

Collaborative Graph Neural Networks for Attributed Network Embedding

Qiaoyu Tan, Xin Zhang, Xiao Huang, Hao Chen, Jundong Li, and Xia Hu

Abstract—Graph neural networks (GNNs) have shown prominent performance on attributed network embedding. However, existing efforts mainly focus on exploiting network structures, while the exploitation of node attributes is rather limited as they only serve as node features at the initial layer. This simple strategy impedes the potential of node attributes in augmenting node connections, leading to limited receptive field for inactive nodes with few or even no neighbors. Furthermore, the training objectives (i.e., reconstructing network structures) of most GNNs also do not include node attributes, although studies have shown that reconstructing node attributes is beneficial. Thus, it is encouraging to deeply involve node attributes in the key components of GNNs, including graph convolution operations and training objectives. However, this is a nontrivial task since an appropriate way of integration is required to maintain the merits of GNNs. To bridge the gap, in this paper, we propose COllaborative graph Neural Networks—CONN, a tailored GNN architecture for attribute network embedding. It improves model capacity by 1) selectively diffusing messages from neighboring nodes and involved attribute categories, and 2) jointly reconstructing node-to-node and node-to-attribute-category interactions via cross-correlation. Experiments on real-world networks demonstrate that CONN excels state-of-the-art embedding algorithms with a great margin.

Index Terms—Attributed Network Embedding, Graph Neural Networks, Collaborative Aggregation, Cross-Correlation



1 INTRODUCTION

ATTRIBUTED networks [1], [2], [3], [4] are ubiquitous in a myriad of real-world information systems, such as academic networks and social medial systems. Unlike plain networks in which only node-to-node interactions are available, each node in attributed network is associated with a rich set of attributes, describing its distinctive characteristics. For example, in social networks, users connect with others as friends, share opinions, and post comments as attributes. In academic citation networks, different articles are connected via citation links, and each article has substantial text information, such as an abstract sentence to describe its own topic. Several studies [5], [6] in social science have revealed that attributes of nodes can reflect and affect their community structures [7] in practice. Thus, it is encouraging and important to study attributed networks. To this end, attributed network embedding [1], [8], [9], targeting at leveraging both network proximity and node attribute affinity to learn low-dimensional node representations, has attracted great attention in recent years, and many efforts have been devoted from both academia and industry [10], [11], [12], [13], [14], [15], [16].

Among them, embedding paradigm based on graph neural networks (GNNs) [15], [17], [18], [18], [19], [20], [21], [22], [23], [24], [25], [26], [27], [28], [29], [30] has achieved remarkable success over a variety of downstream graph analytical tasks, including node classification [31], [32], [33],

graph classification [34], [35], link prediction [36], [37], [38], [39], node clustering [40], [41], and anomaly detection [42], [43], [44], [45]. The design recipes for GNNs based methods include two major components: 1) a GNN encoder that takes node attributes and node-to-node interaction network as input and outputs low-dimensional node representations; 2) the training objective, which is derived to reconstruct the input data (e.g., network structure), so as to train the model unsupervised. Thanks to the critical message propagation mechanism in GNNs, the GNN encoder is naturally applicable for attributed networks. Therefore, existing GNNs-based embedding efforts mainly focus on advancing model performance with more expressive message passing schema, such as adaptively aggregating messages from neighbors via the attention layer [46], [47], [48].

Despite their popularity and recent advances in refining GNNs architectures to effectively model the topological structures [49], [50], [51], [52], [53], little attention was paid to the node attributes. Prior GNNs studies focus on updating the representation of each node by aggregating the representations of itself and its neighbors [54] recursively. In this learning process, node attributes are merely employed as the representations of nodes in the initial layer [55]. They would be blocked from message propagation if the network structure is incomplete or missing, which is quite common in real-world applications where graphs exhibit long-tail node degree distribution [56], [57]. Besides, even when we design the training objectives of GNNs, the node attributes are seldom used. For instance, the common practice to train GNNs based embedding algorithm for attributed network embedding is to reconstruct the observed node interactions [15], by either employing a negative sampling based objective or directly recovering the whole input network structure [14]. As a summary, node attributes are not well exploited in the existing works.

Qiaoyu Tan is with the Department of Computer Science and Engineering, Texas A&M University, Texas, USA. E-mail: qytan@tamu.edu.

Xin Zhang, Xiao Huang, and Hao Chen are with the Department of Computing, The Hong Kong Polytechnic University, Hung Hom, Hong Kong. E-mail: xin12.zhang@connect.polyu.hk; xiaohuang@comp.polyu.edu.hk; sundaychenhao@gmail.com.

Jundong Li is with the Department of Electrical and Computer Engineering, University of Virginia, Virginia, USA. E-mail: jundong@virginia.edu.

Xia Hu is with the Department of Computer Science, Rice University, Texas, USA. E-mail: xia.hu@rice.edu.

Recently, there has been a revolution to rethink the value of node attributes in random-walk based embedding approach [58], [59]. The core idea is to redesign the crucial component—random-walk in an attribute aware fashion. Specifically, ANRL [60] refines the conditional probability that estimates the propensity score between the anchor node and its context by utilizing node attributes. FeatWalk [59] conducts an attribute-aware joint random walk to increase the diversity of generated walks. According to their empirical experiments, both of them have shown significant performance gains compared with their vanilla counterparts. Nevertheless, they are tailored for random-walk based approaches and cannot be directly applied to GNN models with trivial efforts. Motivated by this, in this paper, we propose to explore whether node attributes can be effectively employed to advance the essential building blocks of GNN architectures (i.e., message aggregation mechanism and training objective).

However, it is a non-trivial and challenging task to integrate node attributes into GNN architectures mainly because of two reasons: (i) GNNs already show promising results by integrating node attributes as the initial node representations. It is difficult to further incorporate node attributes into the key components of GNNs, while maintaining the existing benefits and prominent performance [54], [61]; (ii) real-world node attributes, such as comments of users, abstracts of papers, and descriptions of products, are distinct from the network topological structures and not in line with the graph convolutional operation in GNNs. Specifically, the values in node attributes are often multi-categorical or continuous variables, while these in the network structures are binary. They are not compatible with each other. Thus a tailored operation is required to learn from them jointly. For example, employing autoencoders as the training objective of GNNs would achieve sub-optimal performance [14], [62].

To address the aforementioned challenges, we propose a novel unsupervised representation learning model, dubbed **COL**laborative Graph Neural Network (CONN). It aims to develop a tailored GNN architecture for attributed networks, such that node attributes can be explicitly fused into the message aggregation process as well as the training objective. Specifically, we aim to investigate two important research questions. (i) How to leverage node attributes to explicitly guide the message propagation of vanilla GNN, and conduct a collaborative aggregation mechanism? (ii) In the training objective, how to effectively model the heterogeneous interactions, so as to jointly reconstruct the network structure and node attributes? We summarize our major contributions as follows.

- We focus on unsupervised representation learning on attributed networks and propose an effective GNN framework - CONN, to leverage node attributes in the two aforementioned key components of vanilla GNNs.
- By conducting a bipartite graph on node attributes, we develop a collaborative aggregation mechanism for node embedding. It not only helps to enrich or rebuild node connections through attribute category, but also provides a principled way to update node representation from both neighboring nodes and involved attribute categories.

- Based on the node and attribute category representations, we design a novel cross-correlation layer to effectively model the complex node or node attribute category interactions. It highlights the similarity of two anchor nodes from their multi-granularity features and significantly boosts the reconstruction capability of vanilla GNNs.
- We evaluate CONN on node classification and link prediction tasks. Empirical results on benchmark datasets show that CONN performs consistently better than other state-of-the-art embedding methods. Moreover, we also analyze the robustness and convergence speed of CONN in Section 5.8.

2 RELATED WORK

There are three types of related works, based on whether node attributes are modeled for network embedding or not. In this section, we briefly review some related works in these two fields. Please refer to [63], [64], [65], [66], [67] for comprehensive review.

2.1 Network embedding

The first class of approach can be tracked back to traditional graph machine learning problem, which aims to learn node embedding while preserving the local manifold structure, such as the LPP [68] and Laplacian Eigenmaps [69]. Nevertheless, these methods suffer from scalability challenges to large-scale networks, due to the time-expensive eigendecomposition operation on the adjacency matrix, whose time complexity is $O(n^3)$ with n is the number of nodes. Inspired by the recent breakthrough of distributed representation learning [70] in natural language processing, a lot of scalable network embedding methods [58], [71], [72], [73], [74], [75], [76], [77] have been developed. For example, DeepWalk [58] and Node2vec [73] conduct truncated random walks on the network to generate node sequences, and then feed them into the Skip-Gram algorithm [70] for node embedding. LINE [71] optimizes the first and second order neighborhood associations to learn representations. GraRep [72] extends LINE to capture high-order neighborhood relationships. SDNE [74] applies deep learning for node embedding and targets to capture the non-linear graph structure as well as preserve the global and local structures of the graph. [78] attempts to incorporate the community structure of the graph for representation learning. Some other studies aim to deal with large-scale graphs [79], [80], [81]. Despite their simplicity, the aforementioned methods may be limited in practice, since they cannot exploit side information, such as user profiles, worlds in posts, and contexts in photos on social media.

2.2 Attributed Network Embedding

Different from traditional network embedding, attributed network embedding [31], [82], [83], [84], [85], [86], [87], [88], which aims to learn node representations by leveraging both the network structure and node attributes, has attracted substantial attention in recent years. For instance, ANRL [60] incorporates the node attributes into the conditional probability function, which predicts the propensity

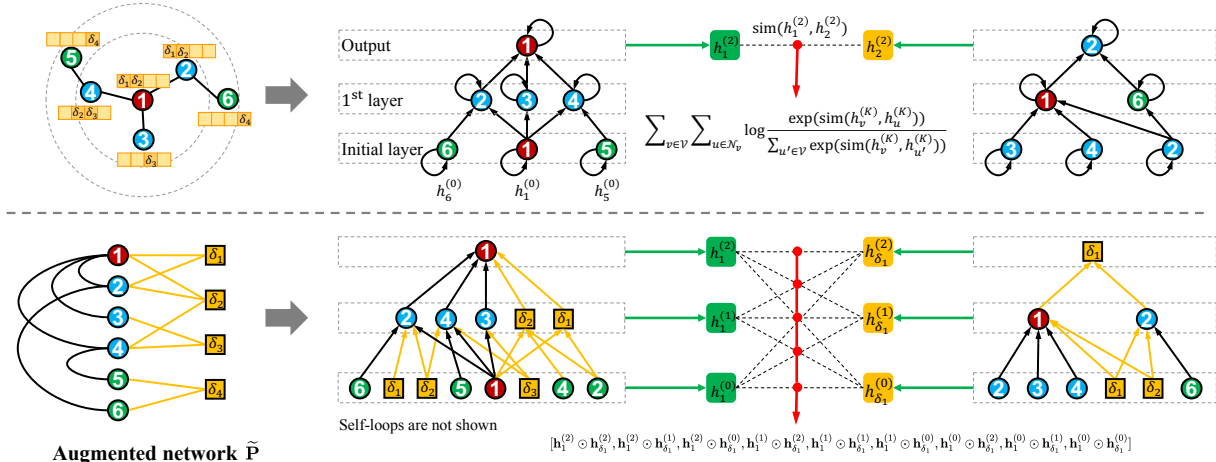


Fig. 1. Instead of inferring the node-to-node links by using node attributes as side features, CONN targets at jointly learning the network structure and node attributes in a unified latent space under graph neural network framework.

scores between an anchor node and its context, of Skip-gram to capture structure correlation. MUSAE [13] advances Skip-gram model by considering multi-scale neighborhood relationships based on node attributes. FeatWalk [59] aims to conduct attribute-aware random walks to increase the diversity of generated walks, so as to boost the Skip-gram model. PANE [89] is another random walk-based method for scalable training. TADW [83] incorporates node attributes into DeepWalk under matrix factorization framework. PTE [90] employs different orders of world co-occurrence relationships and node label to generate predictive text representations. ProGAN [91] aims to preserve the underlying proximities in the hidden space based on a generative adversarial network. Although the aforementioned methods are capable of employing the network structure and node attributes for joint representation learning, they are limited in capturing the structure information.

More recently, graph neural networks [23] based attributed network embedding has attracted increasing attention, due to its ability in capturing structure information, incorporating node attributes, and modeling non-linear relationships. As a pioneering GNN work, GCN [23] formally suggests to update node presentation by recursively aggregating representations from its adjacent neighbors, and achieves significant performance improvement on semi-supervised classification task. To follow up, GAE [14] makes the first effort to extend GCN to an autoencoder framework to learn node representations unsupervisedly. Meanwhile, GraphSage [15] improves the scalability of GCN by sampling and aggregating features from a node’s local neighborhood. Their results show that GCN is naturally suitable for attributed networks, since it can directly utilize node attributes as initial node features in the first layer for training. Motivated by this, many follow-up studies [46], [50], [92] have been proposed to advance model capability by developing more expressive GNN architectures. Recently, some efforts have also been devoted to redefine the training objective of GCN to learn more effective node representations. CAN [62] and [93] advocate to jointly reconstruct the network structure and node attributes under variational autoencoder framework. DGI [94] and GIC [95] suggests to

learn node representations by maximizing the mutual information between local node representation and the graph summary representation. GCA [96] and MVGRL [97] seek to update node representations by maximizing the agreement of representations between different views generated by data augmentation. However, they can only explore the node attributes implicitly in the message propagation process of GNN. In this paper, we propose a principled GNN variant for attributed network embedding, which allows node attributes to explicitly guide the message propagation and training objectives.

2.3 Attributed Network Embedding with Extra Knowledge

In addition to plain attributed networks, including node features and homogeneous graph structure, some studies also utilize extra knowledge, such as knowledge graph (KG) [98], [99], [100], [101], to boost the model’s performance. For example, PGE [102] studies how to incorporate additional edge features. KGCN [103] and KGAT [104] investigate how to leverage external knowledge graphs (e.g., item knowledge graphs) to improve recommendation quality. However, these methods require additional efforts to generate high-quality information sources, such as edge features and knowledge graphs. In contrast, in this work, we focus on standard attributed network embedding, modeling on homogeneous graphs with pure node attributes, such as continuous features and categorical features.

3 PROBLEM STATEMENT

We assume that an attributed network denoted by $\mathcal{G} = (\mathcal{V}, \mathbf{A}, \mathbf{X})$ is given. It has n vertexes collected in set \mathcal{V} . These n nodes are connected by an undirected network with its adjacency matrix denoted as $\mathbf{A} \in \mathbb{R}^{n \times n}$. If there is an edge between nodes v_i and v_j , then $\mathbf{A}_{ij} = 1$; otherwise, $\mathbf{A}_{ij} = 0$. Besides the network \mathbf{A} , each node v also has a descriptive feature vector $\mathbf{x}_v \in \mathbb{R}^m$, known as node attributes, where m is the total number of attribute categories. We denote the neighbor set of node v as \mathcal{N}_v , and represent the neighbor set of attribute category δ_j as \mathcal{N}_{δ_j} ,

TABLE 1
Summary of notations in this paper.

Notation	Description
$\mathbf{A} \in \mathbb{R}^{n \times n}$	adjacency matrix
$\mathbf{X} \in \mathbb{R}^{n \times m}$	a matrix collects all node attributes
\mathcal{V}	a set collects all n nodes
\mathcal{U}	a set collects all m attribute categories
$\delta_j \in \mathcal{U}$	the j^{th} node attribute category
\mathcal{N}_v	a set collects adjacent neighbors of v
K	the number of graph convolutional layers
$\tilde{\mathbf{P}} \in \mathbb{R}^{(n+m) \times (n+m)}$	adjacency matrix of augmented network
$\mathbf{h}_v^{(k)} \in \mathbb{R}^{1 \times d}$	representation of node v at the k^{th} layer
$\mathbf{h}_{\delta_j}^{(k)} \in \mathbb{R}^{1 \times d}$	representation of δ_j at the k^{th} layer

i.e., $\mathcal{N}_{\delta_j} = \{u | \mathbf{X}_{uj} > 0, \text{ for } u \in \mathcal{V}\}$. We employ a diagonal matrix $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$ to denote the degree matrix, where $d_i = \sum_j \mathbf{A}_{ij}$. The main symbols are listed in Table 1. To study the GNNs in an unsupervised setting, we follow the literature [23], [105] and formally define the problem of attributed network embedding in Definition 1.

Definition 1. Attributed Network Embedding. Given an attributed network $\mathcal{G} = (\mathcal{V}, \mathbf{A}, \mathbf{X})$, the goal is to learn a d -dimensional continuous vector $\mathbf{h}_v \in \mathbb{R}^d$ for each node $v \in \mathcal{V}$, such that the topological structures in \mathbf{A} and side information characterized by \mathbf{X} could be preserved in the embedding representations \mathbf{H} . The performance of this learning task is evaluated by applying \mathbf{H} to various downstream tasks such as node classification and link prediction.

To perform attributed network embedding, GNN models [15], [23] learn the embedding representation of each node v by aggregating the representations of itself and its neighbors \mathcal{N}_v in the previous layer. A typical neighborhood aggregation mechanism [23] is expressed as below,

$$\mathbf{h}_v^{(k)} = \frac{1}{1 + d_v} \mathbf{h}_v^{(k-1)} + \sum_{u \in \mathcal{N}_v} \frac{\mathbf{A}_{vu}}{\sqrt{(1 + d_v)(1 + d_u)}} \mathbf{h}_u^{(k-1)}, \quad (1)$$

where $\mathbf{h}_v^{(k)}$ denotes the hidden representation of node v at the k^{th} layer of GNNs, and $\mathbf{h}_v^{(0)} = \mathbf{x}_v$. The final output is $\mathbf{h}_v^{(K)}$, where K is the maximum number of layers considered. The top subfigure in Figure 1 illustrates this traditional neighborhood aggregation by using a toy example. We could see that, by stacking two layers, the representations of two-hop neighbors 5 and 6 could be accessed by node $v = 1$. To train this model in an unsupervised manner, a widely-adopted [15] graph-based loss function is defined as,

$$\mathcal{L} = - \sum_{v \in \mathcal{V}} \sum_{u \in \mathcal{N}_v} \log \frac{\exp(\text{sim}(\mathbf{h}_v^{(K)}, \mathbf{h}_u^{(K)}))}{\sum_{u' \in \mathcal{V}} \exp(\text{sim}(\mathbf{h}_v^{(K)}, \mathbf{h}_{u'}^{(K)}))}. \quad (2)$$

$\text{sim}(\cdot, \cdot)$ is a similarity function, e.g., inner product. The goal is to make the representations of connected nodes similar to each other, while enforcing the representations of unconnected nodes to be distinct. We observe that, node attributes are employed only as the initial representations $\mathbf{h}_v^{(0)}$. They have not been further exploited, especially being integrated into the core mechanisms of GNNs.

4 COLLABORATIVE GRAPH CONVOLUTION

Node attributes are informative, and significantly correlated with and complementary to the network [60], [106], [107], [108], [109]. Since they have not been fully exploited in GNNs, we explore to deeply integrate node attributes into the core mechanisms of GNNs, and develop a novel framework named Collaborative graph Neural Network (CONN). Figure 1 depicts the two major components of CONN, i.e., collaborative neighborhood aggregation and collaborative training objective. First, we redefine the graph convolutions by considering node attribute categories \mathcal{U} as another set of nodes. As illustrated in Figure 1, we augment the original network \mathbf{A} to a new one $\tilde{\mathbf{P}}$ with $n + m = 6 + 4$ nodes. It preserves all edges in \mathbf{A} , and contains a link from nodes v to δ_j if v has a non-zero value in its node attributes \mathbf{X}_{vj} . Notice that \mathbf{X}_{vj} can be both positive or negative, where negative value means the neighbor has negative impact towards the anchor node. Based on $\tilde{\mathbf{P}}$, we perform neighborhood aggregation. For example, the first-order neighbors of node 1 have been augmented from $\{2, 3, 4\}$ to $\{2, 3, 4, \delta_1, \delta_2\}$, while the second-order neighbors of node 1 have been augmented from $\{1, 5, 6\}$ to $\{1, 2, 4, 5, 6, \delta_1, \delta_2, \delta_3\}$. We observe that our collaborative neighborhood aggregation could not only capture node-to-node interactions, but also node-to-attribute-category interactions. Second, to train our model, we design a collaborative loss. The goal is to collaboratively predict all links in the augmented network $\tilde{\mathbf{P}}$, which incorporates \mathbf{A} and \mathbf{X} . Additionally, we design a novel cross correlation mechanism to model the complex interactions between any pair of nodes in $\tilde{\mathbf{P}}$ (e.g., $1 \in \mathcal{V}$ and $\delta_1 \in \mathcal{U}$ in Figure 1). It employs not only the node representations in the last layer $\mathbf{h}_v^{(K)}$, but also all the remaining ones $\{\mathbf{h}_v^{(k)}, \text{ for } k = 0, 1, \dots, K - 1\}$. We now introduce the details in the following subsections.

4.1 Collaborative Neighborhood Aggregation

Given an attributed graph $\mathcal{G} = (\mathcal{V}, \mathbf{A}, \mathbf{X})$, existing GCN architectures mainly define the multiple-hop neighbors of node $v \in \mathcal{V}$ purely based on the network structure \mathbf{A} . The first-order neighbors of node v is represented by $\mathcal{N}_v = \{u | \mathbf{A}_{uv} = 1\}$, while the second-order neighbors is denoted by $\mathcal{N}_v^{(2)} = \{u | \mathbf{A}_{uv}^2 > 0\}$, so on and so forth. \mathbf{A}^k indicates the k^{th} power of \mathbf{A} . As a result, other than serving as the initial node representations, i.e., $\mathbf{h}_v^{(0)} = \mathbf{x}_v$, node attributes are excluded from the graph convolutions, i.e., the core operation of GNNs. As discussed in Introduction section, this could be suboptimal for inactive nodes with few or no neighbors in practice, since they don't have sufficient neighborhood information for GNNs to effectively learn their embedding representations [56].

4.1.1 Augmented network

To tackle the aforementioned issue, we propose to leverage the geometrical property of node attributes. We regard each node attribute category $\delta_j \in \mathcal{U}$ as a new node and node attributes \mathbf{X} as a weighted bipartite graph. By adding it into the original network \mathbf{A} , we would get an augmented

network, denoted as \mathbf{P} . Mathematically, its adjacency matrix is written as,

$$\mathbf{P} = \begin{bmatrix} \mathbf{A} & \mathbf{X} \\ \mathbf{X}^\top & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{(n+m) \times (n+m)}. \quad (3)$$

Eq.(9) is applicable for both categorical and continuous node attributes. For continuous features, however, it would generate a dense bipartite graph based on \mathbf{X} , which may significantly increase the computation costs when performing convolution on it. To save the computation, we empirically simplify the dense bipartite graph into a sparse one by only preserving the top- N values in each row of \mathbf{X} . An appropriate N value acts as a trade-off between the efficiency and accuracy, we analyze its impact in Section 5.8. In summary, we directly use the feature matrix \mathbf{X} of categorical attributes to construct the augmented graph, while adopting top- N values in each row of \mathbf{X} to generate a sparse graph for continuous features.

4.1.2 Augmented multi-hop neighbors

By using \mathbf{P} , the first-order and second-order neighbors of node v can be expanded as:

$$\begin{aligned} \mathcal{N}_v &= \{u | \mathbf{A}_{vu} = 1\} + \{\delta_j | \mathbf{X}_{vj} \neq 0\}, \\ \mathcal{N}_v^{(2)} &= \{u | \mathbf{A}_{vu}^2 > 0 \text{ or } (\mathbf{X}\mathbf{X}^\top)_{vu} > 0\} + \{\delta_j | (\mathbf{A}\mathbf{X})_{vj} > 0\}, \end{aligned} \quad (4)$$

where $(\mathbf{X}\mathbf{X}^\top) \in \mathbb{R}^{n \times n}$ collects all node pairs that share at least one node attribute category, i.e., all $v \rightarrow \delta_j \rightarrow u$ paths. $(\mathbf{A}\mathbf{X}) \in \mathbb{R}^{n \times m}$ implies node-to-attribute-category interactions reflected by $v \rightarrow u \rightarrow \delta_j$ paths. Compared with traditional GNNs, we explicitly model the node-to-attribute-category interactions within K hops. Original node-to-node interactions have been enriched by paths passing through $\delta_j \in \mathcal{U}$. Similarly, the first and second-order neighbors of $\delta_j \in \mathcal{U}$ could be computed as,

$$\begin{aligned} \mathcal{N}_{\delta_j} &= \{u | \mathbf{X}_{uj} \neq 0\}, \\ \mathcal{N}_{\delta_j}^{(2)} &= \{u | (\mathbf{A}\mathbf{X})_{uj} > 0\} + \{\delta_i | (\mathbf{X}^\top \mathbf{X})_{ji} \geq 1\}, \end{aligned} \quad (5)$$

where $(\mathbf{X}^\top \mathbf{X}) \in \mathbb{R}^{m \times m}$ denotes the attribute category correlations estimated by their common nodes. It collects $\delta_j \rightarrow u \rightarrow \delta_i$ paths. We enable nodes in \mathcal{V} to propagate messages to attribute categories, and correlations among attribute categories to affect the graph convolution. Based on this, we can not only enrich node interactions using attribute-category as intermediate to improve model performance (see Table 3 and 4), but also enhance model robustness *w.r.t.* missing edges as shown in Section 5.8.

Another thing we want to remark is that different from previous efforts [110], [111] that define a node-to-node similarity network based on feature distance, we directly build a node-to-attribute-category bipartite graph on feature matrix by using attribute values as edge weights. As analyzed before, a bipartite graph between node and attribute-category can not only enrich or rebuild node interactions using attribute as intermedia, but also preserve feature information as much as possible.

4.1.3 Collaborative aggregation

We now illustrate how to learn node representations based on the augmented network \mathbf{P} . Given the recent advances

in heterogeneous graph embedding [112], [112], [113], the intuitive solution is to apply these well-established heterogeneous GNNs to learn embeddings from \mathbf{P} (it consists of two objects (node and attribute-category) and two relations (node-to-node and node-to-attribute-category)). However, our preliminary experiments show that these methods perform not good on our *synthetic heterogeneous graph* (See discussion in 5.8.4), since they may over-emphasize the heterogeneity between node and its attributes, making the learning process substantially complex. Since standard GNNs architectures [46], [50], [62], [92] could already achieve high performance on attributed networks by looking them as homogeneous graph, we propose to follow the tradition and regard the augmented network \mathbf{P} as homogeneous resource with simple weight schema. We leave attributed network embedding from heterogeneous GNNs perspective as our future work.

Specifically, our essential idea is to treat node vertex and attribute-category vertex as identical vertices but do provide a weight hyperparameter α to control the information diffusion between the network \mathbf{A} and node attributes \mathbf{X} . It is because the importance of node-to-node interactions and node attributes are not explicitly available. In our setting, the binary values in \mathbf{A} might not be compatible with the feature values in \mathbf{X} (e.g., continuous features), we define a refined transition probability matrix as follow.

$$\tilde{\mathbf{P}} = \begin{bmatrix} \alpha \tilde{\mathbf{A}} & (1 - \alpha) \tilde{\mathbf{X}} \\ (1 - \alpha) \tilde{\mathbf{X}}^\top & \alpha \mathbf{I} \end{bmatrix}, \quad (6)$$

where $\tilde{\mathbf{A}}$ and $\tilde{\mathbf{X}}$ denote the normalization of $(\mathbf{A} + \mathbf{I})$ and \mathbf{X} after applying ℓ_1 norm to normalize each row respectively. $\alpha \in [0, 1]$ is a trade-off hyper-parameter to impose our inductive bias about the importance of network structure and node attributes. Specifically, when $\alpha = 1$, it yields to a vanilla graph convolutional operation purely based on the network structure. As α increases, node representations would be more dependent on node attributes, and node-to-attribute-category interactions and attribute category correlations will be gradually incorporated into the graph convolution process.

Based on $\tilde{\mathbf{P}}$, standard graph convolutional layers can be directly applied to update node representations. We follow a simple GNN model [54], and update the corresponding embedding matrix with a simple sparse matrix multiplication. We use $\mathbf{H}^{(k)} \in \mathbb{R}^{(n+m) \times d}$ to denote the intermediate representation of all $(n + m)$ nodes in the k^{th} layer. Mathematically, it could be written as,

$$\mathbf{H}^{(k)} \leftarrow \tilde{\mathbf{P}}^K \mathbf{H}^{(0)}. \quad (7)$$

Our initial node representations $\mathbf{H}^{(0)} \in \mathbb{R}^{(n+m) \times d}$ are not based on \mathbf{X} . Instead, $\mathbf{H}^{(0)}$ is a trainable embedding matrix that is randomly initialized following common protocols [58], [114]. To further illustrate the correlation between $\mathbf{h}_v^{(k)}$ and $\{\alpha, \mathbf{A}, \mathbf{X}\}$, we rewrite the corresponding update rule as follows,

$$\begin{aligned} \mathbf{h}_v^{(k)} &= \alpha \tilde{\mathbf{A}}_{vv} \mathbf{h}_v^{(k-1)} + \alpha \sum_{u \in \mathcal{N}_v} \tilde{\mathbf{A}}_{vu} \mathbf{h}_u^{(k-1)} \\ &\quad + (1 - \alpha) \sum_{\delta_j \in \mathcal{N}_v} \tilde{\mathbf{X}}_{vj} \mathbf{h}_{\delta_j}^{(k-1)}. \end{aligned} \quad (8)$$

Eq. (8) provides a principled solution to utilize and control node attributes. On the one hand, it can explicitly enrich or replenish node interactions by treating attribute-category as additional node. On the other hand, neighborhood information from node and attribute-category are selectively combined via trade-off parameter α .

4.2 Collaborative Training Objective

Given the updated representations $\mathbf{H}^{(k)}$ for $k = 0, 1, \dots, K$, we need an unsupervised objective to train the GNN model. A widely-adopted approach [14], [15], [62] is to employ the node representations at the last layer $\mathbf{H}^{(K)}$ to reconstruct all edges. As illustrated in Eq. (2), it estimates the probability of two nodes v and u being connected based on the similarity between their vectors $\mathbf{h}_v^{(K)}$ and $\mathbf{h}_u^{(K)}$. This approach has been demonstrated to be effective in plain networks, but node attributes are often available in practice. Eq. (2) could not directly incorporate node attributes. Another intuitive solution [14] is to employ autoencoders to reconstruct both \mathbf{A} and \mathbf{X} . It would achieve suboptimal performance because the topological structures \mathbf{X} are heterogeneous with node attributes \mathbf{A} .

4.2.1 Cross correlation mechanism

To cope with aforementioned issues, we propose a novel cross correlation mechanism to model and predict the complex node-to-node and node-to-attribute-category interactions. It has two major steps. First, we remove all weights in \mathbf{P} and convert it into a binary matrix. We target at integrating the network and node attribute into our training objective. To make them compatible with each other, we define a binary adjacency matrix as,

$$\mathbf{P}^* = \begin{bmatrix} \mathbf{A} & \mathbf{X}^* \\ \mathbf{X}^{*\top} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{(n+m) \times (n+m)}, \quad (9)$$

where $\mathbf{X}_{v\delta_j}^* = 1$ if $\mathbf{X}_{v\delta_j} > 0$. Our goal is to recover the node-to-node and node-to-attribute-category interactions in \mathbf{P}^* . Second, the node representations $\mathbf{H}^{(k)}$ for $k = 0, 1, \dots, K$ at all layers are learned based on different orders of neighbors. One merit of our model is that we have integrated node attributes \mathbf{X} into the collaborative neighborhood aggregation and no longer need to employ \mathbf{X} as the initial node representations $\mathbf{H}^{(0)}$. Since we are flexible to define $\mathbf{H}^{(0)}$, we could make the dimensions of all $\{\mathbf{H}^{(k)}\}$ the same. In such a way, we could easily take full advantage of them, and model the complex interaction from node v to node u or node v to node attribute category δ_j as,

$$\begin{aligned} \mathbf{y}_{vu} &= \text{MLP}(\|_{k=0}^K \|_{i=0}^K \mathbf{h}_v^{(k)} \odot \mathbf{h}_u^{(i)}), \\ \mathbf{y}_{v\delta_j} &= \text{MLP}(\|_{k=0}^K \|_{i=0}^K \mathbf{h}_v^{(k)} \odot \mathbf{h}_{\delta_j}^{(i)}) \end{aligned} \quad (10)$$

where $\|$ indicates the concatenation operation. $\mathbf{h}_v^{(k)} \odot \mathbf{h}_u^{(i)} \in \mathbb{R}^{n \times d}$ denotes the correlation feature between the under- $(k+1)$ th-order neighborhood of node v and the under- $(i+1)$ th-order neighborhood of node u . \odot represents the element-wise multiplication. $\mathbf{y}_{vu} \in \mathbb{R}$ and $\mathbf{y}_{v\delta_j} \in \mathbb{R}$ are the predicted scores for pairs (v, u) and (v, δ_j) , respectively. MLP denotes a three-layer multilayer perceptron with a Relu activation function.

It is worth noting that the element-wise operation between node representations has been explored in KGAT [104] and NGCF [115] for the recommendation. Our method differs in two ways. First, KGAT and NGCF utilize this technique to facilitate message propagation in each GNN layer (for node-level embedding). Yet, CONN applies it to generate edge representations of two end nodes from different granularities. Second, the proposed cross-correlation mechanism is far more element-wise. Its novelty lies in the proposal to obtain an informative edge representation of end nodes by integrating their cross-correlations from different combinations of the GNN layers. Therefore, the proposed cross-correlation layer is different from KGAT and NGCF because it aims to improve the quality of edge representation while the referenced methods work on boosting the representations of each node per GNN layer, from which they are orthogonal to us and can be incorporated as the base GNN backbone.

4.2.2 Analysis

The subfigure on the bottom of Figure 1 illustrates the key idea of the proposed cross correlation mechanism. Given any two entities (e.g., node $v = 1$ and δ_1 in Figure 1), we aim to capture the second-order correlations across all $K+1$ embedding representations of the two entities. So, in total, Eq. (10) has concatenated $(K+1)^2$ correlation features (e.g., 9 in Figure 1), which is a small number. It should be noted that the element-wise multiplication [116] would highlight the shared patterns within the two input node representations, e.g., $\mathbf{h}_v^{(k)}$ and $\mathbf{h}_u^{(i)}$, while eliminating some inconsistent and noisy information.

4.2.3 Optimization and final representation

In our collaborative training objective, the goal is to reconstruct the node-to-node interactions in \mathbf{A} by using \mathbf{y}_{vu} and the node-to-attribute-category interactions in \mathbf{X}^* by using $\mathbf{y}_{v\delta_j}$. The corresponding objective function is defined as follows.

$$\begin{aligned} \mathcal{L} = & - \sum_{v \in \mathcal{V}} \sum_{z \in \mathcal{N}_v} \log \frac{\exp(\mathbf{y}_{vz})}{\sum_{z' \in \mathcal{V} \cup \mathcal{U}} \exp(\mathbf{y}_{vz'})} \\ & - \sum_{\delta_j \in \mathcal{U}} \sum_{u \in \mathcal{N}_{\delta_j}} \log \frac{\exp(\mathbf{y}_{\delta_j u})}{\sum_{u' \in \mathcal{V} \cup \mathcal{U}} \exp(\mathbf{y}_{\delta_j u'})}. \end{aligned} \quad (11)$$

Eq. (11) is usually intractable in practice, because the sum operation of the denominator is computationally prohibitive. Therefore, we employ the negative sampling [15], [71] strategy to accelerate the optimization.

After we have trained the model CONN by using the collaborative objective in Eq. (11), we need to define the final embedding representation of nodes to perform the downstream tasks. For link prediction task, we directly predict the probability of linking based on \mathbf{y}_{vu} for a node pair (v, u) . For node classification, we adopt the output from the last layer $\mathbf{H}^{(K)}$ as node embedding, in which off-the-shelf classification algorithms could directly use it to classify.

4.3 Comparison with Prior Work

To the best of our knowledge, few efforts have been devoted to leverage node attributes to redefine the core components

of embedding methods for graphs. We roughly divide them into two categories and analyze the difference below.

Random-walk based approach. Random-walk based methods focus on conducting truncated random walks to generate node sequences, and then applying Skip-gram algorithm to learn node representations based on the sequences. This approach is initially not applicable for node attributes. ANRL [60] addresses this issue by modifying the loss function of Skip-gram to depend on node attributes. Featwalk [59] suggests to conduct an attribute-aware random walks to inject node attributes in the random walk generation process. Compared with these methods, our model belongs to graph neural network approach that takes advantage of graph convolutional networks to explicitly model local graph structure.

Graph neural network based approach. This line of methods focuses on exploiting graph convolutional networks [23] to model the local subgraph of an anchor node for node representation. It is natural to cope with attributed graphs by using node attributes as initial node features in the first layer. However, such approach is rather limited in exploiting node attributes as they are excluded from the two crucial components of GCN, i.e., neighborhood message propagation and the training objective. CAN [62] attempts to alleviate this problem by jointly reconstructing node attributes and network structure under variational autoencoder. In contrast, our model focuses on exploiting node attributes to impact the two building blocks jointly.

Besides, we also want to remark that the utilization of node attributes in our work is different from skip connection trick [117]. Skip connection targets to skip some higher-order neighbors in deeper layers by looking back initial features, but our focus is to redefine the neighborhood set of nodes in different orders by regarding attribute categories as "neighbor". Namely, our model introduces additional node-to-attribute-category relations in the message aggregation process of each layer. In a summary, skip connection is complementary with our model and can be added to our GNN backbone to avoid over-smoothing.

Heterogeneous network embedding. Thanks to our proposal in constructing the augmented network between nodes and attribute categories in Section 4.1.1, we can also regard the resultant embedding task as a special heterogeneous network embedding (HNE) problem [112], [113], [118], [119], where we have two object types (node and attribute category) and two relation types (node-to-node and node-to-attribute-category). Under this principle, state-of-the-art heterogeneous models might be applied. However, we found that such approach may make the learning task unnecessarily complex, since HNE will over-emphasize the heterogeneity of attributed networks. This may be suboptimal because attributed networks are not truly "heterogeneous" graphs. Moreover, it's nontrivial to train the model well, because no features are available for the second object-attribute-category. Also, it's hard to control the importance of two relations or node types in a mixed manner, such as mixed random-walk [112]. In contrast, by regarding the augmented network as homogeneous graph in our setting, the problem itself is substantially simplified and arbitrary standard GNNs architectures can be used in the plug-and-play fashion, equipping our proposal with broader applicability and practicability.

TABLE 2
Statistics of the datasets.

	$ \mathcal{V} $	# Edge	$ \mathcal{U} $	#Label
Pubmed	19,717	44,338	500	3
ACM	48,579	119,974	10,000	9
BlogCatalog	5,196	171,743	8,189	6
ogbn-arxiv	169,343	1,166,243	128	40
Reddit	232,965	11,606,919	602	41
ogbl-collab	235,868	1,285,465	128	—

bility and practicability.

To summarize, we propose a new alternative approach for GNNs based attributed network embedding by regarding attribute category as additional nodes and converting attributed network into an augmented graph. The augmented graph offers flexible information diffusion between nodes and attribute categories without information loss. To effectively learn node representations from the new graph, we develop two tailored designs: collaborative aggregation and cross-correlation mechanism, where the former helps to explicitly control the information propagation between nodes and node attributes in convolution, while the latter improves model's reconstruction capacity using multi-granularity features. Although our proposal seems general and is applicable for both homogeneous GNNs and heterogeneous GNNs, we empirically found that it's nontrivial to learn node representations based on heterogeneous GNNs (See discussion in 5.8.4). Therefore, we focus on the simple case and leave attributed network embedding based on heterogeneous GNNs as the future work.

5 EXPERIMENTS

We analyze the effectiveness of CONN on multiple real-world datasets with various scales and types. Specifically, our evaluation centers around three questions.

- **Q1:** Compared with the state-of-the-art embedding methods, can CONN achieve better performance in terms of node classification and link prediction tasks?
- **Q2:** There are three crucial components, i.e., graph mixing convolutional layer, graph correlation layer, and collaborative optimization, in CONN, how much does each component contribute?
- **Q3:** What are the impacts of hyperparameters: the trade-off parameter α and embedding dimension d , on CONN?

5.1 Datasets

We conduct experiments on six publicly available attributed networks of various scales and types. Their statistical information is summarized in Table 2.

Pubmed [120]. It is the biggest benchmark citation network used in [23]. Nodes correspond to documents and edges correspond to citations. Each node has a bag-of-words feature vector according to the paper abstract. Labels are defined as the academic topics.

ACM [121]. It is a large-scale citation network consisting of 48,579 papers published in ACM. Words in the paper abstracts are adopted as node attributes based on the bag-of-words model. Citation links are treated as edges. Each

TABLE 3
Node classification performance.

Method	Pubmed		ACM		BlogCatalog		Reddit		ogbn-arxiv	
	F1-micro	F1-macro	F1-micro	F1-macro	F1-micro	F1-macro	F1-micro	F1-macro	F1-micro	F1-macro
GAE	0.825	0.819	0.710	0.611	0.636	0.632	0.624	0.467	0.614	0.401
ARGE	0.837	0.833	0.736	0.674	0.606	0.601	0.645	0.601	0.633	0.423
DSGC	0.844	0.841	0.759	0.705	0.644	0.639	0.693	0.611	0.643	0.426
DGI	0.857	0.856	0.768	0.709	0.753	0.750	0.605	0.418	0.646	0.415
GIC	0.859	0.857	0.754	0.698	0.773	0.780	0.622	0.433	0.628	0.427
GCA	0.864	0.860	0.757	0.702	0.815	0.836	0.752	0.639	0.655	0.431
FeatWalk	0.843	0.843	0.760	0.703	0.935	0.934	0.663	0.503	0.637	0.424
CAN	0.841	0.835	0.721	0.657	0.652	0.648	0.820	0.726	0.650	0.428
CONN	0.866	0.862	0.775	0.723	0.945	0.944	0.913	0.879	0.674	0.456

paper is published under a specific area which serves as label for classification.

BlogCatalog [1]. It is a social network collected from a blog community. Nodes are web users and edges indicate the user interactions. Node attributes denote the keywords of their blogs. Each user could register his/her blogs into six different predefined classes, which are considered as class labels for node classification.

Reddit [15]. It is another social network dataset collected from the online discussion forum-Reddit. Nodes correspond to the Reddit posts and edges represent the co-comment relationships. The posts are preprocessed into 602-dimensional feature vectors via Glove CommonCrawl world embedding [122]. Hence, node attributes refer to the 602 latent dimensions. We use the communities or ‘subreddit’ that the post belongs to a target label.

ogbl-collab [123]. It is a challenging author collaboration network from KDD Cup 2021. Each node is an author and edges indicate the collaboration between authors. All nodes come with 128-dimensional features, obtained by averaging the word embeddings of papers that are published by the authors. It is widely used to conduct link prediction task.

ogbn-arxiv [123]. It is a large-scale paper citation network of arXiv papers from KDD Cup 2021. Each node is an arXiv paper and the edge indicates that one paper cites another one. Each paper is represented by a 128-dimensional feature vector obtained by averaging the embeddings of words in its title and abstract. The target is to predict 40 subject areas.

5.2 Baseline Methods

To validate the effectiveness of CONN, we include four categories of unsupervised baselines as follows. First, to study why we need tailored framework to incorporate node attributes into GCN-architectures, we compare with vanilla GCN methods GAE [14]. Second, to investigate how effective is CONN compared with other tailored solutions, we include two recent works, i.e., CAN [62] and FeatWalk [59]. Third, to have a comprehensive evaluation with state-of-the-art unsupervised models, we include three popular self-supervised learning based GNN methods, DGI [94], GIC [95], and GCA [96]. Note that, other non-GCN based

embedding methods are not included, i.e., DeepWalk [58], LINE [71] and ANRL [60], since they are outperformed by CAN and FeatWalk in their experiments [59], [62]. Besides, other classical GNNs architectures, i.e., GAT [46], SGC [54], and APPNP [61], are not included for comparison, since they are initially dedicated for supervised learning while we focus on unsupervised representation learning.

- **GAE** [14]. It learns node embeddings by reconstructing the network structure under the autoencoder approach. Specifically, it employs graph convolutional network to encode a subgraph into latent space.
- **ARGE** [124]. It is an adversarially regularized GAE. We do not consider the variational version since the two variants perform quick similar in most cases.
- **DGI** [94]. It learns node embeddings by maximizing the mutual information between the local patch representation and the global graph representation.
- **GIC** [95]. It updates DGI by leveraging cluster-level node representation for unsupervised representation learning.
- **GCA** [96]. It learns node representations by minimizing the contrastive loss between the original graph and its augmented forms.
- **CAN** [62]. It learns node embeddings by reconstructing both the network structure and attribute matrix under the variational autoencoder framework.
- **FeatWalk** [59]. It advances vanilla random-walk based methods via introducing an attribute-enhanced random walk strategy, which helps to generate diversified random walks for representation learning.
- **DSGC** [125]. It defines a k -NN graph based on node features and then uses it as an attribute-aware graph filter for network embedding.

Besides, we also introduce three variants to validate the effectiveness of core components in CONN.

- **CONN-gcn**. It replaces the graph mixing convolutional layer with vanilla graph convolutional layer, to verify the effectiveness of modeling mixed neighbors, i.e., nodes and attributes.
- **CONN-inner**. It excludes the graph correlation layer and utilizes the inner-product to estimate the similarity between two nodes based on their last layer representations similar to [14], [15]. We use it to verify the

usefulness of the proposed correlation layer.

- **CONN-ncoll**. It only considers the node-to-node interactions in the objective function. This variant is used to certify the contribution of jointly optimizing the node-to-node and node-to-attribute interactions.

5.3 Experimental Settings

We follow the common protocol [14], [62] to evaluate the performance of CONN. The effectiveness of the learned latent representations is evaluated over two popular downstream tasks, i.e., link prediction and node classification. For link prediction task, we randomly split 85%, 10% and 5% edges in the network to form training, testing and validation sets similar to [14]. The link prediction task aims to estimate whether a missing edge in the network should be connected or not, based on its embedding representations. The performance is measured by two standard metrics, i.e., area under the ROC curve (AUC) and average precision (AP) scores.

For node classification task, it targets to classify a new instance into one or multiple categories, based on the obtained node representations and the trained classifier. Specifically, we apply 5-fold cross-validation on all datasets to construct the training and test sets. To perform classification, we build an SVM classifier based on scikit-learn package and train the classifier based on the nodes in the training group and corresponding labels. Then we apply the learned classifier to predict the labels of instances in the test groups. The averaged results of five-fold is reported.

We use the official released codes of baselines for experiments and use the validation set to tune their parameters. For our method, we train CONN for 100 epochs using Adam optimizer with learning rate 0.01 and early stopping with a patience of 20 epochs. If it is not specified, d is set as 128, $K = 2$, and α equals to 0.2 and 0.8 for node classification and link prediction tasks, respectively. For continuous datasets (Reddit, ogbn-arxiv, and ogbl-collab), we use top-50 values in each row of \mathbf{X} to construct the bipartite graph by default. The source code of CONN is available at <https://github.com/Qiaoyut/CONN>.

5.4 Node Classification

We start to evaluate the performance of CONN in node classification task (Q1). Table 3 summarizes the results on five datasets in terms of micro-average and macro-average scores.

From Table 3, we observe that CONN performs consistently better than other baselines across two evaluation metrics on datasets with categorical attributes (PubMed, ACM, and BlogCatalog) and continuous feature embeddings (Reddit and ogbn-arxiv). It demonstrates the effectiveness of CONN. To be specific, compared with vanilla GCN variants, CONN improves 48.6% and 9.8% over GAE on BlogCatalog and ogbn-arxiv datasets in terms of micro-average score, respectively. CONN improves 46.7% and 4.8% over ARGE on BlogCatalog and ogbn-arxiv datasets under the macro-average score, respectively. This improvement validates the necessity to design tailored GCN architecture for attributed networks. CAN is also proposed to model on node attributes, but it loses to CONN significantly on five

scenarios. The major difference between them is that CAN targets to jointly reconstruct node attributes and network under auto-encoder framework, while CONN aims to revise GCN architecture by explicitly leveraging node attributes to guide message propagation. This comparison certifies that a tailored GCN solution for attributed networks is more promising and effective. FeatWalk incorporates node attributes to random-walk based models, but it is outperformed by CONN in all cases. Taking Reddit dataset for example, CONN improves FeatWalk over 37.7%. It is reasonable since our model can leverage structure information for node embedding while random-walk based methods fail to use that. Although DSGC also constructs an attribute-aware graph for network embedding, it loses to CONN in almost all cases. This is because DSGC requires to construct a k -NN graph based on node features, which may delete a lot of important attribute information.

Another promising observation is that CONN performs significantly better than state-of-the-art self-supervised competitors (DGI, GIC and GCA) in general. This results indicates the effectiveness of reconstructing original network structure for unsupervised representation learning. Besides, the performance gap between CONN and FeatWalk increases on continuous-value datasets (reddit and orgn-arxiv). This is mainly because FeatWalk reduces to random choice among all attribute categories as the sampling graph is fully-connected, which damages the quality of random walks.

5.5 Link Prediction

We now evaluate the performance of CONN in terms of link prediction task (Q1). Since DGI, GIC, GCA, and FeatWalk are not originally tested on this task, we delete the test edges from the adjacent matrix and use the resulting adjacent matrix to train them. Then, we estimate the similarity scores for test edges based on the inner product of corresponding node representations similar to [62]. Table 4 reports the results on three median size (PubMed, ACM, and BlogCatalog) and two large-scale (Reddit and ogbl-collab) datasets in terms of AUC and AP scores.

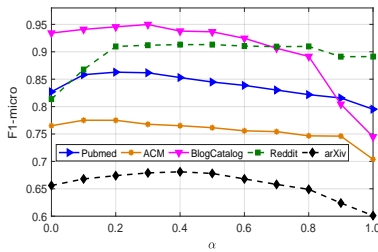
From the table, we can see that CONN performs significantly better than other baselines. Specifically, it improves 11.4%, 21.5%, 16.5%, 25.2%, 22.7%, 11.1%, 48.8%, and 9.7% over GAE, ARGE, DSGC, DGI, GIC, GCA, FeatWalk, and CAN on BlogCatalog in terms of AUC value, respectively. CAN loses to FeatWalk on node classification task in most cases, but outperforms FeatWalk on predicting missing edges. It indicates the importance of capturing structure information for link prediction. Both CAN and CONN target to leverage node attributes for node embedding, but CONN outperforms CAN with a great margin in all cases. The main difference is that CAN focuses on using node attributes to enrich the objective function, with the hope to enhance GCN encoder optimization. In contrast, our model directly merges useful node attributes into the GCN building blocks, such that a more powerful GCN encoder could be explicitly achieved. Although DGI, GIC, and GCA are trained based on the advanced contrastive loss function, our model performs substantially better than them with a wide margin. These results indicate the insufficiency of

TABLE 4
Link prediction results.

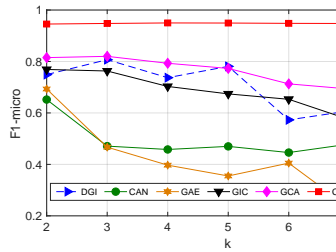
	Pubmed		ACM		BlogCatalog		Reddit		ogbl-collab	
	AUC	AP	AUC	AP	AUC	AP	AUC	AP	AUC	AP
GAE	0.920	0.911	0.957	0.956	0.824	0.822	0.578	0.565	0.821	0.737
ARGE	0.968	0.971	0.964	0.969	0.755	0.723	0.603	0.597	0.847	0.814
DSGC	0.964	0.970	0.961	0.965	0.788	0.774	0.593	0.586	0.842	0.808
DGI	0.942	0.927	0.582	0.644	0.733	0.737	0.594	0.586	0.818	0.728
GIC	0.937	0.935	0.674	0.775	0.748	0.745	0.693	0.677	0.892	0.804
GCA	0.955	0.956	0.756	0.820	0.826	0.808	0.717	0.741	0.886	0.833
FeatWalk	0.940	0.941	0.963	0.963	0.617	0.617	0.852	0.870	0.828	0.797
CAN	0.980	0.977	0.896	0.899	0.837	0.837	0.909	0.903	0.901	0.886
CONN	0.994	0.993	0.986	0.987	0.918	0.906	0.986	0.985	0.930	0.912

TABLE 5
Ablation study of CONN on node classification task.

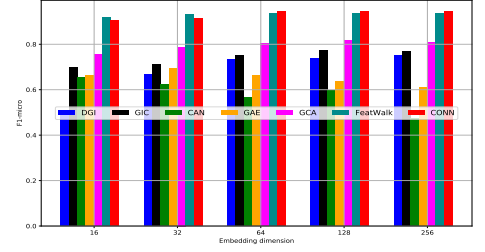
		CONN-gcn	CONN-inner	CONN-ncoll	CONN
F1- micro	Pubmed	0.795	0.701	0.851	0.866
	ACM	0.704	0.643	0.756	0.775
	BlogCatalog	0.745	0.794	0.900	0.945
	Reddit	0.891	0.571	0.891	0.913
	ogbn-arxiv	0.614	0.553	0.635	0.674
F1- macro	Pubmed	0.795	0.708	0.846	0.862
	ACM	0.633	0.551	0.694	0.723
	BlogCatalog	0.740	0.787	0.898	0.944
	Reddit	0.849	0.464	0.844	0.879
	ogbn-arxiv	0.408	0.388	0.415	0.456



(a) Trade-off parameter α



(b) The number of layers K



(c) Embedding dimension d

Fig. 2. Hyper-parameter analysis of CONN.

existing GCN-based architectures in exploiting useful node attributes. Based on these observations, we believe that our proposed framework is more suitable for the link prediction task.

5.6 Ablation Study

We now investigate the second question (Q2), i.e., how much could CONN’s three major components, i.e., graph mixing convolutional layer, graph correlation layer, and collaborative optimization contribute? Three variants, i.e., CONN-gcn, CONN-inner, and CONN-ncoll, that are introduced at the beginning, are used for this ablation study. Table 5 records the node classification performance on five datasets in terms of micro-and macro-average scores.

Based on Table 5, we have three major observations. First, without node attributes to guide message propaga-

tion in the graph convolution process, the performance of CONN-gcn decreases. It validates the effectiveness of explicitly incorporating node attributes into GCN architecture. Second, without the graph correlation layer, CONN-inner loses to CONN with great margin. For instance, CONN achieves 23.5% improvements than CONN-inner on Pubmed in terms of micro-average. The major difference between CONN-inner and CONN is that the former one adopts simple inner product to estimate edge similarity, while CONN devises deep correlation layer, which is capable of capturing the complex correlations between two nodes. This comparison verifies the effectiveness of the proposed graph correlation layer. Third, CONN performs slightly better than CONN-ncoll across five datasets. It indicates that jointly optimizing node-to-node and node-to-

attribute interactions is beneficial. Given that CONN outperforms three variants, it verifies that we propose a principled framework to reinforce the reciprocal effects among the three components.

5.7 Parameter Sensitivity Analysis

We now study the impact of parameters α and embedding dimension d over CONN on BlogCatalog (Q3). α controls the importance of node-to-node interaction and node-to-attribute interaction for message passing. We plot the performance of CONN when α varies from 0 to 1 with step size 0.1 on Figure 2-a. From the results, we observe that the performance of CONN increases as α increases from 0.0 to 0.2 on five datasets. CONN obtains the best results when $\alpha = 0.2$ in general. Notice that, when $\alpha = 0$, CONN only utilizes the node attribute interactions for graph convolution, and the network structure is excluded. When $\alpha = 1$, node-to-attribute interactions are not considered. CONN reduces to the vanilla GCN, except that both node attribute and node interactions are jointly optimized.

K is the number of layers for GCN backbones. Large K means that high-order neighbors are included. Since FeatWalk has no such parameter, we omit it for comparison. We vary K from 2 to 8, and the performance on BlogCatalog is shown in Figure 2-b. From the results, we observe that our model CONN consistently outperforms all baselines on the different number of layers. Another interesting observation is that when K varies from 2 to 8, the performance of CAN, GAE, GIC, GCA, and DGI decreases while CONN performs stable. It validates the superiority of our model in capturing high-order dependencies. It is worth noting that although CONN achieves relatively stable results than standard GNN models in Figure 2 (b) when K is small, e.g., 8, we do observe an obvious performance drop when K is larger, i.e., 15. Therefore, by modeling attribute-aware relationships, our method can alleviate the over-smoothing problem when the GNN model is relatively shallow to some extent. Still, tailored efforts are required to thoroughly tackle the over-smoothing issues as done in [52].

In CONN, the dimension d presents the output latent space for downstream applications, i.e., node classification and link prediction. We search it from $\{16, 32, 64, 128, 256\}$, and the performance of all methods are depicted in Figure 2-c. We observe that different methods have different optimal d , and our model CONN could achieve satisfactory performance when $d = 128$. Similar observations are made on other datasets. To make it fair, we tune the best embedding dimension for all methods and report their best results.

5.8 Further Analysis

With performance comparison with SOTA methods completed, we delve into our proposal for even more insights.

5.8.1 Constructing a sparse bipartite graph from continuous features is a good trade-off.

The first question we would like to answer is: what is the impact of top- N values for graphs with continuous features? Figure 3 shows the results of our model *w.r.t.* different N values on two graphs (Reddit and ogbl-collab) with embedding vectors (e.g., continuous values) as node

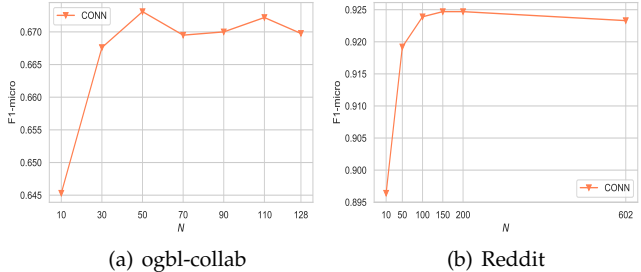


Fig. 3. CONN performance *w.r.t.* different top- N values.

attributes. From the figures, we can see that the performance of CONN increases with the increasing of N values, and it achieves good results when N equals the dimension of node attributes, i.e., 128 and 602 for collab and reddit, respectively.

This observation sheds light on the following insights. (i) Top- N strategy is a good way to handle continuous feature vectors in our model, since $N = 50$ can already achieve satisfactory results in both cases. (ii) Furthermore, it indicates that our model can be used in high-dimensional data, since we can first adopt classical dimensional reduction techniques to reduce the dimension.

To check the scalability of our model on high-dimensional data, we first adopt Deepwalk to obtain 128-d embedding vectors for nodes in ACM and BlogCatalog, and then use the resultant node features as input to train our model. We observe comparable results in two downstream tasks. Specifically, the AUC results for ACM and BlogCatalog on link prediction are 0.987 and 0.916; while the Micro-average results on node classification are 0.769 and 0.941.

5.8.2 Modeling node attributes as a bipartite graph enhances the robustness of GNNs for representation learning.

Next, we explore whether modeling node attributes as an augmented graph can improve model’s robustness towards noises, i.e., missing links. We analyze the robustness of our model and the vanilla GNNs competitor (GAE) towards edge perturbation, i.e., randomly masking some edges with ratios, ranging from 0.0 to 0.9 with step size 0.1. We show the results in terms of link prediction in Figure 4. Similar sensitive tendencies are observed in other datasets.

The results in two figures show that our model performs relatively stable when the masking ratio is less than 50% across two datasets, while the performance of GAE drops in general when the ratio increases. We believe this stable merit is attributed to the proposed augmented bipartite graph because it can replenish missing connections between nodes via using node categories as intermedia.

5.8.3 Jointly reconstructing node attributes and node interactions does not impact the convergence speed.

Moreover, we would like to examine the empirical convergence of our model. We show the training curves of CONN under different GNNs layers over two representative datasets (BlogCatalog and Reddit for categorical and continuous attributes, respectively) in Figure 5. We can

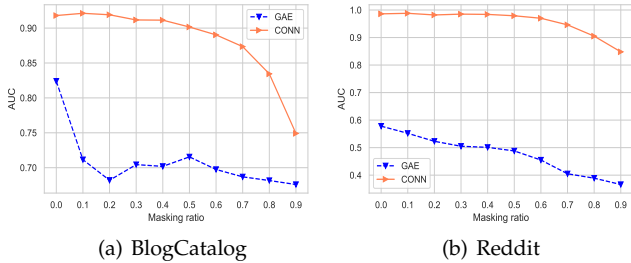


Fig. 4. Robustness analysis of CONN vs. GAE with different edge perturbation ratios.

TABLE 6

Performance of CONN vs. heterogeneous GNNs methods. "NC" stands for node classification results in terms of F1-micro metric; "LP" stands for link prediction results in terms of AUC.

		HetGNN	MAGNN	CONN
NC	Pubmed	0.836	0.839	0.866
	ACM	0.724	0.726	0.775
	BlogCatalog	0.785	0.816	0.945
	Reddit	0.646	0.659	0.913
	ogbn-arxiv	0.625	0.637	0.674
LP	Pubmed	0.936	0.942	0.994
	ACM	0.869	0.877	0.986
	BlogCatalog	0.808	0.812	0.918
	Reddit	0.646	0.681	0.986
	ogbn-collab	0.834	0.856	0.930

observe that the training loss decreases very fast in the first 10 epochs, and our model tends to converge within 50 epochs in general.

5.8.4 Regarding the augmented network as heterogeneous graph is suboptimal.

Finally, we investigate the generalization of our proposal in terms of heterogeneous GNNs. In particular, we want to explore whether the well-established heterogeneous GNNs efforts can be directly applied for attributed network embedding by using our proposed augmented network. To this end, we select two popular unsupervised heterogeneous GNN embedding methods as backbone: HetGNN [112] and MAGNN [126]. For a fair comparison, we initialize learnable embeddings for nodes and attribute categories similar to our model. Table 6 reports the results on both node classification and link prediction tasks.

We observe a clear performance gap between our CONN and two heterogeneous GNNs methods (HetGNN and MAGNN) on all evaluation scenarios. Jointly considering the results in Table 3 and 4, the performance of two heterogeneous variants ranks in the middle against all baselines. The possible explanation is that regarding augmented networks as a heterogeneous graph may over-emphasize the heterogeneity between node interactions and node attributes. It is hard to control the importance between them via either random walk [112] or meta-path [126]. These comparisons demonstrate our motivation to treat the augmented network as a homogeneous graph. We believe nontrivial

efforts are needed to unleash the power of heterogeneous GNNs on attributed network embedding.

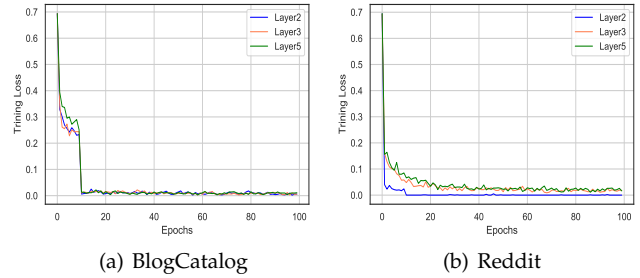


Fig. 5. Empirical training curves of CONN on two datasets with different GNN layers.

6 CONCLUSION

In this paper, we study the problem of unsupervised node representation learning on attributed graphs under graph neural networks (GNNs). Existing GNNs efforts mainly focus on exploiting topological structures, while only using node attributes as initial node presentations in the first layer. Here, we argue that such GNN architecture is suboptimal to modeling real-world attributed graphs, since node attributes are totally excluded from the key factors of GNNs, i.e., the message aggregation mechanism and the training objectives. To tackle this problem, we propose a novel collaborative graph neural network termed CONN. It allows node attributes to determine the message-passing process for neighborhood aggregation and enables the node-to-node and node-to-attribute-category interactions to be jointly recovered. Empirical results on node classification and link prediction tasks over social and citation graphs demonstrate the superiority of CONN against state-of-the-art embedding methods. Our future work is to explore its applicability for dynamic graphs and further improve the robustness of our model by integrating adversarial training. Moreover, we are interested in exploring similar ideas in recommendation scenarios [127], such as sequential recommendation [128], [129] and session recommendation [130], [131].

7 ACKNOWLEDGEMENT

We thank the anomalous reviewers for the feedback. The work is, in part, supported by NSF (IIS-2224843 and IIS-1849085).

REFERENCES

- [1] X. Huang, J. Li, and X. Hu, "Label informed attributed network embedding," in *Proceedings of the tenth ACM international conference on web search and data mining*, 2017, pp. 731–739.
- [2] H. Gao and H. Huang, "Deep attributed network embedding," in *International Joint Conference on Artificial Intelligence*, vol. 18. New York, NY, 2018, pp. 3364–3370.
- [3] R. Hong, Y. He, L. Wu, Y. Ge, and X. Wu, "Deep attributed network embedding by preserving structure and attribute information," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 51, no. 3, pp. 1434–1445, 2019.

- [4] H. Wang, D. Lian, H. Tong, Q. Liu, Z. Huang, and E. Chen, "Decoupled representation learning for attributed networks," *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [5] M. McPherson, L. Smith-Lovin, and J. M. Cook, "Birds of a feather: Homophily in social networks," *Annual review of sociology*, pp. 415–444, 2001.
- [6] P. V. Marsden and N. E. Friedkin, "Network studies of social influence," *Sociological Methods & Research*, vol. 22, no. 1, pp. 127–151, 1993.
- [7] Z. Yang, J. Guo, K. Cai, J. Tang, J. Li, L. Zhang, and Z. Su, "Understanding retweeting behaviors in social networks," in *Proceedings of the 19th ACM international conference on Information and knowledge management*, 2010, pp. 1633–1636.
- [8] L. Liao, X. He, H. Zhang, and T.-S. Chua, "Attributed social network embedding," *IEEE Transactions on Knowledge and Data Engineering*, vol. 30, no. 12, pp. 2257–2270, 2018.
- [9] S. Wang, C. Aggarwal, J. Tang, and H. Liu, "Attributed signed network embedding," in *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, 2017, pp. 137–146.
- [10] H. Wang, J. Wang, J. Wang, M. Zhao, W. Zhang, F. Zhang, W. Li, X. Xie, and M. Guo, "Learning graph representation with generative adversarial nets," *IEEE Transactions on Knowledge and Data Engineering*, vol. 33, no. 8, pp. 3090–3103, 2019.
- [11] H. Yang, S. Pan, P. Zhang, L. Chen, D. Lian, and C. Zhang, "Binarized attributed network embedding," in *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2018, pp. 1476–1481.
- [12] Y. Gao, M. Gong, Y. Xie, and H. Zhong, "Community-oriented attributed network embedding," *Knowledge-Based Systems*, vol. 193, p. 105418, 2020.
- [13] B. Rozemberczki, C. Allen, and R. Sarkar, "Multi-scale attributed node embedding," *Journal of Complex Networks*, vol. 9, no. 2, p. cnab014, 2021.
- [14] T. N. Kipf and M. Welling, "Variational graph auto-encoders," *arXiv preprint arXiv:1611.07308*, 2016.
- [15] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Advances in Neural Information Processing Systems*, 2017, pp. 1024–1034.
- [16] J. Li, H. Dani, X. Hu, J. Tang, Y. Chang, and H. Liu, "Attributed network embedding for learning in a dynamic environment," in *International Conference on Information and Knowledge Management*, 2017, pp. 387–396.
- [17] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, "Neural message passing for quantum chemistry," in *International Conference on Machine Learning*, 2017, pp. 1263–1272.
- [18] X. Han, Z. Jiang, N. Liu, Q. Song, J. Li, and X. Hu, "Geometric graph representation learning via maximizing rate reduction," in *Proceedings of the ACM Web Conference 2022*, 2022, pp. 1226–1237.
- [19] H. Gao, Z. Wang, and S. Ji, "Large-scale learnable graph convolutional networks," in *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2018, pp. 1416–1424.
- [20] Z. Jiang, X. Han, C. Fan, Z. Liu, N. Zou, A. Mostafavi, and X. Hu, "Fmp: Toward fair graph message passing against topology bias," *arXiv preprint arXiv:2202.04187*, 2022.
- [21] X. Li and Y. Cheng, "Understanding the message passing in graph neural networks via power iteration clustering," *Neural Networks*, vol. 140, pp. 130–135, 2021.
- [22] Z. Jiang, X. Han, C. Fan, Z. Liu, X. Huang, N. Zou, A. Mostafavi, and X. Hu, "Topology matters in fair graph learning: a theoretical pilot study," 2022.
- [23] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *International Conference on Learning Representations*, 2017.
- [24] Q. Tan, N. Liu, X. Huang, S.-H. Choi, L. Li, R. Chen, and X. Hu, "S2gae: Self-supervised graph autoencoders are generalizable learners with graph masking," in *ACM International Conference on Web Search and Data Mining*, 2023, pp. 787–795.
- [25] Z. Jiang, X. Han, C. Fan, Z. Liu, N. Zou, A. Mostafavi, and X. Hu, "Fair graph message passing with transparency," 2022.
- [26] Q. Li, Z. Han, and X.-M. Wu, "Deeper insights into graph convolutional networks for semi-supervised learning," in *AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [27] 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.
- [28] Q. Tan, D. Zha, S.-H. Choi, L. Li, R. Chen, and X. Hu, "Double wins: Boosting accuracy and efficiency of graph neural networks by reliable knowledge distillation," 2022.
- [29] X. Miao, W. Zhang, Y. Shao, B. Cui, L. Chen, C. Zhang, and J. Jiang, "Lasagne: A multi-layer graph convolutional network framework via node-aware deep architecture," *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [30] Q. Tan, N. Liu, and X. Hu, "Deep representation learning for social network analysis," *Frontiers in Big Data*, vol. 2, p. 2, 2019.
- [31] X. Huang, J. Li, and X. Hu, "Accelerated attributed network embedding," in *Proceedings of the 2017 SIAM international conference on data mining*. SIAM, 2017, pp. 633–641.
- [32] Q. Tan, S. Ding, N. Liu, S.-H. Choi, L. Li, R. Chen, and X. Hu, "Graph contrastive learning with model perturbation," 2022.
- [33] X. Gao, W. Hu, and G.-J. Qi, "Self-supervised graph representation learning via topology transformations," *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [34] X. Han, Z. Jiang, N. Liu, and X. Hu, "G-mixup: Graph data augmentation for graph classification," in *International Conference on Machine Learning*. PMLR, 2022, pp. 8230–8248.
- [35] X. Zhang, Q. Tan, X. Huang, and B. Li, "Graph contrastive learning with personalized augmentation," *arXiv preprint arXiv:2209.06560*, 2022.
- [36] Q. Tan, X. Zhang, N. Liu, D. Zha, L. Li, R. Chen, S.-H. Choi, and X. Hu, "Bring your own view: Graph neural networks for link prediction with personalized subgraph selection," in *ACM International Conference on Web Search and Data Mining*, 2023, pp. 625–633.
- [37] M. Zhang and Y. Chen, "Link prediction based on graph neural networks," in *Advances in Neural Information Processing Systems*, 2018, pp. 5165–5175.
- [38] M. Zhang, P. Li, Y. Xia, K. Wang, and L. Jin, "Labeling trick: A theory of using graph neural networks for multi-node representation learning," *Advances in Neural Information Processing Systems*, vol. 34, pp. 9061–9073, 2021.
- [39] Q. Tan, N. Liu, X. Zhao, H. Yang, J. Zhou, and X. Hu, "Learning to hash with graph neural networks for recommender systems," in *International Conference on World Wide Web*, 2020, pp. 1988–1998.
- [40] Y. He, "Gnns for node clustering in signed and directed networks," in *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*, 2022, pp. 1547–1548.
- [41] C. Fattal, L. Labiod, and M. Nadif, "Efficient graph convolution for joint node representation learning and clustering," in *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*, 2022, pp. 289–297.
- [42] S. Zhou, Q. Tan, Z. Xu, X. Huang, and F.-I. Chung, "Subtractive aggregation for attributed network anomaly detection," in *ACM International Conference on Information & Knowledge Management*, 2021, pp. 3672–3676.
- [43] S. Zhou, X. Huang, N. Liu, Q. Tan, and F.-L. Chung, "Unseen anomaly detection on networks via multi-hypersphere learning," in *SIAM International Conference on Data Mining*. SIAM, 2022, pp. 262–270.
- [44] Q. Zhang, J. Dong, Q. Tan, and X. Huang, "Integrating entity attributes for error-aware knowledge graph embedding," *IEEE Transactions on Knowledge and Data Engineering*, 2023.
- [45] J. Dong, Q. Zhang, X. Huang, Q. Tan, D. Zha, and Z. Zihao, "Active ensemble learning for knowledge graph error detection," in *ACM International Conference on Web Search and Data Mining*, 2023, pp. 877–885.
- [46] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.
- [47] J. Gao, J. Gao, X. Ying, M. Lu, and J. Wang, "Higher-order interaction goes neural: a substructure assembling graph attention network for graph classification," *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [48] Z. Luo, Y. Cui, S. Zhao, and J. Yin, "g-inspector: Recurrent attention model on graph," *IEEE Transactions on Knowledge and Data Engineering*, 2020.
- [49] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How powerful are graph neural networks?" *arXiv preprint arXiv:1810.00826*, 2018.
- [50] H. Gao and S. Ji, "Graph representation learning via hard and channel-wise attention networks," in *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2019, pp. 741–749.
- [51] W. Feng, J. Zhang, Y. Dong, Y. Han, H. Luan, Q. Xu, Q. Yang, E. Kharlamov, and J. Tang, "Graph random neural networks for semi-supervised learning on graphs," *Advances in Neural Information Processing Systems*, vol. 33, 2020.

- [52] M. Chen, Z. Wei, Z. Huang, B. Ding, and Y. Li, "Simple and deep graph convolutional networks," in *International conference on machine learning*. PMLR, 2020, pp. 1725–1735.
- [53] K. Stachenfeld, J. Godwin, and P. Battaglia, "Graph networks with spectral message passing," *arXiv preprint arXiv:2101.00079*, 2020.
- [54] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, and K. Weinberger, "Simplifying graph convolutional networks," in *International conference on machine learning*, 2019, pp. 6861–6871.
- [55] V. P. Dwivedi, C. K. Joshi, T. Laurent, Y. Bengio, and X. Bresson, "Benchmarking graph neural networks," *arXiv preprint arXiv:2003.00982*, 2020.
- [56] Z. Liu, T.-K. Nguyen, and Y. Fang, "Tail-gnn: Tail-node graph neural networks," in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021, pp. 1109–1119.
- [57] W. Zheng, E. W. Huang, N. Rao, S. Katariya, Z. Wang, and K. Subbian, "Cold brew: Distilling graph node representations with incomplete or missing neighborhoods," *arXiv preprint arXiv:2111.04840*, 2021.
- [58] B. Perozzi, R. Al-Rfou, and S. Skiena, "Deepwalk: Online learning of social representations," in *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2014, pp. 701–710.
- [59] X. Huang, Q. Song, F. Yang, and X. Hu, "Large-scale heterogeneous feature embedding," in *AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 3878–3885.
- [60] Z. Zhang, H. Yang, J. Bu, S. Zhou, P. Yu, J. Zhang, M. Ester, and C. Wang, "Anrl: Attributed network representation learning via deep neural networks." in *International Joint Conference on Artificial Intelligence*, vol. 18, 2018, pp. 3155–3161.
- [61] J. Klicpera, A. Bojchevski, and S. Günnemann, "Predict then propagate: Graph neural networks meet personalized pagerank," *International Conference on Learning Representations*, 2019.
- [62] Z. Meng, S. Liang, H. Bao, and X. Zhang, "Co-embedding attributed networks," in *ACM International Conference on Web Searching and Data Mining*, 2019, pp. 393–401.
- [63] Z. Zhang, P. Cui, and W. Zhu, "Deep learning on graphs: A survey," *IEEE Transactions on Knowledge and Data Engineering*, 2020.
- [64] J. Zhou, G. Cui, Z. Zhang, C. Yang, Z. Liu, L. Wang, C. Li, and M. Sun, "Graph neural networks: A review of methods and applications," *arXiv preprint arXiv:1812.08434*, 2018.
- [65] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and S. Y. Philip, "A comprehensive survey on graph neural networks," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [66] I. Makarov, D. Kiselev, N. Nikitinsky, and L. Subelj, "Survey on graph embeddings and their applications to machine learning problems on graphs," *PeerJ Computer Science*, vol. 7, p. e357, 2021.
- [67] D. Zhang, J. Yin, X. Zhu, and C. Zhang, "Network representation learning: A survey," *IEEE transactions on Big Data*, vol. 6, no. 1, pp. 3–28, 2018.
- [68] X. He and P. Niyogi, "Locality preserving projections," *Advances in neural information processing systems*, vol. 16, 2003.
- [69] M. Belkin and P. Niyogi, "Laplacian eigenmaps and spectral techniques for embedding and clustering," *Advances in Neural Information Processing Systems*, vol. 14, pp. 585–591, 2001.
- [70] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Advances in Neural Information Processing Systems*, 2013, pp. 3111–3119.
- [71] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, and Q. Mei, "Line: Large-scale information network embedding," in *International World Wide Web Conference*, 2015, pp. 1067–1077.
- [72] S. Cao, W. Lu, and Q. Xu, "Grarep: Learning graph representations with global structural information," in *International Conference on Information and Knowledge Management*, 2015, pp. 891–900.
- [73] A. Grover and J. Leskovec, "node2vec: Scalable feature learning for networks," in *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2016, pp. 855–864.
- [74] D. Wang, P. Cui, and W. Zhu, "Structural deep network embedding," in *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2016, pp. 1225–1234.
- [75] R. A. Rossi, R. Zhou, and N. K. Ahmed, "Deep inductive graph representation learning," *IEEE transactions on knowledge and data engineering*, vol. 32, no. 3, pp. 438–452, 2018.
- [76] C. Tu, X. Zeng, H. Wang, Z. Zhang, Z. Liu, M. Sun, B. Zhang, and L. Lin, "A unified framework for community detection and network representation learning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 31, no. 6, pp. 1051–1065, 2018.
- [77] N. Liu, Q. Tan, Y. Li, H. Yang, J. Zhou, and X. Hu, "Is a single vector enough? exploring node polysemy for network embedding," in *ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 932–940.
- [78] X. Wang, P. Cui, J. Wang, J. Pei, W. Zhu, and S. Yang, "Community preserving network embedding," in *AAAI Conference on Artificial Intelligence*, vol. 31, no. 1, 2017.
- [79] R. Yang, J. Shi, X. Xiao, Y. Yang, and S. S. Bhowmick, "Homogeneous network embedding for massive graphs via reweighted personalized pagerank," *arXiv preprint arXiv:1906.06826*, 2019.
- [80] J. Qiu, Y. Dong, H. Ma, J. Li, C. Wang, K. Wang, and J. Tang, "Netsmf: Large-scale network embedding as sparse matrix factorization," in *The World Wide Web Conference*, 2019, pp. 1509–1520.
- [81] W. Lin, "Large-scale network embedding in apache spark," in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021, pp. 3271–3279.
- [82] S. Zhao, Z. Du, J. Chen, Y. Zhang, J. Tang, and P. Yu, "Hierarchical representation learning for attributed networks," *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [83] C. Yang, Z. Liu, D. Zhao, M. Sun, and E. Y. Chang, "Network representation learning with rich text information." in *International Joint Conference on Artificial Intelligence*, vol. 2015, 2015, pp. 2111–2117.
- [84] W. Wang, X. Wei, X. Suo, B. Wang, H. Wang, H.-N. Dai, and X. Zhang, "Hgate: heterogeneous graph attention auto-encoders," *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [85] X. Zhao, Q. Dai, J. Wu, H. Peng, M. Liu, X. Bai, J. Tan, S. Wang, and P. Yu, "Multi-view tensor graph neural networks through reinforced aggregation," *IEEE Transactions on Knowledge and Data Engineering*, 2022.
- [86] X. Zhang, H. Liu, Q. Li, and X.-M. Wu, "Attributed graph clustering via adaptive graph convolution," *arXiv preprint arXiv:1906.01210*, 2019.
- [87] Z. Liu, C. Huang, Y. Yu, and J. Dong, "Motif-preserving dynamic attributed network embedding," in *Proceedings of the Web Conference 2021*, 2021, pp. 1629–1638.
- [88] W. Huang, Y. Li, Y. Fang, J. Fan, and H. Yang, "Biane: Bipartite attributed network embedding," in *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*, 2020, pp. 149–158.
- [89] R. Yang, J. Shi, X. Xiao, Y. Yang, J. Liu, and S. S. Bhowmick, "Scaling attributed network embedding to massive graphs," *arXiv preprint arXiv:2009.00826*, 2020.
- [90] J. Tang, M. Qu, and Q. Mei, "Pte: Predictive text embedding through large-scale heterogeneous text networks," in *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2015, pp. 1165–1174.
- [91] H. Gao, J. Pei, and H. Huang, "Progan: Network embedding via proximity generative adversarial network," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 1308–1316.
- [92] M. Liu, H. Gao, and S. Ji, "Towards deeper graph neural networks," in *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, 2020, pp. 338–348.
- [93] J. Fang, S. Liang, Z. Meng, and M. De Rijke, "Hyperspherical variational co-embedding for attributed networks," *ACM Transactions on Information Systems (TOIS)*, vol. 40, no. 3, pp. 1–36, 2021.
- [94] P. Velickovic, W. Fedus, W. L. Hamilton, P. Liò, Y. Bengio, and R. D. Hjelm, "Deep graph infomax." in *International Conference on Learning Representations*, 2019.
- [95] C. Mavromatis and G. Karypis, "Graph infoclust: Maximizing coarse-grain mutual information in graphs," in *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 2021, pp. 541–553.
- [96] Y. Zhu, Y. Xu, F. Yu, Q. Liu, S. Wu, and L. Wang, "Graph contrastive learning with adaptive augmentation," in *Proceedings of the Web Conference 2021*, 2021, pp. 2069–2080.
- [97] K. Hassani and A. H. Khasahmadi, "Contrastive multi-view representation learning on graphs," in *International Conference on Machine Learning*. PMLR, 2020, pp. 4116–4126.
- [98] N. Park, A. Kan, X. L. Dong, T. Zhao, and C. Faloutsos, "Estimating node importance in knowledge graphs using graph

- neural networks,” in *Proceedings of the ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 596–606.
- [99] Z. Zhang, F. Zhuang, H. Zhu, Z. Shi, H. Xiong, and Q. He, “Relational graph neural network with hierarchical attention for knowledge graph completion,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 05, 2020, pp. 9612–9619.
- [100] M. Yasunaga, H. Ren, A. Bosselut, P. Liang, and J. Leskovec, “Qaggn: Reasoning with language models and knowledge graphs for question answering,” in *Proceedings of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2021, pp. 535–546.
- [101] J. Dong, Q. Zhang, X. Huang, K. Duan, Q. Tan, and Z. Jiang, “Hierarchy-aware multi-hop question answering over knowledge graphs,” in *International World Wide Web Conference*, 2023, pp. 2519–2527.
- [102] Y. Hou, H. Chen, C. Li, J. Cheng, and M.-C. Yang, “A representation learning framework for property graphs,” in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 65–73.
- [103] H. Wang, M. Zhao, X. Xie, W. Li, and M. Guo, “Knowledge graph convolutional networks for recommender systems,” in *The world wide web conference*, 2019, pp. 3307–3313.
- [104] X. Wang, X. He, Y. Cao, M. Liu, and T.-S. Chua, “Kgat: Knowledge graph attention network for recommendation,” in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 950–958.
- [105] L. Yao, C. Mao, and Y. Luo, “Graph convolutional networks for text classification,” in *AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 7370–7377.
- [106] C. Park, D. Kim, J. Han, and H. Yu, “Unsupervised attributed multiplex network embedding,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 5371–5378.
- [107] G. Cui, J. Zhou, C. Yang, and Z. Liu, “Adaptive graph encoder for attributed graph embedding,” in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 976–985.
- [108] J. Chen, S. Zheng, Y. Song, J. Rao, and Y. Yang, “Learning attributed graph representations with communicative message passing transformer,” *arXiv preprint arXiv:2107.08773*, 2021.
- [109] Z. Zhou, S. Zhou, B. Mao, X. Zhou, J. Chen, Q. Tan, D. Zha, C. Wang, Y. Feng, and C. Chen, “Opengsl: benchmarking graph structure learning,” in *arxiv*, 2023.
- [110] X. Wang, M. Zhu, D. Bo, P. Cui, C. Shi, and J. Pei, “Amgcn: Adaptive multi-channel graph convolutional networks,” in *Proceedings of the 26th ACM SIGKDD International conference on knowledge discovery & data mining*, 2020, pp. 1243–1253.
- [111] C. Liu, L. Wen, Z. Kang, G. Luo, and L. Tian, “Self-supervised consensus representation learning for attributed graph,” in *Proceedings of the 29th ACM International Conference on Multimedia*, 2021, pp. 2654–2662.
- [112] C. Zhang, D. Song, C. Huang, A. Swami, and N. V. Chawla, “Heterogeneous graph neural network,” in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 793–803.
- [113] D. Jin, C. Huo, C. Liang, and L. Yang, “Heterogeneous graph neural network via attribute completion,” in *Proceedings of the Web Conference 2021*, 2021, pp. 391–400.
- [114] J. Qiu, Y. Dong, H. Ma, J. Li, K. Wang, and J. Tang, “Network embedding as matrix factorization: Unifying deepwalk, line, pte, and node2vec,” in *ACM International Conference on Web Searching and Data Mining*, 2018, pp. 459–467.
- [115] X. Wang, X. He, M. Wang, F. Feng, and T.-S. Chua, “Neural graph collaborative filtering,” in *Proceedings of the 42nd international ACM SIGIR conference on Research and development in Information Retrieval*, 2019, pp. 165–174.
- [116] H. Zhu, F. Feng, X. He, X. Wang, Y. Li, K. Zheng, and Y. Zhang, “Bilinear graph neural network with neighbor interactions,” in *International Joint Conference on Artificial Intelligence*, vol. 5, 2020.
- [117] Y. Li and I. King, “Autograph: Automated graph neural network,” in *International Conference on Neural Information Processing*. Springer, 2020, pp. 189–201.
- [118] C. Yang, Y. Xiao, Y. Zhang, Y. Sun, and J. Han, “Heterogeneous network representation learning: A unified framework with survey and benchmark,” *IEEE Transactions on Knowledge and Data Engineering*, 2020.
- [119] X. Wang, H. Ji, C. Shi, B. Wang, Y. Ye, P. Cui, and P. S. Yu, “Heterogeneous graph attention network,” in *The world wide web conference*, 2019, pp. 2022–2032.
- [120] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassi-Rad, “Collective classification in network data,” *AI magazine*, vol. 29, no. 3, pp. 93–93, 2008.
- [121] J. Tang, J. Zhang, L. Yao, J. Li, L. Zhang, and Z. Su, “Arnetminer: extraction and mining of academic social networks,” in *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2008, pp. 990–998.
- [122] J. Pennington, R. Socher, and C. D. Manning, “Glove: Global vectors for word representation,” in *EMNLP*, 2014, pp. 1532–1543.
- [123] K. Wang, Z. Shen, C. Huang, C.-H. Wu, Y. Dong, and A. Kanakia, “Microsoft academic graph: When experts are not enough,” *Quantitative Science Studies*, vol. 1, no. 1, pp. 396–413, 2020.
- [124] S. Pan, R. Hu, G. Long, J. Jiang, L. Yao, and C. Zhang, “Adversarially regularized graph autoencoder for graph embedding,” *arXiv preprint arXiv:1802.04407*, 2018.
- [125] Q. Li, X. Zhang, H. Liu, Q. Dai, and X.-M. Wu, “Dimensionwise separable 2-d graph convolution for unsupervised and semi-supervised learning on graphs,” in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021, pp. 953–963.
- [126] X. Fu, J. Zhang, Z. Meng, and I. King, “Magnn: Metapath aggregated graph neural network for heterogeneous graph embedding,” in *Proceedings of The Web Conference 2020*, 2020, pp. 2331–2341.
- [127] D. Zha, L. Feng, Q. Tan, Z. Liu, K.-H. Lai, B. Bhushanam, Y. Tian, A. Kejariwal, and X. Hu, “Dreamshard: Generalizable embedding table placement for recommender systems,” in *Advances in Neural Information Processing Systems*, 2022.
- [128] Q. Tan, J. Zhang, J. Yao, N. Liu, J. Zhou, H. Yang, and X. Hu, “Sparse-interest network for sequential recommendation,” in *ACM International Conference on Web Search and Data Mining*, 2021, pp. 598–606.
- [129] Q. Tan, J. Zhang, N. Liu, X. Huang, H. Yang, J. Zhou, and X. Hu, “Dynamic memory based attention network for sequential recommendation,” in *AAAI Conference on Artificial Intelligence*, vol. 35, no. 5, 2021, pp. 4384–4392.
- [130] H. Zhou, Q. Tan, X. Huang, K. Zhou, and X. Wang, “Temporal augmented graph neural networks for session-based recommendations,” in *International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2021, pp. 1798–1802.
- [131] H. Zhou, S. Zhou, K. Duan, X. Huang, Q. Tan, and Z. Yu, “Interest driven graph structure learning for session-based recommendation,” in *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 2023, pp. 284–296.