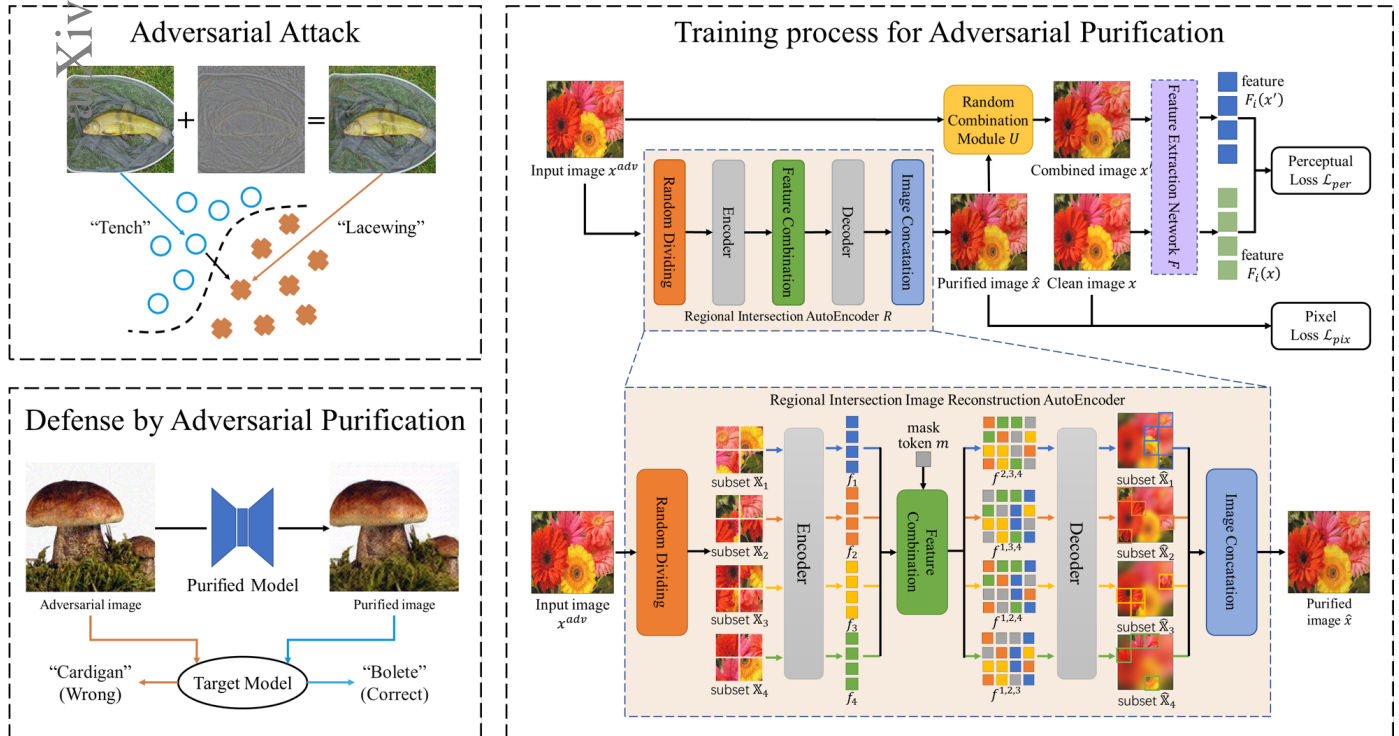


Graphical Abstract

Adversarial Purification of Information Masking

Sitong Liu, Zhichao Lian, Shuangquan Zhang, Liang Xiao



Highlights

Adversarial Purification of Information Masking

Sitong Liu, Zhichao Lian, Shuangquan Zhang, Liang Xiao

- Establishes a quantitative correlation between residual adversarial perturbation scale and attack capability.
- Introduces various information masking strategies to resist the same-position adversarial perturbation.
- Advances a regional intersection reconstruction approach aimed at diminishing the scale of residual same-position adversarial perturbations.
- Proposes a simulated technique of residual perturbation to fortify defenses against content-similar adversarial perturbations.
- Explores a joint constraint form encompassing pixel loss and perceptual loss to augment the flexibility in clean sample generation.

Adversarial Purification of Information Masking

Sitong Liu^a, Zhichao Lian^{a,*}, Shuangquan Zhang^a, Liang Xiao^{a,*}

^aorganization=School of Computer Science and Engineering, Nanjing University of Science and Technology, addressline=No.200 Xiaolingwei Street, city=Nanjing, postcode=210094, state=Jiangsu, country=China

Abstract

Adversarial attacks meticulously generate minuscule, imperceptible perturbations to images to deceive neural networks. Counteracting these, adversarial purification methods seek to transform adversarial input samples into clean output images to defend against adversarial attacks. Nonetheless, extent generative models fail to effectively eliminate adversarial perturbations, yielding less-than-ideal purification results. We emphasize the potential threat of residual adversarial perturbations to target models, quantitatively establishing a relationship between perturbation scale and attack capability. Notably, the residual perturbations on the purified image primarily stem from the same-position patch and similar patches of the adversarial sample. We propose a novel adversarial purification approach named Information Mask Purification (IMPure), aims to extensively eliminate adversarial perturbations. To obtain an adversarial sample, we first mask part of the patches information, then reconstruct the patches to resist adversarial perturbations from the patches. We reconstruct all patches in parallel to obtain a cohesive image. Then, in order to protect the purified samples against potential similar regional perturbations, we simulate this risk by randomly mixing the purified samples with the input samples before inputting them into the feature extraction network. Finally, we establish a combined constraint of pixel loss and perceptual loss to augment the model’s reconstruction adaptability. Extensive experiments on the ImageNet dataset with three classifier models demonstrate that our approach achieves state-of-the-art results against nine adversarial attack methods. Implementation code and pre-trained weights can be accessed at <https://github.com/NoWindButRain/IMPure>.

Keywords: adversarial purification, adversarial attacks and defenses, adversarial machine learning, image reconstruction

1. Introduction

Neural network-based machine learning methodologies have become instrumental in a multitude of domains, including autonomous driving [36] and access control systems [18]. Nonetheless, emerging research has revealed the vulnerability of neural networks to adversarial attacks [13]. As can be seen from Fig. 1, even slight perturbations can induce erroneous outputs from the target model, diverging from original sample predictions. Adversarial attack strategies possess not only the capability to deceive neural networks but also transferable applicability, capable of inflicting damage on black-box models without the requisite knowledge of the target model’s architecture or parameters [35, 5, 2]. Given the significant security vulnerabilities these adversarial attacks introduce, the development of efficient defense mechanisms becomes vital, aiming to enhance the adversarial robustness of neural networks.

One of the primary objectives of defensive methods is to ensure the target model’s immunity to adversarial

perturbations, maintaining accurate outputs, especially in high-stakes applications like autonomous driving and access control. Defenses can generally be categorized into model-specific and model-agnostic strategies [15]. Model-specific methods, including adversarial training [13] and gradient masking [9], augment the robustness against adversarial attacks by modifying the training strategy of the target model. However, Model-specific methods have notable downsides like a negative robustness-accuracy correlation [45] and the necessity of retraining, which can be resource-intensive. Model-agnostic methods, such as JPEG compression [12], random padding [47], and high-level representation guided denoising (HGD) [23], operate at the input stage, preprocessing input samples to diminish adversarial perturbations’ impact. Model-agnostic methods, while intuitive and seamlessly integrated, often struggle to entirely clean adversarial perturbations and preserve image quality, while also being potentially susceptible to compromise in a white-box setting.

Generative model-based adversarial purification methods are an important part of model-agnostic methods. Generative models have seen extensive deployment across various computer vision and deep learning tasks, including image denoising [49] and image restoring [4]. As shown

*Corresponding author

Email addresses: 1stnjust@163.com (Sitong Liu), lzcts@163.com (Zhichao Lian), zhangsq@njust.edu.cn (Shuangquan Zhang), xiaoliang@njust.edu.cn (Liang Xiao)

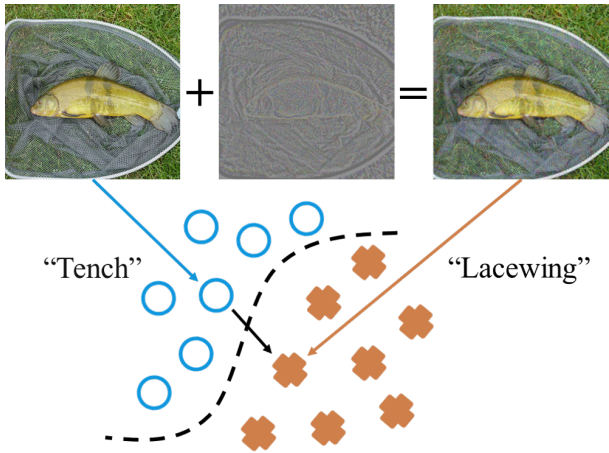


Figure 1: Adversarial Attack Process: The illustration depicts a pair of images—the pristine image on the left is accurately classified as a tench by the target model. Conversely, the image on the right showcases an adversarial example, derived post the infusion of a minor perturbation. This adversarial perturbation alters the feature mapping, propelling it beyond the decision boundary delineated by the target model, and consequently, the adversarial sample is misclassified as a lacewing.

in Fig. 2, Viewing adversarial perturbations as a distinct type of noise, researchers have proposed adversarial purification methods using generative models to denoise or reconstruct adversarial samples prior to input into the target model. However, the pixel-level loss functions, typically employed in traditional generative tasks, fall short in adversarial purification due to the error amplification effect. Some approaches [23, 50] employ perceptual loss functions, calculating feature differences in target model intermediate layers as a remedy. It’s noteworthy that the noise scale in purified images might exceed that of the original adversarial images, not only degrading image quality but also potentially undermining the accuracy of the target model. In summary, our adversarial purification methodology is structured around three primary objectives: 1) To extensively eliminate adversarial perturbations present in the input image; 2) To effectively defend against any residual adversarial perturbations; 3) To minimize the discrepancy between the purified image and its original, clean counterpart.

In this paper, we delve into the detrimental impact of residual adversarial perturbations on purified images, furnishing rigorous proof that articulates the quantitative relationship between the scale of residual perturbations and the capabilities of the attack. Consequently, we emphasize the importance of limiting the scale of adversarial perturbations for adversarial purification methods. The residual perturbations observed on the purified image patches come from same-position patches and content-similar patches in the adversarial image. We propose a novel adversarial purification method, termed Information Mask Purification (IMPure), devised to resist same-position and content-similar perturbations through

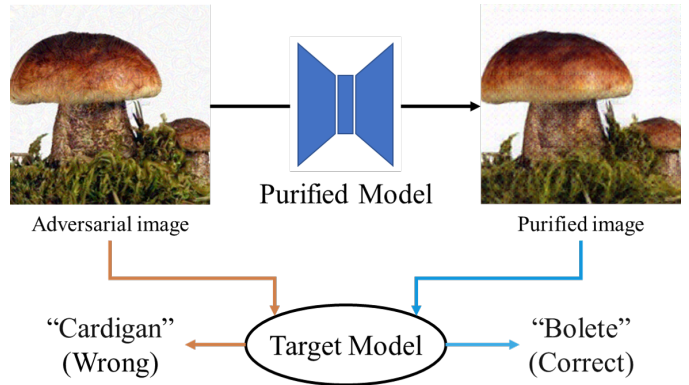


Figure 2: Adversarial Attack Process: The illustration depicts a pair of images—the pristine image on the left is accurately classified as a tench by the target model. Conversely, the image on the right showcases an adversarial example, derived post the infusion of a minor perturbation. This adversarial perturbation alters the feature mapping, propelling it beyond the decision boundary delineated by the target model, and consequently, the adversarial sample is misclassified as a lacewing.

the meticulous design of a Regional Intersection image reconstruction Autoencoder (RIAE) and a Random Combination Module (RCM), respectively. RIAE is based on the Masked Autoencoder (MAE) [16] and showcases robust reconstruction capabilities through a reconfigured and parallelized architecture. We selectively mask part of the adversarial sample patch information, using the remaining patches to reconstruct the masked patch and parallelizing the reconstruction process on all patches to obtain a complete purified image. Given the minute pixel values of the adversarial perturbation, the influence of the perturbation can be effectively mitigated through information masking. The image information loss of the masked patches can be reconstructed by leveraging the information from the remaining intact patches through the reconstruction network. Identifying content-similar patches poses a challenge, yet RCM achieves resistance to content-similar perturbations indirectly. Operational during the training phase, RCM strategically merges the purified image with the adversarial image, deliberately preserving some adversarial perturbations. Thus, to counteract the impact of preservational perturbations on the target model, the reconstruction network needs to amplify the overall defense capability of the purified image. Furthermore, our proposed method proffers enhanced flexibility in selecting feature extraction networks relative to existing adversarial purification methods. We employ joint constraints of pixel loss and perceptual loss, facilitating the integration of both low-level pixel information and high-level perceptual information throughout the training phase, thereby helping the model on a virtuous cycle.

Our work has made the following main contributions:

1. We meticulously explore the perils posed by residual adversarial perturbations within adversarial purification strategies, elucidating a quantitative nexus be-

tween perturbation scale and offensive capabilities.

2. We propose IMPure, an innovative adversarial purification method, which consists of RIAE that maximizes the elimination of same-position perturbations and RCM that encourages the network to improve the defense capabilities to resist content-similar perturbations. Additionally, we provide a flexible joint-constrained loss to promote the consistency of purified images in both low-level information and high-level representation.
3. A suite of comprehensive experiments demonstrates that our method achieves excellent defense performance on three target networks and nine adversarial attack methods in the classification task of the ImageNet dataset.

2. Related works

In this section, we introduce representative methods for adversarial attack and defense. To simplify the problem, we focus only on the image classification task. First, we designate some notations used in this paper. Let \mathbf{x} denote the original clean image from a given dataset, and y denotes the corresponding label. A neural network $f : \mathbf{x} \rightarrow y$ is called the target model. $\mathcal{L}(\mathbf{x}, y)$ denotes the loss function of the network. $p(y|\mathbf{x})$ predicted probability of class y output by the network f . $y_{\mathbf{x}} = \arg \max_y p(y|\mathbf{x})$ is the predicted class of \mathbf{x} . \mathbf{x}^{adv} denotes an adversarial example of \mathbf{x} .

2.1. Adversarial attack

The adversarial attack method generates an adversarial sample \mathbf{x}^{adv} on a clean image by carefully constructing an adversarial perturbation that causes the classifier f to give wrong predictions $y_{\mathbf{x}^{adv}}$. Based on their purpose and environment setting respectively, adversarial attack methods are categorised, and some representative methods are introduced in this section.

Adversarial attacks can be divided into targeted and untargeted attacks according to their purposes. The goal of an untargeted attack is $y_{\mathbf{x}^{adv}} \neq y$. The target attack specifies a specific class y' and makes $y_{\mathbf{x}^{adv}} = y'$.

Goodfellow et al. [13] suggested the cumulative effects of high dimensional model weights and the deep model is more linear in high dimensions. They proposed an untargeted method called the Fast Gradient Sign Method (FGSM) to generate adversarial examples only by computing the gradient once:

$$\mathbf{x}^{adv} = \mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, y)) \quad (1)$$

where ϵ is the strength of the perturbation, $\text{sign}(\cdot)$ represents a symbolic function. By modifying the FGSM to maximize the probability of a specific class y^{target} , a version of the target attack can be obtained:

$$\mathbf{x}^{adv} = \mathbf{x} - \epsilon \text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, y^{target})) \quad (2)$$

Kurakin et al. [22] proposed a Basic iterative method (BIM) which is an iterative FGSM attack to achieve stronger attack effects through multiple rounds of small steps:

$$\mathbf{x}_{t+1}^{adv} = \text{Clip}_{\mathbf{x}, \epsilon} \{ \mathbf{x}_t^{adv} + \alpha \text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_t^{adv}, y)) \} \quad (3)$$

where \mathbf{x}_t^{adv} represents the adversarial sample in the t_{th} iteration, α represents the step size, and Clip represents the truncation method.

Adversarial attacks can be broadly categorized based on the attack environment into white-box and black-box attacks. In white-box attacks, adversaries possess complete knowledge of the target model, while black-box attacks operate under the premise that the attackers lack direct access to the model’s specific details.

Athalye et al. [1] introduced a gradient approximation method that adeptly circumvented the obfuscated gradients defense in a white-box context. They classified obfuscated gradient defenses into three types: broken gradients, random gradients, and vanishing/exploding gradients. Additionally, they proposed two distinctive attack strategies: Backward Pass Differentiable Approximation (BPDA) and Expectation Over Transformation (EOT).

A salient feature of adversarial samples is their inherent transferability across varying model architectures and parameterizations. Transferability paves the way for the feasibility of black-box attacks [25], even facilitating their extension from digital settings to real-world scenarios [44].

2.2. Defense against adversarial attack

For the deployment of intelligent models in real-world scenarios, it’s imperative that models maintain accurate predictions even when confronted with adversarial attacks. Ilyas et al. [19] referred to the model’s ability to resist adversarial attacks as adversarial robustness.

Adversarial training [13], a model-specific defense, has garnered considerable attention. It seeks to bolster a model’s adversarial robustness by retraining on the training set with adversarial samples added. However, recent studies have highlighted its limitations. Training on large datasets like ImageNet, coupled with the generation of ample adversarial samples, can be resource-intensive [46]. Wong et al. [38] pointed out the potential detriment of excessive adversarial training, while Schmidt et al. [39] observed that achieving pronounced adversarial robustness demands a significantly larger sample size compared to standard high-accuracy models. Furthermore, Tsipras et al. [45] noted that model accuracy and adversarial robustness are negatively correlated. Such challenges impede the widespread adoption of adversarial training.

Model-agnostic defenses primarily revolve around pre-processing the input image to neutralize adversarial perturbations. Pioneering efforts revealed that basic image transformations, such as random padding [47] and JPEG compression [12], could effectively thwart FGSM attacks.

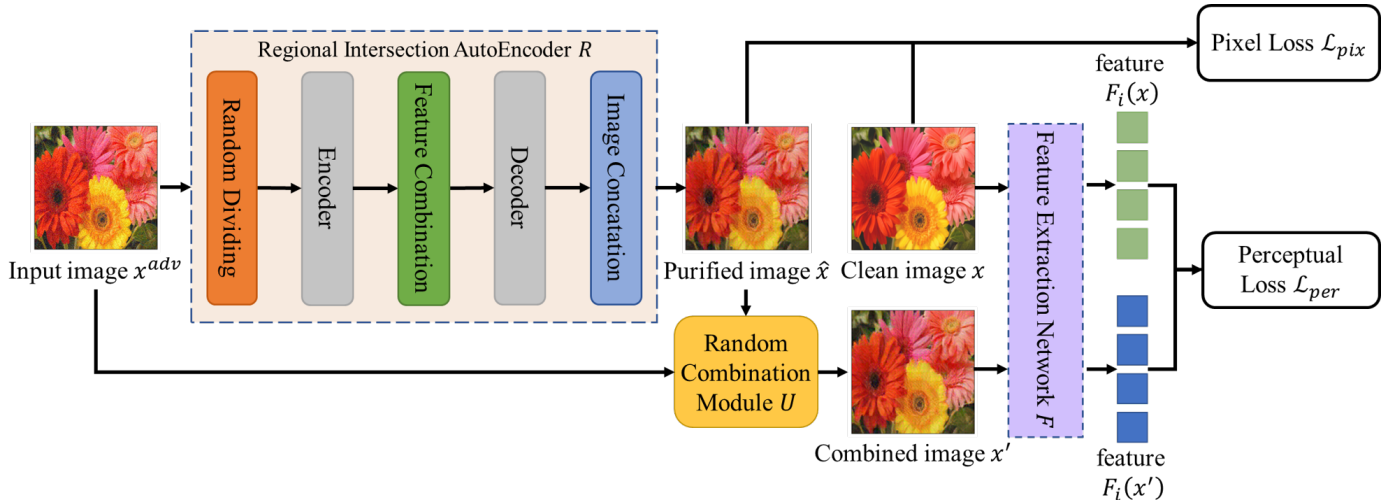


Figure 3: The framework of our adversarial purification method IMPure. We input the adversarial image into the image reconstruction network to get the purified image and calculate the pixel loss with the clean image stable low-level information. Then the adversarial image and the purified image are input to RCM to obtain the combined image. The feature extraction network extracts high-level features from combined and clean images and computes the loss to stabilize the semantic information.

However, these rudimentary methods falter against sophisticated adversarial techniques and can also degrade image quality. Subsequent research shifted towards sophisticated approaches like denoising and reconstruction to restore clean images. For instance, Gu et al. [14] first proposed the use of denoising autoencoders to remove adversarial perturbations. Liao et al. [23] pointed out the disadvantages of pixel loss and introduced the High-Level Representation Guided Denoiser (HGD). Mustafa et al. [32] combined wavelet denoising with super-resolution to enhance image quality, though its efficacy against potent strong adversarial attacks remains questionable. Jie et al. [20] rolled out ComDefend, an end-to-end image compression technique, that sidesteps the need for a vast adversarial sample collection, albeit with limited defensive capabilities. Meng et al. [30] put forth MagNet, encompassing detectors and reformers to stave off adversarial attacks while preserving accuracy on standard inputs. Zhang et al. [50] designed an image reconstruction network to eliminate the effects of adversarial perturbations. Nie et al. [33] proposed that DiffPure purifies adversarial samples through forward and reverse processes of diffusion models, and residual perturbations may still attack successfully. Model-agnostic methods offer commendable versatility, yet several challenges persist.

In this paper, we present a method aimed at addressing the challenges inherent in the aforementioned defense mechanisms. Firstly, the employment of perceptual loss often leads to the loss of fine-grained details, subsequently compromising model accuracy. To mitigate this, we introduce a dual constraint encompassing both pixel loss and perceptual loss, striving to retain essential image information. Secondly, fully eradicating adversarial perturbations remains challenging, and purifying images have potential security risks. In response, we devise RCM to emulate po-

tential disturbance risks, fostering an active response from reconstruction models. Lastly, the defense of strong adversarial attacks is more challenging. We construct RIAE to resist the impact of strong adversarial perturbations through information masking and regional intersection reconstruction.

3. Proposed Method

In this section, we first probe the impact of residual adversarial perturbations on the purified image and furnish a quantitative analysis linking perturbation scale to attack capability (Section 3.1). Recognizing the pivotal role of curtailing residual perturbation scales in enhancing adversarial purification efficacy, we propose IMPure (Section 3.2). A schematic representation of our proposed defense model is depicted in Fig. 3. Given the difficulty of diminishing the aggregate perturbation scale of a purified image, our approach pivots on image patch processing. The perturbations in an image patch predominantly originate from the same-position patch and content-similar patch of the adversarial sample. To tackle perturbations from identical position patches, we put forth RIAE, visualized in Fig. 4. Adversarial perturbations are mitigated by masking part of the image patches information. Then, the reconstruction model uses the information from the remaining image patches to reconstruct the masked image patches. The reconstruction process is executed in parallel on all patches to obtain an overall purified image. Addressing content-similar patches, RCM is introduced, bypassing intricate content computations. Operational during the training phase, RCM intentionally retains part of the adversarial perturbations on the purified image to simulate perturbations from content-similar patches. RIAE generates the purified image with stronger defense capabilities to resist

the adverse effects of preservation perturbations on the target model. Conclusively, we design a joint constrained loss integrating pixel loss and perceptual loss, and explore a more flexible perceptual loss design scheme. The joint loss is designed to encourage the model to achieve a virtuous cycle between low-level information and high-level representation.

3.1. Theoretical Analysis

We begin with a brief overview of adversarial attack principles and the mechanisms of adversarial purification. Subsequently, we delve into the potential risks posed by residual adversarial perturbations on purified images to the target model. We establish a quantifiable correlation between perturbation scale and attack efficacy using a streamlined target model. Additionally, we discuss the influence of different residual conditions of perturbations on attack capability, leading us to assert the importance of constraining the scale of residual disturbances.

Adversarial attack methodologies are orchestrated with the primary intent of manipulating the target models to yield incorrect predictions. These techniques can generate small perturbations, which are added to clean images to obtain adversarial samples. It is difficult for the human eyes to detect the difference between clean images and adversarial samples, but adversarial samples can induce the target model to make wrong predictions. To formalize, consider an efficient target model T and an input image x , yielding a prediction $y = T(x)$. An adversarial perturbation η is derived via an adversarial attack methodology and generates an adversarial sample $x^{adv} = x + \eta$, subsequently causing the target model T to produce incorrect predictions, i.e., $T(x^{adv}) \neq y$.

Adversarial purification serves as a defensive mechanism against adversarial attacks, wherein the image is subjected to a ‘‘purification’’ process prior to being fed into the target model. An adversarial purification model P receives an adversarial sample x^{adv} and yields a purified sample $\hat{x} = P(x^{adv})$, ensuring that $T(x) = T(\hat{x})$.

Ideally, we hope to completely remove the adversarial perturbation η , making the purified sample \hat{x} exactly the same as the clean image x . Moreover, current research [23] shows that the remaining adversarial perturbations can still successfully attack, which inspired us to design adversarial purification methods to remove the adversarial perturbations and add defensive perturbations.

We formally discuss why simply reducing adversarial perturbations cannot achieve good defense results. We define a simple linear model $f(x) = \omega^T x$, a clean input x and an adversarial sample $x^{adv} = x + \eta$, where ω is a weight vector, ω and x follow the standard normal distribution $N(0, 1)$, η is a small adversarial perturbation subject to the constraint $\|\eta\|_\infty < \epsilon$, and ϵ is a small constant representing the perturbation threshold used to ensure the concealment of the adversarial sample. When the input of the linear model f is an adversarial sample x^{adv} , the

output of model f is

$$\begin{aligned} f(x^{adv}) &= \omega^T x^{adv} \\ &= \omega^T x + \omega^T \eta. \end{aligned} \quad (4)$$

The difference caused by the adversarial perturbation is

$$\begin{aligned} \Delta f &= f(x^{adv}) - f(x) \\ &= \omega^T \eta. \end{aligned} \quad (5)$$

We maximize the difference by assigning $\eta = \epsilon \text{sign}(\omega)$. At this time, the attack caused by the perturbation is the strongest, and the difference $\Delta f = \epsilon \|\omega\|_1$. If ω has n dimensions and the average magnitude of an element of the weight vector is m , then the expectation of the maximum difference $\mathbb{E}(\Delta f) = \epsilon mn$. n is usually very large in neural networks so even a small perturbation can cause a large difference in the prediction. Although many neural networks add nonlinearity through such activation functions ReLU, they still retain linearity as a whole, making it difficult for the network to resist adversarial perturbations.

Next, we discuss what happens after the adversarial purification method reduces adversarial perturbations. The previous discussion demonstrated that the attack capability of the adversarial sample varies linearly with the scale of the perturbation. Because the pixel value of the image is a discrete value of $[0, 255]$, if a perturbation threshold against perturbation is 16 unless we can completely remove the perturbation, even if the remaining perturbation is only 1, its attack capability remains 1/16. Since the powerful attack capability against perturbations mainly comes from the high dimensionality of the weight vector, limited linear reduction still retains considerable attack capability. Therefore, simply reducing confrontational disturbances cannot achieve the reliable defense effects.

Since global reduction has a limited effect against disturbance, are other reduction methods effective? For example, compared with reducing the overall perturbation to half, would it be better to completely remove the purification in half of the area to achieve a better defensive effect? We define the scale of an adversarial perturbation as the sum of the absolute values of the perturbation and design two methods to reduce the adversarial perturbation: global reduction $R_g(\eta, s)$ and local reduction $R_l(\eta, s)$, where $s \in (0, 1)$ represents the multiple of the scale. The global method is to multiply the adversarial perturbation by s :

$$R_g(\eta, s) = s\eta. \quad (6)$$

The local method is to retain s of the area of the adversarial perturbation and assign 0 to the rest:

$$R_l(\eta, s) = v \odot \eta, \quad (7)$$

where v is a mask vector with the same shape as η in which sn elements are 1 and the remaining elements are 0. \odot specifies the element-wise multiplication. Two methods reduce adversarial perturbations to the same scale and

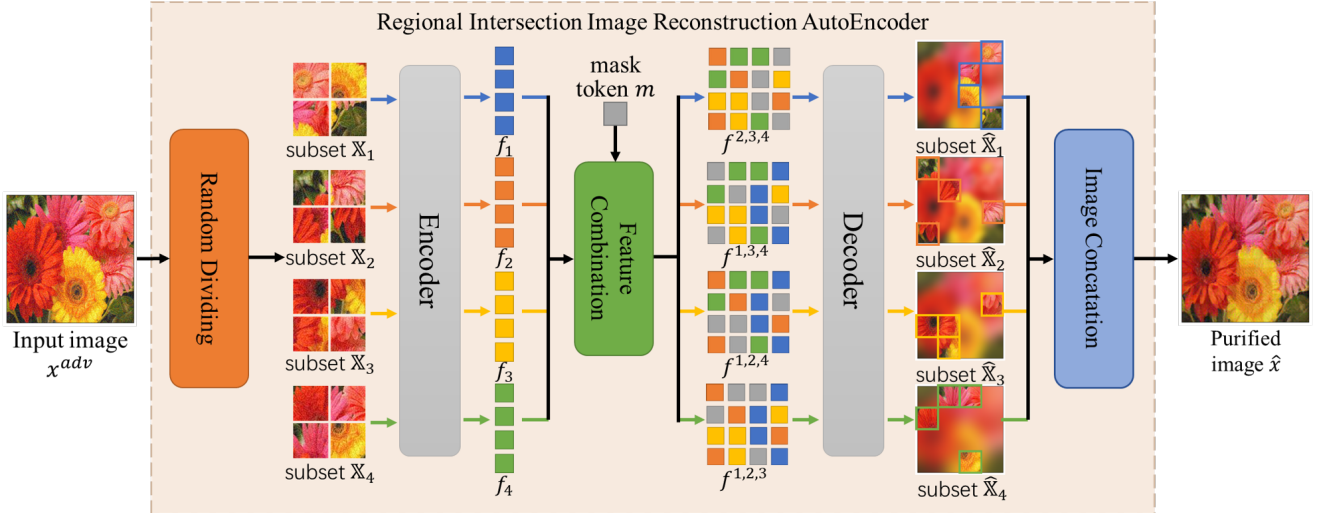


Figure 4: An overview of RIAE. The input image first is randomly divided into several subsets and input into the Encoder to output features. The features are combined and replaced with mask tokens to replace the parts that need to be masked, and then input into the Decoder respectively to output the purified subset. Finally, the patches in the subsets are combined to obtain the purified image.

then compare the attack capabilities of their processed adversarial samples. When the global method R_g and the local method R_l are used to reduce the adversarial perturbation η respectively, expectations of the maximum difference Δf_g and Δf_l are:

$$\mathbb{E}(\Delta f_g) = \mathbb{E}(\omega^T s \eta) = semn, \quad (8)$$

$$\mathbb{E}(\Delta f_l) = \mathbb{E}(\omega^T (v \odot \eta)) = semn. \quad (9)$$

It is observed that two different reduction methods yield equivalent outcomes, with the attack capability against the perturbation varying linearly with scale. This shows that for a generative model with a limited ability to eliminate perturbations, the ability to defend against adversarial attacks is closely related to the scale of perturbations that the model can eliminate. Combined with the previous results, we consider that the complete elimination of adversarial perturbations is necessary for adversarial purification methods.

3.2. Information Mask Purification

In this section, we propose IMPure based on the theory in Section 3.1. IMPure is architecturally composed of three principal components: RIAE, RCM, and a Feature Extraction Module. RIAE is responsible for transforming the input adversarial samples x^{adv} into purified samples \hat{x} without attack capabilities. The region intersection design we propose allows the network to reconstruct image patches without being interfered with by the same-position patches' adversarial perturbations. We propose RCM to encourage the generation of images with stronger defense capabilities. The adversarial samples x^{adv} and purified samples \hat{x} are spliced into a combined image x' which are fed into the feature extraction network to obtain features for calculating the perceptual loss. To improve the effectiveness of the reconstruction network, we use joint

constraints of pixel loss and feature loss to stabilize the underlying information and semantic information respectively. To this end, we extract high-level features of the purified image \hat{x} and the clean image x respectively through pre-trained feature extraction network for calculating the feature loss. The structure of IMPure is illustrated in Fig. 3.

3.2.1. Regional Intersection AutoEncoder

Given the discrete nature of image pixel values, limiting the overall perturbation scale of an image is difficult. Therefore we choose to reconstruct partial image patches at a time to simplify the task, by reconstructing in parallel on all patches to get the complete image. As shown in Fig. 5, to further reduce the perturbation scale, we mask the information of the reconstructed patches. One straightforward method to remove perturbations in patches is to erase pixels, leading to a dual loss: both the adversarial perturbations and the underlying image information. In the experimental part, we also evaluated two other information masking methods, image noise and feature noise. We refer to this image reconstruction technique as ‘‘Regional Intersection’’ and have developed a dedicated network RIAE for its implementation.

The architecture of the image reconstruction network is introduced below. In order to conveniently the operation of the image area, we build an autoencoder with reference to the Masked AutoEncoder (MAE) [16] structure. An overview of the autoencoder is shown in Fig. 4. First, like Vision Transformer (ViT) [11], an image $x \in \mathbb{R}^{H \times W \times C}$ is divided into patches $x_p \in \mathbb{R}^{N \times (P^2 \cdot C)}$ and mapped into embedded patches $z \in \mathbb{R}^{N \times D}$ through Patch embedding and Position embedding, where (H, W) is the resolution of the image, C is the number of the image channels, (P, P) is the resolution of each image patch, $N = HW/P^2$ is

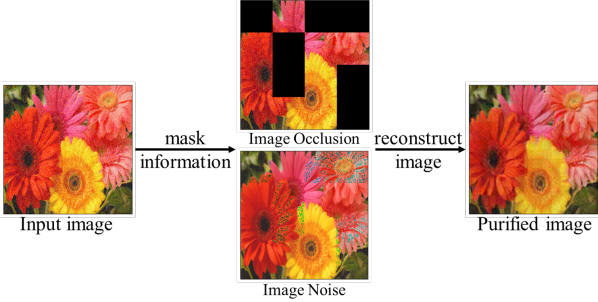


Figure 5: Information masking process: Get an input image, first use the information masking method on the specified area, and then reconstruct the masked image to obtain the purified image. Two information masking methods are illustrated respectively: Image Occlusion and Image Noise.

the number of the patches, D is the constant latent vector size for the transformer, Patch embedding is a trainable linear projection, and Position embedding is a fixed 2D sine-cosine encoding [6] to simplify training, unlike ViT.

Second, we treat the embedding patch as a set $\mathbb{X} = \{\mathbf{z}_i, |i \in \{1, \dots, N\}\}$ and divide it into several subsets $\{\mathbb{X}_i | i \in \{1, \dots, S\}\}$, where S is the number of subsets we specify. Randomly shuffle the order of patches and divide them into subsets evenly. The subsets are equivalent to uniform sampling of the image, and there are no repeated patches between subsets. Then just like standard ViT, we process these subsets with a series of transformer blocks to obtain the feature set $\mathbb{F} = \{\mathbf{f}_i \in \mathbb{R}^{L \times D}\}$, where $L = N/S$ is the number of patches in the subset. Each subset is encoded separately, there is no data interaction between subsets, and adversarial perturbations in different subsets will not affect each other.

Third, for the decoder, we occlude a subset and input the other subsets to the decoder to reconstruct the mask patches. We perform this step for each subset and finally combine all reconstructed parts together to obtain the complete reconstructed image. In the specific implementation, we combine all subsets into complete features according to their original positions and replace the features of the subset to be occluded with a mask token $\mathbf{m} \in \mathbb{R}^D$, then add position embedding and use a series of transformer blocks to complete the reconstruction of the subset set $\{\hat{\mathbb{X}}_i | i \in \{1, \dots, S\}\}$. Finally, we combine the patches in the subset set by location to obtain the purified image $\hat{\mathbf{x}}$.

3.2.2. Random Combination Module

Through the utilization of the aforementioned RIAE, we have procured purified images capable of counteracting adversarial attacks. While the autoencoder remains ostensibly impervious to the adversarial perturbations within the corresponding regions, it's imperative to acknowledge that in realistic imagery, the presence of regions with analogous features is not uncommon, and the adversarial perturbations within such regions often exhibit simi-

larity. This congruency in perturbations potentially renders the autoencoder susceptible to adversarial influence. To enhance defensive capabilities and mitigate this susceptibility, we introduce RCM. This module is meticulously designed to retain adversarial perturbations within random regions during the training phase intentionally. Such strategic retention is aimed at incentivizing the image reconstruction network to synthesize purified images endowed with fortified defensive attributes. This augmentation in the training paradigm serves to bolster the resilience of the autoencoder against adversarial perturbations, even those exhibiting subtle similarities within different regions, thereby refining the robustness of the overall model.

Our combination strategy is to randomly generate a mask matrix $\mathbf{u} \in \mathbb{R}^{H \times W}$ with the same resolution as the image, with each pixel having a grayscale value between 0.0 and 1.0, and obtain the combined image \mathbf{x}' through calculation:

$$\mathbf{x}' = \hat{\mathbf{x}} \odot \mathbf{u} + \mathbf{x} \odot (1 - \mathbf{u}). \quad (10)$$

In order to echo the transformer structure used by the image reconstruction network, the mask matrix also takes the same patch size to simulate the situation where the adversarial perturbations of some patches are not cleared.

3.2.3. Loss function

We utilize a combination of reconstruction loss and perceptual loss to train our image reconstruction network R . θ represent the parameters of our network. Our network receives an input image \mathbf{x}^{adv} and outputs a purified image $\hat{\mathbf{x}}$:

$$\hat{\mathbf{x}} = R(\mathbf{x}^{adv}, \theta). \quad (11)$$

The pixel-level loss as the reconstruction loss is defined as:

$$\mathcal{L}_{pix}(\hat{\mathbf{x}}, \mathbf{x}, \theta) = \frac{1}{HW} \|\hat{\mathbf{x}} - \mathbf{x}\|_1, \quad (12)$$

where \mathbf{x} is a clean image paired with \mathbf{x}^{adv} . We utilize L_1 loss which has been demonstrated to better than L_2 loss for image restoration [51].

Adversarial images and clean images are often very similar at the pixel level but have great differences in the feature representation of the target model, which leads to different predictions by the target model. In order to make the purified image consistent with the clean image in terms of feature representation, we use perceptual loss for constraints. Specifically, we record a convolutional network used for feature extraction as F , and $F_i(\mathbf{x})$ is the feature map output by the image \mathbf{x} at the i -th layer of the network F . The perceptual loss is defined as the difference in feature representation between the combined image \mathbf{x}' output by the image combination module and the clean image \mathbf{x} :

$$\begin{aligned} \mathcal{L}_{per-i}(\mathbf{x}', \mathbf{x}, \theta) &= D(F_i(\mathbf{x}), F_i(\mathbf{x}')) \\ &= \frac{1}{H_{F_i} W_{F_i}} \|F_i(\mathbf{x}) - F_i(\mathbf{x}')\|_1, \end{aligned} \quad (13)$$

where D is the distance function used to measure feature differences. In our method, the L_1 norm is used as the distance function, and (H_{F_i}, W_{F_i}) is the resolution of the feature map.

The key to the success of perceptual loss is the structure of the feature extraction network [24]. The model’s feature representation usually has certain characteristics, hence it may be effective to choose the corresponding model as the feature extraction network. However, single-model selection can diminish the generalization of the defense method, while the overhead of multi-model integrated training remains substantial. Consequently, we elect to employ the 19-layer VGG network [41] for feature extraction. First, the VGG network has proven its superiority in tasks using perceptual loss in other fields. Second, the convolutional structure of the VGG network is widely used by many CNNs. Finally, the network’s depth plays a pivotal role in image feature extraction. To capture the correlation between multi-layer statistics extracted by multi-layer CNN, we integrate multi-layer feature maps in the perceptual loss. The perceptual loss is defined as:

$$\mathcal{L}_{per}(\mathbf{x}', \mathbf{x}, \theta) = \sum_K \lambda_i \mathcal{L}_{per-i}(\mathbf{x}', \mathbf{x}, \theta), \quad (14)$$

where K represents the set of layers of VGG feature maps involved in the operation. λ_i represents the weight of the loss. VGG19 contains 5 similar modules, each module contains convolution layers, pooling layers, and nonlinear activation layers. We select the feature map output after the ReLU activation layer in each module to participate in the calculation of perceptual loss.

The overall loss we use to train our network is defined as:

$$\mathcal{L}_{overall}(\mathbf{x}^{adv}, \mathbf{x}, \theta) = \mathcal{L}_{pix}(\hat{\mathbf{x}}, \mathbf{x}, \theta) + \sum_K \lambda_i \mathcal{L}_{per-i}(\mathbf{x}', \mathbf{x}, \theta), \quad (15)$$

where we obtain the combination of losses through experiments select the feature maps of the 4-th and 5-th modules and the previous layer output of the softmax layer. We set $\lambda_4 = 30$, $\lambda_5 = 10$ and $\lambda_{pre_soft} = 5$.

4. Experiments

In this section, we illustrate the experiments and use large-scale results to prove the superiority of our method.

4.1. Experimental Settings

Target models: The target model architectures we choose are Inception V3 (IncV3) [43], Inception-ResNet-v2 (IncRes) [42] and ResNet-v2-101 (ResNet) [17], and the task is image classification. All models are pre-trained on ImageNet dataset, training settings and model weights can be found on this page¹. We do not make any modifications

to the structure and weights of the model, and the testing of the model will be conducted according to the original settings.

Attack methods: We select a series of adversarial attack methods to generate adversarial examples, including FGSM [13], C&W (CW) [3], Deepfool (Df) [31], PGD [29], MIFGSM (MI) [10], DIFGSM (DI) [48], MD-FGSM (MD) [48], APGD-ce (A-ce) [7] and APGD-dlr (A-dlr) [7]. The code implementation of the first five methods refers to the CleverHans library² [34], and the remaining four come from the Torchattacks library³ [21]. For methods that support targeted attacks and non-targeted attacks, we uniformly choose non-targeted attacks. To generate strong adversarial samples, the perturbation threshold ϵ is set to 16/255, and other parameters of the attack method will use the default settings of the source code author.

Dataset: We train and evaluate our defense method on the ImageNet [8] dataset with top-1 classification accuracy, which is widely used in adversarial attack [10, 48, 7] and defense [47, 37, 20] domains. For each target model, we only select images that the model can correctly classify, and modify the resolution of the image according to the input requirements of the model. We randomly select 10 images from each category in the Imagenet training set and then use all adversarial attack methods to generate adversarial samples as the training set for our defense method. We randomly select 5 images from each category in the Imagenet verification set to generate adversarial samples as the verification set. Use the same method to select different images to generate adversarial samples as the test set.

Implementation details: The structure and parameter settings of the transformer part of Regional Intersection AutoEncoder refer to MFFAE [26]. The patch size is set to 16×16 . For the transformer block of the encoder part, the embedding dimension is set to 768, the depth is set to 12, and the number of multi-nodes inside is set to 12. For the transformer block of the decoder part, the embedding dimension is set to 512, the depth is set to 8, and the number of multi-nodes inside is set to 16. Feature extraction network selection VGG19. We train our model using AdamW [28] with $\beta_1 = 0.9$, $\beta_2 = 0.95$. The learning rate is initially set to $1e^{-4}$ and use CosineAnnealing [27] with 100 epochs. The training procedure is described in the Algorithm. 1.

4.2. Black-box evaluation

We assess the efficacy of our proposed method in a black-box setting. The black-box scenario where it is presumed that the attacker has access to the structure and parameters of the target network, yet remains unaware of the defense mechanism. Table 1 delineates the top-1

¹<https://pytorch.org/vision/stable/models.html#classification>

²<https://github.com/cleverhans-lab/cleverhans>

³<https://github.com/Harry24k/adversarial-attacks-pytorch>

Table 1: Top-1 classification accuracy of the three target models on the test set without defense and with defense in a black-box setting.

Methods		Clean	FGSM [13]	CW [3]	Df [31]	PGD [29]	MI [10]	DI [48]	MD [48]	A-ce [7]	A-dlr [7]
ResNet [17]	No defense	100.00	53.70	1.20	0.12	0.74	0.42	2.24	2.26	0.56	0.36
	Defense	87.54	74.42	83.96	83.74	75.70	66.04	68.44	57.52	74.38	74.76
IncV3 [43]	No defense	100.00	26.40	0.22	0.00	0.04	0.06	0.16	0.08	0.08	0.84
	Defense	87.30	66.02	84.28	78.18	75.00	64.18	70.04	58.52	74.94	73.78
IncRes [42]	No defense	100.00	48.32	2.24	0.08	1.34	1.26	2.84	3.04	1.52	1.78
	Defense	85.86	70.88	82.40	80.34	76.44	66.52	71.26	60.74	75.20	74.38

Table 2: Top-1 classification accuracy of IncV3 against different defense methods on the test set in a black-box setting.

Method	Clean	FGSM [13]	CW [3]	Df [31]	PGD [29]	MI [10]	DI [48]	MD [48]	A-ce [7]	A-dlr [7]
Attack	100.00	26.40	0.22	0.00	0.04	0.06	0.16	0.08	0.08	0.84
Random [47]	95.56	36.18	37.86	25.30	3.70	1.34	0.34	0.42	4.52	15.30
PD+WD [37]	80.34	36.78	58.22	37.64	15.26	3.00	7.54	1.54	12.48	24.56
WD+SR [32]	84.84	48.34	78.50	69.72	42.40	29.40	32.84	20.44	40.30	49.70
ComDefend [20]	89.44	36.92	83.40	54.98	38.78	10.10	17.30	2.84	45.74	51.30
Recon [50]	94.02	64.58	52.48	42.26	60.84	55.36	43.82	43.24	65.16	67.66
DIR [52]	64.02	53.56	63.02	61.58	58.92	52.32	55.88	48.96	57.40	57.16
Ours	87.30	66.02	84.28	78.18	75.00	64.18	70.04	58.52	74.94	73.78

Algorithm 1 Regional Intersection Adversarial Purification Method Training Algorithm.

Input: Adversarial image \mathbf{x}^{adv} , clean image \mathbf{x} , regional intersection autoencoder R parameterized by θ , random combine module U , feature extraction network F with map layers K , learning rate l_r and number of training epochs of T .

- 1: Initialize θ with random values;
- 2: **for** $t \leftarrow 0$ to T **do**
- 3: $\hat{\mathbf{x}} \leftarrow R(\mathbf{x}^{adv}, \theta)$;
- 4: $\mathbf{x}' \leftarrow U(\hat{\mathbf{x}}, \mathbf{x}^{adv})$;
- 5: Calculate \mathcal{L}_{pix} , \mathcal{L}_{per} and $\mathcal{L}_{overall}$ using Eq. 12, Eq. 14 and Eq. 15, respectively;
- 6: $\theta \leftarrow \theta - l_r \nabla_{\theta} (\mathcal{L}_{overall}(\mathbf{x}^{adv}, \mathbf{x}, \theta))$;
- 7: **end for**.

classification accuracy of three target models on the test dataset, both in the absence and presence of defense. Each attack methodology is specified under the “Methods” row, with “Clean” denoting the original unaltered sample. The “No defense” row reflects the classification accuracy of the target model when subjected to adversarial samples; a lower value here signifies a more potent attack methodology. The “Defense” column represents the classification accuracy of the purified images, with higher values indicating enhanced adversarial robustness.

The adversarial robustness of the target model can be ascertained through metrics presented in the “No defense” row. IncRes exhibits the best adversarial robustness in the face of adversarial attacks, while IncV3 performs the worst. Nevertheless, the substantial dip in classification accuracy witnessed across all three target models underscores the imperativeness of defensive strategies.

Post-defense application, a marginal decline in the “Clean” metrics across all three target models is observed. This emanates from the information masking our method enacts during image reconstruction, inevitably leading to some loss of image information. In juxtaposition with the “No defense” scenario, a marked improvement is discernible across various attack methodologies once the defense is activated. Overall, IncRes continues to exhibit the most robust adversarial resilience.

Different attack samples perform differently under defense strategies. Attacks via C&W and DeepFool are relatively easier to thwart as these techniques prioritize stealth and employ L2 norm constraints to mitigate perturbations, thereby moderating their attack potency to just the threshold of effectiveness, rendering them more easily defensible. In stark contrast, warding off attacks from MIFGSM, DIFGSM, and MDFGSM proves more challenging due to their emphasis on attack migration. APGD, an enhanced variant of PGD, presents an even tougher defense challenge. It’s noteworthy to mention that given Inception’s dismal accuracy against FGSM attack samples, the defense metrics too, are correspondingly low.

4.3. Comparison with previous state-of-the-art methods

In order to prove the effectiveness of our proposed defense model, we selected four model-agnostic methods for comparison, including random resizing and padding (Random) ⁴[47], pixel deflection(PD+WD)⁵ [37], wavelet denoising & super-resolution(WD+SR)⁶ [32], ComDe-

⁴https://github.com/cihangxie/NIPS2017_adv_challenge_defense

⁵<https://github.com/iamaaditya/pixel-deflection>

⁶<https://github.com/aamir-mustafa/super-resolution-adversarial-defense>

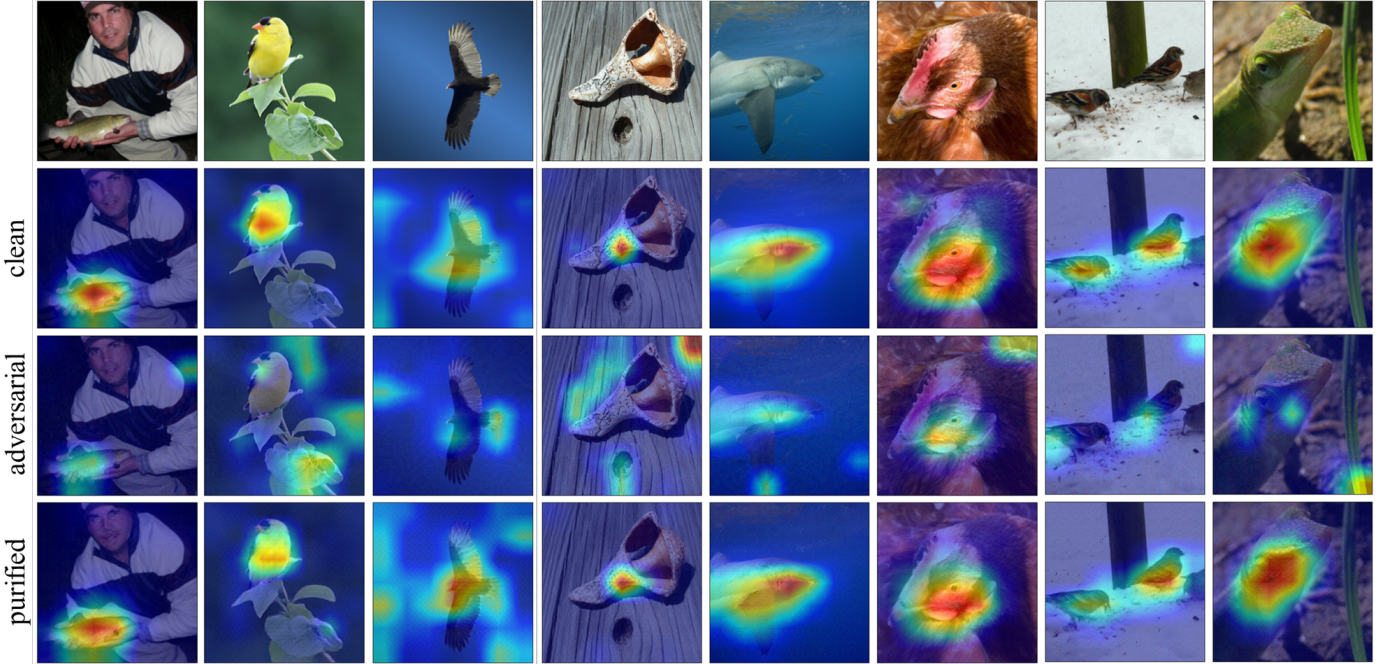


Figure 6: Class activation maps of GradCAM on resnet101 for clean images, adversarial examples and sanitized images.

fund⁷ [20], Recon⁸ [50] and DIR⁹ [52]. Among them, Random and PD are static defense strategies, and WD+SR, ComDefend, Recon and DIR include deep models. We used the weights provided in the WD+SR project. ComDefend, Recon and DIR did not provide appropriate weights, so we trained the weights ourselves to participate in the test set in strict accordance with the project requirements. All testing procedures are carried out strictly in accordance with the requirements in their projects. We compared the classification accuracy of various adversarial examples of IncV3 on the test set. The higher the classification accuracy, the more effective the defense against this type of attack is.

The “Clean” metric illuminates the capability of the defense technique to retain the fundamental information within the image. Random emerges as the top performer, attributing to the fact that resizing and padding operations preserve most of the image information, and are frequently employed in diverse data augmentation techniques. On the other hand, Recon has obvious advantages in maintaining image information based on the Unet-like model, rendering its performance slightly inferior to Random. ComDefend executes compression operations, consequently sacrificing some level of detailed information. Similarly, our approach induces some information loss due to information masking, aligning its performance closely with that of ComDefend. Given that CNN lacks scale invariance, the target model is difficult to classify the super-resolution images generated by WD+SR. The

PD method, which involves pixel swapping, disrupts local information, exerting a significant impact on the classification task. The image erasure used by DIR loses a lot of information when facing high-resolution images, which greatly affects the quality of the purified image and limits the defense performance.

The ensuing metrics unveil the defense method’s ability to resist specific adversarial attacks. Our defense method achieves the best defensive performance against all attacks, achieving conspicuous advantages, especially in the three strong adversarial attack methods of MIFGSM, DIFGSM, and MDFGSM. Static methods, Random and PD, exhibit a semblance of defensive potential against adversarial samples characterized by weaker attack capabilities, however, their efficacy dwindles against robust attack modalities. Both WD+SR and ComDefend provide defense against all attacks, albeit their performance is circumscribed by their exclusive reliance on pixel loss.

4.4. Effect of the adversarial perturbation on the feature map

Adversarial perturbations are usually invisible to the human eye, but when added to a clean image, they cause very noticeable changes in the intermediate feature maps of the target model. This perturbation usually amplifies as the network propagates layer by layer, eventually causing the target model to output incorrect predictions. By visualizing and comparing the feature maps of clean images, adversarial samples, and purified images in the target model, the defense effect of the adversarial purified model can be intuitively demonstrated. As shown in Fig. 6, we used GradCAM [40] to visualize the class activation map

⁷<https://github.com/jiaxiaojunQAQ/Comdefend>

⁸<https://github.com/ZOMIN28/Reconstructing-Images>

⁹<https://github.com/dwDavidxd/DIR>

of the image on resnet101. Adversarial examples are generated by the PGD method. We can clearly see that the areas of focus for the clean image target model are closely related to the targets in the image. For adversarial samples, there is no obvious correlation between the area of interest and the target. The clean images generated by our adversarial cleansing method are basically the same as the clean images on the class activation map.

5. Ablation experiment

To further optimize the network structure, we try various combinations and optimization strategies to construct our methods and compare their performance. In order to simplify the problem, we use ResNet as the target model, the training epoch is set to 30, and other experimental settings remain consistent with the previous section.

5.1. Effects of difference information mask methods

In our proposed framework, the information masking technique is pivotal. An optimal masking method ought to retain the innate information of the image to the maximum extent while neutralizing a majority of the adversarial perturbations. We have instituted three distinct information masking methods: Image Occlusion (ImgOcc), Image Noise (ImgNoise), and Feature Noise (FeaNoise). ImgOcc epitomizes the total elimination of image information, standing as the most radical among the trio, as it retains neither the adversarial perturbation nor the original information. ImgNoise involves the injection of Gaussian noise with random intensity into the image. Considering that the pixel value that resists disturbance is small, a noise of suitable intensity can resist most perturbations, preserving the original information to a certain degree, thereby facilitating image reconstruction. FeaNoise operates on the features yielded by the RIAE encoder, introducing Gaussian noise of random intensity to these features. As depicted in Fig. 7, we conduct an evaluation to gauge the efficacy of the three information masking methods alongside the amalgam of ImgNoise and FeaNoise (Img+Fea). To render a lucid comparative insight via a singular graph, we normalized the results.

The Clean indicator succinctly illustrates the image reconstruction capability of the model and echoes the information loss induced by the information masking method. Our observations reveal the most subpar performance with ImgOcc, given its failure to transmit any information. ImgNoise exhibits a marginally better performance as the Gaussian noise retains some low-frequency information. FeaNoise outshone owing to its model’s proficiency in learning anti-noise features, thus being minimally impacted. Reflecting upon the defense performance, ImgOcc is markedly inferior across most attack methods. However, it matched FeaNoise on DIFGSM, outperformed FeaNoise on MIFGSM, and outstripped both ImgNoise and FeaNoise on MDFGSM. This underscores the superior perturbation resistance of ImgOcc over ImgNoise and

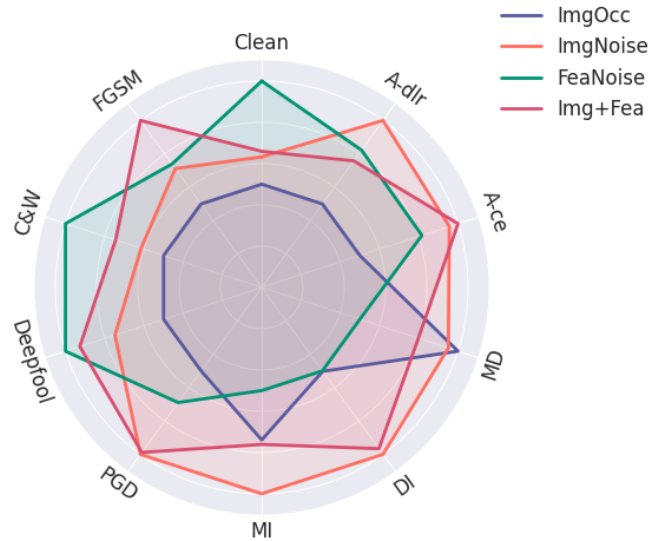


Figure 7: Performance of defense models against various adversarial attack methods under different information mask methods.

FeaNoise. In the battles of FGSM, CW, and DeepFool, FeaNoise trumps ImgNoise, although it falls significantly behind ImgNoise in other attack methodologies, revealing a better disturbance resistance of ImgNoise over FeaNoise. Finally, the composite method of Img+Fea manifests the pinnacle of overall performance, demonstrating that both masking methods are individually resistant to unique perturbations. ImgNoise+FeaNoise is thereby employed in our defense method.

5.2. Effects of difference perceptual losses

We evaluate the resulting performance of various loss combinations on lightweight experimental designs to obtain optimal losses. The output of the RELU layer in each convolution module of VGG19 will participate in the calculation of perceptual loss. We mark the layer in which the i -th module participates in the operation as layer i . As shown in Fig. 8, we set the initial perceptual loss to the previous layer output of the softmax layer of VGG19, use the pixel loss plus the perceptual loss as the benchmark, and then add one layer of feature maps to the perceptual loss each time and observe the experimental results.

On the three indicators of Clean, FGSM, and APGD-dlr, the performance gets better as the number of feature maps increases. The shallower the feature map, the lower the level of information it represents. This low-level information is obviously helpful for the reconstruction of clean images. At the same time, the experimental results also show that these shallow features are also of positive significance for weak adversarial attacks. In contrast, on MIFGSM, DIFGSM, and MDFGSM, the performance decreases as the number of feature maps increases. These three adversarial attacks are the three with the worst defense effects among the various defense methods currently available. For such adversarial samples with robust attack

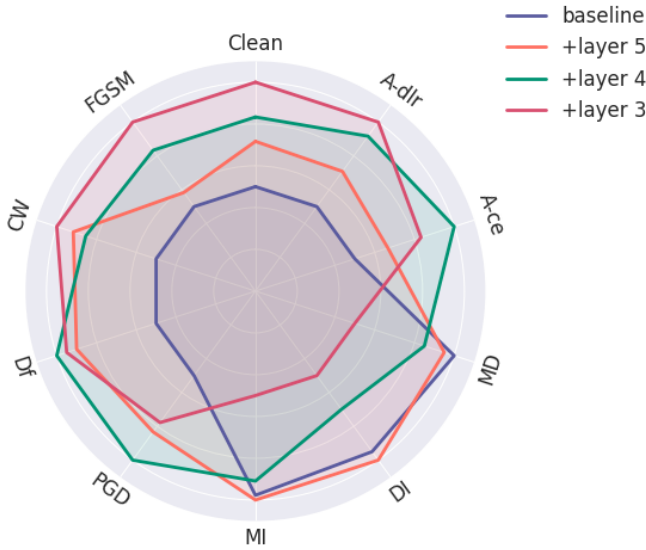


Figure 8: Performance of defense models against various adversarial attack methods under different perceptual loss settings.

capabilities, shallow features can no longer have a positive effect. Based on various indicators, our defense method finally chose to add layer5+layer4 to the perceptual loss.

5.3. Effects of difference dividing methods

The image is partially occluded and then reconstructed after the input area Regional Intersection image reconstruction AutoEncoder. In the training phase, in order to prevent overfitting, we randomly select some areas for occlusion. During the testing phase, we experimentally determined which method of selecting occlusion areas would achieve the best performance. The two dividing methods are: a uniform method is used to select the area, and the same random selection as in the training stage. We show in Fig. 9 the performance of the defense model in resisting various adversarial attack methods using two methods. In order to visually demonstrate the performance gap between different methods, we normalized the results. The uniform method has obvious advantages over the random method in most indicators. The random method’s defense performance against MIFGSM, DIFGSM, and MDFGSM is slightly better than the uniform method. This may be due to the fact that the adversarial perturbations generated by these robust adversarial attack methods can also take effect in adjacent areas, making it difficult for uniform methods to resist. We consider the gap in comprehensive performance between the two methods and choose to use the uniform dividing method during testing.

5.4. Effects of defense against agnostic attacks

Agnostic attacks denote adversarial assault mechanisms that remain external to the model training regimen. Given the incessant evolution of adversarial attack methodologies, fortifying models against agnostic attacks

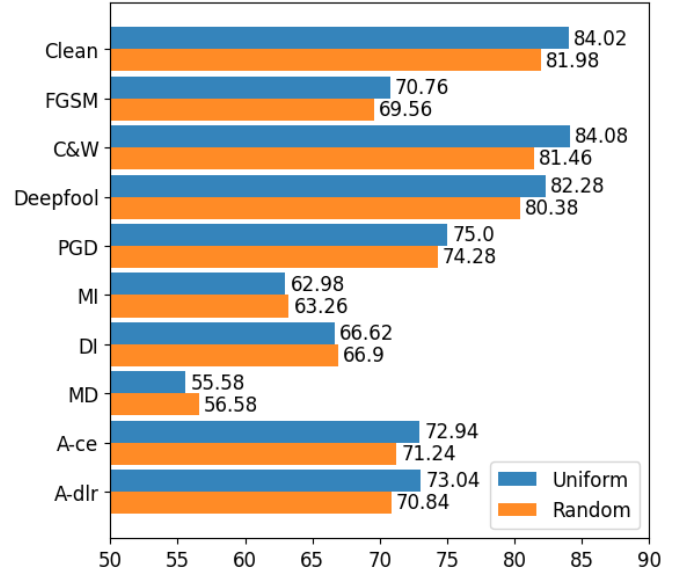


Figure 9: Top-1 classification accuracy of defense models against various adversarial attack methods for two dividing methods.

has emerged as a crucial endeavor. We conduct experiments to appraise the robustness of our defense strategy against gradient-based agnostic onslaughts. Early gradient-based attack variants such as FGSM, PGD, and MIFGSM are integrated into the training dataset, while DIFGSM, DMFGSM, and APGD are employed as agnostic attack exemplars to evaluate the defense efficacy of the model. The outcomes, delineated in Table 3, affirm that our defense paradigm continues to exhibit commendable resilience when confronted with agnostic assaults.

Table 3: Top-1 classification accuracy of ResNet against against agnostic attacks on the test set in a black-box setting.

Method	No Defense	Defense
Clean	100.00	86.12
FGSM [13]	53.70	86.12
PGD [29]	0.74	73.78
MI [10]	0.42	62.04
DI [48]	2.24	62.28
MD [48]	2.26	47.46
A-ce [7]	0.56	72.64
A-dlr [7]	0.36	72.42

6. Conclusion

In this paper, we demonstrate the hazards of residual adversarial perturbations and advocate for adversarial purification methods that eliminates adversarial perturbations wherever possible. Based on this theory, we propose a novel adversarial purification method MIPure that defend against adversarial attacks by maximizing the elimination of same-position perturbations and resisting

content-similar perturbations. We construct RIAE to constrain the scale of the same-position perturbations on the purified image by masking the image patch information to destroy the perturbations and reconstruct the patch. Then we also propose RCM to encourage our model to resist the influence of content-similar perturbations by simulating residual perturbations. Finally, we propose a joint constraint on pixel loss and perceptual loss to make the generation of purified images more flexible. Experiments show that our adversarial purification model is very effective in defending against adversarial attacks and also exhibits good robustness in the face of agnostic attacks. In future work, we will continue to optimize the structure of the autoencoder and solve problems such as training instability and adversarial sample augmentation.

Acknowledgments

This work was supported in part by the National Natural Science Foundation of China under Grant 61871226, Grant 61571230, Grant 61802190, and Grant 61906093; in part by the Jiangsu Provincial Social Developing Project under Grant BE2018727; and in part by the Open Research Fund in 2021 of Jiangsu Key Laboratory of Spectral Imaging and Intelligent Sense under Grant JSGP202101 and Grant JSGP202204.

References

- [1] Athalye, A., Carlini, N., Wagner, D.A., 2018. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples, in: Dy, J.G., Krause, A. (Eds.), Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, Stockholm, Sweden, July 10-15, 2018, PMLR. pp. 274–283. URL: <http://proceedings.mlr.press/v80/athalye18a.html>.
- [2] Brendel, W., Rauber, J., Bethge, M., 2018. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models, in: 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings, OpenReview.net. URL: <https://openreview.net/forum?id=SyZIOGWZ>.
- [3] Carlini, N., Wagner, D.A., 2016. Towards evaluating the robustness of neural networks. 2017 IEEE Symposium on Security and Privacy (SP), 39–57. URL: <https://api.semanticscholar.org/CorpusID:2893830>.
- [4] Chen, L., Chu, X., Zhang, X., Sun, J., 2022. Simple baselines for image restoration, in: European Conference on Computer Vision. URL: <https://api.semanticscholar.org/CorpusID:248085491>.
- [5] Chen, P.Y., Zhang, H., Sharma, Y., Yi, J., Hsieh, C.J., 2017. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security URL: <https://api.semanticscholar.org/CorpusID:2179389>.
- [6] Chen, X., Xie, S., He, K., 2021. An empirical study of training self-supervised vision transformers, in: 2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021, IEEE. pp. 9620–9629. URL: <https://doi.org/10.1109/ICCV48922.2021.00950>, doi:10.1109/ICCV48922.2021.00950.
- [7] Croce, F., Hein, M., 2020. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks, in: Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event, PMLR. pp. 2206–2216. URL: <http://proceedings.mlr.press/v119/croce20b.html>.
- [8] Deng, J., Dong, W., Socher, R., Li, L., Li, K., Li, F., 2009. Imagenet: A large-scale hierarchical image database, in: 2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2009), 20-25 June 2009, Miami, Florida, USA, IEEE Computer Society. pp. 248–255. URL: <https://doi.org/10.1109/CVPR.2009.5206848>, doi:10.1109/CVPR.2009.5206848.
- [9] Dhillon, G.S., Azizzadenesheli, K., Lipton, Z.C., Bernstein, J., Kossaiji, J., Khanna, A., Anandkumar, A., 2018. Stochastic activation pruning for robust adversarial defense, in: 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings, OpenReview.net. URL: <https://openreview.net/forum?id=H1uR4GZRZ>.
- [10] 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://openaccess.thecvf.com/content_cvpr_2018/html/Dong_Boosting_Adversarial_Attacks_CVPR_2018_paper.html, doi:10.1109/CVPR.2018.00957.
- [11] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N., 2021. An image is worth 16x16 words: Transformers for image recognition at scale, in: 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021, OpenReview.net. URL: <https://openreview.net/forum?id=YicbFdNTTy>.
- [12] Dziugaite, G.K., Ghahramani, Z., Roy, D.M., 2016. A study of the effect of jpg compression on adversarial images. ArXiv preprint abs/1608.00853. URL: <https://arxiv.org/abs/1608.00853>.
- [13] Goodfellow, I.J., Shlens, J., Szegedy, C., 2015. Explaining and harnessing adversarial examples, in: Bengio, Y., LeCun, Y. (Eds.), 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings. URL: <http://arxiv.org/abs/1412.6572>.
- [14] Gu, S.S., Rigazio, L., 2014. Towards deep neural network architectures robust to adversarial examples. CoRR abs/1412.5068. URL: <https://api.semanticscholar.org/CorpusID:15538683>.
- [15] Guo, C., Rana, M., Cissé, M., van der Maaten, L., 2018. Countering adversarial images using input transformations, in: 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings, OpenReview.net. URL: <https://openreview.net/forum?id=SyJ7C1WCb>.
- [16] He, K., Chen, X., Xie, S., Li, Y., Doll'ar, P., Girshick, R.B., 2021. Masked autoencoders are scalable vision learners. 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 15979–15988. URL: <https://api.semanticscholar.org/CorpusID:243985980>.
- [17] He, K., Zhang, X., Ren, S., Sun, J., 2016. Identity mappings in deep residual networks, in: European Conference on Computer Vision. URL: <https://api.semanticscholar.org/CorpusID:6447277>.
- [18] Huang, Z., Zhang, J., Shan, H., 2021. When age-invariant face recognition meets face age synthesis: A multi-task learning framework. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 7278–7287. URL: <https://api.semanticscholar.org/CorpusID:232092666>.
- [19] Ilyas, A., Santurkar, S., Tsipras, D., Engstrom, L., Tran, B., Madry, A., 2019. Adversarial examples are not bugs, they

- are features, in: Wallach, H.M., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E.B., Garnett, R. (Eds.), Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pp. 125–136. URL: <https://proceedings.neurips.cc/paper/2019/hash/e2c420d928d4bf8ce0ff2ec19b371514-Abstract.html>.
- [20] Jia, X., Wei, X., Cao, X., Foroosh, H., 2019. Comdefend: An efficient image compression model to defend adversarial examples, in: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019, Computer Vision Foundation / IEEE. pp. 6084–6092. URL: http://openaccess.thecvf.com/content_CVPR_2019/html/Jia_ComDefend_An_Efficient_Image_Compression_Model_to_Defend_Adversarial_Examples_CVPR_2019_paper.html, doi:10.1109/CVPR.2019.00624.
- [21] Kim, H., 2020. Torchattacks: A pytorch repository for adversarial attacks. ArXiv preprint abs/2010.01950. URL: <https://arxiv.org/abs/2010.01950>.
- [22] Kurakin, A., Goodfellow, I.J., Bengio, S., 2016. Adversarial examples in the physical world. ArXiv preprint abs/1607.02533. URL: <https://arxiv.org/abs/1607.02533>.
- [23] Liao, F., Liang, M., Dong, Y., Pang, T., Hu, X., Zhu, J., 2018. Defense against adversarial attacks using high-level representation guided denoiser, 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. 1778–1787. URL: http://openaccess.thecvf.com/content_cvpr_2018/html/Liao_Defense_Against_Adversarial_CVPR_2018_paper.html, doi:10.1109/CVPR.2018.00191.
- [24] Liu, Y., Chen, H., Chen, Y., Yin, W., Shen, C., 2021. Generic perceptual loss for modeling structured output dependencies. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 5420–5428 URL: <https://api.semanticscholar.org/CorpusID:232290506>.
- [25] Liu, Y., Chen, X., Liu, C., Song, D., 2017. Delving into transferable adversarial examples and black-box attacks, in: 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings, OpenReview.net. URL: <https://openreview.net/forum?id=Sys6GJqxl>.
- [26] Liu, Y., Zhang, S., Chen, J., Yu, Z., Chen, K., Lin, D., 2023. Improving pixel-based mim by reducing wasted modeling capability. ArXiv preprint abs/2308.00261. URL: <https://arxiv.org/abs/2308.00261>.
- [27] Loshchilov, I., Hutter, F., 2017. SGDR: stochastic gradient descent with warm restarts, in: 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings, OpenReview.net. URL: <https://openreview.net/forum?id=Skq89Scxx>.
- [28] Loshchilov, I., Hutter, F., 2019. Decoupled weight decay regularization, in: 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019, OpenReview.net. URL: <https://openreview.net/forum?id=Bkg6RiCqY7>.
- [29] 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, Conference Track Proceedings, OpenReview.net. URL: <https://openreview.net/forum?id=rJzIBfZAb>.
- [30] Meng, D., Chen, H., 2017. Magnet: A two-pronged defense against adversarial examples. Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security URL: <https://api.semanticscholar.org/CorpusID:3583538>.
- [31] Moosavi-Dezfooli, S., 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>, doi:10.1109/CVPR.2016.282.
- [32] Mustafa, A., Khan, S.H., Hayat, M., Shen, J., Shao, L., 2019. Image super-resolution as a defense against adversarial attacks. IEEE Transactions on Image Processing 29, 1711–1724. URL: <https://api.semanticscholar.org/CorpusID:57573757>.
- [33] Nie, W., Guo, B., Huang, Y., Xiao, C., Vahdat, A., Anandkumar, A., 2022. Diffusion models for adversarial purification, in: Chaudhuri, K., Jegelka, S., Song, L., Szepesvári, C., Niu, G., Sabato, S. (Eds.), International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA, PMLR. pp. 16805–16827. URL: <https://proceedings.mlr.press/v162/nie22a.html>.
- [34] Papernot, N., Faghri, F., Carlini, N., Goodfellow, I., Feinman, R., Kurakin, A., Xie, C., Sharma, Y., Brown, T., Roy, A., Matyasko, A., Behzadan, V., Hambardzumyan, K., Zhang, Z., Juang, Y.L., Li, Z., Sheatsley, R., Garg, A., Uesato, J., Gierke, W., Dong, Y., Berthelot, D., Hendricks, P., Rauber, J., Long, R., 2016a. Technical report on the cleverhans v2.1.0 adversarial examples library. ArXiv preprint abs/1610.00768. URL: <https://arxiv.org/abs/1610.00768>.
- [35] Papernot, N., Mcdaniel, P., Goodfellow, I.J., Jha, S., Celik, Z.B., Swami, A., 2016b. Practical black-box attacks against machine learning. Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security URL: <https://api.semanticscholar.org/CorpusID:1090603>.
- [36] Prakash, A., Chitta, K., Geiger, A., 2021. Multi-modal fusion transformer for end-to-end autonomous driving. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 7073–7083 URL: <https://api.semanticscholar.org/CorpusID:233148602>.
- [37] Prakash, A., Moran, N., Garber, S., DiLillo, A., Storer, J.A., 2018. Deflecting adversarial attacks with pixel deflection. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 8571–8580 URL: <https://api.semanticscholar.org/CorpusID:4528012>.
- [38] Rice, L., Wong, E., Kolter, J.Z., 2020. Overfitting in adversarially robust deep learning, in: Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event, PMLR. pp. 8093–8104. URL: <http://proceedings.mlr.press/v119/rice20a.html>.
- [39] Schmidt, L., Santurkar, S., Tsipras, D., Talwar, K., Madry, A., 2018. Adversarially robust generalization requires more data, in: Bengio, S., Wallach, H.M., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (Eds.), Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, pp. 5019–5031. URL: <https://proceedings.neurips.cc/paper/2018/hash/f708f064faaf32a43e4d3c784e6af9ea-Abstract.html>.
- [40] Selvaraju, R.R., Das, A., Vedantam, R., Cogswell, M., Parikh, D., Batra, D., 2016. Grad-cam: Visual explanations from deep networks via gradient-based localization. International Journal of Computer Vision 128, 336–359. URL: <https://api.semanticscholar.org/CorpusID:15019293>.
- [41] Simonyan, K., Zisserman, A., 2015. Very deep convolutional networks for large-scale image recognition, in: Bengio, Y., LeCun, Y. (Eds.), 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings. URL: <http://arxiv.org/abs/1409.1556>.
- [42] Szegedy, C., Ioffe, S., Vanhoucke, V., Alemi, A.A., 2017. Inception-v4, inception-resnet and the impact of residual connections on learning, in: Singh, S.P., Markovitch, S. (Eds.), Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA, AAAI Press. pp. 4278–4284. URL: <http://aaai.org/ocs/index.php/AAAI/AAAI17/paper/view/14806>.
- [43] 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: <https://doi.org/10.1109/CVPR.2016.308>, doi:10.1109/CVPR.2016.308.
- [44] Thys, S., Ranst, W.V., Goedemé, T., 2019. Fooling automated surveillance cameras: Adversarial patches to attack person detection. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW) , 49–55 URL: <https://api.semanticscholar.org/CorpusID:121124946>.
- [45] Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., Madry, A., 2019. Robustness may be at odds with accuracy, in: 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019, OpenReview.net. URL: <https://openreview.net/forum?id=SyxAb30cY7>.
- [46] Wong, E., Rice, L., Kolter, J.Z., 2020. Fast is better than free: Revisiting adversarial training, in: 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020, OpenReview.net. URL: <https://openreview.net/forum?id=BJx040EFvH>.
- [47] Xie, C., Wang, J., Zhang, Z., Ren, Z., Yuille, A.L., 2018. Mitigating adversarial effects through randomization, in: 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings, OpenReview.net. URL: <https://openreview.net/forum?id=Sk9yuq10Z>.
- [48] Xie, C., Zhang, Z., Zhou, Y., Bai, S., Wang, J., Ren, Z., Yuille, A.L., 2019. Improving transferability of adversarial examples with input diversity, in: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019, Computer Vision Foundation / IEEE. pp. 2730–2739. URL: http://openaccess.thecvf.com/content_CVPR_2019/html/Xie_Improving_Transferability_of_Adversarial_Examples_With_Input_Diversity_CVPR_2019_paper.html, doi:10.1109/CVPR.2019.00284.
- [49] Zhang, K., Li, Y., Liang, J., Cao, J., Zhang, Y., Tang, H., Timofte, R., Gool, L.V., 2022. Practical blind image denoising via swin-conv-unet and data synthesis. Machine Intelligence Research URL: <https://api.semanticscholar.org/CorpusID:247748724>.
- [50] Zhang, S., Gao, H., Rao, Q., 2021. Defense against adversarial attacks by reconstructing images. IEEE Transactions on Image Processing 30, 6117–6129. URL: <https://api.semanticscholar.org/CorpusID:235710687>.
- [51] Zhao, H., Gallo, O., Frosio, I., Kautz, J., 2017. Loss functions for image restoration with neural networks. IEEE Transactions on Computational Imaging 3, 47–57. URL: <https://api.semanticscholar.org/CorpusID:5334482>.
- [52] Zhou, D., Chen, Y., Wang, N., Liu, D., Gao, X., Liu, T., 2023. Eliminating adversarial noise via information discard and robust representation restoration, in: International Conference on Machine Learning. URL: <https://api.semanticscholar.org/CorpusID:260847065>.