

Decouple Graph Neural Networks: Train Multiple Simple GNNs Simultaneously Instead of One

Hongyuan Zhang, Yanan Zhu, and Xuelong Li*, *Fellow, IEEE*

Abstract—Graph neural networks (GNN) suffer from severe inefficiency due to the exponential growth of node dependency with the increase of layers. It extremely limits the application of stochastic optimization algorithms so that the training of GNN is usually time-consuming. To address this problem, we propose to decouple a multi-layer GNN as multiple simple modules for more efficient training, which is comprised of classical *forward training* (FT) and designed *backward training* (BT). Under the proposed framework, each module can be trained efficiently in FT by stochastic algorithms without distortion of graph information owing to its simplicity. To avoid the only unidirectional information delivery of FT and sufficiently train shallow modules with the deeper ones, we develop a backward training mechanism that makes the former modules perceive the latter modules, inspired by the classical backward propagation algorithm. The backward training introduces the reversed information delivery into the decoupled modules as well as the forward information delivery. To investigate how the decoupling and greedy training affect the representational capacity, we theoretically prove that the error produced by linear modules will not accumulate on unsupervised tasks in most cases. The theoretical and experimental results show that the proposed framework is highly efficient with reasonable performance, which may deserve more investigation.

Index Terms—Graph Neural Network, Backward Training, Efficient Training.

1 INTRODUCTION

In recent years, neural networks [1], [2], due to the impressive performance, have been extended to graph data, known as graph neural networks (GNNs) [3]. As GNNs significantly improve the results of graph tasks, it has been extensively investigated from different aspects, such as graph convolution network (GCN) [4], [5], graph attention networks (GATs) [6], [7], spatial-temporal GNN (STGNN) [8], graph auto-encoder [9], [10], graph contrastive learning [11], etc.

Except for the variants that originate from different perspectives, an important topic is motivated by the well-known inefficiency of GNN. In classical neural networks [2], the optimization is usually based on stochastic algorithms with limited batch [12], [13] since samples are independent of each other. However, the aggregation-like operations defined in [14] result in the dependency of each node on its neighbors and the amount of dependent nodes for one node increases exponentially with the growth of layers, which results in the unexpected increases of batch size. Some works are proposed based on neighbor sampling [14]–[17] and graph approximation [18] to limit the batch size, while some methods [19], [20] attempt to directly apply high-order graph operation and sacrifice the most non-linearity. The training stability is a

problem for neighbor sampling methods [14], [15], [17] though VRGCN [16] has attempted to control the variance via improving sampling. Note that the required nodes may still grow (slowly) with the increase of depth. Cluster-GCN [18] finds an approximate graph with plenty of connected components so that the batch size is strictly upper-bounded. The major challenge of these methods is the information missing during sampling. The simplified methods [19], [20] are efficient but the limited non-linearity may be the bottleneck of these methods. These methods may incorporate the idea of GIN [21] to improve the capacity [20].

To apply stochastic optimization while retaining the exact graph structure, we propose a framework, namely stacked graph neural network (SGNN), which decouples a multi-layer GNN as multiple simple GNN modules and then trains them simultaneously rather than connecting them with the increase of the depth. Inspired by the backward propagation algorithm, we find that the main difference between stacked networks [22] and classical networks is *no training information propagated from the latter modules to the former ones*. The lack of backward information delivery may be the main reason of the performance limitation of stacked models. The contributions are concluded as: (1) We accordingly propose a backward training strategy to let the former modules receive the information from the final loss and latter modules, which leads to a cycled training framework to control bias and train shallow modules correctly. (2) Under this framework, a multi-layer GNN can be decoupled into multiple simple GNNs, named as separable GNNs in this paper, so that every training step could use the stochastic optimization without any samplings or changes on graph. Therefore, SGNN could take both non-linearity and high efficiency into account. (3) We investigate how the decoupling and greedy training affect the representational capacity of the linear SGNN. It is proved that the error would not accumulate in most cases when the final objective is graph reconstruction.

* Corresponding author

This work is supported by The National Natural Science Foundation of China (No. 61871470).

Hongyuan Zhang and Yanan Zhu are with the School of Artificial Intelligence, Optics and ElectroNics (iOPEN), Northwestern Polytechnical University, Xi'an 710072, P.R. China. Hongyuan Zhang and Xuelong Li are also with the Institute of Artificial Intelligence (TeleAI), China Telecom Corp Ltd, 31 Jinrong Street, Beijing 100033, P. R. China.

E-mail: hyzhang98@gmail.com, xuelong_li@ieee.org.

©2022 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

The source codes are available at <https://github.com/hyzhang98/SGNN>.

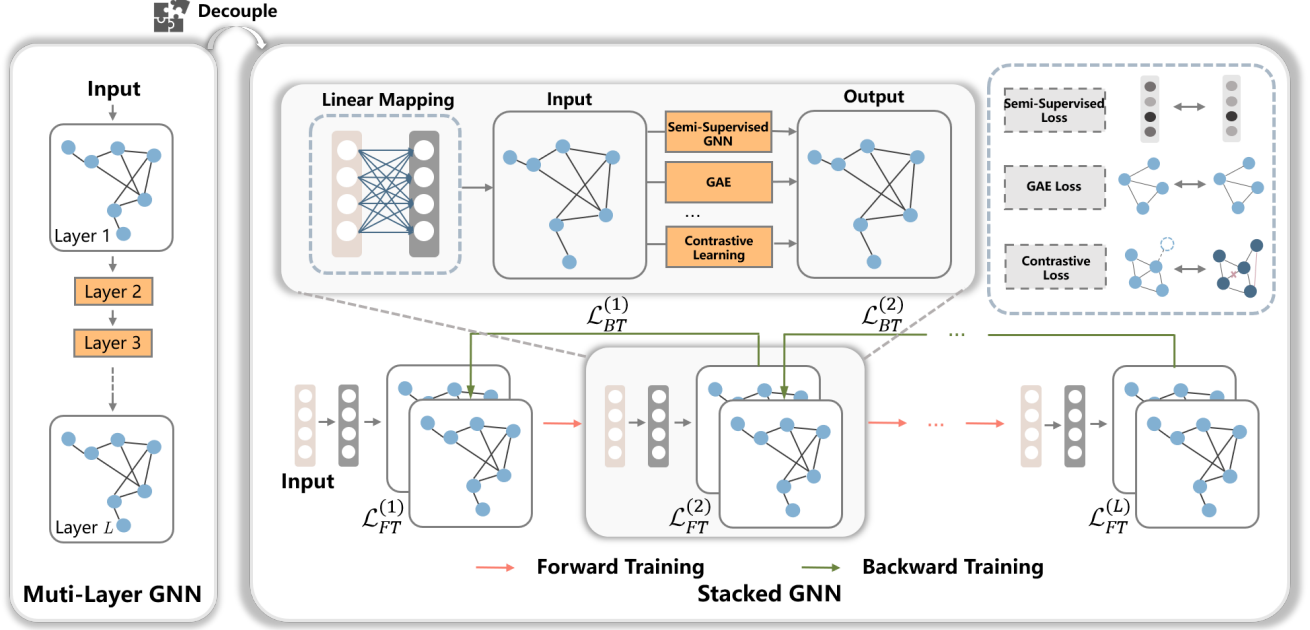


Fig. 1: Illustration of a stacked graph neural network decoupled from an L -layer GNN. To train each module individually, some loss (e.g., semi-supervised loss, unsupervised loss, contrastive loss) is required and it is denoted by $\mathcal{L}_{FT}^{(t)}$. To let the shallow modules perceive the deeper ones, \mathcal{M}_t passes back the *expected* input features to \mathcal{M}_{t-1} during the backward training. The divergence between the features output by \mathcal{M}_t and the expected features of \mathcal{M}_{t-1} formulates the BT loss $\mathcal{L}_{BT}^{(t)}$.

2 BACKGROUND

Graph Neural Networks: In the past years, graph neural networks [4]–[6], [21], [23], [24] have attracted more and more attention. GNNs are applied to not only graph tasks [25] (e.g., recommend systems [26]) but also other applications [27], [28] (e.g., computer vision [29]). In particular, graph convolution network (GCN) [4] has become an important baseline. By introducing self-attention techniques [30], graph attention networks (GAT) [6], [7] are proposed and applied to other applications [29], [31]. As [32] claimed that GNNs suffer from the over-smoothing problem, GALA [10] develops the graph sharpening and ResGCN [33] attempts to designs a deeper architecture. The theoretical works [32], [34], [35] have different views towards the depth of GNNs. Some works [32], [34] claimed that the expressive power of GNN decreases with the increase of layers, while the others argue that the assumptions in [34] may not hold and deeper GNNs have stronger power [35]. Moreover, some works [21], [36] investigate the expressive capability by showing the connection between Weisfeiler-Lehman test [37] and GNNs. Nevertheless, *most of them neglect the inefficiency problem of GNNs.*

Efficient Graph Neural Networks: To accelerate the optimization through batch gradient descent to GNN without too much deviation, several models [14]–[17] propose to sample data points according to graph topology. These models propose different sampling strategies to obtain stable results. GraphSAGE [14] produces a subgraph with limited neighbors for each node while FastGCN [15] samples fixed nodes for each layer with the importance sampling. The variance of sampling is further controlled in [16]. Cluster-GCN [18] aims to generate an approximate graph with plenty of connected components so that each component can be used as a batch per step. AnchorGAE [38] proposes to accelerate the graph operation by introducing the anchors to convert the original graph into a bipartite one so that the complexity can be

reduced to $\mathcal{O}(n)$ compared with the existing models [39], [40]. SGc [19] simplifies GCN by setting all activations of middle layers as linear functions and SSGC [20] further improves it. In summary, *SGNN proposed in this paper retains the non-linearity and requires no node sampling or sub-graph sampling.* L2-GCN [41] attempts to extend the idea of the classical stacked auto-encoder into the popular GCN while DGL-GNN [42] further develops a parallel version. They both fail to train all GNN modules jointly but *SGNN firstly offers a novel framework to train them like training layers in a conventional neural network.*

Connections to Existing Models: Stacked Auto-Encoder (SAE) [22] is a model applied to the pre-training of neural networks. It trains the current two-layer auto-encoder [43] and then feeds the latent features output by the middle layer to the next auto-encoder. The model is often used as a pre-training model instead of a formal model. MGAE [44] is an extension of SAE and its fundamental module is graph auto-encoder [9]. The main difference compared with the proposed model is whether each module could be perceived by modules from both forward and backward directions. The stack paradigm is similar to the classical boosting models [45]–[47] while some works [48], [49] also investigated the boosting algorithm of neural networks. In recent years, some boosting GNN models [50], [51] are also developed. The most boosting algorithms (e.g., [48], [50]) aim to learn a prediction function gradually while the proposed SGNN aims to learn ideal embeddings gradually. Note that AdaGCN [51] is also trained gradually and the features are combined using AdaBoost [45]. More importantly, *all these boosting methods for GNNs are only trained forward and the backward training is missing.* Deep neural interface [52] proposes to decouple neural networks to asynchronously accelerate the computation of gradients. *The decoupling is an acceleration trick to compute the gradients of L -layer networks,* while SGNN proposed in this paper explicitly separates an L -layer GNN into L simple

modules. In other words, *the ultimate goal of SGNN is not to optimize an L -layer GNN.*

3 PROPOSED METHOD

Motivated by SAE [22] and the fact that the simplified models [19], [20] are highly efficient for GNN, we therefore rethink the substantial difference between the stacked networks and multi-layer GNNs. To sum up, we attempt to answer the following two questions in this paper:

- Q1: *How to decouple a complex GNN into multiple simple GNNs and train them jointly?*
 Q2: *How does the decoupling affect the representational capacity and final performance?*

We will discuss the first question in this section and then elaborate on another one in Section 4.

3.1 Preliminary

Each decoupled GNN model of the proposed model is named as a module and the t -th module is denoted by \mathcal{M}_t for simplicity. The vector and matrix are denoted by lower-case and upper-case letters in bold, respectively. $\|\cdot\|$ represents the Frobenius norm. Given a graph, let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be adjacency matrix and $\mathbf{X} \in \mathbb{R}^{n \times d}$ be node features. A typical GNN layer can be usually defined as

$$\mathbf{H} = f(\mathbf{A}, \mathbf{X}, \mathbf{W}) = \varphi(\mathbf{P}\mathbf{X}\mathbf{W}), \quad (1)$$

where \mathbf{W} is projection coefficient and $\mathbf{P} = \phi(\mathbf{A})$ is a function of \mathbf{A} . When we discuss each individual module, we assume that $\mathbf{W} \in \mathbb{R}^{d \times k}$ for simplicity. For example, GCN [4] defines \mathbf{P} as $\mathbf{P}_{\text{GCN}} = \mathbf{D}^{-\frac{1}{2}}(\mathbf{A} + \mathbf{I})\mathbf{D}^{-\frac{1}{2}}$ and \mathbf{D} is the degree matrix of $\mathbf{A} + \mathbf{I}$. When multiple layers are integrated, the learned representation given by multiple GNN layers can be written as

$$\begin{aligned} \mathbf{H} &= f_1 \circ f_2 \circ \dots \circ f_L(\mathbf{A}, \mathbf{X}, \mathbf{W}_1, \dots, \mathbf{W}_L) \\ &= \varphi_L(\mathbf{P}_{\varphi_{L-1}}(\dots \varphi_1(\mathbf{P}\mathbf{X}\mathbf{W}_1) \dots) \mathbf{W}_L), \end{aligned} \quad (2)$$

where L is the amount of layers. Assume that the average number of neighbors is c . To compute \mathbf{H} , each sample will need $\mathcal{O}(c^L)$ extra samples. If the depth is large and the graph is connected, then all nodes have to be engaged to compute for one node. The uncontrolled batch size results in the time-consuming training. In vanilla GNNs, the computational complexity is $\mathcal{O}(KL\|\mathbf{A}\|_0 \sum_{i=0}^{m-1} d_i d_{i+1})$ on sparse graphs where K is the number of iterations and d_i is the dimension of \mathbf{W}_i . For large-scale datasets, both time and space complexity are too expensive.

3.2 Stacked Graph Neural Networks

Although the stacked networks usually have more parameters than multi-layer networks, which frequently indicates that the stacked networks may be more powerful, they only serve as a technique for pre-training. Specifically speaking, they simply transfer the representations learned by the current network to the next one but *no feedback is passed back*. It causes the invisibility of the succeeding modules and the final objective. As a result of the unreliability of the former modules, the stacked model is conventionally used as an unsupervised pre-training model.

Rethinking the learning process of a network, multiple layers are optimized simultaneously by gradient-based methods where the gradient is calculated by the well-known backward propagation algorithm [53]. The algorithm consists of forward propagation (FP)

and backward propagation (BP). FP computes the required values for BP, which can be viewed as an information delivery process. Note that FP is similar to the training of the stacked networks. Specifically, transferring the output of the current module to the next one in the stacked network is like the computation of neurons layer by layer during FP. Inspired by this, we aim to design a BP-like training strategy, namely *backward training (BT)*, so that the former modules could be tuned according to the feedback. The core idea of our stacked graph neural network (SGNN) is shown in Figure 1.

3.2.1 Separability: Crucial Concept for Efficiency

Before introducing SGNN in detail, we formally introduce the key concept and core motivation of how to accelerate GNN via SGNN.

Definition 3.1. *If a GNN model can be formulated as*

$$f(\mathbf{A}, \mathbf{X}, \mathbf{W}) = f_1(f_0(\mathbf{A}, \mathbf{X}), \mathbf{W}),$$

then it is a separable GNN. If it can be further formulated as

$$f(\mathbf{A}, \mathbf{X}, \mathbf{W}) = f_1(f_0(\mathbf{A}, \mathbf{X}), \mathbf{W}) = g_1(\mathbf{A}, g_0(\mathbf{X}, \mathbf{W})),$$

then it is a fully-separable GNN.

To keep simplicity, define the set of separable GNNs as

$$\begin{aligned} \mathcal{F}_k &= \{f : \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times d} \times \mathbb{R}^{d \times k} \mapsto \mathbb{R}^{n \times k} \\ &\quad | f(\mathbf{A}, \mathbf{X}, \mathbf{W}) = f_1(f_0(\mathbf{A}, \mathbf{X}), \mathbf{W})\} \end{aligned} \quad (3)$$

and the set of fully-separable GNNs as

$$\begin{aligned} \mathcal{F}_k^* &= \{f : \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times d} \times \mathbb{R}^{d \times k} \mapsto \mathbb{R}^{n \times k} \\ &\quad | f(\mathbf{A}, \mathbf{X}, \mathbf{W}) = f_1(f_0(\mathbf{A}, \mathbf{X}), \mathbf{W}) = g_1(\mathbf{A}, g_0(\mathbf{X}, \mathbf{W}))\}. \end{aligned} \quad (4)$$

Note that most single-layer GNN models are separable. For instance, SGC [19] is fully-separable where $f_0(\mathbf{A}, \mathbf{X}) = \mathbf{P}^m \mathbf{X}$ and $f_1(f_0(\mathbf{A}, \mathbf{X}), \mathbf{W}) = \varphi(f_0(\mathbf{A}, \mathbf{X}) \cdot \mathbf{W})$, while the single-layer GIN [3] is separable but not fully-separable since $\mathbf{P} \cdot \text{MLP}(\mathbf{X}) \neq \text{MLP}(\mathbf{P}\mathbf{X})$ usually holds. JKNet [54] consisting of one layer is also separable but not fully-separable. However, a single-layer GAT [6] is not separable since the graph operation is relevant to \mathbf{W} .

The separable property actually factorizes a GNN model to 2 parts, graph operation f_0 and neural operation f_1 . Since all dependencies among nodes in GNNs are caused by the graph operation, one can compute $\mathbf{X}' = f_0(\mathbf{A}, \mathbf{X})$ once (like **preprocessing**) in separable GNNs and then the GNN is converted into a typical network. After computing \mathbf{X}' , the information contained in graph has been passed into \mathbf{X}' and the succeeding sampling would not affect the topology of graph. Therefore, we can obtain a highly efficient GNN model that can be optimized by SGD, provided that each module is separable. On the other hand, *the fully-separable condition is essential for the backward training to pass back the information over multiple modules.* Since most single-layer GNNs are separable but not fully-separable, we show how to revise separable GNNs to **introduce the fully-separability**.

Then, we formally clarify the core idea of SGNN by showing how to handle **Q1**.

3.2.2 Forward Training (FT)

The first challenge is how to set the training objective for each module \mathcal{M}_t . It is crucial to apply SGNN to both supervised and unsupervised scenes. Suppose that we have a separable GNN module \mathcal{M} and let $\mathbf{H} = f(\mathbf{A}, \mathbf{X}, \mathbf{W})$ be the features learned by

the separable GNN module. For the unsupervised cases, if \mathcal{M}_t is a GAE, then the loss of FT is formulated as

$$\mathcal{L}_{FT} = \ell(\mathbf{A}, \mathbf{X}, \mathbf{W}) = d(\mathbf{A}, \kappa(\mathbf{H})), \quad (5)$$

where $d(\cdot, \cdot)$ represents the metric function and $\kappa : \mathbb{R}^{n \times k} \mapsto \mathbb{R}^{n \times n}$ is a mapping function. For instance, a simple loss introduced by [9] is $d(\mathbf{A}, \kappa(\mathbf{H})) = KL(\mathbf{A} \parallel \sigma(\mathbf{H}\mathbf{H}^T))$ where $\sigma(\cdot)$ is the sigmoid function, and $KL(\cdot \parallel \cdot)$ is the Kullback-Leibler divergence. The other options include but not limited to symmetric content reconstruction [10] and graph contrastive learning [11]. For modules with supervision information, a projection matrix $\mathbf{R} \in \mathbb{R}^{k \times c}$ is introduced to map the k -dimension embedding vector into soft labels with c classes. For the node classification, the loss can be simply set as

$$\mathcal{L}_{FT} = KL(\mathbf{Y} \parallel \text{softmax}(\mathbf{H}\mathbf{R})), \quad (6)$$

where $\mathbf{Y} \in \mathbb{R}^{n \times c}$ is the supervision information for supervised tasks. Note that the above loss is equivalent to the classical softmax regression if \mathbf{H} is constant. The loss could also be link prediction, graph classification, etc. Although base modules can utilize diverse losses, we only discuss the situation that all modules use the same kind of loss in this paper for simplicity.

Algorithm 1: Procedure of Stacked Graph Neural Networks Composed of L Modules

Input : Adjacency matrix \mathbf{A} , feature matrix \mathbf{X} , balance coefficient η , the number of epochs K, L , separable GNN modules $\{\mathcal{M}_t\}_{t=1}^L$, $\mathbf{H}_0 \leftarrow \mathbf{X}$.

Output : Features output by \mathcal{M}_L .

for $i = 1, \dots, K$ **do**

 # *Forward Train Stacked Graph Neural Networks*

for $t = 1, \dots, L - 1$ **do**

 Feed the current features to \mathcal{M}_t and reset \mathbf{U}_t :

$\mathbf{X}_t \leftarrow \mathbf{H}_{t-1}, \mathbf{U}_t \leftarrow \mathbf{I}$.

 # *Train with only \mathcal{L}_{FT} at the first forward training*

$\mathbf{X}'_t = f_0^{(t)}(\mathbf{A}, \mathbf{X}_t)$ # *Preprocessing for mini-batch algorithms*

 Compute loss $\mathcal{L}^{(t)} \leftarrow \mathcal{L}_{FT}^{(t)}$ if $i == 1$ else

$\mathcal{L}_{FT}^{(t)} + \eta \mathcal{L}_{BT}^{(t)}$

 Train \mathcal{M}_t by optimizing $\min_{\mathbf{W}_t} \mathcal{L}^{(t)}$ based on mini-batch algorithms.

 Obtain the features: $\mathbf{H}_t \leftarrow f_1^{(t)}(\mathbf{X}'_t, \mathbf{W}_t)$.

$\mathbf{X}_L \leftarrow \mathbf{H}_t, \mathbf{U}_L \leftarrow \mathbf{I}, \mathbf{X}'_L = f_0^{(L)}(\mathbf{A}, \mathbf{X}_L)$.

 Train \mathcal{M}_L by optimizing $\min_{\mathbf{W}_L, \mathbf{U}_L} \mathcal{L}_{FT}^{(L)}$ based on mini-batch algorithms.

 # *Backward Train Stacked Graph Neural Networks*

for $t = L - 1, L - 2, \dots, 1$ **do**

 Compute the expected output feature of \mathcal{M}_t :

$\mathbf{Z}_{t+1} \leftarrow (g_0^*)_{t+1}(\mathbf{X}_{t+1}, \mathbf{U}_{t+1})$.

 Train \mathcal{M}_t by optimizing $\min_{\mathbf{W}_t, \mathbf{U}_t} \mathcal{L}_{FT}^{(t)} + \eta \mathcal{L}_{BT}^{(t)}$ based on mini-batch algorithms.

by FP and BP. *BP lets the shallow layers perceive the deep ones through the feedback. In SGNN, the tail modules are invisible to the head ones in FT.* We accordingly design the *backward training (BT)* for SGNN. To achieve the reverse information delivery, the core idea is to **introduce the fully-separability** via defining **expected features**. For a separable GNN layer serving as a module in SGNN

$$\mathbf{H} = f(\mathbf{A}, \mathbf{X}, \mathbf{W}) = f_1(f_0(\mathbf{A}, \mathbf{X}), \mathbf{W}), \forall f \in \mathcal{F}_k, \quad (7)$$

where $\mathbf{W} \in \mathbb{R}^{d \times k}$, we aim at tuning its input so that the previous module can be aware of what kind of representations are required by the current module. However, \mathbf{X} is not a learnable parameter. A direct scheme is to use a transformation ($\mathbb{R}^d \mapsto \mathbb{R}^d$) to tune the input of the module

$$\mathbf{H} = f_1(f_0(\mathbf{A}, \psi(\mathbf{X}\mathbf{U})), \mathbf{W}) \quad (8)$$

where $\mathbf{U} \in \mathbb{R}^{d \times d}$ is a learnable square matrix and $\psi : \mathbb{R}^d \mapsto \mathbb{R}^d$ is a non-parametric function. Clearly, $f_0(\mathbf{A}, \mathbf{X}, \mathbf{U})$ can be regarded as a parametric GNN. To retain the high efficiency owing to f_0 , we therefore further constrain the revised $f_0(\mathbf{A}, \psi(\mathbf{X}\mathbf{U}))$ as a fully-separable GNN, i.e.,

$$\begin{aligned} f_0(\mathbf{A}, \mathbf{X}\mathbf{U}) &= f^*(\mathbf{A}, \mathbf{X}, \mathbf{U}) = f_1^*(f_0^*(\mathbf{A}, \mathbf{X}), \mathbf{U}) \\ &= g_1^*(\mathbf{A}, g_0^*(\mathbf{X}, \mathbf{U})), \end{aligned} \quad (9)$$

where $f^* \in \mathcal{F}_d^*$. Note that if \mathbf{U} is fixed as \mathbf{I} (or other constant matrices), then the modified layer is equivalent to original separable GNN in Eq. (7). To sum up, a separable GNN layer f is modified by introducing the fully-separability, which is shown as

$$\begin{aligned} \forall f \in \mathcal{F}_k, \mathbf{H} &= f(\mathbf{A}, \mathbf{X}, \mathbf{W}) = f_1(f_0(\mathbf{A}, \mathbf{X}), \mathbf{W}) \\ \Rightarrow \forall f^* \in \mathcal{F}_d^*, \mathbf{H} &= F(\mathbf{A}, \mathbf{X}, \mathbf{U}, \mathbf{W}) = f_1(f^*(\mathbf{A}, \mathbf{X}, \mathbf{U}), \mathbf{W}), \end{aligned} \quad (10)$$

where F represents a function from $\mathbb{R}^{n \times n} \times \mathbb{R}^{n \times d} \times \mathbb{R}^{d \times d} \times \mathbb{R}^{d \times k}$ to $\mathbb{R}^{n \times k}$. Denote $\mathbf{Z} = g_0^*(\mathbf{X}, \mathbf{U})$ and \mathbf{Z} is the *expected features*. Specifically, let \mathbf{Z}_t be the learned expected features during the backward training of \mathcal{M}_t and it serves as the expected input of \mathcal{M}_t from \mathcal{M}_{t-1} . In the forward training, the delivery of information is based on the learned features \mathbf{H}_t , and \mathbf{Z}_t plays the similar role in the backward training. The loss of backward training attempts to shrink the difference between the output feature \mathbf{H}_t of \mathcal{M}_t and expected input \mathbf{Z}_{t+1} of \mathcal{M}_{t+1} ,

$$\mathcal{L}_{BT}^{(t)} = d(\mathbf{H}_t, \mathbf{Z}_{t+1}) = d(F_t(\mathbf{A}, \mathbf{X}_t, \mathbf{U}_t, \mathbf{W}_t), \mathbf{Z}_{t+1}). \quad (11)$$

Note that \mathcal{L}_{BT} is only activated after the first forward training leading to the final loss of \mathcal{M}_t as

$$\mathcal{L}^{(t)} = \mathcal{L}_{FT}^{(t)} + \eta \mathcal{L}_{BT}^{(t)}, \quad (12)$$

and it is updated during each backward training. The introduction of \mathbf{Z} will not limit the application of stochastic optimization since the expected features can also be sampled at each iteration without restrictions. The procedure is summarized as Algorithm 1.

Remark that \mathbf{U}_t remains as the identity matrix during FT. This setting leads to each forward computation across L base modules being equivalent to a forward propagation L -layer GNN. In other words, an SGNN with L modules can be regarded as a decomposition of an L -layer GNN. One may concern that why not to learn \mathbf{U}_t and \mathbf{W}_t together in FT. In this case, we prefer to use \mathbf{U}_t only for learning the expected features of \mathcal{M}_{t+1} and the capability improvement from the co-learning \mathbf{U}_t in FT could also be implemented by \mathbf{W} , which is equivalent to use GIN [21] as base modules.

3.2.3 Backward Training (BT)

The second challenge is how to train multiple separable GNNs simultaneously in order to ensure performance. Roughly speaking, the gradients of all layers in classical multi-layer neural networks are computed exactly due to the repeated delivery of information

TABLE 1: Comparison of different efficient GNN models. “# Activations” represents the number of activation functions that can be used for each model at most. Preprocessing complexity is the required computational complexity before applying mini-batch gradient descent. For time complexity, we report the computational complexity per update. For simplicity, we assume that the original features and hidden features are all d -dimension. For sampling-based methods, r represents the number of sampled nodes. $\|\mathbf{A}_B\|_0$ represents the average node number of the partitioned sub-graphs.

	GraphSAGE	FastGCN	Cluster-GCN	SGC	S ² GC	SGNN
Separability	\times	\times	\times	\checkmark	\checkmark	\checkmark
Sampling	\checkmark	\checkmark	\checkmark	\times	\times	\times
# Activations	L	L	L	1	1	L
Preprocessing complexity	0	$\mathcal{O}(\ \mathbf{A}\ _0)$	$\mathcal{O}(n)$	$\mathcal{O}(L\ \mathbf{A}\ _0 d)$	$\mathcal{O}(L(L+1)\ \mathbf{A}\ _0 d + nd)$	$\mathcal{O}(LK\ \mathbf{A}\ _0 d)$
Time complexity per update	$\mathcal{O}(r^L md^2)$	$\mathcal{O}(rLmd^2)$	$\mathcal{O}(\ \mathbf{A}_B\ _0 Ld^2)$	$\mathcal{O}(md^2)$	$\mathcal{O}(md^2)$	$\mathcal{O}(md^2)$

3.3 Complexity

As each base module is assumed as a separable GNN, both FT and BT of \mathcal{M}_t can be divided into two steps, the preprocessing step for graph operation and the training step for parameters learning. Denote the output dimension of \mathcal{M}_t as d_t and the dimension of original content feature as $d_0 = d$. The preprocessing to compute $\mathbf{X}'_t = f_0(\mathbf{A}, \mathbf{X}_t)$ requires $\mathcal{O}(\|\mathbf{A}\|_0 d_t)$ cost. Suppose that each module is trained E iterations and the batch size is set as m . Then the computation cost of the training step is $\mathcal{O}(Emd_{t-1}d_t)$. Note that only the GNN mapping is considered and the computation of the loss is ignored. Overall, the computational complexity of an SGNN with L modules is approximately $\mathcal{O}(K\|\mathbf{A}\|_0 \sum_{t=0}^{L-1} d_t + Em \sum_{t=1}^L d_{t-1}d_t)$. Remark that the graph is only used once during every epoch and no sampling is processed on the graph such that the graph structure is completely retained which is unavailable in the existing fast GNNs. The coefficient of $\|\mathbf{A}\|_0$ is only $K \sum_{t=1}^{L-1} d_t$. The space complexity is only $\mathcal{O}(\|\mathbf{A}\|_0 + md_{t-1}d_t)$. Therefore, the growth of graph scale will not affect the efficiency of SGNN. Due to the efficiency, most experiments of SGNN can be conducted on a PC with an NVIDIA 1660 (6GB) and 16GB RAM. To better clarify the advantages of the proposed SGNN, we summarize the characteristics of diverse efficient GNNs in Table 1.

4 THEORETICAL ANALYSIS

To answer Q2 raised in the beginning of Section 3, we discuss the impact of the decoupling in this section. Intuitively speaking, if an L -layer GNN achieves satisfactory results, then there exists $\{\mathbf{W}_t\}_{t=1}^L$ such that an SGNN with L modules could achieve the same results. However, each \mathcal{M}_t is trained greedily according to the forward training loss \mathcal{L}_{FT} , while middle layers of a multi-layer GNN are trained according to the same objective. The major concern is whether the embedding learned by a greedy strategy leads to an irreversible deviation in the forward training.

In this section, we investigate the possible side effects on a specific SGNN comprised of unsupervised modules defined in Eq. (5) with linear activations. The conclusion is not apparent since simply setting \mathbf{W} as an identity matrix does not prove it for GNN due to the existence of \mathbf{P} . Remark that a basic premise is that the previous module has achieved a reasonable result.

Given a linear separable-GNN module \mathcal{M} which is defined as $\mathbf{H} = \mathbf{P}\mathbf{X}\mathbf{W} \in \mathbb{R}^{n \times k}$, suppose that the forward training uses the reconstruction error $\|\mathbf{P} - \mathbf{H}\mathbf{H}^T\|$ as \mathcal{L}_{FT} . We first introduce the matrix angle to better understand whether the preconditions of Theorem 4.1 are practicable.

Definition 4.1. Given two matrices $\mathbf{B}_1, \mathbf{B}_2 \in \mathbb{R}^{n \times n}$, we define the matrix angle as $\theta(\mathbf{B}_1, \mathbf{B}_2) = \arccos\langle \mathbf{B}_1, \mathbf{B}_2 \rangle / (\|\mathbf{B}_1\| \cdot \|\mathbf{B}_2\|)$.

Before elaborating on theorems, we introduce the following assumption, which separates the discussions into two cases.

Assumption 4.1. $\mathbf{X}\mathbf{X}^T$ does not share the same eigenspace with $\mathbf{P} - \mathbf{X}\mathbf{X}^T$.

Note that the above assumption is weak and frequently holds in most cases. For simplicity, $\mathbf{U}_r \in \mathbb{R}^{n \times r}$ ($r \in \mathbb{N}_+$) is the eigenvectors associated with r leading eigenvalues. Under this assumption, we find that the error of \mathcal{M} is upper-bounded by ε .

Theorem 4.1. Let $\delta = 1 - \cos(\theta_*/2)$ and $\theta_* = \theta((\mathbf{P} - \mathbf{X}\mathbf{X}^T)\mathbf{Q}, \mathbf{Q}(\mathbf{P} - \mathbf{X}\mathbf{X}^T))$ where $\mathbf{Q} = \mathbf{U}_o \mathbf{U}_o^T - \mathbf{I}/2$ and $o = \min(\text{rank}(\mathbf{X}), k)$. Under Assumption 4.1, if $\|\mathbf{P} - \mathbf{H}\mathbf{H}^T\| = \varepsilon \leq \mathcal{O}(\delta)$ and $\sigma_* \leq \mathcal{O}(\sqrt{\delta})$ where σ_* is the $(o+1)$ -th largest singular value of \mathbf{X} , then there exists $\mathbf{W} \in \mathbb{R}^{d \times k}$ so that $\|\mathbf{P} - \mathbf{H}\mathbf{H}^T\| \leq \varepsilon$. In other words, if ε is small enough, then $\mathbf{H}\mathbf{H}^T$ could be a better approximation than $\mathbf{X}\mathbf{X}^T$.

Specially, if $\text{rank}(\mathbf{X}) \leq k$ or $k = d$, $\sigma_* = 0$ so that $\sigma_* \leq \mathcal{O}(\sqrt{\delta})$ holds. From the above theorem, we claim that the error through \mathcal{M} will not accumulate (i.e., bound by ε) provided that the input \mathbf{X} , the output of the previous modules, is well-trained. We further provide an upper-bound of error if Assumption 4.1 does not hold. The following theorem shows the increasing speed of error is at most linear with the tail singular-values.

Theorem 4.2. If Assumption 4.1 does not hold, then there exists $\mathbf{W} \in \mathbb{R}^{d \times k}$ so that $\|\mathbf{P} - \mathbf{H}\mathbf{H}^T\| \leq \varepsilon + \mathcal{O}(\sigma_*^2)$.

Corollary 4.1. Given an SGNN with L linear modules $\{\mathcal{M}_t\}_{t=1}^L$ with $\mathcal{L}_{FT}^{(t)} = \|\mathbf{P} - \mathbf{H}_t \mathbf{H}_t^T\|$, if $\mathbf{P} - \mathbf{H}_1 \mathbf{H}_1^T$ and $\mathbf{H}_1 \mathbf{H}_1^T$ share the same eigenspace, then $\mathcal{L}_{FT}^{(L)} \leq \mathcal{L}_{FT}^{(1)} + \sum_{t=1}^{L-1} \mathcal{O}(\sigma_*^2(\mathbf{H}_t))$.

Based on Theorem 4.2, we conclude that the residual would not accumulate rapidly when Assumption 4.1 does not hold. All proofs are put in appendix.

5 EXPERIMENTS

In this section, we conduct experiments to investigate: (1) whether the performance of SGNN could approach the performance of the

TABLE 2: Data Information

Dataset	Nodes	Edges	Classes	Features	Train / Val / Test Nodes
Cora	2,708	5,429	7	1,433	140 / 500 / 1,000
Citeseer	3,327	4,732	6	3,703	120 / 500 / 1,000
Pubmed	19,717	44,338	3	500	60 / 500 / 1,000
Reddit	233K	11.6M	41	602	152K / 24K / 55K

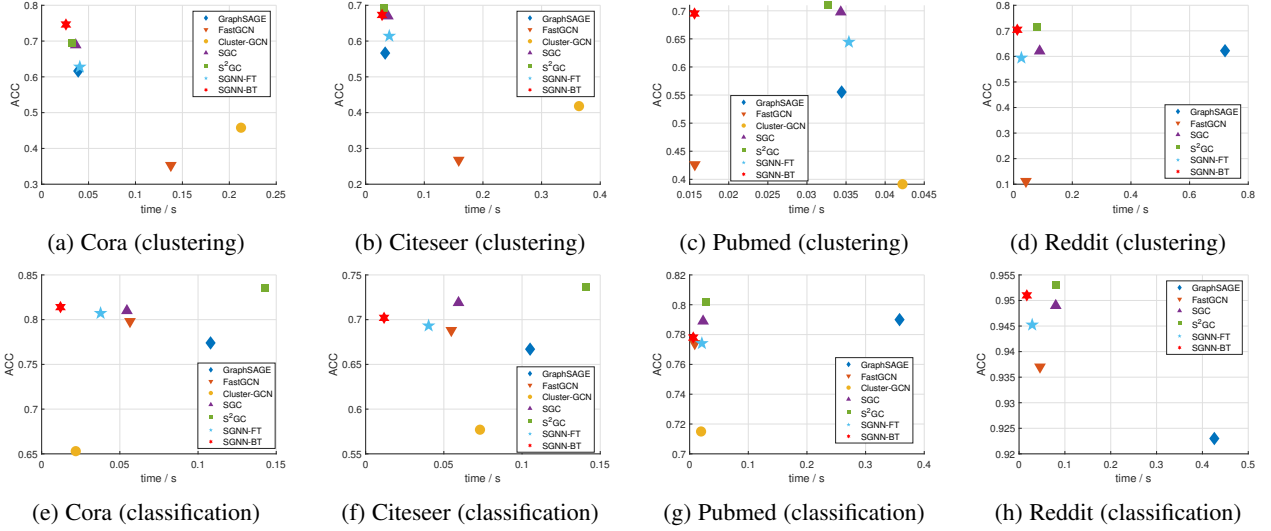


Fig. 2: Performance and training efficiency of several scalable GNNs. The efficiency metric is computed by “Consuming Time / # Iterations”. The consuming time begins from loading data into RAM. The first line shows the result of node clustering and the second line is the result of node classification. Note that SGC is slower than SGNN since SGNN updates parameters more times than SGC per graph preprocessing owing to the design of BT.

TABLE 3: Node clustering results on 4 datasets: The boldface and underline results indicate the optimal and sub-optimal results respectively. “N/A” means the Out-Of-Memory (OOM) exception.

Datasets	Cora		Citeseer		PubMed		Reddit	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
K-Means	0.4922	0.3210	0.5401	0.3054	0.5952	0.2780	0.1927	0.2349
ARGA	0.6400	0.4490	0.5730	0.3500	0.6807	0.2757	N/A	N/A
MGAE	0.6806	0.4892	0.6691	0.4158	0.5932	0.2957	N/A	N/A
GraphSAGE	0.6163	0.4826	0.5664	0.3425	0.5554	0.0943	0.6225	0.7291
FastGAE	0.3527	0.1553	0.2672	0.1178	0.4262	0.0442	0.1115	0.0715
ClusterGAE	0.4579	0.2261	0.4182	0.1767	0.3913	0.0001	N/A	N/A
GAE	0.5960	0.4290	0.4080	0.1760	0.6861	0.2957	N/A	N/A
AGC (SGC)	0.6892	0.5368	0.6700	0.4113	0.6978	0.3159	0.5833	0.6894
SGNN-FT	0.6278	0.5075	0.6141	0.3776	0.6444	0.2312	0.5943	0.7156
SGNN-BT	0.7463	0.5546	<u>0.6730</u>	<u>0.4159</u>	0.6951	0.3337	0.7042	0.7601
S ² GC	0.6960	<u>0.5471</u>	0.6911	0.4287	0.7098	<u>0.3321</u>	0.7011	0.7509
GAE-S ² GC	0.6976	0.5317	0.6435	0.3969	0.6528	0.2452	0.6272	0.7158
SGNN-S ² GC	<u>0.7223</u>	0.5404	<u>0.6822</u>	<u>0.4243</u>	<u>0.7084</u>	0.3302	<u>0.7023</u>	<u>0.7575</u>

original L -layer GNN in a highly-efficient way; (2) what the impact of the non-linearity and flexibility brought by the decoupling is.

To sufficiently answer the above 2 problems, both node clustering and semi-supervised node classification are used. It should be emphasized that SGNN can be used for both transductive and inductive tasks. We mainly conduct experiments on datasets for transductive learning since the inefficiency problem is more likely to appear in transductive learning scenes. There are 4 common datasets, including Cora, Citeseer, Pubmed [25], and Reddit [14], used for our experiments. Cora and Citeseer contain thousands of nodes. Pubmed has nearly 20 thousands nodes and Reddit contains more than 200 thousands nodes, which are middle-scale and large-scale datasets respectively. The details of four common datasets are shown in Table 2. Experiments on two OGB datasets can be

found in Section 5.4.

5.1 Node Clustering

5.1.1 Experimental Settings

We first testify the effectiveness of SGNN on the node clustering. We compare our method against 10 methods, including a baseline clustering model K-means, three GCN models without considering training efficiency (GAE [9], ARGA [55], MGAE [44]), and six fast GCN models with GAE-loss (GraphSAGE [14], FastGAE [15], ClusterGAE [18], AGC [56] (an unsupervised extension of SGC [19]), S²GC [20], and GAE-S²GC).

The used codes are based on the publicly released implementation. To ensure fairness, all multi-layer GNN models consist of two layers and SGNN is comprised of two modules. For models that

TABLE 4: Investigation about the impact of L on SGNN and GAE regarding node clustering.

L	Cora			Citeseer			Pubmed		
	SGNN	GAE	FastGAE	SGNN	GAE	FastGAE	SGNN	GAE	FastGAE
2	0.75	0.60	0.35	0.67	0.41	0.27	0.70	0.69	0.43
3	0.66	0.63	0.33	0.65	0.58	0.25	0.64	0.64	0.42
4	0.68	0.65	0.33	0.59	0.58	0.24	0.64	0.60	0.41
5	0.69	0.62	0.33	0.53	0.45	0.24	0.64	0.60	0.41
6	0.69	0.53	0.32	0.44	0.32	0.24	0.64	0.48	0.41
7	0.68	0.52	0.32	0.44	0.31	0.24	0.64	0.46	N/A

could be trained by stochastic algorithms, the size of mini-batch is set as 128. The learning rate is set as 0.001 and the number of epochs is set as 100. We set the size first layer as 128 and the second layer size as 64. \mathbf{U} is initialized as an identity matrix, and η is set as 10^3 by default. The number of backward training is set as 5 or 10. To investigate the effectiveness of backward training, we report the experimental results with sufficient training for each module, which is denoted by SGNN-FT, while SGNN with the proposed backward training is represented by SGNN-BT.

To study the performance of SGNN with different base models, we choose a fully-separable GNN, S^2GC , as the base model and this method is marked as SGNN- S^2GC . Note that SGNN- S^2GC is also trained with the proposed BT strategy, which is same as SGNN-BT. The original S^2GC does not apply any activation functions while we add a activation functions, which are same as the SGNN-FT, to the embeddings of each S^2GC base models. As S^2GC does not use the GAE framework, we also added a competitor, namely GAE- S^2GC , which uses S^2GC as the encoder, to ensure fairness. For SGNN- S^2GC and GAE- S^2GC , all extra hyper-parameters are simply set according to the original paper of S^2GC for both node clustering and node classification. We do not tune the hyper-parameters of S^2GC manually.

All methods are run five times and the average scores are recorded. The result are summarized in Table 3.

5.1.2 Performance

From Table 3, we find that SGNN outperforms in most datasets. If the released codes could not run on Reddit due to out-of-memory (OOM), we put the notation “N/A” instead of results. In particular, SGNN-BT obtains good improvements on Reddit with high efficiency. Specifically speaking, it is about 8% higher than the well-known GraphSAGE. SGNN-FT performs above the average on some datasets. It usually outperforms GraphSAGE but fails to exceed SGC. Due to the deeper structure caused by multiple modules, the performance of SGNN excels the simple GAE. It also outperforms SGC due to more non-linearity brought by multiple modules. Note that S^2GC and SGC are strong competitors, while **SGNN can easily employ them as a base module since they are separable** which is shown by SGNN- S^2GC . It is easy to find that SGNN- S^2GC usually achieves similar results compared with S^2GC . As it is slower than SGNN-BT and the performance improvement is not stable, we recommend to use simple SGNN-FT instead of SGNN- S^2GC in practice. From the table, we find that it is unnecessary to modify S^2GC to GAE- S^2GC since the GAE architecture does not improve the performance of S^2GC at all. From the ablation experiments, SGNN-BT works better than SGNN-FT, which indicates the necessity of the backward training.

We also investigate how the number of modules L affects the node clustering accuracy and the results averaged over 5 runs are reported in Table 4. To ensure fairness, we also show the performance of GAE with the same depth though deeper GCN and GAE usually return unsatisfied results. Note that the neurons of each layers are set as [256, 128, 64, 32, 16, 16, 16], respectively. It is not hard to find that SGNN is robust to the depth compared with the conventional GNNs.

5.1.3 Efficiency

Figure 2 shows the consuming time of several GNNs with higher efficiency on Pubmed and Reddit. Instead of neglecting the preprocessing operation, we measure the efficiency through a more rational way. We record the totally consuming time after loading data into RAM and then divide the total number of updating parameters of GNNs. The measurement could reflect the real difference of diverse training techniques aiming to apply batch-based algorithms to GNN. It should be emphasized the reason why SGC is worse than SGNN regarding the consuming time. The key point is the different costs of their *preprocessing* operation. For an L -order SGC, the computation cost of $\mathbf{P}^L \mathbf{X}$ is at least $\mathcal{O}(\|\mathbf{A}\|_0 L d)$ while SGNN with L first-order modules totally requires $\mathcal{O}(\|\mathbf{A}\|_0 \sum_{i=1}^L d_i)$ for the same preprocessing operation. The metric also provides a fair comparison between SGC and other models since the stopping criteria are always different.

5.2 Node Classification

5.2.1 Experimental Setting

We also conduct experiments of semi-supervised classification on four datasets. The split of datasets follows [19] which is shown in Table 2. We compare SGNN against GCN [4], GAT [6], DGI [57], APPNP [58], L2-GCN [41] FastGCN [15], GraphSAGE [14], Cluster-GCN [18], SGC [19], GCNII [59], and S^2GC [20]. Similarly, we testify SGNN with two different base models, namely SGNN-BT and SGNN- S^2GC . The experimental settings are same as the experiments of node clustering. For GraphSAGE, we use the mean operator by default and some notations are added if the extra operators are used. On citation networks, the learning rate is set as 0.01, while it is 10^{-4} on Reddit. Since the nodes for training are less than 200 on citation networks, we use all training points in each iteration for all methods while we sample 256 points as a mini-batch for approaching expected features during backward training of SGNN. On Reddit, the batch size of all batch-based models is set as 512. We do not apply the early stopping criterion used in [4] and the max iteration follows the setting of SGC. The embedding dimensions of each module are the same as the setting in node clustering. For the sake of fairness, we report the results obtained by

TABLE 5: Test accuracy (%) of supervised SGNN averaged over 10 runs on 3 citation datasets

Datasets	Cora	Citeseer	Pubmed
GAT	83.0 \pm 0.7	72.5 \pm 0.7	79.0 \pm 0.3
DGI	82.3 \pm 0.6	71.8 \pm 0.7	76.8 \pm 0.6
GCNII	85.5 \pm 0.5	<u>73.4 \pm 0.6</u>	80.2 \pm 0.4
FastGCN [†]	79.8 \pm 0.3	68.8 \pm 0.6	77.4 \pm 0.3
GraphSAGE [†]	77.4 \pm 0.6	66.7 \pm 0.3	79.0 \pm 0.3
Cluster-GCN [†]	65.3 \pm 3.9	57.7 \pm 2.3	71.5 \pm 1.9
GCN	81.4 \pm 0.4	70.9 \pm 0.5	79.0 \pm 0.4
SGC	81.0 \pm 0.0	71.9 \pm 0.1	78.9 \pm 0.0
SGNN-FT	81.0 \pm 0.6	69.3 \pm 0.4	77.4 \pm 0.3
SGNN-BT	81.4 \pm 0.5	70.2 \pm 0.4	77.8 \pm 0.3
S ² GC	83.5 \pm 0.0	73.6 \pm 0.1	80.2 \pm 0.0
SGNN-S ² GC	<u>83.8 \pm 0.2</u>	73.2 \pm 0.3	80.2 \pm 0.4

TABLE 6: Test accuracy (%) averaged over 5 runs on Reddit

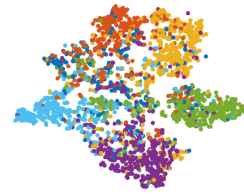
Methods	Reddit
GAT	N/A
DGI	94.0
SAGE-mean	95.0
SAGE-GCN	93.0
FastGCN	93.7
Cluster-GCN	96.6
GCN	N/A
SGC	94.9
SGNN-FT	94.54 \pm 0.01
SGNN-BT	95.10 \pm 0.02
S ² GC	<u>95.3</u>
SGNN-S ² GC	95.28 \pm 0.03



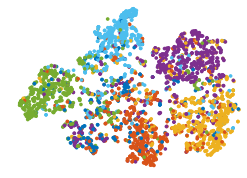
(a) GCN on Cora



(b) SGNN on Cora



(c) GCN on Citeseer



(d) SGNN on Citeseer

Fig. 3: Visualization of SGNN comprised of 3 modules and 3-layer GCN on node classification. For SGNN, the output of \mathcal{M}_3 is visualized. For GCN, the output of the final GCN-layer is visualized.

SGNN with two modules using first-order operation. The forward training loss is defined in Eq. (6). Moreover, all compared models share an identical implementation of their mini-batch iterators, loss function and neighborhood sampler (when applicable). The balance coefficient of \mathcal{L}_{FT} and \mathcal{L}_{BT} is set as 1 by default. We report the results averaged over 10 runs on citation datasets and 5 runs on Reddit in Table 5 and Table 6. The hyper-parameters are shared for different datasets which are optimized on Cora.

5.2.2 Performance

The results of compared methods in Table 5 are taken from the corresponding papers. When the experimental results are missed, we run the publicly released codes and the corresponding records are superscripted by [†]. From Tables 5 and 6, we conclude that SGNN outperforms the models with neighbor sampling such as GraphSAGE, FastGCN, and Cluster-GCN on citation networks and the performance of SGNN exceeds most models on Reddit. On simple citation networks, SGNN **loses the least accuracy compared with other batch-based models**, which is close to GCN. Owing to the separability of each module, the batch sampling requires no neighbor sampling and causes no loss of graph information. Note that we simply employ the single-layer GCN as the base modules in our experiments, while some high-order methods that obtain competitive results are also available for SGNN, e.g., SGNN-S²GC. As shown by the results of SGNN-S²GC, SGNN can be indeed improved by employing more complicated separable GNNs as base models.

TABLE 7: Node Classification Results on Large Datasets

Datasets	Products		Arxiv	
	Test Acc	Val Acc	Test Acc	Val Acc
MLP	61.06	75.54	55.50	57.65
Softmax	47.70	N/A	52.77	N/A
GraphSAGE	<u>78.50</u>	92.24	71.49	<u>72.77</u>
Cluster-GCN	78.97	<u>92.12</u>	N/A	N/A
GCN	75.64	92.00	71.74	73.00
SGC	68.87	N/A	68.78	N/A
S ² GC	70.22	N/A	70.15	N/A
SGNN-FT	68.10	86.26	65.25	64.00
SGNN-BT	74.44	91.13	<u>71.57</u>	71.66

Although some methods achieve preferable results, they either fail to run or obtain unsatisfactory results on large-scale datasets. Besides, we also show the comparison of efficiency on node classification task in Figure 2.

5.2.3 Visualization to Show Impact of the Decoupling

In Figure 3, we visualize the output of a 3-module SGNN and a 3-layer GCN to directly show that *the decoupling would not cause the trivial features*, which corresponds to the theoretical conclusion in Section 4. To show the benefit of the non-linearity brought by SGNN and the backward training, the convergence curves of SGC,

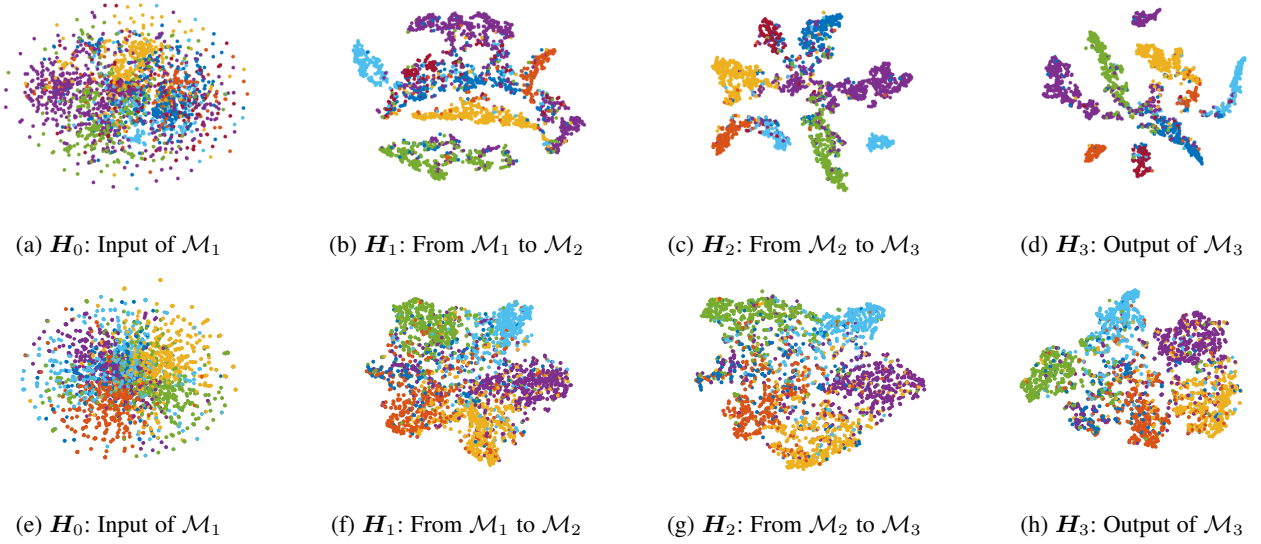


Fig. 4: Visualization of a trained SGNN comprised of 3 modules on node classification of Cora and Citeseer. The first line shows the visualization of Cora and the bottom line shows the visualization of Citeseer.

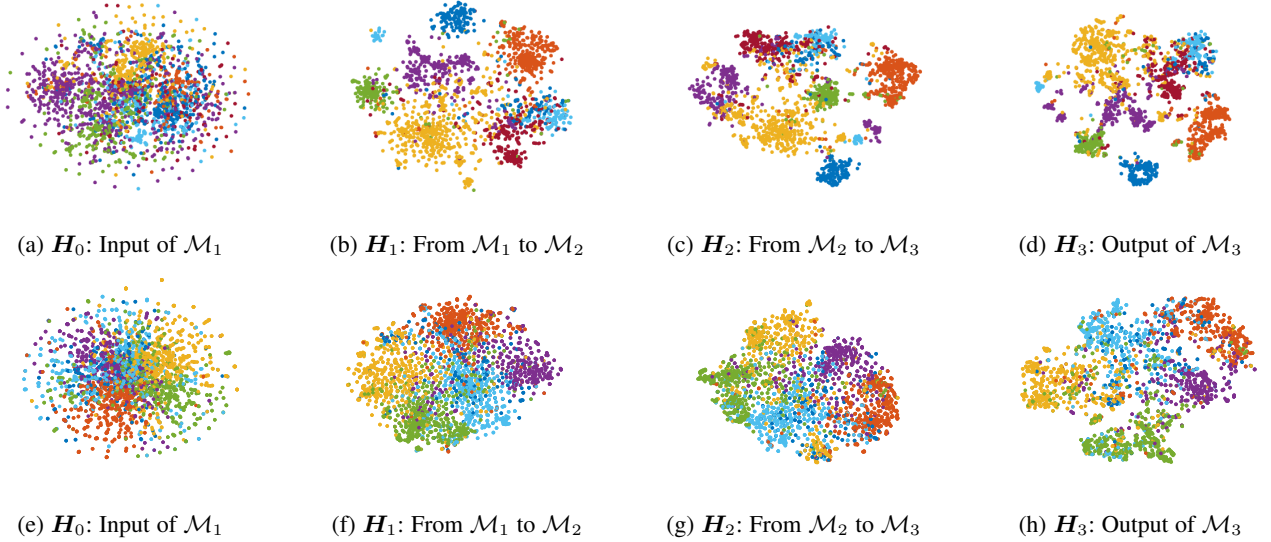


Fig. 5: Visualization of a trained SGNN comprised of 3 modules on node clustering of Cora and Citeseer. The first line is visualization on Cora and the second line is visualization on Citeseer.

SGNN-FT, and SGNN-BT are shown in Figure 7. Note that the figure shows the variation of the *final loss*. In SGNN, the final loss is the loss of \mathcal{M}_L , while it is the unique training loss in SGC. SGC with L -order graph operation is used. From this figure, we can conclude that: (1) The non-linearity does lead to a better loss value; (2) The backward training significantly decreases the loss. In summary, *the decoupling empirically does not cause the negative impact*.

5.3 Visualization

We also provide more visualizations in Figures 4 and 5. We run SGNN with 3 GNN modules and visualize the input and output of \mathcal{M}_1 , \mathcal{M}_2 , and \mathcal{M}_3 through t -SNE on Cora and Citeseer, for node clustering and node classification. The purpose of these two figures is to empirically investigate *whether the decoupling would cause*

the accumulation of residuals and errors. The experimental results support the theoretical results that are provided in Section 4. One may concern the impact of η (trade-off coefficient between \mathcal{L}_{FT} and \mathcal{L}_{BT}) on the performance. We testify SGNN with different η from $\{10^{-5}, 10^{-3}, 10^{-1}, 10^1, 10^3, 10^5\}$ and find that $\eta = 10^3$ **usually leads to good results**. Accordingly, we only report results SGNN with $\eta = 10^3$ in this paper. Moreover, we show the impact of η to node clustering on Cora and Citeseer in Figure 8.

Moreover, we show the output of \mathcal{M}_1 of different periods in Figure 6, in order to show the impact of the backward training. From the figure, we find that BT indeed affects the latent features, which is particularly apparent in Figure 6d.

5.4 Experiments on OGB Datasets

We further show some experiments of node classification on two OGB datasets, OGB-Products and OGB-Arxiv, which are

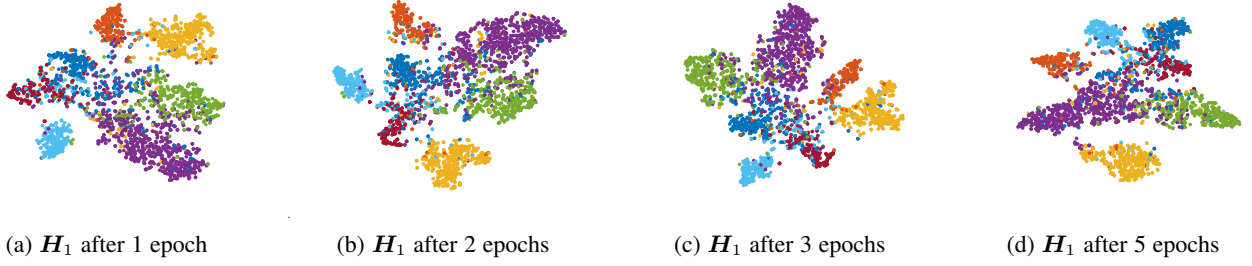


Fig. 6: t-SNE Visualization of the output of \mathcal{M}_1 from a trained SGNN comprised of 3 modules on node classification of Citeseer.

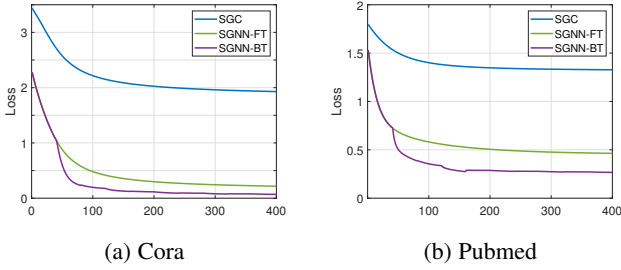


Fig. 7: Convergence curve of the *final loss*. The order of SGC is set as L . The backward training significantly decreases the final loss and the non-linearity also plays an important role.

downloaded from <https://ogb.stanford.edu/docs/nodeprop/>. The OGB-Products contains more than 2 million nodes and OGB-Arxiv contains more than 150 thousand nodes.

It should be emphasized that we only use the simple single-layer GCN as the base module of SGNN. The performance can be further improved by incorporating different models such as GCNII, GIN, *etc.* In particular, we only tune hyper-parameters on Arxiv, and we simply report results of SGNN with settings from Reddit.

Remark: One may concern that the proposed SGNN cannot achieve the state-of-the-art results like node clustering. It should be emphasized that **an SGNN with L modules should be regarded as a variant of an L -layer GNN**. So it is more fair to compare SGNN with GCN. From Tables 5, 6, and 7, we can find that the performance of SGNN can approach GCN with high efficiency. A main reason why SGNN cannot outperforms other models on supervised tasks like node clustering is the difficulty of designing proper training losses for middle modules. It essentially originates from the black-box property of neural networks, *i.e.*, **what kind of latent features is preferable for deeper layers**. From the experimental results, we find that it is not the optimal scheme to simply set the final supervised loss as the training loss of middle module. Instead, on node clustering, we prove that the greedy strategy would not accumulate the error and the experimental results validate the theoretical conclusion.

6 CONCLUSION AND FUTURE WORKS

In this paper, we propose the Stacked Graph Neural Networks (SGNN). We first decouple a multi-layer GNN into multiple simple GNNs, which is formally defined as separable GNNs in our paper to ensure the availability of batch-based optimization without loss of graph information. The bottleneck of the existing stacked models is that the information delivery is only unidirectional, and therefore a backward training mechanism is developed to make the former

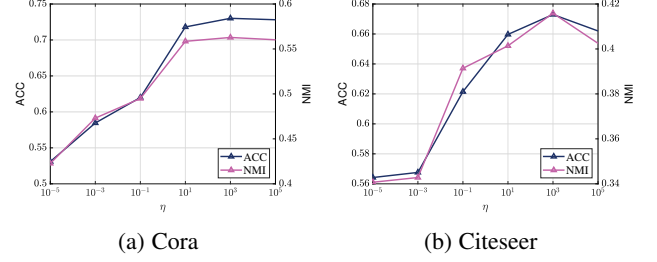


Fig. 8: Impact of η to node clustering on Cora and Citeseer.

modules perceive the latter ones. We also theoretically prove that the residual of linear SGNN would not accumulate in most cases for unsupervised graph tasks. The theoretical and experimental results show that the proposed framework is more than an efficient method and it may deserve further investigation in the future. The theoretical analysis focuses on linear SGNN and the generalization bound is also not investigated in this paper. Therefore, they will be the core of our future work. Moreover, since \mathcal{L}_{FT} could be any losses and SGNN fails to achieve state-of-the-art results, how to set the most appropriate training loss, especially on supervised tasks, for each module will be also a crucial topic in our future works. It may help us understand how neural networks work.

7 PROOFS

7.1 Lemma for Proofs

Lemma 7.1. For two given symmetric matrices A and B , A and B share the same eigenspace if and only if A and B commute.

Proof. First, if A and B share the same eigenspace, then there exists P such that

$$A = P\Lambda_A P^T, B = P\Lambda_B P^T. \quad (13)$$

Accordingly, we have $AB = P\Lambda_A\Lambda_B P^T = BA$.

Then, we turn to prove the converse. If A and B commute, suppose that $Au = \lambda u$ and then

$$ABu = BAu = \lambda Bu. \quad (14)$$

Apply eigenvalue decomposition, we have $A = U\Lambda U^T$. Note that if $Au_1 = \lambda_1 u_1$, $Au_2 = \lambda_2 u_2$, and $\lambda_1 \neq \lambda_2$, then $u_2^T B u_1 = 0$ since Bu_1 is also an eigenvector associated with λ_1 . Therefore, $U^T B U$ is a block-diagonal matrix, *i.e.*,

$$U^T B U = \begin{bmatrix} C_1 & & & \\ & C_2 & & \\ & & \ddots & \\ & & & C_k \end{bmatrix}. \quad (15)$$

Apply eigendecomposition to C_i ,

$$C_i = U_i \Lambda_i U_i^T. \quad (16)$$

Denote

$$T = \begin{bmatrix} U_1 & & & \\ & U_2 & & \\ & & \ddots & \\ & & & U_k \end{bmatrix}, \quad (17)$$

and $V = UT$, which leads to

$$T^T U^T B U T = \Lambda_B \text{ and } T^T U^T A U T = \Lambda_A \quad (18)$$

where $V^T V = I$. Hence, the lemma is proved. \square

7.2 Proof of Theorem 4.1

Proof. Use the notation $\ell(\cdot, \cdot, \cdot)$ as the reconstruction loss

$$\ell(P, X, W) = \|P - P X W W^T X^T P^T\|. \quad (19)$$

According to the conditions, we define

$$E = P - X X^T \Rightarrow \|E\| = \varepsilon. \quad (20)$$

Apply SVD, we can factorize X as

$$X = U_o \Sigma_o V_o^T + U_e \Sigma_e V_e^T, \quad (21)$$

where $e = d - o$. Clearly, we have $V_o^T V_e = 0$ and thus

$$P = X X^T + E = U_k \Sigma_k^2 U_k^T + U_e \Sigma_e^2 U_e^T + E. \quad (22)$$

Therefore, HH^T can be written as

$$P(U_o \Sigma_o V_o^T + U_e \Sigma_e V_e^T) W W^T (V_o \Sigma_o U_o^T + V_e \Sigma_e U_e^T) P. \quad (23)$$

Let W_0 be a valid solution as

$$W_0 = \begin{cases} W_0 & \text{s.t. } V_k^T W_0 = \Sigma_k^{-2} \text{ if } \text{rank}(X) = k \\ [W_r; 0] & \text{s.t. } V_r^T W_r = \Sigma_r^{-2} \text{ if } \text{rank}(X) = r \leq k \end{cases}. \quad (24)$$

By the above definition, $V_e^T W = 0$. Therefore, with $\text{rank}(X) > k$,

$$\begin{aligned} & \|HH^T - P\| \\ &= \|U_k \Sigma_k^3 V_k W_0 W_0^T V_k \Sigma_k^3 U_k^T \\ & \quad + EU_k \Sigma_k V_k W_0 W_0^T V_k \Sigma_k^3 U_k^T \\ & \quad + U_k \Sigma_k^3 V_k W_0 W_0^T V_k \Sigma_k U_k^T E \\ & \quad + EU_k \Sigma_k V_k W_0 W_0^T V_k \Sigma_k U_k^T E - P\| \\ &= \|U_k \Sigma_k^2 U_k^T + EU_k U_k^T + U_k U_k^T E + EU_k \Sigma_k^{-2} U_k^T E \\ & \quad - U_k \Sigma_k^2 U_k^T - U_e \Sigma_e^2 U_e^T - E\| \\ &= \|EU_k U_k^T + U_k U_k^T E + EU_k \Sigma_k^{-2} U_k^T E \\ & \quad - U_e \Sigma_e^2 U_e^T - E\| \\ &\leq \|EU_k U_k^T + U_k U_k^T E - E\| + \|EU_k \Sigma_k^{-2} U_k^T E\| \\ & \quad + \|U_e \Sigma_e^2 U_e^T\|. \end{aligned} \quad (25)$$

Similarly, if $\text{rank}(X) = r \leq k$,

$$\|HH^T - P\| = \|EU_r U_r^T + U_r U_r^T E + EU_r \Sigma_r^{-2} U_r^T E - E\| \quad (26)$$

$$\leq \|EU_r U_r^T + U_r U_r^T E - E\| + \|EU_r \Sigma_r^{-2} U_r^T E\|. \quad (27)$$

Now we focus on the general case, $\text{rank}(X) > k$ and the conclusion can be easily extended into the low-rank case. Note that

$$\|E(U_k U_k^T - \frac{1}{2}I)\| = \|(U_k U_k^T - \frac{1}{2}I)E\|, \quad (28)$$

and the first term can be written as

$$\begin{aligned} & \|EU_k U_k^T + U_k U_k^T E - E\|^2 \\ &= \|E(U_k U_k^T - \frac{1}{2}I) + (U_k U_k^T - \frac{1}{2}I)E\|^2 \\ &= \|E(U_k U_k^T - \frac{1}{2}I)\|^2 + \|(U_k U_k^T - \frac{1}{2}I)E\|^2 \\ & \quad + 2\langle E(U_k U_k^T - \frac{1}{2}I), (U_k U_k^T - \frac{1}{2}I)E \rangle \\ &= 2\|E(U_k U_k^T - \frac{1}{2}I)\|^2 + 2s\|E(U_k U_k^T - \frac{1}{2}I)\|^2 \\ &= 2(1+s)\|E(U_k U_k^T - \frac{1}{2}I)\|^2, \end{aligned}$$

where

$$s = \frac{\langle E(U_k U_k^T - \frac{1}{2}I), (U_k U_k^T - \frac{1}{2}I)E \rangle}{\|E(U_k U_k^T - \frac{1}{2}I)\|^2}. \quad (29)$$

Due to that

$$\begin{aligned} & \|E(U_k U_k^T - \frac{1}{2}I)\|^2 = \text{tr}(E^2(U_k U_k^T - \frac{1}{2}I)) \\ &= \text{tr}(E^2(U_k U_k^T U_k U_k^T - 2 \times \frac{1}{2}U_k U_k^T + \frac{1}{4}I)) \\ &= \frac{1}{4}\text{tr}(E^2) = \frac{1}{4}\varepsilon^2, \end{aligned}$$

we have

$$\begin{aligned} & \|EU_k U_k^T + U_k U_k^T E - E\| \\ &= \sqrt{2(1+s)}\|E(U_k U_k^T - \frac{1}{2}I)\| = \sqrt{\frac{1+s}{2}}\varepsilon. \end{aligned} \quad (30)$$

Let $Q = (U_k U_k^T - I/2)$ and s can be reformulated as

$$s = \cos(EQ, QE), \quad (31)$$

and we have the following definition

$$\theta_* = \arccos(s) = \theta(EQ, QE). \quad (32)$$

According to Lemma 7.1, Assumption 4.1 indicates that

$$E(U_k U_k^T - \frac{1}{2}I) = (U_k U_k^T - \frac{1}{2}I)E \Rightarrow s < 1. \quad (33)$$

And therefore, $\sqrt{(1+s)/2} \in [0, 1]$. Let $\delta = 1 - \sqrt{(1+s)/2} = 1 - \cos(\theta_*/2) > 0$ and the above equation can be reformulated as

$$\|EU_k U_k^T + U_k U_k^T E - E\| = (1-\delta)\varepsilon. \quad (34)$$

The second term can be formulated as

$$\|EU_k \Sigma_k^{-2} U_k^T E\| \leq \|U_k \Sigma_k^{-2} U_k^T\| \varepsilon^2 \quad (35)$$

$$= \varepsilon^2 \sqrt{\text{tr}(U_k \Sigma_k^{-4} U_k^T)} \quad (36)$$

$$= \varepsilon^2 \sqrt{\text{tr}(\Sigma^{-4} U_k^T U_k)} \quad (37)$$

$$= \varepsilon^2 \sqrt{\sum_{i=1}^k \frac{1}{\sigma_i^4}} \leq \sqrt{r} \frac{\varepsilon^2}{\sigma_k^2}, \quad (38)$$

while the third term is

$$\|U_e \Sigma_e^2 U_e^T\| = \|\Sigma_e^2\| = (\sum_{i=k+1}^n \sigma_i^4)^{1/2} \leq \sqrt{n-k} \sigma_*^2. \quad (39)$$

To sum up, the error of \mathcal{M} is bounded as

$$\|\mathbf{H}\mathbf{H}^T - \mathbf{P}\| \leq (1 - \delta)\varepsilon + \sqrt{k}\frac{\varepsilon^2}{\sigma_k^2} + \sqrt{e}\sigma_*^2. \quad (40)$$

If

$$\sqrt{k}\frac{\varepsilon}{\sigma_k^2} \leq \frac{\delta}{2} \Rightarrow \varepsilon \leq \frac{\delta\sigma_k^2}{2\sqrt{k}} \quad (41)$$

and

$$\sqrt{n - k}\sigma_*^2 \leq \frac{\delta}{2} \Rightarrow \sigma_* \leq \sqrt{\frac{\delta}{2\sqrt{n - k}}}, \quad (42)$$

then $\|\mathbf{H}\mathbf{H}^T - \mathbf{P}\| \leq \varepsilon$. In other words, when $\varepsilon \leq \mathcal{O}(\delta)$ and $\sigma_* \leq \mathcal{O}(\sqrt{\delta})$, the error will be bounded by ε .

For the case that $\text{rank}(\mathbf{X}) = r < k$, it is not hard to verify that

$$\begin{aligned} & \|\mathbf{H}\mathbf{H}^T - \mathbf{P}\| \\ & \leq \|\mathbf{E}\mathbf{U}_r\mathbf{U}_r^T + \mathbf{U}_r\mathbf{U}_r^T\mathbf{E} - \mathbf{E}\| + \|\mathbf{E}\mathbf{U}_r\mathbf{\Sigma}_r^{-2}\mathbf{U}_r^T\mathbf{E}\| \end{aligned} \quad (43)$$

$$\leq (1 - \delta)\varepsilon + \sqrt{r}\frac{\varepsilon^2}{\sigma_r^2}. \quad (44)$$

As

$$\sqrt{r}\frac{\varepsilon}{\sigma_r^2} \leq \frac{\delta}{2} \Rightarrow \varepsilon \leq \frac{\delta\sigma_r^2}{2\sqrt{r}} = \mathcal{O}(\delta) \quad (45)$$

and $\sigma_* = 0 \leq \mathcal{O}(\sqrt{\delta})$, we get $\|\mathbf{H}\mathbf{H}^T - \mathbf{P}\| \leq \varepsilon$. Hence, we have

$$\min_{\mathbf{W}} \ell(\mathbf{P}, \mathbf{X}, \mathbf{W}) \leq \ell(\mathbf{P}, \mathbf{X}, \mathbf{W}_0) \leq \varepsilon, \quad (46)$$

and the theorem is proved. \square

7.3 Proof of Theorem 4.2

Proof. According to Ineq. (25),

$$\begin{aligned} \|\mathbf{H}\mathbf{H}^T - \mathbf{P}\| & \leq \|\mathbf{E}\mathbf{U}_k\mathbf{U}_k^T + \mathbf{U}_k\mathbf{U}_k^T\mathbf{E} - \mathbf{E}\| \\ & \quad + \|\mathbf{E}\mathbf{U}_k\mathbf{\Sigma}_k^{-2}\mathbf{U}_k^T\mathbf{E}\| + \|\mathbf{U}_e\mathbf{\Sigma}_e^2\mathbf{U}_e^T\|. \end{aligned} \quad (47)$$

Suppose that $\mathbf{E} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ so that $\mathbf{P} = \mathbf{U}(\mathbf{S} + \mathbf{\Lambda})\mathbf{U}^T$ where $\mathbf{S} = \mathbf{\Sigma}\mathbf{\Sigma}^T = \text{diag}(\mathbf{\Sigma}_r^2; \mathbf{0})$. So, we have

$$\begin{aligned} \mathbf{P} & = \mathbf{U}(\mathbf{S} + \mathbf{\Lambda})\mathbf{U}^T \\ & = [\mathbf{U}_r, \mathbf{U}_e] \begin{bmatrix} \mathbf{\Sigma}_r^2 + \mathbf{\Lambda}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{\Sigma}_e^2 + \mathbf{\Lambda}_e \end{bmatrix} \begin{bmatrix} \mathbf{U}_r^T \\ \mathbf{U}_e^T \end{bmatrix} \\ & = \mathbf{U}_r(\mathbf{\Sigma}_r^2 + \mathbf{\Lambda}_r)\mathbf{U}_r^T + \mathbf{U}_e(\mathbf{\Sigma}_e^2 + \mathbf{\Lambda}_e)\mathbf{U}_e^T. \end{aligned}$$

Let $\mathbf{V}_k\mathbf{W} = (\mathbf{\Sigma}_k^3 + \mathbf{\Lambda}_k\mathbf{\Sigma}_k)(\mathbf{\Sigma}_k^2 + \mathbf{\Lambda}_k)^{1/2}$ and we have

$$\begin{aligned} & \|\mathbf{P} - \mathbf{H}\mathbf{H}^T\| \\ & = \|\mathbf{U}(\mathbf{S} + \mathbf{\Lambda})\mathbf{\Sigma}\mathbf{V}^T\mathbf{W}\mathbf{W}^T\mathbf{V}\mathbf{\Sigma}^T(\mathbf{S} + \mathbf{\Lambda})\mathbf{U}^T - \mathbf{P}\| \\ & = \|\mathbf{U}_k(\mathbf{\Sigma}_k^3 + \mathbf{\Lambda}_k\mathbf{\Sigma}_k)\mathbf{V}_k^T\mathbf{W}\mathbf{W}^T\mathbf{V}_k(\mathbf{\Sigma}_k^3 + \mathbf{\Lambda}_k\mathbf{\Sigma}_k)\mathbf{U}_k^T - \mathbf{P}\| \\ & = \|\mathbf{U}_k\hat{\mathbf{I}}(\mathbf{\Sigma}_k^2 + \mathbf{\Lambda}_k)\hat{\mathbf{I}}\mathbf{U}_k^T - \mathbf{U}_k(\mathbf{\Sigma}_k^2 + \mathbf{\Lambda}_k)\mathbf{U}_k^T \\ & \quad - \mathbf{U}_e(\mathbf{\Sigma}_e^2 + \mathbf{\Lambda}_e)\mathbf{U}_e^T\| \\ & = \|\mathbf{U}_e(\mathbf{\Sigma}_e^2 + \mathbf{\Lambda}_e)\mathbf{U}_e^T\| \leq \|\mathbf{\Sigma}_e^2\| + \|\mathbf{\Lambda}_e\| \\ & \leq \varepsilon + \mathcal{O}(\sigma_*^2), \end{aligned}$$

where $\hat{\mathbf{I}} = (\mathbf{\Sigma}_k^3 + \mathbf{\Lambda}_k\mathbf{\Sigma}_k)(\mathbf{\Sigma}_k^3 + \mathbf{\Lambda}_k\mathbf{\Sigma}_k)^\dagger = \mathbf{\Sigma}_k(\mathbf{\Sigma}_k^2 + \mathbf{\Lambda}_k)(\mathbf{\Sigma}_k^2 + \mathbf{\Lambda}_k)^\dagger\mathbf{\Sigma}_k^{-1}$. Clearly, $\hat{\mathbf{I}} = \mathbb{1}\{\mathbf{\Sigma}_r^2 + \mathbf{\Lambda}_r \neq \mathbf{0}\}$. Therefore, Hence, the theorem is proved. \square

Corollary 7.1. *If Assumption 4.1 does not hold and $\text{rank}(\mathbf{X}) \leq k$, then there exists $\mathbf{W} \in \mathbb{R}^{d \times k}$ so that $\|\mathbf{P} - \mathbf{H}\mathbf{H}\| \leq \varepsilon$.*

Proof. Suppose that $\mathbf{E} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ so that $\mathbf{P} = \mathbf{U}(\mathbf{S} + \mathbf{\Lambda})\mathbf{U}^T$ where $\mathbf{S} = \mathbf{\Sigma}\mathbf{\Sigma}^T = \text{diag}(\mathbf{\Sigma}_r^2; \mathbf{0})$. Then we have

$$\begin{aligned} \mathbf{P} & = \mathbf{U}(\mathbf{S} + \mathbf{\Lambda})\mathbf{U}^T = [\mathbf{U}_r, \mathbf{U}_e] \begin{bmatrix} \mathbf{\Sigma}_r^2 + \mathbf{\Lambda}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{\Lambda}_e \end{bmatrix} \begin{bmatrix} \mathbf{U}_r^T \\ \mathbf{U}_e^T \end{bmatrix} \\ & = \mathbf{U}_r(\mathbf{\Sigma}_r^2 + \mathbf{\Lambda}_r)\mathbf{U}_r^T + \mathbf{U}_e\mathbf{\Lambda}_e\mathbf{U}_e^T. \end{aligned} \quad (48)$$

Let $\mathbf{W} = [\mathbf{W}_r; \mathbf{0}]$ subjected to $\mathbf{V}_r^T\mathbf{W}_r = (\mathbf{\Sigma}_r^3 + \mathbf{\Lambda}_r\mathbf{\Sigma}_r)^\dagger(\mathbf{\Sigma}_r^2 + \mathbf{\Lambda}_r)^{1/2}$. Then

$$\begin{aligned} & \|\mathbf{H}\mathbf{H}^T - \mathbf{P}\| \\ & = \|\mathbf{U}(\mathbf{S} + \mathbf{\Lambda})\mathbf{\Sigma}\mathbf{V}^T\mathbf{W}\mathbf{W}^T\mathbf{V}\mathbf{\Sigma}^T(\mathbf{S} + \mathbf{\Lambda})\mathbf{U}^T - \mathbf{P}\| \\ & = \|\mathbf{U}_r(\mathbf{\Sigma}_r^3 + \mathbf{\Lambda}_r\mathbf{\Sigma}_r)\mathbf{V}_r^T\mathbf{W}\mathbf{W}^T\mathbf{V}_r(\mathbf{\Sigma}_r^3 + \mathbf{\Lambda}_r\mathbf{\Sigma}_r)\mathbf{U}_r^T - \mathbf{P}\| \\ & = \|\mathbf{U}_r\hat{\mathbf{I}}(\mathbf{\Sigma}_r^2 + \mathbf{\Lambda}_r)\hat{\mathbf{I}}\mathbf{U}_r^T - \mathbf{U}_r(\mathbf{\Sigma}_r^2 + \mathbf{\Lambda}_r)\mathbf{U}_r^T - \mathbf{U}_e\mathbf{\Lambda}_e\mathbf{U}_e^T\|, \end{aligned}$$

where $\hat{\mathbf{I}} = (\mathbf{\Sigma}_r^3 + \mathbf{\Lambda}_r\mathbf{\Sigma}_r)(\mathbf{\Sigma}_r^3 + \mathbf{\Lambda}_r\mathbf{\Sigma}_r)^\dagger = \mathbf{\Sigma}_r(\mathbf{\Sigma}_r^2 + \mathbf{\Lambda}_r)(\mathbf{\Sigma}_r^2 + \mathbf{\Lambda}_r)^\dagger\mathbf{\Sigma}_r^{-1}$. Clearly, $\hat{\mathbf{I}} = \mathbb{1}\{\mathbf{\Sigma}_r^2 + \mathbf{\Lambda}_r \neq \mathbf{0}\}$. Therefore,

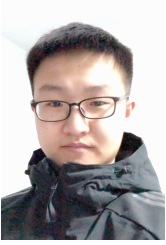
$$\|\mathbf{H}\mathbf{H}^T - \mathbf{P}\| = \|\mathbf{U}_e\mathbf{\Lambda}_e\mathbf{U}_e^T\| = \|\mathbf{\Lambda}_e\| \leq \varepsilon. \quad (49)$$

Hence, the corollary is proved. \square

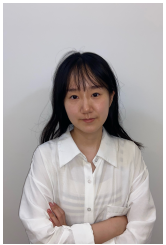
REFERENCES

- [1] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [3] 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.
- [4] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *ICLR*, 2017.
- [5] M. Niepert, M. Ahmed, and K. Kutzkov, "Learning convolutional neural networks for graphs," in *International Conference on Machine Learning*, 2016, pp. 2014–2023.
- [6] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in *6th International Conference on Learning Representations*, 2018.
- [7] D. Kreuzer, D. Beaini, W. Hamilton, V. Létourneau, and P. Tossou, "Rethinking graph transformers with spectral attention," in *Advances in Neural Information Processing Systems*, vol. 34, 2021, pp. 21 618–21 629.
- [8] X. Wang, Y. Ma, Y. Wang, W. Jin, X. Wang, J. Tang, C. Jia, and J. Yu, "Traffic flow prediction via spatial temporal graph neural network," in *Proceedings of The Web Conference 2020*, 2020, pp. 1082–1092.
- [9] T. N. Kipf and M. Welling, "Variational graph auto-encoders," *arXiv preprint arXiv:1611.07308*, 2016.
- [10] J. Park, M. Lee, H. J. Chang, K. Lee, and J. Y. Choi, "Symmetric graph convolutional autoencoder for unsupervised graph representation learning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6519–6528.
- [11] K. Hassani and A. H. Khasahmadi, "Contrastive multi-view representation learning on graphs," in *International Conference on Machine Learning*, PMLR, 2020, pp. 4116–4126.
- [12] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization," *Journal of machine learning research*, vol. 12, no. 7, 2011.
- [13] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015.
- [14] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Advances in Neural Information Processing Systems*, 2017, pp. 1024–1034.
- [15] J. Chen, T. Ma, and C. Xiao, "Fastgcn: Fast learning with graph convolutional networks via importance sampling," *arXiv preprint arXiv:1801.10247*, 2018.

- [16] J. Chen, J. Zhu, and L. Song, "Stochastic training of graph convolutional networks with variance reduction," in *International Conference on Machine Learning*, 2018, pp. 942–950.
- [17] H. Zeng, H. Zhou, A. Srivastava, R. Kannan, and V. K. Prasanna, "Graphsaint: Graph sampling based inductive learning method," in *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- [18] W.-L. Chiang, X. Liu, S. Si, Y. Li, S. Bengio, and C.-J. Hsieh, "Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 257–266.
- [19] F. Wu, A. H. S. Jr., T. Zhang, C. Fifty, T. Yu, and K. Q. Weinberger, "Simplifying graph convolutional networks," in *Proceedings of the 36th International Conference on Machine Learning*, vol. 97, 2019, pp. 6861–6871.
- [20] H. Zhu and P. Koniusz, "Simple spectral graph convolution," in *9th International Conference on Learning Representations*, 2021.
- [21] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How powerful are graph neural networks?" in *International Conference on Learning Representations*, 2019.
- [22] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P. Manzagol, "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion," *Journal of Machine Learning Research*, vol. 11, pp. 3371–3408, 2010.
- [23] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, "Spectral networks and locally connected networks on graphs," in *2nd International Conference on Learning Representations*, 2014.
- [24] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *Advances in Neural Information Processing Systems*, 2016, pp. 3844–3852.
- [25] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Gallagher, and T. Eliassi-Rad, "Collective classification in network data," *AI Mag.*, vol. 29, no. 3, pp. 93–106, 2008.
- [26] C. Zhang, Y. Wang, L. Zhu, J. Song, and H. Yin, "Multi-graph heterogeneous interaction fusion for social recommendation," *ACM Transactions on Information Systems (TOIS)*, vol. 40, no. 2, pp. 1–26, 2021.
- [27] B. Qian, Y. Wang, H. Yin, R. Hong, and M. Wang, "Switchable online knowledge distillation," in *European Conference on Computer Vision*, 2022, pp. 449–466.
- [28] B. Qian, Y. Wang, R. Hong, and M. Wang, "Adaptive data-free quantization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 7960–7968.
- [29] D. Guo, Y. Shao, Y. Cui, Z. Wang, L. Zhang, and C. Shen, "Graph attention tracking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 9543–9552.
- [30] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017*, 2017, pp. 5998–6008.
- [31] B.-H. Kim, J. C. Ye, and J.-J. Kim, "Learning dynamic graph representation of brain connectome with spatio-temporal attention," in *Advances in Neural Information Processing Systems*, vol. 34. Curran Associates, Inc., 2021, pp. 4314–4327.
- [32] Q. Li, Z. Han, and X.-M. Wu, "Deeper insights into graph convolutional networks for semi-supervised learning," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [33] G. Li, M. Muller, A. Thabet, and B. Ghanem, "Deepgcns: Can gcns go as deep as cnns?" in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9267–9276.
- [34] K. Oono and T. Suzuki, "Graph neural networks exponentially lose expressive power for node classification," in *International Conference on Learning Representations*, 2020.
- [35] W. Cong, M. Ramezani, and M. Mahdavi, "On provable benefits of depth in training graph convolutional networks," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [36] C. Morris, M. Ritzert, M. Fey, W. L. Hamilton, J. E. Lenssen, G. Rattan, and M. Grohe, "Weisfeiler and leman go neural: Higher-order graph neural networks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019, pp. 4602–4609.
- [37] B. Weisfeiler and A. Leman, "The reduction of a graph to canonical form and the algebra which appears therein," *Nauchno-Tekhnicheskaya Informatsia*, vol. 2, no. 9, pp. 12–16, 1968.
- [38] H. Zhang, J. Shi, R. Zhang, and X. Li, "Non-graph data clustering via o(n) bipartite graph convolution," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 7, pp. 8729–8742, 2023.
- [39] X. Li, H. Zhang, and R. Zhang, "Adaptive graph auto-encoder for general data clustering," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 12, pp. 9725–9732, 2022.
- [40] H. Zhang, Y. Zhu, and X. Li, "Toward projected clustering with aggregated mapping," *IEEE Transactions on Image Processing*, vol. 32, pp. 4103–4113, 2023.
- [41] Y. You, T. Chen, Z. Wang, and Y. Shen, "L2-GCN: layer-wise and learned efficient training of graph convolutional networks," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, 2020, pp. 2124–2132.
- [42] Y. Wang, J. Tang, Y. Sun, and G. Wolf, "Decoupled greedy learning of graph neural networks," in *Optimization for Machine Learning*, 2020.
- [43] G. E. Hinton and R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [44] C. Wang, S. Pan, G. Long, X. Zhu, and J. Jiang, "Mgae: Marginalized graph autoencoder for graph clustering," in *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, 2017, pp. 889–898.
- [45] Y. Freund and R. E. Schapire, "A decision-theoretic generalization of on-line learning and an application to boosting," *J. Comput. Syst. Sci.*, vol. 55, no. 1, pp. 119–139, 1997.
- [46] J. H. Friedman, "Greedy function approximation: a gradient boosting machine," *Annals of statistics*, pp. 1189–1232, 2001.
- [47] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, B. Krishnapuram, M. Shah, A. J. Smola, C. C. Aggarwal, D. Shen, and R. Rastogi, Eds., pp. 785–794.
- [48] H. Schwenk and Y. Bengio, "Boosting neural networks," *Neural Comput.*, vol. 12, no. 8, pp. 1869–1887, 2000.
- [49] Z. Zhou, J. Wu, and W. Tang, "Ensembling neural networks: Many could be better than all," *Artif. Intell.*, vol. 137, no. 1-2, pp. 239–263, 2002.
- [50] S. Ivanov and L. Prokhorenkova, "Boost then convolve: Gradient boosting meets graph neural networks," in *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- [51] K. Sun, Z. Zhu, and Z. Lin, "Adagcn: Adaboosting graph convolutional networks into deep models," in *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- [52] M. Jaderberg, W. M. Czarnecki, S. Osindero, O. Vinyals, A. Graves, D. Silver, and K. Kavukcuoglu, "Decoupled neural interfaces using synthetic gradients," in *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, vol. 70, 2017, pp. 1627–1635.
- [53] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *nature*, vol. 323, no. 6088, pp. 533–536, 1986.
- [54] K. Xu, C. Li, Y. Tian, T. Sonobe, K. Kawarabayashi, and S. Jegelka, "Representation learning on graphs with jumping knowledge networks," in *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, vol. 80, 2018, pp. 5449–5458.
- [55] S. Pan, R. Hu, G. Long, J. Jiang, L. Yao, and C. Zhang, "Adversarially regularized graph autoencoder for graph embedding," in *IJCAI*, 2018, pp. 2609–2615.
- [56] X. Zhang, H. Liu, Q. Li, and X.-M. Wu, "Attributed graph clustering via adaptive graph convolution," in *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, 2019, pp. 4327–4333.
- [57] P. Velickovic, W. Fedus, W. L. Hamilton, P. Liò, Y. Bengio, and R. D. Hjelm, "Deep graph infomax," in *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, 2019.
- [58] J. Klicpera, A. Bojchevski, and S. Günnemann, "Predict then propagate: Graph neural networks meet personalized pagerank," in *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, 2019.
- [59] M. Chen, Z. Wei, Z. Huang, B. Ding, and Y. Li, "Simple and deep graph convolutional networks," in *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, vol. 119, 2020, pp. 1725–1735.



Hongyuan Zhang received the B.E. degree in software engineering from Xidian University, Xi'an, China in 2019 and received the Ph.D. degree from the School of Computer Science and the School of Artificial Intelligence, Optics and Electronics (iOPEN), Northwestern Polytechnical University, Xi'an, China in 2024.



Yanan Zhu received the B.E. degree in computer science and technology from Nanjing University of Aeronautics and Astronautics, Nanjing, China, in 2021 and received the master's degree with School of Computer Science and School of Artificial Intelligence, Optics and Electronics (iOPEN), Northwestern Polytechnical University, Xi'an, China in 2024.

Xuelong Li (M'02-SM'07-F'12) is the Chief Technology Officer (CTO) and Chief Scientist of the China Telecom.