
Prompt Learning with Optimal Transport for Vision-Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 With the increasing attention to large vision-language models such as CLIP, there
2 has been a significant amount of effort dedicated to building efficient prompts.
3 Unlike conventional methods of only learning one single prompt, we propose
4 to learn multiple comprehensive prompts to describe diverse characteristics of
5 categories such as intrinsic attributes or extrinsic contexts. However, directly
6 matching each prompt to the same visual feature is problematic, as it pushes the
7 prompts to converge to one point. To solve this problem, we propose to apply
8 optimal transport to match the vision and text modalities. Specifically, we first
9 model images and the categories with visual and textual feature sets. Then, we
10 apply a two-stage optimization strategy to learn the prompts. In the inner loop, we
11 optimize the optimal transport distance to align visual features and prompts by the
12 Sinkhorn algorithm, while in the outer loop, we learn the prompts by this distance
13 from the supervised data. Extensive experiments are conducted on the few-shot
14 recognition task and the significant improvement demonstrates the superiority of
15 our method.

16 1 Introduction

17 In the past few years, large-scale vision-language pre-trained (VLP) models, such as CLIP [39],
18 ALIGN [17], and BLIP [23] have achieved remarkable success in open-world visual concept learning.
19 These methods have brought new light but also pose a new question: how to efficiently adapt the
20 knowledge from pretraining to the downstream tasks since these models are typically of massive sizes
21 which are not feasible for normal users to re-train.

22 One of the conventional
23 paradigms of utilizing pretrained
24 knowledge is “pre-training,
25 fine-tuning”, which fixes the
26 architecture of the pre-trained
27 neural network and tunes its
28 parameters using task-specific
29 objective functions. Beyond

30 fine-tuning the parameters,
31 many recent methods [62, 63]
32 introduce the concept of prompt
33 learning from the field of NLP to the vision domain and achieve striking performance gain for the
34 few-shot visual classification. They fix the model parameters and instead learn suitable prompts
35 by turning a template sentence into a set of learnable vectors. Then, these prompts are learned by
36 minimizing the distance between the visual features and prompt-based language features.

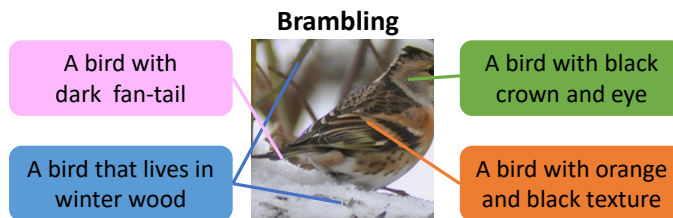


Figure 1: The motivation that one category can be complementarily described in different views (An example of “Brambling”).

37 Despite significant improvements over manual prompts, learning only a sentence is intuitively
38 insufficient to represent a class. One class can be described by many intrinsic characteristics and
39 even extrinsic context relations. Thus, for one object, we may have multiple prompt candidates
40 which focus on different attributes. As shown in Figure 1, we can describe the class “Brambling” in
41 different views: such as the color of the wing, the color of the crown and eyes, the shape and color of
42 the tail, and even the living environment information. It motivates us to learn multiple prompts to
43 comprehensively represent the class and thus facilitate classification.

44 The most natural solution is to directly learn multiple prompts by respectively matching each prompt
45 with the visual features. However, it is the same as matching the mean of prompt features and the
46 visual features. This solution is problematic since all prompts are encouraged to be closer to one single
47 point and thus tend to learn the same characteristics. It contradicts our purpose to learn comprehensive
48 prompts. To solve this problem, we tested adding some constraints to push away the prompt from
49 each other, but found that this solution still fails to learn representative and comprehensive prompts.
50 This solution treats the visual representation as one single point, and such a unified view of visual
51 features ignores the fact that different prompts may only focus on one or a subset of characteristics.

52 To address this problem, in this paper, we propose Prompt Learning with Optimal Transport (PLOT),
53 which applies optimal transport (OT) to align the local visual features and multiple textual prompts.
54 Optimal transport can calculate the distance between two distributions under the form of multiple
55 sampling. In our prompts learning framework, we formulate local visual features and multiple prompts
56 as the samplings of two discrete distributions and use OT to encourage fine-grained cross-modal
57 matching. Specifically, to obtain the local visual features with different semantic clues, we extract all
58 feature maps as the visual representation instead of the single global representation. Fortunately, we
59 can easily obtain the visual feature maps from the visual encoder of CLIP by using all outputs of the
60 multi-head self-attention layer [42]. Then the problem comes down to how to calculate the distance
61 between two feature sets.

62 We solve this problem by introducing the optimal transport theory [50] and formulate the feature sets
63 as a discrete probability distribution where each feature has an equal probability value. Furthermore,
64 to reduce the computational cost and avoid the extra model parameters, we learn the prompts with
65 a two-stage optimization strategy. At the first stage in the inner loop, we fix both visual and text
66 features and optimize the optimal transport problem by a fast Sinkhorn distances algorithm [6]. Then,
67 in the outer loop, we fix all parameters of optimal transport and back-propagate the gradient to learn
68 the prompts with different characteristics. We conduct comprehensive experiments on 11 datasets
69 following the standard setting of CLIP [39] and CoOp [62] to evaluate our method. These experiments
70 span the visual classification on generic objects, scenes, actions, fine-grained categories and so on.
71 The significant result improvement demonstrates that PLOT can effectively learn representative and
72 comprehensive prompts.

73 2 Related Work

74 **Optimal Transport** The Optimal Transport [30] is initially introduced to solve the problem of how
75 to reduce the cost when moving simultaneously several items. Recently, OT theory has drawn wide
76 attention in the machine learning and computer vision community by comparing distributions readily
77 available to them under the form of feature sets [37]. Due to the brilliant property of distribution
78 matching, OT has been applied in many theoretic and application tasks including generative models [1,
79 44, 59], structural matching [4, 56, 60, 55] (e.g. sequence matching [4] and graph matching [55]),
80 and other distribution-based tasks (such as clustering [22], distribution estimation [2], and causal
81 discovery [49]). In this paper, we use OT distance to align the features of vision and language
82 modalities and propose a two-stage learning strategy to guide the learning of multiple prompts.

83 **Vision-Language Pre-trained Models** Vision-Language Pre-trained (VLP) models aim to explore
84 the semantic correspondence between the vision and language modalities through large-scale pre-
85 training. Recently, VLP models have achieved an exciting performance improvement in the zero-shot
86 and few-shot visual recognition [39, 10, 62, 63, 58], which shows the great potential to promote
87 open-world visual understanding with the help of language. One key part of learning VLP models is
88 the self-supervised learning objective on two modalities. The popular VLP objectives can be divided
89 into reconstruction [25, 15, 8, 20], contrastive matching [39, 17, 16], or the combination of both
90 two [24, 53, 19]. Besides, recent progresses in the field of VLP also benefit a lot from large scale

91 pair-wised datasets. For example, CLIP [39] apply 400 million image-text pairs for the contrastive
 92 learning, while ALIGN even exploits 1.8 billion data pairs. Beyond recognition, these VLP models
 93 also show great potential for other downstream applications, such as dense prediction [42, 61], image
 94 generation [31, 41, 35], and action understanding [52, 47].

95 **Prompt Learning** Prompt learning is introduced from the field of NLP to efficiently adapt the large
 96 language model to downstream tasks. Different from the conventional “pre-training, fine-tuning”
 97 paradigm which initializes the pre-trained model and tunes the parameters of the network using
 98 downstream task-specific objective functions, prompt learning applies textual prompt to reformulate
 99 the downstream tasks as the original pretrained task [27, 36]. By the prompt, the domain shift between
 100 pretrained task and downstream application is reduced and thus the pretrained knowledge can be
 101 easier adapted to downstream tasks. The concept of prompt learning [36, 40, 38] begins from the
 102 success of GPT [40] series. Early prompt learning methods (such as Petroni *et al.* [36] and Pörner *et*
 103 *al.* [38]) always manually create templates based on human prior knowledge. Furthermore, some
 104 mining-based methods [18] and gradient-based methods [45] are proposed to automatically search for
 105 appropriate templates. Beyond search in the discrete space, some methods [26, 48, 28] remove the
 106 constraint that the prompts are “words” and instead learn prompts in the continuous embedding space.
 107 Recently, CoOp [62] and its extended version [63] introduce prompt learning into open-world visual
 108 understanding to adapt the knowledge from the large-scale visual-language pretrained models and
 109 achieve great performance improvement on the few-shot visual recognition. Compared with CoOp,
 110 our PLOT method further improve prompt learning by introducing the optimal transport distance to
 111 learn multiple local prompts and achieve fine-grained vision-language matching.

112 3 Approach

113 In this section we will first revisit the baseline method CoOp 3.1, review the preliminaries of
 114 optimal transport 3.2, and then introduce our proposed PLOT 3.3 to show how we can learn multiple
 115 comprehensive prompts.

116 3.1 A Revisit of CoOp

117 CoOp [62] is one of the pioneering methods to learn the prompts for using vision language pretrained
 118 knowledge (such as CLIP [39]) for downstream open-world visual recognition. Different from CLIP
 119 that manually designs the prompt templates, CoOp sets a part of context words in the template as
 120 continuous learnable parameters which can be learned from the few-shot data. Then the classification
 121 weights can be represented by the distance between the learned prompt and visual feature.

122 Specifically, given an image x , a visual feature $\mathbf{f} = f(x)$ is obtained by the visual encoder f of
 123 CLIP. Then, the textual prompt can be formulated as $\mathbf{t}_k = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_L, \mathbf{c}_k\}$, where \mathbf{c}_k is the
 124 word embedding of the class name, $\{\mathbf{v}_l\}_{l=1}^L$ are learnable vectors with the same dimension as the
 125 original word embedding and L is the length of context words. With prompt \mathbf{t}_k as the input, the text
 126 encoder g outputs the textual feature as $\mathbf{g}_k = g(\mathbf{t}_k)$. The final prediction probability is computed by
 127 the matching score as follows:

$$p(y = k|x) = \frac{\exp(\text{sim}(\mathbf{f}, \mathbf{g}_k)/\tau)}{\sum_{k'=1}^K \exp(\text{sim}(\mathbf{f}, \mathbf{g}_{k'})/\tau)}, \quad (1)$$

128 where $\text{sim}(\cdot, \cdot)$ denotes a metric function such as cosine similarity, and τ stands for the temperature
 129 of Softmax. Then we can optimize the parameters of $\{\mathbf{v}_l\}_{l=1}^L$ with the cross-entropy loss between
 130 the prediction and the labeled target.

131 3.2 Optimal Transport

132 Optimal transport (OT) distance is a widely used metric for the comparison of distributions. Here, we
 133 only focus on the discrete situation which is more related to our framework. Assuming we have two
 134 sets of points (features), the discrete distributions are formulated as:

$$U = \sum_{m=1}^M u_m \delta_{\mathbf{f}_m} \quad \text{and} \quad V = \sum_{n=1}^N v_n \delta_{\mathbf{g}_n}, \quad (2)$$

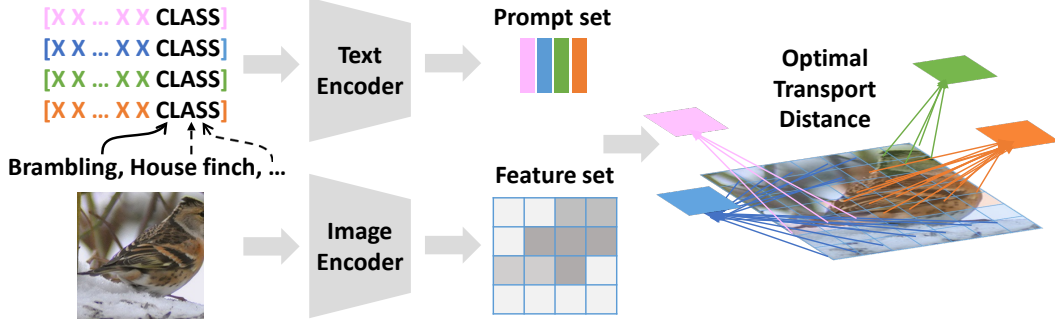


Figure 2: The framework of PLOT. PLOT first describes each category with multiple prompts and obtains a set of prompt features by text encoder. The image is also encoded as a set of local features. Then the optimal transport is used as the metric between prompts and visual features.

135 where \mathbf{u} and \mathbf{v} are the discrete probability vectors that sum to 1, and $\delta_{\mathbf{f}}$ is a Dirac delta function
 136 placed at support point \mathbf{f} in the embedding space. Then, the total distance of these two distributions
 137 are written as:

$$\langle \mathbf{T}, \mathbf{C} \rangle = \sum_{m=1}^M \sum_{n=1}^N T_{m,n} C_{m,n}. \quad (3)$$

138 We call \mathbf{C} the cost matrix in which each point denotes the cost between \mathbf{f}_m and \mathbf{g}_n , such as
 139 $C_{m,n} = 1 - \text{sim}(\mathbf{f}_m, \mathbf{g}_n)$. While the \mathbf{T} is called the transport plan, which is learned to minimize the
 140 total distance. The optimization problem of optimal transport is formulated as:

$$\begin{aligned} d_{OT}(\mathbf{u}, \mathbf{v} | \mathbf{C}) &= \underset{\mathbf{T}}{\text{minimize}} \langle \mathbf{T}, \mathbf{C} \rangle \\ \text{subject to} \quad & \mathbf{T}\mathbf{1} = \mathbf{u}, \mathbf{T}^T \mathbf{1} = \mathbf{v}, \mathbf{T} \geq 0. \end{aligned} \quad (4)$$

141 As directly optimizing the above objective is always time-consuming, we apply the Sinkhorn dis-
 142 tance [6] to use an entropic constraint for fast optimization. The optimization problem with a
 143 Lagrange multiplier of the entropy constraint is:

$$\begin{aligned} d_{OT,\lambda}(\mathbf{u}, \mathbf{v} | \mathbf{C}) &= \underset{\mathbf{T}}{\text{minimize}} \langle \mathbf{T}, \mathbf{C} \rangle - \lambda h(\mathbf{T}) \\ \text{subject to} \quad & \mathbf{T}\mathbf{1} = \mathbf{u}, \mathbf{T}^T \mathbf{1} = \mathbf{v}, \end{aligned} \quad (5)$$

144 where $h(\cdot)$ is entropy and $\lambda \geq 0$ is a hyper-parameter. Then we can have a fast optimization solution
 145 with few iterations as:

$$\mathbf{T}^* = \text{diag}(\mathbf{u}^t) \exp(-\mathbf{C}/\lambda) \text{diag}(\mathbf{v}^t), \quad (6)$$

146 where t denotes iteration and in each iteration $\mathbf{u}^t = \mathbf{u} / ((\exp(-\mathbf{C}/\lambda) \mathbf{v}^{t-1}) \mathbf{1})$ and $\mathbf{v}^t =$
 147 $\mathbf{v} / ((\exp(-\mathbf{C}/\lambda)^T \mathbf{u}^t) \mathbf{1})$, with the initiation $\mathbf{v}^0 = \mathbf{1}$.

148 3.3 Prompt Learning with Optimal Transport

149 In this subsection, we introduce the details of our PLOT, which learns multiple prompts to describe
 150 different characteristics of the category by minimizing the OT distance.

151 Specifically, as shown in Figure 2, given an image \mathbf{x} , we first feed it to the visual encoder branch of
 152 CLIP. Apart from the global visual feature \mathbf{f} , we can also obtain a set of local features $\{\mathbf{f}_m\}_{m=1}^M$.
 153 The visual encoder has a multi-head attention pooling layer in which the input is the combination of
 154 the global feature and a set of local features (feature map) and the output is a tensor with the shape
 155 $\mathbb{R}^{(H \times W + 1) \times C}$, where H and W is the height and width of feature map and C is the feature dimension.
 156 Therefore, we can obtain $M = H \times W$ local features and a global feature. At the same time, for
 157 class k , we can initialize N local prompts as $\{\mathbf{t}_{k,n}\}_{n=1}^N$ with learnable vectors $\{\mathbf{v}_{l,n}\}_{l=1, n=1}^{L,N}$, where
 158 each is same as the prompt in CoOp. With both visual and textual encoders, we can obtain local
 159 visual features $\mathbf{F} = \{\mathbf{f}_m\}_{m=1}^M \in \mathbb{R}^{M \times C}$ and prompt features $\mathbf{G}_k = \{\mathbf{g}_n\}_{n=1}^N \in \mathbb{R}^{N \times C}$.

160 In the inner loop, we learn the transport plan \mathbf{T} with these fixed support sets \mathbf{F}, \mathbf{G}_k , by minimizing
 161 the following OT distance to push \mathbf{G}_k to \mathbf{F} :

$$d_{OT}(k) = d_{OT}(\mathbf{u}, \mathbf{v} | \mathbf{1} - \mathbf{F}^T \mathbf{G}_k), \quad (7)$$

162 where $C = F^T G_k$ denotes that we use the cosine similarity between F and G_k as the cost matrix.
 163 Then we can obtain the solution of transport plan T^* as Eq (6) and the final OT distance $d_{OT}(k)$.

164 Given the OT distance between G_k and F , we reformulate the prediction probability as:

$$p_{ot}(y = k|\mathbf{x}) = \frac{\exp((1 - d_{OT}(k))/\tau)}{\sum_{k'=1}^K \exp((1 - d_{OT}(k'))/\tau)}. \quad (8)$$

165 In the outer loop, we fix the transport plan T^* and apply the cross entropy loss to optimize the
 166 $\{\mathbf{v}_{l,n} |_{l=1,n=1}^{L,N}\}$ as:

$$L_{CE} = -\frac{1}{|\mathcal{X}|} \sum_{\mathbf{x} \in \mathcal{X}} \sum_{k=1}^K y_{\mathbf{x},k} p_{ot}(y = k|\mathbf{x}), \quad (9)$$

167 where $\mathbf{y}_{\mathbf{x}}$ is a one-hot label vector. The detail algorithm can be found in the supplementary materials.

168 4 Experiments

169 Extensive experiments are conducted to evaluate our method, including comparison with CoOp,
 170 ablation studies, parameter analysis extensibility analysis, computing cost analysis and visualization.

171 4.1 Datasets

172 We followed the experimental settings in the CoOp [62] for the few-shot learning evaluation. The
 173 experiments are conducted on the 11 visual recognition datasets, including Caltech101 [9], DTD [5],
 174 EuroSAT [12], FGVCAircraft [29], Flowers102 [32], Food101 [3], ImageNet [7], OxfordPets [33],
 175 StanfordCars [21], SUN397 [54], and UCF101 [46]. These datasets span visual classification on
 176 generic objects, scenes, actions, fine-grained categories and so on, which constitutes a comprehensive
 177 evaluation of our method. All experiments adopted the few-shot evaluation protocol used in CLIP [39]
 178 and CoOp [62], where we respectively choose 1, 2, 4, 8 and 16 shots for model training and use
 179 original test set for evaluation. Besides, we also evaluated the robustness of our method with domain
 180 shift. Following CoOp, we used the ImageNet as the source domain and evaluate our method with
 181 ImageNet-based robustness evaluation datasets including ImageNetV2 [43], ImageNet-Sketch [51],
 182 ImageNet-A [14], and ImageNet-R [13]. The detailed introduction of each dataset can be found in
 183 the supplementary materials.

184 4.2 Implementation details

185 We chose CoOp [62] as our main competitor to evaluate our method. Compared with CoOp which
 186 only learns a global prompt for one class, our PLOT method learns multiple local prompts and applies
 187 the OT distance for better fine-grained alignment. Besides, we also reported the performance of
 188 training a linear classifier with the CLIP [39] features. It is also a widely-used strategy to adapt the
 189 pretrained knowledge for the downstream task [45]. We reproduced the performance of CoOp and
 190 CLIP linear probe with the released official code.

191 The original CoOp method has different versions with different class token positions and parameter
 192 initialization strategies. We applied the default model that fixes the class token positions in the end due
 193 to the limited performance gap between two different ways of positioning the class token. Besides, we
 194 used the random parameter initialization strategy but not the class-specific context version. Following
 195 the widely used setting in [62, 63, 10, 57], we also chose RN50 [11] as the backbone network of the
 196 visual branch and set the length of learnable context tokens as 16. All the code of our method is based
 197 on CoOp, which adopted the SGD optimizer with 0.002 initial learning rate, CosineAnnealingLR
 198 schedule, and a warmup trick with 1e-5 learning rate. Besides, we also followed the epoch strategy to
 199 train more epochs for more shots.

200 For the optimal transport distance, we apply $N = 4$ prompts for each category and use $M = 7 \times 7$
 201 due to the feature map size. We set the hyper-parameters in the Sinkhorn distances algorithm [6] as
 202 $\lambda = 0.1$ for all the datasets. We set the maximum iteration number of the inner loop as 100 and will
 203 early stop the iteration when the average absolute update value $\Lambda < 0.01$. We initialize all values in
 204 the vector v and μ as $1/N$ and $1/M$ respectively. All models are conducted on the Pytorch [34] 1.7.1
 205 and trained on 4 NVIDIA A100 GPUs. We repeated the experiments three times with different seeds
 206 and reported the average.

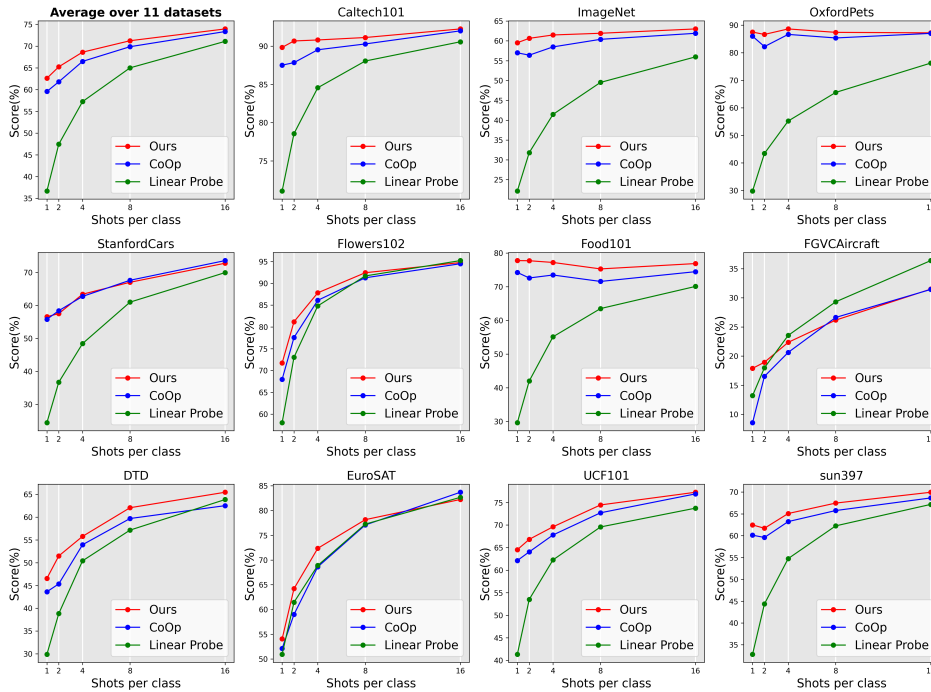


Figure 3: The few-shot learning results on 11 datasets. We compare our PLOT with CoOp and Linear Probe method, and observe the consistent and significant performance improvement on most of datasets. (The average accuracy on all datasets is shown in the left top.)

207 4.3 Comparison With CoOp

208 In this subsection, we compare our PLOT with the baseline CoOp on the few-shot recognition and
 209 domain generalization tasks.

210 4.3.1 Few-Shot Learning

211 We summarized the experimental results in Figure 3 where the red line denotes our PLOT method,
 212 the blue one denotes CoOp, and the green one is CLIP linear probe. The detailed accuracy can be
 213 found in the supplementary materials. We observed that both prompt learning methods (PLOT and
 214 CoOp) outperform the linear probe method by a large margin. Besides, PLOT can further improve
 215 the performance of CoOp on most of the datasets. Taking the average accuracy (at the left top) as
 216 the example, Plot respectively gained 3.03%, 3.45%, 2.13%, 1.38%, 0.61% performance boost over
 217 CoOp at 1, 2, 4, 8, 16 shots. We found the performance gap will reduce when shots increase. It is not
 218 surprising since both CoOp and PLOT focus on utilizing the pre-trained knowledge, and the effect
 219 of pre-training diminishes given more training data. Among all datasets, PLOT achieves a larger
 220 improvement over CoOp on the FOOD101 and DTD datasets and achieves comparable performance
 221 only on the StanfordCars datasets. For the FGVCAircraft dataset in which the CoOp only obtains
 222 7.77% accuracy, our PLOT can achieve an accuracy of 17.79%, twice as high as that of CoOp. Note
 223 that we don't use the class-specific context, thus the performance on the fine-grained classification
 224 datasets is lower, e.g. the performance of both CoOp and PLOT without class-specific context is
 225 lower than the linear probing on FGVCAircraft. All these performance comparisons can serve as
 226 experimental evidence to demonstrate that multiple local prompts and optimal transport distance
 227 facilitate the prompt learning of vision-language models.

228 4.3.2 Domain generalization

229 The robustness also plays a critical role in model applications since the real-world environment may
 230 have large domain shifts with the training data. Therefore, we conducted a robustness evaluation to
 231 investigate the transferability of models learned by PLOT.

Table 1: Comparison with CoOp on robustness to domain shift.

| Method | Source | Target | | | |
|-----------------------|--------------|--------------|--------------|--------------|--------------|
| | ImageNet | -V2 | -Sketch | -A | -R |
| CLIP + CoOp | 61.91 | 54.26 | 32.47 | 21.78 | 54.21 |
| CLIP + PLOT ($N=4$) | 63.01 | 55.11 | 33.00 | 21.86 | 55.61 |

Table 2: **Ablation studies on few-shot recognition.** PLOT is our defined model with $N = 4$, CoOp is the baseline method, M denotes that we respectively match the global visual feature and multiple textual prompts, V denotes that we apply a constraint to add the variance of prompts, M indicates using the visual feature map instead of the global visual feature. The comparison with different prompt number is shown in white, while the comparison with different ablation versions are highlighted in gray.

| Dataset | Settings | 1 shot | 2 shots | 4 shots | 8 shots | 16 shots |
|------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|
| Caltech101 | PLOT | 89.83 ± 0.33 | 90.67 ± 0.21 | 90.80 ± 0.20 | 91.54 ± 0.33 | 92.24 ± 0.38 |
| | CoOp | 87.51 ± 1.02 | 87.84 ± 1.10 | 89.52 ± 0.80 | 90.28 ± 0.42 | 91.99 ± 0.31 |
| | G | 88.13 ± 0.36 | 86.98 ± 1.25 | 88.45 ± 0.79 | 90.16 ± 0.22 | 90.72 ± 0.18 |
| | G+V | 88.28 ± 0.43 | 87.72 ± 1.25 | 88.45 ± 0.30 | 89.82 ± 0.20 | 92.00 ± 0.13 |
| | M | 69.78 ± 1.75 | 71.57 ± 1.59 | 77.18 ± 2.16 | 81.77 ± 0.47 | 86.21 ± 0.20 |
| | M+V | 66.11 ± 8.29 | 71.45 ± 3.98 | 79.30 ± 3.96 | 86.96 ± 0.78 | 89.80 ± 0.17 |
| | N=1 | 88.47 ± 1.15 | 89.19 ± 0.39 | 89.70 ± 0.38 | 90.45 ± 0.24 | 91.56 ± 0.14 |
| | N=2 | 88.86 ± 0.51 | 89.60 ± 0.10 | 90.60 ± 0.17 | 91.25 ± 0.65 | 91.89 ± 0.36 |
| | N=4 | 89.83 ± 0.33 | 90.67 ± 0.21 | 90.80 ± 0.20 | 91.54 ± 0.33 | 92.24 ± 0.38 |
| N=8 | 89.74 ± 0.30 | 90.18 ± 0.46 | 91.02 ± 0.18 | 91.28 ± 0.28 | 92.04 ± 0.29 | |
| DTD | PLOT | 46.55 ± 2.62 | 51.24 ± 1.95 | 56.03 ± 0.43 | 61.70 ± 0.35 | 65.60 ± 0.82 |
| | CoOp | 43.62 ± 1.96 | 45.35 ± 0.31 | 53.94 ± 1.37 | 59.69 ± 0.13 | 62.51 ± 0.25 |
| | G | 45.12 ± 1.69 | 48.39 ± 2.08 | 54.75 ± 0.48 | 60.15 ± 0.70 | 63.59 ± 0.76 |
| | G+V | 45.90 ± 2.00 | 48.50 ± 0.99 | 53.96 ± 0.48 | 59.69 ± 1.01 | 63.51 ± 0.66 |
| | M | 13.18 ± 4.57 | 12.25 ± 3.86 | 13.00 ± 4.73 | 20.76 ± 5.42 | 26.99 ± 1.98 |
| | M+V | 12.61 ± 5.93 | 15.11 ± 1.81 | 20.35 ± 1.33 | 44.13 ± 2.39 | 56.85 ± 0.54 |
| | N=1 | 43.91 ± 0.65 | 48.21 ± 2.20 | 53.69 ± 1.10 | 58.90 ± 0.19 | 62.85 ± 0.74 |
| | N=2 | 45.59 ± 2.46 | 48.06 ± 1.92 | 55.58 ± 1.71 | 61.56 ± 0.17 | 64.60 ± 0.92 |
| | N=4 | 46.55 ± 2.62 | 51.24 ± 1.95 | 56.03 ± 0.43 | 61.70 ± 0.35 | 65.60 ± 0.82 |
| N=8 | 46.89 ± 1.94 | 51.87 ± 2.06 | 54.45 ± 0.48 | 62.20 ± 0.56 | 65.25 ± 0.38 | |
| FOOD101 | PLOT | 77.74 ± 0.47 | 77.70 ± 0.02 | 77.21 ± 0.43 | 75.31 ± 0.30 | 77.09 ± 0.18 |
| | CoOp | 74.25 ± 1.52 | 72.61 ± 1.33 | 73.49 ± 2.03 | 71.58 ± 0.79 | 74.48 ± 0.15 |
| | G | 74.63 ± 0.11 | 70.15 ± 0.49 | 70.41 ± 0.46 | 70.72 ± 0.98 | 73.68 ± 0.46 |
| | G+V | 74.83 ± 0.31 | 70.09 ± 0.85 | 70.86 ± 0.22 | 70.80 ± 0.68 | 73.93 ± 0.35 |
| | M | 52.02 ± 4.86 | 46.12 ± 1.46 | 46.86 ± 1.39 | 53.43 ± 0.88 | 61.28 ± 0.23 |
| | M+V | 46.52 ± 1.15 | 45.95 ± 2.66 | 53.57 ± 0.83 | 62.95 ± 0.37 | 67.63 ± 1.11 |
| | N=1 | 75.96 ± 0.48 | 76.12 ± 0.59 | 77.11 ± 0.41 | 76.56 ± 0.69 | 77.43 ± 0.80 |
| | N=2 | 77.12 ± 0.49 | 76.89 ± 0.23 | 76.16 ± 0.52 | 75.23 ± 0.69 | 76.81 ± 0.50 |
| | N=4 | 77.74 ± 0.47 | 77.70 ± 0.02 | 77.21 ± 0.43 | 75.31 ± 0.30 | 77.09 ± 0.18 |
| N=8 | 78.05 ± 0.15 | 78.19 ± 0.07 | 78.12 ± 0.17 | 76.63 ± 0.22 | 77.48 ± 0.12 | |

232 Table 1 summarize the results of our PLOT method and CoOp on four ImageNet-based robustness
 233 evaluation datasets. For both methods, we trained the models on ImageNet with 16 shots per class.
 234 For PLOT, we set the number of prompts as $N = 4$. We can observe that PLOT outperforms CoOp
 235 consistently on both source and target domains. These experimental results demonstrate that the
 236 performance improvement of our learning multiple prompts doesn't rely on single-domain overfitting.

237 4.4 Ablation Studies and More Analysis

238 In this subsection, we conducted the ablation studies to investigate the effectiveness of different
 239 components, in order to answer the following questions.

240 **Q: Can we directly learn multiple prompts by respectively matching each prompt with the**
 241 **global visual feature? A: No.** As shown in Table 2, we report the performance of directly matching
 242 the global visual feature (notated as "G") and compare it with the baseline CoOp and our PLOT on
 243 three datasets including Caltech101, DTD, and FOOD101. We observe that there is no improvement

Table 3: Comparison with Adapter-based method.

| Dataset | Methods | 1 shot | 2 shots | 4 shots | 8 shots | 16 shots |
|----------|----------------------|--------------|--------------|--------------|--------------|--------------|
| ImageNet | Tip-Adapter-F | 61.32 | 61.69 | 62.52 | 64.00 | 65.51 |
| | Tip-Adapter-F + OT | 61.44 | 61.98 | 62.86 | 64.13 | 65.76 |
| | Tip-Adapter-F + PLOT | 62.27 | 64.31 | 63.89 | 65.04 | 66.17 |

244 over the baseline on some datasets (such as Caltech101 and FOOD101) if we only directly match
 245 prompts and global features. Though “G” obtained the improvement on the DTD dataset, this
 246 improvement is still less than that of PLOT. It is because this “G” method is incentivized to learn the
 247 indistinguishable prompts, which contradicts our purpose to learn multiple comprehensive prompts.
 248 We further add some constraints to push away the prompt from each other. For example, we add an
 249 objective function to add the distance between every two prompts as a regularization term, which
 250 is notated as “V”. However, comparing “G” and “G+V”, we do not find significant and consistent
 251 improvement when using variance loss.

252 **Q: Does the improvement mainly come from using all feature maps? A: No.** In PLOT, we apply
 253 all feature maps of the visual encoder branch, where each feature is a local embedding at one spatial
 254 position. Compared with the global feature, these local features are more informative and contain
 255 fine-grained clues. However, we demonstrate that the improvement of PLOT does not only rely on
 256 using all feature maps. On the contrary, directly using the feature map to replace the global feature
 257 causes a large performance drop. For example, on all three datasets, directly using the feature map
 258 (“M” or “M+V”) has an around 20% 1 shot accuracy drop over using the global visual feature. It
 259 is not surprising since the original CLIP model is trained by matching the global visual feature and
 260 language feature. Without using the OT method, the distance between the feature map and multiple
 261 textual prompts degenerates to the mean distance of each feature-prompt pair. Besides, when using
 262 the feature map, adding the variance loss works well, especially for more shots. For example, the
 263 accuracy on 16 shots DTD is improved by a large margin (from 26.99 to 56.85).

264 **Q: How many prompts are needed? A: 4 prompts are enough** One important hyper-parameter
 265 in PLOT is the number of prompts. To analyze the effect of the number of prompts, we conducted
 266 the experiments on three datasets with 1, 2, 4, 8 prompts. The results are summarized in the white
 267 part of Table 2. We can observe that the performance obviously increases when adding the number
 268 of prompts from 1 to 4. For example, PLOT (N=4) respectively obtains 1.36%, 2.64%, and 1.68%
 269 1-shot accuracy improvement over PLOT (N=1) on three datasets. Besides, when we further increase
 270 the number of prompts, the improvement is not consistent. To balance the improvement and cost,
 271 we set $N = 4$ as the default configuration of our PLOT model. In the experiments, we tuned this
 272 hyper-parameter on the Caltech101 dataset and applied it to other datasets.

273 **Q: Can PLOT benefit zero shot learning? A: No.** CLIP [39] shows that manually designing the
 274 prompts can still achieve good performance. we obtain 7 prompts by prompt engineering on the
 275 ImageNet dataset and can further ensemble them to obtain **60.38%** top 1 accuracy. In this section,
 276 we replace the cosine distance between the global visual feature and prompt ensemble with the OT
 277 distance between the feature map and all 7 prompts. However, without any learning, the OT distance
 278 only obtains **58.78%** accuracy. We argue there are two reasons why the OT distance does not work
 279 without learning: 1) prompt engineering selects prompts based on the global feature and cosine
 280 distance, instead of OT distance with feature map; 2) all these selected prompts are closed to the
 281 global feature and lack the complementarity.

282 **Q: Can PLOT benefit Adapter-based methods? A: Yes.** Adapter-based methods [10, 57] is another
 283 research direction of the efficient adaptation of pre-trained vision-language models. Different from
 284 the prompt learning that fixes the model parameters and tunes the language prompt, adapter-based
 285 methods [10, 57] allow for fine-tuning a part of the network or adding an extra model for training.
 286 Recently, adapter-based methods also achieve good performance on few-shot visual recognition.
 287 Therefore, we want to explore whether our PLOT method can benefit them, and how.

288 We apply the Tip-adapter-F [57] as our baseline method, which learns a $Linear(d, N_{cls} \times K_{shots})$
 289 model to describe one image by the similarity with all training samples, where d is the dimension of
 290 visual feature, N_{cls} is the number of categories (e.g. 1000 in ImageNet), and K_{shots} is the number
 291 of shots. Then, the final similarity consists of the original distance between the visual feature and
 292 prompt ensembling and the new distance calculated by the learned feature and one-hot vector of

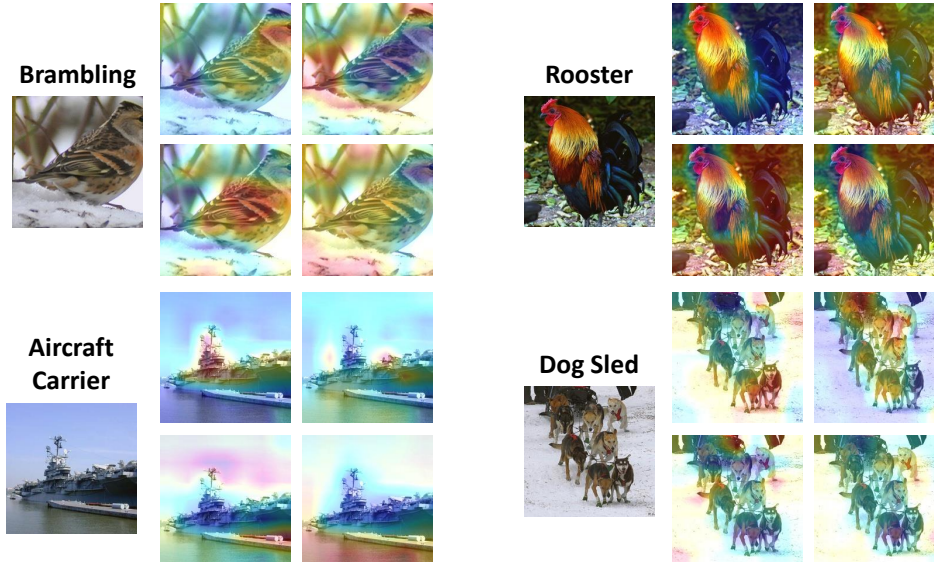


Figure 4: Visualizations. We provide the heatmaps of transport plan T related to each prompt on 4 categories in ImageNet. Different transport plans focus on different attributes of the object.

293 labels (whose dimension is $(N_{cls} \times K_{shots}, N_{cls})$). Please find details in Tip-adapter-F [57]. To
 294 introduce PLOT to this framework, we first used the feature map to replace the global feature and
 295 then learned multiple linear models. As a result, with different local features and different linear
 296 models, we can obtain a $M \times N$ distance matrix and applied the Sinkhorn algorithm [6] to calculate
 297 the OT distance. Furthermore, we can apply the learned prompts as co-partner of the ensembling
 298 prompt to refine the final similarity.

299 Table 3 summarizes the few-shot recognition results of the original Tip-Adapter-F method and our
 300 adapter-based PLOT methods on ImageNet. From this table, We observe that using the OT distance
 301 can improve the performance of the adapter-based method. Using the learned prompts, we can further
 302 promote the accuracy of all settings.

303 **Q: What is the extra computation time cost of PLOT over CoOp baseline? A: Around 10%**
 304 **inference speed and 5% training time.** Please see the detailed comparisons and analysis in the
 305 supplementary materials.

306 4.5 Visualization

307 In this subsection, we provide some visualization examples of the transport plans T related to different
 308 prompts (N=4). We translate each transport plan into colorful heatmaps and resize them into their
 309 original size and combine them with the raw image. As shown in Figure 4, we provide the heatmaps
 310 of 4 categories in ImageNet. We observe that different transport plans highlight different regions of
 311 the image, which demonstrates that the learned multiple prompts are complementary. For the class
 312 “Brambling”, the prompts respectively focus on the head, tail, wing, and environment. For “Dog
 313 Sled”, the prompts are related to dogs, the sled, some ties, and the snow environment.

314 5 Conclusion

315 In this paper, we present a method, named PLOT, to learn multiple comprehensive prompts to
 316 describe diverse characteristics of one category. To avoid convergence to one point, we propose to
 317 apply the optimal transport to achieve the fine-grained alignment between both vision and language
 318 domains. We apply a two-stage optimization strategy where the inner loop fixes the prompts and
 319 learns the transport plan to calculate the cross-modality distance, and the outer loop uses this distance
 320 to optimize the prompt learner. We build our method on the base of CoOp and achieve significant
 321 improvement on the few-shot recognition task in various datasets, which demonstrates the advantage
 322 to learn multiple prompts instead of a single one.

323 References

- 324 [1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In
325 *ICML*, pages 214–223, 2017.
- 326 [2] Emmanuel Boissard, Thibaut Le Gouic, and Jean-Michel Loubes. Distribution’s template estimate with
327 wasserstein metrics. *Bernoulli*, 21(2):740–759, 2015.
- 328 [3] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101–mining discriminative components
329 with random forests. In *ECCV*, pages 446–461, 2014.
- 330 [4] Liqun Chen, Yizhe Zhang, Ruiyi Zhang, Chenyang Tao, Zhe Gan, Haichao Zhang, Bai Li, Dinghan Shen,
331 Changyou Chen, and Lawrence Carin. Improving sequence-to-sequence learning via optimal transport.
332 *arXiv preprint arXiv:1901.06283*, 2019.
- 333 [5] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing
334 textures in the wild. In *CVPR*, pages 3606–3613, 2014.
- 335 [6] Marco Cuturi. Sinkhorn distances: lightspeed computation of optimal transport. In *NeurIPS*, volume 2,
336 page 4, 2013.
- 337 [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical
338 image database. In *CVPR*, pages 248–255, 2009.
- 339 [8] Zi-Yi Dou, Yichong Xu, Zhe Gan, Jianfeng Wang, Shuohang Wang, Lijuan Wang, Chenguang Zhu, Zicheng
340 Liu, Michael Zeng, et al. An empirical study of training end-to-end vision-and-language transformers.
341 *arXiv preprint arXiv:2111.02387*, 2021.
- 342 [9] Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training examples:
343 An incremental bayesian approach tested on 101 object categories. In *CVPRW*, pages 178–178, 2004.
- 344 [10] Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li,
345 and Yu Qiao. Clip-adapter: Better vision-language models with feature adapters. *arXiv preprint*
346 *arXiv:2110.04544*, 2021.
- 347 [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
348 In *CVPR*, pages 770–778, 2016.
- 349 [12] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep
350 learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied*
351 *Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.
- 352 [13] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai,
353 Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of
354 out-of-distribution generalization. *arXiv preprint arXiv:2006.16241*, 2020.
- 355 [14] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial
356 examples. *arXiv preprint arXiv:1907.07174*, 2019.
- 357 [15] Yicong Hong, Qi Wu, Yuankai Qi, Cristian Rodriguez-Opazo, and Stephen Gould. Vln bert: A recurrent
358 vision-and-language bert for navigation. In *CVPR*, pages 1643–1653, 2021.
- 359 [16] Aashi Jain, Mandy Guo, Krishna Srinivasan, Ting Chen, Sneha Kudugunta, Chao Jia, Yinfei Yang, and
360 Jason Baldridge. Mural: multimodal, multitask retrieval across languages. *arXiv preprint arXiv:2109.05125*,
361 2021.
- 362 [17] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung,
363 Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text
364 supervision. In *ICML*, pages 4904–4916, 2021.
- 365 [18] Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. How can we know what language models
366 know? *TACL*, 8:423–438, 2020.
- 367 [19] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion.
368 Mdetr-modulated detection for end-to-end multi-modal understanding. In *ICCV*, pages 1780–1790, 2021.
- 369 [20] Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or
370 region supervision. In *ICML*, pages 5583–5594, 2021.
- 371 [21] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained
372 categorization. In *ICCVW*, pages 554–561, 2013.
- 373 [22] Charlotte Laclau, Ievgen Redko, Basarab Matei, Younes Bennani, and Vincent Brault. Co-clustering
374 through optimal transport. In *ICML*, pages 1955–1964, 2017.
- 375 [23] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training
376 for unified vision-language understanding and generation. *arXiv preprint arXiv:2201.12086*, 2022.

- 377 [24] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong
378 Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *NeurIPS*,
379 34, 2021.
- 380 [25] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and
381 performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- 382 [26] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv*
383 *preprint arXiv:2101.00190*, 2021.
- 384 [27] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train,
385 prompt, and predict: A systematic survey of prompting methods in natural language processing. *arXiv*
386 *preprint arXiv:2107.13586*, 2021.
- 387 [28] Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. Gpt understands,
388 too. *arXiv preprint arXiv:2103.10385*, 2021.
- 389 [29] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual
390 classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- 391 [30] Gaspard Monge. Mémoire sur la théorie des déblais et des remblais. *Histoire de l'Académie Royale des*
392 *Sciences de Paris*, 1781.
- 393 [31] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya
394 Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided
395 diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- 396 [32] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of
397 classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, pages
398 722–729, 2008.
- 399 [33] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *CVPR*, pages
400 3498–3505, 2012.
- 401 [34] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen,
402 Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep
403 learning library. *NeurIPS*, 2019.
- 404 [35] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven
405 manipulation of stylegan imagery. In *ICCV*, pages 2085–2094, 2021.
- 406 [36] Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H Miller, and
407 Sebastian Riedel. Language models as knowledge bases? *arXiv preprint arXiv:1909.01066*, 2019.
- 408 [37] Gabriel Peyre and Marco Cuturi. Computational optimal transport. *Foundations and Trends in Machine*
409 *Learning*, 11(5-6):355–607, 2019.
- 410 [38] Nina Poerner, Ulli Waltinger, and Hinrich Schütze. Bert is not a knowledge base (yet): Factual knowledge
411 vs. name-based reasoning in unsupervised qa. *arXiv preprint arXiv:1911.03681*, 2019.
- 412 [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish
413 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from
414 natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021.
- 415 [40] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
416 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 417 [41] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional
418 image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 419 [42] Yongming Rao, Wenliang Zhao, Guangyi Chen, Yansong Tang, Zheng Zhu, Guan Huang, Jie Zhou, and
420 Jiwen Lu. Densclip: Language-guided dense prediction with context-aware prompting. *arXiv preprint*
421 *arXiv:2112.01518*, 2021.
- 422 [43] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishal Shankar. Do imagenet classifiers
423 generalize to imagenet? *arXiv preprint arXiv:1902.10811*, 2019.
- 424 [44] Tim Salimans, Han Zhang, Alec Radford, and Dimitris Metaxas. Improving gans using optimal transport.
425 *ICLR*, 2018.
- 426 [45] Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Autoprompt: Eliciting
427 knowledge from language models with automatically generated prompts. *arXiv preprint arXiv:2010.15980*,
428 2020.
- 429 [46] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions
430 classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- 431 [47] Guy Tevet, Brian Gordon, Amir Hertz, Amit H Bermano, and Daniel Cohen-Or. Motionclip: Exposing
432 human motion generation to clip space. *arXiv preprint arXiv:2203.08063*, 2022.

- 433 [48] Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Multimodal
434 few-shot learning with frozen language models. *NeurIPS*, 34:200–212, 2021.
- 435 [49] Ruibo Tu, Kun Zhang, Hedvig Kjellström, and Cheng Zhang. Optimal transport for causal discovery.
436 *ICLR*, 2022.
- 437 [50] Cédric Villani. *Optimal transport: old and new*, volume 338. Springer, 2009.
- 438 [51] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by
439 penalizing local predictive power. *NeurIPS*, 32, 2019.
- 440 [52] Mengmeng Wang, Jiazheng Xing, and Yong Liu. Actionclip: A new paradigm for video action recognition.
441 *arXiv preprint arXiv:2109.08472*, 2021.
- 442 [53] Wenhui Wang, Hangbo Bao, Li Dong, and Furu Wei. Vlmo: Unified vision-language pre-training with
443 mixture-of-modality-experts. *arXiv preprint arXiv:2111.02358*, 2021.
- 444 [54] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-
445 scale scene recognition from abbey to zoo. In *CVPR*, pages 3485–3492, 2010.
- 446 [55] Hongteng Xu, Dixin Luo, Hongyuan Zha, and Lawrence Carin Duke. Gromov-wasserstein learning for
447 graph matching and node embedding. In *ICML*, pages 6932–6941, 2019.
- 448 [56] Jingjing Xu, Hao Zhou, Chun Gan, Zaixiang Zheng, and Lei Li. Vocabulary learning via optimal transport
449 for neural machine translation. *arXiv preprint arXiv:2012.15671*, 2020.
- 450 [57] Renrui Zhang, Rongyao Fang, Peng Gao, Wei Zhang, Kunchang Li, Jifeng Dai, Yu Qiao, and Hong-
451 sheng Li. Tip-adapter: Training-free clip-adapter for better vision-language modeling. *arXiv preprint*
452 *arXiv:2111.03930*, 2021.
- 453 [58] Renrui Zhang, Longtian Qiu, Wei Zhang, and Ziyao Zeng. Vt-clip: Enhancing vision-language models
454 with visual-guided texts. *arXiv preprint arXiv:2112.02399*, 2021.
- 455 [59] He Zhao, Dinh Phung, Viet Huynh, Trung Le, and Wray Buntine. Neural topic model via optimal transport.
456 *ICLR*, 2021.
- 457 [60] Wenliang Zhao, Yongming Rao, Ziyi Wang, Jiwen Lu, and Jie Zhou. Towards interpretable deep metric
458 learning with structural matching. In *ICCV*, pages 9887–9896, 2021.
- 459 [61] Chong Zhou, Chen Change Loy, and Bo Dai. Denseclip: Extract free dense labels from clip. *arXiv preprint*
460 *arXiv:2112.01071*, 2021.
- 461 [62] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language
462 models. *arXiv preprint arXiv:2109.01134*, 2021.
- 463 [63] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for
464 vision-language models. In *CVPR*, 2022.

465 Checklist

- 466 1. For all authors...
- 467 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contribu-
468 tions and scope? [Yes]
- 469 (b) Did you describe the limitations of your work? [Yes] See our analysis in Section 4.4. 1) Our
470 method is still need few-shot data for optimization, which cannot be applied in zero-shot setting.
471 2) The method needs more computing cost than CoOp.
- 472 (c) Did you discuss any potential negative societal impacts of your work? [N/A] We propose
473 a general framework for using the vision-language pre-trained model. It is not for specific
474 applications, which does not directly involve societal issues.
- 475 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 476 2. If you are including theoretical results...
- 477 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 478 (b) Did you include complete proofs of all theoretical results? [N/A]
- 479 3. If you ran experiments...
- 480 (a) Did you include the code, data, and instructions needed to reproduce the main experimental
481 results (either in the supplemental material or as a URL)? [Yes]
- 482 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)?
483 [Yes]
- 484 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
485 multiple times)? [Yes]

- 486 (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs,
487 internal cluster, or cloud provider)? [Yes] See Section 4.2.
- 488 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 489 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 490 (b) Did you mention the license of the assets? [N/A]
- 491 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 492 (d) Did you discuss whether and how consent was obtained from people whose data you're us-
493 ing/curating? [N/A]
- 494 (e) Did you discuss whether the data you are using/curating contains personally identifiable informa-
495 tion or offensive content? [N/A]
- 496 5. If you used crowdsourcing or conducted research with human subjects...
- 497 (a) Did you include the full text of instructions given to participants and screenshots, if applicable?
498 [N/A]
- 499 (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB)
500 approvals, if applicable? [N/A]
- 501 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on
502 participant compensation? [N/A]