
Tree of Attributes Prompt Learning for Vision-Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Prompt learning has proven effective in adapting vision language models for
2 downstream tasks. However, existing methods usually append learnable prompt
3 tokens solely with the category names to obtain textual features, which fails to fully
4 leverage the rich context indicated in the textual category name. To address this
5 issue, we propose the Tree of Attributes Prompt learning (TAP), which first instructs
6 LLMs to generate a tree of attributes with a “concept - attribute - description”
7 structure for each associated category name, and then learn the hierarchy with
8 vision and text prompt tokens. Unlike existing methods that merely augment
9 category names with a set of unstructured descriptions, our approach essentially
10 distills structured knowledge graphs associated with class names from LLMs.
11 Furthermore, our approach introduces text and vision prompts designed to explicitly
12 learn the corresponding visual attributes, effectively serving as domain experts.
13 Additionally, the general and diverse descriptions generated based on the class
14 names may be wrong or absent in the specific given images. To address this
15 misalignment, we further introduce a vision-conditional pooling module to extract
16 instance-specific text features. Extensive experimental results demonstrate that
17 our approach outperforms state-of-the-art methods on the zero-shot base-to-novel
18 generalization as well as few-shot classification across 11 diverse datasets.

19 1 Introduction

20 Recent advancements in vision-language models (VLMs) like CLIP [33] and ALIGN [13] merge
21 the capabilities of visual perception with linguistic understanding, which have revolutionized the
22 landscape with their zero-shot learning abilities. They proficiently handle tasks on unseen data,
23 bypassing the conventional requirement for task-specific training. This feature has enabled a plethora
24 of applications, ranging from content-based image retrieval to complex visual question answering,
25 setting new benchmarks in the domain. A crucial development in this domain is the concept of
26 prompt learning, which has significantly influenced both natural language processing (NLP) [20–22]
27 and vision-only models [14, 43, 44, 51]. This approach leverages learnable prompts to guide model
28 understanding, tailoring responses to specific tasks or datasets.

29 Prompt learning, particularly in vision-language models, has garnered considerable interest due
30 to its parameter efficiency and rapid convergence [54, 53, 55, 8, 23]. Techniques like CoOp [54]
31 optimize learnable continuous prompts for few-shot image recognition, enhancing model performance
32 significantly. Recent efforts have expanded to multimodal prompt learning, optimizing prompts
33 in both visual and language domains [15, 16, 38, 19]. Despite their success, these models rely on
34 simplistic text prompts, typically formatted as “a photo of a {class}”, illustrated in Fig. 1 (a). While
35 functional, this approach lacks depth, failing to encapsulate the intricacies and finer details inherent in

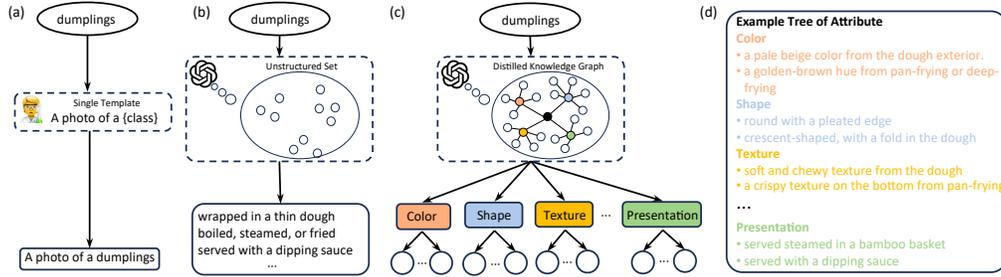


Figure 1: Illustration of the methods for CLIP text prompts formation. (a) Manually created prompt with the single “a photo of a {class}” template; (b) A unstructured set of detailed descriptions generated by LLMs; (c) The proposed Tree of Attribute that organizes the descriptions in a “concept - attribute - descriptions” structure, essentially distilling knowledge graphs from LLMs; (d) An example Tree of Attribute for “dumplings”.

36 visual data. Such limitations hinder the model’s ability to fully leverage the rich, descriptive potential
 37 offered by more detailed and contextually relevant textual information.

38 In parallel, another stream of research has been exploring the utilization of large language models
 39 (LLMs) to generate more elaborate and descriptive text prompts for enhancing zero-shot learning
 40 capabilities [26, 32, 35, 17, 30, 48, 49, 36, 52, 40]. These LLM-generated descriptions offer a wealth
 41 of detail and context, potentially enriching the model’s interpretative capabilities. However, current
 42 methodologies in integrating these descriptions often do not exploit the full potential of this richness.
 43 As shown in Fig. 1 (b), most of these approaches lack a structured framework to organize and utilize
 44 these descriptions effectively, leading to a scattergun approach where not all generated descriptions
 45 are contextually relevant or optimally aligned with the visual content. In addition, as noted in [35],
 46 descriptions generated by such paradigms are usually diverse, which covers most possibilities of the
 47 class, but include descriptions that are either likely not co-occurring, e.g. “steamed” and “fried”, or
 48 absent in the input image, e.g. “long tail” for a cat shot from the front, necessitating the need for a
 49 selective pooling mechanism for clearer image-text alignments.

50 In response to these challenges, our work introduces “Tree of Attribute Prompt learning (TAP),”
 51 a method that redefines the integration and utilization of detailed descriptions within VLMs. As
 52 indicated in Fig. 1 (c), unlike existing methods that merely augment category names with a set of
 53 unstructured descriptions, our approach essentially distills structured knowledge graphs associated
 54 with class names from LLMs. Specifically, we adopt a hierarchical, tree-like structure to systematic-
 55 ally generate and integrate descriptions, ensuring a layered and comprehensive understanding of
 56 visual content. Each branch of this tree represents a specific attribute, with finer details fleshed out in
 57 the subsequent leaves, ensuring that every aspect of the visual content is captured and represented.
 58 Furthermore, we reimagine the learnable prompt tokens as “domain experts”, each specializing in
 59 different aspects of the image, supplemented by the CLS token’s global perspective. In addition, we
 60 introduce vision-conditional layers for each expert-attribute pair, which pool the most applicable
 61 descriptions from each of the attribute sets with condition on the input image content, ensuring
 62 optimal image-text alignment. This setup not only provides a detailed, attribute-focused analysis but
 63 also harmonizes these insights with the overall context.

64 Extensive experiments in both base-to-novel generalization and few-shot classification across 11
 65 diverse datasets demonstrate the effectiveness of our method. On base-to-novel generalization, TAP
 66 achieves average performance gains of 1.07% in harmonic mean over the state-of-the-art methods,
 67 and 9.34% over the vanilla CLIP. Competitive results are also observed in few-shot classification.

68 2 Related Work

69 **Prompt Learning for Vision-Language Models.** Prompt learning bridges linguistic understanding
 70 and visual perception by guiding VLMs with text prompts, a concept originated in NLP [20–22]
 71 and adapted to vision-only [14, 43, 44, 51] and multimodal contexts [54, 53, 15, 16, 38, 19, 40, 34,
 72 36, 52, 55, 4, 23]. In the textual domain, CoOp [54] optimizes learnable continuous prompts in
 73 CLIP’s language branch for few-shot image recognition, while CoCoOp [53] addresses CoOp’s

74 overfitting issues by conditioning prompts on visual features. In the visual domain, Visual Prompt
 75 Tuning (VPT) [1] and Dual-modality Prompt Tuning (DPT) [47] enhance CLIP’s vision encoder by
 76 learning visual prompts in pixel space and dynamically generating prompts through cross-attention,
 77 respectively. TransHP [42] leverages category hierarchy for prompt learning to improve classification
 78 performance. LoGoPrompt [38] enhances classification by incorporating synthetic images with class
 79 name text as auxiliary visual prompts. MaPLe [15] explores multimodal prompt learning, jointly
 80 optimizing prompts in both vision and language branches. Other recent works have focused on
 81 regularizing prompt learning to leverage the knowledge from base VLMs effectively, demonstrating
 82 enhanced generalization in varied downstream visual tasks [16, 4, 36]. PromptSRC, for instance,
 83 introduced a self-regulating method that restricts both the vision and text prompt, demonstrating
 84 improved generalization. Distinct from these approaches, PLOT [5] and ALIGN [41] leverage
 85 Optimal Transport to align multiple prompts with local visual features, either from the multi-head
 86 self-attention layer or at a token level. Our work diverges from these methods by introducing a
 87 hierarchical "Tree of Attribute" framework derived from LLMs to structure textual descriptions and
 88 guide the learning of specialized "domain expert" tokens for attribute-level understanding.

89 **Image classification by descriptions.** There’s a growing emphasis on using visual descriptions for
 90 zero-shot recognition, moving beyond generic prompts [54, 53]. These descriptions, like the “fur
 91 pattern” or “tail shape” of a cat, provide fine-grained and distinctive characteristics. The use of LLMs
 92 like GPT-3 [3], allows for efficient generation of a broad spectrum of class-specific descriptions,
 93 offering an advantage over manually crafted templates. While this approach has been extensively
 94 researched in zero-shot contexts [17, 26, 30, 35, 48, 49, 10, 32, 28], its application in conjunction
 95 with prompt learning for few-shot tasks remains relatively unexplored[25, 19, 40, 52, 50]. Previ-
 96 ous methodologies, however, have largely utilized unstructured descriptions, lacking an organized
 97 framework for effective utilization. Our approach diverges by structuring these descriptions into a
 98 “Tree of Attribute” model, coupled with learnable visual prompts as domain experts. Additionally,
 99 LLM-generated descriptions often cover a wide range of potential class descriptions, of which not
 100 all may be pertinent to a given image, pointing to the need for a selective pooling mechanism to
 101 ensure optimal image-text alignment. We further introduce a vision-conditional pooling layer for
 102 refined image-text alignment. This structured approach not only enhances the interpretability of the
 103 model’s learning process but also significantly improves alignment accuracy between image content
 104 and descriptive text.

105 3 Methodology

106 3.1 Preliminary

107 **CLIP.** Our approach is built on the pre-trained vision-language model, CLIP [33]. Formally, let (x, c)
 108 denote the dataset, where x is an image and $c \in \{1, \dots, C\}$ are the class labels. For an image x , the
 109 vision encoder $h_I(\cdot)$ transforms it into a feature vector $\mathbf{f}_x^v = h_I(x)$. Simultaneously, each class label
 110 c is mapped to a text prompt $t_c = \text{a photo of a } \{c\}$, and converted into textual feature vectors
 111 $\mathbf{f}_c^t = h_T(t_c)$. The predicted class \hat{y} is given by:

$$\hat{y} = \operatorname{argmax}_c \cos(\mathbf{f}_x^v, \mathbf{f}_c^t) \quad (1)$$

112 where $\cos(\cdot)$ denotes cosine similarity.

113 **Image classification with class descriptions.** To improve the model’s understanding of the categories
 114 in the transfer datasets, previous works [26, 35] use more detailed descriptions from Large Language
 115 Models (LLMs) instead of the simple "a photo of a {c}" to prompt the CLIP text encoder.
 116 Under this approach, a convoluted set of descriptions is generated for a class c as $\mathcal{D}_c : \{ "c, \text{ which}$
 117 $\text{is/has/etc description.} " \}$, e.g. $c = \text{"television"}$ and $\text{description} = \text{"black or grey"}$.
 118 This classification is reformulated as

$$\hat{y} = \operatorname{argmax}_c \frac{1}{|\mathcal{D}_c|} \sum_{d \in \mathcal{D}_c} \cos(\mathbf{h}_I(x), \mathbf{h}_T(d)) \quad (2)$$

119 3.2 Overall Framework

120 We rethink the descriptions by LLM \mathcal{D}_c as nodes in knowledge graphs. While previous methods
 121 generate an unstructured set of descriptions, we distill structured knowledge graphs for each class c

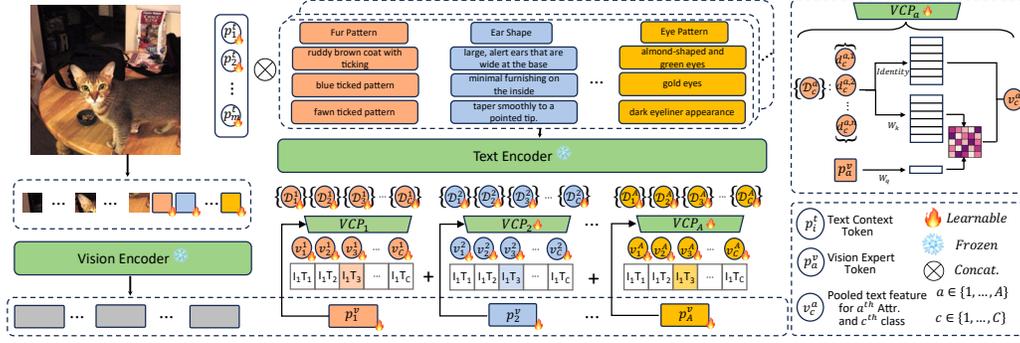


Figure 2: Overview of the proposed TAP method. TAP utilizes fine-grained descriptions from LLMs and organizes them in a Tree of Attribute. Vision expert tokens are added to the vision encoder to learn from specific attributes such as color and shape. A vision-conditional pooling layer is introduced to ensure optimal image-text alignment. Textual context tokens are also incorporated to the textual branch, shared across descriptions.

122 from LLM, in which the root node is the class name c , capturing the highest level semantics, and the
 123 leaf nodes are the detailed descriptions capturing fine-grained details. In this framework, previous
 124 paradigms only generate the leaf nodes of the graph, with the edges and graph structure missing,
 125 where the rich and inherent structure from the descriptions is overlooked. To address this limitation,
 126 we formulate our approach as a Tree of Attribute, which follows the “concept - attribute - description”
 127 structures, as illustrated in Fig. 1 (c).

128 Besides weighting the descriptions equally, previous works typically align descriptions that describe
 129 images from different aspects and at different granularities with a singular CLS token from the image
 130 encoder. However, while the use of a single CLS token is effective in certain contexts, we note that
 131 the CLS token is designed to capture the global information of an input image x [9]. As a result, even
 132 though this helps to further inform global understanding, it may fail to effectively capture the nuances
 133 and variances at the attribute level. This leads to suboptimal use of the rich descriptions. We address
 134 this by introducing a set of learnable prompt tokens that serve as domain experts in the vision branch,
 135 each of which aligns with a specific attribute-level textual embedding.

136 Additionally, close inspection of the LLM-generated descriptions indicates limited contextual rele-
 137 vance and a high degree of diversity. Previous works [35] reflect the issue of descriptions that are
 138 likely not co-occurring e.g. “steam” and “fried”. We further identify cases where the descriptions are
 139 technically correct but irrelevant to certain images, such as describing “long tail” in frontal images
 140 of cats, underscoring the need for a selective pooling mechanism. Thus, we introduce a vision-
 141 conditional pooling layer to extract instance-specific text features for each attribute for selecting the
 142 most applicable descriptions.

143 Overall, our approach utilizes fine-grained descriptions and organizes them in a Tree of Attribute
 144 following the “concept - attributes - descriptions” structure. Learnable vision expert tokens are
 145 appended to the input image embedding to learn from specific fine-grained attributes such as color
 146 and shape. A vision-conditional pooling layer is further added for each attribute to ensure optimal
 147 image-text alignment. Inspired by CoOP [54], we also incorporate textual contextual tokens in the
 148 text encoder. The overall framework is presented in Fig. 2.

149 3.3 Tree of Attribute generation by LLMs

150 We redefine the process of integrating LLM-generated descriptions by introducing a knowledge graph
 151 $\mathcal{G}_c = \{\mathcal{V}_c, \mathcal{E}_c\}$ for each class c , where \mathcal{V}_c denotes the set of nodes, and \mathcal{E}_c denotes the edges that
 152 capture the semantic relationship between nodes. In previous works, \mathcal{V}_c is the set of descriptions
 153 \mathcal{D}_c , while \mathcal{E}_c is missing. We argue that such methods overlook the inherent structure among the
 154 descriptions and thus do not exploit the richness of these descriptions effectively. To better leverage
 155 knowledge from LLMs, we introduce an attribute layer to link the root node class name, and the leaf
 156 node descriptions. The attribute nodes include visual attributes generated by LLMs, such as color and
 157 shape, for systematically guiding description generation as illustrated in Fig. 1 (c). Each branch of
 158 this “tree” represents a specific attribute, with the subsequent “leaves” fleshing out the descriptions

159 with finer details. In this framework, \mathcal{V}_c includes the class name which is the root node, the set of
 160 attributes such as color and shape being the intermediate layer, and lastly the set of descriptions
 161 under each attribute node. \mathcal{E}_c includes the edges that build up the hierarchy. This structure allows
 162 for a nuanced representation of class information, spanning from general concepts down to specific
 163 attributes and detailed descriptions.

164 To this end, we introduce the Tree of Attribute (ToA), where we use a tree structure to model the
 165 relationship and structure of the descriptions. Let \mathcal{A}_c denote the set of attributes, and for each attribute
 166 $a_c \in \mathcal{A}_c$, we denote its leaf nodes as \mathcal{D}_c^a . Each set \mathcal{D}_c^a contains descriptions that specifically pertain
 167 to attribute a for class c , which is denoted as

$$\mathcal{D}_c^a = \{d_c^{a,1}, d_c^{a,2}, \dots, d_c^{a,n}\}, \quad (3)$$

168 where $d_c^{a,i}$ represents the i -th description for attribute a of class c and n is the number of descriptions
 169 per attribute.

170 The process of generating a Tree of Attribute (ToA) unfolds in three steps: 1) **Attribute Generation:**
 171 We first query LLMs with the dataset information and ask it to generate a set of attributes \mathcal{A} which are
 172 considered relevant and characteristic of the dataset. 2) **Example Generation:** We then ask LLMs to
 173 generate descriptions for a randomly sampled class in the dataset, using the attributes \mathcal{A} identified
 174 in the previous step. Each description takes the format of “class, which {is/has/etc} {description}”.
 175 Human review is performed to ensure the quality of the example. 3) **Description Generation for**
 176 **All Classes:** Building upon the Q&A template from the previous step, the LLM is then tasked with
 177 generating descriptions for all classes in the dataset.

178 Additionally, we incorporate a “global context” attribute which is aligned with the CLS token in the
 179 vision encoder. The descriptions are the 7 standard templates provided in [33].

180 3.4 Learning TAP with Learnable Expert Tokens

181 To fully exploit the structured Tree of Attribute, we introduce learnable visual expert tokens \mathbf{p}_a^v in the
 182 vision branch to learn from each of the attribute nodes $a \in \mathcal{A}$. Unlike traditional methods that rely
 183 on a single CLS token for alignment, these expert tokens enable focused learning on specific image
 184 attributes, such as color or shape, enhancing the model’s performance and interpretability.

185 We denote the set of introduced visual expert tokens as $\mathcal{P}^v = \{\mathbf{p}_a^v | a \in \mathcal{A}\}$. Akin to the idea
 186 of visual prompt tuning (VPT) [14], we insert \mathcal{P}^v into the input sequence of the vision encoder,
 187 forming the prompted input sequences $\tilde{\mathbf{X}}_{\mathbf{p}} = \{\mathbf{e}_{\text{CLS}}, \mathcal{P}^v, \mathbf{E}_{\text{patch}}\}$, where \mathbf{e}_{CLS} is the input CLS
 188 token, and $\mathbf{E}_{\text{patch}}$ denotes the embedded patch tokens. To further boost the model’s capacity for
 189 nuanced attribute representation, we employ deep prompting by introducing a zero-initialized layer
 190 residual for each prompt token across transformer layers, which provides more explicit attribute
 191 guidance across transformer layers. In parallel, we adopt a set of m learnable context tokens
 192 $\mathcal{P}^t = \{\mathbf{p}_j^t | j \in \{1, 2, \dots, m\}\}$ for the text encoder shared across all descriptions, similar to [54].

193 3.5 Vision-Conditional Pooling

194 To mitigate issues of misalignment and potential misleading information from the broad spectrum of
 195 LLM-generated descriptions, we proposed an adaptive vision-conditional pooling layer, applicable to
 196 each set of attribute descriptions \mathcal{D}_a shared across all classes to dynamically pool the most applicable
 197 descriptions based on the visual content of the image x using its corresponding visual expert token
 198 denoted as $\mathbf{p}_{a,x}^v$. For ease of expression, we will proceed without explicitly mentioning x , though it’s
 199 important to note that both the expert token and the resulting attribute-level embeddings are dependent
 200 on the visual information. Intuitively, VCP uses attention to calculate the similarity between \mathbf{p}_a^v and
 201 all embedded descriptions in attribute \mathcal{D}_a , which are then used as weights for a weighted sum of the
 202 original description embeddings. Formally, for each attribute a and its associated expert token \mathbf{p}_a^v ,
 203 the pooled attribute-level embedding \mathbf{v}_c^a for class c and attribute a is:

$$\begin{aligned} \text{Query} &= W_q \cdot \mathbf{p}_a^v, \\ \text{Key} &= W_k \cdot \text{Emb}(\mathcal{D}_c^a), \\ \text{Attention Score} &= \text{softmax}(\text{Query} \cdot \text{Key}^T), \\ \mathbf{v}_c^a &= \text{Attention Score} \cdot \text{Emb}(\mathcal{D}_c^a), \end{aligned} \quad (4)$$

204 where W_q and W_k are learnable weights $\in \mathbb{R}^{d \times d}$, $\text{Emb}(\cdot)$ denotes the embedding function, and
 205 $\text{softmax}(\cdot)$ is the Softmax function. This layer mirrors cross-attention but omits W_v to maintain the
 206 output within the CLIP V-L space.

207 3.6 Training and Inference

208 **Training objective.** During training, each visual expert token \mathbf{p}_a^v is aligned with its associated
 209 attribute-level embedding \mathbf{v}_c^a , trained with the following contrastive objective:

$$L_{con}(\mathbf{p}_a^v, \mathbf{v}_c^a) = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(\cos(\mathbf{p}_a^v, \mathbf{v}_c^a)/\tau)}{\sum_{c=1}^C \exp(\cos(\mathbf{p}_a^v, \mathbf{v}_c^a)/\tau)}, \quad (5)$$

210 where N represents the number of training samples, and τ is the learned temperature of CLIP. The
 211 total classification loss L_{class} is the average of the contrastive loss from each expert token as well as
 212 the CLS token, defined as:

$$L_{class} = \frac{1}{|\mathcal{A}|} \left(\sum_{a \in \mathcal{A}} L_{con}(\mathbf{p}_a^v, \mathbf{v}_c^a) \right), \quad (6)$$

213 Similar to [16] and [4], we regularize the vision CLS token, text feature, and the prediction logits
 214 from each attribute using the vanilla CLIP model. We denote the regularization loss as L_{reg} , where
 215 the details can be found in Appendix. The overall training objective is $L_{total} = L_{class} + L_{reg}$.

216 **Prediction fusion.** During inference, we integrate the prediction by each attribute expert pair by a
 217 weighted sum, formulated as follows:

$$\tilde{y} = \underset{c}{\operatorname{argmax}} \left(\alpha \cos(\mathbf{f}_{CLS}^v, \mathbf{v}_c^{CLS}) + \frac{1 - \alpha}{|\mathcal{A}| - 1} \sum_{a \in \mathcal{A} \setminus \{CLS\}} \cos(\mathbf{p}_a^v, \mathbf{v}_c^a) \right) \quad (7)$$

218 where α is a hyperparameter that signifies the weight assigned to the global context provided by the
 219 CLS token, balancing its contribution with that of the attribute-specific expert prompts.

220 4 Experiments

221 We extensively evaluate our method in two settings: 1) Base-to-novel class generalization, where the
 222 datasets are equally split into base and novel classes. We train the model on the base classes only and
 223 evaluate on both base and novel classes; and 2) Few-shot classification with 16 shots per class.

224 **Datasets and baselines.** For both base to novel class generalization and few-shot setting, we follow
 225 previous works [54, 53], using 11 image recognition datasets. The datasets span a range of recog-
 226 nition tasks: ImageNet [7] and Caltech101 [11] for generic object recognition; OxfordPets [30],
 227 StanfordCars [18], Flowers102 [27], Food101 [2], and FGVC Aircraft [24] for fine-grained classifica-
 228 tion; SUN397 [46] for scene recognition; UCF101 [39] for action recognition; DTD [6] for texture
 229 classification; and EuroSAT [12] for satellite image analysis. We benchmark against several leading
 230 methods, including CLIP [33], CoOp [54], Co-CoOP [53], ProGrad [55], RPO [19], LoGoPrompt
 231 [38], and the state-of-the-art PromptSRC [16].

232 **Implementation details.** A pre-trained CLIP model with a ViT-B/16 vision backbone is used in all
 233 of our experiments and results are averaged over 3 runs. We use GPT-3.5-turbo [29] for attribute and
 234 description generation. We initialize the text context tokens with the word embedding of a photo
 235 of a. For both settings, we iteratively train the vision and text encoders with 5 epochs for vision
 236 and 1 epoch for text schedule. We set $\alpha = 0.4$, $\mu_1 = 10$, and $\mu_2 = 2.5$ for all datasets. We train
 237 the vision encoder for 50 and 100 epochs, and text encoder for 10 and 20 epochs for base-to-novel
 238 generalization and few-shot experiments, respectively. For DTD, Oxford Flowers, Stanford Cars,
 239 UCF101, and Caltech101 datasets, we use a learning rate of 0.002 for the text encoder and 0.006 for
 240 the vision encoder, with $\mu_3 = 3$. For the remaining 6 datasets, the learning rates for both text and
 241 vision encoders are set as 0.004, with $\mu_3 = 1.5$. We also use a Gaussian Prompt Weighting (GPA)
 242 following [16], with a mean of 45, std of 10 for base-to-novel generalization, and 80, 20 for few-shot
 243 experiments. Refer to the Appendix for additional implementation details.

Table 1: Comparison of TAP in base-to-novel generalization. HM: harmonic mean [45].

(a) Average			(b) ImageNet			(c) Caltech101			(d) OxfordPets		
	Base	Novel HM		Base	Novel HM		Base	Novel HM		Base	Novel HM
CLIP	69.34	74.22 71.70	CLIP	72.43	68.14 70.22	CLIP	96.84	94.00 95.40	CLIP	91.17	97.26 94.12
CoOp	82.69	63.22 71.66	CoOp	76.47	67.88 71.92	CoOp	98.00	89.81 93.73	CoOp	93.67	95.29 94.47
Co-CoOp	80.47	71.69 75.83	Co-CoOp	75.98	70.43 73.10	Co-CoOp	97.96	93.81 95.84	Co-CoOp	95.20	97.69 96.43
ProGrad	82.48	70.75 76.16	ProGrad	77.02	66.66 71.46	ProGrad	98.02	93.89 95.91	ProGrad	95.07	97.63 96.33
RPO	81.13	75.00 77.78	RPO	76.60	71.57 74.00	RPO	97.97	94.37 96.03	RPO	94.63	97.50 96.05
LoGoPrompt	84.47	74.24 79.03	LoGoPrompt	76.74	70.83 73.66	LoGoPrompt	98.19	93.78 95.93	LoGoPrompt	96.07	96.31 96.18
PromptSRC	84.26	76.10 79.97	PromptSRC	77.60	70.73 74.01	PromptSRC	98.10	94.03 96.02	PromptSRC	95.33	97.30 96.30
TAP	84.75	77.63 81.04	TAP	77.97	70.40 73.99	TAP	98.90	95.50 97.17	TAP	95.80	97.73 96.76
(e) StanfordCars			(f) Flowers102			(g) Food101			(h) FGVC Aircraft		
	Base	Novel HM		Base	Novel HM		Base	Novel HM		Base	Novel HM
CLIP	63.37	74.89 68.65	CLIP	72.08	77.80 74.83	CLIP	90.10	91.22 90.66	CLIP	27.19	36.29 31.09
CoOp	78.12	60.40 68.13	CoOp	97.60	59.67 74.06	CoOp	88.33	82.26 85.19	CoOp	40.44	22.30 28.75
Co-CoOp	70.49	73.59 72.01	Co-CoOp	94.87	71.75 81.71	Co-CoOp	90.70	91.29 90.99	Co-CoOp	33.41	23.71 27.74
ProGrad	77.68	68.63 72.88	ProGrad	95.54	71.87 82.03	ProGrad	90.37	89.59 89.98	ProGrad	40.54	27.57 32.82
RPO	73.87	75.53 74.69	RPO	94.13	76.67 84.50	RPO	90.33	90.83 90.58	RPO	37.33	34.20 35.70
LoGoPrompt	78.36	72.39 75.26	LoGoPrompt	99.05	76.52 86.34	LoGoPrompt	90.82	91.41 91.11	LoGoPrompt	45.98	34.67 39.53
PromptSRC	78.27	74.97 76.58	PromptSRC	98.07	76.50 85.95	PromptSRC	90.67	91.53 91.10	PromptSRC	42.73	37.87 40.15
TAP	80.70	74.27 77.35	TAP	97.90	75.57 85.30	TAP	90.97	91.83 91.40	TAP	44.40	36.50 40.06
(i) SUN397			(j) DTD			(k) EuroSAT			(l) UCF101		
	Base	Novel HM		Base	Novel HM		Base	Novel HM		Base	Novel HM
CLIP	69.36	75.35 72.23	CLIP	53.24	59.90 56.37	CLIP	56.48	64.05 60.03	CLIP	70.53	77.50 73.85
CoOp	80.60	65.89 72.51	CoOp	79.44	41.18 54.24	CoOp	92.19	54.74 68.69	CoOp	84.69	56.05 67.46
Co-CoOp	79.74	76.86 78.27	Co-CoOp	77.01	56.00 64.85	Co-CoOp	87.49	60.04 71.21	Co-CoOp	82.33	73.45 77.64
ProGrad	81.26	74.17 77.55	ProGrad	77.35	52.35 62.45	ProGrad	90.11	60.89 72.67	ProGrad	84.33	74.94 79.35
RPO	80.60	77.80 79.18	RPO	76.70	62.13 68.61	RPO	86.63	68.97 76.79	RPO	83.67	75.43 79.34
LoGoPrompt	81.20	78.12 79.63	LoGoPrompt	82.87	60.14 69.70	LoGoPrompt	93.67	69.44 79.75	LoGoPrompt	86.19	73.07 79.09
PromptSRC	82.67	78.47 80.52	PromptSRC	83.37	62.97 71.75	PromptSRC	92.90	73.90 82.32	PromptSRC	87.10	78.80 82.74
TAP	82.87	79.53 81.17	TAP	84.20	68.00 75.24	TAP	90.70	82.17 86.22	TAP	87.90	82.43 85.08

Table 2: Few shot classification results with 16 shots.

16-Shot Classification												
	Average	ImageNet	Caltech101	Pets	Cars	Flowers	Food101	Aircraft	SUN397	DTD	EuroSAT	UCF101
CLIP	78.79	67.31	95.43	85.34	80.44	97.37	82.90	45.36	73.28	69.96	87.21	82.11
CoOp	79.89	71.87	95.57	91.87	83.07	97.07	84.20	43.40	74.67	69.87	84.93	82.23
CoCoOp	74.90	70.83	95.16	93.34	71.57	87.84	87.25	31.21	72.15	63.04	73.32	78.14
MaPLe	81.79	72.33	96.00	92.83	83.57	97.00	85.33	48.40	75.53	71.33	92.33	85.03
PSRC	82.87	73.17	96.07	93.67	83.83	97.60	87.50	50.83	77.23	72.73	92.43	86.47
TAP	83.37	73.76	96.73	93.90	85.37	98.10	87.53	50.43	77.30	74.90	91.90	87.17

244 **4.1 Base-to-Novel Generalization**

245 In base-to-novel generalization, we equally split the classes into base and novel classes. Initial
 246 training and evaluations are conducted on the seen base classes, followed by evaluation on the unseen
 247 novel classes in a zero-shot manner. TAP surpasses prior state-of-the-art models in terms of the
 248 base and novel class accuracy, as well as their harmonic mean across most of the 11 datasets, with
 249 an average increase of 1.53% in the zero-shot novel class prediction, and a 1.07% increase in the
 250 overall harmonic mean in average, as detailed in Table 1. Notably, our method improves unseen class
 251 prediction without compromising base class performance, exhibiting an average performance boost
 252 of 0.49%. In the challenging fine-grained tasks such as DTD, EuroSAT, and UCF101, TAP achieves
 253 significant improvements in novel class prediction by 5.03%, 8.27%, and 3.63% respectively. These
 254 results underscore the robust generalizability and efficacy of our method across diverse scenarios.

255 **4.2 Few-Shot Classification**

256 In few-shot classification, TAP also outperforms existing methods in 9 out of the 11 datasets. Detailed
 257 in Table 2, we achieve an average accuracy of 83.37 across the 11 datasets, surpassing the previous
 258 state-of-the-art methods by 0.5%, further demonstrating the effectiveness of our method.

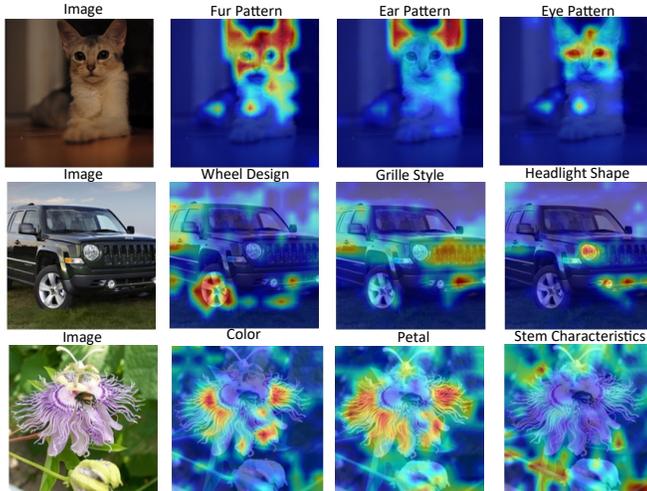


Figure 3: Visualization of the class activation maps.

Table 3: Effects of the Tree of Attributes.

Des. Org.	Unstructured	Ours
Base	82.89	84.75
Novel	75.32	77.63
HM	78.93	81.04

Table 4: Effects of domain experts.

Align. Token	CLS	Ours
Base	83.89	84.75
Novel	76.85	77.63
HM	80.22	81.04

Table 5: Effects of the number of experts.

Attrs. Num.	1	2	3	4	5	6	7	8	Ours
Base Acc.	83.20	83.97	84.1	84.41	84.45	84.62	84.66	84.74	84.75
Novel Acc.	74.90	76.20	76.35	77.06	77.13	77.17	77.35	76.67	77.63
HM	78.83	79.90	80.04	80.57	80.63	80.72	80.84	80.50	81.04

259 4.3 Ablation Study

260 **Effects of Tree of Attribute.** A core inquiry is whether structuring descriptions into a Tree of
 261 Attribute (ToA) offers advantages over an unstructured aggregation of LLM-generated descriptions.
 262 To evaluate, we revert to aligning a mixed, unstructured set of descriptions with the CLS token
 263 - a common practice in prior studies [25, 19, 40, 52], while keeping the same number of visual
 264 prompt tokens. According to Table 3, substituting the ToA with an unstructured set results in
 265 significant performance decreases of 1.86%, 2.31%, and 2.11% across the average base, novel, and
 266 their harmonic mean performances, respectively. This stark contrast underscores the ToA’s critical
 267 role in enhancing model efficacy.

268 **Effects of Learning through Domain Experts.** Further, we examine the impact of substituting the
 269 CLS token with visual expert tokens for learning fine-grained attributes, commonly adopted in in
 270 previous works [25, 19, 40, 52]. Our findings (Table 4) reveal improvements of 0.89%, 0.78%, and
 271 0.82% in the average base, novel, and harmonic mean accuracies, respectively, upon integrating visual
 272 expert tokens. These results support the notion that domain-specific, learnable tokens enhance the
 273 model’s ability to grasp fine-grained details by focusing on distinct aspects of the image, as opposed
 274 to the CLS token’s global focus.

275 **Effects of Number of Attributes.** In our framework, the selection of attributes is dynamically
 276 determined by LLMs, leading to variability across different datasets. This adaptability stands in
 277 contrast to a static approach where the number of attributes is uniformly set across all datasets. To
 278 understand the impact of this variability, we explore how altering the number of attributes from 1 to 8
 279 influences model performance. Our findings, detailed in Table 5, reveal a performance improvement
 280 trend as the number of attributes increases, with an optimal peak at 7 attributes before a slight decline
 281 at 8. However, crucially, across all fixed-attribute scenarios, none matched the performance achieved
 282 through our method’s dynamic attribute determination. These results underscore the importance of
 283 an adaptive approach to attribute selection, as opposed to a one-size-fits-all strategy.

284 **Design choice of the vision-conditional pooling layer.** Lastly, we ablate the design of the pooling
 285 layer, starting from the naive training-free average pooling, to the attention-based pooling mechanism
 286 with condition on the input image. Compared to average pooling, VCP demonstrates a performance
 287 gain of 1.08% in the average harmonic mean. Furthermore, when compared with attention-based max
 288 pooling, which selects a single description per attribute according to the attention score in Eq. (4),

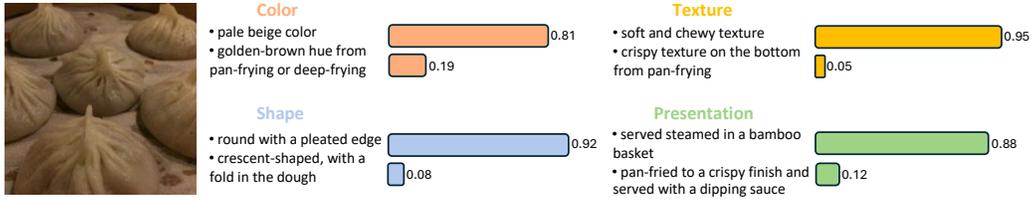


Figure 4: Visualization of the attention weights in the VCP layer for an example “dumplings” image.

Table 6: Design choice of the pooling layer.

Pooling Method	Base Acc.	Novel Acc.	HM
Attn. Max Pooling	82.90	76.36	79.49
Average Pooling	83.18	76.98	79.96
VCP (Ours)	84.75	77.63	81.04

VCP maintains a superior advantage of 1.55% in average harmonic mean. These outcomes attest to the VCP layer’s integral role in finetuning attribute relevance to the visual context, substantiating its design and implementation within our model.

4.4 Visualization

Expert tokens focus on attribute-related regions. We further investigate the effects of vision domain experts by visualizing their class activation maps from three illustrative examples using GradCAM [37], as shown in Fig. 3. These visualizations underscore the precision with which each expert token concentrates on the image regions pertinent to its designated attribute. Take the first cat image as an example. The “fur pattern” expert distinctly highlights the animal’s fur texture, whereas the “ear” and “eye” experts focus precisely on the respective anatomical features. This pattern of attribute-specific attention is consistent across the evaluated examples, reinforcing the conceptualization of expert tokens as dedicated “domain experts” within the visual field.

VCP layer pools the most applicable descriptions. The inherently interpretable nature of the VCP layer, thanks to its attention mechanism, allows for insightful visualizations of its operational process. Through the examination of attention weights assigned by the VCP layer to different attributes in a given image, we elucidate the layer’s capability to discern and prioritize the most applicable descriptions. As illustrated in Fig. 4 with a “dumplings” image, the VCP layer adeptly allocates higher attention weights to descriptions accurately reflecting the observed instance (e.g., assigning weights of 0.92 to “round with a pleated edge” under the “Shape” attribute and 0.95 to “soft and chewy texture” under the Texture”). In contrast, less relevant descriptions for the specific image context (e.g., “crenant-shaped” for Shape and “crispy texture from pan-frying” for Texture) receive significantly lower weights. This discernment is crucial, given the class “dumplings” encompasses a broad variety of appearances based on cooking methods, yet not all descriptions are fitting for every instance. These visualizations compellingly demonstrate the VCP layer’s effectiveness in refining description relevance, thereby enhancing the model’s interpretative alignment with the visual content.

5 Conclusion

This paper introduces Tree of Attribute Prompt learning (TAP), a novel method that integrates detailed, LLM-generated descriptions within VLMs, achieving state-of-the-art performance in both base-to-novel generalization and few-shot image classification tasks across 11 diverse datasets. TAP leverages a hierarchical “Tree of Attribute” framework, distilling structured knowledge graphs from LLMs for nuanced representation of visual concepts, and employs learnable “domain expert” tokens and a vision-conditional pooling module for optimal image-text alignment. While promising, we note that the reliance on LLMs presents challenges in fine-grained datasets where similar classes require nuanced differentiation, in which cases LLMs generate identical descriptions for distinct classes, impacting novel class prediction performance. It highlights the current limitations of LLMs in discerning highly fine-grained distinctions. Addressing this challenge through enhanced LLM capabilities or alternative strategies will be a key focus of future research.

References

- 326
- 327 [1] Hyojin Bahng, Ali Jahanian, Swami Sankaranarayanan, and Phillip Isola. Visual prompting: Modifying
328 pixel space to adapt pre-trained models. *arXiv preprint arXiv:2203.17274*, 3:11–12, 2022.
- 329 [2] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative components
330 with random forests. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland,
331 September 6–12, 2014, Proceedings, Part VI 13*, pages 446–461. Springer, 2014.
- 332 [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
333 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.
334 *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 335 [4] Adrian Bulat and Georgios Tzimiropoulos. Lasp: Text-to-text optimization for language-aware soft
336 prompting of vision & language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision
337 and Pattern Recognition*, pages 23232–23241, 2023.
- 338 [5] Guangyi Chen, Weiran Yao, Xiangchen Song, Xinyue Li, Yongming Rao, and Kun Zhang. Prompt learning
339 with optimal transport for vision-language models. In *ICLR*, 2023.
- 340 [6] Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing
341 textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
342 pages 3606–3613, 2014.
- 343 [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical
344 image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255.
345 Ieee, 2009.
- 346 [8] Mohammad Mahdi Derakhshani, Enrique Sanchez, Adrian Bulat, Victor G Turrissi da Costa, Cees GM
347 Snoek, Georgios Tzimiropoulos, and Brais Martinez. Bayesian prompt learning for image-language model
348 generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages
349 15237–15246, 2023.
- 350 [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
351 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,
352 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In
353 *International Conference on Learning Representations*, 2021. URL [https://openreview.net/forum?
354 id=YicbFdNTTy](https://openreview.net/forum?id=YicbFdNTTy).
- 355 [10] Zalan Fabian, Zhongqi Miao, Chunyuan Li, Yuanhan Zhang, Ziwei Liu, Andrés Hernández, Andrés
356 Montes-Rojas, Rafael Escucha, Laura Siabatto, Andrés Link, et al. Multimodal foundation models for
357 zero-shot animal species recognition in camera trap images. *arXiv preprint arXiv:2311.01064*, 2023.
- 358 [11] Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training examples:
359 An incremental bayesian approach tested on 101 object categories. In *2004 conference on computer vision
360 and pattern recognition workshop*, pages 178–178. IEEE, 2004.
- 361 [12] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep
362 learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied
363 Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.
- 364 [13] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung,
365 Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text
366 supervision. In *International conference on machine learning*, pages 4904–4916. PMLR, 2021.
- 367 [14] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and
368 Ser-Nam Lim. Visual prompt tuning. In *European Conference on Computer Vision*, pages 709–727.
369 Springer, 2022.
- 370 [15] Muhammad Uzair Khattak, Hanoona Rasheed, Muhammad Maaz, Salman Khan, and Fahad Shahbaz Khan.
371 Maple: Multi-modal prompt learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision
372 and Pattern Recognition*, pages 19113–19122, 2023.
- 373 [16] Muhammad Uzair Khattak, Syed Talal Wasim, Muzammal Naseer, Salman Khan, Ming-Hsuan Yang, and
374 Fahad Shahbaz Khan. Self-regulating prompts: Foundational model adaptation without forgetting. In
375 *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15190–15200,
376 October 2023.

- 377 [17] Jae Myung Kim, A Koepke, Cordelia Schmid, and Zeynep Akata. Exposing and mitigating spurious
378 correlations for cross-modal retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision
379 and Pattern Recognition*, pages 2584–2594, 2023.
- 380 [18] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained
381 categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pages
382 554–561, 2013.
- 383 [19] Dongjun Lee, Seokwon Song, Jihee Suh, Joonmyeong Choi, Sanghyeok Lee, and Hyunwoo J Kim.
384 Read-only prompt optimization for vision-language few-shot learning. In *Proceedings of the IEEE/CVF
385 International Conference on Computer Vision*, pages 1401–1411, 2023.
- 386 [20] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning,
387 2021.
- 388 [21] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation, 2021.
- 389 [22] Xiao Liu, Kaixuan Ji, Yicheng Fu, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-tuning v2: Prompt tuning
390 can be comparable to fine-tuning universally across scales and tasks. *CoRR*, abs/2110.07602, 2021. URL
391 <https://arxiv.org/abs/2110.07602>.
- 392 [23] Yuning Lu, Jianzhuang Liu, Yonggang Zhang, Yajing Liu, and Xinmei Tian. Prompt distribution learning.
393 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5206–
394 5215, 2022.
- 395 [24] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual
396 classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- 397 [25] Chengzhi Mao, Revant Teotia, Amrutha Sundar, Sachit Menon, Junfeng Yang, Xin Wang, and Carl
398 Vondrick. Doubly right object recognition: A why prompt for visual rationales. In *Proceedings of the
399 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2722–2732, 2023.
- 400 [26] Sachit Menon and Carl Vondrick. Visual classification via description from large language models. *ICLR*,
401 2023.
- 402 [27] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of
403 classes. In *2008 Sixth Indian conference on computer vision, graphics & image processing*, pages 722–729.
404 IEEE, 2008.
- 405 [28] Zachary Novack, Julian McAuley, Zachary Lipton, and Saurabh Garg. Chils: Zero-shot image classification
406 with hierarchical label sets. In *International Conference on Machine Learning (ICML)*, 2023.
- 407 [29] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang,
408 Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with
409 human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- 410 [30] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012 IEEE
411 conference on computer vision and pattern recognition*, pages 3498–3505. IEEE, 2012.
- 412 [31] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and
413 A. Lerer. Automatic differentiation in PyTorch. In *NeurIPS Autodiff Workshop*, 2017.
- 414 [32] Sarah Pratt, Ian Covert, Rosanne Liu, and Ali Farhadi. What does a platypus look like? generating
415 customized prompts for zero-shot image classification. In *Proceedings of the IEEE/CVF International
416 Conference on Computer Vision*, pages 15691–15701, 2023.
- 417 [33] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish
418 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from
419 natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR,
420 2021.
- 421 [34] Hanoona Rasheed, Muhammad Uzair Khattak, Muhammad Maaz, Salman Khan, and Fahad Shahbaz
422 Khan. Fine-tuned clip models are efficient video learners. In *Proceedings of the IEEE/CVF Conference on
423 Computer Vision and Pattern Recognition*, pages 6545–6554, 2023.
- 424 [35] Karsten Roth, Jae Myung Kim, A. Sophia Koepke, Oriol Vinyals, Cordelia Schmid, and Zeynep Akata.
425 Waffling around for performance: Visual classification with random words and broad concepts. In
426 *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15746–15757,
427 October 2023.

- 428 [36] Shuvendu Roy and Ali Etemad. Consistency-guided prompt learning for vision-language models. In *The*
429 *Twelfth International Conference on Learning Representations*, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=wsRXwlwx4w)
430 [forum?id=wsRXwlwx4w](https://openreview.net/forum?id=wsRXwlwx4w).
- 431 [37] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and
432 Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In
433 *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- 434 [38] Cheng Shi and Sibeil Yang. Logoprompt: Synthetic text images can be good visual prompts for vision-
435 language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages
436 2932–2941, 2023.
- 437 [39] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions
438 classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- 439 [40] Xinyu Tian, Shu Zou, Zhaoyuan Yang, and Jing Zhang. Argue: Attribute-guided prompt tuning for
440 vision-language models. *arXiv preprint arXiv:2311.16494*, 2023.
- 441 [41] Dongsheng Wang, Miaoge Li, Xinyang Liu, MingSheng Xu, Bo Chen, and Hanwang Zhang. Tuning
442 multi-mode token-level prompt alignment across modalities. In *Thirty-seventh Conference on Neural*
443 *Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=A253n2EXCd>.
- 444 [42] Wenhao Wang, Yifan Sun, Wei Li, and Yi Yang. TransHP: Image classification with hierarchical
445 prompting. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL
446 <https://openreview.net/forum?id=vpQuCsZXz2>.
- 447 [43] Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren, Guolong
448 Su, Vincent Perot, Jennifer Dy, et al. Dualprompt: Complementary prompting for rehearsal-free continual
449 learning. In *European Conference on Computer Vision*, pages 631–648. Springer, 2022.
- 450 [44] Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent
451 Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *Proceedings of the*
452 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 139–149, 2022.
- 453 [45] Yongqin Xian, Bernt Schiele, and Zeynep Akata. Zero-shot learning—the good, the bad and the ugly. In
454 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4582–4591, 2017.
- 455 [46] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-
456 scale scene recognition from abbey to zoo. In *2010 IEEE computer society conference on computer vision*
457 *and pattern recognition*, pages 3485–3492. IEEE, 2010.
- 458 [47] Yinghui Xing, Qirui Wu, De Cheng, Shizhou Zhang, Guoqiang Liang, Peng Wang, and Yanning Zhang.
459 Dual modality prompt tuning for vision-language pre-trained model. *IEEE Transactions on Multimedia*,
460 pages 1–13, 2023. doi: 10.1109/TMM.2023.3291588.
- 461 [48] An Yan, Yu Wang, Yiwu Zhong, Chengyu Dong, Zexue He, Yujie Lu, William Yang Wang, Jingbo Shang,
462 and Julian McAuley. Learning concise and descriptive attributes for visual recognition. In *Proceedings of*
463 *the IEEE/CVF International Conference on Computer Vision*, pages 3090–3100, 2023.
- 464 [49] Yue Yang, Artemis Panagopoulou, Shenghao Zhou, Daniel Jin, Chris Callison-Burch, and Mark Yatskar.
465 Language in a bottle: Language model guided concept bottlenecks for interpretable image classification.
466 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19187–
467 19197, 2023.
- 468 [50] Yi Zhang, Ce Zhang, Ke Yu, Yushun Tang, and Zhihai He. Concept-guided prompt learning for gener-
469 alization in vision-language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38
470 (7):7377–7386, Mar. 2024. doi: 10.1609/aaai.v38i7.28568. URL [https://ojs.aaai.org/index.php/](https://ojs.aaai.org/index.php/AAAI/article/view/28568)
471 [AAAI/article/view/28568](https://ojs.aaai.org/index.php/AAAI/article/view/28568).
- 472 [51] Yuanhan Zhang, Kaiyang Zhou, and Ziwei Liu. Neural prompt search. *arXiv preprint arXiv:2206.04673*,
473 2022.
- 474 [52] Zhaoheng Zheng, Jingmin Wei, Xuefeng Hu, Haidong Zhu, and Ram Nevatia. Large language models are
475 good prompt learners for low-shot image classification. *arXiv preprint arXiv:2312.04076*, 2023.
- 476 [53] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for
477 vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
478 *Recognition*, pages 16816–16825, 2022.

479 [54] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language
 480 models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022.

481 [55] Beier Zhu, Yulei Niu, Yucheng Han, Yue Wu, and Hanwang Zhang. Prompt-aligned gradient for prompt
 482 tuning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15659–
 483 15669, 2023.

484 A Appendix

485 A.1 Model regularization

486 Denote the frozen image feature from CLIP vision encoder as \mathbf{f}^v , the frozen text feature for description d from
 487 CLIP text encoder as \mathbf{f}_d^t , and the zero-shot logit prediction from CLIP as \hat{y} . Additionally, denote the trained
 488 image feature as $\tilde{\mathbf{f}}^v$, the trained text feature for description d as $\tilde{\mathbf{f}}_d^t$, and the logit prediction from attribute a after
 489 training as \tilde{y}_a . The losses are as follows:

$$L_{L_1-V} = \|\mathbf{f}^v - \tilde{\mathbf{f}}^v\|_1 \quad (8)$$

$$L_{con-T} = - \sum_{d \in \mathcal{D}} \left(\frac{1}{2} \log \frac{\exp(\cos(\mathbf{f}_d^t, \tilde{\mathbf{f}}_d^t))}{\sum_{k \in \mathcal{D}_s} \exp(\cos(\mathbf{f}_d^t, \tilde{\mathbf{f}}_k^t))} + \frac{1}{2} \log \frac{\exp(\cos(\tilde{\mathbf{f}}_d^t, \mathbf{f}_d^t))}{\sum_{k \in \mathcal{D}_s} \exp(\cos(\tilde{\mathbf{f}}_k^t, \mathbf{f}_d^t))} \right) \quad (9)$$

490

$$L_{KL-attr} = \frac{1}{|\mathcal{A}|} \left(\sum_{a \in \mathcal{A}} \mathcal{D}_{KL}(\hat{y}, \tilde{y}_a) \right) \quad (10)$$

491 The regularization loss is then:

$$L_{reg} = \mu_1 L_{L_1-V} + \mu_2 L_{KL-attr} + \mu_3 L_{con-T}, \quad (11)$$

492 Our overall training objective is thus given by:

$$L_{total} = L_{class} + L_{reg} \quad (12)$$

493 A.2 Additional implementation details

494 We use PyTorch [31] to implement all experiments on a single NVIDIA A100-80GB GPU. Our code is developed
 495 based on the implementation of CoOp [54], which is available at <https://github.com/KaiyangZhou/CoOp> and
 496 released under the MIT license. Our code is also released under the MIT license. Baseline results for both
 497 base-to-novel generalization and few-shot classification are taken from their respective publications. For the
 498 “global context” attribute which is aligned with the CLS token in the vision encoder, we use the following 7
 499 selected templates provided in [33].

500 "itap of a {class}."
 501 "a bad photo of the {class}."
 502 "a origami {class}."
 503 "a photo of the large {class}."
 504 "a {class} in a video game."
 505 "art of the {class}."
 506 "a photo of the small {class}."

507 A.3 Prompts for Tree-of-Attribute generation

508 As introduced in Section 3.3, we generate the Tree-of-Attribute with the following three steps: 1) Attribute
 509 Generation, 2) In-Context Example Generation, and 3) Description Generation for All Classes. The prompts for
 510 each step are as follows:

511 1) Attribute Generation:

512 {Dataset Description.}

513 *Visual attributes refer to observable, describable features of the images that can include color, shape, size,*
 514 *texture, and any specific patterns or markings, which can help differentiate between classes for the dataset. They*

515 *should be consistently observable across multiple images of the same class. Your task is to generate a list of*
516 *visual attributes (less than 10) for the {Dataset Name} dataset. Ensure this list is clear, concise, and specific to*
517 *the dataset's needs. Avoid generic attributes that do not contribute to distinguishing between classes.*

518 **2) In-Context Example Generation**

519 *Describe describe what a "{Random Class Name}" class in the {Dataset Name} dataset look like using the*
520 *generated visual attributes.*

521 *You must follow the following rules:*

522 *1. For each visual attribute, describe all possible variations as separate sentences. This approach allows for a*
523 *detailed and clear presentation of each attribute's range.*

524 *2. Provide a maximum of five descriptions for each visual attribute to maintain focus and relevance. Also, aim to*
525 *provide at least two descriptions to ensure a comprehensive overview of the attribute.*

526 *3. The descriptions should provide clear, distinguishable features of each class to support image classification*
527 *tasks.*

528 *4. Descriptions for each attribute are independent from each other, and they should not serve as context for each*
529 *other.*

530 *5. Each description describes an image independently. If certain description is possible for a class, please just*
531 *list that description, and do not use words like "may have" or "sometimes have".*

532 *6. Reply descriptions only. Do not include any explanation before and after the description.*

533 *7. The descriptions should follow the format of "classname, which ...", where "..." is the description of the visual*
534 *attribute.*

535 **3) Description Generation for All Classes**

536 *{Dataset Description.}*

537 *Your task is to write detailed descriptions for various classes within the {Dataset Name} dataset, using the*
538 *provided visual attributes such as color and shape. These descriptions will help in accurately classifying and*
539 *understanding the unique features of each class.*

540 *You must follow the following rules:*

541 *1. For each visual attribute, describe all possible variations as separate sentences. This approach allows for a*
542 *detailed and clear presentation of each attribute's range.*

543 *2. Provide a maximum of five descriptions for each visual attribute to maintain focus and relevance. Also, aim to*
544 *provide at least two descriptions to ensure a comprehensive overview of the attribute.*

545 *3. The descriptions should provide clear, distinguishable features of each class to support image classification*
546 *tasks.*

547 *4. Descriptions for each attribute are independent from each other, and they should not serve as context for each*
548 *other.*

549 *5. Each description describes an image independently. If certain description is possible for a class, please just*
550 *list that description, and do not use words like "may have" or "sometimes have".*

551 *6. Reply descriptions only. Do not include any explanation before and after the description.*

552 *7. The descriptions should follow the format of "classname, which ...", where "..." is the description of the visual*
553 *attribute.*

554 *Q: Describe what a "{Random Class Name}" in the {Dataset Name} look like using the following visual attributes:*
555 *{Visual Attributes from Step 1.}*

556 *A: {Answer from Step 2.}*

557 *Q: Describe what a "{Target Class Name}" in the {Dataset Name} look like using the following visual attributes:*
558 *{Visual Attributes from Step 1.}*

559 *A:*

560 *In the prompt templates, "Dataset Description" is the description of the dataset from their official website,*
561 *"Random Class Name" is a randomly sampled class name in the dataset for in-context example generation, and*
562 *"Target Class Name" is the class name of interest for the current query. While step 1 and 2 are made in two*
563 *consecutive calls to provide contexts which are queried once per dataset, step 3 is queried independently for*

564 each of the remaining classes in the dataset. Human review is performed after step 2 to ensure a high-quality set
565 of attributes and in-context example.

566 **A.4 Potential societal impacts**

567 While our work primarily focuses on advancing prompt learning in vision-language models, it's crucial to
568 acknowledge the potential broader societal implications of such advancements. On the positive side, TAP could
569 lead to more efficient and accurate image understanding systems, benefiting various domains. For instance, it
570 could enhance accessibility for visually impaired individuals by providing more detailed descriptions of visual
571 content. Furthermore, improved visual understanding could contribute to more effective content moderation,
572 mitigating the spread of harmful online materials. However, these advancements also present potential risks.
573 LLMs used for description generation can perpetuate existing societal biases present in their training data, leading
574 to biased outcomes in image recognition. Moreover, sophisticated VLMs could be misused to create misleading
575 visual content, contributing to misinformation and manipulation. The enhanced ability to analyze and understand
576 images also raises privacy concerns, particularly in surveillance contexts where personal information could be
577 extracted from visual data. Addressing these potential negative impacts necessitates careful consideration of bias
578 mitigation techniques during LLM training, promoting transparency and explainability in VLM decision-making,
579 and establishing ethical guidelines for responsible development and deployment of such technologies.

580 **NeurIPS Paper Checklist**

581 **1. Claims**

582 Question: Do the main claims made in the abstract and introduction accurately reflect the paper's
583 contributions and scope?

584 Answer: [Yes]

585 Justification: The abstract and introduction clearly state the problem of limited context in existing
586 prompt learning methods, propose TAP as a solution using structured knowledge graphs and domain
587 experts, and highlight the strong experimental results in both base-to-novel generalization and few-shot
588 classification. This accurately reflects the paper's contributions and scope.

589 Guidelines:

- 590 • The answer NA means that the abstract and introduction do not include the claims made in the
591 paper.
- 592 • The abstract and/or introduction should clearly state the claims made, including the contributions
593 made in the paper and important assumptions and limitations. A No or NA answer to this
594 question will not be perceived well by the reviewers.
- 595 • The claims made should match theoretical and experimental results, and reflect how much the
596 results can be expected to generalize to other settings.
- 597 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are not
598 attained by the paper.

599 **2. Limitations**

600 Question: Does the paper discuss the limitations of the work performed by the authors?

601 Answer: [Yes]

602 Justification: The paper includes a discussion of the limitations associated with relying on LLMs for
603 generating descriptions, particularly in fine-grained datasets where similar classes require nuanced
604 differentiation. This discussion can be found in "Conclusion".

605 Guidelines:

- 606 • The answer NA means that the paper has no limitation while the answer No means that the paper
607 has limitations, but those are not discussed in the paper.
- 608 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 609 • The paper should point out any strong assumptions and how robust the results are to violations of
610 these assumptions (e.g., independence assumptions, noiseless settings, model well-specification,
611 asymptotic approximations only holding locally). The authors should reflect on how these
612 assumptions might be violated in practice and what the implications would be.
- 613 • The authors should reflect on the scope of the claims made, e.g., if the approach was only tested
614 on a few datasets or with a few runs. In general, empirical results often depend on implicit
615 assumptions, which should be articulated.
- 616 • The authors should reflect on the factors that influence the performance of the approach. For
617 example, a facial recognition algorithm may perform poorly when image resolution is low or
618 images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide
619 closed captions for online lectures because it fails to handle technical jargon.
- 620 • The authors should discuss the computational efficiency of the proposed algorithms and how
621 they scale with dataset size.
- 622 • If applicable, the authors should discuss possible limitations of their approach to address problems
623 of privacy and fairness.
- 624 • While the authors might fear that complete honesty about limitations might be used by reviewers
625 as grounds for rejection, a worse outcome might be that reviewers discover limitations that
626 aren't acknowledged in the paper. The authors should use their best judgment and recognize
627 that individual actions in favor of transparency play an important role in developing norms that
628 preserve the integrity of the community. Reviewers will be specifically instructed to not penalize
629 honesty concerning limitations.

630 **3. Theory Assumptions and Proofs**

631 Question: For each theoretical result, does the paper provide the full set of assumptions and a complete
632 (and correct) proof?

633 Answer: [NA]

634 Justification: The paper focuses on proposing a novel method for prompt learning in VLMs and
635 evaluating its empirical performance. It doesn't introduce any new theoretical results or theorems
636 requiring formal proofs.

637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. **Experimental Result Reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The paper provides all the necessary information for reproducing the experimental results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. **Open access to data and code**

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: Our codebase is built based on the CoOP and CoCoOP [54, 53], and can be reproduced based on our Methods, Implementation details in main text and appendix. Our code will be released upon acceptance.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.

- 695 • While we encourage the release of code and data, we understand that this might not be possible,
696 so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless
697 this is central to the contribution (e.g., for a new open-source benchmark).
- 698 • The instructions should contain the exact command and environment needed to run to reproduce
699 the results. See the NeurIPS code and data submission guidelines ([https://nips.cc/public/
700 guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 702 • The authors should provide instructions on data access and preparation, including how to access
the raw data, preprocessed data, intermediate data, and generated data, etc.
- 703 • The authors should provide scripts to reproduce all experimental results for the new proposed
704 method and baselines. If only a subset of experiments are reproducible, they should state which
705 ones are omitted from the script and why.
- 706 • At submission time, to preserve anonymity, the authors should release anonymized versions (if
707 applicable).
- 708 • Providing as much information as possible in supplemental material (appended to the paper) is
709 recommended, but including URLs to data and code is permitted.

710 6. Experimental Setting/Details

711 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters,
712 how they were chosen, type of optimizer, etc.) necessary to understand the results?

713 Answer: [Yes]

714 Justification: The training and test details can be found in section Experiments and Appendix.

715 Guidelines:

- 716 • The answer NA means that the paper does not include experiments.
- 717 • The experimental setting should be presented in the core of the paper to a level of detail that is
718 necessary to appreciate the results and make sense of them.
- 719 • The full details can be provided either with the code, in appendix, or as supplemental material.

720 7. Experiment Statistical Significance

721 Question: Does the paper report error bars suitably and correctly defined or other appropriate informa-
722 tion about the statistical significance of the experiments?

723 Answer: [No]

724 Justification: We follow previous works [54, 53] to report results averaged over 3 runs. Error bars are
725 not reported.

726 Guidelines:

- 727 • The answer NA means that the paper does not include experiments.
- 728 • The authors should answer "Yes" if the results are accompanied by error bars, confidence
729 intervals, or statistical significance tests, at least for the experiments that support the main claims
730 of the paper.
- 731 • The factors of variability that the error bars are capturing should be clearly stated (for example,
732 train/test split, initialization, random drawing of some parameter, or overall run with given
733 experimental conditions).
- 734 • The method for calculating the error bars should be explained (closed form formula, call to a
735 library function, bootstrap, etc.)
- 736 • The assumptions made should be given (e.g., Normally distributed errors).
- 737 • It should be clear whether the error bar is the standard deviation or the standard error of the
738 mean.
- 739 • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report
740 a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is
741 not verified.
- 742 • For asymmetric distributions, the authors should be careful not to show in tables or figures
743 symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- 744 • If error bars are reported in tables or plots, The authors should explain in the text how they were
745 calculated and reference the corresponding figures or tables in the text.

746 8. Experiments Compute Resources

747 Question: For each experiment, does the paper provide sufficient information on the computer
748 resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

749 Answer: [Yes]

750 Justification: The type of compute used is provided in Appendix.

751 Guidelines:

752 • The answer NA means that the paper does not include experiments.

753 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud

754 provider, including relevant memory and storage.

755 • The paper should provide the amount of compute required for each of the individual experimental

756 runs as well as estimate the total compute.

757 • The paper should disclose whether the full research project required more compute than the

758 experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into

759 the paper).

760 **9. Code Of Ethics**

761 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code

762 of Ethics <https://neurips.cc/public/EthicsGuidelines>?

763 Answer: [Yes]

764 Justification: We conform with the NeurIPS Code of Ethics in every aspect.

765 Guidelines:

766 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.

767 • If the authors answer No, they should explain the special circumstances that require a deviation

768 from the Code of Ethics.

769 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due

770 to laws or regulations in their jurisdiction).

771 **10. Broader Impacts**

772 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts

773 of the work performed?

774 Answer: [Yes]

775 Justification: The potential societal impacts are discussed in Appendix.

776 Guidelines:

777 • The answer NA means that there is no societal impact of the work performed.

778 • If the authors answer NA or No, they should explain why their work has no societal impact or

779 why the paper does not address societal impact.

780 • Examples of negative societal impacts include potential malicious or unintended uses (e.g.,

781 disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deploy-

782 ment of technologies that could make decisions that unfairly impact specific groups), privacy

783 considerations, and security considerations.

784 • The conference expects that many papers will be foundational research and not tied to particular

785 applications, let alone deployments. However, if there is a direct path to any negative applications,

786 the authors should point it out. For example, it is legitimate to point out that an improvement in

787 the quality of generative models could be used to generate deepfakes for disinformation. On the

788 other hand, it is not needed to point out that a generic algorithm for optimizing neural networks

789 could enable people to train models that generate Deepfakes faster.

790 • The authors should consider possible harms that could arise when the technology is being used

791 as intended and functioning correctly, harms that could arise when the technology is being used

792 as intended but gives incorrect results, and harms following from (intentional or unintentional)

793 misuse of the technology.

794 • If there are negative societal impacts, the authors could also discuss possible mitigation strategies

795 (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitor-

796 ing misuse, mechanisms to monitor how a system learns from feedback over time, improving the

797 efficiency and accessibility of ML).

798 **11. Safeguards**

799 Question: Does the paper describe safeguards that have been put in place for responsible release of

800 data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or

801 scraped datasets)?

802 Answer: [NA]

803 Justification: This paper primarily focuses on a novel prompt learning method and doesn't involve

804 the release of a new pre-trained LLM, image generator, or scraped dataset. Therefore, this question

805 doesn't directly apply in this context. We leverage an existing pre-trained LLM (GPT-3.5-turbo), and

806 any ethical considerations regarding its release and potential misuse fall under the responsibility of its

807 creators.

808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We credited the creators of the CoOp codebase [54] by including the attribution statement in appendix.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [No]

Justification: Code will be released upon acceptance.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: Our paper doesn't involve crowdsourcing or research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.

- 863 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other
864 labor should be paid at least the minimum wage in the country of the data collector.

865 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

866 Question: Does the paper describe potential risks incurred by study participants, whether such
867 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an
868 equivalent approval/review based on the requirements of your country or institution) were obtained?

869 Answer: [NA]

870 Justification: This paper focuses on developing a novel prompt learning method and evaluating its
871 performance on established image recognition datasets. It doesn't involve any form of crowdsourcing,
872 human subject research, or data collection that would necessitate IRB approval or ethical considerations
873 related to study participants. Therefore, this question doesn't apply to our research.

874 Guidelines:

- 875 • The answer NA means that the paper does not involve crowdsourcing nor research with human
876 subjects.
- 877 • Depending on the country in which research is conducted, IRB approval (or equivalent) may be
878 required for any human subjects research. If you obtained IRB approval, you should clearly state
879 this in the paper.
- 880 • We recognize that the procedures for this may vary significantly between institutions and
881 locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for
882 their institution.
- 883 • For initial submissions, do not include any information that would break anonymity (if applica-
884 ble), such as the institution conducting the review.