# PTTA:
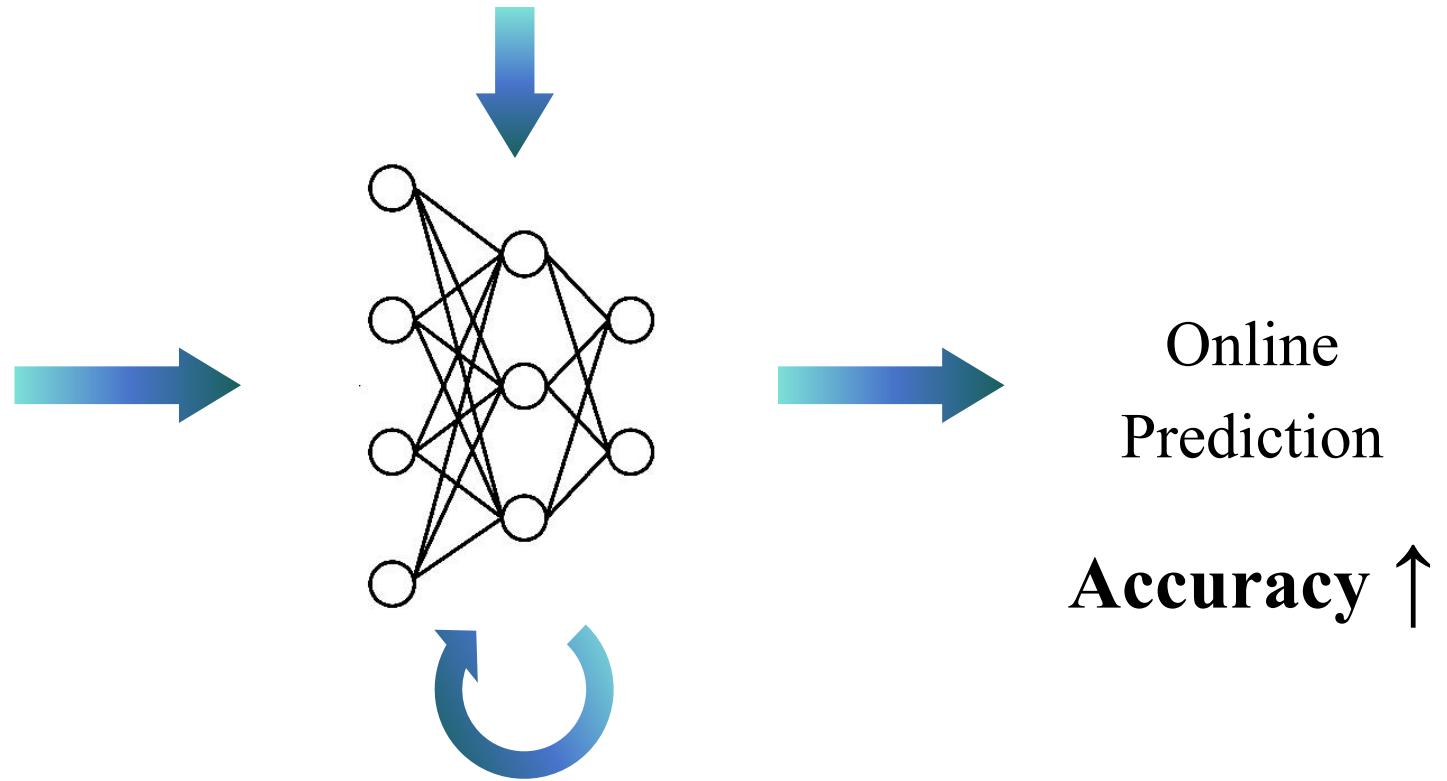# Purifying Malicious Samples for Test-Time Model Adaptation

Jing Ma,    Hanlin Li,    Xiang Xiang

# Introduction

Model Zoo (e.g., Hugging Face 🤗)

Test Samples

Online Prediction

**Accuracy ↑**

**Test-Time Model Adaptation**

# Malicious Sample Hazards

# Introduction

Sample selection wastes our limited test samples.

## Question:

Rather than <u>selecting and discarding</u> malicious samples, why not **purify** them into benign ones?

# Introduction

# Method

> 1. Logit-Saliency Indicator

$$\nabla_z \mathcal{L}_{\text{Ent}}(f_\theta(x)) = -(\mathbf{z} - \mathbf{p} \cdot \mathbf{z}) \odot \mathbf{p} \qquad (4)$$

$f_\theta$     The model with parameters θ     $\mathbf{z}$     The output logits

$x$     A test sample     $\mathbf{p}$     The predicted probabilities

$\mathcal{L}_{\text{Ent}}$     Entropy minimization objective

# Method

➤ 2. Benign Sample Retrieval



**Saliency Indicator & Distance:**

$$\mathcal{D}_{sa} = 1 - \mathrm{Cosine}(\nabla_x \mathcal{L}_{\mathrm{Ent}}(f_\theta(x_i)), \ \nabla_x \mathcal{L}_{\mathrm{Ent}}(f_\theta(x_j)))$$

**Benign Sample Retrieval:**

$$x_j^* = \underset{1 \leq j \leq N_{mb}}{\mathrm{argmax}} \mathcal{D}_{sa}(x_i, x_j)$$

# Method

➢ 3. Malicious Sample Purification



Saliency Indicator & Distance:
$$\mathcal{D}_{sa} = 1 - \text{Cosine}(\nabla_x \mathcal{L}_{\text{Ent}}(f_\theta(x_i)), \ \nabla_x \mathcal{L}_{\text{Ent}}(f_\theta(x_j)))$$

Benign Sample Retrieval:
$$x_j^* = \underset{1 \le j \le N_{mb}}{\text{argmax}} \mathcal{D}_{sa}(x_i, x_j)$$

Malicious Sample Purification:
$$x_{ij}^* = \lambda x_i + (1 - \lambda)x_j^*, \ y_{ij}^* = \lambda \hat{y}_i + (1 - \lambda)\hat{y}_j^*$$

Test Samples: Test Samples Forward

Purified Samples: Purified Samples Forward

Benign Sample Selection

Benign Sample Retrieval

Memory Bank or In-Batch

Loss Function
$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{tta}} + \alpha \mathcal{L}_{\text{pur}}$$

Sample Store & Retrieval

Purification

Purification Loss:
$$\mathcal{L}_{\text{pur}} = -\frac{1}{N_{bs}} \sum_i^{N_{bs}} \sum_c^C (y_{ij}^*)_c \log \sigma_c(f_\theta(x_{ij}^*))$$

# Main Experiments

Table 2: Experimental results (top-1 classification accuracy (%)) on the lifelong TTA task.

| METHODS | ROUND R-1 | R-2 | R-3 | R-4 | R-5 | R-6 | R-7 | R-8 | R-9 | R-10 | AVERAGE |
|---|---|---|---|---|---|---|---|---|---|---|---|
| NoAdapt | $31.6_{\pm0.00}$ | $31.6_{\pm0.00}$ | $31.6_{\pm0.00}$ | $31.6_{\pm0.00}$ | $31.6_{\pm0.00}$ | $31.6_{\pm0.00}$ | $31.6_{\pm0.00}$ | $31.6_{\pm0.00}$ | $31.6_{\pm0.00}$ | $31.6_{\pm0.00}$ | $31.6_{\pm0.00}$ |
| Tent | $8.1_{\pm0.06}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $0.9_{\pm0.01}$ |
| + PTTA | $60.0_{\pm0.05}$ | $31.8_{\pm0.11}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $0.1_{\pm0.00}$ | $9.3_{\pm0.02}$ |
| CoTTA | $42.9_{\pm0.25}$ | $40.6_{\pm0.57}$ | $37.0_{\pm1.33}$ | $34.8_{\pm0.82}$ | $33.5_{\pm1.02}$ | $32.0_{\pm0.67}$ | $31.0_{\pm0.50}$ | $30.5_{\pm0.61}$ | $30.7_{\pm0.89}$ | $30.6_{\pm0.90}$ | $34.4_{\pm0.53}$ |
| + PTTA | $52.1_{\pm0.28}$ | $46.3_{\pm0.25}$ | $42.8_{\pm0.04}$ | $40.3_{\pm0.06}$ | $39.2_{\pm0.09}$ | $38.8_{\pm0.01}$ | $38.4_{\pm0.14}$ | $38.2_{\pm0.03}$ | $38.0_{\pm0.05}$ | $37.6_{\pm0.08}$ | $41.2_{\pm0.02}$ |
| SoTTA | $59.5_{\pm0.22}$ | $60.7_{\pm0.15}$ | $61.0_{\pm0.15}$ | $61.3_{\pm0.16}$ | $61.5_{\pm0.13}$ | $61.5_{\pm0.19}$ | $61.6_{\pm0.15}$ | $61.7_{\pm0.15}$ | $61.8_{\pm0.22}$ | $61.9_{\pm0.20}$ | $61.3_{\pm0.24}$ |
| + PTTA | $61.3_{\pm0.39}$ | $62.6_{\pm0.42}$ | $63.1_{\pm0.34}$ | $63.4_{\pm0.31}$ | $63.5_{\pm0.27}$ | $63.7_{\pm0.26}$ | $63.8_{\pm0.27}$ | $63.9_{\pm0.25}$ | $63.9_{\pm0.25}$ | $64.0_{\pm0.24}$ | $63.3_{\pm0.43}$ |
| SAR | $60.0_{\pm0.02}$ | $61.1_{\pm0.02}$ | $61.4_{\pm0.02}$ | $61.6_{\pm0.02}$ | $61.7_{\pm0.02}$ | $61.8_{\pm0.03}$ | $61.7_{\pm0.02}$ | $61.6_{\pm0.08}$ | $59.3_{\pm0.26}$ | $60.4_{\pm0.06}$ | $61.1_{\pm0.04}$ |
| + PTTA | $61.6_{\pm0.01}$ | $63.0_{\pm0.00}$ | $63.4_{\pm0.00}$ | $63.6_{\pm0.03}$ | $63.8_{\pm0.02}$ | $63.9_{\pm0.01}$ | $64.0_{\pm0.02}$ | $64.0_{\pm0.02}$ | $64.1_{\pm0.00}$ | $64.1_{\pm0.02}$ | $63.5_{\pm0.00}$ |
| ETA | $62.2_{\pm0.07}$ | $59.5_{\pm0.13}$ | $55.4_{\pm0.39}$ | $46.6_{\pm7.84}$ | $31.5_{\pm27.2}$ | $29.5_{\pm25.4}$ | $28.2_{\pm24.3}$ | $26.2_{\pm22.6}$ | $26.1_{\pm22.6}$ | $25.1_{\pm21.6}$ | $39.0_{\pm15.2}$ |
| + PTTA | $\underline{65.3}_{\pm0.06}$ | $\underline{65.4}_{\pm0.03}$ | $\underline{65.3}_{\pm0.06}$ | $\underline{65.1}_{\pm0.01}$ | $64.9_{\pm0.03}$ | $64.7_{\pm0.05}$ | $64.6_{\pm0.03}$ | $64.4_{\pm0.06}$ | $64.3_{\pm0.04}$ | $\underline{64.1}_{\pm0.04}$ | $64.8_{\pm0.00}$ |
| EATA | $62.5_{\pm0.40}$ | $62.2_{\pm0.39}$ | $61.9_{\pm0.37}$ | $61.7_{\pm0.44}$ | $61.6_{\pm0.48}$ | $61.4_{\pm0.44}$ | $61.3_{\pm0.53}$ | $61.2_{\pm0.38}$ | $61.0_{\pm0.42}$ | $61.0_{\pm0.35}$ | $61.6_{\pm0.41}$ |
| + PTTA | $64.7_{\pm0.31}$ | $65.0_{\pm0.43}$ | $65.0_{\pm0.39}$ | $65.1_{\pm0.43}$ | $\underline{65.1}_{\pm0.40}$ | $\underline{65.0}_{\pm0.38}$ | $\underline{65.0}_{\pm0.45}$ | $\underline{65.0}_{\pm0.43}$ | $\mathbf{65.0}_{\pm0.39}$ | $\mathbf{65.0}_{\pm0.41}$ | $\underline{65.0}_{\pm0.40}$ |
| DeYO | $62.0_{\pm0.50}$ | $43.1_{\pm18.9}$ | $32.7_{\pm29.0}$ | $25.0_{\pm27.8}$ | $17.9_{\pm30.7}$ | $16.8_{\pm28.8}$ | $15.9_{\pm27.4}$ | $3.5_{\pm5.94}$ | $0.1_{\pm0.01}$ | $0.1_{\pm0.00}$ | $21.7_{\pm15.9}$ |
| + PTTA | $\mathbf{65.8}_{\pm0.01}$ | $\mathbf{66.0}_{\pm0.05}$ | $\mathbf{66.0}_{\pm0.04}$ | $\mathbf{65.9}_{\pm0.05}$ | $\mathbf{65.9}_{\pm0.05}$ | $\mathbf{65.8}_{\pm0.07}$ | $\mathbf{65.7}_{\pm0.05}$ | $\mathbf{65.5}_{\pm0.07}$ | $\underline{64.9}_{\pm0.91}$ | $61.6_{\pm4.30}$ | $\mathbf{65.3}_{\pm0.53}$ |
| CPL | $55.5_{\pm0.13}$ | $57.3_{\pm0.11}$ | $57.8_{\pm0.10}$ | $58.1_{\pm0.15}$ | $58.4_{\pm0.27}$ | $58.5_{\pm0.21}$ | $58.5_{\pm0.20}$ | $58.7_{\pm0.11}$ | $58.8_{\pm0.10}$ | $58.8_{\pm0.21}$ | $58.0_{\pm0.13}$ |
| + PTTA | $59.7_{\pm0.02}$ | $61.5_{\pm0.02}$ | $62.0_{\pm0.03}$ | $62.3_{\pm0.02}$ | $62.6_{\pm0.02}$ | $62.7_{\pm0.00}$ | $62.9_{\pm0.03}$ | $63.0_{\pm0.00}$ | $63.1_{\pm0.01}$ | $63.2_{\pm0.03}$ | $62.3_{\pm0.01}$ |

# Main Experiments

➢ Why logit-saliency indicator works?



(a) Logit-saliency indicator      (b) Feature-saliency indicator      (c) Pixel-saliency indicator
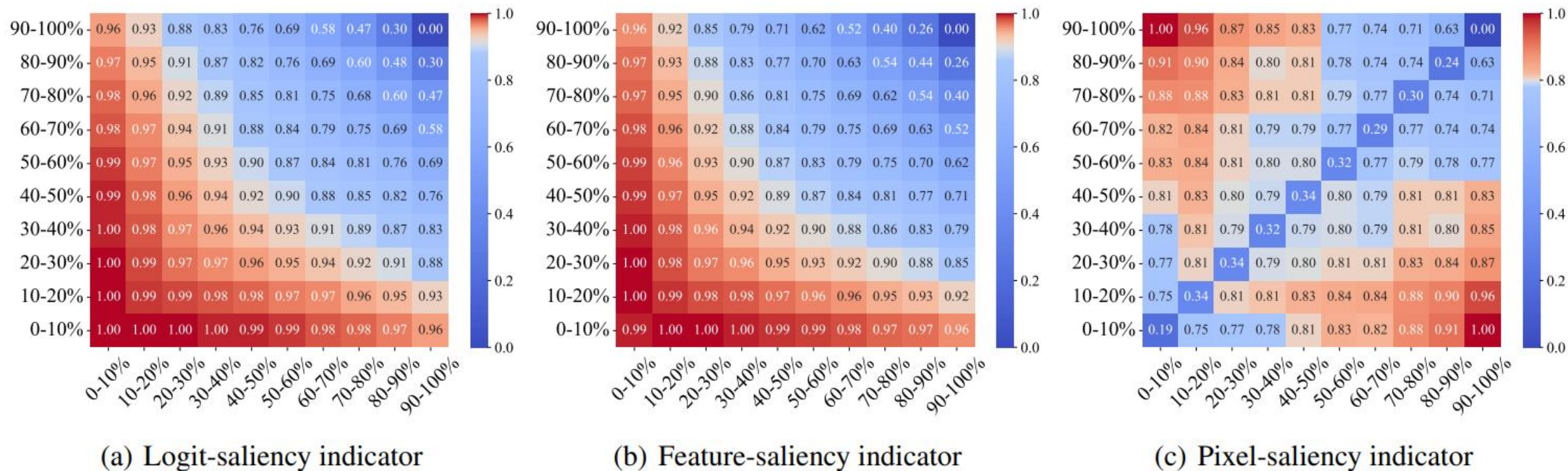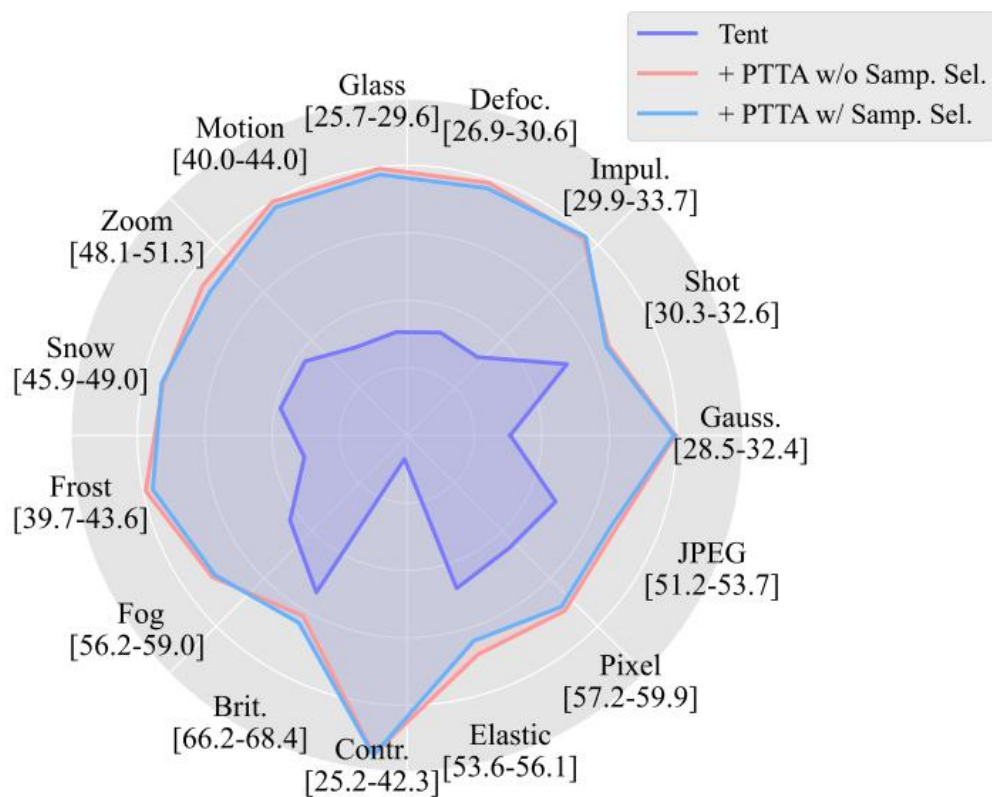
Figure 6: The Saliency Distance (normalized to the range of $0 \sim 1$) among test samples sorted in the ascending order of prediction entropy and split by percentages. A good indicator satisfies $\mathcal{D}_{sa}(x^+, x^+) > \mathcal{D}_{sa}(x^+, x^-) > \mathcal{D}_{sa}(x^-, x^-)$.

# Main Experiments

> **Unnecessary** of Sample Selection
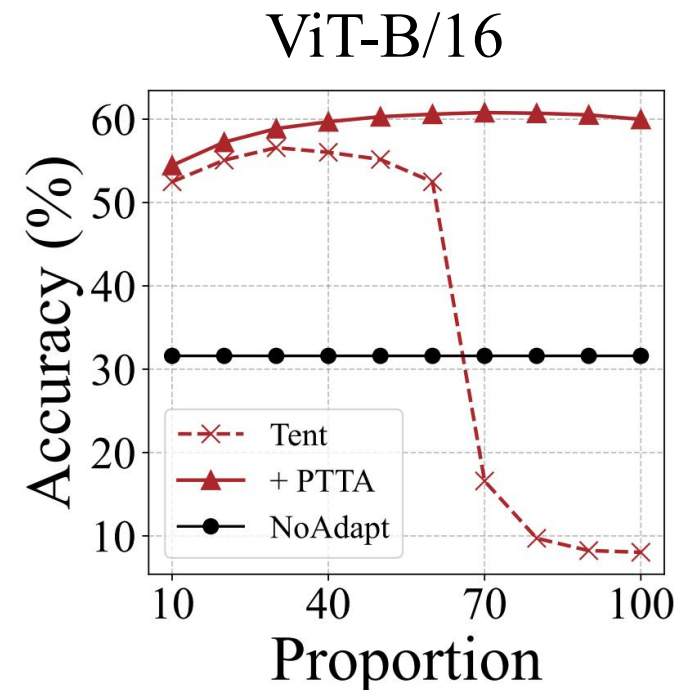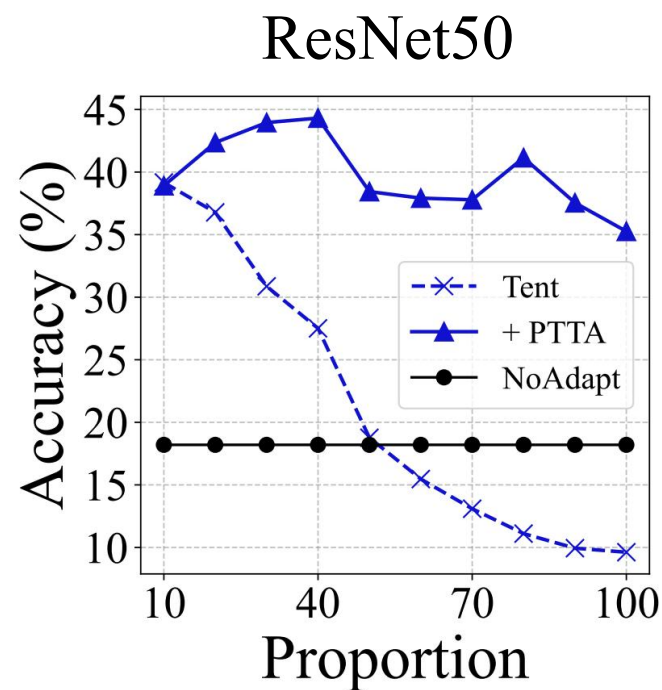
> **Insensitivity** to thresholds



Figure 5: Comparison of selecting benign samples and using all samples as candidates for purifying malicious samples.

# Main Experiments

➤ Efficiency of PTTA

Table 8: The running time (seconds) per batch for different TTA algorithms and their PTTA-applied versions. The batch size is set to 64. Δ denotes the runtime increase ratio of PTTA-applied versions compared to the base TTA methods.

| METHODS | RESNET50 | VIT-B/16 | METHODS | RESNET50 | VIT-B/16 |
|---|---|---|---|---|---|
| TENT | 0.157 | 0.264 | CPL | 0.157 | 0.310 |
| + PTTA | 0.219 | 0.376 | + PTTA | 0.225 | 0.400 |
| Δ | 39.5% | 42.4% | Δ | 43.3% | 29.0% |
| ETA | 0.158 | 0.270 | SAR | 0.225 | 0.540 |
| + PTTA | 0.225 | 0.384 | + PTTA | 0.328 | 0.770 |
| Δ | 42.4% | 42.2% | Δ | 45.8% | 42.6% |
| EATA | 0.171 | 0.291 | COTTA | 0.530 | 1.440 |
| + PTTA | 0.236 | 0.410 | + PTTA | 0.673 | 1.597 |
| Δ | 38.0% | 40.9% | Δ | 27.0% | 10.9% |
| DEYO | 0.177 | 0.328 | SOTTA | 0.794 | 1.379 |
| + PTTA | 0.254 | 0.420 | + PTTA | 0.872 | 1.597 |
| Δ | 43.5% | 28.0% | Δ | 9.80% | 15.8% |

Table 9: The storage overhead of the memory bank for a first-in-first-out queue with a maximum length of $1,000$.

| LOGIT-SALIENCY INDICATOR | PREDICTED PROBABILITIES | RAW IMAGES |
|---|---|---|
| 28.0 KB | 28.0 KB | 35.0 MB |

# Thank You



More interesting analyses in our paper!



Code available on Github!