# Enabling Homogeneous GNNs to Handle Heterogeneous Graphs via Relation Embedding

Junfu Wang, Yuanfang Guo, *Senior Member, IEEE,*
Liang Yang, Yunhong Wang, *Fellow, IEEE*

**Abstract**—Graph Neural Networks (GNNs) have been generalized to process the heterogeneous graphs by various approaches. Unfortunately, these approaches usually model the heterogeneity via various complicated modules. This paper aims to propose a simple yet effective framework to assign adequate ability to the homogeneous GNNs to handle the heterogeneous graphs. Specifically, we propose Relation Embedding based Graph Neural Network (RE-GNN), which employs only one parameter per relation to embed the importance of distinct types of relations and node-type-specific self-loop connections. To optimize these relation embeddings and the model parameters simultaneously, a gradient scaling factor is proposed to constrain the embeddings to converge to suitable values. Besides, we interpret the proposed RE-GNN from two perspectives, and theoretically demonstrate that our RE-GCN possesses more expressive power than GTN (which is a typical heterogeneous GNN, and it can generate meta-paths adaptively). Extensive experiments demonstrate that our RE-GNN can effectively and efficiently handle the heterogeneous graphs and can be applied to various homogeneous GNNs.

**Index Terms**—Graph Neural Network, Heterogeneous Graph Neural Network, Heterogeneous Graph.

◆

## 1 INTRODUCTION

G RAPH Neural Networks (GNNs) have shown great expressive power in graph representation learning [1], [2], [3], [4], [5]. They have been widely adopted in various downstream applications, such as social network analysis [6], [7], protein prediction [8], [9], traffic prediction [10], [11], drug discovery [12], etc.

Unfortunately, the majorities of traditional GNNs are designed for homogeneous graphs, which disregard the variations in both the node type and edge type. The potential of GNNs in handling the data with complex relations has not been fully exploited. Currently, heterogeneous graph, a.k.a, heterogeneous information network, which possesses more than one type of nodes or edges, has attracted more attentions from researchers, due to its ability to represent data with more complicated relations.

Recently, several literatures have generalized GNNs to handle heterogeneous graphs. Existing heterogeneous GNNs can be classified into two categories, based on their mechanisms for handling the heterogeneity.

The first type of approaches utilizes the composite relations, i.e., meta-paths, to convert the heterogeneous graph to several meta-path based homogeneous graphs, as shown in Fig. 1b. For example, in the ACM dataset, the papers, which are not directly connected in the original heterogeneous graph yet connected with the same authors, are connected via a composite relation (*paper-author-paper*). These meta-path based heterogeneous GNNs [13], [14], [15], [16] utilize the handcrafted or learned meta-paths as shortcuts to boost the efficiency of the message passing among the target types



(a) Heterogeneous graph

(b) Meta-path based homogeneous graphs

(c) Relation based subgraphs

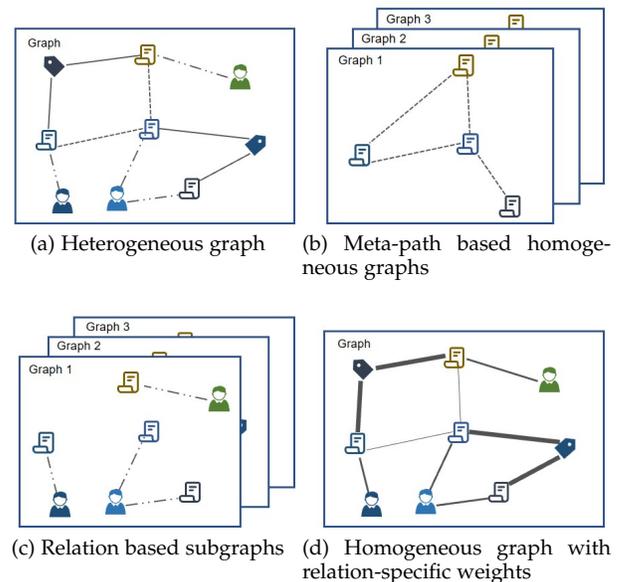(d) Homogeneous graph with relation-specific weights

Fig. 1: Different mechanisms to handle the heterogeneous graphs. A heterogeneous academic graph is utilized as an example.

of nodes. However, they require certain prior knowledge to either design the meta-paths [13], [14] or generate the meta-paths with excessive computations [15], [16]. Besides, the performances of these methods is highly correlated to the quality of the constructed meta-paths.

The other kind of approaches directly models the heterogeneity, especially the heterogeneous relations. They usually handle the subgraph of each type of relations (as shown in Fig. 1c) separately, via constructing different modules [17],

• J. Wang, Y. Guo, and Y. Wang are with the School of Computer Science and Engineering, Beihang University, Beijing 100191, China.
E-mail: {wangjunfu, andyguo, yhwang}@buaa.edu.cn.
• L. Yang is with the School of Artificial Intelligence, Hebei University of Technology, Tianjin 300401, China. E-mail: yangliang@vip.qq.com.

[18], [19], [20], [21]. Since their complexities are correlated to the number of relation types, they are less efficient when there exist many types of relations. Besides, [22] generates the relation representations and utilizes them as a part of the neighbor messages, i.e., concatenate them with the node representations. Then, the node attributes and relation representations are aggregated in a homogeneous manner.

Recently, [23] indicates that the homogeneous GNNs possess more potential to handle the heterogeneous graphs than previously reported [13], [14], [15]. It firstly reduces the complexity of the heterogeneous graph by converting its nodes and edges to a single type. Then, a specific homogeneous GNN, such as Graph Attention Network (GAT) [2], can achieve a comparable performance with the well-designed feature initialization strategies on certain datasets. Although [23] then proposes an improvement to make GAT more suitable for processing heterogeneous graphs, it has not provided a general mechanism for handling heterogeneity with various homogeneous GNNs. However, different GNNs [1], [2], [3], [4], [24] usually possess unique advantages in handling homogeneous graphs with different properties, which may be suitable for different applications. Thus, it is necessary to enable various homogeneous GNNs to process different heterogeneous graphs while preserving their advantages.

This paper aims to assign adequate abilities to the homogeneous GNNs for handling the heterogeneous graphs, by proposing a simple yet effective framework. The key of our work is to effectively convert the entire to-be-processed heterogeneous graph to one homogeneous graph, i.e., to homogenize the various types of nodes and relations. Generally, the attributes in different types of nodes are extracted from different perspectives. To homogenize the graph, the features in different types of nodes are firstly projected into the same feature space, by assuming that different types of nodes are correlated implicitly. For example, in a typical heterogeneous academic network, all types of nodes, e.g., papers, authors, conferences and subjects, are actually correlated to a certain extent. Then, each type of nodes can pass messages according to various relations. Considering the heterogeneity of relations, in the neighbourhood aggregation step, the messages from the nodes with different relations should not possess the same importance. Therefore, we must determine the appropriate importance for each relation in the neighbourhood aggregation step.

Specifically, we propose Relation Embedding based Graph Neural Network (RE-GNN), where only one embedding parameter is employed for each relation to model the aggregation importance. Then, the heterogeneous graph is adaptively converted to a weighted homogeneous graph by the proposed relation embeddings, as shown in Fig. 1d. Similar to the proposed relation embeddings, node-type-specific self-loop embeddings are exploited to add the self-loop connections. Then, the weighted graph with self-loops can be directly processed by the traditional homogeneous GNNs, such as [1], [2], [24]. Since the heterogeneity is embedded into the weighted homogeneous graph, the importance of a neighbor in the aggregation step is highly correlated with the type of its relation to the current node. With this simple yet effective framework, homogeneous GNNs can possess adequate abilities to handle the heterogeneous graphs.

To effectively model the heterogeneity, the learned weights of distinct relations are expected to be highly distinguishable. However, since the numerical values of relation embeddings are much larger than the other parameters, a straightforward simultaneous optimization of the relation embeddings and the other parameters cannot fully exploit the potential of relation embeddings. To tackle this *numerical inconsistency*, a gradient scaling factor is proposed to enlarge/suppress the updating modifications to the embeddings in each iteration, which enables the relation embeddings to gradually converge to the optimal values, to generate effective weights for the relations.

The proposed RE-GNN can be interpreted from two perspectives. Although RE-GNN is proposed from the perspective of spatial aggregations, it also possesses a meta-path based explanation. Specifically, when utilizing GCN [1] as backbone, we reveal that our RE-GCN can degenerate to the Graph Transformer Network (GTN) [15], a typical meta-path based heterogeneous GNN, which automatically generates composite relations. Besides, we theoretically demonstrate that our RE-GCN possesses more expressive power than GTN, without explicitly generating the meta-paths.

Extensive experiments demonstrate that our RE-GNN can effectively and efficiently handle the heterogeneous graphs and can be applied to various homogeneous GNNs.

Our contributions are summarized as follows:

- We propose a simple yet effective framework, named RE-GNN, to assign adequate ability to the homogeneous GNNs for handling heterogeneous graphs.
- To tackle the *numerical inconsistency*, we propose a gradient scaling factor to effectively optimize the relation embeddings and other model parameters simultaneously.
- We interpret the proposed RE-GNN from two perspectives, and theoretically demonstrate that our RE-GCN possesses more expressive power than GTN, a typical heterogeneous GNN.
- Extensive experiments demonstrate that our RE-GNN can effectively and efficiently handle the heterogeneous graphs and can be easily applied to various homogeneous GNNs.

## 2 RELATED WORK

### 2.1 Graph Neural Networks

Inspired by the traditional deep neural networks, Graph Neural Networks [25] are designed to handle the irregular graph data. Typical GNNs are usually designed from either the spectral or spatial perspectives. Spectral methods employ the graph spectral theory [26] to define the graph convolution operation [1], [27], [28], [29]. Spatial methods design the graph convolution operation directly on the graph, and aggregate the message from the spatial neighbors [2], [3], [30], [31]. Unfortunately, most of them are only designed to handle the homogeneous graphs.

### 2.2 Heterogeneous GNNs

To process the heterogeneous graph, researchers generalize traditional GNNs to heterogeneous graphs. Currently, there exist two types of approaches to model the heterogeneity.

The first type of approaches is developed based on meta-path constructions, which is firstly introduced by HAN [13]. It converts a heterogeneous graph to multiple homogeneous graphs by various manually-designed meta-paths. For each type of meta-paths, the meta-path based neighbors are connected in the corresponding homogeneous graph. Then, the results of each meta-path based homogeneous graph are fused by an attention scheme. Based on the above mechanism, MAGNN [14] utilizes the intermediate nodes along the meta-path, instead of only considering the meta-path based neighbors. These two methods require certain prior domain knowledge to design the meta-paths for each heterogeneous graph. GTN [15] firstly generates meta-paths via a soft selection of edge types. Then, an ensemble of GCNs is utilized to process the learned composite relations. Subsequently, the neural architecture search (NAS) technique [32] is utilized in DiffMG [16] to seek for a suitable meta-graph to model the more complicated composite relations.

The other type of approaches intends to model the heterogeneity directly, via different kinds of nodes and relations. R-GCN [17] models the relations in knowledge graphs by employing specialized parameter matrices, which separately constructs GNNs on each relation graph. R-GSN [18] further improves it by leveraging attention-based intra- and inter-relation aggregations. HetGNN [19] encodes node heterogeneous content and aggregates the neighbors from the same node type by using Bi-LSTMs, and combines the node embeddings from different node types with an attention scheme. HGT [20] utilizes a transformer network for each relation to model the importance of each relation. Besides, RSHN [22] constructs an edge-centric coarsened line graph to generate the relation representations, and then transfers the content of the node and relation representations to the target nodes, in the message passing process. R-HGNN [21] learns the relation representations and fuse the relation-aware representations of the target nodes semantically. All of these methods directly model the heterogeneity via complicated mechanism and possess excessive parameters. HGB [23] generalizes GAT [2] by adding the edge type attention to the original pair-wise self-attention mechanism. However, it does not provide a general mechanism to enable the homogeneous GNNs to handle the heterogeneity.

# 3 PRELIMINARIES

## 3.1 Heterogeneous Graph

A graph is defined as $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{F}, \mathcal{R}, \phi, \varphi)$, where $\mathcal{V}$ and $\mathcal{E}$ stand for the collections of nodes and edges, respectively. $\phi : \mathcal{V} \rightarrow \mathcal{F}$ and $\varphi : \mathcal{E} \rightarrow \mathcal{R}$ are the node type and edge type mapping functions, respectively. $\mathcal{F} = \{\phi(v) : \forall v \in \mathcal{V}\}$ represents the dynamic range of the node type projections and $\mathcal{R} = \{\varphi(e) : \forall e \in \mathcal{E}\}$ is the dynamic range of the edge type projections. A typical heterogeneous graph [33] contains more than one type of nodes or edges, i.e., $|\mathcal{F}| + |\mathcal{R}| > 2$. $A_i$ denotes the corresponding adjacency matrix of the edge type $i \in \mathcal{R}$. On the contrary, in homogeneous graphs, both the node and edge types only possess one valid value, i.e., $|\mathcal{F}| = |\mathcal{R}| = 1$.

## 3.2 Meta-path

Meta-path is widely adopted in the learning of heterogeneous graphs. A meta-path models a path passing through multiple relations (which may belong to different types of relations), e.g., $v_1 \xrightarrow{r_1} v_2 \xrightarrow{r_2} \dots \xrightarrow{r_l} v_{l+1}$, where $v_i \in \mathcal{V}$ and $r_j \in \mathcal{R}$. It describes a pair of nodes $v_1$ and $v_{l+1}$, which is connected by a composite relation. Note that $R = r_1 \odot r_2 \odot \dots \odot r_l$, where $\odot$ denotes the composite operator on relations. $v_1$ and $v_{l+1}$ are the meta-path based neighbors, and the corresponding composite adjacency matrix is $A_P = A_{r_1} A_{r_2} \dots A_{r_l}$.

## 3.3 Graph Convolutional Network

Graph Convolutional Network (GCN) [1] has become the most popular GNN in the past few years. Given a homogeneous graph $\mathcal{G}$, its graph convolution operation can be described as

$$H^{(l+1)} = \sigma(\tilde{A} H^{(l)} W^{(l)}), \tag{1}$$

where $\tilde{A} = \hat{D}^{-\frac{1}{2}} \hat{A} \hat{D}^{-\frac{1}{2}}$ is the normalized adjacency matrix, $\hat{A} = A + I$ denotes the adjacency matrix with self-loops, and $\hat{D}$ represents the corresponding degree matrix. $W^{(l)} \in \mathbb{R}^{d_{in}^{(l)} \times d_{out}^{(l)}}$ contains the learnable parameters. $H^{(l+1)}$ is the output of the $l$-th layer and the input of the $(l+1)$-th layer, and $H^{(0)} = X$. $\sigma$ stands for the non-linear activation function, e.g., ReLU. For a directed graph (i.e., asymmetric adjacency matrix), $\hat{A}$ can be normalized by the inverse of the degree matrix, i.e., $\hat{D}^{-1}$, as $\tilde{A} = \hat{D}^{-1} \hat{A}$.

## 3.4 Graph Transformer Network

Graph Transformer Network (GTN) [15] generates meta-paths in an end-to-end manner, which possesses the ability to search for the task-specific meta-paths from all the possible meta-paths. Since GTN is a typical meta-path based heterogeneous GNN and it will be compared to our approach in latter sections, we give a brief review here.

Firstly, GTN learns an adjacency matrix of $l$-length meta-paths via a GT layer

$$A_P = \left( \sum_{r_1 \in \mathcal{R}} \alpha_{r_1}^{(1)} A_{r_1} \right) \cdots \left( \sum_{r_l \in \mathcal{R}} \alpha_{r_l}^{(l)} A_{r_l} \right), \tag{2}$$

where $\alpha_{r_i}^{(j)} = \text{softmax}_j (w_{jr_i})$ is the soft weight for the edge type $r_i \in \mathcal{R}$ and $w$ is the learnable parameters.

Then, the learned adjacency matrix is fed into the standard GCN as

$$Z^{(l)} = \sigma(\tilde{D}_P^{-1} \tilde{A}_P^{(l)} H^{(l)} W^{(l)}), \tag{3}$$

where $\tilde{A}_P^{(l)} = A_P^{(l)} + I$ and $\tilde{D}_P$ denotes the corresponding degree matrix.

At last, the representations learned by multiple generated meta-paths are concatenated as

$$H^{(l+1)} = \|_{i=1}^{C} Z_i^{(l)}, \tag{4}$$

where $\|$ denotes the concatenation operator and $C$ represents the number of generated meta-paths. This architecture can be viewed as an ensemble of GCNs on multiple generated meta-path relations.
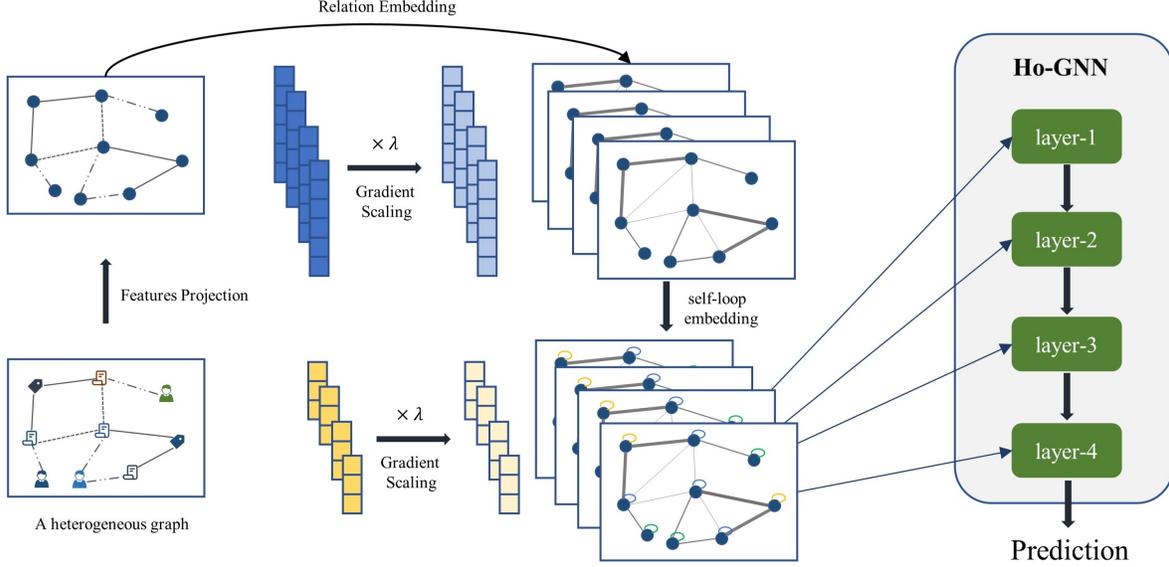
Fig. 2: The illustration of the Relation Embedding based Graph Neural Network (RE-GNN).

## 4 METHODOLOGY

To assign adequate abilities to the homogeneous GNNs for handling the heterogeneous graphs, in this section, we propose a simple yet effective framework, named Relation Embedding based Graph Neural Network (RE-GNN). Specifically, by assuming that all types of nodes are associated in particular relations, the node features are projected via a node-type-specific transformation matrix. Then, for the heterogeneous topology, one embedding parameter per relation is exploited to learn the importance of different types of relations and self-loop connections. Then, the heterogeneous graph can be converted to a weighted homogeneous graph, which can be handled by typical homogeneous GNNs. Besides, to ensure that the proposed relation embeddings can acquire effective information, a gradient scaling factor is proposed to adjust the updates in each optimizing iteration, which constrains the weights of the relations to converge to a proper value. A detailed illustration of our RE-GNN is presented in Fig. 2.

### 4.1 Feature Projection

Usually, the features in different types of nodes are extracted via various schemes, while these nodes are always associated in certain conditions. Therefore, the features in each node are firstly projected via a type-specific linear transformation [13], [15]. The projected features are represented as

$$\tilde{x}_i = x_i W_{\phi(i)}, \tag{5}$$

where $x_i$ is the original features in node $i$ and $W_{\phi(i)}$ represents the learnable projection matrix for the node type $\phi(i)$. Let $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, ..., \tilde{x}_n)$ be the matrix form of the features after transformation.

### 4.2 Relation Embeddings

After projecting the attributes of each type of nodes into the same feature space, nodes can exchange message with all their neighbors from different types of relations. In the aggregation step, the messages from different types of edges should possess different weights (importances). Then, one relation embedding value is utilized for each edge type to learn its importance. The learned weighted adjacency matrix is formulated as

$$A_H^{(l)} = \sum_{i \in \mathcal{R}} \tau(e_i^{(l)}) A_i, \tag{6}$$

where $e_i^{(l)} \in \mathbb{R}$ stands for the learnable relation embedding of edge type $i \in \mathcal{R}$ for the $l$-th layer, and $\tau(\cdot)$ is a function to generate the aggregation importances from the embeddings. For simplicity, we set $\tau(\cdot)$ to be the LeakyReLU function to ensure the importances of different relations to be non-negative. In different layers, the relation embeddings can learn different importances for various relation types.

Similarly, we consider the self-loop connections, which connect the node with itself, as special types of relations. Here, similar embedding values are employed to learn the importances of node-type-specific self-loop connections. The adjacency matrix with self-loops is formulated as

$$\hat{A}_H^{(l)} = A_H^{(l)} + \sum_{j \in \mathcal{F}} \tau(s_j^{(l)}) I_j, \tag{7}$$

where $I_j$ is the diagonal matrix and $(I_j)_{ii} = 1$, if and only if $\phi(i) = j$.

Then, the learned weighted adjacency matrix can be utilized by the homogeneous GNN layers, i.e.,

$$H^{(l+1)} = GNNLayer(\hat{A}_H^{(l)}, H^{(l)}), \tag{8}$$

where $H^{(0)} = \tilde{X}$ represents the projected node features and $GNNLayer(\cdot, \cdot)$ can be the layers in any homogeneous GNNs, such as GCN and GAT. If GCN is employed as the baseline model, the RE-GCN layers can then be represented by

$$H^{(l+1)} = \sigma(\tilde{A}_H^{(l)} H^{(l)} W^{(l)}), \tag{9}$$

where $\tilde{A}_H^{(l)}$ is the normalized version of $\hat{A}_H^{(l)}$. Note that the symmetric normalized version is $\tilde{A}_H = \hat{D}_H^{-\frac{1}{2}} \hat{A}_H \hat{D}_H^{-\frac{1}{2}}$,

which is usually utilized in processing the undirected graphs. Meanwhile, the asymmetric normalized version is $\tilde{A}_H = \hat{D}_H^{-1} \hat{A}_H$, which is usually applied to the directed graphs. Here, $\hat{D}_H$ denotes the degree diagonal matrix corresponding to $\hat{A}$.

## 4.3 Gradient Scaling

In the optimization process, the relation embeddings are firstly initialized as ones, where each type of relations has identical importance. As the model training progresses, the weights of the relations are expected to diversify. Since the absolute values of the relation embeddings are much larger than the other parameters, the learned embeddings are not differentiable enough. For example, a regular parameter $w$ is initialized via the Xavier Uniform Initialization [34] method by setting $|w| < a$, where $a$ is a small value correlated to the input and output dimensions and the relation embeddings are set to ones. Due to the above *numerical inconsistency*, it is difficult for the popular optimizers to simultaneously optimize the relation embeddings and model parameters. To tackle this problem, a gradient scaling factor is proposed to enlarge the gradients of the weights of the relations.

Specifically, a pre-defined scaling factor $\lambda > 0$ is exploited on the original embeddings as

$$\alpha_i = \lambda e_i. \tag{10}$$

Then, the scaled weights are utilized to generate the adjacency matrix as

$$A_H = \sum_{i \in \mathcal{R}} \tau(\alpha_i) A_i. \tag{11}$$

To initialize the scaled weights as ones, each embedding parameter $e_i$ is initially set to $\frac{1.0}{\lambda}$. The value of the scaling factor $\lambda$ will affect the updating modification to $\alpha$. In general, if $\lambda \geq 1$, Eq. (10) serves as a gradient enlarging machine and vice versa. To better illustrate the proposed gradient scaling factor, here we reveal how this scaling factor $\lambda$ functions in the optimization process.

By denoting the loss function as $\mathcal{L}$, the gradient of each relation embedding $e_i$ in Eq. (6) is $g_i = \frac{\partial \mathcal{L}}{\partial e_i}$. For the gradient based optimizers, the parameters are updated via $e_i^{next} = e_i - \Delta e_i$, where $\Delta e_i = \kappa(g_i)$. By utilizing the scaling factor in Eq. (10), the gradient to the embedding parameter $e_i$ is $\lambda g_i$. Then, the updating modification becomes $\kappa(\lambda g_i)$. For the common gradient optimizers, such as SGD, Momentum [35], and Nesterov [36], $\kappa(\lambda g_i) = \lambda \kappa(g_i)$. For the adaptive optimizers like Adagrad [37] and Adam [38], $\kappa(\lambda g_i) = \kappa(g_i)$, due to their adaptive learning strategies. By considering the modifications of the weights, since the scaling factor $\lambda > 0$, the modification $\Delta \alpha_i$ becomes $\lambda^2 \Delta \alpha_i$ for the common optimizers like SGD, and $\lambda \Delta \alpha_i$ for the adaptive optimizers like Adam. The detailed proof can be found in Appendix A.

In general, with the proposed scaling factor $\lambda$, the weights of the relations can be optimized within a suitable numerical interval. Similarly, the adjacency matrix with self-loops can be obtained via

$$\hat{A}_H = \sum_{i \in \mathcal{R}} \tau(\alpha_i) A_i + \sum_{j \in \mathcal{F}} \tau(\beta_j) I_j, \tag{12}$$

where $\beta_j = \lambda s_j$.

## 4.4 Memory and Computational Complexities

Each RE-GNN layer only employs one parameter for each relation type to model the importance of the corresponding relation. Generally, our RE-GNN only introduces $|\mathcal{F}| + |\mathcal{R}|$ more parameters than the corresponding homogeneous GNN in each layer. For the computational complexity, each RE-GNN layer only introduces $O(|\mathcal{E}| + |\mathcal{V}|)$ more calculations to compute the weighted adjacency matrix by Eq. (12), which can usually be neglected by the computational complexity of the corresponding homogeneous GNN layer, such as [1], [2], [39]. Considering both the memory and computational complexities, our RE-GNN is indeed efficient.

## 5 ANALYSIS

In this section, we analyze the effectiveness of our RE-GNN from two distinct perspectives. Firstly, we provide a spatial analysis which intuitively illustrates the effectiveness of our framework. Then, we reveal that our RE-GNN can be interpreted as a meta-path based heterogeneous GNN. Specifically, when employing GCN as the backbone, our RE-GCN can degenerate to Graph Transformer Network (GTN) [15], which is a powerful meta-path based heterogeneous GNN and can generate meta-paths adpatively. At last, we theoretically demonstrate that our RE-GCN possesses more expressive power than GTN.

### 5.1 Spatial Analysis

Under the assumption that all types of nodes are implicitly correlated, our RE-GNN firstly projects the raw features in different types of nodes into the same feature space. Then, RE-GNN learns a weighted relation adjacency matrix for the homogeneous GNN, to model the importance of different relations. For example, for a paper node in the ACM network, different messages from different relations (e.g., Paper-Paper, Conference-Paper, Subject-Paper and node-type-specific self-loop (Paper)) are aggregated based on the corresponding relational importances. Similarly, for a conference, the message passed by Paper-Conference relation and self-loop (Conference) connection are aggregated based on the weights of different relations. A detailed case study can be found in Sec. 6.8.

### 5.2 Meta-Path based Analysis

Besides of explaining RE-GNN from the perspective of spatial aggregations, RE-GNN can also be interpreted as a meta-path based heterogeneous GNN, implicitly. For convenience, Graph Transformer Network (GTN) [15] is selected as the compared heterogeneous GNN in this subsection. Our RE-GNN framework also employ GCN [1] as the backbone, similar to GTN. Here, we reveal that our RE-GCN can degenerate to GTN in certain conditions. Besides, we also theoretically demonstrate that RE-GCN possesses more expressive power than GTN.

#### 5.2.1 Simplifying RE-GCN

As stated in Sec. 3.4, GTN can be regarded as an ensemble of GCNs on multiple generated composite relations. By excluding the ensemble trick in Eq. (4), we consider a simple

case, where GTN only learns a 2-length composite relation, as below.

$$Z_P = \sigma\left(\tilde{A}_2 H^{(0)} W_P^{(0)}\right), \qquad (13)$$

where $\tilde{A}_2 = \hat{D}_2^{-1}\hat{A}_2$, $\hat{A}_2 = A_2 + I$, and $\hat{D}_2$ represents the degree matrix of $\hat{A}_2$. Note that $A_2 = \left(\sum_{i\in\mathcal{R}}\alpha_{1,i}A_i\right)\left(\sum_{i\in\mathcal{R}}\alpha_{2,i}A_i\right)$ is the learned adjacency matrix of 2-length meta-path. By stacking two RE-GCN layers, we can obtain

$$Z_H = \sigma\left(\tilde{A}_H^{(1)}\sigma\left(\tilde{A}_H^{(0)}H^{(0)}W_H^{(0)}\right)W_H^{(1)}\right), \qquad (14)$$

where $\tilde{A}_H^{(i)} = (\hat{D}_H^{(i)})^{-1}\hat{A}_H^{(i)}, i = 0, 1$, and $\hat{A}^{(i)}$ is calculated via Eq. (12). If we simplify Eq. (14) by removing the nonlinearity and coarsening the weight matrices (like SGC [29]), it will degenerate to

$$Z_{H,SGC} = \sigma\left(\tilde{A}_{H,2}H^{(0)}W_H\right), \qquad (15)$$

where $\tilde{A}_{H,2} = \hat{D}_{H,2}^{-1}\hat{A}_{H,2}$ and $\hat{A}_{H,2} = \hat{A}_H^{(1)}\hat{A}_H^{(0)}$. Note that $\hat{D}_{H,2} = \hat{D}_H^{(1)}\hat{D}_H^{(0)}$ is the degree matrix of $\hat{A}_{H,2}$. After simplifying 2 RE-GCN layers, we actually obtain a 2-length GTN layer. Note that the only difference between RE-GCN and GTN is $\hat{A}_H$ (RE-GCN) and $\hat{A}_P$ (GTN), which are generated from different formulas. $\hat{A}_H$ (RE-GCN) contains extra self-loop relation embeddings, while $\hat{A}_P$ (GTN) only adds an identity self-loop matrix. The above observation indicates that RE-GCN can implicitly learn the composite relations.

Since RE-SGC is a simplified version of RE-GCN, which will give superior expressive power than RE-SGC, we can intuitively conclude that the 2-layered RE-GCN in Eq. (14) possesses more expressive powers than the 1-layered GTN in Eq. (13). Then, we will prove it theoretically and extend it to the general case.

### 5.2.2 Theoretical Analysis

Due to limited computing and storing resources, the values of the input features and parameters are always finite for any neural network. Usually, to ensure the stability and robustness of a neural network, the input and parameters are normalized to small values. Under such circumstance, we firstly introduce the bounded set to quantitatively analyze the expressive power of neural network. Note that all the proofs are provided in the appendices.

**Definition 1** (Bounded set). *A bounded set is a set $\mathcal{D}$ where $\forall d \in \mathcal{D}, |d| < k$. $k > 0$ is a real number, which is the bound of $\mathcal{D}$. $\mathcal{D}$ can be a set of value, vector or matrix and $|\cdot|$ is a corresponding norm.*

As stated above, we assume that the dynamic ranges of the network inputs and parameters are both bounded sets. For example, a multilayer perceptron (MLP) layer is a function $f : \mathcal{H}_{in} \rightarrow \mathcal{H}_{out}$, which is defined on a bounded set $\mathcal{H}_{in} \subset \mathbb{R}^{d_{in}}$ and $f$'s range is also a bounded set $\mathcal{H}_{out} \subset \mathbb{R}^{d_{out}}$. $\forall h_{in} \in \mathcal{H}_{in}$ and $\forall h_{out} \in \mathcal{H}_{out}$, we can obtain $||h_{in}||_2^2 < k_{in}$ and $||h_{out}||_2^2 < k_{out}$, respectively. $\forall h_{in} \in \mathcal{H}_{in}$, $f$ is defined as

$$f(h_{in}) = \sigma(h_{in}W + b), \qquad (16)$$

where $W \in \mathcal{W} \subset \mathbb{R}^{d_{in}\times d_{out}}$ is the learnable weight matrix and $\sigma(\cdot)$ is a ReLU function. The bound of $W$ is defined as $k_w = |W| \stackrel{\text{def}}{=} \max||w_j||_2^2$, where $w_j$ is the column vector of

$W$. Note that $b$ is a bias vector and its bound is defined as $k_b = \max|b_i|$.

**Definition 2** (Bound of Multilayer Perceptron Layer). *A multilayer perceptron (MLP) layer $f(W, b) : \mathcal{H}_{in} \rightarrow \mathcal{H}_{out}$ is defined on a bounded set $\mathcal{H}_{in} \subset \mathbb{R}^{d_{in}}$. $k_{in}, k_w, k_b$ are the bounds of the input $h_{in}$, parameter $W$ and bias $b$, respectively. The bound of $f$ is defined as $k = \max(k_{in}, k_w, k_b)$.*

Since it is usually necessary to constrain the bound of input data when training or designing a neural network, we take the bound of the input into account for an MLP layer. According to the definition above, $k$ is also the bound of the input, parameters, and bias. Subsequently, we can infer the rationality of an MLP layer based on its bound. Here, we quantitatively analyze the expressive power of a composite of two MLP layers, compared to that of only one MLP layer.

**Lemma 3.**

*Given an MLP layer $f : \mathcal{H}_{in} \rightarrow \mathcal{H}_{out}$ bounded by $k$, there exists a composite of two MLP layers $(f_1 \circ f_2)$ which equals to $f$, where $f_1$ is also bounded by $k$ and $f_2$ is bounded by $\max\{1, 2k\}$.*

Since a particular solution of the layer $f_2$ is $W_2$ being an identity matrix when $b_2$ is bounded by $2k$, the layer $f_2$ is bounded by $\max(1, 2k)$ in Lemma 3. By considering that the bias plays a complementary role to parameter $W$, we do not separately restrict the bound of $b$, in practice. Thus, it is acceptable that the bias of $f_2$, i.e., $b_2$, has a bit larger bound.

According to Lemma 3, we can conclude that a composite of two MLP layers possesses no less expressive power than one MLP layer.

Now, let us consider the scenario of GNN. Similarly, for a GCN layer, we still utilize $k = \max(k_{in}, k_w, k_b)$ as its bound, where $k_{in}, k_w$ and $k_b$ are the bounds of the input, parameter and bias, respectively. In practice, GCN layer usually utilizes a bias vector. For example, the GTN layer in Eq. (13) can be rewritten as

$$Z_P = \sigma\left(\tilde{A}_2 HW_P + B_P\right), \qquad (17)$$

where $B = \begin{pmatrix} b \\ \vdots \\ b \end{pmatrix}$ is the *broadcast* of a bias vector $b$ (stacking multiple rows with a row vector $b$). The GTN layer actually adds a bias vector to each sample. Similarly, a composite of two RE-GCN layers in Eq. (14) can be rewritten as

$$Z_H = \sigma\left(\tilde{A}_H^{(1)}\sigma\left(\tilde{A}_H^{(0)}XW_H^{(0)} + B_H^{(0)}\right)W_H^{(1)} + B_H^{(1)}\right). \quad (18)$$

**Lemma 4.**

*If a vector set $\mathcal{H}$ is bounded by $k$. Denote $o = p_1 h_1 + p_2 h_2 + ... + p_r h_r$, where $\sum_{i=1}^r p_i = 1$, $h_i \in \mathcal{H}, i = 1, 2, ..., r$. Then, $o$ is also bounded by $k$, i.e., $||o||_2^2 < k$.*

Then, according to Lemma 3 and 4, Corollary 5 can be obtained.

**Corollary 5.**

*If $f_{GTN} : \mathcal{G} \rightarrow \mathbb{R}^{N\times C}$ is a 2-lengthed GTN layer, which is bounded by $\xi$, there exists a composite of two RE-GCN layers, where the first layer is bounded by $\xi$ and the second layer is bounded by $\max\{1, 2\xi\}$, which is equivalent to $f_{GTN}$.*

Corollary 5 theoretically reveals that a composite of two RE-GCN layers possesses no less expressive powers than a 2-lengthed GTN layer. Then, we can extend it to the general case.

**Theorem 6.**

TABLE 1: Statistics of the datasets.

| Datasets | Nodes | Edges |
|---|---|---|
| DBLP | # author(A): 4,057<br># paper(P): 14,328<br># term(T): 7,723<br># venue(V): 20 | # A-P:19,645<br># P-T:85,810<br># P-V:14,328 |
| ACM | # paper(P): 4,019<br># author(A): 7,167<br># subject(S): 60 | # P-P:9,615<br># P-A:13,407<br># P-S:4,019 |
| IMDB | # movie(M): 4,278<br># director(D): 2,081<br># actor(A): 5,257 | # M-D:4,278<br># M-A:12,828 |
| OGBN-MAG | # paper (P): 736,389<br># author (A): 1,134,649<br># institutions (I): 8,740<br># fields(F): 59,965 | # A-I: 1,043,998<br># A-P: 7,145,660<br># P-P: 7,505,078<br># P-F: 10,792,672 |

*Let $G \in \mathcal{G}$ be a heterogeneous graph, where the node features $X \in \mathcal{X}$ are normalized, i.e., $\forall X \in \mathcal{X}, |X| < \xi$. If a non-ensembled GTN, $m : \mathcal{G} \to \mathbb{R}^{N \times C}$, maps the nodes in $G$ to any node embeddings $Z \in \mathbb{R}^{N \times C}$, there exists a RE-GCN which is equivalent to GTN.*

Theorem 6 further demonstrates that RE-GCN possesses no less expressive power than the GTN. Then, Theorem 7 is obtained as follows.

**Theorem 7.**

*Let $G \in \mathcal{G}$ be a heterogeneous graph, where the node features $X \in \mathcal{X}$ are normalized, i.e., $\forall X \in \mathcal{X}, |X| < \xi$. There exists a RE-GCN, $r : \mathcal{G} \to \mathbb{R}^{N \times C}$, which can map $G$ to the node embeddings $Z \in \mathbb{R}^{N \times C}$ that GTN cannot map $G$ to.*

Theorems 6 and 7 jointly prove that our RE-GCN possesses more expressive power than the GTN. According to the above analysis, we can conclude that RE-GCN can be interpreted as an implicit meta-path based heterogeneous GNN, which indicates that it may be unnecessary to design or generate the meta-path explicitly. By simply stacking multiple RE-GNN layers, the heterogeneity may also be efficiently handled.

# 6 EVALUATIONS

## 6.1 Datasets

Four widely utilized heterogeneous datasets, including three heterogeneous academic networks (i.e., DBLP, ACM, and OGBN-MAG) and one heterogeneous movie network (i.e., IMDB), are employed to demonstrate the effectiveness of our RE-GNNs. Their details are shown in Tab. 1.

In DBLP, four types of nodes, i.e., papers, authors, venues, and terms, and three types of edges (relations), i.e., A-P, P-T, and P-V, are constructed. According to the conferences they are submitted to, authors are categorized into four research areas, i.e., Database, Data Mining, Artificial Intelligence and Information Retrieval. Each paper is described by a bag-of-words (BOW) representation of its keywords, and each author is described by the BOW embeddings of its published papers. Each term is represented by the Glove word vectors [40], and the attributes of venues are formed as one-hot vectors.

In ACM, three types of nodes, i.e., papers, authors, and subjects, and three types of directed edges (relations) are constructed. According to the conference they published,

papers are classified into three categories, i.e., Database, Wireless Communication, and Data Mining. Each paper is described by a BOW representation of keywords. The attributes of authors and subjects are formed as one-hot vectors.

In IMDB, three types of nodes, i.e., movies, directors, and actors, and two types of directed edges (relations) are formed. According to their genre, the movies are classified into three categories, i.e., Action, Comedy, and Drama. Movies are described by the BOW embeddings of their plots. The attributes of actors and directors are represented by the averaged attributes of their related movies.

In OGBN-MAG, four types of nodes, i.e., papers, authors, institutions, and fields of study, and four types of directed edges (relations), i.e., A-I, A-P, P-P, P-F, are constructed. According to the venue (conference or journal) they published, papers are divided into 349 categories. Each paper is associated with a 128-dimensional word2vec feature vector. We employ metapath2vec [41] to generate the features in other types of nodes.

We process the relations with directions for all the datasets and consider their reversed ones. For DBLP, ACM, and IMDB, we adopt the same data division strategies as [14], [42], where 400/400 nodes are randomly selected as training/validation sets, and the remaining nodes (approximately 80%) are selected as the test set. For OGBN-MAG, we adopt its official splits [43].

## 6.2 Baselines

We compare our method with various kinds of popular GNN methods, including five homogeneous GNNs (i.e., GCN [1] and GAT [2], GIN [24], GATv2 [39], GraphSAGE [3], and MixHop [44]), one heterogeneous graph embedding method (i.e., Metapath2vec [41]), four meta-path based heterogeneous GNNs (i.e., HAN [13], GTN [15], MAGNN [14], and MAGNN-AC [42]), and four relation based heterogeneous GNNs (i.e., R-GCN [17], HGT [20], HGB [23], and R-HGNN [21]). Note that the homogeneous GNNs, such as GCN and GAT, are constructed on the homogeneous graph, where the heterogeneity is directly eliminated. Here, GCN-M and GAT-M respectively represent the GCN and GAT models which are constructed on several meta-path based homogeneous graphs, and the best scores are reported.

## 6.3 Implementation Details

For the homogeneous GNNs, four-layered GNNs with 64 hidden neurons and 4 heads (only for GAT and GATv2), are employed as our baselines. Note that we employ GAT and GATv2 with 128 hidden neurons and 1 heads on IMDB, which are obtained via hyperparameter search. For the corresponding RE-GNNs, the gradient scaling factor $\lambda$ is set to 100. In the training process, both the homogeneous GNNs and RE-GNNs are trained for a maximum of 200 epochs with an early stopping condition at 50 epochs. The cross-entropy loss is utilized as the loss function. Adam [38] optimizer is employed with the learning rate of 0.001, the weight decay rate of 0.001 for ACM and DBLP, and 0.005 for IMDB. The dropout layers are utilized for the input of each GNN layer in the training process with a dropout rate of 0.6.

TABLE 2: Performances on the node classification task.

| Methods | DBLP | | ACM | | IMDB | |
|---|---|---|---|---|---|---|
| | Macro-F1 | Micro-F1 | Macro-F1 | Micro-F1 | Macro-F1 | Micro-F1 |
| Metapath2vec | 91.71 ± 0.00 | 92.39 ± 0.00 | 75.34 ± 0.00 | 76.79 ± 0.00 | 48.70 ± 0.00 | 50.40 ± 0.00 |
| HAN | 92.45 ± 0.87 | 92.96 ± 0.83 | 91.80 ± 0.38 | 91.67 ± 0.31 | 57.69 ± 0.91 | 58.15 ± 0.93 |
| GTN | 93.55 ± 0.22 | 94.29 ± 0.21 | 91.92 ± 0.67 | 91.81 ± 0.65 | 58.63 ± 1.45 | 60.19 ± 1.66 |
| MAGNN | 93.39 ± 0.33 | 93.89 ± 0.30 | 91.49 ± 1.18 | 91.46 ± 1.08 | 58.75 ± 2.04 | 59.48 ± 1.45 |
| MAGNN-AC | 93.00 ± 0.29 | 93.58 ± 0.21 | 88.34 ± 3.55 | 88.60 ± 2.98 | 56.98 ± 0.38 | 57.80 ± 0.44 |
| R-GCN | 92.26 ± 0.47 | 92.87 ± 0.42 | 93.36 ± 0.21 | 93.36 ± 0.21 | 58.96 ± 0.29 | 59.34 ± 0.27 |
| HGT | 91.43 ± 1.15 | 92.31 ± 0.91 | 92.40 ± 0.60 | 92.30 ± 0.54 | 58.16 ± 0.84 | 58.37 ± 0.70 |
| HGB | 93.20 ± 0.33 | 93.69 ± 0.29 | 93.15 ± 0.30 | 93.13 ± 0.27 | 58.11 ± 0.86 | 58.42 ± 0.71 |
| R-HGNN | 92.89 ± 0.59 | 93.43 ± 0.53 | 92.16 ± 0.71 | 92.10 ± 0.67 | 56.79 ± 1.01 | 57.27 ± 0.95 |
| GAT | 90.19 ± 0.89 | 90.94 ± 0.77 | 93.16 ± 0.16 | 93.16 ± 0.16 | 57.62 ± 1.20 | 58.92 ± 0.91 |
| GCN | 87.39 ± 0.23 | 88.31 ± 0.19 | 93.60 ± 0.29 | 93.58 ± 0.27 | 58.37 ± 1.26 | 59.59 ± 1.10 |
| GAT-M | 90.57 ± 0.70 | 91.15 ± 0.63 | 89.98 ± 0.52 | 89.94 ± 0.47 | 54.20 ± 0.84 | 54.80 ± 0.59 |
| GCN-M | 89.04 ± 0.71 | 89.99 ± 0.64 | 90.20 ± 0.24 | 90.15 ± 0.23 | 53.99 ± 0.95 | 54.72 ± 0.61 |
| RE-GAT | **95.06 ± 0.16** | **95.41 ± 0.15** | **94.04 ± 0.23** | **93.99 ± 0.22** | **60.01 ± 1.09** | **60.53 ± 0.85** |
| RE-GCN | **95.46 ± 0.29** | **95.80 ± 0.27** | 93.95 ± 0.16 | 93.93 ± 0.16 | **60.88 ± 0.95** | **61.51 ± 0.64** |

TABLE 3: Performances on the node clustering task.

| Methods | DBLP | | ACM | | IMDB | |
|---|---|---|---|---|---|---|
| | NMI | ARI | NMI | ARI | NMI | ARI |
| metapath2vec | 78.61 ± 0.04 | 83.64 ± 0.04 | 41.80 ± 0.01 | 34.71 ± 0.01 | 5.49 ± 0.24 | 5.16 ± 0.27 |
| HAN | 76.30 ± 0.68 | 82.12 ± 0.56 | 70.71 ± 0.91 | 75.16 ± 0.88 | 12.96 ± 1.09 | 14.34 ± 1.28 |
| MAGNN | 79.89 ± 0.85 | 84.94 ± 0.77 | 70.70 ± 1.08 | 74.34 ± 1.42 | 15.70 ± 0.77 | **15.73 ± 1.19** |
| R-GCN | 76.85 ± 1.70 | 82.73 ± 1.70 | 75.00 ± 0.80 | 77.82 ± 0.80 | 13.63 ± 0.16 | 15.59 ± 0.16 |
| HGB | 77.09 ± 1.93 | 82.86 ± 1.55 | 75.91 ± 0.84 | 80.56 ± 0.82 | 11.32 ± 2.48 | 11.40 ± 4.09 |
| R-HGNN | 75.16 ± 3.53 | 79.14 ± 6.78 | 63.86 ± 0.04 | 63.42 ± 0.08 | 9.67 ± 0.91 | 10.03 ± 1.72 |
| GCN | 60.59 ± 0.02 | 63.80 ± 0.03 | 77.28 ± 0.42 | 80.82 ± 0.56 | 14.15 ± 2.23 | 13.01 ± 2.43 |
| RE-GCN | **83.57 ± 1.78** | **88.13 ± 1.88** | **77.49 ± 0.58** | **81.67 ± 0.56** | **15.83 ± 1.10** | 15.65 ± 2.28 |

## 6.4 Node Classification

### 6.4.1 Comparisons

The node classification results of three heterogeneous datasets, i.e., DBLP, ACM, and IMDB, are presented in Tab. 2. For the employed baselines, both the meta-path based and relation based heterogeneous GNNs can achieve decent performances. Besides, though homogeneous GNNs, i.e., GCN and GAT, can perform well on ACM and IMDB datasets, they are not suitable for DBLP, compared to other heterogeneous GNNs. On the contrary, our RE-GNNs, which only utilize one-dimensional embeddings, can allow the homogeneous GNNs to process the heterogeneous graphs effectively. For example, in DBLP, RE-GCN performs much better than the original GCN, which has 8.07 and 7.49 improvements in Macro-F1 and Micro-F1 scores, respectively. A similar trend happens on GAT and RE-GAT. Meanwhile, for the ACM dataset, though our RE-GCN still achieves the best performance, the performance gain is relatively small, e.g., 0.35 in the Macro-F1 score. The reason may be that the weights of the relations with the current initialization are already close to the optimal weights in some particular cases. Under such circumstances, the performances of homogeneous GNNs may be close to or even approximately equal to the corresponding RE-GNNs. Generally, our RE-GCN and RE-GAT outperform the other well-designed heterogeneous GNNs, which usually possess complicated modules.
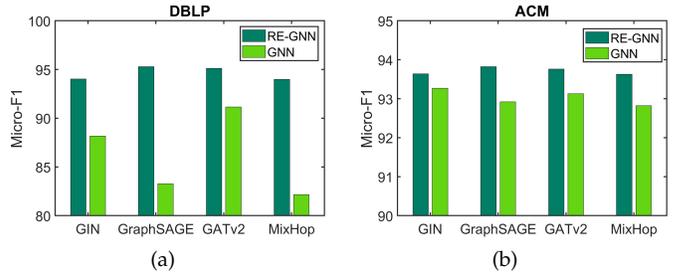


Fig. 3: Results of our RE-GNNs with other homogeneous GNNs as backbones.

### 6.4.2 Applied to other GNNs

The proposed framework can be easily applied to other homogeneous GNNs. Besides of GCN and GAT, four other typical homogeneous GNNs, i.e., GIN [24], GraphSAGE [3], GATv2 [39], and MixHop [44], are employed for further validations. As can be observed from Fig. 3, though these homogeneous GNNs cannot perform well on DBLP, their corresponding RE-GNNs can achieve superior yet decent results. On the ACM dataset, these homogeneous GNNs can achieve good results, similar to GCN and GAT. By introducing our relation embeddings, the performances of the corresponding RE-GNNs also possess certain improvements. In general, these results demonstrate that our RE-GNNs can effectively assign adequate abilities to the homogeneous

(a) Metapath2vec          (b) HAN          (c) MAGNN          (d) R-GCN

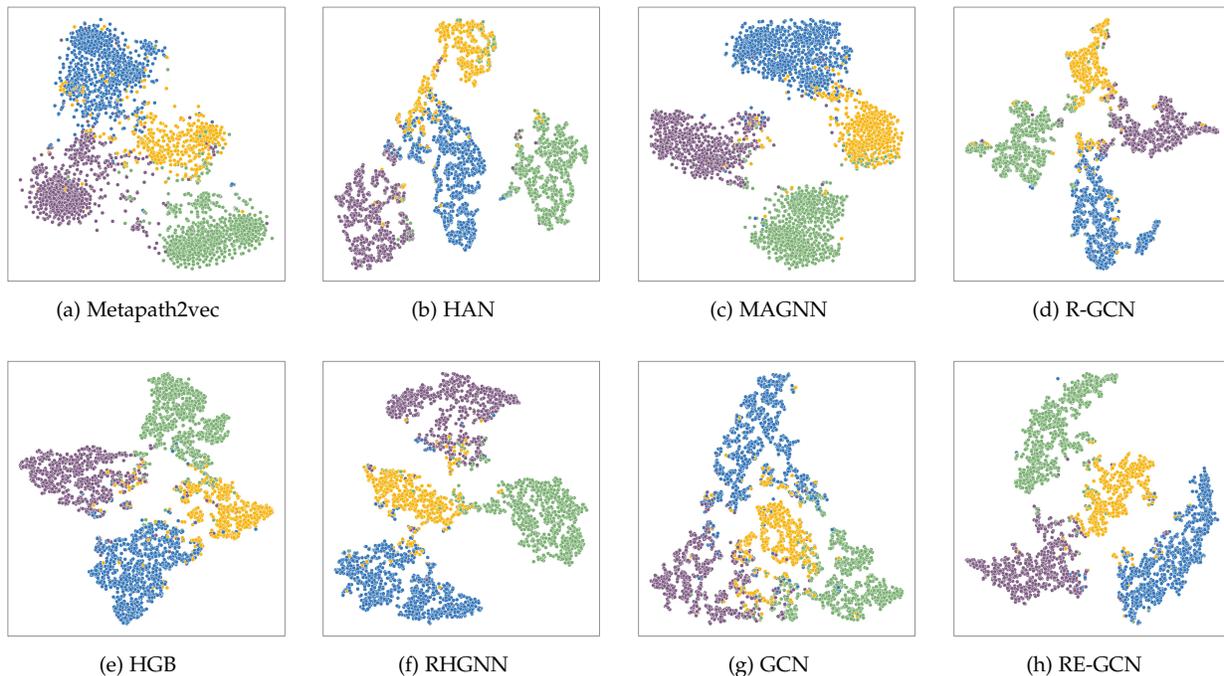(e) HGB          (f) RHGNN          (g) GCN          (h) RE-GCN

Fig. 4: Visualizations of the node representations on DBLP. Each data point represents an author, and its color represents its category.

GNNs to handle heterogeneous graphs.

## 6.5 Node Representation

### 6.5.1 Node Clustering

The node clustering task is conducted to validate the effectiveness of the learned node representations. The representations of testing nodes, which is generated by each trained model, are fed to the K-Means algorithm. Normalized Mutual Information (NMI) and Adjusted Rand Index (ARI) are employed as the evaluation metrics, where higher scores correspond to better models. By following [14], we repeat K-Means 10 times for each run of the models, and each model is tested for 10 runs. As shown in Tab. 3, our RE-GCN achieves superior performances on all the employed datasets, which demonstrates that our relation embeddings can enable GCN to learn meaningful node representations for heterogeneous graphs.

### 6.5.2 Visualizations

For visual comparisons, we utilize T-SNE [45] to project the learned embeddings of the authors in DBLP into a 2-dimensional space. As can be observed from Fig. 4g, for the visualization of GCN, authors with the same category tend to be decentralized, and authors with different categories tend to be mixed. On the contrary, after introducing the relation embeddings to model the heterogeneity, our RE-GCN can effectively learn suitable embeddings. As shown in Fig. 4h, the authors with identical research interests are more compact, and authors with different research interests are more distinguishable. Besides, the heterogeneous baseline methods either fail to gather the authors with

identical research interests or cannot provide clear boundaries for authors belonging to different categories. These results further validate that our RE-GNN can learn effective representations and thus handle the heterogeneous graphs.

## 6.6 Efficiency Analysis

As stated in Sec. 4, our RE-GNN only introduces one parameter for each relation type in each layer compared with the corresponding homogeneous GNN, which is an efficient framework to handle the heterogeneous graphs. Here, we analyze its experimental efficiency compared to the homogeneous GNNs. All the experiments are conducted on the Nvidia GeForce RTX 2080Ti GPU with 12 GB memory. Fig. 5 presents the training times and the model sizes for different methods trained on DBLP. For fair comparisons, all of these methods contain 64 hidden units for each hidden layer. An exception is R-HGNN, which possesses 16 hidden units (obtained via a hyper-parameter search).

As can be observed, our RE-GCN has an approximately equivalent model size to GCN, which is substantially smaller than the other heterogeneous GNNs. For example, HGB, a simple heterogeneous GNNs, possesses a 2.6x larger model size than our RE-GCN. Besides, as can be observed from Fig. 5b, our RE-GCN outperforms all the heterogeneous GNNs in terms of training speed. HGB possesses 2.9x training times compared to our RE-GCN. In addition, two complicated models, i.e., R-HGNN and MAGNN, require more than 17x training times. In general, compared to these heterogeneous GNNs, our RE-GCN provides the fastest training speed and the fewest model size with the best performances.
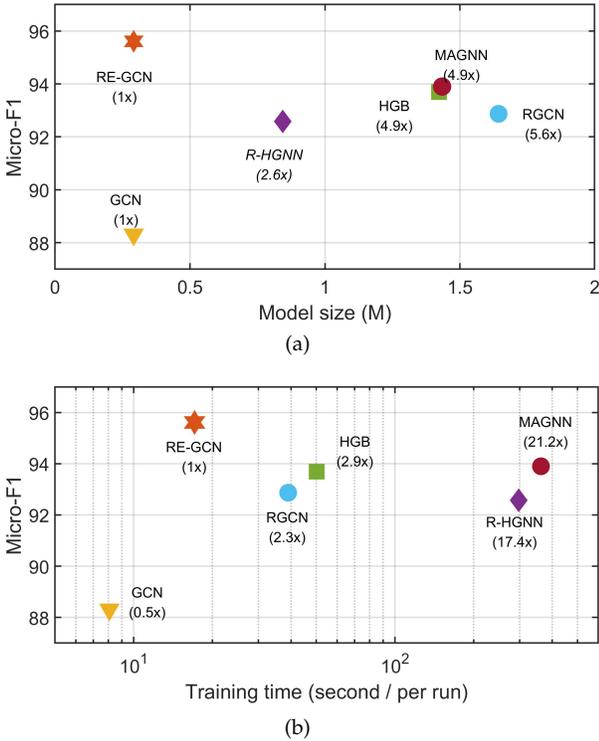
Fig. 5: Comparisons of the model size and average training time on DBLP. Note that all the model is trained with 64 hidden units except R-HGNN is trained with 16 hidden units.

TABLE 4: Performances on OGBN-MAG.

| | Validation Acc. | Test Acc. | Parameters |
|---|---|---|---|
| MetaPath2vec | 35.06 ± 0.17 | 35.44 ± 0.36 | 94,479,069 |
| HGT | 49.89 ± 0.47 | 49.27 ± 0.61 | 21,173,389 |
| R-GCN | 49.12 ± 0.50 | 47.89 ± 0.53 | 154,366,772 |
| R-GSN | 51.33 ± 0.35 | 50.10 ± 0.42 | 154,373,028 |
| R-HGNN | 53.61 ± 0.22 | 52.04 ± 0.26 | 5,638,053 |
| RE-GCN | 51.94 ± 0.12 | 50.82 ± 0.18 | **1,037,171** |
| RE-GAT | 53.08 ± 0.26 | **52.10 ± 0.17** | 1,039,373 |

### 6.7 Scalability for Large-scale Graph

Here, we employ the OGBN-MAG dataset to validate the scalability of our RE-GNN for handling large-scale heterogeneous graphs. Note that many training techniques can boost the performance on OGBN-MAG, such as multi-stage training [46] and adversarial training [47], which are not considered in our comparison, because these techniques are orthogonal to model design. The results of employed baselines in Tab. 4, except for R-GSN, are borrowed from the OGB leaderboards [43]. For R-GSN, we report its official results when removing the FALG training technique [47]. For utilizing both GCN and GAT as backbones, we construct 2-layered RE-GNNs followed by a fully connected layer as the output layer. The hidden dimension is set to 512 for RE-GCN and 64 for RE-GAT with 8 heads. We utilize the Neighbor Sampling method [3] to train our RE-GNNs with mini-batch.

As can be observed from Tab. 4, both RE-GCN and RE-GAT can achieve outstanding performances on OGBN-

TABLE 5: Ablation study on the DBLP dataset. Note that *ET emb* represents the edge type relation embeddings, and *SL emb* represents the node-type-specific self-loop embeddings.

| | ET emb | SL emb | Macro-F1 | Micro-F1 |
|---|---|---|---|---|
| GAT | | | 90.19 | 90.94 |
| GCN | | | 87.39 | 88.31 |
| GAT-S | | ✓ | 90.63 | 91.29 |
| GCN-S | | ✓ | 88.77 | 89.43 |
| GAT-E | ✓ | | 94.91 | 95.28 |
| GCN-E | ✓ | | 95.15 | 95.48 |
| RE-GAT | ✓ | ✓ | 95.06 | 95.41 |
| RE-GCN | ✓ | ✓ | **95.46** | **95.80** |

MAG. When utilizing GAT as the backbone, our RE-GAT achieves the best performance and outperforms the complicated heterogeneous GNNs, e.g., R-GSN and R-HGNN. Besides, when utilizing GCN as the backbone, our RE-GCN can still outperform R-GCN. Meanwhile, our RE-GCN and RE-GAT possess much fewer number of parameters than other heterogeneous methods. These results demonstrate our RE-GNNs are still effective and efficient on the large-scale heterogeneous graphs.

### 6.8 Case Study

Here, we analyze the learned embeddings of a 4-layered RE-GCN on the DBLP dataset to further study the mechanism of our RE-GNN framework. As presented in Tab. 1, DBLP contains four types of nodes and three types of edges. The task is designed to predict the category of authors. By considering the self-loop connections and reversing the directed relations, our RE-GNN considers 10 types of relations.

As can be observed from Fig. 6, the importances of P-* (i.e., P-A, P-T, P-V, and P) relations are much larger than the others in the first layer. It indicates that RE-GCN believes the paper attributes to be more critical and utilizes them to initialize the features in other types of nodes. In the second layer, nodes exchange messages according to different types of relations. Then, in the third layer, there are two groups of dominate messages: 1) P-A; 2) *-P (i.e., A-P, T-P, V-P, and P). It indicates that the third RE-GCN layer utilizes the messages from papers to update the states of authors, as well as utilizes the messages from other types of nodes to update the states of papers. In the last layer, the states of authors are updated by the messages from paper nodes and itself. In general, we can conclude that our RE-GCN can adaptively make the paper nodes to dominate in this case, which is consistent with the common understandings of academic networks.

### 6.9 Ablation Study

Here, we verify the effectiveness of the edge relation embeddings and node-type-specific self-loop embeddings. Two variants of our RE-GNN are given: 1) GNN-S: homogeneous GNN with node-type-specific self-loop embeddings; 2) GNN-E: homogeneous GNN with edge relation embeddings.

As shown in Tab. 5, the F1 scores of the GNN-E variants are significantly higher than these of corresponding GNNs. Meanwhile, the node-type-specific self-loop embeddings

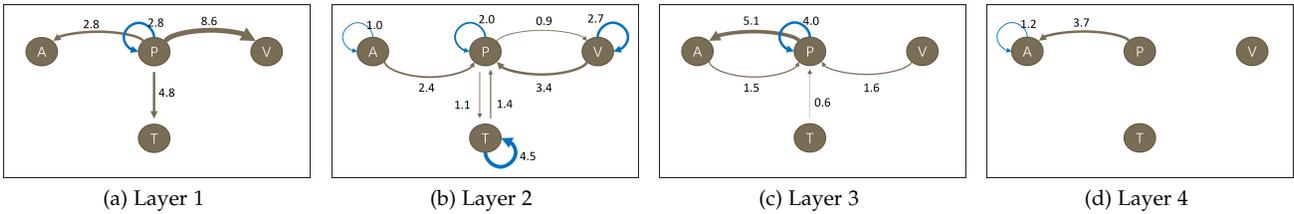| (a) Layer 1 | (b) Layer 2 | (c) Layer 3 | (d) Layer 4 |

Fig. 6: Visualizations of the learned relation embeddings for a 4-layered RE-GCN on the DBLP dataset. Note that we remove the relation when its weight is lower than 0.4, for a clear visualization.



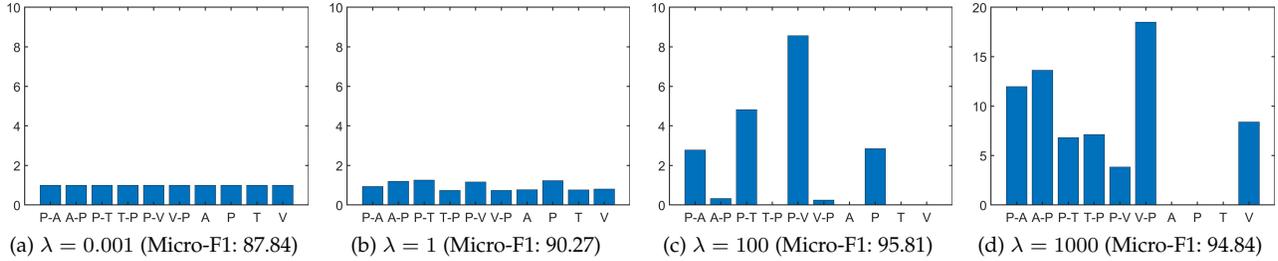| (a) $\lambda = 0.001$ (Micro-F1: 87.84) | (b) $\lambda = 1$ (Micro-F1: 90.27) | (c) $\lambda = 100$ (Micro-F1: 95.81) | (d) $\lambda = 1000$ (Micro-F1: 94.84) |

Fig. 7: The learned weights of the relations for the first layer of RE-GCN on the DBLP dataset.

give relatively minor improvements. It indicates that the relation types are more critical than the node-type-specific self-loop types. Note that the GNN-E variants can also implicitly exploit the general self-loop embeddings, because the normalizations in the aggregation step of GNNs will jointly normalize the message from different types of relations and self-loops. Besides, RE-GNNs, which explicitly learn the node-type-specific self-loop embeddings, outperform their GNN-E variants. Therefore, the effectiveness of our node-type-specific self-loop embeddings has also been verified.

### 6.10 Impacts of Gradient Scaling Factor

As stated in Sec. 4.3, the update of the weights of the relations in each optimization iteration can be scaled via the gradient scaling factor $\lambda$, theoretically. By using the Adam optimizer, the update of the weights of the relations changes $\lambda$ times. Here, we experimentally verify the effectiveness of the proposed gradient scaling factor. Fig. 7 presents the results of the first layer of RE-GCN, which is learned on the DBLP dataset. When $\lambda$ is set to 0.001, the update of the weights of the relations is negligible. Under such circumstance, RE-GCN becomes an approximation of GCN with the fixed weights of the relations being ones. When $\lambda$ is set to 1.0, RE-GCN possesses the original relation embeddings in Eqs. (6) and (7). As shown in Fig. 7b, the final weights of the relations are similar. When we set a proper scaling factor, e.g., $\lambda = 100$, RE-GCN can learn distinguishable weights of the relations. However, when $\lambda$ is too large, e.g., 1000, the weights of the relations are very sensitive in the training process and are easy to fall into polarization, i.e., part of the relations are dominating the others. According to the above results, our gradient scaling factor can help RE-GNN to learn proper relation embeddings and process the heterogeneous graphs.

## 7 CONCLUSION

This paper proposes a simple yet effective framework, named Relation Embedding based Graph Neural Network (RE-GNN), to assign adequate ability to the homogeneous GNNs for handling the heterogeneous graphs. Specifically, we exploit only one parameter per relation to model the importance of distinct types of relations and node-type-specific self-loop connections. To optimize the relation embeddings and the model parameters simultaneously, a gradient scaling factor is proposed to enable the embeddings to converge to appropriate values. Besides, we interpret the proposed RE-GNN from two perspectives, and theoretically demonstrate that our RE-GCN possesses more expressive power than GTN. Extensive experiments demonstrate that our RE-GNN can effectively and efficiently handle the heterogeneous graphs and can be applied to various homogeneous GNNs.

## REFERENCES

[1] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *ICLR*, 2017.

[2] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in *ICLR*, 2018.

[3] W. L. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *NIPS*, 2017, pp. 1024–1034.

[4] J. Wang, Y. Wang, Z. Yang, L. Yang, and Y. Guo, "Bi-gcn: Binary graph convolutional network," in *IEEE CVPR*, 2021, pp. 1561–1570.

[5] B. Jiang, Z. Zhang, D. Lin, J. Tang, and B. Luo, "Semi-supervised learning with graph learning-convolutional networks," in *IEEE CVPR*, 2019, pp. 11 313–11 320.

[6] J. Qiu, J. Tang, H. Ma, Y. Dong, K. Wang, and J. Tang, "Deepinf: Social influence prediction with deep learning," in *ACM SIGKDD*, 2018, pp. 2110–2119.

[7] C. Li and D. Goldwasser, "Encoding social information with graph convolutional networks for political perspective detection in news media," in *ACL*, 2019, pp. 2594–2604.

[8] A. Fout, J. Byrd, B. Shariat, and A. Ben-Hur, "Protein interface prediction using graph convolutional networks," in *NIPS*, 2017, pp. 6533–6542.

[9] J. Shang, C. Xiao, T. Ma, H. Li, and J. Sun, "Gamenet: Graph augmented memory networks for recommending medication combination," in *AAAI*, 2019, pp. 1126–1133.

[10] S. Guo, Y. Lin, N. Feng, C. Song, and H. Wan, "Attention based spatial-temporal graph convolutional networks for traffic flow forecasting," in *AAAI*, 2019, pp. 922–929.

[11] J. Li, Z. Han, H. Cheng, J. Su, P. Wang, J. Zhang, and L. Pan, "Predicting path failure in time-evolving graphs," in *ACM SIGKDD*, 2019, pp. 1279–1289.

[12] Z. Xiong, D. Wang, X. Liu, F. Zhong, X. Wan, X. Li, Z. Li, X. Luo, K. Chen, H. Jiang *et al.*, "Pushing the boundaries of molecular representation for drug discovery with the graph attention mechanism," *Journal of medicinal chemistry*, vol. 63, no. 16, pp. 8749–8760, 2019.

[13] X. Wang, H. Ji, C. Shi, B. Wang, Y. Ye, P. Cui, and P. S. Yu, "Heterogeneous graph attention network," in *WWW*, 2019, pp. 2022–2032.

[14] X. Fu, J. Zhang, Z. Meng, and I. King, "Magnn: Metapath aggregated graph neural network for heterogeneous graph embedding," in *WWW*, 2020, pp. 2331–2341.

[15] S. Yun, M. Jeong, R. Kim, J. Kang, and H. J. Kim, "Graph transformer networks," in *NIPS*, 2019, pp. 11 983–11 993.

[16] Y. Ding, Q. Yao, H. Zhao, and T. Zhang, "Diffmg: Differentiable meta graph search for heterogeneous graph neural networks," in *ACM SIGKDD*, 2021, pp. 279–288.

[17] M. Schlichtkrull, T. N. Kipf, P. Bloem, R. Van Den Berg, I. Titov, and M. Welling, "Modeling relational data with graph convolutional networks," in *ESWC*, 2018, pp. 593–607.

[18] X. Wu, M. Jiang, and G. Liu, "R-gsn: The relation-based graph similar network for heterogeneous graph," *arXiv preprint arXiv:2103.07877*, 2021.

[19] C. Zhang, D. Song, C. Huang, A. Swami, and N. V. Chawla, "Heterogeneous graph neural network," in *ACM SIGKDD*, 2019, pp. 793–803.

[20] Z. Hu, Y. Dong, K. Wang, and Y. Sun, "Heterogeneous graph transformer," in *WWW*, 2020, pp. 2704–2710.

[21] L. Yu, L. Sun, B. Du, C. Liu, W. Lv, and H. Xiong, "Heterogeneous graph representation learning with relation awareness," *IEEE Transactions on Knowledge and Data Engineering*, 2022.

[22] S. Zhu, C. Zhou, S. Pan, X. Zhu, and B. Wang, "Relation structure-aware heterogeneous graph neural network," in *ICDM*, 2019, pp. 1534–1539.

[23] Q. Lv, M. Ding, Q. Liu, Y. Chen, W. Feng, S. He, C. Zhou, J. Jiang, Y. Dong, and J. Tang, "Are we really making much progress?: Revisiting, benchmarking and refining heterogeneous graph neural networks," in *ACM SIGKDD*, 2021, pp. 1150–1160.

[24] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How powerful are graph neural networks?" in *ICLR*, 2019.

[25] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," *IEEE Transactions on Neural Networks*, vol. 20, no. 1, pp. 61–80, 2008.

[26] D. K. Hammond, P. Vandergheynst, and R. Gribonval, "Wavelets on graphs via spectral graph theory," *Applied and Computational Harmonic Analysis*, pp. 129–150, 2011.

[27] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, "Spectral networks and locally connected networks on graphs," in *ICLR*, 2014.

[28] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *NIPS*, 2016, pp. 253–261.

[29] F. Wu, A. H. S. Jr., T. Zhang, C. Fifty, T. Yu, and K. Q. Weinberger, "Simplifying graph convolutional networks," in *ICML*, 2019, pp. 6861–6871.

[30] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, "Neural message passing for quantum chemistry," in *ICML*, 2017, pp. 1263–1272.

[31] J. Chen, T. Ma, and C. Xiao, "Fastgcn: Fast learning with graph convolutional networks via importance sampling," in *ICLR*, 2018.

[32] T. Elsken, J. H. Metzen, and F. Hutter, "Neural architecture search: A survey," *JMLR*, vol. 20, no. 1, pp. 1997–2017, 2019.

[33] Y. Sun and J. Han, "Mining heterogeneous information networks: principles and methodologies," *Synthesis Lectures on Data Mining and Knowledge Discovery*, vol. 3, pp. 1–159, 2012.

[34] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *AISTATS*, 2010, pp. 249–256.

[35] N. Qian, "On the momentum term in gradient descent learning algorithms," *Neural networks*, vol. 12, no. 1, pp. 145–151, 1999.

[36] B. He and X. Yuan, "On the o(1/n) convergence rate of the douglas–rachford alternating direction method," *SIAM Journal on Numerical Analysis*, vol. 50, no. 2, pp. 700–709, 2012.

[37] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization." *JMLR*, vol. 12, no. 7, 2011.

[38] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR*, 2015.

[39] S. Brody, U. Alon, and E. Yahav, "How attentive are graph attention networks?" in *ICLR*, 2022.

[40] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *EMNLP*, 2014, pp. 1532–1543.

[41] Y. Dong, N. V. Chawla, and A. Swami, "metapath2vec: Scalable representation learning for heterogeneous networks," in *ACM SIGKDD*, 2017, pp. 135–144.

[42] D. Jin, C. Huo, C. Liang, and L. Yang, "Heterogeneous graph neural network via attribute completion," in *WWW*, 2021, pp. 391–400.

[43] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, and J. Leskovec, "Open graph benchmark: Datasets for machine learning on graphs," in *NIPS*, 2020.

[44] S. Abu-El-Haija, B. Perozzi, A. Kapoor, N. Alipourfard, K. Lerman, H. Harutyunyan, G. V. Steeg, and A. Galstyan, "Mixhop: Higher-order graph convolutional architectures via sparsified neighborhood mixing," in *ICML*, 2019, pp. 21–29.

[45] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne." *JMLR*, vol. 9, no. 11, 2008.

[46] Q. Li, Z. Han, and X. Wu, "Deeper insights into graph convolutional networks for semi-supervised learning," in *AAAI*, 2018, pp. 3538–3545.

[47] K. Kong, G. Li, M. Ding, Z. Wu, C. Zhu, B. Ghanem, G. Taylor, and T. Goldstein, "Flag: Adversarial data augmentation for graph neural networks," *arXiv preprint arXiv:2010.09891*, 2020.

# APPENDIX A
## THEORETICAL IMPACTS OF GRADIENT SCALING FACTOR

Here, we discuss the theoretical impacts of the proposed gradients scaling factor, i.e., $\alpha_i = \lambda e_i$, where the relation embedding $e_i$ is scaled by a factor $\lambda > 0$. The gradient of object function $\mathcal{L}$ to the weight of each relation $\alpha_i$ is $\frac{\partial \mathcal{L}}{\partial \alpha_i}$. Then the gradient to each relation embedding $e_i$ is $\lambda \frac{\partial \mathcal{L}}{\partial \alpha_i}$. In the original case, where the relation embedding is utilized directly as the weight of the relation ($\alpha_i = e_i$), the gradient to $e_i$ is denoted as $g_i$. By employing the gradient scaling factor, the scaled gradient to $e_i$ is denoted as $g_i' = \lambda g_i$.

Now, we consider the updating modification for each iteration in the optimization process. For the gradient descent based optimizers, each relation embedding is updated via $e_i^{next} = e_i - \Delta e_i$, where $\Delta e_i = \kappa(g_i)$ is correlated to the gradient $g_i$. With the scaling factor $\lambda$, the updating modification becomes $\kappa(\lambda g_i)$.

We argue that $\kappa(\lambda g_i) = \lambda \kappa(g_i)$, for the common gradient optimizers such as SGD, Momentum [35], and Nesterov Momentum [36]. For the adaptive optimizers like Adagrad [37] and Adam [38], $\kappa(\lambda g_i) = \kappa(g_i)$. Here, we present the detailed proofs of SGD and Adam optimizers as examples. The proofs of other commonly utilized optimizers can be generalized by these proofs.

**Stochastic Gradient Descent (SGD).** For a Stochastic Gradient Descent (SGD) optimizer, a parameter $\theta$ is updated by

$$\theta_t = \theta_{t-1} - \eta g_t, \tag{S1}$$

where $\eta$ stands for the learning rate and $g_t = \frac{\partial \mathcal{L}}{\partial \theta_{t-1}}$ represents the gradient of a batch of the input data. Thus, the updating modification function is $\kappa(g) = -\eta g$. Then, $\kappa(g') = \kappa(\lambda g) = -\eta \lambda g = \lambda(-\eta g) = \lambda \kappa(g)$.

**Adaptive Moment Estimation (Adam).** Adaptive Moment Estimation (Adam) is a typical gradient descent method which computes an adaptive learning rate for each parameter. In each iteration $t$, it firstly computes the exponential averages of the gradient $m_t$ and the squared gradient $v_t$ respectively as

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t, \tag{S2}$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2. \tag{S3}$$

where $\beta_1$ and $\beta_2$ are the two pre-defined hyperparameters. Both $m_t$ and $v_t$ are initialized as vectors of zeros. Then, the bias elimination is utilized to obtain

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \tag{S4}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}. \tag{S5}$$

At last, the parameter is updated accordingly.

$$\theta_t = \theta_{t-1} - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}, \tag{S6}$$

where $\epsilon > 0$ is an extremely small number which can be neglected when $v_t \neq 0$. Then, the updating modification function is $\kappa(g_t) = -\eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$. In the following proof, induction is employed to prove $\kappa(\lambda g_t) = \kappa(g_t)$.

*Proof.* Initially, we have $m_0 = 0$ and $v_0 = 0$. Let $g_t' = \lambda g_t$. Then, $\kappa(\lambda g) = \kappa(g')$.

In the first iteration, since $m_0 = v_0 = 0$, it is easy to validate $m_1' = \lambda m_1$, $v_1' = \lambda^2 v_1$, $\hat{m}_1' = \lambda \hat{m}_1$ and $\hat{v}_1' = \lambda^2 \hat{v}_1$. Then, the updating function is $\kappa(\lambda g_t) = -\eta \frac{\lambda \hat{m}_1}{\sqrt{\lambda^2 \hat{v}_1} + \epsilon} = \kappa(g_t)$. Since $\epsilon$ is a very small number, which is only utilized to prevent a zero denominator in practice, we ignore its effect here.

Assuming that in the $(t-1)$-th iteration, we have obtained $\hat{m}_{t-1}' = \lambda \hat{m}_{t-1}$, $\hat{v}_{t-1}' = \lambda^2 \hat{v}_{t-1}$ and $\kappa(g_{t-1}') = \kappa(g_{t-1})$. Then, in the $t$-th iteration, we can obtain

$$\begin{aligned} m_t' &= \beta_1 m_{t-1}' + (1 - \beta_1) g_t' \\ &= \beta_1 \lambda m_{t-1} + (1 - \beta_1) \lambda g_t \\ &= \lambda(\beta_1 m_{t-1} + (1 - \beta_1) g_t) \\ &= \lambda m_t. \end{aligned} \tag{S7}$$

Similarly, we can obtain

$$\begin{aligned} v_t' &= \beta_2 v_{t-1}' + (1 - \beta_2) g_t'^2 \\ &= \beta_2 \lambda^2 v_{t-1} + (1 - \beta_2) \lambda^2 g_t^2 \\ &= \lambda^2(\beta_2 m_{t-1} + (1 - \beta_2) g_t^2) \\ &= \lambda^2 v_t. \end{aligned} \tag{S8}$$

Then, for the bias eliminated $\hat{m}_t'$ and $\hat{v}_t'$, we can compute $\hat{m}_t' = \lambda \hat{m}_t$ and $\hat{v}_t' = \lambda^2 \hat{v}_t$ respectively. At last, the updating modification is $\kappa(\lambda g_t) = -\eta \frac{\hat{m}_1'}{\sqrt{\hat{v}_1'} + \epsilon} = \kappa(g_t)$. $\square$

Then, we consider the modification for the weight of each relation $\Delta \alpha_i$. Since, $\alpha_i$ is not a parameter, its modification is correlated to $\Delta e_i$. In the original case, $\Delta \alpha = \Delta e_i = \kappa(g_i)$. By employing the gradient scaling factor, the modification of the weight of each relation is $\Delta \alpha' = \lambda \Delta e_i = \lambda \kappa(g_i')$. As stated above, for the common gradient optimizers, such as SGD, Momentum and Nesterov, $\Delta \alpha' = \lambda^2 \Delta \alpha$. For the adaptive optimizers like Adagrad and Adam, $\Delta \alpha' = \lambda \Delta \alpha$.

# APPENDIX B
## PROOF OF LEMMA 3

*Proof.* $\forall h \in \mathcal{H}_{in}$, we can obtain

$$(f_1 \circ f_2)(h) = \sigma(\sigma(hW_1 + b_1)W_2 + b_2), \tag{S10}$$

where $\sigma$ is a ReLU function, $h$, $W_1$, and $W_2$ are bounded by $k$, $k_{w_1}$, $k_{w_2}$, respectively. By letting $t = hW_1$, we can obtain $t_i = hw_i^T$. Since $hh^T = \sum_j h_j^2 < k$ and $w_i w_i^T = \sum_j w_{ij}^2 < k_{w_1}$, we can calculate $|t_i| = |hw_i^T| = |\sum_j h_j w_{ij}| \leq \sum_j |h_j w_{ij}| \leq \frac{1}{2} \sum_j (h_j^2 + w_{ij}^2) = \frac{k + k_{w_1}}{2}$. Then, there is a $|b_i| \leq \frac{k + k_{w_1}}{2}$, when $b_i + t_i \geq 0$, i.e., there exists a $b_1'$ with a bound of $\frac{k + k_1}{2}$, when $hW_1 + b_1' > 0$. Thus, Eq. (S10) can be rewritten as

$$(f_1 \circ f_2)(h) = \sigma(hW_1 W_2 + b_1' W_2 + b_2). \tag{S11}$$

Then, when $W_1 W_2 = W$ and $b_1' W_2 + b_2 = b$, $(f_1 \circ f_2)$ equals to $f$. There exists a simple solution to the above formulas, i.e., $W_1 = W, W_2 = I, b_2 = b - b_1'$.

$\square$

## APPENDIX C
## PROOF OF LEMMA 4

*Proof.* According to the definition, the bound of $o$ is

$$||o||_2^2 = \sum_j o_j^2 = \sum_j \left( \sum_i p_i h_{ij} \right)^2. \qquad \text{(S12)}$$

Since $\sum_{i=1}^r p_i = 1$, and $(\cdot)^2$ is a convex function,

$$\sum_j \left( \sum_i p_i h_{ij} \right)^2 \leq \sum_j \left( \sum_i p_i h_{ij}^2 \right), \qquad \text{(S13)}$$

according to the Jensen's inequality. Then,

$$
\begin{aligned}
\sum_j \left( \sum_i p_i h_{ij}^2 \right) &= \sum_i p_i \left( \sum_j h_{ij}^2 \right) \\
&< \sum_i p_i k \\
&= k.
\end{aligned}
\qquad \text{(S14)}
$$

At last, we can obtain $||o||_2^2 < k$. □

## APPENDIX D
## PROOF OF COROLLARY 5

*Proof.* For Eq. (14), X is the collection of $x$. Each row of $X$, which is denoted as $x_i$, is bounded by $\xi$. $\tilde{A} = \hat{D}^{-1}\hat{A}$ is a non-negative matrix, where the summation of each row equals to 1, i.e., $\sum_j \tilde{a}_{ij} = 1$ and $\tilde{a}_{ij} \geq 0$. Let $H = \tilde{A}X$. According to Lemma 4, the bound of each row $h_i$ is $\xi$.

Then, according to Lemma 3, we can conclude that there exists a $b'$ with a bound of $\frac{k+k_1}{2}$, $hW_1 + b'_1 > 0$, i.e., $HW_1 + B'_1 > 0$. Thus, Eq. (14) can be rewritten as

$$Z_H = \sigma(\tilde{A}_H^{(1)} \tilde{A}_H^{(0)} X W_H^{(0)} W_H^{(1)} + \tilde{A}_H^{(1)} B_H^{(0)'} W_H^{(1)} + B_H^{(1)}). \qquad \text{(S17)}$$

Since $\tilde{A}_H^{(1)}$ is a matrix, which is normalized in each row, and $B_H^{(0)'}$ is a column equivalent matrix, we can obtain

$$\tilde{A}_H^{(1)} B_H^{(0)'} = B_H^{(0)'}. \qquad \text{(S18)}$$

Besides,

$$
\begin{aligned}
\tilde{A}_H^{(1)} \tilde{A}_H^{(0)} &= \hat{D}^{(1)} \hat{A}_H^{(1)} \hat{D}^{(0)} \hat{A}_H^{(0)} \\
&= \hat{D}^{(1)} \hat{D}^{(0)} \hat{A}_H^{(1)} \hat{A}_H^{(0)} \\
&= \hat{D}_{H2} \hat{A}_{H2} \\
&= \tilde{A}_{H2}.
\end{aligned}
\qquad \text{(S19)}
$$

With Eqs. (S18) and (S19), Eq. (S17) can be reformed to

$$Z_H = \sigma \left( \tilde{A}_{H2} X W_H^{(0)} W_H^{(1)} + B_H^{(0)'} W_H^{(1)} + B_H^{(1)} \right). \qquad \text{(S20)}$$

Similar to the proof of Lemma 3, we can easily construct a solution that $W_H^{(0)} = W_P$, $W_H^{(1)} = I$, $B_H^{(1)} = B_P - B_H^{(0)'}$. □

## APPENDIX E
## PROOF OF THEOREM 6

*Proof.* Theorem 1 can be proved in two steps.

- T1(1). For any $L$-lengthed one-layered GTN, there exists an $L$-layered RE-GCN which is equivalent to it.
- T1(2). For any $L$-lengthed $K$-layered GTN, there exists an $(LK)$-layered RE-GCN which is equivalent to it, where $K > 1$.

According to Corollary 5, for an $L$-lengthed GTN layer, we can also obtain a stack of $L$ RE-GCN layers which is equivalent to it. This equivalency can be achieved via removing the ReLU function and choosing a proper bias vector $b^{(L)}$ in each layer (except the last layer). Note that the bias vector $b^{(l)}$ is bounded by $\frac{k_{in}^{(l)}+k_w^{(l)}}{2}$. $k_{in}^{(l)}$ and $k_w^{(l)}$ are the bounds of the input features and parameters in the $L$-th layer, respectively. Therefore, T1(1) is proved.

For T1(2), we can employ a composite of $L$ RE-GCN layers which is equivalent to each $L$-lengthed GTN layer. Then, we can stack these $(KL)$ RE-GCN layers, which is equivalent to the $L$-lengthed $K$-layered GTN. Therefore, T1(2) is proved. □

## APPENDIX F
## PROOF OF THEOREM 7

*Proof.* Theorem 2 can be proved in two steps.

- T2(1). There exists a 2-layered RE-GCN, which cannot be represented by any $L$-lengthed one-layered GTN.
- T2(2). There exists a 2-layered RE-GCN, which cannot be represented by any $L$-lengthed $K$-layered GTN, where $K > 1$.

Consider a graph with only one node and a self-loop connection. Then, a GCN layer is degenerated to an MLP layer. Since the composite of two MLP layers can be a non-linear mapping while one MLP layer can only be a linear mapping (by ignoring the last ReLU function), there exists a composite of two MLP layers that one MLP layer cannot be equivalent to. Therefore, T2(1) is proved.

The adjacency matrix enables the nodes (samples) to exchange messages with others. Then, the effects of neighbourhood aggregation with $A$ cannot be replaced by the effects of weight projection with $W$. For a $K$-layered GTN ($K > 1$) with an arbitrary learned adjacency matrix $A_P$, it should satisfy that $A_P^{k_1} = A_{H_1}$, $A_P^{k_2} = A_{H_2}$, where $k_1 + k_2 = K$. However, this cannot be satisfied at all the time for any possible $A_{H_1}$ and $A_{H_2}$. For example, if $A_{H_1}$ is a non-singular matrix, while $A_{H_2}$ is a singular matrix, no proper $A_P$ can be obtained. The reason is that $A^n$ is a non-singular matrix, if and only if the matrix $A$ is non-singular. Therefore, T2(1) is proved.

□