

# Robust Heterogeneous Graph Neural Networks against Adversarial Attacks

Mengmei Zhang<sup>1</sup>, Xiao Wang<sup>1</sup>, Meiqi Zhu<sup>1</sup>, Chuan Shi<sup>1\*</sup>, Zhiqiang Zhang<sup>2</sup>, Jun Zhou<sup>2</sup>

<sup>1</sup>Beijing University of Posts and Telecommunications

<sup>2</sup>Ant Group, Hangzhou, China

{zhangmm, xiaowang, zhumeiqi, shichuan}@bupt.edu.cn, {lingyao.zzq, jun.zhoujun}@antfin.com

## Abstract

Heterogeneous Graph Neural Networks (HGNNs) have drawn increasing attention in recent years and achieved outstanding performance in many tasks. However, despite their wide use, there is currently no understanding of their robustness to adversarial attacks. In this work, we first systematically study the robustness of HGNNs and show that they can be easily fooled by adding the adversarial edge between the target node and large-degree node (i.e., hub). Furthermore, we show two key reasons for such vulnerabilities of HGNNs: one is *perturbation enlargement effect*, i.e., HGNNs, failing to encode transiting probability, will enlarge the effect of the adversarial hub in comparison of GCNs, and the other is *soft attention mechanism*, i.e., such mechanism assigns positive attention values to obviously unreliable neighbors. Based on the two facts, we propose a novel robust HGNN framework *RoHe* against topology adversarial attacks by equipping an attention purifier, which can prune malicious neighbors based on topology and feature. Specifically, to eliminate the perturbation enlargement, we introduce the metapath-based transiting probability as the prior criterion of the purifier, restraining the confidence of malicious neighbors from adversarial hub. Then the purifier learns to mask out neighbors with low confidence, thus can effectively alleviate the negative effect of malicious neighbors in the soft attention mechanism. Extensive experiments on different benchmark datasets for multiple HGNNs are conducted, where the considerable improvement of HGNNs under adversarial attacks will demonstrate the effectiveness and generalization ability of our defense framework.

## Introduction

Many real-world datasets are usually modeled with Heterogeneous Graphs (HGs) (Shi et al. 2017), which contain diverse types of objects and relations. An example of ACM citation network characterized by HG is given in Figure 1(a), consisting of three types of objects (Author (A), Paper (P), Subject (S)), and two types of relations (P-A and P-S). Since HGs contain rich high-order structural information, metapath (sequence of relation types between two node types) is widely used as a basic tool to capture such information (Shi et al. 2017), such as P-A-P (papers written by the

same author) and P-S-P (papers attached to the same subject). In recent years, with deep learning employed on HGs, there is a surge of Heterogeneous Graph Neural Networks (HGNNs) (Wang et al. 2019b; Yun et al. 2019; Fu et al. 2020), which often adopt a hierarchical aggregation (including node-level and semantic-level) to capture the information from metapath-based neighbors, and have achieved state-of-the-art performance on a wide range of tasks, e.g., node classification and link prediction.

Despite the great success of HGNNs, there is no systematic understanding of the adversarial robustness of HGNNs, i.e., *whether the HGNNs can be easily fooled by slight perturbations of the input topology*. This is especially important for HGNN models since they are widely applied to many real-world applications, e.g., e-commerce (Zhang et al. 2019a; Hu et al. 2019) and cyber security (Zhang et al. 2019c; Zhong et al. 2020). So far, most works focus on adversarial vulnerabilities of homogeneous GNNs (Sun et al. 2018), but the robustness for HGNNs is indeed not foreseeable due to unique metapath-based aggregation.

To answer this question, we introduce the first study of adversarial robustness of HGNNs through evaluating their performance on dataset ACM under the same evasion adversarial attacks<sup>1</sup>, which perturb topology in the test phase, and the attack results are shown in Figure 1(b). Surprisingly, compared to the drop of GCN by about 3 points, HGNNs, i.e., HAN, MAGNN, and GTN, dramatically decrease by an average of 28 points. Obviously, HGNNs have significantly different adversarial robustness from GCNs, which motivates us to further investigate the differences of architectures between GCNs and HGNNs.

In the further analysis of attack results, we observe that the attackers tend to maliciously link the target node to the large-degree node (i.e., hub). Taking HAN as an example, the attacker injects an adversarial edge ( $p_1, a_3$ ) in Figure 1(a), which will lead malicious (red) papers  $p_4 \cdots p_{66}$  to be the direct neighbors of  $p_1$  under metapath PAP. And even they are assigned small attention values, they can still dominate the receptive field of HAN in Figure 1(c). We argue such vulnerabilities of HGNNs can be attributable to two key reasons: (1) **Perturbation enlargement effect**.

\*Corresponding author

Copyright © 2022, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

<sup>1</sup>Attack algorithm and the implementation details of the compared GCN (Kipf and Welling 2017) can be found in Appendix.

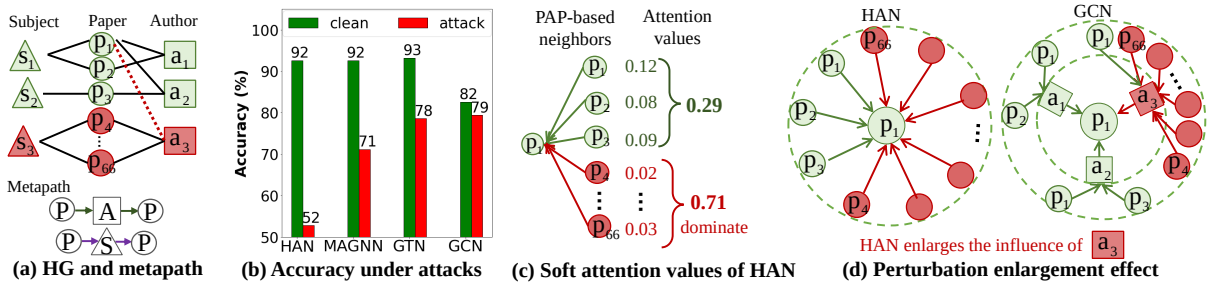


Figure 1: The illustrative example for adversarial attack against HGNNs on ACM dataset. (a) Basic concepts of HG (HG, metapath). (b) The a robustness evaluation of HGNNs and GCN. (c) A toy example of the soft attention values of HAN under adversarial edge  $(p_1, a_3)$ . (d) The comparison of the influence of adversarial link  $(p_1, a_3)$  to HAN and GCN.

We will prove that HGNNs will rapidly enlarge the effect of the adversarial hub. While GCNs will not enlarge it, since compared to HGNNs, as shown in Figure 1(d), GCNs will not regard the malicious  $p_4 \cdots p_{66}$  as the direct neighbors of  $p_1$ , and thus the malicious two-hop neighbors  $p_4 \cdots p_{66}$  can only influence  $p_1$  through one-hop neighbor  $a_3$ . However, HAN directly aggregates all neighbors under PAP with equal weights  $\frac{1}{66}$ , and thus enlarges the effect of adversarial  $(p_1, a_3)$  to  $\frac{63}{66}$  (i.e., the total weights of malicious  $p_4 \cdots p_{66}$ ). (2) **Soft attention mechanism.** Conventional attention mechanisms assume all neighbors are reliable and aggregate them with soft (i.e., positive) values. This soft attention mechanism may hurt the performance of GNNs when existing adversarial/noisy/disassortative neighbors (Zhang and Zitnik 2020; Bo et al. 2021). It will cause more serious damage to HGNN when injecting adversarial hub as shown in Figure 1 (c).

Once the vulnerabilities of HGNNs are identified, there is a strong need for further improving the adversarial robustness of HGNNs. Thus, in this paper, we propose a **Robust Heterogeneous GNN framework (RoHe)** against topology adversarial attacks by designing an attention purifier, which can prune malicious neighbors based on topology and feature. More specifically, for the problem of perturbation enlargement, we introduce the metapath-based transiting probability as the prior criterion of the purifier, restraining the confidence of malicious neighbors from the adversarial hub. Then the purifier learns a differentiable mask vector to remove the unreliable neighbors in the soft attention mechanism.

The contributions of this work are three folds:

- We introduce the first systematic exploration and assessment of the robustness of HGNNs, and point out that the HGNNs are highly fragile to adversarial link to the hub, which can be attributed to the problems of perturbation enlargement effect and soft attention mechanism.
- Based on the above findings, we propose a novel robust HGNN framework (RoHe) against adversarial attacks by designing an attention purifier, which can constrain the enlargement perturbations by transiting probability and eliminate the negative impact of malicious neighbors through mask operation.
- We perform experiments on different benchmark datasets

for multiple HGNNs. The effectiveness and generalization ability of our defense framework RoHe are well demonstrated by the considerable improvement of HGNNs under adversarial attacks.

## Preliminaries

**Definition 1 Heterogeneous Graph.** A heterogeneous graph, defined as  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , consists of an object set  $\mathcal{V}$  and an edge set  $\mathcal{E}$ .  $\mathcal{G}$  is also associated with a node type mapping function  $\phi : \mathcal{V} \rightarrow \mathcal{A}$  and an edge type mapping function  $\psi : \mathcal{E} \rightarrow \mathcal{R}$ .  $\mathcal{A}$  and  $\mathcal{R}$  denote the predefined sets of node types and edge types, where  $|\mathcal{A}| + |\mathcal{R}| > 2$ . For each type  $R \in \mathcal{R}$ ,  $\mathbf{M}^R$  represents the corresponding binary adjacency matrix.

**Definition 2 Metapath.** A metapath  $\Phi$  is defined as a path in the form of  $\Phi = A_1 \xrightarrow{R_1} A_2 \xrightarrow{R_2} \cdots \xrightarrow{R_l} A_{l+1}$ , which describes a composite relation  $R = R_1 \odot R_2 \odot \cdots \odot R_l$  between node types  $A_1$  and  $A_{l+1}$ .

**Definition 3 Metapath based Neighbors.** Given a node  $v$  and a metapath  $\Phi$  in a heterogeneous graph, the metapath based neighbors  $\mathcal{N}_v^\Phi$  are defined as the set of nodes that connect with  $v$  via metapath  $\Phi$ .

**Metapath-based transiting probability.** In this paper, we consider the metapath-based transiting probability denoted by  $\mathbf{P}_{vu}^\Phi$  (from node  $v$  to neighbor  $u$  along metapath  $\Phi = A_1 \xrightarrow{R_1} A_2 \xrightarrow{R_2} \cdots \xrightarrow{R_l} A_{l+1}$ ), which can be used to guide a metapath-based random walk for learning heterogeneous graph embedding, e.g., metapath2vec (Dong, Chawla, and Swami 2017). And the matrix  $\mathbf{P}^\Phi$  can be calculated by

$$\mathbf{P}^\Phi = \mathbf{P}^{R_1} \cdots \mathbf{P}^{R_l}, \quad (1)$$

where  $\mathbf{P}^{R_i} = (\mathbf{D}^{R_i})^{-1} \mathbf{M}^{R_i}$  for  $i \in \{1 \cdots l\}$ . This shows that given a metapath  $\Phi$ ,  $\mathbf{P}_{vu}^\Phi$  is defined in terms of two parts: (1) their connectivity defined by the number of paths between  $v$  and  $u$  following  $\Phi$ ; and (2) the degree information of all nodes along paths.

**Heterogeneous graph neural networks.** HGNNs often adopt a hierarchical aggregation: the node-level one aims to merge the neighbors based on a specific metapath, and the semantic-level one can fuse the information of different metapaths. In this paper, we focus on three representative

HGNNs, i.e., HAN (Wang et al. 2019b), MAGNN (Fu et al. 2020) and GTN (Yun et al. 2019). Taking HAN as an example, the metapath-based embedding of node  $v$  can be aggregated as follows:

$$\mathbf{z}_v^\Phi = \sigma \left( \sum_{u \in \mathcal{N}_v^\Phi} a_{vu}^\Phi \cdot \mathbf{h}_u \right), \quad (2)$$

where  $a_{vu}^\Phi$  is the attention value for neighbor  $u$ ,  $\mathbf{h}_u$  is the projected feature of  $u$ ,  $\mathcal{N}_v^\Phi$  is the metapath-based neighbors.

To facility analysis, we provide more preliminaries about the mechanism of (asymmetrically normalized) GCN and MAGNN/GTN in Appendix. For clarity, we also formally provide the simplified structure-based weights of  $u \in \mathcal{N}_v^\Phi$  for these models as shown in Table 1, by giving metapath  $\Phi = A_1 \xrightarrow{R_1} A_2 \xrightarrow{R_2} A_3$  and supposing the features of nodes are the same.

**Adversarial attacks on GNNs.** In this paper, we focus on evasion attack, a typical type of adversarial attack that perturbs graph in the test phase and guides the model to misclassify the target node  $v$ . Specifically, given a homogeneous graph with adjacency matrix  $\mathbf{M}$  and node features  $\mathbf{X}$ , the goal of an attacker is to find the optimal perturbed adjacency matrix  $\mathbf{M}_\Delta$ :

$$\operatorname{argmax}_{\mathbf{M}_\Delta} \mathcal{L}(f_{GNN}^*(\mathbf{M}_\Delta, \mathbf{X})_v, c_v), \quad (3)$$

where  $f_{GNN}^*(\mathbf{M}_\Delta, \mathbf{X})_v$  is the prediction of trained GNN model  $f_{GNN}^*$  for node  $v$ ,  $c_v$  is the label of  $v$ ,  $\Delta$  (named budget) is the maximum number of the perturbed edges, and  $\mathcal{L}$  is the classification loss in this paper. The yielding optimal  $\mathbf{M}_\Delta$  will lead to minimum test accuracy.

### Adversarial Vulnerability Analysis

We perform adversarial attacks on HGNNs and GCN (details of attack method are in Appendix), and the results presented in Figure 1 (b) clearly show that compared to GCN, HGNNs are highly vulnerable to adversarial attacks, especially HAN. Here we further analyze the key reasons for such vulnerabilities.

### Perturbation Enlargement Effect

We discover that HGNNs exist the phenomenon of perturbation enlargement, i.e., HGNNs will enlarge the effect of the adversarial hub. As shown in Figure 1 (d), the influence of adversarial hub  $a_3$ , expected to be less than  $\frac{1}{3}$  (the inverse of the  $p_1$ 's author neighbors), is enlarged to  $\frac{63}{66}$  for HAN. Specifically, when the attacker injects an adversarial hub  $a_3$  as the direct neighbor of  $p_1$ , the influence of  $a_3$  to  $p_1$  should be proportional to the inverse of the degree of target node  $p_1$  (i.e.,  $\frac{1}{3}$ ) from the perspective of network science (i.e., transiting probability). While HGNNs can not satisfy it and enlarge the total weights to  $\frac{63}{66}$  for HAN and  $\frac{63}{68}$  for MAGNN/GTN, since they skip the intermediate layers (e.g., the layer for author in PAP), and directly aggregate multi-hop neighbors  $p_1 \cdots p_{66}$ , failing to encode transiting probability  $\mathbf{P}^{R_1} \mathbf{P}^{R_2}$  in structural weight of  $u$ .

Model	Weight of $u$	$\mathbf{p}_1$	$\mathbf{p}_2/\mathbf{p}_3$	$\mathbf{p}_4\text{-}66$
TransP	$(\mathbf{P}^{R_1} \mathbf{P}^{R_2})_{vu}$	$\frac{2}{6} + \frac{1}{3 \times 64}$	$\frac{1}{6}$	$\frac{1}{3 \times 64}$
HAN	$\frac{1}{ \mathcal{N}_v^\Phi }$	$\frac{1}{66}$	$\frac{1}{66}$	$\frac{1}{66}$
MAGNN	$(\mathbf{D}_v^\Phi)^{-1} \mathbf{M}_{vu}^\Phi$	$\frac{3}{68}$	$\frac{1}{68}$	$\frac{1}{68}$
GTN	$(\mathbf{D}_v^\Phi)^{-1} \mathbf{M}_{vu}^\Phi$	$\frac{3}{68}$	$\frac{1}{68}$	$\frac{1}{68}$

Table 1: The structural weights of  $u \in \mathcal{N}_v^\Phi$  and their examples in HGNNs, TransP (short for Transiting Probability). Here  $\mathbf{P}^{R_i} = (\mathbf{D}^{R_i})^{-1} \mathbf{M}^{R_i}$ .

We also find that the perturbation enlargement effect is more significant in HAN than MAGNN and GTN. Taking Figure 1(d) as an example, we can see that  $p_1$  is connected to  $p_1$  more densely (by 3 paths) than  $p_4$  (by 1 path). Thus the  $p_1$  is expected to have larger weights than  $p_4$  in transiting probability. MAGNN and GTN can satisfy it and assign higher weight to itself  $p_1$  ( $\frac{3}{68}$ ) than malicious  $p_4$  ( $\frac{1}{68}$ ), by encoding the total number of path instances (i.e.,  $(\mathbf{D}_v^\Phi)^{-1} \mathbf{M}_{vu}^\Phi$ ). While HAN equally treats all neighbors with same weights  $\frac{1}{66}$ , thus yields the larger total weights of malicious  $p_4 \cdots p_{66}$  ( $\frac{63}{66}$ ) than MAGNN/GTN ( $\frac{63}{68}$ ).

### Soft Attention Mechanism

We argue the soft attention mechanism will especially hurt the generalization performance on adversarial attacks for HGNNs. As shown in Figure 1 (c), vast malicious neighbors  $p_4 \cdots p_{66}$  can accumulate the smaller but positive attention values and finally dominate the receptive field of HGNNs, misleading the classification of  $p_1$ . Based on this fact, the power of assigning zero attention values to obviously unreliable neighbors is significant for HGNNs.

## The Proposed Robust Heterogeneous GNN

This section depicts our proposed **Robust Heterogeneous GNNs (RoHe)** against topology adversarial attacks. HGNNs often adopt a hierarchical aggregation (including node-level and semantic-level), and our RoHe is applied to purify the node-level aggregation. Figure 2 illustrates the overall architecture of RoHe. The node-level attention for metapath-based neighbors will be equipped by our purifier, which can constrain the enlargement perturbations by transiting prior and eliminate the negative impact of malicious neighbors through mask operation. Then the purified attention will be used for node-level aggregation, yielding the node embeddings for different metapath, which can be fused in semantic-level aggregation. Note that here we present our general framework based on HAN (Wang et al. 2019b) for simplicity.

### Node-level Aggregation

Here we first detail the node-level attention mechanism and show that our purifier can eliminate the problems of perturbation enlargement and soft attention mechanism by transiting probability and purification mask.

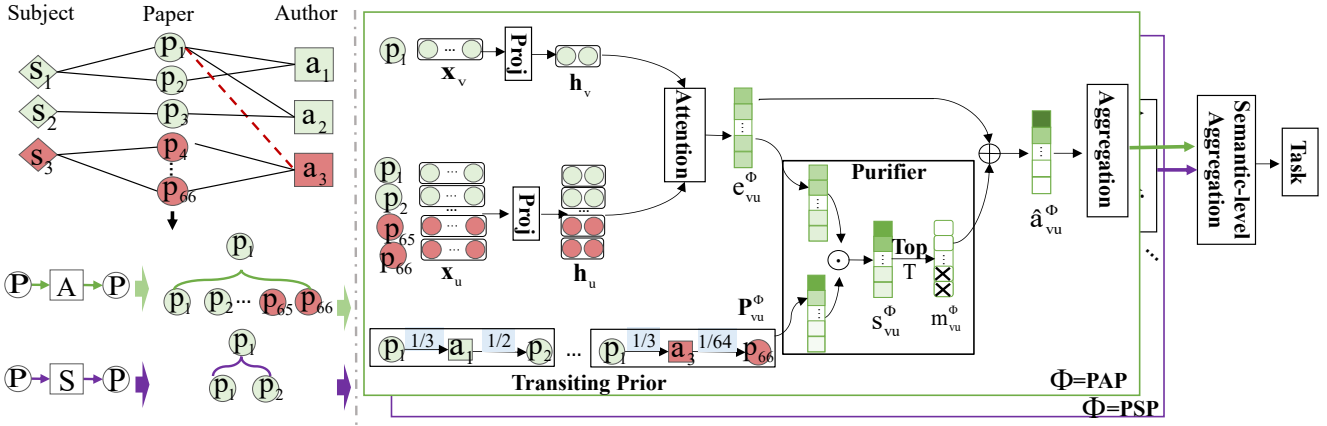


Figure 2: The overall framework of RoHe.

**Node feature transformation.** Since different node types may have unequal dimensions of feature vectors or lie in different feature spaces, HGNNs usually project the features of different types of nodes into the common space. Specifically, for the target node  $v$  with type  $A \in \mathcal{A}$ , we use a type-specific transformation matrix  $\mathbf{W}_A$  to obtain the projected features  $\mathbf{h}_v$  as follows:

$$\mathbf{h}_v = \mathbf{W}_A \mathbf{x}_v. \quad (4)$$

**Feature-based importance.** Given a metapath  $\Phi$ , based on the hypothesis that the nodes with similar features are more likely to be important than dissimilar ones, we estimate the importance  $e_{vu}^\Phi$  of neighbors  $u$  to target node  $v$  under  $\Phi$  by dot-product similarity of features:

$$e_{vu}^\Phi = \mathbf{h}_v \cdot \mathbf{h}_u. \quad (5)$$

In conventional node-level attention mechanism, the feature-based importance  $e_{vu}^\Phi$  will be directly normalized across  $\mathcal{N}_v^\Phi$  with the softmax function, yielding the final soft attention values  $a_{vu}^\Phi$ . We argue  $a_{vu}^\Phi$  only considers the feature information of nodes, while equally treats the multi-hop neighbors  $\mathcal{N}_v^\Phi$  from the perspective of topology, which will lead to enlarging the effect of adversarial hub neighbor as proved in Section . Besides, all neighbors in  $\mathcal{N}_v^\Phi$  are assigned positive values after softmax function. Such soft attention mechanism has excellent differentiability in back propagation (Chaudhari et al. 2019), but fails to assign zero value to obviously malicious neighbors as showed in Section .

To solve above problems, we introduce a differentiable purifier to mask out the neighbor  $u \in \mathcal{N}_v^\Phi$  with low confidence score  $s_{vu}^\Phi$ . Specifically, we first utilize metapath-based transiting probability  $\mathbf{P}_{vu}^\Phi$  as the prior of confidence of  $u$  to eliminate perturbation enlargement problem.

**Transiting prior.** Given metapath  $\Phi = A_1 \xrightarrow{R_1} A_2 \xrightarrow{R_2} \dots \xrightarrow{R_l} A_{l+1}$ , to encode the probability of transiting along metapath  $\Phi$  as a prior, we first calculate the transiting probability matrix  $\mathbf{P}^{R_i} = (\mathbf{D}^{R_i})^{-1} \mathbf{M}^{R_i}$  for relation  $R_i \in \{R_1, \dots, R_l\}$ . Each element  $\mathbf{P}_{vu}^{R_i}$  represents the probability of transiting from node  $v$  to  $u$  in relation  $R_i$ . And

powering it along  $\Phi$  will lead to the metapath-based transiting probability  $\mathbf{P}_{vu}^\Phi$  as introduced in Section . Then we use the element  $\mathbf{P}_{vu}^\Phi$  as the prior confidence of neighbor  $u$  for target node  $v$  in metapath  $\Phi$ . We can see that the neighbor  $u$  is expected to obtain a small  $\mathbf{P}_{vu}^\Phi$  for confidence, if  $u$  is indirectly connected to  $v$  passing through the hub node, which can solve the enlargement of adversarial hub as described in Section .

**Confidence score.** Based on the transiting prior  $\mathbf{P}_{vu}^\Phi$ , to determine the unreliable neighbors, we can calculate the confidence score vector  $\mathbf{s}_v^\Phi \in \mathbb{R}^{|\mathcal{N}_v^\Phi|}$  for neighbors  $\mathcal{N}_v^\Phi$  based on both feature and topology, by incorporating feature similarity  $e_{vu}^\Phi$  and  $\mathbf{P}_{vu}^\Phi$ :

$$s_{vu}^\Phi = \sigma(\mathbf{P}_{vu}^\Phi \cdot e_{vu}^\Phi). \quad (6)$$

The notation  $s_{vu}^\Phi$ , as an element of  $\mathbf{s}_v^\Phi$ , is the confidence score for neighbor  $u \in \mathcal{N}_v^\Phi$ , indicating that neighbors with similar features and high transiting probabilities are regarded to be reliable.

For the problem of soft attention, we design a mask operation, which can mask out neighbors with low confidence in a differentiable way.

**Purification mask.** We model the mask operation by constructing a mask vector  $\mathbf{m}_v^\Phi \in \{1, -\infty\}^{|\mathcal{N}_v^\Phi|}$  for all the neighbors of target node  $v$  by

$$m_{vu}^\Phi = \begin{cases} 0 & \text{if } u \in \text{Top}(\mathbf{s}_v^\Phi, T), \\ -\infty & \text{otherwise,} \end{cases} \quad (7)$$

where  $T$  is the number of neighbors to be kept, and  $\text{Top}(\cdot)$  returns the set of the  $T$  most reliable neighbors based on their confidence scores  $\mathbf{s}_v^\Phi$ . Then the other neighbors will be removed by setting their mask values as  $-\infty$ . When a softmax is applied to a sum of  $e_{vu}^\Phi$  and  $m_{vu}^\Phi = -\infty$ , the node  $u$  will be effectively masked out, since the output of softmax for  $-\infty$  is zero.

Thus, we can use  $\mathbf{m}_v^\Phi$  to mask out the abundant adversarial/noisy neighbors, yielding purified attention  $\hat{a}_{vu}^\Phi$  via softmax function:

$$\hat{a}_{vu}^\Phi = \frac{\exp(m_{vu}^\Phi + e_{vu}^\Phi)}{\sum_{i \in \mathcal{N}_v^\Phi} \exp(m_{vi}^\Phi + e_{vi}^\Phi)}. \quad (8)$$

In this way, node-level attention is enhanced to encode the transiting probability of metapath-based neighbors and only aggregate top- $T$  reliable neighbors, alleviating the problems of perturbation enlargement and soft attention mechanism.

$$\hat{a}_{vu}^\Phi = \text{softmax}_u(e_{vu}^\Phi) = \frac{\exp(e_{vu}^\Phi)}{\sum_{i \in \mathcal{N}_v^\Phi} \exp(e_{vi}^\Phi)}. \quad (9)$$

**Aggregation of neighbors.** Finally, the final purified attention  $\hat{a}_{vu}^\Phi$  will be used to aggregate neighbors for semantic-specific embedding  $\mathbf{z}_v^\Phi$  as

$$\mathbf{z}_v^\Phi = \sum_{u \in \mathcal{N}_v^\Phi} (\hat{a}_{vu}^\Phi \cdot \mathbf{h}_u). \quad (10)$$

### Semantic-level Aggregation

Since different metapaths capture different semantics of the HG, HGNNs usually adopt semantic-level attention to calculate the importance of each metapath. Given the metapath set  $\{\Phi_0, \Phi_1, \dots, \Phi_P\}$ , after node-level aggregation, we can obtain a group of semantic-specific node embeddings of  $v$ , denoted as  $\{\mathbf{z}_v^{\Phi_0}, \mathbf{z}_v^{\Phi_1}, \dots, \mathbf{z}_v^{\Phi_P}\}$ . HAN further calculates the importance of metapath  $\Phi \in \{\Phi_1, \dots, \Phi_P\}$  by

$$w^\Phi = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \mathbf{q}^T \cdot \text{tanh}(\mathbf{W} \cdot \mathbf{z}_v^\Phi + \mathbf{b}), \quad (11)$$

where  $\mathbf{W}$  and  $\mathbf{b}$  denote the weight matrix and bias of the MLP, respectively.  $\mathbf{q}$  is the semantic-level attention vector. Then HAN uses the softmax function to normalize the importance  $w^\Phi$  to yield the attention value  $\beta^\Phi$  for  $\Phi$ . Hence, the final embedding  $\mathbf{z}_v$  of  $v$  can be obtained by semantic-level aggregation:

$$\mathbf{z}_v = \sum_{\Phi \in \{\Phi_1, \dots, \Phi_P\}} \beta^\Phi \cdot \mathbf{z}_v^\Phi. \quad (12)$$

Finally, the overall proposed model can be optimized by minimizing following loss:

$$\mathcal{L} = - \sum_{v \in \mathcal{V}_L} \ln(\mathbf{W}_{clf} \cdot \mathbf{z}_{v,c_v}), \quad (13)$$

where  $\mathbf{W}_{clf}$  is the parameter of the classifier,  $c_v$  is the class of training node  $v \in \mathcal{V}_L$ . The overall process of our proposed RoHe is summarized in Algorithm 1.

## Experiments

### Experimental Setup

**Datasets.** RoHe is evaluated on three widely used HG datasets: (1) **ACM** (Wang et al. 2019a) consists of {Paper (P), Author (A), Subject (S)} and we employ metapath set {PAP, PSP} for paper classification. (2) **DBLP** (Fu et al. 2020) consists of {Author (A), Paper (P), Term (T), Conference (C)} and we use metapath set {APA, APCPA, APTPA} for author classification. (3) **Aminer** (Hu, Fang, and Shi 2019) consists of {Paper (P), Author (A), Reference (R)} and we employ metapaths set {PAP, PRP} for paper classification. Note that the features of ACM and DBLP are based on bag-of-words representations, and the features of Aminer are assigned one-hot id vectors to nodes. Details are in Appendix.

---

### Algorithm 1: RoHe: Robust heterogeneous HAN

---

**Require:** The heterogeneous graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ ,  
The node features  $\{\mathbf{x}_v, v \in \mathcal{V}\}$ ,  
The metapath set  $\{\Phi_0, \Phi_1, \dots, \Phi_P\}$ ,  
The mask threshold  $T$ .

**Ensure:** The final node embeddings  $\{\mathbf{z}_v, v \in \mathcal{V}\}$ .

- 1: Pre-process transiting matrix via Eq. (1);
  - 2: **for** node type  $A \in \mathcal{A}$  **do**
  - 3:   Type-specific transformation to obtain  $\{\mathbf{h}_v, v \in \mathcal{V}\}$ ;
  - 4: **end for**
  - 5: **for**  $\Phi \in \{\Phi_0, \Phi_1, \dots, \Phi_P\}$  **do**
  - 6:   **for**  $v \in \mathcal{V}$  **do**
  - 7:     Find the metapath-based neighbors  $\mathcal{N}_v^\Phi$
  - 8:     **for**  $u \in \mathcal{N}_v^\Phi$  **do**
  - 9:       Calculate the feature-based importance  $e_{vu}$  for  $u \in \mathcal{N}_v^\Phi$  via Eq. (5);
  - 10:       Calculate confidence score via Eq. (6);
  - 11:       Obtain purification mask vector  $\mathbf{m}_v^\Phi$  via Eq. (7);
  - 12:       Obtain the purified attention  $\hat{a}_{vu}^\Phi$  via Eq. (9);
  - 13:     **end for**
  - 14:     Obtain the node embedding  $\mathbf{z}_v^\Phi$  for  $\Phi$  via Eq. (10);
  - 15:   **end for**
  - 16: **end for**
  - 17: Calculate the semantic-level attention values  $\{\beta^\Phi\}$  for  $\Phi \in \{\Phi_0, \Phi_1, \dots, \Phi_P\}$ ;
  - 18: Obtain final node embeddings  $\{\mathbf{z}_v, v \in \mathcal{V}\}$  by fusing the embeddings from different metapath via Eq. (12);
- 

**Setup.** (1) **HGNNs:** We mainly evaluate the effectiveness of our RoHe on HAN, and we also generalize RoHe to MAGNN (Fu et al. 2020) and GTN (Yun et al. 2019). (2) **Baselines:** Since there are no existing robust HGNNs methods, we compare with the direct adaptations of following strategies: Jaccard (Wu et al. 2019), GGCL (Zhu et al. 2019) and SimP (Jin et al. 2021), and the variants of our RoHe: **RoHe<sub>T</sub>** (only keeping transiting probability) and **RoHe<sub>P</sub>** (only keeping pruning operation). (3) **Generating adversarial attack:** We employ FGSM-based attacks (Goodfellow, Shlens, and Szegedy 2015) to generate perturbation edges in experiments. Given a target node, we limit adversarial edges with budget  $\Delta = \{1, 2, 3\}$  and edge types as P-A for ACM/DBLP and P-R for Aminer. We evaluate the performance with Micro-F1 metric over 500 target nodes, which are randomly sampled from the test set. More details about experimental settings are in Appendix.

### Defense Effectiveness of RoHe

Here we evaluate the effectiveness of RoHe on HAN (i.e., HAN-RoHe) against all baselines, under two scenarios (Clean and Attack). The overall results are presented in Table 2, and results of more metrics are in Appendix. Here we have the following observations:

(1) Attacker can dramatically decrease the performance of HAN by about 43% by adding one edge. However, the proposed HAN-RoHe successfully restores the performance of GNNs to the level comparable to when there is no attack. For example, with the increase of budget  $\Delta$ , the per-

Data	Model	Clean	Attack		
			$\Delta = 1$	$\Delta = 3$	$\Delta = 5$
ACM	HAN	0.926	0.528	0.330	0.240
	Jaccard	0.918	0.892	0.860	0.848
	SimP	0.898	0.746	0.476	0.358
	GGCL	0.902	0.260	0.084	0.084
	HAN-RoHe <sub>P</sub>	0.924	0.780	0.868	0.870
	HAN-RoHe <sub>T</sub>	<b>0.940</b>	0.900	0.564	0.304
DBLP	HAN	0.942	0.332	0.096	0.060
	Jaccard	0.934	0.816	0.812	0.802
	SimP	0.942	0.790	0.670	0.600
	GGCL	0.914	0.684	0.464	0.344
	HAN-RoHe <sub>P</sub>	0.862	0.686	0.714	0.702
	HAN-RoHe <sub>T</sub>	<b>0.944</b>	0.760	0.360	0.220
Aminer	HAN	<b>0.882</b>	0.346	0.134	0.102
	GGCL	0.808	0.276	0.056	0.042
	HAN-RoHe <sub>P</sub>	0.840	0.772	0.772	0.774
	HAN-RoHe <sub>T</sub>	0.842	0.788	0.668	0.562
	HAN-RoHe	0.838	<b>0.840</b>	<b>0.812</b>	<b>0.802</b>

Table 2: Results (Micro-F1) of HAN-RoHe. A higher value indicates better robustness.

formance of HAN-RoHe only drops by about 5% for ACM and Aminer. The reason is that HAN can greatly benefit from RoHe by equipping an attention purifier, which filters adversarial neighbors and retains the essential neighbors.

(2) The proposed HAN-RoHe consistently outperforms all defense methods in the Attack scenario. 1) Jaccard and SimP, pruning unreliable neighbors based on the feature similarity, will only alleviate the problem of soft attention mechanism and thus have limited improvement. 2) The Gaussian layer of GGCL also cannot completely absorb the vast adversarial neighbors, failing to defend against such attacks. 3) HAN-RoHe<sub>T</sub> and HAN-RoHe<sub>P</sub> only solve one of the problems of HGNNs respectively, and thus fail to achieve best adversarial robustness. In summary, the above observations prove the reasons for the adversarial vulnerabilities of HGNNs.

(3) HAN-RoHe also successfully defends the non-attributed HG Aminer. The defense model SimP and Jaccard, relying on the original feature (i.e., attribute), hence cannot be directly applied to Aminer. While our RoHe relies on node embedding rather than original features, thus can still enhance the robustness of HAN on Aminer.

### Generalization Performance of RoHe

**Generalization performance on random noise.** We evaluate the robustness of the proposed RoHe under random noise by linking the target node to random nodes. The results are shown in Figure 3. Results of more metrics and datasets are in Appendix. We can see HAN-RoHe achieves the best performance on most metrics and its performance only slightly drops with the increase of budget  $\Delta$ , meanwhile, Jaccard and SimP can also achieve comparable performance in comparison with RoHe on some metrics. The

Data	HGNNs	Clean	Attack		
			$\Delta = 1$	$\Delta = 3$	$\Delta = 5$
ACM	HAN	<b>0.926</b>	0.528	0.330	0.240
	HAN-RoHe	0.920	<b>0.904</b>	<b>0.902</b>	<b>0.882</b>
	MAGNN	<b>0.926</b>	0.711	0.647	0.589
	MAGNN-RoHe	0.916	<b>0.901</b>	<b>0.907</b>	<b>0.909</b>
	GTN	0.932	0.786	0.466	0.302
	GTN-RoHe <sub>T</sub>	<b>0.932</b>	<b>0.892</b>	<b>0.772</b>	<b>0.656</b>
DBLP	HAN	0.942	0.332	0.096	0.060
	HAN-RoHe	<b>0.942</b>	<b>0.936</b>	<b>0.864</b>	<b>0.808</b>
	MAGNN	<b>0.920</b>	0.620	0.494	0.416
	MAGNN-RoHe	0.898	<b>0.798</b>	<b>0.740</b>	<b>0.682</b>
	GTN	0.946	0.564	0.200	0.128
	GTN-RoHe <sub>T</sub>	<b>0.950</b>	<b>0.644</b>	<b>0.334</b>	<b>0.172</b>

Table 3: Results (Micro-F1) of RoHe on different HGNNs.

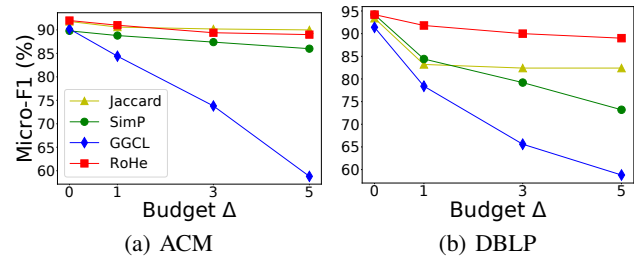


Figure 3: Results of HAN-RoHe under random noise.

reason is that these noise neighbors possibly have different features with the target node, and can be filtered well by the defense models based on feature similarity.

**Generalization performance on different HGNNs.** To demonstrate that our proposed defense framework is generic to other HGNNs, we generalize RoHe to MAGNN and GTN. The results are presented in Table ??, and results of more metrics are in Appendix. We first observe that the performance of all HGNNs dramatically drops under adversarial attacks, which demonstrates their common limitations. And RoHe can significantly improve the robustness of diverse HGNNs, especially for HAN and MAGNN. The reason is that the memory-consuming GTN can only be equipped by variant RoHe<sub>T</sub>, yielding limited improvement. Additionally, GTN and MAGNN show better robustness than HAN, since they can better encode structural information as explained in Section . We also find that all HGNNs are more vulnerable on DBLP than ACM, since the perturbations on P-A can be relieved by metapath PSP in ACM, which can be shown in Figure 4. But all metapaths in DBLP (APA, APCPA and APTPA) contain the perturbed relation type P-A, leading to less robustness.

### Robustness of Aggregations

**Analysis of node-level aggregation.** To verify whether RoHe can learn robust node-level attention values, here we take a paper node P3143 about “Wireless Communication” in ACM as an example. For clean data (Clean), the P3143 is connected to 6 neighbors with PAP metapath. For perturbed data (Attack), the attacker just injects one perturbation edge



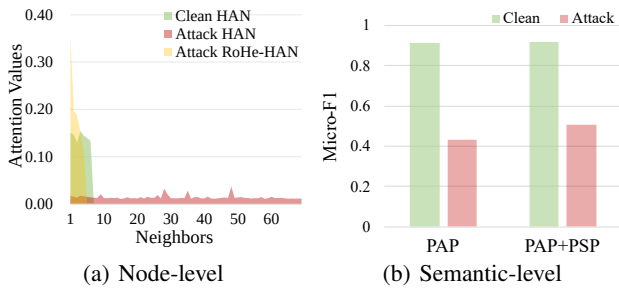


Figure 4: Analysis of node/semantic-level aggregations in ACM under Clean/Attack. (a) Node-level attention values of HAN(-RoHe) for paper P3143. (b) Results of different metapaths (only attacking P-A edges).

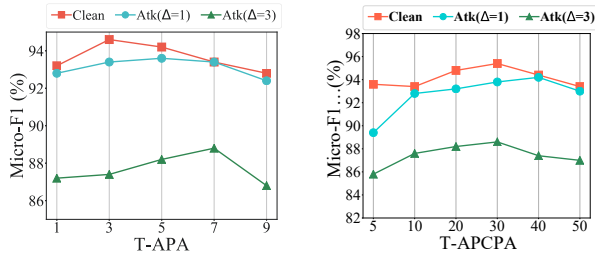


Figure 5: Analysis (Micro-F1) of parameter  $T$ .

by linking P3143 with author “Jiawei Han” who published 63 papers mainly about “Data Mining”. And the node-level attention values are shown in Figure 4 (a). Obviously, under attack, the original soft attention mechanism of HAN has to aggregate the 63 adversarial neighbors with positive values, leading to the distortion of P3143 embedding. While RoHe successfully filters these perturbations and assigns high confidence scores for true neighbors.

**Analysis of semantic-level aggregation.** To evaluate whether the rich semantics conveyed by metapaths can enhance the adversarial robustness of HGNNs by semantic-level aggregation, as shown in Figure 4 (b), we report the performance of different metapath sets in ACM dataset under Clean and Attack, where the type of adversarial edges are constrained to P-A. Obviously, the HAN under full metapaths achieves better robustness, since the perturbations on P-A can be relieved by the information within P-S for ACM.

### Parameter Study

We analyze hyper-parameters  $T$  which is the number of neighbors to be kept in purifier of HAN-RoHe, under scenarios of Clean and Attack with  $\Delta = \{1, 3\}$ . Here, we take metapaths APA and APCPA in DBLP dataset as examples. Figure 5 demonstrates how performance responds when threshold  $T$  increases. There exists an optimal  $T$  that delivers the best performance. When  $T$  is small, RoHe can only make use of little relevant neighbor information, which leads to inferior performance. When  $T$  increases, the purified receptive field involves more noise, leading to a higher chance of incorporating harmful neighbors, which negatively impacts the classification performance.

## Related Work

**Heterogeneous graph neural networks.** Recently, HGNNs showed outstanding performance in various tasks. Roughly speaking, HGNNs fell into two categories: (1) Directly aggregating metapath-based neighbors. HAN (Wang et al. 2019b) proposed directly aggregated metapath-based neighbors with node-level attention. Then MAGNN (Fu et al. 2020) extended HAN by considering the intermediate nodes along metapath. GTN (Yun et al. 2019) further automatically identified the useful metapaths in the process of learning node embeddings. (2) Indirectly aggregating multi-hop neighbors. HGT (Hu et al. 2020) and R-GCN (Schlichtkrull et al. 2018) indirectly incorporated long-range neighbors through message passing across layers. Here we focus on the former type, which is widely used in many safety-related tasks (Hu et al. 2019; Zhong et al. 2020; Zhang et al. 2019c).

**Adversarial robust on graphs.** Recently, a magnitude of adversarial attacks were introduced for homogeneous graphs (Zügner, Akbarnejad, and Günnemann 2018; Li et al. 2020; Ma, Ding, and Mei 2020; Sun et al. 2018), pointing out their sensitivity regarding such attacks. However, there are few existing investigations on the adversarial attacks for HGs (Hou et al. 2019; Pezeshkpour, Tian, and Singh 2019; Zhang et al. 2019b), and they all focused on non-GNN based methods (e.g., metapath2vec (Hou et al. 2019)). This paper sheds the first light on this important problem. On the other side, these adversarial attacks works also triggered the research on adversarial defenses on GNNs (Wu et al. 2019; Jin et al. 2020). With a unique aggregation mechanism, HGNNs show different adversarial vulnerabilities from GCNs and need additional specially designed defense solutions.

## Conclusion

In this paper, we introduce the first study on the adversarial robustness of HGNNs. Our extensive experiments show that HGNNs are highly fragile to topology adversarial attacks in comparison with GCNs, which can be attributed to the facts of perturbation enlargement and soft attention values. To address them, we propose an effective robust HGNN framework RoHe by equipping an attention purifier, which can prune unreliable neighbors based on topology and feature, alleviating the above vulnerabilities of HGNNs. Experiments on various datasets and multiple HGNNs show the effectiveness of RoHe. In future work, we will explore how to make full use of multiple aspects of information based on metapath to further improve robustness.

## Acknowledgments

This work is supported in part by the National Natural Science Foundation of China (No. U20B2045, 62172052, 61772082, 61702296, 62002029), the Fundamental Research Funds for the Central Universities 2021RC28, and BUPT Excellent Ph.D. Students Foundation (No. CX2021202). This work is also sponsored by CCF-Ant Group Research Fund.

## References

- Bo, D.; Wang, X.; Shi, C.; and Shen, H. 2021. Beyond Low-frequency Information in Graph Convolutional Networks. 3950–3957. AAAI Press.
- Chaudhari, S.; Polatkan, G.; Ramanath, R.; and Mithal, V. 2019. An Attentive Survey of Attention Models. *CoRR*, abs/1904.02874.
- Dong, Y.; Chawla, N. V.; and Swami, A. 2017. meta-path2vec: Scalable Representation Learning for Heterogeneous Networks. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Halifax, NS, Canada, August 13 - 17, 2017*, 135–144. ACM.
- Fu, X.; Zhang, J.; Meng, Z.; and King, I. 2020. MAGNN: Metapath Aggregated Graph Neural Network for Heterogeneous Graph Embedding. In *WWW*, 2331–2341.
- Goodfellow, I. J.; Shlens, J.; and Szegedy, C. 2015. Explaining and Harnessing Adversarial Examples. In *ICLR*.
- Hou, S.; Fan, Y.; Zhang, Y.; Ye, Y.; Lei, J.; Wan, W.; Wang, J.; Xiong, Q.; and Shao, F. 2019.  $\alpha$ Cyber: Enhancing Robustness of Android Malware Detection System against Adversarial Attacks on Heterogeneous Graph based Model. In *CIKM*, 609–618.
- Hu, B.; Fang, Y.; and Shi, C. 2019. Adversarial Learning on Heterogeneous Information Networks. In *KDD*, 120–129.
- Hu, B.; Zhang, Z.; Shi, C.; Zhou, J.; Li, X.; and Qi, Y. 2019. Cash-Out User Detection Based on Attributed Heterogeneous Information Network with a Hierarchical Attention Mechanism. In *AAAI*, 946–953.
- Hu, Z.; Dong, Y.; Wang, K.; and Sun, Y. 2020. Heterogeneous Graph Transformer. In *WWW*, 2704–2710.
- Jin, W.; Derr, T.; Wang, Y.; Ma, Y.; Liu, Z.; and Tang, J. 2021. Node Similarity Preserving Graph Convolutional Networks. In *WSDM*.
- Jin, W.; Li, Y.; Xu, H.; Wang, Y.; and Tang, J. 2020. Adversarial Attacks and Defenses on Graphs: A Review and Empirical Study. *CoRR*.
- Kipf, T. N.; and Welling, M. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *ICLR*.
- Li, J.; Zhang, H.; Han, Z.; Rong, Y.; Cheng, H.; and Huang, J. 2020. Adversarial Attack on Community Detection by Hiding Individuals. In *WWW*, 917–927.
- Ma, J.; Ding, S.; and Mei, Q. 2020. Towards More Practical Adversarial Attacks on Graph Neural Networks. In *NeurIPS*.
- Pezeshkpour, P.; Tian, Y.; and Singh, S. 2019. Investigating Robustness and Interpretability of Link Prediction via Adversarial Modifications. In *NAACL-HLT*, 3336–3347.
- Schlichtkrull, M. S.; Kipf, T. N.; Bloem, P.; van den Berg, R.; Titov, I.; and Welling, M. 2018. Modeling Relational Data with Graph Convolutional Networks. In *ESWC*, 593–607.
- Shi, C.; Li, Y.; Zhang, J.; Sun, Y.; and Yu, P. S. 2017. A Survey of Heterogeneous Information Network Analysis. *IEEE Trans. Knowl. Data Eng.*, 29(1): 17–37.
- Sun, L.; Dou, Y.; Yang, C.; Wang, J.; Yu, P. S.; and Li, B. 2018. Adversarial Attack and Defense on Graph Data: A Survey. *CoRR*.
- Wang, M.; Zheng, D.; Ye, Z.; Gan, Q.; Li, M.; Song, X.; Zhou, J.; Ma, C.; Yu, L.; Gai, Y.; Xiao, T.; He, T.; Karypis, G.; Li, J.; and Zhang, Z. 2019a. Deep Graph Library: A Graph-Centric, Highly-Performant Package for Graph Neural Networks. *CoRR*.
- Wang, X.; Ji, H.; Shi, C.; Wang, B.; Ye, Y.; Cui, P.; and Yu, P. S. 2019b. Heterogeneous Graph Attention Network. In *WWW*, 2022–2032.
- Wu, H.; Wang, C.; Tyshetskiy, Y.; Docherty, A.; Lu, K.; and Zhu, L. 2019. Adversarial Examples for Graph Data: Deep Insights into Attack and Defense. In *IJCAI*, 4816–4823.
- Yun, S.; Jeong, M.; Kim, R.; Kang, J.; and Kim, H. J. 2019. Graph Transformer Networks. In *NeurIPS*, 11960–11970.
- Zhang, C.; Song, D.; Huang, C.; Swami, A.; and Chawla, N. V. 2019a. Heterogeneous Graph Neural Network. In *KDD*, 793–803.
- Zhang, H.; Zheng, T.; Gao, J.; Miao, C.; Su, L.; Li, Y.; and Ren, K. 2019b. Data Poisoning Attack against Knowledge Graph Embedding. In *IJCAI*, 4853–4859.
- Zhang, X.; and Zitnik, M. 2020. GNNGuard: Defending Graph Neural Networks against Adversarial Attacks. In *NeurIPS*.
- Zhang, Y.; Fan, Y.; Ye, Y.; Zhao, L.; and Shi, C. 2019c. Key Player Identification in Underground Forums over Attributed Heterogeneous Information Network Embedding Framework. In *CIKM*, 549–558.
- Zhong, Q.; Liu, Y.; Ao, X.; Hu, B.; Feng, J.; Tang, J.; and He, Q. 2020. Financial Defaulter Detection on Online Credit Payment via Multi-view Attributed Heterogeneous Information Network. In *WWW*, 785–795.
- Zhu, D.; Zhang, Z.; Cui, P.; and Zhu, W. 2019. Robust Graph Convolutional Networks Against Adversarial Attacks. In *KDD*, 1399–1407.
- Zügner, D.; Akbarnejad, A.; and Günnemann, S. 2018. Adversarial Attacks on Neural Networks for Graph Data. In *KDD*, 2847–2856.