



THE UNIVERSITY OF
MELBOURNE

Leveraging Distributional Discrepancies For Accuracy-robustness Trade-off

Jiacheng Zhang

School of Computing and Information Systems

The University of Melbourne

25 June 2025



Outline

- ❑ Background
- ❑ **ICML 2025:** Sample-specific Noise Injection for Diffusion-based Adversarial Purification
- ❑ **ICML 2025:** One Stone, Two Birds: Enhancing Adversarial Defense Through the Lens of Distributional Discrepancy

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



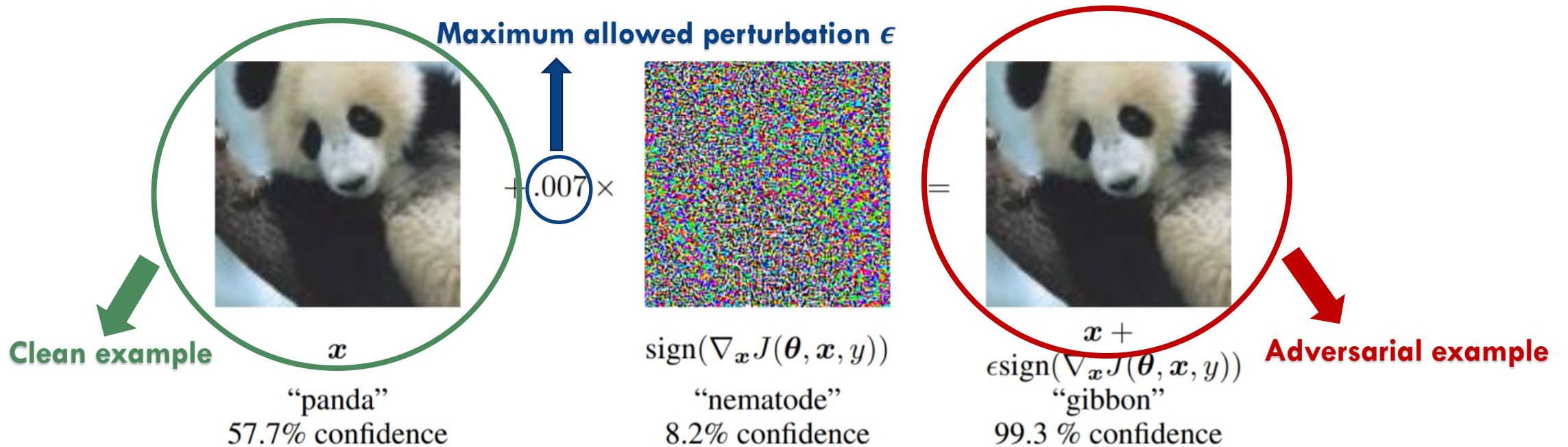
What is an adversarial example (attack)?

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

How it works?

What is an adversarial example (attack)?

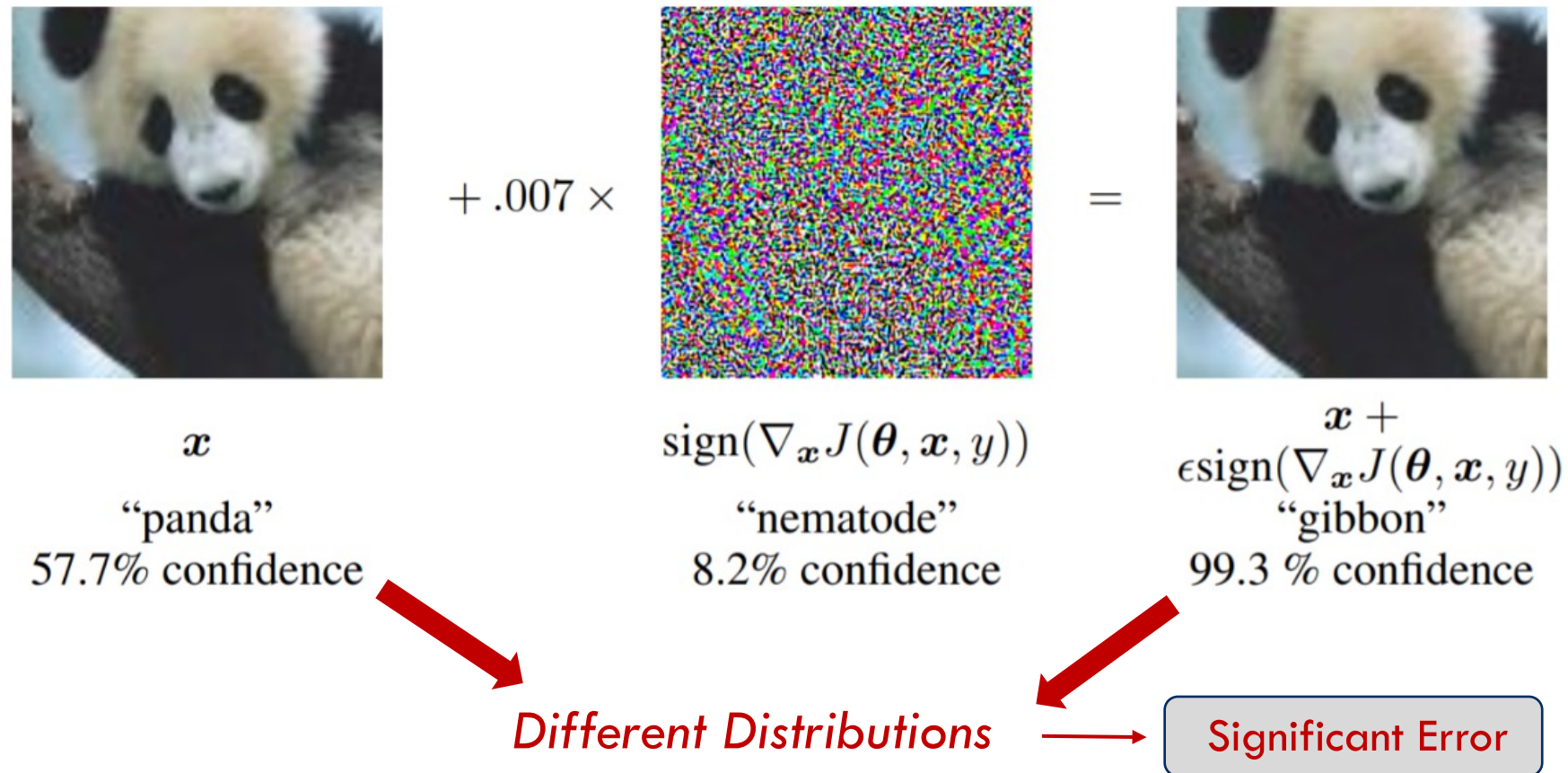
- Adding **imperceptible, non-random perturbations** to input data.



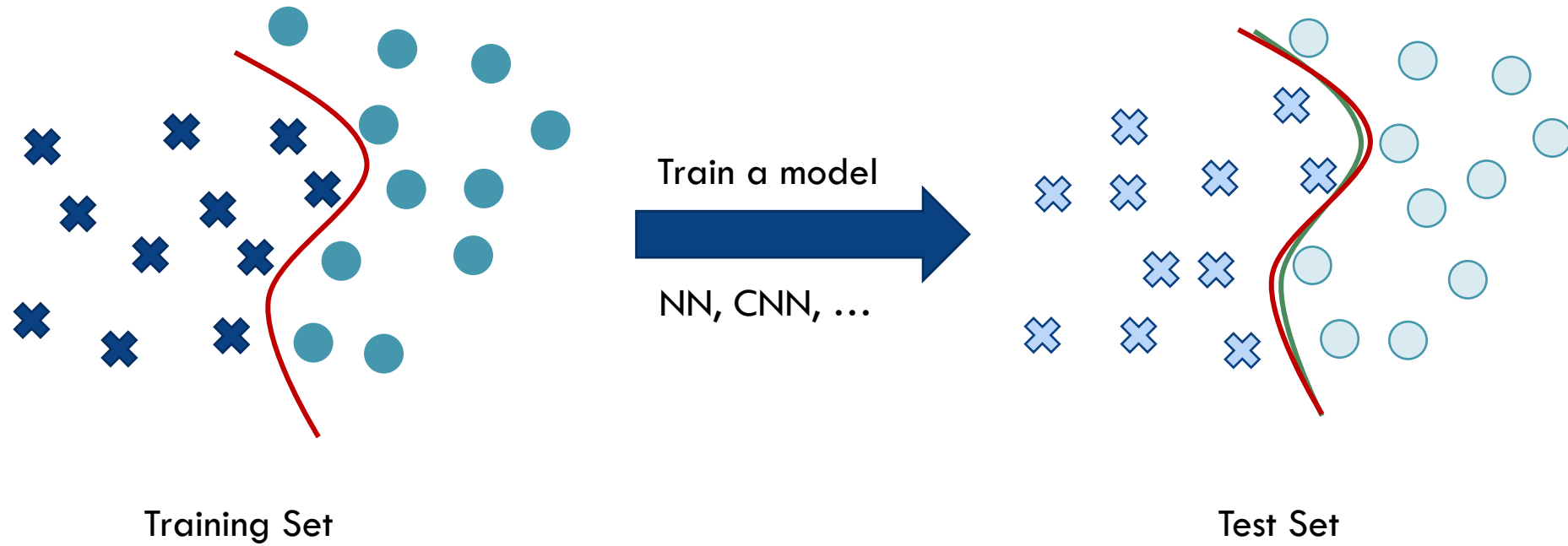
- Cannot fool human eyes but **can easily fool** state-of-the-art neural networks.

Why it works?

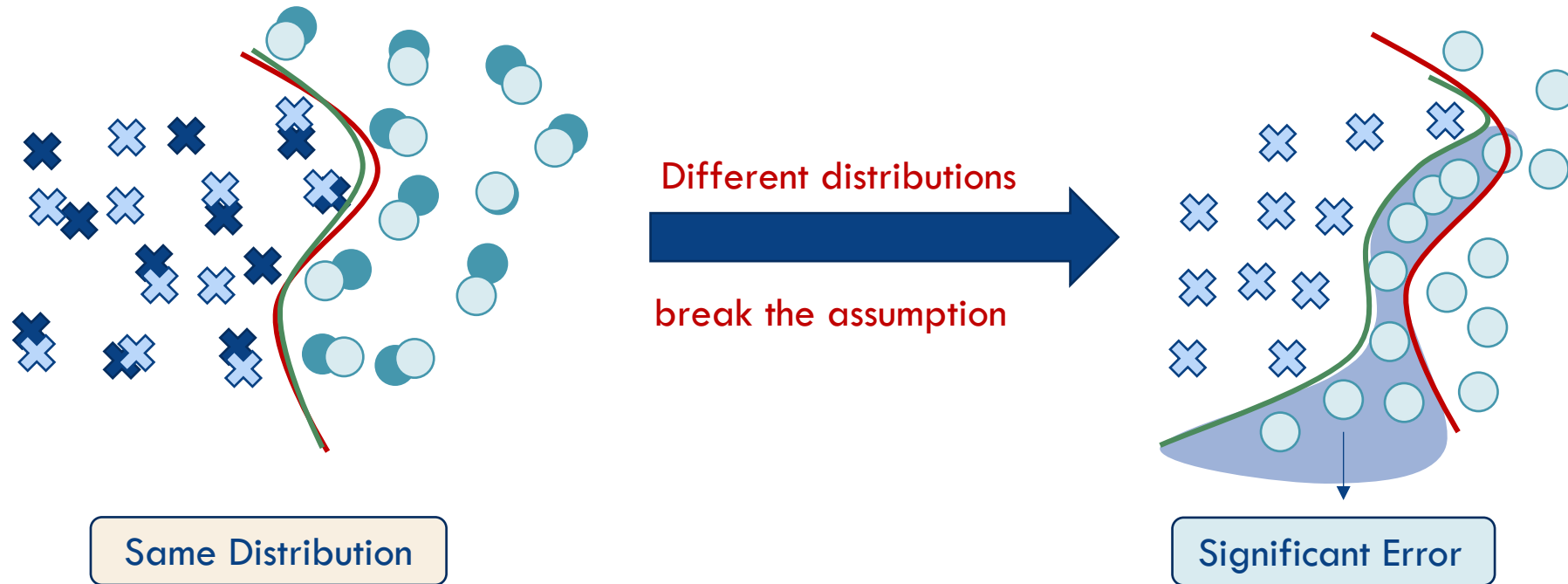
Why adversarial attack can be successful?



Basic assumption in machine learning



Basic assumption in machine learning

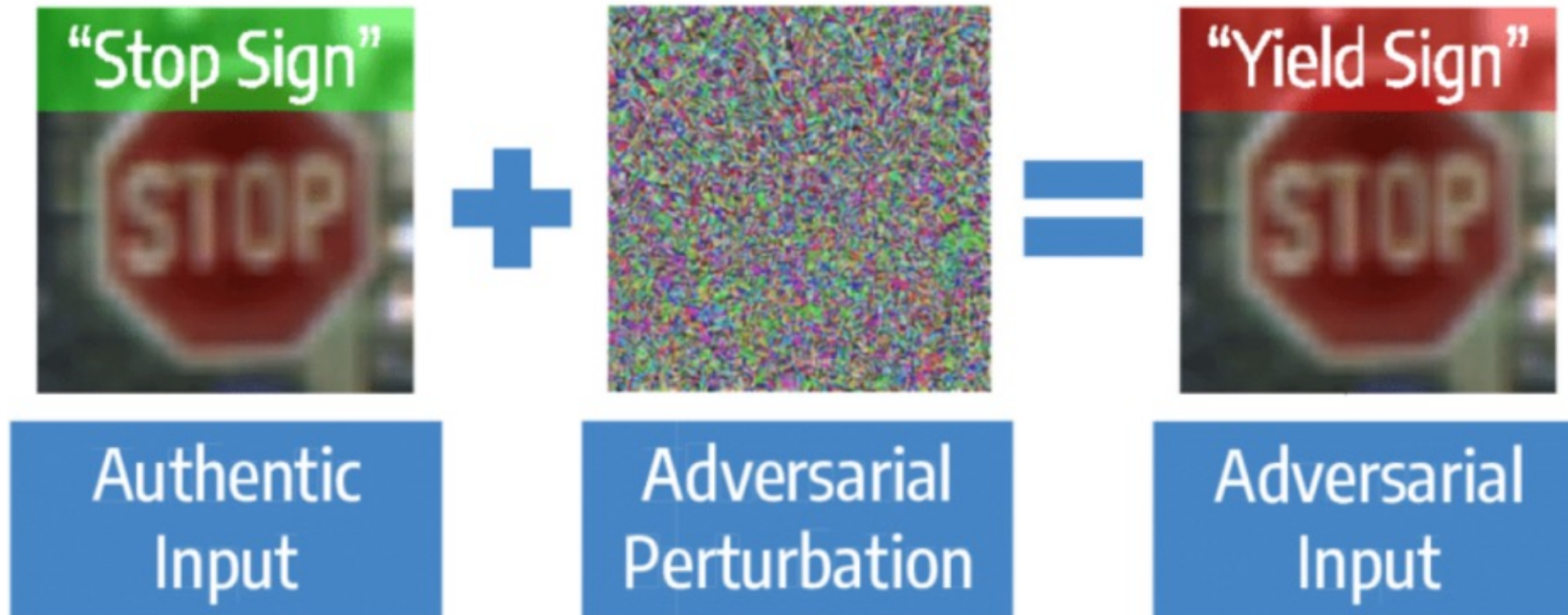


Basic assumption in machine learning

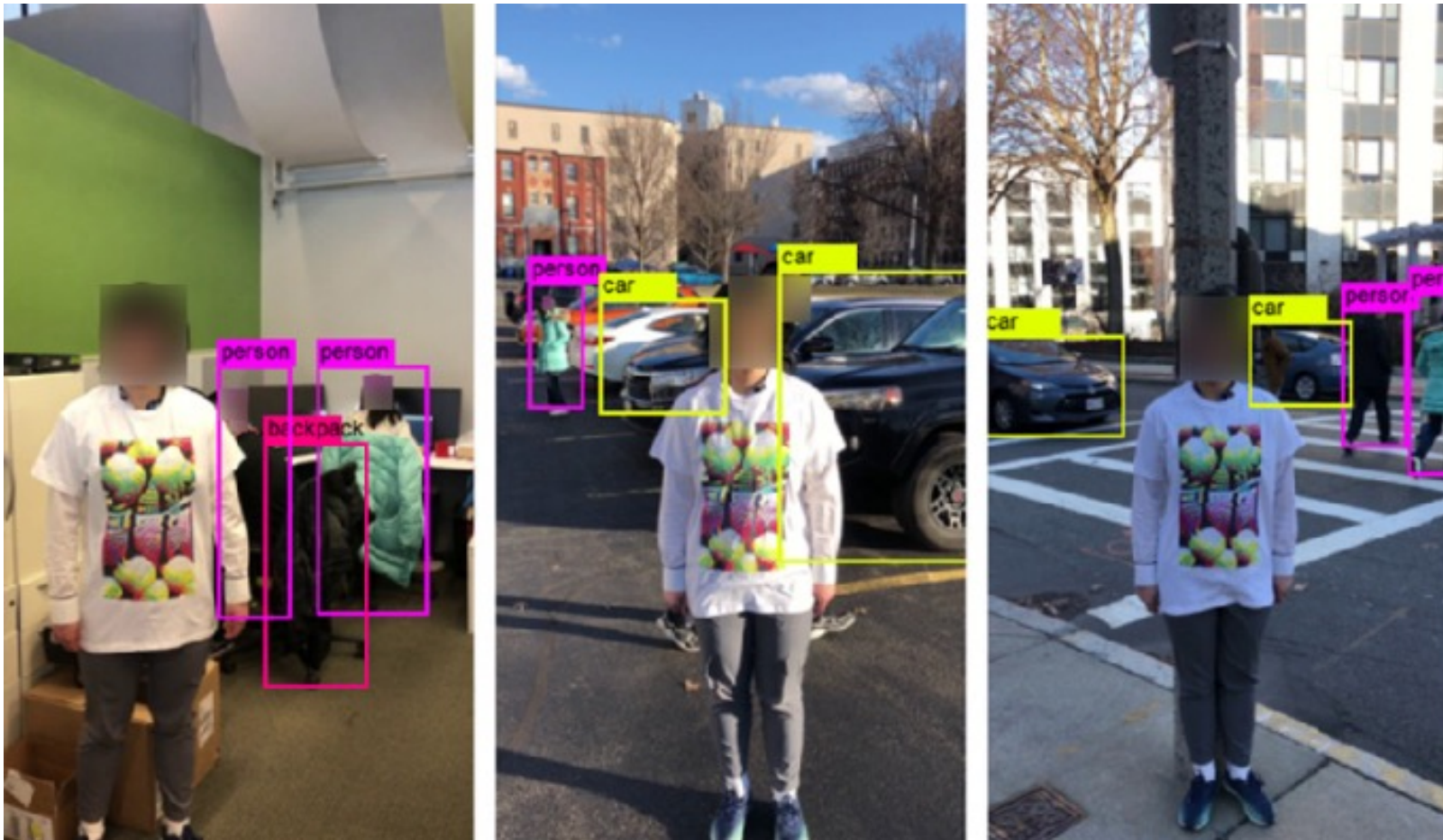
Why do we care?

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**.



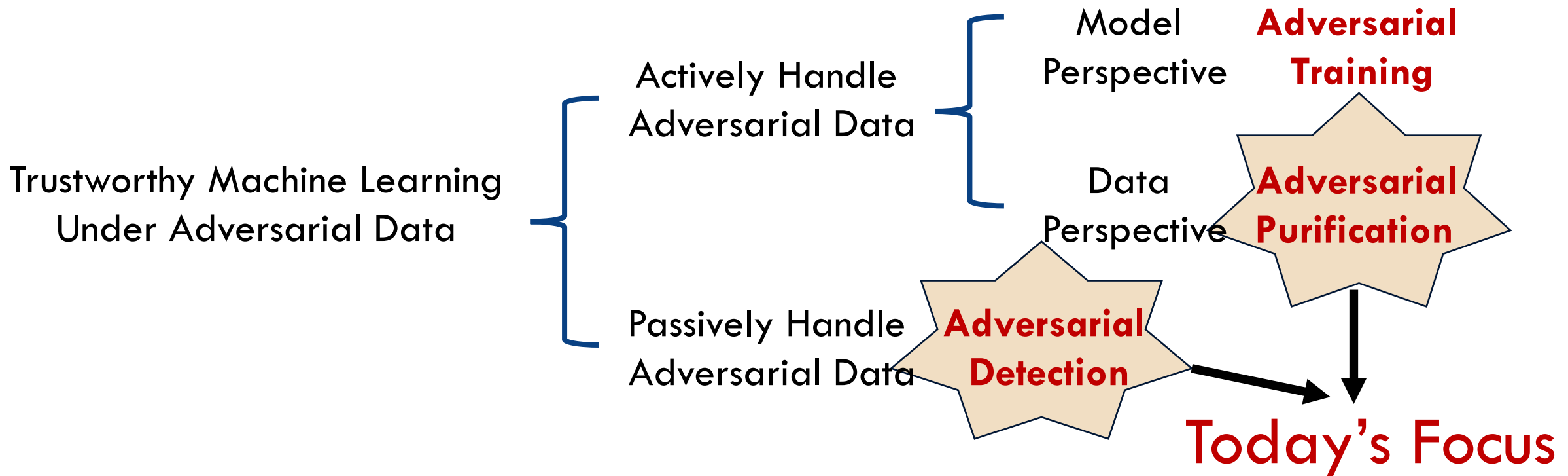
Why do we care?



- ❑ 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 it?

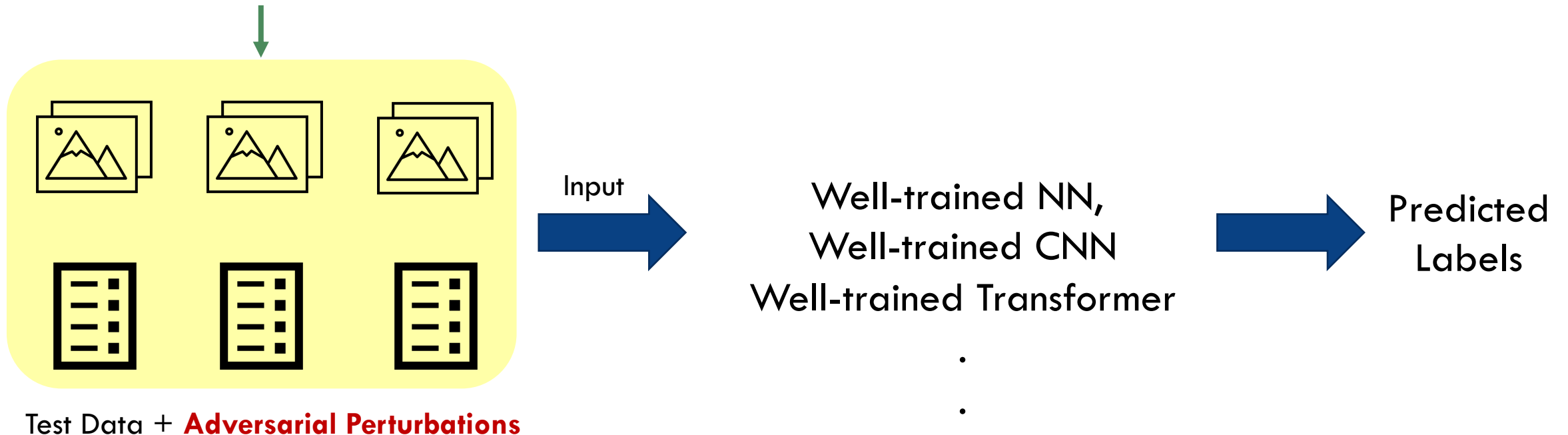
Defend against adversarial attacks



Adversarial detection

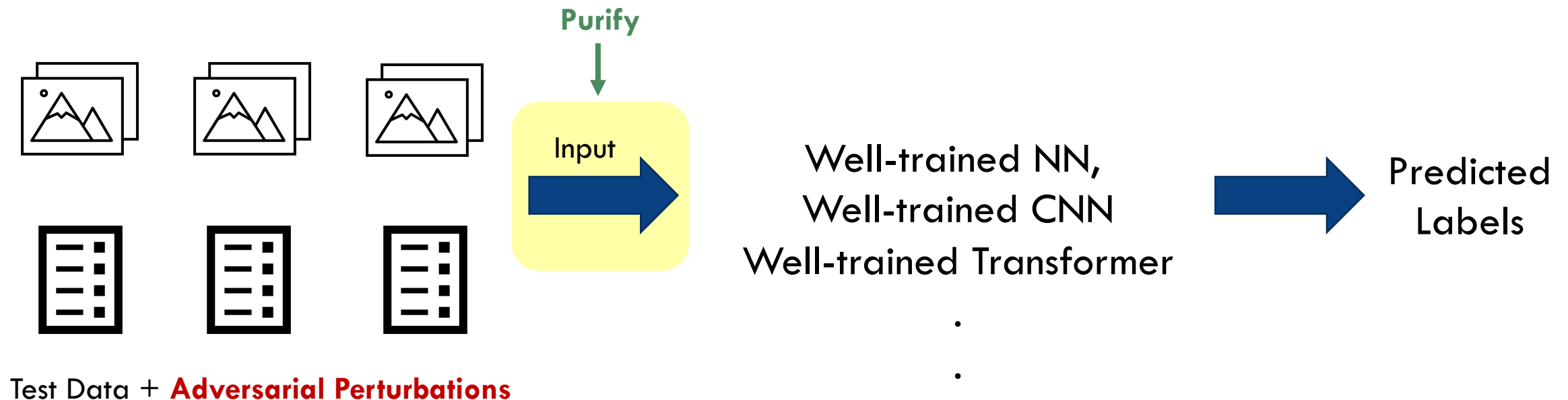
❑ *Adversarial Detection (AD)*: aims to detect and discard AEs.

Discard the adversarial data



Adversarial purification

□ *Adversarial Purification* (AP): aims to shift AEs back towards their natural counterparts.



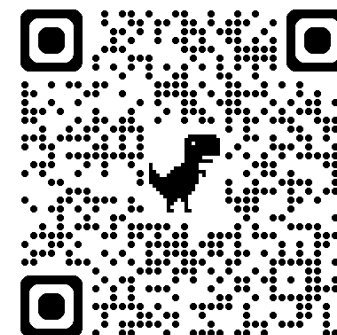
Sample-specific Noise Injection for Diffusion-based Adversarial Purification

Yuhao Sun[^], Jiacheng Zhang[^], Zesheng Ye[^], Chaowei Xiao, Feng Liu^{*}

([^] Co-first authors, ^{*} Corresponding authors)

In *ICML*, 2025.

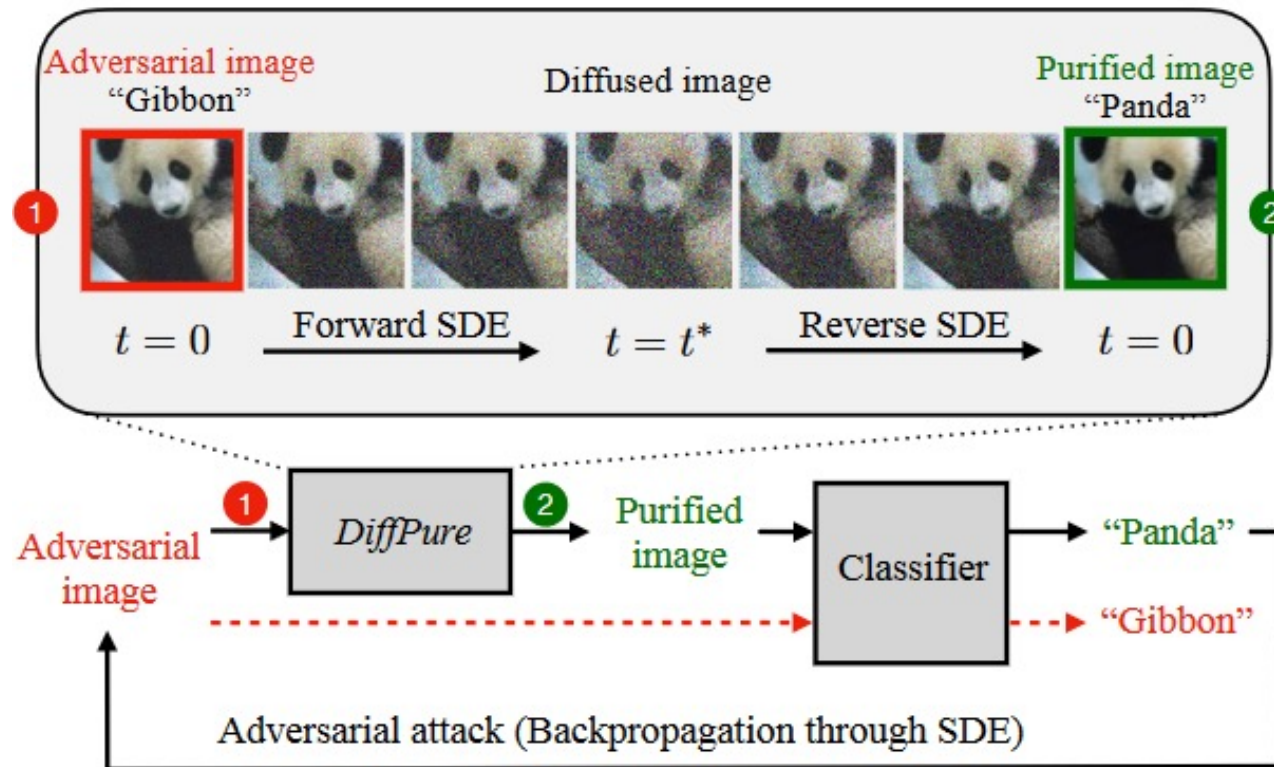
Paper



Code



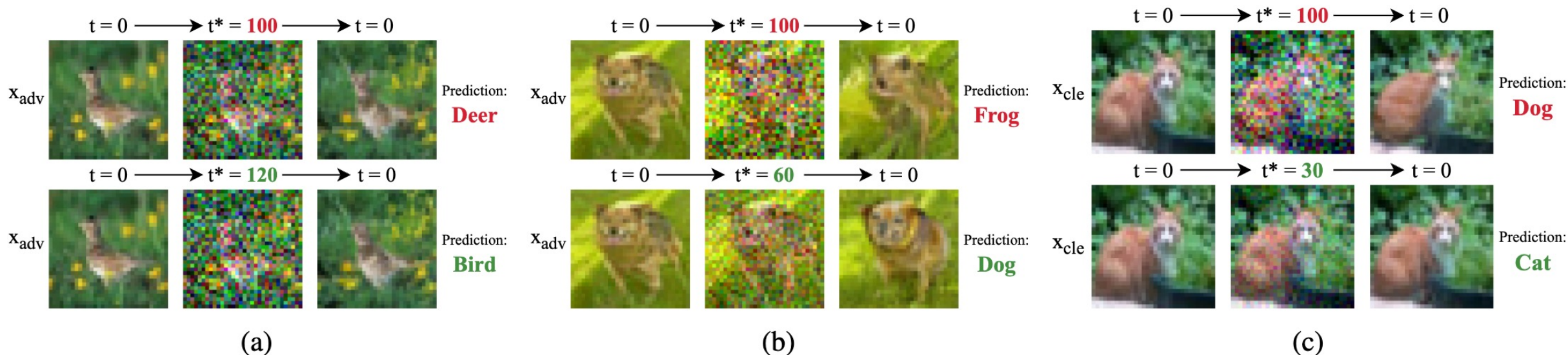
Preliminary: diffusion-based adversarial purification



A Key Challenge: The Choice of t

- ❑ If t is too small, then adversarial noise cannot be fully removed.
- ❑ If t is too large, then the purified image may have a different semantic meaning.
- ❑ **Research gap:** current methods empirically select a *fixed* timestep t for all images, which is *counterintuitive*.

Motivation



- ❑ Sample-shared noise level *fail* to address diverse adversarial perturbations.
- ❑ These findings *highlight* the need for sample-specific noise injection levels.

What is the metric?

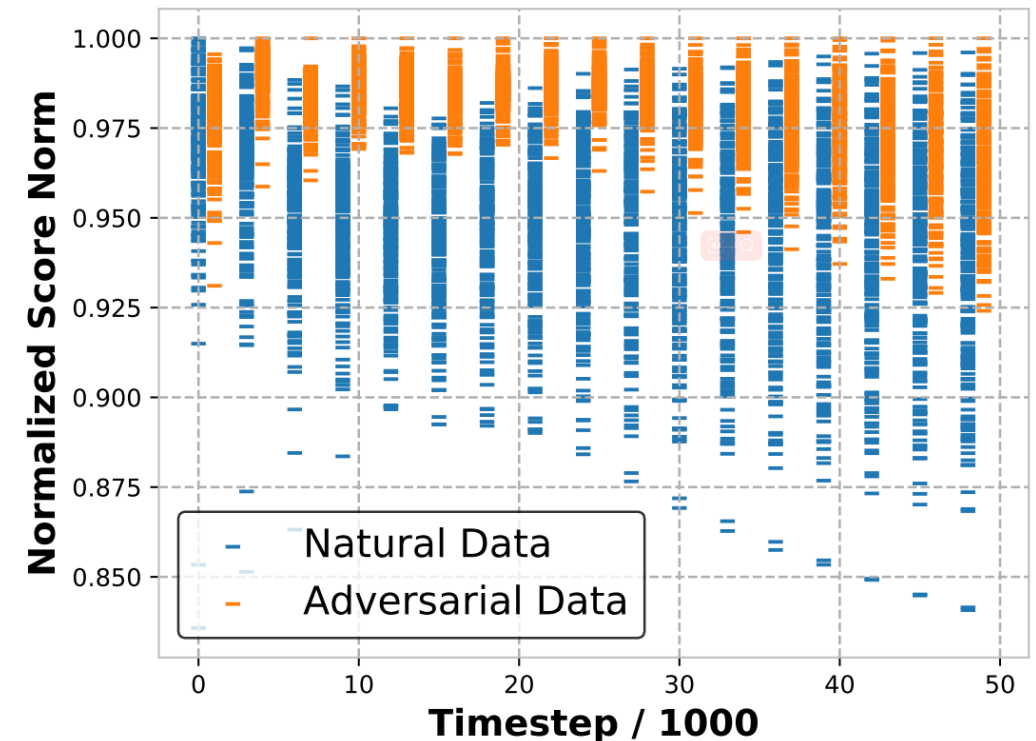
Intuition from score function

□ Intuition from score function $\nabla_{\mathbf{x}} \log p(\mathbf{x})$

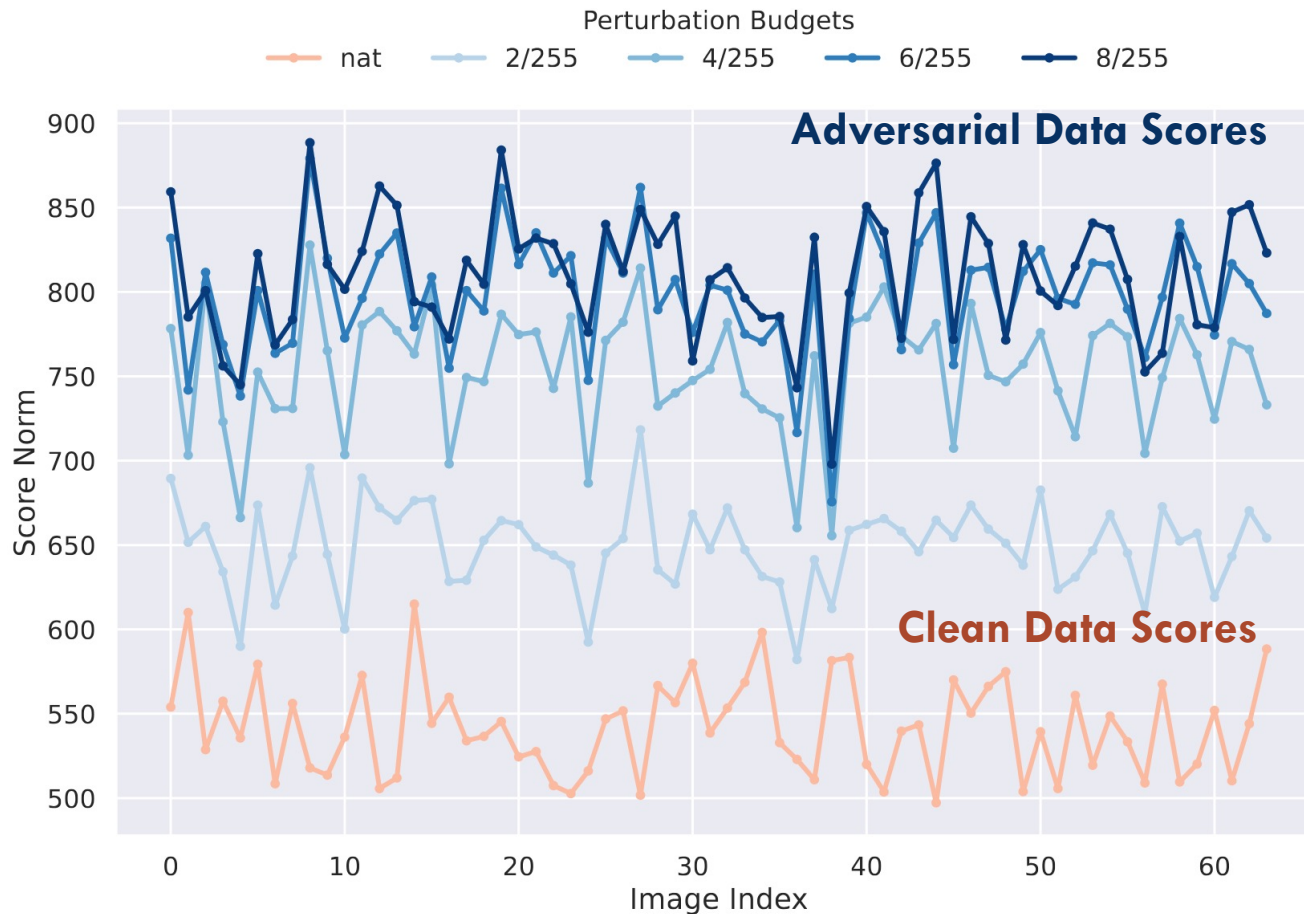
- Score $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ represents the momentum of the sample towards **high density areas** of natural data distribution (Song et al., 2019)



- A **lower score norm** $\|\nabla_{\mathbf{x}} \log p(\mathbf{x})\|$ indicates the sample is **closer** to the high-density areas of natural data distribution

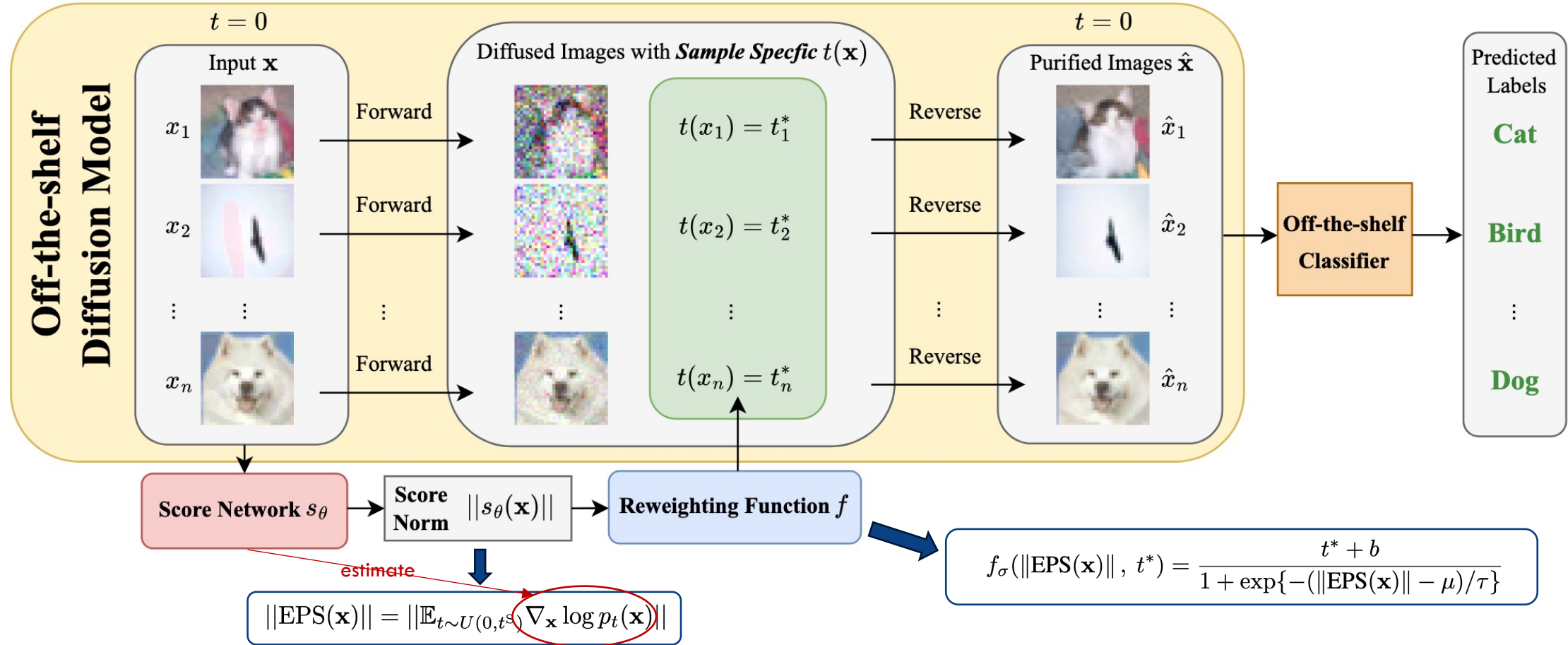


Score norms vs perturbation budgets



- We further find that score norms *scale directly* with perturbation budgets.
- Score norms can act as *proxies* for estimating the sample-specific noise level.

Sample-specific Score-aware Noise Injection (SSNI)



Main results: CIFAR10

PGD+EOT ℓ_∞ ($\epsilon = 8/255$)		
DBP Method	Standard	Robust
WRN-28-10	Nie et al. (2022)	89.71 \pm 0.72
	+ SSNI-N	93.29\pm0.37 (+3.58)
	Wang et al. (2022)	92.45 \pm 0.64
	+ SSNI-N	94.08\pm0.33 (+1.63)
	Lee & Kim (2023)	90.10 \pm 0.18
	+ SSNI-N	93.55\pm0.55 (+2.66)
WRN-70-16	Nie et al. (2022)	90.89 \pm 1.13
	+ SSNI-N	94.47\pm0.51 (+3.58)
	Wang et al. (2022)	93.10 \pm 0.51
	+ SSNI-N	95.57\pm0.24 (+2.47)
	Lee & Kim (2023)	89.39 \pm 1.12
	+ SSNI-N	93.82\pm0.24 (+4.44)

PGD+EOT ℓ_2 ($\epsilon = 0.5$)		
DBP Method	Standard	Robust
WRN-28-10	Nie et al. (2022)	91.80 \pm 0.84
	+ SSNI-N	93.95\pm0.70 (+2.15)
	Wang et al. (2022)	92.45 \pm 0.64
	+ SSNI-N	94.08\pm0.33 (+1.63)
	Lee & Kim (2023)	90.10 \pm 0.18
	+ SSNI-N	93.55\pm0.55 (+3.45)
WRN-70-16	Nie et al. (2022)	92.90 \pm 0.40
	+ SSNI-N	95.12\pm0.58 (+2.22)
	Wang et al. (2022)	93.10 \pm 0.51
	+ SSNI-N	95.57\pm0.24 (+2.47)
	Lee & Kim (2023)	89.39 \pm 1.12
	+ SSNI-N	93.82\pm0.24 (+4.43)

Main results: ImageNet-1K

PGD+EOT ℓ_∞ ($\epsilon = 4/255$)		
DBP Method	Standard	Robust
RN-50	Nie et al. (2022)	68.23 \pm 0.92
	+ SSNI-N	70.25 \pm 0.56 (+2.02)
	Wang et al. (2022)	30.34 \pm 0.72
	+ SSNI-N	33.66 \pm 1.04 (+3.32)
	Lee & Kim (2023)	0.39 \pm 0.03
	+ SSNI-N	5.21 \pm 0.24 (+4.82)
	Lee & Kim (2023)	42.45 \pm 0.92
	+ SSNI-N	43.48 \pm 0.25 (+1.03)

Limitations of DBP framework & SSNI

- ❑ **Limitation 1:** Having a pre-trained diffusion model is not always feasible, training a diffusion model is resource-consuming.
- ❑ **Limitation 2:** The inference speed of DBP-based methods is slow.
- ❑ **Limitation 3:** SSNI still injects noise to clean samples, which cannot fully preserve the utility (i.e., clean accuracy) of the model.

Can we do better?

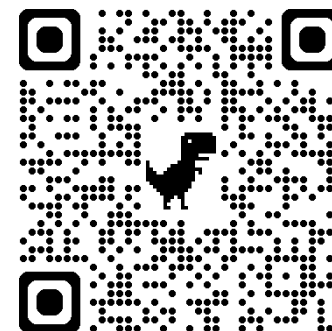
One Stone, Two Birds: Enhancing Adversarial Defense Through the Lens of Distributional Discrepancy

Jiacheng Zhang, Benjamin I. P. Rubinstein, Jingfeng Zhang, Feng Liu*

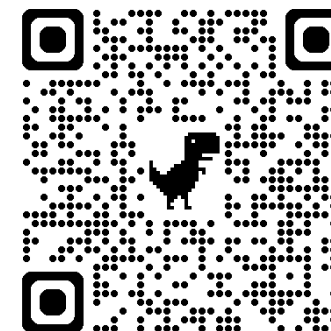
(* Corresponding authors)

In *ICML*, 2025.

Paper



Code



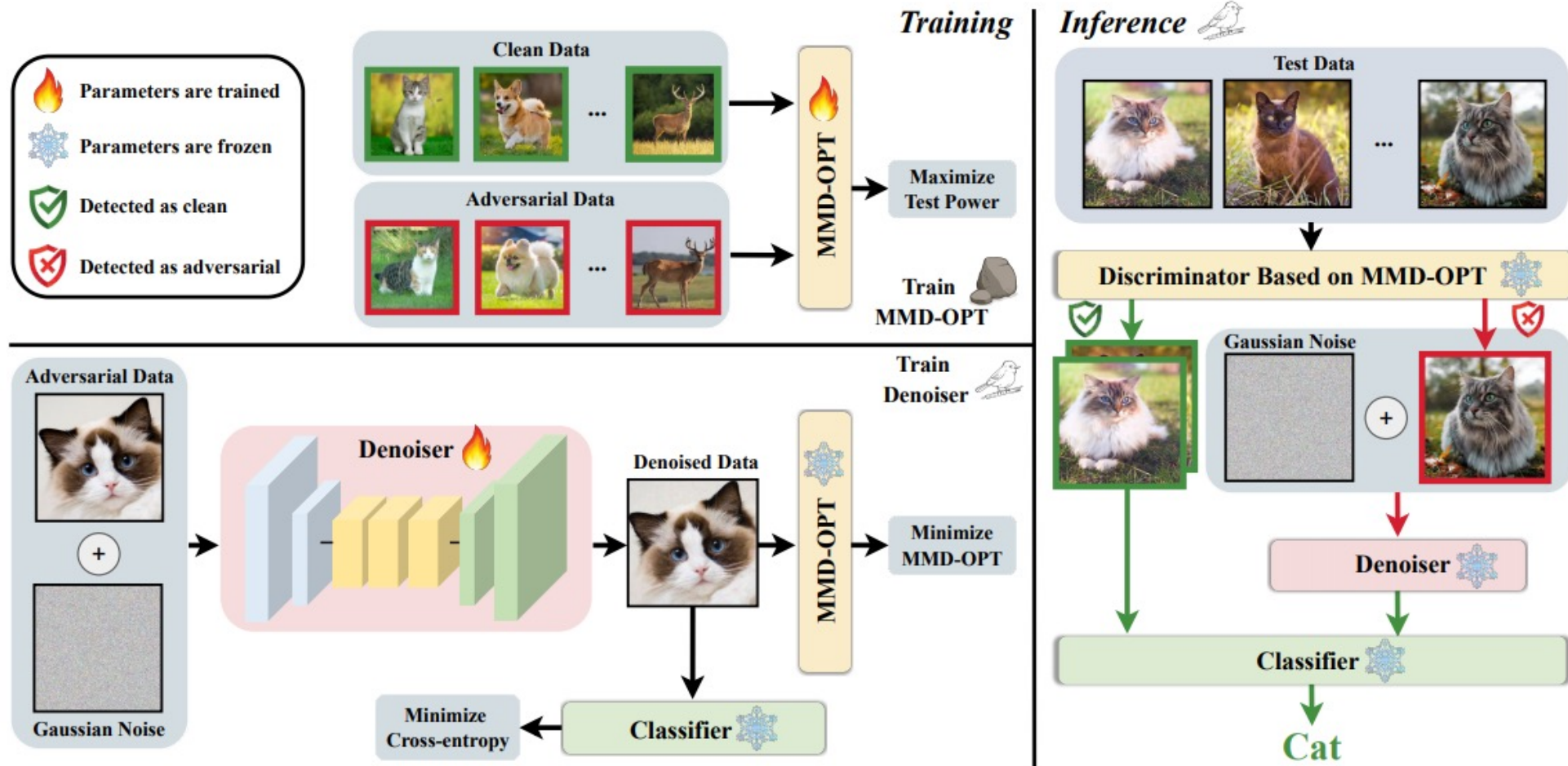
Distributional discrepancy minimization improves robustness

Theorem 1. *For a hypothesis $h \in \mathcal{H}$ and a distribution $\mathcal{D}_A \in \mathbb{D}$:*

$$R(h, f_A, \mathcal{D}_A) \leq R(h, f_c, \mathcal{D}_c) + d_1(\mathcal{D}_c, \mathcal{D}_A).$$

risk on adversarial data **risk on clean data** **distributional discrepancy**

Distributional-discrepancy-based Adversarial Defense (DAD)

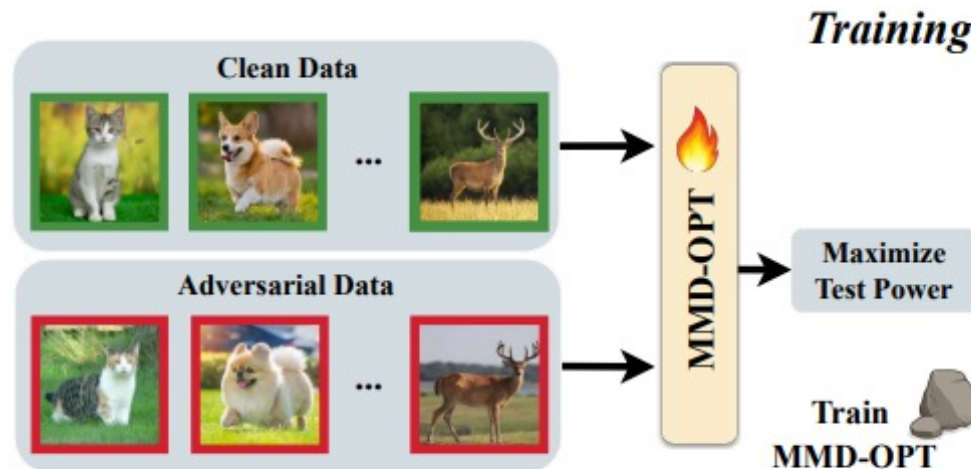


One stone: optimized MMD

$$\text{MMD-OPT}(S'_X, S'_Z) = \widehat{\text{MMD}}_u^2(S'_X, S'_Z; k_\omega^*)$$

MMD values

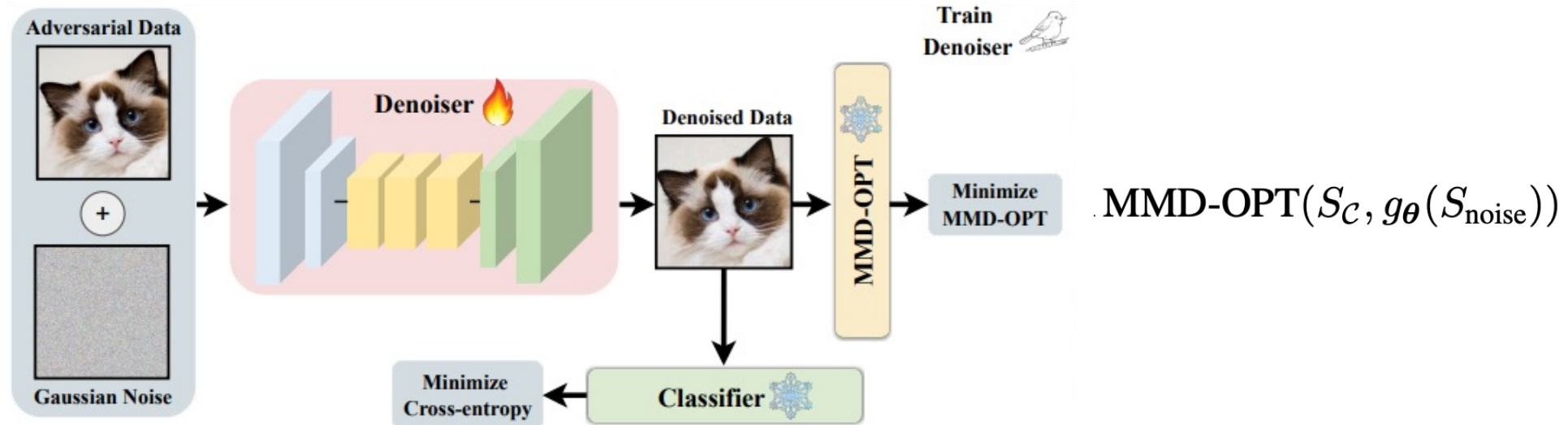
**0 if same distribution
1 if different**



Algorithm 1 Optimizing MMD (Liu et al., 2020).

- 1: **Input:** clean data S_C^{train} , adversarial data S_A^{train} , learning rate η , epoch T ;
- 2: Initialize $\omega \leftarrow \omega_0$; $\lambda \leftarrow 10^{-8}$;
- 3: **for** epoch = 1, ..., T **do**
- 4: $S'_C \leftarrow$ minibatch from S_C^{train} ;
- 5: $S'_A \leftarrow$ minibatch from S_A^{train} ;
- 6: $k_\omega \leftarrow$ kernel function with parameters ω using Eq. (3);
- 7: $M(\omega) \leftarrow \widehat{\text{MMD}}_u^2(S'_C, S'_A; k_\omega)$ using Eq. (2);
- 8: $V_\lambda(\omega) \leftarrow \hat{\sigma}_\lambda(S'_C, S'_A; k_\omega)$ using Eq. (5);
- 9: $\hat{J}_\lambda(\omega) \leftarrow M(\omega) / \sqrt{V_\lambda(\omega)}$ using Eq. (4);
- 10: $\omega \leftarrow \omega + \eta \nabla_{\text{Adam}} \hat{J}_\lambda(\omega)$;
- 11: **end for**
- 12: **Output:** k_ω^*

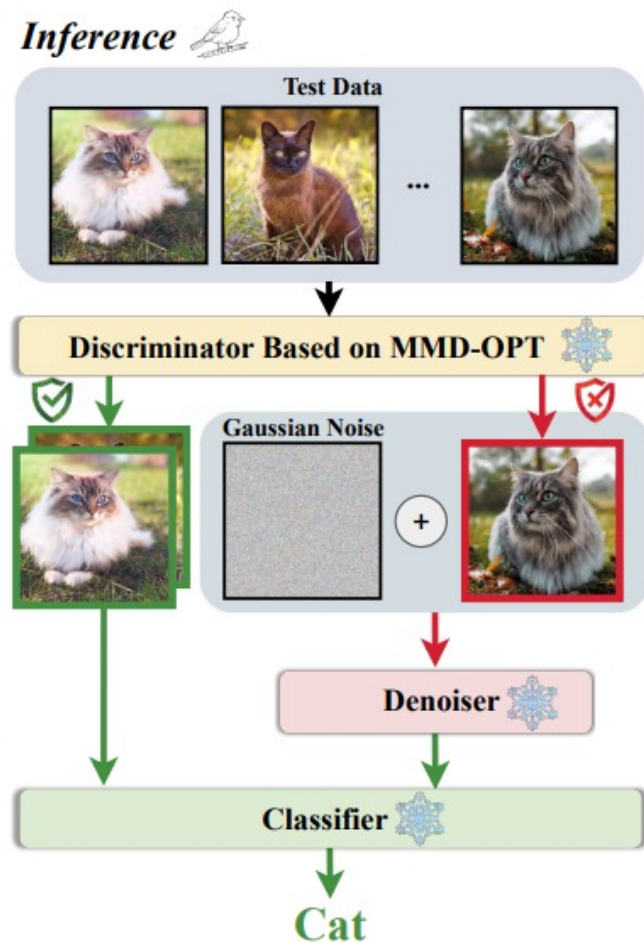
First bird: MMD-OPT-based denoiser



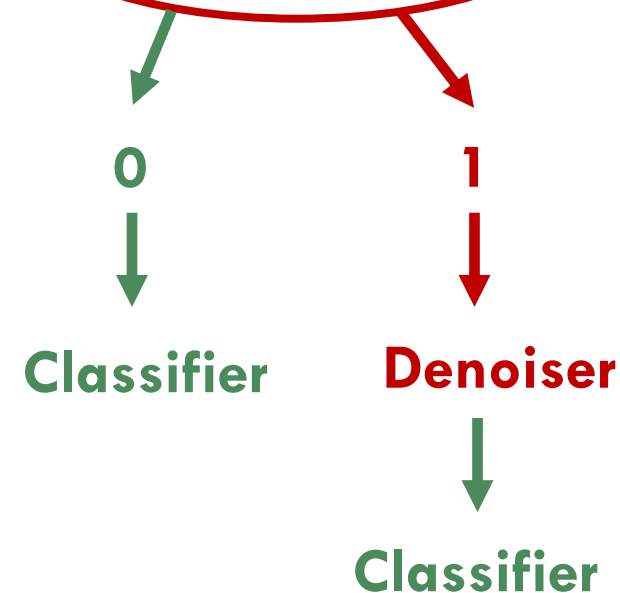
$$\mathcal{L}_{\text{ce}}(\widehat{h}_C^*(g_\theta(S_{\text{noise}})), Y_C)$$

$$g_{\theta^*} = \arg \min_{\theta} \text{MMD-OPT}(S_C, g_\theta(S_{\text{noise}})) + \alpha \cdot \mathcal{L}_{\text{ce}}(\widehat{h}_C^*(g_\theta(S_{\text{noise}})), Y_C)$$

Second bird: MMD-OPT-based discriminator



$$\text{MMD-OPT}(S_{\mathcal{V}}, S_{\mathcal{T}}) = \widehat{\text{MMD}}_u^2(S_{\mathcal{V}}, S_{\mathcal{T}}; k_{\omega}^*)$$



Main results: CIFAR-10

ℓ_∞ ($\epsilon = 8/255$)			
Type	Method	Clean	Robust
WRN-28-10			
AT	Gowal et al. (2021)	87.51	63.38
	Gowal et al. (2020)*	88.54	62.76
	Pang et al. (2022a)	88.62	61.04
AP	Yoon et al. (2021)	85.66	33.48
	Nie et al. (2022)	90.07	46.84
	Lee & Kim (2023)	90.16	55.82
Ours	DAD	94.16 \pm 0.08	67.53 \pm 1.07
WRN-70-16			
AT	Rebuffi et al. (2021)*	92.22	66.56
	Gowal et al. (2021)	88.75	66.10
	Gowal et al. (2020)*	91.10	65.87
AP	Yoon et al. (2021)	86.76	37.11
	Nie et al. (2022)	90.43	51.13
	Lee & Kim (2023)	90.53	56.88
Ours	DAD	93.91 \pm 0.11	67.68 \pm 0.87

ℓ_2 ($\epsilon = 0.5$)			
Type	Method	Clean	Robust
WRN-28-10			
AT	Rebuffi et al. (2021)*	91.79	78.80
	Augustin et al. (2020) [†]	93.96	78.79
	Schwag et al. (2022) [†]	90.93	77.24
AP	Yoon et al. (2021)	85.66	73.32
	Nie et al. (2022)	91.41	79.45
	Lee & Kim (2023)	90.16	83.59
Ours	DAD	94.16 \pm 0.08	84.38 \pm 0.81
WRN-70-16			
AT	Rebuffi et al. (2021)*	95.74	82.32
	Gowal et al. (2020)*	94.74	80.53
	Rebuffi et al. (2021)	92.41	80.42
AP	Yoon et al. (2021)	86.76	75.66
	Nie et al. (2022)	92.15	82.97
	Lee & Kim (2023)	90.53	83.57
Ours	DAD	93.91 \pm 0.11	84.03 \pm 0.75

Main results: ImageNet-1K

ℓ_∞ ($\epsilon = 4/255$)			
Type	Method	Clean	Robust
RN-50			
AT	Salman et al. (2020a)	64.02	34.96
	Engstrom et al. (2019)	62.56	29.22
	Wong et al. (2020)	55.62	26.24
AP	Nie et al. (2022)	71.48	38.71
	Lee & Kim (2023)	70.74	42.15
Ours	DAD	78.61 \pm 0.04	53.85 \pm 0.23

- ❑ **Strength 1:** DAD can largely preserve the original utility (i.e., clean accuracy of the classifier).
- ❑ **Strength 2:** Compared to DBP methods that rely on density estimation, learning distributional discrepancies is a simpler and more feasible task.
- ❑ **Strength 3:** DAD is efficient in both training and inferencing.

Thank You!