# **GD**<sup>2</sup>: Robust Graph Learning under Label Noise via Dual-View Prediction Discrepancy

Kailai Li $^1$ †, Jiong Lou $^1$ †, Jiawei Sun $^1$ , Honghong Zeng $^1$ , Wen Li $^5$ ,\*, Chentao Wu $^{1,3,4}$ , Yuan Luo $^1$ , Wei Zhao $^2$ , Shouguo Du $^6$ , Jie Li $^{1,3,4,*}$ 

<sup>1</sup>Department of Computer Science and Engineering, Shanghai Jiao Tong University

<sup>2</sup>Shenzhen University of Advanced Technology

<sup>3</sup>Yancheng Blockchain Research Institute

<sup>4</sup>Shanghai Key Laboratory of Trusted Data Circulation and Governance and Web3

<sup>5</sup>Shanghai University of International Business and Economics

<sup>6</sup>Shanghai Municipal Big Data Center

{kailai\_li, lj1994, noelsjw, zhhhhh5}@sjtu.edu.cn,
wen.li.sh@outlook.com, {wuct,luoyuan}@cs.sjtu.edu.cn,
{weizhao86, shouguo.du.sh}@outlook.com, lijiecs@sjtu.edu.cn

#### **Abstract**

Graph Neural Networks (GNNs) achieve strong performance in node classification tasks but exhibit substantial performance degradation under label noise. Despite recent advances in noise-robust learning, a principled approach that exploits the node-neighbor interdependencies inherent in graph data for label noise detection remains underexplored. To address this gap, we propose GD<sup>2</sup>, a noise-aware Graph learning framework that detects label noise by leveraging Dual-view prediction Discrepancies. The framework contrasts the *ego-view*, constructed from node-specific features, with the *structure-view*, derived through the aggregation of neighboring representations. The resulting discrepancy captures disruptions in semantic coherence between individual node representations and the structural context, enabling effective identification of mislabeled nodes. Building upon this insight, we further introduce a view-specific training strategy that enhances noise detection by amplifying prediction divergence through differentiated view-specific supervision. Extensive experiments on multiple datasets and noise settings demonstrate that GD<sup>2</sup> achieves superior performance over state-of-the-art baselines.

## 1 Introduction

Graph Neural Networks (GNNs) have recently achieved significant advancements in node classification, with extensive applicability across domains such as social network analysis [9], recommender systems [38], and biomedical research [36]. Nevertheless, obtaining high-quality labeled data remains challenging in real-world scenarios, and noisy labels are inevitably present. Prior research indicates that noisy labels can mislead the training of GNNs, substantially degrading predictive accuracy and generalization capability [45]. Furthermore, due to the message-passing mechanism inherent in GNNs, the negative impact of label noise is exacerbated [5, 28]. Consequently, developing robust GNN models under label noise represents a crucial and pressing research objective.

To alleviate negative impacts associated with noisy labels, general noise-learning approaches commonly exploit the memorization effect exhibited by neural networks [1, 42] or rely on handcrafted

<sup>\*</sup>Corresponding authors: Wen Li, Jie Li.

<sup>†</sup>Equal contribution.

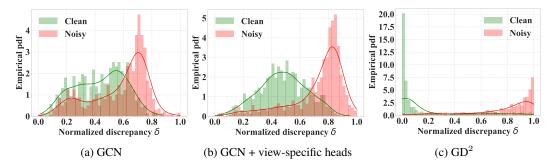


Figure 1: Normalized prediction discrepancy between the ego-view and structure-view on the Computer dataset. (a) Standard GCN yields observable but limited separation in the discrepancy distribution between clean and noisy nodes, hindering effective noise detection. (b) GCN with a shared encoder and view-specific classifier heads produces a more distinguishable discrepancy distribution. (c) GD<sup>2</sup> enhances the distributional separability in prediction discrepancies between clean and noisy nodes, yielding a well-separated distribution and facilitating noise detection.

heuristic signals [27, 10, 14] for filtering or rectifying mislabeled samples. However, these methods implicitly assume that the data are independent and identically distributed (i.i.d.), which does not hold in graph data due to interdependencies among connected nodes. Consequently, existing noise-learning techniques perform suboptimally when directly applied to node classification tasks [32]. Recently, several methods specifically tailored for graph data have emerged, incorporating structural information into noise-robust learning. Representative methods include graph structure refinement [5, 6], structure-based curriculum learning [39, 35], label correction via propagation [4, 2], and incorporation of auxiliary supervision [28, 8, 43]. A common trait among these approaches is the implicit use of semantic alignment between node features and the structural context. Such alignment captures the node—neighbor interdependencies inherent to graph data, contributing to robust representation learning and label noise mitigation. However, most existing techniques inherit noise detection mechanisms originally developed for i.i.d. data and do not explicitly leverage such relational patterns. The potential of interdependencies for identifying noisy labels remains underexplored. Consequently, an open question naturally arises: can the intrinsic dependencies between node features and structural context be explicitly modeled to enable effective label noise detection in graph data?

To address the aforementioned research question, this paper presents GD<sup>2</sup>, a noise-aware graph learning framework that explicitly leverages the semantic discrepancies between node features and structural context to detect label noise. The core component is a noise detection module based on dual-view prediction discrepancy. Specifically, node representations are decoupled into two distinct views based on information source. The ego-view is constructed solely from node-specific features and captures individual semantic information. The structure-view is derived by aggregating neighbor representations and encodes contextual semantics from local graph topology. The prediction discrepancy between the two views is exploited as an intrinsic signal for identifying potential label errors. Figure 1 illustrates the prediction discrepancy between the ego-view and structure-view under label noise. The key insight is that label noise disrupts semantic coherence between the two views, leading to divergent predictions. The ego-view, relying solely on individual features, tends to overfit spurious correlations introduced by incorrect labels, resulting in semantically misaligned representations. In contrast, the structure-view captures topological patterns that remain relatively stable under noise. Therefore, prediction divergence between the two views serves as an indicator for identifying mislabeled nodes. This insight is further supported by theoretical analysis, which reveals that clean-labeled nodes inherently exhibit bounded cross-view representation discrepancies. Building upon the above noise detection mechanism, we also introduce a view-specific training strategy, which provides differentiated supervision to each view to further enhance prediction discrepancies and facilitate noise detection. The main contributions are summarized as follows:

- (1) We propose a novel perspective for noisy label detection in graph data, that prediction discrepancy between ego-view and structure-view can be an indicator for identifying label errors.
- (2) We introduce a graph learning framework named GD<sup>2</sup>, which incorporates a noise detection module and a view-specific training strategy. By quantifying predictive discrepancies between views, the proposed framework effectively detects mislabeled nodes and facilitates robust learning.

(3) Extensive experiments on multiple real-world graph datasets under varying noise levels demonstrate that GD<sup>2</sup> consistently outperforms state-of-the-art baselines, highlighting its effectiveness in addressing label noise in graph learning.

#### 2 Related Work

## 2.1 Noisy Label Detection

Noisy label detection aims to identify and remove incorrectly labeled samples from the training data. The main challenge lies in defining a suitable metric to quantify the label quality. A widely-adopted approach leverages the memorization effect of deep neural networks [1], utilizing criteria such as training loss or prediction confidence. Representative methods include Decoupling [23], Co-teaching [13], Co-teaching+ [42], and ProMix [41]. More advanced criteria have also been developed, such as AUM [27], which identifies mislabeled samples by tracking the prediction margin during training, TopoFilter [37], which isolates noisy samples based on topological properties in the representation space, and OT-Filter [10], which selects samples from the perspective of optimal transport theory.

Recently, some studies have extended sample selection methods to graph-structured data. For instance, TSS [39] employs class-conditional betweenness centrality based on graph topology to select nodes; UnionNet [19] aggregates labels to perform sample re-weighting and label correction; and CLN-ode [35] introduces a multi-perspective difficulty measurer for detecting label noise. Nevertheless, these methods typically inherit noise detection mechanisms originally developed for i.i.d. data, failing to explicitly model intrinsic dependencies between node features and structural context for effective noisy label detection.

## 2.2 Noise Robust Learning

Besides detecting noisy labels, another line of research focuses on developing robust training algorithms and models. Such methods include designing noise-resistant network architectures [17, 11], robust loss functions [46], and regularization [34, 24], as well as training strategies [15, 20, 40].

In the context of Graph Neural Networks (GNNs), previous studies [5, 8] have demonstrated that label noise can propagate through the graph structure, significantly impairing the performance of GNNs. Several methods have mitigated this issue by refining graph structure. For instance, NRGNN [5] connects unlabeled nodes to trustworthy labeled neighbors to alleviate the effects of label noise, while RS-GNN [6] learns a denoised graph structure to reduce noise propagation. Other methods introduce auxiliary supervisory signals to strengthen robustness, including self-training and consistency regularization [28], information-theoretic objectives [49], and pairwise consistency modeling [8]. Graph contrastive learning has also been leveraged, with representative methods including RNCGLN [50], ALEX [43], and CGNN [44]. Label propagation approaches have likewise been explored: ERASE [2] integrates label propagation with coding rate minimization, while R<sup>2</sup>LP [4] applies propagation on reconstructed graphs to correct noisy labels. DND-NET [7] mitigates noise by decoupling message passing in GNNs.

Although these methods have improved the robustness of GNNs against noisy labels, their identification of mislabeled samples still predominantly relies on loss-based or confidence-based metrics. Different from existing approaches, we propose a graph learning framework that identifies noisy labels based on dual-view prediction discrepancy.

# 3 Preliminary

This work focuses on node classification under label noise. Consider a graph  $\mathcal{G}=(\mathcal{V},\mathcal{E})$ , where  $\mathcal{V}$  denotes the set of nodes and  $\mathcal{E}\subseteq\mathcal{V}\times\mathcal{V}$  represents the set of edges. The graph  $\mathcal{G}$  is assumed to be undirected and without self-loops. Node features are represented by a matrix  $\mathbf{X}\in\mathbb{R}^{|\mathcal{V}|\times d}$ , where  $|\mathcal{V}|$  is the number of nodes and d is the feature dimension. Let  $\mathbf{Y}\in\mathbb{R}^{|\mathcal{V}|\times C}$  denote the one-hot encoded ground-truth label matrix, and  $\tilde{\mathbf{Y}}\in\mathbb{R}^{|\mathcal{V}|\times C}$  the one-hot encoded noisy label matrix, where C is the number of classes. The adjacency matrix is denoted by  $\mathbf{A}\in\mathbb{R}^{|\mathcal{V}|\times |\mathcal{V}|}$ , where  $\mathbf{A}_{ij}=1$  if nodes  $v_i$  and  $v_j$  are connected, and  $\mathbf{A}_{ij}=0$  otherwise. The objective is to train a Graph Neural Network

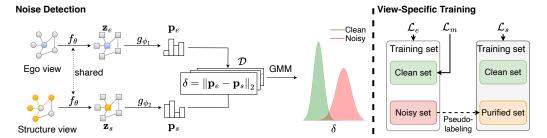


Figure 2: Overview of the proposed framework GD<sup>2</sup>, which consists of two core modules: noise detection and view-specific training. In the noise detection module, the cross-view prediction discrepancy is quantified and modeled using a Gaussian Mixture Model (GMM) to identify label noise. For model training, pseudo-labels are assigned to confident noisy nodes to refine supervision. The framework then jointly optimizes three view-specific objectives. The two modules are mutually reinforcing: the denoised training set improves model robustness, which in turn sharpens the discrepancy for more effective noise detection.

(GNN) on the training set  $\mathcal{V}_{train}$  with noisy labels  $\tilde{\mathbf{Y}}$ , such that accurate predictions can be made for the true labels  $\tilde{\mathbf{Y}}$  on the unlabeled nodes.

Graph Neural Networks (GNNs) typically follow a message-passing architecture, where node representations are updated by aggregating information from neighboring nodes. At the  $\ell$ -th layer, the representation of node v is computed as:

$$\mathbf{h}_v^{(\ell)} = \text{UPDATE}^{(\ell)} \left( \mathbf{h}_v^{(\ell-1)}, \text{AGG}^{(\ell)} \left( \left\{ \mathbf{h}_u^{(\ell-1)} \mid u \in \mathcal{N}(v) \right\} \right) \right), \tag{1}$$

where  $\mathcal{N}(v)$  denotes the set of neighbors of node v,  $AGG^{(\ell)}$  is the aggregation function, and  $\text{UPDATE}^{(\ell)}$  is the update function. Initial representations are given by node features, i.e.,  $\mathbf{h}_v^{(0)} = \mathbf{x}_v \in \mathbb{R}^d$ . For Graph Convolutional Networks [16], the representation matrix at the  $\ell$ -th layer is updated as:

$$\mathbf{H}^{(\ell)} = \sigma \left( \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{H}^{(\ell-1)} \mathbf{W}^{(\ell)} \right), \tag{2}$$

where  $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$  is the adjacency matrix with self-loops,  $\mathbf{I}$  is the identity matrix,  $\tilde{\mathbf{D}}$  is the degree matrix of  $\tilde{\mathbf{A}}$ ,  $\mathbf{W}^{(\ell)}$  is the learnable weight matrix, and  $\sigma$  denotes a non-linear activation function.

## 4 Method

This section details the proposed framework  $GD^2$ . A standard approach to learning under label noise typically involves two stages: first identifying mislabeled samples, and then training the model using the samples estimated to be clean. As illustrated in Figure 2,  $GD^2$  follows this paradigm and consists of two core components: a noise detection module and a view-specific training module. In the noise detection module, cross-view prediction discrepancies are used to estimate the likelihood of label correctness via a Gaussian Mixture Model. In the training module, pseudo-labels are generated from a mixed-view, and then distinct training objectives are designed for different views. The two components are trained in a mutually reinforcing manner: the detection module provides a cleaner training set for the model, while the training process strengthens view discrepancies, further improving noise detection. The pseudo code of  $GD^2$  is presented in the Appendix E.

# 4.1 Noise Detection via Prediction Discrepancy

Identification of mislabeled nodes is achieved by explicitly modeling prediction discrepancies between different views. The distinct views differ in information sources: the ego-view relies exclusively on intrinsic node features, whereas the structure-view aggregates neighborhood information. Node representations of the ego-view and structure-view, denoted as  $\mathbf{Z}_e$  and  $\mathbf{Z}_s$ , respectively, are first computed. Specifically, a Graph Convolutional Network (GCN)-based encoder  $f_\theta$  with parameters  $\theta$  is employed. Representations at the  $\ell$ -th layer for the ego-view and structure-view, denoted by  $\mathbf{H}_e^{(\ell)}$ 

and  $\mathbf{H}_{s}^{(\ell)}$ , are updated according to:

$$\mathbf{H}_{e}^{(\ell)} = \sigma \left( \mathbf{H}_{e}^{(\ell-1)} \mathbf{W}^{(\ell)} \right), \quad \mathbf{H}_{s}^{(\ell)} = \sigma \left( \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{H}_{s}^{(\ell-1)} \mathbf{W}^{(\ell)} \right), \tag{3}$$

where  $\mathbf{W}^{(\ell)}$  denotes learnable weights at layer  $\ell$ ,  $\sigma$  represents a non-linear activation function,  $\mathbf{A}$  is the adjacency matrix without self-loops, and  $\mathbf{D}$  is the corresponding degree matrix of  $\mathbf{A}$ . The final representations for both views are:

$$\mathbf{Z}_e = f_{\theta}(\mathbf{X}, \mathbf{I}), \quad \mathbf{Z}_s = f_{\theta}(\mathbf{X}, \mathbf{A}), \tag{4}$$

where X is the node feature matrix, I denotes the identity matrix.

To enhance the discrepancy between view-specific predictions and avoid mutual interference between the different training objectives (further discussed in Section 4.2.2), distinct linear classifiers are employed for each view. The ego-view prediction matrix  $\mathbf{P}_e \in \mathbb{R}^{|\mathcal{V}| \times C}$  and the structure-view prediction  $\mathbf{P}_s \in \mathbb{R}^{|\mathcal{V}| \times C}$  are computed using classifiers  $g_{\phi_1}$  and  $g_{\phi_2}$ , respectively:

$$\mathbf{P}_e = g_{\phi_1}(\mathbf{Z}_e), \quad \mathbf{P}_s = g_{\phi_2}(\mathbf{Z}_s). \tag{5}$$

The prediction discrepancy between the two views for node  $i \in \mathcal{V}_{\text{train}}$  is quantified using the  $\ell^2$ -norm:

$$\delta_i = \left\| \mathbf{P}_e^{[i]} - \mathbf{P}_s^{[i]} \right\|_2,\tag{6}$$

where  $\mathbf{P}^{[i]}$  represents the prediction vector of node i.

To adapt flexibly across different graphs and noise scenarios, the discrepancy set  $\mathcal{D} = \{\delta_i \mid i \in \mathcal{V}_{\text{train}}\}$  is modeled using a two-component Gaussian Mixture Model (GMM) via the Expectation-Maximization (EM) algorithm. Each discrepancy value  $\delta_i$  is mapped to a posterior probability  $\mathbb{P}_{\text{clean}}(s \mid \delta_i)$ , indicating the likelihood of belonging to the clean-label component s with the smaller mean. A threshold  $\tau$  is then applied to partition the training set into clean and noisy subsets:

$$\mathcal{V}_{\text{clean}} = \{ i \mid \mathbb{P}_{\text{clean}}(s \mid \delta_i) > \tau, i \in \mathcal{V}_{\text{train}} \}, \quad \mathcal{V}_{\text{noisy}} = \mathcal{V}_{\text{train}} \setminus \mathcal{V}_{\text{clean}}.$$
 (7)

# 4.2 Model Training

#### 4.2.1 Label Purification

After identifying nodes with noisy labels, pseudo-labels are generated to purify corrupted supervision. However, predictions obtained from either the ego-view or the structure-view exhibit inherent bias, as each is derived from a single source of information. To mitigate this issue and improve pseudo-label reliability, a mixed-view representation is introduced. The mixed-view integrates node features with local structural context to provide more accurate semantic estimates. Formally, the mixed-view representation is computed as:

$$\mathbf{Z}_m = f_{\theta}(\mathbf{X}, \tilde{\mathbf{A}}), \tag{8}$$

where  $\tilde{\bf A}$  denotes the adjacency matrix with added self-loops. Representations are updated following the standard GCN formulation [16] as described in Eq. 2. Mixed-view predictions are then obtained via the classifier  $g_{\phi_2}$  as  ${\bf P}_m = g_{\phi_2}({\bf Z}_m)$ . To ensure pseudo-label quality, only high-confidence predictions are retained. Specifically, noisy nodes with prediction confidence exceeding a predefined threshold  $\gamma$  are selected to construct the purified training set:

$$\mathcal{V}_{\text{purified}} = \{ i \mid \max(\mathbf{P}_{m}^{[i]}) > \gamma, \ i \in \mathcal{V}_{\text{noisy}} \}. \tag{9}$$

For each node selected into  $V_{purified}$ , the pseudo-label is determined by the class with the highest predicted probability:

$$\hat{\mathbf{Y}}^{[i]} = \arg \max \mathbf{P}_m^{[i]}, \quad i \in \mathcal{V}_{\text{purified}}. \tag{10}$$

## 4.2.2 Training Objective

As training progresses, predictions from different views tend to overfit noisy labels, resulting in a reduced discrepancy between views, especially for mislabeled nodes. The diminished discrepancy degrades the effectiveness of noise detection. To address this issue, view-specific training objectives are introduced to encourage distinct learning for each view. Specifically, the overall training objective

comprises three loss terms corresponding to the ego-view, structure-view, and mixed-view, each optimized on dedicated training nodes and labels.

First, the ego-view loss employs the original noisy labels  $\tilde{\mathbf{Y}}$  to supervise predictions over all training nodes  $\mathcal{V}_{\text{train}}$ , allowing the model to closely capture potential label inconsistencies. Letting  $L(\cdot,\cdot)$  denote the cross-entropy loss function, the ego-view loss is formulated as:

$$\mathcal{L}_e = \frac{1}{|\mathcal{V}_{\text{train}}|} \sum_{i \in \mathcal{V}_{\text{train}}} L(\mathbf{P}_e^{[i]}, \tilde{\mathbf{Y}}^{[i]}). \tag{11}$$

For the structure-view, pseudo-labels  $\hat{\mathbf{Y}}$  generated from the mixed-view serve as the supervisory signals. The mixed-view integrates both node features and local structural context, producing more reliable pseudo-labels that mitigate the effect of noisy labels. Notably, nodes estimated to be mislabeled receive divergent supervision: the structure-view is trained with pseudo-labels, while the ego-view relies on the original noisy labels. This asymmetry amplifies the prediction discrepancy between views, which in turn facilitates the identification of mislabeled nodes. The structure-view loss is defined as:

$$\mathcal{L}_{s} = \frac{1}{|\mathcal{V}_{\text{clean}} \cup \mathcal{V}_{\text{purified}}|} \left( \sum_{i \in \mathcal{V}_{\text{clean}}} L(\mathbf{P}_{s}^{[i]}, \tilde{\mathbf{Y}}^{[i]}) + \sum_{i \in \mathcal{V}_{\text{purified}}} L(\mathbf{P}_{s}^{[i]}, \hat{\mathbf{Y}}^{[i]}) \right). \tag{12}$$

The mixed-view is trained using the original noisy labels  $\tilde{\mathbf{Y}}$ , but restricted to the confidently identified clean nodes  $\mathcal{V}_{clean}$ . This design aims to mitigate confirmation bias, which refers to the risk of reinforcing incorrect predictions when the model is trained on unreliable pseudo-labels. The mixed-view loss is given by:

$$\mathcal{L}_{m} = \frac{1}{|\mathcal{V}_{\text{clean}}|} \sum_{i \in \mathcal{V}_{\text{clean}}} L(\mathbf{P}_{m}^{[i]}, \tilde{\mathbf{Y}}^{[i]}). \tag{13}$$

Finally, these three losses are jointly optimized, resulting in the overall training objective:

$$\mathcal{L} = \mathcal{L}_e + \mathcal{L}_s + \mathcal{L}_m. \tag{14}$$

During the evaluation stage, predictions  $\mathbf{P}_m$  from the mixed-view are employed as the final classification outputs.

## 5 Theoretical Justification

To theoretically justify the use of prediction discrepancy as an indicator of label noise, we analyze the deviation between ego-view and structure-view representations in graph neural networks. We begin by introducing assumptions on graphs and then propose a theorem that characterizes the probabilistic behavior of the discrepancy.

**Assumptions on Graphs.** Following previous works [22, 47], we consider a graph  $\mathcal{G}$ , where each node i has features  $\mathbf{x}_i \in \mathbb{R}^d$  and label  $y_i$ . We assume that (1) The features of node i are sampled from feature distribution  $\mathcal{F}_{y_i}$ , i.e.,  $\mathbf{x}_i \sim \mathcal{F}_{y_i}$ ; (2) Dimensions of  $\mathbf{x}_i$  are independent to each other; (3) The features in  $\mathbf{X}$  are bounded by a positive scalar B, i.e.,  $\max_{i,j} |\mathbf{X}[i,j]| \leq B$ ; (4) For node i, its neighbor's labels are independently sampled from neighbor distribution  $\mathcal{D}_{y_i}$ . The sampling is repeated for  $\deg(i)$  times to sample the labels for  $\deg(i)$  neighbors.

We denote a graph following these assumptions (1)–(4) as  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}, \{\mathcal{F}_c, c \in \mathcal{C}\}, \{\mathcal{D}_c, c \in \mathcal{C}\}\}$ . Note that we use the subscripts in  $\mathcal{F}_{y_i}$  and  $\mathcal{D}_{y_i}$  to indicate that these two distributions are shared by all nodes with the same label as node i. Define  $\mu_e := \mathbb{E}_{\mathbf{x} \sim \mathcal{F}_{y_i}}[\mathbf{x}]$  as the expected feature of node i given the ground-truth label  $y_i$ , and define  $\mu_s := \mathbb{E}_{c \sim \mathcal{D}_{y_i}, \ \mathbf{x} \sim \mathcal{F}_c}[\mathbf{x}]$  as the expected feature of neighbors. We analyze the embeddings obtained after a GCN operation. Following previous works [18, 3, 22, 47], we drop the non-linearity in the analysis.

**Theorem 1.** Consider a graph  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}, \{\mathcal{F}_c, c \in \mathcal{C}\}, \{\mathcal{D}_c, c \in \mathcal{C}\}\}$  that satisfies Assumptions (1)–(4). For any node  $i \in \mathcal{V}$ , let  $\mathbf{h}_i = \mathbf{W}\mathbf{x}_i$  denote the pre-activation output of a single-layer GCN using only the node's own features (ego-view), where  $\mathbf{x}_i \sim \mathcal{F}_{y_i}$ . Let  $\overline{\mathbf{h}}_i = \mathbf{W}\left(\frac{1}{|\mathcal{N}_i|}\sum_{j\in\mathcal{N}_i}\mathbf{x}_j\right)$  denote the corresponding structure-view representation that aggregates features from neighbors. Let

Table 1: Node classification accuracy (%, mean  $\pm$  std) on noisy datasets. p denotes the noise rate for each setting. Best results are highlighted in **bold**, and second-best results are underlined.

		Uniform			Pair			
Dataset	Method	p = 0.2	p = 0.4	p = 0.6	p = 0.2	p = 0.3	p = 0.4	
	GCN NRGNN RTGNN PI-GNN	$85.33 \pm 0.45$ $86.86 \pm 0.67$ $83.25 \pm 1.14$ $83.25 \pm 1.94$	$82.68 \pm 0.66$ $83.09 \pm 0.90$ $82.44 \pm 1.63$ $81.00 \pm 1.89$	$80.78 \pm 0.85$ $80.04 \pm 4.16$ $79.94 \pm 2.52$ $79.10 \pm 2.24$	$84.11 \pm 0.51$ $85.75 \pm 1.60$ $82.54 \pm 0.65$ $81.42 \pm 2.41$	$82.05 \pm 0.90$ $83.59 \pm 0.80$ $80.14 \pm 1.09$ $79.81 \pm 1.03$	$77.85 \pm 1.89$ $76.20 \pm 1.72$ $72.24 \pm 1.92$ $75.80 \pm 2.89$	
Computer	TSS ERASE GD <sup>2</sup>	$86.87 \pm 2.14$ $86.87 \pm 0.93$ $86.91 \pm 0.56$	$81.00 \pm 1.89$ $83.62 \pm 1.73$ $84.84 \pm 0.83$ $86.26 \pm 0.71$	$     \begin{array}{r}       79.10 \pm 2.24 \\       82.11 \pm 1.02 \\       \hline       83.41 \pm 1.13 \\       \hline       84.76 \pm 0.81 \\     \end{array} $	$81.42 \pm 2.41$ $83.80 \pm 1.78$ $86.68 \pm 0.75$ $86.88 \pm 0.55$	$82.56 \pm 1.75$ $84.60 \pm 1.03$ $85.75 \pm 0.64$	$79.80 \pm 2.89$ $79.89 \pm 2.03$ $80.13 \pm 4.42$ $83.16 \pm 1.75$	
Photo	GCN NRGNN RTGNN PI-GNN TSS ERASE GD <sup>2</sup>	$\begin{array}{c} 90.10 \pm 0.72 \\ 89.29 \pm 1.99 \\ 89.65 \pm 1.10 \\ 90.83 \pm 0.71 \\ 91.66 \pm 1.08 \\ \underline{91.86 \pm 0.47} \\ \mathbf{92.45 \pm 1.30} \end{array}$	$86.38 \pm 1.27$ $83.25 \pm 5.20$ $87.35 \pm 2.58$ $87.07 \pm 0.81$ $87.27 \pm 0.95$ $88.00 \pm 0.53$ $89.72 \pm 0.98$	$72.55 \pm 5.13$ $69.64 \pm 6.85$ $73.49 \pm 3.38$ $73.33 \pm 1.44$ $73.07 \pm 4.46$ <b>75.10</b> ± <b>1.57</b> $74.32 \pm 5.85$	$\begin{array}{c} 89.03 \pm 0.49 \\ 88.36 \pm 1.21 \\ \underline{91.66 \pm 1.06} \\ 89.58 \pm 0.41 \\ 89.03 \pm 0.85 \\ 91.41 \pm 0.75 \\ \underline{92.83 \pm 1.10} \end{array}$	$87.26 \pm 0.78$ $83.10 \pm 3.94$ $89.90 \pm 1.89$ $88.18 \pm 1.54$ $88.00 \pm 2.13$ $87.67 \pm 1.47$ $91.84 \pm 1.00$	$85.84 \pm 2.21$ $82.16 \pm 5.35$ $87.03 \pm 6.18$ $87.79 \pm 4.81$ $86.04 \pm 5.75$ $85.59 \pm 2.33$ $89.00 \pm 1.23$	
CS	GCN NRGNN RTGNN PI-GNN TSS ERASE GD <sup>2</sup>	$\begin{array}{c} 91.51 \pm 0.13 \\ 92.07 \pm 0.50 \\ \underline{93.23 \pm 0.33} \\ \underline{92.40 \pm 0.18} \\ 89.80 \pm 0.52 \\ \underline{91.52 \pm 0.25} \\ \mathbf{93.43 \pm 0.34} \end{array}$	$90.57 \pm 0.29$ $91.40 \pm 0.46$ $91.91 \pm 0.55$ $\overline{91.68 \pm 0.34}$ $87.23 \pm 1.04$ $91.18 \pm 0.33$ $92.37 \pm 0.45$	$88.07 \pm 0.72$ $89.37 \pm 1.23$ $\mathbf{90.62 \pm 0.94}$ $89.57 \pm 0.88$ $82.13 \pm 1.71$ $89.00 \pm 0.46$ $89.61 \pm 0.54$	$\begin{array}{c} 91.07 \pm 0.27 \\ 89.50 \pm 0.14 \\ \underline{92.44 \pm 0.52} \\ \overline{91.17 \pm 0.76} \\ 87.70 \pm 0.94 \\ 90.23 \pm 0.35 \\ \underline{92.75 \pm 0.38} \end{array}$	$88.54 \pm 0.46$ $86.44 \pm 1.10$ $89.69 \pm 1.33$ $89.75 \pm 0.93$ $84.69 \pm 2.01$ $86.92 \pm 1.34$ $90.54 \pm 0.43$	$80.67 \pm 2.12$ $75.51 \pm 6.00$ $76.16 \pm 8.44$ $82.68 \pm 5.65$ $75.83 \pm 7.20$ $81.66 \pm 2.66$ $84.68 \pm 1.92$	
WikiCS	GCN NRGNN RTGNN PI-GNN TSS ERASE GD <sup>2</sup>	$78.35 \pm 0.49$ $74.52 \pm 0.97$ $76.75 \pm 0.57$ $74.11 \pm 1.66$ $78.23 \pm 0.28$ $79.85 \pm 0.62$ $79.45 \pm 0.36$	$75.06 \pm 0.52$ $69.48 \pm 1.16$ $75.46 \pm 0.70$ $71.48 \pm 1.33$ $77.39 \pm 0.66$ $77.27 \pm 0.66$ $78.21 \pm 0.60$	$72.36 \pm 0.86$ $64.73 \pm 2.81$ $72.35 \pm 1.44$ $69.85 \pm 6.87$ $69.98 \pm 0.85$ $74.08 \pm 0.76$ $74.47 \pm 0.80$	$78.84 \pm 0.67$ $74.89 \pm 1.66$ $76.68 \pm 0.67$ $65.05 \pm 8.78$ $77.76 \pm 0.40$ $78.86 \pm 0.52$ $79.84 \pm 0.29$	$74.33 \pm 1.00$ $71.01 \pm 1.65$ $75.12 \pm 1.93$ $\overline{63.40 \pm 7.91}$ $74.67 \pm 1.03$ $74.43 \pm 1.52$ $78.09 \pm 0.42$	$69.15 \pm 1.75$ $67.74 \pm 2.83$ $71.56 \pm 3.47$ $59.28 \pm 8.92$ $65.43 \pm 3.49$ $69.74 \pm 4.07$ $75.88 \pm 1.68$	
Roman-Empire	GCN NRGNN RTGNN PI-GNN TSS ERASE GD <sup>2</sup>	$64.65 \pm 0.80$ $61.85 \pm 0.46$ $64.85 \pm 0.47$ $63.64 \pm 0.48$ $58.29 \pm 0.35$ $63.26 \pm 0.41$ $66.19 \pm 0.69$	$\begin{array}{c} 57.86 \pm 2.00 \\ 59.53 \pm 1.01 \\ \underline{60.13 \pm 0.98} \\ \overline{57.49 \pm 1.51} \\ 54.18 \pm 0.66 \\ 51.66 \pm 0.31 \\ \mathbf{60.26 \pm 0.80} \end{array}$	$50.50 \pm 1.02$ $49.60 \pm 1.36$ $51.80 \pm 1.31$ $47.74 \pm 2.15$ $47.66 \pm 1.10$ $48.74 \pm 0.44$ $51.24 \pm 1.68$	$\begin{array}{c} 66.78 \pm 1.05 \\ \hline 62.41 \pm 0.68 \\ 62.28 \pm 0.71 \\ 62.45 \pm 0.99 \\ 66.35 \pm 0.48 \\ 61.53 \pm 0.76 \\ \textbf{68.25} \pm \textbf{1.01} \end{array}$	$64.16 \pm 1.04$ $60.55 \pm 0.69$ $59.82 \pm 0.92$ $57.87 \pm 1.32$ $63.60 \pm 0.52$ $57.42 \pm 1.48$ $66.22 \pm 0.86$	$58.34 \pm 2.62$ $53.71 \pm 1.39$ $51.61 \pm 1.68$ $50.06 \pm 1.81$ $58.52 \pm 1.04$ $51.81 \pm 2.27$ <b>62.16</b> ± <b>1.14</b>	
Amazon-Ratings	GCN NRGNN RTGNN PI-GNN TSS ERASE GD <sup>2</sup>	$36.82 \pm 0.65$ $37.90 \pm 0.26$ $36.65 \pm 0.01$ $37.22 \pm 0.42$ $39.56 \pm 0.55$ $38.85 \pm 0.62$ $39.80 \pm 0.73$	$35.85 \pm 2.16$ $37.03 \pm 0.84$ $36.62 \pm 0.04$ $35.10 \pm 0.80$ $37.62 \pm 0.97$ $36.53 \pm 0.48$ $37.58 \pm 1.10$	$34.95 \pm 2.10$ $35.33 \pm 1.06$ $\overline{34.59 \pm 0.18}$ $33.17 \pm 2.25$ $34.73 \pm 2.84$ $33.26 \pm 0.74$ $35.72 \pm 1.04$	$37.58 \pm 1.07$ $37.15 \pm 1.32$ $36.59 \pm 0.12$ $36.60 \pm 0.49$ $39.65 \pm 0.52$ $37.11 \pm 1.99$ $39.97 \pm 0.55$	$35.09 \pm 3.06$ $36.98 \pm 0.62$ $34.58 \pm 3.64$ $33.35 \pm 1.42$ $38.99 \pm 0.84$ $35.47 \pm 2.87$ $38.76 \pm 0.74$	$32.99 \pm 2.12$ $36.52 \pm 0.85$ $30.87 \pm 3.39$ $31.98 \pm 2.73$ $36.45 \pm 1.46$ $31.67 \pm 5.13$ $38.28 \pm 0.90$	

d denote the feature dimensionality, and  $\rho(\mathbf{W})$  denote the largest singular value of the weight matrix  $\mathbf{W}$ . Define the discrepancy between the two views as  $\delta_i = \|\mathbf{h}_i - \overline{\mathbf{h}}_i\|_2$ . Then, for any t > 0, the probability that  $\delta_i$  exceeds t is bounded by

$$\mathbb{P}(\delta_i \ge t) \le 2d \cdot \exp\left(-\frac{\deg(i) + 1}{2B^2d} \left(\frac{t}{\rho(\mathbf{W})} - \|\boldsymbol{\mu}_e - \boldsymbol{\mu}_s\|_2\right)^2\right). \tag{15}$$

The proof is provided in Appendix D. Theorem 1 shows that clean-label nodes inherently exhibit bounded representation discrepancies between the ego-view and structure-view. This behavior arises from the statistical alignment between node-specific features  $\mu_e$  and neighborhood semantics  $\mu_s$ . Furthermore, the bound strengthens with increasing node degree  $\deg(i)$ , implying that nodes with richer neighborhood contexts exhibit more stable discrepancy patterns. This explains the robustness of structure-view representations under label noise, as neighborhood aggregation smooths out individual perturbations. In contrast, the ego-view representation of a mislabeled node is constructed independently of neighbor information and is therefore more susceptible to corrupted supervision. As a result, the representation deviates from the structure-view and violates the concentration behavior observed in clean-label cases.

Table 2: Ablation of GD<sup>2</sup> on the Computer and Roman-Empire datasets under the Pair-0.4 noise.

Ablation	$\mathcal{L}_e$	$\mathcal{L}_s$	Separate classifier	Computer	Roman-Empire
$\mathrm{GD}^2$	$\mathcal{V}_{ ext{train}}$	$\mathcal{V}_{clean} \cup \mathcal{V}_{purified}$	√	$\textbf{83.16} \pm \textbf{1.75}$	$62.16 \pm 1.14$
$\mathrm{GD}^2$ with $\phi_1 \equiv \phi_2$	$\mathcal{V}_{ ext{train}}$	$\mathcal{V}_{clean} \cup \mathcal{V}_{purified}$	×	$81.89 \pm 1.76$	$60.87 \pm 2.52$
GD <sup>2</sup> w/o label purification	$\mathcal{V}_{train}$	$\mathcal{V}_{ ext{clean}}$	$\checkmark$	$82.35 \pm 1.29$	$60.76 \pm 1.70$
GD <sup>2</sup> w/o distinct objectives	$egin{aligned} \mathcal{V}_{ ext{train}} \ \mathcal{V}_{ ext{clean}} \ \mathcal{V}_{ ext{clean}} \cup \mathcal{V}_{ ext{purified}} \end{aligned}$	$egin{aligned} \mathcal{V}_{ ext{train}} \ \mathcal{V}_{ ext{clean}} \ \mathcal{V}_{ ext{clean}} \cup \mathcal{V}_{ ext{purified}} \end{aligned}$	√ √ √	$81.13 \pm 1.26$ $82.69 \pm 0.88$ $82.35 \pm 2.06$	$61.62 \pm 1.48$ $60.97 \pm 2.73$ $61.13 \pm 1.06$

# 6 Experiments

## 6.1 Setup

**Datasets and Settings.** We conduct experiments on six benchmark datasets, including four homophilous graphs: Computer, Photo, CS [29], and WikiCS [25], and two heterophilous graphs: Roman-Empire, and Amazon-Ratings [26]. Dataset statistics are provided in Appendix A. Following previous works [12, 48], we randomly select 10% of the nodes for training, 10% for validation, and use the remaining nodes for testing. All results are averaged over 10 runs with different random seeds, and we report both the mean and standard deviation.

Following prior studies [5, 2, 28, 8], we synthetically corrupt labels of the training and validation sets. We consider two types of label noise:

- 1) Uniform-p: Each label is flipped to a randomly selected incorrect class with probability p.
- 2) Pair-p: Each label is flipped to a corresponding pair class with probability p, simulating realistic mislabeling patterns.

We set  $p \in \{0.2, 0.4, 0.6\}$  for uniform noise and  $p \in \{0.2, 0.3, 0.4\}$  for pair noise. More implementation details, including model architectures and training configurations, are provided in Appendix B.

**Baselines.** We compare GD<sup>2</sup> with several representative methods designed for learning with noisy labels on graph data. The baselines include: (1) NRGNN [5], which refines the graph structure to alleviate the effects of label noise; (2) ERASE [2], which performs label denoising via label propagation; (3) RTGNN [28] and PI-GNN [8], which incorporate auxiliary supervision; and (4) TSS [39], which leverages graph topology for node selection under a curriculum learning framework. All hyperparameters are tuned following the original publications. Additional comparisons with general noisy label learning methods are provided in Appendix C.

## 6.2 Main Results

Table 1 reports the node classification performance across diverse noise settings and datasets. As shown in the table,  $\mathrm{GD}^2$  achieves the best or second-best performance across a wide range of datasets and noise settings. Compared to representative graph-specific methods such as NRGNN and ERASE,  $\mathrm{GD}^2$  achieves substantial performance improvements, particularly under high noise rates. For example, on Computer,  $\mathrm{GD}^2$  outperforms ERASE by more than 3.03% under pair-0.4 noise. These consistent improvements indicate that  $\mathrm{GD}^2$  offers a general solution for robust node classification under label noise, delivering strong performance across diverse graphs and noise settings.

## **6.3** Ablation Study

Effect of separate classifiers. We begin by evaluating the necessity of using separate classifiers for ego-view and structure-view representations. The variant  $GD^2$  with  $\phi_1 \equiv \phi_2$  employs a shared classifier head for both views while keeping all other settings unchanged. As shown in Table 2, this modification results in a noticeable performance drop (-1.27% on Computer and -1.29% on Roman-Empire), indicating that separate classifiers are crucial for preventing interference between views. The decoupled heads allow each view to independently learn view-specific predictions and avoid mutual supervision contamination.

**Effect of label purification.** To evaluate the effect of label purification, we compare  $GD^2$  with the variant  $GD^2$  w/o label purification, where  $\mathcal{L}_s$  is trained only on the clean set while the purified

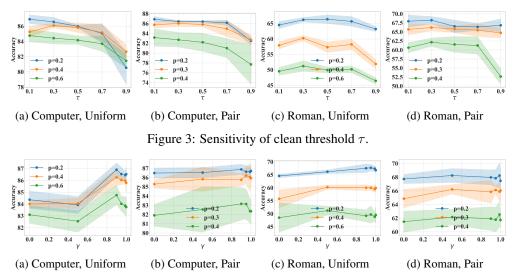


Figure 4: Sensitivity of pseudo-label confidence threshold  $\gamma$ .

set is excluded. As reported in Table 2, this variant performs worse than GD<sup>2</sup>, confirming that label purification is beneficial not only for supervision but also for enhancing the prediction discrepancy signal. Although the purified set is not directly used for training the mix-view branch, it plays an essential role in amplifying cross-view disagreement, thereby improving noise detection quality.

Effect of distinct training objectives. Finally, we assess the impact of using distinct training objectives for the ego and structure views. The variant  $GD^2$  w/o distinct objectives uses the same training set for both  $\mathcal{L}_e$  and  $\mathcal{L}_s$ , with three configurations considered: (i) both views are trained on the full noisy set  $\mathcal{V}_{\text{train}}$ , (ii) both are trained on the clean set  $\mathcal{V}_{\text{clean}}$ , and (iii) both are trained on the union of clean and purified nodes  $\mathcal{V}_{\text{clean}} \cup \mathcal{V}_{\text{purified}}$ . As shown in Table 2, all three variants exhibit performance degradation compared to  $GD^2$ , with the largest drop reaching 2.03% on Computer and 1.19% on Roman-Empire. These results highlight the importance of assigning different supervision to the views, which amplify prediction discrepancies among views.

## 6.4 Hyperparameter Sensitivity

Effect of clean threshold  $\tau$ . The clean threshold  $\tau$  determines the minimum agreement required between ego-view and structure-view predictions for a node to be considered clean. Figure 3 shows the accuracy of  $GD^2$  as  $\tau$  varies from 0.1 to 0.9. Overall, increasing  $\tau$  leads to a gradual decline in performance, as stricter thresholds reduce the coverage of clean samples. However, when  $\tau \leq 0.7$ , the degradation remains limited, suggesting that  $GD^2$  is relatively robust to moderate threshold choices. A sharp performance drop is observed at  $\tau = 0.9$ , indicating that overly conservative filtering discards too much supervision and impairs learning. These results highlight the importance of avoiding excessively large thresholds, while also demonstrating stability under moderate settings.

Effect of pseudo-label confidence threshold  $\gamma$ . The pseudo-label confidence threshold  $\gamma$  controls the trade-off between the quantity and quality of pseudo-labels used in label purification. Figure 4 reports the accuracy of GD<sup>2</sup> as  $\gamma$  varies from 0.0 to 1.0. As  $\gamma$  increases, accuracy generally improves, suggesting that high-confidence pseudo-labels provide more reliable supervision and enhance model robustness. In the low- $\gamma$  regime, performance exhibits greater variance and tends to degrade. In contrast, when  $\gamma$  is sufficiently large (e.g.,  $\gamma \geq 0.9$ ), accuracy becomes both higher and more stable, indicating that GD<sup>2</sup> is relatively insensitive to the threshold within the high-confidence range. Notably, setting  $\gamma = 1.0$ , which disables pseudo-labeling entirely, results in inferior performance compared to slightly relaxed thresholds (e.g.,  $\gamma = 0.9$ ). This confirms that incorporating a small number of highly confident pseudo-labels is more beneficial than discarding them altogether. These results highlight the importance of maintaining both precision and minimal coverage in the label purification process.

## 7 Conclusion

In this paper, we propose a novel graph learning framework  $\mathrm{GD}^2$  to mitigate the detrimental effects of label noise in Graph Neural Networks (GNNs). Our key insight is leveraging the prediction discrepancy derived from two complementary views—the ego-view, capturing individual node semantics, and the structure-view, encoding neighborhood structural semantics. This discrepancy serves as an intrinsic signal for noisy-label detection. Furthermore, we introduce a view-specific training strategy that amplifies prediction discrepancies through differentiated supervision, enhancing the noise detection module. Extensive experiments across multiple real-world datasets under various noise conditions demonstrated that  $\mathrm{GD}^2$  outperforms state-of-the-art noise-robust GNN baselines. These results demonstrate the effectiveness of modeling node—neighbor dependencies for label noise detection and highlight the potential of graph-specific relational patterns in improving robustness. Future work includes extending this framework to dynamic and heterogeneous graphs, where structural variations and multiple relation types pose additional challenges for noise detection.

# Acknowledgments and Disclosure of Funding

This work was supported in part by National Key R&D Program of China No. 2024YFB2705300, NSFC Grant 62232011, 62402315, and the Shanghai Science and Technology Innovation Action Plan Grant 24BC3201200.

#### References

- [1] Devansh Arpit, Stanisław Jastrzębski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S. Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, and Simon Lacoste-Julien. A closer look at memorization in deep networks. In *Proceedings of the International Conference on Machine Learning*, pages 233–242, 2017.
- [2] Ling-Hao Chen, Yuanshuo Zhang, Taohua Huang, Liangcai Su, Zeyi Lin, Xi Xiao, Xiaobo Xia, and Tongliang Liu. Erase: Error-resilient representation learning on graphs for label noise tolerance. In *Proceedings of the ACM International Conference on Information and Knowledge Management*, pages 270–280, 2024.
- [3] Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph convolutional networks. In *Proceedings of the 37th International Conference on Machine Learning*, pages 1725–1735, 2020.
- [4] Yao Cheng, Caihua Shan, Yifei Shen, Xiang Li, Siqiang Luo, and Dongsheng Li. Resurrecting label propagation for graphs with heterophily and label noise. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 433–444, 2024.
- [5] Enyan Dai, Charu Aggarwal, and Suhang Wang. Nrgnn: Learning a label noise resistant graph neural network on sparsely and noisily labeled graphs. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 227–236, 2021.
- [6] Enyan Dai, Wei Jin, Hui Liu, and Suhang Wang. Towards robust graph neural networks for noisy graphs with sparse labels. In *Proceedings of the ACM International Conference on Web Search and Data Mining*, pages 181–191, 2022.
- [7] Kaize Ding, Xiaoxiao Ma, Yixin Liu, and Shirui Pan. Divide and denoise: Empowering simple models for robust semi-supervised node classification against label noise. In *Proceedings of the* ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pages 574–584, 2024.
- [8] Xuefeng Du, Tian Bian, Yu Rong, Bo Han, Tongliang Liu, Tingyang Xu, Wenbing Huang, Yixuan Li, and Junzhou Huang. Noise-robust graph learning by estimating and leveraging pairwise interactions. *Transactions on Machine Learning Research*, 2023.
- [9] Wenqi Fan, Yao Ma, Qing Li, Yuan He, Eric Zhao, Jiliang Tang, and Dawei Yin. Graph neural networks for social recommendation. In *The World Wide Web Conference*, pages 417–426, 2019.

- [10] Chuanwen Feng, Yilong Ren, and Xike Xie. Ot-filter: An optimal transport filter for learning with noisy labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16164–16174, 2023.
- [11] Jacob Goldberger and Ehud Ben-Reuven. Training deep neural-networks using a noise adaptation layer. In *International Conference on Learning Representations*, 2017.
- [12] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Ávila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent A new approach to self-supervised learning. In *Advances in Neural Information Processing Systems*, pages 21271–21284, 2020.
- [13] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor W. Tsang, and Masashi Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. In *Advances in Neural Information Processing Systems*, pages 8536–8546, 2018.
- [14] Suyeon Kim, Dongha Lee, SeongKu Kang, Sukang Chae, Sanghwan Jang, and Hwanjo Yu. Learning discriminative dynamics with label corruption for noisy label detection. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 22477–22487, 2024.
- [15] Taehyeon Kim, Jaehoon Oh, Nak Yil Kim, Sangwook Cho, and Se-Young Yun. Comparing kullback-leibler divergence and mean squared error loss in knowledge distillation. In *Pro*ceedings of the International Joint Conference on Artificial Intelligence, pages 2628–2635, 2021.
- [16] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *Proceedings of the International Conference on Learning Representations*, 2017.
- [17] Kimin Lee, Sukmin Yun, Kibok Lee, Honglak Lee, Bo Li, and Jinwoo Shin. Robust inference via generative classifiers for handling noisy labels. In *International Conference on Machine Learning*, pages 3763–3772, 2019.
- [18] Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional networks for semi-supervised learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- [19] Yayong Li, Jie Yin, and Ling Chen. Unified robust training for graph neural networks against label noise. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 528–540, 2021.
- [20] Michal Lukasik, Srinadh Bhojanapalli, Aditya Krishna Menon, and Sanjiv Kumar. Does label smoothing mitigate label noise? In *Proceedings of the International Conference on Machine Learning*, pages 6448–6458, 2020.
- [21] Xingjun Ma, Hanxun Huang, Yisen Wang, Simone Romano, Sarah Erfani, and James Bailey. Normalized loss functions for deep learning with noisy labels. In *Proceedings of the 37th International Conference on Machine Learning*, pages 6543–6553, 2020.
- [22] Yao Ma, Xiaorui Liu, Neil Shah, and Jiliang Tang. Is homophily a necessity for graph neural networks? In *International Conference on Learning Representations*, 2022.
- [23] Eran Malach and Shai Shalev-Shwartz. Decoupling "when to update" from "how to update". In *Advances in Neural Information Processing Systems*, pages 960–970, 2017.
- [24] Aditya Krishna Menon, Ankit Singh Rawat, Sashank J. Reddi, and Sanjiv Kumar. Can gradient clipping mitigate label noise? In *Proceedings of the International Conference on Learning Representations*, 2020.
- [25] Péter Mernyei and Cătălina Cangea. Wiki-cs: A wikipedia-based benchmark for graph neural networks. *arXiv preprint arXiv:2007.02901*, 2020.

- [26] Oleg Platonov, Denis Kuznedelev, Michael Diskin, Artem Babenko, and Liudmila Prokhorenkova. A critical look at the evaluation of GNNs under heterophily: Are we really making progress? In *International Conference on Learning Representations*, 2023.
- [27] Geoff Pleiss, Tianyi Zhang, Ethan Elenberg, and Kilian Q. Weinberger. Identifying mislabeled data using the area under the margin ranking. In *Advances in Neural Information Processing Systems*, pages 17044–17056, 2020.
- [28] Siyi Qian, Haochao Ying, Renjun Hu, Jingbo Zhou, Jintai Chen, Danny Z. Chen, and Jian Wu. Robust training of graph neural networks via noise governance. In *Proceedings of the ACM International Conference on Web Search and Data Mining*, pages 607–615, 2023.
- [29] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls of graph neural network evaluation. *arXiv preprint arXiv:1811.05868*, 2018.
- [30] Haobo Wang, Ruixuan Xiao, Yixuan Li, Lei Feng, Gang Niu, Gang Chen, and Junbo Zhao. Pico+: Contrastive label disambiguation for robust partial label learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 3183–3198, 2024.
- [31] Yisen Wang, Xingjun Ma, Zaiyi Chen, Yuan Luo, Jinfeng Yi, and James Bailey. Symmetric cross entropy for robust learning with noisy labels. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 322–330, 2019.
- [32] Zhonghao Wang, Danyu Sun, Sheng Zhou, Haobo Wang, Jiapei Fan, Longtao Huang, and Jiajun Bu. NoisyGL: A comprehensive benchmark for graph neural networks under label noise. In *Advances in Neural Information Processing Systems*, pages 38142–38170, 2024.
- [33] Hongxin Wei, Lei Feng, Xiangyu Chen, and Bo An. Combating noisy labels by agreement: A joint training method with co-regularization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13726–13735, 2020.
- [34] Hongxin Wei, Lue Tao, Renchunzi Xie, and Bo An. Open-set label noise can improve robustness against inherent label noise. In *Advances in Neural Information Processing Systems*, pages 7978–7992, 2021.
- [35] Xiaowen Wei, Xiuwen Gong, Yibing Zhan, Bo Du, Yong Luo, and Wenbin Hu. Clnode: Curriculum learning for node classification. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*, pages 670–678, 2023.
- [36] Oliver Wieder, Stefan Kohlbacher, Mélaine Kuenemann, Arthur Garon, Pierre Ducrot, Thomas Seidel, and Thierry Langer. A compact review of molecular property prediction with graph neural networks. *Drug Discovery Today: Technologies*, pages 1–12, 2020.
- [37] Pengxiang Wu, Songzhu Zheng, Mayank Goswami, Dimitris Metaxas, and Chao Chen. A topological filter for learning with label noise. In *Advances in Neural Information Processing Systems*, pages 21382–21393, 2020.
- [38] Shiwen Wu, Fei Sun, Wentao Zhang, Xu Xie, and Bin Cui. Graph neural networks in recommender systems: a survey. *ACM Computing Surveys*, pages 1–37, 2022.
- [39] Yuhao Wu, Jiangchao Yao, Xiaobo Xia, Jun Yu, Ruxin Wang, Bo Han, and Tongliang Liu. Mitigating label noise on graphs via topological sample selection. In *Proceedings of the International Conference on Machine Learning*, pages 53944–53972, 2024.
- [40] Xiaobo Xia, Tongliang Liu, Bo Han, Chen Gong, Nannan Wang, Zongyuan Ge, and Yi Chang. Robust early-learning: Hindering the memorization of noisy labels. In *Proceedings of the International Conference on Learning Representations*, 2021.
- [41] Ruixuan Xiao, Yiwen Dong, Haobo Wang, Lei Feng, Runze Wu, Gang Chen, and Junbo Zhao. Promix: combating label noise via maximizing clean sample utility. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pages 4442–4450, 2023.
- [42] Xingrui Yu, Bo Han, Jiangchao Yao, Gang Niu, Ivor W. Tsang, and Masashi Sugiyama. How does disagreement help generalization against label corruption? In *Proceedings of the International Conference on Machine Learning*, pages 7164–7173, 2019.

- [43] Jingyang Yuan, Xiao Luo, Yifang Qin, Zhengyang Mao, Wei Ju, and Ming Zhang. Alex: Towards effective graph transfer learning with noisy labels. In *Proceedings of the ACM International Conference on Multimedia*, pages 3647–3656, 2023.
- [44] Jingyang Yuan, Xiao Luo, Yifang Qin, Yusheng Zhao, Wei Ju, and Ming Zhang. Learning on graphs under label noise. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, pages 1–5, 2023.
- [45] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, pages 107–115, 2021.
- [46] Zhilu Zhang and Mert R. Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. In *Advances in Neural Information Processing Systems*, pages 8792–8802, 2018.
- [47] Zhiwei Zhang, Minhua Lin, Junjie Xu, Zongyu Wu, Enyan Dai, and Suhang Wang. Robustness inspired graph backdoor defense. In *International Conference on Learning Representations*, 2025.
- [48] Tong Zhao, Yozen Liu, Leonardo Neves, Oliver Woodford, Meng Jiang, and Neil Shah. Data augmentation for graph neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 11015–11023, 2021.
- [49] Zhanke Zhou, Jiangchao Yao, Jiaxu Liu, Xiawei Guo, Quanming Yao, Li He, Liang Wang, Bo Zheng, and Bo Han. Combating bilateral edge noise for robust link prediction. In *Advances in Neural Information Processing Systems*, pages 21368–21414, 2023.
- [50] Yonghua Zhu, Lei Feng, Zhenyun Deng, Yang Chen, Robert Amor, and Michael Witbrock. Robust node classification on graph data with graph and label noise. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 17220–17227, 2024.

# **NeurIPS Paper Checklist**

#### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We have claimed that this paper focuses on developing a graph learning framework under label noise for the node classification task in the abstract and introduction. Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

#### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We have discussed in the appendix F that the current framework is limited to static, homogeneous graphs and may not generalize to dynamic or heterogeneous settings where temporal variations and multi-relational structures complicate noise detection.

#### Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

## 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: We have provided the full set of assumptions in Section 5 and the complete proof in Appendix D.

#### Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

## 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: In the experiment section and Appendix B, we have provided the detailed experimental settings for results reproducing.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We have provided open access to the datasets and detailed experimental settings for reproducing results in the Appendix B. The code will be released after the paper has been accepted.

#### Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We have provided the detailed experimental settings in the experiment section and appendix.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Following previous studies, we use the average accuracy and standard deviation to evaluate the model performance in the node classification task.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We have provided detailed information on hardware and software environments in the appendix.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: Our paper is conducted with the NeurIPS Code of Ethics.

## Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: This paper focuses on foundational research of graph learning under label noise, which does not involve societal impacts.

## Guidelines:

• The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

#### 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This paper poses no such risks.

## Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

#### 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We use the public datasets for experiments in this paper and have cited the original paper that produced the dataset.

#### Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: This paper does not release new assets.

#### Guidelines:

- The answer NA means that the paper does not release new assets.
- · Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- · At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

## 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This paper focuses on graph learning under label noise, which does not involve crowdsourcing nor research with human subjects.

#### Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

## 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

## 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components.

## Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

## A Dataset Statistics

Dataset	# Nodes	# Edges	# Classes	# Features	Type
Computer	13,752	245,861	10	767	Homophily
Photo	7,650	119,081	8	745	Homophily
CS	18,333	81,894	15	6,805	Homophily
WikiCS	11,701	431,206	10	300	Homophily
Roman-Empire	22,662	32,927	18	300	Heterophily
Amazon-Ratings	24,492	93,050	5	300	Heterophily

Table 3: Statistics of used graph datasets in this paper.

# **B** Implementation Details

The GNN encoder  $f_{\theta}$  is implemented as a Graph Convolutional Network (GCN) optionally equipped with residual connections, Batch Normalization, or Layer Normalization, depending on the specific dataset. The classifiers  $g_{\phi_1}$  and  $g_{\phi_2}$  are implemented as single-layer linear networks. We use the AdamW optimizer for training. For graph data augmentations, we adopt a combination of random node feature masking and edge masking, where node features and edges are independently masked with probabilities  $p_f$  and  $p_e$ , respectively.

**Hyperparameters.** We select the hyperparameters through a grid search, with the search ranges detailed in Table 4. All hyperparameters are tuned using validation sets. For baseline methods, hyperparameter tuning is conducted within the ranges recommended in their original publications.

Search range			
500, 1000			
0.001, 0.005, 0.01			
64, 256, 512			
1, 2, 3, 4, 5			
Batch Norm, Layer Norm			
True, False			
0.2, 0.3, 0.5, 0.7			
0.0, 0.1, 0.2, 0.3, 0.4			
0.0, 0.1, 0.2, 0.3, 0.4			
0.1, 0.3, 0.5, 0.7, 0.9			
0.0, 0.5, 0.9, 0.95, 0.99, 1.0			

Table 4: Search ranges for hyperparameters.

**Label Noise Generation** Following prior work [5, 8, 30], label noise is simulated by corrupting the ground-truth labels. The noise generation process is characterized using a noise transition matrix  $Q \in \mathbb{R}^{C \times C}$ , where C denotes the number of classes. Each entry  $Q_{ij} = \Pr(y_{\text{noisy}} = j \mid y_{\text{true}} = i)$  represents the probability of a true label i being flipped to a noisy label j.

For Uniform-p noise, the transition matrix Q takes the following form:

$$Q = \begin{bmatrix} 1 - p & \frac{p}{C - 1} & \cdots & \frac{p}{C - 1} \\ \frac{p}{C - 1} & 1 - p & \cdots & \frac{p}{C - 1} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{p}{C - 1} & \frac{p}{C - 1} & \cdots & 1 - p \end{bmatrix}$$

For Pair-p noise, a typical transition matrix Q is defined as:

$$Q = \begin{bmatrix} 1 - p & p & 0 & \cdots & 0 \\ 0 & 1 - p & p & \cdots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & 1 - p & p \\ p & 0 & \cdots & 0 & 1 - p \end{bmatrix}$$

Note that GD<sup>2</sup> does not rely on knowledge of the transition matrix or the noise rate to operate effectively.

**Infrastructures.** We implement the proposed GD<sup>2</sup> with PyTorch 2.4.0 and torch\_geometric 2.6.1. All experiments are conducted on a Linux server with a NVIDIA GeForce RTX 4090 GPU with 24GB memory.

# C Performance Comparison against General Noisy Label Learning Methods

We further compare GD<sup>2</sup> with representative general noisy label learning methods that are originally designed for i.i.d. data rather than graph-structured data. The baselines include: (i) robust loss function approaches such as SCE [31] and APL [21]; (ii) small-loss based sample selection methods such as Co-teaching [13] and JoCoR [33]; and (iii) confidence-based sample selection and semi-supervised learning methods, exemplified by ProMix [41]. These methods represent different paradigms of noise-robust learning and provide a complementary evaluation to validate the effectiveness of GD<sup>2</sup>.

Table 5 presents the comparison results. It can be observed that general noisy label learning methods, although effective on i.i.d. data, exhibit limited robustness when applied to graph-structured data. In contrast,  $\mathrm{GD}^2$  achieves consistently superior performance, demonstrating its effectiveness in addressing label noise by explicitly leveraging the structural information inherent in graphs. These results highlight the necessity of designing noise-robust methods tailored to graph domains rather than directly transferring techniques developed for i.i.d. settings.

# **D** Proof of Theorem

*Proof.* The expectation of  $\mathbf{h}_i$  is given by:

$$\mathbb{E}[\mathbf{h}_i] = \mathbb{E}[\mathbf{W}\mathbf{x}_i] = \mathbf{W}\mathbb{E}[\mathbf{x}_i] = \mathbf{W}\boldsymbol{\mu}_e$$

and the aggregated neighbor representation is:

$$\overline{\mathbf{h}}_i = \mathbf{W}\left(rac{1}{|\mathcal{N}_i|}\sum_{j\in\mathcal{N}_i}\mathbf{x}_j
ight), \quad \mathbb{E}[\overline{\mathbf{h}}_i] = \mathbf{W}oldsymbol{\mu}_s.$$

Define the discrepancy:

$$\delta_i = \left\| \mathbf{h}_i - \overline{\mathbf{h}}_i \right\|_2 = \left\| \mathbf{W} \left( \mathbf{x}_i - \frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} \mathbf{x}_j \right) \right\|_2 = \| \mathbf{W} \boldsymbol{\xi}_i \|_2,$$

where we let:

$$\boldsymbol{\xi}_i := \mathbf{x}_i - \frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} \mathbf{x}_j.$$

By the sub-multiplicativity of matrix norms, we have:

$$\delta_i \leq \|\mathbf{W}\|_2 \cdot \|\boldsymbol{\xi}_i\|_2 = \rho(\mathbf{W}) \cdot \|\boldsymbol{\xi}_i\|_2.$$

Thus, for any t > 0,

$$\mathbb{P}(\delta_i \ge t) \le \mathbb{P}\left(\|\boldsymbol{\xi}_i\|_2 \ge \frac{t}{\rho(\mathbf{W})}\right).$$

Table 5: Node classification accuracy (%, mean  $\pm$  std) on noisy datasets. p denotes the noise rate for each setting. Best results are highlighted in **bold**.

		Uniform			Pair			
Dataset	Method	p = 0.2	p = 0.4	p = 0.6	p = 0.2	p = 0.3	p = 0.4	
Computer	SCE	$86.15 \pm 1.80$	$84.41 \pm 0.92$	$79.91 \pm 2.10$	$86.05 \pm 0.36$	$83.57 \pm 1.31$	$76.51 \pm 2.16$	
	APL	$83.86 \pm 5.95$	$79.22 \pm 2.25$	$78.80 \pm 2.29$	$77.91 \pm 3.14$	$76.90 \pm 5.62$	$74.95 \pm 4.51$	
	Co-teaching	$84.57 \pm 1.33$	$83.59 \pm 1.38$	$78.18 \pm 2.15$	$82.62 \pm 1.16$	$82.22 \pm 1.08$	$74.38 \pm 1.47$	
	JoCoR	$81.10 \pm 1.29$	$81.03 \pm 1.20$	$78.14 \pm 1.02$	$79.24 \pm 0.26$	$77.57 \pm 0.78$	$68.63 \pm 5.16$	
	ProMix	$84.66 \pm 1.01$	$82.26 \pm 1.06$	$80.61 \pm 1.52$	$81.60 \pm 0.75$	$80.11 \pm 0.92$	$78.57 \pm 3.04$	
	$\mathrm{GD}^2$	$86.91 \pm 0.56$	$86.26 \pm 0.71$	$84.76 \pm 0.81$	$86.88 \pm 0.55$	$85.75 \pm 0.64$	$83.16 \pm 1.75$	
	SCE	$91.20 \pm 0.64$	$88.86 \pm 0.93$	$73.54 \pm 1.21$	$89.90 \pm 1.08$	$86.79 \pm 1.51$	$76.21 \pm 3.54$	
	APL	$91.72 \pm 1.05$	$89.73 \pm 0.73$	$72.09 \pm 2.05$	$89.80 \pm 2.55$	$88.30 \pm 1.11$	$83.05 \pm 3.06$	
TN .	Co-teaching	$85.60 \pm 7.02$	$77.11 \pm 6.62$	$61.98 \pm 5.02$	$86.43 \pm 5.64$	$84.67 \pm 5.67$	$78.78 \pm 9.57$	
Photo	JoCoR	$88.69 \pm 3.36$	$83.33 \pm 4.87$	$68.32 \pm 4.67$	$83.30 \pm 3.62$	$81.41 \pm 2.96$	$74.49 \pm 2.42$	
	ProMix	$91.76 \pm 1.83$	$88.36 \pm 0.11$	$68.59 \pm 2.40$	$89.03 \pm 2.14$	$88.60 \pm 1.11$	$83.86 \pm 2.39$	
	$\mathrm{GD}^2$	$92.45 \pm 1.30$	$89.72 \pm 0.98$	$74.32 \pm 5.85$	$92.83 \pm 1.10$	$91.84 \pm 1.00$	$89.00 \pm 1.23$	
	SCE	$92.12 \pm 0.29$	$91.25\pm0.39$	$88.22 \pm 1.38$	$90.10 \pm 0.66$	$87.23 \pm 1.53$	$75.89 \pm 2.60$	
	APL	$92.18 \pm 0.22$	$90.48 \pm 0.96$	$87.00 \pm 0.56$	$90.04 \pm 0.89$	$86.74 \pm 1.90$	$79.73 \pm 5.30$	
~~	Co-teaching	$91.66 \pm 0.37$	$90.52 \pm 0.48$	$86.38 \pm 1.96$	$90.17 \pm 0.89$	$87.26 \pm 0.80$	$73.81 \pm 6.66$	
CS	JoCoR	$91.30 \pm 0.82$	$91.17 \pm 0.34$	$87.44 \pm 1.05$	$88.45 \pm 1.65$	$84.15 \pm 0.88$	$70.32 \pm 6.13$	
	ProMix	$92.56 \pm 0.51$	$90.33 \pm 0.59$	$88.91 \pm 1.26$	$90.98 \pm 0.61$	$87.83 \pm 1.13$	$77.31 \pm 3.37$	
	$\mathrm{GD}^2$	$93.43 \pm 0.34$	$92.37 \pm 0.45$	$89.61 \pm 0.54$	$92.75 \pm 0.38$	$90.54 \pm 0.43$	$84.68 \pm 1.92$	
	SCE	$72.21 \pm 0.63$	$70.68 \pm 0.47$	$61.89 \pm 1.94$	$71.22\pm0.70$	$68.90 \pm 0.85$	$60.48 \pm 3.16$	
	APL	$71.61 \pm 2.16$	$68.36 \pm 2.84$	$60.82 \pm 2.17$	$69.90 \pm 0.45$	$69.09 \pm 0.73$	$63.21 \pm 5.68$	
*****	Co-teaching	$76.15 \pm 1.46$	$71.69 \pm 2.97$	$60.65 \pm 1.58$	$74.45 \pm 1.54$	$71.53 \pm 1.81$	$62.23 \pm 5.45$	
WikiCS	JoCoR	$76.09 \pm 1.70$	$74.84 \pm 1.90$	$65.19 \pm 4.33$	$73.50 \pm 0.79$	$68.79 \pm 2.57$	$60.39 \pm 3.32$	
	ProMix	$73.74 \pm 2.29$	$72.52 \pm 2.81$	$68.69 \pm 2.19$	$74.45 \pm 0.85$	$72.42 \pm 0.23$	$67.69 \pm 5.71$	
	$\mathrm{GD}^2$	$79.45 \pm 0.36$	$\textbf{78.21} \pm \textbf{0.60}$	$74.47 \pm 0.80$	$79.84 \pm 0.29$	$78.09 \pm 0.42$	$75.88 \pm 1.68$	
Roman-Empire	SCE	$57.78 \pm 1.13$	$54.83 \pm 1.86$	$50.46 \pm 3.35$	$66.84 \pm 1.22$	$54.14 \pm 0.60$	$52.20\pm0.78$	
	APL	$63.29 \pm 1.46$	$50.67 \pm 1.09$	$45.82 \pm 2.61$	$61.49 \pm 2.31$	$58.29 \pm 0.98$	$54.65 \pm 1.55$	
	Co-teaching	$53.92 \pm 0.43$	$52.16 \pm 0.46$	$48.53 \pm 0.82$	$61.59 \pm 0.63$	$59.13 \pm 0.57$	$55.47 \pm 0.62$	
	JoCoR	$59.52 \pm 0.71$	$55.25 \pm 0.75$	$50.41 \pm 1.05$	$66.91 \pm 0.22$	$57.60 \pm 0.78$	$53.67 \pm 1.70$	
	ProMix	$53.51 \pm 0.17$	$52.27 \pm 0.37$	$48.90 \pm 0.83$	$61.68 \pm 0.80$	$59.43 \pm 0.88$	$55.50 \pm 0.49$	
	$\mathrm{GD}^2$	$66.19 \pm 0.69$	$60.26 \pm 0.80$	$51.24 \pm 1.68$	$68.25 \pm 1.01$	$66.22 \pm 0.86$	$62.16 \pm 1.14$	
Amazon-Ratings	SCE	$36.66 \pm 0.17$	$36.52 \pm 0.14$	$36.29 \pm 0.33$	$36.38 \pm 0.77$	$37.25 \pm 0.40$	$35.32 \pm 2.09$	
	APL	$38.28 \pm 1.49$	$35.85 \pm 0.73$	$34.88 \pm 2.19$	$37.70 \pm 0.26$	$36.89 \pm 0.75$	$34.97 \pm 0.80$	
	Co-teaching	$38.86 \pm 1.91$	$36.12 \pm 1.46$	$35.02 \pm 1.00$	$38.63 \pm 1.64$	$37.71 \pm 0.75$	$35.87 \pm 1.09$	
	JoCoR	$38.15 \pm 0.47$	$36.93 \pm 0.31$	$35.22 \pm 1.44$	$38.57 \pm 0.38$	$37.40 \pm 0.50$	$36.50 \pm 0.63$	
	ProMix	$38.55 \pm 1.10$	$37.08 \pm 1.54$	$35.38 \pm 1.74$	$38.89 \pm 1.49$	$37.54 \pm 0.84$	$35.33 \pm 1.84$	
	$\mathrm{GD}^2$	$39.80 \pm 0.73$	$37.58 \pm 1.10$	$35.72 \pm 1.04$	$39.97 \pm 0.55$	$38.76 \pm 0.74$	$38.28 \pm 0.90$	

We now bound  $\|\boldsymbol{\xi}_i\|_2$ . First note that:

$$\mathbb{E}[\boldsymbol{\xi}_i] = \mathbb{E}[\mathbf{x}_i] - \frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} \mathbb{E}[\mathbf{x}_j] = \boldsymbol{\mu}_e - \boldsymbol{\mu}_s.$$

We utilize Hoeffding's inequality to bound the deviation of each dimension of  $\xi_i$  from its expectation. Let  $\mathbf{x}_i[k]$  denote the k-th coordinate of  $\mathbf{x}$ , for  $k=1,\ldots,d$ . Then, for each fixed k, the set  $\{\mathbf{x}_j[k]\mid j\in\mathcal{N}(i)\}$  consists of i.i.d. bounded random variables (by assumption), and by directly applying Hoeffding's inequality, we have:

$$\mathbb{P}(|\boldsymbol{\xi}_{i}[k] - \mathbb{E}[\boldsymbol{\xi}_{i}[k]]| \geq t_{1}) = \mathbb{P}\left(\left|\mathbf{x}_{i}[k] - \frac{1}{|\mathcal{N}_{i}|} \sum_{j \in \mathcal{N}_{i}} \mathbf{x}_{j}[k] - \mathbb{E}\left[\mathbf{x}_{i}[k] - \frac{1}{|\mathcal{N}_{i}|} \sum_{j \in \mathcal{N}_{i}} \mathbf{x}_{j}[k]\right]\right) \geq t_{1}\right)$$

$$= \mathbb{P}\left(\left|(\mathbf{x}_{i}[k] - \mathbb{E}[\mathbf{x}_{i}[k]]) - \left(\frac{1}{|\mathcal{N}_{i}|} \sum_{j \in \mathcal{N}_{i}} \mathbf{x}_{j}[k] - \mathbb{E}\left[\frac{1}{|\mathcal{N}_{i}|} \sum_{j \in \mathcal{N}_{i}} \mathbf{x}_{j}[k]\right]\right)\right| \geq t_{1}\right)$$

$$\leq 2 \exp\left(-\frac{\deg(i) + 1}{2B^{2}}t_{1}^{2}\right)$$
(16)

We now use the following implication: if the  $d_2$  norm of the deviation vector satisfies

$$\|\boldsymbol{\xi}_i - \mathbb{E}[\boldsymbol{\xi}_i]\|_2 \ge \sqrt{d}t_1,$$

then at least one coordinate  $k \in \{1, \dots, d\}$  must satisfy

$$|\boldsymbol{\xi}_i[k] - \mathbb{E}[\boldsymbol{\xi}_i[k]]| \geq t_1$$

Therefore, by the union bound over all d dimensions:

$$\mathbb{P}\left(\|\boldsymbol{\xi}_{i} - \mathbb{E}[\boldsymbol{\xi}_{i}]\|_{2} \geq \sqrt{d}t_{1}\right)$$

$$\leq \sum_{k=1}^{d} \mathbb{P}\left(|\boldsymbol{\xi}_{i}[k] - \mathbb{E}[\boldsymbol{\xi}_{i}[k]]| \geq t_{1}\right)$$

$$\leq 2d \cdot \exp\left(-\frac{\deg(i) + 1}{2B^{2}}t_{1}^{2}\right).$$
(17)

Now, set  $t_1 = \frac{t_2}{\sqrt{d}}$ , we obtain the vector norm deviation bound:

$$\mathbb{P}\left(\left\|\boldsymbol{\xi}_{i} - \mathbb{E}[\boldsymbol{\xi}_{i}]\right\|_{2} \ge t_{2}\right) \le 2d \cdot \exp\left(-\frac{\deg(i) + 1}{2B^{2}d}t_{2}^{2}\right).$$

By triangle inequality:

$$\|\boldsymbol{\xi}_i\|_2 \leq \|\boldsymbol{\xi}_i - \mathbb{E}[\boldsymbol{\xi}_i]\|_2 + \|\mathbb{E}[\boldsymbol{\xi}_i]\|_2 = \|\boldsymbol{\xi}_i - \mathbb{E}[\boldsymbol{\xi}_i]\|_2 + \|\boldsymbol{\mu}_e - \boldsymbol{\mu}_s\|_2.$$

Hence,

$$\mathbb{P}\left(\|\boldsymbol{\xi}_i\|_2 \geq \frac{t}{\rho(\mathbf{W})}\right) \leq \mathbb{P}\left(\|\boldsymbol{\xi}_i - \mathbb{E}[\boldsymbol{\xi}_i]\|_2 \geq \frac{t}{\rho(\mathbf{W})} - \|\boldsymbol{\mu}_e - \boldsymbol{\mu}_s\|_2\right),$$

as long as  $\frac{t}{\rho(\mathbf{W})} - \|\boldsymbol{\mu}_e - \boldsymbol{\mu}_s\|_2 > 0$ . Finally, we conclude:

$$\mathbb{P}(\delta_i \ge t) \le 2d \cdot \exp\left(-\frac{\deg(i) + 1}{2B^2d} \left(\frac{t}{\rho(\mathbf{W})} - \|\boldsymbol{\mu}_e - \boldsymbol{\mu}_s\|_2\right)^2\right).$$

## E Pseudo-code

We describe the overall pipeline of GD<sup>2</sup> in Algorithm 1.

## **F** Limitation

One limitation of the current framework lies in its focus on static, homogeneous graphs. While the proposed discrepancy-based noise detection mechanism demonstrates strong performance under these settings, its applicability to more complex graph structures remains unexplored. In particular, dynamic graphs introduce temporal evolution in both node features and structure patterns, which may obscure discrepancy signals or introduce new sources of semantic drift. Similarly, heterogeneous graphs involve multiple node and edge types, where semantic alignment across views can be relation-specific and context-dependent. These additional layers of structural and semantic complexity pose significant challenges for robust noise detection. Future work will explore how to adapt the proposed framework to these more general graph scenarios, potentially incorporating temporal modeling techniques or relation-aware discrepancy estimation to maintain effectiveness under such conditions.

## **Algorithm 1** Pseudo code of GD<sup>2</sup>

**Input**: Node feature matrix X, adjacency matrix A, noisy labels Y, training sets  $V_{\text{train}}$ **Parameter**: Training epochs E, clean probability threshold  $\tau$ , pseudo-label confidence threshold  $\gamma$ ; GNN encoder  $f_{\theta}$ , linear classifiers  $g_{\phi_1}, g_{\phi_2}$ 

- 1: Initialize parameters  $\theta$ ,  $\phi_1$ ,  $\phi_2$
- 2: **for** epoch = 1 to E **do**
- 3: // Noise Detection
- 4:
- 5:
- Compute view-specific representations:  $\mathbf{Z}_e = f_{\theta}(\mathbf{X}, \mathbf{I}), \mathbf{Z}_s = f_{\theta}(\mathbf{X}, \mathbf{A})$ Compute predictions:  $\mathbf{P}_e = g_{\phi_1}(\mathbf{Z}_e), \quad \mathbf{P}_s = g_{\phi_2}(\mathbf{Z}_s)$ Measure prediction discrepancy  $\delta_i = \|\mathbf{P}_e^{[i]} \mathbf{P}_s^{[i]}\|_2, \quad i \in \mathcal{V}_{\text{train}}$ Fit GMM on  $\{\delta_i\}$  and compute clean probability  $\mathbb{P}_{\text{clean}}(s \mid \delta_i)$ 6:
- 7:
- Determine clean set  $V_{\text{clean}}$  and noisy set  $V_{\text{noisy}}$  using threshold  $\tau$  (Eq. (7)) 8:
- 9: // Label Purification
- Compute mixed-view prediction  $\mathbf{P}_m = g_{\phi_2}(f_{\theta}(\mathbf{X}, \tilde{\mathbf{A}}))$ 10:
- Construct  $V_{purified}$  and assign pseudo-labels  $\hat{\mathbf{Y}}$  (Eq. (9) and Eq. (10)) 11:
- 12: // Model Training
- 13:
- Compute loss terms  $\mathcal{L}_e$ ,  $\mathcal{L}_s$ , and  $\mathcal{L}_m$ Joint optimization of total loss  $\mathcal{L} = \mathcal{L}_e + \mathcal{L}_s + \mathcal{L}_m$ 14:
- 15: Update model parameters  $\theta$ ,  $\phi_1$  and  $\phi_2$
- 16: **end for**

**Output**: Trained encoder  $f_{\theta}$  and classifiers  $g_{\phi_1}, g_{\phi_2}$