irregular in nature. To adequately exploit the complex and multifaceted relationships between nodes, we further design an *embedding transformation function* to map the structure-only (or task-agnostic) node representations to the task-specific ones for different few-shot classification tasks. To some extent, the meta-learning module implicitly encodes the shared propagation patterns of different labels through learning a variety of tasks. Finally, the optimization module is proposed to train the preceding two modules with a simple yet effective scheduling strategy in order to ensure the training stability and the effectiveness. One advantage of MetaTNE is that, after training, it is natural to directly apply the learned meta-learning module to few-shot novel labels.

Our main contributions are summarized as follows:

- To the best of our knowledge, this is the first work that only uses the graph structure and some known labels to study the problem of NCFNL. Compared with previous graph convolution based works [39, 34] that rely on high-quality node content for feature propagation and aggregation, our work is more challenging and at the same time more applicable to content-less scenarios.
- We propose an effective framework to solve NCFNL in a meta-learning manner. Our framework is able to generalize to classifying emerging novel labels with only a few support nodes. In particular, we design a transformation function that captures the multifaceted relationships between nodes to facilitate applying meta-learning to the graph data.
- We conduct extensive experiments on four publicly available real-world datasets, and empirical results show that MetaTNE achieves up to 150.93% and 47.58% performance improvement over the state-of-the-art methods in terms of Recall and F_1 , respectively.

2 Related Work

Unsupervised Network Embedding. This line of works focus on learning node embeddings that preserve various structural relations between nodes, including skip-gram based methods [20, 26, 7, 23], deep learning based methods [3, 31], and matrix factorization based methods [2, 21]. A straightforward way to adapt these methods for NCFNL is to simply train a new classifier (e.g., logistic regression) when novel labels emerge, while the learned node embeddings hold constant. However, this does not incorporate the guidance from node labels into the process of network embedding, which dramatically degrades the performance in the few-shot setting.

Semi-Supervised Network Embedding. These approaches typically formulate a unified objective function to jointly optimize the learning of node embeddings and the classification of nodes, such as combining the objective functions of DeepWalk and support vector machines [13, 28], as well as regarding labels as a kind of context and using node embeddings to simultaneously predict structural neighbors and node labels [5, 32]. Another line of works [11, 8, 30, 9] explore graph neural networks

to solve semi-supervised node classification as well as graph classification. Two recent works [15, 38] extend graph convolutional network (GCN) [11] to accommodate to the few-shot setting. However, the above methods are limited to a fixed set of labels and the adaptation of them to NCFNL requires to train the corresponding classification models or parameters from scratch when a novel label appears, which is not a well-designed solution to the few-shot novel labels and usually cannot reach satisfactory performance. Recently, Chauhan et al. [4] study few-shot graph classification with unseen novel labels based on graph neural networks. Zhang et al. [37] propose a few-shot knowledge graph completion method that essentially performs link prediction in a novel graph given a few training links. In comparison, we study node classification with respect to few-shot novel labels in the same graph and their methods are not applicable.

In addition, GCN based methods **heavily rely on** high-quality node content for feature propagation and aggregation, while in some networks (e.g., online social networks), some nodes (e.g., users) may not expose or expose noisy (low-quality) content, or even all node content is unavailable due to privacy concerns, which would limit the practical use of these methods. In contrast, our focus is to solve the problem of NCFNL by exploiting the relationships between the graph structure and the node labels, without involving node content.

Meta-Learning on Graphs. Zhou et al. [39] propose Meta-GNN that applies MAML [6] to GCN in a meta-learning way. More recently, Yao et al. [34] propose a method that combines GCN with metric-based meta-learning [25]. To some extent, all methods could handle novel labels emerging in a graph. However, they are built upon GCN and thus need high-quality node content for better performance, while in this paper we are interested in graphs without node content.

Few-Shot Learning on Images. Recently, few-shot learning has received considerable attention. Most works [22, 6, 25, 33, 17] focus on the problem of few-shot image classification in which there are no explicit relations between images. Some works also introduce task-specific designs for better generalization and learnability, such as task-specific null-space projection [35] and infinite mixture prototypes [1]. However, graph-structured data exhibits complex relations between nodes (i.e., the graph structure) which are the most fundamental and important information in a graph, making it difficult to directly apply these few-shot methods to graphs.

3 Problem Formulation

Throughout the paper, we use lowercase letters to denote scalars (e.g., ℓ), boldface lowercase letters to denote vectors (e.g., \mathbf{u}), and boldface uppercase letters to denote matrices (e.g., \mathbf{W}).

We denote a graph of interest by $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{Y})$, where $\mathcal{V} = \{v_1, v_2, \dots, v_{|\mathcal{V}|}\}$ is the set of nodes, $\mathcal{E} = \{e_{ij} = (v_i, v_j)\} \subseteq \mathcal{V} \times \mathcal{V}$ is the set of edges, and \mathcal{Y} is the set of labels associated with nodes in the graph. Here, we consider the multi-label setting where each node may have multiple labels. Let $\ell_{v_i, y} \in \{0, 1\}$ be the label indicator of the node v_i in terms of the label $y \in \mathcal{Y}$, where $\ell_{v_i, y} = 1$ suggests that the node v_i holds the label y and $\ell_{v_i, y} = 0$ otherwise. We use $\mathcal{D}_y^+ = \{v_i | \ell_{v_i, y} = 1\}$ to denote nodes that hold the label y , and $\mathcal{D}_y^- = \{v_i | \ell_{v_i, y} = 0\}$ to denote nodes that do not hold the label y . In this paper, we assume \mathcal{G} is undirected for ease of presentation.

Known Labels and Novel Labels. We divide the labels into two categories: the known labels $\mathcal{Y}_{\text{known}}$ and the novel labels $\mathcal{Y}_{\text{novel}}$. The former are given before we start any kind of learning process (e.g., semi-supervised network embedding), while the latter emerge after we have learned a model.

We assume that each known label is complete, namely $|\mathcal{D}_y^+| + |\mathcal{D}_y^-| = |\mathcal{V}|$ for $y \in \mathcal{Y}_{\text{known}}$. To some extent, the known labels refer to relatively stable labels (e.g., an interest group that has existed and evolved for a long time in online social networks). Although for some nodes, inevitably we are not sure whether they hold specific known labels or not, we simply assume that the corresponding label indicators equal 0 (i.e., not holding) like many other node classification works [20, 26]. In practice, a more principled way is to additionally consider the case of uncertain node-label pairs and define the label indicator as 1, 0, and -1 for the cases of holding the label, uncertain label, and not holding the label, respectively, which we leave as future work.

On the other hand, a novel label has only a few support nodes (e.g., 10 positive nodes and 10 negative nodes). By leveraging the known labels that have sufficiently many positive and negative nodes, we aim to explore the propagation patterns of labels along the graph structure and learn a model that generalizes well to classifying emerging novel labels with only a few support nodes.

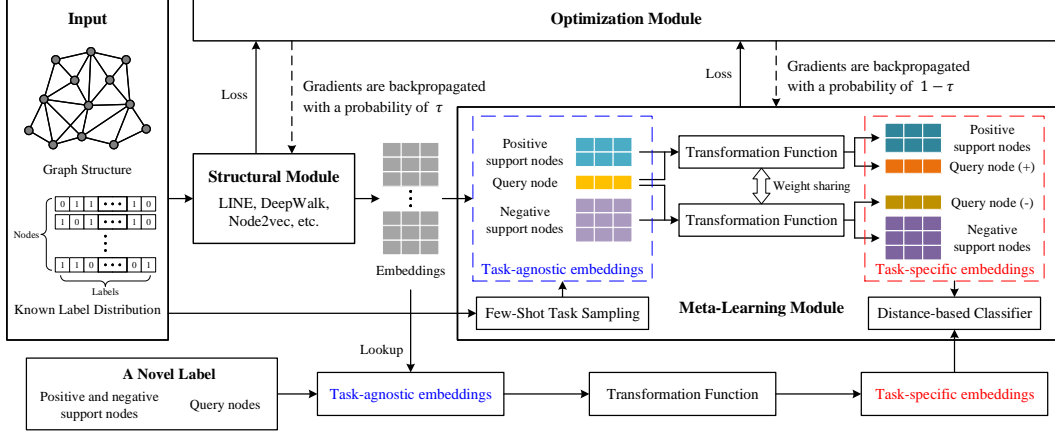


Figure 1: A schematic depiction of our MetaTNE. In the meta-learning module, we use 2 positive and 3 negative support nodes for simplicity of illustration. The threshold τ gradually decreases from 1 to 0 during training. The flow of applying MetaTNE to a novel label is shown at the bottom.

Our Problem. Given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{Y}_{\text{known}}, \mathcal{Y}_{\text{novel}})$, the problem of NCFNL aims to explore the relationships between the graph structure and the known labels $\mathcal{Y}_{\text{known}}$ and learn a generalizable model for classifying novel labels $\mathcal{Y}_{\text{novel}}$. Specifically, for each $y \in \mathcal{Y}_{\text{novel}}$, after observing only a few corresponding support nodes, the model should be able to generate or act as a good classifier to determine whether other nodes hold the label y or not.

4 Algorithm

In this section, we present our proposed MetaTNE in detail, which consists of three modules: the structural module, the meta-learning module, and the optimization module, as shown in Fig. 1. Given a graph and some known labels, the structural module learns an embedding for each node based on the graph structure. Then, the meta-learning module learns a transformation function that adapts the structure-only node embeddings for each few-shot node classification task sampled from the known labels and performs few-shot classification using a distance-based classifier. Finally, to optimize our model, we propose a learning schedule that optimizes the structural and meta-learning modules with probabilities that gradually decrease and increase from 1 to 0 and from 0 to 1, respectively.

4.1 Structural Module

The structural module aims to learn a representation or embedding in the latent space for each node while preserving the graph structure (i.e., the connections between nodes). Mathematically, for each node $v_i \in \mathcal{V}$, we maximize the log-probability of observing its neighbors by optimizing the following objective function: $\min \sum_{v_i \in \mathcal{V}} \sum_{v_j \in \mathcal{N}(v_i)} \log \mathbb{P}(v_j | v_i)$, where $\mathcal{N}(v_i)$ denotes the neighboring nodes of v_i . We optimize the above objective function following the skip-gram architecture [19]. Regarding the construction of the neighboring set $\mathcal{N}(\cdot)$, although there are many choices such as 1-hop neighbors based on the connectivity of nodes [26] and the random walk strategy [20, 7], in this paper we adopt the 1-hop neighbors for the sake of simplicity. By optimizing the above objective, we are able to obtain an embedding matrix $\mathbf{U} \in \mathbb{R}^{|\mathcal{V}| \times d}$, of which the i -th row \mathbf{u}_i represents the embedding vector of v_i .

4.2 Meta-Learning Module

As alluded before, we cast the problem of NCFNL into the meta-learning framework [22, 6] and simulate the few-shot setting with $\mathcal{Y}_{\text{known}}$ during training. In what follows, we first describe how to organize the graph structure and the known labels in the meta-learning scenario. Then, we give a metric-based meta-learning paradigm for solving NCFNL. In particular, we propose a transformation function that transforms the task-agnostic embeddings to the task-specific ones in order to better deal with the multi-label setting where each node may be associated with multiple labels.

4.2.1 Data Organization

Instead of directly optimizing over the entire set of known labels like traditional semi-supervised learning methods [32], we propose to construct a pool of few-shot node classification tasks according to the known labels $\mathcal{Y}_{\text{known}}$. Analogous to few-shot image classification tasks in the literature of meta-learning [22], a few-shot node classification task $\mathcal{T}_i = (\mathcal{S}_i, \mathcal{Q}_i, y_i)$ is composed of a support set \mathcal{S}_i , a query set \mathcal{Q}_i , and a label identifier y_i randomly sampled from $\mathcal{Y}_{\text{known}}$. The support set $\mathcal{S}_i = \mathcal{S}_i^+ \cup \mathcal{S}_i^-$ contains the set \mathcal{S}_i^+ of randomly sampled positive nodes and the set \mathcal{S}_i^- of randomly sampled negative nodes, where $\mathcal{S}_i^+ \subset \mathcal{D}_{y_i}^+$ and $\mathcal{S}_i^- \subset \mathcal{D}_{y_i}^-$. The query set $\mathcal{Q}_i = \mathcal{Q}_i^+ \cup \mathcal{Q}_i^-$ is defined in the same way but does not intersect with the support set, namely $\mathcal{Q}_i^+ \subset \mathcal{D}_{y_i}^+ \setminus \mathcal{S}_i^+$ and $\mathcal{Q}_i^- \subset \mathcal{D}_{y_i}^- \setminus \mathcal{S}_i^-$. The task is, given the support set of node-label pairs, finding a classifier $f_{\mathcal{T}_i}$ which is able to predict the probability $\hat{\ell}_{v_q, y_i} \in [0, 1]$ for each query node v_q with a low misclassification rate. We denote by $\mathcal{T}_i \sim p(\mathcal{T}|\mathcal{Y}_{\text{known}})$ sampling a few-shot node classification task from $\mathcal{Y}_{\text{known}}$.

4.2.2 Meta-Learning with Embedding Transformation for NCFNL

To facilitate learning to classify for a label with few associated nodes in a graph, we apply a meta-learning flavored learning scheme. Following the above definition of few-shot node classification tasks, for each task $\mathcal{T}_i = (\mathcal{S}_i, \mathcal{Q}_i, y_i) \sim p(\mathcal{T}|\mathcal{Y}_{\text{known}})$, we aim to construct a classifier $f_{\mathcal{T}_i}$ for the label y_i given the support set \mathcal{S}_i , which is able to classify the query nodes in the set \mathcal{Q}_i . Formally, for each $(v_q, \ell_{v_q, y_i}) \in \mathcal{Q}_i$, the classification loss is defined as follows:

$$\mathcal{L}(\hat{\ell}_{v_q, y_i}, \ell_{v_q, y_i}) = -\ell_{v_q, y_i} \log \hat{\ell}_{v_q, y_i} - (1 - \ell_{v_q, y_i}) \log(1 - \hat{\ell}_{v_q, y_i}), \quad (1)$$

where $\hat{\ell}_{v_q, y_i}$ denotes the predicted probability that v_q holds label y_i . Here, to calculate the probability, we adopt a distance-based classifier which is commonly used in the metric-based meta-learning literature [25]. Specifically, for each task \mathcal{T}_i , the classifier $f_{\mathcal{T}_i}$ is parametrized by two d -dimensional latent representations, $\mathbf{c}_+^{(i)}$ (called positive prototype) and $\mathbf{c}_-^{(i)}$ (called negative prototype), that correspond to the cases of holding and not holding label y_i , respectively. The predictions are made based on the distances between the node representations and these two prototypes. Mathematically, given the embedding vector \mathbf{u}_q of each query node v_q , we have the predicted probability as

$$\hat{\ell}_{v_q, y_i} = f_{\mathcal{T}_i}(v_q | \mathbf{c}_+^{(i)}, \mathbf{c}_-^{(i)}) = \frac{\exp(-\text{dist}(\mathbf{u}_q, \mathbf{c}_+^{(i)}))}{\sum_{m \in \{+, -\}} \exp(-\text{dist}(\mathbf{u}_q, \mathbf{c}_m^{(i)}))}, \quad (2)$$

where $\text{dist}(\cdot, \cdot) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow [0, +\infty)$ is the squared Euclidean distance function and the positive or negative prototype is usually calculated as the mean vector of node representations in the corresponding support set [25].

Why do we need Embedding Transformation? Equation (2) makes predictions under the condition that each node is represented by the same or *task-agnostic* embedding vector regardless of which label or task we are concerned about. Technically, this scheme makes sense for few-shot image classification in prior works [25] where each image is assigned to the same one and only one label. However, this is problematic in the multi-label scenario where each node could be assigned to multiple labels. Here is an illustrating example. In social networks, suppose we have two classification tasks \mathcal{T}_1 and \mathcal{T}_2 with respect to different labels, namely “Sports” from $\mathcal{Y}_{\text{known}}$ and “Music” from $\mathcal{Y}_{\text{novel}}$, and two users A and B are involved in these two tasks. Both users A and B could give positive feedback to “Sports”, while on the other hand, they could give positive and negative feedback to “Music” respectively. Intuitively, the task-agnostic scheme may provide similar embeddings after fitting well on the task \mathcal{T}_1 , which is not appropriate for the task \mathcal{T}_2 .

High-Level Module Design. To mitigate the above problem, we propose to learn a transformation function $Tr(\cdot)$ which transforms the task-agnostic embeddings to some task-specific ones for each task. First, we argue that different query nodes have different correlation patterns with the nodes in the support set. To fully explore how a query node correlates with the support nodes, we propose to tailor the embeddings of the support nodes for each query node. Second, to classify a query node, we are more interested in characterizing the distance relationship between the query node and either positive or negative support nodes rather than the relationship between the positive and negative support nodes. Thus, during the transformation, we propose to adapt the query node with the positive and the negative nodes in the support set separately.

Based on the above two principles, for each query node, we first construct two sets: one containing the task-agnostic embeddings of the query node and the positive support nodes, and the other containing the task-agnostic embeddings of the query node and the negative support nodes. Then, we separately feed the two sets into the transformation function. The meta-learning module in Fig. 1 illustrates this process. Formally, given a task $\mathcal{T}_i = (\mathcal{S}_i, \mathcal{Q}_i, y_i)$, for each query node $v_q \in \mathcal{V}_{\mathcal{Q}_i}$, we have

$$\{\tilde{\mathbf{u}}_{q,m}^{(i)}\} \cup \{\tilde{\mathbf{u}}_{k,q}^{(i)} | v_k \in \mathcal{V}_{\mathcal{S}_i^m}\} = Tr(\{\mathbf{u}_q\} \cup \{\mathbf{u}_k | v_k \in \mathcal{V}_{\mathcal{S}_i^m}\}), \quad m \in \{+, -\}, \quad (3)$$

where $\tilde{\mathbf{u}}_{q,m}^{(i)}$ denotes the adapted embedding of the query node v_q in relation to the positive or negative support nodes, and $\tilde{\mathbf{u}}_{k,q}^{(i)}$ denotes the adapted embedding of the support node v_k tailored for the query node v_q . As a result, each query node has two different adapted embeddings $\tilde{\mathbf{u}}_{q,+}^{(i)}$ and $\tilde{\mathbf{u}}_{q,-}^{(i)}$ that are further used for comparisons with the adapted embeddings of the positive and negative support nodes, respectively. A consequential benefit is that the transformation function is more flexible to capture the multifaceted relationships between nodes in the multi-label scenario. Imagine that even if the task-specific embeddings of the positive and negative support nodes or prototypes are distributed close, we are still able to make right predictions through altering $\tilde{\mathbf{u}}_{q,+}^{(i)}$ and $\tilde{\mathbf{u}}_{q,-}^{(i)}$. The ablation study in Section Experiments and the visualization in § C.5 confirm the superiority of this design.

Instantiation. As per the above discussions, we propose to implement $Tr(\cdot)$ using the self-attention architecture with the scaled dot-product attention mechanism [29], which separately takes as input the two sets $\{\mathbf{u}_q\} \cup \{\mathbf{u}_k | v_k \in \mathcal{V}_{\mathcal{S}_i^m}\}$ where $m \in \{+, -\}$. Mathematically, we have

$$\tilde{\mathbf{X}}_m^{(i)} = \text{SelfAttention}(\mathbf{X}_m^{(i)}) = \text{softmax}\left(\frac{\mathbf{X}_m^{(i)} \mathbf{W}_Q (\mathbf{X}_m^{(i)} \mathbf{W}_K)^\top}{\sqrt{d'}}\right) \mathbf{X}_m^{(i)} \mathbf{W}_V, \quad (4)$$

where $\mathbf{X}_m^{(i)}$ is the concatenation of $[\mathbf{u}_q]$ and $[\mathbf{u}_k; \forall v_k \in \mathcal{V}_{\mathcal{S}_i^m}]$, $\tilde{\mathbf{X}}_m^{(i)}$ is the concatenation of the transformed representations $[\tilde{\mathbf{u}}_{q,m}^{(i)}]$ and $[\tilde{\mathbf{u}}_{k,q}^{(i)}; \forall v_k \in \mathcal{V}_{\mathcal{S}_i^m}]$, $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$ are trainable projection matrices, and d' is the dimension after projection. We refer readers to § A.1 for more details on the instantiation of the transformation function.

With the transformed embeddings, we further calculate the positive and negative prototypes tailored for v_q as well as the predicted probability as follows:

$$\tilde{\mathbf{c}}_{m,q}^{(i)} = \frac{1}{|\mathcal{S}_i^m|} \sum_{v_k \in \mathcal{V}_{\mathcal{S}_i^m}} \tilde{\mathbf{u}}_{k,q}^{(i)}, \quad m \in \{+, -\}, \quad \text{and} \quad \hat{\ell}_{v_q, y_i} = \frac{\exp(-\text{dist}(\tilde{\mathbf{u}}_{q,+}^{(i)}, \tilde{\mathbf{c}}_{+,q}^{(i)}))}{\sum_{m \in \{+, -\}} \exp(-\text{dist}(\tilde{\mathbf{u}}_{q,m}^{(i)}, \tilde{\mathbf{c}}_{m,q}^{(i)}))}. \quad (5)$$

The final meta-learning objective is formulated as:

$$\min_{\mathbf{U}, \Theta} \sum_{\mathcal{T}_i} \sum_{(v_q, \ell_{v_q, y_i}) \in \mathcal{Q}_i} \mathcal{L}(\hat{\ell}_{v_q, y_i}, \ell_{v_q, y_i}) + \lambda \sum \|\Theta\|_2^2, \quad (6)$$

where $\mathcal{T}_i \sim p(\mathcal{T} | \mathcal{Y}_{\text{known}})$, $\hat{\ell}_{v_q, y_i}$ is calculated through Eqn. (5), Θ refers to the set of parameter matrices (e.g., $\mathbf{W}_Q, \mathbf{W}_K$, and \mathbf{W}_V) contained in $Tr(\cdot)$, and $\lambda > 0$ is a balancing factor.

4.3 Optimization and Using the Learned Model for Few-Shot Novel Labels

For optimization, one typical way is to minimize the (weighted) sum of the structural loss and the meta loss. However, the structure information of the graph is still not properly embedded at the beginning of the training stage, and the node representations are somewhat random which make no sense for the few-shot classification tasks. Therefore, a training procedure that focuses on optimizing the structural module at the beginning and then gradually pays more attention to optimizing the meta-learning module is preferably required. To satisfy this requirement, we take inspiration from learning rate annealing [12] and introduce a probability threshold τ , and in each training step the structural and meta modules are optimized with probabilities of τ and $1 - \tau$, respectively. The probability threshold τ is gradually decayed from 1 to 0 in a staircase manner, namely $\tau = 1 / (1 + \gamma \lfloor \frac{\text{step}}{N_{\text{decay}}} \rfloor)$ where γ is the decay rate, step is the current step number, and N_{decay} indicates how often the threshold is decayed. The complete optimization procedure is outlined in Algorithm 1 in the appendices. In addition, the time complexity is analyzed in § A.3.

Our ultimate goal is to, after observing a few support nodes associated with a novel label $y \in \mathcal{Y}_{\text{novel}}$, predict whether other (or some query) nodes have the label y or not. In effect, this can be regarded as a few-shot node classification task $\mathcal{T} = (\mathcal{S}, \mathcal{Q}, y)$. After optimization, we have obtained the task-agnostic node representations \mathbf{U} , and the transformation function $Tr(\cdot)$ parameterized by Θ . Thus, to classify a query node $v_q \in \mathcal{Q}$, we simply look up the representations of the query node and the support nodes from \mathbf{U} , adapt their representations using the transformation function as formulated in Eqn. (3) and (4), and compute the predicted probability according to Eqn. (5). The detailed procedure is presented in Algorithm 2 in the appendices.

5 Experiments

Four publicly available real-world benchmark datasets are used to validate the effectiveness of our method. The statistics of these datasets are summarized in Table 1. For each dataset, we split the labels into training, validation, and test labels according to a ratio of 6:2:2. In the training stage, we regard the training labels as the known labels and sample few-shot node classification tasks from them. For validation and test purposes, we regard the validation and test labels as the novel labels and sample 1,000 tasks from them, respectively.

Table 1: Statistics of the datasets.

Dataset	#Nodes	#Edges	#Labels
BlogCatalog	10,312	333,983	39
Flickr	80,513	5,899,882	195
PPI	3,890	76,584	50
Mashup	16,143	300,181	28

We use the average classification performance on the test tasks for comparisons of different methods. For ease of presentation, we use $K_{\mathcal{S},+}$, $K_{\mathcal{S},-}$, $K_{\mathcal{Q},+}$, and $K_{\mathcal{Q},-}$ to indicate the respective numbers of positive support, negative support, positive query, and negative query nodes in a task. We compare MetaTNE with Label Propagation [40], unsupervised network embedding methods (LINE [26] and Node2vec [7]), semi-supervised network embedding methods (Planetoid [32] and GCN [11]), and Meta-GNN [39]. For detailed experimental settings including dataset and baseline descriptions, baseline evaluation procedure, and parameter settings, please refer to § B.

Overall Comparisons. Following the standard evaluation protocol of meta-learning [6], we first compare different methods with $K_{\mathcal{S},+} = K_{\mathcal{Q},+}$ and $K_{\mathcal{S},-} = K_{\mathcal{Q},-}$ (hereafter using $K_{*,+}$ and $K_{*,-}$ for simplicity), and these numbers are the same for both training and test tasks. Considering that negative samples are usually easier than positive samples to acquire we report the overall performance with $K_{*,+}$ set to 10 and $K_{*,-}$ set to 20 and 40, respectively. The comparison results on the four datasets are presented in Table 2. Since in our application scenarios we prefer to discover proteins with new functions in biological networks and find users who are interested in the latest advertisements on online social networks rather than predict negative samples accurately, we report Recall in addition to AUC and F_1 . To eliminate randomness, all of the results here and in the following quantitative experiments are averaged over 50 different trials.

From Table 2, we observe that MetaTNE consistently and significantly outperforms all other methods in terms of the three metrics across all the four datasets except the AUC scores on Flickr dataset. By jointly analyzing the F_1 and Recall scores, MetaTNE predicts positive nodes from imbalanced data more effectively than the baselines, with little loss of precision. In particular, MetaTNE achieves 44.22% and 150.93% gains over the strongest baseline (i.e., Planetoid) with respect to Recall on BlogCatalog dataset when $K_{*,-}$ equals 20 and 40, respectively.

Compared with the unsupervised methods, Planetoid reaches better performance owing to the use of training labels. On the other hand, GCN also uses training labels as supervision, while does not show satisfactory performance and even worse performance than Node2vec, which is due to that the graph convolution relies heavily on node attributes for feature propagation and aggregation as mentioned before and the lack of node attributes limits its representativeness and thus classification capacity.

Besides, Meta-GNN underperforms the unsupervised methods and GCN in some cases, which seems to contradict the published results in the original paper. The reasons are twofold: (1) Meta-GNN is built upon GCN and the predictive ability is also limited due to the lack of node attributes, while the original paper focuses on attributed graphs; (2) Meta-GNN simply applies MAML to GCN and is originally used for the multi-class setting (e.g., each document has the same and only one label in Cora [24]). However, we consider the multi-label setting and the same pair of nodes may have opposite relations in different tasks, which will introduce noisy and contradictory signals in the optimization process of MAML and further degrade the performance in some cases.

Table 2: Results on few-shot node classification tasks with novel labels. OOM means out of memory (16 GB GPU memory). The standard deviation is provided in § C.1.

(a) $K_{*,+} = 10$ and $K_{*,-} = 20$.

Method	BlogCatalog			Flickr			PPI			Mashup		
	AUC	F ₁	Recall	AUC	F ₁	Recall	AUC	F ₁	Recall	AUC	F ₁	Recall
LP	0.6422	0.1798	0.2630	0.8196	0.4321	0.4989	0.6285	0.2147	0.2769	0.6488	0.3103	0.4535
LINE	0.6690	0.2334	0.1595	0.8593	0.6194	0.5418	0.6372	0.2147	0.1456	0.6926	0.2970	0.2142
Node2vec	0.6697	0.3750	0.2940	0.8504	0.6664	0.6147	0.6273	0.3545	0.2860	0.6575	0.3835	0.3147
Planetoid	0.6850	0.4657	0.4301	0.8601	0.6638	0.6331	0.6791	0.4672	0.4411	0.7056	0.4825	0.4218
GCN	0.6102	0.2730	0.2194	OOM	OOM	OOM	0.6544	0.3379	0.2721	0.6895	0.3052	0.2390
Meta-GNN	0.4805	0.2375	0.2141	OOM	OOM	OOM	0.5466	0.3289	0.3081	0.7078	0.4576	0.4176
MetaTNE	0.6986	0.5380	0.6203	0.8462	0.7118	0.7700	0.6865	0.5188	0.5621	0.7645	0.5764	0.5566
%Improv.	1.99	15.53	44.22	-1.62	6.81	21.62	1.09	11.04	27.43	8.01	19.46	22.73

(b) $K_{*,+} = 10$ and $K_{*,-} = 40$.

Method	BlogCatalog			Flickr			PPI			Mashup		
	AUC	F ₁	Recall	AUC	F ₁	Recall	AUC	F ₁	Recall	AUC	F ₁	Recall
LP	0.6421	0.0554	0.0727	0.8253	0.3055	0.3040	0.6298	0.0773	0.0748	0.6534	0.1156	0.1284
LINE	0.6793	0.0529	0.0328	0.8644	0.4154	0.3485	0.6423	0.0496	0.0300	0.7009	0.0956	0.0617
Node2vec	0.6792	0.1982	0.1340	0.8558	0.5295	0.4602	0.6309	0.1894	0.1306	0.6643	0.2070	0.1447
Planetoid	0.6981	0.2980	0.2319	0.8728	0.5040	0.4461	0.6879	0.3100	0.2523	0.7095	0.3279	0.2551
GCN	0.6198	0.1011	0.0704	OOM	OOM	OOM	0.6608	0.1403	0.0957	0.6979	0.0813	0.0531
Meta-GNN	0.4811	0.1042	0.0859	OOM	OOM	OOM	0.5399	0.2085	0.1867	0.7050	0.3279	0.2768
MetaTNE	0.7139	0.4398	0.5819	0.8505	0.6220	0.7460	0.7039	0.4298	0.5327	0.7684	0.4814	0.4816
%Improv.	2.26	47.58	150.93	-2.55	17.47	62.10	2.33	38.65	111.14	8.30	46.81	73.99

Ablation Study. In what follows, to gain deeper insight into the contributions of different components involved in our approach, we conduct ablation studies by considering the following variants: (1) a variant without the transformation function; (2) a variant that produces task-specific embeddings by simply feeding all support and query node representations into the self-attention network instead of according to

Eqn. (3); (3) a variant that optimizes the total loss of the two modules with the meta-learning loss scaled by a balancing factor searched over $\{10^{-2}, 10^{-1}, \dots, 10^2\}$. We refer to these variants as V1, V2, and V3. The results are summarized in Table 3. We see that MetaTNE consistently outperforms its three ablated variants. Especially, we observe that V1 performs the worst in most cases, which confirms the necessity to introduce the transformation function. The comparison with V2 demonstrates the effectiveness of our special design in Eqn (3). Moreover, the results of V3 indicate that our proposed scheduling strategy can boost the performance of MetaTNE with a better balance between the two modules during optimization.

Additional Experiments. In § C, we present more analytical experiments on the numbers of support and query nodes, and illustrate the effect of the proposed transformation function through a visualization experiment.

6 Conclusion and Future Work

This paper studies the problem of node classification on graphs with few-shot novel labels. To address this problem, we propose a new semi-supervised framework MetaTNE that integrates network embedding and meta-learning. Benefiting from utilizing known labels in a meta-learning manner, MetaTNE is able to automatically capture the relationships between the graph structure and the node labels as prior knowledge and make use of the prior knowledge to help recognize novel labels with only a few support nodes. Extensive experiments on four real-world datasets demonstrate the superiority of our proposed method. In the future, to improve the interpretability, we plan to extend our approach to quantify the relationships between different labels (e.g., the weight that one label contributes to another) during meta-learning. Another interesting idea is to explicitly incorporate the graph structure information into the meta-learning module, such as developing a more principled way to construct few-shot tasks according to the graph structure instead of random sampling.

Table 3: Results of ablation study in terms of F₁.

Method	$K_{*,+} = 10, K_{*,-} = 20$		$K_{*,+} = 10, K_{*,-} = 40$	
	BlogCatalog	PPI	BlogCatalog	PPI
MetaTNE	0.5380	0.5188	0.4398	0.4298
V1	0.5028	0.4851	0.3998	0.3721
V2	0.5020	0.5011	0.4141	0.4078
V3	0.5205	0.4980	0.4039	0.4074

Broader Impact

In general, this work has potential positive impact on graph-related fields that need to deal with the classification problem with respect to few-shot novel labels. For instance, our work is beneficial for social networking service providers such as Facebook and Twitter. These providers can obtain quick and effective feedback on newly developed features through distributing surveys among a small group of users on social networks. In addition, our work can also help biologists, after discovering a new function of certain existing proteins, quickly understand whether other proteins in a protein-protein interaction network have the new function, which improves the efficiency of wet laboratory experimentation. Moreover, many recommender systems model users and items as a graph and enhance the recommendation performance with the aid of network embedding. To some extent, our work is potentially useful to alleviate the cold-start problem as well.

At the same time, our model could be biased towards the few-shot setting after training and not provide superior performance on those labels with many support nodes. In practice, if the original few-shot label gradually has enough support nodes (e.g., biologists identify more proteins with and without the new function through laboratory experiments), we recommend using general unsupervised or semi-supervised methods (e.g., Node2vec [7] or Planetoid [32]) to recognize the label.

References

- [1] Kelsey R. Allen, Evan Shelhamer, Hanul Shin, and Joshua B. Tenenbaum. Infinite mixture prototypes for few-shot learning. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *ICML*, 2019.
- [2] Shaosheng Cao, Wei Lu, and Qiongkai Xu. Grarep: Learning graph representations with global structural information. In *CIKM*, 2015.
- [3] Shaosheng Cao, Wei Lu, and Qiongkai Xu. Deep neural networks for learning graph representations. In *Thirtieth AAAI conference on artificial intelligence*, 2016.
- [4] Jatin Chauhan, Deepak Nathani, and Manohar Kaul. Few-shot learning on graphs via super-classes based on graph spectral measures. In *ICLR*, 2020.
- [5] Jifan Chen, Qi Zhang, and Xuanjing Huang. Incorporate group information to enhance network embedding. In *CIKM*, 2016.
- [6] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017.
- [7] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *SIGKDD*, 2016.
- [8] Will Hamilton, Zhitaoying, and Jure Leskovec. Inductive representation learning on large graphs. In *NIPS*, 2017.
- [9] Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure Leskovec. Strategies for pre-training graph neural networks. In *ICLR*, 2020.
- [10] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [11] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017.
- [12] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [13] Juzheng Li, Jun Zhu, and Bo Zhang. Discriminative deep random walk for network classification. In *ACL*, 2016.
- [14] Liam Li, Kevin Jamieson, Afshin Rostamizadeh, Ekaterina Gonina, Moritz Hardt, Benjamin Recht, and Ameet Talwalkar. Massively parallel hyperparameter tuning. *arXiv preprint arXiv:1810.05934*, 2018.

- [15] Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional networks for semi-supervised learning. In *AAAI*, 2018.
- [16] Richard Liaw, Eric Liang, Robert Nishihara, Philipp Moritz, Joseph E Gonzalez, and Ion Stoica. Tune: A research platform for distributed model selection and training. *arXiv preprint arXiv:1807.05118*, 2018.
- [17] Yanbin Liu, Juho Lee, Minseop Park, Saehoon Kim, Eunho Yang, Sung Ju Hwang, and Yi Yang. Learning to propagate labels: Transductive propagation network for few-shot learning. In *ICLR*, 2019.
- [18] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *JMLR*, 2008.
- [19] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, 2013.
- [20] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social representations. In *SIGKDD*, 2014.
- [21] Jiezhong Qiu, Yuxiao Dong, Hao Ma, Jian Li, Kuansan Wang, and Jie Tang. Network embedding as matrix factorization: Unifying deepwalk, line, pte, and node2vec. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, pages 459–467, 2018.
- [22] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *ICLR*, 2017.
- [23] Leonardo Filipe Rodrigues Ribeiro, Pedro H. P. Saverese, and Daniel R. Figueiredo. *struc2vec*: Learning node representations from structural identity. In *SIGKDD*, pages 385–394, 2017.
- [24] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.
- [25] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *NIPS*, 2017.
- [26] Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. Line: Large-scale information network embedding. In *WWW*, 2015.
- [27] Lei Tang and Huan Liu. Relational learning via latent social dimensions. In *KDD*, 2009.
- [28] Cunchao Tu, Weicheng Zhang, Zhiyuan Liu, Maosong Sun, et al. Max-margin deepwalk: Discriminative learning of network representation. In *IJCAI*, 2016.
- [29] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017.
- [30] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. In *ICLR*, 2018.
- [31] Daixin Wang, Peng Cui, and Wenwu Zhu. Structural deep network embedding. In *SIGKDD*, 2016.
- [32] Zhilin Yang, William W Cohen, and Ruslan Salakhutdinov. Revisiting semi-supervised learning with graph embeddings. In *ICML*, 2016.
- [33] Huaxiu Yao, Xian Wu, Zhiqiang Tao, Yaliang Li, Bolin Ding, Ruirui Li, and Zhenhui Li. Automated relational meta-learning. In *ICLR*, 2020.
- [34] Huaxiu Yao, Chuxu Zhang, Ying Wei, Meng Jiang, Suhang Wang, Junzhou Huang, Nitesh V Chawla, and Zhenhui Li. Graph few-shot learning via knowledge transfer. In *AAAI*, 2020.
- [35] Sung Whan Yoon, Jun Seo, and Jaekyun Moon. Tapnet: Neural network augmented with task-adaptive projection for few-shot learning. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *ICML*, 2019.

- [36] Xiang Yue, Zhen Wang, Jingong Huang, Srinivasan Parthasarathy, Soheil Moosavinasab, Yungui Huang, Simon M Lin, Wen Zhang, Ping Zhang, and Huan Sun. Graph embedding on biomedical networks: Methods, applications, and evaluations. *arXiv preprint arXiv:1906.05017*, 2019.
- [37] Chuxu Zhang, Huaxiu Yao, Chao Huang, Meng Jiang, Zhenhui Li, and Nitesh V. Chawla. Few-shot knowledge graph completion. *CoRR*, abs/1911.11298, 2019.
- [38] Shengzhong Zhang, Ziang Zhou, Zengfeng Huang, and Zhongyu Wei. Few-shot classification on graphs with structural regularized GCNs, 2019. URL <https://openreview.net/forum?id=r1znKiAcY7>.
- [39] Fan Zhou, Chengtai Cao, Kunpeng Zhang, Goce Trajcevski, Ting Zhong, and Ji Geng. Meta-gnn: On few-shot node classification in graph meta-learning. In *CIKM*, 2019.
- [40] Xiaojin Zhu, Zoubin Ghahramani, and John D Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *ICML*, 2003.

A Additional Algorithm Details

A.1 Details of the Transformation Function

For the transformation function, we stack multiple computation blocks as shown in Fig. 2. The stacking mechanism helps the function capture comprehensive relationships between nodes such that the performance is boosted. In each computation block, there are mainly two modules. The first is a self-attention module used to capture the relationships between input nodes, and the second is a node-wise fully-connected feed-forward network used to introduce nonlinearity. In addition, following [29], we employ a residual connection around each of the self-attention module and the feed-forward network and then perform layer normalization, in order to make the optimization faster and more stable.

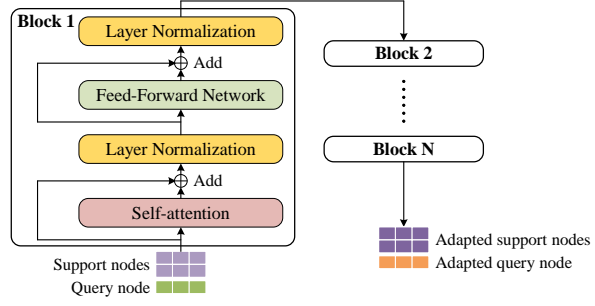


Figure 2: Illustration of the transformation function. The support nodes are either positive or negative.

The detailed architecture of the self-attention module is illustrated in Fig. 3. Following [29], we extend the self-attention with multiple parallel attention heads using multiple sets of trainable matrices (i.e., $\mathbf{W}_Q^h, \mathbf{W}_K^h, \mathbf{W}_V^h \in \mathbb{R}^{d'/H \times d}$ where $h = 1, \dots, H$). In each attention head (i.e., each scaled dot-product attention), for any two nodes $v_i, v_j \in \{v_q\} \cup \mathcal{V}_{S_i^m}$ (v_i and v_j could be the same and $m \in \{+, -\}$) within task \mathcal{T}_i , we first calculate the attention ω_{ij}^h that v_i pays to v_j as follows:

$$\omega_{ij}^h = \frac{\exp((\mathbf{W}_Q^h \mathbf{u}_i) \cdot (\mathbf{W}_K^h \mathbf{u}_j) / \sqrt{d'/H})}{\sum_{v_k \in \{v_q\} \cup \mathcal{V}_{S_i^m}} \exp((\mathbf{W}_Q^h \mathbf{u}_i) \cdot (\mathbf{W}_K^h \mathbf{u}_k) / \sqrt{d'/H})}, \quad (7)$$

where “ \cdot ” denotes the dot product operator. Then, we compute the output vector of the query node v_q as

$$\tilde{\mathbf{u}}_{q,m}^{i,h} = \omega_{qq}^h \mathbf{W}_V^h \mathbf{u}_q + \sum_{v_k \in \mathcal{V}_{S_i^m}} \omega_{qk}^h \mathbf{W}_V^h \mathbf{u}_k, \quad (8)$$

and compute the output vector of each support node $v_k \in \mathcal{V}_{S_i^m}$ tailored for the query node v_q as

$$\tilde{\mathbf{u}}_{k,q}^{i,h} = \omega_{kk}^h \mathbf{W}_V^h \mathbf{u}_k + \sum_{v_j \in (\mathcal{V}_{S_i^m} \setminus \{v_k\}) \cup \{v_q\}} \omega_{kj}^h \mathbf{W}_V^h \mathbf{u}_j. \quad (9)$$

Finally, we concatenate the output vectors of all attention heads and use a trainable matrix $\mathbf{W}_O \in \mathbb{R}^{d \times d'}$ to project the concatenated vectors into the original space with the input dimension:

$$\tilde{\mathbf{u}}_{q,m}^{(i)} = \mathbf{W}_O(\tilde{\mathbf{u}}_{q,m}^{i,1} \oplus \dots \oplus \tilde{\mathbf{u}}_{q,m}^{i,H}), \quad \text{and} \quad \tilde{\mathbf{u}}_{k,q}^{(i)} = \mathbf{W}_O(\tilde{\mathbf{u}}_{k,q}^{i,1} \oplus \dots \oplus \tilde{\mathbf{u}}_{k,q}^{i,H}), \quad \forall v_k \in \mathcal{V}_{S_i^m}. \quad (10)$$

The multiple parallel attention heads allow the function to jointly attend to information from different input nodes for each input node, and thus help the function better exploit the relationships between input nodes.

A.2 Pseudo Codes

The optimization procedure is outlined in Algorithm 1. The procedure of using the learned model for few-shot novel labels is presented in Algorithm 2.

A.3 Time Complexity Analysis

For the structural module, we optimize the objective function in a way similar to [26] and the time complexity is $O(kd|\mathcal{E}|)$ where k is the number of negative nodes at each iteration, d is the dimension

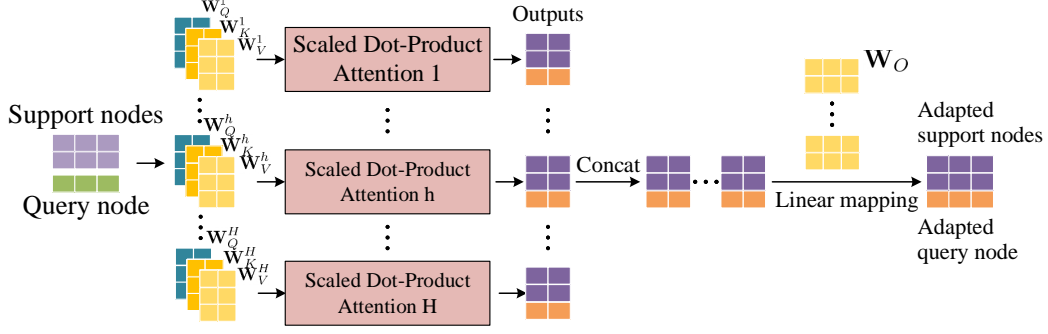


Figure 3: Illustration of the self-attention module. The support nodes are either positive or negative.

of node embeddings, and $|\mathcal{E}|$ is the number of edges. For the meta-learning module, the time cost mainly comes from the embedding transformation through the self-attention architecture [29]. Specifically, let m be the number of query nodes and n be the number of positive or negative support nodes. Calculating the *query*, *key*, and *value* vectors takes $O(mndd')$, where d' is the dimension of the *query*, *key*, and *value* vectors. Calculating the attention weights and the weighted sum of *value* vectors takes $O(mn^2d')$. Calculating the final output vectors takes $O(mndd')$. Overall, the time complexity of MetaTNE is $O(kd|\mathcal{E}| + mndd' + mn^2d')$. Note that we can take advantage of GPU acceleration for optimization in practice.

B Details of the Experimental Settings

B.1 Datasets

Four datasets are used in our experiments.

BlogCatalog [27]: This dataset is the friendship network crawled from the BlogCatalog website. The friendships and group memberships are encoded in the edges and labels, respectively.²

Flickr [27]: This dataset is the friendship network among the bloggers crawled from the Flickr website. The friendships and group memberships are encoded in the edges and the labels, respectively.³

PPI [7]: This dataset is a protein-protein interaction network for Homo Sapiens. Different labels represent different function annotations of proteins.⁴

Mashup [36]: This dataset is a protein-protein interaction network for human. Different labels represent different function annotations of proteins.⁵

B.2 Baselines

The following baselines are considered:

Label Propagation (LP) [40]: This method is a semi-supervised learning algorithm that estimates labels by propagating label information through a graph. It assigns a node the label which most of its neighborhoods have and propagates until no label is changing.

LINE [26]: This method first separately learns node embeddings by preserving 1- and 2-step neighborhood information between nodes and then concatenates them as the final node embeddings.

Node2Vec [7]: This method converts graph structure to node sequences by mixing breadth- and depth-first random walk strategies and learns node embeddings with the skip-gram model [19].

²<http://socialcomputing.asu.edu/datasets/BlogCatalog3>

³<http://socialcomputing.asu.edu/datasets/Flickr>

⁴<https://snap.stanford.edu/node2vec/>

⁵<https://github.com/xiangyue9607/BioNEV>

Algorithm 1 The Optimization Procedure of MetaTNE

Input: Graph G , total number of steps N , decay rate γ , decay period N_{decay}

Output: The embedding matrix $\mathbf{U} \in \mathbb{R}^{|V| \times M}$, the function $Tr(\cdot)$

- 1: Randomly initialize \mathbf{U} and the parameters Θ of $Tr(\cdot)$
- 2: **for** $step = 0$ **to** N **do**
- 3: Calculate the threshold $\tau = 1/(1 + \gamma \lfloor \frac{step}{N_{\text{decay}}} \rfloor)$
- 4: Draw a random number $r \sim \text{Uniform}(0, 1)$
- 5: **if** $r < \tau$ **then** ▷ Optimize the structural module
- 6: Sample a batch of pairs $\{(v_i, v_j) | v_i \in \mathcal{V}, v_j \in \mathcal{N}(v_i)\}$
- 7: Update \mathbf{U} to optimize the objective function:

$$\min \sum_{v_i \in \mathcal{V}} \sum_{v_j \in \mathcal{N}(v_i)} \log \mathbb{P}(v_j | v_i) \quad (11)$$

- 8: **else** ▷ Optimize the meta-learning module
- 9: Sample a batch of tasks \mathcal{T}_i from $\mathcal{Y}_{\text{known}}$
- 10: **for all** $\mathcal{T}_i = (\mathcal{S}_i, \mathcal{Q}_i, y_i)$ **do**
- 11: **for all** $v_q \in \mathcal{Q}_i$ **do**
- 12: Calculate the adapted embeddings $\{\tilde{\mathbf{u}}_{k,q}^{(i)} | v_k \in \mathcal{V}_{\mathcal{S}_i^m}\}$ and $\tilde{\mathbf{u}}_{q,m}^{(i)}$, where $m \in \{+, -\}$, via Eqn. (10)
- 13: Calculate the prototypes $\tilde{\mathbf{c}}_{+,q}^{(i)}$ and $\tilde{\mathbf{c}}_{-,q}^{(i)}$:

$$\tilde{\mathbf{c}}_{m,q}^{(i)} = \frac{1}{|\mathcal{S}_i^m|} \sum_{v_k \in \mathcal{V}_{\mathcal{S}_i^m}} \tilde{\mathbf{u}}_{k,q}^{(i)}, \quad m \in \{+, -\} \quad (12)$$

- 14: Calculate the predicted probability that v_q holds y_i :

$$\hat{\ell}_{v_q, y_i} = \frac{\exp(-\text{dist}(\tilde{\mathbf{u}}_{q,+}^{(i)}, \tilde{\mathbf{c}}_{+,q}^{(i)}))}{\sum_{m \in \{+, -\}} \exp(-\text{dist}(\tilde{\mathbf{u}}_{q,m}^{(i)}, \tilde{\mathbf{c}}_{m,q}^{(i)}))} \quad (13)$$

- 15: **end for**
- 16: **end for**
- 17: Update \mathbf{U} and Θ to optimize the objective function:

$$\min_{\mathbf{U}, \Theta} \sum_{\mathcal{T}_i} \sum_{(v_q, \ell_{v_q, y_i}) \in \mathcal{Q}_i} \mathcal{L}(\hat{\ell}_{v_q, y_i}, \ell_{v_q, y_i}) + \lambda \sum \|\Theta\|_2^2, \quad (14)$$

- 18: **end if**
 - 19: **end for**
-

Algorithm 2 Applying MetaTNE to Few-Shot Novel Labels

Input: The embedding matrix \mathbf{U} , the function $Tr(\cdot)$, a novel label $y \in \mathcal{Y}_{\text{novel}}$, associated positive support nodes $\mathcal{V}_{\mathcal{S}+}$ and negative support nodes $\mathcal{V}_{\mathcal{S}-}$, query nodes $\mathcal{V}_{\mathcal{Q}}$

Output: The predicted probability $\hat{\ell}_{v_q, y}$ for each query node v_q

- 1: Look up in \mathbf{U} to get the support and query embeddings $\mathbf{u}_k, \mathbf{u}_q$.
 - 2: **for** v_q in $\mathcal{V}_{\mathcal{Q}}$ **do**
 - 3: Adapt v_q together with $\mathcal{V}_{\mathcal{S}+}$ according to Eqn. (10) and obtain adapted embeddings $\{\tilde{\mathbf{u}}_{q,+}\} \cup \{\tilde{\mathbf{u}}_{k,q} | v_k \in \mathcal{V}_{\mathcal{S}+}\}$.
 - 4: Adapt v_q together with $\mathcal{V}_{\mathcal{S}-}$ according to Eqn. (10) and obtain adapted embeddings $\{\tilde{\mathbf{u}}_{q,-}\} \cup \{\tilde{\mathbf{u}}_{k,q} | v_k \in \mathcal{V}_{\mathcal{S}-}\}$.
 - 5: Calculate the positive and negative prototypes $\tilde{\mathbf{c}}_{m,q}, m \in \{+, -\}$ for classification according to Eqn. (12).
 - 6: Calculate the predicted probability with $\tilde{\mathbf{c}}_{m,q}$ and $\tilde{\mathbf{u}}_{q,m}$ according to Eqn. (13).
 - 7: **end for**
-

GCN [11]: This method is a semi-supervised method that uses a localized first-order approximation of spectral graph convolutions to exploit the graph structure. Here we use the identity matrix as the input feature matrix of GCN as suggested in [11].

Planetoid [32]: This is a semi-supervised method that learns node embeddings by using them to jointly predict node labels and node neighborhoods in the graph.

Meta-GNN [39]: This method directly applies MAML [6] to train GCN [11] in a meta-learning manner. Similarly, we use the identity matrix as the input feature matrix of GCN as suggested in [11].

Baseline Evaluation Procedure. We assess the performance of the baselines on the node classification tasks sampled from the test labels as follows: (1) For LP, we propagate the labels of the support nodes over the entire graph and inspect the predicted labels of the query nodes for each test tasks; (2) For each unsupervised network embedding method, we take the learned node embeddings as features to train a logistic regression classifier with L2 regularization for each test task. We use the support set to train the classifier and then predict the labels of the query nodes; (3) For each semi-supervised network embedding method, we first use the training labels to train the model for multi-label node classification. Then, for each test task, we fine-tune the model by substituting the final classification layer with a binary classification layer. Analogous to (2), we use the support set to train the new layer and then predict the labels of the query nodes; (4) For Meta-GNN, we first employ MAML [6] to learn a good initialization of GCN on the training tasks (binary node classification tasks). Then, for each test task, we use the support set to update the GCN from the learned initialization and apply the adapted GCN to the query nodes.

B.3 Parameter Settings

For LP, we use an open-source implementation⁶ and set the maximum iteration number to 30. For fair comparisons, we set the dimension of node representations to 128 for LINE, Node2vec, and Planetoid. For LINE, we set the initial learning rate to 0.025 and the number of negative samples to 5. For Node2vec, we set the window size to 10, the length of each walk to 40, and the number of walks per node to 80. The best in-out and return hyperparameters are tuned on the validation tasks with a grid search over $p, q \in \{0.25, 0.5, 1, 2, 4\}$. For Planetoid, we use the variant Planetoid-G since there are no input node features in our datasets. We tune the respective batch sizes and learning rates used for optimizing the supervised and the structural objectives based on the performance on the validation tasks. For GCN, we use a two-layer GCN with the number of hidden units as 128 and ReLU nonlinearity, and tune the dropout rate, learning rate, and weight decay based on the performance on the validation tasks. and set other hyperparameters as the original paper. For Meta-GNN⁷, we also use a two-layer GCN with 128 hidden units and ReLU nonlinearity. We set the number of inner updates to 2 due to the limitation of GPU memory and tune the fast and meta learning rates based on the performance on the validation tasks. For Planetoid, GCN, and Meta-GNN, we apply the best performing models on the validation tasks to the test tasks.

For our proposed MetaTNE, there are three parts of hyperparameters. In the structural module, we need to set the size d of node representations and sample N_1 node pairs at each training step. We also sample N_{neg} negative nodes per pair to speed up the calculation as in [26]. In the meta-learning module, we sample N_2 training tasks at each training step. The hyperparameters involved in the transformation function include the number H of parallel attention heads, the size d'/H of the query, key, and value vectors, the size d_{ff} of the hidden layer in the two-layer feed-forward network, the number L of stacked computation blocks. Besides, we apply dropout to the output of each of the self-attention modules and the feed-forward networks before it is added to the corresponding input and normalized, and the dropout rate is denoted by P_{drop} . Another hyperparameter is the weight decay coefficient λ . In the optimization module, we use the Adam optimizer [10] to optimize the structural and the meta-learning modules with learning rates of α_1 and α_2 , respectively. In addition, we have the decay rate γ and the decay period N_{decay} to control the optimization of the structural and meta-learning modules.

For all four datasets, we set $d = 128$, $N_{\text{neg}} = 5$, $P_{\text{drop}} = 0.1$, and $\gamma = 0.1$. We tune other hyperparameters on the validation tasks over the search space shown in Table 4. We utilize the Ray

⁶https://github.com/yamaguchiuto/label_propagation

⁷Since the authors do not provide the implementation that uses GCN as the learner, we implement it on the basis of the released code at <https://github.com/ChengtaiCao/Meta-GNN> to perform experiments.

Table 4: The hyperparameter search space.

Hyperparameter	Values	Hyperparameter	Values
N_1	{512, 1024, 2048}	L	{1, 2, 3}
N_2	{32, 64, 128}	λ	{0.001, 0.01, 0.1}
H	{1, 2, 4}	α_1	{0.0001, 0.001}
d'	{128, 256}	α_2	{0.0001, 0.001}
d_{ff}	{256, 512}	N_{decay}	{500, 1000, 1500, 2000}

Tune library [16] with asynchronous HyperBand scheduler [14] to accelerate the searching process. Note that, for each dataset, we only search the best hyperparameters with $K_{*,+} = 10$ and $K_{*,-} = 20$ for both training and test tasks, and directly apply these hyperparameters to other experimental scenarios. The resulting hyperparameters are available in our attached code.

C Additional Experiments

C.1 Full Results of Overall Comparisons

The full results of overall comparisons in our original paper are presented in Table 5 in the form of mean \pm std. Overall, we observe that our proposed MetaTNE achieves comparable or even lower standard deviation, which demonstrates the statistical significance of the superiority of MetaTNE.

C.2 The Performance w.r.t. the Numbers of Positive and Negative Nodes

To further investigate the performance under different combinations of $K_{*,+}$ and $K_{*,-}$, we conduct experiments with $K_{*,+}$ fixed at either 10 or 20 while varying $K_{*,-}$ from 10 to 50 for both training and test tasks. Figure 4 gives the performance comparisons of MetaTNE and the best performing baseline (i.e., Planetoid) in terms of F_1 on BlogCatalog dataset. We observe that Planetoid and MetaTNE achieve comparable performance when $K_{*,+}$ is the same as or larger than $K_{*,-}$, while the performance gap between MetaTNE and Planetoid gradually increases as the ratio of $K_{*,+}$ to $K_{*,-}$ decreases, which demonstrates the practicability of our method since the positive nodes are relatively scarce compared with the negative ones in many realistic applications.

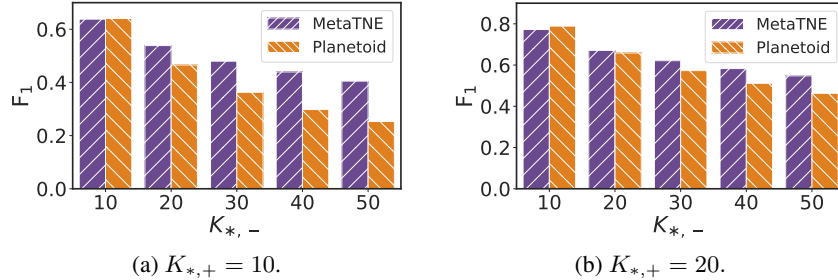


Figure 4: The performance w.r.t. the numbers of positive and negative nodes on BlogCatalog dataset.

C.3 The Performance w.r.t. the Number of Query Nodes

In the above experiments, we presume that, for each few-shot node classification task, the support and the query sets have the same numbers of positive and negative nodes following the standard protocol of meta-learning (called the *standard-setting*). However, in practice, the query set could have different numbers of positive and negative nodes as well as a different ratio of the number of positive nodes to the number of negative nodes compared to the support set. Thus, we further examine how the number of query nodes influences the performance. Towards this end, we sample additional test tasks by varying the numbers of positive and negative nodes in the query set (i.e., $K_{Q,+}^{\text{test}}$ and $K_{Q,-}^{\text{test}}$), with the numbers of positive and negative nodes in the support set fixed at 10 and 30, respectively (i.e., $K_{S,+}^{\text{test}} = 10$ and $K_{S,-}^{\text{test}} = 30$), and then compare the performance on these tasks. This setting is

Table 5: Results with standard deviation on few-shot node classification tasks with novel labels. OOM means out of memory (16 GB GPU memory).

(a) $K_{*,+} = 10$ and $K_{*,-} = 20$.

Method	BlogCatalog			Flickr		
	AUC	F ₁	Recall	AUC	F ₁	Recall
LP	0.6422 \pm 0.0289	0.1798 \pm 0.0198	0.2630 \pm 0.0309	0.8196 \pm 0.0175	0.4321 \pm 0.0392	0.4989 \pm 0.0492
LINE	0.6690 \pm 0.0323	0.2334 \pm 0.0499	0.1595 \pm 0.0403	0.8593 \pm 0.0145	0.6194 \pm 0.0334	0.5418 \pm 0.0382
Node2vec	0.6697 \pm 0.0325	0.3750 \pm 0.0478	0.2940 \pm 0.0432	0.8504 \pm 0.0151	0.6664 \pm 0.0284	0.6147 \pm 0.0332
Planetoid	0.6850 \pm 0.0320	0.4657 \pm 0.0437	0.4301 \pm 0.0451	0.8601 \pm 0.0360	0.6638 \pm 0.0796	0.6331 \pm 0.0821
GCN	0.6102 \pm 0.0285	0.2730 \pm 0.0415	0.2194 \pm 0.0392	OOM	OOM	OOM
Meta-GNN	0.4805 \pm 0.0364	0.2375 \pm 0.0365	0.2141 \pm 0.0392	OOM	OOM	OOM
MetaTNE	0.6986 \pm 0.0305	0.5380 \pm 0.0342	0.6203 \pm 0.0375	0.8462 \pm 0.0164	0.7118 \pm 0.0223	0.7700 \pm 0.0227
%Improv.	1.99	15.53	44.22	-1.62	6.81	21.62

Method	PPI			Mashup		
	AUC	F ₁	Recall	AUC	F ₁	Recall
LP	0.6285 \pm 0.0221	0.2147 \pm 0.0384	0.2769 \pm 0.0630	0.6488 \pm 0.0258	0.3103 \pm 0.0414	0.4535 \pm 0.0991
LINE	0.6372 \pm 0.0270	0.2147 \pm 0.0373	0.1456 \pm 0.0280	0.6926 \pm 0.0354	0.2970 \pm 0.0602	0.2142 \pm 0.0537
Node2vec	0.6273 \pm 0.0258	0.3545 \pm 0.0350	0.2860 \pm 0.0326	0.6575 \pm 0.0303	0.3835 \pm 0.0413	0.3147 \pm 0.0396
Planetoid	0.6791 \pm 0.0251	0.4672 \pm 0.0314	0.4411 \pm 0.0328	0.7056 \pm 0.0223	0.4825 \pm 0.0287	0.4218 \pm 0.0334
GCN	0.6544 \pm 0.0211	0.3379 \pm 0.0338	0.2721 \pm 0.0324	0.6895 \pm 0.0250	0.3052 \pm 0.0424	0.2390 \pm 0.0404
Meta-GNN	0.5466 \pm 0.0311	0.3289 \pm 0.0349	0.3081 \pm 0.0411	0.7078 \pm 0.0323	0.4576 \pm 0.0393	0.4176 \pm 0.0381
MetaTNE	0.6865 \pm 0.0205	0.5188 \pm 0.0209	0.5621 \pm 0.0311	0.7645 \pm 0.0251	0.5764 \pm 0.0291	0.5566 \pm 0.0337
%Improv.	1.09	11.04	27.43	8.01	19.46	22.73

(b) $K_{*,+} = 10$ and $K_{*,-} = 40$.

Method	BlogCatalog			Flickr		
	AUC	F ₁	Recall	AUC	F ₁	Recall
LP	0.6421 \pm 0.0288	0.0554 \pm 0.0118	0.0727 \pm 0.0158	0.8253 \pm 0.0156	0.3055 \pm 0.0413	0.3040 \pm 0.0485
LINE	0.6793 \pm 0.0320	0.0529 \pm 0.0316	0.0328 \pm 0.0216	0.8644 \pm 0.0139	0.4154 \pm 0.0471	0.3485 \pm 0.0471
Node2vec	0.6792 \pm 0.0314	0.1982 \pm 0.0516	0.1340 \pm 0.0398	0.8558 \pm 0.0150	0.5295 \pm 0.0381	0.4602 \pm 0.0420
Planetoid	0.6981 \pm 0.0315	0.2980 \pm 0.0550	0.2319 \pm 0.0507	0.8728 \pm 0.0382	0.5040 \pm 0.0790	0.4461 \pm 0.0741
GCN	0.6198 \pm 0.0297	0.1011 \pm 0.0345	0.0704 \pm 0.0265	OOM	OOM	OOM
Meta-GNN	0.4811 \pm 0.0405	0.1042 \pm 0.0589	0.0859 \pm 0.0558	OOM	OOM	OOM
MetaTNE	0.7139 \pm 0.0309	0.4398 \pm 0.0401	0.5819 \pm 0.0451	0.8505 \pm 0.0154	0.6220 \pm 0.0245	0.7460 \pm 0.0523
%Improv.	2.26	47.58	150.93	-2.55	17.47	62.10

Method	PPI			Mashup		
	AUC	F ₁	Recall	AUC	F ₁	Recall
LP	0.6298 \pm 0.0228	0.0773 \pm 0.0231	0.0748 \pm 0.0277	0.6534 \pm 0.0259	0.1156 \pm 0.0276	0.1284 \pm 0.0509
LINE	0.6423 \pm 0.0268	0.0496 \pm 0.0193	0.0300 \pm 0.0122	0.7009 \pm 0.0345	0.0956 \pm 0.0489	0.0617 \pm 0.0348
Node2vec	0.6309 \pm 0.0264	0.1894 \pm 0.0373	0.1306 \pm 0.0286	0.6643 \pm 0.0311	0.2070 \pm 0.0417	0.1447 \pm 0.0333
Planetoid	0.6879 \pm 0.0250	0.3100 \pm 0.0368	0.2523 \pm 0.0323	0.7095 \pm 0.0223	0.3279 \pm 0.0298	0.2551 \pm 0.0278
GCN	0.6608 \pm 0.0223	0.1403 \pm 0.0357	0.0957 \pm 0.0264	0.6979 \pm 0.0241	0.0813 \pm 0.0231	0.0531 \pm 0.0162
Meta-GNN	0.5399 \pm 0.0316	0.2085 \pm 0.0337	0.1867 \pm 0.0405	0.7050 \pm 0.0346	0.3279 \pm 0.0574	0.2768 \pm 0.0666
MetaTNE	0.7039 \pm 0.0218	0.4298 \pm 0.0242	0.5327 \pm 0.0420	0.7684 \pm 0.0244	0.4814 \pm 0.0318	0.4816 \pm 0.0393
%Improv.	2.33	38.65	111.14	8.30	46.81	73.99

called the *generalized-setting*. Note that here we only alter the sampling of test tasks as described above and the training tasks are always sampled under the condition that both the support and query sets contain 10 positive and 30 negative nodes (i.e., $K_{*,+}^{\text{train}} = 10$ and $K_{*,-}^{\text{train}} = 30$). Figure 5 shows the experimental results on PPI dataset.

We observe that MetaTNE consistently yields better performance than Planetoid under different combinations of $K_{*,+}^{\text{test}}$ and $K_{*,-}^{\text{test}}$. In particular, jointly analyzing Table 5 and Fig. 5a, MetaTNE achieves almost the same performance in both the standard- and generalized-settings when the query

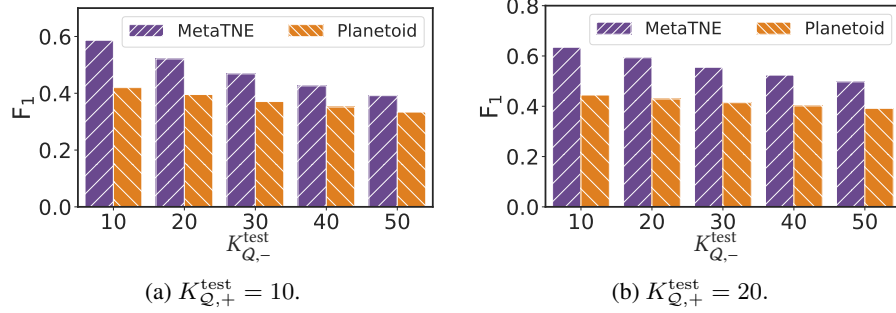


Figure 5: The performance w.r.t. the number of query nodes on PPI dataset.

set contains 10 positive nodes as well as 20 or 40 negative nodes, which indicates that to some extent MetaTNE is not sensitive to the choice of $K_{*,+}$ and $K_{*,,-}$ for sampling training tasks as well as $K_{S,+}^{\text{test}}$ and $K_{Q,+}^{\text{test}}$ and demonstrates the robustness of MetaTNE. On the other hand, it essentially becomes easier to classify the query nodes as the ratio of $K_{Q,+}^{\text{test}}$ to $K_{Q,-}^{\text{test}}$ increases, whereas the performance of Planetoid does not change markedly as $K_{Q,-}^{\text{test}}$ decreases in Fig. 5, which evidences that Planetoid tends to overfit the training tasks (e.g., the ratio of the number of positive nodes to the number of negative nodes).

C.4 The Performance with Fewer Positive Nodes

We further examine the performance of different methods by using fewer positive nodes and conduct experiments with $K_{*,+}$ set to 5 and $K_{*,,-}$ set to 10 or 20. Table 6 reports the experimental results on BlogCatalog dataset. From Table 6, we observe similar results to Table 5 and MetaTNE still significantly outperforms all other methods in the case that there are fewer positive nodes.

Table 6: Results of fewer positive nodes on BlogCatalog dataset.

Method	$K_{*,+} = 5, K_{*,,-} = 10$			$K_{*,+} = 5, K_{*,,-} = 20$		
	AUC	F1	Recall	AUC	F1	Recall
LP	0.6231 ± 0.0284	0.1753 ± 0.0168	0.2831 ± 0.0279	0.6226 ± 0.0288	0.0567 ± 0.0101	0.0930 ± 0.0159
LINE	0.6355 ± 0.0295	0.1296 ± 0.0379	0.0884 ± 0.0291	0.6432 ± 0.0300	0.0116 ± 0.0141	0.0076 ± 0.0098
Node2vec	0.6384 ± 0.0299	0.2912 ± 0.0440	0.2267 ± 0.0387	0.6451 ± 0.0305	0.1017 ± 0.0372	0.0689 ± 0.0273
Planetoid	0.6473 ± 0.0303	0.4221 ± 0.0408	0.4052 ± 0.0437	0.6583 ± 0.0318	0.2305 ± 0.0509	0.1853 ± 0.0470
GCN	0.5879 ± 0.0262	0.2176 ± 0.0336	0.1790 ± 0.0316	0.5971 ± 0.0283	0.0643 ± 0.0231	0.0464 ± 0.0178
Meta-GNN	0.4671 ± 0.0343	0.2673 ± 0.0342	0.2772 ± 0.0422	0.4580 ± 0.0394	0.0714 ± 0.0601	0.0630 ± 0.0573
MetaTNE	0.6546 ± 0.0286	0.4523 ± 0.0371	0.4842 ± 0.0469	0.6756 ± 0.0295	0.3730 ± 0.0387	0.4539 ± 0.0505
%Improv.	1.13	7.15	19.50	2.63	61.82	144.95

C.5 Visualization

To better demonstrate the effectiveness of the transformation function, we select two typical query nodes from the test tasks on Flickr dataset and visualize the relevant node embeddings before and after adaptation with t-SNE [18] in Fig. 6. Note that “Query (+)” and “Query (-)”, respectively, indicate the adapted embeddings of the query node in relation to the positive and negative support nodes in Eqn. (8). From Fig. 6a where the label of the query node is negative, we see that, before adaptation, the embedding of the query node is closer to the positive prototype than the negative prototype and thus misclassification occurs. After adaptation, the distance between “Query (-)” and the negative prototype is smaller than that between “Query (+)” and the positive prototype and hence the query node is classified correctly. The similar behavior is observed in Fig. 6b. Moreover, we observe that the transformation function is capable of either (1) gathering the positive and negative support nodes into two separate regions as shown in Fig. 6a or (2) adjusting “Query (+)” and “Query (-)” to make the right prediction when the positive and negative prototypes are close as shown in Fig. 6b. Another observation is that the transformation function has the tendency of enlarging the distances between node embeddings to facilitate classification.

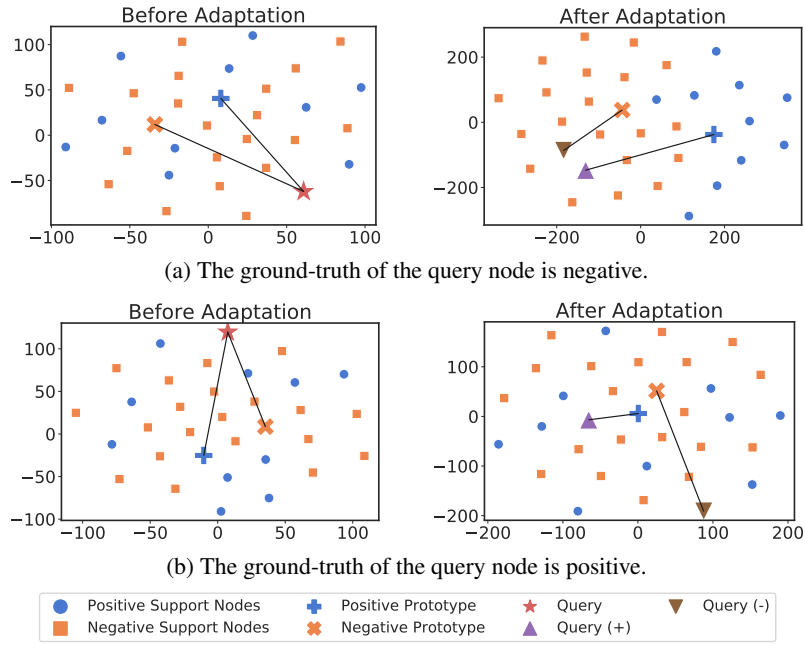


Figure 6: t-SNE visualization of embedding adaptation.