# Salient Feature Extractor for Adversarial Defense on Deep Neural Networks

Jinyin Chen[a,b,*], Ruoxi Chen[b], Haibin Zheng[b], Zhaoyan Ming[c], Wenrong Jiang[d], Chen Cui[e]

[a]*Institute of Cyberspace Security, Zhejiang University of Technology, Hangzhou, China*
[b]*College of Information Engineering, Zhejiang University of Technology, Hangzhou, China*
[c]*Institute of Computing Innovation, Zhejiang Univeristy, Hangzhou, China*
[d]*College of Computer Science, Hangzhou Dianzi University,Hangzhou, China*
[e]*the Big Data and Cyber Security Research Institute, Zhejiang Police College, Hangzhou, China*

## Abstract

Recent years have witnessed unprecedented success achieved by deep learning models in the field of computer vision. However, their vulnerability towards carefully crafted adversarial examples has also attracted the increasing attention of researchers. Motivated by the observation that adversarial examples are due to the non-robust feature learned from the original dataset by models, we propose the concepts of salient feature(SF) and trivial feature(TF). The former represents the class-related feature, while the latter is usually adopted to mislead the model. We extract these two features with coupled generative adversarial network model and put forward a novel detection and defense method named salient feature extractor (SFE) to defend against adversarial attacks. Concretely, detection is realized by separating and comparing the difference between SF and TF of the input. At the same time, correct labels are obtained by re-identifying SF to reach the purpose of defense. Extensive experiments are carried out on MNIST, CIFAR-10, and ImageNet datasets where SFE shows state-of-the-art results in effectiveness and efficiency compared with baselines. Furthermore, we provide an interpretable understanding of the defense and detection process. The code of SFE could be downloaded from `https://github.com/haibinzheng/SFE`.

*Keywords:* Adversarial attack, defense, generative adversarial network, salient feature.

## 1. Introduction

Deep learning enjoys great popularity in both academic and industrial application for its superior performance, ranging from the field of image classification to object detection, natural language processing and bioinformatics analysis. However, deep neural networks (DNNs) are vulnerable to adversarial perturbations, imperceptible to humans, easily lead to misclassification, as proved by Szegedy [43]. In the process of independent decision-making, the vulnerability of deep models to adversarial examples has

---

*Corresponding author
*Email addresses:* chenjinyin@zjut.edu.cn (Jinyin Chen), 2112003149@zjut.edu.cn (Ruoxi Chen), haibinzheng320@gmail.com (Haibin Zheng), mingzhaoyan@gmail.com (Zhaoyan Ming), jiangwenrong@zjjcxy.com (Wenrong Jiang), cuichen@zjjcxy.com (Chen Cui)

posed non-negligible threat to data and information security. This problem, furthermore, impedes the application in mission-critical areas, such as face recognition and auto pilot. Therefore, it is crucial to study defense against adversarial attacks and further improve the robustness of deep models.

In the field of image classification, numerous adversarial attacks have been proposed to discover vulnerabilities of DNNs. Based on the knowledge degree of the target model, adversarial attacks are categorized into white-box attacks and black-box attacks. One of the typical white-box attacks is gradient-based attack, e.g., fast gradient sign method (FGSM) [11], basic iterative method (BIM) [22], momentum-based iterative FGSM (MI-FGSM) [7] and decision boundary-based attack DeepFool [29]. Those attacks require less time to compute perturbations but a thorough knowledge of the targeted model should be given in advance. Black-box attacks are mainly decision-based and score-based, such as the zeroth-order optimization attack (ZOO) [5], point wise attack (PWA) [39] and local search attack (LSA) [31]. Only knowing the output of the targeted model, black-box attacks can still be carried out, usually with a larger size of perturbation.

Meanwhile, the profound implications of DNN's vulnerability have motivated a wide range of investigations into defense for DNNs. Depending on different defense purposes, they can be categorized as re-identification defense and adversarial detector. Re-identification defenses, known as complete defenses, are now developed along three main directions. Training/input modification includes adversarial training [11] and input pre-processing; model modification includes defense distillation [34] and reverse cross-entropy loss [32]; defense of network add-on contains methods based on generative adversarial network(GAN) [10] and autoencoder [13]. Re-identification defense can provide correct class labels of adversarial examples while adversarial detector only determines whether the input is benign or adversarial. The ACT-detector [2], perturbation detection [27], GAT [51] are three of those available defenses.

Another novel angle to study the adversarial example is from MIT [15], they first claimed that the adversarial example is not a bug, concluding that the existence of adversarial examples arise from the non-robust features learned from the original dataset by the model. They gave interpretation of adversarial example from the machine's angle instead of man. From that point of view, the adversarial example takes advantage of non-robust feature to fool DNN while the robust feature is still working for man since the adversarial perturbation is imperceptible. Inspired by this, we propose the concepts of salient feature(SF) and trivial feature(TF), then compare visualization of these features for benign and adversarial examples via Grad-CAM [40], heatmaps on ImageNet VGG19 [41] model shown in Figure1. From red to blue, the weight that the model allocates decreases. By comparing heatmaps, we find that the red area of SF in the benign example is similar to the heatmap of original map, but different from that of the adversarial. Correct labels, consistent with human prediction, can be gained when the model pays attention to SF. On the contrary, misclassification occurs when TF is focused.

Since the SF captures the salient features of the adversarial examples, we proposed a novel defense method, salient feature extractor, SFE for short. It extracts and reconstructs SF and TF to defend against

2

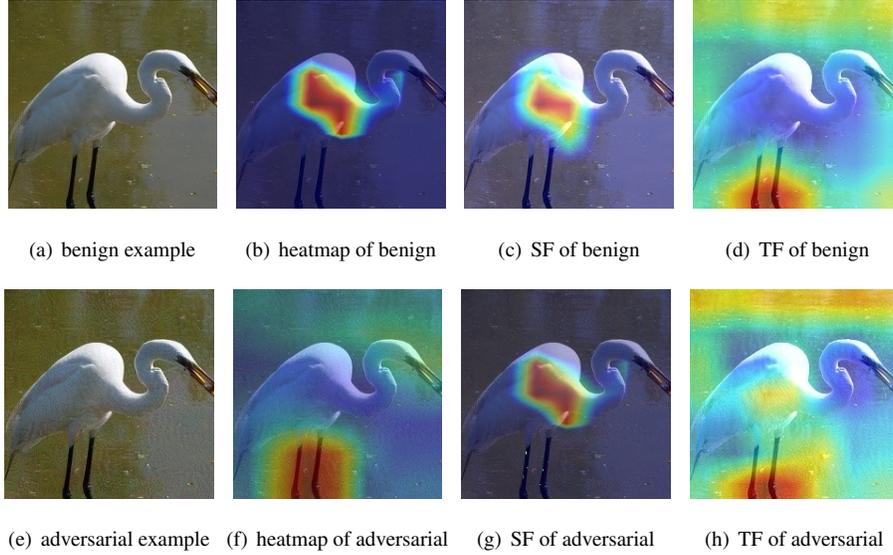|     |     |     |     |
| (a) benign example | (b) heatmap of benign | (c) SF of benign | (d) TF of benign |
| (e) adversarial example | (f) heatmap of adversarial | (g) SF of adversarial | (h) TF of adversarial |

Figure 1: Visualization of SF and TF for VGG19 model on ImageNet via Grad-CAM. Red area in SF of adversarial looks similar to that in benign heatmaps, leading to successful defense. Four images from left to right are the benign example (classified as "American egret") and its heatmap, visualization of SF and TF reconstructed by SFE. The second row represents heatmaps of adversarial examples (misclassified as " cheetah") generated by MI-FGSM and its corresponding SF and TF results after reconstruction. Heatmaps in the second column uses the logits of benign input for calculation while the last two columns use that of the extracted SF and TF.

adversarial attacks. It has been proved that generative adversarial network (GAN) [10] is a well-designed optimizer on basis of game player. We adopt a coupled GAN model to extract SF and TF respectively. Adversarial examples are distinguished by calculating difference between SF and TF. Consequently, correct classification labels are obtained by recognition of reconstructed and strengthened SF in adversarial examples.

Our main contributions are summarized as follows:

1. We propose concepts of salient features(SF) and trivial features(TF). Observation is obtained that the correct classification of examples attributes to the model's concentration on SF. Moreover, consistent with human semantics, SF and TF could be adopted for detection and defense strategy.

2. We design a feature extraction method of SF and TF, namely SFE, which can effectively achieve the detection of benign and adversarial examples by feature separation and calculation of distribution difference. By reconstructing SF and TF of adversarial examples on the basis of coupled GAN framework, both detection and re-identification defense of adversarial examples can be completed.

3. Extensive experiments are implemented on various models and datasets to demonstrate the detection and defense effect of SFE, including within and among-class distance, transferability, parameter sensitivity and time complexity. Besides, interpretable defense via visualization is provided for better understanding.

The rest of the paper is organized as follows. The related works are discussed in Section 2, the SFE

method and critical techniques are introduced in Section 3. Experiments and analysis are detailed in Section 4 respectively. At last, we conclude limitations and future works.

## 2. Related Work

In this section, we review the related literature and briefly summarize the techniques of the attacks and defenses used in the experiment. Moreover, we introduce GAN structure and its related attack and defense methods as well.

### 2.1. Attack Method

We briefly introduce the classic and lately proposed adversarial attacks, including seven white-box attacks and four black-box attacks, which adopted as baselines in experiments. They are implemented to generate adversarial examples of diverse perturbations, which can fairly verify the defense and detection effect of different methods.

### 2.1.1. White-box Attack

White-box attacks mainly use gradient information to determine the direction of perturbation, and increase the value of loss function by adding perturbation until the input is misclassified.

Goodfellow et al. [11] proposed FGSM to find the direction where the gradient of the model changes the most. The algorithm adds perturbation along that direction, which leads to the flip of the predicted label. On the basis of FGSM, Kurakin et al. [22] proposed BIM method, which expanded the operation of increasing the loss function of the classifier to several small steps. As a result, BIM generates smaller perturbation with less transferability. Dong et al. [7] proposed a kind of momentum based iterative algorithm MI-FGSM, which integrates momentum term into the iterative process of attack to generate more transferable adversarial examples. MI-FGSM does boost the effectiveness of adversarial attacks. projected gradient descent (PGD) [26] is an iterative white-box attack. Compared with FGSM of one iteration, it does several iterations with one small step at a time, and each iteration will project the perturbation to the specified range. Consequently it conducts a stronger attack than other previous iterative methods like BIM. Papernot et al. [33] proposed Jacobian-based saliency map attack (JSMA) method to calculate the partial derivative of each output of the last layer of neural network to each input feature and find the part which has the greatest impact on the specific output of the classifier. Based on the critical pixels in the saliency map, the algorithm crafts perturbation on the input, generating adversarial examples. Moosavi et al. [29] generated the minimum norm perturbations by iterative calculation method, namely DeepFool, and gradually pushed the image within the classification boundary to the outside until the wrong predicted label occurred. The "universal" perturbation calculated by Moosavi et al. [28] can fool the network on "any" image with high probability, which is called universal adverse perturbations (UAP). Similar to DeepFool, UAP uses the principle of the classification boundary for all images and

crafts perturbation invisible to humans. Chen et al. [3] focused on the object contours and performed FineFool on the basis of attention mechanism. To reduce image distortion caused by large perturbations, Xiao et al. [48] proposed an adaptive gradient-based adversarial attack method named Adaptive Iteration Fast Gradient Method (AI-FGM). By adaptively seeking gradient, attacks are performed with fewer pixel modifications.

### 2.1.2. Black-box Attack

Black-box attacks do not require the specific parameters of the model; they tend to use the output label or confidence of the model to calculate the perturbation.

LSA [31] is a black-box attack based on greedy local search, which adds perturbation to a single pixel or a small part of randomly selected pixels. Adopting the idea of greedy local search, the algorithm constructs a small set of pixels to perturb, improving the effectiveness of the attack. PWA attack [39] starts with an adversarial and performs a binary search between the adversarial and the original for each dimension of the input individually until the input is misclassified. Contrast reduction attack (CRA) reduces the contrast of the input until it is misclassified. Similarly, additive uniform noise attack (AUNA) adds uniform noise to the input, gradually increasing the standard deviation until the model is fooled. In addition, Chen et al. [17] adopted the genetic algorithm and proposed POBA-GA, which achieves white-box comparable attack performances. Wei et al. [46] put forward the adversarial attributes and use it for the generation of black-box perturbations, only with the knowledge of predicted probabilities from the model.

Adversarial attacks could be found in the real-world scenario as well. In the field of face recognition [37], autonomous vehicles [9] and license plate recognition [36], malicious manipulation of inputs will cause disastrous results.

### 2.2. Defense Method

Numerous defense strategies have been proposed to deal with adversarial attacks. Based on different purposes, defense against adversarial attacks could be classified into two categories: complete defense and adversarial detector, the former also named re-identification defense.

### 2.2.1. Complete Defense

Complete defense mainly develops in the following three directions: using modified input for training or testing and modifying network parameters or structure, e.g, adding more layers or changing loss function, and using add-on network for defense.

In terms of modified training/input, Goodfellow et al. [11] and Huang et al. [14] proposed adversarial training, in which adversarial examples were injected into the training set, enhancing the robustness of neural network towards adversarial examples. On this basis, Mummadi et al. [30] adopted a new training set to finetune the target model by mixing a large number of adversarial and benign example

pairs. Although adversarial training does effectively improve the robustness of the model, it only takes effect on specific attack method with high computation and time cost. Besides, image preprocessing also belongs to defense of input modification. Xie et al. [49] found that random resizing and padding of adversarial examples can weaken the strength of the attack. Prakash et al. [35] redistribute the pixel values in adversarial examples by pixel deflection, and then denoise them based on a wavelet, so as to effectively correct class labels. Image preprocessing methods are easy to operate but hard to deal with large perturbed adversarial examples. Zhang et al. [52] measured the distance between feature distribution of adversarial and benign examples using an optimal transport-based Wasserstein distance. By aligning feature representations, models are no longer easily fooled by a diversity of adversaries.

Network modification mainly provides robustness by changing model structure, loss or activation function. Papernot et al. [34] put forward defense distillation, which uses the knowledge of network to shape its own robustness, and proved that it can resist adversarial examples of small perturbation. Guneet S. Dhillon et al. [6] developed stochastic activation pruning (SAP) method, which prunes the random subset of the activation function and enlarges remaining ones to compensate. As a result, the model gains a certain defense ability against adversarial attacks while maintaining high classification accuracy. Pang et al. [32] minimized the reverse cross entropy in the process of training and proposed RCE, which improves the robustness of the model in the adversarial setting. Most of the network modification defense need to retrain the model, which decreases computational efficiency. Besides, training parameters closely related to defense effect should be carefully chosen to achieve expected results.

The method of using network add-on enables the model to cope with adversarial examples with the help of one or more external models such as autoencoder, GAN or ensemble models. Hlihor et al. [13] adopted the DAE method to train the autoencoder to minimize the distance between adversarial and benign examples, so as to remove perturbations. Ju et al. [20] studied the neural network ensemble method Ens-D for image recognition task. The ensemble of multiple models can still make robust classification when one of them is hacked. Besides, GAN is introduced to defend against adversarial attacks as well, which will be detailed in Section 2.3.

### 2.2.2. Adversarial Detector

Adversarial detectors can distinguish adversarial examples from benign ones, which serves as an alarming bell in the system. Meng and Chen put forward Magnet [8], which uses one or more separated detectors and networks to discriminate adversarial examples by approximating the manifold of benign ones. The detection methods designed by Madry [26], Tramèr [45] and Yin [51] are all based on adversarial training, achieving convincing performance in adversarial detection. Metzen et al. [27] proposed an extended subnetwork detector to distinguish the real data from adversarial ones, so as to make the network itself more robust. These detection methods need to retrain the subnetwork or classifier, which increases computation burden. Tian et al. [44] found that adversarial images are sensitive to transformation operations such as rotation and translation, while benign ones are not. They implemented 45

different image transformation methods to detect adversarial examples, achieving a fair good detection ratio. Moreover, another light-weighted detector named ACT-detector [2] is designed to achieve better detection rate with much less channels, i.e., 5 channels at least.

### 2.3. GAN for Adversarial Attacks and Defenses

GAN, first proposed by Goodfellow et al. [10], has been widely used in image generation, video prediction, object detection and semantic segmentation. Until now, variant structures have derived from GAN, e.g. WGAN, DCGAN, BEGAN [12], dualgan [50] and CoGAN [23]. In general, GAN is a two-player network structure composed of a generator and a discriminator. The generator is designed to imitate, model and learn the distribution of real data as much as possible, reconstruct the random noise or latent variables, and finally generate realistic examples. The discriminator is to differentiate the real data from generated data. Through the continuous competition between the two internal models, the generation and discrimination ability of them are enhanced until a Nash equilibrium is reached.

In aspect of DNN security, GAN is used to generate malicious examples, posing invisible threats to mission-critical field, and also provides a more powerful way out for defense mechanisms.

For GAN based attacks, Xiao et al. [47] trained a conditional GAN to implement attacks, namely AdvGAN, which could generate diverse adversarial examples without accessing the targeted model itself. Chen et al. [4] designed MAG-GAN to generate large-scale adversarial examples, and used it as an effective tool to explore the vulnerability and improve the defense capability of DNNs. Liu et al. [24] introduced adversarial examples in the process of training and proposed Rob-GAN, which leads to the enhancement of convergence speed of GAN training and the quality of generated adversarial examples.

On the other hand, Defense-GAN [38] and APE-GAN [16] introduced GAN for adversarial defense. As input, the mixture of benign examples and adversarial examples is used to train GAN model until it eliminates adversarial perturbation. The defense of network add-on requires to train the extended network with adversarial examples, and generate the output close to the benign examples through reconstruction. In this way, defensive effect is obtained at the cost of more computation times and lower implementation efficiency. However, these GAN-based defense methods directly deal with the whole images, increasing the complexity and difficulty of training. Besides, the defense effect cannot be guaranteed when faced with unknown attacks since they are usually attack-dependent.

## 3. Methodology

We first give the definitions of salient feature and trivial features, and then describe SFE in detail. Adversarial example is detected by extraction, separation and comparison of SF and TF. The reconstructed SF of adversarial examples is re-identified by the model to get the correct classification label. Finally, we analyze the convergence to guarantee the stability and robustness of SFE.

## 3.1. Preliminary

A DNN model consists of input, hidden and output layer. For an input example $x$, the hidden layer of DNN gives the output of the feature, denoted as $h(x)$. The result of the output layer is the classification label, denoted as $f_{DNN}(x)$. Here, we adopt the output of the last fully connected layer as feature.

**Definition of salient feature (SF) and trivial feature (TF)**. For an example pair $< X, X* >$, where $x \in X$ and $x* \in X*$ represent benign and its corresponding adversarial example, respectively. The SF and TF of the benign example $x$ are both $h(x)$, where $f_{DNN}(x) \neq f_{DNN}(x*)$. For the adversarial example $x*$, the SF and TF will be separated as $h(x)$ and $h(x*)$.

Table 1: The definition of SF and TF.

| example pairs | input | feature $x_F$ | salient feature $x_{SF}$ | trivial feature $x_{TF}$ |
|---|---|---|---|---|
| $< X, X^* >$ | $x \in X$ | $h(x)$ | $h(x)$ | $h(x)$ |
| | $x^* \in X^*$ | $h(x^*)$ | $h(x)$ | $h(x^*)$ |

Table 1 clearly denotes SF and TF, where SF and TF of the input $x$ are recorded as $x_{SF}$ and $x_{TF}$, and $x_F$ denotes the output feature of the hidden layer. By learning the difference between SF and TF, the detector obtains the capability of discriminating adversarial examples. The output label of defense is obtained when the reconstructed SF of adversarial examples is input to the model for re-identification.

## 3.2. Framework

The framework of our proposed method is shown in Figure 2, which includes feature extraction and separation of SF and TF, adversarial example detection via training adversarial detector (AdvD) and defense against adversarial examples via re-identification of SF.
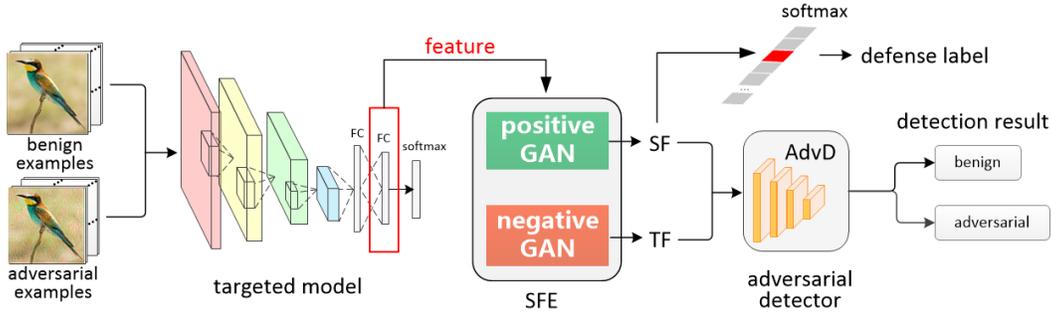


Figure 2: The framework of our proposed method, which includes feature extraction, adversarial example detection and re-identification defense.

Concretely, the targeted model is fed with benign examples and their corresponding adversarial examples at the beginning. The output of the last fully connected layer of the model, considered as high-dimensional feature, is input to SFE. SF and TF are then extracted and separated on the basis of coupled GAN structure. Next, AdvD is trained using the reconstructed SF and TF. Adversarial examples are
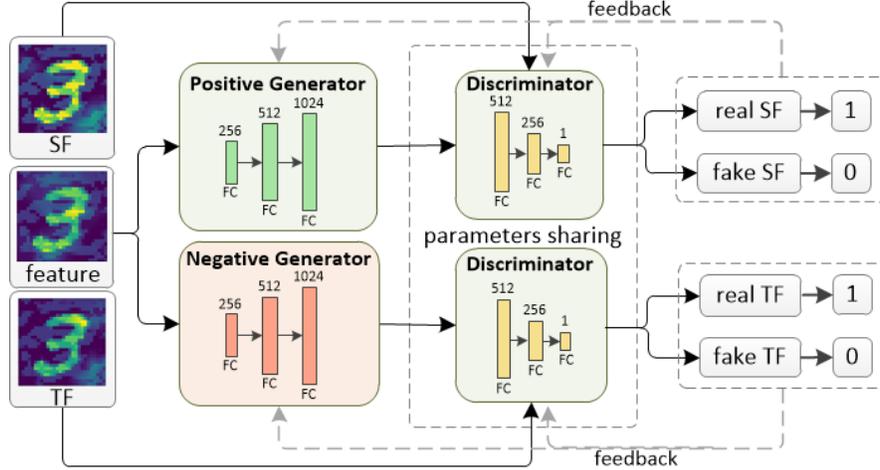
Figure 3: The structure of salient feature extractor. SFE consists of coupled GAN structures called positive GAN and negative GAN, which are painted in green and red. Generators and discriminators are all composed of fully connected layers, FC for short, sizes marked above. By reconstructing SF and TF of adversarial examples on the basis of coupled GAN framework, detection and re-identification defense of adversarial examples can be completed at the same time. SF, TF and feature are visualization heatmaps are obtained via Grad-CAM on adversarial examples generated by MI-FGSM on MNIST-CNN1 model. From yellow to blue, the weight that model allocates decreases.

distinguished by calculating the difference between SF and TF while correct classification labels are obtained by recognition of reconstructed SF, achieving successful defense.

### 3.3. Salient Feature Extractor

SFE consists of two coupled GAN structures called positive GAN and negative GAN, as shown in Fig 3. The former includes a positive generator (PG) and discriminator (D), responsible for learning and generating salient features, while the latter is composed of a negative generator (NG) and D, responsible for trivial features. After inputting the high-dimensional feature of the output of the last fully connected layer of the target model, SFE remaps the features through a coupled GAN structure. Reconstructed SF and TF are provided by PG and NG respectively, which are then input to D for classification. The binary result of D is fed back to the generator and discriminator for optimization of the model parameters. In detail, the output of D is 1 when generated SF and TF are classified as true, vice versa.

The specific structures of generators and discriminators used in the experiment are provided in Table 2. PG and NG are structurally identical with different functions, so they are trained with different data. To reduce complexity, D of positive GAN and negative GAN share parameters during training with the same parameters. This constraint forces high-level features to decode in the same way in both discriminators, thus better capturing the relationship between SF and TF, separating and generating features similar to that distribution of ground truth. In our experiment, G has the same structure in three datasets with stacked fully connected layers, which are sufficient to deal with high-dimensional image features. The activation function of D on CIFAR-10 and ImageNet datasets is slightly different from that

Table 2: The specific structures of G, D and AdvD in our experiment.

| Generator | | Discriminator | | Adversarial example detector | |
|---|---|---|---|---|---|
| Datasets | structure | Datasets | structure | Datasets | structure |
| MNIST CIFAR-10 ImageNet | Dense 256 | MNIST | Dense 512 | MNIST CIFAR-10 ImageNet | Dense 512 activation='ReLU' |
| | LeakyReLU(alpha=0.2) | | LeakyReLU(alpha=0.2) | | BatchNormalization |
| | BatchNormalization(momentum=0.8) | | Dense 256 | | Dropout 0.25 |
| | Dense 512 | | LeakyReLU(alpha=0.2) | | Dense 256 activation='ReLU' |
| | LeakyReLU(alpha=0.2) | | Dense 1 activation='tanh' | | BatchNormalization |
| | BatchNormalization(momentum=0.8) | CIFAR-10 ImageNet | Dense 512 | | Dropout 0.25 |
| | Dense 1024 | | LeakyReLU(alpha=0.2) | | Dense 128 activation='ReLU' |
| | LeakyReLU(alpha=0.2) | | Dense 256 | | BatchNormalization |
| | BatchNormalization(momentum=0.8) | | LeakyReLU(alpha=0.2) | | Dropout 0.25 |
| | Dense activation='tanh' | | Dense 1 activation='sigmoid' | | Dense 64 activation='ReLU' |
| | | | | | BatchNormalization |
| | | | | | Dropout 0.125 |
| | | | | | Dense 1 activation='sigmoid' |

on MNIST. The size of the input layer of G is [$H$,$W$,$C$], the same as that of the image. The size of the output layer of G and that of D is [$H \times W \times C$,1], and the output layer of D is [1,1].

The training process of SFE is shown in Algorithm 1. Benign and adversarial examples will be input into the targeted model, and the output of the last fully connected layer will be taken out as high-dimensional feature for training SFE.

We use mean square error(MSE) as the optimization objective to minimize the distance between the input features and the corresponding generated features, approximating the generated data to real data. The parameters of PG and D are updated alternately during the training process, and loss function is defined as:

$$loss_{PG} = MSE(\text{PG}(\text{x}_\text{F}), \text{x}_\text{SF}) + \text{CE}(\text{D}(\text{PG}(\text{x}_\text{F})), 1) \tag{1}$$

$$loss_{D_{PG}} = CE(\text{D}(\text{PG}(\text{x}_\text{F})), 0) + \text{CE}(\text{D}(\text{x}_\text{SF}), 1) \tag{2}$$

where $MSE(\cdot, \cdot)$ denotes mean square error, $CE(\cdot, \cdot)$denotes cross-entropy of binary classification. $PG(\cdot)$ and $D(\cdot)$ represents the output of the generator and discriminator in positive GAN respectively. $x_F$ denotes feature output by the last fully connected layers, while $x_{SF}$ denotes salient feature.

Similarly, when training negative GAN, the parameters of NG and D are updated alternately, and the loss functions are defined as follows:

$$loss_{NG} = MSE(\text{NG}(\text{x}_\text{F}), \text{x}_\text{TF}) + \text{CE}(\text{D}(\text{NG}(\text{x}_\text{F})), 1) \tag{3}$$

$$loss_{D_{NG}} = CE(\text{D}(\text{NG}(\text{x}_\text{F})), 0) + \text{CE}(\text{D}(\text{x}_\text{TF}), 1) \tag{4}$$

where $NG(\cdot)$ and $D(\cdot)$ denotes output of generator and discriminator in negative GAN separately. $x_{TF}$ denotes trivial feature.

The calculation formula of the average square error is defined as follows:

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 \tag{5}$$

where $y_i$ represents the true class label and $\hat{y}_i$ denotes the predicted value of the model.

Cross-entropy of binary classification is calculated as follows:

$$CE = -(y \cdot log(\hat{y}) + (1 - y) \cdot log(1 - \hat{y})) \tag{6}$$

where $\hat{y}$ is the probability of correct prediction while $y$ represents ground truth label. Concretely, the value is 1 if the example is positive, vice versa.

The pseudo-code of SFE is presented in Algorithm 1.

---

**Algorithm 1**: The training process of SFE

**Input**: Benign examples $X = \{x_1, x_2, \ldots, x_m\}$. Adversarial examples $X* = \{x_1^*, x_2^*, \ldots, x_m^*\}$. High-dimensional feature of the model $x_F = H(X) = \{h(x_1), h(x_2), \ldots, h(x_m)\}$. $x_F^* = H(X^*) = \{h(x_1^*), h(x_2^*), \ldots, h(x_m^*)\}$. The minibatch size $m_b$. Input parameters, the number of steps to pre-train the discriminator $k_D$.

**Output**: salient feature $x_{SF}$ and trivial feature $x_{TF}$.

---

1. **for** $k_D$ steps **do**

2.     Sample minibatch of $m$ benign feature from $x_F$.

3.     Sample minibatch of $m$ adversarial feature from $x_F^*$.

4.     Obtaining generated data $\{x_{SF}^{\sim 1}, x_{SF}^{\sim 2}, \ldots, x_{SF}^{\sim m}\}$, $x_{SF}^{\sim i} = \text{PG}(\text{x}_F^*)$.

5.     Obtaining generated data $\{x_{TF}^{\sim 1}, x_{TF}^{\sim 2}, \ldots, x_{TF}^{\sim m}\}$, $x_{TF}^{\sim i} = \text{NG}(\text{x}_F^*)$.

6.     Update the discriminator parameter $\theta_D$ to minimize $\{loss_{D(PG)} + loss_{D(NG)}\}$.

7. **end for**

8. **for** number of training iterations **do**

9.     Sample new minibatch of $m$ adversarial feature from $x_F^*$.

10.    Update PG parameter $\theta_{\text{PG}}$ to minimize $\{loss_{\text{PG}}\}$.

11.    Update NG parameter $\theta_{\text{PG}}$ to minimize $\{loss_{\text{NG}}\}$.

12.    Update the discriminator parameter $\theta_D$ to minimize $\{loss_{D(\text{PG})} + loss_{D(\text{NG})}\}$.

13. **end**

14. The minimization can use any standard optimization learning rule. We used Adam optimizer in our experiments.

---

### 3.4. Adversarial Examples Detection via AdvD

By using coupled GAN structure, we can separate and extract SF and TF in the input image. As defined in Section 3.1, SF in benign examples is similar to TF, while there is a difference between SF and TF in adversarial examples. Based on that definition, we design the adversarial example detector (AdvD) to determine whether the input data is adversarial.

AdvD consists of five fully connected layers, whose specific structure is shown in Table 2. The input layer size of AdvD is $[H \times W \times C, 1]$, the same as the output of the generator in SFE. And the size of its output layer is $[1,1]$.

As Figure 2 shows, after the training of SFE finishes, its output will be used to train AdvD. The output of PG, the SF of benign and adversarial examples, and the output of NG, TF, will be concatenated to generate the training set of AdvD. During training, the input of AdvD is the concatenated feature, and the output is the detection result. "benign" means benign example, marked with 0, and "adversarial" means adversarial example, marked with 1. During detection, The loss function of AdvD is defined as follows:

$$loss_{\text{AdvD}} = CE(\text{AdvD}(\text{Concat}(\text{PG}(\text{h(x)}), \text{NG}(\text{h(x)}))), 0) +$$
$$CE(\text{AdvD}(\text{Concat}(\text{PG}(\text{h(x}^*)), \text{NG}(\text{h(x}^*)))), 1) \tag{7}$$

where $x$ denotes benign examples while $x^*$ denotes adversarial examples. $\text{AdvD}(\cdot)$ represents the output of AdvD. $Concat(\cdot)$ is concatenate function, which does concatenate operation on the last dimension of the matrix. PG and NG represent generators in positive GAN and negative GAN respectively, whose inputs are $h(x)$, the high-dimensional feature of the last fully connected layer of the targeted model.

For the well-trained SFE, The parameters of targeted model, PG, and NG are fixed and the parameters of AdvD are updated by $min\ loss_{AdvD}$. After training, AdvD will give detection results when the mixture of benign and adversarial examples is input to it.

## 3.5. Adversarial Examples Re-identification via SF

As defined in Section 3.1, SF of the adversarial examples are the same as that of its corresponding benign examples and the hidden layer of the model. In the original model, the high-dimensional image features contain important information closely related to the label, which are fed into the next layer for image classification. Positive GAN in SFE reconstructs the high-dimensional feature and strengthens that important features. Therefore, for well-trained SFE, SF, output from PG, still retains critical information for classification, which can be adopted to correct labels of adversarial examples. Similar to the detection process, we only use generators for the well-trained SFE model. Correct classification results will be given when reconstructed SF is input to the targeted model.

## 3.6. Convergence Analysis

The optimization objective function of GAN is a minimax game corresponding to two players, generator G and discriminator D. To prove the convergence of SFE, we first consider the optimal of any given G. Take PG as an example for the following mathematical proof.

**Proposition 1**: For a fixed PG, D is optimized by PG, that is $\exists D(x)^* = \frac{P_{\text{data}}(x)}{P_{\text{data}(x)} + P_{\text{PG}}(x)}$.

**Proof**: Given PG, D is maximized for

$$V(PG, D) = E_{x \sim P_{data}}[log\mathrm{D(x)}] + \mathrm{E}_{x \sim P_{PG}}[\log(1 - \mathrm{D(x)})]$$
$$= \int_x P_{data}(x) log\mathrm{D(x)}\mathrm{dx} + \int_x \mathrm{P_{PG}(x)}\log(1 - \mathrm{D(x)})\mathrm{dx} \qquad (8)$$
$$= \int_x [P_{data}(x) log\mathrm{D(x)} + \mathrm{P_{PG}(x)}\log(1 - \mathrm{D(x)})]\mathrm{dx}$$

This formula can be simplified as $f(\mathrm{D}) = \mathrm{a}log(\mathrm{D}) + \mathrm{b}log(1 - \mathrm{D})$. If and only if $P_{\mathrm{PG}}(x) = P_{SF}(x)$, D reaches the optimal value $\mathrm{D(x)}^* = \frac{\mathrm{P_{data}(x)}}{\mathrm{P_{data}(x)} + \mathrm{P_{PG}(x)}}$. The same procedure may be easily adapted to obtain the optimum result of D in NG.

**Theorem 1**: if and only if $P_G = P_{data}$, the global minimum of training criterion $C(\mathrm{G})$ = max $V(\mathrm{G}, \mathrm{D})$ can be reached.

**Proof**: According to GAN theory, if and only if $P_G = P_{data}$. At that point, $V(\mathrm{G}, \mathrm{D})$ achieves the value $-log4$. For any G, we can substitute the optimal discriminator $D^*$ obtained in the previous step into $C(\mathrm{G})$ = max $V(\mathrm{G}, \mathrm{D})$ and obtain that $C(\mathrm{G}) = -log4 + 2\mathrm{JS}(\mathrm{P_{data}|P_G})$ where $JS(\cdot)$ denotes Jensen-Shannon divergence. $P_G = P_{data}$ is the possible value of $C(\mathrm{G})$. This concludes the proof.

If PG, NG and D have enough capacity, D can reach the optimal value of given PG and NG.

## 4. Experiments and Analysis

To illustrate general significance of the results, extensive experiments have been carried out to testify the state-of-the-art performance of SFE, including the following parts:

- **RQ1**: Is SFE capable of shielding the model from adversarial manipulation of inputs? And is it competitive enough when compared with baselines about detection and defense results?

- **RQ2**: Does SFE have good transferability over attacks?

- **RQ3**: Can SFE meet the need of efficiency with low time complexity?

- **RQ4**: Can SFE provide a visual understanding of interpretable defense?

*4.1. Setup*

**Platform**: i7-7700K 4.20GHzx8 (CPU), TITAN Xp 12GiB x2 (GPU), 16GBx4 memory (DDR4), Ubuntu 16.04 (OS), Python 3.6, Tensorflow-gpu-1.3, Tflearn-0.3.2. [1].

**Datasets**: We verify the effectiveness of the proposed method on the MNIST [2], CIFAR10[3], and ImageNet[4] datasets. MNIST includes 60000 training examples and 10000 testing examples. Each of it

---

[1]Tflearn can be downloaded at *https://github.com/tflearn/tflearn.*

[2]MNIST can be download at *http://yann.lecun.com/exdb/mnist/*

[3]CIFAR10 can be download at *https://www.cs.toronto.edu/ kriz/cifar.html*

[4]ImageNet can be download at *http://www.image-net.org/*

Table 3: The network structure for MNIST.

| Layer Types | CNN1 (acc=99.79%) | CNN2 (acc=99.76%) |
|---|---|---|
| Conv+ReLU | 5*5*32 | 5*5*16 |
| Max Pooling | - | 2*2 |
| Conv+ReLU | 5*5*64 | 5*5*32 |
| Max Pooling | 2*2 | 2*2 |
| Flatten | 1 | - |
| Dropout | 0.5 | 0.25 |
| Flatten | - | 1 |
| Dense | 128 | 128 |
| Dropout | 0.5 | 0.5 |
| Softmax | 10 | 10 |

is a $28 \times 28$ pixel gray handwritten digital image, marked with ten classes range from 0 to 9. CIFAR-10 dataset consists of 60000 $32 \times 32$ color images of 10 classes such as airplane, bird, ship and frog, with 6000 in each class. ImageNet project is a large visual database for visual object recognition research. It contains more than 14 million images, covering up to 20000 categories. In our experiment, we selected 10 classes with a total of 13500 images for detection and defense test.

**DNNs**: A Variety of different models are adopted in our experiments. For MNIST, we crafted two self-trained ConvNets, whose network structure is shown in Table 3. For CIFAR-10, we trained AlexNet [21] and VGG19 [41], whose classification rate are 99.82% and 99.88% respectively. Besides, Inception v3 (Inc-v3) [42] and VGG19 are used in the experiment of ImageNet dataset. Their recognition accuracy go to 99.48% and 99.50%.

**Attack methods**: We used eleven attack methods, which adopt various algorithms to generate adversarial examples of different perturbation size and distribution, so as to prove the effectiveness of defense and detection of SFE. The white-box attacks include FGSM [11], DeepFool (L2) [29], PGD [26], MI-FGSM [7], BIM [22], JSMA [33] and UAP [28] while black-box attacks contain CRA [19], AUNA [18], PWA [39] and LSA [31]. Attack success rate and average perturbation sizes are shown in Table 4 and Table 5, where $\rho_{adv}$ denotes the average perturbation size of each pixel after normalization to $[0, 1]$.

Table 4: The attack success rate of adversarial examples. The pixel values of each image are normalized to [0,1]. The data in the table represents the attack success rate.

| Datasets | Models | Attack Methods | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | DeepFool | FGSM | PGD | UAP | BIM | MI-FGSM | JSMA | PWA | CRA | AUNA | LSA |
| MNIST | CNN1 | 86.71% | 84.16% | 91.44% | 90.07% | 99.99% | 99.99% | 99.93% | 99.82% | 91.37% | 90.72% | 99.26% |
| | CNN2 | 90.96% | 86.79% | 92.41% | 89.76% | 100.00% | 99.98% | 99.98% | 99.87% | 89.87% | 96.23% | 95.18% |
| CIFAR-10 | AlexNet | 88.18% | 89.03% | 86.71% | 89.83% | 96.22% | 89.57% | 99.85% | 87.45% | 92.95% | 99.99% | 95.37% |
| | VGG19 | 96.66% | 88.29% | 99.95% | 88.82% | 99.87% | 99.76% | 100.00% | 97.60% | 91.67% | 100.00% | 90.03% |
| ImageNet | Inc-v3 | 89.80% | 89.70% | 95.03% | 99.15% | 98.58% | 96.70% | 99.25% | 100.00% | 94.50% | 99.60% | 94.50% |
| | VGG19 | 98.56% | 89.66% | 100.00% | 100.00% | 99.75% | 99.50% | 99.75% | 96.15% | 100.00% | 99.50% | 90.00% |

**Defense baselines**: Nine defense methods and five detection baselines were applied in our experi-

Table 5: The average perturbation size and attack parameters of adversarial examples.

| Attack | $\rho_{adv}$ | Parameters |
|--------|------|------------|
| DeepFool | 0.002 | epsilon=1e-6,maxiter=100 |
| FGSM | 0.008 | epsilon=0.3,stepsize=0.05,iterations=10. |
| PGD | 0.001 | epsilon=0.3,stepsize=0.01,iterations=100 |
| UAP | 0.002 | delta=0.2,maxiter=20 |
| BIM | 0.001 | epsilon=0.3,stepsize=0.05,iterations=10, |
| MI-FGSM | 0.001 | epsilon=0.3,stepsize=0.06,iterations=10,decayfactor=1 |
| JSMA | 0.002 | maxiter=2000,num random targets=0,fast=True,theta=0.1, max perturbations per pixel=7 |
| PWA | 0.006 | |
| CRA | 0.031 | epsilons=1000 |
| AUNA | 0.012 | |
| LSA | 0.061 | perturbation parameter r=1.5,p=10.0,d=5,max perturbations per pixel=5,maxiter=150 |

Table 6: The operation setting of different defenses.

| Defense | Operation |
|---------|-----------|
| resize (RS) | $(H, W) \rightarrow (H/2, W/2) \rightarrow (H, W)$ |
| random resize (RRS) | $(H, W) \rightarrow (H', W') \rightarrow (H, W)$, where $H' \in [H/2, H]$, $W' \in [W/2, W]$ |
| rotate (RT) | rotate $-45^o \rightarrow$ rotate $45^o$, fill in missing values caused by rotation |
| random rotate (RRT) | rotate rotate $-r \rightarrow$ rotate $r$, fill in missing values caused by rotation, where $r \in [0^o, 45^o]$ |

ment. Defense methods includes resize(RS), rotate(RT), random resize(RRS),random rotate(RRT) [49], RCE [32], Ens-D [20], DAE [13], Defense-GAN [38] and APE-GAN [16]. The settings of different defenses are detailed in Table 6. As for detection, 45C-Detector [44],Perturbation detection (Per-D) [27], adversarial training-based detection (AdvT-D) [45], PGD-based detection (PGD-D) [26] and GAT [51] are used to make fair comparison with SFE.

**Metrics**: The metrics used in the experiments are defined as follows:

① Classification accuracy: $acc = \frac{n_{true}}{N}$, where $n_{true}$is the number of clean examples correctly classified by the targeted model and $N$ denotes the total number of benign images.

② Attack success rate: $ASR = \frac{N_{adv}}{N}$, where $N_{adv}$ denotes the number of adversarial examples misclassified by the targeted model after attacks.

③ Perturbation L2-norm: $\rho_{adv} = ||x_{adv} - x||_2$, where $x$ and $x_{adv}$ are benign and its corresponding adversarial example respectively, and $|| \cdot ||_2$ represents L2 norm.

④ Detection rate, $DR = \frac{N_{det}}{N}$,where $N_{det}$ is the number of input examples that are detected by model.

⑤ Defense success rate, $DSR = \frac{N_{def}}{N_{adv}}$,where $N_{def}$ denotes the number of adversarial examples correctly classified by the target model after defense.

⑥ Within-class distance, $FSA_k = 1 - \frac{norm(\sum j=1^{n_k} dist(f_{x_j,k}, f_{ck}))}{n_k}, FSA = \frac{\sum i=1^K d^2(k)}{K}$ [1], where $n_k$ denotes the number of examples belonging to class $k$. $f_{x_j,k}$ denotes the feature vector of the example

$x_j$ belonging to the $k-th$ class in the high-dimensional feature space. $f_{ck}$ represents the center of class $k$ in the feature subspace. $K$ denotes the number of total class and $dist(\cdot)$ represents the normalization function. For each class, we sum $FSA_k$ and average it to gain the final $FSA$. The model with smaller with-in class distance always has higher classification accuracy.

⑦ Among-class distance, $FSD_{k,k+1} = dist(f_{ck}, f_{ck+1}), FSD = \frac{1}{K(K+1)/2} \sum_{i=1}^{K-1} \sum_{j=i+1}^{K} FSD_{i,j}$ [1], where $f_{ck}$ and $f_{ck+1}$ represent the center of class $k$ and $k+1$. The final $FSD$ is calculated by averaging the sum of $FSD_{i,j}$ in each class. Similar to with-in class distance, the model with larger among-class distance performs better in identification and recognition.

*4.2. Comparison of Detection Results*

In this section, we will focus on the aspect of detection in question **RQ1**, and use detection rate (DR) to measure the detection effect of SFE on adversarial examples. Comparison will be made with SFE and detection baselines. Fig. 4 shows the detection results of SFE and its baselines on three different datasets and six models, where the ordinate represents DR and the last column mean denotes the average of it. In the experiment, we mix the randomly selected adversarial and benign examples then input them to the model for detection. The image of training and testing set is 7:3.

According to the experiment results shown in Fig. 4, SFE has high detection accuracy for a variety of datasets and models, and is superior to five baselines in most cases. Reaching almost 100% on MNIST and CIFAR-10, DR of SFE slightly decreased on ImageNet dataset. The reason goes to the fact that image size on the ImageNet is larger, so the difference between the salient and trivial features after reconstruction of SFE are obscure, which somewhat reduces the detection performance. As indicated in Fig. 4, the detection effects of SFE remain stable. Meanwhile, DR of black-box attacks are practically as good as that of white-box attacks, which indicates that SFE can accurately distinguish differences between features for comparison and detection. In this way, a good transferability between various attacks is guaranteed. In contrast to other detection algorithms, their DR fluctuates greatly among different attacks. This can be attributed to the fact that there exist differences in the image transformation and training of the adversarial examples generated by various attacks, leading to unstable and unreliable performance. However, SFE shows its vulnerability to attacks with large perturbation, such as CRA and LSA, and has a decline in DR. This is because large perturbation destroys important features in some cases, so the reconstructed SF and TF can not fully reflect the classification information of the original image.

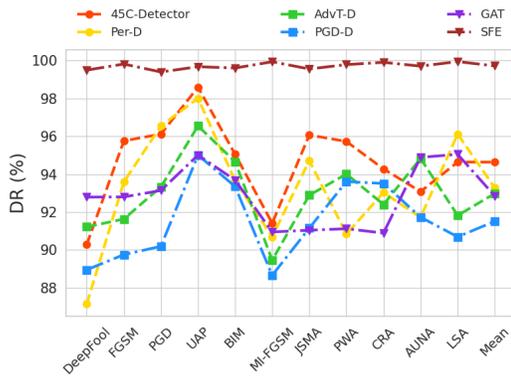*4.3. Comparison of Defense Results*

In the previous section, we have discussed the detection capability of SFE against adversarial examples. And in this section, we conduct the defense experiment to answer the defense aspects of **RQ1** by using DSR to measure the effect of re-identification defense.
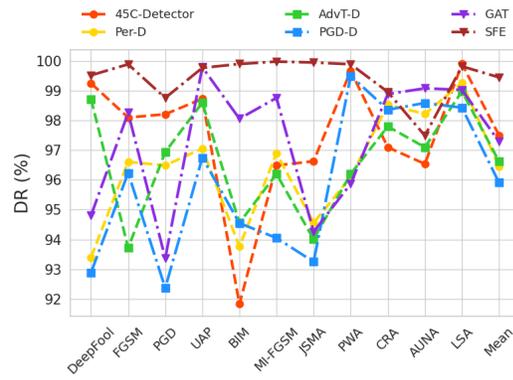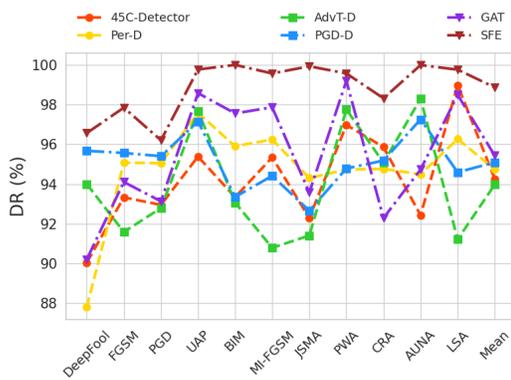
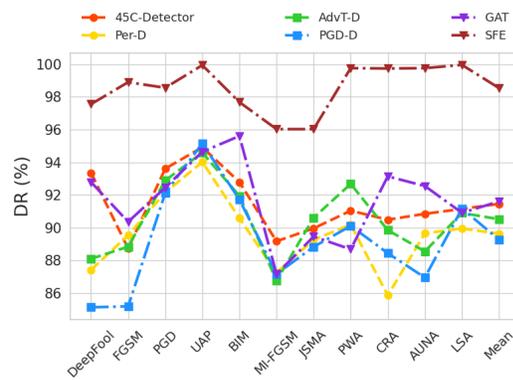(a) MNIST-CNN1

(b) MNIST-CNN2

(c) CIFAR-10 AlexNet

(d) CIFAR-10 VGG19

(e) ImageNet Inc-v3

(f) ImageNet VGG19

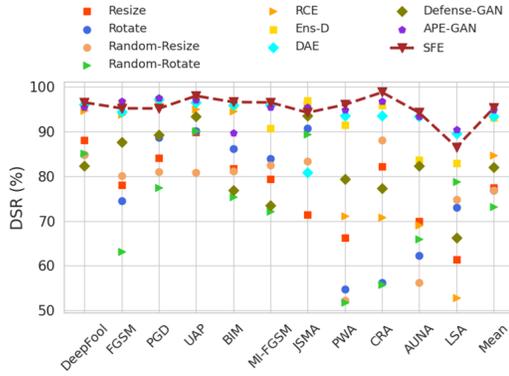Figure 4: Comparison of detection results against various adversarial attacks.

17

Fig. 5 shows the experimental results of SFE and nine other defense methods on different datasets, models, and attacks, where the ordinate represents DSR, and the last column is the average of it. The results of SFE are denoted by broken brown lines while other defense methods are represented by scattered points with different shapes.

Generally speaking, our proposed method significantly enhances the robustness of the model against adversarial attacks. It can be clearly observed that the average DSR on three datasets is about 90%, which reaches the highest among all baselines. This indicates that SFE can successfully reconstruct the salient features of different kinds of adversarial examples, making them similar to benign images. Therefore, when salient features are input into the model for re-identification, the correct classification results can be obtained. With regard to different attacks, SFE has a certain defense ability against black-box attacks, but the effect is inferior to white box, opposite to detection results. One possible reason is that black-box attacks tend to have larger perturbation, which may pose a negative impact in the process of extracting and reconstructing salient features. Similarly, for white-box attacks like FGSM with large perturbation, the effect of SFE on feature extraction and reconstruction is slightly reduced, which leads to the decline of DSR. As for different datasets, the defense capability of SFE is affected by data distribution as well. For large-scale images and complex data sets such as ImageNet, the difficulty of feature extraction increases, which affects the defense effect to a certain extent. However, with the increase of the complexity of data sets, DSR of SFE remains around 85%. In comparison with baseline, we can find that image transformation operations such as resize and rotate do weaken the attack strength, but only when the parameters are selected properly can it get satisfactory defense effect. Although defense results of RCE and Ens-D are close to SFE in DSR, SFE outperforms them in most situations, as well as APE-GAN and defense-GAN.
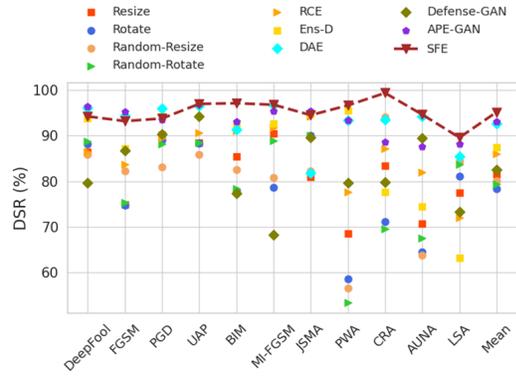
For **RQ1**, we can conclude that SFE has quite competitive performance in both detection and re-identification defense. It is slightly superior when compared with baselines. When a malicious perturbation is input to the model, SFE can play a good alarming role while correct classification labels for most of the adversarial examples.

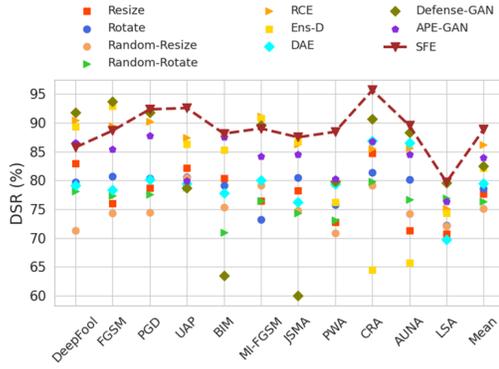*4.4. Analysis of Within and Among-class Distance*

In order to better verify the effectiveness of SFE, we pay attention to the feature manifold during the process, and further calculate the within-class and among-class distance of benign example, its corresponding adversarial example and the defense model output. As shown in the formula in Section 4.1, within-class distance with a smaller value indicates a more concentrated feature distribution of the same class in the image. Similarly, among-class distance with larger value shows that the feature distribution of different classes is more dispersed, which is easy for the model to distinguish.
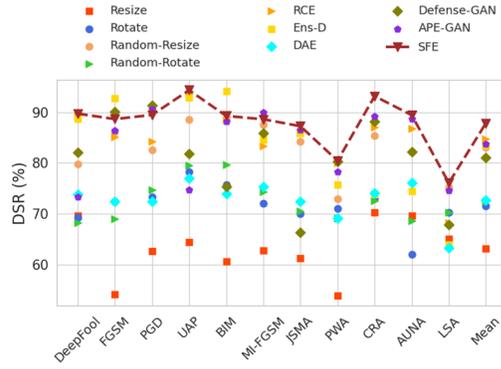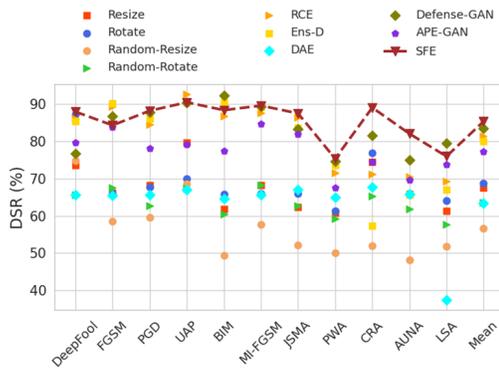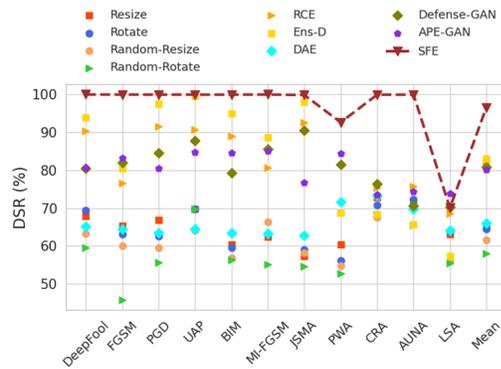
18

(a) MNIST-CNN1

(b) MNIST-CNN2

(c) CIFAR-10 AlexNet

(d) CIFAR-10 VGG19

(e) ImageNet Inc-v3

(f) ImageNet VGG19

Figure 5: Comparison of defense results for different classifiers on different datasets.
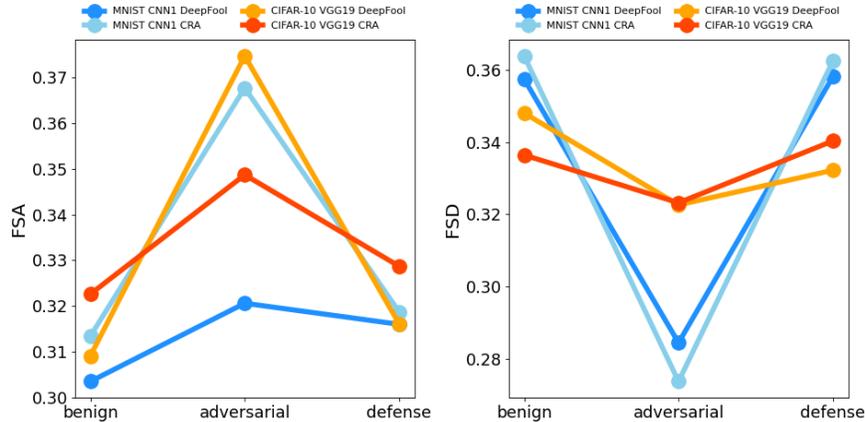
Figure 6: FSA and FSD for benign and adversarial examples before and after defense, where all data in the figure is normalizated.

Fig. 6 shows the changes of within-class and among-class distance of four models in MNIST and CIFAR-10 on white-box attack DeepFool and black-box attack CRA. The ordinate represents the value of two distances, and the abscissa represents the output of the benign and adversarial example before and after defense.

As for within-class distance, eight broken lines show the same trend: all rise first and then decline, with close values in benign and defense columns. The curve of among-class distance showed a "V" shape, which decreases first and then increases. The reason lies in that adversarial attacks increase the distance of within-class, that is to say, scatter the feature distribution of the same class. On the other hand, it also reduces the distance between different classes, making the feature distribution of different classes more concentrated and hard to distinguish. Consequently, misclassification of the model appears. SFE reconstructs the features of adversarial examples to approximate that of benign examples, reduces the distance within class and increases the distance among classes. As the result, it is easier for model to distinguish examples from different classes, and finally guarantees the correct label output by the model after defense.

*4.5. Defense Impact on Benign Examples*

A robust defense should reduce or minimize the misclassification rate caused by various adversarial attacks while maintaining a high classification accuracy. In order to further demonstrate the reliability of SFE, we randomly selected 1000 benign examples from MNIST,CIFAR-10 and ImageNetand input them to the model after defense for re-identification, so as to measure the impact of defense on benign examples. At the same time, we also tested the computation time and obtained results in Table 7.

It can be observed from the table that SFE does not sacrifice the classification accuracy of benign examples while completing efficient defense. It indicates that SFE can accurately find most of the critical pixels that affect the classification results when reconstructing salient features. In this way, no obvious decline in accuracy could be found during the experiment. Besides, SFE consumes less time than that on

Table 7: The classification accuracy and time complexity of SFE on benign examples before and after defense, where "benign" denotes classification and time results of benign examples while "defense" denotes results of benign examples after SFE defense.

| Datasets | Models | acc | | time/s | |
|---|---|---|---|---|---|
| | | benign | defense | benign | defense |
| MNIST | CNN1 | 99.79% | 98.57% | 5.36 | 5.42 |
| | CNN2 | 99.76% | 98.77% | 5.60 | 5.54 |
| CIFAR-10 | AlexNet | 99.92% | 98.87% | 6.04 | 5.95 |
| | VGG19 | 99.88% | 98.53% | 6.15 | 6.02 |
| ImageNet | Inc-v3 | 99.48% | 98.35% | 7.72 | 7.60 |
| | VGG19 | 99.50% | 91.71% | 7.35 | 7.28 |

original setting. This is because SFE extracts features of benign examples for re-identification, decreases the dimension and size of the data to be processed, thus reducing the time cost in the defense process.
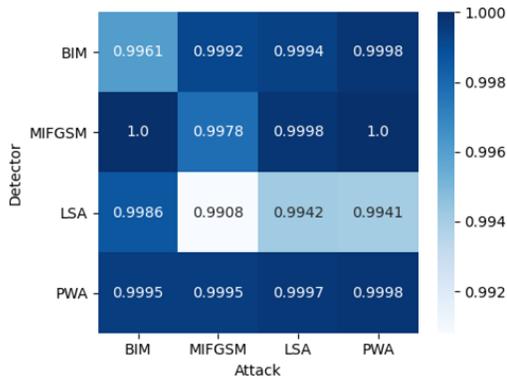
*4.6. Defense Transferability of SFE*

In the real-world scene, defense is carried out without the knowledge of the attack algorithm implemented by the attacker. Therefore, the defense transferability among attacks is particularly important. In this section, we discuss the transferability of SFE. Under this setting, GAN is trained with adversarial examples crafted by a certain attack, but reconstructed SF and TF are used to detect and defend against other attacks. Under this setting, we performed a transferability experiment on SFE to address **RQ2**.
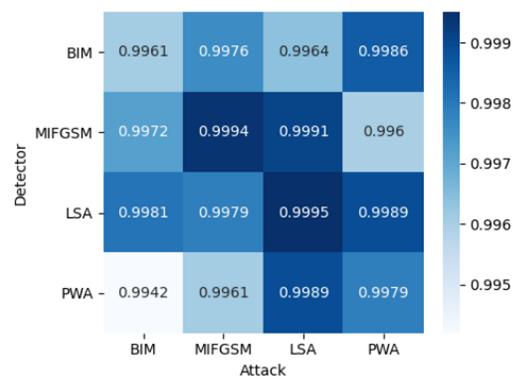
In order to analyze the detection ability of SFE between various attacks, we input adversarial examples into detectors trained by the features of adversarial examples generated by other kinds of attacks for testing. Experiments are carried out in CNN1 of MNIST and AlexNet of CIFAR-10. BIM and MI-FGSM are randomly selected as white-box attacks, with LSA and PWA as black-box attacks. Other experimental settings are the same as that in Section 4.2. The visualization results of detection are shown in Fig. 7, where the color intensity of the squares is proportional to DR. The horizontal line represents the adversarial examples used for testing, and the vertical line represents those used to train the detector. The result on the diagonal is the non-transferable data in Section 4.2.

As the figure suggests, SFE shows detection success rate of more than 98% regardless of whether the salient and trivial features are trained by the detected attack or not. Compared with the adversarial examples of training and transferability testing, they are quite the same in the detection accuracy between them. More specifically, some results of transferability testing even exceed that of the non-transferable experiment. In addition, SFE has good transferability between white-box attacks and black-box attacks, which indicates that the detection performance is fairly independent on adversarial examples used for training.

Meanwhile, we also implement similar experiments to investigate the transferability of defense. PWA and CRA are selected as black-box attacks this time, and other experiment settings are consistent with Section 4.3. Defense results are shown in Fig. 8.
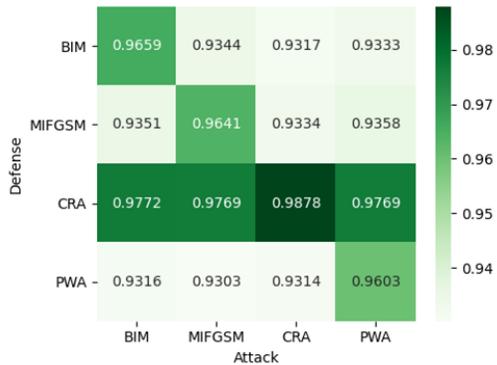
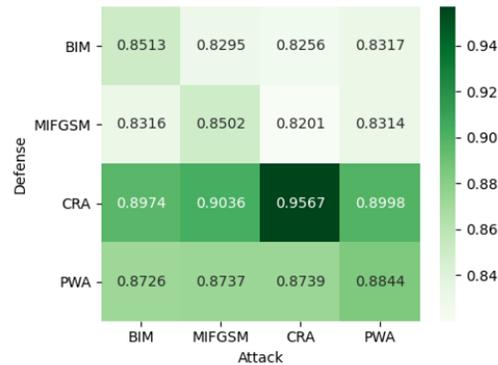(a) DR of MNIST-CNN1        (b) DR of CIFAR-10 AlexNet

Figure 7: The visualization of heatmap of SFE in detection transferable setting on MNIST and CIFAR-10 dataset, where blue bars denote the results of DR.

It can be observed from the figure that the results of defense are roughly similar to those of detection. Although the defense effect of the transferability experiment was slightly lower than that of original experiment, no obvious difference between them could be seen from the figure. The same result can be observed as well between black and white-box attacks. No matter what kind of adversarial examples are used for training, SFE can maintain a stable defense effect among a variety of attack methods.

For **RQ2**, the experimental results demonstrate that SFE shows a certain transferability among a variety of attack methods in detection and defense. Moreover, its defensive capability is unrelated to attack algorithms.



(a) DSR of MNIST-CNN1        (b) DSR of CIFAR-10 AlexNet

Figure 8: The visualization of heatmap of SFE in detection transferable setting on MNIST and CIFAR-10 dataset, where blue bars denote the results of DR.

Table 8: The structure of G model after different operations.

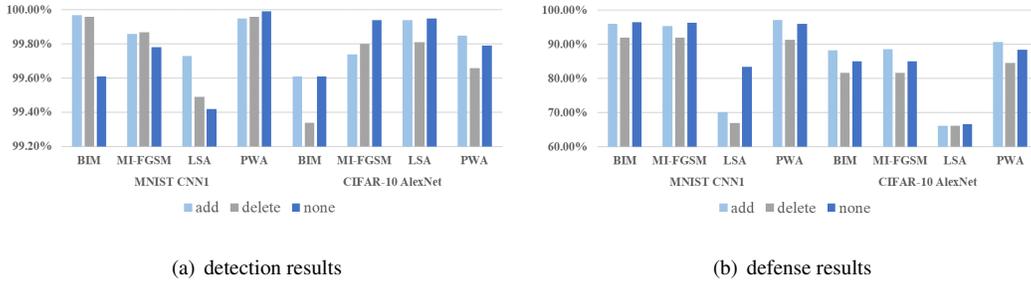| Layers | Operations | |
|---|---|---|
| | add | delete |
| Dense | 256 | 256 |
| LeakyReLU | alpha=0.2 | alpha=0.2 |
| BatchNormalization | momentum=0.8 | momentum=0.8 |
| Dense | 512 | - |
| LeakyReLU | alpha=0.2 | - |
| BatchNormalization | momentum=0.8 | - |
| Dense | 1024 | 1024 |
| LeakyReLU | alpha=0.2 | alpha=0.2 |
| BatchNormalization | momentum=0.8 | momentum=0.8 |
| Dense | 1024 | - |
| LeakyReLU | alpha=0.2 | - |
| BatchNormalization | momentum=0.8 | - |
| Dense | 1024 | - |
| LeakyReLU | alpha=0.2 | - |
| BatchNormalization | momentum=0.8 | - |



(a) detection results

(b) defense results

Figure 9: The results of detection and defense on different operations of G structure. This curve is based on the CNN1 on MNIST dataset and AlexNet on CIFAR-10 dataset.

### 4.7. Parameter Sensitivity Analysis

In this section, to further illustrate the reliability of SFE, we study the influence of hyper-parameter on defense and detection, especially the generator(G) structure of GAN.

We add or delete the layer of G in GAN. The structure of the original G is the same as that described in Section 3. We carried out the experiment in CNN1 of MNIST and AlexNet of CIFAR-10. Adversarial examples we used are crafted by BIM, MI-FGSM, LSA, PWA with the same other experimental parameters. The results of detection and defense are shown in Fig. 9, where add means adding the layer of G, delete means deleting the layer, and none means unchanged. The structure of the G model after operations are shown in Table 8.

According to the figure, we can conclude that adding or deleting the structure of G has little effect on the detection effect on MNIST and CIFAR-10 data sets with all DR stable at around more than 99%. This shows that the difference between salient and trivial feature are still obvious for model to differentiate, which is less susceptible to structure of G. On the other hand, Increasing the structure of G may improve

23

the defense effect to some extent, but deleting the structure of it may slightly reduce the defense effect. This is not hard to understand: By adding the structure of G, the fully connected layer maps the learned distribution features in the high-dimensional space for more times, which improves the learning ability of GAN for salient and trivial features, and vice versa. Meanwhile, we also notice that the simple model like CNN1 is more vulnerable to change of G structure than more complex model like AlexNet in defense, especially in attacks with large perturbation such as LSA.

### 4.8. Comparison of Algorithm Complexity

In the previous section, we have confirmed that SFE has competitive defense capability, good transferability, and stable parameter sensitivity.

In the design process of the algorithm, effectiveness and efficiency should all be taken into consideration with the fast development of big data and cloud computing technology. Therefore, in this section, the efficiency of the algorithm will be verified. We present a comparison of detection and defense times between SFE and baselines. Combining with the experimental results, we will give an answer to question **RQ3**.

In the experiment, We chose 10000 images for training and 1000 for testing. Adversarial examples are crafted by DeepFool. Fig. 10 shows computation time of training and testing during detection on MNIST and ImageNet.

For MNIST, the training time of SFE is almost four times that of Per-D, mainly because GAN needs to learn and imitate the salient and trivial features during training. Compared with the other three detection methods, the training time of SFE is less than 5 seconds, much shorter than others. Besides, the test speed of SFE is the fastest among these detection methods.

Similar pattern could be concluded in the experiments results for ImageNet. But for large data sets, the training complexity of the detector increases, extending the training time. The training time of SFE is still less than 12 seconds, superior to other methods. In terms of testing time, SFE still ranks top, with the shortest time consumption.

We implement the same experiment for defense, the result of which is shown in Fig. 11.

As the figure suggests, resize and rotate belong to input transformation defense, which does not involve training operation, so the training time is 0. Ens-D needs to train two ensemble models, so it is time-consuming. RCE, DAE, Defense-GAN and APE-GAN all have to train a new model to filter out perturbations while SFE only needs to reconstruct the features extracted from the model, so the training cost is much lower than these baselines. In general, the test time of all defense methods is quite close, but SFE only needs to re-identify reconstructed features, which proves to be the most efficient algorithm.

Here, we can answer **RQ3**: SFE shows the low cost in terms of time complexity. Concretely, the training time of SFE is less than most of the detection and defense methods that require training.
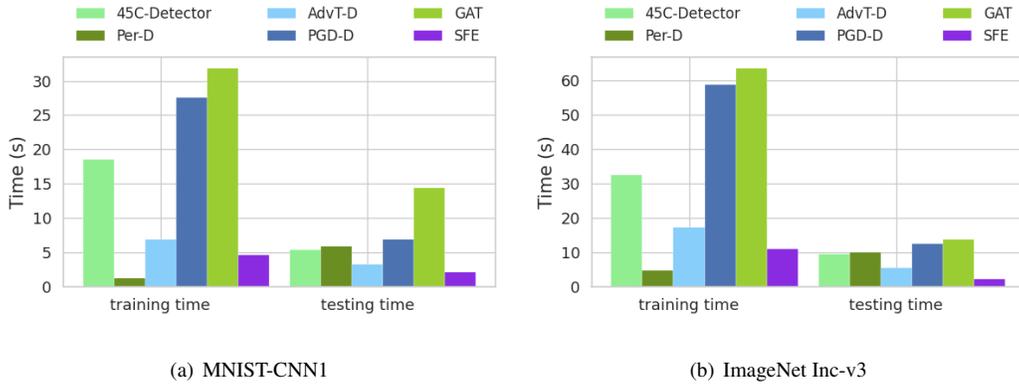
24

(a) MNIST-CNN1          (b) ImageNet Inc-v3

Figure 10: Time comparison among detection algorithms during training and testing for CNN1 classifier on MNIST dataset and Inc-v3 classifier on ImageNet.


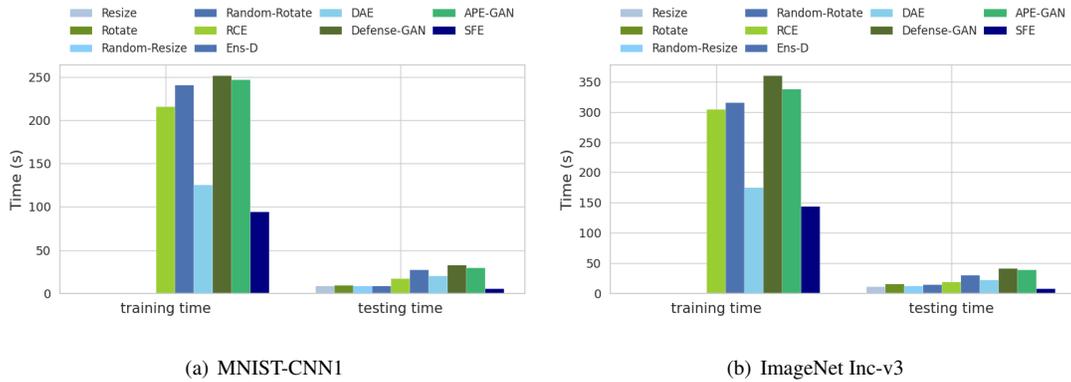
(a) MNIST-CNN1          (b) ImageNet Inc-v3

Figure 11: Time comparison among defense algorithms during training and testing for CNN1 classifier on MNIST dataset and Inc-v3 classifier on ImageNet.

> Moreover, the test time is much lower than baselines as well. By comparing and re-identifying features, SFE achieves low complexity and high efficiency, which helps it applicable to various scenarios.

### 4.9. Detection and Defense Results of Adaptive Attack

In this section, we discuss the detection and defense effectiveness of SFE under adaptive settings, where the attacker knows our defense methods in advance.

To prevent SFE from extracting and separating SF and TF from the input, we increased the perturbation size added to the benign examples to 0.08, almost 100 times that of the previous experiment. The specific adversarial examples on MNIST and CIFAR-10 are shown in Fig. 12, where a human can hardly tell the contents of the image. The results of detection and defense, as shown in Table 9, are obtained with the same experimental parameters as Section 4.2 and Section 4.3.
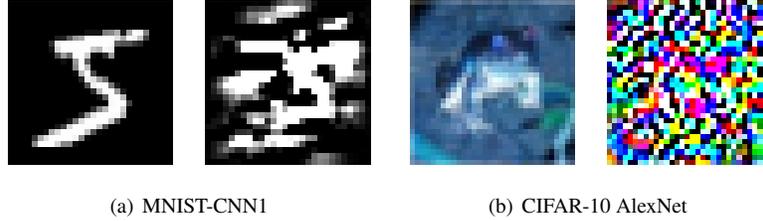
(a) MNIST-CNN1        (b) CIFAR-10 AlexNet

Figure 12: Adaptive adversarial examples of MNIST and CIFAR-10, where the left is benign example while the right is its adversarial example.

From the data in table 9, SFE still has certain detection and defense capabilities for adversarial examples with such exaggerated perturbations: more than half of adversarial examples can be detected and finally re-identified. It is indicated that SFE shows robustness under adaptive attack settings. Generally speaking, for adaptive attacks, detection results of SFE are better than defense. Reasons lie in that detection only calculates the difference between SF and TF, and does not require the label-related information in extracted features. On the contrary, re-identification defense directly use that information of important features for classification, which is more sensitive to the large perturbation added by adaptive attacks in the image.

Table 9: The detection and defense performance of SFE under adaptive attack settings.

| Datasets | MNIST | | CIFAR-10 | | ImageNet | |
|---|---|---|---|---|---|---|
| **Model** | CNN1 | CNN2 | AlexNet | VGG19 | Inc-v3 | VGG19 |
| **DR** | 70.37% | 71.52% | 70.40% | 70.20% | 68.60% | 69.20% |
| **DSR** | 56.10% | 54.96% | 55.92% | 52.04% | 55.06% | 58.64% |

### 4.10. Visualization based Interpretation

In this section, we will analyze and verify the effectiveness of defense and detection from high-dimensional feature space via t-SNE [25] and image pixel features via Grad-CAM [40]. Interpretable understanding of the whole process is detailed and **RQ4** is also answered later.

### 4.10.1. Visualization of t-SNE

t-SNE is adopted to visualize the process of attack, defense and detection from the perspective of cluster distribution in high-dimensional feature space.

① Visualization of Attacks

Fig. 13 shows the t-NSE visualization results of benign examples before and after attack on MNIST-CNN1 model. From left to right, are corresponding t-NSE results of benign examples, and adversarial examples generated by MI-FGSM, JSMA and PWA. Different colour blocks represent different class clusters.

The well-trained model can accurately classify images from different class, so all kinds of clusters in t-SNE visualization of benign examples are separated from each other. Comparing the visualization
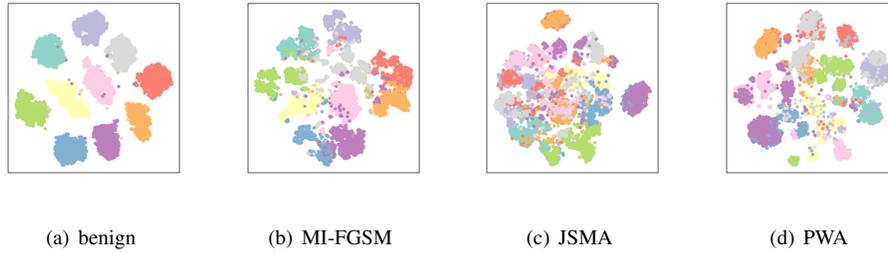
(a) benign  (b) MI-FGSM  (c) JSMA  (d) PWA

Figure 13: The t-NSE visualization of different attacks on CNN1 of MNIST dataset, where the last three images are high-dimensional feature distributions of the adversarial examples crafted by MI-FGSM, JSMA, PWA respectively.

result of the benign and adversarial example, we can find that the attacks break up the distribution of classes while cutting down the distance between different classes. Whether it is the gradient-based white-box attack MI-FGSM or decision- based black-box attack PWA, they all fool the model by increasing the distance between classes.

② SF and TF of benign example



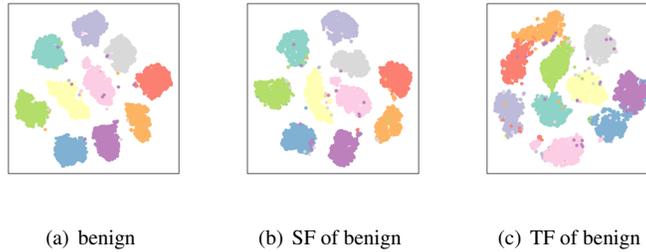(a) benign  (b) SF of benign  (c) TF of benign

Figure 14: The t-NSE visualization of SF and TF of benign examples on MNIST CNN1.

Fig. 14 shows the t-SNE visualization of benign examples and their corresponding SF and TF on MNIST-CNN1. From left to right are benign examples classified by the original model, SF and TF reconstructed by GAN.

It can be seen from the figure that the distribution of SF and TF of benign examples after reconstruction of GAN are similar, both scattered, which is consistent with our definition in Section 3.1. The distribution of the reconstructed SF resembles that of the original model, which indicates that the defense poses little negative impact on classification accuracy of benign examples.

③ The effectiveness of detection

Fig. 15 shows SF and TF of benign and adversarial examples on the MNIST-CNN1 model. As can be seen from the figure, the distribution of SF and TF of benign examples looks similar, while that of adversarial examples are quite different, which conforms to our previous definition as well. SFE detectors learn the differences between SF and TF, thus successfully detecting adversarial examples.

④ The effectiveness of defense

Fig. 16 shows the visualization results of VGG19 model on ImageNet before and after attack and
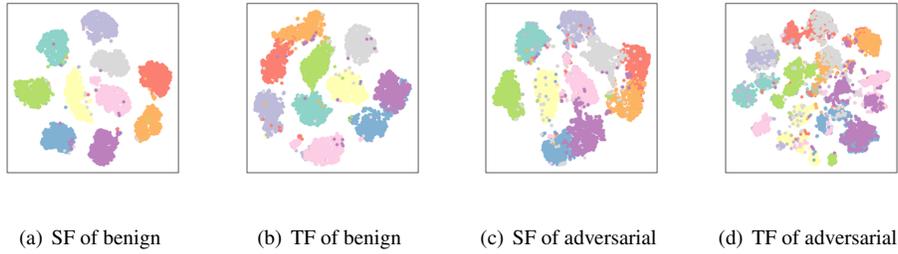
| (a) SF of benign | (b) TF of benign | (c) SF of adversarial | (d) TF of adversarial |

Figure 15: Reconstructed SF and TF of benign and adversarial examples on MNIST-CNN1. Adversarial examples are generated by PWA.



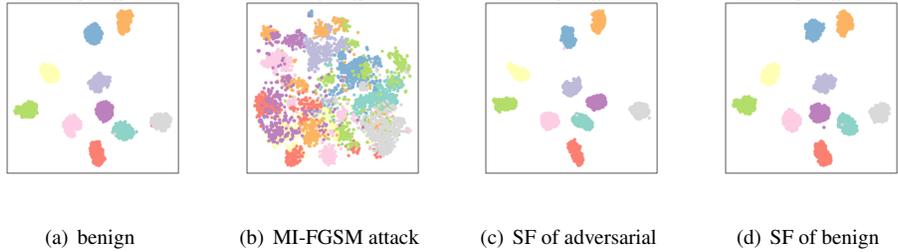| (a) benign | (b) MI-FGSM attack | (c) SF of adversarial | (d) SF of benign |

Figure 16: t-NSE visualization of ImageNet VGG 19 during the defense process. Adversarial examples are crafted by MI-FGSM.

after defense. From left to right are benign and adversarial examples crafted by MI-FGSM attack, SF of adversarial examples after GAN reconstruction and SF of reconstructed benign image.

After reconstructing by GAN, the distribution of t-SNE of benign and adversarial examples are similar to that of original model. The visualization result demonstrates the GAN's competitive performance of learning in terms of feature distribution: it will study the features of the adversarial examples to resemble SF of benign examples, so that the reconstructed SF of adversarial examples will be classified correctly when input to the model for re-identification defense. Besides, the SF of adversarial examples after reconstruction of GAN is similar to that of the benign example, consistent with the definition proposed above, which also verifies the validity of re-identification of SF during defense.

⑤ The effectiveness of defense

From the results of t-SNE visualization above, it can be concluded that images of a certain class are misclassified to another class by attacks. During that process, the cluster distribution is disrupted and the distance between classes is reduced. SFE remaps the distribution of different features, widens the distance between classes, and scatters the distribution, approximating the original model. Consequently, defense and detection of SFE are based on the characteristics of class distribution and independent of attack methods, which contributes to transferability among different attack algorithms.

### 4.10.2. Visualization of Grad-CAM

Grad-CAM provides heatmap visualization from the perspective of pixel-level features. Here, we provide more heatmaps generated by Grad-CAM on ImageNet for detailed indication, as shown in Figure

17. Areas more closely related to classification label and drawn more attention by the model are painted in red. From red to blue, the weight that the model allocates decreases. More visualization results will be shown in Appendix A.

It can be easily observed from the figure that red areas in SF of benign are quite similar to that in heatmap of benign example. This demonstrates that SF is relevant to class label, consistent with our definition and human semantics. Therefore, defense effect is guaranteed by the high similarity of highlight areas between reconstructed SF of adversarial examples and that in benign ones. On the contrary, TF are class-independent features. When attention is paid to those areas irrelevant to the label, misclassification happens, as shown in heatmap of adversarial.



(a) benign example     (b) heatmap of benign     (c) SF of benign     (d) TF of benign

(e) adversarial example     (f) heatmap of adv     (g) SF of adv     (h) TF of adv

(i) benign example     (j) heatmap of benign     (k) SF of benign     (l) TF of benign

(m) adversarial example     (n) heatmap of adv     (o) SF of adv     (p) TF of adv

Figure 17: Visualization of SF and TF on ImageNet VGG19 model via Grad-CAM. From left to right are the benign example and its heatmap, visualization of SF and TF reconstructed by SFE. The second row represents heatmaps of adversarial examples and their corresponding SF and TF results after reconstruction. From red to blue, the weight that the model allocates decreases. Adversarial examples are generated by DeepFool and BIM, respectively.

29

Based on the visualization results of t-SNE and Grad-CAM, we can answer **RQ4** with the following conclusions: 1) In the visualization of the high-dimensional features of the model, SFE successfully achieved the defense against the adversarial examples by reconstructing the feature distribution and increasing the distance between classes; 2) As for pixel-level feature in an image, SF is quite class-related while TF may cause misclassification. By reconstructing SF of adversarial exmaples, the defense effect is reached when it resembles a benign one.

## 5. Conclusions

In this paper, we propose the concepts of salient features and trivial features. We use GAN framework to extract these two features, and put forward the SFE method, which can detect and re-identify adversarial examples at the same time. Comprehensive experiments have shown that, compared with baselines, SFE achieves fairly good detection and defense effect for a variety of attack algorithms on different datasets and models, and has fairly good defense transferability among different attacks.

To improve the robustness of DNNs against attacks, as a future work, we plan further to improve the defense accuracy of SFE on complex datasets by optimizing the network structure and training strategy. Meanwhile, we would cut down the complexity of model structure and computation time to meet large-scale applications and computing needs.

## Acknowledgment

## References

[1] Chen, J., Wang, Z., Zheng, H., Xiao, J., Ming, Z., 2020a. Roby: Evaluating the robustness of a deep model by its decision boundaries. arXiv preprint arXiv:2012.10282 .

[2] Chen, J., Zheng, H., Shangguan, W., Liu, L., Ji, S., 2021a. Act-detector: Adaptive channel transformation-based light-weighted detector for adversarial attacks. Information Sciences .

[3] Chen, J., Zheng, H., Xiong, H., Chen, R., Du, T., Hong, Z., Ji, S., 2021b. Finefool: A novel dnn object contour attack on image recognition based on the attention perturbation adversarial technique. Computers & Security , 102220.

[4] Chen, J., Zheng, H., Xiong, H., Shen, S., Su, M., 2020b. Mag-gan: Massive attack generator via gan. Information Sciences 536, 67–90. URL: `http://dx.doi.org/10.1016/j.ins.2020.04.019`.

[5] Chen, P., Zhang, H., Sharma, Y., Yi, J., Hsieh, C., 2017. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models, in: Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security, AISec@CCS 2017, Dallas, TX, USA, November 3, 2017, ACM. pp. 15–26. URL: `https://doi.org/10.1145/3128572.3140448`, doi:`10.1145/3128572.3140448`.

[6] Dhillon, G.S., Azizzadenesheli, K., Lipton, Z.C., Bernstein, J., Kossaifi, J., Khanna, A., Anandkumar, A., 2018. Stochastic activation pruning for robust adversarial defense. arXiv preprint arXiv:1803.01442 .

[7] Dong, Y., Liao, F., Pang, T., Su, H., Zhu, J., Hu, X., Li, J., 2018. Boosting adversarial attacks with momentum, in: 2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018, IEEE Computer Society. pp. 9185–9193. URL: `http://dx.doi.org/10.1109/CVPR.2018.00957`.

[8] Dongyu, M., Hao, C., 2017. Magnet: A two-pronged defense against adversarial examples, in: Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, CCS 2017, Dallas, TX, USA, October 30 - November 03, 2017, ACM. pp. 135–147. URL: `https://doi.org/10.1145/3133956.3134057`.

[9] Eykholt, K., Evtimov, I., Fernandes, E., Li, B., Rahmati, A., Xiao, C., Prakash, A., Kohno, T., Song, D., 2018. Robust physical-world attacks on deep learning visual classification, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1625–1634.

[10] Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y., 2014. Generative adversarial networks. arXiv preprint arXiv:1406.2661 .

[11] Goodfellow, I.J., Shlens, J., Szegedy, C., 2015. Explaining and harnessing adversarial examples, in: 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, pp. 1–10. URL: `http://arxiv.org/abs/1412.6572`.

[12] Gui, J., Sun, Z., Wen, Y., Tao, D., Ye, J., 2020. A review on generative adversarial networks: Algorithms, theory, and applications. arXiv preprint arXiv:2001.06937 .

[13] Hlihor, P., Volpi, R., Malagò, L., 2020. Evaluating the robustness of defense mechanisms based on autoencoder reconstructions against carlini-wagner adversarial attacks, in: Proceedings of the Northern Lights Deep Learning Workshop, pp. 6–6.

[14] Huang, R., Xu, B., Schuurmans, D., Szepesvári, C., 2015. Learning with a strong adversary. arXiv preprint arXiv:1511.03034 .

[15] Ilyas, A., Santurkar, S., Tsipras, D., Engstrom, L., Tran, B., Madry, A., 2019. Adversarial examples are not bugs, they are features. arXiv preprint arXiv:1905.02175 .

[16] Jin, G., Shen, S., Zhang, D., Dai, F., Zhang, Y., 2019. Ape-gan: Adversarial perturbation elimination with gan, in: IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2019, Brighton, United Kingdom, May 12-17, 2019, IEEE. pp. 3842–3846. URL: https://doi.org/10.1109/ICASSP.2019.8683044, doi:10.1109/ICASSP.2019.8683044.

[17] Jinyin, C., Mengmeng, S., Shijing, S., Hui, X., Haibin, Z., 2019. Poba-ga: Perturbation optimized black-box adversarial attacks via genetic algorithm. Computers and Security 85, 89–106. URL: https://doi.org/10.1016/j.cose.2019.04.014.

[18] Jonas, R., Wieland, B., Behar, V., Evgenia, R., a. Additive uniform noise attack in foolbox tool. https://foolbox.readthedocs.io/en/v1.8.0/modules/attacks/decision.html#foolbox.attacks.AdditiveUniformNoiseAttack.

[19] Jonas, R., Wieland, B., Behar, V., Evgenia, R., b. Contrast reduction attack in foolbox tool. https://foolbox.readthedocs.io/en/v1.8.0/modules/attacks/decision.html#foolbox.attacks.ContrastReductionAttack.

[20] Ju, C., Bibaut, A., van der Laan, M., 2018. The relative performance of ensemble methods with deep convolutional neural networks for image classification. Journal of Applied Statistics 45, 2800–2818. URL: http://dx.doi.org/10.1080/02664763.2018.1441383.

[21] Krizhevsky, A., Sutskever, I., Hinton, G.E., 2012. Imagenet classification with deep convolutional neural networks, in: Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States, pp. 1106–1114. URL: http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.

[22] Kurakin, A., Goodfellow, I.J., Bengio, S., 2017. Adversarial examples in the physical world, in: 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Workshop Track Proceedings, OpenReview.net. pp. 1–14. URL: https://openreview.net/forum?id=HJGU3Rodl.

[23] Liu, M.Y., Tuzel, O., 2016. Coupled generative adversarial networks , 469–477URL: https://dl.acm.org/doi/abs/10.5555/3157096.3157149.

[24] Liu, X., Hsieh, C.J., 2019. Rob-gan: Generator, discriminator, and adversarial attacker, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 11234–11243. URL: `https://openaccess.thecvf.com/content_CVPR_2019/html/Liu_Rob-GAN_Generator_Discriminator_and_Adversarial_Attacker_CVPR_2019_paper.html`.

[25] Van der Maaten, L., Hinton, G., 2008. Visualizing data using t-sne. Journal of machine learning research 9.

[26] Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A., 2018. Towards deep learning models resistant to adversarial attacks, in: 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30-May 3, 2018, OpenReview.net. pp. 1–28. URL: `https://openreview.net/forum?id=rJzIBfZAb`.

[27] Metzen, J.H., Genewein, T., Fischer, V., Bischoff, B., 2017. On detecting adversarial perturbations, in: 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings, OpenReview.net. pp. 1–12. URL: `https://openreview.net/forum?id=SJzCSf9xg`.

[28] Moosavi-Dezfooli, S.M., Fawzi, A., Fawzi, O., Frossard, P., 2017. Universal adversarial perturbations, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017, IEEE Computer Society. pp. 86–94. URL: `http://dx.doi.org/10.1109/CVPR.2017.17`.

[29] Moosavi-Dezfooli, S.M., Fawzi, A., Frossard, P., 2016. Deepfool: a simple and accurate method to fool deep neural networks, in: 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016, IEEE Computer Society. pp. 2574–2582. URL: `https://doi.org/10.1109/CVPR.2016.282`.

[30] Mummadi, C.K., Brox, T., Metzen, J.H., 2019. Defending against universal perturbations with shared adversarial training, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 4928–4937. URL: `https://openaccess.thecvf.com/content_ICCV_2019/html/Mummadi_Defending_Against_Universal_Perturbations_With_Shared_Adversarial_Training_ICCV_2019_paper.html`.

[31] Narodytska, N., Kasiviswanathan, S.P., 2016. Simple black-box adversarial perturbations for deep networks. arXiv preprint arXiv:1612.06299 .

[32] Pang, T., Du, C., Dong, Y., Zhu, J., 2018. Towards robust detection of adversarial examples, in: Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December

2018, Montréal, Canada, pp. 4579–4589. URL: `http://papers.nips.cc/paper/7709-towards-robust-detection-of-adversarial-examples`.

[33] Papernot, N., McDaniel, P., Jha, S., Fredrikson, M., Celik, Z.B., Swami, A., 2016a. The limitations of deep learning in adversarial settings, in: 2016 IEEE European symposium on security and privacy (EuroS&P), IEEE. pp. 372–387. URL: `https://10.1109/EuroSP.2016.36`.

[34] Papernot, N., McDaniel, P., Wu, X., Jha, S., Swami, A., 2016b. Distillation as a defense to adversarial perturbations against deep neural networks, in: IEEE Symposium on Security and Privacy, SP 2016, San Jose, CA, USA, May 22-26, 2016, IEEE Computer Society. pp. 582–597. URL: `http://dx.doi.org/10.1109/SP.2016.41`, doi:`10.1109/SP.2016.41`.

[35] Prakash, A., Moran, N., Garber, S., DiLillo, A., Storer, J., 2018. Deflecting adversarial attacks with pixel deflection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 8571–8580. URL: `https://openaccess.thecvf.com/content_cvpr_2018/html/Prakash_Deflecting_Adversarial_Attacks_CVPR_2018_paper.html`.

[36] Qian, Y., Ma, D., Wang, B., Pan, J., Wang, J., Gu, Z., Chen, J., Zhou, W., Lei, J., 2020. Spot evasion attacks: Adversarial examples for license plate recognition systems with convolutional neural networks. Computers & Security 95, 101826.

[37] Rozsa, A., Günther, M., Rudd, E.M., Boult, T.E., 2019. Facial attributes: Accuracy and adversarial robustness. Pattern Recognition Letters 124, 100–108.

[38] Samangouei, P., Kabkab, M., Chellappa, R., 2018. Defense-gan: Protecting classifiers against adversarial attacks using generative models, in: International Conference on Learning Representations. URL: `https://openreview.net/forum?id=BkJ3ibb0-`.

[39] Schott, L., Rauber, J., Bethge, M., Brendel, W., 2018. Towards the first adversarially robust neural network model on mnist. arXiv preprint arXiv:1805.09190 .

[40] Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D., 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization, in: IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017, IEEE Computer Society. pp. 618–626. URL: `http://dx.doi.org/10.1109/ICCV.2017.74`, doi:`10.1109/ICCV.2017.74`.

[41] Simonyan, K., Zisserman, A., 2015. Very deep convolutional networks for large-scale image recognition, in: 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, pp. 1–14. URL: `http://arxiv.org/abs/1409.1556`.

[42] Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z., 2016. Rethinking the inception architecture for computer vision, in: 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016, IEEE Computer Society. pp. 2818–2826. URL: http://dx.doi.org/10.1109/CVPR.2016.308.

[43] Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R., 2013. Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199 .

[44] Tian, S., Yang, G., Cai, Y., 2018. Detecting adversarial examples through image transformation, in: Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018, AAAI Press. pp. 4139–4146. URL: https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/17408.

[45] Tramèr, F., Kurakin, A., Papernot, N., Goodfellow, I., Boneh, D., McDaniel, P., 2018. Ensemble adversarial training: Attacks and defenses, in: 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings, OpenReview.net. pp. 1–20. URL: https://openreview.net/forum?id=rkZvSe-RZ.

[46] Wei, X., Guo, Y., Li, B., 2021. Black-box adversarial attacks by manipulating image attributes. Information Sciences 550, 285–296.

[47] Xiao, C., Li, B., Zhu, J.Y., He, W., Liu, M., Song, D., 2018. Generating adversarial examples with adversarial networks, in: Proceedings of the 27th International Joint Conference on Artificial Intelligence, pp. 3905–3911. URL: https://dl.acm.org/doi/abs/10.5555/3304222.3304312.

[48] Xiao, Y., Pun, C.M., Liu, B., 2020. Adversarial example generation with adaptive gradient search for single and ensemble deep neural network. Information Sciences 528, 147–167.

[49] Xie, C., Wang, J., Zhang, Z., Ren, Z., Yuille, A., 2017. Mitigating adversarial effects through randomization. arXiv preprint arXiv:1711.01991 .

[50] Yi, Z., Zhang, H., Tan, P., Gong, M., 2017. Dualgan: Unsupervised dual learning for image-to-image translation, in: Proceedings of the IEEE international conference on computer vision, pp. 2849–2857.

[51] Yin, X., Kolouri, S., Rohde, G.K., 2020. Gat: Generative adversarisl training for adversarial example detection and robust classification, in: Proceedings of the Eighth International Conference on Learning Representations, pp. 1–26.

[52] Zhang, X., Wang, J., Wang, T., Jiang, R., Xu, J., Zhao, L., 2021. Robust feature learning for adversarial defense via hierarchical feature alignment. Information Sciences 560, 256–270.

# Appendix A. More visualization of SF and TF

*Appendix A.1. Visualization of CIFAR-10*

Here we show visualizations of CIFAR-10 on VGG19 model. From left to right are the benign example and its heatmap of SF and TF in adversarial exmaples. From red to blue, the weight that the model allocates decreases. Adversarial examples are generated by all attacks used in Section 4.1.



| (a) benign example | (b) SF of adv | (c) TF of adv | (d) benign example | (e) SF of adv | (f) TF of adv |
| (g) benign example | (h) SF of adv | (i) TF of adv | (j) benign example | (k) SF of adv | (l) TF of adv |
| (m) benign example | (n) SF of adv | (o) TF of adv | (p) benign example | (q) SF of adv | (r) TF of adv |
| (s) benign example | (t) SF of adv | (u) TF of adv | (v) benign example | (w) SF of adv | (x) TF of adv |
| (y) benign example | (z) SF of adv | (aa) TF of adv | (ab) benign example | (ac) SF of adv | (ad) TF of adv |

Figure A.18: Visualization of SF and TF on CIFAR-10 VGG19 model via Grad-CAM.

The supplementary visualization of SF and TF on ImageNet VGG19 model via Grad-CAM are shown as follows. From left to right are the benign example and its heatmap of SF and TF in adversarial exmaples. From red to blue, the weight that the model allocates decreases. Adversarial examples are generated by all attacks used in Section 4.1.
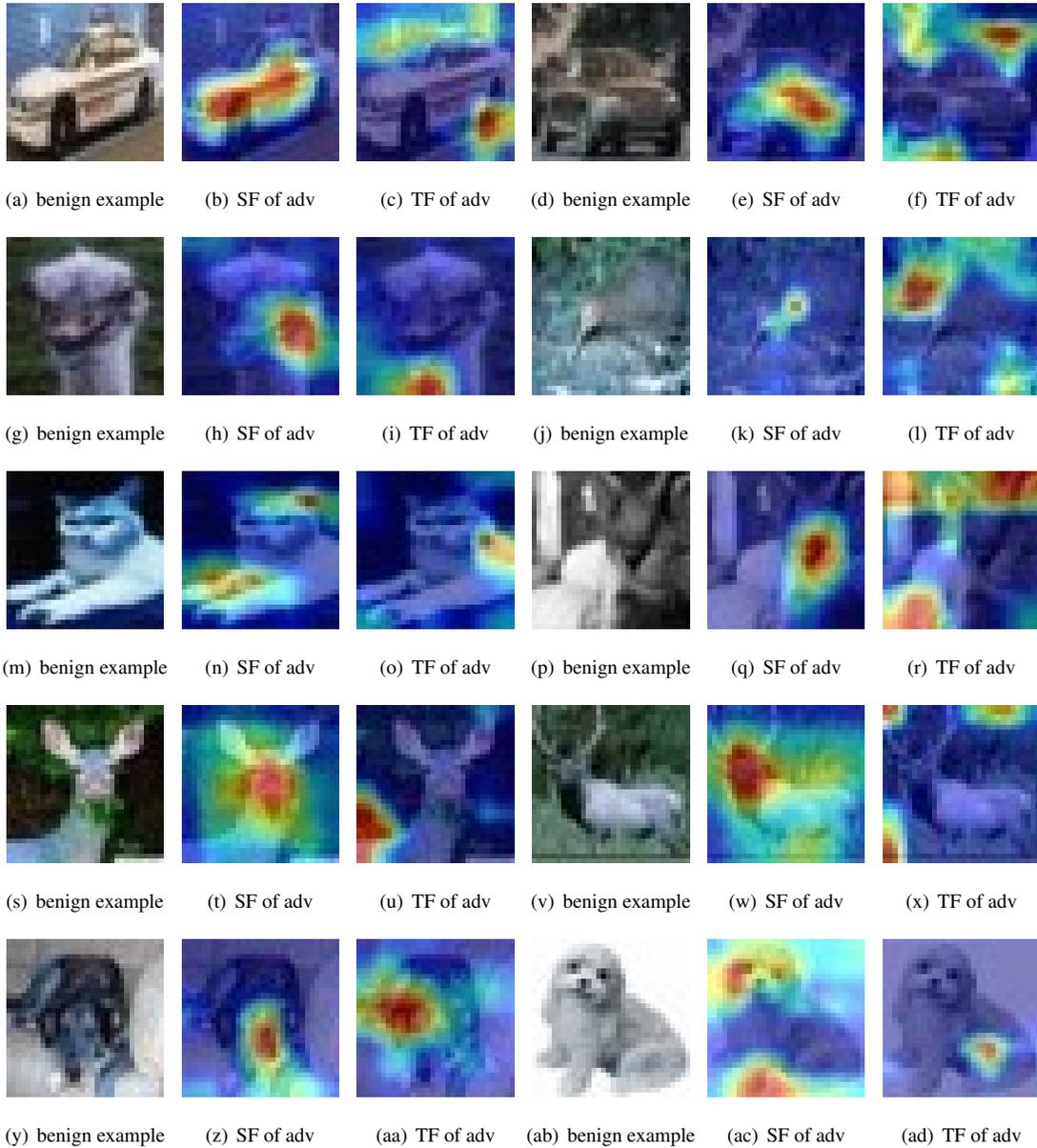


(a) benign example    (b) SF of adv    (c) TF of adv    (d) benign example    (e) SF of adv    (f) TF of adv

(g) benign example    (h) SF of adv    (i) TF of adv    (j) benign example    (k) SF of adv    (l) TF of adv

(m) benign example    (n) SF of adv    (o) TF of adv    (p) benign example    (q) SF of adv    (r) TF of adv

(s) benign example    (t) SF of adv    (u) TF of adv    (v) benign example    (w) SF of adv    (x) TF of adv

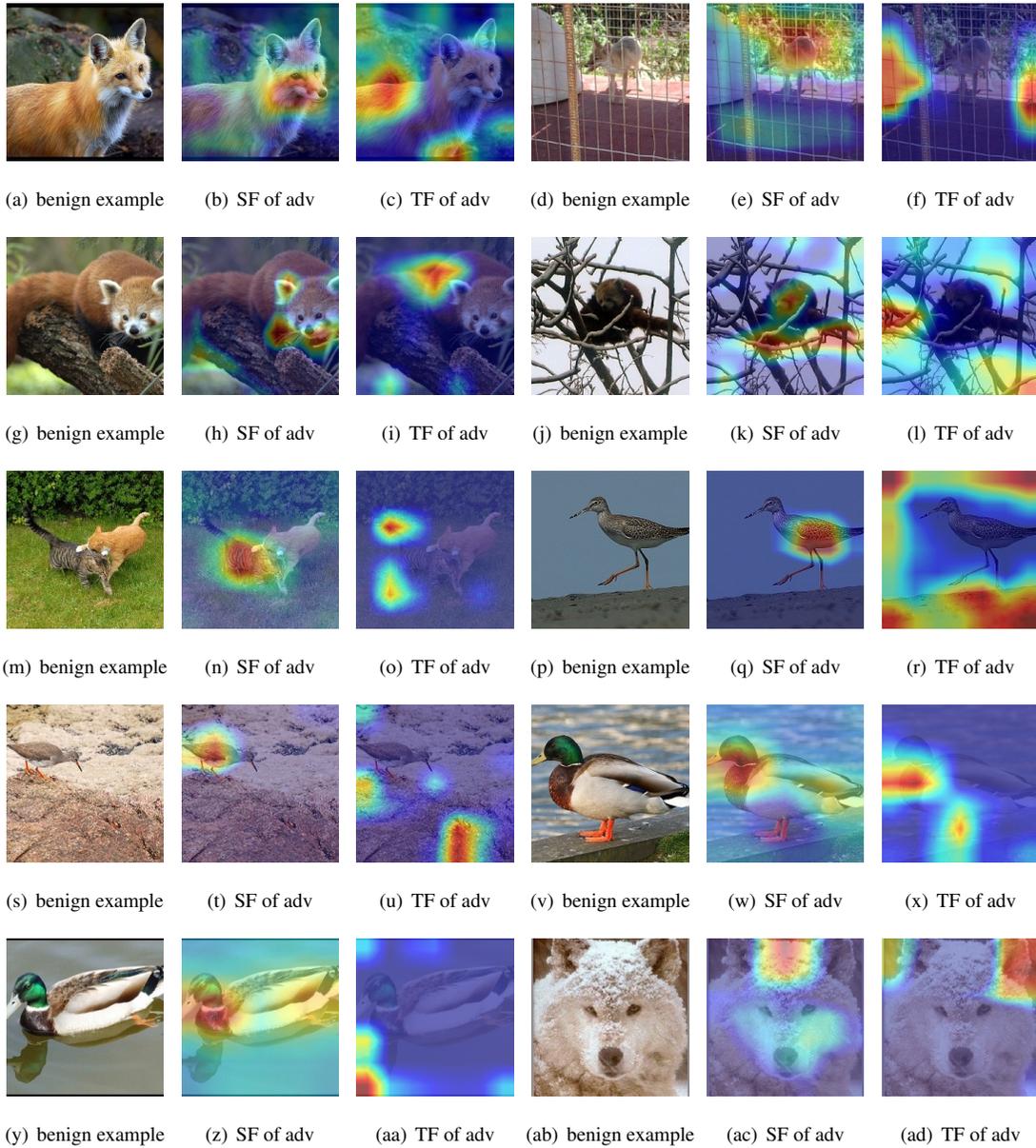(y) benign example    (z) SF of adv    (aa) TF of adv    (ab) benign example    (ac) SF of adv    (ad) TF of adv

Figure A.19: Visualization of SF and TF on ImageNet VGG19 model via Grad-CAM.