

# Securing Visually-Aware Recommender Systems: An Adversarial Image Reconstruction and Detection Framework

MINGLEI YIN, West Virginia University, USA

BIN LIU, West Virginia University, USA

NEIL ZHENQIANG GONG, Duke University, USA

XIN LI, State University of New York at Albany, USA

With rich visual data, such as images, becoming readily associated with items, visually-aware recommendation systems (VARS) have been widely used in different applications. Recent studies have shown that VARS are vulnerable to item-image adversarial attacks, which add human-imperceptible perturbations to the clean images associated with those items. Attacks on VARS pose new security challenges to a wide range of applications such as e-Commerce and social networks where VARS are widely used. How to secure VARS from such adversarial attacks becomes a critical problem. Currently, there is still a lack of systematic study on how to design secure defense strategies against visual attacks on VARS. In this paper, we attempt to fill this gap by proposing an *adversarial image reconstruction and detection* framework to secure VARS. Our proposed method can simultaneously (1) secure VARS from adversarial attacks characterized by *local* perturbations by image reconstruction based on *global* vision transformers; and (2) accurately detect adversarial examples using a novel contrastive learning approach. Meanwhile, our framework is designed to be used as both a filter and a detector so that they can be *jointly* trained to improve the flexibility of our defense strategy to a variety of attacks and VARS models. We have conducted extensive experimental studies with two popular attack methods (FGSM and PGD). Our experimental results on two real-world datasets show that our defense strategy against visual attacks is effective and outperforms existing methods on different attacks. Moreover, our method can detect adversarial examples with high accuracy.

CCS Concepts: • **Information systems** → **Recommender systems**; • **Security and privacy** → *Web application security*.

Additional Key Words and Phrases: Recommendation systems, visual features, adversarial machine learning, attack detection, contrastive learning

## ACM Reference Format:

Minglei Yin, Bin Liu, Neil Zhenqiang Gong, and Xin Li. 2023. Securing Visually-Aware Recommender Systems: An Adversarial Image Reconstruction and Detection Framework. 1, 1 (June 2023), 25 pages. <https://doi.org/10.1145/1122445.1122456>

## 1 INTRODUCTION

In the era of data explosion, people often face overwhelming information overload problems. Recommender systems play an important role in helping users find the information they are interested in more easily [1]. As of today, recommender systems have become essential components in a wide range of Internet services—from E-commerce to social networks—to help users deal with

---

Authors' addresses: Minglei Yin, my0033@mix.wvu.edu, West Virginia University, USA; Bin Liu, bin.liu1@mail.wvu.edu, West Virginia University, USA; Neil Zhenqiang Gong, neil.gong@duke.edu, Duke University, USA; Xin Li, xin.li@ieee.org, State University of New York at Albany, USA.

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2022 Association for Computing Machinery.

XXXX-XXXX/2023/6-ART \$15.00

<https://doi.org/10.1145/1122445.1122456>

information overload, engage users, and improve user experience. In a typical recommender system setting, we are given a set of users, a set of items, and a record of the users’ historical interactions (e.g., ratings, likes, or clicks) with the items, and our goal is to model user preferences from the user-item interactions and then recommend each user a list of new items that the users have not experienced yet. Meanwhile, rich visual data, such as images, are becoming more widely associated with items in various web services such as e-commerce sites (e.g., Amazon), image sharing and social media services (e.g., Pinterest) and online marketplaces (e.g., Airbnb). The emerging field of visually aware recommender systems (VARS) [2, 20, 28] is expected to have a significant impact in several application domains such as e-Commerce [20, 28], image sharing and social networks [44, 65], fashion [24], food [12] and real estate and tourism [48]. As the cliché says, “a picture is worth a thousand words” – the visual appearance of a product image (e.g., the picture of an outfit or an apartment) could affect the final decision of an online consumer [19, 20].

Recently, some studies have shown that visually aware recommender systems (VARS) are vulnerable to item image adversarial attacks, which add human-imperceptible perturbations to clean images associated with these items [8, 36, 51]. Tang et al. [51] found that a small but intentional perturbation in the input image will severely decrease the accuracy of the recommendation, implying the vulnerability of VARS to untargeted attacks. A black-box attack model for VARS was studied [8] and it shows that the visual attack model can effectively influence the preference scores and classifications of items without knowing the parameters of the model to promote the push of certain items. Liu et al. [36] studied adversarial item promotion attacks at VARS in the top N item generation stage in the cold start setting. Although there exist many studies [6, 14, 25, 32, 34, 35, 42, 49, 53, 60, 61, 69] on attacks on general recommendation systems that manipulate user-item interaction data in different ways, attacks against VARS are different in that they only manipulate images associated with items. The security aspect of VARS systems is much less explored.

Attacks on VARS pose new security challenges to a wide range of applications such as e-Commerce and social networks where VARS are widely used. However, to our knowledge, there has been a limited amount of research on defending VARS from adversarial attacks. Tang et al. [51] proposed to apply adversarial training to improve the robustness of VARS. Anelli et al. [3] conducted a study on the effectiveness of adversarial training methods to improve the robustness of VARS to different adversarial image manipulations including the Fast Gradient Sign Method (FGSM) [17], Projected Gradient Descent (PGD) [31], and Carlini & Wagner (CW) [5] attack. Despite the fact that there are various studies on attacks and defenses on vision learning systems [16, 64], there is still no systematic study on the defense of VARS against increasingly more powerful attacks. Generally speaking, *robust model construction* and *attack detection* are two popular strategies for defending against attacks against recommender systems [10]. Robust model construction approaches, such as robust statistics-based methods [38] and more recently adversarial training-based methods [3, 22, 51, 54, 63], aim to design recommender systems proactively that are more secure against attacks. Attack detection approaches aim to detect malicious user profiles [4, 33, 39, 56] and recover from attacks.

In this paper, we propose a novel *adversarial image reconstruction and detection* framework to protect visually aware recommender systems from adversarial attacks that manipulate images associated with items. The proposed framework is designed to simultaneously (1) mitigate the impacts of adversarial attacks on recommendation performance and (2) detect adversarial attacks. Specifically, the adversarial image reconstruction component reconstructs clean item images by denoising the adversarial perturbations in attacked images. To this end, we develop an image reconstruction network composed of several residual blocks and a global vision transformer [45]. Since adversarial attacks such as FGSM and PGD are local perturbations, the proposed transformer-based global filtering strategy can effectively make the reconstructed images as close to clean images

as possible, alleviating the adversarial impact on the performance of VARS models. Moreover, we also design a novel contrastive learning-based component to accurately detect adversarial examples. A novel strategy is proposed to construct positive and negative pairs for contrastive learning. By pushing the reconstructed images toward clean images and away from adversarial ones, our detection network can detect adversarial examples with high accuracy. Furthermore, our framework is designed to be a flexible filter and detector plug-in, which can defend against the adversarial attack without modifying the original recommender system. We build our defense framework on Visual Bayesian Personalized Ranking (VBPR) [20], the most popular VARS model, as our baseline model.

We evaluate our defense framework with two popular attack methods, *i.e.*, FGSM and PGD. Our experimental results on two real-world datasets (Amazon Men and Amazon Fashion) show that our framework can effectively defend against visual attacks and outperforms existing methods on different attacks. Furthermore, our method can detect adversarial examples with high accuracy.

A summary of key contributions is listed below.

- We provide a systematic study on securing visually-aware recommender systems by unifying two defense strategies—robust model construction and attack detection—from adversarial attacks that manipulate images associated with items.
- We design a novel *adversarial image reconstruction and detection* framework to simultaneously mitigate the impacts of adversarial attacks and detect the attacks. We demonstrate that end-to-end training of reconstruction and detection networks can significantly improve the robustness of VARS to a variety of adversarial attacks.
- Extensive experiments on two real-world datasets demonstrate that the proposed framework can effectively defend the recommender system model from attacked images with varying strengths. Moreover, our proposed framework can detect adversarial examples with high accuracy.

## 2 RELATED WORK

We organize the related work into two categories: security of general recommender systems, and attacks and defenses in visually-aware recommender systems.

### 2.1 Security of General Recommender Systems

In the era of data explosion, recommender systems play an important role in reducing information overload by helping users find the information they are interested in more easily [1]. In a *general recommender system* setting, we are given a set of users, a set of items, and a record of the users' historical interactions (e.g., ratings, likes, or clicks) with the items, and our goal is to model user preferences from the user-item interactions and then recommend to each user a list of new items that the users have not experienced yet. Many algorithms have been developed for recommender systems in recent decades, for example, neighborhood-based [47], graph-based [15], matrix-factorization-based [30, 40], and deep-learning-based [23, 55, 67].

Recommendation systems have been successfully applied in different applications and many methods have been proposed to improve recommendation performance, the security aspect of recommender systems is much less explored but has received increasingly more attention in recent years [10]. Due to the nature of openness, where user-item interaction data is used to train a recommendation system, a body of studies has shown that recommender systems are vulnerable to various adversarial attacks [6, 14, 25, 32, 34, 35, 42, 49, 53, 60, 61, 69]. By injecting carefully crafted fake data into a recommender system, an attacker can spoof the recommender system to recommend attacker-chosen items or arbitrary incorrect items to many genuine users. Depending on the intent

of the attack, attacks on recommender systems can be categorized into *targeted attacks* with the goal of promoting or demoting targeted items and *untargeted attacks* whose goal is to dysfunction a recommender system. To attack such general recommender systems, attackers manipulate user-item interaction data to deceive the recommendation model. There are two categories of attacks. The first category is *data poisoning attacks* (a.k.a. shilling attacks), in which attackers inject fake users with carefully crafted ratings into a recommender system [6, 13, 14, 25, 32, 34, 42, 49, 53, 61, 69]. As such, different attacks have been proposed to different types of recommender system models, e.g., neighborhood-based [6], graph-based [14], matrix factorization-based [13, 34], deep learning based [25], and graph neural network-based [42, 61] recommender systems. Meanwhile, model-agnostic data poisoning attacks to recommender systems have been studied [49, 53, 69], i.e., the attacks are designed without knowing specific recommendation models used by the recommender system. For example, when prior knowledge about the target recommender system is unavailable or limited, reinforcement learning-based attacks [49], surrogate model-based attacks [69], and generative adversarial network-based [53] have been proposed. The second category of attack to recommender systems is *profile pollution attacks* with the goal of polluting users' profiles [35, 60] to perturb the results of a recommender system. For example, Xing *et al.* [60] studied user profile pollution by generating false item clicks through cross-site request forgery (XSRF) and demonstrated the compromised user profiles leaning to promotion of certain items on three platforms—YouTube, Google, and Amazon; Liu *et al.* [35] studied the impacts of profile pollution attacks on recommender systems under factorization machines (FMs) framework and showed that FMs were vulnerable to such adversarial perturbations.

*Robust model construction* and *attack detection* are two major strategies to defend against attacks on recommendation systems [10]. The first strategy is to proactively design robust recommender systems so that they are more secure against attacks. Along this line, Mehta *et al.* [38] proposed a robust matrix factorization approach to attack-resistant collaborative filtering for recommender systems by leveraging robust statistics such as M-estimators. More recently adversarial training [3, 22, 54, 63] has been applied to improve the robustness of recommender systems. The basic idea of adversarial training [17] is to train a recommender system model on a training dataset that is augmented with adversarial examples so that the adversarially trained recommendation model is resistant to adversarial attacks. For example, He *et al.* [22] proposed an adversarial personalized ranking framework that applied adversarial training to the widely used Bayesian Personalized Ranking (BPR) model [46] by introducing adversarial perturbations in the embedding vectors of users and items. Yuan *et al.* [63] studied adversarial training on collaborative denoising auto-encoder recommendation model. Unlike the adversarial training framework in [22, 63] that add perturbations to the model parameters (e.g., embedding vectors of users and items), Wu *et al.* [54] proposed an adversarial poisoning training method to counteract data poisoning attacks to recommender systems. The second strategy, the attack detection-based method, aims to detect malicious user profiles and then remove compromised user profiles in the data processing stage. Attack detection-based method assumes that malicious users and genuine users have different user-item interaction patterns. Different attack detection methods have been proposed, including classification [4, 41] by extracting attributes derived from user profiles and items, semi-supervised learning [56], and unsupervised learning such as clustering in the user-item rating matrix [33, 39] and graph-based methods [62] on user-item graph. Zhang *et al.* [68] proposed a method of unifying the robust recommendation task and fraudster detection task by combining a graph convolutional network (GCN) model to predict user-item ratings and a neural random forest model to predict the mean square of all ratings per user. They assumed that if a user's rating is largely deviated from the predicted ratings, this user is most likely to be a fraudster.



## 2.2 Attacks and Defenses in Visually-Aware Recommender Systems

With rich visual data, such as images, becoming readily associated with items, visually aware recommendation systems (VARS) have been widely used to improve users’ decision-making process and support online sales [19, 20, 28]. Typically, VARS models [20, 28] first extract image features using deep neural networks and then combine the extracted features with existing recommendation models [9]. Visual Bayesian Personalized Ranking (VBPR) [20] is one of the most widely used models in VARS. Specifically, VBPR uses neural networks such as pre-trained CNN such as *ResNet50* [18] to extract image features from images associated with items and then fuses the extracted features into the widely used Bayesian personalized ranking (BPR) model [46]. More recently, visual-aware deep Bayesian personalized ranking (DVBPR) [28] was developed to simultaneously extract task-guided visual features and learn user latent factors, leading to improved recommendation performance by directly learning “fashion-aware” image representations through joint training of image representation and recommender systems.

Recent studies have shown that VARS are vulnerable to item image adversarial attacks, which only add human-imperceptible perturbations to clean images associated with those items [8, 11, 36, 51]. Tang et al. [51] studied the vulnerability of VARS in an untargeted attack environment and found that a small but intentional perturbation in the input image will severely decrease the precision of the recommendation. Noia *et al.* [11] proposed a targeted attack to VARS, and formulated the attack goal as to spoof the recommender system to misclassify the images of a category of low recommended products towards the class of more recommended products. Since the item catalogs of VARS are usually large, recommender systems often count on the item providers to provide images as supplementary information. This reliance on external sources has inspired the design of a black-box attack on VARS in [8]. An attacker was shown to unfairly promote targeted items by modifying the item scores and pushing their rankings. By systematically creating human-imperceptible perturbations of the images of the pushed item, the attackers manage to incrementally increase the item score. Liu et al. [36] studied adversarial item promotion attacks in VARS in the top-N item generation stage under the cold-start recommendation setting. Note that attacks on VARS are different from existing work on attacks against general recommendation systems [6, 13, 14, 25, 32, 34, 42, 49, 53, 61]. While studies on attacks against general recommender systems focus on manipulating user-item records to achieve attack goals, attacks against VARS only manipulate images associated with items.

There is a limited amount of research on the defenses against item image adversarial attacks to VARS. After demonstrating the vulnerability of VARS in an untargeted attack environment, Tang et al. [51] proposed to apply adversarial training to improve the robustness of VARS. Specifically, adversarial perturbations were applied to the item image’s deep feature vector and the adversarial training was formulated as an adversarial regularizer added to the BPR loss. Anelli et al. [3] conducted a study on the effectiveness of adversarial training methods to improve the robustness of VARS to different adversarial image manipulations, including the Fast Gradient Sign Method (FGSM) [17], Projected Gradient Descent (PGD) [31], and Carlini & Wagner (CW) [5] attack. They assumed that adversaries were aware of recommendation lists, and then formulated the adversarial attack as to misclassify the images of a category of low-recommended products towards the class of more recommended products. Instead of focusing on the item-level raking performances of recommender systems, models were evaluated in terms of the fraction of compromised items in the top-N recommendations.

Our work is different from previous research [3, 51] on defenses against item image adversarial attacks to VARS in several different ways. First, both our work and [51] consider untargeted attacks to VARS with the goal of dysfunctioning the recommender system; however, the ways to generate

adversarial samples are different. In [51] perturbations are added to model parameters (*i.e.*, item image’s deep feature vector), but our attack is to add human-imperceptible perturbations to clean images associated with items. Second, our attack model is also different from the targeted attack in [3], which assumes that the adversaries are aware of the recommendation lists, and adversarial samples are generated by attacking an image classifier to misclassify an item from low-rank category to high-rank category. Third, from the defense perspective, different from [3, 51] which only applies adversarial training to boost the robustness of the recommendation, we propose an adversarial image reconstruction and detection framework in this paper, which not only can secure VARS from adversarial attacks via image reconstruction, but also can accurately detect adversarial examples. We note that orthogonal to VARS, attacks and their defenses have also been studied in various image classification algorithms in the computer vision community [16, 64].

### 3 PROPOSED METHOD

In this section, we first introduce item image adversarial attacks to visually aware recommendation systems in this study. We then introduce our proposed framework, which takes into account both *robust model construction* and *attack detection* into account, to defend against such attacks and elaborate details about our proposed defense method.

#### 3.1 Preliminaries and Problem Formulation

Our goal is to secure visually-aware recommendation systems (VARS) from adversarial attacks. Specifically, we consider a typical VARS setting where we have a set of  $M$  users  $\mathcal{U} = \{1, 2, \dots, M\}$  and a set of  $N$  items  $\mathcal{I} = \{1, 2, \dots, N\}$ , and we are given a record of user-item interactions  $\mathcal{D} = \{\langle u, i, r_{ui} \rangle\}$ , where  $r_{ui}$  denotes the preference score of user  $u$  for item  $i$ . We assume that an image  $x_i$  is associated with each item  $i$ , and the images are denoted as  $\mathcal{X} = \{x_i\}_{i \in \mathcal{I}}$ . In this paper, without loss of generality, we focus on the widely used VARS model, Visual Bayesian Personalized Ranking (VBPR) [20], as our baseline model, although the methodology could be generalized to other VARS. To avoid overfitting, following [51] we define the user preference score of user  $u$  for the item  $i$  as  $r_{ui} = \gamma_u^T (\gamma_i + \mathbf{E}f_i)$ , where  $\gamma_u$  and  $\gamma_i$  are latent factors of the user and the item, respectively,  $f_i$  is the vector of visual features extracted by a pre-trained CNN such as *ResNet50* [18] from the image  $x_i$  associated with the item  $i$ , and  $\mathbf{E}$  is a transformation matrix that maps the visual feature  $f_i$  into the latent factor space of the item. Then the task of VARS is to build a model  $r_{ui} = F(u, i, x_i | \Theta)$ , where  $u \in \mathcal{U}$ ,  $i \in \mathcal{I}$ , and  $\Theta$  indicate the model parameters, to infer the preference scores for the items that the users have not yet experienced. Recommendations are made based on the inferred preference scores. We use Bayesian Personalized Ranking (BPR) optimization framework, which is a pairwise ranking loss, to train the VBPR model. Specifically, we first construct a training dataset  $\mathcal{D}_s$  as follows:

$$\mathcal{D}_s = \{(u, i, j, x_i, x_j) \mid u \in \mathcal{U} \wedge i \in \mathcal{I}_u^+ \wedge j \in \mathcal{I} / \mathcal{I}_u^+\} \quad (1)$$

where,  $\mathcal{I}_u^+$  represents the set of interacted items of user  $u$ . The triplet  $(u, i, j)$  indicates that the user  $u$  prefers item  $i$  over item  $j$ .  $x_i$  and  $x_j$  are images of items  $i$  and  $j$ . BPR loss is defined as follows:

$$\mathcal{L}_{BPR} = \underset{\Theta}{\operatorname{argmin}} \left\{ - \sum_{(u,i,j) \in \mathcal{D}_s} \ln \sigma(r_{ui} - r_{uj}) + \lambda \|\Theta\|^2 \right\}, \quad (2)$$

where  $\sigma(\cdot)$  is the Sigmoid function and  $\lambda$  is a regularization hyperparameter.

We consider untargeted attacks to VARS with the goal of dysfunctioning the recommender system. The objective of adversarial attacks to VARS is to disrupt the output of VARS, for example, to alter the ranking of user recommendations, by adding human-imperceptible perturbations  $\delta_i$  to clean images  $x_i$  associated with each item  $i \in \mathcal{I}$ , namely,  $x_i^* = x_i + \delta_i$ . We consider a white-box

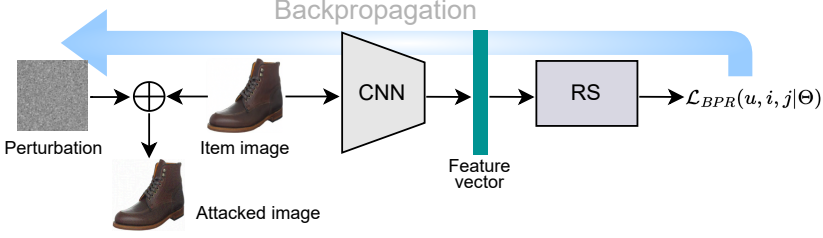


Fig. 1. Illustration of adversary example generation in attacks to visually-aware recommender systems. Adversary examples are generated based on the VBPR model [20], which includes a pre-trained CNN (e.g., ResNet50) for image feature extraction and a latent factor recommendation model for user preference prediction. The whole pipeline is differentiable. Therefore, the attack signal can be computed by the backpropagation of the loss value.

attack setting in which we assume that an attacker has access to the user-item interaction data  $\mathcal{D}$  and the images  $x_i$  associated with the item  $i$ , the recommender system model (the VBPR model in this work) and the neural network architecture to derive image features. Specifically, when the user  $u$  prefers item  $i$  over item  $j$ , then adversarial attacks aim to alter preference scores  $F(u, i, x_i | \Theta)$  and  $F(u, j, x_j | \Theta)$  so that  $F(u, i, x_i | \Theta) < F(u, j, x_j | \Theta)$ . Meanwhile, to make perturbations imperceptible, the attacked image  $x_i^*$  ( $x_j^*$ ) should be similar to the clean image  $x_i$  ( $x_j$ ). The generation of adversarial samples can be modeled as a constrained minimization problem:

$$\begin{aligned} & \text{minimize } \|x_i^* - x_i\|_2^2 + \|x_j^* - x_j\|_2^2 \\ & \quad + F(u, i, x_i | \Theta) - F(u, j, x_j | \Theta) \\ & \text{such that } x_i^*, x_j^* \in [0, 1]^n \end{aligned} \quad (3)$$

where  $i$  and  $j$  indicate positive and negative samples in the VBPR model, respectively. The first part  $\|x_i^* - x_i\|_2^2 + \|x_j^* - x_j\|_2^2$  aims to keep the attacked images  $x_i^*$  and  $x_j^*$  to stay close to the original images  $x_i$  and  $x_j$  in pixel space, respectively, the second part  $F(u, i, x_i | \Theta) - F(u, j, x_j | \Theta)$  disrupts the output of VARS by flipping the preference scores, and  $x_i^*, x_j^* \in [0, 1]^n$  are image value range constraints to keep the perturbed images within the original image value range.

Given the adversarial attack goal as shown in Equation (3), different attack methods, such as FGSM [17] and PGD [31], can be used to generate adversarial samples, and the perturbed images are denoted as  $\mathcal{X}^* = \{x_i^*\}_{i \in \mathcal{I}}$ . Specifically, the preference score flipping part  $F(u, i, x_i | \Theta) - F(u, j, x_j | \Theta)$  in Equation (3) can be achieved by maximizing the BPR loss function  $\mathcal{L}_{BPR}$  as defined in Equation (2), namely,  $\arg\max_{\Theta} \{-\sum_{(u,i,j) \in \mathcal{D}_s} \ln \sigma(r_{ui} - r_{uj}) + \lambda \|\Theta\|^2\}$ . Figure 1 illustrates the procedure of generating adversary examples in attacks to visually-aware recommender systems, and different attack methods, such as FGSM and PGD, can be applied.

- **Fast Gradient Sign Method (FGSM)** [17] has the advantage of generating adversarial perturbations quickly over other methods. It needs only one step to generate the attacked image. Given a clean input image  $x$ , the adversarial example can be computed by a local perturbation computed by the fast gradient sign method.

$$x^* = x + \epsilon \text{ sign}(\nabla_x \mathcal{L}_{BPR}) \quad (4)$$

where  $\epsilon$  corresponds to the magnitude of adversarial signals,  $\text{sign}(\cdot)$  is the sign function and  $\nabla_x \mathcal{L}_{BPR}$  is the gradient of the attack loss function.

- **Projected Gradient Descent (PGD)** [31] is an iterative version of FGSM. The attack algorithm iterates FGSM with a smaller step size. After each completed perturbation step, the intermediate attacked image is clipped to a  $\epsilon$ -neighborhood of the original image  $x$ .

Our goal is then to design an effective defense strategy that can simultaneously (1) mitigate the impacts of adversarial attacks on recommendation performance and (2) detect adversarial attacks. We assume that we have the user-item interaction data  $\mathcal{D} = \{\langle u, i, r_{ui} \rangle\}$ , the images associated with the items  $\mathcal{X} = \{x_i\}_{i \in \mathcal{I}}$ , and the perturbed images  $\mathcal{X}^* = \{x_i^*\}_{i \in \mathcal{I}}$  generated according to Equation (3). Given data  $\{\mathcal{D}, \mathcal{X}, \mathcal{X}^*\}$ , our objective is to build a reconstruction network to reconstruct image items for boosting the robustness of the recommendation, and a detection network to detect adversarial samples. To our knowledge, this formulation of joint *robust model construction* and *attack detection* for VARS has not been considered in the open literature.

### 3.2 Overview of Proposed Adversarial Image Reconstruction and Detection Framework

To secure visually-aware recommendation systems (VARS) from adversarial attacks, we propose a framework that takes into account the robustness of the model to adversarial attack and the detection of adversarial examples. Meanwhile, our framework is designed to be used as a filter and detector prior to the recommendation system model. Therefore, our model can be trained first and then used as a pre-processing step without affecting the architecture and parameters of the current recommender system. Figure 2 shows the entire pipeline of our framework, which can be divided into three parts: the *reconstruction network* to denoise adversarial signals, the *detection network* to detect adversarial examples, and the *recommendation system* (RS) model to predict the final preference scores of the items for each user.

For the reconstruction part, we first sample the inputs  $(u, i, j)$  from the training set  $\mathcal{D}_s = \{(u, i, j)\}$ , where the triplet  $(u, i, j)$  indicates that the user  $u$  prefers item  $i$  over item  $j$ . Beyond the user-item interaction records, we have item images, clear images  $\mathcal{X} = \{x_i\}_{i \in \mathcal{I}}$  or perturbed images  $\mathcal{X}^* = \{x_i^*\}_{i \in \mathcal{I}}$ , associated with the items. In our work, the triplet  $(u, i, j)$  is used as input of the recommendation model, and the corresponding image of the items is input of our proposed defense model. We assume that the input images can be clean ( $\mathcal{X} = \{x_i\}_{i \in \mathcal{I}}$ ) or perturbed ( $\mathcal{X}^* = \{x_i^*\}_{i \in \mathcal{I}}$ ) via adversarial perturbation. Our objective is to denoise/remove the adversarial signals while preserving the clean images. Thus, we devised a reconstruction network in our framework to reconstruct the images from their perturbed images to increase the robustness of the recommendation. The reconstruction network  $T(\cdot)$  is a neural network that transforms an input image, a clean image  $x$  or a perturbed image  $x^*$ , into a reconstructed image  $\hat{x}$ , namely  $\hat{x} = T(x^*)$  for a perturbed image and  $\hat{x} = T(x)$  for a clean image. To train the reconstruction network, we take a pair of clean and perturbed images  $\langle x_i, x_i^* \rangle_{i \in \mathcal{I}}$  as input, and apply the perceptual loss to enforce the reconstructed image  $\hat{x}_i = T(x_i^*)$  is similar to the corresponding clean image  $x_i$ . The trained reconstruction network is then used to reconstruct item images whether they have been attacked or not, and we hope that the reconstructed images can be as clean as possible so that the image features extracted through vision models (e.g., ResNet50) would not have a negative effect on the following recommender system models.

For the detection part, based on the reconstruction network, we compare the reconstructed image with the input image. If they are similar, the input image is originally clean; if they are different, the input image might have experienced adversarial perturbation. Then the question boils down to a reliable evaluation of the similarity or differences between two images. Along this line of reasoning, we turn to metric learning [29], which aims to automatically construct domain-specific distance metrics from the training dataset. We then use the learned distance metric for other tasks,

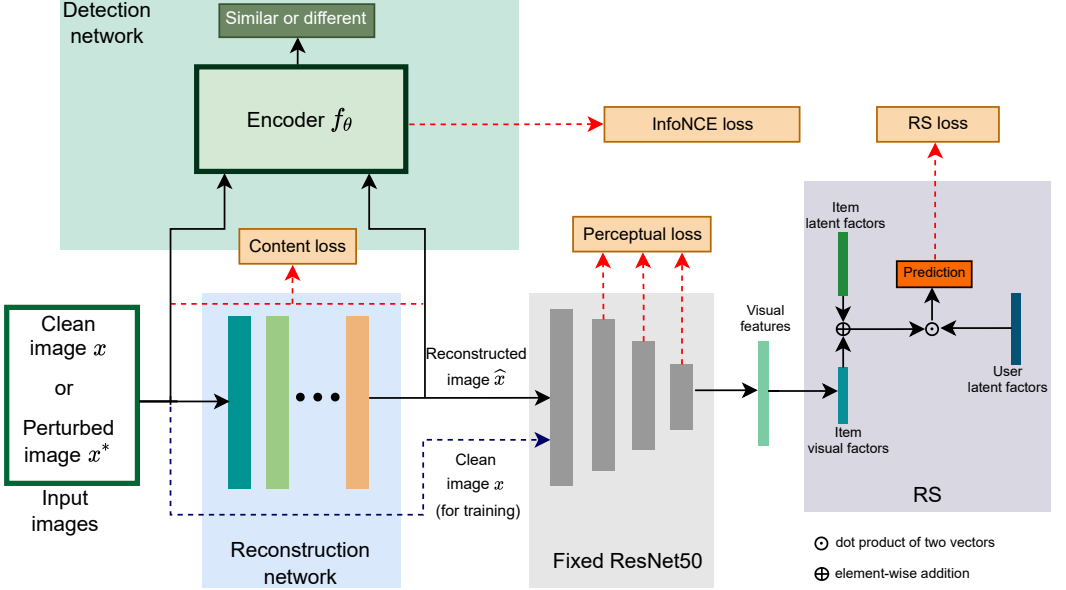


Fig. 2. The pipeline of our proposed framework. The defense architecture mainly consists of two parts: (1) the reconstruction network to reconstruct item images whether the input images have been attacked or not to boost the recommendation robustness. Note that pairs of clean and perturbed images  $\langle x, x^* \rangle$  are needed as input to train the reconstruction network; and (2) the detection network for detecting adversarial samples. Other parts in this figure are basically the same as the components of the VBPR model, which includes a pre-trained network (ResNet50 in this work) for image feature extraction and a latent factor model for user preference prediction.

such as adversarial detection. In our case, the distance metric is used as a judgment of similarity; our detection network is based on the distance obtained by metric learning.

A notable feature of the proposed system is the *joint training* of the detection and reconstruction networks. Note that the objectives of reconstruction and detection are mutually beneficial to each other. Reconstruction can help detection by filtering out adversarial attack signals; on the other hand, detection can facilitate reconstruction by pushing clean images away from noisy ones. Joint optimization of the reconstruction and detection modules in an end-to-end manner allows them to interact with each other for improved generalization performance, as will be experimentally verified later. To our knowledge, such end-to-end optimization of detection and reconstruction modules has not been proposed before. In the next three sections, we will first elaborate on the reconstruction network in Sec. 3.3 and the detection network in Sec. 3.4. Then we present the total loss function for joint reconstruction and detection in Sec. 3.5.

### 3.3 Image Reconstruction based on Residual and Transformer Blocks

The objective of image reconstruction is to generate clean images through denoising adversarial perturbations that might impair the performance of the recommendation model. Note that VARS models, for example, VBPR model, first extract image features from item images using a convolutional neural network (CNN) and then fuse the extracted features into existing recommendation models such as the widely used BPR model. The reason for the degradation of recommendation performance is that the introduced adversarial signals can activate semantically irrelevant regions

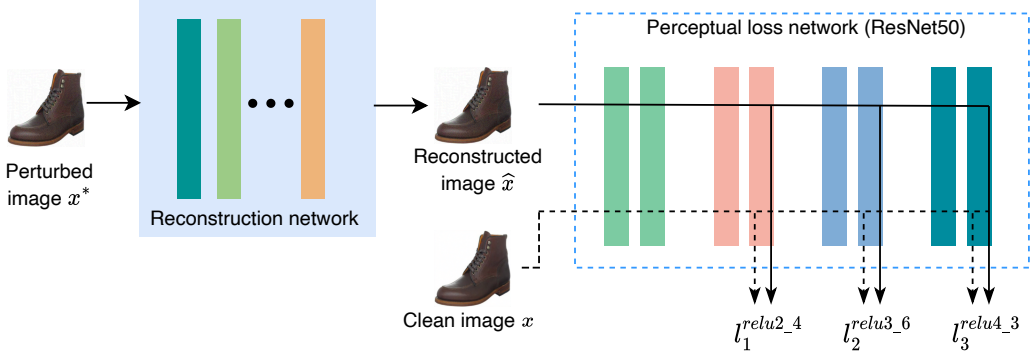


Fig. 3. Illustration of the reconstruction network. The reconstruction network transforms a perturbed image  $x^*$  into a reconstructed image  $\hat{x} = T(x^*)$ . Perceptual loss that measures the feature differences between the clean image  $x$  and reconstructed image  $\hat{x}$  at different intermediate feature maps is used to supervise the training of the reconstruction network.

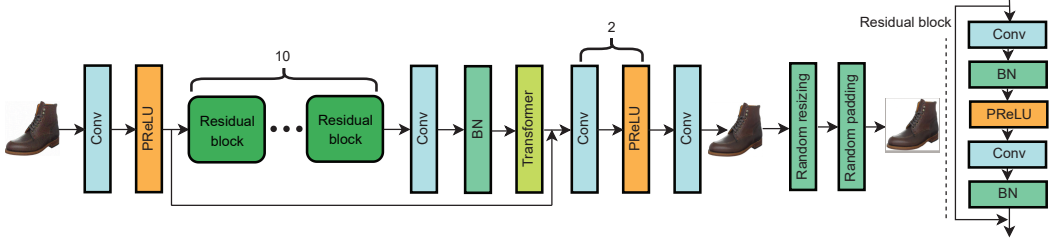


Fig. 4. The reconstruction network architecture. It primarily consists of a series of 10 residual blocks (dark green) and transform blocks (light green) along with the basic convolution operations with kernel size 3 (except for the first and last one whose kernel size is 9). For the activation function, we have adopted both parametric ReLU (PReLU) and Batch Normalization (BN). Note that we can append one random layer that contains random resizing and padding after the final layer at inference time.

in the intermediate feature maps, which interfere with the features of innocent signals that focus on task-related content in CNN networks (e.g., ResNet50) used to extract image features from item images [59]. Furthermore, perturbation signals will be amplified in the high-level features of the CNN network, although they may be human-imperceptible. Operation through the layers of the network can gradually increase the magnitude of the perturbation of misleading characteristics, which negatively affects the performance of the recommendation model.

Now the problem of defense against adversarial attacks boils down to how to denoise the intermediate feature maps generated by the layers of the neural network used for image extraction. Inspired by [59, 66], we devised a reconstruction network to reconstruct the images from their perturbed images to increase the robustness of the recommendation. The reconstruction network  $T(\cdot)$  is a neural network that transforms a perturbed image  $x^*$  into a reconstructed image  $\hat{x}$  through  $\hat{x} = T(x^*)$ . As shown in Figure 3, to train the reconstruction network, we take pairs of clean and perturbed images  $\langle x, x^* \rangle$  as input and apply perceptual loss to ensure that the reconstructed image  $\hat{x}_i = T(x^*)$  is similar to the corresponding clean image  $x$ .

The primary components of the reconstruction network are shown in Figure 4. The residual blocks are the core components of the image reconstruction network. Because the goal of the

reconstruction network is to denoise adversarial perturbations, a residual block is a good choice as it excels in learning the difference (or the residual) between input and output. Residual blocks are able to keep features for the clean image that focus primarily on semantically informative content in the image, and remove feature maps for the adversarial image that are activated across semantically irrelevant regions. Furthermore, recent advances in visual recognition (e.g., bottleneck transformer [50]) have shown that it is better to replace spatial convolutions with *global* self-attention in the final bottleneck block of the reconstruction network for the task of visual recognition. Based on the observation that adversarial attacks are local perturbations, we advocate the inclusion of a global vision transformer block in image reconstruction. The entire network consists of nine residual blocks followed by a transformer block and some separated convolution operations. The first and last convolutions are equipped with a kernel size of 9, and the others with 3. To generate clean reconstructed images, the entire pipeline does not involve any downsampling operators. Following [58], in this work, random operations at inference time are also used to mitigate the adversarial effect. Specifically, two random operations are imposed on the reconstructed images. The reconstructed image of size  $224 \times 224$  is first randomly resized to a smaller image  $n \times n, n \in [212, 224]$ , and then the resized image is randomly padded with zeros to the size of  $224 \times 224$  with zeros.

We use the perceptual loss of the image [27] to supervise the task of denoising feature maps. As shown in Figure 3, the perceptual loss function is calculated by comparing high-level differences based on intermediate features extracted from pre-trained networks. Note that the perturbation signals will be amplified in the high-level features of the CNN network. We expect the intermediate features derived from the reconstructed images to be as close to the features extracted from clean images as possible. Specifically, given a pre-trained neural network  $\phi$ , let  $\phi_l(x)$  be the features of the  $l$ -th convolution layer of the network  $\phi$ ,  $x$  and let  $\hat{x}$  be the input image and the reconstructed image. The perceptual loss in the  $l$ -th convolution layer can be defined as

$$\mathcal{L}_{perc}^l = \sum_{x \in \mathcal{X}} \text{Dist}(\phi_l(x), \phi_l(\hat{x})) \quad (5)$$

where  $\text{Dist}(\cdot)$  denotes the  $L_2$  distance function. As studied in [9], ResNet50 has shown quantitatively and qualitatively to produce the best recommended products. Therefore, we use ResNet50 [18] as our pre-trained network for the computation of perceptual loss. In particular, we use the outputs of stages 2, 3, and 4 of the ResNet50, denoted as  $l_1^{\text{relu2\_4}}$ ,  $l_2^{\text{relu3\_6}}$  and  $l_3^{\text{relu4\_3}}$ , to retrieve the high-level features to compute the distance differences. Then the perceptual loss  $\mathcal{L}_{perc}$  is the sum of the perceptual losses in the layers of  $l_1^{\text{relu2\_4}}$ ,  $l_2^{\text{relu3\_6}}$  and  $l_3^{\text{relu4\_3}}$ .

In addition to perceptual loss, in the reconstruction network, we also use content loss to supervise the image generation task. The content loss is defined as the  $L_2$  of the different between the input image  $x$  and the constructed image  $\hat{x}$ ,

$$\mathcal{L}_{pix} = \sum_{x \in \mathcal{X}} \|x - \hat{x}\|_2^2. \quad (6)$$

Content loss is employed to obtain a slightly blurred image (i.e., filter out adversarial perturbation), while perceptual loss can also help the reconstruction model preserve rich details in an image.

### 3.4 Attack Detection Based on Contrastive Learning

Beyond building a more robust model, attack detection is also an important strategy to defend against attacks on recommendation systems. In particular, the purpose of attack detection in this paper is to distinguish an innocent image from a malicious image with adversarial perturbations (adversarial example), which can degrade the performance of the recommendation model. We propose a contrastive learning-based approach to the detection of adversarial examples. Contrastive

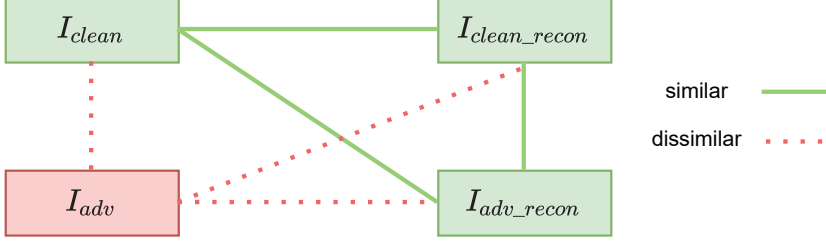


Fig. 5. Positive and negative pairs for training the detection network via contrastive learning. The pair connected by a green line means positive pair (expected to be pulled close), and the pair connected by the red dot line means negative pair (expected to be pushed away).

learning [7] aims to learn representations of images that push away dissimilar examples and bring similar examples closer. As shown in Figure 2, we will reconstruct the image by  $\hat{x} = T(x^*)$  when the item image is perturbed. We can detect the perturbed image  $x^*$  if its representation is away from the representation of the reconstructed image  $\hat{x}$ . Similarly, when the item image is innocent, namely, the input is a clean image  $x$ , we would expect its representation to be close to that of the reconstructed image  $\hat{x} = T(x)$ .

To this end, we build the detection network  $f_\theta(\cdot)$  which includes two parts: a neural network encoder  $g(\cdot)$  to extract image features and a small projection head  $h(\cdot)$  to project the extracted image space to a common embedding space for attack detection purpose. In particular, we use ResNet50 for  $g(\cdot)$  and a multi-layer perceptron (MLP) for  $h(\cdot)$ . Without causing ambiguity, let  $x$  be the input image for the detection network<sup>1</sup>, then we have its image representation in the embedding space as  $z = f_\theta(x)$ . Specifically, we have:

$$\begin{aligned} s &= g(x) = \text{ResNet50}(x) \\ z &= h(s) = W^{(2)}\sigma(W^{(1)}s), \end{aligned} \tag{7}$$

where  $s \in \mathbb{R}^{2048}$  is the result of the average pooling layer in ResNet,  $\sigma$  indicates the nonlinear activation function of ReLU and  $W^{(1)}$ ,  $W^{(2)}$  are the weights in the MLP. The final output  $z \in \mathbb{R}^{128}$  is the feature vector in the embedding space, and will be used to calculate the similarity to the other. The feature vectors of clean images will be pulled close in the embedding space, while those of adversarial images will be pushed away.

We then train the encoder with contrastive learning, which encourages the clustering of different classes of samples around their centroids. Specifically, for each item with a corresponding image  $x$ , we first construct its positive and negative pairs as training data samples for encoder training  $f_\theta$  based on inputs  $I_{clean}$ ,  $I_{adv}$  and the corresponding outputs  $I_{clean\_recon}$ ,  $I_{adv\_recon}$  from the reconstruction network. The feature vectors of the pairs are then placed close to each other in the embedding space if two samples in the pair are similar; otherwise, the feature vectors of dissimilar samples in the pair are separated by a large distance from each other. Figure 5 shows the construction of positive and negative pairs. Ideally, we assume that our reconstruction network will produce clean images without adversarial perturbations for both natural and adversarial examples. For each example  $I_{clean}$ , we will have a corresponding adversarial version  $I_{adv}$ . When our reconstruction network is completed, the reconstructed versions  $I_{clean\_recon}$  and  $I_{adv\_recon}$  will

<sup>1</sup>Note that the input image for the detection network can be a perturbed image  $x^*$  or a clean image  $x$  along with the reconstructed image  $\hat{x} = T(x^*)$  or  $\hat{x} = T(x)$  respectively.



be generated. Based on the above observation, we expect that different distances will be learned between similar and dissimilar pairs. Positive and negative pair sets can be constructed as follows.

$$S_{pos} = \left\{ \begin{array}{l} (I_{clean}, I_{clean\_recon}), \\ (I_{clean}, I_{adv\_recon}), \\ (I_{clean\_recon}, I_{adv\_recon}) \end{array} \right\} \quad S_{neg} = \left\{ \begin{array}{l} (I_{adv}, I_{clean}), \\ (I_{adv}, I_{adv\_recon}), \\ (I_{clean\_recon}, I_{adv}) \end{array} \right\} \quad (8)$$

The similarity of two samples in the pair can be measured by cosine similarity. Let  $sim(\mathbf{u}, \mathbf{v}) = \mathbf{u}^T \mathbf{v} / \|\mathbf{u}\| \|\mathbf{v}\|$  denote cosine similarity. For training the encoder  $f_\theta$ , we propose to minimize the following contrastive loss similar to the InfoNCE loss [43]:

$$\mathcal{L}(x) = -\log \frac{\sum_{x_i, x_j \in S_{pos}} \exp(sim(f_\theta(x_i), f_\theta(x_j))/\tau)}{\sum_{x_i, x_j \in S} \exp(sim(f_\theta(x_i), f_\theta(x_j))/\tau)} \quad (9)$$

where  $x_i, x_j$  are derived from item image  $x$ ,  $S = S_{pos} \cup S_{neg}$ , and  $\tau$  is a temperature hyperparameter [57]. Given a dataset  $X$ , the total contrastive loss will be calculated as

$$\mathcal{L}_{contr} = \mathbb{E}_{x \sim X} [\mathcal{L}(x)]. \quad (10)$$

After obtaining a well-trained feature encoder  $f_\theta$ , its output can be used to calculate the similarity between two samples. Let  $I_{in}$  (perturbed image  $x^*$  or clean image  $x$ ) and  $I_{out}$  (reconstructed image  $\hat{x} = T(x^*)$  or  $\hat{x} = T(x)$ ) be the input and output of our reconstruction network. The feature vectors can be obtained by  $\mathbf{z}_{in} = f_\theta(I_{in})$  and  $\mathbf{z}_{out} = f_\theta(I_{out})$ . Then we use the Euclidean distance between  $\mathbf{z}_{in}$  and  $\mathbf{z}_{out}$  as the dissimilarity score. A high score indicates that  $I_{in}$  and  $I_{out}$  are dissimilar, which implies that  $I_{in}$  is an adversarial example because  $I_{out}$  is purified after going through our reconstruction network. In contrast, a small score indicates that  $I_{in}$  and  $I_{out}$  are similar and both are likely to be clean examples. The threshold for the decision boundary is set to 0.2 in our experiments.

### 3.5 Total Loss Functions for Joint Reconstruction and Detection

As shown in Figure 2, there are four kinds of loss function in total in our framework. The detection network is optimized by contrastive loss, InfoNCE loss [43] in our case, which helps the detection network to learn an embedding function that maps the image  $\mathbf{x}$  to the feature  $\mathbf{z}$ . This embedding will pull the two clean samples close and push a clean sample and an adversarial sample away in the embedding space. For image quality enhancement and removal of adversarial signals, we equip the reconstruction network with content loss and perceptual loss. For the recommender system, we use the same loss function as VBPR because we use VBPR as our recommendation model in the backend. The loss function of VBPR is defined as [20]:

$$\mathcal{L}_{rs} = \argmin_{\Theta} \sum_{(u,i,j) \in \mathcal{D}_s} -\ln \sigma(\hat{r}_{u,i} - \hat{r}_{u,j}) + \lambda_\Theta \|\Theta\|^2 \quad (11)$$

where  $(u, i, j)$  is the sampled pair of items of user  $u$ ,  $\hat{r}_{u,i}$  is the preference score of user  $u$  for item  $i$ ,  $\Theta$  represents all parameters of the model,  $\lambda_\Theta$  is the weight of the regularization term and  $\sigma$  is the Sigmoid function. To avoid overfitting [51], we define the user preference score  $u$  for item  $i$  as  $\hat{r}_{u,i} = \gamma_u^T (\gamma_i + \mathbf{E} \mathbf{f}_i)$ , where  $\gamma_u$  and  $\gamma_i$  are vectors  $K$ -dimensional (e.g.,  $K = 64$ ) that represent the latent factors of the user  $u$  and the item  $i$ , respectively. Furthermore,  $\mathbf{f}_i \in \mathbb{R}^{2048}$  is a feature vector of the item  $i$  extracted by a pre-trained neural network, and  $\mathbf{E} \in \mathbb{R}^{K \times 2048}$  is an embedding matrix that maps  $\mathbf{f}_i$  into a  $K$  dimensional latent space.

In summary, we add these loss functions together to jointly optimize our reconstruction and detection networks as

$$\mathcal{L} = \mathcal{L}_{pix} + \alpha \mathcal{L}_{perc} + \beta \mathcal{L}_{contr} + \gamma \mathcal{L}_{rs}, \quad (12)$$

where  $\alpha, \beta, \gamma$  are hyperparameters. It should be noted that the detection and reconstruction networks with their corresponding loss are optimized jointly. Based on our constructed multiple types of positive and negative pairs, the detection network can provide useful feedback signals to help the reconstruction network suppress adversarial perturbations and vice versa.

## 4 EXPERIMENTS

In this section, we report our experimental results on two real-world datasets to show the effectiveness of our proposed framework on improving the robustness to adversarial attacks and the detection of adversarial examples.

### 4.1 Datasets

Our experiments are conducted on two real-world datasets: Amazon Men and Amazon Fashion, both of which are derived from the Amazon Web store [37]. The original dataset contains over 180 million relationships among almost 6 million objects, which are the result of recording the product recommendations of more than 20 million users. The visual features of these two datasets have been shown to provide meaningful information for recommendation. User review histories are viewed as implicit feedback, which can be used to sample the triplet  $(u, i, j)$ , each item paired with an image to extract visual features. The statistics of these two datasets are summarized in Table 1. For data processing, we first convert each user’s review into a binary-valued interaction vector, and then filter out cold users by a rule of thumb:  $|\mathcal{I}_u^+| < 5$ . Following the same evaluation protocol [51], we have used a standard leave-one-out method to generate the test sets. That is, we randomly select one interaction vector for testing and use the remaining data for training.

Table 1. Statistics of the datasets we used in this work.

Dataset	User#	Item#	Interaction#
Amazon Men	34,244	110,636	254,870
Amazon Fashion	45,184	166,270	358,003

### 4.2 Evaluation Metrics

Since our recommender system generates the top- $N$  list based on the computed preference scores, we use the Hit Ratio (HR) and Normalized Discounted Cumulative Gain (NDCG) [51] to measure the quality of the top- $N$  lists generated.

- **Hit Ratio (HR@ $N$ ):** given the top- $N$  recommendation list, we check if the groundtruth item is in the list. If yes, we mark 1 for this user and 0 otherwise.
- **Normalized Discounted Cumulative Gain (NDCG@ $N$ ):** given the top- $N$  recommendation list, we consider the rank position of the groundtruth item in the list. The score decrease as the groundtruth’s rank goes lower.

The difference between these two metrics is that HR only considers whether the recommended item exists in the top- $N$  list, while NDCG further takes into account the position of the recommended item in the top- $N$  list. Since it is time-consuming to rank all items for every user during the evaluation, we have followed the common strategy [11, 23] that randomly samples 20 items with which the user does not interact, ranking the test item among these 20 items. More specifically,

we first generate a list of size 20 for each user in the test set. The generated list contains items that do not have any interactions with the user. Then we predict the preference scores with the recommender system model. Based on the computed preference scores, we rank the items interacted with in the list and calculate the HR and NDCG metrics. We separately calculate the metrics when  $N$  is set to 5, 10, and 20, respectively. For the detection part, we measure the detection performance according to the classification accuracy.

### 4.3 Training Procedures

As shown in Figure 2, the complete pipeline consists of three parts: the detection network, the reconstruction network, and the recommender system (RS). There are three hyperparameters in the loss function as shown in Equation (12), in our experiments, we choose  $\alpha = 1.0$ ,  $\beta = 100.0$ , and  $\gamma = 0.1$ . To optimize these three components in a more computationally efficient manner, we opt to process them separately before fine-tuning. Specifically, we first train our model on the Amazon Men dataset. We first spent about 100 epochs training the VBPR model, which can be used as the baseline RS model by both AMR [51] and our framework. The learning rate when training the baseline VBPR model is set to  $1e - 3$  in the first 80 epochs and decreases to  $1e - 4$  in the last 20 epochs. We train the VBPR model in the same way on two datasets to produce a well-trained model. After having a pre-trained VBPR model, we spent another five epochs fine-tuning the model with a learning rate of  $1e - 4$  in an adversarial training way as described in [51] to obtain the well-trained AMR model.

Based on the pre-trained VBPR model, we then build our framework and train the detection network and reconstruction network, respectively. We first link our reconstruction network with the pre-trained VBPR model. To enforce the reconstruction network’s learning how to recover clean images and remove adversarial signals from perturbed images, we have fixed the parameters of the VBPR model, i.e., there is no updating of the VBPR model when training the reconstruction network. During training, we do not activate the random layers of the reconstruction network because it does not help the training. We use them only in the testing phase to disturb the adversarial signals. To speed up reconstruction network training, we only use clean input and adversarial input generated by the FGSM method. Under this condition, it takes about 10 epochs with a learning rate of  $1e - 5$  to complete the training process. Then it will take another three epochs to retrain our model with clean input, FGSM-adversarial input, and PGD-adversarial input, respectively. The proportion of these three inputs is controlled to be close to 1 : 1 : 1. We use the random number generator to determine which types of input should be used in one iteration. The learning rate is maintained at  $1e - 5$ .

The well-trained reconstruction network can be transferred to another dataset by fine-tuning for other small-number (e.g., 2-3) epochs. Finally, we train the detection network based on the well-trained reconstruction network. It takes about two epochs to complete the training. When training the detection network, the adversarial inputs we used are generated by the FGSM method with  $\epsilon = 16$  only. As shown in our analysis (Table 3) later, the accuracy of the detection network is robust to different choices of  $\epsilon$  values. The proportion of clean and adversarial inputs is close to 1 : 1. We note that the well-trained detection network can be effortlessly transferred to other datasets without extra fine-tuning.

### 4.4 Defense Performance and Analysis

Our reconstruction network first strives to generate clean images regardless of the input images, whether they are clean or adversarial. The reconstructed images will then be fed into the following pre-trained neural network, which extracts visual features to be used by the recommender system models for predicting user preferences. To gain a deeper understanding, we will analyze the defense

Table 2. Recommendation performance on two real-world datasets. HR and NDCG metrics are used for comparison (the higher, the better). The clean suffix in the leftmost column means that all models run on clean data without adversarial perturbations. And rand indicates that the reconstructed images by our model are transformed with two random operations. The attack perturbation level  $\epsilon$  is set to 16 in this table.

Metrics	HR@N			NDCG@N		
	N = 5	N = 10	N = 20	N = 5	N = 10	N = 20
Amazon Men						
VBPR-Clean	0.5074	0.7127	0.9806	0.3495	0.4155	0.4830
AMR-Clean	0.5113	0.7123	0.9799	0.3519	0.4166	0.4840
OURS-Clean	0.5069	0.7129	0.9806	0.3488	0.4146	0.4825
(rand)OURS-Clean	0.5022	0.7110	0.9805	0.3461	0.4131	0.4806
VBPR-FGSM	0.0159	0.0465	0.3804	0.0087	0.0183	0.0978
AMR-FGSM	0.0150	0.0426	0.3576	0.0082	0.0170	0.0918
OURS-FGSM	0.4580	0.6669	0.9730	0.3119	0.3791	0.4555
(rand)OURS-FGSM	0.4550	0.6646	0.9736	0.3058	0.3730	0.4510
VBPR-PGD	0.0000	0.0000	0.0030	0.0000	0.0000	0.0007
AMR-PGD	0.0000	0.0000	0.0031	0.0000	0.0000	0.0008
OURS-PGD	0.0512	0.1449	0.6590	0.0268	0.0562	0.1806
(rand)OURS-PGD	0.2406	0.4420	0.9098	0.1402	0.2043	0.3206
Amazon Fashion						
VBPR-Clean	0.5504	0.7602	0.9845	0.3784	0.4462	0.5029
AMR-Clean	0.5508	0.7606	0.9851	0.3800	0.4476	0.5044
OURS-Clean	0.5487	0.7598	0.9846	0.3780	0.4459	0.5027
(rand)OURS-Clean	0.5412	0.7578	0.9860	0.3713	0.4406	0.4990
VBPR-FGSM	0.0042	0.0188	0.3013	0.0022	0.0068	0.0736
AMR-FGSM	0.0042	0.0178	0.2839	0.0023	0.0065	0.0694
OURS-FGSM	0.4949	0.7170	0.9782	0.3308	0.4025	0.4683
(rand)OURS-FGSM	0.4791	0.7036	0.9792	0.3161	0.3888	0.4586
VBPR-PGD	0.0000	0.0000	0.0040	0.0000	0.0000	0.0009
AMR-PGD	0.0000	0.0000	0.0044	0.0000	0.0000	0.0011
OURS-PGD	0.0216	0.0917	0.5818	0.0106	0.0326	0.1510
(rand)OURS-PGD	0.2063	0.4311	0.9043	0.1137	0.1859	0.3033

performance of our reconstruction network from the following perspectives: *Recommendation performance* and *transferability study*. We also present the visual quality of the reconstructed images.

**Recommendation performance evaluation.** Table 2 shows the quantitative results of our defense model in terms of HR@N and NDCG@N metrics on Amazon Men and Amazon Fashion datasets [21]. Similarly to AMR [51], the performance of our defense model can be compared to VBPR (baseline), which shows that the new reconstruction network does not affect the existing RS model in clean data. As expected, both the performance of VBPR and AMR drop dramatically in the presence of FGSM and PGD attacks on the RS model. A plausible reason for the failure of AMR is that adversarial perturbations are added to the positive image  $i$  and the negative image  $j$  based on the input sample triplet  $(u, i, j)$ . Perturbations on two-item images can not only decrease the rank position of the positive item, but also increase the possibility of recommending a negative item. In contrast, our defense model appears to be more robust to adversarial attacks than both the VBPR and AMR baselines. The reconstruction network can partially remove adversarial signals from

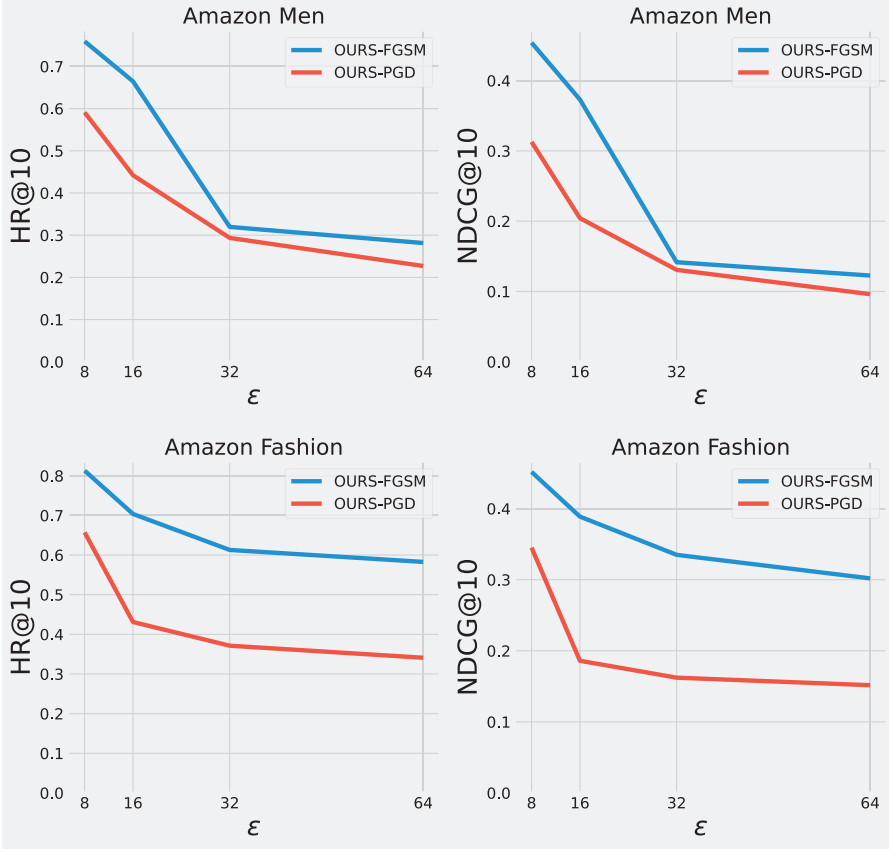


Fig. 6. Performance of the reconstruction network w.r.t. different values of attack perturbation level  $\epsilon$  (our model is trained with  $\epsilon = 16$  only).

the input images and effectively defend the RS model from FGSM attacks. When faced with more aggressive PGD attacks, our model alone becomes inadequate; however, with random operations (marked by “(rand)”), it still achieves decent performance in challenging PGD attacks. We conclude that our proposed reconstruction network is an effective defense strategy against adversarial attacks without altering the original recommendation model.

**Transferability study.** We want to demonstrate the transferability performance of our defense model when faced with attack signals with different strengths. Transferability is a desirable property for machine learning algorithms. As shown in the first column of Figure 6, attack signals with a smaller value  $\epsilon = 8$  do not greatly affect the performance of the RS model. However, for  $\epsilon = 32$  and  $\epsilon = 64$ , the impact of adversarial perturbation on top-10 HR performance is more observable. Our reconstruction network manages to remove these adversarial perturbations while simultaneously improving visual quality. It has been empirically verified that the preservation of visual quality can be observed for different values  $\epsilon$ , although our reconstruction network is trained with  $\epsilon = 16$  only, which justifies the good transferability of our defense model.

**Visual quality comparison.** Figure 7 compares the images generated by our reconstruction network with different inputs (clean vs. adversarial). The figure shows the inputs and outputs of our reconstruction network under three circumstances and four different values of  $\epsilon$ . We manually



Fig. 7. Reconstructed images from clean and adversarial images. FGSM and PGD are used to generate adversarial images. And we also investigate the effects of the magnitude of the attack perturbation magnitude  $\epsilon$  on generated images. This figure is better viewed when zoomed in.

adjusted this value to strike the best trade-off between attack efficiency and the imperceptibility of adversarial signals. In the presence of adversarial signals, there is a conflict between the objectives of preventing degradation of model performance and generating imperceptible perturbations. For example, if one wants the model performance to degrade gracefully, the introduced attacked signals cannot be too large to become easily observable based on the visual appearance. Meanwhile, if the attacked signals remain imperceptible, they may be too weak to have an adversarial impact on

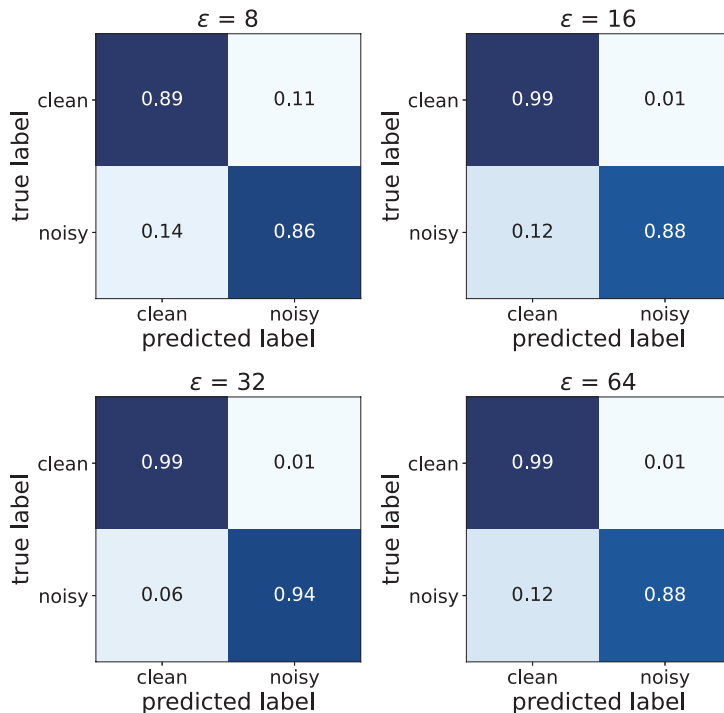


Fig. 8. The detection confusion matrix with different values of the attack perturbation level  $\epsilon$  in the Amazon Men dataset.

the performance of the model. Based on the above analysis, we have handpicked  $\epsilon = 16$  so that our model can reach a good balance between attack efficiency and visual quality degradation. As shown in the second column of Figure 7, the visual quality of the reconstructed images from clean images is almost identical to that from adversarial input from FGSM and PGD attacks.

#### 4.5 Detection Performance and Analysis

As shown in Figure 2, we can train our detection network based on the output images of the reconstruction network. We feed both the original input images and the corresponding reconstructed images into the detection network to obtain two feature vectors and compute the distance of these two feature vectors. By analyzing the distribution of the computed distance, we can learn the representations by maximizing feature consistency under differently augmented views. This way of separating adversarial examples from clean images is conceptually similar to the Adversarial-To-Standard (A2S) model in adversarial contrastive learning [26]. More specifically, the dimension of the output of the feature vector of the detection network is 1024 in our current implementation. When evaluating our detection network, we have used 1000 items in the cold list as negative examples that have never been sampled to train the reconstruction network. Then we sample 1000 positive examples together with 1000 negative examples to obtain the testing set for evaluation. We control the percentage of adversarial examples in the original inputs to be around 50%. Table 3 shows the high accuracy of the detection results in the two datasets, even with different values of  $\epsilon$ . In most cases, our detection network has high accuracy in distinguishing adversarial samples from clean ones.

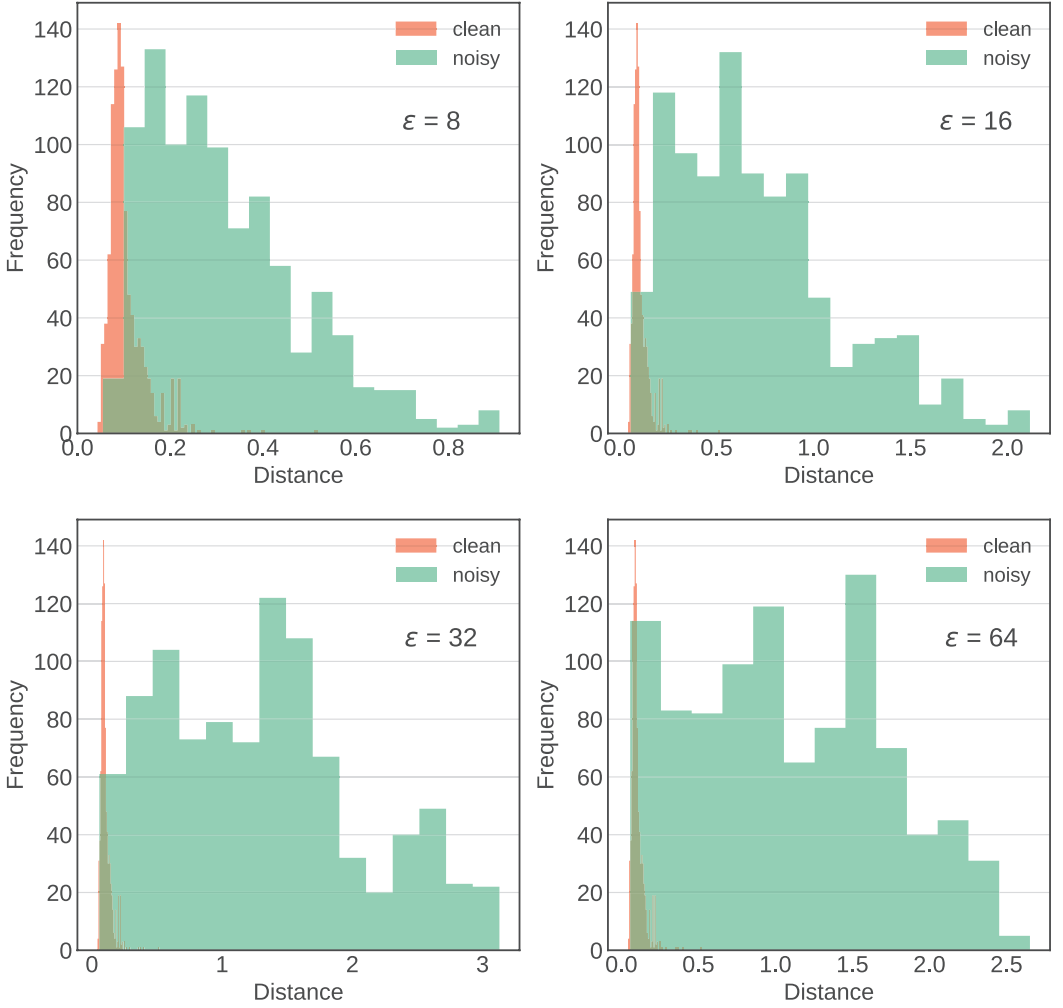


Fig. 9. Distribution of two classes (clean vs. noisy) of distances learned by detection network with different  $\epsilon$  values on Amazon Men dataset.

Figure 8 shows the confusion matrix of the detection results. We can see that our detection network can correctly divide the testing set into two parts: the adversarial set and the clean set in an ideal way. As mentioned in the last section, the defense performance of our reconstruction network decreases with increasing value  $\epsilon$ . However, this figure tells us that our detection network can still have a high detection accuracy in the case of large  $\epsilon$ . Although the performance of our detection network degrades slightly with decreasing  $\epsilon$  values, our reconstruction network still performs well and maintains the performance of the original recommender systems. Through the experimental results, we find that our proposed detection network and reconstruction network complement each other in defense against adversarial attack. With a high  $\epsilon$ , the detection network predominates, while with a small  $\epsilon$ , the reconstruction network is more important.

Figure 9 compares the distributions of the distance profiles learned by the detection network on the Amazon Men data set. It can be seen that after training, the detection network has successfully



Table 3. Detection accuracy under different  $\epsilon$  values.

	$\epsilon$	8	16	32	64
Amazon Fashion	0.9523	0.9665	0.9636	0.9608	
Amazon Men	0.8826	0.9461	0.9660	0.9406	

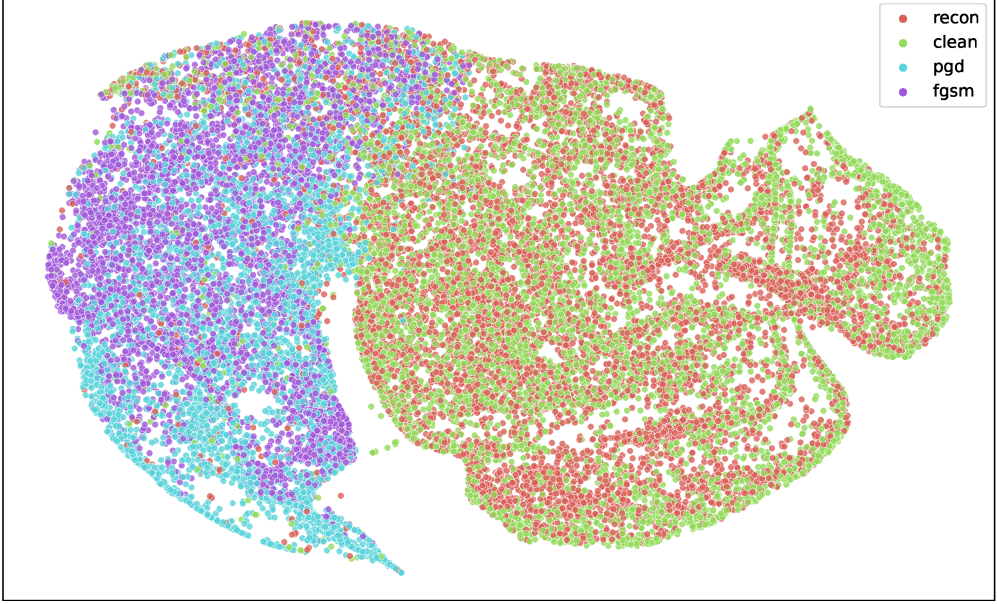


Fig. 10. t-SNE visualization of four classes of images: clean images without perturbations (green), reconstructed images by our reconstruction network (red), and adversarial images with FGSM and PGD perturbations (blue and purple).

separated clean samples from adversarial ones in the 128-dimensional feature space extracted. Additionally, we can calculate appropriate distance thresholds for binary classification (noisy vs. clean). In our implementation, we use 0.2 to classify the test set. We have trained the detection network with the pre-trained reconstruction network with  $\epsilon = 16$  only and the adversarial perturbations generated by the FGSM attack method with  $\epsilon = 8$  only. The reason we use a different value  $\epsilon$  in training the detection network is that we aim to find a more refined and accurate decision boundary because a lower value  $\epsilon$  makes it more difficult for the detection network to separate the two classes apart. Training is completed on the Amazon Men dataset and testing on the Amazon Fashion dataset. We still get satisfactory results in detection accuracy, as shown in Table 3, demonstrating that our detection network has good generalizability.

#### 4.6 Interaction between Reconstruction and Detection Networks

The analysis in the previous subsection suggests that our detection network can have a high detection accuracy in the case of large  $\epsilon$  values. Although the performance of our detection network decreases slightly with decreasing values of  $\epsilon$ , our reconstruction network still performs well and maintains the performance of the original recommender systems. Through empirical studies, we find that our proposed detection network and reconstruction network complement each other in

defense against adversarial attack. With a high  $\epsilon$ , the detection network predominates because the performance of the reconstruction network degrades; while with a small  $\epsilon$ , the role played by the reconstruction network becomes more important because it is more challenging to detect the presence of subtle adversarial perturbations.

To better illustrate the interaction between reconstruction and detection networks, we report the visualization result of t-SNE [52] in Figure 10. It can be clearly seen that the clusters of clean and reconstructed images (marked with green and red) are separated from those of adversarial images with FGSM and PGD perturbations (marked with blue and purple). This clear separation echoes the observation we made in Fig. 9. When combined with the performance of the recommendations reported in Table 2, we conclude that the detection network and the reconstruction network mutually help each other by the joint training proposed with contrast loss.

## 5 CONCLUSION AND FUTURE WORK

In this work, we present a defense framework against adversarial attacks in the recommender system. Our framework can be trained and used as a pre-processing step without affecting the currently running recommender system. The proposed framework consists mainly of two components: the detection network for separating adversarial inputs from clean inputs and the reconstruction network for generating clean images. Two networks play different roles in the defense against different magnitudes of adversarial attack signals. Through extensive experiments, we conclude that these two networks complement each other and collectively succeed in defending the recommender systems from adversarial attacks.

*Limitations and Future Research.* There are a few limitations in this work and interesting future research directions. First, the proposed approach is designed to defend against untargeted attacks under the assumption of a white box attack. An interesting future work is to investigate defenses against targeted attacks under black-box attack assumptions. Second, in this study, we have only experimentally tested two most well-known attack models, namely FGSM and PGD on two popular datasets (Amazon Man and Fashion). It will be interesting to test the performance of our system on other attack models, such as CW attacks on other datasets. Third, the rich interaction between reconstruction and the detection network can be further exploited by joint optimization of two modules. On the one hand, reconstruction can benefit detection because a defense system can strategically manipulate adversarial perturbation so that the objectives of cooperating with both the recommender system and the detection systems can be met. On the other hand, adversarial contrastive learning can also facilitate the task of image reconstruction, because they share the common interest of distinguishing adversarial attacks (aiming at recommender systems) from innocent perturbations such as JPEG compression. The topic of how to achieve self-supervised training of the joint reconstruction and detection network is left for our future study.

## REFERENCES

- [1] Gediminas Adomavicius and Alexander Tuzhilin. 2005. Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE transactions on knowledge and data engineering* 17, 6 (2005), 734–749.
- [2] Vito Walter Anelli, Alejandro Bellogín, Antonio Ferrara, Daniele Malitesta, Felice Antonio Merra, Claudio Pomo, Francesco Maria Donini, and Tommaso Di Noia. 2021. V-Elliot: Design, Evaluate and Tune Visual Recommender Systems. In *Fifteenth ACM Conference on Recommender Systems*. 768–771.
- [3] Vito Walter Anelli, Yashar Deldjoo, Tommaso Di Noia, Daniele Malitesta, and Felice Antonio Merra. 2021. A Study of Defensive Methods to Protect Visual Recommendation Against Adversarial Manipulation of Images. In *The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21)*. 10.
- [4] Robin Burke, Bamshad Mobasher, Chad Williams, and Runa Bhauumik. 2006. Classification features for attack detection in collaborative recommender systems. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*. 542–547.

- [5] Nicholas Carlini and David Wagner. 2017. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE, 39–57.
- [6] Liang Chen, Yangjun Xu, Fenfang Xie, Min Huang, and Zibin Zheng. 2020. Data poisoning attacks on neighborhood-based recommender systems. *Transactions on Emerging Telecommunications Technologies* (2020).
- [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*. PMLR, 1597–1607.
- [8] Rami Cohen, Oren Sar Shalom, Dietmar Jannach, and Amihod 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 (WSDM '21)*. 94–102.
- [9] Yashar Deldjoo, Tommaso Di Noia, Daniele Malatesta, and Felice Antonio Merra. 2021. A Study on the Relative Importance of Convolutional Neural Networks in Visually-Aware Recommender Systems. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*. 3961–3967.
- [10] Yashar Deldjoo, Tommaso Di Noia, and Felice Antonio Merra. 2021. A survey on adversarial recommender systems: from attack/defense strategies to generative adversarial networks. *ACM Computing Surveys (CSUR)* 54, 2 (2021), 1–38.
- [11] Tommaso Di Noia, Daniele Malatesta, 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.
- [12] David Elsweiler, Christoph Trattner, and Morgan Harvey. 2017. Exploiting food choice biases for healthier recipe recommendation. In *Proceedings of the 40th international acm sigir conference on research and development in information retrieval*. 575–584.
- [13] Minghong Fang, Neil Zhenqiang Gong, and Jia Liu. 2020. Influence function based data poisoning attacks to top-n recommender systems. In *Proceedings of The Web Conference 2020*. 3019–3025.
- [14] Minghong Fang, Guolei Yang, Neil Zhenqiang Gong, and Jia Liu. 2018. Poisoning attacks to graph-based recommender systems. In *Proceedings of the 34th Annual Computer Security Applications Conference*. 381–392.
- [15] Francois Fouss, Alain Pirotte, Jean-Michel Renders, and Marco Saerens. 2007. Random-walk computation of similarities between nodes of a graph with application to collaborative recommendation. *IEEE Transactions on knowledge and data engineering* 19, 3 (2007), 355–369.
- [16] Micah Goldblum, Dimitris Tsipras, Chulin Xie, Xinyun Chen, Avi Schwarzschild, Dawn Song, Aleksander Mądry, Bo Li, and Tom Goldstein. 2022. Dataset security for machine learning: Data poisoning, backdoor attacks, and defenses. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45, 2 (2022), 1563–1580.
- [17] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572* (2014).
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [19] Ruining He and Julian McAuley. 2016. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *proceedings of the 25th international conference on world wide web*. 507–517.
- [20] 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.
- [21] Xiangnan He, Tao Chen, Min-Yen Kan, and Xiao Chen. 2015. Trirank: Review-aware explainable recommendation by modeling aspects. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*. 1661–1670.
- [22] Xiangnan He, Zhankui He, Xiaoyu Du, and Tat-Seng Chua. 2018. Adversarial personalized ranking for recommendation. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*. 355–364.
- [23] 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.
- [24] Yang Hu, Xi Yi, and Larry S Davis. 2015. Collaborative fashion recommendation: A functional tensor factorization approach. In *Proceedings of the 23rd ACM international conference on Multimedia*. 129–138.
- [25] Hai Huang, Jiaming Mu, Neil Zhenqiang Gong, Qi Li, Bin Liu, and Mingwei Xu. 2021. Data Poisoning Attacks to Deep Learning Based Recommender Systems. In *Network and Distributed System Security Symposium (NDSS) 2021*.
- [26] Ziyu Jiang, Tianlong Chen, Ting Chen, and Zhangyang Wang. 2020. Robust Pre-Training by Adversarial Contrastive Learning. In *NeurIPS*.
- [27] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. 2016. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*. Springer, 694–711.
- [28] 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.
- [29] Martin Koestinger, Martin Hirzer, Paul Wohlhart, Peter M Roth, and Horst Bischof. 2012. Large scale metric learning from equivalence constraints. In *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 2288–2295.

- [30] Yehuda Koren, Robert Bell, and Chris Volinsky. 2009. Matrix factorization techniques for recommender systems. *Computer* 42, 8 (2009), 30–37.
- [31] Alexey Kurakin, Ian J Goodfellow, and Samy Bengio. 2018. Adversarial examples in the physical world. In *Artificial intelligence safety and security*. Chapman and Hall/CRC, 99–112.
- [32] Shyong K Lam and John Riedl. 2004. Shilling recommender systems for fun and profit. In *Proceedings of the 13th international conference on World Wide Web*. 393–402.
- [33] Jong-Seok Lee and Dan Zhu. 2012. Shilling attack detection—a new approach for a trustworthy recommender system. *INFORMS Journal on Computing* 24, 1 (2012), 117–131.
- [34] Bo Li, Yining Wang, Aarti Singh, and Yevgeniy Vorobeychik. 2016. Data poisoning attacks on factorization-based collaborative filtering. *Advances in neural information processing systems* 29 (2016), 1885–1893.
- [35] Yang Liu, Xianzhao Xia, Liang Chen, Xiangnan He, Carl Yang, and Zibin Zheng. 2020. Certifiable robustness to discrete adversarial perturbations for factorization machines. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 419–428.
- [36] 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 (WWW '21)*.
- [37] 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.
- [38] Bhaskar Mehta, Thomas Hofmann, and Wolfgang Nejdl. 2007. Robust collaborative filtering. In *Proceedings of the 2007 ACM conference on Recommender systems*. 49–56.
- [39] Bhaskar Mehta and Wolfgang Nejdl. 2009. Unsupervised strategies for shilling detection and robust collaborative filtering. *User Modeling and User-Adapted Interaction* 19, 1 (2009), 65–97.
- [40] Andriy Mnih and Russ R Salakhutdinov. 2007. Probabilistic matrix factorization. *Advances in neural information processing systems* 20 (2007), 1257–1264.
- [41] Bamshad Mobasher, Robin Burke, Runa Bhaumik, and Chad Williams. 2007. Toward trustworthy recommender systems: An analysis of attack models and algorithm robustness. *ACM Transactions on Internet Technology (TOIT)* 7, 4 (2007), 23–es.
- [42] Toan Nguyen Thanh, Nguyen Duc Khang Quach, Thanh Tam Nguyen, Thanh Trung Huynh, Viet Hung Vu, Phi Le Nguyen, Jun Jo, and Quoc Viet Hung Nguyen. 2023. Poisoning GNN-based recommender systems with generative surrogate-based attacks. *ACM Transactions on Information Systems* 41, 3 (2023), 1–24.
- [43] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748* (2018).
- [44] Aditya Pal, Chantat Eksombatchai, Yitong Zhou, Bo Zhao, Charles Rosenberg, and Jure Leskovec. 2020. Pinnorsage: Multi-modal user embedding framework for recommendations at pinterest. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2311–2320.
- [45] Sayak Paul and Pin-Yu Chen. 2021. Vision transformers are robust learners. *arXiv preprint arXiv:2105.07581* (2021).
- [46] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2009. BPR: Bayesian Personalized Ranking from Implicit Feedback. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence* (Montreal, Quebec, Canada) (UAI '09). AUAI Press, Arlington, Virginia, USA, 452–461.
- [47] Badrul Sarwar, George Karypis, Joseph Konstan, and John Riedl. 2001. Item-based collaborative filtering recommendation algorithms. In *Proceedings of the 10th international conference on World Wide Web*. 285–295.
- [48] Mete Sertkan, Julia Neidhardt, and Hannes Werthner. 2020. PicTouRe-A Picture-Based Tourism Recommender. In *Fourteenth ACM Conference on Recommender Systems*. 597–599.
- [49] Junshuai Song, Zhao Li, Zehong Hu, Yucheng Wu, Zhenpeng Li, Jian Li, and Jun Gao. 2020. Poisonrec: an adaptive data poisoning framework for attacking black-box recommender systems. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*. IEEE, 157–168.
- [50] Aravind Srinivas, Tsung-Yi Lin, Niki Parmar, Jonathon Shlens, Pieter Abbeel, and Ashish Vaswani. 2021. Bottleneck transformers for visual recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 16519–16529.
- [51] Jinhui Tang, Xiaoyu Du, Xiangnan He, Fajie Yuan, Qi Tian, and Tat-Seng Chua. 2020. Adversarial Training Towards Robust Multimedia Recommender System. *IEEE Transactions on Knowledge and Data Engineering* 32, 5 (2020), 855–867.
- [52] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, 11 (2008).
- [53] Chenwang Wu, Defu Lian, Yong Ge, Zhihao Zhu, and Enhong Chen. 2021. Triple Adversarial Learning for Influence based Poisoning Attack in Recommender Systems. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 1830–1840.

- [54] Chenwang Wu, Defu Lian, Yong Ge, Zhihao Zhu, Enhong Chen, and Senchao Yuan. 2021. Fight Fire with Fire: Towards Robust Recommender Systems via Adversarial Poisoning Training. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1074–1083.
- [55] Shiwen Wu, Fei Sun, Wentao Zhang, Xu Xie, and Bin Cui. 2022. Graph neural networks in recommender systems: a survey. *Comput. Surveys* 55, 5 (2022), 1–37.
- [56] Zhiang Wu, Junjie Wu, Jie Cao, and Dacheng Tao. 2012. HySAD: A semi-supervised hybrid shilling attack detector for trustworthy product recommendation. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. 985–993.
- [57] Zhirong Wu, Yuanjun Xiong, Stella Yu, and Dahua Lin. 2018. Unsupervised Feature Learning via Non-Parametric Instance-level Discrimination. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3733–3742.
- [58] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan Yuille. 2018. Mitigating Adversarial Effects Through Randomization. In *International Conference on Learning Representations*.
- [59] Cihang Xie, Yuxin Wu, Laurens van der Maaten, Alan L Yuille, and Kaiming He. 2019. Feature denoising for improving adversarial robustness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 501–509.
- [60] Xingyu Xing, Wei Meng, Dan Doozan, Alex C Snoeren, Nick Feamster, and Wenke Lee. 2013. Take this personally: Pollution attacks on personalized services. In *22nd USENIX Security Symposium (USENIX Security 13)*. 671–686.
- [61] Senrong Xu, Liangyue Li, Zenan Li, Yuan Yao, Feng Xu, Zulong Chen, Quan Lu, and Hanghang Tong. 2023. On the Vulnerability of Graph Learning-based Collaborative Filtering. *ACM Transactions on Information Systems* 41, 4 (2023), 1–28.
- [62] Zhihai Yang, Zhongmin Cai, and Xiaohong Guan. 2016. Estimating user behavior toward detecting anomalous ratings in rating systems. *Knowledge-Based Systems* 111 (2016), 144–158.
- [63] Feng Yuan, Lina Yao, and Boualem Benatallah. 2019. Adversarial collaborative neural network for robust recommendation. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1065–1068.
- [64] Xiaoyong Yuan, Pan He, Qile Zhu, and Xiaolin Li. 2019. Adversarial Examples: Attacks and Defenses for Deep Learning. *IEEE Transactions on Neural Networks and Learning Systems* 30, 9 (2019), 2805–2824.
- [65] Andrew Zhai, Dmitry Kislyuk, Yushi Jing, Michael Feng, Eric Tzeng, Jeff Donahue, Yue Li Du, and Trevor Darrell. 2017. Visual discovery at pinterest. In *Proceedings of the 26th International Conference on World Wide Web Companion*. 515–524.
- [66] Shudong Zhang, Haichang Gao, and Qingxun Rao. 2021. Defense against adversarial attacks by reconstructing images. *IEEE Transactions on Image Processing* 30 (2021), 6117–6129.
- [67] 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), 5.
- [68] Shijie Zhang, Hongzhi Yin, Tong Chen, Quoc Viet Nguyen Hung, Zi Huang, and Lizhen Cui. 2020. Gcn-based user representation learning for unifying robust recommendation and fraudster detection. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*. 689–698.
- [69] Yihe Zhang, Xu Yuan, Jin Li, Jiadong Lou, Li Chen, and Nian-Feng Tzeng. 2021. Reverse Attack: Black-box Attacks on Collaborative Recommendation. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*. 51–68.