

RETHINKING THE SECURITY OF SKIP CONNECTIONS IN RESNET-LIKE NEURAL NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Skip connections are an essential component of current state-of-the-art deep neural networks (DNNs) such as ResNet, WideResNet, DenseNet, and ResNeXt. Despite their huge success in building deeper and more powerful DNNs, we identify a surprising *security weakness* of skip connections in this paper. Use of skip connections *allows easier generation of highly transferable adversarial examples*. Specifically, in ResNet-like (with skip connections) neural networks, gradients can backpropagate through either skip connections or residual modules. We find that using more gradients from the skip connections rather than the residual modules according to a decay factor, allows one to craft adversarial examples with high transferability. Our method is termed *Skip Gradient Method* (SGM). We conduct comprehensive transfer attacks against 10 state-of-the-art DNNs including ResNets, DenseNets, Inceptions, Inception-ResNet, Squeeze-and-Excitation Network (SENet) and robustly trained DNNs. We show that employing SGM on the gradient flow can greatly improve the transferability of crafted attacks in almost all cases. Furthermore, SGM can be easily combined with existing black-box attack techniques, and obtain high improvements over state-of-the-art transferability methods. Our findings not only motivate new research into the architectural vulnerability of DNNs, but also open up further challenges for the design of secure DNN architectures.

1 INTRODUCTION

In deep neural networks (DNNs), a skip connection builds a short-cut from a shallow layer to a deep layer by connecting the input of a convolutional block (also known as the residual module) directly to its output. While different layers of a neural network learn different “levels” of features, skip connections can help preserve low-level features and avoid performance degradation when adding more layers. This has been shown to be crucial for building very deep and powerful DNNs such as ResNet (He et al., 2016a;b), WideResNet (Zagoruyko & Komodakis, 2016), DenseNet (Huang et al., 2017) and ResNeXt (Xie et al., 2017). In the meantime, despite their superior performance, DNNs have been found extremely vulnerable to adversarial examples (or attacks), which are input examples slightly perturbed with an intention to fool the network to make a wrong classification (Szegedy et al., 2013; Goodfellow et al., 2014). Adversarial examples often appear imperceptible to human observers, and are transferable across different models (Liu et al., 2017).

Adversarial examples can be crafted following either a white-box setting (the adversary has full access to the target model) or a black-box setting (the adversary has no information of the target model). White-box methods such as Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2014), Basic Iterative Method (BIM) (Kurakin et al., 2016), Projected Gradient Decent (PGD) (Madry et al., 2018) and Carlini and Wagner (CW) (Carlini & Wagner, 2017) often suffer from low transferability in a black-box setting, thus posing only limited threats to DNN models which are usually kept secret in practice (Dong et al., 2018; Xie et al., 2019). Several techniques have been proposed to improve the transferability of black-box attacks crafted on a surrogate model, such as momentum boosting (Dong et al., 2018), diverse input (Xie et al., 2019) and translation invariance (Dong et al., 2019). Although these techniques are effective, they (as well as white-box methods) all treat the entire network (either the target model or the surrogate model) as a single component while ignore its inner architectural characteristics. *The question of whether or not the DNN architecture itself can expose more security weaknesses to adversarial attacks is an unexplored problem.*

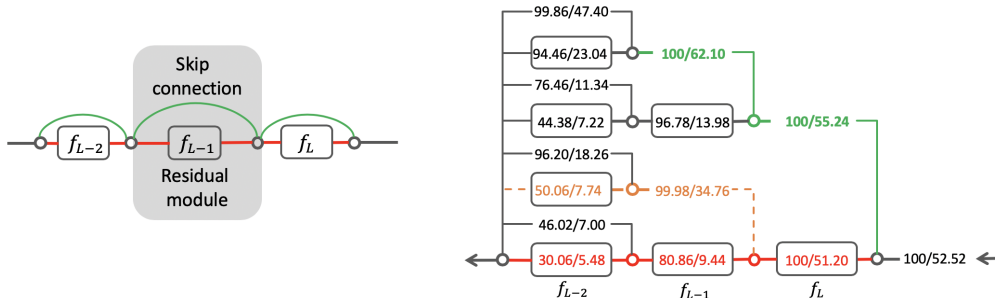


Figure 1: *Left*: Illustration of the last 3 skip connections (green lines) and residual modules (black boxes) of a ImageNet-trained ResNet-18. *Right*: The success rate (in the form of “white-box/black-box”) of adversarial attacks crafted using gradients flowing through either a skip connection (going upwards) or a residual module (going leftwards) at each junction point (circle). Three example backpropagation paths are highlighted in different colors, with the green path skipping over the last two residual modules is the best while the red path through all 3 residual modules is the worst. The attacks are crafted by BIM on 5000 ImageNet validation images under maximum L_∞ perturbation $\epsilon = 16$ (pixel values are in $[0, 255]$). The black-box success rate is tested against a VGG19 target model.

In this paper, we identify one such weakness about the skip connections used by many state-of-the-art DNNs. We first conduct a toy experiment with the BIM attack and ResNet-18 on the ImageNet validation dataset (Deng et al., 2009) to investigate how skip connections affect the adversarial strength of attacks crafted on the network. At each of the last 3 skip connections and residual modules of ResNet-18, we illustrate the success rate of attacks crafted using gradients backpropagate through either the skip connection or the residual module in Figure 1. As can be observed, the success rate drops more drastically whenever using gradients from a residual module instead of the skip connection. This implies that gradients from the skip connections are more vulnerable (high success rate). In addition, we surprisingly find that skip connections expose more transferable information. For example, the black-box success rate was even improved from 52.52% to 62.10% when the attack skips the last two residual modules (following the path in green color).

Motivated by the above observations, in this paper, we propose the *Skip Gradient Method* (SGM) to generate adversarial examples using gradients more from the skip connections rather than the residual modules. In particular, SGM utilizes a decay factor to reduce gradients from the residual modules. We find that this simple adjustment on the gradient flow can generate highly transferable adversarial examples, and the more skip connections in a network, the more transferable are the crafted attacks. This is in sharp contrast to the design principles (eg. “going deeper” with skip connections) underpinning many modern DNNs. In particular, our main contributions are:

- We identify an important security weakness of skip connections in ResNet-like neural networks, *i.e.*, they allow an easy generation of highly transferable adversarial examples.
- We propose the *Skip Gradient Method* (SGM) to craft adversarial examples using gradients more from the skip connections. Using a single decay factor on gradients, SGM is an appealingly simple and generic technique that can be used by any existing gradient-based attack methods.
- We provide comprehensive transfer attack experiments, from 8 different source models against 10 state-of-the-art DNNs, showing that SGM can greatly improve the transferability of crafted adversarial examples. When combined with existing transfer techniques, SGM improves the state-of-the-art transferability benchmarks by a large margin.

2 RELATED WORK

Existing adversarial attacks can be categorized into two groups: 1) white-box attacks and 2) black-box attacks. In the white-box setting, the adversary has full access to the parameters of the target model, while in the black-box setting, the target model is kept secret from the adversary.

2.1 WHITE-BOX ATTACKS

Given a clean example \mathbf{x} with class label y and a target DNN model f , the goal of an adversary is to find an adversarial example \mathbf{x}_{adv} that fools the network into making an incorrect prediction (eg. $f(\mathbf{x}_{adv}) \neq y$), while still remaining in the ϵ -ball centered at \mathbf{x} (eg. $\|\mathbf{x}_{adv} - \mathbf{x}\|_\infty \leq \epsilon$).

Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2014). FGSM perturbs clean example \mathbf{x} for one step by the amount of ϵ along the gradient direction:

$$\mathbf{x}_{adv} = \mathbf{x} + \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \ell(f(\mathbf{x}), y)). \quad (1)$$

The Basic Iterative Method (BIM) (Kurakin et al., 2016) is an iterative version of FGSM that perturbs for T steps with step size ϵ/T .

Projected Gradient Descent (PGD) (Madry et al., 2018). PGD perturbs normal example \mathbf{x} for T steps with smaller step size. After each step of perturbation, PGD projects the adversarial example back onto the ϵ -ball of \mathbf{x} , if it goes beyond the ϵ -ball:

$$\mathbf{x}_{adv}^{t+1} = \Pi_\epsilon(\mathbf{x}_{adv}^t + \alpha \cdot \text{sign}(\nabla_{\mathbf{x}} \ell(f(\mathbf{x}_{adv}^t), y))), \quad (2)$$

where $\Pi_\epsilon(\cdot)$ is the projection operation. Different to BIM, PGD allows step size $\alpha > \epsilon/T$.

There are also other types of white-box attacks including sparsity-based methods such as Jacobian-based Saliency Map Attack (JSMA) (Papernot et al., 2016), sparse attack (Modas et al., 2019), one-pixel attack (Su et al., 2019), and optimization-based methods such as Carlini and Wagner (CW) (Carlini & Wagner, 2017) and elastic-net (EAD) (Chen et al., 2018).

2.2 BLACK-BOX ATTACKS

Black-box attacks can be generated by either attacking a surrogate model or using gradient estimation methods in combination with queries to the target model. Gradient estimation methods estimate the gradients of the target model using black-box optimization methods such as Finite Differences (FD) (Chen et al., 2017; Bhagoji et al., 2018) or Natural Evolution Strategies (NES) (Ilyas et al., 2018; Jiang et al., 2019). These methods all require a large number of queries to the target model, which not only reduces efficiency but also potentially exposes the attack. Alternatively, black-box adversarial examples can be crafted on a surrogate model then applied to attack the target model. Although the white-box methods can be directly applied on the surrogate model, they are far less effective in the black-box setting (Dong et al., 2018; Xie et al., 2019). Several transfer techniques have been proposed to improve the transferability of black-box attacks.

Momentum Iterative boosting (MI) (Dong et al., 2018). MI incorporates a momentum term into the gradient to boost the transferability:

$$\mathbf{x}_{adv}^{t+1} = \Pi_\epsilon(\mathbf{x}_{adv}^t + \alpha \cdot \text{sign}(\mathbf{g}^{t+1})), \quad \mathbf{g}^{t+1} = \mu \cdot \mathbf{g}^t + \frac{\nabla_{\mathbf{x}} \ell(f(\mathbf{x}_{adv}^t), y)}{\|\nabla_{\mathbf{x}} \ell(f(\mathbf{x}_{adv}^t), y)\|_1}, \quad (3)$$

where \mathbf{g}^t is the adversarial gradient at the t -th step, $\alpha = \epsilon/T$ is the step size for a total of T steps, μ is a decay factor, and $\|\cdot\|_1$ is the L_1 norm.

Diverse Input (DI) (Xie et al., 2019). DI proposes to craft adversarial examples using gradient with respect to the randomly-transformed input example:

$$\mathbf{x}_{adv}^{t+1} = \Pi_\epsilon(\mathbf{x}_{adv}^t + \alpha \cdot \text{sign}(\nabla_{\mathbf{x}} \ell(f(H(\mathbf{x}_{adv}^t; p)), y))), \quad (4)$$

where $H(\mathbf{x}_{adv}^t; p)$ is a stochastic transformation function on \mathbf{x}_{adv}^t for a given probability p .

Translation Invariant (TI) (Dong et al., 2019). TI targets to evade robustly trained DNNs by generating adversarial examples that are less sensitive to the discriminative regions of the surrogate model. More specifically, TI computes the gradients with respect to a set of translated versions of the original input:

$$\mathbf{x}_{adv}^{t+1} = \Pi_\epsilon(\mathbf{x}_{adv}^t + \alpha \cdot \text{sign}(\mathbf{W} * \nabla_{\mathbf{x}} \ell(f(\mathbf{x}_{adv}^t), y))), \quad (5)$$

where \mathbf{W} is a predefined kernel (e.g., uniform, linear, and Gaussian) matrix of size $(2k+1)(2k+1)$ (k being the maximal number of pixels to shift). This kernel convolution is equivalent to the weighted sum of gradients over $(2k+1)^2$ number of shifted input examples.

Furthermore, there are other studies focusing on intermediate feature representations. For example, Activation Attack (Inkawhich et al., 2019) drives the activation of a specified layer on a given image towards the layer of a target image, to yield a highly transferable targeted example. Intermediate Level Attack (Huang et al., 2019) attempts to fine-tune an existing adversarial example for greater black-box transferability by increasing its perturbation on a pre-specified layer of the source model.

Although the above transfer techniques are effective, they (including white-box attacks) either 1) treat the network (either the surrogate model or the target model) as a single component or 2) only use the intermediate layer output of the network. In other words, they do not directly consider the effects of different DNN architectural characteristics. In the following, we exploits an architectural security weakness about the skip connection.

3 PROPOSED SKIP GRADIENT ATTACK

In this section, we first introduce the gradient decomposition of skip connection and residual module. Following that, we propose our Skip Gradient Method (SGM), then demonstrate the security weakness of skip connection via a case study.

3.1 GRADIENT DECOMPOSITION WITH SKIP CONNECTIONS

In ResNet-like neural networks, a skip connection uses identity mapping to bypass residual layers, allowing data flow from a shallow layer directly to subsequent deep layers. Thus, we can decompose the network into a collection of paths of different lengths (Veit et al., 2016). We denote a skip connection together with its associated residual module as a building block (residual block) of a network. Considering three successive building blocks (eg. $z_{i+1} = z_i + f_{i+1}(z_i)$) in a residual network from input z_0 to output z_3 , the output z_3 can be expanded as:

$$\begin{aligned} z_3 &= z_2 + f_3(z_2) = [z_1 + f_2(z_1)] + f_3(z_1 + f_2(z_1)) \\ &= [z_0 + f_1(z_0) + f_2(z_0 + f_1(z_0))] + f_3((z_0 + f_1(z_0)) + f_2(z_0 + f_1(z_0))). \end{aligned} \quad (6)$$

According to the chain rule in calculus, the gradient of a loss function ℓ with respect to input z_0 can then be decomposed as,

$$\frac{\partial \ell}{\partial z_0} = \frac{\partial \ell}{\partial z_3} \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial z_0} = \frac{\partial \ell}{\partial z_3} (1 + \frac{\partial f_3}{\partial z_2}) (1 + \frac{\partial f_2}{\partial z_1}) (1 + \frac{\partial f_1}{\partial z_0}). \quad (7)$$

Extending this toy example to a network with L residual blocks, the gradient can be decomposed from L -th to the $(l + 1)$ -th ($0 \leq l < L$) residual block as,

$$\frac{\partial \ell}{\partial x} = \frac{\partial \ell}{\partial z_L} \prod_{i=l}^{L-1} \left(\frac{\partial f_{i+1}}{\partial z_i} + 1 \right) \frac{\partial z_l}{\partial x}. \quad (8)$$

The example illustrated in Figure 1 is a the above decomposition of a ResNet-18 at the last 3 building blocks ($l = L - 3$).

3.2 SKIP GRADIENT METHOD (SGM)

In order to use more gradient from the skip connections, here, we introduce a decay parameter into the decomposed gradient to reduce the gradient from the residual modules. Following the decomposition in Equation (8), the ‘‘skipped’’ gradient is,

$$\nabla_x \ell = \frac{\partial \ell}{\partial z_L} \prod_{i=0}^{L-1} \left(\gamma \frac{\partial f_{i+1}}{\partial z_i} + 1 \right) \frac{\partial z_0}{\partial x}, \quad (9)$$

where $z_0 = x$ is the input of the network, and $\gamma \in (0, 1]$ is the decay parameter. Accordingly, given a clean example x and a DNN model f , an adversarial example can be crafted iteratively by,

$$x_{adv}^{t+1} = \Pi_\epsilon \left(x_{adv}^t + \alpha \cdot \text{sign} \left(\frac{\partial \ell}{\partial z_L} \prod_{i=0}^{L-1} \left(\gamma \frac{\partial f_{i+1}}{\partial z_i} + 1 \right) \frac{\partial z_0}{\partial x} \right) \right). \quad (10)$$

SGM is a generic technique that can be easily implemented on any neural network that has skip connections. During the backpropagation process, SGM simply multiplies the decay parameter to

Table 1: The success rates (%) of black-box attacks (untargeted) crafted by FGSM, PGD and their “skip gradient” (+SGM) versions, against a VGG19 target model. The best results are in **bold**.

| Attack \ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
|-----------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| FGSM | 54.92 | 52.20 | 45.28 | 42.12 | 41.36 | 51.32 | 48.24 | 49.84 |
| FGSM+SGM | 56.60 | 57.54 | 49.90 | 49.40 | 47.46 | 54.96 | 55.12 | 56.66 |
| PGD | 59.88 | 54.50 | 52.34 | 43.24 | 45.26 | 61.04 | 54.98 | 57.68 |
| PGD+SGM | 64.52 | 71.48 | 68.32 | 62.44 | 66.14 | 81.74 | 79.28 | 81.24 |

the gradient whenever it passes a residual module. Therefore, SGM does not require any computation overhead, and works efficiently even on densely connected networks such as DenseNets. The reduction of residual gradients is accumulated along the backpropagation path, that is, the residual gradients at lower layers will be reduced more times than those at higher layers. This is because, compared to high-level features, low-level features have already been well preserved by skip connections (see feature decompositions in Equation (6)).

3.3 SECURITY WEAKNESS OF SKIP CONNECTIONS: A CASE STUDY

To demonstrate the security weakness of skip connections, we conduct a case study, on FGSM, 10-step PGD, and their corresponding SGM versions, to investigate the success rates of black-box attacks crafted with or without manipulating the skip connections. The black-box attacks are generated on 8 different source (surrogate) models ResNet(RN)-18/34/50/101/152 and DenseNet(DN)-121/169/201, then applied to attack a VGG19 target model. All models were trained on ImageNet training set. We randomly select 5000 ImageNet validation images that are correctly classified by all source models, and craft untargeted attacks under maximum L_∞ perturbation $\epsilon = 16$, which is a typical black-box setting (Dong et al., 2018; Xie et al., 2019; Dong et al., 2019). The step size of PGD was set to $\alpha = 2$, and the decay parameter of SGM was set to $\gamma = 0.6$ (results with varying γ can be found in Appendix C).

The success rates (transferability) of different methods are reported in Table 1. As can be seen, when the skip connections are manipulated with our SGM, the transferability of FGSM and PGD is greatly improved across all source models. On some source models such as DN169 and DN201, the improvements are even more than 23%. Without SGM, the best transferability against the VGG19 target model is 61.04% which is achieved by PGD on DN121, however, this is improved further by our proposed SGM to 81.74% ($> 20\%$ gain). This not only highlights the security weakness of skip connections in terms of the generation of highly transferable attacks, but also indicates the severeness of this weakness, as such a huge boost in transferability only takes a single decay factor. Another important observation is that when there are more skip connections in a network (*e.g.*, DenseNet $>$ ResNet), the the crafted attacks become more transferable, especially when the skip connections are manipulated by our SGM. This raises questions about the design principle behind many state-of-the-art DNNs: “going deeper” with techniques like skip connection and 1×1 convolution.

4 COMPARISON TO EXISTING TRANSFER ATTACKS

In this section, we compare the transferability of adversarial examples crafted by our proposed SGM and existing methods on ImageNet dataset (Deng et al., 2009) against both unsecured and secured target models.

Baselines. We compare SGM with FGSM, PGD, and 3 state-of-the-art transfer attacks: (1) Momentum Iterative (MI) (Dong et al., 2018), (2) Diverse Input (DI) (Xie et al., 2019), and (3) Transition Invariant (TI) (Dong et al., 2019). Note that the TI attack was originally proposed to attack secured models, although here we include TI to attack both unsecured models and secured models. For TI and our SGM, we test both the one-step and the iterative version, however, the other methods DI and MI only have an iterative version. The iteration step is set to 10 and 20 for unsecured and secured target models respectively. For all iterative methods PGD, TI and our SGM, the step size is set to $\alpha = 2$. For our proposed SGM, the decay parameter is set to $\gamma = 0.2$ (0.5) and $\gamma = 0.5$ (0.7) on ResNet and DenseNet source models in PGD (FGSM) respectively. Other parameters of existing methods are configured as in their original papers.

Threat Model. We adopt a black-box threat model in which adversarial examples are generated by attacking a source model and then applied to attack the target model. The target model is of

either a different architecture to the source model, or the same architecture but trained separately. The attacks are crafted on 5000 randomly selected ImageNet validation images that are classified correctly by all source models. For all attack methods, we follow the standard setting (Dong et al., 2018; Xie et al., 2019) to craft untargeted attacks under maximum L_∞ perturbation $\epsilon = 16$ with respect to pixel values in $[0, 255]$.

Target Models. We consider two types of target models: 1) unsecured models that are trained on ImageNet training set using traditional training; and 2) secured models trained using adversarial training. For unsecured target model, we choose 7 state-of-the-art DNNs: VGG19 (with batch normalization) (Simonyan & Zisserman, 2015), ResNet-152 (RN152) (He et al., 2016a), DenseNet-201 (DN152), 154 layer Squeeze-and-Excitation network (SE154) (Hu et al., 2018), Inception V3 (IncV3) (Szegedy et al., 2016), Inception V4 (IncV4) (Szegedy et al., 2017) and Inception-ResNet V2 (IncResV2) (Szegedy et al., 2017). For secured target models, we consider 3 robustly trained DNNs using ensemble adversarial training (Tramèr et al., 2018): IncV3_{ens3} (ensemble of 3 IncV3 networks), IncV3_{ens4} (ensemble of 4 IncV3 networks) and IncResV2_{ens3} (ensemble of 3 IncResV2 networks).

Source Models. We choose 8 different source models from the ResNet family: ResNet(RN)-18/34/50/101/152 and DenseNet(DN)-121/169/201. Whenever the input size of the source model does not match the target model, we resize the crafted adversarial images to the input size of the target model. For VGG19, ResNet and DenseNet models, images are cropped and resized to 224×224 , while for Inception/Inception-ResNet models, images are cropped and resized to 299×299 .

4.1 TRANSFERABILITY AGAINST UNSECURED MODELS

We first investigate the transferability of all attack methods against the 7 unsecured models. The goal is to find the best method that can generate the most transferable attacks on *one* source model against *all* target models.

One-step Transferability. The one-step transferability is measured by the success rate of one-step attacks, as reported in Table 2. Here, we only show the results on two source models: 1) RN152 which is the best ResNet source model with the highest success rate on average against all target models, and 2) DN201 which is the best DenseNet source model. Also note that, when the source and target models are the same, the result represents the white-box success rate. Overall, adversarial examples crafted on DN201 have significantly better transferability than those crafted on RN152, especially for our SGM method. This is because there are $\sim 30\times$ more skip connections that can be manipulated by our SGM in DN201 compared to RN152. In comparison to both FGSM and TI, transferability is improved considerably by SGM in almost all test scenarios, except when transferring from RN152 to VGG19/IncV3/IncV4 where SGM is outperformed by TI. This implies that, when transferring across different architectures (eg. ResNet \rightarrow VGG/Inception), translation adaptation may help increase the transferability of one-step perturbations. However, this advantage of TI disappears when there are more skip connections, as is the case for the DN201 source model.

Table 2: One-step transferability: the success rates (%) of black-box attacks crafted by different methods on 2 source models against 7 unsecured target models. The best results are in **bold**.

| Source | Attack | VGG19 | RN152 | DN201 | SE154 | IncV3 | IncV4 | IncResV2 |
|--------|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| RN152 | FGSM | 41.36 | 71.84 | 37.40 | 30.00 | 25.74 | 21.62 | 20.46 |
| | TI | 49.52 | 49.66 | 36.54 | 30.18 | 33.86 | 29.06 | 20.64 |
| | SGM | 47.70 | 77.54 | 43.56 | 31.16 | 29.18 | 25.00 | 22.80 |
| DN201 | FGSM | 49.84 | 39.10 | 81.44 | 35.46 | 31.74 | 26.82 | 24.10 |
| | TI | 54.00 | 33.62 | 57.72 | 34.02 | 34.58 | 30.12 | 20.74 |
| | SGM | 56.70 | 47.38 | 87.72 | 42.84 | 38.36 | 32.56 | 29.92 |

Multi-step Transferability. First we provide a detailed study about the transferability of all attack methods from the 8 source models to the 7 unsecured target models. We then compare different attack methods on two best source models: the best ResNet source model and the best DenseNet source model. The multi-step (e.g., 10 step) transferability from all source models to three representative target models (VGG19, SE154 and IncV3) is illustrated in Figure 2 (see Appendix B for more results). In all transfer scenarios, our proposed SGM outperforms existing methods consistently on almost all source models except RN18. Adversarial attacks crafted by SGM become

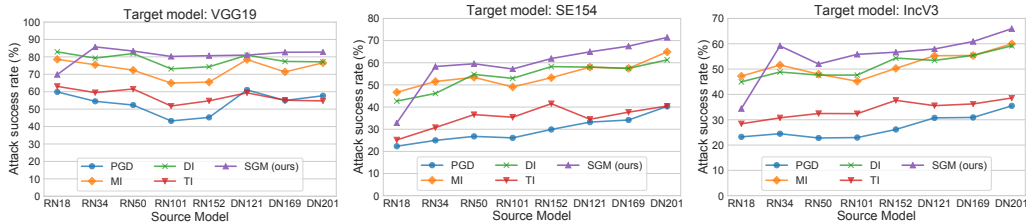


Figure 2: The attack success rates of black-box attacks crafted by different attack methods on 8 source models against 3 *unsecured* target models: VGG19 (left), SE154 (middle) and IncV3 (right).

Table 3: Multi-step transferability: the success rates (%) of black-box attacks crafted by different methods on 2 source models against 7 *unsecured* target models. The best results are in **bold**.

| Source | Attack | VGG19 | RN152 | DN201 | SE154 | IncV3 | IncV4 | IncRes |
|--------|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| RN152 | PGD | 45.26 | 99.92 | 50.90 | 29.90 | 26.16 | 20.84 | 18.78 |
| | TI | 54.78 | 99.62 | 62.94 | 41.54 | 37.66 | 34.94 | 28.30 |
| | MI | 65.52 | 99.78 | 75.36 | 53.26 | 50.26 | 43.04 | 41.58 |
| | DI | 74.32 | 99.90 | 78.04 | 58.26 | 54.34 | 47.16 | 43.26 |
| | SGM | 80.68 | 99.90 | 81.70 | 61.90 | 56.66 | 48.50 | 45.86 |
| DN201 | PGD | 57.68 | 59.40 | 99.86 | 40.32 | 35.48 | 31.82 | 26.04 |
| | TI | 54.72 | 49.46 | 99.56 | 40.46 | 38.58 | 36.06 | 28.40 |
| | MI | 75.14 | 76.58 | 99.80 | 64.92 | 59.96 | 54.28 | 49.78 |
| | DI | 77.68 | 77.14 | 99.76 | 61.30 | 59.18 | 55.80 | 48.08 |
| | SGM | 82.82 | 86.16 | 99.58 | 71.72 | 65.38 | 59.12 | 55.24 |

more transferable when there are more skip connections in the source model (*e.g.*, from RN18 to DN201). An interesting observation is that, when the target model is shallow such as VGG19 (left figure in Figure 2), shallow source models transfer better, however, when the target model is deep such as SE154 and IncV3 (middle and right figures in Figure 2), deeper source models tend to have better transferability. We suspect this is due to the architectural similarities shared by the target and source models. Note that against the VGG19 target model, the success rate of baseline methods all drop significantly when the ResNet source models become more complex (from RN18 to RN152). The small variations at RN50 and DN121 source models may be caused by the architectural difference between RN18/34 which consist of normal residual blocks, RN50/101/152 which consist of “bottleneck” residual blocks and DN121/169/201 which has dense skip connections.

Results for the best source models RN152 and DN201 against all target models are reported in Table 3. The proposed SGM attack outperforms existing methods by a large margin consistently against all target models. Particularly, for transfer DN201 \rightarrow SE154 (a recent state-of-the-art DNN with only 2.251% top-5 error on ImageNet), SGM achieves a success rate of 71.4%, which is $> 6\%$ and $> 10\%$ higher than MI and DI respectively.

Combining with Existing Methods. We further demonstrate that the weakness of skip connections can be exploited in combination with existing techniques. The experiments are conducted on DN201 (the best source model in the above multi-step experiments), and TI attack is excluded as it was originally proposed against secured models and demonstrates limited improvement over PGD against unsecured models. The results are reported in Table 4. The transferability of MI and DI is improved remarkably of 11.98% \sim 21.98% by SGM. When combined with both MI and DI, SGM improves the state-of-the-art (MI+DI) transferability by a huge margin consistently against all target models. In particular, SGM pushes the new state-of-the-art to at least 80.52% which previously was only 71%. This illustrates that the security weakness of skip connections can be easily manipulated to craft highly transferable attacks against many state-of-the-art DNN models.

4.2 TRANSFERABILITY AGAINST ROBUSTLY TRAINED MODELS

The success rates of our SGM and other baseline methods against the 3 secured target models are reported in Table 5. Overall, with translation adaptation specifically designed for evading adversarially trained models, TI achieves the best standalone transferability, while SGM is the second best

Table 4: Combined with existing methods: the success rates (%) of attacks crafted on source model DN201 against 7 *unsecured* target models. The best results are in **bold** and + indicates improvement.

| Attack \ Target | VGG19 | RN152 | DN201 | SE154 | IncV3 | IncV4 | IncRes |
|------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| MI | 75.14 | 76.58 | 99.80 | 64.92 | 59.96 | 54.28 | 49.78 |
| MI+SGM | +11.98 | +13.24 | 99.52 | +17.26 | +21.98 | +15.62 | +18.38 |
| DI | 77.68 | 77.14 | 99.76 | 61.30 | 59.18 | 55.80 | 48.08 |
| DI+SGM | +12.38 | +13.86 | 99.52 | +20.82 | +17.76 | +15.78 | +20.20 |
| MI+DI | 87.06 | 87.30 | 99.76 | 79.90 | 76.68 | 75.18 | 71.00 |
| MI+DI+SGM | 93.02 | 93.90 | 99.52 | 89.66 | 85.68 | 81.22 | 80.52 |

Table 5: Transferability against secured models: the success rates of multi-step attacks crafted on RN152 and DN201 source models against 3 secured models. The best results are in **bold**.

| Source | Attack | IncV3 _{ens3} | IncV _{ens4} | IncRes _{ens3} |
|--------|---------------|-----------------------|----------------------|------------------------|
| RN152 | PGD | 11.24 | 9.22 | 6.21 |
| | TI | 44.28 | 44.80 | 37.42 |
| | MI | 22.84 | 20.98 | 15.55 |
| | DI | 27.29 | 22.88 | 16.78 |
| | SGM | 31.18 | 27.44 | 19.56 |
| | TI+SGM | 52.62 | 52.80 | 43.96 |
| DN201 | PGD | 17.64 | 14.69 | 10.18 |
| | TI | 41.75 | 41.10 | 33.72 |
| | MI | 31.07 | 28.04 | 20.73 |
| | DI | 33.29 | 28.46 | 21.02 |
| | SGM | 41.25 | 37.87 | 29.42 |
| | TI+SGM | 44.70 | 46.35 | 38.41 |

with higher success rates than either PGD, MI or DI. When combined with TI, SGM also improves the TI attack by a considerable margin across all transfer scenarios. This indicates that, although manipulating the skip connections alone may not sufficient to attack secured models, it still can make existing attacks more powerful. One interesting observation is that attacks crafted here on RN152 are more transferable than those crafted on DN201, which is quite the opposite to attacking unsecured models.

4.3 A CLOSER LOOK AT SGM

In this part, we conduct more experiments to investigate the gradient decay factor of our proposed SGM, and explore the potential use of SGM for white-box attacks.

Effect of Residual Gradient Decay. We test the transferability of our proposed SGM with varying decay parameter $\gamma \in [0.1, 1.0]$, where $\gamma = 1.0$ means no decay on the residual gradients. The attacks are crafted by 10-step SGM on 5000 random ImageNet validation images. The results against 3 target models (VGG19, SE154 and IncV3) are illustrated in Figure 3 (See Appendix C for more results). As can be observed, the trends are very consistent against different target models. On DenseNet source models, decreasing decay parameter (increasing decay strength) tends to improve transferability until it exceeds a certain threshold, *e.g.*, $\gamma = 0.5$. This is because the decay encourages the attack to focus on more transferable low-level information, however, it becomes less sufficient if all high-level class-relevant information is ignored. On ResNet source models, decreasing decay parameter can constantly improve transferability for $\gamma \geq 0.2$. Compared to DenseNet source models, ResNets require more decay on the residual gradients. Recalling that skip connections reveal more transferable information of the source model, ResNets require more penalty on the residual gradients to increase the importance of skip gradients that reveal more transferable information of the source model.

Improving Weak White-box Attacks. In addition to the black-box transferability, we next show that SGM can also improve the weak (one-step) white-box attack FGSM. Note that the one-step version of SGM is equivalent to FGSM plus residual gradient decay. Our experiments are conducted on the 8 source models, and the white-box success rates under maximum L_∞ perturbation $\epsilon = 8$

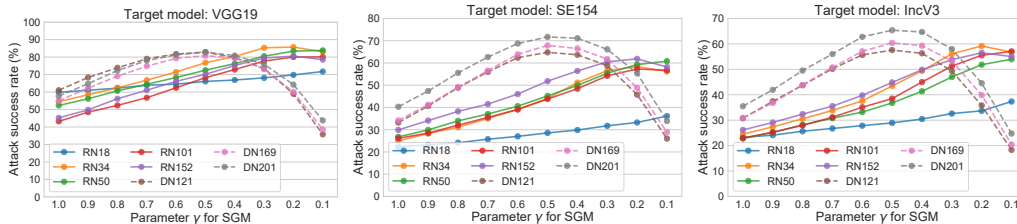
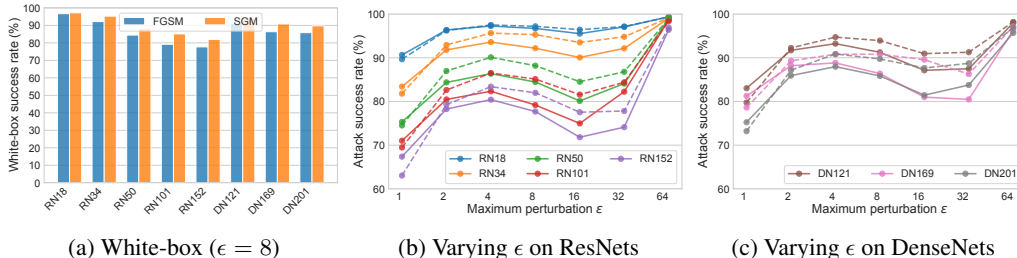


Figure 3: Parameter tuning: the success rates of black-box attacks crafted by 10-step SGM with varying decay parameter $\gamma \in [0.1, 1.0]$. The solid and dash curves represent results on ResNet and DenseNet source models respectively.



(a) White-box ($\epsilon = 8$) (b) Varying ϵ on ResNets (c) Varying ϵ on DenseNets

Figure 4: White-box success rate for FGSM versus SGM. In (b) and (c), each color corresponds to one model, with FGSM is represented by solid curve and SGM is represented by dashed curve.

(a typical white-box setting) are shown in Figure 4a. As can be observed, using SGM can help improve the adversarial strength (*i.e.*, higher success rate). We then vary the maximum perturbation $\epsilon \in [1, 64]$, and show the results on ResNet and DenseNet models separately in Figure 4b and Figure 4c. Compared to FGSM, SGM can always give better adversarial strength, except when ϵ is extremely small ($\epsilon \leq 2$). When the perturbation space becomes infinitely small, the loss landscape within the space becomes flat and the gradient points to the optimal perturbation direction. However, when the perturbation space expands, one-step gradient becomes less accurate due to changes in the loss landscape (success rate decreases as ϵ increases from 4 to 16), and in this case, the skip gradient which contains more low-level information is more reliable than the residual gradient (the improvement is more significant for $\epsilon \in [4, 16]$). Another interesting observation is that adversarial strength decreases when the model becomes more complex from RN18 to RN152, or DN121 to DN201. This is likely because the loss landscape of complex models is steeper than shallow models, making one-step gradient less reliable.

5 CONCLUSION

In this paper, we have identified a surprising security weakness of the skip connections used by many state-of-the-art ResNet-like neural networks, that is, they can be easily used to generate highly transferable adversarial examples. To demonstrate this architectural weakness, we proposed the *Skip Gradient Method* (SGM) to craft adversarial examples using more gradients from the skip connections rather than the residual ones, via a decay factor on gradients. We conducted a series of transfer attack experiments with 8 source models and 10 target models including 7 unsecured and 3 secured models, and showed that attacks crafted by SGM have significantly better transferability than those crafted by existing methods. When combined with existing techniques, SGM can also boost state-of-the-art transferability by a huge margin. We believe this security weakness of skip connections is due to the fact that they expose extra low-level information which is more transferable across different DNNs. Our findings in this paper not only reminds researcher to pay attention to the architectural vulnerability of DNNs, but also raises new challenges for secure DNN architecture design.

REFERENCES

Arjun Nitin Bhagoji, Warren He, Bo Li, and Dawn Song. Practical black-box attacks on deep neural networks using efficient query mechanisms. In *ECCV*, 2018.

- Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *S&P*, 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 *AISeC*, 2017.
- Pin-Yu Chen, Yash Sharma, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. Ead: elastic-net attacks to deep neural networks via adversarial examples. In *AAAI*, 2018.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- 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.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2014.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016a.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *ECCV*, 2016b.
- Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, 2018.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *CVPR*, 2017.
- Qian Huang, Isay Katsman, Horace He, Zeqi Gu, Serge Belongie, and Ser-Nam Lim. Enhancing adversarial example transferability with an intermediate level attack. In *ICCV*, 2019.
- Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks with limited queries and information. In *ICML*, 2018.
- Nathan Inkawhich, Wei Wen, Hai Helen Li, and Yiran Chen. Feature space perturbations yield more transferable adversarial examples. In *CVPR*, 2019.
- Linxi Jiang, Xingjun Ma, Shaoxiang Chen, James Bailey, and Yu-Gang Jiang. Black-box adversarial attacks on video recognition models. *arXiv preprint arXiv:1904.05181*, 2019.
- Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In *ICLR*, 2016.
- 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.
- Apostolos Modas, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. Sparsefool: a few pixels make a big difference. In *CVPR*, 2019.
- Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z Berkay Celik, and Ananthram Swami. The limitations of deep learning in adversarial settings. In *EuroS&P*, 2016.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai. One pixel attack for fooling deep neural networks. In *IEEE Transactions on Evolutionary Computation*. IEEE, 2019.

- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *ICLR*, 2013.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking 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.
- Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian Goodfellow, Dan Boneh, and Patrick McDaniel. Ensemble adversarial training: Attacks and defenses. In *ICLR*, 2018.
- Andreas Veit, Michael J Wilber, and Serge Belongie. Residual networks behave like ensembles of relatively shallow networks. In *NeurIPS*, 2016.
- 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*, 2019.
- Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *CVPR*, 2017.
- Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.

A VISUALIZATION OF ADVERSARIAL EXAMPLES CRAFTED BY SGM

In this section, we visualize 6 clean images and their corresponding adversarial examples crafted using our SGM on either a ResNet-152 or a DenseNet201 in Figure 5. These visualization results show that the generated adversarial perturbations are human imperceptible.



Figure 5: Visualization of 6 clean images and their corresponding adversarial examples. The clean images are shown in the top row, adversarial images crafted on ResNet-152 are shown in the middle row, while those crafted on DenseNet-201 are shown at the bottom. All adversarial images are crafted using our proposed SGM (10-step) under maximum perturbation $\epsilon = 16$.

B MORE RESULTS FOR MULTI-STEP TRANSFERABILITY

In addition to Figure 2 and Table 3, here, we show more multi-step transferability results for the remaining 5 source models (against all 7 target models) in Figure 6 and Table 6 respectively. The proposed SGM attack outperforms existing methods by a large margin consistently in all transfer scenarios except ResNet-18. We believe this is because there are only 5 “standard” residual blocks in ResNet-18, which might limit the performance of the proposed SGM. For comparison, there are 13 “standard” residual blocks in ResNet-34.

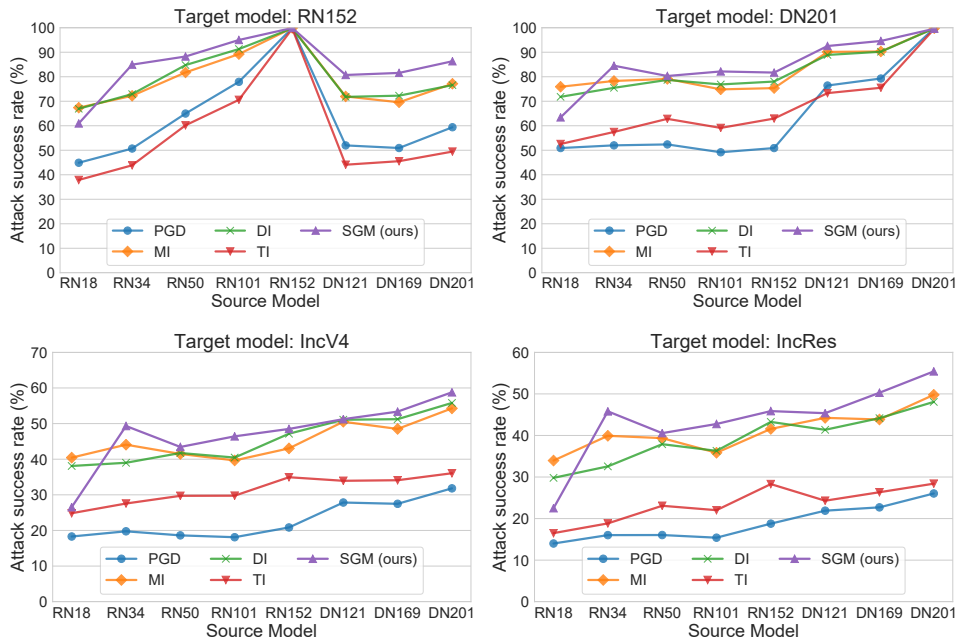


Figure 6: The attack success rates (%) of black-box attacks crafted by different attack methods on 8 source models against all 7 *unsecured* target models.

Table 6: Multi-step transferability: the success rates (%) of black-box attacks crafted by different methods on total 8 source models against 7 *unsecured* target models. The best results are in **bold**.

| | Attack | VGG19 | RN152 | DN201 | SN154 | IncV3 | IncV4 | IncRes |
|-------|--------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| RN18 | PGD | 59.88 | 44.90 | 50.90 | 22.34 | 23.24 | 18.30 | 14.00 |
| | TI | 63.06 | 37.88 | 52.56 | 25.18 | 28.40 | 24.84 | 16.50 |
| | MI | 78.58 | 67.38 | 75.92 | 46.68 | 47.22 | 40.44 | 33.96 |
| | DI | 82.88 | 66.94 | 71.80 | 42.68 | 44.96 | 38.10 | 29.80 |
| | SGM | 69.84 | 60.94 | 63.42 | 32.82 | 34.40 | 26.54 | 22.50 |
| RN34 | PGD | 54.50 | 50.70 | 52.00 | 24.98 | 24.48 | 19.76 | 16.02 |
| | TI | 59.50 | 43.86 | 57.46 | 30.76 | 30.76 | 27.56 | 18.84 |
| | MI | 75.50 | 72.22 | 78.28 | 51.54 | 51.60 | 44.08 | 39.92 |
| | DI | 79.28 | 72.92 | 75.48 | 46.16 | 48.84 | 39.00 | 32.56 |
| | SGM | 85.72 | 84.94 | 84.50 | 58.30 | 59.14 | 49.32 | 45.80 |
| RN50 | PGD | 52.34 | 64.98 | 52.38 | 26.78 | 22.80 | 18.62 | 16.04 |
| | TI | 61.56 | 60.18 | 62.82 | 36.60 | 32.46 | 29.72 | 23.08 |
| | MI | 72.36 | 81.76 | 79.10 | 53.42 | 47.98 | 41.48 | 39.36 |
| | DI | 81.90 | 84.76 | 78.70 | 54.72 | 47.62 | 41.74 | 37.90 |
| | SGM | 83.34 | 88.26 | 80.30 | 59.54 | 52.00 | 43.44 | 40.56 |
| RN101 | PGD | 43.24 | 77.92 | 49.18 | 26.10 | 22.96 | 18.10 | 15.40 |
| | TI | 51.80 | 70.58 | 59.12 | 35.42 | 32.40 | 29.76 | 22.02 |
| | MI | 64.96 | 89.24 | 74.84 | 49.06 | 45.18 | 39.66 | 35.84 |
| | DI | 73.14 | 91.30 | 76.86 | 52.96 | 47.64 | 40.46 | 36.30 |
| | SGM | 80.26 | 95.00 | 82.16 | 57.18 | 55.82 | 46.42 | 42.76 |
| RN152 | PGD | 45.26 | 99.92 | 50.90 | 29.90 | 26.16 | 20.84 | 18.78 |
| | TI | 54.78 | 99.62 | 62.94 | 41.54 | 37.66 | 34.94 | 28.30 |
| | MI | 65.52 | 99.78 | 75.36 | 53.26 | 50.26 | 43.04 | 41.58 |
| | DI | 74.32 | 99.90 | 78.04 | 58.26 | 54.34 | 47.16 | 43.26 |
| | SGM | 80.68 | 99.90 | 81.70 | 61.90 | 56.66 | 48.50 | 45.86 |
| DN121 | PGD | 61.04 | 51.98 | 76.42 | 33.22 | 30.74 | 27.84 | 21.90 |
| | TI | 59.36 | 44.10 | 73.32 | 34.50 | 35.56 | 33.94 | 24.30 |
| | MI | 78.44 | 71.96 | 90.14 | 58.02 | 55.26 | 50.52 | 44.24 |
| | DI | 80.98 | 71.78 | 88.90 | 58.02 | 53.42 | 51.08 | 41.36 |
| | SGM | 81.00 | 80.72 | 92.54 | 64.92 | 57.94 | 51.24 | 45.36 |
| DN169 | PGD | 54.98 | 50.92 | 79.28 | 34.16 | 30.90 | 27.48 | 22.70 |
| | TI | 55.12 | 45.54 | 75.50 | 37.70 | 36.24 | 34.10 | 26.34 |
| | MI | 71.38 | 69.58 | 90.32 | 57.46 | 55.36 | 48.50 | 43.84 |
| | DI | 77.38 | 72.28 | 90.20 | 57.46 | 55.42 | 51.24 | 44.16 |
| | SGM | 82.66 | 81.56 | 94.60 | 67.48 | 60.88 | 53.36 | 50.28 |
| DN201 | PGD | 57.68 | 59.40 | 99.86 | 40.32 | 35.48 | 31.82 | 26.04 |
| | TI | 54.72 | 49.46 | 99.56 | 40.46 | 38.58 | 36.06 | 28.40 |
| | MI | 75.14 | 76.58 | 99.80 | 64.92 | 59.96 | 54.28 | 49.78 |
| | DI | 77.68 | 77.14 | 99.76 | 61.30 | 59.18 | 55.80 | 48.08 |
| | SGM | 82.74 | 86.32 | 99.58 | 71.42 | 65.94 | 58.78 | 55.42 |

C MORE RESULTS FOR PARAMETER TUNING

In this section, we demonstrate the results in black-box attacks with varying decay parameter γ in Figure 7 (PGD) and Figure 8 (FGSM). The numerical results against different target models are shown in Table 7 - 13 (PGD) and Table 14 - 20 (FGSM). It is noteworthy that the performance curve with varying γ has startling similarity within DenseNet or ResNet.

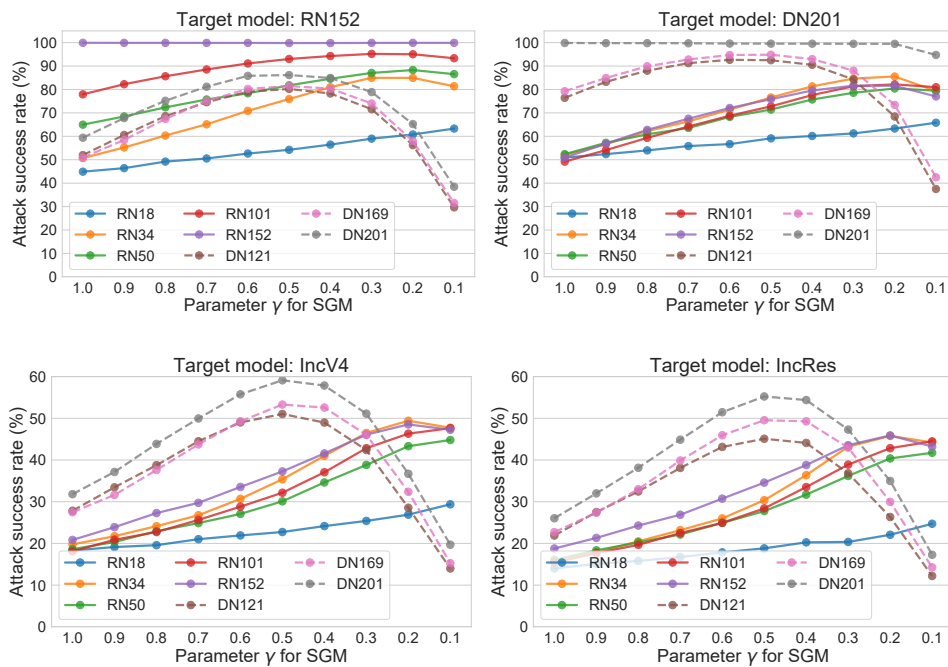


Figure 7: Parameter tuning: the success rates of black-box attacks crafted by 10-step SGM with varying decay parameter $\gamma \in [0.1, 1.0]$. The solid and dash curves represent results on ResNet and DenseNet source models respectively.

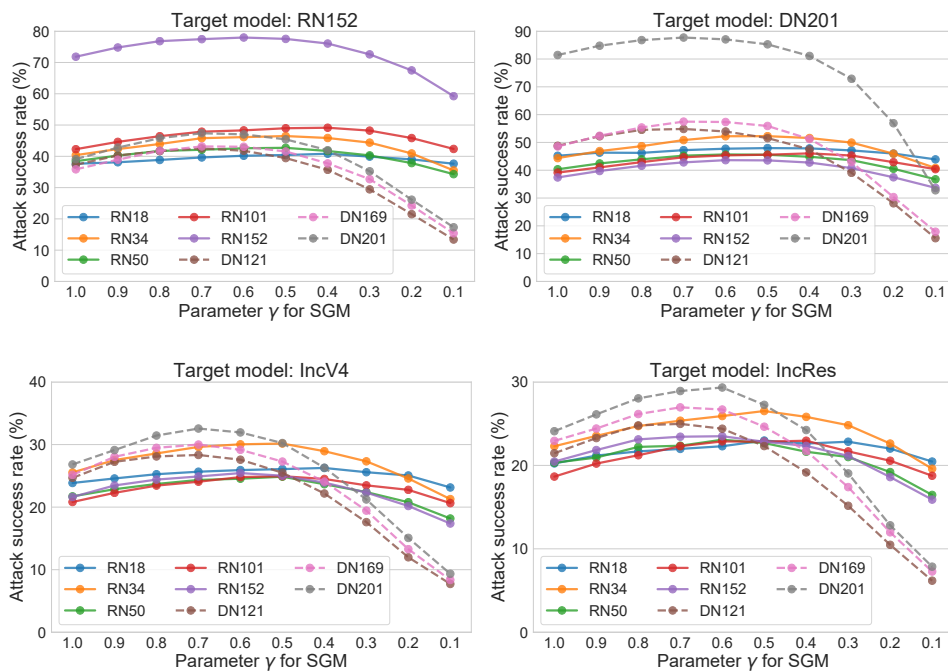


Figure 8: Parameter tuning: the success rates of black-box attacks crafted by single-step SGM with varying decay parameter $\gamma \in [0.1, 1.0]$. The solid and dash curves represent results on ResNet and DenseNet source models respectively.

Table 7: Multi-step (10-step) SGM with different decay parameter γ , against VGG19.

| Target | | VGG19 | | | | | | |
|---------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| $\gamma \setminus$ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
| 1.00 | 59.88 | 54.50 | 52.34 | 43.24 | 45.26 | 61.04 | 54.98 | 57.68 |
| 0.90 | 61.04 | 58.48 | 56.16 | 48.52 | 49.84 | 68.34 | 62.28 | 64.76 |
| 0.80 | 62.24 | 62.40 | 60.68 | 52.36 | 56.16 | 73.90 | 69.06 | 72.32 |
| 0.70 | 63.90 | 66.80 | 64.44 | 56.78 | 61.06 | 79.00 | 74.82 | 78.00 |
| 0.60 | 64.52 | 71.48 | 68.32 | 62.44 | 66.14 | 81.74 | 79.28 | 81.24 |
| 0.50 | 66.08 | 76.78 | 72.50 | 68.22 | 70.18 | 82.86 | 80.82 | 82.82 |
| 0.40 | 66.88 | 80.32 | 76.26 | 72.86 | 75.18 | 79.72 | 79.46 | 80.86 |
| 0.30 | 68.10 | 85.32 | 80.50 | 77.74 | 79.46 | 73.50 | 73.12 | 75.90 |
| 0.20 | 69.82 | 85.72 | 83.40 | 80.02 | 80.54 | 58.98 | 60.14 | 64.26 |
| 0.10 | 71.74 | 83.04 | 83.82 | 80.10 | 78.56 | 35.78 | 38.58 | 43.76 |

Table 8: Multi-step (10-step) SGM with different decay parameter γ , against ResNet-152.

| Target | | RN152 | | | | | | |
|---------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| $\gamma \setminus$ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
| 1.00 | 44.90 | 50.70 | 64.98 | 77.92 | 99.92 | 51.98 | 50.92 | 59.40 |
| 0.90 | 46.40 | 55.22 | 68.48 | 82.24 | 99.90 | 60.58 | 58.44 | 67.78 |
| 0.80 | 49.20 | 60.32 | 72.36 | 85.64 | 99.90 | 68.60 | 67.42 | 75.18 |
| 0.70 | 50.54 | 65.12 | 75.72 | 88.54 | 99.88 | 74.56 | 75.12 | 81.18 |
| 0.60 | 52.64 | 70.88 | 78.36 | 91.10 | 99.88 | 79.02 | 80.24 | 85.82 |
| 0.50 | 54.24 | 75.92 | 81.78 | 93.02 | 99.82 | 80.22 | 81.44 | 86.16 |
| 0.40 | 56.46 | 80.98 | 84.58 | 94.28 | 99.82 | 78.28 | 80.38 | 84.90 |
| 0.30 | 59.00 | 84.96 | 87.06 | 95.20 | 99.88 | 71.58 | 74.02 | 78.88 |
| 0.20 | 60.78 | 84.94 | 88.30 | 95.04 | 99.90 | 56.22 | 57.90 | 65.24 |
| 0.10 | 63.30 | 81.38 | 86.54 | 93.34 | 99.90 | 29.66 | 31.52 | 38.46 |

Table 9: Multi-step (10-step) SGM with different decay parameter γ , against DenseNet-201.

| Target | | DN201 | | | | | | |
|---------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| $\gamma \setminus$ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
| 1.00 | 50.90 | 52.00 | 52.38 | 49.18 | 50.90 | 76.42 | 79.28 | 99.86 |
| 0.90 | 52.40 | 56.62 | 57.22 | 54.12 | 56.86 | 83.26 | 84.90 | 99.78 |
| 0.80 | 53.98 | 62.22 | 60.86 | 59.32 | 62.76 | 88.02 | 89.92 | 99.78 |
| 0.70 | 55.84 | 66.38 | 63.56 | 64.20 | 67.52 | 91.24 | 92.80 | 99.72 |
| 0.60 | 56.68 | 71.38 | 68.28 | 68.80 | 72.10 | 92.66 | 94.72 | 99.60 |
| 0.50 | 59.12 | 76.58 | 71.44 | 72.78 | 75.84 | 92.50 | 94.84 | 99.56 |
| 0.40 | 60.14 | 81.34 | 75.72 | 77.62 | 79.56 | 90.52 | 92.98 | 99.54 |
| 0.30 | 61.24 | 84.66 | 78.58 | 81.38 | 81.62 | 84.30 | 88.04 | 99.50 |
| 0.20 | 63.32 | 85.54 | 80.42 | 82.18 | 81.58 | 68.44 | 73.42 | 99.48 |
| 0.10 | 65.80 | 79.36 | 79.50 | 80.98 | 77.02 | 37.50 | 42.50 | 94.74 |

Table 10: Multi-step (10-step) SGM with different decay parameter γ , against SENet-154.

| Target | | SE154 | | | | | | |
|---------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| $\gamma \setminus$ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
| 1.00 | 22.34 | 24.98 | 26.78 | 26.10 | 29.90 | 33.22 | 34.16 | 40.32 |
| 0.90 | 23.18 | 28.24 | 29.90 | 28.50 | 34.10 | 40.64 | 41.22 | 47.40 |
| 0.80 | 24.16 | 31.14 | 33.98 | 32.04 | 38.24 | 48.90 | 49.08 | 55.54 |
| 0.70 | 25.72 | 35.04 | 37.06 | 35.52 | 41.50 | 55.96 | 56.56 | 62.66 |
| 0.60 | 27.00 | 39.14 | 40.56 | 39.10 | 46.02 | 62.38 | 63.88 | 68.74 |
| 0.50 | 28.56 | 44.10 | 45.20 | 43.78 | 51.82 | 64.78 | 67.88 | 71.72 |
| 0.40 | 29.84 | 51.16 | 49.90 | 48.42 | 56.40 | 63.68 | 66.50 | 71.12 |
| 0.30 | 31.76 | 56.54 | 55.38 | 54.24 | 60.54 | 58.80 | 61.78 | 66.12 |
| 0.20 | 33.25 | 58.38 | 59.28 | 57.28 | 61.84 | 45.76 | 48.84 | 55.34 |
| 0.10 | 36.12 | 56.18 | 60.82 | 56.72 | 58.18 | 26.00 | 28.88 | 33.90 |

Table 11: Multi-step (10-step) SGM with different decay parameter γ , against InceptionV3.

| Target | | IncV3 | | | | | | |
|---------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| $\gamma \setminus$ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
| 1.00 | 23.24 | 24.48 | 22.80 | 22.96 | 26.16 | 30.74 | 30.90 | 35.48 |
| 0.90 | 24.10 | 27.38 | 25.14 | 25.28 | 29.08 | 37.44 | 36.76 | 41.94 |
| 0.80 | 25.64 | 30.54 | 28.16 | 27.96 | 32.38 | 43.74 | 43.82 | 49.52 |
| 0.70 | 26.74 | 33.84 | 30.76 | 31.14 | 35.52 | 50.20 | 50.74 | 56.02 |
| 0.60 | 27.84 | 37.60 | 33.18 | 35.08 | 39.76 | 55.66 | 57.08 | 62.78 |
| 0.50 | 28.92 | 43.44 | 36.78 | 38.48 | 44.84 | 57.58 | 60.40 | 65.38 |
| 0.40 | 30.44 | 49.38 | 41.36 | 45.00 | 49.82 | 56.30 | 59.30 | 64.68 |
| 0.30 | 32.60 | 56.08 | 46.96 | 51.26 | 53.80 | 49.38 | 53.30 | 57.98 |
| 0.20 | 33.66 | 59.18 | 51.84 | 55.52 | 56.50 | 35.76 | 39.84 | 44.54 |
| 0.10 | 37.34 | 56.76 | 53.96 | 57.04 | 55.02 | 18.32 | 20.44 | 24.74 |

Table 12: Multi-step (10-step) SGM with different decay parameter γ , against InceptionV4.

| Target | | IncV3 | | | | | | |
|---------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| $\gamma \setminus$ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
| 1.00 | 18.30 | 19.76 | 18.62 | 18.10 | 20.84 | 27.84 | 27.48 | 31.82 |
| 0.90 | 19.14 | 21.74 | 20.34 | 20.88 | 23.90 | 33.46 | 31.64 | 37.14 |
| 0.80 | 19.58 | 24.14 | 22.94 | 22.74 | 27.28 | 38.76 | 37.60 | 43.86 |
| 0.70 | 21.02 | 26.76 | 24.88 | 25.64 | 29.74 | 44.50 | 43.74 | 49.96 |
| 0.60 | 21.92 | 30.70 | 27.08 | 28.82 | 33.54 | 49.02 | 49.28 | 55.76 |
| 0.50 | 22.74 | 35.38 | 30.10 | 32.16 | 37.26 | 51.04 | 53.30 | 59.12 |
| 0.40 | 24.16 | 41.00 | 34.62 | 37.06 | 41.58 | 48.98 | 52.58 | 57.86 |
| 0.30 | 25.40 | 46.48 | 38.82 | 42.84 | 46.08 | 42.40 | 46.04 | 51.12 |
| 0.20 | 26.90 | 49.42 | 43.32 | 46.32 | 48.54 | 28.56 | 32.38 | 36.70 |
| 0.10 | 29.38 | 47.74 | 44.80 | 47.64 | 47.26 | 14.00 | 15.26 | 19.70 |

Table 13: Multi-step SGM with different decay parameter γ , against InceptionResnet-V2.

| Target | | IncV3 | | | | | | |
|---------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| $\gamma \setminus$ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
| 1.00 | 14.00 | 16.02 | 16.04 | 15.40 | 18.78 | 21.90 | 22.70 | 26.04 |
| 0.90 | 15.10 | 17.88 | 18.34 | 17.68 | 21.32 | 27.52 | 27.38 | 32.00 |
| 0.80 | 15.80 | 20.50 | 20.24 | 19.66 | 24.28 | 32.46 | 33.00 | 38.14 |
| 0.70 | 16.74 | 23.18 | 22.14 | 22.48 | 26.88 | 38.10 | 39.88 | 44.88 |
| 0.60 | 17.86 | 26.00 | 24.92 | 24.94 | 30.74 | 43.14 | 45.94 | 51.46 |
| 0.50 | 18.82 | 30.36 | 27.76 | 28.40 | 34.58 | 45.10 | 49.54 | 55.24 |
| 0.40 | 20.26 | 36.38 | 31.70 | 33.52 | 38.80 | 44.10 | 49.28 | 54.38 |
| 0.30 | 20.36 | 43.20 | 36.18 | 38.90 | 43.56 | 36.82 | 42.98 | 47.30 |
| 0.20 | 22.08 | 45.78 | 40.38 | 42.84 | 45.90 | 26.34 | 29.96 | 34.98 |
| 0.10 | 24.72 | 44.22 | 41.74 | 44.48 | 43.24 | 12.20 | 14.26 | 17.28 |

Table 14: Single-step SGM with different decay parameter γ , against VGG19.

| Target | | VGG19 | | | | | | |
|---------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| $\gamma \setminus$ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
| 1.0 | 54.92 | 52.20 | 45.28 | 42.12 | 41.36 | 51.32 | 48.24 | 49.84 |
| 0.9 | 55.62 | 54.44 | 46.70 | 44.48 | 42.86 | 53.54 | 51.06 | 52.58 |
| 0.8 | 55.68 | 55.60 | 48.24 | 45.86 | 44.86 | 55.16 | 53.26 | 54.84 |
| 0.7 | 56.50 | 56.50 | 49.34 | 47.82 | 45.70 | 55.80 | 54.06 | 56.70 |
| 0.6 | 56.60 | 57.54 | 49.90 | 49.40 | 47.46 | 54.96 | 55.12 | 56.66 |
| 0.5 | 56.62 | 58.22 | 50.70 | 50.26 | 47.66 | 52.78 | 53.56 | 55.86 |
| 0.4 | 56.66 | 57.24 | 50.60 | 50.64 | 48.20 | 48.68 | 50.40 | 53.30 |
| 0.3 | 56.32 | 56.84 | 50.12 | 51.12 | 47.16 | 42.26 | 45.46 | 46.90 |
| 0.2 | 55.46 | 54.38 | 48.72 | 51.36 | 45.60 | 34.90 | 37.10 | 40.22 |
| 0.1 | 53.60 | 50.52 | 47.70 | 49.94 | 43.48 | 25.56 | 27.04 | 29.50 |

Table 15: Single-step SGM with different decay parameter γ , against ResNet-152.

| Target | | RN152 | | | | | | |
|---------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| $\gamma \setminus$ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
| 1.00 | 37.62 | 40.38 | 38.44 | 42.28 | 71.84 | 37.30 | 35.84 | 39.10 |
| 0.90 | 38.10 | 42.32 | 40.24 | 44.62 | 74.80 | 40.28 | 38.92 | 42.78 |
| 0.80 | 38.84 | 43.94 | 41.72 | 46.44 | 76.82 | 41.66 | 41.78 | 45.82 |
| 0.70 | 39.64 | 45.80 | 42.10 | 47.88 | 77.46 | 42.54 | 43.16 | 47.38 |
| 0.60 | 40.18 | 46.14 | 42.58 | 48.30 | 77.98 | 41.78 | 43.06 | 47.02 |
| 0.50 | 40.42 | 46.48 | 42.72 | 48.98 | 77.54 | 39.44 | 41.56 | 45.42 |
| 0.40 | 40.82 | 45.86 | 41.70 | 49.14 | 76.06 | 35.68 | 37.74 | 41.96 |
| 0.30 | 39.90 | 44.38 | 40.30 | 48.22 | 72.60 | 29.44 | 32.70 | 35.22 |
| 0.20 | 39.02 | 40.92 | 37.84 | 45.84 | 67.48 | 21.52 | 24.26 | 26.14 |
| 0.10 | 37.64 | 35.50 | 34.30 | 42.36 | 59.24 | 13.40 | 15.40 | 17.32 |

Table 16: Single-step SGM with different decay parameter γ , against DenseNet-201.

| Target | | DN201 | | | | | | |
|---------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| $\gamma \setminus$ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
| 1.00 | 45.20 | 44.36 | 40.34 | 39.14 | 37.40 | 48.80 | 48.54 | 81.44 |
| 0.90 | 46.32 | 46.88 | 42.44 | 41.02 | 39.72 | 52.20 | 52.40 | 84.78 |
| 0.80 | 46.24 | 48.66 | 43.92 | 42.92 | 41.60 | 54.52 | 55.38 | 86.80 |
| 0.70 | 47.22 | 50.84 | 45.24 | 44.66 | 42.82 | 54.84 | 57.48 | 87.72 |
| 0.60 | 47.74 | 52.22 | 45.68 | 45.34 | 43.62 | 53.94 | 57.34 | 87.06 |
| 0.50 | 47.96 | 52.26 | 45.56 | 45.58 | 43.54 | 51.56 | 55.92 | 85.28 |
| 0.40 | 47.86 | 51.58 | 44.84 | 46.06 | 42.76 | 47.50 | 51.22 | 81.10 |
| 0.30 | 47.12 | 49.94 | 43.62 | 45.28 | 40.76 | 39.06 | 42.94 | 72.90 |
| 0.20 | 46.00 | 45.94 | 40.52 | 42.90 | 37.46 | 28.10 | 30.34 | 56.90 |
| 0.10 | 43.90 | 40.80 | 36.74 | 40.38 | 33.66 | 15.56 | 17.88 | 32.78 |

Table 17: Single-step SGM with different decay parameter γ , against SENet-154.

| Target | | SE154 | | | | | | |
|---------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| $\gamma \setminus$ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
| 1.00 | 28.22 | 29.80 | 26.90 | 25.72 | 26.74 | 31.00 | 31.16 | 35.46 |
| 0.90 | 28.98 | 31.38 | 28.54 | 27.08 | 28.30 | 32.82 | 34.34 | 38.42 |
| 0.80 | 29.50 | 33.08 | 29.04 | 28.98 | 29.84 | 35.02 | 37.46 | 40.98 |
| 0.70 | 29.24 | 34.62 | 30.02 | 29.98 | 30.74 | 36.20 | 38.04 | 42.84 |
| 0.60 | 29.60 | 35.30 | 30.56 | 30.74 | 31.00 | 35.38 | 38.72 | 42.48 |
| 0.50 | 30.06 | 35.12 | 30.82 | 30.58 | 31.12 | 33.46 | 37.28 | 40.18 |
| 0.40 | 29.74 | 34.18 | 30.66 | 30.58 | 29.78 | 30.04 | 34.14 | 37.10 |
| 0.30 | 29.34 | 32.36 | 29.34 | 29.42 | 28.00 | 24.78 | 28.48 | 31.14 |
| 0.20 | 28.44 | 29.74 | 27.68 | 28.14 | 25.78 | 17.40 | 20.86 | 23.38 |
| 0.10 | 26.94 | 26.26 | 24.88 | 26.02 | 22.40 | 10.30 | 12.90 | 14.54 |

Table 18: Single-step SGM with different decay parameter γ , against Inception-V3.

| Target | | IncV3 | | | | | | |
|---------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| $\gamma \setminus$ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
| 1.00 | 29.22 | 30.22 | 26.40 | 24.34 | 25.74 | 28.50 | 29.32 | 31.74 |
| 0.90 | 30.06 | 31.80 | 27.58 | 25.90 | 27.46 | 31.04 | 32.00 | 34.70 |
| 0.80 | 30.82 | 33.80 | 28.18 | 27.94 | 28.58 | 33.16 | 33.98 | 37.10 |
| 0.70 | 31.02 | 34.74 | 29.12 | 28.94 | 29.30 | 33.26 | 35.34 | 38.36 |
| 0.60 | 31.52 | 35.42 | 29.32 | 29.60 | 29.70 | 32.34 | 35.70 | 37.80 |
| 0.50 | 32.16 | 35.48 | 29.16 | 30.46 | 29.18 | 30.70 | 33.38 | 36.14 |
| 0.40 | 32.08 | 34.90 | 29.00 | 30.30 | 28.56 | 27.66 | 30.24 | 32.06 |
| 0.30 | 31.46 | 33.62 | 28.16 | 29.80 | 27.48 | 22.88 | 24.84 | 27.26 |
| 0.20 | 30.70 | 31.14 | 26.90 | 27.90 | 25.06 | 15.74 | 18.16 | 19.80 |
| 0.10 | 29.50 | 27.36 | 24.46 | 26.08 | 22.14 | 10.16 | 11.86 | 12.90 |

Table 19: Single-step SGM with different decay parameter γ , against Inception-V4.

| Target | | Inc4 | | | | | | |
|-------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| γ \ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
| 1.00 | 23.86 | 25.54 | 21.72 | 20.82 | 21.62 | 24.68 | 25.16 | 26.82 |
| 0.90 | 24.58 | 27.52 | 22.86 | 22.28 | 23.46 | 27.24 | 28.04 | 29.12 |
| 0.80 | 25.26 | 28.58 | 23.70 | 23.44 | 24.40 | 28.10 | 29.48 | 31.44 |
| 0.70 | 25.66 | 29.62 | 24.34 | 24.06 | 24.90 | 28.34 | 30.00 | 32.56 |
| 0.60 | 25.92 | 30.04 | 24.54 | 24.80 | 25.44 | 27.58 | 29.16 | 31.94 |
| 0.50 | 26.08 | 30.16 | 24.86 | 24.94 | 25.02 | 25.50 | 27.28 | 30.24 |
| 0.40 | 26.26 | 28.94 | 23.66 | 24.50 | 23.92 | 22.18 | 24.08 | 26.30 |
| 0.30 | 25.58 | 27.32 | 22.40 | 23.48 | 22.32 | 17.60 | 19.44 | 21.22 |
| 0.20 | 25.06 | 24.58 | 20.78 | 22.76 | 20.18 | 11.98 | 13.28 | 15.08 |
| 0.10 | 23.14 | 21.28 | 18.18 | 20.64 | 17.38 | 7.72 | 8.36 | 9.34 |

Table 20: Single-step SGM with different decay parameter γ , against InceptionResnet-V2.

| Target | | IncRes | | | | | | |
|-------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| γ \ Source | RN18 | RN34 | RN50 | RN101 | RN152 | DN121 | DN169 | DN201 |
| 1.00 | 20.24 | 22.30 | 20.32 | 18.66 | 20.46 | 21.46 | 22.94 | 24.10 |
| 0.90 | 21.18 | 23.54 | 20.94 | 20.22 | 21.84 | 23.28 | 24.42 | 26.12 |
| 0.80 | 21.68 | 24.72 | 22.22 | 21.22 | 23.12 | 24.80 | 26.16 | 28.04 |
| 0.70 | 21.98 | 25.36 | 22.36 | 22.30 | 23.44 | 24.98 | 26.96 | 28.92 |
| 0.60 | 22.30 | 25.92 | 23.08 | 22.88 | 23.50 | 24.40 | 26.70 | 29.34 |
| 0.50 | 22.98 | 26.52 | 22.68 | 22.84 | 22.84 | 22.32 | 24.64 | 27.26 |
| 0.40 | 22.62 | 25.82 | 21.64 | 22.96 | 22.34 | 19.16 | 21.72 | 24.24 |
| 0.30 | 22.84 | 24.82 | 21.00 | 21.68 | 21.14 | 15.16 | 17.40 | 19.04 |
| 0.20 | 22.02 | 22.62 | 19.20 | 20.56 | 18.60 | 10.48 | 11.98 | 12.82 |
| 0.10 | 20.46 | 19.60 | 16.46 | 18.76 | 15.90 | 6.18 | 7.28 | 7.86 |