



Australian  
National  
University

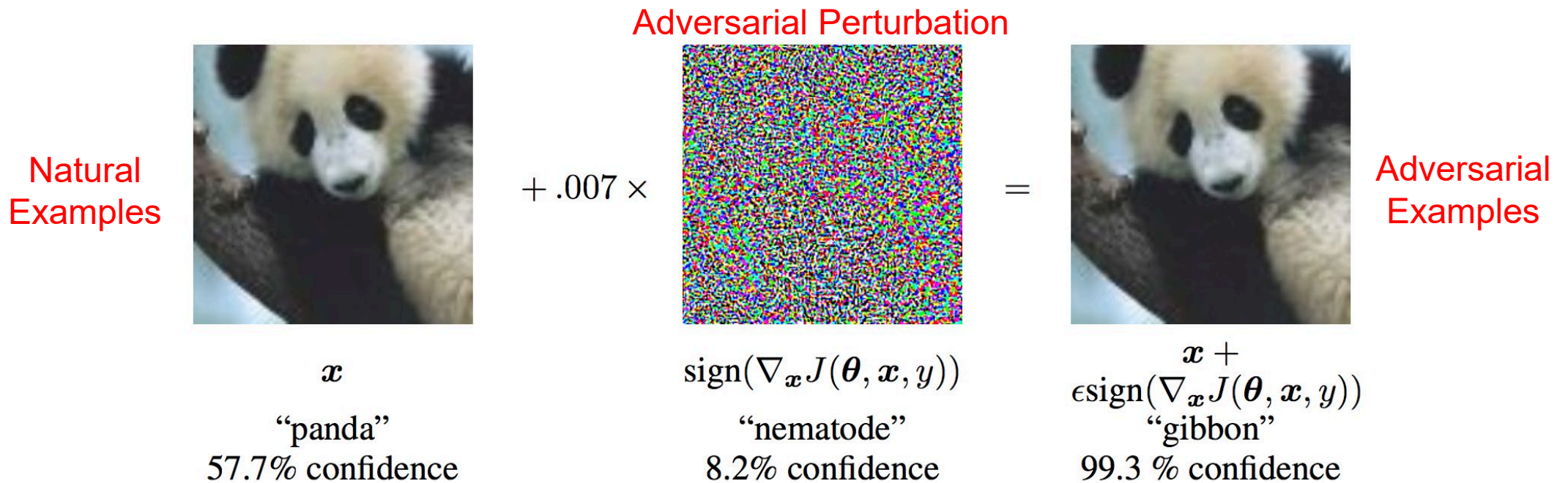


# Robust Distillation via Untargeted and Targeted Intermediate Adversarial Samples

Junhao Dong, Piotr Koniusz, Junxi Chen, Z. Jane Wang, Yew-Soon Ong

Reporter: Junhao Dong

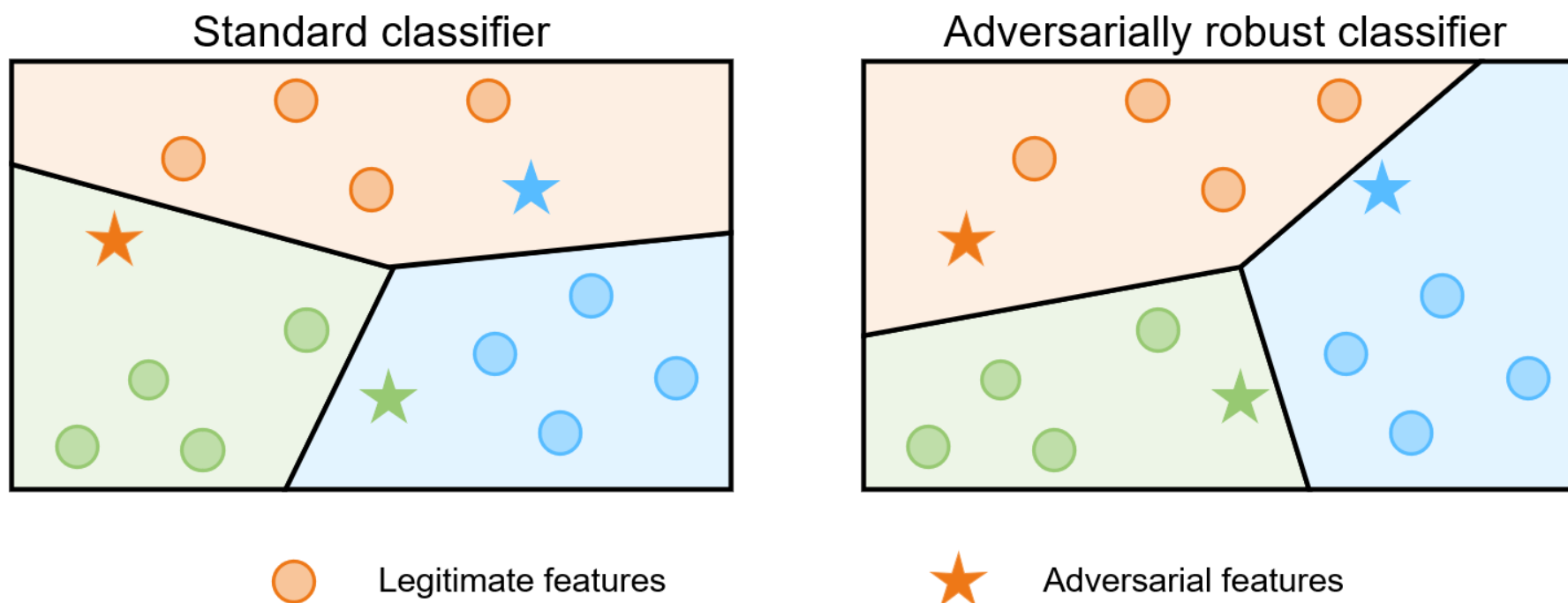
**Adversarial examples** are tailored inputs with the purpose of confusing neural networks. (Visually similar to natural examples)



Introducing gradient ascent at the **image level**.

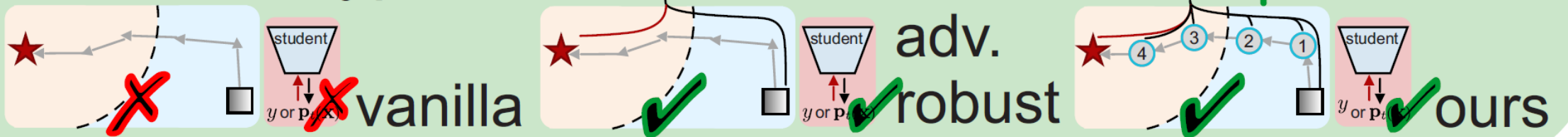
## Adversarial Training (min-max optimization):

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[ \mathcal{L}_{\text{CE}}(f_{\theta}(\mathbf{x}), y) + \max_{\|\delta\|_{\infty} < \epsilon} \mathcal{L}_{\text{KL}}(f_{\theta}(\mathbf{x}) \| f_{\theta}(\mathbf{x} + \delta)) \right]$$

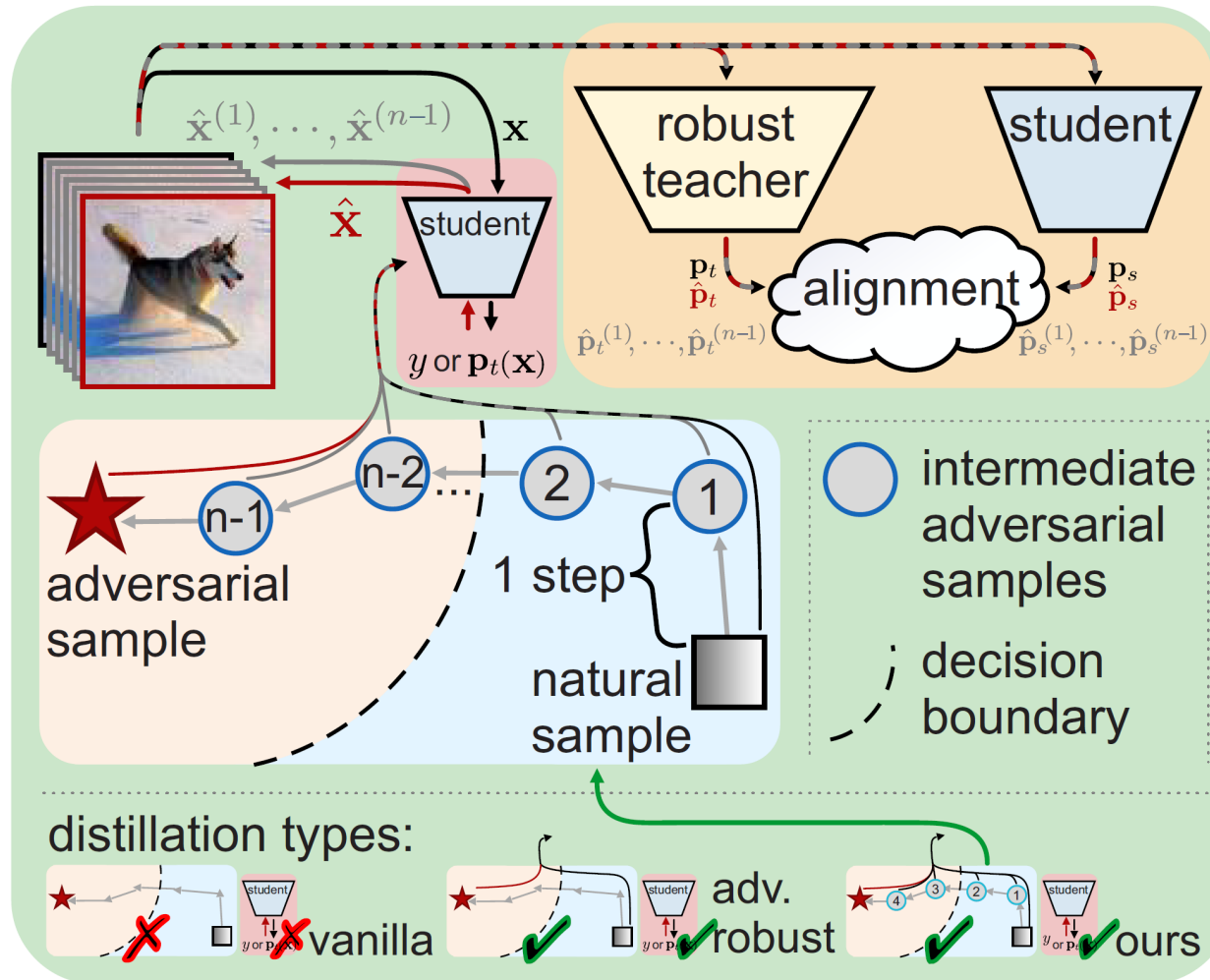


## Adversarially Robust Knowledge Distillation

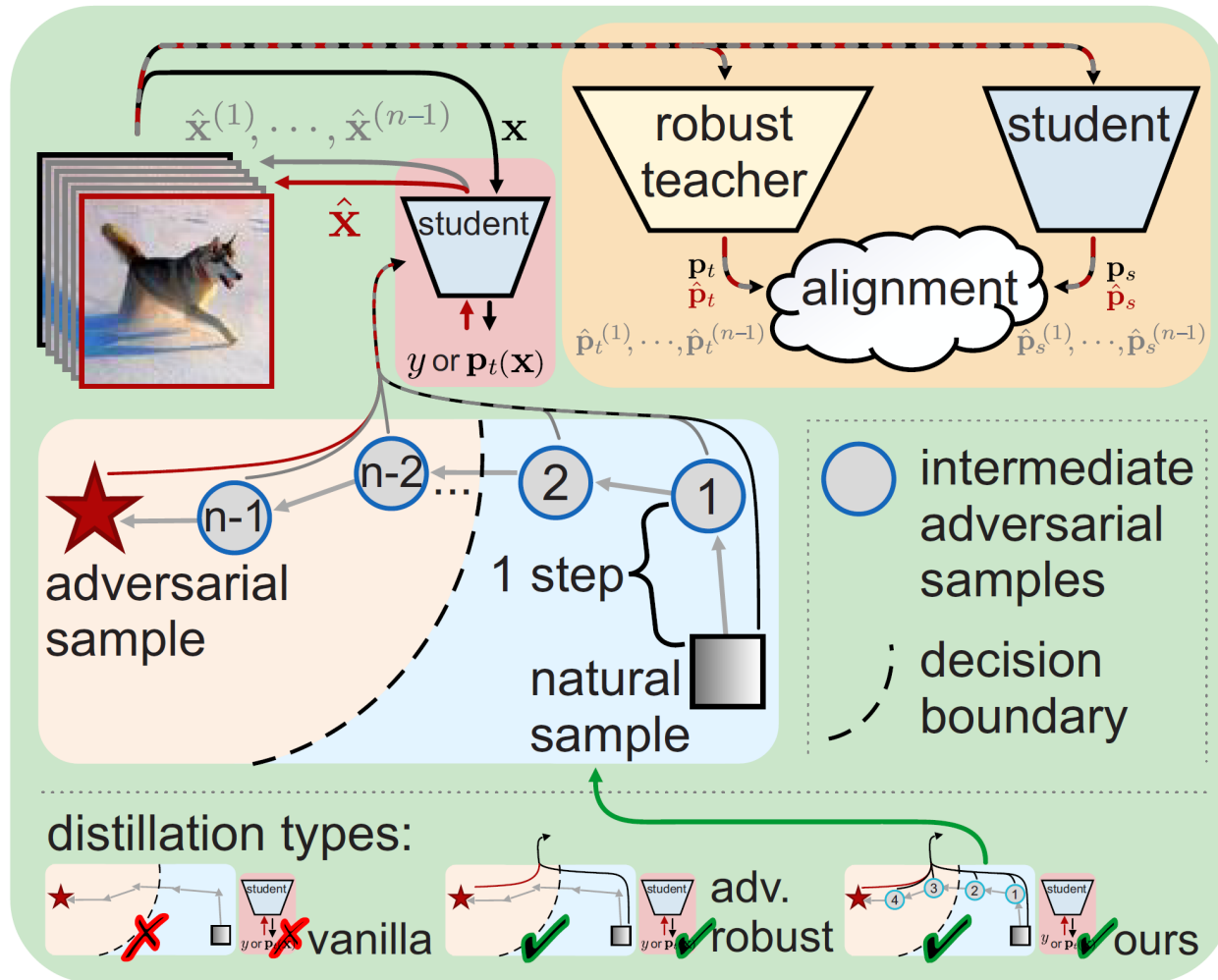
distillation types:



## Adversarially Robust Knowledge Distillation



## Adversarially Robust Knowledge Distillation



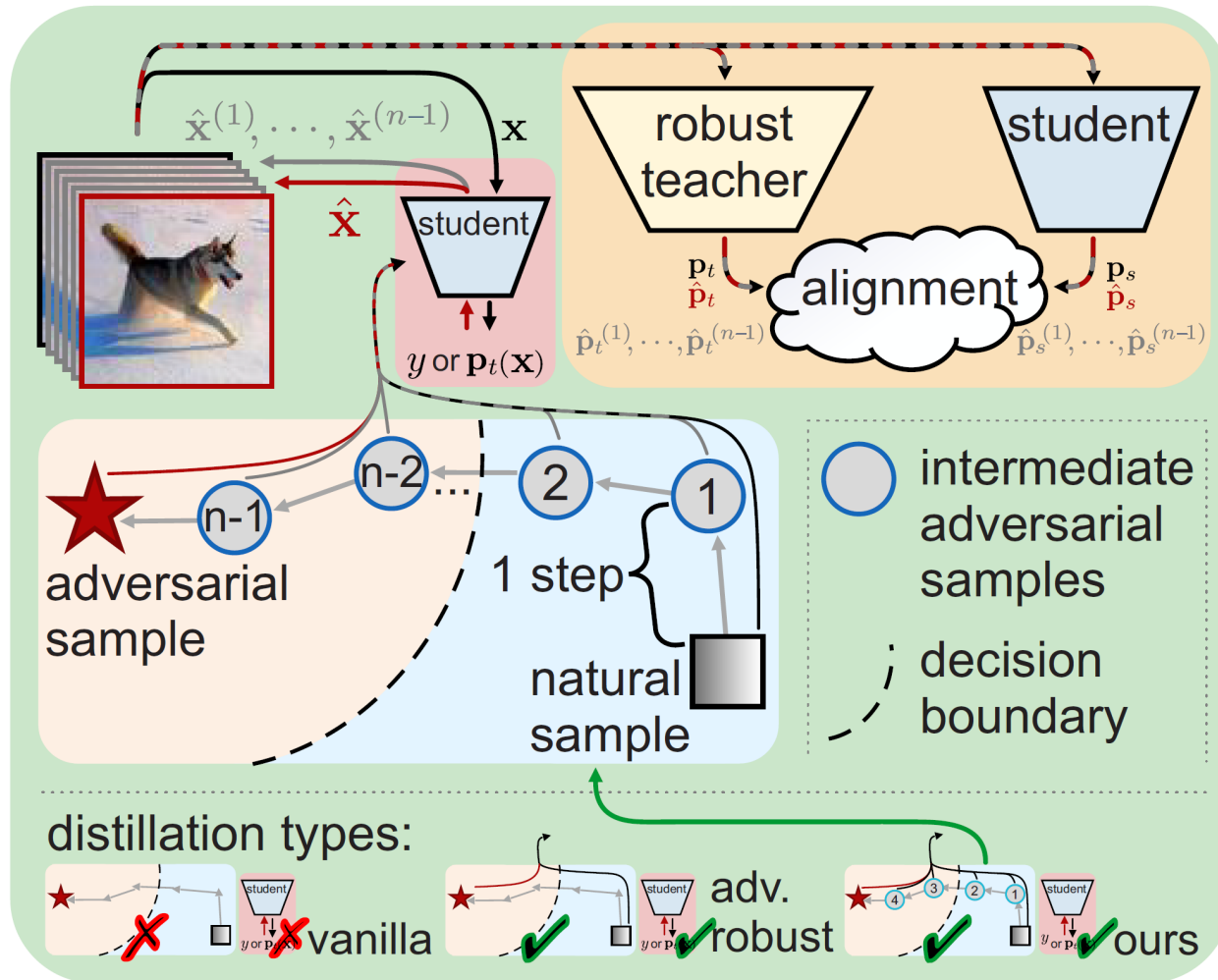
### Untargeted Adversary Generation

$$\begin{aligned}\hat{\mathbf{x}}^{(i+1)} &= \vartheta_{\alpha}(\hat{\mathbf{x}}^{(i)}, y) \\ &= \Pi_{\mathbb{B}(\mathbf{x}, \epsilon)} \left( \hat{\mathbf{x}}^{(i)} + \alpha \cdot \text{sign} \left( \nabla_{\hat{\mathbf{x}}^{(i)}} \mathcal{L}_{\text{CE}}(f_{\theta}(\hat{\mathbf{x}}^{(i)}), y) \right) \right)\end{aligned}$$

### Prediction Alignment of Clean and Adversarial Samples

$$\mathcal{L}_{\text{ARKD}} = \underbrace{\mathcal{L}_{\text{KL}}(f_{\theta_t}(\mathbf{x}) \| f_{\theta_s}(\mathbf{x}))}_{\text{alignment of "natural distributions"}} + \beta \cdot \underbrace{\mathcal{L}_{\text{KL}}(f_{\theta_t}(\hat{\mathbf{x}}) \| f_{\theta_s}(\hat{\mathbf{x}}))}_{\text{alignment of "adversarial distributions"}}$$

## Intermediate Adversarial Knowledge Distillation



### Objective Function

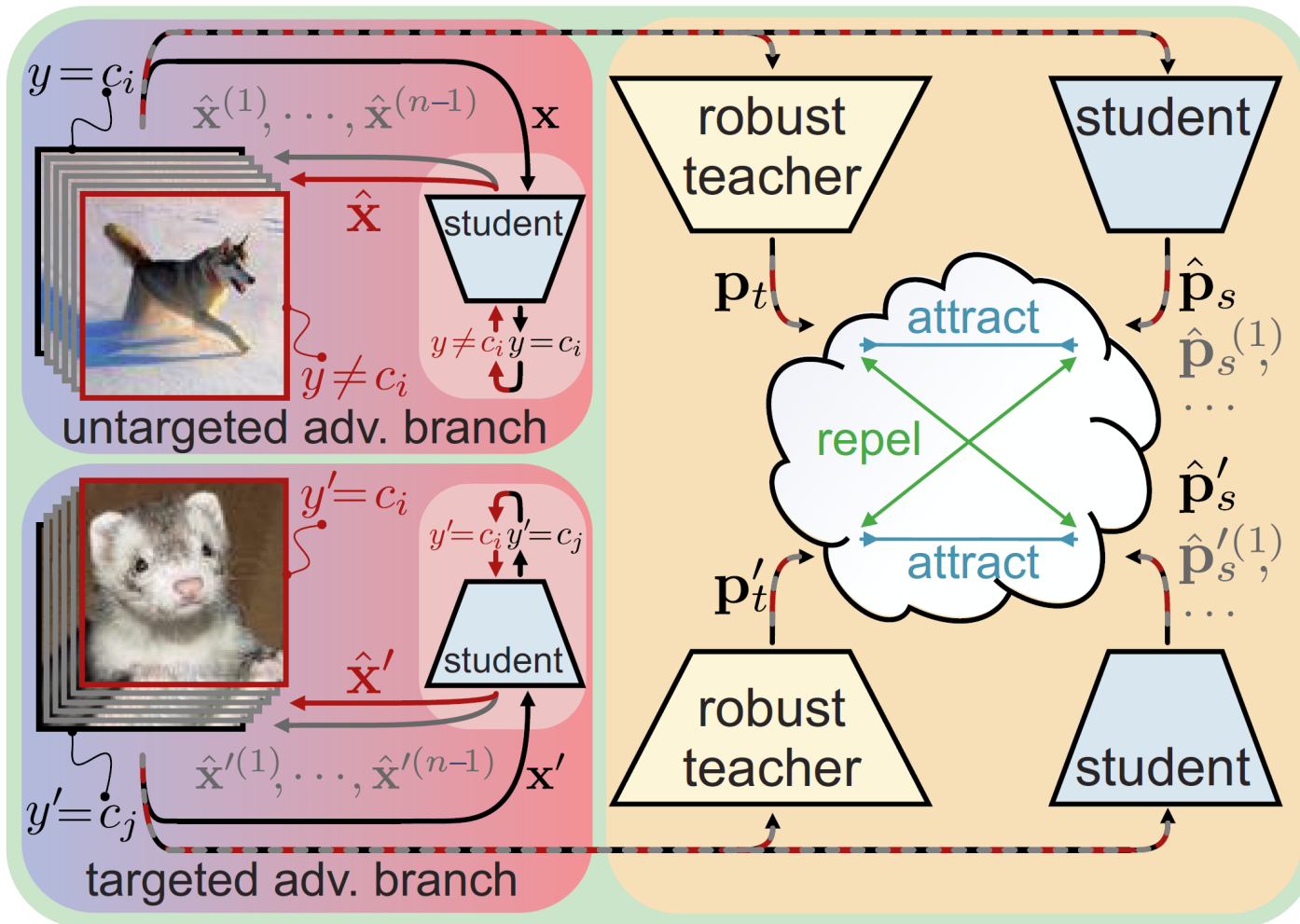
$$\mathcal{L}_{\text{IAKD}} = \sum_{i=1}^{n-1} w(\hat{\mathbf{x}}^{(i)} | \mathbf{x}) \cdot \mathcal{L}_{\text{KL}}(f_{\theta_t}(\hat{\mathbf{x}}^{(i)}) \| f_{\theta_s}(\hat{\mathbf{x}}^{(i)}))$$

### Re-weighting

$$w(\hat{\mathbf{x}}^{(i)} | \mathbf{x}) = \frac{(1-\gamma) i}{n} + \frac{\gamma |(f_{\theta_t}(\mathbf{x}))_y - (f_{\theta_s}(\hat{\mathbf{x}}^{(i)}))_y|}{\max_{j \in \mathcal{B}} |(f_{\theta_t}(\mathbf{x}_j))_{y_j} - (f_{\theta_s}(\hat{\mathbf{x}}_j^{(i)}))_{y_j}|}$$

We theoretically prove that **such a weighting mechanism captures the localized  $\kappa$ -Lipschitz smoothness** of the student model.

## ■ Dual-branch Adversarially Robust knowledge distillation (DARWIN)



Attraction

$$\text{attract} \begin{cases} \{f_{\theta_s}(\hat{\mathbf{x}}^{(i)})\}_{i=1}^n \rightarrow f_{\theta_t}(\mathbf{x}) \\ \{f_{\theta_s}(\hat{\mathbf{x}}'^{(i)})\}_{i=1}^n \rightarrow f_{\theta_t}(\mathbf{x}') \end{cases}$$

Repulsion

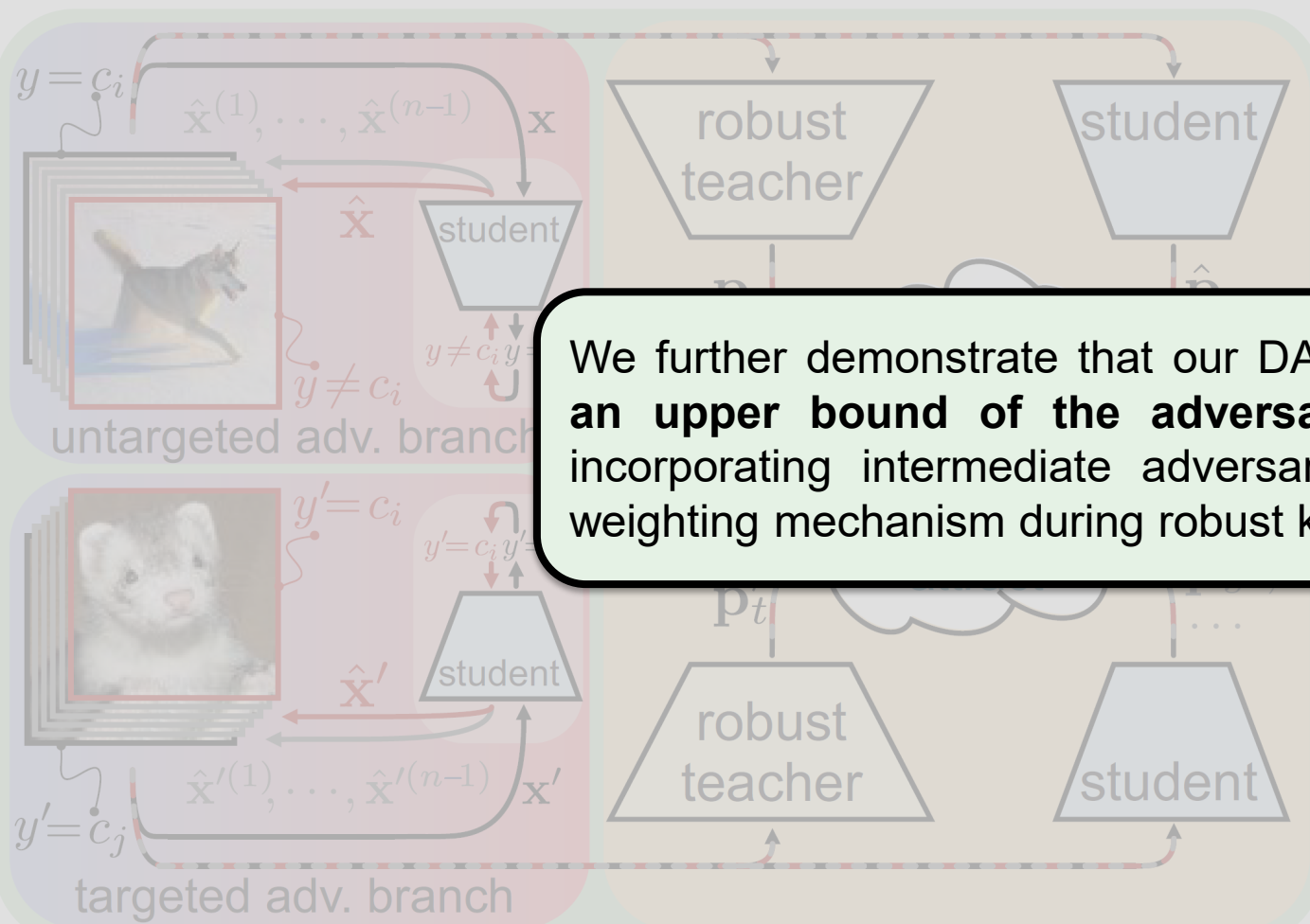
$$\text{repel} \begin{cases} \{f_{\theta_s}(\hat{\mathbf{x}}^{(i)})\}_{i=1}^n \leftrightarrow f_{\theta_t}(\mathbf{x}') \\ \{f_{\theta_s}(\hat{\mathbf{x}}'^{(i)})\}_{i=1}^n \leftrightarrow f_{\theta_t}(\mathbf{x}) \end{cases}$$

Dual-Branch Knowledge Distillation

$$\mathcal{L}_{\text{DBKD}} = \sum_{i=1}^n \left[ w(\hat{\mathbf{x}}^{(i)}|\mathbf{x}) \mathcal{L}_{\text{tri}}(f_{\theta_t}(\mathbf{x}), f_{\theta_s}(\hat{\mathbf{x}}^{(i)}), f_{\theta_s}(\hat{\mathbf{x}}'^{(i)})) \right. \\ \left. + w(\hat{\mathbf{x}}'^{(i)}|\mathbf{x}') \mathcal{L}_{\text{tri}}(f_{\theta_t}(\mathbf{x}'), f_{\theta_s}(\hat{\mathbf{x}}'^{(i)}), f_{\theta_s}(\hat{\mathbf{x}}^{(i)})) \right]$$



## ■ Dual-branch Adversarially Robust knowledge distillation (DARWIN)



We further demonstrate that our DARWIN method **minimizes an upper bound of the adversarially robust risk** when incorporating intermediate adversarial samples with the re-weighting mechanism during robust knowledge distillation.

Attraction

$$\text{attract} \begin{cases} \{f_{\theta_s}(\hat{x}^{(i)})\}_{i=1}^n \rightarrow f_{\theta_t}(\mathbf{x}) \\ \{f_{\theta_s}(\hat{x}'^{(i)})\}_{i=1}^n \rightarrow f_{\theta_t}(\mathbf{x}') \end{cases}$$

$$\begin{aligned} & \leftarrow f_{\theta_t}(\mathbf{x}') \\ & \leftarrow f_{\theta_t}(\mathbf{x}) \end{aligned}$$

Dual-Branch Knowledge Distillation

$$\mathcal{L}_{\text{DBKD}} = \sum_{i=1}^n \left[ w(\hat{x}^{(i)} | \mathbf{x}) \mathcal{L}_{\text{tri}}(f_{\theta_t}(\mathbf{x}), f_{\theta_s}(\hat{x}^{(i)}), f_{\theta_s}(\hat{x}'^{(i)})) \right. \\ \left. + w(\hat{x}'^{(i)} | \mathbf{x}') \mathcal{L}_{\text{tri}}(f_{\theta_t}(\mathbf{x}'), f_{\theta_s}(\hat{x}'^{(i)}), f_{\theta_s}(\hat{x}^{(i)})) \right]$$

## ■ Standard Comparison (Distillation from a Large Model):

Type	Architecture	Method	y	CIFAR-10				CIFAR-100			
				Natural	PGD-20	CW	AA	Natural	PGD-20	CW	AA
<b>Teacher</b>	WRN-34	TRADES [58]	✓	84.92	55.34	54.21	52.55	60.04	31.56	28.64	27.38
<b>Student</b>	ResNet-18	ARD [15]	✓	82.95	52.26	51.69	49.46	57.46	30.14	27.11	25.30
		IAD [63]	✓	82.41	53.06	51.79	49.78	56.38	30.61	27.35	25.51
		RSLAD [64]	✗	83.12	53.91	52.84	51.19	57.23	31.08	28.29	26.62
		CRDND [54]	✗	83.92	52.70	50.95	49.05	58.03	30.16	27.02	25.68
		GACD [1]	✗	82.76	53.42	52.26	50.07	56.82	31.19	27.81	26.12
		<b>DARWIN</b>	✓	<b>84.48</b>	<b>55.07</b>	53.85	52.24	<b>59.12</b>	<b>32.30</b>	<b>28.95</b>	<b>27.26</b>
	<b>DARWIN-LF</b>	✗	84.35	55.02	<b>53.99</b>	<b>52.33</b>	59.04	32.18	28.62	27.13	
	MNV2	ARD [15]	✓	82.44	51.91	50.64	48.40	55.28	30.23	27.05	25.28
		IAD [63]	✓	81.61	52.30	50.19	48.34	54.26	30.46	27.13	25.50
		RSLAD [64]	✗	82.89	52.72	52.04	50.04	57.31	30.48	27.86	25.89
		CRDND [54]	✗	82.77	52.57	50.11	49.28	56.24	29.65	26.68	25.61
		GACD [1]	✗	82.90	52.49	51.40	49.55	56.10	30.49	27.18	25.33
<b>DARWIN</b>		✓	84.06	<b>53.94</b>	<b>53.11</b>	<b>51.28</b>	<b>58.45</b>	<b>31.53</b>	<b>28.36</b>	26.55	
<b>DARWIN-LF</b>	✗	<b>84.08</b>	53.76	52.80	51.09	58.41	31.44	28.33	<b>26.58</b>		

## Robust Distillation from ViTs:

Type	Architecture	Method	Natural	PGD-20	AA
<b>Teacher</b>	ViT-B	AT-PRM [33]	83.98	53.10	49.66
		ARD [15]	82.76	52.95	49.03
<b>Student</b>	ResNet-18	RSLAD [64]	82.33	54.89	49.74
		IAD [63]	82.27	53.42	49.48
		CRDND [54]	82.19	53.16	48.98
		GACD [1]	81.64	54.24	49.95
		<b>DARWIN</b>	<b>83.75</b>	54.80	51.42
		<b>DARWIN-LF</b>	83.73	<b>54.95</b>	<b>51.49</b>
<b>Teacher</b>	DeiT-S	AT-PRM [33]	82.68	52.47	49.27
		ARD [15]	81.59	53.45	49.20
<b>Student</b>	MN2	RSLAD [64]	80.86	53.91	50.18
		IAD [63]	80.41	54.12	49.62
		CRDND [54]	80.27	52.21	48.46
		GACD [1]	79.97	54.00	48.91
		<b>DARWIN</b>	83.02	54.46	<b>51.19</b>
		<b>DARWIN-LF</b>	<b>83.15</b>	<b>54.62</b>	51.13

## Self-Distillation w/ Generated Data:

Dataset	Type	DDPM	Method	Natural	Robust
CIFAR-10	<b>Teacher</b>	$\times$	TRADES [58]	82.45	48.90
		$\checkmark$	ARD [15]	82.89	53.41
	<b>Student</b>	$\checkmark$	RSLAD [64]	82.05	52.60
		$\checkmark$	IAD [63]	82.95	53.47
		$\checkmark$	<b>DARWIN</b>	84.13	55.92
		$\checkmark$	<b>DARWIN-LF</b>	<b>84.68</b>	<b>56.41</b>
CIFAR-100	<b>Teacher</b>	$\times$	TRADES [58]	56.37	23.78
		$\checkmark$	ARD [15]	56.07	26.92
	<b>Student</b>	$\checkmark$	RSLAD [64]	53.40	26.00
		$\checkmark$	IAD [63]	55.82	26.77
		$\checkmark$	<b>DARWIN</b>	58.18	28.24
		$\checkmark$	<b>DARWIN-LF</b>	<b>58.74</b>	<b>28.45</b>

## Ablations:

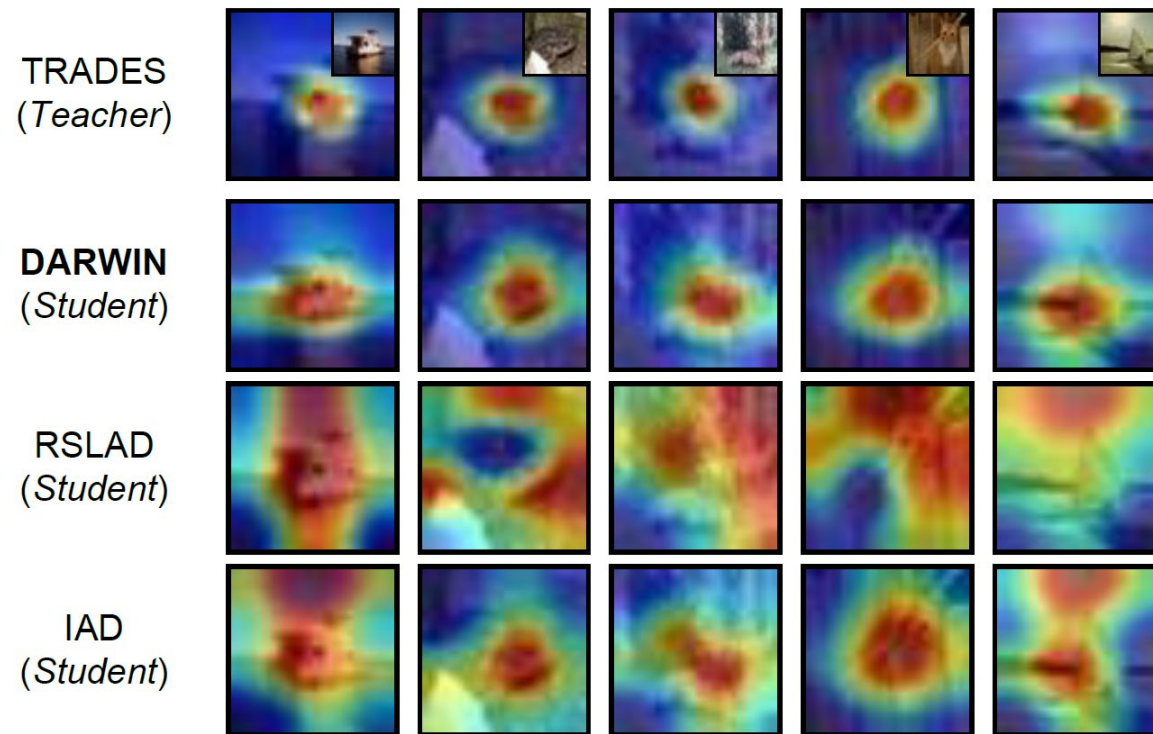
### Impact of Each Module

	ARKD	IAKD	DBKD	Natural	PGD-20	AA
1	✓			83.09/57.34	53.95/30.59	50.66/25.32
2	✓	✓		82.74/56.68	54.52/32.11	51.70/26.95
3	✓		✓	84.68/59.23	54.36/31.53	51.21/26.18
	✓	✓	✓	84.48/59.12	55.07/32.30	52.24/27.26

### Diverse Weighting Strategies

Weighting Strategies	CIFAR-10			CIFAR-100		
	Natural	PGD	AA	Natural	PGD	AA
Uniform Weighting (no weights)	83.32	53.84	50.92	57.81	30.95	25.70
First Term switched on ( $\gamma=0$ )	84.13	54.18	51.58	58.53	31.44	26.20
Second Term switched on ( $\gamma=1$ )	83.85	54.45	51.82	58.16	31.82	26.47
Both Terms switched on ( $\gamma=0.5$ )	<b>84.48</b>	<b>55.07</b>	<b>52.24</b>	<b>59.12</b>	<b>32.30</b>	<b>27.26</b>

## Attention Visualizations:



## ■ Contributions:

- We propose a novel robust knowledge distillation method that **integrates intermediate adversaries along the adversarial path**. An **adaptive weighting mechanism** is proposed to calibrate the influence of each intermediate sample to facilitate the distillation of adversarial paths. Our strategy also leads to **minimizing an upper bound of the adversarially robust risk**.
- To capture relations between decision boundaries, we devise a **dual-branch mechanism** by harnessing the **complementary characteristics of untargeted and targeted adversarial samples**. This inter-class relational learning facilitates a more effective robustness transfer.
- Extensive experiments showcase the superiority of our method compared with the state-of-the-art approaches across various settings, including diverse backbones, auxiliary data, and cross-dataset distillation.



Australian  
National  
University



# Thank you!

**Robust Distillation via Untargeted and Targeted  
Intermediate Adversarial Samples**

Junhao Dong, Piotr Koniusz, Junxi Chen, Z. Jane Wang, Yew-Soon Ong

Reporter: Junhao Dong