

# DISCRETE LATENT FEATURES ABLATE ADVERSARIAL ATTACK: A ROBUST PROMPT TUNING FRAMEWORK FOR VLMS

006 **Anonymous authors**

007 Paper under double-blind review

## ABSTRACT

013 While adversarial fine-tuning can enhance the robustness of vision-language models  
 014 (VLMs), such approaches are computationally expensive. Adversarial prompt  
 015 tuning has emerged as a practical alternative. However, existing methods are  
 016 limited by their reliance on vulnerable continuous image features. To mitigate  
 017 the vulnerability in the feature representation, we propose **DEFEAT** (Discrete  
 018 LatEnt FeaturE based Adversarial Training), a robust prompt tuning framework  
 019 for VLMs. Specifically, the DEFEAT method introduces a perturbation discrete  
 020 shield module that reconstructs discrete latent features and designs a logits fusion  
 021 strategy, substantially reducing the discrepancy between clean and adversarial  
 022 image representations. Moreover, the DEFEAT method integrates prompt tuning  
 023 with adversarial training while applying regularization from learnable prompts  
 024 to hand-crafted prompts, further enhancing the adversarial robustness. Extensive  
 025 experiments across 15 datasets validate the effectiveness of the proposed DEFEAT  
 026 method among existing adversarial prompt tuning methods.

## 1 INTRODUCTION

029 In recent years, vision-language models (VLMs), particularly CLIP  
 030 (Radford et al., 2021) and its successors (Sun et al., 2023; 2024;  
 031 Zhang et al., 2024a; Li et al., 2022; Jia et al., 2021), have demon-  
 032 strated remarkable performance across various visual tasks. As these  
 033 models are deployed for real-world applications, it is crucial to ad-  
 034 dress their vulnerabilities. Recent studies (Touvron et al., 2023; Zou  
 035 et al., 2023; Schlar mann & Hein, 2023; Zhao et al., 2023; Carlini  
 036 et al., 2023) have shown that VLMs are highly susceptible to ad-  
 037 versarial examples (Szegedy et al., 2013), which can lead to significant  
 038 security issues by possessing imperceptible perturbations.

039 Adversarial training (Madry et al., 2018) is widely regarded as the  
 040 most effective defense against such attacks. While many studies  
 041 (Mao et al., 2023; Schlar mann et al., 2024; Hossain & Imteaj, 2024;  
 042 Gong et al., 2025) have improved the robustness of CLIP-like models  
 043 by adversarially fine-tuning the entire model, this approach is im-  
 044 practical for large-scale deployment due to the immense parameter  
 045 counts of modern VLMs. To improve efficiency, Parameter-Efficient  
 046 Fine-Tuning (PEFT) techniques (Ding et al., 2023) have emerged as a promising solution to enable  
 047 the adaptation of large models by tuning only a small subset of parameters. Building on this, Li et al.  
 048 (2024); Zhou et al. (2024); Zhang et al. (2024b) have integrated adversarial training with prompt  
 049 tuning (Zhou et al., 2022b), a PEFT method highly effective for VLMs, to efficiently enhance the  
 robustness on downstream tasks.

050 Among those existing methods, most of them rely on continuous image features for adversarial  
 051 defense and mainly focus on how to train robust models. In contrast, we mainly focus on the feature  
 052 representation, a different direction. After exploring this direction, we have a finding that *discretizing*  
 053 *the latent image feature could effectively mitigate adversarial attacks*. To see that, as illustrated in  
 Figure 2, latent features processed with a discretization exhibit a significantly smaller shift between

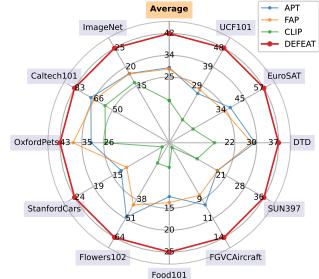


Figure 1: The harmonic mean accuracy of the robustness ( $\epsilon = 4/255$ ) and accuracy of baselines and the proposed DEFEAT method under the adversarial few-shot classification setting.

054 clean and adversarial examples than continuous image features, where the experimental settings are  
 055 introduced in Section 3.2. This indicates that the latent feature discretization could reduce the impact  
 056 of feature shifts induced by adversarial attacks.  
 057

058 Based on this finding, we propose a **Discrete**  
 059 **LatEnt Feature** based **Adversarial Training**  
 060 (**DEFEAT**) method, a robust prompt tuning  
 061 framework for VLMs. DEFEAT designs a per-  
 062 turbation discrete shield (PerturbShield) mod-  
 063 ule that mitigates adversarial attacks through  
 064 grid-based discrete feature reconstruction and  
 065 semantic alignment projection. To optimize the  
 066 trade-off between the robustness and accuracy,  
 067 the DEFEAT method uses a logits fusion strat-  
 068 egy that combines predictions from PerturbShield and original image features. Additionally, DEFEAT  
 069 applies regularization to the learnable prompt using hand-crafted prompts and demonstrates that this  
 070 regularizer enhances the adversarial robustness. Extensive experiments across 15 datasets demon-  
 071 strate that the proposed DEFEAT method achieves state-of-the-art performance compared to existing  
 072 methods and offers an improved trade-off between the robustness and performance (see Figure 1).  
 073 Specifically, the proposed DEFEAT method achieves an average improvement of 13.76% in terms of  
 074 the harmonic mean of the robustness and accuracy compared to the previous state-of-the-art method  
 075 under the adversarial few-shot classification setting.  
 076

077 Our contributions are three-fold. (i) To the best of our knowledge, we are the first to use latent  
 078 feature discretization to mitigate adversarial attacks. (ii) We analyze the mitigating effect of feature  
 079 discretization on visual adversarial perturbation and propose the DEFEAT method as a robust prompt  
 080 tuning framework for VLMs; (iii) Extensive experiments on 15 datasets demonstrate the effectiveness  
 081 of the proposed DEFEAT method under various settings, including adversarial few-shot classification,  
 082 adversarial domain generalization, and adversarial cross-dataset generalization.

## 083 2 RELATED WORK

084 **CLIP-based VLMs.** VLMs have significantly boosted cognitive capabilities by merging visual  
 085 and textual modalities, excelling in real-world vision tasks (Liu et al., 2023; Zhu et al., 2024). The  
 086 introduction of CLIP (Radford et al., 2021), trained on about 400 million image-text pairs, was  
 087 particularly transformative, establishing a new paradigm for vision-language representation learning.  
 088 Numerous subsequent works have followed this paradigm, proposing a broad family of CLIP-like  
 089 models, including ALIGN (Jia et al., 2021), EVA-CLIP (Sun et al., 2023), OpenCLIP (Ilharco et al.,  
 090 2021), and LongCLIP (Zhang et al., 2024a). Given the trailblazing role and widespread adoption of  
 091 this architecture, we focus our work on CLIP as the representative model to investigate robust prompt  
 092 tuning within this foundational VLM paradigm.

093 **Prompt tuning for VLMs.** Prompt tuning (Zhou et al., 2022b; Khattak et al., 2023a; Chen et al.,  
 094 2025; Zhou et al., 2022a; Khattak et al., 2023b) has emerged as a popular lightweight model  
 095 adaptation technique, particularly in VLMs. Unlike conventional fine-tuning approaches that involve  
 096 updating all parameters of a pre-trained model, prompt tuning focuses on optimizing only the  
 097 prompt representations, which significantly reduces computational overhead. CoOp (Zhou et al.,  
 098 2022b) introduces an efficient adaptation strategy by optimizing learnable prompt vectors that replace  
 099 manually crafted textual prompts (e.g., “a photo of a panda”) in the textual branch of CLIP, marking a  
 100 significant milestone in the vision-language alignment. However, CoOp tends to overfit base classes  
 101 with limited generalization ability. To address this, several studies (Zhu et al., 2023; Yao et al., 2023;  
 102 Chen et al., 2024; Khattak et al., 2023b) use the knowledge captured in pre-trained CLIP through  
 103 hand-crafted prompts to enhance the generalization of learnable prompts. In this paper, our study  
 104 builds on prompt tuning of CoOp to explore robust prompt tuning for VLMs.

105 **Adversarial training of VLMs.** Adversarial attacks (Madry et al., 2018; Goodfellow et al., 2014)  
 106 fool models by adding imperceptible perturbations to input images. Several studies (Touvron et al.,  
 107 2023; Zou et al., 2023; Schlarbmann & Hein, 2023; Zhao et al., 2023; Carlini et al., 2023) have shown  
 108 that VLMs are highly vulnerable to adversarial attacks. Among various defense strategies, adversarial  
 109 training (Madry et al., 2018; Zhang et al., 2019; Goodfellow et al., 2014; Zhang et al., 2020) is one  
 110 of the most effective defense methods against such attacks. It strengthens the model robustness by

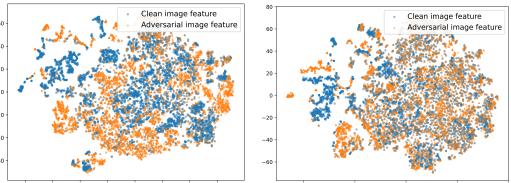


Figure 2: The left and right images illustrate the latent image features without and with the feature discretization, respectively.

incorporating adversarial examples into the training data and optimizing the model to classify both clean and adversarial examples correctly. Building on this idea, TeCoA (Mao et al., 2023) enhances adversarial robustness in CLIP through supervised adversarial fine-tuning. FARE (Schlarmann et al., 2024) introduces an unsupervised adversarial fine-tuning framework to obtain a robust CLIP, which preserves nominal performance while transferring robustness to downstream tasks. However, fine-tuning the entire model is not so economical. Alternatively, AdvPT (Zhang et al., 2024b) and APT (Li et al., 2024) combine prompt tuning with adversarial training, introducing adversarial prompt tuning for CLIP to enhance adversarial robustness through a PEFT way. Building on multi-modal prompt tuning (Khattak et al., 2023a), FAP (Zhou et al., 2024) improves the consistency of multimodal features and encourages differentiation in unimodal features between clean and adversarial examples, thereby enhancing adversarial robustness. Unlike previous adversarial prompt tuning methods that focus on learning robust prompts, the proposed DEFEAT method integrates VQ-VAE to mitigate the impact of adversarial attacks, achieving a better trade-off between robustness and accuracy.

### 3 METHODOLOGY

In this section, we begin with a brief overview of CLIP, adversarial attacks and training in CLIP, and VQ-VAE. Then we analyze the impact of VQ-VAE on visual adversarial perturbation. After that, we introduce the proposed DEFEAT in detail.

#### 3.1 PRELIMINARIES

**CLIP.** CLIP facilitates zero-shot image classification by evaluating the similarity between visual and textual embeddings. It comprises a visual encoder and a text encoder. Let  $\mathbf{I}$  represent the class feature extracted by CLIP’s visual encoder for an image, which serves as the image representation and is generated by a class token.  $\mathbf{t} = \{\mathbf{t}_c\}_{c=1}^C$  denotes a collection of textual embeddings produced by the text encoder for textual prompts, where each  $\mathbf{t}_c$  corresponds to the textual embedding for class  $c$ , and  $C$  is the number of classes. Those prompts typically follow the format “a photo of a [CLS]”, where “[CLS]” is a placeholder for class tokens that are specific class names like “panda”, “dog”, and “bird”. Based on those embeddings, the probability of classifying an image to belong to a class can be calculated as  $p^{\text{clip}}(y|\mathbf{x}) = \frac{\exp(f_y(\mathbf{x})/\tau)}{\sum_{c=1}^C \exp(f_c(\mathbf{x})/\tau)}$ , where  $\tau$  is the temperature parameter,  $f_c(\mathbf{x}) = \cos(\mathbf{t}_c, \mathbf{I})$ , and  $\cos(\cdot, \cdot)$  denotes the cosine similarity.

**Adversarial attack and training in CLIP.** An adversarial attack (Madry et al., 2018) fools the model by learning an imperceptible perturbation  $\delta$  added to a clean example  $\mathbf{x}$ , generating an adversarial example  $\mathbf{x}_a$ . Specifically, in CLIP, the adversarial example  $\mathbf{x}_a$  is optimized by maximizing the image-to-text contrastive loss (Mao et al., 2023; Li et al., 2024; Zhou et al., 2024) as

$$\mathbf{x}_a = \arg \max_{\mathbf{x}_a} \mathcal{L}(\mathbf{x}_a, \mathbf{t}, \mathbf{y}), \text{ s.t. } \mathbf{x}_a = \mathbf{x} + \delta, \|\delta\|_p \leq \epsilon, \quad (1)$$

where  $\mathcal{L}(\mathbf{x}_a, \mathbf{t}, \mathbf{y}) = -\sum_{i=1}^C y_i \log p^{\text{clip}}(y_i|\mathbf{x}_a)$ , label  $\mathbf{y}$  is a one-hot vector with  $y_i$  equal to 1 when class  $i$  is the ground-truth label for  $\mathbf{x}$  and 0 otherwise,  $\epsilon$  denotes the perturbation size, and  $\|\cdot\|_p$  denotes the  $\ell_p$  norm of a vector. Here we focus on the  $\ell_\infty$  threat model (i.e.,  $p = \infty$ ).

After generating adversarial examples, Mao et al. (2023) employs adversarial training to optimize the parameters  $\theta$  of CLIP’s vision encoder, thereby developing a robust CLIP. The objective of adversarial training is formulated as  $\theta = \arg \min_{\theta} \mathcal{L}(\mathbf{x}_a, \mathbf{t}, \mathbf{y})$ .

**VQ-VAE.** VQ-VAE is a generative model based on discrete latent representations, with its core mechanism to convert continuous representations into discrete codes. Specifically, VQ-VAE consists of an encoder  $E$ , a decoder  $D$ , and a codebook  $\mathbf{e}$ . The encoder maps input data  $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$  into a continuous latent vector  $\mathbf{z}_e = E(\mathbf{x}) \in \mathbb{R}^{n \times d}$ , where  $d$  is the dimensionality of the latent embedding,  $n$  is the number of latent positions after encoding. The learnable codebook is defined as  $\mathbf{e} = \{\mathbf{e}_k\}_{k=1}^K$ , where  $K$  is the size of the discrete latent space. The latent vector  $\mathbf{z}_e$  is then quantized into a discrete latent representation,  $\mathbf{z}_q \in \mathbb{R}^{n \times d}$ . This is achieved by replacing each vector in  $\mathbf{z}_e$  with its nearest neighbor from the codebook:

$$\mathbf{z}_q^{(i)} = \arg \min_{\mathbf{e}_k \in \mathbf{e}} \|\mathbf{z}_e^{(i)} - \mathbf{e}_k\|_2, \forall i \in \{1, \dots, n\}, \quad (2)$$

where  $\mathbf{z}_q^{(i)}$  and  $\mathbf{z}_e^{(i)}$  are the vectors at the  $i$ -th position of  $\mathbf{z}_q$  and  $\mathbf{z}_e$ , respectively, and  $\|\cdot\|_2$  denotes the  $\ell_2$  norm. The quantized discrete representation  $\mathbf{z}_q$  is then used to reconstruct the input by the decoder

162 as  $\hat{\mathbf{x}} = D(\mathbf{z}_q)$ . Since there is no real gradient defined for Eq. (2), VQ-VAE copies gradients from the  
 163 decoder input  $\mathbf{z}_q$  to the encoder output  $\mathbf{z}_e$ . The entire training objective in VQ-VAE is formulated as  
 164

$$165 \quad \mathcal{L}_{\text{VQ-VAE}}(\mathbf{x}) = \alpha \underbrace{\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2}_{\text{reconstruction loss}} + \beta \underbrace{\|\text{sg}[\mathbf{z}_e] - \mathbf{z}_q\|_2^2}_{\text{codebook loss}} + \gamma \underbrace{\|\mathbf{z}_e - \text{sg}[\mathbf{z}_q]\|_2^2}_{\text{commitment loss}}, \quad (3)$$

167 where  $\text{sg}[\cdot]$  is the stop gradient operator. By following the original implementation (van den Oord  
 168 et al., 2017),  $\gamma$  is set to be  $0.25\beta$ . Since VQ-VAE is an end-to-end, trainable, and commonly used  
 169 discretization model, we use it to discretize and reconstruct latent image features.  
 170

### 171 3.2 MITIGATING EFFECT OF VQ-VAE ON VISUAL ADVERSARIAL PERTURBATION

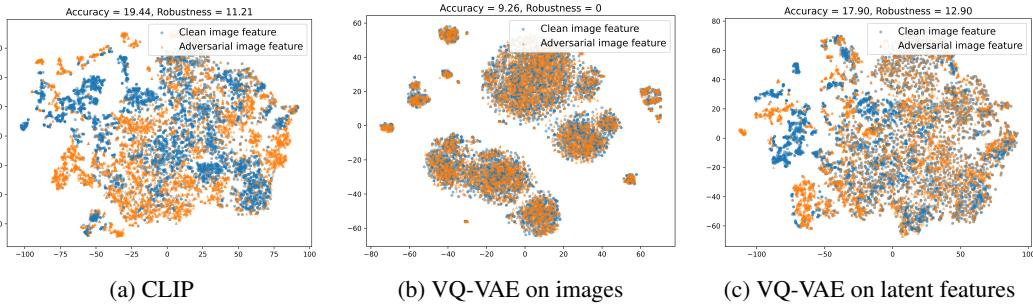


Figure 3: Visualization of clean and adversarial examples on the EuroSAT dataset. The VQ-VAE is trained with 16-shot samples per class under the adversarial few-shot classification setting.

In this section, we explore how VQ-VAE mitigates visual adversarial perturbations. We first visualize the clean and adversarial image features output by the CLIP’s image encoder in Figure 3a. We observe that imperceptible perturbations ( $\epsilon = 4/255$ ) cause a significant shift of latent features between clean and adversarial examples, leading to a sharp decline in classification accuracy for adversarial samples (i.e., from 19.44% to 11.21%).

Inspired by the quantization process of VQ-VAE, we hypothesize that the complex discretization process could mitigate adversarial attacks by reducing the distribution shift between clean and adversarial examples. To validate this hypothesis, we integrate VQ-VAE as a defense module within the VLM framework against adversarial attacks. Firstly, we attempt to use VQ-VAE to reconstruct input images as did in (Mao et al., 2022). As shown in Figure 3b, the shift of latent features between clean and adversarial examples is significantly reduced, demonstrating that the discretization process can suppress adversarial perturbations. However, pixel-level reconstruction with VQ-VAE introduces significant information loss and structural distortions that compromise semantic content preservation when training data is limited. To see that, we can find that the feature distributions in Figure 3b differ greatly from the original image feature distributions shown in Figure 3a. The poor performance (i.e., Accuracy=9.26%, Robustness=0%) could verify this and also make this approach impractical.

Alternatively, we use VQ-VAE to reconstruct the latent feature output by ViT. As shown in Figure 3c, we find that the distribution shift of latent features between adversarial and clean examples is significantly reduced. Moreover, the distribution of latent features after VQ-VAE (shown in Figure 3c) is more similar to the feature distribution by the CLIP (shown in Figure 3a). This suggests that reconstructing latent features can also effectively suppress adversarial perturbations and achieve better performance (i.e., Accuracy=17.90%, Robustness=12.90%). One possible reason is that when adversarial perturbations are introduced into clean image features, the process of finding the nearest codes often makes both clean and adversarially perturbed features be assigned to the same discrete representation (i.e., codes). Hence, the discretization of VQ-VAE effectively neutralizes the perturbations within its discrete representation space, thus hindering the updates of adversarial perturbations and diminishing the distributional differences between clean and adversarial examples.

### 3.3 DEFEAT: DISCRETE LATENT FEATURE BASED ADVERSARIAL TRAINING

Motivated by the above observation, we propose the DEFEAT method for CLIP-based VLMs. Besides a frozen image encoder and a frozen text encoder in CLIP, DEFEAT consists of a learnable prompt and a perturbation discrete shield module. The architecture of DEFEAT is illustrated in Figure 4.

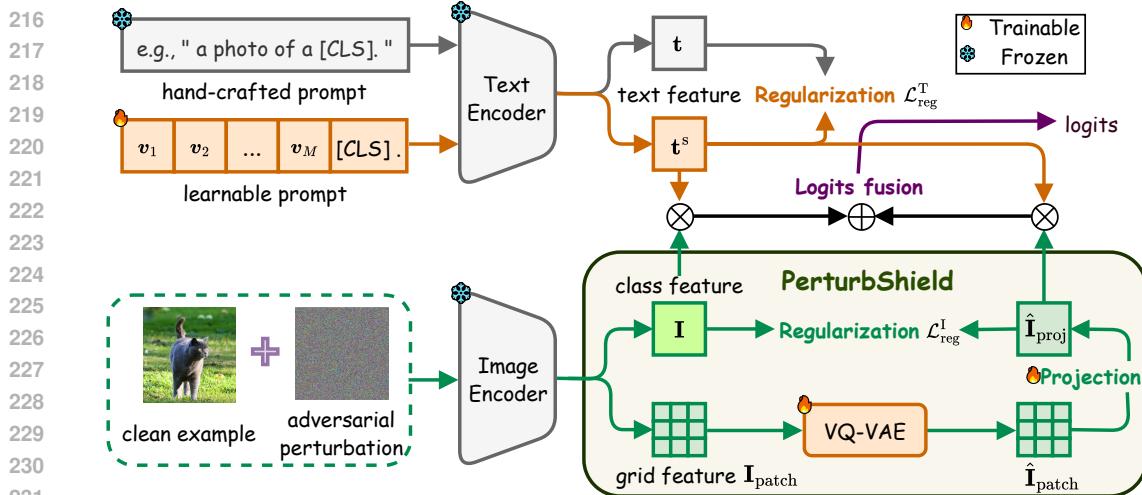


Figure 4: Overview of the proposed DEFEAT method. Note that both image and text encoders remain frozen, with only the learnable prompt and the perturbation discrete shield module being trained. In our experiments, we exclusively use adversarial examples for training.

**Prompt tuning for text encoder.** We introduce learnable prompt vectors  $\mathbf{v} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_M\}$  to replace the textual prompt (e.g., “a photo of a”), where each  $\mathbf{v}_i$  has the same dimension as the word embeddings. The learnable prompt is subsequently constructed by concatenating  $\mathbf{v}$  with the class token. Formally, the learnable prompts are formed as  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_M, [\text{CLS}]\}$ . Let  $\mathbf{t}^s = \{\mathbf{t}_c^s\}_{c=1}^C$  be a set of textual embeddings generated by the text encoder for the learnable prompts, where each  $\mathbf{t}_c^s$  corresponds to the textual embedding of learnable prompt for class  $c$ . The probability of classifying an image can be calculated as

$$p^s(y|\mathbf{x}) = \frac{\exp(f_y^s(\mathbf{x})/\tau)}{\sum_{c=1}^C \exp(f_c^s(\mathbf{x})/\tau)}, \quad (4)$$

where  $f_c^s(\mathbf{x}) = \cos(\mathbf{t}_c^s, \mathbf{I})$  and recall that  $\mathbf{I}$  denotes the class feature.

**Perturbation discrete shield (PerturbShield).** As analyzed in Section 3.2, we propose integrating VQ-VAE as a defense module within the prompt tuning framework to mitigate the impact of visual adversarial perturbations. Given the outstanding performance of the Vision Transformer (ViT) (Dosovitskiy et al., 2020) in visual tasks, we use CLIP’s ViT as the image encoder. The last Transformer layer of ViT can output two types of image features: the class feature  $\mathbf{I} \in \mathbb{R}^{d_v}$  generated by the [class] token and the grid feature  $\mathbf{I}_{\text{patch}} \in \mathbb{R}^{N \times d_v}$  from image patches, where  $d_v$  denotes the dimension (i.e., 512 for CLIP) and  $N$  is the number of patches.

Since CLIP uses the class feature  $\mathbf{I}$  to represent an image, we initially consider reconstructing  $\mathbf{I}$  with VQ-VAE. However, this approach has some limitations. That is, the class feature essentially functions as a global patch representation, meaning that we would be restricted to using only a single code from the codebook to represent the entire image. However, the representation capacity of a single code is far inferior to that of multiple codes combined, ultimately resulting in severe information loss. In this case, if an adversarial attack successfully maps the class feature to an incorrect code, it will lead to a wrong prediction.

To address this problem, we propose Perturbation Discrete Shield (PerturbShield), an end-to-end defense module that utilizes discrete representation learning to mitigate adversarial perturbations in VLMs. Specifically, PerturbShield operates through two designed stages: 1) grid-based discrete feature reconstruction; 2) semantic alignment projection, which are introduced in the following.

1) *Grid-based discrete feature reconstruction.* Instead of reconstructing the class feature, we reconstruct the grid feature  $\mathbf{I}_{\text{patch}}$  through the VQ-VAE. This process aims to represent the grid feature using a combination of multiple codes to preserve more information. Additionally, the discretization process can help mitigate minor perturbations and enhance the model’s robustness as analyzed in Section 3.2. The reconstructed grid feature can be obtained as  $\hat{\mathbf{I}}_{\text{patch}} = \text{VQ-VAE}(\mathbf{I}_{\text{patch}})$ , where  $\text{VQ-VAE}(\cdot)$  denotes a VQ-VAE and  $\hat{\mathbf{I}}_{\text{patch}} \in \mathbb{R}^{N \times d_v}$ .

270 2) *Semantic alignment projection*. We design a learnable matrix  $\mathbf{W}$  that performs two functionalities.  
 271 First, since  $\hat{\mathbf{I}}_{\text{patch}}$  is more robust than  $\mathbf{I}$ , we could use it to learn a robust representation of the  
 272 whole image by a linear transformation:  $\hat{\mathbf{I}}_{\text{proj}} = \mathbf{W} \cdot \hat{\mathbf{I}}_{\text{patch}}$ , where  $\mathbf{W} \in \mathbb{R}^{1 \times N}$  and  $\hat{\mathbf{I}}_{\text{proj}} \in \mathbb{R}^{d_v}$ .  
 273 Another functionality is to semantically align the transformed latent feature with CLIP's pre-trained  
 274 embedding space through a feature alignment regularization between  $\hat{\mathbf{I}}_{\text{proj}}$  and the class feature  $\mathbf{I}$  as  
 275

$$\mathcal{L}_{\text{reg}}^{\mathbf{I}} = \|\hat{\mathbf{I}}_{\text{proj}} - \mathbf{I}\|_1, \quad (5)$$

276 where  $\|\cdot\|_1$  denotes the  $\ell_1$  norm. These dual functionalities ensure that the adversarial defense  
 277 process remains aligned with the embedding space of the VLM. Then the prediction probability for  
 278 an input image  $\mathbf{x}$  can be computed as  
 279

$$p^{\text{vq}}(y|\mathbf{x}) = \frac{\exp(f_y^{\text{vq}}(\mathbf{x})/\tau)}{\sum_{c=1}^C \exp(f_c^{\text{vq}}(\mathbf{x})/\tau)}, \quad (6)$$

280 where  $f_c^{\text{vq}}(\mathbf{x}) = \cos(\mathbf{t}_c^s, \hat{\mathbf{I}}_{\text{proj}})$ .  
 281

282 **Logits fusion.** There is usually some information loss during the reconstruction process. Con-  
 283sequently, using  $\hat{\mathbf{I}}_{\text{proj}}$  alone to compute logits (i.e., the prediction probability) reduces the clean  
 284accuracy (e.g., from 19.44% to 17.90% in Figure 3). Since PerturbShield effectively mitigates  
 285adversarial attacks, we propose a logits fusion strategy that uses logits from  $\hat{\mathbf{I}}_{\text{proj}}$  and  $\mathbf{t}^s$  to counter  
 286adversarial attacks, while using logits from  $\mathbf{I}$  and  $\mathbf{t}^s$  to maintain the clean accuracy, thus achieving a  
 287better trade-off between the robustness and accuracy. Specifically, the former logits are defined in  
 288Eq. (6), and the latter ones are defined in Eq. (4). Then, the fused logits are calculated as  
 289

$$p(y|\mathbf{x}) = (p^{\text{vq}}(y|\mathbf{x}) + p^s(y|\mathbf{x}))/2. \quad (7)$$

290 **Prompt alignment regularization.** To enhance the generalization of learnable prompts and prevent  
 291overfitting (Zhou et al., 2022a), we use prompt alignment regularization for learnable prompts. This  
 292ensures the consistency between text features generated by learnable prompts and hand-crafted  
 293prompts, thereby improving adaptability while preserving CLIP's inherent zero-shot inference capa-  
 294bilities. Specifically, the prompt alignment regularization is formulated as  
 295

$$\mathcal{L}_{\text{reg}}^{\mathbf{T}} = \|\mathbf{t}^s - \mathbf{t}\|_1. \quad (8)$$

296 Moreover, we find that using hand-crafted prompts to regularize learnable prompts as in Eq. (8) can  
 297enhance the model's robustness, which is detailed in Section 4.3.  
 298

299 **Training objective.** Based on Eq. (7), the cross-entropy loss is defined as  $\mathcal{L}_{\text{ce}}(\mathbf{x}, \mathbf{t}^s, \mathbf{y}) =$   
 300 $-\sum_{i=1}^C y_i \log p(y_i|\mathbf{x})$ . To learn a robust prompt, Li et al. (2024); Zhou et al. (2024) combines  
 301prompt tuning with adversarial training. Following this paradigm, we use the learnable prompts  
 302updated in the previous epoch to generate adversarial examples dynamically in each epoch as  
 303

$$\mathbf{x}_a = \arg \max_{\mathbf{x}_a} \mathcal{L}_{\text{ce}}(\mathbf{x}_a, \mathbf{t}^s, \mathbf{y}), \text{ s.t. } \|\mathbf{x}_a - \mathbf{x}\|_p \leq \epsilon. \quad (9)$$

304 Finally, based on those adversarial examples, the overall training objective of the DEFEAT method is  
 305formulated as  
 306

$$\mathcal{L}(\mathbf{x}_a, \mathbf{t}^s, \mathbf{t}, \mathbf{y}) = \mathcal{L}_{\text{ce}}(\mathbf{x}_a, \mathbf{t}^s, \mathbf{y}) + \mathcal{L}_{\text{VQ-VAE}}(\mathbf{I}_{\text{patch}}) + \lambda \mathcal{L}_{\text{reg}}^{\mathbf{I}} + \mu \mathcal{L}_{\text{reg}}^{\mathbf{T}}. \quad (10)$$

307 Due to page limit, the whole algorithm for the proposed DEFEAT method is shown in Appendix A.  
 308

## 315 4 EXPERIMENTS

### 316 4.1 SETUPS

317 **Datasets.** By following APT (Li et al., 2024) and CoOp (Zhou et al., 2022b), experiments are  
 318conducted under three settings, including adversarial few-shot classification, adversarial cross-dataset  
 319generalization, and adversarial domain generalization. The first two settings are performed on 11  
 320image classification datasets, including ImageNet (Deng et al., 2009) and Caltech101 (Fei-Fei et al.,  
 3212004) for generic object classification, OxfordPets (Parkhi et al., 2012), StanfordCars (Krause et al.,  
 3222013), Flowers102 (Nilsback & Zisserman, 2008), Food101 (Bossard et al., 2014), and FGVCAircraft  
 323

(Maji et al., 2013) for fine-grained visual categorization, SUN397 (Xiao et al., 2010) for scene recognition, DTD (Cimpoi et al., 2014) for texture classification, EuroSAT (Helber et al., 2019) for satellite image classification, and UCF101 (Soomro et al., 2012) for action recognition. Adversarial domain generalization is conducted on the ImageNet dataset and its variants, including ImageNetV2 (Recht et al., 2019), ImageNet-Sketch (Wang et al., 2019), ImageNet-A (Hendrycks et al., 2021b), and ImageNet-R (Hendrycks et al., 2021a).

**Baselines.** To demonstrate the effectiveness of the proposed DEFEAT method, we compare it against three baseline methods, including zero-shot CLIP (Radford et al., 2021), the textual prompt tuning method (i.e., Adversarial Prompt Tuning (APT) (Li et al., 2024)), and the multi-modal prompt tuning method (i.e., Few-shot Adversarial Prompt learning (FAP) (Zhou et al., 2024)).

**Implementation details.** Our implementation is based on CoOp (Zhou et al., 2022b) and APT (Li et al., 2024). For all baseline methods, we maintain the same experimental settings (e.g., training epochs, training schedules, and data augmentation settings) as specified in their original implementations. All experiments are conducted with a ViT-B/32 CLIP model. Following APT (Li et al., 2024), the image encoder weights for all baselines were pre-trained using TeCoA (Mao et al., 2023). The length of the prompt vectors is fixed to 16. For the DEFEAT methods,  $\alpha$ ,  $\beta$ ,  $\lambda$ , and  $\mu$  are set to 0.5, 0.1, 10, and 20, respectively, across all experiments. For adversarial training and evaluation, we employ the PGD attack (Madry et al., 2018) under the  $\ell_\infty$  threat model. Following Li et al. (2024); Mao et al. (2023); Schlarbmann et al. (2024), two perturbation sizes,  $\epsilon = 1/255$  and  $\epsilon = 4/255$ , are used. During training, we use 3 steps with a step size of  $2\epsilon/3$ , and for evaluation, 100 steps with a step size of  $\epsilon/4$  and a random start. More details are provided in Appendix B.

## 4.2 MAIN RESULTS

In this section, we report experimental results under different settings.

**Adversarial few-shot classification.** In this scenario, we assess the model’s capability to develop robust representations from limited labeled data. Specifically, models are tuned using 1, 4, 16 shots per class and evaluated on the remaining samples. The average performance of DEFEAT and baselines on 11 datasets is shown in Table 1, where ‘H’ denotes the harmonic mean accuracy of robustness and accuracy, measuring the trade-off between these metrics. DEFEAT consistently outperforms zero-shot CLIP, even with 1-shot tuning, achieving significant improvements in accuracy, robustness, and ‘H’. For example, under  $\epsilon = 1/255$ , DEFEAT provides gains of 6.01%, 6.01%, and 6.15% in accuracy, robustness, and ‘H’, respectively. Furthermore, these improvements scale with the number of shots. With 16-shot tuning, DEFEAT achieves gains of 18.97% (17.31%), 27.79% (25.26%), and 24.39% (25.89%) in accuracy, robustness, and ‘H’ for  $\epsilon = 1/255$  (4/255), respectively.

Existing adversarial prompt tuning methods (i.e., APT and FAP) improve performance over CLIP, confirming the effect of combining adversarial training with prompt tuning. Compared to those methods, DEFEAT performs comparably to the current state-of-the-art method FAP in terms of the accuracy, while consistently outperforms both APT and FAP in terms of the robustness and ‘H’ across all shot and perturbation sizes. Notably, with 16-shot training under  $\epsilon = 4/255$ , DEFEAT surpasses APT (FAP) by 15.79% (16.23%) in robustness, and 13.22% (13.76%) in ‘H’, respectively. Detailed results for each dataset are provided in Appendix C.1.

Table 1: The average performance on the 11 datasets for different  $\epsilon$  and shots under the adversarial few-shot classification setting. ‘H’ denotes the harmonic mean accuracy.

| $\epsilon$ | Method | 1 shot       |              |              | 4 shots      |              |              | 16 shots     |              |              |
|------------|--------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|            |        | Acc.         | Rob.         | H            | Acc.         | Rob.         | H            | Acc.         | Rob.         | H            |
| 1/255      | CLIP   | 46.06        | 32.98        | 38.44        | 46.06        | 32.98        | 38.44        | 46.06        | 32.98        | 38.44        |
|            | APT    | 46.99        | 33.36        | 39.02        | 58.19        | 41.34        | 48.34        | 65.41        | 47.88        | 55.29        |
|            | FAP    | 52.72        | 37.15        | 43.59        | 59.77        | 42.21        | 49.48        | 66.40        | 48.86        | 56.30        |
|            | DEFEAT | 52.07        | <b>38.99</b> | <b>44.59</b> | 58.37        | <b>53.04</b> | <b>55.58</b> | 65.03        | <b>60.77</b> | <b>62.83</b> |
| 4/255      | CLIP   | 33.67        | 10.79        | 16.34        | 33.67        | 10.79        | 16.34        | 33.67        | 10.79        | 16.34        |
|            | APT    | 32.97        | 11.62        | 17.18        | 42.29        | 14.40        | 21.48        | 51.08        | 20.26        | 29.01        |
|            | FAP    | <b>38.16</b> | 13.44        | 19.88        | <b>44.63</b> | 15.34        | 22.83        | 50.50        | 19.82        | 28.47        |
|            | DEFEAT | 35.78        | <b>14.73</b> | <b>20.87</b> | 44.13        | <b>25.28</b> | <b>32.15</b> | <b>50.98</b> | <b>36.05</b> | <b>42.23</b> |

**Adversarial domain generalization.** We assess the model’s capability to generalize to out-of-distribution (OOD) data in the adversarial domain generalization setting. Specifically, models are tuned using 16-shot and 100-shot samples from each of the 1000 classes on ImageNet (source), and then evaluated on four different target domains (i.e., ImageNetV2, ImageNet-Sketch, ImageNet-A, and ImageNet-R). Table 2 shows that all baselines achieve comparable accuracy on ImageNet, as they use the same robust image encoder (Mao et al., 2023) pre-trained on the full ImageNet training set. Nevertheless, with 16-shot training, DEFEAT surpasses CLIP, APT, and FAP in terms of the robustness by 8.92%, 7.04%, and 7.0%, respectively. Across the four target domains, DEFEAT

achieves good robustness and ‘H’, outperforming the strongest baseline (i.e., APT) by 3.31% in robustness and 3.36% in ‘H’. Those results confirm that the robustness learned by DEFEAT can be effectively transferred to OOD data. Furthermore, as the number of training samples increases from 16 shots to 100 shots, DEFEAT demonstrates consistent improvements in terms of the robustness, achieving +3.31% on ImageNet and an average of +1.94% across the target domains, which highlights the scalability of our approach. In contrast, the two prompt tuning baselines (i.e., APT and FAP) show little to no improvement.

Table 2: Performance under adversarial domain generalization setting.  $\epsilon = 4/255$ .

| Method    | Source   |              |              | Target       |              |              |                 |              |              |              |             |             |             |              |              | Average      |              |              |              |
|-----------|----------|--------------|--------------|--------------|--------------|--------------|-----------------|--------------|--------------|--------------|-------------|-------------|-------------|--------------|--------------|--------------|--------------|--------------|--------------|
|           | ImageNet |              |              | ImageNet-V2  |              |              | ImageNet-Sketch |              |              | ImageNet-A   |             |             | ImageNet-R  |              |              | Average      |              |              |              |
|           | Acc.     | Rob.         | H            | Acc.         | Rob.         | H            | Acc.            | Rob.         | H            | Acc.         | Rob.        | H           | Acc.        | Rob.         | H            | Acc.         | Rob.         | H            |              |
| zero shot | CLIP     | 40.11        | 10.14        | 16.19        | 33.11        | 7.49         | 12.22           | 17.59        | 7.24         | 10.26        | 4.04        | 0.28        | 0.52        | 37.52        | 12.51        | 18.76        | 23.07        | 6.88         | 10.60        |
| 16 shots  | APT      | <b>41.06</b> | 12.02        | 18.60        | <b>33.67</b> | 9.09         | 14.32           | 18.22        | 7.87         | 10.99        | <b>4.19</b> | 0.36        | 0.66        | 37.04        | 13.29        | 19.56        | <b>23.28</b> | 7.65         | 11.52        |
|           | FAP      | 40.32        | 12.06        | 18.57        | 32.81        | 9.17         | 14.33           | 16.42        | 7.11         | 9.92         | 3.87        | 0.43        | 0.77        | 36.04        | 13.55        | 19.70        | 22.29        | 7.57         | 11.30        |
|           | DEFEAT   | 40.68        | <b>19.06</b> | <b>25.96</b> | 33.12        | <b>14.96</b> | <b>20.61</b>    | <b>18.27</b> | <b>9.87</b>  | <b>12.82</b> | 3.97        | <b>0.88</b> | <b>1.44</b> | <b>37.35</b> | <b>18.14</b> | <b>24.42</b> | 23.18        | <b>10.96</b> | <b>14.88</b> |
| 100 shots | APT      | 40.35        | 12.11        | 18.63        | 33.38        | 9.19         | 14.41           | <b>18.86</b> | 8.16         | 11.39        | <b>4.21</b> | 0.40        | 0.73        | <b>37.81</b> | 13.70        | 20.11        | <b>23.57</b> | 7.86         | 11.79        |
|           | FAP      | <b>40.76</b> | 12.35        | 18.96        | <b>33.60</b> | 9.14         | 14.37           | 16.41        | 6.98         | 9.79         | 3.65        | 0.39        | 0.70        | 36.36        | 13.60        | 19.80        | 22.51        | 7.53         | 11.28        |
|           | DEFEAT   | 40.23        | <b>22.71</b> | <b>29.03</b> | 32.77        | <b>16.44</b> | <b>21.90</b>    | 18.49        | <b>11.53</b> | <b>14.20</b> | 4.20        | <b>1.37</b> | <b>2.07</b> | 37.34        | <b>21.45</b> | <b>27.25</b> | 23.20        | <b>12.70</b> | <b>16.41</b> |

**Adversarial cross-dataset generalization.** We assess the model’s zero-shot adversarial robustness in the adversarial cross-dataset generalization setting. Models are tuned using 16-shot and 100-shot samples from each of the 1000 classes on ImageNet and evaluated on the other 10 datasets. Table 3 shows that DEFEAT achieves competitive accuracy and the highest robustness and ‘H’ compared to all baseline methods across all datasets, including both source and target domains. Specifically, with 16-shot tuning, DEFEAT outperforms the strongest baseline (i.e., FAP) by an average of 5.61% in terms of robustness and 5.37% in terms of ‘H’. Similar to the adversarial domain generalization setting, APT and FAP show no significant performance improvements when the number of training samples increases from 16 shots to 100 shots. In contrast, DEFEAT demonstrates scalability, achieving gains of 3.65% in terms of robustness on ImageNet and 3.18% in terms of average robustness on 10 target datasets. Those results demonstrate the good zero-shot adversarial robustness of DEFEAT.

Table 3: Performance under adversarial cross-dataset generalization setting.  $\epsilon = 4/255$ .

| 16 shots  | Method | Source       |              |              | Target       |              |              |             |              |              |              |              |              |      |  |  |
|-----------|--------|--------------|--------------|--------------|--------------|--------------|--------------|-------------|--------------|--------------|--------------|--------------|--------------|------|--|--|
|           |        | ImageNet     | Caltech101   | OxfordPets   | StanfordCars | Flowers102   | Food101      | FGVC        | Aircraft     | SUN397       | DTD          | EuroSAT      | UCF101       | Avg. |  |  |
| Acc.      | CLIP   | 40.11        | 78.78        | 66.26        | 10.32        | 30.13        | 23.52        | <b>7.17</b> | <b>33.24</b> | <b>24.29</b> | <b>19.56</b> | <b>36.98</b> | <b>33.03</b> |      |  |  |
|           | APT    | <b>41.06</b> | <b>79.63</b> | 64.92        | <b>11.55</b> | 27.24        | <b>23.58</b> | 4.77        | 30.87        | 22.22        | 16.31        | 33.68        | 31.48        |      |  |  |
|           | FAP    | 40.32        | 79.19        | 63.67        | 9.64         | 30.25        | 22.82        | 5.43        | 31.92        | 23.23        | 16.79        | 32.25        | 31.52        |      |  |  |
|           | DEFEAT | 40.68        | 78.22        | <b>66.69</b> | 9.40         | <b>31.10</b> | 22.01        | 6.09        | 32.28        | 23.94        | 17.95        | 33.15        | 32.08        |      |  |  |
| Rob.      | CLIP   | 10.14        | 43.61        | 15.56        | 0.99         | 8.93         | 3.27         | 0.36        | 6.20         | 11.35        | 11.22        | 7.01         | 10.85        |      |  |  |
|           | APT    | 12.02        | 46.86        | 21.10        | 1.49         | 8.89         | 3.75         | 0.60        | 6.85         | 9.93         | 11.15        | 7.56         | 11.82        |      |  |  |
|           | FAP    | 12.06        | 46.25        | 20.71        | 1.55         | 9.74         | 3.67         | 0.66        | 7.15         | 10.82        | 11.38        | 7.45         | 11.94        |      |  |  |
|           | DEFEAT | <b>19.06</b> | <b>54.32</b> | <b>31.97</b> | <b>3.84</b>  | <b>18.80</b> | <b>8.72</b>  | <b>2.34</b> | <b>13.47</b> | <b>15.60</b> | <b>12.89</b> | <b>13.59</b> | <b>17.55</b> |      |  |  |
| H         | CLIP   | 16.19        | 56.14        | 25.20        | 1.81         | 13.78        | 5.74         | 0.69        | 10.45        | 15.47        | 14.26        | 11.79        | 16.33        |      |  |  |
|           | APT    | 18.60        | 59.00        | 31.85        | 2.64         | 13.41        | 6.47         | 1.07        | 11.21        | 13.73        | 13.25        | 12.35        | 17.18        |      |  |  |
|           | FAP    | 18.57        | 58.40        | 31.25        | 2.67         | 14.74        | 6.32         | 1.18        | 11.68        | 14.76        | 13.57        | 12.10        | 17.32        |      |  |  |
|           | DEFEAT | <b>25.96</b> | <b>64.12</b> | <b>43.22</b> | <b>5.45</b>  | <b>23.43</b> | <b>12.49</b> | <b>3.38</b> | <b>19.01</b> | <b>18.89</b> | <b>15.00</b> | <b>19.28</b> | <b>22.69</b> |      |  |  |
| 100 shots | Method | ImageNet     | Caltech101   | OxfordPets   | StanfordCars | Flowers102   | Food101      | FGVC        | Aircraft     | SUN397       | DTD          | EuroSAT      | UCF101       | Avg. |  |  |
|           | CLIP   | 40.11        | 78.78        | <b>66.26</b> | 10.32        | 30.13        | 23.52        | <b>7.17</b> | <b>33.24</b> | 24.29        | 19.56        | <b>36.98</b> | <b>33.03</b> |      |  |  |
|           | APT    | 40.35        | <b>80.28</b> | 64.60        | <b>11.76</b> | 29.03        | <b>23.87</b> | 5.85        | 32.91        | <b>25.71</b> | 14.17        | 34.87        | 32.31        |      |  |  |
|           | FAP    | <b>40.76</b> | 75.54        | 64.30        | 10.21        | 29.76        | 23.37        | 3.93        | 30.49        | 24.94        | <b>22.01</b> | 31.51        | 31.61        |      |  |  |
|           | DEFEAT | 40.23        | 76.63        | 65.74        | 9.46         | <b>30.25</b> | 20.02        | 6.51        | 31.31        | 23.94        | 16.67        | 31.69        | 31.22        |      |  |  |
| Rob.      | CLIP   | 10.14        | 43.61        | 15.56        | 0.99         | 8.93         | 3.27         | 0.36        | 6.20         | 11.35        | 11.22        | 7.01         | 10.85        |      |  |  |
|           | APT    | 12.11        | 46.82        | 21.18        | 1.52         | 9.66         | 3.76         | 0.69        | 7.20         | 12.41        | 11.06        | 8.33         | 12.26        |      |  |  |
|           | FAP    | 12.35        | 44.50        | 20.47        | 1.47         | 9.50         | 3.89         | 0.39        | 6.91         | 10.17        | 12.19        | 7.35         | 11.68        |      |  |  |
|           | DEFEAT | <b>22.71</b> | <b>59.35</b> | <b>42.95</b> | <b>5.58</b>  | <b>21.52</b> | <b>11.91</b> | <b>3.78</b> | <b>16.29</b> | <b>16.49</b> | <b>12.98</b> | <b>16.44</b> | <b>20.73</b> |      |  |  |
| H         | CLIP   | 16.19        | 56.14        | 25.20        | 1.81         | 13.78        | 5.74         | 0.69        | 10.45        | 15.47        | 14.26        | 11.79        | 16.33        |      |  |  |
|           | APT    | 18.63        | 59.15        | 31.90        | 2.69         | 14.50        | 6.50         | 1.23        | 11.82        | 16.74        | 12.42        | 13.45        | 17.78        |      |  |  |
|           | FAP    | 18.96        | 56.01        | 31.05        | 2.57         | 14.40        | 6.67         | 0.71        | 11.27        | 14.45        | <b>15.69</b> | 11.92        | 17.06        |      |  |  |
|           | DEFEAT | <b>29.03</b> | <b>66.89</b> | <b>51.96</b> | <b>7.02</b>  | <b>25.15</b> | <b>14.94</b> | <b>4.78</b> | <b>21.43</b> | <b>19.53</b> | 14.60        | <b>21.65</b> | <b>24.92</b> |      |  |  |

### 4.3 ABLATION STUDIES

In this section, we conduct ablation studies to analyze the effectiveness of the PerturbShield module, logits fusion strategy, and hyperparameter (i.e.,  $\mu$ ). Experiments are conducted using 16 examples per class on the EuroSAT dataset under the adversarial few-shot classification setting.

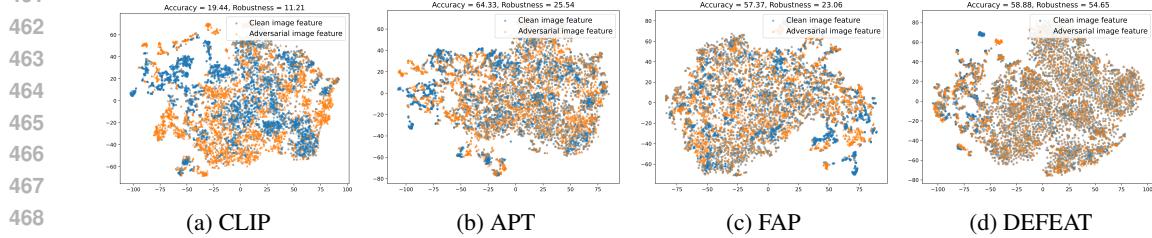
**Effects of components.** Table 4 shows the ablation study for PerturbShield and logits fusion strategy. Due to information loss during the reconstruction process, using PerturbShield alone without logits

432 fusion results in limited performance improvements on CLIP and can negatively impact DEFEAT.  
 433 Unlike CLIP, which computes logits using hand-crafted prompts, DEFEAT relies on learnable  
 434 prompts. The information loss can impair prompt learning, thereby causing adverse effects.  
 435

436 By combining PerturbShield with the logits fusion  
 437 strategy, adversarial attacks can be mitigated effectively  
 438 while maintaining good clean accuracy, thus  
 439 achieving a better trade-off between accuracy and  
 440 robustness. Specifically, CLIP combining Perturb-  
 441 Shield and logits fusion strategy outperforms the orig-  
 442 inal CLIP by 9.91% in ‘H’ (24.13% v.s. 14.22%), and  
 443 DEFEAT with this combination outperforms its coun-  
 444 terpart without it by 20.38% in ‘H’ (56.69% v.s. 36.31%). Therefore, we combine PerturbShield with  
 445 the logits fusion strategy.

446 **Effect of  $\mu$ .**  $\mu$  controls the weight of  $\mathcal{L}_{\text{reg}}^T$ , which apply regularization  
 447 from learnable prompts to hand-crafted prompts. According to Figure 5,  
 448 we can see that applying regularization (i.e.,  $\mu > 0$ ) of learnable prompts  
 449 can enhance the model’s robustness. As  $\mu$  increases, the robustness of  
 450 DEFEAT improves, achieving the best overall performance at  $\mu = 20$ .  
 451 Therefore, we set  $\mu = 20$  in our experiments.

452 **How does DEFEAT enhance adversarial robustness?** To answer this  
 453 question, we visualize the image feature (i.e., class feature) output by the  
 454 image encoder using t-SNE. Figure 6 shows that CLIP exhibits a significant  
 455 distribution shift between clean and adversarial examples, while adversarial  
 456 prompt tuning methods (i.e., APT and FAP) reduce this shift. Notably,  
 457 the image encoder remains frozen during training and testing for all methods, demonstrating that  
 458 adversarial prompt tuning can mitigate the backward adversarial gradient, thereby enhancing model  
 459 robustness. DEFEAT not only incorporates adversarial prompt tuning but also further mitigates  
 460 the backward adversarial gradient through PerturbShield. So the distribution shift of latent features  
 461 between clean and adversarial examples is further reduced compared to APT and FAP, providing  
 462 additional evidence of the effectiveness of the proposed DEFEAT method.



463 Figure 6: Visualization of clean and adversarial examples on the EuroSAT dataset. Models are trained  
 464 with 16-shot samples per class.  $\epsilon = 4/255$ .

465 *Moreover, due to page limit, more experiments can be found in Appendix C, including evaluations  
 466 against stronger attacks and a custom adaptive attack designed for DEFEAT, generalization to  
 467 alternative backbones, detailed hyperparameter analysis, and computational cost analysis.*

## 471 5 CONCLUSION

472 This work presents a robust prompt tuning framework for VLMs. We begin by analyzing the  
 473 mitigating effect of feature discretization on visual adversarial perturbation. Building on this, we  
 474 propose the **Discrete LatEnt FeaturE** based Adversarial Training (**DEFEAT**) method, which ablates  
 475 adversarial attacks in the feature representation. Specifically, DEFEAT introduces a perturbation  
 476 discrete shield module that reconstructs discrete latent features, and designs a logits fusion strategy  
 477 to improve the trade-off between robustness and accuracy. Moreover, DEFEAT integrates prompt  
 478 tuning with adversarial training and applies prompt alignment regularization, further enhancing the  
 479 adversarial robustness. Extensive experiments on 15 datasets demonstrate that DEFEAT achieves  
 480 state-of-the-art performance.

Table 4: Ablation study.  $\epsilon = 4/255$ .

| Baselines | Adversarial prompt tuning | PerturbShield | Logits fusion | Acc.  | Rob.  | H              |
|-----------|---------------------------|---------------|---------------|-------|-------|----------------|
| CLIP      | ✗                         | ✗             | ✗             | 19.44 | 11.21 | 14.22          |
|           | ✗                         | ✓             | ✗             | 17.90 | 12.90 | 14.99 (+0.77)  |
|           | ✗                         | ✓             | ✓             | 25.53 | 22.88 | 24.13 (+9.91)  |
| DEFEAT    | ✓                         | ✗             | ✗             | 64.12 | 25.33 | 36.31          |
|           | ✓                         | ✓             | ✗             | 22.21 | 17.48 | 19.56 (-16.75) |
|           | ✓                         | ✓             | ✓             | 58.88 | 54.65 | 56.69 (+20.38) |

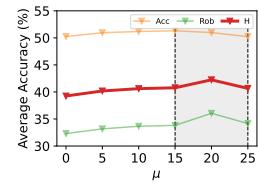


Figure 5: The average performance of DEFEAT on all datasets w.r.t. different  $\mu$  values.

486 ETHICS STATEMENT  
487488 This paper presents work whose goal is to advance the field of Machine Learning by improving the  
489 safety and reliability of AI systems. The authors have read and comply with the ICLR Code of Ethics.  
490 The research did not involve human subjects, animal experiments, or personally identifiable data. All  
491 experiments were conducted on publicly available benchmarks and open-source models. We have  
492 carefully considered the broader impacts and believe that this work poses no foreseeable risks of  
493 harm while contributing to the development of robust and secure vision-language models.  
494495 REPRODUCIBILITY STATEMENT  
496497 The authors have made significant efforts to ensure the reproducibility of results. Section 4.1  
498 details the experimental setup, including datasets, model configurations, and hyperparameter settings.  
499 Additional ablations in Section 4.3 and Appendix C.2 further analyze the effect of each module and  
500 hyperparameters. During the reviewing process, the source code is supplied anonymously as part  
501 of the supplementary materials. Additionally, upon the acceptance of the paper, this code will be  
502 publicly released.  
503504 REFERENCES  
505506 Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of  
507 security: Circumventing defenses to adversarial examples. In *International conference on machine  
508 learning*, pp. 274–283. PMLR, 2018.509 Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative compo-  
510 nents with random forests. In *ECCV*, 2014.512 Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017  
513 ieee symposium on security and privacy (sp)*, pp. 39–57. Ieee, 2017.514 Nicholas Carlini, Milad Nasr, Christopher A Choquette-Choo, Matthew Jagielski, Irena Gao, Pang  
515 Wei W Koh, Daphne Ippolito, Florian Tramer, and Ludwig Schmidt. Are aligned neural networks  
516 adversarially aligned? *Advances in Neural Information Processing Systems*, 36:61478–61500,  
517 2023.519 Yang Chen, Shuai Fu, and Yu Zhang. Mopd: Mixture-of-prompts distillation for vision-language  
520 models, 2024. URL <https://arxiv.org/abs/2412.19087>.521 Yang Chen, Yanbin Wei, Ke Jin, Yi Kong, James Kwok, and Yu Zhang. Alternating training-based  
522 label smoothing enhances prompt generalization, 2025. URL <https://arxiv.org/abs/2508.17846>.525 Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. De-  
526 scribing textures in the wild. In *CVPR*, 2014.527 Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble  
528 of diverse parameter-free attacks. In *International conference on machine learning*, pp. 2206–2216.  
529 PMLR, 2020.531 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
532 hierarchical image database. In *CVPR*, 2009.533 Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin  
534 Chen, Chi-Min Chan, Weize Chen, et al. Parameter-efficient fine-tuning of large-scale pre-trained  
535 language models. *Nature Machine Intelligence*, 5(3):220–235, 2023.537 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
538 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An  
539 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint  
arXiv:2010.11929*, 2020.

540 Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training  
 541 examples: An incremental bayesian approach tested on 101 object categories. In *CVPRW*, 2004.  
 542

543 Shizhan Gong, Haoyu LEI, Qi Dou, and Farzan Farnia. Boosting the visual interpretability of CLIP via  
 544 adversarial fine-tuning. In *The Thirteenth International Conference on Learning Representations*,  
 545 2025.

546 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial  
 547 examples. *arXiv preprint arXiv:1412.6572*, 2014.

548 Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset  
 549 and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected  
 550 Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.

551

552 Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul  
 553 Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical  
 554 analysis of out-of-distribution generalization. In *ICCV*, 2021a.

555 Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial  
 556 examples. In *CVPR*, 2021b.

557

558 Md Zarif Hossain and Ahmed Imteaj. Sim-clip: Unsupervised siamese adversarial fine-tuning for  
 559 robust and semantically-rich vision-language models. *arXiv preprint arXiv:2407.14971*, 2024.

560 Gabriel Ilharco, Mitchell Wortsman, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar,  
 561 Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, et al. Openclip. *Zenodo*,  
 562 2021.

563 Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung,  
 564 Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with  
 565 noisy text supervision. In *Proceedings of the International Conference on Machine Learning*, pp.  
 566 4904–4916, 2021.

567

568 Muhammad Uzair Khattak, Hanoona Rasheed, Muhammad Maaz, Salman Khan, and Fahad Shahbaz  
 569 Khan. Maple: Multi-modal prompt learning. In *CVPR*, pp. 19113–19122, 2023a.

570 Muhammad Uzair Khattak, Syed Talal Wasim, Muzammal Naseer, Salman Khan, Ming-Hsuan  
 571 Yang, and Fahad Shahbaz Khan. Self-regulating prompts: Foundational model adaptation without  
 572 forgetting. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp.  
 573 15190–15200, 2023b.

574

575 Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained  
 576 categorization. In *CVPRW*, 2013.

577

578 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-  
 579 training for unified vision-language understanding and generation. In *International Conference on  
 Machine Learning*, 2022.

580 Lin Li, Haoyan Guan, Jianing Qiu, and Michael Spratling. One prompt word is enough to boost  
 581 adversarial robustness for pre-trained vision-language models. In *Proceedings of the IEEE/CVF  
 582 Conference on Computer Vision and Pattern Recognition*, pp. 24408–24419, 2024.

583

584 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Advances  
 585 in Neural Information Processing Systems*, 2023.

586 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.  
 587 Towards deep learning models resistant to adversarial attacks. In *International Conference on  
 588 Learning Representations*, 2018.

589

590 Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained  
 591 visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.

592

593 Chengzhi Mao, Scott Geng, Junfeng Yang, Xin Wang, and Carl Vondrick. Understanding zero-  
 shot adversarial robustness for large-scale models. In *The Eleventh International Conference on  
 Learning Representations*, 2023.

594 Xiaofeng Mao, Yuefeng Chen, Ranjie Duan, Yao Zhu, Gege Qi, Xiaodan Li, Rong Zhang, Hui Xue,  
 595 et al. Enhance the visual representation via discrete adversarial training. *Advances in Neural*  
 596 *Information Processing Systems*, 35:7520–7533, 2022.

597

598 Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number  
 599 of classes. In *Indian conference on computer vision, graphics & image processing*, 2008.

600 Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *CVPR*,  
 601 2012.

602 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
 603 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
 604 models from natural language supervision. In *International Conference on Machine Learning*,  
 605 2021.

606

607 Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers  
 608 generalize to imagenet? In *ICML*, 2019.

609 Christian Schlarmann and Matthias Hein. On the adversarial robustness of multi-modal foundation  
 610 models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp.  
 611 3677–3685, 2023.

612

613 Christian Schlarmann, Naman Deep Singh, Francesco Croce, and Matthias Hein. Robust clip:  
 614 Unsupervised adversarial fine-tuning of vision embeddings for robust large vision-language models.  
 615 In *International Conference on Machine Learning*, pp. 43685–43704. PMLR, 2024.

616 Lijun Sheng, Jian Liang, Zilei Wang, and Ran He. R-pt: Improving adversarial robustness of  
 617 vision-language models through test-time prompt tuning. In *Proceedings of the Computer Vision  
 618 and Pattern Recognition Conference*, pp. 29958–29967, 2025.

619

620 Manli Shu, Weili Nie, De-An Huang, Zhiding Yu, Tom Goldstein, Anima Anandkumar, and  
 621 Chaowei Xiao. Test-time prompt tuning for zero-shot generalization in vision-language  
 622 models. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Ad-*  
 623 *vances in Neural Information Processing Systems*, volume 35, pp. 14274–14289. Curran Asso-  
 624 ciates, Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/5bf2b802e24106064dc547ae9283bb0c-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/5bf2b802e24106064dc547ae9283bb0c-Paper-Conference.pdf).

625

626 Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions  
 627 classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.

628

629 Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training  
 630 techniques for clip at scale, 2023. URL <https://arxiv.org/abs/2303.15389>.

631

632 Quan Sun, Jinsheng Wang, Qiying Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, and Xinlong Wang.  
 633 Eva-clip-18b: Scaling clip to 18 billion parameters, 2024. URL <https://arxiv.org/abs/2402.04252>.

634

635 Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow,  
 636 and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.

637

638 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
 639 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

640

641 Aaron van den Oord, Oriol Vinyals, and koray kavukcuoglu. Neural discrete representation learning.  
 642 In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett  
 (eds.), *Advances in Neural Information Processing Systems*, volume 30, 2017.

643

644 Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations  
 645 by penalizing local predictive power. In *NIPS*, 2019.

646

647 Xin Wang, Kai Chen, Jiaming Zhang, Jingjing Chen, and Xingjun Ma. Tapt: Test-time adversarial  
 648 prompt tuning for robust inference in vision-language models. In *Proceedings of the Computer  
 Vision and Pattern Recognition Conference*, pp. 19910–19920, 2025.

648 Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database:  
 649 Large-scale scene recognition from abbey to zoo. In *CVPR*, 2010.  
 650

651 Hantao Yao, Rui Zhang, and Changsheng Xu. Visual-language prompt tuning with knowledge-guided  
 652 context optimization. In *CVPR*, 2023.

653 Beichen Zhang, Pan Zhang, Xiaoyi Dong, Yuhang Zang, and Jiaqi Wang. Long-clip: Unlocking the  
 654 long-text capability of clip. In *European conference on computer vision*, pp. 310–325. Springer,  
 655 2024a.

656 Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan.  
 657 Theoretically principled trade-off between robustness and accuracy. In *International conference on  
 658 machine learning*, pp. 7472–7482. PMLR, 2019.

660 Jiaming Zhang, Xingjun Ma, Xin Wang, Lingyu Qiu, Jiaqi Wang, Yu-Gang Jiang, and Jitao Sang.  
 661 Adversarial prompt tuning for vision-language models. In *European Conference on Computer  
 662 Vision*, pp. 56–72. Springer, 2024b.

663 Jingfeng Zhang, Xilie Xu, Bo Han, Gang Niu, Lizhen Cui, Masashi Sugiyama, and Mohan Kankanhalli.  
 664 Attacks which do not kill training make adversarial learning stronger. In *International  
 665 conference on machine learning*, pp. 11278–11287. PMLR, 2020.

666 Yunqing Zhao, Tianyu Pang, Chao Du, Xiao Yang, Chongxuan Li, Ngai-Man Man Cheung, and Min  
 667 Lin. On evaluating adversarial robustness of large vision-language models. In *Advances in Neural  
 668 Information Processing Systems*, volume 36, pp. 54111–54138, 2023.

669 Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for  
 670 vision-language models. In *CVPR*, 2022a.

672 Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-  
 673 language models. *IJCV*, 130(9):2337–2348, 2022b.

675 Yiwei Zhou, Xiaobo Xia, Zhiwei Lin, Bo Han, and Tongliang Liu. Few-shot adversarial prompt  
 676 learning on vision-language models. In *Advances in Neural Information Processing Systems*,  
 677 volume 37, pp. 3122–3156, 2024.

678 Beier Zhu, Yulei Niu, Yucheng Han, Yue Wu, and Hanwang Zhang. Prompt-aligned gradient for  
 679 prompt tuning. In *ICCV*, 2023.

681 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. MiniGPT-4: Enhancing  
 682 vision-language understanding with advanced large language models. In *ICLR*, 2024.

683

684 Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. Universal  
 685 and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*,  
 686 2023.

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

---

702 **Contents of the Appendix**  
703

704 1. Appendix A — Algorithm  
705 2. Appendix B — Implementation Details  
706 3. Appendix C — Additional Experimental Results
 


707 - C.1 — Comprehensive Results of Adversarial Few-shot Classification
708 - C.2 — Analysis on More Hyperparameter Sensitivity
709 - C.3 — Analysis of Logits Fusion Strategy
710 - C.4 — Adversarial Robustness Evaluation under Various Attacks
711 - C.5 — Adaptive Attack on DEFEAT
712 - C.6 — Generalization to Stronger Perturbation Budgets
713 - C.7 — Generalization to Alternative Robust Backbones
714 - C.8 — Generalization to Alternative VLM Backbones
715 - C.9 — Robust Visual Backbone Reliance
716 - C.10 — Computational Cost Analysis
717 - C.11 — [Performance Comparison under the  \$\ell\_2\$  Threat Model](#)
718 - C.12 — [Combination with Test-time Defense Strategy](#)
719 - C.13 — [Theoretical Justification](#)
720 - C.14 — [Stability Analysis with Multi-Seed Evaluation](#)
721 - C.15 — [Stability Analysis against Prompt Variations](#)

722 4. Appendix D — Use of Large Language Models  
723

724 **A ALGORITHM**  
725

726 The whole algorithm for the proposed DEFEAT methods is shown in Algorithm 1.  
727

728 **Algorithm 1** Discrete Latent Feature based Adversarial Training (DEFEAT)  
729

730

---

731 1: **Input:** Dataset  $\mathbb{D}$ , learnable prompt vectors  $\mathbf{v}$ , textual embeddings of hand-crafted prompts  $\mathbf{t}$ ,  
732 CLIP pre-trained image encoder  $\mathcal{I}$  and text encoder  $\mathcal{T}$ , VQ-VAE model  $\theta_{\text{VQ-VAE}}$ , learnable  
733 matrix  $\mathbf{W}$ , weight parameters  $\lambda, \mu$ , learning rate  $\eta$ , and perturbation size  $\epsilon$ .  
734 2: **for all** training epochs **do**  
735 3:   **for all**  $\mathbf{x}, \mathbf{y} \in$  minibatch **do**  
736 4:     Use  $\mathbf{v}$  and [CLS] to generate textual embeddings of learnable prompts  $\mathbf{t}^s$ ;  
737 5:     Generate adversarial examples  
738        $\mathbf{x}_a \leftarrow \arg \max_{\mathbf{x}_a} \mathcal{L}_{\text{ce}}(\mathbf{x}_a, \mathbf{t}^s, \mathbf{y})$ , s.t.  $\|\mathbf{x}_a - \mathbf{x}\|_p \leq \epsilon$ ;  
739 6:     Feed  $\mathbf{x}_a$  to  $\mathcal{I}$  to generate class feature  $\mathbf{I}$  and grid feature  $\mathbf{I}_{\text{patch}}$ ;  
740 7:     Reconstruct  $\mathbf{I}_{\text{patch}}$  using a VQ-VAE, obtaining  $\hat{\mathbf{I}}_{\text{patch}}$ ;  
741 8:     Obtain a robust representation by a transformation  
742        $\hat{\mathbf{I}}_{\text{proj}} \leftarrow \mathbf{W} \cdot \hat{\mathbf{I}}_{\text{patch}}$   
743 9:     Fuse the logits from  $\hat{\mathbf{I}}_{\text{proj}}$  in Eq. (6) and  $\mathbf{I}$  in Eq. (4)  
744        $p(y|\mathbf{x}_a) \leftarrow (p^{\text{vq}}(y|\mathbf{x}_a) + p^s(y|\mathbf{x}_a))/2$ ;  
745 10:   Train with the objective  
746       $\mathcal{L}(\mathbf{x}_a, \mathbf{t}^s, \mathbf{t}, \mathbf{y}) \leftarrow \mathcal{L}_{\text{ce}}(\mathbf{x}_a, \mathbf{t}^s, \mathbf{y}) + \mathcal{L}_{\text{VQ-VAE}}(\mathbf{I}_{\text{patch}}) + \lambda \mathcal{L}_{\text{reg}}^{\mathbf{I}} + \mu \mathcal{L}_{\text{reg}}^{\mathbf{T}}$ ,  
747      where  $\mathcal{L}_{\text{ce}}(\mathbf{x}_a, \mathbf{t}^s, \mathbf{y}) \leftarrow -\sum_{i=1}^C y_i \log p(y_i|\mathbf{x}_a)$ ,  $\mathcal{L}_{\text{reg}}^{\mathbf{I}}$  is defined in Eq. (5),  $\mathcal{L}_{\text{reg}}^{\mathbf{T}}$   
748      is defined in Eq. (8), and  $\mathcal{L}_{\text{VQ-VAE}}(\cdot)$  is defined in Eq. (3);  
749 11:   Upadate the learnable parameters  
750       $\mathbf{v} \leftarrow \mathbf{v} - \eta \nabla_{\mathbf{v}} \mathcal{L}(\mathbf{x}_a, \mathbf{t}^s, \mathbf{t}, \mathbf{y})$ ;  
751       $\mathbf{W} \leftarrow \mathbf{W} - \eta \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{x}_a, \mathbf{t}^s, \mathbf{t}, \mathbf{y})$ ;  
752       $\theta_{\text{VQ-VAE}} \leftarrow \theta_{\text{VQ-VAE}} - \eta \nabla_{\theta_{\text{VQ-VAE}}} \mathcal{L}(\mathbf{x}_a, \mathbf{t}^s, \mathbf{t}, \mathbf{y})$ .  
753 12:   **end for**  
754 13: **end for**  
755

---

756 **B IMPLEMENTATION DETAILS**  
757

758 All experiments are conducted on NVIDIA GeForce RTX 3090, except for the ImageNet dataset,  
 759 which is on Quadro RTX 8000. Training is conducted using SGD with an initial learning rate of  
 760 0.002, which is decayed using the cosine annealing rule. The maximum epoch is set to 200, 100, and  
 761 50 for 16, 4, and 1 shots, respectively. For ImageNet, they are 50, 20, and 20. A warm-up strategy is  
 762 used by fixing the learning rate to  $10^{-5}$  during the first epoch. **The VQ-VAE within the PerturbShield**  
 763 **module is configured with an input dimension of 512 to align with CLIP ViT-B/32 grid features,**  
 764 **projecting them into a latent dimension of 256 using a codebook of 512 unique embedding vectors.**  
 765 **The encoder of VQ-VAE consists of two convolutional layers followed by a residual block, while the**  
 766 **decoder of VQ-VAE mirrors this structure with a residual block and a transposed convolutional layer**  
 767 **for feature reconstruction.**

768 The hand-crafted prompts for different datasets used in the prompt alignment regularization follow  
 769 Radford et al. (2021); Zhou et al. (2022b) and are shown below:

```
770     ImageNet: "a photo of a [CLS]."  

  771     Caltech101: "a photo of a [CLS]."  

  772     OxfordPets: "a photo of a [CLS], a type of pet."  

  773     StanfordCars: "a photo of a [CLS]."  

  774     OxfordFlowers: "a photo of a [CLS], a type of flower."  

  775     Food101: "a photo of [CLS], a type of food."  

  776     FGVCAircraft: "a photo of a [CLS], a type of aircraft."  

  777     SUN397: "a photo of a [CLS]."  

  778     DTD: "a photo of a [CLS], a type of texture."  

  779     EuroSAT: "a centered satellite photo of [CLS]."  

  780     UCF101: "a photo of a person doing [CLS]."
```

781 Note that [CLS] denotes the placeholder for the class name.

782 **C ADDITIONAL EXPERIMENTAL RESULTS**783 **C.1 COMPREHENSIVE RESULTS OF ADVERSARIAL FEW-SHOT CLASSIFICATION**

784 In this section, we provide the full results of adversarial few-shot classification performance on 11  
 785 datasets. The results in Tables A5 and A6 show that the proposed DEFEAT method can consistently  
 786 outperform the baseline methods in terms of the robustness and harmonic mean accuracy across  
 787 nearly all datasets and shot settings. For each specific dataset, DEFEAT consistently outperforms  
 788 CLIP across all shots. For adversarial prompt tuning methods, DEFEAT outperforms APT and FAP  
 789 across almost all the datasets and shot settings, with particularly notable improvements on EuroSAT,  
 790 UCF101, etc.. For example, with 16-shot training, DEFEAT surpasses APT and FAP by 20.13%  
 791 and 23.79% in terms of 'H' (i.e., 56.69% vs. 36.56% vs. 32.90%) on EuroSAT, and by 23.21% and  
 792 21.35% (i.e., 48.45% vs. 25.24% vs. 27.10%) on UCF101.

793 Figure A7 shows that as the number of shots increases, DEFEAT achieves greater improvements  
 794 compared to APT and FAP. Those results demonstrate DEFEAT's effectiveness and scalability.

795 **C.2 ANALYSIS ON MORE HYPERPARAMETER SENSITIVITY**

796 In this section, we conduct ablation studies to analyse the effectiveness of hyperparameters (i.e.,  $\alpha$ ,  $\beta$ ,  
 797 and  $\lambda$ ).

798 **Effect of  $\alpha$ .** In Eq. (3),  $\alpha$  controls the weight of the reconstruction loss, which is used to optimize  
 799 the decoder and the encoder of VQ-VAE. According to Figure A8a, we can see that DEFEAT is  
 800 insensitive within the range [0,5]. We set  $\alpha = 0.5$  in our experiments.

801 **Effect of  $\beta$ .** In Eq. (3),  $\beta$  controls the weight of the codebook loss, which is used to learn the  
 802 embedding space of the codebook. According to Figure A8b, applying the codebook loss (i.e.,  
 803  $\beta > 0$ ) results in better overall performance compared to excluding it (i.e.,  $\beta = 0$ ). Failure to  
 804 update the codebook reduces the representational capacity of the embeddings, thereby increasing  
 805 the reconstruction quantization error of the VQ-VAE. Additionally, the model may collapse if the  
 806 codebook is poorly initialized. The overall performance of DEFEAT is insensitive within the range  
 807  $\beta \in [0.01, 0.1]$ . In our experiments, we set  $\beta = 0.1$ .

810  
811 Table A5: The performance on the 11 datasets for different shots under the adversarial few-shot  
812 classification setting when  $\epsilon = 1/255$ , where ‘H’ denotes the harmonic mean accuracy.

| Dataset      | $\epsilon = 1/255$ | 1 shot |              |              | 4 shot       |              |              | 16 shot      |              |              |              |
|--------------|--------------------|--------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|              |                    | CLIP   | APT          | FAP          | DEFEAT       | APT          | FAP          | DEFEAT       | APT          | FAP          | DEFEAT       |
| Average      | Acc.               | 46.06  | 46.99        | <b>52.72</b> | 52.07        | 58.19        | <b>59.77</b> | 58.37        | 65.41        | <b>66.40</b> | 65.03        |
|              | Rob.               | 32.98  | 33.36        | 37.15        | <b>38.99</b> | 41.34        | 42.21        | <b>53.04</b> | 47.88        | 48.86        | <b>60.77</b> |
|              | H                  | 38.44  | 39.02        | 43.58        | <b>44.59</b> | 48.34        | 49.48        | <b>55.57</b> | 55.29        | 56.29        | <b>62.83</b> |
| ImageNet     | Acc.               | 55.34  | <b>55.48</b> | 55.01        | 54.33        | 56.31        | <b>57.39</b> | 56.11        | 50.73        | <b>58.46</b> | 57.54        |
|              | Rob.               | 38.48  | 39.00        | 37.79        | <b>39.68</b> | 39.30        | 39.63        | <b>41.36</b> | 33.37        | 40.37        | <b>53.07</b> |
|              | H                  | 45.40  | 45.80        | 44.80        | <b>45.86</b> | 46.29        | 46.88        | <b>47.62</b> | 40.26        | 47.76        | <b>55.21</b> |
| Caltech101   | Acc.               | 85.92  | 86.46        | 88.15        | <b>89.66</b> | 89.17        | 90.87        | <b>91.16</b> | 92.90        | 92.21        | <b>93.39</b> |
|              | Rob.               | 75.21  | 77.32        | 78.86        | <b>80.73</b> | 79.68        | 80.81        | <b>87.99</b> | 83.98        | 84.30        | <b>91.12</b> |
|              | H                  | 80.21  | 81.63        | 83.25        | <b>84.96</b> | 84.16        | 85.55        | <b>89.55</b> | 88.22        | 88.08        | <b>92.24</b> |
| OxfordPets   | Acc.               | 79.48  | 78.33        | <b>79.56</b> | 77.19        | 83.48        | 83.21        | <b>84.08</b> | 83.95        | <b>85.55</b> | 84.71        |
|              | Rob.               | 63.26  | 61.57        | 60.56        | <b>62.36</b> | 66.18        | 66.34        | <b>73.37</b> | 65.14        | 69.15        | <b>79.53</b> |
|              | H                  | 70.45  | 68.95        | 68.77        | <b>68.99</b> | 73.83        | 73.82        | <b>78.36</b> | 73.36        | 76.48        | <b>82.04</b> |
| StanfordCars | Acc.               | 25.21  | 41.09        | <b>45.55</b> | 42.87        | 49.51        | <b>51.14</b> | 47.69        | <b>60.20</b> | 57.54        | 56.42        |
|              | Rob.               | 12.50  | 22.44        | 24.05        | <b>26.81</b> | 26.94        | 27.47        | <b>44.05</b> | 37.21        | 32.79        | <b>51.65</b> |
|              | H                  | 16.71  | 29.03        | 31.48        | <b>32.99</b> | 34.89        | 35.74        | <b>45.80</b> | 45.99        | 41.77        | <b>53.93</b> |
| Flowers102   | Acc.               | 48.15  | 23.63        | <b>59.60</b> | 53.39        | <b>76.61</b> | 70.36        | 72.72        | <b>86.80</b> | 82.46        | 84.90        |
|              | Rob.               | 32.68  | 17.46        | <b>42.14</b> | 40.80        | 59.68        | 51.48        | <b>69.02</b> | 73.57        | 66.63        | <b>81.81</b> |
|              | H                  | 38.93  | 20.08        | <b>49.37</b> | 46.25        | 67.09        | 59.46        | <b>70.82</b> | 79.64        | 73.70        | <b>83.33</b> |
| Food101      | Acc.               | 46.58  | 41.18        | <b>55.36</b> | 49.83        | 45.37        | <b>59.24</b> | 51.29        | 54.68        | <b>64.48</b> | 56.31        |
|              | Rob.               | 27.38  | 22.12        | <b>33.32</b> | 32.07        | 24.95        | 36.20        | <b>46.84</b> | 33.02        | 41.19        | <b>49.45</b> |
|              | H                  | 34.49  | 28.78        | <b>41.60</b> | 39.02        | 32.20        | 44.94        | <b>48.96</b> | 41.18        | 50.27        | <b>52.66</b> |
| FGVCAircraft | Acc.               | 12.51  | 2.01         | <b>18.69</b> | 17.55        | 14.64        | 21.00        | <b>21.27</b> | <b>28.74</b> | 26.55        | 25.74        |
|              | Rob.               | 6.09   | 1.08         | 10.41        | <b>11.16</b> | 7.44         | 11.43        | <b>19.98</b> | 16.53        | 14.97        | <b>23.46</b> |
|              | H                  | 8.19   | 1.41         | 13.37        | <b>13.64</b> | 9.87         | 14.80        | <b>20.60</b> | 20.99        | 19.15        | <b>24.55</b> |
| SUN397       | Acc.               | 48.67  | 51.17        | <b>53.50</b> | 52.63        | 56.35        | <b>58.90</b> | 58.75        | 62.56        | 62.69        | <b>62.83</b> |
|              | Rob.               | 31.70  | 32.44        | 34.90        | <b>35.73</b> | 36.48        | 38.71        | <b>54.01</b> | 42.31        | 42.67        | <b>57.53</b> |
|              | H                  | 38.39  | 39.71        | 42.24        | <b>42.56</b> | 44.29        | 46.72        | <b>56.28</b> | 50.48        | 50.78        | <b>60.06</b> |
| DTD          | Acc.               | 31.97  | 34.10        | <b>36.83</b> | 38.48        | 47.04        | <b>48.58</b> | 47.93        | 55.26        | 56.68        | <b>57.15</b> |
|              | Rob.               | 23.88  | 23.40        | 24.23        | <b>28.55</b> | 33.33        | 35.87        | <b>44.33</b> | 39.54        | 41.67        | <b>52.66</b> |
|              | H                  | 27.34  | 27.75        | 29.23        | <b>32.78</b> | 39.02        | 41.27        | <b>46.06</b> | 46.10        | 48.03        | <b>54.81</b> |
| EuroSAT      | Acc.               | 23.62  | <b>51.44</b> | 31.80        | 42.00        | <b>61.43</b> | 55.38        | 47.98        | 74.15        | <b>74.89</b> | 67.20        |
|              | Rob.               | 16.48  | <b>35.37</b> | 21.89        | 31.80        | 39.10        | 31.88        | <b>43.86</b> | 51.53        | 52.48        | <b>62.80</b> |
|              | H                  | 19.41  | <b>41.92</b> | 25.93        | 36.20        | <b>47.75</b> | 40.47        | 45.83        | 60.80        | 61.71        | <b>64.93</b> |
| UCF101       | Acc.               | 49.22  | 52.02        | <b>55.86</b> | 54.80        | 60.22        | 61.43        | <b>63.05</b> | 69.55        | 68.89        | <b>69.10</b> |
|              | Rob.               | 35.13  | 34.73        | <b>40.47</b> | 39.15        | 41.63        | 44.49        | <b>58.58</b> | 50.52        | 51.23        | <b>65.40</b> |
|              | H                  | 41.00  | 41.65        | <b>46.94</b> | 45.67        | 49.23        | 51.61        | <b>60.73</b> | 58.53        | 58.76        | <b>67.20</b> |

837  
838 Table A6: The performance on the 11 datasets for different shots under the adversarial few-shot  
839 classification setting when  $\epsilon = 4/255$ , where ‘H’ denotes the harmonic mean accuracy.

| Dataset      | $\epsilon = 4/255$ | 1 shot |              |              | 4 shot       |              |              | 16 shot      |              |              |              |
|--------------|--------------------|--------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|              |                    | CLIP   | APT          | FAP          | DEFEAT       | APT          | FAP          | DEFEAT       | APT          | FAP          | DEFEAT       |
| Average      | Acc.               | 33.67  | 32.97        | <b>38.16</b> | 35.78        | 42.29        | <b>44.63</b> | 44.13        | <b>51.08</b> | 50.50        | 50.98        |
|              | Rob.               | 10.79  | 11.62        | 13.44        | <b>14.73</b> | 14.40        | 15.34        | <b>25.28</b> | 20.26        | 19.82        | <b>36.05</b> |
|              | H                  | 16.34  | 17.19        | 19.88        | <b>20.87</b> | 21.49        | 22.83        | <b>32.14</b> | 29.02        | 28.47        | <b>42.24</b> |
| ImageNet     | Acc.               | 40.11  | <b>37.88</b> | 36.95        | 37.32        | 39.13        | 38.94        | <b>39.28</b> | 41.06        | 40.45        | <b>41.09</b> |
|              | Rob.               | 10.14  | 10.79        | 11.41        | <b>12.48</b> | 11.49        | 11.86        | <b>12.66</b> | 12.02        | 12.03        | <b>18.42</b> |
|              | H                  | 16.19  | 16.80        | 17.44        | <b>18.70</b> | 17.76        | 18.18        | <b>19.15</b> | 18.60        | 18.54        | <b>25.44</b> |
| Caltech101   | Acc.               | 78.78  | 77.89        | 76.67        | <b>79.55</b> | 81.62        | 77.61        | <b>83.16</b> | 86.29        | 80.04        | <b>87.63</b> |
|              | Rob.               | 43.61  | 45.84        | 48.36        | <b>51.03</b> | 47.59        | 51.81        | <b>58.74</b> | 56.75        | 55.42        | <b>78.30</b> |
|              | H                  | 56.14  | 57.71        | 59.31        | <b>62.18</b> | 60.12        | 62.14        | <b>68.85</b> | 68.47        | 65.49        | <b>82.70</b> |
| OxfordPets   | Acc.               | 66.26  | 59.58        | <b>63.01</b> | 57.51        | 65.99        | 66.86        | <b>69.45</b> | 67.29        | 69.09        | <b>72.17</b> |
|              | Rob.               | 15.56  | 15.54        | 20.31        | <b>22.16</b> | 17.80        | 22.59        | <b>27.42</b> | 19.98        | 26.17        | <b>30.77</b> |
|              | H                  | 25.20  | 24.65        | 30.72        | <b>31.99</b> | 28.04        | 33.77        | <b>39.32</b> | 30.81        | 37.96        | <b>43.14</b> |
| StanfordCars | Acc.               | 10.32  | 19.49        | <b>27.41</b> | 22.96        | 23.93        | <b>33.96</b> | 25.90        | 31.60        | <b>41.06</b> | 32.45        |
|              | Rob.               | 0.99   | 2.90         | 3.27         | <b>4.74</b>  | 4.44         | 4.76         | <b>8.20</b>  | 7.70         | 6.33         | <b>19.51</b> |
|              | H                  | 1.81   | 5.05         | 5.84         | <b>7.86</b>  | 7.49         | 8.35         | <b>12.46</b> | 12.38        | 10.97        | <b>24.37</b> |
| Flowers102   | Acc.               | 30.13  | 35.93        | <b>42.55</b> | 34.59        | <b>60.25</b> | 55.58        | 56.11        | <b>76.41</b> | 66.10        | 71.13        |
|              | Rob.               | 8.93   | 11.49        | <b>14.21</b> | 14.01        | 22.45        | 20.46        | <b>48.60</b> | 37.52        | 30.82        | <b>58.26</b> |
|              | H                  | 13.78  | 17.41        | <b>21.30</b> | 19.94        | 32.71        | 29.91        | <b>52.09</b> | 50.33        | 42.04        | <b>64.05</b> |
| Food101      | Acc.               | 23.52  | 20.73        | <b>33.46</b> | 26.96        | 21.97        | <b>38.84</b> | 29.75        | 30.39        | <b>45.60</b> | 33.57        |
|              | Rob.               | 3.27   | 3.42         | 5.46         | <b>6.01</b>  | 4.01         | 6.39         | <b>8.29</b>  | 7.90         | 8.23         | <b>20.50</b> |
|              | H                  | 5.74   | 5.87         | 9.39         | <b>9.83</b>  | 6.78         | 10.97        | <b>12.97</b> | 12.54        | 13.94        | <b>25.46</b> |
| FGVCAircraft | Acc.               | 7.17   | 1.41         | <b>11.94</b> | 11.67        | 2.31         | <b>15.27</b> | 14.58        | <b>20.31</b> | 18.57        | 18.03        |
|              | Rob.               | 0.36   | 0.42         | 2.52         | <b>2.73</b>  | 0.51         | 2.73         | <b>10.26</b> | 6.15         | 5.07         | <b>11.94</b> |
|              | H                  | 0.69   | 0.65         | 4.16         | <b>4.42</b>  | 0.84         | 4.63         | <b>12.04</b> | 9.44         | 7.97         | <b>14.37</b> |
| SUN397       | Acc.               | 33.24  | 32.32        | <b>36.44</b> | 35.68        | 39.06        | 39.68        | <b>41.29</b> | 45.21        | 43.66        | <b>46.26</b> |
|              | Rob.               | 6.20   | 5.76         | 8.63         | <b>9.14</b>  | 7.92         | 9.69         | <b>19.99</b> | 11.27        | 11.44        | <b>28.88</b> |
|              | H                  | 10.45  | 9.78         | 13.96        | <b>14.55</b> | 13.17        | 15.58        | <b>26.94</b> | 18.04        | 18.13        | <b>35.56</b> |
| DTD          | Acc.               | 24.29  | 23.29        | <b>27.78</b> | 25.24        | 36.17        | 37.41        | <b>38.18</b> | <b>45.86</b> | 42.85        | 44.15        |
|              | Rob.               | 11.35  | 9.46         | 10.87        | <b>11.88</b> | 14.13        | 15.60        | <b>26.48</b> | 21.51        | 20.98        | <b>32.33</b> |
|              | H                  | 15.47  | 13.45        | 15.63        | <b>16.16</b> | 20.32        | 22.02        | <b>31.27</b> | 29.28        | 28.17        | <b>37.33</b> |
| EuroSAT      | Acc.               | 19.56  | 21.09        | <b>28.02</b> | 25.49        | <b>49.98</b> | 40.58        | 39.47        | <b>64.33</b> | 57.37        | 58.88        |
|              | Rob.               | 11.22  | 14.53        | 13.94        | <b>15.77</b> | 16.44        | 8.57         | <b>28.07</b> | 25.54        | 23.06        | <b>54.65</b> |
|              | H                  | 14.26  | 17.21        | 18.62        | <b>19.49</b> | 24.74        | 14.15        | <b>32.81</b> | 36.56        | 32.90        | <b>56.69</b> |
| UCF101       | Acc.               | 36.98  | 33.10        | 35.58        | <b>36.56</b> | 44.75        | 46.21        | <b>48.27</b> | 53.16        | 50.67        | <b>55.43</b> |
|              | Rob.               | 7.01   | 7.72         | 8.86         | <b>12.11</b> | 11.66        | 14.27        | <b>29.32</b> | 16.55        | 18.50        | <b>43.03</b> |
|              | H                  | 11.79  | 12.52        | 14.19        | <b>18.19</b> | 18.50        | 21.81        | <b>36.48</b> | 25.24        | 27.10        | <b>48.45</b> |

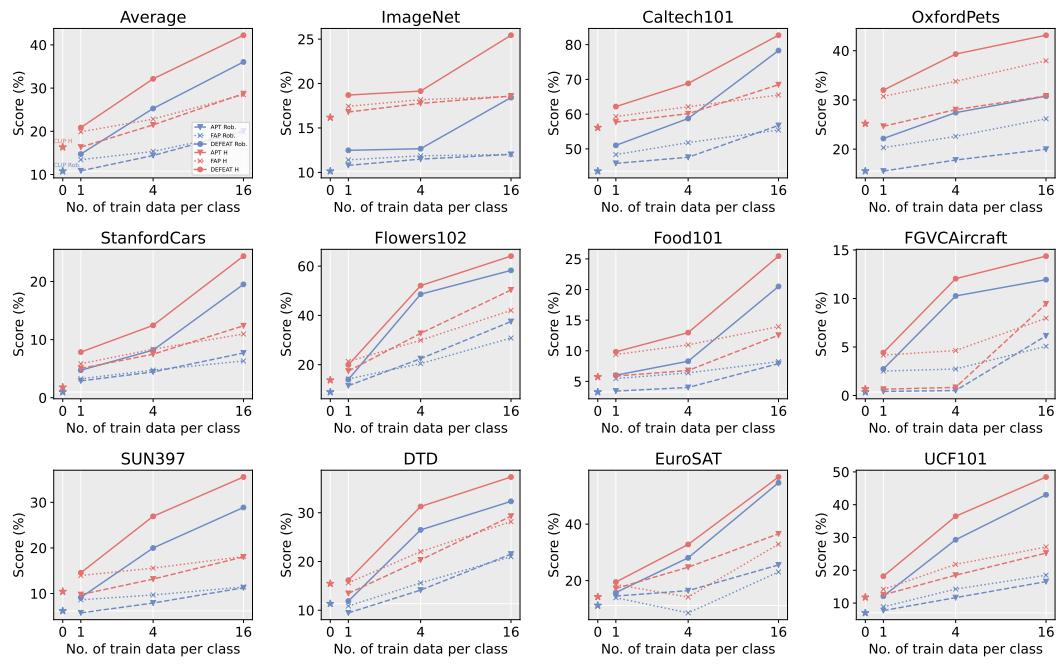


Figure A7: Performance on the 11 datasets under the adversarial few-shot classification setting when  $\epsilon = 4/255$ . As the accuracies among prompt tuning methods are comparable, we only plot the curves for robustness and ‘H’ to improve readability.

**Effect of  $\lambda$ .** In Eq. (10),  $\lambda$  controls the weight of  $\mathcal{L}_{\text{reg}}^{\mathbf{I}}$ , which is designed to ensure consistency between the  $\hat{\mathbf{I}}_{\text{proj}}$  and the CLIP pre-trained class feature  $\mathbf{I}$ . According to Figure A8c, as  $\lambda$  increases, the overall performance of DEFECT improves, and then remains stable thereafter (i.e.,  $\lambda > 10$ ), making it easy to choose an appropriate  $\lambda$  in practice. Those results highlight the importance of regularizing  $\hat{\mathbf{I}}_{\text{proj}}$ . In our experiments, we set  $\lambda = 10$ .

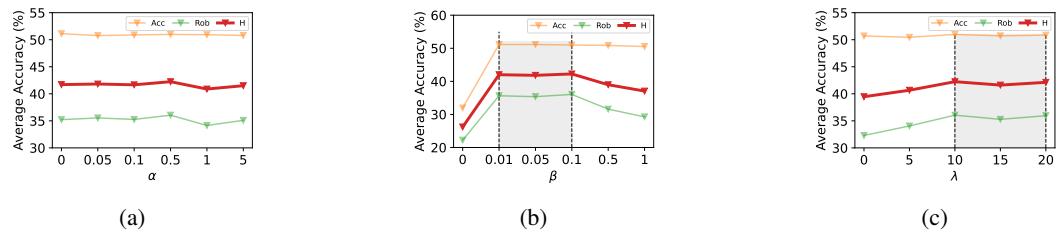


Figure A8: The average performance of DEFECT on all datasets w.r.t. hyperparameters (i.e.,  $\alpha$ ,  $\beta$ , and  $\lambda$ ) under the adversarial few-shot classification setting.

### C.3 ANALYSIS OF LOGITS FUSION STRATEGY

In this section, we analyse a variant of the logits fusion strategy. Specifically, to control the trade-off between robustness and accuracy flexibility, we can modify the fused logits in Eq. (7) to the following form:

$$p(y|\mathbf{x}) = (1 - \omega) p^{\text{vq}}(y|\mathbf{x}) + \omega p^{\text{s}}(y|\mathbf{x}) \quad (11)$$

where  $\omega$  is the hyperparameter. We analyze how different  $\omega$  values impact DEFECT’s performance. Figure A9 shows that as  $\omega$  increases, accuracy gradually improves because the logits of  $\mathbf{I}$  and  $\mathbf{t}^{\text{s}}$ , which contribute to clean accuracy, have a larger weight in the fused logits. Robustness and ‘H’ initially increase with  $\omega$ , then plateau, and eventually decline. Although the range of  $\omega$  values where ‘H’ stabilizes varies slightly across different datasets,  $\omega = 0.5$  consistently falls within or near this range. Therefore, we set  $\omega = 0.5$  for all datasets in our experiments.

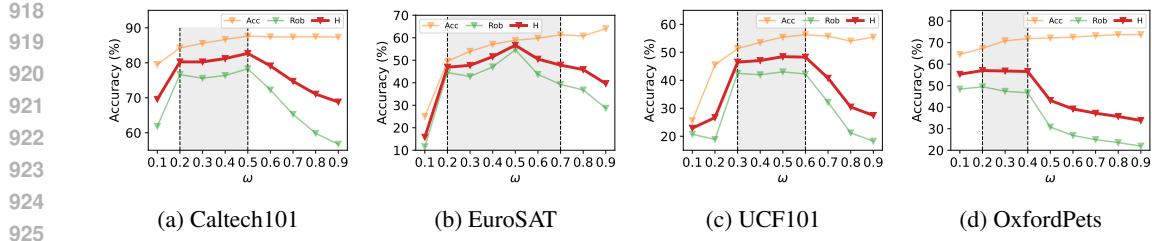


Figure A9: Performance of DEFEAT on Caltech101, EuroSAT, UCF101, and OxfordPets datasets w.r.t. different  $\omega$  values under the adversarial few-shot classification setting.

#### C.4 ADVERSARIAL ROBUSTNESS EVALUATION UNDER VARIOUS ATTACKS

In this section, we evaluate the adversarial robustness of prompt tuning methods using a wider variety of attacks, including Carlini & Wagner (CW) attack (Carlini & Wagner, 2017) and AutoAttack (AA) (Croce & Hein, 2020). CW is a stronger optimization-based adversarial attack that finds the minimal perturbation required to cause a model to misclassify an input. AutoAttack (Croce & Hein, 2020) is a standardized benchmark suite that combines multiple attacks to evaluate a model’s adversarial robustness reliably. Specifically, we use the standard AutoAttack setting, which executes a suite of four attacks: APGD-CE (untargeted), APGD-T (targeted), FAB-T (targeted), and Square (black-box). All component attacks utilize the default setting of 1 restart. We apply non-deterministic seeding to each component attack. As shown in Table A7, both CW and AA generate more potent adversarial examples than PGD, resulting in lower robustness for all methods (i.e., APT, FAP, and DEFEAT). Despite this, DEFEAT consistently outperforms the APT and FAP baselines across all datasets. This consistent superiority against a diverse and powerful suite of attacks provides strong evidence that DEFEAT’s robustness is genuine and not an artifact of attack-specific overfitting or gradient obfuscation.

Table A7: Performance of adversarial prompt tuning methods under the adversarial few-shot classification setting with various attacks.

| 16shot, $\epsilon = 4/255$ |              |              |              | Average      |              |              |              | Caltech101   |              |              |              | OxfordPets   |              |              |              | StanfordCars |              |              |              | Flowers102   |  |  |  |
|----------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--|--|--|
| Method                     | Acc.         | PGD          | AA           | CW           |  |  |  |
| TeCoA                      | 33.00        | 11.37        | 9.64         | 11.07        | 78.78        | 43.61        | 40.93        | 44.30        | 66.26        | 15.56        | 11.28        | 15.26        | 10.32        | 0.99         | 0.62         | 0.99         | 30.13        | 8.93         | 6.54         | 7.59         |  |  |  |
| APT                        | 52.85        | 22.18        | 17.27        | 20.07        | 86.29        | 56.75        | 53.75        | 56.51        | 67.29        | 19.98        | 13.33        | 16.98        | 31.60        | 7.70         | 3.69         | 5.19         | <b>76.41</b> | 37.52        | 30.86        | 33.46        |  |  |  |
| FAP                        | <b>53.17</b> | 21.62        | 16.02        | 17.58        | 80.04        | 55.42        | 50.87        | 53.35        | 69.09        | 26.17        | 16.46        | 18.34        | <b>41.06</b> | 6.33         | 2.43         | 4.00         | 66.10        | 30.82        | 23.02        | 24.81        |  |  |  |
| DEFEAT                     | 52.60        | <b>39.92</b> | <b>27.02</b> | <b>34.89</b> | <b>87.63</b> | <b>78.30</b> | <b>70.30</b> | <b>77.32</b> | <b>72.17</b> | <b>30.77</b> | <b>18.86</b> | <b>23.63</b> | 32.45        | <b>19.51</b> | <b>9.63</b>  | <b>16.71</b> | 71.13        | <b>58.26</b> | <b>45.11</b> | <b>55.38</b> |  |  |  |
|                            |              |              |              |              |              |              |              | Food101      |              |              |              | FGVCAircraft |              |              |              | DTD          |              |              |              | EuroSAT      |  |  |  |
| Method                     | Acc.         | PGD          | AA           | CW           |  |  |  |
| TeCoA                      | 23.52        | 3.27         | 1.81         | 2.57         | 7.17         | 0.36         | 0.06         | 0.21         | 24.29        | 11.35        | 9.75         | 10.17        | 19.56        | 11.22        | 10.49        | 11.30        | 36.98        | 7.01         | 5.26         | 7.22         |  |  |  |
| APT                        | 30.39        | 7.90         | 4.48         | 5.54         | <b>20.31</b> | 6.15         | 3.06         | 4.05         | <b>45.86</b> | 21.51        | 15.07        | 17.14        | <b>64.33</b> | 25.54        | 18.35        | 25.04        | 53.16        | 16.55        | 12.82        | 16.68        |  |  |  |
| FAP                        | <b>45.60</b> | 8.23         | 3.92         | 5.37         | 18.57        | 5.07         | 2.58         | 3.06         | 42.85        | 20.98        | 17.32        | 18.09        | 64.51        | 23.06        | 15.19        | 16.33        | 50.67        | 18.50        | 12.42        | 14.83        |  |  |  |
| DEFEAT                     | 33.57        | <b>20.50</b> | <b>11.08</b> | <b>18.23</b> | 18.03        | <b>11.94</b> | <b>6.48</b>  | <b>10.41</b> | 44.15        | <b>42.33</b> | <b>25.95</b> | <b>30.85</b> | 58.88        | <b>54.65</b> | <b>27.35</b> | <b>40.00</b> | <b>55.43</b> | <b>43.03</b> | <b>28.42</b> | <b>41.48</b> |  |  |  |
|                            |              |              |              |              |              |              |              | UCF101       |              |              |              |              |              |              |              |              |              |              |              |              |  |  |  |

#### C.5 ADAPTIVE ATTACK ON DEFEAT

To provide a more rigorous evaluation and to demonstrate DEFEAT’s robustness, we design a defense-aware adaptive attack specifically tailored to circumvent its core architectural components. Following the principles outlined by Athalye et al. (2018), an adaptive attack must take the specific defense mechanism into account. Our attack targets the logits fusion strategy, which is a unique characteristic of the DEFEAT framework.

An attacker with knowledge of the DEFEAT architecture would realize that the final prediction comes from a standard, continuous feature branch and a defended, discretized feature branch. Based on this insight, we hypothesize that the most effective adaptive attack strategy is to identify the more vulnerable of these two branches and concentrate the full adversarial pressure on it.

To validate this hypothesis and identify the worst-case attack, we design a logits fusion-aware adaptive attack. Instead of maximizing the standard cross-entropy loss on the fused output, we maximize a weighted sum of the individual cross-entropy losses from each branch. The adaptive attack loss

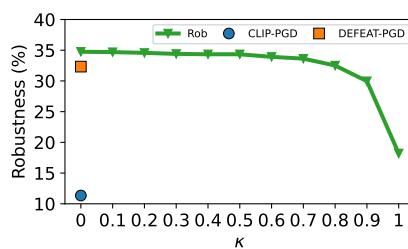


Figure A10: Robustness of DEFEAT on the DTD dataset under the logits fusion-aware adaptive attack with varying  $\kappa$  values. Standard PGD attacks on CLIP and DEFEAT are provided as baselines.

$\mathcal{L}_{\text{adapt}}$  is defined as:

$$\mathcal{L}_{\text{ce}}^s(\mathbf{x}, \mathbf{t}^s, \mathbf{y}) = \sum_{i=1}^C y_i \log p^s(y|\mathbf{x}) \quad (12)$$

$$\mathcal{L}_{\text{ce}}^{\text{vq}}(\mathbf{x}, \mathbf{t}^s, \mathbf{y}) = \sum_{i=1}^C y_i \log p^{\text{vq}}(y|\mathbf{x}) \quad (13)$$

$$\mathcal{L}_{\text{adapt}} = \kappa * \mathcal{L}_{\text{ce}}^s(\mathbf{x}, \mathbf{t}^s, \mathbf{y}) + (1 - \kappa) * \mathcal{L}_{\text{ce}}^{\text{vq}}(\mathbf{x}, \mathbf{t}^s, \mathbf{y}) \quad (14)$$

where  $p^s(y|\mathbf{x})$  and  $p^{\text{vq}}(y|\mathbf{x})$  are defined in Eq. (4) and Eq. (6),  $\kappa$  is a hyperparameter controls the focus of the attack. By varying the value of  $\kappa$ , we can explore the effect of attacking different branches. We then use a PGD-based optimizer to find the perturbation that maximizes  $\mathcal{L}_{\text{adapt}}$ .

We conduct experiments on the DTD dataset under the adversarial few-shot (16-shot) classification setting with  $\epsilon = 4/255$  to analyze DEFEAT’s performance under this adaptive attack. The results are presented in Figure A10.

The experimental results clearly validate our hypothesis. When we concentrate the full attack pressure on the standard, non-discretized feature branch ( $\kappa = 1$ ), DEFEAT’s robustness drops to its lowest point of 18.20%. This demonstrates that, with knowledge of DEFEAT’s defense mechanism, the most effective adaptive attack is indeed to focus entirely on the branch unprotected by the PerturbShield module.

Conversely, when the attack exclusively targets the defended PerturbShield branch ( $\kappa = 0$ ), the model’s robustness reaches its highest point of 34.75%, even surpassing its performance against a standard PGD attack on the fused logits (32.33%). This provides strong evidence for the effectiveness of our proposed discretization module.

In summary, this experiment validates the design of our adaptive attack. Even under this custom-designed, worst-case adaptive attack, DEFEAT’s robustness (18.20%) remains significantly higher than the CLIP baseline (11.35%), providing evidence for the good robustness of DEFEAT.

## C.6 GENERALIZATION TO STRONGER PERTURBATION BUDGETS

A critical measure of a defense mechanism’s effectiveness is its ability to generalize to attacks stronger than those seen during training. To evaluate this, we conduct a challenging experiment where all models are trained using 16 examples per class under a standard PGD attack with a perturbation budget of  $\epsilon = 4/255$ , but are then evaluated against a much stronger, unseen PGD attack with  $\epsilon = 8/255$ .

As shown in Table A8, the results demonstrate a consistent advantage for DEFEAT. While the stronger attack significantly degrades the robustness of all baselines (compared to the 16-shot results in Table A6), DEFEAT maintains a substantially higher level of robustness across every dataset. On average, DEFEAT achieves a robust accuracy of 23.32%, which is over 4 times higher than FAP (5.81%) and 5 times higher than APT (4.62%).

Crucially, this significant gain in robustness is achieved with no sacrifice in clean accuracy. DEFEAT remains competitive with the baselines in this regard. This experiment provides compelling evidence that DEFEAT learns a more fundamental and generalizable defense mechanism. Instead of merely overfitting to a specific perturbation budget, its discrete feature reconstruction allows it to more

effectively ablate perturbations, making it a more reliable defense for real-world scenarios where attack strengths are unknown.

Table A8: Performance against larger perturbation budgets (trained on  $\epsilon = 4/255$ , tested on  $\epsilon = 8/255$ ).

| train ( $\epsilon = 4/255$ ) |        | Average |              | Caltech101 |              | OxfordPets |             | EuroSAT |              | UCF101 |              | DTD   |              | Flowers102 |              | StanfordCars |             | FGVCAircraft |             |
|------------------------------|--------|---------|--------------|------------|--------------|------------|-------------|---------|--------------|--------|--------------|-------|--------------|------------|--------------|--------------|-------------|--------------|-------------|
| test ( $\epsilon = 8/255$ )  | Method | Acc.    | Rob.         | Acc.       | Rob.         | Acc.       | Rob.        | Acc.    | Rob.         | Acc.   | Rob.         | Acc.  | Rob.         | Acc.       | Rob.         | Acc.         | Rob.        | Acc.         | Rob.        |
| PGD                          | APT    | 55.66   | 4.62         | 86.29      | 20.89        | 67.29      | 1.31        | 64.33   | 1.62         | 53.16  | 1.43         | 45.86 | 6.68         | 76.41      | 4.26         | 31.60        | 0.52        | 20.31        | 0.27        |
|                              | FAP    | 53.22   | 5.81         | 80.04      | 25.03        | 69.09      | 2.02        | 57.37   | 2.41         | 50.67  | 2.06         | 42.85 | 8.33         | 66.10      | 5.68         | 41.06        | 0.29        | 18.57        | 0.66        |
|                              | DEFEAT | 54.98   | <b>23.32</b> | 87.63      | <b>58.46</b> | 72.17      | <b>2.94</b> | 58.88   | <b>31.09</b> | 55.43  | <b>21.46</b> | 44.15 | <b>19.56</b> | 71.13      | <b>36.09</b> | 32.45        | <b>7.61</b> | 18.03        | <b>9.33</b> |

### C.7 GENERALIZATION TO ALTERNATIVE ROBUST BACKBONES

To demonstrate that the effectiveness of DEFEAT is not limited to a specific robust backbone, we replace the TeCoA-trained image encoder (Mao et al., 2023) with a more recent and powerful one pre-trained via the FARE method (Schlarmann et al., 2024). We then evaluate all adversarial prompt tuning baselines under the adversarial 16-shot classification setting.

The results in Table A9 show that DEFEAT maintains its superiority. Under the weaker PGD attack ( $\epsilon = 1/255$ ), DEFEAT achieves an average robustness of 69.33%, outperforming the next-best baseline (FAP) by nearly 15%. This advantage becomes even more pronounced under the stronger attack ( $\epsilon = 4/255$ ), where DEFEAT’s average robustness of 52.46% is over 20% higher than the strongest baseline (FAP at 30.53%). Those results validate that DEFEAT is a generalizable framework that enhances robustness independently of the underlying pre-training method, and its benefits are magnified when paired with stronger robust encoders.

Table A9: Performance on a CLIP ViT-B/32 image encoder pre-trained with FARE (Schlarmann et al., 2024).

| $\epsilon = 1/255$ | Method  | Average |              | Caltech101 |              | OxfordPets |              | EuroSAT |              | UCF101 |              | DTD   |              | Flowers102 |              | StanfordCars |              | FGVCAircraft |              |
|--------------------|---------|---------|--------------|------------|--------------|------------|--------------|---------|--------------|--------|--------------|-------|--------------|------------|--------------|--------------|--------------|--------------|--------------|
|                    |         | Acc.    | Rob.         | Acc.       | Rob.         | Acc.       | Rob.         | Acc.    | Rob.         | Acc.   | Rob.         | Acc.  | Rob.         | Acc.       | Rob.         | Acc.         | Rob.         | Acc.         | Rob.         |
| PGD                | FARE    | 54.87   | 34.09        | 91.12      | 76.63        | 86.92      | 61.60        | 24.89   | 13.83        | 59.08  | 33.81        | 41.13 | 26.60        | 62.16      | 33.78        | 56.97        | 21.43        | 16.65        | 5.04         |
|                    | +APT    | 75.57   | 50.60        | 95.25      | 84.06        | 87.87      | 63.29        | 79.58   | 42.44        | 78.91  | 51.86        | 63.95 | 41.31        | 93.10      | 73.08        | 72.89        | 37.71        | 32.97        | 11.07        |
|                    | +FAP    | 71.73   | 54.79        | 93.67      | 86.17        | 88.53      | 72.31        | 80.49   | 52.46        | 72.83  | 55.80        | 61.94 | 45.80        | 85.06      | 72.07        | 63.41        | 37.78        | 27.93        | 15.96        |
|                    | +DEFEAT | 74.75   | <b>69.33</b> | 95.17      | <b>93.02</b> | 89.78      | <b>83.02</b> | 76.63   | <b>68.88</b> | 79.57  | <b>74.23</b> | 64.13 | <b>58.39</b> | 91.64      | <b>87.09</b> | 70.71        | <b>63.11</b> | 30.33        | <b>26.91</b> |
| $\epsilon = 4/255$ | Method  | Average |              | Caltech101 |              | OxfordPets |              | EuroSAT |              | UCF101 |              | DTD   |              | Flowers102 |              | StanfordCars |              | FGVCAircraft |              |
|                    |         | Acc.    | Rob.         | Acc.       | Rob.         | Acc.       | Rob.         | Acc.    | Rob.         | Acc.   | Rob.         | Acc.  | Rob.         | Acc.       | Rob.         | Acc.         | Rob.         | Acc.         | Rob.         |
|                    | FARE    | 42.84   | 15.87        | 86.00      | 54.69        | 77.60      | 24.64        | 16.28   | 10.88        | 43.01  | 8.96         | 31.26 | 15.19        | 39.22      | 7.75         | 39.17        | 3.56         | 10.20        | 1.32         |
|                    | +APT    | 64.86   | 29.04        | 92.41      | 66.61        | 81.36      | 30.36        | 57.42   | 24.96        | 67.96  | 24.45        | 53.31 | 23.05        | 84.49      | 42.55        | 57.08        | 14.99        | 24.84        | 5.37         |
| PGD                | +FAP    | 58.59   | 30.53        | 85.48      | 66.41        | 77.41      | 39.22        | 59.53   | 23.72        | 59.19  | 26.38        | 46.99 | 25.18        | 70.44      | 41.41        | 47.62        | 13.52        | 22.05        | 8.43         |
|                    | +DEFEAT | 65.53   | <b>52.46</b> | 92.54      | <b>86.57</b> | 84.19      | <b>52.79</b> | 63.77   | <b>52.02</b> | 69.10  | <b>55.30</b> | 52.19 | <b>41.84</b> | 80.63      | <b>68.49</b> | 58.03        | <b>42.72</b> | 23.79        | <b>19.98</b> |

### C.8 GENERALIZATION TO ALTERNATIVE VLM BACKBONES

To demonstrate that the effectiveness of DEFEAT is not limited to the ViT-B/32 CLIP backbone used in our main paper, we conduct additional experiments on two distinct and powerful backbones: a larger CLIP model (ViT-L/14) and a model from a different pre-training family (EVA-CLIP (Sun et al., 2023)). The following results validate that DEFEAT is a broadly applicable framework.

First, to assess the scalability of our method, we replace the TeCoA-trained ViT-B/32 (Mao et al., 2023) backbone with the much larger TeCoA-trained ViT-L/14. As shown in Table A10, DEFEAT continues to provide significant robustness gains on this more powerful model. Those results demonstrate that our method scales effectively with model size.

Next, we evaluate DEFEAT’s adaptability by applying it to EVA-CLIP (Sun et al., 2023), which represents a different pre-training paradigm compared to the original OpenAI CLIP. The results in Table A11 show that DEFEAT successfully enhances the robustness of this distinct model family. Those results validate that DEFEAT can serve as a general-purpose robust prompt tuning framework.

Table A10: Performance on the TeCoA-trained CLIP ViT-L/14 (Mao et al., 2023) backbone against PGD attack ( $\epsilon = 4/255$ ) under the adversarial few-shot (16-shot) classification setting.

| $\epsilon = 4/255$ | Method  | Caltech101 |              |              | DTD   |              |              | Flowers102 |              |              | Average |              |              |
|--------------------|---------|------------|--------------|--------------|-------|--------------|--------------|------------|--------------|--------------|---------|--------------|--------------|
|                    |         | Acc.       | Rob.         | H            | Acc.  | Rob.         | H            | Acc.       | Rob.         | H            | Acc.    | Rob.         | H            |
| PGD                | CLIP    | 94.81      | 17.85        | 30.04        | 53.72 | 0.59         | 1.17         | 79.46      | 0.45         | 0.89         | 76.00   | 6.30         | 11.63        |
|                    | TeCoA   | 87.38      | 75.98        | 81.28        | 34.63 | 28.19        | 31.08        | 39.02      | 25.38        | 30.76        | 53.68   | 43.18        | 47.86        |
|                    | +APT    | 95.42      | 86.77        | 90.89        | 57.21 | 44.92        | 50.33        | 88.67      | 76.53        | 82.15        | 80.43   | 69.41        | 74.51        |
|                    | +DEFEAT | 95.70      | <b>91.81</b> | <b>93.71</b> | 57.57 | <b>46.81</b> | <b>51.64</b> | 85.95      | <b>79.01</b> | <b>82.33</b> | 79.74   | <b>72.54</b> | <b>75.97</b> |

Table A11: Performance on the EVA-CLIP (Sun et al., 2023) backbone against PGD attack ( $\epsilon = 1/255$ ) under the adversarial few-shot (16-shot) classification setting.

| $\epsilon = 1/255$ | Method   | Caltech101 |              |              | DTD   |             |              | Flowers102 |             |             | Average |             |              |
|--------------------|----------|------------|--------------|--------------|-------|-------------|--------------|------------|-------------|-------------|---------|-------------|--------------|
|                    |          | Acc.       | Rob.         | H            | Acc.  | Rob.        | H            | Acc.       | Rob.        | H           | Acc.    | Rob.        | H            |
| PGD                | EVA-CLIP | 97.20      | 3.81         | 7.33         | 49.41 | 2.78        | 5.26         | 75.92      | 0.81        | 1.60        | 74.18   | 2.47        | 4.77         |
|                    | +APT     | 97.85      | 5.23         | 9.93         | 74.70 | 2.60        | 5.03         | 98.05      | 2.64        | 5.14        | 90.20   | 3.49        | 6.72         |
|                    | +DEFEAT  | 96.80      | <b>11.70</b> | <b>20.88</b> | 71.28 | <b>5.61</b> | <b>10.40</b> | 95.86      | <b>4.47</b> | <b>8.54</b> | 87.98   | <b>7.26</b> | <b>13.41</b> |

### C.9 ROBUST VISUAL BACKBONE RELIANCE

In this section, we analyse the reliance of the robust visual backbone on the proposed DEFEAT method. While DEFEAT significantly enhances adversarial robustness under few-shot prompt tuning conditions, it relies heavily on pre-trained robust visual backbones (e.g., TeCoA (Mao et al., 2023)) like other adversarial prompt tuning methods (e.g., APT (Li et al., 2024) and FAP (Zhou et al., 2024)). The results in Table A12 show that APT failed to improve robustness when not using TeCoA. DEFEAT and FAP show less dependence on TeCoA when the perturbation size is small (i.e.,  $\epsilon = 1/255$ ). However, using the robust image encoder TeCoA still substantially enhances the robustness of both FAP and DEFEAT. When increasing perturbation size to  $\epsilon = 4/255$ , we find that all methods, i.e., APT, FAP, and DEFEAT) exhibit near-zero robustness or even collapse during training, a phenomenon also reported in Li et al. (2024); Zhou et al. (2024). The reasons may be the considerably smaller parameter space available for tuning compared to fine-tuning the entire model. Those results show that a robust visual backbone is necessary for adversarial prompt tuning methods.

Table A12: Performance of adversarial prompt tuning methods with (w/) or without (w/o) TeCoA on Flowers102 and DTD datasets for 16 shots under the adversarial few-shot classification setting when  $\epsilon = 1/255$ . The numbers in the parentheses denote the improvement or decline compared to the methods without TeCoA.

| Dataset    | $\epsilon = 1/255$ | TeCoA | APT            |                | FAP            |                | DEFEAT                |                       |
|------------|--------------------|-------|----------------|----------------|----------------|----------------|-----------------------|-----------------------|
|            |                    |       | w/o TeCoA      | w/ TeCoA       | w/o TeCoA      | w/ TeCoA       | w/o TeCoA             | w/ TeCoA              |
| Flowers102 | Acc.               | 48.15 | 85.06 (+36.91) | 86.80 (+38.65) | 80.15 (+32.00) | 82.46 (+34.31) | <b>91.72 (+43.57)</b> | 84.90 (+36.75)        |
|            | Rob.               | 32.68 | 0.85 (-31.83)  | 73.57 (+40.89) | 51.36 (+18.68) | 66.63 (+33.95) | 55.99 (+23.31)        | <b>81.81 (+49.13)</b> |
|            | H                  | 38.93 | 1.68 (-37.25)  | 79.64 (+40.71) | 62.60 (+23.67) | 73.70 (+34.77) | 69.53 (+30.60)        | <b>83.33 (+44.40)</b> |
| DTD        | Acc.               | 31.97 | 60.11 (+28.14) | 55.26 (+23.29) | 55.50 (+23.53) | 56.68 (+24.71) | <b>65.13 (+33.16)</b> | 57.15 (+25.18)        |
|            | Rob.               | 23.88 | 2.60 (-21.28)  | 39.54 (+15.66) | 28.67 (+4.79)  | 41.67 (+17.79) | 30.91 (+7.03)         | <b>52.66 (+28.78)</b> |
|            | H                  | 27.34 | 4.98 (-22.35)  | 46.10 (+18.76) | 37.81 (+10.47) | 48.03 (+20.69) | 41.92 (+14.58)        | <b>54.81 (+27.47)</b> |

### C.10 COMPUTATIONAL COST ANALYSIS

To evaluate the practicality and efficiency of DEFEAT, we provide a comprehensive analysis of its computational overhead compared to APT and FAP. For a fair comparison, we adopt the same experimental settings as APT (e.g., training epochs, schedules, and data augmentation), while for FAP, we maintain its original implementation to ensure a faithful reproduction of its results. All experiments were conducted on the ImageNet dataset under the adversarial few-shot (16-shot) classification setting with  $\epsilon = 4/255$ .

The results in Table A13 show that DEFEAT achieves a highly favorable trade-off between computational cost and robustness.

1134     **Training Memory.** The training memory required by DEFEAT (21240M) is nearly identical to  
 1135     that of APT (21204M) and significantly lower than that of FAP (37500M). This indicates that our  
 1136     PerturbShield module adds minimal GPU memory overhead.

1137     **Training Time.** The training time of DEFEAT (19.21h) is negligibly longer than APT’s (18.99h) and  
 1138     is less than half that of FAP (39.04h), demonstrating a clear efficiency advantage.

1140     **Inference Time.** During inference, the per-image processing time for DEFEAT (3.1431ms) is only a  
 1141     minor increase of 0.15ms over APT (2.9959ms), which is insignificant for practical applications.

1142     In summary, the computational overhead introduced by our method is minimal during both training  
 1143     and inference. This slight cost is a highly favorable trade-off for the substantial improvement in  
 1144     robustness that DEFEAT provides over both APT and FAP, showing the practicality and efficiency of  
 1145     our method.

1146     Table A13: Computational cost and performance comparison on ImageNet under the adversarial  
 1147     few-shot (16-shot) classification setting with  $\epsilon = 4/255$ .

| Method        | Train.Memory | Train.time (h) | Infer.time (ms) | Acc.         | Rob.         | H            |
|---------------|--------------|----------------|-----------------|--------------|--------------|--------------|
| APT           | 21204M       | 18.99          | 2.9959          | 41.06        | 12.02        | 18.60        |
| DEFEAT (ours) | 21240M       | 19.21 (+0.22)  | 3.1431          | <b>41.09</b> | <b>18.42</b> | <b>25.44</b> |
| FAP           | 37500M       | 39.04          | 3.0020          | 40.45        | 12.03        | 18.54        |

### C.11 PERFORMANCE COMPARISON UNDER THE $\ell_2$ THREAT MODEL

We extend our experiments to the  $\ell_2$  threat model, evaluating all methods against the PGD- $\ell_2$  ( $\epsilon = 0.5$ ) attack across 10 downstream datasets.

As shown in Table A14, DEFEAT performs well under the  $\ell_2$  threat model. Our method achieves an average robustness of 50.24%, surpassing both APT (47.34%) and FAP (47.07%). This result demonstrates that the robustness gains provided by DEFEAT are generalizable and not restricted to the  $\ell_\infty$  norm.

Table A14: Performance Comparison under the adversarial few-shot (16-shot) classification setting with the  $\ell_2$  threat model.

| PGD- $\ell_2$ | Average |              |              | Caltech101   |              |              | OxfordPets   |              |              | StanfordCars |              |              | Flowers102   |              |              | Food101      |              |              | FGVCAircraft |              |              | SUN397       |              |      | DTD  |      |      | EuroSAT |      |  | UCF101 |  |  |
|---------------|---------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|------|------|------|------|---------|------|--|--------|--|--|
| Method        | Acc.    | Rob.         | H            | Acc.         | Rob.         | Acc. | Rob. | Acc. | Rob. | Acc.    | Rob. |  |        |  |  |
| TeCoA         | 33.03   | 29.17        | 30.98        | 78.78        | 74.40        | 66.26        | 59.14        | 10.32        | 7.80         | 30.13        | 26.59        | 23.52        | 19.34        | 7.17         | 5.07         | 33.24        | 28.13        | 24.29        | 22.51        | 19.56        | 17.73        | 36.98        | 31.03        |      |      |      |      |         |      |  |        |  |  |
| APT           | 52.09   | 47.34        | 49.60        | 86.29        | 84.42        | 67.29        | 61.43        | 31.60        | 29.30        | <b>76.41</b> | 68.69        | 30.39        | 26.08        | <b>20.31</b> | 15.99        | 45.21        | 39.85        | <b>45.86</b> | 40.01        | <b>64.33</b> | 58.90        | 53.16        | 48.72        |      |      |      |      |         |      |  |        |  |  |
| FAP           | 52.22   | 47.07        | 49.51        | 80.04        | 77.36        | 69.09        | 61.43        | <b>41.06</b> | <b>33.33</b> | 66.10        | 61.63        | <b>45.60</b> | <b>38.26</b> | 18.57        | 15.48        | 43.66        | 38.17        | 42.85        | 39.36        | 64.51        | 59.78        | 50.67        | 45.94        |      |      |      |      |         |      |  |        |  |  |
| DEFEAT        | 51.97   | <b>50.24</b> | <b>51.09</b> | <b>87.63</b> | <b>86.41</b> | <b>72.17</b> | <b>67.10</b> | 32.45        | 30.93        | 71.13        | <b>69.43</b> | 33.57        | 31.99        | 18.03        | <b>16.65</b> | <b>46.26</b> | <b>44.14</b> | 44.15        | <b>42.61</b> | 58.88        | <b>58.85</b> | <b>55.43</b> | <b>54.27</b> |      |      |      |      |         |      |  |        |  |  |

### C.12 COMBINATION WITH TEST-TIME DEFENSE STRATEGY

To address the scenario facing strong attacks without a pre-trained robust backbone, we propose **DEFEAT-T**, a variant of the DEFEAT method. DEFEAT-T successfully adapts our core idea to the non-robust backbone scenario by combining it with a test-time defense strategy (Shu et al., 2022; Wang et al., 2025; Sheng et al., 2025).

The core of DEFEAT-T is to pivot our method from an “active defender” (via adversarial training) to a “passive purifier” (via clean training), combined with an efficient, optimization-free test-time ensemble. The mechanism is structured into two phases:

**Phase 1: Clean Training.** We first recognized that adversarial training on a non-robust backbone could be a source of the training collapse. Therefore, we discard adversarial training. Instead, we train our DEFEAT framework (including the learnable prompt and PerturbShield module) on a standard (non-robust) backbone using clean data only. In this phase, our PerturbShield module learns to reconstruct the clean feature distribution of the non-robust backbone.

**Phase 2: Test-Time Ensemble Defense.** At test time, we execute forward-pass ensemble defense, avoiding per-sample backpropagation or optimization, which are major bottlenecks for test-time

prompt tuning methods (Shu et al., 2022; Wang et al., 2025; Sheng et al., 2025). First, we generate  $N^{\text{aug}}$  (e.g.,  $N^{\text{aug}}=64$ ) augmented views (e.g., random crop, random horizontal flip) of the test sample. Second, all  $N^{\text{aug}}$  views are fed through our model trained from clean data, to get  $N^{\text{aug}}$  fused logits by Eq. (7) as shown in Figure 4. Next, to obtain the final prediction, one can directly average these  $N^{\text{aug}}$  fused logits. Here, to have better performance, we perform a reliability-based weighted average. Specifically, we use the  $N^{\text{aug}}$  class features (which are collected alongside the grid features) to compute a similarity matrix. We then average the top- $k$  (e.g.,  $k=20$ ) values in each row to get a reliability score for each augmented view. Finally, we use these reliability scores as weights for a weighted average of the corresponding fused logits to obtain the final prediction.

The rationale for this design is our hypothesis that this view augmentation and weighted-average prediction could mitigate the impact of unseen adversarial perturbations, and that combining it with our DEFEAT framework provides a second stage of purification.

Table A15: Combination adversarial prompt tuning methods with test-time defense strategy under the adversarial few-shot classification setting.

| 16shot, $\epsilon = 4/255$                   | Method               | Average      |              |              | Caltech101   |              | OxfordPets   |              | Flowers102   |              | DTD          |              |
|--|----------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|  |                      | Acc.         | Rob.         | H            | Acc.         | Rob.         | Acc.         | Rob.         | Acc.         | Rob.         | Acc.         | Rob.         |
| zero-shot                                    | CLIP                 | 72.38        | 2.28         | 4.42         | 91.03        | 7.79         | 87.38        | 1.20         | 66.91        | 0.12         | 44.21        | 0.00         |
|  | TeCoA                | 49.87        | 19.86        | 28.41        | 78.78        | 43.61        | 66.26        | 15.56        | 30.13        | 8.93         | 24.29        | 11.35        |
| Test-time prompt tuning                      | R-TPT                | 70.23        | 50.86        | 59.00        | 90.83        | 77.00        | 85.12        | 57.32        | 62.28        | 39.46        | 42.67        | 29.67        |
|  | APT w/o TeCoA        | 14.05        | 0.07         | 0.14         | 43.77        | 0.16         | 3.38         | 0.00         | 3.45         | 0.00         | 5.61         | 0.12         |
|  | DEFEAT w/o TeCoA     | 18.45        | 0.16         | 0.32         | 51.89        | 0.41         | 3.05         | 0.00         | 9.62         | 0.00         | 9.22         | 0.24         |
| Train-time prompt tuning                     | APT (w/ TeCoA)       | 68.96        | 33.94        | 45.49        | 86.29        | 56.75        | 67.29        | 19.98        | 76.41        | 37.52        | 45.86        | 21.51        |
|  | DEFEAT (w/ TeCoA)    | 68.77        | 49.92        | 57.84        | 87.63        | 78.30        | 72.17        | 30.77        | 71.13        | 58.26        | 44.15        | 32.33        |
| Train-time prompt tuning + test-time defense | APT-T (w/o TeCoA)    | 86.72        | 43.33        | 57.78        | 95.01        | 71.08        | 90.19        | 32.54        | <b>94.76</b> | 39.18        | 66.90        | 30.50        |
|  | DEFEAT-T (w/o TeCoA) | <b>87.76</b> | <b>75.96</b> | <b>81.43</b> | <b>96.23</b> | <b>93.63</b> | <b>91.88</b> | <b>86.54</b> | 94.60        | <b>66.38</b> | <b>68.32</b> | <b>57.27</b> |

**Empirical Verification.** As shown in Table A15, DEFEAT-T (w/o TeCoA) achieves good average performance (87.76% Acc., 75.96% Rob., 81.43% H), surpassing test-time prompt tuning methods like R-TPT (70.23% Acc., 50.86% Rob., 59.00% H), and the standard DEFEAT (w/ TeCoA) (68.77% Acc., 49.92% Rob., 57.84% H) even without a robust backbone.

Simultaneously, we also extend the baseline APT to APT-T (w/o TeCoA) (i.e., APT trained on clean data, evaluated with the same test-time ensemble but without the logits fusion strategy and the PerturbShield module). We find that APT-T (w/o TeCoA) (43.33% Rob.) achieves performance comparable to the original DEFEAT (w/ TeCoA) (49.92% Rob.). However, its robustness is still significantly weaker than our DEFEAT-T (75.96% Rob.). This comparison further demonstrates the effectiveness of our proposed logits fusion strategy and PerturbShield.

Crucially, our DEFEAT-T method has a decisive efficiency advantage over existing test-time prompt tuning methods (Shu et al., 2022; Wang et al., 2025; Sheng et al., 2025). Existing methods require per-sample optimization (batch size=1) involving backpropagation, making them slow. Our DEFEAT-T is an inference-only scheme with no gradient calculations, allowing for batch parallelism and much faster inference. However, even so, the time cost of DEFEAT-T’s inference is still much higher than that of the original DEFEAT. As shown in Table A16, the inference time of DEFEAT-T is approximately 2x slower than that of DEFEAT. On the other hand, the training time for DEFEAT-T is much shorter, as it does not require adversarial training.

In summary, DEFEAT and DEFEAT-T can be adopted in different scenarios. The original DEFEAT (w/ TeCoA) is ideal for low-latency inference scenarios, whereas DEFEAT-T is a perfectly viable and highly robust alternative when a robust backbone is not available.

Table A16: Computational cost and performance comparison between DEFEAT, DEFEAT-T, and R-TPT (Sheng et al., 2025) on Caltech101 under the adversarial few-shot (16-shot) classification setting with  $\epsilon = 4/255$ .

| Method               | Train.Time (s) | Infer.Memory    | Infer.Time (s/per img) | Acc.         | Rob..        | H            |
|----------------------|----------------|-----------------|------------------------|--------------|--------------|--------------|
| DEFEAT (w/ TeCoA)    | 3597           | 8481M           | <b>0.09</b>            | 93.39        | 91.12        | 92.24        |
| DEFEAT-T (w/o TeCoA) | 1079           | 10420M          | 0.19                   | <b>96.19</b> | <b>95.58</b> | <b>95.88</b> |
| R-TPT                | 0              | 5248M (batch=1) | 6.63                   | 90.83        | 86.13        | 88.42        |

1242 C.13 THEORETICAL JUSTIFICATION  
1243

1244 Here we provide a theoretical justification for why clean and adversarially perturbed features can be  
1245 mapped to the same discrete representation within the VQ-VAE framework.

1246 Let  $\text{Enc}(\cdot)$  denote the composite visual encoder, which consists of the CLIP image encoder followed  
1247 by the trainable VQ-VAE encoder. The clean latent features output by this composite encoder are  
1248 denoted by  $z_e^c = \text{Enc}(x) \in \mathbb{R}^{n \times d}$  (where  $n$  is the number of latent positions after encoding,  $d$  is  
1249 the dimensionality of the latent embedding). Under an adversarial attack with perturbation  $\delta$ , the  
1250 adversarial latent features become:

$$1251 \quad z_e^a = \text{Enc}(x + \delta) = z_e^c + \Delta z, \quad (15)$$

1252 where  $\Delta z \in \mathbb{R}^{n \times d}$  is the mapped perturbation in the latent space.

1253 The core of VQ-VAE is to discretize continuous latent features through a codebook  $\mathcal{E} = \{e_k\}_{k=1}^K$   
1254 (where  $K$  is the codebook size). The covering radius  $r$  of the codebook is defined as

$$1255 \quad r = \max_{z \in \mathcal{Z}} \min_{e_k \in \mathcal{E}} \|z - e_k\|_2, \quad (16)$$

1256 where  $\mathcal{Z}$  is the set of latent features output by the encoder from all clean images. The covering radius  
1257  $r$  characterizes the “coverage capability” of the codebook over the clean feature space: the distance  
1258 from any clean feature to its nearest codebook element does not exceed  $r$ , and  $r$  is minimized through  
1259 the codebook loss of VQ-VAE.

1260 Based on the above definition, if the perturbation magnitude in the latent space satisfies  $\|\Delta z\|_2 \leq r$ ,  
1261 it is statistically highly probable that both the clean feature  $z_e^c$  and the adversarial feature  $z_e^a$  map to  
1262 the same codebook element  $e_k \in \mathcal{E}$ . Notably, the training objective in DEFEAT encourages images  
1263 that are more similar at the pixel level to be closer in feature space. So the feature shift  $\Delta z$  induced  
1264 by minor perturbations is typically much smaller than the codebook covering radius  $r$ , effectively  
1265 mitigating such adversarial attacks during the quantization process.

1266 **Empirical Verification.** To validate this, we conduct a statistical analysis on the latent space using  
1267 the Caltech101 dataset (16-shot,  $\epsilon = 4/255$ ). We measure two key metrics across all image patches:

- 1268 • The average magnitude of adversarial perturbation per patch:  $\|\Delta z\|_2 \approx 126.0$ .
- 1269 • The average distance between clean latent features and their assigned codebook vectors  
1270 (serving as a proxy for the safety radius):  $r \approx 432.3$ .

1271 The experimental results show that the covering radius is approximately 3.43 times larger than the  
1272 perturbation magnitude (432.3 vs. 126.0). This substantial margin implies that the vast majority  
1273 of adversarial perturbations are insufficient to push the latent features out of their original discrete  
1274 region defined by their assigned code. Consequently, the clean and adversarial patches are quantized  
1275 to the same discrete codes, confirming the effectiveness of our defense.

1276 C.14 STABILITY ANALYSIS WITH MULTI-SEED EVALUATION  
1277

1278 To ensure the reproducibility of our results and confirm that the reported performance is not an  
1279 artifact of a single run, we conduct a stability analysis by re-running our primary adversarial few-shot  
1280 classification experiments. Specifically, we evaluate the models under the 16-shot setting with PGD  
1281 attack ( $\epsilon = 4/255$ ) using three random seeds (0, 1, and 2).

1282 Table A17 compares our originally reported results (derived from Table 1 and Table A6) with the  
1283 new averages calculated across three seeds. As observed, the multi-seed averages align closely with  
1284 our reported single-seed results. Crucially, the robustness advantage of DEFEAT remains consistent  
1285 and stable. For instance, the multi-seed average harmonic mean for DEFEAT is 41.66%, which still  
1286 substantially surpasses the strongest baselines, APT (28.90%) and FAP (28.07%). This consistency  
1287 confirms that our results are reproducible and statistically robust.

1288 C.15 STABILITY ANALYSIS AGAINST PROMPT VARIATIONS  
1289

1290 As shown in Table A19, DEFEAT exhibits remarkable stability. The robustness fluctuates only  
1291 slightly (between 31.09% and 32.33%) across different prompt variations.

1296 Table A17: Comparison between originally reported results and averages over 3 seeds  
1297

| 1298<br>1299<br>1300<br>1301<br>1302<br>1303 | 16shot, $\epsilon = 4/255$<br>Method | Average (reported) |              |              | Average on 3 seeds |              |   |
|--|--------------------------------------|--------------------|--------------|--------------|--------------------|--------------|---|
|  |                                      | Acc.               | Rob.         | H            | Acc.               | Rob.         | H |
| TeCoA  | 33.67                                | 10.79              | 16.34        | 33.67        | 10.78              | 16.33        |   |
| APT  | <b>51.08</b>                         | 20.26              | 29.02        | <b>51.27</b> | 20.12              | 28.90        |   |
| FAP  | 50.50                                | 19.82              | 28.47        | 50.68        | 19.41              | 28.07        |   |
| DEFEAT                                       | 50.98                                | <b>36.05</b>       | <b>42.24</b> | 50.85        | <b>35.28</b>       | <b>41.66</b> |   |

1304 Table A18: Complete results calculated by averaging three seeds across all datasets. (16shot,  
1305  $\epsilon = 4/255$ )

| 1307<br>Method | ImageNet     | Caltech101   | OxfordPets   | StanfordCars | Flowers102   | Food101      | FGVCAircraft | SUN397       | DTD          | EuroSAT      | UCF101       |              |
|----------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| 1308<br>Acc.   | 1309<br>Acc. | 1310<br>Acc. | 1311<br>Acc. | 1312<br>Acc. | 1313<br>Acc. | 1314<br>Acc. | 1315<br>Acc. | 1316<br>Acc. | 1317<br>Acc. | 1318<br>Acc. | 1319<br>Acc. |              |
| 1308<br>Rob.   | 1309<br>Rob. | 1310<br>Rob. | 1311<br>Rob. | 1312<br>Rob. | 1313<br>Rob. | 1314<br>Rob. | 1315<br>Rob. | 1316<br>Rob. | 1317<br>Rob. | 1318<br>Rob. | 1319<br>Rob. |              |
| TeCoA          | 40.11        | 10.15        | 78.78        | 43.57        | 66.26        | 15.49        | 10.32        | 0.94         | 30.13        | 8.94         | 23.52        | 3.26         |
| APT            | <b>40.98</b> | 12.22        | 86.38        | 56.79        | 68.25        | 19.82        | 33.08        | 7.74         | <b>76.24</b> | 37.18        | 30.19        | 7.87         |
| FAP            | 40.53        | 12.10        | 80.36        | 55.50        | 69.27        | 25.59        | <b>40.91</b> | 6.60         | 66.11        | 30.96        | <b>45.18</b> | 8.30         |
| DEFEAT         | 40.79        | <b>18.02</b> | <b>86.96</b> | <b>77.61</b> | <b>72.12</b> | <b>32.31</b> | 32.94        | <b>20.38</b> | 71.00        | <b>58.13</b> | 33.60        | <b>20.50</b> |
|                |              |              |              |              |              |              |              |              | 16.82        | <b>11.32</b> | 46.23        | <b>29.61</b> |
|                |              |              |              |              |              |              |              |              | <b>44.41</b> | <b>32.57</b> | 58.85        | <b>45.35</b> |
|                |              |              |              |              |              |              |              |              |              |              | <b>55.59</b> | <b>42.34</b> |

1312 Crucially, the performance trends remain consistent across all variations. Even under the least optimal  
1313 prompt setting (“{} texture.”), DEFEAT (31.09% Rob.) significantly outperforms the strongest  
1314 baseline, APT (20.26% Rob.), by a margin of +10.83%. This confirms that our method’s superiority  
1315 is intrinsic to its design rather than an artifact of specific prompt engineering.

1317 Table A19: Stability analysis of DEFEAT against different hand-crafted prompt templates on the  
1318 DTD dataset.

| 1320<br>Method                                    | 16shot, $\epsilon = 4/255$ | DTD   |       |       |
|---|----------------------------|-------|-------|-------|
|   |                            | Acc.  | Rob.  | H     |
| TeCoA   |                            | 33.67 | 10.79 | 16.34 |
| APT   |                            | 51.08 | 20.26 | 29.02 |
| FAP   |                            | 50.50 | 19.82 | 28.47 |
| <b>DEFEAT with different hand-crafted prompts</b> |                            |       |       |       |
| “a photo of a {} , a type of texture.”            |                            | 44.15 | 32.33 | 37.33 |
| “{} texture.”                                     |                            | 44.74 | 31.09 | 36.69 |
| “a photo of {}.”                                  |                            | 44.03 | 31.80 | 36.93 |

## D USE OF LARGE LANGUAGE MODELS

1333 The LLM served as an assistive tool for language editing and manuscript polishing. Its usage was  
1334 confined to improving grammar, refining phrasing for clarity, and ensuring idiomatic English.

1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349