# Manipulating Visually-aware Federated Recommender Systems and Its Countermeasures

WEI YUAN, The University of Queensland, Australia

SHILONG YUAN, Nanjing University, China

CHAOQUN YANG, Griffith University, Australia

QUOC VIET HUNG NGUYEN, Griffith University, Australia

HONGZHI YIN*, The University of Queensland, Australia

Federated recommender systems (FedRecs) have been widely explored recently due to their ability to protect user data privacy. In FedRecs, a central server collaboratively learns recommendation models by sharing model public parameters with clients, thereby offering a privacy-preserving solution. Unfortunately, the exposure of model parameters leaves a backdoor for adversaries to manipulate FedRecs. Existing works about FedRec security already reveal that items can easily be promoted by malicious users via model poisoning attacks, but all of them mainly focus on FedRecs with only collaborative information (i.e., user-item interactions). We argue that these attacks are effective because of the data sparsity of collaborative signals. In practice, auxiliary information, such as products' visual descriptions, is used to alleviate collaborative filtering data's sparsity. Therefore, when incorporating visual information in FedRecs, all existing model poisoning attacks' effectiveness becomes questionable. In this paper, we conduct extensive experiments to verify that incorporating visual information can beat existing state-of-the-art attacks in reasonable settings.

However, since visual information is usually provided by external sources, simply including it will create new security problems. Specifically, we propose a new kind of poisoning attack for visually-aware FedRecs, namely image poisoning attacks, where adversaries can gradually modify the uploaded image to manipulate item ranks during FedRecs' training process. Furthermore, we reveal that the potential collaboration between image poisoning attacks and model poisoning attacks will make visually-aware FedRecs more vulnerable to being manipulated. To safely use visual information, we employ a diffusion model in visually-aware FedRecs to purify each uploaded image and detect the adversarial images. Extensive experiments conducted with two FedRecs on two datasets demonstrate the effectiveness and generalization of our proposed attacks and defenses.

CCS Concepts: • **Information systems → Recommender systems**.

Additional Key Words and Phrases: federated learning, poisoning attack, multimodal recommendation, image pollution and purification

---

*Corresponding author.

---

Authors' addresses: Wei Yuan, The University of Queensland, Brisbane, QLD, Australia, w.yuan@uq.edu.au; Shilong Yuan, Nanjing University, Nanjing, Jiangsu, China, shilongyuan@nju.edu.cn; Chaoqun Yang, Griffith University, Gold Coast, QLD, Australia, chaoqun.yang@griffith.edu.au; Quoc Viet Hung Nguyen, Griffith University, Gold Coast, QLD, Australia, henry.nguyen@griffith.edu.au; Hongzhi Yin, The University of Queensland, Brisbane, QLD, Australia, db.hongzhi@gmail.com.

---

## 1 INTRODUCTION

Recommender systems have become an integral part of web applications (e.g., e-commerce [5, 53] and social media [58]), during the era of information explosion, since they are effective in reducing information overload by discovering users' potential interests. Traditionally, recommender systems are trained in a centralized server using a vast collection of user data [66]. However, with the growing awareness of privacy and the release of privacy protection regulations, such as the General Data Protection Regulation (GDPR) [49] in the European Union and the California Consumer Privacy Act (CCPA) [11] in the United States, collecting and storing user data has become more challenging.

Federated learning, as a privacy-preserving paradigm, allows for training models on decentralized data [32]. Consequently, an increasing number of researchers are exploring the potential of federated learning in recommender systems, resulting in the emergence of federated recommender systems (FedRecs). In FedRecs, the central server and users/clients[1] collaborate to learn a recommendation model by sharing model public parameters instead of user private data. Due to the significant advantage of data privacy protection, after the first FedRec framework proposed by Ammad et al. [2], several extended versions have sprung up to enhance the effectiveness and efficiency of FedRecs [18, 26, 34].

Due to the exposure of model parameters to all participants, and some of them may have malicious intentions, the security issues of FedRecs have raised concerns among researchers. Adversarial item rank manipulation is one of the most studied security problems in FedRecs, driven by financial incentives, which can lead to strong unfairness and even reduce the validity and usability of recommender systems. In [68], the first model poisoning attack, PipAttack, was introduced to demonstrate the vulnerability of FedRecs to being controlled by malicious participants who upload poisoned model updates. After that, all existing works about item rank manipulation in FedRecs are based on model poisoning attacks. For example, FedRecAttack [42] argues that PipAttack requires too many malicious users, which is not practical. It, on the other hand, achieves item promotion with fewer malicious users but requires more prior knowledge, such as a small proportion of user interaction data, which even violates the FedRec learning protocol. [41] proposed the first model poisoning attacks without any prior knowledge assumptions. Nevertheless, its performance is unstable and undesirable because it randomly samples vectors from a Gaussian distribution to act as the proxy of user embeddings. In our previous work [62], we proposed PSMU, which achieves state-of-the-art attack performance without relying on any prior knowledge and with fewer malicious users and training epochs, revealing the severe threats of model poisoning attacks to FedRecs.

However, all existing model poisoning attacks, including our previous work [62], only verify the threat in FedRecs with collaborative data (i.e., user-item interaction data). We argue that due to the inherent sparsity of collaborative information, many items, especially cold ones, lack sufficient descriptions. Consequently, these attacks can easily manipulate the rank order of items by uploading poisoned gradients. In real-life scenarios, item visual descriptions are used to alleviate collaborative data's sparsity problem. Intuitively, incorporating these visual signals makes the item features more comprehensive and robust, and as a result, existing state-of-the-art model poisoning attacks may fail

---

[1]In this paper, client and user are equivalent, since a client is only responsible for one user considering privacy protection requirements.
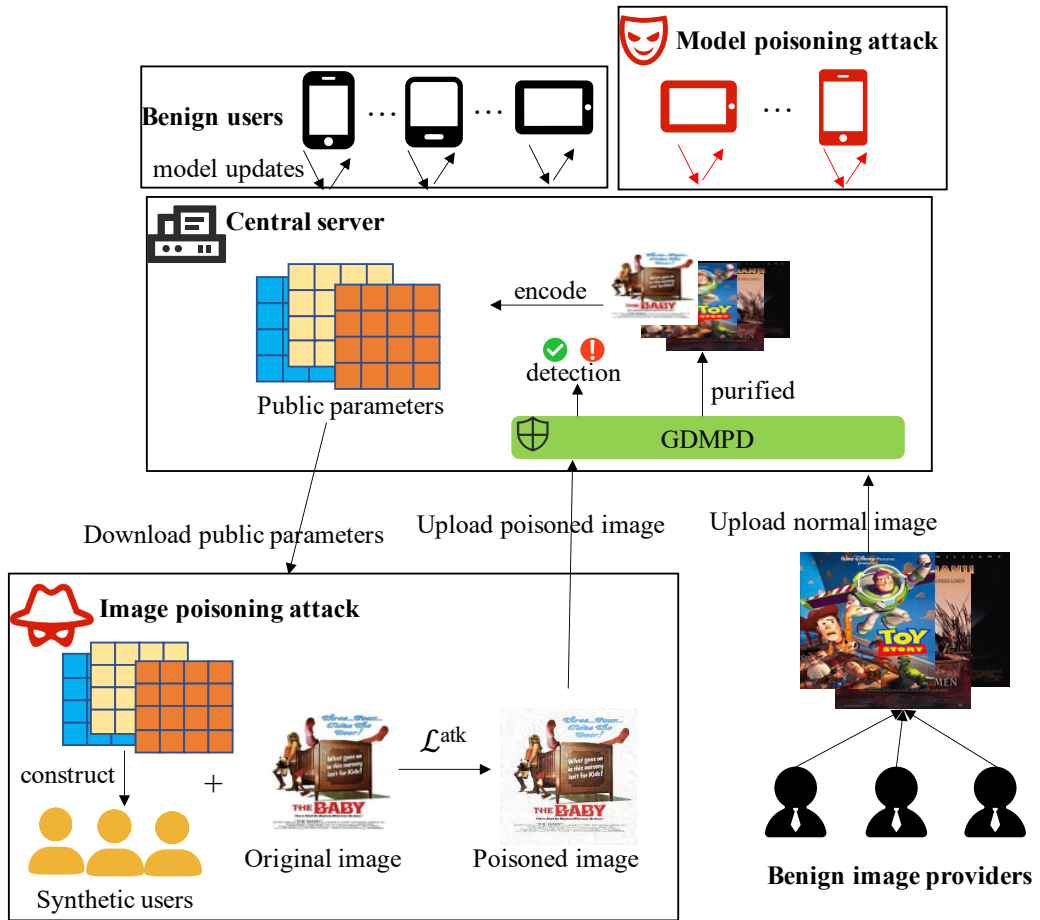
Fig. 1. Overview of the new threats "image poisoning attack" and our diffusion model based defense mechanism.

to promote items adversarially. In this paper, we empirically show that all existing state-of-the-art model poisoning attacks fail the adversarial item promotion in visually-aware FedRecs. i.e., visual information can mitigate the adversarial promotion threat caused by model poisoning attacks.

While visual signals can alleviate the model poisoning problem, incorporating them may leave another backdoor for adversaries to manipulate item ranks, as the product visual descriptions are typically provided by external sources that are not always trustworthy. In other words, the adversaries can be item image providers who promote the target items by uploading images with human-imperceptible noise. In this paper, we refer to such attacks as *image poisoning attacks*. Fig. 1 presents an overview of image poisoning attacks. It is worth noting that some research has used polluted images to change item ranks in centralized recommender systems (i.e., visual attacks [6, 29]). However, our image poisoning attacks have many different settings. Specifically, in centralized recommender systems, the model parameters are not accessible unless a "white-box" assumption is made. In contrast, the public parameters of FedRecs are apparent, but the user's private parameters are strictly out of reach. Moreover, all research in centralized recommender

systems [6, 29] assumes adversaries can obtain benign users' feedback, which is not valid in FedRecs. Furthermore, previous visual attacks for centralized recommendations can only be launched after the recommender system is well-trained, as adversaries cannot participate in the training process. While the image poisoning attacks in FedRecs are continually executed during the training process.

In this paper, we propose the first image poisoning attack, namely PSMU(V) (*poisoning with synthetic malicious users via visual* information), to disclose the risk of directly using visual information in FedRecs as shown in the bottom left of Fig. 1. Specifically, PSMU(V) is an image poisoning version of our previous work PSMU which is a model poisoning attack. The same as PSMU, PSMU(V) constructs a group of synthetic users with randomly selected interactions. It then calculates image perturbations with attack objectives guided by these synthetic users. Finally, the adversary uploads the poisoned image to the central server to influence the target item's feature representation. The above steps are iteratively executed with the training process of FedRecs. Through experimental results, we demonstrate the effectiveness of our image poisoning attacks and reveal the risks of using images provided by external sources directly. Additionally, since PSMU and PSMU(V) can have a consistent attack objective and can be launched simultaneously, we propose PSMU++ (i.e., PSMU+PSMU(V)) to reveal a more severe threat caused by the potential collaboration of PSMU and PSMU(V). That is, by launching both model poisoning attacks and image poisoning attacks, the target items will be more easily exposed to users than executing only one of these attacks.

The threats posed by image poisoning attacks underscore the need for a safer mechanism to use visual information. However, the defense against image poisoning attacks is still under-explored. In centralized recommendation, [48] attempted to employ adversarial training to improve visually-aware recommender systems' robustness, but it can only defend against untargeted attacks that aim to destroy a recommender system, and [29] indicated that adversarial training cannot effectively prevent item promotion attacks. Inspired by the great achievement of the Denoising Diffusion Probabilistic Model (DDPM) [16, 45] in image generation, we propose our novel image poisoning defender, *Guided Diffusion Model for Purification and Detection* (GDMPD), as shown in the middle part of Fig. 1. GDMPD can achieve two functions: purification and detection. The purification function aims to prevent adversarial images from achieving their malicious purpose. Particularly, the purification is based on DDPM which includes two processes: diffusion process and reverse process. During the diffusion process, the model gradually adds noise to the image, which can submerge the adversarial perturbations. Then, the reverse process purifies these noises to recover the image, which can remove both added noise and adversarial perturbations. In FedRecs, besides reducing the effectiveness of attacks, detecting malicious behavior is also necessary since it can provide the system manager with valuable insights for conducting further processes. Therefore, our GDMPD provides the detection function to further indicate which image is adversarial.

To support the proposed attack and defense methods, we extend the base FedRecs (Fed-NCF and Fed-LightGCN) used in our previous work [62] to visually-aware FedRecs. Then, we conduct extensive experiments with these FedRecs on two recommendation datasets (MovieLens-1M and Amazon Cell Phone). The experimental results demonstrate that incorporating visual signals can alleviate model poisoning attacks but simply using visual information provided by untrusty sources will leave a backdoor for image poisoning attacks, and our novel defense method can fix such a backdoor.

To sum up, our major new contributions are listed as follows:

- Our previous work [62] only studied the threat of model poisoning attacks for federated recommender systems on collaborative data. In this paper, we make the exploration of model poisoning attacks in visually-aware federated recommender systems. The empirical results

demonstrate that visually-aware federated recommender systems are robust to existing state-of-the-art model poisoning attacks, since visual signals can alleviate the data sparsity problem of collaborative information.

- Although visual information can defend against model poisoning attacks, we propose the first *image poisoning attack*, PSMU(V), to reveal a new backdoor for adversaries to promote items if visual information is directly used. To the best of our knowledge, this is the first work to reveal such threats in visually-aware FedRecs. Furthermore, we propose PSMU++ to investigate the potential hazard of collaboration between model poisoning attacks and image poisoning attacks.
- To fix the security hole of image poisoning attacks, we propose the Guided Diffusion Model for Purification and Detection (GDMPD), which is a diffusion model based defense mechanism in the central server of FedRecs to purify each uploaded image and detect the adversarial images. So far as we know, this is the first work that utilizes the diffusion model as defense method in federated recommender systems.
- We have performed comprehensive experiments using two visually-aware FedRecs that we extended from our previous work [62] on two widely-used recommendation datasets. The experimental results demonstrate the effectiveness and generalizability of our proposed methods.

The remainder of the paper is organized as follows. Related work is reviewed in Section 2, followed by the introduction of the visually-aware federated recommender systems in Section 3. Section 4 presents the technical details of our attacks extended from our previous work [62]. Then, in Section 5, we show how to fix the security problem revealed by image poisoning attacks. Section 6 exhibits a comprehensive analysis of experimental results. Finally, Section 7 gives a brief conclusion of this paper.

## 2 RELATED WORK

In this section, we briefly review the literature on four related topics: federated recommender systems, attacks and defense mechanisms for federated recommender systems, visually-aware recommender systems, and diffusion models. Other involved topics such as the development of general recommender systems and federated learning can be referred to corresponding surveys [24, 60, 70].

### 2.1 Federated Recommender Systems

Federated Recommender Systems (FedRecs) have gained increasing attention in recent years due to their ability to protect user privacy. Ammad et al [2] presented the first FedRec framework that applies federated learning with a collaborative filtering model. Based on this basic framework, many extended versions have been proposed in a short time [1, 47]. Some works attempt to reduce the performance gap between FedRecs and centralized recommender systems. For example, [54] utilizes Graph Neural Network (GNN) [43] to achieve fairly good recommendation accuracy. Wu et al. [56] employed contrastive learning in FedRecs. Other works focus on the efficiency of FedRecs. [4, 34] investigate fast convergence of FedRecs, while Zhang et al. [65] proposed a lightweight communication strategy based on learning to hash (L2H) [51]. ReFRS [18] learns dynamic and diversified user preferences on resource-constrained devices. Some work transplant FedRecs to specific recommendations, such as news recommendation [38], social recommendation [30], POI recommendation [10], and so on [27].

In addition to enhancing the effectiveness and efficiency of FedRecs, privacy concerns are also a crucial research direction in this area. Chai et al. [3] demonstrated that even if user data is not directly shared, adversaries can still recover sensitive information from the model updates sent

by the target user. To address this issue, a central server can apply a differential privacy (DP) mechanism to perturb the global model, as proposed in [52]. However, this approach assumes that the central server is in a sterile environment, which is not applicable in real-life scenarios where the server may be curious about clients' private information. To protect user privacy further, local differential privacy (LDP) is equipped on the client side [54]. Zhang et al. [67] introduced adaptive LDP, which can protect privacy with less impact on recommendation performance. Nevertheless, Yuan et al. [63] discovered that LDP alone cannot safeguard user-item interaction information, which is known as the user-item interaction leakage problem. They provided a regularization-based method to tackle such a problem. Additionally, [64] enables FedRecs to comply with the privacy regulations of the "right to be forgotten".

## 2.2 Attacks and Defenses for Federated Recommender Systems

Given the significant advancements achieved by FedRecs, many researchers are now investigating the security concerns associated with these systems. They have proposed several effective attack methods against FedRecs, exposing the vulnerabilities of current FedRecs under specific conditions. In general, the attacks in FedRecs can be divided into two categories: inference attacks and poisoning attacks. Inference attacks aim to detect certain information (e.g., user attributes [67], user private date [63]) from FedRecs to reveal certain privacy problems, as introduced in Section 2.1.

In this paper, our topic is closer to poisoning attacks. According to the attack's goal, there are targeted attacks and untargeted attacks. Untargeted attacks aim to cause a loss of recommendation accuracy and undermine the validity of the target model. FedAttack [55] attempts to compromise FedRecs using hard negative samples. Yu et al. [61] introduced a cluster-based attack to disrupt recommender systems. Targeted attacks aim to make specific items to be recommended to as many users as possible. Compared to untargeted attacks, targeted attacks are more stealthy and are more common due to financial incentives. Therefore, in this paper, we focus on targeted attacks. PipAttack [68] demonstrates that malicious users can manipulate the order of item ranks by uploading poisoned gradients. However, their attack requires a large proportion of compromised clients, which may be unaffordable in real applications. Rong et al. [42] reduced the number of malicious users by incorporating more prior knowledge. [41] is the first model poisoning attack that does not rely on any prior knowledge, but its performance is unstable. Our previous work [62] proposed a more effective model poisoning attack, PSMU, which achieves state-of-the-art performance. Besides, our previous work is the first to provide a defense mechanism based on gradient clipping to defend against existing model poisoning attacks.

However, all existing targeted poisoning attacks have been launched in FedRecs that rely on collaborative filtering data. These data suffer from severe sparsity problems, leading to strong biases that make it easier for poisoning attacks to manipulate items, especially for the cold ones. When visual signals are incorporated, data sparsity can be alleviated. Therefore, the performance of existing poisoning attacks for visually-aware FedRecs is unknown.

## 2.3 Visually-aware Recommender Systems and Attacks

Visual signals are important for making accurate recommendations [5, 17]. Some models solely rely on the feature vectors extracted from images to provide recommendations [19, 20]. However, the feature vectors are not optimized for making a good recommendation. VBPR [13] is the first work that fuses both visual information and collaborative signals based on BPR [40]. After that, many works [21, 28] are proposed based on the general framework: a pre-trained model is used to extract visual features, and then, certain fusion methods are used to aggregate visual features with collaborative signals (or other modality signals [39, 59, 69]) to feed a recommendation model. In

this paper, based on such a framework, we extend the basic FedRecs used in our previous work [62] to visually-aware FedRecs.

Due to the large scale of item catalogues, product visual descriptions are usually provided by external sources. Since visual information can influence the ranking of items, image providers may have a chance to adversarially manipulate item ranks. [8] attempts to change the popularity of items with the same categories. [29] is the first work research item promotion via image pollution in centralized recommendation with white-box settings. [6] further investigates such attacks under black-box settings. However, all of the above works are based on centralized recommender systems, and they assume that users' recommendation list is available, which is infeasible in FedRecs. Therefore, in this paper, we explore image poisoning attacks in FedRecs to reveal the threats of incorporating visual information, and then, we propose defense solution to prevent the threat.

## 2.4 Diffusion Models

Motivated by non-equilibrium thermodynamics [45], the diffusion model has shown a strong ability to generate high-quality images [7, 46]. Different from other commonly used generative models such as GANs [9] and VAEs [23], diffusion models generate samples by predicting the noise. Therefore, they are naturally suitable for adversarial image purification [57]. Nie et al. [36] was the first to use the diffusion model to purify adversarial images. Wang et al. [50] utilized guidance to further improve the quality and fidelity of purified images. This paper takes the first time to integrate diffusion models into the visually-aware federated recommender systems to allow the secure using of visual information.

## 3 VISUALLY-AWARE FEDERATED RECOMMENDER SYSTEMS

In this part, we provide the fundamental settings of our visually-aware federated recommender systems. The federated recommendation framework used in this paper is the same as all previous FedRecs attack works [41, 63, 68], which was originally proposed by [2].

Let $\mathcal{U}$ and $\mathcal{V}$ denote the set of benign users and items, respectively. $|\mathcal{U}|$ and $|\mathcal{V}|$ are the sizes of users and items. In FedRec, each user $u_i$ is a client who manages its' own training dataset $\mathcal{D}_i$. $\mathcal{D}_i$ consists of many user-item interactions $(u_i, v_j, r_{ij})$, where $r_{ij}$ is a binary rating denoting whether user $u_i$ has interacted with item $v_j$. That is, $r_{ij} = 1$ means $u_i$ has interacted with $v_j$, while $r_{ij} = 0$ indicates no interaction between $u_i$ and $v_j$. In addition, a single image $i_j$ is available for each item $v_j$ as an auxiliary description, which is uploaded by the item provider and managed by the central server. $\mathcal{V}_i^+$ and $\mathcal{V}_i^-$ are the sets of interacted items and non-interacted items for user $u_i$. Using the above data, FedRec aims to predict $\hat{r}_{ij}$ between user $u_i$ and non-interacted item $v_j$ and recommend items according to top-K highest prediction scores.

To ensure privacy protection, the parameters of the recommendation model are divided into private and public parameters. Private parameters are generally user embeddings U, which are maintained by corresponding users and are never shared with others. The public parameters, on the other hand, include item embeddings V, visual feature extractor $\Phi$, visual feature transform matrix E and other parameters $\Theta$, are transmitted between a central server and clients to achieve collaborative learning.

**Federated learning protocol.** In FedRecs, a central server coordinates the learning process. Initially, the central server initializes all public parameters, meanwhile, the clients initialize their corresponding private parameters locally. Then, a recommender system is trained by iteratively repeating the following steps. First, the central server randomly selects a set of users $\mathcal{U}_{t-1}$ to participate in the training process and sends the public parameters to these users. The selected users combine the received public parameters with the private parameters to form a local recommendation model. The local recommender system is trained on local dataset $\mathcal{D}_i$ by optimizing certain objective

functions, such as:

$$\mathcal{L}^{rec} = -\sum_{(u_i, v_j, r_{ij}) \in \mathcal{D}_i} r_{ij} \log \hat{r}_{ij} + (1 - r_{ij}) \log(1 - \hat{r}_{ij}) \tag{1}$$

After local training, the selected user $u_i$ updates its private parameters (E.q. 2) and transmits public parameters' gradients $\nabla \Theta_i^{t-1}$, $\nabla E_i^{t-1}$ and $\nabla V_i^{t-1}$ to the central server:

$$u_i^t = u_i^{t-1} - lr \nabla u_i^{t-1} \tag{2}$$

Then, the central server aggregates all received public parameter updates. The following formula takes item embeddings as an example:

$$V^t = V^{t-1} - lr \sum_{u_i \in \mathcal{U}_{t-1}} \nabla V_i^{t-1} \tag{3}$$

where $lr$ is the learning rate. Note that the visual extractor $\Phi$ can be trainable or freeze. In this paper, we freeze the visual extractor $\Phi$, since it is a large pretrained model and retraining it will dramatically increase the difficulty of convergence.

**Base recommendation model.** Generally, the above federated recommendation framework is compatible with most existing deep learning-based recommendation models [66]. In our previous work [62], we choose two classical and widely used recommenders, Neural Collaborative Filtering (NCF) [15] and LightGCN [14], as the base model. In this paper, we extend these two models to consider visual information when making recommendations, namely VNCF and LightVGCN. And then, we integrate these two models into the above FedRec framework to form Fed-VNCF and Fed-LightVGCN.

In Fed-VNCF, the local recommendation model in client $u_i$ predicts $\hat{r}_{ij}$ using the following formula:

$$\hat{r}_{ij} = \sigma(h^\top FFN([u_i, v_j, E\Phi(i_j)])) \tag{4}$$

where $h$ and $E$ are trainable public parameters, $u_i$ and $v_j$ are embeddings of user $u_i$ and item $v_j$, and $[\cdot]$ is concatenation operation.

For Fed-LightVGCN, the user-item interactions are viewed as a bipartite graph and all users and items are treated as distinct nodes. Then, user and item embeddings are learned by propagating their neighbour nodes' embeddings:

$$
\begin{aligned}
u_i^l &= \sum_{j \in \mathcal{N}_{u_i}} \frac{1}{\sqrt{|\mathcal{N}_{u_i}|}\sqrt{|\mathcal{N}_{v_j}|}} (v_j^{l-1} + E\Phi(i_j)) \\
v_j^l &= \sum_{i \in \mathcal{N}_{v_j}} \frac{1}{\sqrt{|\mathcal{N}_{v_j}|}\sqrt{|\mathcal{N}_{u_i}|}} u_i^{l-1}
\end{aligned} \tag{5}
$$

where $\mathcal{N}_{u_i}$ and $\mathcal{N}_{v_j}$ denote the sets of $u_i$'s and $v_j$'s neighbors. $l$ is the propagation layer. In order to protect privacy, each user can only perform the above calculation on its local bipartite graph. After propagation, we aggregate all layers' embeddings as the final user and item embeddings:

$$u_i = \sum_{l=0}^{L} u_i^l, \quad v_j = \sum_{l=0}^{L} v_j^l \tag{6}$$

Finally, the same as VNCF, we use E.q. 4 to compute the predicted preference scores.

---

**Algorithm 1** Visually-aware Federated Recommender Systems.

---

**Input:** global epoch $T$; local epoch $L$; learning rate $lr$, visual extractor $\Phi \ldots$
**Output:** public parameter V, E and $\Theta$, local client embedding $\mathrm{u}_i|_{i \in \mathcal{U}}$
 1: Initialize public parameter $\mathrm{V}^0$, $\mathrm{E}^0$, and $\Theta^0$
 2: Initialize item image set $\mathcal{I}^0 = \{\}$
 3: **for** each round t =1, ..., $T$ **do**
 4:     **if** new images uploaded by item providers **then**
 5:         $\mathcal{I}^t \leftarrow$ update $\mathcal{I}^{t-1}$ // The threats of image poisoning attack
 6:     **end if**
 7:     sample a fraction of clients $\mathcal{U}_{t-1}$ from $\mathcal{U}$ // The threats of model poisoning attack
 8:     **for** $u_i \in \mathcal{U}_{t-1}$ **in parallel do**
 9:         // run on client $u_i$
10:         calculate $\nabla \mathrm{u}_i^{t-1}$, $\nabla \mathrm{V}_i^{t-1}$, and $\nabla \mathrm{E}_i^{t-1}$, and $\nabla \Theta_i^{t-1}$ using E.q. 1
11:         $\mathrm{u}_i^t \leftarrow$ update local private parameters using E.q. 2
12:         upload $\nabla \mathrm{V}_i^{t-1}$, $\nabla \mathrm{E}_i^{t-1}$ and $\nabla \Theta_i^{t-1}$ to the central server
13:     **end for**
14:     $\mathrm{V}^t, \mathrm{E}^t, \Theta^t \leftarrow$ aggregate gradients using E.q. 3
15: **end for**

---

## 4 ADVERSARIAL ITEM PROMOTION VIA IMAGE AND MODEL POISONING ATTACKS

In this section, we present the details of our attacks, including the model poisoning attack (PSMU), image poisoning attack (PSMU(V)), and the combination of model poisoning attack and image poisoning attack (PSMU++).

### 4.1 Attack Task Formulation

Manipulating recommender systems includes promoting and demoting the rank order of items. In this work, we mainly discuss item promotion, since the demotion can be achieved by reversing the attack objective or promoting all other items. Adversarial item promotion has been widely studied in the poisoning attacks for FedRecs, which aims to increase the target item's exposure chances motivated by financial incentives [68]. However, all previous works only consider item promotion in FedRecs with collaborative data, none of them investigates the threat in visually-aware FedRecs. In this section, we present preliminaries and basic settings of adversarial item promotion in this paper.

**Attack goal.** Obviously, the goal of adversaries is to promote target items to as many users as possible. Formally, given that a recommender system recommends $K$ items $\hat{\mathcal{V}}_i$ to user $u_i$, the adversaries would like to improve the target item's exposure rate at rank $K$ (ER@K):

$$ER@K = \frac{1}{\left|\widetilde{\mathcal{V}}\right|} \sum_{v_j \in \widetilde{\mathcal{V}}} \frac{\left|\left\{u_i \in \mathcal{U} | v_j \in \hat{\mathcal{V}}_i \wedge v_j \in \mathcal{V}_i^-\right\}\right|}{\left|\left\{u_i \in \mathcal{U} | v_j \in \mathcal{V}_i^-\right\}\right|} \tag{7}$$

$\widetilde{\mathcal{V}}$ is the set of target items.

**Attack approach.** We explore two kinds of poisoning attacks in this work: model poisoning attacks (a.k.a., gradient poisoning attacks) and image poisoning attacks. For model poisoning attacks, as shown in the upper part of Fig. 1, the attacker will employ a group of malicious users $\widetilde{\mathcal{U}}$

to upload poisoned gradients to optimize E.q. 7:

$$\underset{\{\nabla \widetilde{V}^t, \nabla \widetilde{E}^t, \nabla \widetilde{\Theta}^t\}_{t=s}^{T-1}}{argmax} ER@K(U^T, V^T, E^T, \Theta^T) \qquad (8)$$

where $s$ is the epoch when the attacks are launched. $\nabla \widetilde{V}^t$, $\nabla \widetilde{E}^t$ and $\nabla \widetilde{\Theta}^t$ are gradients generated by malicious users at epoch $t$. These poisoned gradients will be aggregated in the central server using E.q. 3, since the central server is unaware of poisoned attacks.

For image poisoning attacks, we assume that the adversary is the target item image provider. It increases the target item's ER@K by uploading an image with human-unaware perturbations as follows:

$$\underset{\{\widetilde{I}^t\}_{t=s}^{T-1}}{argmax} ER@K(U^T, V^T, E^T, \Theta^T)$$
$$\widetilde{i}_j^t = i_j + \boldsymbol{\delta}_j^t, \text{ for } \widetilde{i}_j^t \in \widetilde{I}^t \qquad (9)$$

where $\widetilde{I}^t$ is the set of poisoned images for target items $\widetilde{V}$ at epoch $t$. Note that the model poisoning attacks and image poisoning attacks can be launched simultaneously to promote the same target items, as follows:

$$\underset{\{\nabla \widetilde{V}^t, \nabla \widetilde{E}^t, \nabla \widetilde{\Theta}^t, \widetilde{I}^t\}_{t=s}^{T-1}}{argmax} ER@K(U^T, V^T, E^T, \Theta^T) \qquad (10)$$

**Attack prior knowledge.** Following our previous work [62], in this paper, we assume that for both model poisoning attacks and image poisoning attacks, the attacker knows public parameters V, E, and $\Theta$ received from the central server, which is consistent with FedRecs' protocol. Besides, we assume the image poisoning attacker already knows the visual extractor $\Phi$, which is reasonable since $\Phi$ is an open-source pretrained model and it can be easily inferred by comparing the image feature vectors generated from the system and from guessed extractors.

### 4.2 Poisoning Attack

**Formulate the optimization problem.** The goal of all our attacks, including PSMU, PSMU(V), and PSMU++, is to promote target items $\widetilde{V}$ to as many users as possible. To achieve that, these attacks use different approaches to maximize E.q. 7, such as E.q. 8, E.q. 9, and E.q. 10. However, for all these attacks, it is challenging to directly optimize their objective function because of the following two problems: (1) The complex dependence of model parameter updates [62]; (2) ER@K is not differentiable; For the first problem, instead of finding a globally optimal solution, we greedily calculate the optimal results at each epoch, which will simplify the optimization problem:

$$\text{PSMU:} \quad \underset{\{\nabla \widetilde{V}^{t-1}, \nabla \widetilde{\Theta}^{t-1}\}}{argmax} ER@K(U^{t-1}, V^{t-1} - lr\nabla \widetilde{V}^{t-1}, E^{t-1} - lr\nabla \widetilde{E}^{t-1}, \Theta^{t-1} - lr\nabla \widetilde{\Theta}^{t-1})$$

$$\text{PSMU(V):} \quad \underset{\{\widetilde{I}^{t-1}\}}{argmax} ER@K(U^{t-1}, V^{t-1}, E^{t-1}, \Theta^{t-1})$$

$$\text{PSMU++:} \quad \underset{\{\nabla \widetilde{V}^{t-1}, \nabla \widetilde{\Theta}^{t-1}, \widetilde{I}^{t-1}\}}{argmax} ER@K(U^{t-1}, V^{t-1} - lr\nabla \widetilde{V}^{t-1}, E^{t-1} - lr\nabla \widetilde{E}^{t-1}, \Theta^{t-1} - lr\nabla \widetilde{\Theta}^{t-1})$$

$$(11)$$

For the second problem, we approximately optimize ER@K by forcing the target items' predicted preference scores to be higher than other recommended items' scores:

$$\mathcal{L}^{att} = \sum_{u_i \in \mathcal{U}} \sum_{v_t \in \widetilde{V} \wedge v_t \notin V_i^+} \sum_{v_j \in \widehat{V}_i \wedge v_j \notin \widetilde{V}} \sigma(\hat{r}_{ij} - \hat{r}_{it}) \qquad (12)$$

We omit the time index for the prediction scores in E.q. 12 to make the formula clear. To compute E.q. 12, we need the feedback from benign users (i.e., the recommended items $\hat{\mathcal{V}}_i$ and the scores of target items $\hat{r}_{it}$). The calculation of $\hat{\mathcal{V}}_i$ and $\hat{r}_{it}$ are based on user $u_i$'s private parameters, which is strictly not accessible in FedRecs. Following our previous work [62], we construct a group of synthetic users to replace benign users. Specifically, we randomly select a group of items as fake user $\widetilde{u}_i$'s interacted items $\widetilde{\mathcal{V}}_i^+$. Based on $\widetilde{\mathcal{V}}_i^+$, we build the synthetic dataset $\widetilde{\mathcal{D}}_i$. Then, we fix all public parameters and train the fake user's user embeddings with the recommendation objective:

$$\widetilde{U}^{t-1} = \underset{\widetilde{U}^{t-1}}{argmin} \, \mathcal{L}^{rec}(\widetilde{U}^{t-1}, V^{t-1}, E^{t-1}, \Theta^{t-1}, \widetilde{\mathcal{D}}^{t-1}) \tag{13}$$

Note that at different epochs, we reconstruct different $\widetilde{\mathcal{V}}_i^+$ and $\widetilde{\mathcal{D}}_i$, so that even with a small size of malicious users, we can still simulate many synthetic users.

After constructed synthetic users, E.q. 12 is transformed to:

$$\widetilde{\mathcal{L}}^{att} = \sum_{\widetilde{u}_i \in \widetilde{\mathcal{U}}} \sum_{v_t \in \widetilde{\mathcal{V}} \wedge v_t \notin \widetilde{\mathcal{V}}_i^+} \sum_{v_j \in \widetilde{\mathcal{V}}_i \wedge v_j \notin \widetilde{\mathcal{V}}} \sigma(\hat{r}_{ij} - \hat{r}_{it}) \tag{14}$$

where $\widetilde{\mathcal{V}}_i$ is the set of items that have the highest prediction scores for malicious user $\widetilde{u}_i$.

**Optimizing via poisoned gradients (PSMU).** As shown in E.q. 11, PSMU promotes target items by uploading poisoned gradients to the central server. Specifically, given current public parameters $V^{t-1}$, $E^{t-1}$, $\Theta^{t-1}$, and the calculated fake user embeddings $\widetilde{U}^{t-1}$, we compute the poisoned gradients as follows:

$$\nabla \widetilde{V}^{t-1} = \frac{\partial}{\partial V^{t-1}} \widetilde{\mathcal{L}}^{att}(\widetilde{U}^{t-1}, V^{t-1}, E^{t-1}, \Theta^{t-1})$$

$$\nabla \widetilde{E}^{t-1} = \frac{\partial}{\partial E^{t-1}} \widetilde{\mathcal{L}}^{att}(\widetilde{U}^{t-1}, V^{t-1}, E^{t-1}, \Theta^{t-1}) \tag{15}$$

$$\nabla \widetilde{\Theta}^{t-1} = \frac{\partial}{\partial \Theta^{t-1}} \widetilde{\mathcal{L}}^{att}(\widetilde{U}^{t-1}, V^{t-1}, E^{t-1}, \Theta^{t-1})$$

It is worth noting that compared to original poison FedRecs, in visually-aware FedRecs, we also poison public parameters related to visual signals, such as E, for a fair comparison. Naturally, poisoning more items will make the attack goal easier to achieve, however, it will also cause too many side effects on FedRec performance. Therefore, for item embeddings, PSMU only uploads poisoned gradients for target items' embeddings.

$$\nabla \widetilde{V}^{t-1} = \begin{cases} 0 & v_m \notin \widetilde{V} \\ \nabla \widetilde{V}_m^{t-1} & v_m \in \widetilde{V} \end{cases} \quad m = 0, 1, \ldots, |\mathcal{V}| \tag{16}$$

**Optimizing via poisoned images (PSMU(V)).** In visually-aware FedRecs, item images are provided by external sources, which are usually item providers. These item providers may provide images with slight pollutions to mislead recommender systems to give higher preference scores to their items for as many users as possible. Formally, given $V^{t-1}$, $E^{t-1}$, $\Theta^{t-1}$ and $\widetilde{U}^{t-1}$, the attacker calculates perturbations as follows:

$$\delta^{t-1} = \underset{\delta^{t-1}}{argmin} \, \widetilde{\mathcal{L}}^{att}(\widetilde{U}^{t-1}, V^{t-1}, E^{t-1}, \Theta^{t-1}), \quad \|\delta\| \leq \epsilon \tag{17}$$

To avoid the perturbations being aware by normal users, the attacker should restrict the size of noise $\delta^{t-1}$ with the bound $\epsilon$ at each epoch. Previous visual attacks in centralized recommender systems usually set $\epsilon$ to at least 32 for a 255 pixel value range. In this paper, we only use $\epsilon = 4$ then

we can achieve ER@5=1.0 results, which shows the severe threats of image poisoning attacks in FedRecs.

**PSMU++.** Since PSMU and PSMU(V) have consistent objectives and there is no conflict between their poisoning implementations, the adversary can combine these two attacks together to form a more effective item promotion attack, named PSMU++. In PSMU++, we assume the item providers are the adversaries and they not only upload poisoned product images, but also upload poisoned gradients using a group of compromised users. Specifically, the malicious users upload poisoned gradients by launching the PSMU algorithm, meanwhile, the item provider launch PSMU(V) based on these malicious users' synthetic user embeddings (i.e., Algorithm 3 Line 4).

---

**Algorithm 2** PSMU: Poisoning with Synthetic Malicious Users

---

**Input:** public parameters $V^{t-1}$, $\Theta^{t-1}$, $E$
**Output:** public parameter poisoned gradients $\nabla\widetilde{V}_i^{t-1}$, $\nabla\widetilde{E}_i^{t-1}$, $\nabla\widetilde{\Theta}_i^{t-1}$
1: // run on malicious client $\widetilde{u}_i$
2: randomly construct training set $\widetilde{\mathcal{D}}_i^{t-1}$
3: calculate synthetic user embedding $\widetilde{u}_i^{t-1}$ using E.q. 13
4: calculate $\nabla\widetilde{V}_i^{t-1}$, $\nabla\widetilde{E}_i^{t-1}$, $\nabla\widetilde{\Theta}_i^{t-1}$ using E.q. 15
5: $\nabla\widetilde{V}_i^{t-1} \leftarrow$ constraint $\nabla\widetilde{V}_i^{t-1}$ using E.q. 16
6: upload $\nabla\widetilde{V}_i^{t-1}$, $\nabla\widetilde{E}_i^{t-1}$, $\nabla\widetilde{\Theta}_i^{t-1}$ to the central server

---

**Algorithm 3** PSMU(V): Poisoning with Synthetic Malicious Users via Visual Information

---

**Input:** public parameters $V^{t-1}$, $\Theta^{t-1}$, $E$
**Output:** image adversarial perturbation $\delta^{t-1}$
1: // run on target items' image provider
2: **if** malicious users $\widetilde{\mathcal{U}}$ exist **then**
3:     // PSMU++
4:     request for malicious user embeddings $\widetilde{U}^{t-1}$
5: **else**
6:     randomly construct training set $\widetilde{\mathcal{D}}^{t-1}$ and calculate user embeddings $\widetilde{U}^{t-1}$ using E.q. 13
7: **end if**
8: calculate $\delta^{t-1}$ using E.q. 17
9: construct and upload adversarial images $\widetilde{I}^{t-1}$

---

## 5   GUIDED DIFFUSION MODEL FOR PURIFICATION AND DETECTION

The experimental results in Section 6 indicate that traditional model poisoning attacks, such as PSMU, are ineffective in visually-aware FedRecs. However, the presence of visual information creates another backdoor, which provides an opportunity for corporations to promote items effectively through PSMU and PSMU(V), underscoring the urgent need for image poisoning defense. GDMPD leverages a pretrained Denoising Diffusion Probabilistic Model (DDPM) [16, 45] to purify all uploaded images with guidance, eliminating the need for additional computation resources for training. After purification, adversarial images have a high probability of losing their delicate perturbations. To maintain the recommender system, it is essential to detect which images are adversarial. Based on the purified image, our GDMPD achieves an adversarial image detection
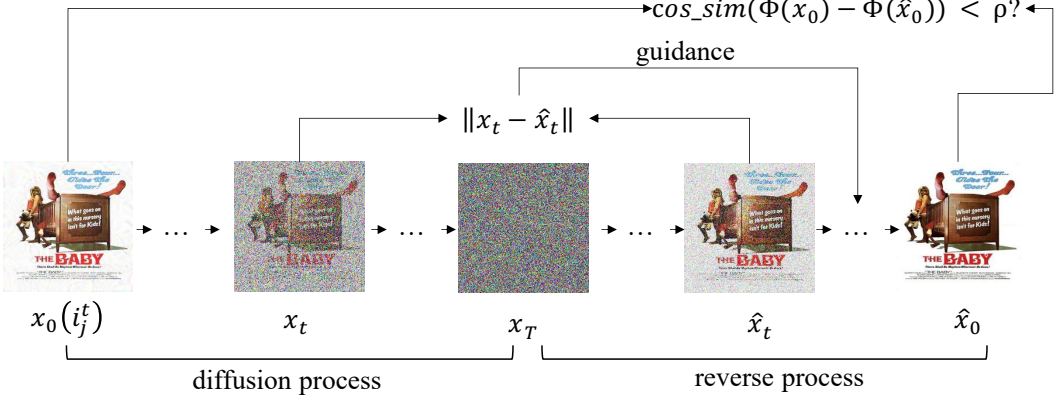
Fig. 2. Overview of the proposed image poisoning defender: GDMPD.

function. Notably, our detection method is also training-free. Fig. 2 gives an overview illustration of our GDMPD.

## 5.1 Diffused Image Guided Purification

In Section 4, we use $i_j^t$ to represent the visual discription of item $j$ uploaded at epoch $t$. For clarity of description, when using the image as input for GDMPD, we denote it as $x$. Generally, DDPM consists of two Markov processes: the diffusion process and the reverse process. In the diffusion process, DDPM adds noise to the input image at each time step until it becomes Gaussian noise. Then, the reverse process gradually removes this noise to recover the input image.

**Diffusion process.** Formally, assume $x_0$ be the input image where $t = 0$. Note that to avoid misunderstanding with FedRec's global epochs $t$ denoted by superscript, here we use subscript $t$ to denote the diffusion time step. $T$ is the length of diffusion steps. DDPM incrementally corrupts the input image $x_0$ into Gaussian noise as follows:

$$q(x_1, x_2, \ldots, x_T | x_0) = \prod_{t=1}^{T} q(x_t | x_{t-1})$$

$$q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I) \tag{18}$$

where $\mathcal{N}(x, \mu, \sigma^2)$ means $x$ is sampled from a Gaussian distribution with a mean $\mu$ and variance $\sigma$. $\beta_t$ is generated from a predefined noise adding schedule $\beta$. Common settings of $\beta$ include consine [16], square-root [25], and linear schedule [50]. In this paper, following [7, 50], we use the linear schedule and define $\beta_t$ as: $\beta_t = \frac{t-1}{T-1}(\beta_T - \beta_1)$, where $\beta_T = 2 \times 10^{-2}$ and $\beta_1 = 1 \times 10^{-4}$ are hyper-parameters. According to [16], we can directly calculate $x_t$ at an arbitrary diffusion step directly conditioned on $x_0$ with following euqation:

$$q(x_t | x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) I)$$

$$\bar{\alpha}_t = \prod_{i=1}^{t} \alpha_i, \quad \alpha_i = 1 - \beta_i \tag{19}$$

Then, with reparameter trick, we can generate $x_t$ as follows:

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \zeta \tag{20}$$

where $\zeta$ is noise sampled from standard Gaussian distribution, i.e., $\zeta \sim \mathcal{N}(0, I)$.

Considering the input image is adversarial image $x_0^{adv}$ (a.k.a., $\widetilde{i}$ in Section 4) and $x_0^{adv} = x_0 + \boldsymbol{\delta}$, then, after $t$ steps diffusion, the image equals to:

$$x_t^{adv} = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{\bar{\alpha}_t}\boldsymbol{\delta} + \sqrt{1 - \bar{\alpha}_t}\zeta \tag{21}$$

When $t$ increases, $\sqrt{\bar{\alpha}_t}$ will gradually decrease while $\sqrt{1 - \bar{\alpha}_t}$ will gradually increase. Since $\|\boldsymbol{\delta}\|$ should be small (lower than $\epsilon$) to guarantee the perturbations' unawareness, after an appropriate length of diffusion, the magnitude of Gaussian noise $\zeta$ will be large enough to submerge the delicately calculated perturbation $\boldsymbol{\delta}$. Meanwhile, the semantic meaning of original image $x_0$ still can be largely preserved, since $\boldsymbol{\delta}$ is negligible compared with $x_0$.

**Reverse process.** The reverse process is a Markov process that denoises the diffused $x_t$ to approximate the original input $x_0$ by predicting the noise added in the diffusion process. Formally, the reverse process from step $T$ to 0 is as follows:

$$p_w(\hat{x}_0, \hat{x}_1, \ldots, \hat{x}_{T-1}|x_T) = \prod_{t=1}^{T} p_w(\hat{x}_{t-1}|\hat{x}_t)$$
$$p_w(\hat{x}_{t-1}|\hat{x}_t) = \mathcal{N}(\hat{x}_{t-1}; \boldsymbol{\mu}_w(\hat{x}_t, t), \Sigma_w(\hat{x}_t, t)I) \tag{22}$$

where the mean $\boldsymbol{\mu}_w(x_t, t)$ is a neural network parameterized by w, the variance $\Sigma_w(x_t, t)$ can be either neural network or predefined time step dependent constant [16, 35]. As results, the reverse process iteratively samples $\hat{x}_{t-1}$ using $p_w(\hat{x}_{t-1}|\hat{x}_t)$ to get the predicted input image $\hat{x}_0$. Assume the input image is adversarial image $x_0^{adv}$, after the process of E.q. 21, the adversarial perturbations are corrupt by gradually added Gaussian noise. Then, E.q. 22 is used to eliminate the Gaussian noise and is very likely to simultaneously remove the perturbations. This is because: (1) The normalization of perturbation is small and the adversarial information is destroyed by adding noise in $x_t^{adv}$; (2) We use a pretrained DDPM which is learned on normal image datasets, therefore, it tends to recover the image to the domain of clean images in the reverse process.

**Guided reverse process.** However, simply using DDPM to purify uploaded images in FedRecs will have the following challenge: how to largely recover the original semantic of the input image meanwhile remove most perturbations. Specifically, if the diffusion steps $t$ are too large, the original information of $x_0$ in $x_t$ will be destroyed and the results of the reverse process will tend to be random [7]. As a result, the FedRecs' recommendation performance will be compromised since all items' visual information will be altered by the defense mechanism. On the contrary, when the diffusion steps are too small, the diffusion and reverse process may not be strong enough to purify all perturbations. In FedRecs, since perturbations are usually small, when we diffuse the image to Gaussian distribution, the perturbations will be largely be submerged. Therefore, the challenge is mainly about how to recover an image with high quality.

To improve the fidelity of images generated by the reverse process, we propose to add guidance during the reverse process. Concretely, we use the counterpart image $x_t$ in the diffusion steps to guide $\hat{x}_t$'s generation as shown in Fig. 2. To achieve this, we modify the reverse process $p_w(\hat{x}_{t-1}|\hat{x}_t)$ in E.q. 22 to condition on $x_t$, i.e., $p_w(\hat{x}_{t-1}|\hat{x}_t, x_t)$. According to [7, 45], we can further get the following approximation:

$$\log p_w(\hat{x}_{t-1}|\hat{x}_t, x_t) \approx \log p_w(\hat{x}_{t-1}|\hat{x}_t)p(x_t|\hat{x}_t)$$
$$\approx \log p(z) \tag{23}$$

$$z \sim \mathcal{N}(\boldsymbol{\mu}_w(\hat{x}_t, t) + \Sigma_w(\hat{x}_t, t)\nabla_{\hat{x}_t}\log p(x_t|\hat{x}_t), \Sigma_w(\hat{x}_t, t)I) \tag{24}$$

where $p_w(\hat{x}_{t-1}|\hat{x}_t)$ is the probability from unconditional DDPM and $p(x_t|\hat{x}_t)$ can be interpreted as "how close $x_t$ and $\hat{x}_t$ are". In this paper, $p(x_t|\hat{x}_t)$ is designed as follows:

$$p(x_t|\hat{x}_t) = \exp(\lambda \|x_t - \hat{x}_t\|) \tag{25}$$

$\lambda$ is the factor that controls the scale of guidance. $\|\cdot\|$ is mean squared error. Combined with E.q. 24, we can get:

$$z \sim \mathcal{N}(\boldsymbol{\mu}_w(\hat{x}_t, t) + \lambda \Sigma_w(\hat{x}_t, t) \nabla_{\hat{x}_t} \|x_t - \hat{x}_t\|, \Sigma_w(\hat{x}_t, t) I) \tag{26}$$

Finally, we can use the pretrained DDPM to infer purified image $\hat{x}_0$ given the input image $x_0$. Algorithm 4 illustrates how to purify uploaded item visual information using our diffused image-guided DDPM.

---

**Algorithm 4** Diffused Image Guided Purification

---

**Input:** pretrained DDPM $\boldsymbol{\mu}_w(\hat{x}_t, t)$ and $\Sigma_w(\hat{x}_t, t)$, guidance factor $\lambda$, input image $x_0$, ...
**Output:** purified image $\hat{x}_0$
 1: // diffusion process
 2: **for** each round t=1, ..., $T$ **do**
 3:     calculate diffused image $x_t$ with E.q. 19 and E.q. 20
 4: **end for**
 5: // reverse process
 6: **for** each round t=$T$, ..., 1 **do**
 7:     $\boldsymbol{\mu}, \Sigma \leftarrow \boldsymbol{\mu}_w(\hat{x}_t, t), \Sigma_w(\hat{x}_t, t)$
 8:     sample $\hat{x}_{t-1}$ from $\mathcal{N}(\boldsymbol{\mu} + \lambda \Sigma \nabla_{\hat{x}_t} \|x_t - \hat{x}_t\|, \Sigma I)$
 9: **end for**

---

## 5.2 Adversarial Image Detection

The image purification function guarantees that all uploaded images used for recommendations are unlikely to contain adversarial perturbations. However, in real-life scenarios, it is also essential to detect adversarial images as it provides insights into system maintenance. For example, the system manager can collect detected images to analyze potential attacks and even punish the detected adversarial image providers directly.

Detecting adversarial images in FedRecs is non-trivial since we only have normal images and we cannot get adversarial images for training before we can detect them. As we know, image poisoning attacks achieve adversarial goals by adding imperceptible perturbations to the image. These perturbations can cause remarkable changes for the image feature vectors which are encoded by the extractor $\Phi$. Based on this characteristic, we propose a training-free method to detect adversarial images. Specifically, we assume that the image purified by Algorithm 4 will be a "safe" image. In other words, the purified image will not cause remarkable changes to its encoding feature, since the purifications are destroyed. Therefore, we employ the feature extractor to encode the image before and after purification and compare the difference between these two feature vectors. If the similarity of these two feature vectors is smaller than a threshold $\rho$, the image will be detected as an adversarial image:

$$cos\_sim(\Phi(x_0) - \Phi(\hat{x}_0)) < \rho \tag{27}$$

We use cosine similarity to measure the difference of image vectors[2]. $\rho$ is a preset hyper-parameter and its value will directly influence the accuracy of the detector. In this paper, we set $\rho$ as follows: first, we use our DDPM to purify a large number of clean images, such as those from public datasets

---

[2]We also tried Euclidean distance and get equivalent experimental results.

like ImageNet. Next, we calculate the difference between purified and original images and use statistics to define $\rho$. If the difference between purified and original images is smaller than $\rho$, such an image is likely different from clean images from the image extractor perspective, indicating that it is an adversarial image.

## 6 EXPERIMENTS

In this section, we take extensive experiments to answer the following research questions (RQs):

- **RQ1.** Are current state-of-the-art model poisoning attacks still effective for visually-aware federated recommender systems?
- **RQ2.** Are there any risks of using visual information in federated recommender systems? i.e., Are our PSMU(V) and PSMU++ effective for visually-aware federated recommendations?
- **RQ3.** How is the effectiveness of our diffusion model-based defense method?

### 6.1 Dataset

In this paper, we leverage two popular federated recommendation datasets for evaluation: MovieLens-1M (ML) [12] and Amazon Cell Phone (AZ) [31]. ML contains $6,040$ users and $3,706$ items with $1,000,208$ feedback, and $3,301$ items have image descriptions. AZ includes $103,593$ interactions with $13,174$ users, $5,970$ products, and $5,877$ visual signals. All users have at least $5$ interactions with different products. Following [62, 68], we binarize the user-item ratings, where all ratings are transformed to $r_{ij} = 1$ and negative instances are sampled with $1 : 4$ ratio. Table 1 illustrates the basic statistics of these two datasets. It is worth pointing out that we choose two datasets with very different data sparsity to show the data sparsity problem's impacts on our poisoning attacks.

Table 1. Statistics of recommendation datasets

| Dataset | #users | #items | #interactions | Avg. | Density |
|---------|--------|--------|---------------|------|---------|
| ML      | 6,040  | 3,706  | 1,000,208     | 166  | 4.46%   |
| AZ      | 13,174 | 5,970  | 103,593       | 8    | 0.13%   |

### 6.2 Evaluation Protocol

For both model poisoning attacks and image poisoning attacks, the evaluation protocol is consistent with our previous work [62]. Specifically, FedRecs are trained without attacks for a few epochs and then we launch the attacks. The FedRecs are trained until convergence or reaching the pre-defined maximum global epochs. We select the most unpopular items as target items. ER@5 is used to evaluate both model poisoning attacks and image poisoning attacks. The purification of the defense method is evaluated from two aspects: (1) whether it can reduce target items' ER@5 to normal level; (2) whether it deteriorates recommendation performance (NDCG@20). The detection of the defense method is evaluated by accuracy.

### 6.3 Parameter Settings

All the experiments are implemented using PyTorch [37]. For both Fed-NCF, Fed-VNCF, Fed-LightGCN, and Fed-LightVGCN, the dimension of user and item embeddings are set to 32 following [62]. We use the deep pretrained CNN model from [44] as our visual extractor $\Phi^3$. Then, the visual feature extracted by CNN model is transformed to 32 dimension size vector by the

---

[3]We also tried other visual extractors such as ResNet and get similar results and trends, since our attack and defense do not make any special assumption on the visual extractor.
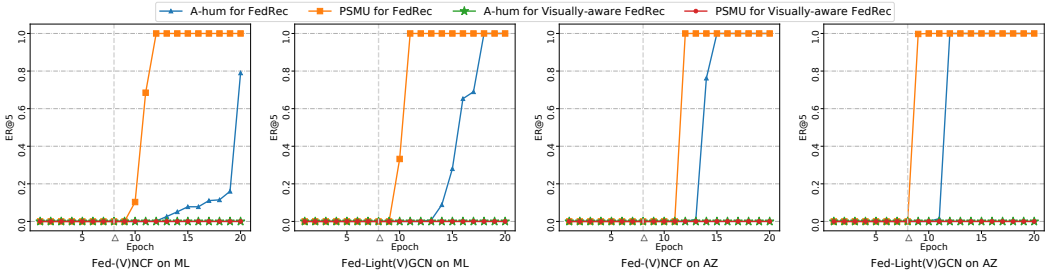
Fig. 3. The performance comparison of existing state-of-the-art (SOTA) model poisoning attacks for general federated recommender systems and visually-aware federated recommender systems. All current SOTA model poisoning attacks are ineffective in visually-aware FedRecs.

visual feature transform matrix E. Three layers of feedforward layers are utilized to process the concatenated user, item, and visual feature vectors (optional) with sizes of $[96, 32, 16]$ for Fed-VNCF and Fed-LightVGCN and $[64, 32, 16]$ for Fed-NCF and Fed-LightGCN, respectively. The layer of LightGCN propagation is 1 for both Fed-LightGCN and Fed-LightVGCN. Adam [22] with a learning rate of $0.001$ is adopted as the optimizer. The same as in [62], all poisoning attacks are launched at 8th global epoch. The number of selected items for synthetic users $\left|\widetilde{\mathcal{V}}_i^+\right|$ is 30. For model poisoning attacks, the proportion of malicious users $\xi$ equals $0.1\%$ without specific mention. For image poisoning attacks, $\epsilon$ is 4 as default, which is much smaller than visual attacks in centralized recommendation [6, 29]. For the diffusion model, we use an unconditional $256 * 256$ DDPM[4] pretrained by [7]. The number of diffusion steps is set to 1000 according to [7]. The guidance factor $\lambda$ is 1000.

## 6.4 Effectiveness of Model Poisoning Attacks for Visually-aware FedRecs (RQ1)

All existing model poisoning attacks [41, 42, 68], including our previous work [62], have demonstrated their effectiveness only in FedRecs with collaborative data. In this paper, we argue that the effectiveness of these attacks is due to the sparsity of collaborative information, which results in less robust item embeddings (especially for cold items) due to insufficient data description. Therefore, when additional item auxiliary information such as product visual description is incorporated, these poisoning attacks may become ineffective.

To support our argument, we conduct experiments with both general FedRecs and visually-aware FedRecs using A-hum and PSMU. A-hum is the earlier state-of-the-art model poisoning attack proposed by Rong et al [41]. PSMU is the current state-of-the-art model poisoning attack proposed by our previous work [62]. We choose these two attacks to do experiments since our previous work [62] already showed that other model poisoning attacks have very poor performance with limited malicious users ($\xi = 0.1\%$).

Fig. 3 presents the performance comparison of these two model poisoning attacks on FedRecs with or without visual information. For Fed-NCF and Fed-LightGCN, both PSMU and A-hum have the ability to influence the value of exposure rate on all datasets, meanwhile, our PSMU achieves better performance than A-hum (i.e., achieving higher ER@5 values or achieving ER@5=1.0 with fewer epochs.) This result proves the effectiveness of these two model poisoning attacks for FedRecs with only collaborative information. Besides, by comparing the same attack's performance in the same FedRec across datasets, we can observe that the sparser the dataset is, the better performance the attack has. For example, PSMU obtains $1.0$ ER@5 scores using about 4 and 3 epochs on the ML

---

[4]https://openaipublic.blob.core.windows.net/diffusion/jul-2021/256x256_diffusion.pt

(a) Fed-VNCF on ML.          (b) Fed-VNCF on AZ.          (c) Fed-LightVGCN on ML.(d) Fed-LightVGCN on AZ.
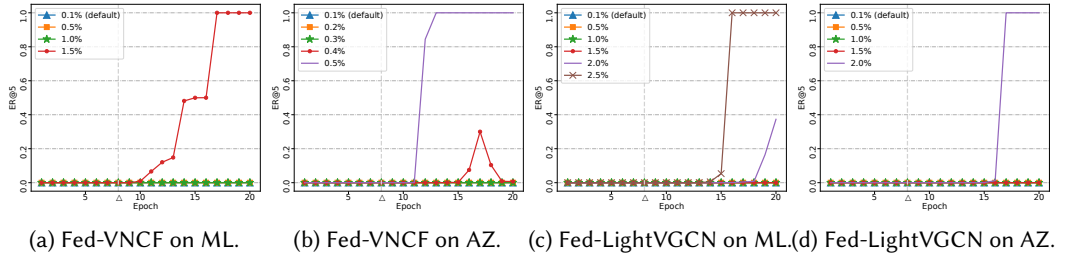
Fig. 4. The effectiveness of PSMU with more malicious user proportion in visually-aware FedRecs.

dataset for Fed-NCF and Fed-LightGCN, while it only costs 3 and 1 epochs on AZ. However, when FedRecs were equipped with product visual information, all these attacks' ER@5 scores dropped to 0 in all cases (the line of "A-hum for Visually-aware FedRec" and "PSMU for Visually-aware FedRec" in Fig. 3), which indicates that incorporating visual information can make Fed-NCF and Fed-LightGCN more robust to malicious poisoning attacks.

In Fig. 3, we have already demonstrated that PSMU, our state-of-the-art model poisoning attack, is unable to promote target items in visually-aware FedRecs when using the default setting of our previous work [62]. However, the effectiveness of PSMU may increase with more malicious users, although the cost of launching such attacks will also increase. Therefore, we investigate the effectiveness of PSMU with more malicious users in Fig. 4. The results show that PSMU requires a higher proportion of malicious users to be effective in visually-aware FedRecs. Specifically, PSMU cannot be effective until the proportion of malicious users increased to $1.5\%$ and $0.5\%$ for Fed-VNCF on ML and AZ, respectively, which are 15 and 5 times higher than the settings for FedRecs with only collaborative data. To manipulate Fed-LightVGCN, PSMU requires at least 25 and 20 times more malicious users on ML and AZ than the default settings. As a result, the costs of utilizing model poisoning attacks to compromise visually-aware FedRecs are much higher than the original FedRecs. Furthermore, by comparing different models on the same dataset (i.e., Fig. 4a and Fig. 4c, Fig. 4b and Fig. 4d), we find that Fed-LightVGCN is relatively more robust than Fed-VNCF when facing model poisoning attacks. This is because visual information is fully utilized in Fed-LightVGCN compared to Fed-NCF: Fed-LightVGCN not only uses visual information for directly predicting the preference scores (E.q. 4) but also utilizes it during LightGCN propagation (E.q. 5).

By combining Fig. 3 and Fig. 4 we can conclude that incorporating visual information can improve the robustness of FedRecs for current state-of-the-art model poisoning attacks.

## 6.5 Effectiveness of PSMU(V) and PSMU++ (RQ2)

Although Section 6.4 manifests that using visual information can defend against model poisoning attacks, in this subsection, we disclose that visual information will create new backdoors for adversaries to promote items by presenting the effectiveness of PSMU(V) (RQ2). Besides, we further reveal that the backdoor of visual information gives adversaries an opportunity to simultaneously launch image and model poisoning attacks (PSMU++) to manipulate item ranks.

As mentioned in Section 1, we are the first to present image poisoning attacks in visually-aware FedRecs. Most previous visual attacks [6, 29] in the centralized recommendation are not applicable in FedRecs settings since they depend on the feedback of benign users. For comparison purposes, we construct the following baselines: No Attack and Popularity Attack. No Attack displays the original exposure rate of target items. Popularity Attack is similar to the EXPA attack proposed by [29], but
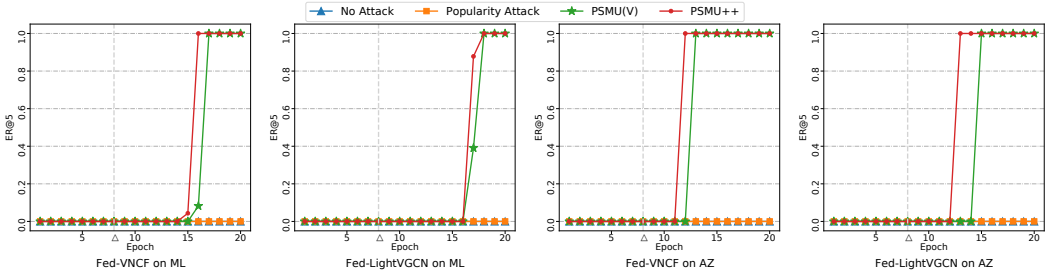
Fig. 5. The trend of exposure rate of target items for different image poisoning attacks with $\epsilon = 4$.

it differs in that it gradually changes item images during the training process. In Popularity Attack, we assume that adversaries have knowledge of the popularity information of items. At each global epoch, the attacker tries to add noise to make the target item's visual vector close to the feature vector of popular items.



Fig. 6. Example of adversarial images with different perturbation constraints ($\epsilon$) and corresponding purified images generated by different purification methods. All these adversarial images are generated by PSMU(V) and can achieve ER@5=1.0. The original image is from the ML dataset.

Fig. 5 presents the results of different image poisoning attacks. Both PSMU(V) and PSMU++ use $0.1\%$ synthetic users which is the same as the setting in model poisoning attacks. In Fig. 5, we can see that No Attack and Popularity Attack cannot create any changes in the exposure rate of the target item. In contrast, PSMU(V) promotes target items to all users with 9, 10, 5, and 7 global epochs in the cases of "Fed-VNCF on ML", "Fed-LightVGCN on ML", "Fed-VNCF on AZ", and "Fed-LightVGCN on AZ", respectively. Moreover, when incorporating PSMU (i.e., PSMU++), the item promotion process is accelerated as shown by the red line in Fig. 5. Besides, comparing the results from different datasets, we can get a consistent conclusion with Fig. 3: Promoting items

is easier on AZ than on ML since AZ is sparser. Additionally, Fed-LightVGCN is relatively more reliable than Fed-NCF under poisoning attacks since the visual information is fused by not only concatenation but also with LightGCN propagation. It is worth noting that we set the normalization of perturbations to be less than 4, making the polluted image be human-imperceptible. The first line of Fig. 6 provides an example of adversarial images with different perturbation scales. $\epsilon = 0$ represents the original image. The adversarial image with $\epsilon = 4$ perturbations is indistinguishable from the original image to humans, ensuring the stealthiness of our image poisoning attacks.

## 6.6 The Effectiveness of GDMPD

The effectiveness of PSMU(V) and PSMU++ reveals the backdoor created by incorporating visual information from external sources. In this paper, we propose a safe way to utilize images from untrustworthy sources through GDMPD, which can purify images and detect adversarial images. In this section, we conduct experiments to demonstrate the effectiveness of our GDMPD. First, we show the purification effectiveness, followed by the accuracy of GDMPD detection.
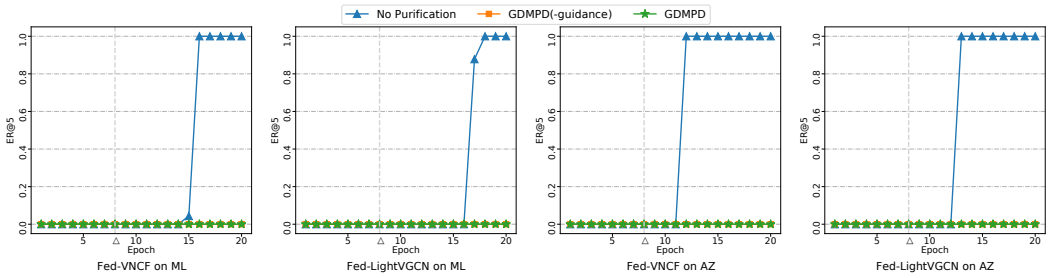


Fig. 7. The trend of exposure rate of PSMU++ for FedRecs with purification mechanism.

Table 2. The comparison of recommendation performance (NDCG@20) for different purified methods in visually-aware FedRecs.

| FedRec Model | Purification Method | ML | AZ |
|---|---|---|---|
| Fed-VNCF | original | 0.03985 | 0.02786 |
| | GDMPD(-guidance) | 0.04833 | 0.02806 |
| | GDMPD | **0.05032** | **0.02849** |
| Fed-LightVGCN | original | 0.03831 | 0.02065 |
| | GDMPD(-guidance) | 0.04494 | 0.02135 |
| | GDMPD | **0.04623** | **0.02246** |

Table 3. The standard deviation of Blur and Brisque scores of images with different purification methods. Lower values indicate that the image quality difference is less.

| Purification Method | ML | | AZ | |
|---|---|---|---|---|
| | Blur | Brisque | Blur | Brisque |
| original | 2444.68 | 11.82 | 1682.55 | 18.138 |
| GDMPD(-guidance) | 821.41 | 8.83 | 552.84 | 16.80 |
| GDMPD | **801.02** | **8.62** | **523.75** | **16.59** |

Table 4. The effectiveness of GDMPD for defending against PSMU++ with different perturbation scales. The value of each cell $(x, y)$ represents: $x$ is the highest ER@5 scores that PSMU++ achieves in FedRecs without defense method, and $y$ is the highest ER@5 scores that PSMU++ achieves when equipped with GDMPD.

| Dataset | Fed-VNCF | | | | Fed-VLightGCN | | | |
|---|---|---|---|---|---|---|---|---|
| $(x, y)$ | $\epsilon = 4$ | $\epsilon = 8$ | $\epsilon = 16$ | $\epsilon = 32$ | $\epsilon = 4$ | $\epsilon = 8$ | $\epsilon = 16$ | $\epsilon = 32$ |
| **ML** | (1.0, 0.0) | (1.0, 0.0) | (1.0, 0.0) | (1.0, 0.0) | (1.0, 0.0) | (1.0, 0.0) | (1.0, 0.0) | (1.0, 0.0) |
| **AZ** | (1.0, 0.0) | (1.0, 0.0) | (1.0, 0.0) | (1.0, 0.0) | (1.0, 0.0) | (1.0, 0.0) | (1.0, 0.0) | (1.0, 0.0) |

To safely use external images, we incorporate purification mechanisms in visually-aware FedRecs. Fig. 7 shows PSMU++'s attack results for different purification methods. GDMPD(-guidance) is GDMPD without diffused image guidance. In Fig. 7, all purification methods can reduce the attack's ER@5 to 0, which demonstrates that by adding Gaussian noises during the diffusion process, the perturbations have been diluted. An effective defense method should not only prevent the attacker's achieving its malicious goals, but also consume less recommendation performance. Table 2 presents the recommendation performance of visually-aware FedRecs with different purification methods. "original" is the visually-aware FedRecs that leverage original images. According to the results in Table 2, FedRecs with purified images even have better performance than using original images. This is because the diffusion model can shrink the variance of original images, where images are provided by different providers and the quality of them is different. Table 3 provides a proof-of-concept. We calculate the standard deviation of Brisque [33] and Blur scores (Laplacian operator value) to evaluate the quality deviation of images generated by different methods. In Table 3, the standard deviations of original images in both ML and AZ are much higher than purified images, indicating the original images' large quality difference.

When adding more perturbations, the attacker's goal will be easier to achieve. Therefore, we evaluate our defense method's effectiveness with increased $\epsilon$ from 4 to 32 in Table 4. The results show that before utilizing our purification mechanism, the attack's ER@5 can reach to 1.0 with all different $\epsilon$ scales. However, when incorporating our defense methods, the scores of ER@5 reduce to 0.0 in all cases, which implies that our defense methods can at least tolerate PSMU++ with $\epsilon < 32$. Fig. 6 presents a case study of adversarial images with different perturbation scales purified by different methods. The comparison of purified images and adversarial images with large-scale perturbations (e.g., $\epsilon = 16$ or 32) shows that the diffusion model can remove abnormal noise. Furthermore, by comparing the images generated by our GDMPD and GDMPD(-guidance), we can see that after adding guidance, the generated images are more consistent with the original ones, demonstrating the effectiveness of our guidance.

To improve the maintenance of recommender systems, we have implemented a detection function in GDMPD that is based on its purification ability. Since our detector is training-free, the detection results for ML and AZ are almost identical, so we present the overall results for the union set of ML and AZ adversarial images. Specifically, we use PSMU++ to generate 100 images that can promote cold items to 1.0 on ML and AZ respectively. Then, we tested whether GDMPD can detect these adversarial images when mixed with all other normal item images. As described in E.q. 27, the accuracy of adversarial image detection mainly depends on the setting of $\rho$. In our experiments, we set $\rho$ as follows. First, we randomly sample a subset (10, 000 in our experiments) of images from a publicly available image dataset, ImageNet. These images are normal images and we use our GDMPD to purify them. After purification, we calculate the cosine similarity between purified and original images. Finally, we naively use the minimal value of the similarity scores as $\rho$. If the original image from FedRecs has smaller scores than $\rho$ with its corresponding purified
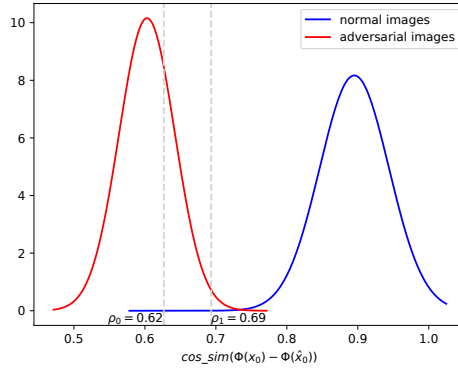
Fig. 8. Visualization of the density distribution of normal and adversarial images similarity scores fitting with normal distribution. $\rho_0$ is the minimal similarity score from normal images. $\rho_1$ is the best setting of $\rho$ that can filter out all adversarial images.

image, GDMPD will mark it as an adversarial image. Based on this setting, we get $0.72$ accuracy for detecting adversarial images and no normal images are falsely predicted as adversarial images. To further analyze our detection, we visualize the distribution of normal images' similarity scores and adversarial images' similarity scores in Fig. 8. $\rho_1 = 0.69$ is the "best" setting of $\rho$ that can achieve $1.0$ accuracy to detect adversarial images from FedRecs but we cannot directly get $\rho_1$ in practice since we have limited prior knowledge of adversarial images. In this paper, we simply set $\rho$ to $\rho_0$ according to the minimal value of normal images' similarity scores. How to estimate a better $\rho$ can be explored in future work.

## 7 CONCLUSION

Recently, numerous studies have exposed the threat of model poisoning attacks on federated recommender systems (FedRecs) that rely on collaborative data. We argue that these attacks are effective due to the sparsity of user-item interactions in the data. In this paper, we propose the incorporation of visual information to alleviate the data sparsity problem and demonstrate that existing model poisoning attacks cannot easily promote target items in visually-aware FedRecs. Subsequently, we propose PSMU(V) image poisoning attacks that exploit the newly created backdoor in visually-aware FedRecs. These attacks can work in tandem with model poisoning attacks, posing a greater threat and highlighting the need for a secure visual information usage mechanism. To address this gap, we propose a novel image poisoning defender based on DDPM that can not only purify adversarial images but also detect them. Extensive experiments conducted on two real-world datasets using two visually-aware FedRecs demonstrate the effectiveness of our proposed attacks and defenses.

## REFERENCES

[1] Zareen Alamgir, Farwa K Khan, and Saira Karim. 2022. Federated recommenders: methods, challenges and future. *Cluster Computing* 25, 6 (2022), 4075–4096.

[2] Muhammad Ammad-Ud-Din, Elena Ivannikova, Suleiman A Khan, Were Oyomno, Qiang Fu, Kuan Eeik Tan, and Adrian Flanagan. 2019. Federated collaborative filtering for privacy-preserving personalized recommendation system. *arXiv preprint arXiv:1901.09888* (2019).

[3] Di Chai, Leye Wang, Kai Chen, and Qiang Yang. 2020. Secure federated matrix factorization. *IEEE Intelligent Systems* 36, 5 (2020), 11–20.

[4] Fei Chen, Mi Luo, Zhenhua Dong, Zhenguo Li, and Xiuqiang He. 2018. Federated meta-learning with fast convergence and efficient communication. *arXiv preprint arXiv:1802.07876* (2018).

[5] Tong Chen, Hongzhi Yin, Guanhua Ye, Zi Huang, Yang Wang, and Meng Wang. 2020. Try this instead: Personalized and interpretable substitute recommendation. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*. 891–900.

[6] Rami Cohen, Oren Sar Shalom, Dietmar Jannach, and Amihood Amir. 2021. A black-box attack model for visually-aware recommender systems. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 94–102.

[7] Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems* 34 (2021), 8780–8794.

[8] Tommaso Di Noia, Daniele Malitesta, and Felice Antonio Merra. 2020. Taamr: Targeted adversarial attack against multimedia recommender systems. In *2020 50th Annual IEEE/IFIP international conference on dependable systems and networks workshops (DSN-W)*. IEEE, 1–8.

[9] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative adversarial networks. *Commun. ACM* 63, 11 (2020), 139–144.

[10] Yeting Guo, Fang Liu, Zhiping Cai, Hui Zeng, Li Chen, Tongqing Zhou, and Nong Xiao. 2021. PREFER: Point-of-interest REcommendation with efficiency and privacy-preservation via Federated Edge leaRning. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 5, 1 (2021), 1–25.

[11] Elizabeth Liz Harding, Jarno J Vanto, Reece Clark, L Hannah Ji, and Sara C Ainsworth. 2019. Understanding the scope and impact of the california consumer privacy act of 2018. *Journal of Data Protection & Privacy* 2, 3 (2019), 234–253.

[12] F Maxwell Harper and Joseph A Konstan. 2015. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)* 5, 4 (2015), 1–19.

[13] Ruining He and Julian McAuley. 2016. VBPR: visual bayesian personalized ranking from implicit feedback. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 30.

[14] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*. 639–648.

[15] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *Proceedings of the 26th international conference on world wide web*. 173–182.

[16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851.

[17] Nguyen Quoc Viet Hung, Huynh Huu Viet, Nguyen Thanh Tam, Matthias Weidlich, Hongzhi Yin, and Xiaofang Zhou. 2017. Computing crowd consensus with partial agreement. *IEEE Transactions on Knowledge and Data Engineering* 30, 1 (2017), 1–14.

[18] Mubashir Imran, Hongzhi Yin, Tong Chen, Quoc Viet Hung Nguyen, Alexander Zhou, and Kai Zheng. 2023. ReFRS: Resource-efficient federated recommender system for dynamic and diversified user preferences. *ACM Transactions on Information Systems* 41, 3 (2023), 1–30.

[19] Vignesh Jagadeesh, Robinson Piramuthu, Anurag Bhardwaj, Wei Di, and Neel Sundaresan. 2014. Large scale visual recommendations from street fashion images. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. 1925–1934.

[20] Yannis Kalantidis, Lyndon Kennedy, and Li-Jia Li. 2013. Getting the look: clothing recognition and segmentation for automatic product suggestions in everyday photos. In *Proceedings of the 3rd ACM conference on International conference on multimedia retrieval*. 105–112.

[21] Wang-Cheng Kang, Chen Fang, Zhaowen Wang, and Julian McAuley. 2017. Visually-aware fashion recommendation and design with generative image models. In *2017 IEEE international conference on data mining (ICDM)*. IEEE, 207–216.

[22] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[23] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).

[24] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. 2020. Federated learning: Challenges, methods, and future directions. *IEEE signal processing magazine* 37, 3 (2020), 50–60.

[25] Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. 2022. Diffusion-lm improves controllable text generation. *Advances in Neural Information Processing Systems* 35 (2022), 4328–4343.

[26] Feng Liang, Weike Pan, and Zhong Ming. 2021. Fedrec++: Lossless federated recommendation with explicit feedback. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 35. 4224–4231.

[27] Guanyu Lin, Feng Liang, Weike Pan, and Zhong Ming. 2020. Fedrec: Federated recommendation with explicit feedback. *IEEE Intelligent Systems* 36, 5 (2020), 21–30.

[28] Qiang Liu, Shu Wu, and Liang Wang. 2017. Deepstyle: Learning user preferences for visual recommendation. In *Proceedings of the 40th international acm sigir conference on research and development in information retrieval*. 841–844.

[29] Zhuoran Liu and Martha Larson. 2021. Adversarial item promotion: Vulnerabilities at the core of top-n recommenders that use images to address cold start. In *Proceedings of the Web Conference 2021*. 3590–3602.

[30] Zhiwei Liu, Liangwei Yang, Ziwei Fan, Hao Peng, and Philip S Yu. 2022. Federated social recommendation with graph neural network. *ACM Transactions on Intelligent Systems and Technology (TIST)* 13, 4 (2022), 1–24.

[31] Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. 2015. Image-based recommendations on styles and substitutes. In *Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval*. 43–52.

[32] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*. PMLR, 1273–1282.

[33] Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. 2012. No-reference image quality assessment in the spatial domain. *IEEE Transactions on image processing* 21, 12 (2012), 4695–4708.

[34] Khalil Muhammad, Qinqin Wang, Diarmuid O'Reilly-Morgan, Elias Tragos, Barry Smyth, Neil Hurley, James Geraci, and Aonghus Lawlor. 2020. Fedfast: Going beyond average for faster training of federated recommender systems. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*. 1234–1242.

[35] Alexander Quinn Nichol and Prafulla Dhariwal. 2021. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*. PMLR, 8162–8171.

[36] Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Anima Anandkumar. 2022. Diffusion models for adversarial purification. *arXiv preprint arXiv:2205.07460* (2022).

[37] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019).

[38] Tao Qi, Fangzhao Wu, Chuhan Wu, Yongfeng Huang, and Xing Xie. 2020. Privacy-preserving news recommendation model learning. *arXiv preprint arXiv:2003.09592* (2020).

[39] Ruihong Qiu, Zi Huang, Hongzhi Yin, and Zijian Wang. 2022. Contrastive learning for representation degeneration problem in sequential recommendation. In *Proceedings of the fifteenth ACM international conference on web search and data mining*. 813–823.

[40] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2012. BPR: Bayesian personalized ranking from implicit feedback. *arXiv preprint arXiv:1205.2618* (2012).

[41] Dazhong Rong, Qinming He, and Jianhai Chen. 2022. Poisoning Deep Learning based Recommender Model in Federated Learning Scenarios. *arXiv preprint arXiv:2204.13594* (2022).

[42] Dazhong Rong, Shuai Ye, Ruoyan Zhao, Hon Ning Yuen, Jianhai Chen, and Qinming He. 2022. Fedrecattack: Model poisoning attack to federated recommendation. In *2022 IEEE 38th International Conference on Data Engineering (ICDE)*. IEEE, 2643–2655.

[43] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. 2008. The graph neural network model. *IEEE transactions on neural networks* 20, 1 (2008), 61–80.

[44] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).

[45] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*. PMLR, 2256–2265.

[46] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. 2021. Maximum likelihood training of score-based diffusion models. *Advances in Neural Information Processing Systems* 34 (2021), 1415–1428.

[47] Zehua Sun, Yonghui Xu, Yong Liu, Wei He, Yali Jiang, Fangzhao Wu, and Lizhen Cui. 2022. A Survey on Federated Recommendation Systems. *arXiv preprint arXiv:2301.00767* (2022).

[48] Jinhui Tang, Xiaoyu Du, Xiangnan He, Fajie Yuan, Qi Tian, and Tat-Seng Chua. 2019. Adversarial training towards robust multimedia recommender system. *IEEE Transactions on Knowledge and Data Engineering* 32, 5 (2019), 855–867.

[49] Paul Voigt and Axel Von dem Bussche. 2017. The eu general data protection regulation (gdpr). *A Practical Guide, 1st Ed., Cham: Springer International Publishing* 10, 3152676 (2017), 10–5555.

[50] Jinyi Wang, Zhaoyang Lyu, Dahua Lin, Bo Dai, and Hongfei Fu. 2022. Guided diffusion model for adversarial purification. *arXiv preprint arXiv:2205.14969* (2022).

[51] Jingdong Wang, Ting Zhang, Nicu Sebe, Heng Tao Shen, et al. 2017. A survey on learning to hash. *IEEE transactions on pattern analysis and machine intelligence* 40, 4 (2017), 769–790.

[52] Qinyong Wang, Hongzhi Yin, Tong Chen, Junliang Yu, Alexander Zhou, and Xiangliang Zhang. 2021. Fast-adapting and privacy-preserving federated recommender system. *The VLDB Journal* (2021), 1–20.

[53] Kangning Wei, Jinghua Huang, and Shaohong Fu. 2007. A survey of e-commerce recommender systems. In *2007 international conference on service systems and service management*. IEEE, 1–5.

[54] Chuhan Wu, Fangzhao Wu, Yang Cao, Yongfeng Huang, and Xing Xie. 2021. Fedgnn: Federated graph neural network for privacy-preserving recommendation. *arXiv preprint arXiv:2102.04925* (2021).

[55] Chuhan Wu, Fangzhao Wu, Tao Qi, Yongfeng Huang, and Xing Xie. 2022. FedAttack: Effective and covert poisoning attack on federated recommendation via hard sampling. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 4164–4172.

[56] Chuhan Wu, Fangzhao Wu, Tao Qi, Yongfeng Huang, and Xing Xie. 2022. Fedcl: Federated contrastive learning for privacy-preserving recommendation. *arXiv preprint arXiv:2204.09850* (2022).

[57] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Yingxia Shao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. 2022. Diffusion models: A comprehensive survey of methods and applications. *arXiv preprint arXiv:2209.00796* (2022).

[58] Hongzhi Yin and Bin Cui. 2016. *Spatio-temporal recommendation in social media*. Springer.

[59] Hongzhi Yin, Bin Cui, Zi Huang, Weiqing Wang, Xian Wu, and Xiaofang Zhou. 2015. Joint modeling of users' interests and mobility patterns for point-of-interest recommendation. In *Proceedings of the 23rd ACM international conference on Multimedia*. 819–822.

[60] Junliang Yu, Hongzhi Yin, Xin Xia, Tong Chen, Jundong Li, and Zi Huang. 2022. Self-supervised learning for recommender systems: A survey. *arXiv preprint arXiv:2203.15876* (2022).

[61] Yang Yu, Qi Liu, Likang Wu, Runlong Yu, Sanshi Lei Yu, and Zaixi Zhang. 2022. Untargeted Attack against Federated Recommendation Systems via Poisonous Item Embeddings and the Defense. *arXiv preprint arXiv:2212.05399* (2022).

[62] Wei Yuan, Quoc Viet Hung Nguyen, Tieke He, Liang Chen, and Hongzhi Yin. 2023. Manipulating Federated Recommender Systems: Poisoning with Synthetic Users and Its Countermeasures. *arXiv preprint arXiv:2304.03054* (2023).

[63] Wei Yuan, Chaoqun Yang, Quoc Viet Hung Nguyen, Lizhen Cui, Tieke He, and Hongzhi Yin. 2023. Interaction-level Membership Inference Attack Against Federated Recommender Systems. *arXiv preprint arXiv:2301.10964* (2023).

[64] Wei Yuan, Hongzhi Yin, Fangzhao Wu, Shijie Zhang, Tieke He, and Hao Wang. 2023. Federated unlearning for on-device recommendation. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*. 393–401.

[65] Honglei Zhang, Fangyuan Luo, Jun Wu, Xiangnan He, and Yidong Li. 2022. LightFR: Lightweight Federated Recommendation with Privacy-preserving Matrix Factorization. *ACM Transactions on Information Systems* (2022).

[66] Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. 2019. Deep learning based recommender system: A survey and new perspectives. *ACM computing surveys (CSUR)* 52, 1 (2019), 1–38.

[67] Shijie Zhang and Hongzhi Yin. 2022. Comprehensive Privacy Analysis on Federated Recommender System against Attribute Inference Attacks. *arXiv preprint arXiv:2205.11857* (2022).

[68] Shijie Zhang, Hongzhi Yin, Tong Chen, Zi Huang, Quoc Viet Hung Nguyen, and Lizhen Cui. 2022. Pipattack: Poisoning federated recommender systems for manipulating item promotion. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*. 1415–1423.

[69] Bolong Zheng, Kai Zheng, Xiaokui Xiao, Han Su, Hongzhi Yin, Xiaofang Zhou, and Guohui Li. 2016. Keyword-aware continuous knn query on road networks. In *2016 IEEE 32Nd international conference on data engineering (ICDE)*. IEEE, 871–882.

[70] Ruiqi Zheng, Liang Qu, Bin Cui, Yuhui Shi, and Hongzhi Yin. 2023. AutoML for Deep Recommender Systems: A Survey. *ACM Transactions on Information Systems* (2023).