
Graph Few-shot Learning via Knowledge Transfer

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Towards the challenging problem of semi-supervised node classification, there have
2 been extensive studies. As a frontier, Graph Neural Networks (GNNs) have aroused
3 great interest recently, which update the representation of each node by aggregating
4 information of its neighbors. However, most GNNs have shallow layers with a
5 limited receptive field and may not achieve satisfactory performance especially
6 when the number of labeled nodes is quite small. To address this challenge, we
7 innovatively propose a graph few-shot learning (GFL) algorithm that incorporates
8 prior knowledge learned from auxiliary graphs to improve classification accuracy
9 on the target graph. Specifically, a transferable metric space characterized by a node
10 embedding and a graph-specific prototype embedding function is shared between
11 auxiliary graphs and the target, facilitating the transfer of structural knowledge.
12 Extensive experiments and ablation studies on four real-world graph datasets
13 demonstrate the effectiveness of our proposed model.

14 1 Introduction

15 Classifying a node (e.g., predicting interests of a user) in a graph in a semi-supervised manner
16 has been challenging but imperative, inasmuch as only a small fraction of nodes have access to
17 annotations which are usually costly. Recently, graph neural networks (GNN) [7, 14] have attracted
18 considerable interest and demonstrated promising performance. To their essential characteristics,
19 GNNs recursively update the feature of each node through aggregation (or message passing) of its
20 neighbors, by which the patterns of graph topology and node features are both captured. Nevertheless,
21 considering that adding more layers increases the difficulty of training and over-smoothens node
22 features [7], most of existing GNNs have shallow layers with a restricted receptive field. Therefore,
23 GNNs are inadequate to characterize the global information, and work not that satisfactorily when
24 the number of labeled nodes is especially small.

25 Inspired by recent success of few-shot learning, from an innovative perspective, we are motivated to
26 *leverage the knowledge learned from auxiliary graphs to improve semi-supervised node classification*
27 *in the target graph of our interest.* The intuition behind lies in that auxiliary graphs and the target
28 graph likely share local topological structures as well as class-dependent node features [11, 8]. Yet
29 it is even more challenging to achieve few-shot learning on graphs than on i.i.d. data (e.g., images)
30 which existing few shot learning algorithms focus on. The two lines of recent few-shot learning
31 works, including gradient-based methods [4, 10] and metric-based methods [12, 15], formulate
32 the transferred knowledge as parameter initializations (or a meta-optimizer) and a metric space,
33 respectively. None of them, however, meets the crucial prerequisite of graph few-shot learning to
34 succeed, i.e., transferring underlying structures across graphs.

35 To this end, we propose a novel **Graph Few-shot Learning (GFL)** model. Built upon metric-based
36 few-shot learning, the basic idea of GFL is to learn a transferable metric space in which the label of a
37 node is predicted as the class of the nearest prototype to the node. The metric space is practically
38 characterized with two embedding functions, which embed a node and the prototype of each class,
39 respectively. Specifically, first, GFL learns the representation of each node using a graph autoencoder
40 whose backbone is GNNs. Second, to better capture global information, we establish a relational
41 structure of all examples belonging to the same class, and learn the prototype of this class by applying

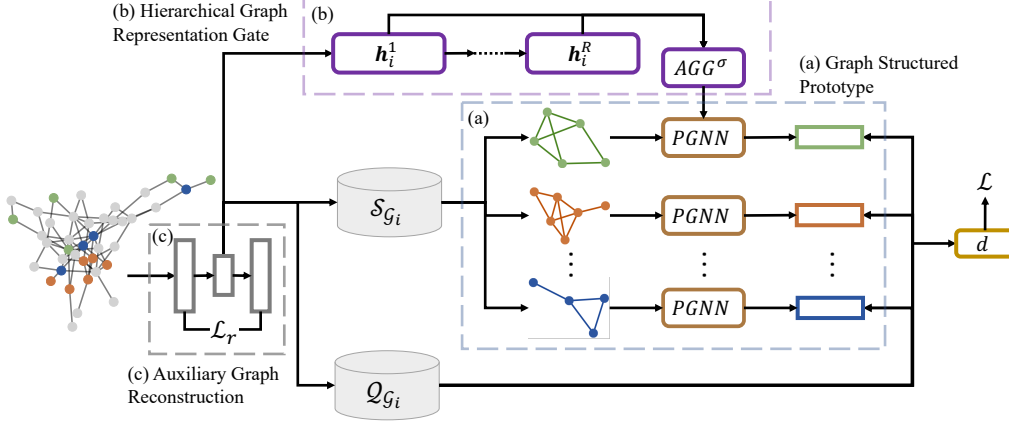


Figure 1: The framework of proposed GFL.

42 a prototype GNN to the relational structure. Most importantly, both embedding functions encrypting
 43 structured knowledge are transferred from auxiliary graphs to the target one, to remedy the lack
 44 of labeled nodes. Besides the two node-level structures, note that we also craft the graph-level
 45 representation via a hierarchical graph representation gate, to enforce that similar graphs have similar
 46 metric spaces. The experiments on node classification problem demonstrate the effectiveness of GFL.

47 2 Preliminaries

48 **Graph Neural Network** A graph \mathcal{G} is represented as (\mathbf{A}, \mathbf{X}) , where $\mathbf{A} \in \{0, 1\}^{n \times n}$ is the adjacent
 49 matrix, and $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \in \mathbb{R}^{n \times h}$ is the node feature matrix. To learn the node representation for
 50 graph \mathcal{G} , an embedding function f with parameter θ are defined as:

$$\mathbf{H}^{(l+1)} = \mathcal{M}(\mathbf{A}, \mathbf{H}^{(l)}; \mathbf{W}^{(l)}), \quad (1)$$

51 where \mathcal{M} is the message passing function. After stacking L graph neural network layers, we can get
 52 the final representation $\mathbf{Z} = \text{GNN}(\mathbf{A}, \mathbf{X}) = f_\theta(\mathbf{A}, \mathbf{X}) = \mathbf{H}^{(L+1)} \in \mathbb{R}^{h'}$.

53 **The Graph Few-Shot Learning Problem** In graph few-shot learning, we are given a sequence
 54 of graphs $\{\mathcal{G}_1, \dots, \mathcal{G}_{N_t}\}$ sampled from a probability distribution \mathcal{E} over tasks [3]. For each graph
 55 $\mathcal{G}_i \sim \mathcal{E}$, we are provided with a small set of n^{s_i} labeled *support* nodes set $\mathcal{S}_i = \{(\mathbf{x}_{i,j}^{s_i}, y_{i,j}^{s_i})\}_{j=1}^{n^{s_i}}$
 56 and a *query* nodes set $\mathcal{Q}_i = \{(\mathbf{x}_{i,j}^{q_i}, y_{i,j}^{q_i})\}_{j=1}^{n^{q_i}}$. For each node j in query set \mathcal{Q}_i , we are supposed to
 57 predict its corresponding label by associating its embedding $f_\theta(\mathbf{A}, \mathbf{x}_{i,j}^{q_i}) : \mathbb{R}^h \rightarrow \mathbb{R}^{h'}$ with represen-
 58 tation $(f_\theta(\mathbf{A}, \mathbf{x}_{i,j}^{s_i}), y_{i,j}^{s_i})$ in support set \mathcal{S}_i via the similarity measure d . Specifically, in prototypical
 59 network [12], the prototype \mathbf{c}_i^k for each class k is defined as $\mathbf{c}_i^k = \sum_{\mathbf{x}_{i,j}^{s_i} \in \mathcal{S}_i^k} f_\theta(\mathbf{A}, \mathbf{x}_{i,j}^{s_i}) / |\mathcal{S}_i^k|$, where
 60 \mathcal{S}_i^k denotes the sample set in \mathcal{S}_i of class k and $|\mathcal{S}_i^k|$ means the number of samples in \mathcal{S}_i^k . For each
 61 graph \mathcal{G}_i , the effectiveness on query set \mathcal{Q}_i is evaluated by the loss $\mathcal{L}_i = \sum_k \mathcal{L}_i^k$, where:

$$\mathcal{L}_i^k = - \sum_{(\mathbf{x}_{i,j}^{q_i}, y_{i,j}^{q_i}) \in \mathcal{Q}_i^k} \log \frac{\exp(-d(f_\theta(\mathbf{A}, \mathbf{x}_{i,j}^{q_i}), \mathbf{c}_i^k))}{\sum_{k'} \exp(-d(f_\theta(\mathbf{A}, \mathbf{x}_{i,j}^{q_i}), \mathbf{c}_i^{k'}))}, \quad (2)$$

62 where \mathcal{Q}_i^k is the query set of class k from \mathcal{Q}_i . To achieve this goal, few-shot learning often includes
 63 two-steps, i.e., meta-training and meta-testing. In meta-training, the parameter θ of embedding
 64 function f_θ is optimized to minimize the expected empirical loss over all historical training graphs,
 65 i.e., $\min_\theta \sum_{i=1}^{N_t} \mathcal{L}_i$. Once trained, given a new graph \mathcal{G}_t , the learned embedding function f_θ can be
 66 used to improve the learning effectiveness with a few support nodes.

67 3 Methodology

68 In this section, we elaborate our proposed GFL whose framework is illustrated in Figure 1. We will de-
 69 tail the three components of GFL, i.e., *graph structured prototype*, *hierarchical graph representation*
 70 *gate* and *auxiliary graph reconstruction*.

71 **Graph Structured Prototype** In most of the cases, a node plays two important roles in a graph: one
 72 is locally interacting with the neighbors that may belong to different classes; the other is interacting

73 with the nodes of the same class in relatively long distance, which can be globally observed. It
 74 is non-trivial to model the relational structure among support nodes and learn their corresponding
 75 prototype, we thus propose a prototype GNN model denoted as PGNN to tackle this challenge.

76 Given the representation of each node, we first extract the relational structure of samples belong
 77 to class k . For each graph \mathcal{G}_i , the relational structure \mathcal{R}_i^k of the sample set \mathcal{S}_i^k can be constructed
 78 based on the number of k -hop common neighbors. Then, the PGNN is used to model the interactions
 79 between samples in the \mathcal{S}_i^k , i.e., $\text{PGNN}_\phi(\mathcal{R}_i^k, f_\theta(\mathcal{S}_i^k))$, where PGNN is parameterized by ϕ and then
 80 we use j to indicate the j -th node representation (see part (a) in Figure 1).

$$\mathbf{c}_i^k = \text{Pool}_{j=1}^{n_i^{s_i^k}}(\text{PGNN}_\phi(\mathcal{R}_i^k, f_\theta(\mathcal{S}_i^k))[j]), \quad (3)$$

81 where Pool operator denotes a max or mean pooling operator over support nodes and $n_i^{s_i^k}$ represents
 82 the number of nodes in support set \mathcal{S}_i^k .

83 **Hierarchical Graph Representation Gate** The above prototype construction process is highly
 84 determined by the PGNN with the globally shared parameter ϕ . However, different graphs have their
 85 own topological structures, motivating us to tailor the globally shared information to each graph.
 86 Thus, we learn a hierarchical graph representation for extracting graph-specific information and
 87 incorporate it with the parameter of PGNN through a gate function (see part (b) in Figure 1 and more
 88 detailed structure is in Appendix A). Following [18], the hierarchical graph representation for each
 89 level is accomplished by alternating between two level-wise stages:

90 *I. Node Assignment* In the assignment step, each low-level node is assigned to high-level community.
 91 In level r , we denote the number of nodes as K^r , the adjacency matrix as \mathbf{A}_i^r , the feature matrix as
 92 \mathbf{X}_i^r . The assignment matrix $\mathbf{P}_i^{r \rightarrow r+1} \in \mathbb{R}^{K^r \times K^{r+1}}$ from level r to level $r+1$ is calculated by applying
 93 softmax function on the output of an assignment GNN (AGNN) as follows:

$$\mathbf{P}_i^{r \rightarrow r+1} = \text{Softmax}(\text{AGNN}(\mathbf{A}_i^r, \mathbf{X}_i^r)), \quad (4)$$

94 *II. Representation Fusion* After getting the assignment matrix, for level $r+1$, the adjacent matrix is
 95 defined as $\mathbf{A}_i^{r+1} = (\mathbf{P}_i^{r \rightarrow r+1})^T \mathbf{A}_i^r \mathbf{P}_i^{r \rightarrow r+1}$ and the feature matrix is calculated by applying assignment
 96 matrix on the output of a fusion GNN (FGNN), i.e., $\mathbf{X}_i^{r+1} = (\mathbf{P}_i^{r \rightarrow r+1})^T \text{FGNN}(\mathbf{A}_i^r, \mathbf{X}_i^r)$. Then, the
 97 feature representation \mathbf{h}_i^{r+1} of level $r+1$ can be calculated as:

$$\mathbf{h}_i^{r+1} = \text{Pool}_{k^{r+1}=1}^{K^{r+1}}((\mathbf{P}_i^{r \rightarrow r+1})^T \text{FGNN}(\mathbf{A}_i^r, \mathbf{X}_i^r)[k^{r+1}]), \quad (5)$$

98 Then, to get the whole graph-specific representation \mathbf{h}_i , the representation of each level is aggregated
 99 via an aggregator AGG, i.e., $\mathbf{h}_i = \text{AGG}(\mathbf{h}_i^1, \dots, \mathbf{h}_i^R)$. Inspired by previous findings [16]: similar
 100 graphs may activate similar parameters (i.e., parameter ϕ of the PGNN), we introduce a gate function
 101 $\mathbf{g}_i = \mathcal{T}(\mathbf{h}_i)$ to tailor graph structure specific information. Then, the global transferable knowledge
 102 (i.e., ϕ) is adapted to the structure-specific parameter via the gate function, i.e., $\phi_i = \mathbf{g}_i \circ \phi = \mathcal{T}(\mathbf{h}_i) \circ \phi$
 103 where \circ represents element-wise multiplication. $\mathbf{g}_i = \mathcal{T}(\mathbf{h}_i) = \sigma(\mathbf{W}_g \mathbf{h}_i + \mathbf{b}_g)$ maps the graph-specific
 104 representation \mathbf{h}_i to the same space of parameter ϕ . Thus, PGNN_ϕ in Eqn. (3) would be PGNN_{ϕ_i} .

105 **Auxiliary Graph Reconstruction** In practice, it is difficult to learn an informative node represen-
 106 tation using only the signal from the matching loss, which motivates us to design a new constraint for
 107 improving the training stability and the quality of node representation (see part (c) in Figure 1). Thus,
 108 for the node embedding function, we refine it by using a graph autoencoder and the reconstruction
 109 loss is defined as follows,

$$\mathcal{L}_r(\mathbf{A}_i, \mathbf{X}_i) = \|\mathbf{A}_i - \text{GNN}_{dec}(\mathbf{Z}_i) \text{GNN}_{dec}^T(\mathbf{Z}_i)\|_F^2, \quad (6)$$

110 where $\mathbf{Z}_i = \text{GNN}_{enc}(\mathbf{A}_i, \mathbf{H}_i)$ is the representation for each node. Recalling the objective in Section 2,
 111 we reach the optimization problem of GFL as $\min_{\Theta} \sum_{i=1}^{N_t} \mathcal{L}_i + \gamma \mathcal{L}_r(\mathbf{A}_i, \mathbf{X}_i)$, where Θ represents all
 112 learnable parameters. The whole training process (algorithm) of GFL is detailed in Appendix B.

113 4 Experiments

114 **Dataset and Experimental Settings** We evaluate our model performance on four tasks: (1) classi-
 115 fying 4 research domains of authors using AMiner collaboration data [2]; (2) classifying 5 communi-
 116 ties of posts using Reddit posts data [6]; (3) classifying 3 categories of papers using AMiner citation
 117 data [2]; (4) classifying 3 subjects of papers using PubMed data [14]. The details of these datasets are
 118 summarized in Appendix C. In this work, we follow the traditional K -way N -shot few-shot learning
 119 settings [4, 12]. For each graph, N labeled nodes for each class are provided. The rest nodes are used
 120 as query set. Like [7], the embedding structure is a two-layer graph convolutional structure (GCN)
 121 with 32 neurons in each layer. The distance metric d is defined as the inner product distance (see
 122 Appendix D for detailed hyperparameter settings).

Table 1: Comparison between GFL and other node classification methods on four graph datasets. Performance of Accuracy \pm 95% confidence intervals on 10-shot classification are reported.

Model	Collaboration	Reddit	Citation	Pubmed
LP [19]	61.09 \pm 1.36%	23.40 \pm 1.63%	67.00 \pm 4.50%	48.55 \pm 6.01%
Planetoid [17]	62.95 \pm 1.23%	50.97 \pm 3.81%	61.94 \pm 2.14%	51.43 \pm 3.98%
Deepwalk [9]	51.74 \pm 1.59%	34.81 \pm 2.81%	56.56 \pm 5.25%	44.33 \pm 4.88%
node2vec [5]	59.77 \pm 1.67%	43.57 \pm 2.23%	54.66 \pm 5.16%	41.89 \pm 4.83%
Base-GCN [7]	63.16 \pm 1.47%	46.21 \pm 1.43%	63.95 \pm 5.93%	54.87 \pm 3.60%
Finetune	76.09 \pm 0.56%	54.13 \pm 0.57%	88.93 \pm 0.72%	83.06 \pm 0.72%
K-NN	67.53 \pm 1.33%	56.06 \pm 1.36%	78.18 \pm 1.70%	74.33 \pm 0.52%
Matchingnet [15]	80.87 \pm 0.76%	56.21 \pm 1.87%	94.38 \pm 0.45%	85.65 \pm 0.21%
MAML [4]	79.37 \pm 0.41%	59.39 \pm 0.28%	95.71 \pm 0.23%	88.44 \pm 0.46%
Protonet [12]	80.49 \pm 0.55%	60.46 \pm 0.67%	95.12 \pm 0.17%	87.90 \pm 0.54%
GFL-mean (Ours)	83.51 \pm 0.38%	62.66 \pm 0.57%	96.51 \pm 0.31%	89.37 \pm 0.41%
GFL-att (Ours)	83.79 \pm 0.39%	63.14 \pm 0.51%	95.85 \pm 0.26%	88.96 \pm 0.43%

123 **Baseline Methods** For performance comparison of node classification, we consider three types
 124 of baselines: (1) *Graph-based semi-supervised methods* including Label Propagation (LP) [19] and
 125 Planetoid [17]; (2) *Graph representation learning methods* including Deepwalk [9], node2vec [5]
 126 and Base-GCN [7]. Note that, for Base-GCN, we train GCN on each meta-testing graph with limited
 127 labeled data rather than transferring knowledge from meta-training graphs; (3) *Transfer/few-shot*
 128 *methods* including finetune, K-nearest-neighbor (K-NN), Matching Network (Matchingnet) [15],
 129 MAML [4], Prototypical Network (Protonet) [12]. Each transfer/few-shot learning method uses the
 130 same embedding structure as GFL. Detailed description of baselines can be found in Appendix E.

131 **Overall Results** For each dataset, we reported the averaged accuracy with 95% confidence interval
 132 over meta-testing graphs of 10-shot node classification in Table 1. Both GFL-mean and GFL-
 133 att achieve the best performance than all three types of baselines on four datasets, indicating the
 134 effectiveness by incorporating graph prototype and hierarchical graph representation. In addition, as
 135 a metric distance based meta-learning algorithm, GFL not only outperforms other algorithms from
 136 this research line (i.e., Matchingnet, Protonet), but also achieves better performance than MAML,
 137 a popular gradient-based meta-learning algorithm. We also conduct extensive ablation studies and
 138 sensitivity analysis, which are reported in Appendix F and G, respectively.

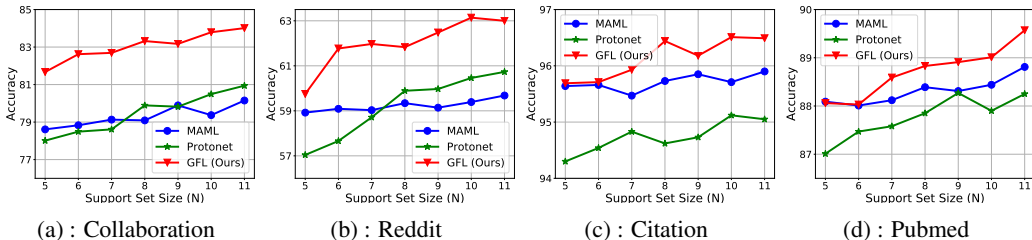


Figure 2: Effect of support set size, which is represented by the shot number N

139 **Effect of Support Set Size** In addition, we analyze the effect of support set size, which is repre-
 140 sented by the shot number N . We select two representative few-shot learning methods: Protonet
 141 (metric-learning based model) and MAML (gradient-based model). The results of each dataset are
 142 shown in Figure 2a-2d. We can see when the support set size is small, Protonet performs worse than
 143 MAML. The potential reason is that calculating prototype by averaging values over samples may be
 144 sensitive to outliers. Thus, it requires more data to reduce the effect of outliers. By extracting the
 145 relational structure among samples of the same class, our algorithm is more robust and achieves best
 146 performance in all scenarios.

147 **5 Conclusion**

148 In this paper, we introduce a new framework GFL to improve the learning effectiveness on a new
 149 graph by transferring knowledge from previous learned graphs. GFL improves metric-based few-shot
 150 learning by integrating graph structure from local node-level to global graph-level. We conduct
 151 extensive experiments and the results demonstrate the effectiveness of our proposed model on four
 152 node classification tasks.

References

- 153 [1] Lada A Adamic and Eytan Adar. Friends and neighbors on the web. *Social networks*, 25(3):211–
154 230, 2003.
- 155 [2] Aminer. <https://aminer.org/>, 2019.
- 156 [3] Jonathan Baxter. Theoretical models of learning to learn. In *Learning to learn*, pages 71–94.
157 Springer, 1998.
- 158 [4] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation
159 of deep networks. In *ICML*, pages 1126–1135, 2017.
- 160 [5] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *KDD*,
161 pages 855–864. ACM, 2016.
- 162 [6] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs.
163 In *NIPS*, pages 1024–1034, 2017.
- 164 [7] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional
165 networks. In *ICLR*, 2017.
- 166 [8] Danai Koutra, Vogelstein Joshua T., and Christos Faloutsos. Deltacon: A principled massive-
167 graph similarity function. In *SDM*, 2013.
- 168 [9] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social represen-
169 tations. In *KDD*, pages 701–710, 2014.
- 170 [10] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. *ICLR*, 2016.
- 171 [11] Nino Shervashidze, Pascal Schweitzer, Erik Jan van Leeuwen, Kurt Mehlhorn, and Karsten M
172 Borgwardt. Weisfeiler-lehman graph kernels. *JMLR*, 12(Sep):2539–2561, 2011.
- 173 [12] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In
174 *NIPS*, pages 4077–4087, 2017.
- 175 [13] Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Kelvin Xu, Ross Goroshin,
176 Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, and Hugo Larochelle. Meta-dataset:
177 A dataset of datasets for learning to learn from few examples. *arXiv preprint arXiv:1903.03096*,
178 2019.
- 179 [14] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua
180 Bengio. Graph attention networks. In *ICLR*, 2018.
- 181 [15] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks
182 for one shot learning. In *NIPS*, pages 3630–3638, 2016.
- 183 [16] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov,
184 Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with
185 visual attention. In *ICML*, pages 2048–2057, 2015.
- 186 [17] Zhilin Yang, William Cohen, and Ruslan Salakhudinov. Revisiting semi-supervised learning
187 with graph embeddings. In *ICML*, pages 40–48, 2016.
- 188 [18] Zhitao Ying, Jiaxuan You, Christopher Morris, Xiang Ren, Will Hamilton, and Jure Leskovec.
189 Hierarchical graph representation learning with differentiable pooling. In *NIPS*, pages 4805–4815,
190 2018.
- 191 [19] Xiaojin Zhu and Zoubin Ghahramani. Learning from labeled and unlabeled data with label
192 propagation. Technical report, Citeseer, 2002.
- 193

194 **A Detailed Framework of Hierarchical Graph Representation Gate**

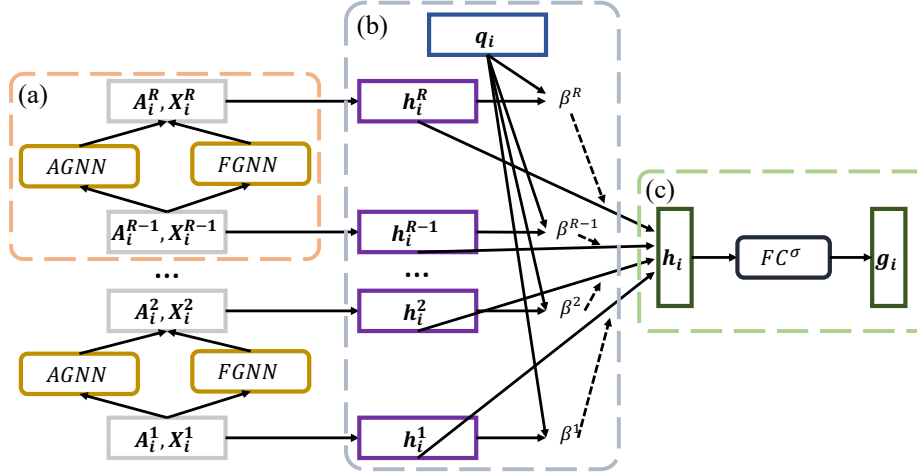


Figure 3: The detailed framework of hierarchical graph representation gate: (a) basic block for representation learning; (b) aggregator to aggregate hierarchical representations; (c) graph representation gate construction.

195 In Figure 3, we illustrate the detailed framework of hierarchical graph representation gate (detailed
 196 description of this component is in Section 4.2). Part (a) shows the basic block for learning hierarchical
 197 representation $\{\mathbf{h}_i^1, \dots, \mathbf{h}_i^R\}$. The aggregator is illustrated in part (b), where the learnable query
 198 vector \mathbf{q}_i is introduced to calculate the attention weight $\beta^1 \dots \beta^R$. Note that, we only illustrate the
 199 attention aggregator. For mean pooling aggregator, we calculate the average value of $\mathbf{h}_i^1, \dots, \mathbf{h}_i^R$ as
 200 the graph representation \mathbf{h}_i . Then, the graph representation \mathbf{h}_i is used to calculate gate \mathbf{g}_i by using a
 201 fully connected layer with sigmoid activation in part (c).

202 **B Meta-training Process of GFL**

Algorithm 1 Training Process of GFL

Require: \mathcal{E} : distribution over graphs; L : # of layers in hierarchical structure; α : stepsize; γ :
 balancing parameter for loss

- 1: Randomly initialize Θ
- 2: **while** not done **do**
- 3: Sample a batch of graphs $\mathcal{G}_i \sim \mathcal{E}$ and its corresponding adjacent matrices \mathbf{A}_i and feature
 matrices \mathbf{X}_i
- 4: **for all** \mathcal{G}_i **do**
- 5: Sample support set \mathcal{S}_i and query set \mathcal{Q}_i
- 6: Compute the embedding $f_\theta(\mathbf{A}_i, \mathbf{X}_i)$ and its reconstruction error $\mathcal{L}_r(\mathbf{A}_i, \mathbf{X}_i)$ in Eqn. (6)
- 7: Compute the hierarchical representation $\{\mathbf{h}_i^1, \dots, \mathbf{h}_i^R\}$ in Eqn. (5) and gate \mathbf{g}_i in Eqn. (??)
- 8: Construct relational graphs $\{\mathcal{R}_i^1, \dots, \mathcal{R}_i^K\}$ for samples in \mathcal{S}_i and compute graph prototype
 $\{\mathbf{c}_i^1, \dots, \mathbf{c}_i^K\}$ in Eqn. (3).
- 9: Compute the matching score using the query set \mathcal{Q}_i and evaluate loss in Eqn. (2)
- 10: **end for**
- 11: Update $\Theta \leftarrow \Theta - \alpha \nabla_{\Theta} \sum_{i=1}^{N_t} \mathcal{L}_i(\mathbf{A}_i, \mathbf{X}_i) + \gamma \mathcal{L}_r(\mathbf{A}_i, \mathbf{X}_i)$
- 12: **end while**

203 **C Detailed Dataset Description**

204 We use four datasets of different kinds of graphs: Collaboration, Reddit, Citation and Pubmed. (1):
 205 *Collaboration data*: Our first task is to predict research domains of different academic authors. We

206 use the collaboration graphs extracted from the AMiner data [2]. Each author is assigned with a
 207 computer science category label according to the majority of their papers’ categories. (2): *Reddit*
 208 *data*: In the second task, we predict communities of different Reddit posts. We construct post-to-post
 209 graphs from Reddit community data [6], where each edge denotes that the same user comments on
 210 both posts. Each post is labeled with a community id. (3): *Citation data*: The third task is to predict
 211 paper categories. We derive paper citation graphs from the AMiner data and each paper is labeled
 212 with a computer science category label. (4): *Pubmed data*: Similar to the third task, the last task is to
 213 predict paper class labels. The difference is that the citation graphs are extracted from the PubMed
 214 database [14] and each node is associated with diabetes class id. The statistics of these datasets are
 reported in Table 2.

Table 2: Data Statistics.

Dataset	Collaboration	Reddit	Citation	Pubmed
# Nodes (avg.)	4,496	5,469	2,528	2,901
# Edges (avg.)	14,562	7,325	14,710	5,199
# Features/Node	128	600	100	500
# Classes	4	5	3	3
# Graphs (Meta-training)	100	150	30	60
# Graphs (Meta-validation)	10	15	3	5
# Graphs (Meta-testing)	20	25	10	15

215

216 D Detailed Hyperparameter Settings

217 For more detailed hyperparameter settings, the learning rate is set as 0.01, the dimension of hier-
 218 archical graph representation \mathbf{h}_i is set as 32. For the hierarchical graph representation learning
 219 structure, we follow the strategy in [18] that the number of nodes in a higher layer is half of that in its
 220 consecutive lower layer. The specific values for the number of nodes in the 1st layer and the number
 221 of layers are determined by the tuning process on the validation set. In our experiments, the number
 222 of hierarchical layer R is set as 3 and the number of nodes in layer 2 and layer 3 are set as 64 and 32,
 223 respectively. The reconstruction loss weight γ is set as 1.0. Additionally, in order to construct the
 224 relational graph of few-shot labelled nodes for each class, we compute the similarity score between
 225 each two nodes by counting the number of k -hop ($k=3$) common neighbors and further smooth this
 226 similarity value by a sigmoid function. The computed similarity matrix is further fed into GNN for
 227 generating prototype embedding. We implement all experiments using Pytorch¹.

228 E Detailed Descriptions of Baselines

229 We detail three types of baselines in this section. Note that, all GCN used in baselines are two-layers
 230 GCN with 32 neurons each layer. The descriptions are as follows:

231 • Graph-based semi-supervised methods:

- 232 – **Label Propagation (LP)** [19]: Label Propagation is a traditional semi-supervised learning
 233 methods.
- 234 – **Planetoid** [17] Planetoid is a semi-supervised learning method based on graph embeddings.
 235 We use transductive formulation of Planetoid in this paper.

236 • Graph representation learning methods:

- 237 – **Deepwalk** [9]: Deepwalk learns the node embedding in an unsupervised way. We concatenate
 238 the learned node embedding and the node features, then feed them to the multiclass classifier
 239 (using Scikit-Learn²). Few-shot node labels in each graph are available for training classifier.
- 240 – **node2vec** [5]: This method is similar to the Deepwalk while we use node2vec model to learn
 241 node embedding.
- 242 – **Non-transfer-GCN** [7] In Non-transfer-GCN, we only train GCN on each meta-testing
 243 network without transferring knowledge from meta-training networks.

¹<https://pytorch.org>

²<https://scikit-learn.org/stable/>

244 • **Transfer/Few-shot methods**

- 245 – **All-Graph-Finetune (AGF)** In AGF, we train GCN by feeding each meta-training graph
246 one-by-one. Then, we can learn the initialization of GCN parameters from meta-training
247 graphs. In meta-testing process, we finetune the learned initialization on every meta-testing
248 graphs. The hyperparameters (e.g., learning rate) of AGF are the same as GFL.
- 249 – **K-nearest-neighbor (K-NN)** Similar as the settings of [13], we first learn the initialization
250 of GCN parameters by using all meta-training graphs. Then, in the testing process, we use the
251 learned embedding function (i.e., GCN) to learn the representation of support nodes and each
252 query node. Finally, we use k-NN for classification.
- 253 – **Matching Network (Matchingnet)** [15]: Matching network is a metric-based few-shot
254 learning method. Since traditional matching network focuses on one-shot learning, in our
255 scenario, we still use each query node representation to match the most similar on in support
256 set and use its label as the predicted label for the query node.
- 257 – **MAML** [4]: MAML is a representative gradient-based meta-learning method, which learn
258 a well-generalized model initialization which can be adapted with a few gradient steps. For
259 MAML, we set the learning rate of inner loop as 0.001.
- 260 – **Prototypical Network (Protonet)** [12] Prototypical network is a representative metric-based
261 few-shot learning method, which constructs the prototype by aggregating nodes belong to the
262 same class using mean pooling.

263 **F Ablation Studies**

264 Since GFL integrates three essential components (i.e., graph structured prototype, hierarchical
265 graph representation gate, auxiliary graph reconstruction), we conduct extensive ablation studies
266 to understand the contribution of each component. Table 3 shows the results of ablation studies
267 on each dataset, where the best results among GFL-att and GFL-mean are reported as GFL results.
268 Performance of accuracy are reported in this table. For the graph structured prototype, in (M1a), we
269 first report the performance of protonet for comparison since Protonet use mean pooling of node
270 embedding instead of constructing and exploiting relational structure for each class.

271 To show the effectiveness of hierarchical graph representation gate, we first remove this component
272 and report the performance in (M2a). The results are inferior, demonstrating that the effectiveness of
273 graph-level representation. In addition, we only use the flat representation structure (i.e., $R = 1$) in
274 (M2b). The results show the effectiveness of hierarchical representation.

275 For auxiliary graph reconstruction, we remove the decoder GNN and only use the encoder GNN to
276 learn the node representation in (M3). GFL outperforms (M3) as the graph reconstruction loss refines
277 the learned node representation and enhance the stability of training.

Table 3: Results of Ablation Studies. Accuracy scores on 10-shot node classification are reported. We select the best performance of GFL-mean and GFL-att for GFL in this table.

Ablation Model	Collab.	Reddit	Citation	Pubmed
(M1a): use the mean pooling prototype (i.e., protonet)	80.49%	60.46%	95.12%	87.90%
(M1b): replace the gate with the Film modulation	82.86%	62.66%	95.42%	88.92%
(M2a): remove the hierarchical representation gate	82.63%	61.99%	95.33%	88.15%
(M2b): use flat representation rather than hierarchical way	83.45%	62.55%	95.76%	89.08%
(M3): remove the graph reconstruction loss	82.98%	62.58%	95.63%	89.11%
GFL (Ours)	83.79%	63.14%	96.51%	89.37%

278 **G Sensitivity Analysis**

279 **G.1 Effect of Distance Functions d**

280 In addition, we replace the distance function d in Eqn. (2) from inner product (used in the original
281 setting) for cosine distance. The results are reported in Table 4. Compared with inner product, the
282 similar results show that GFL is not very sensitivity to the distance function d .

Table 4: Effect of different distance functions d .

Method	Collab.	Reddit	Cita.	Pubmed
GFL (inner product)	83.79%	63.14%	96.51%	89.37%
GFL (cosine)	84.02%	62.95%	96.02%	89.25%

283 G.2 Hyperparameter Sensitivity of Hierarchical Graph Representation Learning Structure

284 We investigate the performances of GFL with different number of nodes in the 1st layer and different
 285 number of layers in hierarchical graph representation structure. We show the results in Table 5. The
 286 number of nodes from bottom level to top level are represented in a tuple. The results demonstrate
 287 that GFL is not sensitive to such hyperparameters as long as the number of nodes in the 1st layer is not
 too small (e.g., 16).

Table 5: Effect of the number of nodes and layers in hierarchical graph representation learning structure. In the first column, the number of nodes from bottom level to top level are listed in a tuple.

Method	Collab.	Reddit	Cita.	Pubmed
(16,8,1)	83.47%	62.78%	95.89%	89.13%
(32,16,1)	83.61%	63.01%	95.92%	89.28%
(64,32,1)	83.79%	63.14%	96.51%	89.37%
(128,64,1)	83.73%	63.16%	96.43%	89.35%
(32,16,8,1)	83.67%	63.12%	96.45%	89.31%
(64,32,16,1)	83.78%	63.07%	96.38%	89.25%

288

289 G.3 Effect of Different Similarity Functions for Constructing Relational Structure \mathcal{R}_i^k

290 We further analyze the effect of different similarity functions for constructing relational structure of
 291 the graph prototype. Jaccard Index, Adamic-Adar [1], PageRank and Top-k Common Neighbors
 292 (Top-k CN) are selected and the results are reported in Table 6. Note that, in previous results, we all
 293 use Top-k CN as similarity function. The results show that GFL is not very sensitive to the similarity
 on a dataset may be achieved by different functions.

Table 6: Effect of different similarity functions for calculating relational structure \mathcal{R}_i^k . Results of Accuracy are reported.

Method	Collab.	Reddit	Cita.	Pubmed
Jaccard	82.98%	62.71%	95.18%	88.91%
Adamic-Adar	83.70%	62.87%	95.49%	89.21%
PageRank	84.14%	63.08%	95.93%	90.02%
Top-k CN	83.79%	63.14%	96.51%	89.37%

294