Test-time Loss Landscape Adaptation for Zero-Shot Generalization in Vision-Language Models

Aodi Li¹, Liansheng Zhuang¹, Xiao Long¹, Minghong Yao¹ and Shafei Wang²

¹University of Science and Technology of China, Hefei 230026, China ²Peng Cheng Laboratory, Shenzhen 518000, China aodili@mail.ustc.edu.cn, lszhuang@ustc.edu.cn

Abstract

Test-time adaptation of pre-trained vision-language models has emerged as a technique for tackling distribution shifts during the test time. Although existing methods, especially those based on Testtime Prompt Tuning (TPT), have shown promising results, their high computational cost associated with parameter optimization presents challenges for scalability and practical application. This paper unveils the unnecessary nature of backpropagation in existing methods from a loss landscape perspective. Building on this insight, this paper proposes a simple yet effective framework called Test-time Loss Landscape Adaptation (TLLA). TLLA leverages the relative position between the training minimum and test loss landscapes to guide the adaptation process, avoiding the update of model parameters at test time. Specifically, it mainly consists of two main stages: In the prompt tuning stage, a Sharpness-Aware Prompt Tuning (SAPT) method is introduced to identify the training flat minimum, setting the foundation for the subsequent test-time adaptation; In the test stage, a Sharpness-based Test Sample Selection (STSS) approach is utilized to ensure the alignment of flat minima within the training loss landscape and each augmented test sample's loss landscape. Extensive experiments on both domain generalization and cross-dataset benchmarks demonstrate that TLLA achieves stateof-the-art performances while significantly reducing computational overhead. Notably, TLLA surpasses TPT by an average of 5.32% and 6.98% on four ImageNet variant datasets when employing ResNet50 and ViT-B/16 image encoders, respectively. The code will be available soon.

1 Introduction

Recent advancements in vision-language (V-L) pretraining, such as CLIP [Radford *et al.*, 2021], have generated new opportunities for developing foundational models in vision tasks [Jia *et al.*, 2021; Yang *et al.*, 2022]. These models, trained on extensive collections of image-text pairs, can learn and represent a diverse range of visual concepts. By

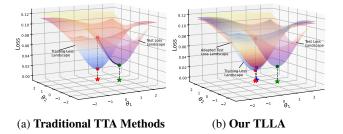


Figure 1: Comparison of conventional TTA methods and our Test-time Loss Landscape Adaptation (TLLA). (a) Traditional TTA methods treat the test landscape as static, aiming to optimize model parameters to achieve the test flat minimum (\star) , using the training minimum (\star) as an initialization. (b) Our TLLA keeps model parameters unchanged during testing. Instead, it adjusts the test landscape to a position where the training minimum (\star) is already very close to the minimum (\star) of the adapted test landscape.

means of well-designed prompts, they can be applied to downstream tasks in a zero-shot manner without requiring task-specific data [Li *et al.*, 2022; Ramesh *et al.*, 2022; Patashnik *et al.*, 2021]. Consequently, various prompt tuning methods [Zhou *et al.*, 2022b; Zhou *et al.*, 2022a] are proposed to directly learn prompts using training data from downstream tasks. Though these methods can find better prompts compared to hand-crafted ones, the learned prompts are limited to the training distribution and may have limited generalization beyond that.

Recently, several studies [Chen et al., 2022; Boudiaf et al., 2022; Wang et al., 2020] have attempted to develop test-time adaptation methods to address this issue at test time. As a typical TTA method for V-L models, Test-time Prompt Tuning (TPT) [Shu et al., 2022] learns adaptive model parameters (i.e., prompts) on the fly using only the given test sample. Along with this paradigm, several variants [Feng et al., 2023; Yoon et al., 2024] have been developed to improve the performance of TPT. Despite the effectiveness of these methods, the computational overhead associated with prompt optimization limits their applicability in many real-world scenarios.

This paper reanalyses the prompt tuning and test-time adaptation of CLIP models from the perspective of loss land-scapes [Li *et al.*, 2018]. From this view, the objective of prompt tuning is to identify the minimum of the training loss landscape, while test-time adaptation aims to seek the mini-

mum in the test loss landscape, using the training minimum as an initialization. Due to constraints on computational resources and inference time, models often struggle to effectively locate the minimum of the test loss landscape, as only a limited number of gradient descent steps are permitted. Furthermore, since training data is unavailable during testing, existing methods fail to account for the similarity between the training and test data. When test samples closely resemble the training data, backpropagations and parameter updates are still necessary for existing methods, leading to a substantial increase in unnecessary computational costs. Recognizing the relative relationship between the well-trained model parameters (i.e., prompts) and the test loss landscape, as illustrated in Figure 1, this paper proposes adjusting the test loss landscape so that its minimum coincides with that of the training loss landscape, avoiding the alteration of prompts.

To achieve this goal, there are several challenges to be addressed. First, altering loss landscapes may require adjustments to the dataset or the model architecture [Bai et al., 2021; Li et al., 2020]. To preserve the model's performance on the original training set, the only feasible solution is to modify test samples. However, employing gradient descent to derive suitable inputs still demands extensive backpropagations, which entail considerable computational complexity and may distort the semantic information of the original test samples. Therefore, efficiently modifying test samples to change the shapes of loss landscapes becomes a critical issue. Furthermore, the ability to quantitatively assess how effectively the training flat minimum can serve as a flat minimum in the test loss landscapes will directly affect the efficacy of test-time adaptation. Consequently, when designing a measurement criterion, it is crucial to consider its potential impact on the classification accuracy of test samples.

Inspired by the above analysis, the paper proposes a novel Test-time Loss Landscape Adaptation (TLLA) framework for efficient test-time adaptation of V-L models. Different from existing methods that typically update prompts for each test sample, TLLA leverages the relative position between the training minimum and test loss landscapes to guide the adaptation process. It mainly consists of two main stages: In the prompt tuning stage, the Sharpness-Aware Prompt Tuning (SAPT) method is introduced to fine-tune the prompts with the downstream training dataset, aiming at seeking the training flat minimum. In the test-time stage, data augmentations are utilized to generate numerous augmented versions for each test sample, each version with its own loss landscape. A Sharpness-based Test Sample Selection (STSS) method is then proposed, using a sharpness-based score to select the augmented versions that align the training flat minimum with those of the augmented loss landscapes. To efficiently modify the loss landscapes, TLLA avoids directly using gradient descent for altering test samples and instead changes the shape of the loss by selecting specific augmented versions of a certain test sample. To design an effective metric for the selection in the test-time stage, TLLA first optimizes for the smallest possible sharpness at the minima while minimizing the loss value in the prompt tuning stage. Since flatter minima generally indicate better model generalization than sharper ones [Keskar et al., 2016; Dziugaite and Roy, 2017; Jiang et al., 2019; Foret et al., 2020], the minimization of sharpness not only improves model generalization but also provides a test-time criterion to measure the alignment of flat minima within the training and test loss landscapes. Theoretical analysis suggests that using the sharpness-based metric will help distinguish the proximity of test samples to the training distribution. The closer an augmented sample is to the training distribution, the smaller its sharpness-based score is likely to be. Since models tend to generate more reliable results for data closer to the training distribution, TLLA significantly improves the zero-shot generalization of vision-language models. Extensive experiments on domain generalization [Hendrycks et al., 2021a] and cross-dataset [Zhou et al., 2022a] benchmarks demonstrate the state-of-the-art performances of TLLA.

Our main contributions can be summarized as follows.

- A novel Test-time Loss Landscape Adaptation (TLLA) framework is proposed for zero-shot generalization in vision-language models. Different from previous works, it leverages the alignment between training and test flat minima to identify reliable predictions to avoid heavy computation on model updates at test time.
- Theoretical analysis is presented to offer a clearer insight into how sample selection at test time improves the reliability of predictions.
- Extensive experiments on domain generalization and cross-dataset benchmarks demonstrate the superior performances of TLLA over other prevalent TTA methods.

2 Related Work

Test-time adaptation (TTA) of vision-language models. Vision-language models, such as CLIP [Radford *et al.*, 2021], trained on large image and text datasets, have shown great promise in semantic representation learning by linking visual and textual representations, enabling exceptional zeroshot reasoning across multiple downstream tasks. To improve CLIP's transfer learning capability for downstream classification tasks, methods combining language prompt learners (e.g., CoOp [Zhou et al., 2022b] and CoCoOp [Zhou et al., 2022a]) with visual adapters (e.g., Tip-Adapter [Zhang et al., 2022]) have been proposed. While effective, these methods typically require large amounts of downstream training data, which is challenging for real-world applications, and they struggle with the distribution misalignment between CLIP's pre-training data and downstream test data. Test-time adaptation (TTA) aims to bridge this gap by adjusting the model at the test time. It generally falls into two categories: the first uses self-supervised proxy tasks like image rotation prediction, modifying the training process (e.g., Test-Time Training [Sun et al., 2020] and TTT++ [Liu et al., 2021]), while the second adapts models without altering the training stage, as seen in approaches like TPT [Shu et al., 2022] and DiffTPT [Feng et al., 2023], which fine-tune prompts for each test sample. Although TPT and DiffTPT are successful in test-time adaptation, test-time prompt tuning is computationally expensive and time-consuming. This paper introduces a novel test-time adaptation method based on the loss landscape, which improves the generalizability of CLIP on downstream tasks while also improving inference efficiency.

Generalization from the view of loss landscape. In recent years, optimization techniques aimed at flat minima in loss landscapes have surged to improve the generalization of deep models [Keskar et al., 2016; Dziugaite and Roy, 2017; Jiang et al., 2019]. One notable method, SAM [Foret et al., 2020], which focuses on finding parameters located in regions of the loss landscape with consistently low loss values, has gained significant attention for its effectiveness and scalability. This concept of flat minima has also been extended to improve the out-of-domain generalization of deep models [Zou et al., 2024]. However, most of these methods concentrate on the training stage and do not investigate the impact of flat minima during the test stage. To the best of our knowledge, this is the first study to reveal the beneficial effects of flatness on test-time adaptation of vision-language models.

3 Methodology

3.1 Preliminaries

Contrastive Language-Image Pre-training. CLIP [Radford et al., 2021] primarily comprises a Text Encoder E_t and an Image Encoder E_v . The Image Encoder is available in two architectures: one based on ResNet [He et al., 2016] and the other using the popular Vision Transformer (ViT) [Dosovitskiy et al., 2020]. This encoder transforms an input image x into its feature representation, i.e., $e_i =$ $E_v(x)$. For a classification task with K classes, the corresponding class labels are formatted into a text template, "a photo of a [cls]", which is then mapped to tokens $y_k =$ $(SOS, t_1, t_2, \dots, t_L, c_k, EOS)$. Here, SOS and EOS represent the embeddings of the start and end tokens, while t_1, t_2, \dots, t_L corresponds to the phrase "a photo of a", and the token c_k denotes the specific description of the k-th class. The text encoder of CLIP, designed as a Transformer architecture, processes these tokens to generate text features: $e_{t,k} = E_t(y_k)$. During the pre-training stage, CLIP is trained on the WIT dataset [Radford et al., 2021] through a contrastive learning approach. In this setup, each image is paired with its corresponding text sentence as a positive sample, while all other image-text combinations are treated as negative samples. The goal of the contrastive learning objective is to enhance the cosine similarity of positive pairs while reducing that of negative pairs. In the zero-shot classification stage, all classes in the dataset are converted to text, and the cosine similarity between image embeddings and text embeddings is computed to determine the probability of an image belonging to each category:

$$p(y_i \mid \boldsymbol{x}) = \frac{\exp(\sin(\boldsymbol{e}_{t,i} \cdot \boldsymbol{e}_i) \tau)}{\sum_{j=1}^{K} \exp(\sin(\boldsymbol{e}_{t,j} \cdot \boldsymbol{e}_i) \tau)},$$
 (1)

where τ is the temperature of the softmax.

Prompt tuning. Prompt tuning has emerged as a popular tuning method for Transformer-based models in downstream tasks. This approach does not modify the model parameters; rather, it changes the input to the model, making it highly efficient. Specifically, instead of using the template "a photo of a [cls]", it replaces the tokens associated with the

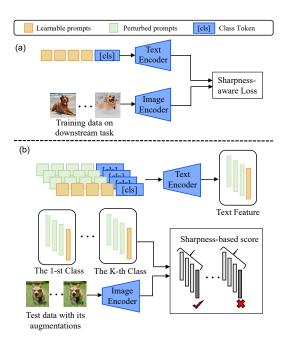


Figure 2: Overview of the proposed Test-time Loss Landscape Adaptation (TLLA). (a) Sharpness-aware Prompt Tuning: It optimizes the model parameters to reduce the sharpness of the loss landscape, enabling stable and effective adaptations during test time without the need for large-scale retraining. (b) Sharpness-based Test Sample Selection: It introduces a selection mechanism to identify augmented test samples that ensure the training flat minimum aligns with those in their loss landscapes, enabling more confident and accurate predictions, particularly for out-of-distribution data.

hand-crafted prompts ("a photo of a") with learnable parameters $p = (p_1, \dots, p_L)$, which are then updated based on the dataset used for downstream tasks.

Test-time prompt tuning. To prevent overfitting that may arise from prompts learned on the downstream training set—which may not perform effectively on test data with distribution shifts—test-time prompt tuning (TPT) [Shu *et al.*, 2022] fine-tunes a specific prompt for each test sample. During testing, multiple augmented views of the test samples are generated. Then, predictions with entropy below a predetermined threshold are kept, while others are discarded using a confidence filter. The averaged entropy of selected predictions is then used as a loss function to update the prompts:

$$\ell_{\text{ENT}}(\boldsymbol{p}) = \arg\min_{\boldsymbol{p}} - \sum_{i=1}^{K} \tilde{p}_{\boldsymbol{p}} (y_i \mid \boldsymbol{x}) \log \tilde{p}_{\boldsymbol{p}} (y_i \mid \boldsymbol{x}), \quad (2)$$

where \tilde{p} denotes the average probability of selected predictions.

3.2 Test-time Loss Landscape Adaptation

This section presents a novel Test-time Loss Landscape Adaptation (TLLA) framework for efficient test-time adaptation of V-L models. Its main idea is to adapt the test landscape such that the well-trained prompts from the downstream training dataset coincide with the flat minimum in the adapted test landscape. The TLLA framework is structured in two main

stages, as depicted in Figure 2. In the first stage (the prompt tuning stage), we utilize the Sharpness-Aware Prompt Tuning (SAPT) method to fine-tune the prompts, positioning them at the flat minimum of the training landscape. The flatness around the minimum not only improves model generalization but also acts as a test-time criterion to align the training minimum with those of test landscapes; In the second stage (the test-time stage), a Sharpness-based Test Sample Selection (STSS) method is proposed, which ensures the alignment of flat minima within the training and selected test loss landscapes through selecting specific augmented versions for each test sample. Next, we'll provide a detailed description of the two-stage TLLA framework.

Sharpness-Aware Prompt Tuning

Prompt tuning is crucial for enabling the effective transfer of the pre-trained CLIP model to downstream tasks. Traditional prompt tuning methods, such as CoOp [Zhou *et al.*, 2022b], typically use cross-entropy loss to fine-tune the prompt *p*:

$$\ell_{\text{CE}}(\boldsymbol{p}) = -\sum_{i=1}^{n} \log p_{\boldsymbol{p}}(y_i|\boldsymbol{x}_i), \tag{3}$$

where $p_p(y_i|x_i)$ represents the predictive probability that x_i belongs to the its true label class.

However, optimizing solely based on training loss can often lead to suboptimal model performance [Foret *et al.*, 2020]. To overcome this limitation, we introduce Sharpness-aware Prompt Tuning (SAPT), which jointly minimizes both the loss and its "sharpness". Sharpness is defined as the sensitivity of the training loss to small perturbations ϵ (with a norm less than ρ) added to the prompts \boldsymbol{p} . This is quantified by evaluating the worst-case increase in loss:

$$\ell_{\mathrm{S}}(\boldsymbol{p}) = \max_{||\boldsymbol{\epsilon}|| \le \rho} \ell_{\mathrm{CE}}(\boldsymbol{p} + \boldsymbol{\epsilon}) - \ell_{\mathrm{CE}}(\boldsymbol{p}). \tag{4}$$

Here, $\ell_{\rm S}$ represents the sharpness measure for the prompts p. Since the perturbation strength ρ is small enough, we can apply a Taylor expansion to approximately solve for the optimal perturbation ϵ^* :

$$\epsilon^{\star} = \rho \frac{\nabla_{\theta} \ell_{\text{CE}}(\boldsymbol{p})}{||\nabla_{\boldsymbol{p}} \ell_{\text{CE}}(\boldsymbol{p})||}.$$
 (5)

Then, the final loss function for SAPT integrates both loss value and sharpness:

$$\ell_{\text{SAPT}}(\boldsymbol{p}) = \ell_{\text{CE}}(\boldsymbol{\theta} + \boldsymbol{\epsilon}^{\star}).$$
 (6)

During training via (stochastic) gradient descent, the contribution from $\nabla_{\boldsymbol{p}} \boldsymbol{\epsilon}^{\star}$ can be disregarded due to the minor perturbation strength ρ . In this way, SAPT not only yields robust prompts that enhance generalization but also provides the sharpness measure as additional information for adaptation during testing.

Sharpness-based Test Sample Selection

Through sharpness-aware prompt tuning, the prompts are positioned at a flat minimum within the training loss landscape. At test time, we aim to maintain the prompts as flat minima in the test loss landscapes for each test sample. This can be accomplished by either fine-tuning the prompts individually

for each sample or adjusting the test loss landscape to work with the current prompts. While fine-tuning can lead prompts away from the training landscape's flat minimum and requires a computationally expensive backward pass, we focus on adjusting the test loss landscape in this paper.

Instead of modifying test samples directly via gradient descent, which is resource-intensive, we propose a Sharpness-based Test Sample Selection (STSS) method. STSS utilizes data augmentations to create multiple test loss landscapes for each sample. By selecting augmented samples that align the current prompts with flat minima in their respective loss landscapes, we ensure current prompts remain optimal. Given that such alignment typically corresponds to small loss values and reduced loss sharpness in these test landscapes, STSS introduces a sharpness-based score as a metric. To mitigate the computational burden of backpropagation in calculating sharpness, we redefine it as follows:

$$\ell'_{S}(\boldsymbol{p}) = \max_{\boldsymbol{\epsilon}_{1},...,\boldsymbol{\epsilon}_{M} \sim \mathcal{N}} \ell_{ENT} \left(\boldsymbol{p} + \rho' \frac{\boldsymbol{\epsilon}_{i}}{\|\boldsymbol{\epsilon}_{i}\|} \right) - \ell_{ENT}(\boldsymbol{p}). \quad (7)$$

Here, $\ell_{\rm ENT}$ denotes entropy, which, according to previous studies [Goyal et al., 2022; Wang et al., 2020], serves as a robust surrogate for the standard cross-entropy loss in the absence of test labels. Note that the entropy loss here applies to individual augmented samples, differing slightly from that in equation (2). The variable ϵ_i refers to random samples drawn from the standard normal distribution \mathcal{N} . Consequently, the maximum variation in the loss resulting from M random perturbations serves as an approximate measure of sharpness. By executing M forward passes, we can efficiently compute sharpness. Since ϵ_i only influences the forward pass through the text encoder and not the image encoder, it requires only M forward passes per test category. This significantly reduces the computational cost compared to traditional test-time adaptation methods. The sharpness is subsequently combined with the loss value to produce the final sharpness-based score, which is employed to select the most reliable augmented test samples. Finally, the prediction for each test sample is determined by a voting mechanism, based on the selected top r augmented versions with the lowest sharpness-based scores.

4 Theoretical Analysis

This section provides a theoretical explanation of how our method improves test-time classification. Let's begin with the following problem: During training, the model learns from data sampled independently and identically from distribution S; During testing, however, data is drawn from two distinct distributions \mathcal{T}_1 and \mathcal{T}_2 . Then, the question is: How can we distinguish between these test distributions and determine on which one the model will perform more reliably?

To address this, we first derive an upper bound for the generalization error, which quantifies the model's performance on unseen data from \mathcal{T}_1 and \mathcal{T}_2 . We will then explore how, when the test distributions are sufficiently distinguishable, TLLA can effectively distinguish between them. This is crucial because, as we will show, when the test distribution closely resembles the training distribution, the generalization error bound decreases, leading to more accurate predictions.

Table 1: **Results on datasets with natural distribution shifts.** We report top-1 accuracy (%) for each method across five target datasets, using two backbones (CLIP-ResNet50 and CLIP-ViT-B/16). We highlight the best results in **bold** and underline the second best results.

Method	ImageNet	ImageNet-A	ImageNet-V2	ImageNet-R	ImageNet-Sk.	Avg.	OOD Avg.
CLIP-ResNet50 [Radford et al., 2021]	59.81	23.24	52.91	60.72	35.48	46.43	43.09
CoOp [Zhou et al., 2022b]	63.33	23.06	55.40	56.60	34.67	46.61	42.43
CoCoOp [Zhou et al., 2022a]	62.81	23.32	55.72	57.74	34.48	46.81	42.82
Tip-Adapter [Zhang et al., 2022]	62.03	23.13	53.97	60.35	35.74	47.04	43.30
TPT [Shu et al., 2022]	60.74	26.67	54.70	59.11	35.09	47.26	43.89
C-TPT [Yoon et al., 2024]	-	25.60	54.80	59.70	35.70	-	44.00
DiffTPT [Feng et al., 2023]	60.80	31.06	55.80	58.80	<u>37.10</u>	48.71	45.69
CoOp [Zhou et al., 2022b]+TPT [Shu et al., 2022]	<u>64.73</u>	30.32	57.83	58.99	35.86	49.55	45.75
CoOp [Zhou et al., 2022b]+ DiffTPT [Feng et al., 2023]	64.70	<u>32.96</u>	61.70	58.20	36.80	50.87	47.42
TLLA (Ours)	65.90	37.87	<u>59.03</u>	62.32	37.63	52.55	49.21
CLIP-ViT-B/16 [Radford et al., 2021]	68.34	49.89	61.88	77.65	48.24	61.20	59.42
CoOp [Zhou et al., 2022b]	71.51	49.71	64.20	75.21	47.99	61.72	59.28
CoCoOp [Zhou et al., 2022a]	71.02	50.63	64.07	76.18	48.75	62.13	59.91
Tip-Adapter [Zhang et al., 2022]	70.75	51.04	63.41	77.76	48.88	62.37	60.27
TPT [Shu et al., 2022]	69.70	53.67	64.30	73.90	46.40	61.59	59.57
C-TPT [Yoon et al., 2024]	-	52.90	63.40	78.00	48.50	-	60.70
DiffTPT [Feng et al., 2023]	70.30	55.68	65.10	75.00	46.80	62.28	60.52
CoOp [Zhou et al., 2022b]+TPT [Shu et al., 2022]	73.30	56.88	66.60	73.80	49.40	64.00	61.67
CoOp [Zhou et al., 2022b]+ DiffTPT [Feng et al., 2023]	75.00	58.09	66.80	73.90	49.50	64.12	61.97
ZERO [Farina et al., 2024]	69.06	61.35	64.13	77.28	48.29	64.02	62.76
PromptAlign [Abdul Samadh et al., 2024]	-	59.37	65.29	<u>79.33</u>	50.23	-	63.55
TLLA (Ours)	74.01	65.90	67.23	81.24	51.81	68.04	66.55

Theorem 1 (Generalization Bound). Consider real-valued function class $\mathcal{F} = \{f_{\theta}(\cdot)\}$, and a bounded loss function $\ell : \mathbb{R} \times \mathbb{R} \to [0, M]$. Define ℓ^{ρ} as:

$$\ell^{\rho}(f_{\theta}(\boldsymbol{x}), y) = \max_{\|\boldsymbol{\epsilon}\|_{2} \le \rho} \ell(f_{\theta+\boldsymbol{\epsilon}}(\boldsymbol{x}), y). \tag{8}$$

Assume that ℓ^{ρ} is μ -Lipschitz with respect to f:

$$|\ell^{\rho}(f,y) - \ell^{\rho}(f',y)| \le \mu |f - f'|.$$
 (9)

Denote the training and test distribution as S and T, respectively. Then, with probability at least $1 - \delta$, the following inequality holds:

$$\ell^{\rho}(f_{\theta}(\boldsymbol{X}_{\mathcal{T}}), Y_{\mathcal{T}}) \leq \frac{M}{2} d_{\mathcal{F}\Delta\mathcal{F}}(\mathcal{S}; \mathcal{T}) + \sqrt{\frac{M^{2}}{2} \log \frac{2}{\delta}} + \hat{\ell}^{\rho}_{\mathcal{S}}(f_{\theta}) + 2\mu R_{n}(\mathcal{F}, \mathcal{S}) + M\sqrt{\frac{\log(2/\delta)}{2n}}.$$
 (10)

Here, (X_T, Y_T) represents the random vector that follows the distribution \mathcal{T} . $R_n(\mathcal{F}, \mathcal{S})$ represents the Rademacher complexity [Zhang, 2023]. The term $d_{\mathcal{F}\Delta\mathcal{F}}(\mathcal{S}; \mathcal{T})$ quantifies the discrepancy between distributions \mathcal{S} and \mathcal{T} .

In the following, we will show that when the two test distributions are sufficiently distinguishable—compared with the tightness of the above upper bound—we can effectively differentiate between them. To proceed with this analysis, we first introduce the concepts of bound tightness and distribution separability.

Definition 2 (β -tightness). Let α be an upper bound for the variable x such that $p(x \leq \alpha) \geq 1 - \delta$. If there exists an oracle upper bound α^* for which $p(x \leq \alpha^*) = 1 - \delta$, we say that the upper bound is β -tight, where $\beta = |\alpha - \alpha^*|$.

Definition 3 (γ -separability). Let \mathcal{T}_1 and \mathcal{T}_2 be two test distributions. We say that they are γ -separable if the condition $|d_{\mathcal{F}\Delta\mathcal{F}}(\mathcal{T}_1;\mathcal{S})-d_{\mathcal{F}\Delta\mathcal{F}}(\mathcal{T}_2;\mathcal{S})| > \gamma$ holds. Here, \mathcal{S} represents the training distribution.

Theorem 4. Consider two test distributions \mathcal{T}_1 and \mathcal{T}_2 that are γ -separable, where $\gamma > \beta$, with β measuring the tightness of the bound established in Theorem 1. Define \mathcal{S} as the training distribution. If the discrepancy measure satisfies $d_{\mathcal{F}\Delta\mathcal{F}}(\mathcal{S};\mathcal{T}_1) < d_{\mathcal{F}\Delta\mathcal{F}}(\mathcal{S};\mathcal{T}_2)$, then there exists a threshold ξ such that:

$$p(\ell^{\rho}(f_{\boldsymbol{\theta}}(\boldsymbol{X}_{\tau_1}), Y_{\tau_1}) < \xi) > p(\ell^{\rho}(f_{\boldsymbol{\theta}}(\boldsymbol{X}_{\tau_2}), Y_{\tau_2}) < \xi), \tag{11}$$

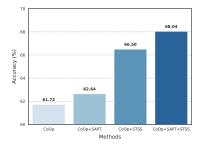
where $(\boldsymbol{X}_{\mathcal{T}_1}, Y_{\mathcal{T}_1})$ and $(\boldsymbol{X}_{\mathcal{T}_2}, Y_{\mathcal{T}_2})$ represent random variables that follow distributions \mathcal{T}_1 and \mathcal{T}_2 , respectively.

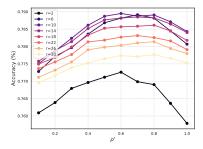
This inequality indicates that a test distribution further from the training distribution tends to exhibit a higher sharpness score. By comparing sharpness scores across test distributions, we can identify which one is closer to the training distribution, thus yielding more reliable predictions. Notably, the tunable parameter ρ controls the tightness of the upper bound, facilitating a precise differentiation between test distributions and improving the model performance. It is important to note that in the theoretical analysis presented in this section, we do not distinguish between ρ and ρ' (which are utilized to calculate the sharpness of the training and test loss landscapes, respectively). However, in practical implementation, we may opt to use different values for ρ and ρ' for better performance. Due to space limitations, detailed proofs and further discussions are provided in the supplementary file.

5 Experiments

5.1 Experimental Setup

Datasets. We conduct two types of experiments to evaluate the model's robustness to natural distribution shifts and its cross-dataset generalization capabilities, following previous research such as TPT [Shu *et al.*, 2022]. To assess the model's robustness to natural distribution shifts, we apply prompt tuning on the ImageNet [Deng *et al.*, 2009] dataset, and evalu-





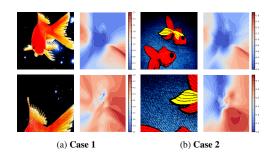


Figure 3: Ablation study on main components of the proposed TLLA.

Figure 4: The influence of the key hyperparameter ρ' on the test accuracy.

Figure 5: 2D Visualization of loss landscapes associated with different augmented test samples.

ate its performance on four ImageNet variants—ImageNet-A [Hendrycks et al., 2021c], ImageNet-V2 [Recht et al., 2019], ImageNet-R [Hendrycks et al., 2021b] and ImageNet-Sketch [Wang et al., 2019]—which is also known as the domain generalization task. In addition, we perform cross-dataset evaluations for image classification across multiple distinct datasets, each from a distinct domain with different classes: including plant and animal species (Flower102 [Nilsback and Zisserman, 2008], OxfordPets [Parkhi et al., 2012]), scenes (SUN397 [Xiao et al., 2010]), textures (DTD [Cimpoi et al., 2014]), transportation (StanfordCars [Krause et al., 2013], Aircraft [Maji et al., 2013]), and general objects (Caltech101 [Fei-Fei et al., 2004]). In this experiment, ImageNet serves as the source dataset, while the remaining fine-grained datasets are used as target datasets for evaluation.

Implementation details. Our experiments are based on pretrained CLIP [Radford et al., 2021] models, specifically CLIP-ResNet50 (using a ResNet50 [He et al., 2016] image encoder) and CLIP-ViT-B/16 (using a Vision Transformer [Dosovitskiy et al., 2020] image encoder). the prompt tuning stage, our experiments are built on the CoOp [Zhou et al., 2022b] framework. We set the number of prompts to 4 and employ SGD as the optimizer, with a learning rate of 0.002. The prompts are trained in a 16-shot manner on the source dataset. For cross-dataset and domain generalization tasks, the prompts were trained for 5 and 50 epochs, with batch sizes of 4 and 32, respectively. During testing, existing TPT-based methods usually leverage the input image along with its 63 augmented views. To ensure a fair comparison, we apply the same data augmentation strategy across all experiments. The key hyperparameters for calculating the sharpness of training and test loss landscapes, ρ and ρ' , are determined through a grid search, with the values ranging from [0.05, 0.1, 0.3, 0.5, 0.7].

5.2 Main Results

Robustness to natural distribution shifts. We first compare the proposed TLLA with state-of-the-art techniques on ImageNet and its variant OOD datasets. The results, presented in Table 1, highlight the superior performance of TLLA across several ImageNet-based OOD datasets. Specifically, TLLA outperforms TPT on both the ResNet-50 and ViT-B/16 architectures, boosting OOD accuracy by 5.32% $(43.89\% \rightarrow 49.21\%)$ and 6.98% $(59.57\% \rightarrow 66.55\%)$, respectively. Moreover, when compared to CoOp+TPT, TLLA shows an average accuracy improvement of 3.46% (45.75%)

 \rightarrow 49.21%) and 4.58% (61.97% \rightarrow 66.55%) for ResNet-50 and ViT-B/16, respectively, on the OOD benchmark. Notably, our method consistently outperforms state-of-the-art (SOTA) methods on both architectures, achieving an average accuracy gain of 1.79% for ResNet-50 and 3.00% for ViT-B/16. These results strongly demonstrate the effectiveness of TLLA in enhancing the out-of-distribution generalization of CLIP across diverse OOD datasets.

Cross-dataset generalization. We also observe superior performance of the TLLA compared to state-of-the-art (SOTA) techniques on cross-domain benchmarks (as shown in Table 3). Specifically, when using ResNet-50 and CLIP-ViT-B/16 as the backbone networks, TLLA achieves an average accuracy improvement of 1.62% (59.37% \rightarrow 60.99%) and 1.44% (66.46% \rightarrow 67.90%) over SOTA methods, respectively. These improvements further validate the effectiveness of the TLLA in adapting to diverse datasets during testing. This capability is particularly valuable for vision-language models like CLIP, as it enables these models to recognize arbitrary categories in image classification tasks without the need for additional training.

5.3 Ablation Study

Main components analysis. To investigate the necessity of each component of the TLLA algorithm, we conduct an ablation study on the domain generalization benchmark dataset, utilizing the CLIP-ViT-B/16 model architecture. The experimental results are presented in Figure 3. Based on the analysis of these results, we can draw several key conclusions: (1) Sharpness-aware prompt tuning enhances gen-The sharpness-aware prompt tuning method eralization. (CoOp+SAPT) significantly improves the model's generalization ability compared to the traditional CoOp method, yielding an average accuracy improvement of 0.92% (from 61.72% to 62.64%) on ImageNet and its out-of-distribution (OOD) datasets. (2) Sharpness-based sample selection during test-time adaptation leads to notable performance gains. During the test-time stage, the CoOp+STSS method, which incorporates sharpness-based sample selection, achieves a substantial performance boost of 4.78% in average accuracy (from 61.72% to 66.50%) over CoOp. This demonstrates the effectiveness of test-time adaptation strategies over sharpness-aware prompt tuning alone. (3) Sharpnessaware prompt tuning facilitates improved sample selection in test-time adaptation. When sharpness-aware prompt tuning (SAPT) is applied during the prompt tuning stage, it

Table 2: Cross-dataset generalization from ImageNet to fine-grained classification datasets. The models are fine-tuned on ImageNet with 16-shot training data per category. We emphasize the best results in **bold** and mark the second-best results with underlining.

Method	Caltech101	Pets	Cars	Flowers102	Aircraft	SUN397	DTD	Avg.
CLIP-ResNet50 [Radford et al., 2021]	87.26	82.97	55.89	62.77	16.11	60.85	40.37	58.03
CoOp [Zhou et al., 2022b] CoCoOp [Zhou et al., 2022a]	86.53 87.38	87.00 88.39	55.32 56.22	61.55 65.57	15.12 14.61	58.15 59.61	37.29 38.53	57.28 58.62
TPT [Shu <i>et al.</i> , 2022] CoOp [Zhou <i>et al.</i> , 2022b]+TPT [Shu <i>et al.</i> , 2022] DiffTPT [Feng <i>et al.</i> , 2023]	87.02 86.90 86.89	84.49 <u>87.54</u> 83.40	58.46 57.65 60.71	62.69 58.83 <u>63.53</u>	17.58 15.84 <u>17.60</u>	61.46 60.00 <u>62.72</u>	40.84 38.06 40.72	58.93 57.83 <u>59.37</u>
TLLA (Ours)	89.74	88.50	<u>59.98</u>	63.78	19.74	63.21	41.97	60.99
CLIP-ViT-B/16 [Radford et al., 2021]	93.35	88.25	65.48	67.44	23.67	62.59	44.27	63.58
CoOp [Zhou et al., 2022b] CoCoOp [Zhou et al., 2022a]	93.70 94.43	89.14 90.14	64.51 65.32	68.71 71.88	18.47 22.94	64.15 67.36	41.92 45.73	62.94 65.40
TPT [Shu et al., 2022] CoOp [Zhou et al., 2022b]+TPT DiffTPT [Feng et al., 2023] PromptAlign [Abdul Samadh et al., 2024]	94.16 93.75 92.49 94.01	87.79 88.93 88.22 90.76	66.87 67.06 67.01 <u>68.50</u>	68.98 68.25 70.10 72.39	24.78 25.89 25.60 24.80	65.50 66.40 65.74 <u>67.54</u>	47.75 47.15 47.00 47.24	65.12 65.35 65.17 66.46
TLLA (Ours)	96.96	91.28	68.93	<u>72.11</u>	26.97	69.29	49.76	67.90

further improves the sample selection during the test stage. Specifically, CoOp+SAPT+STSS (TLLA) improves average accuracy by 5.40% (from 62.64% to 68.04%) compared to CoOp+SAPT, which is a larger performance gain than that without SAPT (+4.78%). This suggests that optimizing for flatter minima during the prompt tuning stage not only improves model generalization but also has a beneficial impact on test-time adaptation performance. Interestingly, the weakened version of our TLLA algorithm, which excludes sharpness-aware prompt tuning (CoOp+STSS), already outperforms the traditional TPT+CoOp algorithm (64.00%) and its variant DiffTPT+CoOp (64.12%) by a noticeable margin, highlighting the efficacy of the proposed test-time loss land-scape adaptation framework even in the absence of sharpness-aware prompt tuning.

Ablative analysis on key hyperparameters. As mentioned earlier, ρ and ρ' are two critical hyperparameters in our algorithm, playing significant roles during the training and testing stages, respectively. Since the impact of ρ on generalization is well-established in prior work [Foret et al., 2020], this paper focuses on how ρ' affects the model's adaptation during testing. Theoretical analysis in Section 4 suggests that ρ' may control the distinguishability between different test distributions. Properly tuning ρ' allows more sensitive discrimination of the test distributions, thereby improving prediction reliability. To empirically validate this, we incrementally adjust ρ' and observe its effect on test accuracy. The curves in Figure 4 show that increasing ρ' initially improves accuracy but eventually causes a decline, similar to ρ 's effect on generalization. When ρ' is too small or too large, the sharpness values become either too small or too large, failing to reflect the true relative sharpness, which reduces accuracy. Notably, when $\rho' = 0$, the sample selection criterion degenerates to entropy, highlighting the positive effect of sharpness in test-time adaptation. Additionally, as the number of retained test samples increases, test accuracy initially improves but eventually decreases. This is because test samples closer to the training distribution tend to have smaller sharpness scores. However, this relationship is probabilistic, not deterministic. Retaining test samples with smaller sharpness scores generally improves accuracy, but keeping too many samples, especially those with larger sharpness scores, harms test accuracy.

5.4 Visualization of Loss Landscapes

To intuitively validate the algorithm's effectiveness, we visualize the test data's loss surface, offering a clearer view of how sample selection enhances prediction reliability. The loss surface shows how the model's loss changes with parameter adjustments. Using a two-dimensional visualization technique [Li et al., 2018], as shown in Figure 5, we present a more intuitive view of the model's loss landscape. When the model's parameters lie at a flat minimum (Figure 5, first row), the corresponding augmented samples accurately reflect the semantic information of the categories. However, when the parameters fall outside this flat minimum (second row), the augmented samples lose semantic integrity, leading to unreliable predictions and reduced generalization to unseen data. This visualization highlights how TLLA can filter out unreliable test samples, mitigating their negative impact on predictions and improving the generalization ability of V-L models.

6 Conclusion

This paper reanalyses the test-time adaptation of V-L models from a loss landscape view, highlighting the redundancy of backpropagation in existing methods. Inspired by this insight, a novel and efficient framework called Test-time Loss Landscape Adaptation (TLLA) is proposed for the test-time adaptation of V-L models. Unlike existing TTA methods that fine-tune prompts for each test sample, TLLA leverages the relative position between the training minimum and test loss landscapes to guide the adaptation process. Specifically, TLLA identifies the training flat minimum during the prompt tuning stage and ensures the alignment of flat minima within the training and test loss landscapes through a test sample selection process. Theoretical analysis and extensive experiments demonstrate the effectiveness and superiority of the proposed TLLA. We hope this work will foster a deeper understanding of loss landscapes and inspire the development of more advanced test-time adaptation methods in the future.

References

- [Abdul Samadh et al., 2024] Jameel Abdul Samadh, Mohammad Hanan Gani, Noor Hussein, Muhammad Uzair Khattak, Muhammad Muzammal Naseer, Fahad Shahbaz Khan, and Salman H Khan. Align your prompts: Test-time prompting with distribution alignment for zero-shot generalization. Advances in Neural Information Processing Systems, 36, 2024.
- [Bai et al., 2021] Haoyue Bai, Fengwei Zhou, Lanqing Hong, Nanyang Ye, S-H Gary Chan, and Zhenguo Li. Nas-ood: Neural architecture search for out-of-distribution generalization. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 8320–8329, 2021.
- [Boudiaf et al., 2022] Malik Boudiaf, Romain Mueller, Ismail Ben Ayed, and Luca Bertinetto. Parameter-free online testtime adaptation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 8344–8353, 2022.
- [Cha et al., 2021] Junbum Cha, Sanghyuk Chun, Kyungjae Lee, Han-Cheol Cho, Seunghyun Park, Yunsung Lee, and Sungrae Park. Swad: Domain generalization by seeking flat minima. Advances in Neural Information Processing Systems, 34:22405– 22418, 2021.
- [Chen et al., 2022] Dian Chen, Dequan Wang, Trevor Darrell, and Sayna Ebrahimi. Contrastive test-time adaptation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 295–305, 2022.
- [Cimpoi et al., 2014] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 3606–3613, 2014.
- [Deng et al., 2009] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pages 248–255. Ieee, 2009.
- [Dosovitskiy et al., 2020] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.
- [Dziugaite and Roy, 2017] Gintare Karolina Dziugaite and Daniel M Roy. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. *arXiv preprint arXiv:1703.11008*, 2017.
- [Farina et al., 2024] Matteo Farina, Gianni Franchi, Giovanni Iacca, Massimiliano Mancini, and Elisa Ricci. Frustratingly easy test-time adaptation of vision-language models. arXiv preprint arXiv:2405.18330, 2024.
- [Fei-Fei et al., 2004] Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In 2004 conference on computer vision and pattern recognition workshop, pages 178–178. IEEE, 2004.
- [Feng et al., 2023] Chun-Mei Feng, Kai Yu, Yong Liu, Salman Khan, and Wangmeng Zuo. Diverse data augmentation with diffusions for effective test-time prompt tuning. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 2704–2714, 2023.

- [Foret et al., 2020] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. arXiv preprint arXiv:2010.01412, 2020.
- [Goyal et al., 2022] Sachin Goyal, Mingjie Sun, Aditi Raghunathan, and J Zico Kolter. Test time adaptation via conjugate pseudo-labels. Advances in Neural Information Processing Systems, 35:6204–6218, 2022.
- [He et al., 2016] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016.
- [Hendrycks et al., 2021a] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In Proceedings of the IEEE/CVF international conference on computer vision, pages 8340–8349, 2021.
- [Hendrycks et al., 2021b] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In Proceedings of the IEEE/CVF international conference on computer vision, pages 8340–8349, 2021.
- [Hendrycks et al., 2021c] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15262–15271, 2021.
- [Jia et al., 2021] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International con*ference on machine learning, pages 4904–4916. PMLR, 2021.
- [Jiang et al., 2019] Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. arXiv preprint arXiv:1912.02178, 2019.
- [Keskar *et al.*, 2016] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.
- [Krause et al., 2013] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In Proceedings of the IEEE international conference on computer vision workshops, pages 554–561, 2013.
- [Li et al., 2018] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. Advances in neural information processing systems, 31, 2018.
- [Li et al., 2020] Yanxi Li, Zhaohui Yang, Yunhe Wang, and Chang Xu. Adapting neural architectures between domains. Advances in Neural Information Processing Systems, 33:789–798, 2020.
- [Li et al., 2022] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image pre-training. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 10965–10975, 2022.
- [Liu et al., 2021] Yuejiang Liu, Parth Kothari, Bastien Van Delft, Baptiste Bellot-Gurlet, Taylor Mordan, and Alexandre Alahi. Ttt++: When does self-supervised test-time training fail or

- thrive? Advances in Neural Information Processing Systems, 34:21808–21820, 2021.
- [Maji et al., 2013] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. arXiv preprint arXiv:1306.5151, 2013.
- [Nilsback and Zisserman, 2008] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In 2008 Sixth Indian conference on computer vision, graphics & image processing, pages 722–729. IEEE, 2008
- [Parkhi et al., 2012] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In 2012 IEEE conference on computer vision and pattern recognition, pages 3498– 3505. IEEE, 2012.
- [Patashnik et al., 2021] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In Proceedings of the IEEE/CVF international conference on computer vision, pages 2085–2094, 2021.
- [Radford et al., 2021] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [Ramesh *et al.*, 2022] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- [Recht et al., 2019] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? In *International conference on machine* learning, pages 5389–5400. PMLR, 2019.
- [Shu et al., 2022] Manli Shu, Weili Nie, De-An Huang, Zhiding Yu, Tom Goldstein, Anima Anandkumar, and Chaowei Xiao. Test-time prompt tuning for zero-shot generalization in visionlanguage models. Advances in Neural Information Processing Systems, 35:14274–14289, 2022.
- [Sun et al., 2020] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training with self-supervision for generalization under distribution shifts. In *International conference on machine learning*, pages 9229–9248. PMLR, 2020.
- [Wang et al., 2019] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. Advances in Neural Information Processing Systems, 32, 2019.
- [Wang et al., 2020] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. arXiv preprint arXiv:2006.10726, 2020.
- [Xiao et al., 2010] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In 2010 IEEE computer society conference on computer vision and pattern recognition, pages 3485–3492. IEEE, 2010.
- [Yang et al., 2022] Jinyu Yang, Jiali Duan, Son Tran, Yi Xu, Sampath Chanda, Liqun Chen, Belinda Zeng, Trishul Chilimbi, and Junzhou Huang. Vision-language pre-training with triple contrastive learning. In *Proceedings of the IEEE/CVF Conference on*

- Computer Vision and Pattern Recognition, pages 15671–15680, 2022
- [Yoon et al., 2024] Hee Suk Yoon, Eunseop Yoon, Joshua Tian Jin Tee, Mark Hasegawa-Johnson, Yingzhen Li, and Chang D Yoo. C-tpt: Calibrated test-time prompt tuning for vision-language models via text feature dispersion. arXiv preprint arXiv:2403.14119, 2024.
- [Zhang et al., 2022] Renrui Zhang, Wei Zhang, Rongyao Fang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. Tip-adapter: Training-free adaption of clip for few-shot classification. In European conference on computer vision, pages 493–510. Springer, 2022.
- [Zhang, 2023] Tong Zhang. *Mathematical analysis of machine learning algorithms*. Cambridge University Press, 2023.
- [Zhao et al., 2018] Han Zhao, Shanghang Zhang, Guanhang Wu, José MF Moura, Joao P Costeira, and Geoffrey J Gordon. Adversarial multiple source domain adaptation. Advances in neural information processing systems, 31, 2018.
- [Zhou et al., 2022a] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for visionlanguage models. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 16816–16825, 2022.
- [Zhou et al., 2022b] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337– 2348, 2022.
- [Zou et al., 2024] Yingtian Zou, Kenji Kawaguchi, Yingnan Liu, Jiashuo Liu, Mong-Li Lee, and Wynne Hsu. Towards robust out-of-distribution generalization bounds via sharpness. arXiv preprint arXiv:2403.06392, 2024.

A Proof to Theorem 3.1

We begin by defining the $\mathcal{F}\Delta\mathcal{F}$ distance, which is crucial for bounding the difference in the losses between two distributions. This definition will play a pivotal role in the subsequent analysis.

Definition 5 ($\mathcal{F}\Delta\mathcal{F}$ distance [Zhao et al., 2018]). Let \mathcal{F} be a hypothesis class for instance space \mathcal{X} , and $\mathcal{A}_{\mathcal{F}}$ be the collection of subsets of \mathcal{X} that are the support of some hypothesis in \mathcal{F} , i.e., $\mathcal{A}_{\mathcal{F}} := \{f^{-1}(\{1\}) \mid f \in \mathcal{F}\}$. The distance between two distributions \mathcal{S} and \mathcal{T} is defined as:

$$d_{\mathcal{F}\Delta\mathcal{F}}\left(\mathcal{S};\mathcal{T}\right) := 2 \sup_{\mathcal{A}(f)\in\mathcal{A}_{\mathcal{F}\Delta\mathcal{F}}} \left| p_{\mathcal{S}}(\mathcal{A}(f)) - p_{\mathcal{T}}(\mathcal{A}(f)) \right|,$$
(12)

where $\mathcal{F}\Delta\mathcal{F}$ is defined as:

$$\mathcal{F}\Delta\mathcal{F} := \left\{ f(x) \oplus f'(x) : f, f' \in \mathcal{F} \right\},\,$$

and \oplus is the XOR operator e.g. $\mathbb{I}(f'(x) \neq f(x))$.

With this definition in hand, we now move to a lemma that connects the distance between distributions to the difference in the losses, setting the stage for bounding the generalization error.

Lemma 1 ([Cha et al., 2021]). Consider an bounded loss function $\ell : \mathbb{R} \times \mathbb{R} \to [0, M]$. Given two distributions, S and T, the difference between the loss with S and the loss with T is bounded by the distance between S and T:

$$|\mathcal{L}_{\mathcal{T}}(f_{\theta}) - \mathcal{L}_{\mathcal{S}}(f_{\theta})| \le \frac{M}{2} d_{\mathcal{F}\Delta\mathcal{F}}(\mathcal{S}; \mathcal{T}).$$
 (13)

Proof.

$$\begin{aligned} &|\mathcal{L}_{\mathcal{T}}\left(f_{\boldsymbol{\theta}}\right) - \mathcal{L}_{\mathcal{S}}\left(f_{\boldsymbol{\theta}}\right)| \\ &= \left|\mathbb{E}_{\boldsymbol{z}' \sim \mathcal{T}}\ell\left(f_{\boldsymbol{\theta}}, \boldsymbol{z}'\right) - \mathbb{E}_{\boldsymbol{z} \sim \mathcal{S}}\ell\left(f_{\boldsymbol{\theta}}, \boldsymbol{z}\right)\right| \\ &\leq \int_{0}^{\infty} \left|p_{\mathcal{T}}\left[\ell\left(f_{\boldsymbol{\theta}}\left(\boldsymbol{x}'\right), y'\right) > t\right] \\ &- p_{\mathcal{S}}\left[\ell\left(f_{\boldsymbol{\theta}}\left(\boldsymbol{x}\right), y\right) > t\right]\right| dt \\ &\leq M \sup_{t \in [0, M]} \sup_{f_{\boldsymbol{\theta}} \in \mathcal{F}} \left|p_{\mathcal{T}}\left(\ell\left(f_{\boldsymbol{\theta}}\left(\boldsymbol{x}'\right), y'\right) > t\right) \right. \\ &- p_{\mathcal{S}}\left(\ell\left(f_{\boldsymbol{\theta}}\left(\boldsymbol{x}\right), y\right) > t\right)\right| \\ &\leq M \sup_{\mathcal{A}(f) \in \mathcal{A}_{\mathcal{F}\Delta\mathcal{F}}} \left|p_{\mathcal{T}}\left(\mathcal{A}(f)\right) - p_{\mathcal{S}}\left(\mathcal{A}(f)\right)\right| \\ &= \frac{M}{2} d_{\mathcal{F}\Delta\mathcal{F}}\left(\mathcal{S}; \mathcal{T}\right). \end{aligned}$$

Next, we incorporate the Rademacher complexity bound for the generalization error. This bound will provide a crucial link between the expected loss and the sample complexity.

Lemma 2 (Theorem 6.31 [Zhang, 2023]). Consider a real-valued function class $\mathcal{F} = \{f_{\theta}(\cdot)\}$ and a bounded loss function $\ell : \mathbb{R} \times \mathbb{R} \to [0, M]$. Assume that the loss function $\ell(f, y)$ is μ -Lipschitz with respect to f:

$$|\ell(f,y) - \ell(f',y)| \le \mu |f - f'|.$$
 (14)

Let $\{x_i, y_i\}_{i=1}^n$ be n i.i.d. samples from S. With probability at least $1 - \delta$:

$$\mathbb{E}_{\mathcal{S}}\ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y) \leq \frac{1}{n} \sum_{i=1}^{n} \ell\left(f_{\boldsymbol{\theta}}\left(\boldsymbol{x}_{i}\right), y_{i}\right) + 2\mu R_{n}(\mathcal{F}, \mathcal{S}) + M\sqrt{\frac{\log(1/\delta)}{2n}}.$$

Here, $R_n(\mathcal{F}, \mathcal{S})$ represents the expected Rademacher complexity:

$$R_{n}(\mathcal{F}, \mathcal{S}) = \mathbb{E}_{(\boldsymbol{x}_{i}, y_{i}) \sim \mathcal{S}} \mathbb{E}_{\sigma} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \ell\left(f_{\boldsymbol{\theta}}(\boldsymbol{x}_{i}), y_{i}\right).$$
(15)

where $\sigma_1, \ldots, \sigma_n$ are independent uniform $\{\pm 1\}$ -valued Bernoulli random variables.

Now, we are ready to prove Theorem 3.1, incorporating the insights from the previous lemmas.

Proof to Theorem 3.1. Define ℓ^{ρ} as:

$$\ell^{\rho}(f_{\theta}(\boldsymbol{x}), y) = \max_{\|\boldsymbol{\epsilon}\|_{2} \le \rho} \ell(f_{\theta+\boldsymbol{\epsilon}}(\boldsymbol{x}), y). \tag{16}$$

Applying Lemma 1 and Lemma 2, we find that with probability at least $1 - \delta/2$,

$$\mathbb{E}_{\mathcal{T}}\ell^{\rho}(f_{\boldsymbol{\theta}}(\boldsymbol{X}_{\mathcal{T}}), Y_{\mathcal{T}}) \leq \frac{1}{n} \sum_{i=1}^{n} \ell^{\rho} \left(f_{\boldsymbol{\theta}} \left(\boldsymbol{x}_{i} \right), y_{i} \right) + \frac{M}{2} d_{\mathcal{F}\Delta\mathcal{F}} \left(\mathcal{S}; \mathcal{T} \right) + 2\gamma R_{n}(\mathcal{F}, \mathcal{S}) + M \sqrt{\frac{\log(2/\delta)}{2n}}.$$
 (17)

Since ℓ^{ρ} is sub-Gaussian for a bounded loss function, we also have that with probability at least $1 - \delta/2$,

$$\ell^{\rho}(f_{\boldsymbol{\theta}}(\boldsymbol{X}_{\mathcal{T}}), Y_{\mathcal{T}}) \leq \mathbb{E}_{\mathcal{T}}\ell^{\rho}(f_{\boldsymbol{\theta}}(\boldsymbol{X}_{\mathcal{T}}), Y_{\mathcal{T}}) + \sqrt{\frac{M^2}{2}\log\frac{2}{\delta}}.$$
(18)

By combining these two inequalities above, we complete the proof of Theorem 3.1.

Remarks. The generalization error bound above does not explicitly include ρ , which could be a significant parameter for mitigating generalization error. However, it is important to note that ρ can enhance the smoothness of ℓ^{ρ} , thereby potentially reducing the Lipschitz constant μ , which is a contributing factor to the expected Rademacher complexity $R_n(\mathcal{F},\mathcal{S})$.

B Proof to Theorem 3.4

Proof to Theorem 3.4. To begin with, we define two quantities, ξ_1 and ξ_2 , as follows:

$$\xi_{1} := \frac{M}{2} d_{\mathcal{F}\Delta\mathcal{F}}\left(\mathcal{S}; \mathcal{T}_{1}\right) + \sqrt{\frac{M^{2}}{2} \log \frac{2}{\delta}} + \frac{1}{n} \sum_{i=1}^{n} \ell^{\rho}\left(f_{\theta}\left(\boldsymbol{x}_{i}\right), y_{i}\right) + 2\mu R_{n}(\mathcal{F}, \mathcal{S}) + M\sqrt{\frac{\log(2/\delta)}{2n}},$$

$$(19)$$

and

$$\xi_{2} := \frac{M}{2} d_{\mathcal{F}\Delta\mathcal{F}}(\mathcal{S}; \mathcal{T}_{2}) + \sqrt{\frac{M^{2}}{2} \log \frac{2}{\delta}} + \frac{1}{n} \sum_{i=1}^{n} \ell^{\rho} \left(f_{\boldsymbol{\theta}}(\boldsymbol{x}_{i}), y_{i} \right) + 2\mu R_{n}(\mathcal{F}, \mathcal{S}) + M \sqrt{\frac{\log(2/\delta)}{2n}}.$$
(20)

Next, we consider the relationship between the distances $d_{\mathcal{F}\Delta\mathcal{F}}(\mathcal{S};\mathcal{T}_1)$ and $d_{\mathcal{F}\Delta\mathcal{F}}(\mathcal{S};\mathcal{T}_2)$. If $d_{\mathcal{F}\Delta\mathcal{F}}(\mathcal{S};\mathcal{T}_1) < d_{\mathcal{F}\Delta\mathcal{F}}(\mathcal{S};\mathcal{T}_2)$, then we immediately conclude that $\xi_1 < \xi_2$. From Theorem 3.1, we can further derive the following two inequalities:

$$p(\ell^{\rho}(f_{\theta}(\boldsymbol{X}_{\mathcal{T}_1}), Y_{\mathcal{T}_1}) > \xi_1) < \delta, \tag{21}$$

and

$$p(\ell^{\rho}(f_{\theta}(\boldsymbol{X}_{\mathcal{T}_2}), Y_{\mathcal{T}_2}) > \xi_2) < \delta. \tag{22}$$

To proceed, let us introduce ξ^* , the oracle upper bound of $\ell^{\rho}(f_{\theta}(X_{\mathcal{T}_2}), Y_{\mathcal{T}_2})$, which satisfies:

$$p(\ell^{\rho}(f_{\theta}(\boldsymbol{X}_{\mathcal{T}_2}), Y_{\mathcal{T}_2}) > \xi^{\star}) = \delta. \tag{23}$$

Next, we define the difference β as the absolute difference between ξ_2 and the oracle bound ξ^\star , i.e., $\beta:=|\xi_2-\xi^\star|$. Given the separability of \mathcal{T}_1 and \mathcal{T}_2 , we can infer that $|\xi_2-\xi_1|>\beta$. This implies that: $\xi_1<\xi_2-\beta\leq\xi^\star$. Consequently, we can conclude that:

$$p(\ell^{\rho}(f_{\theta}(\boldsymbol{X}_{\mathcal{T}_{2}}), Y_{\mathcal{T}_{2}}) > \xi_{1}) > p(\ell^{\rho}(f_{\theta}(\boldsymbol{X}_{\mathcal{T}_{2}}), Y_{\mathcal{T}_{2}}) > \xi^{\star}).$$
(24)

Now, let us set $\xi=\xi_1$, which leads us to the following inequalities:

$$p(\ell^{\rho}(f_{\boldsymbol{\theta}}(\boldsymbol{X}_{\mathcal{T}_{1}}), Y_{\mathcal{T}_{1}}) > \xi) < \delta < p(\ell^{\rho}(f_{\boldsymbol{\theta}}(\boldsymbol{X}_{\mathcal{T}_{2}}), Y_{\mathcal{T}_{2}}) > \xi_{1}), \tag{25}$$

which completes the proof of Theorem 3.2.

To facilitate understanding, we provide several clarifications regarding this theorem as follows.

Remarks 1. In the absence of true labels during test-time adaptation, we must employ entropy as a surrogate loss for cross entropy, which is equivalent to using cross-entropy with conjugate pseudo-labels and is considered the best alternative to cross-entropy[Goyal et al., 2022]. However, using a surrogate loss introduces an additional error term, resulting in a β' -tight bound where $\beta' > \beta$.

Remarks 2. By adjusting the key parameter ρ , we effectively control the tightness of the upper bound. When ρ is tuned appropriately, it allows a more sensitive discrimination of the test distributions, potentially improving the accuracy of identifying which distribution the model performs better on.