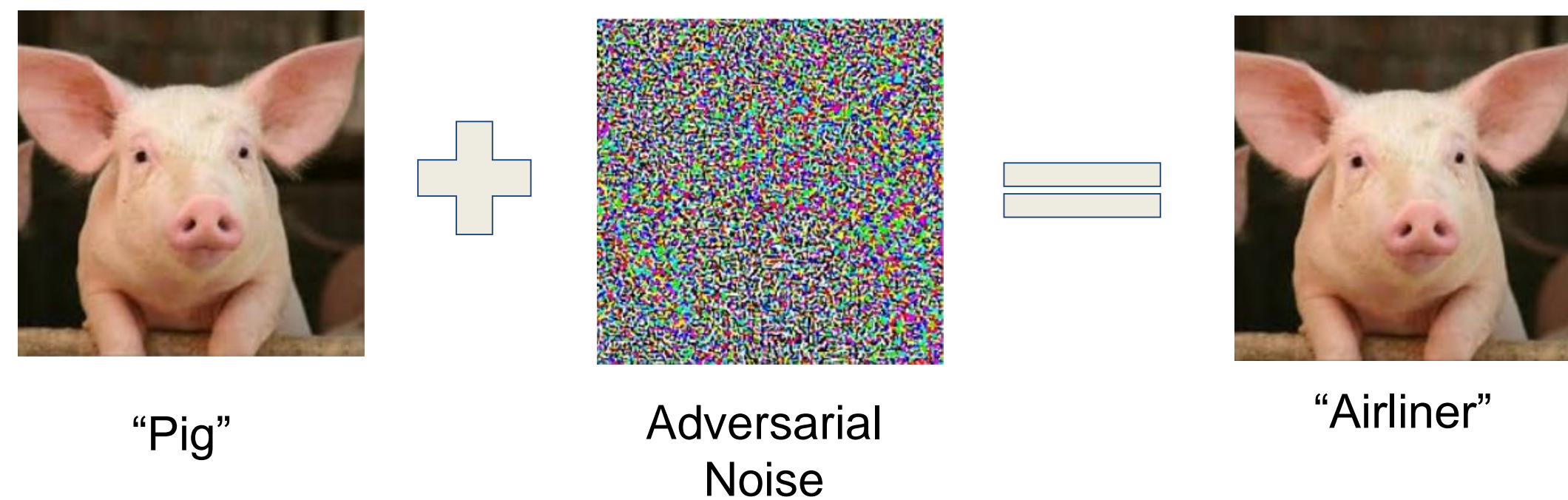


Metric Learning for Adversarial Robustness



Adversarial Examples

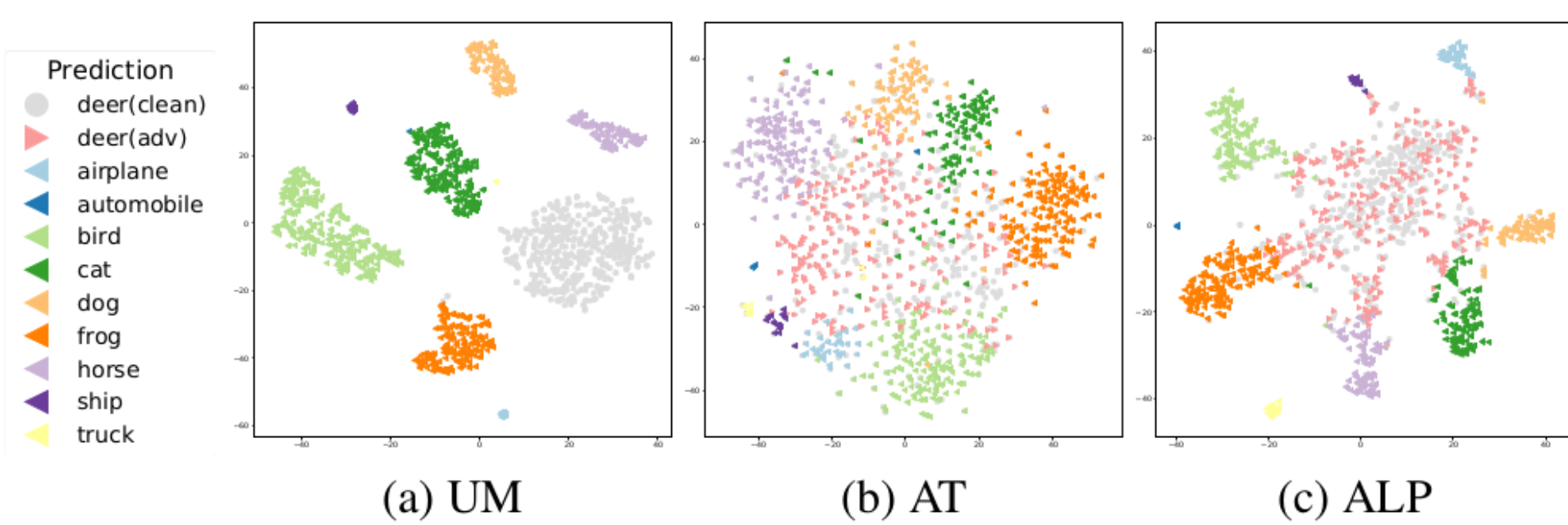
Neural networks are vulnerable under adversarial attacks.



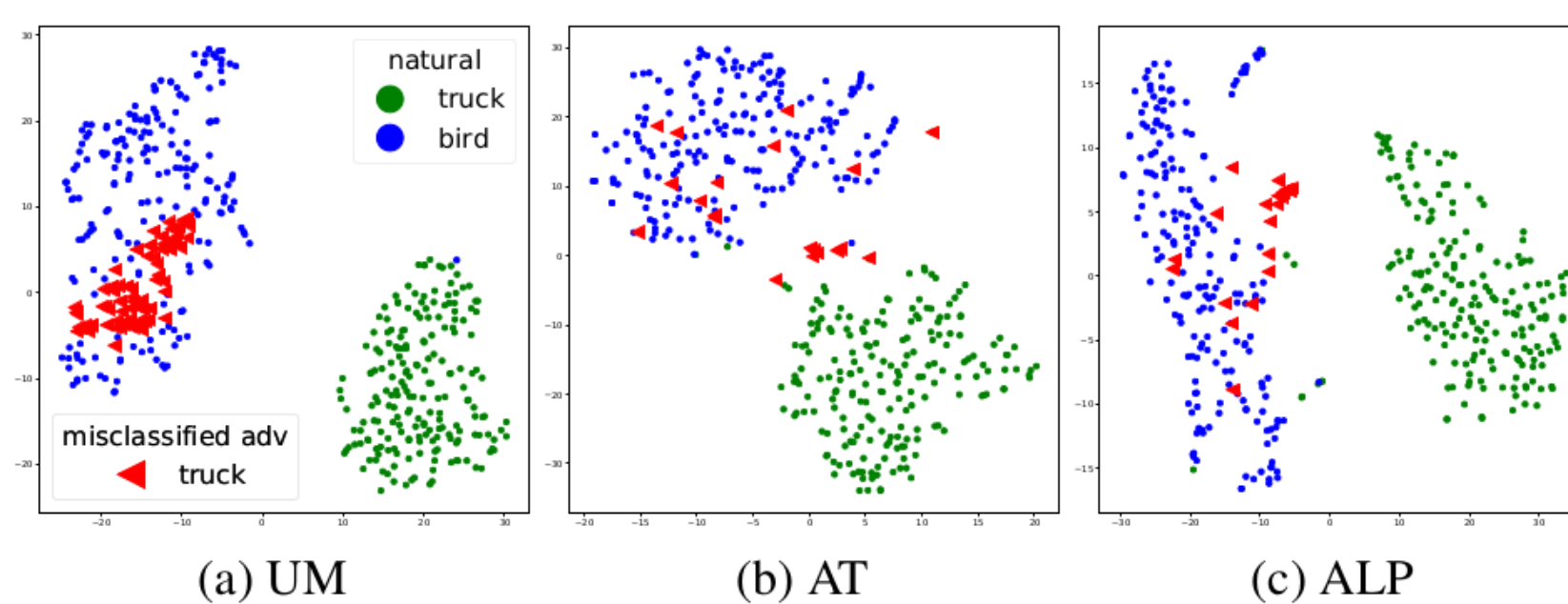
Design Motivations

Hidden representations of DNNs under attack have not been systematically studied.

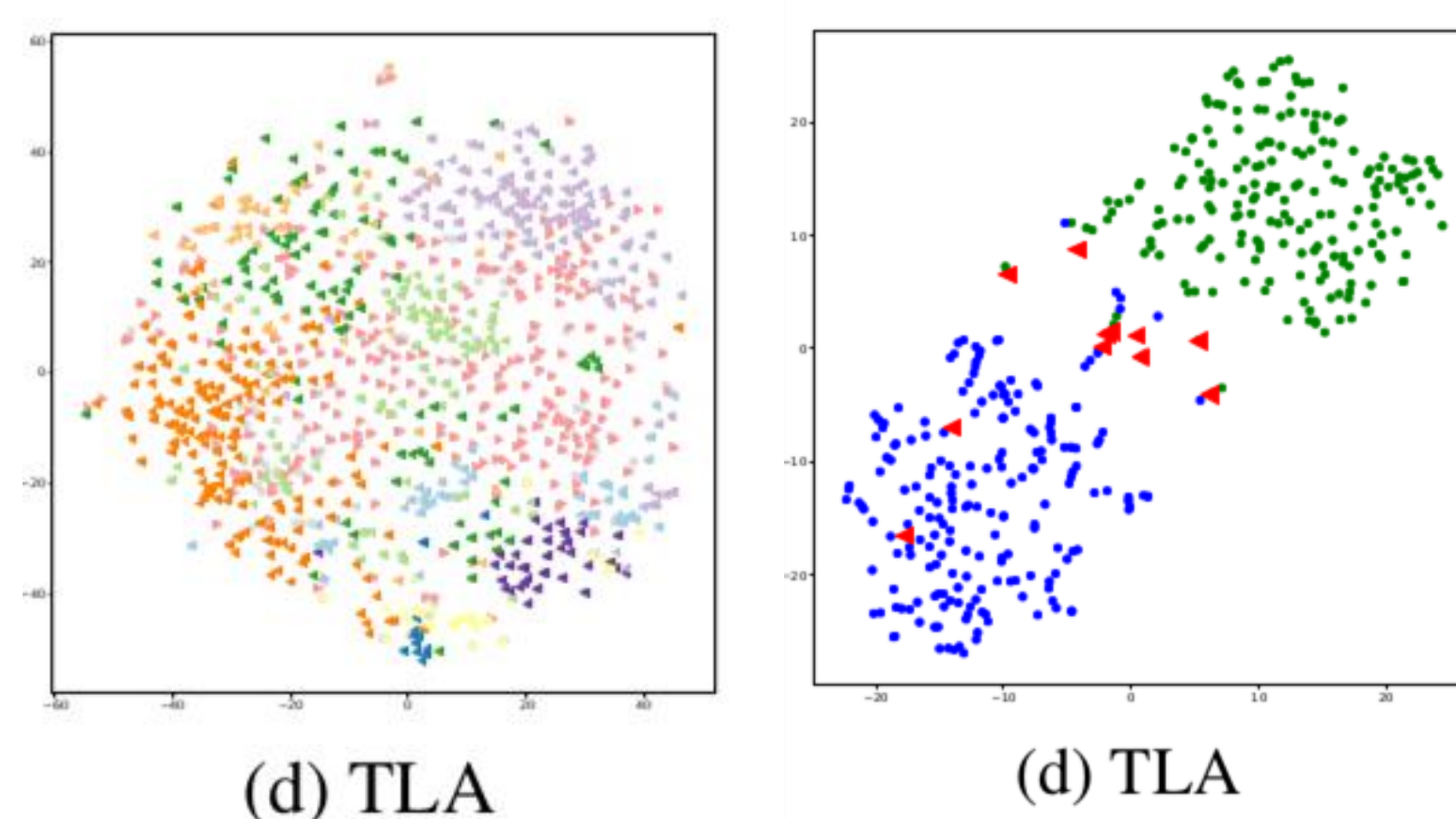
- Representations from the same class are mistakenly classified to different classes.



- Representations of the adversarial examples are shifted inside the false class.



After TLA training



- The examples from the same class are clustered together.
- The adversarial examples are separated from the corresponding false class by a margin.

Methodology



Triplet Loss

- Pull the representations of the same class together
- Push the representations of different classes away up to a margin
- Anchor $x_a^{(i)}$ and positive $x_p^{(i)}$ are from the same class, negative $x_n^{(i)}$ is from a different class

$$\sum_i^N \mathcal{L}_{trip}(x_a^{(i)}, x_p^{(i)}, x_n^{(i)}) = \sum_i^N [D(h(x_a^{(i)}), h(x_p^{(i)})) - D(h(x_a^{(i)}), h(x_n^{(i)})) + \alpha]_+$$

Triplet Loss for Adversarial Robustness (TLA)

$$\mathcal{L}_{all} = \sum_i^N \mathcal{L}_{ce}(f(x_a^{(i)}), y^{(i)}) + \lambda_1 \mathcal{L}_{trip}(h(x_a^{(i)}), h(x_p^{(i)}), h(x_n^{(i)})) + \lambda_2 \mathcal{L}_{norm}$$

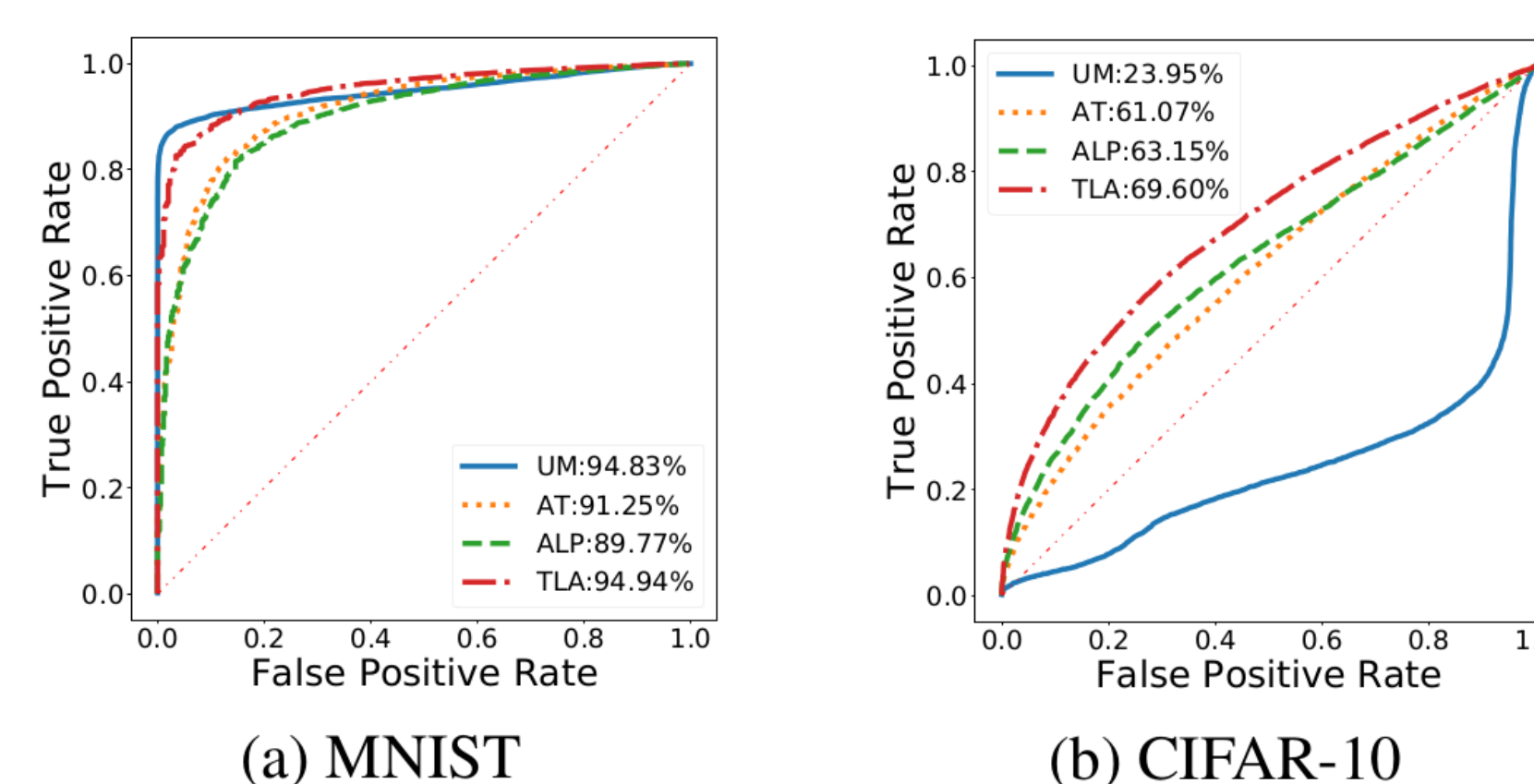
Semi-hard negative example retrieval

- Negative examples randomly sampled from another class are too easy.
- Retrieving the closest examples from the total training set is slow and too hard as a negative example.
- We propose to retrieve the closest negative example from a mini-batch of training data.

- We propose to regularize the distorted representation under attack with metric learning, to produce more robust classifier.
- Our method, TLA, produces desirable internal representation that improves adversarial robust accuracy and misclassification detection efficiency.

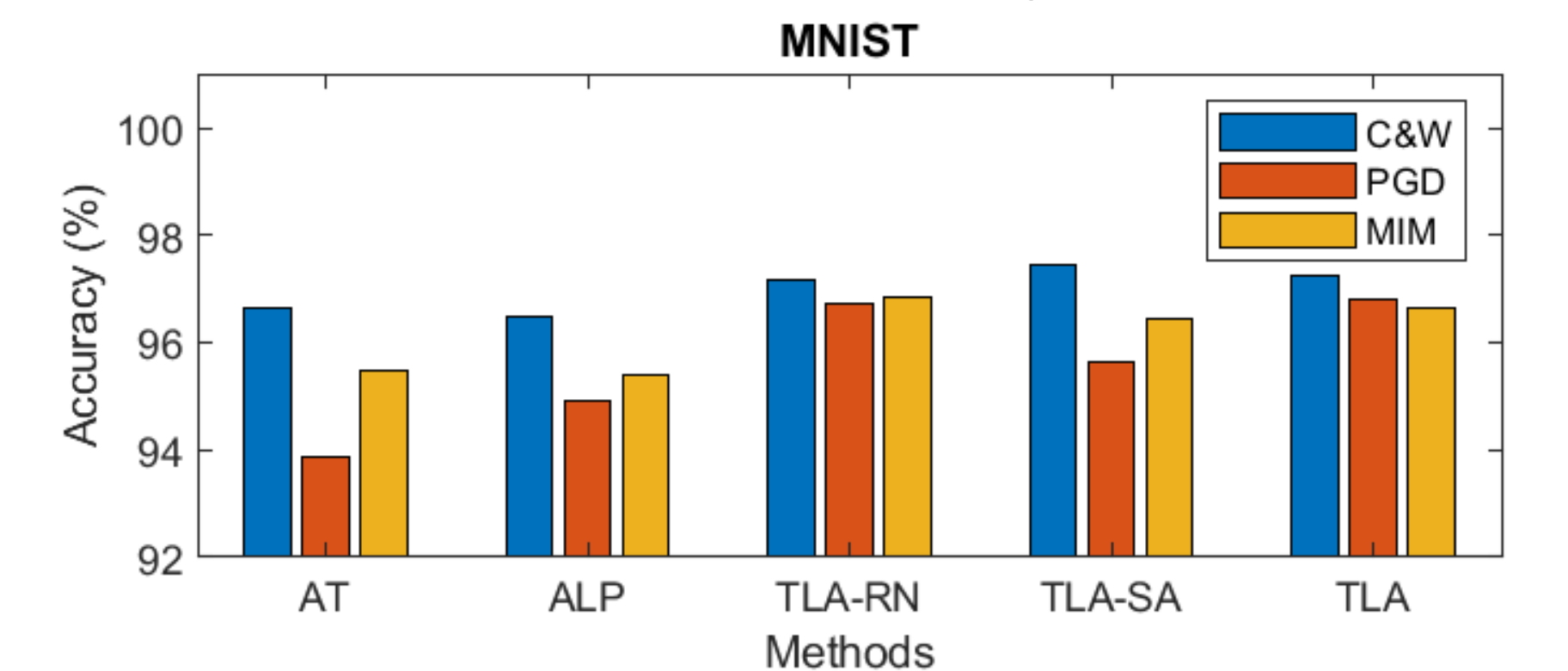
Results

Effect of TLA on mis-classification detection: ROC Curve

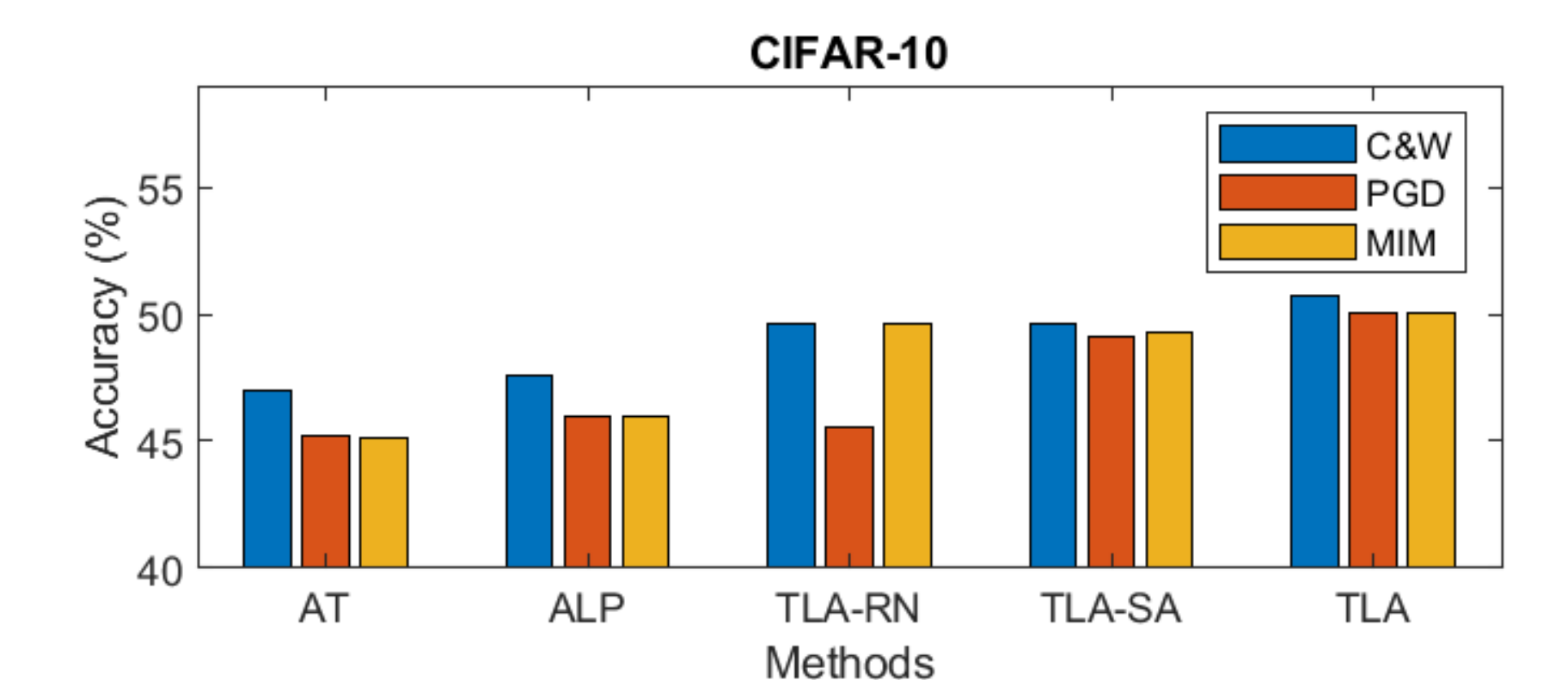


- PGD attack under larger perturbation norm bound compared with training time
- TLA Improves the AUC score for detection by
 - 3.69% on MNIST
 - 6.46% on CIFAR-10.

Effect of TLA on classification accuracy under attacks

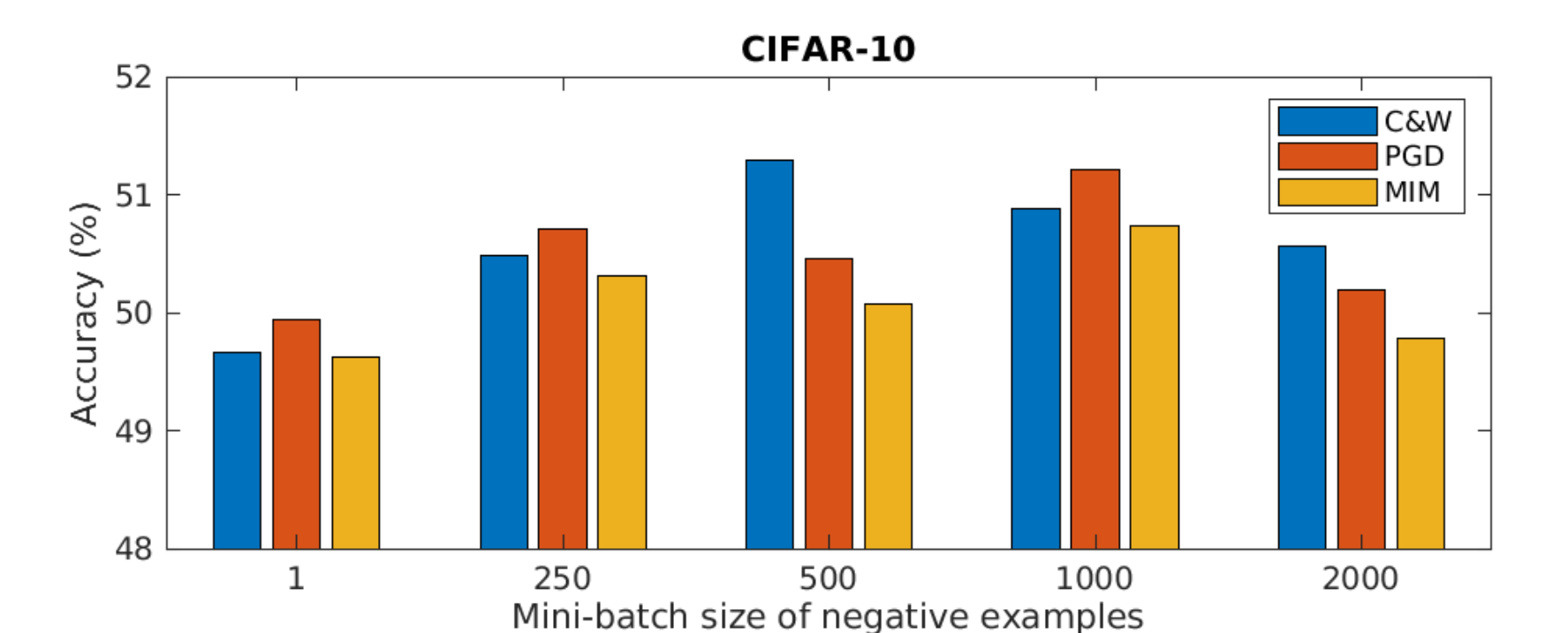


- With 2000 steps of PGD
- Improve the accuracy by
 - 1.86% compared with Adversarial Logit Pairing
 - 2.92% compared with Adversarial Training.



- With 400 steps of PGD
- Improve the accuracy by
 - 4.05% compared with Adversarial Logit Pairing
 - 4.82% compared with Adversarial Training.

Semi-hard negative example retrieval using mini-batch



- By choosing the proper mini-batch size, performance is improved by up to 1.64%

Nearest Neighbors Retrieval Using Learned Metric

