



Delft University of Technology

## Private Graph Extraction via Feature Explanations

Olatunji, Iyiola E ; Rathee, Mandeep; Funke, Thorben ; Khosla, M.

### DOI

[10.56553/popets-2023-0041](https://doi.org/10.56553/popets-2023-0041)

### Publication date

2023

### Document Version

Final published version

### Published in

Proceedings on Privacy Enhancing Technologies 2023(2)

### Citation (APA)

Olatunji, I. E., Rathee, M., Funke, T., & Khosla, M. (2023). Private Graph Extraction via Feature Explanations. In *Proceedings on Privacy Enhancing Technologies 2023(2)* (pp. 59-78) <https://doi.org/10.56553/popets-2023-0041>

### Important note

To cite this publication, please use the final published version (if applicable). Please check the document version above.

### Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

### Takedown policy

Please contact us and provide details if you believe this document breaches copyrights. We will remove access to the work immediately and investigate your claim.

# Private Graph Extraction via Feature Explanations

Iyiola E. Olatunji  
L3S Research Center  
Hannover, Germany  
iyiola@l3s.de

Thorben Funke  
L3S Research Center  
Hannover, Germany  
tfunke@l3s.de

Mandeep Rathee  
L3S Research Center  
Hannover, Germany  
rathee@l3s.de

Megha Khosla  
TU Delft  
Delft, Netherlands  
m.khosla@tudelft.nl

## ABSTRACT

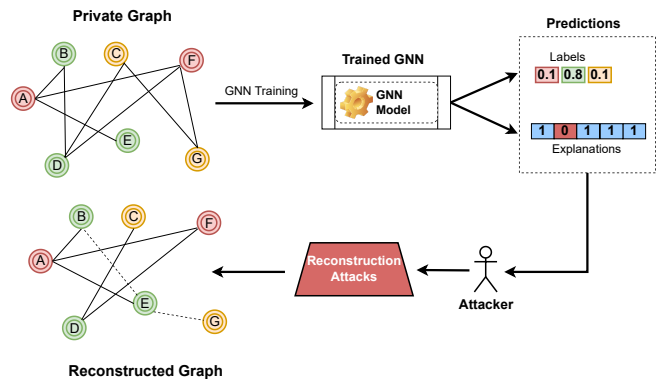
Privacy and interpretability are two important ingredients for achieving trustworthy machine learning. We study the interplay of these two aspects in graph machine learning through graph reconstruction attacks. The goal of the adversary here is to reconstruct the graph structure of the training data given access to model explanations. Based on the different kinds of auxiliary information available to the adversary, we propose several graph reconstruction attacks. We show that additional knowledge of post-hoc feature explanations substantially increases the success rate of these attacks. Further, we investigate in detail the differences between attack performance with respect to three different classes of explanation methods for graph neural networks: gradient-based, perturbation-based, and surrogate model-based methods. While gradient-based explanations reveal the most in terms of the graph structure, we find that these explanations do not always score high in utility. For the other two classes of explanations, privacy leakage increases with an increase in explanation utility. Finally, we propose a defense based on a randomized response mechanism for releasing the explanations, which substantially reduces the attack success rate. Our code is available at <https://github.com/iyempissy/graph-stealing-attacks-with-explanation>.

## KEYWORDS

privacy risk, model explanations, graph reconstruction attacks, private graph extraction, graph neural networks, attacks

## 1 INTRODUCTION

Graphs are highly informative, flexible, and natural ways of representing data in various real-world domains. Graph neural networks (GNNs) [19, 23, 37] have emerged as the standard tool to analyze graph data that is non-euclidean and irregular in nature. GNNs have gained state-of-the-art results in various graph analytical tasks ranging from applications in biology, and healthcare [2, 7] to recommending friends in a social network [12]. GNNs' success can be attributed to their ability to extract powerful latent features via complex aggregation of neighborhood aggregations [16, 46].



**Figure 1: The importance scores for the features, as provided by an explanation, can be exploited to infer the graph structure. Here, we show an example of a binary explanation where a score of 1 indicates that the corresponding feature is part of the explanation.**

However, these models are inherently black-box and complex, making it extremely difficult to understand the underlying reasoning behind their predictions. With the growing adoption of these models in various sensitive domains, efforts have been made to explain their decisions in terms of feature as well as neighborhood attributions. Model explanations can offer insights into the internal decision-making process of the model, which builds the trust of the users. Moreover, owing to the current regulations [28] and guidelines for designing trustworthy AI systems, several proposals advocate for deploying (automated) model explanations [17, 31].

Nevertheless, releasing additional information, such as explanations, can have adverse effects on the privacy of the training data. While the risk to privacy due to model explanations exists for machine learning models in general [34], it can have more severe implications for graph neural networks. For instance, several works [8, 27] have established the increased vulnerability of GNNs to privacy attacks due to the additional encoding of graph structure in the model itself. We initiate the *first investigation* of the effect of releasing feature explanations for graph neural networks on *the leakage of private information* in the training data.

To analyze the information leakage due to explanations, we take the perspective of an adversary whose goal is to infer the hidden connections among the training nodes. Consider a setting where

This work is licensed under the Creative Commons Attribution 4.0 International License. To view a copy of this license visit <https://creativecommons.org/licenses/by/4.0/> or send a letter to Creative Commons, PO Box 1866, Mountain View, CA 94042, USA.  
Proceedings on Privacy Enhancing Technologies 2023(2), 59–78  
© 2023 Copyright held by the owner/author(s).  
<https://doi.org/10.56553/popets-2023-0041>

the user has access to node features and labels but not the graph structure among the nodes. For example, node features could be part of the public profile of various individuals. The graph structure among the nodes could be some kind of social network which is private. Now, let's say that the user wants to obtain a trained GNN on her data while providing the node features and labels to a central authority that has access to the private graph structure. We ask here the question: *how much do the feature-based explanations leak information about the graph structure used to train the GNN model?* We quantify this information leakage via several graph reconstruction attacks. A visual illustration is shown in Figure 1. Specifically, our threat model consists of two main settings: first, where the adversary has access only to feature explanations, and second, where the adversary has additional information on node features/labels. Note that we only focus on feature-based explanations for GNNs in this work. For other explanation types such as node or edge explanations, as the returned nodes and edges belong to the node's neighborhood, it becomes trivial to reconstruct the graph structure.

### 1.1 Our Contributions and Findings

Ours is the first work to analyze the risks of releasing feature explanations in GNNs on the privacy of the relationships/connections in the training nodes. We quantify the information leakage via explanations by the success rate of several graph reconstruction attacks. Our attacks range from a simple explanation similarity-based attack to more complex attacks exploiting graph structure learning techniques. Besides, we provide a thorough analysis of information leakage via feature-based explanations produced by three classes of GNN post-hoc explanation methods, including *gradient-based*, *perturbation-based*, and *surrogate model-based*.

To analyze the differences in the robustness of the explanation methods to our privacy attacks, we investigate the explanation utility in terms of *faithfulness* and *sparsity*. We find that the gradient-based methods are the most susceptible to graph reconstruction attacks even though the corresponding explanations are the least faithful to the model. In other words, they have low utility. This is an important finding as the corresponding explanations could release a large amount of private information without offering any utility to the user. The perturbation-based approach ZORRO and its variant show the highest explanation utility as well as a high success rate of the attack, pointing to the expected trade-off between privacy and explanation utility.

We perform our study over three types of datasets with varying properties. For instance, our first dataset (CORA) has a large number of binary features, but the feature space is very sparse. Our second dataset (CORAML) has fewer but denser features. Our final dataset (BITCOIN) has a very small number of features. We find that the information leakage varies with explanation techniques as well as the feature size of the dataset. The dataset (BITCOIN) with the smallest feature size (8) is the most difficult to attack. Through experiments on three additional datasets (CITESEER, PUBMED, and CREDIT), we later show that although a small feature size might limit the information leakage, the correlation between node features with node connections plays a significant role. All baseline attacks which rely on knowledge of only features and labels perform no better than a random guess. In such a case, explanation-based attacks provide

an improvement, though not very huge, in inferring private graph structure information. We also perform additional experiments to evaluate the effectiveness of the attacks on three other datasets. Also, we perform a subgraph reconstruction attack when only a subset of the explanations is available to the attacker.

Finally, we develop a perturbation-based defense for releasing feature-based explanations. Our defense employs a randomized response mechanism to perturb the individual explanation bits. We show that our defense reduces the attack to a random guess with a small drop in the explanation utility. Our code is available at <https://github.com/iyempissy/graph-stealing-attacks-with-explanation>.

## 2 PRELIMINARIES

### 2.1 Graph Neural Networks

Graph neural networks (GNNs) [19, 23, 37] are a special family of deep learning models designed to perform inference on graph-structured data. The variants of GNNs, such as graph convolutional network (GCN), compute the representation of a node by utilizing its feature representation and that of the neighboring nodes. That is, GNNs compute node representations by recursive aggregation and transformation of feature representations of its neighbors.

Let  $G = (V, E)$  denote a graph where  $V$  is the node-set, and  $E$  represents the edges or links among the nodes. Furthermore, let  $\mathbf{x}_i^{(\ell)}$  be the feature representation of node  $i$  at layer  $\ell$ ,  $\mathcal{N}_i$  denote the set of its 1-hop neighbors and  $\theta$  is a learnable weight matrix. Formally, the  $\ell$ -th layer of a graph convolutional operation can be described as

$$\mathbf{z}_i^{(\ell)} = \text{AGGREGATION}^{(\ell)} \left( \left\{ \mathbf{x}_i^{(\ell-1)}, \left\{ \mathbf{x}_j^{(\ell-1)} \mid j \in \mathcal{N}_i \right\} \right\} \right) \quad (1)$$

$$\mathbf{x}_i^{(\ell)} = \text{TRANSFORMATION}^{(\ell)} \left( \mathbf{z}_i^{(\ell)} \right) \quad (2)$$

Then a softmax layer is applied to the node representations at the last layer (say  $L$ ) for the final prediction of the node classes  $C$

$$\mathbf{y} \leftarrow \text{argmax}(\text{softmax}(\mathbf{z}_i^{(L)}\theta)), \quad (3)$$

GNNs have been shown to possess increased vulnerability to privacy attacks due to encoding of additional graph structure in the model [8, 27]. We further investigate the privacy risks of releasing post-hoc explanations for GNN models.

### 2.2 Explaining Graph Neural Networks

GNNs are deep learning models which are inherently black-box or non-interpretable. This black-box behavior becomes more critical for applying them in sensitive domains like medicine, crime, and finance. Consequently, recent works have proposed post-hoc explainability techniques to explain the decisions of an already-trained model. In this work, we are concerned with the task of node classification, where the model is trained to predict node labels in a graph. For such a task, an instance is a single node. An explanation for a node usually consists of a subset of the most important features as well as a subset of its neighboring nodes/edges responsible for the model's prediction. Depending on the explanation method, the importance is usually quantified either as a continuous score (also referred to as a soft mask) or a binary score (also called a hard

mask). In this work, we consider three popular classes of explanation methods: *gradient-based* [3, 29, 32], *perturbation-based* [16, 46], and *surrogate* [21] methods.

**Gradient-based Methods.** These approaches usually employ the gradients of the target model’s prediction with respect to input features as importance scores for the node features. We use two gradient-based methods in our study, namely **GRAD**[35] and Grad-Input (**GRAD-I**) [36]. For a given graph  $G$  and the trained GNN model  $f(X, G, \theta)$  (where  $\theta$  is the set of parameters and  $X$  is the features matrix), GRAD generates an explanation  $\mathcal{E}_X$  by assigning continuous valued importance scores to the features. For node  $i$  and  $\mathbf{x}_i$  is its features vector, The score is calculated by  $\frac{\partial f}{\partial \mathbf{x}_i}$ . GRAD-I transforms GRAD explanation by an element-wise multiplication with the input features ( $\mathbf{x}_i \odot \frac{\partial f}{\partial \mathbf{x}_i}$ ).

**Perturbation-based Methods.** Perturbation-based methods obtain soft or hard masks over the features/nodes/edges as explanations by monitoring the change in prediction with respect to different input perturbations. We use two methods from this class: **ZORRO** [16] and **GNNEXP** [46]. ZORRO learns discrete masks over input nodes and node features as explanations using a greedy algorithm. It optimizes a fidelity-based objective that measures how the new predictions match the original predictions of the model by fixing the selected nodes/features and replacing the others with random noise values. This returns a hard mask for the explanations. GNNEXP learns soft masks over edges and features by minimizing the cross-entropy loss between the predictions of the original graph and the predictions of the newly obtained (masked) graph. We also utilize the ZORRO variant that provides soft explanation masks called **ZORRO-S**. ZORRO-S relaxes the argmax in ZORRO’s objective with a softmax, such that the masked are retrievable with standard gradient-based optimization. Together with the regularization terms of GNNEXP, ZORRO-S learns sparse soft masks.

**Surrogate Methods.** Surrogate methods fit a simple and interpretable model to a sampled local dataset corresponding to the query node. For example, the sampled dataset can be generated from the neighbors of the given query node. The explanations from the surrogate model are then used to explain the original predictions. We use GraphLime[21], which we denote as **GLIME** from this class. As an interpretable model, it uses the global feature selection method HSIC-Lasso [44]. The sampled dataset consists of the node and its neighborhood. The set of the most important features returned by HSIC-Lasso is used as an explanation for the GNN model.

*Remark:* Please note that in this work, we assume that only feature-based explanations are released to the user. In the presence of node/edge explanations, an adversary can trivially reconstruct large parts of the neighborhood as the returned nodes/edges are part of the node’s original neighborhood.

**2.2.1 Measuring Explanation Quality.** We measure the quality of the explanation by its ability to approximate the model’s behavior which is referred to as *faithfulness*. As the groundtruth for explanations is not available, we use the *RDT-Fidelity* proposed by [16] to measure faithfulness. The corresponding fidelity score measure how the original and new predictions match by fixing the

selected nodes/features and replacing the others with random noise values. Formally, the *RDT-Fidelity* of explanation  $\mathcal{E}_X$  corresponding to explanation mask  $M(\mathcal{E}_X)$  with respect to the GNN  $f$  and the noise distribution  $\mathcal{N}$  is given by

$$\mathcal{F}(\mathcal{E}_X) = \mathbb{E}_{Y_{\mathcal{E}_X} | Z \sim \mathcal{N}} \left[ \mathbb{1}_{f(X) = f(Y_{\mathcal{E}_X})} \right], \quad (4)$$

where the perturbed input is given by

$$\tilde{X}_{\mathcal{E}_X} = X \odot M(\mathcal{E}_X) + Z \odot (\mathbb{1} - M(\mathcal{E}_X)), Z \sim \mathcal{N}, \quad (5)$$

where  $\odot$  denotes an element-wise multiplication, and  $\mathbb{1}$  is a matrix of ones with the corresponding size.

**Sparsity.** We further note that by definition, the complete input is faithful to the model. Therefore, we further measure the sparsity of the explanation. A meaningful explanation should be sparse and only contain a small subset of features most predictive of the model decision. We use the entropy-based sparsity definition from [16] as it is applicable for both soft and hard explanation masks. Let  $p$  be the normalized distribution of explanation (feature) masks. Then the sparsity of an explanation is given by the entropy  $H(p)$  over the mask distribution,

$$H(p) = - \sum_{f \in M} p(f) \log p(f).$$

Note that the entropy here is bounded by  $\log(|M|)$ , where  $M$  corresponds to the size of the feature set. The lower the entropy, the sparser the explanation. We will utilize these two metrics used in measuring explanation quality to argue about the differences in the attack performance.

## 3 THREAT MODEL AND ATTACK METHODOLOGY

### 3.1 Motivation

We consider the setting in which different data holders hold the features and graph (adjacency matrix) in practice. Specifically, the central trusted server has access to the graph structure and trains a GNN model using the node features and labels provided by the user. The user can further query the trained GNN model by providing node features. As the features are already known to the user, revealing *feature explanations* for the prediction might be considered a safe way to increase the user’s trust in the model. We investigate such a scenario and uncover the increased privacy risks of releasing *feature explanations even if the original node features/labels are already known to the adversary*.

Our setting is inspired by previous works and practical scenarios [4, 10, 18, 45, 48]. We elaborate on one such scenario in the following. Consider a multinational company with several subdivisions. For simplicity, we consider two subdivisions (the marketing and product departments). The marketing department collects the interaction between users (edge information), while the product department collects data on user behavior (features). Under privacy law(e.g., GDPR), such user interactions cannot be shared because they are collected for different purposes. However, they aim to train a GNN model jointly since it would benefit both teams. The product department then sends the node features to the marketing department, which in addition, uses their edge information to train a predictive model and releases an API (similar to MLaaS). The API

returns both the prediction for the features and the corresponding feature explanations for such predictions. The product department, having observed how useful the exploratory analysis is, aims to recover the private edges using the information returned by the API (explanations) and their node features.

### 3.2 Threat Model

We consider two scenarios: **first**, in which the adversary has access to node features and/or labels and obtains additional access to feature explanations, and **second**, in which the adversary has access only to feature explanations. The goal of the adversary is to infer the private connections of the graph used to train the GNN model.

Corresponding to the above settings, we categorize our attacks as *explanation augmentation* (when explanations are augmented with other information to launch the attack) and *explanation-only* attacks. Through our five proposed attacks, we investigate the privacy leakage of the explanations generated by six different explanation methods.

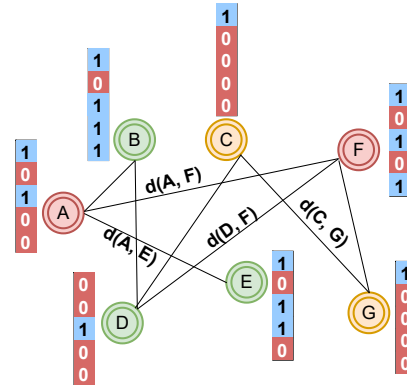
Our simple similarity-based explanation-only attack already allows us to quantify the additional information that the feature-based explanation encodes about the graph structure. Our explanation augmentation attacks are based on the graph structure learning paradigm, which allows us to effectively integrate additional known information in learning the private graph structure. Besides, our explanation augmentation attacks also result in a successfully trained GNN model without the knowledge of the true graph structure, offering an additional advantage to the adversary.

### 3.3 Attack Methodologies

Here, we provide a detailed description of our two types of attacks. We commence with the *explanation-only attack* in which we utilize only the provided explanation to launch the attack, followed by the *explanation augmentation* attacks in which more information such as node labels or/and features are exploited in addition to the explanation. The taxonomy of our attacks based on the attacker’s knowledge is presented in Table 1.

**Explanation-only Attack.** This is an unsupervised attack in which the attacker only has access to the explanations and does not have access to the features or the labels. The attacker measures the distance between each pair of explanation vectors and assigns an edge between them if their distance is small. The intuition is that the model might assign similar labels to connected nodes, which leads to similar explanations. We experimented with various distance metrics, but cosine similarity performs best across all datasets. We refer to this similarity-based attack as **EXPLAINSIM**. This attack is illustrated in Figure 2. We report the performance using the area under the receiver operating characteristic curve (AUC) and average precision (AP). Rather than choosing one threshold which could be domain-dependent, these metrics allow us to provide a holistic evaluation across the complete range of thresholds.

**Explanation Augmentation Attacks.** Towards explanation augmentation attacks, we leverage the graph structure learning paradigm of [13]. In particular, we employ two *generator modules* for generating graph edges corresponding to features and explanations, respectively. The generators are trained using feature/explanation



**Figure 2: EXPLAINSIM attack.** Each node is assigned a feature explanation vector where blue (1) and red (0) indicate whether the feature is part of the explanation or not. The attacker then assigns an edge by computing the pairwise similarity represented as  $d(\text{node}_{e_i}, \text{node}_{e_j})$  between nodes’ explanation vectors. We show the representation of ZORRO’s explanation for easier visualization.

reconstruction-based losses and node classification loss. We commence by describing the common architecture of the attack model followed by its concrete usage in the four attack variations in Section 3.3.1.

**Generators.** The two generators take the node features/explanations as input and output two adjacency matrices. We employ the full parameterization (FP) approach to model the generators similar to those used in [13, 15]. In other words, each element of the adjacency matrix  $\tilde{A}$  is treated as a separate learnable parameter, i.e., the adjacency matrix is fully parameterized. We initialize these learnable parameters by the cosine similarity among the feature/explanation vectors of the corresponding nodes. Let the output adjacency matrix function be given as  $\tilde{A} = G_{FP}(X; \theta_G) = \theta_G$  where  $\theta_G \in \mathbb{R}^{n \times n}$  and  $G_{FP}(\cdot; \cdot)$  denotes the generator function. To obtain a symmetric adjacency matrix with all positive elements, we perform the following transformation

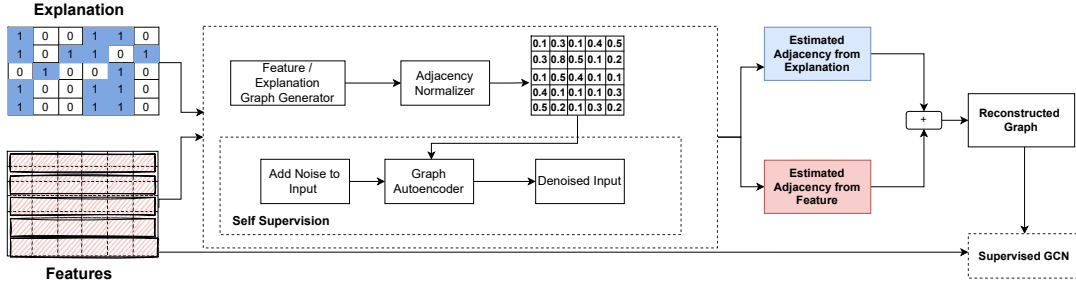
$$A = D^{-\frac{1}{2}} \left( \frac{P_{[0,1]}(\tilde{A}) + P_{[0,1]}(\tilde{A})^T}{2} \right) D^{-\frac{1}{2}},$$

where P is a non-negative function defined by

$$P_{[0,1]}[x] = \begin{cases} 0 & x < 0, \\ 1 & x > 1, \\ x & \text{otherwise.} \end{cases}$$

The final adjacency matrix is computed by adding the matrices corresponding to the two generators. Any element greater than one is trimmed to 1. In other words, a cell in the adjacency matrix follows a Bernoulli distribution. The final graph is then generated by sampling from the learned Bernoulli distributions. More precisely, we select edge  $(i, j)$  with probability  $p$  where  $p \in [0, 1]$  is the learnt weight for cell  $(i, j)$ .

**Training With Self-supervision and Node Classification Loss.** We remark that only some nodes used to train the attack



**Figure 3: Overview of GSEF.** The generator takes node features and explanations as input and outputs an adjacency matrix which may be non-normalized by the adjacency normalizer. The normalized adjacency matrix is used in predicting both the class labels and reconstructing the node features (explanations) by the denoising autoencoders. The final reconstructed adjacency is the one that minimizes the reconstruction error on each of the node features and explanations and the loss of the class label prediction.

model are labeled. To learn effective representations of both labeled and unlabeled nodes, we employ a self-supervision module in addition to training with the supervised classification loss. In particular, for the self-supervision module, we use *denoising graph autoencoders*, which takes noisy features/explanations and the graph sampled from the generator as input. The goal is to reconstruct the true node features and explanations.

Let  $DAE$  and  $DAE_{\mathcal{E}_X}$  denote the feature and explanation reconstruction modules respectively. They take the noisy node features(explanations) as input and output a denoised version of the features (explanations) with the same dimension. To generate the noisy features/explanations, we first select random indices. We then add random noise to feature/explanation values corresponding to the selected indices. Concretely, for binary features/explanations, we randomly flip  $r = 20\%$  of the index with original values 1 to 0. For features and explanations with continuous values, we add independent Gaussian noise to  $r = 20\%$  percent of feature/explanation elements (corresponding to selected indices). To train our self-supervision module, we employ binary cross entropy loss for binary features/explanations and mean squared error for continuous features/explanations. Let  $\mathcal{L}_{DAE}$  and  $\mathcal{L}_{DAE_{\mathcal{E}_X}}$  denote the loss functions corresponding to feature and explanation reconstruction, respectively.

For supervised training, we employ a graph convolution network (GCN) which takes as input the node features and the graph sampled from the generator to predict node labels. Note that the GCN used here is not the target model. We use the cross-entropy loss ( $CE$ ) between the predicted labels ( $\hat{Y}$ ) and the ground-truth labels ( $Y$ ),  $\mathcal{L}_C = CE(Y, \hat{Y})$ .

The final training loss for private graph extraction is then given by

$$\mathcal{L} = \mathcal{L}_{DAE} + \mathcal{L}_{DAE_{\mathcal{E}_X}} + \mathcal{L}_C.$$

We refer to the above attack framework as **Graph Stealing with Explanations and Features (GSEF)**, and the schematic diagram is given in Figure 3. Besides, we have three attack variations that employ a single generator module, as described below.

**3.3.1 Attack Variations.** Besides GSEF, we have three attack variations that employ explanations (i) **GSEF-CONCAT** in which we

**Table 1: Attack taxonomy based on attacker’s knowledge of node features ( $X$ ), labels ( $Y$ ) and feature explanations ( $\mathcal{E}_X$ ).**

ATTACK	$X$	$Y$	$\mathcal{E}_X$
EXPLAINSIM	✗	✗	✓
GSEF	✓	✓	✓
GSEF-CONCAT	✓	✓	✓
GSEF-MULT	✓	✓	✓
GSE	✗	✓	✓

concatenate the node features and explanations and feed the concatenated input to a single generator module (ii) **GSEF-MULT** in which we perform element-wise multiplication between the features and the explanations and feed them into a single graph generator module. This is equivalent to assigning importance to the node features, which emphasize the essential characteristics of the nodes. Similar to GSEF-CONCAT, we reconstruct the adjacency matrix using one generator of Figure 3 and (iii) **GSE** in which the attacker only has access to the explanations and labels. Here, we also employ only one generator with explanations as input.

## 4 EXPERIMENTS

In this section, we present the experimental results to show the effectiveness of explanation-based attacks. Specifically, our experiments are designed to answer the following research questions:

RQ 1. *How does the knowledge of feature explanations influence the reconstruction of private graph structure?*

RQ 2. *What are the differences between explanation methods with respect to privacy leakage?*

RQ 3. *What is the additional advantage of the adversary (for example, in terms of the utility of the inferred information on a downstream task) on explanation augmentation attacks?*

RQ 4. *How much does the lack of knowledge about groundtruth node labels affect attack performance?*

## 4.1 Experimental Settings

### 4.1.1 Attack Baselines Without Explanations.

**FEATURESIM.** In this unsupervised attack, the attacker computes the pairwise similarity between pairs of the actual features to reconstruct the graph. Specifically, an edge exists between two nodes if the distance between their feature representation is low. We use cosine similarity as a measure of similarity because it performs better than other distance metrics.

**LSA [20].** LSA (Link stealing attack) is a black-box attack that assumes that the attacker has access to a dataset drawn from a similar distribution to the target data (shadow dataset). Additionally, LSA knows the architecture of the target model and can train a corresponding shadow model that replicates the behavior of the target model. The goal of the attack is to infer sensitive links between nodes of interest. We compare our results with their proposed attack-2, where an attacker has access to the node features and labels. They trained a separate MLP model (reference model) using the available target attributes and their corresponding labels to obtain posteriors. Then, LSA computes the pairwise distance between posteriors obtained from the target model and that of the reference model for the nodes of interest. We use cosine similarity as the distance metric.

**GRAPHMI [48].** GRAPHMI is a white-box attack in which the attacker has access to the parameters of the target model, node features, all the node labels, and other auxiliary information like edge density. The goal of an attacker is to reconstruct the sensitive links or connections between nodes. The attack model uses the cross entropy loss between the true labels and the output posterior distribution of the target model, along with feature smoothness and adjacency sparsity constraints, to train a fully parameterized adjacency matrix. The graph is then reconstructed using the graph autoencoder module in which the encoder is replaced by learned parameters of the target model, and the decoder is a logistic function.

**SLAPS [13].** Since our attack model is built on top of the graph structure learning framework of SLAPS, we performed an experiment using the vanilla SLAPS. Given node features and labels, SLAPS aims to reconstruct the graph that works best for the node classification task.

**4.1.2 Target GNN Model.** We employ a 2-layer graph convolution network (GCN) [23] as our target GNN model. We used a learning rate of 0.001 and trained for 200 epochs.

**4.1.3 Evaluation Metrics.** Following the existing works [20, 48], we use the area under the receiver operating cost curve (AUC) and average precision (AP) to evaluate our attack. For all experiments (including baselines), we randomly sample an equal number of pairs of connected and unconnected nodes from the original graph to construct the test set. We measure the AUC and the AP on these randomly selected node pairs. We elaborate more in Appendix B. All our experiments were conducted for 10 different instantiations using PyTorch Geometric library [14] on 11GB GeForce GTX 1080 Ti GPU, and we report the mean values across all runs. The standard deviation of the results is in Appendix D.

**Table 2: Dataset statistics.**  $|V|$  and  $|E|$  denotes the number of nodes and edges respectively,  $C$ ,  $X_d$ , and  $\text{deg}$  denotes the number of classes, size of feature dimension and the average degree of the corresponding graph dataset.

	CORA	CORAML	BITCOIN
$ V $	2708	2995	3783
$ E $	10858	8226	28248
$X_d$	1433	300	8
$C$	7	7	2
$\text{deg}$	4.00	2.75	7.50

**4.1.4 Datasets.** We use three commonly used datasets chosen based on varying graph properties, such as their feature dimensions and structural properties. The task on all datasets is node classification. We present the data statistics in Table 2. We also performed additional experiments using PUBMED, CITESEER, and CREDIT datasets (See Appendix E.1).

**CORA.** The CORA dataset [33] is a citation dataset where each research article is a node, and there exists an edge between two articles if one article cites the other. Each node has a label that shows the article category. The features of each node are represented by a 0/1-valued word vector, which indicates the word’s presence or absence from the article’s abstract.

**CORAML.** In CORAML dataset [33], each node is a research article, and its abstract is available as raw text. In contrast to the above dataset, the raw text of the abstract is transformed into a dense feature representation. We preprocess the text by removing stop words, web page links, and special characters. Then, we generate Word2Vec [25] embedding of each word. Finally, we create the feature vector by taking the average over the embedding of all words in the abstract.

**BITCOIN.** The BITCOIN-Alpha dataset[24] is a signed network of trading accounts. Each account is represented as a node, and there is a weighted edge between any two accounts, which represents the trust between accounts. The maximum weight value is +10, indicating total trust, and the lowest is -10, which means total distrust. Each node is assigned a label indicating whether or not the account is trustworthy. The feature vector of each node is based on the rating by other users, such as the average positive or negative rating. We follow the procedures in [38] for generating the feature vectors.

## 4.2 Result Analysis

**4.2.1 Analysing the Information Leakage by Explanations (RQ 1).** The detailed results of different attacks are provided in Table 3. Our results show that the explanation-only (EXPLAINSIM) and explanation augmentation (GSEF) attacks for all explanation methods other than GLIME and GNNEXP outperform all baseline methods and by far reveal the most information about the private graph structure. We attribute the superior performance of GSEF to the multi-task learning paradigm that aims to reconstruct both the features and explanations. Our results also support our assumption

**Table 3: Attack performance and baselines. The best performing attack(s) on each explanation method is(are) highlighted in bold, and the second best attack(s) is(are) underlined.**

Exp	Attack	CORa		CORaML		BITCOIN	
		AUC	AP	AUC	AP	AUC	AP
Baseline	FEATURESIM	0.799	0.827	0.706	0.753	0.535	0.478
	Lsa [20]	0.795	0.810	0.725	0.760	0.532	0.500
	GRAPHMI [48]	0.856	0.830	0.808	0.814	0.585	0.518
	SLAPS [13]	0.736	0.776	0.649	0.702	0.597	0.577
GRAD	GSEF-CONCAT	0.734	0.773	0.640	0.705	0.527	0.515
	GSEF-MULT	0.678	0.737	0.666	0.730	0.264	0.383
	GSEF	<u>0.948</u>	<u>0.953</u>	<b>0.902</b>	<u>0.833</u>	<b>0.700</b>	<b>0.715</b>
	GSE	0.924	0.939	0.699	0.768	0.229	0.365
	EXPLAINSIM	<b>0.984</b>	<b>0.978</b>	<u>0.890</u>	<b>0.891</b>	<u>0.681</u>	<u>0.644</u>
GRAD-I	GSEF-CONCAT	0.734	0.775	0.674	0.734	0.525	0.527
	GSEF-MULT	0.691	0.742	0.717	0.756	0.252	0.380
	GSEF	<u>0.949</u>	<u>0.950</u>	<u>0.887</u>	<u>0.832</u>	<b>0.709</b>	<b>0.723</b>
	GSE	0.903	0.923	0.717	0.781	0.256	0.380
ZORRO	EXPLAINSIM	<b>0.984</b>	<b>0.979</b>	<b>0.903</b>	<b>0.899</b>	<u>0.681</u>	<u>0.644</u>
	GSEF-CONCAT	0.823	0.860	0.735	0.786	<u>0.575</u>	0.529
	GSEF-MULT	0.723	0.756	0.681	0.697	0.399	0.449
	GSEF	<b>0.884</b>	<b>0.880</b>	<u>0.776</u>	<u>0.820</u>	0.537	<u>0.527</u>
	GSE	0.779	0.810	0.722	0.777	<b>0.596</b>	<b>0.561</b>
ZORRO-S	EXPLAINSIM	<u>0.871</u>	<u>0.873</u>	<b>0.806</b>	<b>0.829</b>	0.427	0.485
	GSEF-CONCAT	0.907	0.922	<u>0.747</u>	<u>0.791</u>	<b>0.601</b>	<b>0.590</b>
	GSEF-MULT	0.794	0.815	0.712	0.740	0.490	0.491
	GSEF	<b>0.918</b>	<u>0.923</u>	<b>0.776</b>	<b>0.819</b>	<u>0.598</u>	<u>0.565</u>
	GSE	0.893	0.915	0.742	0.784	0.571	0.564
GLIME	EXPLAINSIM	<u>0.908</u>	<b>0.934</b>	0.732	0.787	0.484	0.496
	GSEF-CONCAT	<u>0.643</u>	<u>0.710</u>	<u>0.610</u>	<u>0.652</u>	<u>0.473</u>	<u>0.493</u>
	GSEF-MULT	0.516	0.522	0.517	0.528	0.264	0.371
	GSEF	<b>0.730</b>	<b>0.773</b>	<b>0.681</b>	<b>0.740</b>	<b>0.542</b>	<b>0.525</b>
	GSE	0.558	0.571	0.540	0.555	0.236	0.361
GNNEXP	EXPLAINSIM	0.505	0.524	0.520	0.523	0.504	0.512
	GSEF-CONCAT	0.614	0.650	0.653	0.705	0.467	0.489
	GSEF-MULT	<u>0.724</u>	<u>0.760</u>	<u>0.637</u>	<u>0.692</u>	0.390	0.454
	GSEF	<b>0.762</b>	<b>0.796</b>	<b>0.700</b>	<b>0.695</b>	<b>0.590</b>	<b>0.563</b>
	GSE	0.517	0.552	0.490	0.508	0.386	0.451
GNNEXP	EXPLAINSIM	0.537	0.541	0.484	0.508	<u>0.551</u>	<u>0.543</u>

**SUMMARY OF ATTACK COMPARISONS**

- The amount of information in the explanation alone for the graph structure can be quantified using the EXPLAINSIM attack.
- Explanation only (EXPLAINSIM) and explanation augmentation (GSEF) attacks for all explanation methods other than GLIME and GNNEXP outperform all baseline methods.
- Among the baseline approaches, the white-box access-based attack, GRAPHMI, performs the best, followed by FEATURESIM.
- The relatively good performance of FEATURESIM in datasets points to a high correlation of node features with node connections in all datasets other than BITCOIN.
- The information leakage for BITCOIN is limited by small feature size. Our results on additional datasets with varying feature dimensions, as discussed in Section 4.2.2, support our observation about the effect of small feature size on the attack’s success.
- For GLIME and GNNEXP, we observe that the explanation contains little information about the graph structure. The reason behind this is further revealed in the fidelity-sparsity analysis of the obtained explanations.

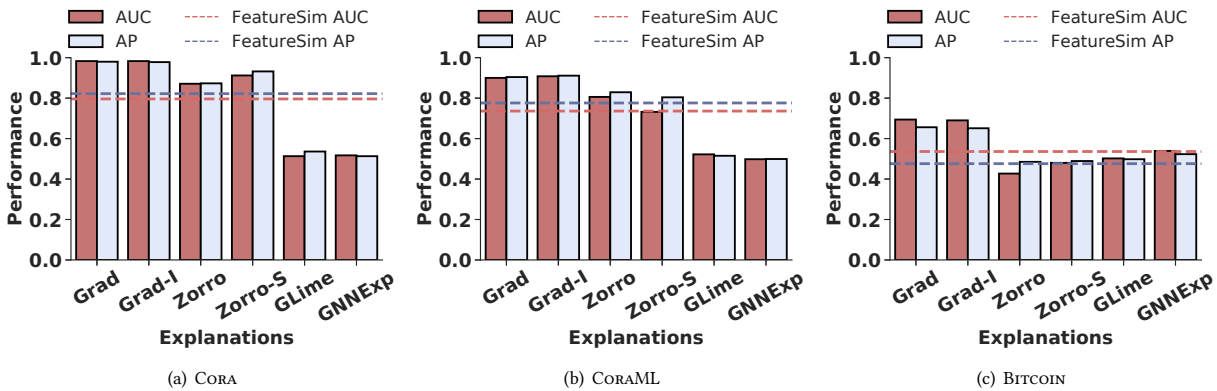
that a graph structure that is good for predicting the node labels is also suitable for predicting the node features and explanations.

In the following, we provide an in-depth result analysis.

- **Baseline Methods.** Among the baseline methods, GRAPHMI is the best-performing attack. This is not very surprising as GRAPHMI has white-box access to the target GNN in addition to access to node features and labels. On the contrary, the FEATURESIM attack, which only uses node features, shows competitive performance. This highlights the fact that the features alone are very informative of the graph structure of the studied datasets (except BITCOIN in which all baseline attacks almost fail with an AUC score of close to 0.5).

- **Comparison of the Privacy Leakage via Explanations With That of Features.** Figure 4 compares the performance of two similarity-based attacks using explanations (EXPLAINSIM) and features (FEATURESIM), respectively. Note that these attacks do not use any other information except explanations and features, respectively. Hence, allowing us to compare their information content. We note that except for GLIME and GNNEXP, EXPLAINSIM outperforms FEATURESIM for all datasets except BITCOIN. Moreover, for BITCOIN, both attacks fail (with AUC close to 0.5) except for gradient-based explanation methods.
- **Explanation Augmentation With Features and Labels.** Next, we compare the explanation augmentation attack GSEF with the





**Figure 4: Performance of explanation-only attack (EXPLAINSIM) on the different datasets. The adopted baseline is FEATURESIM which performs the pairwise similarities using the true node features.**

vanilla graph structure learning approach SLAPS, which only uses node features and labels. GSEF outperforms SLAPS (Figure 6), which points to the added utility of using explanations to reconstruct the graph structure.

- **Explanation Augmentation Attack Variants.** In Figure 5, we compare the performance of GSE, which have no access to the true features, with SLAPS, which utilizes the true features. We observe that on most datasets, the attack significantly outperforms SLAPS. This emphasizes that explanations encode the feature information albeit the importance. Comparing the performance of GSEF-CONCAT and GSEF-MULT with GSEF on all datasets shows that independently extracting the adjacency matrix from the features and explanations respectively and then combining the two adjacency matrices is better than combining the features and explanations at the input stage.
- **Attack on BITCOIN.** We observe that all baseline attacks fail for BITCOIN. We attribute this to the very small feature size (=8) and the small number of labels (=2). The attacks are not able to exploit the little information exposed by a small set of features/labels. Explanation-based attacks, especially in the case of gradient-based explanations, are more successful than baselines but less than for other datasets.

**4.2.2 Results on Additional Datasets.** We further investigate our observation of the impact of feature dimension and disasno on the attack success by performing experiments on the additional datasets (namely CITESEER, PUBMED, and CREDIT) with varying graph properties (details on datasets are provided in Section E.1). The detailed result is in the Appendix in Table 12. We observe that the result on CITESEER and PUBMED supports our observation of the effect of the feature size on the attack success on CORA and CORAML, respectively. First, a dataset with a high feature dimension is most vulnerable to the attack, as seen in the relative performance of CITESEER on all explanation methods and attacks. Secondly, on the GRAD and GRAD-I explanations, simply augmenting backpropagation-based explanations with the actual feature via concatenation (GSEF-CONCAT) or element-wise multiplication

(GSEF-MULT) does not encode the necessary information for launching a successful attack. Instead, it introduces new noise to the reconstructed graph because the explanations are generated such that larger gradients, which only reflect the sensitivity between input and output, may not accurately indicate the importance of the feature. We also observe this phenomenon on other explanations methods, including ZORRO, ZORRO-S, and GLIME, where GSE, which uses only explanations as input, performs better than GSEF-CONCAT and GSEF-MULT on the CORA, CORAML, CITESEER, and PUBMED datasets. Thirdly, explanation methods that are robust to input perturbation, such as ZORRO and ZORRO-S, achieve better attack performance than backpropagation-based explanations on all explanation augmentation attacks. Finally, GLIME and GNNEXP explanations are more robust to the attacks except on GSEF attack, which achieves a high AUC and AP scores.

Among all the additional datasets, the results of the baseline attacks (FEATURESIM, LSA, GRAPHMI, and SLAPS, which all rely only on feature information) on the CREDIT dataset are the least. Recall that the CREDIT dataset has only 13 features. This supports our observation that small feature size affects the attack’s success. Notably, the best performing baseline (GRAPHMI), which utilizes white-box access to the trained model, has a very low AUROC of 0.59 and 0.51 on the BITCOIN and CREDIT datasets, respectively. This is in contrast to the performance of GRAPHMI on CORA, CITESEER, PUBMED, and CORAML with AUROC of 0.86, 0.79, 0.80, and 0.81, respectively.

Furthermore, in comparison with BITCOIN, which has only 8 features, we observe that the performance of FEATURESIM on the CREDIT dataset is 34% higher (0.54 on BITCOIN and 0.72 on CREDIT), which points to high connectivity among similar nodes on the CREDIT dataset. This, in turn, boosts our attack performance and allows better reconstruction of the private graph. Therefore, the poor performance of our attacks on the BITCOIN dataset is due to disassociativity on the BITCOIN dataset. Moreover, it is worth noting that our explanation augmentation attack is highly successful on the CREDIT dataset.

**4.2.3 Differences in Privacy Leakage (RQ 2).** All explanation method leaks significant information via the reconstruction attacks except for GLIME and GNNEXP. We observe that for GLIME and GNNEXP, the explanation-based attacks do not perform better than the baselines, which do not utilize any explanation. Moreover, gradient-based methods are most vulnerable to privacy leakage. To understand the reason behind these observations, we investigate the explanation quality. We measure the goodness of the explanation by its ability to approximate the model’s behavior which is also referred to as *faithfulness*. As the groundtruth for explanations is not available, we use the RDT-Fidelity proposed by [16] to measure faithfulness. The results are shown in Table 4. We further note that by definition, the complete input is faithful to the model. Therefore, in addition, we measure the sparsity of the explanation. A meaningful explanation should be sparse and only contain a small subset of features most predictive of the model decision. We use the entropy-based sparsity definition from [16] as it is applicable for both soft and hard explanation masks. The results are shown in Table 4. We analyze the tradeoffs of privacy and explanation goodness in the following.

- **First**, we observe that GNNEXP has the lowest sparsity (the higher the entropy, the more uniform the explanation mask distribution, i.e., higher the explanation density). In other words, almost all features are marked equally important. Hence, it is not surprising that it shows a high fidelity. This is the main reason why EXPLAINSIM fails because there is no distinguishing power contained in the explanations.
- **Second**, we observe that gradient-based explanations (GRAD and GRAD-I) contain the most information about the graph structure though they have low fidelity, i.e., they do not reflect the model’s decision process. It appears that these two methods provide the most similar explanations for connected nodes. This is really the worst case when the explanations are not useful but leak maximum private information about the graph structure.
- **Third**, GLIME has the highest sparsity and lowest fidelity. GLIME runs the HSIC-Lasso feature selection method over a local dataset created using the node and its neighborhood. HSIC-Lasso is known to output a very small set of most predictive features [6] when used for global feature selection. But for the current setting of instance-wise feature selection, i.e., finding the most predictive features for decision over an instance/node, GLIME’s explanation turns out to be too short, which is neither faithful to the model nor contains any predictive information about the neighborhood.
- **Finally**, the explanations of ZORRO and ZORRO-S show the highest fidelity and intermediate sparsity, pointing to their high quality. The GSEF attack also obtains high AUC scores for two datasets pointing to the expected increased privacy risk with an increase in explanation utility.

**4.2.4 Adversary’s Advantage in Terms of Trained GNN Model (RQ 3).** Here, we formalize a quantitative advantage measure that captures the privacy risk posed by the different attacks. The attacker is at an advantage if she can train a well-performing model (on a downstream task) using the reconstructed graph. As the attack models based on graph structure learning implicitly train a GNN on the reconstructed graph, we quantify the attacker’s

**Table 4: RDT-Fidelity and sparsity (entropy) of different explanation methods. For fidelity, the higher the better. For sparsity, the lower the better**

Exp	CORA		CORAML		BITCOIN	
	Fidelity	Sparsity	Fidelity	Sparsity	Fidelity	Sparsity
<b>GRAD</b>	0.23	3.99	0.22	5.24	0.83	0.64
<b>GRAD-I</b>	0.19	3.99	0.20	5.30	0.82	0.64
<b>ZORRO</b>	0.89	1.83	0.96	3.33	0.99	0.37
<b>ZORRO-S</b>	0.98	2.49	0.84	2.75	0.95	0.96
<b>GLIME</b>	0.19	0.88	0.20	0.98	0.82	0.13
<b>GNNEXP</b>	0.74	7.27	0.55	5.70	0.90	2.05

advantage by the performance on the downstream task of node classification.

**HYPOTHESIS 1.** *If the explanations and the reconstructed graph can perform better on a downstream classification task with high confidence, then the reconstructed adjacency is a valid representation of the graph structure. Hence, the attacker has an advantage quantified by Equation 6.*

We define the attacker’s advantage as

$$\text{Advantage} = \mathcal{R}(f(\mathcal{E}_X; Adj_{rec}; \theta_W), y), \quad (6)$$

where  $f$  is a 2-layer GCN model parameterized by  $\theta_W$ ,  $\mathcal{E}_X$  is the explanation matrix,  $Adj_{rec}$  is the reconstructed graph by the attacker,  $y$  is the groundtruth label and  $\mathcal{R}$  is an advantage measure that compares the predictions on  $f$  and the groundtruth label. We use accuracy as the choice of  $\mathcal{R}$ .

We compare the results to that of SLAPS that uses the actual features and groundtruth label for learning the graph structure and the original performance (denoted by  $Max$ ) in which the model is trained with true features, labels, and graph structure. We analyze the attacker’s advantage corresponding to four attacks; GSEF-CONCAT, GSEF-MULT, GSEF, and GSE. The intuition is that if the attacker’s advantage is not better than SLAPS, then the best advantage an attacker can have is similar to having the actual feature and performing graph structure learning. Also, if the attacker’s advantage is greater or equal to  $Max$ , then the attacker has an equivalent advantage as she would have by possessing the actual feature and graph. An example use case of the attacker’s advantage is shown in Figure 7. Specifically, if a model trained with, say, Jane’s full data (true features and graph) and another trained only with her explanations (no graph or true features), both models will make the same prediction about Jane. The detailed results for the attacker’s advantage are plotted in Figure 8. We observe that on CORA, the attacker obtains the highest advantage for GRAD, GRAD-I, ZORRO, and ZORRO-S explanations. On CORAML, the highest advantage is obtained with GRAD and ZORRO-S explanations. Usually, the attacker’s advantage is positively correlated with the success rate of corresponding attacks for both CORA and CORAML. On BITCOIN, the attacker’s advantage for all explanation methods is usually high. This is surprising as the success rate for attacks is relatively lower than for other datasets. This might imply that the reconstructed

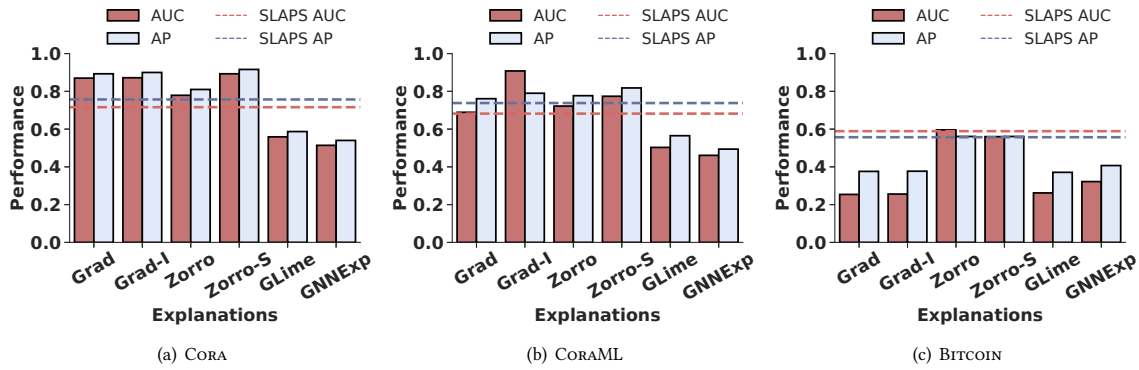


Figure 5: Average AUC and AP of GSE attack on the different datasets. The adopted baseline is SLAPS which use the true node features.

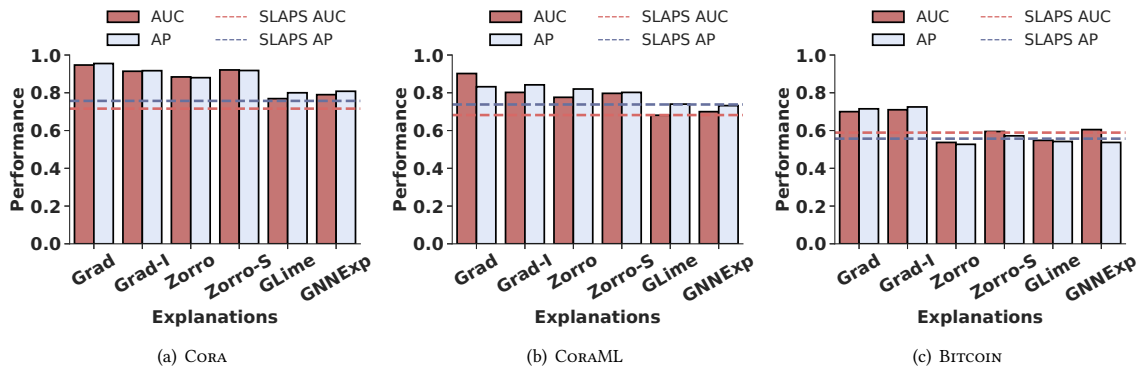


Figure 6: Average AUC and AP of GSEF attack on the different datasets. The adopted baseline is SLAPS which use the true node features.

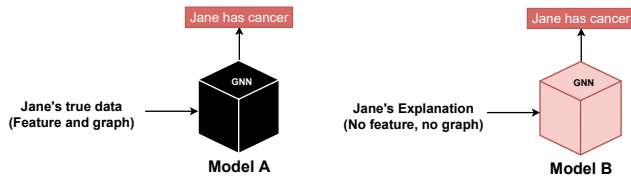


Figure 7: An example use case of the attacker's advantage

graph has the same semantics as the true graph, if not the exact structure.

**4.2.5 Lack of Groundtruth Labels (RQ 4).** We relax the assumption that the attacker has access to groundtruth labels. Instead, she has black-box access to the target model. This is made possible with the popularity of machine learning as a service (MLaaS), where a user can input a query and get the predictions as output. Therefore, the "groundtruth" label is the one obtained from the target model. As representative explanations, we show the performance on GRAD and ZORRO on all datasets.

As shown in Table 5, on the GRAD explanation, we observe a 3% gain in attack performance on GSEF-CONCAT when the attacker has access to the target model on the CORA and CORAML dataset. On the BITCOIN dataset, we observe a decrease of about 3% in AUC and an increase of 6% in AP. The corresponding performance on ZORRO follows the same with no significant change in the performance on CORAML and a 5% decrease in performance on CORA.

The performance of GSEF-MULT and GSEF attack decreases across all datasets, with BITCOIN having the worst performance reduction of up to 29% on GRAD. However, on ZORRO, there is a 2% gain on CORA, a 4% decrease on CORAML, and up to 36% decrease on BITCOIN. We observe performance drop on GSEF and GSE on all datasets across both explanations except for GSE on GRAD, which has up to 5% gain in performance on the CORA and CORAML datasets. It is important to note that the BITCOIN dataset has the least performance across all attacks, even when the true groundtruth label is used. Therefore, the large disparity in performance when the label is generated from the trained black-box model is not surprising.

**Summary.** In the absence of groundtruth labels, we used model predictions. Intuitively, for incorrectly predicted nodes, the model explanations which are based on incorrectly predicted labels should

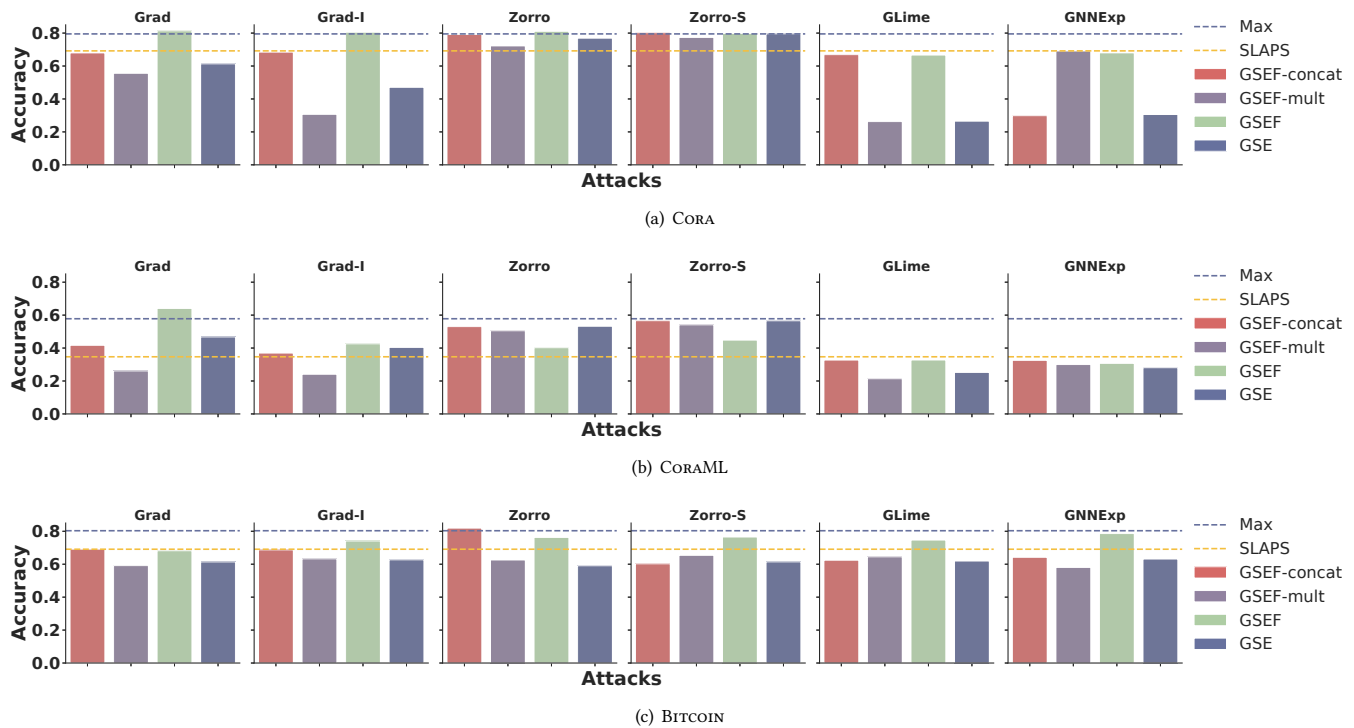


Figure 8: Accuracy of the reconstructed graph on a downstream node classification task by all models on different datasets. The blue line is the original accuracy using the true features and edges, while the yellow line is the SLAPS accuracy.

lead to lower performance in the presence of groundtruth labels. For CORA and CORAML datasets, on GSEF-CONCAT and GSE attacks, having black-box access to the target model performs better than the attacker having access to the groundtruth label. For GSEF-MULT and GSEF attacks, it is better to have access to groundtruth label to achieve the best attack success rate. For the BITCOIN dataset, the groundtruth labels perform better than the black-box access to the target model on all attacks.

### 4.3 Attack Performance When Explanations for Partial Node-set Are Available

In this section, we assume that the attacker has access to explanations of only a subset of the nodes. The attacker is interested in reconstructing the subgraph induced on the partial node set. We perform the experiment on representative datasets CORA and CREDIT. For each of the datasets, we only use 30% of the nodes to obtain explanations. The result is shown in Table 6.

Unsurprisingly, we observe a drop in attack performance scores when explanations for only a node subset are available. However, we observe that explanation augmentation attacks perform best on all explainers and datasets. Recall that when the full explanations are available on GRAD, the augmentation attacks rank second on the CORA dataset. However, with explanations for the partial node set, augmentation attacks rank first with AUROC and AP of 0.91 and 0.92. Moreover, on the CORA and CREDIT datasets, GSEF-CONCAT and GSEF-MULT perform competitively with other attacks. Lastly,

as observed in the experiments, when the full explanations are available, attacks on the GRAD and ZORRO-S explanations are highly successful, while attacks on the GNNEXP are less successful. Nonetheless, we observe AUC and AP  $> 0.76$  on both datasets for GNNEXP.

### 4.4 When Do Augmentation Attacks Work?

To summarize, augmentation attacks perform best on all explanation models except on GRAD and GRAD-I across all datasets. Nonetheless, augmentation attack GSEF on the GRAD and GRAD-I explanations perform comparable to or better than GSE attack, which only uses explanations as input on CITESEER, CORAML, BITCOIN, and CREDIT datasets.

On explainers other than GRAD and GRAD-I, augmentation attacks perform the best. For instance, on the CITESEER dataset, augmentation attack GSEF achieved AUROC and AP of 0.94 and 0.95 respectively on ZORRO-S explanation, which is far better than all baselines and explanation-only attacks. The multi-task learning paradigm that GSEF employs gives it an edge over other augmentation attacks, baselines, and explanation-only attacks.

Overall, the strength of augmentation attacks depends on the choice of the explanation method and the strategy of using the additional feature information. For instance, GSEF-CONCAT and GSEF-MULT augmentation attacks, which combine the additional feature information by concatenation and element-wise multiplication, perform worst on the GRAD, GRAD-I, ZORRO, and ZORRO-S explanations. However, they perform best on surrogate method

**Table 5: Performance comparison of relaxing the availability of groundtruth labels ( $Y$ ) assumption. Here, the attacker has black-box access to the target model ( $M$ ). We perform the experiment on all datasets’ GRAD and ZORRO explanation methods.  $\Delta$  is the percentage difference. A negative value implies that the groundtruth labels are preferred over black-box access.**

Dataset	Attack	$Y_{\text{GRAD}}$		$M_{\text{GRAD}}$		$\Delta_{\text{GRAD}}$		$Y_{\text{ZORRO}}$		$M_{\text{ZORRO}}$		$\Delta_{\text{ZORRO}}$	
		AUC	AP	AUC	AP	AUC	AP	AUC	AP	AUC	AP	AUC	AP
CORA	GSEF-CONCAT	0.734	0.773	0.757	0.784	3.1	1.4	0.823	0.860	0.779	0.810	-5.3	-5.8
	GSEF-MULT	0.678	0.737	0.658	0.700	-2.9	-5.0	0.723	0.756	0.740	0.772	2.4	2.0
	GSEF	0.948	0.953	0.927	0.932	-2.2	-2.2	0.884	0.880	0.871	0.881	-1.5	0.1
	GSE	0.924	0.939	0.946	0.966	2.4	2.9	0.779	0.810	0.814	0.849	4.5	4.8
CORAML	GSEF-CONCAT	0.640	0.705	0.661	0.734	3.3	4.1	0.735	0.786	0.738	0.792	0.4	0.8
	GSEF-MULT	0.666	0.730	0.650	0.693	-2.4	-5.1	0.681	0.697	0.653	0.692	-4.1	-0.7
	GSEF	0.902	0.833	0.808	0.852	-10.4	2.3	0.776	0.820	0.751	0.796	-3.2	-2.9
	GSE	0.699	0.768	0.735	0.795	5.2	3.5	0.722	0.777	0.713	0.759	-1.2	-2.3
BITCOIN	GSEF-CONCAT	0.527	0.515	0.512	0.546	-2.7	6.0	0.575	0.529	0.523	0.517	-9.0	-2.3
	GSEF-MULT	0.264	0.383	0.241	0.367	-8.7	-4.2	0.399	0.449	0.255	0.369	-36.1	-17.8
	GSEF	0.700	0.715	0.500	0.560	-28.6	-21.7	0.537	0.527	0.352	0.438	-34.5	-16.9
	GSE	0.229	0.365	0.210	0.352	-8.3	-3.6	0.596	0.561	0.491	0.503	-17.6	-10.3

**Table 6: Attack performance with explanations over partial node set for CORA and CREDIT datasets. The best performing attack(s) on each explanation method is(are) highlighted in bold, and the second best attack(s) is(are) underlined.**

Exp	Attack	CORA		CREDIT	
		AUC	AP	AUC	AP
GRAD	GSEF-CONCAT	0.683 ± 0.02	0.700 ± 0.03	0.821 ± 0.03	0.851 ± 0.03
	GSEF-MULT	0.674 ± 0.02	0.705 ± 0.02	0.806 ± 0.03	0.839 ± 0.02
	GSEF	<b>0.916 ± 0.03</b>	<b>0.920 ± 0.03</b>	<b>0.842 ± 0.04</b>	<b>0.875 ± 0.04</b>
	GSE	0.865 ± 0.02	0.871 ± 0.02	<u>0.840 ± 0.02</u>	<u>0.870 ± 0.03</u>
	EXPLAINSIM	<u>0.881 ± 0.01</u>	<u>0.913 ± 0.01</u>	0.835 ± 0.01	0.854 ± 0.01
ZORRO-S	GSEF-CONCAT	0.863 ± 0.04	0.896 ± 0.04	<u>0.825 ± 0.03</u>	<u>0.842 ± 0.02</u>
	GSEF-MULT	0.720 ± 0.03	0.757 ± 0.04	0.819 ± 0.02	0.840 ± 0.02
	GSEF	<b>0.901 ± 0.03</b>	<b>0.905 ± 0.03</b>	<b>0.835 ± 0.02</b>	<b>0.847 ± 0.03</b>
	GSE	0.881 ± 0.04	0.893 ± 0.01	0.762 ± 0.03	0.801 ± 0.02
	EXPLAINSIM	<u>0.887 ± 0.02</u>	<u>0.901 ± 0.02</u>	0.701 ± 0.01	0.721 ± 0.01
GNNEXP	GSEF-CONCAT	0.585 ± 0.04	0.600 ± 0.04	<b>0.829 ± 0.03</b>	<b>0.840 ± 0.02</b>
	GSEF-MULT	<u>0.680 ± 0.03</u>	<u>0.699 ± 0.03</u>	<u>0.794 ± 0.02</u>	<u>0.802 ± 0.02</u>
	GSEF	<b>0.762 ± 0.04</b>	<b>0.791 ± 0.06</b>	0.782 ± 0.02	0.799 ± 0.02
	GSE	0.500 ± 0.06	0.520 ± 0.05	0.688 ± 0.04	0.715 ± 0.04
	EXPLAINSIM	0.508 ± 0.02	0.515 ± 0.02	0.615 ± 0.02	0.660 ± 0.02

GLIME and perturbation-based method GNNEXP. Moreover, augmentation attacks show the best performance when explanations only for a partial node set are available.

## 5 DEFENSE

**Explanation Perturbation.** To limit the information leakage by the explanation, we perturb each explanation bit using a randomized response mechanism [22, 40]. Specifically for 0/1 feature (explanation) mask as in ZORRO, we flip each bit of the explanation

with a probability that depends on the privacy budget  $\epsilon$  as follows

$$Pr(\mathcal{E}'_{x_i} = 1) = \begin{cases} \frac{e^\epsilon}{e^\epsilon + 1}, & \text{if } \mathcal{E}_{x_i} = 1, \\ \frac{1}{e^\epsilon + 1}, & \text{if } \mathcal{E}_{x_i} = 0, \end{cases} \quad (7)$$

where  $\mathcal{E}_{x_i}$  and  $\mathcal{E}'_{x_i}$  are true and perturbed  $i^{th}$  bit of explanation  $\mathcal{E}_x$  respectively. Note that our defense mechanism satisfies  $de$ -local differential privacy.

LEMMA 1. For an explanation with  $d$  dimensions, the explanation perturbation defense mechanism in Equation 7 satisfies  $de$ -local differential privacy.

PROOF. Note that for an explanation corresponding to two graph datasets  $D$  and  $D'$  differing in a single edge, the ratio of probabilities of obtaining a certain explanation can be bounded as follows.

$$\frac{\Pr[\mathcal{E}_X(D) = S]}{\Pr[\mathcal{E}_X(D') = S]} = \prod_{i=1}^d \frac{\Pr[\mathcal{E}_{x_i}(D) = S_i]}{\Pr[\mathcal{E}_{x_i}(D') = S_i]} \leq \prod_{i=1}^d \frac{e^{\epsilon}}{e^{\epsilon} + 1} = e^{d\epsilon}.$$

□

**Defense Evaluation.** We evaluate our defense mechanism on EXPLAINSIM attack as it best quantifies the information leakage due to the explanations alone. All other attacks assume the availability of other information, such as features and labels. We use two datasets: CORA and CORAML. As evaluation metrics, we use the AUC score and AP to compute the attack success rate after the defense. Besides, we measure the utility of the perturbed explanation in terms of fidelity, sparsity, and the percentage of 1 bit that is retained from the original explanation (intersection).

**Defense Results.** As shown in Figure 9, the explanation perturbation based on the randomized response mechanism clearly defends against the attack. For instance, at a very high privacy level  $\epsilon = 0.0001$ , which gives  $d\epsilon = 0.14$ , the attack performance drastically dropped to 0.56 in AUC and 0.59 in AP, which is about 36% decrease over the non-private released explanation. As expected, the attack performance decreases significantly with increase in the amount of noise ( $\epsilon$  decreases). In Table 7, we analyze the change in explanation utility due to our perturbation mechanism. We observe that on CORA, with the lowest privacy loss level, there is a drop of 5.61% in the fidelity when the attack is already reduced to a random guess. The entropy of the mask distribution increases. In other words, the explanation sparsity decreases. For ZORRO, this implies that more bits are set to 1 than in the true explanation mask. Even though this decreases explanation utility to some extent, we point out that 74.68% of true explanation is still retained. Moreover, the sparsity is still lower than achieved by GNNExp explanations, even without any perturbations.

While quantitatively, the change in explanation sparsity is acceptable, more application-dependent qualitative studies are required to evaluate the change in the utility of explanations. Nevertheless, we provide a promising first defense for future development and possible improvements. We obtain similar results for CORAML, which are provided in Appendix C.

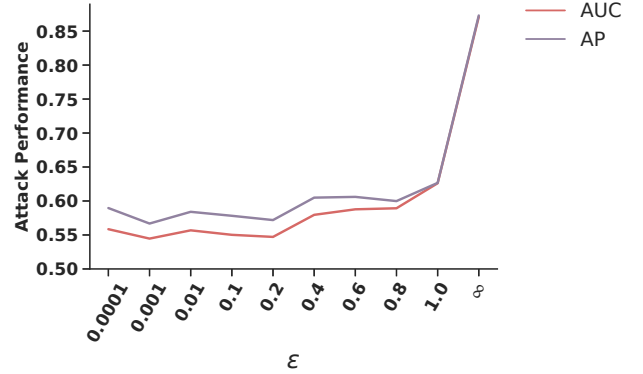
**Defense Variant for Soft Explanation Masks.** Note that Equation 7 only applies to explanations that return binary values. For explanations with continuous values, we can adapt the defense as follows. We keep the original value ( $\mathcal{E}_{x_i}$ ) when the flipped coin lands heads, but when it lands tail, we replace ( $\mathcal{E}_{x_i}$ ) with  $\mathcal{E}_{x'_i}$  where  $\mathcal{E}_{x'_i}$  is a random number drawn from a normal distribution ( $\mathcal{E}_{x'_i} \sim \mathcal{N}(0, 1)$ ).

## 6 CONCLUSION

We initiate the first investigation on the privacy risks of releasing post-hoc explanations of graph neural networks. Concretely, we quantify the information leakage of explanations via our proposed *five* graph reconstruction attacks. The goal of the attacker is to reconstruct the private graph structure information used to train a GNN model. Our results show that even when the explanations

**Table 7: Fidelity, sparsity and percentage of 1 bits in the true explanation that is retained in the perturbed explanation (intersection) after defense for different  $\epsilon$  on the CORA dataset for ZORRO explanation.  $\infty$  implies no privacy.**

$\epsilon$	Fidelity	Sparsity	Intersection
<b>0.0001</b>	0.84	5.91	74.68
<b>0.001</b>	0.84	5.91	74.70
<b>0.01</b>	0.84	5.89	75.03
<b>0.1</b>	0.84	5.80	75.10
<b>0.2</b>	0.83	5.71	75.60
<b>0.4</b>	0.82	5.49	76.45
<b>0.6</b>	0.81	5.25	77.16
<b>0.8</b>	0.81	5.00	78.66
<b>1</b>	0.81	4.73	80.10
$\infty$	0.89	1.83	100



**Figure 9: Privacy budget and corresponding attack performance of EXPLAINSIM for ZORRO explanation on the CORA dataset.  $\infty$  implies that no perturbation is performed.**

alone are available without any additional auxiliary information, the attacker can reconstruct the graph structure with an AUC score of more than 90%. Our explanation-based attacks outperform all baseline methods pointing to the additional privacy risk of releasing explanations. We propose a perturbation-based defense mechanism that reduces the attack to a random guess. The defense leads to a slight decrease in fidelity. At the lowest privacy loss, the perturbed explanation still contains around 75% of the true explanation. While quantitatively, the change in explanation sparsity seems to be acceptable, more application-dependent qualitative studies would be required to evaluate the change in the utility of explanations.

We emphasize that we strongly believe in the transparency of graph machine learning and acknowledge the need to explain trained models. At the same time, our work points out the associated privacy risks which cannot be ignored. We believe that our work would encourage future work on finding solutions to balance the complex trade-off between privacy and transparency.

## ACKNOWLEDGMENTS

This work is partly funded by the Lower Saxony Ministry of Science and Culture under grant number ZN3491 within the Lower Saxony "Vorab" of the Volkswagen Foundation and supported by the Center for Digital Innovations (ZDIN), and the Federal Ministry of Education and Research (BMBF), Germany under the project LeibnizKILabor (grant number 01DD20003). The authors are grateful to the anonymous reviewers for providing insights that further improved our paper.

## REFERENCES

- [1] Chirag Agarwal, Himabindu Lakkaraju, and Marinka Zitnik. 2021. Towards a unified framework for fair and stable graph representation learning. In *Uncertainty in Artificial Intelligence*. PMLR, 2114–2124.
- [2] David Ahmedt-Aristizabal, Mohammad Ali Armin, Simon Denman, Clinton Fookes, and Lars Petersson. 2021. Graph-Based Deep Learning for Medical Diagnosis and Analysis: Past, Present and Future. *Sensors (Basel, Switzerland)* 21, 14 (July 2021). <https://doi.org/10.3390/s21144758>
- [3] Federico Baldassarre and Hossein Azizpour. 2019. Explainability techniques for graph convolutional networks. *arXiv preprint arXiv:1905.13686* (2019).
- [4] Anitesh Barua, Suryanarayanan Ravindran, and Andrew B Whinston. 2007. Enabling information sharing within organizations. *Information Technology and Management* 8, 1 (2007), 31–45.
- [5] Hanjun Dai, Hui Li, Tian Tian, Xin Huang, Lin Wang, Jun Zhu, and Le Song. 2018. Adversarial attack on graph structured data. *arXiv preprint arXiv:1806.02371* (2018).
- [6] Ngan Thi Dong and Megha Khosla. 2020. Revisiting Feature Selection with Data Complexity. In *2020 IEEE 20th International Conference on Bioinformatics and Bioengineering (BIBE)*. 211–216. <https://doi.org/10.1109/BIBE50027.2020.00042>
- [7] Thi Ngan Dong, Stefanie Mucke, and Megha Khosla. 2022. MuCoMiD: A Multi-task graph Convolutional Learning Framework for miRNA-Disease Association Prediction. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* (2022).
- [8] Mengnan Du, Ninghao Liu, Qingquan Song, and Xia Hu. 2018. Towards explanation of dnn-based prediction with guided feature inversion. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1358–1367.
- [9] Vasisht Duddu, Antoine Boutet, and Virat Shejwalkar. 2020. Quantifying Privacy Leakage in Graph Embedding. *arXiv preprint arXiv:2010.00906* (2020).
- [10] Aleksander Fabijan, Helena Holmström Olsson, and Jan Bosch. 2016. The Lack of Sharing of Customer Data in Large Software Organizations: Challenges and Implications. In *Agile Processes, in Software Engineering, and Extreme Programming*, Helen Sharp and Tracy Hall (Eds.). Springer International Publishing, Cham, 39–52.
- [11] Wenqi Fan, Wei Jin, Xiaorui Liu, Han Xu, Xianfeng Tang, Suhang Wang, Qing Li, Jiliang Tang, Jianping Wang, and Charu Aggarwal. 2021. Jointly Attacking Graph Neural Network and its Explanations. *arXiv preprint arXiv:2108.03388* (2021).
- [12] Wenqi Fan, Yao Ma, Qing Li, Yuan He, Eric Zhao, Jiliang Tang, and Dawei Yin. 2019. Graph neural networks for social recommendation. In *The World Wide Web Conference*. 417–426.
- [13] Bahare Fatemi, Layla El Asri, and Seyed Mehran Kazemi. 2021. SLAPS: Self-Supervision Improves Structure Learning for Graph Neural Networks. *Advances in Neural Information Processing Systems* 34 (2021).
- [14] Matthias Fey and Jan Eric Lenssen. 2019. Fast graph representation learning with PyTorch Geometric. *arXiv preprint arXiv:1903.02428* (2019).
- [15] Luca Franceschi, Mathias Niepert, Massimiliano Pontil, and Xiao He. 2019. Learning discrete structures for graph neural networks. In *International conference on machine learning*. PMLR, 1972–1982.
- [16] Thorben Funke, Megha Khosla, Mandeep Rathee, and Avishek Anand. 2022. ZORRO: Valid, Sparse, and Stable Explanations in Graph Neural Networks. *IEEE Transactions on Knowledge and Data Engineering* (2022).
- [17] Bryce Goodman and Seth Flaxman. 2017. European Union regulations on algorithmic decision-making and a “right to explanation”. *AI magazine* 38, 3 (2017), 50–57.
- [18] Danielle Greenstein. 2020. Why you should share data with other in your organization. <https://www.lotame.com/sharing-data-within-organization/>
- [19] William L. Hamilton, Rex Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In *NIPS*.
- [20] Xinlei He, Jinyuan Jia, Michael Backes, Neil Zhenqiang Gong, and Yang Zhang. 2021. Stealing links from graph neural networks. In *30th USENIX Security Symposium (USENIX Security 21)*. 2669–2686.
- [21] Qiang Huang, Makoto Yamada, Yuan Tian, Dinesh Singh, Dawei Yin, and Yi Chang. 2020. Graphlime: Local interpretable model explanations for graph neural networks. *arXiv:2001.06216* (2020).
- [22] Peter Kairouz, Keith Bonawitz, and Daniel Ramage. 2016. Discrete distribution estimation under local privacy. In *International Conference on Machine Learning*. PMLR, 2436–2444.
- [23] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *International Conference on Learning Representations (ICLR)*.
- [24] Srijan Kumar, Francesca Spezzano, VS Subrahmanian, and Christos Faloutsos. 2016. Edge weight prediction in weighted signed networks. In *2016 IEEE 16th International Conference on Data Mining (ICDM)*. IEEE, 221–230.
- [25] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013).
- [26] Iyiola E Olatunji, Thorben Funke, and Megha Khosla. 2021. Releasing Graph Neural Networks with Differential Privacy Guarantees. *arXiv preprint arXiv:2109.08907* (2021).
- [27] Iyiola E Olatunji, Wolfgang Nejdl, and Megha Khosla. 2021. Membership Inference Attack on Graph Neural Networks. In *IEEE International Conference on Trust, Privacy and Security in Intelligent Systems and Applications, TPS-ISA 2021*.
- [28] European Parliament and Council of the European Union. 2016. General Data Protection Regulation. *Official Journal of the European Union* (2016). <https://eur-lex.europa.eu/legal-content/EN/ALL/?uri=CELEX:32016R0679>
- [29] Phillip E Pope, Soheil Kolouri, Mohammad Rostami, Charles E Martin, and Heiko Hoffmann. 2019. Explainability methods for graph convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10772–10781.
- [30] Sina Sajadmanesh and Daniel Gatica-Perez. 2020. Locally Private Graph Neural Networks. *arXiv preprint arXiv:2006.05535* (2020).
- [31] Andrew Selbst and Julia Powles. 2018. “Meaningful Information” and the Right to Explanation. In *Conference on Fairness, Accountability and Transparency*. PMLR, 48–48.
- [32] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*. 618–626.
- [33] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. 2008. Collective classification in network data. *AI magazine* 29, 3 (2008), 93–93.
- [34] Reza Shokri, Martin Strobel, and Yair Zick. 2021. On the privacy risks of model explanations. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Analytics*. 231–241.
- [35] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2013. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv:1312.6034* (2013).
- [36] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *PMLR*.
- [37] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. *International Conference on Learning Representations* (2018).
- [38] Minh Vu and My T Thai. 2020. Pgm-explainer: Probabilistic graphical model explanations for graph neural networks. *Advances in neural information processing systems* 33 (2020), 12225–12235.
- [39] Binghui Wang and Neil Zhenqiang Gong. 2019. Attacking graph-based classification via manipulating the graph structure. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*. 2023–2040.
- [40] Yue Wang, Xintao Wu, and Donghui Hu. 2016. Using Randomized Response for Differential Privacy Preserving Data Collection. In *EDBT/ICDT Workshops*, Vol. 1558. 0090–6778.
- [41] Bang Wu, Xiangwen Yang, Shirui Pan, and Xingliang Yuan. 2020. Model Extraction Attacks on Graph Neural Networks: Taxonomy and Realization. *arXiv preprint arXiv:2010.12751* (2020).
- [42] Fan Wu, Yunhui Long, Ce Zhang, and Bo Li. 2021. LinkTeller: Recovering Private Edges from Graph Neural Networks via Influence Analysis. *arXiv preprint arXiv:2108.06504* (2021).
- [43] Huijun Wu, Chen Wang, Yuriy Tyshetskiy, Andrew Docherty, Kai Lu, and Liming Zhu. 2019. Adversarial examples on graph data: Deep insights into attack and defense. *arXiv preprint arXiv:1903.01610* (2019).
- [44] Makoto Yamada, Wittawat Jitkrittum, Leonid Sigal, Eric P Xing, and Masashi Sugiyama. 2014. High-dimensional feature selection by feature-wise kernelized lasso. *Neural computation* 26, 1 (2014), 185–207.
- [45] Tung-Mou Yang and Terrence A Maxwell. 2011. Information-sharing in public organizations: A literature review of interpersonal, intra-organizational and inter-organizational success factors. *Government information quarterly* 28, 2 (2011), 164–175.
- [46] Rex Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. 2019. GNN Explainer: A tool for post-hoc explanation of graph neural networks. *arXiv:1903.03894* (2019).

- [47] Zaixi Zhang, Jinyuan Jia, Binghui Wang, and Neil Zhenqiang Gong. 2020. Backdoor attacks to graph neural networks. *arXiv preprint arXiv:2006.11165* (2020).
- [48] Zaixi Zhang, Qi Liu, Zhenya Huang, Hao Wang, Chengqiang Lu, Chuanren Liu, and Enhong Chen. 2021. GraphMI: Extracting Private Graph Data from Graph Neural Networks. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, Zhi-Hua Zhou (Ed.). 3749–3755.
- [49] Daniel Zügner, Amir Akbarnejad, and Stephan Günnemann. 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.
- [50] Daniel Zügner and Stephan Günnemann. 2019. Adversarial attacks on graph neural networks via meta learning. *arXiv preprint arXiv:1902.08412* (2019).

## APPENDIX

**Organization.** The Appendix is organized as follows. We first present related works in private graph extraction attacks and other attacks on GNN models. Next, we elaborated on the data sampling used for evaluation in Section B followed by the results for defense on the CORAML dataset in Section C. The performance and the corresponding standard deviation of our main results are provided in Section D. In Section E, we present the results of the additional experiments on the CITESEER, PUBMED, and CREDIT datasets. Finally, we list the hyperparameters used for our attacks and target model in Section F.

### A RELATED WORKS

**Private Graph Extraction Attacks.** Given a black-box access to a GNN model that is trained on a target dataset and the adversary’s background knowledge, He et al. [20] proposed link stealing attacks to infer whether there is a link between a given pair of nodes in the target dataset. Their attacks are specific to the adversary’s background knowledge which range from simply exploiting the node feature similarities to a shadow model-based attack. Wu et al. [42] proposed an edge re-identification attack for vertically partitioned graph learning. Their attack setting is different from ours and is applicable in scenarios where the high-dimensional features and high-order adjacency information are usually heterogeneous and held by different data holders. GraphMI [48] aims at reconstructing the adjacency matrix of the target graph given white-box access to the trained model, node features, and labels.

**Other Inference Attacks and Defenses on GNNs.** Several other attacks, such as membership inference [9, 27] and model extraction attacks [41] have been proposed to quantify privacy leakage in GNNs. In a membership inference attack, the goal of the attacker is to infer whether a node was part of the data used in training the GNN model via black-box access to the trained model. In a model extraction attack, the attacker aims to steal the trained model’s parameter and hyperparameters to duplicate or mimic the functionality of the target model via the predictions returned from querying the model [41]. Recently, several defenses against these attacks have been proposed, which are mainly based on differential privacy. Olatunji et al. [26] proposed a method for releasing GNN models by combining a knowledge-distillation framework with two noise mechanisms, random subsampling and noisy labeling. Their centralized setting approach trains a student model using a public graph and private labels obtained from a teacher model trained exclusively for each query node (personalized teacher models). Their method, by design, defends against membership inference attacks and model extraction attacks since only the student model (which

has limited and perturbed information) is released. Sajadmanesh and Gatica-Perez [30] proposed a locally differentially private GNN model by considering a distributed setting where nodes and labels are private, and the graph structure is known to the central server. Their approach perturbs both the node features and labels to ensure a differential privacy guarantee. However, all attacks and defenses are not applicable to explanations.

**Membership Inference Attack and Explanations.** On euclidean data such as images, Shokri et al. [34] analyzed the privacy risks of feature-based model explanations using membership inference attacks which quantifies the extent to which model predictions and their explanations leak information about the presence of a datapoint in the training set of a model. We emphasize that the goal of Shokri et al. [34] differs from ours in that we focus on reconstructing the entire graph structure from feature-based explanations. Also, their investigations are limited to non-graph data and the corresponding target and explanation models.

**Adversarial Attacks and GNNs.** Another line of research focuses on the vulnerability of GNNs to adversarial attacks [5, 39, 43, 47, 49, 50]. The goal of the attacker is to fool the GNN model into making a wrong prediction by manipulating node features or the structural information of nodes. A recent work [11] used explanation method such as GNNExplainer as a method for detecting adversarial perturbation on graphs. Hence, acting as a tool for inspecting adversarial attacks on GNN models. They further proposed an adversarial attack framework (GEAttack) that exploits the vulnerabilities of explanation methods and the GNN model. This allows the attacker to simultaneously fool the GNN model and misguide the inspection from the explanation method. Our work differs significantly from this work in that first, we aim to reconstruct the graph from the explanations and, secondly, to quantify the privacy leakage of explanations on GNN models.

### B DATA SAMPLING FOR EVALUATION

In this section, we elaborate on the data sampling method used for our evaluation as briefly explained in Section 4.1.3. To generate a test set, we randomly select 10% of all nodes. We use the set of all edges which are incident on at least one of the nodes in our selected set. We then sample an equal number of negative edges (pairs of nodes such that no edge exists between them; one of the nodes in this pair is from the selected set of 10% nodes). In total, we generate 10 random test sets.

We employ 10 random instantiations of the attack models. Corresponding to each instantiation of the attack model, we test the model with one of our random test sets. In total, we have 10 observations of the final scores. We evaluated all methods on the same test sets. As we wanted to keep the number of observations the same for all compared methods, including the non-parameterized methods based on feature and explanation similarity (note that they do not use a neural model, and there is no initialization required in their cases), we tested each attack model (corresponding to each initialization) on one test set. We computed the average precision (AP) and AUROC result on the balanced pairs. For the partial explanation experiment, we randomly sample 30% of the nodes only once to generate a fixed subgraph. Evaluation is performed as in



the previous case. Here, the 10 test sets are generated from this subgraph.

We remark that using AUC and AP scores allows us to evaluate all methods across all ranges of cutoff decision thresholds. Our evaluation is domain and dataset-independent. Nevertheless, in real scenarios, an attacker might need to choose a single decision threshold for which domain-specific knowledge or an auxiliary validation dataset might be required.

### C RESULTS FOR DEFENSE ON CORAML

The attack performance for different values of  $\epsilon$  is plotted in Figure 10. For the lowest privacy budget, we observe that the attack is reduced to a random guess (with an AUC score close to 0.55). The variation in explanation utility and intersection with true explanation is shown in Table 8. Here also, the perturbed explanation is able to retain around 75% of the 1 bit of the true explanation. While there is a drop in fidelity and the explanation becomes denser, we note that the perturbed explanation still shows higher fidelity and sparsity than other explanation methods (c.f. Table 4).

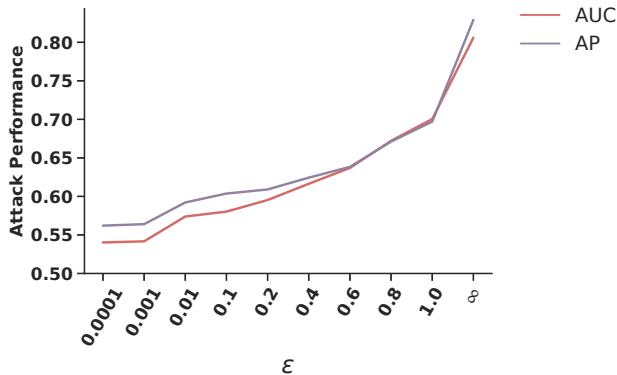


Figure 10: Privacy budget and corresponding attack performance of EXPLAINSIM for ZORRO explanation on the CORAML dataset.  $\infty$  implies no privacy.

### D DETAILED RESULTS

In Table 9, we present the mean and standard deviation of the results in Table 3. Since our results have a low standard deviation across the 10 runs of the experiment, our method is stable.

### E ADDITIONAL EXPERIMENTS

We perform additional experiments on 3 additional datasets to verify our observations and the effectiveness of our attacks. We start by describing the additional datasets. The performance scores are provided in Table 12, and results are discussed in the main paper.

#### E.1 Additional Datasets

We used CITESEER and PUBMED, which are conventional datasets for node classification tasks and the CREDIT dataset. The details about the datasets are in Table 10. The fidelity and sparsity for the corresponding explanations are provided in Table 11.

Table 8: Fidelity, sparsity and percentage of 1 bits in the true explanation that is retained in the perturbed explanation (intersection) after defense for different  $\epsilon$  on the CORAML dataset for ZORRO explanation.  $\infty$  implies no privacy.

$\epsilon$	Fidelity	Sparsity	Intersection
<b>0.0001</b>	0.86	4.53	74.96
<b>0.001</b>	0.86	4.53	74.98
<b>0.01</b>	0.86	4.53	75.12
<b>0.1</b>	0.87	4.46	75.08
<b>0.2</b>	0.87	4.39	75.20
<b>0.4</b>	0.87	4.24	75.70
<b>0.6</b>	0.87	4.08	77.07
<b>0.8</b>	0.88	3.95	78.70
<b>1</b>	0.89	3.81	80.75
$\infty$	0.96	3.33	100

**PUBMED and CITESEER.** Both PUBMED and CITESEER [33] are citation datasets, where each node is a research article, and there is an edge between two articles if one cites the other. In CITESEER, features are binary values like CORA, while in PUBMED, the features are represented by the TF/IDF weighted word vector of the unique words in the dictionary.

**CREDIT.** The Credit defaulter graph [1] is a financial network where an edge exists between individuals (nodes) if there is a similarity in their spending and payment pattern. The task is to determine whether an individual will default on their credit card payment. We used a subset of this credit dataset which has 3000 nodes.

#### E.2 Ensuring the Stability of Results

To ensure the stability of the results, we perform an additional experiment by slightly modifying the evaluation procedure in Section B (hereafter referred to as *10-runs*). Instead of testing the attacks on one of our random test sets for each instantiation which resulted in 10 observations of the final scores, we modify the evaluation as follows. For each instantiation, we evaluate the attacks on *all* 10 held out balanced test sets. This results in 100 observations of the final score (recall that we have 10 instantiations). We report the mean and the standard deviation. This modification which we call *100-runs*, further mitigates any sampling biases from the test sets and any randomness. Similar to 10-runs, we evaluated all methods on the same test sets. The EXPLAINSIM attack is excluded since there is no randomness involved, resulting in the same results as in 10-runs.

We observe that in most cases, the difference between all previous experiments, which were evaluated using the 10-runs procedure and the modified 100-runs is 0. This shows that our experiments are stable and the results are reliable. Although we observe some minor variations in some results, however, such variations are relatively small and do not invalidate our observations or experiments. For example, on the PUBMED dataset, GSE attack on the GRAD explanation for the 10-runs resulted in an AUC of 0.770 while that of 100-runs is 0.767. The results of the 100-runs are in Table 13.

**Table 9: Results with standard deviation of Table 3. The best performing attack(s) on each explanation method is(are) highlighted in bold, and the second best attack(s) is(are) underlined.**

Exp	Attack	CORA		CORAML		BITCOIN	
		AUC	AP	AUC	AP	AUC	AP
Baseline	FEATURESIM	0.799±0.04	0.827±0.04	0.706±0.08	0.753±0.07	0.535±0.03	0.478±0.02
	Lsa [20]	0.795±0.03	0.810±0.02	0.725±0.04	0.760±0.01	0.532±0.05	0.500±0.06
	GRAPHMI [48]	0.856±0.01	0.830±0.01	0.808±0.02	0.814±0.03	0.585±0.05	0.518±0.06
	SLAPS [13]	0.736±0.05	0.776±0.05	0.649±0.06	0.702±0.07	0.597±0.09	0.577±0.07
GRAD	GSEF-CONCAT	0.734±0.05	0.773±0.04	0.640±0.05	0.705±0.04	0.527±0.04	0.515±0.03
	GSEF-MULT	0.678±0.04	0.737±0.03	0.666±0.06	0.730±0.05	0.264±0.07	0.383±0.04
	GSEF	<u>0.948 ± 0.02</u>	<u>0.953 ± 0.01</u>	<b>0.902 ± 0.08</b>	<u>0.833 ± 0.07</u>	<b>0.700 ± 0.05</b>	<b>0.715 ± 0.04</b>
	GSE	0.924±0.03	0.939±0.02	0.699±0.07	0.768±0.05	0.229±0.03	0.365±0.02
	EXPLAINSIM	<b>0.984 ± 0.01</b>	<b>0.978 ± 0.01</b>	<u>0.890 ± 0.04</u>	<b>0.891 ± 0.04</b>	<u>0.681 ± 0.03</u>	<u>0.644 ± 0.03</u>
GRAD-I	GSEF-CONCAT	0.734±0.06	0.775±0.04	0.674±0.05	0.724±0.04	0.525±0.09	0.527±0.05
	GSEF-MULT	0.691±0.02	0.742±0.02	0.717±0.05	0.756±0.06	0.252±0.03	0.380±0.02
	GSEF	<u>0.949 ± 0.02</u>	<u>0.950 ± 0.02</u>	<u>0.787 ± 0.08</u>	<u>0.832 ± 0.07</u>	<b>0.709 ± 0.04</b>	<b>0.723 ± 0.03</b>
	GSE	0.903±0.04	0.923±0.04	0.717±0.08	0.781±0.06	0.256±0.03	0.380±0.02
	EXPLAINSIM	<b>0.984 ± 0.01</b>	<b>0.979 ± 0.01</b>	<b>0.903 ± 0.04</b>	<b>0.899 ± 0.04</b>	<u>0.681 ± 0.03</u>	<u>0.644 ± 0.03</u>
ZORRO	GSEF-CONCAT	0.823±0.04	0.860±0.05	0.735±0.02	0.786±0.01	<u>0.575 ± 0.03</u>	0.529±0.05
	GSEF-MULT	0.723±0.03	0.756±0.03	0.681±0.02	0.697±0.04	0.399±0.07	0.449±0.05
	GSEF	<b>0.884 ± 0.03</b>	<b>0.880 ± 0.04</b>	<u>0.776 ± 0.03</u>	<u>0.820 ± 0.02</u>	0.537±0.05	<u>0.527 ± 0.04</u>
	GSE	0.779±0.04	0.810±0.01	0.722±0.02	0.777±0.02	<b>0.596 ± 0.03</b>	<b>0.561 ± 0.03</b>
	EXPLAINSIM	<u>0.871 ± 0.02</u>	<u>0.873 ± 0.02</u>	<b>0.806 ± 0.02</b>	<b>0.829 ± 0.03</b>	0.427±0.06	0.485±0.05
ZORRO-S	GSEF-CONCAT	0.907±0.03	0.922±0.02	<u>0.747 ± 0.06</u>	<u>0.791 ± 0.05</u>	<b>0.601 ± 0.06</b>	<b>0.590 ± 0.05</b>
	GSEF-MULT	0.794±0.06	0.815±0.06	0.712±0.06	0.740±0.06	0.490±0.08	0.491±0.05
	GSEF	<b>0.918 ± 0.02</b>	<u>0.923 ± 0.02</u>	<b>0.776 ± 0.06</b>	0.819±0.05	<u>0.598 ± 0.03</u>	<u>0.565 ± 0.03</u>
	GSE	0.893±0.04	0.915±0.02	0.742±0.06	<b>0.784 ± 0.05</b>	0.571±0.03	0.564±0.04
	EXPLAINSIM	<u>0.908 ± 0.03</u>	<b>0.934 ± 0.02</b>	0.732±0.05	0.787±0.03	0.484±0.04	0.496±0.03
GLIME	GSEF-CONCAT	<u>0.643 ± 0.05</u>	<u>0.710 ± 0.04</u>	<u>0.610 ± 0.05</u>	<u>0.652 ± 0.04</u>	<u>0.473 ± 0.07</u>	<u>0.492 ± 0.05</u>
	GSEF-MULT	0.516±0.06	0.522±0.04	0.517±0.05	0.528±0.04	0.264±0.03	0.371±0.01
	GSEF	<b>0.730 ± 0.05</b>	<b>0.774 ± 0.03</b>	<b>0.681 ± 0.05</b>	<b>0.740 ± 0.05</b>	<b>0.542 ± 0.05</b>	<b>0.525 ± 0.03</b>
	GSE	0.558±0.06	0.571±0.05	0.540±0.06	0.555±0.05	0.236±0.04	0.361±0.02
	EXPLAINSIM	0.505±0.04	0.524±0.04	0.520±0.04	0.523±0.04	0.504±0.05	0.512±0.03
GNNEXP	GSEF-CONCAT	0.614±0.04	0.650±0.04	0.653±0.05	0.705±0.04	0.467±0.11	0.489±0.06
	GSEF-MULT	<u>0.724 ± 0.05</u>	<u>0.760 ± 0.05</u>	<u>0.637 ± 0.05</u>	<u>0.692 ± 0.05</u>	0.390±0.10	0.454±0.05
	GSEF	<b>0.762 ± 0.04</b>	<b>0.796 ± 0.03</b>	<b>0.700 ± 0.07</b>	<b>0.796 ± 0.06</b>	<b>0.590 ± 0.04</b>	<b>0.563 ± 0.03</b>
	GSE	0.517±0.04	0.552±0.03	0.490±0.05	0.508±0.04	0.386±0.11	0.451±0.07
	EXPLAINSIM	0.537±0.05	0.541±0.04	0.484±0.05	0.508±0.04	<u>0.551 ± 0.04</u>	<u>0.545 ± 0.03</u>

## F HYPERPARAMETER SETTINGS

In Table 14, we provide the hyperparameters used in training the different modules of our attack for reproducibility. We used hyperparameters from [13], which are already fine-tuned for these datasets.

## G LIMITATIONS AND FUTURE DIRECTION

This section highlights some limitations and future directions for our work.

**Partial Explanations.** In the current work, the 30% nodes we sampled to generate the fixed subgraph for the partial explanation

experiment were only sampled once. Such a sampling approach might be subject to sampling bias. Hence, sampling the subgraph multiple times may be desirable to avoid such bias.

**More Work on Attacks and Defenses.** A promising direction will be to launch more attacks and defenses on other explanation outputs, such as nodes or subgraph-level explanations. These are natural derivatives of graph-based explanation methods, and it is worth investigating to quantify the risk involved in releasing explanations. More defenses that optimize the usefulness of the explanations (fidelity and sparsity) are desirable.

**Fidelity-sparsity Analysis.** In this work, we analyzed the effect of the usefulness of explanations as measured by their fidelity and

**Table 10: Dataset statistics for the additional datasets.  $|V|$  and  $|E|$  denotes the number of nodes and edges respectively,  $C$ ,  $X_d$ , and  $\text{deg}$  denotes the number of classes, size of feature dimension and the average degree of the corresponding graph dataset.**

	CITeseer	PUBMED	CREDIT
$ V $	3327	19717	3000
$ E $	9228	88651	28854
$X_d$	3703	500	13
$C$	6	3	2
$\text{deg}$	2.77	4.50	9.62

**Table 11: RDT-Fidelity and sparsity (entropy) of different explanation methods on the additional datasets. For fidelity, the higher the better. For sparsity, the lower the better.**

$Exp$	CITeseer		PUBMED		CREDIT	
	Fidelity	Sparsity	Fidelity	Sparsity	Fidelity	Sparsity
<b>GRAD</b>	0.24	4.15	0.44	4.38	0.55	0.94
<b>GRAD-I</b>	0.19	4.16	0.37	4.39	0.59	0.95
<b>ZORRO-S</b>	0.90	2.62	0.96	2.17	0.94	1.23
<b>GLIME</b>	0.18	0.65	0.38	1.59	0.60	0.49
<b>GNNEXP</b>	0.65	8.21	0.79	6.21	0.65	2.56

sparsity. More work can be done in this direction. In general, our fidelity and sparsity metrics can be complemented with other domain and dataset-specific metrics of explanation utility.

**Table 12: Attack performance and baselines for additional datasets (CITeseer, PubMed, and CREDIT). The best performing attack(s) on each explanation method is(are) highlighted in bold, and the second best attack(s) is(are) underlined.**

Exp	Attack	CITeseer		PubMed		CREDIT	
		AUC	AP	AUC	AP	AUC	AP
Baseline	FEATURESIM	0.890 ± 0.04	0.909 ± 0.03	0.888 ± 0.01	0.889 ± 0.01	0.720 ± 0.01	0.753 ± 0.01
	LSA [20]	0.859 ± 0.03	0.893 ± 0.02	0.783 ± 0.05	0.801 ± 0.04	0.655 ± 0.03	0.671 ± 0.04
	GRAPHMI [48]	0.789 ± 0.02	0.746 ± 0.02	0.795 ± 0.01	0.758 ± 0.02	0.510 ± 0.04	0.569 ± 0.05
	SLAPS [13]	0.834 ± 0.04	0.875 ± 0.03	0.579 ± 0.02	0.626 ± 0.02	0.723 ± 0.02	0.767 ± 0.02
GRAD	GSEF-CONCAT	0.811 ± 0.03	0.860 ± 0.03	0.672 ± 0.02	0.732 ± 0.01	0.825 ± 0.02	0.868 ± 0.02
	GSEF-MULT	0.748 ± 0.05	0.808 ± 0.04	0.579 ± 0.03	0.627 ± 0.03	0.820 ± 0.02	0.861 ± 0.01
	GSEF	<u>0.969 ± 0.02</u>	<u>0.971 ± 0.02</u>	<u>0.788 ± 0.02</u>	<u>0.856 ± 0.02</u>	<b>0.881 ± 0.02</b>	<b>0.918 ± 0.02</b>
	GSE	0.945 ± 0.03	0.955 ± 0.02	0.770 ± 0.03	0.846 ± 0.02	0.861 ± 0.02	0.899 ± 0.01
	EXPLAINSIM	<b>0.991 ± 0.01</b>	<b>0.985 ± 0.01</b>	<b>0.987 ± 0.01</b>	<b>0.981 ± 0.01</b>	<u>0.861 ± 0.01</u>	<u>0.867 ± 0.01</u>
GRAD-I	GSEF-CONCAT	0.810 ± 0.04	0.857 ± 0.04	0.583 ± 0.02	0.632 ± 0.02	0.826 ± 0.02	0.867 ± 0.02
	GSEF-MULT	0.745 ± 0.05	0.804 ± 0.04	0.574 ± 0.01	0.625 ± 0.01	0.813 ± 0.02	0.863 ± 0.01
	GSEF	<u>0.970 ± 0.02</u>	<u>0.974 ± 0.02</u>	<u>0.782 ± 0.02</u>	<u>0.853 ± 0.02</u>	<b>0.881 ± 0.02</b>	<b>0.918 ± 0.01</b>
	GSE	0.947 ± 0.02	0.960 ± 0.02	0.740 ± 0.02	0.836 ± 0.01	0.861 ± 0.02	0.893 ± 0.02
	EXPLAINSIM	<b>0.992 ± 0.001</b>	<b>0.986 ± 0.001</b>	<b>0.988 ± 0.001</b>	<b>0.981 ± 0.001</b>	<u>0.864 ± 0.001</u>	<u>0.868 ± 0.001</u>
ZORRO-S	GSEF-CONCAT	0.932 ± 0.02	0.944 ± 0.02	0.663 ± 0.02	0.729 ± 0.01	<u>0.825 ± 0.03</u>	<u>0.868 ± 0.02</u>
	GSEF-MULT	0.828 ± 0.06	0.855 ± 0.06	0.554 ± 0.03	0.601 ± 0.02	0.828 ± 0.03	0.868 ± 0.02
	GSEF	<b>0.953 ± 0.02</b>	<b>0.960 ± 0.02</b>	<b>0.901 ± 0.02</b>	<b>0.914 ± 0.02</b>	<b>0.822 ± 0.02</b>	<b>0.869 ± 0.02</b>
	GSE	0.932 ± 0.03	0.947 ± 0.03	0.753 ± 0.02	0.772 ± 0.02	0.785 ± 0.01	0.840 ± 0.01
	EXPLAINSIM	<u>0.950 ± 0.01</u>	<u>0.957 ± 0.01</u>	<u>0.874 ± 0.02</u>	<u>0.904 ± 0.01</u>	0.734 ± 0.02	0.768 ± 0.02
GLIME	GSEF-CONCAT	<u>0.757 ± 0.05</u>	<u>0.801 ± 0.05</u>	0.614 ± 0.03	0.641 ± 0.02	<u>0.825 ± 0.03</u>	<u>0.867 ± 0.02</u>
	GSEF-MULT	0.567 ± 0.03	0.637 ± 0.04	0.507 ± 0.10	0.509 ± 0.08	0.791 ± 0.01	0.855 ± 0.01
	GSEF	<b>0.862 ± 0.05</b>	<b>0.894 ± 0.04</b>	<b>0.697 ± 0.02</b>	<b>0.712 ± 0.02</b>	<b>0.845 ± 0.02</b>	<b>0.890 ± 0.01</b>
	GSE	0.571 ± 0.04	0.618 ± 0.04	0.602 ± 0.03	0.591 ± 0.03	0.798 ± 0.03	0.843 ± 0.03
	EXPLAINSIM	0.703 ± 0.06	0.705 ± 0.05	<u>0.661 ± 0.03</u>	<u>0.724 ± 0.02</u>	0.765 ± 0.03	0.810 ± 0.02
GNNEXP	GSEF-CONCAT	0.639 ± 0.04	0.713 ± 0.04	0.552 ± 0.03	0.590 ± 0.03	<b>0.826 ± 0.03</b>	<b>0.869 ± 0.02</b>
	GSEF-MULT	<u>0.812 ± 0.05</u>	<u>0.862 ± 0.04</u>	<u>0.570 ± 0.01</u>	<u>0.618 ± 0.01</u>	<u>0.821 ± 0.02</u>	<u>0.865 ± 0.02</u>
	GSEF	<b>0.847 ± 0.05</b>	<b>0.881 ± 0.05</b>	<b>0.612 ± 0.02</b>	<b>0.648 ± 0.02</b>	0.817 ± 0.02	0.869 ± 0.02
	GSE	0.594 ± 0.07	0.640 ± 0.06	0.507 ± 0.02	0.512 ± 0.02	0.760 ± 0.04	0.804 ± 0.05
	EXPLAINSIM	0.618 ± 0.05	0.654 ± 0.05	0.495 ± 0.01	0.499 ± 0.01	0.787 ± 0.02	0.845 ± 0.02

**Table 13: 100-runs results with standard deviation. The best performing attack(s) on each explanation method is(are) highlighted in bold, and the second best attack(s) is(are) underlined.**

Exp	Attack	CORA		CORAML		BITCOIN		CITESEER		PUBMED		CREDIT	
		AUC	AP	AUC	AP	AUC	AP	AUC	AP	AUC	AP	AUC	AP
GRAD	GSEF-CONCAT	0.740 ± 0.05	0.780 ± 0.04	0.639 ± 0.07	0.703 ± 0.05	<u>0.474 ± 0.10</u>	<u>0.503 ± 0.07</u>	0.814 ± 0.04	0.862 ± 0.03	0.666 ± 0.03	0.728 ± 0.02	0.825 ± 0.03	0.868 ± 0.02
	GSEF-MULT	0.682 ± 0.05	0.742 ± 0.04	0.667 ± 0.07	0.727 ± 0.06	0.283 ± 0.03	0.392 ± 0.02	0.750 ± 0.05	0.806 ± 0.04	0.576 ± 0.02	0.627 ± 0.02	0.820 ± 0.02	0.861 ± 0.01
	GSEF	<b>0.949 ± 0.02</b>	<b>0.955 ± 0.01</b>	<b>0.801 ± 0.07</b>	<b>0.836 ± 0.06</b>	<b>0.563 ± 0.15</b>	<b>0.552 ± 0.14</b>	<b>0.969 ± 0.02</b>	<b>0.972 ± 0.02</b>	<b>0.787 ± 0.02</b>	<b>0.855 ± 0.02</b>	<b>0.878 ± 0.02</b>	<b>0.916 ± 0.01</b>
	GSE	<u>0.919 ± 0.04</u>	<u>0.935 ± 0.02</u>	<u>0.697 ± 0.07</u>	<u>0.767 ± 0.05</u>	0.265 ± 0.04	0.380 ± 0.02	<u>0.946 ± 0.03</u>	<u>0.955 ± 0.02</u>	<u>0.767 ± 0.03</u>	<u>0.845 ± 0.02</u>	<u>0.868 ± 0.02</u>	<u>0.902 ± 0.01</u>
GRAD-I	GSEF-CONCAT	0.727 ± 0.05	0.767 ± 0.04	0.666 ± 0.06	0.719 ± 0.05	<u>0.523 ± 0.08</u>	<u>0.528 ± 0.05</u>	0.813 ± 0.04	0.860 ± 0.04	0.578 ± 0.02	0.628 ± 0.02	0.824 ± 0.03	0.867 ± 0.02
	GSEF-MULT	0.689 ± 0.02	0.741 ± 0.02	0.699 ± 0.06	0.747 ± 0.06	0.289 ± 0.06	0.403 ± 0.04	0.738 ± 0.04	0.797 ± 0.04	0.573 ± 0.02	0.625 ± 0.02	0.810 ± 0.02	0.860 ± 0.01
	GSEF	<b>0.949 ± 0.02</b>	<b>0.953 ± 0.02</b>	<b>0.792 ± 0.07</b>	<b>0.830 ± 0.06</b>	<b>0.557 ± 0.05</b>	<b>0.547 ± 0.03</b>	<b>0.969 ± 0.02</b>	<b>0.973 ± 0.02</b>	<b>0.784 ± 0.02</b>	<b>0.854 ± 0.02</b>	<b>0.879 ± 0.02</b>	<b>0.918 ± 0.01</b>
	GSE	<u>0.900 ± 0.04</u>	<u>0.921 ± 0.03</u>	<u>0.706 ± 0.06</u>	<u>0.772 ± 0.05</u>	0.242 ± 0.05	0.370 ± 0.02	<u>0.945 ± 0.03</u>	<u>0.958 ± 0.02</u>	<u>0.734 ± 0.03</u>	<u>0.834 ± 0.02</u>	<u>0.858 ± 0.02</u>	<u>0.892 ± 0.02</u>
ZORRO-S	GSEF-CONCAT	<b>0.909 ± 0.03</b>	<b>0.924 ± 0.02</b>	<u>0.750 ± 0.06</u>	<u>0.790 ± 0.05</u>	<u>0.573 ± 0.04</u>	<b>0.569 ± 0.03</b>	<u>0.935 ± 0.02</u>	0.947 ± 0.02	0.663 ± 0.02	<u>0.730 ± 0.02</u>	<u>0.826 ± 0.03</u>	<u>0.868 ± 0.02</u>
	GSEF-MULT	0.796 ± 0.05	0.821 ± 0.05	0.717 ± 0.06	0.744 ± 0.06	0.512 ± 0.09	0.512 ± 0.06	0.826 ± 0.06	0.853 ± 0.05	0.559 ± 0.03	0.606 ± 0.02	<b>0.829 ± 0.03</b>	<b>0.869 ± 0.02</b>
	GSEF	<u>0.902 ± 0.02</u>	<u>0.922 ± 0.02</u>	<u>0.772 ± 0.06</u>	<b>0.811 ± 0.05</b>	<b>0.602 ± 0.03</b>	<u>0.568 ± 0.03</u>	<b>0.953 ± 0.03</b>	<b>0.961 ± 0.03</b>	<b>0.686 ± 0.02</b>	<b>0.754 ± 0.02</b>	0.818 ± 0.03	0.867 ± 0.02
	GSE	0.889 ± 0.04	0.915 ± 0.02	0.745 ± 0.07	0.787 ± 0.06	0.547 ± 0.04	0.546 ± 0.04	0.933 ± 0.03	<u>0.948 ± 0.03</u>	<u>0.664 ± 0.02</u>	0.730 ± 0.01	0.785 ± 0.02	0.841 ± 0.01
GLIME	GSEF-CONCAT	<u>0.655 ± 0.06</u>	<u>0.716 ± 0.05</u>	<u>0.614 ± 0.06</u>	<u>0.655 ± 0.05</u>	<b>0.537 ± 0.07</b>	<b>0.520 ± 0.05</b>	<u>0.762 ± 0.05</u>	<u>0.802 ± 0.05</u>	<u>0.614 ± 0.04</u>	<u>0.643 ± 0.03</u>	<u>0.827 ± 0.03</u>	<u>0.869 ± 0.02</u>
	GSEF-MULT	0.518 ± 0.05	0.540 ± 0.05	0.559 ± 0.06	0.564 ± 0.05	0.232 ± 0.03	0.358 ± 0.01	0.566 ± 0.04	0.635 ± 0.05	0.529 ± 0.07	0.513 ± 0.06	0.791 ± 0.01	0.855 ± 0.01
	GSEF	<b>0.743 ± 0.04</b>	<b>0.784 ± 0.03</b>	<b>0.649 ± 0.06</b>	<b>0.697 ± 0.05</b>	<u>0.516 ± 0.05</u>	<u>0.512 ± 0.03</u>	<b>0.853 ± 0.05</b>	<b>0.887 ± 0.05</b>	<b>0.680 ± 0.03</b>	<b>0.698 ± 0.03</b>	<b>0.844 ± 0.03</b>	<b>0.889 ± 0.02</b>
	GSE	0.534 ± 0.05	0.565 ± 0.05	0.523 ± 0.05	0.542 ± 0.05	0.274 ± 0.05	0.373 ± 0.02	0.581 ± 0.05	0.620 ± 0.04	0.578 ± 0.05	0.569 ± 0.05	0.790 ± 0.03	0.830 ± 0.04
CNNEXP	GSEF-CONCAT	0.628 ± 0.04	0.648 ± 0.05	0.642 ± 0.06	<u>0.695 ± 0.05</u>	<u>0.476 ± 0.09</u>	<u>0.497 ± 0.05</u>	0.628 ± 0.05	0.704 ± 0.04	0.555 ± 0.02	0.596 ± 0.02	<b>0.825 ± 0.03</b>	<u>0.868 ± 0.02</u>
	GSEF-MULT	<u>0.730 ± 0.05</u>	<u>0.768 ± 0.05</u>	<b>0.645 ± 0.06</b>	<b>0.699 ± 0.05</b>	0.393 ± 0.11	0.453 ± 0.05	<u>0.805 ± 0.05</u>	<u>0.854 ± 0.04</u>	<u>0.564 ± 0.02</u>	<u>0.613 ± 0.02</u>	0.819 ± 0.02	0.863 ± 0.02
	GSEF	<b>0.762 ± 0.04</b>	<b>0.794 ± 0.03</b>	<u>0.644 ± 0.07</u>	0.684 ± 0.06	<b>0.601 ± 0.04</b>	<b>0.574 ± 0.03</b>	<b>0.856 ± 0.05</b>	<b>0.887 ± 0.04</b>	<b>0.616 ± 0.02</b>	<b>0.644 ± 0.02</b>	<u>0.818 ± 0.02</u>	<b>0.868 ± 0.02</b>
	GSE	0.515 ± 0.05	0.542 ± 0.04	0.500 ± 0.06	0.518 ± 0.05	0.311 ± 0.08	0.407 ± 0.04	0.593 ± 0.05	0.647 ± 0.05	0.499 ± 0.02	0.507 ± 0.02	0.766 ± 0.03	0.816 ± 0.03

**Table 14: Hyperparameters of the different modules of the attack and target model**

Param/Module	Self-supervision	Classification	Target Model
num_layers	2	2	2
learning rate	0.01	0.001	0.01
epochs	2000	200	200
hidden_size	512	32	32
dropout	0.5	0.5	0.5
noisy_mask_ratio	20	-	-
weight_decay	-	-	5e-4