# Learning Domain Invariant Prompt for Vision-Language Models

Cairong Zhao, Yubin Wang, Xinyang Jiang, Yifei Shen, Kaitao Song, Dongsheng Li, and Duoqian Miao

*Abstract*—Prompt learning is one of the most effective and trending ways to adapt powerful vision-language foundation models like CLIP to downstream datasets by tuning learnable prompt vectors with very few samples. However, although prompt learning achieves excellent performance over in-domain data, it still faces the major challenge of generalizing to unseen classes and domains. Some existing prompt learning methods tackle this issue by adaptively generating different prompts for different tokens or domains but neglecting the ability of learned prompts to generalize to unseen domains. In this paper, we propose a novel prompt learning paradigm that directly generates *domain invariant* prompt that can be generalized to unseen domains, called MetaPrompt. Specifically, a dual-modality prompt tuning network is proposed to generate prompts for input from both image and text modalities. With a novel asymmetric contrastive loss, the representation from the original pre-trained vision-language model acts as supervision to enhance the generalization ability of the learned prompt. More importantly, we propose a meta-learning-based prompt tuning algorithm that explicitly constrains the task-specific prompt tuned for one domain or class to also achieve good performance in another domain or class. Extensive experiments on 11 datasets for base-to-new generalization and 4 datasets for domain generalization demonstrate that our method consistently and significantly outperforms existing methods.

*Index Terms*—Prompt learning, meta-learning, few-shot learning, domain generalization.

## I. INTRODUCTION

**R**ECENT research in pre-training large Vision-Language Models (VLM) using web-scale data has shown remarkable progress in learning transferable representations [13], [15]. Compared with traditional supervised learning methods, which learn close-set visual concepts from discrete labels, these models align images in a joint embedding space via contrastive learning, providing a promising opportunity to leverage human language for guiding visual recognition tasks. Benefiting from this paradigm, pre-trained vision-language models can conduct zero-shot or few-shot transfer to downstream tasks with open-set visual concepts learned from natural language supervision. As a result, how to effectively leverage these powerful foundation models becomes an important research direction. Recent studies [20], [21] apply a simple yet effective way to adapt pre-trained vision-language models to downstream tasks, called prompting. Manually designing a proper prompt is a non-trivial task due to its ambiguity, which makes automatic prompt tuning the current mainstream approach. Drawing inspiration from recent advances of prompt learning [17], [18], [19] in NLP, methods like CoOp [20], CoCoOp [21] and MaPLe [69] learn a set of continuous vectors as the context in a prompt (i.e., prompt vector) with the pre-trained parameters fixed, which achieve significant improvement with very few training samples.

Although showing promising performance in i.i.d samples, as discussed by previous works [21], prompt learning still faces a substantial challenge of domain generalization. Like other machine learning methods, conventional prompt tuning approaches [20] tend to overfit the distribution of the training set. When transferred to unseen domains, the good generalization ability of foundation models is compromised, and the learned prompt vectors suffer a significant accuracy drop when transferred to unseen domains. Even with massive tuning, we could not yet guarantee an optimal prompt for the downstream tasks. Recently, several methods [21], [23] have tackled this issue by adaptively generating different prompts for different tokens or domains, known as conditional prompt learning. However, they fail to exploit the generalization ability of the prompt generator or learned prompt, and do not explicitly enforce the prompt to generalize to unseen domains.

In this paper, our goal is to explicitly learn domain invariant prompt for vision-language models. Such prompts are independent of the input instance and should have a low bias toward the visual representations of the target task. In essence, with the distribution shift in text and image modality, both two can be categorized as cross-domain tasks, as the test samples are out-of-domain. As discussed in previous literature [64], [65], [66], input samples are composed of attributes (i.e., factors of variation), such as color, shape, texture, etc., and different domains are defined by different distributions of each attribute. As a result, there exists a unified meta-domain containing all possible attributes, where data domains are attribute distributions sampled from this meta-domain. Under this assumption, our theoretical analysis following [43] shows that tuning prompts with an episodic training strategy has a strong generalization guarantee. Specifically, this type of method has the generalization bound of $O(1/\sqrt{N})$ with $N$ being the number of tasks, independent of the sample size in each domain, which motivates us to propose an episodic
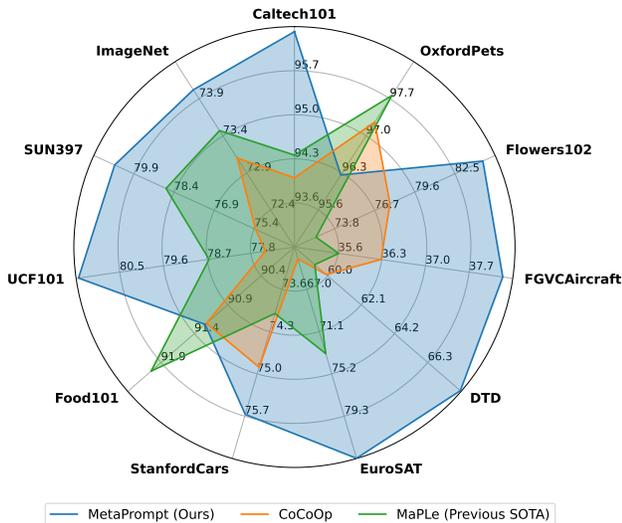
Fig. 1. Comprehensive comparison of the harmonic mean of previous methods CoCoOp, MaPLe, and our method MetaPrompt on 11 diverse image recognition datasets for base-to-new generalization. MetaPrompt surpasses state-of-the-art methods on 9 of 11 datasets, proving that learning domain invariant prompt achieves a good trade-off between in-domain and out-of-domain data.

prompt tuning method in the few-shot setting.

As a result, leveraging the power of episodic training, we propose MetaPrompt, a simple but effective few-shot approach that generates domain invariant prompt for vision-language models. To explicitly enforce the learned prompt to generalize to unseen domains, a novel episodic prompt tuning algorithm is proposed, where we optimize the prompt when trained on a certain domain that can produce good results on samples out of this domain. Furthermore, a dual-modality prompt tuning network is proposed, where two sets of prompt vectors are learned for input from both image and text modalities, respectively. To overcome overfitting on in-domain data and leverage the strong generalization ability of the pre-trained vision-language model, we propose an asymmetric contrastive loss where the representation from the unprompted stream of text modality acts as supervision to learn visual prompt, and vice versa, aiming to learn a set of domain invariant prompts for both modalities.

In this paper, the ability of domain generalization is evaluated from two perspectives, new image domains and new class domains. Our MetaPrompt is applicable for both out-of-domain images (i.e., conventional domain generalization task) and classes (i.e., base-to-new generalization task). As shown in Fig. 1, for base-to-new generalization, MetaPrompt obtains an overall improvement of harmonic mean accuracy by an average gain of 0.54% over the previous state-of-the-art method MaPLe on 11 image recognition datasets. For domain generalization, although based on a few-shot setting, our method achieves comparable performance over other methods training on full datasets with only 1-shot and 5-shot settings. These experimental results demonstrate the effectiveness of our method and show superiority in generalization ability to other prompt tuning approaches.

## II. RELATED WORK

### A. Prompt Learning

Prompt learning emerges from recent advances in natural language processing. The core idea of prompt learning is to formalize various tasks [14], [13], [16] to masked language modeling problems with different prompt templates. A prompt can be seen as a function of the input tokens, providing instruction for pre-trained language models such as BERT [14] or GPT [16] to adapt to downstream tasks. Earlier work [71] enabled the model to understand the task and make better predictions by manually designing discrete natural language prompts. However, some hand-crafted prompt templates are inappropriate in many cases due to their ambiguity, and the recognition performance of such methods is susceptible to the form of the provided content. Recent methods [17], [18], [19] learn continuous contexts to model prompts, called prompt tuning, to automate prompt engineering and explore optimal prompts. This paradigm can also be applied to vision-language models [13], [15]. Specifically, CoOp [20] demonstrates that the performance of CLIP is sensitive to prompts, and a suitable prompt can be learned with very few samples for image recognition. CoCoOp [21] extends CoOp by learning an input-conditional token for each image to obtain generalizable representations. ProDA [70] captures the distribution of diverse prompts to handle the varying visual representations and provides high-quality task-related content for facilitating recognition.

Although these approaches consider prompt learning for text modality, they neglect to tune prompts for generating visual features. To fill this gap, Visual Prompt Tuning (VPT) [22] achieves significant performance gains with only a small amount of trainable parameters as a prompt while keeping the model backbone frozen. MaPLe [69] proposes a prompt tuning method for both vision and language branches to improve alignment between the vision and language representations. In contrast, with an explicit constraint on prompt tuning, our method learns domain invariant prompt for both modalities, resulting in better generalization.

### B. Domain Generalization

Domain generalization refers to learning a robust model generalized to unseen domains. In this paper, the generalization ability of a model is evaluated from the perspectives of both out-of-domain images and classes, corresponding to the conventional domain generalization task and base-to-new generalization task. Conventional domain generalization mainly evaluates the generalization ability on unseen image domains. Many approaches [9], [10], [11], [12] attempt to measure the domain gap and learn domain invariant features. In order to learn a set of parameters that can generalize to unseen domains, several methods [67], [58] adopt meta-learning to simulate domain shift during training. In this paper, we provide a theoretical analysis on the generalization guarantee of meta-learning based on episodic training and incorporate episodic training in prompt tuning for the first time. In contrast to previous methods, we treat the regularization process for generalization as a constraint after regular steps separately,

aiming to create various episodes within one single batch to optimize the domain invariant prompt with less computational complexity.

Recently, another type of generalization task emerges called base-to-new generalization, which aims to exploit the generalization ability on unseen classes [1], [2], [3], [4]. Conventional methods [5], [6], [7], [8] learn a semantic space based on auxiliary information. Compared with supervised learning, CLIP-based methods achieve high performance in generalization due to the more vital transferring ability. CoCoOp [21] attempts to tackle this generalization problem with conditional prompt learning. We investigate the feasibility of learning domain invariant prompts for the pre-trained vision-language model CLIP [13] and propose a training strategy to implement this goal.

### C. Meta-Learning

Most existing meta-learning approaches focus on few-shot learning, which can be divided into metric learning methods, memory network methods, and optimization-based methods. Metric learning methods [25], [26], [27], [28] learn a similarity space to extract discriminative meta-features for new classes efficiently. Memory network methods [29], [30], [31], [32] store meta-knowledge by memory models when learning seen tasks and then generalize it to unseen tasks. Optimization-based methods [24], [33], [34], [35] train meta-optimizer that enable fast adaption for new tasks. Works like MAML [24], [36], [37], [38] focus on learning meta-initial parameters of a deep model so that it would perform well on new tasks after only a small number of gradient updates. Drawing on the advances, our work is the first to propose learning meta-prompt for visual-language foundation models generalizing to unseen domains. Instead of learning the initial parameters of the model, we regularize parameters after every conventional update to learn robust representations. In concrete, we utilize gradients on meta-test subtasks to regularize parameters, i.e., prompts. By imposing this constraint, our model learns robust representations and performs better on base-to-new generalization and domain generalization tasks.

### III. GENERALIZATION BOUND OF EPISODIC TRAINING

Following previous literature [64], our theoretical analysis is based on the assumption that data is composed of attributes (i.e., factors of variation), such as color, shape, texture, etc., and different domains can be defined by different distributions of attributes. For example, as shown in Fig. 2, a sketch domain corresponds to a color distribution with only two values, black and white. In contrast, a cartoon or natural image domain may correspond to a color distribution with more color values. As a result, we assume that there exists a unified meta-domain distribution $\tau$ containing all possible attributes, where data domains $\mathcal{P} = \{\mathcal{P}_i\}_{i=1}^N$ are distributions sampled from this meta-domain with different attribute distributions. Under this assumption, we expect a training strategy to learn invariant features from seen domains and be able to generalize to unseen domains. Specifically, given a training algorithm $\mathbf{F}$ trained on a dataset $\mathbf{D} = \{D_i = D_i^s\}_{i=1}^N$ drawn from a domain distribution
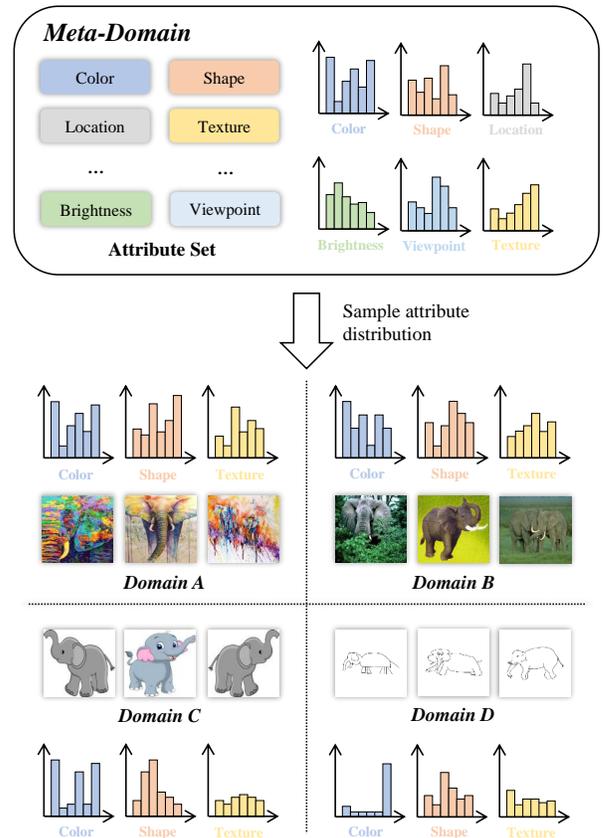


Fig. 2. Input samples are composed of attributes (i.e., factors of variation), such as color, shape, texture, etc., and different domains can be defined by different distributions of attributes.

$\mathcal{P}_i^M$ containing $M$ training samples (i.e., $D_i^s \overset{i.i.d.}{\sim} \mathcal{P}_i^M$), the generalization error of the model obtained by $\mathbf{F}(\mathbf{D})$ is as follows:

$$\mathcal{R}(\mathbf{F}(\mathbf{D}), \tau) = \mathbb{E}_{\mathcal{P} \sim \tau, D^s \sim \mathcal{P}^M, z \sim \mathcal{P}} L(\mathbf{F}(\mathbf{D})(D^s), z). \quad (1)$$

To improve the generalization ability of meta-learning algorithms, the pioneering work [25] proposes a training strategy – episodic training strategy, which treats each task as a training instance and updates the inner-task algorithm by episode (task by task). In this paper, we transfer episodic training to the domain generalization scenario by treating each data domain as a training instance and update the inner-domain algorithm by episode (domain by domain). Specifically, we first update the model on a support domain (i.e., in-domain error). Then the performance of the updated model is measured and optimized on another query domain (i.e., out-of-domain error or episodic training error). As a result, the training error of the episodic training strategy is as follows:

$$\hat{\mathcal{R}}_{epi}(\mathbf{F}(\mathbf{D}), \mathbf{D}) = \frac{1}{N} \sum_{i=1}^N \frac{1}{N_i^q} \sum_{z_i \in D_i^q} \hat{L}(\mathbf{F}(D_i^s), z_i), \quad (2)$$

where $D_i^q$ is the set of data sampled from a query domain; and $N_i^q$ is the number of samples in $D_i^q$. From Eq. 2 we can
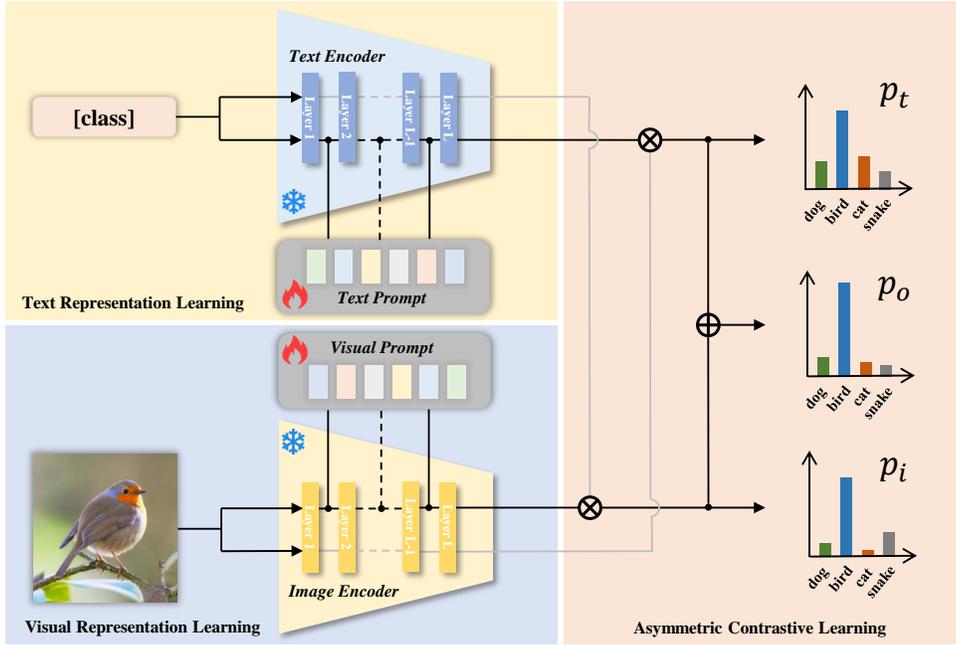
Fig. 3. Our dual-modality prompt tuning network consists of a text encoder for textual representation learning and an image encoder for visual representation learning. Each encoder has two streams, where the unprompted stream (marked by the gray line) from one modality guides the prompted stream (marked by the black line) from another modality by tuning corresponding prompts. The asymmetric contrastive learning module outputs three probability distributions for the end-to-end training to make better predictions.

see that episodic training strategy directly minimizes the out-of-domain testing error, and hence intuitively the in-domain sample number $M$ in the generalization bound vanishes, with the generalization bound only depending on the domain number $N$.

Based on this paradigm, we naturally associate episodic training with domain generalization task, aiming to learn invariance from various distributions by creating meta-tasks with domain gap as episodes. By applying this strategy, the distribution shift between the meta-train and meta-test subtask can be approximately equivalent to that between the original training and test task. The error of the parameter over the meta-test task is exactly the test error of generalization tasks and thereby is an unbiased estimate of the generalization error on unseen domains. Theoretically, following [43], we derive the bound of the generalization gap between these two errors only depending on the domain number $N$, which is formulated by:

$$\mathbb{E}_{\mathbf{F}}[\mathcal{R}(\mathbf{F}(\mathbf{D}), \tau)] \leq \mathbb{E}_{\mathbf{F}}\left[\hat{\mathcal{R}}_{epi}(\mathbf{F}(\mathbf{D}), \mathbf{D})\right] + O\left(\frac{1}{\sqrt{N}}\right). \quad (3)$$

The generalization bound implies a strong generalization guarantee for episodic training algorithms in the few-shot regime, which motivates this paper to adopt episodic training to learn domain invariant prompt with very few samples.

## IV. METHODOLOGY

In this section, we elaborate on our prompt tuning method, MetaPrompt. Sec. 4.1 provides a brief overview of existing prompt tuning approaches for text and image modalities. Sec. 4.2 introduces our dual-modalities prompt tuning network

and demonstrates a novel asymmetric contrastive loss for prompt tuning. Sec. 4.3 introduces our batch-wise episodic training paradigm for learning invariant prompt.

### A. Dual-Modality Prompt Tuning

This paper focuses on prompt tuning for vision-language foundation models, such as CLIP, which are usually composed of a text encoder and an image encoder. The text encoder adopts a transformer [68] to encode textual information while the image encoder can either be a CNN model like ResNet [41] or a vision transformer like ViT [42] to encode visual concepts. Among recent works on prompt tuning, prompt vectors can be learned for both text [20] and image encoder [22]. We describe the setting of prompt tuning for text and image modalities as follows:

*a) Textual Prompt Tuning:* [20] automatically learns a set of tunable continuous vectors as context tokens that are fed into the text encoder together with the class tokens. Given the textual prompt composed of $P$ vectors for the $i$-th class denoted as $\boldsymbol{t}_i$, the prediction probability of the $i$-th class can be calculated by:

$$p_t(y = i \mid \boldsymbol{x}) = \frac{\exp\left(\text{sim}\left(\boldsymbol{x}, g\left(\boldsymbol{t}_i\right)\right) / \tau\right)}{\sum_{j=1}^{K} \exp\left(\text{sim}\left(\boldsymbol{x}, g\left(\boldsymbol{t}_j\right)\right) / \tau\right)}, \quad (4)$$

where $\boldsymbol{x}$ represents the image feature extracted from the image encoder and $g(\cdot)$ denotes the text encoder.

*b) Visual Prompt Tuning:* [22] adopts a similar idea as textual prompt, where extra prompt vectors are automatically learned to be fed into the image encoder. The image patches are firstly embedded into a latent space as the input of the first transformer layer, and then $P$ learnable vectors are introduced

at $L$ transformer layers' input space as prompt. The output of the transformer head is considered the final visual feature $\widetilde{\boldsymbol{x}}$. The prediction probability of the $i$-th class can be calculated by:

$$p_i(y = i \mid \boldsymbol{x}) = \frac{\exp\left(\operatorname{sim}\left(\widetilde{\boldsymbol{x}}, g\left(\boldsymbol{h}_i\right)\right)/\tau\right)}{\sum_{j=1}^{K} \exp\left(\operatorname{sim}\left(\widetilde{\boldsymbol{x}}, g\left(\boldsymbol{h}_j\right)\right)/\tau\right)}, \quad (5)$$

where $g(\cdot)$ denotes the text encoder and $\boldsymbol{h}_i$ denotes the textual feature of the $i$-th class token.

### B. Asymmetric Contrastive Learning

Motivated by previous works on textual and visual prompt tuning, we propose a dual-modality prompt tuning network that jointly learns visual and textual prompts for each transformer layer. As shown in Fig. 3, two sets of learned prompt vectors are fed into the text and image encoder of a foundation model together with the image and text input, respectively.

To prevent the learned prompt vectors from overfitting the in-domain training samples (especially in a few-shot learning setting), we propose to leverage the strong generalization ability of the pre-trained vision language model by a novel Asymmetric Contrastive Loss (AC Loss). AC loss utilizes representations from the unprompted pre-trained vision-language model as supervision to enhance the generalization ability of prompts and learn a set of domain invariant prompts for both modalities. Specifically, instead of training both textual and visual prompts simultaneously with a single contrastive loss, we train them separately, where prompted representations from one modality are aligned with unprompted ones from another modality. For example, as shown in Fig. 3, for learning of visual prompt, the prediction probability is obtained based on the similarity between unprompted textual features and prompted visual features.

With our asymmetric contrastive learning module, we have two probability distributions $p_t$ and $p_i$, corresponding to textual and visual prompts with Eq. 4 and Eq. 5. We sum them up to obtain an overall probability distribution $p_o$, which is formulated by:

$$p_o(y \mid \boldsymbol{x}) = \frac{1}{2}(p_t(y \mid \boldsymbol{x}) + p_i(y \mid \boldsymbol{x})). \quad (6)$$

During training, the cross-entropy loss is adopted to minimize the distance between the ground-truth label $y$ and the three probability distributions $p_t$, $p_i$ and $p_o$, where the losses are denoted as $\mathcal{L}_t$, $\mathcal{L}_i$ and $\mathcal{L}_o$. As a result, AC Loss can be expressed as the sum of the above three losses:

$$\mathcal{L}_{AC}(y, \boldsymbol{x}) = \mathcal{L}_o(y, \boldsymbol{x}) + \mathcal{L}_t(y, \boldsymbol{x}) + \mathcal{L}_i(y, \boldsymbol{x}) \quad (7)$$

### C. Batch-wise Episodic Training

Motivated by the analysis from Section 3, we propose a batch-wise episodic training paradigm for prompt learning. Given a batch of training data containing samples from different domains, we separate the batch into a support set and a query set based on domains. Our proposed episodic training strategy aims to regularize learnable prompts that narrow the gap between training errors on the support set and query set.
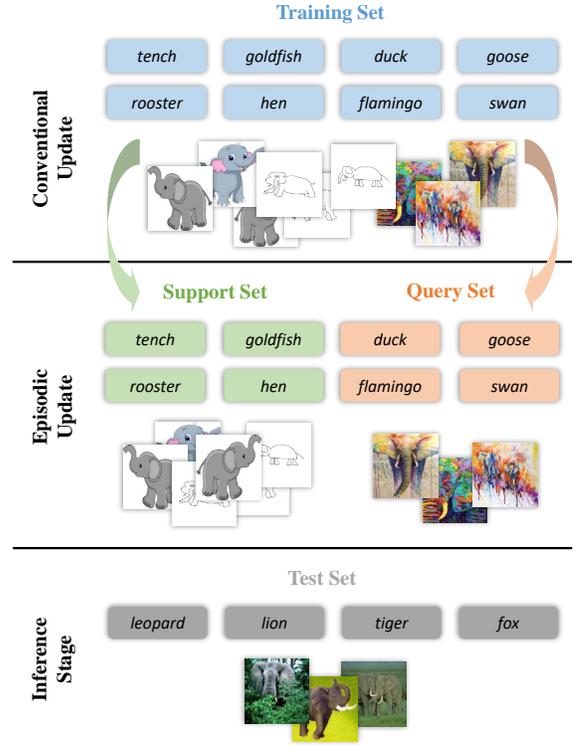


Fig. 4. Given a batch of training data containing samples from different domains, we conduct a conventional update and several episodic updates. In the inference stage, with the domain invariant prompt, we evaluate the generalization ability on unseen domains.

Specifically, given a set of $N$ datasets sampled from $N$ domains denoted as $\mathbf{D} = \{D_i\}_{i=1}^{N}$ where $D_i \sim \mathcal{P}_i$, we split the set by grouping samples from one of the datasets as the query set $D_i^q$, and samples from the rest as the support set $D_i^s$. Note that, for domain generalization, since it is clear which domain each sample belongs to, the query and support set can be easily split. However, for base-to-new generalization, there is no explicit definition of which domain each sample belongs to. Hence, we randomly split the query and support set based on the class label of each sample, as shown in Fig. 4.

To maintain good in-domain performance, which is an essential metric for base-to-new generalization, we adopt an alternated update strategy, where the query/support episodic update and the conventional gradient descent update are conducted alternately. After a conventional update, the learnable prompt $\theta$ is updated with the samples on the support set $D_i^s$ to get the updated prompt $\theta_i'$. Then the generalization error of the updated prompt $\theta_i'$ is measured by the cross-entropy loss on the query set $D_i^q$, whose corresponding gradients are back-propagated to update the original prompt $\theta$. Since this update involves second-order gradient computation with high complexity, in our implementation, we design a first-order approximating method. The parameter $\theta$ is updated as follows:

$$\theta \leftarrow \theta - \alpha\eta \sum_i \nabla_{\theta_i'} \mathcal{L}(\theta_i'; D_i^q), \quad (8)$$

where $\alpha$ is the meta-step size, and $\eta$ is the learning rate of the

---

**Algorithm 1** Batch-wise Episodic Training

---

**Require:** Domain size $N$, learning rate $\eta$, meta-step size $\alpha$, dataset $\mathcal{D}$, loss functions $\mathcal{L}_{AC}$, $\mathcal{L}_t$, $\mathcal{L}_i$

**Ensure:** Prompt parameters $\Theta$

1: Randomly initialize $\Theta_0 = \{\theta^I, \theta^T\}$
2: **for** $t$ in iterations **do**
3:     Randomly sample a batch $\mathcal{D}_t$ from $\mathcal{D}$
4:     **Conventional Update:**
5:     Update $\Theta_t$ w.r.t. $\mathcal{L}_{AC}$:
        $\Theta_t \leftarrow \Theta_t - \eta \nabla_{\Theta_t} \mathcal{L}_{AC}(\Theta_t; \mathcal{D}_t)$
6:     **Episodic Update:**
7:     **if** base-to-new generalization **then**
8:         $\theta \leftarrow \theta_t^T$, $\mathcal{L} \leftarrow \mathcal{L}_t$
9:         $\mathbf{D} = \{D_i\}_{i=1}^N \leftarrow group\_by\_class(\mathcal{D}_t)$
10:     **else if** conventional domain generalization **then**
11:         $\theta \leftarrow \theta_t^I$, $\mathcal{L} \leftarrow \mathcal{L}_i$
12:         $\mathbf{D} = \{D_i\}_{i=1}^N \leftarrow group\_by\_domain(\mathcal{D}_t)$
13:     **end if**
14:     **for** $i = 1$ to $N$ **do**
15:         $D_i^q \leftarrow D_i$
16:         $D_i^s \leftarrow \bigcup_{j=1, j \neq i}^N D_j$
17:         $\theta_i' \leftarrow \theta - \alpha \nabla_\theta \mathcal{L}(\theta; D_i^s)$
18:         $g_i \leftarrow \nabla_{\theta_i'} \mathcal{L}(\theta_i'; D_i^q)$
19:     **end for**
20:     Update $\theta$ with gradients:
        $\theta \leftarrow \theta - \alpha \eta \sum_{i=1}^N g_i$
21:     **if** base-to-new generalization **then**
22:         $\Theta_{t+1} \leftarrow \{\theta_t^I, \theta\}$
23:     **else if** conventional domain generalization **then**
24:         $\Theta_{t+1} \leftarrow \{\theta, \theta_t^T\}$
25:     **end if**
26: **end for**

---

conventional training. $\mathcal{L}$ indicates the generalization loss on the query set for calculating gradients. To simplify the training process, our paradigm treats one batch-wise iteration in Eq. 8 as a series of training episodes and conducts several splits of the query and support set within each batch iteration.

Furthermore, based on our dual-modality prompt tuning network, we propose a modality-specific optimization strategy where the prompt of only one specific modality is tuned during the episodic update step. For example, since base-to-new generalization focuses on generalizing to unseen classes, only textual prompts are tuned with the loss function $\mathcal{L}_t$. Similarly, for conventional domain generalization, only visual prompts are updated with the loss function $\mathcal{L}_i$ to generalize to unseen images. The detailed implementation of the episodic training is shown in Alg. 1.

## V. EXPERIMENTS

We evaluate our approach mainly in the two generalization settings, i.e. base-to-new generalization and conventional domain generalization. In our experiments, we use the open-source CLIP [13] as the foundation vision-language model. Here we elaborate on the experimental configurations.

*a) Datasets:* For base-to-new generalization, we follow Zhou et al. [20] and evaluate the performance of our method using 11 image recognition datasets, which cover a wide range of recognition tasks. Specifically, the benchmark includes ImageNet [44] and Caltech101 [45] for classification on generic objects; OxfordPets [46], StanfordCars [47], Flowers102 [39], Food101 [48] and FGVCAircraft [49] for fine-grained classification; SUN397 [50] for scene recognition; UCF101 [51] for action recognition; DTD [40] for texture classification; and finally EuroSAT [52] for satellite imagery recognition. For each dataset, we split the classes equally into two groups as base and new classes. We train the model only on base classes in a few-shot setting, while evaluation is conducted separately on base and new classes.

For conventional domain generalization experiments, we select four real-world datasets from DomainBed benchmark, including VLCS [53], PACS [54], OfficeHome [55], Domain-Net [56]. We conduct experiments with the leave-one-out strategy. For a dataset, one of the domains is selected as the target domain at a time, and other domains are used as the source domains. We train the model on the source domains in a few-shot setting and evaluated on the target domain.

*b) Implementation Details:* We apply prompt tuning on the pre-trained CLIP model with ViT-B/16 as the visual backbone. Both prompts are randomly initialized from Gaussian distribution with mean of 0 and a standard deviation of 0.02. We adopt SGD optimization with an initial learning rate of 0.002, decayed by the cosine annealing rule, and the meta-step size $\alpha$ is set to 0.2. The warming-up trick is adopted during the first epoch with a fixed learning rate of $10^{-5}$. During inference, we use the overall distribution $p_o$ for prediction.

For base-to-new generalization, the maximum epoch is set to 10 for all datasets with a batch size of 16. The prompt length $P$ and the prompt layer $L$ of visual and textual prompts are set to 2 and 12. Following Zhou et al. [20], we use the few-shot evaluation protocol selecting 16 shots for training and the whole test set for evaluation. For conventional domain generalization, the maximum epoch is set to 5 for all datasets with a batch size of 32. The prompt length $P$ and the prompt layer $L$ of visual and textual prompts are set to 4 and 10. We adopt 1 and 5 shots for each source domain combining them as the training set and test on the target domain. For the hyper-parameter selection of our implementation, we share the same set of hyper-parameters instead of searching for each dataset.

### A. Base-to-New Generalization

The performance of our MetaPrompt in base-to-new generalization setting on 11 image recognition datasets is shown in Table I. We compare its performance with zero-shot CLIP with hand-crafted prompts and recent prompt learning methods, including CoOp, CoCoOp and MaPLe.

*a) Generalization to Unseen Classes:* In comparison with the state-of-the-art prompt tuning method MaPLe, MetaPrompt obtains an overall improvement to 75.48% in terms of the average accuracy of new classes over 11 datasets with our episodic training strategy that explicitly constrains the prompt to generalize to out-of-domain classes. When considering both base and new classes, MetaPrompt shows an

TABLE I
**COMPARISON OF CLIP, CoOp, CoCoOp, MaPLe, AND OUR METAPROMPT ON BASE-TO-NEW GENERALIZATION.** OUR EXPERIMENTS ARE REPEATED THREE TIMES USING DIFFERENT RANDOM SEEDS. METAPROMPT OUTPERFORMS ALL OTHER METHODS ON BOTH BASE AND NEW CLASSES AND DEMONSTRATES STRONG GENERALIZATION PERFORMANCE ON 11 IMAGE RECOGNITION DATASETS. H: HARMONIC MEAN (TO HIGHLIGHT THE GENERALIZATION TRADE-OFF).

**(a) Average over 11 datasets.**

| | Base | New | H |
|---|---|---|---|
| CLIP | 69.34 | 74.22 | 71.70 |
| CoOp | 82.69 | 63.22 | 71.66 |
| CoCoOp | 80.47 | 71.69 | 75.83 |
| MaPLe | 82.28 | 75.14 | 78.55 |
| MetaPrompt | **83.65** | **75.48** | **79.09** |
| vs. MaPLe | +1.37 | +0.34 | +0.54 |

(b) ImageNet.

| | Base | New | H |
|---|---|---|---|
| CLIP | 72.43 | 68.14 | 70.22 |
| CoOp | 76.47 | 67.88 | 71.92 |
| CoCoOp | 75.98 | 70.43 | 73.10 |
| MaPLe | 76.66 | 70.54 | 73.47 |
| MetaPrompt | **77.52** | **70.83** | **74.02** |
| vs. MaPLe | +0.86 | +0.29 | +0.55 |

(c) Caltech101.

| | Base | New | H |
|---|---|---|---|
| CLIP | 96.84 | 94.00 | 95.40 |
| CoOp | 98.00 | 89.81 | 93.73 |
| CoCoOp | 97.96 | 93.81 | 95.84 |
| MaPLe | 97.74 | 94.36 | 96.02 |
| MetaPrompt | **98.13** | **94.58** | **96.32** |
| vs. MaPLe | +0.39 | +0.22 | +0.30 |

(d) OxfordPets.

| | Base | New | H |
|---|---|---|---|
| CLIP | 91.17 | 97.26 | 94.12 |
| CoOp | 93.67 | 95.29 | 94.47 |
| CoCoOp | 95.20 | 97.69 | 96.43 |
| MaPLe | 95.43 | **97.76** | **96.58** |
| MetaPrompt | **95.53** | 97.00 | 96.26 |
| vs. MaPLe | +0.10 | -0.76 | -0.32 |

(e) StanfordCars.

| | Base | New | H |
|---|---|---|---|
| CLIP | 63.37 | 74.89 | 68.65 |
| CoOp | **78.12** | 60.40 | 68.13 |
| CoCoOp | 70.49 | 73.59 | 72.01 |
| MaPLe | 72.94 | 74.00 | 73.47 |
| MetaPrompt | 76.34 | **75.01** | **75.48** |
| vs. MaPLe | +3.40 | +1.01 | +2.01 |

(f) Flowers102.

| | Base | New | H |
|---|---|---|---|
| CLIP | 72.08 | **77.80** | 74.83 |
| CoOp | 97.60 | 59.67 | 74.06 |
| CoCoOp | 94.87 | 71.75 | 81.71 |
| MaPLe | 95.92 | 72.46 | 82.56 |
| MetaPrompt | **97.66** | 74.49 | **84.52** |
| vs. MaPLe | +1.74 | +2.03 | +1.96 |

(g) Food101.

| | Base | New | H |
|---|---|---|---|
| CLIP | 90.10 | 91.22 | 90.66 |
| CoOp | 88.33 | 82.26 | 85.19 |
| CoCoOp | 90.70 | 91.29 | 90.99 |
| MaPLe | 90.71 | **92.05** | **91.38** |
| MetaPrompt | **90.74** | 91.85 | 91.29 |
| vs. MaPLe | +0.03 | -0.20 | -0.09 |

(h) FGVCAircraft.

| | Base | New | H |
|---|---|---|---|
| CLIP | 27.19 | 36.29 | 31.09 |
| CoOp | **40.44** | 22.30 | 28.75 |
| CoCoOp | 33.41 | 23.71 | 27.74 |
| MaPLe | 37.44 | 35.61 | 36.50 |
| MetaPrompt | 40.14 | **36.51** | **38.24** |
| vs. MaPLe | +2.70 | +0.90 | +1.74 |

(i) SUN397.

| | Base | New | H |
|---|---|---|---|
| CLIP | 69.36 | 75.35 | 72.23 |
| CoOp | 80.60 | 65.89 | 72.51 |
| CoCoOp | 79.74 | 76.86 | 78.27 |
| MaPLe | 80.82 | 78.70 | 79.75 |
| MetaPrompt | **82.26** | **79.04** | **80.62** |
| vs. MaPLe | +1.44 | +1.34 | +0.87 |

(j) DTD.

| | Base | New | H |
|---|---|---|---|
| CLIP | 53.24 | **59.90** | 56.37 |
| CoOp | 79.44 | 41.18 | 54.24 |
| CoCoOp | 77.01 | 56.00 | 64.85 |
| MaPLe | 80.36 | 59.18 | 68.16 |
| MetaPrompt | **83.10** | 58.05 | **68.35** |
| vs. MaPLe | +2.74 | -1.13 | +0.19 |

(k) EuroSAT.

| | Base | New | H |
|---|---|---|---|
| CLIP | 56.48 | 64.05 | 60.03 |
| CoOp | 92.19 | 54.74 | 68.69 |
| CoCoOp | 87.49 | 60.04 | 71.21 |
| MaPLe | **94.07** | 73.23 | 82.35 |
| MetaPrompt | 93.53 | **75.21** | **83.38** |
| vs. MaPLe | -0.54 | +1.98 | +1.03 |

(l) UCF101.

| | Base | New | H |
|---|---|---|---|
| CLIP | 70.53 | 77.50 | 73.85 |
| CoOp | 84.69 | 56.05 | 67.46 |
| CoCoOp | 82.33 | 73.45 | 77.64 |
| MaPLe | 83.00 | **78.66** | 80.77 |
| MetaPrompt | **85.33** | 77.72 | **81.35** |
| vs. MaPLe | +2.33 | -0.94 | +0.58 |

absolute average gain of 0.54% on the harmonic mean over MaPLe. The results strongly prove that our method of learning domain invariant prompt improves the generalization ability.

*b) Performance Gain in Seen Classes:* While MetaPrompt achieves excellent performance on generalizing to unseen classes, it still maintains high accuracy on seen classes compared with other methods optimized to fit in-domain data, even better than MaPLe by 1.37%. MetaPrompt achieves a good trade-off between in-domain and out-of-domain data for two reasons. Firstly, our dual-modality prompts improve the recognition accuracy from two modalities simultaneously. With the unprompted pre-trained vision-language model as supervision, we obtain a stable boost in fitting both in-domain and out-of-domain data. Secondly, from the perspective of training strategies, MaPLe does not explicitly consider the in-domain and out-domain trade-off and achieving good generalization at the expense of lower in-domain accuracy,

while our approach proposes an explicit episodic training strategy to learn domain invariant prompt for both seen and unseen classes.

### B. Conventional Domain Generalization

The performance of our MetaPrompt in conventional domain generalization setting on four benchmarks is shown in Table II. We compare its performance with different categories of domain generalization methods, including the non-ensemble methods like ERM [57], MLDG [58], Fish [59], CORAL [60], the ensemble methods like SWAD [61], EoA [62], SEDGE [63], as well as zero-shot CLIP and CoCoOp in domain generalization setting. Since extracting domain invariant features is the mainstream idea in traditional domain generalization tasks, we follow this idea for CLIP-based learning to train domain invariant prompt.

TABLE II
**COMPARISON OF DOMAIN GENERALIZATION METHODS AND OUR METAPROMPT ON FOUR DOMAIN GENERALIZATION BENCHMARKS.** CLIP
(TEMPLATE) INDICATES USING 'A PHOTO OF A {CLASS NAME}' PROMPT. 'ENSEMBLE', 'CLIP' INDICATE ENSEMBLE AND CLIP-BASED METHODS. OUR
EXPERIMENTS ARE REPEATED THREE TIMES USING DIFFERENT RANDOM SEEDS. ALTHOUGH OUR METHOD IS BASED ON *few-shot* SETTING, IT ACHIEVES
COMPETITIVE RESULTS AGAINST FULL-TRAINING METHODS AND DEMONSTRATES STRONG PERFORMANCE ON DOMAIN GENERALIZATION
BENCHMARKS.

| Method | Setting | | Category | | Accuracy(%) | | | |
|---|---|---|---|---|---|---|---|---|
| | Zero-shot | Few-shot | Ensemble | CLIP | PACS | VLCS | OfficeHome | DomainNet |
| ERM [57] | | | | | $84.2 \pm 0.1$ | $77.3 \pm 0.1$ | $67.6 \pm 0.2$ | $44.0 \pm 0.1$ |
| MLDG [58] | | | | | $84.8 \pm 0.6$ | $77.1 \pm 0.4$ | $68.2 \pm 0.1$ | $41.8 \pm 0.4$ |
| Fish [59] | | | | | $85.5 \pm 0.3$ | $77.8 \pm 0.3$ | $68.6 \pm 0.4$ | $42.7 \pm 0.2$ |
| CORAL [60] | | | | | $86.2 \pm 0.3$ | $78.8 \pm 0.6$ | $68.7 \pm 0.3$ | $41.5 \pm 0.1$ |
| SWAD [61] | | | ✓ | | $88.1 \pm 0.1$ | $79.1 \pm 0.1$ | $70.6 \pm 0.2$ | $46.5 \pm 0.1$ |
| EoA [62] | | | ✓ | | $95.8 \pm 0.0$ | $81.1 \pm 0.0$ | $83.9 \pm 0.0$ | $60.9 \pm 0.0$ |
| SEDGE [63] | | | ✓ | | $96.1 \pm 0.0$ | $82.2 \pm 0.0$ | $80.7 \pm 0.2$ | $54.7 \pm 0.1$ |
| CLIP [13] | ✓ | | | ✓ | $95.7 \pm 0.0$ | $75.9 \pm 0.0$ | $79.4 \pm 0.0$ | $57.9 \pm 0.0$ |
| CLIP (template) | ✓ | | | ✓ | $96.1 \pm 0.0$ | $\mathbf{82.3 \pm 0.0}$ | $82.1 \pm 0.0$ | $57.8 \pm 0.0$ |
| CoCoOp [21] (5-shot) | | ✓ | | ✓ | $96.7 \pm 0.4$ | $78.3 \pm 1.0$ | $84.1 \pm 0.1$ | $61.1 \pm 0.2$ |
| MetaPrompt (1-shot) | | ✓ | | ✓ | $96.7 \pm 0.6$ | $81.7 \pm 0.6$ | $84.0 \pm 0.5$ | $61.5 \pm 0.2$ |
| MetaPrompt (5-shot) | | ✓ | | ✓ | $\mathbf{96.9 \pm 0.3}$ | $82.0 \pm 0.9$ | $\mathbf{85.1 \pm 0.4}$ | $\mathbf{61.8 \pm 0.2}$ |

In comparison with traditional domain generalization methods, CLIP-based methods show excellent generalization performance due to the strong transfer learning ability of CLIP. Although training with very few samples, our MetaPrompt provides competitive results in domain generalization benchmarks, by outperforming all other methods on three of four benchmarks on average accuracy in the 5-shot setting and achieving comparable performance even with the 1-shot setting. Furthermore, our method outperforms the conditional prompt tuning method CoCoOp on all datasets, showing a much better generalization ability to unseen domains. By simulating the generalization error between different domains during training, our domain invariant prompt is more generalizable than a conditional-based prompt generator training separately with domains.

But it is worth noting that our approach suffers some performance degradation in the few-shot regime on datasets with a large domain distribution shift, such as VLCS, which indicates that the current strategy still conducts an approximated estimate of generalization error.

## C. Further Analysis

*a) Influence of Model Components:* We analyze the influence of components in our model and conduct an ablation study on various combinations of them, as shown in Table III. The baseline method simultaneously trains both textual and visual prompts with a conventional gradient descent optimizer. The results show that both batch-wise episodic training strategy and asymmetric contrastive loss positively affect generalization to unseen domains. Among them, AC loss achieves an absolute performance gain on new class domains and an overall boost on new image domains, which shows the effectiveness of leveraging unprompted representations of pre-trained vision-language foundation models. Our episodic training strategy also plays an important role in boosting the ability of generalization, which will be analyzed in the subsequent section. In addition, our modality-specific optimization strategy further improves performance on both tasks.

TABLE III
**ABLATION ON DIFFERENT COMPONENTS.** 'EPISODIC' DENOTES OUR BATCH-WISE EPISODIC TRAINING STRATEGY. 'MOS' INDICATES USING OUR MODALITY-SPECIFIC OPTIMIZATION STRATEGY INSTEAD OF REGULARIZING PROMPTS FOR BOTH MODALITIES IN BOTH TASKS DURING EPISODIC UPDATES. FOR DOMAIN GENERALIZATION, WE USE AN AVERAGE OF 1-SHOT AND 5-SHOT ACCURACY AS EVALUATION METRICS, THE SAME BELOW.

(a) Base-to-New Generalization.

| Episodic | AC Loss | MOS | Base | New | H |
|---|---|---|---|---|---|
| | | | 82.73 | 73.05 | 77.24 |
| ✓ | | | 82.88 | 74.96 | 78.62 |
| | ✓ | | 83.19 | 74.87 | 78.68 |
| ✓ | ✓ | | 83.60 | 75.22 | 78.91 |
| ✓ | ✓ | ✓ | **83.65** | **75.48** | **79.09** |

(b) Domain Generalization.

| Episodic | AC Loss | MOS | P | V | O | D |
|---|---|---|---|---|---|---|
| | | | 96.5 | 76.7 | 83.7 | 61.2 |
| ✓ | | | 96.7 | 79.5 | 84.2 | 61.3 |
| | ✓ | | 96.6 | 78.7 | 84.1 | 61.4 |
| ✓ | ✓ | | 96.7 | 81.0 | 84.4 | **61.6** |
| ✓ | ✓ | ✓ | **96.8** | **81.8** | **84.5** | **61.6** |

*b) Visualization of Image Embeddings:* We randomly select four datasets to analyze the t-SNE plots of image embedding, as shown in Fig. 5. Our MetaPrompt shows better inter-class separability and intra-class cohesiveness in both base and new classes. We attribute the good performance of our method to the fact that the visual prompts are learned under the supervision of textual representations of the unprompted pre-trained model. Since the representations remain fixed throughout the tuning process, visual concepts can be better aligned with corresponding textual labels. With the same class name, image embeddings are usually clustered together. On the other hand, the pre-trained CLIP model has a strong capability of representing semantics. With the supervision of distinguished textual semantics, image embeddings with various classes can be separated.
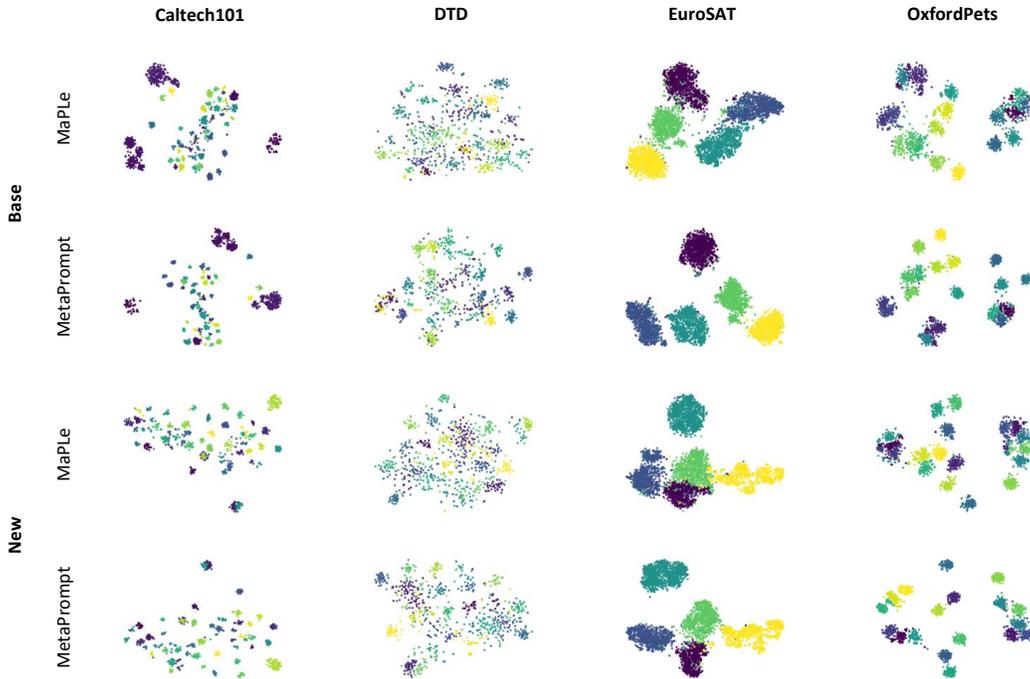
Fig. 5.  T-SNE plots of image embeddings in previous methods MaPLe and our method MetaPrompt on four diverse image recognition datasets. Points with the same color represent image embeddings of the same class. MetaPrompt shows better inter-class separability and intra-class cohesiveness in both base and new classes.

TABLE IV
ABLATION ON DIFFERENT PROMPT LENGTHS.

(a) Base-to-New Generalization.

| Length | Base | New | H |
|--------|------|-----|---|
| 1 | 82.96 | 74.88 | 78.44 |
| 2 | 83.65 | **75.48** | **79.09** |
| 4 | 83.81 | 75.18 | 79.01 |
| 8 | 84.17 | 75.04 | 79.05 |
| 16 | 84.10 | 74.98 | 79.03 |
| 32 | **84.35** | 74.42 | 78.77 |

(b) Domain Generalization.

| Length | P | V | O | D | Average |
|--------|---|---|---|---|---------|
| 1 | 96.61 | 79.91 | 84.38 | 61.45 | 80.59 |
| 2 | 96.69 | 80.28 | 84.52 | 61.53 | 80.76 |
| 4 | 96.80 | **81.84** | **84.54** | 61.63 | **81.20** |
| 8 | 96.92 | 80.49 | 84.35 | 61.50 | 80.82 |
| 16 | **96.93** | 79.75 | 84.49 | 61.60 | 80.69 |
| 32 | 96.87 | 79.28 | 84.48 | **61.67** | 80.57 |

TABLE V
ABLATION ON DIFFERENT LAYERS OF PROMPT.

(a) Base-to-New Generalization.

| Layer | Base | New | H |
|-------|------|-----|---|
| 2 | 80.32 | 74.87 | 77.19 |
| 4 | 81.09 | 74.27 | 77.21 |
| 6 | 81.89 | **75.53** | 78.35 |
| 8 | 82.57 | 75.01 | 78.32 |
| 10 | 83.22 | 75.51 | 78.90 |
| 12 | **83.65** | 75.48 | **79.09** |

(b) Domain Generalization.

| Layer | P | V | O | D | Average |
|-------|---|---|---|---|---------|
| 2 | 96.66 | 81.32 | 83.96 | 60.67 | 80.65 |
| 4 | 96.65 | 81.51 | 84.11 | 60.93 | 80.80 |
| 6 | 96.59 | 81.31 | 84.30 | 61.10 | 80.82 |
| 8 | 96.47 | 80.54 | 84.62 | 61.31 | 80.73 |
| 10 | **96.80** | **81.84** | 84.54 | **61.63** | **81.20** |
| 12 | 96.64 | 78.35 | **84.85** | 61.51 | 80.34 |

*c) Influence of Prompt Length:* The ablation study on the prompt length is carried out in both generalization settings. We study 1, 2, 4, 8, 16, and 32 prompt vectors each layer for both modalities with the same random initialization. Table IV summarizes the performance over both tasks. For base-to-new generalization, we can draw a conclusion that the model with a longer prompt length performs better in base classes. On the other hand, by applying our training strategy, the difference in new classes is relatively small except for 32 prompt vectors with a dramatic drop in performance. The result suggests that using 2 prompt vectors is a better choice when taking into account the accuracy of both base and new classes. For conventional domain generalization, a shorter prompt is not enough to recognize visual concepts well, while a longer prompt seems to overfit in-domain samples. With a prompt length of 4, our method shows promising results considering the overall performance.

*d) Influence of Layer of Prompt :* The ablation study on the layer of prompt is also conducted in both generalization settings. We study 2, 4, 6, 8, 10, and 12 layers of prompts
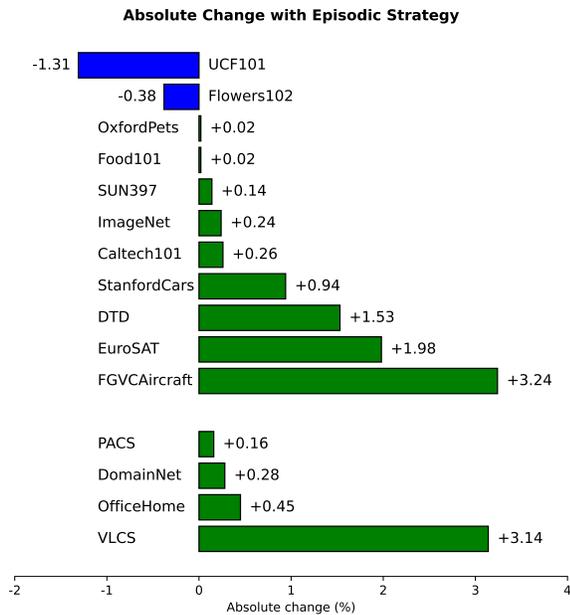
Fig. 6. Performance change on unseen domains with our proposed episodic training over datasets for base-to-new generalization and domain generalization. It shows a general improvement in both generalization tasks, which proves the effectiveness of our episodic training.

for both modalities with the same random initialization. Table V demonstrates the results for these settings. For base-to-new generalization, we can draw a similar conclusion that the accuracy of base classes improves as the layer of prompts increases, while the results on new classes show instability when it changes. The result indicates that applying prompts to all 12 layers has a strong performance when taking into account both base and new classes. For conventional domain generalization, it is clear that with 10 layers of prompt, our method shows excellent performance on all datasets.

*e) Influence of Episodic Training:* We investigate the influence of our proposed batch-wise episodic training strategy. Fig. 6 demonstrates an overall performance boost on datasets for both generalization tasks. The performance of our training strategy remains robust on new classes, which reflects excellent generalization ability. For FGVCAircraft in base-to-new generalization and VLCS in conventional generalization, our training strategy improves the accuracy by more than 3%, which prevents catastrophic failures on out-domain data, highlighting the significance of learning domain invariant prompt.

## VI. CONCLUSION

We propose MetaPrompt for learning domain invariant prompt with CLIP to tackle the problem of generalization. Our theoretical analysis shows that meta-learning applying the episodic training strategy has a strong generalization guarantee. Based on this analysis, we design a dual-modality prompt tuning network with asymmetric contrastive loss and impose a batch-wise episodic training strategy as an explicit constraint on prompt tuning. Prompt can be learned on few-shot data with a high generalization ability to unseen classes

and domains. Extensive experiments on base-to-new generalization and domain generalization demonstrate that our method consistently outperforms existing methods.

While traditional prompt learning approaches often degrade the performance on generalization, our method provides timely insights on how to access the intrinsic association between domains and proposes a feasible solution for learning invariant prompts, which alleviates poor performance on unseen tasks. We reveal that MetaPrompt performs much better in many generalization tasks than input-conditional approaches and show evidence that learning domain invariant prompts has the potential for large pre-trained foundation models. We hope the empirical findings could inspire future research on invariant prompt learning for efficient generalization.

## REFERENCES

[1] W.-L. Chao, S. Changpinyo, B. Gong, and F. Sha, "An empirical study and analysis of generalized zero-shot learning for object recognition in the wild," in *European conference on computer vision*. Springer, 2016, pp. 52–68.

[2] W. Wang, V. W. Zheng, H. Yu, and C. Miao, "A survey of zero-shot learning: Settings, methods, and applications," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 10, no. 2, pp. 1–37, 2019.

[3] Y. Xian, B. Schiele, and Z. Akata, "Zero-shot learning-the good, the bad and the ugly," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4582–4591.

[4] Y. Xian, T. Lorenz, B. Schiele, and Z. Akata, "Feature generating networks for zero-shot learning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 5542–5551.

[5] D. Huynh and E. Elhamifar, "Fine-grained generalized zero-shot learning via dense attribute-based attention," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 4483–4493.

[6] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, M. Ranzato, and T. Mikolov, "Devise: A deep visual-semantic embedding model," *Advances in neural information processing systems*, vol. 26, 2013.

[7] X. Wang, Y. Ye, and A. Gupta, "Zero-shot recognition via semantic embeddings and knowledge graphs," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6857–6866.

[8] M. Kampffmeyer, Y. Chen, X. Liang, H. Wang, Y. Zhang, and E. P. Xing, "Rethinking knowledge graph propagation for zero-shot learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 11 487–11 496.

[9] H. Li, S. J. Pan, S. Wang, and A. C. Kot, "Domain generalization with adversarial feature learning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 5400–5409.

[10] Y. Ganin and V. Lempitsky, "Unsupervised domain adaptation by backpropagation," in *International conference on machine learning*. PMLR, 2015, pp. 1180–1189.

[11] S. Motiian, M. Piccirilli, D. A. Adjeroh, and G. Doretto, "Unified deep supervised domain adaptation and generalization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 5715–5725.

[12] M. Arjovsky, L. Bottou, I. Gulrajani, and D. Lopez-Paz, "Invariant risk minimization," *arXiv preprint arXiv:1907.02893*, 2019.

[13] 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*. PMLR, 2021, pp. 8748–8763.

[14] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[15] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig, "Scaling up visual and vision-language representation learning with noisy text supervision," in *International Conference on Machine Learning*. PMLR, 2021, pp. 4904–4916.

[16] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.

[17] X. L. Li and P. Liang, "Prefix-tuning: Optimizing continuous prompts for generation," *arXiv preprint arXiv:2101.00190*, 2021.

[18] B. Lester, R. Al-Rfou, and N. Constant, "The power of scale for parameter-efficient prompt tuning," *arXiv preprint arXiv:2104.08691*, 2021.

[19] X. Liu, K. Ji, Y. Fu, Z. Du, Z. Yang, and J. Tang, "P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks," *arXiv preprint arXiv:2110.07602*, 2021.

[20] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, "Learning to prompt for vision-language models," *International Journal of Computer Vision*, vol. 130, no. 9, pp. 2337–2348, 2022.

[21] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, "Conditional prompt learning for vision-language models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16 816–16 825.

[22] M. Jia, L. Tang, B.-C. Chen, C. Cardie, S. Belongie, B. Hariharan, and S.-N. Lim, "Visual prompt tuning," *arXiv preprint arXiv:2203.12119*, 2022.

[23] Z. Zheng, X. Yue, K. Wang, and Y. You, "Prompt vision transformer for domain generalization," *arXiv preprint arXiv:2208.08914*, 2022.

[24] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *International conference on machine learning*. PMLR, 2017, pp. 1126–1135.

[25] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra *et al.*, "Matching networks for one shot learning," *Advances in neural information processing systems*, vol. 29, 2016.

[26] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," *Advances in neural information processing systems*, vol. 30, 2017.

[27] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. H. Torr, and T. M. Hospedales, "Learning to compare: Relation network for few-shot learning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1199–1208.

[28] P. Rodríguez, I. Laradji, A. Drouin, and A. Lacoste, "Embedding propagation: Smoother manifold for few-shot classification," in *European Conference on Computer Vision*. Springer, 2020, pp. 121–138.

[29] N. Mishra, M. Rohaninejad, X. Chen, and P. Abbeel, "A simple neural attentive meta-learner," *arXiv preprint arXiv:1707.03141*, 2017.

[30] T. Munkhdalai and H. Yu, "Meta networks," in *International Conference on Machine Learning*. PMLR, 2017, pp. 2554–2563.

[31] B. Oreshkin, P. Rodríguez López, and A. Lacoste, "Tadam: Task dependent adaptive metric for improved few-shot learning," *Advances in neural information processing systems*, vol. 31, 2018.

[32] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap, "Meta-learning with memory-augmented neural networks," in *International conference on machine learning*. PMLR, 2016, pp. 1842–1850.

[33] S. Ravi and H. Larochelle, "Optimization as a model for few-shot learning," in *International conference on learning representations*, 2017.

[34] C. Finn, K. Xu, and S. Levine, "Probabilistic model-agnostic meta-learning," *Advances in neural information processing systems*, vol. 31, 2018.

[35] A. Rajeswaran, C. Finn, S. M. Kakade, and S. Levine, "Meta-learning with implicit gradients," *Advances in neural information processing systems*, vol. 32, 2019.

[36] A. Antoniou, H. Edwards, and A. Storkey, "How to train your maml," *arXiv preprint arXiv:1810.09502*, 2018.

[37] S. Flennerhag, A. A. Rusu, R. Pascanu, F. Visin, H. Yin, and R. Hadsell, "Meta-learning with warped gradient descent," *arXiv preprint arXiv:1909.00025*, 2019.

[38] P. Zhou, X. Yuan, H. Xu, S. Yan, and J. Feng, "Efficient meta learning via minibatch proximal update," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[39] M.-E. Nilsback and A. Zisserman, "Automated flower classification over a large number of classes," in *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*. IEEE, 2008, pp. 722–729.

[40] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi, "Describing textures in the wild," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 3606–3613.

[41] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[42] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[43] J. Chen, X.-M. Wu, Y. Li, Q. Li, L.-M. Zhan, and F.-l. Chung, "A closer look at the training strategy for modern meta-learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 396–406, 2020.

[44] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

[45] 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 *2004 conference on computer vision and pattern recognition workshop*. IEEE, 2004, pp. 178–178.

[46] O. M. Parkhi, A. Vedaldi, A. Zisserman, and C. Jawahar, "Cats and dogs," in *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 2012, pp. 3498–3505.

[47] J. Krause, M. Stark, J. Deng, and L. Fei-Fei, "3d object representations for fine-grained categorization," in *Proceedings of the IEEE international conference on computer vision workshops*, 2013, pp. 554–561.

[48] L. Bossard, M. Guillaumin, and L. V. Gool, "Food-101–mining discriminative components with random forests," in *European conference on computer vision*. Springer, 2014, pp. 446–461.

[49] S. Maji, E. Rahtu, J. Kannala, M. Blaschko, and A. Vedaldi, "Fine-grained visual classification of aircraft," *arXiv preprint arXiv:1306.5151*, 2013.

[50] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba, "Sun database: Large-scale scene recognition from abbey to zoo," in *2010 IEEE computer society conference on computer vision and pattern recognition*. IEEE, 2010, pp. 3485–3492.

[51] K. Soomro, A. R. Zamir, and M. Shah, "Ucf101: A dataset of 101 human actions classes from videos in the wild," *arXiv preprint arXiv:1212.0402*, 2012.

[52] 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*, vol. 12, no. 7, pp. 2217–2226, 2019.

[53] C. Fang, Y. Xu, and D. N. Rockmore, "Unbiased metric learning: On the utilization of multiple datasets and web images for softening bias," in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 1657–1664.

[54] D. Li, Y. Yang, Y.-Z. Song, and T. M. Hospedales, "Deeper, broader and artier domain generalization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 5542–5550.

[55] H. Venkateswara, J. Eusebio, S. Chakraborty, and S. Panchanathan, "Deep hashing network for unsupervised domain adaptation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5018–5027.

[56] X. Peng, Q. Bai, X. Xia, Z. Huang, K. Saenko, and B. Wang, "Moment matching for multi-source domain adaptation," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 1406–1415.

[57] I. Gulrajani and D. Lopez-Paz, "In search of lost domain generalization," *arXiv preprint arXiv:2007.01434*, 2020.

[58] D. Li, Y. Yang, Y.-Z. Song, and T. Hospedales, "Learning to generalize: Meta-learning for domain generalization," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, 2018.

[59] Y. Shi, J. Seely, P. H. Torr, N. Siddharth, A. Hannun, N. Usunier, and G. Synnaeve, "Gradient matching for domain generalization," *arXiv preprint arXiv:2104.09937*, 2021.

[60] B. Sun and K. Saenko, "Deep coral: Correlation alignment for deep domain adaptation," in *European conference on computer vision*. Springer, 2016, pp. 443–450.

[61] J. Cha, S. Chun, K. Lee, H.-C. Cho, S. Park, Y. Lee, and S. Park, "Swad: Domain generalization by seeking flat minima," *Advances in Neural Information Processing Systems*, vol. 34, pp. 22 405–22 418, 2021.

[62] D. Arpit, H. Wang, Y. Zhou, and C. Xiong, "Ensemble of averages: Improving model selection and boosting performance in domain generalization," *arXiv preprint arXiv:2110.10832*, 2021.

[63] Z. Li, K. Ren, X. Jiang, B. Li, H. Zhang, and D. Li, "Domain generalization using pretrained models without fine-tuning," *arXiv preprint arXiv:2203.04600*, 2022.

[64] O. Wiles, S. Gowal, F. Stimberg, S. Alvise-Rebuffi, I. Ktena, T. Cemgil *et al.*, "A fine-grained analysis on distribution shift," *arXiv preprint arXiv:2110.11328*, 2021.

[65] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, "beta-vae: Learning basic visual concepts with a constrained variational framework," in *International conference on learning representations*, 2017.

[66] H. Kim and A. Mnih, "Disentangling by factorising," in *International Conference on Machine Learning*. PMLR, 2018, pp. 2649–2658.

[67] Y. Balaji, S. Sankaranarayanan, and R. Chellappa, "Metareg: Towards domain generalization using meta-regularization," *Advances in neural information processing systems*, vol. 31, 2018.

[68] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[69] M. U. Khattak, H. Rasheed, M. Maaz, S. Khan, and F. S. Khan, "Maple: Multi-modal prompt learning," *arXiv preprint arXiv:2210.03117*, 2022.

[70] Y. Lu, J. Liu, Y. Zhang, Y. Liu, and X. Tian, "Prompt distribution learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5206–5215.

[71] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig, "Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing," *ACM Computing Surveys*, vol. 55, no. 9, pp. 1–35, 2023.

**Yifei Shen** (Graduate Student Member, IEEE) received the Ph.D. degree from the Hong Kong University of Science and Technology and currently works at Microsoft.

**Cairong Zhao** is currently a Professor of College of Electronic and Information Engineering at Tongji University. He received a Ph.D. degree from Nanjing University of Science and Technology, an M.S. degree from Changchun Institute of Optics, Fine Mechanics and Physics, Chinese Academy of Sciences and a B.S. degree from Jilin University, in 2011, 2006 and 2003, respectively. He works on visual and intelligent learning, including computer vision, pattern recognition and visual surveillance. He has published over 40 top-rank international conferences and journals in the field, including CVPR, ICCV, ICLR, AAAI, ACM MM, TIP, TIFS, TMM, TCSVT, and PR. He holds prestigious positions such as the deputy secretary-general of the Pattern Recognition and Machine Intelligence Committee of the Chinese Association of Automation, the chairman of the Computer Vision Special Committee of the Shanghai Computer Society, and an outstanding member of the China Computer Federation, and a senior member of the China Graphics Society. He also serves as the reviewer of more than ten AI-related international journals and conferences, including TPAMI, TIP, CVPR, ICCV, NIPS, ICML, AAAI, etc.

**Kaitao Song** received the B.S. degree and Ph.D. degree in computer science and technology from Nanjing University of Science and Technology, China, in 2015 and 2021. His current research interests include natural language processing, multimodal analysis, deep learning, speech recognition and machine learning. He has published more than 20 academic papers on the top-tier international journals and conferences, such as IEEE TIP, ICML, NeurIPS, ACL, KDD, ICCV, AAAI, IJCAI, InterSpeech, ICASSP, etc. He has served as a PC member of ICML, NeurIPS, ICLR, ACL, EMNLP and etc.

**Yubin Wang** received the B.E. degree in data science and big data technology from Tongji University in 2022. He is currently pursuing the master's degree with Tongji University. His main research interests include prompt learning, multi-modal learning, and person re-identification.

**Dongsheng Li** received B.E. from University of Science and Technology of China in 2007 and Ph.D. from Fudan University in 2012. He is now a principal research manager with Microsoft Research Asia (MSRA) since February 2020. Before joining MSRA, he was a research staff member with IBM Research – China since April 2015. He is also an adjunct professor with School of Computer Science, Fudan University, Shanghai, China. His research interests include recommender systems and machine learning applications. His work on cognitive recommendation engine won the 2018 IBM Corporate Award.

**Xinyang Jiang** received B.E. from Zhejiang University in 2012 and Ph.D. from Zhejiang University in 2017. He is now a researcher from Microsoft Research Asia. Before joining MSRA, he was a researcher from Tencent Youtu Lab. His main research field is computer vision, including person Re-identification, vector graphics recognition and video enhancement and recognition.

**Duoqian Miao** was born in 1964. He is currently a Professor and the Ph.D. Tutor with the College of Electronics and Information Engineering, Tongji University, Shanghai, China. He serves as the Vice President for the International Rough Set Society, the Executive Manager of the Chinese Association for Artificial Intelligence, the Chair of the CAAI Granular Computing Knowledge Discovery Technical Committee, a Distinguished Member of Chinese Computer Federation, the Vice President of the Shanghai Computer Federation, and the Vice President of the Shanghai Association for Artificial Intelligence. He serves as Associate Editor for the International Journal of Approximate Reasoning and an Editor of the Journal of Computer Research and Development (in Chinese).