

# On Success and Simplicity: A Second Look at Transferable Targeted Adversarial Images

(对有目标对抗图像迁移性的反思)

Zhengyu Zhao (赵正宇), Zhuoran Liu, Martha Larson  
Radboud University, The Netherlands (NeurIPS 2021)

# About Me



**CISPA**  
HELMHOLTZ CENTER FOR  
INFORMATION SECURITY



**Radboud  
University**

- 个人经历

2021.12 - present: **Postdoc @ CISPA** (亥姆霍兹信息安全中心), Germany

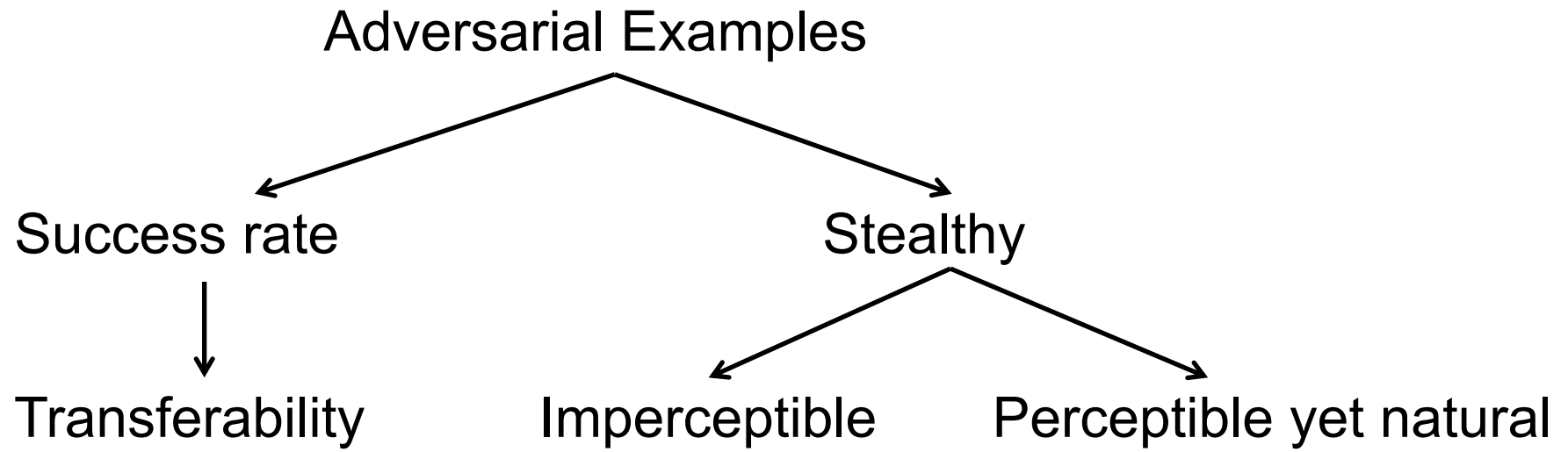
2017.09 - 2022.02: **PhD @ Radboud University**, Netherlands

- 研究兴趣

Security & Privacy in Computer Vision: Adversarial (Image) Examples,  
Data Poisoning (训练数据投毒), Membership Inference (训练成员推理).

zhengyu.zhao@cispa.de  
zhengyuzhao.github.io





# Adversarial Examples

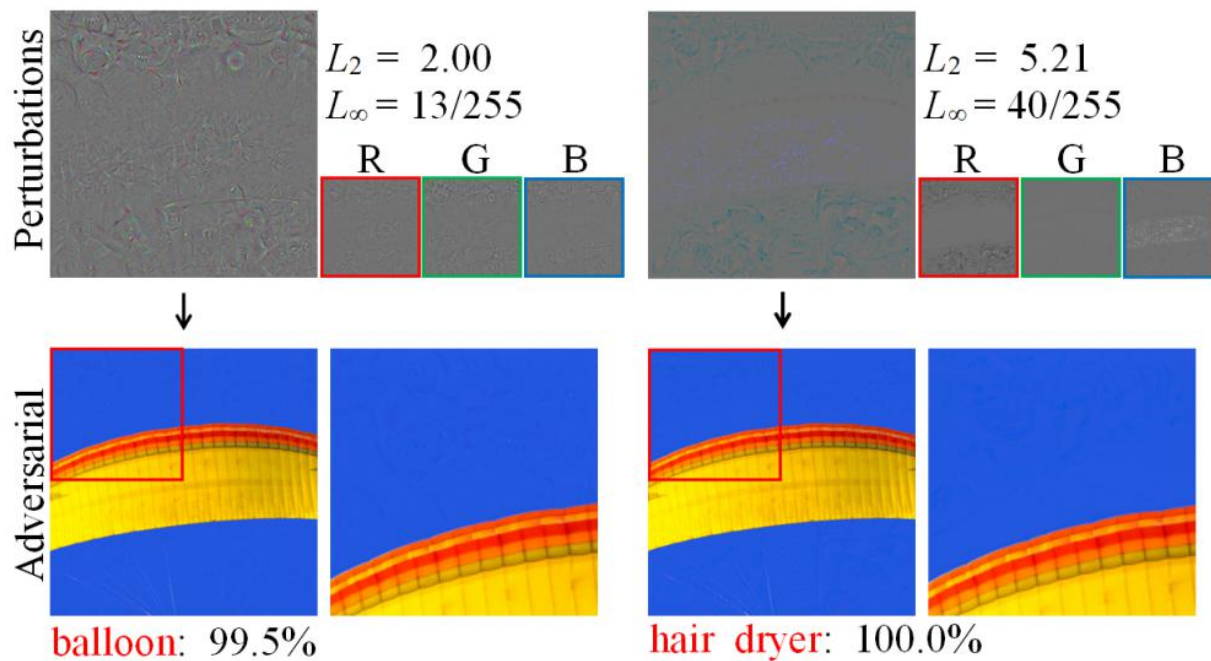
Success rate

Stealthy

Transferability

**Imperceptible**

Perceptible yet natural



(a) C&W

(b) PerC-C&W (ours)

change perspective

Radboud University



# Adversarial Examples

Success rate

Stealthy

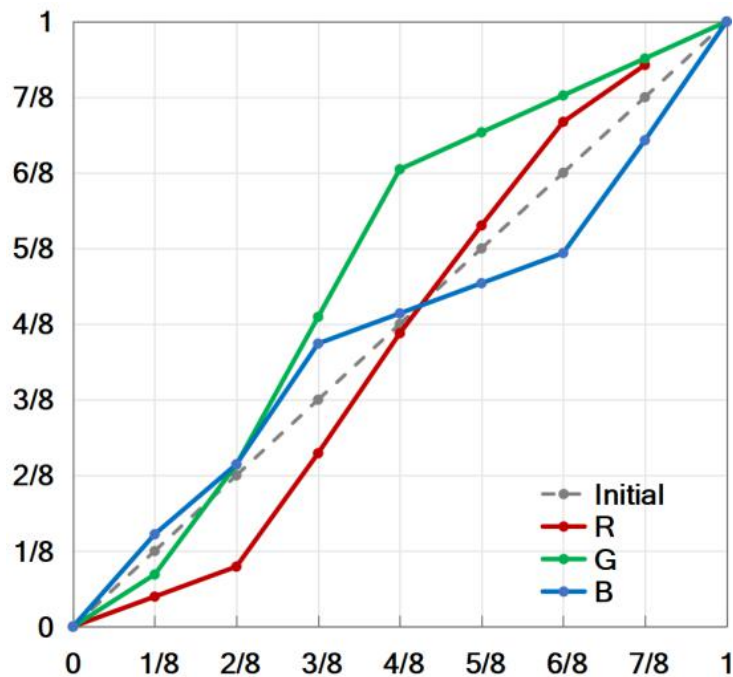
Transferability

Imperceptible

**Perceptible yet Natural**

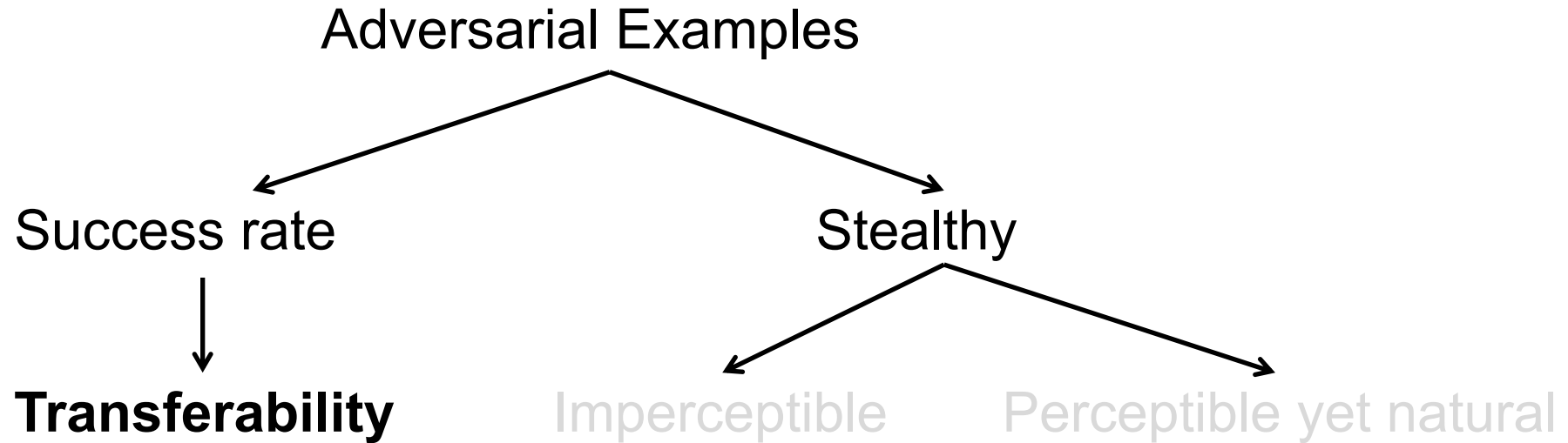


×



=





# On Success and Simplicity: A Second Look at Transferable Targeted Adversarial Images

Zhengyu Zhao (赵正宇), Zhuoran Liu, Martha Larson.  
NeurIPS 2021

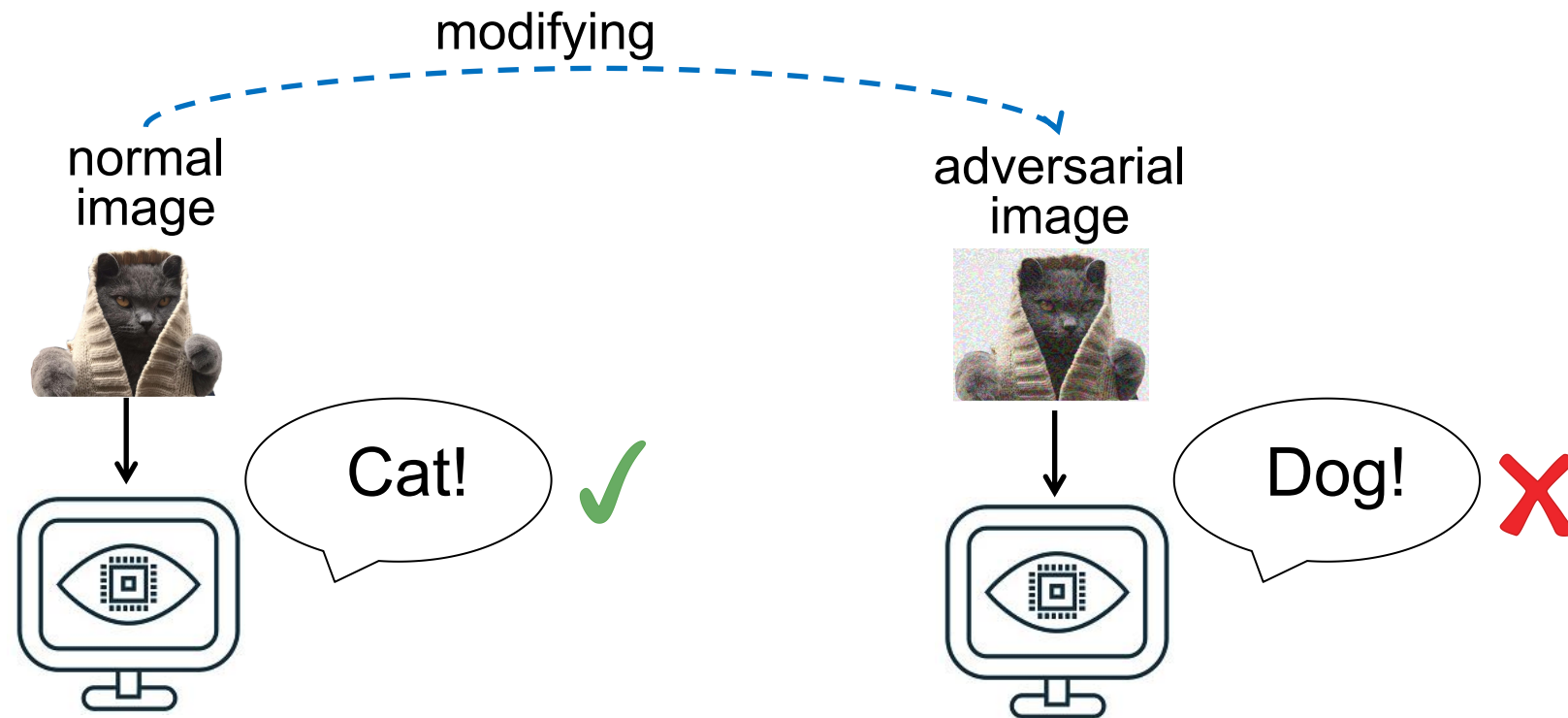


# Outline

- Recap adversarial image and targeted transferability
- Our work on revisiting targeted transferability
- Summary & future challenges

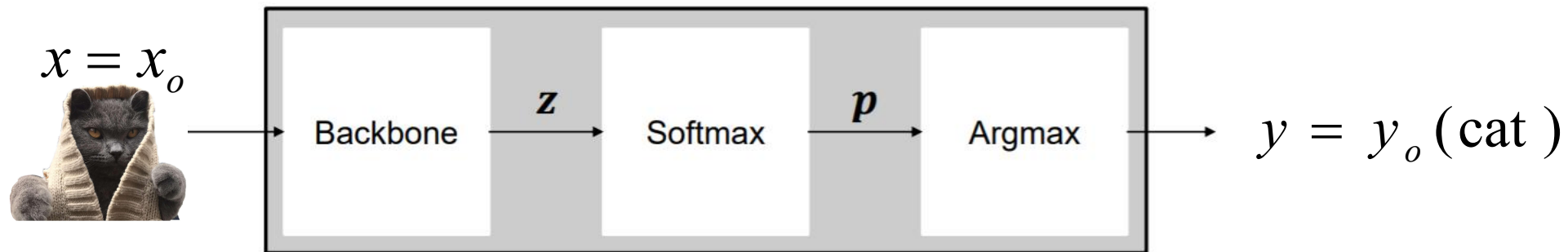


# Adversarial Images: Definition





# Adversarial Images: Formulation



Loss function

$$\theta^* = \arg \min_{\theta} J(x_0, y_0)$$

Learned parameters

Ground-truth label of  $X$

$$x' = \arg \min_x J(x, \underline{y}_t)$$

Target class

$$\|x' - x_0\|_{\infty} \leq \epsilon$$

change perspective

# Adversarial Images: Iterative Optimization

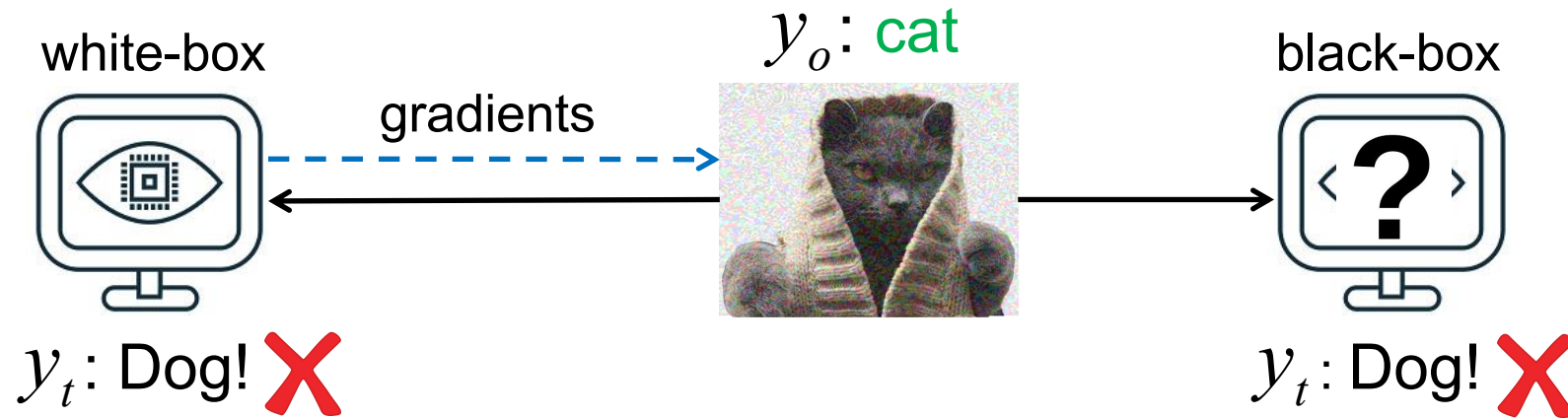
Objective function:  $x' = \arg \min_x J(x, y_t)$  s.t.  $\|x - x_o\|_\infty \leq \varepsilon$

Optimization: Projected Gradient Descent (PGD)

Iterative-Fast Gradient Sign Method (I-FGSM)<sup>[1]</sup>

$$x'_0 = x_o, \quad x'_{i+1} = x'_i - \alpha \cdot \text{sign}(\nabla_x J(x'_i, y_t))$$
$$x'_{i+1} \leftarrow \text{clip}(x'_{i+1} - x_o, -\varepsilon, \varepsilon)$$

# Adversarial Images: (Targeted) Transferability



# Targeted Transferability via Iterative Approach

Iterative-Fast Gradient Sign Method (I-FGSM):  $x'_0 = x_o$ ,  $x'_{i+1} = x'_i - \alpha \cdot \text{sign}(\nabla_x J(x'_i, y_t))$

Improve  
transferability

- Gradient stabilization<sup>[1,2]</sup>  
e.g. momentum-based<sup>[1]</sup>:

$$\mathbf{g}_{i+1} = \mu \cdot \mathbf{g}_i + \frac{\nabla_x J(\mathbf{x}'_i, y_t)}{\|\nabla_x J(\mathbf{x}'_i, y_t)\|_1}$$

$$\mathbf{x}'_{i+1} = \mathbf{x}'_i - \alpha \cdot \text{sign}(\mathbf{g}_i)$$

- Input augmentation<sup>[3,4,5]</sup>  
e.g. random resizing & padding<sup>[4]</sup>:

$$\mathbf{x}'_{i+1} = \mathbf{x}'_i - \alpha \cdot \text{sign}(\nabla_x J(T(\mathbf{x}'_i, p), y_t))$$

1. Dong et al. *Boosting Adversarial Attacks with Momentum*. CVPR 2018.

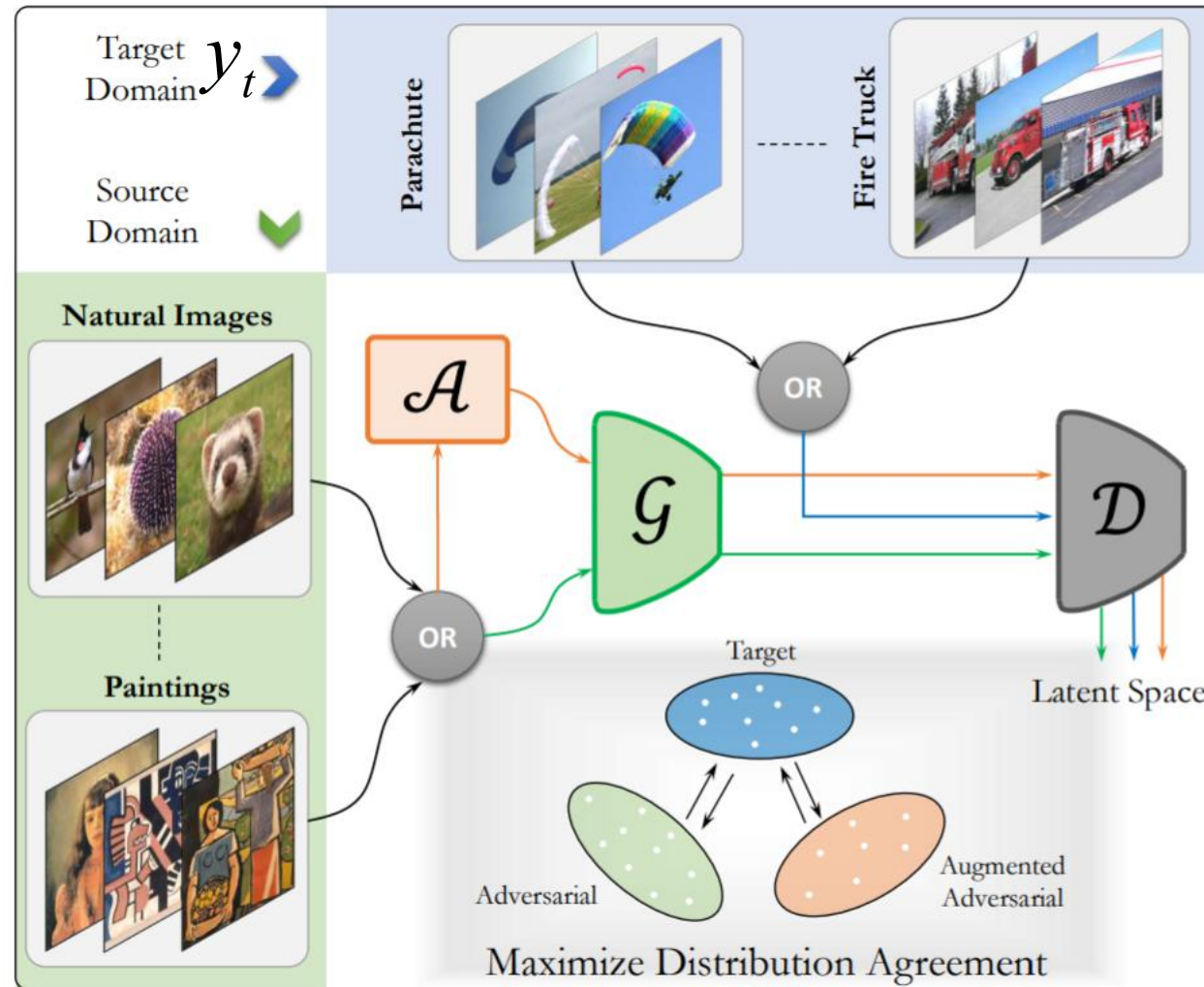
2. Lin et al. *Nesterov Accelerated Gradient and Scale Invariance for Adversarial Attacks*. ICLR 2020

3. Dong et al. *Evading Defenses to Transferable Adversarial Examples by Translation-Invariant Attacks*. CVPR 2019

4. Xie et al. *Improving Transferability of Adversarial Examples with Input Diversity*. CVPR 2019

5. Wang et al. *Admix: Enhancing the transferability of adversarial attacks*. ICCV, 2021.

# Targeted Transferability via Generative Approach



$\mathcal{A}$  : Augmenter

$\mathcal{G}$  : Generator

$\mathcal{D}$  : Discriminator

*change perspective*

# Iterative vs. Generative Approaches

## Iterative

- Data: Single Input image
- Model: 1 × target-agnostic classifier

vs

## Generative

- Massive training data
- 1000 × target-specific generators

Targeted Transferability: Iterative approach << Generative approach



# Revisiting Targeted Transferability: Key Message

Targeted Transferability: Iterative approach ~~X~~ Generative approach



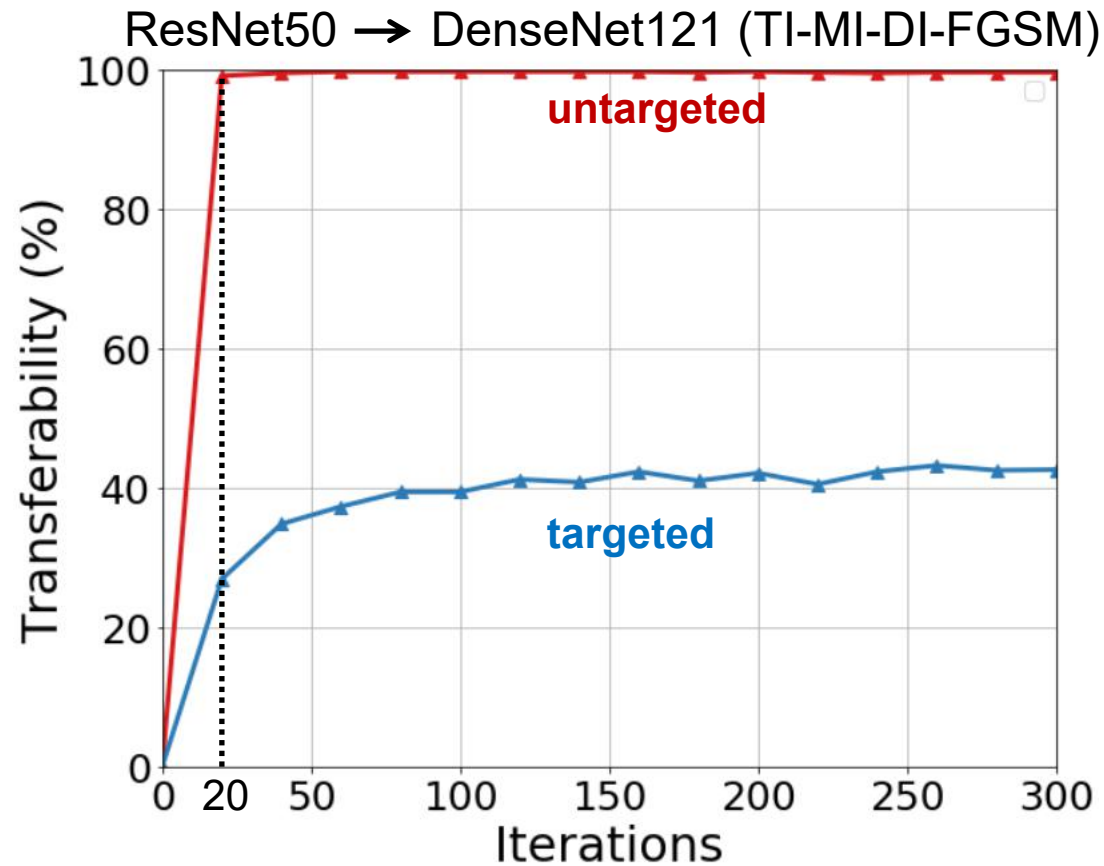
$$\left\| \begin{array}{c} x' \\ \text{[Image of cat in hood]} \\ x_0 \\ \text{[Image of cat in hood]} \end{array} \right\|_{\infty} \leq \epsilon$$

Targeted Transferability (%)

Bound	Attack	D121	V16	D121-ens	V16-ens
$\epsilon = 16$	TTP [8]	<b>79.6</b>	<b>78.6</b>	92.9	89.6
	ours	75.9	72.5	<b>99.4</b>	<b>97.7</b>
$\epsilon = 8$	TTP [8]	37.5	46.7	63.2	66.2
	ours	<b>44.5</b>	<b>46.8</b>	<b>92.6</b>	<b>87.0</b>



# Revisiting Targeted Transferability: More Iterations



Few ( $\leq 20$ ) iterations in the literature:

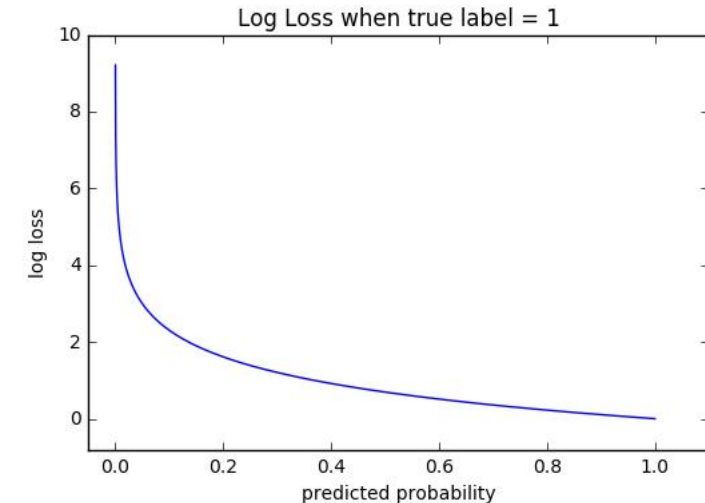
- not converge to optimal performance
- unrealistic iteration budget



# Revisiting Targeted Transferability: Better Loss

Cross-Entropy Loss ( $L_{CE}$ ) causes **decreasing gradient** problem:

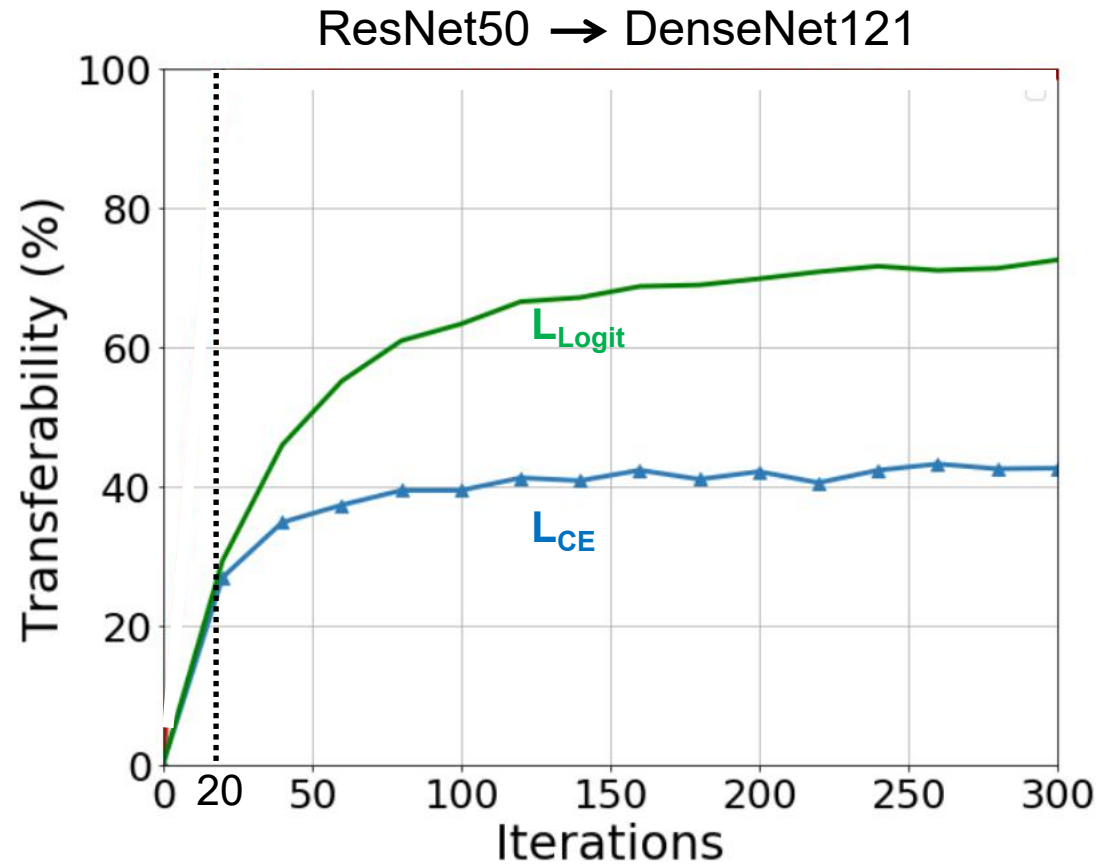
$$L_{CE} = -1 \cdot \log(p_t) = -\log\left(\frac{e^{z_t}}{\sum e^{z_j}}\right) = -z_t + \log\left(\sum e^{z_j}\right),$$
$$\frac{\partial L_{CE}}{\partial z_t} = -1 + \frac{\partial \log(\sum e^{z_j})}{\partial e^{z_t}} \cdot \frac{\partial e^{z_t}}{\partial z_t} = -1 + \frac{e^{z_t}}{\sum e^{z_j}} = -1 + p_t.$$



Logit Loss ( $L_{Logit}$ ) is better:

$$L_{Logit} = -z_t, \quad \frac{\partial L_{Logit}}{\partial z_t} = -1.$$

# Revisiting Targeted Transferability: Better Loss



# Revisiting Targeted Transferability: Better Evaluation

More challenging&realistic scenarios:

1. Model diversity
2. Target class diversity

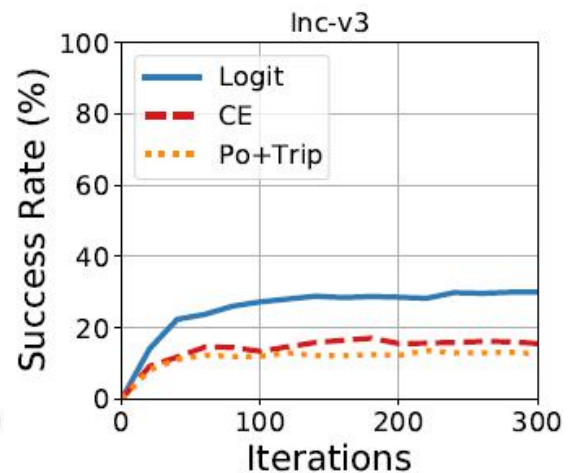
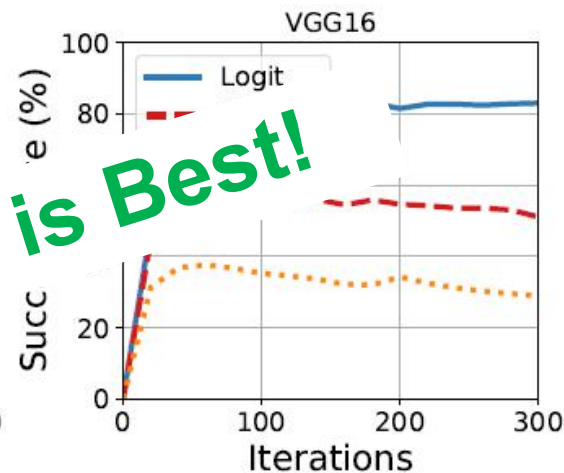
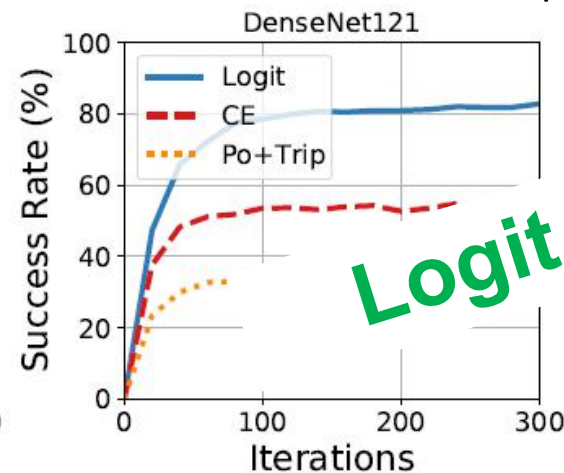
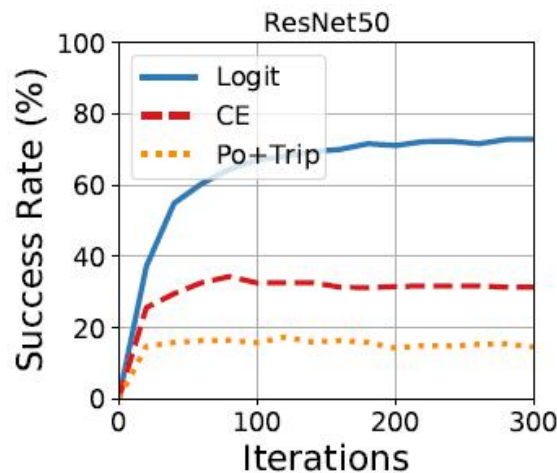
# Revisiting Targeted Transferability: Better Evaluation

1. Model diversity

2. Target class diversity

Attack	-Inc-v3	-Inc-v4	-IncRes-v2	-D	-Res101	-Res152	Average
CE	85.3	83.3	90.1	93.2	93.2	90.7	87.7
Po+Trip	84.4	82.4	85.0	87.9	87.9	85.7	84.4
Logit	<b>85.5</b>	<b>85.8</b>	<b>90.1</b>	<b>90.0</b>	91.4	<b>90.8</b>	<b>88.1</b>

**Saturation!**



**Logit is Best!**



# Revisiting Targeted Transferability: Better Evaluation

1. Model diversity

2. Target class diversity

Results on Google Cloud Vision API (100 images)

	CE	Po+Trip	Logit
<b>Transferability (%)</b>	7	8	<b>18</b>

Cloud Vision API

Labels

- Sky 96%
- Chinese Architecture 88%
- Travel 81%
- Temple 78%
- Composite Material 75%
- Facade 74%
- Building 73%
- Shade 72%

Labels

- Boat 93%
- Sky 92%
- Vehicle 86%
- Watercraft 86%
- Naval Architecture 81%
- Art 75%
- Water 72%
- Ship 72%

change perspective

$y_t$  is “yawl” (a type of boat)

# Revisiting Targeted Transferability: Better Evaluation

1. Model diversity

2. Target class diversity



$\mathcal{Y}_t$

1st: tabby cat 66%  
2nd: tiger cat 28%  
3rd: Egyptian cat 5%  
...  
1000th: airplane 0.1%

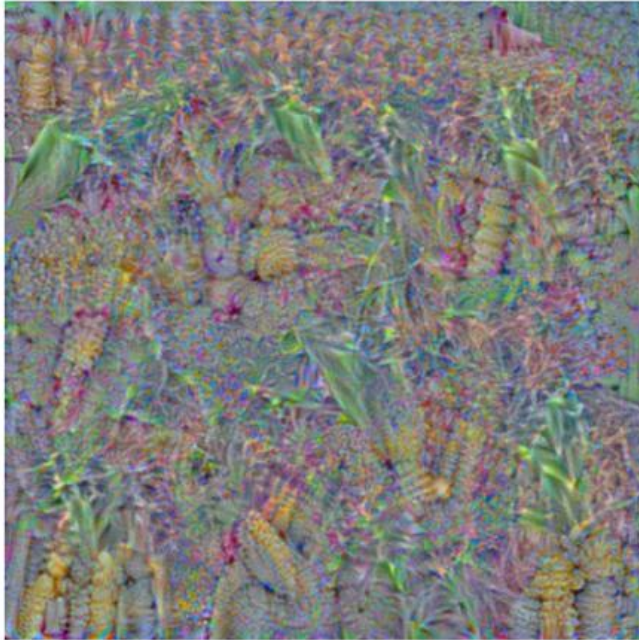
Targeted transferability (%) from best to worst case

Attack	2nd	10th	200th	500th	800th	1000th
CE	<b>89.9</b>	76.7	49.7	43.1	37.0	25.1
Po+Trip	82.6	77.6	58.4	53.6	49.1	38.2
Logit	83.8	<b>81.3</b>	<b>75.0</b>	<b>71.0</b>	<b>65.1</b>	<b>52.8</b>

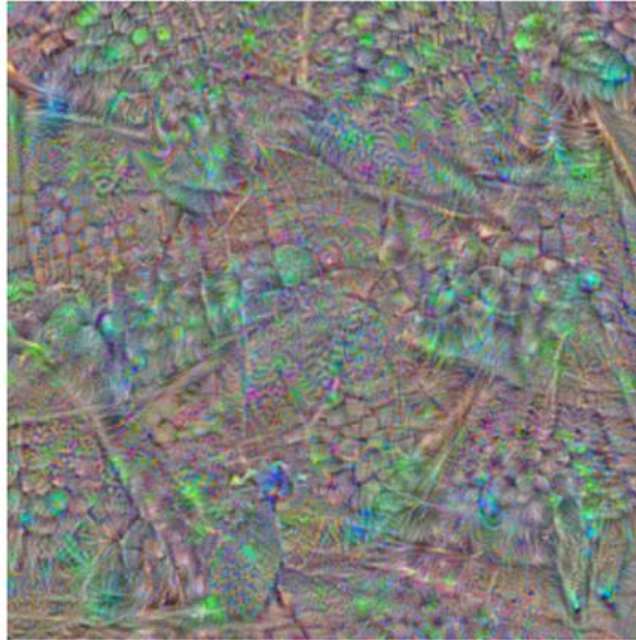


# Perturbation Semantics

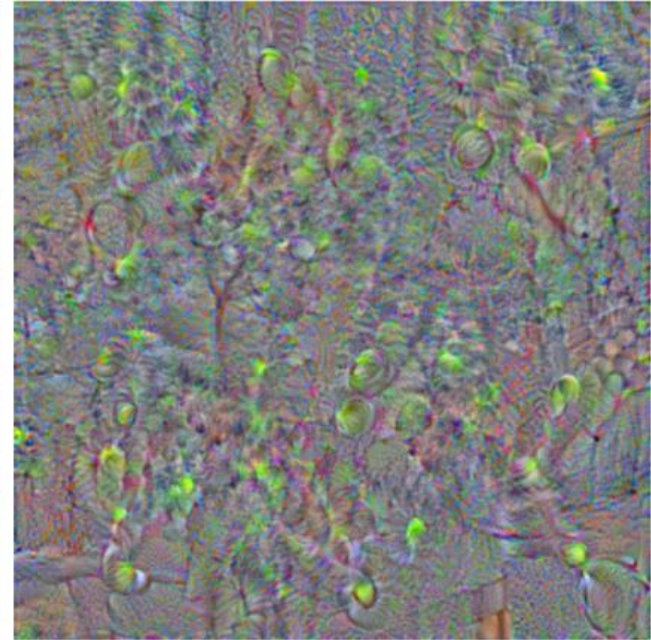
“corn”



“peacock”



“tennis ball”

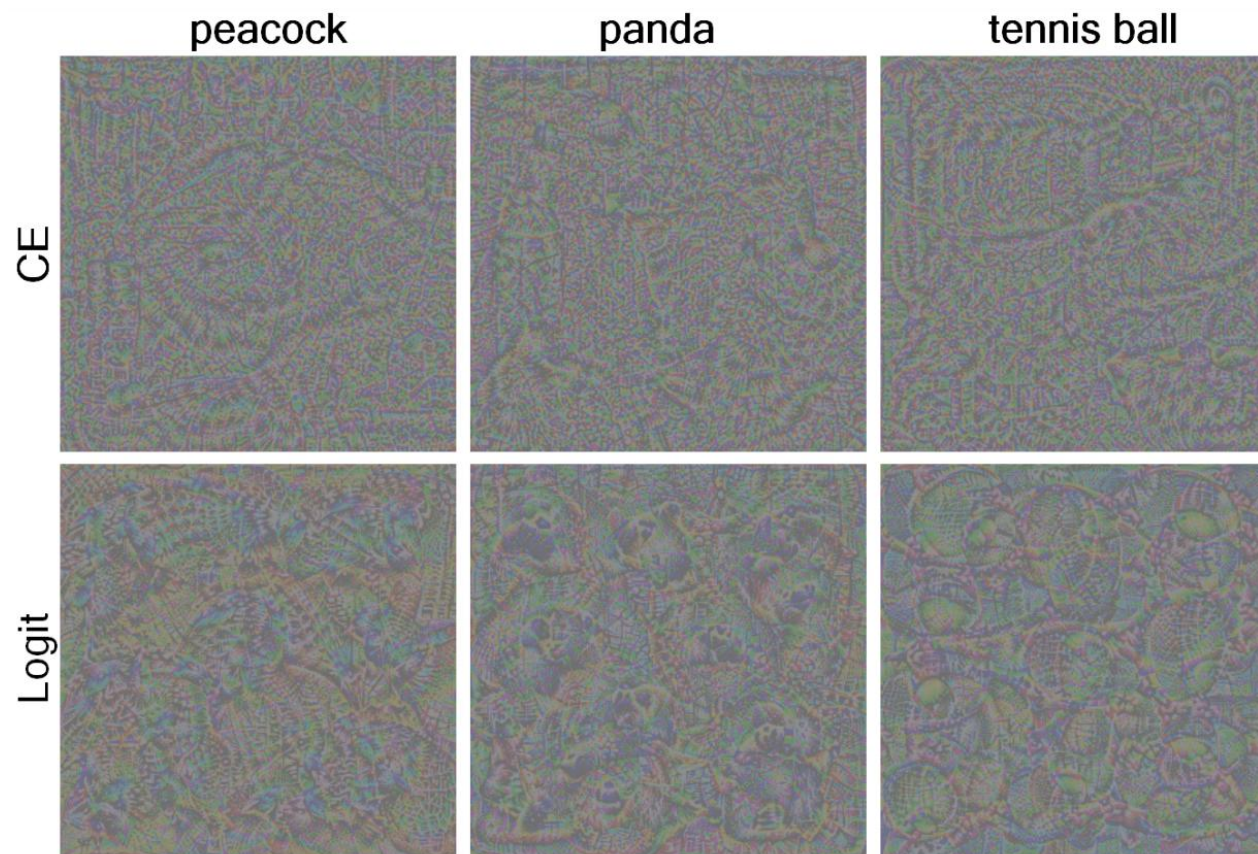


Unbounded Adversarial perturbations

# Data/Training-free Targeted UAPs

Success rates (%) of Targeted UAPs ( $\epsilon=16$ )

Attack	Inc-v3	Res50	Dense121	VGG16
CE	2.6	9.2	8.7	20.1
Logit	<b>4.7</b>	<b>22.8</b>	<b>21.8</b>	<b>65.9</b>





# Summary

- Targeted Transferability: Iterative approach ~~✗~~ Generative approach  
- More iterations >  
- Better (Logit) loss
- Better evaluation with more challenging scenarios  
- Model diversity  
- Target class diversity
- Semantic perturbations for UAPs

# Follow-ups using our iterative baseline

## 1. Slightly robust source model; Other target models (e.g. transformer, CLIP)

Springer et al. *A Little Robustness Goes a Long Way: Leveraging Robust Features for Targeted Transfer Attacks*. NeurIPS 2021.

## 2. Data/Training-free targeted UAPs with a better initialization



Li et al. Learning Universal Adversarial Perturbation by Adversarial Example. AAAI 2022



## 3. Seeking local worst-case perturbations

Qin et al. Boosting the Transferability of Adversarial Attacks with Reverse Adversarial Perturbation. OpenReview 2022

...

# Future Challenge 1

Iterative: lightweight (data-free and model-free) but slow (many iterations).  
 

Generative: heavy (additional data and models) but fast (single forward pass).  
 

## Iterative + Generative: Better source classifier with fewer iterations

1. Wu et al. *Skip Connections Matter: On the Transferability of Adversarial Examples Generated with ResNets*. ICLR 2020.
2. Guo et al. *Backpropagating Linearly Improves Transferability of Adversarial Examples*. NeurIPS 2020.
3. Zhang et al. *Backpropagating Smoothly Improves Transferability of Adversarial Examples*. CVPRW 2021
4. Zhu et al. *Rethinking Adversarial Transferability from a Data Distribution Perspective*. ICLR 2022.

# Future Challenge 2

Interpreting & improving targeted transferability on specific architectures:

Attack	Source Model: Res50			Source Model: Dense121		
	→Dense121	→VGG16	→Inc-v3	→Res50	→VGG16	→Inc-v3
Logit	<b>29.3/63.3/72.5</b>	<b>24.0/55.7/62.7</b>	<b>3.0/7.2/9.4</b>	<b>17.2/39.7/43.7</b>	<b>13.5/35.3/38.7</b>	<b>2.7/6.9/7.6</b>
Attack	Source Model: VGG16			Source Model: Inc-v3		
	→Res50	→Dense121	→Inc-v3	→Res50	→Dense121	→VGG16
Logit	<b>3.3/8.7/11.2</b>	<b>3.6/11.7/13.2</b>	<b>0.2/0.7/0.9</b>	<b>0.8/1.6/2.9</b>	<b>1.2/2.8/5.3</b>	<b>0.7/2.2/3.7</b>

