

ADVFLYP: ADVERSARIALLY FINETUNE LIKE YOU PRETRAIN FOR ZERO-SHOT ROBUSTNESS OF CLIP

Anonymous authors

Paper under double-blind review

ABSTRACT

Pretrained vision-language models (VLMs) like CLIP are shown to be highly susceptible to adversarial perturbations. Adversarial finetuning (AFT) approaches have been proposed to improve the zero-shot adversarial robustness of CLIP on various downstream tasks, based on finetuning the vision encoder on adversarial images generated from a proxy classification dataset, such as TinyImageNet. However, we demonstrate that existing AFT approaches have largely overlooked the important role of the training recipe, particularly the training data and objective. To this end, we propose *Adversarially Finetune Like You Pretrain* (AdvFLYP), which practically retains the training recipe of CLIP’s pretraining during AFT. We finetune CLIP based on adversarial images generated from web-scale image-text data with a contrastive loss. Experiments validate the superiority of AdvFLYP on various downstream datasets. For example, AdvFLYP outperforms existing AFT approaches finetuned on TinyImageNet (ImageNet) by 19.1% (3.1%), averaged on 14 downstream datasets. Further analyses show that sufficiently large training data amounts and batch sizes are crucial for the contrastive learning of AdvFLYP. Our code and model checkpoints will be released.

1 INTRODUCTION

Pre-trained vision-language models (VLMs) have been trained to align images with their descriptive texts over significant amounts of image-text pairs, with CLIP (Radford et al., 2021) and ALIGN (Jia et al., 2021) being notable representatives. These models have exhibited remarkable abilities to perform image classification in a zero-shot manner (Pratt et al., 2023; Saha et al., 2024; Sammani & Deligiannis, 2024). However, recent studies have revealed the alarming vulnerabilities of CLIP to adversarial attacks (Mao et al., 2023; Li et al., 2024; Schlarmann et al., 2024): An imperceptible maliciously manipulated noise added to a test image suffices to substantially reduce the model’s recognition accuracy.

To enhance CLIP’s robustness to adversarial attacks, recent studies introduce adversarial images generated from a proxy dataset such as TinyImageNet (Le & Yang, 2015) or ImageNet (Deng et al., 2009) into the training set, and finetune the vision encoder (Mao et al., 2023; Wang et al., 2024; Yu et al., 2024; Schlarmann et al., 2024) based on adversarial training (AT) (Madry et al., 2018; Zhang et al., 2019). This practice proves effective in improving the model’s adversarial robustness on diverse downstream datasets without further training, which is termed *zero-shot robustness* (Mao et al., 2023). However, these methods incur a significant degradation in the model’s generalization on clean data of downstream tasks. We hypothesize that such degradation is largely due to the misaligned training recipe between CLIP’s pretraining and the finetuning process of existing AFT methods. Intuitively, there is a fundamental difference between adversarially finetuning CLIP and robustly training a model from scratch. CLIP has been pre-trained over web-scale image-text pairs and learned real-world knowledge, and updating its model weights on a specific domain can already lead to noticeable generalization loss (Radford et al., 2021), which further complicates the analysis of this loss in adversarial finetuning. In this work, we investigate the generalization degradation in the adversarial finetuning of CLIP, and identify two important factors: **(1)** the training data distribution that differs from CLIP’s pretraining data. This is evidenced by the following observations: (i) Finetuning CLIP on the clean data of a proxy dataset lowers the accuracy on downstream datasets, and (ii) Finetuning CLIP on adversarial images of a proxy dataset (Mao et al., 2023; Wang et al., 2024; Yu et al., 2024) results in higher accuracy on the clean test set of the proxy dataset than the

original CLIP, which indicates that the model has overfit to the distribution of the proxy dataset, even finetuned with adversarial images; (2) the training objective for AFT. Minimizing the cross-entropy loss between the generated adversarial images and their correct labels on a classification dataset effectively aligns multiple images from a class to the same textual class name, which causes the loss of knowledge.

Built upon our analysis, we propose a simple yet effective paradigm termed *Adversarially Finetune Like You Pretrain* (AdvFLYP). The main idea of AdvFLYP is to finetune CLIP with adversarial images while maximally retaining the same training recipe as employed in the pretraining phase. Specifically, to imitate the distribution of CLIP’s training data, we randomly sample a certain number of web-scale image-text pairs. During AFT, we employ the same contrastive loss for pretraining CLIP. The difference is that we align adversarial images, rather than clean images, with their corresponding texts. To further alleviate the robustness-generalization trade-off, we propose to impose logit- and feature-level regularization during finetuning, which we show to improve robustness transfer and generalization on clean images, respectively. Through extensive experiments, we show that when finetuned with the same amounts of training data, AdvFLYP outperforms existing AFT methods finetuned on TinyImageNet and ImageNet by an average relative improvement of 19.1% and 3.1%, respectively. To facilitate understanding of this paradigm, we vary the training conditions (e.g., batch size, training data amount) and provide insights into contrastive learning of CLIP in the context of adversarial finetuning. We summarize the contributions of this work as follows:

- We investigate the generalization degradation in existing AFT methods for CLIP, and identify two major sources, which are training data distribution and the training objective.
- We introduce *Adversarially Finetune Like You Pretrain* (AdvFLYP), a simple yet effective paradigm to achieve zero-shot adversarial robustness, which resumes contrastive learning of CLIP by aligning adversarial images with their texts.
- Extensive experiments on 14 downstream datasets show that AdvFLYP outperforms mainstream AFT methods. We also vary the training setting of AdvFLYP and show that sufficiently large training data amounts and batch size for AdvFLYP are crucial for robustness and accuracy.

2 RELATED WORK

Adversarial robustness of neural networks. Deep neural networks (Krizhevsky et al., 2012) are vulnerable to adversarial attacks (Carlini & Wagner, 2017; Szegedy et al., 2014): an imperceptible pixel-level perturbation added to the test image can mislead a well-trained model to make a wrong prediction. Adversarial attacks (Carlini & Wagner, 2017; Croce & Hein, 2020) and defences (Madry et al., 2018) have been extensively studied. Among defence methods, adversarial training (AT) (Madry et al., 2018; Zhang et al., 2019; Rice et al., 2020) has been established as the *de-facto* standard to train an adversarially robust model. More recent research finetunes a standardly trained model on adversarial samples to enhance its adversarial robustness, instead of training a robust model from scratch (Suzuki et al., 2023).

Adversarial robustness of vision-language models (VLMs) has also attracted significant research attention (Zhao et al., 2023). In this paper, we focus on *zero-shot adversarial robustness* of CLIP (Radford et al., 2021). Existing methods are largely based on AT, introducing adversarial images into the training set and adapt the CLIP models. There are two types of categories in this regard: *adversarial finetuning* (AFT) (Mao et al., 2023; Wang et al., 2024; Yu et al., 2024; Schlarmann et al., 2024), which finetunes the vision encoder of CLIP; and *adversarial prompt tuning* (Li et al., 2024; Zhang et al., 2024), which learns tunable prompt at the text encoder side to align with adversarial images. More recently, test-time methods for defending CLIP have started to garner interests (Wang et al., 2025; Sheng et al., 2025; Tong et al., 2025; Zhang et al., 2025; Xing et al., 2025), which achieves inference-time robustness without the need for training. We focus on *adversarial finetuning* methods in this work, which is still the most effective approach. Mao et al. (2023) first propose to generate adversarial images on ImageNet (Deng et al., 2009) by maximizing the cross-entropy loss *w.r.t.* the ground-truth label, which are then leveraged for finetuning vision encoder f_θ by minimizing the cross-entropy loss of these adversarial images *w.r.t.* the labels. Subsequent work introduces regularization based on this loss. Wang et al. (2024) impose logit-level regularization terms guided by frozen CLIP models to improve robustness on downstream datasets. Yu et al. (2024)

introduce regularization formulated by aligning text-guided attention of the model with the original CLIP. More recently, Dong et al. (2025) focus on improving the adversarial candidates in adversarial finetuning by forming consecutive vertices and sampling simplices. Aside from supervised AFT, recent work proposes unsupervised AFT (Schlarmann et al., 2024). Gong et al. (2025) employ unsupervised AFT as a novel tool to improve interpretability of visual models. These methods employ a proxy dataset and a training objective that differ from CLIP’s pretraining. We propose a novel AFT paradigm, AdvFLYP, which challenges the common practice of current AFT methods.

3 METHOD

In this section, we first introduce preliminaries regarding CLIP (Radford et al., 2021) and existing finetuning-based methods to achieve *zero-shot adversarial robustness*, and elaborate on our paradigm, termed *Adversarially Finetune Like You Pretrain* (AdvFLYP).

3.1 PRELIMINARIES

CLIP (Radford et al., 2021) is a dual-encoder architecture with a vision encoder $f_\theta(\cdot) \in \mathbb{R}^d$ and a text encoder $g_\phi(\cdot) \in \mathbb{R}^d$, which map an image x and a text t into the same d -dimensional latent space, respectively. In the pretraining phase of CLIP, the vision and text encoders are trained over 400 million web-scale image-text pairs via a contrastive loss (Oord et al., 2018), which maximizes the cosine similarity of an image embedding with its corresponding text embedding. In a single batch $\{(x_i, t_i)\}_{i=1}^N$, the contrastive loss is formulated as follows:

$$\mathcal{L}_{CLIP}(\{(x_i, t_i)\}_{i=1}^N) = -\frac{1}{2N} \sum_{i=1}^N \left[\log \frac{\exp(s_{ii}/\tau)}{\sum_{j=1}^N \exp(s_{ij}/\tau)} + \log \frac{\exp(s_{ii}/\tau)}{\sum_{j=1}^N \exp(s_{ji}/\tau)} \right] \quad (1)$$

where τ is the temperature, and $s_{ij} = \frac{f_\theta(x_i)^\top g_\phi(t_j)}{\|f_\theta(x_i)\| \|g_\phi(t_j)\|}$ is the cosine similarity between x_i and t_j . After the pretraining phase, given an image x_{test} and a set of pre-trained textual categories $\{c_1, \dots, c_K\}$ at inference time, CLIP is able to perform zero-shot classification by classifying it as the category with the highest similarity $\hat{y} = \arg \max_k \frac{f_\theta(x_{test})^\top g_\phi(T[c_k])}{\|f_\theta(x_{test})\| \|g_\phi(T[c_k])\|}$, where $T[\cdot]$ is a textual template, which is usually ‘a photo of a [CLS]’.

Adversarial attacks. A pixel-level perturbation $\delta \in \mathbb{R}^{C \times H \times W}$ bounded by a L_∞ -radius ball, when maliciously designed to maximize the loss of a given image x w.r.t. its label c_{GT} , can cause CLIP to misclassify the sample:

$$\delta_{adv} = \arg \max_{\delta} \mathcal{L}(f_\theta(x + \delta), c_{GT}), \quad s.t. \|\delta\|_\infty \leq \epsilon \quad (2)$$

where \mathcal{L} is cross-entropy loss, and ϵ is the attack budget controlling the attack strength.

Adversarial finetuning of CLIP typically finetunes the pre-trained vision encoder f_θ by generating adversarial images on the fly and aligning them with their correct labels on a proxy dataset. To this end, Mao et al. (2023) propose TeCoA, which is a conventional cross-entropy loss of adversarial images w.r.t. ground-truth labels:

$$\theta' = \arg \min_{\theta} \mathcal{L}(f_\theta(x + \delta_{adv}), c_{GT}) \quad (3)$$

Subsequent finetuning-based methods (Wang et al., 2024; Yu et al., 2024) introduce regularization terms to improve robustness on downstream datasets and generalization on top of this loss. Specifically, Wang et al. (2024) employ a frozen original CLIP model to guide the finetuning process, while Yu et al. (2024) propose text-guided attention and regularize the finetuning and encourage the model to attend to informative areas in adversarial images.

3.2 LIMITATIONS OF EXISTING METHODS

Despite their effectiveness in boosting zero-shot adversarial robustness, these methods incur a significant generalization decline on clean data (Mao et al., 2023). To look into this degradation, we perform introductory experiments by leveraging two proxy datasets, ImageNet (Deng et al., 2009)

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

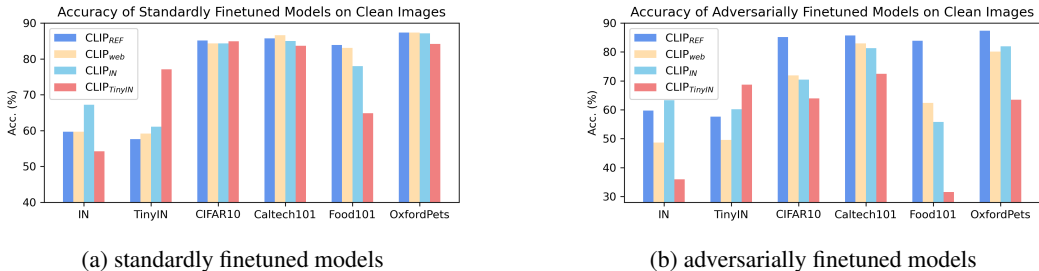


Figure 1: Behaviour of finetuned models. The subscripts in the legends indicate the dataset used for finetuning. CLIP_{web} represents CLIP finetuned on a toy dataset of 100k web-scale image-text pairs. CLIP_{REF} denotes the original pre-trained CLIP without any weight updates.

and TinyImageNet (Le & Yang, 2015), to finetune f_θ of CLIP with either clean images (a.k.a. standard finetuning) or adversarial images (a.k.a. adversarial finetuning). We also collect a toy web-scale dataset of 100k image-text pairs to imitate CLIP’s pre-training data distribution, denoted as *web*. We test these models on the clean images of 6 selected datasets, which include two proxy datasets involved, two general object recognition datasets CIFAR10 (Krizhevsky et al., 2009) and Caltech101 (Fei-Fei et al., 2006), and two fine-grained datasets Food101 (Bossard et al., 2014) and OxfordPets (Parkhi et al., 2012). We report the results in Fig.1 and make the following observations: (1) The model still suffers a noticeable loss in generalization ability even when it is finetuned with clean images on a proxy dataset. (2) Adversarial finetuning of CLIP on a proxy dataset leads to higher accuracy on clean images from this dataset than the original CLIP. Both observations indicate that apart from the inherent robustness-generalization trade-off, the finetuned CLIP overfits to the data distribution of the dataset it has been finetuned on, therefore compromising generalization further.

Furthermore, when finetuning f_θ on a classification dataset via a cross-entropy loss, it equivalently aligns a large number of semantically-rich images from the same category with a single textual prompt. Intuitively, this forces the model to adapt to the classification task instead of retaining its capability of matching images and texts.

3.3 AdvFLYP

To address the limitations discussed above, we propose a simple yet effective adversarial finetuning paradigm, which we term *Adversarially Fine-tune Like You Pretrain* (AdvFLYP)¹, to achieve *zero-shot adversarial robustness*. The idea is intuitive, and can be viewed as resuming the training of CLIP with adversarial images while maximally maintaining the training recipe.

Data preparation. Since the pre-training data of CLIP is not publicly available, to imitate the distribution of CLIP’s pre-training data, we collect 1M webscale image-text pairs. To this end, we randomly sample one million entries with reachable URLs

Algorithm 1 PyTorch-style pseudocode for AdvFLYP

```

# target vision encoder f_theta
# frozen original vision encoder F_theta0
# frozen text encoder g_phi
# collected data: web image-text pairs D
for (X, T) in D: # one batch
    # generate adversarial perturbations
    delta=PGD(f, g, (X, T), l_clip)
    # obtain embeddings
    X=f(X+delta), X_c=f(X), X_ori=F(X+delta), T=g(T)
    # compute probability logit
    P=X@T.t(), P_c=X_c@T.t(), P_ori=X_ori@T.t()
    # logit-level regularization
    l_logit=P*(P/P_c).log()+P*(P/P_ori).log()
    # feature-level regularization
    l_feat=(X-X_c).norm(-1)+(X-X_ori).norm(-1)
    # update theta w.r.t. final loss
    L=(l_clip(f, g, (X_B+delta, T_B))+l_logit+l_feat)
    L.backward()
    optimizer.step()
return theta

```

from LAION-400M (Schuhmann et al., 2021). Following the original work of CLIP (Radford et al., 2021), we utilize these noisy web-scale data without further data cleansing.

¹This paradigm is named after a work on robust finetuning of CLIP, *Finetune Like You Pretrain* (FLYP) (Goyal et al., 2023), which finds that CLIP finetuned with the same contrastive objective as in pretraining compares favourably to CLIP typically finetuned with a cross-entropy loss.

Adversarial finetuning. As in general adversarial training (AT) frameworks (Madry et al., 2018), this process involves a min-max optimization. Instead of employing a cross-entropy loss on adversarial images from a classification dataset (Mao et al., 2023; Wang et al., 2024; Yu et al., 2024), we propose to employ the same contrastive loss as in CLIP’s pretraining (Eq. 1) in our adversarial finetuning paradigm. Specifically, given a batch of image-text data $\{(x_i, t_i)\}_{i=1}^N$, in the inner maximization process, we optimize a L_∞ -bounded perturbation δ_i for each sample x_i in this batch, such that this contrastive loss (Eq. 1) is maximized:

$$\delta = \arg \max_{\{\delta_1, \dots, \delta_N\}} \mathcal{L}_{CLIP}(\{(x_i + \delta_i, t_i)\}_{i=1}^N), \quad s.t. \|\delta\|_\infty \leq \epsilon \quad (4)$$

Note that $\delta \in \mathbb{R}^{N \times C \times H \times W}$ is optimized at the same time, instead of being optimized individually, by employing PGD algorithm (Carlini & Wagner, 2017). This is in stark contrast to existing adversarial finetuning methods (Mao et al., 2023), where a perturbation is optimized independently for each image to maximize its cross-entropy loss against a pre-defined set of categories (Eq. 2). In the experiment section, we will investigate the impact of the contrastive learning setting, *e.g.*, batch size, in the context of adversarial finetuning. In the outer minimization loop, we finetune the model weights θ of the vision encoder to minimize the contrastive loss of this batch of adversarial samples. To further alleviate generalization loss, we also incorporate regularization guided by the frozen original CLIP F_{θ_0} during finetuning. Specifically, we obtain the normalized embeddings of adversarial images, $X_\theta^{adv} = \left[\frac{f_\theta(x_i + \delta_i)}{\|f_\theta(x_i + \delta_i)\|} \right]_{i=1}^N \in \mathbb{R}^{N \times d}$ and $X_{\theta_0}^{adv} = \left[\frac{f_{\theta_0}(x_i + \delta_i)}{\|f_{\theta_0}(x_i + \delta_i)\|} \right]_{i=1}^N \in \mathbb{R}^{N \times d}$, which are output by the target model f_θ and the original model F_θ , respectively. We also feed the clean images to the target model and obtain their embeddings $X_\theta^{clean} = \left[\frac{f_\theta(x_i)}{\|f_\theta(x_i)\|} \right]_{i=1}^N \in \mathbb{R}^{N \times d}$. We compute the probability logits of X_θ^{adv} , $X_{\theta_0}^{adv}$ and X_θ^{clean} w.r.t. the text features $T_\phi = \left[\frac{g_\phi(t_i)}{\|g_\phi(t_i)\|} \right]_{i=1}^N \in \mathbb{R}^{N \times d}$:

$$P_\theta^{adv} = \text{softmax}(X_\theta^{adv} T^\top) \in \mathbb{R}^{N \times N} \quad (5)$$

$$P_{\theta_0}^{adv} = \text{softmax}(X_{\theta_0}^{adv} T^\top) \in \mathbb{R}^{N \times N} \quad (6)$$

$$P_\theta^{clean} = \text{softmax}(X_\theta^{clean} T^\top) \in \mathbb{R}^{N \times N} \quad (7)$$

The logit-level regularization is formulated following Wang et al. (2024):

$$\mathcal{L}_{logit} = \frac{1}{N} [\text{KL}(P_\theta^{adv} \| P_{\theta_0}^{adv}) + \text{KL}(P_\theta^{adv} \| P_\theta^{clean})] \quad (8)$$

where $\text{KL}(\cdot \| \cdot)$ denotes KL divergence. In this work, we additionally introduce feature-level regularization, which we find to benefit generalization on clean images:

$$\mathcal{L}_{feat} = \frac{1}{N} [\|X_\theta^{adv} - X_{\theta_0}^{adv}\|_F + \|X_\theta^{adv} - X_\theta^{clean}\|_F] \quad (9)$$

where $\|\cdot\|_F$ denotes Frobenius norm. To sum up, in the outer minimization loop, the weights of the vision encoder f_θ are updated as follows:

$$\theta' = \arg \min_\theta \{\mathcal{L}_{CLIP}(\{(x_i + \delta_i, t_i)\}_{i=1}^N) + \mathcal{L}_{logit} + \mathcal{L}_{feat}\} \quad (10)$$

We summarize the paradigm of AdvFLYP in Alg. 1.

4 EXPERIMENTS

We conduct extensive experiments to evaluate the adversarial robustness of our proposed paradigm. We implement state-of-the-art finetuning-based methods with available code on two common proxy datasets, TinyImageNet and ImageNet, and compare AdvFLYP with these baselines under various attack scenarios. We also implement AdvFLYP under multiple finetuning settings to understand the behaviour of this contrastive learning framework in an adversarial finetuning (AFT) context.

4.1 IMPLEMENTATION DETAILS

Following previous AFT-based methods, we finetune the pre-trained CLIP’s ViT-B/32 vision encoder. To prepare web-scale image text-pairs that closely follow CLIP’s pre-training data distribution, we collect 1M data pairs. Specifically, we randomly sample one million data points with reachable URLs from LAION-400M (Schuhmann et al., 2021) and denote it as `small-LAION`. Following the data preprocessing of CLIP, we crop and resize the raw images to the size of 224×224 . In the finetuning process, we set the batch size to 256, unless otherwise specified. To generate adversarial images, we employ the PGD algorithm (Carlini & Wagner, 2017) with 2 iterations to update the batch-wise perturbations $\delta \in \mathbb{R}^{N \times C \times H \times W}$ (Eq. 4). The attack strength and step size during finetuning are set to $\epsilon = 1/255$, $\alpha = 1/255$, respectively. We leverage an SGD optimizer and retain the initial learning rate from CLIP’s pre-training stage at $1e - 4$ (Radford et al., 2021), which is dynamically adjusted with cosine scheduling. We finetune the model for 20 epochs on a single NVIDIA RTX A6000 GPU device.

4.2 BASELINES AND DATASETS

We implement TeCoA (Mao et al., 2023), PMG-AFT (Wang et al., 2024) and TGA-ZSR (Yu et al., 2024) based on their released code and finetune f_θ on two proxy datasets, TinyImageNet (Le & Yang, 2015) and ImageNet (Deng et al., 2009), which include 100k and roughly 1.2M training images, respectively. The main training objective of these methods is the cross-entropy loss of adversarial images *w.r.t.* their true labels on a classification dataset. To ensure fair comparison, we use their original hyperparameters in our implementation while keeping other settings such as dataset pre-processing strictly identical. We further implement FARE (Schlarmann et al., 2024), which is an unsupervised adversarial finetuning method that alternately generates adversarial images by enlarging their the L_2 distance to the original embeddings in the latent space, and updates the encoder weights to minimize their distance. When comparing to baselines finetuned on TinyImageNet, we randomly sample a subset of 100k training image-text pairs from `small-LAION` to maintain the same training data amount, which we denote as `tiny-LAION`.

After the finetuning process, we evaluate the *zero-shot adversarial robustness* of all baselines on 14 downstream datasets spanning diverse domains, which include general object recognition datasets CIFAR10 (Krizhevsky et al., 2009), CIFAR100 (Krizhevsky et al., 2009), STL10 (Coates et al., 2011), Caltech101 (Fei-Fei et al., 2006) and Caltech256 (Griffin et al., 2007); fine-grained recognition datasets OxfordPets (Parkhi et al., 2012), Flowers102 (Nilsback & Zisserman, 2008), Food101 (Bossard et al., 2014), StanfordCars (Krause et al., 2013); scene recognition datasets SUN397 (Xiao et al., 2010) and Country211 (Radford et al., 2021); domain-specific datasets FGVCaircraft (Maji et al., 2013), EuroSAT (Helber et al., 2019), DTD (Cimpoi et al., 2014).

4.3 RESULTS AND DISCUSSION

AdvFLYP *v.s.* AFT methods finetuned on TinyIN. Although recent work (Wang et al., 2024; Yu et al., 2024) employs TinyImageNet as a common proxy dataset due to its small size, we argue that it is not an ideal dataset for AFT. In our preliminary experiments (Fig. 1a in Sec. 3.2), we find that when performing standard finetuning on TinyImageNet, it already causes a significant accuracy degradation on downstream tasks. As can be seen in Table 1, in AFT, this degradation is further worsened by introducing adversarial images, with TeCoA and PMG-AFT losing over 20 average points on downstream datasets compared to the original CLIP. This degradation is largely attributed to the fact that TinyImageNet, with very limited semantics in low-resolution (64×64) training images, drastically differs from the distribution of CLIP’s pre-training data, which are noisy web-scale image-text pairs. All AFT baselines exhibit substantially higher accuracy of clean images on TinyImageNet, which indicates that they heavily overfit to the data distribution of TinyImageNet, even when finetuned with adversarial images. Despite weaker zero-shot adversarial robustness compared to supervised AFT methods, FARE retains the best clean accuracy among other baselines, showing that unsupervised AFT better preserves CLIP’s zero-shot capabilities than supervised counterparts. AdvFLYP, finetuned with a contrastive loss on the same amount of images collected from the web, achieves the same level of clean accuracy with unsupervised AFT, while showing higher adversarial robustness than supervised AFT methods under PGD-10 (Carlini & Wagner, 2017) and AutoAttack (Croce & Hein, 2020) at attack strength $\epsilon = 1/255$. When evaluated under PGD-10 at $\epsilon = 4/255$,

AdvFLYP v.s. AFT methods finetuned on ImageNet. As can be seen from Table 2, AFT methods finetuned on ImageNet incur lesser loss of generalization on clean images, compared to when finetuned on TinyImageNet. This is due to the larger dataset size and better image quality of ImageNet, which effectively alleviates overfitting. Nonetheless, all baselines finetuned on ImageNet still exhibit signs of overfitting, with higher test accuracy on ImageNet, especially TGA-ZSR. Additionally, employing ImageNet as a proxy dataset for AFT benefits certain datasets that share more similar classes than others. For example, AFT baselines finetuned on ImageNet, which include a considerable amount of animal classes, transfer better to OxfordPets (Parkhi et al., 2012) and ImageNet-like general classification datasets. Among AFT baselines, the unsupervised method FARE achieves the best clean accuracy. However, FARE exhibits noticeably lower robustness levels, compared to supervised AFT. TGA-ZSR is shown to largely overfit to both the data distribution of ImageNet and the attack type during AFT. Our AdvFLYP performs consistently better than supervised methods in terms of generalization on clean images and different adversarial attack scenarios.

For each downstream dataset, we average the accuracy of a model under all scenarios including PGD-10 ($\epsilon = 1/255$), PGD-10 ($\epsilon = 4/255$), AutoAttack ($\epsilon = 1/255$) and clean images, for a comprehensive evaluation. When finetuned on a tiny portion of web-scale image-text data (tiny-LAION), our AdvFLYP achieves best overall results on 11 out of 14 downstream datasets, with an improvement of 4.86 points (19.1%). Additionally, AdvFLYP finetuned on small-LAION performs best on 12 out of 14 datasets, with an improvement of 1.01 points (3.1%). Results show that our AdvFLYP paradigm steadily enhances zero-shot adversarial robustness of CLIP across datasets of various domains, without overfitting to the distribution of any dataset, proving the importance of following the pre-training data and training objective of CLIP in AFT.

4.4 ANALYSIS ON ADVFLYP

This work proposes an AFT paradigm that practically resumes the pretraining process of CLIP, except that the training data are swapped for adversarial images. However, the behaviour of contrastive learning in a context of adversarial finetuning is understudied. This section explores other training settings of the AdvFLYP paradigm to provide insights into its working.

Ablation studies. On tiny-LAION, we ablate the default AdvFLYP to reveal contributions of each term of the formulated loss (Eq. 10) to its effectiveness. From Table 3, it can be seen that the contrastive loss \mathcal{L}_{CLIP} as employed in CLIP’s pretraining significantly improves its *zero-shot adversarial robustness* from the original CLIP’s 3.04 to 33.06, which plays a major role in AdvFLYP. Adding logit-level regularization \mathcal{L}_{logit} further improves transferability of robustness on downstream data, which is in line with the findings of PMG-AFT (Wang et al., 2024). In this work, we find that feature-level regularization \mathcal{L}_{feat} is highly effectively in retaining generalization on downstream clean images. We introduce logit- and feature-level regularization into our AdvFLYP to reach a sweet spot between robustness and clean accuracy without further tuning their respective weights.

(%)	Rob. Acc.	Clean Acc.	Avg.
\mathcal{L}_{CLIP}	30.41	53.02	41.72
$\mathcal{L}_{CLIP} + \mathcal{L}_{logit}$	33.44	52.64	<u>43.04</u>
$\mathcal{L}_{CLIP} + \mathcal{L}_{feat}$	30.74	55.11	42.93
AdvFLYP	<u>33.06</u>	<u>53.28</u>	43.17

Table 3: Ablation of training objective in AdvFLYP. The reported robust accuracy is tested under PGD-10 ($\epsilon = 1/255$). We report average accuracy over 14 downstream datasets.

Amount of training data. Results in Table 1 and Table 2 show that a larger size of webscale data for AFT benefits both robustness and accuracy. We investigate the impact of training data amount on AdvFLYP in Fig. 2 (left). Interestingly, when the training data amount is scarce, both robustness and accuracy are data-limited and therefore, are not in conflict. Increasing the amount of collected web-scale data for AdvFLYP improves both robustness and accuracy considerably. When there are sufficient image-text pairs, the improvement plateaus and does not noticeably benefit from more training data.

Batch size. Different from existing AFT methods, which finetune f_θ through a cross-entropy *w.r.t.* the true labels of images, AdvFLYP employs the same contrastive loss in AFT as in the pretraining

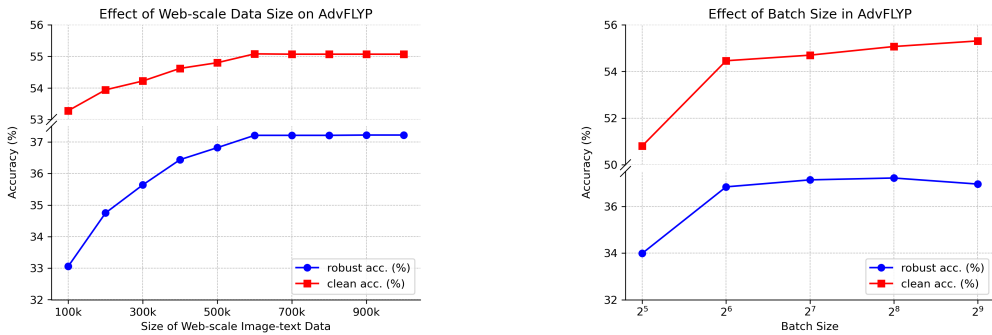


Figure 2: Adversarial robustness (PGD-10, $\epsilon = 1/255$) and clean accuracy of AdvFLYP paradigm under various training settings, averaged on 14 downstream datasets.

phase. As suggested in (Radford et al., 2021), a larger batch size is beneficial for contrastive learning, because it provides more negative samples in a single batch. In this experiment, we vary the batch size from 32 (2^5) to 512 (2^9) for AdvFLYP and report its performance in Fig. 2 (right). We adjust the number of epochs for each batch size to ensure a similar number of updates to model weights. We show that increasing the batch size for AdvFLYP improves clean accuracy significantly, which is in line with the findings of CLIP (Radford et al., 2021). In comparison, enlarging the batch size boosts the robustness at first. However, when the batch size increases to 512, robustness slightly declines, with additional gains of clean accuracy. This shows that both robustness and accuracy benefit from a sufficiently large batch size in contrastive learning in an AFT context. Further enlarging the batch size would trade robustness off for accuracy. We experiment with the maximum batch size of 512 due to hardware constraints.

Unfreeze other components. Existing AFT methods invariably finetune the vision encoder f_θ of CLIP, while keeping the text encoder g_ϕ frozen. Intuitively, resuming the pretraining recipe for AFT would result in whole finetuning of CLIP. In this experiment, we unfreeze more modules of CLIP to investigate the impact on AdvFLYP, and report the results in Table 4. It can be seen that unfreezing more modules of CLIP in AdvFLYP does not lead to better robustness, indicating that f_θ is still the major component in adversarially robust CLIP. Unfreezing g_ϕ and more modules such as layer-wise normalization leads to slightly better clean accuracy. We also find that employing only the image-to-text loss, *i.e.*, the first term of \mathcal{L}_{CLIP} (Eq. 1), leads to best clean accuracy. Utilizing the full contrastive loss \mathcal{L}_{CLIP} and finetuning only f_θ achieves the best overall performance for AdvFLYP.

(%)	Rob. Acc.	Clean Acc.	Avg.
f_θ, g_ϕ	36.04	<u>55.04</u>	45.54
$f_\theta, g_\phi, \text{others}$	35.78	54.95	45.36
$f_\theta (\mathcal{L}_{CLIP} \rightarrow \mathcal{L}_{I2T})$	36.29	55.30	<u>45.80</u>
$f_\theta (\mathcal{L}_{CLIP} \rightarrow \mathcal{L}_{T2I})$	<u>36.32</u>	54.36	45.34
f_θ	36.82	54.80	45.81

Table 4: Impact of trainable CLIP modules on AdvFLYP finetuned on 500k web data.

5 CONCLUSION

In this work, we propose a simple yet paradigm for adversarial finetuning, *Adversarially Finetune Like You Pretrain* (AdvFLYP), which practically resumes the training of CLIP with adversarial images, while retaining the training recipe as much as possible. This paradigm addresses the limitations of existing AFT methods, which invariably finetune the CLIP model on adversarial images on a proxy classification dataset through a cross-entropy loss, compromising the generalization of CLIP. Our AdvFLYP paradigm employs web-scale image-text data that follows the distribution of CLIP’s pretraining data, and finetunes the same contrastive loss as employed during CLIP’s pretraining. We additionally find that logit- and feature-level regularization benefits robustness and clean accuracy, respectively. AdvFLYP outperforms existing AFT methods commonly finetuned on TinyImageNet and ImageNet by 19.1% and 3.1%, respectively, while alleviating overfitting to any proxy dataset distribution. We analyse the behaviour of AdvFLYP and show that a sufficient large training data size and training-time batch size are crucial to both downstream robustness and accuracy, throwing light on the behaviour of contrastive learning in an adversarial finetuning context of CLIP.

REFERENCES

- 486
487
488 Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative compo-
489 nents with random forests. In *European conference on computer vision*, pp. 446–461. Springer,
490 2014.
- 491 Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017*
492 *IEEE symposium on security and privacy (sp)*, pp. 39–57. Ieee, 2017.
- 493
494 Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. De-
495 scribing textures in the wild. In *Proceedings of the IEEE conference on computer vision and*
496 *pattern recognition*, pp. 3606–3613, 2014.
- 497 Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised
498 feature learning. In *Proceedings of the fourteenth international conference on artificial intelli-*
499 *gence and statistics*, pp. 215–223. JMLR Workshop and Conference Proceedings, 2011.
- 500
501 Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble
502 of diverse parameter-free attacks. In *International conference on machine learning*, pp. 2206–
503 2216. PMLR, 2020.
- 504 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-
505 erarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
506 pp. 248–255. Ieee, 2009.
- 507 Junhao Dong, Piotr Koniusz, Yifei Zhang, Hao Zhu, Weiming Liu, Xinghua Qu, and Yew-Soon Ong.
508 Improving zero-shot adversarial robustness in vision-language models by closed-form alignment
509 of adversarial path simplices. In *Forty-second International Conference on Machine Learning*,
510 2025. URL <https://openreview.net/forum?id=WR0ahlhOoy>.
- 511
512 Li Fei-Fei, Robert Fergus, and Pietro Perona. One-shot learning of object categories. *IEEE trans-*
513 *actions on pattern analysis and machine intelligence*, 28(4):594–611, 2006.
- 514 Shizhan Gong, Haoyu LEI, Qi Dou, and Farzan Farnia. Boosting the visual interpretability of CLIP
515 via adversarial fine-tuning. In *The Thirteenth International Conference on Learning Representa-*
516 *tions*, 2025. URL <https://openreview.net/forum?id=khuIvzxPRp>.
- 517
518 Sachin Goyal, Ananya Kumar, Sankalp Garg, Zico Kolter, and Aditi Raghunathan. Finetune like
519 you pretrain: Improved finetuning of zero-shot vision models. In *Proceedings of the IEEE/CVF*
520 *Conference on Computer Vision and Pattern Recognition*, pp. 19338–19347, 2023.
- 521 Gregory Griffin, Alex Holub, Pietro Perona, et al. Caltech-256 object category dataset. Technical
522 report, Technical Report 7694, California Institute of Technology Pasadena, 2007.
- 523
524 Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset
525 and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected*
526 *Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.
- 527 Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan
528 Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning
529 with noisy text supervision. In *International conference on machine learning*, pp. 4904–4916.
530 PMLR, 2021.
- 531 Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained
532 categorization. In *Proceedings of the IEEE international conference on computer vision work-*
533 *shops*, pp. 554–561, 2013.
- 534
535 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
536 2009.
- 537 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-
538 lutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- 539
Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. 2015.

- 540 Lin Li, Haoyan Guan, Jianing Qiu, and Michael Spratling. One prompt word is enough to boost
541 adversarial robustness for pre-trained vision-language models. In *Proceedings of the IEEE/CVF*
542 *Conference on Computer Vision and Pattern Recognition*, pp. 24408–24419, 2024.
- 543 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
544 Towards deep learning models resistant to adversarial attacks. In *International Conference on*
545 *Learning Representations*, 2018.
- 546 Subhansu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained
547 visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- 548 Chengzhi Mao, Scott Geng, Junfeng Yang, Xin Wang, and Carl Vondrick. Understanding zero-
549 shot adversarial robustness for large-scale models. In *The Eleventh International Confer-*
550 *ence on Learning Representations*, 2023. URL <https://openreview.net/forum?id=P4bXCawRi5J>.
- 551 Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number
552 of classes. In *2008 Sixth Indian conference on computer vision, graphics & image processing*, pp.
553 722–729. IEEE, 2008.
- 554 Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predic-
555 tive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- 556 Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012*
557 *IEEE conference on computer vision and pattern recognition*, pp. 3498–3505. IEEE, 2012.
- 558 Sarah Pratt, Ian Covert, Rosanne Liu, and Ali Farhadi. What does a platypus look like? gener-
559 ating customized prompts for zero-shot image classification. In *Proceedings of the IEEE/CVF*
560 *international conference on computer vision*, pp. 15691–15701, 2023.
- 561 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
562 Girish Sastry, Amanda Aspell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
563 models from natural language supervision. In *International conference on machine learning*, pp.
564 8748–8763. PMLR, 2021.
- 565 Leslie Rice, Eric Wong, and Zico Kolter. Overfitting in adversarially robust deep learning. In
566 *International conference on machine learning*, pp. 8093–8104. PMLR, 2020.
- 567 Oindrila Saha, Grant Van Horn, and Subhansu Maji. Improved zero-shot classification by adapting
568 vlms with text descriptions. In *Proceedings of the IEEE/CVF conference on computer vision and*
569 *pattern recognition*, pp. 17542–17552, 2024.
- 570 Fawaz Sammani and Nikos Deligiannis. Interpreting and analysing clip’s zero-shot image classifi-
571 cation via mutual knowledge. *Advances in Neural Information Processing Systems*, 37:39597–
572 39631, 2024.
- 573 Christian Schlarmann, Naman Deep Singh, Francesco Croce, and Matthias Hein. Robust clip: Un-
574 supervised adversarial fine-tuning of vision embeddings for robust large vision-language models.
575 In *International Conference on Machine Learning*, pp. 43685–43704. PMLR, 2024.
- 576 Christoph Schuhmann, Robert Kaczmarczyk, Aran Komatsuzaki, Aarush Katta, Richard Vencu,
577 Romain Beaumont, Jenia Jitsev, Theo Coombes, and Clayton Mullis. Laion-400m: Open dataset
578 of clip-filtered 400 million image-text pairs. In *NeurIPS Workshop Datacentric AI*, number FZJ-
579 2022-00923. Jülich Supercomputing Center, 2021.
- 580 Lijun Sheng, Jian Liang, Zilei Wang, and Ran He. R-tpt: Improving adversarial robustness of vision-
581 language models through test-time prompt tuning. In *Proceedings of the Computer Vision and*
582 *Pattern Recognition Conference*, pp. 29958–29967, 2025.
- 583 Satoshi Suzuki, Shin’ya Yamaguchi, Shoichiro Takeda, Sekitoshi Kanai, Naoki Makishima, Atsushi
584 Ando, and Ryo Masumura. Adversarial finetuning with latent representation constraint to mitigate
585 accuracy-robustness tradeoff. In *2023 IEEE/CVF International Conference on Computer Vision*
586 *(ICCV)*, pp. 4367–4378. IEEE Computer Society, 2023.

- 594 Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfel-
595 low, and Rob Fergus. Intriguing properties of neural networks. In *International Conference on*
596 *Learning Representations*, 2014.
- 597
- 598 Baoshun Tong, Hanjiang Lai, Yan Pan, and Jian Yin. On the zero-shot adversarial robustness of
599 vision-language models: A truly zero-shot and training-free approach. In *Proceedings of the*
600 *Computer Vision and Pattern Recognition Conference*, pp. 19921–19930, 2025.
- 601
- 602 Sibow Wang, Jie Zhang, Zheng Yuan, and Shiguang Shan. Pre-trained model guided fine-tuning for
603 zero-shot adversarial robustness. In *Proceedings of the IEEE/CVF conference on computer vision*
604 *and pattern recognition*, pp. 24502–24511, 2024.
- 605
- 606 Xin Wang, Kai Chen, Jiaming Zhang, Jingjing Chen, and Xingjun Ma. Tapt: Test-time adversarial
607 prompt tuning for robust inference in vision-language models. In *Proceedings of the Computer*
608 *Vision and Pattern Recognition Conference*, pp. 19910–19920, 2025.
- 609
- 610 Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database:
611 Large-scale scene recognition from abbey to zoo. In *2010 IEEE computer society conference on*
612 *computer vision and pattern recognition*, pp. 3485–3492. IEEE, 2010.
- 613
- 614 Songlong Xing, Zhengyu Zhao, and Nicu Sebe. Clip is strong enough to fight back: Test-time
615 counterattacks towards zero-shot adversarial robustness of clip. In *Proceedings of the Computer*
616 *Vision and Pattern Recognition Conference*, pp. 15172–15182, 2025.
- 617
- 618 Lu Yu, Haiyang Zhang, and Changsheng Xu. Text-guided attention is all you need for zero-shot
619 robustness in vision-language models. *Advances in Neural Information Processing Systems*, 37:
620 96424–96448, 2024.
- 621
- 622 Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan.
623 Theoretically principled trade-off between robustness and accuracy. In *International conference*
624 *on machine learning*, pp. 7472–7482. PMLR, 2019.
- 625
- 626 Jiaming Zhang, Xingjun Ma, Xin Wang, Lingyu Qiu, Jiaqi Wang, Yu-Gang Jiang, and Jitao Sang.
627 Adversarial prompt tuning for vision-language models. In *European conference on computer*
628 *vision*, pp. 56–72. Springer, 2024.
- 629
- 630 Mingkun Zhang, Keping Bi, Wei Chen, Jiafeng Guo, and Xueqi Cheng. CLIPure: Purification in
631 latent space via CLIP for adversarially robust zero-shot classification. In *The Thirteenth Interna-*
632 *tional Conference on Learning Representations*, 2025. URL [https://openreview.net/](https://openreview.net/forum?id=TQ2ZOy6miT)
633 [forum?id=TQ2ZOy6miT](https://openreview.net/forum?id=TQ2ZOy6miT).
- 634
- 635 Yunqing Zhao, Tianyu Pang, Chao Du, Xiao Yang, Chongxuan Li, Ngai-Man Man Cheung, and Min
636 Lin. On evaluating adversarial robustness of large vision-language models. *Advances in Neural*
637 *Information Processing Systems*, 36:54111–54138, 2023.
- 638
- 639
- 640
- 641
- 642
- 643
- 644
- 645
- 646
- 647