

RAIN: A Simple Framework for Robust and Accurate Image Classification

Jiawei Du*, Hanshu Yan*, Vincent Y. F. Tan, *Senior Member, IEEE*, Joey Tianyi Zhou, Rick Siow Mong Goh, Jiashi Feng *Member, IEEE*

Abstract—It has been shown that the majority of existing adversarial defense methods achieve robustness at the cost of sacrificing prediction accuracy. The undesirable severe drop in accuracy adversely affects the reliability of machine learning algorithms and prohibits their deployment in realistic applications. This paper aims to address this dilemma by proposing a novel preprocessing framework, which we term *Robust and Accurate Image classificatioN* (RAIN), to improve the robustness of given CNN classifiers and, at the same time, preserve their high prediction accuracies. RAIN introduces a new randomization-enhancement scheme. It applies randomization over inputs to break the ties between the model forward prediction path and the backward gradient path, thus improving the model robustness. However, similar to existing preprocessing-based methods, the randomized process will degrade the prediction accuracy. To understand why this is the case, we compare the difference between original and processed images, and find it is the loss of high-frequency components in the input image that leads to accuracy drop of the classifier. Based on this finding, RAIN enhances the input’s high-frequency details to retain the CNN’s high prediction accuracy. Concretely, RAIN consists of two novel randomization modules: randomized small circular shift (RdmSCS) and randomized down-upsampling (RdmDU). The RdmDU module randomly downsamples the input image, and then the RdmSCS module circularly shifts the input image along a randomly chosen direction by a small but random number of pixels. Finally, the RdmDU module performs upsampling with a detail-enhancement model, such as deep super-resolution networks. We conduct extensive experiments on the STL10 and ImageNet datasets to verify the effectiveness of RAIN against various types of adversarial attacks. Our numerical results show that RAIN outperforms several state-of-the-art methods in both robustness and prediction accuracy. The proposed RAIN framework works in a plug-in manner and does not require to modify the given CNNs, thus easy and efficient for implementation.¹

Index Terms—Adversarial Robustness, Randomized Input Pre-processing,

I. INTRODUCTION

In recent years, convolution neural network (CNN)-based image classification models have been successfully applied to a variety of areas such as robot vision [1], biometrics authentication [2], and autonomous vehicles [3]. In these real-world applications, the CNN models not only should make

Jiawei Du is with the Department of Electrical and Computer Engineering, National University of Singapore and the Institute of High Performance Computing, A*STAR, Singapore.

Hanshu Yan, Vincent Y. F. Tan, and Jiashi Feng are with the Department of Electrical and Computer Engineering, National University of Singapore.

Joey Tianyi Zhou and Rick Siow Mong Goh are with the Institute of High Performance Computing, A*STAR, Singapore.

* Equal Contribution, e-mail: hanshu.yan@u.nus.edu.

¹The source codes of RAIN project are available at <https://github.com/dyjw9/RAIN>.

accurate predictions, but also need to be robust against small perturbations over input, i.e., the models should make consistent decisions for examples that have been contaminated by a small amount of noise. However, CNN classifiers are known to be highly vulnerable to adversarial examples [4], [5]—even certain visually imperceptible perturbations to inputs can easily fool the CNNs, resulting in grossly incorrect predictions. To ameliorate this problem, researchers have developed a variety of adversarial defense methods to protect CNN classifiers from adversarial attacks. These methods can be roughly divided into two categories: adversarial training [6]–[8] and input preprocessing [9]–[11].

The adversarial training (AT)-based methods [6], [7], [12], [13] train the CNN classifiers on the clean images as well as their adversarial versions generated on-the-fly during the model training procedure. Although they can effectively improve the model robustness, generating adversarial examples via gradient descent per training iteration is highly computationally demanding. Thus, the implementation of AT-based methods is inefficient in practice. Unlike the AT-based methods, the input preprocessing methods develop randomized or non-differentiable transformations to block the path of gradient backpropagation from the prediction to input, so that the attackers cannot access inner gradients to craft adversarial examples, thus, fail to fool the classifiers. These methods do not need to re-train the classifiers, and thus are more computationally efficient [9], [11], [14]–[17]. In real-world applications, classifiers are usually uninformed about whether each input is clean or adversarially perturbed. Thus, to ensure reliability, the classifiers have to treat clean and perturbed images on the same footing and should make correct predictions for both types of images. However, existing methods in both of these two categories suffer from the severe classification accuracy drop on clean inputs when improving the robustness of CNN classifiers [9], [18]. Mitigating the trade-off between accuracy and robustness is very challenging; existing works that target this problem are still few and far between.

In this work, we follow the strategy of input processing and aim to develop an efficient and effective defense framework, which also can alleviate the problem of the drop in accuracy. Our methodology is surprisingly straightforward: On the one hand, since adversarial attack methods mostly leverage the gradients with respect to the input to generate adversarial examples, to defend against adversarial examples, we can use randomization to decorrelate the inference forward path and the gradient backward path [9]. On the other hand, to maintain high prediction accuracy, we should utilize image transformations

to which CNN classifiers are (almost) invariant.

Taking into account these two aspects, we develop two novel randomized preprocessing modules—the randomized small circular shift (RdmSCS) and the randomized down-upsampling (RdmDU). In the RdmSCS module, input images are circularly shifted along a randomly chosen direction by a small and random number of pixels. Since CNNs are shown to be invariant to a small shift [19], the proposed RdmSCS will not degrade the natural accuracy by much. To further improve the robustness, we propose the RdmDU module, in which the input images are first downsampled through randomly sampling one pixel from each corresponding small patch. We then resize the low-resolution image into the original size with certain upsampling method. If using the widely-used bicubic interpolation for upsampling, RdmDU would also incur a clear drop in accuracy as the existing preprocessing-based methods do. To understand why the accuracy degrades, we conduct an empirical study which reveals that the loss of high-frequency details is a critical reason for the drop in accuracy. Based on this observation, we choose a detail-enhancement model, such as a deep super-resolution (SR) network, to reconstruct (i.e., obtain upsampled) high-quality images with rich (i.e., high-frequency) details. Consequently, the detail-enhancement RdmDU module manages to preserve high natural accuracy.

Combining these two randomized modules, we propose the *Robust and Accurate Image classificatioN* (RAIN) framework. The input image is processed sequentially as follows: (i) Randomized downsampling, (ii) RdmSCS, (iii) SR upsampling. We conduct extensive sets of experiments on the STL10 [20] and ImageNet [21] datasets to verify the effectiveness of the proposed RAIN. The results demonstrate that RAIN significantly enhances the robustness of given CNN classifiers (e.g., from 0% to 52.3% against C&W attack on the ImageNet) and does not lead to a severe drop in the natural accuracy (e.g., from 100% to 93.3% on ImageNet). In summary, the main contributions of this work are as follows:

- We propose the RAIN framework, which enhances the robustness and retains the high prediction accuracy of a given CNN classifier. This method achieves better or almost equal robustness against adversarial attacks and clearly better accuracy on clean images when compared to several state-of-the-art defense approaches [6], [7], [9], [15].
- We introduce two simple yet effective components within the RAIN framework: RdmSCS and RdmDU. Each of these modules can independently improve the robustness with a minor degradation to the prediction accuracy.
- We reveal that the loss or inadvertent removal of high-frequency components in input images is a critical factor in the drop of the accuracy of a CNN classifier. Hence, we are inspired to use a detail-enhancement model for upsampling within the RdmDU module.
- RAIN works in a plug-in manner and is easy to implement on given any CNN classifier.

The rest of this paper is organized as follows. In Section II, we review related existing adversarial defense methods. In Section III, we first elaborate on the two dedicated randomized

preprocessing modules, namely RdmSCS and RdmDU. Then, we explain why the use of a detail-enhancement model in the RdmDU helps to preserve good classification accuracy of the given CNN. Finally, we formally introduce the complete RAIN framework. In Section IV, we conduct extensive experiments to evaluate the ability of RAIN to robustify CNN classifiers and to preserve high natural accuracies for them. In Section V, we further study how the scale of shifting in the RdmSCS and the order of combination of two modules affect RAIN’s performance. We also illustrate how the RAIN framework purifies the feature maps of adversarial examples. We briefly summarize our work in Section VI and discuss several topics for future research.

II. RELATED WORK

Robustness measures the sensitivity of the accuracy of a classification model to perturbations of the inputs. Many approaches have been proposed to estimate the most harmful adversarial examples within a given neighborhood, such as FGSM [22], DeepFool [23], and PGD [6]. These methods compute adversarial perturbations based on available gradients evaluated at neighborhoods of the given input images. They can access the gradients in each layer and are thus called *white-box* attacks. In contrast, several other methods [24]–[27] propose to estimate the adversarial perturbation by using a number of queries. These approaches, without information about the gradients in each layer, are called *black-box* attacks. In this work, we use five white-box and two black-box attack methods to evaluate the robustness of classifiers.

To enhance the adversarial robustness of given classifiers, researchers have recently proposed a wide variety of defense methods. Here, we review two main classes of such methods, namely adversarial training and input preprocessing methods. We also discuss the connections of these methods to our work.

1) *Adversarial Training*: Adversarial training-based methods train CNN classifiers from scratch on both clean images and their corresponding adversarial examples to enhance their robustness. In each training epoch, the adversarial perturbations corresponding to the clean images are generated in real-time. The works in [6], [28] show that adversarial training is generally an effective defense mechanism against certain types of adversarial perturbations. However, training CNNs from scratch and computing adversarial examples in real-time has an adverse effect on the training cost; in fact, a study [29] has shown that it increases the training cost by 3-to-30 times. This computational burden restricts the applicability of adversarial training methods on large-scale datasets. To ameliorate the problem, researchers have proposed several methods for efficiently estimating adversarial perturbations [12], [30]. Besides, another severe drawback of adversarial training methods is that the resultant classifiers tend to over-fit a certain class of adversarial examples [29], [31], [32]. This results in a sharp drop in the prediction accuracy, and the classifiers are still vulnerable to other types of attacks.

2) *Input Preprocessing*: Preprocessing-based defense methods remove the adversarial information from input images by transforming images in a dedicated way before they are fed into

given CNN classifiers. One straightforward idea is to map input images onto the manifold generated by the clean data. To this end, several works [14], [16] propose to train generative models on the original dataset and use them to reconstruct images that are close to the adversarial examples for making predictions. Similarly, the work of [33] first impairs small patches of the adversarial examples and then performs reconstruction via matrix estimation. Traditional image processing operations have also been shown to be effective in robustifying given CNN classifiers. For example, Mustafa *et al.* [11] propose to use wavelet denoising and SR techniques to deactivate the adversarial perturbations. However, the enhanced robustness that results from these methods is due to the fact that the attackers are not provided with the internal gradients of the preprocessing modules [16], [28]. Athalye *et al.* [28] proposes the Backward Pass Differentiable Approximation technique and can successfully bypass this problem. The works of [9], [15] propose to randomize image processing operations such as padding and resizing. The randomization introduced effectively breaks the connection between the forward inference path and backward gradient path, and consequently, it improves robustness against adversaries. However, these randomization-based methods also suffer from deteriorations in prediction accuracies because CNN models are not invariant to the random padding and resizing operations. Athalye *et al.* [28] also demonstrates that the Expectation over Transformation technique is able to find the worst perturbation that can successfully mislead the CNN classifiers.

Our work develops an effective and efficient preprocessing-based defense framework, which amalgamates randomization concepts into two image processing modules that preserve the overall classification accuracy. RAIN is easy to implement and can defend EoT-based attacks if combined with adversarially trained CNN classifiers. The work that is most related to RAIN is presented in [11]. The difference on the usage of SR models is that the authors of [11] upsample input images with SR networks without performing downsampling. This will leads to a severe performance drop when the upsampling factor becomes large (e.g., 3 or 4) because the given CNN classifier is trained with images that are of the original size and vanilla CNNs are not scale-invariant [34].

III. RAIN: ROBUST AND ACCURATE IMAGE CLASSIFICATION NETWORKS

In this section, we first introduce the background on robustness evaluation and formally define the problem to solve. Then, we elaborate the proposed randomization-enhancement scheme that contains two novel preprocessing modules. Finally, we present the full RAIN framework and its implementation details.

A. Preliminaries

Given a CNN classifier $C(\cdot)$ and a perturbation budget ϵ , the ideal untargeted adversarial example x_ϵ^{adv} of an input x is defined as the data point in its ϵ -neighborhood that maximizes certain classification loss. Formally,

$$x_\epsilon^{\text{adv}} = \arg \max_{x' \in \mathcal{B}(x, \epsilon)} l(C(x'), y), \quad (1)$$

where $l(\cdot, \cdot)$ is the loss function, y is the label of x , and $\mathcal{B}(x, \epsilon)$ is the ϵ ball around x , which defined based on certain norm, i.e., $\mathcal{B}(x, \epsilon) = \{x' : \|x' - x\| \leq \epsilon\}$. Most of existing attack methods [6], [22] are based on the gradient descent to solve the optimization problem in Equation (1). In the following, we use $A_\epsilon(\cdot, \cdot)$ to represent certain adversarial attack method and the ideal adversarial example x_ϵ^{adv} is approximated by $A_\epsilon(C, x)$.

The robustness of the given CNN $C(\cdot)$ can be quantified by the classification accuracy over the adversarial examples of input images. Since it is less meaningful to attack inputs that are already classified incorrectly, we follow the principle in [9], [15] and randomly collect a set of labelled images $\mathcal{T} = \{(x_j, y_j)\}_{j=1}^M$ for robustness evaluation, where (x_j, y_j) denotes the j th sample *correctly classified* by $C(\cdot)$ before being attacked, i.e., $C(x_j) = y_j$. For a particular adversarial attack method $A_\epsilon(\cdot, \cdot)$, the robustness of $C(\cdot)$ is evaluated as

$$R_{\mathcal{T}}(C, A_\epsilon) = \frac{1}{M} \sum_{j=1}^M \mathbb{1}[C(A_\epsilon(C, x_j)) = y_j]. \quad (2)$$

To enhance the robustness of CNN classifiers, many defense methods have been proposed. Let $C'(\cdot)$ be the robustified version of $C(\cdot)$. The natural accuracy of $C'(\cdot)$ is quantified by the classification accuracy on the collected set \mathcal{T} , i.e.,

$$P_{\mathcal{T}}(C') = \frac{1}{M} \sum_{j=1}^M \mathbb{1}[C'(x_j) = y_j]. \quad (3)$$

Although most of the existing defense methods manage to robustify given classifiers, they suffer from severe drops in accuracy [10], [18]. For example, on the ImageNet, the Pixel Deflection method improves the robustness of the given CNN from 0 to 32.0%, but the natural accuracy on \mathcal{T} drops from 100% to 85.8% (refer to Table III). To ensure the reliability of classification in real applications, the CNN classifiers have to make correct predictions for both of the clean images and their perturbed examples. In this work, we aim to develop a preprocessing-based method to mitigate this trade-off, i.e., for the robustified model C' , $R_{\mathcal{T}}(C', A_\epsilon)$ is significantly higher than $R_{\mathcal{T}}(C, A_\epsilon)$ and the accuracy of $C'(\cdot)$ over \mathcal{T} should remain above a high value.

B. Randomized Input Preprocessing

Inspired by the invariance of CNNs to slight shifts and scalings [19], [35], we aim to utilize these properties to design robust and accurate CNN models. In particular, we consider two randomized transformation modules: RdmSCS and RdmDU. One the one hand, randomization can result in a mismatch between the forward prediction path and the gradient back-propagation path [9]. Hence, these two transformations are expected to be able to improve the adversarial robustness of classifiers. On the other hand, CNNs are capable of learning translation-invariant representations due to the gradual increase in the size of the receptive field [36] and the pooling operations [37]. The downsampling and subsequent upsampling operations also do not change the size of the images and the positions of the contents. Thus, the two transformations of interest would likely not degrade the prediction accuracy significantly.

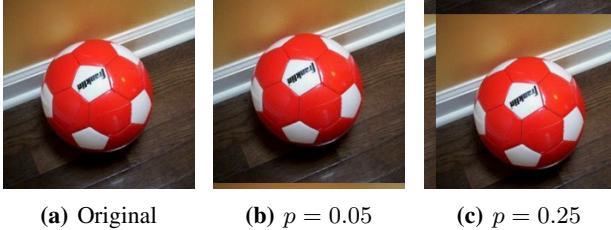


Fig. 1: An illustration of the RdmSCS module. Left: an input image; Mid: the input shifted with $p = 0.05$; Right: the input shifted with $p = 0.25$.

1) *Randomized Small Circular Shift*: We first consider a random shift operation of the input. Before it is fed to a given CNN model, the input is shifted by a random number of pixels along a randomly selected direction (horizontal or vertical). More precisely, two shifting parameters, Δh and Δw , are randomly sampled from uniform distributions, i.e., $\Delta h \sim \mathcal{U}(-hp, hp)$, $\Delta w \sim \mathcal{U}(-wp, wp)$, where h and w are the height and width of the input image respectively, and p is a predefined positive constant. The sign of Δh (or Δw) represents the vertical (or horizontal) direction for shifting, while the magnitude quantifies the number of pixels to be shifted. Algorithm 1 elaborates on the implementation of the proposed random shifting operation. We use the superscript s to denote the shifted version of an image.

Algorithm 1 Randomized Small Circular Shift

```

1: function RDMSCS( $x, p$ )
2:   for  $i_s$  in  $(0, h - 1)$  and  $j_s$  in  $(0, w - 1)$  do
3:     Randomly initialize
        $\Delta h \sim \mathcal{U}(-hp, hp)$ ,  $\Delta w \sim \mathcal{U}(-wp, wp)$ 
4:      $i \equiv (i^s + \Delta h) \bmod h$ 
5:      $j \equiv (j^s + \Delta w) \bmod w$ 
6:      $x^s(i^s, j^s) = x(i, j)$ 
return  $x^s$ 
```

Conventional translation operations pad zeros into the empty regions [38]. Instead, our proposed operation splices back the shifted part of the input image circularly from the opposite boundary (see Figure 1). The aim of this operation is to retain the frequency components and other statistics (such as the means and variances of pixel values) of original images as much as possible so that the preprocessing operation would not have an obvious negative impact on the prediction accuracy. For this reason, we usually set p to be a small value in practice. Accordingly, we dub the proposed shifting operation as *Randomized Small Circular Shift* (RdmSCS).

To understand how the value of p affects the efficacy of RdmSCS, we evaluate the performance of the given CNN classifier equipped with RdmSCS. From Table I, we can see that the proposed RdmSCS clearly improves on the robustness of CNN classifiers when p is greater or equal to 0.05. Choosing a large value of p (e.g., 0.15) can significantly robustify the CNN classifiers, but will degrade the accuracy on clean images because CNNs are not invariant to a large degree of shifting. To achieve a good balance between accuracy and robustness,

TABLE I: Evaluation of RdmSCS on the robustness and prediction accuracy on STL10 dataset (Refer to Section IV-A for the setting of attack methods).

STL10	Clean Images	Accuracy			Robustness		
		FGSM	C&W	Deep Fool	FGSM	C&W	Deep Fool
CNN-only	1.000	0.074	0.000	0.000	0.074	0.000	0.000
$p = 0.01$	0.975	0.075	0.117	0.000	0.075	0.117	0.000
$p = 0.05$	0.958	0.215	0.421	0.093	0.215	0.421	0.093
$p = 0.1$	0.925	0.311	0.523	0.252	0.311	0.523	0.252
$p = 0.15$	0.836	0.299	0.671	0.422	0.299	0.671	0.422

the value of p should be relatively small (less or equal to 0.1), so that the natural accuracy can be retained at a high value (above 90%).

2) *Randomized Down-Upsampling*: In addition to the RdmSCS operation, we also introduce a *Randomized Down-Upsampling* (RdmDU) module to further improve the robustness which simultaneously preserves the classification accuracy. Downsampling and subsequent upsampling do not change the size of the input image and its semantic contents, thus, they are expected to maintain a good classification performance of the given CNN over clean images. The RdmDU module consists of a randomized downsampling operation followed by an upsampling operation. Given an input image $x \in \mathbb{R}^{h \times w \times 3}$, the random downsampling operation of RdmDU partitions the input into non-overlapping 2×2 patches and randomly picks one pixel from the four in each patch with the same probability $1/4$. Consequently, the resultant image, denoted as x^\downarrow , has size of $\frac{1}{2}h \times \frac{1}{2}w \times 3$. As the CNN classifier is trained with images of the original size, the upsampling operation of RdmDU reconstructs a high-resolution image of size $h \times w \times 3$, denoted as $x^{\downarrow\uparrow}$, from x^\downarrow . Algorithm 2 shows the implementation of the RdmDU module in detail.

Algorithm 2 Randomized Down-Upsampling

Input: Input image x ;

- $x = \text{RANDOM_DOWNSAMPLING}(x)$
- $x = \text{Upsampling}(x, \text{factor} = 2)$

Output: x

- function** RANDOM_DOWNSAMPLING(x)
- for** i in $(0, \lceil \frac{1}{2}h \rceil - 1)$ and j in $(0, \lceil \frac{1}{2}w \rceil - 1)$ **do**
- Randomly initialize
 $\Delta i \sim \mathcal{U}([0, 1])$, $\Delta j \sim \mathcal{U}([0, 1])$
- $x^\downarrow(i, j) = x(\min(2i + \Delta i, h - 1),$
 $\min(2j + \Delta j, w - 1))$

return x^\downarrow

For the upsampling operation in the RdmDU module, we first consider the bicubic interpolation method [39], which is efficient and widely used in image processing tasks. We conduct experiments on the STL10 dataset to evaluate the effectiveness of the RdmDU-Bicubic module. From Table II, we find that both of the RdmDU-Bicubic can enhance the robustness of the CNN classifier against various kinds of attacks. However, the prediction accuracy of RdmDU-Bicubic modification drops significantly to 81.2%, which is undesirable in high-stakes applications, such as the autonomous vehicles. To understand why upsampling with bicubic interpolation in

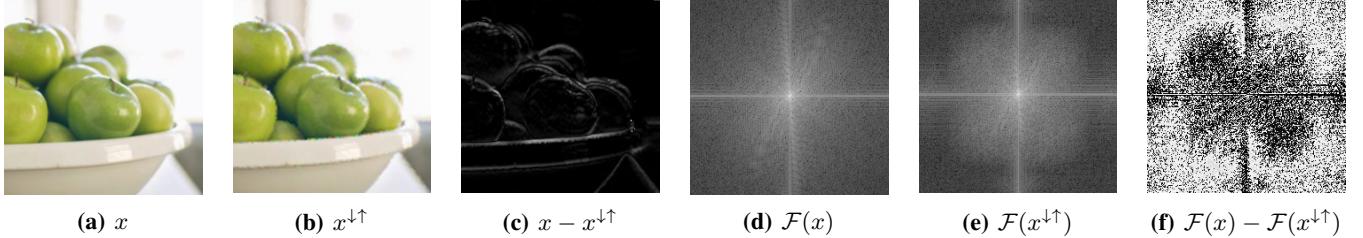


Fig. 2: Comparison between an input and the output of RdmDU-bic. From left to right: input x , output $x^{\downarrow\uparrow}$, difference $x - x^{\downarrow\uparrow}$; spectrum $\mathcal{F}(x)$, spectrum $\mathcal{F}(x^{\downarrow\uparrow})$, and $\mathcal{F}(x) - \mathcal{F}(x^{\downarrow\uparrow})$.

RdmDU degrades the prediction accuracy, we then conduct an empirical study in Section III-C and find that the loss of high-frequency components of original images leads to the drop in the classification accuracy. Based on the findings, we propose to replace the bicubic with a detail-enhancement model, so that the RdmDU module can reconstruct an image with rich high-frequency details and retain the CNN’s accuracy at a high value (above 90%). Finally, we use a deep learning-based super-resolution model, namely the EDSR [40] network, as the upsampling operation. The experimental results show that RdmDU-SR can preserve 93.5% natural accuracy, but also achieves consistently better adversarial robustness.

TABLE II: Evaluation of RdmDU modules on the robustness and prediction accuracy on STL10 dataset (Refer to Section IV-A for the setting of attack methods).

	Accuracy		Robustness		
	STL10	Clean	FGSM	C&W	Deep Fool
CNN-only	1.000	0.074	0.000	0.000	
w/ RdmDU-Bicubic	0.812	0.485	0.625	0.539	
w/ RdmDU-SR	0.935	0.565	0.742	0.664	

C. What is Lacking in Preserving Classification Accuracy?

To discover why the bicubic interpolation leads to a severe drop in natural accuracy, we first compare the original input and the output images of RdmDU-Bicubic. From Figure 2, we see that visually, the difference, if any, between a certain image x and its processed version $x^{\downarrow\uparrow}$ is imperceptible. However, if we transform the image into the frequency domain using the 2D-Fast Fourier Transform (FFT) $\mathcal{F}(\cdot)$, we can see that the RdmDU-Bicubic tends to suppress the high-frequency components of the original image. Therefore, we hypothesize that the accuracy drop arises because of the loss of high-frequency components. To corroborate this hypothesis, we then conduct experiments to study the contribution of high-frequency components to the accuracy of a well-trained CNN model.

Setup: Given a CNN classifier, we remove an increasing amount of high frequency components from input images and examine how the prediction accuracy changes with respect to the loss of high-frequency components. Specifically, we preprocess input images in two steps, namely low-pass filtering and energy normalizing:

a). The low-pass filter removes the high frequency components above a certain threshold. We first compute the spectrum z of an input image x via the 2D-FFT, i.e., $z = \mathcal{F}(x)$. Next,

we remove the components with frequencies above a given threshold r , and obtain the resultant spectrum z' , where

$$z'(\mu, \nu) = z(\mu, \nu) \cdot \mathbb{1} \left(\frac{d((\mu, \nu), (c_h, c_w))}{1/2\sqrt{h^2 + w^2}} \leq r \right). \quad (4)$$

Here $d(\cdot, \cdot)$ denotes the Euclidean distance between two points and (c_h, c_w) is the position of the center, which corresponds to the zero-frequency component. The larger the value of r is, the more high-frequency components are retained.

b). Removing high-frequency components reduces the energy of the spectrum $E(z)$, where

$$E(z) \equiv \int_{\mu} \int_{\nu} |z(\mu, \nu)|^2 d\mu d\nu. \quad (5)$$

To ensure that the comparisons are fair, we normalize the spectrums of the images before and after processing so that they have the same energy. Thus, we multiply the resultant spectrum z' with $\sqrt{E(z)/E(z')}$ and reconstruct the RGB image for making subsequent predictions.

Findings: We conduct experiments on the STL10 and ImageNet datasets. For each dataset, we plot a curve of the performance of the classifier as a function of the threshold on the frequency r . From Figure 4, we find that, on the two datasets, the prediction accuracies drop rapidly as r decreases, corresponding to an increase in the amount of high-frequency components being filtered out. Thus, the high-frequency components of the input images are critical for maintaining the accuracy of predictions. Given this observation, to preserve high natural accuracy, we propose to select an upsampling method that can super-resolve images while generating rich details.

RdmDU-SR: Deep-learning-based SR methods have been shown to achieve impressive performances on SR tasks [40], [41]. They are capable of generating high-frequency details such as rich textures and sharp edges. Here, we adopt the EDSR [40] model, which is among the most effective implementations for SR tasks and has been widely used as the backbone of various SR algorithms. We transform images processed by the EDSR into the frequency domain and compare the spectrums with those of images processed by the bicubic upsampling. The spectrums are shown in Figure 3, and we also calculate the heatmaps of the difference between the spectrums of processed images and original images.

From Figure 3, we can observe that, in each row, the four corners of $\mathcal{F}(x) - \mathcal{F}(x_{\text{Bicubic}}^{\downarrow\uparrow})$ are obviously bright, while the corners and the center of $\mathcal{F}(x) - \mathcal{F}(x_{\text{EDSR}}^{\downarrow\uparrow})$ are all relatively

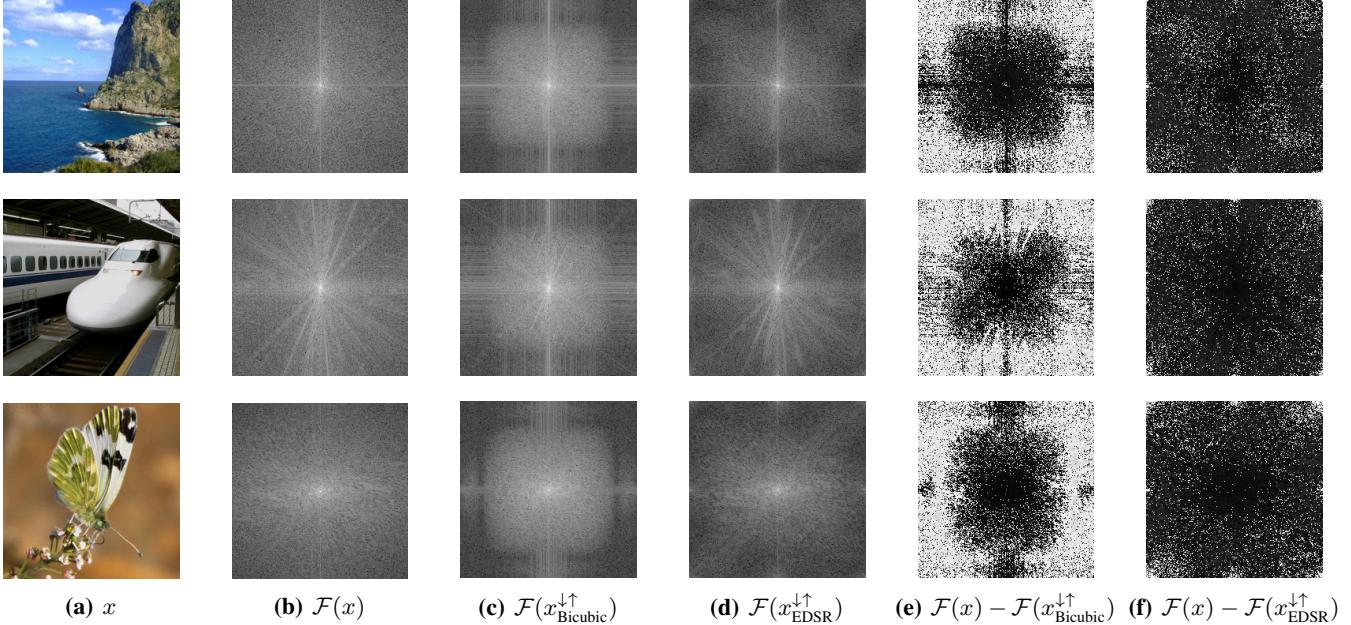


Fig. 3: Comparison between images processed by RdmDU-Bicubic and RdmDU-SR. From left to right, column (a), original images x ; column (b), spectrum of x ; column (c), spectrums of images processed by RdmDU-Bicubic $\mathcal{F}(x_{\text{Bicubic}}^{\downarrow\uparrow})$; column (d), spectrums of images processed by RdmDU-SR $\mathcal{F}(x_{\text{EDSR}}^{\downarrow\uparrow})$; column (e), difference maps $\mathcal{F}(x) - \mathcal{F}(x_{\text{Bicubic}}^{\downarrow\uparrow})$; column (f), difference maps $\mathcal{F}(x) - \mathcal{F}(x_{\text{EDSR}}^{\downarrow\uparrow})$.

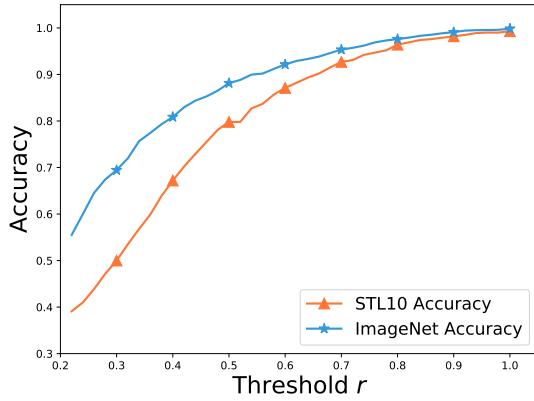


Fig. 4: Accuracy v.s. the amount of high-frequency components removed. The smaller r is, the more high-frequency components are removed. On both two datasets, we can see even the energy of the filtered spectrum is controlled to be the same with the origin, the loss of high-freq components still impairs the test accuracy.

dark. It means that, in comparison to the bicubic interpolation, EDSR can effectively reconstruct high-frequency components of the original images. According to the previous findings, we can deduce that RdmDU-SR preserves the high natural accuracies of CNN classifiers because of its ability to generate rich details and textures.

D. RAIN Framework

In the previous sections, we described two randomized preprocessing operations, both of which were shown to improve

the robustness of CNN classifiers with minor degradations to the prediction accuracies. Here, we formally introduce our *Robust and Accurate Image Classification* framework, which combines the RdmSCS and the RdmDU.

The whole pipeline of the RAIN is as follows (see Figure 5). For a given well-trained classification model $C(\cdot)$, RAIN first downsamples a certain image x through the randomized downsampling operation within the RdmDU module and shifts the resultant image via the RdmSCS module successively. Then, a well-trained EDSR model upsamples the image so that the resulting size is the same as the original image, and the details are also enriched. Lastly, the resultant image is fed to the given CNN classifier $C(\cdot)$ for subsequent prediction. To generate image details that are most useful for $C(\cdot)$, we fix the CNN's parameters and fine-tune the EDSR model's parameters to minimize the loss of the classification task over a few training epochs.

IV. EXPERIMENTS

In this section, we conduct extensive sets of experiments to evaluate the effectiveness of RAIN in terms of the enhancement of the CNNs' robustness and the preservation of their accuracies. We first examine how well the prediction accuracies of the CNN equipped with RAIN on clean images are retained, and evaluate the robustness against various white-box attacks and black-box attacks. Then, we also compare the proposed RAIN framework with other state-of-the-art methods to show the advantages of RAIN on robustness and accuracy. RAIN is a preprocessing-based method, we finally show that RAIN can work in conjunction with adversarial training. When dealing with the Expectation over Transformation (EoT) attacks,

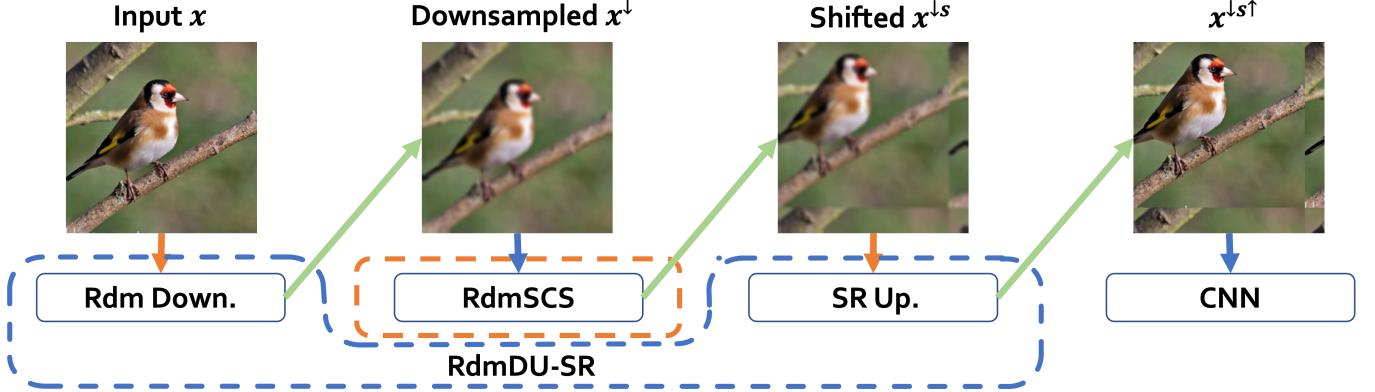


Fig. 5: The pipeline of our proposed RAIN framework: given an input, it is sequentially downsampled via RdmDU, circularly shifted via RdmSCS, and upsampled by an SR model within the RdmDU.

equipped with RAIN, the adversarially trained CNN classifiers achieve better natural accuracies.

A. Experimental Setting

We use five white-box attack methods [6], [22], [23], [42], [43] and three black-box methods [24], [25], [44] to evaluate the effectiveness of the proposed RAIN. Since experiments are interspersed throughout Section IV and V, here we describe the common experimental settings.

1) *Models and Datasets*:: We conduct experiments on the STL10 [20] and / or the ImageNet [21] datasets. We use a seven-layer CNN classifier for the STL10 dataset and the Resnet-50 [45] model for the ImageNet. Details about the network architectures can be found in Section VI-A. For each dataset, we randomly select 5,000 images from the test set for robustness evaluation, i.e., $M = 5,000$ in Equation (2). Besides, the performance of defense methods on original nature images is also evaluated on \mathcal{T} .

2) *White-box Attack Setting*:: For the white-box attacks, we consider one single-step attack method (FGSM [22]) and four iterative attack methods (PGD [6], C&W [42], Deepfool [23] and EAD [43]). The FGSM and Deepfool attacks are based on the l_∞ norm with $\epsilon = 8/255$. The PGD attack is based on the l_∞ norm with $\epsilon = 8/255$ and learning rate $2/255$. The C&W attack is based on l_2 norm with learning rate $1/255$. The EAD attack is implemented in EN-rule with learning rate $1/255$ and $\beta = 0.01$. All the four iterative attacks work over 40 iterations of gradient descent.

3) *Black-box Attack Setting*:: For the black-box attacks, we consider the ZOO method [24] with a maximum of 10,240 queries, the NES [25] with a maximum of 2,000 queries, and the RayS [44] with a maximum 10,000 queries. All of these methods are based on the l_∞ norm with $\epsilon = 8/255$. All the experiments follow the same settings discussed above unless otherwise specified.

4) *EoT-based Attack Setting*:: We evaluate the adversarial robustness against FGSM, PGD, and C&W attacks under the EoT framework, in which each type of attack estimates the expected worst perturbation by averaging the obtained adversaries over a fixed number of queries. For the FGSM attack, the perturbation budget ϵ is $8/255$ based on the l_∞ norm;

The C&W attack is based on l_2 norm with rate $1/255$; The PGD-1 attack based on l_∞ with $\epsilon = 8/255$ and rate $2/255$. Both of the PGD-1 and C&W work in 40 iterations. All of these three attacks query the classifier for 40 times. We also use the PGD attack in another setting: based on l_∞ , PGD-2 works in 10 iterations with $\epsilon = 8/255$ and rate $2/255$. The number of queries is set to be 10.

B. Adversarial Robustness of RAIN against White- and Black-box attacks

We evaluate the robustness of our proposed RAIN on both of the STL10 [20] and ImageNet [21] datasets under various white-box and black-box attacks. The experimental results demonstrate that the proposed RAIN framework improves the adversarial robustness of given CNNs, and, at the same time, manages to preserve high natural accuracies (>90% on both the two datasets).

We also compare RAIN with other four recent adversarial defense methods: two of them are preprocessing-based defense methods, namely Random Resizing & Padding (RandR&P) [9] and Pixel Deflection (PixelDef) [15]; the other two are adversarial training-based methods, namely Madry *et al.* [6] and Feature Denoising (FeatDn) [7] (both are trained with PGD examples of $\epsilon = 16/255$, step = $1/255$ and 40 iterations).

White-box Attack: From Table III, we observe that RAIN enhances the adversarial robustness of given CNNs, from 0% to 78.1% on STL10 against the C&W attack and from 0% to 86.7% on ImageNet against the Deep Fool attack. In comparison to the two preprocessing-based competitors, RAIN achieves better adversarial robustness (e.g., 52.3% v.s. 43.8% against C&W attack on ImageNet) but also maintains higher natural accuracies (e.g., 92.9% v.s. 90.7% of RandP&R on STL10). For the two AT-based methods, RAIN also shows remarkable advantages regarding both robustness and natural accuracy. In specific, the accuracies of RAIN on both two datasets are at least 20 points higher than those of the two AT-based methods (e.g., 92.9% v.s. 70.5% on the STL10).

Black-box Attacks: We compare the proposed RAIN framework to the baseline methods under black-box adversarial attacks. Table IV shows that RAIN framework achieves the best black-box adversarial robustness against the ZOO and RayS

TABLE III: Robustness comparison of different defense methods against white-box attacks. We report the classification accuracies on clean images and various types of adversarial examples. The results show that the proposed RAIN framework achieves the best performance in prediction accuracy and robustness. The methods of [6] and [7] train classifiers with PGD adversarial examples. Thus, we do not evaluate their robustness with respect to PGD attacks for comparison.

		Accuracy		Robustness				
STL10		Clean		FGSM	PGD	C&W	Deep Fool	EAD
CNN-only		1.000		0.074	0	0	0	0
RAIN		0.929		0.681	0.045	0.781	0.734	0.828
PixelDef [15]		0.883		0.397	0	0.320	0.117	0.407
RandP&R [9]		0.907		0.376	0.025	0.699	0.305	0.789
Madry [6]		0.705		0.607	-	0.289	0.234	0.028
FeatDn [7]		0.696		0.615	-	0.300	0.227	0.034
ImageNet		Clean		FGSM	PGD	C&W	Deep Fool	EAD
CNN-only		1.000		0.125	0	0	0	0
RAIN		0.933		0.625	0.012	0.523	0.867	0.857
PixelDef [15]		0.858		0.357	0	0.320	0.210	0.365
RandP&R [9]		0.928		343	0	0.438	0.703	0.759
Madry [6]		0.623		0.598	-	0.355	0.426	0.063
FeatDn [7]		0.653		0.603	-	0.398	0.410	0.115

attacks (91.2% on STL10 and 88.5% on ImageNet against ZOO, 88.2% on STL10 and 92.4% on ImageNet against RayS). In terms of the NES attack, the performance of RAIN is very close to that of the Random Resizing & Padding method [9] (87.1% v.s. 87.3% on STL10 and 88.2% v.s. 88.1% on ImageNet), and is clearly better than the other three baselines competitors.

TABLE IV: Evaluation of the robustness of the proposed RAIN framework and the baseline methods under black-box attacks on STL10 and ImageNet datasets.

		Accuracy		Robustness		
STL10		Clean		ZOO	NES	RayS
CNN-only		1		0	0.02	0
RAIN		0.929		0.912	0.871	0.882
PixelDef [15]		0.883		0.679	0.650	0.743
RandP&R [9]		0.907		0.854	0.873	828
Madry [6]		0.705		0.705	0.663	0.265
FeatDn [7]		0.696		0.621	0.594	0.304
ImageNet		Clean		ZOO	NES	RayS
CNN-only		1		0	0.02	0
RAIN		0.933		0.885	0.882	0.924
PixelDef [15]		0.858		0.846	0.841	0.848
RandP&R [9]		0.928		0.867	0.881	0.913
Madry [6]		0.623		0.620	0.611	0.396
FeatDn [7]		0.653		0.663	0.641	0.421

C. Working in Conjunction with AT against EoT-based Attacks

Defense methods [9], [10], [15], [16] based on gradient obfuscation can block the back-propagation path of gradients from the predictions to inputs. Consequently, the attacker cannot accurately approximate the gradients to generate adversarial perturbations. However, Athalye *et al.* [28] demonstrated that attackers can circumvent the gradient obfuscation problem through the Expectation over Transformation (EoT) technique for randomization methods or the Backward Pass Differentiable Approximate (BPDA) technique for non-differentiable transformations. Since our RAIN framework is built with two randomized transformations, RdmSCS and RdmDU, here we will discuss the effectiveness of RAIN against EoT-based attacks.

Given an input x , it is firstly preprocessed by some random transformation T before fed into the classifier $C(\cdot)$. Assume the random transformation T is i.i.d. sampled from the distribution of transformation P_T , the EoT-based attack calculates the gradients of the expected prediction, $\nabla_x \mathbb{E}_{T \sim P_T}[C(T(x))]$. Since

$$\nabla_x \mathbb{E}_{T \sim P_T}[C(T(x))] = \mathbb{E}_{T \sim P_T}[\nabla_x C(T(x))], \quad (6)$$

the gradients of the expected prediction can be estimated by averaging the gradients of predictions of x over multiple queries. When the number of query times is sufficiently large, the randomization-based defense will fail [28].

RAIN-AT: To deal with the EoT-based attack, it is imperative to perform adversarial training (AT) [6] on CNN classifiers. Within the RAIN framework, an input image is first modified with the RdmSCS and RdmDU-SR modules arranged in sequence; then, the processed image is sent to the *adversarially trained CNN* classifier for making a prediction. To achieve the best performance, we then fine-tune the parametric EDSR model *only* with natural clean samples (denoted as RAIN-AT).

Conventional adversarial training suffers from the undesired trade-off between the prediction accuracy on natural images and the robustness against adversarial attacks [18]. Equipped with our RAIN framework, we find that the adversarially trained CNN classifiers achieve higher natural accuracies but with no degradation to the robustness. Besides, since the pre-trained EDSR only needs to be fine-tuned with one epoch (on both the STL10 dataset and the ImageNet dataset), the implementation of RAIN does not incur too much extra computational cost.

Evaluation: From Table V, we observe that, based on normally trained CNN classifiers, RAIN and the other two preprocessing-based methods, namely PixelDef [15] and RandP&R [9], are totally misled by EoT-based attacks (The accuracy is almost equal to 0% against PGD attacks). In contrast, the adversarially trained CNNs [6] can effectively defend EoT-based attacks. More specifically, Madry *et al.* [6] achieves 36.7% accuracy against the PGD attack on the STL10 dataset and 36.1% against PGD attack on the ImageNet dataset.

Based on the adversarially trained CNNs (Madry [6]), the RAIN framework further enhances adversarial robustness and

TABLE V: Robustness comparison of different defense methods against EoT-based attacks. We report the classification accuracies on clean images and various types of adversarial examples. The results show that the proposed RAIN framework with adversarially trained CNNs (RAIN-AT) achieves the best robustness against various attacks.

	Accuracy		Robustness		
	Clean	FGSM	PGD-1	PGD-2	C&W
PixelDef [15]	0.883	0.082	0	0	0.283
RandP&R [9]	0.907	0.125	0.015	0.024	0.365
RAIN	0.929	0.188	0.01	0.023	0.402
Madry [6]	0.705	0.607	0.367	0.363	0.294
RAIN-AT	0.745	0.647	0.361	0.375	0.605
RandP&R-AT	0.619	0.575	0.364	0.374	0.428
ImageNet	Accuracy		Robustness		
	Clean	FGSM	PGD-1	PGD-2	C&W
PixelDef [15]	0.858	0.022	0	0	0.247
RandP&R [9]	0.928	0.067	0	0	0.276
RAIN	0.933	0.125	0	0	0.255
Madry [6]	0.623	0.598	0.361	0.361	0.349
RAIN-AT	0.656	0.630	0.363	0.377	0.473
RandP&R-AT	0.569	0.574	0.347	0.356	0.430

improves the prediction accuracy on natural images. More precisely, RAIN-AT achieves 65.6% natural accuracy against 62.3% of Madry [6] on the ImageNet dataset and 74.5% against 70.5% on the STL10 dataset. The Random Padding & Resizing [9] method also can be combined with adversarially trained CNNs (RandP&R-AT). We can see that our RAIN-AT outperforms RandP&R-AT consistently on all types of attacks. Unlike the dedicated preprocessing modules in our RAIN, the random padding and resizing operations results in a severe drop of accuracy (around 10 points lower than ours on both two datasets).

In summary, our RAIN and the adversarial training can work together to defend EoT-based attacks. Equipped with RAIN, the adversarially trained CNN can achieve better accuracy on clean images.

V. ABALTAION STUDY

A. On the Combination Order of RdmSCS and RdmDU

The two randomized preprocessing modules can be combined in three different orders: RdmSCS+RdmDU, RdmDU+RdmSCS, and RAIN. Here, we conduct experiments on the STL10 dataset to show why we choose the RAIN framework with the modules arranged as in Figure 5 instead of the other two variants. The experimental settings have been described in Section IV-A.

From Table VI, we see that all of the three variations are clearly more robust than the original CNN model. In comparison to the results in Table I and Table II, the combination of RdmSCS and RdmDU-SR clearly outperforms the separate utilization of these two modules for enhancing robustness. In real-world applications, since the probability that input images are naturally clean is much higher than that of adversarial examples, we choose the order of ‘Downsampling-Shifting-Upsampling’, i.e., RAIN, which achieves the highest natural accuracy (92.9%) among the three variants.

B. On the Selection the Value of p

In the proposed RAIN framework, the value of p is an important hyperparameter; this parameter affects the robustness

TABLE VI: Comparison of the three type of combinations of RdmSCS and RdmDU-SR (with $p = 0.05$). Among the three types, the type RAIN achieves the best natural accuracy.

	Accuracy		Robustness	
	Clean	FGSM	DeepFool	C&W
CNN-only	1.000	0.074	0	0
RdmSCS-RdmDU	0.909	0.675	0.737	0.843
RdmDU-RdmSCS	0.899	0.666	0.768	0.797
RAIN	0.929	0.680	0.734	0.781

as well as the accuracy on natural images. We conduct experiments on the STL10 dataset to examine the effectiveness of RdmSCS with different values of p .

From Table VII and Table I, we observe that, for each value of p , the combination of RdmSCS and RdmDU-SR (RAIN) results in significantly better adversarial robustness in comparison to using RdmSCS only. However, introducing an extra preprocessing module also deteriorates the natural accuracy on clean images. To preserve a high natural accuracy (above 90%), we should set the value of p to be less or equal to 0.1. Regarding the adversarial robustness, we choose $p = 0.05$, which achieves the best robustness against FGSM and Deep Fool attacks and also maintains a 92.9% natural accuracy.

TABLE VII: Evaluation of RAIN with various values of p on the robustness and prediction accuracy on STL10 dataset.

	Accuracy		Robustness	
	Clean	FGSM	DeepFool	C&W
$p = 0.01$	0.938	0.623	0.695	0.726
$p = 0.05$	0.929	0.683	0.734	0.781
$p = 0.1$	0.909	0.64	0.677	0.820
$p = 0.15$	0.774	0.59	0.703	0.671

C. Effects of Adversarial Noise on Feature Maps

Adversarial attacks add imperceptible noise to input images and results in a catastrophic failure of CNN classifiers. Even though the magnitude of perturbations on input images are small, the negative effects will be amplified across the hidden layers. In CNN classifiers, features generated at the last convolutional layer can reflect the semantic information and

indicate pixels that are critical for subsequent prediction [7], [11]. To investigate the efficacy of our RAIN framework, we extract and compare the feature maps of clean images and their adversarial examples.

We train a ResNet50 model on the ImageNet dataset and collect two images from the ImageNet for illustration. We use FGSM attack with $\epsilon = 8/255$ to generate the adversarial examples, x_ϵ^C , of a certain input x . Here, the superscript ‘C’ means the target model is the CNN only. We extract the features of x and x_ϵ^C in the last convolutional layer for comparison. From Figure 6, we see that the feature maps $F(x)$ focus on the region corresponding to the class objects in the original images (e.g., in the second row, the highlighted pixels corresponds to the position of airplane). We subtract the feature map of the adversarial example x_ϵ^C from that of x and obtain the difference map $F(x_\epsilon^C) - F(x)$. It is obvious that the adversarial feature map highlights many other irrelevant pixels in addition to the critical ones, which could lead to wrong predictions downstream. When the CNN classifier is equipped with RAIN, the adversarial examples are crafted in an end-to-end manner, i.e., the gradients of prediction are backpropagated through the defense modules. We denote the adversarial example in the RAIN framework as x_ϵ^R . From the difference map $F(x_\epsilon^R) - F(x)$ in Figure 6, we observe that the feature map of x_ϵ^R is close to the feature map of clean image x . In comparison to results in column (c), the RAIN framework manages to suppress the malicious effects of the adversarial noise in the hidden layer.

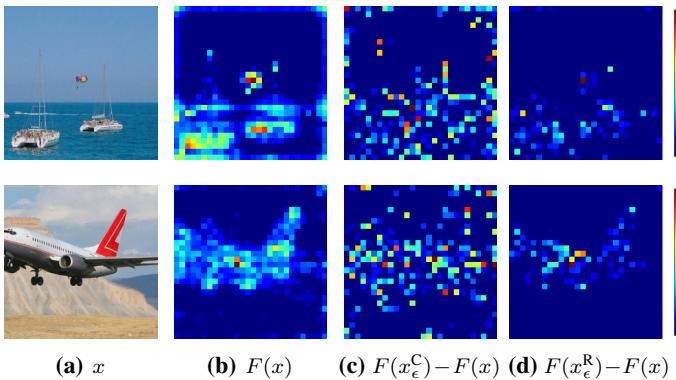


Fig. 6: A comparison on the feature maps of clean images and their adversarial examples. In each row, column (a) shows the natural image x ; In column (b), $F(x)$ denotes the feature map of x extracted from the given CNN; column (c) shows the difference between the feature map of x and that of the adversarial example x_ϵ^C generated w.r.t. the CNN; column (d) shows the difference between the feature map of x and that of the adversarial example x_ϵ^R crafted w.r.t. the CNN equipped with our RAIN.

VI. CONCLUSIONS

This paper proposes a novel adversarial defense framework, RAIN, which combines two randomization modules, namely RdmSCS and RdmDU. As shown through our extensive experiments on two benchmark datasets, , RAIN can significantly enhance the robustness of given CNN classifiers and

simultaneously preserve their high classification accuracies. RAIN preserves the reliability of CNN classifiers when the set of inputs is a combination of clean images and their adversarially perturbed versions.

In the future, we aim to develop a defense framework that is applicable to the classification of low-resolution and low-quality images because the quality of such images can be severely degraded by preprocessing operations, such as downsampling. The degradation may be so severe that effective SR models cannot reconstruct images to a level that predictions can be reliably made. We also plan to analyze the fundamental limits of the trade-off between prediction accuracy and adversarial robustness in a theoretically-grounded manner. This may lead to the proposal of a defense mechanism against adversarial attacks that has some theoretical guarantees.

APPENDIX

A. Networks Used on the STL10 and ImageNet datasets

CNN Classifier on the STL10: The CNN architecture used on the STL10 dataset is shown in Table VIII. Our implementation builds on the open-source codes.²

TABLE VIII: The architecture of the CNN classifier on the STL10 dataset. The classifier consists of either convolutional (Conv) layers or fully connected (FC) layers.

STL10		Layer
Conv Layers	$\times 1$	Conv(3, 32, 3, 1) + MaxPooling2d + ReLU
	$\times 1$	Conv(32, 64, 3, 1) + MaxPooling2d + ReLU
	$\times 1$	Conv(64, 128, 3, 1) + MaxPooling2d + ReLU
	$\times 1$	Conv(128, 128, 3, 1) + MaxPooling2d + ReLU
	$\times 2$	Conv(128, 256, 3, 0) + ReLU
FC Layer	$\times 1$	MaxPooling2d + Linear(256, 10)

CNN Classifier on the ImageNet: We use the Resnet-50 [45] as the classification model on the Imagenet dataset. The pretrained version of the Resnet-50 can be found from the Torchvision repository.³

EDSR: Our RAIN uses the same EDSR model for both of the two datasets. We use the single-scale baseline version of EDSR,⁴ which consists of only 16 residual blocks (see Table 1 in [40]).

REFERENCES

- [1] A. Giusti, J. Guzzi, D. C. Cireşan, F.-L. He, J. P. Rodríguez, F. Fontana, M. Faessler, C. Forster, J. Schmidhuber, G. Di Caro *et al.*, “A machine learning approach to visual perception of forest trails for mobile robots,” *IEEE Robotics and Automation Letters*, vol. 1, no. 2, pp. 661–667, 2015.
- [2] Y. Sun, X. Wang, and X. Tang, “Deep learning face representation from predicting 10,000 classes,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 1891–1898.
- [3] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami, “Practical black-box attacks against deep learning systems using adversarial examples,” *arXiv preprint arXiv:1602.02697*, 2016.
- [4] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, “Intriguing properties of neural networks,” *arXiv preprint arXiv:1312.6199*, 2013.

²<https://github.com/aaron-xichen/pytorch-playground>

³<https://github.com/pytorch/vision/blob/master/torchvision/models>

⁴<https://github.com/thstkdgus35/EDSR-PyTorch>

- [5] B. Biggio, I. Corona, D. Maiorca, B. Nelson, N. Šrndić, P. Laskov, G. Giacinto, and F. Roli, “Evasion attacks against machine learning at test time,” in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2013.
- [6] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, “Towards deep learning models resistant to adversarial attacks,” *arXiv preprint arXiv:1706.06083*, 2017.
- [7] C. Xie, Y. Wu, L. v. d. Maaten, A. L. Yuille, and K. He, “Feature denoising for improving adversarial robustness,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 501–509.
- [8] J. Zhang, X. Xu, B. Han, G. Niu, L. Cui, M. Sugiyama, and M. Kankanhalli, “Attacks which do not kill training make adversarial learning stronger,” *arXiv preprint arXiv:2002.11242*, 2020.
- [9] C. Xie, J. Wang, Z. Zhang, Z. Ren, and A. Yuille, “Mitigating adversarial effects through randomization,” *arXiv preprint arXiv:1711.01991*, 2017.
- [10] G. K. Dziugaite, Z. Ghahramani, and D. M. Roy, “A study of the effect of jpg compression on adversarial images,” *arXiv preprint arXiv:1608.00853*, 2016.
- [11] A. Mustafa, S. H. Khan, M. Hayat, J. Shen, and L. Shao, “Image super-resolution as a defense against adversarial attacks,” *IEEE Transactions on Image Processing*, vol. 29, pp. 1711–1724, 2019.
- [12] D. Zhang, T. Zhang, Y. Lu, Z. Zhu, and B. Dong, “You only propagate once: Accelerating adversarial training via maximal principle,” in *Advances in Neural Information Processing Systems*, 2019, pp. 227–238.
- [13] M. Mosbach, M. Andriushchenko, T. Trost, M. Hein, and D. Klakow, “Logit pairing methods can fool gradient-based attacks,” *arXiv preprint arXiv:1810.12042*, 2018.
- [14] Y. Song, T. Kim, S. Nowozin, S. Ermon, and N. Kushman, “Pixeldefend: Leveraging generative models to understand and defend against adversarial examples,” *arXiv preprint arXiv:1710.10766*, 2017.
- [15] A. Prakash, N. Moran, S. Garber, A. DiLillo, and J. Storer, “Deflecting adversarial attacks with pixel deflection,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8571–8580.
- [16] P. Samangouei, M. Kabkab, and R. Chellappa, “Defense-gan: Protecting classifiers against adversarial attacks using generative models,” *arXiv preprint arXiv:1805.06605*, 2018.
- [17] C. Kou, H. K. Lee, E.-C. Chang, and T. K. Ng, “Enhancing transformation-based defenses against adversarial attacks with a distribution classifier,” in *International Conference on Learning Representations*, 2020.
- [18] H. Zhang, Y. Yu, J. Jiao, E. P. Xing, L. E. Ghaoui, and M. I. Jordan, “Theoretically principled trade-off between robustness and accuracy,” *arXiv preprint arXiv:1901.08573*, 2019.
- [19] E. Kauderer-Abrams, “Quantifying translation-invariance in convolutional neural networks,” *ArXiv*, 2018.
- [20] A. Coates, A. Ng, and H. Lee, “An analysis of single-layer networks in unsupervised feature learning,” in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, 2011, pp. 215–223.
- [21] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2009, pp. 248–255.
- [22] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” *arXiv preprint arXiv:1412.6572*, 2014.
- [23] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, “Deepfool: A simple and accurate method to fool deep neural networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2574–2582.
- [24] P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh, “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*. ACM, 2017, pp. 15–26.
- [25] A. Ilyas, L. Engstrom, A. Athalye, and J. Lin, “Black-box adversarial attacks with limited queries and information,” *arXiv preprint arXiv:1804.08598*, 2018.
- [26] C.-C. Tu, P. Ting, P.-Y. Chen, S. Liu, H. Zhang, J. Yi, C.-J. Hsieh, and S.-M. Cheng, “Autozoom: Autoencoder-based zeroth order optimization method for attacking black-box neural networks,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019, pp. 742–749.
- [27] J. Du, H. Zhang, J. T. Zhou, Y. Yang, and J. Feng, “Query-efficient meta attack to deep neural networks,” *arXiv preprint arXiv:1906.02398*, 2019.
- [28] A. Athalye, N. Carlini, and D. Wagner, “Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples,” *arXiv preprint arXiv:1802.00420*, 2018.
- [29] L. Schmidt, S. Santurkar, D. Tsipras, K. Talwar, and A. Madry, “Adversarially robust generalization requires more data,” in *Advances in Neural Information Processing Systems*, 2018, pp. 5014–5026.
- [30] A. Shafahi, M. Najibi, M. A. Ghiasi, Z. Xu, J. Dickerson, C. Studer, L. S. Davis, G. Taylor, and T. Goldstein, “Adversarial training for free!” in *Advances in Neural Information Processing Systems*, 2019, pp. 3353–3364.
- [31] L. Rice, E. Wong, and J. Z. Kolter, “Overfitting in adversarially robust deep learning,” *arXiv preprint arXiv:2002.11569*, 2020.
- [32] Y. Balaji, T. Goldstein, and J. Hoffman, “Instance adaptive adversarial training: Improved accuracy tradeoffs in neural nets,” *arXiv preprint arXiv:1910.08051*, 2019.
- [33] Y. Yang, G. Zhang, D. Katabi, and Z. Xu, “Me-net: Towards effective adversarial robustness with matrix estimation,” *arXiv preprint arXiv:1905.11971*, 2019.
- [34] Y. Xu, T. Xiao, J. Zhang, K. Yang, and Z. Zhang, “Scale-invariant convolutional neural networks,” *arXiv preprint arXiv:1411.6369*, 2014.
- [35] O. S. Kayhan and J. C. van Gemert, “On translation invariance in cnns: Convolutional layers can exploit absolute spatial location,” *ArXiv*, 2020.
- [36] R. Gens and P. M. Domingos, “Deep symmetry networks,” in *Advances in Neural Information Processing Systems*, 2014.
- [37] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu, “Spatial transformer networks,” in *Advances in Neural Information Processing Systems*, 2015.
- [38] A. Azulay and Y. Weiss, “Why do deep convolutional networks generalize so poorly to small image transformations,” *arXiv preprint arXiv:1805.12177*, 2018.
- [39] R. G. Keys, “Cubic convolution interpolation for digital image processing,” 1981.
- [40] B. Lim, S. Son, H. Kim, S. Nah, and K. Mu Lee, “Enhanced deep residual networks for single image super-resolution,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2017, pp. 136–144.
- [41] C. Dong, C. C. Loy, K. He, and X. Tang, “Image super-resolution using deep convolutional networks,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 2, pp. 295–307, 2015.
- [42] N. Carlini and D. Wagner, “Towards evaluating the robustness of neural networks,” in *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2017.
- [43] P.-Y. Chen, Y. Sharma, H. Zhang, J. Yi, and C.-J. Hsieh, “Ead: elastic-net attacks to deep neural networks via adversarial examples,” in *Thirty-second AAAI Conference on Artificial Intelligence*, 2018.
- [44] J. Chen and Q. Gu, “Rays: A ray searching method for hard-label adversarial attack,” *arXiv preprint arXiv:2006.12792*, 2020.
- [45] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.