

R-TPT: Improving Adversarial Robustness of Vision-Language Models through Test-Time Prompt Tuning

Lijun Sheng^{1,2}, Jian Liang^{2,3*}, Zilei Wang¹, Ran He^{2,3}

¹ University of Science and Technology of China

² NLPR & MAIS, Institute of Automation, Chinese Academy of Sciences

³ University of Chinese Academy of Sciences

slj0728@mail.ustc.edu.cn, liangjian92@gmail.com

Abstract

Vision-language models (VLMs), such as CLIP, have gained significant popularity as foundation models, with numerous fine-tuning methods developed to enhance performance on downstream tasks. However, due to their inherent vulnerability and the common practice of selecting from a limited set of open-source models, VLMs suffer from a higher risk of adversarial attacks than traditional vision models. Existing defense techniques typically rely on adversarial fine-tuning during training, which requires labeled data and lacks of flexibility for downstream tasks. To address these limitations, we propose robust test-time prompt tuning (R-TPT), which mitigates the impact of adversarial attacks during the inference stage. We first reformulate the classic marginal entropy objective by eliminating the term that introduces conflicts under adversarial conditions, retaining only the pointwise entropy minimization. Furthermore, we introduce a plug-and-play reliability-based weighted ensembling strategy, which aggregates useful information from reliable augmented views to strengthen the defense. R-TPT enhances defense against adversarial attacks without requiring labeled training data while offering high flexibility for inference tasks. Extensive experiments on widely used benchmarks with various attacks demonstrate the effectiveness of R-TPT. The code is available in <https://github.com/TomSheng21/R-TPT>.

1. Introduction

Vision-language models (VLMs) [5, 32, 57] are multimodal models pretrained on large-scale paired image-text data. Their powerful zero-shot inference capability and broad applicability across a range of downstream tasks have made them a foundational tool in the research community. CLIP [32], a milestone work, aligns features from the text and

*To whom correspondence should be addressed.

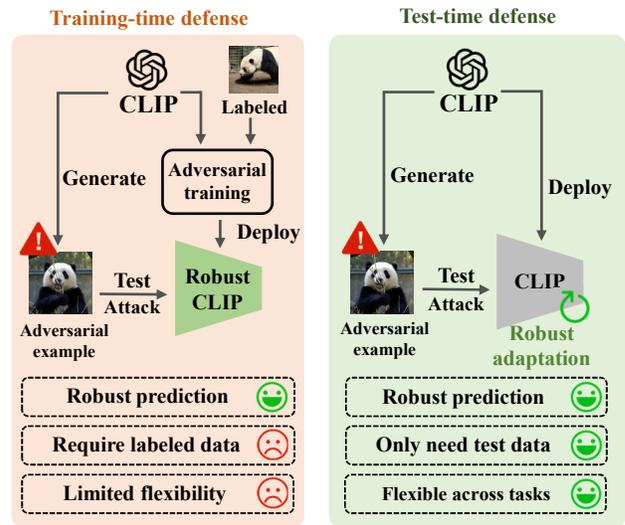


Figure 1. Comparison between training-time and test-time defense for CLIP. Our test-time defense paradigm provides robust prediction as the conventional training-time methods and requires no annotated dataset or adversarial training.

visual modalities using two specialized feature extractors trained with a contrastive loss function. Due to its concise architecture and impressive performance, CLIP has become the most widely used VLM across diverse research topics. For classification tasks, CLIP extracts features from both images and category descriptions, then chooses the category whose features exhibit the highest similarity to the image’s feature representation. Beyond classification [48, 49, 61, 62], CLIP has been successfully applied to various vision tasks, such as semantic segmentation [37, 60], object detection [56, 59], and image clustering [3, 22].

The impressive performance on downstream tasks and the broad range of applications of CLIP not only highlight its powerful capabilities but also expose it to potential vul-

nerabilities, particularly under adversarial attacks [19, 26]. While adversarial attacks and defenses [13, 24, 55] have been widely explored for conventional visual models, the situation for CLIP is more complex. The pre-training of CLIP requires the collection of vast amounts of image-text pairs and substantial computational resources that most deployers cannot afford. As a result, many deployers choose to adopt open-source versions of CLIP from a small range of candidates. This fact introduces a relatively high risk of adversarial attacks targeting CLIP-based applications.

Recent works [21, 26] explore adversarial prompt tuning using annotated data to enhance the robustness of CLIP. However, their reliance on labeled data and limited flexibility across tasks pose challenges for real-world deployment. To address this, we choose to defend adversarial attacks in the inference stage, which is applicable across various scenarios and requires no labeled dataset or prior knowledge of the downstream task, as shown in Figure 1. Existing test-time adaptation methods [38, 54] primarily focus on improving accuracy for clean test samples, while overlooking the potential risks posed by adversarial attacks. Moreover, deploying defense at test time necessitates short inference time and avoiding using additional resources, such as large language models or diffusion models, to ensure versatility.

To address the above challenges and achieve successful defense against potential attacks, we propose robust test-time prompt tuning (R-TPT). First, we revisit and refine the widely used optimization objective for instance adaptation. Many previous works [38, 42], following MEMO [58], augment the test instance and aim to minimize marginal entropy, which is defined as the entropy of the mean prediction. We decompose marginal entropy into two components: a pointwise entropy term and the Kullback–Leibler (KL) divergence, which measures the divergence between predictions from each augmented view and the mean prediction. We observe that when adapting to adversarial examples with high-confidence inaccurate predictions, the KL divergence term tends to pull the augmented views toward the misleading mean prediction, which does not exist in the clean scenario. This meaningless operation introduces conflicts into the optimization process. To mitigate the influence of adversarial samples and preserve simplicity, we discard the KL divergence term, retaining only the pointwise entropy minimization for tuning textual prompts. This straightforward modification not only defends against adversarial attacks but also maintains clean performance.

In order to effectively leverage knowledge from augmented views, we propose a reliability-based weighted ensembling strategy. To assign a larger weight to reliable augmented views during ensembling, we introduce a similarity-based metric to assess the reliability of samples. We hypothesize that samples with higher similarity to their neighbors are farther away from outliers and contain more reli-

able information. Thus, we calculate the average similarity of each sample with its neighbors to estimate its reliability. Using this metric, outliers such as adversarial examples and noisy augmented views are assigned lower reliability scores, which means less participation in the ensembling. Finally, we obtain the final prediction by ensembling the individual model predictions, weighted according to their reliability. Extensive experiments on fine-grained classification benchmarks and distribution shift benchmarks validate the effectiveness of our method in both adaptation and defense against adversarial attacks. Our contributions are summarized as follows:

- We propose R-TPT, which is the first to explore test-time paradigms for defending against potential adversarial attacks in CLIP.
- We discard the KL divergence term from the marginal entropy objective to eliminate optimization conflicts and propose a reliability-based weighted ensembling strategy to integrate knowledge from reliable augmented views.
- Extensive experiments demonstrate the effectiveness of our method in both adaptation and adversarial defense.

2. Related Work

2.1. Adversarial Attack and Defense

A lot of research [13, 24, 43] has been devoted to studying neural network’s vulnerability to adversarial noise. A pioneering work [43] introduces the concept of adversarial examples and finds small-amplitude noise that humans cannot recognize leads to misclassifications. Since then, a series of attack methods to generate adversarial samples [4, 8, 13, 24] have been proposed. FGSM [13] proposes to utilize the gradient sign to generate an adversarial example. Researchers [24] generate the adversarial noise by projected gradient descent (PGD) operation, which becomes the standard measurement of model robustness. Moreover, a variety of works explore adversarial attacks in many restricted conditions, ranging from one-pixel attack [41], universal perturbation [28] to more realistic black-box setting [2, 18].

In parallel, numerous defense strategies [24, 36, 50, 55] have been proposed to mitigate adversarial attacks. Adversarial training [13, 24] improves the robustness of the model by incorporating the adversarial samples into the training set. TRADES [55] provides a theoretical analysis of adversarial error to trade adversarial robustness against accuracy. AWP [50] perturbs the model’s weights with small adversarial noise during training to enhance robustness. Reconstruction [29, 34] with generative models [12, 39] is also a commonly used technique in test-time defense methods.

As for the vision-language models, TeCoA [26] and APT [21] employ adversarial training to pretrained CLIP [32] to improve the adversarial robustness. However, their training-time solution requires an annotated dataset and lacks flexi-

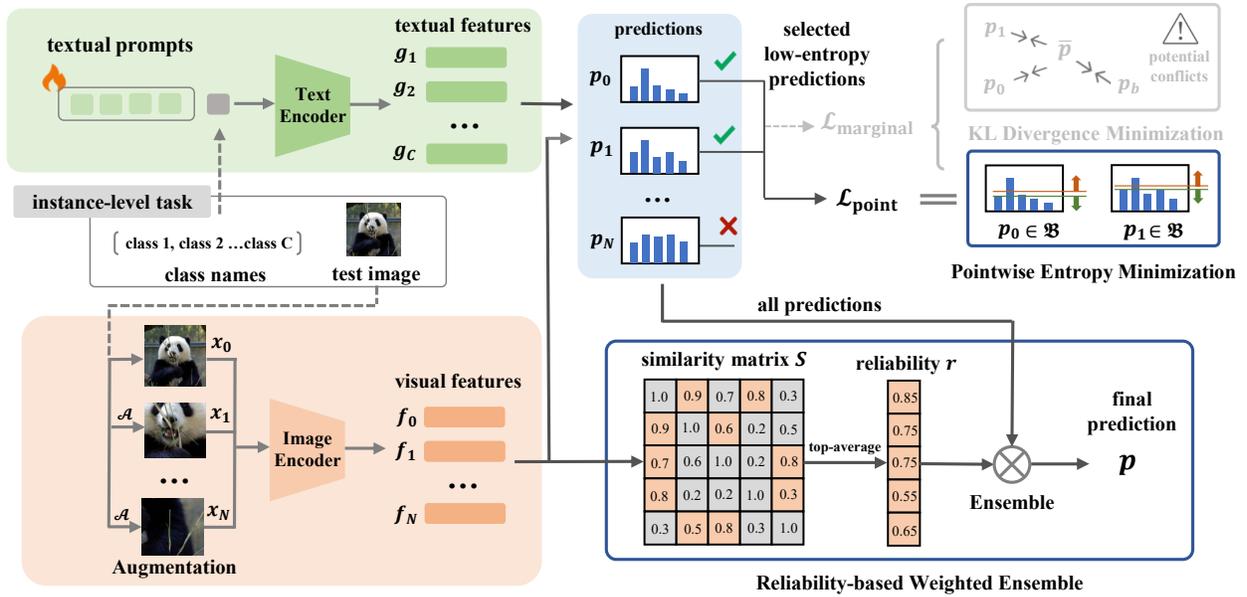


Figure 2. The pipeline of R-TPT framework. Given an instance-level task, we deploy augmentation on the test image and build a classifier with CLIP’s text branch. After selecting a low-entropy batch, R-TPT discards the KL divergence minimization term which potentially introduces conflicts in the marginal entropy [58] and optimizes textual prompts with pointwise entropy minimization. To effectively utilize the knowledge of the augmented views, R-TPT applies a reliability-based weighted ensembling mechanism in the final inference process.

bility across tasks. In this work, we choose to deploy adversarial defense in the test time that is more effective and flexible and can collaborate with their training-time defense.

2.2. Test-Time Adaptation for VLMs

Test-time adaptation [23, 38, 52] aims to adapt pre-trained models to the test data at inference time to improve the performance further. According to the test data form, test-time adaptation is divided into streaming data adaptation [44, 46, 53] and single instance adaptation [38, 58], and our work focuses on the latter. Recent research [38, 42, 54] has been devoted to exploring the instance adaptation methods for VLMs [32, 47, 57]. TPT [38] employs the marginal entropy minimization to augmentation views of the test instance to correct the prediction. DiffTPT [11] utilizes the diffusion technique to obtain diverse views which is helpful for the adaptation. PromptAlign [1] aligns the statistics of the test instance and collected natural images to make the model’s parameters adapt to test samples. To take advantage of more prompt templates [32], TPS [42] chooses to optimize feature shift and utilizes prompt ensembling for initialization. Moreover, researchers [54] propose a training-free adaptation method by ensembling the augmented views with the MeanShift algorithm [7].

Besides accuracy, researchers [51] also focus on improving calibration performance through higher text feature dispersion. In this work, we first utilize the test-time paradigm to defend against adversarial attacks for CLIP due to its se-

vere vulnerability to adversarial examples.

3. Methodology

In this paper, we improve CLIP’s adversarial robustness through a test-time paradigm, motivated by its inherent vulnerabilities and the high resource demands of train-time defense methods. Our proposed robust test-time prompt tuning (R-TPT), as illustrated in Figure 2, requires no annotated data and is equipped with greater flexibility. In Sec. 3.1, we begin with reviewing the foundational concepts of CLIP and the test-time prompt tuning approach. We then introduce the two key components of R-TPT: pointwise entropy minimization in Sec. 3.2, and the reliability-based weighted ensembling strategy in Sec. 3.3.

3.1. Preliminary

Contrastive language-image pre-training. CLIP is a popular language-vision model with a double-tower architecture, consisting of an image encoder $F(\cdot)$ and a text encoder $G(\cdot)$. It is pretrained by the contrastive learning objective with a large amount of image-text pairs. Benefiting from rich pretraining knowledge, CLIP has a strong zero-shot generalization ability. Take a C -way classification task with class names $\{t_c\}_{c=1}^C$ as an example, CLIP obtains textual features g_c by the text encoder $G(\cdot)$ with a prompt template (e.g., “a photo of a []”) and the class name t_c as the input. Also, each test image x_i queries the image encoder

to calculate the image feature $f_i = F(x_i)$. The probability that x_i belongs to category c is calculated by a softmax operation with the cosine similarity of those features:

$$p_c(x_i) = \frac{\exp(\cos(f_i, g_c)/\tau)}{\sum_{j=1}^C \exp(\cos(f_i, g_j)/\tau)}, \quad (1)$$

where $\cos(\cdot)$ represents the cosine similarity operation and τ refers to the temperature default set to 0.01.

Test-time prompt tuning. Although CLIP exhibits strong classification performance, it is sensitive to distribution shifts. To address this issue, test-time prompt tuning (TPT) [38] improves the model’s performance on individual test instances, without requiring additional training data for adaptation. During the test time, TPT deploys augmentation operations on the test instance x_0 to obtain augmented views $\{x_i\}_{i=1}^N$ and tunes the textual prompts with low-entropy views. The core objective of TPT is to minimize marginal entropy, which is formulated as:

$$\mathcal{L}_{\text{marginal}} = \mathcal{H}(\bar{p}) = \mathcal{H}\left(\frac{1}{|\mathcal{B}|} \sum_{x \in \mathcal{B}} p(x)\right), \quad (2)$$

where $\mathcal{H}(\cdot)$ is the Shannon entropy and \mathcal{B} represents the sample set selected from all views $\{x_i\}_{i=0}^N$ based on the low entropy. When augmentation provides valuable information, the classification boundary adapts according to the new prompt, leading to more accurate predictions.

3.2. Refining Marginal Entropy Minimization

In the optimization of a single test sample, marginal entropy is the default choice for both visual models [58] and multi-modal models [38]. Minimizing marginal entropy $\mathcal{H}(\bar{p})$ encourages the model to produce a consistent output across a set of selected low-entropy augmented views \mathcal{B} . However, when the test sample is adversarially perturbed, it can easily be selected into the set \mathcal{B} . While random augmentations can weaken the adversarial noise, enforcing consistency across all augmented outputs may mislead the optimization. Since our goal is to utilize test time adaptation to strengthen the model’s defense ability to adversarial examples, we propose refining the marginal entropy objective.

We decompose the marginal entropy objective into two items as follows:

$$\begin{aligned} \mathcal{H}(\bar{p}) &= -\sum_{c=1}^C \bar{p}_c \log \bar{p}_c = -\frac{1}{|\mathcal{B}|} \sum_{b=1}^{|\mathcal{B}|} \sum_{c=1}^C p_c^b \log \bar{p}_c \\ &= \frac{1}{|\mathcal{B}|} \sum_{b=1}^{|\mathcal{B}|} \left(-\sum_{c=1}^C p_c^b \log p_c^b + \sum_{c=1}^C p_c^b \log \frac{p_c^b}{\bar{p}_c} \right) \\ &= \frac{1}{|\mathcal{B}|} \sum_{b=1}^{|\mathcal{B}|} \left(\mathcal{H}(p^b) + \mathcal{KL}(p^b \parallel \bar{p}) \right), \end{aligned} \quad (3)$$

where $\mathcal{H}(\cdot)$ denotes the Shannon entropy, $\mathcal{KL}(\cdot \parallel \cdot)$ refers to Kullback–Leibler (KL) divergence. Eq. 3 shows that marginal entropy is a combination of a pointwise entropy term and a KL divergence term. Minimizing pointwise entropy helps move the classification boundary away from low-entropy points, which is the main composition of the marginal entropy. The KL divergence term, on the other hand, encourages consistent predictions across the augmented views. In the case of clean test samples, the differences between low-entropy augmented views are small, and the KL divergence term has a minimal effect. However, in adversarial scenarios, the mean prediction is distorted by the original perturbed sample, which differs significantly from the augmented views. As a result, enforcing consistency across the augmented views leads to conflicts and neglects the valuable information in these views.

To improve performance under both clean and adversarial conditions, we discard the KL divergence term and focus solely on minimizing pointwise entropy, as follows:

$$\text{minimize } \mathcal{L}_{\text{point}} = \frac{1}{|\mathcal{B}|} \sum_{b=1}^{|\mathcal{B}|} \mathcal{H}(p^b). \quad (4)$$

For clean test instances, minimizing pointwise entropy functions similarly to the original marginal entropy objective. Also, for adversarial examples, this approach focuses more on augmented views, which have relatively high entropy and information, and ignores the original adversarial inputs. Therefore, our objective can handle both natural conditions and adversarial attacks well.

3.3. Reliability-based Weighted Ensembling

Unsupervised prompt tuning typically makes no change to image features and a small adjustment to the classification boundary. For adversarial examples, it is challenging to correct incorrect predictions since their feature representations are far from the correct decision boundary. Fortunately, augmentation at the pixel level can help mitigate the adversarial noise, as it weakens the effect of perturbations. Thus, we can leverage diverse augmented views of the test image, which provide valuable knowledge. To integrate information from augmented views effectively and protect the ensembling from the noise of lower-quality views, we propose a reliability-based weighted ensembling mechanism.

Since augmentation inherently involves randomness (e.g., background areas left after cropping), we introduce a reliability metric to represent the quality of each augmented view. Reliability is defined as the similarity between the sample and its nearest neighbors. High reliability indicates that the sample is densely clustered in the feature space, containing useful information. Conversely, low reliability suggests that the sample is an outlier, likely containing augmentation or adversarial noise, and should be down-weighted or ignored during ensembling.

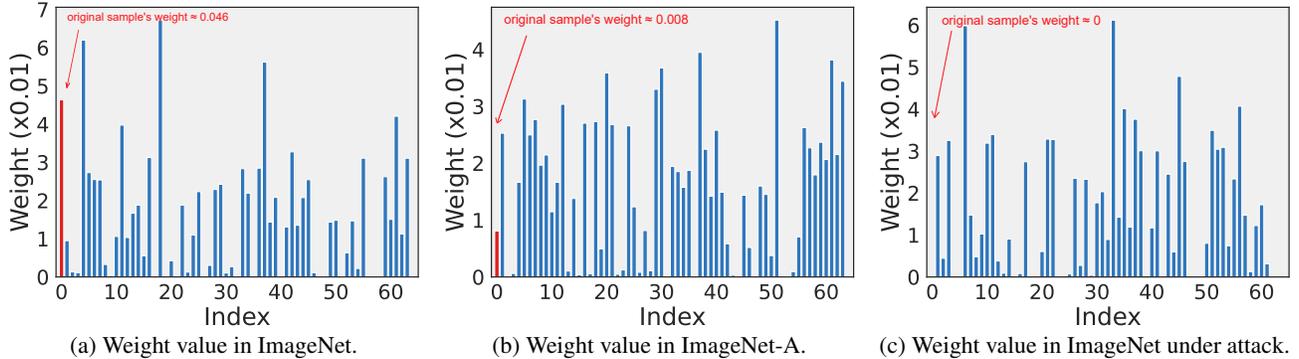


Figure 3. Visualization of assigned reliability-based weight value under (a) clean, (b) distribution shift, and (c) adversarial attack scenarios. The red and blue bars denote the weight value of the original test instance and its augmented views, respectively. The weight assigned to the clean instance is higher than under the distribution shift. The weight of the adversarial test sample is close to 0.

Given an augmented batch, which includes the test image and its N augmented views, we calculate the similarity matrix $\{S_{i,j}\}_{i,j=0}^N$ of the visual features $\{f_i\}_{i=0}^N$:

$$S_{i,j} = \cos(f_i, f_j), \quad i, j = 0, 1, \dots, N, \quad (5)$$

where $\cos(\cdot, \cdot)$ denotes the cosine similarity operation. The i -th row of the matrix S represents the similarity between x_i and all samples in the batch. We then select the K closest neighbor of x_i (excluding x_i itself) to form the neighboring set \mathcal{N}_i . The reliability of x_i is calculated as the average similarity within the neighboring set \mathcal{N}_i :

$$r_i = \frac{1}{K} \sum_{k \in \mathcal{N}_i} \cos(f_i, f_k), \quad i = 0, 1, \dots, N. \quad (6)$$

The reliability score reflects the degree to which x_i is surrounded by similar samples in the feature space.

To obtain the final prediction, we ensemble the predictions from all augmented views, weighted by their respective reliability scores. The weight assigned to each view is calculated by applying a softmax operation on the reliabilities $\{r_i\}_{i=0}^N$. This weight assignment mechanism is flexible and adapts well to various scenarios. We present the weight values for clean samples (ImageNet), samples with distribution shift (ImageNet-A), and adversarial samples (ImageNet) in Figure 3. As shown, clean samples are assigned higher weights, and the ensembling mechanism focuses on augmented views when the test instance shows significant distribution shifts from natural images. For adversarial samples, the weight assigned to the poisoned instance is close to zero, which indicates that our mechanism protects the ensembling process from being misled by adversarial noise. This approach not only improves the robustness of predictions on adversarial samples but also maintains strong performance on clean samples, thus enhancing the overall reliability of the inference process.

4. Experiment

4.1. Setup

Datasets. To evaluate our proposed test-time defense method for VLMs, we conduct experiments on eight fine-grained classification datasets. These databases cover general objects (**Caltech101** [10]), animals (**Pets** [31]), plants (**Flower102** [30]), vehicles (**Cars** [20], **Aircraft** [25]), textures (**DTD** [6]), satellite images (**EuroSAT** [14]), and actions from videos (**UCF101** [40]). Moreover, we evaluate on **ImageNet** [9] and four ImageNet-out-of-distribution (OOD) benchmarks with distribution shift. **ImageNet-A** [17] contains 200 classes and 7,500 natural adversarial examples which are collected with adversarial filtration technique. **ImageNet-V2** [33] consists of 10,000 natural images across 1,000 categories from a different source than ImageNet. **ImageNet-R** [16] is a dataset containing 30,000 images with various renditions, leading to different textures and local statistics from ImageNet. **ImageNet-S** [45] consists of 50,889 sketch-style images and shares the same category space with ImageNet. Since our method is to defend against potential adversarial attacks during the test time, we do not need access to the training set of the above datasets.

Evaluation metrics. Since our task focuses on instance-level test-time adaptation, the model update and prediction of each test sample can not utilize the information of other samples. Following the previous works [38, 54], we report the average classification accuracy (**Acc.**) to measure the method’s adaptation ability on clean samples. To evaluate the adversarial defense performance, we provide the average accuracy (**Rob.**) on adversarial examples generated by the PGD algorithm [24] with various noise radii. It is worth noting that the adversarial examples are calculated on CLIP before adaptation, which is more suitable for real-world applications since attackers rely on open-source models and have no idea about the victim model’s algorithm.

Method	Caltech101		Pets		Cars		Flower102		Aircraft		DTD		EuroSAT		UCF101		Avg.	
	Acc.	Rob.																
CLIP [32]	85.9	2.6	83.5	0.0	55.7	0.0	61.7	0.0	15.7	15.7	40.4	0.8	23.7	0.0	58.9	0.0	53.2	2.4
Ensemble	83.5	74.8	82.3	69.9	57.1	36.2	58.0	46.6	16.4	16.4	37.1	29.5	16.7	13.7	53.9	43.0	50.6	41.3
TPT [38]	87.9	7.0	84.7	0.1	58.4	0.0	62.1	0.0	17.3	17.3	42.4	4.3	28.4	0.0	60.6	0.3	55.2	3.6
C-TPT [51]	87.7	3.7	83.6	0.0	56.6	0.0	64.8	0.0	16.7	16.7	41.5	1.3	27.0	0.0	60.1	0.1	54.8	2.7
MTA [54]	87.3	65.9	84.8	59.8	58.7	17.8	61.0	31.5	18.1	18.1	40.3	18.8	22.5	1.6	60.6	31.3	54.1	30.6
R-TPT	86.7	79.8	84.6	74.2	58.1	42.9	60.6	51.9	17.5	17.5	41.3	33.5	21.2	15.9	59.7	50.9	53.7	45.8

Table 1. Results (%) of clean accuracy (Acc.) and adversarial accuracy (Rob.) of various adaptation methods on **fine-grained classification datasets** with pre-trained CLIP-ResNet50 ($\epsilon = 1.0$). Best clean accuracies are (**bold**), best adversarial accuracies are (**bold red**).

Method	Caltech101		Pets		Cars		Flower102		Aircraft		DTD		EuroSAT		UCF101		Avg.	
	Acc.	Rob.	Acc.	Rob.	Acc.	Rob.	Acc.	Rob.										
CLIP [32]	94.0	0.0	88.3	0.0	65.5	0.0	67.4	0.0	23.9	23.9	44.4	0.0	42.2	0.0	65.2	0.0	61.4	3.0
Ensemble	91.9	74.7	86.2	51.2	65.7	26.0	65.9	36.3	23.4	23.4	43.2	25.1	28.2	2.2	63.0	30.6	58.4	33.7
TPT [38]	94.1	0.0	87.4	0.0	66.5	0.0	69.1	0.0	23.4	23.4	46.9	0.0	42.6	0.0	67.9	0.0	62.2	2.9
C-TPT [51]	93.9	0.0	88.2	0.0	65.8	0.0	69.6	0.0	23.9	23.9	45.9	0.0	42.3	0.0	65.5	0.0	61.9	3.0
MTA [54]	94.3	72.1	88.0	51.8	67.7	18.5	67.4	27.9	25.0	25.0	46.5	16.2	42.5	1.2	67.5	27.5	62.3	30.0
R-TPT	93.7	82.0	87.2	60.2	67.0	34.7	68.7	44.6	23.9	23.9	46.4	32.8	34.7	8.5	67.2	43.2	61.1	41.2

Table 2. Results (%) of various adaptation methods on **fine-grained classification datasets** with pre-trained CLIP-ViT-B/16 ($\epsilon = 4.0$).

Baselines. We compare R-TPT on the above benchmarks with existing test-time adaptation methods for CLIP, including TPT [38], C-TPT [51], and MTA [54], as well as zero-shot prediction from CLIP [32]. Also, we regard Ensemble as an additional baseline method, which employs simple average operation on predictions of all augmented views. Note that both our method and the compared baseline methods rely only on CLIP and AugMix [15] augmentation, without any additional foundation models (*e.g.*, LLM, diffusion models) or knowledge. All results of baselines are reproduced with the official code.

Implementation details. For all experiments, we utilize official pre-trained CLIP-ResNet50 and CLIP-ViT-B/16 [32] as our base model. As for adversarial example generation, we utilize the PGD algorithm [24] with $\epsilon = 1.0$ and 7 steps for ResNet, while $\epsilon = 4.0$ and 100 steps for ViT. In the defense stage at test time, the parameter optimized in all experiments is a textual prompt with a context length of 4 and is initiated with “a photo of a”. We adopt the Adam optimizer with weight decay and a single step. Following previous work [38], the learning rate is set to 0.005 and the augmented batch size is 64. All experiments use the PyTorch framework and run on RTX3090 GPUs.

4.2. Experimental Results

Results on fine-grained datasets. We evaluate the adaptation and adversarial defense ability of R-TPT on eight fine-grained benchmarks and report the results in Table 1. It is shown that CLIP with strong zero-shot generalization ability is suffering from the adversarial attack with a small radius. TPT can steadily improve the accuracy of clean im-

ages but has weak defense capability. We observe that methods with an ensembling strategy (*e.g.*, Ensemble, MTA, R-TPT) achieve significantly better adversarial accuracy. In particular, R-TPT achieves a 45.8% of adversarial accuracy which outperforms all baseline methods. Besides the attractive defense performance, R-TPT also improves the clean accuracy of CLIP from 53.2% to 53.7%. In contrast, although Ensemble has strong defense capabilities, it suffers from negative transfer in the clean scenario.

Results on ImageNet and ImageNet-OOD datasets. We provide the experiential results on ImageNet and four ImageNet-OOD datasets in Table 3. CLIP can overcome distribution shifts, but fails to defend against adversarial attacks. R-TPT performs the best in adversarial defense under both ImageNet and its related out-of-distribution benchmarks. Especially, R-TPT achieves an adversarial accuracy of 47.7% on ImageNet, 7.6% higher than the second-best method Ensemble, while CLIP’s adversarial accuracy on this benchmark is 0.1%. Moreover, our method obtains a similar clean accuracy with TPT, indicating that R-TPT is also effective for distribution shifts.

Results of CLIP-ViT backbone. We evaluate our method on CLIP-ViT-B/16 [32] model and provide the results in Table 2. We observe that R-TPT outperforms all baseline methods regarding adversarial defense performance. However, the performance gain of all adaptation methods for clean samples is small, and R-TPT and the remaining two ensembling-based methods make weak negative transfers. This illustrates that the space for improvement of clean accuracy for strong pre-training models is small, and R-TPT can greatly enhance its weak robustness.

Method	ImageNet		ImageNet-A		ImageNet-V2		ImageNet-R		ImageNet-S		Avg.	
	Acc.	Rob.										
CLIP [32]	58.2	0.1	21.8	0.0	51.5	0.1	56.1	0.8	33.3	0.5	44.2	0.3
Ensemble	58.0	40.1	22.6	10.1	52.0	37.2	51.3	39.3	29.5	20.7	42.7	29.5
TPT [38]	60.7	0.3	26.5	0.0	54.8	0.3	58.9	1.8	35.0	1.4	47.2	0.7
C-TPT [51]	60.4	0.1	24.1	0.0	54.3	0.1	57.7	1.0	34.7	0.9	46.2	0.4
MTA [54]	60.4	30.0	27.5	5.6	54.2	24.6	58.4	29.8	35.2	11.3	47.1	20.3
R-TPT	60.9	47.7	28.4	14.4	54.9	41.6	57.6	46.9	34.0	26.2	47.1	35.4

Table 3. Results (%) of various adaptation methods on **ImageNet and ImageNet-OOD benchmarks** with pre-trained CLIP-ResNet50. OOD Avg. refers to the average results among four ImageNet-OOD benchmarks ($\epsilon = 1.0$).

Method	Caltech101		Pets		Cars		Flower102		Aircraft		DTD		EuroSAT		UCF101		Avg.	
	Acc.	Rob.	Acc.	Rob.	Acc.	Rob.	Acc.	Rob.	Acc.	Rob.	Acc.	Rob.	Acc.	Rob.	Acc.	Rob.	Acc.	Rob.
CLIP-TeCoA [26]	79.3	44.3	66.9	15.8	10.2	1.0	30.8	9.0	6.6	6.6	24.5	10.7	14.5	10.8	34.6	6.7	33.4	13.1
Ensemble	72.7	55.1	59.9	38.9	5.6	2.7	26.6	16.0	4.2	4.2	23.5	16.2	12.5	11.0	26.4	14.0	28.9	19.8
TPT [38]	79.3	52.7	65.2	27.4	9.6	2.0	27.9	12.3	6.7	6.7	25.5	14.6	12.2	11.2	34.9	10.2	32.7	17.1
C-TPT [51]	79.8	47.3	66.1	19.5	10.6	1.3	29.4	10.7	6.4	6.4	26.2	12.4	13.0	11.1	36.4	8.1	33.5	14.6
MTA [54]	79.7	55.7	66.2	31.2	9.0	2.5	29.1	14.0	6.5	6.5	24.4	13.5	13.3	11.2	34.6	12.5	32.9	18.4
R-TPT	76.1	60.5	63.2	40.1	7.7	3.5	26.6	16.5	6.1	6.1	25.2	17.7	11.5	11.3	31.1	17.4	30.9	21.7

Table 4. Results (%) of adaptation methods on **fine-grained classification datasets** with TeCoA pre-trained CLIP-ViT-B/32 ($\epsilon = 4.0$).

4.3. More Analysis

Results under robust-pretrained models. we report the results of deploying adaptation methods on CLIP-ViT-B/32 with TeCoA [26] robust pretrained models in Table 4. It is shown that TeCoA significantly improves the robustness of CLIP, and the defense effect against adversarial attacks can benefit from our method during the testing stage, which increases adversarial accuracy from 13.1% to 21.7%. At the same time, we also find that R-TPT’s defense effect on the TeCoA pretrained model is weaker than the vanilla CLIP. The reason is that introducing adversarial learning during the training time decreases the clean accuracy, which is the upper bound of defense. Robust pre-trained models also generate more difficult adversarial examples.

Analysis under various attacks. To demonstrate the versatility, we investigate the defense performance of R-TPT and baseline methods under various attack methods. We employ CW [4], DeepFool [27], and FGSM [13] as new attacks and report the adversarial accuracy on three benchmarks in Table 5. It is shown that our method improves CLIP’s defense capability under each attack, which is consistent with the trend of PGD. In particular, R-TPT has reached 51.8% adversarial accuracy on Flowers. The excellent defense performance among various attack methods proves the versatility of R-TPT.

Analysis of inference efficiency. We study the inference efficiency of R-TPT. Experiments of training-time defense method APT and R-TPT on UCF101 are provided in Table 6. It is shown that R-TPT spends a certain amount of time on each test sample, which is different from training-time defense which spends a lot of resources at one time

Method	Flowers				DTD			
	CW	DF	FGSM	Avg.	CW	DF	FGSM	Avg.
CLIP [32]	0.8	0.4	4.8	2.0	2.3	7.6	13.4	7.8
Ensemble	50.1	52.2	46.6	49.7	31.1	32.9	29.7	31.2
TPT [38]	13.8	10.8	14.2	12.9	21.3	24.4	22.2	22.6
C-TPT [51]	6.6	5.5	6.2	6.1	11.9	15.8	17.5	15.1
MTA [54]	34.5	35.4	36.6	35.5	23.6	23.5	23.9	23.7
R-TPT	51.6	54.7	49.2	51.8	34.2	35.9	32.5	34.2

Table 5. Adversarial accuracies (%) of adaptation methods against different attacks on two fine-grained datasets (DF = DeepFool).

Method	Stage	Running time	Rob.
APT+TeCoA (4shots)	Training time	208s/50epochs	39.4
R-TPT (64 views)	Test time	0.58s/image	41.0
R-TPT (32 views)	Test time	0.28s/image	40.8
R-TPT (16 views)	Test time	0.20s/image	40.0

Table 6. Running time and adversarial accuracies (%) of adaptation methods against adversarial attack on UCF101 dataset.

during training. Moreover, we find that inference can be accelerated by appropriately reducing the number of augmented views to obtain better real-time performance. When we choose to reduce the number of views from 64 to 16, the time required for each test sample is reduced to less than half of the original.

Analysis under different prompt templates. Since in the realistic scenario, the attacker can not access the prompts of the model textual branch, we evaluate the defense performance of the adversarial examples generated by

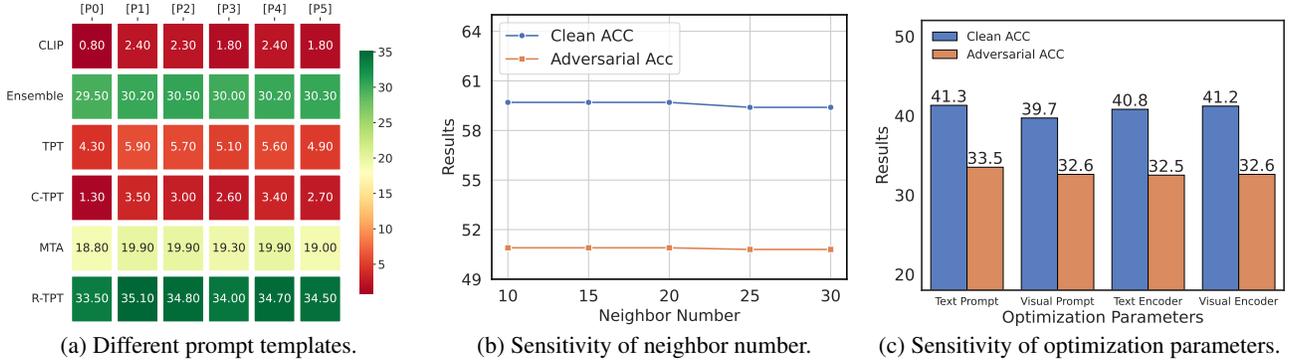


Figure 4. Results of different prompt templates in adversarial examples generation, different neighbor numbers, and different choices of optimization parameters in DTD dataset (CLIP-ResNet50, $\epsilon = 1.0$).

different prompt templates. The result of DTD benchmark with six textual templates is provided in Figure 4(a) (P0=‘a photo of a []’, P1=‘a bad photo of the []’, P2=‘a origami []’, P3=‘a photo of the large []’, P4=‘a toy []’, P5=‘art of the []’). It is shown that CLIP is suffering from adversarial samples even if attackers use different prompt templates. Also, R-TPT outperforms all baseline methods under each template.

Analysis of sensitivity of hyperparameters and optimization parameters. We study the impact of the neighbor number in calculating reliability and the optimization parameter to further demonstrate the effectiveness of R-TPT. We provide the results of R-TPT with different neighbor numbers (10, 15, 20, 25, 30) in Figure 4(b). The results demonstrate that R-TPT performs stable across different neighbor numbers in terms of both clean and adversarial accuracy. Please note that if the neighbor number equals to augmented batch size, it will assign all samples with equal weights and the reliability-based ensembling mechanism will degenerate into the vanilla ensembling. Besides, we report the results with different optimization parameters in Figure 4(c). Compared to other choices of the parameter space, optimizing textual prompts owns better defense and adaptation effects and fewer parameters, which make it always a popular solution for parameter-efficient fine-tuning.

4.4. Ablation Study

To study the contribution of terms in our method, we investigate the effectiveness of ensembling, reliability-based weighted mechanism, and pointwise entropy objective in R-TPT. The results of the ablation study are reported in Table 7. From the table, we find that entropy minimization focuses on enhancing the accuracy of adaptation on the clean samples and slightly improves the defense performance on CLIP. However, when cooperating with weighted ensembling, the improvement of entropy minimization on both metrics is significant. The ensembling strategy greatly

Ensemble	Weighted	EntMin	Fine-grained		ImageNet-X	
			Acc.	Rob.	Acc.	Rob.
✗	✗	✗	53.2	2.4	44.2	0.3
✗	✗	✓	55.2	3.6	47.2	0.7
✓	✗	✗	50.6	41.3	42.7	29.5
✓	✗	✓	53.3	44.3	46.7	34.2
✓	✓	✗	51.6	44.3	44.8	33.4
✓	✓	✓	53.7	45.8	47.1	35.4

Table 7. **Ablation study.** Clean and adversarial accuracies (%) on fine-grained datasets and ImageNet dataset (CLIP-ResNet50).

strengthens the defense capabilities of the model, but it also reduces the adaptation performance on clean samples. The weighted mechanism gives reliable samples more important roles during prediction, further improving the defense capabilities and mitigating the clean performance drop.

5. Conclusion

In this paper, we for the first time explore the adversarial defense for CLIP with a test-time paradigm and propose robust test-time prompt tuning (R-TPT). We first review the classic test time adaptation method and decompose its marginal entropy objective into a pointwise entropy term and a KL divergence term. We find that minimizing KL divergence will introduce conflicts when meeting the adversarial test instance, thus we discard KL divergence terms and only optimize the textual prompt with the pointwise entropy. We also introduce a reliability-based weighted ensembling strategy to utilize knowledge of the augmented views, which contain more knowledge than the risky original input under adversarial attacks. Extensive results show that R-TPT achieves the best defense against various adversarial attacks among all baseline methods and maintains a clean adaptation performance. We believe that our work will provide a new perspective on the defense and shed light on the safety issues of VLMs.

Acknowledgements

This work was funded by the National Natural Science Foundation of China under Grants (62276256, U2441251) and the Young Elite Scientists Sponsorship Program by CAST (2023QNRC001). The authors would like to thank Professor Tan Tieniu for his valuable guidance and contribution to this work.

References

- [1] Jameel Abdul Samadh, Mohammad Hanan Gani, Noor Hussein, Muhammad Uzair Khattak, Muhammad Muzammal Naseer, Fahad Shahbaz Khan, and Salman H Khan. Align your prompts: Test-time prompting with distribution alignment for zero-shot generalization. In *Proc. NeurIPS*, 2024. 3
- [2] Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: a query-efficient black-box adversarial attack via random search. In *Proc. ECCV*, 2020. 2
- [3] Shaotian Cai, Liping Qiu, Xiaojun Chen, Qin Zhang, and Longteng Chen. Semantic-enhanced image clustering. In *Proc. AAAI*, 2023. 1
- [4] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *Proc. S&P*, 2017. 2, 7
- [5] Fei-Long Chen, Du-Zhen Zhang, Ming-Lun Han, Xiu-Yi Chen, Jing Shi, Shuang Xu, and Bo Xu. Vlp: A survey on vision-language pre-training. *Machine Intelligence Research*, 20(1):38–56, 2023. 1
- [6] Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Proc. CVPR*, 2014. 5
- [7] Dorin Comaniciu and Peter Meer. Mean shift analysis and applications. In *Proc. ICCV*, 1999. 3
- [8] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *Proc. ICML*, 2020. 2
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proc. CVPR*, 2009. 5
- [10] Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In *Proc. CVPR Workshops*, 2004. 5
- [11] Chun-Mei Feng, Kai Yu, Yong Liu, Salman Khan, and Wangmeng Zuo. Diverse data augmentation with diffusions for effective test-time prompt tuning. In *Proc. ICCV*, 2023. 3
- [12] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Proc. NeurIPS*, 2014. 2
- [13] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *Proc. ICLR*, 2015. 2, 7
- [14] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019. 5
- [15] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. Augmix: A simple data processing method to improve robustness and uncertainty. In *Proc. ICLR*, 2020. 6
- [16] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kada-vath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proc. ICCV*, 2021. 5
- [17] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proc. CVPR*, 2021. 5
- [18] Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks with limited queries and information. In *Proc. ICML*, 2018. 2
- [19] Dehong Kong, Siyuan Liang, Xiaopeng Zhu, Yuansheng Zhong, and Wenqi Ren. Patch is enough: naturalistic adversarial patch against vision-language pre-training models. *Visual Intelligence*, 2(1):1–10, 2024. 2
- [20] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proc. ICCV Workshops*, 2013. 5
- [21] Lin Li, Haoyan Guan, Jianing Qiu, and Michael Spratling. One prompt word is enough to boost adversarial robustness for pre-trained vision-language models. In *Proc. CVPR*, 2024. 2, 11, 12
- [22] Yunfan Li, Peng Hu, Dezhong Peng, Jiancheng Lv, Jianping Fan, and Xi Peng. Image clustering with external guidance. In *Proc. ICML*, 2024. 1
- [23] Jian Liang, Ran He, and Tieniu Tan. A comprehensive survey on test-time adaptation under distribution shifts. *International Journal of Computer Vision*, pages 1–34, 2024. 3
- [24] Aleksander Madry. Towards deep learning models resistant to adversarial attacks. In *Proc. ICLR*, 2018. 2, 5, 6
- [25] Subhansu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013. 5
- [26] Chengzhi Mao, Scott Geng, Junfeng Yang, Xin Wang, and Carl Vondrick. Understanding zero-shot adversarial robustness for large-scale models. In *Proc. ICLR*, 2023. 2, 7, 11, 12
- [27] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and accurate method to fool deep neural networks. In *Proc. CVPR*, 2016. 7
- [28] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. In *Proc. CVPR*, 2017. 2
- [29] Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Anima Anandkumar. Diffusion models for adversarial purification. In *Proc. ICML*, 2022. 2
- [30] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *Proc. ICVGIP*, 2008. 5

- [31] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *Proc. CVPR*, 2012. 5
- [32] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *Proc. ICML*, 2021. 1, 2, 3, 6, 7, 11, 12
- [33] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? In *Proc. ICML*, 2019. 5
- [34] P Samangouei. Defense-gan: protecting classifiers against adversarial attacks using generative models. In *Proc. ICLR*, 2018. 2
- [35] Christian Schlarman, Naman Deep Singh, Francesco Croce, and Matthias Hein. Robust clip: Unsupervised adversarial fine-tuning of vision embeddings for robust large vision-language models. In *Proc. ICML*, 2024. 11, 12
- [36] Lijun Sheng, Jian Liang, Ran He, Zilei Wang, and Tieniu Tan. Adaptguard: Defending against universal attacks for model adaptation. In *Proc. ICCV*, 2023. 2
- [37] Gyungin Shin, Weidi Xie, and Samuel Albanie. Reco: Retrieve and co-segment for zero-shot transfer. In *Proc. NeurIPS*, 2022. 1
- [38] Manli Shu, Weili Nie, De-An Huang, Zhiding Yu, Tom Goldstein, Anima Anandkumar, and Chaowei Xiao. Test-time prompt tuning for zero-shot generalization in vision-language models. In *Proc. NeurIPS*, 2022. 2, 3, 4, 5, 6, 7, 12
- [39] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *Proc. ICLR*, 2021. 2
- [40] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. A dataset of 101 human action classes from videos in the wild. *Center for Research in Computer Vision*, 2(11):1–7, 2012. 5
- [41] Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai. One pixel attack for fooling deep neural networks. *IEEE Transactions on Evolutionary Computation*, 23(5):828–841, 2019. 2
- [42] Elaine Sui, Xiaohan Wang, and Serena Yeung-Levy. Just shift it: Test-time prototype shifting for zero-shot generalization with vision-language models. In *Proc. WACV*, 2025. 2, 3
- [43] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *Proc. ICLR*, 2014. 2
- [44] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. In *Proc. ICLR*, 2021. 3
- [45] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. In *Proc. NeurIPS*, 2019. 5
- [46] Qin Wang, Olga Fink, Luc Van Gool, and Dengxin Dai. Continual test-time domain adaptation. In *Proc. CVPR*, 2022. 3
- [47] Xiao Wang, Guangyao Chen, Guangwu Qian, Pengcheng Gao, Xiao-Yong Wei, Yaowei Wang, Yonghong Tian, and Wen Gao. Large-scale multi-modal pre-trained models: A comprehensive survey. *Machine Intelligence Research*, 20(4):447–482, 2023. 3
- [48] Zhengbo Wang, Jian Liang, Ran He, Nan Xu, Zilei Wang, and Tieniu Tan. Improving zero-shot generalization for clip with synthesized prompts. In *Proc. ICCV*, 2023. 1
- [49] Zhengbo Wang, Jian Liang, Lijun Sheng, Ran He, Zilei Wang, and Tieniu Tan. A hard-to-beat baseline for training-free clip-based adaptation. In *Proc. ICLR*, 2024. 1
- [50] Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust generalization. In *Proc. NeurIPS*, 2020. 2
- [51] Hee Suk Yoon, Eunseop Yoon, Joshua Tian Jin Tee, Mark Hasegawa-Johnson, Yingzhen Li, and Chang D Yoo. C-tpt: Calibrated test-time prompt tuning for vision-language models via text feature dispersion. In *Proc. ICLR*, 2024. 3, 6, 7, 12
- [52] Yongcan Yu, Lijun Sheng, Ran He, and Jian Liang. Benchmarking test-time adaptation against distribution shifts in image classification. *arXiv preprint arXiv:2307.03133*, 2023. 3
- [53] Yongcan Yu, Lijun Sheng, Ran He, and Jian Liang. Stamp: Outlier-aware test-time adaptation with stable memory replay. In *Proc. ECCV*, 2024. 3
- [54] Maxime Zanella and Ismail Ben Ayed. On the test-time zero-shot generalization of vision-language models: Do we really need prompt learning? In *Proc. CVPR*, 2024. 2, 3, 5, 6, 7, 12
- [55] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theoretically principled trade-off between robustness and accuracy. In *Proc. ICML*, 2019. 2
- [56] Hao Zhang, Feng Li, Xueyan Zou, Shilong Liu, Chunyuan Li, Jianwei Yang, and Lei Zhang. A simple framework for open-vocabulary segmentation and detection. In *Proc. ICCV*, 2023. 1
- [57] Jingyi Zhang, Jiaying Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 1, 3
- [58] Marvin Zhang, Sergey Levine, and Chelsea Finn. Memo: Test time robustness via adaptation and augmentation. In *Proc. NeurIPS*, 2022. 2, 3, 4
- [59] Shiyu Zhao, Zhixing Zhang, Samuel Schulter, Long Zhao, BG Vijay Kumar, Anastasis Stathopoulos, Manmohan Chandraker, and Dimitris N Metaxas. Exploiting unlabeled data with vision and language models for object detection. In *Proc. ECCV*, 2022. 1
- [60] Chong Zhou, Chen Change Loy, and Bo Dai. Extract free dense labels from clip. In *Proc. ECCV*, 2022. 1
- [61] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *Proc. CVPR*, 2022. 1
- [62] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022. 1

R-TPT: Improving Adversarial Robustness of Vision-Language Models through Test-Time Prompt Tuning

Supplementary Material

6. Algorithm

Here, we provide the pseudocode algorithm of R-TPT to show the process of our proposed defense method clearly.

Algorithm 1 Pseudocode of R-TPT.

Require: Test sample x_t , CLIP model.

- ▷ Augment x_t via AugMix to obtain views $\{x_i\}_{i=0}^N$ and select low-entropy views \mathcal{B} .
 - ▷ Update textual prompts via minimizing pointwise entropy of selected views \mathcal{B} via Eq.4.
 - ▷ Obtain the reliability $\{r_i\}_{i=0}^N$ of all views via Eq.6.
 - ▷ Obtain the robust prediction by ensembling the predictions $\{p_i\}_{i=0}^N$ of all views weighted by the reliability $\{r_i\}_{i=0}^N$.
-

7. Datasets

We provide the content, number of categories and number of images of all datasets involved in the experimental section in Table 8.

Dataset	Description	# Classes	# Test
Caltech101	Object images	100	2,465
Pets	Pet images	37	3,669
Cars	Car images	196	8,041
Flower102	Flower images	102	2,463
Aircraft	Aircraft images	100	3,333
DTD	Describable textures dataset	47	1,692
EuroSAT	Sentinel-2 satellite images	10	8,100
UCF101	Human action images	101	3,783
ImageNet	Object and scene images	1,000	50,000
ImageNet-A	Adversarially filtered images	200	7,500
ImageNet-V2	New test images	1,000	10,000
ImageNet-R	Rendered images	200	30,000
ImageNet-S	Sketch-style images	1,000	50,889

Table 8. Introduction of all datasets involved in experiments.

8. Experimental Results

8.1. Results of Larger Backbone.

We evaluate our method using the CLIP-ViT_L/14 model [32] and present the results in Table 9. Our experiments demonstrate that R-TPT outperforms all baseline methods in terms of defense performance, highlighting its robustness even when applied to large-scale backbone architectures. Also, we observe that, in terms of clean adaptation

performance, only TPT and C-TPT exhibit positive gains, whereas the remaining methods suffer from negative transfer.

8.2. Results Compared with Training-time defense Methods.

Training-time defense methods [21, 26, 35] typically rely on labeled data and robust pre-trained checkpoints to achieve their performance. To ensure a fair comparison, we have focused our main text on test-time baselines that utilize the same resources as our proposed method. Here, we provide a comprehensive evaluation of training-time methods on fine-grained datasets in Tables 10, 11 to highlight the competitive performance of R-TPT, even in the absence of external data and pre-trained robust checkpoints. It is shown that R-TPT not only remains competitive with training-time methods but also achieves significantly better performance on clean samples. More importantly, R-TPT can further improve the robustness of training-time methods.

Method	Caltech101		Pets		Cars		Flower102		Aircraft		DTD		EuroSAT		UCF101		Avg.	
	Acc.	Rob.																
CLIP [32]	95.2	0.1	93.1	0.0	76.8	0.0	76.2	0.0	30.0	30.0	52.4	0.0	55.1	0.0	73.7	0.0	69.1	3.8
Ensemble	94.9	83.6	93.4	63.5	76.3	40.5	75.0	48.6	31.7	31.7	51.3	31.3	38.7	11.1	71.7	48.3	66.6	44.8
TPT [38]	95.9	0.2	93.8	0.0	78.0	0.0	76.9	0.0	31.6	31.6	55.1	0.0	51.8	0.0	74.7	0.0	69.7	4.0
C-TPT [51]	95.6	0.1	94.3	0.0	77.4	0.0	76.3	0.0	30.4	30.4	55.4	0.0	54.0	0.0	75.1	0.0	69.8	3.8
MTA [54]	95.8	83.1	93.7	64.9	78.4	36.6	76.1	44.2	32.7	32.7	53.4	27.2	47.8	7.5	74.7	47.5	69.1	43.0
R-TPT	95.7	88.2	93.7	72.9	77.2	49.1	76.2	55.6	31.7	31.7	54.0	38.0	44.3	20.4	74.3	55.6	68.4	51.4

Table 9. Results (%) of various adaptation methods on **fine-grained classification datasets** with pre-trained CLIP-ViT-L/14 ($\epsilon = 4.0$).

Method	Caltech101		Pets		Cars		Flower102		Aircraft		DTD		EuroSAT		UCF101		Avg.	
	Acc.	Rob.																
CLIP [32]	85.9	2.6	83.6	0.0	55.7	0.0	61.7	0.0	15.7	15.7	40.4	0.8	23.7	0.0	59.0	0.0	53.2	2.4
TeCoA ¹ [26]	78.3	78.3	76.0	75.8	22.4	22.3	33.5	33.4	5.8	5.8	26.2	26.0	16.5	16.6	38.4	38.1	37.1	37.0
APT ¹ [21]	2.9	1.7	31.9	3.8	8.5	0.6	2.6	1.1	0.9	0.9	16.6	7.9	17.0	4.0	11.2	0.9	11.4	2.6
APT ¹ +TeCoA ¹ [21]	82.8	82.8	79.3	79.0	33.9	33.6	42.7	42.6	9.9	9.9	39.2	39.0	32.9	32.9	51.5	51.4	46.5	46.4
R-TPT	86.7	79.8	84.6	74.2	58.1	42.9	60.6	51.9	17.5	17.5	41.3	33.5	21.2	15.9	59.7	50.9	53.7	45.8

Table 10. Results (%) of training-time defense methods on **fine-grained classification datasets** with pre-trained ResNet50 ($\epsilon = 1.0$).

Method	Caltech101		Pets		Cars		Flower102		Aircraft		DTD		EuroSAT		UCF101		Avg.	
	Acc.	Rob.																
CLIP [32]	91.4	0.2	85.1	0.0	60.1	0.0	64.0	0.0	18.1	18.1	43.0	0.0	35.8	0.0	61.6	0.0	57.4	2.3
TeCoA ⁴ [26]	79.3	78.0	66.9	63.7	10.2	9.1	30.8	28.9	6.6	6.6	24.5	24.0	14.5	14.3	34.6	33.4	33.4	32.2
FARE ⁴ [35]	86.3	85.4	76.7	73.8	39.2	34.4	37.0	34.0	9.5	9.5	28.3	27.3	16.6	16.3	44.2	41.9	42.2	40.3
APT ⁴ [21]	10.7	0.4	10.0	0.2	1.5	0.1	0.9	0.2	2.6	2.6	9.0	0.1	7.8	6.7	3.7	0.2	5.8	1.3
APT ⁴ +TeCoA ⁴ [21]	81.4	80.2	66.7	63.9	20.8	18.9	42.5	40.4	5.2	5.2	35.2	33.7	29.3	29.2	40.2	39.4	40.2	38.9
R-TPT	90.6	76.4	84.5	55.8	63.1	28.4	62.6	37.6	19.1	19.1	42.1	29.1	32.0	5.1	62.8	41.0	57.1	36.6

Table 11. Results (%) of training-time defense methods on **fine-grained classification datasets** with pre-trained ViT-B/32 ($\epsilon = 4.0$).