

Correlation-Aware Graph Convolutional Networks for Multi-Label Node Classification

Yuanchen Bei
Zhejiang University
Hangzhou, China
yuanchenbei@zju.edu.cn

Weizhi Chen
Zhejiang University
Hangzhou, China
chenweizhi@zju.edu.cn

Hao Chen
The Hong Kong Polytechnic
University
HongKong SAR, China
sundaychenhao@gmail.com

Sheng Zhou*
Zhejiang University
Hangzhou, China
zhousheng_zju@zju.edu.cn

Carl Yang
Emory University
Atlanta, USA
j.carlyang@emory.edu

Jiapei Fan
Alibaba Group
Hangzhou, China
jiapei.fjp@@alibaba-inc.com

Longtao Huang
Alibaba Group
Hangzhou, China
kaiyang.hlt@@alibaba-inc.com

Jiajun Bu
Zhejiang University
Hangzhou, China
bjj@zju.edu.cn

ABSTRACT

Multi-label node classification is an important yet under-explored domain in graph mining as many real-world nodes belong to multiple categories rather than just a single one. Although a few efforts have been made by utilizing Graph Convolution Networks (GCNs) to learn node representations and model correlations between multiple labels in the embedding space, they still suffer from the ambiguous feature and ambiguous topology induced by multiple labels, which reduces the credibility of the messages delivered in graphs and overlooks the label correlations on graph data. Therefore, it is crucial to reduce the ambiguity and empower the GCNs for accurate classification. However, this is quite challenging due to the requirement of retaining the distinctiveness of each label while fully harnessing the correlation between labels simultaneously. To address these issues, in this paper, we propose a **Correlation-aware Graph Convolutional Network (CorGCN)** for multi-label node classification. By introducing a novel Correlation-Aware Graph Decomposition module, CorGCN can learn a graph that contains rich label-correlated information for each label. It then employs a Correlation-Enhanced Graph Convolution to model the relationships between labels during message passing to further bolster the classification process. Extensive experiments on five datasets demonstrate the effectiveness of our proposed CorGCN.

*Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-XXXX-X/18/06
<https://doi.org/XXXXXXXX.XXXXXXX>

CCS CONCEPTS

• **Information systems** → **Web mining**; *Social networks*; • **Human-centered computing** → *Social network analysis*.

KEYWORDS

multi-label node classification, graph learning, graph neural networks, graph data mining

ACM Reference Format:

Yuanchen Bei, Weizhi Chen, Hao Chen, Sheng Zhou, Carl Yang, Jiapei Fan, Longtao Huang, and Jiajun Bu. 2018. Correlation-Aware Graph Convolutional Networks for Multi-Label Node Classification. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 14 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

1 INTRODUCTION

Node classification serves as a cornerstone in the field of graph mining [7, 45]. Over the past decade, Graph Convolutional Networks (GCNs) have achieved remarkable success in this area by aggregating information from neighboring nodes, where edges in the graph often indicate similarities in the single-label space [20, 34, 42]. However, graphs in real-world scenarios often entail multi-label nodes. [For example, users in social networks usually exhibit broad interests and embody multiple labels [32, 51], and protein nodes in protein interaction networks usually carry several relevant gene ontology annotations [12, 50]. The diverse labels offer varied perspectives for delineating node characteristics while also introducing both challenges and opportunities for GCNs in the realm of multi-label node classification.

Inspired by the success of multi-label learning on non-relational data [24, 38, 40], such as images and text, there has been a growing interest in applying multi-label learning to relational data, including graphs [11, 44, 55]. Among the few existing methods for multi-label node classification, most employ GCNs to map nodes to low-dimensional representations. They then follow the traditional multi-label learning paradigm to model the relationships

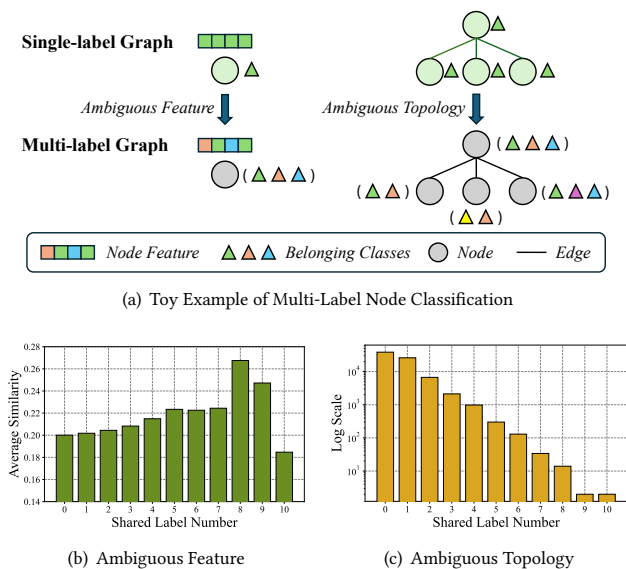


Figure 1: (a) A toy example of multi-label node classification challenges. (b)-(c) Illustrations of the ambiguous feature and ambiguous topology issues on the PCG dataset.

between nodes and labels, as well as the inter-label relationships in the representation space. However, this paradigm overlooks the characteristics of graph data in multi-label scenarios, as Figure 1-(a), leading to challenges for GCNs. (1) **Ambiguous Feature**: In the single-label setting, nodes can aggregate patterns specific to a particular type of label from one neighborhood node by transforming the features. However, in the multi-label setting, the feature of one node may be related to multiple labels and the patterns represented by the feature are ambiguous. As a result, aggregating information from such ambiguous feature will significantly impair the discriminative ability of the representation. Figure 1-(b) illustrates that nodes that allocate similar features may share a diverse number of labels in real-world graphs. (2) **Ambiguous Topology**: In the single-label setting, connected nodes typically share the same single label (also known as the homophily assumption), so the pattern propagated along the edge is usually deterministic. However, in the multi-label setting, both the connected nodes have multiple labels, and the patterns propagated along the edge are often ambiguous. This ambiguity makes it difficult to determine from which connected nodes we should aggregate specific label information. Aggregating information directly from all neighbors would further accumulate ambiguity, compromising the discriminative power of the representations learned by GCNs and eventually impacting the inferring of specific labels. Figure 1-(c) illustrates that connected nodes may share diverse numbers of labels in real-world graphs. *It is essential to reduce the ambiguity so that the potential of GCNs in multi-label node classification can be fully released.*

Although important, addressing the above issues is non-trivial and meets the following challenges: (1) **Label Distinctiveness**: As previously discussed, node attributes and edges in a multi-label graph may be influenced by multiple labels simultaneously. Directly extracting information from such a graph with a mixture of node

labels can result in the loss of label distinctiveness and lead to inadequate exploration of each label. (2) **Label Correlation**: In the multi-label setting, a node's association with multiple labels implies a correlation between these labels. The success of existing multi-label learning methods also demonstrates that fully leveraging these inter-label correlations can significantly enhance the quality of the representation [33, 54]. Consequently, merely extracting messages or neighbors for a single label only could forfeit these vital correlations, leading to suboptimal outcomes. This raises a pertinent question: *How can we retain the distinctiveness of each label while fully harnessing the correlation between labels to achieve more accurate multi-label learning?*

To address the aforementioned challenges, in this paper, we present (**CorGCN**), a **Correlation-Aware Graph Convolutional Network** for multi-label node classification. To tackle the first challenge, we propose a novel graph decomposition strategy where we learn an individual graph for each label while preserving its unique characteristics. To tackle the second challenge, we fully consider the correlation between labels during the graph decomposition process. Instead of indiscriminately discarding all information from other categories, we retain the information of related labels within each category's graph. Lastly, based on the multiple graphs generated, we further exploit the correlation between labels using a newly designed correlation-enhanced GCN that takes category correlation into account, thereby enhancing the final multi-label classification results. We conduct extensive experiments on five datasets with comprehensive metrics. The results demonstrate the effectiveness of the proposed method. The primary contributions of this paper can be summarized as follows:

- We highlight the overlooked impact of ambiguous feature and ambiguous topology on GCNs in the multi-label node classification, which is fundamental in this critical task.
- We propose a Correlation-aware Graph Convolutional Network (CorGCN) to simultaneously retain label-distinct characters and label correlation information for multi-label node classification.
- Extensive experiments on five datasets demonstrate CorGCN significantly outperforms nine state-of-the-art baselines across seven metrics. Further in-depth analysis from diverse perspectives also demonstrates the strengths of CorGCN.

2 RELATED WORKS

2.1 Multi-Label Learning

Multi-label learning, where each instance is associated with a set of correlated labels rather than a single label, is widely applicable in various domains such as object detection [13] and text classification [43]. For example, an image may contain multiple relevant objects and a text can encompass various topics [5, 24, 25, 40]. The key to multi-label learning is to model the multi-label correlation. Early studies convert the multi-label learning problem into multiple single-label learning problems, which has limitations in its ability to explicitly model correlations [6, 29]. Recently, various studies have been proposed with deep representation learning. Among them, one type of method learns correlated representations with statistical analysis [37, 48]. The second type of method learns the label correlation on label embeddings with label sequence/relational modeling techniques, such as RNN and GNN [10, 36, 47]. Another type of

model proposes to co-learn label and instance representations with the auto-encoder architecture, which models the correlation during the latent encoding [1, 2, 54].

Despite significant progress, these methods are all focused on Euclidean data. In recent years, due to the development of GNNs, the problem of multi-label learning on non-Euclidean graphs has also shown its importance [53]. However, thus far, multi-label learning on graph data has not received much attention.

2.2 Multi-Label Node Classification

Node classification is a well-known fundamental task in graph data mining. Typically, node classification refers to the assignment of a unique label to each node [4, 14]. In recent years, multi-label node classification (MLNC) has posed significance due to the increasing realization that nodes on graphs often exhibit multiple associated categories simultaneously [11, 53]. For example, in social network analysis, a user often belongs to multiple interest groups [8, 31, 41]. This involves assigning non-unique and variable numbers of labels to each node. Currently, MLNC is still in its infancy, and only a few studies have focused on MLNC [31, 53]. Representatively, ML-GCN [11] learns distinct embeddings for each label and constructs node-label correlations to augment the GCN-encoded node embeddings. Subsequently, LANC [55] and LARN [44] further incorporate attention mechanisms to integrate label embeddings with GCN-encoded node embeddings more effectively.

Nonetheless, current approaches continue to utilize a unified graph convolution process across neighborhoods under various labels, which ignores the possibility that information among different multi-label nodes may become intermixed in a unified message passing process, including both correlated and uncorrelated labels.

2.3 Graph Structure Learning

Since the natural graph structure often contains noise information, it's not always directly applicable to a wide range of downstream tasks. This limitation has brought the field of Graph Structure Learning (GSL) into the spotlight, garnering increasing attention in recent years [3, 17, 21]. Representatively, GRCN [49] adjusts edge weights in the graph, optimizing with downstream tasks. IDGL [9] iteratively improves the graph structure through node embeddings. And SUBLIME [26] uses a contrastive-based objective function to guide the learned graph. These GSL models generally strive to refine the graph structure by enhancing graph homophily or by augmenting task-specific useful information in the single-label space. Despite these advances, there is a noticeable absence of GSL approaches specifically tailored for MLNC, leaving a significant gap in the field.

In light of these challenges, our paper is dedicated to introducing a GSL-based methodology, uniquely designed for MLNC. This method aims to address the existing limitations by integrating the understanding of multi-label correlations and adapting the graph structure accordingly.

3 PRELIMINARIES

Notations. Given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, consisting of two sets: a set of nodes \mathcal{V} and a set of edges \mathcal{E} . Here, \mathcal{V} represents the set of nodes in the graph, while $\mathcal{E} \subseteq \{(u, v) | u, v \in \mathcal{V}\}$ represents the

set of edges between the nodes. Then, the adjacency matrix of \mathcal{G} can be defined as $A \in \mathbb{R}^{n \times n}$ with m non-zero element in the matrix based on the edge set, where $n = |\mathcal{V}|$ and $m = |\mathcal{E}|$ are the node number and edge number, respectively. Each node $v_i \in \mathcal{V}$ in the graph \mathcal{G} is associated with a feature vector $\mathbf{x}_i \in \mathbb{R}^f$ and thus all the node feature vectors form the overall feature matrix $X \in \mathbb{R}^{n \times f}$ of the graph \mathcal{G} . To distinguish from the notation of decomposed graphs introduced in this paper, we represent the original graph adjacent matrix with the symbol A^0 .

Definition 1 (Graph with Multi-Label Nodes). In a multi-label graph, each node $v_i \in \mathcal{V}$ is associated with a set of labels $L_i \subseteq L$, where L is the set of all possible labels and $|L| = K$ is the total number of labels. Thus, each node can have multiple labels, i.e., $|L_i| \geq 1$. In practice, each node v_i will be assigned a multi-hot label vector $\mathbf{y}_i \in \mathbb{R}^{1 \times K}$, where each element $y_{i,j} \in [0, 1]$ in \mathbf{y}_i indicates whether node v_i belongs to class j .

Definition 2 (Multi-Label Node Classification). In the multi-label node classification scenario, we typically have a subset of nodes with known labels and another subset with unknown labels [53]. The node set \mathcal{V} in \mathcal{G} can be divided into two disjoint subsets: a labeled node set \mathcal{V}_L and an unlabeled node set \mathcal{V}_U . Thus, $\mathcal{V} = \mathcal{V}_L \cup \mathcal{V}_U$ and $\mathcal{V}_L \cap \mathcal{V}_U = \emptyset$. The goal of multi-label node classification is to learn a model $f : \mathcal{V} \rightarrow \{0, 1\}^K$ using the labeled nodes $v_i \in \mathcal{V}_L$ and the graph topology A . This model aims to predict the label set \mathbf{y}_j for unlabeled nodes $v_j \in \mathcal{V}_U$.

4 METHODOLOGY

Our proposed model for multi-label node classification comprises the correlation-aware graph learning module with node feature and graph topology decomposition (Sec. 4.1) and the correlation-aware graph convolution (Sec. 4.2 and Sec. 4.3). Further, we illustrate how CorGCN can be expanded to the scenario with large label space (Sec. 4.4). Figure 2 provides an overview of the architecture of our proposed CorGCN.

4.1 Correlation-Aware Graph Decomposition

Due to the ambiguity in node features and topological structures of graphs in multi-label scenarios, we aim to decompose them into multiple graphs. However, directly performing the decomposition would result in the loss of the critically important multi-label correlation property. Therefore, we need to carry out the decomposition based on representations that have encapsulated label correlation.

4.1.1 Correlation-Aware Feature Decomposition. This component primarily adopts an approach that first models the node-label correlation and label-label correlation, and then, based on this, performs the decomposition of node features across different label spaces. Summarizing these, we model the correlations with contrastive mutual information estimator ϕ and likelihood maximizing decoder θ , which can be defined as,

$$\arg \max_{\phi, \theta} \underbrace{I_{\phi}(E^X, E^L)}_{\text{Mutual Information}} + \underbrace{\mathcal{L}_{\theta}(\theta; E^X) + \mathcal{L}_{\theta}(\theta; E^L)}_{\text{Likelihood}}. \quad (1)$$

Specifically, to fully utilize the knowledge from labeled nodes for generalization to all nodes in the graph, we first assign a trainable

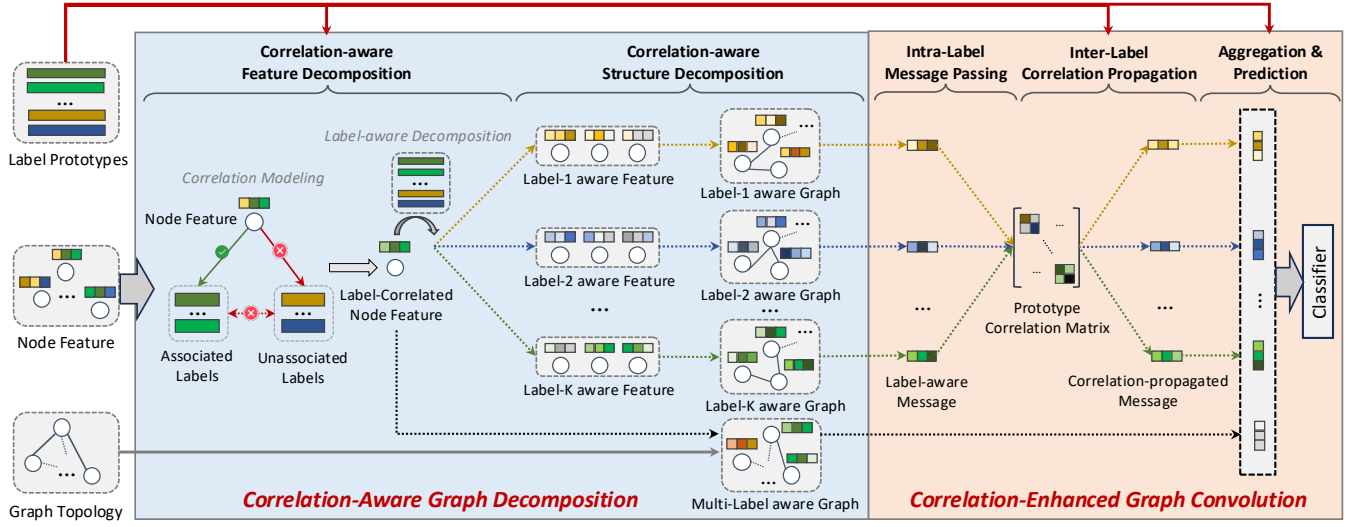


Figure 2: The overall architecture of CorGCN. (a) Correlation-Aware Graph Decomposition: it first learns label-correlated node features and decomposes them into multiple label-aware features. Then, based on the decomposed features, it decomposes multiple label-aware graphs. **(b) Correlation-Enhanced Graph Convolution:** each layer includes intra-label message passing of the neighborhood in each label-aware graph view and inter-label correlation propagation between label-aware messages.

prototype embedding to each label to describe the characteristics of this label:

$$E^l = [E_1^l, \dots, E_k^l, \dots, E_K^l], \quad (2)$$

where $1 \leq k \leq K$, and $E_k^l \in \mathbb{R}^d$ is the k -th label prototype. We then transform the original node features X to the same d -dimensional space using the linear transformation: $E^x = XW_t$, where $W_t \in \mathbb{R}^{f \times d}$ is a learnable transformation matrix, and each vector $E_i^x \in \mathbb{R}^d$ in E^x denotes the transformed feature of node v_i :

$$E^x = [E_1^x, \dots, E_i^x, \dots, E_n^x]. \quad (3)$$

Then, the mutual information estimator ϕ captures both node-label and label-label correlations based on those supervised nodes. Previous works for multi-label learning have successfully adopted contrastive learning to model multi-label correlations and inspired by [2], we utilize contrastive learning as ϕ here. Given the labeled nodes \mathcal{V}_L , the contrastive-based mutual information estimator with \mathcal{L}_{cmi} can be written as follows:

$$\mathcal{L}_{cmi} = -\frac{1}{|\mathcal{V}_L|} \sum_{i \in \mathcal{V}_L} \frac{1}{|\text{Pos}(\mathbf{y}_i)|} \sum_{p \in \text{Pos}(\mathbf{y}_i)} \log \frac{\exp(E_i^x \cdot E_p^l)}{\sum_{t \in \mathbf{y}_i} \exp(E_i^x \cdot E_t^l)}, \quad (4)$$

where $\text{Pos}(\mathbf{y}_i)$ denotes the label set that node v_i belongs to. Through this formulation, we directly map node features into a compact space that reflects multi-label correlations, while implicitly capturing the nuances of multi-label prototype relationships. Unlike single-label scenarios where classes are mutually exclusive, multi-label classification acknowledges the complexity of label correlations. Therefore, we do not enforce strict contrastive relations

among labels, maintaining their intrinsic associations [2, 54]. Instead, it utilizes feature embeddings as anchors and label embeddings as positive or negative examples, allowing frequently co-occurring labels to become more similar over time.

Subsequently, due to the above process mainly focuses on embedding space without explicit guidance, in order to ensure that different node features and label prototypes retain their own classification characteristics, we further constrain this through likelihood maximizing decoder with corresponding loss \mathcal{L}_{lm} as follows:

$$\rho_k = \frac{\sqrt{\frac{1}{\sum_{v_i \in \mathcal{V}_L} \mathbf{y}_{i,k}}}}{\sum_{j=1}^K \sqrt{\frac{1}{\sum_{v_i \in \mathcal{V}_L} \mathbf{y}_{i,j}}}}, \quad (5)$$

$$\mathcal{L}_{lm} = -\frac{1}{2|\mathcal{V}_L|} \sum_{v_i \in \mathcal{V}_L} \sum_{k=1}^K \rho_k [(1 - p_k(E_i^x | \theta))^\gamma \cdot \log p_k(E_i^x | \theta) + (1 - p_k(\mathbf{y}_i \cdot E^l | \theta))^\gamma \cdot \log p_k(\mathbf{y}_i \cdot E^l | \theta)], \quad (6)$$

where ρ_k is a statistical parameter to enhance the learning for classes with a small number of node samples and γ is a hyper-parameter to control the observation of hard node/label samples [23]. $\hat{p}_k(E_i^x | \theta)$ is the predicted result of the decoder θ , and if $y_{i,k} = 1$, $p_k(E_i^x | \theta) = \hat{p}_k(E_i^x | \theta)$. Otherwise, $p_k(E_i^x | \theta) = 1 - \hat{p}_k(E_i^x | \theta)$.

In the above process, we obtain the correlated node features E^x and label prototypes E^l . Each label prototype is associated with a specific label and its potentially correlated labels while each correlated node feature is still ambitious due to the mixture of multiple label information. Based on this, we decompose the correlated node features to each label-friendly feature space that *considers both the label-specific and its correlated label-specific information* with the

guidance of label prototypes. Specifically, given the label-correlated node feature E_i^x of node v_i and the label prototype of label k , the projection process can be obtained by:

$$w_{i,k} = \text{sim}(E_i^x, E_k^l) = \frac{E_i^x \cdot E_k^l}{\|E_i^x\| \|E_k^l\|}, \quad (7)$$

$$\begin{aligned} E_i^{proj} &= [E_{i,1}^{proj}, \dots, E_{i,k}^{proj}, \dots, E_{i,K}^{proj}] \\ &= [w_{i,1} E_i^x, \dots, w_{i,k} E_i^x, \dots, w_{i,K} E_i^x], \end{aligned} \quad (8)$$

where $w_{i,k}$ is the projection coefficient of node i to label k , $E_{i,k}^{proj}$ is the projected representation of node v_i towards the k -th label.

4.1.2 Correlation-Aware Structure Decomposition. Based on the correlation-aware decomposed node features E^{proj} , we aim to decompose the graph structure (message passing path) for each label with its correlated labels.

To preserve the graph topology patterns, we first aggregate the neighborhood projected features of each center node v_i in each label view k for structure decomposition:

$$E_{i,k}^{sd} = \text{Agg}(E^{proj}, A^0) = \frac{1}{|\mathcal{N}_i^0| + 1} \sum_{j \in \{\mathcal{N}_i^0 \cup \{v_i\}\}} E_{j,k}^{proj}, \quad (9)$$

where \mathcal{N}_i^0 is the neighborhood of node v_i in original graph adjacent matrix A^0 . We leverage the aggregated representations to measure the node similarity in each label view:

$$S_{i,j}^k = \text{sim}(E_{i,k}^{sd}, E_{j,k}^{sd}) = \frac{E_{i,k}^{sd} \cdot E_{j,k}^{sd}}{\|E_{i,k}^{sd}\| \|E_{j,k}^{sd}\|}, \quad (10)$$

where S^k is the similarity score matrix in label view- k , $S_{i,j}^k$ is the score between node v_i and node v_j .

Then, for each node v_i , the structure decomposing for different label-aware graphs can be described as:

$$A_{i,j}^k = \begin{cases} 1, & S_{i,j}^k \in \text{Top-}\lambda(S_i^k); \\ 0, & \text{otherwise.} \end{cases} \quad (11)$$

$$\mathcal{G}^k = (A^k, E_k^{proj}), \quad (12)$$

where A^k is the adjacent matrix of the k -th label view, and λ is the hyperparameter that controls the density of the decomposed graph.

Furthermore, with the correlated label-aware node feature and the original graph topology, the multi-label aware graph $\mathcal{G}^0 = (A^0, E^x)$ can be obtained to capture correlated structure patterns by message passing, and thus the learned correlation-aware decomposed graphs CDG can be obtained as follows:

$$CDG = \{\mathcal{G}^0, \mathcal{G}^1, \dots, \mathcal{G}^K\}. \quad (13)$$

4.2 Correlation-Enhanced Graph Convolution

Previous studies for multi-label node classification mainly conduct unified neighborhood message passing [11, 44, 55]. We contend this approach has two primary limitations: (1) passing ambiguous messages from the neighborhood in a unified manner; and (2) correlation ignorance, which overlooks the label correlation when passing the message. Therefore, we further equip the CDG with Correlation-Enhanced graph convolution.

4.2.1 Intra-Label Message Passing. Firstly, Correlation-Enhanced graph convolution conducts intra-label message passing within the graph \mathcal{G}^k of each label view k :

$$\begin{aligned} \hat{Z}_{[:,k]}^{(l)} &= \text{GCN}(Z_{[:,k]}^{(l-1)}, \tilde{A}^k) \\ &= \sigma(\tilde{A}^k Z_{[:,k]}^{(l-1)} \mathbf{W}^{(l)}), \end{aligned} \quad (14)$$

where $0 \leq k \leq K$, $\tilde{A}^k = \hat{D}^{k-\frac{1}{2}} A^{*k} \hat{D}^{k-\frac{1}{2}} \in \mathbb{R}^{n \times n}$ is the normalized adjacency matrix of label view- k , $\hat{D}^k \in \mathbb{R}^{n \times n}$ is the degree matrix of $A^{*k} = A^k + I$ where I is the identity matrix. $\hat{Z}^{(l)}$, $\mathbf{W}^{(l)}$ is the output features, trainable parameters in l -th message passing layer, respectively. The input of the first layer $Z^{(0)} = [E^x, E_1^{proj}, \dots, E_K^{proj}] \in \mathbb{R}^{n \times (K+1) \times d}$ is the feature matrix. Note that this is a GCN-like [20] message passing function, any other graph message passing functions [16, 34] can also be adopted here.

4.2.2 Inter-Label Correlation Propagation. After the above intra-label message passing, we obtain the neighborhood messages $\hat{Z}^{(l)} \in \mathbb{R}^{n \times K \times d}$ from each label-aware view. Then, to model their correlations, we further propose the inter-label correlation propagation between each label-aware graph view:

$$\text{Cor}_i = \text{Softmax}\left(\frac{(E^l \mathbf{W}_1)(\hat{Z}_{[i,1]}^{(l)} \mathbf{W}_2)^T}{\sqrt{d}}\right), \quad (15)$$

$$Z_{[i,1]}^{(l)} = \text{Cor}_i \hat{Z}_{[i,1]}^{(l)} \mathbf{W}_3, \quad (16)$$

where $\text{Cor}_i \in \mathbb{R}^{K \times K}$ is the label prototype correlation matrix for node v_i , $\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3$ are the trainable parameters, and $Z_{[i,1]}^{(l)} \in \mathbb{R}^{K \times d}$ is the output of v_i in l -th graph convolution layer containing K views of inter-label correlated representations.

4.2.3 Aggregation & Prediction. After the multi-layer graph message passings, here we aggregate these representations to obtain the final representation with both the node representations and label prototypes. The aggregation process can be written as:

$$Z_i^{cls} = [Z_{i,0} || (\sum_{k=1}^K \text{sim}(Z_{i,k}, E_k^l) \cdot Z_{i,k})], \quad (17)$$

$$\hat{y}_i = \sigma(\mathbf{W}^{cls} Z_i^{cls} + \mathbf{b}^{cls}), \quad (18)$$

where $Z_{i,k}$ is output features from final layers of correlation propagation for node v_i in the label view k , $Z_i^{cls} \in \mathbb{R}^{n \times 2d}$ is the final node representation for multi-label node classification of node v_i and $||$ is the concatenation operation, \mathbf{W}^{cls} and \mathbf{b}^{cls} are the trainable parameters in the classifier, and σ is the Sigmoid function.

4.3 Objective Function

To train and optimize the model parameters, we apply the binary cross-entropy loss as the model objective function. Formally, for each node $v_i \in \mathcal{V}_L$, the classification objective function can be expressed as:

$$\mathcal{L}_{cls} = -\frac{1}{|\mathcal{V}_L|} \sum_{i \in \mathcal{V}_L} \frac{1}{K} \sum_{j=1}^K y_{i,j} \log(\hat{y}_{i,j}) + (1-y_{i,j}) \log(1-\hat{y}_{i,j}), \quad (19)$$

where $\hat{y}_{i,j}$ is the predicted logits of node v_i towards label j and $y_{i,j}$ is the corresponding ground-truth label. Then, considering

the objective losses in each CorGCN module, the overall objective function of CorGCN is as follows:

$$\mathcal{L} = \mathcal{L}_{cls} + \alpha \cdot \mathcal{L}_{cmi} + \beta \cdot \mathcal{L}_{lm}, \quad (20)$$

where α and β are the controllable hyper-parameters.

4.4 Extension to Large Label Space

In some real-world graph structures, nodes may exist in a very large multi-label space, such as in large protein interaction networks, where individual protein nodes may have hundreds of characteristic labels simultaneously [18, 30]. In this subsection, we demonstrate that our proposed CorGCN can be readily extended to multi-label node classification with large label space for efficient learning.

Macro Label Prototypes. For the large label space, directly applying multi-label graph learning with K prototypes will obtain hundreds of label-aware graphs, which is not efficient enough for applications. Therefore, we aim to refine and cluster the original prototypes into macro-label prototypes, with the principle that correlated labels will form a macro prototype.

Specifically, we first pre-train the K label prototypes in Eq.(1) with \mathcal{L}_{cmi} and \mathcal{L}_{lm} (as Sec. 4.1). Then, we adopt K-means-based clustering to refine the K pre-trained label prototypes into K' macro prototypes [15]. Given the pre-trained label prototypes E^l , we randomly initialize K' macro prototype centroids $\{\mu_1, \dots, \mu_k, \dots, \mu_{K'}\}$, where $\mu_k \in \mathbb{R}^d$ is the centroids of macro label prototype C_k , $K' \ll K$ is the hyperparameter set as the macro prototype number. Then, we explore and assign the original label prototypes to the appropriate macro label prototype based on their representation similarity, and update the centroid of the macro label prototype iteratively. The process can be expressed as:

$$\mu_k = \frac{1}{|C_k|} \sum_{x_v=k, E_v^l \in C_k} E_v^l, \quad (21)$$

where $|C_k|$ is the number of labels within C_k , and x_v is the macro label prototype index that v is assigned to. Further, the optimization target of the macro label prototype generation is:

$$\begin{aligned} & \min_{x_1, \dots, x_K; \mu_1, \dots, \mu_{K'}} J(x_1, \dots, x_K; \mu_1, \dots, \mu_{K'}) \\ & \triangleq \sum_{k=1}^{K'} \sum_{x_v=k, E_v^l \in C_k} \sqrt{(E_v^l - \mu_k)(E_v^l - \mu_k)^\top}, \end{aligned} \quad (22)$$

where J is the objective function of macro label prototype generation. In this way, we obtain the macro label prototype set:

$$E^{ml} = [\mu_1, \dots, \mu_k, \dots, \mu_{K'}], \quad (23)$$

where $1 \leq k \leq K' \ll K$. Therefore, we can reduce the original number of label prototypes to meet our needs. With the obtained E^{ml} , we can flexibly replace E^l with E^{ml} and procedure Sec. 4.1 - Sec. 4.3 in the same way. Model analyses are detailed in Appendix A.

5 EXPERIMENTS

In this section, we conduct comprehensive experiments on benchmark datasets with CorGCN, aiming to answer the following questions: **RQ1:** How does CorGCN perform compared to representative and state-of-the-art models? **RQ2:** What is the effect of different components in CorGCN? **RQ3:** Can CorGCN be generalized and

adapted to different GNN message passing backbones? **RQ4:** How does the impact of CorGCN on each fine-grained class? **RQ5:** What is the efficiency of CorGCN on training and inference? **RQ6:** How do key hyper-parameters impact the performance of CorGCN?

5.1 Experimental Setup

5.1.1 Datasets. We conduct experiments on five widely used benchmark datasets: **Humloc** [53], **PCG** [53], **Blogcatalog** [55], **PPI** [50] (a dataset with large label space), and **Delve** [44] to verify the effectiveness of CorGCN. The statistics of these datasets are shown in Table 1. We randomly divided the labeled data of these datasets into training, validation, and test sets, following a 6:2:2 ratio, respectively [53]. The performance on PPI and Delve are conducted on the extension of CorGCN, with 20 and 10 macro label prototypes, respectively. Detailed descriptions of these datasets are illustrated in Appendix B.1.

Table 1: Statistics of the experimental datasets.

Dataset	# Nodes	# Edges	# Features	# Classes	Density
Humloc	3,106	18,496	32	14	0.3836%
PCG	3,233	37,351	32	15	0.7149%
Blogcatalog	10,312	333,983	100	39	0.6282%
PPI	14,755	225,270	50	121	0.2070%
Delve	1,229,280	4,322,275	300	20	0.0006%

5.1.2 Baselines. To evaluate the effectiveness of CorGCN, we compare it with nine representative and state-of-the-art models. (i) Traditional GNN for node classification: **GCN** [20]. Due to the traditional workflow of GCN being unsuitable for multi-label node classification, we replace the last activation function from softmax to sigmoid. (ii) Graph structure learning methods: **GRCN** [49], **IDGL** [9], and **SUBLIME** [26]. We equip them with the GCN backbone for fair comparisons. (iii) Multi-label node classification methods: **GCN-LPA** [35], **ML-GCN** [11], **LANC** [55], **SMLG** [31], and **LARN** [44]. The details of these baselines are left in Appendix B.2.

5.1.3 Implementation Setting. For all models, the embedding size is fixed to 64 for fair node classification. The learning rate of CorGCN is searched from $\{1 \times 10^{-2}, 5 \times 10^{-3}, 1 \times 10^{-3}\}$, the regularization term is searched from $\{1 \times 10^{-4}, 5 \times 10^{-5}, 1 \times 10^{-5}\}$. The batch size is set to 1024 for all models and the Adam optimizer is used. The detailed experiment implementations are posed in Appendix B.3. Source code is available at <https://github.com/YuanchenBei/CorGCN>.

5.1.4 Evaluation Metrics. We evaluate the model under the setting for semi-supervised multi-label learning on graphs as previous works [31]. We evaluate the models with seven widely-adopted metrics including *Ranking Loss*, *Hamming Loss*, *Macro-AUC*, *Micro-AUC*, *Macro-AP*, *Micro-AP*, and *Label Ranking AP* (abbreviated as *Ranking*, *Hamming*, *Ma-AUC*, *Mi-AUC*, *Ma-AP*, *Mi-AP*, and *LPAP* in the following subsections, respectively). All these multi-label node classification metrics are clearly defined in [52]. Note that we run all the experiments *five* times with different random seeds and report the average results with standard deviation.

Table 2: Multi-label node classification comparison results in percentage over five trial runs (\uparrow : the higher, the better; \downarrow : the lower, the better). The best and second-best results in each column are highlighted in bold font and underlined. OOM denotes out-of-memory during the model training. The average ranking is calculated based on the numerical result of each model.

Dataset	Metrics	GCN	GRCN	IDGL	SUBLIME	GCN-LPA	ML-GCN	LANC	SMLG	LARN	CorGCN
Humloc	Ranking (\downarrow)	13.29 \pm 1.01	16.84 \pm 0.09	19.04 \pm 0.13	14.44 \pm 0.10	19.58 \pm 4.96	<u>13.28 \pm 0.70</u>	14.66 \pm 0.36	17.82 \pm 0.43	13.62 \pm 0.46	12.57 \pm 0.31
	Hamming (\downarrow)	7.52 \pm 0.13	8.37 \pm 0.03	8.40 \pm 0.02	8.53 \pm 0.03	8.30 \pm 0.13	<u>7.46 \pm 0.22</u>	7.82 \pm 0.16	8.55 \pm 0.26	8.04 \pm 0.29	7.37 \pm 0.17
	Ma-AUC (\uparrow)	69.10 \pm 2.25	56.27 \pm 0.23	50.26 \pm 0.94	72.42 \pm 0.40	57.91 \pm 2.39	72.41 \pm 2.47	68.73 \pm 1.76	<u>72.64 \pm 1.65</u>	71.63 \pm 3.30	77.31 \pm 1.58
	Mi-AUC (\uparrow)	85.39 \pm 1.30	82.85 \pm 0.02	78.48 \pm 0.17	85.63 \pm 0.07	80.97 \pm 5.29	<u>87.63 \pm 0.52</u>	86.06 \pm 0.47	79.06 \pm 0.38	86.63 \pm 0.50	88.57 \pm 0.37
	Ma-AP (\uparrow)	<u>24.65 \pm 1.29</u>	14.11 \pm 0.12	9.65 \pm 0.35	20.14 \pm 0.13	13.70 \pm 1.72	23.78 \pm 1.17	20.94 \pm 1.69	21.17 \pm 0.61	20.60 \pm 0.34	26.91 \pm 1.67
	Mi-AP (\uparrow)	46.46 \pm 2.06	35.11 \pm 0.10	26.31 \pm 0.34	36.61 \pm 0.13	31.88 \pm 7.37	<u>47.76 \pm 1.59</u>	42.10 \pm 2.61	35.53 \pm 0.18	41.30 \pm 1.22	48.97 \pm 1.30
LRAP (\uparrow)	64.66 \pm 0.80	55.45 \pm 0.19	52.04 \pm 0.13	59.71 \pm 0.25	52.50 \pm 3.97	<u>64.91 \pm 1.14</u>	61.02 \pm 1.45	54.27 \pm 0.38	62.46 \pm 1.46	65.40 \pm 0.83	
PCG	Ranking (\downarrow)	27.50 \pm 0.58	27.30 \pm 0.04	27.97 \pm 0.10	27.31 \pm 0.16	28.55 \pm 0.34	<u>27.11 \pm 0.44</u>	28.26 \pm 0.51	28.15 \pm 0.56	28.45 \pm 0.53	25.97 \pm 0.57
	Hamming (\downarrow)	13.22 \pm 0.30	12.96 \pm 0.00	13.20 \pm 0.01	13.16 \pm 0.01	12.59 \pm 0.20	<u>12.56 \pm 0.31</u>	12.59 \pm 0.18	16.28 \pm 0.16	12.65 \pm 0.14	12.41 \pm 0.23
	Ma-AUC (\uparrow)	62.92 \pm 0.84	48.03 \pm 0.36	45.60 \pm 0.27	57.32 \pm 0.41	53.84 \pm 0.90	61.48 \pm 0.84	58.75 \pm 1.05	<u>63.02 \pm 0.97</u>	57.00 \pm 1.14	64.86 \pm 1.05
	Mi-AUC (\uparrow)	71.74 \pm 0.38	67.61 \pm 0.13	62.86 \pm 0.13	70.89 \pm 0.09	70.09 \pm 0.38	<u>72.03 \pm 0.40</u>	71.31 \pm 0.56	71.08 \pm 0.62	70.02 \pm 0.47	74.16 \pm 0.61
	Ma-AP (\uparrow)	<u>23.49 \pm 0.38</u>	13.12 \pm 0.05	12.68 \pm 0.12	18.45 \pm 0.29	16.01 \pm 0.69	20.95 \pm 0.91	19.49 \pm 0.85	22.73 \pm 0.75	16.23 \pm 0.66	24.64 \pm 0.90
	Mi-AP (\uparrow)	<u>30.04 \pm 0.54</u>	23.43 \pm 0.09	21.32 \pm 0.28	28.04 \pm 0.29	25.50 \pm 0.95	29.33 \pm 1.21	27.54 \pm 0.98	24.77 \pm 0.85	24.87 \pm 0.59	31.91 \pm 1.07
LRAP (\uparrow)	48.03 \pm 1.25	46.42 \pm 0.13	46.76 \pm 0.24	<u>48.23 \pm 0.26</u>	45.37 \pm 1.03	48.21 \pm 0.57	46.62 \pm 1.27	47.58 \pm 0.86	45.76 \pm 0.99	49.04 \pm 0.83	
Blogcatalog	Ranking (\downarrow)	25.67 \pm 0.20	25.69 \pm 0.02	25.95 \pm 0.01	<u>25.48 \pm 0.01</u>	42.19 \pm 3.87	25.68 \pm 0.25	<u>25.48 \pm 0.12</u>	26.61 \pm 0.86	25.67 \pm 0.27	25.42 \pm 0.23
	Hamming (\downarrow)	3.58 \pm 0.03	3.57 \pm 0.00	<u>3.56 \pm 0.00</u>	3.59 \pm 0.00	3.58 \pm 0.03	3.58 \pm 0.03	3.55 \pm 0.03	3.66 \pm 0.04	3.58 \pm 0.03	3.58 \pm 0.03
	Ma-AUC (\uparrow)	50.59 \pm 0.51	48.19 \pm 0.04	47.08 \pm 0.01	50.52 \pm 0.02	50.19 \pm 2.06	50.94 \pm 1.17	<u>52.39 \pm 0.52</u>	51.70 \pm 0.32	50.52 \pm 1.07	54.48 \pm 0.52
	Mi-AUC (\uparrow)	73.80 \pm 0.20	69.22 \pm 0.05	65.28 \pm 0.01	72.37 \pm 0.01	56.74 \pm 2.32	73.85 \pm 0.17	<u>74.10 \pm 0.07</u>	71.75 \pm 0.86	73.73 \pm 0.47	74.15 \pm 0.21
	Ma-AP (\uparrow)	4.16 \pm 0.14	3.67 \pm 0.01	3.67 \pm 0.01	3.91 \pm 0.02	4.10 \pm 0.24	4.16 \pm 0.14	5.07 \pm 0.20	4.13 \pm 0.13	4.15 \pm 0.48	<u>4.61 \pm 0.15</u>
	Mi-AP (\uparrow)	9.41 \pm 0.15	7.97 \pm 0.01	7.13 \pm 0.02	8.89 \pm 0.02	4.72 \pm 0.53	9.42 \pm 0.15	<u>9.60 \pm 1.08</u>	5.54 \pm 1.33	9.38 \pm 1.22	9.65 \pm 0.21
LRAP (\uparrow)	27.88 \pm 0.15	27.30 \pm 0.00	27.82 \pm 0.04	27.81 \pm 0.00	17.58 \pm 3.93	27.87 \pm 0.18	<u>28.19 \pm 1.01</u>	25.37 \pm 0.27	28.10 \pm 0.92	28.32 \pm 0.29	
PPI	Ranking (\downarrow)	<u>18.32 \pm 0.16</u>	25.73 \pm 0.01	OOM	25.44 \pm 0.01	25.33 \pm 0.08	20.05 \pm 0.62	18.44 \pm 0.14	18.38 \pm 0.30	19.85 \pm 0.18	16.17 \pm 0.22
	Hamming (\downarrow)	22.63 \pm 0.14	26.28 \pm 0.00	OOM	26.24 \pm 0.00	25.60 \pm 0.18	23.83 \pm 0.14	<u>21.85 \pm 0.28</u>	31.06 \pm 0.32	23.41 \pm 0.13	20.79 \pm 0.29
	Ma-AUC (\uparrow)	73.06 \pm 0.21	44.65 \pm 0.02	OOM	51.82 \pm 0.12	58.86 \pm 0.29	70.26 \pm 1.12	<u>74.24 \pm 0.42</u>	74.04 \pm 0.41	70.21 \pm 0.24	77.54 \pm 0.42
	Mi-AUC (\uparrow)	80.19 \pm 0.14	67.84 \pm 0.01	OOM	70.24 \pm 0.02	72.48 \pm 0.08	78.37 \pm 0.72	<u>81.15 \pm 0.26</u>	80.46 \pm 0.26	78.50 \pm 0.13	83.35 \pm 0.29
	Ma-AP (\uparrow)	54.70 \pm 0.42	29.34 \pm 0.01	OOM	34.07 \pm 0.01	39.06 \pm 0.30	50.69 \pm 1.54	<u>57.19 \pm 0.15</u>	55.71 \pm 0.70	51.16 \pm 0.61	61.33 \pm 0.67
	Mi-AP (\uparrow)	67.64 \pm 0.43	55.66 \pm 0.01	OOM	56.52 \pm 0.01	58.68 \pm 0.37	64.98 \pm 1.32	<u>69.59 \pm 0.27</u>	63.11 \pm 0.62	65.57 \pm 0.54	72.35 \pm 0.58
LRAP (\uparrow)	68.82 \pm 0.39	62.19 \pm 0.01	OOM	62.28 \pm 0.01	62.18 \pm 0.40	66.68 \pm 1.01	<u>68.87 \pm 0.35</u>	63.68 \pm 0.61	67.96 \pm 0.36	71.40 \pm 0.34	
Delve	Ranking (\downarrow)	3.46 \pm 0.04	OOM	OOM	6.10 \pm 0.53	OOM	3.33 \pm 0.01	OOM	5.07 \pm 0.02	1.46 \pm 0.02	<u>2.41 \pm 0.02</u>
	Hamming (\downarrow)	3.16 \pm 0.01	OOM	OOM	30.26 \pm 0.51	OOM	2.99 \pm 0.01	OOM	7.50 \pm 0.01	1.67 \pm 0.07	<u>2.48 \pm 0.01</u>
	Ma-AUC (\uparrow)	94.98 \pm 0.14	OOM	OOM	58.37 \pm 0.74	OOM	95.71 \pm 0.04	OOM	95.26 \pm 0.09	97.85 \pm 0.05	<u>96.84 \pm 0.10</u>
	Mi-AUC (\uparrow)	96.55 \pm 0.06	OOM	OOM	57.58 \pm 0.04	OOM	97.29 \pm 0.02	OOM	95.11 \pm 0.04	98.39 \pm 0.07	<u>98.05 \pm 0.02</u>
	Ma-AP (\uparrow)	64.28 \pm 0.39	OOM	OOM	11.60 \pm 0.07	OOM	65.20 \pm 0.07	OOM	65.67 \pm 0.31	80.16 \pm 0.32	<u>72.80 \pm 0.26</u>
	Mi-AP (\uparrow)	78.72 \pm 0.07	OOM	OOM	10.12 \pm 0.07	OOM	81.78 \pm 0.09	OOM	68.97 \pm 0.25	90.44 \pm 0.12	<u>86.12 \pm 0.03</u>
LRAP (\uparrow)	85.57 \pm 0.10	OOM	OOM	29.44 \pm 0.40	OOM	86.12 \pm 0.05	OOM	76.47 \pm 0.27	92.09 \pm 0.21	<u>89.20 \pm 0.01</u>	
Average Rank		3.86	7.79	8.52	6.17	7.79	<u>3.57</u>	<u>3.57</u>	5.94	4.46	1.31

5.2 Main Results (RQ1)

In this subsection, we compare our proposed CorGCN with nine state-of-the-art baseline models on the five experimental datasets. The comparison results on seven metrics are reported in Table 2, corresponding Friedman statistics are shown in Table 3, and Bonferroni-Dunn test can be found in Appendix C.1, with the following observations:

- **CorGCN can achieve significant improvements over state-of-the-art methods on all experimental datasets.** From the table, we can observe that our proposed CorGCN can generally achieve the best performance among the seven metrics on average. Specifically, CorGCN outperforms the best baseline by 6.43%, 2.92%, 3.99%, and 4.14% for Macro-AUC on Humloc, PCG, Blogcatalog, and PPI, respectively. The significance of model gains was further validated by Friedman statistics and the Bonferroni-Dunn test. These results verify that designing the correlation

decomposed graph and the equipped CorGCN contributes to achieving better multi-label node classification performance.

- **Graph structure learning models do not bring about performance gains.** Comparing three graph structure learning models (GRCN, IDGL, SUBLIME) with the simply adopted GCN, we can find that these methods fail to yield substantial benefits in multi-label node classification scenarios and, in many instances, may even lead to negative outcomes. For example, all three models negatively impact the GCN model performance in terms of Micro-AUC, Macro-AP, and Micro-AP. This indicates that directly learning graph structures without considering multi-label correlations is less useful for the multi-label node classification downstream task.
- **The models with multi-label modeling generally become the best baselines.** We can observe from the results that LANC and ML-GCN achieve the best performance among the other baseline models. Their shared advantage lies in the explicit modeling

Table 3: Friedman statistics F_F on each evaluation metric with corresponding p -value (all < 0.05 level).

Evaluation Metric	F_F	p -value
Ranking Loss	17.59	0.0015
Hamming Loss	21.97	0.0002
Macro-AUC	12.14	0.0163
Micro-AUC	13.41	0.0095
Macro-AP	14.85	0.0050
Micro-AP	16.54	0.0024
Label Ranking AP	16.54	0.0024

of multi-labels. However, due to the unified message passing with ambitious information, they are less optimal compared to the CorGCN. Our proposed CorGCN explicitly considers the label correlation and addresses the ambitious message passing with information aggregation over the correlation-aware decomposed graph, which achieves the best performance.

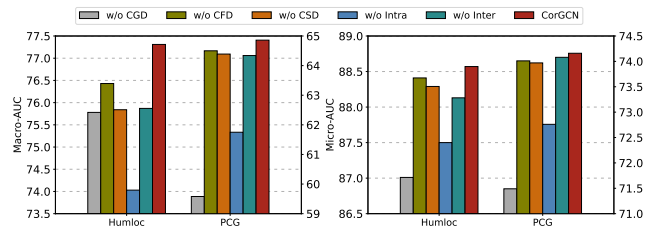
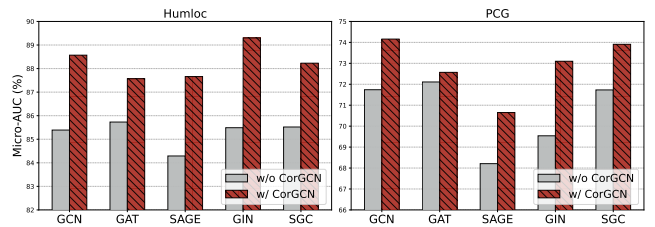
We have also assessed the convergence of CorGCN, with the results detailed in Appendix C.2. The results indicate that our model adeptly balances effectiveness and efficiency, also exhibiting faster convergence compared to the pure backbone.

5.3 Ablation Study (RQ2)

To verify the effectiveness of the key components in CorGCN, we conduct ablation studies by comparing CorGCN with its five variants: (1) *w/o CGD* replaces the decomposed graphs by the origin input graph. (2) *w/o CFD* removes the correlation-aware decomposition of node features, replacing them with original features for graph learning. (3) *w/o CSD* excludes the correlation-aware structure decomposition, replacing them with correlated node features and origin structure for CorGCN message passing. (4) *w/o Intra* ignores the intra-label message passing and merges all views of graphs for unified message passing. (5) *w/o Inter* ignore the inter-label correlation propagation without considering message correlations. From Figure 3, we can have following observations:

Effectiveness of the correlation decomposed graph. The severe performance decline of *w/o CGD*, along with the diminished efficacy of *w/o CFD* and *w/o CSD*, demonstrates that the original graph topology and node features are suboptimal to multi-label node classification. The proposed correlation decomposed graph (CDG) with both correlation decomposed node features and topology structure can effectively benefit the classification.

Effectiveness of the correlation decomposed graph convolution. Based on the effective correlation decomposed graph, the lack of carefully designed decomposed graph convolution modules (*w/o intra* and *w/o inter*) both results in inferior classification performance than CorGCN, which demonstrates that modeling intra-label message and passing inter-label correlation is meaningful. Further, removing the intra-label message passing with unified message passing has a greater impact on CorGCN, indicating that the issue of feature and topology ambitiousness indeed brings negative effects during the message passing in GNNs.

**Figure 3: Ablation study on CorGCN with its five variants.****Figure 4: Generalization study results of CorGCN in different GNN message passings.**

5.4 Generalization Study (RQ3)

To investigate the generalization of the proposed CorGCN method, we evaluate the performance of CorGCN with different GNN backbones. Specifically, we utilize the representative GCN-like [20], GAT-like [34], SAGE-like [16], GIN-like [46], and SGC-like [39] for CorGCN in Eq.(10) of message passing, respectively. The study results are illustrated in Figure 4.

From the results, we can find that equipped with CorGCN (*w/ CorGCN*) is consistently significantly better than message passing with vanilla backbones (*w/o CorGCN*), which brings 3.50% and 3.15% average Micro-AUC improvements on Humloc and PCG, respectively. The improvements brought by CorGCN to other backbones are greater compared to GAT. This is because GAT implicitly differentiates the neighborhoods by considering the importance of different nodes. However, the message passing paradigms in these backbones all remain ambitious in multi-label node classification. CorGCN explicitly addresses this issue and thus exhibits better performance. The results demonstrate both the effectiveness and generalization of CorGCN in different message passing functions. The complete backbone generalization results with other metrics can be found in Appendix C.3.

5.5 Case Study (RQ4)

To further analyze the performance of our CorGCN with the graph decomposition, we conduct a case study that analyzes the performance of each fine-grained class. The comparison between the original GCN and the top-performed baseline MLGCN over the Humloc dataset is as follows. From the results in Figure 5, we can have the following findings: (i) Under the decomposition paradigm, CorGCN can achieve significantly better performance among all the fine-grained classes than the pure backbone. (ii) Although MLGCN is the top-performing baseline on Humloc, due to the lack of

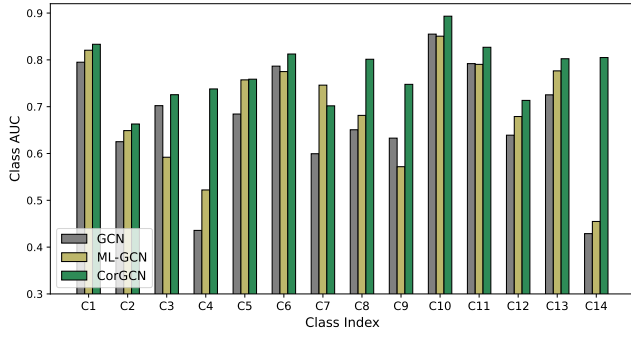


Figure 5: Case study of class AUC performance on Humloc.

Table 4: Efficiency study results on Humloc and PCG datasets each epoch in seconds.

Dataset	Stage	GCN	ML-GCN	LANC	CorGCN
Humloc	Training	0.37	15.79	1.18	0.76
	Inference	0.37	0.49	0.48	0.57
PCG	Training	0.39	18.15	1.21	0.84
	Inference	0.30	0.62	0.39	0.64

ambiguous discrimination ability with a unified convolution process, it will lead to 5/14 classes with performance dropping than the original GCN.

Therefore, the decomposition paradigm of CorGCN is helpful in better understanding each class characteristic over the multi-label graph with ambiguous features and topology.

5.6 Efficiency Study (RQ5)

To further evaluate whether CorGCN can balance the effectiveness and efficiency, we conduct the efficiency study of CorGCN compared with the top-3 performed baselines on both training and inference stage in Table 4.

From the results, we can find that our CorGCN can balance effectiveness and efficiency. For the training stage, the end-to-end label correlation modeling strategy employed by CorGCN yields an efficiency that is second to that of pure GCN. During the inference phase, CorGCN achieves computational efficiency that is slightly lower than top-3 performing baselines due to the multi-label structure decomposition while demonstrating a good enhancement in performance compared to these models.

5.7 Parameter Study (RQ6)

5.7.1 Effect of Parameter λ . We investigate the effect of λ for density controlling in correlation decomposed structure learning with the range from 1 to 13 with a step size of 2 as Figure 6. From the results, we can observe that a too-small value of λ will cause poor performance, which indicates the information is too sparse and insufficient. Furthermore, the suitable value of λ for Humloc is larger than the value for PCG, one possible reason is that the λ of CorGCN should be set larger on more sparse graphs.

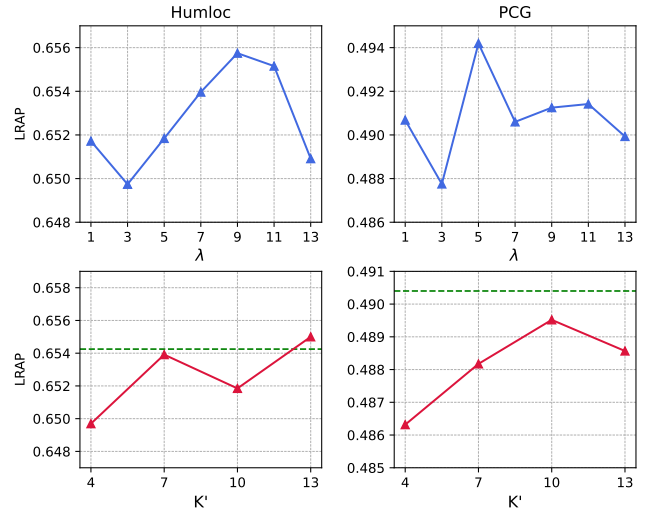


Figure 6: Parameter study results of λ and K' .

5.7.2 Effect of Macro Label Prototype Number K' . We also evaluate the impact of different macro label numbers K' in Sec. 4.4 to CorGCN. From Figure 6, we can find that the too-small cluster number will lead to too coarse label segmentation and result in relatively poor results. In addition, choosing a relatively appropriate number of macro labels, such as 13 for Humloc and 10 for PCG, can achieve model performance close to or even equivalent to the performance of original label numbers (the green dashed line).

More hyperparameter studies on learning rate and γ for CorGCN can be found in C.4.

6 CONCLUSION

In this paper, we propose CorGCN for multi-label node classification. Specifically, aiming to address the ambitious feature and topology in the original graph in multi-label scenarios, we introduce more suitable Correlation-aware Decomposed Graphs (CDG). Then, based on the CDG, we further design a novel Correlation-enhanced Graph Convolution with inter-label message passing and intra-label correlation propagation. Extensive experiments on five benchmark datasets and in-depth analyses in multiple perspectives illustrate the effectiveness and strengths of CorGCN.

REFERENCES

- [1] Junwen Bai, Shufeng Kong, and Carla Gomes. 2021. Disentangled variational autoencoder based multi-label classification with covariance-aware multivariate probit model. In *Proceedings of the Twenty-Ninth International Conference on Artificial Intelligence*. 4313–4321.
- [2] Junwen Bai, Shufeng Kong, and Carla P Gomes. 2022. Gaussian mixture variational autoencoder with contrastive learning for multi-label classification. In *International Conference on Machine Learning*. PMLR, 1383–1398.
- [3] Yuanchen Bei, Sheng Zhou, Qiaoyu Tan, Hao Xu, Hao Chen, Zhao Li, and Jiajun Bu. 2023. Reinforcement neighborhood selection for unsupervised graph anomaly detection. In *2023 IEEE International Conference on Data Mining (ICDM)*. IEEE, 11–20.
- [4] Smriti Bhagat, Graham Cormode, and S Muthukrishnan. 2011. Node classification in social networks. *Social network data analytics* (2011), 115–148.
- [5] Jasmin Bogatinovski, Ljupčo Todorovski, Sašo Žderoski, and Dragi Kocev. 2022. Comprehensive comparative study of multi-label classification methods. *Expert Systems with Applications* 203 (2022), 117215.

- [6] Matthew R Boutell, Jiebo Luo, Xipeng Shen, and Christopher M Brown. 2004. Learning multi-label scene classification. *Pattern recognition* 37, 9 (2004), 1757–1771.
- [7] Hongyun Cai, Vincent W Zheng, and Kevin Chen-Chuan Chang. 2018. A comprehensive survey of graph embedding: Problems, techniques, and applications. *IEEE transactions on knowledge and data engineering* 30, 9 (2018), 1616–1637.
- [8] Hao Chen, Yuanchen Bei, Qijie Shen, Yue Xu, Sheng Zhou, Wenbing Huang, Feiran Huang, Senzhang Wang, and Xiao Huang. 2024. Macro graph neural networks for online billion-scale recommender systems. In *Proceedings of the ACM on Web Conference 2024*. 3598–3608.
- [9] Yu Chen, Lingfei Wu, and Mohammed Zaki. 2020. Iterative deep graph learning for graph neural networks: Better and robust node embeddings. *Advances in neural information processing systems* 33 (2020), 19314–19326.
- [10] Thibaut Durand, Nazanin Mehrasa, and Greg Mori. 2019. Learning a deep convnet for multi-label classification with partial labels. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 647–657.
- [11] Kaisheng Gao, Jing Zhang, and Cangqi Zhou. 2019. Semi-supervised Graph Embedding for Multi-label Graph Node Classification. In *Web Information Systems Engineering – WISE 2019*. Springer International Publishing, Cham, 555–567.
- [12] Ziqi Gao, Chenran Jiang, Jiawen Zhang, Xiaosen Jiang, Lanqing Li, Peilin Zhao, Huanming Yang, Yong Huang, and Jia Li. 2023. Hierarchical graph learning for protein–protein interaction. *Nature Communications* 14, 1 (2023), 1093.
- [13] Weifeng Ge, Sibe Yang, and Yizhou Yu. 2018. Multi-evidence filtering and fusion for multi-label classification, object detection and semantic segmentation based on weakly supervised learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1277–1286.
- [14] Aditya Grover and Jure Leskovec. 2016. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*. 855–864.
- [15] Greg Hamerly and Charles Elkan. 2003. Learning the k in k-means. *Advances in neural information processing systems* 16 (2003).
- [16] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In *Advances in Neural Information Processing Systems*, Vol. 30.
- [17] Wei Jin, Yao Ma, Xiaorui Liu, Xianfeng Tang, Suhang Wang, and Jiliang Tang. 2020. Graph structure learning for robust graph neural networks. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*. 66–74.
- [18] Xiaoyang Jing and Jinbo Xu. 2021. Fast and effective protein model refinement using deep graph neural networks. *Nature computational science* 1, 7 (2021), 462–469.
- [19] Samuel Kerrien, Bruno Aranda, Lionel Breuza, Alan Bridge, Fiona Broackes-Carter, Carol Chen, Margaret Duesbury, Marine Dumousseau, Marc Feuermann, Ursula Hinz, et al. 2012. The IntAct molecular interaction database in 2012. *Nucleic acids research* 40, D1 (2012), D841–D846.
- [20] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *International Conference on Learning Representations*.
- [21] Zhixun Li, Xin Sun, Yifan Luo, Yanqiao Zhu, Dingshuo Chen, Yingtao Luo, Xiangxin Zhou, Qiang Liu, Shu Wu, Liang Wang, et al. 2024. GSLB: the graph structure learning benchmark. *Advances in Neural Information Processing Systems* 36 (2024).
- [22] Arthur Liberzon, Chet Birger, Helga Thorvaldsdóttir, Mahmoud Ghandi, Jill P Mesirov, and Pablo Tamayo. 2015. The molecular signatures database hallmark gene set collection. *Cell systems* 1, 6 (2015), 417–425.
- [23] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. 2017. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*. 2980–2988.
- [24] Jingzhou Liu, Wei-Cheng Chang, Yuexin Wu, and Yiming Yang. 2017. Deep learning for extreme multi-label text classification. In *Proceedings of the 40th international ACM SIGIR conference on research and development in information retrieval*. 115–124.
- [25] Weiwei Liu, Haobo Wang, Xiaobo Shen, and Ivor W Tsang. 2021. The emerging trends of multi-label learning. *IEEE transactions on pattern analysis and machine intelligence* 44, 11 (2021), 7955–7974.
- [26] Yixin Liu, Yu Zheng, Daokun Zhang, Hongxu Chen, Hao Peng, and Shirui Pan. 2022. Towards unsupervised deep graph structure learning. In *Proceedings of the ACM Web Conference 2022*. 1392–1403.
- [27] Zhibin Pan, Yidi Wang, and Weiping Ku. 2017. A new general nearest neighbor classification based on the mutual neighborhood information. *Knowledge-Based Systems* 121 (2017), 142–152.
- [28] Janet Piñero, Juan Manuel Ramírez-Anguita, Josep Saüch-Pitarch, Francesco Ronzano, Emilio Centeno, Ferran Sanz, and Laura I Furlong. 2020. The DisGeNET knowledge platform for disease genomics: 2019 update. *Nucleic acids research* 48, D1 (2020), D845–D855.
- [29] Jesse Read, Bernhard Pfahringer, Geoff Holmes, and Eibe Frank. 2009. Classifier chains for multi-label classification. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2009, Bled, Slovenia, September 7–11, 2009, Proceedings, Part II 20*. Springer, 254–269.
- [30] Yu Rong, Yatao Bian, Tingyang Xu, Weiyang Xie, Ying Wei, Wenbing Huang, and Junzhou Huang. 2020. Self-supervised graph transformer on large-scale molecular data. *Advances in Neural Information Processing Systems* 33 (2020), 12559–12571.
- [31] Zixing Song, Ziqiao Meng, Yifei Zhang, and Irwin King. 2021. Semi-supervised multi-label learning for graph-structured data. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 1723–1733.
- [32] Peijie Sun, Le Wu, and Meng Wang. 2018. Attentive recurrent social recommendation. In *The 41st international ACM SIGIR conference on research & development in information retrieval*. 185–194.
- [33] Adane Nega Tarekn, Mohib Ullah, and Faouzi Alaya Cheikh. 2024. Deep learning for multi-label learning: A comprehensive survey. *arXiv preprint arXiv:2401.16549* (2024).
- [34] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2018. Graph attention networks. In *The International Conference on Learning Representations*.
- [35] Hongwei Wang and Jure Leskovec. 2021. Combining graph convolutional neural networks and label propagation. *ACM Transactions on Information Systems (TOIS)* 40, 4 (2021), 1–27.
- [36] Jiang Wang, Yi Yang, Junhua Mao, Zhiheng Huang, Chang Huang, and Wei Xu. 2016. Cnn-rnn: A unified framework for multi-label image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2285–2294.
- [37] Kaixiang Wang, Ming Yang, Wanqi Yang, and YiLong Yin. 2018. Deep correlation structure preserved label space embedding for multi-label classification. In *Asian Conference on Machine Learning*. PMLR, 1–16.
- [38] Ya Wang, Dongliang He, Fu Li, Xiang Long, Zhichao Zhou, Jinwen Ma, and Shilei Wen. 2020. Multi-label classification with label graph superimposing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 12265–12272.
- [39] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. 2019. Simplifying graph convolutional networks. In *International conference on machine learning*. PMLR, 6861–6871.
- [40] Jian Wu, Victor S Sheng, Jing Zhang, Hua Li, Tetiana Dadakova, Christine Leon Swisher, Zhiming Cui, and Pengpeng Zhao. 2020. Multi-label active learning algorithms for image classification: Overview and future promise. *ACM Computing Surveys (CSUR)* 53, 2 (2020), 1–35.
- [41] Shiwen Wu, Fei Sun, Wentao Zhang, Xu Xie, and Bin Cui. 2022. Graph neural networks in recommender systems: a survey. *Comput. Surveys* 55, 5 (2022), 1–37.
- [42] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. 2020. A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems* 32, 1 (2020), 4–24.
- [43] Lin Xiao, Xin Huang, Boli Chen, and Liping Jing. 2019. Label-specific document representation for multi-label text classification. In *Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP)*. 466–475.
- [44] Lin Xiao, Pengyu Xu, Liping Jing, Uchenna Akujobi, and Xiangliang Zhang. 2022. Semantic guide for semi-supervised few-shot multi-label node classification. *Information Sciences* 591 (2022), 235–250.
- [45] Shunxin Xiao, Shipping Wang, Yuanfei Dai, and Wenzhong Guo. 2022. Graph neural networks in node classification: survey and evaluation. *Machine Vision and Applications* 33 (2022), 1–19.
- [46] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2019. How Powerful are Graph Neural Networks?. In *International Conference on Learning Representations*.
- [47] Vacit Oгуз Yazici, Abel Gonzalez-Garcia, Arnau Ramisa, Bartłomiej Twardowski, and Joost van de Weijer. 2020. Orderless recurrent models for multi-label classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13440–13449.
- [48] Chih-Kuan Yeh, Wei-Chieh Wu, Wei-Jen Ko, and Yu-Chiang Frank Wang. 2017. Learning deep latent space for multi-label classification. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 31.
- [49] Donghan Yu, Ruohong Zhang, Zhengbao Jiang, Yuexin Wu, and Yiming Yang. 2021. Graph-revised convolutional network. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2020, Ghent, Belgium, September 14–18, 2020, Proceedings, Part III*. Springer, 378–393.
- [50] Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor Prasanna. 2019. GraphSAINT: Graph Sampling Based Inductive Learning Method. In *International Conference on Learning Representations*.
- [51] Jiawei Zhang and Philip S Yu. 2018. Broad learning: An emerging area in social network analysis. *ACM SIGKDD Explorations Newsletter* 20, 1 (2018), 24–50.
- [52] Min-Ling Zhang and Zhi-Hua Zhou. 2013. A review on multi-label learning algorithms. *IEEE transactions on knowledge and data engineering* 26, 8 (2013), 1819–1837.
- [53] Tianqi Zhao, Thi Ngan Dong, Alan Hanjalic, and Megha Khosla. 2023. Multi-label Node Classification On Graph-Structured Data. *Transactions on Machine Learning Research* (2023).

- [54] Wenting Zhao, Shufeng Kong, Junwen Bai, Daniel Fink, and Carla Gomes. 2021. Hot-vae: Learning high-order label correlation for multi-label classification via attention-based variational autoencoders. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 15016–15024.
- [55] Cangqi Zhou, Hui Chen, Jing Zhang, Qianmu Li, Dianming Hu, and Victor S Sheng. 2021. Multi-label graph node classification with label attentive neighborhood convolution. *Expert Systems with Applications* 180 (2021), 115063.

A METHODOLOGY DETAILS

A.1 Analysis of the Objective Function on your Multi-Label Estimator

The mutual information estimator ϕ and likelihood maximizing decoder θ work collaboratively for node-label correlation and label-label correlation modeling as

$$\arg \max_{\phi, \theta} \underbrace{I_{\phi}(E^x, E^l)}_{\text{Mutual Information}} + \underbrace{\mathcal{L}_{\theta}(\theta; E^x) + \mathcal{L}_{\theta}(\theta; E^l)}_{\text{Likelihood}}. \quad (24)$$

Analysis of the effectiveness of mutual information estimator ϕ for multi-label correlation modeling can be found in previous works [2, 54]. Then, The likelihood maximizing decoder θ can be utilized as a likelihood estimator of node feature and label prototype with the entropy KL divergence with focal enhancement [23] for target labels as follows:

$$\mathcal{L}_{lm} = \frac{1}{2^{|\mathcal{V}_L|}} \sum_{v_i \in \mathcal{V}_L} \sum_{k=1}^K \rho_k \left[\underbrace{(1 - p_k(E_i^x | \theta))^y \log p_k(E_i^x | \theta)}_{\text{Likelihood Estimation with Node Feature}} + \underbrace{(1 - p_k(\mathbf{y}_i \cdot E^l | \theta))^y \log p_k(\mathbf{y}_i \cdot E^l | \theta)}_{\text{Likelihood Estimation with Label Prototype}} \right]. \quad (25)$$

A.2 Complexity Analysis

The overall complexity of CorGCN arises from the following stages. During the correlation-aware graph decomposition, the calculation involves the projection process, incurring a computational cost of $\Theta(K \cdot n^2)$. For the Correlation-Enhanced graph convolution, we postulate that the message passing process within a graph comprising n nodes entails a cost of $\Theta_{mp}(n)$, which varies according to the specific GCN backbone employed. During the intra-label message passing, due to the process can be paralleled across K label views, the overall expense of message passing is quantified as $\Theta_{mp}(n)$. Then, the cost associated with inter-label correlation propagation is $\Theta(K^2 \cdot n)$. Thus, the overall time complexity of the process can be expressed as $T(n) = \Theta(K(n^2 + \Theta_{mp}(n) + K \cdot n))$.

As outlined in Sec. 4.4, this complexity can be reduced to $T(n) = \Theta(K'(n^2 + \Theta_{mp}(n) + K'n))$ for the large label space, where $K' \ll K$, denoting a significant reduction in computational cost. In practice, the entire process can be batch-processed, where the similarity is calculated among the nodes within each batch, thus making the time complexity of each batch B to $T(B) = \Theta(K'(B^2 + \Theta_{mp}(B) + K'B))$.

A.3 Pseudocode of CorGCN

The overall pipeline of CorGCN includes (1) correlation-aware graph learning with feature and topology decompositions, and (2) correlation-aware graph convolution with intra-label message

Algorithm 1 The overall workflow of CorGCN.

Input: Feature matrix X ; Label prototype E^l ; Graph topology A^0 ; Number of GNN layers N_{gnn} .

Output: Predicted labels $\hat{\mathbf{y}}$.

- 1: $E^x \leftarrow X$; // Compute transformed features
 - 2: $\mathcal{L}_{cmi} \leftarrow \{E^x, E^l\}$ using Eq.(4); // Compute contrastive learning loss
 - 3: $\mathcal{L}_{lm} \leftarrow \{E^x, E^l, \mathbf{y}\}$ using Eq.(5) and Eq.(6); // Compute likelihood loss
 - 4: $E^{proj} \leftarrow \{E^x, E^l\}$ using Eq.(7), (8); // Compute projected features as node representations
 - 5: $E^{sd} \leftarrow \{E^{proj}, A\}$ using Eq.(9); // Aggregate the node representations from neighborhood
 - 6: $CDG \leftarrow \{E^{sd}, E^{proj}, A\}$ using Eq.(10) - Eq.(13); // Learn the graphs by aggregated representations
 - 7: **for** $epoch \in 1, 2, \dots, N_{gnn}$ **do**
 - 8: $\hat{Z}^{(l)} \leftarrow \{Z^{(l-1)}, CDG\}$ using Eq.(14) // Update the node representations by message passing
 - 9: $Z^{(l)} \leftarrow \{\hat{Z}^{(l)}, E^l\}$ using Eq.(15) - (16) // Update the node representations by correlation propagation
 - 10: **end for**
 - 11: $\hat{\mathbf{y}} \leftarrow \{Z^{(E)}, E^l\}$ using Eq.(17), (18) // Predict the probabilities
 - 12: $\mathcal{L}_{cls} \leftarrow \{\hat{\mathbf{y}}, \mathbf{y}\}$ using Eq.(19); // Compute the classification loss
 - 13: $\mathcal{L} \leftarrow \mathcal{L}_{cls} + \alpha \mathcal{L}_{cmi} + \beta \mathcal{L}_{lm}$;
 - 14: **if** *Training* **then**
 - 15: Back-propagate \mathcal{L} to update model parameters.
 - 16: **end if**
-

passing and inter-label correlation propagation. The general process of this pipeline is described as the pseudocode in Algorithm 1.

B EXPERIMENTAL DETAILS

B.1 Dataset Details

We adopt five publicly available datasets for comprehensive evaluation of our proposed CorGCN. Among these datasets, the included PPI dataset is with large label space and the Delve dataset is a large-scale dataset in multi-label node classification. The detailed description is as follows:

Humloc Dataset¹ [53] is a human protein subcellular location prediction dataset, consisting of 3,106 nodes and 18,496 edges. Each node may have one or more labels in 14 possible locations. The edge data is derived from protein-protein interactions obtained from the IntAct database [19]. An edge exists between two nodes in the graph if there is an interaction between the respective proteins in IntAct.

PCG Dataset¹ [53] is a protein phenotype prediction dataset with 3,233 nodes into 15 classes and 37,351 edges between them. Each node represents a protein with 32-dimensional features. Each edge between a pair of protein nodes is their functional interaction. The multi-label that each node is associated with is its correlated phenotypes, of which a phenotype is any observable characteristic or trait of a disease. The correspondence between protein and phenotype is retrieved from the DisGeNET database [28].

¹<https://github.com/Tianqi-py/MLGNC>

Blogcatalog Dataset² [55] is a social network dataset with 10,312 nodes and 333,983 edges, which each node represents a blogger at the platform and each edge represents the contact relationship between a pair of nodes. The labels denote the social interest groups a blogger node is a part of, forming a total of 39 interest groups.

PPI Dataset³ [50] serves as a benchmark dataset with a large label space for protein-protein interaction networks, comprising 14,755 nodes (representing proteins) and 225,270 edges (indicating functional interactions). The dataset is annotated with gene ontology categories as labels, encompassing a total of 121 unique sets, sourced from the Molecular Signatures Database [22].

Delve Dataset⁴ [44] is a large-scale citation network dataset with 1,229,280 nodes and 4,322,275 edges, which each node is a paper and each edge represents a citation relationship between a pair of papers. The labels denote the research fields that papers belong to, forming a total of 20 categories.

B.2 Baseline Details

We compare our proposed CorGCN with nine representative state-of-the-art GNN models as follows.

- **GCN** [20] is a widely recognized model in the domain of graph neural networks, grounded in spectral theory.
- **GRCN** [49] incorporates a module based on GCN specifically engineered for the prediction of missing edges and the adjustment of edge weights, which is optimized with downstream tasks.
- **IDGL** [9] mutually enhances the graph structure and node embeddings: the premise is to iteratively improve the graph structure through enhanced node embeddings and, reciprocally, to refine node embeddings leveraging an optimized graph structure.
- **SUBLIME** [26] uses a contrastive loss to maximize the agreement between the anchor graph and the learned graph, which is further enhanced with a bootstrapping mechanism.
- **GCN-LPA** [35] unifies GCNs and label propagation in the same learning framework, which we utilized to propagate the multi-label simultaneously.
- **ML-GCN** [11] employs a strategy to generate and utilize a label matrix for capturing node-label and label-label correlations through a relaxed skip-gram model in a unified vector space.
- **LANC** [55] develops a label attention module that employs an additive attention mechanism to synergize input and output contextual representations.
- **SMLG** [31] comprises a joint variational representation module that generates node and label embeddings and a confidence-rated margin ranking module that detects second-order label correlations and refines the decision boundaries.
- **LARN** [44] designs a novel label correlation scanner that adaptively captures label relationships, extracting vital encoded information for the generation of comprehensive representations.

²https://figshare.com/articles/dataset/BlogCatalog_dataset/11923611

³<https://github.com/GraphSAINT/GraphSAINT>

⁴<https://www.dropbox.com/sh/tg8nclx5gpcbtbmo/AADQc0YFpuVeqXhaNuUK4MMba?dl=0>

B.3 Implementation Details

We explore the combination space of the following hyper-parameters lr , λ , K' , γ and list the values finally chosen in Table 5 as the setting of performing the main results. Specifically, the learning rate lr is selected from $\{0.0001, 0.0005, 0.001, 0.005, 0.01, 0.02, 0.05\}$; λ that controls the density of graphs is tuned from 1 to 20; the macro label prototype number K' is searched from 1 to 20; and γ that controls the focusing factor in \mathcal{L}_{lm} is chosen from $\{1.0, 1.5, 2.0, 2.5, 3.0\}$. Besides, to compare fairly, the hidden dim d is set to 64, and the dropout rate of the network is set to 0.3.

Table 5: The values of hyper-parameters used in CorGCN.

Dataset	lr	λ	K'	γ
Humloc	0.001	7	\	2.0
PCG	0.001	19	\	2.0
Blogcatalog	0.001	5	\	2.0
PPI	0.02	5	20	2.0
Delve	0.001	5	10	2.0

C ADDITIONAL EXPERIMENTAL RESULTS

C.1 Bonferroni-Dunn Test Results

The Bonferroni-Dunn test is a post-hoc procedure used to control the family-wise error rate when multiple comparisons are being made, such as in the case of comparing multiple models or treatments against a control in statistical analysis [27]. This formula is used after performing a Friedman statistic, in which we have obtained the statistics all show significance with p -value less than 0.05 level in Table 3.

The critical difference (CD) for the Bonferroni-Dunn test, which is used to determine if the difference between two means is statistically significant, can be expressed as:

$$CD = q_{\alpha} \sqrt{\frac{k(k+1)}{6N}}, \quad (26)$$

Where q_{α} is the critical value based on the Studentized range statistic q , for α significance level (we set 0.05) adjusted for the number of comparisons. k is the number of groups or treatments being compared (we use five datasets). N is the total number of observations across all groups ($N = 5 \times 10 = 50$ in our experiments).

The visualization results of the Bonferroni-Dunn on the evaluation metrics are presented in Figure 7 as follows all demonstrate that the performance advantage of our model is significant.

C.2 Convergence Analysis

To further analyze the convergence of our CorGCN, we provide the figures of training classification loss changing over epochs and validation metrics changing over epochs as Figure 8.

From the results, we can find that CorGCN can achieve a faster and better convergence than the pure backbone for multi-label node classification with lower training loss and higher validation performance. This is due to the design of CorGCN with the correlation-aware multi-label graph decompositon and convolution can help the backbone better and faster understand the multi-label relationships and important information during message passing.

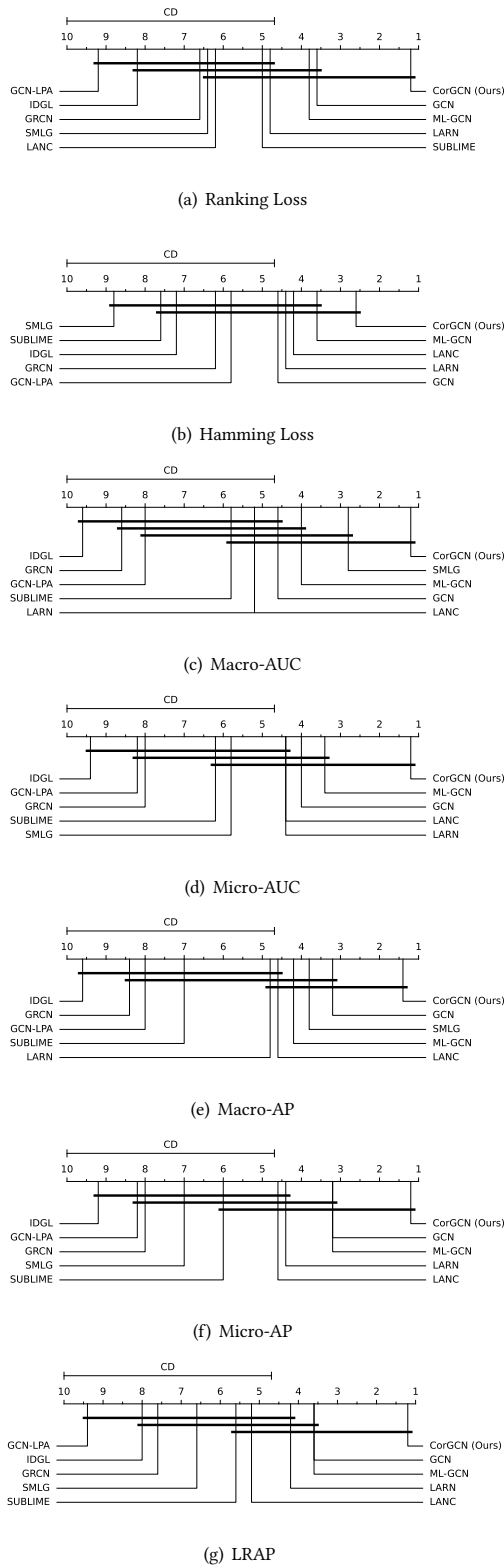


Figure 7: Bonferroni-Dunn test results of our proposed model against other baselines on all the evaluation metrics, where $CD = 5.31$ at 0.05 significance level.

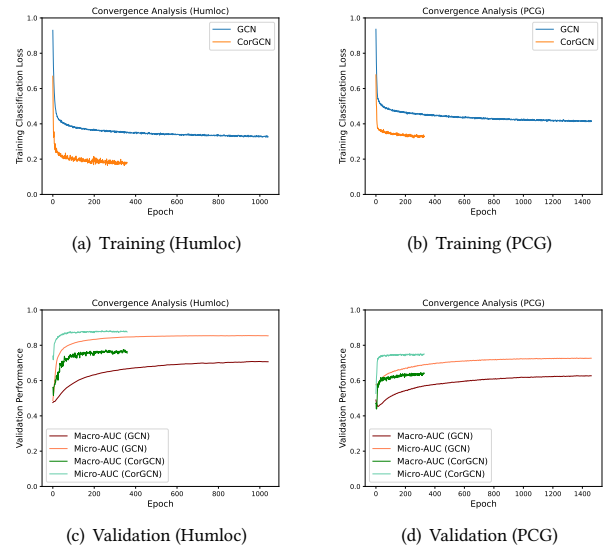


Figure 8: (a)-(b) Convergence analysis on the training classification loss. (c)-(d) Convergence analysis on the validation set performance.

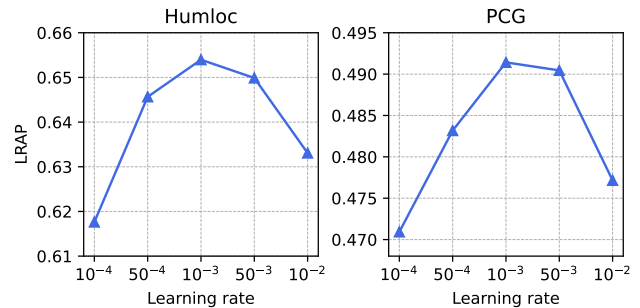


Figure 9: Parameter study results on learning rate of CorGCN.

C.3 More Backbone Generalization Study Results

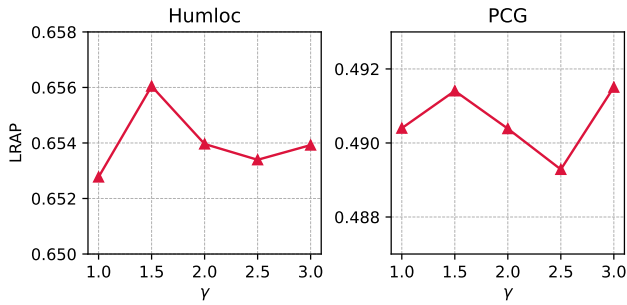
In Table 6, we have shown the complete backbone generalization study results. From the overall results, we can observe that CorGCN can generally positively improve the backbone performance on multi-label node classification, which illustrates that CorGCN can serve as a flexible plug-in framework to enhance the performance of traditional graph convolution backbone models.

C.4 More Parameter Study Results

C.4.1 Effect of Learning Rate. We tested the impact of the learning rate on the performance of CorGCN. As shown in Figure 9, it can be observed that when the learning rate is too low, the model’s performance is poor. It’s possibly the small step size of gradient descent that makes it difficult for the model to escape local optima. On the other hand, when the learning rate is too high, the model’s

Table 6: Backbone study results on Humloc and PCG datasets of all seven evaluation metrics (\uparrow : the higher, the better; \downarrow : the lower, the better).

Dataset	Metrics	GCN	CorGCN (GCN)	GAT	CorGCN (GAT)	SAGE	CorGCN (SAGE)	GIN	CorGCN (GIN)	SGC	CorGCN (SGC)
Humloc	Ranking (\downarrow)	13.29 \pm 1.01	12.57 \pm 0.31	13.61 \pm 0.47	13.27 \pm 0.31	14.60 \pm 0.46	12.77 \pm 0.27	13.59 \pm 0.22	12.17 \pm 0.51	14.92 \pm 0.35	12.87 \pm 0.30
	Hamming (\downarrow)	7.52 \pm 0.13	7.37 \pm 0.17	7.67 \pm 0.19	7.59 \pm 0.13	8.08 \pm 0.18	7.73 \pm 0.12	7.80 \pm 0.05	7.10 \pm 0.16	8.19 \pm 0.08	7.31 \pm 0.12
	Ma-AUC (\uparrow)	69.10 \pm 2.25	77.31 \pm 1.58	68.76 \pm 1.53	73.51 \pm 2.19	65.83 \pm 2.97	73.59 \pm 2.02	68.59 \pm 2.55	79.88 \pm 2.16	68.56 \pm 1.56	77.17 \pm 1.12
	Mi-AUC (\uparrow)	85.39 \pm 1.30	88.57 \pm 0.37	85.73 \pm 0.37	87.57 \pm 0.31	84.29 \pm 0.48	87.66 \pm 0.15	85.49 \pm 0.31	89.30 \pm 0.79	85.52 \pm 0.28	88.23 \pm 0.34
	Ma-AP (\uparrow)	24.65 \pm 1.29	26.91 \pm 1.67	24.63 \pm 1.25	25.04 \pm 1.69	20.21 \pm 1.34	24.47 \pm 0.85	23.50 \pm 2.13	31.76 \pm 2.77	19.15 \pm 0.34	27.62 \pm 1.83
	Mi-AP (\uparrow)	46.46 \pm 2.06	48.97 \pm 1.30	44.55 \pm 2.35	46.00 \pm 1.48	38.68 \pm 0.84	44.74 \pm 1.84	41.10 \pm 0.38	52.70 \pm 1.47	36.20 \pm 0.67	48.35 \pm 1.19
	LRAP (\uparrow)	64.66 \pm 0.80	65.40 \pm 0.83	64.38 \pm 1.23	64.16 \pm 1.11	62.12 \pm 1.11	64.61 \pm 1.02	63.39 \pm 0.78	66.56 \pm 0.66	59.74 \pm 0.35	65.21 \pm 0.32
PCG	Ranking (\downarrow)	27.50 \pm 0.58	25.97 \pm 0.57	27.60 \pm 0.37	27.34 \pm 0.65	29.84 \pm 0.51	28.10 \pm 0.41	28.64 \pm 0.72	27.09 \pm 0.76	27.54 \pm 0.10	26.39 \pm 0.33
	Hamming (\downarrow)	13.22 \pm 0.30	12.41 \pm 0.23	12.62 \pm 0.33	12.47 \pm 0.17	12.95 \pm 0.14	12.61 \pm 0.16	12.64 \pm 0.11	12.37 \pm 0.08	12.75 \pm 0.05	12.47 \pm 0.21
	Ma-AUC (\uparrow)	62.92 \pm 0.84	64.86 \pm 1.05	62.13 \pm 1.18	62.33 \pm 1.03	55.17 \pm 0.80	57.95 \pm 0.55	58.06 \pm 1.32	62.80 \pm 0.67	61.20 \pm 0.43	64.15 \pm 1.46
	Mi-AUC (\uparrow)	71.74 \pm 0.38	74.16 \pm 0.61	72.11 \pm 0.41	72.57 \pm 0.61	68.21 \pm 0.61	70.65 \pm 0.29	69.54 \pm 1.08	73.10 \pm 0.65	71.73 \pm 0.24	73.91 \pm 0.38
	Ma-AP (\uparrow)	23.49 \pm 0.38	24.64 \pm 0.90	21.33 \pm 0.74	21.87 \pm 1.14	15.99 \pm 0.26	17.83 \pm 0.78	17.37 \pm 0.81	23.21 \pm 0.60	19.22 \pm 0.49	24.25 \pm 1.07
	Mi-AP (\uparrow)	30.04 \pm 0.54	31.91 \pm 1.07	29.19 \pm 0.94	29.91 \pm 1.66	24.04 \pm 0.63	25.98 \pm 0.45	25.94 \pm 0.93	31.13 \pm 1.00	26.68 \pm 0.38	31.73 \pm 0.51
	LRAP (\uparrow)	48.03 \pm 1.25	49.04 \pm 0.83	48.02 \pm 1.06	47.53 \pm 1.18	45.23 \pm 0.83	46.40 \pm 0.73	45.94 \pm 1.52	47.70 \pm 0.96	48.72 \pm 0.21	48.82 \pm 0.90

**Figure 10: Parameter study results on γ of CorGCN.**

parameters may miss the minimum points during gradient descent, preventing convergence to the optimal solution.

C.4.2 Effect of Parameter γ . We explored the impact of the γ on the performance of CorGCN. Parameter γ amplifies the loss weight of misclassified samples, thereby increasing the model’s attention to them. As shown in Figure 10, it can be observed that the model achieves the best performance on the LARP metric when γ is set to 1.5. However, as γ increases, it may cause the model to overly focus on misclassified samples, neglecting the well-classified samples, resulting in a slight decrease in model performance.

C.4.3 Adaptive Parameters α and β . To balance the proportions between \mathcal{L}_{cls} , \mathcal{L}_{cmi} and \mathcal{L}_{le} , we calculate the loss balancing parameters $\alpha = \left| \frac{\mathcal{L}_{cls}}{3\mathcal{L}_{cmi}} \right|$ and $\beta = \left| \frac{\mathcal{L}_{cls}}{3\mathcal{L}_{le}} \right|$ without gradients. With these adaptive parameters, We ensure that the classification loss \mathcal{L}_{cls} remains the primary loss while avoiding excessive differences among contrastive loss \mathcal{L}_{cmi} and likelihood loss \mathcal{L}_{le} , which will neither dominate excessively nor become too small to effectively constrain the model parameters.