

CRAFT: Cross-modal Aligned Features Improve Robustness of Prompt Tuning

Jingchen Sun* Rohan Sharma Vishnu Suresh Lokhande* Changyou Chen
 University at Buffalo, State University of New York, USA

Abstract

*Prompt Tuning has emerged as a prominent research paradigm for adapting vision-language models to various downstream tasks. However, recent research indicates that prompt tuning methods often lead to overfitting due to limited training samples. In this paper, we propose a **Cross-modal Aligned Feature Tuning (CRAFT)** method to address this issue. Cross-modal alignment is conducted by first selecting anchors from the alternative domain and deriving relative representations of the embeddings for the selected anchors. Optimizing for a feature alignment loss over anchor-aligned text and image modalities creates a more unified text-image common space. Overfitting in prompt tuning also deteriorates model performance on out-of-distribution samples. To further improve the prompt model’s robustness, we propose minimizing Maximum Mean Discrepancy (MMD) over the anchor-aligned feature spaces to mitigate domain shift. The experiment on four different prompt tuning structures consistently shows the improvement of our method, with increases of up to 6.1% in the Base-to-Novel generalization task, 5.8% in the group robustness task, and 2.7% in the out-of-distribution tasks. The code is available at <https://github.com/Jingchensun/Craft>.*

1. Introduction

Large-scale Vision Language (VL) models, such as CLIP [42], ALIGN [22] and LLaVA [29], demonstrate remarkable capabilities in recognition and representation. To leverage the powerful zero-shot recognition capabilities of vision-language models while avoiding the significant computational resource consumption of fine-tuning, various parameter-efficient fine-tuning methods have been proposed. Among these, prompt tuning methods [2, 15, 23, 24, 30, 55, 70, 71] have gained considerable attention in the research community due to their simple structure and relatively intuitive interpretability.

While prompt tuning methods can effectively adapt the vision-language (VL) models to various downstream tasks, a significant challenge remains unresolved: *Existing prompt*

tuning method often lead to overfitting [8, 33, 34, 40, 72]. Existing methods optimize prompts using text-based cross-entropy loss. Although this loss function is commonly used in many image recognition tasks, however, the cross-entropy loss alone make the model more prone to overfit the training data, especially when the available training data is limited in many downstream tasks. While the prompt model is trained on only a few samples, the soft prompts concentrate on task-specific knowledge, easily over-fitting to the target task. E.g, In Table 2, When CoOp [71] or MaPLe [55] is fine-tuned on the Base class, the recognition accuracy on the Novel class is even lower than that of the pre-trained CLIP model.

We attribute this overfitting to the lack of regularization in the latent space. To address these challenges, we propose a cross-modal feature alignment method that leverages relative representations [21, 37, 49] to construct the static anchor and stochastic anchors in the dual model branch. While the stochastic anchors are dynamic samples that are stochastic and selected during every training batch, all the image stochastic samples and the text stochastic samples share a common semantic space. We show that combining the static anchors and the stochastic anchors can provide effective regularization for the learnable prompt and improve the generalization ability of the prompt model.

Additionally, prompt overfitting is exacerbated by domain discrepancies in downstream tasks [20, 58]. As existing methods primarily focus on acquiring task-specific knowledge from the source domains [24, 55, 70], they often neglect the domain distribution discrepancies from the source to target domains. To address this issue, we proposed an image-text joint distribution mean discrepancy minimization to enhance the cross-domain consistency of optimized prompts. Our empirical results show that the modified MMD method efficiently mitigates overfitting and improves the model’s out-of-distribution recognition capability.

We conducted experiments on four different prompt tuning structures: Visual Prompt Tuning (VPT) [2], Language Prompt Tuning (LPT) [71], and Multi-modal Prompt Tuning with MaPLe [55] and PromptSRC [24]. Our method consistently improved the average classification accuracy across three tasks. In the Base-to-Novel generalization experiments, applying our method to MaPLe improved accuracy by up to

*Corresponding authors: Jingchen Sun and Vishnu Suresh Lokhande

6.1 points. Additionally, applying our method to PromptSRC reduced the domain gap by 5.4 points in the Waterbird group robustness task. Furthermore, when applied to ImageNet and its four variant datasets, accuracy improved by up to 2.7 points. These experiments demonstrate the effectiveness of our proposed method in enhancing the generalization and robustness of prompt tuning.

Our contributions are summarized as follows: (1) We propose a novel relative anchor-based cross-modal feature alignment method, which effectively regularizes prompt tuning and mitigates overfitting issues. (2) We design a modified Maximum Mean Discrepancy (MMD) loss based on our aligned feature space, enhancing the model’s ability to handle out-of-distribution data. (3) We validate our approach on four different prompt tuning structures across three image classification tasks. The experimental results demonstrate the effectiveness and robustness of our method.

2. Related Work

Prompt Tuning. CoOp [71] pioneered the use of learnable prompts in the language branch to adapt vision-language models. Building on this concept, Visual Prompt Tuning (VPT) [23] employs customized visual prompts tailored to image embeddings to enhance the distinguishability of image features. Unified Prompt Tuning (UPT) [64] learns a unified prompt for both image and text branches, aligning them in latent space. MaPLE [55] extends multimodal prompt tuning to deeper network structures, embedding learnable prompts into each layer of the transformer. Despite various prompt structure designs proposed for adaptation, the cross-entropy loss used in these methods will make the model prone to causing overfitting.

Out-of-Distribution Adaption. Many works show that the overfitting is harmful for the out-of-distribution generalization task [10, 13, 35, 47, 48, 50, 61, 62, 73]. To solve this problem, PromptSRC [24] proposes a Gaussian weighted sampling of prompts learned at different epochs, DePT [65] decouples base-specific knowledge during tuning to preserve task-shared knowledge. Our method differs from these approaches by assuming that a feature space that is more aligned across modalities and domains can more effectively alleviate the problems of overfitting and out-of-distribution.

3. Method

Notation. In the prompt-tuning of the vision-language models, we consider images denoted by x and text by y . Transpose of the image vector is denoted as x^T . Images are drawn from a probability distribution \mathbb{P}_x defined over measurable space $(\Omega_x, \mathcal{F}_x)$, while the text is drawn from a probability distribution \mathbb{P}_y defined on a measurable space $(\Omega_y, \mathcal{F}_y)$. Each image-text pair is associated with a label c , which belongs to one of the K classes. In our approach, we

utilize a CLIP [42] model, which provides both an image encoder f_θ and a text encoder g_ϕ . These encoders effectively align the feature spaces of images and text [37].

Incorporating Prompt Tuning Parameters into $\{\theta, \phi\}$.

In prompt-tuning methods, certain learnable prompts are prepended to the text prompts onto the text branch, denoted by $\Lambda_{\text{text}} \in \{\lambda_{\text{text}}^1, \lambda_{\text{text}}^2, \dots, \lambda_{\text{text}}^L\}$. Similarly, learnable prompts $\Lambda_{\text{visual}} \in \{\lambda_{\text{visual}}^1, \lambda_{\text{visual}}^2, \dots, \lambda_{\text{visual}}^L\}$, are prepended into the image branch. For more details regarding the process of prepend, refer to MaPLE [55] architecture or Figure 1 in the paper. In an image classification task, the tokenized image embeddings and learnable visual prompts are fed into the image encoder [55], resulting in image features $f_\theta(x, \Lambda_{\text{visual}})$. Similarly, for the text encoder, tokenized class embeddings and learnable text prompts are processed in the text encoder, yielding text embeddings $g_\phi(y, \Lambda_{\text{text}})$. To maintain simplicity in this paper, the notation $\{\theta, \phi\}$ will be used to represent the parameters that include $\{\Lambda_{\text{visual}}, \Lambda_{\text{text}}\}$ respectively.

Our Feature Alignment Method. Existing methods mainly use text-based cross-entropy loss, which leads to overfitting, as elaborated in Section 1. We propose an anchor-based Cross-modal Feature Alignment method. Our methodology involves generating a distribution over the query point x (or y) by applying a softmax function to the distances from certain anchors in the alternate modality. Cross-modal alignment is possible by comparing samples from the alternative domain (say image) for every sample in the given domain (say text). A distribution of anchors \mathbb{A}_X from \mathbb{P}_X and another discrete distribution \mathbb{A}_Y from \mathbb{P}_Y can be drawn to define the anchors. Comparisons made on these anchors help in reducing computational complexity; instead of comparing them to all, we only compare them to a few. Having defined anchors, let us write down the distribution of classes for every query image-text pair (x, y) . For $a_x \sim \mathbb{A}_X$ and $a_y \sim \mathbb{A}_Y$ and using $\langle \cdot, \cdot \rangle$ as the inner product measure, we have the discrete probability density function over the image/text domains defined as

$$p_x(c = k | x; \theta) = \frac{\exp\langle f_\theta(x), a_y^k \rangle}{\sum_{k'} \exp\langle f_\theta(x), a_y^{k'} \rangle} \quad (1)$$

$$p_y(c = k | y; \phi) = \frac{\exp\langle g_\phi(y), a_x^k \rangle}{\sum_{k'} \exp\langle g_\phi(y), a_x^{k'} \rangle} \quad (2)$$

Minimizing the negative log-probability of Equation 1 and Equation 2, we obtain the feature alignment loss as follows:

$$\mathcal{L}_{\text{Aligned}}(\theta, \phi) = -\log p_x(c = k | x; \theta)p_y(c = k | y; \phi) \quad (3)$$

Minimizing negative log probability often involves transitioning to expectations, while the objective is ideally expressed

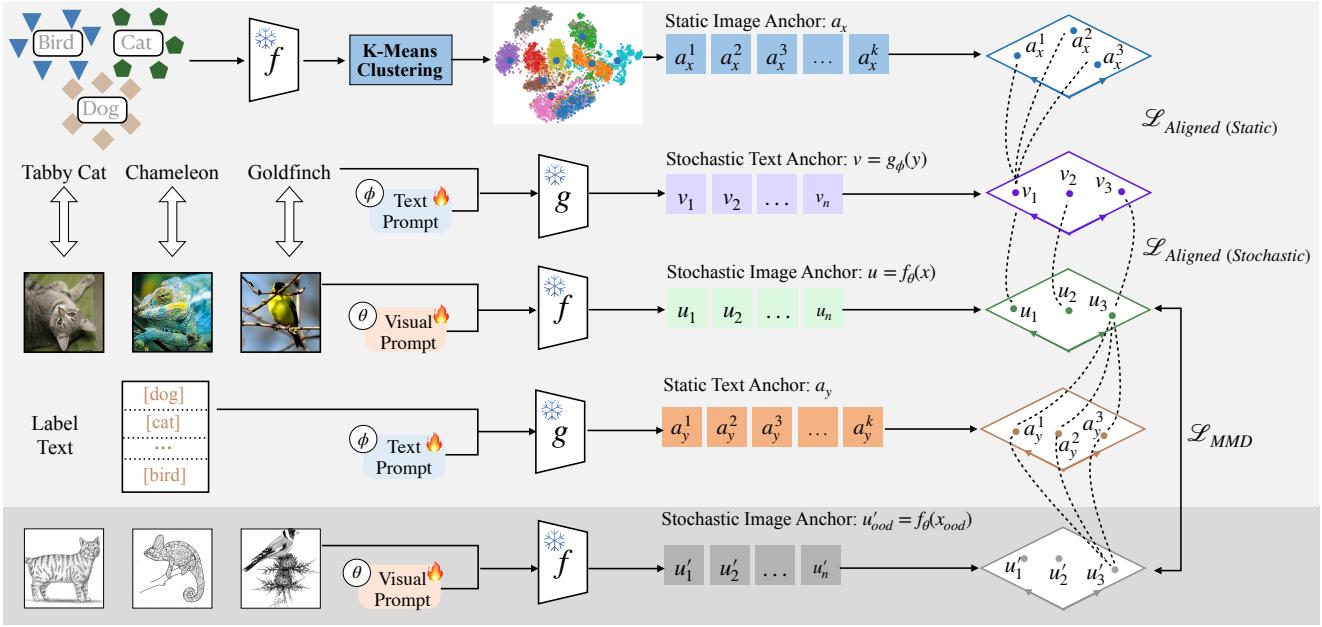


Figure 1. Illustration of Our Proposed Cross-Modal Feature Alignment Method. Firstly k-means clustering is conducted image embeddings to obtain static image anchors. Static text anchors are derived from the class-text labels. Simultaneously, we construct batch-level image and text samples to create stochastic image and text anchors. Static image anchors are aligned with stochastic text anchors using Equation 2. Additionally, stochastic text anchors and stochastic image anchors are aligned with each other by Equation 3. To address out-of-distribution samples, we apply the Maximum Mean Discrepancy (MMD) method to the aligned features, ensuring consistency within the latent space.

as an expectation over the true data distribution, in practice, we approximate this using the empirical distribution of observed data. Training proceeds by minimizing $\mathcal{L}_{\text{Aligned}}$ of the true class k through SGD [45]. Such a cross-modal loss function promotes the alignment of feature spaces across both the text and image modalities [39]. Additionally, leveraging anchors from the alternative modality can offer more reliable classification cues [51].

3.1. Anchors selection for feature alignment

In the previous section, we have already two types of anchors - the imaging anchors, denoted as a_x , and the text anchors, denoted as a_y . Based on how we sample a given image or text anchor, we introduce two additional categories.

- 1. Static Anchors.** Anchors are chosen before the training begins. The selected anchors remain the same throughout the training process. For the text modality, these anchors might represent a static template, such as "a photo of {}" where the class-name is inserted within {}. For the image modality, static anchors can be determined using clustering algorithms.
- 2. Stochastic Anchors.** As the name implies, these anchors are chosen randomly in every iteration of training. In a given iteration, a fixed count of anchors are sampled from both the text and image modality.

Using static anchors in (3) yields $\mathcal{L}_{\text{Aligned}}(\text{static})$. We obtain $\mathcal{L}_{\text{Aligned}}(\text{stochastic})$ using stochastic anchors in (3). Overall,

$$\mathcal{L}_{\text{Aligned}} = \mathcal{L}_{\text{Aligned}}(\text{static}) + \mathcal{L}_{\text{Aligned}}(\text{stochastic}) \quad (4)$$

The interplay between static anchors and stochastic anchors is visible in Figure 1.

Necessity of Static and Stochastic Anchors. Static anchors provide reliable reference points throughout training, ensuring consistency during feature comparisons. Stochastic anchors serve as a form of data augmentation. In each iteration, samples from various modalities are presented to the model, allowing for sampling from the complete support of the data distribution. This approach enhances robustness against minor perturbations and outliers, resulting in more stable and generalized performance. [4].

Remark 1. Pre-trained vision-language (V-L) dual model architectures, like MaPLe [55], employ text-based cross-entropy loss. In text cross-entropy loss, a random batch of images are drawn and are compared to static text-labels. This is simply optimizing for $\mathcal{L}_{\text{Aligned}}(\theta) \propto \langle f_\theta(x), a_y \rangle$, where a_y represent the static text anchors defined by text-labels. Thus, our proposed $\mathcal{L}_{\text{Aligned}}(\theta, \phi)$ from (3) subsumes text-cross entropy loss.

3.2. Accounting for the variations in distribution

Although prompt-tuning strategies enhance the performance of vision language models, they have a tendency to overfit the samples, as described in Section 1. Overfitting leads to the generation of inaccurate and uncertain predictions by models when dealing with samples that either have different distributions or belong to Novel classes [5]. This paper will analyze two main distributions of image data, the in-domain samples $x_{id} \sim \mathbb{P}_x^{id}$ and the out-of-domain samples $x_{ood} \sim \mathbb{P}_x^{ood}$. Extending multiple distributions using the concept of group robustness is achievable [52].

Applying a parametric function to any of the two distributions, x_{id} and x_{ood} , and updating the function's parameters to make \mathbb{P}_x^{id} and \mathbb{P}_x^{ood} indistinguishable from each other has been demonstrated to be effective in the literature [16]. Although techniques such as KL divergence, Optimal Transport measures, can be effective in measuring distributional discrepancies, the maximum mean discrepancy (MMD) is widely favored due to its applicability to various types of data, including imaging, and its well-established statistical properties [16] (page 728) and [11] (page 292).

Lemma 2. $(\Omega_x, \mathcal{F}_x)$ is a measurable space, and \mathbb{P}_x^{id} and \mathbb{P}_x^{ood} are two borel probability measures for in-domain and out-of-domain imaging data. Then $\mathbb{P}_x^{id} = \mathbb{P}_x^{ood}$ if and only if $\mathbb{E}_{x_{id}}(f(x_{id})) = \mathbb{E}_{x_{ood}}(f(x_{ood}))$ for all $f \in C(\Omega_x)$ where $C(\Omega_x)$ is the space of bounded continuous functions on Ω_x .

A function class that is sufficiently rich is required to uniquely determine whether $\mathbb{P}_x^{id} = \mathbb{P}_x^{ood}$ as per Lemma 2. The unit balls in Reproducing Kernel Hilbert space (RKHS) are such a function class that are utilized to establish maximum mean discrepancy (MMD) [16].

MMD measures the distributional discrepancy between in-domain and out-of-domain data; however, challenges occur with multi-modal data, involving vision and language. The straightforward use of MMD necessitates optimization across vision-language, language-language, and vision-vision distribution pairs. We propose calculating the Maximum Mean Discrepancy (MMD) through a vision-language similarity measure. An obvious solution is to utilize a product measure $\mathbb{P}_x \times \mathbb{P}_y$. The product measure theorem (Theorem 4.4.6 from [11]) enables the use of product measures because there exists a unique measure on the product of the σ -algebras $\mathcal{F}_x \times \mathcal{F}_y$. It reminds us of the fundamental [11],

Fact 3. The cardinality of the σ -algebra for the product measure is greater than that of the individual probability measures.

Please refer to page 139 of [11] for more details. The above Fact 3 significantly affects MMD computation, as empirical estimates require more samples now to approximate

expectation calculations. As a solution, we propose utilizing an induced measure to evaluate MMD in lieu of evaluating it over product measures. For the image encoder f_θ and a unit-normed text anchor a_y , let $a_y^T f_\theta(\cdot)$ be a measurable function from Ω_x, \mathcal{F}_x to Γ_x, \mathcal{G}_x . Our induced measure, the anchor-aligned probability measure, is

$$\mathbb{P}_x^{a_y} := \mathbb{P}_x \circ (a_y^T f_\theta)^{-1}(B) = \mathbb{P}_x((a_y^T f_\theta)^{-1}(B)), B \in \mathcal{G} \quad (5)$$

Based on the the anchor-aligned probability measure $\mathbb{P}_x^{a_y}$ defined in Equation 5, the anchor-aligned MMD measure defined as follows:

Definition 4. Given a non-negative, characteristic, and bounded kernel k in a Reproducing Kernel Hilbert Space (RKHS), with a bounded search region θ , and a holder-continuous function f_θ with a specific ratio and support, we consider a text anchor a_y that belongs to the same class $c = k$ as the image embedding vector x . We calculate the maximum mean discrepancy between anchor-aligned feature vectors from two distributions, \mathbb{P}_x^{id} and \mathbb{P}_x^{ood} as

$$\mathcal{L}_{MMD} = \|\mathbb{E}_{x \sim \mathbb{P}_x^{id}} k(f_\theta^T(x)a_y, \cdot) - \mathbb{E}_{x \sim \mathbb{P}_x^{ood}} k(f_\theta^T(x)a_y, \cdot)\| \quad (6)$$

Equation 6 defines the MMD loss between two distributions \mathbb{P}_x^{id} (in-domain) and \mathbb{P}_x^{ood} (out-of-domain) over the anchor-aligned feature space. Here, $k(x, \cdot)$ signifies that the kernel has one parameter fixed at x , while the second parameter, represented by \cdot , is free and can accommodate any random variable. In particular, with the feature mapping onto RKHS as Φ , we have $\langle \Phi(x), \Phi(y) \rangle = k(x, y)$. In the empirical form of the MMD, the expectation is replaced with the sample average calculated over a group. A previous study [69] is of particular relevance because it has demonstrated that calculating the average of samples taken from training batches (sub-samples) is feasible while maintaining consistency and asymptotic normality criteria. Moreover, [69] can be seen as an extension of the paper [3] where the in-domain and out-domain samples are parameterized separately. For a detailed explanation of the consistency and asymptotic properties of MMD and the necessary assumptions, we recommend the reader to visit [68]. However, it is crucial that the anchor-aligned feature vectors require holder continuity. We demonstrate the validity of the statement below

Lemma 5. If f_θ exhibits Holder continuity with a constant α and the anchor vectors are normalized to have unit length, then the anchor-aligned feature vectors likewise exhibit Holder continuity with the same constant α .

Proof Outline.

$$\begin{aligned} \|f_\theta(x_{id})^T a_y - f_\theta(x_{ood})^T a_y\| &\leq \|f_\theta(x_{id}) - f_\theta(x_{ood})\| \|a_y\| \\ &\leq \|x_{id} - x_{ood}\|^\alpha \end{aligned} \quad (7)$$

The first inequality in Equation (7) is valid by the Cauchy-Schwarz inequality. The second inequality is valid because of the unit norm constraint on the anchor aligned vectors a_y . The lemma 5 ensures that the consistency properties and the asymptotic guarantees will remain valid for the feature vectors that are aligned with the anchor.

4. Experiment

In order to assess the effectiveness of our proposed method, we conducted three primary experiments: 1) Base-to-Novel Generalization on 11 datasets, 2) Cross-Domain Robustness on Waterbird and CelebA dataset, 3) Out-of-Distribution Task on ImageNet and its four variant datasets.

4.1. Feature Alignment Improves the Base-to-Novel Generalization

Setting. We utilize 11 publicly available image classification datasets for experiments, including ImageNet [9], Caltech101 [12], OxfordPet [41], Flowers102 [38], StanfordCars [27], Food101 [6], FGVC-Aircraft [36], SUN397 [60], DTD [7], EuroSAT [17], and UCF101 [53]. Each dataset is split into non-overlapping Base and Novel classes. We fine-tune the model in Base classes and evaluate it in Novel classes, we call this setting Base-to-Novel Generalization. The results are displayed in Table 1 and Table 2.

The feature alignment method improves the accuracy of both Base and Novel group classes. In Table 1, for the Base group, the final average results were consistently enhanced by applying $\mathcal{L}_{\text{Aligned}}$ to four prompt tuning structures, with a 2.0-point gain achieved by MaPLe [55] and a 1.9-point gain achieved by PromptSRC [24] on average. For some individual datasets, applying our method $\mathcal{L}_{\text{Aligned}}$ to MaPLe improved the results by 6.1 points on StanfordCars. Additionally, our $\mathcal{L}_{\text{Aligned}}$ loss consistently enhanced performance on novel classes. An average gain of 2.3 points is yielded by $\mathcal{L}_{\text{Aligned}}$ on MaPLe, while an average improvement of 1.5 points is resulted on PromptSRC. Although some datasets, such as Flowers102 and Food101, show only minimal improvements, we attribute this to the already high recognition accuracy (over 90%) in these tasks from the original CLIP model, thus the improved space is not as significant as other tasks (like DTD or EuroSAT). In most datasets, better results are achieved by our method compared to the baseline models.

The best model within feature aligned is also very competitive among state-of-the-art methods. From table 2, it can be seen that while existing prompt tuning methods (E.g., CoOp [71], Co-CoOp [70], ProGrad [73], etc.) perform well on the base class, they perform poorly on the novel class, with recognition accuracy even lower than the original CLIP model. In contrast, our best model achieves better accuracy than the original CLIP model, which indicates the proposed feature alignment method can prevent overfitting

on the prompt model. When compared our best approach, namely $\mathcal{L}_{\text{Aligned}}$ applied on PromptSRC from Table 1, with state-of-the-art methods like KgCoOp [61], ProGrad [73], HPT [57], OGEN [62], and DePT [65]. It is seen that our best model maintains the same recognition level on the base group (with a 0.4-point improvement), and outperforms the state-of-the-art prompt-tuning methods by 1.4 points on the novel class. The better performance on most individual datasets is also achieved by our method, with a maximum gain of 2.3 points on the Base group and 2.8 points on the Novel group. These results further demonstrate that our method is highly competitive with state-of-the-art methods.

4.2. Feature Alignment Improves the Group Robustness

We further evaluated the model’s Group Robustness [66] by using the Waterbirds [46] and CelebA [32] datasets. The Waterbirds dataset is divided into four categories based on background (land or water) and bird type (land bird or water bird). The CelebA dataset is also categorized into four groups by hair color (blond or non-blond) and gender (male or female). The model was trained across all groups, aiming to reduce the classification performance gap between the minority group and the overall groups. This setup tests how well the prompt tuning model handles the spurious correlations.

The proposed feature alignment method improves the worst group accuracy and reduces the performance gap. From Table 3, we can see that models incorporating the feature alignment $\mathcal{L}_{\text{Aligned}}$ consistently reduce the accuracy gap between the worst-group and average-group. Specifically, for the Waterbirds dataset, applying our method to MaPLe reduced the accuracy gap by 5.2 points, and applying it to PromptSRC reduced the gap by 5.4 points. Similar trends are observed for the CelebA dataset, where applying $\mathcal{L}_{\text{Aligned}}$ to the four prompt tuning structures consistently reduced the performance gap. These results demonstrate that our approach effectively enhances the model’s group robustness. Furthermore, compared to other state-of-the-art methods, such as Contrastive-Adapter [66], our best model, $\mathcal{L}_{\text{Aligned}}$ applied to PromptSRC, reduced the performance gap from 5.7 to 1.0, showcasing its competitiveness.

4.3. Feature Alignment and MMD Improves the Out-of-Distribution Recognition

In this setting, the model was trained on the ImageNet dataset and evaluated on its four variants: ImageNetV2 [43], ImageNetSketch [56], ImageNet-A [19], and ImageNet-R [18]. Table 4 shows the results. It is observed that the proposed robust prompt tuning method (indicated by $\mathcal{L}_{\text{Aligned}} + \mathcal{L}_{\text{MMD}}$) consistently improves accuracy across different ImageNet variants (V2, S, A, R) when compared to the baseline prompt tuning method. Regarding the maximum gain on in-

Method	OxfordPets	Flowers102	FGVC-Aircraft	DTD	EuroSAT	StanfordCars	Food101	Caltech101	UCF101	SUN397	ImageNet	Average		
	Base	Novel	Base	Novel	Base	Novel	Base	Novel	Base	Novel	Base	Novel	Base	Novel
LPT [55]	94.8	97.3	96.4	73.6	33.0	28.3	79.4	53.1	87.1	68.9	72.6	74.0	90.2	91.2
Ours-LPT	96.1	97.9	97.6	75.2	37.1	33.0	81.1	56.4	89.3	70.9	73.6	75.1	90.5	91.7
Δ	+1.3	+0.6	+1.2	+1.6	+4.1	+4.7	+1.7	+3.3	+2.2	+2.0	+1.0	+1.1	+0.3	+0.5
VPT [2]	94.8	96.1	85.2	68.5	30.8	33.8	77.7	53.3	89.4	69.0	68.8	73.4	89.3	90.1
Ours-VPT	95.3	96.5	90.0	73.2	34.6	34.7	80.2	57.1	92.9	74.5	70.0	74.8	90.4	91.6
Δ	+0.5	+0.4	+4.8	+4.7	+3.8	+0.9	+2.5	+3.8	+3.5	+5.5	+1.2	+1.4	+1.1	+1.5
MaPLe [55]	95.4	96.4	95.9	72.5	37.4	35.6	80.4	59.2	94.1	75.2	72.9	72.2	90.7	90.7
Ours-MaPLe	95.5	98.0	97.8	75.3	42.3	37.1	81.7	63.9	96.3	78.1	79.0	74.5	89.9	92.2
Δ	+0.1	+1.6	+1.9	+2.8	+4.9	+1.5	+1.3	+4.7	+2.2	+2.9	+6.1	+2.3	-0.8	+1.4
PromptSRC [24]	95.0	97.4	98.0	78.0	42.1	37.0	82.4	59.3	91.4	73.9	77.9	75.2	90.6	91.7
Ours-PromptSRC	95.8	98.0	98.7	77.7	47.0	38.2	84.4	64.0	95.5	79.4	80.8	75.6	91.0	92.2
Δ	+0.8	+0.6	+0.7	-0.3	+4.9	+1.2	+2.0	+4.7	+4.1	+5.5	+2.9	+0.4	+0.4	+0.5

Table 1. **Base-to-Novel Generalization on 11 Datasets.** In this setting, the prompt model is trained on Base group data and evaluated on Novel group data. Our method differs from the baseline in that we apply $\mathcal{L}_{\text{Aligned}}$ to the baseline method. ■ represents the max gains in the Base group, ■ represents the max gains in the Novel group. The table shows that our feature alignment method consistently improves the average classification accuracy across four different prompt-tuning structures.

Method	OxfordPets	Flowers102	FGVC-Aircraft	DTD	EuroSAT	StanfordCars	Food101	Caltech101	UCF101	SUN397	ImageNet	Average		
	Base	Novel	Base	Novel	Base	Novel	Base	Novel	Base	Novel	Base	Novel	Base	Novel
CLIP [42]	91.2	97.3	72.1	77.8	27.2	36.3	53.2	59.9	56.5	64.1	63.4	74.9	90.1	91.2
CoOp [71]	93.7	95.3	97.6	59.7	40.4	22.3	79.4	41.2	92.2	54.7	78.1	60.4	88.3	82.3
Co-CoOp [70]	95.2	97.7	94.9	71.8	33.4	23.7	77.0	56.0	87.5	60.0	70.5	73.6	90.7	91.3
KgCoOp [61]	94.7	97.8	95.0	74.7	36.2	33.6	77.6	55.0	85.6	64.3	71.8	75.0	90.5	91.7
ProGrad [73]	95.1	97.6	95.5	71.9	40.5	27.6	77.4	52.4	90.1	60.9	77.7	68.6	90.4	89.6
HPT [57]	95.8	97.7	98.2	78.4	42.7	38.1	83.8	63.3	94.2	77.1	77.0	74.2	90.5	91.6
OGEN [63]	96.0	97.5	97.3	77.7	41.3	40.3	83.8	62.5	93.4	76.7	77.6	75.2	90.7	91.7
DePT [65]	95.4	97.3	98.4	77.1	45.7	36.7	84.8	61.2	93.2	77.9	80.8	75.0	90.9	91.6
Ours-Best Model	95.8	98.0	98.7	77.7	47.0	38.2	84.4	64.0	95.5	79.4	80.8	75.6	91.0	92.2
Δ	+0.4	+0.7	+0.3	+0.6	+1.3	+1.4	-0.4	+2.8	+2.3	+1.5	+0.0	+0.6	+0.1	+0.6

Table 2. **Ours Best Model Compare with Some State-of-art Methods.** Our best model here is the PromptSRC model comes from Table 1. When compared to the state-of-the-art method, our feature alignment method mitigates the overfitting issues on novel classes, thus leading to higher accuracy even than the SOTA method.

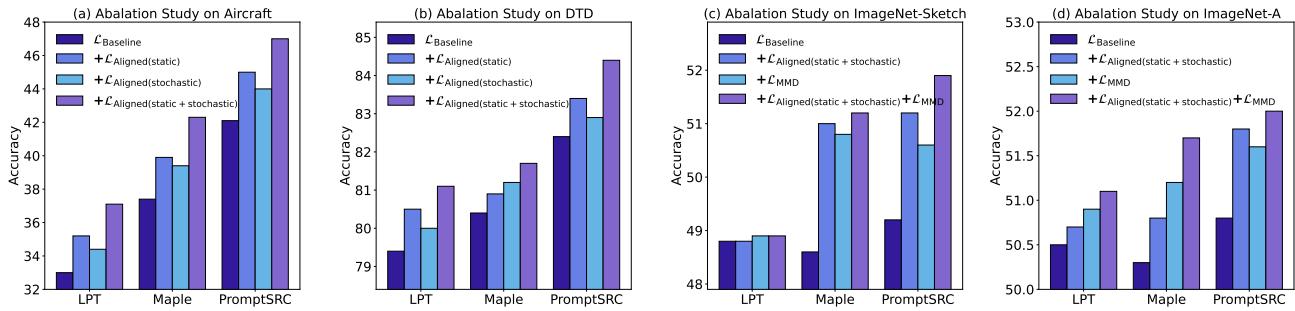


Figure 2. **The ablation study on individual datasets.** $\mathcal{L}_{\text{Baseline}}$ refers to the use of text-based cross-entropy loss in the method. Figures (a) and (b) demonstrate that adding $\mathcal{L}_{\text{Aligned}}(\text{Static})$ or $\mathcal{L}_{\text{Aligned}}(\text{Stochastic})$ can complementarily improve accuracy in in-distribution tasks. Figures (c) and (d) show that adding \mathcal{L}_{MMD} further enhances accuracy across out-of-distribution tasks.

dividual datasets, VPT improved by 2.1 points on ImageNet-A, and PromptSRC by 2.7 points on the ImageNet-Sketch dataset. These improvements illustrate that, compared to traditional cross-entropy methods, our approach narrows the distributional gap between ImageNet and its four variants by utilizing the cross-domain discrepancy match loss, thereby

enhancing its out-of-distribution generalization capabilities on a larger scale. Further analysis of the reasons behind these accuracy enhancements will be conducted in the Ablation Study.

	Waterbirds			CelebA				
	WG	Avg	Gap	△	WG	Avg	Gap	△
ERM Linear Probe [28]	7.9	93.5	85.6		11.9	94.7	82.8	
ERM Adapter [14]	60.8	96.0	35.2		36.1	94.2	58.1	
WiSE-FT [59]	49.8	91.0	41.2		85.6	88.6	3.0	
DFR (Subsample) [26]	63.9	91.8	27.9		76.9	92.5	15.6	
DFR (Upsample) [66]	51.3	92.4	41.1		89.6	91.8	2.2	
Contrastive Adapter [66]	83.7	89.4	5.7		90.0	90.7	0.7	
LPT [55]	78.5	89.6	11.1		78.8	88.6	9.8	
Ours-LPT	79.2	88.8	9.6	-1.5 ↓	80.4	87.9	7.5	-2.3 ↓
VPT [2]	79.3	87.7	8.4		79.9	89.8	9.9	
Ours-VPT	82.1	88.9	6.8	-1.6 ↓	82.2	90.4	8.2	-1.7 ↓
MaPLe [55]	84.6	91.8	7.2		83.6	94.2	10.6	
Ours-MaPLe	85.4	87.4	2.0	-5.2 ↓	85.4	90.1	4.7	-5.8 ↓
PromptSRC [24]	85.7	92.1	6.4		89.4	93.4	3.9	
Ours-PromptSRC	89.6	90.7	1.0	-5.4 ↓	90.2	91.1	0.9	-3.1 ↓

Table 3. **Cross Robustness Evaluation on Waterbird and CelebA Datasets.** Here we applied our $\mathcal{L}_{\text{Aligned}}$ on the baseline method. "WG" refers to the Worst-Group accuracy and "Avg" refers to the Average-Group accuracy. The A downward arrow \downarrow indicates a reduction in the accuracy gap for the groups. The decrease in average accuracy is normal in the group robustness setting, as the goal is to reduce the accuracy gap between groups.

	Source		Target				
	ImageNet	V2	S	A	R	Avg	
ProGrad [73]	70.5	63.4	48.2	49.5	75.2	59.0	
KgCoOp [61]	71.2	64.1	49.0	50.7	76.7	60.1	
CLIPOOD [50]	71.6	64.9	49.3	50.4	77.2	60.5	
HPT [57]	71.7	65.3	49.4	50.9	77.4	60.7	
CoPrompt [44]	70.8	64.3	49.4	50.5	77.5	60.4	
ArGue-N [54]	71.8	65.0	49.3	51.5	77.0	60.7	
LPT [55]	71.7	64.4	48.8	50.5	75.5	59.8	
Ours-LPT	71.9	64.7	48.9	51.1	76.1	60.2	
△	-0.2	+0.3	+0.1	+0.6	+0.6	+0.4	
VPT [2]	69.9	63.1	47.9	43.0	75.9	57.5	
Ours-VPT	70.2	63.3	48.5	45.1	76.0	58.2	
△	+0.3	+0.2	+0.6	+2.1	+0.1	+0.8	
MaPLe [55]	70.2	63.9	48.6	50.3	76.9	59.9	
Ours-MaPLe	70.8	64.1	51.2	51.7	77.8	61.2	
△	-0.6	+0.2	+2.6	+1.3	+0.9	+1.3	
PromptSRC [24]	71.3	64.2	49.2	50.8	77.7	60.5	
Ours-PromptSRC	71.5	65.4	51.9	52.0	77.8	61.8	
△	-0.3	+1.1	+2.7	+1.2	+0.1	+1.3	

Table 4. **Out-of-Distribution Evaluation on ImageNet Variant Datasets.** Here ours method refers to applied both $\mathcal{L}_{\text{Aligned}}$ and \mathcal{L}_{MMD} on the baseline method. In this setting, our method ensures the alignment of ImageNet data from different distributions, resulting in improved accuracy across all four datasets.

$\mathcal{L}_{\text{Baseline}}$	$\mathcal{L}_{\text{Aligned}}(\text{Static})$	$\mathcal{L}_{\text{Aligned}}(\text{Stochastic})$	LPT	VPT	MaPLe	PromptSRC
✓			81.0	78.7	82.3	83.7
✓	✓		82.0	79.9	83.7	84.8
✓		✓	81.6	79.1	82.6	84.0
✓	✓	✓	82.2	80.5	84.2	85.5

Table 5. **Ablation study of $\mathcal{L}_{\text{Aligned}}(\text{Static})$ and $\mathcal{L}_{\text{Aligned}}(\text{Stochastic})$ on Average Performance Across 11 Datasets.** $\mathcal{L}_{\text{Baseline}}$ refers to the use of text-based cross-entropy loss in the method. The results indicate that both $\mathcal{L}_{\text{Aligned}}(\text{Static})$ and $\mathcal{L}_{\text{Aligned}}(\text{Stochastic})$ complementarily improved the final accuracy.

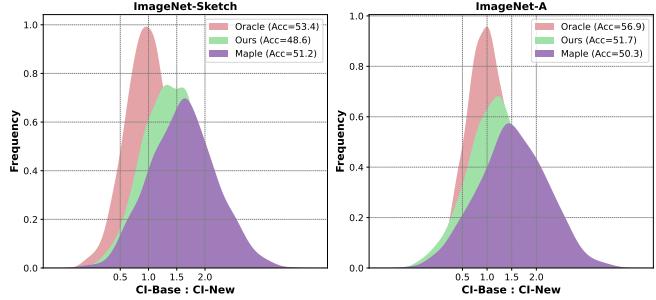


Figure 3. **The effectiveness analysis on channel importance ratio distribution.** The Oracle model is trained on the combination of labeled source and labeled target data, while our model is trained on the labeled source and unlabeled target data. Our \mathcal{L}_{MMD} mitigate the domain shift of the baseline method and push the channel importance distribution close to the Oracle model.

$\mathcal{L}_{\text{Baseline}}$	$\mathcal{L}_{\text{Aligned}}(\text{Static})$	$\mathcal{L}_{\text{Aligned}}(\text{Stochastic})$	\mathcal{L}_{MMD}	LPT	VPT	MaPLe	PromptSRC
✓				59.8	57.5	59.9	60.5
✓	✓			60.1	57.9	60.4	60.7
✓		✓		59.9	57.7	60.1	61.4
✓			✓	60.1	58.0	60.6	61.0
✓	✓	✓	✓	60.2	58.2	61.2	61.8

Table 6. **Ablation study of \mathcal{L}_{MMD} and $\mathcal{L}_{\text{Aligned}}$ on the out-of-distribution task.** The results show that \mathcal{L}_{MMD} significantly enhances classification accuracy, while $\mathcal{L}_{\text{Aligned}}$ provides a slight improvement in accuracy.

4.4. Abalation Studies

Effect of Static and Stochastic Anchors. We first conducted an ablation study on $\mathcal{L}_{\text{Aligned}}(\text{Static})$ and $\mathcal{L}_{\text{Aligned}}(\text{Stochastic})$, presenting the average accuracy across 11 datasets in Table 5. It was evident that combining the static anchor with the baseline method significantly enhanced the accuracy of the four prompt tuning models, with VPT increasing by 1.2 points and MaPLe by 1.4 points. We believe this enhancement was due to the extra supervision information introduced by the image anchor, which helped the model learn better decision boundaries. Furthermore, when combining the two losses, they complemented each other in enhancing model performance, resulting in a 1.9-point improvement in VPT and a 2.0-point improvement in MaPLe. We selected specific datasets and plotted bar charts of their ablation study results. As shown in Figure 2 (a) and (b), applying the Static Anchor or Stochastic Anchor separately with the baseline method led to an accuracy improvement on the Aircraft dataset, while the combination of all three losses yielded the optimal results.

Effect of MMD. Table 6 presents the out-of-distribution ablation experiment on ImageNet and its four variant datasets. We found that applying either the static anchor or stochastic anchor with the baseline method brings limited improvement on the OOD classification accuracy. In

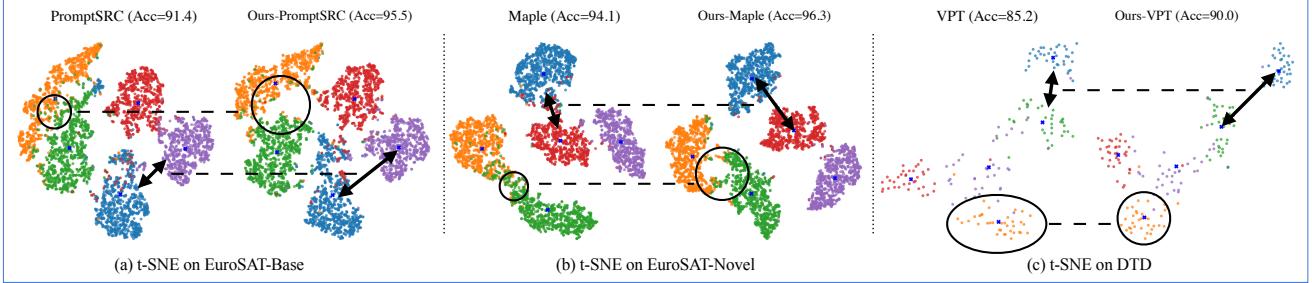


Figure 4. The t-SNE Visualization of Latent Embeddings. The arrows in the three sub-figures illustrate our method can push the boundary between the two categories further apart. The circles in Figures (a) and (b) demonstrate that our method can separate the overlapping features of the two categories away from each other. The circle in Figure (c) shows that our method can achieve a more compact feature space.

contrast, using \mathcal{L}_{MMD} resulted in significant improvements across the four different Prompt Tuning structures. When combining all these loss functions ($\mathcal{L}_{\text{Aligned}} + \mathcal{L}_{\text{MMD}}$), we achieved the best enhancement in final classification performance. The ablation experiments on individual datasets, as shown in Figure 2 (c) and (d), demonstrated similar conclusions, indicating that the combination of all three losses jointly contributes to the improvement of OOD classification accuracy.

We further verify how our \mathcal{L}_{MMD} addresses the domain shift. We train three different models for comparison: 1) the Oracle model, where the prompt is trained on the combination of Target and Source datasets; 2) the Baseline model, where only the text-based cross-entropy loss is used, and only the supervised Source data is used to train the prompt; and 3) our model, where our proposed method is used with both the supervised Source data and the unsupervised Target data. We use the metric of Channel Importance Ratio [65] to show the distribution shift. We conduct experiment on two datasets in Figures 3 (a) and 3 (b). The channel importance distribution shown in these figures demonstrates a clear gap between the Oracle model and the Baseline model. By employing our proposed \mathcal{L}_{MMD} loss, which utilizes cross-domain discrepancy matching to overcome domain shift, we not only bring the CI distribution close to the Oracle model but also increase recognition accuracy.

4.5. Aligned Feature Visualization

Figure 6 presents the results of image feature visualization. In Figure (a), it is evident that applying our $\mathcal{L}_{\text{Aligned}}$ method to MaPLe increases the distance between cluster centers of the same color in the visualized features. This indicates that our method enhances the learned latent space, bringing it closer to real samples, strengthening the model’s decision boundaries, and consequently improving its accuracy. Similar improvements are observed in Figures 6 (b) and (c). Additionally, Figure 6 (c) demonstrates that our method achieves a more compact cluster representation, as

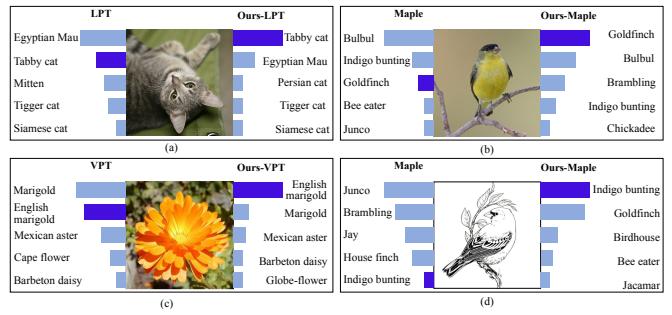


Figure 5. Comparison of Prediction Probabilities With and Without Our Method. Our robust prompt tuning method effectively corrects misclassifications made by the baseline method.

highlighted by the circled areas. These results confirm the effectiveness of our $\mathcal{L}_{\text{Aligned}}$ method.

Figure 5 illustrates the qualitative results of our method on a single sample. In Figure 5 (a), without using the $\mathcal{L}_{\text{Aligned}}$ loss function, an image of a Tabby cat was incorrectly classified as Egyptian Mau. However, after applying the $\mathcal{L}_{\text{Aligned}}$ loss function, the top-1 probability for this image corresponded to the correct category. Figures 5(b), (c) and (d) show similar results.

5. Conclusion

In this paper, we proposed a novel feature alignment method to mitigate the overfitting issue in prompt tuning. Our method combines static and stochastic anchors to learn a more aligned feature space. Based on this aligned space, we applied a modified cross-domain discrepancy matching loss to address domain shift. Experiments demonstrated that our approach outperforms existing methods in Base-to-Novel generalization, group robustness, and out-of-distribution tasks. We hope our work will inspire more research on the adaptation of visual language models.

6. Acknowledgement

This work is partially supported by NSF AI Institute-2229873, NSF RI-2223292, an Amazon research award, and an Adobe gift fund. Prof. Lokhande thanks support provided by University at Buffalo Startup funds. We thank Prof. Won Hwa Kim (POSTECH South Korea) for the insightful discussions on the project.

References

- [1] Dyah Adila, Changho Shin, Linrong Cai, and Frederic Sala. Zero-shot robustification of zero-shot models with foundation models. *arXiv preprint arXiv:2309.04344*, 2023. [14](#)
- [2] Hyojin Bahng, Ali Jahanian, Swami Sankaranarayanan, and Phillip Isola. Exploring visual prompts for adapting large-scale models. *arXiv preprint arXiv:2203.17274*, 1(3):4, 2022. [1, 6, 7](#)
- [3] Mahsa Baktashmotagh, Mehrtash T Harandi, Brian C Lovell, and Mathieu Salzmann. Unsupervised domain adaptation by domain invariant projection. In *Proceedings of the IEEE international conference on computer vision*, pages 769–776, 2013. [4](#)
- [4] Pierre Baldi and Peter J Sadowski. Understanding dropout. *Advances in neural information processing systems*, 26, 2013. [3](#)
- [5] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017. [4](#)
- [6] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. 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*, pages 446–461. Springer, 2014. [5](#)
- [7] 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. [5](#)
- [8] Guillaume Couairon, Matthijs Douze, Matthieu Cord, and Holger Schwenk. Embedding arithmetic of multimodal queries for image retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4950–4958, 2022. [1](#)
- [9] 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. [5](#)
- [10] Xuefeng Du, Yiyou Sun, Jerry Zhu, and Yixuan Li. Dream the impossible: Outlier imagination with diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024. [2](#)
- [11] Richard M Dudley. *Real analysis and probability*. Chapman and Hall/CRC, 2004. [4](#)
- [12] 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. [5](#)
- [13] 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. [2](#)
- [14] Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and Yu Qiao. Clip-adapter: Better vision-language models with feature adapters. *arXiv preprint arXiv:2110.04544*, 2021. [7, 14](#)
- [15] Tianyu Gao, Adam Fisch, and Danqi Chen. Making pre-trained language models better few-shot learners. *arXiv preprint arXiv:2012.15723*, 2020. [1](#)
- [16] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *The Journal of Machine Learning Research*, 13(1):723–773, 2012. [4, 13](#)
- [17] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019. [5](#)
- [18] 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. [5](#)
- [19] 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. [5](#)
- [20] Yu Huang, Chenzhuang Du, Zihui Xue, Xuanyao Chen, Hang Zhao, and Longbo Huang. What makes multi-modal learning better than single (provably). *Advances in Neural Information Processing Systems*, 34:10944–10956, 2021. [1](#)
- [21] Jintian Ji and Songhe Feng. Anchor structure regularization induced multi-view subspace clustering via enhanced tensor rank minimization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19343–19352, 2023. [1](#)
- [22] 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 Conference on Machine Learning*, pages 4904–4916. PMLR, 2021. [1](#)
- [23] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIII*, pages 709–727. Springer, 2022. [1, 2](#)
- [24] Muhammad Uzair Khattak, Syed Talal Wasim, Muzammal Naseer, Salman Khan, Ming-Hsuan Yang, and Fahad Shahbaz Khan. Self-regulating prompts: Foundational model adaptation without forgetting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15190–15200, 2023. [1, 2, 5, 6, 7](#)

- [25] Muhammad Uzair Khattak, Syed Talal Wasim, Muzammal Naseer, Salman Khan, Ming-Hsuan Yang, and Fahad Shahbaz Khan. Self-regulating prompts: Foundational model adaptation without forgetting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15190–15200, 2023. 12
- [26] Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient for robustness to spurious correlations. *arXiv preprint arXiv:2204.02937*, 2022. 7, 14
- [27] 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. 5
- [28] Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution. *arXiv preprint arXiv:2202.10054*, 2022. 7, 14
- [29] Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. *Advances in Neural Information Processing Systems*, 36, 2024. 1
- [30] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021. 1
- [31] Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa, Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training group information. In *International Conference on Machine Learning*, pages 6781–6792. PMLR, 2021. 14
- [32] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015. 5
- [33] Chengcheng Ma, Yang Liu, Jiankang Deng, Lingxi Xie, Weiming Dong, and Changsheng Xu. Understanding and mitigating overfitting in prompt tuning for vision-language models. *IEEE Transactions on Circuits and Systems for Video Technology*, 33(9):4616–4629, 2023. 1
- [34] Chengcheng Ma, Yang Liu, Jiankang Deng, Lingxi Xie, Weiming Dong, and Changsheng Xu. Understanding and mitigating overfitting in prompt tuning for vision-language models. *IEEE Transactions on Circuits and Systems for Video Technology*, 33(9):4616–4629, 2023. 1
- [35] Mohammad Mahdi Derakhshani, Enrique Sanchez, Adrian Bulat, Victor Guilherme Turrisi da Costa, Cees GM Snoek, Georgios Tzimiropoulos, and Brais Martinez. Bayesian prompt learning for image-language model generalization. *arXiv e-prints*, pages arXiv-2210, 2022. 2
- [36] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013. 5
- [37] Luca Moschella, Valentino Maiorca, Marco Fumero, Antonio Norelli, Francesco Locatello, and Emanuele Rodolà. Relative representations enable zero-shot latent space communication. *arXiv preprint arXiv:2209.15430*, 2022. 1, 2
- [38] 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. 5
- [39] Antonio Norelli, Marco Fumero, Valentino Maiorca, Luca Moschella, Emanuele Rodolà, and Francesco Locatello. Asif: Coupled data turns unimodal models to multimodal without training. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [40] Jinyoung Park, Juyeon Ko, and Hyunwoo J Kim. Prompt learning via meta-regularization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26940–26950, 2024. 1
- [41] 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. 5
- [42] 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. 1, 2, 6
- [43] 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. 5
- [44] Shuvendu Roy and Ali Etemad. Consistency-guided prompt learning for vision-language models. *arXiv preprint arXiv:2306.01195*, 2023. 7
- [45] Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016. 3
- [46] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*, 2019. 5
- [47] Cheng Shi and Sibei Yang. Logoprompt: Synthetic text images can be good visual prompts for vision-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2932–2941, 2023. 2
- [48] Manli Shu, Weili Nie, De-An Huang, Zhiding Yu, Tom Goldstein, Anima Anandkumar, and Chaowei Xiao. Test-time prompt tuning for zero-shot generalization in vision-language models. *Advances in Neural Information Processing Systems*, 35:14274–14289, 2022. 2
- [49] Xiangbo Shu, BinQian Xu, Liyan Zhang, and Jinhui Tang. Multi-granularity anchor-contrastive representation learning for semi-supervised skeleton-based action recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(6):7559–7576, 2022. 1
- [50] Yang Shu, Xingzhuo Guo, Jialong Wu, Ximei Wang, Jianmin Wang, and Mingsheng Long. Clipood: Generalizing clip to out-of-distributions. *arXiv preprint arXiv:2302.00864*, 2023. 2, 7

- [51] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017. 3
- [52] Nimit S Sohoni, Maziar Sanjabi, Nicolas Ballas, Aditya Grover, Shaoliang Nie, Hamed Firooz, and Christopher Ré. Barack: Partially supervised group robustness with guarantees. *arXiv preprint arXiv:2201.00072*, 2021. 4
- [53] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. 5
- [54] Xinyu Tian, Shu Zou, Zhaoyuan Yang, and Jing Zhang. Argue: Attribute-guided prompt tuning for vision-language models. *arXiv preprint arXiv:2311.16494*, 2023. 7
- [55] Muhammad Uzair Khattak, Hanoona Rasheed, Muhammad Maaz, Salman Khan, and Fahad Shahbaz Khan. Maple: Multi-modal prompt learning. *arXiv e-prints*, pages arXiv–2210, 2022. 1, 2, 3, 5, 6, 7, 12
- [56] 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. 5
- [57] Yubin Wang, Xinyang Jiang, De Cheng, Dongsheng Li, and Cairong Zhao. Learning hierarchical prompt with structured linguistic knowledge for vision-language models, 2023. 5, 6, 7
- [58] Zehan Wang, Yang Zhao, Haifeng Huang, Jiageng Liu, Aoxiong Yin, Li Tang, Linjun Li, Yongqi Wang, Ziang Zhang, and Zhou Zhao. Connecting multi-modal contrastive representations. *Advances in Neural Information Processing Systems*, 36:22099–22114, 2023. 1
- [59] Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs, Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, et al. Robust fine-tuning of zero-shot models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7959–7971, 2022. 7, 14
- [60] 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. 5
- [61] Hantao Yao, Rui Zhang, and Changsheng Xu. Visual-language prompt tuning with knowledge-guided context optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6757–6767, 2023. 2, 5, 6, 7
- [62] Yuhang Zang, Hanlin Goh, Josh Susskind, and Chen Huang. Overcoming the pitfalls of vision-language model finetuning for ood generalization. *arXiv preprint arXiv:2401.15914*, 2024. 2, 5
- [63] Yuhang Zang, Hanlin Goh, Josh Susskind, and Chen Huang. Overcoming the pitfalls of vision-language model finetuning for ood generalization. *arXiv preprint arXiv:2401.15914*, 2024. 6
- [64] Yuhang Zang, Wei Li, Kaiyang Zhou, Chen Huang, and Chen Change Loy. Unified vision and language prompt learning. *arXiv preprint arXiv:2210.07225*, 2022. 2, 12
- [65] Ji Zhang, Shihan Wu, Lianli Gao, Hengtao Shen, and Jingkuan Song. Dept: Decoupled prompt tuning. *arXiv preprint arXiv:2309.07439*, 2023. 2, 5, 6, 8
- [66] Michael Zhang and Christopher Ré. Contrastive adapters for foundation model group robustness. *Advances in Neural Information Processing Systems*, 35:21682–21697, 2022. 5, 7, 14
- [67] Michael Zhang, Nimit S Sohoni, Hongyang R Zhang, Chelsea Finn, and Christopher Ré. Correct-n-contrast: A contrastive approach for improving robustness to spurious correlations. *arXiv preprint arXiv:2203.01517*, 2022. 14
- [68] Hao Zhou, Vamsi K Ithapu, Sathya Narayanan Ravi, Vikas Singh, Grace Wahba, and Sterling C Johnson. Hypothesis testing in unsupervised domain adaptation with applications in alzheimer’s disease. *Advances in neural information processing systems*, 29, 2016. 4
- [69] Hao Henry Zhou, Vikas Singh, Sterling C Johnson, Grace Wahba, Alzheimer’s Disease Neuroimaging Initiative, Berkeley, and Charles DeCarli. Statistical tests and identifiability conditions for pooling and analyzing multisite datasets. *Proceedings of the National Academy of Sciences*, 115(7):1481–1486, 2018. 4
- [70] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16816–16825, 2022. 1, 5, 6, 12
- [71] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022. 1, 2, 5, 6, 12
- [72] Beier Zhu, Yulei Niu, Yucheng Han, Yue Wu, and Hanwang Zhang. Prompt-aligned gradient for prompt tuning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15659–15669, 2023. 1
- [73] Beier Zhu, Yulei Niu, Yucheng Han, Yue Wu, and Hanwang Zhang. Prompt-aligned gradient for prompt tuning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15659–15669, 2023. 2, 5, 6, 7

A. Overview of Appendix

In addition to the main method and experiments outlined in the paper, we offer supplementary information about our work in the Appendix. In Appendix B, we delve into implementation within the methodology section. We provide more details about our implementation in the Static Anchor and Stochastic Anchor, as well as the Maximum Mean Discrepancy Minimization. Subsequently in Appendix C, we provide more results on dataset setting and adaptation experiments, covering the Base-to-Novel Generalization and group Robustness task. Further in Appendix D, we present additional results of ablation studies along with their visualization outcomes. Lastly in Appendix E and F, we explore the limitations of our work and analyze its broader impact.

B. Method

B.1. Review the Adaption of CLIP

We first review the pretraining and inference stage of the CLIP model, then we discuss the adaptation of CLIP. During the pretraining phase, a large-scale dataset of image-text pairs (x, y) is collected for training the model using a contrastive learning approach. Here x represents an image, and y denotes its corresponding textual description. For each image x , an image encoder model f_θ is parameterized by θ to extract its visual feature vector $u \in \mathbb{R}^{1 \times H}$: $u = f_\theta(x)$. Similarly, for each textual description y , a text encoder g_ϕ is parameterized by ϕ to get its feature embedding $v \in \mathbb{R}^{1 \times H}$: $v = g_\phi(y)$. For the i -th image x_i and the j -th language description y_j in a batch \mathcal{B} , we normalize their feature vectors to a hyper-sphere using: $u_i = \frac{f_\theta(x_i)}{\|f_\theta(x_i)\|}$ and $v_j = \frac{g_\phi(y_j)}{\|g_\phi(y_j)\|}$.

Test phase of CLIP. In this phase, a predefined prompt "a photo of a " is commonly employed for inference. Let's consider a single test image x_{test} of class C , where $x_{\text{test}} \in \mathbb{R}^{C \times H \times W}$ and $C \in \mathbb{R}^K$ for a K -class classification problem. The predefined prompt is prepended to each class label in C to construct the language description. The zero-shot prediction probability for the test image is determined by:

$$Pr(c = k | x_{\text{test}}) = \frac{\exp(\text{sim}(u, v_k) \tau)}{\sum_{i=1}^K \exp(\text{sim}(u, v_i) \tau)} \quad (8)$$

B.2. Introduction of the Prompt Tuning Method

Language Prompt Tuning involves introducing learnable parameters into the text branch. We follow the same notation in [55] and [25]. In the text branch, the class label c is formatted as a language description within a text template as "a photo of a {label}", which can be further transferred as $\tilde{y} = \{t_{SOS}, t_1, t_2, \dots, t_L, c_n, t_{EOS}\}$. Here t_l are the word embeddings of the text template, and c_n are the class label. The

t_{SOS} and t_{EOS} are the learnable start and end token embeddings. The text encoder g encodes the input tokens \tilde{y} through multiple transformer blocks to generate a latent text feature representation $\tilde{g} = g(\tilde{y}, \theta_g)$. In Language prompt tuning, learnable text prompts $\Lambda_{\text{text}} \in \{\lambda_{\text{text}}^1, \lambda_{\text{text}}^2, \dots, \lambda_{\text{text}}^L\}$ are appended to the input \tilde{y} . In CoOp [71] or CoCoOp [70], the learnable text prompts λ_{text} are only added to input of text encoder. While in Maple [55], the learnable text prompts are appended to multiple transformer layers as $[\dots, \tilde{y}_i] = \mathcal{L}_i([\Lambda_{\text{text}}, \tilde{y}_{i-1}]) \quad i = 1, 2, \dots, J$, where \mathcal{L}_i represent the i layer number in the transformer, J represent the prompt depth.

Visual Prompt Learning involves the integration of learnable prompts within the image branch. The input image x is divided into M patches, and these patches are projected to generate patch embeddings $\tilde{x} = \{e_{cls}, e_1, e_2, \dots, e_M\}$, where e_{cls} is the learnable class token. Subsequently, learnable visual prompts are introduced as $\Lambda_{\text{visual}} \in \{\lambda_{\text{visual}}^1, \lambda_{\text{visual}}^2, \dots, \lambda_{\text{visual}}^L\}$. The learnable visual prompts are appended to multiple transformer layers as $[c_i, \tilde{x}_i, \dots] = \mathcal{V}_i([c_{i-1}, \tilde{x}_{i-1}, \Lambda_{\text{visual}}]) \quad i = 1, 2, \dots, J$, where \mathcal{V}_i represent the i layer in vision transformer, J represent the prompt depth.

Multi-modal Prompt Learning integrates language prompt learning and visual prompt learning, combining them synergistically. Simply combining text prompt learning and visual prompt learning is called independent prompt learning, which is used in UPT [64]. To foster interaction between the image and text branches, multi-modal prompt learning employs projection layers $L_t = \{\tilde{l}^1, \tilde{l}^2, \dots, \tilde{l}^J\}$ for projecting the learnable language prompts onto the visual prompts, defined as $\Lambda_{\text{visual}} = \{\tilde{l}^1(\lambda_{\text{text}}^1), \tilde{l}^2(\lambda_{\text{text}}^2), \dots, \tilde{l}^J(\lambda_{\text{text}}^L)\}$. This formulation facilitates interaction between the visual and language prompts. Such unified prompt tuning is a key feature of the Maple [55] framework.

B.3. Static Anchor Implementation

To address the overfitting issues of text-based cross-entropy loss in scenarios with limited data, we propose a symmetrical static anchor alignment method, analogous to an image-based cross-entropy loss. This method involves two primary steps:

Step 1: Construction of Image Anchors. We use a pre-trained CLIP image encoder to extract features for each category in the source dataset, followed by K-means clustering to identify the centroid of each category's features, denoted as a_x^k . It is important to note that the dimensionality of a_x^k differs from that of the batch image features $f_\theta(x)$.

Step 2: Alignment with Text Samples. For each image batch, corresponding text labels represented as language descriptions are aligned, with batch text embeddings $v' = g_\phi(y, \Lambda_{\text{txt}})$, where $v' \in \mathbb{R}^{B \times H}$, also differing in feature

dimensions from class labels.

B.4. Stochastic Anchors Implementation

Stochastic anchors, selected during each batch, can be implemented as cross-modal contrastive learning process. Traditional supervised learning, which models relationships between images and discrete labels, often neglects textual concepts associated with labels. In contrast, stochastic anchors learning fosters understanding of visual concepts through enforcing batch-level text-image alignment.

We construct a contrastive similarity matrix $s' = \text{sim}(u, v')$, where $s' \in \mathbb{R}^{N \times N}$. This matrix supports the formulation of both image-to-text and text-to-image contrastive losses, averaging these to derive the final text-image contrastive loss. In this matrix, only diagonal entries are treated as positive examples, enhancing the robustness of the latent space by introducing a larger set of negative samples.

B.5. Maximum Mean Discrepancy Implementation

Maximum Mean Discrepancy (MMD) is a kernel-based method primarily used to test the equality of two distributions from samples. Introduced in [16], MMD compares the mean embeddings in a feature space, facilitating its use as a loss function in various machine learning tasks, including density estimation, generative modeling, and inverse problems tackled with invertible neural networks. Its simplicity and robust theoretical foundations make MMD particularly advantageous.

To compute MMD for multi-modal data, a product measure is constructed to create a new probability space. Combining two probability spaces increases the complexity of the resulting σ -algebra, necessitating additional samples to characterize the probability space, as noted in Fact 3. We propose Equation 5 to define an induced measure, specifically the anchor-aligned probability measure, as a replacement for the traditional product measure in MMD computation. Equation 5 is essential for the application of MMD in anchor-aligned feature spaces. The transformation of the original probability measure P_x via an anchor-aligned mapping is demonstrated. This equation defines a new probability measure, $P_x^{a_y}$, corresponding to the anchor a_y .

Equation 6 specifies the MMD loss between two distributions, \mathbb{P}_x^{id} (in-domain) and \mathbb{P}_x^{ood} (out-of-domain), within the anchor-aligned feature space. This equation quantifies the discrepancy between two probability distributions in the anchor-aligned feature space. Here is how to use Equation 6 for the current task: The first term computes the expectation over all x_{id} samples, while the second term computes the expectation over all x_{ood} samples. Here, k is the kernel function, and in our experiments, we use the Gaussian kernel. f is the image encoder, and θ is being updated during training. In practical implementations, the expectations are replaced with sample averages.

C. Experiments Protocol

C.1. Base-to-Novel Dataset Split

In our experiment, we partition all class samples into two distinct groups, as outlined in the tables: the Base group (Table 7) and the Novel group (Table 8).

Consider the ImageNet dataset illustrated in Table 9, which consists of 1,000 classes. We divide the training set into two subsets, each containing 500 non-overlapping classes. For instance, one subset may include classes such as ["trench", "goldfish", "great white shark", "tiger shark", ...], and the other might feature ["spindle", "sports car", "spotlight", ...]. This separation ensures that no class from one group appears in the other, thereby preventing the model from encountering unknown classes during training and enhancing the fairness and credibility of our out-of-distribution evaluations. The test set follows a similar bifurcation, maintaining correspondence with the class labels from the training set.

	Classes	Train-Samples	Val-Samples	Test-Samples
OxfordPets	18	288	368	1881
Flowers102	51	816	817	1053
FGVCAircraft	50	800	1667	1666
DTD	23	368	564	864
EuroSAT	5	80	2700	4200
StanfordCars	98	1568	818	4002
Food101	50	800	10100	15300
SUN397	198	3168	1985	9950
Caltech101	50	800	825	1549
UCF101	50	800	949	1934
ImageNet	500	8000	25000	25000

Table 7. **Base class samples statistics.** The first column "Classes" denotes the total number of classes for each category. The columns "Train-Samples", "Val-Samples", and "Test-Samples" represent the respective number of images allocated for model training, validation, and testing purposes.

	Classes	Train-Samples	Val-Samples	Test-Samples
OxfordPets	19	304	368	1788
Flowers102	51	816	816	1410
FGVCAircraft	50	800	1,666	1667
DTD	24	384	564	828
EuroSAT	5	80	2,700	3900
StanfordCars	98	1568	817	4039
Food101	51	816	10,100	15000
SUN397	199	3184	1,985	9900
Caltech101	50	800	824	916
UCF101	51	816	949	1849
ImageNet	500	8000	25000	25000

Table 8. **Novel class samples statistics.** The first column "Classes" denotes the total number of classes for each category. The columns "Train-Samples", "Val-Samples", and "Test-Samples" represent the respective number of images allocated for model training, validation, and testing purposes.

	Classes	Train-Samples	Val-Samples	Test-Samples	Task
OxfordPets	37	2944	736	3669	Fine-Grained
Flowers102	102	4093	1633	2463	Fine-Grained
FGVCAircraft	100	3334	3333	3333	Fine-Grained
DTD	47	2820	1128	1692	Textures
EuroSAT	10	13500	5400	8100	Satellite Images
StanfordCars	196	6509	1635	8041	Fine-Grained
Food101	101	50500	20200	30300	Food
SUN397	397	15880	3970	19850	Scene
Caltech101	100	4128	1649	2465	Object
UCF101	101	7639	1898	3783	Action
ImageNet	1000	12800000	N/A	50000	Object

Table 9. All class samples statistics from the original datasets. The last column "task" provides a broad categorization of these image classification tasks, such as fine-grained classification or texture classification.

C.2. Group Robustness Baseline

For the group robustness experiment described in Section 4.4, we give a more comprehensive introduction about the baseline method that we compared with.

We evaluate our method against several methods in group robustness experiments, including zero-shot classification, ERM linear probing [28], and ERM adapter training [14]. Additionally, we compare against recent approaches tailored to enhancing downstream transfer in analogous scenarios, all while utilizing only pretrained model embeddings [1].

One such method is Weight space ensembling (WiSEFT) [59], which initially trains a linear classifier using standard ERM and then combines the classifier outputs with the initial zero-shot predictions. Although originally proposed for training linear classifiers and fine-tuning the original weights of a foundational model, we focus on the prompt tuning in the extra parameter in our context.

Another approach is Deep feature reweighting (DFR) [26], which entails training a linear probe on embeddings computed from a pretrained model over group-balanced data. Similar to previous studies [31, 67], we treat incorrectly and correctly classified samples as proxies for distinct groups.

Lastly, we consider the Contrastive Adapter approach [66], which introduces contrastive adapting. This method trains adapters with contrastive learning to bring sample embeddings closer to both their ground-truth class embeddings and other sample embeddings within the same class. While our method differs from this work, as we apply Contrastive learning to Prompt Tuning instead of Adapters.

C.3. Training Details

We utilized SGD as the optimizer optimizer with an initial learning rate of 0.0025 for Batch size 4, and a learning rate of 0.01 for batch size 128. The cosine annealing strategy is chosen to schedule the learning rate. For the Base to Novel Generalization setting, we use a few-shot training of 16 shots with a training duration of 20 epochs, while for Group Roubustness, we train 10 epochs on Waterbird and 5 epochs

		Method	Pets	Flowers	Aircraft	DTD	EuroSAT	Cars	Food	Caltech	UCF	Avg
Base	G-Means	95.4	97.6	43.2	82.9	91.9	78.6	90.7	98.2	86.9	85.0	
	H-Cluster	95.2	97.5	43.1	83.0	92.2	78.7	90.5	98.0	86.7	85.0	
	K-Means	95.3	97.5	43.0	83.3	92.4	78.8	90.6	98.1	86.5	85.1	
Novel	G-Means	97.5	77.2	38.2	63.6	68.1	75.4	91.7	94.5	78.2	76.0	
	H-Cluster	97.4	77.1	36.7	64.1	72.5	74.9	90.5	94.2	78.5	76.2	
	K-Means	97.5	77.7	36.9	63.9	79.4	75.2	91.7	94.1	79.1	77.3	

Table 10. **Anchor Selection Method Comparison.** K-means is the default anchor selection method used in this paper. G-Means represents the group means anchor selection method. H-Cluster means hierarchical clustering anchor selection method. The 'Avg' represents the average accuracy over all the datasets.

on CelebA for the full dataset. All images were resized to 224x224 pixels, utilizing the same image preprocessing technique for the CLIP image encoder. All CLIP models adopted the ViT-B/16 backbone. We maintained consistency across all other settings as the baseline work, making modifications solely to the loss function to ensure a fair comparison between our method and the standard cross-entropy loss.

D. More Experiments Results

D.1. Anchor Selections Comparison

To evaluate whether different static anchor selections affect the final results, we conducted the ablation study on the anchor selection experiment, with the results shown in Table 10. We used the pre-trained CLIP model with a ViT-B/16 backbone as the feature extractor. All training images from each dataset were fed into the model's image encoder, and the resulting features were stored. The features are grouped by the ground truth label, then we use different anchor selection methods to choose the most representative one as the static anchor. The anchor selection methods we have are (1) K-means: we utilize the cluster center of K-means as the static anchor; (2) Hierarchical clustering: also the cluster center is utilized as the static anchor; (3) Group means, we direct calculation of the mean features for all the samples in each group. Table 10 shows that K-means and other methods do not have significant differences, while K-means yield better results compared to the hierarchical clustering method.

D.2. t-SNE Visualization

We show more t-SNE visualization results in Figure 6. In Figure (a), it is evident that applying our $\mathcal{L}_{\text{Aligned}}$ method to LPT increases the distance between cluster centers of the green color point and the orange color points. This indicates that our method enhances the learned latent space, bringing it closer to real samples, strengthening the model's decision

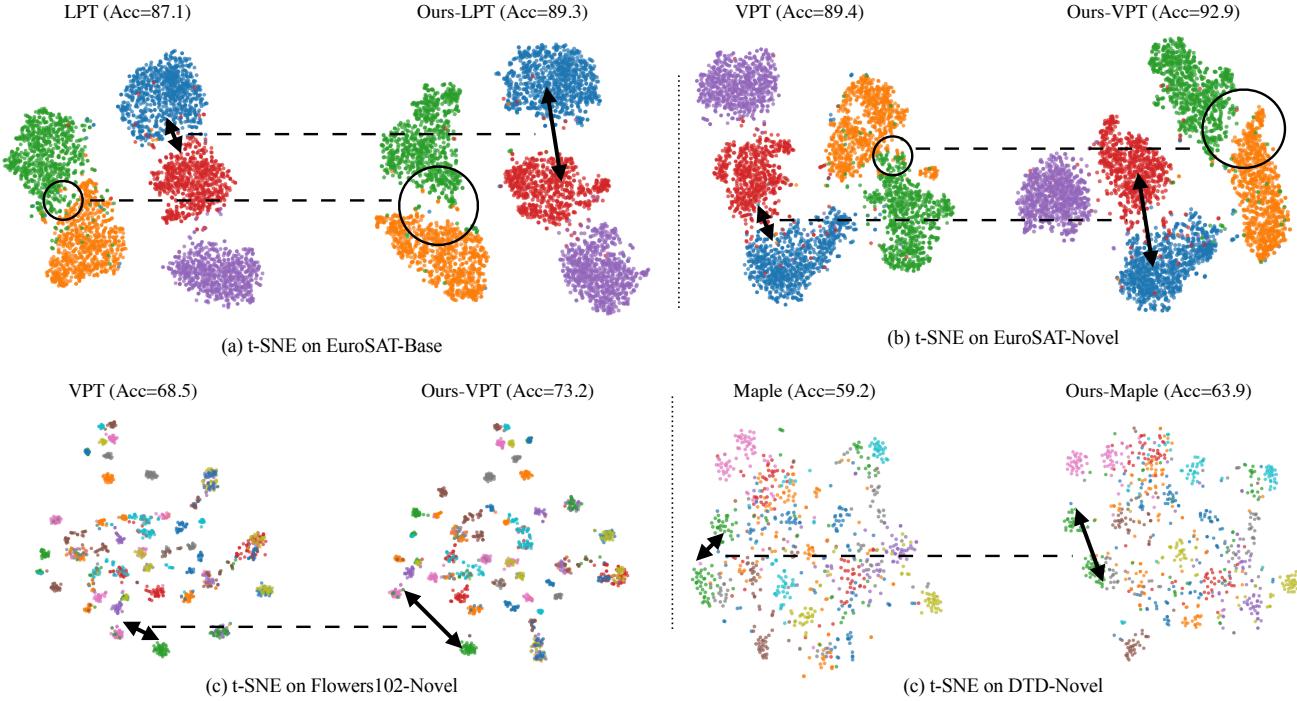


Figure 6. The t-SNE Visualization of Latent Embeddings. The arrows in the figures illustrate our method can push the boundary between the two categories further apart. The circles in Figures (a) and (b) demonstrate that our method can separate the overlapping features of the two categories away from each other.

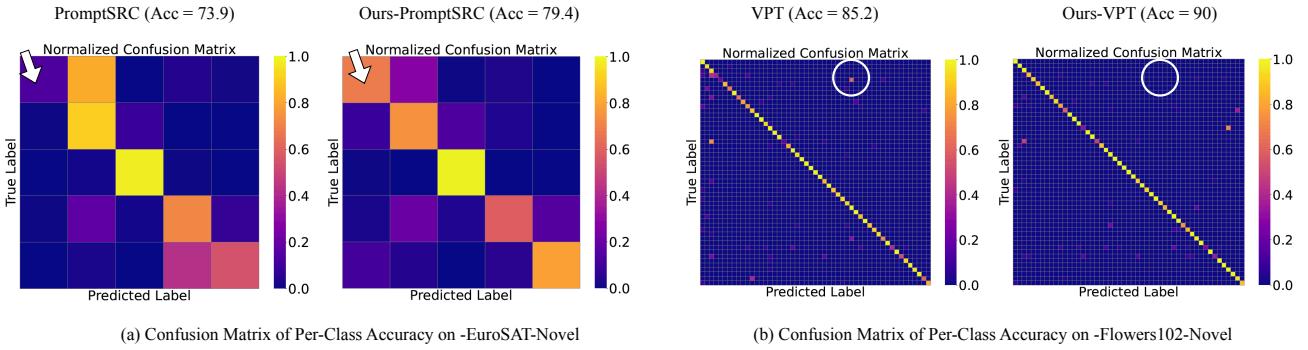


Figure 7. The Confusion Matrix for Per-Class Accuracy. For Figure (a), without our method, the category in the first row is misclassified as the second category. After using our method, the first category classification is successfully made to achieve the highest accuracy. For Figure (b), our method also significantly improves one misclassified subclass, thereby improving the overall accuracy on the entire task.

boundaries, and consequently improving its accuracy. Similar improvements are observed in Figures 6 (b) (c) and (d). Additionally, the circle in Figure 6 (a) and Figure (b) shows that by using our method, we separate the overlap clusters to no-overlap clusters, which also confirms the effectiveness of our $\mathcal{L}_{\text{Aligned}}$ method.

D.3. Confusion Matrix Comparison

To conduct a more granular analysis of the performance improvements brought about by our method, we visualized the confusion matrices representing the accuracy for each category. The experimental results are illustrated in Figure 7. In Figure 7 (a), in the PromptSRC classification experiment on the EuroSAT dataset, the highest value in the first row of the baseline confusion matrix deviated from the diagonal, representing the Pasture Land category, with an accuracy of

only 13.2%. Upon utilizing our $\mathcal{L}_{\text{Aligned}}$ loss function, the first row of the confusion matrix aligned with the diagonal, and the classification accuracy for Pasture Land improved to 68%, which lead to the all-class accuracy improved to 79.4%. Similarly in Figure 7 (b), the figure shows the classification experiments of VPT on the Oxford Flowers dataset. In the confusion matrix of the baseline model, we observed that the classification accuracy for the fifth category, English Marigold, deviated significantly from the diagonal, with an accuracy of only 20%. After applying our proposed $\mathcal{L}_{\text{Aligned}}$ loss function, the classification accuracy for English Marigold increased to 90%.

E. Limitations

Our method aims to construct relative representations in the latent space for cross-modal alignment between image and text modalities, utilizing both static and stochastic anchors. A significant limitation of this method is its high dependency on the selection of the Anchor. For instance, if the static anchor selected does not accurately capture the clustering characteristics of the targeted image category, it may result in biased cross-modal alignment, thereby adversely affecting the learning performance of the model. Additionally, in complex or non-standardized scenarios, finding a suitable static anchor point can be challenging, which constrains the general applicability of our approach

F. Broader Impact

Our proposed approach offers an effective technique applicable to visual language models characterized by an Image-Text dual-branch architecture, which is plug-and-play and can be integrated with many existing prompt tuning methods. Consequently, applying our method to the more sophisticated Prompt Tuning framework could yield further enhancements in performance. We leave these explorations for future research.