

THU-AM-276

# Open-set Fine-grained Retrieval via Prompting Vision-Language Evaluator

Shijie Wang<sup>1</sup>, Jianlong Chang<sup>2</sup>, Haojie Li<sup>1,3\*</sup>, Zihui Wang<sup>1</sup>, Wanli Ouyang<sup>4</sup>, Qi Tian<sup>2</sup>

<sup>1</sup>International School of Information Science & Engineering, Dalian University of Technology, China

<sup>2</sup> Huawei Cloud & AI, China

<sup>3</sup>College of Computer and Engineering, Shandong University of Science and Technology, China

<sup>4</sup> Sense Time Computer Vision Research Group, The University of Sydney, Australia

JUNE 18-22, 2023

CVPR

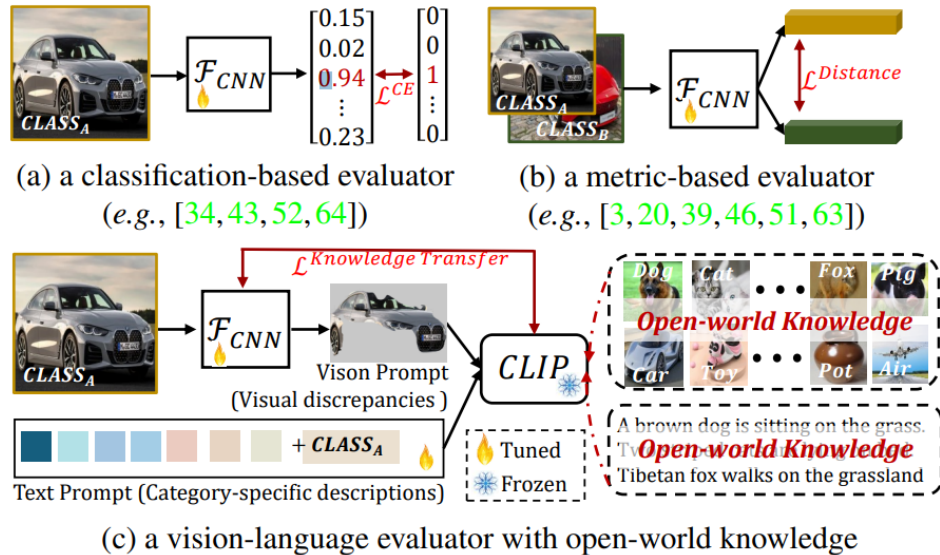


VANCOUVER, CANADA



HUAWEI

# Motivation



## Problem

Current works focus on **close-set visual concepts**, where all the subcategories are pre-defined, and make it hard to **capture discriminative knowledge from unknown subcategories**, consequently failing to handle unknown subcategories in open-world scenarios.

# Method

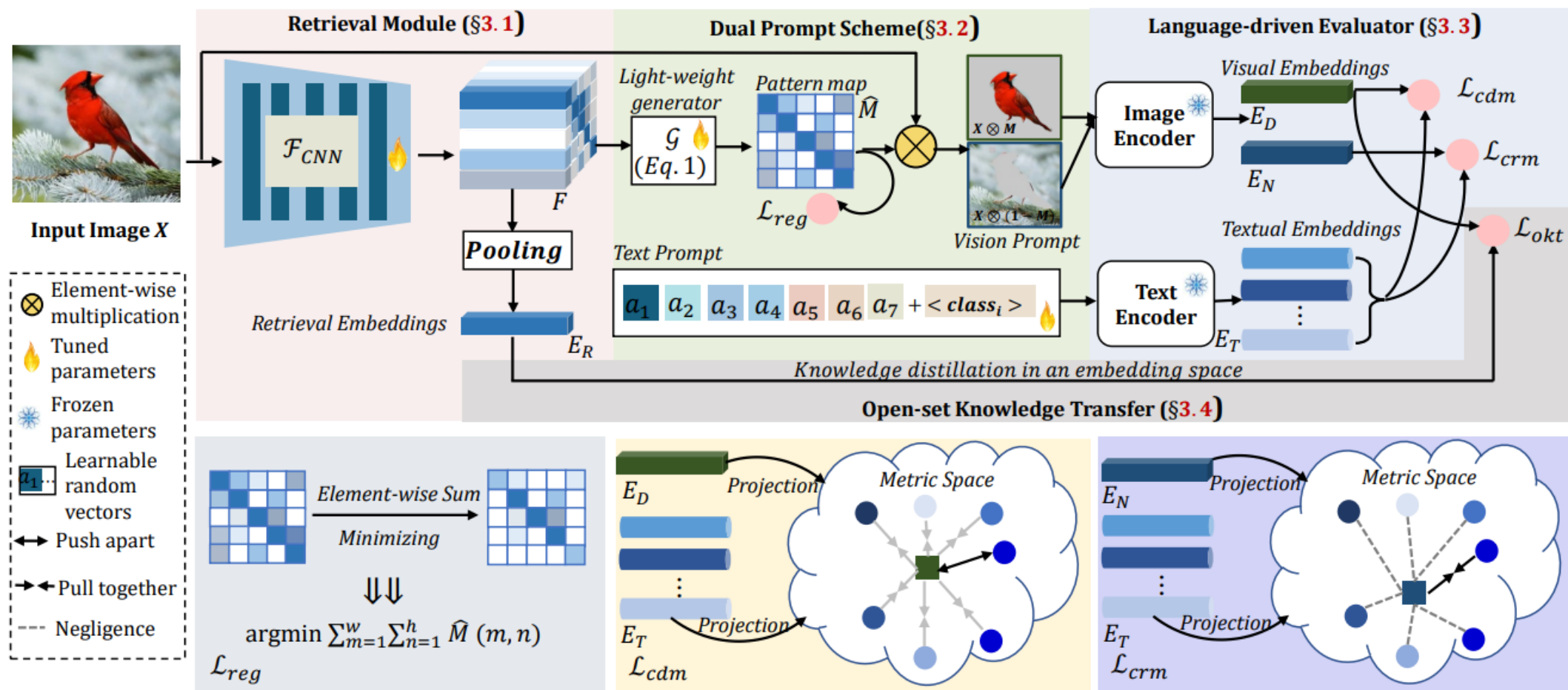
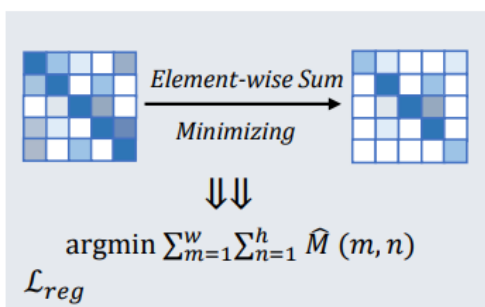
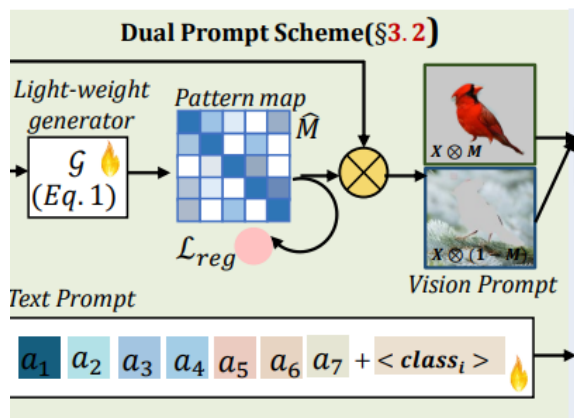


Figure 2. Detailed illustration of **prompting vision-language evaluator**. See §3 for more details.

# Experiments

Method	Arch	CUB-200-2011				Stanford Cars 196				FGVC Aircraft			
		1	2	4	8	1	2	4	8	1	2	4	8
SCDA <small>TIP<sub>17</sub></small> [52]	R50	57.3	70.2	81.0	88.4	48.3	60.2	71.8	81.8	56.5	67.7	77.6	85.7
CRL <small>IJCAI<sub>18</sub></small> [64]	R50	62.5	74.2	82.9	89.7	57.8	69.1	78.6	86.6	61.1	71.6	80.9	88.2
CEP <small>ECCV<sub>20</sub></small> [3]	R50	69.2	79.2	86.9	91.6	89.3	93.9	96.6	98.1	-	-	-	-
HDCL <small>IJON<sub>21</sub></small> [59]	R50	69.5	79.6	86.8	92.4	84.4	90.1	94.1	96.5	71.1	81.0	88.3	93.3
DGCRL <small>AAAI<sub>19</sub></small> [65]	R50	67.9	79.1	86.2	91.8	75.9	83.9	89.7	94.0	70.1	79.6	88.0	93.0
DCML <small>CVPR<sub>21</sub></small> [62]	R50	68.4	77.9	86.1	91.7	85.2	91.8	96.0	98.0	-	-	-	-
DAS <small>ECCV<sub>22</sub></small> [27]	R50	69.2	79.3	87.1	92.6	87.8	93.2	96.0	97.9	-	-	-	-
IBC <small>ICML<sub>21</sub></small> [41]	R50	70.3	80.3	87.6	92.7	88.1	93.3	96.2	98.2	-	-	-	-
NIA <small>CVPR<sub>22</sub></small> [40]	R50	70.5	80.6	-	-	89.1	93.4	-	-	-	-	-	-
Proxy <small>CVPR<sub>21</sub></small> [17]	BN	71.1	80.4	87.4	92.5	88.3	93.1	95.7	97.5	-	-	-	-
HIST <small>CVPR<sub>22</sub></small> [25]	R50	71.4	81.1	88.1	-	89.6	93.9	96.4	-	-	-	-	-
ETLR <small>CVPR<sub>21</sub></small> [18]	BN	72.1	81.3	87.6	-	89.6	94.0	96.5	-	-	-	-	-
PNCA++ <small>ECCV<sub>20</sub></small> [46]	R50	72.2	82.0	89.2	93.5	90.1	94.5	97.0	98.4	-	-	-	-
<b>Our PLEor</b>	R50	<b>74.8</b>	<b>84.5</b>	<b>91.3</b>	<b>94.9</b>	<b>94.4</b>	<b>96.9</b>	<b>98.3</b>	<b>98.9</b>	<b>86.3</b>	<b>91.7</b>	<b>95.1</b>	<b>96.7</b>

# Dual Prompt Scheme



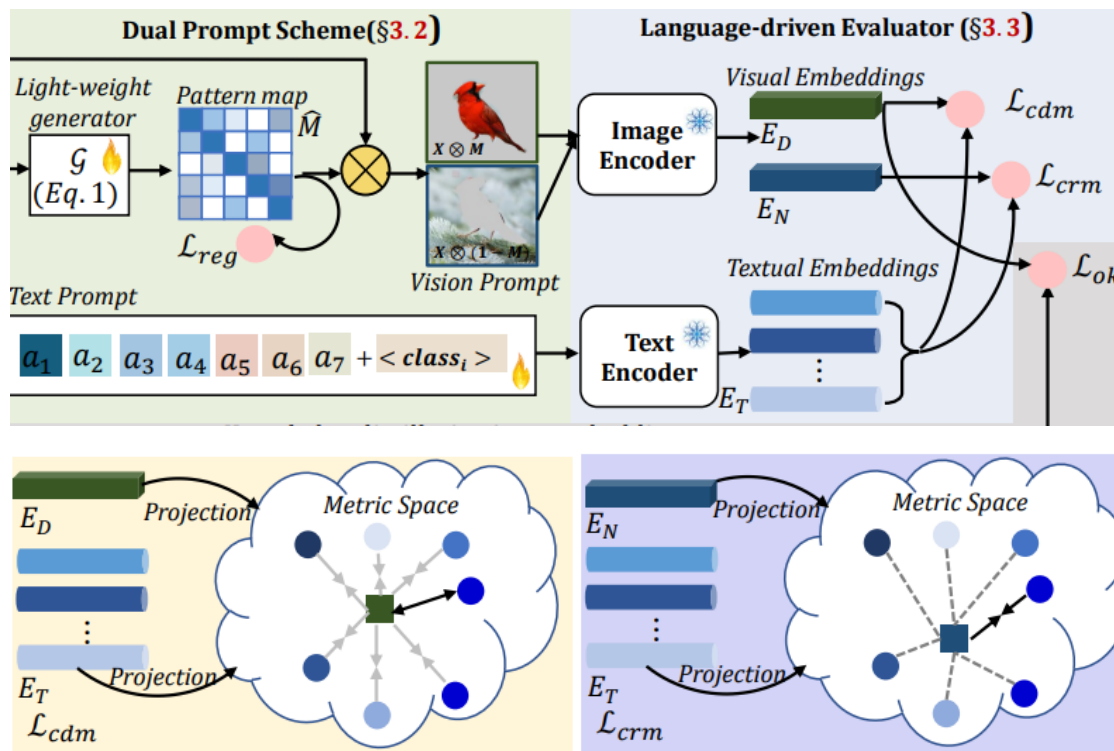
- **Vision Prompt:** To obtain the category-specific discrepancies, the vision prompt aims to project the semantic features into a new space where the location, scale and intensity of these discrepancies are specified.

$$\hat{M} = \sigma(\mathcal{G}(\mathbf{F})), \quad \mathbf{V}_D = \mathbf{X} \otimes \mathbf{M}, \quad \mathbf{V}_N = \mathbf{X} \otimes (1 - \mathbf{M}),$$

- **Text Prompt:** A text prompt is designed to generate appropriate text descriptions automatically via keeping semantically coherent with the category-specific vision prompt.

$$\mathcal{P}_{\text{class}} = (a_1, a_2, \dots, a_i, \dots, a_k, \langle class \rangle),$$

# Vision-Language Evaluator



- **Evaluator:** the contrastive objective of vision-language evaluator encourages the pre-trained CLIP model to locate the category-specific descriptions in vision prompt and generate the category-specific semantics into text prompt.

$$\mathcal{L}_{cdm} = - \sum_{i=1}^N y_i \cdot \log \frac{\exp(\cos \langle \mathbf{E}_D, \mathbf{E}_T^i \rangle / \tau)}{\sum_{i=1}^N \exp(\cos \langle \mathbf{E}_D, \mathbf{E}_T^i \rangle / \tau)},$$

$$\mathcal{L}_{crm} = - \sum_{i=1}^N y_i \cdot \log \left( 1 - \frac{\exp(\cos \langle \mathbf{E}_N, \mathbf{E}_T^i \rangle / \tau)}{\sum_{i=1}^N \exp(\cos \langle \mathbf{E}_N, \mathbf{E}_T^i \rangle / \tau)} \right). \quad (7)$$

# Ablation Experiments

Table 1. Comparison of performance and efficiency on CUB-200-2011 using different combinations of constraints. The first row indicates that we use the traditional classification-based classifier (*i.e.*, ResNet-50) as supervision, to replace the proposed PLEor for comparison. "Time" is the time of extracted retrieval embeddings.

$\mathcal{L}_{cdm}$	$\mathcal{L}_{crm}$	$\mathcal{L}_{reg}$	$\mathcal{L}_{okt}$	Recall@1	Time
				66.3%	21.1ms
✓				72.1%	42.3ms
✓	✓			74.4%	42.3ms
✓	✓	✓		75.1%	42.3ms
✓	✓	✓	✓	<b>74.8%</b>	<b>21.1ms</b>

Table 2. Evaluation results of retrieval performance on CUB-200-2011 dataset with/without the prompt learning. Hand-craft prompt denotes that we use the handcrafted prompt template ("a photo of a [·].") in text prompt.

Prompt	Recall@1
CLIP + Hand-craft prompt	71.5%
CLIP + Text Prompt	73.3% <sub>+1.8</sub>
CLIP + Vision&Hand-craft Prompt	72.4% <sub>+0.9</sub>
CLIP + Vision&Text Prompt	<b>74.8%</b> <sub>+3.3</sub>

# Visualization Experiments

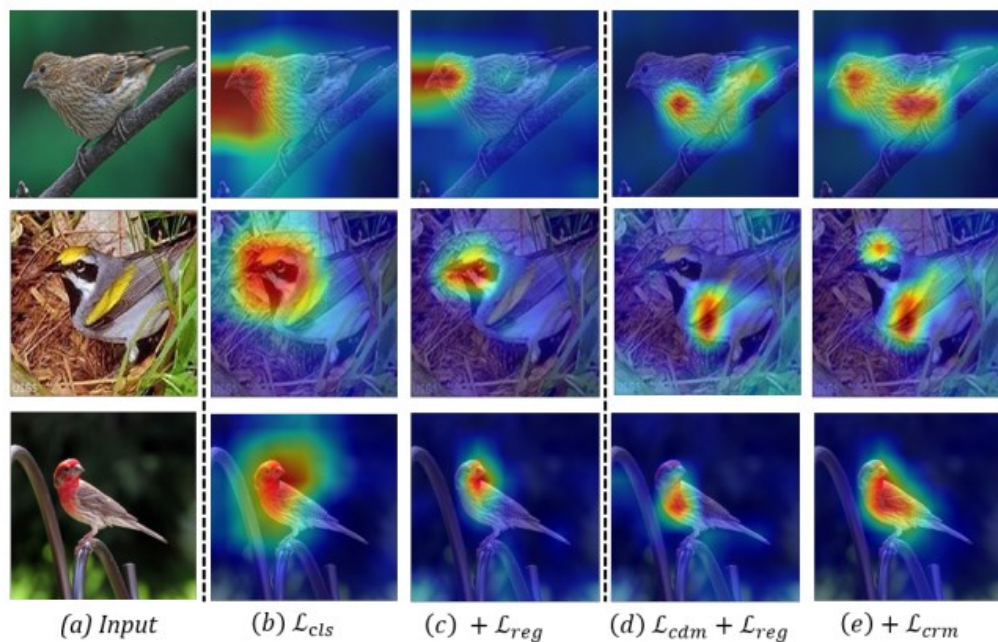


Figure 4. Visualization of vision prompt based on classification-based evaluator (b)(c) and our vision-language evaluator (d)(e), respectively.  $+\mathcal{L}$  means that we successively add this constraint, *i.e.*,  $+\mathcal{L}_{reg} = \mathcal{L}_{cls} + \mathcal{L}_{reg}$ ,  $+\mathcal{L}_{crm} = \mathcal{L}_{cdm} + \mathcal{L}_{reg} + \mathcal{L}_{crm}$ .

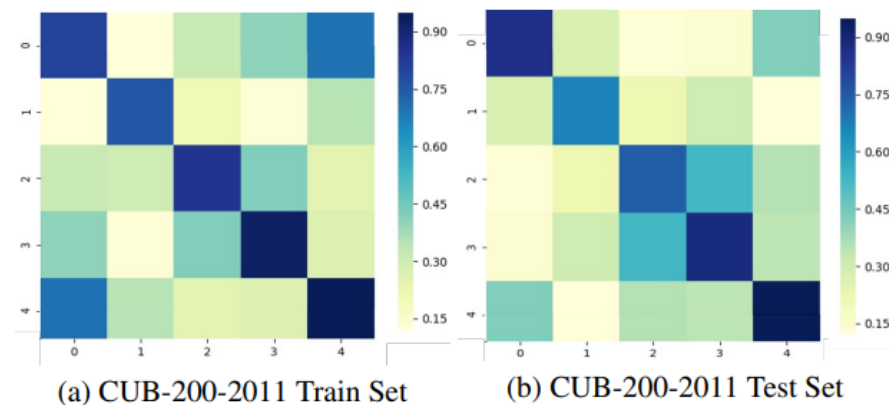


Figure 5. The nearest description for text prompt learned by PLEor, with their similarity shown in grids.



# Conclusion

- A prompting vision-language evaluator, i.e., PLEor, is proposed. It can distill the knowledge with open-world visual concepts from CLIP model to alleviate the problems behind open-set scenarios. To our best knowledge, we are the first to regard CLIP model as an evaluator specifically for OSFR task.
- PLEor provides timely insights into the adaptation of pre-trained CLIP model adopting prompt learning, and crucially, demonstrates the effectiveness of a simple modification for inputs of CLIP model in OSFR.
- PLEor achieves new state-of-the-art results compared with classification-based and metric-based evaluators, which is significant gains of 8.0% average retrieval accuracy on three widely-used OSFR datasets.