

# UAG: Uncertainty-aware Attention Graph Neural Network for Defending Adversarial Attacks

Boyuan Feng, Yuke Wang, Yufei Ding

University of California, Santa Barbara  
{boyuan, yuke\_wang, yufeidng}@ucsb.edu

## Abstract

With the increasing popularity of graph-based learning, graph neural networks (GNNs) emerge as the essential tool for gaining insights from graphs. However, unlike the conventional CNNs that have been extensively explored and exhaustively tested, people are still worrying about the GNNs' robustness under the critical settings, such as financial services. The main reason is that existing GNNs usually serve as a black-box in predicting and do not provide the uncertainty on the predictions. On the other side, the recent advancement of Bayesian deep learning on CNNs has demonstrated its success of quantifying and explaining such uncertainties to fortify CNN models. Motivated by these observations, we propose UAG, the first systematic solution to defend adversarial attacks on GNNs through identifying and exploiting hierarchical uncertainties in GNNs. UAG develops a Bayesian uncertainty technique to explicitly capture uncertainties in GNNs and further employs an uncertainty-aware attention technique to defend adversarial attacks on GNNs. Intensive experiments show that our proposed defense approach outperforms the state-of-the-art solutions by a significant margin.

## Introduction

As the emerging trend of extending deep learning from Euclidean data (*e.g.*, images) to non-Euclidean data (*e.g.*, graphs), graph neural network (GNN) (Xu et al. 2019b; Kipf and Welling 2017; Veličković et al. 2018) wins lots of attentions from both research and industrial domains. Compared with the conventional graph-learning approaches (*e.g.*, random walk (Grover and Leskovec 2016; Perozzi, Al-Rfou, and Skiena 2014), and graph Laplacians (Luo et al. 2011, 2009; Cheng et al. 2018)), GNNs excel at both computation efficiency and runtime performance for various tasks, such as the node classification (Kaspar and Horst 2010; Gibert, Valveny, and Bunke 2012; Duran and Niepert 2017) and link prediction (Chen, Li, and Huang 2005; Kunegis and Lommatzsch 2009; Tylenda, Angelova, and Bedathur 2009). Despite the stunning success, people still concern about the robustness of GNNs, especially in some safety-critical domains (*e.g.*, financial services and medicinal chemistry). Existing work (Zügner and Günnemann 2019; Xu et al. 2019a; Waniek et al.

2018) has shown that GNNs are sensitive to small perturbations on the topology and the node features, which motivates our work for defending adversarial attacks.

The most recent work, RGCN (Zhu et al. 2019), improves GNN robustness with a simple strategy that replaces deterministic GNN features with a Gaussian distribution and measures the variance in the intermediate feature vectors. However, it assumes fixed GNN weights without quantifying the uncertainty from GNN models and does not consider the uncertainty from the graph topology, leading to unsatisfactory accuracy under severe attacks. We believe the key to improve the GNN robustness is to develop a powerful technique to quantify and exploit uncertainties from various sources to absorb the effect of adversarial attacks.

In this paper, we focus on exploring the benefits of explicitly quantifying GNN uncertainty to defend GNNs against adversarial attacks. In particular, there are two major types of uncertainties in GNNs – the *model* uncertainty and the *data* uncertainty. The former refers to the uncertainty in model parameters to tell whether the selected parameters can best suit the distribution of the collected data. The latter refers to the uncertainty in the noisy data collection, coming from either the noises in the data collection process or the adversarial attacks. However, exploring these uncertainties is non-trivial since there are several challenges to overcome:

1. **Uncertainty Measurement:** How to explicitly measure the uncertainty of GNNs?
2. **Robustness:** How to effectively incorporate the measured uncertainty into existing GNNs for defending adversarial attacks?

To tackle these challenges, we propose the first Bayesian-based uncertainty guided approach to defend the GNN effectively. First, we develop a *Bayesian uncertainty technique* based on the powerful Bayesian framework to capture these uncertainties from different sources. Intuitively, we measure the uncertainty value for individual nodes where a higher uncertainty usually indicates a lower prediction accuracy. Then, we design an *uncertainty-aware attention technique* to dynamically adjust the impact of one node towards its neighboring nodes according to its uncertainty. In particular, for nodes with high uncertainty that may have been attacked,

we restrict its feature propagation towards neighboring nodes in order to absorb the attack impact.

In short, we summarize our contributions as follows:

- We identify two types of uncertainties in GNNs (*i.e.*, model uncertainty and data uncertainty) and propose a Bayesian uncertainty technique to explicitly capture both types of uncertainties.
- We introduce an uncertainty-aware attention technique to defend the adversarial attack by assigning less impact (weights) on nodes with high uncertainty, thus, mitigating their impact on the final prediction.
- Rigorous experiments and studies across various datasets on mainstream GNNs show that our proposed defense approach outperforms the state-of-the-art RGCN by a significant margin.

## Related Work

**Graph Neural Network** Graph Neural Networks (GNNs) are now becoming a major way of gaining insights from the graph structures. It generally includes several graph convolutional layers, each of which consists of a neighbor aggregation and a node update step. The most common graph convolutional layer (Kipf and Welling 2017) computes the embedding for node  $v$  at layer  $k+1$  based on node embedding at layer  $k$ , where  $k \geq 0$ .

$$h_v^{(k+1)} = \sigma\left(\sum_{u \in \bar{N}(v)} \frac{1}{c_u c_v} h_u^{(k)} \cdot W^{(k)}\right) \quad (1)$$

As shown in Equation 1,  $h_v^{(k)}$  is the embedding vector for node  $v$  at layer  $k$ ,  $W^{(k)}$  is the GNN weight at layer  $k$ , and  $\bar{N}(v) = N(v) \cup v$  is the set of node  $v$  and its neighboring nodes.  $c_u$  and  $c_v$  are fixed values determined by the degree of node  $u$  and  $v$  and will be omitted in following sections for notation simplicity. Intuitively, the graph convolution layer aggregates information across nodes by averaging features in nearby nodes. More advanced GNNs utilize different aggregation methods. For example, GAT (Veličković et al. 2018) aggregates node features with weighted average based on the cosine similarity between node features.

**Graph Adversarial Attacks and Defense** Existing works have explored the robustness of the GNNs in two opposite but closely related directions, GNN attacks, and GNN defense. On the attack side, existing GNN attacks can be broadly classified into two major categories, *poisoning* (Zügner, Akbarnejad, and Günnemann 2018; Zügner and Günnemann 2019) and *evasion* (Dai et al. 2018), depending on the time they happen. The former (poisoning attack) happens during the training time of the GNNs through modifying training data and the latter (evasion attack) takes place during the GNN inference time by changing test data samples. Our work is orthogonal and complementary to these existing GNN attack research, since 1) our goal is to minimize the impact of these GNN attacks by incorporating model and data uncertainties during the GNN computation; 2) our defense-oriented research may potentially motivate more diverse adversarial attacks tailored for GNNs.

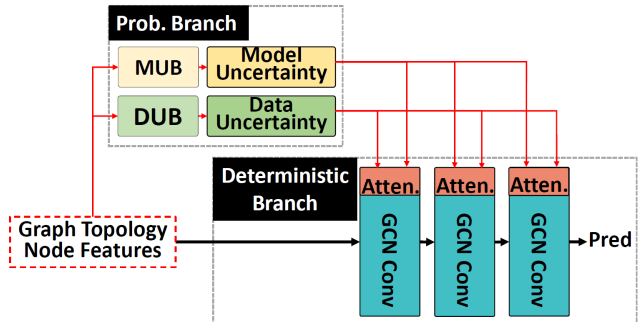


Figure 1: Overview of UAG.

On the defense side, RGCN (Zhu et al. 2019) proposes a novel model to make GCN immune from adversarial attacks by leveraging Gaussian distributions to reduce the impact of GNN attacks. Different from RGCN, our UAG is the first work to identify and quantify how adversarial attacks affect GNN’s performance – model and data uncertainties that take both model (weight) and data (topology and embedding features) into consideration. And we further exploit such uncertainty information by incorporating novel techniques to facilitate the defense.

**Bayesian Neural Network and Uncertainty** Many research efforts have been made towards developing Bayesian Neural Network to measure uncertainty (Gal 2016) in computer vision (Kendall and Gal 2017; Gal and Ghahramani 2016a,b; Alex Kendall and Cipolla 2017), natural language processing (Xiao and Wang 2019), and time series analysis (Zhu and Laptev 2017). These works usually focus on convolutional neural networks and use Bayesian Neural Network as a regularization technique. Some recent contributions (Zhang et al. 2019; Hasanzadeh et al. 2020) extend the Bayesian Neural Network to the graph domain as a stochastic regularization technique. These works aim to solve the over-smoothing problem in GNNs and do not explicitly quantify the uncertainty. To the best of our knowledge, we are the first to exploit the uncertainty in the graph domain to defend adversarial attacks.

## Methodology

The overview of UAG is presented in Figure 1. UAG takes two inputs – the adjacency matrix  $A \in \mathcal{R}^{N \times N}$  and the node features  $X \in \mathcal{R}^{N \times D}$ , where  $N$  is the number of nodes and  $D$  is the dimension of node features. There are two main branches in UAG – the probabilistic branch (including the Model Uncertainty Branch (MUB) and the Data Uncertainty Branch (DUB)) and the deterministic branch, where the architecture and weights are different across branches. Given the graph data  $(A, X)$ , the probabilistic branch measures the node-wise uncertainty  $U = [U_M, U_D] \in \mathcal{R}^{N \times 2}$  from the GNN model weights and the graph data. Here, the probabilistic branch adopts a novel Bayesian uncertainty technique to measure the uncertainty for each node due to the adversarial attacks. The deterministic branch takes the measured node-wise uncertainty  $U$  and the graph data  $(A, X)$  to generate the node classification results  $Y \in \mathcal{R}^N$ . It contains an uncertainty-aware attention technique to adaptively ad-

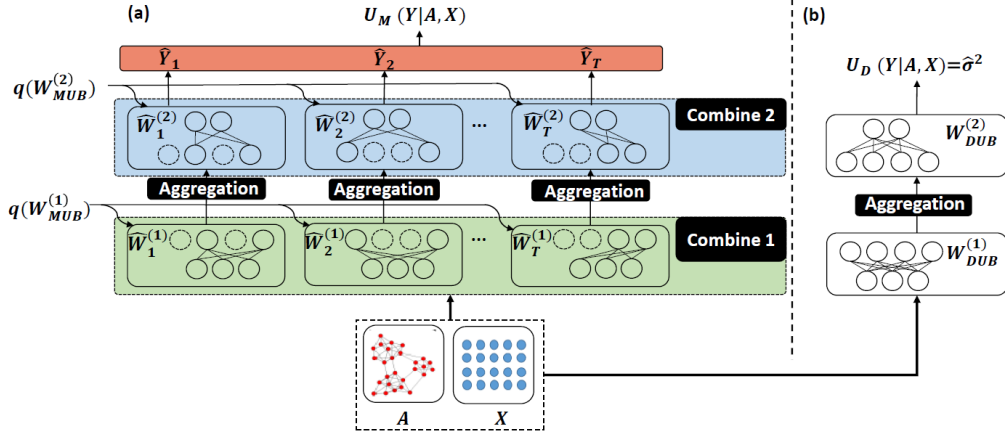


Figure 2: Bayesian Uncertainty Technique Overview. (a) Model Uncertainty Branch (MUB); (b) Data Uncertainty Branch (DUB)

just the edge attention during the inference to defend against adversarial attacks.

### Bayesian Uncertainty Technique

Bayesian uncertainty technique measures two sources of uncertainties – the model uncertainty  $U_M \in \mathcal{R}^N$  and the data uncertainty  $U_D \in \mathcal{R}^N$ . Formally, following the law of total variance, the uncertainty in the prediction  $\hat{Y}$  is

$$\begin{aligned} \text{Var}(Y) &= \text{Var}(E[Y|A, X]) + E[\text{Var}(Y|A, X)] \\ &= U_M(Y|A, X) + U_D(Y|A, X) \end{aligned} \quad (2)$$

Here, we use a model uncertainty branch (MUB) and a data uncertainty branch (DUB) to access  $U_M(Y|A, X)$  and  $U_D(Y|A, X)$ , respectively.

**Model Uncertainty.** The model uncertainty  $U_M$  measures the uncertainty in the mapping process  $E[Y|A, X]$  due to model weight selection. We use a 2-layer GCN to quantify the model uncertainty. Instead of using fixed weights, we utilize a probability distribution to describe the uncertainty from the model weights, as illustrated in Figure 2(a). Formally, given the graph data  $(A, X)$  and the partial label  $Y \in \mathcal{R}^{N_L}$  with  $N_L$  as the number of labeled nodes, we first train the weight posterior distribution  $p(W|A, X)$ . Then, we conduct the prediction mapping procedure as

$$p(Y|X, A) = \int_W p(Y|W, A, X) p(W|A, X) dW \quad (3)$$

Since the exact inference is intractable, we adopt the MC dropout variational inference (Gal and Ghahramani 2016a) method by multiplying a deterministic model weight  $W_{MUB}$  with a random variable  $B$  following the Bernoulli distribution. This provides  $q(W)$  as an approximation to the true posterior  $p(W|A, X)$ . In particular, the model weights  $W$  follows

$$\begin{aligned} q(W) &\sim B \odot W_{MUB} \\ P(B) &\sim \text{Bernoulli}(p) \end{aligned} \quad (4)$$

where  $\odot$  is the Hadamard product, and  $p$  is a hyperparameter ( $=0.8$  by default in our evaluation). During training, we can train the weight by minimizing the cross-entropy loss

$$L_{model} = -\frac{1}{T} \sum_{t=1}^T \log p(\hat{Y}_t | \hat{W}_t, A, X) + \frac{1-p}{2T} \|W_{MUB}\|^2 \quad (5)$$

where  $\hat{W}_t$  is sampled from  $q(W)$ ,  $\hat{Y}_t$  is the prediction under sampled weight  $\hat{W}_t$ , and  $T$  is the number of samples during the MC dropout variational inference. During inference, we perform the Monte Carlo integration:

$$E(Y|A, X) = \frac{1}{T} \sum_{t=1}^T \hat{Y}_t \quad (6)$$

where  $\hat{Y}_t \in \mathcal{R}^L$  is the prediction after the softmax layer, and  $L$  is the number of classes in the graph data.

Here, we observe that applying Bernoulli distribution at different granularities leads to different probabilistic interpretation. To provide a comprehensive measurement on the model uncertainty, we apply dropout independently for individual GNN layers, channels, nodes, and edges

$$W_{uv}^{(k)} = B_{uv}^{(k)} \odot B_u^{(k)} \odot W^{(k)} \quad (7)$$

where  $W^{(k)} \in \mathcal{R}^{f_k \times f_{k+1}}$  is the GNN weights at the  $k^{th}$  layer,  $B_{uv}^{(k)} \in \mathcal{R}$  determines the dropout on the edge-level,  $B_u^{(k)} \in \{0, 1\}^{f_k}$  drops the weight at the channel level, and  $f_k$  is the number of feature channels at layer  $k$ . From the perspective of individual nodes, we have

$$\begin{aligned} h_v^{(k+1)} &= \sigma\left(\sum_{u \in \bar{N}(v)} h_u^{(k)} \cdot W_{uv}^{(k)}\right) \\ &= \sigma\left(\sum_{u \in \bar{N}(v)} h_u^{(k)} \cdot (B_{uv}^{(k)} \odot B_u^{(k)} \odot W^{(k)})\right) \end{aligned} \quad (8)$$

where  $\bar{N}(v) = N(v) \cup v$ . Noting that this dropout also provides a Bayesian view of dropping edges or nodes when either  $B_{uv}^{(k)} = 0$  or  $B_u^{(k)} = 0$ .

Given the Bayesian framework on GNNs, we can measure the model uncertainty as the variance in predictions

$$\begin{aligned} U_M(Y|A, X) &= \text{Var}(Y|A, X) \\ &= E(Y^2|A, X) - [E(Y|A, X)]^2 \\ &= \frac{1}{T} \sum_{t=1}^T \hat{Y}_t^2 - [E(Y|A, X)]^2 \end{aligned} \quad (9)$$

Here, we additionally apply a reduce operation to transform the  $L$ -dimension vector  $U_M$  to a scalar value as the model uncertainty.

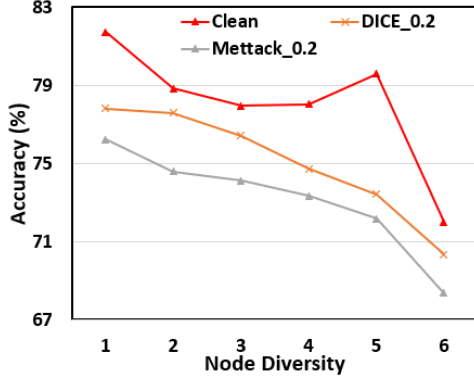


Figure 3: Relationship between accuracy and node diversity. *DICE\_0.2* and *Mettack\_0.2* indicates perturbing 20% edges with DICE and Mettack, respectively.

**Data Uncertainty** The data uncertainty  $U_D$  measures the prediction noise intrinsic to the data inputs. There are two standard approaches to identify the data uncertainty on each node. One approach utilizes the maximum predicted probability to measure confidence in the prediction. However, this confidence comes as a side effect of the model training and lacks sophisticated probabilistic interpretation. The other approach from the CNN domain predicts the uncertainty value based on the image inputs. However, naively borrowing this approach into the GNN domain focuses only on the node features and fails to exploit the important topology information in the graph data.

Instead, we aim to capture the data uncertainty that is intrinsic to the graph topology in terms of node diversity  $Div_{node}^{(k)}$ , defined as the number of different labels in the node’s  $k$ -hop neighbors. Our key observation is that adversarial attacks on graph data usually increase the node diversity and add edge connections between nodes with different labels. For example, DICE attack (delete edges internally, connect externally) (Wanik et al. 2018) exploits the node label information to increase node diversity by deleting edges between nodes with the same label and adding edges between nodes with different labels. Figure 3 shows the accuracy among nodes that have node diversity larger than various thresholds. We observe that the accuracy usually decreases significantly as the node diversity increases, which holds for both the clean graph data and the attacked graph data from various attacking algorithms.

To this end, we explicitly measure the data uncertainty by treating the prediction as a Gaussian distribution and setting the variance to be the node diversity, as illustrated in Figure 2(b). Formally, we have

$$Y \sim N(\hat{\mu}(A, X), \hat{\sigma}^2(A, X)) \quad (10)$$

where  $\hat{\mu}$  and  $\hat{\sigma}^2$  are the predicted label and node diversity, respectively. Here, we parameterize the  $\hat{\mu}$  and  $\hat{\sigma}^2$  with the adjacency matrix  $A$  and the node feature  $X$  and use a 2-layer GCN to predict their values. During inference, we will use the  $\hat{\sigma}^2(A, X)$  as the data uncertainty

$$U_D(Y|A, X) = \hat{\sigma}^2(A, X) \quad (11)$$

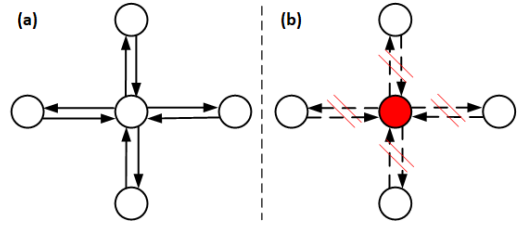


Figure 4: Illustration of UAT. (a) Aggregation on the clean graph; (b) aggregation when the red node is attack.

To train the data uncertainty, we have two losses on the labeled nodes and unlabeled nodes, respectively. On the labeled nodes, we focus on the ground truth labels and have a KL-divergence that requires the predicted distribution to match with the ground truth distribution

$$L_1 = KL(N(\hat{\mu}(A, X), \hat{\sigma}^2(A, X))|N(Y, \sigma^2)) \quad (12)$$

where the  $Y$  comes from the ground truth label and  $\sigma^2$  measures the node diversity  $Div_{node}^{(k)}$  in the graph data. When computing the node diversity, we utilize only labeled node in the clean graph data. In particular, for a given node, we first collect all labeled 2-hop neighboring nodes and then count the number of distinct labels. On the unlabeled nodes, similar to the unsupervised learning (Bojchevski and Günnemann 2018) on graph data, we focus on the graph topology and adopt an energy-based unsupervised loss

$$L_2 = \sum_i \sum_{k < l} \sum_{j_k \in N_{ik}} \sum_{j_l \in N_{il}} (E_{ij_k}^2 + \exp^{-E_{ij_l}}) \quad (13)$$

$$E_{ij} = D_{KL}(N(\hat{Y}_j, \hat{\sigma}_j^2) || N(\hat{Y}_i, \hat{\sigma}_i^2))$$

Assuming that a node tends to have a similar label with neighboring nodes, this loss implicitly captures the node diversity by forcing higher feature similarities in neighboring nodes.

### Uncertainty-aware Attention Technique

The Uncertainty-aware Attention Technique (UAT) adaptively adjusts the edge attention during the inference to defend against adversarial attacks, as illustrated in Figure 4. We equip a 2-layer GCN with our UAT. On a clean graph (Figure 4a), we adopt edge aggregation similar to existing GNNs that each node aggregates and propagates information across neighboring nodes. On an attacked graph (e.g., the red node in Figure 4b), UAT adaptively limits the information propagation between the attacked node and other nodes. While existing works (Zügner, Akbarnejad, and Günnemann 2018) have shown that attacking one node in the graph can also lead to the wrong prediction on other nodes, UAT mitigates it by reducing the impact from attacked nodes to remaining nodes.

Formally, given the feature  $h_u^{(k+1)}$  for each node  $u$  at the  $k + 1$  layer, we compute an attention  $Att_\tau(u)$  for each node  $u$  and compute each GNN layer as

$$h_v^{(k+1)} = \sigma\left(\sum_{u \in \tilde{N}(v)} Att_\tau^{uv} \cdot h_u^{(k)} \cdot W^{(k)}\right) \quad (14)$$

$$Att_\tau^{uv} = \min(Att_\tau(u), Att_\tau(v))$$

where  $\tau \in \{M, D\}$  indicates whether we are using model uncertainty or the data uncertainty, each node embedding  $h_u^k$  is weighted by an attention. Here, we use attention from both nodes  $u$  and  $v$  to decide the attention value on the edge. Note that the deterministic branch focuses on improving accuracy and utilizes independent weight from the probabilistic branch, which focuses on capturing uncertainties. We design two attentions to measure the model uncertainty and the data uncertainty, respectively

$$\begin{aligned} Att_\tau(u) &= \exp(-\zeta \cdot U_{\tau,u}) \\ \zeta &= \alpha_\tau \cdot \exp(-\beta_\tau \cdot Range(U_\tau)) \end{aligned} \quad (15)$$

where  $U_{\tau,u}$  measures the uncertainty on node  $u$ , and  $\alpha_\tau > 0$  and  $\beta_\tau > 0$  are two hyperparameters to adjust the impact from uncertainty. Intuitively, a larger uncertainty  $U_{\tau,u}$  on a node  $u$  leads to lower weight in the information propagation along with the graph topology. Here, we additionally utilize a  $Range(U_\tau)$  operation to measure the global uncertainty diversity in order to absorb the uniform uncertainty scale change on all nodes under diverse attacks. In particular, we collect  $U_\tau$  for all nodes and measure the  $Range(U_\tau)$  as the absolute difference between the first and the third quartiles. We have investigated several functions to combine the data uncertainty and the model uncertainty, and find out a simple minimal combination can already lead to good performance

$$Att_{Both}(u) = \min(Att_M, Att_D) \quad (16)$$

Intuitively, we restrict the information propagation from one node when it shows either a high model uncertainty or a high data uncertainty.

## Evaluation

In this section, we evaluate UAG on three popular datasets and compare with three baselines to show its effectiveness.

### Experiment Environments

**Datasets** We select the most typical datasets (*Cora*, *Citeseer*, and *Pubmed*) used by many GNN papers (Kipf and Welling 2017; Xu et al. 2019b; Hamilton, Ying, and Leskovec 2017) to evaluate our UAG. In these datasets, the node represents documents, edge refers to citations, and each node has its own associated bag-of-word features. Table 1 summarizes the details of these datasets. We follow the common data split by selecting 10% nodes as the training dataset, 10% nodes as the validation dataset, and 80% nodes as the testing dataset.

**Baselines** *Graph Convolutional Network (GCN)* (Kipf and Welling 2017) is one of the most popular GNN architectures. It has been widely adopted in node classification, graph classification, and link prediction tasks. Besides, it is also the key backbone network for many other GNNs, such as *GraphSage* (Hamilton, Ying, and Leskovec 2017), and differentiable pooling (*Diffpool*) (Ying et al. 2018). *Graph Attention Network (GAT)* (Xu et al. 2019b), another typical type of GNN, aims to distinguish the graph-structure that cannot be identified by GCN. GAT differs from GCN in its aggregation function, which assigns different weights to different nodes during the aggregation. *Robust GCN*

Dataset	#Vertex	#Edge	#Dim	#Class
Citeseer	3,327	9,464	3,703	6
Cora	2,708	10,858	1,433	7
Pubmed	19,717	88,676	500	3

Table 1: Datasets for Evaluation.

(**RGCN**) (Zhu et al. 2019), leverages the Gaussian distributions for node representations to amortize the effects of adversarial attacks.

**Attack Methods** *Random Attack* is a popular attack method that randomly adds fake edges into the graph dataset without considering the label of nodes. *DICE Attack* (delete edges internally, connect externally) (Wanick et al. 2018) exploits the node label information to increase node diversity by deleting edges between nodes with the same label and adding edges between nodes with different labels. *Metattack* (Zügner and Günnemann 2019) is another representative attack that adopts a meta-learning approach to reason the loss change by iteratively perturbing individual edges and features.

**Platforms.** We implement UAG based on PyTorch Geometric (Fey and Lenssen 2019). We evaluate UAG on a Dell Workstation T7910 (Ubuntu 18.04) with an Intel Xeon CPU E5-2603, 64 GB memory, and an NVIDIA 1080Ti GPU with 12GB memory.

### Overall Performance

In this section, we demonstrate the effectiveness of our proposed UAG approach (*UAG-both*) by comparing accuracy (under different attack methods and different attack ratio) with the original unoptimized GCN and GAT model as well as the one equipped with the state-of-the-art *RGCN* (Zhu et al. 2019) defense method. Besides, to gain more design insights, we add two more baselines for comparison, including *UAG-Data* (only considering data uncertainty) and *UAG-Model* (only considering data uncertainty). We show the accuracy under diverse attack ratio, defined as the number of attacked edges over the number of total edges. For *Metattack* on Pubmed dataset, since running the adversarial attack for all nodes is very time-consuming, we randomly sample 10% of them. For other datasets and attack method, we attack on all nodes.

Figure 5 shows accuracy performance comparisons for **Random Attack** among different implementations on the Cora, Citeseer, and Pubmed datasets. As we can easily tell from those three sub-plots, UAG method and its variants (UAG-Data, UAG-Model, and UAG-Both) consistently outperform the state-of-the-art *RGCN* defense approach. The major source of such performance improvements is that UAG effectively captures the data uncertainty and the model uncertainty, based on which UAG adaptively adjusts the edge weights and the amount of information propagation between nodes. For individual dataset settings, with the increase of the attack ratio, we see the overall trend of accuracy decreasing among these implementations. We also notice that on Cora and Citeseer dataset, *RGCN* offers notable accuracy improvement over the original GAT and GCN. However, it is still inferior compared with our UAG approach, since we explicitly capture the uncertainties, instead of relying on a Gaussian distribution to implicitly defend adversarial attacks.

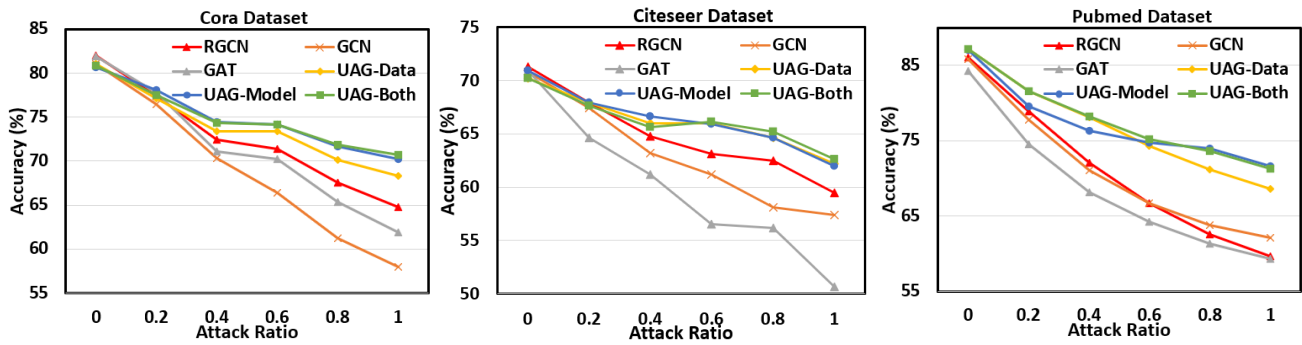


Figure 5: Results of different methods when adopting Random Attack as the attack method.

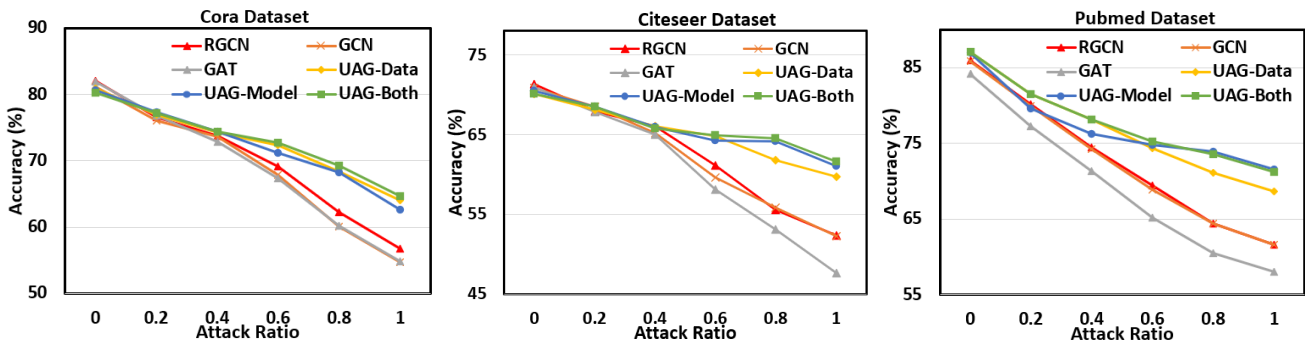


Figure 6: Results of different methods when adopting DICE Attack as the attack method.

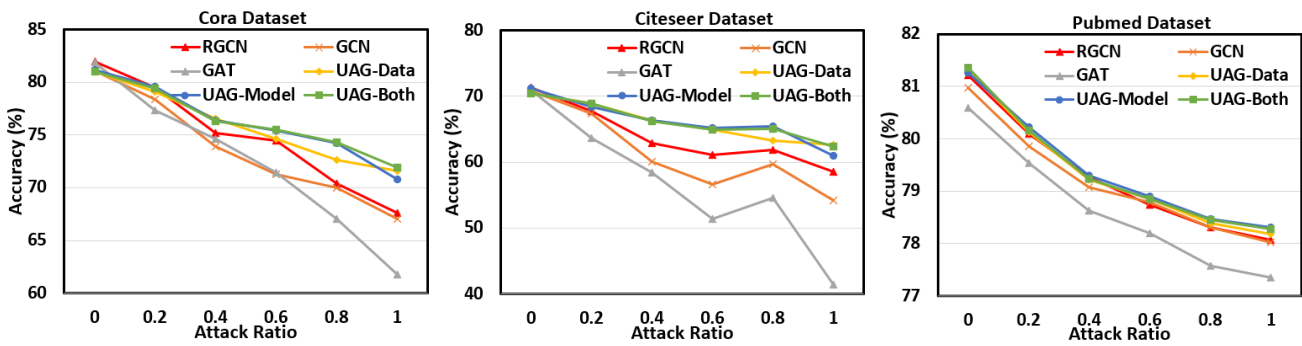


Figure 7: Results of different methods when adopting Mettack as the attack method.

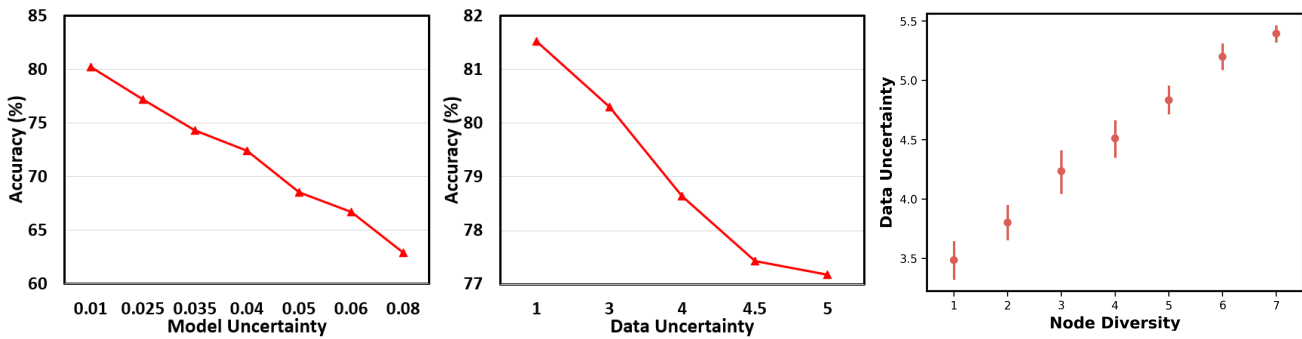


Figure 8: Relationship between Accuracy and Uncertainty. Left: Model Uncertainty v.s. Accuracy. Mid: Data Uncertainty v.s. Accuracy. Right: Data Uncertainty v.s. True Diversity.

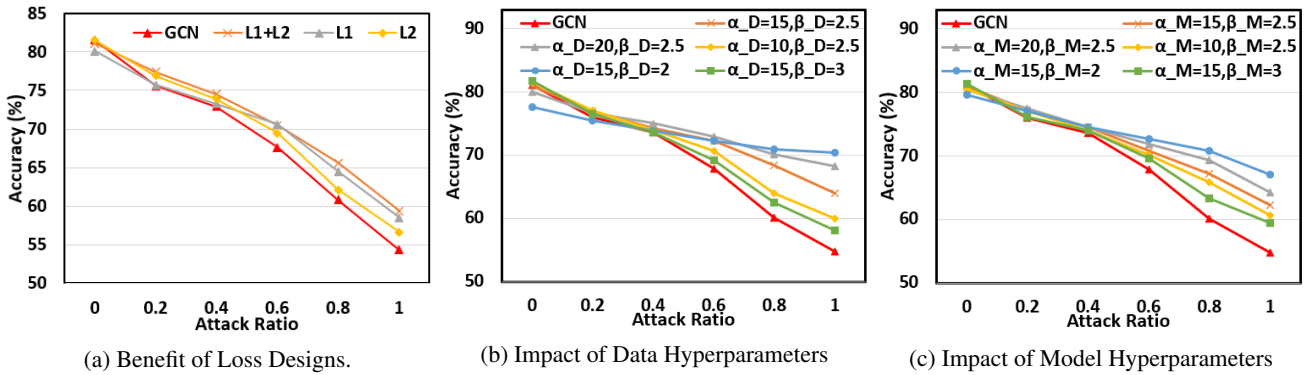


Figure 9: Loss Design Benefits and Parameter Analysis.

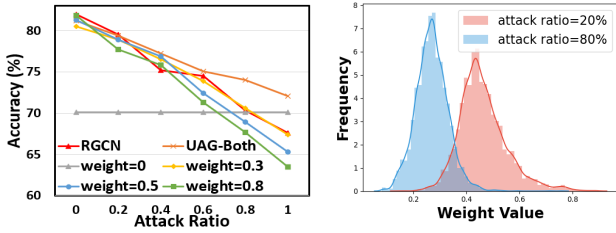


Figure 10: Left: Accuracy of Static Edge Weights. Right: Edge weight distribution under Random Attack.

Another observation is that UAG-Both usually outperforms UAG-Data and UAG-Model under diverse datasets and attack ratios. The reason is that the data uncertainty and the model uncertainty capture uncertainties from different sources and combining these two values usually offers a more comprehensive measurement on the prediction uncertainties.

On **DICE Attack** (Figure 6), besides the similar observations as the above attack setting, we have more observations. GAT and RGCN are sensitive towards DICE attack as demonstrated with significant accuracy drop with the increase of attack ratio, because DICE intentionally increases the node diversity by adding cross-community connections (*i.e.*, edges between nodes with different labels). In contrast, our UAG approach can handle this attack effectively, because it explicitly captures the node diversity as the data uncertainty.

We also observe a similar performance trend on **Metattack** (Figure 7) as the previous two types of attacks and further demonstrate the advantage of our UAG approach in terms of higher accuracy under diverse attack ratios.

## Ablation Studies

In this section, we conduct a set of ablation studies for in-depth analysis.

**Accuracy and Uncertainty.** As shown in Figure 8(a) and (b), the increase of the model uncertainty and data uncertainty would lead to the decrease of the accuracy. This also strengthens our initial assumption that the relationship between uncertainties and accuracy can be explored to defend adversarial attacks. Figure 8(c) exhibits the relation between data uncertainty and node diversity, showing the effectiveness of our Loss design (Eq 12) in learning the node diversity.

**Loss Design Benefits.** Figure 9(a) validates the effectiveness of the two loss designs in the data uncertainty. Here, L1

(Eq 12) forces UAG to learn the node diversity in labeled nodes, while L2 (Eq 13) is an unsupervised loss that implicitly encodes the node diversity in unlabeled nodes. We observe that the UAG with L1+L2 outperforms only L1 or L2, since it fully exploits both the labeled and unlabeled node.

**Parameter Analysis.** Figure 9(b) and (c) shows the impact from different values of the data-uncertainty-related hyperparameters and the model-uncertainty-related hyperparameters on the UAG performance. Intuitively, larger values of  $\alpha_\tau$  and lower value of  $\beta_\tau$  lead to a stronger impact from the uncertainties and higher accuracy under large attack ratios, where  $\tau \in \{M, D\}$ . However, setting  $\alpha_\tau$  too large or  $\beta_\tau$  too small may also impose too many constraints on the information propagation. In our experiments, we can achieve satisfying results by using  $\alpha_\tau = 15$  and  $\beta_\tau = 2.5$ .

**Uncertainty on Attention Values.** Figure 10 visualizes the weight changes under different ratios of Random Attack. We can notice that a higher random attack ratio would lead to denser weight distribution towards zero. This is because our UAG approach would try to amortize the impact of such attacks by changing weight values towards zero that can minimize the value propagation between neighboring nodes. Furthermore, we consider pre-assigning different weights to show the key importance of weight value for adversarial attacks. We can see that the higher the static weight value the poorer the performance in maintaining model accuracy under the attack. The major reason is more “attack impacts” will be propagated to different nodes through node aggregation, thus, lowering the model overall performance. Our UAG can adaptively determine the weight value for different nodes based on the model and data uncertainty factor, thus, largely absorbing the influence of the adversarial attack.

## Conclusion

In this paper, we propose UAG, the first systematic defense solution for adversarial attacks on GNNs by considering hierarchical uncertainty in GNNs. UAG incorporates a Bayesian uncertainty technique to explicitly capture uncertainty in GNNs and further employs an uncertainty-aware attention technique to fortify GNNs. Extensive experiments further demonstrate UAG’s advantages over the state-of-the-art solutions. Overall, our work paves a new way of exploring uncertainty benefits in GNN research.

## Acknowledgements

We thank all anonymous reviewers for their valuable comments. This work was supported in part by NSF 1925717.

## References

- Alex Kendall, V. B.; and Cipolla, R. 2017. Bayesian SegNet: Model Uncertainty in Deep Convolutional Encoder-Decoder Architectures for Scene Understanding. In Tae-Kyun Kim, Stefanos Zafeiriou, G. B.; and Mikolajczyk, K., eds., *Proceedings of the British Machine Vision Conference (BMVC)*, 57.1–57.12. BMVA Press. ISBN 1-901725-60-X. doi: 10.5244/C.31.57. URL <https://dx.doi.org/10.5244/C.31.57>.
- Bojchevski, A.; and Günnemann, S. 2018. Deep Gaussian Embedding of Graphs: Unsupervised Inductive Learning via Ranking. In *International Conference on Learning Representations*. URL <https://openreview.net/forum?id=r1ZdKJ-0W>.
- Chen, H.; Li, X.; and Huang, Z. 2005. Link prediction approach to collaborative filtering. In *Proceedings of the 5th ACM/IEEE-CS Joint Conference on Digital Libraries (JCDL)*, 141–142. IEEE.
- Cheng, D.; Gong, Y.; Chang, X.; Shi, W.; Hauptmann, A.; and Zheng, N. 2018. Deep feature learning via structured graph Laplacian embedding for person re-identification. *Pattern Recognition* 82: 94–104.
- Dai, H.; Li, H.; Tian, T.; Huang, X.; Wang, L.; Zhu, J.; and Song, L. 2018. Adversarial attack on graph structured data. *arXiv preprint arXiv:1806.02371*.
- Duran, A. G.; and Niepert, M. 2017. Learning graph representations with embedding propagation. In *Advances in neural information processing systems (NIPS)*, 5119–5130.
- Fey, M.; and Lenssen, J. E. 2019. Fast Graph Representation Learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds (ICLR)*.
- Gal, Y. 2016. *Uncertainty in Deep Learning*. Ph.D. thesis, University of Cambridge.
- Gal, Y.; and Ghahramani, Z. 2016a. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48, ICML'16*, 1050–1059. JMLR.org.
- Gal, Y.; and Ghahramani, Z. 2016b. A Theoretically Grounded Application of Dropout in Recurrent Neural Networks. In *Proceedings of the 30th International Conference on Neural Information Processing Systems, NIPS'16*, 1027–1035. Red Hook, NY, USA: Curran Associates Inc. ISBN 9781510838819.
- Gibert, J.; Valveny, E.; and Bunke, H. 2012. Graph embedding in vector spaces by node attribute statistics. *Pattern Recognition* 45(9): 3072–3083.
- Grover, A.; and Leskovec, J. 2016. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM international conference on Knowledge discovery and data mining (SIGKDD)*, 855–864.
- Hamilton, W.; Ying, Z.; and Leskovec, J. 2017. Inductive representation learning on large graphs. In *Advances in neural information processing systems (NIPS)*, 1024–1034.
- Hasanzadeh, A.; Hajiramezanali, E.; Boluki, S.; Zhou, M.; Duffield, N.; Narayanan, K.; and Qian, X. 2020. Bayesian Graph Neural Networks with Adaptive Connection Sampling. In *ICML 2020: International Conference on Machine Learning*. URL <https://arxiv.org/abs/2006.04064>.
- Kaspar, R.; and Horst, B. 2010. *Graph classification and clustering based on vector space embedding*, volume 77. World Scientific.
- Kendall, A.; and Gal, Y. 2017. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in neural information processing systems (NeurIPS)*, 5574–5584.
- Kipf, T. N.; and Welling, M. 2017. Semi-supervised classification with graph convolutional networks. *International Conference on Learning Representations (ICLR)*.
- Kunegis, J.; and Lommatzsch, A. 2009. Learning spectral graph transformations for link prediction. In *Proceedings of the 26th Annual International Conference on Machine Learning (ICML)*, 561–568.
- Luo, D.; Ding, C.; Huang, H.; and Li, T. 2009. Non-negative laplacian embedding. In *2009 Ninth IEEE International Conference on Data Mining (ICDM)*, 337–346. IEEE.
- Luo, D.; Nie, F.; Huang, H.; and Ding, C. H. 2011. Cauchy graph embedding. In *Proceedings of the 28th International Conference on Machine Learning (ICML)*, 553–560.
- Perozzi, B.; Al-Rfou, R.; and Skiena, S. 2014. DeepWalk: Online Learning of Social Representations. In *Proceedings of the 20th ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, 701–710. New York, NY, USA: Association for Computing Machinery. ISBN 9781450329569. doi:10.1145/2623330.2623732. URL <https://doi.org/10.1145/2623330.2623732>.
- Tylenda, T.; Angelova, R.; and Bedathur, S. 2009. Towards time-aware link prediction in evolving social networks. In *Proceedings of the 3rd workshop on social network mining and analysis*, 1–10.
- Veličković, P.; Cucurull, G.; Casanova, A.; Romero, A.; Liò, P.; and Bengio, Y. 2018. Graph Attention Networks. In *International Conference on Learning Representations (ICLR)*. URL <https://openreview.net/forum?id=rJXmpikCZ>.
- Waniek, M.; Michalak, T. P.; Wooldridge, M. J.; and Rahwan, T. 2018. Hiding individuals and communities in a social network. *Nature Human Behaviour* 2: 139.
- Xiao, Y.; and Wang, W. Y. 2019. Quantifying uncertainties in natural language processing tasks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, 7322–7329.
- Xu, K.; Chen, H.; Liu, S.; Chen, P.-Y.; Weng, T.-W.; Hong, M.; and Lin, X. 2019a. Topology attack and defense for graph neural networks: An optimization perspective. *arXiv preprint arXiv:1906.04214*.



Xu, K.; Hu, W.; Leskovec, J.; and Jegelka, S. 2019b. How Powerful are Graph Neural Networks? In *International Conference on Learning Representations (ICLR)*. URL <https://openreview.net/forum?id=ryGs6iA5Km>.

Ying, R.; You, J.; Morris, C.; Ren, X.; Hamilton, W. L.; and Leskovec, J. 2018. Hierarchical Graph Representation Learning with Differentiable Pooling. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems (NIPS)*, 4805–4815. Red Hook, NY, USA.

Zhang, Y.; Pal, S.; Coates, M.; and Ustebay, D. 2019. Bayesian Graph Convolutional Neural Networks for Semi-Supervised Classification. *Proceedings of the AAAI Conference on Artificial Intelligence* 33: 5829–5836. doi: 10.1609/aaai.v33i01.33015829.

Zhu, D.; Zhang, Z.; Cui, P.; and Zhu, W. 2019. Robust graph convolutional networks against adversarial attacks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 1399–1407.

Zhu, L.; and Laptev, N. 2017. Deep and Confident Prediction for Time Series at Uber. In *2017 IEEE International Conference on Data Mining Workshops (ICDMW)*, 103–110.

Zügner, D.; Akbarnejad, A.; and Günnemann, S. 2018. Adversarial attacks on neural networks for graph data. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2847–2856.

Zügner, D.; and Günnemann, S. 2019. Adversarial attacks on graph neural networks via meta learning. *arXiv preprint arXiv:1902.08412* .