

# GLINKX: A SCALABLE UNIFIED FRAMEWORK FOR HOMOPHILOUS AND HETEROPHILOUS GRAPHS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

In graph learning, there have been two predominant inductive biases regarding graph-inspired architectures: On the one hand, higher-order interactions and message passing work well on homophilous graphs and are leveraged by GCNs and GATs. Such architectures, however, cannot easily scale to large real-world graphs. On the other hand, shallow (or node-level) models using ego features and adjacency embeddings work well in heterophilous graphs. In this work, we propose a novel scalable shallow method – GLINKX – that can work both on homophilous and heterophilous graphs. GLINKX leverages (i) novel monophilous label propagations (ii) ego/node features, (iii) knowledge graph embeddings as positional embeddings, (iv) node-level training, and (v) low-dimensional message passing. Formally, we prove novel error bounds and justify the components of GLINKX. Experimentally, we show its effectiveness on several homophilous and heterophilous datasets.

## 1 INTRODUCTION

In recent years, graph learning methods have emerged with a strong performance for various ML tasks. Graph ML methods leverage the topology of graphs underlying the data (Battaglia et al., 2018) to improve their performance. Two very important design options for proposing graph ML based architectures in the context of *node classification* are related to whether the data is *homophilous* or *heterophilous*.

For homophilous data – where neighboring nodes share similar labels (McPherson et al., 2001; Altenburger & Ugander, 2018a) – Graph Neural Network (GNN)-based methods are able to achieve high accuracy. Specifically, a broad subclass successful GNNs are Graph Convolutional Networks (GCNs) (e.g., GCN, GAT, etc.) (Kipf & Welling, 2016; Veličković et al., 2017; Zhu et al., 2020). In the GCN paradigm, *message passing* and *higher-order interactions* help node classification tasks in the homophilous setting since such inductive biases tend to bring the learned representations of linked nodes close to each other. However, GCN-based architectures suffer from *scalability issues*. Performing (higher-order) propagations during the training stage are hard to scale in large graphs because the number of nodes grows exponentially with the increase of the filter receptive field. Thus, for practical purposes, GCN-based methods require *node sampling*, substantially increasing their training time. For this reason, architectures (Huang et al., 2020; Zhang et al., 2022b; Sun et al., 2021; Maurya et al., 2021; Rossi et al., 2020) that leverage propagations outside of the training loop (as a preprocessing step) have shown promising results in terms of scaling to large graphs.

In *heterophilous* datasets (Rogers et al., 2014), the nodes that are connected tend to have different labels. Currently, many works that address heterophily can be classified into two categories concerning scale. On the one hand, recent successful architectures (in terms of accuracy) (Jin et al., 2022a; Di Giovanni et al., 2022; Zheng et al., 2022b; Luan et al., 2021; Chien et al., 2020; Lei et al., 2022) that address heterophily resemble GCNs in terms of design and thus suffer from the same scalability issues. On the other hand, *shallow or node-level models* (see, e.g., (Lim et al., 2021; Zhong et al., 2022)), i.e., models that are treating graph data as tabular data and do not involve propagations during training, has shown a lot of promise for large heterophilous graphs. In (Lim et al., 2021), it is shown that combining ego embeddings (node features) and adjacency embeddings works in the heterophilous setting. One element that LINKX exploits via the adjacency embeddings is monophily (Altenburger & Ugander, 2018a;b), namely the similarity of the labels of a node’s neighbors. However, their design is still impractical in real-world data since the method (LINKX) is *not* inductive (see Section 2),

and embedding the adjacency matrix directly requires many parameters in a model. In LINKX, the adjacency embedding of a node can alternatively be thought of as a *positional embedding* (PE) of the node in the graph, and recent developments (Kim et al., 2022; Dwivedi et al., 2021; Lim et al., 2021) have shown the importance of PEs in both homophilous and heterophilous settings. However, most of these works suggest PE parametrizations that are difficult to compute in large-scale settings. We argue that PEs can be obtained in a scalable manner by utilizing *knowledge graph embeddings*, which, according to prior work, can (El-Kishky et al., 2022; Lerer et al., 2019; Bordes et al., 2013; Yang et al., 2014) be trained in very large networks.

**Goal & Contribution:** In this work, we develop a scalable method for node classification that: (i) works both on homophilous and heterophilous graphs, (ii) is *simpler and faster* than conventional message passing networks (by avoiding the neighbor sampling and message passing overhead during training), and (iii) can work in both a *transductive* and an *inductive*<sup>1</sup> setting. For a method to be scalable, we argue that it should: (i) run models on node-scale (thus leveraging i.i.d. minibatching), (ii) avoid doing message passing during training and do it a constant number of times before training, and (iii) transmit small messages along the edges. Our proposed method – GLINKX (see Section 3) – combines all the above desiderata. GLINKX has three components: (i) ego embeddings<sup>2</sup>, (ii) PEs inspired by architectures suited for heterophilous settings, and (iii) *scalable* 2nd-hop-neighborhood propagations inspired by architectures suited for monophilous settings (Section 2.5). We provide novel theoretical error bounds and justify components of our method (Section 3.4). Finally, we evaluate GLINKX’s empirical effectiveness on several homophilous and heterophilous datasets (Section 4).

## 2 PRELIMINARIES

### 2.1 NOTATION

We denote scalars with lower-case, vectors with bold lower-case letters, and matrices with bold uppercase. We consider a directed graph  $G = G(V, E)$  with vertex set  $V$  with  $|V| = n$  nodes, and edge set  $E$  with  $|E| = m$  edges, and adjacency matrix  $\mathbf{A}$ .  $\mathbf{X} \in \mathbb{R}^{n \times d_X}$  represents the  $d_X$ -dimensional features and  $\mathbf{P} \in \mathbb{R}^{n \times d_P}$  represent the  $d_P$ -dimensional PE matrix. A node  $i$  has a feature vector  $\mathbf{x}_i \in \mathbb{R}^{d_X}$  and a positional embedding  $\mathbf{p}_i \in \mathbb{R}^{d_P}$  and belongs to a class  $y_i \in \{1, \dots, c\}$ . The training set is denoted by  $G_{\text{train}}(V_{\text{train}}, E_{\text{train}})$ , the validation set by  $G_{\text{valid}}(V_{\text{valid}}, E_{\text{valid}})$ , and test set by  $G_{\text{test}}(V_{\text{test}}, E_{\text{test}})$ .  $\mathbb{I}\{\cdot\}$  is the indicator function.  $\Gamma_c$  is the  $c$ -dimensional simplex.

### 2.2 GRAPH CONVOLUTIONAL NEURAL NETWORKS

In homophilous datasets, GCN-based methods have been used for node classification. GCNs (Kipf & Welling, 2016) utilize feature propagations together with non-linearities to produce node embeddings. Specifically, a GCN consists of multiple layers where each layer  $i$  collects  $i$ -th hop information from the nodes through propagations and forwards this information to the  $i + 1$ -th layer. More specifically, if  $G$  has a symmetrically-normalized adjacency matrix  $\mathbf{A}'_{\text{sym}}$  (with self-loops) (ignoring the directionality of edges), then a GCN layer has the form

$$\mathbf{H}^{(0)} = \mathbf{X}, \mathbf{H}^{(i+1)} = \sigma \left( \mathbf{A}'_{\text{sym}} \mathbf{H}^{(i)} \mathbf{W}^{(i)} \right) \forall i \in [L], \mathbf{Y} = \text{softmax} \left( \mathbf{H}^{(L)} \right).$$

Here  $\mathbf{H}^{(i)}$  is the embedding from the previous layer,  $\mathbf{W}^{(i)}$  is a learnable projection matrix and  $\sigma(\cdot)$  is a non-linearity (e.g. ReLU, sigmoid, etc.).

### 2.3 LINKX

In heterophilous datasets, the simple method of LINKX has been shown to perform well. LINKX combines two components – MLP on the node features  $\mathbf{X}$  and LINK regression (Altenburger & Ugander, 2018a) on the adjacency matrix – as follows:

<sup>1</sup>For this paper, we operate in the *transductive setting*. See App. B for the inductive setting.

<sup>2</sup>We use ego embeddings and node features interchangeably.

**Algorithm 1** GLINKX Algorithm**Input:** Graph  $G(V, E)$  with train set  $V_{\text{train}} \subseteq V$ , node features  $\mathbf{X}$ , labels  $\mathbf{Y}$ **Output:** Node Label Predictions  $\mathbf{Y}_{\text{final}}$ **1st Stage (KGEs).** Pre-train knowledge graph embeddings  $\mathbf{P}$  with Pytorch Biggraph.**2nd Stage (MLaP).** Propagate labels and predict neighbor distribution

1. **MLaP Forward:** Calculate the distribution of each training node’s neighbors, i.e.

$$\hat{\mathbf{y}}_i = \frac{\sum_{j \in V_{\text{train}}: (j,i) \in E_{\text{train}}} \mathbf{y}_j}{|\{j \in V_{\text{train}}: (j,i) \in E_{\text{train}}\}|} \text{ for all } i \in V_{\text{train}}$$

2. **Learn distribution of a node’s neighbors:**

- (a) For each epoch, calculate  $\tilde{\mathbf{y}}_i = f_1(\mathbf{x}_i, \mathbf{p}_i; \theta_1)$  for  $i \in V_{\text{train}}$

- (b) Update the parameters s.t. the negative cross-entropy

$$\mathcal{L}_{\text{CE},1}(\theta_1) = \sum_{i \in V_{\text{train}}} \text{CE}(\hat{\mathbf{y}}_i, \tilde{\mathbf{y}}_i; \theta_1) \text{ is maximized in order to bring } \tilde{\mathbf{y}}_i \text{ statistically close to } \hat{\mathbf{y}}_i.$$

- (c) Let  $\theta_1^*$  be the parameters at the end of the training that correspond to the epoch with the best validation accuracy.

3. **MLaP Backward:** Calculate  $\mathbf{y}'_i = \frac{\sum_{j \in V: (i,j) \in E} \tilde{\mathbf{y}}_j}{|\{j \in V: (i,j) \in E\}|}$  for all  $i \in V_{\text{train}}$ , where

$$\tilde{\mathbf{y}}_j = f_1(\mathbf{x}_j, \mathbf{p}_j; \theta_1^*).$$

**3rd Stage (Final Model).** Predicting node’s own label distribution:

1. For each epoch, calculate  $y_{\text{final},i} = f_2(\mathbf{x}_i, \mathbf{p}_i, \mathbf{y}'_i; \theta_2)$ .

2. Update the parameters s.t. the negative cross-entropy

$$\mathcal{L}_{\text{CE},2}(\theta_2) = \sum_{i \in V_{\text{train}}} \text{CE}(y_i, y_{\text{final},i}; \theta_2) \text{ is maximized.}$$

Return  $\mathbf{Y}_{\text{final}}$ 

$$\mathbf{H}_X = \text{MLP}_X(\mathbf{X}), \mathbf{H}_A = \text{MLP}_A(\mathbf{A}), \mathbf{Y} = \text{ResNet}(\mathbf{H}_X, \mathbf{H}_A).$$

## 2.4 NODE CLASSIFICATION

In node classification problems on graphs, we have a model  $f(\mathbf{X}, \mathbf{Y}_{\text{train}}, \mathbf{A}; \theta)$  that takes as an input the node features  $\mathbf{X}$ , the training labels  $\mathbf{Y}_{\text{train}}$  and the graph topology  $\mathbf{A}$  and produces a prediction for each node  $i$  of  $G$ , which corresponds to the probability that a given node belongs to any of  $c$  classes (with the sum of such probabilities being one). The model is trained with back-propagation. Once trained, the model can be used for the prediction of labels of nodes in the test set.

There are two training regimes: *transductive* and *inductive*. In the transductive training regime, we have full knowledge of the graph topology (for the train, test, and validation sets) and the node features, and the task is to predict the labels of the validation and test set. In the inductive regime, only the graph induced by  $V_{\text{train}}$  is known at the time of training, and then the full graph is revealed for prediction on the validation and test sets. In real-world scenarios, such as online social networks, the dynamic nature of problems makes the inductive regime particularly useful.

## 2.5 HOMOPHILY, HETEROHILY &amp; MONOPHILY

**Homophily and Heterophily:** There are various measures of homophily in the GNN literature like node homophily and edge homophily Lim et al. (2021). Intuitively, homophily in a graph implies that nodes with similar labels are connected. GNN-based approaches like GCN, GAT, etc., leverage this property to improve the node classification performance. Alternatively, if a graph has low homophily – namely, nodes that connect tend to have different labels – it is said to be *heterophilous*. In other words, a graph is heterophilous if neighboring nodes do not share similar labels.

**Monophily:** Generally, we define a graph to be monophilous if the label of a node is similar to that of its neighbors’ neighbors<sup>3</sup>. Etymologically, the word “monophily” is derived from the Greek words “*monos*” (unique) and “*philos*” (friend), which in our context means that a node – regardless of its label – has neighbors of primarily one label. In the context of a directed graph, monophily can be

<sup>3</sup>A similar definition of monophily has appeared in (Altenburger & Ugander, 2018a), whereby many nodes have extreme preferences for connecting to a certain class.

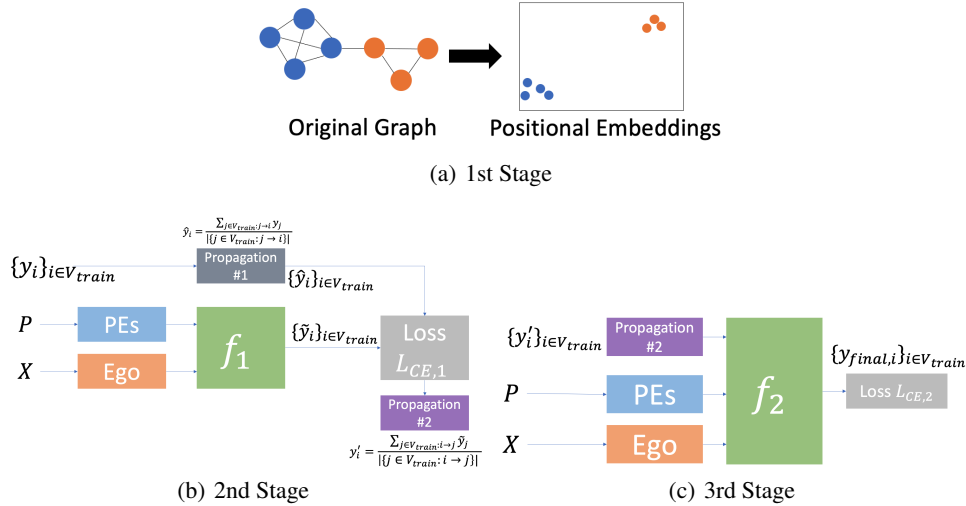


Figure 1: Block Diagrams of GLINKX stages.

thought of as a structure that resembles Fig. 2(a) where similar nodes (in this case, three green nodes connected to a yellow node) are connected to a node with a different label.

We argue that encoding monophily into a model can be helpful for both heterophilous and homophilous graphs (see Figs. 3(b) and 3(c)), which is one of the main motivators behind our work. In homophilous graphs, monophily will fundamentally encode the 2nd-hop neighbor’s label information, and since in such graphs, neighboring nodes have similar labels, it can provide a helpful signal for node classification. In heterophily, neighboring nodes have different labels, but the 2nd-hop neighbors may share the same label, providing helpful information for node classification. Monophily is effective for heterophilous graphs (Lim et al., 2021). Therefore, an approach encoding monophily has an advantage over methods designed specifically for homophilous and heterophilous graphs, especially when varying levels of homophily can exist between different sub-regions in the same graph (see Section 3.3). It may also not be apparent if the (sub-)graph is purely homophilous/heterophilous (since these are not binary constructs), which makes a unified architecture that can leverage graph information for both settings all the more important.

### 3 METHOD

#### 3.1 COMPONENTS & MOTIVATION

The desiderata we laid down on Section 1 can be realized by three components: (i) PEs, (ii) ego embeddings, and (iii) label propagations that encode monophily. More specifically, ego embeddings and PEs are used as primary features, which have been shown to work for both homophilous and heterophilous graphs for the models we end up training. Finally, the propagation step is used to encode monophily to provide additional information to our final prediction.

**Positional Embeddings:** We use PEs to provide our model information about the position of each node and hypothesize that PEs are an important piece of information in the context of large-scale node classification. PEs have been used to help discriminate isomorphic graph (sub)-structures (Kim et al., 2022; Dwivedi et al., 2021; Srinivasan & Ribeiro, 2019). This is useful for both homophily (Kim et al., 2022; Dwivedi et al., 2021) and heterophily (Lim et al., 2021) because isomorphic (sub)-structures can exist in both the settings. In the homophilous case, adding positional information can help distinguish nodes that have the same neighborhood but distinct position (Dwivedi et al., 2021; Morris et al., 2019; Xu et al., 2019), circumventing the need to do higher-order propagations (Dwivedi et al., 2021; Li et al., 2019; Bresson & Laurent, 2017) which are prone to over-squashing (Alon & Yahav, 2021). In heterophily, structural similarity among nodes is important for classification, as in the case of LINKX – where adjacency embedding can be considered a PE. However, in large graphs, using

adjacency embeddings or Laplacian eigenvectors (as methods such as (Kim et al., 2022) suggest) can be a computational bottleneck and may be infeasible.

In this work, we leverage *knowledge graph embeddings* (KGEs) to encode positional information about the nodes, **and embed the graph**. Using KGEs has two benefits: Firstly, KGEs can be trained quickly for large graphs. This is because KGEs compress the adjacency matrix into a fixed-sized embedding, and adjacency matrices have been shown to be effective in heterophilous cases. Further, KGEs are lower-dimensional than the adjacency matrix (e.g.,  $d_P \sim 10^2$ ), allowing for faster training and inference times. Secondly, KGEs can be pre-trained efficiently on such graphs (Lerer et al., 2019) and can be used off-the-shelf for other downstream tasks, including node classification (El-Kishky et al., 2022)<sup>4</sup>. So, in the 1st Stage of our methods in Alg. 1 (Fig. 1(a)) we train KGEs model on the available graph structure. Here, we fix this positional encoding once they are pre-trained for downstream usage. Finally, we note that this step is transductive but we can easily make it inductive (El-Kishky et al., 2022; Albooyeh et al., 2020).

**Ego Embeddings:** We get ego embeddings from the node features. Such embeddings have been used in homophilous and heterophilous settings (Lim et al., 2021; Zhu et al., 2020). Node embeddings are useful for tasks where the graph structure provides little/no information about the task.

**Monophilous Label Propagations:** We now propose a novel monophily (refer Section 2.5) inspired label propagation which we refer to as Monophilous Label Propagation (MLaP). MLaP has the advantage that we can use it both for homophilous and heterophilous graphs or in a scenario with varying levels of graph homophily (see Section 3.3) as it encodes monophily (Section 2.5).

To understand how MLaP encodes monophily, we consider the example in Fig. 2. In this example, we have three green nodes connected to a yellow node and two nodes of different colors connected to the yellow node. Then, one way to encode monophily in Fig. 2(a) while predicting label for  $j_\ell, \ell \in [5]$ , is to get a *distribution* of labels of nodes connected to node  $i$  thus encoding its neighbors’ distribution. The fact that there are more nodes with green color than other colors can be used by the model to make a prediction. But this information may only sometimes be present, or there may be few labeled nodes around node  $i$ . Consequently, we propose to use a model that predicts the label distribution of nodes connected to  $i$ . We use the node features ( $\mathbf{x}_i$ ) and PE ( $\mathbf{p}_i$ ) of node  $i$  to build this model since nodes that are connected to node  $i$  share similar labels, and thus, the features of node  $i$  must be predictive of its neighbors. So, in Fig. 2(a), we train a model to predict a distribution of  $i$ ’s neighbors. Next, we provide  $j_\ell$  the learned distribution of  $i$ ’s neighbors by propagating the learned distribution from  $i$  back to  $j_\ell$ . Eqs. (1) to (3) correspond to MLaP. We train a final model that leverages this information together with node features and PEs (Fig. 2(b)).

### 3.2 OUR METHOD: GLINKX

We put the components discussed in Section 3.1 together into three stages. In the first stage, we pre-train the PEs by using KGEs. Next, encode monophily into our model by training a model that predicts a node’s neighbors’ distribution and by propagating the soft labels from the fitted model. Finally, we combine the propagated information, node features, and PEs to train a final model. GLINKX is described in Alg. 1 and consists of three main components detailed as block diagrams in Fig. 1. Fig. 2 shows the GLINKX stages from Alg. 1 on a toy graph:

**1st Stage (KGEs):** We train DistMult KGEs with Pytorch-Biggraph (Yang et al., 2014) treating  $G$  as a knowledge graph with only one relation (see App. A.4 for more details). Here we have decided to use DistMult, but one can use their method of choice to embed the graph.

**2nd Stage (MLaP):** First (2nd Stage in Alg. 1, Fig. 1(b), and Fig. 2(a)), for a node we want to learn the distribution of *its neighbors*. To achieve this, we propagate the labels from a node’s neighbors (we call this step MLaP Forward), i.e. calculate

$$\hat{\mathbf{y}}_i = \frac{\sum_{j \in V_{\text{train}} : (j, i) \in E_{\text{train}}} \mathbf{y}_j}{|\{j \in V_{\text{train}} : (j, i) \in E_{\text{train}}\}|} \quad \forall i \in V_{\text{train}}. \quad (1)$$

<sup>4</sup>Positional information can also be provided by other methods such as node2vec (Grover & Leskovec, 2016), however, most of such methods are less scalable.

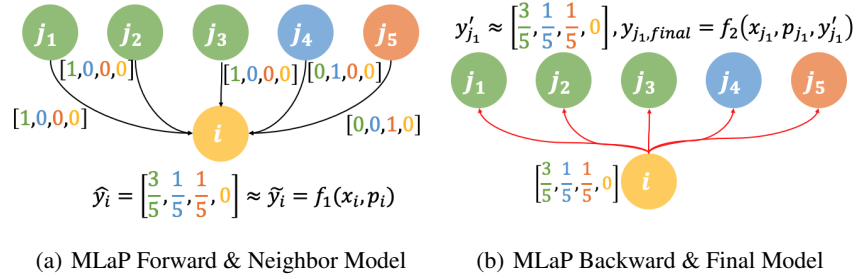


Figure 2: Example. For node  $i$  we want to learn a model that takes  $i$ 's features  $\mathbf{x}_i \in \mathbb{R}^{d_x}$ , and PEs  $\mathbf{p}_i \in \mathbb{R}^{d_p}$  and predict a value  $\tilde{y}_i \in \mathbb{R}^c$  that matches the label distribution of its neighbors neighbors  $\hat{\mathbf{y}}_i$  using a shallow model. Next, we want to propagate (outside the training loop) the (predicted) distribution of a node back to its neighbors and use it together with the ego features and the PEs to make a prediction about a node's own label. We propagate  $\tilde{\mathbf{y}}_i$  to its neighbors  $j_1$  to  $j_5$ . For example, for  $j_1$ , we encode the propagated distribution estimate  $\tilde{\mathbf{y}}_i$  from  $i$  to form  $\mathbf{y}'_{j_1}$ . We predict the label by using  $\mathbf{y}'_{j_1}, \mathbf{x}_{j_1}, \mathbf{p}_{j_1}$ .

Then, we train a model that predicts the distribution of neighbors, which we denote with  $\tilde{\mathbf{y}}_i$  using the ego features  $\{\mathbf{x}_i\}_{i \in V_{\text{train}}}$  and the PEs  $\{\mathbf{p}_i\}_{i \in V_{\text{train}}}$  and maximize the **negative** cross-entropy with treating  $\{\hat{\mathbf{y}}_i\}_{i \in V_{\text{train}}}$  as ground truth labels, namely we maximize

$$\mathcal{L}_{\text{CE}, 1}(\theta_1) = \sum_{i \in V_{\text{train}}} \sum_{l \in [c]} \hat{y}_{i,l} \log(\tilde{y}_{i,l}), \quad (2)$$

where  $\tilde{\mathbf{y}}_i = f_1(\mathbf{x}_i, \mathbf{p}_i; \theta_1)$  and  $\theta_1 \in \Theta_1$  is a learnable parameter vector. Although in this paper we assume to be in the *transductive setting*, this step allows us to be inductive (see App. B). In Section 3.4 we give a theoretical justification of this step, namely “*why is it good to use a parametric model to predict the distribution of neighbors?*”.

Finally, we propagate the predicted soft-labels  $\tilde{\mathbf{y}}_i$  back to the original nodes, i.e. calculate

$$\mathbf{y}'_i = \frac{\sum_{j \in V: (i,j) \in E} \tilde{\mathbf{y}}_j}{|\{j \in V : (i,j) \in E\}|} \quad \forall i \in V_{\text{train}}, \quad (3)$$

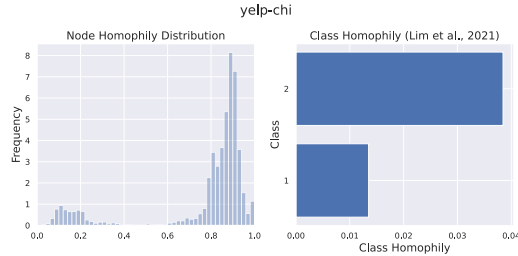
where the soft labels  $\{\tilde{\mathbf{y}}_i\}_{i \in V_{\text{train}}}$  have been computed with the parameter  $\theta_1^*$  of the epoch with the best validation accuracy from model  $f_1(\cdot | \theta_1)$ . We call this step **MLaP Backward**.

**3rd Stage (Final Model):** We make the final predictions  $\mathbf{y}_{\text{final}, i} = f_2(\mathbf{x}_i, \mathbf{p}_i, \mathbf{y}'_i; \theta_2)$  by combining the ego embeddings, PEs, and the (back)-propagated soft labels ( $\theta_2$  is a learnable parameter vector). We use the soft-labels  $\tilde{\mathbf{y}}_i$  instead of the actual labels one-hot ( $y_i$ ) in order to avoid label leakage, which hurts performance (see also (Shi et al., 2020) for a different way to combat label leakage). Finally, we maximize the **negative** cross-entropy with respect to a node's own labels,

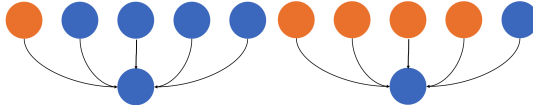
$$\mathcal{L}_{\text{CE}, 2}(\theta_2) = \sum_{i \in V_{\text{train}}} \sum_{l \in [c]} \mathbb{I}\{y_i = l\} \log(\mathbf{y}_{\text{final}, i, l}), \quad (4)$$

Overall, Stage 2 corresponds to learning the neighbor distributions and propagating these distributions, and Stage 3 uses these distributions to train a new model which predicts a node's labels. In Section 3.4, we prove that such a two-step procedure incurs lower errors than directly using the features to predict a node's labels.

**Scalability:** GLINKX is highly scalable as it performs message passing a constant number of times by paying an  $O(mc)$  cost, where the dimensionality of classes  $c$  is usually small (compared to  $d_X$  that GCNs rely on). In both Stages 2 and 3 of Alg. 1, we train node-level MLPs, which allow us to leverage i.i.d. (row-wise) mini-batching, like tabular data, and thus our complexity is similar to other shallow methods (LINKX, FSGNN) (Lim et al., 2021; Maurya et al., 2021). This, combined with the



(a) Node and class homophily



(b) Homophilous example (c) Heterophilous example

Figure 3: Top: Node and class homophily distributions for the yelp-chi dataset. Bottom: Examples of a homophilous (Fig. 3(b)) and a heterophilous (Fig. 3(c)) region in the same graph that are both monophilous, namely they are connected to many neighbors of the same kind. In a spam network, the homophilous region corresponds to many non-spam reviews connecting to non-spam reviews (which is the expected behaviour of a non-spammer user), and the heterophilous region corresponds to spam reviews targeting non-spam reviews (which is the expected behaviour of spammers), thus, yielding a graph with both homophilous and heterophilous regions such as in Fig. 3(a).

propagation outside the training loops, circumvents the scalability issues of GCNs. For more details, refer App. A.3.

### 3.3 VARYING HOMOPHILY

Graphs with monophily experience homophily, heterophily, or both. For instance, in the yelp-chi dataset – where we classify a review as spam/non-spam (see Fig. 3) – we observe a case of monophily together with varying homophily. Specifically in this dataset, spam reviews are linked to non-spam reviews, and non-spam reviews usually connect to other non-spam reviews, which makes the node homophily distribution bimodal. Here the 2nd-order similarity makes the MLaP mechanism effective.

### 3.4 THEORETICAL ANALYSIS

**Justification of Stage 2:** In Stage 2, we train a parametric model to learn the distribution of a node’s neighbors from the node features  $\xi_i$ <sup>5</sup>. Arguably, we can learn such a distribution naïvely by counting the neighbors  $i$  that belong to each class. This motivates our first theoretical result. In summary, we show that training a parametric model for learning the distribution of a node’s neighbors (as in Stage 2) yields a lower error than the naïve solution. Below we present the Thm. 1 (proof in App. F) for undirected graphs (the case of directed graphs is the same, but we omit it for simplicity of exposition):

**Theorem 1.** *Let  $G([n], E)$  be an undirected graph of minimum degree  $K > c^2$  and let  $Q_i \in \Gamma_c$  be the likelihood, from the viewpoint of node  $i$ , of any node in its neighborhood  $\mathcal{N}(i)$  to be assigned to different classes for every node  $i \in [n]$ . The following two facts are true (under standard assumptions for SGD and the losses):*

1. Let  $\hat{Q}_i$  be the sample average of  $Q_i$ , i.e.  $\hat{Q}_{i,j} = \frac{1}{|\mathcal{N}(i)|} \sum_{k \in \mathcal{N}(i)} \mathbb{I}\{y_k = j\}$ . Then, for every  $i \in [n]$ , we have that  $\max_{j \in [c]} \mathbb{E}[|Q_{i,j} - \hat{Q}_{i,j}|] \leq \mathbb{E}[\|Q_i - \hat{Q}_i\|_\infty] \leq O\left(\sqrt{\frac{\log(Kc)}{K}}\right)$ .
2. Let  $q(\cdot|\xi_i; \theta)$  be a model parametrized by  $\theta \in \mathbb{R}^D$  that uses the features  $\xi_i$  of each node  $i$  to predict  $Q_i$ . We estimate the parameter  $\theta_1$  by running SGD for  $t = n$  steps to maximize  $\mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^c Q_{i,j} \log q(j|\xi_i; \theta)$ . Then, for every  $i \in [n]$ , we have that  $\max_{j \in [c]} \mathbb{E}[|q(j; \xi_i; \theta_1) - Q_{i,j}|] \leq O\left(\sqrt{\frac{\log n}{n}}\right)$ .

<sup>5</sup>In Section 3.1,  $\xi_i$ s correspond to the augmented features  $\xi_i = [x_i; p_i]$

It is evident here that if the minimum degree  $K$  is much smaller than  $n$ , then the parametric model has lower error than the naïve approach, namely  $\hat{O}(n^{-1/2})$  compared to  $\hat{O}(K^{-1/2})$ .

**Justification of Stages 2 and 3:** We now provide theoretical foundations for the two-stage approach. Specifically, we argue that a two-stage procedure involving learning the distribution of a node’s 2nd-hop neighbor distributions (we assume for simplicity, again, that the graph is undirected) first with a parametric model such as in Thm. 1, and then running a two-phase algorithm to learn a parametric model that predicts a node’s label, yields a lower error than naïvely training a shallow parametric model to learn a node’s labels. The first phase of the two-phase algorithm involves training the model first by minimizing the cross-entropy between the predictions and the 2nd-hop neighborhood distribution. Then the model trains a joint objective that uses the learned neighbor distributions and the actual labels starting from the model learned in the previous phase.

**Theorem 2.** *Let  $G([n], E)$  be an undirected graph of minimum degree  $K > c^2$  and, let  $P_i$  be the likelihood of node  $i$  to be assigned to a different class, and let  $Q_i, q(\cdot|\xi_i; \theta_1)$  defined as in Thm. 1. Let  $p(\cdot|\xi_i; \mathbf{w})$  be a model parametrized by  $\mathbf{w} \in \mathbb{R}^D$  that is used to predict the class assignments  $y_i \sim p(\cdot|\xi_i; \mathbf{w})$ . Let  $\mathbf{w}_*$  be the optimal parameter. The following are true (under standard assumptions for SGD and the losses):*

1. *The naïve optimization scheme that runs SGD to maximize  $\mathcal{G}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^c P_{i,j} \log p(j|\xi_i; \mathbf{w})$  for  $n$  steps has error  $\mathbb{E}[\mathcal{G}(\mathbf{w}_{n+1}) - \mathcal{G}(\mathbf{w}_*)] \leq O\left(\frac{\log n}{n}\right)$ .*
2. *The two-phase optimization scheme that runs SGD to maximize  $\hat{\mathcal{G}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^c \left(\frac{1}{|\mathcal{N}(i)|} \sum_{k \in \mathcal{N}(i)} q(j|\xi_k; \theta_1)\right) \log p(j|\xi_i; \mathbf{w})$  for  $n_1$  steps, to estimate a solution  $\mathbf{w}'$  and then runs SGD on the objective  $\lambda \hat{\mathcal{G}}(\mathbf{w}) + (1 - \lambda) \mathcal{G}(\mathbf{w})$  for  $n$  steps starting from  $\mathbf{w}'$ , achieves error  $\mathbb{E}[\mathcal{G}(\mathbf{w}_{n+1}) - \mathcal{G}(\mathbf{w}_*)] \leq O\left(\frac{\sqrt{\log n \log \log n}}{n}\right)$ .*

You can find the proof in App. F. We observe that the two-phase optimization scheme can reduce the error by a factor of  $\sqrt{\log n / \log \log n}$  highlighting the importance of using the distribution of the 2nd-hop neighbors of a node to predict its label **and holds regardless of the homophily properties of the graph**. Also, note that the above two-phase optimization scheme differs from the description of the method we gave in Alg. 1. The difference is that the distribution of a node’s neighbors is embedded into the model in the case of Alg. 1, and the distribution of a node’s neighbors are embedded into the loss function in Thm. 2 as a regularizer. In Alg. 1, we chose to incorporate this information in the model because using multiple losses harms scalability and makes training harder in practice. In the same spirit, the conception of GCNs (Kipf & Welling, 2016) replaces explicit regularization with the graph Laplacian with the topology into the model (see also (Hamilton et al., 2017; Yang et al., 2016)).

### 3.5 COMPLEMENTARITY

Different components of GLINKX provide a *complementary* signal to components proposed in the GNN literature (Maurya et al., 2021; Zhang et al., 2022b; Rossi et al., 2020). One can combine GLINKX with existing architectures (e.g. feature propagations (Maurya et al., 2021; Rossi et al., 2020), label propagations (Zhang et al., 2022b)) for potential metric gains. For example, SIGN computes a series of  $r \in \mathbb{N}$  feature propagations  $[X, \Phi X, \Phi^2 X, \dots, \Phi^r X]$  where  $\Phi$  is a matrix (e.g., normalized adjacency or normalized Laplacian) as a preprocessing step. We can include this complementary signal, namely, embed each of the propagated features and combine them in the 3rd Stage to GLINKX. Overall, although in this paper we want to keep GLINKX simple to highlight its main components, we conjecture that adding more components to GLINKX would improve its performance on datasets with highly variable homophily (see Section 3.3).

## 4 EXPERIMENTS & CONCLUSIONS

**Comparisons:** We experiment with homophilous and heterophilous datasets (see Tab. 1 and App. D.3). We train KGEs with Pytorch-Bigraph (Lerer et al., 2019; Yang et al., 2014). For homophilous datasets, we compare with vanilla GCN and GAT, FSGNN, and Label Propagation (LP). For a fair comparison, we compare with one-layer GCN/GAT/FSGNN/LP since our method is one-hop. We also compare with higher-order (h.o.) GCN/GAT/FSGNN/LP with 2 and 3 layers. In the heterophilous



Table 1: Experimental results. (\*) = results from the OGB leaderboard.

	Homophilous Datasets		Heterophilous Datasets		
	PubMed	ogbn-arxiv	squirrel	yelp-chi	arxiv-year
$n$	19.7K	169.3K	5.2K	169.3K	45.9K
$m$	44.3K	1.16M	216.9K	7.73M	1.16M
Edge-insensitive homophily (Lim et al., 2021)	0.66	0.41	0.02	0.05	0.27
$d_X/c$	500 / 27	128 / 40	2089 / 5	32 / 2	128 / 5
GLINKX w/ KGEs	87.95 $\pm$ 0.30	<b>69.27 <math>\pm</math> 0.25</b>	45.83 $\pm$ 2.89	87.82 $\pm$ 0.20	<b>54.09 <math>\pm</math> 0.61</b>
GLINKX w/ Adjacency	<b>88.03 <math>\pm</math> 0.30</b>	69.09 $\pm$ 0.13	<b>69.15 <math>\pm</math> 1.87</b>	<b>89.32 <math>\pm</math> 0.45</b>	53.07 $\pm$ 0.29
Label Propagation (1-hop)	83.02 $\pm$ 0.35	<b>69.59 <math>\pm</math> 0.00</b>	32.22 $\pm$ 1.45	85.98 $\pm$ 0.28	43.71 $\pm$ 0.22
LINKX (from (Lim et al., 2021))	87.86 $\pm$ 0.77	67.32 $\pm$ 0.24	61.81 $\pm$ 1.80	85.86 $\pm$ 0.40	<b>56.00 <math>\pm</math> 1.34</b>
LINKX (our runs)	87.55 $\pm$ 0.37	63.91 $\pm$ 0.18	61.46 $\pm$ 1.60	88.25 $\pm$ 0.24	53.78 $\pm$ 0.06
GCN w/ 1 Layer	86.43 $\pm$ 0.74	50.76 $\pm$ 0.20		N/A	
GAT w/ 1 Layer	86.41 $\pm$ 0.53	54.42 $\pm$ 0.10		N/A	
FSGNN w/ 1 Layer	88.93 $\pm$ 0.31	61.82 $\pm$ 0.84	64.06 $\pm$ 2.69	86.36 $\pm$ 0.36	42.86 $\pm$ 0.22
Higher-order GCN	86.29 $\pm$ 0.46	71.18 $\pm$ 0.27 (*)		N/A	
Higher-order GAT	86.64 $\pm$ 0.40	<b>73.66 <math>\pm</math> 0.11</b> (*)		N/A	
Higher-order FSGNN	<b>89.37 <math>\pm</math> 0.49</b>	69.26 $\pm$ 0.36	68.04 $\pm$ 2.19	86.33 $\pm$ 0.30	44.89 $\pm$ 0.29
Label Propagation (2-hop)	83.44 $\pm$ 0.35	69.78 $\pm$ 0.00	43.41 $\pm$ 1.44	85.95 $\pm$ 0.26	46.30 $\pm$ 0.27
Label Prop. on $\mathbb{I}[A^2 - A - I \geq 0]$	82.14 $\pm$ 0.33	9.87 $\pm$ 0.00	24.43 $\pm$ 1.18	85.68 $\pm$ 0.32	23.08 $\pm$ 0.13

Table 2: Ablation Study. We use the hyperparameters of the best run from Tab. 1 with KGEs.

	Ablation Type	Stages	All	Remove ego embeddings	Remove propagation	Remove PEs
Heterophilous	arxiv-year	All Stages	54.09 $\pm$ 0.61	53.52 $\pm$ 0.77	50.83 $\pm$ 0.24	39.06 $\pm$ 0.35
	arxiv-year	3rd Stage	54.09 $\pm$ 0.61	53.69 $\pm$ 0.65	50.83 $\pm$ 0.24	49.13 $\pm$ 1.10
Homophilous	ogbn-arxiv	All Stages	69.27 $\pm$ 0.25	61.26 $\pm$ 0.33	62.70 $\pm$ 0.34	65.64 $\pm$ 0.18
	ogbn-arxiv	3rd Stage	69.27 $\pm$ 0.25	67.60 $\pm$ 0.39	62.70 $\pm$ 0.34	69.62 $\pm$ 0.15

case, we compare with LINKX<sup>6</sup> because it is scalable and is shown to work better than other baselines (e.g., H2GCN, MixHop, etc.) and with FSGNN. Note that we do not compare GLINKX with other more complex methods because GLINKX is complementary to them (see Section 3.5), and we can incorporate these designs into GLINKX. We use a *ResNet* module to combine our algorithm’s components from Stages 2 and 3. Details about the hyperparameters we use are in App. C.

In the heterophilous datasets, GLINKX outperforms LINKX (except for arxiv-year, where we are within the confidence interval). Moreover, the performance gap between using KGEs and adjacency embeddings shrinks as the dataset grows. In the homophilous datasets, GLINKX outperforms 1-layer GCN/GAT/LP/FSGNN and LINKX. In PubMed, GLINKX beats h.o. GCN/GAT and in arxiv-year GLINKX is very close to the performance of GCN/GAT.

Finally, we note that our method produces consistent results across regime shifts. In detail, in the heterophilous regime, our method performs on par with LINKX; however, when we shift to the homophilous regime, LINKX’s performance drops, whereas our method’s performance continues to be high. Similarly, while FSGNN performs similarly to GLINKX on the homophilous datasets, we observe a significant performance drop on the heterophilous datasets (see arxiv-year).

**Ablation Study:** We ablate each component of Alg. 1 to see each component’s performance contribution. We use the hyperparameters of the best model from Tab. 1. We perform two types of ablations: (i) we remove each of the components from all stages of the training, and (ii) we remove the corresponding components only from the 3rd Stage. Except for removing the PEs from the 3rd Stage only on ogbn-arxiv, all components contribute to increased performance on both datasets. Note that adding PEs in the 1st Stage does improve performance, suggesting the primary use case of PEs.

**Conclusion:** We present GLINKX, a scalable method for node classification in homophilous and heterophilous graphs that combines three components: (i) ego embeddings, (ii) PEs, and (iii) monophilous propagations. As future work, (i) GLINKX can be extended in heterogeneous graphs, (ii) use more expressive methods such as attention or Wasserstein barycenters (Cuturi & Doucet, 2014) for averaging the low-dimensional messages, and (iii) add complementary signals.

<sup>6</sup>We have run our method with hyperparameter space that is a subset of the sweeps reported in (Lim et al., 2021) due to resource constraints. A bigger hyperparameter search would improve our results.

## REFERENCES

- Marjan Albooyeh, Rishab Goel, and Seyed Mehran Kazemi. Out-of-sample representation learning for knowledge graphs. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 2657–2666, 2020.
- Uri Alon and Eran Yahav. On the bottleneck of graph neural networks and its practical implications. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=i800PhOCVH2>.
- Kristen M Altenburger and Johan Ugander. Monophily in social networks introduces similarity among friends-of-friends. *Nature human behaviour*, 2(4):284–290, 2018a.
- Kristen M Altenburger and Johan Ugander. Node attribute prediction: An evaluation of within-versus across-network tasks. In *NeurIPS Workshop on Relational Representation Learning*, 2018b.
- Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- K. Bhatia, K. Dahiya, H. Jain, P. Kar, A. Mittal, Y. Prabhu, and M. Varma. The extreme classification repository: Multi-label datasets and code, 2016. URL <http://manikvarma.org/downloads/XC/XMLRepository.html>.
- Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. *Advances in neural information processing systems*, 26, 2013.
- Xavier Bresson and Thomas Laurent. Residual gated graph convnets. *arXiv preprint arXiv:1711.07553*, 2017.
- Eli Chien, Jianhao Peng, Pan Li, and Olga Milenkovic. Adaptive universal generalized pagerank graph neural network. *arXiv preprint arXiv:2006.07988*, 2020.
- Marco Cuturi and Arnaud Doucet. Fast computation of wasserstein barycenters. In Eric P. Xing and Tony Jebara (eds.), *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pp. 685–693, Beijing, China, 22–24 Jun 2014. PMLR. URL <https://proceedings.mlr.press/v32/cuturi14.html>.
- Francesco Di Giovanni, James Rowbottom, Benjamin P Chamberlain, Thomas Markovich, and Michael M Bronstein. Graph neural networks as gradient flows. *arXiv preprint arXiv:2206.10991*, 2022.
- Yingtong Dou, Zhiwei Liu, Li Sun, Yutong Deng, Hao Peng, and Philip S Yu. Enhancing graph neural network-based fraud detectors against camouflaged fraudsters. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pp. 315–324, 2020.
- Vijay Prakash Dwivedi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. Graph neural networks with learnable structural and positional representations. *arXiv preprint arXiv:2110.07875*, 2021.
- Ahmed El-Kishky, Thomas Markovich, Serim Park, Chetan Verma, Baekjin Kim, Ramy Eskander, Yury Malkov, Frank Portman, Sofía Samaniego, Ying Xiao, et al. Twhin: Embedding the twitter heterogeneous information network for personalized recommendation. *arXiv preprint arXiv:2202.05387*, 2022.
- Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 855–864, 2016.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. *Advances in neural information processing systems*, 30, 2017.

- Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *Advances in neural information processing systems*, 33:22118–22133, 2020.
- Qian Huang, Horace He, Abhay Singh, Ser-Nam Lim, and Austin R Benson. Combining label propagation and simple models out-performs graph neural networks. *arXiv preprint arXiv:2010.13993*, 2020.
- Di Jin, Rui Wang, Meng Ge, Dongxiao He, Xiang Li, Wei Lin, and Weixiong Zhang. Raw-gnn: Random walk aggregation based graph neural network. *arXiv preprint arXiv:2206.13953*, 2022a.
- Wei Jin, Lingxiao Zhao, Shichang Zhang, Yozen Liu, Jiliang Tang, and Neil Shah. Graph condensation for graph neural networks. In *International Conference on Learning Representations*, 2022b. URL <https://openreview.net/forum?id=WLEx3Jo4QaB>.
- Jinwoo Kim, Tien Dat Nguyen, Seonwoo Min, Sungjun Cho, Moontae Lee, Honglak Lee, and Seunghoon Hong. Pure transformers are powerful graph learners. *arXiv preprint arXiv:2207.02505*, 2022.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- Jon M Kleinberg. Authoritative sources in a hyperlinked environment. *Journal of the ACM (JACM)*, 46(5):604–632, 1999.
- Runlin Lei, Zhen Wang, Yaliang Li, Bolin Ding, and Zhewei Wei. Evennet: Ignoring odd-hop neighbors improves robustness of graph neural networks. *arXiv preprint arXiv:2205.13892*, 2022.
- Adam Lerer, Ledell Wu, Jiajun Shen, Timothee Lacroix, Luca Wehrstedt, Abhijit Bose, and Alex Peysakhovich. Pytorch-biggraph: A large scale graph embedding system. *Proceedings of Machine Learning and Systems*, 1:120–131, 2019.
- Guohao Li, Matthias Muller, Ali Thabet, and Bernard Ghanem. Deepgcns: Can gcns go as deep as cnns? In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9267–9276, 2019.
- Derek Lim, Felix Hohne, Xiuyu Li, Sijia Linda Huang, Vaishnavi Gupta, Omkar Bhalerao, and Ser Nam Lim. Large scale learning on non-homophilous graphs: New benchmarks and strong simple methods. *Advances in Neural Information Processing Systems*, 34:20887–20902, 2021.
- Sitao Luan, Chenqing Hua, Qincheng Lu, Jiaqi Zhu, Mingde Zhao, Shuyuan Zhang, Xiao-Wen Chang, and Doina Precup. Is heterophily a real nightmare for graph neural networks to do node classification? *arXiv preprint arXiv:2109.05641*, 2021.
- Sunil Kumar Maurya, Xin Liu, and Tsuyoshi Murata. Improving graph neural networks with simple architecture design. *arXiv preprint arXiv:2105.07634*, 2021.
- Miller McPherson, Lynn Smith-Lovin, and James M Cook. Birds of a feather: Homophily in social networks. *Annual review of sociology*, pp. 415–444, 2001.
- Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen, Gaurav Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pp. 4602–4609, 2019.
- Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, and Bo Yang. Geom-gcn: Geometric graph convolutional networks. *arXiv preprint arXiv:2002.05287*, 2020.
- Everett M Rogers, Arvind Singhal, and Margaret M Quinlan. Diffusion of innovations. In *An integrated approach to communication theory and research*, pp. 432–448. Routledge, 2014.
- Emanuele Rossi, Fabrizio Frasca, Ben Chamberlain, Davide Eynard, Michael Bronstein, and Federico Monti. Sign: Scalable inception graph neural networks. *arXiv preprint arXiv:2004.11198*, 7:15, 2020.

- Benedek Rozemberczki, Carl Allen, and Rik Sarkar. Multi-scale attributed node embedding. *Journal of Complex Networks*, 9(2):cnab014, 2021.
- 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.
- Yunsheng Shi, Zhengjie Huang, Shikun Feng, Hui Zhong, Wenjin Wang, and Yu Sun. Masked label prediction: Unified message passing model for semi-supervised classification. *arXiv preprint arXiv:2009.03509*, 2020.
- Arnab Sinha, Zhihong Shen, Yang Song, Hao Ma, Darrin Eide, Bo-June Hsu, and Kuansan Wang. An overview of microsoft academic service (mas) and applications. In *Proceedings of the 24th international conference on world wide web*, pp. 243–246, 2015.
- Balasubramaniam Srinivasan and Bruno Ribeiro. On the equivalence between positional node embeddings and structural graph representations. *arXiv preprint arXiv:1910.00452*, 2019.
- Chuxiong Sun, Hongming Gu, and Jie Hu. Scalable and adaptive graph neural networks with self-label-enhanced training. *arXiv preprint arXiv:2104.09376*, 2021.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=ryGs6iA5Km>.
- Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. Embedding entities and relations for learning and inference in knowledge bases. *arXiv preprint arXiv:1412.6575*, 2014.
- Zhilin Yang, William Cohen, and Ruslan Salakhudinov. Revisiting semi-supervised learning with graph embeddings. In *International conference on machine learning*, pp. 40–48. PMLR, 2016.
- Shichang Zhang, Yozen Liu, Yizhou Sun, and Neil Shah. Graph-less neural networks: Teaching old MLPs new tricks via distillation. In *International Conference on Learning Representations*, 2022a. URL [https://openreview.net/forum?id=4p6\\_5HBWPCw](https://openreview.net/forum?id=4p6_5HBWPCw).
- Wentao Zhang, Ziqi Yin, Zeang Sheng, Yang Li, Wen Ouyang, Xiaosen Li, Yangyu Tao, Zhi Yang, and Bin Cui. Graph attention multi-layer perceptron. *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2022b.
- Wenqing Zheng, Edward W Huang, Nikhil Rao, Sumeet Katariya, Zhangyang Wang, and Karthik Subbian. Cold brew: Distilling graph node representations with incomplete or missing neighborhoods. In *International Conference on Learning Representations*, 2022a. URL <https://openreview.net/forum?id=1ugNpm7W6E>.
- Xin Zheng, Yixin Liu, Shirui Pan, Miao Zhang, Di Jin, and Philip S Yu. Graph neural networks for graphs with heterophily: A survey. *arXiv preprint arXiv:2202.07082*, 2022b.
- Zhiqiang Zhong, Sergey Ivanov, and Jun Pang. Simplifying node classification on heterophilous graphs with compatible label propagation. *arXiv preprint arXiv:2205.09389*, 2022.
- Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, and Danai Koutra. Beyond homophily in graph neural networks: Current limitations and effective designs. *Advances in Neural Information Processing Systems*, 33:7793–7804, 2020.
- Jiong Zhu, Ryan A Rossi, Anup Rao, Tung Mai, Nedim Lipka, Nesreen K Ahmed, and Danai Koutra. Graph neural networks with heterophily. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 11168–11176, 2021.