

# Boosting the Transferability of Adversarial Attacks with Reverse Adversarial Perturbation

Zeyu Qin<sup>\*1</sup>, Yanbo Fan<sup>\*2</sup>, Yi Liu<sup>1</sup>, Li Shen<sup>3</sup>, Yong Zhang<sup>2</sup>, Jue Wang<sup>2</sup>, Baoyuan Wu<sup>1</sup>

<sup>1</sup> The Chinese University of Hong Kong, Shenzhen

<sup>2</sup>Tencent AI Lab

<sup>3</sup> JD Explore Academy

NeurIPS 2022

# Outline

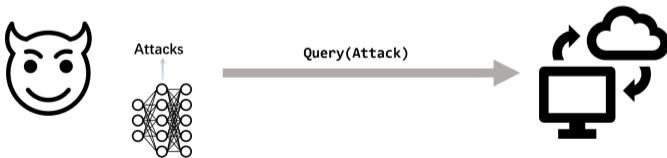
Introduction

Reverse Adversarial Perturbation

Experimental Evaluation

## Severe Threats from Black-box Attacks

- ▶ **Transfer attacks: using  $x^{adv}$  from  $\mathcal{M}^S$  to attack  $\mathcal{M}^T$** 
  - The attackers can **utilize same dataset to train the surrogate model  $\mathcal{M}^S$**
  - Generating  $x^{adv}$  (white-box attacks) on  $\mathcal{M}^S$ , then attacking  $\mathcal{M}^T$ .
  - Don't need to iteratively query but it is not practical and performs poor attack performance.



## Transfer attacks overfits to $\mathcal{M}^S$

- ▶ Taking the target attack as example, the general formulation of many existing transfer attack methods can be written as follows:

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

where  $\mathcal{L}$  is the adversarial loss function,  $y_t$  is target label.

- ▶ The existing transfer attack methods exhibit poor transferability on  $\mathcal{M}^T$  (not successfully attacking  $\mathcal{M}^T$ )
- ▶  $\mathbf{x}^{adv}$  severely depends on (overfits to) the decision boundaries of  $\mathcal{M}^S$  and there are huge difference of decision boundaries between  $\mathcal{M}^S$  and  $\mathcal{M}^T$ . [1,2,3]

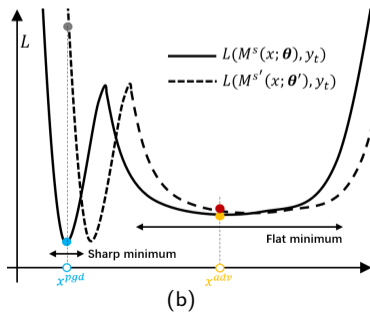
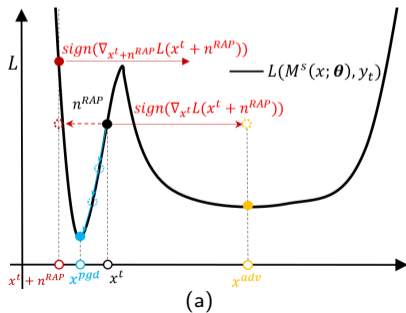
[1] Tramer et al., Ensemble Adversarial Training: Attacks and Defenses, ICLR 2018.

[2] Demontis et al., Why Do Adversarial Attacks Transfer? Explaining Transferability of Evasion and Poisoning Attacks, ACM CCS 2019.

[3] Dong et al., Boosting adversarial attacks with momentum, CVPR 2018.

## New Perspective to Interpret Adversarial Transferability

- ▶ We propose a new perspective to interpret the adversarial transferability, the flatness of (adversarial) loss landscape of  $x^{adv}$  on  $\mathcal{M}^S$ .
- ▶ The  $x^{adv}$  located at the flat local minimum is less sensitive to the changes of decision boundary (the difference of  $\mathcal{M}^S$  and  $\mathcal{M}^T$ ). Therefore, it could have the better adversarial transferability.



# Outline

Introduction

Reverse Adversarial Perturbation

Experimental Evaluation

## Finding $x^{adv}$ located at a local flat region

- ▶ We encourage that **not only**  $x^{adv}$  itself has low loss value, **but also** the points in the **vicinity of**  $x^{adv}$  have similarly low loss values.
- ▶ We propose to **minimize the maximal loss value within a local neighborhood region** around  $x^{adv}$ .
- ▶ The maximal loss is implemented by perturbing  $x^{adv}$  (adding perturbation  $\mathbf{n}^{adv}$ ) to maximize the adversarial loss, named **Reverse Adversarial Perturbation (RAP)**. So, we aim to solve this problem,

$$\min_{\mathbf{x}^{adv} \in \mathcal{B}_\epsilon(\mathbf{x})} \mathcal{L}(\mathcal{M}^S(\mathbf{x}^{adv} + \mathbf{n}^{adv}; \boldsymbol{\theta}), y_t)$$

Where,

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

## RAP with Late-Start

- ▶ In our preliminary experiments, we find that RAP requires more iterations to converge and the performance is slightly lower during the initial iterations.
- ▶ Hence, we further propose a better initialization, **late-start** (LS in following content) which first only runs the outer loop for several iterations then conducts the min-max loop.



# Outline

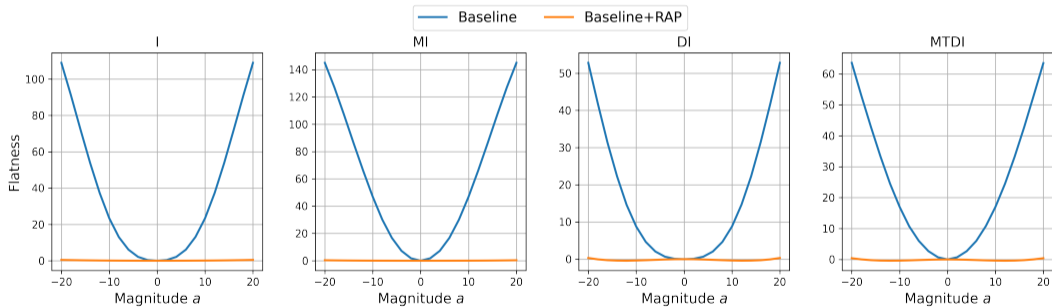
Introduction

Reverse Adversarial Perturbation

Experimental Evaluation

## A Closer Look at RAP

- ▶ First we **visualize the loss landscape around  $x^{adv}$  on  $\mathcal{M}^S$**  by plotting the loss variations. We can observe that RAP could help find  $x^{adv}$  located at the flat region.



## RAP achieves the better attack performance

- Combined with existing attacks, RAP **further boosts their transferability for both untargeted and targeted attacks.**

The below tables show the transfer targeted attack performance ( $\mathcal{M}^S \Rightarrow \mathcal{M}^T$ ).

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

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

## RAP achieves the better attack performance on diverse architectures

- ▶ We also conduct experiments on diverse network architectures, ViT and ensemble models. Our RAP-LS achieves the better attack performance.

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

## RAP achieves the SOTA attack performance stronger defense models

- ▶ We also take a comparison on stronger defense models. Our methods also achieve the SOTA performance on defense models. Our RAP-LS achieves the better attack performance.

Attack	Untargeted		
	Res-50 AT ( $\ell_2$ )	Res-50 AT ( $\ell_\infty$ )	Feature Denoising
MTDI	42.5	32.4	44.1
MTDI+RAP-LS	<b>59.5</b>	<b>34.4</b>	<b>44.4</b>
MTDSI	56.6	35.8	45.0
MTDSI+RAP-LS	<b>70.3</b>	<b>36.6</b>	<b>45.7</b>
MTDAI	62.1	35.6	44.2
MTDAI+RAP-LS	<b>73.7</b>	<b>37.7</b>	<b>45.2</b>

**Thank Y'all!**