








# Towards Long-Tailed Recognition for Graph Classification via Collaborative Experts

Si-Yu Yi , Zhengyang Mao , Wei Ju , *Member, IEEE*, Yong-Dao Zhou , Luchen Liu , Xiao Luo ,  
and Ming Zhang 

**Abstract**—Graph classification, aiming at learning the graph-level representations for effective class assignments, has received outstanding achievements, which heavily relies on high-quality datasets that have balanced class distribution. In fact, most real-world graph data naturally presents a long-tailed form, where the head classes occupy much more samples than the tail classes, it thus is essential to study the graph-level classification over long-tailed data while still remaining largely unexplored. However, most existing long-tailed learning methods in visions fail to jointly optimize the representation learning and classifier training, as well as neglect the mining of the hard-to-classify classes. Directly applying existing methods to graphs may lead to sub-optimal performance, since the model trained on graphs would be more sensitive to the long-tailed distribution due to the complex topological characteristics. Hence, in this paper, we propose a novel long-tailed graph-level classification framework via Collaborative Multi-expert Learning (CoMe) to tackle the problem. To equilibrate the contributions of head and tail classes, we first develop balanced contrastive learning from the view of representation learning, and then design an individual-expert classifier training based on hard class mining. In addition, we execute gated fusion and disentangled knowledge distillation among the multiple experts to promote the collaboration in a multi-expert framework. Comprehensive experiments are performed on seven widely-used benchmark datasets to demonstrate the superiority of our method CoMe over state-of-the-art baselines.

**Index Terms**—Balanced contrastive learning, class-imbalanced learning, hard class extraction, multi-expert learning.

## I. INTRODUCTION

GRAPH-STRUCTURED data [1], [2], [3] is ubiquitous in a variety of domains, such as social networks, protein-protein interaction networks, and citation networks. Graph classification [4], [5], [6], [7], as one of the most fundamental tasks

in data mining for graphs, has attracted significant attention. It attempts to predict the class label and the property of each graph in a dataset. Promising applications include property prediction for quantum mechanics [8] and functional assessment of chemical compounds [9].

Graph neural networks (GNNs) [1], [2], [10], [11] have become one of the most prominent approaches in graph representation learning and achieved remarkable progress for graph classification. The key idea of GNNs is to iteratively propagate and aggregate the information from the neighbors of each node in a graph for generating node-level representations [12]. Afterward, a readout function [13], [14] integrates all of the node representations into a graph-level representation, which is then fed into a classifier to predict the graph label. Hence, the learned graph-level representation can incorporate the characteristics of the topological structure and reveal the whole semantic information of the graph, which works well for the downstream graph classification task. In spite of this, the huge success of GNNs is typically built on high-quality datasets of having a roughly balanced distribution. Nevertheless, real-world datasets commonly exhibit long-tailed class distributions, where a large portion of tail classes occupy a limited number of data whereas few head classes have most of the data. For example, in the Cora citation network, the proportion of the instances in head class *Neural Network* is about 26.8%, while those in the tail classes *Rule Learning* and *Reinforcement Learning* are only about 7.9% and 4.8% respectively. In such scenarios, directly adopting GNNs on long-tailed graph datasets may lead to notorious prediction bias and significant performance degradation, since the model is prone to being dominated by the head class while attention to tail classes is easily overlooked. Consequently, the high-frequency head classes may achieve impressive predictive performance, while the low-resource tail classes receive unsatisfactory accuracy, thereby hindering the strength of GNNs learned from long-tailed graph data.

To tackle the problem caused by long-tailed data distribution, existing studies have proposed lots of methods [15], [16], [17], [18], [19], [20] in computer vision, which fall into three categories: re-sampling, re-weighting, and ensemble learning. The re-sampling strategy [15], [21] aims to balance the data distribution in the head and tail classes, including over- and under-sampling. Over-sampling replicates existing samples in the tail classes, whereas under-sampling discards some samples from the head classes to reduce the imbalance. However, re-sampling can lead to over-fitting or under-fitting problems,

Manuscript received 10 March 2023; revised 10 July 2023; accepted 29 August 2023. Date of publication 7 September 2023; date of current version 13 November 2023. This work is partially supported by the National Natural Science Foundation of China (NSFC) under Grants 62306014, 62106008, 62276002, and 12131001, in part by the China Postdoctoral Science Foundation with under Grant 2023M730057, and in part by the Fundamental Research Funds for the Central Universities, LPMC, and KLMDASR. Recommended for acceptance by H. Kitagawa. (Si-Yu Yi and Zhengyang Mao contribute equally to this paper.) (Corresponding authors: Wei Ju; Ming Zhang.)

Si-Yu Yi and Yong-Dao Zhou are with the School of Statistics and Data Science, Nankai University, Tianjin 300071, China (e-mail: siyuyi@mail.nankai.edu.cn; ydzhou@nankai.edu.cn).

Zhengyang Mao, Wei Ju, Luchen Liu, and Ming Zhang are with the School of Computer Science, Peking University, Beijing 100871, China (e-mail: mao.zhengyang.cn@gmail.com; juwei@pku.edu.cn; liuluchen@pku.edu.cn; mzhang\_cs@pku.edu.cn).

Xiao Luo is with the Department of Computer Science, University of California, Los Angeles, CA 90095 USA (e-mail: xiaoluo@cs.ucla.edu).

Digital Object Identifier 10.1109/TBDDATA.2023.3313029

since improper re-sampling may excessively utilize the minority samples or abandon useful information in the majority samples. The re-weighting strategy [22], [23] modifies the weights of the losses from different classes to make the long-tailed data contribute properly to training. But plain re-weighting strategy could benefit classifier learning while hurting representation learning, because it may under-represent the head classes and cause unstable training [24], [25]. Different from the formers, ensemble learning [18], [26] combines multiple expert networks from a complementary perspective to obtain reliable and robust predictions, whose works have achieved satisfactory progress. Nevertheless, most current ensemble learning methods lack mutual supervision among different experts and knowledge transfer is also deficient.

Despite the wealth of the researches in visions, however, to the best of our knowledge, the long-tailed graph-level classification is yet rarely explored. Thus, it is essential to pay insight into the task for meeting the demand of the practical applications. When dealing with long-tailed graphs, because of the complex topology structure and the ubiquitous over-smoothing issue caused by the message passing in GNN-based encoder [27], the classification model trained on graphs would be more sensitive to the long-tailed class distribution and it thus is necessary to improve all the components of the model that can be affected by the long-tailed distribution, such as the representation learning and classifier training. Nevertheless, due to the shortcomings of the existing methods in visions discussed above, directly applying them to the long-tailed graph data cannot solve such problems roundly and effectively, which may result in sub-optimal performance. Moreover, many existing methods focus on the training of the hard samples by re-sampling or re-weighting while ensemble learning lays emphasis on aggregating different perspectives by multiple experts for more comprehensive data exploration. These strategies overlook the selection of the hard-to-classify classes, which leads to many samples being indistinguishable from the non-target classes (i.e., hard classes) but still having high confidence. Instead, the selection of the hard classes for all samples is actually a more refined and comprehensive way to distinguish samples that are difficult to classify accurately and acquire higher classification accuracy. Hence, it is highly desirable to proposed a novel method to jointly optimize representation learning and classifier learning as well as effectively capture hard classes for long-tailed graphs.

To address these issues, based on **Collaborative Multi-expert Learning**, we propose a novel framework named CoMe for long-tailed graph classification. The key idea of CoMe is to propose tailored balanced contrastive learning along with individual-expert classifier training to jointly optimize the representation learning and classifier learning, and then fuse and distill the multiple expert networks from both global and local views for stronger collaboration capability. Specifically, CoMe first introduces tailored balanced contrastive learning to alleviate the class imbalance for representation learning, which can effectively balance the contributions of the head and tail classes on the contrastive loss. Then, we propose the balanced predicted

probability in the classifier learning for each expert from both global and local perspectives to alleviate the influence of the sample sizes in different classes on embedded representations and enhance the hard class mining. Moreover, to fully benefit from multi-expert/ensemble learning, we fuse different expert models by gating functions to increase the diversity of the whole training network and meanwhile, we perform knowledge distillation among experts in a disentangled manner to encourage them learn extra knowledge from others and decouple the effects of the predictions for the target class and non-target classes. Comprehensive experiments are conducted to show that the proposed method can greatly improve the performance of the long-tailed graph classification compared with existing state-of-the-art methods over multiple benchmark datasets. Moreover, the combination of tailored representation learning and classifier training is highly effective for dealing with imbalanced settings. To summarize, the main contributions of our work are as follows:

- This paper studies long-tailed graph-level classification under the multi-expert learning framework, which jointly optimizes the representation learning and classifier training, and explores the organic fusion and effective cooperation among expert networks.
- We propose a framework to balance the contributions of the head and tail classes on both representation learning and classifier learning. We explicitly capture the hard-to-classify classes for all samples and equip the multi-expert learning with gated fusion and disentangled knowledge distillation to enhance the long-tailed learning.
- Extensive experiments are performed on various commonly used datasets to validate the superiority of the proposed approach against existing state-of-the-art models.

## II. RELATED WORK

In this section, we briefly review the related works in three aspects, namely graph-level classification, long-tailed learning, and contrastive learning.

### A. Graph-Level Classification

Graph-level classification is one of the most critical problems in the graph domain, which aims at predicting the class label of the entire graph. Existing algorithms for graph classification can be broadly categorized into graph kernel-based methods and GNN-based methods. The core of classic graph kernels is to decompose each graph into substructures (e.g., graphlets [28], subtrees [29], or shortest paths [30]) to measure the similarity between two graphs. Recent works have focused on designing expressive GNNs [1], [10], [11], [31] and achieved remarkable success. The key idea of GNNs [1], [2], [10], [32], [33] is to iteratively update the node feature according to its neighbor nodes with pooling methods [13], [14] to integrate all node representations and characterize meaningful representation of the whole graph. Our paper goes further and explores a challenging and under-explored scenario, i.e., long-tailed graph-level classification.

### B. Long-Tailed Learning

Long-tailed learning aims to alleviate the impact of class imbalance on model training, and there are currently three main strategies to address this practical problem: re-sampling [15], [16], [34], [35], [36], [37], re-weighting [22], [38], [39], [40], [41], and ensemble learning [18], [26], [42], [43], [44]. For the re-sampling group, SMOTE [15] aims to rebalance the data distribution by generating new samples and performing interpolation in the tail classes, which belongs to an over-sampling approach. GraphSMOTE [35] extends SMOTE to the graph domain by encoding and synthesizing new samples based on the similarity between nodes, and training an edge generator to model relationship information. For the re-weighting group, LDAM [22] proposes a label-distribution-aware margin loss inspired by minimizing a margin-based generalization bound. DisAlign [39] develops an adaptive calibration function that enables the adjustment of classification scores for individual data points. As for ensemble learning, RIDE [18] designs an effective strategy to reduce model variance and bias as well as mitigate computational costs via dynamic expert routing. LFME [42] aggregates knowledge from multiple experts and proposes a knowledge distillation framework to learn a unified student model. Our work inherits the advantages of ensemble learning and is dedicated to organic cooperation and supervision among experts to promote effective long-tailed learning.

### C. Contrastive Learning

Contrastive learning (CL) learns the common and discriminative attributes by a contrasting principle among the positive and negative pairs. Many approaches have been proposed with competitive performance [45], [46], [47], [48], [49]. SimCLR [45] combines the data augmentations and a learnable nonlinear transformation in CL to learn the representations. MoCo [46] builds a dynamic dictionary with a queue and a moving-averaged encoder to facilitate contrastive self-supervised learning. SupCon [48] extends the contrastive approach to the fully-supervised setting, which allows for effective leverage of the label information. MoCov2 [47] improves MoCo by using a multi-layer perceptron (MLP) projection head and more data augmentations to ease the burden on the batch size in training. SACC [50] incorporates strong and weak augmentations into instance- and cluster-level CL for deep clustering. Moreover, there are many recent methods extending CL to graph domains [51], [52], [53], [54], [55], [56], [57]. GraphCL [51] designs four types of graph augmentations to incorporate various priors for effective CL. CGCN [52] proposes a semi-supervised contrastive loss to maximize the homogeneity of the original topology graph and the self-adaptive one. For long-tailed graph classification, our paper focuses on designing tailored supervised contrastive learning that balances the head and tail classes to learn effective graph-level representations.

## III. PRELIMINARIES

In this section, we first briefly present the basic notations and formal terminologies for the long-tailed graph dataset. Then, we

provide the problem definition for the long-tailed graph-level classification task and give an introduction to the GNN-based classifier.

*Notations:* Given a graph dataset  $\mathcal{G} = \{G_i, y_i\}_{i=1}^N$  with  $M$  classes, where  $G_i$  is the  $i$ -th graph and  $y_i \in \{1, \dots, M\}$  is the ground-truth class label of  $G_i$ . Let  $N_j$  denote the number of graphs in the  $j$ -th class and assume that  $N_1 \geq N_2 \geq \dots \geq N_M$  without loss of generality. The imbalance factor (IF) of the dataset is defined as  $N_1/N_M$  to measure the extent of class imbalance. Let  $p(G|y)$  be the probability density function of graph  $G$  conditioned on the class  $y$ . A dataset is long-tailed if the following relationships hold:

$$\begin{cases} \int p(G|y = j')dG \geq \int p(G|y = j'')dG, & \forall j' \leq j'', \\ \lim_{j \rightarrow \infty} \int p(G|y = j)dG = 0, \end{cases} \quad (1)$$

where  $j', j'' \in \{1, \dots, M\}$  are indices of class labels. Equation (1) indicates that for any given class index  $j' \leq j''$ , the number of samples in the  $j'$ -th class is larger than the number of samples in the  $j''$ -th class. In other words, the above formula reflects that the class size successively decays with the class index increasing and the probability finally approaches zero in the last few classes. Under the long-tailed setting in (1), the classes can be divided into head, medium, and tail classes based on different numbers of graphs. The number of graphs in head classes is far more than that of tail classes.

*Problem Definition:* The target of long-tailed graph-level classification is to train an unbiased classification model based on the long-tailed graph dataset. The model needs to learn effective and discriminative representations for all the graphs, such that it is not overwhelmed by the abundant head classes and is able to correctly classify the graphs in all classes. Further, the trained classifier should have a powerful generalization ability on a balanced test dataset.

*GNN-based Classifier:* To obtain effective probability assignments on the classes for each graph, GNN is the most popular strategy in graph representation learning. We begin with leveraging GNN as the encoder to acquire the whole graph representation, which reasonably captures the node attribute and structure information. Specifically, the propagation rule in a layer of GNN is

$$\mathbf{h}_v^{(l+1)} = \mathcal{C}_\theta^{(l)} \left( \mathbf{h}_v^{(l)}, \mathcal{A}_\theta^{(l)}(\{\mathbf{h}_{v'}^{(l)}\}_{v' \in \mathcal{N}(v)}) \right),$$

where  $\mathbf{h}_v^{(l)}$  is the learned representation of node  $v$  in Graph  $G$  at the  $l$ -th layer,  $\mathcal{N}(v)$  is the neighbors of node  $v$ ,  $\mathcal{C}_\theta^{(l)}$  and  $\mathcal{A}_\theta^{(l)}$  are the combination and aggregation functions at the  $l$ -th layer, respectively. Based on the node-level representations obtained by a  $L$ -layer GNN, a pooling operation *READOUT* is performed to merge them into a graph-level representation of graph  $G$ , which is formulated as

$$\mathbf{h} = \text{READOUT}(\{\mathbf{h}_{v'}^{(L)} : v' \in V\}), \quad (2)$$

where  $V$  is the node set of graph  $G$ . Then, the graph representation of fusing attributive and topological information can be used for graph-level classification, where the representation is fed into a classifier, such as a multi-layer perceptron (MLP) followed



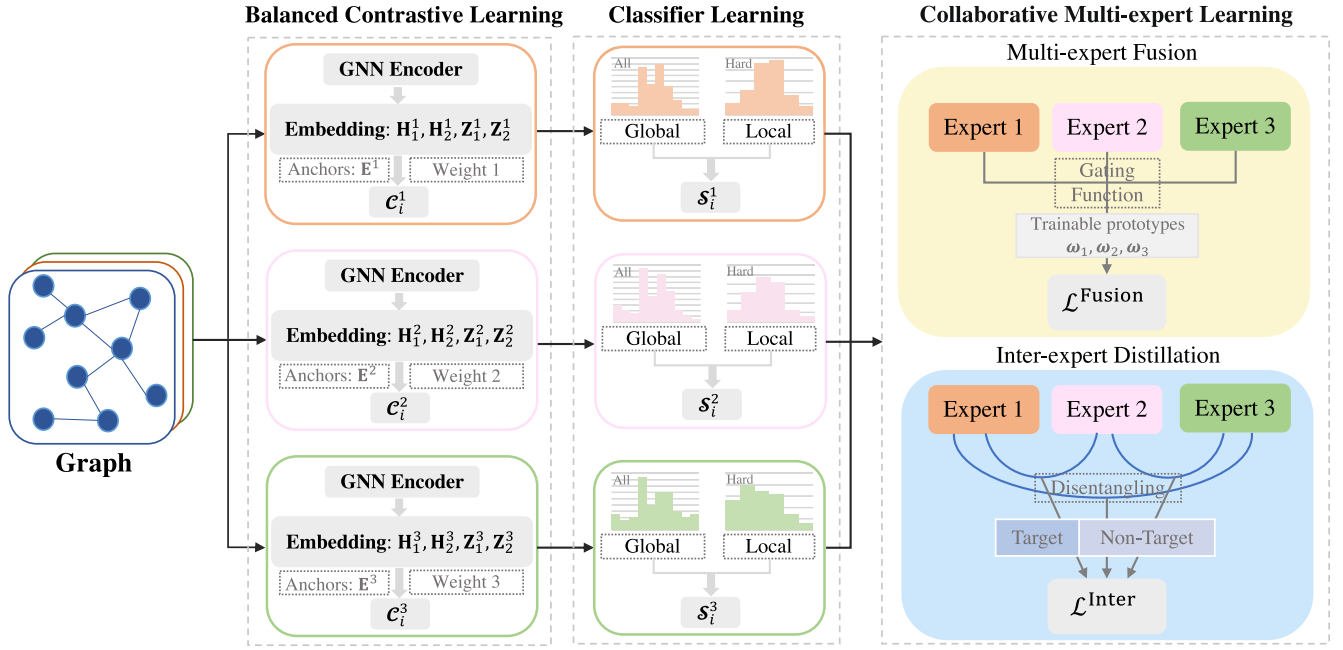


Fig. 1. Framework overview of the proposed CoMe. Balanced contrastive learning and classifier learning from both global and local views are performed to jointly guide the individual-expert training. Then the experts are integrated by gated multi-expert fusion and disentangled inter-expert knowledge distillation for long-tailed graph classification.

by a softmax function, to derive the probability assignment of graph  $G$ .

However, direct employment of the GNN-based classifier on a long-tailed dataset easily results in prediction bias, since the model may be dominated by the high-frequency head classes. Hence, ensuring the high accuracy of the GNN-based classifier for both head and tail classes is what we focus on.

#### IV. THE PROPOSED METHOD

In this section, we introduce our proposed framework named CoMe, which explores the challenging long-tailed graph classification in the context of multi-expert learning. Our CoMe mainly contains four modules, i.e., balanced contrastive learning of individual expert, individual-expert classifier learning, multi-expert fusion module, and inter-expert distillation module. Fig. 1 presents the framework overview of the proposed method. In the following, we show the four parts of our framework CoMe in detail.

##### A. Balanced Contrastive Learning of Individual Expert

The high performance of collaborative multi-expert learning is based on the fact that each expert network has excellent learning ability. Thus, we first discuss the representation learning of the individual expert. To accommodate the long-tailed dataset, we propose a tailored balanced contrastive learning (BCL) to generate discriminative representations for each expert, which can facilitate a balance between the contributions of head and tail classes and alleviate the problem caused by insufficient samples in the tail classes.

Contrastive learning [45] is a self-supervised learning that extracts meaningful representations by maximizing the similarity

of positive pairs and minimizing the correlation of negative pairs based on two augmented view of the data. To achieve this, data augmentation plays a vital role in contrastive learning to enrich the dataset and acquire the positive/negative pairs. To increase the diversity of the experts, we leverage different graph augmentation strategies in different expert networks. Specifically, we consider four strategies following [51]:

- *Attribute Masking*: We randomly mask a few entries of the node attributes in the graph. This strategy slightly disturbs the node features, which would not essentially change the semantic information.
- *Node Dropping*: We randomly delete a certain portion of nodes in each graph along with the connected edges removed, the probability of a node being removed follows a uniform distribution.
- *Edge Perturbation*: We randomly add or delete several edges from each graph. It is premised on the assumption that the semantic information is resistant to variations in edge connection patterns.
- *Subgraph*: We utilize the random walk algorithm to sample a subgraph from the original graph with the assumption that the main semantics can be preserved in the local structure to some extent.

Assume that a total of  $K$  experts are employed and in the  $k$ -th expert model, we obtain two augmented datasets  $\{\mathcal{G}_1^k, \mathcal{G}_2^k\}$  through different augmentations on the original  $N$  graphs, where the subscripts (i.e.,  $a = 1, 2$ ) denote the  $a$ -th augmented view. The graphs in  $\mathcal{G}_1^k$  and  $\mathcal{G}_2^k$  are fed into the shared GNN encoder to learn the graph-level representations as (2), which are denoted by  $\mathbf{H}_1^k = \{\mathbf{h}_1^{1,k}, \dots, \mathbf{h}_N^{1,k}\}$  and  $\mathbf{H}_2^k = \{\mathbf{h}_1^{2,k}, \dots, \mathbf{h}_N^{2,k}\}$ , respectively. As the traditional contrastive learning [45], the learned representations are then mapped into a shared space

by an MLP to compute the contrastive loss, where the mapped embeddings are denoted by  $\mathbf{Z}_1^k = \{\mathbf{z}_1^{1,k}, \dots, \mathbf{z}_N^{1,k}\}$  and  $\mathbf{Z}_2^k = \{\mathbf{z}_1^{2,k}, \dots, \mathbf{z}_N^{2,k}\}$ , respectively. To balance the contributions of the head and tail classes in contrastive learning without sacrificing the accuracies of the head classes, we introduce the trainable anchors  $E(k) = \{\mathbf{e}_1^k, \dots, \mathbf{e}_M^k\}$  as learnable class representations. Specifically, we define the balanced contrastive loss for the  $i$ -th graph and the  $k$ -th expert as

$$\mathcal{C}_i^k = - \sum_{\mathbf{z}_+ \in P(i) \cup \mathbf{e}_{y_i}^k} w(\mathbf{z}_+) \log \frac{\exp(\mathbf{z}_+ \cdot T(\mathbf{h}_i^{1,k})/\tau)}{\sum_{\mathbf{z}_j \in A(i) \cup E(k)} \exp(\mathbf{z}_j \cdot T(\mathbf{h}_i^{1,k})/\tau)}, \quad (3)$$

with

$$A(i) = \{\mathbf{z}_j \in \mathbf{Z}_1^k \cup \mathbf{Z}_2^k \setminus \mathbf{z}_i^{1,k},$$

$$P(i) = \{\mathbf{z}_j \in A(i) : y_j = y_i\},$$

$$\mathbf{z} \cdot T(\mathbf{h}_i^{1,k}) = \begin{cases} \mathbf{z} \cdot \mathbf{z}_i^{1,k}, & \mathbf{z} \in A(i), \\ \mathbf{z} \cdot \mathbf{h}_i^{1,k}, & \mathbf{z} \in E(k), \end{cases}$$

$$w(\mathbf{z}_+) = \begin{cases} \alpha, & \mathbf{z}_+ \in P(i), \\ 1, & \mathbf{z}_+ = \mathbf{e}_{y_i}^k, \end{cases}$$

where  $\tau$  is the temperature hyper-parameter,  $w(\mathbf{z}_+)$  is the weight function,  $y_i$  is the true label of the  $i$ -th graph, and  $\alpha \in [0, 1]$  is a weight parameter. The set  $A(i)$  contains all the graph representations from the  $k$ -th expert except the representation  $\mathbf{z}_i^{1,k}$  for the  $i$ -th graph from the first augmented view, and the set  $P(i)$  includes the representations of the graphs that belong to class  $y_i$  in  $A(i)$ .

Next, we theoretically illustrate the usefulness of the defined loss  $\mathcal{C}_i^k$  for representation learning on long-tailed data. Denote  $W_i$  as the cardinality of  $P(i)$ ,  $P(i) = \{\mathbf{z}_1^+, \mathbf{z}_2^+, \dots, \mathbf{z}_{W_i}^+\}$ ,  $p_q^+ = \exp(\mathbf{z}_q^+ \cdot T(\mathbf{h}_i^{1,k})/\tau) / \sum_{\mathbf{z}_j \in A(i) \cup E(k)} \exp(\mathbf{z}_j \cdot T(\mathbf{h}_i^{1,k})/\tau)$  for  $q = 1, 2, \dots, W_i$ ,  $p_e^+ = \exp(\mathbf{e}_{y_i}^k \cdot T(\mathbf{h}_i^{1,k})/\tau) / \sum_{\mathbf{z}_j \in A(i) \cup E(k)} \exp(\mathbf{z}_j \cdot T(\mathbf{h}_i^{1,k})/\tau)$ , and  $p^+ = p_1^+ + \dots + p_{W_i}^+ + p_e^+ \in [0, 1]$ . Then the balanced contrastive loss in (3) can be rewritten as

$$\mathcal{C}_i^k = -\alpha(\log p_1^+ + \dots + \log p_{W_i}^+) - \log p_e^+.$$

To achieve the minimum of  $\mathcal{C}_i^k$ , according to the Lagrange multiplier method [58], the Lagrange equation is formed as

$$F = -\alpha(\log p_1^+ + \dots + \log p_{W_i}^+) - \log p_e^+ + \lambda(p_1^+ + \dots + p_{W_i}^+ + p_e^+ - p^+),$$

where  $\lambda$  is the Lagrange multiplier. We take the first-order derivatives of  $F$  with respect to  $\lambda$ ,  $p_q^+$ ,  $p_e^+$  and set the derivatives equal to 0, which are formulated as

$$\begin{aligned} \frac{\partial F}{\partial \lambda} &= p_1^+ + \dots + p_{W_i}^+ + p_e^+ - p^+ = 0, \\ \frac{\partial F}{\partial p_q^+} &= -\alpha \frac{1}{p_q^+} + \lambda = 0, \end{aligned}$$

$$\frac{\partial F}{\partial p_e^+} = -\frac{1}{p_e^+} + \lambda = 0.$$

By jointly solving the above equations, we obtain that  $p_q^+ = \alpha p^+ / (\alpha W_i + 1)$  and  $p_e^+ = p^+ / (\alpha W_i + 1)$ . Hence, if  $p^+ = 1$ ,  $\mathcal{C}_i^k$  achieves its minimum with

$$\frac{\exp(\mathbf{z}_+ \cdot T(\mathbf{h}_i^{1,k})/\tau)}{\sum_{\mathbf{z}_j \in A(i) \cup E(k)} \exp(\mathbf{z}_j \cdot T(\mathbf{h}_i^{1,k})/\tau)} = \frac{1}{W_i + 1/\alpha}$$

satisfied for any  $\mathbf{z}_+ \in P(i)$ , which is the probability that two graphs in the same class are a positive pair.

As for the original supervised contrastive loss (SupCon) in [48], which is defined as

$$\bar{\mathcal{C}}_i^k = - \sum_{\mathbf{z}_+ \in P(i)} \log \frac{\exp(\mathbf{z}_+ \cdot T(\mathbf{h}_i^{1,k})/\tau)}{\sum_{\mathbf{z}_j \in A(i)} \exp(\mathbf{z}_j \cdot T(\mathbf{h}_i^{1,k})/\tau)}.$$

It can be similarly obtained that when  $\bar{\mathcal{C}}_i^k$  reaches its minimum, for any  $\mathbf{z}_+ \in P(i)$ , we have

$$\frac{\exp(\mathbf{z}_+ \cdot T(\mathbf{h}_i^{1,k})/\tau)}{\sum_{\mathbf{z}_j \in A(i)} \exp(\mathbf{z}_j \cdot T(\mathbf{h}_i^{1,k})/\tau)} = \frac{1}{W_i}. \quad (4)$$

We can observe in (4) that a larger  $W_i$  yields a smaller probability that two graphs in the same class are a positive pair, which causes a greater contribution from the negative log of (4) to  $\bar{\mathcal{C}}_i^k$ . It implies that when the training in SupCon converges, the high-frequency head classes have a greater impact on the loss than the low-resource tail classes since the samples in head classes possess larger  $W_i (= 2N_{y_i} - 1)$ , which inevitably inhibits the performance of the tail classes.

Let  $W_h$  and  $W_t$  be the cardinalities of  $P(i_h)$  and  $P(i_t)$ , where  $i_h$  and  $i_t$  are the head- and tail-class graphs, respectively. By introducing the class anchors  $E(k)$  and the weight  $w(\mathbf{z}_+)$  in (3), our BCL reduces the difference of the probabilities that two graphs are a positive pair in the head and tail classes from  $1/W_t - 1/W_h$  to  $1/(W_t + 1/\alpha) - 1/(W_h + 1/\alpha)$ . When  $\alpha$  is smaller, the difference could be smaller, which balances the contributions of the head and tail classes to the loss in (3). In addition, with the decrease of  $\alpha$ , the weight of the contrast between the graph and the corresponding class representation in (3) would be higher, which improves the contrast intensity between them, thereby bearing the ability to push graphs in the same class close to each other and implicitly benefits hard class learning. Hence, the proposed balanced contrastive loss enhances the performance of the learned representation over long-tailed graph data.

## B. Individual-Expert Classifier Learning

Based on the learned embedding  $\mathbf{H}^k$  from the original graphs in the aforementioned balanced contrastive learning, we feed it into a classifier network such as an MLP, to obtain the logit vectors  $\mathbf{O}^k = \{\mathbf{o}_1^k, \dots, \mathbf{o}_N^k\}$ . Then we use the softmax-like function to derive the predicted probability assignment for each graph sample. If the naive softmax function is adopted, for the  $i$ -th graph and the  $k$ -th expert, the predicted probability of assigning

to the  $j$ -th class is defined as

$$\bar{p}_{i,j}^k = \frac{\exp(o_{i,j}^k)}{\sum_{m=1}^M \exp(o_{i,m}^k)}, \quad (5)$$

where  $o_{i,j}^k$  is the  $j$ -th entry in the logit vector  $\mathbf{o}_i^k$ . Apparently, the naive softmax function cannot mitigate the impact of the long-tailed distribution on the classifier, which inevitably leads to the learned classification model dominated by the high-frequency head classes. It implies that the learned logits from the classifier implicitly incorporate the effects of the sample sizes for different classes, which is not beneficial to the testing graph whose ground-truth label belongs to the tail class. Hence, we need to eliminate the effects of the sample sizes for different classes on the logits, which can indirectly equilibrate their impacts on the trained classifier.

To this end, we leverage the idea of Bayesian inference and incorporate class frequency into the predicted probability by treating it as the prior information. Specifically, we propose the balanced predicted probability as follows,

$$\begin{aligned} p_{i,j}^k &= p(y = j | \mathbf{x} = \mathbf{o}_i^k) = \frac{p(\mathbf{x} = \mathbf{o}_i^k | y = j) p(y = j)}{p(\mathbf{x} = \mathbf{o}_i^k)} \\ &= \frac{N_j \exp(o_{i,j}^k)}{\sum_{m=1}^M N_m \exp(o_{i,m}^k)}, \end{aligned} \quad (6)$$

where  $y$  is the label of  $\mathbf{x}$  and  $N_j$  is the number of graphs in the  $j$ -th class of the dataset. The probability  $p(\mathbf{x} = \mathbf{o}_i^k | y = j)$  is defined as  $\bar{p}_{i,j}^k$  in (5) and  $p(y = j) = N_j/N$ . When the graph dataset is completely balanced, i.e.,  $N_1 = \dots = N_M$ , the defined  $p_{i,j}^k$  naturally degenerates into the vanilla softmax function. Intuitively, the incorporation of the prior information  $p(y = j)$  in  $p_{i,j}^k$  decouples the effects of sample size  $N_j$  and the logit  $\mathbf{o}_i^k$  on the prediction result. Moreover, in the testing stage, a common principle is to conduct softmax operation on the learned logit vector for the testing sample to acquire the predicted probability. Thus, accompanied by a supervised loss (e.g., cross entropy) based upon the true labels and the balanced probabilities defined in (6), the training process can be well supervised to alleviate the influence of the sample sizes for different classes on the logits, which effectively avoids the prediction performance being overwhelmed by sufficient head-class samples.

**Hard-to-classify Class Mining:** Further, under the long-tailed setting, it is ubiquitous that the predicted class for the sample in the tail class is not the ground-truth class but with a high predicted score [59]. Hence, motivated by [60], [61], additional attention shall be paid to hard class mining (HCM). For the  $i$ -th graph and the  $k$ -th expert, we explicitly extract the hard classes  $\Omega_i^k$  by selecting classes with the top- $M_{\text{hard}}$  largest logits of the  $i$ -th row in  $\mathbf{O}^k$  and the ground-truth class, which is formulated as

$$\Omega_i^k = \text{TopHC}\{j : j \neq y_i\} \cup \{y_i\}. \quad (7)$$

Then, we define the balanced predicted probability focusing on the hard classes by integrating  $\Omega_i^k$  into (6), which derives

$$\tilde{p}_{i,j}^k = \frac{N_j \exp(o_{i,j}^k)}{\sum_{m \in \Omega_i^k} N_m \exp(o_{i,m}^k)} \text{ for any } j \in \Omega_i^k. \quad (8)$$

Equation (8) presents a fine-grained characterization of the probability assignment among the hard classes. It enables more targeted and effective supervised classifier learning, which ultimately leads to improved prediction accuracy for tail classes.

Based on the above discussion, to derive the expert network with strong discrimination ability for graph classification, we conduct classifier learning for each expert from both global and local perspectives. From the global view, we resort to  $p_{i,j}$  in (6) to empower the ability of balancing the contributions of head and tail classes in training, whereas from the local view, we enhance the hard class learning with the help of  $\tilde{p}_{i,j}^k$  in (8). Formally, for the  $i$ -th graph and the  $k$ -th expert, the supervised loss of individual classifier is defined as

$$\mathcal{S}_i^k = -(\log(p_{i,y_i}^k) + \log(\tilde{p}_{i,y_i}^k)).$$

### C. Multi-Expert Fusion Module

In addition to individual learning for each expert, how to organically combine the losses from different experts is also the focus of multi-expert learning since the graphs in different classes may have heterogeneous interactions with different expert networks. Here we employ gating functions to control the fusion process, which can be regarded as learnable weights and trained along with expert networks. By the law of total probability, the joint prediction mechanism with multiple experts is formulated as

$$p(y = y_i | \mathbf{O}) = \sum_{k=1}^K p(e_i = k | \mathbf{o}_i^k) p(y = y_i | e_i = k, \mathbf{o}_i^k),$$

where  $\mathbf{o}_i^k = (o_{i,1}^k, \dots, o_{i,M}^k)^\top$ ,  $\mathbf{O}$  is formed by stacking the logit vectors of  $K$  experts,  $e_i$  is a latent variable indicating the expert index for the  $i$ -th graph,  $p(e_i = k | \mathbf{o}_i^k)$  is the gating function of the  $k$ -th expert with  $\sum_{k=1}^K p(e_i = k | \mathbf{o}_i^k) = 1$  satisfied, and  $p(y = y_i | e_i = k, \mathbf{o}_i^k)$  represents the predicted probability of the  $k$ -th expert. For the gating function, we parameterize it as an input-dependent soft assignment based on the cosine similarity between the embedded logit  $\mathbf{o}_i^k$  and the trainable gating prototype  $\boldsymbol{\omega}^k$ , which has the form as

$$p(e_i = k | \mathbf{o}_i^k) = \frac{\exp(\mathbf{o}_i^k \cdot \boldsymbol{\omega}^k / \kappa)}{\sum_{k=1}^K \exp(\mathbf{o}_i^k \cdot \boldsymbol{\omega}^k / \kappa)},$$

where  $\kappa$  is the temperature hyper-parameter. The value  $p(y = y_i | e_i = k, \mathbf{o}_i^k)$  can be substituted by the balanced predicted probability  $p_{i,y_i}^k$  or its variant  $\tilde{p}_{i,y_i}^k$ . Hence, with the adoption of both  $p_{i,y_i}^k$  and  $\tilde{p}_{i,y_i}^k$  from global and local views, the fused supervised loss (FSL) of classifiers can be equipped as

$$\mathcal{L}^{\text{FSL}} = \sum_{i=1}^N \sum_{k=1}^K p(e_i = k | \mathbf{o}_i^k) \mathcal{S}_i^k.$$

In addition, following an analogous pipeline, we can similarly fuse the balanced contrastive loss  $\mathcal{C}_i^k$  in (3) under different experts, which derives the fused contrastive loss (FCL) as

$$\mathcal{L}^{\text{FCL}} = \sum_{i=1}^N \sum_{k=1}^K p(e_i = k | \mathbf{o}_i^k) \mathcal{C}_i^k.$$

Thoroughly, for the  $N$  graphs and  $K$  experts, the total fusion loss for multiple experts is defined as

$$\begin{aligned} \mathcal{L}^{\text{Fusion}} &= \mathcal{L}^{\text{FSL}} + \eta \mathcal{L}^{\text{FCL}} \\ &= \sum_{i=1}^N \sum_{k=1}^K p(e_i = k | \mathbf{o}_i^k) (\mathcal{S}_i^k + \eta \mathcal{C}_i^k), \end{aligned} \quad (9)$$

where  $\eta$  is the pre-defined hyper-parameter to adjust the weight between  $\mathcal{S}_i^k$  and  $\mathcal{C}_i^k$ . With the adaptive fusion of multiple experts, the diversity of the whole training network is increased, which can jointly boost the performance of the graph-level classification task on long-tailed data.

#### D. Inter-Expert Distillation Module

Besides the fusion of the experts, the knowledge distillation among the experts is also significant to allow each expert network to learn extra signals from others and achieve information sharing. We employ the Kullback-Leibler (KL) divergence-like metric to support the inter-expert distillation.

Here we define the balanced predicted probability among the non-target classes as

$$\check{p}_{i,j}^k = \frac{N_j \exp(o_{i,j}^k)}{\sum_{m \neq y_i} N_m \exp(o_{i,m}^k)}, j \in \{1, \dots, M\} \setminus \{y_i\}.$$

We easily have  $\check{p}_{i,j}^k = p_{i,j}^k / (1 - p_{i,y_i}^k)$ . In the standard KL divergence, the distance between the probability assignments of two experts for the  $i$ -th graph can be formulated as

$$\begin{aligned} \text{KL}(\mathbf{p}_i^k || \mathbf{p}_i^q) &= p_{i,y_i}^k \log \left( \frac{p_{i,y_i}^k}{p_{i,y_i}^q} \right) + \sum_{j \neq y_i} p_{i,j}^k \log \left( \frac{p_{i,j}^k}{p_{i,j}^q} \right) \\ &= p_{i,y_i}^k \log \left( \frac{p_{i,y_i}^k}{p_{i,y_i}^q} \right) + (1 - p_{i,y_i}^k) \log \left( \frac{1 - p_{i,y_i}^k}{1 - p_{i,y_i}^q} \right) \\ &\quad + (1 - p_{i,y_i}^k) \sum_{j \neq y_i} \check{p}_{i,j}^k \log \left( \frac{\check{p}_{i,j}^k}{\check{p}_{i,j}^q} \right) \\ &\triangleq \text{KL}(\mathbf{b}_i^k || \mathbf{b}_i^q) + (1 - p_{i,y_i}^k) \text{KL}(\check{\mathbf{p}}_i^k || \check{\mathbf{p}}_i^q), \end{aligned} \quad (10)$$

where  $\mathbf{p}_i^k = (p_{i,1}^k, \dots, p_{i,M}^k)^\top$ ,  $\mathbf{b}_i^k = (p_{i,y_i}^k, 1 - p_{i,y_i}^k)^\top$  is the binary probability assignment vector of whether it belongs to the target class (i.e., the ground-truth class), and  $\check{\mathbf{p}}_i^k = (\check{p}_{i,1}^k, \dots, \check{p}_{i,y_i-1}^k, \check{p}_{i,y_i+1}^k, \dots, \check{p}_{i,M}^k)^\top$  represents the probability assignment among the non-target classes.  $\text{KL}(\mathbf{b}_i^k || \mathbf{b}_i^q)$  measures the similarity of the binary probability assignments of the target class in different expert networks, while  $\text{KL}(\check{\mathbf{p}}_i^k || \check{\mathbf{p}}_i^q)$  represents the similarity of the probability assignments among the non-target classes in different expert networks. The weight  $1 - p_{i,y_i}^k$  is negatively correlated with the prediction performance. Under the long-tailed setting, we can find from (10)

that due to the natural weight, the high prediction confidence on the graphs in the head classes would inevitably suppress the information sharing among the non-target classes between two experts; whereas, the low prediction confidence on the samples in the tail classes would spontaneously hinder the mutual learning of target knowledge between the experts, which further inhibits the prediction performance of the tail classes under the framework of collaborative multi-expert learning. Hence, to better leverage the advantage of multi-expert learning, we propose a disentangled metric (DKL) for knowledge distillation, which is defined as

$$\text{DKL}(\mathbf{p}_i^k || \mathbf{p}_i^q; \beta_1, \beta_2) = \beta_1 \text{KL}(\mathbf{b}_i^k || \mathbf{b}_i^q) + \beta_2 \text{KL}(\check{\mathbf{p}}_i^k || \check{\mathbf{p}}_i^q),$$

where  $\beta_1, \beta_2$  are weight hyper-parameters. With the flexible weight, DKL effectively avoids mutual inhibition between the distillations of the target and non-targets to promote high-performance cooperation.

Moreover, we analogously analyze the KL divergence  $\text{KL}(\check{\mathbf{p}}_i^k || \check{\mathbf{p}}_i^q)$  between two different experts' probability assignments that focus on the hard classes with  $\check{\mathbf{p}}_i^k = (\check{p}_{i,1}^k, \dots, \check{p}_{i,M}^k)^\top$  and generate a disentangled version denoted by  $\text{DKL}(\check{\mathbf{p}}_i^k || \check{\mathbf{p}}_i^q; \beta_1, \beta_2)$ . Based on  $\text{DKL}(\mathbf{p}_i^k || \mathbf{p}_i^q; \beta_1, \beta_2)$  and  $\text{DKL}(\check{\mathbf{p}}_i^k || \check{\mathbf{p}}_i^q; \beta_1, \beta_2)$ , we incorporate both the global and local perspectives as Section IV-C to further promote the distillation performance for the long-tailed data. Accordingly, the total inter-expert distillation loss is defined as

$$\begin{aligned} \mathcal{L}^{\text{Inter}} &= \sum_{k=1}^K \sum_{q \neq k}^K \sum_{i=1}^N (\text{DKL}(\mathbf{p}_i^k || \mathbf{p}_i^q; \beta_1, \beta_2) \\ &\quad + \text{DKL}(\check{\mathbf{p}}_i^k || \check{\mathbf{p}}_i^q; \beta_1, \beta_2)). \end{aligned} \quad (11)$$

#### E. Joint Optimization Module for Graph Classification

Graph classification is essentially a supervised task and we focus on the long-tailed data in this paper. Toward this end, we incorporate the multi-expert learning framework to increase the diversity of the whole network, where individual learning, multi-expert fusion and inter-expert distillation are discussed to jointly promote the performance of the long-tailed graph classification task.

Formally, we unite the multi-expert fusion loss  $\mathcal{L}^{\text{Fusion}}$  in (9) and the inter-expert distillation loss  $\mathcal{L}^{\text{Inter}}$  in (11) to jointly optimize our proposed framework CoMe, where the total loss  $\mathcal{L}$  is

$$\mathcal{L} = \mathcal{L}^{\text{Fusion}} + \epsilon \mathcal{L}^{\text{Inter}}, \quad (12)$$

where  $\epsilon$  is the pre-defined hyper-parameter to adjust the influence of  $\mathcal{L}^{\text{Fusion}}$  and  $\mathcal{L}^{\text{Inter}}$  on training.

After converges, we feed each test instance into the network to obtain the logit vectors  $\{\mathbf{o}_{\text{test}}^k\}_{k=1}^K$  and output the predicted probability vector  $\mathbf{p}_{\text{test}}$  by performing the softmax operation on the averaged logit  $\bar{\mathbf{o}}_{\text{test}} = \sum_{k=1}^K \mathbf{o}_{\text{test}}^k / K$ . The whole process of our proposed CoMe is summarized in Algorithm 1.



**Algorithm 1:** The Pseudo-Code of the Proposed CoMe.

---

**Input:** Graph dataset  $\mathcal{G} = \{G_i, y_i\}_{i=1}^N$ ; Class number  $M$ ;  
Expert number  $K$ ; Maximum iterations  $T_{max}$ ;  
**Output:** Classification result  $y$ ;

- 1: Initialize the parameters in balanced contrastive learning;
- 2: **for**  $t = 1$  to  $T_{max}$  **do**
- 3: Obtain  $\{\mathcal{G}_1^k, \mathcal{G}_2^k\}_{k=1}^K$  by graph augmentation;
- 4: Update the embeddings  $\{\mathbf{H}_1^k, \mathbf{H}_2^k, \mathbf{Z}_1^k, \mathbf{Z}_2^k\}_{k=1}^K$  and the logits  $\{\mathbf{O}^k\}_{k=1}^K$  by the GNN-based encoder;
- 5: Extract the hard classes  $\Omega_i^k$  in (7) traversing all the graphs and experts;
- 6: Calculate the losses  $\mathcal{L}^{\text{Fusion}}$  and  $\mathcal{L}^{\text{Inter}}$  in (9) and (11), respectively;
- 7: Conduct backpropagation and update the whole network in CoMe by minimizing  $\mathcal{L}$  in (12);
- 8: **end for**
- 9: Obtain the predicted probability  $\mathbf{p}_{\text{test}}$  with the averaged logit  $\bar{\mathbf{o}}_{\text{test}}$  on the test instance;
- 10: **return**  $y$ ;

---

*F. Computational Complexity Analysis*

In training, we adopt the mini-batch stochastic gradient descent to optimize our method. Assume that the batch size is  $B$  and the complexity of producing embeddings and logits from the GNN-based encoder is  $O(W)$ . For  $K$  experts, we calculate the balanced contrastive loss in  $O(Kd(\sum_{i=1}^B |A(i)| + BM))$  time, where  $A(i)$  is defined in (3),  $d$  is the dimension of the embeddings and  $M$  is the number of classes. The complexity of classifier learning is  $O(KBM)$ . Moreover, the time complexities of the multi-expert fusion and inter-expert distillation are  $O(K(M+B))$  and  $O(KB(M-1))$ , respectively. Hence, the total computational complexity of our approach is  $O(KW + Kd\sum_{i=1}^B |A(i)| + KBM)$ .

**V. EXPERIMENTS**

In this section, we first introduce the experimental settings which include benchmark datasets, compared baselines, and implementation details of the proposed method. Then we conduct experiments to validate the effectiveness of CoMe. We aim to answer the following research questions.

- *RQ1*: Does our proposed CoMe outperform baseline methods in long-tailed graph classification?
- *RQ2*: How do different components of CoMe contribute to the overall classification performance?
- *RQ3*: How do the hyperparameters in CoMe affect the final classification performance?
- *RQ4*: How does disentangled knowledge distillation affect the classification performance of CoMe?

*A. Experimental Setup*

**Datasets:** We evaluate the proposed CoMe against several baselines on seven publicly accessible datasets from various fields, including a) social networks: COLLAB [62], b) synthetic:

Synthetic [63], c) bioinformatics: ENZYMES [64], and d) computer vision: MNIST [65], Letter-high [66], Letter-low [66], and COIL-DEL [66]. Among the seven datasets utilized in the experiments, COLLAB and ENZYMES datasets are natural graph datasets derived from real-world data, while the Synthetic dataset is a synthetic graph dataset. Additionally, the visual-world MNIST dataset is represented as graphs by setting the superpixels as nodes and their spatial relations as edges, while the Letter-high, Letter-low, and COIL-DEL datasets are transformed into graph structures by representing the lines in letters as undirected edges and considering the endpoints of these lines as nodes, all originating from real-world image datasets. To ensure that the datasets follow Zipf's law exactly [67], we processed the original training sets into standard long-tailed datasets with different imbalance factors (IFs), while the validation and test sets remained to be balanced. We choose distinct IFs for each dataset to ensure the number of training samples in the tail class with the least samples falls within the range of 2 to 4.

**Baselines:** To demonstrate the efficacy of our proposed framework, we compare our CoMe with a range of competitive long-tailed learning baselines. Below we give a brief introduction to nine baseline models, which can be divided into four main categories, i.e., , data re-balancing, loss re-weighting, information augmentation, and contrastive learning based methods. Among the baseline methods considered, Graph augmentation, G<sup>2</sup> GNN, and GraphCL are specifically designed for long-tailed learning on graphs, whereas CB loss, LACE loss, SupCon, and SBCL are adapted from approaches originally developed for long-tailed classification tasks in visual world.

- GraphSAGE [10] serves as a basic GNN encoder in our implementation, where we utilize the mean aggregator to aggregate feature information. This method is used as the base encoder for CoMe and other baselines.
- Over-sampling [68] technique often incorporates repeating samples from tail classes randomly as a means of making datasets balanced.
- CB loss [38] is a loss re-weighting approach at the class level. To tackle the training problem for imbalanced data, CB loss adds class-balanced weighting to the loss function for class  $i$  inversely, which is proportional to the effective number of samples.
- LACE loss [69] is another re-weighting technique that adjusts the prediction probabilities using the label frequencies, the LACE loss function employs adjustment to the logits during the model inference phase.
- Graph augmentation [70] is a popular technique in graph representation learning that enhances model generalization and generates extra training data. We use over-sampling to increase training data and select either edge permutation or node dropping as one of the two fundamental topological augmentations.
- G<sup>2</sup>GNN [71] measures the graph kernel-based similarity between different graph samples to construct a Graph-of-Graph (GoG), which links graphs with their



TABLE I  
OVERALL PERFORMANCE (%) WITH VARIOUS IFs ON SEVEN BENCHMARK DATASETS FOR LONG-TAILED GRAPH CLASSIFICATION

Model	COLLAB		Synthie		ENZYMES		MNIST		Letter-high		Letter-low		COIL-DEL	
	IF=10	IF=20	IF=15	IF=30	IF=15	IF=30	IF=50	IF=100	IF=25	IF=50	IF=25	IF=50	IF=10	IF=20
GraphSAGE	63.07	53.33	34.74	30.25	30.66	25.16	68.67	63.46	51.06	42.16	86.00	84.32	38.80	31.32
Over-sampling	72.33	70.25	35.25	33.50	32.33	28.50	64.69	59.78	53.62	44.20	88.48	86.72	39.20	26.96
CB loss	68.78	65.85	34.75	30.75	32.19	26.83	68.85	63.40	53.76	45.06	87.46	85.44	41.72	32.34
LACE loss	68.33	64.77	33.25	30.85	31.16	25.50	69.72	64.59	47.46	38.94	87.89	84.69	41.96	32.18
Augmentation	72.85	71.14	39.37	35.37	32.08	26.75	72.18	68.17	49.28	42.36	88.32	86.40	38.18	30.80
G <sup>2</sup> GNN <sub>n</sub>	73.94	71.89	38.08	27.94	35.00	29.17	70.91	66.73	<u>58.91</u>	51.12	89.49	87.98	38.32	27.98
G <sup>2</sup> GNN <sub>e</sub>	<u>74.50</u>	<u>72.76</u>	40.19	37.53	35.83	29.50	73.69	70.31	58.85	49.96	89.84	87.80	39.18	31.06
GraphCL	69.33	67.36	40.25	36.25	36.66	29.83	69.37	65.12	57.34	48.93	89.28	87.89	42.02	33.19
SupCon	69.25	67.14	40.34	37.25	37.08	30.67	69.76	64.88	57.29	48.93	89.12	87.36	42.93	34.20
SBCL	71.63	69.12	<u>42.08</u>	<u>38.19</u>	<u>37.63</u>	<u>32.41</u>	<u>75.06</u>	<u>72.12</u>	<u>57.73</u>	<u>51.38</u>	<u>90.54</u>	<u>89.03</u>	<u>44.78</u>	<u>37.64</u>
CoMe	<b>76.88</b>	<b>74.40</b>	<b>42.25</b>	<b>38.50</b>	<b>38.00</b>	<b>33.50</b>	<b>76.24</b>	<b>72.80</b>	<b>63.42</b>	<b>54.12</b>	<b>91.67</b>	<b>90.56</b>	<b>46.56</b>	<b>38.88</b>

The best results are shown in boldface and the second-best results are underlined.

k-nearest neighbors. After constructing the kNN graph, neighboring graph representations are aggregated together via the GoG propagation on the established kNN graph.

- GraphCL [51] proposes four distinct strategies to augment input graphs and learn graph-level representations, which aims to maximize the mutual information between the original graph and its augmented variants.
- SupCon [48] is an extended version of the contrastive loss, which leverages label information effectively. SupCon allows more than one view to be positive so that views of the same label can be attracted to each other in the embedding space.
- SBCL [72] tackles the limitations of contrastive learning by clustering each head class into multiple subclasses of comparable sizes to the tail classes, thereby achieving subclass balance and learning more balanced representation space for long-tailed data.

*Implementation Details:* In all experiments, we used GraphSAGE [10] as the GNN backbone encoder, with a two-layer MLP classifier. To optimize all models, we utilized the Adam optimizer with a fixed learning rate of 0.0001 and a batch size of 32. For our proposed CoMe, we set the expert number  $K$  to 3 to balance performance and efficiency. We also set the temperature hyper-parameter  $\tau$  and the contrast weight  $\alpha$  for the BCL module to 0.2 and 0.05, respectively. For joint training, we set the contrast weight  $\eta$  to 1.0 and the inter-expert distillation weight  $\epsilon$  to 0.6. Moreover, we tuned the hyper-parameters of HCM number  $M_{\text{hard}}$ , and the DKL hyper-parameters  $\beta_1$  and  $\beta_2$  for each dataset. All baseline models are implemented in PyTorch using the open-sourced code provided by the original paper. The top-1 average accuracy of 10 run times is employed for evaluation.

## B. Overall Comparison (RQ1)

In this section, we evaluate the performance of CoMe on long-tailed graph classification and compare it with various baseline methods. Table I presents the experimental results on

seven benchmarks with different IFs. Based on the quantitative results, the following observations can be made:

- Based on the classification accuracies on all seven datasets, we observed a sharp decrease in the performance of all methods as the long-tailedness between head and tail classes increases. This indicates that GNNs are highly susceptible to long-tailed distribution and degrade severely in such long-tailed settings.
- For all four types of baseline methods, information augmentation and contrastive learning approaches mainly focus on learning better representations, while loss re-balancing methods aim to balance the classifier during training. From the results, it can be observed that existing baselines fail to balance the representation learning and classifier training in long-tailed learning, which leads to sub-optimal classification performance. Overall, information augmentation approaches outperform re-balancing approaches on most datasets as they incorporate additional knowledge to enrich the tail classes. Moreover, contrastive learning baselines demonstrate relatively stable performance across all datasets.
- From the table, it can be concluded that our CoMe achieves the best performance on all seven datasets compared with all the baselines. It is mainly attributed to the fact that for effectively addressing the long-tailed problem in graph data, our CoMe not only enhances representation learning with balanced contrastive learning but also promotes classifier balancing with balanced predicted probability. Compared to methods that directly adapt from image long-tailed classification approaches, which mainly concentrate on rebalancing the classifier (CB loss, LACE loss) or improving representations (SupCon, SBCL), our approach excels in effectively addressing the specific long-tailed problem encountered in graph data. In addition, CoMe outperforms other competitors with a large margin under strict imbalance settings (COIL-DEL with IF = 10 and 20), where most classes have fewer than 5 training instances. This highlights the importance of employing hard class mining to derive the expert network with strong discrimination ability.

TABLE II

COMPARISON WITH SEVERAL VARIANTS FOR ABLATION STUDY (%) ON LETTER-HIGH AND ENZYMES DATASETS (MULTI-EXPERT FRAMEWORK: MeF; BALANCED CONTRASTIVE LEARNING: BCL; HARD CLASS MINING: HCM; GATED MULTI-EXPERT FUSION: GMeF; DISENTANGLED INTER-EXPERT DISTILLATION: DIeD)

	MeF	BCL	HCM	GMeF	DIeD	Letter-high Accuracy	$\Delta$	ENZYMES Accuracy	$\Delta$
$M_1$	✓					54.92	-	35.00	-
$M_2$	✓		✓			57.47	+2.55	35.67	+0.67
$M_3$	✓	✓				62.34	+7.42	36.93	+1.93
$M_4$	✓	✓	✓			62.54	+7.62	37.00	+2.00
$M_5$	✓	✓	✓	✓		62.85	+7.93	37.13	+2.13
$M_6$	✓	✓	✓		✓	62.65	+7.73	37.33	+2.33
$M_7$	✓	✓	✓	✓	✓	63.42	+8.50	38.00	+3.00

TABLE III

EFFECTIVENESS ANALYSIS (%) OF BALANCED CONTRASTIVE LEARNING AND BALANCED PREDICTED PROBABILITY ON THE SYNTHIE AND LETTER-HIGH DATASETS (UNSUPERVISED CONTRASTIVE LEARNING: UCL; SUPERVISED CONTRASTIVE LEARNING: SCL; BALANCED CONTRASTIVE LEARNING: BCL; BALANCED PREDICTED PROBABILITY: BPP)

UCL	SCL	BCL	BPP	Synthie Accuracy	$\Delta$	Letter-high Accuracy	$\Delta$
				36.25	-	53.60	-
✓			✓	37.35	+1.10	60.37	+6.77
	✓		✓	40.23	+3.98	61.65	+8.05
		✓		39.62	+3.37	60.69	+7.09
		✓	✓	42.25	+6.00	63.42	+9.82

### C. Ablation Study (RQ2)

*Ablation study on major components:* We conduct an ablation analysis for the key components in our CoMe, i.e., multi-expert framework (MeF), balanced contrastive learning (BCL), hard class mining (HCM), gated multi-expert fusion (GMeF), and disentangled inter-expert distillation (DIeD). We analyze the effect of each component by adding them to the base model until obtaining the complete method. In all cases, the classification results in the training process are generated by the balanced predicted probability. Table II demonstrates the results of the ablation study. The base model ( $M_1$ ) includes only the model ensemble, without any of the other key components. Adding HCM ( $M_2$ ) and BCL ( $M_3$ ) to the baseline model significantly improves accuracy on both datasets. Particularly,  $M_3$  shows a substantial improvement, emphasizing the importance of balanced contrastive learning for effective representation learning. Furthermore, the combination of HCM and BCL ( $M_4$ ) leads to further accuracy improvements. The introduction of gating prototypes ( $M_5$ ) enhances discrimination ability and enables the dynamic fusion of multiple experts, resulting in higher accuracy. Incorporating disentangled inter-expert distillation ( $M_6$ ) also contributes to accuracy improvement by facilitating effective knowledge sharing between experts through disentangled knowledge distillation. Finally, when all key components are combined ( $M_7$ ), the model achieves the best performance, underscoring the collective impact of the key components on overall accuracy.

*Ablation study on BCL and BPP:* To analyze the effect of balanced contrastive learning, we carry out an experiment by substituting BCL with unsupervised contrastive learning (UCL) and supervised contrastive learning (SCL). Table III shows the experimental results on two datasets. From the table, we can observe that the classification accuracy of using both BCL and BPP is higher than combining UCL or SCL with BPP, indicating that BCL learns better representation in the long-tailed settings. However, employing BCL without BPP results in a minor drop in classification accuracy. These results suggest that the combination of BCL and BPP can lead to significant improvements in both representation learning and classifier learning,

which can ultimately enhance the overall performance of the model.

### D. Hyper-Parameters Analysis (RQ3)

In this section, we study the key hyper-parameters in our CoMe, i.e., the hard class mining number  $M_{\text{hard}}$ , the balanced contrastive learning weight  $\eta$ , and the inter-expert distillation weight  $\epsilon$ . We tune  $M_{\text{hard}}$  separately by fixing other hyper-parameters, and jointly tune the loss weight hyper-parameters  $\eta$  and  $\epsilon$ . Fig. 2 demonstrates how these hyper-parameters affect classification performance.

The impact of  $M_{\text{hard}}$  on the classification performance of CoMe is illustrated in Fig. 2(a) and (b). The experimental results on the COIL-DEL and Letter-high datasets show a similar trend, where the classification accuracy initially improves with an increase in  $M_{\text{hard}}$ . Our model achieves the best performance when the hard class ratio ( $M_{\text{hard}}/M$ ) is around 0.3, which corresponds to 30 out of 100 classes for the COIL-DEL dataset and 5 out of 15 classes for the Letter-high dataset. We also conduct more fine-grained adjustments to  $M_{\text{hard}}$  around 30 on the COIL-DEL dataset in Fig. 3. When  $M_{\text{hard}}$  is set around 30, the differences in performance are minimal, which indicates a relatively consistent and steady trend. However, setting  $M_{\text{hard}}$  to a value that is either too small or too large brings limited gains due to the under- or over-exploration of hard categories. Further, when dealing with a new dataset, the optimal value of  $M_{\text{hard}}$  may vary depending on the dataset's specific characteristics. We recommend conducting a small-scale hyper-parameter tuning around  $0.3 * M$  to select a suitable  $M_{\text{hard}}$  for the new dataset.

The loss weight parameters are critical in adjusting the importance of different loss components during training. The impact of  $\eta$  and  $\epsilon$  is visualized in Fig. 2(c) and (d). In our experiments, we vary  $\eta \in \{0.1, 0.5, 1.0, 1.5, 2.0\}$  and  $\epsilon \in \{0.2, 0.4, 0.6, 0.8, 1.0\}$  on the Synthie dataset (IF = 15) and Letter-high dataset (IF = 25). The results indicate that the model is not significantly impacted by the change of  $\epsilon$ , and the optimal performance is achieved when  $\epsilon$  is set to 0.4 for Synthie and 0.6 for Letter-high. Moreover, the classification accuracy increases when  $\eta$  is raised from 0.1 to 0.5, suggesting that balanced contrastive learning significantly contributes to better representation learning.

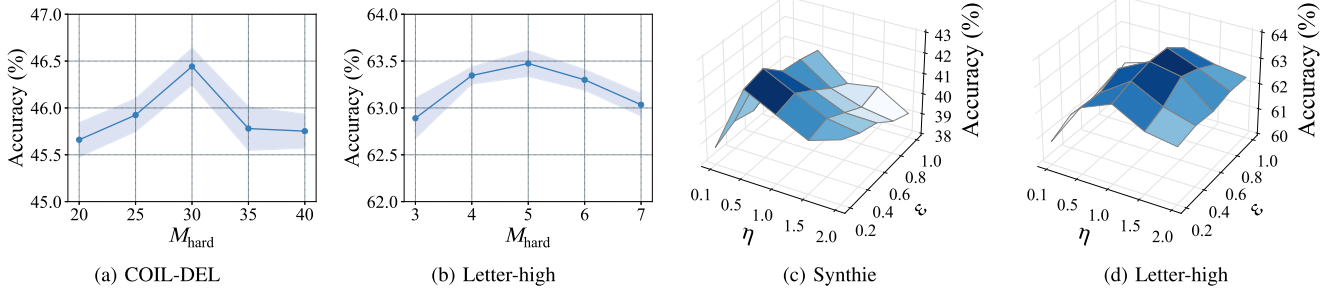


Fig. 2. Sensitivity analysis of (a)–(b): hard class number and (c)–(d): loss weight parameters.

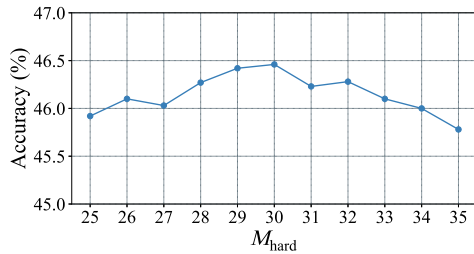


Fig. 3. Fine-grained sensitivity analysis of hard class number on COIL-DEL (IF = 10).

TABLE IV  
EFFECTIVENESS ANALYSIS (%) OF DISENTANGLED KNOWLEDGE DISTILLATION ON THE ENZYMES, LETTER-HIGH AND LETTER-LOW DATASETS (TARGET CLASS DISTILLATION: TCD; NON-TARGET CLASS DISTILLATION: NTCD; DISENTANGLED DISTILLATION: DD)

TCD	NTCD	DD	ENZYMES Accuracy	$\Delta$	Letter-high Accuracy	$\Delta$	Letter-low Accuracy	$\Delta$
			37.13	-	62.85	-	90.88	-
✓	✓		37.67	+0.54	63.18	+0.33	91.20	+0.32
✓			37.31	+0.18	62.76	-0.09	91.11	+0.23
	✓		37.57	+0.44	62.95	+0.10	91.31	+0.43
✓	✓	✓	38.00	+0.87	63.42	+0.57	91.67	+0.79

#### E. Effect of Disentangled Knowledge Distillation (RQ4)

*Performance gain of each distillation component:* Here we individually study the effects of knowledge distillations on target and non-targets over the ENZYMES, Letter-high, and Letter-low datasets. The accuracy and performance gain are reported in Table IV. For each dataset, we present results for five variants of the method, including 1) the vanilla training baseline without inter-expert distillation, 2) classical knowledge distillation using both target and non-target class distillation, 3) target class distillation only, 4) non-target class distillation only, and 5) disentangled distillation.

From the table, we can observe that the adoption of classical knowledge distillation leads to an improvement in the overall performance of the model. However, we notice that applying target class distillation solely leads to less improvement or even degrades the overall performance (e.g., 0.09% drop on Letter-high dataset). It can also be observed that the performances of using non-target class distillation only are comparable and

even better than the classical knowledge distillation (e.g., 0.43% accuracy gain on Letter-low dataset). These experimental results illustrate that knowledge related to non-target classes could be more critical than target-class-specific knowledge. Finally, the application of disentangled distillation provides the most significant improvement to our model. This finding suggests that separating different knowledge components in a disentangled manner helps alleviate the issue of information-sharing suppression between target class and non-target classes of different experts.

*Influence of distillation weight:* To demonstrate the effectiveness of the distillation weight of target and non-target class knowledge distillation, we carry out an experiment to evaluate the performance under varying distillation weight settings. The experimental results on three datasets are shown in Fig. 4. We vary the target class distillation weight  $\beta_1$  across  $\{0.5, 1.0, 1.5, 2.0\}$  for all three datasets. Additionally, we tune the non-target class distillation weight  $\beta_2$  within the range of  $\{0.5, 1.0, 2.0, 4.0\}$  for ENZYMES,  $\{1.0, 2.0, 3.0, 4.0\}$  for Letter-high, and  $\{1.0, 2.0, 4.0, 8.0\}$  for Letter-low.

Overall, the classification performance remains stable when increasing  $\beta_1$  from 0.5 to 2.0. Additionally, we observed that the highest improvement from target class knowledge distillation is achieved when  $\beta_1$  is set to approximately 1.0. These results suggest that the model is relatively insensitive to the setting of different  $\beta_1$  values. However, different settings of  $\beta_2$  can have a significant impact on the classification performance, and the optimal value of  $\beta_2$  varies among the different datasets. Specifically, for the ENZYMES dataset, the best performance is achieved when  $\beta_2$  is set to 0.5, and increasing the contribution of non-target knowledge distillation results in a drop in accuracy. For the Letter-high and Letter-low datasets, we observed an initial increase in accuracy as  $\beta_2$  becomes larger, but an excessively large value of  $\beta_2$  can also lead to performance degradation. These experimental results demonstrate that the optimal selection of the non-target distillation weight is closely related to the confidence of the experts. When the experts are confident, i.e., the predicted probability of the target class is much higher than all non-target classes, the non-target knowledge should be valued more. Specifically, the experts are less confident in the ENZYMES dataset, and the value of  $\beta_2$  should be set lower than those in the other two datasets. Moreover, a larger value of  $\beta_2$  could increase the gradient contributed by non-target classes, potentially leading to a degradation of the model's accuracy.



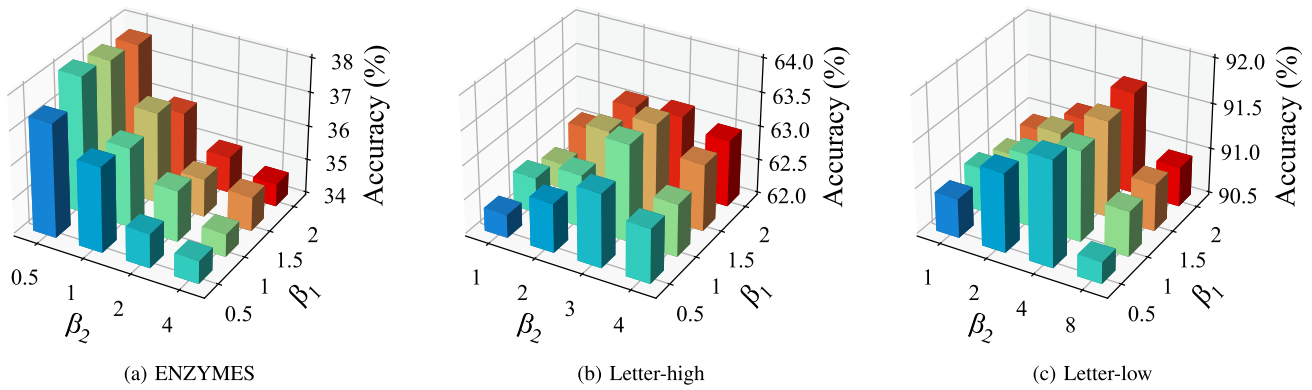


Fig. 4. Performance comparison w.r.t. different settings of the distillation weight hyper-parameters.

## VI. CONCLUSION

In this paper, we study long-tailed graph-level classification and propose a novel method termed Collaborative Multi-expert Learning (CoMe) to handle the imbalanced setting of the distribution in graph datasets. CoMe incorporates the balanced contrastive learning for representation learning, accompanied by the classifier learning of each expert to alleviate the influence of the sample sizes in different classes and enhance the hard class mining. To combine the advantages of multiple experts, we design a mechanism to fuse their diversities in a multi-expert framework, thus enhancing the cooperation. Furthermore, we develop a disentangled knowledge distillation to encourage the knowledge transfer and mutual supervision among multiple experts. Extensive experiments demonstrate that our proposed CoMe consistently outperforms the competitive baseline methods on various benchmark graph datasets.

In our future work, there are several aspects of our proposed model that deserve further investigation: i) the distribution of the test set is unknown in real-world scenarios, and better mechanisms need to be designed to overcome the unknown distribution rather than the balanced distribution; ii) extending our framework to more challenging settings such as noisy labels or noisy graphs; iii) exploring more fundamental theoretical research related to generalization beyond the training distribution such as out-of-distribution problem.

## ACKNOWLEDGMENTS

The authors are grateful to the anonymous reviewers for critically reading the manuscript and for giving important suggestions to improve their paper.

## REFERENCES

- [1] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," 2016, *arXiv:1609.02907*.
- [2] Y. Xie, S. Lv, Y. Qian, C. Wen, and J. Liang, "Active and semi-supervised graph neural networks for graph classification," *IEEE Trans. Big Data*, vol. 8, no. 4, pp. 920–932, Aug. 2022.
- [3] W. Ju et al., "A comprehensive survey on deep graph representation learning," 2023, *arXiv:2304.05055*.
- [4] Y. Xie, Y. Liang, M. Gong, A. Qin, Y.-S. Ong, and T. He, "Semisupervised graph neural networks for graph classification," in *IEEE Trans. Cybern.*, to be published, doi: [10.1109/TCYB.2022.3164696](https://doi.org/10.1109/TCYB.2022.3164696).
- [5] J. Li, Y. Huang, H. Chang, and Y. Rong, "Semi-supervised hierarchical graph classification," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2022.
- [6] W. Ju et al., "Tgnn: A joint semi-supervised framework for graph-level classification," 2023, *arXiv:2304.11688*.
- [7] X. Luo, Y. Zhao, Y. Qin, W. Ju, and M. Zhang, "Towards semi-supervised universal graph classification," *IEEE Trans. Knowl. Data Eng.*, vol. 45, no. 5, pp. 6265–6276, May 2023.
- [8] Z. Hao et al., "ASGN: An active semi-supervised graph neural network for molecular property prediction," in *Proc. 26th ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining*, 2020, pp. 731–752.
- [9] R. Kojima, S. Ishida, M. Ohta, H. Iwata, T. Honma, and Y. Okuno, "kGCN: A graph-based deep learning framework for chemical structures," *J. Cheminformatics*, vol. 12, pp. 1–10, 2020.
- [10] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Proc. Adv. Neural Inf. Process. Syst.*, 2017, pp. 1024–1034.
- [11] Y. Zhang, X. Wang, C. Shi, X. Jiang, and Y. Ye, "Hyperbolic graph attention network," *IEEE Trans. Big Data*, vol. 8, no. 6, pp. 1690–1701, Dec. 2022.
- [12] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, "Neural message passing for quantum chemistry," in *Proc. Int. Conf. Mach. Learn.*, PMLR, 2017, pp. 1263–1272.
- [13] J. Lee, I. Lee, and J. Kang, "Self-attention graph pooling," in *Proc. Int. Conf. Mach. Learn.*, PMLR, 2019, pp. 3734–3743.
- [14] Z. Ying, J. You, C. Morris, X. Ren, W. Hamilton, and J. Leskovec, "Hierarchical graph representation learning with differentiable pooling," in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, pp. 4805–4815.
- [15] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic minority over-sampling technique," *J. Artif. Intell. Res.*, vol. 16, pp. 321–357, 2002.
- [16] H. He and E. A. Garcia, "Learning from imbalanced data," *IEEE Trans. Knowl. Data Eng.*, vol. 21, no. 9, pp. 1263–1284, Sep. 2009.
- [17] M. Ren, W. Zeng, B. Yang, and R. Urtasun, "Learning to reweight examples for robust deep learning," in *Proc. Int. Conf. Mach. Learn.*, PMLR, 2018, pp. 4334–4343.
- [18] X. Wang, L. Lian, Z. Miao, Z. Liu, and S. X. Yu, "Long-tailed recognition by routing diverse distribution-aware experts," 2020, *arXiv: 2010.01809*.
- [19] J. Cui, Z. Zhong, S. Liu, B. Yu, and J. Jia, "Parametric contrastive learning," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2021, pp. 715–724.
- [20] C. K. Maurya, D. Toshniwal, and G. V. Venkatarao, "Distributed sparse class-imbalance learning and its applications," *IEEE Trans. Big Data*, vol. 7, no. 5, pp. 832–844, Nov. 2021.
- [21] H. Guo and S. Wang, "Long-tailed multi-label visual recognition by collaborative training on uniform and re-balanced samplings," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 15089–15098.
- [22] K. Cao, C. Wei, A. Gaidon, N. Arechiga, and T. Ma, "Learning imbalanced datasets with label-distribution-aware margin loss," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, pp. 1565–1576.
- [23] J. Tan et al., "Equalization loss for long-tailed object recognition," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 11662–11671.
- [24] B. Kang et al., "Decoupling representation and classifier for long-tailed recognition," 2019, *arXiv: 1910.09217*.

- [25] S. Alshammari, Y.-X. Wang, D. Ramanan, and S. Kong, "Long-tailed recognition via weight balancing," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2022, pp. 6897–6907.
- [26] B. Li, Z. Han, H. Li, H. Fu, and C. Zhang, "Trustworthy long-tailed classification," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2022, pp. 6970–6979.
- [27] J. Park, J. Song, and E. Yang, "GraphENS: Neighbor-aware ego network synthesis for class-imbalanced node classification," in *Proc. Int. Conf. Learn. Representations*, 2021.
- [28] N. Shervashidze, S. Vishwanathan, T. Petri, K. Mehlhorn, and K. Borgwardt, "Efficient graphlet kernels for large graph comparison," in *Proc. Int. Conf. Artif. Intell. Statist.*, 2009, pp. 488–495.
- [29] N. Shervashidze, P. Schweitzer, E. J. Van Leeuwen, K. Mehlhorn, and K. M. Borgwardt, "Weisfeiler-Lehman graph kernels," *J. Mach. Learn. Res.*, vol. 12, no. 9, pp. 2539–2561, 2011.
- [30] H. Kashima, K. Tsuda, and A. Inokuchi, "Marginalized kernels between labeled graphs," in *Proc. Int. Conf. Mach. Learn.*, 2003, pp. 321–328.
- [31] F. Chen and G. Long, "FedGE: Break the scalability limitation of graph neural network with federated graph embedding," in *IEEE Trans. Big Data*, to be published, doi: [10.1109/TBDATA.2022.3216747](https://doi.org/10.1109/TBDATA.2022.3216747).
- [32] W. Ju, X. Luo, Z. Ma, J. Yang, M. Deng, and M. Zhang, "GHNN: Graph harmonic neural networks for semi-supervised graph-level classification," *Neural Netw.*, vol. 151, pp. 70–79, 2022.
- [33] W. Ju, J. Yang, M. Qu, W. Song, J. Shen, and M. Zhang, "KGNN: Harnessing kernel-based networks for semi-supervised graph classification," in *Proc. 15th ACM Int. Conf. Web Search Data Mining*, 2022, pp. 421–429.
- [34] M. Buda, A. Maki, and M. A. Mazurowski, "A systematic study of the class imbalance problem in convolutional neural networks," *Neural Netw.*, vol. 106, pp. 249–259, 2018.
- [35] T. Zhao, X. Zhang, and S. Wang, "GraphSMOTE: Imbalanced node classification on graphs with graph neural networks," in *Proc. 14th ACM Int. Conf. Web Search Data Mining*, 2021, pp. 833–841.
- [36] L. Qu, H. Zhu, R. Zheng, Y. Shi, and H. Yin, "ImGAGN: Imbalanced network embedding via generative adversarial graph networks," in *Proc. 27th ACM SIGKDD Conf. Knowl. Discov. Data Mining*, 2021, pp. 1390–1398.
- [37] S.-Y. Yi and Y.-D. Zhou, "Model-free global likelihood subsampling for massive data," *Statist. Comput.*, vol. 33, no. 1, 2023, Art. no. 9.
- [38] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Longie, "Class-balanced loss based on effective number of samples," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 9268–9277.
- [39] C. Huang, Y. Li, C. C. Loy, and X. Tang, "Deep imbalanced learning for face recognition and attribute prediction," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 11, pp. 2781–2794, Nov. 2020.
- [40] Z. Liu, T.-K. Nguyen, and Y. Fang, "Tail-GNN: Tail-node graph neural networks," in *Proc. 27th ACM SIGKDD Conf. Knowl. Discov. Data Mining*, 2021, pp. 1109–1119.
- [41] S. Zhang, Z. Li, S. Yan, X. He, and J. Sun, "Distribution alignment: A unified framework for long-tail visual recognition," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 2361–2370.
- [42] L. Xiang, G. Ding, and J. Han, "Learning from multiple experts: Self-paced knowledge distillation for long-tailed classification," in *Proc. 16th Eur. Conf. Comput. Vis.*, Glasgow, U.K., Springer, 2020, pp. 247–263.
- [43] J. Cai, Y. Wang, and J.-N. Hwang, "ACE: Ally complementary experts for solving long-tailed recognition in one-shot," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2021, pp. 112–121.
- [44] F. Hu, W. Liping, L. Qiang, S. Wu, L. Wang, and T. Tan, "GraphDIVE: Graph classification by mixture of diverse experts," in *Proc. 31st Int. Joint Conf. Artif. Intell.*, 2022, pp. 2080–2086.
- [45] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *Proc. Int. Conf. Mach. Learn.*, PMLR, 2020, pp. 1597–1607.
- [46] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 9729–9738.
- [47] X. Chen, H. Fan, R. Girshick, and K. He, "Improved baselines with momentum contrastive learning," 2020, *arXiv: 2003.04297*.
- [48] P. Khosla et al., "Supervised contrastive learning," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 33, pp. 18661–18673, 2020.
- [49] Y. Ren, J. Bai, and J. Zhang, "Label contrastive coding based graph neural network for graph classification," in *Proc. Database Syst. Adv. Appl.: 26th Int. Conf.*, Taipei, Taiwan, Springer, 2021, pp. 123–140.
- [50] X. Deng, D. Huang, D.-H. Chen, C.-D. Wang, and J.-H. Lai, "Strongly augmented contrastive clustering," *Pattern Recognit.*, vol. 139, 2023, Art. no. 109470.
- [51] Y. You, T. Chen, Y. Sui, T. Chen, Z. Wang, and Y. Shen, "Graph contrastive learning with augmentations," *Adv. Neural Inf. Process. Syst.*, vol. 33, pp. 5812–5823, 2020.
- [52] L. Zhong, J. Yang, Z. Chen, and S. Wang, "Contrastive graph convolutional networks with generative adjacency matrix," *IEEE Trans. Signal Process.*, vol. 71, pp. 772–785, 2023.
- [53] W. Ju et al., "Unsupervised graph-level representation learning with hierarchical contrasts," *Neural Netw.*, vol. 158, pp. 359–368, 2023.
- [54] X. Luo et al., "Clear: Cluster-enhanced contrast for self-supervised graph representation learning," in *IEEE Trans. Neural Netw. Learn. Syst.*, to be published, doi: [10.1109/TNNLS.2022.3177775](https://doi.org/10.1109/TNNLS.2022.3177775).
- [55] X. Luo et al., "DualGraph: Improving semi-supervised graph classification via dual contrastive learning," in *Proc. IEEE 38th Int. Conf. Data Eng.*, 2022, pp. 699–712.
- [56] W. Ju et al., "GLCC: A general framework for graph-level clustering," in *Proc. AAAI Conf. Artif. Intell.*, 2023, pp. 4391–4399.
- [57] J. Yuan, X. Luo, Y. Qin, Y. Zhao, W. Ju, and M. Zhang, "Learning on graphs under label noise," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process.*, 2023, pp. 1–5.
- [58] D. P. Bertsekas, *Constrained Optimization and Lagrange Multiplier Methods*. New York, NY, USA: Academic Press, 2014.
- [59] J. L. Leevy, T. M. Khoshgoftaar, R. A. Bauder, and N. Seliya, "A survey on addressing high-class imbalance in Big Data," *J. Big Data*, vol. 5, no. 1, pp. 1–30, 2018.
- [60] J. Li, Z. Tan, J. Wan, Z. Lei, and G. Guo, "Nested collaborative learning for long-tailed visual recognition," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2022, pp. 6949–6958.
- [61] Z. Tan, J. Li, J. Du, J. Wan, Z. Lei, and G. Guo, "NCL++: Nested collaborative learning for long-tailed visual recognition," 2023, *arXiv:2306.16709*.
- [62] P. Yanardag and S. Vishwanathan, "Deep graph kernels," in *Proc. 21th ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining*, 2015, pp. 1365–1374.
- [63] C. Morris, N. M. Kriege, K. Kersting, and P. Mutzel, "Faster kernels for graphs with continuous attributes via hashing," in *Proc. IEEE 16th Int. Conf. Data Mining*, 2016, pp. 1095–1100.
- [64] I. Schomburg et al., "Brenda, the enzyme database: Updates and major new developments," *Nucleic Acids Res.*, vol. 32, no. suppl\_1, pp. D431–D433, 2004.
- [65] V. P. Dwivedi, C. K. Joshi, T. Laurent, Y. Bengio, and X. Bresson, "Benchmarking graph neural networks, 2020," *arXiv:2003.00982*.
- [66] K. Riesen and H. Bunke, "Iam graph database repository for graph based pattern recognition and machine learning," in *Proc. Joint IAPR Int. Workshops Statist. Techn. Pattern Recognit. Struct. Syntactic Pattern Recognit.*, Springer, 2008, pp. 287–297.
- [67] W. J. Reed, "The pareto, zipf and other power laws," *Econ. Lett.*, vol. 74, no. 1, pp. 15–19, 2001.
- [68] N. V. Chawla, "C4. 5 and imbalanced data sets: Investigating the effect of sampling method, probabilistic estimate, and decision tree structure," in *Proc. Int. Conf. Mach. Learn.*, CIBC Toronto, ON, Canada, 2003, Art. no. 66.
- [69] A. K. Menon, S. Jayasumana, A. S. Rawat, H. Jain, A. Veit, and S. Kumar, "Long-tail learning via logit adjustment," 2020, *arXiv: 2007.07314*.
- [70] S. Yu, H. Huang, M. N. Dao, and F. Xia, "Graph augmentation learning," in *Proc. Companion Proc. Web Conf.*, 2022, pp. 1063–1072.
- [71] Y. Wang, Y. Zhao, N. Shah, and T. Derr, "Imbalanced graph classification via graph-of-graph neural networks," in *Proc. 31st ACM Int. Conf. Inf. Knowl. Manage.*, 2022, pp. 2067–2076.
- [72] C. Hou, J. Zhang, H. Wang, and T. Zhou, "Subclass-balancing contrastive learning for long-tailed recognition," 2023, *arXiv:2306.15925*.



**Si-Yu Yi** received the BS and MS degrees in mathematics from Sichuan University, Sichuan, China, in 2017 and 2020, respectively. She is currently working toward the PhD degree in statistics with Nankai University, Tianjin, China. Her research interests focus on graph representation learning, design of experiments, statistical sampling, and subsampling in Big Data.



**Zhengyang Mao** is currently working toward the master's degree with the School of Computer Science, Peking University. His research interests include graph representation learning and long-tailed learning.



**Luchen Liu** received the PhD degree in computer science from Peking University, in 2020. He is currently a post-doctoral research fellow in computer science with Peking University. His current research interests lie primarily in the area of deep learning for temporal graph data and interdisciplinary applications such as intelligent healthcare and quantitative investment.



top-tier venues and has won the best paper finalist in IEEE ICDM 2022.

**Wei Ju** (Member, IEEE) received the BS degree in mathematics from Sichuan University, Sichuan, China, in 2017, and the PhD degree in computer science from Peking University, Beijing, China, in 2022. He is currently a postdoc research fellow in computer science with Peking University. His current research interests lie primarily in the area of machine learning on graphs including graph representation learning and graph neural networks, and interdisciplinary applications such as drug discovery and recommender systems. He has published more than 20 papers in



**Xiao Luo** received the BS degree in mathematics from Nanjing University, Nanjing, China, in 2017, and the PhD degree from the School of Mathematical Sciences, Peking University, Beijing, China. He is a postdoctoral researcher with the Department of Computer Science, University of California, Los Angeles, USA. His research interests include machine learning on graphs, image retrieval, statistical models and bioinformatics.



**Yong-Dao Zhou** received the BS degree in mathematics, the MS and PhD degrees in statistics from Sichuan University, China, in 2002, 2005, and 2008, respectively. After graduation, he joined Sichuan University and was a professor after 2015. In 2017, he then joined Nankai University, where he is presently a professor in statistics. His research agenda focuses on design of experiments and Big Data analysis. He published more than 60 papers and five monographs. His research publications have won best paper awards in WCE 2009 and Sci Sin Math in 2023.



**Ming Zhang** received the BS, MS, and PhD degrees in computer science from Peking University, respectively. She is a full professor with the School of Computer Science, Peking University. He is a member of Advisory Committee of Ministry of Education in China and the chair of ACM SIGCSE China. She is one of the fifteen members of ACM/IEEE CC2020 Steering Committee. She has published more than 200 research papers on Text Mining and Machine Learning in the top journals and conferences. She won the best paper of ICML 2014 and best paper nominee of WWW 2016. He is the leading author of several textbooks on Data Structures and Algorithms in Chinese, and the corresponding course is awarded as the National Elaborate Course, National Boutique Resource Sharing Course, National Fine-designed Online Course, National First-Class Undergraduate Course by MOE China.