

# LOREAL: Mitigating Low-Resolution Challenges in Vision-Language Models with Attribute-driven Prompt Self-Distillation

Anonymous CVPR submission

## Abstract

001 *Prompt Learning (PL) has emerged as a parameter-efficient*  
 002 *technique for adapting Vision-Language Models (VLMs)*  
 003 *to downstream tasks. However, almost all existing PL*  
 004 *methods are primarily designed and evaluated on well-*  
 005 *curated datasets, overlooking a critical post-deployment*  
 006 *phenomenon, i.e., the intrinsic connection between input*  
 007 *resolution and storage-memory consumption. Specifically,*  
 008 *to satisfy the stringent storage-memory constraints on edge*  
 009 *devices, models are often limited to low-resolution inputs*  
 010 *(e.g.,  $\leq 224 \times 224$  for CLIP-ViT/B-16) and generate fewer*  
 011 *tokens (with the position embedding resized), which poses*  
 012 *a unique challenge in performance robustness. To tackle*  
 013 *this issue, we propose LOREAL, an efficient prompt self-*  
 014 *distillation framework that learns resolution-invariant rep-*  
 015 *resentations by excavating attribute semantics. At the heart*  
 016 *of LOREAL is a dual-student architecture, i.e., two student*  
 017 *models fed with inputs at different resolutions synergisti-*  
 018 *cally learn from each other. Building upon this, we con-*  
 019 *textualize the students' prompt with resolution-invariant a-*  
 020 *tributes queried from the LLM, then leverage cross-modality*  
 021 *meta-nets to generate attribute semantics. These meta-*  
 022 *nets are bridged between the different encoders of two stu-*  
 023 *denters, wherein we introduce Low-Level Distillation (LLD)*  
 024 *and High-Level Distillation (HLD) to facilitate the learn-*  
 025 *ing of more cross-resolution representations. Extensive ex-*  
 026 *periments show that LOREAL significantly improves VLMs'*  
 027 *performance and robustness under varied resolution set-*  
 028 *tings, **underscoring significant practical utilities**. Code is*  
 029 *in the Supplementary Material.*

## 030 1. Introduction

031 The research community has witnessed a remarkable ad-  
 032 vancement of Vision-Language Models (VLMs) [24, 35,  
 033 36, 40, 53], which demonstrates promising zero-shot gen-  
 034 eralization ability enabled by their joint modeling of multi-  
 035 modal semantics. As a typical representative, CLIP [53]

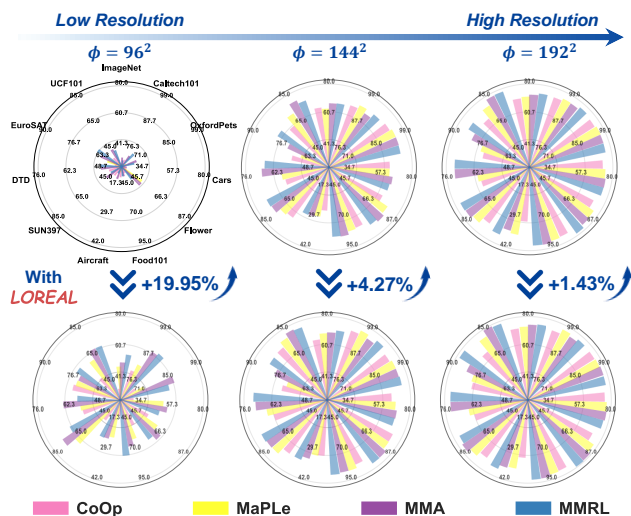


Figure 1. The Harmonic Mean (HM) of four SOTAs (CoOp [83], MaPLe [27], MMA [75], MMRL [15]) w/ or w/o our LOREAL on 11 datasets under varied resolution settings  $\phi$ . LOREAL substantially mitigates their performance degradations under LR settings.

employs contrastive pre-training over separate modality encoders to learn a unified manifold space. Its prominent cross-modal alignment capabilities have facilitated a wide range of downstream applications [3, 56, 67].

To better adapt VLMs for downstream tasks, an emerging trend focuses on conducting Parameter-Efficient Fine-Tuning (PEFT) [42, 78, 83] over VLMs. Among these techniques, Prompt Learning (PL) [25, 27, 43, 48, 82, 83] has gained significant traction: PL freezes the entire pre-trained VLM and fine-tunes only a small set of learnable prompts inserted into its layers, demonstrating remarkable parameter efficiency and strong compatibility. As a brief revisit, we categorize existing PL methods over VLMs along two dimensions: **(a)** according to the prompt types, there are (a): uni-modality methods, which employs either pure textual prompts [34, 77, 82, 83] as that for Large Language Models (LLMs), or pure visual prompts [73] for vision encoders; (b): bi-modality methods [15, 16, 27, 75], in which prompts are inserted into both encoders with explicit inter-

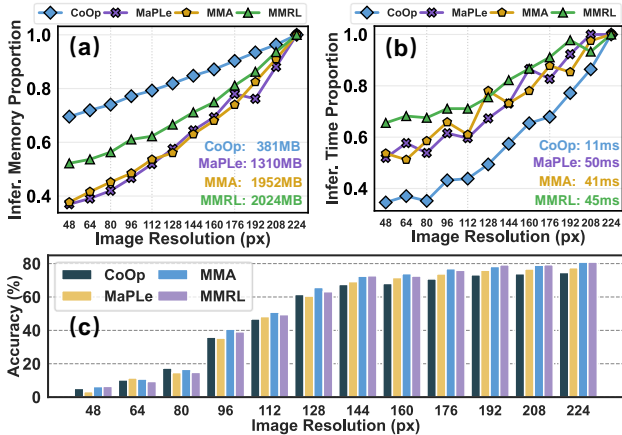


Figure 2. The (a) proportion of inference memory respective to that when  $\phi = 224^2$  (marked in **Lower-right**), (b) proportion of inference time respective to that when  $\phi = 224^2$  (marked in **Lower-right**), and (c) Harmonic Mean (HM) of four existing SOTAs under different resolution settings averaged over 11 datasets.

actions. **②** according to the learning scheme, there are (a): consistency-based methods [15, 47], which employs extra regularization loss to mitigate forgetting or overfitting; (b): procedure-based methods [34, 38, 39], which introduces refined tuning procedures like dual-prompt collaboration [34] or prompt-based knowledge distillation [38].

Despite their promising advancements, these methods typically maintain a fixed input resolution (typically,  $224 \times 224$  for CLIP-ViT-B/16), contradicting real-world applications where edge-deployed VLMs are required to process low-resolution inputs [14, 37, 66, 81] and generate fewer tokens due to constrained computational resources. We name this setting as the Low-Resolution (LR) setting. **The rationality of LR can be explained from two perspectives:** **①** the capturing, transmission, and storage overhead required for low-resolution images is nearly orders of magnitude lower than that for standard-resolution images, thus naturally adapting to the realities of edge devices. **②** since the complexity of the visual encoder accounts for a considerable portion of the overall complexity and polynomially increases with token number, reducing visual tokens for low-resolution images significantly cuts off the inference memory, thereby naturally meeting the requirements of real-time tasks. Furthermore, we collect the inference cost of four SOTAs under varied resolutions  $\phi$  in Figure 2 (a,b); the results unravel that reducing  $\phi$  can achieve up to 62% memory savings and 64% speed improvement during inference.

**Given that LR closely mirrors real-world requirements, how do existing models perform under it?** As shown in Figure 2 (c), we observe an overall degradation in existing SOTAs’ performance as  $\phi$  decreases; We also find that when  $\phi$  is reduced by approximately half, the model’s accuracy shows a dramatic plunge since discriminative visual features are largely blurred, posing unique challenges.

As the first attempt to handle this problem, this paper introduces **LOW-RE**olution **A**tttribute-guided prompt **L**earning (LOREAL), a prompt self-distillation scheme designed to equip VLMs with enhanced robustness against resolution shifts. Specifically, inspired by recent progress of automatic attribute-searching [39, 60], we first leverage the Large-Language Models (LLMs) to generate generic attributes which remain salient under varied resolutions. These attributes serve to refine the model’s awareness to a portion of resolution-invariant visual features. Then, we construct a contextualized prompt with attributes incorporated: ‘A photo of a [CLASS] with  $S_1$  [Attr1]  $S_2$  [Attr2], ...’, where  $S$  are the learnable attribute contents. Furthermore, we introduce cross-modality meta-nets, where each one is dedicated to one attribute and transfers the output visual embeddings into attribute contents. To enable the model to capture resolution-invariant semantics, we further introduce a self-distillation scheme, where two student networks that pin on different input resolutions share the meta-nets to generate attribute contents for each other. These two students further learn from each other via a Low-Level Distillation (LLD) that aligns the prompt attributes and a High-Level Distillation (HLD) that aligns the output logits. Extensive experiments show the superiority of our model (primarily shown in Figure 1). An architecture comparison between our LOREAL and other distillation methods is in Figure 3. Our contribution can be summarized as follows:

- We first study the Low-Resolution (LR) challenges in prompt learning of VLMs, and unveil the vulnerability of existing methods to this realistic setting.
- We propose LOREAL, a prompt self-distillation framework, as the first endeavor towards this problem. LOREAL is guided by resolution-invariant attributes and shared meta-nets to refine the resolution robustness.
- Extensive experiments on various datasets and benchmarks demonstrate the superiority and compatibility of our LOREAL to the realistic LR challenges.

## 2. Related Work

**Prompt Learning of Vision-Language Models.** Prompt Learning (PL) originates in natural language processing to steer the Language Model toward producing optimal responses. The pioneer work CoOp first leverages PL to adapt Vision-Language Models (VLMs) with a learnable template combined with classnames. Further advancement focuses on **①** strengthening the modality interactions by generating text prompts conditioned on visual priors [73, 77, 82]; **②** incorporating multi-modal prompts in intermediate layers [15, 27, 28, 73, 75], typically represented by MaPLe [27] which leverages layer-wise entangled prompts for both encoders; **③** employing distribution alignment to mitigate bias and avoid forgetting in finetuning [15, 28, 73, 84]. For example, PromptSRC [28] and MMRL [28] introduce the

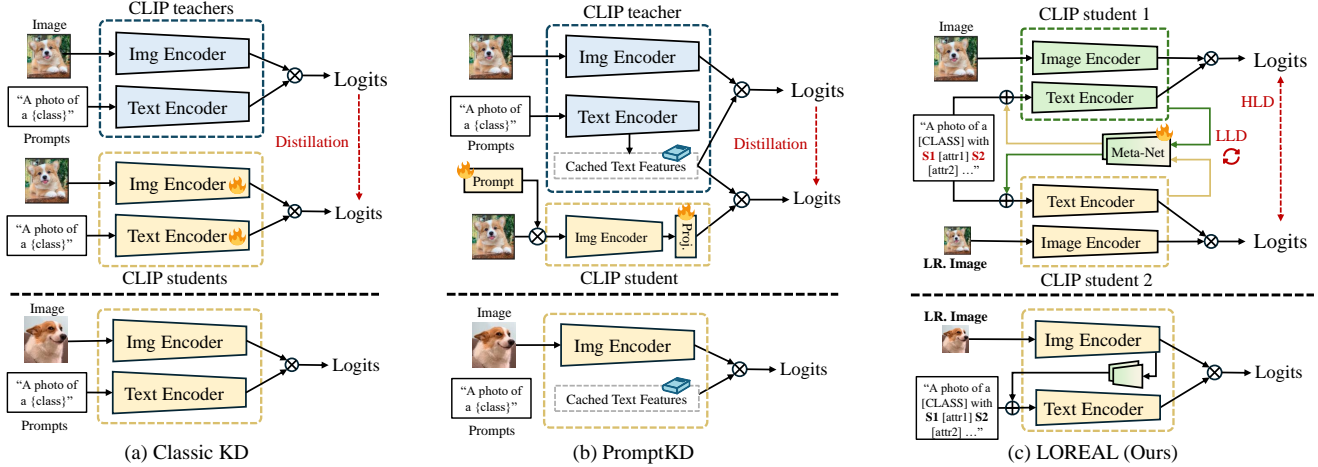


Figure 3. Comparisons of Classic KD [74], PromptKD [38] and our LOREAL. The upper/lower part is the training/inference stage. TE means the Text Encoder. LR means Low-Resolution. (a) Classic KD of VLMs, where students are fully-tuned. (b) PromptKD, which leverages prompts to learning from teachers. Both (a) and (b) are designed for non-LR inference. (c) The proposed LOREAL, a prompt self-distillation scheme to solve the LR challenges. Here, two students are the same models but fed with different inputs. LOREAL leverages fine-grained attribute guidance and simultaneously distills from two levels to boost the model’s robustness to data resolutions.

141 prediction from zero-shot VLM to guide each training step. 142 Recently, Skip-Tuning [71] caches intermediate visual features 143 and drops gradients of unrelated classes for more efficient 144 finetuning. PL for VLMs has also shown wide applications 145 in knowledge distillation [38, 70], federated learning [52], 146 and interpretability analysis [10, 11]. *Notably, although both 147 Skip-Tuning [71] and our work aim to explore lightweighting, 148 it fails to discuss inference memory savings induced by LR 149 settings during inference.*

150 **Knowledge Distillation (KD).** Knowledge Distillation (KD) [21] 151 aims to learn a lightweight student network under the guidance 152 of the bulky teacher, thereby achieving comparable performance 153 with greater efficiency and empowering the latent-sensitive 154 applications. Based on the type of distilling sources, KD can 155 be categorized into Logit Distillation (LD) [26, 30, 45, 80] (i.e., 156 only to align the final logits), and Feature Distillation (FD) [4, 5, 17, 49, 55, 61–63, 76] 157 which employ extra intermediate features for alignment. Self-Distillation (SD) 158 is a special variant of KD, where the students act as the teachers 159 of themselves and learn from each other. Some recent KD methods 160 [29, 32, 38, 74] turn their gaze to the VLMs, for example, Clip-KD [74] 161 systematically studies the adaptability of distillation techniques 162 to VLMs; PromptKD [38] only leverages prompts to adapt 163 knowledge from teachers. COSMOS [29] introduces a cross-modality 164 self-distillation for VLM pretraining. 165 166

167 **Low-Resolution Recognition (LRR).** Low-Resolution 168 Recognition [69] (LRR) aims at maintaining the model’s 169 performance when only Low-Resolution (LR) images are 170 available in inference. Early LR methods are designed for 171 face recognition [57, 68, 85]; Subsequent works [13, 51, 65] 172 have explored the integration of Knowledge Distillation

(KD) to obtain compact yet high-performing models suitable 173 for LR inference: for example, PixelDistillation [14] 174 introduces a cost-flexible heterogeneous framework to distill 175 small convolution models for fast inference; When it comes to 176 LRR on large pre-trained models, existing studies [1, 8, 33, 41, 59] 177 mostly focus on designing position embeddings that enables 178 flexible input resolution for ViTs; for example, ResFormer [59] 179 proposes a multi-resolution training and global-local embedding 180 strategy to learn resolution-invariant semantics; MSPE [41] 181 designs a plug-and-play hierarchical position embedding to adapt 182 different patch sizes. *Despite these advancements, addressing the 183 LR challenge through prompt learning for VLMs remains 184 unexplored.* 185

### 3. Methodology 186

#### 3.1. Preliminaries 187

**Prompt Learning (PL) of Vision-Language (VLMs).** 188 Consider a typical VLM, i.e. CLIP; it consists of a visual 189 encoder  $\mathcal{E}_v$  and text encoder  $\mathcal{E}_t$  where each comprises  $L$  layers, 190 i.e.  $\mathcal{E}_v = \{\mathcal{E}_{v,i}\}_{i=1}^L$ ,  $\mathcal{E}_t = \{\mathcal{E}_{t,i}\}_{i=1}^L$ ; For image branch 191  $\mathcal{V}$ , denote  $\mathcal{B}$  as an image batch; for one image  $x$  from  $\mathcal{B}$ , 192 a patch embedding module  $\text{emb}$  first project it into  $M$ -length 193 visual embeddings  $\mathbf{v}_0$  ( $m$  is the fixed token number), then, 194 the process of each layer  $\mathcal{E}_{v,i}$  can be formalized as: 195

$$[\text{cls}_i, \mathbf{v}_i] = \mathcal{E}_{v,i}([\text{cls}_{i-1}, \mathbf{v}_{i-1}]), \quad i \in \{1, 2, \dots, L\}, \quad (1) \quad 196$$

where  $\text{cls}_i \in \mathbb{R}^{D_v}$  means the class embeddings of layer  $i$ . 197  $D_v$  means the visual hidden dimension. We apply another 198 projector  $M_v$  over  $\text{cls}_L$  to generate the final visual embeddings 199  $\mathbf{f}_v$ , i.e.  $\mathbf{f}_v = \text{cls}_L M_v$ . For the text branch, we fill all 200 classnames into the prompt  $P$  that is initialized from either 201

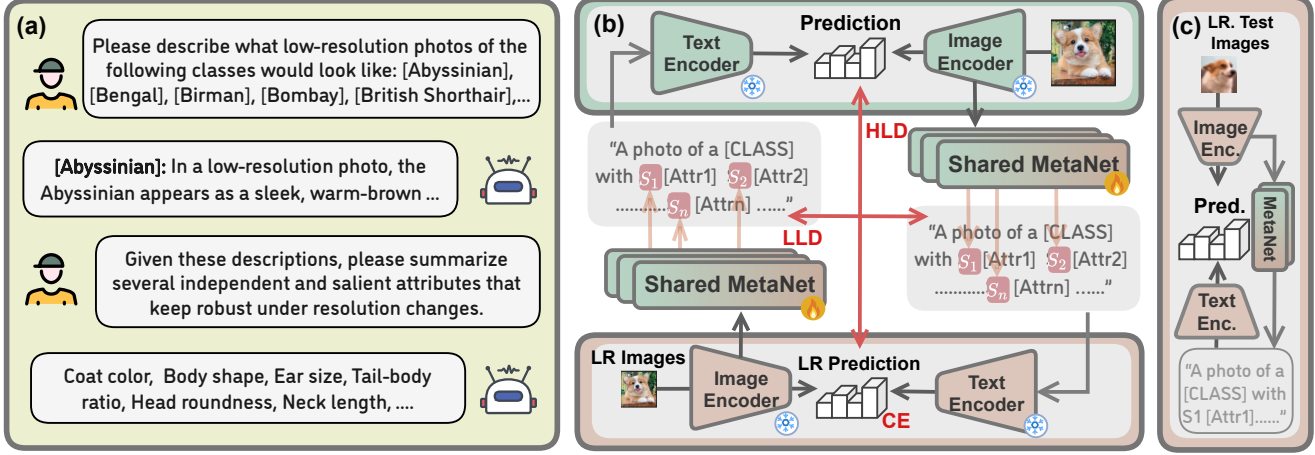


Figure 4. Our LOREAL framework. (a): We leverage the LLM to generate several resolution-invariant attributes. (b): Self-distillation framework. We utilize the visual embeddings to fill the prompt attributes via meta-nets, then leverage Low-Level Distillation (LLD) and High-Level Distillation (HLD) for self-distillation. Only the meta-nets are learnable, and the parameters of two illustrated meta-nets are shared. LR represents Low-Resolution. (c): Inference stage. The model takes LR images and contextualizes prompts with the meta-nets.

202 pre-defined (i.e. “A photo of a [CLASS]”) [82, 83] or random  
 203 values  $\mathcal{P}^0$  to generate the text input. Denote  $\mathcal{C}$  as all  
 204 classes, and  $C$  as the number of classes. The text input is  
 205 then tokenized and converted into text embeddings  $\mathbf{t}_0$ . Finally,  
 206  $\mathbf{t}_0$ , together with a Begin-Of-Text (BoT) token  $\text{bot}$  and  
 207 End-Of-Text (EOT) token  $\text{eot}$ , is fed into each  $\mathcal{E}_{t,i}$ :

$$208 \quad [\text{bot}_i, \mathbf{t}_i, \text{eot}_i] = \mathcal{E}_{t,i}([\text{bot}_{i-1}, \mathbf{t}_{i-1}, \text{eot}_{i-1}]). \quad (2)$$

209 We pick up the EOT token  $\text{eot}_L$  from the last layer  $L$ , and  
 210 project it to the modality-shared space with projector  $M_t$ ,  
 211 i.e.  $\mathbf{f}_t = \text{eot}_L M_t$ . We may also denote  $\mathbf{f}_t(P)$  as the textual  
 212 outputs conditioned on prompt  $P$ . For the classification  
 213 stage, we calculate the cosine similarity  $\text{Sim}(\cdot)$  between  $\mathbf{f}_v$   
 214 and each entry in  $\mathbf{f}_t \in \mathbb{R}^{C \times d}$  as the raw prediction, i.e.:

$$215 \quad \hat{\mathbf{y}}_c = \frac{\exp(\text{Sim}(\mathbf{f}_v, \mathbf{f}_{t,c})/\tau)}{\sum_{c'} \exp(\text{Sim}(\mathbf{f}_v, \mathbf{f}_{t,c'})/\tau)}, \quad (3)$$

216 where  $\mathbf{f}_{t,c}$  is the  $c$ -th entry of  $\mathbf{f}_t$ .  $\tau$  is the temperature. We  
 217 also denote  $\hat{\mathbf{y}}(\mathbf{f}_v, \mathbf{f}_{t,c})$  as  $\hat{\mathbf{y}}$  generated from  $\mathbf{f}_v, \mathbf{f}_{t,c}$ .

218 **Low-Resolution Recognition (LRR) with VLMs.** As  
 219 we introduced above, LRR generally represents that models  
 220 are directly fed with low-resolution inputs and generate  
 221 fewer tokens accordingly. For LRR on VLMs, considering  
 222 that the length  $J$  of visual position embedding  $\mathcal{P} \in \mathbb{R}^{J \times D_v}$   
 223 is fixed, we follow baseline implementations [23] and re-  
 224 shape it to dynamically meet the shape requirements  $(\phi, \phi)$   
 225 (assume  $\phi$  is the resolution setting of the input image):

$$226 \quad \mathcal{P} = \text{Cat}(\mathcal{P}[:1], \text{Intp}(\mathcal{P}[1:], \phi, \phi)), \quad (4)$$

227 where  $[a:b]$  means slicing, and  $[:1]$  separates the location  
 228 of  $\text{cls}$ -token apart.  $\text{Cat}$  means concatenation,  $\text{Intp}$  is the  
 229 interpolation function (bicubic mode by default).

### 230 3.2. Generating Resolution-invariant Attributes

231 While the Low-Resolution (LR) settings may induce con-  
 232 siderate vague in class distinctions, a natural idea is diving  
 233 into attribute-level semantics in the image and focusing on  
 234 those that tend to remain robust under varied resolution. In-  
 235 spired by recent progress in attribute searching [6, 39, 79],  
 236 we design an LLM-based pipeline with explicit Chain-of-  
 237 Thought (CoT) to search resolution-invariant attributes.

238 Concretely, as shown in Fig. 4, we first instruct the LLM  
 239 to generate extensive descriptions of how LR photos would  
 240 typically appear for all [CLASSES]. Leveraging these de-  
 241 scriptions as contextual background, we then instruct the  
 242 LLM to carefully summarize a set of generic and invari-  
 243 ant attributes that remain perceptible across varying res-  
 244 olutions. This process favors macro-level attributes (e.g.,  
 245 coat-color and body-shape for OxfordPets [50];  
 246 window and roofline for StanfordCars [31]) while fil-  
 247 tering out more subtle ones (e.g., texture, eye-color  
 248 for OxfordPets; logo or badge for StanfordCars). Note  
 249 that both macro and subtle attributes are considered generic  
 250 in standard non-LR settings [39], and *our key insight is to*  
 251 *identify and utilize the subset of these generic features that*  
 252 *demonstrably retain their saliency under LR.*

253 Subsequently, with the LLM-generated raw attributes  
 254  $\{A_k\}_{k=1}^K$ , we build contextualized prompts  $\mathbf{p} = \{\mathbf{p}_c\}_{c=1}^C$   
 255 by filling them into the attribute slot [A.] as the following  
 256 to actively refine model’s awareness:

257 A photo of a [CLS] with S<sub>1</sub> [A<sub>1</sub>] S<sub>2</sub> [A<sub>2</sub>]  $\cdots$  S<sub>K</sub> [A<sub>K</sub>],

258 where the learnable tokens are underlined,  $S_i \in \mathbb{R}^{M \times D_t}$   
 259 represents the learnable contents of attributes, and  $M$  is the  
 260 learnable token number in each attribute. Notably, the or-  
 261 der of attributes is not taken into specific consideration as  
 262 it yields minor effects on the performance. Details of the

263 attributes and attribute order are in **Sup. Mat. C.** To incor-  
 264 porate visual semantics into the prompt, we contextualize  
 265  $S_i$  using cross-modality meta-nets in the following.

### 266 3.3. Cross-modality Meta-Nets

267 In correspondence with the attributes, we introduce  $K$  sep-  
 268 arate meta-nets, denoted as  $\{\mathcal{M}_k\}_{k=1}^K$  to dynamically con-  
 269 textualize prompts with visual semantics. Specifically, the  
 270 meta-nets are built upon a LoRA-like architecture, and the  
 271 output visual embeddings  $\mathbf{f}_v$  are fed into each meta-net to  
 272 generate their respective  $\{S_k\}_{k=1}^K$ :

$$273 \quad S_k = M_k(\mathbf{f}_v) = W_{\uparrow,k}(W_{\downarrow,k}(\mathbf{f}_v)), \quad (5)$$

274 where  $W_{\downarrow,k} \in \mathbb{R}^{D_v \times D_s}$ ,  $W_{\uparrow,k} \in \mathbb{R}^{D_s \times D_t}$  are LoRA pro-  
 275 jection matrices,  $D_s$  represents LoRA intermediate dimen-  
 276 sion. We denote  $p(S_k)$  as prompts conditioned on  $S_k$ .

### 277 3.4. Attribute-driven Self-Distillation

278 Despite the straightforward idea of finetuning attribute-  
 279 contextualized prompts to focus on invariant properties, the  
 280 model’s exposure solely to LR images during inference nec-  
 281 essitates a more sophisticated approach. We therefore intro-  
 282 duce a self-distillation method that guides the model in  
 283 reconciling semantics across different resolutions.

284 Our self-distillation incorporates two student VLMs  
 285  $\{\mathcal{E}_t^\alpha, \mathcal{E}_v^\alpha\}$   $\{\mathcal{E}_t^\beta, \mathcal{E}_v^\beta\}$  with all backbone frozen and only the  
 286 meta-nets learnable. One student  $\{\mathcal{E}_t^\alpha, \mathcal{E}_v^\alpha\}$  is fed with stan-  
 287 dard image  $\mathbf{x}$  and the other one  $\{\mathcal{E}_t^\beta, \mathcal{E}_v^\beta\}$  is fed with low-  
 288 resolution image  $\mathbf{x}'$  (when combined with other methods  
 289 besides zero-shot CLIP, the students are pre-tuned with im-  
 290 ages of their separate resolution). Recall the cross-modality  
 291 meta-nets in the last subsection, we bridge these meta-nets  
 292 between the cross-modal encoders of different students:

$$293 \quad \begin{aligned} \mathbf{f}_v^\alpha &= \mathcal{E}_v^\alpha(\mathbf{x}); & \mathbf{f}_t^\beta &= \mathcal{E}_t^\beta(p(S(\mathbf{f}_v^\alpha))) \\ \mathbf{f}_v^\beta &= \mathcal{E}_v^\beta(\mathbf{x}'); & \mathbf{f}_t^\alpha &= \mathcal{E}_t^\alpha(p(S(\mathbf{f}_v^\beta))) \end{aligned} \quad (6)$$

294 Then, to enable inputs of different resolution to generate  
 295 similar attribute contexts, we introduce Low-Level Distilla-  
 296 tion (LLD), where we employ the contrastive learning be-  
 297 tween the prompt contexts  $S(\mathbf{f}_v^\alpha)$  and  $S(\mathbf{f}_v^\beta) \in \mathbb{R}^{K \times D_t}$ :

$$298 \quad \mathcal{L}_{\text{LLD}} = - \sum_{k=1}^K \log \frac{\exp(\text{Sim}(S_k(\mathbf{f}_v^\alpha), S_k(\mathbf{f}_v^\beta))/\tau)}{\sum_{k'}^K \exp(\text{Sim}(S_k(\mathbf{f}_v^\alpha), S_{k'}(\mathbf{f}_v^\beta))/\tau)} \quad (7)$$

299 Meanwhile, we generate the predictions  $\hat{\mathbf{y}}^\alpha \hat{\mathbf{y}}^\beta$  for both stu-  
 300 dents, i.e:  $\hat{\mathbf{y}}^\alpha = \hat{\mathbf{y}}(\mathbf{f}_v^\alpha, \mathbf{f}_{t,c}^\alpha)$ ,  $\hat{\mathbf{y}}^\beta = \hat{\mathbf{y}}(\mathbf{f}_v^\beta, \mathbf{f}_{t,c}^\beta)$ . These sep-  
 301 arate predictions are further aligned with Kullback-Leibler  
 302 (KL)-Divergence as our High-Level Distillation (HLD):

$$303 \quad \mathcal{L}_{\text{HLD}} = \sum_{c=1}^C (\hat{y}_c^\alpha \log \hat{y}_c^\alpha - \hat{y}_c^\alpha \log \hat{y}_c^\beta) \quad (8)$$

304 These designs enable the textual prompts to be dynamically  
 305 contextualized by visual features from another resolution,

thereby progressively aligning them in the multi-modal  
 manifold space. Finally, our two distillation losses are har-  
 monized with the task loss  $\mathcal{L}_{\text{CE}} = - \sum_c y_c \log \hat{y}_c^\beta$  via two  
 coefficients  $\lambda_1 \lambda_2$ , i.e,  $\mathcal{L} = \mathcal{L}_{\text{CE}} + \lambda_1 \mathcal{L}_{\text{HLD}} + \lambda_2 \cdot 1/K \cdot \mathcal{L}_{\text{LLD}}$ .  
 In the inference stage, we use the meta-nets to transfer vi-  
 sual embeddings into attribute contents, then calculate the  
 predictions leveraging the attribute-structured prompt.

## 4. Experiments

### 4.1. Settings

**Datasets.** We evaluate our LOREAL on three benchmarks:

316 **1 Low-Resolution Base-to-New (LR-B2N) benchmark,**  
 317 which integrates LR with Base-to-New tasks to further evalu-  
 318 ate the models’ generalization robustness. The datasets are  
 319 first split into disjoint base and new classes, then models  
 320 are trained on base classes and evaluated on both base and  
 321 new classes. The test-sets are down-sampled to the resolu-  
 322 tion  $\phi \in \{96^2, 144^2, 192^2\}$ , and datasets employed for this  
 323 task are: ImageNet [9], OxfordPets [50], SUN397 [72], Eu-  
 324 roSAT [18], Caltech101 [12], StanfordCars [31], DTD [7],  
 325 UCF101 [58], Flowers102 [46], Food101 [2] and FGVC-  
 326 Aircraft [44]. **2 Low-Resolution Cross-dataset Evaluation**  
 327 **(LR-CE) benchmark,** which integrates LR with a cross-  
 328 dataset task to evaluate models’ robustness to LR test-sets  
 329 of new datasets. We train the models on ImageNet and evalu-  
 330 ate them on the above ten datasets down-sampled to  $\phi$ . **3**  
 331 **Low-Resolution Domain Generalization (LR-DG) bench-**  
 332 **mark,** which integrates LR with domain-generalization to  
 333 evaluate models’ robustness to LR test-sets of new domains.  
 334 We train models on ImageNet and evaluate them on four  
 335 variants (ImageNet-A [20], ImageNet-V2 [54], ImageNet-  
 336 R [19], and ImageNet-S [64]) down-sampled to  $\phi$ . All ex-  
 337 periments adopt a 16-shot setting, i.e., 16 samples per class.

**Implementation Details.** Following previous studies, we  
 338 adopt pretrained CLIP-VIT-B/16 as the backbone.  $\tau = 1$ .  
 339 We use the SGD optimizer with a learning rate of 0.002.  
 340 According to the ablation studies, the intermediate hidden  
 341 dimension  $D_s$  of meta-nets is 32;  $M = 2, \lambda_1 = 1, \lambda_2 = 2$ .  
 342 The employed LLM is GPT-4o [22], and we generate 5 at-  
 343 tributes for each class. Notably, when combining our model  
 344 with methods that cannot perform separate visual or text  
 345 embeddings [15, 27], we proactively store their visual em-  
 346 beddings offline before distillation and leverage the stored  
 347 embeddings to fill in the prompts. The training epoch and  $\phi$   
 348 are unique for each method. More implementation specifics  
 349 and experiment results are in **Sup. Mat. A / B.**

### 4.2. Experiment Results

**Results on the LR-B2N Benchmark.** As demonstrated in  
 352 Table 1, integrating our LOREAL with four existing SOTAs  
 353 under different  $\phi$  settings consistently yields a substantial  
 354 performance gains. Specifically, we have three key find-  
 355

Table 1. Results on the LR-B2N Benchmark over three resolution settings  $\phi \in \{96^2, 144^2, 192^2\}$  w/ or w/o our LOREAL. HM means the harmonic mean.  $\eta$  is the ratio of patch tokens relative to the  $224^2$  input. All results are averaged over 3 runs. **Bold** marks the best results.

	Method	Average			ImageNet			Caltech101			OxfordPets			StanfordCars			Flowers102		
		Base	Novel	HM	Base	Novel	HM	Base	Novel	HM	Base	Novel	HM	Base	Novel	HM	Base	Novel	HM
		$\phi : (96 \times 96), \eta \approx 18.37\%$																	
	CoOp	38.40	33.74	35.54	24.92	23.30	24.08	70.95	65.16	67.93	58.72	57.45	58.08	17.37	14.19	15.62	43.97	21.42	28.81
	+LOREAL	<b>54.81</b>	<b>42.77</b>	<b>47.48</b>	<b>35.96</b>	<b>40.84</b>	<b>38.24</b>	<b>83.21</b>	<b>72.60</b>	<b>77.54</b>	<b>66.92</b>	<b>64.20</b>	<b>65.53</b>	<b>27.61</b>	<b>18.16</b>	<b>21.91</b>	<b>76.73</b>	<b>35.04</b>	<b>48.11</b>
	MaPLe	37.17	33.71	34.85	21.24	26.74	23.67	69.34	62.88	65.95	53.32	66.44	59.16	12.42	12.82	12.62	50.90	37.52	43.20
	+LOREAL	<b>65.02</b>	<b>52.18</b>	<b>57.25</b>	<b>46.07</b>	<b>42.76</b>	<b>44.35</b>	<b>85.48</b>	<b>76.14</b>	<b>80.54</b>	<b>69.17</b>	<b>70.43</b>	<b>69.79</b>	<b>52.22</b>	<b>41.89</b>	<b>46.49</b>	<b>86.32</b>	<b>43.40</b>	<b>57.76</b>
	MMA	41.55	39.91	40.50	30.99	25.88	28.21	73.98	66.92	70.27	60.77	69.30	64.76	15.57	16.17	15.86	51.66	38.58	44.17
	+LOREAL	<b>67.81</b>	<b>59.74</b>	<b>63.14</b>	<b>51.80</b>	<b>44.29</b>	<b>47.75</b>	<b>92.32</b>	<b>84.61</b>	<b>88.30</b>	<b>83.47</b>	<b>87.47</b>	<b>85.42</b>	<b>54.10</b>	<b>46.01</b>	<b>49.73</b>	<b>84.90</b>	<b>51.49</b>	<b>64.10</b>
	MMRL	41.64	36.92	38.90	29.57	23.44	26.15	75.98	68.36	71.97	58.75	68.67	63.32	14.02	11.54	12.66	40.36	33.83	36.81
	+LOREAL	<b>66.97</b>	<b>57.80</b>	<b>61.71</b>	<b>52.54</b>	<b>45.34</b>	<b>48.68</b>	<b>93.29</b>	<b>85.26</b>	<b>89.09</b>	<b>81.61</b>	<b>82.27</b>	<b>81.94</b>	<b>46.53</b>	<b>42.63</b>	<b>44.49</b>	<b>79.36</b>	<b>48.72</b>	<b>60.38</b>
$\phi : (144 \times 144), \eta \approx 41.5\%$																			
	CoOp	23.03	30.37	26.20	10.20	9.30	9.73	35.24	30.86	32.90	47.45	36.47	41.24	54.36	52.05	53.18	36.14	30.56	33.12
	+LOREAL	<b>50.92</b>	<b>38.96</b>	<b>44.14</b>	<b>13.87</b>	<b>10.14</b>	<b>11.72</b>	<b>52.83</b>	<b>36.03</b>	<b>42.84</b>	<b>62.85</b>	<b>47.45</b>	<b>54.07</b>	<b>79.71</b>	<b>67.41</b>	<b>73.05</b>	<b>52.33</b>	<b>39.64</b>	<b>45.11</b>
	MaPLe	19.67	22.83	21.13	7.02	7.86	7.42	26.33	28.03	27.15	50.20	29.71	37.33	64.60	51.13	57.08	33.87	24.82	28.65
	+LOREAL	<b>56.91</b>	<b>55.35</b>	<b>56.12</b>	<b>23.65</b>	<b>19.98</b>	<b>21.66</b>	<b>68.60</b>	<b>65.71</b>	<b>67.12</b>	<b>73.84</b>	<b>43.24</b>	<b>54.54</b>	<b>82.52</b>	<b>64.23</b>	<b>72.24</b>	<b>70.48</b>	<b>50.89</b>	<b>59.10</b>
	MMA	27.36	30.29	28.75	10.38	11.10	10.73	37.57	38.27	37.92	49.88	40.94	44.97	57.48	64.56	60.81	41.37	36.99	39.06
	+LOREAL	<b>61.90</b>	<b>61.55</b>	<b>61.72</b>	<b>22.63</b>	<b>22.80</b>	<b>22.71</b>	<b>74.50</b>	<b>75.39</b>	<b>74.94</b>	<b>72.80</b>	<b>54.83</b>	<b>62.55</b>	<b>74.90</b>	<b>68.92</b>	<b>71.79</b>	<b>72.54</b>	<b>59.76</b>	<b>65.53</b>
	MMRL	30.47	27.06	28.66	7.56	9.06	8.24	34.49	32.66	33.55	57.25	37.44	45.27	70.05	58.23	63.60	39.55	35.88	37.63
	+LOREAL	<b>70.46</b>	<b>69.43</b>	<b>69.94</b>	<b>22.93</b>	<b>22.62</b>	<b>22.77</b>	<b>68.55</b>	<b>63.46</b>	<b>65.91</b>	<b>69.56</b>	<b>46.74</b>	<b>55.91</b>	<b>81.74</b>	<b>68.08</b>	<b>74.29</b>	<b>70.10</b>	<b>61.28</b>	<b>65.39</b>
$\phi : (192 \times 192), \eta \approx 73.47\%$																			
	CoOp	81.36	66.68	73.29	75.20	65.24	69.87	98.19	91.48	94.72	94.42	96.14	95.27	76.06	66.50	70.96	97.72	64.54	77.74
	+LOREAL	<b>82.55</b>	<b>68.60</b>	<b>74.93</b>	<b>77.10</b>	<b>66.18</b>	<b>71.22</b>	<b>98.32</b>	<b>91.70</b>	<b>94.89</b>	<b>95.84</b>	<b>96.42</b>	<b>96.13</b>	<b>78.81</b>	<b>69.52</b>	<b>73.87</b>	<b>97.15</b>	<b>65.62</b>	<b>78.33</b>
	MaPLe	82.19	70.75	76.04	75.41	66.28	70.55	98.05	93.56	95.75	95.25	97.04	96.14	74.71	71.30	72.97	97.34	70.99	82.10
	+LOREAL	<b>83.32</b>	<b>72.60</b>	<b>77.59</b>	<b>76.95</b>	<b>67.98</b>	<b>72.19</b>	<b>98.46</b>	<b>94.76</b>	<b>96.57</b>	<b>94.95</b>	<b>97.48</b>	<b>96.20</b>	<b>76.39</b>	<b>74.53</b>	<b>75.45</b>	<b>97.87</b>	<b>72.44</b>	<b>83.26</b>
	MMA	81.28	75.47	78.27	76.25	68.45	72.14	98.39	92.90	95.57	94.36	97.54	95.92	75.89	72.25	74.03	<b>97.63</b>	71.56	82.59
	+LOREAL	<b>82.61</b>	<b>76.98</b>	<b>79.70</b>	<b>78.32</b>	<b>69.28</b>	<b>73.52</b>	<b>98.44</b>	<b>93.01</b>	<b>95.65</b>	<b>95.52</b>	<b>98.06</b>	<b>96.77</b>	<b>77.41</b>	<b>75.12</b>	<b>76.25</b>	97.24	<b>76.24</b>	<b>85.47</b>
	MMRL	83.32	75.55	79.25	76.18	69.89	72.90	98.84	94.32	96.53	95.22	96.39	95.80	76.26	71.75	73.94	98.39	74.18	84.59
	+LOREAL	<b>84.39</b>	<b>76.67</b>	<b>80.34</b>	<b>77.60</b>	<b>70.85</b>	<b>74.07</b>	<b>99.10</b>	<b>94.43</b>	<b>96.71</b>	<b>96.01</b>	<b>97.09</b>	<b>96.55</b>	<b>76.89</b>	<b>73.00</b>	<b>74.89</b>	<b>98.48</b>	<b>76.11</b>	<b>85.86</b>
$\phi : (96 \times 96), \eta \approx 18.37\%$																			
	CoOp	23.03	30.37	26.20	10.20	9.30	9.73	35.24	30.86	32.90	47.45	36.47	41.24	54.36	52.05	53.18	36.14	30.56	33.12
	+LOREAL	<b>50.92</b>	<b>38.96</b>	<b>44.14</b>	<b>13.87</b>	<b>10.14</b>	<b>11.72</b>	<b>52.83</b>	<b>36.03</b>	<b>42.84</b>	<b>62.85</b>	<b>47.45</b>	<b>54.07</b>	<b>79.71</b>	<b>67.41</b>	<b>73.05</b>	<b>52.33</b>	<b>39.64</b>	<b>45.11</b>
	MaPLe	19.67	22.83	21.13	7.02	7.86	7.42	26.33	28.03	27.15	50.20	29.71	37.33	64.60	51.13	57.08	33.87	24.82	28.65
	+LOREAL	<b>56.91</b>	<b>55.35</b>	<b>56.12</b>	<b>23.65</b>	<b>19.98</b>	<b>21.66</b>	<b>68.60</b>	<b>65.71</b>	<b>67.12</b>	<b>73.84</b>	<b>43.24</b>	<b>54.54</b>	<b>82.52</b>	<b>64.23</b>	<b>72.24</b>	<b>70.48</b>	<b>50.89</b>	<b>59.10</b>
	MMA	27.36	30.29	28.75	10.38	11.10	10.73	37.57	38.27	37.92	49.88	40.94	44.97	57.48	64.56	60.81	41.37	36.99	39.06
	+LOREAL	<b>61.90</b>	<b>61.55</b>	<b>61.72</b>	<b>22.63</b>	<b>22.80</b>	<b>22.71</b>	<b>74.50</b>	<b>75.39</b>	<b>74.94</b>	<b>72.80</b>	<b>54.83</b>	<b>62.55</b>	<b>74.90</b>	<b>68.92</b>	<b>71.79</b>	<b>72.54</b>	<b>59.76</b>	<b>65.53</b>
	MMRL	30.47	27.06	28.66	7.56	9.06	8.24	34.49	32.66	33.55	57.25	37.44	45.27	70.05	58.23	63.60	39.55	35.88	37.63
	+LOREAL	<b>70.46</b>	<b>69.43</b>	<b>69.94</b>	<b>22.93</b>	<b>22.62</b>	<b>22.77</b>	<b>68.55</b>	<b>63.46</b>	<b>65.91</b>	<b>69.56</b>	<b>46.74</b>	<b>55.91</b>	<b>81.74</b>	<b>68.08</b>	<b>74.29</b>	<b>70.10</b>	<b>61.28</b>	<b>65.39</b>
$\phi : (144 \times 144), \eta \approx 41.5\%$																			
	CoOp	74.21	61.88	67.00	67.71	59.10	63.11	94.13	89.25	91.63	91.44	93.06	92.24	65.39	57.04	60.93	89.36	56.74	69.41
	+LOREAL	<b>77.30</b>	<b>67.37</b>	<b>71.73</b>	<b>69.22</b>	<b>63.92</b>	<b>66.46</b>	<b>97.32</b>	<b>91.66</b>	<b>94.41</b>	<b>92.77</b>	<b>93.96</b>	<b>93.36</b>	<b>67.89</b>	<b>59.98</b>	<b>63.69</b>	<b>94.40</b>	<b>68.73</b>	<b>79.55</b>
	MaPLe	74.59	64.54	68.80	67.79	57.55	62.25	94.25	91.48	92.84	92.40	<b>94.74</b>	93.56	63.39	59.10	61.17	93.45	64.96	76.64
	+LOREAL	<b>79.46</b>	<b>69.32</b>	<b>73.72</b>	<b>72.79</b>	<b>63.26</b>	<b>67.69</b>	<b>97.81</b>	<b>94.21</b>	<b>95.98</b>	<b>94.26</b>	94.69	<b>94.47</b>	<b>70.89</b>	<b>63.53</b>	<b>67.01</b>	<b>95.25</b>	<b>69.03</b>	<b>80.05</b>
	MMA	75.21	69.90	72.30	70.13	62.90	66.32	94.77	87.81	91.16	91.23	95.64	93.38	65.39	63.21	64.28	90.88	69.43	78.72
	+LOREAL	<b>78.68</b>	<b>73.08</b>	<b>75.60</b>	<b>72.01</b>	<b>64.78</b>	<b>68.20</b>	<b>97.35</b>	<b>91.59</b>	<b>94.38</b>	<b>92.93</b>	<b>95.97</b>	<b>94.43</b>	<b>70.21</b>	<b>67.18</b>	<b>68.66</b>	<b>95.54</b>	<b>72.36</b>	<b>82.35</b>
	MMRL	76.92	68.94	72.43	70.67	62.93	66.58	95.35	90.67	92.95	92.24	94.07	93.15	63.92	60.58	62.21	93.45	65.60	77.09
	+LOREAL	<b>80.93</b>	<b>73.03</b>	<b>76.56</b>	<b>73.74</b>	<b>66.38</b>	<b>69.87</b>	<b>97.81</b>	<b>94.23</b>	<b>95.99</b>	<b>94.21</b>	<b>95.53</b>	<b>94.87</b>	<b>69.74</b>	<b>66.28</b>	<b>67.97</b>	<b>97.06</b>	<b>71.99</b>	<b>82.67</b>
$\phi : (192 \times 192), \eta \approx 73.47\%$																			
	CoOp	81.36	66.68	73.29	75.20	65.24	69.87	98.19	91.48	94.72	94.42	96.14	95.27	76.06	66.50	70.96	97.72	64.54	77.74
	+LOREAL	<b>82.55</b>	<b>68.60</b>	<b>74.93</b>	<b>77.10</b>	<b>66.18</b>	<b>71.22</b>	<b>98.32</b>	<b>91.70</b>	<b>94.89</b>	<b>95.84</b>	<b>96.42</b>	<b>96.13</b>	<b>78.81</b>	<b>69.52</b>	<b>73.87</b>	<b>97.15</b>	<b>65.62</b>	<b>78.33</b>
	MaPLe	82.19	70.75	76.04	75.41	66.28	70.55	98.05	93.56	95.75	95.25	97.04	96.14	74.71	71.30	72.97	97.34	70.99	82.10
	+LOREAL	<b>83.32</b>	<b>72.60</b>	<b>77.59</b>	<b>76.95</b>	<b>67.98</b>	<b>72.19</b>	<b>98.46</b>	<b>94.76</b>	<b>96.57</b>	<b>94.95</b>	<b>97.48</b>	<b>96.20</b>	<b>76.39</b>	<b>74.53</b>	<b>75.45</b>	<b>97.87</b>	<b>72.44</b>	<b>83.26</b>
	MMA	81.28	75.47	78.27	76.25	68.45	72.14	98.39	92.90	95.57	94.36	97.54	95.92	75.89	72.25	74.03	<b>97.63</b>	71.56	82.59
	+LOREAL	<b>82.61</b>	<b>76.98</b>	<b>79.70</b>	<b>78.32</b>	<b>69.28</b>	<b>73.52</b>	<b>98.44</b>	<b>93.01</b>	<b>95.65</b>	<b>95.52</b>	<b>98.06</b>	<b>96.77</b>	<b>77.41</b>	<b>75.12</b>	<b>76.25</b>	97.24	<b>76.24</b>	<b>85.47</b>
	MMRL	83.32</																	

Table 2. Results on the LR-CE Benchmark over three resolution settings  $\phi \in \{96^2, 144^2, 192^2\}$  w/ or w/o our LOREAL.  $\eta$  is the ratio of patch tokens relative to the  $224^2$  input. All results are averaged over 3 runs. **Bold** marks the best results.

Method	ImageNet (Source)			Caltech101			OxfordPets			StanfordCars			Flowers102			Food101		
	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>
CoOp	22.70	57.97	67.64	47.63	93.73	94.49	48.74	91.33	93.08	13.73	58.23	66.20	24.12	63.87	71.00	28.86	80.16	87.75
+LOREAL	<b>37.65</b>	<b>63.28</b>	<b>70.92</b>	<b>62.41</b>	<b>96.90</b>	<b>95.51</b>	<b>55.13</b>	<b>93.22</b>	<b>93.88</b>	<b>22.40</b>	<b>61.52</b>	<b>69.64</b>	<b>42.06</b>	<b>74.51</b>	<b>78.30</b>	<b>45.23</b>	<b>83.87</b>	<b>88.08</b>
MaPLe	20.28	57.30	68.92	52.82	93.07	78.24	49.69	93.92	95.40	14.81	59.52	64.44	30.94	66.77	72.59	33.09	77.26	88.12
+LOREAL	<b>43.79</b>	<b>64.76</b>	<b>71.60</b>	<b>79.59</b>	<b>94.56</b>	<b>95.38</b>	<b>60.37</b>	<b>94.09</b>	<b>95.33</b>	<b>39.95</b>	<b>67.87</b>	<b>68.25</b>	<b>52.28</b>	<b>72.37</b>	<b>83.17</b>	<b>44.74</b>	<b>87.46</b>	<b>90.05</b>
MMA	24.15	58.33	70.89	47.55	93.69	92.78	58.79	94.36	95.26	10.82	61.60	64.79	35.35	60.78	78.68	35.53	73.81	88.57
+LOREAL	<b>46.53</b>	<b>65.91</b>	<b>73.26</b>	<b>79.38</b>	<b>96.68</b>	<b>94.80</b>	<b>67.58</b>	<b>95.20</b>	<b>95.65</b>	<b>36.67</b>	<b>69.83</b>	<b>71.60</b>	<b>55.84</b>	<b>74.89</b>	<b>79.12</b>	<b>48.00</b>	<b>80.83</b>	<b>90.02</b>
MMRL	24.35	62.06	71.87	69.98	93.44	94.89	60.75	94.92	<b>95.87</b>	12.11	63.36	70.79	42.05	69.36	<b>78.94</b>	37.35	80.25	89.96
+LOREAL	<b>47.50</b>	<b>68.04</b>	<b>73.71</b>	<b>78.86</b>	<b>94.94</b>	<b>96.36</b>	<b>70.10</b>	<b>96.03</b>	95.60	<b>37.58</b>	<b>65.64</b>	<b>73.84</b>	<b>58.28</b>	<b>75.18</b>	78.73	<b>50.65</b>	<b>85.95</b>	<b>90.48</b>
Method	FGVCAircraft			SUN397			DTD			EuroSAT			UCF101			Target Avg		
	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>
CoOp	8.25	13.42	22.02	31.50	62.30	70.87	30.74	41.97	51.08	50.23	63.22	63.95	33.40	62.77	71.61	31.72	63.10	69.20
+LOREAL	<b>13.96</b>	<b>20.59</b>	<b>28.15</b>	<b>46.14</b>	<b>69.35</b>	<b>76.63</b>	<b>54.55</b>	<b>61.80</b>	<b>66.02</b>	<b>69.15</b>	<b>67.15</b>	<b>77.89</b>	<b>46.74</b>	<b>68.89</b>	<b>72.59</b>	<b>45.77</b>	<b>69.78</b>	<b>74.67</b>
MaPLe	10.59	13.72	25.54	44.21	64.05	70.95	33.46	47.28	61.73	54.36	62.46	68.49	36.71	62.13	70.61	36.07	64.02	69.61
+LOREAL	<b>14.26</b>	<b>23.11</b>	<b>29.27</b>	<b>48.02</b>	<b>70.76</b>	<b>78.03</b>	<b>60.58</b>	<b>64.00</b>	<b>64.25</b>	<b>68.49</b>	<b>75.16</b>	<b>88.72</b>	<b>59.38</b>	<b>76.72</b>	<b>79.50</b>	<b>52.77</b>	<b>72.61</b>	<b>77.19</b>
MMA	11.52	16.95	26.51	44.65	71.89	74.41	34.32	46.86	62.69	55.97	63.26	65.05	38.23	70.04	<b>74.04</b>	37.27	65.32	72.28
+LOREAL	<b>15.04</b>	<b>24.84</b>	<b>32.59</b>	<b>50.51</b>	<b>73.15</b>	<b>78.08</b>	<b>64.35</b>	<b>62.17</b>	<b>65.10</b>	<b>75.38</b>	<b>79.49</b>	<b>87.60</b>	<b>53.80</b>	<b>72.85</b>	73.55	<b>54.66</b>	<b>72.99</b>	<b>76.81</b>
MMRL	11.82	23.55	32.54	48.22	71.40	76.99	37.85	54.83	54.09	54.67	72.74	81.23	32.94	73.23	77.72	40.77	69.71	75.30
+LOREAL	<b>17.88</b>	<b>29.81</b>	<b>37.97</b>	<b>52.34</b>	<b>74.10</b>	<b>78.35</b>	<b>62.15</b>	<b>67.49</b>	<b>72.68</b>	<b>70.77</b>	<b>76.69</b>	<b>89.31</b>	<b>52.17</b>	<b>75.34</b>	<b>78.04</b>	<b>55.08</b>	<b>74.12</b>	<b>79.14</b>

Table 3. Results on the LR-DG Benchmark over three resolution settings  $\phi \in \{96^2, 144^2, 192^2\}$  w/ or w/o our LOREAL.  $\eta$  is the ratio of patch tokens relative to the  $224^2$  input. All results are averaged over 3 runs. **Bold** marks the best results.

Method	ImageNet (Source)			ImageNet-V2			ImageNet-S			ImageNet-A			ImageNet-R			Target Avg.		
	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>	96 <sup>2</sup>	144 <sup>2</sup>	192 <sup>2</sup>
CoOp	22.70	57.97	67.64	17.68	53.38	60.12	10.27	36.74	42.48	10.51	22.29	38.71	21.63	63.06	71.32	15.02	43.87	53.16
+LOREAL	<b>37.65</b>	<b>63.28</b>	<b>70.92</b>	<b>31.32</b>	<b>56.55</b>	<b>63.84</b>	<b>18.84</b>	<b>38.47</b>	<b>42.88</b>	<b>17.29</b>	<b>28.51</b>	<b>43.24</b>	<b>25.63</b>	<b>69.74</b>	<b>73.77</b>	<b>23.27</b>	<b>48.32</b>	<b>55.93</b>
MaPLe	20.28	57.30	68.92	12.08	50.61	61.28	9.58	35.66	42.15	7.29	20.84	37.18	13.29	59.13	69.75	10.56	41.56	52.59
+LOREAL	<b>43.79</b>	<b>64.76</b>	<b>71.60</b>	<b>36.83</b>	<b>57.25</b>	<b>64.16</b>	<b>21.02</b>	<b>40.46</b>	<b>42.93</b>	<b>15.91</b>	<b>28.97</b>	<b>39.19</b>	<b>33.68</b>	<b>66.98</b>	<b>72.85</b>	<b>26.86</b>	<b>48.42</b>	<b>54.78</b>
MMA	24.15	58.33	70.89	17.98	53.87	61.81	10.24	37.21	43.46	9.53	23.32	38.25	20.11	60.37	69.21	14.47	43.69	53.18
+LOREAL	<b>46.53</b>	<b>65.91</b>	<b>73.26</b>	<b>35.10</b>	<b>57.00</b>	<b>62.59</b>	<b>19.98</b>	<b>38.31</b>	<b>43.98</b>	<b>14.22</b>	<b>27.69</b>	<b>39.06</b>	<b>32.77</b>	<b>66.41</b>	<b>71.00</b>	<b>25.52</b>	<b>47.35</b>	<b>54.16</b>
MMRL	24.35	62.06	71.87	18.10	55.29	63.61	10.84	39.57	46.19	9.80	28.53	44.92	22.94	65.98	74.75	15.42	47.34	57.37
+LOREAL	<b>47.50</b>	<b>68.04</b>	<b>73.71</b>	<b>42.56</b>	<b>58.65</b>	<b>65.50</b>	<b>21.95</b>	<b>41.59</b>	<b>46.55</b>	<b>12.89</b>	<b>32.55</b>	<b>45.95</b>	<b>45.15</b>	<b>69.32</b>	<b>75.89</b>	<b>30.64</b>	<b>50.53</b>	<b>58.47</b>

LOREAL shows more pronounced effectiveness for multi-modal approaches like MaPLe, MMA and MMRL ( $\phi=96^2$ : an average of +22.81% HM) compared to CoOp ( $\phi=96^2$ : 11.3% HM), possibly attributed to their cross-modal design which facilitates the meta-network’s convergence to local minima during resolution-aware semantic refinement.

**Results on the LR-CE Benchmark.** Table 2 presents the results of models (trained on the standard ImageNet training set) being evaluated on LR test-sets of ImageNet and 10 additional datasets. First, when augmented with LOREAL, all methods achieve average improvements of (21.0%, 6.58%, and 2.54%) on the source ImageNet; Furthermore, LOREAL substantially enhances performance on the diverse set of target datasets, especially under low  $\phi$  settings, underscoring its remarkable effectiveness to cross-dataset tasks.

**Results on the LR-DG Benchmark.** Table 3 summarizes the results on the domain generalization benchmark, where models trained on ImageNet are evaluated on LR test-sets of four ImageNet variants. LOREAL confers consistent improvements to all backbone methods across varied  $\phi$ : Specifically, at  $\phi = 96^2$ , all methods exhibit a sharp average performance gain of 12.71%, further validating the effectiveness of LOREAL for domain generalization.

**Efficiency Study.** Integrating our LOREAL with existing approaches only requires additional meta-nets, resulting in high efficiency. In Table 4, we report the tunable param-

Table 4. Efficiency study. Tra./Infer. time is tested for a sample.

$\phi$	Method	Tunable Para.	Tra. Time	Infer. Time	Infer. Mem.	HM $\uparrow$
96 <sup>2</sup>	MaPLe	3555K	107ms	32ms	612MB	34.85
	+LOREAL	+104.5K	+4ms	+1ms	+33MB	57.25
	MMRL	4992K	113ms	29ms	1236MB	38.90
	+LOREAL	+104.5K	+5ms	+1ms	+33MB	61.71
144 <sup>2</sup>	MaPLe	3555K	126ms	38ms	844	68.80
	+LOREAL	+104.5K	+4ms	+2ms	+33MB	73.72
	MMRL	4992K	145ms	37ms	1440MB	72.43
	+LOREAL	+104.5K	+5ms	+1ms	+33MB	76.56

eters, training, and inference cost of MaPLe and MMRL with/without LOREAL under two  $\phi$ . As observed, the additional computation cost is minimal or even negligible; In contrast, LOREAL brings substantial gains to the HM performance, demonstrating remarkable efficiency.

### 4.3. Ablation Studies and Discussions

**The Proposed Modules.** We ablate each proposed module of LOREAL and report the results in Table 6. “LR→St.” represents using LR image embeddings, to contextualize prompt attributes of the student which accepts standard resolution. Our key findings are: (a): both “LR→St.” and “St.→LR” contribute to the performance, where the former exhibits a marginally superior effect (0.84% HM) compared to the latter. This indicates that achieving image-text alignment for LR inputs in training is more conducive to generalization in LR inference. (b): Removing either LLD or HLD leads to a considerable performance drop across all

Table 5. Ablation studies on  $D_s$  (left) and  $\tau$  (right).

$D_s$	Base	New	HM	$\tau$	Base	New	HM
8	73.19	62.18	66.98	0.5	75.99	66.34	70.58
16	76.07	63.75	69.11	1	76.38	67.06	71.16
32	77.30	67.37	71.73	2	76.55	67.63	71.55
64	76.54	66.96	71.17	4	77.30	67.37	71.73
128	76.35	67.58	71.44	6	76.80	66.85	71.22
256	75.78	67.49	71.14	8	76.10	66.18	70.53
512	75.97	61.20	67.53	12	75.99	66.11	70.45
1024	75.21	64.53	69.20	16	76.07	66.03	70.44
2048	73.88	62.74	67.60	20	75.95	65.72	70.21

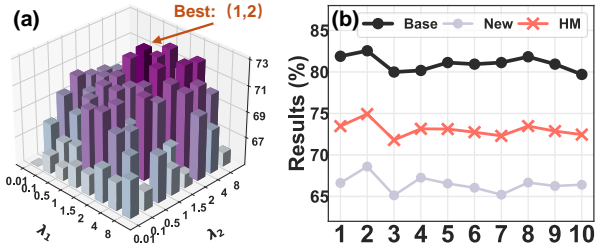


Figure 5. Ablation studies on (a)  $\lambda_1$ ,  $\lambda_2$  and (b)  $M$ .

415 metrics, confirming that both distillation losses are essential  
 416 for bridging the semantic gap across resolutions. Between  
 417 them, LLD contributes more significantly than HLD, likely  
 418 because its direct attribute-level alignment is more readily  
 419 optimized compared to the logit-level alignment of HLD.

420 **Balance Coefficients  $\lambda_1$  and  $\lambda_2$ .** We perform a grid-search  
 421 over these parameters in Figure 5 (a). Although increasing  
 422 either hyperparameter generally improves HM performance  
 423 initially, a dramatic degradation occurs when they continually  
 424 grow. This is likely because an excessively large distillation  
 425 loss overwhelms the task objective, thereby hindering  
 426 overall convergence. The optimal value for  $\lambda_2$  is higher than  
 427 that for  $\lambda_1$ , further indicating the effectiveness of LLD.

428 **Length of Attribute Prompts  $M$ .** We vary  $M$  from 1 to 10,  
 429 with results presented in Figure 5(b). We find that increasing  
 430  $M$  generally leads to a slight performance degradation  
 431 ( $\sim 2\%$ ), potentially because excessive tokens introduce additional  
 432 meta-networks and increase optimization difficulty.  
 433 We fix  $M$  to 2 as indicated by the results.

434 **Intermediate Dimension  $D_s$  of Meta-nets.** The results are  
 435 shown in Table 5 (left). Results suggest that while increasing  
 436  $D_s$  brings significant improvement initially, this trend  
 437 rapidly reverses as  $D_s$  undergoes excessive expansion. This  
 438 may be attributed to the expanded optimization difficulties  
 439 of meta-nets. According to the results, the optimal  $D_s$  is 32.

440 **The Temperature  $\tau$ .** We vary  $\tau$  from 0.5 to 20 and report  
 441 the results in Table 5 (right). The results indicate that  
 442 performance degrades with either excessively small or large  
 443 values, likely because such values over-sharpen or over-  
 444 smooth the output distributions, impairing the knowledge  
 445 transfer. We keep  $\tau = 4$  as suggested by the results.

446 **Further Analysis.** In this part, we discover the semantic  
 447 meanings of attribute contents we learned for interpretability  
 448 analysis. Concretely, we first construct an attribute vo-

Table 6. Ablation studies on proposed modules.

LR $\rightarrow$ St.	St. $\rightarrow$ LR	LLD	HLD	Base	New	HM
✓				75.45	63.42	68.91
	✓			74.78	62.46	68.07
✓	✓			75.70	64.72	69.78
✓	✓	✓		76.10	66.35	70.89
✓	✓	✓	✓	77.30	67.37	71.73

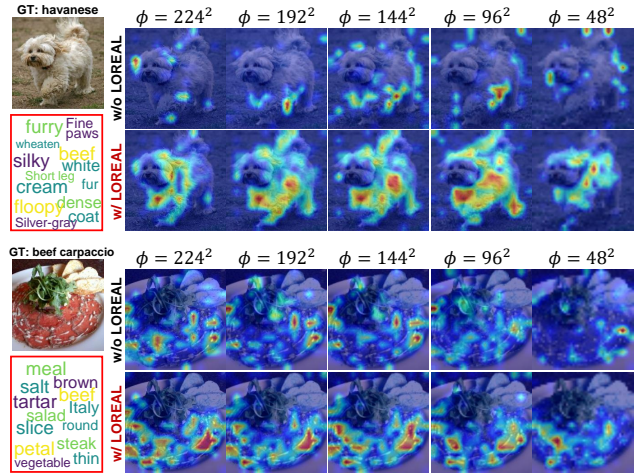


Figure 6. Visualizations of the learned attribute tokens (Lower-Left), and attention heatmaps (Right) across varied resolution settings with or without our LOREAL on zero-shot CLIP. Note that here attention maps of different  $\phi$  are resized for better illustration.

449 cabulary (Details are in **Sup. Mat. D**) by querying an LLM.  
 450 We then compute the similarity between each vocabulary  
 451 word’s embedding and the learned output text embeddings  
 452 from the student. In Figure 6 (lower-left), we showcase the  
 453 word cloud of several most similar attributes in the vocabu-  
 454 lary. Our results reveal that the learned attribute contents ac-  
 455 curately correspond to visual semantics, demonstrating the  
 456 validity of our approach. Furthermore, we plot the atten-  
 457 tion heatmaps with or without LOREAL in Figure 6 (right).  
 458 With LOREAL, the model maintains focus on semantically  
 459 central regions even at very low resolutions. Furthermore,  
 460 it shifts attention toward holistic, resolution-robust features  
 461 (e.g., the havanese’s entirety) rather than details. This  
 462 demonstrates that our approach guides the model to priori-  
 463 tize invariant semantic concepts.

## 5. Conclusion

464 Edge-deployed VLMs are often constrained by computa-  
 465 tional resources, necessitating inference on Low-Resolution  
 466 (LR) inputs and presenting a significant challenge to model  
 467 robustness. To this end, this paper introduces LOREAL,  
 468 a self-distillation framework that endows models with ro-  
 469 bustness to LR via resolution-invariant attributes. LOREAL  
 470 comprises LLM-based attribute generation and attribute-  
 471 driven self-distillation. Experiments under various LR set-  
 472 tings demonstrate the superiority of our model.  
 473

474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529

## References

- [1] Lucas Beyer, Pavel Izmailov, Alexander Kolesnikov, Mathilde Caron, Simon Kornblith, Xiaohua Zhai, Matthias Minderer, Michael Tschannen, Ibrahim Alabdulmohsin, and Filip Pavetic. Flexivit: One model for all patch sizes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14496–14506, 2023. 3
- [2] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative components with random forests. In *Computer vision—ECCV 2014: 13th European conference, zurich, Switzerland, September 6–12, 2014, proceedings, part VI 13*, pages 446–461. Springer, 2014. 5
- [3] Haoxing Chen, Zizheng Huang, Yan Hong, Yanshuo Wang, Zhongcai Lyu, Zhuoer Xu, Jun Lan, and Zhangxuan Gu. Efficient transfer learning for video-language foundation models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 29129–29138, 2025. 1
- [4] Liqun Chen, Dong Wang, Zhe Gan, Jingjing Liu, Ricardo Henao, and Lawrence Carin. Wasserstein contrastive representation distillation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16296–16305, 2021. 3
- [5] Pengguang Chen, Shu Liu, Hengshuang Zhao, and Jiaya Jia. Distilling knowledge via knowledge review. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5008–5017, 2021. 3
- [6] Shiming Chen, Bowen Duan, Salman Khan, and Fahad Shahbaz Khan. Interpretable zero-shot learning with locally-aligned vision-language model. *arXiv preprint arXiv:2506.23822*, 2025. 4
- [7] Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3606–3613, 2014. 5
- [8] Mostafa Dehghani, Basil Mustafa, Josip Djolonga, Jonathan Heek, Matthias Minderer, Mathilde Caron, Andreas Steiner, Joan Puigcerver, Robert Geirhos, Ibrahim M Alabdulmohsin, et al. Patch n’pack: Navit, a vision transformer for any aspect ratio and resolution. *Advances in Neural Information Processing Systems*, 36:2252–2274, 2023. 3
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 5
- [10] Tong Ding, Wanhua Li, Zhongqi Miao, and Hanspeter Pfister. Tree of attributes prompt learning for vision-language models. *arXiv preprint arXiv:2410.11201*, 2024. 3
- [11] Yingjun Du, Wenfang Sun, and Cees Snoek. Ipo: Interpretable prompt optimization for vision-language models. *Advances in Neural Information Processing Systems*, 37: 126725–126766, 2024. 3
- [12] Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In *2004 conference on computer vision and pattern recognition workshop*, pages 178–178. IEEE, 2004. 5
- [13] Shiming Ge, Shengwei Zhao, Chenyu Li, and Jia Li. Low-resolution face recognition in the wild via selective knowledge distillation. *IEEE Transactions on Image Processing*, 28(4):2051–2062, 2018. 3
- [14] Guangyu Guo, Dingwen Zhang, Longfei Han, Nian Liu, Ming-Ming Cheng, and Junwei Han. Pixel distillation: Cost-flexible distillation across image sizes and heterogeneous networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12):9536–9550, 2024. 2, 3
- [15] Yuncheng Guo and Xiaodong Gu. Mmrl: Multi-modal representation learning for vision-language models. *arXiv preprint arXiv:2503.08497*, 2025. 1, 2, 5
- [16] Yuncheng Guo and Xiaodong Gu. Mmrl++: Parameter-efficient and interaction-aware representation learning for vision-language models. *arXiv preprint arXiv:2505.10088*, 2025. 1
- [17] Zhiwei Hao, Jianyuan Guo, Kai Han, Yehui Tang, Han Hu, Yunhe Wang, and Chang Xu. One-for-all: Bridge the gap between heterogeneous architectures in knowledge distillation. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [18] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019. 5
- [19] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8340–8349, 2021. 5
- [20] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15262–15271, 2021. 5
- [21] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 3
- [22] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 5
- [23] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hananeh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, 2021. 4
- [24] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, pages 4904–4916. PMLR, 2021. 1
- [25] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *European conference on computer vision*, pages 709–727. Springer, 2022. 1

- 588 [26] Ying Jin, Jiaqi Wang, and Dahua Lin. Multi-level logit distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24276–24285, 2023. 3
- 589
- 590
- 591
- 592 [27] Muhammad Uzair Khattak, Hanoona Rasheed, Muhammad Maaz, Salman Khan, and Fahad Shahbaz Khan. Maple: Multi-modal prompt learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19113–19122, 2023. 1, 2, 5
- 593
- 594
- 595
- 596
- 597 [28] Muhammad Uzair Khattak, Syed Talal Wasim, Muzammal Naseer, Salman Khan, Ming-Hsuan Yang, and Fahad Shahbaz Khan. Self-regulating prompts: Foundational model adaptation without forgetting. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 15190–15200, 2023. 2
- 598
- 599
- 600
- 601
- 602
- 603 [29] Sanghwan Kim, Rui Xiao, Mariana-Iuliana Georgescu, Stephan Alaniz, and Zeynep Akata. Cosmos: Cross-modality self-distillation for vision language pre-training. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 14690–14700, 2025. 3
- 604
- 605
- 606
- 607
- 608 [30] Youmin Kim, Jinbae Park, YounHo Jang, Muhammad Ali, Tae-Hyun Oh, and Sung-Ho Bae. Distilling global and local logits with densely connected relations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6290–6300, 2021. 3
- 609
- 610
- 611
- 612
- 613 [31] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 554–561, 2013. 4, 5
- 614
- 615
- 616
- 617 [32] Jungsoo Lee, Debasmit Das, Munawar Hayat, Sungha Choi, Kyuwoong Hwang, and Fatih Porikli. Customkd: Customizing large vision foundation for edge model improvement via knowledge distillation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 25176–25186, 2025. 3
- 618
- 619
- 620
- 621
- 622
- 623 [33] Kenton Lee, Mandar Joshi, Iulia Raluca Turc, Hexiang Hu, Fangyu Liu, Julian Martin Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. Pix2struct: Screenshot parsing as pretraining for visual language understanding. In *International Conference on Machine Learning*, pages 18893–18912. PMLR, 2023. 3
- 624
- 625
- 626
- 627
- 628
- 629 [34] Haoyang Li, Liang Wang, Chao Wang, Jing Jiang, Yan Peng, and Guodong Long. Dpc: Dual-prompt collaboration for tuning vision-language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 25623–25632, 2025. 1, 2
- 630
- 631
- 632
- 633
- 634 [35] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pages 12888–12900. PMLR, 2022. 1
- 635
- 636
- 637
- 638
- 639 [36] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 1
- 640
- 641
- 642
- 643
- 644 [37] Runze Li, Dahun Kim, Bir Bhanu, and Weicheng Kuo. Reclip: Resource-efficient clip by training with small images. *arXiv preprint arXiv:2304.06028*, 2023. 2
- 645
- 646
- 647 [38] Zheng Li, Xiang Li, Xinyi Fu, Xin Zhang, Weiqiang Wang, Shuo Chen, and Jian Yang. Promptkd: Unsupervised prompt distillation for vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26617–26626, 2024. 2, 3
- 648
- 649
- 650
- 651
- 652 [39] Zheng Li, Yibing Song, Penghai Zhao, Ming-Ming Cheng, Xiang Li, and Jian Yang. Atprompt: Textual prompt learning with embedded attributes. *arXiv preprint arXiv:2412.09442*, 2024. 2, 4
- 653
- 654
- 655
- 656 [40] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023. 1
- 657
- 658
- 659 [41] Wenzhuo Liu, Fei Zhu, Shijie Ma, and Cheng-Lin Liu. Mspe: multi-scale patch embedding prompts vision transformers to any resolution. *Advances in Neural Information Processing Systems*, 37:29191–29212, 2024. 3
- 660
- 661
- 662
- 663 [42] Wenzhuo Liu, Fei Zhu, Longhui Wei, and Qi Tian. C-clip: Multimodal continual learning for vision-language model. In *The Thirteenth International Conference on Learning Representations*, 2025. 1
- 664
- 665
- 666
- 667 [43] Yuning Lu, Jianzhuang Liu, Yonggang Zhang, Yajing Liu, and Xinmei Tian. Prompt distribution learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5206–5215, 2022. 1
- 668
- 669
- 670
- 671 [44] Subhansu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013. 5
- 672
- 673
- 674
- 675 [45] Roy Miles, Ismail Elezi, and Jiankang Deng. Vkd: Improving knowledge distillation using orthogonal projections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15720–15730, 2024. 3
- 676
- 677
- 678
- 679 [46] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian conference on computer vision, graphics & image processing*, pages 722–729. IEEE, 2008. 5
- 680
- 681
- 682
- 683 [47] Shuaicheng Niu, Jiayang Wu, Yifan Zhang, Yaofu Chen, Shijian Zheng, Peilin Zhao, and Minghui Tan. Efficient test-time model adaptation without forgetting. In *International conference on machine learning*, pages 16888–16905. PMLR, 2022. 2
- 684
- 685
- 686
- 687
- 688 [48] Bikang Pan, Qun Li, Xiaoying Tang, Wei Huang, Zhen Fang, Feng Liu, Jingya Wang, Jingyi Yu, and Ye Shi. Nlprompt: Noise-label prompt learning for vision-language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 19963–19973, 2025. 1
- 689
- 690
- 691
- 692
- 693 [49] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3967–3976, 2019. 3
- 694
- 695
- 696
- 697 [50] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3498–3505. IEEE, 2012. 4, 5
- 698
- 699
- 700



815 models. In *Proceedings of the IEEE/CVF Conference on*  
816 *Computer Vision and Pattern Recognition*, pages 23826–  
817 23837, 2024. 1, 2

818 [76] Zhendong Yang, Zhe Li, Mingqi Shao, Dachuan Shi, Ze-  
819 huan Yuan, and Chun Yuan. Masked generative distillation.  
820 In *European conference on computer vision*, pages 53–69.  
821 Springer, 2022. 3

822 [77] Hantao Yao, Rui Zhang, and Changsheng Xu. Tc-  
823 based class-aware prompt tuning for visual-language model.  
824 In *Proceedings of the IEEE/CVF Conference on Computer*  
825 *Vision and Pattern Recognition*, pages 23438–23448, 2024.  
826 1, 2

827 [78] Renrui Zhang, Rongyao Fang, Wei Zhang, Peng Gao,  
828 Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li.  
829 Tip-adapter: Training-free clip-adapter for better vision-  
830 language modeling. *arXiv preprint arXiv:2111.03930*, 2021.  
831 1

832 [79] Yi Zhang, Ke Yu, Siqi Wu, and Zhihai He. Conceptual code-  
833 book learning for vision-language models. In *European Con-*  
834 *ference on Computer Vision*, pages 235–251. Springer, 2024.  
835 4

836 [80] Borui Zhao, Quan Cui, Renjie Song, Yiyu Qiu, and Jiajun  
837 Liang. Decoupled knowledge distillation. In *Proceedings of*  
838 *the IEEE/CVF Conference on computer vision and pattern*  
839 *recognition*, pages 11953–11962, 2022. 3

840 [81] Li Zhong, Ahmed Ghazal, Jun-Jun Wan, Frederik Zilly,  
841 Patrick Mackens, Joachim Vollrath, and Bogdan Coseriu.  
842 Clip4retrofit: Enabling real-time image labeling on edge de-  
843 vices via cross-architecture clip distillation. In *Proceedings*  
844 *of the Computer Vision and Pattern Recognition Conference*,  
845 pages 3829–3837, 2025. 2

846 [82] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei  
847 Liu. Conditional prompt learning for vision-language mod-  
848 els. In *Proceedings of the IEEE/CVF conference on com-*  
849 *puter vision and pattern recognition*, pages 16816–16825,  
850 2022. 1, 2, 4

851 [83] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei  
852 Liu. Learning to prompt for vision-language models. *In-*  
853 *ternational Journal of Computer Vision*, 130(9):2337–2348,  
854 2022. 1, 4

855 [84] Beier Zhu, Yulei Niu, Yucheng Han, Yue Wu, and Hanwang  
856 Zhang. Prompt-aligned gradient for prompt tuning. In *Pro-*  
857 *ceedings of the IEEE/CVF international conference on com-*  
858 *puter vision*, pages 15659–15669, 2023. 2

859 [85] Wilman WW Zou and Pong C Yuen. Very low resolution  
860 face recognition problem. *IEEE Transactions on image pro-*  
861 *cessing*, 21(1):327–340, 2011. 3