

About Me

Zhengyu Zhao (赵正宇)

✉ zhengyu.zhao@cispa.de 🏠 [zhengyuzhao.github.io](https://github.com/zhengyuzhao)

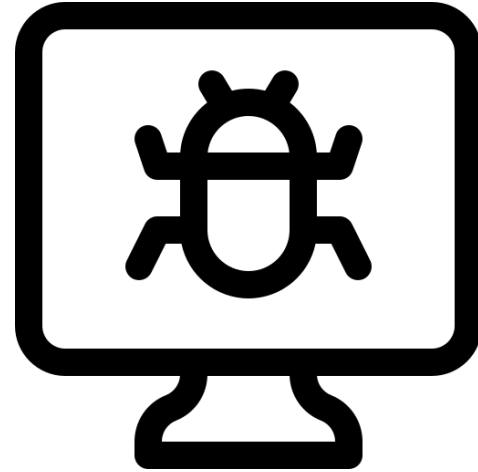
Postdoc @ CISPA Helmholtz Center for Information Security, Germany

PhD @ Radboud University, The Netherlands



Research focus:

Analyzing the vulnerability of deep neural networks to various attacks, e.g., (test-time) adversarial examples and (training-time) data poisons.



Failures of Computer Vision in Adversarial Scenarios



Outline

- Overview of adversarial images in computer vision
- Two recent projects
- Other related projects



Outline

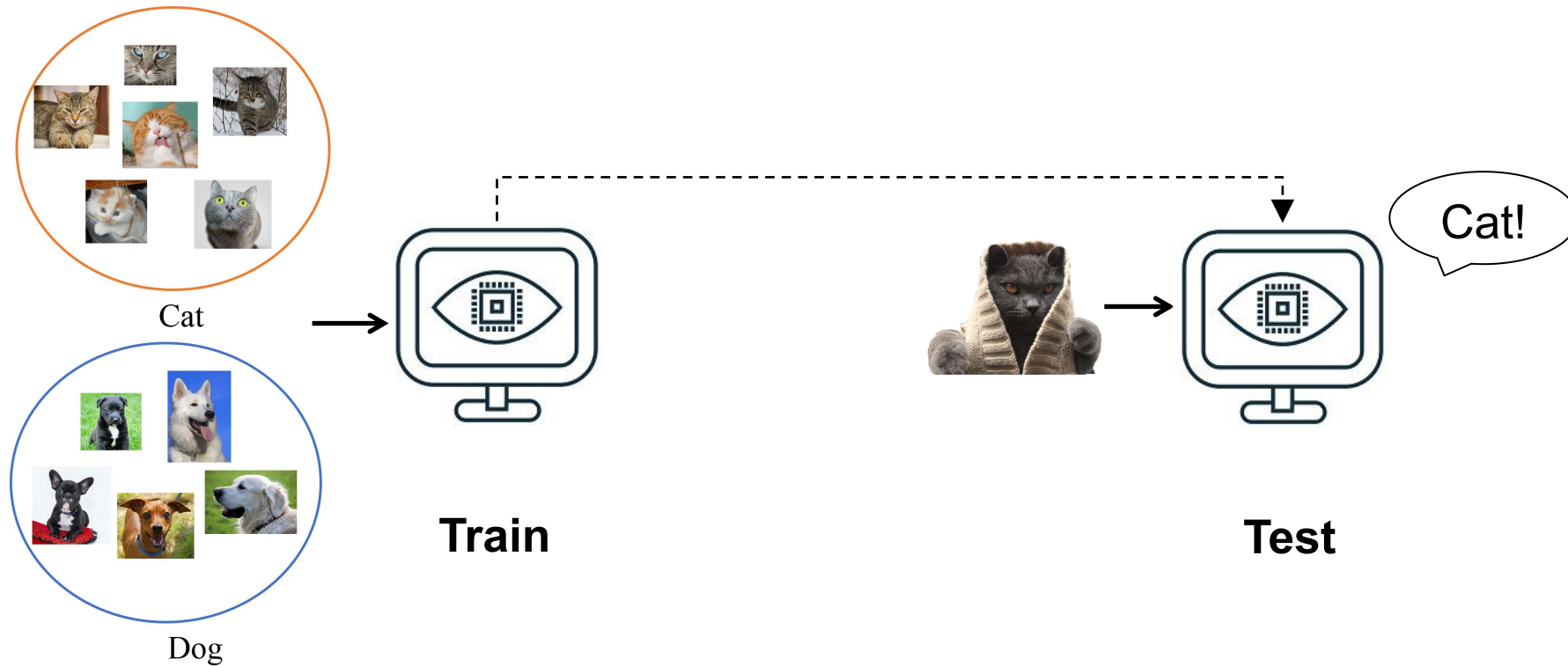
- Overview of adversarial images in computer vision
- Two recent projects
- Other related projects



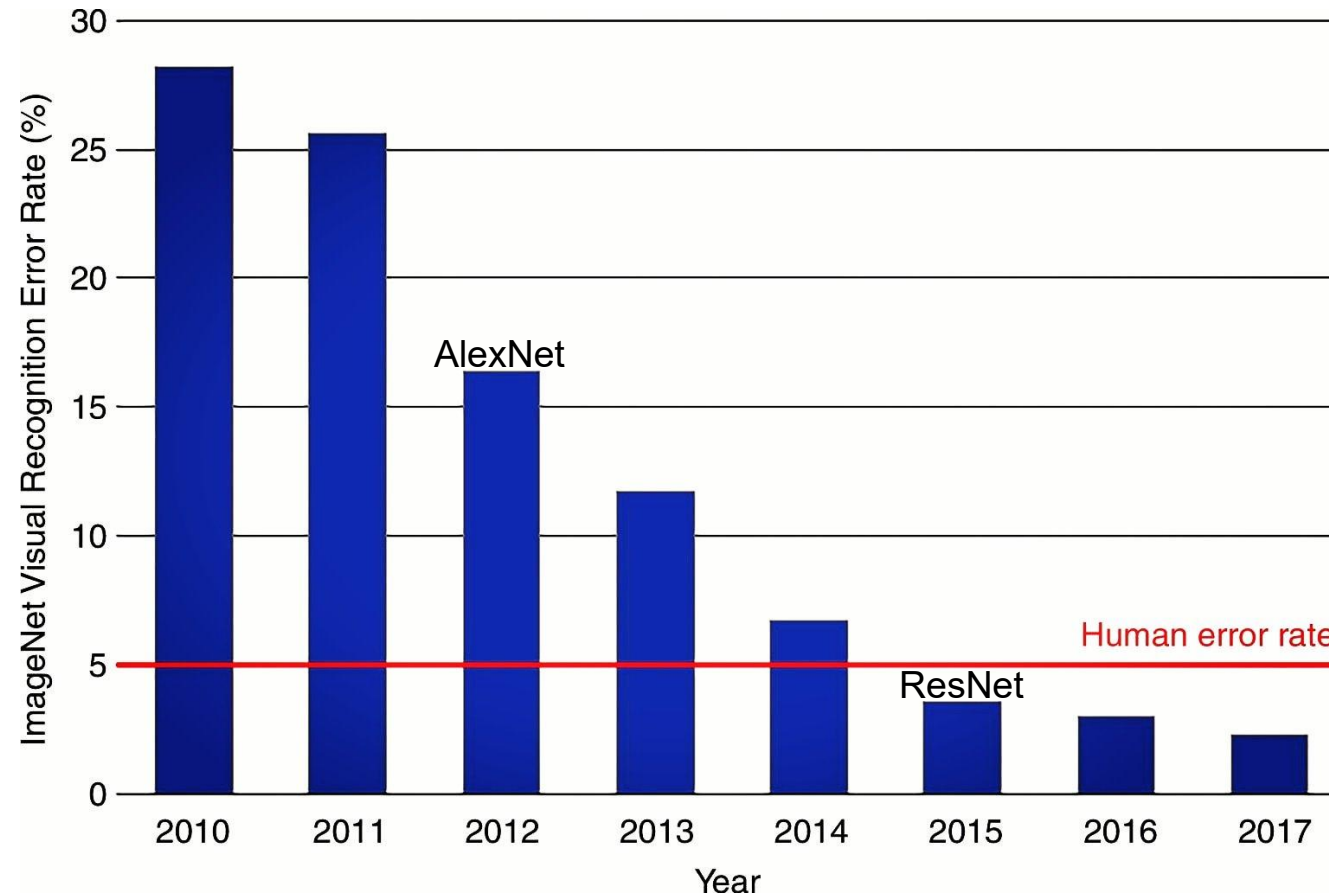
Computer Vision (CV)



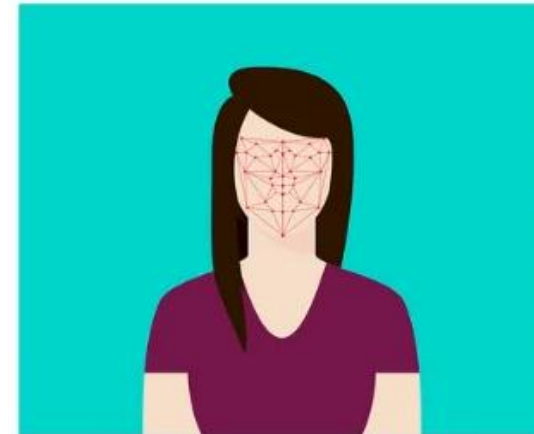
Working pipeline of CV



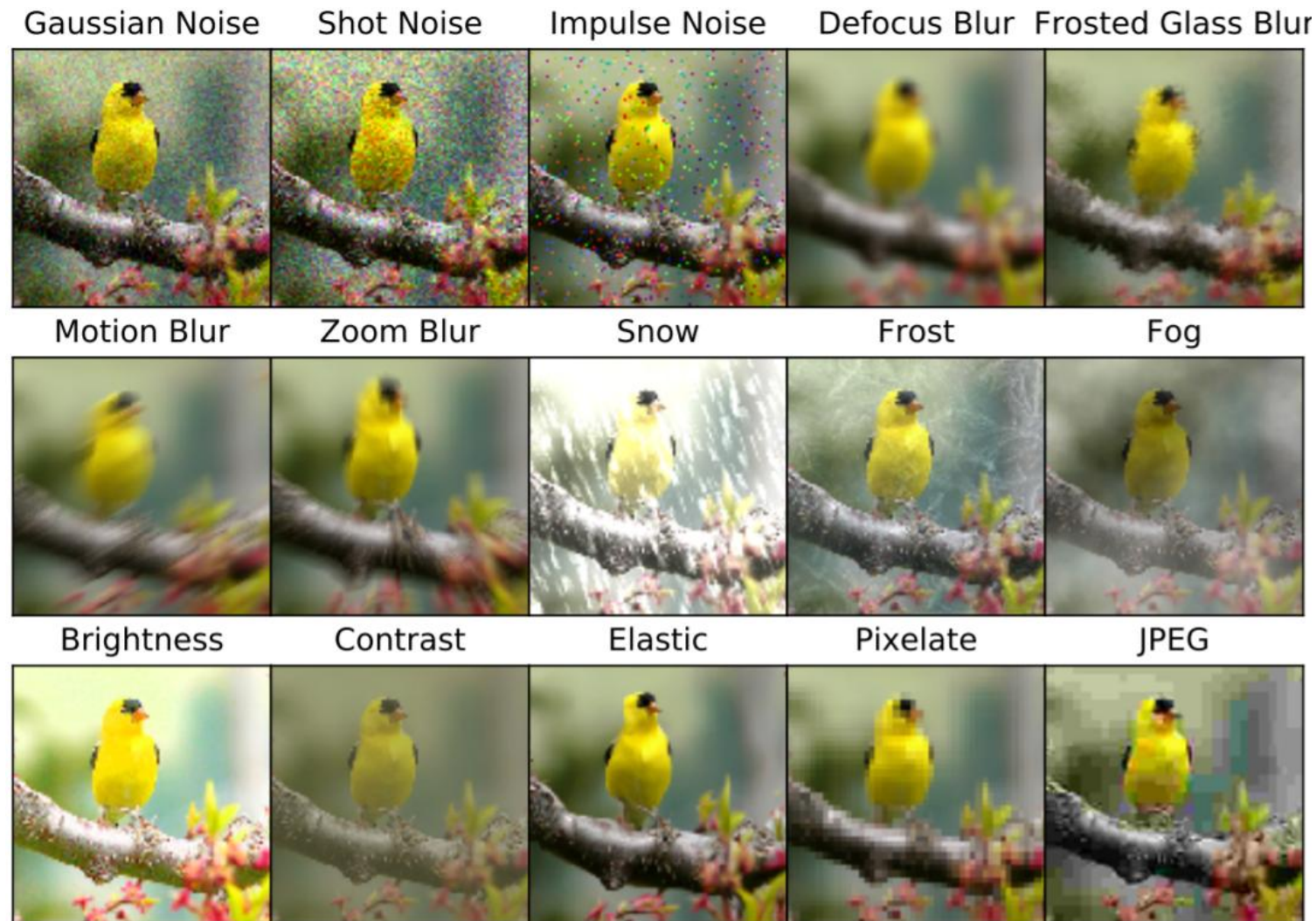
Success of CV



Success of CV



Failure of CV (against Real-world Perturbations)



Failure of CV (against Real-world Perturbations)



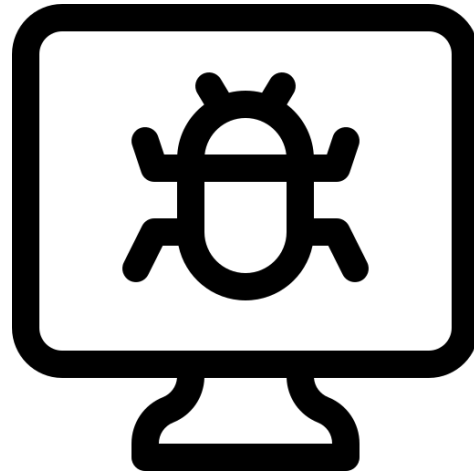
face recognition^[1]



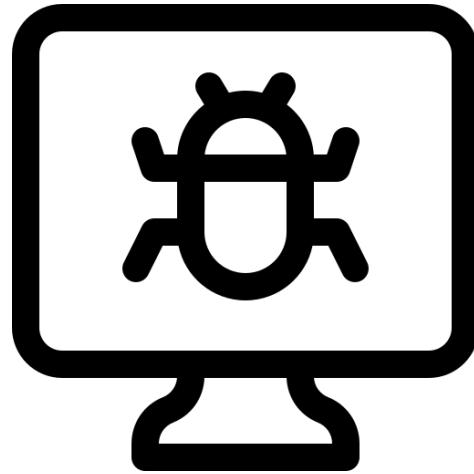
self-driving car^[2]

[1] <https://ipvm.com/reports/face-masks>

[2] <https://www.theguardian.com/technology/2018/mar/22/video-released-of-uber-self-driving-crash-that-killed-woman-in-arizona>



average-case (real-world) Image perturbations?

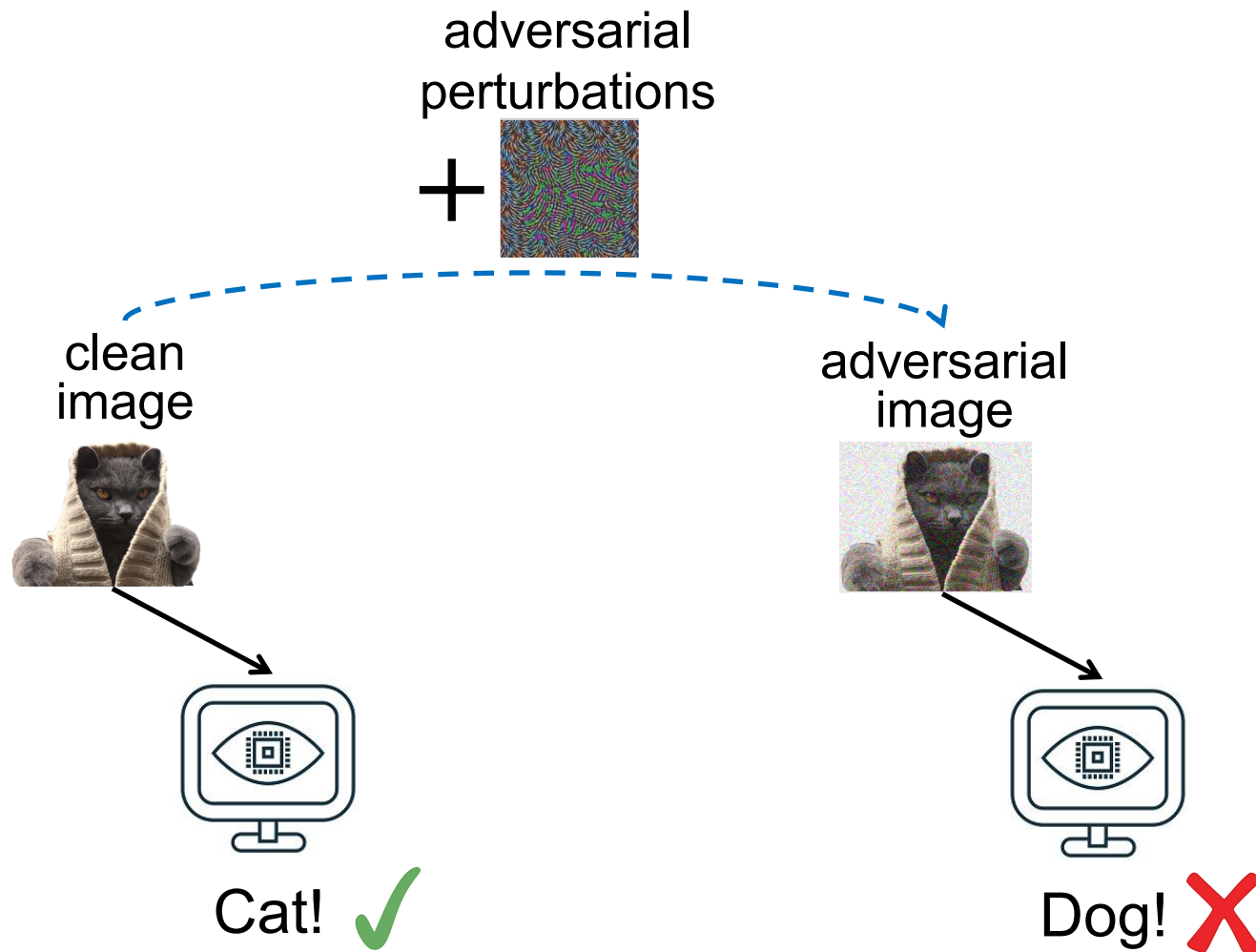


average-case (real-world) Image perturbations?



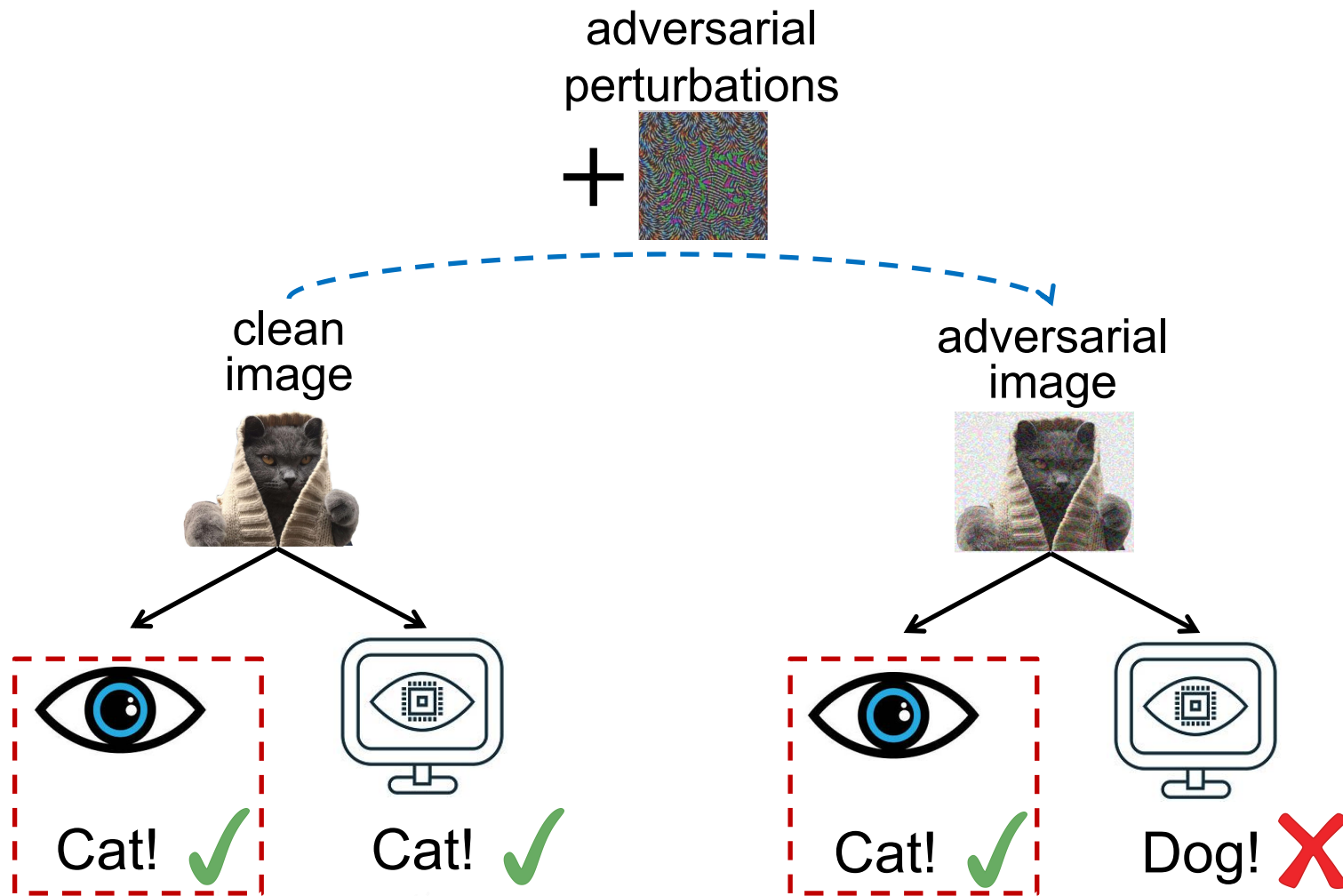
worst-case (adversarial) Image perturbations!

Formalize Adversarial Image Perturbations



change perspective

Stealthy Attacks with Imperceptible Perturbations



$$\left\| \begin{array}{c} x' \\ \text{[Image of cat with perturbations]} \end{array} - \begin{array}{c} x_{\text{cat}} \\ \text{[Image of clean cat]} \end{array} \right\|_{\infty} \leq \epsilon$$

Real-world → Adversarial Image Perturbations



face recognition^[1]



adversarial mask^[2]

[1] <https://ipvm.com/reports/face-masks>

[2] <https://towardsdatascience.com/fooling-facial-detection-with-fashion-d668ed919eb>

Real-world → Adversarial Image Perturbations



self-driving car^[1]



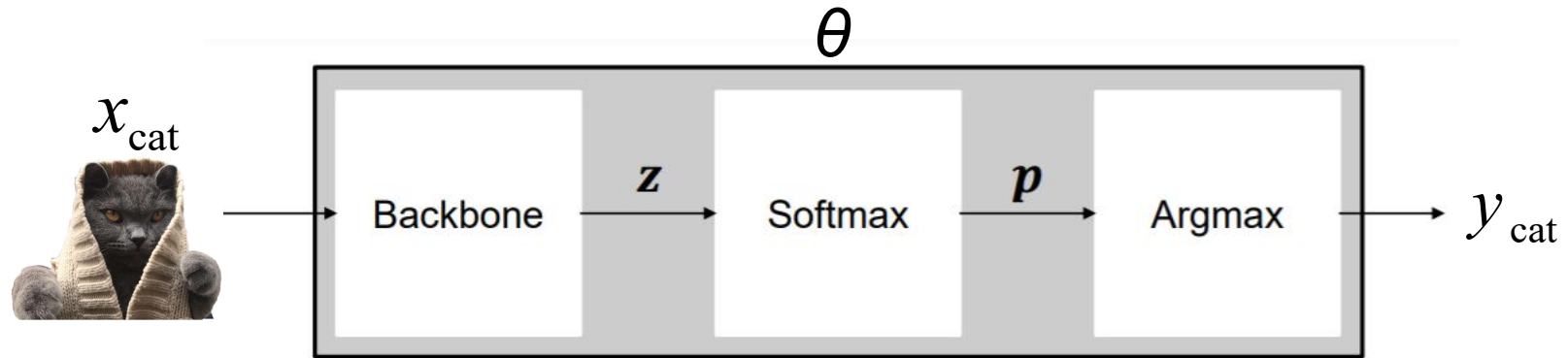
adversarial graffiti^[2]

[1] <https://www.theguardian.com/technology/2018/mar/22/video-released-of-uber-self-driving-crash-that-killed-woman-in-arizona>

[2] Eykholt et al. *Robust physical-world attacks on deep learning visual classification*. CVPR 2018.

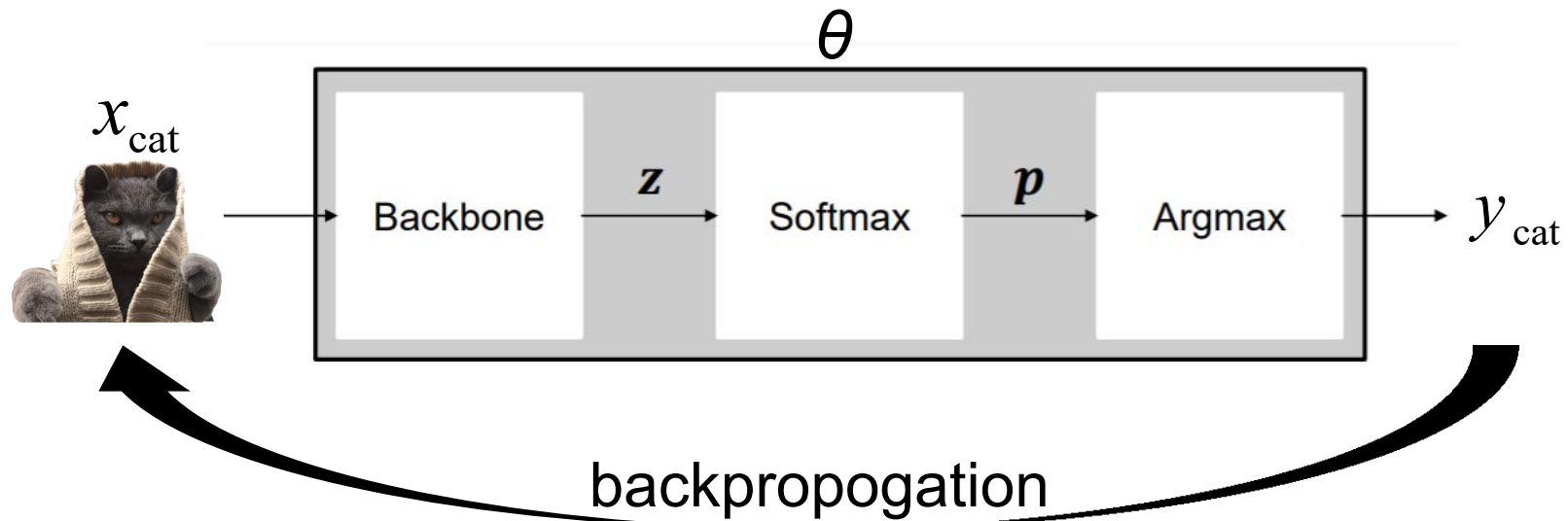


Optimize Adversarial Images x'



$$\theta' = \arg \min_{\theta} J(\theta, x_{\text{cat}}, y_{\text{cat}})$$

Optimize Adversarial Images x'

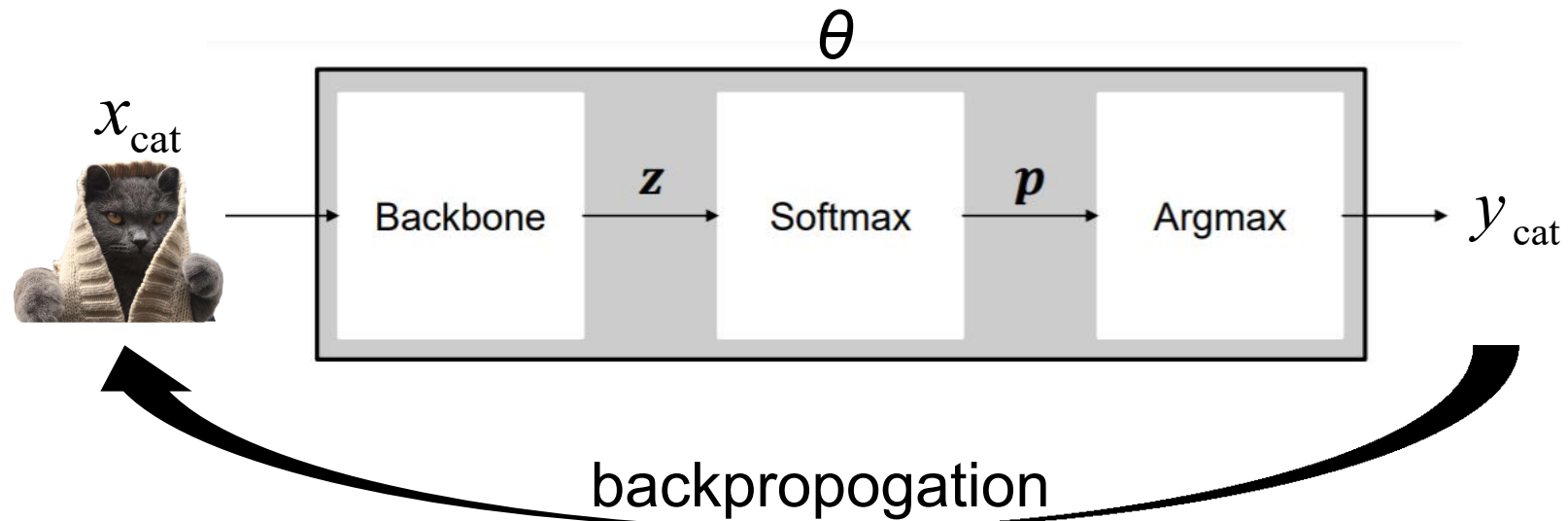


$$\theta' = \arg \min_{\theta} J(\theta, x_{\text{cat}}, y_{\text{cat}})$$



$$x' = \arg \min_x J(\theta_o, x, y_t) \quad \text{targeted}$$

Optimize Adversarial Images x'



$$\theta' = \arg \min_{\theta} J(\theta, x_{\text{cat}}, y_{\text{cat}})$$



$$x' = \arg \min_x J(\theta_o, x, y_t) \quad \text{targeted}$$

$$x' = \arg \max_x J(\theta_o, x, y_{\text{cat}}) \quad \text{non-targeted}$$

Optimize Adversarial Images x'

$$\text{Objective: } x' = \arg \min_x J(\theta_o, x, y_t) \quad \text{s.t.} \quad \|x' - x_{\text{cat}}\|_{\infty} \leq \varepsilon$$

Optimization: Iterative-Fast Gradient **Sign** Method (I-FGSM)^[1]

$$x'_0 = x_{\text{cat}}, \quad x'_{i+1} = x'_i - \text{sign}(\nabla_x J(x'_i, y_t))$$

$$x'_{i+1} \leftarrow \text{clip}(x'_{i+1} - x_{\text{cat}}, -\varepsilon, \varepsilon)$$

Recap

Success of computer vision

↳ Failures against real-world perturbations

↳ ... adversarial images

↳ optimize adversarial images

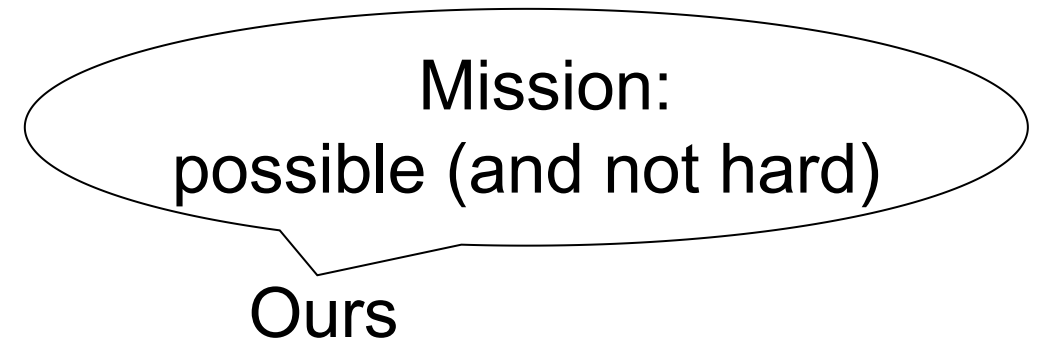
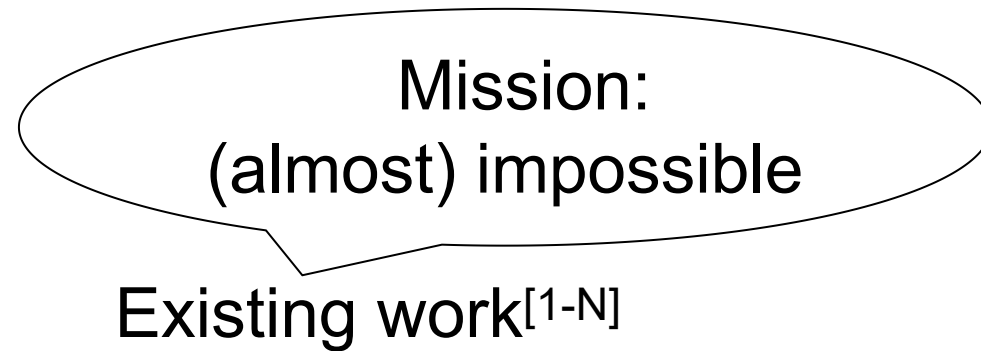


Outline

- Overview of adversarial images in computer vision
- **Two recent projects**
- Other related projects



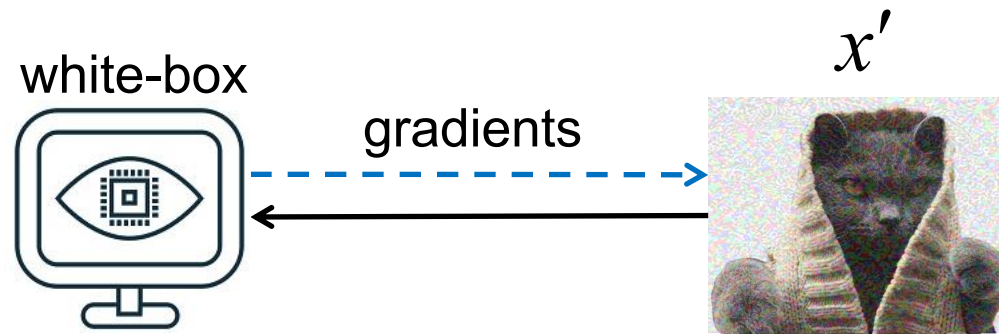
Consensus-Challenging Insights



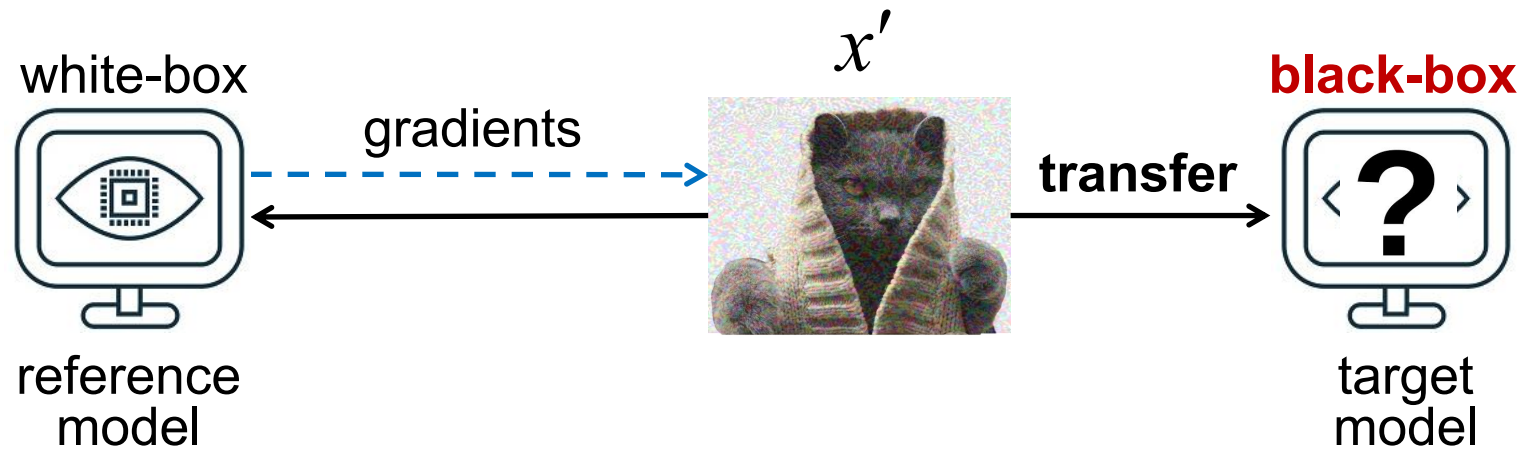
Project 1. Transferable Targeted Attacks



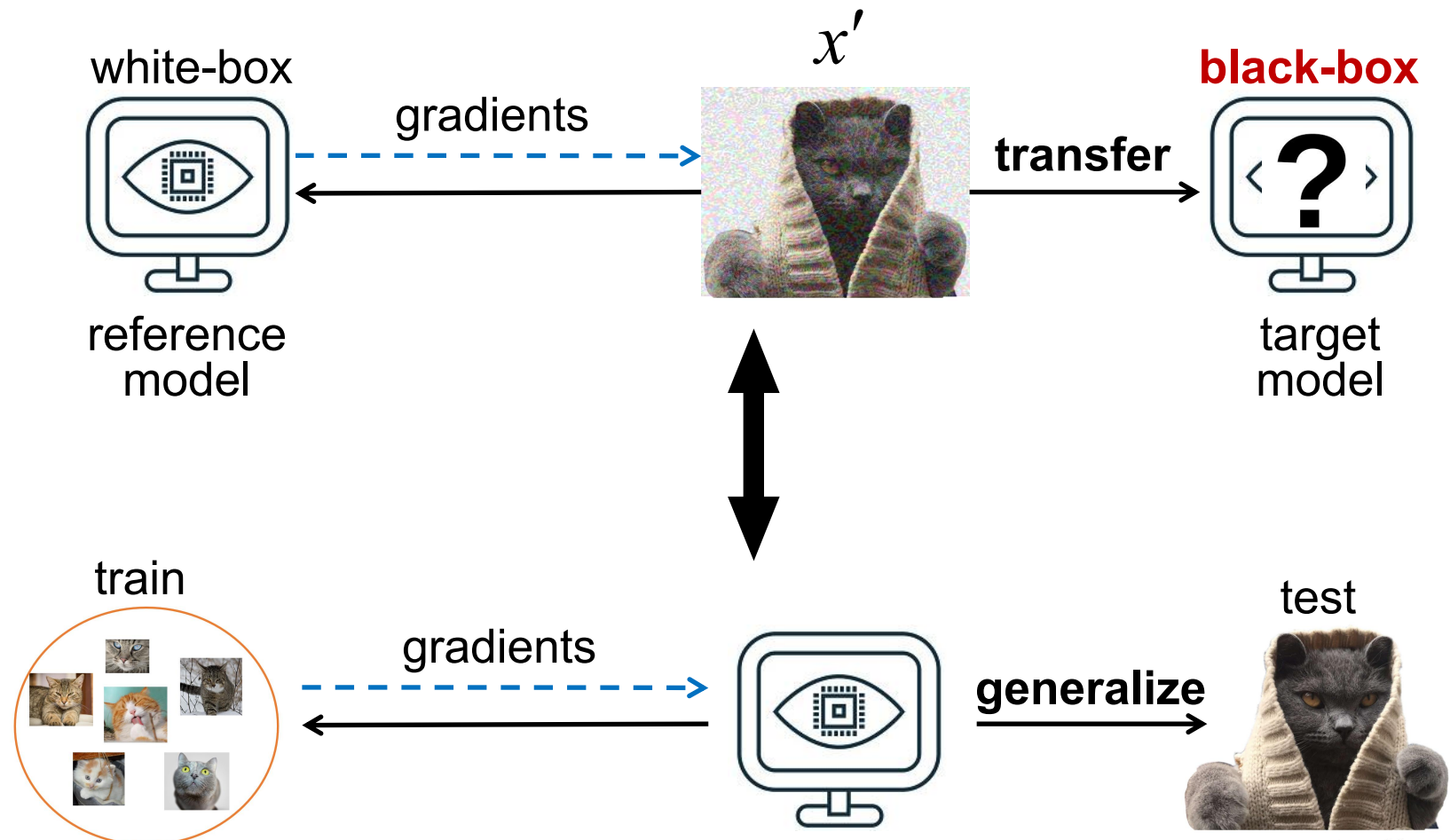
Transferable Targeted Attacks



Transferable Targeted Attacks



Transferable Targeted Attacks



Transfer Techniques

- Gradient stabilization

e.g., momentum-based (MI-FGSM)^[1]:

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

- Data augmentation

e.g., resizing & padding (DI-FGSM)^[2]
translation (TI-FGSM)^[3]:

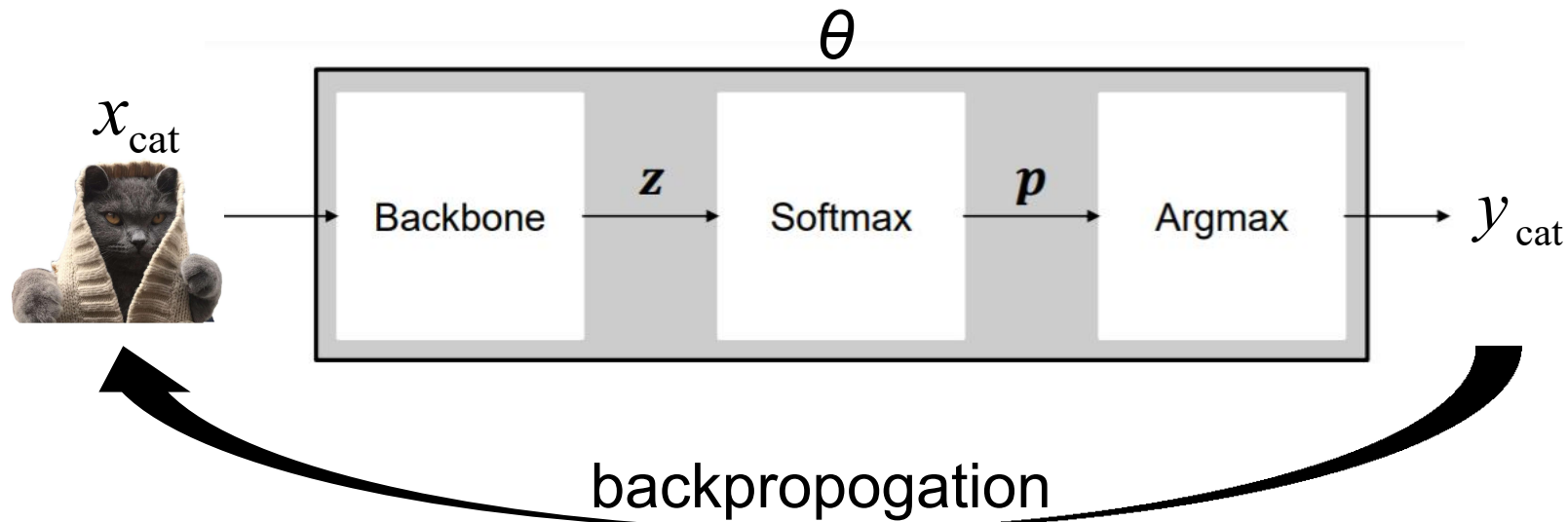
$$\mathbf{x}'_{i+1} = \mathbf{x}'_i - \alpha \cdot \text{sign}(\nabla_{\mathbf{x}} J(T(\mathbf{x}'_i, p), y_t))$$

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

[2] Xie et al. *Improving Transferability of Adversarial Examples with Input Diversity*. CVPR 2019

[3] Dong et al. *Evading defenses to transferable adversarial examples by translation-invariant attacks*. CVPR 2019.

Transferable Targeted Attacks



$$\theta' = \arg \min_{\theta} J(\theta, x_{\text{cat}}, y_{\text{cat}})$$



$$x' = \arg \max_x J(\theta_o, x, y_{\text{cat}}) \quad \text{non-targeted} \text{ 😊}$$

$$x' = \arg \min_x J(\theta_o, x, y_t) \quad \text{targeted} \text{ 😞}$$

Consensus-Challenging Insight

Impossible for I-FGSM
to achieve **targeted**
transferability.

Existing work^[1-6]

Possible and
even SOTA.

Ours

[1] Liu et al. *Delving into transferable adversarial examples and black-box attacks*. ICLR 2017.

[2] Dong et al. *Boosting Adversarial Attacks with Momentum*. CVPR 2018.

[3] Inkawhich et al. *Feature space perturbations yield more transferable adversarial examples*. CVPR 2019.

[4] Inkawhich et al. *Transferable perturbations of deep feature distributions*. ICLR 2020.

[5] Inkawhich et al. *Perturbing across the feature hierarchy to improve standard and strict blackbox attack transferability*. NeurIPS 2020.

[6] Naseer et al. *On generating transferable targeted perturbations*. ICCV 2021.



Fix I-FGSM: Step 1. Ensemble (0% → 15%)

ResNet50 → DenseNet121 (Iter. =10)

I-FGSM: ~0%

MI-FGSM: ~0.5%

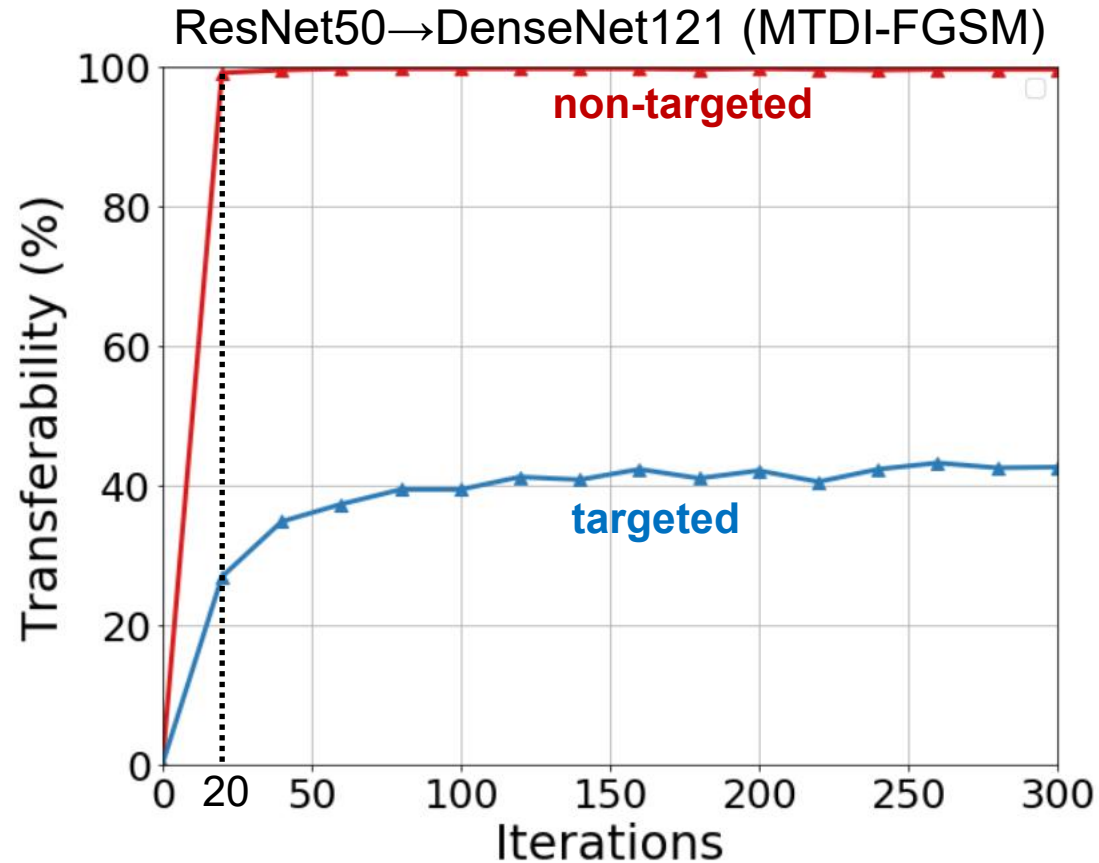
TI-FGSM: ~0.5%

DI-FGSM: ~5%

MTDI-FGSM: ~15%

single technique in existing work

Fix I-FGSM: Step 2. More Iterations (15% → 42%)



<20 iterations in existing work:

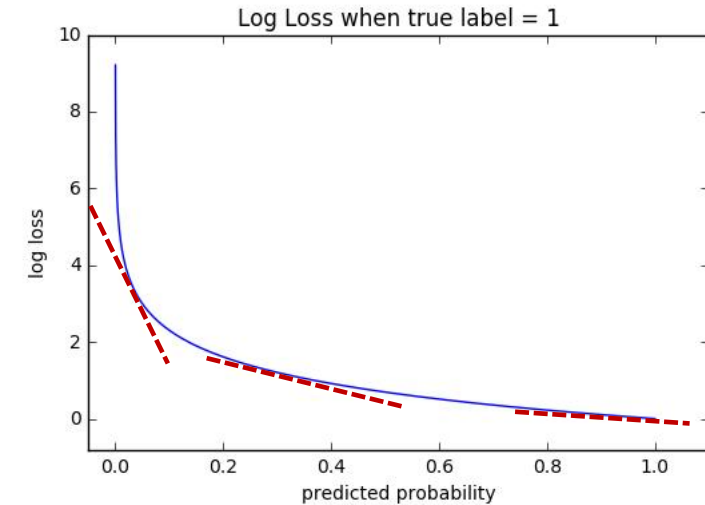
- fail to converge
- efficiency is not important



Fix I-FGSM: Step 3. Suitable Loss

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

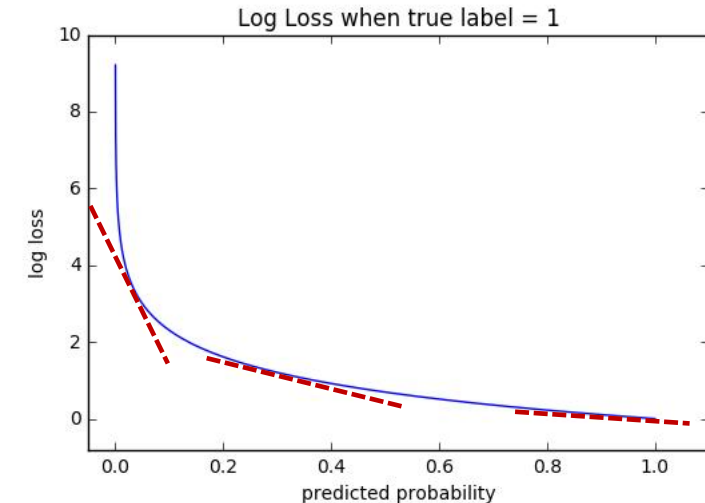
$$L_{CE} = -1 \cdot \log(p_t) = -\log\left(\frac{e^{z_t}}{\sum e^{z_j}}\right) = -z_t + \log\left(\sum e^{z_j}\right),$$
$$\frac{\partial L_{CE}}{\partial z_t} = -1 + \frac{\partial \log(\sum e^{z_j})}{\partial e^{z_t}} \cdot \frac{\partial e^{z_t}}{\partial z_t} = -1 + \frac{e^{z_t}}{\sum e^{z_j}} =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



Fix I-FGSM: Step 3. Suitable Loss

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

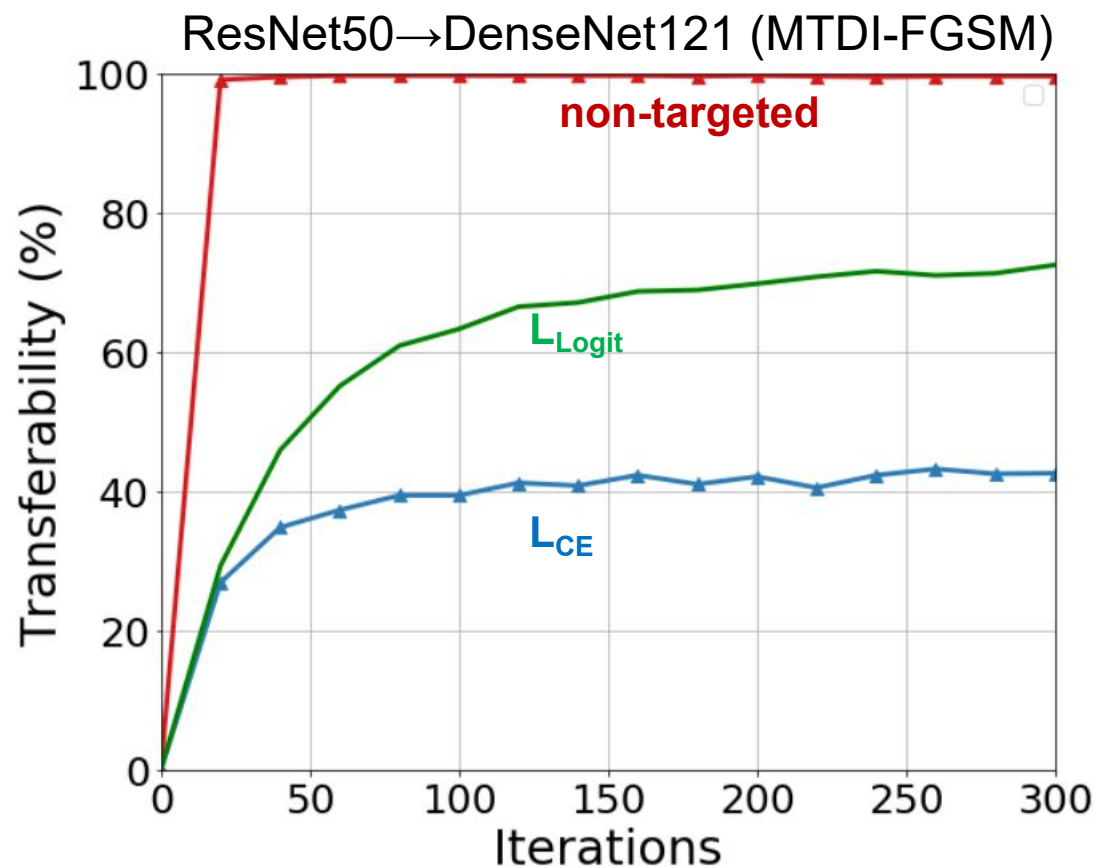
$$L_{CE} = -1 \cdot \log(p_t) = -\log\left(\frac{e^{z_t}}{\sum e^{z_j}}\right) = \underline{-z_t} + \log\left(\sum e^{z_j}\right),$$
$$\frac{\partial L_{CE}}{\partial z_t} = -1 + \frac{\partial \log(\sum e^{z_j})}{\partial e^{z_t}} \cdot \frac{\partial e^{z_t}}{\partial z_t} = -1 + \frac{e^{z_t}}{\sum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



Logit Loss (L_{Logit}):

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

Fix I-FGSM: Step 3. Suitable Loss (42% \rightarrow 72%)



Other Analyses: Real-World Attacks

Services	Evaluation	Ori	CE	Po+Trip	Logit
Object localization	non-targeted	31.50	53.00	51.75	62.50
	targeted	0	9.00	8.50	19.25
Label detection	non-targeted	9.75	34.00	22.50	35.00
	targeted	0	4.50	2.25	6.25

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%

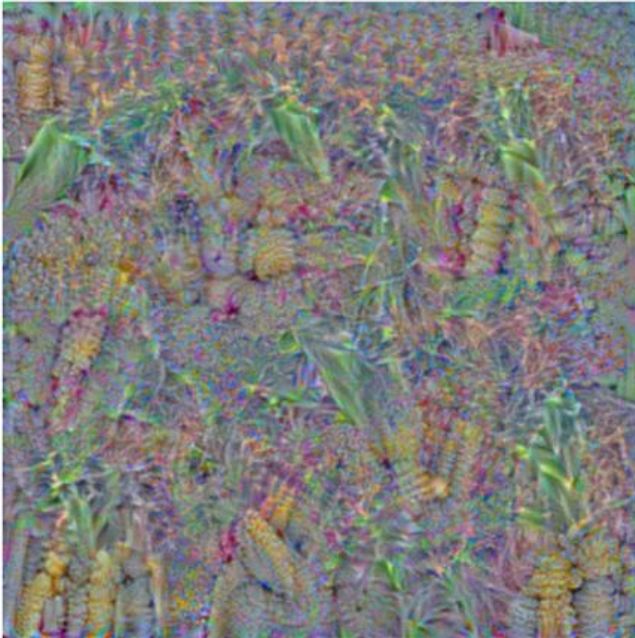
✗

$y_t = \text{“yawl” (a type of boat)}$

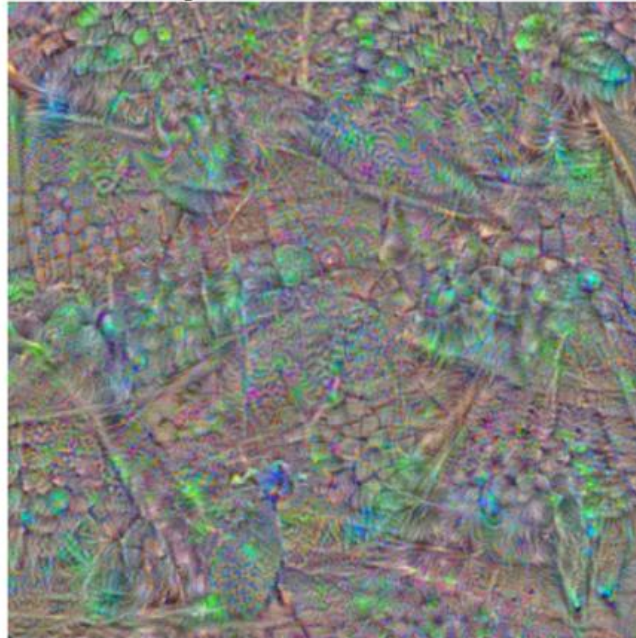
Other Analyses: Perturbation Semantics

target label:

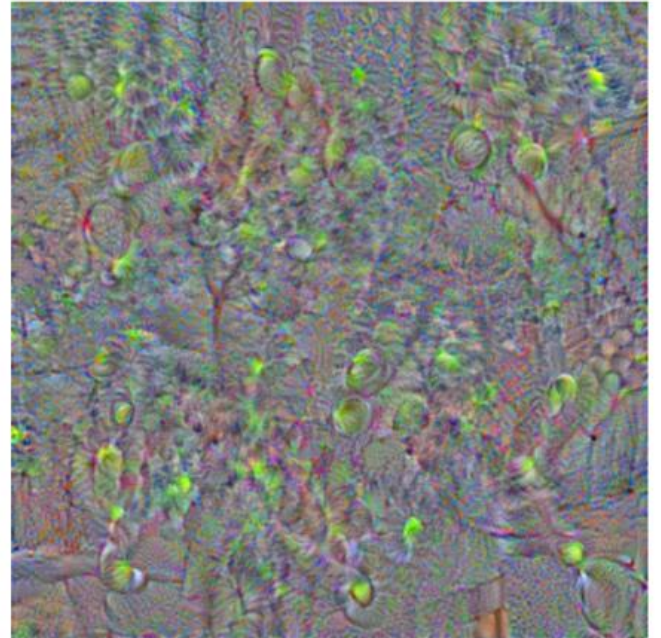
“corn”



“peacock”

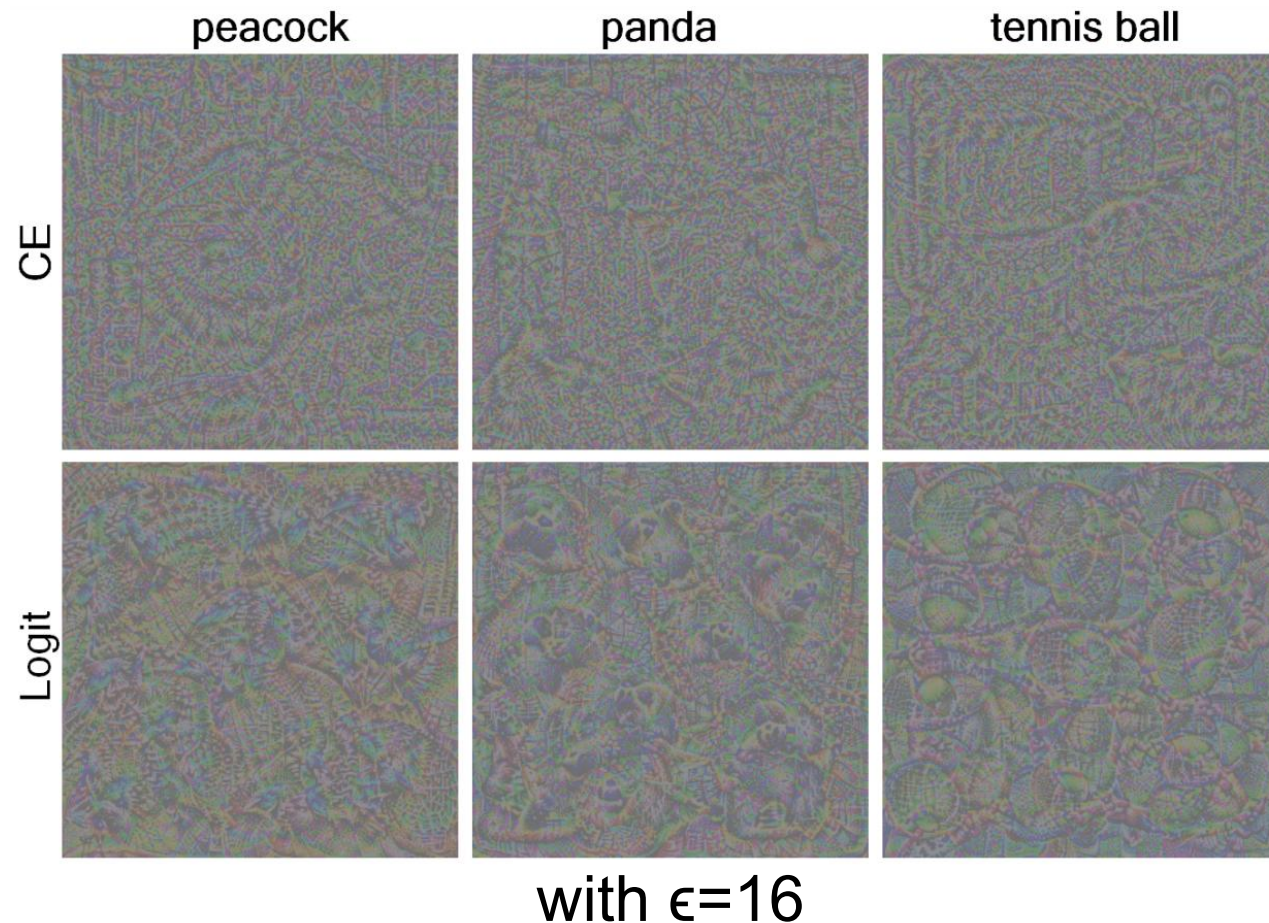


“tennis ball”



without ϵ

Other Analyses: Targeted Universal Perturbations^[1]

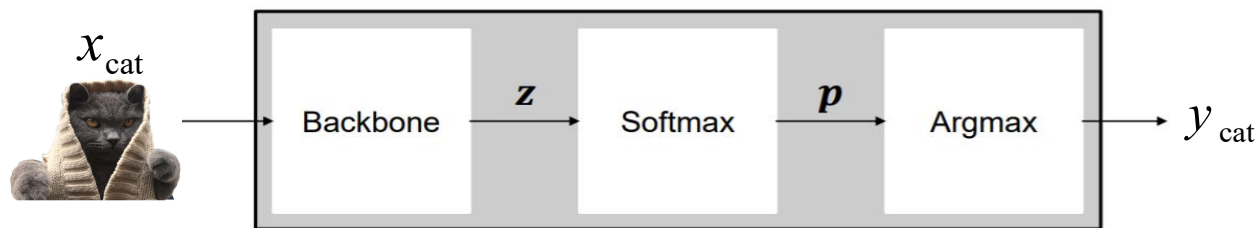


Success rates (%)

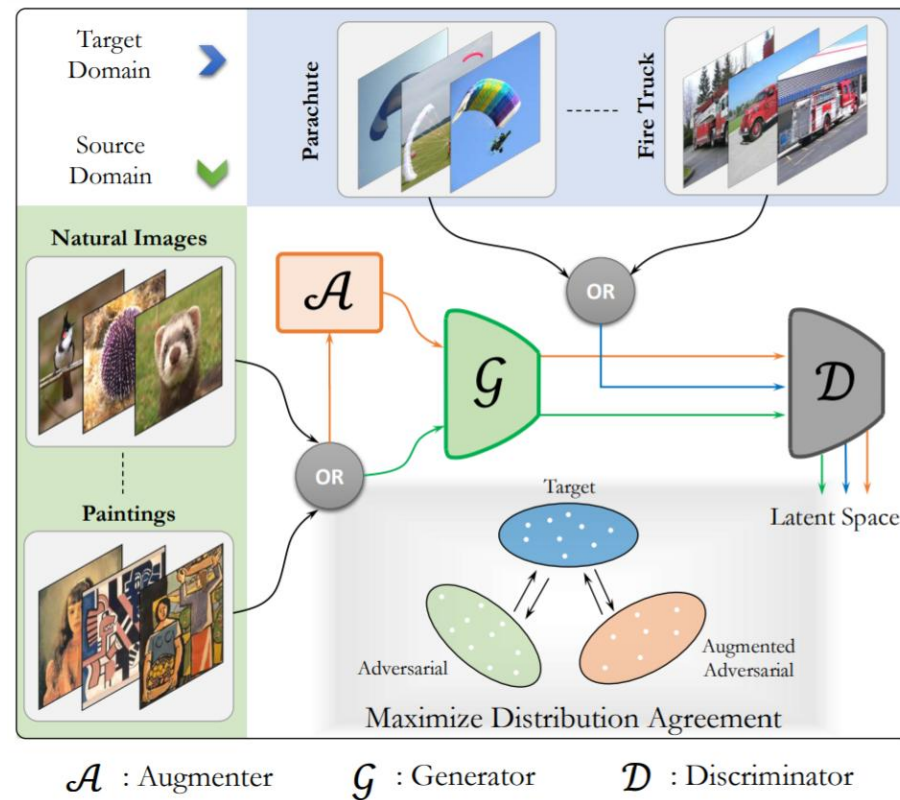
Attack	Inc-v3	Res50	Dense121	VGG16
CE	2.6	9.2	8.7	20.1
Logit	4.7	22.8	21.8	65.9

Other Analyses: I-FGSM (ours) vs. Generative (SOTA)

Ours



Generative^[1]



- Data: Single Input image
- Model: 1 \times surrogate classifier

vs

- Massive training data
- 1000 \times target-specific generators

Other Analyses: I-FGSM (ours) vs. Generative (SOTA)

Targeted Transferability (%)

Bound	Attack	D121	V16	D121-ens	V16-ens
$\epsilon = 16$	SOTA	79.6	78.6	92.9	89.6
	ours	75.9	72.5	99.4	97.7
$\epsilon = 8$	SOTA	37.5	46.7	63.2	66.2
	ours	44.5	46.8	92.6	87.0

Summary of Project 1

- 3 steps to fix I-FGSM
 - ensemble
 - more iterations
 - suitable (logit) loss
- Other Analyses
 - real-world attacks
 - universal perturbations
 - I-FGSM (data/training-free) vs. generative

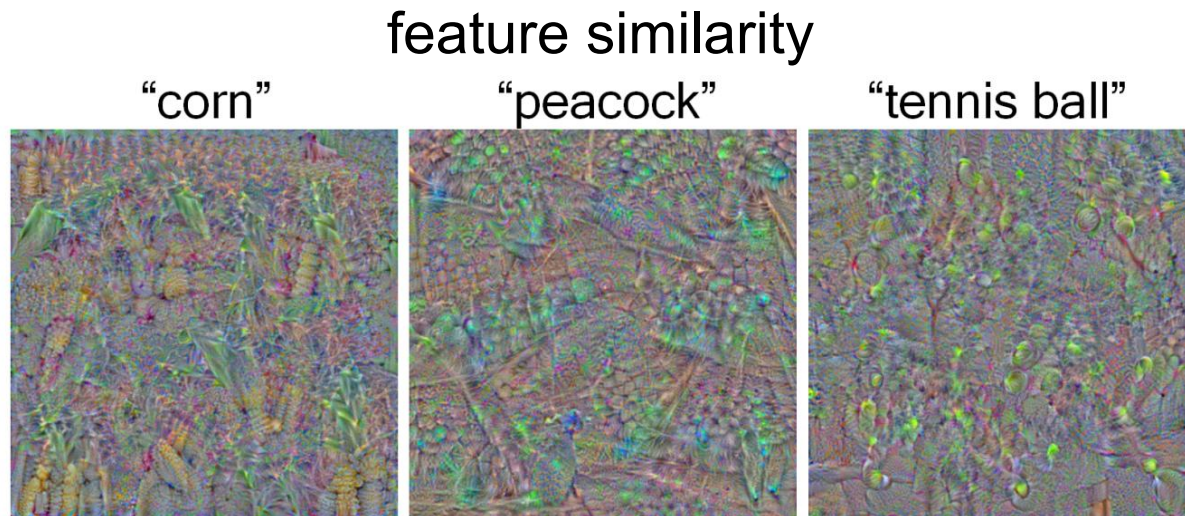
Summary of Project 1

- 3 steps to revive I-FGSM
 - ensemble
 - more iterations
 - suitable (logit) loss
- Other Analyses
 - real-world attacks
 - universal perturbations
 - I-FGSM (data/training-free) vs. generative

"God is in the details"

Future Work

- Explaining transferability



Or model similarity

Res50 → Dense121: ~70% 😊
Res50 → Incv3: ~10% 😭

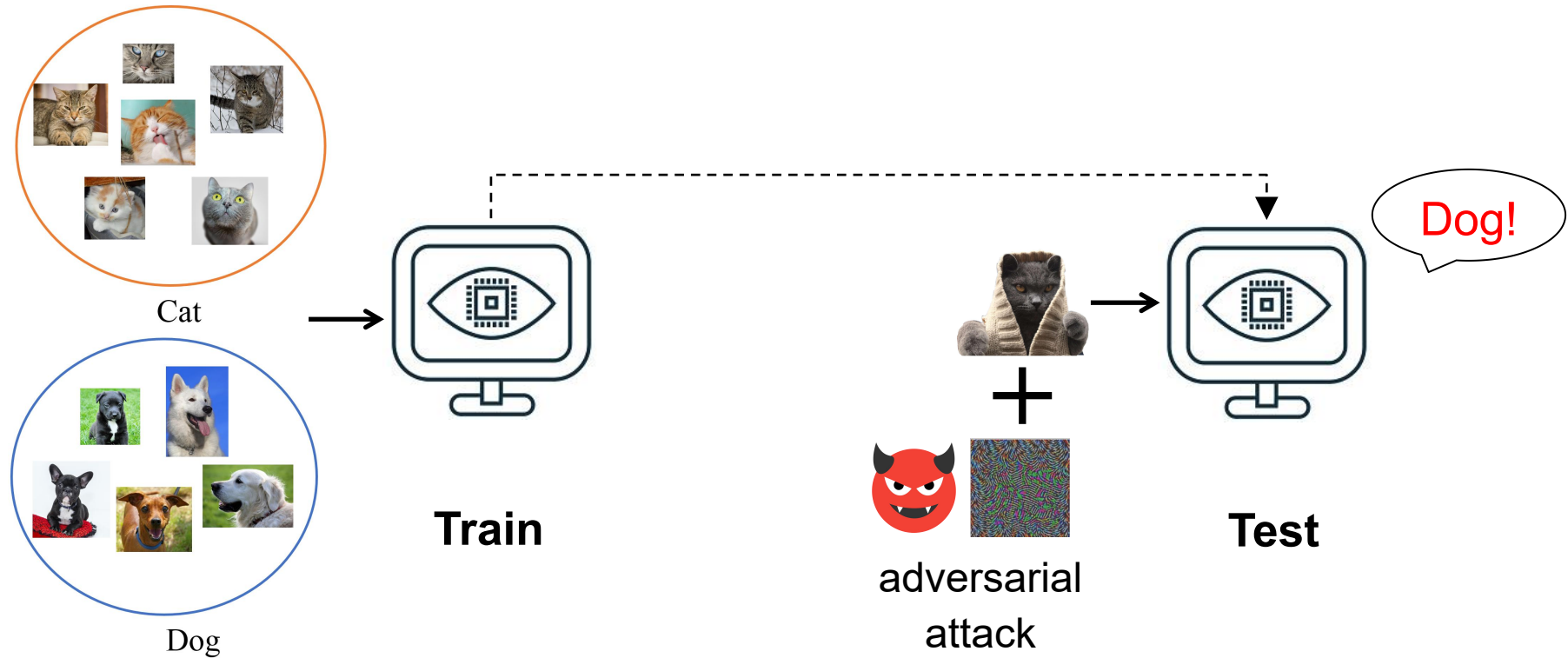
- Benchmarking transferability

📖 Zhao et al. *Towards Good Practices in Evaluating Transfer Adversarial Attacks*. arXiv 2022

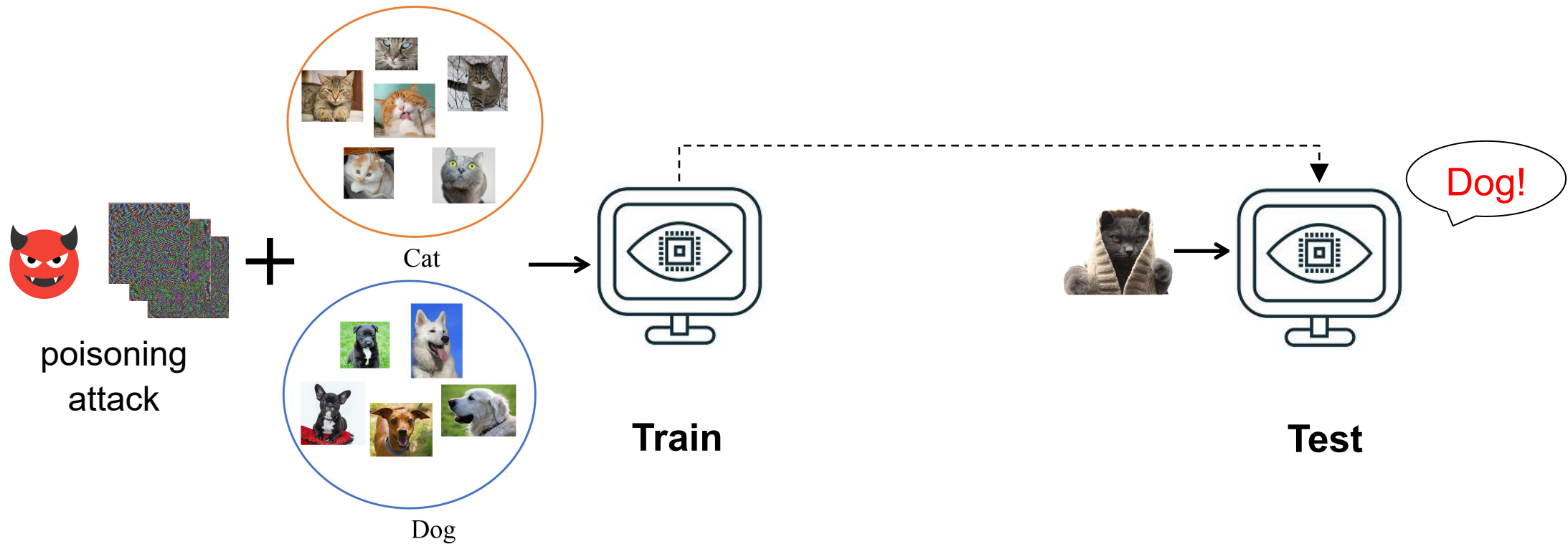
🔗 <https://github.com/ZhengyuZhao/TransferAttackEval>

- Systematic categorization of 40+ transfer attacks
- 23 representative attacks against 9 representative defenses on ImageNet
- Consensus-challenging insights

Testing-Stage Attack



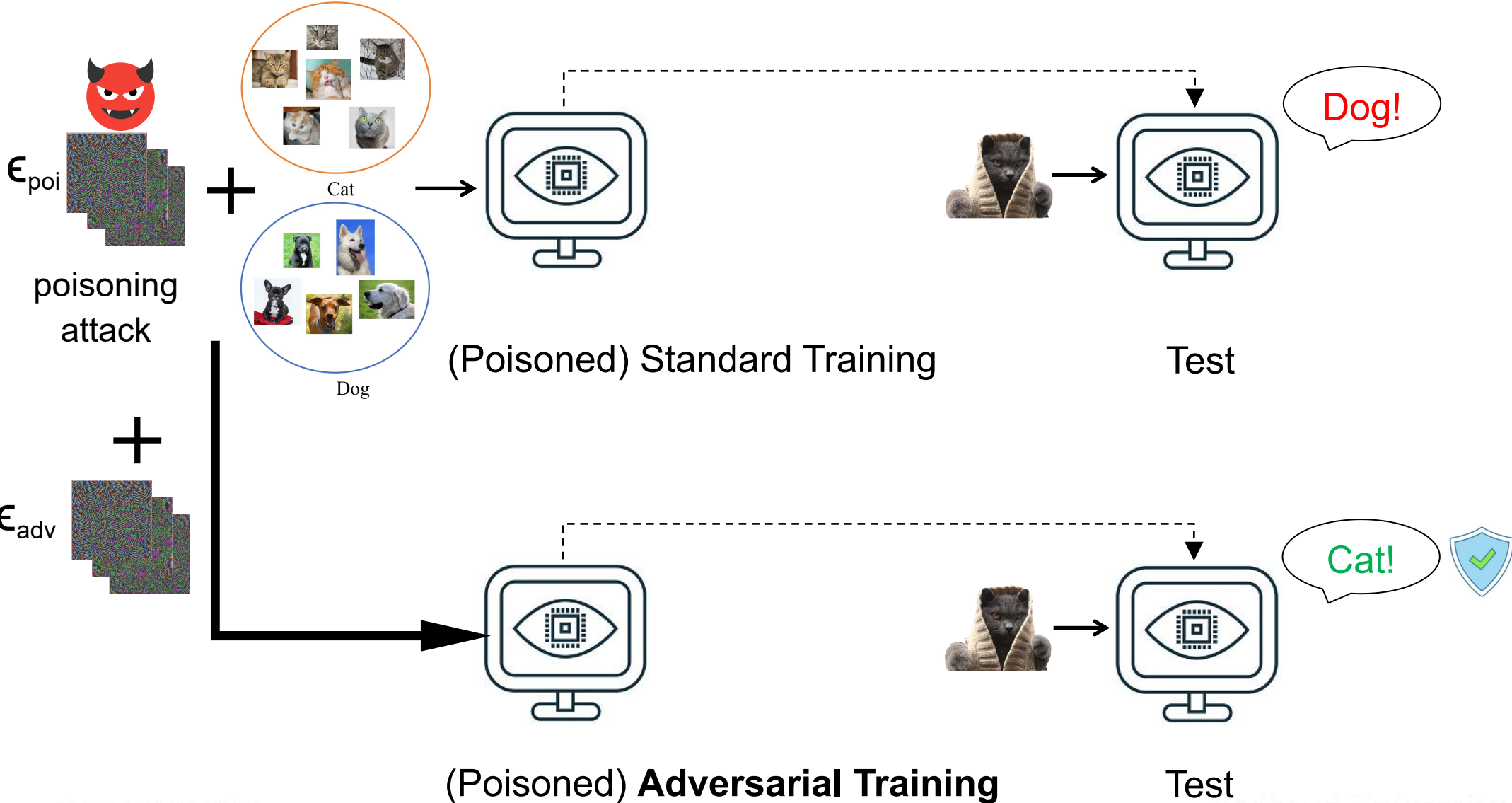
Training-Stage Attack



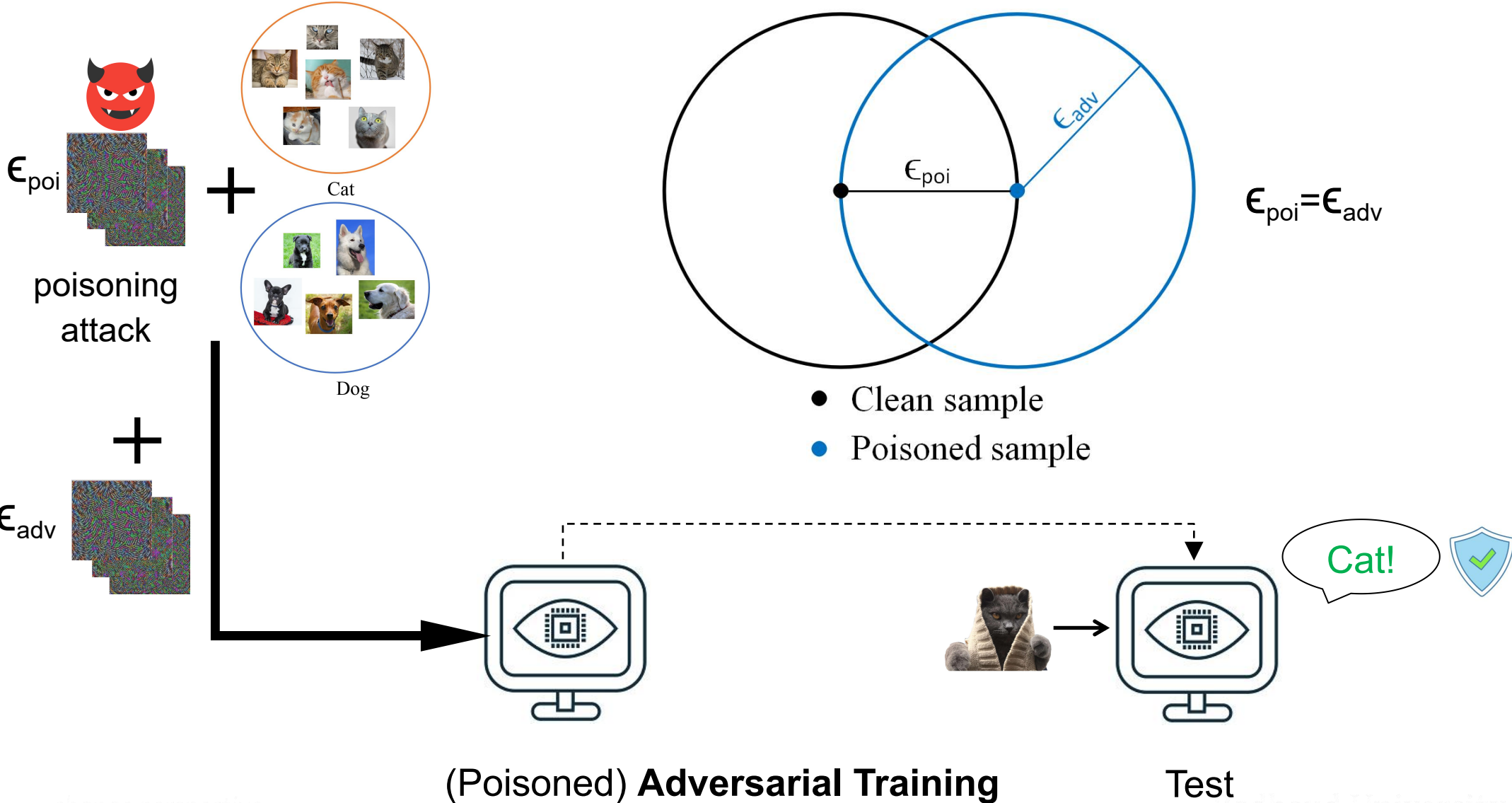
Project 2. Poisoning Against Adversarial Training



Adversarial Training-based Defense



Adversarial Training-based Defense



change perspective

[1] Tao et al. *Better Safe Than Sorry: Preventing Delusive Adversaries with Adversarial Training*. NeurIPS 2021.



Consensus-Challenging Insight

Impossible to poison
AT models

Existing work^[1-6]

Possible (with a new
attack strategy)

Ours

[1] Fowl et al. *Adversarial Examples Make Strong Poisons*. NeurIPS 2021.

[2] Huang et al. *Unlearnable Examples: Making Personal Data Unexploitable*. ICLR 2021.

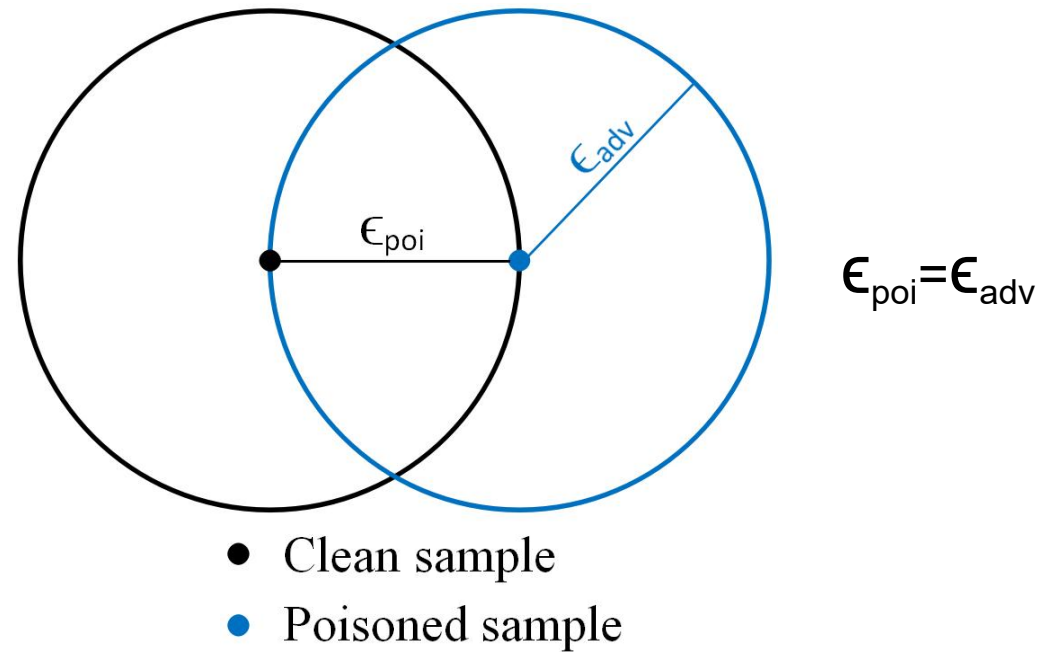
[3] Tao et al. *Better Safe Than Sorry: Preventing Delusive Adversaries with Adversarial Training*. NeurIPS 2021.

[4] Wang et al. *Fooling Adversarial Training with Inducing Noise*. arXiv 2021.

[5] Fu et al. *Robust Unlearnable Examples: Protecting Data Against Adversarial Learning*. ICLR 2022.

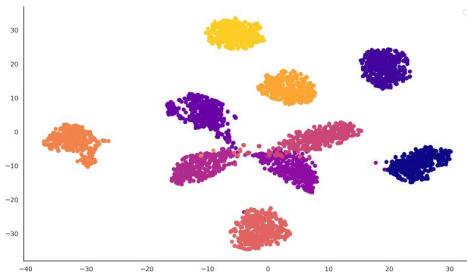
[6] Tao et al. *Can Adversarial Training Be Manipulated By Non-Robust Features?* NeurIPS 2022.

Consensus-Challenging Insight



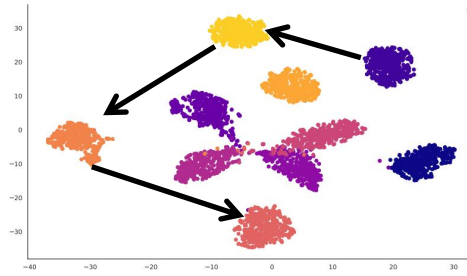
Existing Poisoning

clean training

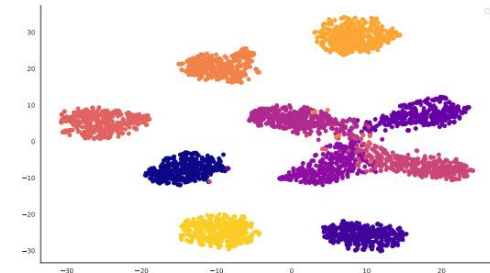


Test Acc: 84.88%

existing poisoning



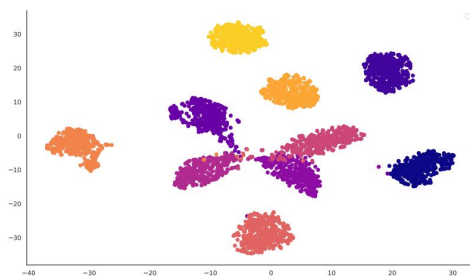
$$x' = \arg \min_x J(x, y_t)$$



Test Acc: 83.11% 😈

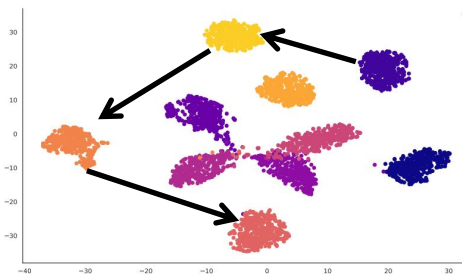
Our Poisoning

clean training

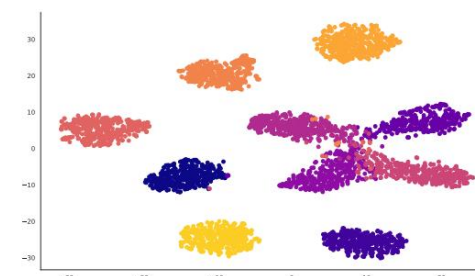


Test Acc: 84.88%

existing poisoning

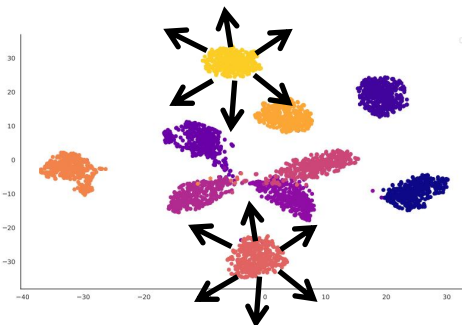


$$x' = \arg \min_x J(x, y_t)$$



Test Acc: 83.11% 😈

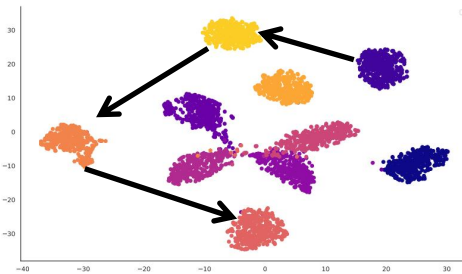
our poisoning



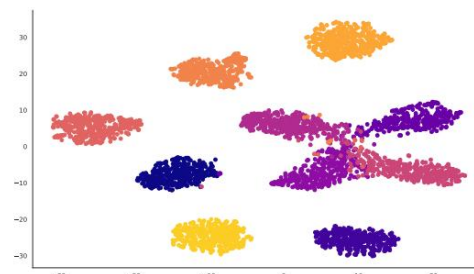
change perspective

Our Poisoning

existing poisoning

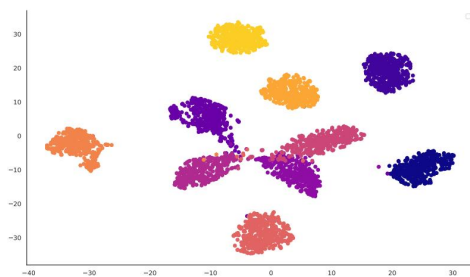


$$x' = \arg \min_x J(x, y_t)$$



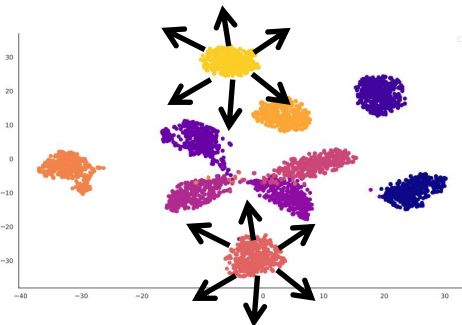
Test Acc: 83.11% 😈

clean training



Test Acc: 84.88%

our poisoning



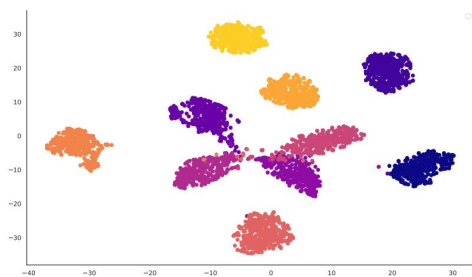
Test Acc: 72.99% 😈

Test Acc: 71.57% 😈

change perspective

Our Poisoning

clean training

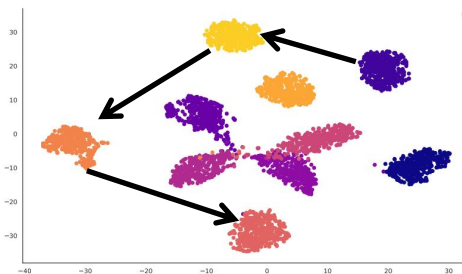


Test Acc: 84.88%

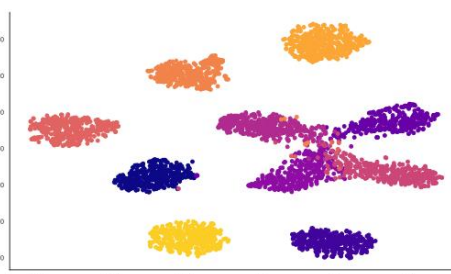
equal to discarding
83% training data!

change perspective

existing poisoning

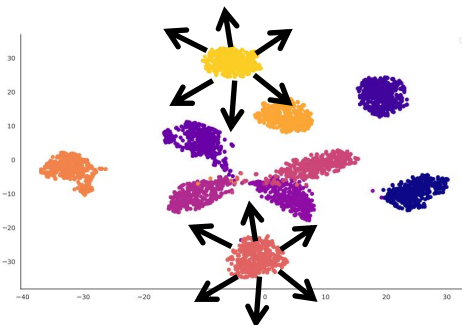


$$x' = \arg \min_x J(x, y_t)$$



Test Acc: 83.11% 😈

our poisoning

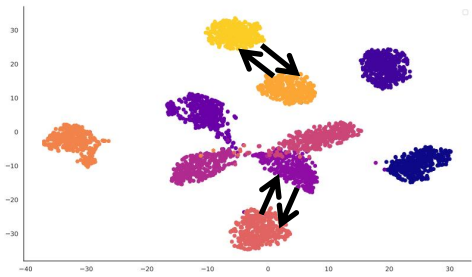


Test Acc: 72.99% 😈

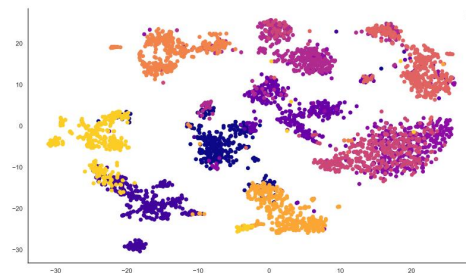
Test Acc: 71.57% 😈

Our Poisoning

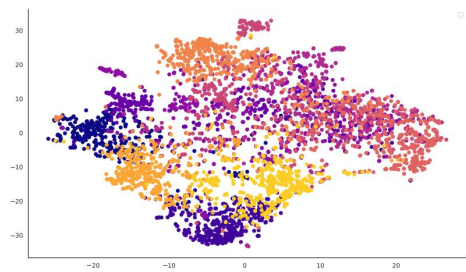
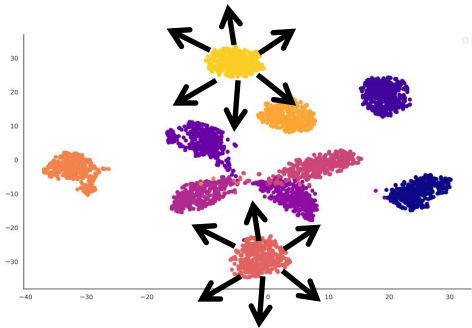
$$\mu = \frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} F_{L-1}^*(x)$$



our poisoning



Test Acc: 72.99% 😈



Test Acc: 71.57% 😈

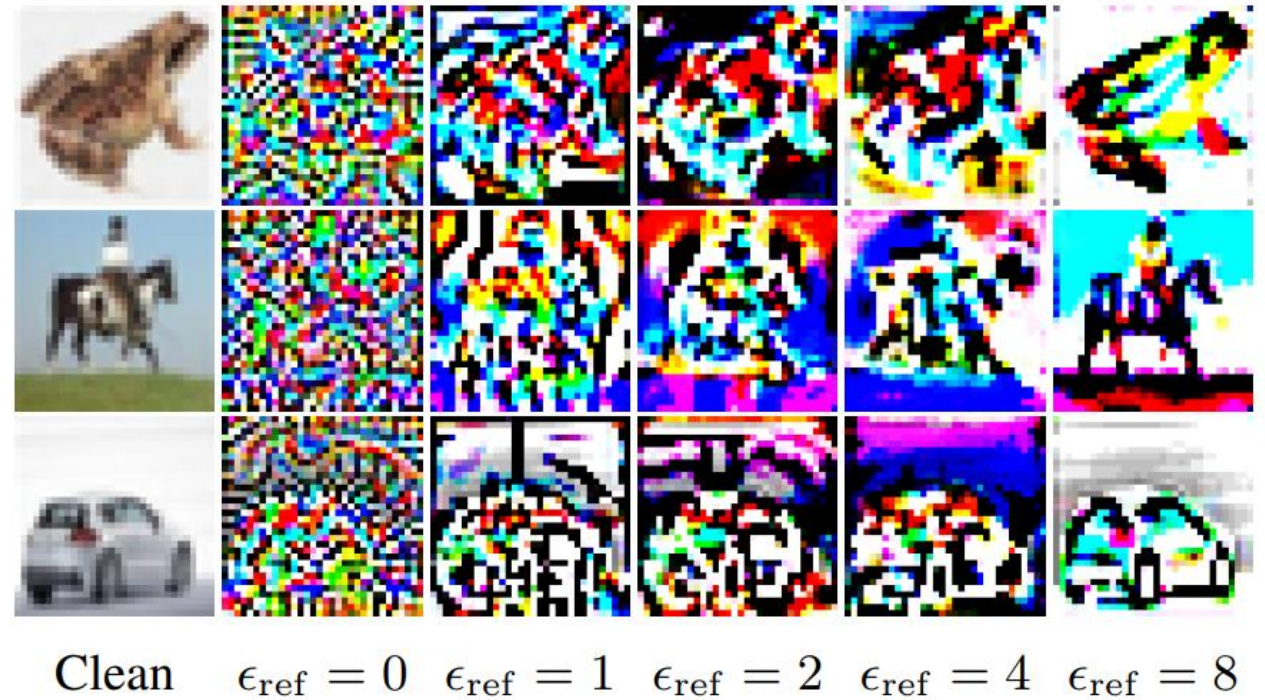
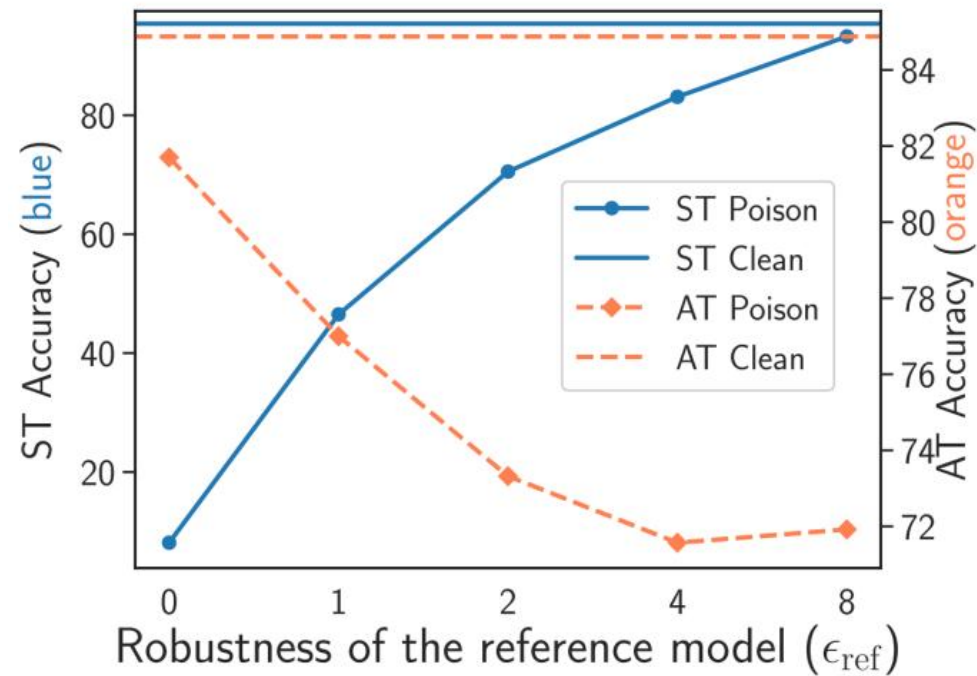
$$\mathcal{L}_{\text{pull}} = \min_{\delta^{\text{poi}}} \|F_{L-1}^*(x + \delta^{\text{poi}}) - \mu_{y'}\|_2$$

$$\mathcal{L}_{\text{push}} = \max_{\delta^{\text{poi}}} \|F_{L-1}^*(x + \delta^{\text{poi}}) - \mu_y\|_2$$

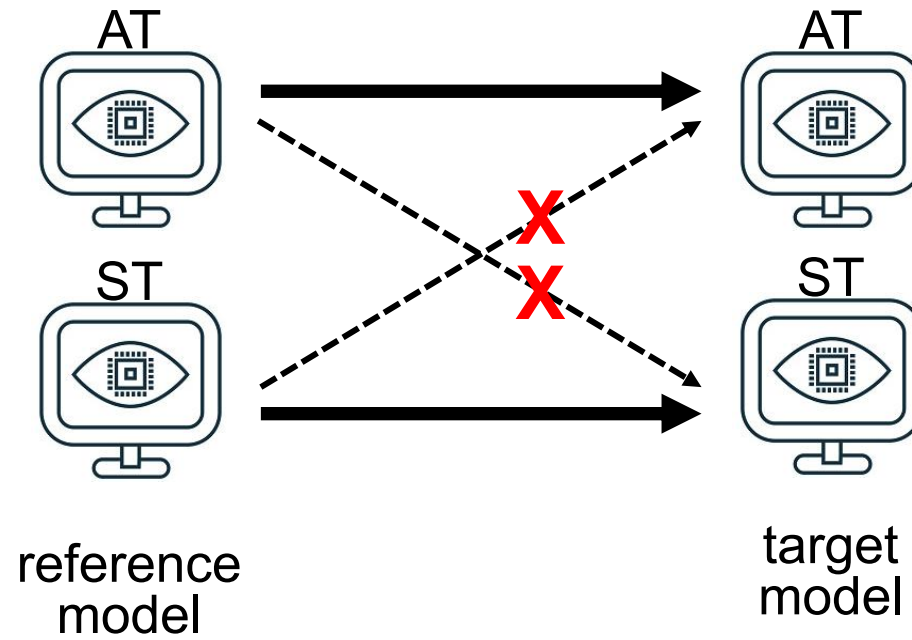
Results

- Different datasets
- Different AT frameworks
- Transferability
- Partial data Poisoning training data
- Ensemble defenses
- Adaptive defenses
- ...

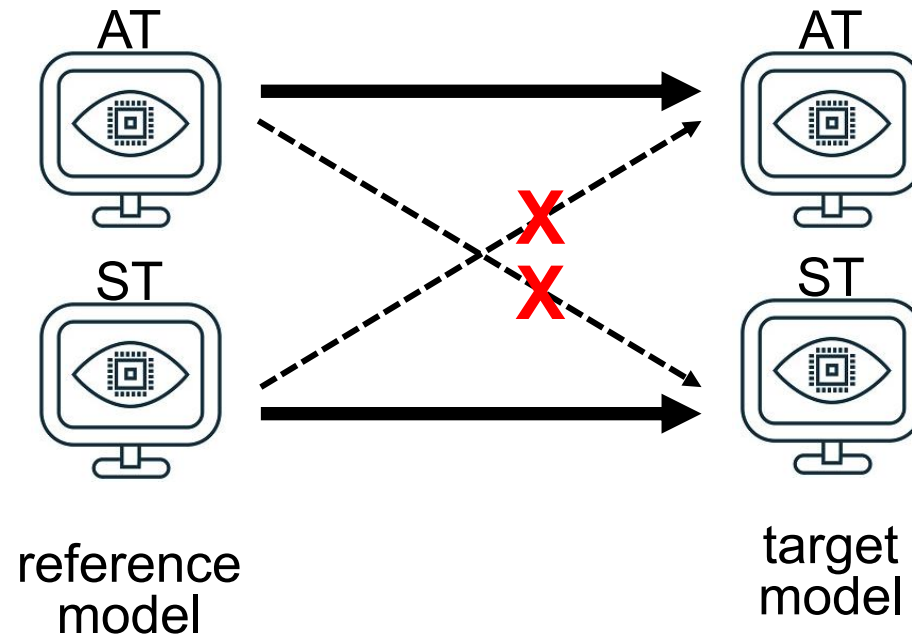
Standard Training (ST) vs. Adversarial Training (AT)



Hybrid Attack



Hybrid Attack



$$\mathcal{L}_{\text{push}} = \max_{\delta^{\text{poi}}} \|F_{L-1}^*(x + \delta^{\text{poi}}) - \mu_y\|_2$$

$$\mathcal{L}_{\text{hybrid}} = \max_{\delta^{\text{poi}}} \|F_{L-1, \text{ST}}^*(x + \delta^{\text{poi}}) - \mu_{y[\text{ST}]}\|_2 + \lambda \|F_{L-1, \text{AT}}^*(x + \delta^{\text{poi}}) - \mu_{y[\text{AT}]}\|_2$$

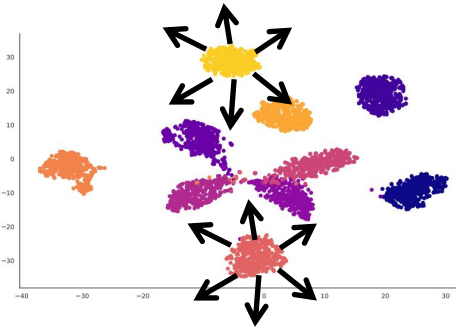
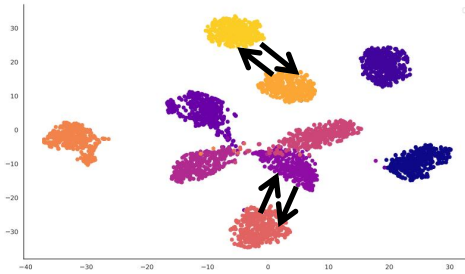
Hybrid Attack

METHOD ($\epsilon_{\text{poi}} = 8/255$) \ ϵ_{adv}	0/255	4/255	8/255	16/255	OPTIMAL TEST ACC.
NONE (CLEAN)	94.59	90.31	84.88	73.78	94.59
ADVPOISON	9.91	88.98	83.11	71.31	88.98
REM	25.59	46.57	84.21	85.76	85.76
ADVIN	77.31	90.08	86.76	72.16	90.08
UNLEARNABLE	25.69	90.47	84.91	79.81	90.47
HYPOCRITICAL	74.06	91.18	84.96	73.33	91.18
HYPOCRITICAL+	75.22	84.82	86.56	82.26	86.56
OURS	83.10	75.39	71.51	63.73	83.10
OURS (HYBRID)	12.93	76.55	74.30	65.75	76.55



Summary of Project 2

- Poisoning AT is possible based on a new attack strategy

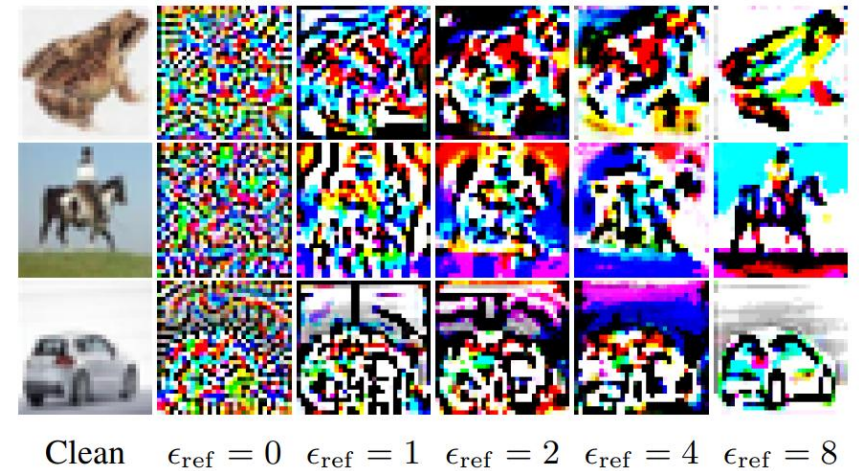


- Poisoning AT vs. ST
- Hybrid attack

Future Directions

- Possible defenses against our new attack
 - generic: training techniques for noisy labels?
 - specific: detecting/pre-filtering our attack?
- More efficient hybrid attack than

$$\mathcal{L}_{\text{hybrid}} = \max_{\delta^{\text{poi}}} \|F_{L-1, \text{ST}}^*(\mathbf{x} + \delta^{\text{poi}}) - \boldsymbol{\mu}_{y, \text{ST}}\|_2 + \lambda \|F_{L-1, \text{AT}}^*(\mathbf{x} + \delta^{\text{poi}}) - \boldsymbol{\mu}_{y, \text{AT}}\|_2$$



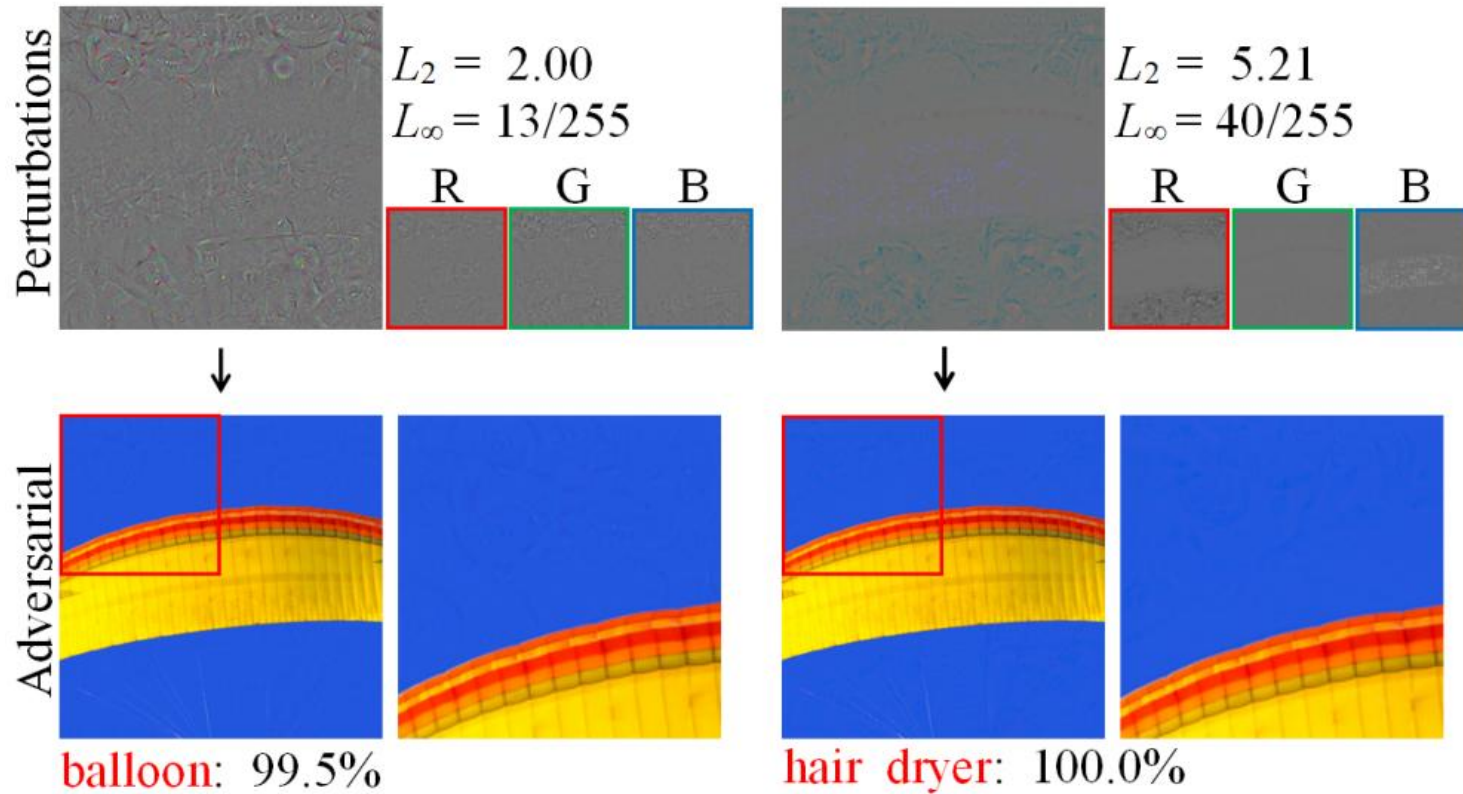
Outline

- Overview of adversarial images in computer vision
- Two recent projects
- **Other related projects**



Imperceptible Perturbations

$$\left\| \begin{array}{c} x' \\ \text{[Image of cat with perturbation]} \end{array} - \begin{array}{c} x_{\text{cat}} \\ \text{[Image of cat]} \end{array} \right\|_{\infty} \leq \epsilon \quad \rightarrow \quad \left\| \begin{array}{c} x' \\ \text{[Image of cat with perturbation]} \end{array} - \begin{array}{c} x_{\text{cat}} \\ \text{[Image of cat]} \end{array} \right\|_{\text{CIEDE2000}} \leq \epsilon$$



(a) C&W

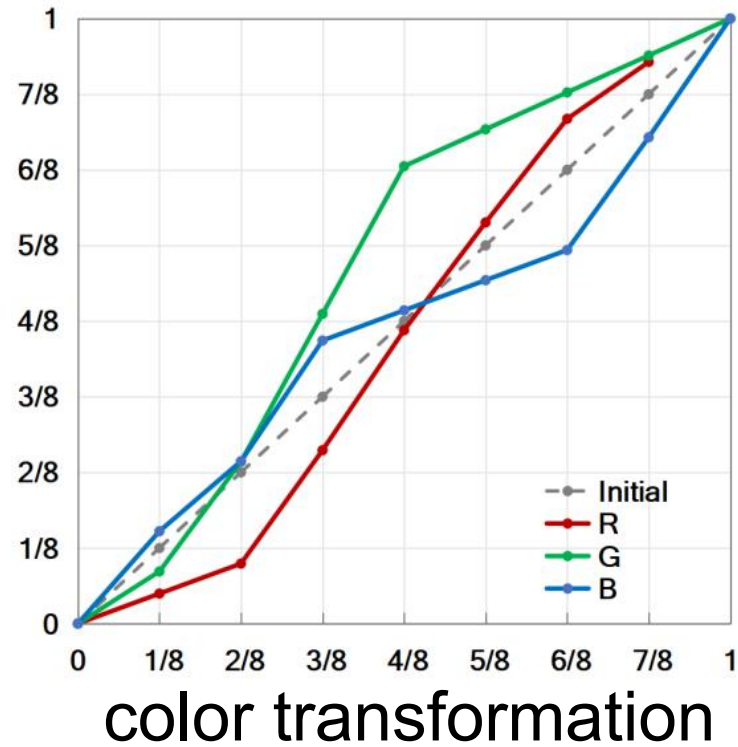
(b) PerC-C&W (ours)

Perceptible yet Stealthy Attacks

x



\times



$=$

x'



Adversarial attacks on Image Retrieval

Query

Top-3 ranking results

Adversary



- **On Success and Simplicity: A Second Look at Transferable Targeted Attacks (Project 1)**
Zhengyu Zhao, Zhuoran Liu, Martha Larson. NeurIPS 2021.
- **Is Adversarial Training Really a Silver Bullet for Mitigating Data Poisoning? (Project 2)**
Rui Wen, Zhengyu Zhao, Zhuoran Liu, Michael Backes, Tianhao Wang, Yang Zhang. ICLR 2023.
- **Towards Good Practices in Evaluating Transfer Adversarial Attacks**
Zhengyu Zhao*, Hanwei Zhang*, Renjue Li*, Ronan Sircé, Laurent Amsaleg, Michael Backes. arXiv 2022.
- **Towards Large yet Imperceptible Adversarial Image Perturbations with Perceptual Color Distance**
Zhengyu Zhao, Zhuoran Liu, Martha Larson. CVPR 2020.
- **Adversarial Image Color Transformations in Explicit Color Filter Space**
Zhengyu Zhao, Zhuoran Liu, Martha Larson. BMVC 2020.
- **Who's Afraid of Adversarial Queries? The Impact of Image Modifications on Content-based Image Retrieval**
Zhuoran Liu, Zhengyu Zhao, Martha Larson. ICMR 2019.

Thank you!

Q&A

