



**AJCAI 2024 Encore Track:**

**Trustworthy AI: Robustness and Ethics**

# **Improving Accuracy-robustness Trade-off via Pixel Reweighted Adversarial Training**

**Presenter: Jiacheng Zhang**

**School of Computing and Information Systems**

**The University of Melbourne**

**Date: 26 November 2024**



## What is an adversarial example (attack)?

☐ **Left-or-right challenge:** Guess which one is the adversarial example?





## What is an adversarial example (attack)?

**99% Guacamole**



**88% Tabby Cat**



## What is an adversarial example (attack)?

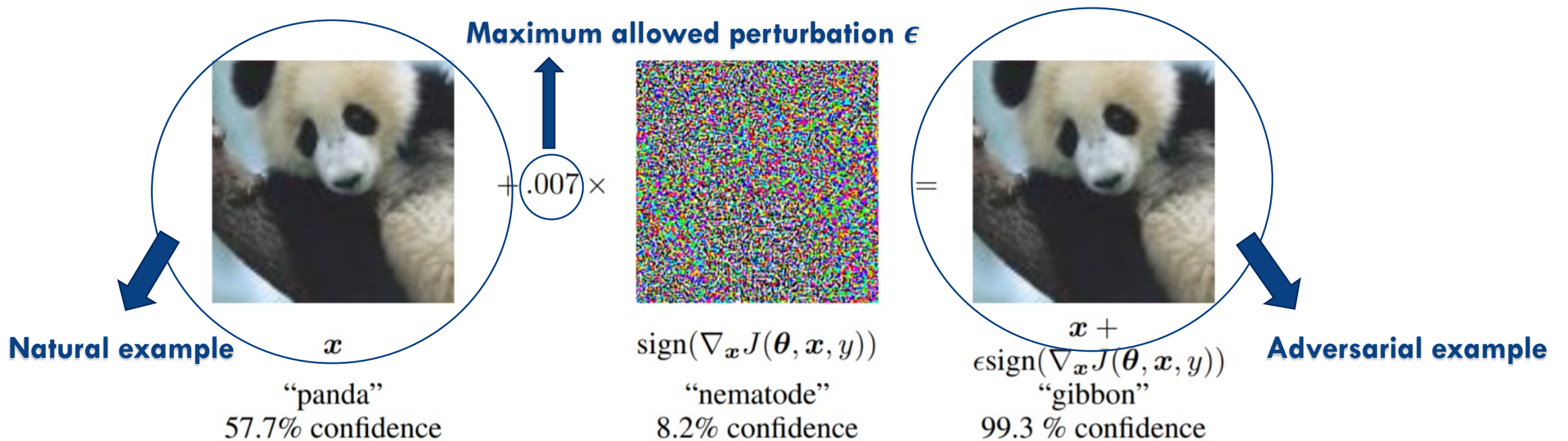
---

**Adversarial examples** can significantly drop the classification accuracy to **0%**.

**How it works?**

# What is an adversarial example (attack)?

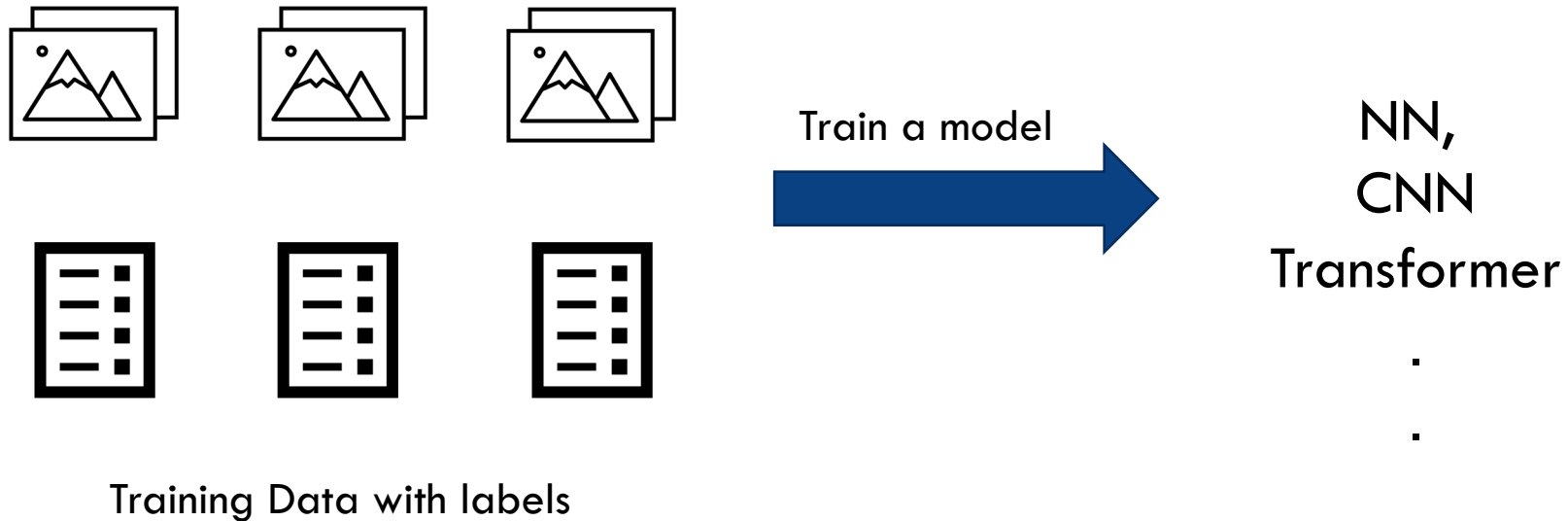
- Adding **imperceptible, non-random perturbations** to input data.



- Cannot fool human eyes, but **can easily fool** state-of-the-art neural networks.

# What is an adversarial example (attack)?

Conventional Machine Learning Pipeline (classification):



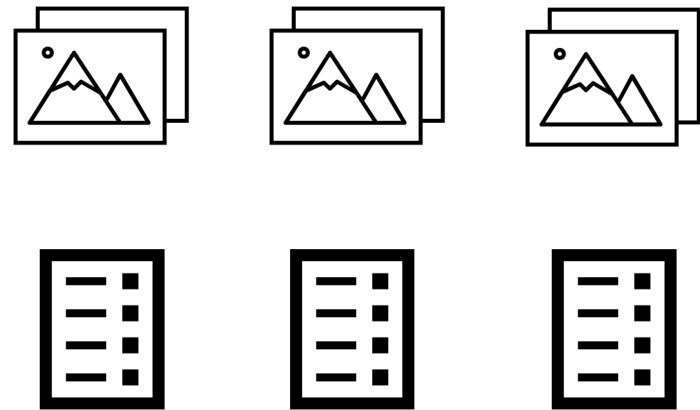
# What is an adversarial example (attack)?

Conventional Machine Learning Pipeline (classification):

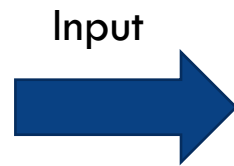


# What is an adversarial example (attack)?

Adversarial attack happens:



Test Data + **Adversarial Perturbations**



Well-trained NN,  
Well-trained CNN  
Well-trained Transformer



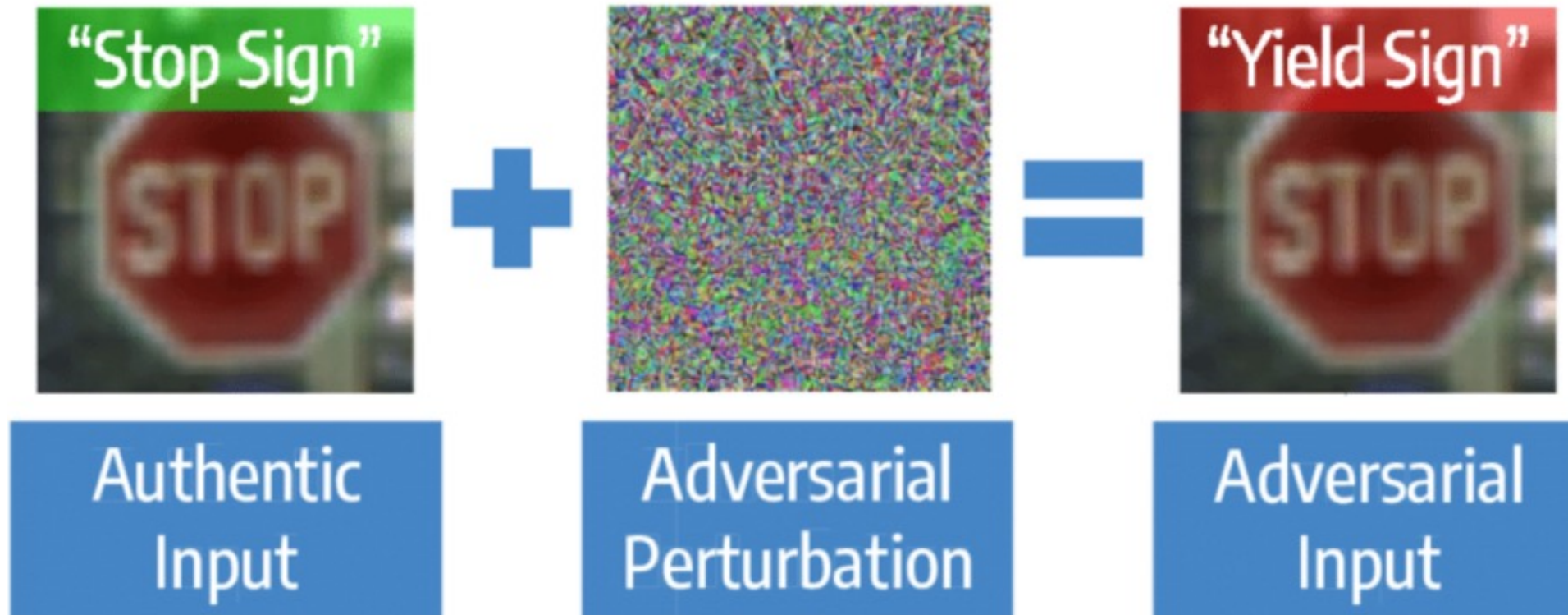
Predicted  
Labels

**Untrusted**

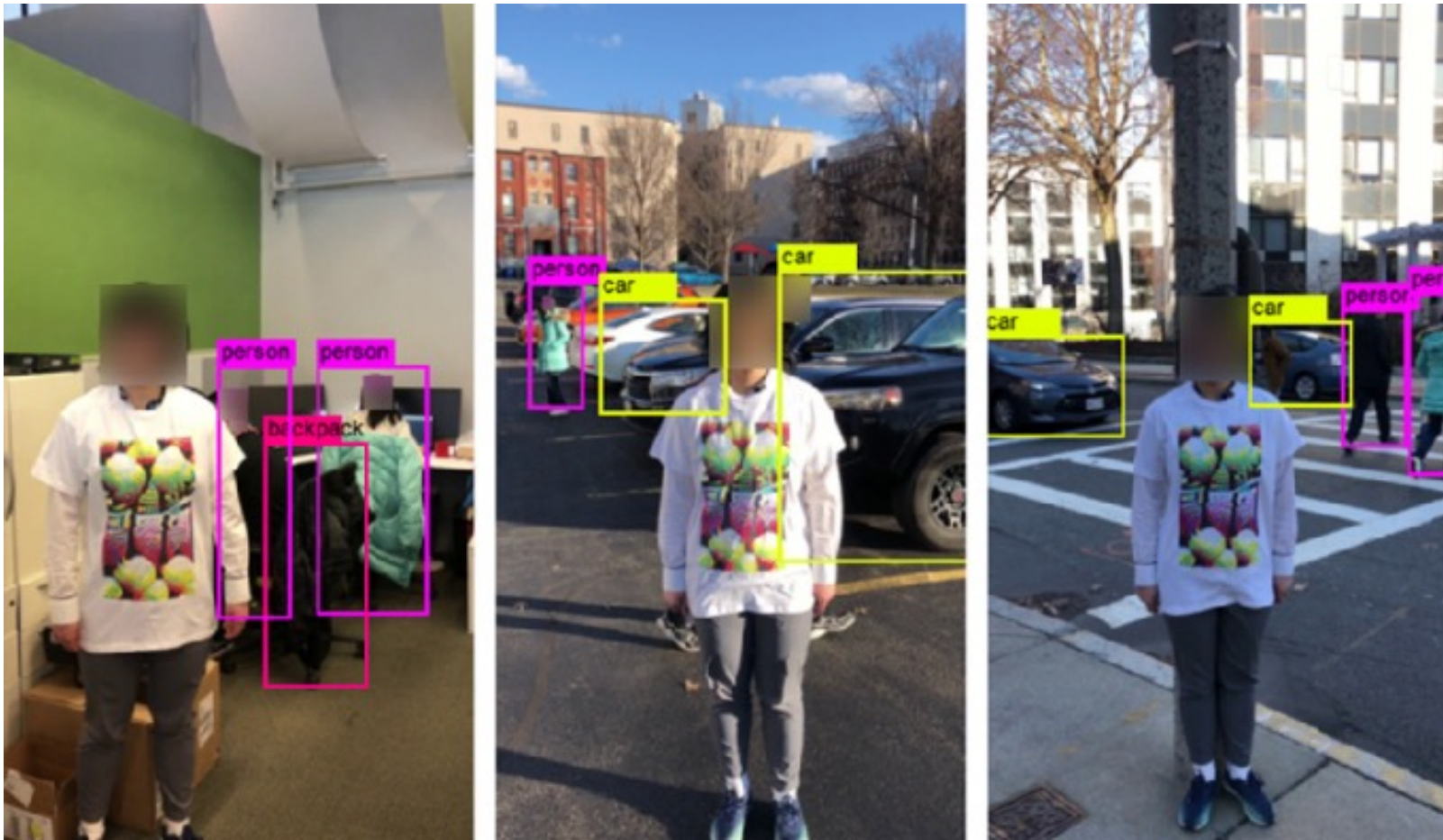


## Why do we care?

- ❑ Cause **security and reliability issues** in the deployment of machine learning systems.
- ❑ E.g., mislead the autonomous driving system to recognize **a stop sign** into **something else**.

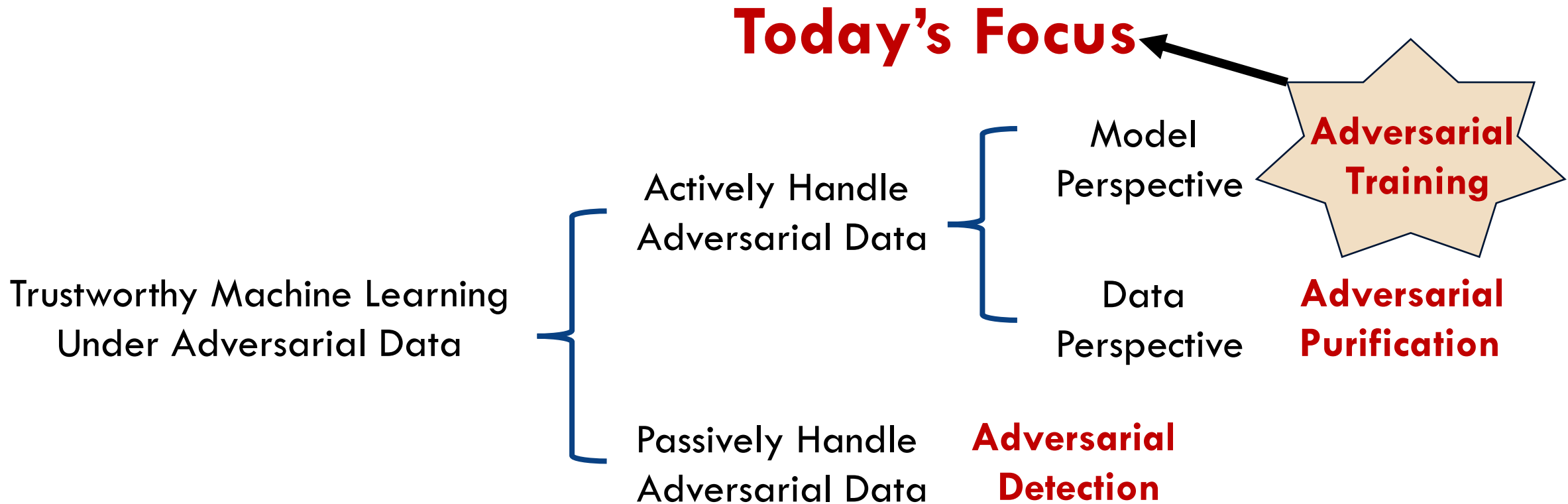


## Why do we care?



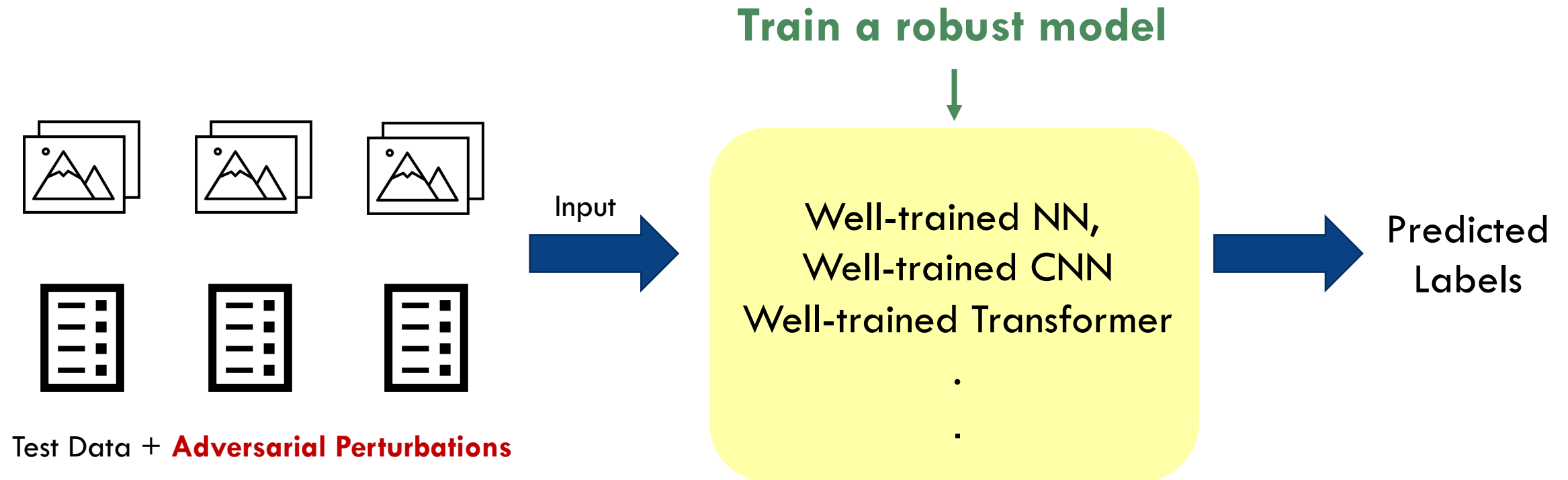
- ❑ Adding **adversarial examples** on T-shirts can bypass the AI detection system.
- ❑ Let you be invisible to the AI detection system!
- ❑ It's cool but it can cause **security and reliability issues.**

# How to defend against adversarial attacks?



# Adversarial training

- ❑ *Adversarial Training (AT)*: aims to train a robust model on *adversarial examples (AEs)*

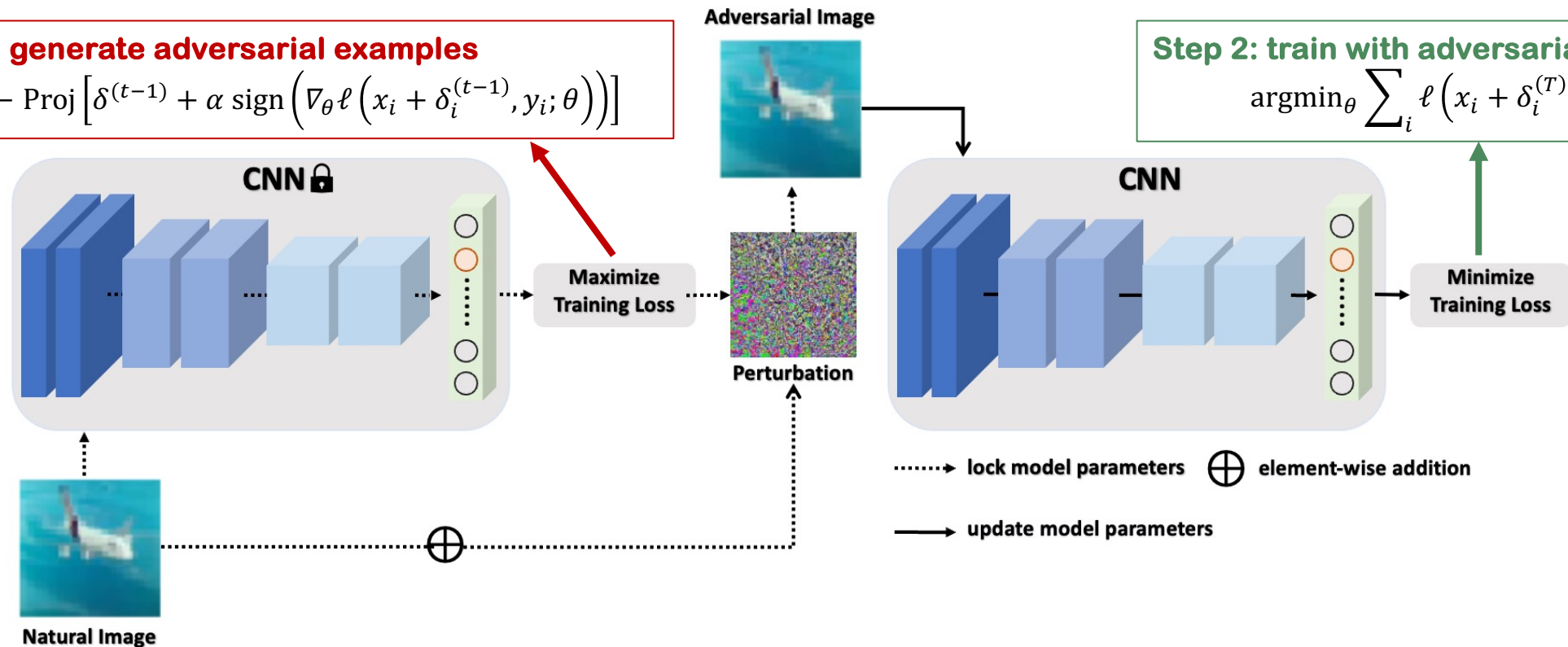


# Adversarial training

□ *Adversarial Training (AT)*: aims to train a robust model on *adversarial examples (AEs)*

## Step 1: generate adversarial examples

$$\delta^{(t)} \leftarrow \text{Proj} \left[ \delta^{(t-1)} + \alpha \text{sign} \left( \nabla_{\theta} \ell \left( x_i + \delta_i^{(t-1)}, y_i; \theta \right) \right) \right]$$





# What's the problem of adversarial training?

---

## What's good about AT?

- ❑ AT improves robustness (i.e., accuracy on adversarial examples).

## What's bad about AT?

- ❑ AT drops natural accuracy (i.e., accuracy on natural examples).

**Problem: there is an accuracy-robustness trade-off!!!**

- ❑ Higher robustness, lower accuracy; higher accuracy, lower robustness.

**What we really want:** a model that has high robustness without sacrificing natural accuracy.

# How to mitigate this problem?

What will happen if we change  $\epsilon$ ?

$$\begin{array}{ccc}
 \begin{array}{c} \text{Panda image} \\ x \\ \text{"panda"} \\ 57.7\% \text{ confidence} \end{array} & \begin{array}{c} \uparrow \\ + \textcircled{.007} \times \\ \text{Noise image} \\ \text{sign}(\nabla_x J(\theta, x, y)) \\ \text{"nematode"} \\ 8.2\% \text{ confidence} \end{array} & = \\
 & & \begin{array}{c} \text{Panda image} \\ x + \epsilon \text{sign}(\nabla_x J(\theta, x, y)) \\ \text{"gibbon"} \\ 99.3\% \text{ confidence} \end{array}
 \end{array}$$

- ❑ **Extreme case:** we decrease  $\epsilon$  to 0, AT will converge to natural training.
- ❑ **Conclusion:** decrease the budget can improve natural accuracy.

## Research gap and research question

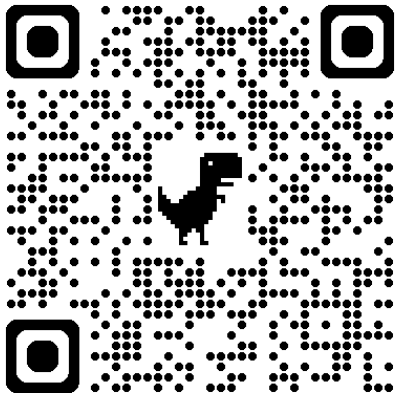
- ❑ **Research gap:** Existing AT methods apply a fixed  $\epsilon$  for all pixels in an image. Therefore, changing  $\epsilon$  must sacrifice natural accuracy or robustness.
- ❑ **Research question:** Can we design an adaptive method to **reweight  $\epsilon$  for only partial pixels in an image** so that we can increase natural accuracy without sacrificing robustness?
- ❑ In our recent work, we show that the answer to this question is **YES**.

# ICML 2024

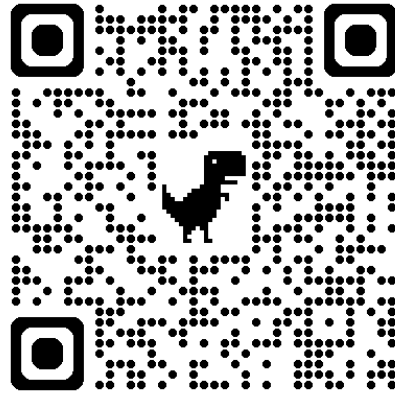
## Improving Accuracy-robustness Trade-off via Pixel Reweighted Adversarial Training

Jiacheng Zhang, Feng Liu\*, Dawei Zhou, Jingfeng Zhang, Tongliang Liu\*

Paper

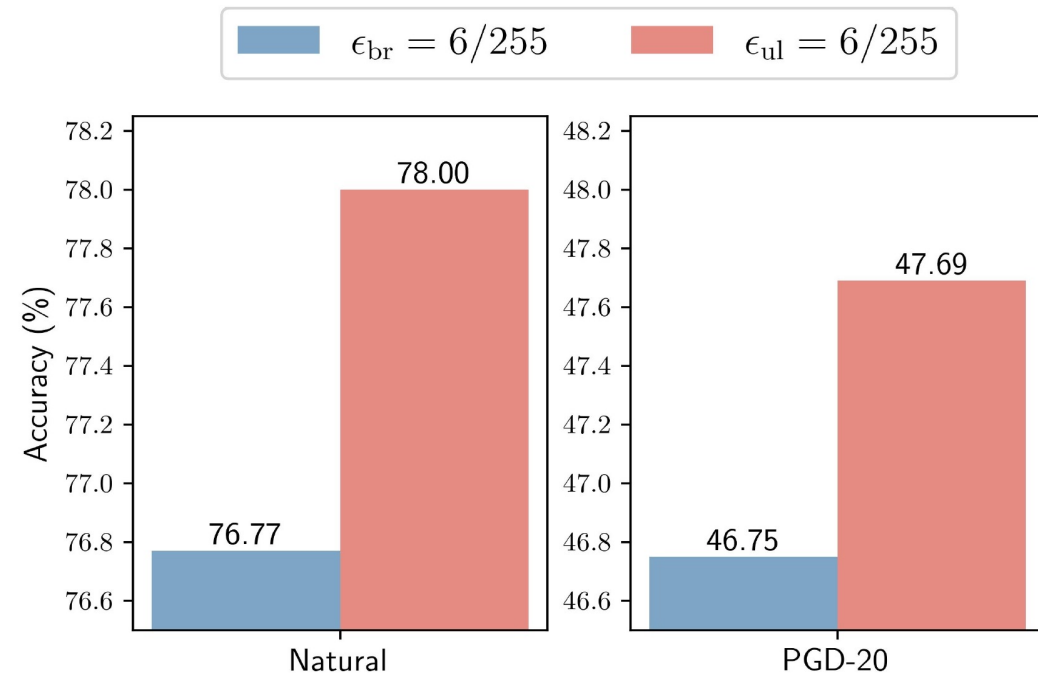
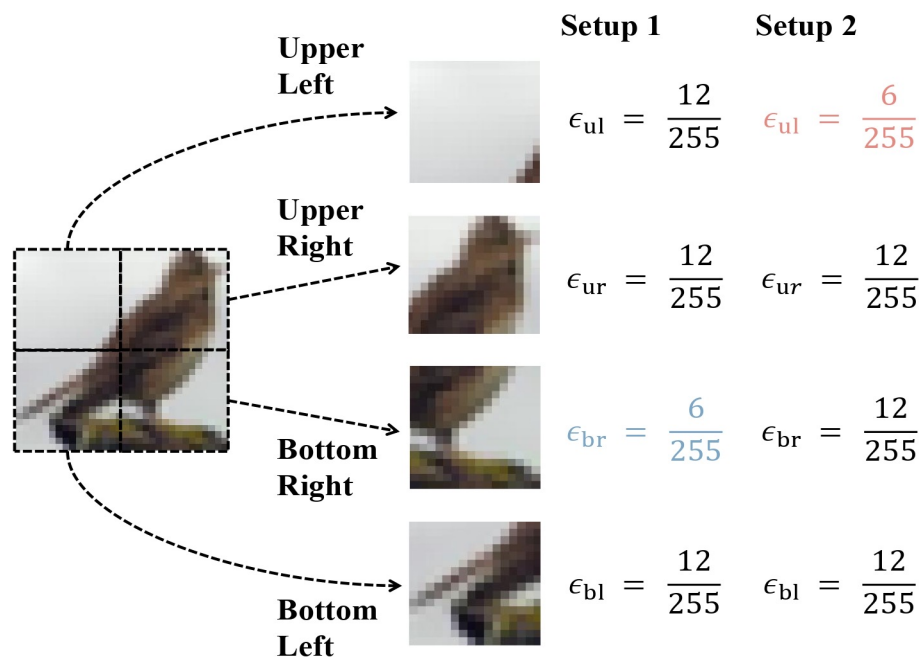


Code



# Are all pixels equally important in robust classification?

- ❑ **Proof-of-concept experiment:** changing the perturbation budgets for different parts of an image has the potential to boost robustness and accuracy **at the same time**.





# Pixel-reweighted Adversarial Training (PART)

- ❑ **Proof-of-concept experiment:** changing the perturbation budgets for different parts of an image has the potential to boost robustness and accuracy **at the same time**.



Pixels in an image have different importances in classification.



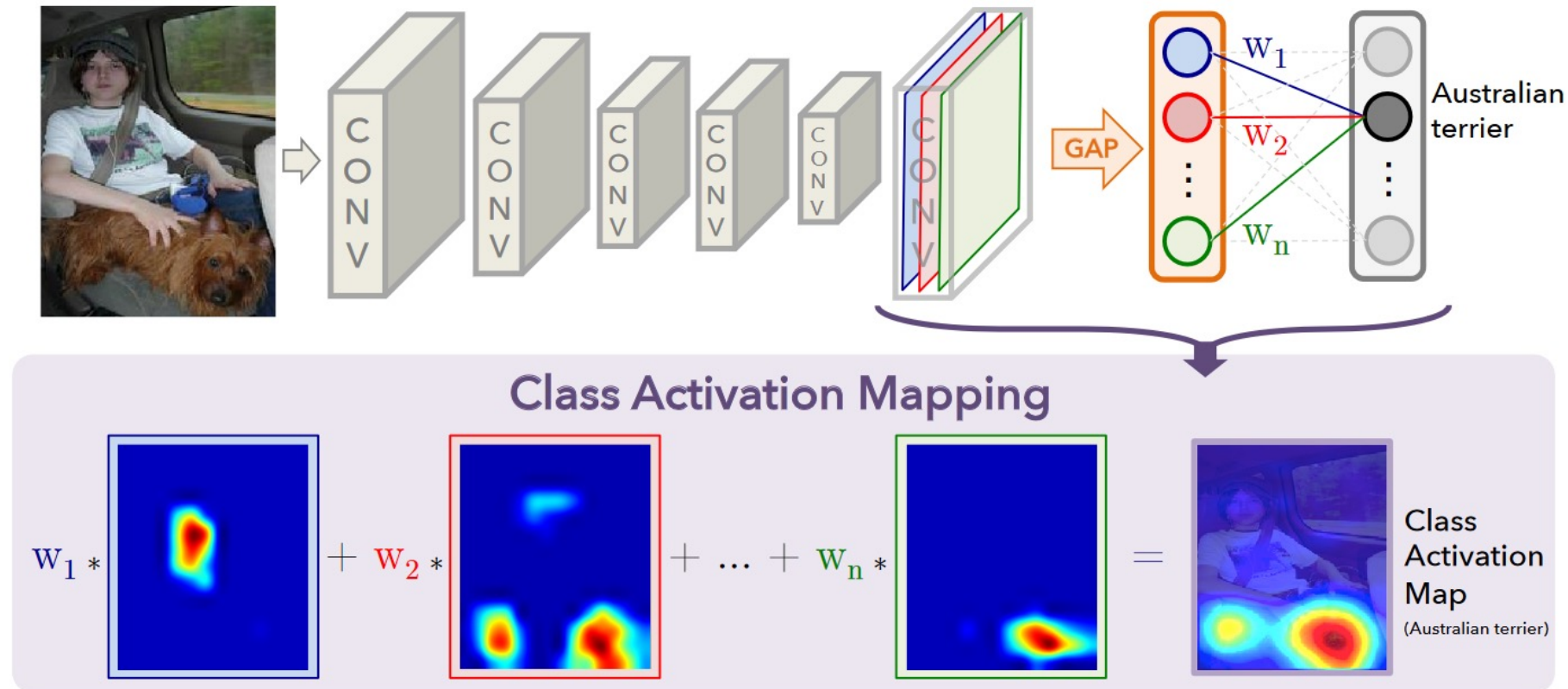
We need to guide the model to **focus more on important pixels** during training.



*Class activation mapping (CAM)* can help.

# Class activation mapping (CAM)

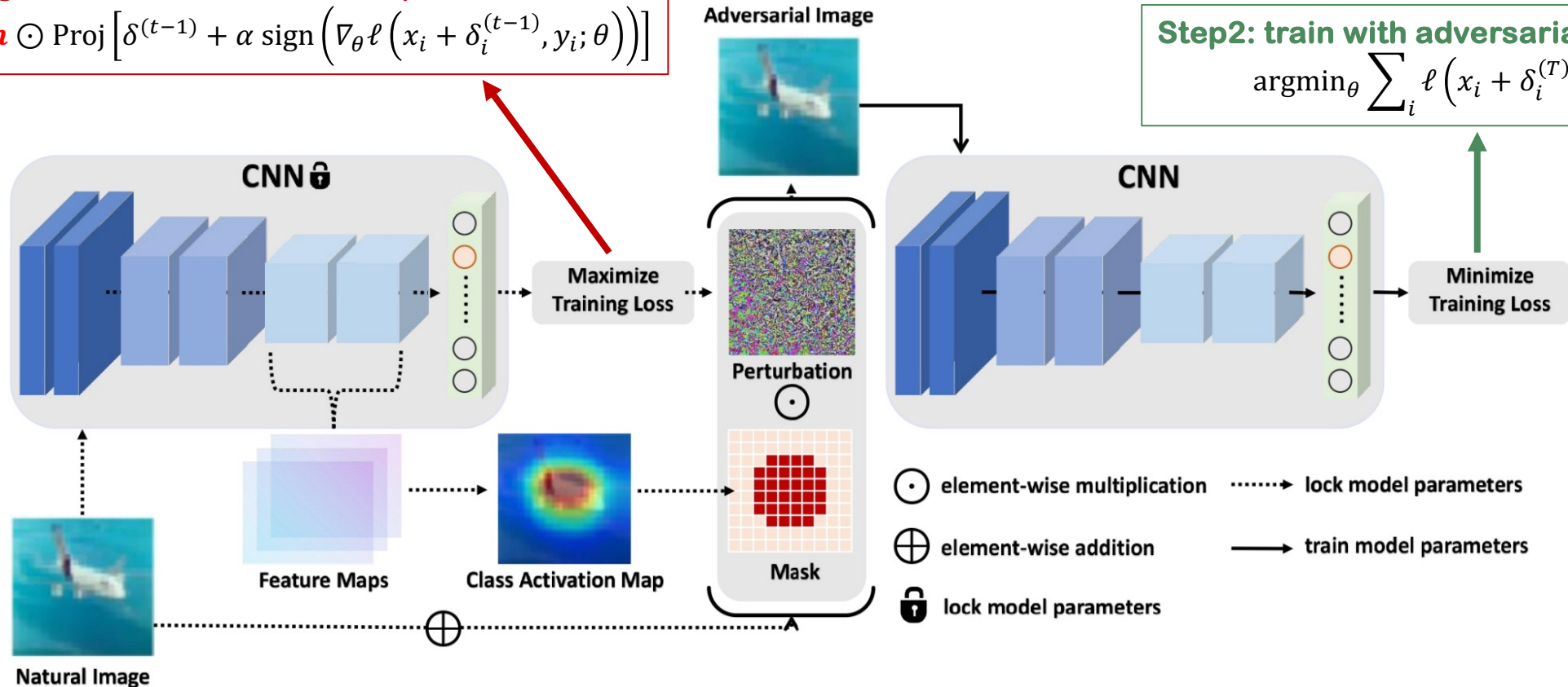
- ❑ **CAM:** the predicted class score is mapped back to the previous convolutional layer to generate the class activation maps. The CAM **highlights the regions that influence the classification result the most.**



# Pixel-reweighted Adversarial Training (PART)

## Step1: generate adversarial examples

$$\delta^{(t)} \leftarrow m \odot \text{Proj} \left[ \delta^{(t-1)} + \alpha \text{sign} \left( \nabla_{\theta} \ell \left( x_i + \delta_i^{(t-1)}, y_i; \theta \right) \right) \right]$$



## Step2: train with adversarial examples

$$\text{argmin}_{\theta} \sum_i \ell \left( x_i + \delta_i^{(T)}, y_i; \theta \right)$$

## Pixel-reweighted Adversarial Training (PART)

**Q: Why can PART improve natural accuracy?**

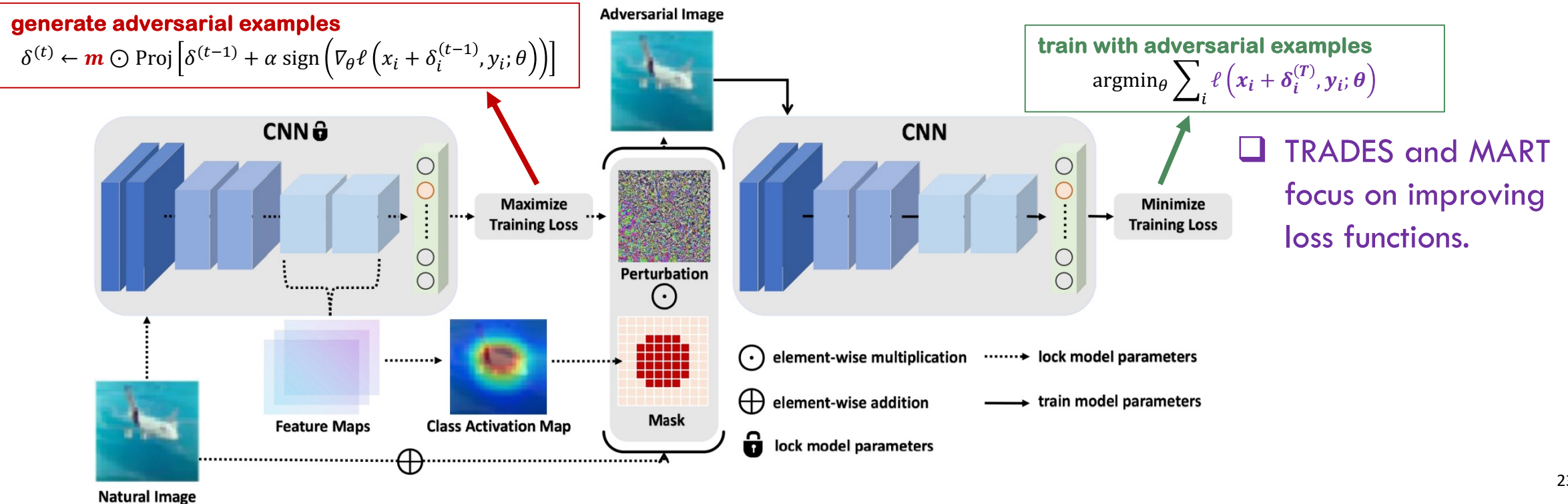
A: Because PART partially reduces  $\epsilon$ .

**Q: Why can PART maintain robust accuracy?**

A: Because PART explicitly guides the model to focus more on important pixel regions during training, making sure these regions stay robust to adversarial attacks.

# Strength of PART

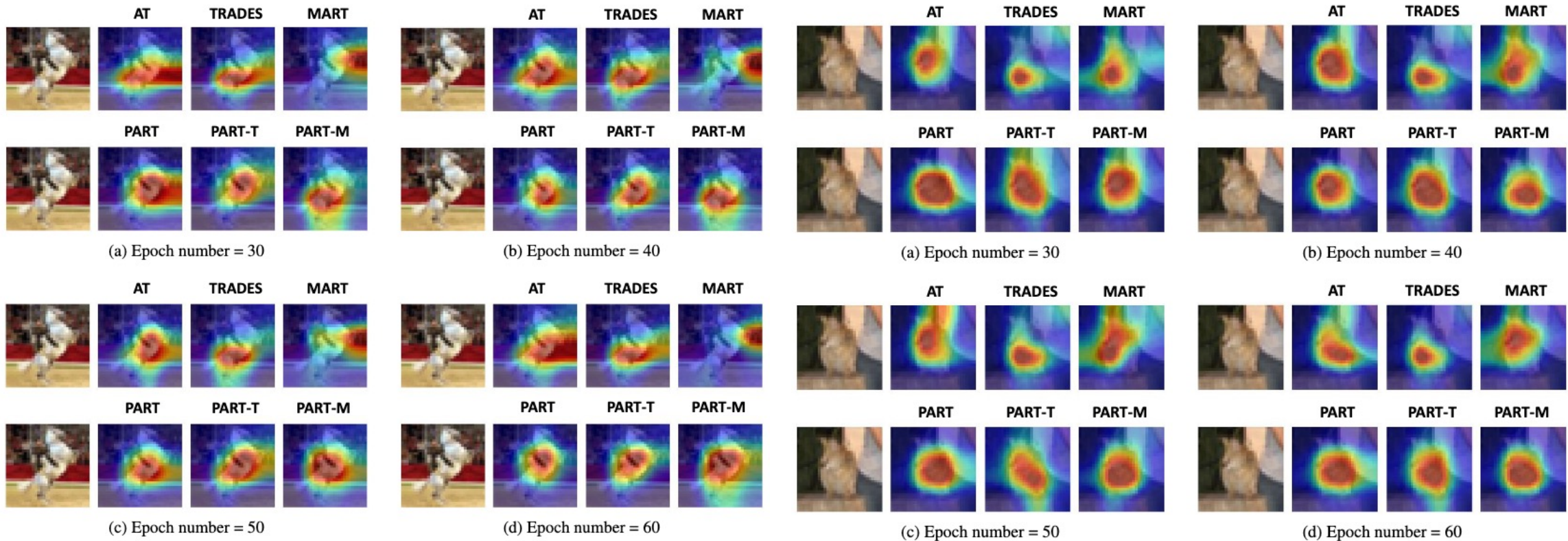
- ❑ PART is **orthogonal** to many SOTA AT methods (e.g., TRADES, MART).
- ❑ PART can be easily integrated into those methods.





# Strength of PART

□ PART-based methods align better with semantic information.



# Main results

Dataset	Method	Natural	PGD-20	MMA	AA
ResNet-18					
CIFAR-10	AT	82.58 $\pm$ 0.14	<b>43.69 <math>\pm</math> 0.28</b>	41.80 $\pm$ 0.10	41.63 $\pm$ 0.22
	PART ( $s = 1$ )	83.42 $\pm$ 0.26 (+ 0.84)	43.65 $\pm$ 0.16 (- 0.04)	<b>41.98 <math>\pm</math> 0.03 (+ 0.18)</b>	<b>41.74 <math>\pm</math> 0.04 (+ 0.11)</b>
	PART ( $s = 10$ )	<b>83.77 <math>\pm</math> 0.15 (+ 1.19)</b>	43.36 $\pm$ 0.21 (- 0.33)	41.83 $\pm$ 0.07 (+ 0.03)	41.41 $\pm$ 0.14 (- 0.22)
	TRADES	78.16 $\pm$ 0.15	48.28 $\pm$ 0.05	45.00 $\pm$ 0.08	45.05 $\pm$ 0.12
	PART-T ( $s = 1$ )	79.36 $\pm$ 0.31 (+ 1.20)	<b>48.90 <math>\pm</math> 0.14 (+ 0.62)</b>	<b>45.90 <math>\pm</math> 0.07 (+ 0.90)</b>	<b>45.97 <math>\pm</math> 0.06 (+ 0.92)</b>
	PART-T ( $s = 10$ )	<b>80.13 <math>\pm</math> 0.16 (+ 1.97)</b>	48.72 $\pm$ 0.11 (+ 0.44)	45.59 $\pm$ 0.09 (+ 0.59)	45.60 $\pm$ 0.04 (+ 0.55)
	MART	76.82 $\pm$ 0.28	49.86 $\pm$ 0.32	45.42 $\pm$ 0.04	45.10 $\pm$ 0.06
	PART-M ( $s = 1$ )	78.67 $\pm$ 0.10 (+ 1.85)	<b>50.26 <math>\pm</math> 0.17 (+ 0.40)</b>	<b>45.53 <math>\pm</math> 0.05 (+ 0.11)</b>	<b>45.19 <math>\pm</math> 0.04 (+ 0.09)</b>
	PART-M ( $s = 10$ )	<b>80.00 <math>\pm</math> 0.15 (+ 3.18)</b>	49.71 $\pm$ 0.12 (- 0.15)	45.14 $\pm$ 0.10 (- 0.28)	44.61 $\pm$ 0.24 (- 0.49)
ResNet-18					
SVHN	AT	91.06 $\pm$ 0.24	49.83 $\pm$ 0.13	47.68 $\pm$ 0.06	45.48 $\pm$ 0.05
	PART ( $s = 1$ )	93.14 $\pm$ 0.05 (+ 2.08)	<b>50.34 <math>\pm</math> 0.14 (+ 0.51)</b>	<b>48.08 <math>\pm</math> 0.09 (+ 0.40)</b>	<b>45.67 <math>\pm</math> 0.13 (+ 0.19)</b>
	PART ( $s = 10$ )	<b>93.75 <math>\pm</math> 0.07 (+ 2.69)</b>	50.21 $\pm$ 0.10 (+ 0.38)	48.00 $\pm$ 0.14 (+ 0.32)	45.61 $\pm$ 0.08 (+ 0.13)
	TRADES	88.91 $\pm$ 0.28	58.74 $\pm$ 0.53	53.29 $\pm$ 0.56	52.21 $\pm$ 0.47
	PART-T ( $s = 1$ )	91.35 $\pm$ 0.11 (+ 2.44)	<b>59.33 <math>\pm</math> 0.22 (+ 0.59)</b>	<b>54.04 <math>\pm</math> 0.16 (+ 0.75)</b>	<b>53.07 <math>\pm</math> 0.67 (+ 0.86)</b>
	PART-T ( $s = 10$ )	<b>91.94 <math>\pm</math> 0.18 (+ 3.03)</b>	59.01 $\pm$ 0.13 (+ 0.27)	53.80 $\pm$ 0.20 (+ 0.51)	52.61 $\pm$ 0.24 (+ 0.40)
	MART	89.76 $\pm$ 0.08	58.52 $\pm$ 0.53	52.42 $\pm$ 0.34	49.10 $\pm$ 0.23
	PART-M ( $s = 1$ )	91.42 $\pm$ 0.36 (+ 1.66)	<b>58.85 <math>\pm</math> 0.29 (+ 0.33)</b>	<b>52.45 <math>\pm</math> 0.03 (+ 0.03)</b>	<b>49.92 <math>\pm</math> 0.10 (+ 0.82)</b>
	PART-M ( $s = 10$ )	<b>93.20 <math>\pm</math> 0.22 (+ 3.44)</b>	58.41 $\pm$ 0.20 (- 0.11)	52.18 $\pm$ 0.14 (- 0.24)	49.25 $\pm$ 0.13 (+ 0.15)
WideResNet-34-10					
TinyImagenet-200	AT	43.51 $\pm$ 0.13	11.70 $\pm$ 0.08	10.66 $\pm$ 0.11	10.53 $\pm$ 0.14
	PART ( $s = 1$ )	44.87 $\pm$ 0.21 (+ 1.36)	<b>11.93 <math>\pm</math> 0.16 (+ 0.23)</b>	<b>10.96 <math>\pm</math> 0.12 (+ 0.30)</b>	<b>10.76 <math>\pm</math> 0.06 (+ 0.23)</b>
	PART ( $s = 10$ )	<b>45.59 <math>\pm</math> 0.14 (+ 2.08)</b>	11.81 $\pm$ 0.10 (+ 0.11)	10.91 $\pm$ 0.08 (+ 0.25)	10.68 $\pm$ 0.10 (+ 0.15)
	TRADES	43.05 $\pm$ 0.15	13.86 $\pm$ 0.10	12.62 $\pm$ 0.16	12.55 $\pm$ 0.09
	PART-T ( $s = 1$ )	44.31 $\pm$ 0.12 (+ 1.26)	<b>14.08 <math>\pm</math> 0.22 (+ 0.22)</b>	<b>13.01 <math>\pm</math> 0.09 (+ 0.39)</b>	<b>12.84 <math>\pm</math> 0.14 (+ 0.29)</b>
	PART-T ( $s = 10$ )	<b>45.16 <math>\pm</math> 0.10 (+ 2.11)</b>	13.98 $\pm$ 0.15 (+ 0.12)	12.88 $\pm$ 0.12 (+ 0.26)	12.72 $\pm$ 0.08 (+ 0.17)
	MART	42.68 $\pm$ 0.22	14.77 $\pm$ 0.18	13.58 $\pm$ 0.13	13.42 $\pm$ 0.16
	PART-M ( $s = 1$ )	43.75 $\pm$ 0.24 (+ 1.07)	<b>14.93 <math>\pm</math> 0.15 (+ 0.16)</b>	<b>13.76 <math>\pm</math> 0.06 (+ 0.18)</b>	<b>13.68 <math>\pm</math> 0.13 (+ 0.24)</b>
	PART-M ( $s = 10$ )	<b>45.02 <math>\pm</math> 0.16 (+ 2.34)</b>	14.65 $\pm$ 0.14 (- 0.12)	13.41 $\pm$ 0.11 (- 0.17)	13.37 $\pm$ 0.15 (- 0.05)

- Our method can achieve a notable improvement in accuracy-robustness trade-off.
- In most cases, our method can improve natural accuracy by a notable margin without sacrificing robustness.

# Main results

ResNet-18						
Corruption	AT	PART ( $s = 10$ )	TRADES	PART-T ( $s = 10$ )	MART	PART-M ( $s = 10$ )
Gaussian Noise	81.05	<b>82.46</b>	76.05	<b>79.42</b>	77.57	<b>79.51</b>
Shot Noise	81.11	<b>82.83</b>	75.91	<b>79.70</b>	77.66	<b>79.74</b>
Impulse Noise	79.42	<b>80.76</b>	74.59	<b>78.03</b>	76.12	<b>78.16</b>
Speckle Noise	81.42	<b>82.68</b>	75.97	<b>79.68</b>	77.63	<b>79.79</b>
Defocus Blur	81.07	<b>82.48</b>	76.06	<b>79.42</b>	77.59	<b>79.50</b>
Glass Blur	77.82	<b>78.20</b>	72.60	<b>76.26</b>	74.17	<b>76.57</b>
Motion Blur	79.30	<b>80.13</b>	74.28	<b>77.43</b>	75.35	<b>77.45</b>
Zoom Blur	78.87	<b>79.30</b>	73.27	<b>76.74</b>	74.10	<b>76.74</b>
Gaussian Blur	81.05	<b>82.46</b>	76.05	<b>79.42</b>	77.57	<b>79.51</b>
Snow	81.09	<b>82.01</b>	76.13	<b>79.43</b>	77.63	<b>79.34</b>
Frost	78.96	<b>80.04</b>	73.90	<b>76.60</b>	75.01	<b>75.80</b>
Fog	79.34	<b>80.18</b>	72.95	<b>77.15</b>	75.29	<b>78.26</b>
Brightness	81.89	<b>83.22</b>	76.87	<b>80.16</b>	78.60	<b>80.26</b>
Spatter	81.03	<b>82.03</b>	75.81	<b>79.31</b>	77.62	<b>79.39</b>
Contrast	77.09	<b>77.67</b>	70.08	<b>75.00</b>	71.90	<b>75.81</b>
Elastic Transform	77.32	<b>78.16</b>	71.99	<b>75.39</b>	73.13	<b>75.68</b>
Pixelate	81.09	<b>82.63</b>	76.05	<b>79.48</b>	77.45	<b>79.43</b>
JPEG Compression	80.50	<b>81.91</b>	75.71	<b>79.09</b>	77.15	<b>79.29</b>
Saturate	78.01	<b>79.34</b>	73.59	<b>76.34</b>	74.70	<b>75.52</b>

- ❑ On corrupted dataset (e.g., CIFAR-10-C), our method can outperform baseline methods by a notable margin.
- ❑ The main reason is that the model can still figure out important parts of an image even if an image is corrupted.



## Conclusion and future work

- ❑ **Key message we want to deliver in this paper:** Guiding the model to focus more on essential pixel regions during training can help improve the generalizability (e.g., accuracy-robustness trade-off) of vision models.

# Questions?



THE UNIVERSITY OF  
MELBOURNE

# Thank you