# CLNode: Curriculum Learning for Node Classification

Xiaowen Wei
School of Computer Science
Wuhan University
Wuhan, China
weixiaowen@whu.edu.cn

Xiuwen Gong
Faculty of Engineering
The University of Sydney
Sydney, Australia
xiuwen.gong@sydney.edu.au

Yibing Zhan
JD Explore Academy
Beijing, China
zhanyibing@jd.com

Bo Du
School of Computer Science
Wuhan University
Wuhan, China
gunspace@163.com

Yong Luo
School of Computer Science
Wuhan University
Wuhan, China
luoyong@whu.edu.cn

Wenbin Hu*
School of Computer Science
Wuhan University
Wuhan, China
hwb@whu.edu.cn

## ABSTRACT

Node classification is a fundamental graph-based task that aims to predict the classes of unlabeled nodes, for which Graph Neural Networks (GNNs) are the state-of-the-art methods. Current GNNs assume that nodes in the training set contribute equally during training. However, the quality of training nodes varies greatly, and the performance of GNNs could be harmed by two types of low-quality training nodes: (1) inter-class nodes situated near class boundaries that lack the typical characteristics of their corresponding classes. Because GNNs are data-driven approaches, training on these nodes could degrade the accuracy. (2) mislabeled nodes. In real-world graphs, nodes are often mislabeled, which can significantly degrade the robustness of GNNs. To mitigate the detrimental effect of the low-quality training nodes, we present CLNode, which employs a selective training strategy to train GNN based on the quality of nodes. Specifically, we first design a multi-perspective difficulty measurer to accurately measure the quality of training nodes. Then, based on the measured qualities, we employ a training scheduler that selects appropriate training nodes to train GNN in each epoch. To evaluate the effectiveness of CLNode, we conduct extensive experiments by incorporating it in six representative backbone GNNs. Experimental results on real-world networks demonstrate that CLNode is a general framework that can be combined with various GNNs to improve their accuracy and robustness.

## CCS CONCEPTS

• **Theory of computation** → **Graph algorithms analysis**; • **Information systems** → *Social networks*; • **Computing methodologies** → Neural networks.

*Corresponding author.

## KEYWORDS

node classification; curriculum learning; graph neural networks

## 1 INTRODUCTION

Node classification is a fundamental graph-based task. Given a graph with limited labeled nodes (training nodes), the task aims to assign labels to unlabeled nodes [23]. The state-of-the-art node classification methods are Graph Neural Networks (GNNs) [35, 40]. Generally, GNNs update the node representations by aggregating the messages passed from their neighbors. Benefiting from this aggregation mechanism, GNNs learn low-dimensional node representations that preserve the topological information and node feature attributes, which are then used to predict the labels. Although many GNN-based node classification works [3, 10, 16, 22, 34] have been proposed, these works usually assume that all training nodes contribute equally. In fact, the quality of training nodes varies widely. Being data-driven approaches, GNNs exhibit degraded performance by training on the low-quality nodes.

To illustrate the quality of nodes, we define training nodes whose representations lack the typical characteristics of their label classes as *difficult nodes*, because it is difficult for GNNs to learn class characteristics from these low-quality nodes. In contrast, *easy nodes* refer to high-quality nodes that have the typical representations of their label classes. We illustrate *difficult nodes* and *easy nodes* using the paper citation network in Figure 1. As illustrated, the cross-field paper $v_1$ connects papers from multiple classes. During neighborhood aggregation, $v_1$ aggregates messages from neighbors $\{v_2, v_3, v_4, v_5, v_6\}$. By aggregating messages $\{v_4 \rightarrow v_1, v_5 \rightarrow v_1, v_6 \rightarrow v_1\}$ from classes $\{c_1, c_2, c_4\}$, $v_1$ obtains an unclear representation that mixes characteristics of different classes, indicating that $v_1$ is a *difficult node*. In contrast, all the aggregated messages of $v_{15}$ are from class $c_4$, which makes it an *easy node*. Therefore, the above observation raises the question of whether these uneven-quality training nodes should be treated equally by GNNs.
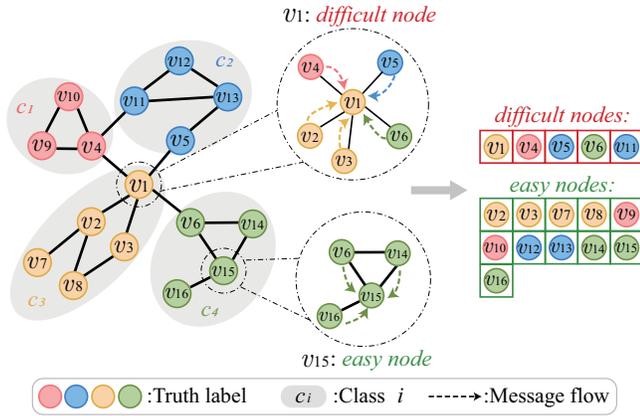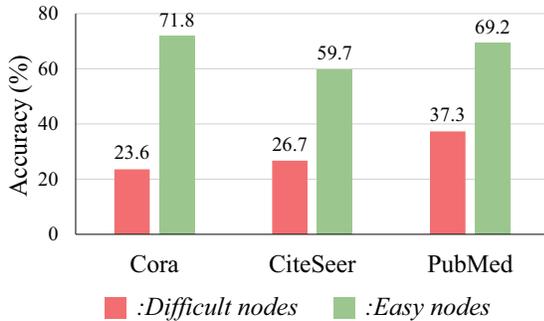
**Figure 1: Illustration of node difficulty.**



**Figure 2: Accuracy of GCN trained on *difficult nodes* or *easy nodes*.**

*Easy nodes* and *difficult nodes* play different roles during training. The representations of *easy nodes* are typical, and training on such nodes helps GNNs find clear decision boundaries; whereas, *difficult nodes* should be used carefully, as their representations lack the typical characteristics of their classes. There are two types of *difficult nodes* that degrade the performance of GNNs: (1) inter-class nodes situated near class boundaries. By aggregating messages from neighbors, these nodes obtain unclear representations; as a result, training on these nodes degrades the accuracy of GNNs. (2) mislabeled nodes. Real-world graphs often contain label noise[5, 17, 20] and current GNNs are easily perturbed by training on these mislabeled nodes. Figure 2 shows the accuracy of GCN [16] on three paper citation networks [23], where the same number of *difficult nodes* or *easy nodes* are utilized for training. The node difficulty is evaluated using Eq.(11), which we will detail in Section 4. From the results, we can see that training on *easy nodes* leads to higher accuracy. For example, on the Cora network, if all training nodes are *easy nodes*, the accuracy is 71.8%, and the accuracy is only 23.6% when only *difficult nodes* are utilized. Based on the above analysis, mitigating the detrimental effect of *difficult nodes* can improve the accuracy and robustness of GNNs. In this paper, we introduce curriculum learning [1] to mitigate the effect of these low-quality training nodes.

In particular, curriculum learning is a training strategy that initially trains the machine learning models using an easier training subset and then gradually introduces more difficult samples. By excluding low-quality difficult samples during initial training, curriculum learning mitigates overfitting to data noise, and thus improves models' accuracy and robustness [19, 30, 41]. The most critical component of curriculum learning is the difficulty measurer, which estimates the difficulty (quality) of samples. In existing works, difficulty measurers are often designed by observing the sample features; for example, sentence length is a popular difficulty measurer in NLP tasks because shorter sentences are often easier for models to learn [21]. However, difficulty cannot be measured directly from node features using a similar approach. One feasible way is to utilize the graph structure, e.g., if a node connects neighbors from multiple classes, it is likely to be an inter-class *difficult node*. However, in the node classification task, this is challenging due to the limited node labels.

In this paper, we attempt to address the above challenging problem by proposing a **C**urriculum **L**earning framework for **Node** Classification, called CLNode. The key idea behind CLNode is to enhance the performance of backbone GNN by incrementally introducing nodes into the training process, starting with *easy nodes* and progressing to harder ones. Specifically, we first propose to assign pseudo-labels to unlabeled nodes. With the help of label information, we design a multi-perspective difficulty measurer, in which two difficulty measurers from local and global perspectives are proposed to measure the difficulty of training nodes. The local difficulty measurer computes local label distribution to identify inter-class *difficult nodes* because their neighbors have diverse labels; the global difficulty measurer identifies mislabeled *difficult nodes* by analyzing the node feature. Based on the measured node difficulty, we propose a continuous training scheduler that selects appropriate training nodes in each epoch to mitigate the negative effect of *difficult nodes*. CLNode is a general framework that can be combined with various GNNs to improve their node classification performance. The key contributions of this paper are summarized as follows:

- We propose CLNode, a novel curriculum learning framework for node classification. CLNode first accurately identifies two types of *difficult nodes*, and then employs a selective training strategy to mitigate the detrimental effect of these nodes.
- We demonstrate that CLNode can be directly plugged into existing GNNs. Without increasing the time complexity, CLNode enhances backbone GNNs by simply feeding nodes to the training process in order from easy to difficult.
- We conduct extensive experiments on five datasets. The results demonstrate that compared with baseline methods without curriculum learning, CLNode effectively improves the accuracies and enhances the robustness to label noise.

## 2 RELATED WORK

### 2.1 Node Classification and GNNs

Node classification [23] aims to predict labels for unlabeled nodes in a given graph. As a fundamental task on graphs, node classification has various applications, including fraud detection [7, 12, 39],

security and privacy analytics [27], and community detection [11, 14].

Recently, GNNs have emerged as promising approaches for analyzing graph data. Due to the long history of GNNs, we refer readers to [35, 40] for a comprehensive review. Based on the definition of graph convolution, GNNs can be broadly divided into two categories, namely spectral-based [2, 16, 26] and spatial-based [10, 37]. Bruna et al. [2] first explore spectral-based GNNs by utilizing a spectral filter on the spectral space. In a follow-up work, GCN [16] simplifies the graph convolution operation. SGC[34] proposes to remove the nonlinearity in GCN and thereby speed up the model. Different from spectral-based methods, spatial-based methods define convolutions directly on graphs by performing operations on spatially close neighbors. GraphSAGE [10] is a general inductive framework that generates representations for nodes by sampling local neighbors. JK-Net [37] devises an alternative graph structure-based strategy to select neighbors for nodes. Although GNNs have achieved great success, they simply assume all training nodes to make equal contributions; consequently, training on the low-quality *difficult nodes* significantly degrades their accuracy and robustness.

## 2.2 Curriculum Learning

Inspired by the learning principle underlying human cognitive processes, curriculum learning [1] is proposed as a training strategy that trains machine learning models from easier samples to harder samples. Previous studies [1, 32, 33] have shown that curriculum learning improves generalization capacity and guides the model towards a better parameter space. Motivated by this, scholars have exploited the power of curriculum learning in a wide range of fields, including computer vision (CV) [9, 13, 38], natural language processing (NLP) [6, 28, 29] and graph classification [31], etc. To the best of our knowledge, however, no work has yet attempted to apply curriculum learning to node classification.

## 3 PRELIMINARIES

### 3.1 Notation

Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, X)$ denote a graph, where $\mathcal{V}$ is the node set, $\mathcal{E}$ is the edge set, and $X$ is the node feature matrix. The input feature of node $i$ is $x_i$, and the neighborhood of node $i$ is $\mathcal{N}(i)$. For the node classification task, a labeled node set $\mathcal{V}_L = \{v_1, ..., v_l\}$ is given with $Y_L$ denoting the input labels. $C$ is the set of classes. The goal of node classification is to predict the labels of unlabeled nodes in the graph.

### 3.2 Graph Neural Networks

Generally, a GNN involves two key computations for each node $i$ at every layer: (1) neighborhood aggregation: aggregating messages passed from $\mathcal{N}(i)$. (2) update representation: updating $i$'s representation from its representation in the previous layer and the aggregated messages. Formally, the $l$-th layer representation of node $i$ is given by:

$$h_i^l = \text{UPDATE}(h_i^{l-1}, \text{AGGREGATE}(\{h_j^{l-1}|j \in \mathcal{N}(i)\})). \quad (1)$$

The final node representation $h_i^L$, i.e., the output of the last layer, is used for various downstream tasks. For the node classification

task, after obtaining node representations, a multilayer perceptron is often used to map them to the predicted labels.

## 3.3 Curriculum Learning

Curriculum learning mitigates the detrimental effect of low-quality samples by using a curriculum to train the model. A curriculum is a sequence of training criteria $< Q_1, ..., Q_t, ..., Q_T >$ over $T$ training epochs. Each criterion $Q_t$ is a training subset. The initial $Q_1$ consists of easier samples; as $t$ increases, more difficult samples are gradually introduced into $Q_t$. In essence, designing such a curriculum for node classification requires us to design a **difficulty measurer** and a **training scheduler**. Here, the difficulty measurer estimates the difficulty of each training node; subsequently, based on the difficulty, the training scheduler generates $Q_t$ at any training epoch $t$ to train the model.

## 4 METHODOLOGY

In this section, we present the details of CLNode. As shown in Figure 3, CLNode comprises two components: (i) multi-perspective difficulty measurer (Figure 3(a)). We first perform a standard node classification to obtain additional label information, then two difficulty measurers from local and global perspectives are proposed to measure the node difficulty. (ii) continuous training scheduler (Figure 3(b)). After determining the node difficulty, we design a training scheduler to train backbone GNN with *easy nodes* initially and continuously introduce harder training nodes. By paying less attention to *difficult nodes*, CLNode improves the accuracy and robustness of backbone GNN. We detail the components of CLNode in the following subsections.

### 4.1 Multi-perspective Difficulty Measurer

In general, neighborhood aggregation benefits from the homophily of graphs, i.e., a node $i$'s neighbors $\mathcal{N}(i)$ tend to have the same label as $i$. However, the *difficult nodes* violate the homophily; for example, the neighbors of an inter-class *difficult node* have diverse labels because they belong to multiple classes. Taking a step further, the difficulty of nodes can be measured with the help of label information. Therefore, the first step is to assign pseudo-labels to unlabeled nodes (see Figure 3(a)). Specifically, we first train a GNN $f_1$ on the whole training set $\mathcal{V}_L$ to perform a standard node classification. After the training process, $f_1$ is used to get the pseudo-labels:

$$H = f_1(\mathcal{G}), \quad (2)$$

$$Y_P = MLP(H), \quad (3)$$

where $H$ is the node representation matrix obtained by GNN $f_1$ and $Y_P$ is the pseudo-labels predicted by a multilayer perceptron. However, directly using $Y_P$ to measure node difficulty may lead to inaccurate results, since $Y_P$ of training nodes may be different from the input labels $Y_L$. Therefore, to better measure node difficulty, we retain the input labels for training nodes:

$$\tilde{Y}[i] = \begin{cases} Y_L[i], & i \in \mathcal{V}_L \\ Y_P[i], & otherwise. \end{cases} \quad (4)$$
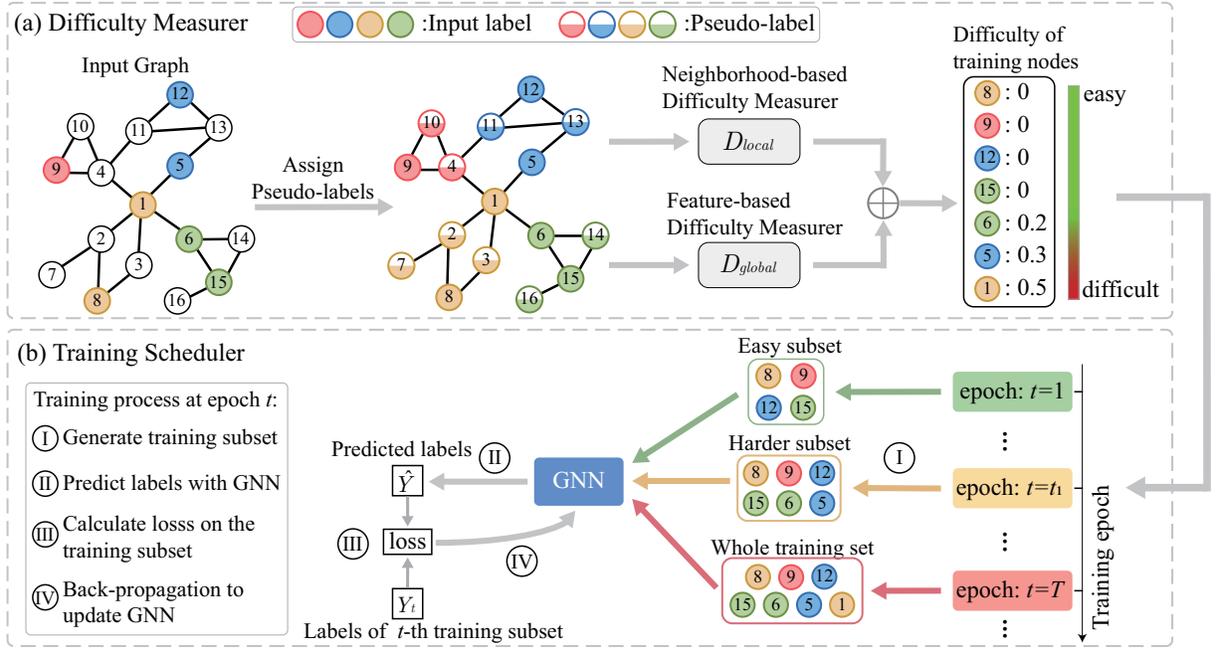
**Figure 3: An overall framework of the proposed CLNode.**

Subsequently, to identify two types of *difficult nodes*, i.e., inter-class nodes and mislabeled nodes, we propose two difficulty measurers to capture both local and global information for measuring the node difficulty.

### 4.1.1 Neighborhood-based Difficulty Measurer.

We first introduce how to identify *difficult nodes* from a local perspective. After obtaining $\tilde{Y}$, for each training node $u$, we calculate its difficulty with reference to the label distribution of its neighborhood. The first type of *difficult nodes* (inter-class nodes) have diverse neighbors that belong to multiple classes. In order to identify these inter-class *difficult nodes*, we calculate the diversity of neighborhood's labels:

$$P_c(u) = \frac{|\{\tilde{Y}[v] = c \mid v \in \hat{\mathcal{N}}(u)\}|}{|\hat{\mathcal{N}}(u)|}, \tag{5}$$

$$D_{local}(u) = -\sum_{c \in C} P_c(u) \, log(P_c(u)), \tag{6}$$

where $\hat{\mathcal{N}}(u)$ denotes $\mathcal{N}(u) \cup \{u\}$ and $P_c(u)$ denotes the proportion of the neighborhood $\hat{\mathcal{N}}(u)$ belonging to class $c$. A larger $D_{local}$ indicates a more diverse neighborhood. Taking Figure 3(a) as an example, the $D_{local}$ of node 1 is 0.54, which is much larger than $D_{local}(8) = 0$, indicating that node 1 has more diverse neighbors than node 8. Nodes with larger $D_{local}$ are more likely to be inter-class nodes. As a result, during neighborhood aggregation, these nodes aggregate neighbors' features to get an unclear representation, making them difficult for GNNs to learn. By paying less attention to these *difficult nodes*, CLNode learns more useful information and effectively improves the accuracy of backbone GNNs.

### 4.1.2 Feature-based Difficulty Measurer.

Because the pseudo-labels could be inaccurate, mislabeled training nodes may not be identified using local information. For instance, consider the training node 7 in Figure 4, whose truth label is $c_3$ but is mislabeled as $c_1$. The label information of node 7 will affect the pseudo-labels of its neighbors. As a result, the pseudo-label of node 2 is likely to be predicted as the mislabeled class $c_1$, thus the local label distribution of node 7 is consistent, from which we cannot identify it as a mislabeled node. Therefore, we propose to use global feature information to identify mislabeled nodes.

Nodes of the same class have similar features, e.g., in a paper citation network, papers in the same field tend to contain the same keywords. However, the mislabeled nodes violate this principle. For instance, in Figure 4, the mislabeled node 7 has low feature similarity to many nodes of its label class (e.g., node 10), since they do not in fact belong to the same class. Conversely, node 7 has high feature similarity to nodes in class $c_3$ (e.g., node 8). Therefore, by exploring the feature similarity, we can deduce that node 7 is likely to be mislabeled. The input feature $X$ is sparse in high-dimensional space, instead, we use $H$ (see Eq.(2)) as the node feature to compute similarity. Let $h_v$ denote the feature of node $v$, then the representative feature of class $c$ is defined as the average of the node features in class $c$:

$$\mathcal{V}_c = \{v \mid \tilde{Y}[v] = c\}, \tag{7}$$

$$h_c = \text{Avg}(h_v \mid v \in \mathcal{V}_c), \tag{8}$$

where $\mathcal{V}_c$ denotes the nodes belonging to class $c$, and $h_c$ is the representative feature of class $c$. To identify mislabeled *difficult nodes*, for each training node $u$, we compute its feature similarity to the label class:

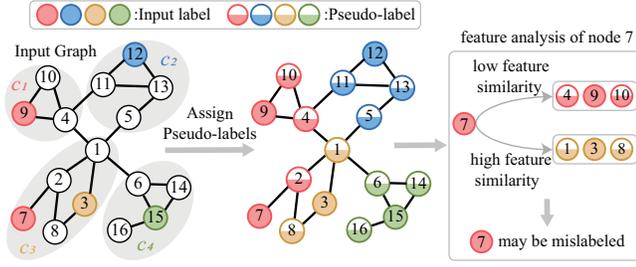$$S(u) = \frac{exp(h_u \cdot h_{c_u})}{\max_{c \in C} exp(h_u \cdot h_c)}, \tag{9}$$

**Figure 4: Illustration of the feature-based difficulty measurer.**



**Figure 5: Visualization of three pacing functions.**

where $c_u$ denotes the label class of node $u$, $S(u)$ calculates the feature similarity between $h_u$ and $h_{c_u}$. Mislabeled nodes tend to have smaller $S(u)$ than correctly labeled nodes. Based on $S(u)$, the feature-based difficulty measurer is defined as:

$$D_{global}(u) = 1 - S(u). \tag{10}$$

$D_{global}$ measures node difficulty from a global perspective. By using $D_{global}$ to identify mislabeled training nodes, CLNode selectively excludes these nodes from the training process, thus improving the robustness of the backbone GNNs to label noise. Considering two difficulty measurers from local and global perspectives, we finally define the difficulty of $u$ as:

$$D(u) = D_{local}(u) + \alpha \cdot D_{global}(u), \tag{11}$$

where $\alpha$ is a hyper-parameter that controls the weight of $D_{global}(u)$.

## 4.2 Continuous Training Scheduler

After measuring the node difficulty, we use a curriculum-based training strategy to train a better GNN model (see Figure 3(b)). To distinguish it from $f_1$, we denote the model trained with curriculum as $f_2$. We propose a continuous training scheduler to generate the easy-to-difficult curriculum. In more detail, we first sort the training set $\mathcal{V}_L$ in ascending order of node difficulty; subsequently, a pacing function $g(t)$ is used to map each training epoch $t$ to a scalar $\lambda_t$ whose range is $(0, 1]$, meaning that a proportion $\lambda_t$ of the easiest training nodes are used as the training subset at the $t$-th epoch. Let $\lambda_0$ denote the initial proportion of the available easiest nodes, while $T$ denotes the epoch when $g(t)$ reaches 1 for the first time. We consider three pacing functions, namely linear, root, and geometric:

- linear:

$$g(t) = min(1, \lambda_0 + (1 - \lambda_0) * \frac{t}{T}). \tag{12}$$

- root:

$$g(t) = min(1, \sqrt{\lambda_0^2 + (1 - \lambda_0^2) * \frac{t}{T}}). \tag{13}$$

- geometric:

$$g(t) = min(1, 2^{log_2 \lambda_0 - log_2 \lambda_0 * \frac{t}{T}}). \tag{14}$$

The visualization of these three pacing functions is presented in Figure 5. As shown in the figure, the linear function increases the difficulty of training nodes at a uniform rate; the root function
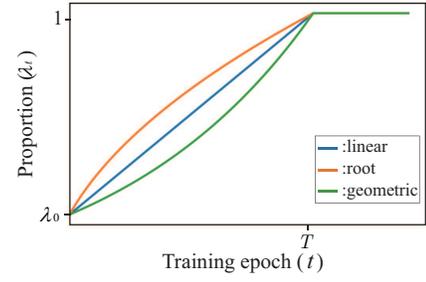
introduces more *difficult nodes* in fewer epochs, while the geometric function trains for a greater number of epochs on the subset of *easy nodes*. By using the pacing function to continuously introduce training nodes into the training process, CLNode assigns appropriate training weights to nodes of different levels of difficulty. Specifically, the more difficult a training node is, the later it is introduced into the training process, meaning it has a smaller training weight.

Moreover, we do not stop training immediately when $t = T$, because at this time, the backbone GNN $f_2$ may not have fully explored the knowledge of nodes which have been recently introduced. Instead, when $t > T$, we use the whole training set to train $f_2$ until the test accuracy on validation set converges.

## 4.3 Pseudo-code and Complexity Analysis

In this subsection, we present the pseudo-code of CLNode and explore its time complexity. The process of CLNode is detailed in Algorithm 1. Lines 2–7 describe the the process of measuring node difficulty and lines 8–17 describe the process of training the backbone GNN $f_2$ with a curriculum. After the training process, $f_2$ is finally used for node classification (see line 18). As the pseudo-code shows, CLNode is easy to be plugged into any backbone GNN, as it only changes the training set in each training epoch.

For the convenience of complexity analysis, we consider GCN as the backbone. The time complexity of an $L$-layer GCN in one epoch is $O(L|\mathcal{E}|F + L|\mathcal{V}|F^2)$, where $F$ is the number of node feature attributes. We assume that GCN converges after $T_1$ epochs, thus its time complexity is $O(T_1 \cdot (L|\mathcal{E}|F + L|\mathcal{V}|F^2))$, which is also the time complexity of training $f_1$. Next, the time complexity of measuring node difficulty is $O(ld + l|C|F)$, where $d$ is the average node degree. The time complexity of sorting $\mathcal{V}_L$ is $O(l \cdot log\, l)$. Finally, we analyze the time complexity of training $f_2$. We first train $T$ epochs using the curriculum, after which we train $f_2$ with the whole $\mathcal{V}_L$ until convergence. The training of the first $T$ epochs can be seen as pre-training $f_2$ with high-quality training nodes. Therefore, $f_2$ will converge before $T + T_1$ epochs. Because $l < |\mathcal{V}| \ll |\mathcal{E}|$, the upper bound on the time complexity of CLNode is $O((2T_1 + T) \cdot (L|\mathcal{E}|F + L|\mathcal{V}|F^2))$. In our experiments, we observe that the running time of CLNode is about twice that of the baseline GNN.

## 5 EXPERIMENTS

In this section, we first evaluate the improvement in accuracy achieved by CLNode over various backbone GNNs. Further experiments are

---

**Algorithm 1:** CLNode

**Input:** A graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, X)$, the labeled node set $\mathcal{V}_L$,
the input labels $Y_L$, the backbone GNN model, the
hyper-parameters $\alpha, \lambda_0, T$.

**Output:** The predicted labels $\hat{Y}$.

1 Initialize parameters of two GNN models $f_1$ and $f_2$;

2 Train $f_1$ on $(\mathcal{G}, \mathcal{V}_L, Y_L)$;

3 Predict pseudo-labels $Y_P$ with $f_1$;

4 $\tilde{Y} \leftarrow$ Eq.(4);

5 **for** $u \in \mathcal{V}_L$ **do**

6 $\quad$ Calculate node difficulty $D(u) \leftarrow$ Eq.(11);

7 **end**

8 Sort $\mathcal{V}_L$ according to node difficulty in ascending order;

9 Let $t = 1$;

10 **while** $t < T$ *or not converge* **do**

11 $\quad$ $\lambda_t \leftarrow g(t)$;

12 $\quad$ Generate training subset $\mathcal{V}_t \leftarrow \mathcal{V}_L[1, ..., \lfloor \lambda_t \cdot l \rfloor]$ ;

13 $\quad$ Use $f_2$ to predict the labels $Y_t$;

14 $\quad$ Calculate loss $\mathcal{L}$ on $\{Y_t[v], Y_L[v] \mid v \in \mathcal{V}_t\}$;

15 $\quad$ Back-propagation on $f_2$ for minimizing $\mathcal{L}$;

16 $\quad$ $t \leftarrow t + 1$;

17 **end**

18 Predict $\hat{Y}$ with $f_2$;

---

conducted on graphs with label noise to demonstrate the robustness of CLNode. Subsequently, we conduct ablation studies to verify the effectiveness of components in CLNode. Finally, we discuss the parameter sensitivity to hyper-parameters.

We conduct experiments on five benchmark datasets: Cora, Citeseer, PubMed [23], Amazon Computers (A-Computers), and Amazon Photo (A-Photo) [24]. Cora, CiteSeer, and PubMed are paper citation networks while A-Computers and A-Photo are product co-purchase networks. Experiments are conducted on these datasets with random splits and standard splits. The random splits follow [25, 36] to randomly label a specific proportion of nodes as the training set, and the label rates are listed in Table 1; the standard splits follow [16, 26] in using 20 labeled nodes per class as the training set. In each dataset, we follow [26, 34] to use 500 nodes for validation and 1000 nodes for testing.

We use six popular GNNs as the backbone models, namely GCN [16], GraphSAGE [10], GAT [26], SuperGAT [15], JK-Net [37] and GCNII [3], which are representative of a broad range of GNNs. In more detail, GCN is a typical convolution-based GNN, GraphSAGE can be applied to inductive learning, GAT and SuperGAT use attention mechanism in neighborhood aggregation, while JK-Net and GCNII are deep GNNs. We use backbone GNNs without curriculum learning as baselines to explore the improvement achieved by CLNode. All models are implemented in PyTorch-geometric [8]. We use the Adam optimizer with a learning rate of 0.01 and the weight decay is $5 \times 10^{-4}$. The hidden unit is fixed at 16 in paper citation networks and 64 in product co-purchase networks. We apply two graph convolutional layers for GCN, GAT, GraphSage, and SuperGAT, 6 layers for JK-Net, and 64 layers for GCNII. To facilitate fair comparison, the backbone GNNs' parameters of CLNode

are the same as the baselines. For CLNode, $\alpha$ is fixed at 1 because we observe good performance at this value. We use the geometric pacing function by default. The hyper-parameter $\lambda_0$ is searched in the range of {0.25, 0.5, 0.75}, while the search space of $T$ is {50, 100, 150}. The code is available at https://github.com/wxwmd/CLNode.

## 5.1 Node Classification

In this subsection, node classification experiments are conducted on five datasets. For each baseline GNN, we compare its original accuracy with the accuracy of being plugged into the CLNode framework. We conduct each experiment for ten trials to report the average test accuracy and standard deviation.

Table 2 reports the experimental results under random splits. The results demonstrate that CLNode can be combined with six backbone GNNs and improve their accuracy on node classification. For example, on the Cora dataset, CLNode improves the test accuracy of backbone GNNs by 3.5% (GCN), 2.0% (GraphSAGE), 2.9% (GAT), 1.1% (SuperGAT), 2.8% (JK-Net), and 1.6% (GCNII). The results prove that CLNode effectively mitigates the detrimental effect of *difficult nodes*, thereby enabling more useful information to be learned from uneven-quality training nodes.
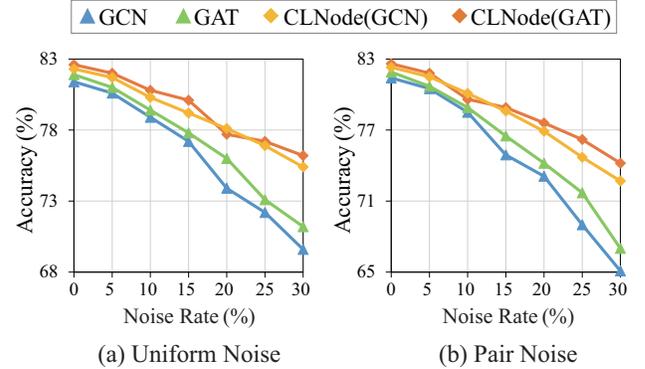


**Figure 6: Accuracy (%) on Cora with two kinds of label noise.**

Moreover, we conduct node classification experiments under different label rates. Table 3 shows the accuracy on Cora dataset at label rates of 1%, 2%, 3%, respectively. We observe that when there are fewer labeled training nodes, the improvement achieved by CLNode is more obvious. This is because when there are more training nodes, the detrimental effect of *difficult nodes* is mitigated by a large number of *easy nodes*; conversely, when there are fewer training nodes, *difficult nodes* easily mislead GNNs to learn the wrong knowledge. Therefore, by excluding *difficult nodes* from initial training, CLNode significantly improves the accuracy of GNNs at a low label rate. For many real-world graphs, the labeling process is tedious and costly, resulting in limited labels, and it would be highly beneficial to use CLNode in these situations.

## 5.2 Robustness to Noise

In this subsection, we investigate whether CLNode enhances the robustness of backbone GNNs to label noise. In a noisily labeled graph, the labels have a probability of $p$ to be flipped to other

**Table 1: Statistics of five benchmark datasets.**

| Dataset | Nodes | Edges | Features | Classes | Label rate |
|---|---|---|---|---|---|
| Cora | 2708 | 5429 | 1433 | 7 | 2% |
| CiteSeer | 3327 | 4732 | 3703 | 6 | 2% |
| PubMed | 19717 | 88648 | 500 | 3 | 0.1% |
| A-Computers | 13381 | 245778 | 767 | 10 | 1% |
| A-Photo | 7487 | 119043 | 745 | 8 | 1% |

**Table 2: Node classification performance (Accuracy (%)±Std) on five datasets.**

| | Method | Cora | CiteSeer | PubMed | A-Computers | A-Photo |
|---|---|---|---|---|---|---|
| GCN | Original | 73.5±0.8 | 62.8±2.6 | 64.3±2.9 | 79.0±3.7 | 89.1±0.8 |
| | +CLNode | **77.0±0.7** | **65.5±2.3** | **65.9±1.3** | **84.7±0.5** | **90.8±1.0** |
| | (Improv.) | 3.5% | 2.7% | 1.6% | 5.7% | 1.7% |
| GraphSAGE | Original | 70.1±2.3 | 57.4±3.7 | 61.3±1.4 | 71.7±2.4 | 83.0±2.6 |
| | +CLNode | **72.1±1.4** | **60.3±3.1** | **64.1±3.8** | **77.5±1.6** | **87.5±1.2** |
| | (Improv.) | 2.0% | 2.9% | 2.8% | 5.8% | 4.5% |
| GAT | Original | 74.2±1.2 | 63.7±2.8 | 64.6±2.5 | 80.2±0.8 | 89.4±1.8 |
| | +CLNode | **77.1±1.1** | **65.3±2.6** | **68.2±2.6** | **82.6±1.1** | **90.1±1.1** |
| | (Improv.) | 2.9% | 1.6% | 3.6% | 2.4% | 0.7% |
| SuperGAT | Original | 74.4±4.3 | **64.8±3.3** | 67.4±4.3 | 81.2±2.0 | 87.3±2.0 |
| | +CLNode | **75.5±2.7** | 63.0±3.2 | **72.2±3.0** | **83.4±2.4** | **88.8±1.2** |
| | (Improv.) | 1.1% | - | 4.8% | 2.2% | 1.5% |
| JK-Net | Original | 74.0±1.5 | 62.1±3.7 | **66.0±1.7** | 83.2±1.3 | 89.2±0.7 |
| | +CLNode | **76.8±0.8** | **63.6±1.2** | **71.5±3.2** | **84.4±1.0** | **90.4±0.9** |
| | (Improv.) | 2.8% | 1.5% | 5.5% | 1.2% | 1.2% |
| GCNII | Original | 76.2±4.0 | 64.5±4.3 | 70.8±6.1 | 79.8±1.8 | 87.4±2.1 |
| | +CLNode | **77.8±2.1** | **66.5±2.2** | **71.3±4.6** | **82.2±1.5** | **89.3±2.0** |
| | (Improv.) | 1.6% | 2.0% | 0.5% | 2.4% | 1.9% |

**Table 3: Accuracy (%) on Cora under different label rates.**

| | Method | 1% | 2% | 3% |
|---|---|---|---|---|
| GCN | Original | 62.4±2.7 | 73.5±0.8 | 78.6±0.6 |
| | +CLNode | **66.9±1.2** | **77.0±0.7** | **79.7±0.6** |
| GraphSage | Original | 54.8±3.0 | 70.1±2.3 | 76.0±0.8 |
| | +CLNode | **61.8±2.6** | **72.1±1.4** | **77.7±1.5** |
| GAT | Original | 65.2±2.4 | 74.2±1.2 | 78.8±1.0 |
| | +CLNode | **68.5±2.0** | **77.1±1.1** | **79.9±0.5** |
| SuperGAT | Original | 65.5±6.0 | 74.4±4.3 | **78.7±1.6** |
| | +CLNode | **67.9±3.2** | **75.5±2.7** | 78.5±2.4 |
| JK-Net | Original | 67.5±1.7 | 74.0±1.5 | 77.4±1.4 |
| | +CLNode | **69.4±1.4** | **76.8±0.8** | **78.8±0.3** |
| GCNII | Original | 68.5±3.9 | 76.2±4.0 | 79.0±2.2 |
| | +CLNode | **71.2±3.8** | **77.8±2.1** | **80.2±2.0** |

- Uniform noise. The label has a probability of $p$ to be mislabeled as any other class.
- Pair noise. We assume that nodes in one class can only be mislabeled as their closest class; that is, labels have a probability $p$ to flip to their pair class.

We conduct experiments on Cora under standard splits and vary $p$ from {0, 5%,..., 30%} to compare the performance of CLNode and the baseline GNNs under different levels of noise. We only report the results using GCN and GAT as backbone GNNs because we have similar observations for other GNNs. CLNode(GCN) and CLNode(GAT) denote the CLNode method using GCN and GAT as backbone GNN, respectively.

The results are shown in Figure 6, from which we observe that as the noise rate increases, the performance of all baselines drops dramatically. CLNode also suffers under conditions of increasing noise rate; however, when there is more noise in the graph, the performance gap between CLNode and the baseline increases. This observation demonstrates that CLNode effectively enhances the robustness of backbone GNNs to two kinds of label noise, since CLNode considers mislabeled training nodes as *difficult nodes* and selectively excludes them from the training process, while the baseline

classes, where $p$ denotes the noise rate. Following [5, 18], we corrupt the labels of the training and validation set with two kinds of label noise:

**Table 4: Comparisons between different difficulty measurers.**

|  | Method | Cora | CiteSeer | PubMed |
|---|---|---|---|---|
| GCN | original | 69.6 | 55.3 | 69.4 |
|  | +CLNode(local) | 74.8 | 61.8 | 74.2 |
|  | +CLNode(global) | 72.3 | 62.5 | 73.2 |
|  | +CLNode | **75.4** | **63.1** | **74.4** |

**Table 5: Comparisons between different pacing functions.**

|  | Pacing Function | Cora | CiteSeer | PubMed |
|---|---|---|---|---|
| CLNode | linear | 74.8 | 62.7 | 74.2 |
|  | root | 74.5 | 62.5 | 73.9 |
|  | geometric | **75.4** | **63.1** | **74.4** |

GNNs treat all training nodes as equal and consequently overfit to noise.
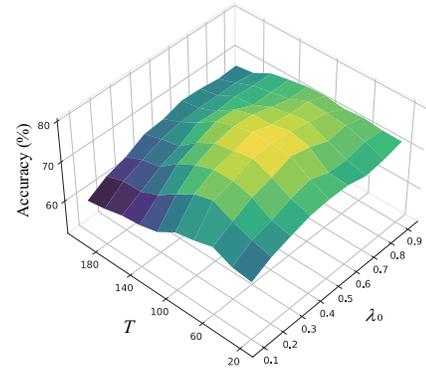
## 5.3 Ablation Study

In this subsection, we conduct ablation studies to explore the effectiveness of the multi-perspective difficulty measurer and the sensitivity of CLNode to different pacing functions. Ablation studies are conducted on three paper citation datasets under standard splits, where the graphs are corrupted by uniform label noise and the noise rate $p$ is set to 30%.

First, to verify the multi-perspective difficulty measurer benefits from combining the local and global information, we design two difficulty measurers to replace it:

- Measuring difficulty only with local information, i.e., we only use $D_{local}$ to measure node difficulty.
- Measuring difficulty only with global information, i.e., we only use $D_{global}$ to measure node difficulty.

We use these two difficulty measurers for ablation studies; in the below, we refer to the ablated methods as CLNode(local) and CLNode(global), respectively. GCN is used as the baseline method. The results are reported in Table 4, from which we observe the following: (1) both CLNode(local) and CLNode(global) outperform the baseline method, which demonstrates that they measure the node difficulty from different perspectives, and thus mitigate the detrimental effect of different types of *difficult nodes*; (2) CLNode achieves the best results in all experiments, proving that by combining local and global perspectives to measure the node difficulty, CLNode effectively identifies two types of *difficult nodes*, thus enhancing the accuracy and robustness of backbone GNNs.

In Table 5, we evaluate the sensitivity of CLNode to three pacing functions: linear, root, and geometric. We find that the geometric pacing function has a slight advantage on all datasets. As shown in Figure 5, the geometric function trains for a greater number of epochs on the subset of *easy nodes* before introducing *difficult nodes*. Therefore, to mitigate the detrimental effect of *difficult nodes*, we believe that the high-confidence knowledge in *easy nodes* should be fully explored before more *difficult nodes* are introduced.



**Figure 7: Parameter sensitivity analysis on Cora.**

## 5.4 Parameter Sensitivity Analysis

Last but not least, we investigate how the hyper-parameters $\lambda_0$ and $T$ affect the performance of CLNode. $\lambda_0$ controls the initial number of training nodes, while $T$ controls the speed at which *difficult nodes* are introduced to the training process. To explore the parameter sensitivity, we alter $\lambda_0$ and $T$ from {0.1, 0.2,..., 0.9} and {20, 40,..., 200}, respectively. We use GCN as the backbone GNN and report the results on Cora under random splits. The results in Figure 7 show the following: (1) Generally, with increasing $\lambda_0$, the performance tends to first increase and then decrease; specifically, the performance is relatively good when $\lambda_0$ is between 0.3 and 0.7. A too small $\lambda_0$ results in few training nodes in the initial training process, meaning that the model cannot learn efficiently. In contrast, an overly large $\lambda_0$ introduces *difficult nodes* during initial training and thus degrades the accuracy. (2) Similarly, as $T$ increases, the test accuracy tends to first increase and then decrease. A too small $T$ will quickly introduce more *difficult nodes*, thus degrading the backbone GNN's performance; conversely, an extremely large $T$ causes the backbone GNN to be trained mainly on the easy subset, causing a loss of the information contained in *difficult nodes*.

## 6 CONCLUSION

In this paper, we study the problem of training GNNs on uneven-quality training nodes. Current GNNs assume that all training nodes contribute equally during training; as a result, *difficult nodes* degrade their accuracy and robustness. To address these issues, we propose a novel framework CLNode to mitigate the detrimental effect of *difficult nodes*. Specifically, we design a multi-perspective difficulty measurer to accurately measure node difficulty using local and global information. Based on these measurements, a continuous training scheduler is proposed to feed nodes to the training progress in an easy-to-difficult curriculum. Extensive experiments on five benchmark datasets demonstrate that CLNode is a general framework that can be combined with six representative backbone GNNs to improve their accuracy. Further experiments are conducted on noisily labeled graphs to prove that CLNode enhances backbone GNNs' robustness. An interesting future direction to expand the current work is to explore the application of curriculum learning to more graph-related tasks, e.g., link prediction.

## 7 ACKNOWLEDGMENTS

## REFERENCES

[1] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*. 41–48.

[2] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. 2013. Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203* (2013).

[3] Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. 2020. Simple and deep graph convolutional networks. In *International Conference on Machine Learning*. 1725–1735.

[4] Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. 2020. Simple and deep graph convolutional networks. In *International Conference on Machine Learning*. 1725–1735.

[5] Enyan Dai, Charu Aggarwal, and Suhang Wang. 2021. Nrgnn: Learning a label noise resistant graph neural network on sparsely and noisily labeled graphs. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 227–236.

[6] Miryam de Lhoneux, Sheng Zhang, and Anders Søgaard. 2022. Zero-Shot Dependency Parsing with Worst-Case Aware Automated Curriculum Learning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. 578–587.

[7] Yingtong Dou. 2022. Robust Graph Learning for Misbehavior Detection. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*. 1545–1546.

[8] Matthias Fey and Jan Eric Lenssen. 2019. Fast graph representation learning with PyTorch Geometric. *arXiv preprint arXiv:1903.02428* (2019).

[9] Guy Hacohen and Daphna Weinshall. 2019. On the power of curriculum learning in training deep networks. In *International Conference on Machine Learning*. 2535–2544.

[10] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. *Advances in neural information processing systems* 30 (2017).

[11] Dongxiao He, Xinxin You, Zhiyong Feng, Di Jin, Xue Yang, and Weixiong Zhang. 2018. A network-specific Markov random field approach to community detection. In *Thirty-Second AAAI Conference on Artificial Intelligence*.

[12] Mengda Huang, Yang Liu, Xiang Ao, Kuan Li, Jianfeng Chi, Jinghua Feng, Hao Yang, and Qing He. 2022. AUC-oriented Graph Neural Network for Fraud Detection. In *Proceedings of the ACM Web Conference 2022*. 1311–1321.

[13] Yuge Huang, Yuhan Wang, Ying Tai, Xiaoming Liu, Pengcheng Shen, Shaoxin Li, Jilin Li, and Feiyue Huang. 2020. Curricularface: adaptive curriculum learning loss for deep face recognition. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 5901–5910.

[14] Di Jin, Ziyang Liu, Weihao Li, Dongxiao He, and Weixiong Zhang. 2019. Graph convolutional networks meet markov random fields: Semi-supervised community detection in attribute networks. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 152–159.

[15] Dongkwan Kim and Alice Oh. 2022. How to find your friendly neighborhood: Graph attention design with self-supervision. *arXiv preprint arXiv:2204.04879* (2022).

[16] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *International Conference on Learning Representations*.

[17] Yayong Li, Jie Yin, and Ling Chen. 2021. Unified robust training for graph neural networks against label noise. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 528–540.

[18] Yangdi Lu, Yang Bo, and Wenbo He. 2022. An Ensemble Model for Combating Label Noise. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*. 608–617.

[19] Yueming Lyu and Ivor W. Tsang. 2020. Curriculum Loss: Robust Learning and Generalization against Label Corruption. In *International Conference on Learning Representations*.

[20] Hoang NT, Choong Jun Jin, and Tsuyoshi Murata. 2019. Learning graph neural networks with noisy labels. *arXiv preprint arXiv:1905.01591* (2019).

[21] Emmanouil Antonios Platanios, Otilia Stretcu, Graham Neubig, Barnabas Poczos, and Tom M Mitchell. 2019. Competence-based curriculum learning for neural machine translation. *arXiv preprint arXiv:1903.09848* (2019).

[22] Meng Qu, Huiyu Cai, and Jian Tang. 2022. Neural Structured Prediction for Inductive Node Classification. In *International Conference on Learning Representations (ICLR)*.

[23] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. 2008. Collective classification in network data. *AI magazine* 29, 3 (2008), 93–93.

[24] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. 2018. Pitfalls of graph neural network evaluation. *arXiv preprint arXiv:1811.05868* (2018).

[25] Ke Sun, Zhouchen Lin, and Zhanxing Zhu. 2020. Multi-stage self-supervised learning for graph convolutional networks on graphs with few labeled nodes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 5892–5899.

[26] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2017. Graph attention networks. *arXiv preprint arXiv:1710.10903* (2017).

[27] Binghui Wang, Jinyuan Jia, and Neil Zhenqiang Gong. 2018. Graph-based security and privacy analytics via collective classification with joint weight learning and propagation. *arXiv preprint arXiv:1812.01661* (2018).

[28] Chengyi Wang, Yu Wu, Shujie Liu, Ming Zhou, and Zhenglu Yang. 2020. Curriculum pre-training for end-to-end speech translation. *arXiv preprint arXiv:2004.10093* (2020).

[29] Peiyi Wang, Liang Chen, Tianyu Liu, Damai Dai, Yunbo Cao, Baobao Chang, and Zhifang Sui. 2022. Hierarchical Curriculum Learning for AMR Parsing. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. 333–339.

[30] Wei Wang, Ye Tian, Jiquan Ngiam, Yinfei Yang, Isaac Caswell, and Zarana Parekh. 2020. Learning a Multi-Domain Curriculum for Neural Machine Translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 7711–7723.

[31] Yiwei Wang, Wei Wang, Yuxuan Liang, Yujun Cai, and Bryan Hooi. 2021. Curgraph: Curriculum learning for graph classification. In *Proceedings of the Web Conference 2021*. 1238–1248.

[32] Daphna Weinshall and Dan Amir. 2020. Theory of curriculum learning, with convex loss functions. *Journal of Machine Learning Research* 21, 222 (2020), 1–19.

[33] Daphna Weinshall, Gad Cohen, and Dan Amir. 2018. Curriculum learning by transfer learning: Theory and experiments with deep networks. In *International Conference on Machine Learning*. 5238–5246.

[34] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. 2019. Simplifying graph convolutional networks. In *International conference on machine learning*. 6861–6871.

[35] 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.

[36] Yiqing Xie, Sha Li, Carl Yang, Raymond Chi Wing Wong, and Jiawei Han. 2020. When do gnns work: Understanding and improving neighborhood aggregation. In *IJCAI International Joint Conference on Artificial Intelligence*.

[37] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. 2018. Representation learning on graphs with jumping knowledge networks. In *International conference on machine learning*. PMLR, 5453–5462.

[38] Bowen Zhang, Yidong Wang, Wenxin Hou, HAO WU, Jindong Wang, Manabu Okumura, and Takahiro Shinozaki. 2021. FlexMatch: Boosting Semi-Supervised Learning with Curriculum Pseudo Labeling. In *Advances in Neural Information Processing Systems*, Vol. 34. 18408–18419.

[39] Ge Zhang, Zhao Li, Jiaming Huang, Jia Wu, Chuan Zhou, Jian Yang, and Jianliang Gao. 2022. efraudcom: An e-commerce fraud detection system via competitive graph neural networks. *ACM Transactions on Information Systems (TOIS)* 40, 3 (2022), 1–29.

[40] Ziwei Zhang, Peng Cui, and Wenwu Zhu. 2020. Deep learning on graphs: A survey. *IEEE Transactions on Knowledge and Data Engineering* (2020).

[41] Tianyi Zhou, Shengjie Wang, and Jeff Bilmes. 2021. Robust Curriculum Learning: from clean label detection to noisy label self-correction. In *International Conference on Learning Representations*.