# GalLoP: Learning Global and Local Prompts for Vision-Language Models

Marc Lafon<sup>\*1</sup><sup>©</sup>, Elias Ramzi<sup>\*1</sup><sup>©</sup>, Clément Rambour<sup>1</sup><sup>©</sup>, Nicolas Audebert<sup>1,2</sup><sup>©</sup>, and Nicolas Thome<sup>3</sup><sup>©</sup>

<sup>1</sup> Conservatoire national des arts et métiers, CEDRIC, F-75141 Paris, France

 $^2~$  Univ. Gustave Eiffel, ENSG, IGN, LASTIG, F-94160 Saint-Mandé, France

<sup>3</sup> Sorbonne Université, CNRS, ISIR, F-75005 Paris, France {marc.lafon, elias.ramzi}@cnam.fr

Abstract. Prompt learning has been widely adopted to efficiently adapt vision-language models (VLMs), e.g. CLIP, for few-shot image classification. Despite their success, most prompt learning methods trade-off between classification accuracy and robustness, e.q. in domain generalization or outof-distribution (OOD) detection. In this work, we introduce Global-Local Prompts (GalLoP), a new prompt learning method that learns multiple diverse prompts leveraging both global and local visual features. The training of the local prompts relies on local features with an enhanced vision-text alignment. To focus only on pertinent features, this local alignment is coupled with a sparsity strategy in the selection of the local features. We enforce diversity on the set of prompts using a new "prompt dropout" technique and a multiscale strategy on the local prompts. GalLoP outperforms previous prompt learning methods on accuracy on eleven datasets in different few shots settings and with various backbones. Furthermore, GalLoP shows strong robustness performances in both domain generalization and OOD detection, even outperforming dedicated OOD detection methods. Code and instructions to reproduce our results will be open-sourced.

Keywords: Vision-language models  $\cdot$  Few shot classification  $\cdot$  Prompt learning  $\cdot$  Local and global prompts  $\cdot$  Robustness  $\cdot$  OOD detection

# 1 Introduction

Vision-Language Models (VLMs), *e.g.* CLIP [35] or ALIGN [20], have shown impressive performances for zero-shot image classification. Prompt learning [3,21,22, 26,33,51,52] has been among the leading approaches to efficiently adapt VLMs to a specific downstream dataset. These methods train a learnable context in the form of *soft prompts* to optimize the text/image alignment. Prompt learning methods

<sup>\*</sup> Equal contribution.



**Fig. 1:** Our GalLoP method demonstrates excellent performances in accuracy plus robustness, *i.e.* out-of-distribution detection (a) and domain generalization (b), while state-of-the-art prompt learning methods compromise between these aspects. Additionally, unlike recent methods utilizing ineffective local zero-shot CLIP features, GalLoP learns discriminative local prompts precisely aligned with sparse image regions at various scales, facilitating the discriminability between classes. GalLoP integrates both global and local prompts, with their diversity explicitly enforced during few-shot learning, which significantly enhances the performance of their combination (c).

benefit from the strong generalization capability of VLMs' textual encoder and are effective even when only a few labeled examples are available.

Despite their success, we observe that these methods trade off between classification accuracy and robustness. This is illustrated on Fig. 1(a), where methods exhibiting the best accuracy sacrifice out-of-distribution (OOD) detection performances, *e.g.* PromptSRC [22], while those excelling in OOD detection often have poor accuracy results, *e.g.* LoCoOp [29]. A similar observation is done in domain generalization (DG), see Fig. 1(b): PromptSRC [22] presents two different versions, one optimized for accuracy (PromptSRC<sup> $\triangleright$ </sup>) and the other for domain generalization (PromptSRC<sup> $\diamond$ </sup>), highlighting the intrinsic conflict between both criteria.

To boost classification accuracy, prompt learning can involve learning multiple prompts [1] to emulate "prompt ensembling", *e.g.* prompts specialized for specific classes [33,45] or Transformer's layers [21,22], or casting multiple prompts learning within a probabilistic framework [26]. The key challenge in prompt ensembling lies in learning diverse prompts to optimize the combination. However, since these approaches only operate on global visual representations, they cannot utilize diverse prompts aligned with specific image regions to maximize their diversity.

Recently, attempts have been made to use local image representations in prompt learning, e.g. LoCoOp [29] or PLOT [3]. Although these approaches are promising, their performances in accuracy/robustness are suboptimal compared to state-of-the-art results, see Fig. 1(a),(b). Their limited performances stem from two main factors: i) they use "dense" (*i.e.* all) local features from CLIP, which includes irrelevant or noisy regions for a given concept, and ii) these local features are not as well aligned with the text due to CLIP's pre-training with the global representation. In consequence, the performance of prompts trained with those local features is much lower than their global counterpart, and this degradation affects performances when combined with global, as illustrated in Fig. 1(c).

In this paper, we introduce Global-Local Prompts (GalLoP), a new method to learn a diverse set of prompts by leveraging both global and local visual representations. GalLoP learns sparse discriminative local features, *i.e.* text prompts are aligned to a sparse subset of regions at multiple scales. This enables fine-grained and accurate text-to-image matching, making GalLoP local prompts highly competitive. Moreover, we train GalLoP with diverse global and local prompts, unlocking the complementarity between both sets and significantly improving their combination, as shown in Fig. 1(c).

To achieve this, GalLoP relies on two main methodological contributions:

- Effective local prompts learning. In GalLoP, we propose to align local prompts with sparse subsets of k image regions, enabling text-to-image matching that captures fine-grained semantics. To adapt visual representations to the downstream dataset, we refine the textual alignment of visual local features by employing a simple linear projection amenable to few-shots learning.
- Enforcing ensemble diversity. We learn both global prompts aligned with the whole image and local spatially-localized prompts, and enforce diversity between them to improve their combination. We induce diversity through randomization using a new "prompt dropout" strategy, which enhances generalization when learning multiple prompts. Additionally, we employ a multiscale strategy to align local prompts with image regions of varying sizes, capturing different visual aspects of a concept's semantics.

We conduct an extensive experimental validation of GalLoP on 11 few-shot image classification datasets and 8 datasets evaluating robustness. We show that GalLoP outperforms state-of-the-art prompt learning methods on classification accuracy, OOD detection, and domain generalization, therefore improving the observed tradeoff in these 3 criteria. We validate that our two main contributions, *i.e.* learning strong local prompts and diverse representations, are essential for reaching excellent performances.

# 2 Related work

**Prompt learning.** Prompt learning has emerged as an efficient way to adapt VLMs to downstream datasets. These methods, *e.g.* CoOp [52], learn *soft prompts* to adapt CLIP textual features to specific labels without the need for a cumbersome step of "prompt engineering" as performed in [35]. Following these seminal works, many variants have been proposed. [51] uses a meta-network to bias the learnable prompt using the global visual representation of the input image. To boost prompt learning performances, recent works have focused on learning multiple prompts [3,21,22,26]. MaPLe [21] introduces prompts in several layers of both textual and visual encoders. PromptSRC [22] builds upon this work by introducing several regularization losses, boosting both accuracy and robustness performances. We note that PromptSRC uses a set of hand-crafted prompts to regularize the learning of the textual prompts,

which is not fully aligned with the initial motivation behind prompt learning. Furthermore, both MaPLe and PromptSRC are limited to the use of vision transformer architectures. ProDA [26] models the distribution over the textual representation of classes using a multivariate Gaussian distribution, and indirectly learns the distribution over prompts using a surrogate loss. PromptStyler [4] learns several prompts that represent different "styles" to perform source-free domain generalization. These two approaches achieve prompt diversity by enforcing orthogonality among the prompts. In GalLoP, we induce diversity with a "prompt dropout" technique, which randomly drops subsets of prompts during training, thus avoiding the introduction of an additional loss, while limiting prompt over-fitting observed in [22]. We further improve diversity by specializing the local prompts on different image scales to align them with different sets of attributes for each class.

**Prompt learning using visual local features.** There has been a growing interest in leveraging CLIP's local features in prompt learning methods [3, 29, 41]. PLOT [3] learns a set of prompts by using the optimal transport (OT) [44] distance between them and the set of local features, which is prohibitive to compute. Furthermore, the OT distance enforces the prompts to use information from all local visual features during training, including possibly detrimental ones. Also, PLOT adds the global visual features to the local features to achieve strong results on the ImageNet dataset. In GalLoP, we use a sparse mechanism to learn localized prompts. This removes the negative influence of background features while being computationally efficient. Finally, GalLoP learns prompts from the local features without any access to CLIP's original global visual feature. LoCoOp [29] introduced an entropy loss leveraging "irrelevant" local visual features in an outlier exposure fashion [18] to improve out-of-distribution detection but at the expense of accuracy. [41] introduces a method specifically designed for multi-label classification, which learns prompts using local visual features. While these methods obtained promising results, we show in this work that their performance is intrinsically limited by the lower discriminative power of CLIP's zero-shot local visual features.

**Prompt learning & robustness.** As VLMs are becoming increasingly prevalent in few-shot classification applications, their robustness capabilities, such are OOD detection or domain generalization (DG), are receiving increasing attention. To address OOD detection for VLMs, [28] proposed the maximum concept matching (MCM). In [30], the authors leveraged global and local visual information to construct the GL-MCM score to improve OOD detection results. In addition to LoCoOp [29], other recent works [31] uses negative prompts to perform OOD detection by "Learning to Say No" (LSN). Furthermore, prompt learning methods are also evaluated on DG tasks where the prompt is learned on a source dataset (*e.g.* Imagenet) and tested on a domain-shifted target dataset (*e.g.* Imagenet-Sketch). CoOp [52] achieves significantly better DG results than CLIP, and more recent prompt learning methods like [21, 22] continue to improve over CoOp's strong performance. By unlocking the potential of local visual features to learn prompts, GalLoP further improves over the state-of-the-art for top-1 accuracy, OOD detection and domain generalization.



**Fig. 2: Illustration of GalLoP.** GalLoP learns a diverse set of global prompts and local prompts. Pertinent local prompts are learned using only the most relevant regions of the image for each class. We further improve the limited text-vision alignment of CLIP's local features using a simple linear layer. The diversity is encouraged using a new "prompt dropout" technique for global prompts, and a multiscale loss for local prompts.

## 3 Combining global and local prompts with GalLoP

In this section, we describe our proposed method, GalLoP, which seeks to learn an ensemble of diverse prompts from both global and local CLIP's visual representations. As illustrated in Fig. 2, GalLoP learns two specialized sets of prompts: the "global prompts" receiving a signal from the global visual representation, and the "local prompts" trained using local features only.

Formally, let us consider a set of n learnable local prompts  $\mathcal{P}_l = (\mathbf{p}_1^l, \cdots, \mathbf{p}_n^l)$  and a set of m learnable global prompts  $\mathcal{P}_g = (\mathbf{p}_1^g, \cdots, \mathbf{p}_m^g)$ . Each of these prompt  $\mathbf{p}$  is composed of V learnable embeddings, *i.e.*  $\mathbf{p} := [p^1, \dots, p^V] \in \mathbb{R}^{V \times d'}$ , and are prepended to the class name embeddings  $\mathbf{c}$  to perform classification. Let  $\mathcal{D} = \{(\mathbf{x}, y)\}$  denote the downstream dataset, where  $\mathbf{x}$  is an image and y its class, and let  $\mathcal{T}$  and  $\mathcal{V}$ denote CLIP's text and vision encoder, respectively. The textual encoder produces a normalized textual representation  $\mathbf{t}_c = \mathcal{T}([\mathbf{p}, \mathbf{c}]) \in \mathbb{R}^d$  of the  $c^{th}$  class. Given the input image  $\mathbf{x}$ , the visual encoder produces a visual representation  $\mathbf{z}$ .  $\mathbf{z}$  can be a global vector for learning global prompts, *i.e.* the global visual feature on which CLIP has been pre-trained. For local prompts,  $\mathbf{z}$  will be a set of localized features outputted by the encoder. From its visual representation  $\mathbf{z}$ , the probability for the image  $\mathbf{x}$  to be classified into the class  $y_c$  can be expressed as:

$$p(y = y_c | \boldsymbol{x}; \boldsymbol{p}) = \frac{\exp(\sin(\boldsymbol{z}, \boldsymbol{t}_c) / \tau)}{\sum_{c'} \exp(\sin(\boldsymbol{z}, \boldsymbol{t}_{c'}) / \tau)},$$
(1)

6 M. Lafon et al.

where  $sim(\cdot, \cdot)$  is a measure of similarity, and  $\tau$  is fixed a temperature scaling parameter. With this general definition of the probability in Eq. (1), we can train a prompt  $\boldsymbol{p}$  using the standard cross-entropy loss  $\mathcal{L}_{CE}(\boldsymbol{p}(\boldsymbol{y}=\boldsymbol{y}_{c}|\boldsymbol{x};\boldsymbol{p}))$ .

In Sec. 3.1, we introduce a relevant similarity measure  $sim(z, t_c)$  for implementing Eq. (1) on local prompts. We rely on a sparsification strategy that only considers a small subset of class-relevant regions of the image. Furthermore, we use a linear projection to improve the vision-text alignment of local features, thus enhancing the quality of the learned prompts. In Sec. 3.2 we describe how we learn a diverse set of global and local prompts, whose combination can improve predictions' performance. We introduce "prompt dropout" to increase the diversity of global prompts by randomly selecting a subset of prompts for each image. Finally, we introduce a multiscale loss by dedicating each local prompt to select different sub-region sizes of the input image.

#### 3.1 Learning prompts from local visual representations

In this section, we temporarily consider a single local prompt  $\mathbf{p}_j^l \in \mathcal{P}_l$  without loss of generality. In this case, the visual representation  $\mathbf{z}$  that we consider is the set of visual local features, *i.e.*  $\mathbf{z} = \mathcal{Z}_l \in \mathbb{R}^{L \times d}$ , obtained following [7] (see details in supplementary Sec. A.1). Here, we can not directly compute the probability of Eq. (1) as we need to define the similarity between the set of vectors  $\mathcal{Z}_l = (\mathbf{z}_1^1, \dots, \mathbf{z}_L^1)$  and the textual representation of the  $c^{th}$  class,  $\mathbf{t}_c = \mathcal{T}([\mathbf{p}_l^i, \mathbf{c}])$ .

**Sparse local similarity.** A naive way to obtain a single similarity for all regions is to average the similarities of each spatial location with the textual representation of the class. However, a substantial portion of the local features are



Fig. 3: GalLoP sparse local similarity  $sim(\mathcal{Z}_l, t_c)$  between class prompt  $t_c$  and visual features  $\mathcal{Z}_l$  is the average of the top-k highest similarities (here, k=3).

irrelevant to the class, *e.g.* features from background areas, which may introduce noise and perturb the learning process. To solve this problem, we adopt a sparse approach, where only local features semantically related to the class are kept to perform classification. As illustrated in Fig. 3, we select the top-k local features with the highest similarities with the prompted class textual representation, and average their similarities to measure  $sim(\mathcal{Z}_l, t_c)$ .

Formally, we define the similarity between a prompt  $t_c$  and the set of visual features  $\mathcal{Z}_l$  as the average similarity for the k most similar regions:

$$\operatorname{sim}_{\operatorname{top-}k}(\mathcal{Z}_l, \boldsymbol{t}_c) \coloneqq \frac{1}{k} \sum_{i=1}^{L} \mathbb{1}_{\operatorname{top-}k}(i) \cdot \langle \boldsymbol{z}_i^l, \boldsymbol{t}_c \rangle$$
where
$$\mathbb{1}_{\operatorname{top-}k}(i) = \begin{cases} 1 & \text{if } \operatorname{rank}_i(\langle \boldsymbol{z}_i^l, \boldsymbol{t}_c \rangle) \leq k, \\ 0 & \text{otherwise.} \end{cases}$$
(2)

which we plug into Eq. (1) to compute the probability for class c. We show in Sec. 4.3 that relying on sparsity is mandatory for local prompt learning, boosting performances by almost 20pt in top-1 accuracy.

Improving local text-vision alignment. While previous works [29, 41, 50] have exploited the text-vision alignment of CLIP's local features, we empirically verified in Sec. 4.3 that using these features leads to poor zero-shots classification results on ImageNet. This is expected, as CLIP is pre-trained to align the global visual features with its textual representation. Local features are thus suboptimal to learn effective prompts for image classification. Motivated by this observation, we propose to improve the discriminative power of CLIP's local visual features by realigning them with the textual representations of the class labels of the downstream dataset. To do so, we propose to use a simple linear projection  $h_{\theta}$ . To ease the learning process, we initialize the linear layer  $h_{\theta}$  to identity, so that the initial features are close to CLIP's representations. Henceforth, we use the set of linearly transformed local visual features  $h_{\theta}(\mathcal{Z}_l)$  to compute the probability of Eq. (1), which becomes:

$$p(y = y_c | \boldsymbol{x}; \boldsymbol{p}_j^l, k, \boldsymbol{\theta}) = \frac{\exp(\sin_{\text{top-}k}(h_{\boldsymbol{\theta}}(\mathcal{Z}_l), \boldsymbol{t}_c) / \tau)}{\sum_{c'} \exp(\sin_{\text{top-}k}(h_{\boldsymbol{\theta}}(\mathcal{Z}_l), \boldsymbol{t}_{c'}) / \tau)}.$$
(3)

Thus, a local prompt can be optimized by maximizing this probability with the cross-entropy loss. These design choices in GalLoP allow us to train a powerful classifier for local features: the sparsity helps to focus on the most relevant regions of an image and to remove potential background noise, while the linear projection enhances the text-vision alignment and boosts the fine-grained discriminating power of the local features. We study these design choices in Sec. 4.3.

#### 3.2 Learning multiple diverse prompts

In this section, we describe how we induce diversity among the learned prompts. Besides exploiting different sources of information – the global and visual ones –, we introduce two mechanisms to increase diversity: "prompt dropout" and multiscale training.

To train our set of global prompts, we can simply train each one of them independently using the cross-entropy loss as in CoOp [52]. However, this strategy will necessarily produce identical global prompts, losing the advantage of using an ensemble of prompts. A possible strategy to avoid this behavior is to use an explicit diversity loss inducing a semantic orthogonality between the different global prompts. Instead of adding a another loss term, which can disturb the training process, we chose to take a different approach.

**Prompt dropout.** Motivated by the success of the "dropout" [10, 40] technique classically used in deep learning, we introduce "prompt dropout" into the prompt learning framework. In "prompt dropout", we randomly mask a subset of prompts for each image of the batch. Alternatively, from the perspective of each prompt, we select a different subset of the batch of images, thus inducing diversity in the learning process of the prompts through input randomization (see Fig. 4(a)).



**Fig. 4:** (a) Prompt dropout induces diversity by randomly selecting different subsets of prompts for each image of the batch. In (a), each image will be used by half the prompts. (b) To learn diverse local prompts, we specialize each one of them using a different number of regions, and therefore a different level of sparsity.

Formally, the loss used to train our global prompts  $\mathcal{P}_g$  with prompt dropout can be expressed as:

$$\mathcal{L}_{\text{global}}(\mathcal{P}_g) = \mathbb{E}_{\boldsymbol{x}, y} \sum_{i=1}^{m} \delta_i(\boldsymbol{x}) \mathcal{L}_{\text{CE}}(p(y|\boldsymbol{x}; \boldsymbol{p}_i^g))$$
(4)

where  $\delta_i(\boldsymbol{x}) \sim \text{Bernoulli}(1-r)$  with r the dropout rate.

**Multiscale training.** To specifically improve the diversity of the local prompts, we specialize each local prompt to select a different number of class-specific visual patches (scales). In this way, prompts dedicated to small scales will get more signals from classes corresponding to small visual concepts, *e.g.* "daisy flower" or "tailed frog", while prompts learned with larger scales will receive more signals from images with wider concepts, *e.g.* "castle" or "valley". More formally, let  $(k_1, k_1 + \Delta_k, \dots, k_1 + (n-1) \cdot \Delta_k)$  denote a set of increasing scales with  $k_1$  the first scale and  $\Delta_k$  the expansion factor. Each local prompt  $p_j^l$  will be learned with its associated scale  $k_j = k_1 + (j-1) \cdot \Delta_k$ .

The training of our n local prompts is performed by optimizing the probability defined in Eq. (3) for each prompt with a different scale, *i.e.* value of k:

$$\mathcal{L}_{\mathbf{x}\text{-scale}}(\mathcal{P}_l, \boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{x}, y} \sum_{j=1}^{n} \mathcal{L}_{\text{CE}}(p(y | \boldsymbol{x}; \boldsymbol{p}_j^l, \boldsymbol{\theta}, k_j)).$$
(5)

The overall loss to train our set of prompts  $\mathcal{P} = \mathcal{P}_l \cup \mathcal{P}_g$  is the sum of the local multiscale and global losses:

$$\mathcal{L}_{\text{total}}(\mathcal{P}, \boldsymbol{\theta}) = \mathcal{L}_{\text{global}}(\mathcal{P}_g) + \mathcal{L}_{\text{x-scale}}(\mathcal{P}_l, \boldsymbol{\theta}) \tag{6}$$

Inference. We perform averaging of the similarities obtained with each prompt to obtain a final similarity,  $\sin(z, t_c)$ , for each class which is used as logit for classification. For OOD detection we use the GL-MCM score [30]. The inference procedures are detailed in Sec. A.2.

# 4 Experimental results

In this section, we present the experimental validation of GalLoP. We first show, in Sec. 4.1, that GalLoP outperforms previous methods on top-1 accuracy on a collection of 11 datasets used in [52] with ViT-B/16 [8]. We also show that GalLoP performs well for different few shot settings on ImageNet and with a ResNet-50 [14]. In Sec. 4.2, we compare robustness performances of GalLoP and other prompts learning methods in domain generalization and OOD detection, and show that GalLoP has better trade-off with top-1 accuracy contrary to previous methods. In Sec. 4.3, we conduct ablation studies of the different components of GalLoP.

Implementation details. We experiment with both ResNet-50 and ViT-B/16 CLIP models. When not specified, we use ViT-B/16. We train for 50 epochs on ImageNet and 200 epochs for other datasets with SGD, a learning rate of 0.002 decayed using cosine annealing and a weight decay of 0.01, following the setting of [52]. Unless specified otherwise, we train the models using 16 shots. Our base parameters for GalLoP are as follows: m=4 global prompts with a dropout of r=75% (in practice we keep a single prompt for each image), n=4 local prompts with scales  $k_1=10$  and  $\Delta_k=10$  for ViT-B/16 and  $k_1=5$  and  $\Delta_k=5$  for ResNet-50 as there are fewer local patches. We keep  $\tau$  fixed from CLIP.

Baselines. We compare GalLoP to recent prompt learning methods. Including, single prompt learning CoOp and Co-CoOp. Multi-prompt learning MaPLe, ProDA, PLOT, PromptSRC. We denote by PromptSRC<sup>b</sup> the version designed for accuracy and PromptSRC<sup>o</sup> the version designed for domain generalization. We also include OOD detection specific methods such as LoCoOp and LSN.

**Table 1:** Top-1 accuracy with ViT-B/16 backbone. Comparison of GalLoP to other prompt learning methods on several standard benchmarks. <sup> $\dagger$ </sup> results based on our own re-implementation.

Dataset	ImageNet for	Caltech101 for	OxfordPets 12.	C <sub>ars</sub> [23]	Flowers102 120	Foodlo1 [2]	$A_{ircraft}$ [27]	SUN397 1471	$^{DTD}_{DT}$	$E_{uroSAT/I51}$	$v_{CF_{101}}$	$A_{Verage}$
CLIP [35]	66.7	92.2	88.4	65.5	70.7	84.8	24.8	62.3	44.1	48.3	64.7	75.7
Linear Probe	67.3	95.4	85.3	80.4	97.4	82.9	45.4	73.3	70.0	87.2	82.1	78.8
CoOp [52]	71.7	95.6	91.9	83.1	97.1	84.2	43.4	74.7	69.9	84.9	82.2	79.9
Co-CoOp [51]	71.0	95.2	93.3	71.6	87.8	87.2	31.2	72.2	63.0	73.3	78.1	74.9
MaPLe [21]	72.3	<u>96.0</u>	92.8	83.6	97.0	85.3	48.4	75.5	71.3	92.3	85.0	81.8
PLOT [3]	72.6	<u>96.0</u>	<u>93.6</u>	84.6	<u>97.6</u>	87.1	46.7	76.0	71.4	92.0	85.3	82.1
PromptSRC <sup>▷</sup> [22]	73.2	96.1	93.7	85.8	97.6	86.5	50.8	77.2	72.7	<b>92.4</b>	86.5	82.9
LoCoOp <sup>†</sup> [29]	71.5	94.9	92.4	79.8	96.3	84.7	40.7	74.2	69.5	86.1	81.6	79.2
$ProDA^{\dagger}$ [26]	71.9	95.5	93.5	79.8	96.8	86.8	40.2	75.7	70.9	85.1	83.3	80.0
GalLoP	75.1	96.7	94.1	89.2	98.8	86.5	58.3	77.2	75.5	<u>90.1</u>	86.9	84.4



Fig. 5: Results on ImageNet with different few shot settings Fig. 5a, and ResNet-50 Fig. 5b.

#### 4.1 Main in-distribution results.

On Tab. 1, we compare GalLoP with a ViT-B/16 backbone on a suite of 11 datasets, a standard benchmark for prompt learning methods. On average, GalLoP outperforms previous methods by a large margin with +1.5pt compared to PromptSRC<sup> $\triangleright$ </sup> the next best performing method. Furthermore, GalLoP performs well on most datasets, achieving state-of-the-art among prompt learning methods. For instance, on the large-scale ImageNet dataset, it outperforms PLOT by +2.5pt and PromptSRC<sup> $\triangleright$ </sup> by +1.9pt. On some datasets, *e.g.* FGVC Aircraft, GalLoP outperforms the next best method by a large margin, with +7.5pt compared to PromptSRC<sup> $\triangleright$ </sup>.

We then compare GalLoP on Fig. 5a to prompt learning methods in different few-shot settings on ImageNet. GalLoP performs well in all configurations, outperforming for each setting the very competitive method, PromptSRC<sup> $\triangleright$ </sup>. Finally, in Fig. 5b we show that GalLoP works well with a ResNet-50, outperforming PLOT and CoOp by +3.1pt. Note that compared to other methods, *e.g.* MaPLe and PromptSRC, GalLoP is amenable to both convolutional and transformer vision backbones. Detailed results for ResNet-50 can be found in the supplementary material B.2

#### 4.2 Robustness results.

In this section, we compare the robustness performances of GalLoP *vs*. other prompt learning methods, see Fig. 6, on domain generalization and OOD detection. For both benchmarks, models are trained on ImageNet (16 shots).

**Domain generalization results.** We compare on Fig. 6a the domain generalization performances of GalLoP vs. other prompt learning methods. After being trained on ImageNet (16 shots), the models are evaluated on top-1 accuracy for different domains with the same classes as ImageNet, *i.e.* ImageNet-V2 [36], ImageNet-Sketch [46], ImageNet-A [19] and ImageNet-R [16]. GalLoP outperforms the domain-generalization specific method PromptSRC<sup> $\diamond$ </sup> by +0.5pt on average,



Fig. 6: GalLoP robustness performances. GalLoP achieves strong performances on domain generalization Fig. 6a and on OOD detection Fig. 6b ImageNet benchmarks, while outperforming prompt learning methods on top-1 accuracy (detailed results in Tab. 7 and Tab. 8).

while outperforming it by +4.9pt on ImageNet. This illustrates the trade-off made by PromptSRC between top-1 accuracy and domain generalization. Indeed, Gal-LoP outperforms PromptSRC<sup> $\triangleright$ </sup>, designed for ImageNet accuracy, by +1.9pt on ImageNet and +1.5pt on average in domain generalization. GalLoP achieves the best trade-off between top-1 performances and domain generalization. The detailed results can be found in supplementary material B.4

**Results on OOD detection.** In OOD detection the models must recognize between in-distribution examples (ImageNet test set) and different OOD datasets, namely iNaturalist [43], SUN [47], Places [49] and Textures [5], a standard benchmark in the OOD detection literature. We plot on Fig. 6b the average results on the ImageNet OOD benchmark of GalLoP and other prompt learning methods measured in FPR95 (lower is better,  $\downarrow$ ). GalLoP outperforms traditional prompt learning methods, *e.g.* CoOp -3pt FPR95, as well as dedicated OOD detection methods, *e.g.*-1.4pt FPR95 *vs.* LoCoOp or -2.9pt FPR95 *vs.* LSN. Meanwhile, GalLoP also outperforms both LSN and LoCoOp by a large margin in top-1 accuracy, *i.e.* +3.2pt and +3.6pt respectively. The detailed results can be found in the supplementary material B.5

#### 4.3 Ablation studies.

In this section, we investigate – on ImageNet 16 shots – the design choices for GalLoP. We first show how GalLoP leverages the complementarity of strong global and local prompts to boost performances Tab. 2. We then demonstrate the benefit of sparsity and local alignment in Fig. 7. Finally, we show the impact of our choice when learning multiple prompts for both global and local features Fig. 8.

#### 12 M. Lafon et al.

Combining global and local features. On Tab. 2, we show that leveraging global and local features requires some important design choices. Indeed, we experiment with a baseline using CoOp on local features ("CoOp<sub>Local</sub>"), learning a single prompt, without sparsity and no alignment. This baseline already outperforms using zero-shot local features, +28.7pt top-1. However, its combination with a standard CoOp<sub>Global</sub>, *i.e.* "CoOp<sub>GL</sub>", is detrimental to final top-1 performances, with -1.9pt top-1 or -3.6pt DG compared to CoOp<sub>Global</sub>. On the other

**Table 2:** Ablation studies for the different components of our GalLoP (ImageNet, 16 shots).

	Top-1	DG	FPR95	↓ AUC
CLIP <sub>Global</sub>	66.6	57.2	42.8	90.8
$\operatorname{CLIP}_{\operatorname{Local}}$	12.5	9.49	73.3	73.7
$\mathrm{CLIP}_{\mathrm{GL}}$	61.1	49.3	35.5	90.8
CoOp <sub>Global</sub>	71.4	59.2	39.1	91.1
$\rm CoOp_{\rm Local}$	41.2	30.1	65.2	78.3
$\rm CoOp_{GL}$	69.5	55.6	33.7	90.5
$\overline{\mathrm{GalLoP}_{\mathrm{Global}}}$	72.0	60.4	37.0	91.7
$\operatorname{GalLoP}_{\operatorname{Local}}$	70.9	54.1	36.0	90.1
GalLoP	75.1	61.3	27.3	93.2

hand, GalLoP enjoys a boost in performances on all metrics when combining the learned global (GalLoP<sub>Global</sub>) and local (GalLoP<sub>Local</sub>) prompts. We can see that the top-1 performances of GalLoP increase by +3.1pt compared to (GalLoP<sub>Global</sub>). Similarly, on OOD detection, GalLoP has a decrease of -8.9pt FPR95 compared to GalLoP<sub>Local</sub>. Tab. 2 illustrates how the resulting performances of GalLoP, in both accuracy and robustness, come from the complementarity of both the local and global features.

The need for sparsity. In Fig. 7 we show how the sparsity when using local features allows achieving higher performances than attending to each local feature, for three regimes: zero-shot CLIP ("zero-shot"), while learning a local prompt ("w/o linear"), and when aligning a local prompt and our linear projection ("w. linear"). On the three regimes, the difference between looking at all local features and the best reported sparsity level is, respectively, +18.4pt, +17.6pt, and +8.5 pt. Furthermore, we can see that when aligning a local prompt and the linear layer, our sparsity ratio works for a wide range of k, with performances above 69pt between k = 5 and k = 50. This shows the



Fig. 7: Impact of our sparsity choice for three regimes, zero-shot CLIP, learning a local prompt ("w/o linear") and aligning our linear projection with a local prompt ("w. linear") (ImageNet, 16 shots).

robustness to the choice of k. Finally, learning a local prompt allows to significantly boost the performances for the local features, e.g. +27.9pt for k=10, and aligning with a linear projection further boosts performances, with +10pt for k=10 compared to learning the prompt only. Fig. 7 shows the interest of both enforcing the sparsity when looking at local features and further aligning the local features with a local prompt.



Fig. 8: Impact of our design choices on learning global Fig. 8a and local prompts Fig. 8b.

**Global prompt learning with prompt dropout.** We display on Fig. 8a how prompt dropout allows learning efficiently multiple prompts for the global features. We display the top-1 accuracy when using more and more prompts, with ("w.") or without ("w/o") prompt dropout. We can observe that adding more prompts does not result in better performances without prompt dropout. For example, performances with 6 prompts decrease compared to using a single prompt . This is due to limited diversity among the learned prompts. In comparison, adding more prompts is always beneficial when using prompt dropout.

**Local prompt learning at multiple scales.** On Fig. 8b, we show the interest of our multiscale approach. We experiment with various number of scales, *i.e.* from 1 to 6 scales with  $k_1 = 10$  and  $\Delta_k = 10$  and report the top-1 accuracy. We can observe a steady increase from 1 scale to 4 scales (+1pt). Performances stabilize afterward for 5 and 6 scales. Fig. 8b shows that learning at different scales is beneficial, but also that GalLoP is not too sensitive to the choice of number of prompts. Furthermore, learning at different scale also reduce the need to select an optimal k, although we show in Fig. 7 that performances are stable with respect to k.

#### 4.4 Qualitative study.

We conduct in this section a qualitative study of GalLoP, by comparing it to CLIP on Fig. 9, and visualizing its different scales on Fig. 10. We show other qualitative results in supplementary material B.6

**Comparison to CLIP.** On Fig. 9, we compare GalLoP and CLIP local features. We can observe that CLIP's local features are not discriminative and do not allow to classify images correctly, which was observed in Sec. 4.3. On the other hand, GalLoP classifies correctly the images, even with a single scale. We can also observe GalLoP accurately segments the object of interest when using all its scales.

**Visualize multiple scales.** Finally, we show the different regions each of the local prompts attend to. We can see that scale # 1 focuses on the most discriminative features, *i.e.* the head and tail of the "Ring tailed lemur". Each scale progressively attends to different parts of the body, leading to an accurate prediction.

#### 14 M. Lafon et al.



Fig. 9: Qualitative comparison of CLIP and GalLoP. From left to right, the original image with its ground truth, CLIP local wrong prediction, one scale (k=10) of GalLoP with correct prediction and GalLoP multiscale, resulting in correct prediction and segmentation.



**Fig. 10:** GalLoP multiscale visualization. Regions observed by the different prompts of GalLoP for a "Ring tailed lemur".

# 5 Conclusion

This paper introduces GalLoP, a new prompt learning method that leverage both global and local visual representations. The key features of GalLoP are the strong discriminability of its local representations and its capacity to produce diverse predictions from both local and global prompts. Extensive experiments show that GalLoP outperforms previous prompt learning methods on top-1 accuracy on average for 11 datasets; that it works in different few shot settings; and for both convolutional and transformer vision-backbones. We show in ablation studies the interest of the design choices that make GalLoP work, *i.e.* complementarity between local and global prompts; sparsity and enhanced alignment; encouraging diversity. Finally, we conduct a qualitative study to show what local prompts focus on when classifying an image. Future works include learning the local feature alignment on a large vision-language dataset.

# Acknowledgements

This work was done under grants from the DIAMELEX ANR program (ANR-20-CE45-0026) and the AHEAD ANR program (ANR-20-THIA-0002). It was granted access to the HPC resources of IDRIS under the allocation AD011012645R1 and AD011013370R1 made by GENCI.

## References

- 1. Agnolucci, L., Baldrati, A., Todino, F., Becattini, F., Bertini, M., Del Bimbo, A.: Eco: Ensembling context optimization for vision-language models. In: ICCV (2023) 2
- Bossard, L., Guillaumin, M., Van Gool, L.: Food-101 mining discriminative components with random forests. In: ECCV (2014) 9, 24
- Chen, G., Yao, W., Song, X., Li, X., Rao, Y., Zhang, K.: Plot: Prompt learning with optimal transport for vision-language models. In: The Eleventh International Conference on Learning Representations (2023) 1, 2, 3, 4, 9
- Cho, J., Nam, G., Kim, S., Yang, H., Kwak, S.: Promptstyler: Prompt-driven style generation for source-free domain generalization. In: IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023. pp. 15656–15666. IEEE (2023) 4, 20
- Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., Vedaldi, A.: Describing textures in the wild. In: Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2014) 9, 11, 24, 25, 26
- Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009) 9, 24
- Dong, X., Bao, J., Zheng, Y., Zhang, T., Chen, D., Yang, H., Zeng, M., Zhang, W., Yuan, L., Chen, D., et al.: Maskclip: Masked self-distillation advances contrastive language-image pretraining. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10995–11005 (2023) 6
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020) 9
- Fei-Fei, L., Fergus, R., Perona, P.: Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. Computer Vision and Pattern Recognition Workshop (2004) 9, 24
- Gal, Y., Ghahramani, Z.: Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In: International Conference on Machine Learning. pp. 1050–1059. PMLR (2016) 7
- Gao, P., Geng, S., Zhang, R., Ma, T., Fang, R., Zhang, Y., Li, H., Qiao, Y.: Clipadapter: Better vision-language models with feature adapters. International Journal of Computer Vision 132(2), 581–595 (2024) 23, 24
- Gondal, M.W., Gast, J., Ruiz, I.A., Droste, R., Macri, T., Kumar, S., Staudigl, L.: Domain aligned clip for few-shot classification. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 5721–5730 (2024) 24
- Goyal, S., Kumar, A., Garg, S., Kolter, Z., Raghunathan, A.: Finetune like you pretrain: Improved finetuning of zero-shot vision models. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023. pp. 19338–19347. IEEE (2023) 23, 24

- 16 M. Lafon et al.
- He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. arxiv e-prints. arXiv preprint arXiv:1512.03385 10 (2015) 9
- 15. Helber, P., Bischke, B., Dengel, A., Borth, D.: Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification (2017) 9, 24
- Hendrycks, D., Basart, S., Mu, N., Kadavath, S., Wang, F., Dorundo, E., Desai, R., Zhu, T., Parajuli, S., Guo, M., Song, D., Steinhardt, J., Gilmer, J.: The many faces of robustness: A critical analysis of out-of-distribution generalization. ICCV (2021) 10, 25
- Hendrycks, D., Gimpel, K.: A baseline for detecting misclassified and out-ofdistribution examples in neural networks. arXiv preprint arXiv:1610.02136 (2016) 20
- 18. Hendrycks, D., Mazeika, M., Dietterich, T.: Deep anomaly detection with outlier exposure. arXiv preprint arXiv:1812.04606 (2018) 4
- Hendrycks, D., Zhao, K., Basart, S., Steinhardt, J., Song, D.: Natural adversarial examples. CVPR (2021) 10, 25
- Jia, C., Yang, Y., Xia, Y., Chen, Y.T., Parekh, Z., Pham, H., Le, Q., Sung, Y.H., Li, Z., Duerig, T.: Scaling up visual and vision-language representation learning with noisy text supervision. In: International conference on machine learning. pp. 4904–4916. PMLR (2021) 1
- Khattak, M.U., Rasheed, H., Maaz, M., Khan, S., Khan, F.S.: Maple: Multi-modal prompt learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 19113–19122 (2023) 1, 2, 3, 4, 9
- Khattak, M.U., Wasim, S.T., Naseer, M., Khan, S., Yang, M.H., Khan, F.S.: Self-regulating prompts: Foundational model adaptation without forgetting. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 15190–15200 (2023) 1, 2, 3, 4, 9, 22
- Krause, J., Stark, M., Deng, J., Fei-Fei, L.: 3d object representations for fine-grained categorization. In: 4th International IEEE Workshop on 3D Representation and Recognition (3dRR-13). Sydney, Australia (2013) 9, 24
- Lafon, M., Ramzi, E., Rambour, C., Thome, N.: Hybrid energy based model in the feature space for out-of-distribution detection. In: International Conference on Machine Learning. pp. 18250–18268. PMLR (2023) 20
- Lee, K., Lee, K., Lee, H., Shin, J.: A simple unified framework for detecting outof-distribution samples and adversarial attacks. Advances in neural information processing systems **31** (2018) 20
- Lu, Y., Liu, J., Zhang, Y., Liu, Y., Tian, X.: Prompt distribution learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5206–5215 (2022) 1, 2, 3, 4, 9, 20
- Maji, S., Kannala, J., Rahtu, E., Blaschko, M., Vedaldi, A.: Fine-grained visual classification of aircraft. Tech. rep. (2013) 9, 24
- Ming, Y., Cai, Z., Gu, J., Sun, Y., Li, W., Li, Y.: Delving into out-of-distribution detection with vision-language representations. Advances in Neural Information Processing Systems 35, 35087–35102 (2022) 4, 20, 25
- Miyai, A., Yu, Q., Irie, G., Aizawa, K.: Locoop: Few-shot out-of-distribution detection via prompt learning. NeurIPS 36 (2023) 2, 4, 7, 9, 26
- Miyai, A., Yu, Q., Irie, G., Aizawa, K.: Zero-shot in-distribution detection in multiobject settings using vision-language foundation models. CoRR (2023) 4, 8, 19, 20
- Nie, J., Zhang, Y., Fang, Z., Liu, T., Han, B., Tian, X.: Out-of-distribution detection with negative prompts. In: The Twelfth International Conference on Learning Representations (2024) 4, 26

- Nilsback, M.E., Zisserman, A.: Automated flower classification over a large number of classes. In: Proceedings of the Indian Conference on Computer Vision, Graphics and Image Processing (Dec 2008) 9, 24
- Parisot, S., Yang, Y., McDonagh, S.: Learning to name classes for vision and language models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 23477–23486 (2023) 1, 2
- Parkhi, O.M., Vedaldi, A., Zisserman, A., Jawahar, C.V.: Cats and dogs. In: IEEE Conference on Computer Vision and Pattern Recognition (2012) 9, 24
- 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: International conference on machine learning. pp. 8748–8763. PMLR (2021) 1, 3, 9
- Recht, B., Roelofs, R., Schmidt, L., Shankar, V.: Do imagenet classifiers generalize to imagenet? In: International conference on machine learning. pp. 5389–5400. PMLR (2019) 10, 25
- Sehwag, V., Chiang, M., Mittal, P.: SSD: A unified framework for self-supervised outlier detection. In: 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021 (2021) 20
- Shu, Y., Guo, X., Wu, J., Wang, X., Wang, J., Long, M.: Clipood: Generalizing CLIP to out-of-distributions. In: ICML (2023) 23, 24
- 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) 9, 24
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: a simple way to prevent neural networks from overfitting. The journal of machine learning research 15(1), 1929–1958 (2014) 7
- Sun, X., Hu, P., Saenko, K.: Dualcoop: Fast adaptation to multi-label recognition with limited annotations. Advances in Neural Information Processing Systems 35, 30569–30582 (2022) 4, 7, 19
- Sun, Y., Ming, Y., Zhu, X., Li, Y.: Out-of-distribution detection with deep nearest neighbors. In: International Conference on Machine Learning. pp. 20827–20840. PMLR (2022) 20
- Van Horn, G., Mac Aodha, O., Song, Y., Cui, Y., Sun, C., Shepard, A., Adam, H., Perona, P., Belongie, S.: The inaturalist species classification and detection dataset. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 8769–8778 (2018) 11, 25, 26
- 44. Villani, C., et al.: Optimal transport: old and new, vol. 338. Springer (2009) 4
- 45. Wang, F., Li, M., Lin, X., Lv, H., Schwing, A., Ji, H.: Learning to decompose visual features with latent textual prompts. In: The Eleventh International Conference on Learning Representations (2022) 2
- Wang, H., Ge, S., Lipton, Z., Xing, E.P.: Learning robust global representations by penalizing local predictive power. In: Advances in Neural Information Processing Systems. pp. 10506–10518 (2019) 10, 25
- Xiao, J., Hays, J., Ehinger, K.A., Oliva, A., Torralba, A.: Sun database: Large-scale scene recognition from abbey to zoo. In: CVPR (2010) 9, 11, 24, 25, 26
- Zhang, R., Zhang, W., Fang, R., Gao, P., Li, K., Dai, J., Qiao, Y., Li, H.: Tip-adapter: Training-free adaption of CLIP for few-shot classification. In: ECCV. pp. 493–510 (2022) 24
- Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., Torralba, A.: Places: A 10 million image database for scene recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence (2017) 11, 25, 26

- 18 M. Lafon et al.
- Zhou, C., Loy, C.C., Dai, B.: Extract free dense labels from CLIP. In: Avidan, S., Brostow, G.J., Cissé, M., Farinella, G.M., Hassner, T. (eds.) Computer Vision -ECCV 2022 - 17th European Conference, Tel Aviv, Israel. Lecture Notes in Computer Science, vol. 13688, pp. 696–712. Springer (2022) 7, 19
- 51. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Conditional prompt learning for visionlanguage models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16816–16825 (2022) 1, 3, 9
- Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Learning to prompt for vision-language models. International Journal of Computer Vision 130(9), 2337–2348 (2022) 1, 3, 4, 7, 9, 21, 22