
Boosting the Transferability of Adversarial Attacks with Reverse Adversarial Perturbation

Zeyu Qin^{1*}, Yanbo Fan^{2*}, Yi Liu¹, Li Shen³, Yong Zhang², Jue Wang², Baoyuan Wu^{1†}

¹School of Data Science, Shenzhen Research Institute of Big Data,
The Chinese University of Hong Kong, Shenzhen

²Tencent AI Lab

³JD Explore Academy

{zeyu6181136, fanyanbo0124, yiliuhk2000}@gmail.com

{mathshenli, zhangyong201303, arphid}@gmail.com

wubaoyuan@cuhk.edu.cn

Abstract

Deep neural networks (DNNs) have been shown to be vulnerable to adversarial examples, which can produce erroneous predictions by injecting imperceptible perturbations. In this work, we study the transferability of adversarial examples, which is significant due to its threat to real-world applications where model architecture or parameters are usually unknown. Many existing works reveal that the adversarial examples are likely to overfit the surrogate model that they are generated from, limiting its transfer attack performance against different target models. To mitigate the overfitting of the surrogate model, we propose a novel attack method, dubbed *reverse adversarial perturbation* (RAP). Specifically, instead of minimizing the loss of a single adversarial point, we advocate seeking adversarial example located at a region with unified low loss value, by injecting the worst-case perturbation (*i.e.*, the reverse adversarial perturbation) for each step of the optimization procedure. The adversarial attack with RAP is formulated as a min-max bi-level optimization problem. By integrating RAP into the iterative process for attacks, our method can find more stable adversarial examples which are less sensitive to the changes of decision boundary, mitigating the overfitting of the surrogate model. Comprehensive experimental comparisons demonstrate that RAP can significantly boost adversarial transferability. Furthermore, RAP can be naturally combined with many existing black-box attack techniques, to further boost the transferability. When attacking a real-world image recognition system, *i.e.*, Google Cloud Vision API, we obtain 22% performance improvement of targeted attacks over the compared method. Our codes are available at: https://github.com/SCLBD/Transfer_attack_RAP.

1 Introduction

Deep neural networks (DNNs) have been successfully applied in many safety-critical tasks, such as autonomous driving, face recognition and verification, *etc.* However, it has been shown that DNN models are vulnerable to adversarial examples [9, 12, 27, 32, 35, 44, 45, 50], which are indistinguishable from natural examples but make a model produce erroneous predictions. For real-world applications, the DNN models are often hidden from users. Therefore, the attackers need to generate the adversarial examples under black-box setting where they do not know any information of the target model [2, 3, 18, 32]. For black-box setting, the adversarial transferability matters since

*Equal contribution. †Corresponding author.

‡This work is done when Zeyu Qin is a research intern at Tencent AI Lab.

it can allow the attackers to attack target models by using adversarial examples generated on the surrogate models. Therefore, learning how to generate adversarial examples with high transferability has gained more attentions in the literature [5, 10, 14, 23, 26, 38, 48].

Under white-box setting where the complete information of the attacked model (*e.g.*, architecture and parameters) is available, the gradient-based attacks such as PGD [27] have demonstrated good attack performance. However, they often exhibit the poor transferability [5, 48], *i.e.*, the adversarial example x^{adv} generated from the surrogate model \mathcal{M}^S performs poorly against different target models \mathcal{M}^T . The previous works attribute that to the overfitting of adversarial examples to the surrogate models [5, 24, 48]. Figure 1 (b) gives an illustration. The PGD attack aims to find an adversarial point x^{pgd} with minimal attack loss, while doesn't consider the attack loss of the neighborhood regions round x^{pgd} . Due to the highly non-convex of deep models, when x^{pgd} locates at a sharply local minimum, a slight change on model parameters of \mathcal{M}^S could cause a large increase of the attack loss, making x^{pgd} fail to attack the perturbed model.

Many techniques have been proposed to mitigate the overfitting and improve the transferability, including input transformation [6, 48], gradient calibration [14], feature-level attacks [17], and generative models [30], etc. However, there still exists a large gap of attack performance between the transfer setting and the ideal white-box setting, especially for targeted attack, requiring more efforts for boosting the transferability.

In this work, we propose a novel attack method called *reverse adversarial perturbation* (RAP) to alleviate the overfitting of the surrogate model and boost the transferability of adversarial examples. We encourage that x^{adv} is not only of low attack loss but also locates at a local flat region, *i.e.*, the points within the local neighborhood region around x^{adv} should also be of low loss values. Figure 1 (b) illustrates the difference between the sharp local minimum and flat local minimum. When the model parameter of \mathcal{M}^S has some slight changes, the variation of the attack loss *w.r.t.* the flat local minimum is less than that of the sharp one. Therefore, the flat local minimum is less sensitive to the changes of decision boundary. To achieve this goal, we formulate a min-max bi-level optimization problem. The inner maximization aims to find the worst-case perturbation (*i.e.*, that with the largest attack loss, and this is why we call it reverse adversarial perturbation) within the local region around the current adversarial example, which can be solved by the projected gradient ascent algorithm. Then, the outer minimization will update the adversarial example to find a new point added with the provided reverse perturbation that leads to lower attack loss. Figure 1 (a) provides an illustration of the optimization process. For t -th iteration and x^t , RAP first finds the point $x^t + n^{rap}$ with max attack loss within the neighborhood region of x^t . Then it updates x^t with the gradient calculated by minimizing the attack loss *w.r.t.* $x^t + n^{rap}$. Compared to directly adopting the gradient at x^t , RAP could help escape from the sharp local minimum and pursue a relatively flat local minimum. Besides, we design a late-start variant of RAP (RAP-LS) to further boost the attack effectiveness and efficiency, which doesn't insert the reverse perturbation into the optimization procedure in the early stage. Moreover, from the technical perspective, since the proposed RAP method only introduces one specially designed perturbation onto adversarial attacks, one notable advantage of RAP is that it can be naturally combined with many existing black-box attack techniques to further boost the transferability. For example, when combined with different input transformations (*e.g.*, the random resizing and padding in Diverse Input [48]), our RAP method consistently outperforms the counterparts by a clear margin.

Our main contributions are three-fold: **1)** Based on a novel perspective, the flatness of loss landscape for adversarial examples, we propose a novel adversarial attack method RAP that encourages both the adversarial example and its neighborhood region be of low loss value; **2)** we present a vigorous experimental study and show that RAP can significantly boost the adversarial transferability on both untargeted and targeted attacks for various networks also containing defense models; **3)** we demonstrate that RAP can be easily combined with existing transfer attack techniques and outperforms the state-of-the-art performance by a large margin.

2 Related Work

The black-box attacks can be categorized into two categories: 1) *query-based attacks* that conduct the attack based on the feedback of iterative queries to target models, and 2) *transfer attacks* that use the adversarial examples generated on some surrogate models to attack the target models. In this work,

we focus on the transfer attacks. For surrogate models, existing attack algorithms such as FGSM [12] and I-FGSM [21] could achieve good attack performance. However, they often overfit the surrogate models and thus exhibit poor transferability. Recently, many works have been proposed to generate more transferable adversarial examples [5, 6, 13, 14, 17, 19, 20, 24, 39, 40, 41, 43, 48], which we briefly summarize as below.

Input transformation: Data augmentation, which has been shown to be effective in improving model generalization, has also been studied to boost the adversarial transferability, such as randomly resizing and padding [48], randomly scaling [24], and adversarial mixup [41]. In addition, the work of Dong et al. [6] uses a set of translated images to compute gradient and get the better performance against defense models. Expectation of Transformation (EOT) method [1] synthesizes adversarial examples over a chosen distribution of transformations to enhance its adversarial transferability. **Gradient modification:** Instead of the I-FGSM, the work of Dong et al. [5] integrates momentum into the updating strategy. And Lin et al. [24] uses the Nesterov accelerated gradient to boost the transferability. The work of Wang and He [40] aims to find a more stable gradient direction by tuning the variance of each gradient step. There are also some model-specific designs to boost the adversarial transferability. For example, Wu et al. [43] found that the gradient of skip connections is more crucial to generate more transferable attacks. The work of Guo et al. [14] proposed LinBP to utilize more gradient of skip connections during the back-propagation. However, these methods tend to be specific to a particular model architecture, such as skip connection, and it is nontrivial to extend the findings to other architectures or modules. **Intermediate feature attack:** Meanwhile, Huang et al. [17], Inkawhich et al. [19, 20] proposed to exploit feature space constraints to generate more transferable attacks. Yet they need to identify the best performing intermediate layers or train one-vs-all binary classifiers for all attacked classes. Recently, Zhao et al. [49] find iterative attacks with much more iterations and logit loss can achieve relatively high targeted transferability and exceed the feature-based attacks. **Generative models:** In addition, there have been some methods utilizing the generative models to generate the adversarial perturbations [28, 30, 31]. For example, the work of Naseer et al. [30] proposed to train a generative model to match the distributions of source and target class, so as to increase the targeted transferability. However, the learning of the perturbation generator is nontrivial, especially on large-scale datasets.

In summary, the current performance of transfer attacks is still unsatisfactory, especially for targeted attacks. In this work, we study adversarial transferability from the perspective of the flatness of adversarial examples. We find that adversarial examples located at flat local minimum will be more transferable than those at sharp local minimum and propose a novel algorithm to find adversarial example that locates at flat local minimum.

3 Methodology

3.1 Preliminaries of Transfer Adversarial Attack

Given a benign sample $(\mathbf{x}, y) \in (\mathcal{X}, \mathcal{Y})$, the procedure of transfer adversarial attack is firstly constructing the adversarial example \mathbf{x}^{adv} within the neighborhood region $\mathcal{B}_\epsilon(\mathbf{x}) = \{\mathbf{x}' : \|\mathbf{x}' - \mathbf{x}\|_p \leq \epsilon\}$ by attacking the white-box surrogate model $\mathcal{M}^s(\mathbf{x}; \boldsymbol{\theta}) : \mathcal{X} \rightarrow \mathcal{Y}$, then transferring \mathbf{x}^{adv} to directly attack the black-box target model $\mathcal{M}^t(\mathbf{x}; \boldsymbol{\phi}) : \mathcal{X} \rightarrow \mathcal{Y}$. The attack goal is to mislead the target model, *i.e.*, $\mathcal{M}^t(\mathbf{x}^{adv}; \boldsymbol{\phi}) \neq y$ (untargeted attack), or $\mathcal{M}^t(\mathbf{x}^{adv}; \boldsymbol{\phi}) = y_t$ (targeted attack) with $y_t \in \mathcal{Y}$ indicting the target label. Taking the target attack as example, the general formulation of many existing transfer attack methods can be written as follows:

$$\min_{\mathbf{x}^{adv} \in \mathcal{B}_\epsilon(\mathbf{x})} \mathcal{L}(\mathcal{M}^s(\mathcal{G}(\mathbf{x}^{adv}); \boldsymbol{\theta}), y_t). \quad (1)$$

The loss function \mathcal{L} is often set as the cross entropy (CE) loss [48] or the logit loss [49], which will be specified in later experiments. Besides, the formulation of untargeted attack can be easily obtained by replacing the loss function \mathcal{L} and y_t by $-\mathcal{L}$ and y , respectively.

Since \mathcal{M}^s is white-box, if $\mathcal{G}(\cdot)$ is set as the identity function, then any off-the-shelf white-box adversarial attack method can be adopted to solve Problem (1), such as I-FSGM [21], MI-FGSM [5], *etc.* Meanwhile, existing works have designed different $\mathcal{G}(\cdot)$ functions and developed the corresponding optimization algorithms, to boost the adversarial transferability between surrogate and target models. For example, $\mathcal{G}(\cdot)$ is specified as random resizing and padding (DI) [48], translation transformation (TI) [6], scale transformation (SI) [24], and adversarial mixup (Admix) [41].

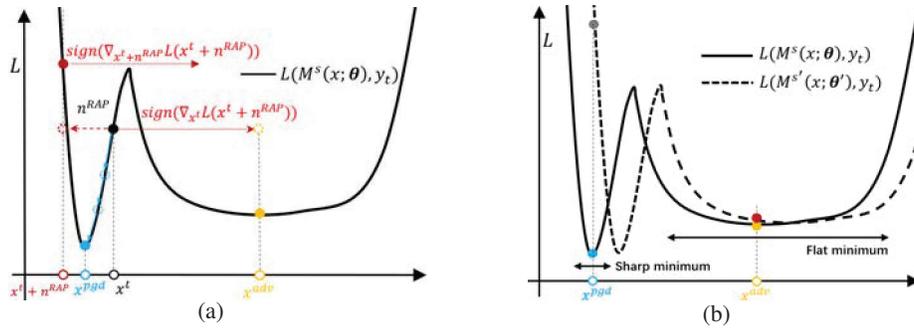


Figure 1: These two plots are schematic diagrams in 1D space. The x-axis means the value of input \mathbf{x} . The y-axis means the value of attack loss function \mathcal{L} . (a) Illustration of our attack method and the original PGD attack. (b) Illustration of attack loss landscape of \mathcal{M}^S and $\mathcal{M}^{S'}$. $\mathcal{M}^{S'}$ denotes a slight change on the model parameters of \mathcal{M}^S . The blue and yellow dots correspond to attacks located at different local minima on \mathcal{M}^S , respectively. The gray and red points are their counterparts on $\mathcal{M}^{S'}$.

3.2 Reverse Adversarial Perturbation

As discussed above, although having good performance for the white-box setting, the adversarial examples generated from \mathcal{M}^S exist poor adversarial transferability on \mathcal{M}^T , especially for targeted attacks. The previous works attribute this issue to the overfitting of adversarial attack to \mathcal{M}^S [4, 5, 6, 38, 46]. As shown in Figure 1 (b), when \mathbf{x}^{pgd} locates at a sharp local minimum, it is not stable and is sensitive to changes of \mathcal{M}^S . When having some changes on model parameters, \mathbf{x}^{pgd} could results in a high attack loss against $\mathcal{M}^{S'}$ and lead to a failure attack.

To mitigating the overfitting to \mathcal{M}^S , we advocate to find \mathbf{x}^{adv} located at flat local region. That means we encourage that not only \mathbf{x}^{adv} itself has low loss value, but also the points in the vicinity of \mathbf{x}^{adv} have similarly low loss values.

To this end, we propose to minimize the maximal loss value within a local neighborhood region around the adversarial example \mathbf{x}^{adv} . The maximal loss is implemented by perturbing \mathbf{x}^{adv} to maximize the attack loss, named *Reverse Adversarial Perturbation* (RAP). By inserting the RAP into the formulation (1), we aim to solve the following problem,

$$\min_{\mathbf{x}^{adv} \in \mathcal{B}_\epsilon(\mathbf{x})} \mathcal{L}(\mathcal{M}^s(\mathcal{G}(\mathbf{x}^{adv} + \mathbf{n}^{rap}); \theta), y_t), \quad (2)$$

where

$$\mathbf{n}^{rap} = \arg \max_{\|\mathbf{n}^{rap}\|_\infty \leq \epsilon_n} \mathcal{L}(\mathcal{M}^s(\mathbf{x}^{adv} + \mathbf{n}^{rap}; \theta), y_t), \quad (3)$$

with \mathbf{n}^{rap} indicating the RAP, and ϵ_n defining its search region. The above formulations Equation (2) and Equation (3) correspond to the targeted attack, and the corresponding untargeted formulations can be easily obtained by replacing the loss function \mathcal{L} and y_t by $-\mathcal{L}$ and y , respectively.

It is a min-max bi-level optimization problem [25], and can be solved by iteratively optimizing the inner maximization and the outer minimization problem. Specifically, in each iteration, given \mathbf{x}^{adv} , the inner maximization *w.r.t.* \mathbf{n}^{rap} is solved by the projected gradient ascent algorithm:

$$\mathbf{n}^{rap} \leftarrow \mathbf{n}^{rap} + \alpha_n \cdot \text{sign}(\nabla_{\mathbf{n}^{rap}} \mathcal{L}(\mathcal{M}^s(\mathbf{x}^{adv} + \mathbf{n}^{rap}; \theta), y_t)). \quad (4)$$

The above update is conducted by T steps, and $\alpha_n = \frac{\epsilon_n}{T}$. Then, given \mathbf{n}^{rap} , the outer minimization *w.r.t.* \mathbf{x}^{adv} can be solved by any off-the-shelf algorithm that is developed for solving Equation (1). For example, it can be undated by one step projected gradient descent, as follows:

$$\mathbf{x}^{adv} \leftarrow \text{Clip}_{\mathcal{B}_\epsilon(\mathbf{x})}[\mathbf{x}^{adv} - \alpha \cdot \text{sign}(\nabla_{\mathbf{x}^{adv}} \mathcal{L}(\mathcal{M}^s(\mathcal{G}(\mathbf{x}^{adv} + \mathbf{n}^{rap}); \theta), y_t))], \quad (5)$$

with $\text{Clip}_{\mathcal{B}_\epsilon(\mathbf{x})}(\mathbf{a})$ clipping \mathbf{a} into the neighborhood region $\mathcal{B}_\epsilon(\mathbf{x})$. The overall optimization procedure is summarized in Algorithm 1. Moreover, since the optimization *w.r.t.* \mathbf{x}^{adv} can be implemented by any off-the-shelf algorithm for solving Problem Equation (1), one notable advantage of the proposed RAP is that it can be naturally combined with any one of them, such as the input transformation methods [6, 24, 41, 48].

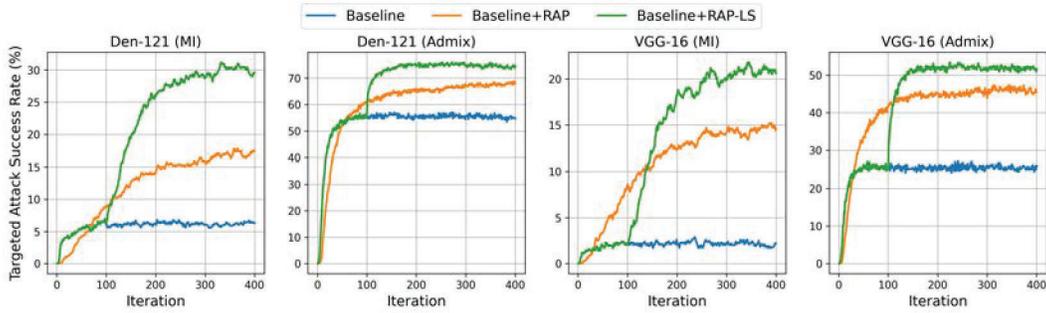


Figure 2: Targeted attack success rate (%) on Dense-121 and VGG-16. We take the Res-50 as the surrogate model and take MI and Admix as baseline methods.

Algorithm 1 Reverse Adversarial Perturbation (RAP) Algorithm

Input: Surrogate model \mathcal{M}^s , benign data (x, y) , target label y_t , loss function \mathcal{L} , transformation \mathcal{G} , the global iteration number K , the late-start iteration number K_{LS} of RAP, as well as hyper-parameters in optimization (specified in later experiments) **Output:** the adversarial example x^{adv}

- 1: Initialize $x^{adv} \leftarrow x$
 - 2: for $k = 1, \dots, K$ do
 - 3: if $k \geq K_{LS}$ then
 - 4: Initialize $n^{rap} \leftarrow 0$
 - 5: for $t = 1, \dots, T$ do
 - 6: Update n^{rap} using Equation (4)
 - 7: Update x^{adv} using Equation (5)
-

A Late-Start (LS) Variant of RAP. In our preliminary experiments, we find that RAP requires more iterations to converge and the performance is slightly lower during the initial iterations, compared to its baseline attack methods. As shown in Figure 2, we combine MI [5] and Admix [41] with RAP, and adopt ResNet-50 as the surrogate model. We take the evaluation on 1000 images from ImageNet (see Sec.4.1). It is observed that the method with RAP (see the orange curves) quickly surpasses its baseline method (see the blue curves) and finally achieves much higher success rate with more iterations, which verify the effect of RAP on enhancing the adversarial transferability. However, it is also observed that the performance of RAP is slightly lower than its baseline method in the early stage. The possible reason is that the early-stage attack is of very weak attack performance to the surrogate model. In this case, it may be waste to pursue better transferable attacks by solving the min-max problem. A better strategy may be only solving the minimization problem Equation (1) in the early stage to quickly achieve the region of relatively high adversarial attack performance, then starting RAP to further enhance the attack performance and transferability simultaneously. This strategy is denoted as RAP with late-start (RAP-LS), whose effect is preliminarily supported by the results shown in Figure 2 (see the green curve) and will be evaluated extensively in later experiments.

3.3 A Closer Look at RAP

To verify whether RAP can help us find a x^{adv} located at the local flat region or not, we use ResNet-50 as surrogate model and conduct the untargeted attacks. We visualize the loss landscape around x^{adv} on \mathcal{M}^S by plotting the loss variations when we move x^{adv} along a random direction with different magnitudes a . The details of the calculation are provided in Appendix. Figure 3 plots the visualizations. We take I-FGSM [21] (denoted as I), MI [5], DI [48], and MI-TI-DI (MTDI) as baselines attacks and combined them with RAP. We can see that comparing to the baselines, RAP could help find x^{adv} located at the flat region.

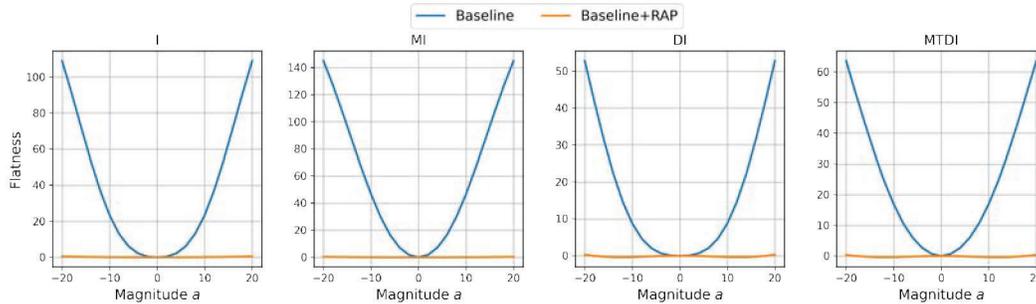


Figure 3: The flatness visualization of untargeted adversarial examples on \mathcal{M}^S . The implementation details are shown in Section B of *Appendix*.

4 Experiments

4.1 Experimental Settings

Dataset and Evaluated Models. We conduct the evaluation on the ImageNet-compatible dataset¹ comprised of 1,000 images. For the surrogate models, we consider the four widely used network architectures: Inception-v3 (Inc-v3) [36], ResNet-50 (Res-50) [15], DenseNet-121 (Dense-121) [16], and VGG-16bn (VGG-16) [34]. For target models, apart from the above models, we also utilize more diverse architectures: Inception-ResNet-v2 (Inc-Res-v2) [37], NASNet-Large (NASNet-L) [51], and ViT-Base/16 (ViT-B/16) [8]. For defense models, we adopt the two widely used ensemble adversarial training (AT) models: adv-Inc-v3 (Inc-v3_{adv}) and ens-adv-Inc-Res-v2 (IncRes-v2_{ens}) [38]. Besides, we also test multi-step AT model [33], input transformation defense [46], feature denoising [47], and purification defense (NRP) [29].

Compared Methods. We adopt I-FGSM [21] (denoted as I), MI [5], TI [6], DI [48], SI [24], Admix [41], VT [40], EMI [42], ILA [17], LinBP [14], Ghost Net [22], and the generative targeted attack method TTP [30]. We also consider the combination of baseline methods, including MI-TI-DI (MTDI), MI-TI-DI-SI (MTDSI), and MI-TI-DI-Admix (MTDAI). Besides, Expectation of Transformation (EOT) method [1] is also a comparable baseline method. We also conduct the comparison of RAP and EOT.

Implementation Details. For untargeted attack, we adopt the Cross Entropy (CE) loss. For targeted attack, apart from CE, we also experiment with the logit loss, where Zhao et al. [49] shows it behaves better for targeted attack. The adversarial perturbation ϵ is restricted by $\ell_\infty = 16/255$. The step size α is set as $2/255$ and number of iteration K is set as 400 for all attacks. In the following, we mainly show the results at $K = 400$ and the results at different value of K are shown in *Appendix*. For RAP, we set K_{LS} as 100 and α_n as $2/255$. We set ϵ_n as $12/255$ for I and TI in untargeted attack and $16/255$ for other attacks in all other settings. The computational cost is shown in Section B of *Appendix*.

Extra Experiments in Appendix. Due to the space limitation, we put extra experiment results in *Appendix*. The comparisons of RAP and EOT, VT, EMI, and Ghost Net methods are shown in Section C in *Appendix*. The evaluation of RAP on stronger defense model, multi-step AT models, NRP, and feature denoising, is shown in Section D. The evaluation of ensemble-model attacks on these diverse network architectures is given in Section E.2.

4.2 The Evaluation of Untargeted Attacks

Baseline Methods. We first evaluate the performance of RAP and RAP-LS with different baseline attacks, including I, MI, DI, TI, SI, and Admix. The results are shown in Table 1. For instance, the ‘MI/ +RAP/ +RAP-LS’ denotes the methods of baseline MI, MI+RAP, and MI+RAP-LS, respectively. RAP achieves the significant improvements for all methods on each target model. For average attack success rate of all target models, RAP outperforms the I and MI by 9.6% and 16.3%, respectively.

¹Publicly available from https://github.com/cleverhans-lab/cleverhans/tree/master/cleverhans_v3.1.0/examples/nips17_adversarial_competition/dataset

Table 1: The **untargeted attack success rate (%) of baseline attacks with RAP**. The results with *CE* loss are reported. The best results are bold and the second best results are underlined.

Attack	ResNet-50 \implies		Inc-v3	DenseNet-121 \implies		
	Dense-121	VGG-16		Res-50	VGG-16	
I / +RAP / +RAP-LS	79.2 / <u>91.5</u> / 91.9	78.0 / <u>91.1</u> / 92.9	34.6 / <u>57.0</u> / 57.2	87.4 / <u>94.2</u> / 94.3	85.1 / <u>91.7</u> / 92.8	46.5 / <u>60.2</u> / 61.1
MI / +RAP / +RAP-LS	85.8 / <u>95.0</u> / 96.1	82.4 / <u>93.9</u> / 94.5	50.3 / <u>75.9</u> / 77.4	90.3 / <u>97.6</u> / 97.9	87.5 / <u>96.0</u> / 97.6	59.3 / <u>80.4</u> / 82.8
TI / +RAP / +RAP-LS	82.0 / <u>94.1</u> / 95.1	81.0 / <u>93.1</u> / 93.3	45.5 / <u>66.1</u> / 67.0	89.6 / <u>94.2</u> / 94.8	87.0 / <u>92.1</u> / 93.3	54.2 / <u>66.7</u> / 70.0
DI / +RAP / +RAP-LS	99.0 / <u>99.6</u> / 99.7	99.0 / <u>99.6</u> / 99.7	57.7 / <u>82.9</u> / 85.0	98.2 / <u>99.6</u> / 99.7	98.1 / <u>99.4</u> / 99.4	67.6 / <u>86.6</u> / 86.9
SI / +RAP / +RAP-LS	94.9 / <u>98.9</u> / 99.7	88.6 / <u>95.7</u> / 97.2	65.9 / <u>79.7</u> / 84.4	95.1 / <u>96.9</u> / 98.8	91.9 / <u>95.0</u> / 97.5	71.6 / <u>83.2</u> / 87.4
Admix / +RAP / +RAP-LS	97.9 / <u>99.6</u> / 99.9	95.8 / <u>97.7</u> / 99.0	77.7 / <u>87.4</u> / 92.6	97.0 / <u>99.0</u> / 99.2	95.6 / <u>97.7</u> / 98.6	82.0 / <u>89.8</u> / 93.8

Attack	VGG-16 \implies		Inc-v3	Inc-v3 \implies		
	Res-50	Dense-121		Res-50	Dense-121	
I / +RAP / +RAP-LS	53.7 / <u>53.0</u> / 54.2	49.1 / <u>50.6</u> / 51.4	22.0 / <u>24.7</u> / 24.9	51.5 / <u>62.1</u> / 62.0	48.7 / <u>60.8</u> / 60.0	55.1 / <u>65.9</u> / 68.0
MI / +RAP / +RAP-LS	62.5 / <u>76.2</u> / 76.4	60.5 / <u>73.0</u> / 73.9	30.0 / <u>42.7</u> / 42.2	62.0 / <u>85.8</u> / 84.8	56.7 / <u>84.6</u> / 84.6	63.1 / <u>84.9</u> / 84.6
TI / +RAP / +RAP-LS	62.8 / <u>64.8</u> / 65.8	55.9 / <u>63.7</u> / 62.1	29.1 / <u>36.2</u> / 37.1	49.3 / <u>63.4</u> / 61.6	49.4 / <u>63.4</u> / 63.8	58.1 / <u>68.6</u> / 69.5
DI / +RAP / +RAP-LS	72.2 / <u>86.0</u> / 88.8	68.8 / <u>85.0</u> / 87.4	29.9 / <u>46.6</u> / 51.6	68.4 / <u>81.7</u> / 81.8	71.9 / <u>85.0</u> / 84.0	76.1 / <u>85.2</u> / 86.4
SI / +RAP / +RAP-LS	80.0 / <u>92.7</u> / 94.7	82.1 / <u>94.8</u> / 95.7	45.8 / <u>74.0</u> / 74.7	66.2 / <u>69.8</u> / 72.8	65.9 / <u>74.9</u> / 77.2	66.0 / <u>69.2</u> / 73.0
Admix / +RAP / +RAP-LS	87.3 / <u>94.6</u> / 96.8	88.2 / <u>96.4</u> / 97.2	55.5 / <u>77.6</u> / 80.8	75.9 / <u>80.2</u> / 84.9	78.5 / <u>83.7</u> / 87.4	74.5 / <u>77.2</u> / 83.5

Table 2: The **untargeted attack success rate (%) of combinational methods with RAP**. The results with *CE* loss are reported. The best results are bold and the second best results are underlined.

Attack	ResNet-50 \implies		Inc-v3	DenseNet-121 \implies		
	Dense-121	VGG-16		Res-50	VGG-16	
MTDI / +RAP / +RAP-LS	99.8 / 100 / 100	99.8 / 100 / <u>99.9</u>	85.7 / <u>96.0</u> / 96.9	99.4 / <u>99.8</u> / 100	99.2 / <u>99.5</u> / 100	89.1 / <u>97.1</u> / 97.1
MTDSI / +RAP / +RAP-LS	100 / 100 / 100	99.7 / 99.9 / <u>99.8</u>	97.0 / 99.1 / <u>99.1</u>	99.8 / 99.9 / 99.9	99.2 / <u>99.3</u> / 99.7	95.1 / <u>98.3</u> / 98.4
MTDAI / +RAP / +RAP-LS	100 / 100 / 100	99.8 / 99.9 / 99.9	98.3 / <u>99.2</u> / 99.8	<u>99.8</u> / <u>99.8</u> / 99.9	99.4 / <u>99.6</u> / 99.8	97.9 / <u>98.8</u> / 98.9

Attack	VGG-16 \implies		Inc-v3	Inc-v3 \implies		
	Res-50	Dense-121		Res-50	Dense-121	
MTDI / +RAP / +RAP-LS	90.0 / <u>97.2</u> / 97.7	88.8 / <u>97.0</u> / 97.3	56.8 / 82.6 / <u>81.4</u>	82.9 / 91.8 / <u>90.6</u>	85.7 / 94.2 / <u>93.3</u>	85.1 / 92.7 / <u>91.0</u>
MTDSI / +RAP / +RAP-LS	97.6 / <u>98.8</u> / 99.4	98.1 / <u>99.2</u> / 99.4	85.0 / <u>94.1</u> / 95.2	89.0 / <u>91.2</u> / 92.3	92.0 / <u>95.2</u> / 95.6	87.6 / <u>90.3</u> / 92.2
MTDAI / +RAP / +RAP-LS	97.8 / <u>99.2</u> / 99.6	98.9 / <u>99.5</u> / 99.6	89.3 / <u>95.0</u> / 95.5	91.5 / <u>94.1</u> / 94.7	95.4 / <u>96.2</u> / 97.6	91.4 / <u>93.2</u> / 94.1

For TI, DI, SI, and Admix, RAP gets the improvements by 10.2%, 10.9%, 9.3%, and 6.3%. With late-start, RAP-LS further enhance the transfer attack performance for almost all methods.

Combinational Methods. Prior works demonstrate the combination of baseline methods could largely boost the adversarial transferability [41, 49]. We also investigate of behavior of RAP when incorporated with the combinational attacks. The results are shown in Table 2. As shown in the table, there exist the clear improvements of the combinational attacks over all baseline attacks shown in Table 1. In addition, our RAP-LS further boosts the average attack success rate of the three combinational attacks by 6.9%, 2.6%, and 1.7% respectively. Combined with the three combinational attacks, RAP-LS achieves 95.4%, 97.6%, and 98.3% average attack success rate, respectively. These results demonstrate RAP can significantly enhance the transferability.

4.3 The Evaluation of Targeted Attacks

We then evaluate the targeted attack performance of the different methods with RAP. The results with logit loss are presented and the results with CE loss are shown in *Appendix*.

Baseline Methods. The results of RAP with baseline attacks are shown in Table 3. From the results, RAP is also very effective in enhancing the transferability in targeted attacks. Taking ResNet-50 and DenseNet-121 as surrogate models for example, the average performance improvements induced by RAP are 5.0% (I), 8.1% (MI), 4.6% (TI), 10.4% (DI), 18.5% (SI), and 15.1% (Admix), respectively. Comparing to the ResNet-50 and DenseNet-121, the baseline attacks generally achieve lower transferability when using the VGG-16 or Inception-v3 as the surrogate models, which has also been verified in existing works [48, 49]. However, for Inception-v3 and VGG-16 as the surrogate models, RAP also consistently boosts the transferability under all cases. With late-start, RAP-LS could further improve the transferability of RAP for most attacks. The average attack success rate under all attack cases of RAP-LS is 2.6% higher than that of RAP.

Combinational Methods. As did in the untargeted attacks, we also evaluate the performance of combinational methods. The results are shown in Table 4. Similar to the findings in untargeted attacks, the combinational methods obtain significantly improvements over baseline methods. The RAP-LS outperforms all combinational methods by a significantly margin. For example, taking the average attack success rate of all target models as evaluation metric, RAP-LS obtains 14.2%, 11.8%, 9.3% improvements over the MTDI, MTDSI and MTDAI, respectively.

Table 3: The **targeted attack success rate (%) of baseline methods with RAP**. The results with logit loss are reported. The best results are bold and the second best results are underlined.

Attack	ResNet-50 \implies		Inc-v3	DenseNet-121 \implies		
	Dense-121	VGG-16		Res-50	VGG-16	Inc-v3
I/+RAP/+RAP-LS	4.5 / 9.5 / 14.3	2.4 / 9.8 / 11.8	0.1 / 0.1 / 0.7	5.0 / 12.8 / 17.9	2.9 / 10.1 / 15.9	0.0 / 0.8 / 1.2
MI/+RAP/+RAP-LS	6.3 / <u>17.5</u> / 29.6	2.2 / <u>14.5</u> / 20.6	0.1 / <u>1.1</u> / 2.4	4.6 / <u>16.2</u> / 26.5	3.1 / <u>13.4</u> / 23.2	0.3 / <u>2.0</u> / 3.4
TI/+RAP/+RAP-LS	7.2 / <u>11.0</u> / 17.3	4.0 / <u>12.9</u> / 15.3	0.1 / <u>0.8</u> / 1.2	8.4 / <u>13.5</u> / 20.8	5.2 / <u>12.4</u> / 16.4	0.2 / <u>2.1</u> / 3.0
DI/+RAP/+RAP-LS	62.6 / <u>64.9</u> / 73.9	57.2 / <u>63.4</u> / 69.3	1.5 / <u>7.9</u> / 10.1	30.2 / <u>52.6</u> / 60.4	32.1 / <u>49.5</u> / 58.9	1.4 / <u>8.8</u> / 10.0
SI/+RAP/+RAP-LS	30.0 / <u>53.2</u> / 61.1	9.5 / <u>32.8</u> / 36.0	1.8 / <u>9.3</u> / 10.5	14.2 / <u>41.5</u> / 43.4	8.4 / <u>31.0</u> / 35.2	1.6 / <u>8.5</u> / 10.4
Admix/+RAP/+RAP-LS	54.6 / <u>68.0</u> / 74.6	26.0 / <u>45.4</u> / 51.6	5.8 / <u>17.1</u> / 19.6	29.3 / <u>53.0</u> / 58.2	21.5 / <u>42.7</u> / 48.2	5.0 / <u>17.1</u> / 17.6

Attack	VGG-16 \implies		Inc-v3	Inc-v3 \implies		
	Res-50	Dense-121		Res-50	Dense-121	VGG-16
I/+RAP/+RAP-LS	0.1 / <u>0.7</u> / 1.4	0.2 / <u>1.4</u> / 1.7	0.0 / <u>0.1</u> / 0.2	0.2 / 0.9 / <u>0.5</u>	0.2 / 0.6 / <u>0.3</u>	0.1 / 0.5 / <u>0.5</u>
MI/+RAP/+RAP-LS	0.5 / <u>1.3</u> / 1.9	0.5 / <u>2.3</u> / 3.0	0.0 / <u>0.0</u> / 0.3	0.2 / <u>1.7</u> / <u>1.5</u>	0.1 / <u>1.6</u> / <u>1.5</u>	0.2 / <u>1.3</u> / <u>1.0</u>
TI/+RAP/+RAP-LS	0.7 / <u>1.2</u> / 3.2	0.8 / <u>1.7</u> / 2.9	0.0 / <u>0.1</u> / 0.4	0.2 / <u>0.5</u> / 0.7	0.1 / 0.7 / <u>0.6</u>	0.2 / 0.8 / <u>0.6</u>
DI/+RAP/+RAP-LS	2.8 / <u>7.3</u> / 9.7	3.8 / <u>8.4</u> / 12.7	0.0 / <u>0.4</u> / 1.1	1.6 / <u>4.6</u> / 6.4	2.8 / <u>5.8</u> / 7.5	2.6 / <u>6.3</u> / 8.1
SI/+RAP/+RAP-LS	3.3 / 9.8 / 9.8	7.2 / <u>16.8</u> / 17.8	0.2 / <u>1.7</u> / 1.8	0.6 / 2.9 / <u>2.5</u>	0.9 / <u>2.7</u> / 3.2	0.5 / <u>1.5</u> / 2.3
Admix/+RAP/+RAP-LS	5.6 / <u>11.1</u> / 11.9	13.0 / <u>20.2</u> / 23.6	0.7 / <u>2.4</u> / 2.8	1.5 / <u>4.9</u> / 5.2	2.0 / <u>6.9</u> / 7.5	1.3 / <u>3.3</u> / 4.4

Table 4: The **targeted attack success rate (%) of combinational methods with RAP**. The results with logit loss are reported. The best results are bold and the second best results are underlined.

Attack	ResNet-50 \implies		Inc-v3	DenseNet-121 \implies		
	Dense-121	VGG-16		Res-50	VGG-16	Inc-v3
MTDI/+RAP/+RAP-LS	74.9 / <u>78.2</u> / 88.5	62.8 / <u>72.9</u> / 81.5	10.9 / <u>28.3</u> / 33.2	44.9 / <u>64.3</u> / 74.5	38.5 / <u>55.0</u> / 65.5	7.7 / <u>23.0</u> / 26.5
MTDSI/+RAP/+RAP-LS	86.3 / <u>88.4</u> / 93.3	70.1 / <u>77.7</u> / 84.7	38.1 / <u>51.8</u> / 58.0	55.0 / <u>71.2</u> / 75.8	42.0 / <u>58.4</u> / 62.3	19.8 / <u>39.0</u> / 39.2
MTDAI/+RAP/+RAP-LS	<u>91.4</u> / 89.4 / 93.6	<u>79.9</u> / 79.0 / 86.3	50.8 / <u>57.1</u> / 64.1	69.1 / <u>74.2</u> / 82.1	54.7 / <u>63.1</u> / 69.3	32.0 / <u>43.5</u> / 49.3

Attack	VGG-16 \implies		Inc-v3	Inc-v3 \implies		
	Res-50	Dense-121		Res-50	Dense-121	VGG-16
MTDI/+RAP/+RAP-LS	11.8 / <u>16.7</u> / 22.9	13.7 / <u>19.4</u> / 27.4	0.7 / <u>3.4</u> / 4.6	1.8 / 8.3 / <u>7.5</u>	4.1 / 14.8 / <u>13.4</u>	2.9 / <u>8.0</u> / 9.8
MTDSI/+RAP/+RAP-LS	31.0 / <u>35.3</u> / 38.7	41.7 / <u>44.4</u> / 49.6	9.6 / <u>15.2</u> / 13.7	5.6 / 11.9 / <u>10.7</u>	10.4 / 21.2 / <u>20.9</u>	4.2 / 8.9 / <u>8.6</u>
MTDAI/+RAP/+RAP-LS	36.2 / <u>39.0</u> / 43.1	48.0 / <u>45.1</u> / 55.2	11.6 / <u>17.1</u> / 17.6	9.6 / <u>13.6</u> / 16.7	17.9 / <u>27.5</u> / 31.6	8.4 / <u>12.0</u> / 12.1

4.4 The Comparison with Other Types of Attacks

Apart from the baseline and the combinational methods, we also experiment with more diverse attack methods, including the model-specific attack LinBP [14], the feature-based attack ILA [17], and the generative targeted attack TTP [30]. The LinBP depends on the skip connection and the authors only provide the source code about ResNet-50. We use their released

code and thus conduct experiments with ResNet-50 as \mathcal{M}^S . The results of LibBP and ILA are shown in Table 5, where we also implement the variants of LinBP following Guo et al. [14], including LinBP-ILA, LinBP-ILA-SGM, LinBP-MI-DI, and LinBP-MI-DI-SGM. We observe that our MI-DI-RAP significantly outperforms the LinBP and ILA, especially for the targeted attacks. Compared with the second-best method (*i.e.*, LinBP-MI-DI-SGM), we obtain a large improvement by 33.5% on average ASR of targeted attacks.

TTP [30] is the state-of-the-art generative method to conduct targeted attack. To compare with it, we adopt the generators based on ResNet-50 provided by the authors. Since TTP needs to train the perturbation generator for each targeted class, we follow their “10-Targets (all-source)” setting, as did in Zhao et al. [49]. The results are shown in Table 6, where our MTDSI+RAP-LS behaves best and outperforms TTP and MTDI by large margins of 14.9% and 25.7%, respectively.

4.5 The Evaluation on Diverse Network Architectures and Defense Models

To further demonstrate the efficacy of RAP, we evaluate our method on more diverse network architectures, including Inception-ResNet-v2, NASNet-Large and ViT-Base/16. We adopt ResNet-50 as the surrogate model and the results are shown in Table 7, *col 2-7*. As shown in the table, the

Table 5: The comparison with ILA and LinBP. We use ResNet-50 as \mathcal{M}^S . The best results are bold.

Attack	Untargeted			Targeted		
	Dense-121	VGG-16	Inc-v3	Dense-121	VGG-16	Inc-v3
ILA	95.0	94.2	77.7	2.8	1.5	0.5
LinBP-ILA	99.5	99.2	89.8	9.4	4.9	2.0
LinBP-ILA-SGM	99.7	99.3	91.1	13.3	7.2	2.8
LinBP-MI-DI	99.5	99.2	89.3	26.1	16.5	3.2
LinBP-MI-DI-SGM	99.8	99.3	90.2	32.6	22.1	4.6
MI-DI+RAP	99.9	100	93.7	75.1	69.7	13.9

Table 6: The comparison with TTP on targeted attack. The best results are bold.

Attack	Dense-121	VGG-16	Inc-v3
TTP	79.6	78.6	40.3
MTDI	78.6	74.6	12.7
MTDI+RAP-LS	90.8	87.2	35.4
MTDSI	93.2	80.0	41.3
MTDSI+RAP-LS	95.7	88.1	59.3

Table 7: The evaluation on diverse network architectures and defense models.

Attack	Untargeted			Targeted			Untargeted		Targeted	
	IncRes-v2	NASNet-L	ViT-B/16	IncRes-v2	NASNet-L	ViT-B/16	Inc-v3 _{adv}	IncRes-v2 _{ens}	Inc-v3 _{adv}	IncRes-v2 _{ens}
MTDI	83.4	89.0	27.9	14.8	32.1	0.4	68.1	50.9	0.8	0.0
MTDI+RAP-LS	95.6	97.5	42.7	43.0	62.5	1.7	86.5	72.3	9.7	4.1
MTDSI	95.7	98.0	43.0	45.5	67.9	2.6	90.0	79.6	12.7	6.7
MTDSI+RAP-LS	98.6	99.7	57.4	64.0	80.4	5.3	96.5	91.5	31.0	22.0
MTDAI	97.3	98.8	45.5	58.4	75.3	3.3	92.1	82.7	17.2	12.2
MTDAI+RAP-LS	99.2	99.8	60.2	70.4	82.6	7.4	96.7	91.6	34.4	26.0

proposed RAP-LS achieves significant improvements for all three combinational methods on all target models, and MTDAI+RAP-LS achieves the best performance for diverse models. For MTDAI, the average performance improvements induced by RAP-LS is 5.9% and 7.8% for untargeted and targeted attacks, respectively. Since ViT is based on the transformer architecture that totally being different from convolution models, the transfer attacks based on Resnet-50 behave relatively poor on it, especially on targeted attacks. Yet our RAP-LS still obtains consistent improvements for all compared methods. We also consider the ensemble-model attack on these diverse network architectures and the results are given in *Appendix*.

Furthermore, we evaluate RAP on attacking defense models. We choose ensemble adversarial training (AT) model, adv-Inc-v3 and ens-adv-Inc-Res-v2 [38], multi-step AT model [33], input transformation defense [46], feature denoising [47], and purification defense (NRP) [29]. We only demonstrate the results of ensemble AT in main submission. The evaluations of other defense models are shown in *Appendix*. Following prior works [41, 48], we adopt the ensemble-model attack by averaging the logits of different surrogate models, including ResNet-50, ResNet-101, Inception-v3, and Inception-ResNet-v2. The transfer attack success rate on defense models are shown in Table 7, *col 8-11*. We can observe that our RAP-LS further boosts transferability of the baseline methods on both targeted and untargeted attacks. For untargeted attacks, RAP-LS achieves average performance improvements of 9.8% and 14.1% on Inc-v3_{adv} and IncRes-v2_{ens}, respectively. For targeted attacks, the average performance improvements of RAP-LS are 14.8% and 11.1%, respectively.

4.6 Ablation Study

We conduct ablation study on the hyper-parameters of the proposed RAP, including the size of neighborhoods ϵ_n , the iteration number of inner optimization T and late-start K_{LS} . We adopt targeted attacks with ResNet-50 as the surrogate model.

We first evaluate the effect of ϵ_n and T . We consider different values of ϵ_n , including 2/255, 4/255, 8/255, 12/255, 16/255, and 20/255. In Figure 4 (a), we plot the tendency curves of the targeted attack success rate under different values of ϵ_n and T . Note that in Sec. 3.2, we set $\alpha_n = \epsilon_n/T$. Thus for a fixed ϵ_n , larger T indicating lower stepsize α_n . The minimum stepsize of α_n is set to 2/255. We have the following observation from the plot: for a fixed ϵ_n , the more iterations T , the better attack performance. Thus, we adopt a relatively smaller $\alpha_n = 2/255$ in our experiments. In Figure 4 (b-d), we further plot the results of different attack methods and target models *w.r.t.* ϵ_n , where $\alpha_n = 2/255$. As shown in the plots, the larger ϵ_n generally improves the attack performance. For Inception-v3 and DenseNet-121, the improvements become mild for even larger ϵ_n . Overall, the value of 12 or 16 could lead to satisfactory result under most cases.

Then we conduct the ablation study of K_{LS} . In Figure 5, we report the targeted attack success rate of I, MI, DI, and MI-TI-DI combined with RAP-LS with $K_{LS} = 0, 25, 50, 100, 150, 200$. Note that the RAP-LS with $K_{LS} = 0$ reduces to RAP. As shown in the plots, the proposed late-start strategy can further boost attack performance of RAP for most cases. In general, the performance improvements increase as K_{LS} increases, and then become mild when K_{LS} is larger than 100. The suitable value of K_{LS} is relatively consistent among different methods and target models.

4.7 The Targeted Attack Against Google Cloud Vision API

Finally, we conduct the transfer attacks to attack a practical and widely used image recognition system, Google Cloud Vision API, and in the more challenging targeted attack scenario. MTDAI-RAP-LS behaves the best performance in above experiments, so we choose it to conduct the attack. We take the evaluation on randomly selected 500 images and use ResNet-50 as surrogate model. As the API

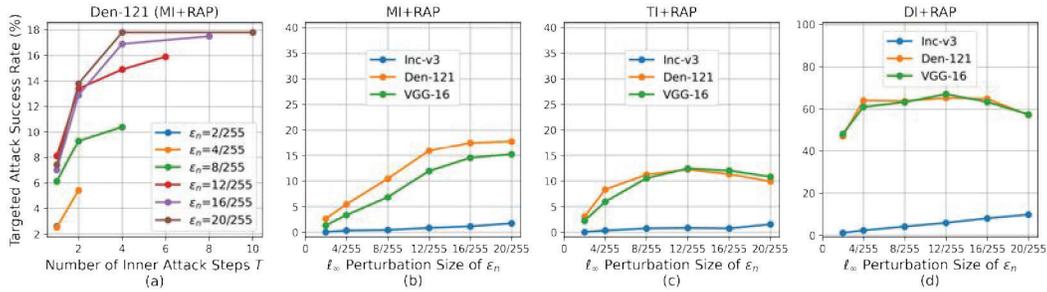


Figure 4: Targeted attack success rate (%) with various T and ϵ_n . Res-50 is set as surrogate model.

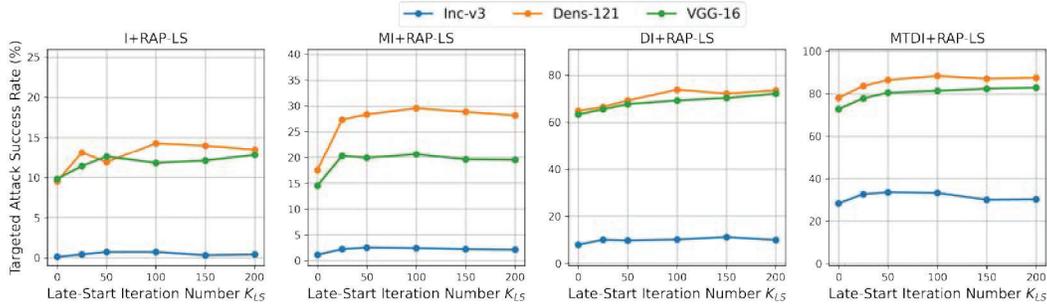


Figure 5: Targeted attack success rate (%) with various K_{LS} . Res-50 is set as surrogate model.

returns 10 predicted labels for each query, to evaluate the attacking performance, we test whether or not the target class appears in the returned predictions. Since the predicted label space of Google Cloud Vision API do not fully correspond to the 1000 ImageNet classes, we manually treat classes with similar semantics to be the same classes. In comparison, the baseline MTDAL successfully attacks 232 images against the Google API. Our RAP-LS achieves a large improvement, successfully attacking 342 images, leading to a 22.0% performance improvements. These demonstrates the high efficacy of our method to improve transferability on real-world system.

5 Conclusion

In this work, we study the transferability of adversarial examples that is significant for black-box attacks. The transferability of adversarial examples is generally influenced by the overfitting of surrogate models. To alleviate this, we propose to seeking adversarial examples that locate at flatter local regions. That is, instead of optimizing the pinpoint attack loss, we aim to obtain a consistently low loss at the neighbor regions of the adversarial examples. We formulate this as a min-max bi-level optimization problem, where the inner maximization aims to inject the worse-case perturbation for the adversarial examples. We conduct a rigorous experimental study, covering untargeted attack and targeted attack, standard and defense models, and a real-world Google Cloud Vision API. The experimental results demonstrate that RAP can significantly boost the transferability of adversarial examples, which also demonstrates that transfer attacks have become serious threats. We need to consider how to effectively defense against them.

Acknowledgments

We want to thank the anonymous reviewers for their valuable suggestions and comments. This work is supported by the National Natural Science Foundation of China under grant No.62076213, Shenzhen Science and Technology Program under grant No.RCYX20210609103057050, and the university development fund of the Chinese University of Hong Kong, Shenzhen under grant No.01001810, and sponsored by CCF-Tencent Open Fund.

References

- [1] Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Synthesizing robust adversarial examples. In *International conference on machine learning*, pages 284–293. PMLR, 2018. 2, 4.1, C, C.3
- [2] Minhao Cheng, Thong Le, Pin-Yu Chen, Huan Zhang, JinFeng Yi, and Cho-Jui Hsieh. Query-efficient hard-label black-box attack: An optimization-based approach. In *International Conference on Learning Representations*, 2019. 1
- [3] Minhao Cheng, Simranjit Singh, Patrick H. Chen, Pin-Yu Chen, Sijia Liu, and Cho-Jui Hsieh. Sign-opt: A query-efficient hard-label adversarial attack. In *International Conference on Learning Representations*, 2020. 1
- [4] Ambra Demontis, Marco Melis, Maura Pintor, Matthew Jagielski, Battista Biggio, Alina Oprea, Cristina Nita-Rotaru, and Fabio Roli. Why do adversarial attacks transfer? explaining transferability of evasion and poisoning attacks. In *28th USENIX security symposium (USENIX security 19)*, pages 321–338, 2019. 3.2
- [5] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting adversarial attacks with momentum. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9185–9193, 2018. 1, 2, 3.1, 3.2, 3.2, 3.3, 4.1
- [6] Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. Evading defenses to transferable adversarial examples by translation-invariant attacks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4312–4321, 2019. 1, 2, 3.1, 3.2, 3.2, 4.1
- [7] Yinpeng Dong, Qi-An Fu, Xiao Yang, Tianyu Pang, Hang Su, Zihao Xiao, and Jun Zhu. Benchmarking adversarial robustness on image classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 321–331, 2020. D
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 4.1
- [9] Yanbo Fan, Baoyuan Wu, Tuanhui Li, Yong Zhang, Mingyang Li, Zhifeng Li, and Yujiu Yang. Sparse adversarial attack via perturbation factorization. In *European conference on computer vision*, pages 35–50. Springer, 2020. 1
- [10] Yan Feng, Baoyuan Wu, Yanbo Fan, Li Liu, Zhifeng Li, and Shu-Tao Xia. Boosting black-box attack with partially transferred conditional adversarial distribution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15095–15104, June 2022. 1
- [11] Lianli Gao, Qilong Zhang, Jingkuan Song, Xianglong Liu, and Heng Tao Shen. Patch-wise attack for fooling deep neural network. In *European Conference on Computer Vision*, pages 307–322. Springer, 2020. B
- [12] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014. 1, 2
- [13] Martin Gubri, Maxime Cordy, Mike Papadakis, Yves Le Traon, and Koushik Sen. Lgv: Boosting adversarial example transferability from large geometric vicinity. *arXiv preprint arXiv:2207.13129*, 2022. 2
- [14] Yiwen Guo, Qizhang Li, and Hao Chen. Backpropagating linearly improves transferability of adversarial examples. *Advances in Neural Information Processing Systems*, 33:85–95, 2020. 1, 2, 4.1, 4.4, B
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4.1

- [16] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017. 4.1
- [17] Qian Huang, Isay Katsman, Horace He, Zeqi Gu, Serge Belongie, and Ser-Nam Lim. Enhancing adversarial example transferability with an intermediate level attack. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4733–4742, 2019. 1, 2, 4.1, 4.4
- [18] Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks with limited queries and information. In *International Conference on Machine Learning*, pages 2137–2146. PMLR, 2018. 1
- [19] Nathan Inkawhich, Kevin Liang, Lawrence Carin, and Yiran Chen. Transferable perturbations of deep feature distributions. In *International Conference on Learning Representations*, 2020. 2
- [20] Nathan Inkawhich, Kevin Liang, Binghui Wang, Matthew Inkawhich, Lawrence Carin, and Yiran Chen. Perturbing across the feature hierarchy to improve standard and strict blackbox attack transferability. *Advances in Neural Information Processing Systems*, 33, 2020. 2
- [21] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *Artificial Intelligence Safety and Security*, page 99–112, Jul 2018. doi: 10.1201/9781351251389-8. 2, 3.1, 3.3, 4.1
- [22] Yingwei Li, Song Bai, Yuyin Zhou, Cihang Xie, Zhishuai Zhang, and Alan Yuille. Learning transferable adversarial examples via ghost networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 11458–11465, 2020. 4.1, C
- [23] Siyuan Liang, Baoyuan Wu, Yanbo Fan, Xingxing Wei, and Xiaochun Cao. Parallel rectangle flip attack: A query-based black-box attack against object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7697–7707, 2021. 1
- [24] Jiadong Lin, Chuanbiao Song, Kun He, Liwei Wang, and John E. Hopcroft. Nesterov accelerated gradient and scale invariance for adversarial attacks. In *International Conference on Learning Representations*, 2020. 1, 2, 3.1, 3.2, 4.1
- [25] Risheng Liu, Jiaxin Gao, Jin Zhang, Deyu Meng, and Zhouchen Lin. Investigating bi-level optimization for learning and vision from a unified perspective: A survey and beyond. *arXiv preprint arXiv:2101.11517*, 2021. 3.2
- [26] Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and black-box attacks. *arXiv preprint arXiv:1611.02770*, 2016. 1
- [27] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018. 1
- [28] Muhammad Muzammal Naseer, Salman H Khan, Muhammad Haris Khan, Fahad Shahbaz Khan, and Fatih Porikli. Cross-domain transferability of adversarial perturbations. *Advances in Neural Information Processing Systems*, 32:12905–12915, 2019. 2
- [29] Muzammal Naseer, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Fatih Porikli. A self-supervised approach for adversarial robustness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 262–271, 2020. 4.1, 4.5, D
- [30] Muzammal Naseer, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Fatih Porikli. On generating transferable targeted perturbations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7708–7717, 2021. 1, 2, 4.1, 4.4, 4.4, B
- [31] Omid Poursaeed, Isay Katsman, Bicheng Gao, and Serge Belongie. Generative adversarial perturbations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4422–4431, 2018. 2
- [32] Zeyu Qin, Yanbo Fan, Hongyuan Zha, and Baoyuan Wu. Random noise defense against query-based black-box attacks. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021. 1

- [33] Hadi Salman, Andrew Ilyas, Logan Engstrom, Ashish Kapoor, and Aleksander Madry. Do adversarially robust imagenet models transfer better? In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 3533–3545. Curran Associates, Inc., 2020. 4.1, 4.5, D
- [34] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 4.1
- [35] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013. 1
- [36] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. 4.1
- [37] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Thirty-first AAAI conference on artificial intelligence*, 2017. 4.1
- [38] Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian Goodfellow, Dan Boneh, and Patrick McDaniel. Ensemble adversarial training: Attacks and defenses. In *International Conference on Learning Representations*, 2018. 1, 3.2, 4.1, 4.5
- [39] JingKang Wang, Tianyun Zhang, Sijia Liu, Pin-Yu Chen, Jiachen Xu, Makan Fardad, and Bo Li. Adversarial attack generation empowered by min-max optimization. *Advances in Neural Information Processing Systems*, 34, 2021. 2
- [40] Xiaosen Wang and Kun He. Enhancing the transferability of adversarial attacks through variance tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1924–1933, June 2021. 2, 4.1, C
- [41] Xiaosen Wang, Xuanran He, Jingdong Wang, and Kun He. Admix: Enhancing the transferability of adversarial attacks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16158–16167, 2021. 2, 3.1, 3.2, 3.2, 4.1, 4.2, 4.5, B
- [42] Xiaosen Wang, Jiadong Lin, Han Hu, Jingdong Wang, and Kun He. Boosting adversarial transferability through enhanced momentum. *arXiv preprint arXiv:2103.10609*, 2021. 4.1, C
- [43] 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 *International Conference on Learning Representations*, 2020. 2
- [44] Jiancong Xiao, Yanbo Fan, Ruoyu Sun, and Zhi-Quan Luo. Adversarial rademacher complexity of deep neural networks, 2022. 1
- [45] Jiancong Xiao, Yanbo Fan, Ruoyu Sun, Jue Wang, and Zhi-Quan Luo. Stability analysis and generalization bounds of adversarial training. *arXiv preprint arXiv:2210.00960*, 2022. 1
- [46] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan Yuille. Mitigating adversarial effects through randomization. In *International Conference on Learning Representations*, 2018. 3.2, 4.1, 4.5, D
- [47] Cihang Xie, Yuxin Wu, Laurens van der Maaten, Alan L Yuille, and Kaiming He. Feature denoising for improving adversarial robustness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 501–509, 2019. 4.1, 4.5, D, D
- [48] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2730–2739, 2019. 1, 2, 3.1, 3.2, 3.3, 4.1, 4.3, 4.5
- [49] Zhengyu Zhao, Zhuoran Liu, and Martha Larson. On success and simplicity: A second look at transferable targeted attacks. *arXiv preprint arXiv:2012.11207*, 2020. 2, 3.1, 4.1, 4.2, 4.3, 4.4, B, E.3

- [50] Xin Zheng, Yanbo Fan, Baoyuan Wu, Yong Zhang, Jue Wang, and Shirui Pan. Robust physical-world attacks on face recognition. *Pattern Recognition*, page 109009, 2022. 1
- [51] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8697–8710, 2018. 4.1

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes] See final part of Section 1.
 - (b) Did you describe the limitations of your work? [Yes] See Appendix
 - (c) Did you discuss any potential negative societal impacts of your work? [Yes] See Appendix
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [N/A]
 - (b) Did you include complete proofs of all theoretical results? [N/A]
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [No]
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 4.1
 - (b) Did you mention the license of the assets? [Yes] See Appendix
 - (c) Did you include any new assets either in the supplemental material or as a URL? [No]
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [No] No, all the used data are well-known public benchmark datasets such as CIFAR10, ImageNet.
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] No, all the used data are well-known public benchmark datasets such as CIFAR10, ImageNet.
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] Not applicable.
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] Not applicable.
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] Not applicable.