056						
057	Method	External Module	Text Template	TTT Iterations	Trainable Parameter	Test Instances
058	ZERO	None	No restriction	0	None	1
059	TDA	Key-value Cache	No restriction	0	None	1
060	CLIPArTT	None	No restriction	10 (default)	Layer Normalization	128 (Multiple required)
061	WATT	None	Diverse set of templates	Multiple required	Layer Normalization	128 (default)
062	ТРТ	None	single template	1	Text Prompt	1
063	LoRA-TTT	None	No restriction	1	LoRA	1

Table 1: Comparison of VLM-focused TTT methods.

### **A. Broader Impact**

064 065 066

067

068

069

070

072

074 075

076

082

Test-Time Training (TTT) for Vision-Language Models (VLMs) is crucial for enhancing their generalization ability and broadening their applicability to real-world AI applications. This study introduces a novel method that achieves strong zero-shot generalization across diverse categories. Our approach enables the development of systems that can adapt to various environments, ranging from memory-constrained edge devices to high-stakes applications, thereby making VLMs more versatile and practical in real-world scenarios. We hope that Parameter-Efficient Fine-Tuning (e.g., LoRA) will play a pivotal role in TTT, inspiring future research aimed at improving the performance of foundation models.

## **B.** Limitations

Our method has limitations that should be addressed in future work. The MEM loss is primarily designed for image 077 classification, making its adaptation to other tasks, such as object detection or segmentation, challenging. In contrast, the 078 MAE loss is task-agnostic, and extending it to such tasks is a promising direction. Additionally, LoRA hyperparameters 079 (e.g., r and  $\gamma$ ) require careful tuning as their optimal values depend on the target domain. Developing a mechanism to dynamically adjust these parameters based on domain characteristics could improve adaptability and performance. 081

#### 083 **C. Related Work**

Test-Time Training (TTT) allows models to adapt to distribution shifts between training and test data during inference 085 through dynamic parameter updates (Liang et al., 2024; Wang et al., 2024; Chen et al., 2022). The challenges in this area lie in designing an effective test-time objective without labels and developing an efficient system suitable for real-world 087 deployment. For example, TENT (Wang et al., 2020) tunes batch normalization statistics at test time using entropy loss; 088 however, this approach requires batch processing rather than instance-level processing, making it challenging to handle 089 sequential data in real-time. In contrast, MEMO (Zhang et al., 2022) computes test loss from a single instance, a strategy we 090 extend to VLMs. Sun et al. (Sun et al., 2020) and Gandelsman et al. (Gandelsman et al., 2022) update the image encoder 091 by introducing auxiliary tasks and applying self-supervision; however, these methods require fine-tuning the model with 092 auxiliary tasks beforehand for TTT. Our approach eliminates this need, allowing for direct adaptation of pre-trained VLMs 093 without additional pre-training steps. We demonstrate that our reconstruction loss enhances performance on foundation 094 models like CLIP, offering a simple yet effective alternative to prior methods. 095

096

097 **TTT for VLMs.** TPT (Shu et al., 2022) focuses on optimizing a text prompt at test time, valued for its simplicity and 098 effectiveness. It demonstrates that augmenting a single test instance and calculating marginal entropy minimization (Zhang 099 et al., 2022) serves as an effective loss for VLMs. DiffTPT (Feng et al., 2023) utilizes stable diffusion to enhance data 100 augmentation quality, while C-TPT (Yoon et al., 2024) is a technique that calibrates TPT to improve reliability. While text prompt tuning remains the predominant approach in TTT for VLMs, some methods instead focus on adapting the image encoder. RLCF (Zhao et al., 2023) tunes the image encoder and demonstrates that CLIP-ViT-B can achieve performance comparable to CLIP-ViT-L but requires CLIP-ViT-L as a feedback source, which poses challenges in memory-constrained 104 environments. As shown in Table 1, WATT (Osowiechi et al., 2024) and CLIPArTT (Hakim et al., 2024) tune the layer 105 normalization parameters of the vision encoder; however, WATT relies on a diverse set of text templates, while CLIPArTT 106 requires multiple test instances, imposing significant constraints on real-world applicability. Moreover, both methods update 107 these parameters across all layers, leading to high computational costs and requiring multiple backpropagation steps. In 108 contrast, our method tunes only the two layers closest to the output, significantly improving computational efficiency and 109

Dataset	# Classes	Test set size
ImageNet (Deng et al., 2009)	1,000	50,000
ImageNet-A (Hendrycks et al., 2021b)	200	7,500
ImageNetV2 (Recht et al., 2019)	1,000	10,000
ImageNet-R (Hendrycks et al., 2021a)	200	30,000
ImageNet-Sketch (Wang et al., 2019)	1,000	50,889
Flowers102 (Nilsback & Zisserman, 2008)	102	2,463
DTD (Cimpoi et al., 2014)	47	1,692
OxfordPets (Parkhi et al., 2012)	37	3,669
StanfordCars (Krause et al., 2013)	196	8,041
UCF101 (Soomro et al., 2012)	101	3,783
Caltech101 (Li et al., 2022)	100	2,465
Food101 (Bossard et al., 2014)	101	30,300
SUN397 (Xiao et al., 2010)	397	19,850
FGVCAircraft (Maji et al., 2013)	100	3,333
EuroSAT (Helber et al., 2019)	10	8,100

Table 2: The details of the datasets used in the experiments.

enabling faster backpropagation. Additionally, lightweight, backpropagation-free methods such as ZERO (Farina et al., 2024) and TDA (Karmanov et al., 2024) have also been proposed. ZERO offers low computational overhead but struggles with generalization performance compared to TPT. While TDA is efficient, it relies on a key-value cache. In contrast, our method adapts to a single test instance in one step, without relying on external modules. This ensures feasibility even in closed, memory-constrained environments such as edge devices, where external resources and cached data are unavailable. 

Application of Low-rank adaptation (LoRA) aims to achieve efficient fine-tuning of large models with vast numbers of parameters in memory-constrained environments by introducing trainable low-rank matrices into each layer of the Transformer architecture, allowing the pre-trained parameters to remain frozen (Hu et al., 2021; Han et al., 2024; Xin et al., 2024). MeLo (Zhu et al., 2024) demonstrates that applying LoRA to vision transformers (ViT) for downstream medical image diagnosis achieves comparable performance to fully fine-tuned ViT models while significantly reducing memory consumption. CLIP-LoRA (Zanella & Ben Ayed, 2024a) demonstrate significant performance improvements in few-shot learning by applying LoRA to the vision encoder of CLIP. However, CLIP-LoRA requires a few labeled samples from the target downstream task.

# **D.** Experiments Details

### **D.1.** Datasets

The evaluation includes out-of-distribution testing on ImageNet and its four variants, as well as fine-grained classification assessments across categories derived from 10 different datasets. The details of the datasets are provided in Table 2.

### **D.2. Detailed Implementation Settings**

Backbone and Optimization. We adopt the pre-trained CLIP-ViT-B/16 as the common backbone architecture. LoRA parameters are optimized in a single step using the AdamW optimizer (Loshchilov & Hutter, 2017) with a learning rate of 0.001 and weight decay of 0.2. All experiments are conducted on a single NVIDIA RTX 3090 GPU with 24GB of memory.

LoRA Configuration. LoRA is applied exclusively to the transformer architecture in layers 11 and 12 of the image encoder with a rank of 16, targeting the key, query, value, and output projection matrices. The scale factor  $\gamma$  is set to 12 for the OOD benchmark and 2 for the fine-grained benchmark. Matrix A is initialized using Kaiming-uniform (He et al., 2015), while matrix **B** is initialized to zero.

Table 3: The 80 hand-crafted text prompts.

"a bad photo of a { class }", "a photo of many { class }", "a sculpture of a { class }", "a photo of the hard to see { class }", "a low resolution photo of the { class }", "a rendering of a { class }", "graffiti of a { class }", "a bad photo of the { class }", "a cropped photo of the { class }", "a tattoo of a { class }", "the embroidered { class }", "a photo of a hard to see { class }", "a bright photo of a { class }", "a photo of a clean { class }", "a photo of a dirty { class }", "a dark photo of the { class }", "a drawing of a { class }", "a photo of my { class }", "the plastic { class }", "a photo of the cool { class }", "a close-up photo of a { class }", "a black and white photo of the { class }", "a painting of the { class }", "a painting of a { class }", "a pixelated photo of the { class }", "a sculpture of the { class }", "a bright photo of the { class }", "a cropped photo of a { class }", "a photo of the dirty { class }", "a jpeg corrupted photo of a { class }", "a class }", "a photo of the dirty { class }", "a photo of a { class }", "a class directly of the dirty { class }", "a photo of a { class }", "a photo of a { class }", "a photo of the dirty { class }", "a photo of a { class }", ", "a blurry photo of the  $\{ class \}$ ", "a photo of the  $\{ class \}$ ", "a good photo of the  $\{ class \}$ ", "a rendering of the a rend class }", "a { class } in a video game", "a photo of one { class }", "a doodle of a { class }", "a close-up photo of the { class }", "a photo of a { class }", "the origami { class }", "the { class } in a video game", "a sketch of a { class }", "a doodle of the { class }", "a origami { class }", "a low resolution photo of a { class }", "the toy { class }", "a rendition of the { class }", "a photo of the clean { class }", "a photo of a large { class }", "a rendition of a { class }", "a photo nice { class }", "a photo of a weird { class }", "a blurry photo of a { class }", "a cartoon { class }", "art of a { class }", "a sketch of the { class }", "a embroidered { class }", "a pixelated photo of a { class }", "itap of the { class }", "a jpeg corrupted photo of the { class }", "a good photo of a { class }", "a plushie { class }", "a photo of the nice { class }", "a photo of the small { class }", "a photo of the weird { class }", "the cartoon { class }", "art of the { class }", "a drawing of the { class }", "a photo of the large { class }", "a black and white photo of a { class }", "the plushie { class }", "a dark photo of a { class }", "itap of a { class }", "graffiti of the { class }", "a toy { class }", "itap of my { class }", "a photo of a cool { class }", "a photo of a small { class }", "a tattoo of the { class }".

**Baselines.** For text prompt tuning baselines, we use 4 trainable text tokens initialized with the hard prompt "a photo of a". We prepare three versions of precomputed prompts: (1) the default hard prompt, (2) an ensemble of 80 hand-crafted prompts (Radford et al., 2021), and (3) CoOp (Zhou et al., 2022b), which uses 4 tokens and is pre-trained on ImageNet with 16-shot supervision. The 80 hand-crafted prompts are listed in Table 3.

**Other Tuning Methods.** For Image Encoder Tuning, we directly tune the key, query, value, and output matrices in layers 11 and 12 of the image encoder using the same optimizer and loss configuration as LoRA-TTT. For Layer Normalization Tuning, we tune only the layer normalization parameters in the same layers with identical settings.

**Test-time Augmentation.** Following TPT (Shu et al., 2022), we expand a single test instance into a batch of 64 using random resized crops (including the original instance). To suppress noise, we select the top 10% of high-confidence samples from the batch for computing the test loss.

#### D.3. MAE loss variants

In this section, we provide details about the variants of the MAE loss. In addition to the MAE loss applied in LoRA-TTT, we explore the following approaches, as illustrated in Figure 1. The loss  $L_{MAE}^{vis, enc}$  calculates the mean squared error of unmasked visual tokens after image encoding. The loss  $L_{MAE}^{cls, dec}$  reconstructs class tokens following the decoding process. The loss  $L_{MAE}^{pix, dec}$  rearranges the visual tokens obtained after decoding back into image pixels, enabling pixel-level reconstruction. This method of calculating  $L_{MAE}^{pix, dec}$  is consistent with traditional TTT approaches based on MAE (Gandelsman et al., 2022; Wang et al., 2023). These methodologies provide diverse perspectives on leveraging MAE loss for effective reconstruction.

### D.4. Evaluation metric

We use the Expected Calibration Error (ECE) (Naeini et al., 2015; Yoon et al., 2024) as a metric to evaluate the calibration performance of the model in image classification. ECE is calculated on a given evaluation dataset by dividing the model's outputs into equally sized bins based on prediction confidence and measuring the discrepancy between the predicted probabilities and the true probabilities within each bin. A well-calibrated model exhibits a smaller gap between predicted



Figure 1: Variants of MAE Loss.

Table 4: **Top1 accuracy of zero-shot image classification on the OOD benchmark** when using the default hard prompt. The results of CoCoOp are obtained from the TPT paper, while others are reproduced with our code. The best results under zero-shot conditions are highlighted in **bold**. Performance improvements over the zero-shot CLIP-ViT-B/16 are indicated with an upward blue arrow (*\triangle blue*) and a downward red arrow (*\triangle red*).

Method	ImageNet	ImageNet-A	ImageNet-V2	ImageNet-R	ImageNet-Sketch	Average	OOD Avg.
CLIP-ViT-B/16	66.71	47.80	60.63	73.99	46.15	59.06	57.14
CoOp (Zhou et al., 2022b)	71.75	50.13	64.51	75.28	47.92	61.92	59.46
CoCoOp (Zhou et al., 2022a)	71.02	50.63	64.07	76.18	48.75	62.13	59.91
TPT (Shu et al., 2022)	69.02	54.73	63.70	77.15	47.99	62.52	60.89
C-TPT (Yoon et al., 2024)	68.50	51.60	62.70	76.00	47.90	61.34	59.55
MTA (Zanella & Ben Ayed, 2024b)	69.23	56.87	63.67	76.88	48.54	63.04	61.49
Image Encoder Tuning	64.26	56.31	59.70	75.89	47.65	60.76	59.89
Layer Normalization Tuning	66.93	48.24	60.94	74.31	46.31	59.35	57.45
LoRA-TTT-M (Ours)	69.21( <sup>12.49</sup> )	60.57 <sub>(†12.77)</sub>	64.28 <sub>(13.65)</sub>	77.53 <sub>(†3.54)</sub>	48.73((2.57)	64.06 <sub>(15.01)</sub>	62.78 <sub>(15.64)</sub>
LoRA-TTT-A (Ours)	66.27 <sub>(10.45)</sub>	52.55(14.75)	60.87( <u>10.24</u> )	75.57 <sub>(1.58)</sub>	47.01( <sup>10.85</sup> )	60.45 <sub>(1.39)</sub>	59.00 <sub>(1.86)</sub>
LoRA-TTT (Ours)	69.40 <sub>(†2.68)</sub>	60.52 <sub>(†12.72)</sub>	64.43 <sub>(†3.80)</sub>	77.84 <sub>(13.85)</sub>	48.94 <sub>(↑2.79)</sub>	64.23 <sub>(15.17)</sub>	62.93 <sub>(15.79)</sub>

confidence and actual accuracy, resulting in a lower ECE value. The ECE is computed as follows:

$$ECE = \sum_{k=1}^{K} \frac{|B_k|}{m} \left| \operatorname{acc}(B_k) - \operatorname{conf}(B_k) \right|, \tag{1}$$

where K represents the number of bins,  $|B_k|$  denotes the number of samples in bin k,  $\operatorname{acc}(B_k)$  is the average accuracy of the samples in bin k, and  $\operatorname{conf}(B_k)$  represents the average prediction confidence of the samples in bin k. In our experiments, the number of bins is set to 20.

#### E. Additional Experiments

#### E.1. Zero-shot classification

Table 4 shows the results on all datasets in the OOD benchmark.

LoRA-TTT: Low-Rank Test-Time Training for Vision-Language Models

Table 5: Error analysis of top-1 accuracy in zero-shot image classification on the OOD benchmark.

Method	ImageNet	ImageNet-A	ImageNet-V2	ImageNet-R	ImageNet-Sketch	Average	OOD Avg.
CLIP-ViT-B/16	66.71	47.80	60.63	73.99	46.15	59.06	57.14
TPT (Shu et al., 2022)	69.02 (±.14)	54.73 (±.11)	63.70 (±.09)	77.15 (±.06)	47.99 (±.04)	62.52 (±.03)	60.89 (±.04)
LoRA-TTT-M (Ours)	69.21 (±.05)	60.57 (±.16)	64.28 (±.08)	77.53 (±.09)	48.73 (±.04)	64.06 (±.02)	62.78 (±.02)
LoRA-TTT-A (Ours)	66.27 (±.11)	52.55 (±.35)	$60.87 (\pm .19)$	75.57 (±.09)	47.01 (±.08)	60.45 (±.06)	59.00 (±.08)
LoRA-TTT (Ours)	$69.40 (\pm .08)$	60.52 (±.19)	64.43 (±.12)	$77.84 (\pm .03)$	$48.94 (\pm .05)$	64.23 (±.01)	62.93 (±.02)

Table 6: Error analysis of top-1 accuracy in zero-shot image classification on the fine-grained benchmark.

Method	Flower102	DTD	Pets	Cars	UCF101	Caltech	Food101	SUN397	Aircraft	EuroSAT	FG Avg.
CLIP-ViT-B/16	67.40	44.39	88.25	65.51	65.24	93.31	83.64	62.56	23.91	42.22	63.64
TPT (Shu et al., 2022)	68.98 (±.18)	45.92 (±.33)	87.27 (±.20)	67.02 (±.14)	68.99 (±.15)	93.55 (±.22)	85.00 (±.06)	65.11 (±.08)	23.76 (±.36)	43.44 (±.08)	64.91 (±.04)
LoRA-TTT-M (Ours)	67.60 (±.33)	46.04 (±.25)	87.11 (±.19)	67.81 (±.16)	68.38 (±.07)	93.59 (±.13)	84.83 (±.13)	64.61 (±.09)	25.68 (±.12)	39.27 (±.23)	64.49 (±.08)
LoRA-TTT-A (Ours)	68.33 (±.02)	45.21 (±.07)	88.72 (±.13)	66.94 (±.02)	66.35 (±.34)	93.71 (±.02)	84.39 (±.05)	63.63 (±.12)	25.38 (±.20)	44.52 (±.15)	64.72 (±.04)
LoRA-TTT (Ours)	67.88 (±.22)	45.86 (±.12)	87.63 (±.06)	67.72 (±.03)	68.38 (±.12)	93.83 (±.16)	84.99 (±.05)	$64.59 (\pm .11)$	25.92 (±.39)	43.23 (±.33)	$65.00 (\pm .06)$

#### E.2. Error analysis

Table 5 and Table 6 present the standard deviation of three random runs with different seeds for zero-shot image classification on the OOD and fine-grained benchmarks, respectively. The randomness of LoRA-TTT-M mainly stems from random data augmentation and one-step optimization, similar to TPT. Additionally, LoRA-TTT-A introduces an additional source of randomness through its masking strategy. Nevertheless, our method achieves an error magnitude comparable to that of TPT.

#### E.3. Expected Calibration Error

Table 7 presents the calibration results on the OOD benchmark, while Table 8 shows the results on the fine-grained benchmark for each dataset. The comparison includes our method, TPT (Shu et al., 2022), and C-TPT (Yoon et al., 2024). The results show that LoRA-TTT-A (*i.e.*, MAE loss) achieves calibration performance comparable to or surpassing that of C-TPT across a wide range of categories, highlighting the effective calibration properties of MAE loss.

#### E.4. Hyper-parameter tuning and sensitivity

Figure 2a shows that adding the MAE loss improves performance on fine-grained datasets without degrading performance on OOD datasets. We chose  $\lambda_1 = 1$  and  $\lambda_2 = 16$  for their consistent strong results across datasets. In Figures 2b to 2d, increasing the number of data augmentations tends to enhance performance; however, for efficiency, we chose 64, aligning with TPT. As  $N_p$  increases, performance tends to decrease, which is consistent with TPT's results. Therefore, following TPT, we select the top 10% ( $N_p = 6$ ). Additionally, AugMix proves to be effective for data augmentation.

We chose different LoRA scales for the two benchmarks because only ImageNet-A exhibited distinct behavior depending on  $\gamma$ , as shown in Table 9. The relationship between this dataset and LoRA parameters requires further investigation. Our method consistently achieves strong performance and outperforms TPT on both benchmarks with r = 16 and  $\gamma = 2$ , while also using the same masking ratio and LoRA layer settings. Hyperparameter sensitivity is inherent in TTT methods. For example, the optimal iteration count in WATT (Osowiechi et al., 2024) and CLIPArTT (Hakim et al., 2024) varies by domain. Our approach generalizes well across multiple benchmarks, highlighting its stability with the fixed configuration.

#### E.5. Scalability Analysis of Our Method

In this section, we evaluate the scalability of our proposed method by applying it to a larger baseline model. Table 10 and Table 11 show the results obtained using the pretrained CLIP-ViT-L/14 on the OOD and fine-grained benchmarks, respectively. LoRA is applied exclusively to the transformer architecture in layers 23 and 24 of the image encoder, targeting the key, query, value, and output matrices with a rank of 16, and the LoRA scale  $\gamma$  is set to 2. All other experimental parameters are consistent with those in the main paper. The results demonstrate that LoRA-TTT consistently outperforms the baseline CLIP-ViT-L/14 across both benchmarks and multiple categories while maintaining the zero-shot setting. It also demonstrates performance improvements when combined with the ensemble text prompts, exhibiting generalization properties to text prompts similar to those observed with CLIP-ViT-B/16. Performance improvements are observed with Table 7: **Expected Calibration Error** (**ECE** $\downarrow$ ) of zero-shot image classification with TTT on the OOD benchmark. The best results, except for the baseline, are highlighted in **bold**.

3	3	2
3	3	3

Method	ImageNet	ImageNet-A	ImageNet-V2	ImageNet-R	ImageNet-Sketch	Average	OOD Avg.
CLIP-ViT-B/16	1.93	8.37	2.51	3.53	4.79	4.23	4.80
TPT (Shu et al., 2022) C-TPT (Yoon et al., 2024)	10.61 3.11	15.35 <b>6.40</b>	11.85 4.64	4.97 2.80	16.14 7.69	11.78 <b>4.93</b>	12.08 5.38
LoRA-TTT-M (Ours)	20.32	25.47	22.66	12.65	30.25	22.27	22.76
LoRA-TTT-A (Ours)	2.97	9.60	4.13	1.35	7.28	5.07	5.59
LoRA-TTT (Ours)	14.04	19.27	16.19	8.08	22.45	16.00	16.49

Table 8: Expected Calibration Error (ECE↓) of zero-shot image classification with TTT on the fine-grained benchmark.
 The best results, except for the baseline, are highlighted in **bold**.

Method	Flower102	DTD	Pets	Cars	UCF101	Caltech	Food101	SUN397	Aircraft	EuroSAT	Average
CLIP-ViT-B/16	3.21	8.23	4.41	4.45	2.93	5.12	2.03	2.24	5.50	7.16	4.53
TPT (Shu et al., 2022)	13.57	23.45	6.18	5.92	11.65	3.60	4.49	11.94	17.81	18.48	11.71
C-TPT (Yoon et al., 2024)	5.24	13.77	1.56	1.56	2.30	3.27	3.31	5.02	4.41	12.47	5.29
LoRA-TTT-M (Ours)	24.27	34.64	10.97	16.96	18.91	4.61	11.96	20.80	25.45	28.70	19.73
LoRA-TTT-A (Ours)	4.10	12.27	3.08	2.20	3.52	4.09	1.83	3.01	6.51	7.34	4.80
LoRA-TTT (Ours)	19.54	26.05	6.68	7.73	11.30	2.31	7.94	13.15	16.76	16.02	12.75

both types of loss (*i.e.*, LoRA-TTT-M and LoRA-TTT-A), highlighting the robustness of our method and its scalability to larger baseline models.

# F. Ablation Study

### F.1. How to apply LoRA for TTT

In this section, we explore the utilization of LoRA for TTT. We investigate the key factors for effectively applying LoRA, including: (1) determining the optimal layers and the extent of LoRA application within the transformer model, (2) understanding the relationship between the appropriate rank and scale, and (3) selecting the attention matrices for tuning.

Which layers should we apply LoRA to? Table 12 presents the zero-shot classification performance when LoRA is applied to specific layers of the image encoder in CLIP-ViT-B/16. Our results indicate that applying LoRA to deeper layers is more effective than to shallower ones, aligning with trends observed in fine-tuning language models (Zhang et al., 2023). Additionally, applying LoRA to more layers does not necessarily improve performance. Limiting its application to the 11th and 12th layers not only outperforms applying it across all layers in terms of performance but also reduces memory consumption and runtime, making our approach more efficient for TTT.

LoRA rank and scale. As shown in Figure 3a, increasing the rank does not directly lead to performance gains. Each
rank has an optimal scale, and as the rank increases, the corresponding optimal scale tends to decrease. When the rank is
small (*e.g.*, rank 4), performance remains stable across different scales, reducing the need for extensive hyperparameter
tuning.

LoRA rank and attention matrices. We investigate the optimal application of LoRA to different attention matrices in CLIP-ViT-B/16. In Figure 3b, we observe that applying LoRA to  $W_v$  at the same rank achieves the best results among the 4 matrices ( $W_o$ ,  $W_v$ ,  $W_q$ , and  $W_k$ ). This trend aligns with previous research (Zhang et al., 2023; Zanella & Ben Ayed, 2024a), even in the context of TTT. Given the same total number of parameters, applying LoRA to  $W_{kvqo}$  shows little difference in performance compared to applying it to  $W_{vq}$  or  $W_{kq}$ .



Table 9: Top-1 accuracy of zero-shot image classification.

Figure 2: Hyper-parameter tuning. Results may slightly vary due to trial randomness, even with the same parameters.

### F.2. Masking strategy

In masked image modeling, the mask strategy plays a crucial role (Hondru et al., 2024; Gao et al., 2024). We examine the effects of the masking ratio, the confidence selection cutoff, the use of an image decoder, and the impact of reconstruction targets. We use a randomly initialized transformer-based decoder with 8 layers, 16 heads, and a 768 embedding size, without prior fine-tuning to ensure a fair evaluation. This decoder allows us to incorporate the pixel-wise reconstruction loss proposed in TTT methods based on MAE (Gandelsman et al., 2022; Wang et al., 2023).

As shown in Table 13, while the masking ratio does not significantly affect the overall performance, we choose a default masking ratio of 50% as it strikes a good balance between performance and computational efficiency. As proposed in TPT, selecting and masking the top 10% of augmented images with the lowest entropy yields better performance than masking all 64 images (*i.e.*, applying a cutoff of 1), with an improvement of over 1% observed in the OOD average. The 10% cutoff not only improves performance but also enhances the computational efficiency of TTT by calculating the loss on only one-tenth of the images. Furthermore, reconstructing the class token is more effective than reconstructing masked visual tokens or image pixels using the decoder. This supports the hypothesis that improving zero-shot image classification performance in VLMs relies more on aligning high-level semantics than on capturing fine-grained features. 

### F.3. Initialization of LoRA weights

LoRA demonstrates high effectiveness and efficiency for TTT, even when initialized with random weights. In this section, we explore the performance gains achieved by fine-tuning the LoRA weights before TTT. We prepare a third dataset, CC3M (Sharma et al., 2018), for LoRA initialization and train only the LoRA weights using the same contrastive loss as in CLIP pre-training (Radford et al., 2021) with image-text pairs. We employ Adam with a learning rate of 1e-6 and a weight decay of 0.05 for optimization, performing one epoch of training with a batch size of 64. 

Table 10: Top-1 accuracy of zero-shot image classification on the OOD benchmark with the CLIP-ViT-L/14 baseline.

Performance improvements over the zero-shot CLIP-ViT-L/14 are indicated with an upward blue arrow (*totlet*) and a downward red arrow (*tred*).

43								
4	Method	ImageNet	ImageNet-A	ImageNet-V2	ImageNet-R	ImageNet-Sketch	Average	OOD Avg.
45	CLIP-ViT-L/14	73.45	68.77	67.75	85.41	57.82	70.64	69.94
46	CLIP-ViT-L/14 + Ensemble	75.53	70.75	69.70	87.85	59.60	72.69	71.97
7	LoRA-TTT-M (Ours)	75.20 <sub>(†1.75)</sub>	73.73 <sub>(†4.96)</sub>	69.74 <sub>(†1.99)</sub>	87.69 <sub>(†2.28)</sub>	59.76 <sub>(†1.94)</sub>	73.22 <sub>(†2.58)</sub>	72.73 <sub>(†2.79)</sub>
F/	LoRA-TTT-A (Ours)	73.88( <u>10.43</u> )	69.91 <sub>(1.13)</sub>	67.98( <sup>10.23</sup> )	85.99 <sub>(10.58)</sub>	58.30( <u></u> 10.49)	71.21(10.57)	70.55( <u>10.61</u> )
3	LoRA-TTT (Ours)	75.07 <sub>(1.61)</sub>	72.56(13.79)	69.24 <sub>(1.49)</sub>	87.27 <sub>(1.86)</sub>	59.48 <sub>(1.67)</sub>	72.72(12.08)	72.14(12.20)
	LoRA-TTT + Ensemble (Ours)	77.03(13.57)	75.63( <u></u> <b>†6.85</b> )	71.86 <sub>(↑4.11)</sub>	89.73 <sub>(14.32)</sub>	61.41 <sub>(†3.59)</sub>	75.13( <u>14.49</u> )	74.66(14.72)

Table 11: **Top-1 accuracy of zero-shot image classification on the fine-grained benchmark with the CLIP-ViT-L/14 baseline**. Performance improvements over the zero-shot CLIP-ViT-L/14 are indicated with an upward blue arrow ( $\uparrow$ blue) and a downward red arrow ( $\downarrow$ red).

454												
455	Method	Flower102	DTD	Pets	Cars	UCF101	Caltech	Food101	SUN397	Aircraft	EuroSAT	FG Avg.
456	CLIP-ViT-L/14 CLIP-ViT-L/14 + Ensemble	76.21 75.92	52.42 54.73	93.05 93.05	76.91 77.78	73.72 75.89	95.17 95.62	88.58 89.20	67.68 70.15	30.03 31.86	55.09 51.70	70.89 71.59
457	LoRA-TTT-M (Ours)	76.45 <sub>(10.24)</sub>	54.14 <sub>(†1.71)</sub>	93.81 <sub>(10.76)</sub>	78.34( <sup>1.43</sup> )	75.23((1.51)	95.05( <u>10.12</u> )	89.32( <u>10.74</u> )	68.97 <sub>(†1.29)</sub>	33.30(13.27)	52.32(12.77)	71.69( <sup>10.81</sup> )
458	LoRA-TTT-A (Ours)	76.65( <sup>10.45</sup> )	52.72(10.30)	93.43(10.38)	77.42( <sup>+0.51</sup> )	74.17( <sup>10.45</sup> )	95.13 <sub>(10.04)</sub>	88.90 <sub>(10.32)</sub>	67.81( <u>10.13</u> )	30.42(10.39)	55.01(10.07)	71.17(10.28)
150	LoRA-TTT (Ours)	76.57( <sup>10.37</sup> )	54.14( <sup>1.71</sup> )	93.87( <sup>10.82</sup> )	78.31( <sup>1.41</sup> )	74.83( <sup>1.11</sup> )	95.54 <sub>(10.37)</sub>	89.34(10.77)	68.72( <sup>1.04</sup> )	33.12(13.09)	53.74(1.35)	71.82( <u></u> <u>0.93</u> )
459	LoRA-TTT + Ensemble (Ours)	75.92 <sub>(↓0.28)</sub>	55.08 <sub>(†2.66)</sub>	93.08 <sub>(↑0.03)</sub>	79.38( <sup>12.47</sup> )	76.79 <sub>(†3.07)</sub>	95.94 <sub>(↑0.77)</sub>	89.79 <sub>(†1.21)</sub>	71.13(13.45)	35.34(15.32)	52.19 <sub>(12.90)</sub>	72.46(11.58)
460												

As shown in Figure 4, LoRA initialization using 21k randomly sampled image-text pairs from CC3M (*i.e.*, only 1% of the total CC3M dataset) improves performance by more than 1% on the fine-grained benchmark and by 0.6% on the OOD benchmark. Furthermore, TTT consistently improves performance on both the benchmarks, regardless of the LoRA initialization. Our experiments demonstrate that fine-tuning LoRA with a small amount of data shows the potential to enhance its performance. While adhering to the constraints of not leveraging domain-specific information or a teacher model, LoRA fine-tuning delivers significant performance improvements in TTT, establishing it as an effective approach for future applications of LoRA in TTT.

# G. Qualitative Analysis

Table 14 shows the t-SNE visualization of image features after the image encoder for various evaluation datasets, comparing the baseline CLIP-ViT-B/16 and our method. The results show that our approach achieves better class separation than the baseline, indicating improved classification performance on the test data. Additionally, the visualizations highlight that the type of test loss affects how class separation is achieved.

# References

- Bossard, L., Guillaumin, M., and Van Gool, L. Food-101–mining discriminative components with random forests. In *Computer vision–ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part VI 13*, pp. 446–461. Springer, 2014.
- Chen, D., Wang, D., Darrell, T., and Ebrahimi, S. Contrastive test-time adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 295–305, 2022.
- Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., and Vedaldi, A. Describing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3606–3613, 2014.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pp. 248–255. Ieee, 2009.
- Farina, M., Franchi, G., Iacca, G., Mancini, M., and Ricci, E. Frustratingly easy test-time adaptation of vision-language
  models. *arXiv preprint arXiv:2405.18330*, 2024.
- Feng, C.-M., Yu, K., Liu, Y., Khan, S., and Zuo, W. Diverse data augmentation with diffusions for effective test-time prompt tuning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2704–2714, 2023.

LoRA Layer	ImageNet	OOD Average	FG Average
12	69.59	62.65	64.68
11-12	69.40	62.93	65.00
9-12	68.97	62.86	64.73
5-8	66.88	61.34	64.83
1-4	67.99	60.69	64.56
All	68.12	62.54	64.62

Table 12: Layers for LoRA application.



Figure 3: **Impact of LoRA application design.** The average top-1 accuracy on the fine-grained benchmark is shown, with LoRA applied to layers 11 and 12 of the image encoder.

- Gandelsman, Y., Sun, Y., Chen, X., and Efros, A. Test-time training with masked autoencoders. *Advances in Neural Information Processing Systems*, 35:29374–29385, 2022.
- Gao, P., Lin, Z., Zhang, R., Fang, R., Li, H., Li, H., and Qiao, Y. Mimic before reconstruct: Enhancing masked autoencoders with feature mimicking. *International Journal of Computer Vision*, 132(5):1546–1556, 2024.
- Hakim, G. A. V., Osowiechi, D., Noori, M., Cheraghalikhani, M., Bahri, A., Yazdanpanah, M., Ayed, I. B., and Desrosiers,C. Clipartt: Adaptation of clip to new domains at test time. *arXiv preprint arXiv:2405.00754*, 2024.
- Han, Z., Gao, C., Liu, J., Zhang, S. Q., et al. Parameter-efficient fine-tuning for large models: A comprehensive survey. *arXiv preprint arXiv:2403.14608*, 2024.
- He, K., Zhang, X., Ren, S., and Sun, J. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pp. 1026–1034, 2015.
- Helber, P., Bischke, B., Dengel, A., and Borth, D. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7): 2217–2226, 2019.
- Hendrycks, D., Basart, S., Mu, N., Kadavath, S., Wang, F., Dorundo, E., Desai, R., Zhu, T., Parajuli, S., Guo, M., 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*, pp. 8340–8349, 2021a.
- Hendrycks, D., Zhao, K., Basart, S., Steinhardt, J., and Song, D. Natural adversarial examples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15262–15271, 2021b.

Reconstruction	Mask Ratio	Cutoff	Decoder	ImageNet	OOD Average	FG Average
Class token	0.25	0.1		67.65 <sub>(10.94)</sub>	<b>58.94</b> ( <sup>1.80</sup> )	64.57 <sub>(↑0.92)</sub>
	0.5	0.1		67.78 <sub>(1.07)</sub>	58.85 <sub>(1.71)</sub>	64.72 <sub>(1.08)</sub>
	0.75	0.1		67.48 <sub>(↑0.76)</sub>	57.87 <sub>(10.72)</sub>	64.35 <sub>(10.71)</sub>
	0.5	0.5		67.52 <sub>(↑0.80)</sub>	58.30 <sub>(11.16)</sub>	64.49 <sub>(10.84)</sub>
	0.5	1		67.20 <sub>(↑0.48)</sub>	57.79 <sub>(10.65)</sub>	64.34 <sub>(10.69)</sub>
	0.5	0.1	1	67.27 <sub>(↑0.55)</sub>	58.28 <sub>(1.14)</sub>	64.14 <sub>(↑0.50)</sub>
Visual tokens	0.5	0.1		66.89 <sub>(↑0.17)</sub>	57.46 <sub>(↑0.32)</sub>	63.79 <sub>(↑0.15)</sub>
Image pixel	0.5	0.1	1	66.67(10.05)	57.05(10.10)	$63.50_{(10,14)}$

Table 13: **Masking strategy.** The LoRA scale  $\gamma$  is set to 2 for both benchmarks. Performance differences from zero-shot CLIP-ViT-B/16 are shown with a blue ( $\uparrow$ ) or red ( $\downarrow$ ) arrow.



Figure 4: Impact of LoRA weight initialization by data size and comparison with TTT.

- Hondru, V., Croitoru, F. A., Minaee, S., Ionescu, R. T., and Sebe, N. Masked image modeling: A survey. *arXiv preprint arXiv:2408.06687*, 2024.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Karmanov, A., Guan, D., Lu, S., El Saddik, A., and Xing, E. Efficient test-time adaptation of vision-language models. In
   *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14162–14171, 2024.
- Krause, J., Stark, M., Deng, J., and Fei-Fei, L. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pp. 554–561, 2013.
- Li, F.-F., Andreeto, M., Ranzato, M., and Perona, P. Caltech 101. CaltechDATA: Pasadena, CA, USA, 2022.
- Liang, J., He, R., and Tan, T. A comprehensive survey on test-time adaptation under distribution shifts. *International Journal of Computer Vision*, pp. 1–34, 2024.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Maji, S., Rahtu, E., Kannala, J., Blaschko, M., and Vedaldi, A. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- Naeini, M. P., Cooper, G., and Hauskrecht, M. Obtaining well calibrated probabilities using bayesian binning. In *Proceedings* of the AAAI conference on artificial intelligence, volume 29, 2015.
- Nilsback, M.-E. and Zisserman, A. Automated flower classification over a large number of classes. In 2008 Sixth Indian conference on computer vision, graphics & image processing, pp. 722–729. IEEE, 2008.
- Osowiechi, D., Noori, M., Hakim, G. A. V., Yazdanpanah, M., Bahri, A., Cheraghalikhani, M., Dastani, S., Beizaee, F.,
   Ayed, I. B., and Desrosiers, C. Watt: Weight average test-time adaptation of clip. *arXiv preprint arXiv:2406.13875*, 2024.

- Parkhi, O. M., Vedaldi, A., Zisserman, A., and Jawahar, C. Cats and dogs. In 2012 IEEE conference on computer vision and pattern recognition, pp. 3498–3505. IEEE, 2012.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Recht, B., Roelofs, R., Schmidt, L., and Shankar, V. Do imagenet classifiers generalize to imagenet? In *International conference on machine learning*, pp. 5389–5400. PMLR, 2019.
- Sharma, P., Ding, N., Goodman, S., and Soricut, R. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pp. 2556–2565, 2018.
- Shu, M., Nie, W., Huang, D.-A., Yu, Z., Goldstein, T., Anandkumar, A., and Xiao, C. Test-time prompt tuning for zero-shot generalization in vision-language models. *Advances in Neural Information Processing Systems*, 35:14274–14289, 2022.
- Soomro, K., Zamir, A. R., and Shah, M. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv* preprint arXiv:1212.0402, 2012.
- Sun, Y., Wang, X., Liu, Z., Miller, J., Efros, A., and Hardt, M. Test-time training with self-supervision for generalization under distribution shifts. In *International conference on machine learning*, pp. 9229–9248. PMLR, 2020.
- Wang, D., Shelhamer, E., Liu, S., Olshausen, B., and Darrell, T. Tent: Fully test-time adaptation by entropy minimization. *arXiv preprint arXiv:2006.10726*, 2020.
- Wang, H., Ge, S., Lipton, Z., and Xing, E. P. Learning robust global representations by penalizing local predictive power. *Advances in Neural Information Processing Systems*, 32, 2019.
- Wang, R., Sun, Y., Gandelsman, Y., Chen, X., Efros, A. A., and Wang, X. Test-time training on video streams. *arXiv* preprint arXiv:2307.05014, 2023.
- Wang, Z., Luo, Y., Zheng, L., Chen, Z., Wang, S., and Huang, Z. In search of lost online test-time adaptation: A survey. *International Journal of Computer Vision*, pp. 1–34, 2024.
- Xiao, J., Hays, J., Ehinger, K. A., Oliva, A., and Torralba, A. Sun database: Large-scale scene recognition from abbey to zoo. In 2010 IEEE computer society conference on computer vision and pattern recognition, pp. 3485–3492. IEEE, 2010.
- Xin, Y., Luo, S., Zhou, H., Du, J., Liu, X., Fan, Y., Li, Q., and Du, Y. Parameter-efficient fine-tuning for pre-trained vision models: A survey. *arXiv preprint arXiv:2402.02242*, 2024.
- Yoon, H. S., Yoon, E., Tee, J. T. J., Hasegawa-Johnson, M., Li, Y., and Yoo, C. D. C-tpt: Calibrated test-time prompt tuning for vision-language models via text feature dispersion. *arXiv preprint arXiv:2403.14119*, 2024.
- Zanella, M. and Ben Ayed, I. Low-rank few-shot adaptation of vision-language models. In *Proceedings of the IEEE/CVF* Conference on Computer Vision and Pattern Recognition, pp. 1593–1603, 2024a.
- Zanella, M. and Ben Ayed, I. On the test-time zero-shot generalization of vision-language models: Do we really need prompt learning? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23783–23793, 2024b.
- Zhang, M., Levine, S., and Finn, C. Memo: Test time robustness via adaptation and augmentation. *Advances in neural information processing systems*, 35:38629–38642, 2022.
- Zhang, Q., Chen, M., Bukharin, A., Karampatziakis, N., He, P., Cheng, Y., Chen, W., and Zhao, T. Adalora: Adaptive budget allocation for parameter-efficient fine-tuning. *arXiv preprint arXiv:2303.10512*, 2023.
- Zhao, S., Wang, X., Zhu, L., and Yang, Y. Test-time adaptation with clip reward for zero-shot generalization in visionlanguage models. arXiv preprint arXiv:2305.18010, 2023.

- Zhou, K., Yang, J., Loy, C. C., and Liu, Z. Conditional prompt learning for vision-language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16816–16825, 2022a.
- Zhou, K., Yang, J., Loy, C. C., and Liu, Z. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022b.
- Zhu, Y., Shen, Z., Zhao, Z., Wang, S., Wang, X., Zhao, X., Shen, D., and Wang, Q. Melo: Low-rank adaptation is better than fine-tuning for medical image diagnosis. In 2024 IEEE International Symposium on Biomedical Imaging (ISBI), pp. 1–5. IEEE, 2024.

