# IMPROVING DEFENSE MECHANISMS FOR SUBGRAPH STRUCTURE MEMBERSHIP INFERENCE ATTACKS

Anonymous authors

Paper under double-blind review

### ABSTRACT

Graph neural networks (GNNs) are of significant importance in diverse real-world applications since they leverage powerful graph learning techniques to solve problems pertaining to social network mining and medical data analysis. Despite their practical relevance, GNNs remain vulnerable to adversarial attacks such as membership inference attacks (MIAs) which pose privacy risks by revealing whether specific data records were part of the training set of the model. While most existing research has focused on designing defense mechanisms for known nodelevel MIAs, and in particular, for determining if a certain node was used during training, only limited attention has been paid to subgraph-structure MIA (SMIA) problems. SMIA methods seek to infer whether a set of nodes forms a particular target structure of interest (such as a graph motif, e.g., clique or multi-hop path) in the training graph. The main contributions of our work are three-fold. The first is a novel robust defense mechanism for GNNs against SMIA attacks. It combines an alternating train-test schedule with a flattening strategy to mitigate the attacks. The second contribution is a new end-to-end SMIA attack model that outperforms existing attacks by using multiset functions to generate learnable embeddings for collections of nodes. Extensive simulations reveal that the new attack model outperforms prior state-of-the-art attack models on GNNs by 12.31%across four datasets when no defense mechanism is present. With the new defense mechanism, one can achieve an average decrease of 14.30% in the attack AU-ROC and an 10.05% improvement in target model utility compared to classical defenses, even when using the improved attack scheme. The third contribution is a study that shows that our defense mechanism extends to node-level MIAs as well, offering similar improvements in attack resistance and utility.

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

### 1 INTRODUCTION

Graph neural networks (GNNs) have emerged as indispensable learning modalities for diverse real-037 world problems, ranging from social network mining and recommendation system design to biological data analysis (Wu et al., 2021b; Zhang et al., 2024). For example, GNNs have been used to improve personalized search and recommendations for customers on e-commerce platforms (e.g., 040 AliGraph at Alibaba (Zhu et al., 2019) and GIANT at Amazon (Chien et al., 2021a)) and to per-041 form inference and prediction on social networks (e.g., PinnerSage at Pinterest (Pal et al., 2020) and 042 LiGNN at LinkedIn (Borisyuk et al., 2024)). GNNs, unlike standard neural networks, make full 043 use of both the discrete graph topology and node and edge features via special embedding methods 044 based on graph convolutions or random walks (Wu et al., 2020). For example, graph convolutions allow GNNs to generate informative embeddings that are well-suited for different downstream tasks.

Despite the successful deployment of GNNs, several weaknesses of GNN models have been pointed out in the literature. One major concern pertains to data privacy, which is becoming an issue of ever increasing importance. GNNs have exhibited privacy vulnerabilities to various attacks (Sun et al., 2023) such as membership inference attacks (MIAs) (Shokri et al., 2017; Hu et al., 2022; Olatunji et al., 2021), whose aim is to determine if a certain sample is used in model training; attribute inference attacks (AIAs), whose focus is on inferring statistical information about the data, such as the number of nodes and edges (Gong & Liu, 2018) and others. Among all existing attacks, the most commonly observed and studied attacks are MIAs. As already pointed out, MIAs have the goal to reveal whether a given record is part of the training dataset used to build a specific target model,

054 and are usually based on the model itself, the record and information about the dataset. Typically, 055 MIAs use prediction logits of shadow models to train attack models, where the shadow training 056 data is obtained either through inference of the target model or through access to a potentially noisy 057 version of the original training dataset. Clearly, successful attacks lead to unacceptable information 058 leakage through the model. For example, if a GNN is trained on nodes belonging to a private group within a large social network - e.g., a support group for sensitive medical issues, successful MIAs on GNNs can reveal both the patient identity and their medical condition. As MIAs exploit the 060 differences of the outputs of target model on training and test datasets, most defense mechanisms 061 work towards suppressing the common patterns that quality attacks rely on (Shokri et al., 2017; Jia 062 et al., 2019; Choquette-Choo et al., 2021; Hayes et al., 2017; Leino & Fredrikson, 2020; Salem et al., 063 2018; Kaya & Dumitras, 2021; Yu et al., 2021; Wang et al., 2020; Chen et al., 2022). 064

Recent studies have focused on node (or edge) MIAs for node (or edge) classification downstream 065 tasks (Olatunji et al., 2021; Wu et al., 2021a; He et al., 2021; Conti et al., 2022). Compared to these 066 standard MIAs on GNNs, little attention has been paid to subgraph-structure MIAs (SMIAs) (Wang 067 & Wang, 2024). SMIAs can lead to an especially problematic form of privacy leakage, where at-068 tackers can infer not only if certain nodes were present in the training graph but also if there were 069 certain relationships between them. These relationships are usually captured via graph motives (triangles, cliques, paths etc). For instance, in a medical data network, cliques may involve members of 071 the same family and indicative of genetic/familial diseases. In this case, SMIA attacks can compro-072 mise not only individual medical histories, but – by association – whole family medical conditions. 073 In comparison, standard MIA methods can only reveal if the nodes were used for training a node 074 classifier, but not what the relationships between themselves and other nodes in the graph are (one 075 may think it plausible to perform individual node MIA attacks and then use link inference attacks to infer the subgraph induced by these nodes, but this process is both ineffective and usually of poor 076 utility). It is also important to point out that SMIA attacks differ from subgraph inference attacks 077 (SIAs) (Zhang et al., 2022) whose goal is to determine if a certain type of subgraph is present in the training dataset (without revealing the nodes that constitute the subgraph). 079

080 The only prior line of work on SMIAs is Wang & Wang (2024). There, the authors propose an 081 attack method based on the use of similarities between posterior vectors of the target nodes which allows them to generate training data for the attack model. On the defense side, the authors calculate the importance of each dimension of the node embeddings using the SHAP algorithm (Scott et al., 083 2017), and add noise to the least important dimensions in an attempt to balance out the quality of the 084 defence against SMIAs and the utility of the model. However, this attack and defense mechanisms 085 have several notable limitations. First, the similarity metrics used for the attack are fixed and cannot be adapted even when the data distribution has changed. Second, the attack performs well only on 087 homophilic graphs, but as we subsequently show, offers poor performance on heterophilic graphs. 880 Third, the similarity calculations split the attack process into two different parts, preventing an end-089 to-end approach, which results in an increase of the complexity of practical implementation. Finally, the defense still relies on the addition of noise to the embeddings, which inevitably compromises the 091 model utility. To address these issues, we propose both a new attack and improved defense system 092 that can counter both the attack of the original SMIA and our improved attack. The gist of our approach is to replace the fixed similarity calculation with learnable multiset functions that allow 093 for an end-to-end attack that dynamically adapts the training data generation process to the dataset, thereby offering excellent performance on both homophilic and heterophilic datasets. Furthermore, we propose a novel two-stage defense (TSD) strategy followed by flattening (Chen et al., 2022) that 096 does not significantly sacrifice model utility but provides stronger defense capabilities compared to noise addition. The key intuition behind the defense is to obfuscate the posterior distributions via 098 controlled overfitting and to add flattening noise that can also obfuscate the graph information (note that flattening noise cannot be directly "translated" to adding noise to node embeddings). Under 100 the TSD approach, the attack approach learned on GNNs without defense mechanisms cannot be 101 easily adapted to those which use defense, leading to a drop in the attack performance. In parallel, 102 flattening (Chen et al., 2022) allows us to increase the variance of the training loss distribution and 103 provides another mechanism for ensuring different posterior distribution for test and trainsets (since 104 different loss distributions lead to different posterior distributions). A detailed description of the approach is delegated to Section 4.

105

106 To demonstrate the utility of our novel SMIA attack and defense, we adopt the following standard 107 modeling assumptions. First, we focus on black-box attacks. Second, unlike some other approaches



Figure 1: (a) Diagram of the defense method SHAP-based noise addition (SHNA) proposed by Wang 123 & Wang (2024). Here GNNs are split into two parts: graph encoder (GE) that embeds node features 124 and graph topology as node embeddings, and fully connected (FC) layer that transform node em-125 beddings into posteriors. (b) Diagram of the two-stage training schedule that leads to an improved 126 defense against SMIAs. The key idea behind the TSD approach is to use the predicted labels of 127 test nodes from the first stage as psudolabels and switch the train and testset in the second stage. 128 During both training stages, flattening is performed as well to alleviate overfitting (see Section 4). (c) Comparison of standard SMIA (similarity-based) and our new end-to-end SMIA approach with 129 improved and generalized attacks. The standard SMIA model uses pairwise similarity of posteriors 130 to generate data for training the attack model. In contrast our end-to-end SMIA replaces similarity 131 metrics with multiset functions that directly deal with posteriors in a permutation invariant fashion. 132 This allows the attacker to be trained end-to-end from data preparation to model training. 133

that only describe attacks for specific GNN architectures or fixed defense mechanisms (Olatunji 135 et al., 2021), we design our defense strategy to be generally applicable to different graph learning 136 paradigms, datasets and resistant to both standard and improved SMIA attacks. Our extensive ex-137 periments, performed on three homophilic (CiteSeer, Facebook, LastFM) and one heterophilic dat-138 set (Chameleon), coupled with GCN, GAT, SGC and GPRGNN networks (Kipf & Welling, 2016; 139 Velickovic et al., 2017; Wu et al., 2019; Chien et al., 2020), reveal that our end-to-end attack model 140 outperforms the standard SMIA by 12.31% when no defense is used. We also demonstrate an aver-141 age 14.30% decrease in attack AUROC and an 10.05% improvement in target model utility trained 142 with our new defense technique compared to classical defenses, and under the improved attack.

143 144 145

146

### 2 RELATED WORKS

Due to space limitations, we only provide a brief summary of related results and delegate a more detailed discussion to Appendix A.

149 MIAs. The concept of MIA was first proposed in Homer et al. (2008) and later extended in various 150 directions, ranging from white-box settings (Nasr et al., 2019; Rezaei & Liu, 2021; Melis et al., 151 2019; Leino & Fredrikson, 2020) to black-box setting (Shokri et al., 2017; Salem et al., 2018; Song 152 & Mittal, 2020; Li & Zhang, 2021; Choquette-Choo et al., 2021; Carlini et al., 2022). Upon iden-153 tification of the informative features (e.g., posterior predictions, loss values, gradient norms, etc.) that reveal the sample membership, the attacker can choose to learn either a binary classifier (Shokri 154 et al., 2017) or metric-based decisions (Yeom et al., 2018; Salem et al., 2018) from a shadow model 155 trained on a shadow dataset to extract patterns within features of the training samples to determine 156 the membership. The description of a standard MIA pipeline is available in Appendix C. 157

Defense Against MIAs. As MIAs exploit the behavioral differences of the target model on trainsets and testsets, most defense mechanisms work by suppressing the common patterns among the
two. Frequently used defense methods include confidence score masking (Shokri et al., 2017; Jia
et al., 2019; Yang et al., 2020; Li et al., 2021; Choquette-Choo et al., 2021; Hanzlik et al., 2021,
regularization (Hayes et al., 2017; Salem et al., 2018; Leino & Fredrikson, 2020; Wang et al., 2020;

162 Choquette-Choo et al., 2021; Kaya & Dumitras, 2021; Yu et al., 2021; Chen et al., 2022), knowledge 163 distillation (Shejwalkar & Houmansadr, 2020; Tang et al., 2022), and differential privacy (Naseri 164 et al., 2020; Saeidian et al., 2021). For example, confidence score masking aims to hide the true 165 prediction vector returned by the target model and it thus mitigates the effectiveness of MIAs by 166 providing only top-k logits per inference (Shokri et al., 2017), or adding noise to the prediction vector in an adversarial manner (Jia et al., 2019). Regularization aims to reduce the degree of over-167 fitting of target models to mitigate MIAs (Choquette-Choo et al., 2021; Hayes et al., 2017; Leino & 168 Fredrikson, 2020; Salem et al., 2018; Kaya & Dumitras, 2021; Yu et al., 2021; Wang et al., 2020; Chen et al., 2022). Knowledge distillation aims to transfer the knowledge from an unprotected 170 model to a protected model (Shejwalkar & Houmansadr, 2020), while differential privacy (Saeidian 171 et al., 2021) protects membership information via noise injection and offers theoretical performance 172 guarantees, at the cost of substantial utility drop. 173

MIAs and Defense Strategies for GNNs. There are a handful of research that focuses on extending
MIA and corresponding defense mechanisms to graph learning framework. Olatunji et al. (2021)
analyzed graph MIA in two settings (train on subgraph, test on subgraph/full), and proposed the
LBP defense based on the confidence score masking idea, He et al. (2021) proposed zero-hop and
two-hop attacks designed for inductive GNNs, Wang & Wang (2023) studied the link membership
inference problem in an unsupervised fashion, and Chen et al. (2024) developed MaskArmor based
on masking and distillation technique. Besides node and edge MIAs, Zhang et al. (2022); Wang &
Wang (2024) also explored the subgraph attacks along with perturbation-based defense mechanism.

k-node Structure Membership Inference (k-SMIA) is a new form of a privacy attack (Wang & 182 Wang, 2024). The aim of this attack is to determine whether a collection of k nodes within the 183 training set belongs to a subgraph structure of interest (e.g., path or clique). k-SMIA hence uses a 184 new definition of membership: members are substructures of k nodes that form a relevant topology. 185 The original SMIA work outlines a novel black-box SMIA attack that combines training a three-186 label classifier with training shadows. This attack outperforms node-level MIAs followed by link 187 prediction. The defense against SMIA relies on perturbing embedding of the nodes with the smallest 188 "contribution" to the accuracy of the model and it results in a performance comparable to that offered 189 by differential privacy. Despite these positive initial findings, one has to point out that a performance matched by differential privacy methods may not be desirable, since the latter is usually significantly 190 reduced compared to defenseless models. Note that k-SMIA still constitutes an attack against node-191 level models, so that the downstream task of the target model remains node classification. 192

193 194

195

### **3 PROBLEM FORMULATION**

196 In this paper we focus on the downstream task as supervised node classifications; nevertheless, our method is applicable to different graph learning scenarios. Let  $\mathcal{G} = (X, A, Y, \mathcal{V}^{\text{Train}}, \mathcal{V}^{\text{Test}})$  denote 197 the graph dataset with node features  $X \in \mathbb{R}^{n \times d}$ , adjacent matrix  $A \in \mathbb{R}^{n \times n}$ , one-hot encoded node labels  $Y \in \mathbb{R}^{n \times C}$ , trainset  $\mathcal{V}^{\text{Train}}$  and testset  $\mathcal{V}^{\text{Train}}$ . Here *n* is the number of nodes, *d* is the feature dimension, *C* is the number of classes,  $\mathcal{V}^{\text{Train}}$  and  $\mathcal{V}^{\text{Test}}$  are disjoint and  $|\mathcal{V}^{\text{Train}}| + |\mathcal{V}^{\text{Test}}| = n$ . We 199 200 later on use  $Y^{\text{Train}}$  to denote the labels of  $\mathcal{V}^{\text{Train}}$ , and  $\hat{Y}^{\text{Test}}$  to denote the predicted labels of  $\mathcal{V}^{\text{Test}}$ . 201 Since X and A are already known during training, the goal of k-SMIA is determing the subgraph-202 substructure membership: given a set of k nodes  $\mathcal{V}_{att} \subset \mathcal{V}^{\text{Train}} \cup \mathcal{V}^{\text{Test}}$ , determine if  $\mathcal{V}_{att}$  forms 203 either a k-clique or a (k-1)-hop path in  $\mathcal{G}$  (member) or not (non-member). Meanwhile, the goal of 204 **MIA** is determing the **label membership**: given a node  $v \in \mathcal{V}^{\text{Train}} \cup \mathcal{V}^{\text{Test}}$ , determine if  $v \in \mathcal{V}^{\text{Train}}$ 205 (member) or not (non-member). Following the practice of SMIA and MIA, we also need a shadow 206 dataset  $\mathcal{G}_s = (X_s, A_s, Y_s, \mathcal{V}_s^{\text{Train}}, \mathcal{V}_s^{\text{Test}})$  to train the shadow model, and  $\mathcal{G}_s$  can be different from  $\mathcal{G}$ . 207

208

### 4 THE TSD METHOD

209 210

To defend against SMIAs, we introduce a *two-stage defense* method (TSD) to train *target* GNNs, depicted in Algorithm 1. The key intuition behind our method is to change the posterior distribution.

The first stage of TSD involves using  $\mathcal{G}$  to train a model checkpoint  $M_1(\theta_1)$  with parameters  $\theta_1$ . We adopt a flattening strategy in the first stage as a form of regularization, inspired by (Chen et al.,

2022). The flattening is implemented by transforming hard labels (one-hot) to soft labels (probability

vectors) when the loss on the trainset falls below a threshold  $\alpha$ . For simplicity, we assign the value

 $\beta$  to the ground-truth class, and  $\frac{1-\beta}{C-1}$  to the other classes. Here  $\alpha, \beta$  are hyperparameters while *C* denotes the number of classes. Note that we only use soft labels to compute the loss when the loss is small enough to keep the model utility as high as possible. The key of flattening is to increase the mean and variance of the training loss distribution while introducing noise to the label distribution.

221	Algorithm 1 TSD Training Procedure
222	<b>Input:</b> Dataset $\mathcal{C} = (X \land Y \lor \mathcal{V}^{\text{Train}} \lor \mathcal{V}^{\text{Test}})$
223	number of training epochs E, learning rate $\gamma$ .
224	number of classes $C$ , loss threshold $\alpha$ , flatten-
225	ing parameter $\beta$ .
226	<b>Output:</b> Second stage target model $M_2(\theta_2)$
227	First stage:
228	Perform random initialization for the first stage
229	model $M_1(\theta_1)$
230	for epoch $\in [1, E]$ do
231	if loss $L(M_1(\theta_1), Y^{\text{Ham}}) \ge \alpha$ then
232	$\theta_1 \leftarrow \theta_1 - \gamma \cdot \nabla_{\theta_1} L(M_1(\theta_1), Y^{\text{man}})$
233	Construct soft labels S <sup>Train</sup> where s <sup>Train</sup>
234	Construct soft labels $S^{\text{construct}}$ where $S_c^{\text{construct}} = (\rho_c; \epsilon_c, \text{Train})$
235	$\begin{cases} \beta, \text{ if } y_c^{\text{num}} = 1; \\ \beta & \beta & \beta \\ \beta & \beta & \beta \\ \beta & \beta & \beta \\ \beta & \beta &$
236	$\left( (1-\beta)/(C-1), \text{ otherwise} \right)$
237	$\theta_1 \leftarrow \theta_1 - \gamma \cdot \nabla_{\theta_1} L(M_1(\theta_1), S^{\text{Train}})$
238	end if
239	end for
240	Second stage: Initialize the second stage model $M(0)$ with
241	Initialize the second stage model $M_2(\theta_2)$ with
242	$M_1(\theta_1)$ , and generate pseudolabels Y for the original testest by information $M_1(\theta_1)$
243	the original testset by interence via $M_1(\theta_1)$
244	in cpoint $\in [1, L]$ in $\hat{V}^{\text{Test}} > \alpha$ there
245	If loss $L(M_2(\theta_2), Y^{-1}) \ge \alpha$ then
246	$\theta_2 \leftarrow \theta_2 - \gamma \cdot \nabla_{\theta_2} L(M_2(\theta_2), Y^{\text{rest}})$
247	eise Construct soft lobala C <sup>Test</sup> where s <sup>Test</sup>
248	Construct soft labels b where $s_c^{-1} = (\rho_c : c \cap \text{Test}_{-1})$
249	$\begin{cases} p, \text{ If } y_c^{cos} = 1; \\ 0, 0, 0, 0 \end{cases}$
250	$\left(\frac{(1-\beta)}{(C-1)}\right)$ , otherwise
251	$\theta_2 \leftarrow \theta_2 - \gamma \cdot \nabla_{\theta_2} L(M_2(\theta_2), S^{\text{lest}})$
252	end if
	end for

The second stage of TSD is similar to the first one, with the main difference that we instead use  $(X, A, \hat{Y}^{\text{Test}})$  for training. The pseudolabels  $\hat{Y}^{\text{Test}}$  are generated via inference of the checkpoint  $M_1(\theta_1)$  on the testset. Instead of random initialization, we use checkpoint initialization for the second stage model  $M_2(\theta_2)$ to resume training. The subsequent training process also proceeds with flattening, and  $M_2(\theta_2)$  is the final output target model. The role of the second stage is to include the testset into training, even without having access to their groundtruth labels  $Y^{\text{Test}}$ . In this case, the testset is also "trained", as it undergoes through the same process as the trainset.

SHAP-based noise addition (SHNA) defense. The SHNA defense (Wang & Wang, 2024) aims to weaken the attack by introducing noise into the node embeddings. To reduce the impact of noise on the model performance, SHNA uses the SHAP algorithm (Scott et al., 2017) to quantify the contribution of each embedding dimension to the final classification accuracy, and selectively adds noise to the  $l = \lfloor r \times d \rfloor$  dimensions with the smallest contributions, where  $r \in (0, 1]$  is a hyperparameter, and d is the embedding dimension. The added noise is Laplacian, with scale b and mean  $\mu$ .

**Comparison with SHNA.** The reasons why TSD can outperform SHNA are two-fold. First, as we explain in Section 5, the attacker targets the posterior distributions directly rather than the node embeddings. Therefore, adding noise to the node embeddings is only an *indirect* method of perturbing the attack model's input,

and the relationship between the noise scale and the attackers performance is unclear. In contrast, TSD directly alters the posterior distribution, as evidenced by the significant change in the loss distribution. Figure 2 illustrates this point: without TSD (standard training), the average training loss is lower than the testing loss. With TSD, however, the discrepancy in the loss distribution is notably reduced. Additionally, we conduct another simulation to examine the correlation between attack performance and loss distribution. Note that when a clique exists in the graph, three scenarios are possible: (1) all nodes belong to the trainset; (2) nodes belong to both the train and testset; and (3) all nodes belong to the testset. The attack accuracy for triangle identification in these cases is 0.9722, 0.9347, and 0.8526, respectively. These results confirm that the attack accuracy is positively correlated with how well the posterior distribution is learned.

262 263

266

253

254

256

257

258

259

260

261

### <sup>264</sup> 265 5

### END-TO-END SMIA: A MULTISET FUNCTION APPROACH

To demonstrate that our defense mechanism is effective against various types of SMIAs, we evaluate
 it using both the similarity-based attacks introduced by Wang & Wang (2024) and a novel attack
 based on multiset functions. We begin by reviewing the similarity-based attacks, highlighting their
 strengths and limitations, before introducing our new end-to-end attack scheme.

**Similarity-Based Attacks.** The core of similarity-based SMIA lies in processing the shadow model's output to generate data for attack model training. For example, for a k = 3 clique, the shadow model provides three posterior vectors for three nodes. Pairwise similarities between these vectors are computed and sorted in ascending order to form a new vector. This process is repeated for three different similarity metrics, namely the dot product, cosine similarity, and  $\ell_2$  distance. The sorted vectors for each metric are concatenated to create the training data for the attack model along with the structure labels (see also Appendix B).

277 While similarity-based attacks have been 278 shown to outperform other link-level attacks for 279 SMIA, they have limitations. First, the choice 280 of similarity metrics is heuristic and may not 281 be effective in all scenarios. Second, the com-282 putational complexity is  $O(k^3)^1$ , which is pro-283 hibitive for large k.

284 Multiset Function-Based Attacks. To ad-285 dress the limitations of similarity-based attacks, we propose using multiset functions to aggre-286 287 gate posterior vectors, rather than relying on heuristic similarity measures. The primary ad-288 vantage of multiset functions is that they can 289 directly process posterior vectors without ad-290 ditional preprocessing. Furthermore, multiset 291



Figure 2: Comparison of training and test loss distribution under (a) standard training; (b) TSD training. The simulation is conducted on Chameleon with NLGCN as the backbone. The attack is our end-to-end SMIA.

functions treat inputs as sets, ensuring that the output remains unaffected by the input order. We 292 implement the multiset function using Deep Sets (Zaheer et al., 2017), which has proven effective 293 in other contexts requiring permutation invariance, such as hypergraph GNNs (Chien et al., 2021b). For example, consider three nodes  $v_1, v_2, v_3$  with corresponding posterior distributions  $p_1, p_2, p_3$ . 295 The final embedding used for training the attacker is computed as  $\rho(\phi(p_1) + \phi(p_2) + \phi(p_3))$ , where 296  $\rho$  and  $\phi$  are neural networks (specifically, two-layer MLPs in our simulations), which makes the 297 attack model trainable in an end-to-end fashion. Based on the performance guarantees of Deep Sets and the universal approximation theorem, it is straightforward to show that the similarity-based 298 approach is a special case of the multiset function-based method. 299

300 301

302 303

304

### 6 EXPERIMENTS

6.1 EXPERIMENTS PERTAINING TO SMIA DEFENSE

Datasets and GNN Baselines. We train four GNN (GCN (Kipf & Welling, 2016), GAT (Velickovic
et al., 2017), SGC (Wu et al., 2019), GPRGNN (Chien et al., 2020)) on three homophilic datasets
(CiteSeer, Facebook (Yang et al., 2023), LastFM (Rozemberczki & Sarkar, 2020)), and four GNN
(NLGCN, NLGAT, NLMLP (Wang et al., 2018), GPRGNN (Chien et al., 2020)) on one heterophilic
datasets (Chameleon (Rozemberczki et al., 2021)). Detailed experimental setups, properties and
statistics of the datasets are available in Appendix F.

Evaluation Metrics. We follow previous literature to use the following two metrics for evaluation.
 We report classification accuracy of the target model on the testset (CA) to measure model utility,
 and AUROC scores of the attack model (AU), which is a widely used approach in the field of
 MIA (Carlini et al., 2022). Note that good defenses should lead to large CAs and small AUs.

315 316

### 6.1.1 COMPARISON OF SMIA DEFENSE METHODS AND ATTACK METHOD

Table 1 and 2 present a comparison of End2end-SMIA with Standard-SMIA, as well as the defense performance of TSD and SHAN. Additional results can be found in Appendix G. Standard-SMIA refers to similarity-based attacks while End2end-SMIA refers to the multiset function-based attack. For simplicity, we refer to them as S-SMIA and E-SMIA, respectively. In the experiments, SHAN uses the noise addition strategy from the original work for homophilic datasets, r = 0.4 and b = 0.3, and for the heterophilic dataset, r = 0.1 and b = 0.01.

<sup>323</sup> 

 $<sup>{}^{1}</sup>O(k^{2})$  pairs with O(k) computational complexity per pair.

Table 1: 3-SMIA attack performance comparison between Standard-SMIA and End2end-SMIA & Defense performance comparison between SHNA and TSD. Compared to SHNA, TSD achieves a decrease in attack AUROC by 16.64% against Standard-SMIA & 14.30% against End2end-SMIA, and increases in utility performance by 10.05% (on average). Full table with variance information is available in Appendix Table 8. 

Dataset	Models	CA(SHNA)	CA(TSD)	S-SMIA AU (SHNA)	S-SMIA AU (TSD)	E-SMIA AU(SHNA)	E-SMIA AU(TSI
	GCN	0.6887	0.7729	0.9157	0.7767	0.9684	0.8703
CitaSaar	GAT	0.7082	0.7548	0.8784	0.7145	0.9336	0.8519
Chescei	SGC	0.6982	0.7454	0.9015	0.7978	0.9862	0.8904
	GPRGNN	0.6736	0.7864	0.8222	0.5969	0.9006	0.7364
	GCN	0.6028	0.7095	0.7143	0.5463	0.7942	0.6127
Eastral	GAT	0.5321	0.6435	0.6854	0.5382	0.8229	0.6068
гасевоок	SGC	0.5011	0.6423	0.7028	0.5253	0.7903	0.5945
	GPRGNN	0.5529	0.6544	0.6512	0.5562	0.8081	0.6213
	GCN	0.8256	0.8523	0.9215	0.8337	0.9825	0.8702
LastEM	GAT	0.8379	0.8662	0.8528	0.8664	0.9274	0.8821
Lastrivi	SGC	0.8154	0.8400	0.9305	0.8835	0.9840	0.9060
	GPRGNN	0.8320	0.8577	0.9004	0.8268	0.9543	0.8687
	NLGCN	0.6373	0.6725	0.8904	0.4538	0.9196	0.7875
Chamalaan	NLGAT	0.6437	0.6835	0.7025	0.5611	0.8848	0.7958
Chameleon	NLMLP	0.5010	0.5484	0.7128	0.5342	0.8152	0.7057
	GPRGNN	0.6637	0.6835	0.7191	0.6259	0.8929	0.7711

Table 2: 4-SMIA attack performance comparison between Standard-SMIA and End2end-SMIA & Defense performance comparison between SHNA and TSD. Compared to SHNA, TSD achieves a decrease in attack AUROC by 16.56% against Standard-SMIA & 13.89% against End2end-SMIA, and increases in utility performance by 10.05% (on average). Full table with variance information is available in Appendix Table 9.

Dataset	Models	CA(SHNA)	CA(TSD)	S-SMIA AU(SHNA)	S-SMIA AU(TSD)	E-SMIA AU(SHNA)	E-SMIA AU(TSD
	GCN	0.6887	0.7729	0.9685	0.9013	0.9786	0.9267
CitaSaar	GAT	0.7082	0.7548	0.9718	0.8868	0.9891	0.9191
Cheseer	SGC	0.6982	0.7454	0.9839	0.9316	0.9955	0.9442
	GPRGNN	0.6736	0.7864	0.9010	0.8523	0.9692	0.8894
	GCN	0.6028	0.7095	0.7377	0.5308	0.7964	0.6033
Faarbaal	GAT	0.5321	0.6435	0.6732	0.5027	0.7436	0.5978
Facebook	SGC	0.5011	0.6423	0.6992	0.5097	0.7520	0.6325
	GPRGNN	0.5529	0.6544	0.6520	0.4978	0.7090	0.5885
	GCN	0.8256	0.8523	0.9649	0.8296	0.9972	0.8561
LoctEM	GAT	0.8379	0.8662	0.8936	0.7442	0.9354	0.7654
Lastrivi	SGC	0.8154	0.8401	0.9574	0.8184	0.9979	0.8270
	GPRGNN	0.8320	0.8577	0.9144	0.8044	0.9844	0.8444
	NLGCN	0.6373	0.6725	0.7663	0.5668	0.9952	0.8346
<b>a</b> 1	NLGAT	0.6437	0.6835	0.7849	0.5901	0.9816	0.8409
Chameleon	NLMLP	0.5010	0.5484	0.7531	0.6076	0.9186	0.8359
	GPRGNN	0.6637	0.6835	0.8358	0.5924	0.9486	0.8127

When comparing the attack AUROC on the same dataset, using the same model and defense method, End2end-SMIA is noticeably better than Standard-SMIA, especially on the heterophilic Chameleon. This is due to the trainable Deep Sets, which extract structural information of the node sets from their posterior vectors, significantly reducing the learning difficulty for the attack model. Additionally, whereas the similarity metrics chosen in Standard-SMIA may not be well-suited to the attacked node sets, Deep Sets can continuously adjust during training to offer improved performance.

When comparing TSD and SHAN under the same conditions, we observe that TSD offers superior defensive capabilities while maintaining higher classification accuracy for the target model. This is because TSD alters the posterior distribution through an additional training phase and the intro-duction of flattening. These two methods are applied only after the target model has already com-pleted learning from the data, so the classification performance of the target model remains largely unaffected. In contrast, SHAN aims to minimize the impact on the target model's classification performance by adding noise only to the least important dimensions of the embeddings. However, this also means that the noise's ability to disrupt the attack decreases due to the reduced importance of those dimensions. Therefore, it is difficult to achieve both high model utility and strong defense ability, regardless of how the noise is added.

### 6.1.2 DATA TRANSFER



Figure 3: Result for 3-SMIA Attack AUROC on GCN under data transfer setting for (a) TSD against Standard-SMIA; (b) SHAN against Standard-SMIA; (c) TSD against End2end-SMIA; and (d) SHAN against End2end-SMIA.

In practical applications, the training data of the target model is typically unknown to the attacker. Therefore, to train a shadow model, attackers often use publicly available datasets as training data, which usually have different distributions from the training data of the target model. To better evaluate the effectiveness of the TSD defense method, we set up an experimental scenario that resembles real-world SMIA, where experiments are conducted for the data transfer setting. Data transfer refers to the situation where the training data for the target model and the shadow model come from different datasets. Evidently, when the training data for the target and shadow models come from the same dataset, the attack on the target model is the strongest, corresponding to the lower bound of the defense method's effectiveness. 

Figure 3 presents the results obtained using GCN as both the target and shadow model under the 3
SMIA attack. The experimental settings remain the same as described in Section 6.1.1. The results
demonstrate the continued effectiveness of the TSD defense method in the data transfer setting.
Across all dataset combinations, TSD consistently outperformed SHAN, indicating that TSD can
maintain excellent performance against SMIA attacks that may occur in real-world scenarios. Additionally, it is worth noting that End2end-SMIA exhibits stronger attack performance compared to
Standard-SMIA, which can be attributed to the powerful generalization capability of DeepSets.

### 

### 6.1.3 ABLATION STUDY OF TSDS

Table 3:	Ablation	study	for the	contributors	to the	utility/attack	gains/mitigation	for	TSD	under
3-SMIA	Attack.									

Method	Dataset	GNN Models	Classify Acc	Attack AUROC
Standard Training	CiteSeer Facebook LastFM Chameleon	GCN GCN GCN NLGCN		$\begin{array}{c} 0.9745 \pm 0.0084 \\ 0.8759 \pm 0.0091 \\ 0.9841 \pm 0.0094 \\ 0.9664 \pm 0.0084 \end{array}$
Flattening (One-Stage)	CiteSeer Facebook LastFM Chameleon	GCN GCN GCN NLGCN		$\begin{array}{c} 0.9428 \pm 0.0081 \\ 0.8249 \pm 0.0084 \\ 0.9520 \pm 0.0075 \\ 0.9187 \pm 0.0093 \end{array}$
Two-Stage (without Flattening)	CiteSeer Facebook LastFM Chameleon	GCN GCN GCN NLGCN		$\begin{array}{c} 0.9013 \pm 0.0081 \\ 0.6644 \pm 0.0091 \\ 0.9008 \pm 0.0091 \\ 0.8324 \pm 0.0074 \end{array}$
TSD (Two-Stage & Flattening)	CiteSeer Facebook LastFM Chameleon	GCN GCN GCN NLGCN		$\begin{array}{c} 0.8703 \pm 0.0075 \\ 0.6127 \pm 0.0096 \\ 0.8702 \pm 0.0088 \\ 0.7875 \pm 0.0098 \end{array}$

432 Our TSD method differs from the conventional training methods in two aspects: (1) the use of flattening operations; and (2) two-stage training. To demonstrate their roles in enhancing defense capability, we conducted the following ablation experiments.

In the experiments, we examined four variants: (1) standard training; (2) two-stage (without flattening); (3) flattening (one-stage); and (4) TSD. Here, standard training refers to training a target model only on the trainset; two-stage trains a target model in a train-test alternate fashion, equivalent to TSD without flattening; flattening is the same as described in Section 4, combined with one-stage training. Clearly, TSD is two-stage combined with flattening. We conduct our experiments using End2end-SMIA.

Table 3 presents the results of the ablation study for these four variants on four datasets and one GNN backbone – GCN. Additional results can be found in Appendix I. The findings indicate that the primary source of improvement for TSD is the two-stage training technique. Moreover, compared to the defense improvement brought by flattening, its impact on model utility is acceptable.

445 446

447

### 6.2 EXPERIMENTS PERTAINING TO MIA DEFENSE

448 TSD is also capable of defending against node-level MIA. Previous studies on MIA (Homer et al., 449 2008; Shokri et al., 2017; Salem et al., 2018; Song & Mittal, 2020) have pointed out that the key to a successful MIA attack is to allow overfitting in the target model. Since TSD has the ability 450 to reduce overfitting, it should intuitively have the capability to defend against MIA: the two-stage 451 training process in TSD can help reduce the gap between the average losses of training and testing 452 nodes, while the use of flattening reduces the difference in the variance of loss distributions between 453 the training and testing nodes. We demonstrate the effectiveness of TSD's defense by comparing it with other defense methods. We choose two representative defense methods, Laplacian Binned Pos-455 terior Perturbation (LBP) (Olatunji et al., 2021) on GNNs and Distillation for Membership Privacy 456 (DMP) (Shejwalkar & Houmansadr, 2020) on graphless models, as our defense baselines. LBP is 457 the state-of-the-art defense method for GNNs and it works by adding Laplacian noise to the poste-458 rior before it is released to the user. To reduce the amount of noise needed to distort the posteriors, 459 LBP does not add noise to each element of the posterior, but to binned posterior. In our experiments, 460 we first randomly shuffle the posteriors and then assign each posterior to a partition/bin. The total number of bins N is predefined based on the number of classes. For each bin, we sample noise with 461 scale b from the Laplace distribution. The sampled noise is added to each element in the bin. After 462 the noise added to each bin, we restore the initial positions of the noisy posteriors and release them. 463 Appendix Table 13 shows the best set of parameters for LBP that we used in our experiments. 464

465 On the other hand, we adapted DMP to the case of GNN training. DMP consists of three phases, namely pre-distillation, distillation and post-distillation. The pre-distillation phase trains an unpro-466 tected model on a private training data without any privacy protection. Next, in the distillation phase, 467 DMP selects reference data and transfers the knowledge of the unprotected model into predictions 468 om the reference data. Notice that private training data and the reference data have no intersec-469 tion. Finally, in the post-distillation phase, DMP uses the predictions to train a protected model. 470 Our experiments use the same model structure for the unprotected and protected models. To follow 471 the DMP procedure, we need to further split the trainset into a private dataset and reference dataset, 472 where the private dataset trains the unprotected models, and the reference dataset trains the protected 473 target model. Compared to DMP, TSD can directly train the target model with the full trainset. 474

In our experiments, the split ratio of trainsets and testsets for TSD and LBP is 1:1, and the split ratio 475 of private datasets, reference datasets and testsets in DMP is 0.45:0.45:0.1. Since MIA is easier to 476 defend against than SMIAs, our experiments are conducted on more complex and diverse datasets. 477 Table 4 and Table 5 shows partial result of our experiments, and the complete results can be found in 478 Appendix J.1 and J.2. The results indicate that our method achieves defense capabilities comparable 479 to LBP and DMP, while achieving better model utility performance. It is important to emphasize 480 that a lower attack AUROC does not necessarily indicate stronger defense. The stronger the defense, 481 the closer the AUROC should be to 0.5, as 0.5 represents random guessing. When the AUROC is 482 less than 0.5, the attackers can flip the prediction results to make the AUROC greater than 0.5. Compared to LBP, TSD significantly improved the model utility. The main reason is that LBP is a 483 perturbation-based method, which can potentially hurt the target model performance significantly. 484 However, TSD achieves defense by alleviating overfitting, which delves deeper into the core issue, 485 instead of adversely affecting target models. In addition, compared to DMP, TSD still achieves

higher model utility. This gain is expected, as TSD not only can make use of the full trainsets, but also utilizes testsets in the second stage, thus enhancing the model's generalization ability.

Table 4: Performance comparison between TSD and LBP. Compared to LBP, TSD has an average increase in utility by 17.28%, with a comparable attack AUROC.

Dataset	Models	Classify Acc (LBP)	Classify Acc (TSD)	Attack AUROC (LBP)	Attack AUROC (7
	GCN	$0.6886 \pm 0.0041$	$\textbf{0.8381} \pm \textbf{0.0023}$	$0.4998 \pm 0.0050$	$0.4990 \pm 0.0048$
DubMad	GAT	$0.7631 \pm 0.0037$	$\textbf{0.8400} \pm \textbf{0.0028}$	$0.5021 \pm 0.0084$	$0.4911 \pm 0.0061$
Publided	SGC	$0.6564 \pm 0.0035$	$\textbf{0.8080} \pm \textbf{0.0020}$	$0.5007 \pm 0.0065$	$\textbf{0.5005} \pm \textbf{0.0057}$
	GPRGNN	$0.7843 \pm 0.0029$	$\textbf{0.8553} \pm \textbf{0.0014}$	$\textbf{0.5003} \pm \textbf{0.0038}$	$0.4967 \pm 0.0034$
	GCN	$0.5195 \pm 0.0041$	$\textbf{0.6778} \pm \textbf{0.0049}$	$0.4912 \pm 0.0021$	$\textbf{0.4993} \pm \textbf{0.0023}$
Easthaalt	GAT	$0.5460 \pm 0.0037$	$\textbf{0.6519} \pm \textbf{0.0039}$	$0.5120 \pm 0.0025$	$\textbf{0.5010} \pm \textbf{0.0028}$
Facebook	SGC	$0.4833 \pm 0.0044$	$\textbf{0.6249} \pm \textbf{0.0043}$	$0.4901 \pm 0.0026$	$\textbf{0.5004} \pm \textbf{0.0031}$
	GPRGNN	$0.4627 \pm 0.0032$	$\textbf{0.5890} \pm \textbf{0.0035}$	$0.4807 \pm 0.0031$	$\textbf{0.5014} \pm \textbf{0.0020}$
	GCN	$0.6509 \pm 0.0037$	$\textbf{0.8378} \pm \textbf{0.0035}$	$0.5118 \pm 0.0024$	$\textbf{0.4971} \pm \textbf{0.0022}$
Lootfm	GAT	$0.7210 \pm 0.0034$	$\textbf{0.8683} \pm \textbf{0.0032}$	$0.5136 \pm 0.0031$	$\textbf{0.4980} \pm \textbf{0.0029}$
Lastini	SGC	$0.6395 \pm 0.0040$	$\textbf{0.8336} \pm \textbf{0.0045}$	$0.5121 \pm 0.0030$	$\textbf{0.4965} \pm \textbf{0.0034}$
	GPRGNN	$0.6875 \pm 0.0038$	$\textbf{0.8443} \pm \textbf{0.0033}$	$0.5101 \pm 0.0035$	$\textbf{0.4999} \pm \textbf{0.0025}$
	NLGCN	$0.5987 \pm 0.0068$	$\textbf{0.6657} \pm \textbf{0.0062}$	$0.5233 \pm 0.0062$	$\textbf{0.4954} \pm \textbf{0.0065}$
Channelson	NLGAT	$0.5926 \pm 0.0072$	$\textbf{0.6585} \pm \textbf{0.0070}$	$0.5246 \pm 0.0065$	$\textbf{0.4902} \pm \textbf{0.0063}$
Chameleon	NLMLP	$0.4281 \pm 0.0078$	$\textbf{0.4824} \pm \textbf{0.0074}$	$0.5623 \pm 0.0057$	$\textbf{0.4848} \pm \textbf{0.0051}$
	GPRGNN	$0.5230 \pm 0.0054$	$\textbf{0.6550} \pm \textbf{0.0058}$	$0.5107 \pm 0.0060$	$\textbf{0.4936} \pm \textbf{0.0049}$

Table 5: Performance comparison between TSD and DMP. Compared to DMP, TSD has an average increase in utility by 4.35%, with a comparable attack AUROC.

Dataset	Models	Classify Acc (DMP)	Classify Acc (TSD)	Attack AUROC (DMP)	Attack AUROC (TSD)
	GCN GAT	$\begin{array}{c} 0.8235 \pm 0.0037 \\ 0.8027 \pm 0.0047 \end{array}$	$\begin{array}{c} \textbf{0.8387} \pm \textbf{0.0034} \\ \textbf{0.8434} \pm \textbf{0.0024} \end{array}$	$0.5026 \pm 0.0039$ 0.5013 ± 0.0036	$\begin{array}{c} \textbf{0.4978} \pm \textbf{0.0042} \\ \textbf{0.5005} \pm \textbf{0.0043} \end{array}$
PubMed	SGC	$0.8013 \pm 0.0041$	$0.8096 \pm 0.0045$	$0.5013 \pm 0.0030$ $0.5024 \pm 0.0042$	$0.5003 \pm 0.0043$ $0.5003 \pm 0.0038$
	GPRGNN	$0.8104 \pm 0.0031$	$0.8423 \pm 0.0036$	$0.5020 \pm 0.0027$	$0.4994 \pm 0.0023$
Easthaalt	GCN GAT	$\begin{array}{c} 0.6896 \pm 0.0046 \\ 0.6420 \pm 0.0040 \end{array}$	$\begin{array}{c} 0.7054 \pm 0.0045 \\ 0.6797 \pm 0.0043 \end{array}$	$\begin{array}{c} 0.4806 \pm 0.0032 \\ 0.4768 \pm 0.0042 \end{array}$	$\begin{array}{c} 0.4964 \pm 0.0030 \\ 0.4910 \pm 0.0037 \end{array}$
Facebook	SGC	$0.6152 \pm 0.0049$ 0.5784 ± 0.0051	$\begin{array}{c} \textbf{0.6351} \pm \textbf{0.0040} \\ \textbf{0.5920} \pm \textbf{0.0046} \end{array}$	$0.4821 \pm 0.0038$ 0.4780 ± 0.0045	$\begin{array}{c} \textbf{0.4969} \pm \textbf{0.0039} \\ \textbf{0.4964} \pm \textbf{0.0032} \end{array}$
		$0.3784 \pm 0.0051$	$0.3920 \pm 0.0040$ 0.8401 ± 0.0055	$0.4780 \pm 0.0043$	$0.4978 \pm 0.0032$
Lastfm	GAT	$\begin{array}{c} 0.8102 \pm 0.0054 \\ 0.8434 \pm 0.0054 \end{array}$	$0.8769 \pm 0.0041$	$0.5142 \pm 0.0028$	$0.4972 \pm 0.0023$
	SGC GPRGNN	$\begin{array}{c} 0.8112 \pm 0.0050 \\ 0.8140 \pm 0.0042 \end{array}$	$\begin{array}{c} \textbf{0.8414} \pm \textbf{0.0044} \\ \textbf{0.8485} \pm \textbf{0.0037} \end{array}$	$\begin{array}{c} 0.5091 \pm 0.0036 \\ 0.5119 \pm 0.0026 \end{array}$	$\begin{array}{c} \textbf{0.4967} \pm \textbf{0.0031} \\ \textbf{0.4979} \pm \textbf{0.0019} \end{array}$
	NLGCN	$0.6681 \pm 0.0064$	$\textbf{0.6963} \pm \textbf{0.0065}$	$0.5210 \pm 0.0064$	$\textbf{0.5182} \pm \textbf{0.0062}$
Chameleon	NLGAT NLMLP	$0.6516 \pm 0.0066$ $0.4643 \pm 0.0079$	$\begin{array}{c} \textbf{0.7082} \pm \textbf{0.0073} \\ \textbf{0.4955} \pm \textbf{0.0070} \end{array}$	$\begin{array}{c} \textbf{0.5116} \pm \textbf{0.0061} \\ \textbf{0.5054} \pm \textbf{0.0053} \end{array}$	$\begin{array}{r} 0.5159 \pm 0.0059 \\ \textbf{0.5017} \pm \textbf{0.0056} \end{array}$
	GPRGNN	$0.6471 \pm 0.0060$	$0.6934 \pm 0.0068$	$0.5172 \pm 0.0067$	$0.5163 \pm 0.0045$

7 CONCLUSIONS AND LIMITATIONS

We proposed a novel two-stage defense method (TSD) against SMIAs for GNNs, and leveraged multiset functions to enhance SMIA attacks (End2end-SMIA). Our evaluation showed that TSD surpasses SHAN in defending against various types of SMIAs, establishing a new state-of-the-art benchmark. Additionally, we compared TSD with LBP and DMP in defending against classical node-level MIAs, demonstrating that TSD also performs effectively in protecting against the normal attacks. We conducted ablation studies and validated the origin of TSD's defense capability. TSD exhibits superior performance and is easy to integrate into various GNNs training processes.

**Limitations and Future Work.** The current form of TSD has the following limitations: (1) It can lead to lower model utility because the labels used in second stage are pseudolabels of test nodes, instead of groud-truth labels; (2) The flattening parameter  $\beta$  is not end-to-end learnable, and uniform flattening may not be the optimal way to counter SMIAs and MIAs. To address these limitations, we are exploring to use only the test nodes with high confidence predictions and change the formula of soft labels to make  $\beta$  learnable.

## 540 REFERENCES

547

554

566

567

568

569

577

578

579 580

- Fedor Borisyuk, Shihai He, Yunbo Ouyang, Morteza Ramezani, Peng Du, Xiaochen Hou, Chengming Jiang, Nitin Pasumarthy, Priya Bannur, Birjodh Tiwana, et al. Lignn: Graph neural networks at linkedin. *arXiv preprint arXiv:2402.11139*, 2024.
- <sup>545</sup> Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramer. Mem <sup>546</sup> bership inference attacks from first principles, 2022.
- Chenyang Chen, Xiaoyu Zhang, Hongyi Qiu, Jian Lou, Zhengyang Liu, and Xiaofeng Chen.
   Maskarmor: Confidence masking-based defense mechanism for gnn against mia. *Information Sciences*, 669:120579, 2024.
- <sup>551</sup> Dingfan Chen, Ning Yu, and Mario Fritz. Relaxloss: Defending membership inference attacks
   <sup>552</sup> without losing utility. In *International Conference on Learning Representations*, 2022. URL
   <sup>553</sup> https://openreview.net/forum?id=FEDfGWVZYIn.
- Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. Adaptive universal generalized pagerank graph neural network. *arXiv preprint arXiv:2006.07988*, 2020.
- Eli Chien, Wei-Cheng Chang, Cho-Jui Hsieh, Hsiang-Fu Yu, Jiong Zhang, Olgica Milenkovic, and
   Inderjit S Dhillon. Node feature extraction by self-supervised multi-scale neighborhood predic tion. arXiv preprint arXiv:2111.00064, 2021a.
- Eli Chien, Chao Pan, Jianhao Peng, and Olgica Milenkovic. You are allset: A multiset function framework for hypergraph neural networks. *arXiv preprint arXiv:2106.13264*, 2021b.
- Christopher A Choquette-Choo, Florian Tramer, Nicholas Carlini, and Nicolas Papernot. Label-only
   membership inference attacks. In *International conference on machine learning*, pp. 1964–1974.
   PMLR, 2021.
  - Mauro Conti, Jiaxin Li, Stjepan Picek, and Jing Xu. Label-only membership inference attack against node-level graph neural networks. In *Proceedings of the 15th ACM Workshop on Artificial Intelligence and Security*, pp. 1–12, 2022.
- Neil Zhenqiang Gong and Bin Liu. Attribute inference attacks in online social networks. ACM
   Trans. Priv. Secur., 21(1), January 2018. ISSN 2471-2566. doi: 10.1145/3154793. URL https:
   //doi.org/10.1145/3154793.
- Lucjan Hanzlik, Yang Zhang, Kathrin Grosse, Ahmed Salem, Maximilian Augustin, Michael Backes, and Mario Fritz. Mlcapsule: Guarded offline deployment of machine learning as a service. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 3300–3309, 2021.
  - Jamie Hayes, Luca Melis, George Danezis, and Emiliano De Cristofaro. Logan: Membership inference attacks against generative models. *arXiv preprint arXiv:1705.07663*, 2017.
  - Xinlei He, Rui Wen, Yixin Wu, Michael Backes, Yun Shen, and Yang Zhang. Node-level membership inference attacks against graph neural networks. *arXiv preprint arXiv:2102.05429*, 2021.
- Nils Homer, Szabolcs Szelinger, Margot Redman, David Duggan, Waibhav Tembe, Jill Muehling, John V Pearson, Dietrich A Stephan, Stanley F Nelson, and David W Craig. Resolving individuals contributing trace amounts of dna to highly complex mixtures using high-density snp genotyping microarrays. *PLoS genetics*, 4(8):e1000167, 2008.
- Hongsheng Hu, Zoran Salcic, Lichao Sun, Gillian Dobbie, Philip S Yu, and Xuyun Zhang. Membership inference attacks on machine learning: A survey. ACM Computing Surveys (CSUR), 54 (11s):1–37, 2022.
- Jinyuan Jia, Ahmed Salem, Michael Backes, Yang Zhang, and Neil Zhenqiang Gong. Memguard:
   Defending against black-box membership inference attacks via adversarial examples. 2019.
- 593 Yigitcan Kaya and Tudor Dumitras. When does data augmentation help with membership inference attacks? In *International conference on machine learning*, pp. 5345–5355. PMLR, 2021.

605

633

- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- Klas Leino and Matt Fredrikson. Stolen memories: Leveraging model memorization for calibrated
   {White-Box} membership inference. In 29th USENIX security symposium (USENIX Security 20),
   pp. 1605–1622, 2020.
- Jiacheng Li, Ninghui Li, and Bruno Ribeiro. Membership inference attacks and defenses in classification models. In *Proceedings of the Eleventh ACM Conference on Data and Application Security and Privacy*, pp. 5–16, 2021.
  - Zheng Li and Yang Zhang. Membership leakage in label-only exposures. In *Proceedings of the* 2021 ACM SIGSAC Conference on Computer and Communications Security, pp. 880–895, 2021.
- Luca Melis, Congzheng Song, Emiliano De Cristofaro, and Vitaly Shmatikov. Exploiting unintended
   feature leakage in collaborative learning. In 2019 IEEE symposium on security and privacy (SP),
   pp. 691–706. IEEE, 2019.
- Mohammad Naseri, Jamie Hayes, and Emiliano De Cristofaro. Toward robustness and privacy
   in federated learning: Experimenting with local and central differential privacy. *arXiv preprint arXiv:2009.03561*, 2020.
- Milad Nasr, Reza Shokri, and Amir Houmansadr. Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning. In 2019 IEEE Symposium on Security and Privacy (SP). IEEE, May 2019. doi: 10.1109/sp.2019.00065.
   URL http://dx.doi.org/10.1109/SP.2019.00065.
- Iyiola E Olatunji, Wolfgang Nejdl, and Megha Khosla. Membership inference attack on graph
   neural networks. In 2021 Third IEEE International Conference on Trust, Privacy and Security in
   Intelligent Systems and Applications (TPS-ISA), pp. 11–20. IEEE, 2021.
- Aditya Pal, Chantat Eksombatchai, Yitong Zhou, Bo Zhao, Charles Rosenberg, and Jure Leskovec. Pinnersage: Multi-modal user embedding framework for recommendations at pinterest. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '20, pp. 2311–2320, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450379984. doi: 10.1145/3394486.3403280. URL https: //doi.org/10.1145/3394486.3403280.
- 627
   628
   629
   629 Shahbaz Rezaei and Xin Liu. On the difficulty of membership inference attacks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7892–7900, 2021.
- Benedek Rozemberczki and Rik Sarkar. Characteristic functions on graphs: Birds of a feather, from statistical descriptors to parametric models, 2020. URL https://arxiv.org/abs/2005.
   07959.
- Benedek Rozemberczki, Carl Allen, and Rik Sarkar. Multi-Scale attributed node embedding. Journal of Complex Networks, 9(2):cnab014, 05 2021. ISSN 2051-1329. doi: 10.1093/comnet/cnab014. URL https://doi.org/10.1093/comnet/cnab014.
- <sup>637</sup> Alexandre Sablayrolles, Matthijs Douze, Yann Ollivier, Cordelia Schmid, and Hervé Jégou. White <sup>638</sup> box vs black-box: Bayes optimal strategies for membership inference. 2019.
- Sara Saeidian, Giulia Cervia, Tobias J Oechtering, and Mikael Skoglund. Quantifying membership
   privacy via information leakage. *IEEE Transactions on Information Forensics and Security*, 16: 3096–3108, 2021.
- Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes.
   Ml-leaks: Model and data independent membership inference attacks and defenses on machine learning models. 2018.
- 647 M Scott, Lee Su-In, et al. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30:4765–4774, 2017.

667

677

684

688

689

690

691

- 648
   649
   650
   650
   Virat Shejwalkar and Amir Houmansadr. Membership privacy for machine learning models through knowledge transfer. 2020.
- Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In 2017 IEEE symposium on security and privacy (SP), pp. 3–18. IEEE, 2017.
- Liwei Song and Prateek Mittal. Systematic evaluation of privacy risks of machine learning models.
   2020.
- Lichao Sun, Yingtong Dou, Carl Yang, Kai Zhang, Ji Wang, Philip S. Yu, Lifang He, and Bo Li.
  Adversarial attack and defense on graph data: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 35(8):7693–7711, 2023. doi: 10.1109/TKDE.2022.3201243.
- Kinyu Tang, Saeed Mahloujifar, Liwei Song, Virat Shejwalkar, Milad Nasr, Amir Houmansadr, and
   Prateek Mittal. Mitigating membership inference attacks by {Self-Distillation} through a novel
   ensemble architecture. In *31st USENIX Security Symposium (USENIX Security 22)*, pp. 1433–1450, 2022.
- Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, Yoshua Ben gio, et al. Graph attention networks. *stat*, 1050(20):10–48550, 2017.
- Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In
   *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7794–7803, 2018.
- Kiuling Wang and Wendy Hui Wang. Link membership inference attacks against unsupervised graph representation learning. In *Proceedings of the 39th Annual Computer Security Applications Conference*, pp. 477–491, 2023.
- Xiuling Wang and Wendy Hui Wang. Subgraph structure membership inference attacks against
   graph neural networks. *Proceedings on Privacy Enhancing Technologies*, 2024.
- Yijue Wang, Chenghong Wang, Zigeng Wang, Shanglin Zhou, Hang Liu, Jinbo Bi, Caiwen Ding, and Sanguthevar Rajasekaran. Against membership inference attack: Pruning is all you need. *arXiv preprint arXiv:2008.13578*, 2020.
- Bang Wu, Xiangwen Yang, Shirui Pan, and Xingliang Yuan. Adapting membership inference attacks
   to gnn for graph classification: Approaches and implications. In 2021 IEEE International Conference on Data Mining (ICDM), pp. 1421–1426, 2021a. doi: 10.1109/ICDM51629.2021.00182.
- Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying graph convolutional networks. In *International conference on machine learning*, pp. 6861–6871. PMLR, 2019.
  - Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems*, 32(1):4–24, 2020.
- Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu. A
   comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 32(1):4–24, 2021b. doi: 10.1109/TNNLS.2020.2978386.
- Renchi Yang, Jieming Shi, Xiaokui Xiao, Yin Yang, Sourav S. Bhowmick, and Juncheng Liu.
   Pane: scalable and effective attributed network embedding. *The VLDB Journal*, 32(6):1237–1262,
   March 2023. ISSN 0949-877X. doi: 10.1007/s00778-023-00790-4. URL http://dx.doi.org/10.1007/s00778-023-00790-4.
- Ziqi Yang, Bin Shao, Bohan Xuan, Ee-Chien Chang, and Fan Zhang. Defending model inversion and membership inference attacks via prediction purification. *arXiv preprint arXiv:2005.03915*, 2020.

702 703 704	Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learn- ing: Analyzing the connection to overfitting. In 2018 IEEE 31st computer security foundations symposium (CSF), pp. 268–282. IEEE, 2018.
705 706 707 708	Da Yu, Huishuai Zhang, Wei Chen, Jian Yin, and Tie-Yan Liu. How does data augmentation affect privacy in machine learning? In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pp. 10746–10753, 2021.
709 710	Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. <i>Advances in neural information processing systems</i> , 30, 2017.
711 712 713 714	He Zhang, Bang Wu, Xingliang Yuan, Shirui Pan, Hanghang Tong, and Jian Pei. Trustworthy graph neural networks: Aspects, methods, and trends. <i>Proceedings of the IEEE</i> , 112(2):97–139, 2024. doi: 10.1109/JPROC.2024.3369017.
715 716 717	Zhikun Zhang, Min Chen, Michael Backes, Yun Shen, and Yang Zhang. Inference attacks against graph neural networks. In <i>31st USENIX Security Symposium (USENIX Security 22)</i> , pp. 4543–4560, 2022.
718 719 720 721 722	Rong Zhu, Kun Zhao, Hongxia Yang, Wei Lin, Chang Zhou, Baole Ai, Yong Li, and Jingren Zhou. Aligraph: a comprehensive graph neural network platform. <i>Proc. VLDB Endow.</i> , 12 (12):2094–2105, aug 2019. ISSN 2150-8097. doi: 10.14778/3352063.3352127. URL https://doi.org/10.14778/3352063.3352127.
723	
724	
725	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	

#### 756 EXTENDED RELATED WORKS А

758

Membership Inference Attacks. MIA on ML models aim to infer whether a data record was used to train a target ML model or not. This concept is firstly proposed by Homer et al. (2008) and later 760 on extended to various directions, ranging from white-box setting where the whole target model is 761 released (Nasr et al., 2019; Rezaei & Liu, 2021; Melis et al., 2019; Leino & Fredrikson, 2020), to 762 black-box setting where only (partial of) output predictions are accessible to the adversary (Shokri et al., 2017; Salem et al., 2018; Song & Mittal, 2020; Li & Zhang, 2021; Choquette-Choo et al., 763 764 2021; Carlini et al., 2022). As a general guideline for MIA, the attacker first need to determine the most informative features that distinguish the sample membership. This feature can be posterior 765 predictions (Shokri et al., 2017; Salem et al., 2018; Jia et al., 2019), loss values (Yeom et al., 2018; 766 Sablayrolles et al., 2019), or gradient norms (Nasr et al., 2019; Rezaei & Liu, 2021). Upon identify-767 ing the informative features, the attacker can choose to learn either a binary classifier (Shokri et al., 768 2017) or metric-based decisions (Yeom et al., 2018; Salem et al., 2018) from shadow model trained 769 on shadow dataset to extract patterns in these features among the training samples for identifying 770 membership. The shadow dataset can be either generated from target model inferences, or a noisy 771 version of the original dataset depending on the assumptions of the attacker. 772

Defense Against Membership Inference Attacks. As MIA exploit the behavioral differences of 773 the target model on trainset and testset, most defense mechanisms work towards suppressing the 774 common patterns that an optimal attack relies on. Popular defense methods include confidence 775 score masking, regularization, knowledge distillation, and differential privacy. Confidence score 776 masking aims to hide the true prediction vector returned by the target model and thus mitigates the 777 effectiveness of MIAs, including only providing top-k logits per inference (Shokri et al., 2017), or 778 add noise to the prediction vector in an adversarial manner (Jia et al., 2019). Regularization aims 779 to reduce the overfitting degree of target models to mitigate MIAs. Existing regularization methods including  $L_2$ -norm regularization (Choquette-Choo et al., 2021; Hayes et al., 2017), dropout (Leino 780 & Fredrikson, 2020; Salem et al., 2018), data argumentation (Kaya & Dumitras, 2021; Yu et al., 781 2021), model compression (Wang et al., 2020), and label smoothing (Chen et al., 2022). Knowledge 782 distillation aims to transfer the knowledge from a unprotected model to a protected model (She-783 jwalkar & Houmansadr, 2020), and differential privacy (Saeidian et al., 2021) naturally protects the 784 membership information with theoretical guarantees at the cost of lower model utility. 785

786 787

788

#### В STANDARD SMIA PROCESS

789 The Standard SMIA contains three phases: shadow GNN model training, attack model training, 790 and membership inference. (1) shadow GNN model training: The adversary trains a set of shadow models  $\Phi_1^s$  to  $\Phi_2^r$  on shadow graphs to mimic the behaviors of the target model. In fact, training a 791 single shadow model is sufficient to achieve performance comparable to training T models. To train  $\Phi_1^s$  to  $\Phi_T^s$ , the attackers will randomly sample the training dataset  $\mathcal{G}_i^s$  for each shadow model from shadow dataset  $\mathcal{G}_s$ . Then the attackers train  $\Phi_1^s$  to  $\Phi_T^s$  by using  $(X_i^s, A_i^s, Y_i^{\text{Train, s}}, \mathcal{V}_i^{\text{Train, s}})$ . (2) attack 793 794 model training: The attackers firstly generate a training dataset  $\mathcal{A}^{\text{Train}}$  from the output of the shadow models on the shadow graph. The generation of  $\mathcal{A}^{\text{Train}}$  follows four steps. First, for each shadow 796 graph sample  $\mathcal{G}_i^s$ , the attackers select a set of k-node sets  $\mathcal{V}_{att}^i$ . Each node set  $\mathcal{V}_{att} \in \mathcal{V}_{att}^i$  comprises 797 k nodes randomly selected from  $\mathcal{G}_i^s$ . Second, for each node set  $\mathcal{V}_{att} \in \mathcal{V}_{att}^i$ , the attackers obtain the 798 posterior probability vector for each node in  $\mathcal{V}_{att}$  output by the shadow model  $\Phi_i^s$ , so there will be k 799 posterior probability vectors for  $V_{att}$ . Third, these k posterior probability vectors are aggregated into 800 a single vector which will serve as the attack feature vector x. Specifically, the attackers measure the 801 pairwise similarity of the k posterior probability vectors and obtains  $\binom{k}{2}$  pairwise similarity values 802 accordingly. Next, the attackers sort these similarity values in ascending order, and concatenates the 803 sorted values as a vector. Following the the same choice in Wang & Wang (2024), we consider three 804 similarity metrics, namely, dot product, cosine similarity, and euclidean distance. Therefore, there 805 will be 3 sorted vectors accordingly, where each vector corresponds to a similarity metric. These 806 3 vectors are further concatenated as one vector, acting as the attack feature x. After the attackers 807 generate the feature x of the node set  $\mathcal{V}_{att}$ , they associates x with its label y. In particular, y = 1if  $\mathcal{V}_{att}$  forms a k-clique in the  $\mathcal{G}_i^s$ , y = 2 if  $\mathcal{V}_{att}$  contains a (k-1)-hop path, and y = 0 otherwise. Finally, the attackers add the newly formed data sample (x, y) to  $\mathcal{A}^{\text{Train}}$ . After  $\mathcal{A}^{\text{Train}}$  is generated, the 808 attackers proceed to train the attack classifier A on  $\mathcal{A}^{\text{Train}}$ . (3) membership inference: The attackers employ the same methodology as the generation of training dataset  $\mathcal{A}^{\text{Train}}$  to derive the feature  $x_a$ for the target node set  $\mathcal{V}_a$ , utilizing the same similarity functions. It is important to note that, unlike the attack features of the training data  $\mathcal{V}_{att}$  that uses the probability output by the shadow model, the attackers employ the posterior probability output of  $\mathcal{V}_a$  by the target model to calculate  $x_a$ . Finally, the attackers feed  $x_a$  into A to obtain predictions. The complete standard SMIA process can be found in Wang & Wang (2024).

816 817

818

831

832

### C STANDARD MIA PROCESS

819 The standard MIA process also has three phases: shadow GNN model training, attack model train-820 ing, and membership inference. (1) shadow GNN model training: shadow GNN model S is a model 821 trained by attackers to replicate the behavior of the target GNN model M, providing training data 822 for the attack model A. To train S, we assume that the shadow dataset  $\mathcal{G}_s$  comes from the same or similar underlying distribution as  $\mathcal{G}_t$ . Then the attackers train S by using  $(X_s, A_s, Y_s^{\text{Train}}, \mathcal{V}_s^{\text{Train}})$  (2) 823 attack model training: To train A, attackers use the trained S to predict all nodes in  $\mathcal{V}_s^{\text{Train}}$  and  $\mathcal{V}_s^{\text{Test}}$ 824 and obtain the corresponding posteriors. For each node, attackers take its posteriors as input of the 825 attack model and assigns a label "1" if the node is from  $\mathcal{V}_s^{\text{train}}$  and "0" if the node is from  $\mathcal{V}_s^{\text{test}}$  to 826 supervise. (3) membership inference: To implement membership inference attack on a given node 827 v, attackers query M with v's feature to obtain its posterior. Then attackers input the posterior into 828 the attack model to obtain the membership information. The complete standard MIA process can be 829 found in Olatunji et al. (2021). 830

D COMPARISON WITH LBP AND DMP

833 Compared to state-of-the-art defense methods based on perturbations and distillation, such as 834 LBP (Olatunji et al., 2021) for GNNs and DMP (Shejwalkar & Houmansadr, 2020) for graphless 835 models, TSD can achieve a better balance between model utility and defense performance. LBP 836 employs noise addition to the posteriors of the target model, grouping the elements randomly and 837 adding noise from the same Laplace distribution to each group to reduce the required amount of 838 noise. While LBP offers strong defense capabilities, the added noise significantly degrades the tar-839 get model utility. DMP, on the other hand, tunes the data used for knowledge transfer to enhance 840 membership privacy. It utilizes an unprotected model trained on private data to guide the training of 841 a protected target model on reference data, optimizing the tradeoff between membership privacy and 842 utility. However, DMP necessitates the collection of an additional dataset for training the protected 843 model, which complicates the whole process. Meanwhile, our method addresses the overfitting 844 problem implicitly by ensuring that both the training set and testset undergo the same procedure. Consequently, our method offers several advantages over LBP and DMP: (1) it avoids explicitly 845 adding noise to the target model predictions, thereby preserving model utility; (2) it does not require 846 additional data; and (3) it fully leverages the whole graph data through a train-test alternate training 847 schedule. 848

849

### E DETAILS OF DMP LOSS FUNCTION

850 851

In our experiments, the post-distillation phase of DMP consists of two parts of loss to train the protected model, with the proportion adjusted by a hyperparameter. One loss is the cross-entropy loss, supervised by the true labels of the reference data. The other loss is the KL divergence between the prediction of the protected model and the unprotected model on the reference data. The former is to ensure that the protected model has a high classification accuracy on the testset, while the latter is to guide the protected model by using the knowledge from the unprotected model. In our experiments, we adjust the hyperparameters to balance the testset classification accuracy and defense capability of the protected model.

859 860 861

862

### F COMPLETE EXPERIMENTAL SETTINGS

Appendix Table 6 contains properties and statistics about benchmark datasets we used in SMIA experiments. For target models and shadow models, we used 2-layer GCN, 2-layer GAT, 2-layer

SGC, NLGCN, NLGAT, NLMLP, and GPRGNN architecture for CiteSeer, lastFM, Chameleon and
3-layer GCN, 3-layer GAT, 3-layer SGC, GPRGNN for Facebook. The attack model is a 3-layer
MLP model. The optimizer we used is Adam. All target and shadow models are trained such
that they achieve comparable performance as reported by the authors in the literature. We used one
NVIDIA GeForce RTX 3090 for training. The time for finishing one experiment is about 10 minutes
to 30 minutes depends on the complexity of datasets.

Table 6: Benchmark dataset properties and statistics of SMIA experiments:  $|\mathcal{V}|$  and  $|\mathcal{E}|$  denote the number of vertices and edges in the corresponding graph dataset.

Dataset	$ \mathcal{V} $	$ \mathcal{E} $	Features	Classes
CiteSeer	3327	4552	3703	6
Facebook	4039	88234	1283	193
LastFM	7624	55612	128	18
Chameleon	2277	31371	2325	5

Table 7: Benchmark dataset properties and statistics of MIA experiments:  $|\mathcal{V}|$  and  $|\mathcal{E}|$  denote the number of vertices and edges in the corresponding graph dataset.

Dataset	$ \mathcal{V} $	$ \mathcal{E} $	Features	Classes
PubMed	19717	44324	500	5
Computers	13752	245861	767	10
Photo	7650	119081	745	8
Ogbn-Arxiv	169343	1166243	128	40
Texas	183	279	1703	5
Squirrel	5201	198353	2089	5

Appendix Table 7 provides additional information on the properties and statistics of the datasets used in the MIA experiments. For target models and shadow models, we used 2-layer GCN, 2-layer GAT, 2-layer SGC, NLGCN, NLGAT, NLMLP, and GPRGNN architecture for PubMed, Computers, Photo, lastFM, Texas, Chameleon, Squirrel and 3-layer GCN, 3-layer GAT, 3-layer SGC, GPRGNN for Facebook and Ogbn-Arxiv. The attack model is a 3-layer MLP model. The optimizer we used is Adam. All target and shadow models are trained such that they achieve comparable performance as reported by the authors in the literature. We used one NVIDIA GeForce RTX 3090 for training. The time for finishing one experiment is about 10 minutes to 5 hours depends on the complexity of datasets. 

### G COMPLETE RESULTS OF SMIA DEFENSE COMPARISON AND SMIA ATTACK COMPARISON

Appendix Table 8 and 9 show the complete results including standard deviation.

### H ADAPTIVE SMIA ATTACK

Adaptive SMIA attack refers to scenarios where the attacker is aware that the target model employs
 the TSD defense method. We consider two scenarios of adaptive attack.

913 The first scenario of an adaptive attack occurs when the attacker knows that the target model uti-914 lizes the TSD defense method and is familiar with the target model's architecture. Additionally, 915 the attacker coincidentally selects the same dataset to train a shadow model. However, they are 916 completely unaware of how the training and testing sets are partitioned within TSD, making the 917 thresholds used for Flattening and the test set in the two-stage training unknown. In this case, it is 918 difficult for the attacker to adjust their attack strategy to better counter the TSD method, even if they

Table 8: 3-SMIA attack performance comparison between Standard-SMIA and End2end-SMIA &
Defense performance comparison between SHNA and TSD. Compared to SHNA, TSD achieves a
decrease in attack AUROC by 16.64% against Standard-SMIA & 14.30% against End2end-SMIA,
and increases in utility performance by 10.05% on average.

Dataset	Models	CA(SHNA)	CA(TSD)	S-SMIA AU (SHNA)	S-SMIA AU (TSD)	E-SMIA AU(SHNA)	E-SMIA AU(TSD)
CiteSeer	GCN GAT SGC GPRGNN		$\begin{array}{c} 0.7729 \pm 0.0080 \\ 0.7548 \pm 0.0092 \\ 0.7454 \pm 0.0084 \\ 0.7864 \pm 0.0079 \end{array}$		$\begin{array}{c} 0.7767 \pm 0.0080 \\ 0.7145 \pm 0.0075 \\ 0.7978 \pm 0.0082 \\ 0.5969 \pm 0.0066 \end{array}$	$\begin{array}{c} 0.9684 \pm 0.0084 \\ 0.9336 \pm 0.0093 \\ 0.9862 \pm 0.0089 \\ 0.9006 \pm 0.0078 \end{array}$	$\begin{array}{c} 0.8703 \pm 0.0075 \\ 0.8519 \pm 0.0094 \\ 0.8904 \pm 0.0103 \\ 0.7364 \pm 0.0084 \end{array}$
Facebook	GCN GAT SGC GPRGNN	$\begin{array}{c} 0.6028 \pm 0.0076 \\ 0.5321 \pm 0.0073 \\ 0.5011 \pm 0.0075 \\ 0.5529 \pm 0.0070 \end{array}$	$\begin{array}{c} 0.7095 \pm 0.0081 \\ 0.6435 \pm 0.0098 \\ 0.6423 \pm 0.0091 \\ 0.6544 \pm 0.0079 \end{array}$	$\begin{array}{c} 0.7143 \pm 0.0061 \\ 0.6854 \pm 0.0057 \\ 0.7028 \pm 0.0052 \\ 0.6512 \pm 0.0069 \end{array}$	$\begin{array}{c} 0.5463 \pm 0.0064 \\ 0.5382 \pm 0.0066 \\ 0.5253 \pm 0.0058 \\ 0.5562 \pm 0.0072 \end{array}$	$\begin{array}{c} 0.7942 \pm 0.0086 \\ 0.8229 \pm 0.0085 \\ 0.7903 \pm 0.0076 \\ 0.8081 \pm 0.0079 \end{array}$	$\begin{array}{c} 0.6127 \pm 0.0096 \\ 0.6068 \pm 0.0084 \\ 0.5945 \pm 0.0129 \\ 0.6213 \pm 0.0072 \end{array}$
LastFM	GCN GAT SGC GPRGNN		$\begin{array}{c} 0.8523 \pm 0.0086 \\ 0.8662 \pm 0.0077 \\ 0.8400 \pm 0.0094 \\ 0.8577 \pm 0.0078 \end{array}$	$\begin{array}{c} 0.9215 \pm 0.0088 \\ 0.8528 \pm 0.0072 \\ 0.9305 \pm 0.0089 \\ 0.9004 \pm 0.0093 \end{array}$	$\begin{array}{c} 0.8337 \pm 0.0068 \\ 0.8664 \pm 0.0087 \\ 0.8835 \pm 0.0065 \\ 0.8268 \pm 0.0073 \end{array}$	$\begin{array}{c} 0.9825 \pm 0.0078 \\ 0.9274 \pm 0.0068 \\ 0.9840 \pm 0.0087 \\ 0.9543 \pm 0.0074 \end{array}$	$\begin{array}{c} 0.8702 \pm 0.0088 \\ 0.8821 \pm 0.0075 \\ 0.9060 \pm 0.0109 \\ 0.8687 \pm 0.0069 \end{array}$
Chameleon	NLGCN NLGAT NLMLP GPRGNN		$\begin{array}{c} 0.6725 \pm 0.0091 \\ 0.6835 \pm 0.0078 \\ 0.5484 \pm 0.0096 \\ 0.6835 \pm 0.0085 \end{array}$	$\begin{array}{c} 0.8904 \pm 0.0076 \\ 0.7025 \pm 0.0077 \\ 0.7128 \pm 0.0096 \\ 0.7191 \pm 0.0087 \end{array}$	$\begin{array}{c} 0.4538 \pm 0.0083 \\ 0.5611 \pm 0.0079 \\ 0.5342 \pm 0.0091 \\ 0.6259 \pm 0.0069 \end{array}$	$\begin{array}{c} 0.9196 \pm 0.0074 \\ 0.8848 \pm 0.0082 \\ 0.8152 \pm 0.0091 \\ 0.8929 \pm 0.0092 \end{array}$	$\begin{array}{c} 0.7875 \pm 0.0098 \\ 0.7958 \pm 0.0074 \\ 0.7057 \pm 0.0130 \\ 0.7711 \pm 0.0097 \end{array}$
DBLP	GCN GAT SAGE	$\begin{array}{c} 0.6201 \pm 0.0131 \\ 0.6146 \pm 0.0145 \\ 0.6254 \pm 0.0116 \end{array}$	$\begin{array}{c} 0.7652 \pm 0.0117 \\ 0.7721 \pm 0.0135 \\ 0.7486 \pm 0.0131 \end{array}$	$\begin{array}{c} 0.5462 \pm 0.0153 \\ 0.5631 \pm 0.0114 \\ 0.5398 \pm 0.0151 \end{array}$	$\begin{array}{c} 0.5312 \pm 0.0134 \\ 0.5116 \pm 0.0137 \\ 0.5267 \pm 0.0115 \end{array}$	$\begin{array}{c} 0.7124 \pm 0.0140 \\ 0.6921 \pm 0.0124 \\ 0.6691 \pm 0.0134 \end{array}$	$\begin{array}{c} 0.6544 \pm 0.0128 \\ 0.5974 \pm 0.0134 \\ 0.6235 \pm 0.0122 \end{array}$
IMDB	GCN GAT SAGE	$\begin{array}{c} 0.4366 \pm 0.0125 \\ 0.4297 \pm 0.0143 \\ 0.4172 \pm 0.0134 \end{array}$	$\begin{array}{c} 0.5167 \pm 0.0131 \\ 0.5283 \pm 0.0121 \\ 0.5244 \pm 0.0135 \end{array}$	$\begin{array}{c} 0.5132 \pm 0.0134 \\ 0.5094 \pm 0.0137 \\ 0.5068 \pm 0.0150 \end{array}$	$\begin{array}{c} 0.5107 \pm 0.0120 \\ 0.5026 \pm 0.0135 \\ 0.5024 \pm 0.0127 \end{array}$	$\begin{array}{c} 0.6518 \pm 0.0153 \\ 0.6324 \pm 0.0129 \\ 0.6063 \pm 0.0128 \end{array}$	$\begin{array}{c} 0.5811 \pm 0.0113 \\ 0.5637 \pm 0.0111 \\ 0.5644 \pm 0.0124 \end{array}$

Table 9: 4-SMIA attack performance comparison between Standard-SMIA and End2end-SMIA & Defense performance comparison between SHNA and TSD. Compared to SHNA, TSD achieves a decrease in attack AUROC by 16.56% against Standard-SMIA & 13.89% against End2end-SMIA, and increases in utility performance by 10.05% on average.

Dataset	Models	CA(SHNA)	CA(TSD)	S-SMIA AU(SHNA)	S-SMIA AU(TSD)	E-SMIA AU(SHNA)	E-SMIA AU(TSD)
CiteSeer	GCN GAT SGC GPRGNN	$\begin{array}{c} 0.6887 \pm 0.0066 \\ 0.7082 \pm 0.0063 \\ 0.6982 \pm 0.0077 \\ 0.6736 \pm 0.0069 \end{array}$	$\begin{array}{c} 0.7729 \pm 0.0080 \\ 0.7548 \pm 0.0092 \\ 0.7454 \pm 0.0084 \\ 0.7864 \pm 0.0079 \end{array}$	$\begin{array}{c} 0.9685 \pm 0.0079 \\ 0.9718 \pm 0.0080 \\ 0.9839 \pm 0.0082 \\ 0.9010 \pm 0.0077 \end{array}$	$\begin{array}{c} 0.9013 \pm 0.0076 \\ 0.8868 \pm 0.0079 \\ 0.9316 \pm 0.0080 \\ 0.8523 \pm 0.0072 \end{array}$	$\begin{array}{c} 0.9786 \pm 0.0081 \\ 0.9891 \pm 0.0073 \\ 0.9955 \pm 0.0068 \\ 0.9692 \pm 0.0091 \end{array}$	$\begin{array}{c} 0.9267 \pm 0.0088 \\ 0.9191 \pm 0.0079 \\ 0.9442 \pm 0.0074 \\ 0.8894 \pm 0.0094 \end{array}$
Facebook	GCN GAT SGC GPRGNN	$\begin{array}{c} 0.6028 \pm 0.0076 \\ 0.5321 \pm 0.0073 \\ 0.5011 \pm 0.0075 \\ 0.5529 \pm 0.0070 \end{array}$	$\begin{array}{c} 0.7095 \pm 0.0081 \\ 0.6435 \pm 0.0098 \\ 0.6423 \pm 0.0091 \\ 0.6544 \pm 0.0079 \end{array}$	$\begin{array}{c} 0.7377 \pm 0.0080 \\ 0.6732 \pm 0.0078 \\ 0.6992 \pm 0.0086 \\ 0.6520 \pm 0.0075 \end{array}$	$\begin{array}{c} 0.5308 \pm 0.0082 \\ 0.5027 \pm 0.0074 \\ 0.5097 \pm 0.0089 \\ 0.4978 \pm 0.0078 \end{array}$	$\begin{array}{c} 0.7964 \pm 0.0069 \\ 0.7436 \pm 0.0085 \\ 0.7520 \pm 0.0073 \\ 0.7090 \pm 0.0076 \end{array}$	$\begin{array}{c} 0.6033 \pm 0.0074 \\ 0.5978 \pm 0.0077 \\ 0.6325 \pm 0.0073 \\ 0.5885 \pm 0.0089 \end{array}$
LastFM	GCN GAT SGC GPRGNN	$\begin{array}{c} 0.8256 \pm 0.0072 \\ 0.8379 \pm 0.0067 \\ 0.8154 \pm 0.0064 \\ 0.8320 \pm 0.0071 \end{array}$	$\begin{array}{c} 0.8523 \pm 0.0086 \\ 0.8662 \pm 0.0077 \\ 0.8401 \pm 0.0094 \\ 0.8577 \pm 0.0078 \end{array}$	$\begin{array}{c} 0.9649 \pm 0.0094 \\ 0.8936 \pm 0.0086 \\ 0.9574 \pm 0.0079 \\ 0.9144 \pm 0.0089 \end{array}$	$\begin{array}{c} 0.8296 \pm 0.0071 \\ 0.7442 \pm 0.0082 \\ 0.8184 \pm 0.0070 \\ 0.8044 \pm 0.0094 \end{array}$	$\begin{array}{c} 0.9972 \pm 0.0094 \\ 0.9354 \pm 0.0087 \\ 0.9979 \pm 0.0092 \\ 0.9844 \pm 0.0076 \end{array}$	$\begin{array}{c} 0.8561 \pm 0.0086 \\ 0.7654 \pm 0.0080 \\ 0.8270 \pm 0.0081 \\ 0.8444 \pm 0.0097 \end{array}$
Chameleon	NLGCN NLGAT NLMLP GPRGNN	$\begin{array}{c} 0.6373 \pm 0.0071 \\ 0.6437 \pm 0.0083 \\ 0.5010 \pm 0.0069 \\ 0.6637 \pm 0.0059 \end{array}$	$\begin{array}{c} 0.6725 \pm 0.0091 \\ 0.6835 \pm 0.0078 \\ 0.5484 \pm 0.0096 \\ 0.6835 \pm 0.0085 \end{array}$	$\begin{array}{c} 0.7663 \pm 0.0074 \\ 0.7849 \pm 0.0080 \\ 0.7531 \pm 0.0070 \\ 0.8358 \pm 0.0075 \end{array}$	$\begin{array}{c} 0.5668 \pm 0.0079 \\ 0.5901 \pm 0.0087 \\ 0.6076 \pm 0.0082 \\ 0.5924 \pm 0.0089 \end{array}$	$\begin{array}{c} 0.9952 \pm 0.0086 \\ 0.9816 \pm 0.0091 \\ 0.9186 \pm 0.0076 \\ 0.9486 \pm 0.0094 \end{array}$	$\begin{array}{c} 0.8346 \pm 0.0079 \\ 0.8409 \pm 0.0063 \\ 0.8359 \pm 0.0078 \\ 0.8127 \pm 0.0088 \end{array}$

are aware that the defense mechanism is TSD. If the attacker attempts to better simulate the target model by training the shadow model using the same approach as TSD, any significant differences in the partitioning of training and testing sets between the shadow and target models would actually weaken the effectiveness of the attack. Under such circumstances, the best strategy for the attacker is to train the shadow model in a standard manner without applying TSD.

The second scenario occurs when the attacker, in addition to the first case, also knows the target model's method for partitioning the dataset (though they still do not know the actual dataset used by the target model and have merely coincidentally chosen the same one). In this situation, if the attacker trains the shadow model in a standard manner, the results will be consistent with those in Tables 1 and 2, since the experimental setup in Section 6.1.1 aligns with the second scenario of the adaptive attack. However, in the case where the attacker employs the same training strategy as TSD, the attack will become stronger. Therefore, we conducted experiments under these circumstances to better test TSD defense method, and the experimental results are presented in Table10. 

971 In experiments, we used both Standard-SMIA and End2end-SMIA. The experimental results indicate that compared to standard attacks, adaptive attacks indeed achieve higher AUROC and present

Table 10: The performance of TSD under Adaptive 3-SMIA attack. *no d* indicates that no defense
method was employed, *no a* indicates that adaptive attacks were not used when the TSD defense
was applied, *a* indicates that adaptive attacks were employed when the TSD defense was applied.

Dataset   Models	S-SMIA AU(no d)	S-SMIA AU (no a)	S-SMIA AU (a)   E-SMIA AU (no d)	E-SMIA AU (no a)	E-SMIA AU (a)
CiteSeer   GCN	$0.9587 \pm 0.0078$	$0.7767 \pm 0.0086$	$0.8316 \pm 0.0074 ~\big ~ 0.9821 \pm 0.0069$	$0.8703 \pm 0.0091$	$0.9245 \pm 0.0086$
Facebook   GCN	$0.8362 \pm 0.0074$	$0.5463 \pm 0.0079$	$0.6581 \pm 0.0082 \ \big  \ 0.9014 \pm 0.0077$	$0.6127 \pm 0.0084$	$0.7061 \pm 0.0081$
LastFM   GCN	$0.9728 \pm 0.0076$	$0.8337 \pm 0.0077$	$0.8986 \pm 0.0078 \ \big  \ 0.9931 \pm 0.0082$	$0.8702 \pm 0.0094$	$0.9211 \pm 0.0080$
Chameleon   NLGC	N 0.9267 $\pm$ 0.0073	$0.4538 \pm 0.0070$	$0.6358 \pm 0.0081 \   \ 0.9747 \pm 0.0088$	$0.7875 \pm 0.0085$	$0.8553 \pm 0.0079$

greater challenges for defense mechanisms. However, when comparing the AUROC of attacks without any defense to that of adaptive attacks employing the TSD defense method, it is evident that TSD still maintains a significant defensive ability. Therefore, TSD is a highly effective defense method against SMIA attacks.

### I COMPLETE ABLATION STUDY RESULTS OF TSD'S DEFENSE PERFORMANCE AGAINST SMIA

Appendix Table 11 and 12 shows the complete ablation study results of TSD's defense performance against SMIA. The experimental results show that the same conclusions as in Section 6.1.3 can be drawn across all datasets and GNN architectures. Additionally, it is worth noting the performance of MLP under Standard Training. MLP exhibits strong defense capabilities but poor classification performance. This aligns with intuition, as MLP does not use message passing and is therefore completely unaware of the structural information between nodes.

999 1000 1001

984

985

986

987 988 989

990

991 992 993

994

995

996

997

998

### J COMPLETE EXPERIMENTS OF MIA DEFENSE

1002 1003

1005

### 1004 J.1 COMPARISON WITH LBP

The parameters we used for LBP is shown in Appendix Table 13. For each experiment, we repeated 5 times and presented the mean and standard deviation of the results in Appendix Table 14.

1008 Our analysis corresponding to different datasets is as follows: For PubMed, Computers, Photo, 1009 Facebook, LastFm and Ogbn-Arxiv, TSD achieves much better classify performance. However, 1010 there is a slight improvement in defense capability. This is because the average degree of nodes in these datasets is relatively large, and similar nodes tend to cluster in greater numbers. The target 1011 model can learn classification capabilities through a large number of similar node features, leading 1012 to more severe overfitting on the testsets, making the attack model more dangerous. Although LBP 1013 defense method can also achieve decent defense capability, it comes at the cost of significant loss in 1014 target model classification capability. 1015

For Texas dataset, TSD shows significant improvements in both classification and defense capabilities. This is because the Texas dataset has a smaller number of nodes, leading to insufficient training data for the target model and severe overfitting. However, TSD converts the testset into training data for the target model, greatly enhancing the model's generalization ability and thus strengthening its defense capabilities. In contrast, the LBP defense method excessively sacrifices the model's classification ability, making it difficult to be utilized effectively.

For Chameleon and Squirrel dataset, it can be seen that even with very low model classification accuracy, the attack model can still achieve membership inference with a probability exceeding random selection. TSD demonstrates significant improvements in model classification on two datasets. We note that NLMLP's defense capability has been greatly enhanced, this is because the improvement in its generalization ability.

Method	Dataset	GNN Models	Classify Acc	Attack Al
		GCN	$0.8012 \pm 0.0093$	$0.9745 \pm$
	CitaSaan	GAT	$0.7855 \pm 0.0086$	$0.9628 \pm$
	CheSeer	SGC	$0.7708 \pm 0.0094$	$0.9743 \pm$
		GPRGNN	$0.8247 \pm 0.0079$	$0.8672 \pm$
		MLP	$0.6988 \pm 0.0095$	$0.6531 \pm$
		GCN	$0.7353 \pm 0.0085$	$0.8759 \pm$
	Faaboolt	GAT	$0.6745 \pm 0.0096$	$0.8641 \pm$
	гасевоок	SGC	$0.6721 \pm 0.0121$	$0.8528 \pm$
Standard Training		GPRGNN	$0.7001 \pm 0.0094$	$0.8740 \pm$
		MLP	$0.3502 \pm 0.0079$	$0.5746 \pm$
		GCN	$0.8823 \pm 0.0088$	$0.9841 \pm$
	LastEM	GAT	$0.8930 \pm 0.0094$	$0.9854 \pm$
	Lastrivi	SGC	$0.8726 \pm 0.0093$	$0.9832 \pm$
		GPRGNN	$0.8991 \pm 0.0081$	0.9897 ±
		MLP	$0.6808 \pm 0.0092$	$0.6539 \pm$
	Chameleon	NLGCN	$0.7052 \pm 0.0085$	$0.9664 \pm$
		NLGAT	$0.7013 \pm 0.0097$	0.9748 ±
		MLMLP	$0.5752 \pm 0.0115$	0.9086 ±
		GPRGNN	$0.7054 \pm 0.0083$	0.9701 ±
		MLP	$0.4569 \pm 0.0096$	0.5603 ±
		GCN	$0.7956 \pm 0.0082$	0.9428 ±
	CiteSeer	GAT	$0.7814 \pm 0.0074$	0.9331 ±
		SGC	$0.7657 \pm 0.0098$	0.9452 ±
		GPRGNN	$0.8194 \pm 0.0068$	$0.8351 \pm$
		GCN	$0.7291 \pm 0.0080$	$0.8249 \pm$
	Facebook	GAT	$0.6684 \pm 0.0093$	0.8121 ±
	1 decoord	SGC	$0.6667 \pm 0.0094$	0.7901 ±
Flattening (One-Stage)		GPRGNN	$0.6974 \pm 0.0083$	$0.8298 \pm$
Flattening (One-Stage)		GCN	$0.8762 \pm 0.0096$	$0.9520 \pm$
	LastFM	GAT	$0.8881 \pm 0.0091$	0.9461 ±
		SGC	$0.8679 \pm 0.0126$	0.9487 ±
		GPRGNN	$0.8924 \pm 0.0089$	0.9511 ±
		NLGCN	$0.7009 \pm 0.0094$	0.9187 ±
	Chameleon	NLGAT	$0.6961 \pm 0.00101$	0.9258 ±
		NLMLP	$0.5691 \pm 0.00134$	0.8609 ±
		GPRGNN	$0.6977 \pm 0.0096$	0.9176 ±

Table 11: Complete ablation study on the source of gains for TSD under 3-SMIA Attack.

1064

1026

1065

### 1066 1067 J.2 COMPARISON WITH DMP

For each experiment, we repeated 5 times and presented the mean and standard deviation of the results in Appendix Table 15.

Our analysis corresponding to different datasets is as follows: For all datasets, compared to the DMP, TSD has a slight lead in both classification accuracy and defense performance. The knowledge distillation of the DMP method is pronounced in guiding the protected target model, and its defense capability is comparable to the TSD. However, the DMP method still results in a reduction in the amount of training data, which still has a significant negative impact on the model's classification ability.

In the experiments, we also observed that controlling the hyperparameters that determine the proportions of the two different losses in the post distillation phase of DMP is crucial. It requires achieving a tradeoff between classification accuracy and defense capability. Adjusting these hyperparameters will increase the implementation cost of the DMP method.

Method	Dataset	GNN Mode	els	Classify Acc	Attack AUROC
	CiteSeer	GCN GAT SGC GPRGNN		$\begin{array}{c} 0.7804 \pm 0.0087 \\ 0.7627 \pm 0.0084 \\ 0.7529 \pm 0.0095 \\ 0.7935 \pm 0.0078 \end{array}$	$ \begin{vmatrix} 0.9013 \pm 0.0081 \\ 0.8832 \pm 0.0086 \\ 0.9130 \pm 0.0090 \\ 0.7628 \pm 0.0079 \end{vmatrix} $
	Facebook	GCN GAT SGC GPRGNN		$\begin{array}{c} 0.7168 \pm 0.0080 \\ 0.6567 \pm 0.0086 \\ 0.6554 \pm 0.0094 \\ 0.6639 \pm 0.0087 \end{array}$	$\begin{array}{c} 0.6644 \pm 0.0091 \\ 0.6621 \pm 0.0082 \\ 0.6510 \pm 0.0086 \\ 0.6715 \pm 0.0084 \end{array}$
Two-Stage (without Flattening)	LastFM	GCN GAT SGC GPRGNN		$\begin{array}{c} 0.8610 \pm 0.0076 \\ 0.8749 \pm 0.0072 \\ 0.8532 \pm 0.0085 \\ 0.8679 \pm 0.0081 \end{array}$	$\begin{array}{c} 0.9008 \pm 0.0091 \\ 0.9082 \pm 0.0087 \\ 0.9336 \pm 0.0104 \\ 0.9017 \pm 0.0094 \end{array}$
	Chameleon	NLGCN NLGAT NLMLP GPRGNN		$\begin{array}{c} 0.6863 \pm 0.0085 \\ 0.6944 \pm 0.0082 \\ 0.5571 \pm 0.0091 \\ 0.6948 \pm 0.0071 \end{array}$	$\begin{array}{c} 0.8324 \pm 0.0074 \\ 0.8396 \pm 0.0082 \\ 0.7689 \pm 0.0118 \\ 0.8302 \pm 0.0068 \end{array}$
	CiteSeer	GCN GAT SGC GPRGNN		$\begin{array}{c} 0.7729 \pm 0.0080 \\ 0.7548 \pm 0.0092 \\ 0.7454 \pm 0.0084 \\ 0.7864 \pm 0.0079 \end{array}$	
	Facebook	GCN GAT SGC GPRGNN GCN GAT SGC GPRGNN NLGCN NLGAT NLMLP GPRGNN		$\begin{array}{c} 0.7095 \pm 0.0081 \\ 0.6435 \pm 0.0098 \\ 0.6423 \pm 0.0091 \\ 0.6544 \pm 0.0079 \end{array}$	$\begin{array}{c} 0.6127 \pm 0.0096 \\ 0.6068 \pm 0.0084 \\ 0.5945 \pm 0.0129 \\ 0.6213 \pm 0.0072 \end{array}$
13D (1wo Stage & Flattenning)	LastFM			$\begin{array}{c} 0.8523 \pm 0.0086 \\ 0.8662 \pm 0.0077 \\ 0.8400 \pm 0.0094 \\ 0.8577 \pm 0.0078 \end{array}$	$\begin{array}{c} 0.8702 \pm 0.0088 \\ 0.8821 \pm 0.0075 \\ 0.9060 \pm 0.0109 \\ 0.8687 \pm 0.0069 \end{array}$
	Chameleon			$\begin{array}{c} 0.6725 \pm 0.0091 \\ 0.6835 \pm 0.0078 \\ 0.5484 \pm 0.0096 \\ 0.6835 \pm 0.0085 \end{array}$	$\begin{array}{c} 0.7875 \pm 0.0098 \\ 0.7958 \pm 0.0074 \\ 0.7057 \pm 0.0130 \\ 0.7711 \pm 0.0097 \end{array}$
	Table 13: Pa	arameters fo	r LB	P	
	Dataset	N	b	-	
	PubMed Compute	l 2 ers 2	1 0.2	-	
	Faceboo LastFM	ok 2 2	0.2 0.2 0.2		
	Ogbn-A Texas Chamele Squirrel	rxiv 2 2 eon 2 2	10 0.2 0.2 0.2		
.3 TSD DEFENSE METHOD F	LEDUCE THE	Generali	ZAT	- Ion Gap	
			_		

Table 12: Complete ablation study on the source of gains for TSD under 3-SMIA Attack. Continu-ation of Appendix Table 11.

In this section, we analyzed the changes in training loss and testing loss distribution before and after TSD training. Experiments are conducted on the heterophilic dataset Chameleon.

1133 Through experiments, we demonstrated that TSD can: (1) reduce the gap between the average losses of training and testing nodes, thereby alleviating overfitting; (2) increase the variance of both

1134	Table 14: Performance comparison between TSD and LBP. Compared to LBP, TSD achieves an
1135	increase in utility performance by 17.28% on average, while achieving comparable attack AUROC
1136	to LBP.

Dataset	Models	Classify Acc(LBP)	Classify Acc(TSD)	Attack AUROC(LBP)	Attack AUROC(TSD)
PubMed	GCN GAT SGC GPRGNN	$\begin{array}{c} 0.6886 \pm 0.0041 \\ 0.7631 \pm 0.0037 \\ 0.6564 \pm 0.0035 \\ 0.7843 \pm 0.0029 \end{array}$	$\begin{array}{c} 0.8381 \pm 0.0023 \\ 0.8400 \pm 0.0028 \\ 0.8080 \pm 0.0020 \\ 0.8553 \pm 0.0014 \end{array}$	$\begin{array}{c} 0.4998 \pm 0.0050 \\ 0.5021 \pm 0.0084 \\ 0.5007 \pm 0.0065 \\ 0.5003 \pm 0.0038 \end{array}$	$\begin{array}{c} 0.4990 \pm 0.0048 \\ 0.4911 \pm 0.0061 \\ 0.5005 \pm 0.0057 \\ 0.4967 \pm 0.0034 \end{array}$
Computers	GCN GAT SGC GPRGNN	$\begin{array}{c} 0.6900 \pm 0.0023 \\ 0.7414 \pm 0.0026 \\ 0.6383 \pm 0.0038 \\ 0.7203 \pm 0.0023 \end{array}$	$\begin{array}{c} 0.8818 \pm 0.0018 \\ 0.9086 \pm 0.0025 \\ 0.8310 \pm 0.0023 \\ 0.8942 \pm 0.0012 \end{array}$	$\begin{array}{c} 0.5128 \pm 0.0039 \\ 0.5165 \pm 0.0053 \\ 0.5112 \pm 0.0041 \\ 0.5154 \pm 0.0033 \end{array}$	$\begin{array}{c} 0.5068 \pm 0.0035 \\ 0.5039 \pm 0.0051 \\ 0.5074 \pm 0.0043 \\ 0.5050 \pm 0.0030 \end{array}$
Photo	GCN GAT SGC GPRGNN		$\begin{array}{c} 0.9299 \pm 0.0023 \\ 0.9453 \pm 0.0029 \\ 0.9001 \pm 0.0028 \\ 0.9430 \pm 0.0015 \end{array}$	$\begin{array}{c} 0.5179 \pm 0.0046 \\ 0.5123 \pm 0.0045 \\ 0.5159 \pm 0.0039 \\ 0.5153 \pm 0.0025 \end{array}$	$\begin{array}{c} 0.5110 \pm 0.0041 \\ 0.5066 \pm 0.0039 \\ 0.5106 \pm 0.0035 \\ 0.5088 \pm 0.0028 \end{array}$
Facebook	GCN GAT SGC GPRGNN	$\begin{array}{c} 0.5195 \pm 0.0041 \\ 0.5460 \pm 0.0037 \\ 0.4833 \pm 0.0044 \\ 0.4627 \pm 0.0032 \end{array}$	$\begin{array}{c} 0.6778 \pm 0.0049 \\ 0.6519 \pm 0.0039 \\ 0.6249 \pm 0.0043 \\ 0.5890 \pm 0.0035 \end{array}$	$\begin{array}{c} 0.4912 \pm 0.0021 \\ 0.5120 \pm 0.0025 \\ 0.4901 \pm 0.0026 \\ 0.4807 \pm 0.0031 \end{array}$	$\begin{array}{c} 0.4993 \pm 0.0023 \\ 0.5010 \pm 0.0028 \\ 0.5004 \pm 0.0031 \\ 0.5014 \pm 0.0020 \end{array}$
Lastfm	GCN GAT SGC GPRGNN		$\begin{array}{c} 0.8378 \pm 0.0035 \\ 0.8683 \pm 0.0032 \\ 0.8336 \pm 0.0045 \\ 0.8443 \pm 0.0033 \end{array}$	$\begin{array}{c} 0.5118 \pm 0.0024 \\ 0.5136 \pm 0.0031 \\ 0.5121 \pm 0.0030 \\ 0.5101 \pm 0.0035 \end{array}$	$\begin{array}{c} 0.4971 \pm 0.0022 \\ 0.4980 \pm 0.0029 \\ 0.4965 \pm 0.0034 \\ 0.4999 \pm 0.0025 \end{array}$
Ogbn-Arxiv	GCN GAT SGC GPRGNN		$\begin{array}{c} 0.7200 \pm 0.0019 \\ 0.7223 \pm 0.0024 \\ 0.7272 \pm 0.0027 \\ 0.7315 \pm 0.0016 \end{array}$	$\begin{array}{c} 0.5005 \pm 0.0037 \\ 0.5001 \pm 0.0045 \\ 0.4998 \pm 0.0040 \\ 0.5017 \pm 0.0031 \end{array}$	$\begin{array}{c} 0.4995 \pm 0.0032 \\ 0.4978 \pm 0.0038 \\ 0.4923 \pm 0.0034 \\ 0.5002 \pm 0.0026 \end{array}$
Texas	NLGCN NLGAT NLMLP GPRGNN		$\begin{array}{c} 0.6152 \pm 0.0031 \\ 0.5882 \pm 0.0026 \\ 0.6686 \pm 0.0043 \\ 0.7224 \pm 0.0029 \end{array}$	$\begin{array}{c} 0.5954 \pm 0.0061 \\ 0.5710 \pm 0.0045 \\ 0.5585 \pm 0.0035 \\ 0.6201 \pm 0.0038 \end{array}$	$\begin{array}{c} 0.4812 \pm 0.0054 \\ 0.4785 \pm 0.0043 \\ 0.5532 \pm 0.0038 \\ 0.4559 \pm 0.0032 \end{array}$
Chameleon	NLGCN NLGAT NLMLP GPRGNN		$\begin{array}{c} 0.6657 \pm 0.0062 \\ 0.6585 \pm 0.0070 \\ 0.4824 \pm 0.0074 \\ 0.6550 \pm 0.0058 \end{array}$	$\begin{array}{c} 0.5233 \pm 0.0062 \\ 0.5246 \pm 0.0065 \\ 0.5623 \pm 0.0057 \\ 0.5107 \pm 0.0060 \end{array}$	$\begin{array}{c} 0.4954 \pm 0.0065 \\ 0.4902 \pm 0.0063 \\ 0.4848 \pm 0.0051 \\ 0.4936 \pm 0.0049 \end{array}$
Squirrel	NLGCN NLGAT NLMLP GPRGNN		$\begin{array}{c} 0.4910 \pm 0.0073 \\ 0.5446 \pm 0.0077 \\ 0.3137 \pm 0.0068 \\ 0.4013 \pm 0.0060 \end{array}$	$\begin{array}{c} 0.5239 \pm 0.0059 \\ 0.5260 \pm 0.0053 \\ 0.5659 \pm 0.0058 \\ 0.5225 \pm 0.0055 \end{array}$	$\begin{array}{c} 0.4913 \pm 0.0060 \\ 0.4920 \pm 0.0049 \\ 0.4746 \pm 0.0054 \\ 0.4922 \pm 0.0050 \end{array}$

1137

1170 1171

member and non-member loss distributions and reduce the disparity between their means; and (3)decrease the distinguishability between member and non-member loss distributions.

1174 Reduce the gap between the average losses of training and testing nodes. Appendix Figure 4 1175 shows the variations of the average losses of training and testing nodes with increasing training 1176 epochs for both standard training and two-stage training on Chameleon dataset. We also recorded 1177 the losses of all models from Appendix Figure 4 at the end of training to Appendix Table 16, and 1178 additionally added the result of comparative experiments on model utility and defense capability. Comparing Appendix Figure 4 (a) with (b), it can be observed that the difference between the aver-1179 age losses of training and testing nodes in the standard training increases as epochs increase, indicat-1180 ing that overfitting exists and becomes worse as training proceeds. However, when using two-stage 1181 training, although overfitting cannot be completely avoided, the difference between training and 1182 testing losses decreases in the second stage as training proceeds, indicating a gradual alleviation of 1183 overfitting. Appendix Table 16 also shows that our method achieved lower average loss gap after the 1184 entire training process. All experimental results demonstrate the capability of our method to reduce 1185 overfitting and the generalization gap. 1186

### 1187 Increase the variance of both member and non-member loss distributions and reduce the disparity between their means. Appendix Figure 5 illustrates the loss distributions of member and

Dataset	Models	Classify Acc (DMP)	Classify Acc (TSD)	Attack AUROC (DMP)	Attack AUROC (TSD)
PubMed	GCN GAT SGC GPRGNN		$\begin{array}{c} 0.8387 \pm 0.0034 \\ 0.8434 \pm 0.0024 \\ 0.8096 \pm 0.0045 \\ 0.8423 \pm 0.0036 \end{array}$	$\begin{array}{c} 0.5026 \pm 0.0039 \\ 0.5013 \pm 0.0036 \\ 0.5024 \pm 0.0042 \\ 0.5020 \pm 0.0027 \end{array}$	$\begin{array}{c} 0.4978 \pm 0.0042 \\ 0.5005 \pm 0.0043 \\ 0.5003 \pm 0.0038 \\ 0.4994 \pm 0.0023 \end{array}$
Computers	GCN GAT SGC GPRGNN	$\begin{array}{c} 0.8756 \pm 0.0038 \\ 0.9071 \pm 0.0044 \\ 0.8323 \pm 0.0042 \\ 0.8683 \pm 0.0029 \end{array}$	$\begin{array}{c} 0.8808 \pm 0.0040 \\ 0.9127 \pm 0.0038 \\ 0.8434 \pm 0.0040 \\ 0.8898 \pm 0.0021 \end{array}$	$\begin{array}{c} 0.5055 \pm 0.0042 \\ 0.5097 \pm 0.0045 \\ 0.5086 \pm 0.0045 \\ 0.5045 \pm 0.0020 \end{array}$	$\begin{array}{c} 0.5004 \pm 0.0045 \\ 0.4933 \pm 0.0051 \\ 0.5023 \pm 0.0047 \\ 0.5036 \pm 0.0032 \end{array}$
Photo	GCN GAT SGC GPRGNN		$\begin{array}{c} 0.9304 \pm 0.0042 \\ 0.9493 \pm 0.0038 \\ 0.8988 \pm 0.0037 \\ 0.9315 \pm 0.0026 \end{array}$	$\begin{array}{c} 0.5061 \pm 0.0052 \\ 0.5072 \pm 0.0058 \\ 0.5105 \pm 0.0060 \\ 0.5043 \pm 0.0044 \end{array}$	$\begin{array}{c} 0.5004 \pm 0.0048 \\ 0.4966 \pm 0.0055 \\ 0.5018 \pm 0.0059 \\ 0.4976 \pm 0.0042 \end{array}$
Facebook	GCN GAT SGC GPRGNN	$\begin{array}{c} 0.6896 \pm 0.0046 \\ 0.6420 \pm 0.0040 \\ 0.6152 \pm 0.0049 \\ 0.5784 \pm 0.0051 \end{array}$	$\begin{array}{c} 0.7054 \pm 0.0045 \\ 0.6797 \pm 0.0043 \\ 0.6351 \pm 0.0040 \\ 0.5920 \pm 0.0046 \end{array}$		$\begin{array}{c} 0.4964 \pm 0.0030 \\ 0.4910 \pm 0.0037 \\ 0.4969 \pm 0.0039 \\ 0.4964 \pm 0.0032 \end{array}$
Lastfm	GCN GAT SGC GPRGNN		$\begin{array}{c} 0.8401 \pm 0.0055 \\ 0.8769 \pm 0.0041 \\ 0.8414 \pm 0.0044 \\ 0.8485 \pm 0.0037 \end{array}$	$\begin{array}{c} 0.5115 \pm 0.0030 \\ 0.5142 \pm 0.0028 \\ 0.5091 \pm 0.0036 \\ 0.5119 \pm 0.0026 \end{array}$	$\begin{array}{c} 0.4978 \pm 0.0027 \\ 0.4972 \pm 0.0023 \\ 0.4967 \pm 0.0031 \\ 0.4979 \pm 0.0019 \end{array}$
Ogbn-Arxiv	GCN GAT SGC GPRGNN	$\begin{array}{c} 0.6876 \pm 0.0039 \\ 0.6869 \pm 0.0037 \\ 0.6798 \pm 0.0023 \\ 0.6903 \pm 0.0025 \end{array}$	$\begin{array}{c} 0.6921 \pm 0.0020 \\ 0.6907 \pm 0.0018 \\ 0.6884 \pm 0.0023 \\ 0.6993 \pm 0.0020 \end{array}$	$\begin{array}{c} 0.4920 \pm 0.0045 \\ 0.4913 \pm 0.0042 \\ 0.4924 \pm 0.0048 \\ 0.5035 \pm 0.0037 \end{array}$	$\begin{array}{c} 0.4952 \pm 0.0049 \\ 0.4934 \pm 0.0048 \\ 0.4991 \pm 0.0043 \\ 0.4972 \pm 0.0040 \end{array}$
Texas	NLGCN NLGAT NLMLP GPRGNN	$\begin{array}{c} 0.6846 \pm 0.0069 \\ 0.6916 \pm 0.0074 \\ 0.6948 \pm 0.0070 \\ 0.6115 \pm 0.0054 \end{array}$	$\begin{array}{c} 0.7027 \pm 0.0044 \\ 0.7263 \pm 0.0038 \\ 0.7282 \pm 0.0042 \\ 0.7445 \pm 0.0031 \end{array}$	$\begin{array}{c} 0.4966 \pm 0.0060 \\ 0.4923 \pm 0.0063 \\ 0.4926 \pm 0.0068 \\ 0.5187 \pm 0.0052 \end{array}$	$\begin{array}{c} 0.4970 \pm 0.0056 \\ 0.4976 \pm 0.0058 \\ 0.4953 \pm 0.0063 \\ 0.5032 \pm 0.0046 \end{array}$
Chameleon	NLGCN NLGAT NLMLP GPRGNN		$\begin{array}{c} 0.6963 \pm 0.0065 \\ 0.7082 \pm 0.0073 \\ 0.4955 \pm 0.0070 \\ 0.6934 \pm 0.0068 \end{array}$	$\begin{array}{c} 0.5210 \pm 0.0064 \\ 0.5116 \pm 0.0061 \\ 0.5054 \pm 0.0053 \\ 0.5172 \pm 0.0067 \end{array}$	$\begin{array}{c} 0.5182 \pm 0.0062 \\ 0.5159 \pm 0.0059 \\ 0.5017 \pm 0.0056 \\ 0.5163 \pm 0.0045 \end{array}$
Squirrel	NLGCN NLGAT NLMLP GPRGNN		$\begin{array}{c} 0.5102 \pm 0.0076 \\ 0.5723 \pm 0.0071 \\ 0.3393 \pm 0.0073 \\ 0.4581 \pm 0.0068 \end{array}$		$\begin{array}{c} 0.5001 \pm 0.0058 \\ 0.4934 \pm 0.0055 \\ 0.5098 \pm 0.0050 \\ 0.5054 \pm 0.0052 \end{array}$

Table 15: Performance comparison between TSD and DMP. Compared to DMP, TSD achieves an increase in utility performance by 4.35% on average, while achieving comparable attack AUROC to DMP.

Table 16: The performance comparison between standard training and TSD training.

Datasat	Models	Avg Trai	n Loss	Avg Te	st Loss	Classif	y Acc	Attack A	UROC
Dataset	Widdels	Standard	TSD	Standard	TSD	Standard	TSD	Standard	TSD
	NLGCN	0.4988	0.6192	1.1278	1.2631	0.6631	0.6657	0.5267	0.4954
Chamalaan	NLGAT	0.1617	0.5782	1.1650	1.28911	0.6603	0.6585	0.5304	0.4902
Chameleon	NLMLP	0.1555	0.9322	3.6396	2.0010	0.4857	0.4824	0.7681	0.4848
	GPRGNN	0.3575	0.8498	1.6416	1.3626	0.6582	0.6550	0.5837	0.4936

non-member nodes on the Chameleon dataset with NLGCN after training with both standard and two-stage methods. In the figure, members refer to the nodes in the trainset of the target model, while non-members refer to the nodes in the testset. Therefore, Figure 5 can also be viewed as the training and testing loss distributions of the target model after using different training methods. Comparing Figure 5 (a) with (b), it can be observed that the loss distributions of members and non-members after standard training have relatively small variances, and their means differ significantly. This conclusion is consistent with the results of the average losses in Appendix Table 16. However, after using two-stage training, significant changes occur in the loss distributions: the variances of both two distributions increase. And combined with the results in Appendix Table 16, it is obvious that their means become closer.



Figure 4: Comparison of average training loss and testing loss on Chameleon for (a) standard training, (b) two-stage training (TSD). In (b), the left half of the orange dashed line indicates the first training stage of our method, while the one on the right indicates the second stage.



Figure 5: Loss distribution histograms for (a) standard training on Chameleon, (b) two-stage training (TSD) on Chameleon

1259

1261

1262 1263

1264 1265

1266

1267 1268

1269

1270

Decrease the distinguishability between member and non-member loss distributions. From
Figure 5, it can be seen that the overlap between the member and non-member loss distributions
of the target model after two-stage training is significantly larger than that of standard training.
Combined with the conclusions obtained above, we can confirm that the distinguishability between
member and non-member distributions has decreased, which will increase the difficulty of MIA.

In summary, the changes of the target model induced by our two-stage training method are significant. Appendix Table 16 also demonstrates that such changes not only substantially enhance defense capability but also result in only subtle decline in downstream classification accuracy.

1283 1284 J.4 DATA TRANSFER

Figure 6, 7 display the experimental results using GCN as backbones. The results demonstrate that TSD also exhibits excellent defense capability in the data transfer setting. Since model utility is only related to the target dataset, combining the classification performance of TSD in Appendix Table 14 and 15, it is evident that TSD can still achieve an outstanding balance between model utility and defense capability in the data transfer setting, which means that TSD can be effectively deployed in real-world applications.

1291

1293

J.5 ABLATION STUDY OF TSD'S DEFENSE PERFORMANCE AGAINST MIA

1294 In the experiments, we set up the same four variants as described in Section 6.1.3. Additionally, 1295 we also considered RelaxLoss (Chen et al., 2022), which is essentially a combination of alternate flattening and gradient ascent when the training loss falls below a predefined threshold. We use

<sup>1274</sup> 







Figure 7: Comparison of TSD and DMP's defense performance under data transfer setting. (a)
Attack AUROC of TSD; (b) Attack AUROC of DMP.

RelaxLoss as an example to show the difference between the defense methods effective for graph and graphless models, and necessities to design defense mechanisms specially for graph models.
Note that the split ratio of trainsets and testsets for TSD is 9:1.

Appendix Table 17, 18,19 presents the results of ablation study regarding five variants. Our analysis is as follow: We first focus on the improvement of defense capability. It can be observed that two-stage (without flattening), flattening, and gradient ascent all enhance the defense capability of the target model compared to standard training. The effect of two-stage (without flattening) on reducing the AUROC of the attack model is the most pronounced, followed by flattening, while gradient as-cent slightly reduces it. These results align with expectations because two-stage (without flattening) directly enables the model to learn the distribution of the testing data and flattening decrease the difference between the loss distributions' variances of training and testing nodes. Surprisingly, gra-dient ascent hardly improves the model's defense capability, suggesting that our method's exclusion of gradient ascent is reasonable. 

Then we focus on the decline of classification accuracy caused by these variants. From the results,
it can be seen that two-stage (without flattening) and gradient ascent hardly lead to a decrease in
classification accuracy, and flattening only results in a slight decline. These results are interpretable:
two-stage (without flattening) only used testing data for an extra training; gradient ascent has a minimal impact on the model's defense capability, which also means that it hardly change the model;

flattening slightly alters the model's mapping when using soft labels, so the classification accuracy decrease. However, the degree of decline caused by flattening is acceptable compared to its enhancement in defense capability. In summary, our TSD method (two stage with flattening) is the best.

1354 1355

Table 17: Ablation Study of TSD's Defense performance against MIA.

Method	Dataset	GNN Models	Classify Acc	Attack AURC
		GCN	$0.8423 \pm 0.0027$	$0.5306 \pm 0.0$
	DubMad	GAT	$0.8421 \pm 0.0025$	$0.5306 \pm 0.0$
	rubivieu	SGC	$0.8129 \pm 0.0024$	$0.5335 \pm 0.0$
		GPRGNN	$0.8637 \pm 0.0026$	$0.5388 \pm 0.0$
		GCN	$0.8865 \pm 0.0025$	$0.5331 \pm 0.0$
	G (	GAT	$0.9145 \pm 0.0026$	$0.5360 \pm 0.0$
	Computers	SGC	$0.8419 \pm 0.0022$	$0.5334 \pm 0.0$
		GPRGNN	$0.8917 \pm 0.0023$	$0.5377 \pm 0.0$
		GCN	$0.9349 \pm 0.0022$	$0.5381 \pm 0.0$
	-	GAT	$0.9471 \pm 0.0023$	$0.5388 \pm 0.0$
	Photo	SGC	$0.9122 \pm 0.0025$	$0.5362 \pm 0.0$
		GPRGNN	$0.9499 \pm 0.0028$	$0.5302 \pm 0.0$ 0.5390 + 0.0
Standard Training		GCN	$0.9797 \pm 0.0020$	$0.5350 \pm 0.0$ 0.5352 ± 0.0
		GAT	$0.6465 \pm 0.0034$	$0.5352 \pm 0.0$ 0.5426 + 0.0
	Facebook	SGC	$0.0405 \pm 0.0050$ $0.6275 \pm 0.0042$	$0.5420 \pm 0.0$ 0.5341 ± 0.0
		GPRGNN	$0.0275 \pm 0.0042$ 0.5635 $\pm 0.0043$	$0.5378 \pm 0.0$
		GCN	$0.3035 \pm 0.0043$	$0.5378 \pm 0.0$
	LastFM	GAT	$0.0333 \pm 0.0033$	$0.5302 \pm 0.0$
		SCC	$0.8092 \pm 0.0030$ 0.8224 $\pm 0.0020$	$0.5375 \pm 0.0$
		CDDCNN	$0.8554 \pm 0.0050$	$0.3300 \pm 0.0$
	Chameleon	OPROININ	$0.8403 \pm 0.0029$	$0.3400 \pm 0.0$
		NLGUN	$0.0817 \pm 0.0054$	$0.3469 \pm 0.0$
			$0.6451 \pm 0.0053$	$0.5790 \pm 0.0$
			$0.5077 \pm 0.0060$	$0.7868 \pm 0.0$
		GPRGININ	$0.0841 \pm 0.0037$	$0.3939 \pm 0.0$
		GCN	$0.8345 \pm 0.0025$	$0.5278 \pm 0.0$
	PubMed	GAT	$0.8328 \pm 0.0023$	$0.5213 \pm 0.0$
	1 uomuu	SGC	$0.8040 \pm 0.0023$	$0.5263 \pm 0.0$
		GPRGNN	$0.8541 \pm 0.0016$	$0.5224 \pm 0.0$
		GCN	$0.8821 \pm 0.0026$	$0.5254 \pm 0.0$
	Computers	GAT	$0.9124 \pm 0.0023$	$0.5236 \pm 0.0$
		SGC	$0.8389 \pm 0.0028$	$0.5267 \pm 0.0$
		GPRGNN	$0.8892 \pm 0.0034$	$0.5276 \pm 0.0$
		GCN	$0.9302 \pm 0.0032$	$0.5297 \pm 0.0$
	Photo	GAT	$0.9435 \pm 0.0035$	$0.5268 \pm 0.0$
	FIIOIO	SGC	$0.9087 \pm 0.0035$	$0.5275 \pm 0.0$
Elettonning (One Sterr)		GPRGNN	$0.9468 \pm 0.0033$	$0.5308 \pm 0.0$
riattenning (One-Stage)		GCN	$0.6751 \pm 0.0037$	$0.5275 \pm 0.0$
	East - 1	GAT	$0.6458 \pm 0.0033$	$0.5335 \pm 0.0$
	Facebook	SGC	$0.6236 \pm 0.0035$	$0.5263 \pm 0.0$
		GPRGNN	$0.5614 \pm 0.0040$	$0.5267 \pm 0.0$
		GCN	$0.8321 \pm 0.0037$	$0.5245 \pm 0.0$
	T I I I	GAT	$0.8659 \pm 0.0039$	$0.5255 \pm 0.0$
	LastFM	SGC	$0.8309 \pm 0.0036$	$0.5282 \pm 0.0$
		GPRGNN	$0.8430 \pm 0.0043$	$0.5325 \pm 0.0$
		NLGCN	$0.6720 \pm 0.0049$	$0.5327 \pm 0.0$
		NLGAT	$0.6364 \pm 0.0066$	$0.5685 \pm 0.0$
	Chameleon	NLMIP	$0.0004 \pm 0.0000$ 0.4954 + 0.0075	$0.5005 \pm 0.0$
		CDDCNN	$0.757 \pm 0.0073$	$0.7707 \pm 0.0$

Table 18: Ablation Study of TSD's Defense performance against MIA. Continuation of Appendix Table 17.

Method	Dataset	GNN Models	Classify Acc	Attack A
		GCN	$0.8310 \pm 0.0043$	0.5282 ±
	DubMad	GAT	$0.8430 \pm 0.0040$	$0.5201 \pm$
	Publied	SGC	$0.8190 \pm 0.0066$	$0.5289 \pm$
		GPRGNN	$0.8657 \pm 0.0041$	$0.5242 \pm$
-		GCN	$0.8798 \pm 0.0033$	$0.5248 \pm$
	~	GAT	$0.9084 \pm 0.0027$	0.5245 +
	Computers	SGC	$0.8378 \pm 0.0024$	$0.5253 \pm$
		GPRGNN	$0.8849 \pm 0.0027$	$0.5257 \pm$
r		GCN	$0.9268 \pm 0.0029$	$0.5264 \pm$
		GAT	$0.9389 \pm 0.0024$	$0.5275 \pm$
	Photo	SGC	$0.9057 \pm 0.0021$	$0.5245 \pm$
		GPRGNN	$0.9399 \pm 0.0025$	$0.5215 \pm 0.5288 \pm$
Flattenning & Gradient Asent (One-Stage)		GCN	$0.9399 \pm 0.0025$ 0.6732 ± 0.0025	$0.5200 \pm$ 0.5243 +
	Facebook	GAT	$0.6752 \pm 0.0025$ $0.6363 \pm 0.0024$	$0.5245 \pm 0.5342 +$
		SGC	$0.0303 \pm 0.0024$ $0.6180 \pm 0.0028$	$0.5342 \pm 0.5251 \pm$
		GPRGNN	$0.0100 \pm 0.0020$ $0.5591 \pm 0.0030$	$0.5251 \pm$ 0.5278 $\pm$
-		GCN	$0.3371 \pm 0.0030$ 0.8284 ± 0.0026	$0.5270 \pm$
	LastFM	GAT	$0.8284 \pm 0.0020$ 0.8592 $\pm 0.0028$	$0.5204 \pm$ 0.5235 $\pm$
		SGC	$0.8392 \pm 0.0028$ 0.8284 $\pm$ 0.0024	$0.5255 \pm 0.5264 \pm$
		GPRCNN	$0.8284 \pm 0.0024$ 0.8272 $\pm 0.0027$	$0.5204 \pm$
	Chameleon	NLCCN	$0.8372 \pm 0.0027$	$0.3342 \pm$
		NLOCN	$0.0707 \pm 0.0008$	$0.3302 \pm$
		NLGAI NLMLD	$0.0402 \pm 0.0074$ 0.4022 $\pm 0.0075$	$0.3301 \pm$
		CDDCNN	$0.4933 \pm 0.0073$	$0.7741 \pm$
		GPRGININ	$0.0733 \pm 0.0004$	0.3829 ±
		GCN	$0.8328 \pm 0.0035$	$0.5049 \pm$
	PubMed	GAT	$0.8489 \pm 0.0031$	$0.5041 \pm$
	Publied	SGC	$0.8055 \pm 0.0028$	$0.5038 \pm$
		GPRGNN	$0.8685 \pm 0.0023$	$0.5044 \pm$
		GCN	$0.8868 \pm 0.0033$	$0.5026 \pm$
	Computers	GAT	$0.9149 \pm 0.0026$	$0.5035 \pm$
		SGC	$0.8393 \pm 0.0031$	$0.5036 \pm$
		GPRGNN	$0.8921 \pm 0.0023$	$0.5045 \pm$
		GCN	$0.9326 \pm 0.0028$	$0.5046 \pm$
	Photo	GAT	$0.9464 \pm 0.0022$	$0.5041 \pm$
	1 11010	SGC	$0.9125 \pm 0.0025$	$0.5056 \pm$
Two-Stage (without Flattening)		GPRGNN	$0.9512 \pm 0.0032$	$0.5034 \pm$
Two Stage (without Flattening)		GCN	$0.6775 \pm 0.0034$	$0.5031 \pm$
	Facebook	GAT	$0.6443 \pm 0.0036$	$0.5026 \pm$
	1 accook	SGC	$0.6280 \pm 0.0037$	0.5045 ±
		GPRGNN	$0.5630 \pm 0.0029$	$0.5039 \pm$
		GCN	$0.8364 \pm 0.0031$	$0.5044 \pm$
	LastEM	GAT	$0.8689 \pm 0.0034$	0.5037 ±
		SGC	$0.8327 \pm 0.0032$	0.5018 ±
		GPRGNN	$0.8472 \pm 0.0037$	0.5040 ±
		NLGCN	$0.6745 \pm 0.0060$	0.5132 ±
	Chamalaar	NLGAT	$0.5965 \pm 0.0069$	0.4849 ±
	Chameleon	NLMLP	$0.4923 \pm 0.0077$	$0.5357 \pm$
		GPRGNN	$0.6541 \pm 0.0064$	$0.5145 \pm$

Table 19: Ablation Study of TSD's Defense performance against MIA. Continuation of Appendix Table 18. 

Method	Dataset	GNN Models	Classify Acc	Attack AUR
		GCN	$0.8381 \pm 0.0023$	$0.4997 \pm 0.1$
	PubMed	GAT	$0.8400 \pm 0.0028$	$0.4981 \pm 0.$
		SGC	$0.8080 \pm 0.0020$	$0.5005 \pm 0.5005$
		GPRGNN	$0.8553 \pm 0.0014$	$0.4987 \pm 0.12$
		GCN	$0.8845 \pm 0.0034$	$0.4996 \pm 0.12$
	Computers	GAT	$0.9122 \pm 0.0036$	$0.5005 \pm 0.5005$
	Computers	SGC	$0.8364 \pm 0.0032$	$0.4989 \pm 0$
		GPRGNN	$0.8918 \pm 0.0029$	$0.5013 \pm 0$
	Photo	GCN	$0.9301 \pm 0.0030$	$0.4987 \pm 0$
		GAT	$0.9432 \pm 0.0032$	$0.4992 \pm 0$
		SGC	$0.9089 \pm 0.0038$	$0.5008 \pm 0$
TSD (Two Stage & Flattenning)		GPRGNN	$0.9491 \pm 0.0027$	$0.5006 \pm 0$
15D (1wo Stage & Flattenning)		GCN	$0.6744 \pm 0.0034$	$0.4981 \pm 0$
	Facebook	GAT	$0.6426 \pm 0.0029$	$0.4979 \pm 0$
		SGC	$0.6256 \pm 0.0037$	$0.4990 \pm 0$
		GPRGNN	$0.5609 \pm 0.0040$	$0.4985 \pm 0$
	LastEM	GCN	$0.8340 \pm 0.0031$	$0.4994 \pm 0$
		GAT	$0.8672 \pm 0.0032$	$0.4999 \pm 0$
	Eusti IVI	SGC	$0.8325 \pm 0.0039$	$0.4983 \pm 0$
		GPRGNN	$0.8441 \pm 0.0027$	$0.5005 \pm 0$
		NLGCN	$0.6657 \pm 0.0062$	$0.4954 \pm 0$
	Chameleon	NLGAT	$0.6585 \pm 0.0070$	$0.4902 \pm 0$
	Chambreon	NLMLP	$0.4824 \pm 0.0074$	$0.5248 \pm 0$
	1	GPRGNN	$0.6550 \pm 0.0058$	$0.4956 \pm 0$
		011101111	0.0000 ± 0.0000	0