

---

# NoisyGL: A Comprehensive Benchmark for Graph Neural Networks under Label Noise

---

Zhonghao Wang<sup>1</sup>, Danyu Sun<sup>1</sup>, Sheng Zhou<sup>1\*</sup>, Haobo Wang<sup>1</sup>,  
Jiapei Fan<sup>2</sup>, Longtao Huang<sup>2</sup>, Jiajun Bu<sup>1</sup>

<sup>1</sup>Zhejiang Key Laboratory of Accessible Perception and Intelligent Systems,  
Collage of Computer Science, Zhejiang University <sup>2</sup>Alibaba Group  
{wangzhonghao, zhousheng\_zju, wanghaobo, bj} @zju.edu.cn,  
danyu.21@intl.zju.edu.cn,  
{jiapei.fjp, kaiyang.hlt} @alibaba-inc.com

## Abstract

Graph Neural Networks (GNNs) exhibit strong potential in node classification tasks through a message-passing mechanism. However, their performance often hinges on high-quality node labels, which are challenging to obtain in real-world scenarios due to unreliable sources or adversarial attacks. Consequently, *label noise* is common in real-world graph data, negatively impacting GNNs by propagating incorrect information during training. To address this issue, the study of Graph Neural Networks under Label Noise (GLN) has recently gained traction. However, due to variations in dataset selection, data splitting, and preprocessing techniques, the community currently lacks a comprehensive benchmark, which impedes deeper understanding and further development of GLN. To fill this gap, we introduce NoisyGL in this paper, the first comprehensive benchmark for graph neural networks under label noise. NoisyGL enables fair comparisons and detailed analyses of GLN methods on noisy labeled graph data across various datasets, with unified experimental settings and interface. Our benchmark has uncovered several important insights missed in previous research, and we believe these findings will be highly beneficial for future studies. We hope our open-source benchmark library will foster further advancements in this field. The code of the benchmark can be found in <https://github.com/eaglelab-zju/NoisyGL>.

## 1 Introduction

Many complex real-world systems can be represented as graph-structured data, including the citation network [19], biological networks [7], traffic networks [6], and social networks [9]. Graph Neural Networks (GNNs) have demonstrated substantial effectiveness in modeling graph data through a message-passing process that aggregates information from neighboring nodes [5]. Among the numerous applications of GNNs, node classification is the most thoroughly studied task, where GNNs are trained with the explicit assistance of semi-supervised node labels [1].

Although GNNs have achieved success, their performance in semi-supervised graph learning tasks is highly dependent on precise node labels, which are difficult to obtain in real-world scenarios [1]. For instance, in online social networks, the process of manually labeling millions of users is costly, and the labels often depend on unreliable user input [18]. Furthermore, graph data is vulnerable to adversarial label-flipping attacks [31]. Consequently, label noise is widespread in graph data. Research has demonstrated that label noise can significantly reduce the generalizability of machine learning models on computer vision and natural language processing scenarios [21]. In GNNs, the

---

\*Sheng Zhou is the Corresponding Author

message-passing mechanism can further exacerbate this negative impact by propagating incorrect supervision from mislabeled nodes throughout the graph, leading to substantial results [18].

To address this challenge, an intuitive solution is to draw on the success of previous Learning with Label Noise (LLN) strategies and apply them to GNNs. However, these approaches are not always applicable to graph learning tasks due to the non-i.i.d nature, sparse labeling of graph data, and message-passing mechanism of GNNs[1]. All these factors make GNNs vulnerable to label noise and hinder traditional LLN methods from being directly applied to graph learning tasks.

In recent years, researchers have developed a series of Graph Neural Networks under Label Noise (GLN) methods to achieve robust graph learning in the presence of label noise. These methods succeeded greatly by adopting Loss regularization [14, 31, 11, 2], Robust training strategy [24], Graph structure augmentation [1, 18, 32], and contrastive learning [30, 10]. Despite the researcher’s claim of the robustness of their proposed GLN methods, the comprehensive benchmark for evaluating these methods remains absent, bringing out the following problems: 1) *Existing works utilize different datasets, noise types, rates, data splitting, and processing strategies, which makes it challenging to achieve a fair comparison.* 2) *Existing work lacks an empirical understanding regarding the impact of the graph structure itself on label noise—a critical distinction between LLN and GLN.* 3) *No existing work has thoroughly examined the applicability of traditional LLN methods to graph learning problems.* These problems hinder us from gaining a comprehensive understanding of the progress in this field.

In this research, we present NoisyGL — the first comprehensive benchmark for graph neural networks under label noise. Our benchmark includes seventeen representative methods: ten GLN methods to assess their effectiveness and robustness on graphs with noisy labels, and seven LLN methods to evaluate their applicability in graph learning tasks. We employ standardized backbones and APIs, consistent data splitting, and processing strategies to ensure a fair comparison and allow users to construct their models or datasets with minimal effort easily. Besides performance and robustness evaluations, our benchmark supports multidimensional analysis, enabling researchers to explore the time efficiency of different methods and understand the influence of graph structure on the handling of label noise.

Through extensive experiments, we have the following key findings: 1) Simply applying LLN methods can’t significantly improve GNNs’ robustness to label noise. 2) Existing GLN methods can alleviate label noise in their applicable scenarios. 3) Pair noise is the most harmful label noise due to its misleading impact. 4) Negative effects of label noise can spread through the graph structure, especially in sparse graphs. 5) GLN methods involving graph structure augmentation effectively mitigate the spread effect of label noise. Our contributions can be summarized as follows:

- **Perform an in-depth review of the current research challenge.** In our study, we revisited and scrutinized the entire progression of GLN. We discovered that the lack of a thorough benchmark in this domain significantly hinders a deeper understanding.
- **Provide a comprehensive and user-friendly benchmark.** We present NoisyGL, the first comprehensive benchmark for GLN. In this benchmark, we have selected and implemented a variety of LLN and GLN methods and evaluated them across eight commonly used datasets under uniform experimental settings. Our benchmark library is available to the public on GitHub, intending to aid future research efforts.
- **Highlight the key findings and future opportunities.** Our study has resulted in several crucial findings that have the potential to greatly advance this field.

## 2 Formulations and Background

**Notations.** Consider a graph denoted by  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ , where  $\mathcal{V}$  is the set of  $N$  nodes and  $\mathcal{E}$  is the set of edges.  $\mathbf{A} \in \mathbb{R}^{N \times N}$  is the adjacency matrix and  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{N \times d}$  denotes node features matrix with dimension  $d$ . Each node has a ground truth label, the set of which is denoted by  $\mathcal{Y} = \{y_1^*, y_2^*, \dots, y_N^*\}$ . We focus on the semi-supervised node classification problem, where only a small set of nodes  $\mathcal{V}_L$  has assigned labels for training procedure, denoted as  $\mathcal{Y}_L = \{y_1^*, y_2^*, \dots, y_l^*\}$ , where  $l$  is the number of labeled nodes. The rest of them are unlabeled nodes, denoted as  $\mathcal{V}_U = \mathcal{V} - \mathcal{V}_L$ . Given  $\mathbf{X}$  and  $\mathbf{A}$ , the goal of node classification is to train a classifier  $f_\theta : (\mathbf{X}, \mathbf{A}) \rightarrow \hat{\mathbf{Y}}^{N \times c} = \{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_N\}$  by minimizing  $\mathcal{L}(f_\theta(\mathbf{X}, \mathbf{A}), \mathcal{Y}_L)$ , where  $c$  is

the number of classes,  $\mathcal{L}$  is a loss function that measures the difference between the predicted labels and the ground truth labels. Typically  $f_\theta$  is a well-designed Graph Neural Network(GNN). In this way, according to the Empirical Risk Minimization (ERM) principle, the well-trained classifier  $f_\theta$  can generalize on unseen data  $\mathcal{V}_U$ .

However, the accessible labels  $\mathcal{Y}_L$  can be corrupted by label noise in the real world, reducing the generalization ability of  $f_\theta$ . We denote the observed noisy labels as  $\mathcal{Y}_N = \{\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_l\}$  and  $\mathcal{Y}_L$  is their corresponding true labels. Typically, we consider two types of label noise, and here are their definitions:

**Uniform noise [21]** or symmetric noise assumes that the true label has a probability of  $\epsilon \in (0, 1)$  to be uniformly flipped to another class. Formally, for  $\forall j \neq i$ , we have  $p(\tilde{y} = j | y^* = i) = \frac{\epsilon}{c-1}$ , where  $c$  represents the number of classes.

**Pair noise [28]** or pair flipping. Assumes that the true label can only be flipped to its corresponding pair class with a probability  $\epsilon$ . Formally, we have  $p(\tilde{y} = y_p | y^* = y_c) = \epsilon$  and  $\forall j \neq y_p, y_c, p(\tilde{y} = j | y^* = y_c) = 0$ , where  $y_p$  is the corresponding pair class of  $y_c$ .

The transition patterns of pair noise and uniform noise are illustrated in the Appendix. B.1. It is important to note that these noise types assume that the transition probability depends only on the observed and true labels, and is independent of instances. In real-world scenarios, label noise can be much more complex. We focus on the most frequently used noise types, leaving the investigation of the other noise types for future studies.

### 3 Benchmark Design

#### 3.1 Datasets and Implementations

**Datasets.** We selected 8 node classification datasets widely used among different studies on graph label noise. These selected datasets come from different domains and exhibit different characteristics, enabling us to evaluate the generalizability of existing methods across a range of scenarios. Specifically, we use three classic citation datasets [19], namely Cora, Citeseer, Pubmed, and one author collaboration network DBLP [15], as well as two representative product co-purchase network datasets Amazon-Computers and Amazon-Photo [20]. Additionally, to analyze the model performance on heterophilous graphs, we include two representative social media network datasets BlogCatalog and Flickr [26]. We present detailed introductions to these datasets in Appendix C.1.

The splitting methods for training, validation, and test sets of the same dataset in different tasks are not always consistent. This necessitates a unified dataset splitting in our work to achieve fair comparisons. For three citation datasets, i.e. Cora Citeseer and Pubmed, we follow the most commonly used split in [31, 24, 11, 8]. For the author collaboration network DBLP, we follow the split as [1, 10]. For two co-purchase datasets Amazon-Computers and Amazon-Photo, we follow the split as [24]. For the social network datasets BlogCatalog and Flickr, we use the same split as [18]. In this study, we assume that the labels of both the training set and validation set have been affected by label noise. A clean test set is used to evaluate the model’s performance.

**Label Corruption.** In each experiment, we first generate a label transition probability matrix based on the given noise rate and the definition of noise. Then, for each clean label in the training and validation set, we draw a noisy label from a categorical distribution according to its corresponding transition probability. These noisy labels are used in the subsequent training procedure.

**Implementations** We consider a collection of state-of-the-art GLN algorithms, including NRGNN [1], RTGNN [18], CP [31], D-GNN [14], RCNGLN [32], CLNode [24], PIGNN [2], UnionNET [11], CGNN [30], and CRGNN [10]; and a set of well-designed LLN methods, including two loss correction methods Forward and Backward correction [16], two robust loss functions APL [13] and SCE [22], two multi-network learning methods Coteaching [4] and JoCoR [23], and one noise adaptation layer method S-model [3]. We have rigorously reproduced all methods according to their papers and source code. More details about these algorithms and implementations can be found in the Appendix C.2.

### 3.2 Research Questions

In this study, we aim to answer the following research questions:

#### **RQ1: Can LLN methods be applied directly to graph learning tasks?**

**Motivation.** While recent studies have suggested that applying traditional Learning with Label Noise (LLN) methods directly to graph learning tasks may not yield the best results [1], a comprehensive analysis of this issue is still lacking. We aim to investigate the suitability of existing LLN methods for graph learning and understand the underlying reasons. By tackling this question, we can gain a clearer insight into the unique challenges posed by graph label noise and identify which LLN techniques remain effective in graph learning contexts.

**Experiment Design.** To investigate this question, we select various LLN methods referenced in the Section 3.1 and implement them on the GCN[8] backbone using unified hyper-parameters. We then perform node classification experiments on the most frequently used datasets, evaluating their effectiveness under various types and levels of label noise. For each method and dataset, we record the mean accuracy metrics and standard deviations over 10 runs. Data splitting is performed randomly with a consistent ratio. By comparing the performance of these LLN methods with GCN, we determine whether they enhance the robustness of the backbone.

#### **RQ2: How much progress has been made by existing GLN methods?**

**Motivation.** While numerous GLN methods have been introduced in the literature, previous studies have used varied datasets, data splits, and preprocessing techniques, complicating the fair comparison of these methods' performance. Furthermore, we notice that the majority of existing approaches have been tested on homophily graphs, leading to concerns about their relevance to heterophily graphs, which are also commonly encountered in practice. By investigating this issue, we seek to determine if current GLN methods effectively address graph label noise and to identify their shortcomings.

**Experiment Design.** To address this question, we select and implement many advanced GLN methods as described in Section 3.1. We then assess the performance of these methods using uniform datasets and experimental settings. For each method and dataset, we record the mean test accuracy and the standard deviation across 10 runs. Since many of these GLN methods use GCN as their foundation, we compare their performance with GCN to evaluate their robustness to label noise.

#### **RQ3: Are existing GLN methods computationally efficient?**

**Motivation.** The efficiency of GNNs in terms of computation is crucial for their use in real-world applications, and considering label noise can lead to higher computational expenses. While previous research has deeply investigated the accuracy, generalization, and robustness of the GLN method, it has failed to address the computational efficiency of these approaches. Therefore, it is important to evaluate the computational efficiency of different methods.

**Experiment design.** To answer this question, we recorded the runtime and test accuracy of various methods on different datasets under 30% uniform noise. Specifically, for each method, we conducted 10 experiments for each method on each dataset. In each experiment, we measured the time when the model achieved the best accuracy on the validation set, considering it as the total runtime for that method. Through these experiments, we can assess whether the GLN methods strike a balance between computational efficiency and test accuracy.

#### **RQ4: Are existing GLN methods sensitive to noise rate?**

**Motivation.** Previous studies utilize different noise rates, making it difficult to fairly compare the performance of various methods. Therefore, it is essential to assess different methods using a consistent set of noise rates and to verify if existing GLN methods maintain stable performance across different noise levels.

**Experiment Design.** To investigate this question, we assess the performance of several GLN methods over varying noise levels using the same datasets and noise types. Specifically, we introduce label contamination with pair noise and uniform noise at rates of 10%, 20%, 30%, 40%, and 50%, while using clean labels as a baseline. We then train the GLN methods on these datasets following the experimental settings described in RQ2 and record the mean test accuracy and standard deviation from 10 runs. This evaluation allows us to determine the robustness of each method.

#### **RQ5: Are existing GLN methods robust to different types of label noise?**

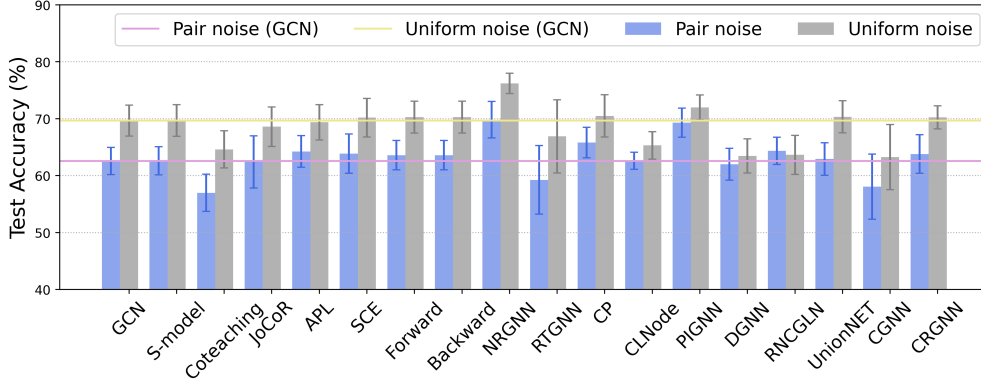


Figure 1: Test accuracy of LLN and GLN methods on DBLP dataset under 30% pair and uniform noise, respectively (10 Runs).

**Motivation.** Existing GLN methods have been developed with a variety of techniques and underlying assumptions, so they have unique strengths and weaknesses in managing different types of noise. It is crucial to identify which type of noise is most detrimental to graph learning and to understand the underlying reasons. Addressing this question will enhance our understanding of the specific scenarios in which each method excels and the distinct characteristics of different label noise types.

**Experiment design.** To tackle this question, we maintain a constant noise rate of 30% and apply both uniform and pair noise to the labels. Subsequently, we train GLN models on these noisy datasets following the experimental settings specified in RQ2 and record the mean test accuracy and standard deviation over 10 runs. This analysis enables us to determine the best method for each type of label noise and to comprehend the characteristics of various label noise types.

**RQ6: Good or bad? Revisiting the role of graph structure in label noise.**

**Motivation.** The graph structure plays a key role in distinguishing graph data from other types of data. The success of graph neural networks has largely relied on the neural message-passing mechanism, which aggregates information from neighboring nodes. However, in the presence of label noise, the messages propagated along the edges can have dual effects: on the one hand, label noise can negatively impact graph learning by spreading incorrect information; on the other hand, it can be alleviated by aligning with the majority label among the neighbors. Therefore, it is crucial to investigate whether the additional graph structure amplifies the effects of label noise and whether existing GLN methods can effectively address this challenge.

**Experiment Design.** To answer this question, we conducted comprehensive experiments on eighteen methods, including one GCN baseline, seven LLN methods, and ten GLN methods. Aiming to figure out how graph structure affects graph learning in the presence of label noise. Specifically, we recorded several metrics, including the Accuracy of Correctly Labeled Training nodes (ACLT), Accuracy of Incorrectly Labeled Training nodes (AILT), Accuracy of Unlabeled Correctly Supervised nodes (AUCS), Accuracy of Unlabeled Unsupervised nodes (AUU), and Accuracy of Unlabeled Incorrectly Supervised nodes (AUIS) under 30% uniform noise. Here, “correctly supervised,” “incorrectly supervised,” and “unsupervised” refer to unlabeled nodes that have a correctly labeled training node, an incorrectly labeled training node, and no labeled node in their neighborhood, respectively.

**4 Experiment Results and Analyses**

We present the performance of the eight methods, including vanilla GCN as a baseline, seven LLN methods with GCN backbone, and 10 GLN methods on eight datasets with different types and rates of label noise in Appendix A. Here are the key findings from the experimental results.

① **(RQ1) Most LLN methods do not significantly improve GNN robustness to label noise.** Table A1, A2, Figure 1 and Figure 2 reveal that most of the selected LLN methods do not substantially improve the performance of the GNN backbone when label noise is present. Mostly, the performance of these LLN methods remains statistically similar to the baseline. In some cases, the application of

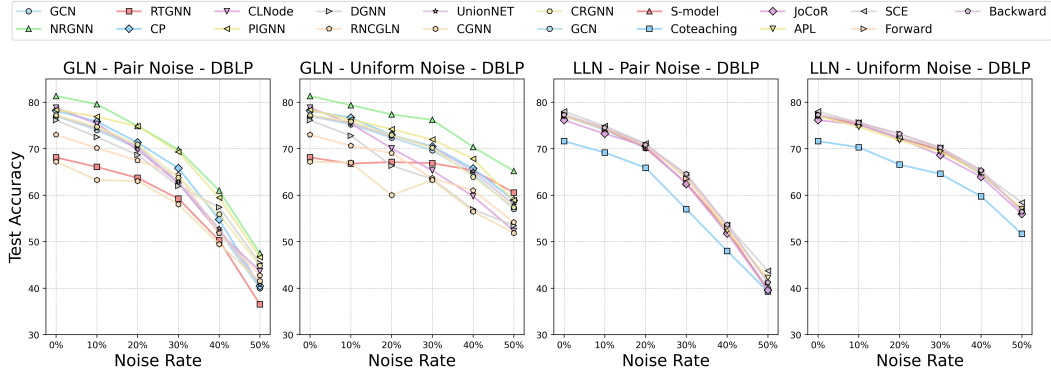


Figure 2: Test accuracy of GLN and LLN methods on DBLP dataset under different rate of pair and uniform noise (10 Runs).

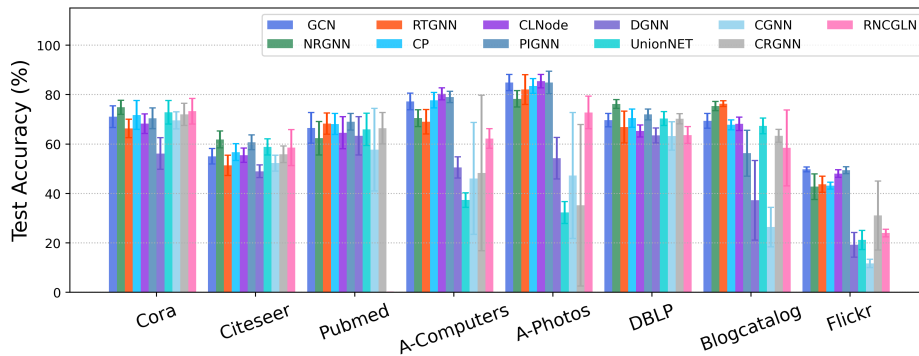


Figure 3: Test accuracy of GLN methods on all dataset under 30% uniform noise (10 Runs).

additional LLN methods can lead to a worse result. Three LLN methods incorporate a noise transition matrix, i.e. S-model, Forward, and Backward correction, demonstrating performance similar to the baseline in most cases. Typically, these transition matrix-based methods learn a diagonal transition matrix, indicating their failure to learn the label transition pattern due to the scarcity of annotations. The multi-network learning methods Coteaching and JoCoR perform similarly to the baseline on sparse graphs but underperform on dense graphs. Notably, we find that two robust loss functions, Active Passive Loss (APL) and Symmetric Cross-Entropy (SCE), slightly enhance the robustness of the baseline model across most datasets. This improvement is likely due to their ability to reduce over-fitting on mislabeled samples, though it is limited by i.i.d. assumptions. Therefore, we conclude that merely applying LLN methods to GNNs does not achieve a label noise-robust graph learning solution. Detailed experimental results are available in Appendix A.

② (RQ2) Existing GLN methods can alleviate label noise in most cases, but this improvement is limited to specific applicable scenarios. As illustrated in Table A4, A5 and Figure 3, for each dataset, there is always at least one GLN method that consistently outperforms the baseline GCN across different types of label noise, indicating that these GLN methods are effective in mitigating the graph label noise problem. However, none of them consistently perform well across all datasets. For example, NRGNN significantly outperforms the baseline GCN in Cora, Citeseer, and DBLP, but not in other datasets. This observation suggests that existing GLN methods cannot generalize across different types of data. Additionally, we observed that on Flickr, all GLN methods fail to achieve better performance than the baseline, highlighting their deficiencies in dealing with highly heterophilous graphs. Detailed experimental results are available in Appendix A.

③ (RQ3) Some GLN methods are computationally inefficient. Table A6 demonstrates that multiple GLN methods, although effective at reducing label noise, often require substantial computational resources. Figure 4 indicates that some modern GLN techniques struggle to balance performance with computational efficiency. For instance, RNCGLN is the slowest, taking 66.8 times longer than

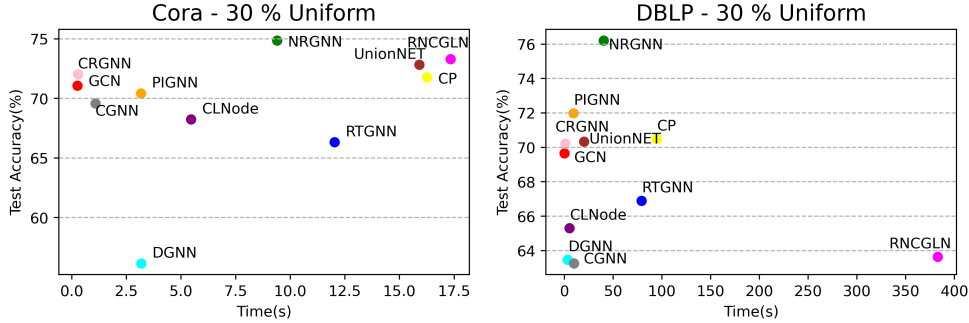


Figure 4: Time consumption and Test accuracy of different GLN methods on Cora and DBLP under 30% uniform noise (10 Runs).

Table 1: Misleading train accuracy of different methods under pair and uniform noise (10 Runs)

Dataset (Avg. # Degree)	Noise type	GCN	JoCoR	APL	NRGNN	CLNode
<b>Cora (3.90)</b>	50% <b>Uniform noise</b>	78.49 ± 9.87	70.82 ± 3.87	77.29 ± 12.35	16.16 ± 7.00	68.46 ± 17.64
	50% <b>Pair Noise</b>	94.53 ± 5.51	84.90 ± 3.49	92.93 ± 6.46	64.05 ± 12.02	88.03 ± 6.22
<b>Citeseer (2.74)</b>	50% <b>Uniform noise</b>	98.78 ± 1.15	82.64 ± 5.51	93.21 ± 5.13	32.74 ± 6.88	79.54 ± 15.19
	50% <b>Pair Noise</b>	98.54 ± 2.46	87.54 ± 5.49	96.19 ± 3.85	60.74 ± 12.91	82.23 ± 8.78
<b>A-Computers (35.76)</b>	50% <b>Uniform noise</b>	19.68 ± 6.59	19.02 ± 5.02	23.49 ± 6.44	13.36 ± 2.89	15.27 ± 5.00
	50% <b>Pair Noise</b>	75.30 ± 7.91	66.85 ± 5.91	72.42 ± 9.90	45.17 ± 8.59	64.34 ± 11.66
<b>Blogcatalog (66.11)</b>	50% <b>Uniform noise</b>	29.13 ± 9.95	27.73 ± 13.31	42.23 ± 16.14	7.93 ± 2.97	31.67 ± 9.33
	50% <b>Pair Noise</b>	72.50 ± 15.54	64.25 ± 7.95	73.39 ± 13.65	56.92 ± 10.29	64.70 ± 11.64

GCN on the Cora dataset and an astounding 2945.8 times longer on the DBLP dataset. Moreover, RNCGLN runs out of memory on the PubMed dataset, underscoring its inefficiency in memory usage. On the other hand, while the NRGNN method also consumes more time than GCN, it achieves a reasonable trade-off between performance and computational efficiency across both datasets. Detailed experimental results can be found in Appendix A.

④ **(RQ4) Most GLN methods can’t ensure a high performance under severe noise.** Figure 2 depicts the performance of different GLN methods on the DBLP dataset under various types and levels of label noise. We observe that, in general, as the noise level increases, the test accuracy of each method decreases. This decrease is most pronounced for pair noise, where the test accuracy of all methods almost halves at 50% pair noise. Additionally, we noticed that RTGNN maintains relatively stable performance under uniform noise. Moreover, two methods, NRGNN and PIGNN, show better results than the baseline GNN over different noise levels and types on the DBLP dataset. Detailed experimental results are provided in Appendix A.

⑤ **(RQ5) Pair noise is more harmful to graph learning.** In our experiments, we consistently observed that pair noise poses the most significant threat to the generalization ability of models. We have an explanation for this finding: Recall the definition provided in Section 2. For uniform noise, the true label has a chance to flip to any other class, incorrect parameter updates caused by mislabeled instances can be partially compensated by other mislabeled instances. Pair noise, however, restricts the flipping class to a specific pair class. For classifiers, this type of pair flipping can be more misleading. After being fully trained, the classifier is more likely to over-fit the pair class. This becomes particularly harmful when node features propagate through message-passing mechanisms, which can lead to a more similar embedding within local neighbors and thus make them have a similar probability of being misclassified to their corresponding pair class. To validate our hypothesis, we conducted an empirical study. Specifically, we recorded the misleading train accuracy of five methods (including 1 GCN baseline, 2 LLN and 2 GLN) on four datasets under 50% pair and uniform noise. Here the misleading train accuracy represents the model’s accuracy in making incorrect predictions to the misclassified classes. The experimental results (shown in Table 1) demonstrate that pair noise has the greatest impact, leading the model to overfit predict the mislabeled classes across different methods and datasets. Detailed experimental results are available in Appendix A.

Table 2: AUCS, AUU, AUIS of different methods on Cora and Amazon-Photos under 30% uniform noise (10 Runs)

Dataset (Avg. # Degree)	Records	GCN	NRGNN	RTGNN	CP	CLNode	RNCGLN
Cora (3.90)	AUCS	80.76 ± 2.95	83.11 ± 3.16	75.53 ± 4.80	80.85 ± 3.46	76.77 ± 3.30	77.06 ± 3.30
	AUU	71.99 ± 3.44	81.33 ± 2.25	73.44 ± 4.95	72.32 ± 4.45	67.11 ± 5.11	75.36 ± 3.33
	AUIS	51.55 ± 6.53	78.81 ± 5.94	69.50 ± 8.06	51.46 ± 12.99	43.86 ± 7.48	72.92 ± 4.66
A-Photos (31.13)	AUCS	92.21 ± 2.44	85.62 ± 2.96	89.98 ± 2.10	90.74 ± 2.98	93.08 ± 1.97	75.71 ± 6.59
	AUU	89.76 ± 2.84	84.29 ± 2.79	89.19 ± 2.13	88.85 ± 3.20	91.15 ± 2.33	74.40 ± 6.73
	AUIS	87.01 ± 4.72	82.46 ± 5.86	88.84 ± 4.76	86.34 ± 3.99	88.71 ± 3.40	71.08 ± 7.89

⑥ **(RQ6) Graph structure can amplify the negative effect of label noise.** From the experimental results in Table 2, we observed that in the sparse graph (Cora), AUIS and AUU exhibit a significant decrease compared to AUCS. Taking the performance of GCN on the Cora dataset as an example, this decrease is 36.17% and 10.85%, respectively. These results highlight the importance of proper supervision of neighboring nodes with correct annotations. Proper supervision of neighboring nodes with correct annotations significantly improves the classification accuracy of unlabeled nodes, while incorrect supervision of neighboring nodes severely reduces the classification accuracy of these nodes, even worse than when no neighborhood supervision is applied. Besides, our investigation also highlights the effectiveness of graph structure augmentation methods in mitigating the spread effect of label noise. According to Table 2, three methods, i.e. NRGNN, RTGNN, and RNCGLN, exhibit the smallest decrease in AUIS compared to AUCS and AUU among all methods. This indicates that they can effectively mitigate the spread effect of label noise. This phenomenon is even more pronounced in sparse graphs like Cora. One possible explanation can be easily drawn from the previous findings: The additional graph structure learning measures they adopted can lead to a denser graph structure used for predictions during the up-sampling process. Consequently, the classifier can rely on more references from the neighborhood, reducing its dependence on a small number of incorrectly labeled samples. Detailed experimental results are available in Appendix A.

⑦ **(RQ6) Sparse graphs are more vulnerable to the spread effect of label noise.** From Table 2 we see that the propagation effect of label noise can be very severe on sparse graphs with a relatively low average degree, like Cora, Citeseer, Pubmed, and DBLP, but not on dense graphs such as Amazon-Computers, Amazon-Photos, Blogcatalog and Flickr. The explanation for this observation is that unlabeled nodes on sparse graphs usually have only a limited number of annotated nodes in their neighborhood available for training. The prediction results of unlabeled nodes rely heavily on the annotated nodes in their neighborhood. However, if these nodes are incorrectly labeled, it will lead to erroneous learning of the embedding for the unlabeled nodes. In contrast, for dense graphs, the neighborhood of unlabeled nodes contains many annotated nodes that can serve as references. As a result, the classifier model is more likely to find correct supervision from these annotated nodes. This hypothesis is further supported by empirical evidence from Table 1, where we observe that compared to sparse graphs (such as Cora, Citeseer, and Pubmed), GCN is less susceptible to misleading on dense graphs like Blogcatalog and Amazon-Computers with a high average degree. Detailed experimental results are available in Appendix A.

## 5 Future directions

Based on the experimental results and analysis, we present several potential directions for the further development of the GLN.

**Designing widely applicable GLN approaches.** Our observations in finding ② reveal that most existing GLN methods cannot ensure consistently high performance across all scenarios. To address this problem, we need to explore three key questions: 1) What are the common properties of different graph datasets? 2) How can these common properties be utilized to enhance the robustness of GNNs against label noise? Our finding ⑥ indicated that enhancing graph structures can reduce the spread of label noise in graphs with varying densities, leading to the third question: 3) If identifying common properties is challenging, can we unify these features through data augmentation?

**Designing GLN approaches for various graph learning tasks.** Previous studies on GLN have predominantly focused on node classification tasks. However, the field of graph learning includes



other important tasks such as link prediction, edge property prediction, and graph classification. However, there is limited work on graph classification [27] and graph transfer learning [29] in the presence of label noise. Overall, research in other areas of graph learning, beyond node classification, is still in its early stages, and warrants further attention and exploration.

**Considering other types of label noise in graph learning.** Previous studies of GLN have mainly focused on pair noise and uniform noise. These noise types are instance-independent, assuming that the label corruption process is conditionally independent of node features when the true labels are given [21]. However, there exists another type of label noise—instance-dependent label noise—that is more realistic. In this case, the corruption probability depends on both the node features and the observed labels. However, none of the previous GLN studies have investigated this problem. Furthermore, unlike traditional machine learning tasks, graph learning involves additional graph structure, so the label noise model on graphs may also depend on graph topology. These issues are worth investigating, as they are more likely to occur in real-world scenarios.

## 6 Conclusions and Future work

In this research, we present NoisyGL, the first comprehensive benchmark designed for Graph Neural Networks under Label Noise (GLN) conditions. NoisyGL includes 7 prominent LLN and 10 GLN methods, allowing the community to fairly evaluate their effectiveness and robustness across various datasets. By using standardized backbones and APIs, consistent data splitting, and processing strategies, NoisyGL ensures a fair comparison and allows users to easily construct their own models or datasets with minimal effort. From this benchmark, we extract several key insights that are highly promising for the progression of this evolving field: Firstly, we point out that simply applying LLN methods cannot significantly improve the robustness of GNNs to label noise. Secondly, we found that existing GLN methods can alleviate label noise in their own applicable scenarios. In particular, pair noise emerges as the most harmful label noise due to its misleading effects. Finally, we discovered that negative effects of label noise can spread through the graph structure, especially in sparse graphs, and graph structure augmentation proves to be effective in mitigating the spread effect of label noise.

**Border Impacts and Limitations.** As NoisyGL provides a comprehensive benchmark for GNNs under label noise, we aim to attract more attention on the quality of graph data from the graph learning community, including the topology, node attributes and labels. However, NoisyGL also has some limitations that we aim to address in future work. Firstly, we aim to include a broader range of datasets to evaluate methods in different scenarios. While our current datasets are predominantly homogeneous graphs, we recognize that most GLN methods struggle with heterogeneous graphs, such as the Flickr network. Secondly, we hope to implement more GLN methods to gain a deeper understanding of the progress in the field. We will continuously update our repository to keep track of the latest advances in the field. We are also open to any suggestions and contributions that will improve the usability and effectiveness of our benchmark.

## 7 Acknowledgement

This work is supported by the National Natural Science Foundation of China (62106221, 62372408), Zhejiang Provincial Natural Science Foundation of China (Grant No: LTGG23F030005), Ningbo Natural Science Foundation (Grant No: 2022J183).

## References

- [1] 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 27th ACM SIGKDD conference on knowledge discovery & data mining*, pages 227–236, 2021.
- [2] 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. *arXiv preprint arXiv:2106.07451*, 2021.
- [3] Jacob Goldberger and Ehud Ben-Reuven. Training deep neural-networks using a noise adaptation layer. In *International conference on learning representations*, 2022.
- [4] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. *Advances in neural information processing systems*, 31, 2018.
- [5] Adrián Javaloy, Pablo Sanchez-Martin, Amit Levi, and Isabel Valera. Learnable graph convolutional attention networks. *arXiv preprint arXiv:2211.11853*, 2022.
- [6] Weiwei Jiang and Jiayun Luo. Graph neural network for traffic forecasting: A survey. *Expert Systems with Applications*, 207:117921, 2022.
- [7] Jeremy Kawahara, Colin J Brown, Steven P Miller, Brian G Booth, Vann Chau, Ruth E Grunau, Jill G Zwicker, and Ghassan Hamarneh. Brainnetcnn: Convolutional neural networks for brain networks; towards predicting neurodevelopment. *NeuroImage*, 146:1038–1049, 2017.
- [8] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- [9] Jure Leskovec and Julian McAuley. Learning to discover social circles in ego networks. *Advances in neural information processing systems*, 25, 2012.
- [10] Xianxian Li, Qiyu Li, Haodong Qian, Jinyan Wang, et al. Contrastive learning of graphs under label noise. *Neural Networks*, 172:106113, 2024.
- [11] Yayong Li, Jie Yin, and Ling Chen. Unified robust training for graph neural networks against label noise. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 528–540. Springer, 2021.
- [12] Yuwen Li, Miao Xiong, and Bryan Hooi. Graphcleaner: Detecting mislabelled samples in popular graph learning benchmarks. In *International Conference on Machine Learning*, pages 20195–20209. PMLR, 2023.
- [13] Xingjun Ma, Hanxun Huang, Yisen Wang, Simone Romano, Sarah Erfani, and James Bailey. Normalized loss functions for deep learning with noisy labels. In *International conference on machine learning*, pages 6543–6553. PMLR, 2020.
- [14] Hoang NT, Choong Jun Jin, and Tsuyoshi Murata. Learning graph neural networks with noisy labels. *arXiv preprint arXiv:1905.01591*, 2019.
- [15] Shirui Pan, Jia Wu, Xingquan Zhu, Chengqi Zhang, and Yang Wang. Tri-party deep network representation. In *International Joint Conference on Artificial Intelligence 2016*, pages 1895–1901. Association for the Advancement of Artificial Intelligence (AAAI), 2016.
- [16] Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. Making deep neural networks robust to label noise: A loss correction approach. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1944–1952, 2017.
- [17] 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? *arXiv preprint arXiv:2302.11640*, 2023.
- [18] 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 Sixteenth ACM International Conference on Web Search and Data Mining*, pages 607–615, 2023.
- [19] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.

- [20] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls of graph neural network evaluation. *arXiv preprint arXiv:1811.05868*, 2018.
- [21] Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. Learning from noisy labels with deep neural networks: A survey. *IEEE transactions on neural networks and learning systems*, 2022.
- [22] 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.
- [23] 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.
- [24] 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.
- [25] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826*, 2018.
- [26] Renchi Yang, Jieming Shi, Xiaokui Xiao, Yin Yang, Juncheng Liu, and Sourav S Bhowmick. Scaling attributed network embedding to massive graphs. *Proceedings of the VLDB Endowment*, 14(1):37–49, 2020.
- [27] Nan Yin, Li Shen, Mengzhu Wang, Xiao Luo, Zhigang Luo, and Dacheng Tao. Omg: Towards effective graph classification against label noise. *IEEE Transactions on Knowledge and Data Engineering*, 35(12):12873–12886, 2023.
- [28] Xingrui Yu, Bo Han, Jiangchao Yao, Gang Niu, Ivor Tsang, and Masashi Sugiyama. How does disagreement help generalization against label corruption? In *International conference on machine learning*, pages 7164–7173. PMLR, 2019.
- [29] 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 31st ACM international conference on multimedia*, pages 3647–3656, 2023.
- [30] Jingyang Yuan, Xiao Luo, Yifang Qin, Yusheng Zhao, Wei Ju, and Ming Zhang. Learning on graphs under label noise. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- [31] Mengmei Zhang, Linmei Hu, Chuan Shi, and Xiao Wang. Adversarial label-flipping attack and defense for graph neural networks. In *2020 IEEE International Conference on Data Mining (ICDM)*, pages 791–800. IEEE, 2020.
- [32] 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, Mar. 2024.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]**
  - (b) Did you describe the limitations of your work? **[Yes]** See Section 6.
  - (c) Did you discuss any potential negative societal impacts of your work? **[No]**
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
  - (b) Did you include complete proofs of all theoretical results? **[N/A]**
3. If you ran experiments (e.g. for benchmarks)...

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[Yes]** See Appendix D
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]** See Section 3, Appendix D and our project URL.
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[Yes]** See Section 4 and Appendix A.
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[Yes]** See Appendix D.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- (a) If your work uses existing assets, did you cite the creators? **[Yes]** See Appendix E.
  - (b) Did you mention the license of the assets? **[Yes]** See Appendix E.
  - (c) Did you include any new assets either in the supplemental material or as a URL? **[No]**
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? **[Yes]** See Appendix E.
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[No]**
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? **[N/A]**
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? **[N/A]**
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? **[N/A]**





Table A3: Additional experiment results for LLN under 30% Uniform noise (10 Runs). ACLT denotes Accuracy of Correct Labeled Training nodes, AILT denotes Accuracy of Incorrect Labeled Training nodes, AUCS denotes Accuracy of Unlabeled Correct Supervised nodes, AUU denotes Accuracy of Unlabeled Unsupervised nodes, AUIS denotes Accuracy of Unlabeled Incorrect Supervised nodes.

Dataset	Records	GCN	S-model	Coteaching	JoCoR	APL	SCE	Forward	Backward
Cora	ACLT	98.53 ± 0.93	98.79 ± 1.34	98.81 ± 1.13	95.00 ± 3.52	97.50 ± 2.09	98.33 ± 1.65	98.70 ± 0.96	98.70 ± 0.96
	AILT	33.77 ± 9.40	30.92 ± 6.95	11.87 ± 5.23	41.59 ± 10.40	38.66 ± 12.00	22.90 ± 13.41	31.19 ± 9.71	31.19 ± 9.71
	AUCS	80.76 ± 2.95	81.12 ± 3.10	77.33 ± 4.31	80.02 ± 3.30	81.11 ± 3.16	79.75 ± 3.35	81.46 ± 3.10	81.46 ± 3.10
	AUU	71.99 ± 3.44	72.04 ± 3.97	64.88 ± 6.03	72.51 ± 5.82	73.41 ± 3.98	69.34 ± 4.94	72.89 ± 3.95	72.94 ± 4.07
	AUIS	51.55 ± 6.53	51.12 ± 7.11	34.72 ± 8.61	55.40 ± 11.08	56.25 ± 7.93	44.72 ± 8.07	53.56 ± 6.87	53.70 ± 7.44
	Time	0.18 ± 0.29	0.10 ± 0.03	2.58 ± 1.02	2.36 ± 1.46	0.13 ± 0.03	1.16 ± 1.03	1.75 ± 0.16	1.71 ± 0.13
Citeseer	ACLT	99.18 ± 0.96	99.42 ± 0.83	98.95 ± 1.02	90.11 ± 2.50	98.23 ± 2.02	99.88 ± 0.37	99.05 ± 1.48	98.93 ± 1.44
	AILT	1.39 ± 2.65	1.68 ± 2.64	1.98 ± 1.87	12.05 ± 7.35	3.10 ± 3.39	1.75 ± 1.51	2.48 ± 2.70	2.48 ± 2.70
	AUCS	74.70 ± 3.49	74.55 ± 3.03	72.39 ± 3.22	73.15 ± 3.06	75.44 ± 3.33	74.12 ± 3.32	75.00 ± 2.56	75.02 ± 3.09
	AUU	57.50 ± 3.60	57.19 ± 3.42	56.25 ± 4.86	57.08 ± 2.90	57.81 ± 3.87	57.19 ± 3.49	58.75 ± 3.11	58.96 ± 3.19
	AUIS	16.45 ± 5.60	15.37 ± 5.67	19.02 ± 6.67	18.53 ± 7.35	16.30 ± 4.76	16.80 ± 5.65	20.57 ± 6.05	20.91 ± 5.57
	Time	0.55 ± 0.44	0.74 ± 0.69	2.26 ± 1.23	1.54 ± 1.25	0.94 ± 1.59	3.79 ± 1.47	2.19 ± 0.77	2.40 ± 0.86
Pubmed	ACLT	97.08 ± 2.61	96.79 ± 2.89	97.67 ± 2.50	84.79 ± 4.05	96.41 ± 2.41	89.38 ± 12.10	97.55 ± 3.24	97.55 ± 3.24
	AILT	30.01 ± 15.38	25.21 ± 18.01	18.08 ± 12.42	30.72 ± 10.11	28.59 ± 15.11	35.51 ± 15.07	12.98 ± 14.48	12.46 ± 13.81
	AUCS	69.49 ± 7.41	70.11 ± 7.61	73.03 ± 8.80	60.31 ± 11.59	71.10 ± 8.60	68.16 ± 5.19	71.72 ± 5.73	71.72 ± 5.73
	AUU	59.47 ± 7.87	59.47 ± 8.61	60.53 ± 9.37	52.63 ± 10.23	61.05 ± 7.53	59.47 ± 7.46	59.47 ± 7.87	59.47 ± 7.87
	AUIS	27.06 ± 21.05	25.95 ± 20.31	34.80 ± 28.88	27.82 ± 20.72	28.97 ± 18.29	41.27 ± 24.98	22.86 ± 15.84	22.86 ± 15.84
	Time	0.21 ± 0.06	0.35 ± 0.30	2.25 ± 1.33	2.22 ± 1.51	0.35 ± 0.13	4.13 ± 1.97	2.12 ± 0.74	2.05 ± 0.65
A-Computers	ACLT	91.05 ± 3.19	92.33 ± 2.22	86.33 ± 3.46	84.85 ± 3.60	92.49 ± 4.26	91.04 ± 4.03	51.54 ± 15.75	37.62 ± 39.19
	AILT	68.90 ± 9.98	74.61 ± 4.68	67.96 ± 6.57	64.28 ± 8.40	72.38 ± 9.26	72.56 ± 5.99	45.95 ± 14.45	31.87 ± 27.46
	AUCS	83.43 ± 2.44	83.97 ± 2.99	80.30 ± 4.47	74.94 ± 8.12	84.59 ± 2.65	84.00 ± 3.22	47.10 ± 11.71	33.83 ± 33.61
	AUU	81.40 ± 3.76	82.77 ± 3.19	78.96 ± 4.28	72.77 ± 8.48	82.95 ± 3.02	82.37 ± 3.29	47.34 ± 12.53	33.87 ± 32.51
	AUIS	77.33 ± 7.57	81.24 ± 5.33	74.86 ± 4.86	69.32 ± 9.02	80.01 ± 5.06	79.58 ± 5.87	48.72 ± 15.15	31.86 ± 29.47
	Time	1.38 ± 0.48	1.25 ± 0.40	3.70 ± 0.95	3.47 ± 1.04	3.16 ± 1.27	3.12 ± 1.09	3.11 ± 0.75	4.97 ± 2.37
A-Photos	ACLT	91.74 ± 2.96	93.33 ± 3.53	89.70 ± 4.48	86.91 ± 6.75	92.50 ± 7.11	91.48 ± 1.91	71.55 ± 15.92	51.17 ± 34.26
	AILT	80.78 ± 5.59	81.99 ± 7.92	77.28 ± 9.37	72.73 ± 9.64	79.08 ± 11.03	83.32 ± 6.83	60.99 ± 15.49	47.84 ± 29.61
	AUCS	92.18 ± 2.42	91.67 ± 3.15	87.03 ± 8.64	83.16 ± 7.40	90.99 ± 6.20	92.92 ± 1.09	69.80 ± 15.91	50.49 ± 36.54
	AUU	89.73 ± 2.81	89.84 ± 3.62	85.44 ± 7.99	80.69 ± 6.92	88.91 ± 6.24	90.91 ± 2.47	67.41 ± 15.81	49.68 ± 35.63
	AUIS	87.01 ± 4.72	86.12 ± 5.98	82.50 ± 7.82	77.67 ± 8.44	85.82 ± 8.87	87.61 ± 4.86	63.56 ± 15.53	49.36 ± 33.39
	Time	0.90 ± 0.32	0.83 ± 0.39	3.40 ± 0.81	2.51 ± 1.42	2.91 ± 1.18	1.82 ± 0.59	2.68 ± 0.46	2.92 ± 1.14
DBLP	ACLT	98.34 ± 1.71	97.93 ± 1.88	98.60 ± 2.26	94.53 ± 3.17	97.59 ± 3.28	96.69 ± 3.50	97.58 ± 1.24	97.58 ± 1.24
	AILT	21.13 ± 8.02	22.55 ± 8.58	12.27 ± 7.26	35.82 ± 7.49	25.92 ± 9.42	28.77 ± 8.09	26.37 ± 10.39	26.37 ± 10.39
	AUCS	83.15 ± 1.43	83.28 ± 1.30	81.90 ± 1.54	82.52 ± 2.03	83.12 ± 1.26	83.11 ± 1.51	83.35 ± 1.64	83.38 ± 1.61
	AUU	74.10 ± 2.64	74.60 ± 2.59	70.47 ± 5.51	75.60 ± 3.29	74.85 ± 2.70	76.80 ± 3.02	75.58 ± 2.40	75.62 ± 2.35
	AUIS	55.65 ± 9.47	57.08 ± 8.88	48.87 ± 11.69	62.24 ± 6.59	57.99 ± 9.69	64.64 ± 7.72	59.80 ± 8.13	59.87 ± 8.04
	Time	0.13 ± 0.03	0.14 ± 0.04	2.41 ± 1.32	1.71 ± 1.46	0.19 ± 0.07	1.32 ± 1.02	1.65 ± 0.14	1.62 ± 0.08
Blogcatalog	ACLT	87.78 ± 3.14	88.25 ± 4.84	58.38 ± 35.27	81.27 ± 5.43	94.92 ± 3.48	63.98 ± 8.80	78.50 ± 6.35	78.11 ± 6.82
	AILT	58.52 ± 6.32	59.95 ± 7.53	37.31 ± 16.42	53.17 ± 4.94	46.59 ± 10.66	55.37 ± 6.46	65.54 ± 5.15	64.26 ± 4.67
	AUCS	71.37 ± 3.05	71.96 ± 2.63	44.65 ± 25.60	65.42 ± 3.40	71.06 ± 4.20	59.45 ± 8.76	68.73 ± 4.84	68.63 ± 6.14
	AUU	70.64 ± 2.99	71.23 ± 2.50	43.95 ± 25.74	64.57 ± 3.45	70.21 ± 4.29	58.64 ± 8.74	68.13 ± 4.73	68.10 ± 5.96
	AUIS	72.32 ± 3.60	72.92 ± 3.15	45.86 ± 24.99	66.18 ± 3.41	71.63 ± 4.87	60.15 ± 8.94	69.36 ± 5.60	69.75 ± 5.52
	Time	0.94 ± 0.26	1.05 ± 0.35	3.87 ± 4.26	3.32 ± 1.03	3.86 ± 1.66	5.95 ± 0.23	3.47 ± 0.47	2.80 ± 0.45
Flickr	ACLT	93.22 ± 3.98	89.85 ± 6.21	77.80 ± 22.21	76.89 ± 1.79	81.19 ± 8.35	35.80 ± 6.44	31.01 ± 4.21	30.41 ± 7.36
	AILT	20.27 ± 5.10	25.32 ± 9.39	21.54 ± 7.61	24.37 ± 3.79	22.35 ± 5.49	22.27 ± 6.43	25.70 ± 4.85	25.02 ± 5.29
	AILMT	68.88 ± 10.70	57.16 ± 16.09	58.87 ± 18.91	55.92 ± 5.74	58.79 ± 11.60	17.57 ± 4.30	15.95 ± 1.97	16.15 ± 2.16
	AUCS	51.31 ± 2.98	50.55 ± 2.72	42.35 ± 10.53	45.96 ± 3.58	45.69 ± 4.07	25.86 ± 4.90	27.27 ± 3.34	26.77 ± 4.94
	AUU	50.72 ± 2.80	50.09 ± 2.60	41.95 ± 10.40	45.45 ± 3.43	45.28 ± 3.81	25.30 ± 5.00	27.02 ± 3.22	26.63 ± 4.73
	AUIS	50.92 ± 3.67	49.94 ± 2.87	42.63 ± 10.01	45.23 ± 3.69	45.09 ± 3.99	24.81 ± 4.58	27.67 ± 2.46	27.09 ± 4.25
Time	1.83 ± 0.72	1.40 ± 0.67	6.48 ± 2.78	4.21 ± 1.01	3.76 ± 2.29	1.13 ± 1.84	3.05 ± 0.74	2.79 ± 0.65	









## B Additional experiment result figures

### B.1 Transition patterns

In this work, we primarily consider two types of label noise: pair noise and uniform noise, as defined in Section 2. Additionally, our code implementation supports random noise. Unlike the fixed transition patterns of the first two types, the label transition pattern for random noise is generated randomly. Below are examples of their label transition probabilities.

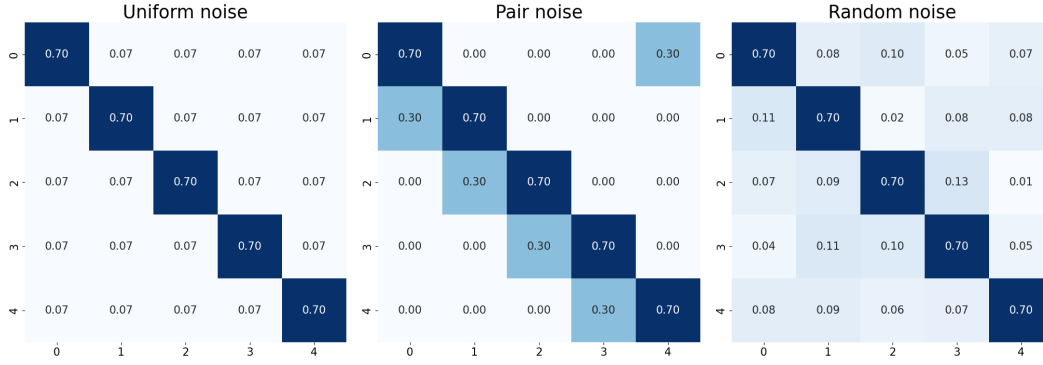


Figure A1: Label transition probability under 30% uniform noise, pair noise and random noise, respectively.

### B.2 Test accuracy of different methods under different noise rate

In Section 4, we investigated the performance of different LLN and GLN methods, here are additional experimental results.

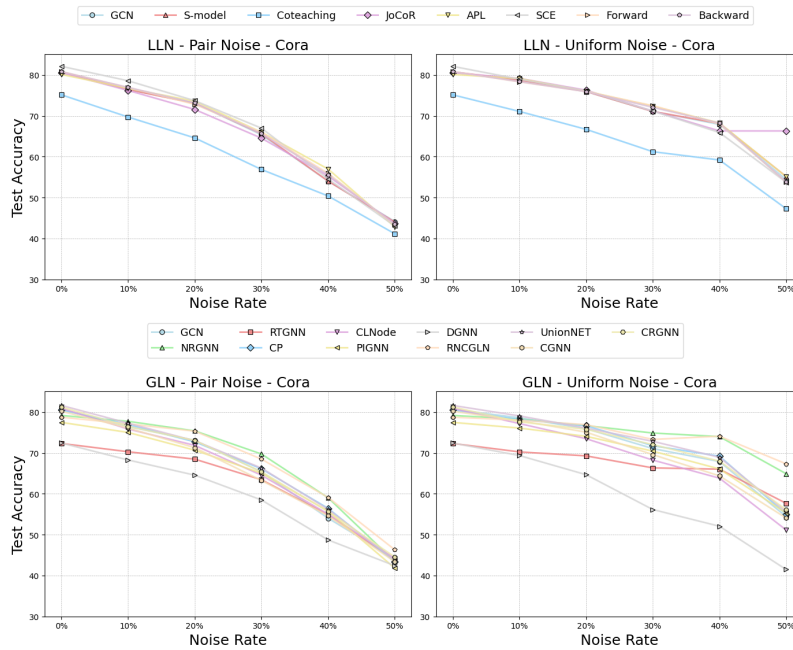


Figure A2: Test accuracy of LLN and GLN methods on Cora dataset under different rate of pair and uniform noise, respectively (10 Runs).

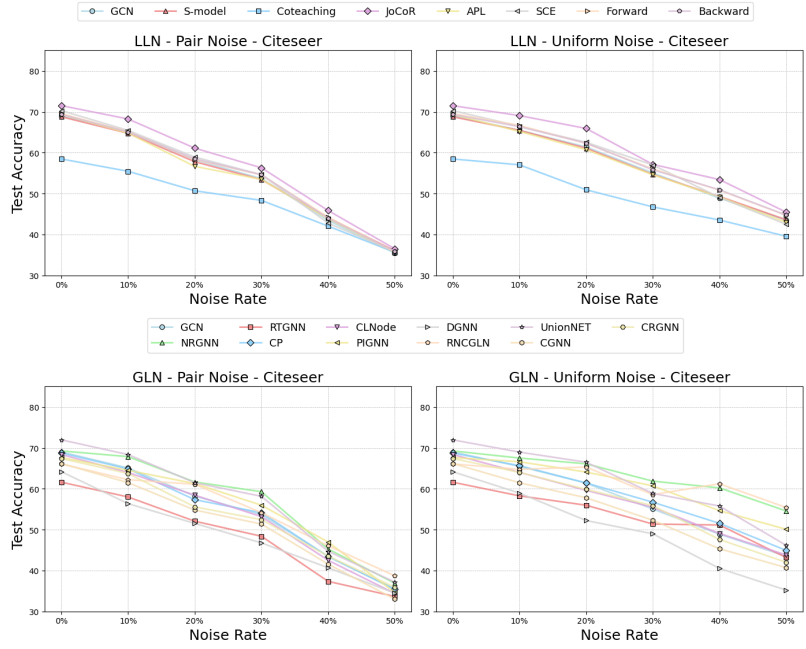


Figure A3: Test accuracy of LLN and GLN methods on Citeseer dataset under different rate of pair and uniform noise, respectively (10 Runs).

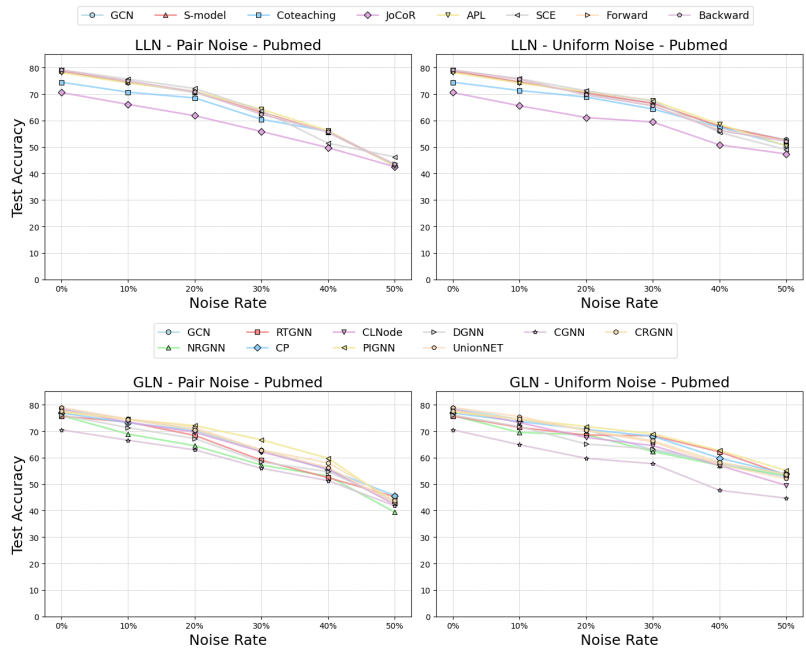


Figure A4: Test accuracy of LLN and GLN methods on Pubmed dataset under different rate of pair and uniform noise, respectively (10 Runs).

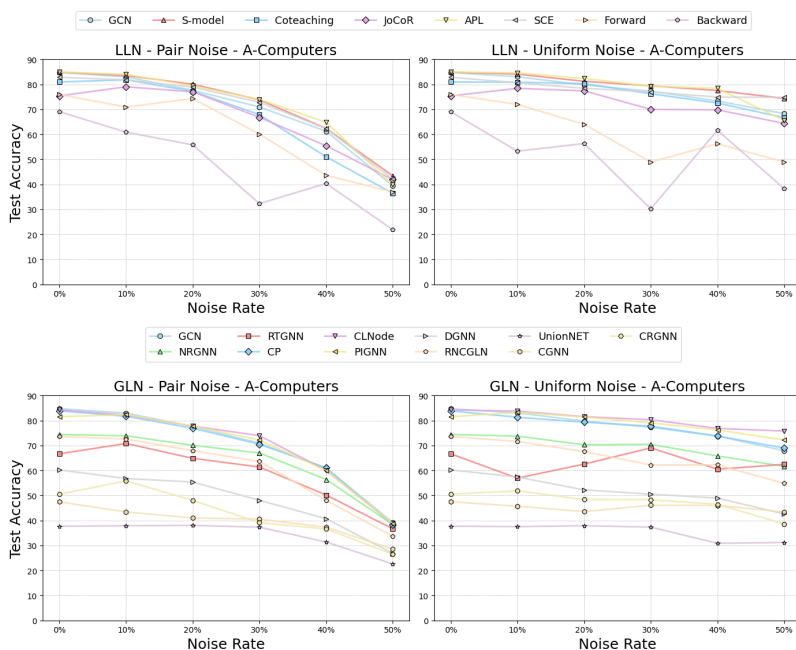


Figure A5: Test accuracy of LLN and GLN methods on Amazon-Computers dataset under different rate of pair and uniform noise, respectively (10 Runs).

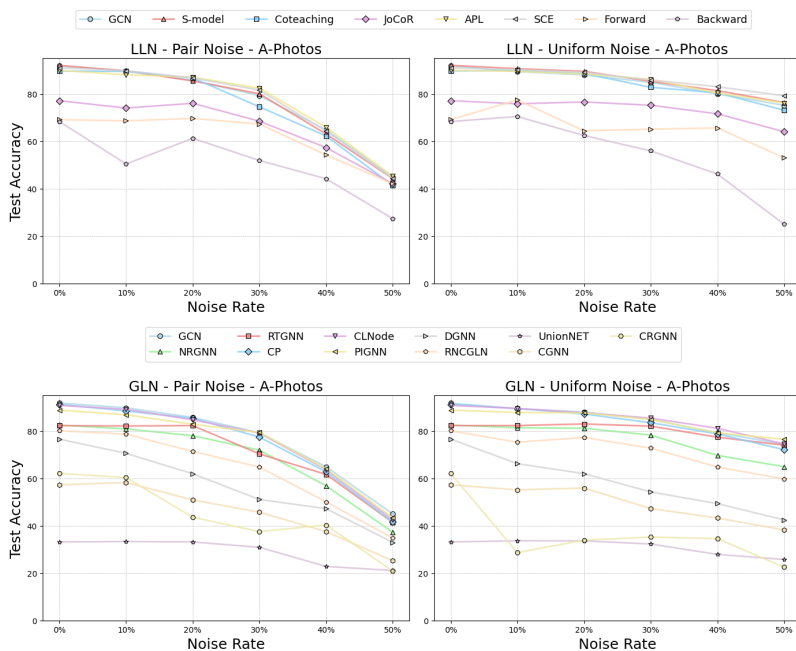


Figure A6: Test accuracy of LLN and GLN methods on Amazon-Photos dataset under different rate of pair and uniform noise, respectively (10 Runs).

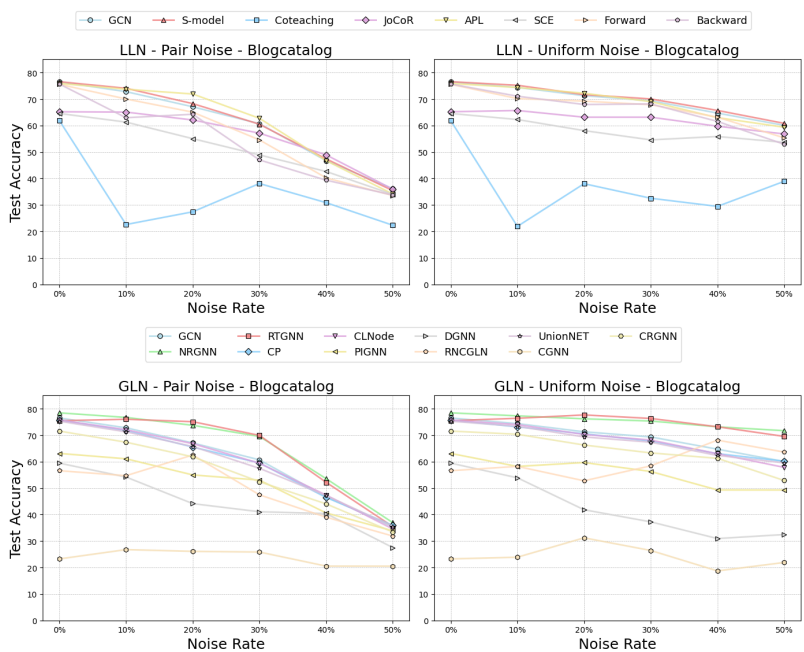


Figure A7: Test accuracy of LLN and GLN methods on Blogcatalog dataset under different rate of pair and uniform noise, respectively (10 Runs).

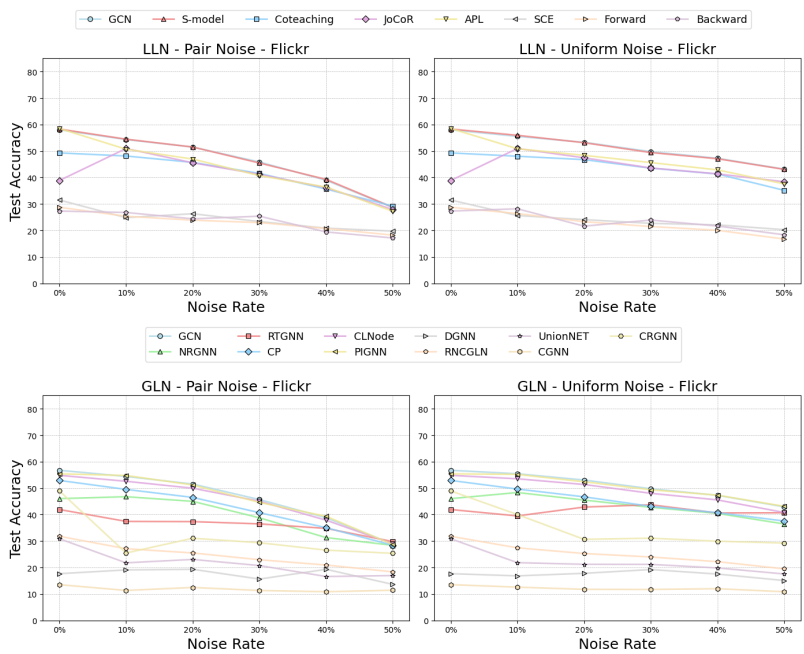


Figure A8: Test accuracy of LLN and GLN methods on Flickr dataset under different rate of pair and uniform noise, respectively (10 Runs).

### B.3 Additional results of time efficiency

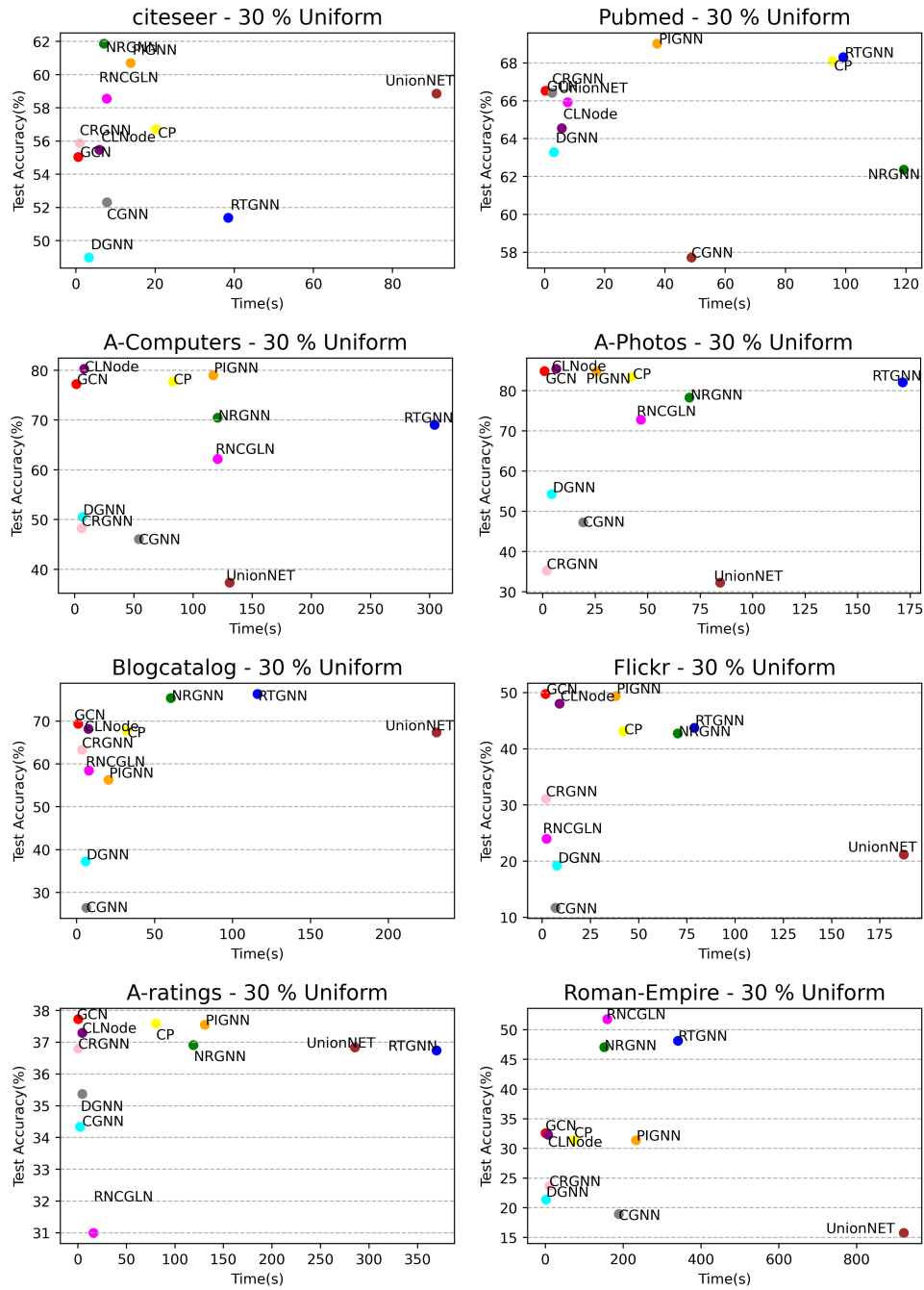


Figure A9: Time consumption and Test accuracy of different GLN methods on different datasets under 30% uniform noise (10 Runs).

## C Additional details of the Benchmark

### C.1 Datasets

Table A7: Overview of the datasets used in this study.

Dataset	# Nodes	# Edges	# Feat.	# Classes	# Homophily	Avg. # degree
Cora	2,708	5,278	1,433	7	0.81	3.90
Citeseer	3,327	4,552	3,703	6	0.74	2.74
Pubmed	19,717	44,324	500	3	0.80	4.50
Amazon-Computers	13,752	491,722	767	10	0.78	35.76
Amazon-Photos	7,650	238,162	745	8	0.83	31.13
DBLP	17,716	105,734	1,639	4	0.83	5.97
BlogCatalog	5,196	343,486	8,189	6	0.40	66.11
Flickr	7575	239738	12047	9	0.24	63.30
Amazon-Ratings	24,492	93,050	300	5	0.38	7.60
Roman-Empire	22,662	32,927	300	18	0.05	2.90

**Cora, Citeseer and Pubmed** [19] are citation networks that most commonly used in previous graph learning under label noise studies [1, 31, 24, 11, 18]. Each node represents a paper and each edge represents citation relationship between papers. Node features are 0/1-valued word vector indicating the absence/presence of the corresponding word from the dictionary. The label of each node is its category of research topic.

**Amazon-Computers and Amazon-Photo** [20] are co-purchase graphs extracted from Amazon, where each node represents a product, edges represent the co-purchased relationships between products. Features are bag-of-words vectors extracted from product reviews, labels of each node is its corresponding product category. These datasets were frequently used in robust graph learning under label noise studies [24, 12].

**Amazon-Ratings** [17] is derived from the Amazon product co-purchasing network metadata from SNAP Datasets. In this dataset, nodes represent products, while edges indicate products frequently bought together. Node features are calculated as the average of FastText embeddings for words in the product descriptions. The product ratings are categorized into five distinct classes as labels.

**Roman-Empire** [17] is derived from the Roman Empire article on English Wikipedia. The text was obtained from the English Wikipedia dump dated 2022.03.01. Each node in the dataset’s graph represents a (non-unique) word in the text. Node features are calculated by FastText embeddings. Two words are connected if they follow each other in the text or are linked in the dependency tree of a sentence. Nodes are labeled based on their syntactic roles and identified using spaCy. The 17 most frequent roles are considered unique classes, while others are grouped into the 18th class.

**DBLP** [15] is an author collaboration network in computer science, each node represents a document and edges represent their citation links. Features are word vectors and labels are category of research topic. This dataset was used in study [1].

**Blogcatalog** [26] is a social network formed by an online community, each node is a blogger and edges represent their relationships. The features of each node are derived from the keywords present in their blog descriptions, and the labels are selected from a collection of established categories that reflect the bloggers’ interests. This dataset was used in study [18].

**Flickr** [26] is a platform where users can share videos and images. User can follow each others thus form a social network. The feature of each node are generated from the user-specified tags, and labels represent the groups they have joined.



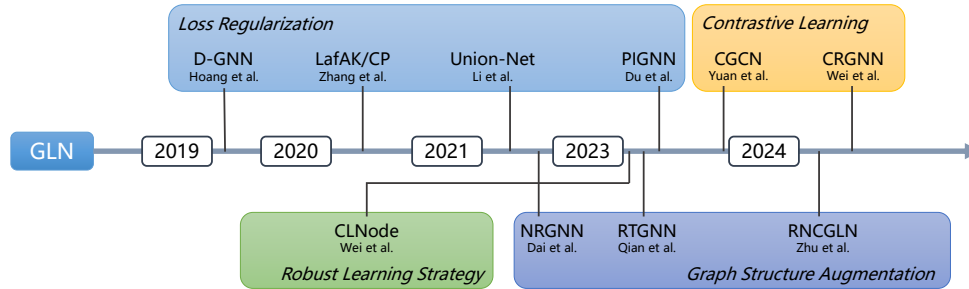


Figure A10: Timeline of GLN research. Existing GLN methods can be categorized into Loss regularization, Robust training strategy, Graph structure augmentation and Contrastive learning.

## C.2 Algorithms

### C.2.1 Graph Neural Networks with Label Noise

**NRGNN** [1] believe that since the labels on the graph are sparse, the falsely labeled nodes may affect the unlabeled nodes in its neighborhood, which make it difficult to receive supervision from correctly labeled nodes. To address these issues, NRGNN first connects nodes with similar features to create a refined graph. Based on this refined graph, precise pseudo-labels are generated, allowing unlabeled samples to receive more supervision from correctly labeled samples, thereby reducing the impact of noisy labels.

**RTGNN** [18] followed the work of NRGNN. The authors point out that although NRGNN emphasizes providing additional supervision for unlabeled nodes through link prediction, it does not distinguish between incorrectly labeled and correctly labeled nodes. Instead, it merely connects nodes with similar features indiscriminately, which may lead to the spread of influence from incorrectly labeled nodes. To solve this problem, the authors propose RTGNN, which, building based on NRGNN, uses the small-loss criterion from Co-teaching [4] to further distinguish between trustworthy and untrustworthy nodes, and corrects the labels of some untrustworthy nodes, mitigates some of the shortcomings of NRGNN.

**CP** [31] studied the impact of adversarial label-flipping attacks on the generalization ability of Graph Convolutional Networks (GCNs). To counteract label-flipping attacks, the authors proposed a defense framework named CP, which uses community labels as high-level signals to guide the node classification task. The CP framework includes a constraint with community information to prevent overfitting to the flipped noisy labels. The use of community labels is motivated by their similarity to the output of GCNs.

**D-GNN** [14] obtains a label noise robust Graph Neural Network by adopting backward loss correction [16] on GIN [25] backbone, which estimates the unbiased loss on clean labels.

**RNCGLN** [32] aim to simultaneously mitigate graph and label noise issues. To achieve this, it first use graph contrastive loss to conduct local graph learning, and adopt multi-head self-attention mechanism to learn node representation from a global perspective. Then utilize pseudo graphs and pseudo labels to deal with graph noise and label noise, respectively.

**CLNode** [24] adopt a curriculum learning strategy to mitigate the impact of label noise. To be specific, it first utilize a multi-perspective difficulty measurer to accurately measure the quality of training nodes. Then employ a training scheduler that selects appropriate training nodes to train GNN in each epoch based on the measured qualities. The authors demonstrated this method enhances the robustness of backbone GNN to label noise.

**PIGNN** [2] enhances the GNN’s resistance to label noise by introducing additional pair-wise labels. The motivation is pair-wise labels are more robust than node-wise labels. In authors’ definition, a pair interaction label is 1 if the nodes have the same label, and 0 otherwise, and they designed a PI label estimation method based on the similarity of node embeddings. During training, the estimated PI label serves as the confidence level for the node classifier’s predictions, thereby constraining the training process of the node classifier. This method performs well on homophilic graphs but poorly on heterophilic graphs.

**Union-NET** [11] tries to limit the gradient passing process of mislabeled samples through neighborhood labeling, like a kind of neighborhood voting with node representation similarity weighting. A GNNs first generates node representations and predicted labels. Context nodes are then aggregated using random walks, and an attention mechanism calculates class probability distributions. This guides a reweighting scheme to minimize the impact of noisy labels. Labels are corrected by aligning them with the most consistent context labels, and a KL-divergence loss maintains alignment with the prior distribution. The training involves pre-training the GNN and updating model parameters with a combined loss function, ensuring robust training and effective label correction.

**CGNN** [30] addresses label noise in GNNs by combining neighborhood-based label correction and contrastive learning. It utilizes message passing neural networks to update node representations, integrating graph contrastive learning for consistent representations across augmented graph views. Finally, CGNN employs an MLP for prediction distributions and iteratively corrects noisy labels by comparing them with their neighbors and choosing the most labels.

**CR-GNN** [10] introduces contrastive learning to enhance GNNs robustness in the face of sparse and noisy node labels. Through techniques like feature masking and edge dropping, CR-GNN preserves node semantics while generating augmented views. Contrastive loss captures local structural information and mitigates noisy label effects, while dynamic cross-entropy loss addresses overfitting and adversarial vulnerabilities. Also, cross-space consistency ensures semantic alignment between embeddings.

### C.2.2 Learning with Label Noise methods

**S-model** [3] adds a noise adaptation layer that models the transition pattern of noisy labels on true labels. In the training procedure, this layer is parameterized by bias terms and allows the network to learn both the classifier and noise model simultaneously. In the test procedure, the noise adaptation layer is removed, which enables the network to predict true labels more effectively.

**Co-teaching** [4] works by simultaneously training two deep neural networks (DNNs), each of which selects a certain number of small-loss samples from them and passes these samples to the other for further training. It assume that mislabeling typically leads to larger losses and thus is less likely to be selected, and then each network selects the samples that perform best on its own with lower loss. This peer-to-peer training mechanism helps to reduce the effect of noisy labels, as both networks focus on more reliable data.

**JoCoR** [23] utilizes consistency maximization to deal with the noisy labels. Instead of using hard sampling, two different classifiers are made to converge in their predictions through explicit regularization. Specifically, the two classifiers are trained by a joint loss function to minimize the differences between them. During the training process, these two classifiers update their parameters at the same time and are jointly trained by means of a pseudo-twin network. The loss function consists of a supervised learning loss and a contrast loss, where the contrast loss is used to maximize the agreement between the two classifiers.

**SCE** [22] enhances the robustness of a model in the presence of noisy labels by combining a noise tolerance term with the standard cross entropy (CCE) loss. Inspired by the Kullback-Leibler scattering symmetry, SCE incorporates the reverse cross entropy loss, a noise tolerance term, and combines it with the standard CCE loss to improve the model's ability to tolerate noisy data. This approach not only retains the advantages of the CCE loss, but also significantly improves the generalization performance in noisy environments through symmetry processing and noise tolerance.

**Forward correction** [16] corrects the sample loss by linearly combining the softmax output of the DNN before applying the loss function. During the forward propagation process, the estimated label transfer probability is multiplied with the softmax output to obtain the corrected loss value. In this way, the softmax output of each sample is first combined with the corresponding transfer probability, and then the loss function is applied, which improves the robustness of the model in noisy labeling environments.

**Backward correction** [16] adjusts the loss for each sample by multiplying the estimated label transfer probability with the output of the specified DNN. The learning of the label transfer probability is decoupled from the learning of the model, and the label transfer matrix is first approximated using the softmax output of the DNN in the uncorrected loss case. Then, when retraining the DNN, the original loss is corrected based on the estimated matrix. The correction loss is computed by linearly combining

the loss values for each observable label, where the coefficients are the transfer probabilities from each observable label to the target label.

## D Package

We have developed an open-source software package NoisyGL, which provides a comprehensive and unbiased platform for evaluating GLN algorithms and advancing future research. The code structure of NoisyGL is well-designed to ensure fair experimental setups for different algorithms, easy reproducibility of experimental results, and support for flexible assembly of models for experiments. NoisyGL includes the following key modules. The Config module consists of the files that define the necessary hyperparameters and settings. The Dataset module is used to load datasets, and the label Contaminator modifies the raw data to create a contaminated graph. The Base-predictor serves as the base class for various reproduced LLN/GLN predictors, and the LLN/GLN Predictor evaluates the contaminated graph to predict performance.



As shown in the Figure A11, the code structure is well-organized to ensure fair experimental settings across algorithms, easy reproduction of experimental results, and convenient trials on flexibly assembled models. Given a specific dataset and config file, a solver will return the learned structure and the task performance. For more details and updated features, please refer to our GitHub repository.

**General Experimental Settings.** We endeavor to follow the original implementations of the various GLN methods in their associated papers or source code. To this end, we integrate the different options into a standardized framework as shown in Figure. In this way, we can ensure consistency and comparability of experiments, allowing the performance of different GLN methods to be fairly evaluated on the same platform. We run most experiments on NVIDIA Geforce RTX 3090 GPU with 24 GB memory, the out-of-memory error during the training is reported as N/A in Appendix A. For the two large datasets, Amazon Ratings and Roman Empire, we run these experiments on NVIDIA A100 with 80GB memory.

**Hyperparameter.** We performed manual hyperparameter tuning to ensure an unbiased evaluation of these GLN methods. The hyperparameter search space for all methods is shown in Table A8. For details on the meaning of these hyperparameters, please refer to their original papers. Through exhaustive tuning and setup, we strive for the best performance of each method under different configurations, thus ensuring the accuracy and fairness of the evaluation.

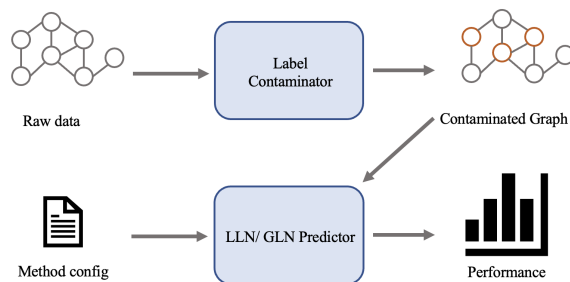


Figure A11: The structure of NoisyGL. The Raw data is processed by the Label Contaminator to introduce label noise, resulting in a Contaminated Graph. This contaminated graph, along with the Method config, is then input to the LLN/GLN Predictor, which evaluates performance metrics based on the specified method configuration.

Table A8: Hyper-parameter search space of all implemented GLN methods.

Algorithm	Hyper-parameter	Search Space
General Settings	learning rate	1e-1, 5e-2, 1e-2, 5e-3, 1e-3, 5e-4, 1e-4, 5e-5
	weight decay	5e-2, 5e-3, 5e-4, 5e-5
	layer number	2, 3, 4, 5
	hidden size	16, 32, 64, 128
NRGNN [1]	$\alpha$	0.01, 0.02, 0.03
	$\beta$	0.01, 0.1, 1, 10
RTGNN [18]	$\tau$	0, 0.05, 0.1, 0.2
	$\lambda$	0.05, 0.1, 0.2
	$\alpha$	0.03, 0.1, 0.3, 1
LAFAK/CP [31]	$\lambda$	0.1, 0.2, 0.3
CLNode [24]	$\lambda_0$	0.25, 0.5, 0.75
	$T$	50, 100, 150
PIGNN [2]	N/A	N/A
DGNN [14]	N/A	N/A
RNCGLN [32]	$\alpha$	$10^{-3}, 10^{-2}, \dots, 10^3$
UnionNET [11]	$\alpha$	0.1, 0.5, 1.0
	$\beta$	0.1, 0.5, 1.0
CGNN [30]	$\gamma$	0.6, 0.7, 0.8, 0.9, 0.95
	$\omega$	0.6, 0.7, 0.8, 0.9, 0.95
CRGNN [10]	$\alpha$	0.1, 0.2, 0.3, $\dots$ , 1
	$\beta$	0.1, 0.2, 0.3, $\dots$ , 1

## E Reproducibility

All of NoisyGL’s experimental results are highly reproducible. We provide more detailed information on the following aspects to ensure the reproducibility of the experiments.

**Accessibility.** You can access all datasets, algorithm implementations, and experimental configurations in our open source project <https://github.com/eaglelab-zju/NoisyGL> without a personal request.

**Dataset.** The datasets used are publicly available. The Cora, Citeseer, and Pubmed datasets are accessible online and are used under the Creative Commons 4.0 license. The BlogCatalog and Flickr datasets were originally published by [26] and further processed in subsequent studies. To the best of our knowledge, these datasets do not have a specific license. The DBLP dataset can be found in [15] and is released under the MIT license. All of these datasets are licensed by the authors for academic research and do not contain any personally identifiable information or offensive content.

**Documentation and uses.** We’ve dedicated ourselves to providing users with comprehensive documentation, guaranteeing a smooth experience with our library. Our code includes ample comments to enhance readability. Furthermore, we furnish all essential files to replicate experimental outcomes, which also serve as illustrative guides on library utilization. Running the code is straightforward; users need only execute the '.py' files with specified arguments like data, method, and GPU.

**License.** We use an MIT license for our open-sourced project.

**Code maintenance.** We are dedicated to maintaining our code through continuous updates, actively engaging with user feedback, and addressing any issues promptly. Additionally, we are eager to receive contributions from the community to improve our library and benchmark algorithms. However, we will uphold rigorous version control measures to uphold reproducibility standards during maintenance procedures.