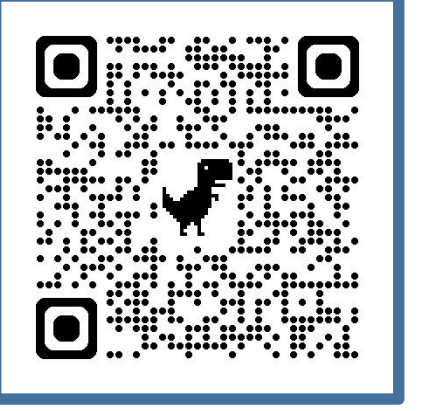# CLIP Is Strong Enough to Fight Back:
# Test-time Counterattacks towards Zero-shot Adversarial Robustness of CLIP

Songlong Xing[1]    Zhengyu Zhao[2]    Nicu Sebe[1]
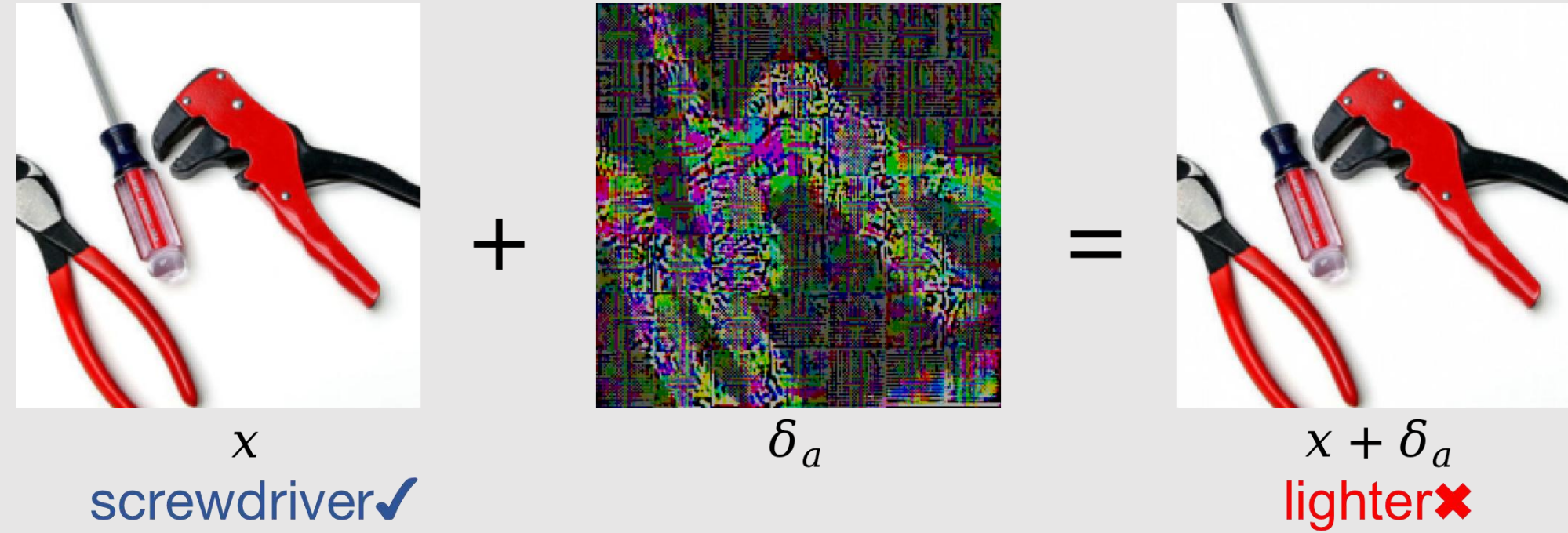
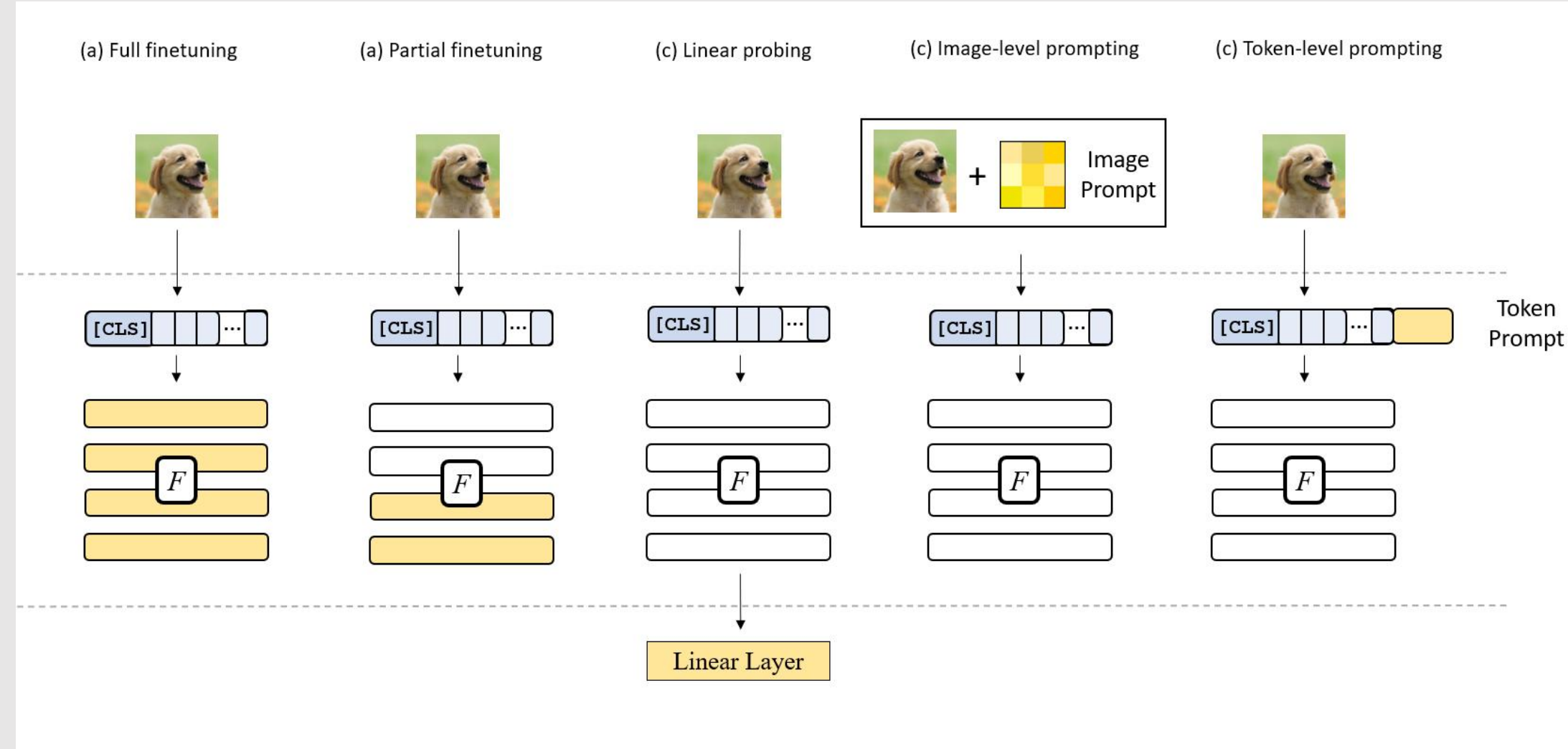[1]University of Trento, Italy    [2]Xi'an Jiaotong University, China

CVPR Nashville JUNE 11-15, 2025

paper & code!

## Motivation

### 1. CLIP is vulnerable to adversarial perturbations



$x$
screwdriver ✔    $\delta_a$    $x + \delta_a$
lighter ✘

### 2. State-of-the-art methods employ finetuning or prompt tuning with adversarial images.



(a) Full finetuning  (a) Partial finetuning  (c) Linear probing  (c) Image-level prompting  (c) Token-level prompting
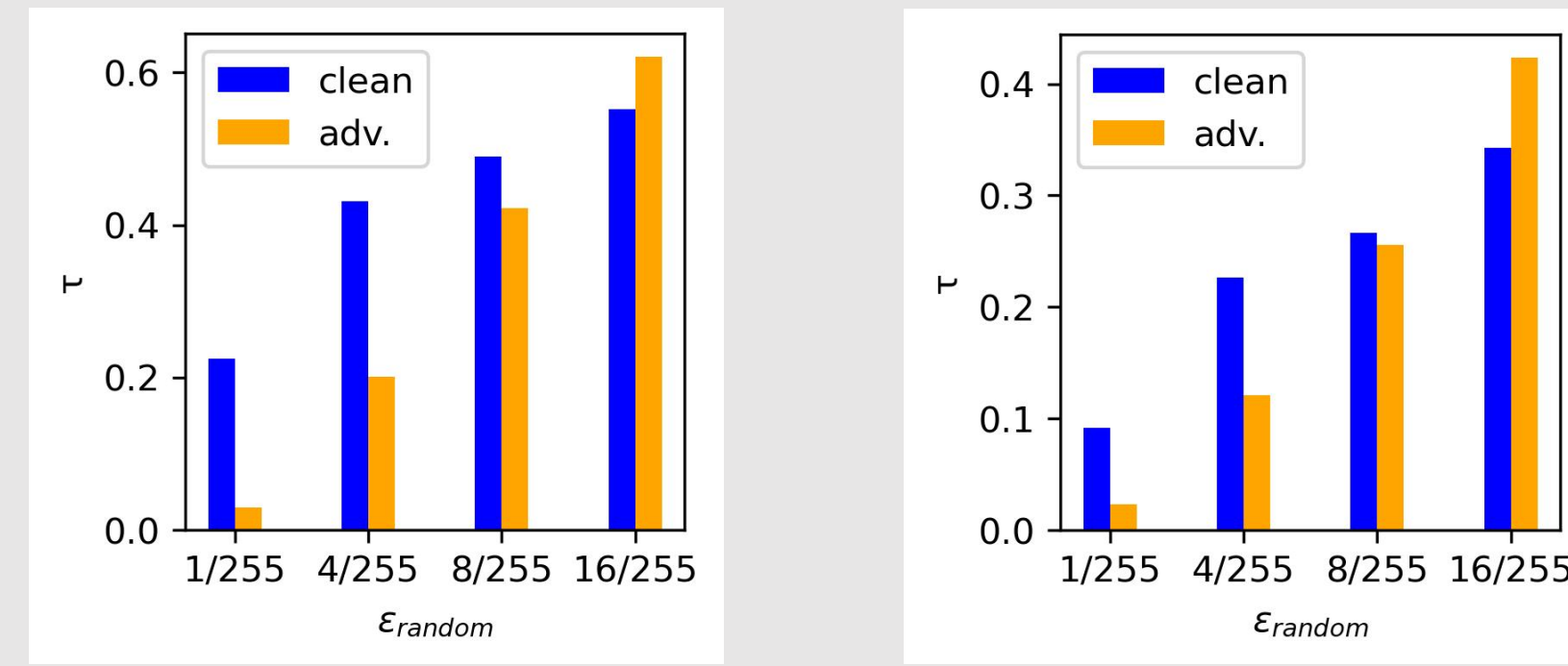
(Figure borrowed from TeCoA [1])

**Limitations are apparent:**
- Expensive training
- Overfitting to adversarial samples
- Significant loss of clean performance

**What if we discard training and counter adversary at test time?**

## Preliminary Experiment

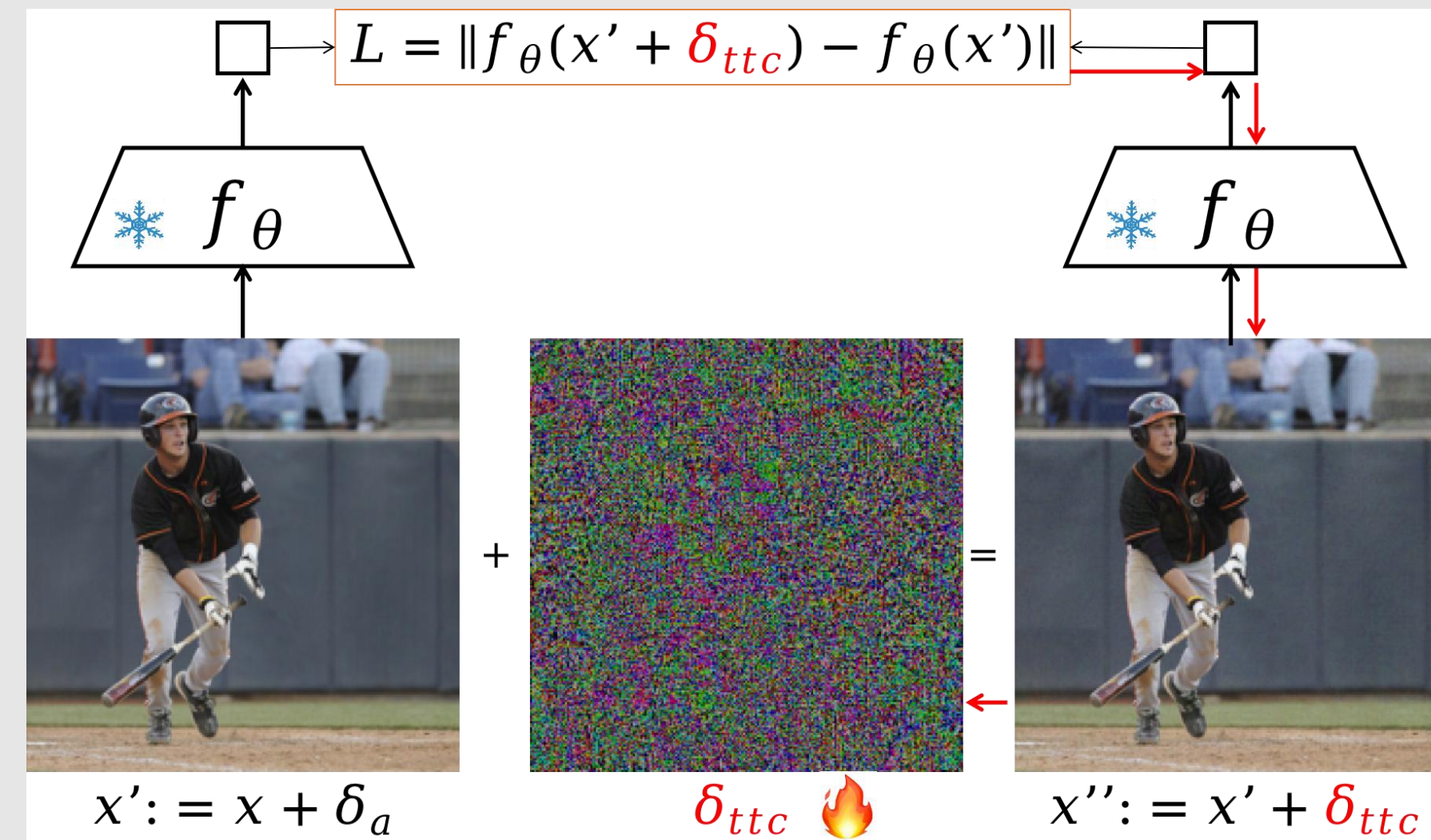**Perturbations that maximize downstream loss cause 'false robustness'.**



CIFAR100          ImageNet

(Check out Appendix for theoretical explanation)

## Methodology

**Leverage the vision encoder to counterattack adversarial images to mitigate 'false robustness'.**



$L = \|f_\theta(x' + \delta_{ttc}) - f_\theta(x')\|$

$x' := x + \delta_a$    $\delta_{ttc}$ 🔥    $x'' := x' + \delta_{ttc}$

$$\delta_{ttc} = \arg\max_\delta \|f_\theta(x + \delta) - f_\theta(x)\|$$

$$s.t. \|\delta\|_\infty \leq \varepsilon_{ttc}$$

## Algorithm

### 1. $\tau$ describes representational variation induced by random noise:

$$\tau = \frac{\|f_\theta(x + n) - f_\theta(x)\|}{f_\theta(x)}$$

### 2. Counterattack test image $x$ depending on $\tau$ at the initial step ($\tau < \tau_{thres}$):

$$\delta_{ttc}^i = \prod (\delta_{ttc}^{i-1} + \alpha \nabla_\delta \|f_\theta(x + \delta_{ttc}^{i-1}) - f_\theta(x)\|)$$

### 3. Weight and sum $\delta_{ttc}^i$ across N steps:

$$\delta_{ttc} = \sum_{i=0}^{N} w_i \bullet \delta_{ttc}^i$$

---

**Algorithm 1** $\tau$-thresholded weighted counterattacks.

**Require:** Test image $x$, pre-trained CLIP vision encoder $f_\theta$, counterattack budget $\epsilon_{ttc}$, stepsize $\alpha$, number of steps $N$, user-defined parameters $\tau_{thres}$ and $\beta$.
1: **procedure** TEST-TIME COUNTERATTACKS
2:     $\delta_{ttc}^0 \sim U(-\epsilon_{ttc}, \epsilon_{ttc})$.
3:     Compute $\tau$ based on Eq. (4) using $\delta_{ttc}^0$.
4:     **if** $\tau \geq \tau_{thres}$ **then**
5:       $w_0 = 1$
6:       **return** $\delta_{ttc} = \delta_{ttc}^0$
7:     **else if** $\tau < \tau_{thres}$ **then**
8:       $\mathbf{w}, \boldsymbol{\delta}_{ttc} := \{\}, \{\}$
9:       **for** $i = 1, 2, \ldots, N$ **do**
10:        $\delta_{ttc}^i = \Pi(\delta_{ttc}^{i-1} + \alpha \nabla_\delta \|f_\theta(x + \delta_{ttc}^{i-1}) - f_\theta(x)\|)$
11:        $w_i = \exp(\beta \cdot i) / \sum_{j=0}^{N} \exp(\beta \cdot j)$ (Eq. (5))
12:        $\mathbf{w} \leftarrow w_i$, $\boldsymbol{\delta}_{ttc} \leftarrow \delta_{ttc}^i$
13:       **end for**
14:       **return** $\delta_{ttc} = \sum_{i=0}^{N} w_i \cdot \delta_{ttc}^i$ (Eq. (6))
15:     **end if**
16: **end procedure**

## Results

| (%) | | CLIP | TeCoA [1] | FARE [2] | RN | TTC (ours) | Δ (w.r.t. CLIP) |
|---|---|---|---|---|---|---|---|
| ImageNet | Rob. | 1.15 | 18.89 | 14.00 | 1.77 | 38.41 | +37.26 |
| | Acc. | 59.69 | 34.89 | 48.79 | 59.34 | 49.39 | -10.30 |
| CIFAR10 | Rob. | 0.74 | 33.61 | 19.65 | 2.01 | 28.75 | +28.01 |
| | Acc. | 85.12 | 64.61 | 74.44 | 81.18 | 81.18 | -3.94 |
| Caltech256 | Rob. | 8.47 | 43.19 | 38.79 | 11.33 | 60.11 | +51.64 |
| | Acc. | 81.72 | 61.14 | 73.32 | 81.25 | 79.66 | -2.06 |
| Cars | Rob. | 0.02 | 8.76 | 6.75 | 0.16 | 33.01 | +32.99 |
| | Acc. | 52.02 | 20.91 | 38.68 | 52.14 | 48.16 | -3.86 |
| Avg. 16 datasets | Rob. | 2.70 | 26.54 | 20.00 | 3.86 | 39.17 | +36.47 |
| | Acc. | 61.51 | 40.25 | 51.02 | 61.61 | 59.75 | -1.76 |

Tab.1 PGD-10 attacks ($\varepsilon_a = 1/255$). Find full table in the paper.

| (%) | C10 | IN | Cal256 | Cars | Rob. | Acc. |
|---|---|---|---|---|---|---|
| TeCoA | 33.61 | 18.89 | 43.19 | 8.76 | 26.54 | 40.25 |
| TeCoA + TTC | 34.68 | 23.14 | 48.49 | 12.09 | 29.02 | 39.85 |
| Δ | 1.07 ↑ | 4.25 ↑ | 5.30 ↑ | 3.33 ↑ | 2.48 ↑ | −0.40 ↓ |
| FARE | 19.65 | 14.00 | 38.79 | 6.85 | 20.00 | 51.02 |
| FARE + TTC | 35.55 | 30.52 | 59.20 | 20.46 | 33.89 | 49.91 |
| Δ | 15.90 ↑ | 16.52 ↑ | 20.41 ↑ | 13.61 ↑ | 13.89 ↑ | −1.11 ↓ |

Tab.2 Applying TTC on finetuned models further improves robustness.

## Conclusion & Discussion

1. CLIP can leverage $f_\theta$ to counterattack
2. First method to defend CLIP at test time

**Limitations to be addressed:**
- limited robustness gains on finetuned models
- Incurs compute expense at test time
- May be circumvented by adaptive attacks

### Reference
[1] Mao et al. Understanding zero-shot adversarial robustness for large-scale models. In ICLR, 2023.

[2] Schlarmann et al. Robust CLIP: Unsupervised adversarial fine-tuning of vision embeddings for robust large vision-language models. In ICML, 2024.