
Prompt Learning with Optimal Transport for Vision-Language Models

Supplementary Material

Anonymous Author(s)

Affiliation

Address

email

1 A. More Details of Optimal Transport

2 The Optimal Transport [11] is initially introduced to find a transportation plan to move simultaneously
3 several items at a minimal cost, such as moving a pile of sand to fill all the holes. Recently, it is widely
4 used for the comparison of distributions. Mathematically, given two probability density function U
5 and V over space \mathcal{X} and \mathcal{Y} , the OT (Wasserstein) distance [16] can be defined as

$$D_{OT}(U, V) = \inf_{\Gamma} \int_{\mathcal{X} \times \mathcal{Y}} C(\mathbf{x}, \mathbf{y}) d\gamma(\mathbf{x}, \mathbf{y}), \quad (1)$$

6 where $C(\mathbf{x}, \mathbf{y})$ is the cost between two points in the space $\mathcal{X} \times \mathcal{Y}$, and Γ denotes the set of transport
7 plans between support points \mathbf{x} and \mathbf{y} (e.g. $\gamma(\mathbf{x}, \mathbf{y})$). We can regard two probability density functions
8 U and V as piles and holes and C is the cost function of moving a unit of sand.

9 In our problem of multiple prompts learning, we formulate the sets of visual features and prompt
10 features as two discrete distributions as

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

11 where \mathbf{u} and \mathbf{v} are the discrete probability vectors that sum to 1, and $\delta_{\mathbf{f}}$ is a Dirac delta function
12 placed at support point \mathbf{f} in the embedding space. Given two support points \mathbf{f}_m and \mathbf{g}_n , the cost
13 function is written as $C(\mathbf{f}_m, \mathbf{g}_n) = 1 - \text{sim}(\mathbf{f}_m, \mathbf{g}_n) = 1 - \frac{\mathbf{f}_m^T \mathbf{g}_n}{\|\mathbf{f}_m\| \|\mathbf{g}_n\|}$. For simply, in this discrete
14 situation, $C \in \mathbb{R}^{M \times N}$ is a cost matrix in which each point denotes the cost between \mathbf{f}_m and \mathbf{g}_n .
15 Then, the total distance of these two distributions is written as:

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

16 where the $\mathbf{T} \in \mathbb{R}^{M \times N}$ is a matrix of transport plan, which is learned to minimize the total distance.
17 Each point $T_{m,n}$ in \mathbf{T} is a weight of local cost $C_{m,n}$.

18 The optimization problem of optimal transport is formulated as:

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

19 These constraints of \mathbf{T} is used to match its marginal distributions and original discrete distributions
20 in Eq (2). In our framework, we treat visual features \mathbf{f}_m and prompt features \mathbf{g}_n equally and thus
21 $\mathbf{u} = \mathbf{1}_{M \times 1} / M$ and $\mathbf{v} = \mathbf{1}_{N \times 1} / N$.

Algorithm 1: Prompt Learning with Optimal Transport

Input: Training few-shot image data: $\mathbf{X} = \{\mathbf{x}\}$, pretrained CLIP model f and g . number of prompts N , entropy parameter λ , maximum number of iterations in inner and outer loops T_{in}, T_{out} .

Output: The parameters of prompts $\{\mathbf{vec}_{l,n}\}_{l=1,n=1}^{L,N}$

```
1: Initialize  $\{\mathbf{vec}_{l,n}\}$ 
2: for  $t_{out} = 1, 2, \dots, T_{out}$  in the outer loop do
3:   Obtain a visual feature set  $\mathbf{F} \in \mathbb{R}^{M \times C}$  with the visual encoder of CLIP;
4:   Generate prompt feature set  $\mathbf{G}_k \in \mathbb{R}^{N \times C}$  of each class with the textual encoder;
5:   Calculate the cost matrix  $\mathbf{C}_k = \mathbf{1} - \mathbf{F}^T \mathbf{G}_k \in \mathbb{R}^{M \times N}$  of each class
6:   Calculate the OT distance with an inner loop: Initialize the  $\mathbf{v}^0 = \mathbf{1}$ ,  $\delta = 0.01$  and  $\Delta_v = \infty$ 
7:   for  $t_{in} = 1, 2, \dots, T_{in}$  do
8:     Update  $\mathbf{u}^{t_{in}} = \mathbf{u} / ((\exp(-\mathbf{C}/\lambda) \mathbf{v}^{t_{in}-1})$ 
9:     Update  $\mathbf{v}^{t_{in}} = \mathbf{v} / ((\exp(-\mathbf{C}/\lambda)^T \mathbf{u}^{t_{in}})$ 
10:    Update  $\Delta_v = \sum |\mathbf{v}^{t_{in}} - \mathbf{v}^{t_{in}-1}| / N$ 
11:    if  $\Delta_v < \delta$  then
12:      break
13:    end if
14:  end for
15:  Obtain optimal transport plan as  $\mathbf{T}_k^* = \text{diag}(\mathbf{u}^t) \exp(-\mathbf{C}_k/\lambda) \text{diag}(\mathbf{v}^t)$ ,
16:  Calculate the OT distance  $d_{OT}(k) = \langle \mathbf{T}_k^*, \mathbf{C}_k \rangle$ 
17:  Calculate the classification probability  $p_{ot}(y = k | \mathbf{x})$  with the OT distance
18:  Update the parameters of prompts with cross-entropy loss  $L_{CE}$ 
19: end for
20: return  $\{\mathbf{vec}_{l,n}\}_{l=1,n=1}^{L,N}$ 
```

22 As directly optimizing the above objective is always time-consuming, we apply the Sinkhorn dis-
23 tance [3] to use an entropic constraint for fast optimization. The optimization problem with a
24 Lagrange multiplier of the entropy constraint is:

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

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

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

27 where t denotes iteration and in each iteration $\mathbf{u}^t = \mathbf{u} / ((\exp(-\mathbf{C}/\lambda) \mathbf{v}^{t-1})$ and $\mathbf{v}^t =$
28 $\mathbf{v} / ((\exp(-\mathbf{C}/\lambda)^T \mathbf{u}^t)$, with the initiation $\mathbf{v}^0 = \mathbf{1}$. The detailed algorithm is shown in Algorithm 1

29 B. Dataset Details

30 The datasets we used in the experiments follow CoOp [19], which include 11 datasets for few-shot
31 visual recognition and 4 ImageNet-based datasets for generalization (robustness) evaluation. The
32 details of each dataset are shown in Table 1, including the number of classes, the sizes of training and
33 testing sets, and original tasks.

34 C. Few-shot Recognition Accuracy

35 In Section 4.3.1, we provide a line chart to show and compare the performance of PLOT and CoOp.
36 In this section, we provide detailed performance results on all 11 few-shot recognition datasets in
37 Table 2, where we use gray for our method and white for CoOp. To highlight, we respectively use
38 dark cyan and light cyan to represent the performance of PLOT and CoOp on the average of all 11
39 datasets. We repeat all experiments 3 times and report the mean and standard deviation in the table.

Table 1: The detailed statistics of datasets used in experiments.

Dataset	Classes	Training size	Testing size	Task
Caltech101 [5]	100	4,128	2,465	Object recognition
DTD [2]	47	2,820	1,692	Texture recognition
EuroSAT [6]	10	13,500	8,100	Satellite image recognition
FGVCAircraft [10]	100	3,334	3,333	Fine-grained aircraft recognition
Flowers102 [12]	102	4,093	2,463	Fine-grained flowers recognition
Food101 [1]	101	50,500	30,300	Fine-grained food recognition
ImageNet [4]	1,000	1.28M	50,000	Object recognition
OxfordPets [13]	37	2,944	3,669	Fine-grained pets recognition
StanfordCars [9]	196	6,509	8,041	Fine-grained car recognition
SUN397 [18]	397	15,880	19,850	Scene recognition
UCF101 [15]	101	7,639	3,783	Action recognition
ImageNetV2 [14]	1,000	-	10,000	Robustness of collocation
ImageNet-Sketch [17]	1000	-	50,889	Robustness of sketch domain
ImageNet-A [8]	200	-	7,500	Robustness of adversarial attack
ImageNet-R [7]	200	-	30,000	Robustness of multi-domains

Table 2: The few-shot visual recognition accuracy on 11 datasets.

Dataset	Methods	1 shot	2 shots	4 shots	8 shots	16 shots
Caltech101	PLOT	89.83 \pm 0.33	90.67 \pm 0.21	90.80 \pm 0.20	91.54 \pm 0.33	92.24 \pm 0.38
	CoOp	87.51 \pm 1.02	87.84 \pm 1.10	89.52 \pm 0.80	90.28 \pm 0.42	91.99 \pm 0.31
DTD	PLOT	46.55 \pm 2.62	51.24 \pm 1.95	56.03 \pm 0.43	61.70 \pm 0.35	65.60 \pm 0.82
	CoOp	43.62 \pm 1.96	45.35 \pm 0.31	53.94 \pm 1.37	59.69 \pm 0.13	62.51 \pm 0.25
EuroSAT	PLOT	54.05 \pm 5.95	64.21 \pm 1.90	72.36 \pm 2.29	78.15 \pm 2.65	82.23 \pm 0.91
	CoOp	52.12 \pm 5.46	59.00 \pm 3.48	68.61 \pm 3.54	77.08 \pm 2.42	83.69 \pm 0.47
FGVCAircraft	PLOT	17.90 \pm 0.09	18.94 \pm 0.44	22.36 \pm 0.42	26.17 \pm 0.29	31.49 \pm 0.89
	CoOp	8.59 \pm 5.79	16.52 \pm 2.38	20.63 \pm 2.46	26.63 \pm 0.86	31.43 \pm 0.96
Flowers102	PLOT	71.72 \pm 0.97	81.19 \pm 0.79	87.82 \pm 0.20	92.43 \pm 0.25	94.76 \pm 0.34
	CoOp	67.98 \pm 1.98	77.58 \pm 1.46	86.10 \pm 1.05	91.27 \pm 0.83	94.49 \pm 0.40
FOOD101	PLOT	77.74 \pm 0.47	77.70 \pm 0.02	77.21 \pm 0.43	75.31 \pm 0.30	77.09 \pm 0.18
	CoOp	74.25 \pm 1.52	72.61 \pm 1.33	73.49 \pm 2.03	71.58 \pm 0.79	74.48 \pm 0.15
ImageNet	PLOT	59.54 \pm 0.16	60.64 \pm 0.06	61.49 \pm 0.23	61.92 \pm 0.09	63.01 \pm 0.13
	CoOp	56.99 \pm 1.03	56.40 \pm 0.87	58.48 \pm 0.47	60.39 \pm 0.57	61.91 \pm 0.17
OxfordPets	PLOT	87.49 \pm 0.57	86.64 \pm 0.63	88.63 \pm 0.26	87.39 \pm 0.74	87.21 \pm 0.40
	CoOp	85.99 \pm 0.28	82.22 \pm 2.15	86.65 \pm 0.97	85.36 \pm 1.00	87.02 \pm 0.89
StanfordCars	PLOT	56.60 \pm 0.36	57.52 \pm 0.71	63.41 \pm 0.29	67.03 \pm 0.50	72.80 \pm 0.75
	CoOp	55.81 \pm 1.67	58.41 \pm 0.43	62.74 \pm 0.16	67.64 \pm 0.06	73.60 \pm 0.19
SUN397	PLOT	62.47 \pm 0.43	61.71 \pm 0.65	65.09 \pm 0.43	67.48 \pm 0.04	69.96 \pm 0.24
	CoOp	60.12 \pm 0.82	59.60 \pm 0.76	63.24 \pm 0.63	65.77 \pm 0.02	68.36 \pm 0.66
UCF101	PLOT	64.53 \pm 0.70	66.83 \pm 0.43	69.60 \pm 0.67	74.45 \pm 0.50	77.26 \pm 0.64
	CoOp	62.13 \pm 1.14	64.05 \pm 0.99	67.79 \pm 0.71	72.71 \pm 0.50	76.90 \pm 0.50
Average	PLOT	62.59 \pm 1.13	65.23 \pm 0.72	68.60 \pm 0.52	71.23 \pm 0.51	73.94 \pm 0.54
	CoOp	59.56 \pm 2.06	61.78 \pm 1.39	66.47 \pm 1.29	69.85 \pm 0.69	73.33 \pm 0.42

40 D. Computation Cost Evaluation

41 As shown in Table 3, we provide the comparison of the training time and inference seed of the
42 baseline method CoOp [19] and our PLOT with the different number of prompts. We report the
43 one-epoch time training on the 1-shot setting of the Food101 [1] dataset and the number of images
44 processed by the model in 1 second. Taking $N = 4$ as an example, PLOT only reduces the 9.2%

Table 3: The training and inference time comparison.

Settings	CoOp	PLOT(N=1)	PLOT(N=2)	PLOT(N=4)	PLOT(N=8)
Training Time (s)	1.127	1.135	1.148	1.182	1.267
Inference Time (images/s)	719.1	714.4	690.7	653.0	519.8

Table 4: The nearest words for 16 context vectors of all $N = 4$ prompts learned by PLOT. N/A means non-Latin characters.

Number	Prompt 1	Prompt 2	Prompt 3	Prompt 4
1	ag	pa	trying	gaz
2	flint	as	field	white
3	leaving	wit	N/A	t
4	sot	l	icons	ario
5	tint	N/A	eclub	safe
6	tar	yl	indiffe	class
7	attn	N/A	ts	represented
8	2	job	cold	attend
9	rollingstones	built	yeah	vie
10	N/A	brought	band	recognized
11	N/A	or	love	old
12	bel	j	late	stel
13	head	ag	industry	awhile
14	artifact	bad	N/A	ded
15	an	chie	across	these
16	5	in	actual	visiting

inference speed and requires an extra 4.9% training time, which is acceptable given the performance improvement.

E. The Interpretation of the Learned Prompts

The learned prompts are difficult to be understood by humans since the parameters are optimized in the continuous space [19]. CoOp proposes to use the word which is nearest to learned prompts in the embedding space to visualize the prompts. Following this manner, we show the nearest words of our learned prompts in Table 4. Similar to CoOp, most words can not be directly understood by human logic. However, we still find the relations between the learned prompts and the corresponding optimal transport plan. As shown in Figure 4 in the main paper, we can observe that the optimal transport plan for Prompt 1 always focuses on the “head”, such as the head of “brambling”, the head of “rooster”, and even the head of “aircraft carrier”. It is because the word “head” is in Prompt 1. Similarly, we can find that Prompt 4 prefers the white part of images, such as the white environment in the image of “brambling” and the snow in the image of “dog sled”. It demonstrates that the learned multiple prompts focus on different characteristics of categories.

F. Visualization of the Failure Cases

To better understand the method and further discover the reason of the failure cases, we visualize the attention maps of some failure cases. As shown in Figure 1, we showed two failure examples with class "2000 AM General Hummer" in the StanfordCars dataset. During the training, we set the number of prompts as 4, but in these visualization results, we found that some of learned prompts remarkably coincide with each other. These prompts can be roughly divided into two classes: Foreground and Background. For example, in both images, prompt 2 (right top) and 3 (left down) focus on the foreground car, while the others focus on the background. It demonstrates that not all classes have multiple complementary attributes, which motivates us to go further to learn the dynamic local prompts numbers to reduce the computational load in the future.

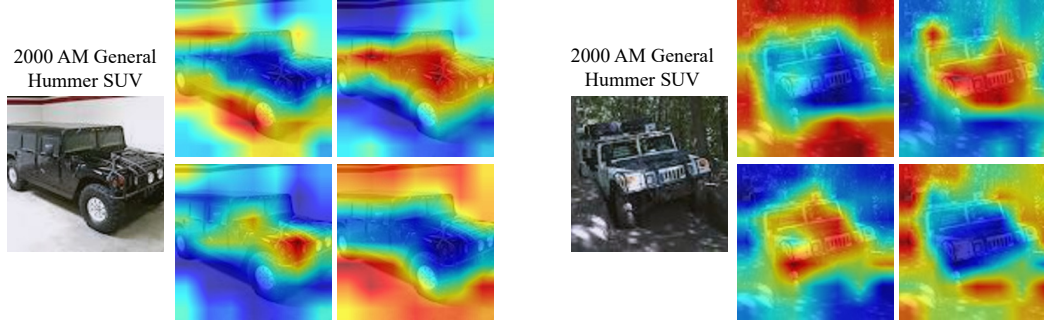


Figure 1: *Failure Visualization. We provide the heatmaps of transport plan T related to each prompt on 2 failure examples in the StanfordCars dataset.*

69 G. Code

70 The code used in our experiments is provided in the supplementary material. We implement the
71 method with Pytorch 1.7. Our implementation is based on the publicly released code of CoOp [19].

72 References

- 73 [1] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative components
74 with random forests. In *ECCV*, pages 446–461, 2014.
- 75 [2] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing
76 textures in the wild. In *CVPR*, pages 3606–3613, 2014.
- 77 [3] Marco Cuturi. Sinkhorn distances: lightspeed computation of optimal transport. In *NeurIPS*, volume 2,
78 page 4, 2013.
- 79 [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical
80 image database. In *CVPR*, pages 248–255, 2009.
- 81 [5] Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training examples:
82 An incremental bayesian approach tested on 101 object categories. In *CVPRW*, pages 178–178, 2004.
- 83 [6] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep
84 learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied*
85 *Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.
- 86 [7] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai,
87 Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of
88 out-of-distribution generalization. *arXiv preprint arXiv:2006.16241*, 2020.
- 89 [8] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial
90 examples. *arXiv preprint arXiv:1907.07174*, 2019.
- 91 [9] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained
92 categorization. In *ICCVW*, pages 554–561, 2013.
- 93 [10] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual
94 classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- 95 [11] Gaspard Monge. Mémoire sur la théorie des déblais et des remblais. *Histoire de l’Académie Royale des*
96 *Sciences de Paris*, 1781.
- 97 [12] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of
98 classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, pages
99 722–729, 2008.
- 100 [13] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *CVPR*, pages
101 3498–3505, 2012.
- 102 [14] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishal Shankar. Do imagenet classifiers
103 generalize to imagenet? *arXiv preprint arXiv:1902.10811*, 2019.
- 104 [15] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions
105 classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- 106 [16] Matthew Thorpe. Introduction to optimal transport. *Lecture Notes*, 2019.

- 107 [17] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by
108 penalizing local predictive power. *NeurIPS*, 32, 2019.
- 109 [18] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-
110 scale scene recognition from abbey to zoo. In *CVPR*, pages 3485–3492, 2010.
- 111 [19] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language
112 models. *arXiv preprint arXiv:2109.01134*, 2021.