

# STAIRCASE SIGN METHOD FOR BOOSTING ADVERSARIAL ATTACKS

Anonymous authors

Paper under double-blind review

## ABSTRACT

Crafting adversarial examples for the transfer-based attack is challenging and remains a research hot spot. Currently, such attack methods are based on the hypothesis that the substitute model and the victim’s model learn similar decision boundaries, and they conventionally apply Sign Method (SM) to manipulate the gradient as the resultant perturbation. Although SM is efficient, it only extracts the sign of gradient units but ignores their value difference, which inevitably leads to a serious deviation. Therefore, we propose a novel Staircase Sign Method ( $S^2M$ ) to alleviate this issue, thus boosting transfer-based attacks. Technically, our method heuristically divides the gradient sign into several segments according to the values of the gradient units, and then assigns each segment with a staircase weight for better crafting adversarial perturbation. As a result, our adversarial examples perform better in both white-box and black-box manner without being more visible. Since  $S^2M$  just manipulates the resultant gradient, our method can be generally integrated into any transfer-based attacks, and the computational overhead is negligible. Extensive experiments on the ImageNet dataset demonstrate the effectiveness of our proposed methods, which significantly improve the transferability (i.e., on average, **5.1%** for normally trained models and **11.2%** for adversarially trained defenses). Our code is available in the supplementary material.

## 1 INTRODUCTION

With the remarkable performance of deep neural networks (DNNs) in various tasks, the robustness of DNNs are becoming a hot spot of the current research. However, DNNs are pretty vulnerable to the adversarial examples (Szegedy et al., 2014; Goodfellow et al., 2015; Biggio et al., 2017; Zhang et al., 2020) which are only added with human-imperceptible perturbations but can fool state-of-the-art DNNs (Szegedy et al., 2016; 2017; He et al., 2016; Huang et al., 2017) successfully. To make the matter worse, attacking in the physical world (Sharif et al., 2016; Kurakin et al., 2017b; Xu et al., 2020) is also practicable, which inevitably raises concerns in real-world applications such as self-driving cars.

To better evaluate the robustness of DNNs, various works have been proposed to seek the vulnerability of DNNs. Specifically, white-box attacks such as Iterative Fast Gradient Sign Method (I-FGSM) (Kurakin et al., 2017a), Deepfool (Moosavi-Dezfooli et al., 2016) and Carlini & Wagner’s (C&W) method (Carlini & Wagner, 2017) can achieve impressive performance with the complete knowledge of the victim’s model, e.g., gradient and structure. However, deployed DNNs are usually transparent to unauthorized users for security, and thus the adversary cannot base on any knowledge of the victim’s model. Therefore, resorting to cross-model transferability (Szegedy et al., 2014; Moosavi-Dezfooli et al., 2017; Naseer et al., 2019; Gao et al., 2021; Naseer et al., 2021) of adversarial examples is a common practice. That is to say, the ad-

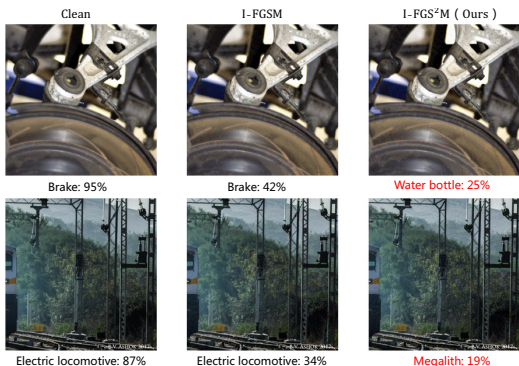


Figure 1: We generate adversarial examples by I-FGSM and our proposed I-FGS<sup>2</sup>M. The Red highlight denotes the pre-set target label.

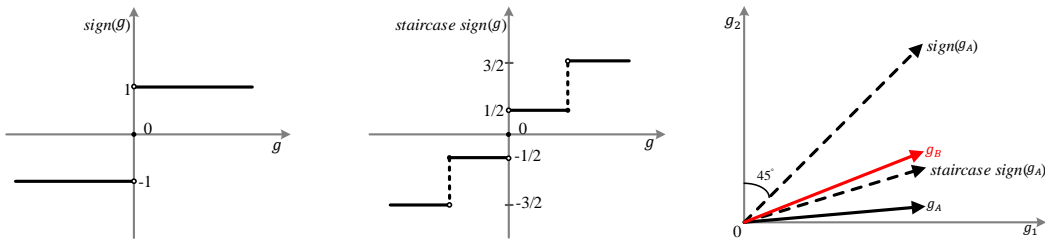


Figure 2: Sign function vs. Staircase sign function (take  $K = 2$  for example). **Left & Middle:** the graph of functions. **Right:** illustration of the gradient direction of the substitute model  $g_A$ , black-box model  $g_B$ , and the resultant update directions of sign and staircase sign function w.r.t  $g_A$ . To express more visually, here we ignore the magnitude of each perturbation. For sign function, each unit of the resultant direction is the same. By contrast, our method reflects the difference between units, and is more close to both  $g_A$  and  $g_B$ .

versarial examples crafted via known white-box models (*a.k.a* substitute model) are also dangerous for other unknown black-box models, which makes the black-box transfer-based attack possible. In this field, Goodfellow *et al.* (Goodfellow et al., 2015) hypothesize that the vulnerability of DNNs is their linear nature. Since raw gradient magnitude is extremely small (*e.g.* the minimal unit in gradient  $\approx 10^{-9}$ ) and digital images usually use 8 bits per pixel, they propose Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2015) so that each pixel can be fully perturbed with only a single step. Conventionally, the following transfer-based iterative attack methods (Dong et al., 2018; Xie et al., 2019b; Lin et al., 2020b; Gao et al., 2020a; Wu et al., 2020) are all based on Sign Method (SM) to boost adversarial attack.

However, there is a limitation in SM, *i.e.*, ignores the difference among each unit in the gradient vector. As illustrated in Figure 2(a), the update direction obtained by the sign function is that whether the partial derivative of loss function at each pixel is positive, negative or zero. Since the transferability phenomenon is mainly due to the fact that decision boundaries around the same data point of different models are similar, a naive application of SM results in a poor gradient estimation (as depicted in Figure 2(c)). Consequently, the adversarial examples especially for targeted ones may deviate from the *global optimal region* where most of DNNs can be fooled, thus decreasing the transferability.

Motivated by this, we propose a **Staircase Sign Method (S<sup>2</sup>M)** to effectively utilize the gradient for the substitute model, thus more closely approximating the gradients of both the black-box and white-box models. Technically, our proposed method first utilizes the sign function to roughly get the gradient direction for the substitute model, then heuristically assigns different weights for each pixel by our staircase sign function. In short, we merely manipulate the perturbation added on the image. Thus, S<sup>2</sup>M can be generally integrated into any transfer-based attacks, *e.g.*, the family of FGSM algorithms. Based on I-FGSM, we propose its variant I-FGS<sup>2</sup>M (Algorithm 14) which can also serve as a iterative attack baseline to be combined with the state-of-the-art approaches, *e.g.*, input diversity (Xie et al., 2019b), Poincaré space loss (Li et al., 2020), and patch-wise++ method (Gao et al., 2020b). To sum up, our main contributions are as follows:

- To the best of our knowledge, we are the first to point out the poor gradient estimation limitation of Sign Method in transfer-based attacks, which causes the adversarial examples to deviate from the *global optimal region*.
- We propose a novel Staircase Sign Method to alleviate the problem. Notably, our method is simple but effective, and can be integrated into any transfer-based attacks.
- Extensive experiments on the ImageNet dataset (Russakovsky et al., 2015) demonstrate the effectiveness of our proposed attacks which consistently outperform vanilla FGSM-based non-targeted & targeted ones in both black-box and white-box manner.

## 2 RELATED WORKS

In this section, we first briefly review the development of transfer-based attack methods in Sec. 2.1, and then introduce several defense methods in Sec. 2.2.

### 2.1 TRANSFER-BASED BLACK-BOX ATTACKS

Unlike the white-box attack, the black-box attack cannot obtain the gradient or parameters of the victim’s model. Although query-based black-box attack (Chen et al., 2017; Ilyas et al., 2018; Andriushchenko et al., 2020) can be applied in this manner, a large number of queries is computationally expensive. Thus, we turn to the transferability of adversarial examples in the remainder of this paper.

For non-targeted attacks, (Goodfellow et al., 2015) quantify the gradient by the sign function and propose single-step FGSM with the step size equal to maximum perturbation. However, perturbing images with single-step attacks usually cannot get a high success rate on the white-box model. Therefore, (Kurakin et al., 2017a) propose I-FGSM which applies FGSM multiple times with a small step size. Considering that iterative methods usually sacrifice transferability to improve the white-box performance, (Dong et al., 2018) integrate momentum term into the iterative process to avoid adversarial examples falling into local optimum. (Xie et al., 2019b) apply random transformations to the input images to alleviate the overfitting problem. (Wu et al., 2020) explore the security weakness of skip connections (He et al., 2016) to boost adversarial attacks. To effectively evade defenses, (Dong et al., 2019) propose a translation-invariant attack method to smooth the perturbation. (Lin et al., 2020a) adapt Nesterov accelerated gradient and leverage scale-invariant property of DNNs to optimize the perturbations. (Naseer et al., 2019) train cross-domain generators by using different training data distribution. (Gao et al., 2020a) craft patch-wise noise to further increase the success rate of adversarial examples.

However, targeted attacks are more challenging which need to guide the victim’s model to predict a specific target class with high confidence rather than just cause misclassification. In this setting, (Li et al., 2020) replace the cross-entropy loss with Poincaré distance and introduce triple loss to make adversarial examples close to the target label. (Gao et al., 2020b) extend non-targeted (Gao et al., 2020a) to targeted version and adopt temperature term to push the adversarial examples into the *global optimal region* of DNNs. Instead of optimizing the output distribution with a few iterations, (Zhao et al., 2020) directly maximize the target logits with more iterations. To make the adversarial examples more transferable, several researchers (Sabour et al., 2016; Inkawhich et al., 2019) turn to directly optimize intermediate features instead of output distribution. (Naseer et al., 2021) train generators to match “distribution” of the target class. Although (Inkawhich et al., 2020a;b; Naseer et al., 2021) have achieved impressive performance in this way, both training specific auxiliary models and generators for each target class is very time-consuming. Therefore, we resort to FGSM-based attacks in this paper.

### 2.2 DEFENSE METHODS

With the development of adversarial examples, researchers pay more and more attention to the robustness of DNNs, and various defense methods are proposed to circumvent potential risks. (Guo et al., 2018) apply multiple input transformation such as JPEG compression (Dziugaite et al., 2016), total variance minimization (Rudin et al., 1992) and image quilting (Efros & Freeman, 2001) to recover from the adversarial perturbations. (Theagarajan et al., 2019) introduce probabilistic adversarial robustness to neutralize adversarial attacks by concentrating sample probability to adversarial-free zones.

Although the above methods are efficient, *i.e.*, do not require a time-consuming training process, the adversarial training defense mechanism is more robust in practice. In this field, (Madry et al., 2018) adopt a natural saddle point formulation to cover the blind spots of DNNs. (Tramèr et al., 2018) introduce ensemble adversarial training which augments training data with perturbations transferred from other models. (Xie et al., 2019a) impose constraints at the feature level by denoising technique.

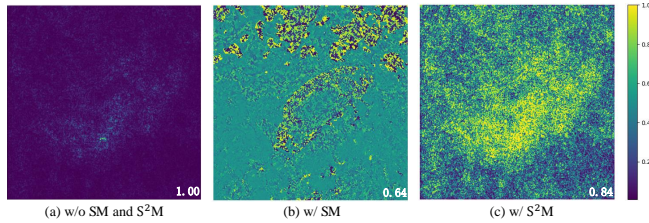


Figure 3: We visualize the perturbations of an image at the first iteration. All perturbations are normalized to  $[0, 1]$ . **(a)**: the gradient of an ensemble of Inc-v4 (Szegedy et al., 2017), IncRes-v2 (Szegedy et al., 2017) and Res-152 (He et al., 2016); **(b)&(c)**: the results of SM and  $S^2M$  w.r.t (a). The numbers in the lower right corner indicate the average cosine similarity between (a) and the other two perturbations on 1,000 images. Notably, our  $S^2M$  not only keeps the perturbation magnitude but also has higher cosine similarity to (a).

### 3 METHODOLOGY

Before introducing our algorithm in detail, we first describe the background knowledge of generating adversarial examples. Given a DNN network  $f(\mathbf{x}) : \mathbf{x} \in \mathcal{X} \rightarrow y \in \mathcal{Y}$ , it takes an input  $\mathbf{x}$  (*i.e.*, a clean image) to predict its label  $y$ . For targeted attacks<sup>1</sup>, it requires us to find a human-imperceptible perturbation  $\delta$  to satisfy  $f(\mathbf{x} + \delta) = y^*$ , where  $\mathbf{x}^* = \mathbf{x} + \delta$  is the generated adversarial example and  $y^*$  is our preset target label.

In order to make the resultant adversarial examples invisible to the clean ones, the adversary usually sets a small perturbation upper bound  $\epsilon$ , and lets  $\|\delta\|_\infty \leq \epsilon$ . By minimizing the loss function  $J(\mathbf{x}^*, y^*)$ , *e.g.*, cross entropy loss, the constrained optimization problem can be denoted as:

$$\arg \min_{\mathbf{x}^*} J(\mathbf{x}^*, y^*), \quad s.t. \quad \|\mathbf{x}^* - \mathbf{x}\|_\infty \leq \epsilon. \quad (1)$$

For targeted attacks a resultant adversarial example at iteration  $t + 1$  can be formally written as:

$$\mathbf{x}_{t+1}^* = \text{clip}_{\mathbf{x}, \epsilon} \{ \mathbf{x}_t^* - \alpha \cdot \text{sign}(\nabla_{\mathbf{x}} J(\mathbf{x}_t^*, y^*)) \}, \quad (2)$$

where  $\text{clip}_{\mathbf{x}, \epsilon}(\cdot)$  keeps the adversarial example  $\mathbf{x}^*$  within the  $\epsilon$ -ball of  $\mathbf{x}$ , and  $\alpha = \epsilon/T$  is the step size.

#### 3.1 OUR METHOD

In this section, we first introduce our motivation in Sec. 3.1.1. Then, we elaborate our method in Sec. 3.1.2 and ensemble strategy in Sec. 3.1.3.

##### 3.1.1 RETHINKING THE SIGN METHOD

To the best of our knowledge, existing transfer-based attack methods are all based on SM. In addition to the linear hypothesis (Goodfellow et al., 2015), the motivation of this method is to modify more information for each pixel than directly adding the gradient, especially for single-step attacks. Besides, manipulating the gradients by SM for iterative attacks can quickly reach the boundary of  $\ell_\infty$ -ball with only a few iterations (Dong et al., 2018; 2019).

However, the transferability of adversarial examples is mainly based on the phenomenon that decision boundaries of different models are similar. Since targeted attacks need to guide the adversarial examples into a specific territory of the target class, directly applying SM inevitably discards significant information of the gradient of the substitute model. As the toy example shown in Figure 2, the direction derived from the SM is less closer to the victim’s model than that of our  $S^2M$ . Consequently, the resultant perturbation badly deviates from the target territory, thus decreases the transferability.

<sup>1</sup>Non-targeted attacks are discussed in Appedix. C.

### 3.1.2 STAIRCASE SIGN METHOD

Motivated by the limitation of SM, we propose a novel **Staircase Sign Method (S<sup>2</sup>M)** to alleviate this problem. Figure 2 depicts the difference between sign function and our proposed staircase sign function. Since our method merely manipulates the gradient at each iteration, it can be generally integrated into any transfer-based attacks, *e.g.*, the family of FGSM algorithms. For simplicity purpose, we only take our variant I-FGS<sup>2</sup>M (summarized in Algorithm 14) as an example to demonstrate the integration process.

Technically, our method can be mainly divided into four steps. Firstly, as with the other methods, *e.g.*, (Dong et al., 2019; Xie et al., 2019b; Gao et al., 2020a), we need to compute the gradient  $\mathbf{G}_t$  at  $t$ -iteration of the substitute model with respect to the input (in line 5):

$$\mathbf{G}_t = \nabla_{\mathbf{x}} \mathbf{J}(\mathbf{x}_t^*, y^*) \quad (3)$$

Secondly, we calculate the  $p$ -th percentile  $g^p$  of  $|\mathbf{G}_t|$  (in line 7) according to the number of staircase  $K$ , where  $p$  ranges from  $100/K$  to 100 with the percentile interval  $\tau = 100/K$ . Thirdly, we assign the staircase weights  $\mathbf{W}_t$  according to  $g_t^p$  by Eq. (4) (in line 8):

$$\mathbf{W}_t = \begin{cases} \frac{\tau}{100}, & g_t^0 \leq |\mathbf{G}_t^{i,j}| \leq g_t^\tau, \\ \frac{3\tau}{100}, & g_t^\tau < |\mathbf{G}_t^{i,j}| \leq g_t^{2\tau}, \\ \vdots & \vdots \\ \frac{(2k+1)\tau}{100}, & g_t^{p-\tau} < |\mathbf{G}_t^{i,j}| \leq g_t^p, \\ \vdots & \vdots \\ \frac{(2K-1)\tau}{100}, & g_t^{100-\tau} < |\mathbf{G}_t^{i,j}| \leq g_t^{100}. \end{cases} \quad (4)$$

where  $k$  ranges from 0 to  $K - 1$ , and also equals to  $p/\tau - 1$ . As a result, our  $\mathbf{W}$  is bounded in  $[1/K, 2 - 1/K] \subseteq [0, 2]$ . Finally, combined with the sign direction of  $\mathbf{G}_t$ , we update our adversarial examples  $\mathbf{x}_t^*$  (in lines 11):

$$\mathbf{x}_{t+1}^* = \text{clip}_{\mathbf{x}, \epsilon} \{ \mathbf{x}_t^* - \alpha \cdot \text{sign}(\mathbf{G}_t) \odot \mathbf{W}_t \}, \quad (5)$$

where  $\odot$  is Hadamard product.

**Proposition 1** Assume that  $\mathbf{G}_t$  is i.i.d. for all  $t \in [0, T - 1]$ , the adversarial examples can reach the boundary of  $\epsilon$ -ball.

*Proof.* Due to the fact that  $\|\cdot\|_\infty = \max(\text{abs}(\cdot))$ , here we only discuss  $\mathbf{W}^{i,j}$  which is also i.i.d. Therefore,  $\left\| \sum_{t=0}^{T-1} \alpha \cdot \text{sign}(\mathbf{G}) \odot \mathbf{W} \right\|_\infty = \left\| \sum_{t=0}^{T-1} \alpha \cdot \mathbf{W}_t^{i,j} \right\|_\infty = \alpha \cdot \sum_{t=0}^{T-1} \mathbf{W}_t^{i,j}$ . Besides,  $\mathbb{E}(\mathbf{W}_t^{i,j}) = \sum_{k=0}^{K-1} \frac{1}{K} \cdot \frac{(2k+1)\tau}{100} = \frac{1}{100K} \cdot (\tau + 3\tau + \dots + (2K-1)\tau) = \frac{K\tau}{100} = 1$ . So we have:

$$\begin{aligned} \mathbb{E} \left( \left\| \sum_{t=0}^{T-1} \alpha \cdot \text{sign}(\mathbf{G}_t) \odot \mathbf{W}_t \right\|_\infty \right) &= \alpha \cdot \sum_{t=0}^{T-1} \mathbb{E}(\mathbf{W}_t^{i,j}) \\ &= \alpha \cdot T \\ &= \epsilon \end{aligned} \quad (6)$$

In fact, S<sup>2</sup>M is equivalent to applying adaptive weight for each pixel in sign noise. With the help of our S<sup>2</sup>M, the poor gradient estimation problem caused by SM can be effectively alleviated. As demonstrated in Figure 3, the cosine similarity between the gradients (a) and the perturbations manipulated by our proposed S<sup>2</sup>M (c) is up to **0.84**, while the result of SM is only 0.64. Please note that our S<sup>2</sup>M does not aim to make the cosine similarity close to 1.0. This is because the victim’s model is only similar to the substitute model, but it cannot be exactly the same. As demonstrated in Figure 2(c), “overfitting” on the substitute model will enlarge the gap with the victim’s model.

The adversarial examples are shown in Figure 1. Compared with I-FGSM which cannot effectively decrease the confidence of true class, our proposed variant successfully misleads the model to classify our resultant human-imperceptible adversarial examples as the pre-set target classes.

### 3.1.3 ATTACKING AN ENSEMBLE OF MODELS

To craft adversarial examples with high transferability, attacking an ensemble of models (Liu et al., 2017; Dong et al., 2018) is a pretty effective strategy, especially for black-box attacks. It is mainly



because crafting adversarial examples on multiple models has the potential to capture the global optimum of “blind spots” easily. In this paper, we follow the ensemble strategy of (Dong et al., 2018), which fuses the logits (the output before the softmax) of an ensemble of  $N$  models:

$$l(\mathbf{x}) = \sum_{n=1}^N u_n l_n(\mathbf{x}), \quad (7)$$

where  $l_n(\cdot)$  is the logits of  $n$ -th model, and  $u_n$  is its ensemble weight with  $u_n > 0$  and  $\sum_{n=1}^N u_n = 1$ .

---

**Algorithm 1: I-FGS<sup>2</sup>M**


---

**Input** : The cross-entropy loss function  $J$  of our substitute model; iterations  $T$ ;  $\ell_\infty$  constraint  $\epsilon$ ; a clean image  $\mathbf{x}$  (Normalized to  $[-1, 1]$ ) and its corresponding true label  $y$ ; the target label  $y^*$ ; the number of staircase  $K$  ( $\geq 2$ );

**Output**: The adversarial example  $\mathbf{x}^*$ ;

```

1:  $\mathbf{x}_0^* = \mathbf{x}$ ;
2:  $\alpha = \epsilon/T$ ;  $\tau = 100/K$ ;
3: Initialize staircase weights  $\mathbf{W}$  to 0, and  $p$  to  $100/K$ ;
4: for  $t \leftarrow 0$  to  $T$  do
5:    $\mathbf{G}_t = \nabla_{\mathbf{x}} J(\mathbf{x}_t^*, y^*)$ ;
6:   for  $k \leftarrow 0$  to  $K$  do
7:     calculate the  $p$ -th percentile  $g_t^p$  of  $|\mathbf{G}_t|$ ;
8:      $\mathbf{W}_t = \begin{cases} \frac{\tau}{100}, & g_t^0 < |\mathbf{G}_t^{i,j}| \leq g_t^\tau, \\ \frac{3\tau}{100}, & g_t^\tau < |\mathbf{G}_t^{i,j}| \leq g_t^{2\tau}, \\ \vdots & \vdots \\ \frac{(2k+1)\tau}{100}, & g_t^{p-\tau} < |\mathbf{G}_t^{i,j}| \leq g_t^p, \\ \vdots & \vdots \\ \frac{(2K-1)\tau}{100}, & g_t^{100-\tau} < |\mathbf{G}_t^{i,j}| \leq g_t^{100}. \end{cases}$ 
9:      $p = p + \tau$ 
10:   end for
11:    $\mathbf{x}_{t+1}^* = \text{clip}_{\mathbf{x}, \epsilon} \{ \mathbf{x}_t^* - \alpha \cdot \text{sign}(\mathbf{G}_t) \odot \mathbf{W}_t \}$ ;
12:    $\mathbf{x}_{t+1}^* = \text{clip}(\mathbf{x}_{t+1}^*, -1, 1)$ ;
13: end for
14: Return  $\mathbf{x}^* = \mathbf{x}_T^*$ ;
```

---

## 4 EXPERIMENTS

To demonstrate the effectiveness of our staircase sign mechanism, we conduct extensive experiments based on the family of FGSM methods. Firstly, we introduce the setup of experiments in Sec. 4.1. Next, we analyze the effect of staircase number in Sec. 4.2. Finally, the attack success rate for normally trained models and defense models is reported in Sec. 4.3 and Sec. 4.4, respectively. Due to the space limitation, non-targeted attacks are discussed in Appendix C. Notably, our non-targeted FGS<sup>2</sup>M variants remarkably outperform vanilla FGSM ones by 19.1% at most.

### 4.1 SETUP

**Networks:** To comprehensively compare the performance between different attack methods, we consider nine state-of-the-art models, including six normally trained models: Inception-v3 (Inc-v3) (Szegedy et al., 2016), Inception V4 (Inc-v4) (Szegedy et al., 2017), Inception-ResNet V2 (IncRes-v2) (Szegedy et al., 2017), ResNet-50 (Res-50), ResNet-101 (Res-101) and ResNet-152 (Res-152) (He et al., 2016), and three ensemble adversarial training models: Inc-v3<sub>ens3</sub>, Inc-v3<sub>ens4</sub> and IncRes-v2<sub>ens</sub> (Tramèr et al., 2018).

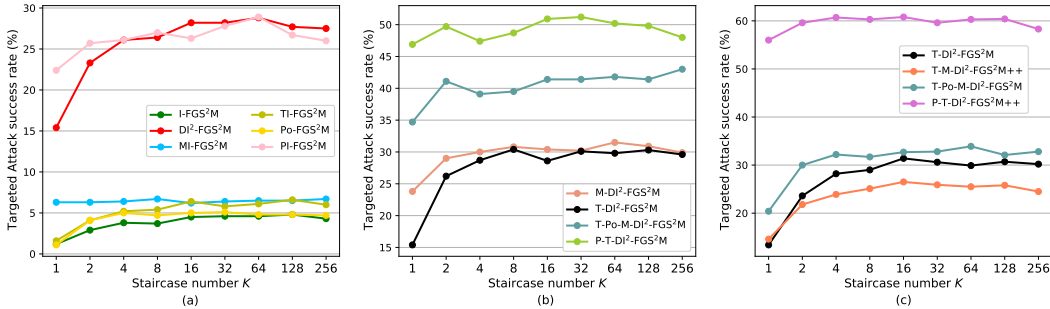


Figure 4: The success rate (%) of targeted black-box attack (Hold-out) for different methods w.r.t staircase number  $K$  ( $K = 1$  denotes SM and  $K \geq 2$  denotes  $S^2M$ ). For (a) and (b), the adversarial examples are crafted via an ensemble of Inc-v4, IncRes-v2 and Res-152, and the hold-out model is Inc-v3. For (c), the white-box models are an ensemble of Inc-v3, Inc-v4, IncRes-v2, Res-152, Res-101, Res-50, Inc-v3<sub>ens4</sub> and IncRes-v2<sub>ens</sub>, and the hold-out model is Inc-v3<sub>ens3</sub>.

**Dataset:** We conduct our experiments on ImageNet-compatible dataset<sup>2</sup>. This dataset is comprised of 1,000 images, and widely used in recent FGSM-based attacks (Dong et al., 2018; Xie et al., 2019b; Dong et al., 2019; Gao et al., 2020a; Li et al., 2020; Gao et al., 2020b). The target label for each image in this dataset is pre-set and usually different.

**Parameters:** Following the previous works (Li et al., 2020; Gao et al., 2020b), in our experiments, the maximum perturbation  $\epsilon$  is set to 16, the iteration  $T$  of all methods is 20, and thus the step size  $\alpha = \epsilon/T = 0.8$ . When attacking an ensemble of  $N$  models simultaneously, the weight for the logits of each model is equal, *i.e.*,  $1/N$ . For MI-FGSM (Dong et al., 2018), the decay factor  $\mu = 1.0$ . For DI-FGSM (Xie et al., 2019b), the transformation probability  $p = 0.7$ . For TI-FGSM (Dong et al., 2019), when the victim’s model is in normally trained models, the Gaussian kernel length is  $5 \times 5$ , and  $15 \times 15$  for defense models. For Po-FGSM (Li et al., 2020), we set  $\lambda = 0.1$ . For PI-FGSM (Gao et al., 2020a) and PI-FGSM++ (Gao et al., 2020b), the amplification factor  $\beta = 10$ , the project factor  $\gamma = 0.8\beta\alpha$ , and the project kernel length is  $3 \times 3$ . The temperature  $\tau$  for PI-FGSM++ is set to 1.5. For our  $S^2M$ , the number of staircase  $K$  is set to 64. Please note that the parameters of each method are fixed no matter what methods are combined.

**Evaluation Metrics:** There are three metrics are used to evaluate the performance of targeted attacks: “Hold-out” is the success rate of the black-box models (*i.e.* transferability), “Ensemble” denotes the white-box success rates for an ensemble of models, and “AoE” (Gao et al., 2020b) averages the white-box success rate of each model.

## 4.2 THE EFFECT OF STAIRCASE NUMBER $K$

In this section, we analyze the effect of the staircase number  $K$  for the state-of-the-art FGSM-based attacks. Here we tune  $K = 2, 4, 8, \dots, 256$ . Please note that our methods only take  $K \geq 2$  as the input.  $K = 1$  denotes their corresponding FGSM-based baseline.

We depict the curve for the black-box attacks (Hold-out) in Figure 4 and the results for white-box attacks can be found in Appendix. B. Compared with the vanilla FGSM implementation, our FGSM variants improve the transferability by a large margin as a whole. Specifically, when transferring adversarial examples to normally trained models (Figure 4(a)), DI-FGSM with  $K = 64$  sharply increases the success rate by 13.4% from 15.4% to 28.8%. Besides, as shown in Figure 4(c), our methods also boost the attack performance on defenses, *i.e.*, consistently outperform the corresponding baseline attacks by 4.3% ~ 16.5%. Considering that almost all curves turn to remain stable when  $K$  is big and reach the peak when  $K = 64$  in Figure 6 and Figure 4, we set the staircase number  $K = 64$  in the following experiments. Note that the computational overhead for this setting is almost negligible compared to the cost of forward pass and backpropagation.

<sup>2</sup>[https://github.com/tensorflow/cleverhans/tree/master/examples/nips17\\_adversarial\\_competition/dataset](https://github.com/tensorflow/cleverhans/tree/master/examples/nips17_adversarial_competition/dataset)

Table 1: The success rate (%) of targeted FGSM-based/FGS<sup>2</sup>M-based attacks. We study four models—Inc-v3, Inc-v4, Res-152 and IncRes-v2, and adversarial examples are crafted via an ensemble of three of them. In each column, “-” denote the hold-out model. We comprehensively report the results from three metrics, *i.e.*, Ensemble, AoE and Hold-out.

	Attacks	-Inc-v3	-Inc-v4	-Res-152	-IncRes	Avg.
Ensemble (White-box)	I	99.9 / <b>100.0</b>	99.9 / <b>100.0</b>	<b>100.0 / 100.0</b>	<b>100.0 / 100.0</b>	<b>100.0 / 100.0</b>
	MI	<b>99.9 / 99.9</b>	99.9 / <b>100.0</b>	<b>100.0 / 100.0</b>	<b>100.0 / 100.0</b>	<b>100.0 / 100.0</b>
	DI <sup>2</sup>	91.8 / <b>98.3</b>	93.3 / <b>98.2</b>	94.5 / <b>98.7</b>	94.8 / <b>98.5</b>	93.6 / <b>98.4</b>
	TI	<b>99.9 / 99.9</b>	99.8 / <b>99.9</b>	97.0 / <b>99.9</b>	<b>100.0 / 100.0</b>	99.2 / <b>99.9</b>
	Po	<b>100.0 / 100.0</b>	99.9 / <b>100.0</b>	<b>100.0 / 100.0</b>	<b>100.0 / 100.0</b>	<b>100.0 / 100.0</b>
	PI	<b>100.0 / 99.4</b>	<b>100.0 / 99.9</b>	<b>100.0 / 99.9</b>	<b>99.8 / 99.8</b>	<b>100.0 / 99.8</b>
	M-DI <sup>2</sup>	88.9 / <b>95.7</b>	90.4 / <b>96.7</b>	91.0 / <b>96.2</b>	92.8 / <b>98.3</b>	90.8 / <b>96.7</b>
	T-DI <sup>2</sup>	92.0 / <b>97.7</b>	92.2 / <b>97.7</b>	93.2 / <b>98.2</b>	94.0 / <b>98.5</b>	92.9 / <b>98.0</b>
	T-Po-M-DI <sup>2</sup>	91.4 / <b>97.4</b>	93.0 / <b>97.4</b>	91.9 / <b>96.5</b>	94.9 / <b>97.5</b>	92.8 / <b>97.2</b>
	P-T-DI <sup>2</sup>	<b>99.3 / 98.8</b>	<b>99.4 / 98.5</b>	<b>99.4 / 99.0</b>	<b>99.6 / 99.4</b>	<b>99.4 / 98.9</b>
AoE (White-box)	I	94.7 / <b>97.2</b>	88.6 / <b>93.3</b>	92.5 / <b>97.0</b>	89.7 / <b>93.1</b>	91.4 / <b>95.2</b>
	MI	94.5 / <b>96.6</b>	90.1 / <b>93.5</b>	93.4 / <b>97.0</b>	90.5 / <b>93.2</b>	92.1 / <b>95.1</b>
	DI <sup>2</sup>	77.8 / <b>89.1</b>	76.3 / <b>86.8</b>	84.4 / <b>93.7</b>	77.8 / <b>86.0</b>	79.1 / <b>88.9</b>
	TI	94.0 / <b>97.0</b>	87.2 / <b>92.3</b>	92.5 / <b>96.6</b>	88.6 / <b>92.3</b>	90.6 / <b>94.6</b>
	Po	88.7 / <b>92.8</b>	82.4 / <b>88.5</b>	78.6 / <b>85.9</b>	87.1 / <b>91.6</b>	84.2 / <b>89.7</b>
	PI	<b>98.1 / 97.9</b>	<b>97.3 / 97.2</b>	<b>98.5 / 97.8</b>	<b>96.6 / 96.9</b>	<b>97.6 / 97.5</b>
	M-DI <sup>2</sup>	75.4 / <b>84.8</b>	74.6 / <b>83.7</b>	80.9 / <b>89.8</b>	76.6 / <b>85.0</b>	76.9 / <b>85.8</b>
	T-DI <sup>2</sup>	78.6 / <b>89.1</b>	75.9 / <b>87.0</b>	83.8 / <b>93.6</b>	76.8 / <b>86.6</b>	78.8 / <b>89.1</b>
	T-Po-M-DI <sup>2</sup>	79.4 / <b>86.3</b>	76.4 / <b>83.8</b>	77.0 / <b>84.6</b>	78.0 / <b>84.6</b>	77.7 / <b>84.8</b>
	P-T-DI <sup>2</sup>	<b>95.1 / 94.0</b>	<b>93.3 / 92.7</b>	<b>97.3 / 96.7</b>	92.7 / <b>93.5</b>	<b>94.6 / 94.2</b>
Hold-out (Black-box)	I	1.2 / <b>4.6</b>	1.2 / <b>3.4</b>	0.0 / <b>1.1</b>	0.9 / <b>1.9</b>	0.8 / <b>2.8</b>
	MI	6.3 / <b>6.5</b>	3.6 / <b>3.7</b>	<b>1.6 / 1.4</b>	3.0 / <b>3.6</b>	3.6 / <b>3.8</b>
	DI <sup>2</sup>	15.4 / <b>28.8</b>	13.8 / <b>27.6</b>	3.2 / <b>8.6</b>	9.4 / <b>20.6</b>	10.5 / <b>21.4</b>
	TI	1.6 / <b>6.1</b>	1.3 / <b>4.2</b>	0.3 / <b>1.3</b>	0.8 / <b>3.2</b>	1.0 / <b>3.7</b>
	Po	1.1 / <b>4.8</b>	0.9 / <b>2.9</b>	0.0 / <b>0.4</b>	0.3 / <b>2.3</b>	0.6 / <b>2.6</b>
	PI	22.4 / <b>28.9</b>	17.2 / <b>23.2</b>	4.2 / <b>6.2</b>	13.9 / <b>20.9</b>	14.4 / <b>19.8</b>
	M-DI <sup>2</sup>	23.8 / <b>31.5</b>	24.1 / <b>30.8</b>	12.3 / <b>14.2</b>	21.3 / <b>28.3</b>	20.4 / <b>26.2</b>
	T-DI <sup>2</sup>	15.5 / <b>29.8</b>	15.9 / <b>30.6</b>	3.9 / <b>9.7</b>	11.3 / <b>25.1</b>	11.6 / <b>23.8</b>
	T-Po-M-DI <sup>2</sup>	34.7 / <b>41.8</b>	32.3 / <b>40.4</b>	17.3 / <b>18.0</b>	28.3 / <b>34.4</b>	28.2 / <b>33.7</b>
	P-T-DI <sup>2</sup>	46.9 / <b>50.2</b>	47.1 / <b>50.8</b>	14.2 / <b>19.4</b>	41.3 / <b>44.7</b>	37.4 / <b>41.3</b>

### 4.3 ATTACKING NORMALLY TRAINED MODELS

In this section, we compare ten FGSM-based attacks including I-FGSM, MI-FGSM, DI<sup>2</sup>-FGSM, TI-FGSM, Po-FGSM, M-DI<sup>2</sup>-FGSM, T-DI<sup>2</sup>-FGSM, T-Po-M-DI<sup>2</sup>-FGSM, P-T-DI<sup>2</sup>-FGSM with our FGS<sup>2</sup>M variants. In this experiment, four models including Inc-v3, Inc-v4, Res-152, and IncRes-v2 are considered. At each iteration, we select one model as the hold-out model to test the transferability, and an ensemble of the rest three with the weight of each model 1/3 serves as the substitute model.

As indicated in Table 1, our proposed FGS<sup>2</sup>M variants effectively boost both the white-box and black-box attacks. On average, they increase the success rate in Ensemble, AoE and Hold-out cases by **2.1%**, **5.2%** and **5.1%**, respectively. This demonstrates that our adversarial examples are more close to the *global optimal region*.

From the results of Table 1, we also observe that several methods, especially for these integrated with diversity input patterns (DI<sup>2</sup>), suffer badly from SM which cannot well utilize the gradient with respect to the random input transformation. Specifically, DI<sup>2</sup>-FGSM only successfully attack 93.6% images against the substitute model (Ensemble). The average success rate of each white-box model (AoE) is even reduced to 79.1%, and merely 10.5% images transfer to the black-box model (Hold-out) on average. With the help of our S<sup>2</sup>M at each iteration, we dramatically alleviate the poor gradient estimation problem, that is, increasing the success rate in Ensemble and AoE cases by **4.8%** and **9.8%**, respectively. Furthermore, in the Hold-out case our DI<sup>2</sup>-FGS<sup>2</sup>M remarkably improves the transferability by **10.9%**.

Another observation from the results is that the staircase sign perturbation seems to be less effective on vanilla MI but more effective for other momentum-based methods in the black-box manner, *e.g.*,



M-DI<sup>2</sup>. It may be because that the update direction derived by our FGS<sup>2</sup>M variants suffers more from the noise curing (Li et al., 2020) caused by vanilla MI, thus causing a lack of diversity and adaptability of noise.

Table 2: The success rate (%) of targeted FGSM-based/FGS<sup>2</sup>M-based attacks. We study nine models—Inc-v3, Inc-v4, Res-152, Res-101, Res-50, IncRes-v2, Inc-v3<sub>ens3</sub>, Inc-v3<sub>ens4</sub> and IncRes-v2<sub>ens</sub>, and adversarial examples are crafted via an ensemble of eight of them. In each column, “-” denote the hold-out model. We comprehensively report the results from three metrics, *i.e.*, Ensemble, AoE and Hold-out.

Attacks		-Inc-v3 <sub>ens3</sub>	-Inc-v3 <sub>ens4</sub>	-IncRes-v2 <sub>ens</sub>	Avg.
Ensemble (White-box)	T-DI <sup>2</sup>	56.5 / <b>78.1</b>	56.5 / <b>77.8</b>	55.2 / <b>76.7</b>	56.1 / <b>77.5</b>
	T-M-DI <sup>2</sup>	44.8 / <b>64.7</b>	45.4 / <b>64.9</b>	48.9 / <b>68.4</b>	46.4 / <b>66.0</b>
	T-Po-M-DI <sup>2</sup>	58.9 / <b>74.3</b>	58.3 / <b>75.3</b>	60.6 / <b>76.9</b>	59.3 / <b>75.5</b>
	P-T-DI <sup>2</sup> ++	<b>94.1</b> / 93.7	<b>94.3</b> / 93.3	94.2 / <b>95.0</b>	<b>94.2</b> / 94.0
AoE (White-box)	T-DI <sup>2</sup>	43.4 / <b>65.1</b>	44.6 / <b>65.5</b>	46.2 / <b>68.1</b>	44.7 / <b>66.2</b>
	T-M-DI <sup>2</sup>	34.7 / <b>52.1</b>	36.2 / <b>52.4</b>	37.8 / <b>54.5</b>	36.2 / <b>53.0</b>
	T-Po-M-DI <sup>2</sup>	48.5 / <b>63.8</b>	48.9 / <b>63.8</b>	50.1 / <b>65.5</b>	49.2 / <b>64.4</b>
	P-T-DI <sup>2</sup> ++	<b>87.5</b> / 87.0	<b>87.3</b> / 86.7	<b>88.4</b> / 87.5	<b>87.7</b> / 87.1
Hold-out (Black-box)	T-DI <sup>2</sup>	13.4 / <b>29.9</b>	12.6 / <b>30.0</b>	10.4 / <b>29.0</b>	12.1 / <b>29.6</b>
	T-M-DI <sup>2</sup>	14.6 / <b>25.5</b>	14.5 / <b>25.2</b>	14.2 / <b>24.3</b>	14.4 / <b>25.0</b>
	T-Po-M-DI <sup>2</sup>	20.4 / <b>33.9</b>	20.0 / <b>32.5</b>	19.2 / <b>30.8</b>	19.9 / <b>32.4</b>
	P-T-DI <sup>2</sup> ++	56.0 / <b>60.3</b>	56.5 / <b>58.1</b>	45.5 / <b>51.7</b>	52.7 / <b>56.7</b>

#### 4.4 ATTACKING ADVERSARIALLY TRAINED DEFENSES

Adversarially trained defenses are shown to effectively withstand the adversarial examples. Therefore, in this experiment, we consider all nine models, which are introduced in Sec. 4.1. To conduct the experiment, we select one of the ensemble adversarial training models (*e.g.* Inc-v3<sub>ens4</sub>) as the hold-out model to evaluate the transferability, and then craft adversarial examples with an ensemble of the rest eight models with the weight of each model 1/8. Here we compare four stronger FGSM-based attacks including T-DI<sup>2</sup>-FGSM, T-M-DI<sup>2</sup>-FGSM, T-Po-M-DI<sup>2</sup>-FGSM, P-T-DI<sup>2</sup>-FGSM++ with our FGS<sup>2</sup>M variants.

As demonstrated in Table 2, regardless of the attacks are white-box or black-box, our methods generally surpass the vanilla FGSM-based methods. Specifically, in the white-box manner, FGS<sup>2</sup>M-based attacks, on average, outperform FGSM-based ones by **14.3%** (Ensemble) and **13.2%** (AoE). This again demonstrates that our method can effectively alleviate the poor gradient estimation problem caused by SM. Besides, even under the more challenging black-box attack manner, our proposed attacks, on average, significantly improves the transferability by an increase of **11.2%** (Hold-out). Notably, compared with T-DI<sup>2</sup>-FGSM, which only successfully transfers 10.4% adversarial examples to IncRes-v2<sub>ens</sub>, our T-DI<sup>2</sup>-FGS<sup>2</sup>M achieves an higher transferability, *i.e.*, **29.0%**, which is about **3×**.

## 5 CONCLUSION

In this paper, we rethink the limitation of Sign Method (SM) applied by state-of-the-art transfer-based attacks and experimentally demonstrate that it causes poor gradient estimation. To address this issue, we propose a simple but effective Staircase Sign Method (S<sup>2</sup>M) to boost transferability. With the help of staircase weights, our methods effectively fool both white-box models and black-box models. Extensive experiments on the ImageNet dataset demonstrate the effectiveness of our FGS<sup>2</sup>M-based attacks, which significantly improves the transferability by **5.1%** for normally trained models and **11.2%** for adversarially trained defenses on average. Therefore, we hope our method can serve as an effective baseline to boost adversarial attacks and evaluate the robustness of various deep neural networks.

## REFERENCES

- Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: A query-efficient black-box adversarial attack via random search. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm (eds.), *ECCV*, 2020.
- Battista Biggio, Iginio Corona, Davide Maiorca, Blaine Nelson, Nedim Srndic, Pavel Laskov, Giorgio Giacinto, and Fabio Roli. Evasion attacks against machine learning at test time. *CoRR*, abs/1708.06131, 2017.
- Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. In *Symposium on Security and Privacy*, 2017.
- Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. ZOO: zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In Bhavani M. Thuraisingham, Battista Biggio, David Mandell Freeman, Brad Miller, and Arunesh Sinha (eds.), *AISec@CCS*, 2017.
- Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting adversarial attacks with momentum. In *CVPR*, 2018.
- Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. Evading defenses to transferable adversarial examples by translation-invariant attacks. In *CVPR*, 2019.
- Dziugaite, Gintare Karolina, Zoubin Ghahramani, and Daniel M. Roy. A study of the effect of jpg compression on adversarial images. *CoRR*, abs/1608.00853, 2016.
- Alexei A. Efros and William T. Freeman. Image quilting for texture synthesis and transfer. In *SIGGRAPH*, 2001.
- Lianli Gao, Qilong Zhang, Jingkuan Song, Xianglong Liu, and Hengtao Shen. Patch-wise attack for fooling deep neural network. In *ECCV*, 2020a.
- Lianli Gao, Qilong Zhang, Jingkuan Song, and Heng Tao Shen. Patch-wise++ perturbation for adversarial targeted attacks. *CoRR*, abs/2012.15503, 2020b.
- Lianli Gao, Yaya Cheng, Qilong Zhang, Xing Xu, and Jingkuan Song. Feature space targeted attacks by statistic alignment. In *IJCAI*, 2021.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In Yoshua Bengio and Yann LeCun (eds.), *ICLR*, 2015.
- Chuan Guo, Mayank Rana, Moustapha Cissé, and Laurens van der Maaten. Countering adversarial images using input transformations. In *ICLR*, 2018.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *CVPR*, 2017.
- Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks with limited queries and information. In Jennifer G. Dy and Andreas Krause (eds.), *ICML*, 2018.
- Nathan Inkawhich, Wei Wen, Hai (Helen) Li, and Yiran Chen. Feature space perturbations yield more transferable adversarial examples. In *CVPR*, 2019.
- Nathan Inkawhich, Kevin J. Liang, Lawrence Carin, and Yiran Chen. Transferable perturbations of deep feature distributions. In *ICLR*. OpenReview.net, 2020a.
- Nathan Inkawhich, Kevin J. Liang, Binghui Wang, Matthew Inkawhich, Lawrence Carin, and Yiran Chen. Perturbing across the feature hierarchy to improve standard and strict blackbox attack transferability. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *NeurIPS*, 2020b.

- Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial machine learning at scale. In *ICLR*, 2017a.
- Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In *ICLR*, 2017b.
- Maosen Li, Cheng Deng, Tengjiao Li, Junchi Yan, Xinbo Gao, and Heng Huang. Towards transferable targeted attack. In *CVPR*, 2020.
- Jiadong Lin, Chuanbiao Song, Kun He, Liwei Wang, and John E. Hopcroft. Nesterov accelerated gradient and scale invariance for adversarial attacks. In *ICLR*, 2020a.
- Jiadong Lin, Chuanbiao Song, Kun He, Liwei Wang, and John E. Hopcroft. Nesterov accelerated gradient and scale invariance for adversarial attacks. In *ICLR*, 2020b.
- Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and black-box attacks. In *ICLR*, 2017.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018.
- Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: A simple and accurate method to fool deep neural networks. In *CVPR*, 2016.
- Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. In *CVPR*, 2017.
- Muzammal Naseer, Salman H. Khan, Muhammad Haris Khan, Fahad Shahbaz Khan, and Fatih Porikli. Cross-domain transferability of adversarial perturbations. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *NeurIPS*, 2019.
- Muzammal Naseer, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Fatih Porikli. On generating transferable targeted perturbations. *CoRR*, abs/2103.14641, 2021.
- Leonid I Rudin, Stanley Osher, and Emad Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: nonlinear phenomena*, 60(1-4):259–268, 1992.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Fei-Fei Li. Imagenet large scale visual recognition challenge. *IJCV*, 2015.
- Sara Sabour, Yanshuai Cao, Fartash Faghri, and David J. Fleet. Adversarial manipulation of deep representations. In *ICLR*, 2016.
- Mahmood Sharif, Sruti Bhagavatula, Lujo Bauer, and Michael K. Reiter. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *SIGSAC*, 2016.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In Yoshua Bengio and Yann LeCun (eds.), *ICLR*, 2014.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Re-thinking the inception architecture for computer vision. In *CVPR*, 2016.
- Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *AAAI*, 2017.
- Rajkumar Theagarajan, Ming Chen, Bir Bhanu, and Jing Zhang. Shieldnets: Defending against adversarial attacks using probabilistic adversarial robustness. In *CVPR*, 2019.
- Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian J. Goodfellow, Dan Boneh, and Patrick D. McDaniel. Ensemble adversarial training: attacks and defenses. In *ICLR*, 2018.

- Dongxian Wu, Yisen Wang, Shu-Tao Xia, James Bailey, and Xingjun Ma. Skip connections matter: On the transferability of adversarial examples generated with resnets. In *ICLR*, 2020.
- Cihang Xie, Yuxin Wu, Laurens van der Maaten, Alan L. Yuille, and Kaiming He. Feature denoising for improving adversarial robustness. In *CVPR*, 2019a.
- Cihang Xie, Zhishuai Zhang, Yuyin Zhou, Song Bai, Jianyu Wang, Zhou Ren, and Alan L. Yuille. Improving transferability of adversarial examples with input diversity. In *CVPR*, 2019b.
- Kaidi Xu, Gaoyuan Zhang, Sijia Liu, Quanfu Fan, Mengshu Sun, Hongge Chen, Pin-Yu Chen, Yanzhi Wang, and Xue Lin. Adversarial t-shirt! evading person detectors in a physical world. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm (eds.), *ECCV*, 2020.
- Chaoning Zhang, Philipp Benz, Tooba Imtiaz, and In So Kweon. Understanding adversarial examples from the mutual influence of images and perturbations. In *CVPR*, 2020.
- Zhengyu Zhao, Zhuoran Liu, and Martha A. Larson. On success and simplicity: A second look at transferable targeted attacks. *CoRR*, abs/2012.11207, 2020.

## A APPENDIX

### A.1 SETUP

**Parameters:** For targeted attacks, we adopt the same parameter setting in our paper. For non-targeted attacks, here we following the previous works Gao et al. (2020a); Lin et al. (2020b). In our experiments, the maximum perturbation  $\epsilon$  is set to 16, the iteration  $T$  of all methods is 10, and thus the step size  $\alpha = \epsilon/T = 1.6$ . For MI-FGSM, the decay factor  $\mu = 1.0$ . For DI-FGSM, the transformation probability  $p = 0.7$ . For TI-FGSM, the Gaussian kernel length is  $15 \times 15$ . For PI-FGSM, the amplification factor  $\beta = 5$ , the project factor  $\gamma = 1.0$ , and the project kernel length is  $3 \times 3$ . For SI-FGSM, the number of scale copies  $m = 5$ . For our S<sup>2</sup>M, the number of staircase  $K$  is set to 64. Please note that the parameters of each method are fixed no matter what methods are combined.

## B EXPERIMENTS FOR TARGETED ATTACKS

### B.1 THE EFFECT OF STAIRCASE NUMBER $K$

Here we show the experimental results of white-box attacks, *i.e.*, Figure 5 for AoE and Figure 6 for Ensemble. In this section, we analyze the effect of the staircase number  $K$  for state-of-the-art FGSM-based attacks. Here we tune  $K = 2, 4, 8, \dots, 256$ . Similar to the observation in Sec. 4.2, our FGS<sup>2</sup>M variants can also improve the success rates by a large margin even when  $K = 2$  and the success rate continues to increase and then remain stable after  $K$  exceeds 64.

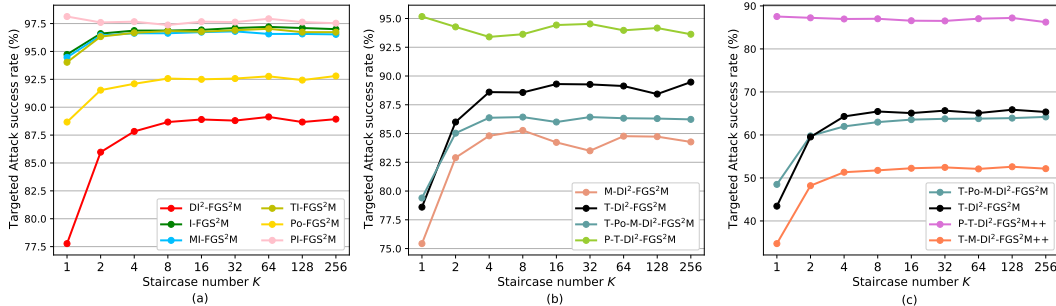


Figure 5: The success rate (%) of targeted white-box attack (AoE) for different methods w.r.t staircase number  $K$  ( $K = 1$  denotes SM and  $K \geq 2$  denotes S<sup>2</sup>M). For (a) and (b), the adversarial examples are crafted via an ensemble of Inc-v4, IncRes-v2 and Res-152, and the hold-out model is Inc-v3. For (c), the white-box models are an ensemble of Inc-v3, Inc-v4, IncRes-v2, Res-152, Res-101, Res-50, Inc-v3<sub>ens4</sub> and IncRes-v2<sub>ens</sub>, and the hold-out model is Inc-v3<sub>ens3</sub>.

## C EXPERIMENTS FOR NON-TARGETED ATTACKS

Due to space limitation, we mainly discuss the more challenging targeted attacks in our paper. Since the poor gradient estimation problem is caused by SM, crafting adversarial perturbations by our proposed S<sup>2</sup>M can also boost non-targeted attacks. In this section, we report our experimental results to demonstrate the effectiveness of our methods.

### C.1 STAIRCASE SIGN METHOD FOR NON-TARGETED ATTACKS

In this section, we compare six FGSM-based attacks including I-FGSM, MI-FGSM, DI<sup>2</sup>-FGSM, TI-FGSM, SI-FGSM and PI-FGSM with our FGS<sup>2</sup>M variants. In this experiment, we study nine models introduced in Sec. A.1. Since non-targeted attacks are less challenging than targeted attacks, here we craft adversarial examples via one model instead of an ensemble of models. More Specifically, we consider four well-known normally trained models including Inc-v3, Inc-v4, IncRes-v2 and Res-152 as our substitute model.



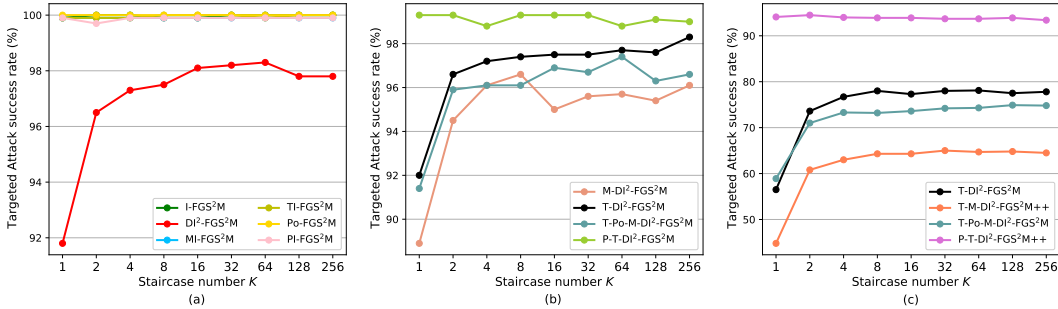


Figure 6: The success rate (%) of targeted white-box attack (Ensemble) for different methods w.r.t staircase number  $K$  ( $K = 1$  denotes SM and  $K \geq 2$  denotes  $S^2M$ ). For (a) and (b), the adversarial examples are crafted via an ensemble of Inc-v4, IncRes-v2 and Res-152. For (c), the white-box models are an ensemble of Inc-v3, Inc-v4, IncRes-v2, Res-152, Res-101, Res-50, Inc-v3<sub>ens4</sub> and IncRes-v2<sub>ens</sub>.

As demonstrated in Table 3, Table 4, Table 5 and Table 6, our FGS<sup>2</sup>M-based attacks consistently outperform FGSM-based ones in both the white-box and black-box manner. For the black-box manner, we improve the significantly transferability by **8.2%** on average. Furthermore, when adversarial examples are crafted via IncRes-v2 by SI-FGS<sup>2</sup>M, we can remarkably transfer an extra **19.1%** adversarial examples to Inc-v3<sub>ens3</sub>. For the white-box manner, we observe that our FGS<sup>2</sup>M variants further increase the success rate of white-box attack toward 100%. As demonstrated in Table 3, I-FGSM only successfully attacks Inc-v3 with a 99.2% success rate. But with the help of our staircase weights, our I-FGS<sup>2</sup>M effectively fool Inc-v3 on all the images, *i.e.*, **100%**.

Attacks	Inc-v3*	Inc-v4	IncRes-v2	Res-152	Res-101	Res-50	Inc-v3 <sub>ens3</sub>	Inc-v3 <sub>ens4</sub>	IncRes-v2 <sub>ens</sub>	Avg.
I	99.2 / <b>100.0</b>	30.0 / <b>40.6</b>	21.5 / <b>34.9</b>	18.9 / <b>26.6</b>	20.8 / <b>29.6</b>	23.3 / <b>32.6</b>	12.1 / <b>16.0</b>	12.1 / <b>17.3</b>	4.9 / <b>8.4</b>	18.0 / <b>25.8</b>
MI	99.2 / <b>100.0</b>	55.7 / <b>57.5</b>	51.2 / <b>55.1</b>	<b>44.0 / 44.0</b>	44.6 / <b>46.9</b>	49.9 / <b>51.8</b>	21.9 / <b>22.7</b>	20.7 / <b>23.1</b>	11.0 / <b>11.6</b>	37.4 / <b>39.1</b>
DI	99.9 / <b>100.0</b>	52.5 / <b>67.0</b>	42.5 / <b>57.8</b>	32.4 / <b>42.9</b>	36.0 / <b>48.4</b>	41.4 / <b>52.4</b>	13.9 / <b>22.0</b>	14.6 / <b>22.7</b>	6.9 / <b>11.6</b>	30.0 / <b>40.6</b>
TI	99.1 / <b>100.0</b>	27.5 / <b>35.9</b>	14.0 / <b>24.5</b>	16.3 / <b>23.0</b>	17.9 / <b>25.6</b>	22.1 / <b>28.1</b>	17.8 / <b>28.0</b>	16.5 / <b>26.6</b>	10.4 / <b>16.9</b>	17.8 / <b>26.1</b>
SI	<b>100.0 / 100.0</b>	53.8 / <b>69.3</b>	47.2 / <b>64.9</b>	39.4 / <b>53.5</b>	45.3 / <b>58.9</b>	48.7 / <b>61.3</b>	21.7 / <b>33.4</b>	22.6 / <b>37.2</b>	10.8 / <b>20.1</b>	36.2 / <b>49.8</b>
PI	<b>100.0 / 100.0</b>	54.5 / <b>62.9</b>	47.4 / <b>55.9</b>	39.7 / <b>47.6</b>	43.0 / <b>48.7</b>	48.2 / <b>52.1</b>	26.3 / <b>31.3</b>	25.4 / <b>29.0</b>	15.5 / <b>19.2</b>	37.5 / <b>43.3</b>

Table 3: The success rate (%) of non-targeted FGSM-based/FGS<sup>2</sup>M-based attacks w.r.t adversarial examples crafted via Inc-v3. We study nine models—Inc-v3, Inc-v4, IncRes-v2, Res-152, Res-101, Res-50, Inc-v3<sub>ens3</sub>, Inc-v3<sub>ens4</sub> and IncRes-v2<sub>ens</sub> here.

Attacks	Inc-v3	Inc-v4*	IncRes-v2	Res-152	Res-101	Res-50	Inc-v3 <sub>ens3</sub>	Inc-v3 <sub>ens4</sub>	IncRes-v2 <sub>ens</sub>	Avg.
I	43.2 / <b>56.9</b>	99.2 / <b>100.0</b>	26.3 / <b>39.2</b>	25.2 / <b>34.6</b>	25.9 / <b>36.1</b>	30.9 / <b>39.8</b>	12.0 / <b>16.9</b>	12.6 / <b>19.0</b>	6.4 / <b>10.5</b>	22.8 / <b>31.6</b>
MI	70.8 / <b>73.1</b>	99.2 / <b>100.0</b>	58.1 / <b>60.4</b>	52.2 / <b>53.5</b>	53.8 / <b>54.8</b>	56.6 / <b>59.1</b>	23.8 / <b>26.0</b>	23.8 / <b>25.0</b>	12.7 / <b>13.9</b>	44.0 / <b>45.7</b>
DI	64.5 / <b>75.3</b>	99.1 / <b>100.0</b>	48.3 / <b>63.5</b>	38.6 / <b>48.5</b>	40.0 / <b>50.4</b>	44.3 / <b>54.5</b>	16.0 / <b>21.5</b>	16.3 / <b>22.9</b>	8.6 / <b>13.7</b>	34.6 / <b>43.8</b>
TI	36.8 / <b>46.6</b>	99.2 / <b>100.0</b>	16.8 / <b>27.9</b>	20.8 / <b>27.9</b>	18.9 / <b>26.9</b>	22.3 / <b>31.9</b>	16.2 / <b>26.8</b>	19.8 / <b>27.7</b>	11.6 / <b>17.9</b>	20.4 / <b>29.2</b>
SI	72.0 / <b>81.7</b>	<b>100.0 / 100.0</b>	57.0 / <b>71.7</b>	51.1 / <b>62.3</b>	52.1 / <b>64.2</b>	56.7 / <b>68.6</b>	26.6 / <b>43.7</b>	28.0 / <b>45.4</b>	16.9 / <b>29.8</b>	45.1 / <b>58.4</b>
PI	68.7 / <b>74.9</b>	<b>100.0 / 100.0</b>	51.4 / <b>60.2</b>	45.4 / <b>53.3</b>	44.5 / <b>52.8</b>	52.2 / <b>57.7</b>	28.2 / <b>35.4</b>	27.8 / <b>33.6</b>	19.7 / <b>23.6</b>	42.2 / <b>48.9</b>

Table 4: The success rate (%) of non-targeted FGSM-based/FGS<sup>2</sup>M-based attacks w.r.t adversarial examples crafted via Inc-v4. We study nine models—Inc-v3, Inc-v4, IncRes-v2, Res-152, Res-101, Res-50, Inc-v3<sub>ens3</sub>, Inc-v3<sub>ens4</sub> and IncRes-v2<sub>ens</sub> here.

Attacks	Inc-v3	Inc-v4	IncRes-v2*	Res-152	Res-101	Res-50	Inc-v3 <sub>ens3</sub>	Inc-v3 <sub>ens4</sub>	IncRes-v2 <sub>ens</sub>	Avg.	
IncRes-v2	I	46.7 / <b>59.5</b>	38.2 / <b>49.0</b>	99.2 / <b>100.0</b>	25.4 / <b>36.5</b>	28.2 / <b>39.9</b>	30.7 / <b>42.1</b>	13.2 / <b>21.4</b>	13.0 / <b>19.5</b>	8.3 / <b>15.2</b>	25.5 / <b>35.4</b>
	MI	<b>76.1</b> / 75.8	67.9 / <b>68.8</b>	99.2 / <b>100.0</b>	<b>57.6</b> / 56.3	57.9 / <b>58.9</b>	61.3 / <b>63.3</b>	32.0 / <b>34.7</b>	28.4 / <b>28.8</b>	20.5 / <b>22.0</b>	50.2 / <b>51.1</b>
	DI	71.4 / <b>79.5</b>	65.3 / <b>76.6</b>	98.5 / <b>99.7</b>	47.8 / <b>58.3</b>	49.6 / <b>59.8</b>	54.38 / <b>64.7</b>	19.5 / <b>31.0</b>	19.1 / <b>28.1</b>	12.2 / <b>22.6</b>	42.5 / <b>52.6</b>
	TI	43.5 / <b>52.2</b>	41.2 / <b>47.8</b>	98.8 / <b>99.9</b>	26.3 / <b>31.4</b>	28.6 / <b>34.5</b>	30.2 / <b>38.0</b>	26.7 / <b>38.6</b>	24.7 / <b>36.4</b>	20.8 / <b>32.4</b>	30.3 / <b>38.5</b>
	SI	74.0 / <b>83.9</b>	64.2 / <b>75.7</b>	<b>99.9</b> / <b>99.9</b>	51.7 / <b>66.1</b>	52.9 / <b>66.5</b>	60.5 / <b>73.3</b>	29.6 / <b>48.7</b>	28.8 / <b>44.3</b>	22.1 / <b>40.7</b>	48.0 / <b>62.4</b>
	PI	72.6 / <b>79.0</b>	64.1 / <b>72.8</b>	<b>100.0</b> / <b>100.0</b>	52.2 / <b>59.0</b>	53.8 / <b>61.2</b>	56.9 / <b>63.9</b>	34.3 / <b>43.4</b>	30.9 / <b>38.5</b>	25.6 / <b>33.2</b>	48.8 / <b>56.4</b>

Table 5: The success rate (%) of non-targeted FGSM-based/FGS<sup>2</sup>M-based attacks w.r.t adversarial examples crafted via IncRes-v2. We study nine models—Inc-v3, Inc-v4, IncRes-v2, Res-152, Res-101, Res-50, Inc-v3<sub>ens3</sub>, Inc-v3<sub>ens4</sub> and IncRes-v2<sub>ens</sub> here.

Attacks	Inc-v3	Inc-v4	IncRes-v2	Res-152*	Res-101	Res-50	Inc-v3 <sub>ens3</sub>	Inc-v3 <sub>ens4</sub>	IncRes-v2 <sub>ens</sub>	Avg.	
IncRes-v2	I	31.3 / <b>43.8</b>	25.9 / <b>35.6</b>	17.7 / <b>31.6</b>	98.7 / <b>99.5</b>	67.3 / <b>80.8</b>	66.1 / <b>78.0</b>	12.2 / <b>17.7</b>	13.3 / <b>19.4</b>	7.6 / <b>12.6</b>	30.2 / <b>39.9</b>
	MI	55.9 / <b>59.5</b>	50.0 / <b>52.2</b>	45.9 / <b>50.3</b>	98.7 / <b>99.5</b>	85.2 / <b>88.0</b>	83.3 / <b>87.6</b>	26.9 / <b>29.7</b>	25.7 / <b>26.8</b>	15.3 / <b>16.4</b>	48.5 / <b>51.3</b>
	DI	60.6 / <b>74.0</b>	56.5 / <b>68.2</b>	51.0 / <b>65.7</b>	98.4 / <b>99.6</b>	86.8 / <b>93.6</b>	84.2 / <b>92.0</b>	21.2 / <b>33.6</b>	20.1 / <b>31.9</b>	13.0 / <b>21.6</b>	49.2 / <b>60.1</b>
	TI	25.5 / <b>32.6</b>	21.9 / <b>28.2</b>	11.0 / <b>19.1</b>	98.2 / <b>99.2</b>	53.3 / <b>62.8</b>	48.2 / <b>56.3</b>	18.4 / <b>26.2</b>	18.7 / <b>25.9</b>	12.6 / <b>20.0</b>	26.2 / <b>33.9</b>
	SI	43.6 / <b>56.5</b>	40.3 / <b>50.5</b>	32.4 / <b>47.1</b>	<b>99.8</b> / <b>99.8</b>	84.7 / <b>91.6</b>	83.7 / <b>89.9</b>	19.0 / <b>33.0</b>	18.9 / <b>31.4</b>	12.5 / <b>22.4</b>	41.9 / <b>52.8</b>
	PI	57.5 / <b>63.9</b>	50.3 / <b>57.8</b>	47.4 / <b>55.0</b>	99.6 / <b>99.7</b>	82.7 / <b>90.6</b>	81.8 / <b>87.9</b>	31.7 / <b>38.1</b>	29.5 / <b>37.4</b>	21.2 / <b>27.1</b>	50.3 / <b>57.2</b>

Table 6: The success rate (%) of non-targeted FGSM-based/FGS<sup>2</sup>M-based attacks w.r.t adversarial examples crafted via Res-152. We study nine models—Inc-v3, Inc-v4, IncRes-v2, Res-152, Res-101, Res-50, Inc-v3<sub>ens3</sub>, Inc-v3<sub>ens4</sub> and IncRes-v2<sub>ens</sub> here.