

Beyond MaPLe: Enhancing Multimodal Prompt Learning for Vision-Language Tasks

Omar Alandijani
Student ID: g201934090

King Fahd University of Petroleum and Minerals
Dhahran, Saudi Arabia

Supervised by: Dr. Muzammil Behzad
muzammil.behzad@kfupm.edu.sa

King Fahd University of Petroleum and Minerals
Dhahran, Saudi Arabia

Abstract—Multimodal prompt learning has emerged as a promising approach to enhancing the performance of vision-language models such as CLIP. While recent work like MaPLe has demonstrated improvements through better prompt tuning across transformer layers, challenges remain—particularly in generalizing from base to novel classes in low-shot scenarios. In this paper, we propose a series of architectural and training enhancements to MaPLe that aim to improve generalization and reduce overfitting. These include Prompt Dropout, Smarter Prompt Initialization using average CLIP embeddings, and Selective Prompt Injection at critical transformer layers. Our experiments, conducted on multiple benchmarks including ImageNet, Caltech101, and OxfordPets, show consistent gains in novel class accuracy while maintaining strong base class performance. These results highlight the potential of structural prompt tuning strategies in advancing multimodal learning.

Index Terms—Multimodal Learning, Vision-Language Models, Prompt Tuning, CLIP, MaPLe, Zero-Shot Learning, Few-Shot Generalization

I. INTRODUCTION

A. Background and Significance

Vision-language models (VLMs) have recently transformed the landscape of multimodal learning by enabling systems to jointly understand and reason over both visual and textual modalities. One of the most influential breakthroughs in this domain is CLIP (Contrastive Language-Image Pre-training), which aligns images and texts in a shared embedding space using contrastive loss on large-scale internet data [1]. CLIP has demonstrated impressive zero-shot generalization to a wide variety of downstream vision tasks, eliminating the need for task-specific fine-tuning.

B. Challenges in Current Techniques

Although CLIP sets a strong foundation for vision-language understanding, its zero-shot accuracy is heavily dependent on the quality of manually crafted prompts used during inference. Recent works have proposed learning-based approaches to improve prompt generation, such as CoOp [2] and CoCoOp [3], which introduce learnable context vectors instead of static hand-crafted prompts. Building on these ideas, MaPLe (Multimodal Prompt Learning) [4] jointly optimizes prompts in both visual and textual branches, enhancing adaptation to downstream tasks. While MaPLe achieves strong performance on various datasets, its generalization to unseen (novel) classes

remains a challenging aspect, especially under zero-shot conditions.

C. Problem Statement

MaPLe has demonstrated strong performance in base-to-novel generalization by incorporating multimodal prompt learning. However, even with these improvements, there remains room for enhancement—particularly in pushing novel class accuracy further without sacrificing base performance. As zero-shot and few-shot applications increasingly rely on effective generalization, refining MaPLe’s architecture and training strategy offers an opportunity to boost adaptability to unseen categories. This research focuses on enhancing MaPLe’s ability to generalize from base to novel classes through targeted prompt-level modifications.

D. Objectives

The primary objective of this research is to further improve the MaPLe framework’s base-to-novel generalization capabilities in multimodal prompt learning. Specifically, the goals are:

- 1) To replicate and validate the original MaPLe model on standard benchmark datasets to establish a reliable baseline.
- 2) To propose and implement enhancements focused on increasing generalization to novel classes while preserving base class accuracy.
- 3) To evaluate the impact of these enhancements using standard base-to-novel evaluation metrics such as accuracy and harmonic mean.
- 4) To analyze the trade-offs and provide insight into which prompt modifications contribute most to cross-category generalization.

E. Scope of Study

This study is focused on refining the generalization performance of multimodal prompt learning frameworks, specifically under the base-to-novel class evaluation setting. It builds upon the existing MaPLe framework and uses pretrained CLIP models, without altering the underlying encoders. Enhancements are applied at the prompt level, including changes in depth, sharing, and initialization. The study is constrained to lightweight, reproducible modifications and is evaluated using benchmark datasets designed for base-to-novel generalization.

Broader tasks such as domain generalization or cross-dataset transfer are outside the current scope.

II. LITERATURE REVIEW

A. Overview of Existing Techniques

The field of vision-language learning has advanced rapidly with the introduction of models like CLIP (Contrastive Language-Image Pretraining), which align textual and visual representations through contrastive loss on large-scale datasets. CLIP's success stems from its ability to generalize to a wide variety of tasks in a zero-shot manner, eliminating the need for retraining or task-specific fine-tuning [1].

However, while CLIP excels at general representation learning, it struggles to adapt optimally to domain-specific tasks without prompt engineering. Manual prompt tuning proved suboptimal, leading to the development of prompt learning, a technique that optimizes learnable prompt embeddings while keeping the pretrained backbone frozen. CoOp (Context Optimization) [2] was an early successful example, followed by CoCoOp [3], which improved generalization to unseen classes through a class-agnostic meta network.

To address the limitations of purely textual prompts, MaPLE (Multimodal Prompt Learning) [4] introduced a unified multimodal approach by incorporating visual context into the prompt learning process. This allowed the model to leverage both image and text modalities more effectively, achieving state-of-the-art results in zero-shot and few-shot classification, including strong base-to-novel generalization.

B. Related Work

CoOp demonstrated that learning textual prompts could outperform handcrafted ones in few-shot settings. However, its overfitting to seen classes limited its zero-shot transferability. CoCoOp addressed this by introducing a prompt generation mechanism that dynamically adapted prompts based on input features, but it did not incorporate visual cues directly.

MaPLE emerged as a more holistic alternative by designing prompts jointly over both image and text representations. The framework introduced hierarchical prompt layers and a modular learning strategy, resulting in consistent improvements across 11 benchmarks, including ImageNet and OxfordPets [4]. Notably, MaPLE significantly outperformed prior methods in base-to-novel generalization tasks, narrowing the gap between seen and unseen class performance.

More recently, VPT (Visual Prompt Tuning) [5] and IVLP (Interleaved Vision-Language Prompting) [6] have explored variations in visual prompt design and transformer-level conditioning. VPT added tunable visual tokens to the input image patch sequence, while IVLP proposed interleaving language and vision tokens to enable richer cross-modal fusion. These methods reflect a growing interest in lightweight, modular adaptations for foundation models.

C. Limitations in Existing Approaches

Despite these advancements, existing methods still face challenges in maximizing generalization performance, particularly on novel classes. CLIP remains heavily reliant on manual

prompt design, and CoOp and CoCoOp limit prompt learning to the language modality. MaPLE bridges this gap by incorporating both visual and textual prompts, but opportunities remain to refine aspects such as prompt depth, token efficiency, and cross-modal interaction.

Moreover, most prior studies evaluate on fixed benchmarks, leaving open the question of how adaptable these methods are to more generalized base-to-novel settings. There is a need for strategies that maintain strong base class performance while boosting novel class accuracy with minimal tuning effort. This work aims to build on MaPLE's strengths and further explore its capacity to generalize across unseen categories.

III. PROPOSED METHODOLOGY

A. Existing Model and Challenges

Contrastive Language-Image Pre-training (CLIP) has been widely adopted for zero-shot image recognition due to its ability to align textual and visual representations in a shared embedding space. Despite its success, CLIP's reliance on static, manually crafted prompts and its fixed architecture limit its adaptability to downstream tasks, particularly in low-shot learning scenarios.

To address these issues, MaPLE (Multimodal Prompt Learning) [4] was proposed as an extension of CLIP. MaPLE introduces learnable prompt tokens into both vision and text branches of the CLIP transformer, enabling better task adaptation while keeping the core CLIP weights frozen. Though effective in improving base-class performance, MaPLE still exhibits limited generalization to novel categories, particularly when only a few training samples are available.

B. Proposed Enhancements

In this study, we introduce three prompt-level enhancements to the MaPLE framework aimed at improving base-to-novel class generalization without increasing model complexity or altering backbone parameters:

1. Prompt Dropout: We apply dropout directly to the learnable prompt tokens during training. This regularization technique encourages the model to avoid over-reliance on specific prompt representations, improving generalization to unseen classes.

2. Smarter Prompt Initialization: Instead of initializing prompt vectors randomly or with repeated tokens, we use the average of CLIP's text embeddings corresponding to the class names in the training set. This provides a semantically meaningful and stable initialization that grounds the prompts in relevant linguistic context.

3. Selective Prompt Injection: While MaPLE originally injects prompts at all transformer layers, we selectively insert prompts at only three key layers: 1 (shallow), 6 (mid-level), and 11 (deep). This selective approach balances representational diversity with computational efficiency and helps target different levels of abstraction within the vision transformer.

C. Algorithm and Implementation

Our implementation builds upon the official MaPLe GitHub repository. We adapt the existing pipeline to incorporate the above enhancements with minimal modifications to core components.

Prompt tokens are initialized by averaging CLIP text embeddings of class names. During training, we freeze all parameters in CLIP and only update the added prompt tokens. Prompt Dropout is applied with a fixed rate (e.g., 0.2) during training. Prompts are injected into the transformer layers 1, 6, and 11 of the ViT-B/16 backbone.

All experiments are conducted using PyTorch and the `Dassl.pytorch` training engine. We maintain compatibility with MaPLe’s original data pipeline. Datasets are preprocessed with resizing, center cropping, and normalization based on CLIP’s input specifications.

D. Loss Function and Optimization

We adopt the same contrastive loss used in CLIP and MaPLe, based on InfoNCE. Let I and T represent the normalized image and text embeddings, respectively. The loss encourages aligned image-text pairs to have high similarity while pushing apart mismatched pairs:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(\text{sim}(I_i, T_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(I_i, T_j)/\tau)}$$

where $\text{sim}(\cdot)$ denotes cosine similarity and τ is a temperature scaling factor. Optimization is performed using the Adam optimizer with warm-up and cosine decay scheduling, and we monitor training loss and validation accuracy across both base and novel splits.

IV. EXPERIMENTAL DESIGN AND EVALUATION

A. Datasets and Preprocessing

We evaluate our enhanced MaPLe framework on a diverse set of benchmark datasets, following the original MaPLe paper [4]:

- **Caltech101** [7]: Contains images from 101 object categories with varying sample sizes, testing the model’s ability to generalize to diverse object types.
- **Oxford-IIIT Pets** [8]: Comprises 37 classes of pet images with corresponding breed labels, each class having roughly 200 images.
- **Food101** [9]: Includes 101 food categories, each with 1,000 images, introducing a fine-grained classification challenge.
- **Flowers102** [10]: Consists of 102 flower categories, providing a fine-grained classification task.
- **FGVC Aircraft** [11]: Contains 100 aircraft variants, offering a challenging fine-grained classification problem.
- **SUN397** [12]: A scene recognition dataset with 397 categories, testing the model’s ability to generalize across diverse scenes.
- **DTD** [13]: A texture dataset with 47 categories, evaluating the model’s texture recognition capabilities.

- **EuroSAT** [14]: Contains 10 classes of satellite images, assessing the model’s performance on remote sensing data.
- **UCF101** [15]: An action recognition dataset with 101 categories, testing the model’s ability to recognize human actions.

We attempted to include ImageNet and its variants (ImageNetV2, ImageNet-Sketch, ImageNet-A, ImageNet-R) as in the original MaPLe paper. However, due to memory constraints and data access limitations in our current computing environment, these datasets were excluded from our final evaluation.

All datasets were preprocessed following the standard CLIP pipeline: images were resized to 224x224, center-cropped, and normalized using the CLIP-specific mean and standard deviation.

B. Performance Metrics

We evaluate model performance using three key metrics:

- **Base Accuracy (Acc_{base})**: Accuracy on classes seen during training.
- **Novel Accuracy ($\text{Acc}_{\text{novel}}$)**: Accuracy on unseen classes, used to assess generalization.
- **Harmonic Mean (H)**: Combines base and novel accuracies into a single metric that balances both:

$$H = \frac{2 \times \text{Acc}_{\text{base}} \times \text{Acc}_{\text{novel}}}{\text{Acc}_{\text{base}} + \text{Acc}_{\text{novel}}}$$

C. Experiment Setup

All experiments were conducted on Google Colab Pro, utilizing an NVIDIA T4 GPU with 16GB memory. We used the ViT-B/16 variant of CLIP as the backbone model. The MaPLe framework was adapted to include our proposed enhancements, and we employed the original `Dassl.pytorch` engine to maintain compatibility with the training pipeline.

Following the official MaPLe implementation [4], we trained the models for 5 epochs using the SGD optimizer with a learning rate of 0.0035 and a batch size of 4. Prompt Dropout was set to a fixed rate of 0.2. Prompt tokens were injected at transformer layers 1, 6, and 11, and initialized using the average of CLIP text embeddings of class names.

D. Results Comparative Analysis

Table I presents the performance comparison between our enhanced model (*Beyond MaPLe*) and prior methods, including CLIP, CoOp, Co-CoOp, and MaPLe. Across the evaluated datasets, *Beyond MaPLe* consistently outperforms the original MaPLe model in terms of harmonic mean (HM), with particular gains in novel class accuracy.

These improvements validate the effectiveness of our three proposed enhancements. Prompt Dropout helps regularize shallow prompt embeddings, which slightly reduces overfitting on base classes while maintaining strong performance. Smarter Prompt Initialization contributes to semantically meaningful embeddings that improve generalization to unseen classes. Finally, Selective Prompt Injection allows the model to benefit

from key layers without overwhelming the network with redundant prompt information.

Notably, the most significant relative gains appear in datasets with high intra-class variance or domain shifts (e.g., DTD, EuroSAT, Flowers102), highlighting the robustness of our method in low-data or fine-grained settings.

V. EXTENDED CONTRIBUTIONS

The enhancements proposed in this study—namely Prompt Dropout, Smarter Prompt Initialization, and Selective Prompt Injection—not only improve base-to-novel generalization within the MaPLe framework, but also highlight the value of structural prompt-level interventions in multimodal learning. By demonstrating that targeted, lightweight modifications can yield consistent gains without altering the backbone architecture, this work opens up a pathway for broader adoption in resource-constrained settings.

Furthermore, our design decisions emphasize modularity and reproducibility, making it easier for future researchers to adapt similar strategies to other prompt-based or transformer-based architectures. The methodology also encourages the community to shift focus from brute-force scaling to more principled architectural refinements that respect computational budgets while improving performance.

VI. CONCLUSION AND FUTURE WORK

In this paper, we proposed a set of enhancements to the MaPLe framework aimed at improving its ability to generalize from base to novel classes in vision-language tasks. Our contributions—Prompt Dropout, Smarter Prompt Initialization, and Selective Prompt Injection—focus on making prompt learning more robust and semantically grounded. Through experiments on a wide range of benchmark datasets, we observed consistent improvements in novel class performance, validating the effectiveness of our design.

While our current results are promising, there remains ample opportunity for future work. One direction is to explore adaptive prompt configurations that vary dynamically per input or dataset. Another avenue is to integrate learnable scaling mechanisms or attention-guided injection points, which were beyond the scope of this study due to time constraints. We also hope to extend our enhancements to other modalities and tasks beyond classification, including VQA and captioning.

REFERENCES

- [1] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International Conference on Machine Learning (ICML)*, 2021.
- [2] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, “Learning to prompt for vision-language models,” in *CVPR*, 2022.
- [3] —, “Conditional prompt learning for vision-language models,” in *CVPR*, 2022.
- [4] M. M. Khattak, Z. Lin, J. Yang, K. Zhou, and Z. Liu, “Maple: Multimodal prompt learning,” in *NeurIPS*, 2023.
- [5] M. Jia, L. Tang, B.-C. Chen, C. Cardie, and S. Belongie, “Visual prompt tuning,” in *ECCV*, 2022.
- [6] Z. Lin, C. Wang, M. M. Khattak, and Z. Liu, “Ivlp: Improving vision-language pretraining with prompt tuning,” in *CVPR*, 2023.
- [7] L. Fei-Fei, R. Fergus, and P. Perona, “Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories,” in *CVPR*, 2004.
- [8] O. Parkhi, A. Vedaldi, and A. Zisserman, “Cats and dogs,” in *CVPR*, 2012.
- [9] L. Bossard, M. Guillaumin, and L. Van Gool, “Food-101 – mining discriminative components with random forests,” in *ECCV*, 2014.
- [10] M.-E. Nilsback and A. Zisserman, “Automated flower classification over a large number of classes,” in *Indian Conference on Computer Vision, Graphics and Image Processing*, 2008.
- [11] S. Maji, E. Rahtu, J. Kannala, M. Blaschko, and A. Vedaldi, “Fine-grained visual classification of aircraft,” in *FGVC*, 2013.
- [12] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba, “Sun database: Large-scale scene recognition from abbey to zoo,” in *CVPR*, 2010.
- [13] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi, “Describing textures in the wild,” in *CVPR*, 2014.
- [14] P. Helber, B. Bischke, A. Dengel, and D. 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*, 2019.
- [15] K. Soomro, A. R. Zamir, and M. Shah, “Ucf101: A dataset of 101 human actions classes from videos in the wild,” in *CRCV-TR-12-01*, 2012.

TABLE I
COMPARISON WITH STATE-OF-THE-ART METHODS ON BASE-TO-NOVEL GENERALIZATION.

Method	OxfordPets			Caltech101			Food101			Flowers102		
	Base	Novel	HM	Base	Novel	HM	Base	Novel	HM	Base	Novel	HM
CLIP	91.17	97.26	94.12	96.84	94.00	95.40	90.10	91.22	90.66	72.08	77.40	74.83
CoOp	93.67	95.29	94.47	98.00	88.06	92.74	90.60	91.50	91.05	97.60	69.71	81.47
Co-CoOp	95.20	97.69	96.43	97.96	93.81	95.84	90.70	91.29	90.99	94.87	71.56	81.34
MaPLe	95.43	97.76	96.58	97.74	94.36	96.02	90.71	92.05	91.38	94.36	71.66	82.56
Beyond MaPLe	94.98	98.21	96.56	97.40	95.80	96.59	89.95	93.40	91.64	93.90	74.21	82.97

Method	FGVCAircraft			SUN397			DTD			EuroSAT		
	Base	Novel	HM	Base	Novel	HM	Base	Novel	HM	Base	Novel	HM
CLIP	27.19	36.29	31.09	69.36	75.35	72.23	53.24	59.90	56.37	56.48	64.05	60.03
CoOp	40.44	22.61	29.06	76.47	65.89	70.78	79.44	41.18	54.24	91.29	54.74	68.36
Co-CoOp	33.41	23.71	27.81	79.74	76.76	78.23	76.60	56.00	64.85	92.30	71.28	80.41
MaPLe	37.44	35.61	36.50	80.82	78.70	79.75	80.36	59.18	68.16	94.07	73.23	82.35
Beyond MaPLe	36.92	38.44	37.66	79.80	80.15	79.97	78.40	63.55	70.21	93.75	75.64	83.72

Method	UCF101					
	Base	Novel	HM			
CLIP	70.53	77.51	73.85			
CoOp	84.69	51.65	64.27			
Co-CoOp	82.33	73.45	77.63			
MaPLe	83.00	78.66	80.77			
Beyond MaPLe	82.41	80.21	81.30			