

Tuning Vision-Language Models with Candidate Labels by Prompt Alignment [★]

Zhifang Zhang¹, Yuwei Niu², Xin Liu³, and Beibei Li¹ ✉

¹ College of Computer Science, Chongqing University, Chongqing 400044, China
libeibeics@cqu.edu.cn

² Hongshen Honor School, Chongqing University, Chongqing 400044, China

³ School of Computer Science and Engineering, Southeast University, Nanjing
210096, China

Abstract. Vision-language models (VLMs) can learn high-quality representations from a large-scale training dataset of image-text pairs. Prompt learning is a popular approach to fine-tuning VLM to adapt them to downstream tasks. Despite the satisfying performance, a major limitation of prompt learning is the demand for labeled data. In real-world scenarios, we may only obtain candidate labels (where the true label is included) instead of the true labels due to data privacy or sensitivity issues. In this paper, we provide the first study on prompt learning with candidate labels for VLMs. We empirically demonstrate that VLMs can learn from candidate labels through prompt learning. Nonetheless, its performance drops when the label ambiguity increases. In order to improve its robustness, we propose a simple yet effective framework that better leverages the prior knowledge of VLMs to guide the learning process with candidate labels. Specifically, our framework disambiguates candidate labels by aligning the model output with the mixed class posterior jointly predicted by both the learnable and the handcrafted prompt. Besides, our framework can be equipped with various training objectives for learning with candidate labels to further improve their performance. Extensive experiments demonstrate the effectiveness of our proposed framework. Our code is available at <https://github.com/zhangzf01/CoOpPLL>.

Keywords: Multimodal Models · Prompt Learning · Candidate Labels.

1 Introduction

Vision-language models (VLMs), such as CLIP [20], have become excellent base models in multiple domains, most of which employ a dual-encoder architecture to align the natural images with descriptive texts. Remarkably, this special training pattern has endowed VLMs with superior zero-shot transfer performance on visual recognition tasks. In specific, during the inference, the pre-trained text encoder receives inputs in the form of man-crafted prompts, *e.g.*, “a photo of

[★] This work was supported by Chongqing Science and Technology Bureau (CSTB2022TIAD-KPX0180).

$\langle CLS \rangle$.". Subsequently, all the generated textual embeddings are matched with the visual embedding obtained from the image encoder to predict the image category. However, the powerful zero-shot ability of VLMs was shown to be heavily dependent on the choice of handcrafted prompts, which needs substantial efforts and professional domain knowledge to design [32]. To avoid the manual design of the prompts, *prompt learning* [32] is proposed, which treats the textual prompt as additional learnable parameters and tunes them while keeping all the original parameters of the pre-trained model fixed. Later, the concept of the prompt is extended to visual prompt [9] and multi-modal prompt [10] in VLMs. Overall, there has been increasing attention paid to prompt learning due to its potential to perform significantly better than zero-shot transfer with a few sets of labeled data [32,31,10,9,30].

While prompt learning has demonstrated effectiveness and efficiency in few-shot supervised learning, the true labels must be provided for the training data used in prompt learning. This is a significant defect and will limit the usage of prompt learning in various real-world scenarios because we may be unable to collect accurate labels due to security issues or labeling difficulties. Fortunately, obtaining a set of candidate labels that includes the true label in these situations is easier. As we see, learning with only candidate labels (also widely known as *partial-label learning (PLL)* [23,26,16,5,7,12,13]) is practically significant, which also has arisen in many vital applications such as web mining [15], online annotation [22] and ecoinformatics [14]. Nevertheless, existing PLL methods primarily focus on training a model from scratch, and the effectiveness of PLL in the new training paradigm called prompt learning remains unconfirmed. To bridge this research gap, we, for the first time, explore the validity and potential approaches for prompt learning with candidate labels.

This paper empirically shows that prompt learning combined with the prevailing PLL training objectives in a vanilla way can learn from candidate labels in Figure 1. However, as experimentally suggested, if the candidate labels become more ambiguous, the model’s performance will drop significantly. Fortunately, prompt learning is still more robust than linear probe [20], another tuning method that trains a linear classifier on top of a frozen pre-trained model. We conjecture that the reason lies in the prior knowledge brought by the fixed class token, which can keep VLMs from over-fitting to the false-positive labels in the candidate label set and provide VLMs with preferred zero-shot ability, thus mitigating the error accumulation problem [27] in PLL with high label ambiguity.

Therefore, to enhance the robustness of prompt learning with candidate labels, we propose a simple yet effective framework incorporating the handcrafted prompt to distill the model with more comprehensive prior knowledge. Concretely, it dynamically mixes the class posteriors predicted by both the handcrafted and learnable prompt, followed by aligning the mixed class posterior with the model output. Besides, due to the simplicity and flexibility of our framework, it can cooperate with any current PLL training objectives. With our framework, the overall performance of various PLL training objectives has improved by a large margin when tuning VLMs with candidate labels.

Tuning VLMs with Candidate Labels by Prompt Alignment

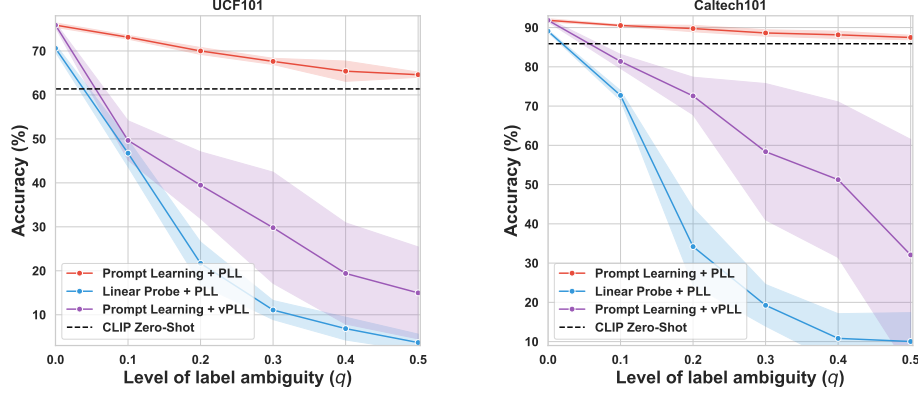


Fig. 1. Performance comparison with multiple fine-tuning approaches combined with PiCO [23] in a vanilla way on UCF101 [21] and Caltech101 [4] with candidate labels of the incremental label ambiguity. vPLL means a simple baseline that treats every candidate label as the ground-truth label and uses cross-entropy loss to learn. We define the level of label ambiguity q as the uniform probability of flipping negative labels $\bar{y}_i \neq y_i$ to false-positive labels inside the candidate label set Y_i : $q = \Pr(\bar{y}_i \in Y_i | \bar{y}_i \neq y_i)$.

Our main contributions can be summarized as follows:

- We provide the first study on the scenario when vision-language models are tuned with only candidate labels.
- We demonstrate empirically and explain that prompt learning combined with PLL training objectives in a vanilla way can learn from candidate labels but is not robust when the label ambiguity is high.
- A framework is proposed to enhance the robustness of prompt learning with candidate labels by aligning the model output with the dynamically mixed prediction by both the handcrafted and the learnable prompt.
- Extensive experiments demonstrate the effectiveness of our framework.

2 Preliminaries of Prompt Learning

CoOp [32] is the first work that migrated prompt learning to vision tasks, which tunes CLIP by optimizing the parameters of the learnable textual prompt (also called soft prompt) while keeping the class token fixed. To be specific, assume CoOp introduces M learnable vectors $\{\mathbf{v}_k\}_{k=1}^M$ and C fixed class tokens $\{\mathbf{c}_l\}_{l=1}^C$. Together, they are usually concatenated to form the full prompt $\mathbf{s}_i = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_M, \mathbf{c}_i\}$ for class i . Let the normalized image embedding be \mathbf{f}^v , then the class posterior is estimated as:

$$p(y = i | \mathbf{x}) = \frac{\exp(\text{sim}(\mathbf{f}^v, \text{TextProj}(\mathbf{s}_i)) / \tau)}{\sum_{j=1}^C \exp(\text{sim}(\mathbf{f}^v, \text{TextProj}(\mathbf{s}_j)) / \tau)}. \quad (1)$$

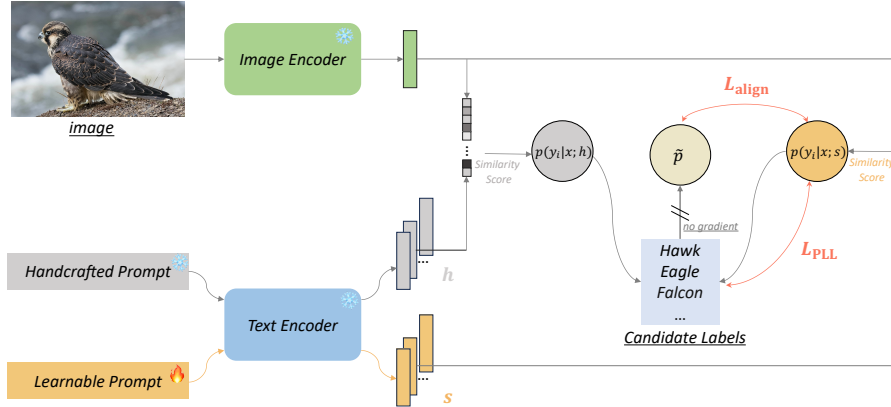


Fig. 2. Illustration of our framework. Our framework includes a prompt alignment module and a PLL method module. All parameters of this framework are **frozen** except for the learnable prompt.

At last, the learnable context vectors $\{\mathbf{v}_k\}_{k=1}^M$ are optimized on a dataset $D = \{(\mathbf{x}_i, y_i)_{i=1}^N\}$ with the cross-entropy loss:

$$\mathcal{L}_{\text{true}} = -\mathbb{E}_{(\mathbf{x}, y) \in D} [\log p(y|\mathbf{x})]. \quad (2)$$

Notably, optimizing this training objective of CoOp requires examples of true labels. But in many realistic scenarios, the true label is not accessible. The former research has focused on prompt learning with noisy labels [25] or no label [8,29] and showed prompt learning is not only robust with label noise but also can effectively learn from unlabelled data with specifically designed algorithms.

In this work, we focus on the scenario where only a candidate label set can be obtained, which we unfortunately cannot utilize to optimize the above training objective. Therefore, we will study prompt learning with candidate labels using PLL training objectives instead of cross-entropy loss.

3 Our Framework

This framework is proposed to improve the robustness of prompt learning with candidate labels. It provides a powerful regularization that aligns the dynamically mixed prediction of the handcrafted and learnable prompt with the current model output using weighted cross-entropy loss. It is shown in Figure 2.

Prompt Alignment Regularization.

Let \mathcal{X} be the input space, and $\mathcal{Y} = \{1, 2, \dots, C\}$ be the label space. The i -th learnable prompt $\mathbf{s}_i = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_M, \mathbf{c}_i\}$, where $\{\mathbf{v}_m\}_{m=1}^M$ denotes M learnable tokens and \mathbf{c}_i is the word embedding for the i -th class name. Similarly, the i -th manually crafted prompt is denoted as \mathbf{h}_i . To clarify, $f_i(\mathbf{x}; \mathbf{s})$ and $g_i(\mathbf{x}; \mathbf{h})$ are the softmax outputs of the i -th label, as predicted by the learnable prompt

and handcrafted prompt separately. When (\mathbf{x}, Y) is drawn from the partialized dataset, the prompt alignment loss can be calculated as:

$$\mathcal{L}_{\text{align}}(\mathbf{x}, Y) = - \sum_{i=1}^C \tilde{p}_i \log f_i(\mathbf{x}; \mathbf{s}), \quad (3)$$

where \tilde{p}_i is mixed linearly with the class posteriors predicted by both the learnable and handcrafted prompts:

$$\tilde{p}_i = \alpha p(y = i \mid \mathbf{x}; \mathbf{s}) + (1 - \alpha) p(y = i \mid \mathbf{x}; \mathbf{h}). \quad (4)$$

Since the non-candidate labels can never be the ground-truth label, the class posteriors are recalculated as:

$$p(y = i \mid \mathbf{x}; \mathbf{s}) = \begin{cases} \frac{f_i(\mathbf{x}; \mathbf{s})}{\sum_{j \in Y} f_j(\mathbf{x}; \mathbf{s})}, & i \in Y, \\ 0, & i \notin Y. \end{cases} \quad (5)$$

$$p(y = i \mid \mathbf{x}; \mathbf{h}) = \begin{cases} \frac{g_i(\mathbf{x}; \mathbf{h})}{\sum_{j \in Y} g_j(\mathbf{x}; \mathbf{h})}, & i \in Y, \\ 0, & i \notin Y. \end{cases} \quad (6)$$

This regularization can be adapted to any prevailing PLL training objective:

$$\mathcal{L}_{\text{total}}(\mathbf{x}, Y) = \mathcal{L}_{\text{PLL}}(\mathbf{x}, Y) + \beta \mathcal{L}_{\text{align}}(\mathbf{x}, Y), \quad (7)$$

where β is the factor controlling the strength of the alignment loss.

Dynamic Mixing Strategy.

In Equation (4), we use a balancing factor α to mix the handcrafted and learnable prompt predictions. However, a fixed balancing factor may be sub-optimal since the prediction quality of the handcrafted prompt will be influenced by several conditions, such as:

- Label Ambiguity: In cases of high label ambiguity, the handcrafted prompt tends to provide more reliable predictions, as soft prompts are difficult to learn effectively in such scenarios.
- Dataset Complexity: For datasets with poor CLIP zero-shot performance, reducing reliance on the handcrafted prompt helps avoid steering the learning process in the wrong direction.
- Training Epoch: If the soft prompt outperforms the handcrafted prompt as training proceeds, it is better to make more use of the soft prompt, ensuring the most effective predictions are leveraged.

Therefore, to further enhance the robustness of our framework against these conditions, a dynamic mixing strategy is adopted to adjust the balancing factor:

$$\tilde{p}_i = \alpha(t) p(y = i \mid \mathbf{x}; \mathbf{s}) + (1 - \alpha(t)) p(y = i \mid \mathbf{x}; \mathbf{h}), \quad (8)$$

$$\alpha(t) = \min\left\{\frac{t}{T'}, \lambda\right\}, \quad (9)$$

where T' controls how quickly the balancing factor grows over time, adjusting the shift between the two outputs' contributions and λ sets the maximum weight that the output of soft prompts can have, ensuring its influence is capped. Together, they dynamically balance the contributions of the two prompts for robustness.

4 Experiments

4.1 Experimental Setting

Datasets. We adopt 8 image recognition datasets: ImageNet [3] and Caltech101 [4] for generic object classification; OxfordPets [18], StanfordCars [11], FGVCAircraft [17], and Food101 [1] for fine-grained classification; UCF101 [21] for action classification; DTD [2] for texture classification.

Implementation Details. Our implementation is based on Pytorch [19]. We apply prompt learning on a pre-trained CLIP whose backbone of the image encoder is ResNet-50 [6]. We use the prompt engineering of CLIP [20] to construct the handcrafted prompts. We use a 16-shot fine-tuning strategy and set the learnable prompt token to be 16 with other hyper-parameters the same as CoOp [32]. For the hyperparameters of our method, we set $\lambda = 0.5$, $T' = 25$, $\beta = 1$. Moreover, if a confidence matrix is required in the PLL method, it will be initialized with the model output before training. We report the average test accuracy and the standard deviation of 4 experiments with the seeds fixed.

Training Objectives. We prove the effectiveness of our framework by incorporating six state-of-the-art PLL methods: PRODEN [16], CC [5], LW [24], PiCO [23], PLLCR [26] and CAVL [28].

4.2 Main Results

In Table 1, we compare our framework with vanilla prompt learning for six PLL training objectives on ten benchmark datasets when $q = \{0.1, 0.3, 0.5\}$. For supervised learning, we use the same settings as vanilla prompt learning, except the cross entropy is used, while for CLIP zero-shot, we use the same prompt as the handcrafted prompt in our framework. It is shown that the performance of vanilla prompt learning drops significantly with higher label ambiguity and sometimes even under-performs zero-shot inference, verifying its lack of robustness. On the contrary, our framework not only shows robustness with highly ambiguous candidate labels but also improves general performance that matches supervised learning. For instance, in OxfordPets with LW at $q = 0.5$, vanilla prompt learning under-performs zero-shot inference by 5.84%, while our framework outperforms supervised learning by 2.21%. Moreover, our framework demonstrates consistent improvement across all levels of label ambiguity, not just in cases of high ambiguity. For example, in the Caltech101 dataset, our framework outperforms vanilla prompt learning in 16 out of 18 cases, regardless of the level of label ambiguity and PLL training objectives. Particularly, the more label ambiguity, the more gain in performance with our framework.

The main result has verified the efficacy of our framework. By aligning with the handcrafted prompt, the model can make more accurate predictions at the initial stage of the learning process. Moreover, the class posterior predicted by the handcrafted prompt is hardly affected by the level of label ambiguity, making the learning process more robust.

Tuning VLMs with Candidate Labels by Prompt Alignment

Table 1. Performance comparison of vanilla prompt learning and our framework. We evaluate the test accuracy of these two methods for different PLL training objectives at different label ambiguities. Supervised means prompt learning with true labels. The standard deviation is shown in parentheses.

	Caltech101			DTD			FGVCAircraft			ImageNet		
Supervised	92.03 (0.22)			62.95 (0.43)			27.27 (3.01)			60.90 (0.64)		
Zero-shot	85.84			42.79			17.07			58.16		
q	0.1	0.3	0.5	0.1	0.3	0.5	0.1	0.3	0.5	0.1	0.3	0.5
PRODEN	91.20 (0.37)	88.78 (0.11)	82.05 (1.37)	60.85 (0.42)	54.70 (1.35)	42.15 (1.23)	21.03 (1.65)	9.15 (2.11)	5.58 (0.97)	60.20 (0.19)	47.03 (2.26)	10.18 (0.19)
+ours	91.98 (0.26)	91.38 (0.39)	90.55 (0.32)	61.95 (1.20)	56.40 (0.49)	52.77 (2.01)	24.83 (1.06)	21.27 (0.93)	19.03 (0.37)	61.90 (0.27)	61.55 (0.30)	60.77 (0.11)
CC	91.68 (0.57)	91.23 (0.38)	90.83 (0.66)	61.10 (0.33)	55.58 (1.25)	48.95 (0.99)	26.00 (0.75)	22.42 (0.66)	18.65 (0.15)	61.40 (0.53)	61.08 (0.72)	60.83 (0.69)
+ours	92.35 (0.27)	91.52 (0.19)	90.80 (0.43)	61.25 (0.56)	56.25 (0.60)	51.77 (0.94)	26.00 (0.89)	22.00 (0.56)	19.85 (0.34)	62.12 (0.11)	61.50 (0.32)	61.08 (0.24)
LW	91.35 (0.26)	89.02 (0.37)	82.97 (1.13)	61.65 (1.11)	54.73 (1.84)	41.77 (1.36)	20.77 (1.15)	9.22 (2.04)	5.85 (0.93)	60.30 (0.82)	45.50 (3.52)	8.80 (1.58)
+ours	92.08 (0.37)	91.38 (0.29)	90.67 (0.16)	61.70 (0.47)	56.38 (0.94)	52.95 (1.65)	25.50 (0.75)	21.65 (0.84)	19.43 (0.22)	62.10 (0.35)	61.65 (0.46)	60.80 (0.14)
PLLCR	91.67 (0.29)	91.68 (0.31)	91.10 (0.25)	62.38 (0.38)	58.25 (0.93)	49.60 (1.64)	24.62 (0.81)	14.60 (1.99)	8.80 (2.43)	60.32 (0.31)	59.60 (0.58)	59.20 (3.32)
+ours	91.95 (0.25)	91.75 (0.34)	91.05 (0.59)	62.03 (1.09)	58.50 (0.94)	52.05 (0.84)	26.27 (0.48)	21.68 (0.20)	19.10 (0.30)	61.87 (0.24)	60.77 (0.21)	60.03 (0.42)
PiCO	90.47 (0.19)	88.75 (0.74)	87.75 (0.76)	61.70 (1.07)	55.08 (1.22)	48.08 (1.58)	23.60 (2.16)	18.27 (2.83)	14.75 (3.11)	58.98 (0.59)	55.35 (0.35)	52.72 (0.93)
+ours	91.92 (0.11)	91.33 (0.19)	90.72 (0.15)	62.67 (0.92)	58.42 (0.41)	55.08 (0.93)	25.12 (1.07)	22.18 (1.00)	20.32 (0.71)	62.03 (0.25)	60.92 (0.19)	59.80 (0.37)
CAVL	91.67 (0.40)	90.97 (0.51)	88.62 (1.30)	60.85 (0.97)	50.38 (1.45)	39.80 (4.97)	24.20 (0.51)	7.05 (3.50)	1.88 (0.79)	60.33 (0.62)	58.88 (0.94)	57.55 (1.18)
+ours	91.82 (0.33)	91.58 (0.43)	90.73 (0.13)	60.80 (0.95)	57.18 (0.45)	51.82 (1.67)	25.25 (0.97)	21.45 (0.47)	18.85 (0.17)	62.00 (0.19)	61.15 (0.25)	60.22 (0.39)

	OxfordPets			StanfordCars			UCF101			Food101		
Supervised	87.82 (0.43)			71.52 (0.78)			76.55 (0.42)			77.15 (0.17)		
Zero-shot	85.69			55.82			61.88			77.39		
q	0.1	0.3	0.5	0.1	0.3	0.5	0.1	0.3	0.5	0.1	0.3	0.5
PRODEN	87.20 (0.45)	86.30 (0.70)	79.08 (1.52)	60.65 (1.03)	12.10 (1.88)	9.53 (0.66)	73.38 (0.77)	63.23 (0.86)	44.83 (0.76)	76.75 (0.54)	76.17 (0.49)	73.58 (0.47)
+ours	88.17 (0.37)	88.25 (0.92)	88.15 (0.74)	68.20 (0.90)	62.25 (1.03)	57.70 (1.06)	74.58 (0.73)	70.72 (0.59)	68.12 (0.98)	77.85 (0.23)	78.40 (0.23)	78.52 (0.31)
CC	87.97 (0.62)	87.88 (0.31)	87.45 (0.27)	70.10 (0.60)	66.73 (0.26)	62.73 (1.16)	74.67 (0.26)	71.47 (0.35)	68.95 (1.46)	76.75 (0.54)	76.17 (0.49)	75.45 (0.92)
+ours	88.55 (0.45)	88.65 (0.68)	88.62 (0.33)	70.03 (0.45)	65.85 (0.95)	62.42 (0.91)	75.15 (0.56)	71.60 (0.48)	69.73 (0.63)	77.60 (0.36)	77.80 (0.23)	78.05 (0.25)
LW	87.42 (0.41)	86.18 (0.91)	79.85 (1.20)	60.27 (0.28)	12.52 (2.11)	10.85 (2.16)	73.05 (1.28)	62.28 (0.68)	44.75 (0.69)	77.47 (0.44)	76.92 (0.61)	73.05 (0.73)
+ours	88.07 (0.46)	88.50 (0.80)	87.90 (0.75)	68.75 (0.81)	62.25 (0.75)	57.95 (0.98)	74.67 (0.59)	70.70 (0.60)	68.10 (0.74)	77.88 (0.26)	78.10 (0.34)	78.45 (0.32)
PLLCR	86.70 (0.76)	86.78 (0.35)	85.62 (0.94)	66.25 (0.70)	61.33 (0.33)	25.05 (0.70)	74.68 (0.58)	70.22 (0.65)	65.20 (1.00)	76.03 (0.50)	75.70 (0.42)	73.95 (0.67)
+ours	86.52 (0.57)	87.27 (0.66)	87.27 (0.47)	67.05 (0.38)	64.55 (0.73)	59.85 (1.57)	74.20 (0.31)	72.22 (0.69)	69.95 (1.04)	76.12 (0.45)	76.22 (0.47)	76.53 (0.48)
PiCO	85.90 (1.49)	81.70 (1.07)	78.40 (1.86)	64.27 (0.73)	55.15 (0.97)	49.95 (1.51)	73.45 (0.55)	67.23 (0.69)	64.35 (0.73)	74.03 (0.40)	70.22 (1.01)	67.97 (1.08)
+ours	87.65 (1.22)	86.78 (1.34)	87.15 (1.42)	68.55 (0.59)	63.67 (0.81)	59.55 (0.67)	74.80 (0.98)	71.65 (0.77)	69.55 (0.38)	76.65 (0.39)	76.05 (0.61)	75.68 (0.53)
CAVL	87.83 (0.70)	87.88 (1.13)	86.03 (0.85)	69.65 (0.26)	63.75 (0.43)	57.35 (1.27)	73.75 (1.02)	69.30 (1.12)	63.77 (2.89)	76.90 (0.32)	76.75 (0.17)	75.53 (0.43)
+ours	88.82 (0.48)	88.85 (0.27)	88.88 (0.25)	69.85 (0.61)	65.47 (1.15)	61.65 (0.78)	74.57 (0.63)	71.78 (0.48)	69.02 (0.78)	77.50 (0.16)	78.00 (0.16)	78.10 (0.23)

4.3 Ablation Studies

In this part, we conduct experiments to assess the effectiveness of our framework with different handcrafted prompts, learnable context lengths, visual backbones, and balancing strategies in 3 datasets: DTD, FGVC Aircraft and Caltech101.

Impact of Different Handcrafted Prompts. Because our framework incorporates handcrafted prompts, it is crucial to determine the performance of our framework with differently crafted prompts. In Table 2, we design some prompts with zero-shot performance lower than the default handcrafted prompts [20] and evaluate our framework with these prompts at different levels of label ambiguity. *Introducing a handcrafted prompt can consistently improve the performance, regardless of the quality of the prompt.* Better handcrafted prompts have resulted in better test accuracy in our experiments, which indicates that the handcrafted prompts indeed guide the fine-tuning process with candidate labels. Nonetheless, the performance is improved by a large margin even when we simply set the handcrafted prompt “<CLS>.”.

Table 2. Performance comparison of our framework for different handcrafted prompts.

Dataset	Handcrafted Prompt	q			CLIP-Zeroshot
		0.1	0.3	0.5	
DTD	without handcrafted prompt (only prompt learning)	61.70±1.07	55.08±1.22	48.08±1.58	-
	“a photo of a <CLS>.”	62.15±0.63	56.85±0.63	53.20±0.91	40.36
	“<CLS>.”	62.45±0.18	56.67±0.98	53.23±1.20	41.13
	“<CLS> texture.”	62.67±0.92	58.42±0.41	55.08±0.93	42.79
FGVC Aircraft	without handcrafted prompt (only prompt learning)	23.60±2.16	18.27±2.83	14.75±3.11	-
	“a photo of a <CLS>.”	24.70±1.16	21.82±1.42	19.60±0.81	15.84
	“<CLS>.”	24.65±0.89	20.85±0.84	19.02±0.71	15.54
	“a photo of a <CLS>, a type of aircraft.”	25.12±1.07	22.18±1.00	20.32±0.71	17.07
Caltech101	without handcrafted prompt (only prompt learning)	90.47±0.19	88.75±0.74	87.75±0.76	-
	“<CLS>.”	91.42±0.19	90.97±0.51	90.05±0.73	81.34
	“a photo of a <CLS>.”	91.92±0.11	91.33±0.19	90.72±0.15	85.84

Effectiveness of the Dynamic Mixing Strategy. This experiment is conducted to evaluate the efficacy of our dynamic and mixing strategy. $\alpha(t)$ is set using the same hyperparameters as the main experiment. The result is shown in Table 3. *The dynamic mixing strategy optimizes the model’s performance.* The model with dynamic technique outperforms all others in 6 and ranked second in 2, out of 9 cases, demonstrating its effectiveness. When only one handcrafted or learnable prompt is leveraged, the model’s performance declines and cannot surpass the performance achieved when $\alpha = 0.7$. This validates the rationality of the mixing strategy. Lastly, when $\alpha = 0$, it outperforms models with only learnable prompt guidance ($\alpha = 1$) or no guidance (vanilla prompt learning) in most

Table 3. Performance comparison of our framework for different α .

	DTD			FGVCAircraft			Caltech101		
q	0.1	0.3	0.5	0.1	0.3	0.5	0.1	0.3	0.5
$\alpha = 0$	62.05 \pm 1.17	57.12 \pm 0.98	53.75 \pm 1.41	24.68 \pm 0.83	22.10 \pm 0.56	19.93 \pm 0.31	91.70 \pm 0.23	91.20 \pm 0.16	90.93\pm0.22
$\alpha = 1$	61.83 \pm 0.40	56.45 \pm 2.29	50.98 \pm 3.08	25.15 \pm 0.98	18.83 \pm 4.24	15.35 \pm 3.76	91.62 \pm 0.25	89.83 \pm 0.68	88.47 \pm 0.37
$\alpha = 0.7$	62.28 \pm 0.33	59.38\pm1.43	54.57 \pm 1.14	25.77\pm0.99	22.02 \pm 0.98	19.95 \pm 0.38	91.85 \pm 0.23	91.03 \pm 0.54	90.35 \pm 0.63
$\alpha(t)$	62.67\pm0.92	58.42 \pm 0.41	55.08\pm0.93	25.12 \pm 1.07	22.18\pm1.00	20.32\pm0.71	91.92\pm0.11	91.33\pm0.19	90.72 \pm 0.15

cases, which, to some extent, justifies the effectiveness of explicitly adopting the handcrafted prompt to guide the learning process.

5 Conclusion

This work, for the first time, investigated the scenario when tuning vision-language models (VLMs) with candidate labels. Throughout a series of experiments, we empirically demonstrated that prompt learning combined with PLL training objectives in a vanilla way can learn from candidate labels. However, as the ambiguity of candidate labels increases, its performance is degraded. To alleviate this issue, we proposed a framework that enhances the robustness by aligning the dynamically mixed class posterior of the handcrafted and learnable prompt with the model’s output to guide the learning process with candidate labels. Comprehensive experimental results and analysis on multiple benchmark datasets demonstrate the effectiveness of our proposed framework.

References

1. Bossard, L., Guillaumin, M., Van Gool, L.: Food-101—mining discriminative components with random forests. In: ECCV. pp. 446–461. Springer (2014)
2. Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., Vedaldi, A.: Describing textures in the wild. In: CVPR. pp. 3606–3613 (2014)
3. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: CVPR. pp. 248–255. Ieee (2009)
4. Fei-Fei, L., Fergus, R., Perona, P.: Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In: CVPR. pp. 178–178. IEEE (2004)
5. Feng, L., Lv, J., Han, B., Xu, M., Niu, G., Geng, X., An, B., Sugiyama, M.: Provably consistent partial-label learning. NeurIPS pp. 10948–10960 (2020)
6. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR. pp. 770–778 (2016)
7. Hu, J., Shu, S., Li, B., Xiang, T., He, Z.: An unbiased risk estimator for partial label learning with augmented classes. ACM Trans. Intell. Syst. Technol. **15**(6) (2024)
8. Huang, T., Chu, J., Wei, F.: Unsupervised prompt learning for vision-language models. arXiv preprint arXiv:2204.03649 (2022)
9. Jia, M., Tang, L., Chen, B.C., Cardie, C., Belongie, S., Hariharan, B., Lim, S.N.: Visual prompt tuning. In: ECCV. pp. 709–727. Springer (2022)

10. Khattak, M.U., Rasheed, H., Maaz, M., Khan, S., Khan, F.S.: Maple: Multi-modal prompt learning. In: CVPR. pp. 19113–19122 (2023)
11. Krause, J., Stark, M., Deng, J., Fei-Fei, L.: 3d object representations for fine-grained categorization. In: ICCV Workshops. pp. 554–561 (2013)
12. Li, B., et al.: Asyco: An asymmetric dual-task co-training model for partial-label learning (2024), <https://arxiv.org/abs/2407.15036>
13. Li, B., et al.: Gemini: A dual-task co-training model for partial label learning. In: AI 2023: Advances in Artificial Intelligence. Springer Nature Singapore (2024)
14. Liu, L., Dietterich, T.: A conditional multinomial mixture model for superset label learning. *NeurIPS* **25** (2012)
15. Luo, J., Orabona, F.: Learning from candidate labeling sets. *NeurIPS* **23** (2010)
16. Lv, J., Xu, M., Feng, L., Niu, G., Geng, X., Sugiyama, M.: Progressive identification of true labels for partial-label learning. In: ICML. pp. 6500–6510 (2020)
17. Maji, S., Rahtu, E., Kannala, J., Blaschko, M., Vedaldi, A.: Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151* (2013)
18. Parkhi, O.M., Vedaldi, A., Zisserman, A., Jawahar, C.: Cats and dogs. In: CVPR. pp. 3498–3505. IEEE (2012)
19. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library (2019)
20. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: ICML. pp. 8748–8763 (2021)
21. Soomro, K., Zamir, A.R., Shah, M.: Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402* (2012)
22. Tang, C.Z., Zhang, M.L.: Confidence-rated discriminative partial label learning. In: AAAI (2017)
23. Wang, H., Xiao, R., Li, Y., Feng, L., Niu, G., Chen, G., Zhao, J.: Pico: Contrastive label disambiguation for partial label learning. In: ICLR (2021)
24. Wen, H., Cui, J., Hang, H., Liu, J., Wang, Y., Lin, Z.: Leveraged weighted loss for partial label learning. In: ICML. pp. 11091–11100. PMLR (2021)
25. Wu, C.E., et al.: Why is prompt tuning for vision-language models robust to noisy labels? In: ICCV. pp. 15488–15497 (2023)
26. Wu, D.D., Wang, D.B., Zhang, M.L.: Revisiting consistency regularization for deep partial label learning. In: ICML. pp. 24212–24225. PMLR (2022)
27. Yao, Y., Gong, C., Deng, J., Yang, J.: Network cooperation with progressive disambiguation for partial label learning. In: ECML PKDD. Springer (2021)
28. Zhang, F., Feng, L., Han, B., Liu, T., Niu, G., Qin, T., Sugiyama, M.: Exploiting class activation value for partial-label learning. In: ICLR (2021)
29. Zhang, J., Wei, Q., Liu, F., Feng, L.: Candidate pseudolabel learning: Enhancing vision-language models by prompt tuning with unlabeled data. In: ICML (2024)
30. Zhang, Z., He, S., Shen, B., Feng, L.: Defending multimodal backdoored models by repulsive visual prompt tuning. *arXiv preprint arXiv:2412.20392* (2024)
31. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Conditional prompt learning for vision-language models. In: CVPR. pp. 16816–16825 (2022)
32. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Learning to prompt for vision-language models. *IJCV* **130**(9), 2337–2348 (2022)