# Enhancing Adversarial Robustness for Deep Metric Learning

Mo Zhou [1]     Vishal M. Patel [1]

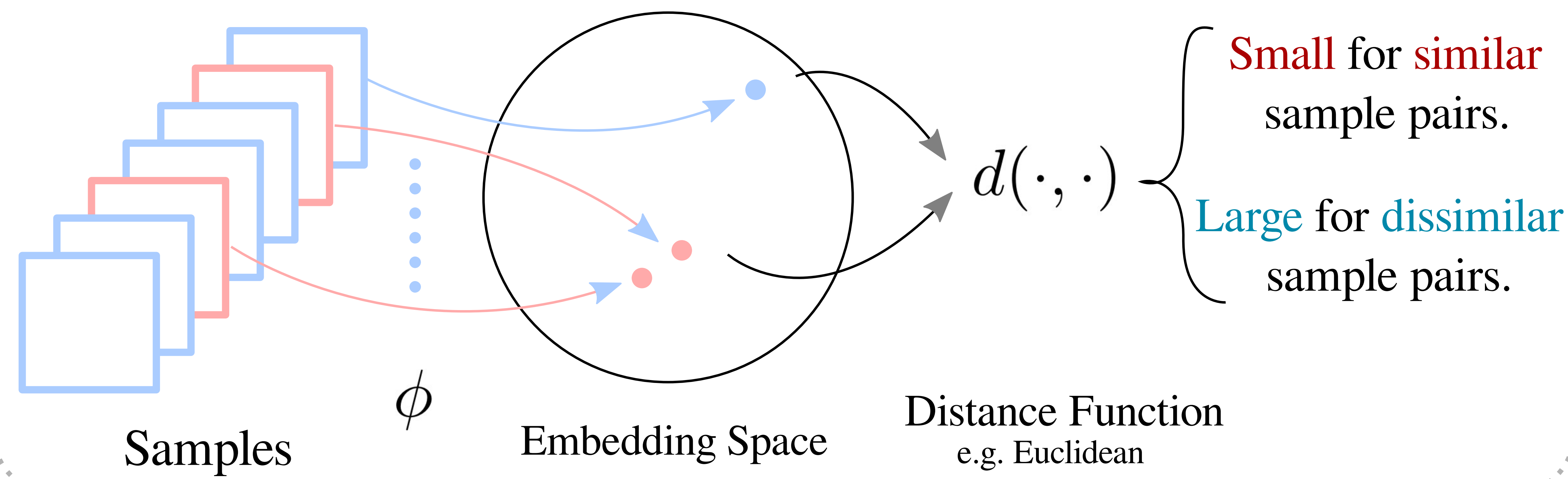[1] Johns Hopkins University

arXiv

Github

CVPR JUNE 19-24 2022 NEW ORLEANS · LOUISIANA

## Deep Metric Learning

Goal: Learn a mapping function $\phi : \mathcal{X} \mapsto \Phi \subseteq \mathbb{R}^D$



Small for similar sample pairs.
Large for dissimilar sample pairs.

Samples

$\phi$

Embedding Space

Distance Function
e.g. Euclidean

## Insights / Motivations

- Deep metric learning (DML) model is vulnerable to adversarial attacks, which results in unexpected retrieval results.

- Leads to safety and security concerns to DML applications.

Embedding-Shifted Triplet (EST) [9] adopts adversarial counterparts of $a, p, n$ with maximum embedding move distance off their original locations, i.e.,

$$L_{EST} = L_T(\tilde{a}, \tilde{p}, \tilde{n}; \gamma) \qquad (1)$$

where $\tilde{a} = \phi(A + r^*)$, and $r^* = \arg\max_r d_\phi(A + r, A)$. The $\tilde{p}$ and $\tilde{n}$ are obtained similarly.

Anti-Collapse Triplet (ACT) [10] "collapses" the embedding vectors of positive and negative sample, and enforces the model to separate them apart, i.e.,

$$L_{ACT} = L_T(a, \overrightarrow{p}, \overleftarrow{n}; \gamma), \qquad (2)$$
$$[\overrightarrow{p}, \overleftarrow{n}] = [\phi(P + r_p^*), \phi(N + r_n^*)] \qquad (3)$$
$$[r_p^*, r_n^*] = \arg\min_{r_p, r_n} d_\phi(P + r_p, N + r_n). \qquad (4)$$
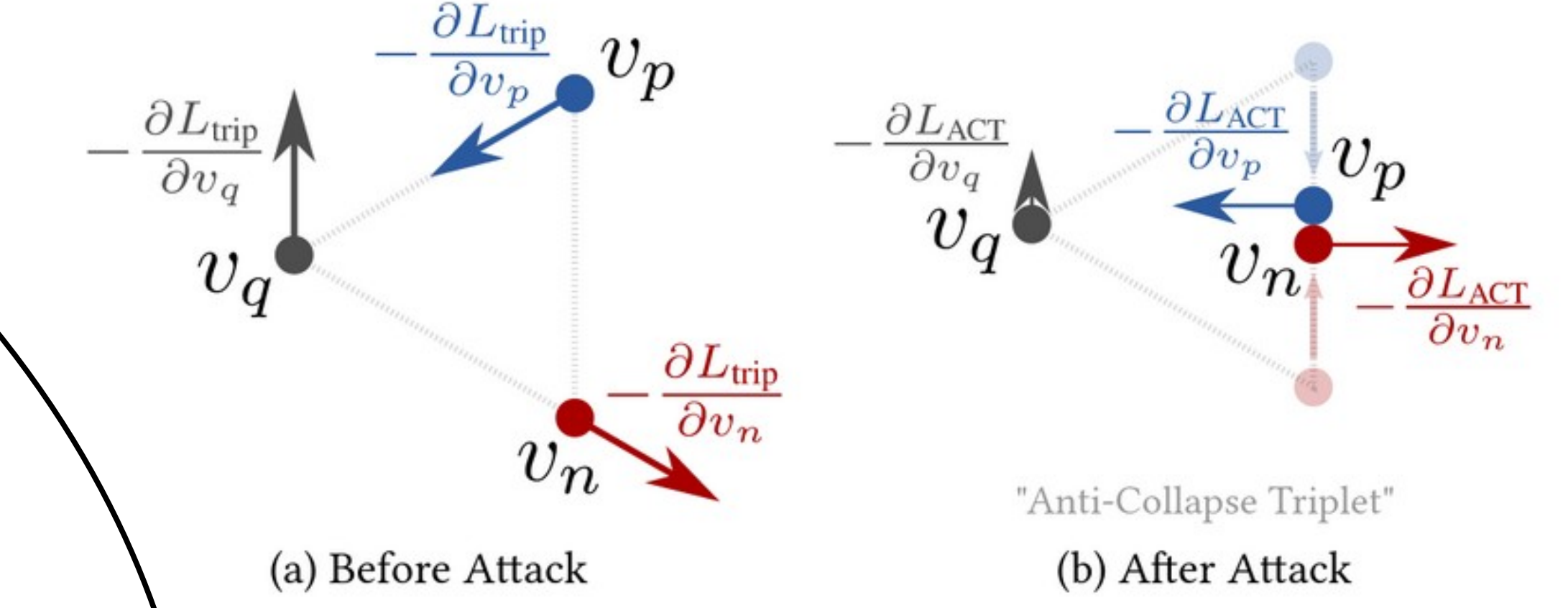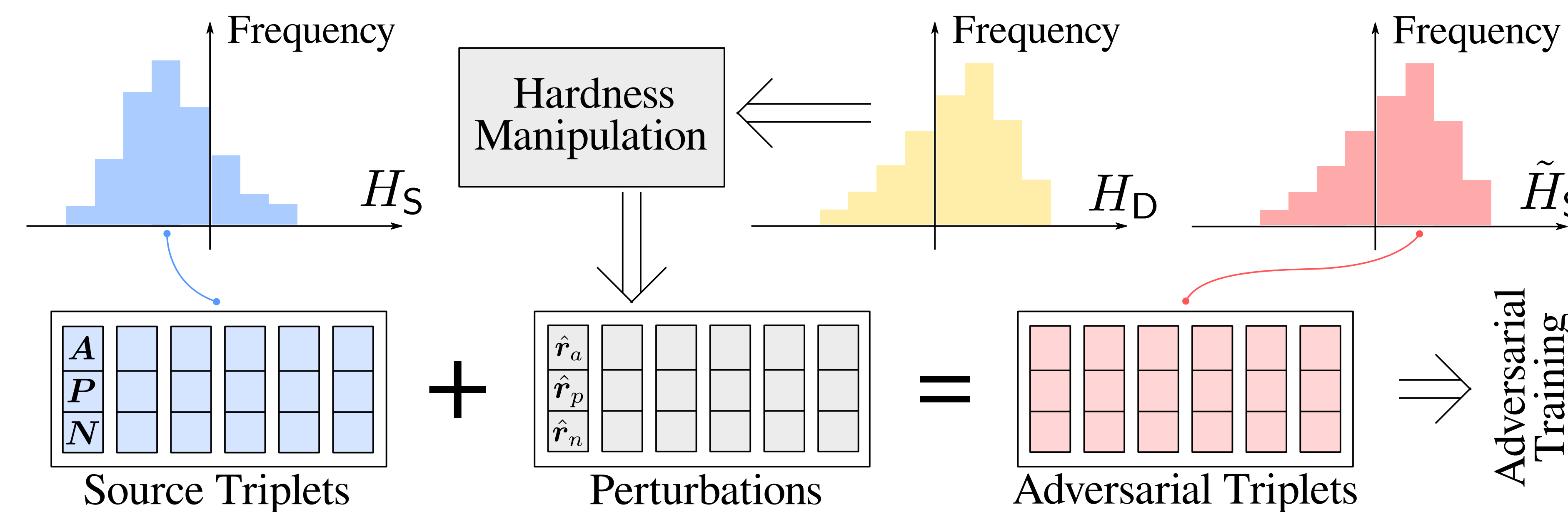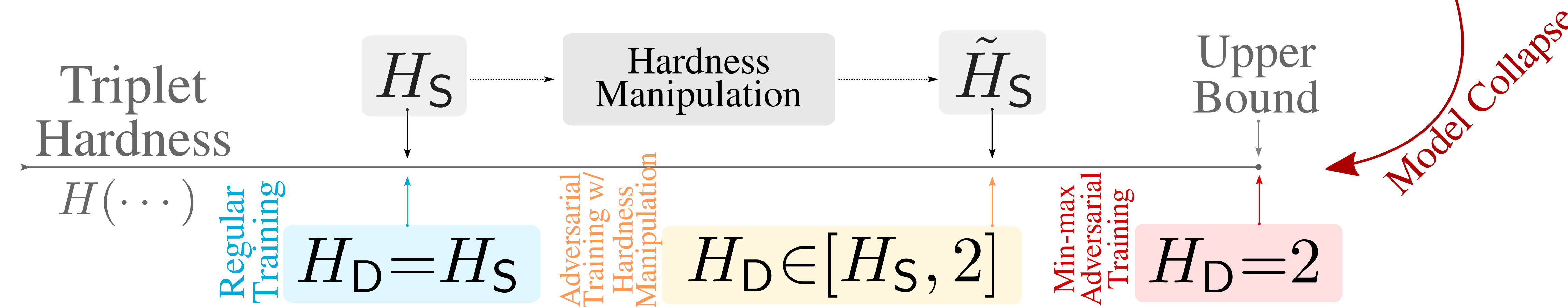


(a) Before Attack   (b) After Attack

Fig. 4: Gradient Direction of Anti-Collapse Triplet (ACT). ACT does not suffer from misleading gradient or inefficient mini-batch exploitation. It does not create excessively hard adversarial example neither.

- Existing defense methods (i.e., EST and ACT) learn from relatively weak adversaries in order to avoid model collapse.

- And suffer from low efficiency in adversarial training, low performance on benign examples, and insufficient robustness.
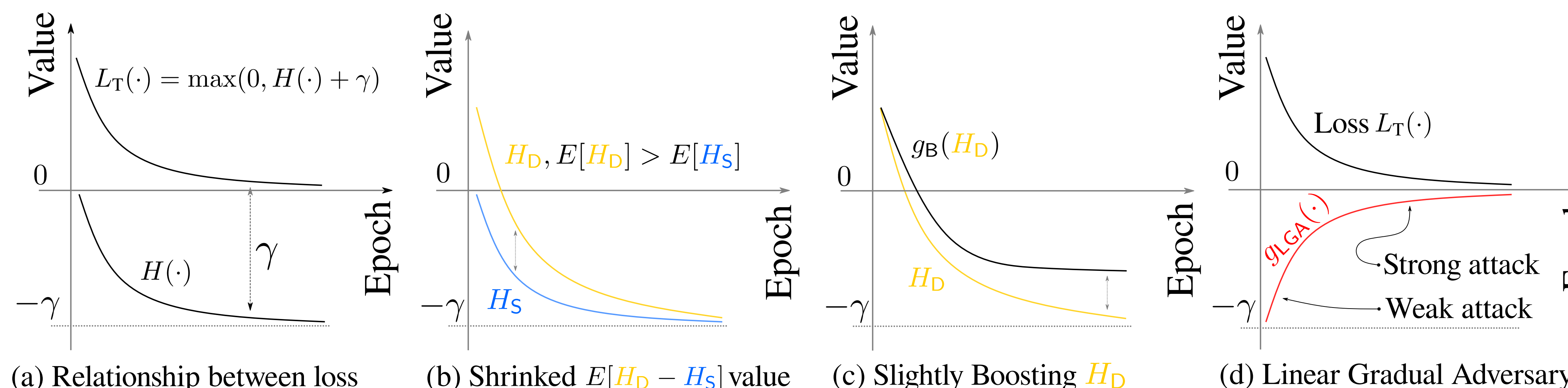
## Contribution 1
## Hardness Manipulation



- Perturb a given image triplet to a specific hardness level, e.g., from semi-hard to soft-hard, instead of as hard as possible.



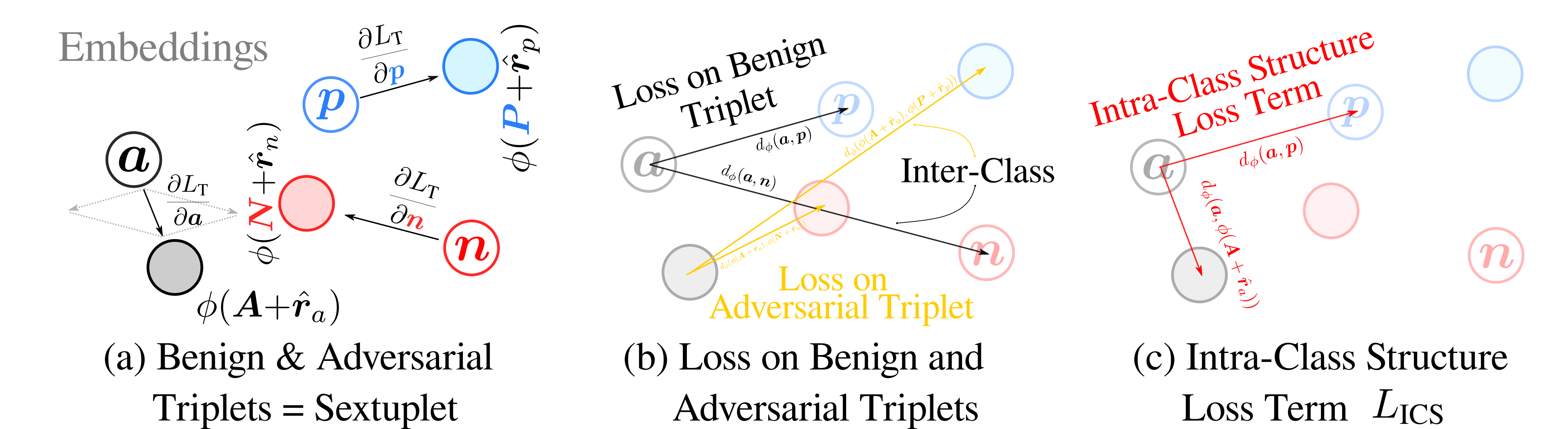| Triplet Hardness $H(\cdots)$ | $H_S$ → Hardness Manipulation → $\tilde{H}_S$ | Upper Bound |
| Regular Training $H_D = H_S$ | Adversarial Training w/ Hardness Manipulation $H_D \in [H_S, 2]$ | Min-max Adversarial Training $H_D = 2$ |

Model Collapse

- Can be seen as a flexible and dynamic "interpolation" between regular training and Madry min-max adversarial training.

- Efficiently create adversarial examples for training and do not easily lead to model collapse.

## Contribution 2
## Linear Gradual Adversary



(a) Relationship between loss and hardness when $L_T(\cdot) > 0$
(b) Shrinked $E[H_D - H_S]$ value leads to weak attack
(c) Slightly Boosting $H_D$ for slightly stronger attack
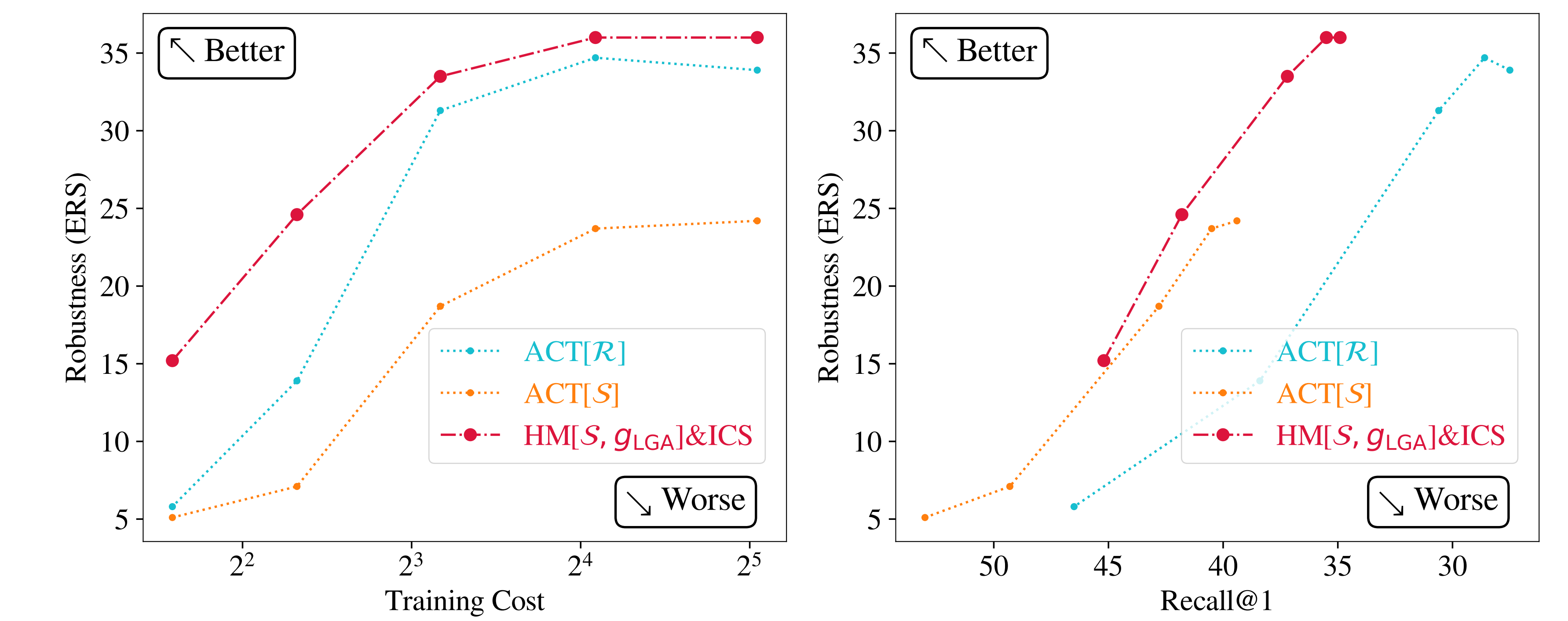(d) Linear Gradual Adversary for balanced training

- Gradually and dynamically increases hardness level to balance the model robustness and performance on benign examples.

## Contribution 3
## Intra-Class Structure Loss



(a) Benign & Adversarial Triplets = Sextuplet
(b) Loss on Benign and Adversarial Triplets
(c) Intra-Class Structure Loss Term $L_{ICS}$

- Enforcing intra-class structure is neglected by existing defenses.
- Some attacks happen within a class.

## Results



- Overwhelmingly outperform SoTA in efficiency of adversarial training, performance on benign examples, and robustness.

- Comparison on:
  (1) CUB-200-2011
  (2) Cars-196
  (2) Stanford Online Product



This table is vector graphics cropped from my PDF instead of bitmap screenshot.
Zoom in and see the numbers.

$\eta$: number of PGD iterations for adversarial example, lower=better
R@1: recall @ 1 metric for deep metric learning, higher=better
ERS: empirical robustness score, higher=better

Table 6. Comparison of our defense with the state-of-the-art methods on commonly used DML datasets.