

TAI++: Text as Image for Multi-Label Image Classification by Co-Learning Transferable Prompt

Xiangyu Wu^{1,2}, Qing-Yuan Jiang¹, Yang Yang^{1*}, Yi-Feng Wu², Qing-Guo Chen², Jianfeng Lu^{1*}

¹Nanjing University of Science and Technology

²Alibaba International Digital Commerce Group

{wxy_yyjhl,yyang,lujf}@njust.edu.cn, qyjiang24@gmail.com,

{yixin.wyf,qingguo.cqg}@alibaba-inc.com

Abstract

The recent introduction of prompt tuning based on pre-trained vision-language models has dramatically improved the performance of multi-label image classification. However, some existing strategies that have been explored still have drawbacks, i.e., either exploiting massive labeled visual data at a high cost or using text data only for text prompt tuning and thus failing to learn the diversity of visual knowledge. Hence, the application scenarios of these methods are limited. In this paper, we propose a pseudo-visual prompt (PVP) module for implicit visual prompt tuning to address this problem. Specifically, we first learn the pseudo-visual prompt for each category, mining diverse visual knowledge by the well-aligned space of pre-trained vision-language models. Then, a co-learning strategy with a dual-adaptor module is designed to transfer visual knowledge from pseudo-visual prompt to text prompt, enhancing their visual representation abilities. Experimental results on VOC2007, MS-COCO, and NUSWIDE datasets demonstrate that our method can surpass state-of-the-art (SOTA) methods across various settings for multi-label image classification tasks. The code is available at <https://github.com/njustkmg/PVP>.

1 Introduction

Multi-label image classification [Wei *et al.*, 2016; Sun *et al.*, 2022; Zhou *et al.*, 2022b; Guo *et al.*, 2023; Mao *et al.*, 2023] has attracted much more attention in many areas including machine learning, computer vision, etc. In recent years, vision-language pre-training models [Radford *et al.*, 2021; Jia *et al.*, 2021; Xu *et al.*, 2023; Bowman *et al.*, 2023; Li *et al.*, 2023; Fu *et al.*, 2024] have exhibited remarkable generalization capabilities by aligning visual and language modalities to shared embedding space. Through utilizing the aligned vision-language embeddings from pre-trained vision-language models, prompt tuning [Zhou *et al.*, 2022a; Gu *et al.*, 2022; Wang *et al.*, 2023b; Sun *et al.*, 2023] has

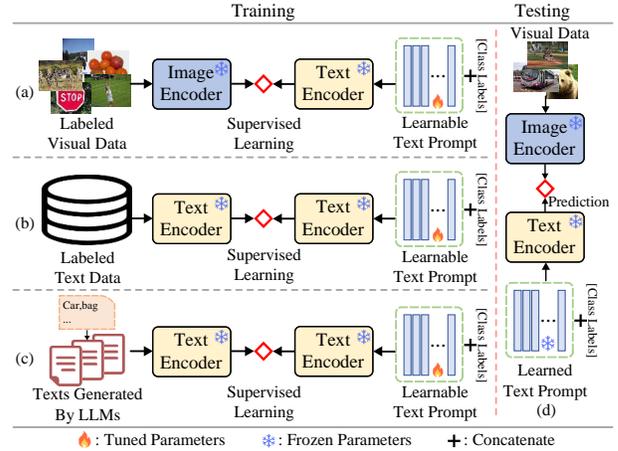


Figure 1: Different prompt tuning paradigms for multi-label image recognition. (a). Tuning text prompt with labeled visual data. (b). Tuning text prompt with labeled text data. (c). Tuning text prompt with texts generated by LLMs. (d). Testing with images and text prompt for image recognition.

emerged as a novel paradigm and significantly enhanced performance in multi-label image classification [Zhou *et al.*, 2022b; Sun *et al.*, 2022; Guo *et al.*, 2023].

As a pioneering work, CoOp [Zhou *et al.*, 2022b], which is illustrated in Figure 1 (a) and Figure 1 (d), ingeniously designed learnable text prompt combined with textual class labels, aligned with images in a supervised manner via frozen encoders. CoCoOp [Zhou *et al.*, 2022a] explored a conditional context of image strategy for unseen classes. To further leverage aligned space of origin CLIP [Radford *et al.*, 2021], DualCoOp [Sun *et al.*, 2022] encoded dual (positive and negative) prompts with partial-labeled images. DualCoOp++ [Hu *et al.*, 2023] separately encoded evidential, positive, and negative prompts to further improve the performance. All of these methods necessitate labeled visual data for model training. However, constructing adequate labeled visual data is costly. And insufficient training data might hinder the learning of robust image recognition networks. Hence, the application scenarios of these methods are limited.

Considering that the visual and text modalities are well-aligned by pre-trained vision-language models, such as CLIP,

*Corresponding author

TAI-DPT [Guo *et al.*, 2023] proposed to enable low-cost text data instead of labeled visual data to learn text prompt. In other words, TAI-DPT leveraged publicly labeled text data for tuning text prompt and directly classified images during testing. The training and testing procedure of TAI-DPT is depicted in Figure 1 (b) and Figure 1 (d), respectively. Apparently, TAI-DPT failed to learn the diversity of visual knowledge because TAI-DPT only utilized text data and textual class labels for text prompt tuning.

There also exist some more challenging scenarios compared with those mentioned above. A typical scenario is that only labeled text data of common categories is available. In this scenario, a reliable option is to utilize the large language models (LLMs) [OpenAI, 2023; Zeng *et al.*, 2023] to generate sufficient pseudo text data for prompt tuning. For example, Guangdong-Hong Kong-Macao International Algorithm Competition¹ advises using LLMs and textual class labels to generate pseudo texts. We can use generated pseudo text data and adopt the same strategy as TAI-DPT, which only uses text data to perform prompt tuning. The training and testing procedure is illustrated in Figure 1 (c) and Figure 1 (d), respectively. However, for multi-label image classification, learning a wide range of discriminative visual-level information for each category is essential and indispensable due to the heterogeneity between images and texts. For example, different images with the same categories, like cars, planes, or bags, encompass various shapes and attributes. To sum up, existing prompt tuning based multi-label classification methods either require a large amount of labeled visual data or fail to learn the diversity of visual knowledge if the algorithm solely adopts text data to perform text prompt tuning.

How to reconcile this problem? In this paper, we design a pseudo-visual prompt module for visual prompt tuning implicitly, mining diverse visual knowledge and explicitly avoiding the usage of massive labeled visual data. More concretely, we propose a co-learning method for both visual and text prompt tuning, where the visual prompt tuning is performed implicitly based on the pseudo-visual prompt module. As visual and text modalities are well-aligned by pre-trained models, we can learn diverse visual knowledge from aligned space of CLIP instead of using massive labeled visual data. Then, we co-learn visual and text prompts with a dual-adaptor module by transferring visual knowledge from learned pseudo-visual prompt to text prompt. Our contributions are summarized as follows:

- We design a pseudo-visual prompt module to perform visual prompt tuning implicitly and propose a novel transferable prompt co-learning method for both visual and text prompt tuning for multi-label image classification. Furthermore, our pseudo-visual prompt can be easily combined with existing methods to improve multi-label image classification performance further.
- We co-learn visual and text prompts by leveraging a dual-adaptor module and contrastive learning for transferring visual knowledge to text prompt, thereby enhancing their visual representation capabilities.

¹<https://iacc.pazhoulab-huangpu.com/>

- Extensive experiments on three widely-used benchmarks, i.e., VOC2007, MS-COCO, and NUSWIDE, show that our proposed method can outperform SOTA methods for multi-label image classification tasks.

2 Related Work

Multi-Label Image Classification. Given an image input, multi-label image classification [Wang *et al.*, 2016; Alfassy *et al.*, 2019; Yang *et al.*, 2018; Zhou *et al.*, 2022b; Lin, 2023; Xi *et al.*, 2023] tries to recognize all the target class labels. Early works [Wang *et al.*, 2016; Wang *et al.*, 2017] focus on utilizing labeled visual data to learn intra-class relationships. However, these methods usually cannot achieve satisfactory performance when labeled visual data is limited. Recent works [Lee *et al.*, 2018; Alfassy *et al.*, 2019; Yang *et al.*, 2021; Simon *et al.*, 2022; Wang *et al.*, 2023a] pay more attention to the data-limited scenarios, including the scenarios with zero-shot, few-shot and partial-labeled image data, in past years.

The partial-labeled data is defined as that some class labels of data are unknown. The authors of [Lee *et al.*, 2018] utilized knowledge graph to mine the relationship between different class labels. SARB [Pu *et al.*, 2022] tried to learn instance-level and prototype-level semantic representations to complement unknown labels. Zero/few-shot multi-label image classification tries to learn novel class labels from limited data. LaSO [Alfassy *et al.*, 2019] leveraged relationships among label sets to extract underlying semantic information for few-shot image classification. In paper [Chen *et al.*, 2023], the authors introduced an embedding matrix with principal embedding vectors trained using a tailored loss function.

Prompt Tuning in Multi-Modal Learning. Prompt tuning [Min *et al.*, 2022; Yan *et al.*, 2023; Yang *et al.*, 2023; Bowman *et al.*, 2023; Yu *et al.*, 2023] has emerged as a promising technique in computer vision and natural language processing, offering a parameter-efficient way to leverage large pre-trained models. KnowPrompt [Chen *et al.*, 2022] involved injecting knowledge into a prompt template and encoding rich semantic knowledge among entities and relations. Pro-Tuning [Nie *et al.*, 2022] learned task-specific visual prompt for downstream input images while keeping the pre-trained model frozen.

Amidst the progress in multi-modal pre-training, researchers have explored the application of prompt tuning in the multi-modal domain. CoOp [Zhou *et al.*, 2022b] modified the pre-trained vision-language models for image recognition tasks by employing learnable prompt context vectors. DualCoOp++ [Hu *et al.*, 2023] efficiently adapted a powerful vision-language model with partial-labeled images by introducing evidence-guided region feature aggregation and winner-take-all modules to improve spatial aggregation and inter-class interaction. These methods necessitated visual modality and textual class labels as default prerequisites in both training and testing. Consequently, TAI-DPT [Guo *et al.*, 2023] extended this paradigm by treating texts as images for zero-shot image recognition, storing the aligned vision-language information from CLIP into text prompt without seeing any image during training.

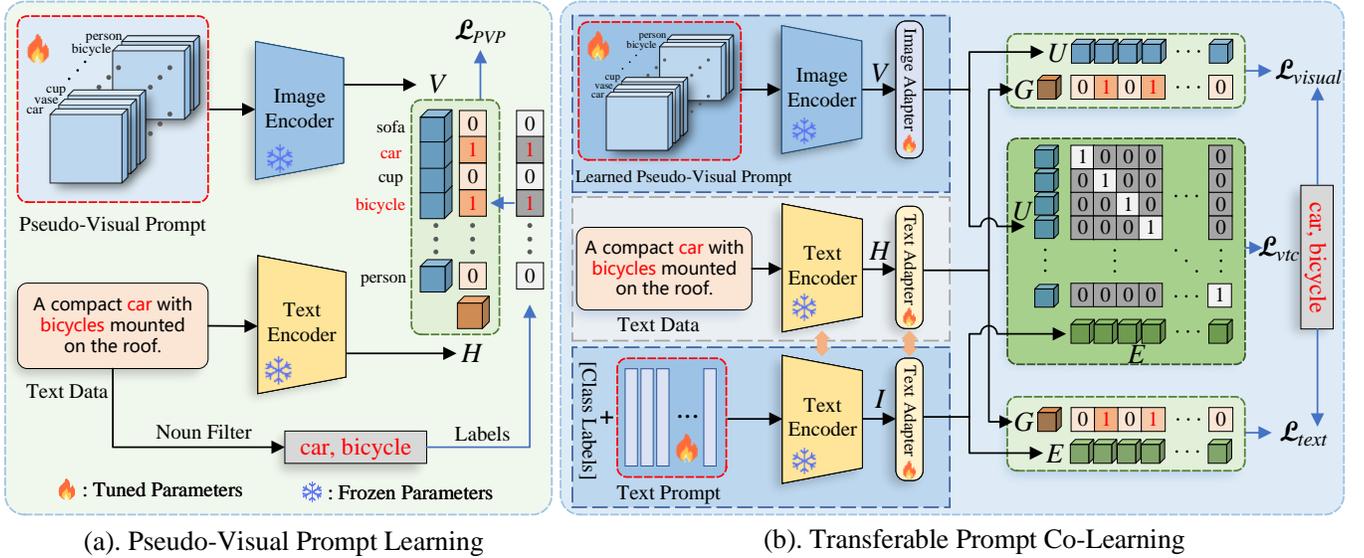


Figure 2: Pseudo-Visual Prompt Learning and Transferable Prompt Co-Learning. Sub-Figure (a) presents the class-specific pseudo-visual prompt module. The global text embedding and pseudo-visual prompt embedding are obtained from the frozen CLIP image and text encoders. The corresponding cosine similarity between the embeddings is guided by the noun-filtered labels with ranking loss. Sub-Figure (b) presents the transferable prompt co-learning module. We perform contrastive learning between the pseudo-visual prompt and the text prompt to enhance the prompts’ visual diversity representation capability.

Prompt tuning based methods can boost the performance of multi-label image classification. However, these methods either require a large amount of labeled visual data or fail to learn the diversity of visual knowledge. In this paper, we propose a novel transferable prompt co-learning method to solve this problem by designing a pseudo-visual prompt module.

3 Methodology

In this section, we introduce the architecture of our proposed pseudo-visual prompt (PVP) method in detail. The whole architecture of our proposed method is illustrated in Figure 2. Our proposed method comprises two key phases: the pseudo-visual prompt learning phase and the transferable prompt co-learning phase. During training, these two phases will be performed sequentially. In practice, we observed that direct co-learning pseudo-visual prompt and text prompt, without pre-learning pseudo-visual prompt, will lead to difficulty in convergence of pseudo-visual prompt. The reason is that the PVP is initialized randomly without textual class labels while the text prompt combines textual class labels. Hence, we adopt the two-stage learning strategy for our proposed method. Furthermore, we introduce two training text construction strategies for the scenarios of labeled visual data and unavailable text data respectively.

3.1 Pseudo-Visual Prompt Learning

The accurate learning of diverse and comprehensive visual knowledge for each class label is pivotal for image classification tasks. Text prompt, combined with textual class labels, can only capture visual information aligned with the textual class labels from the aligned space. Hence, we propose a pseudo-visual prompt module, aiming to construct a class-specific visual prompt for each category and leverage the im-

age encoder, text encoder, and aligned space of CLIP to optimize the pseudo-visual prompt, learning the generic visual knowledge for each category.

As shown in Figure 2 (a), we innovatively design the class-specific pseudo-visual prompt for each category without combining them with any explicit visual or textual labels. Image modality encompasses a wide range of diversity, which has been learned in the well-aligned space of CLIP. Hence, the class-specific pseudo-visual prompt can learn and store the unique visual knowledge of each category. More concretely, the pseudo-visual prompt is defined as:

$$\mathbf{P} = [\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3, \dots, \mathbf{p}_N], \quad (1)$$

where $\mathbf{p}_i \in \mathbb{R}^{H \times W \times 3}$ denotes class-specific pseudo-visual prompt for i -th class, H and W are equal in size, N is the number of class labels. Note that the number of pseudo-visual prompt is batch-size agnostic and equals the number of target categories.

Then, we utilize the well-aligned space of CLIP, collected text training data, and frozen image/text encoder to learn pseudo-visual prompt. The learning procedure of pseudo-visual prompt can be formalized as follows:

$$\langle \Omega_{aligned}, \mathbf{T}, \phi, \psi \rangle \rightarrow \langle \mathbf{P} \rangle,$$

where \mathbf{T} represents the collected labeled text data or pseudo text data generation by LLMs, $\Omega_{aligned}$ represents the origin CLIP’s aligned shared space, $\phi(\cdot)$ and $\psi(\cdot)$ refer to the frozen text and image encoders of CLIP, respectively. As for the input text \mathbf{T} , we directly follow the origin CLIP to obtain the global text embedding by projecting the feature of the last “<EOS>” token. The global visual embedding for each category of \mathbf{P} is obtained by visual attention pooling. Hence, we have:

$$\mathbf{H} = \phi(\mathbf{T}), \quad \mathbf{V} = \psi(\mathbf{P}),$$

where $\mathbf{H} \in \mathbb{R}^{B \times D}$ denotes the extracted normalized global text embeddings of a batch, and $\mathbf{V} \in \mathbb{N}^{N \times D}$ denotes the normalized global visual embeddings of N pseudo-visual prompts. For a specific text $\mathbf{t}_i \in \mathbf{T}$ in a batch, the similarity of text \mathbf{t}_i and pseudo-visual prompt can be computed by:

$$s_{ij} = \langle \mathbf{h}_i, \mathbf{v}_j \rangle, \quad \forall j \in \{1, 2, 3, \dots, N\}. \quad (2)$$

Here, $\mathbf{h}_i \in \mathbf{H}$, $\mathbf{v}_j \in \mathbf{V}$ denote the global text embedding of text \mathbf{t}_i and the global visual embedding of j -th pseudo-visual prompt, respectively. We then perform noun filtering to obtain the positive and negative labels. Specifically, given a text embedding \mathbf{h}_i and a pseudo-visual prompt embeddings \mathbf{v}_j , if the class label filtered from input text by noun filtering is contained in the target category set, \mathbf{h}_i and \mathbf{v}_j are positive pair. Otherwise, they are negative pairs. We employ the ranking loss to measure the discrepancy between similarity scores and text labels following the setting of TAI-DPT [Guo *et al.*, 2023]:

$$\mathcal{L}_{PVP} = \frac{1}{B} \sum_{k=1}^B \sum_{i \in \{c^+\}} \sum_{j \in \{c^-\}} \max(0, m - s_{ki} + s_{kj}), \quad (3)$$

where c^+ and c^- are positive labels and negative labels, s_{ki} and s_{kj} are positive pair and negative pair similarities described in Equation (2), m is the margin used to measure the difference between each pair of predicted values. During training, we fix the text encoder and image encoder and only learn the pseudo-visual prompt by optimizing the objective function in Equation (3).

3.2 Transferable Prompt Co-Learning

After the first learning phase, the diverse class-specific visual knowledge is well aligned with class labels and stored in pseudo-visual prompt. Furthermore, inspired by TAI-DPT [Guo *et al.*, 2023], we design contrastive loss and a dual-adapter module to co-learn the pseudo-visual prompt and text prompt by transferring visual information to the text prompt. Here, the pseudo-visual prompt is initialized by the learned pseudo-visual prompt in the first phase.

Figure 2 (b) illustrates the transferable prompt co-learning procedure. We first adopt the same definition of pseudo-visual prompt \mathbf{P} from Equation (1) during this phase. Then, we define the text prompt as follows:

$$\forall i \in \{1, 2, \dots, N\}, \mathbf{R}_i = [\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_M, \mathbf{g}_i], \\ \mathbf{S} = [\mathbf{R}_1, \mathbf{R}_2, \dots, \mathbf{R}_N],$$

where \mathbf{g}_i denotes word embedding of the i -th class label, \mathbf{r}_i is a learnable context vector of text prompt and M denotes the number of text prompt.

Then, we utilize image and text encoders from CLIP to obtain pseudo-visual prompt, text prompt and global text embeddings.

$$\mathbf{V} = \psi(\mathbf{P}), \mathbf{H} = \phi(\mathbf{T}), \mathbf{I} = \phi(\mathbf{S}).$$

To mine the knowledge of origin CLIP and downstream task, we apply an identical adapter module [Gao *et al.*, 2021] for image encoder and text encoder, a.k.a., a dual-adapter module. Both image and text adapter consist of two fully connected layers, an activation function, and a residual connection. We apply an image adapter for pseudo-visual prompt

and a text adapter for both text prompt and global text embedding extraction. Then, we have:

$$\mathbf{U} = (1 - \lambda)g(\mathbf{V}) + \lambda\mathbf{V}, \\ \mathbf{G} = (1 - \lambda)h(\mathbf{H}) + \lambda\mathbf{H}, \\ \mathbf{E} = (1 - \lambda)h(\mathbf{I}) + \lambda\mathbf{I}. \quad (4)$$

Here, $g(\cdot)$ and $h(\cdot)$ denote the adapter functions for image and text modality, respectively. And $\lambda \in [0, 1]$ denotes the weight between feature from adapter module and feature from image/text encoder. Please note that text prompts and global text utilize the same adapter in Equation (4), which means text adapter is a parameter-sharing network for text prompt and global text learning.

After embedding extraction, we design the objective function for pseudo-visual prompt and text prompt co-learning. We present the objective function based on the given pseudo-visual prompt, text prompt and global text embeddings. Specifically, we first utilize contrastive loss to preserve the similarity between pseudo-visual prompt and text prompt of N class labels. The similarity matrix can be obtained by: $\mathbf{UE}^T \in \mathbb{R}^{N \times N}$. And the ground truth for pseudo-visual prompt and text prompt of N class labels is an identity matrix. Note that the size of similarity matrix is batch-size agnostic and equals the number of target categories N . Thus, the contrastive loss can be written as:

$$p_{ij}^{v2t} = \frac{\exp(s(\mathbf{u}_i, \mathbf{e}_j)/\tau)}{\sum_{k=1}^N \exp(s(\mathbf{u}_i, \mathbf{e}_k)/\tau)}, \\ p_{ij}^{t2v} = \frac{\exp(s(\mathbf{u}_i, \mathbf{e}_j)/\tau)}{\sum_{k=1}^N \exp(s(\mathbf{u}_k, \mathbf{e}_j)/\tau)}, \\ \mathcal{L}_{vtc} = \frac{1}{2} [l_{CE}(y^{v2t}, p^{v2t}) + l_{CE}(y^{t2v}, p^{t2v})],$$

where p_{ij}^{v2t} and p_{ij}^{t2v} denote the softmax-normalized similarity from pseudo-visual prompt to text prompt and from text prompt to pseudo-visual prompt, respectively. τ denotes the temperature scale parameter. l_{CE} denotes cross-entropy loss. $\forall i, j \in \{1, 2, \dots, N\}$, $y_{ij}^{v2t}, y_{ij}^{t2v} \in \{0, 1\}$ denote the similarity ground-truth of text prompt and pseudo-visual prompt. $y_{ij}^{v2t} = 1$ if text prompt \mathbf{t}_i and pseudo-visual prompt \mathbf{p}_j belong to the same category, and $y_{ij}^{v2t} = 0$ otherwise. The definition of y_{ij}^{t2v} is the same with y_{ij}^{v2t} .

Moreover, we utilize ranking loss to obtain \mathcal{L}_{visual} and \mathcal{L}_{text} similar to Equation (3). Specifically, we use \mathcal{L}_{visual} to measure the disparity between global text and pseudo-visual prompt. The \mathcal{L}_{visual} can be formulated as follows:

$$s_{ij}^v = \langle \mathbf{g}_i, \mathbf{u}_j \rangle, \quad \forall j \in \{1, 2, 3, \dots, N\}, \\ \mathcal{L}_{visual} = \frac{1}{B} \sum_{k=1}^B \sum_{i \in \{c^+\}} \sum_{j \in \{c^-\}} \max(0, m - s_{ki}^v + s_{kj}^v).$$

Similarly, the \mathcal{L}_{text} that measures the disparity between global text and text prompt can be formulated as follows:

$$s_{ij}^t = \langle \mathbf{g}_i, \mathbf{e}_j \rangle, \quad \forall j \in \{1, 2, 3, \dots, N\}, \\ \mathcal{L}_{text} = \frac{1}{B} \sum_{k=1}^B \sum_{i \in \{c^+\}} \sum_{j \in \{c^-\}} \max(0, m - s_{ki}^t + s_{kj}^t).$$

Finally, we get the total training loss by combining \mathcal{L}_{vtc} , \mathcal{L}_{visual} and \mathcal{L}_{text} :

$$\mathcal{L} = \mathcal{L}_{vtc} + \mathcal{L}_{visual} + \mathcal{L}_{text}. \quad (5)$$

During training procedure, we fix the image and text encoder from CLIP and learn pseudo-visual prompt, text prompt and dual-adapter by optimizing the objective function in Equation (5).

3.3 Training Text Data Construction

In this section, we discuss the training text data construction for different application scenarios. To obtain training text data in this paper, we utilize two strategies, i.e., human-annotated labeled text data construction [Guo *et al.*, 2023] and LLMs-based pseudo text data construction. The first strategy was introduced by TAI-DPT [Guo *et al.*, 2023], and we follow the setting of TAI-DPT. Specifically, we directly employ public object detection datasets like MS-COCO [Lin *et al.*, 2014] and localized narratives [Krasin *et al.*, 2017] to form labeled text data. For the second strategy, pseudo text data is generated using constructed templates and LLMs [OpenAI, 2023; Zeng *et al.*, 2023] for automatic generation. Concretely, we first combine several class labels with a query template to construct a query prompt. Then, we utilize LLMs to generate pseudo text data. We provide a query prompt example as follows:

PROMPT: Make a sentence to describe a photo. Requirements: Each sentence should be less than 15 words and include keywords: car, dog, cat.

To filter out the unreliable texts generated by LLMs, we re-input the generated text data into LLMs with another query template:

PROMPT: Will the scene described in this text appear in reality? Scene: + "{text}".

Moreover, we judge the reasonableness of the text through the output *likely/unlikely*. For word-level filtered labels in input text, we follow the setting of TAI-DPT [Guo *et al.*, 2023] using NLTK² to perform noun filtering. More details of the query prompts, examples and noun filtering are provided in the supplementary materials.

3.4 Model Inference

During testing procedure, we first replace the input text data with testing images and utilize the image encoder of CLIP to obtain the image embeddings. Then, we utilize image embeddings to compute visual and textual cosine similarities with class label embedding generated by the pseudo-visual prompt and the text prompts, respectively. The final classification score is obtained by fusing the visual and textual cosine similarity.

4 Experiments

4.1 Experiment Setup

Datasets. We evaluate our proposed PVP on VOC2007 [Everingham *et al.*, 2010], MS-COCO [Lin *et al.*, 2014], and NUSWIDE [Chua *et al.*, 2009] datasets. The VOC2007

dataset contains 20 categories with 5,011 images for training and 4,952 images for testing. The MS-COCO dataset contains 80 categories divided into training, testing, and validation sets. We use its training set (82,081 images) for training and validation set (40,504 images) for testing in our experiments because the class labels of origin testing are unavailable. The NUSWIDE dataset contains 81 categories with 161,789 training images and 107,859 testing images to validate our method.

To construct the labeled training text data, we follow the setting of TAI-DPT [Guo *et al.*, 2023]. More concretely, we extract 100K coco-captions as the text training data for VOC2007 and MS-COCO, and localized narratives from OpenImages [Krasin *et al.*, 2017] for NUSWIDE. For the pseudo text data generated by LLMs, we adopt ChatGLM [Zeng *et al.*, 2023] to generate 500k pseudo texts for pseudo-visual prompt learning, and 200k for transferable prompt co-learning. Then we use an identical query prompt and the different number of categories for different datasets, i.e., 20 categories for VOC2007 dataset, 80 categories for MS-COCO dataset, and 81 categories for NUSWIDE dataset.

Implementation Details. For fair comparison, we follow the setting of TAI-DPT [Guo *et al.*, 2023] to choose CLIP ResNet-50 [He *et al.*, 2016] as image encoder and the corresponding CLIP Transformer [Vaswani *et al.*, 2017] as text encoder. We adopt SGD algorithm to perform prompt learning for both two phases. The margin is set to be 1 for the ranking loss in both two phases.

For the pseudo-visual prompt learning, we initialize an identical class-specific visual prompt of size $224 \times 224 \times 3$ for each category. All to be learned prompts are randomly initialized by the same mean and standard value, i.e., $mean = 0, std = 0.02$. To perform pseudo-visual prompt learning, both visual and text encoders are frozen, and only prompts are optimized. The training epoch is set to be 40 for all datasets. The learning rate is empirically initialized with 0.1 and decayed through cosine annealing during training.

For the transferable prompt co-learning, we use the learned PVP in the first phase to initialize the pseudo-visual prompt. We follow the setting of TAI-DPT [Guo *et al.*, 2023] to initialize text prompt by randomly sampling from a Gaussian distribution with mean being 0 and variance being 0.02, and the length of text prompts is set to 16. In this phase, the image and text encoders are also frozen while the pseudo-visual prompt, text prompt, and dual-adapter module are optimized. The hyper-parameter τ is set to be 0.02. The training epoch is set to 20 for all datasets. The learning rate is set to be $1e-4$ and $1e-6$ for text prompt and pseudo-visual prompt learning and decay by cosine annealing, respectively. The hyper-parameter λ of dual-adapter module is set to be 0.5.

4.2 Comparison with SOTA Methods

Results on Zero-Shot Task. To validate the effectiveness of our proposed pseudo-visual prompt, we compare its performance with zero-shot CLIP (ZSCLIP) [Radford *et al.*, 2021] and the current SOTA method³ TAI-DPT [Guo *et al.*, 2023]

³Comparison results with recent work TAI-Adapter [Zhu *et al.*, 2023] submitted to arXiv are reported in supplementary materials.

²<https://www.nltk.org/>

Table 1: The mAP results for zero-shot setting on all datasets. The best performance is shown in boldface. (480) denotes the image resolution during inference

Method	VOC2007		MS-COCO		NUSWIDE	
	Label	Pseudo	Label	Pseudo	Label	Pseudo
ZSCLIP	77.3		49.7		37.4	
TAI-DPT	88.3	88.1	65.1	64.6	46.5	47.3
PVP	88.6	88.9	67.7	67.5	47.6	49.3
TAI-DPT (480)	88.3	88.4	67.2	66.6	42.9	44.1
PVP (480)	89.7	90.0	70.9	70.8	44.3	46.0

Table 2: The mAP results for comparison with few-shot methods on MS-COCO dataset.

Method	0-shot	1-shot	5-shot
LaSO	-	45.3	58.1
ML-FSL	-	54.4	63.6
KGGR	-	52.3	63.5
NLC	-	56.8	64.8
TAI-DPT	59.2	-	-
PVP	62.1	-	-
PVP (480)	64.4	-	-

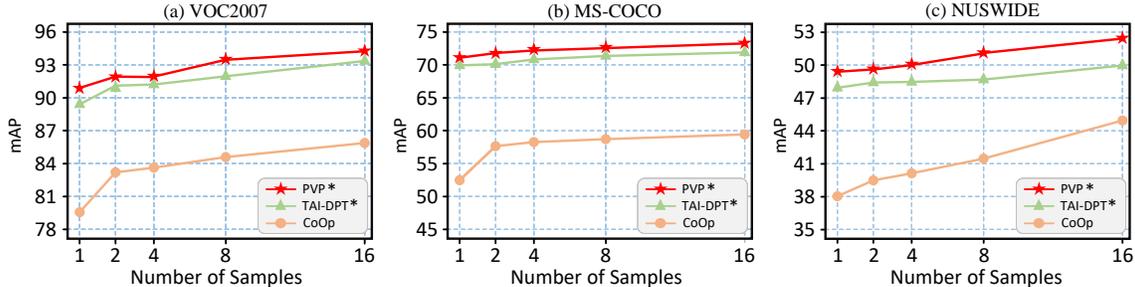


Figure 3: Results for few-shot setting, where the performance of PVP*/TAI-DPT* integrate the predictions of PVP/TAI-DPT and CoOp.

on all datasets. Table 1 presents the results of our proposed method and baselines on all datasets. In Table 1, “Label” and “Pseudo” denote the training with labeled text data (e.g. coco-caption [Lin *et al.*, 2014], localized narratives [Krasin *et al.*, 2017]) and pseudo text data generated by LLMs, respectively. “(480)” denotes that the image resolution is set to be 480×480 during inference.

From Table 1, we can see that our proposed method can achieve the best zero-shot multi-label image recognition performance on all three datasets in all cases. On MS-COCO dataset trained with label text and pseudo text data, our method significantly outperforms TAI-DPT by a large margin of 3.7% and 4.2% points respectively on 480×480 resolution. On VOC2007 and NUSWIDE datasets, our method also improves by 1~2% points over TAI-DPT on both labeled and pseudo text data. The results demonstrate that our method is without relying on any labeled visual and text data, making it more valuable for real-world application scenarios. Notably, we found that on NUSWIDE dataset, all methods perform better on pseudo data than labeled data. This might be due to the texts generated by the LLMs being closer to the original training data of CLIP than the localized narrative [Krasin *et al.*, 2017]. Moreover, the performance on NUSWIDE dataset with higher image resolution is worse than that with lower image resolution. This is due to the resolution of testing images are smaller (about 200×200) than 480.

Results on Few-Shot Task. For few-shot task, we select LaSO [Alfassy *et al.*, 2019], ML-FSL [Misra *et al.*, 2016] and TAI-DPT as baselines and compare the performance of these methods. The LaSO and ML-FSL are used for few-shot setting. Meanwhile, TAI-DPT and our proposed method are used for zero-shot setting. LaSO [Alfassy *et al.*, 2019] and ML-FSL [Misra *et al.*, 2016] require labeled visual data for training. Hence, for the few-shot setting, we select 64 cat-

egories in MS-COCO as normal classes, and the remaining 16 (bicycle, boat, stop sign, bird, backpack, frisbee, snowboard, surfboard, cup, fork, spoon, broccoli, chair, keyboard, microwave, and vase) as novel classes. For few-shot task, the results are reported in Table 2. From Table 2, we can find that our method outperforms TAI-DPT by a large margin of 5.2% points. Furthermore, our proposed method also surpasses ML-FSL trained on 5-shot samples by 0.8% points.

Furthermore, following the same setting of TAI-DPT [Guo *et al.*, 2023], we randomly sample 1, 2, 4, 8, and 16 samples for each class to train the model, and integrate the predictions conveniently with CoOp [Zhou *et al.*, 2022b]. The results are reported in Figure 3. In Figure 3, “PVP*” denotes that the performance is calculated by integrating the predictions of CoOp and our PVP when testing. The definition of “TAI-DPT*” is similar to “PVP*”. From Figure 3, we can find that our proposed method achieves the best performance on various few-shot settings on all datasets without seeing any labeled visual data.

Results on Partial-Label Task. For partial-label task, we select SARb [Pu *et al.*, 2022] and DualCoOp [Sun *et al.*, 2022] as baselines. Following the setting of TAI-DPT [Guo *et al.*, 2023], we use visual training sets with different proportions to complete the training of the SARb and DualCoOp methods. For our proposed method and TAI-DPT, we integrate the predictions of these methods with DualCoOp. The results are shown in Table 3. We can see that PVP* can achieve higher performances than TAI-DPT* in most cases on all datasets after further integrating PVP and TAI-DPT with DualCoOp.

4.3 Visualization

To validate the effectiveness of our proposed method, we conduct several visualization experiments, demonstrating that pseudo-visual prompt are better at focusing on diverse visual

Table 3: The mAP results for partial-label setting on all datasets, where the performance of PVP*/TAI-DPT* integrates the predictions of PVP/TAI-DPT and DualCoOp. The best performance is shown in boldface.

Datasets	Method	10%	20%	30%	40%	50%	60%	70%	80%	90%	Avg.
VOC2007	SARB	83.5	88.6	90.7	91.4	91.9	92.2	92.6	92.8	92.9	90.7
	DualCoOp	91.4	93.8	93.8	94.3	94.6	94.7	94.8	94.9	94.9	94.1
	TAI-DPT*	93.3	94.6	94.8	94.9	95.1	95.0	95.1	95.3	95.5	94.8
	PVP*	93.7	94.4	94.7	95.1	95.1	95.2	95.2	95.3	95.3	94.9
MS-COCO	SARB	71.2	75.0	77.1	78.3	78.9	79.6	79.8	80.5	80.5	77.9
	DualCoOp	81.0	82.3	82.9	83.4	83.5	83.9	84.0	84.1	84.3	83.3
	TAI-DPT*	81.5	82.6	83.3	83.7	83.9	84.0	84.2	84.4	84.5	83.6
	PVP*	81.8	82.8	83.3	83.6	83.9	84.1	84.3	84.6	84.8	83.7
NUSWIDE	DualCoOp	54.0	56.2	56.9	57.4	57.9	57.9	57.6	58.2	58.8	57.2
	TAI-DPT*	56.4	57.9	57.8	58.1	58.5	58.8	58.6	59.1	59.4	58.3
	PVP*	56.9	58.4	58.9	59.3	59.5	59.7	59.9	60.1	60.2	59.2

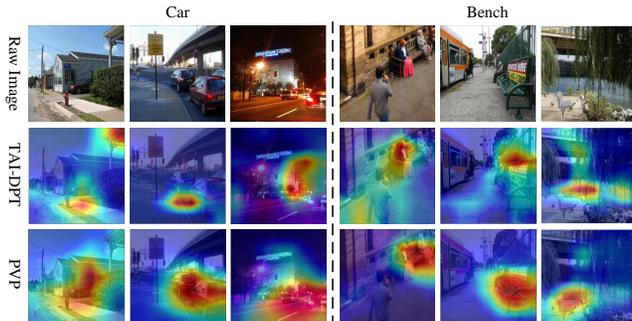


Figure 4: Visualization of PVP and TAI-DPT methods.

information. Specifically, we first select some common class labels, such as car and bench. Then, we randomly select several images of these class labels with different attributes, such as shape, color, size, etc. We then visualize the correlation between the class prompt embeddings of the PVP and TAI-DPT methods and the local image features. The results are shown in Figure 4. From Figure 4, we can find that compared with TAI-DPT, PVP can learn extensive features with higher relevance to the class label in different scenes where the object has different shapes or different attributes, and can accurately identify the position of the object even in dark light, bright light or occlusion scenarios. Hence, we demonstrate that our method can learn diverse and comprehensive visual knowledge for each category through pseudo-visual prompt. More visualization results are provided in the supplementary materials.

4.4 Further Analysis

Ablation Study. To evaluate the effectiveness of our methods, we study the influence of different components, including pseudo-visual prompt, dual-adapter, and contrastive learning (respectively abbreviated as PVP, DA, cLoss in Table 4). From Table 4, we can observe that all components can improve the performance and the PVP module can boost the most significant improvement.

Sensitivity to Hyper-Parameter. In addition, we further explore the impact of the quantity of text training data on the performance of our method on MS-COCO dataset. To eliminate the influence of different data sources, we mix labeled

Table 4: Ablation study for our proposed method.

PVP	DA	cLoss	VOC2007	MS-COCO	NUSWIDE
×	×	×	88.1	64.2	47.0
✓	×	×	89.1	68.6	48.8
✓	✓	×	89.4	69.7	48.9
✓	×	✓	89.7	69.9	48.9
✓	✓	✓	90.0	70.8	49.3

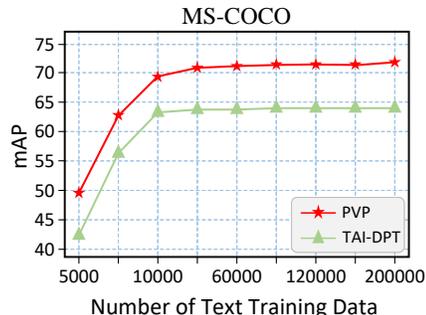


Figure 5: The mAP value with different number of text data on MS-COCO dataset.

text and pseudo text data generated by the LLMs and randomly sample different quantities of text as the training set, as shown in Figure 5. From Figure 5, we can see that as the number of sampled text data increases, the mAP result increases at the beginning and then remains unchanged. In our experiment, we set the number of text data as 200K. More experimental results can be found in supplementary materials.

5 Conclusion

In this paper, we design a novel pseudo-visual prompt module based on pre-trained vision-language models for multi-label image classification tasks. Thus, we can learn diverse visual knowledge from aligned space of CLIP instead of using massive labeled visual data. By leveraging a contrastive loss and dual-adapter module to co-learn the visual and text prompts, our proposed method can enhance the visual representation capabilities. Experiments verify that our PVP method can achieve the best performance compared with the SOTA methods across various datasets.

References

- [Alfassy *et al.*, 2019] Amit Alfassy, Leonid Karlinsky, Amit Aides, Joseph Shtok, Sivan Harary, Rogério Schmidt Feris, Raja Giryes, and Alexander M. Bronstein. Laso: Label-set operations networks for multi-label few-shot learning. In *CVPR*, pages 6548–6557, 2019.
- [Bowman *et al.*, 2023] Benjamin Bowman, Alessandro Achille, Luca Zancato, Matthew Trager, Pramuditha Perera, Giovanni Paolini, and Stefano Soatto. À-la-carte prompt tuning (APT): combining distinct data via composable prompting. In *CVPR*, pages 14984–14993, 2023.
- [Chen *et al.*, 2022] Xiang Chen, Ningyu Zhang, Xin Xie, Shumin Deng, Yunzhi Yao, Chuanqi Tan, Fei Huang, Luo Si, and Huajun Chen. Knowprompt: Knowledge-aware prompt-tuning with synergistic optimization for relation extraction. In *WWW*, pages 2778–2788, 2022.
- [Chen *et al.*, 2023] Wentao Chen, Chenyang Si, Zhang Zhang, Liang Wang, Zilei Wang, and Tieniu Tan. Semantic prompt for few-shot image recognition. *CoRR*, abs/2303.14123, 2023.
- [Chiang *et al.*, 2023] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023.
- [Chua *et al.*, 2009] Tat-Seng Chua, Jinhui Tang, Richang Hong, Haojie Li, Zhiping Luo, and Yantao Zheng. NUS-WIDE: a real-world web image database from national university of singapore. In *CIVR*, 2009.
- [Everingham *et al.*, 2010] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John M. Winn, and Andrew Zisserman. The pascal visual object classes (VOC) challenge. *IJCV*, 88(2):303–338, 2010.
- [Fu *et al.*, 2024] Zhongtian Fu, Kefei Song, Luping Zhou, and Yang Yang. Noise-aware image captioning with progressively exploring mismatched words. In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan, editors, *AAAI*, pages 12091–12099. AAAI Press, 2024.
- [Gao *et al.*, 2021] 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. *CoRR*, abs/2110.04544, 2021.
- [Gu *et al.*, 2022] Yuxian Gu, Xu Han, Zhiyuan Liu, and Minlie Huang. PPT: pre-trained prompt tuning for few-shot learning. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *ACL*, pages 8410–8423, 2022.
- [Guo *et al.*, 2023] Zixian Guo, Bowen Dong, Zhilong Ji, Jinfeng Bai, Yiwen Guo, and Wangmeng Zuo. Texts as images in prompt tuning for multi-label image recognition. In *CVPR*, pages 2808–2817, 2023.
- [He *et al.*, 2016] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [Hu *et al.*, 2023] Ping Hu, Ximeng Sun, Stan Sclaroff, and Kate Saenko. Dualcoop++: Fast and effective adaptation to multi-label recognition with limited annotations. *TPAMI*, 2023.
- [Jia *et al.*, 2021] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In Marina Meila and Tong Zhang, editors, *ICML*, volume 139, pages 4904–4916, 2021.
- [Krasin *et al.*, 2017] Ivan Krasin, Tom Duerig, Neil Alldrin, Vittorio Ferrari, Sami Abu-El-Haija, Alina Kuznetsova, Hassan Rom, Jasper Uijlings, Stefan Popov, Andreas Veit, et al. Openimages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from <https://github.com/openimages>*, 2017.
- [Lee *et al.*, 2018] Chung-Wei Lee, Wei Fang, Chih-Kuan Yeh, and Yu-Chiang Frank Wang. Multi-label zero-shot learning with structured knowledge graphs. In *CVPR*, pages 1576–1585, 2018.
- [Li *et al.*, 2023] Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, volume 202, pages 19730–19742, 2023.
- [Lin *et al.*, 2014] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *ECCV*, volume 8693, pages 740–755, 2014.
- [Lin, 2023] Dekun Lin. Probability guided loss for long-tailed multi-label image classification. In *AAAI*, pages 1577–1585, 2023.
- [Mao *et al.*, 2023] Jun-Xiang Mao, Wei Wang, and Min-Ling Zhang. Label specific multi-semantics metric learning for multi-label classification: Global consideration helps. In *IJCAI*, pages 4055–4063, 2023.
- [Min *et al.*, 2022] Sewon Min, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Noisy channel language model prompting for few-shot text classification. In *ACL*, 2022.
- [Misra *et al.*, 2016] Ishan Misra, C. Lawrence Zitnick, Margaret Mitchell, and Ross B. Girshick. Seeing through the human reporting bias: Visual classifiers from noisy human-centric labels. In *CVPR*, pages 2930–2939, 2016.
- [Nie *et al.*, 2022] Xing Nie, Bolin Ni, Jianlong Chang, Gaomeng Meng, Chunlei Huo, Zhaoxiang Zhang, Shiming Xiang, Qi Tian, and Chunhong Pan. Pro-tuning: Unified prompt tuning for vision tasks. *CoRR*, 2022.
- [OpenAI, 2023] OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023.
- [Pu *et al.*, 2022] Tao Pu, Tianshui Chen, Hefeng Wu, and Liang Lin. Semantic-aware representation blending for multi-label image recognition with partial labels. In *AAAI*, pages 2091–2098, 2022.

- [Radford *et al.*, 2021] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In Marina Meila and Tong Zhang, editors, *ICML*, volume 139, pages 8748–8763, 2021.
- [Simon *et al.*, 2022] Christian Simon, Piotr Koniusz, and Mehrtash Harandi. Meta-learning for multi-label few-shot classification. In *MACV*, pages 346–355, 2022.
- [Sun *et al.*, 2022] Ximeng Sun, Ping Hu, and Kate Saenko. Dualcoop: Fast adaptation to multi-label recognition with limited annotations. In *NeurIPS*, 2022.
- [Sun *et al.*, 2023] Rui Sun, Naisong Luo, Yuwen Pan, Huayu Mai, Tianzhu Zhang, Zhiwei Xiong, and Feng Wu. Appearance prompt vision transformer for connectome reconstruction. In *IJCAI*, pages 1423–1431, 2023.
- [Vaswani *et al.*, 2017] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017.
- [Wang *et al.*, 2016] Jiang Wang, Yi Yang, Junhua Mao, Zhiheng Huang, Chang Huang, and Wei Xu. CNN-RNN: A unified framework for multi-label image classification. In *CVPR*, pages 2285–2294, 2016.
- [Wang *et al.*, 2017] Zhouxia Wang, Tianshui Chen, Guanbin Li, Ruijia Xu, and Liang Lin. Multi-label image recognition by recurrently discovering attentional regions. In *ICCV*, pages 464–472, 2017.
- [Wang *et al.*, 2023a] Haobo Wang, Shisong Yang, Gengyu Lyu, Weiwei Liu, Tianlei Hu, Ke Chen, Songhe Feng, and Gang Chen. Deep partial multi-label learning with graph disambiguation. In *IJCAI*, pages 4308–4316, 2023.
- [Wang *et al.*, 2023b] Henan Wang, Muli Yang, Kun Wei, and Cheng Deng. Hierarchical prompt learning for compositional zero-shot recognition. In *IJCAI*, pages 1470–1478, 2023.
- [Wei *et al.*, 2016] Yunchao Wei, Wei Xia, Min Lin, Junshi Huang, Bingbing Ni, Jian Dong, Yao Zhao, and Shuicheng Yan. HCP: A flexible CNN framework for multi-label image classification. *TPAMI*, 38(9):1901–1907, 2016.
- [Xi *et al.*, 2023] Wenjuan Xi, Xin Song, Weili Guo, and Yang Yang. Robust semi-supervised learning for self-learning open-world classes. In *ICDM*, pages 658–667, 2023.
- [Xu *et al.*, 2023] Xiao Xu, Chenfei Wu, Shachar Rosenman, Vasudev Lal, Wanxiang Che, and Nan Duan. Bridgetower: Building bridges between encoders in vision-language representation learning. In *AAAI*, pages 10637–10647, 2023.
- [Yan *et al.*, 2023] Liqi Yan, Cheng Han, Zenglin Xu, Dongfang Liu, and Qifan Wang. Prompt learns prompt: Exploring knowledge-aware generative prompt collaboration for video captioning. In *IJCAI*, pages 1622–1630, 2023.
- [Yang *et al.*, 2018] Yang Yang, Yi-Feng Wu, De-Chuan Zhan, Zhi-Bin Liu, and Yuan Jiang. Complex object classification: A multi-modal multi-instance multi-label deep network with optimal transport. In *KDD*, pages 2594–2603, 2018.
- [Yang *et al.*, 2021] Yang Yang, Zhao-Yang Fu, De-Chuan Zhan, Zhi-Bin Liu, and Yuan Jiang. Semi-supervised multi-modal multi-instance multi-label deep network with optimal transport. *TKDE*, 33, 2021.
- [Yang *et al.*, 2023] Yang Yang, Ran Bao, Weili Guo, De-Chuan Zhan, Yilong Yin, and Jian Yang. *Sci. China Inf. Sci.*, 66(12), 2023.
- [Yu *et al.*, 2023] Lang Yu, Qin Chen, Jiayu Lin, and Liang He. Black-box prompt tuning for vision-language model as a service. In *IJCAI*, pages 1686–1694, 2023.
- [Zeng *et al.*, 2023] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma, Yufei Xue, Jidong Zhai, Wenguang Chen, Zhiyuan Liu, Peng Zhang, Yuxiao Dong, and Jie Tang. GLM-130B: an open bilingual pre-trained model. In *ICLR*, 2023.
- [Zhou *et al.*, 2022a] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *CVPR*, pages 16795–16804, 2022.
- [Zhou *et al.*, 2022b] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *IJCV*, 130(9):2337–2348, 2022.
- [Zhu *et al.*, 2023] Xuelin Zhu, Jiuxin Cao, Jian Liu, Dongqi Tang, Furong Xu, Weijia Liu, Jiawei Ge, Bo Liu, Qingpei Guo, and Tianyi Zhang. Text as image: Learning transferable adapter for multi-label classification. *CoRR*, abs/2312.04160, 2023.

A Appendix Overview

Here we provide more information about our proposed PVP, more ablation studies, and visualization results. The appendix is organized as follows. In Appendix B, we introduce the construction of text training data in detail and present examples of generated text data and synonym dictionaries. In Appendix C, we compare our method with the recent work TAI-Adapter with different LLMs of text generation. We further present more visualization results on the MSCOCO dataset in Appendix D and conduct hyper-parameter experiments in Appendix E.

B Training Text Data Construction

We present the details of the training text data construction in this section.

Human-Annotated Labeled Text Data Construction. For VOC2007 [Everingham *et al.*, 2010] and MS-COCO [Lin *et al.*, 2014] benchmarks, we obtain the labeled coco-captions from MS-COCO, each text succinctly describes a natural scene, with a maximum length of 25. For NUSWIDE [Chua *et al.*, 2009] dataset, we collect localized narratives from the OpenImages [Krasin *et al.*, 2017]. Each text contains detailed content descriptions, with a maximum length of 60. The examples are shown in Figure 6.

- A young boy stares up at the computer monitor.
- Man in all black doing a trick on his skateboard.
- Men are crowded on back of a overloaded pickup truck.
- The woman in the kitchen is holding a huge pan.
- People riding bicycles down the road approaching a bird.
- A bathroom with a walk in shower currently under repair.
- There is a bathtub and a counter in a bathroom.
- A standing toilet in a bathroom next to a window.
- A large porcelain toilet posed with a tan flower pot.
- ...

(a) COCO-Caption

- The image shows that there are three buses , a car and a lorry on the road. Here we can see buildings , windows and air conditioners. We can also see a street light and this is a tree. Man in all black doing a trick.
- In this picture we can see one person is standing and talking with the microphone in front of the desk, side we can see on a table covered with white cloth, on it we have flower, bottles, glasses. beside the table they is a display board. The woman is holding a huge pan.
- ...

(b) Localized narratives

Figure 6: Labeled text data constructed from COCO-Caption of MS-COCO and localized narratives of OpenImages.

LLMs-based Pseudo Text Data Construction. We describe the LLMs-based pseudo text data construction in detail. Figure 7 illustrates the process of LLMs generating

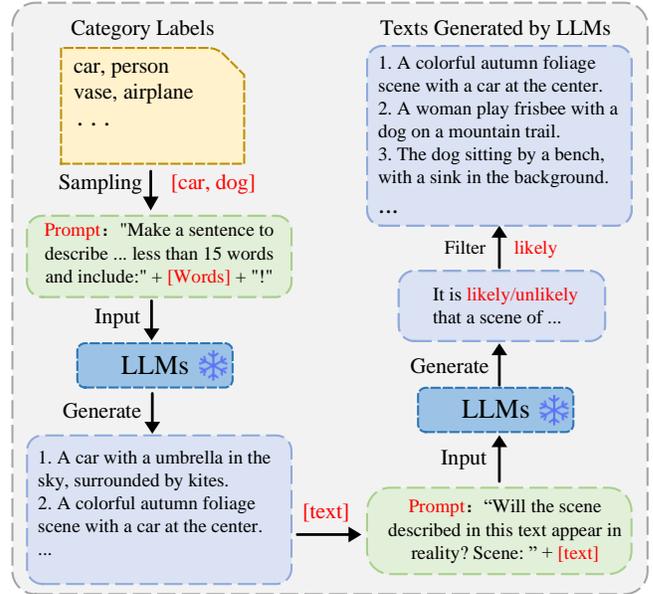


Figure 7: Pseudo text generated by LLMs. Via labels sampling and rationality judgment, obtaining semantically reasonable sentences.

pseudo text data through artificially designed prompt templates and a set of category labels. Given a target category set $\mathcal{C} = \{c_1, c_2, \dots, c_N\}$, where N denotes the number of categories and c_i denotes a particular class, we design a query prompt as follows:

PROMPT: Make a sentence to describe a photo. Requirements: Each sentence should be less than 15 words and include keywords: $\{c_{i_1}, c_{i_2}, \dots, c_{i_l}\}$.

Here $\{c_{i_1}, c_{i_2}, \dots, c_{i_l}\} \subset \mathcal{C}$ and $l \leq 3$. Then, we randomly sample l categories $\{c_{i_1}, c_{i_2}, \dots, c_{i_l}\}$ and input the query prompt to LLMs for generating pseudo text descriptions automatically. To filter out the unreliable sentences generated by LLMs, we re-input the generated text data into LLMs with another query template:

PROMPT: Will the scene described in this text appear in reality? Scene: + $\{\text{text}\}$.

Then we judge the reasonableness of the sentence through the output *likely/unlikely*. For word-level filtered labels in input sentences, we follow the setting of TAI-DPT [Guo *et al.*, 2023], using NLTK⁴ to perform noun filtering. Due to each target category has synonyms with similar meanings, these synonyms also need to be mapped into the class label. Hence, we construct a synonym dictionary, which includes common synonyms of each class in the target dataset. As shown in Figure 8, all the words in each row belong to the same class label. We utilize the synonym dictionary and conduct noun filtration by tokenizing and lemmatizing the words to search for sentences that contain at least one synonym name. The text data that do not match any synonym is discarded. This simple noun filtering strategy ensures that each input text contain at least one class label for prompt tuning.

In Figure 9, we provide an example for labeled text data

⁴<https://www.nltk.org/>

Table 5: Comparison with TAI-Adapter on all datasets. The best performance is shown in boldface.

Method	LLM	VOC2007	MS-COCO	NUSWIDE
TAI-DPT	N/A	88.3	65.1	46.5
TAI-Adapter	Vicuna-33b	89.0	67.7	53.3
PVP	ChatGLM	88.9	67.5	49.3
PVP	Vicuna-33b	89.5	68.9	51.4

['person', 'human', 'people', 'man', 'woman', 'passenger']
 ['bicycle', 'bike', 'cycle']
 ['car', 'taxi', 'auto', 'automobile', 'motor car']
 ['motor bike', 'motor cycle']
 ['aeroplane', 'air craft', 'jet', 'plane', 'air plane']
 ['train', 'rail way', 'railroad']
 ['handbag', 'hand bag', 'pocketbook', 'purse']
 ['baseball glove', 'baseball mitt', 'baseball game glove']
 ['potted plant', 'house plant', 'bonsai', 'pot plant']
 ...

Synonym Dictionary of MS-COCO

Figure 8: Illustration of synonym dictionary.

construction. Given a set of target categories, we randomly sample several categories, such as person, bench, etc. These categories and the first prompt are constructed into a input template that can be processed by LLMs [Zeng *et al.*, 2023], thereby generating a group of short text sentences, such as “A bench in a post office with a person sitting on it”. We then design a rationality judgment template, and concatenate it with the text sentences obtained from the previous dialogue and input it into the LLMs again. Based on the generated results, we retain sentences with reasonable content, thereby forming noisy text training data.

C Comparison with TAI-Adapter

Recent work TAI-Adapter [Zhu *et al.*, 2023] proposed to using a random perturbation mechanism to enhance the transferable capability of the adapter module. TAI-Adapter also required to use massive label visual data during training like TAI-DPT. We compare TAI-Adapter with our method in this section.

TAI-Adapter utilizes the Vicuna-33b-1.3v [Chiang *et al.*, 2023] to construct training text data. For fair comparison, we adopt the same LLMs to construct training text data for PVP. The mAP results⁵ are reported in Table 5. From Table 5, we can find that our proposed PVP can outperform TAI-Adapter in most cases when we utilize Vicuna-33b-1.3v

⁵The results of TAI-Adapter are directly referred from the origin paper.

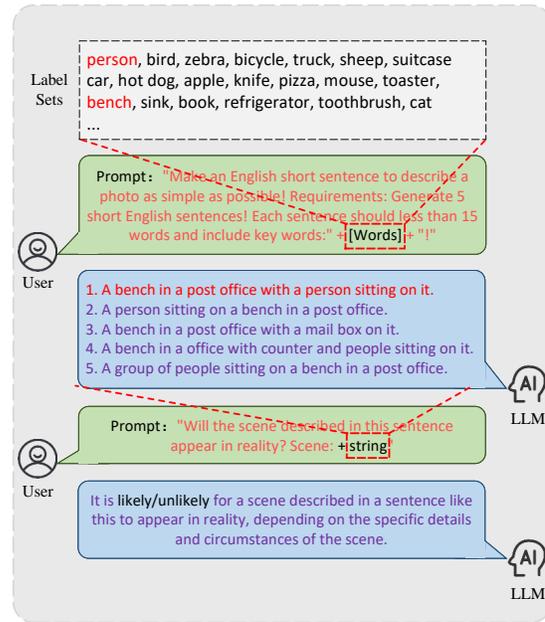


Figure 9: Pseudo text generated by LLMs. Semantically reasonable sentences are obtained via labels sampling and rationality judgment.

to generate training text data. Moreover, by comparing the mAP results based on ChatGLM and Vicuna-33b-1.3v for our proposed method, we can find that the mAP with Vicuna-33b-1.3v is higher than that with ChatGLM. The reason is that the Vicuna-33b-1.3v can extract more diverse features than ChatGLM.

D Visualization Results

In this section, we illustrate more visualization results. In Figure 10, we provide some examples that are randomly selected for visualization. Specifically, we randomly select several common target categories and their corresponding raw images, visualizing the correlation between local image features and class prompt embeddings of TAI-DPT and our method PVP. For raw image, we present the ground-truth category labels. We also compute the similarities between global image feature and class prompt embeddings of TAI-DPT and PVP. And we present the top 5 categories with the highest prediction confidence from different methods. Figure 10 presents the visualization results for the categories including “bicycle”, “hot dog”, “airplane” and “television”. From Figure 10, we can see that our method can learn more diverse and com-

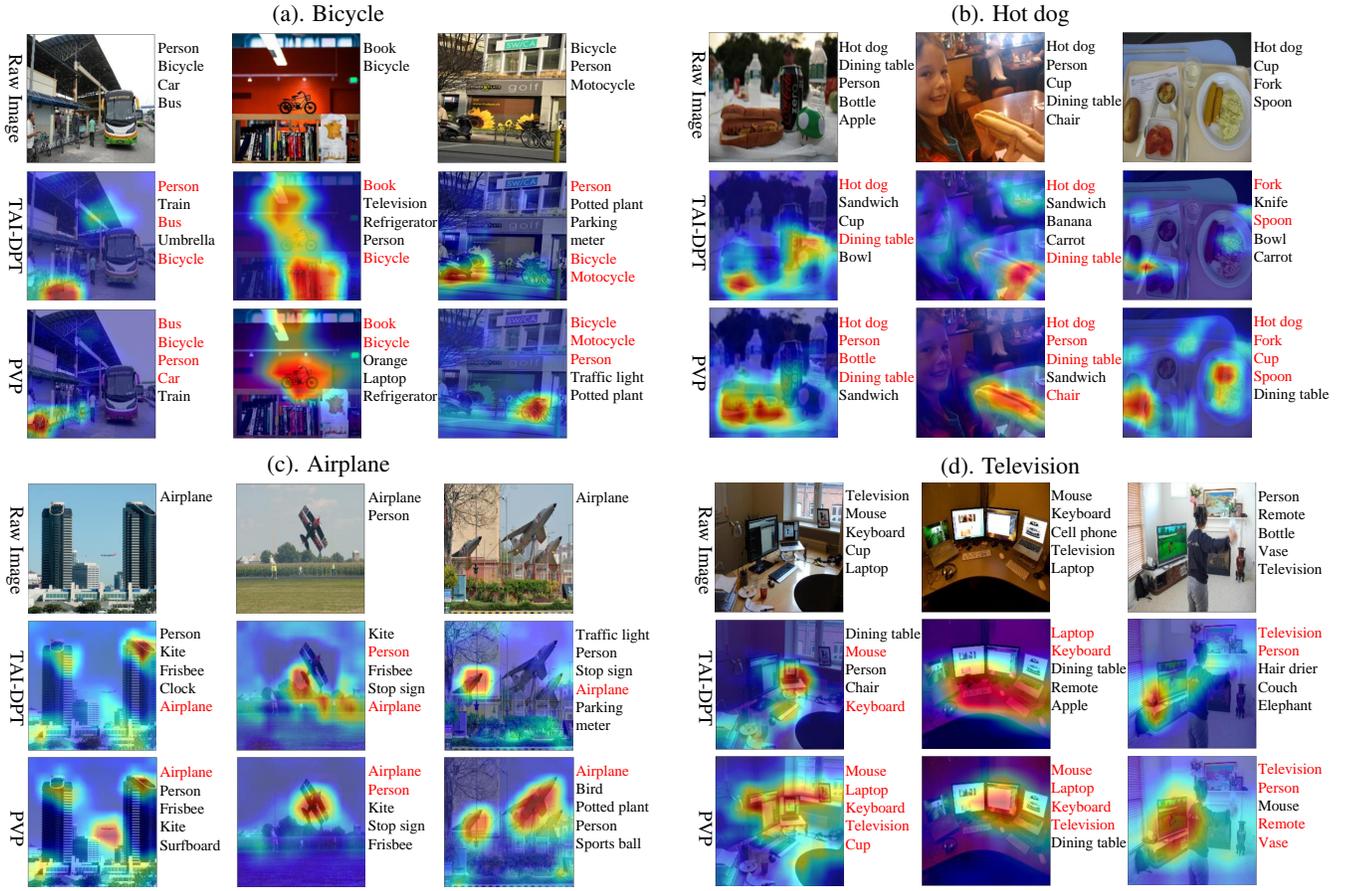


Figure 10: Visualization results for categories “bicycle”, “hot dog”, “airplane” and “television”.

prehensive features. For the images whose shapes and attributes are different but with the same category labels, the accuracy of our method is higher than that of TAI-DPT.

E Further Analysis

E.1 Ablation Study

Training Loss. We discuss the impact of different loss function on multi-label image recognition performance. For the contrastive training of visual and text prompts, we compare contrastive loss and ranking loss. The same loss settings are also used for the contrastive training of prompts and global text features. The results are shown in Table 6. In Table 6, “CE” and “RL” denote the cross-entropy and ranking loss, respectively. As shown in Table 6, applying cross-entropy loss between prompts and ranking loss between prompts and text features achieves the highest multi-label image recognition performance on VOC2007, MS-COCO, and NUSWIDE datasets.

E.2 Sensitivity to Hyper-Parameters

In this section, we present the influence of hyper-parameters, including pseudo-visual prompt size (H and W), training epochs of the first stage.

Table 6: Ablation study for different loss function on all datasets.

$\mathcal{L}_{visual}/\mathcal{L}_{text}$	\mathcal{L}_{vtc}	VOC2007	MS-COCO	NUSWIDE
CE	CE	87.9	68.6	48.3
RL	RL	87.4	67.8	47.5
CE	RL	84.7	65.5	45.6
RL	CE	90.1	70.8	49.3

Prompt Size. Pseudo-visual prompts are processed through an image encoder, with a size of $H \times W \times 3$. Therefore, we explore the impact of different prompt sizes on multi-label image classification performance. In Table 7, for VOC2007, MS-COCO, and NUSWIDE, as the prompt size gradually increases, the image recognition performance on all datasets shows a trend of first rising and then falling, reaching the best performance at the size of $224 \times 224 \times 3$. This indicates that larger prompt size can learn more extensive and diverse visual knowledge to co-learn visual and text prompts. However, an overly large prompt size will overfit the text training data, leading to weaker generalization ability in image testing, thereby affecting the performance of image recognition.

Length of text prompt. In our method, we set the length of the text prompt to 16 following the previous SOTA TAI-DPT. Here, we conduct several experiments to compare the

Table 7: Ablation study for the initialized size of pseudo-visual prompt on all datasets.

PVP	VOC2007	MS-COCO	NUSWIDE
$96 \times 96 \times 3$	89.4	69.9	48.5
$128 \times 128 \times 3$	89.6	70.5	48.9
$160 \times 160 \times 3$	89.8	70.7	48.9
$224 \times 224 \times 3$	90.1	70.8	49.3
$288 \times 288 \times 3$	89.9	70.5	49.2
$324 \times 324 \times 3$	89.8	70.6	48.9

Table 8: Results of the parameter analysis.

M	8	12	16	20	24
mAP	70.40	70.56	70.78	70.71	70.74
λ	0.2	0.4	0.5	0.6	0.8
mAP	70.23	70.62	70.78	70.75	70.64

influence of different prompt lengths on MSCOCO dataset. From Table 8, the prompt length (parameter M) ranges from 8 to 24, we can see the performance of different prompt lengths is similar to others, and the larger size of the prompt can not further improve the image classification results.

Weight of dual-adapter and CLIP. Dual-adapter is designed to both learn the new knowledge of downstream datasets and maintain the origin knowledge of the pretrained CLIP. Therefore, we explore the weight of dual-adapter and origin CLIP to evaluate the importance of both the downstream MSCOCO dataset and CLIP. From Table 8, the weight, denoted as λ , ranges from 0.2 to 0.8, PVP achieves the best performance with the λ being 0.5, demonstrating the knowledge information of MSCOCO and origin CLIP are equally important.

Table 9: Analysis for γ , η , and ν .

γ	η	ν	MSCOCO
1	1	1	70.78
1	2	3	70.62
1	3	2	70.63
2	1	3	70.76
2	3	1	70.80
3	1	2	70.65
3	2	1	70.69

Weights of Loss: Here, we provide the experimental results with different loss weights. We first rewrite the objective function as $\mathcal{L} = \gamma\mathcal{L}_{vtc} + \eta\mathcal{L}_{visual} + \nu\mathcal{L}_{text}$. The experimental results are shown in Table 9. From Table 9, we can see that the mAP almost remains unchanged with different values of weights. We will discuss the influence of the weights within a larger range in the final version.

Training Epochs of The First Stage. Our method consists of two stages: pseudo-visual prompt learning and transferable visual and text prompts co-learning. We explore how the epochs of pseudo-visual prompt learning in the first stage af-

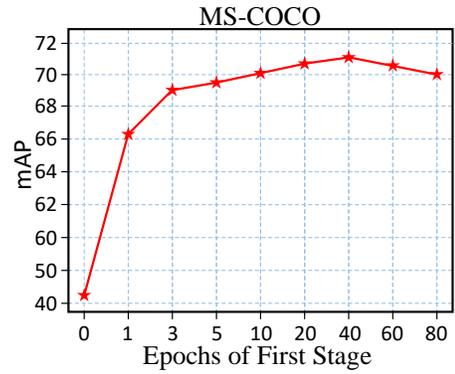


Figure 11: Ablation study for the training epochs of pseudo-visual prompt learning stage on MS-COCO dataset.

fect the multi-label image recognition performance in the second stage. Figure 11 shows that the longer the pseudo-visual prompts are learned, the greater the improvement in image recognition performance on MS-COCO dataset. Moreover, the case where epoch equals 0 denotes that pseudo-visual prompt and text prompt co-learning is performed directly without the pseudo-visual prompt learning in the first stage. This result indicates that the mAP results will be hindered if we only perform transferable prompt co-learning without pseudo-visual prompt learning.