



# Revisiting Adversarial Robustness Distillation from the Perspective of Robust Fairness

Xinli Yue   Ningping Mou   Qian Wang   Lingchen Zhao

Key Laboratory of Aerospace Information Security and Trusted Computing, Ministry  
of Education, School of Cyber Science and Engineering, Wuhan University

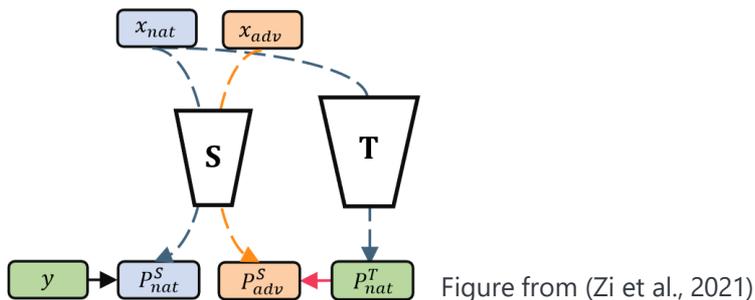
# Overview

- We pioneer research on robust fairness in adversarial robustness distillation (ARD), revealing incomplete fairness inheritance from teacher to student models and uncovering the internal reasons.
- We develop Fair-ARD, a framework enhancing robust fairness in knowledge transfer, adaptable to existing ARD techniques.
- Extensive experiments show Fair-ARD's superiority in improving student models' robust fairness and overall robustness.

# Adversarial Robustness Distillation (ARD)

ARD aims to transfer the robustness of large teacher models to small student models, facilitating the attainment of robust performance on resource-limited devices.

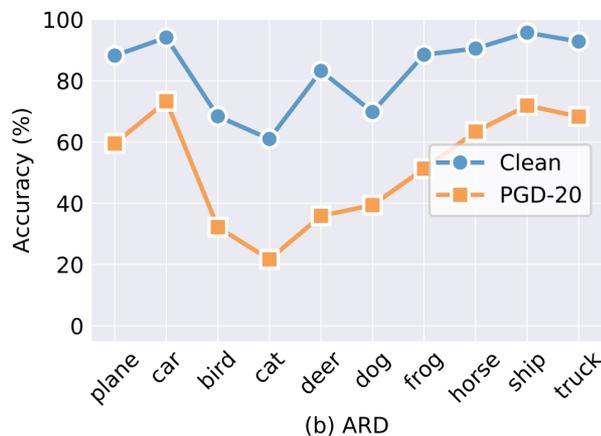
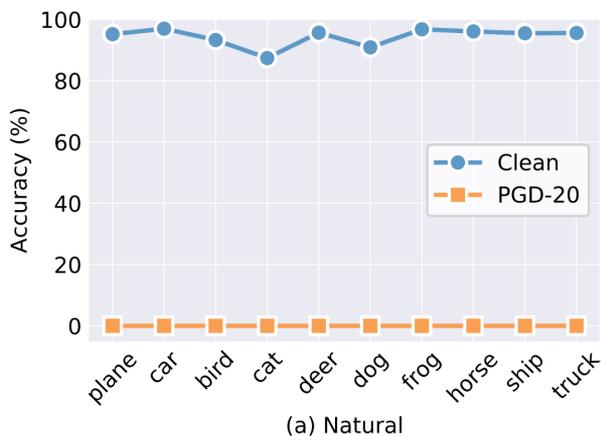
$$\mathcal{L}_{\text{ARD}} = (1 - \alpha)\text{CE}(S^\tau(x_i^j), y_i) + \alpha\tau^2\text{KL}(S^\tau(\tilde{x}_i^j), T^\tau(x_i^j))$$



- Goldblum et al., Adversarially robust distillation, AAAI 2020
- Zi et al., Revisiting Adversarial Robustness Distillation: Robust Soft Labels Make Student Better, ICCV 2021

# Robust Fairness

Models trained by ARD may exhibit high robustness on some classes while demonstrating significantly low robustness on others.



# Problem Statement

## **The Challenge in ARD**

- Traditional ARD often overlooks robust fairness, which often leads to uneven robustness in student models.

## **Research Focus**

- Investigating whether ARD can effectively transfer robust fairness from large teacher models to smaller student models.
- Addressing the gaps in class-wise robustness.

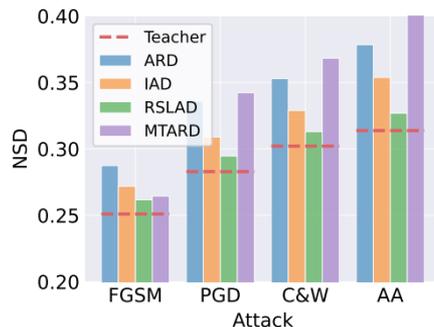
## **Objective**

- To propose a method that addresses the robust fairness issue in ARD.

# Findings on Robust Fairness Inheritance

## Inheritance of Robust Fairness

- The robust fairness of the student models is worse than the teacher model.
- Student models can only partially inherit the teacher model's robust fairness.

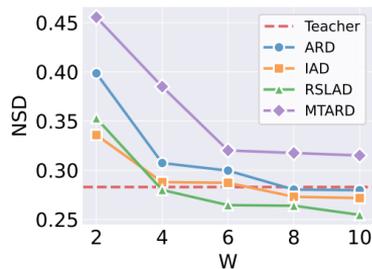


(a) ResNet18

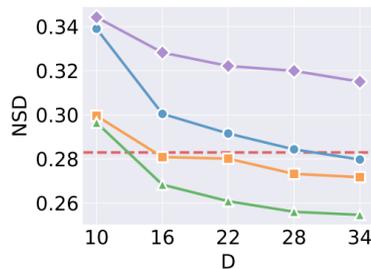
# Findings on Robust Fairness Inheritance

## Impact of Model Capacity Gap

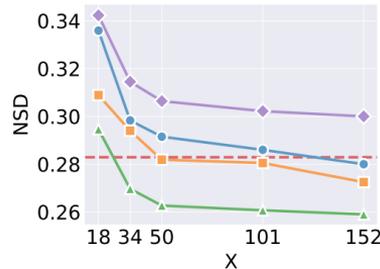
- The gap in capacity between teacher and student models affects the inheritance of robust fairness.



(b) WRN-34-W



(c) WRN-D-10



(d) ResNet-X

# Findings on Robust Fairness Inheritance

## Disparity in Class Robustness Inheritance

- Different classes inherit different proportions of robustness from the teacher model.
- We should assign larger (smaller) weights to harder (easier) classes, allowing student models to focus more on the teacher model's knowledge about hard classes.

Method	FGSM		PGD		C&W		AA	
	cat	car	cat	car	cat	car	cat	car
ARD	85.15%	<b>94.95%</b>	80.60%	<b>92.78%</b>	80.74%	<b>93.59%</b>	72.53%	<b>93.13%</b>
IAD	84.55%	<b>96.39%</b>	81.72%	<b>94.42%</b>	86.89%	<b>95.26%</b>	79.83%	<b>94.69%</b>
RSLAD	90.30%	<b>95.07%</b>	88.81%	<b>93.79%</b>	92.62%	<b>94.36%</b>	89.70%	<b>94.56%</b>
MTARD	89.70%	<b>99.28%</b>	75.37%	<b>96.70%</b>	71.31%	<b>96.41%</b>	63.52%	<b>96.11%</b>

# Fair-ARD

## Overview of Fair-ARD

- A novel framework to address the issue of partial robust fairness inheritance in ARD.
- To ensure a more equitable distribution of class-wise robustness in the student model.

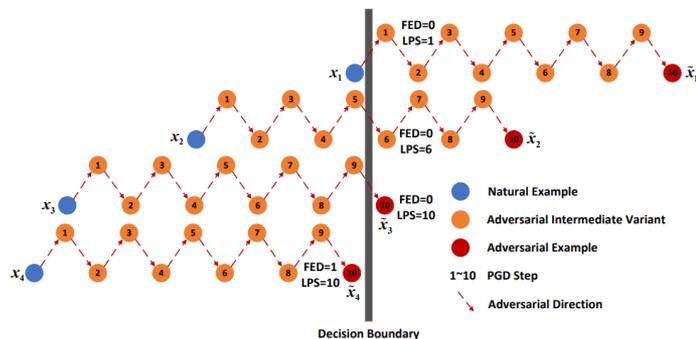
$$\min_{\theta_S} \frac{1}{C} \sum_{i=1}^C \frac{1}{n_i} \sum_{j=1}^{n_i} \omega_i \mathcal{L}_{\text{ARD}}(S, T, x_i^j, y_i, \tau, \alpha)$$

# Fair-ARD

## Metrics to Measure Class Difficulty

- Compared with FED, LPS provides a more fine-grained example difficulty measure.

$$\left\{ \begin{array}{l} \tilde{x}^{(t+1)} = \Pi_{\mathcal{B}_{\epsilon}[x]} \left( \tilde{x}^{(t)} + \gamma \text{sign} \left( \nabla_{\tilde{x}^{(t)}} \ell(S(\tilde{x}^{(t)}), y) \right) \right) \\ d(x, y) = \underset{t \in [0, K]}{\text{argmin}} (S(\tilde{x}^{(t)}) \neq y), \end{array} \right.$$



# Fair-ARD

## The Re-weighting Strategy

- Following the principle of giving larger (smaller) weights to harder (easier) classes, the weight  $\omega_i$  should decrease w.r.t.  $\kappa_i$ .

$$\kappa_i = \frac{1}{n_i} \sum_{j=1}^{n_i} d(x_i^j, y_i)$$
$$\omega_i = \frac{1}{\kappa_i^\beta}$$

# Evaluations

## Worst-class Robustness

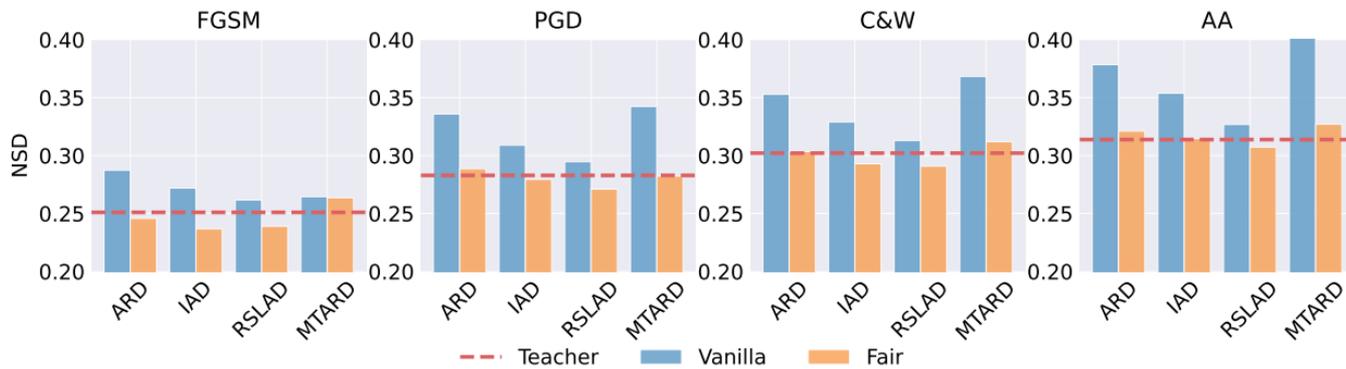
- Fair-ARD shows significant improvement in worst-class robustness under various adversarial attacks. This demonstrates the effectiveness of Fair-ARD in improving robust fairness.

Method	Clean		FGSM		PGD		C&W		AA	
	Avg.	Worst								
Natural	94.35	87.40	16.13	5.30	0.00	0.00	0.00	0.00	0.00	0.00
SAT	84.27	64.10	56.81	28.30	49.11	22.10	48.58	20.50	46.13	17.40
TRADES	82.22	64.80	58.38	31.00	52.35	25.80	50.33	22.30	49.01	21.20
ARD	<b>83.22</b>	61.00	<b>58.77</b>	28.10	51.65	21.60	<b>51.25</b>	19.70	49.05	16.90
Fair-ARD (ours)	82.96	<b>68.10</b>	57.69	<b>39.20</b>	<b>52.05</b>	<b>33.20</b>	50.69	<b>31.00</b>	<b>49.13</b>	<b>29.20</b>
IAD	<b>83.25</b>	60.60	<b>58.90</b>	27.90	52.08	21.90	51.01	21.20	48.95	18.60
Fair-IAD (ours)	83.19	<b>68.70</b>	58.31	<b>36.50</b>	<b>52.27</b>	<b>30.00</b>	<b>51.16</b>	<b>28.10</b>	<b>49.28</b>	<b>25.60</b>
RSLAD	83.04	62.70	60.03	29.80	54.13	23.80	52.76	22.60	51.18	20.90
Fair-RSLAD (ours)	<b>83.59</b>	<b>67.80</b>	<b>60.16</b>	<b>35.50</b>	<b>54.33</b>	<b>29.70</b>	<b>53.07</b>	<b>26.90</b>	<b>51.23</b>	<b>25.10</b>
MTARD	<b>87.14</b>	<b>69.80</b>	<b>60.62</b>	30.10	50.81	20.70	48.85	18.00	46.10	16.10
Fair-MTARD (ours)	81.98	62.10	59.11	<b>30.60</b>	<b>53.96</b>	<b>27.50</b>	<b>52.32</b>	<b>24.20</b>	<b>50.60</b>	<b>22.60</b>

# Evaluations

## NSD

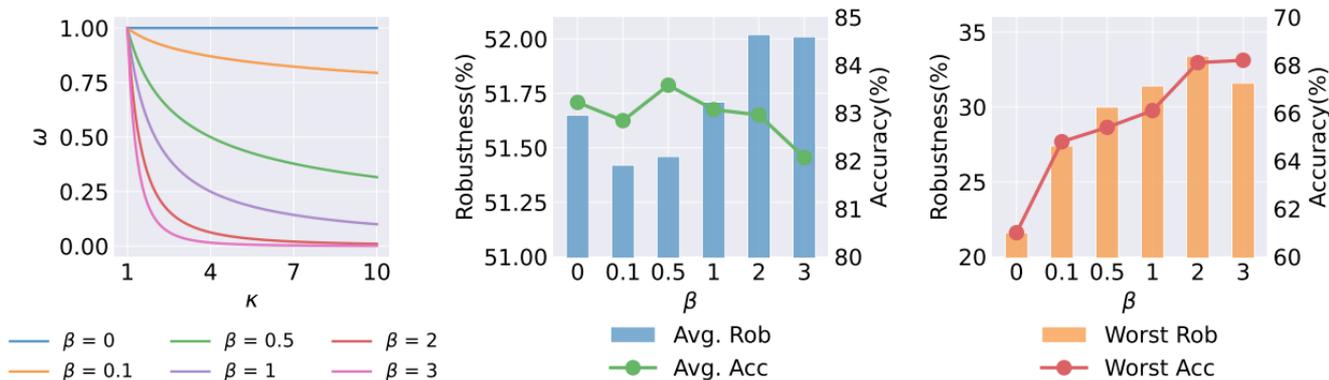
- Fair-ARD achieves lower NSD values, showing its effectiveness in balancing class-wise robustness.



# Ablation Studies

## The Effect of the Hyperparameter $\beta$

- Fair-ARD maintains its effectiveness across a range of  $\beta$ , highlighting its adaptability.



# Discussion

## Comparison with GAIRAT and Robust Fairness Algorithms

- Fair-ARD outperforms others in achieving better worst-class robustness while maintains comparable overall robustness.

Method	Clean		FGSM		PGD		C&W		AA	
	Avg.	Worst								
ARD	<b>83.22</b>	61.00	<b>58.77</b>	28.10	51.65	21.60	<b>51.25</b>	19.70	49.05	16.90
GAIR-ARD	81.09	<u>65.80</u>	54.04	34.60	50.00	<u>32.60</u>	42.69	20.70	40.88	18.90
FRL-ARD	81.47	61.40	56.90	34.30	50.24	28.60	49.76	26.80	47.84	23.90
FAT-ARD	<u>82.98</u>	65.20	57.62	<u>38.00</u>	51.11	31.00	49.99	<u>29.20</u>	48.22	<u>25.90</u>
CFA-ARD	<u>82.00</u>	62.30	57.68	<u>30.60</u>	<b>52.43</b>	25.00	<u>51.23</u>	<u>22.30</u>	<b>49.73</b>	<u>20.00</u>
Fair-ARD (ours)	82.96	<b>68.10</b>	<u>57.69</u>	<b>39.20</b>	<u>52.05</u>	<b>33.20</b>	<u>50.69</u>	<b>31.00</b>	<u>49.13</u>	<b>29.20</b>

# Conclusion

- We explore the issue of robust fairness in ARD, revealing the partial inheritance of robust fairness by student models from teacher models.
- We develop Fair-ARD, a novel framework that employs a refined class difficulty metric and a re-weighting strategy.
- We show that Fair-ARD significantly enhances the robust fairness of student models while maintaining high overall robustness, paving the way for fairer lightweight deep learning models.

# Thanks!



[lczhaocs@whu.edu.cn](mailto:lczhaocs@whu.edu.cn)



<https://github.com/NISP-official/Fair-ARD>