

## Appendix

In this appendix, we present the following details.

- List of notations used in this paper and their descriptions are in § A.
- Overall algorithm of SAP is presented in § B.
- Implementation details are in § C.
- Expanded dataset-wise tables, and additional experiments are presented in § D.
- Examples of class descriptions generated using GPT-3.5 are presented in § E.
- Limitations and Broader Impact in § F.

### A Summary of Notations and Terminology

We use  $\cdot$  (*dot*) to represent various types of multiplication operations – matrix multiplication, matrix-vector or vector-matrix product, and vector dot-product. Detailed descriptions of notations are presented in Tab. 10.

Notation	Description	Dimension
$\theta$	Image Encoder	
$\phi$	Text Encoder	
$\mathcal{Y}$	Classification label space	
$\rho$	Set of all learnable text and visual prompts	
$B$	Batch size	
$N$	Size of the set of descriptions	
$n$	Number of the learnable prompt tokens	
$d$	Dimension of the multimodal space	
$A_y$	LLM generated descriptions for class $y$	
$A$	Union of all descriptions of the classification label space	
$\phi(A)$	Class descriptions features	$\mathbb{R}^{N \times d}$
$\phi(y; A_y)$	Description-guided text features of class $y$	$\mathbb{R}^{N \times d}$
$\theta(x)$	Global image feature	$\mathbb{R}^d$
$\theta^l(x)$	Local image feature	$\mathbb{R}^{M \times d}$
$\theta^{desc}(x)$	Description-guided image features	$\mathbb{R}^{N \times d}$
$\bar{\theta}^{desc}(x)$	Mean Description-guided image features	$\mathbb{R}^d$
$\hat{\theta}(x)$	Fused image features	$\mathbb{R}^d$
$\theta_p(x)$	Prompted Global image feature	$\mathbb{R}^d$
$\theta_p^l(x)$	Prompted Local image feature	$\mathbb{R}^{M \times d}$
$\theta_p^{desc}(x)$	Prompted Description-guided image features	$\mathbb{R}^{N \times d}$
$\mathbf{r}$	Description relevance score for an image	$\mathbb{R}^N$
$\alpha$	average specificity for all descriptions	$\mathbb{R}$

Table 10: Notations used in this paper and their descriptions.

### B SAP: Algorithm

Algorithm 1 outlines the SAP methodology. The algorithm is summarized as follows: In a given dataset, descriptions for each class are acquired by querying the LLM (L1 - L4). Class description features are then derived by passing the descriptions through  $\phi$  (L5). Unprompted and prompted image features are obtained by processing images through  $\theta$  (L7-L8). The description-guided image features are obtained via a parameter-free cross-attention between local features and description features (L9). The local image features are a weighted average of the description-guided features based on the relevance of each description to the

image (L10 - L11). Finally, the mean description-guided image features and global image features are fused to create the fusion image feature (L12). Unprompted and prompted description-guided text features are obtained by passing the description-guided text templates through  $\phi$  (L13-L14).  $L_{ce}$ ,  $L_{steer}^v$ , and  $L_{steer}^t$  loss functions are employed to train the prompts.

---

### Algorithm 1 SAP Algorithm

---

**Require:** Dataset  $D = \{\mathbf{x}_i, y_i\}_{i=1}^B$ ; Classification label space:  $\mathcal{Y}$ ; Vision and Language encoders:  $(\theta, \phi)$ ; LLM: ChatGPT-3.5 model; Hyperparameters: coefficients  $\lambda_1, \lambda_2$ , scaling parameter  $s$ , learning rate  $\delta$ ; Learnable Prompts:  $\rho = \{\rho_t, \rho_v\}$

**Ensure:** Trained parameters  $\hat{\rho}$

```

1: /* Get descriptions for each class by querying LLM */
2: for all  $y \in \mathcal{Y}$  do
3:    $A_y = \text{LLM}(\text{Visual features for distinguishing } y \text{ in a photo?})$ 
4: end for
5:  $A = \bigcup_{y \in \mathcal{Y}} A_y$ 
6:  $\phi(A)$  /* Get class description features */
7: for all epochs do
8:   /* Get unprompted and prompted image features for every image  $\mathbf{x}$  in the batch */
9:    $\theta(\mathbf{x}), \_ = \theta(\mathbf{x})$ 
10:   $\theta_p(\mathbf{x}), \theta_p^t(\mathbf{x}) = \theta(\mathbf{x}; \rho_v)$ 
11:  /* Get description-guided image features using parameter-free cross-attention */
12:   $\theta^{desc}(\mathbf{x}) = \text{Cross\_Attention}(Q = \phi(A), K = \theta^t(x), V = \theta^t(\mathbf{x}))$ 
13:  /* Get mean description-guided image feature using relevance score */
14:   $\mathbf{r} = \text{softmax}(\phi(A) \cdot \theta(\mathbf{x}))$ 
15:   $\bar{\theta}^{desc}(\mathbf{x}) = \theta^{desc}(\mathbf{x})^\top \cdot \mathbf{r}$ 
16:  /* Get fused image feature by fusing global and local feature using description specificity ( $\alpha$ ) */
17:   $\hat{\theta}(\mathbf{x}) = (1 - \alpha) \cdot \theta(\mathbf{x}) + \alpha \cdot \bar{\theta}^{desc}(\mathbf{x})$ 
18:  /* Get unprompted and prompted description guided text features for every class  $y$  */
19:   $\phi(y, A_y) = \phi(y, A_y)$ 
20:   $\phi_p(y, A_y) = \phi(y, A_y; \rho_t)$ 
21:  /* Similarity between an image and a class is the aggregate of similarities over pertinent descriptions of a class */
22:   $\xi(\hat{\theta}_p(\mathbf{x}), \phi_p(y; A_y)) = \frac{1}{|A_y|} \sum_{a \in A_y} \text{sim}(\hat{\theta}_p(\mathbf{x}), \phi_p(y; a))$ 
23:
24:   $L_{ce}(\rho) = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(\xi(\hat{\theta}_p(\mathbf{x}_i), \phi_p(y_i; A_{y_i}))/\tau)}{\sum_{y \in \mathcal{Y}} \exp(\xi(\hat{\theta}_p(\mathbf{x}_i), \phi_p(y; A_y))/\tau)}$ 
25:  /* Compute Steering Losses */
26:   $L_{steer}^v(\rho) = \frac{1}{B} \sum_{i=1}^B \|\theta_p(\mathbf{x}_i) - \theta(\mathbf{x}_i)\|_1$ 
27:   $L_{steer}^t(\rho) = \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \|\phi_p(y; A_y) - \phi(y; A_y)\|_1$ 
28:  /* Perform gradient descent on the total loss */
29:   $\mathcal{L}(\rho) = L_{ce}(\rho) + \lambda_1 L_{steer}^v(\rho) + \lambda_2 L_{steer}^t(\rho)$ 
30:   $\hat{\rho} = \rho - \delta \nabla \mathcal{L}(\rho)$ 
31: end for
32: return  $\hat{\rho}$ 

```

---

## C Implementation Details

**Training Details.** We use the ViT-B/16 (Dosovitskiy et al., 2021)-based CLIP model as our backbone. For the GZS and B2N benchmarks, we fine-tune the model on  $K = 16$  shot training data from the base classes. Prompts are learned in the first three layers for the Cross-dataset benchmark and the first nine layers for the remaining two benchmarks. We introduce a  $d$ -dimensional bias as the sole additional parameter compared to (Khattak et al., 2023). The text prompts in the initial layer are initialized with the word embeddings of ‘a photo of a’, and the rest are randomly initialized from a normal distribution, similar to (Khattak et al., 2023). Our models are trained on a single Tesla V100 GPU with Nvidia driver version 470.199.02. We train for 20 epochs, with a batch size of 4 images,  $\lambda_1 = 10$  and  $\lambda_2 = 25$ . The hyperparameter setup is common across all datasets. We use the SGD optimizer with a momentum of 0.9, a learning rate of 0.0025, and weight decay  $5e - 4$ . A cosine learning rate scheduler is applied with a warmup epoch of 1. We do not tune the temperature, and leave it at the default value of 100, also used by CLIP and PSRC. Image pre-processing involves random crops, random horizontal and vertical flips, and normalization using mean values of [0.48, 0.46, 0.41] and standard deviation values of [0.27, 0.26, 0.27]. All baselines utilize publicly available codes and models. All results are averages over three seeds. We use PyTorch 1.12, CUDA 11.3, and build on the Dassel code repository: <https://github.com/KaiyangZhou/Dassel.pytorch>.

Our code is available at <https://github.com/HariChandana1102/Semantic-Alignment-for-Prompt-Tuning-in-Vision-Language-Models>

## D Expanded Tables and Additional Results

**Using Random Text in place of Class Descriptions.** To study the usefulness of valid descriptions, we replace the descriptions for each class by randomly generated text in Tab. 11. Examples of random descriptions are “Raindrops pattered softly against the roof”, “A solitary figure walked down the empty street”. We observe that descriptions matter for unusual datasets having texture-based images, satellite images, aircraft images and action recognition images. The average HM using random text across 11 datasets on B2N benchmark is **78.27%**, while SAP reports an average HM of **80.94%**. A drop of **2.67%** is noted.

	UCF101	EuroSAT	DTD	OxfordPets	StanfordCars	Flowers102	Food101	FGVCAircraft	SUN397	Caltech101	ImageNet	Average
Base	86.27	95.83	83.1	95.07	78.2	97.5	90.13	41.37	81.87	98.07	76.7	84.01
Novel	76.37	69.23	54.1	95.33	72.33	75.53	89.9	34.8	76.63	94.1	67.7	73.27
HM	81.02	80.39	65.54	95.2	75.15	85.12	90.01	37.8	79.16	96.04	72.17	78.27

Table 11: B2N benchmark results using random text in place of class descriptions. The results show that using irrelevant descriptions hurts model performance.

**Using Class Descriptions of Only Ground Truth Classes** Using class descriptions of the ground-truth class makes sense during training but may lead to noisy local features at inference. Our intention of using class descriptions of all *training classes*, is to construct a generalizable local view of the image, rather than a biased one. Due to the unbiased nature of the feature, it can help with tasks like Classification-without-Classnames. Tab. 12. shows the impact of using just the ground-truth class descriptions during training on three benchmarks. We do not change any hyperparameters. These results corroborate our perspective.

B2N	Base	Novel	HM	GZS	Base	Novel	HM	CwC	Base	Novel	HM
all descriptions (Ours)	84.68	77.51	80.94	all descriptions (Ours)	79.46	69.75	74.29	all descriptions (Ours)	43.30	45.60	44.40
gnd truth descriptions	84.58	76.93	80.58	gnd truth descriptions	79.27	68.96	73.76	gnd truth descriptions	41.76	43.45	42.59

Table 12: Comparison with ground truth class descriptions for B2N, GZS and CwC benchmarks.

**Using class descriptions from other LLMs.** We generate class descriptions from two other LLMs - OpenAI’s GPT4o-mini [OpenAI \(2024\)](#) and Anthropic’s Claude Haiku [Anthropic \(2024\)](#). Both LLMs considered are fast and cheap – for instance generating class descriptions for all classes of all 11 datasets from Claude Haiku takes 40 mins and costs 0.5\$. The results are presented in the Tab. 13. for both B2N and GZS benchmarks:

LLM	Base	Novel	HM	LLM	Base	Novel	HM
GPT-3.5	84.68	77.51	80.94	GPT-3.5	79.47	69.75	74.29
Claude Haiku	84.64	77.05	80.67	Claude Haiku	79.31	69.14	73.88
GPT4o-mini	84.74	77.16	80.77	GPT4o-mini	79.54	69.53	74.2

Table 13: Comparison with class descriptions generated from other LLMs on Base-to-Novel benchmark on the left, and Generalized Zero-Shot benchmark on the right.

The results indicate that we get similar results across varying quality of outputs from different LLMs. We believe that in the future obtaining text semantics is going to be cheaper and easier, which necessitates algorithms that can make use of such cheap semantic information.

**Few-shot Setting.** Our main objective is to train prompts that can generalize effectively to novel classes and datasets. As such, we present results primarily on settings that test generalizability, such as the GZS benchmark, Base-to-Novel benchmark, and the Classification without Class-names benchmark. For completeness, we present results in a few-shot classification setting, where limited training samples are provided for all

classes. Note that there are no novel classes in this setting. We showcase outcomes for  $K = 1, 2, 4, 8,$  and  $16$  shots. As shown in Fig. 6, on average, across 11 datasets, we perform competitively against the best baseline PSRC.

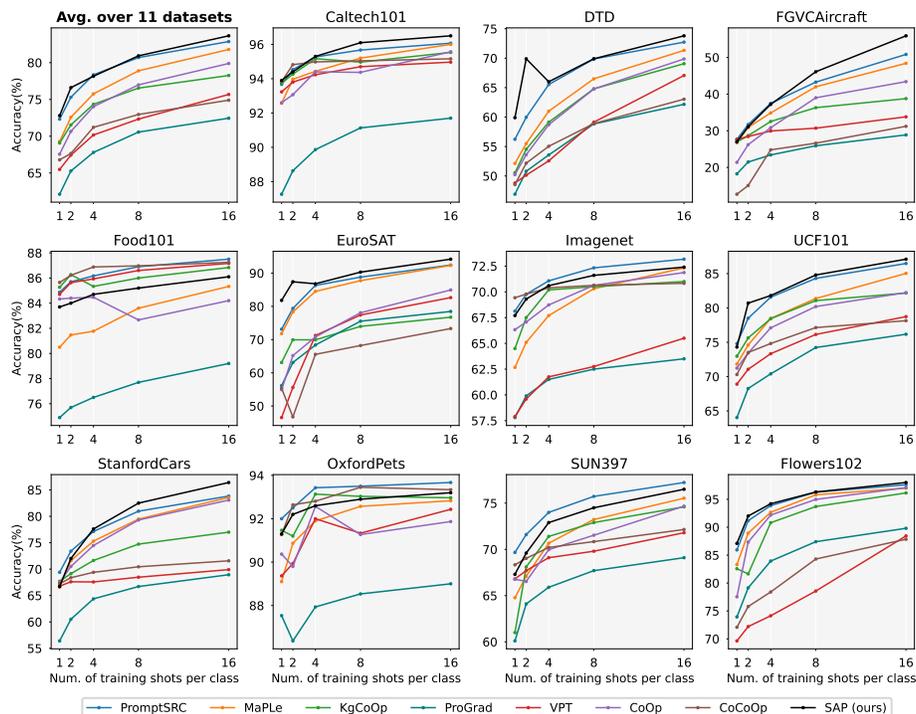


Figure 6: Performance of SAP in the few-shot setting. Our method achieves competitive performance compared to all baselines on average across 11 datasets.

**Domain Generalization.** We show results on Domain Generalization in Tab. 14. We train on  $K = 16$  shot training data from base classes of source dataset ImageNet and evaluation on ImageNetV2, ImageNet-A, ImageNet-Setch, and ImageNet-R target datasets. SAP outperforms two strong baselines PSRC and MaPLe.

	Source		Target				Avg
	ImageNet	-V2	-A	-S	-R		
MaPLe	77.10	71.00	53.70	50.00	77.70	63.10	
PSRC	76.30	71.00	54.10	50.00	77.80	63.22	
SAP	76.40	71.10	55.70	49.80	77.50	<b>63.52</b>	

Table 14: DG benchmark. SAP outperforms baselines on avg.

**ResNet-50 Backbone as Image Encoder.** Here we show the GZS and B2N performance of SAP using the ResNet-50 CLIP model as a backbone. We compare against five baselines which also use the ResNet-50 backbone and present our results in Tab. 15. For all methods including ours, we train the models without tuning any hyperparameters such as prompt-depth, regularization weight, learning rate etc. and use the same values as those of ViT-B/16 CLIP backbone. We observe that PSRC performs particularly poorly with a ResNet backbone. Although we use similar hyperparameters as PSRC, SAP shows good results, indicating that class descriptions help greatly in this setting. We show a gain of  $+0.98\%$  on average gHM for GZS, and  $+2.32\%$  on average HM in the B2N setting.

**Prompt Depth.** Tab. 16 shows the average HM for the B2N benchmark across nine datasets, excluding SUN397 and ImageNet. As seen from the table, adding prompts till depth 9 for image and text encoders is ideal for SAP performance and is used for B2N, GZS and CwC benchmarks.

Depth	1	3	5	7	9	11
HM	76.84	79.35	79.25	80.85	<b>81.76</b>	80.68

Table 16: Prompt depth analysis

**Class Activation Maps (CAMs).** We show additional CAMs for the ResNet-50(He et al., 2015) backbone encoder to visualize image regions that most correlate to a given description. Fig. 7 shows the GradCAM (Sel-

Dataset		CLIP	CoOp	KgCoOp	ProGrad	PSRC	SAP (Ours)
<b>Generalized Zero-Shot Learning Benchmark</b>							
<b>Average on 11 datasets</b>	gBase	57.01	68.65	69.25	<u>69.89</u>	47.41	<b>71.52 (+1.63)</b>
	gNovel	<b>60.73</b>	50.35	59.08	52.26	29.16	<u>59.13 (-1.60)</u>
	gHM	58.81	58.1	<u>63.76</u>	59.81	36.12	<b>64.74 (+0.98)</b>
<b>Base-to-Novel Generalization Benchmark</b>							
<b>Average on 11 datasets</b>	Base	65.27	77.24	75.51	<u>77.98</u>	55.13	<b>78.49 (+0.51)</b>
	Novel	68.14	57.40	<u>67.53</u>	63.41	38.72	<b>69.32 (+1.79)</b>
	HM	66.68	65.86	<u>71.30</u>	69.94	45.49	<b>73.62 (+2.32)</b>

Table 15: Results on GZS and B2N settings using a ResNet-50 backbone. On average, SAP outperforms all the baselines.

varaju et al., 2017) visualizations for base classes “*Floor gymnastics*”, “*Hammering*”, “*Cape Flower*” and “*Highway*”. SAP effectively localizes the text semantics in the image compared to baselines. In Tab. 17, we show quantitative results using an occlusion metric to measure the localization capabilities of our learned prompts. Given a description, we mask out parts of the image which are most activated w.r.t. the description. The occluded image is then classified by the pre-trained CLIP model. A CAM localizes the description well if occluding image regions with the highest activations leads to a large drop in accuracy.

Method	Archery	Baby Crawling	Band Marching	Apply Eye Makeup	Apply Lipstick	Biking	Body Weight Squats
CoOp	57.39	64.42	61.99	75.00	78.66	55.15	53.97
PSRC	47.87	53.69	54.29	50.00	69.33	50.35	50.72
Ours	<b>44.34</b>	<b>49.66</b>	<b>51.58</b>	<b>40.90</b>	<b>62.66</b>	<b>47.96</b>	<b>48.73</b>
	<b>707-320</b>	<b>747-200</b>	<b>737-200</b>	<b>727-200</b>	<b>C-130</b>	<b>CRJ-200</b>	<b>Boeing-717</b>
CoOp	15.21	11.82	23.47	6.13	75.81	38.22	20.63
PSRC	6.14	8.84	21.42	3.06	75.86	32.45	23.58
Ours	<b>3.00</b>	<b>5.92</b>	<b>15.30</b>	<b>0.00</b>	<b>60.61</b>	<b>26.58</b>	<b>14.72</b>

Table 17: Occlusion benchmark (lower number is better): Images are masked at regions of highest activation relevant to a given class description, as identified by prompted image and text encoders, and then evaluated using the pre-trained CLIP model. The lower the accuracy, the better are the localizations. We show results for a few specific classes from the UCF101 dataset (top) and FGVC-Aircraft dataset (bottom). For example, for the class ‘*body weight squats*’, we use the description ‘*person bending knees and hips*’.

For instance, for the text phrase ‘*a photo of a 737-200, which has two engines on the wings*’ we find that masking out important regions given by our prompted image encoder leads to an accuracy of 15.30%. This drop is higher than that of PSRC, whose accuracy drops only to 21.42%. This suggests that regions which are deemed important by SAP are highly correlated to the text phrase. Our parameter-free cross-attention module helps us learn prompts that focus on part-level image information.

**Expanded Dataset-wise Tables.** We present the elaborate tables dataset-wise for the Generalized Zero-Shot setting in Tab. 18 and Base-to-Novel generalization setting in Tab. 21. SAP outperforms the best-performing baseline, PSRC, in 7 of the 11 considered datasets. We perform very well in challenging datasets such as EuroSAT, DTD, and UCF-101. We present dataset-wise results for the Classification without Class-names benchmark in Tab. 19. Tab. 20 has the dataset-wise results for the Cross-Dataset generalization benchmark. In Tab. 15 we show average results on the GZS benchmark and the Base-to-Novel benchmark for the ResNet-50 backbone Image Encoder. We also present detailed, dataset-wise results for the same in Tab. 22.

## E Generation of Class Descriptions

Tab. 23 shows class names sampled from different datasets and their respective descriptions retrieved using GPT-3.5 (Hagendorff et al., 2022). We use the query – “What are useful visual features for distinguishing a [classname] in a photo? Answer concisely.” Class descriptions differ from well-

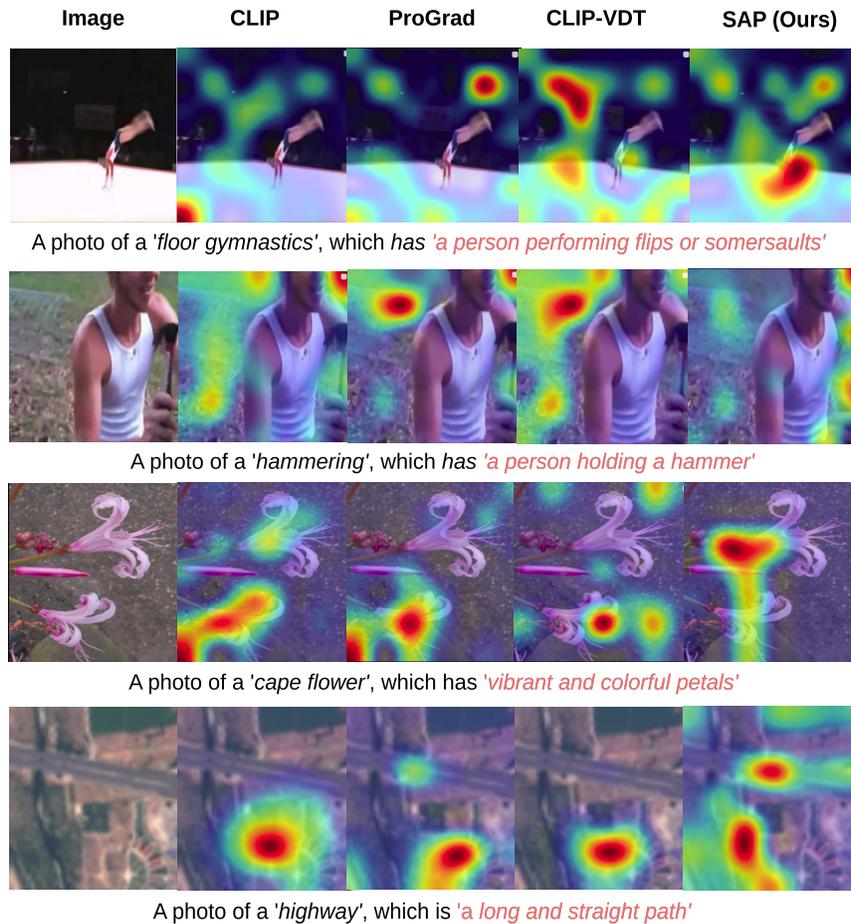


Figure 7: Figure displays GradCAM visualizations that highlight the regions of highest activation relevant to specific text phrases. These visualizations use a ResNet-50 backbone as the image encoder for all baselines, including our. SAP localizes better than the existing baselines.

curated attributes found in datasets with annotated attributes such as AWA (Lampert et al., 2009) and CUB (Wah et al., 2011) in three ways: (i) Our class descriptions may be noisy since no manual curation is used; (ii) They may not necessarily contain class-discriminative information, especially for similar classes; and (iii) Descriptions of a class are generated independently, and may not contain comparative traits w.r.t. other classes. These choices are primarily to keep our approach low-cost while integrating these finer details into fine-tuning of VLMs. It’s important to note that our description generation occurs at the class level, not the image level, making it cost-efficient.

## F Limitations and Broader Impact

A key dependency of our framework is the need for an LLM to provide descriptions at a class level. We however believe that this has become increasingly feasible in recent times, especially since we require at a class level and not at the image level. Our work deals with learning prompts for generalizable image classification by leveraging cheaply available semantic knowledge in the form of class descriptions. We believe that our work can serve as a stepping stone for incorporating semantic information to solve multi-modal tasks like captioning and VQA. To the best of our knowledge, there are no direct detrimental effects of our work.

Dataset		CLIP (ICML '21)	CoOp (IJCV '22)	VPT (ECCV '22)	CoCoOp (CVPR '22)	MaPLe (CVPR '23)	KgCoOp (CVPR '23)	ProGrad (ICCV '23)	PSRC (ICCV '23)	CLIP-VDT (ICCVW '23)	SAP (Ours)
Average on 11 datasets	gBase	60.81	75.19	73.48	73.13	75.47	76.86	70.15	<u>78.81</u>	63.75	<b>79.47 (+0.66)</b>
	gNovel	63.21	60.39	66.62	65.23	67.09	62.12	55.07	<u>68.13</u>	63.89	<b>69.75 (+1.62)</b>
	gHM	61.99	66.99	69.89	68.96	71.04	68.71	61.70	<u>73.08</u>	63.82	<b>74.29 (+1.21)</b>
UCF101	gBase	62.70	80.26	75.76	76.56	76.90	78.96	74.63	<b>82.67</b>	66.19	<u>82.23</u>
	gNovel	64.40	<b>84.76</b>	67.73	64.76	70.40	62.33	51.36	71.40	67.00	<u>76.40</u>
	gHM	63.53	<b>82.45</b>	71.52	70.17	73.51	69.67	60.85	76.62	66.59	<u>79.21</u>
EuroSAT	gBase	51.40	69.26	<u>88.22</u>	70.86	84.06	82.02	76.26	86.60	55.09	<b>94.37</b>
	gNovel	38.90	36.26	53.36	41.03	43.90	31.26	23.43	<u>54.16</u>	50.79	<b>58.53</b>
	gHM	44.28	47.60	66.50	51.97	57.68	45.28	35.85	<u>66.65</u>	52.85	<b>72.25</b>
DTD	gBase	42.70	65.36	58.92	60.29	63.00	66.42	57.19	<b>68.73</b>	55.79	<u>66.47</u>
	gNovel	45.79	34.30	44.26	46.09	47.49	39.73	33.36	<u>47.53</u>	51.00	<b>54.27</b>
	gHM	44.19	44.99	50.55	52.25	54.16	49.72	42.14	<u>56.20</u>	53.28	<b>59.75</b>
Oxford Pets	gBase	84.80	89.56	89.06	91.12	91.69	<u>91.99</u>	88.36	<b>93.00</b>	83.80	91.97
	gNovel	90.19	90.46	93.23	92.50	<b>93.93</b>	<u>92.69</u>	87.76	91.00	90.40	92.30
	gHM	87.41	90.01	91.10	91.81	<b>92.80</b>	<u>92.34</u>	88.06	91.99	86.97	92.13
Stanford Cars	gBase	56.00	74.43	65.13	67.29	69.33	72.56	64.46	74.77	59.50	<b>76.40</b>
	gNovel	64.19	57.16	<u>70.56</u>	68.82	69.86	66.56	55.66	<b>71.23</b>	61.59	69.33
	gHM	59.81	64.67	67.74	68.05	69.61	69.43	59.74	<b>72.96</b>	60.52	<u>72.69</u>
Flowers102	gBase	62.09	93.40	83.12	87.36	91.19	92.80	84.86	95.00	69.90	<b>95.69</b>
	gNovel	69.80	56.92	65.56	65.53	68.29	65.76	62.39	<u>71.00</u>	77.00	<b>71.13</b>
	gHM	65.71	70.74	73.31	74.89	78.10	76.97	71.92	<u>81.27</u>	73.20	<b>81.60</b>
Food101	gBase	79.90	83.59	85.96	86.15	<u>86.76</u>	85.76	78.46	<b>87.07</b>	75.90	86.43
	gNovel	80.90	76.82	84.99	<u>86.50</u>	<b>87.20</b>	83.72	76.23	85.90	77.69	86.09
	gHM	80.39	80.07	85.49	86.33	<b>86.98</b>	84.73	77.33	<u>86.48</u>	76.78	86.26
FGVC Aircraft	gBase	14.50	29.92	25.12	25.90	25.90	32.69	23.93	<u>34.90</u>	16.10	<b>35.00</b>
	gNovel	23.79	22.83	28.03	26.36	<u>28.53</u>	22.06	15.63	28.40	18.60	<b>30.23</b>
	gHM	18.01	25.90	26.50	26.13	27.15	26.35	18.93	<u>31.32</u>	17.59	<b>32.44</b>
SUN397	gBase	60.50	72.56	69.40	71.19	72.76	73.36	67.69	<b>75.63</b>	63.09	<u>75.40</u>
	gNovel	63.70	56.52	67.50	67.26	<u>68.93</u>	61.75	57.00	68.70	66.00	<b>69.80</b>
	gHM	62.05	63.55	68.44	69.17	70.79	67.06	61.89	<u>72.00</u>	64.51	<b>72.30</b>
Caltech101	gBase	91.40	95.92	95.66	95.09	95.83	95.89	91.53	96.20	93.59	<b>96.30</b>
	gNovel	91.69	85.09	<u>92.26</u>	90.93	92.03	92.06	85.26	91.73	86.19	<b>92.82</b>
	gHM	91.54	90.19	93.94	92.97	93.89	<u>93.94</u>	88.29	<u>93.91</u>	89.73	<b>94.53</b>
Imagenet	gBase	63.00	72.80	71.9	72.59	72.80	<u>73.00</u>	64.19	72.30	61.79	<b>73.97</b>
	gNovel	62.00	63.20	65.40	67.80	67.40	65.40	57.70	<b>68.40</b>	56.59	66.66
	gHM	62.49	67.66	68.50	70.11	70.00	68.99	60.77	<b>70.30</b>	59.07	<u>70.13</u>

Table 18: Accuracy comparison on the GZS benchmark. gNovel & gBase indicate the accuracy of the novel classes and base classes respectively under the joint classification label space. gHM is the harmonic mean of gBase and gNovel. The best numbers are in bold, and the second best are underlined. As reported in the first row, SAP outperforms all baselines on average gBase (by +0.66%), gNovel (by +1.62%), and gHM (by 1.21%) computed across all datasets. We indicate the margin of improvement over the corresponding best-performing baseline for each metric in green.

Dataset		CLIP	CoOp	VPT	CoCoOp	MaPLe	KgCoOp	ProGrad	PSRC	SAP
<b>Average on 11 datasets</b>	Base	33.28	36.97	40.28	40.12	<u>41.56</u>	37.95	34.00	40.40	<b>43.31 (+1.75)</b>
	Novel	38.55	<u>43.90</u>	43.72	40.80	43.30	40.69	35.01	43.78	<b>45.66 (+1.76)</b>
	HM	35.72	40.14	41.93	40.46	<u>42.41</u>	39.27	34.50	42.02	<b>44.46 (+2.04)</b>
UCF101	Base	56.60	61.20	61.20	61.70	<u>64.20</u>	62.00	59.70	63.10	<b>64.70</b>
	Novel	62.20	66.80	63.20	<b>70.70</b>	<u>70.40</u>	68.80	63.50	69.40	69.10
	HM	59.27	63.88	62.18	65.89	<b>67.16</b>	65.22	61.54	66.10	<u>66.83</u>
EuroSAT	Base	39.90	47.10	76.50	62.90	<u>84.30</u>	59.70	47.60	71.4	<b>88.70</b>
	Novel	71.10	78.70	<b>83.20</b>	49.00	58.30	57.60	45.80	<u>82.10</u>	80.90
	HM	51.12	58.93	<u>79.71</u>	55.09	68.93	58.63	46.68	76.38	<b>84.62</b>
DTD	Base	40.20	40.90	<u>47.20</u>	44.20	44.90	41.90	39.20	42.70	<b>52.40</b>
	Novel	42.40	44.10	44.30	<u>47.10</u>	42.90	44.40	40.20	44.00	<b>49.00</b>
	HM	41.27	42.44	<u>45.70</u>	45.60	43.88	43.11	39.69	43.34	<b>50.64</b>
Oxford Pets	Base	24.50	32.00	22.30	<b>34.20</b>	<u>32.80</u>	25.40	23.10	27.40	23.60
	Novel	35.20	40.80	40.70	<u>44.10</u>	<b>46.40</b>	39.70	36.00	41.60	<u>44.10</u>
	HM	28.89	35.87	28.81	<b>38.52</b>	<u>38.43</u>	30.98	28.14	33.04	30.75
Stanford Cars	Base	13.50	15.60	17.60	16.30	10.30	12.50	10.00	<u>21.00</u>	<b>22.50</b>
	Novel	15.90	20.70	18.90	11.70	<b>25.80</b>	15.30	8.50	20.40	<u>23.40</u>
	HM	14.60	17.79	18.23	13.62	14.72	13.76	9.19	<u>20.70</u>	<b>22.94</b>
Flowers102	Base	7.40	14.10	12.40	17.70	18.30	12.00	16.40	<u>18.80</u>	<b>19.60</b>
	Novel	9.30	20.40	18.40	17.60	<u>23.20</u>	12.30	13.80	19.30	<b>26.00</b>
	HM	8.24	16.67	14.82	17.65	<u>20.46</u>	12.15	14.99	19.05	<b>22.35</b>
Food101	Base	35.10	42.70	<b>44.00</b>	<u>43.40</u>	35.50	47.10	42.10	41.20	42.20
	Novel	33.80	<b>45.40</b>	<u>44.80</u>	44.40	38.90	44.60	41.80	40.50	44.20
	HM	34.44	<u>44.01</u>	44.40	43.89	37.12	<b>45.82</b>	41.95	40.85	43.18
FGVC Aircraft	Base	6.10	<u>9.50</u>	8.00	7.00	<b>13.40</b>	6.80	5.20	8.30	9.40
	Novel	7.90	<b>15.80</b>	12.80	8.30	<u>15.50</u>	10.70	8.20	12.30	12.30
	HM	6.88	<u>11.87</u>	9.85	7.59	<b>14.37</b>	8.32	6.36	9.91	10.66
SUN397	Base	46.60	49.20	50.50	<u>51.30</u>	50.20	50.10	40.10	50.00	<b>51.40</b>
	Novel	48.30	50.00	51.40	<u>52.50</u>	52.20	<b>53.20</b>	42.90	51.40	51.40
	HM	47.43	49.60	50.95	<b>51.89</b>	51.18	<u>51.60</u>	41.45	50.69	51.40
Caltech101	Base	77.80	76.00	<b>83.00</b>	<b>83.00</b>	82.30	80.80	72.30	81.10	81.70
	Novel	74.80	74.30	<u>75.90</u>	75.80	75.50	<b>76.20</b>	63.20	75.10	75.20
	HM	76.27	75.14	<b>79.29</b>	<u>79.24</u>	78.75	78.43	67.44	77.98	78.32
ImageNet	Base	18.40	18.40	<u>20.40</u>	19.70	<b>21.00</b>	19.20	18.30	19.4	20.30
	Novel	23.20	26.00	<u>27.40</u>	<b>27.60</b>	27.30	24.80	21.30	25.50	26.70
	HM	20.52	21.55	<u>23.39</u>	22.99	<b>23.74</b>	21.64	19.69	22.04	23.06

Table 19: Accuracy comparison in the Classification without Class-names setting. We show average Base, Novel, and HM accuracies over all 11 datasets. During evaluation, descriptions of each class are provided instead of the class name, and visual recognition is conducted based on these descriptions. SAP outperforms baselines by average Base (by +1.75%), Novel (by +1.76%) and HM (by +2.04%) computed over all datasets.

	Source	Target										
	ImageNet	Caltech101	OxfordPets	StanfordCars	Flowers102	Food101	Aircraft	SUN397	DTD	EuroSAT	UCF101	Average
CoOp	71.51	93.70	89.14	64.51	68.71	85.30	18.47	64.15	41.92	46.39	66.55	63.88
CoCoOp	71.02	94.43	90.14	65.32	71.88	86.06	22.94	67.36	45.73	45.37	68.21	65.74
VPT	70.60	91.80	90.40	63.70	67.30	83.10	22.70	66.10	46.10	37.10	65.90	63.42
MaPLe	70.72	93.53	90.49	65.57	72.23	86.20	24.74	67.01	46.49	48.06	68.69	66.30
KgCoOp	69.94	94.08	90.13	65.63	71.21	86.48	23.85	67.47	45.80	41.98	68.33	65.49
ProGrad	62.17	88.30	86.43	55.61	62.69	76.76	15.76	60.16	39.48	28.47	58.70	57.36
PSRC	71.27	93.60	90.25	65.70	70.25	86.15	23.90	67.10	46.87	45.50	68.75	65.81
CLIP-VDT	68.10	85.40	83.50	50.30	56.00	72.50	14.60	56.30	42.70	24.70	53.80	53.98
KAPT	N/A	88.90	89.40	58.15	68.00	79.95	17.95	N/A	44.80	41.35	65.05	61.50
SAP (Ours)	71.40	94.53	90.14	64.58	71.31	86.23	24.47	68.09	48.61	49.10	71.52	<b>66.85</b>

Table 20: Cross-Dataset Generalization benchmark. Models are trained on Imagenet and tested on the entire label space of new datasets without fine-tuning. SAP outperforms all baselines on average. N/A: not available in (Kan et al., 2023).

Dataset		CLIP	CoOp	VPT	CoCoOp	ProDA	MaPLe	KgCoOp	ProGrad	PSRC	L.Prompt	CLIP-VDT	KAPT	SAP
Average on 11 datasets	Base	69.34	82.69	80.81	80.47	81.56	82.28	80.73	82.48	84.26	84.47	82.48	81.10	<b>84.68 (+0.21)</b>
	Novel	74.22	63.22	70.36	71.69	72.30	75.14	73.60	70.75	<u>76.10</u>	74.24	74.50	72.24	<b>77.51 (+1.41)</b>
	HM	71.70	71.66	70.36	75.83	76.65	78.55	77.00	76.16	<u>79.97</u>	79.03	78.28	76.41	<b>80.94 (+0.97)</b>
UCF101	Base	70.53	84.69	82.67	82.33	85.23	83.00	82.89	84.33	<b>87.10</b>	86.19	84.10	80.83	<u>86.60</u>
	Novel	77.50	56.05	74.54	77.64	78.04	<u>80.77</u>	76.67	76.94	78.80	73.07	76.40	67.10	<b>83.90</b>
	HM	73.85	67.46	78.39	77.64	78.04	80.77	79.65	79.35	<u>82.74</u>	79.09	80.07	73.33	<b>85.23</b>
EuroSAT	Base	56.48	92.19	93.01	87.49	83.90	<u>94.07</u>	85.64	90.11	92.90	93.67	88.50	84.80	<b>96.10</b>
	Novel	64.05	54.74	54.89	60.04	66.00	73.23	64.34	60.89	<u>73.90</u>	69.44	70.50	67.57	<b>81.13</b>
	HM	60.03	68.69	69.04	71.21	73.88	<u>82.35</u>	73.48	72.67	82.32	79.75	78.48	75.21	<b>87.98</b>
DTD	Base	53.24	79.44	79.15	77.01	80.67	80.36	77.55	77.35	<u>83.37</u>	82.87	81.80	75.97	<b>84.27</b>
	Novel	59.90	41.18	50.76	56.00	56.48	59.18	54.99	52.35	<u>62.97</u>	60.14	62.30	58.30	<b>67.03</b>
	HM	56.37	54.24	61.85	64.85	66.44	68.16	64.35	62.45	<u>71.75</u>	69.70	70.73	65.97	<b>74.67</b>
Oxford Pets	Base	91.17	93.67	94.81	95.20	<u>95.43</u>	<u>95.43</u>	94.65	95.07	95.33	<b>96.07</b>	94.40	93.13	95.27
	Novel	97.26	95.29	96.00	97.69	<b>97.83</b>	<u>97.76</u>	<u>97.76</u>	97.63	97.30	96.31	97.70	96.53	96.90
	HM	94.12	94.47	95.40	96.43	<b>96.62</b>	<u>96.58</u>	96.18	96.33	96.30	96.18	95.68	94.80	96.08
Stanford Cars	Base	63.37	78.12	72.46	70.49	74.70	72.94	71.76	77.68	78.27	<u>78.36</u>	76.80	69.47	<b>79.70</b>
	Novel	74.89	60.40	73.38	73.59	71.20	74.00	<b>75.04</b>	68.63	<u>74.97</u>	72.39	72.90	66.20	73.47
	HM	68.65	68.13	72.92	72.01	72.91	73.47	73.36	72.88	<b>76.58</b>	75.26	74.80	67.79	<u>76.46</u>
Flowers102	Base	72.08	97.60	95.39	94.87	97.70	95.92	95.00	95.54	<u>98.07</u>	<b>99.05</b>	97.40	95.00	97.83
	Novel	<b>77.80</b>	59.67	73.87	71.75	68.68	72.46	74.73	71.87	<u>76.50</u>	76.52	75.30	71.20	<u>76.50</u>
	HM	74.83	74.06	83.26	81.71	80.66	82.56	83.65	82.03	85.95	<u>86.34</u>	84.94	81.40	<b>86.86</b>
Food101	Base	90.10	88.33	89.88	90.70	90.30	<u>90.71</u>	90.50	90.37	90.67	<b>90.82</b>	90.40	86.13	90.40
	Novel	91.22	82.26	87.76	91.29	88.57	<b>92.05</b>	<u>91.70</u>	89.59	91.53	91.41	91.20	87.06	91.43
	HM	90.66	85.19	88.81	90.99	89.43	<b>91.38</b>	91.09	89.98	91.10	<u>91.11</u>	90.80	86.59	90.91
FGVC Aircraft	Base	27.19	40.44	33.10	33.41	36.90	37.44	36.21	40.54	42.73	<b>45.98</b>	37.80	29.67	<u>42.93</u>
	Novel	36.29	22.30	30.49	23.71	34.13	35.61	33.55	27.57	<u>37.87</u>	34.67	33.00	28.73	<b>38.87</b>
	HM	31.09	28.75	31.74	27.74	35.46	36.50	34.83	32.82	<u>40.15</u>	39.53	35.24	29.19	<b>40.80</b>
SUN397	Base	69.36	80.60	79.66	79.74	78.67	80.82	80.29	81.26	<b>82.67</b>	81.20	81.40	79.40	<u>82.57</u>
	Novel	75.35	65.89	72.68	76.86	76.93	78.70	76.53	74.17	<u>78.47</u>	78.12	76.80	74.33	<b>79.20</b>
	HM	72.23	72.51	79.63	78.27	77.79	79.75	78.36	77.55	<u>80.52</u>	79.63	79.03	76.78	<b>80.85</b>
Caltech101	Base	96.84	98.00	97.86	97.96	<u>98.27</u>	97.74	97.72	98.02	98.10	98.19	<b>98.30</b>	97.10	98.23
	Novel	94.00	89.91	93.76	93.81	93.23	94.36	<u>94.39</u>	93.89	94.03	93.78	<b>95.90</b>	93.53	94.37
	HM	95.40	93.73	95.77	95.84	95.68	96.02	96.03	95.91	96.02	95.93	<b>97.09</b>	95.28	<u>96.26</u>
ImageNet	Base	72.43	76.47	70.93	75.98	75.40	76.66	75.83	77.02	<b>77.60</b>	76.74	76.40	71.10	<b>77.60</b>
	Novel	68.14	67.88	65.90	70.43	70.23	70.54	69.96	66.66	<u>70.73</u>	<b>70.83</b>	68.30	65.20	69.83
	HM	70.22	71.92	73.66	73.10	72.72	73.47	72.78	71.46	<b>74.01</b>	<u>73.66</u>	72.12	68.02	73.51

Table 21: Accuracy comparison on Base-to-Novel Generalization benchmark. The best numbers are in bold, and the second best are underlined. SAP outperforms all baselines on average Base (by +0.21%), Novel (by +1.41%) and HM (by +0.97%) computed over all datasets. We indicate the margin of improvement over the corresponding best-performing baseline for each metric in green.

Dataset		GZS Benchmark					SAP	Base-to-Novel Benchmark					SAP	
		CLIP	CoOp	KgCoOp	Pro-Grad	PSRC		CLIP	CoOp	KgCoOp	Pro-Grad	PSRC		
Average on 11 datasets	gBase	57.01	68.65	69.25	<u>69.89</u>	47.41	<b>71.52 (+1.63)</b>	Base	65.27	77.24	75.51	<u>77.98</u>	55.13	<b>78.49 (+0.51)</b>
	gNovel	<b>60.73</b>	50.35	59.08	52.26	29.16	<u>59.13 (-1.60)</u>	Novel	68.14	57.40	<u>67.53</u>	63.41	38.72	<b>69.32 (+1.79)</b>
	gHM	58.81	58.10	<u>63.76</u>	59.81	36.12	<b>64.74 (+0.98)</b>	HM	66.68	65.86	<u>71.30</u>	69.94	45.49	<b>73.62 (+2.32)</b>
UCF101	gBase	61.20	<u>73.20</u>	71.05	72.75	51.55	<b>74.73</b>	Base	68.40	79.78	77.16	<b>81.04</b>	59.95	<u>80.70</u>
	gNovel	61.79	45.10	56.95	48.05	30.25	<b>63.80</b>	Novel	61.50	48.31	<u>70.13</u>	60.07	38.85	<b>72.67</b>
	gHM	61.49	55.81	63.22	57.87	38.13	<u>68.33</u>	HM	64.77	60.18	<u>73.48</u>	69.00	47.15	<b>76.47</b>
EuroSAT	gBase	32.79	62.70	71.25	<b>73.60</b>	61.15	<u>72.77</u>	Base	55.80	<u>90.25</u>	84.28	88.44	70.35	<b>91.33</b>
	gNovel	<b>46.50</b>	23.45	<u>33.95</u>	19.40	09.00	32.32	Novel	66.90	31.30	<u>53.53</u>	49.49	33.90	<b>67.00</b>
	gHM	38.46	34.13	<b>45.99</b>	30.71	15.69	<u>44.76</u>	HM	60.85	46.48	65.47	<u>63.47</u>	45.75	<b>77.30</b>
DTD	gBase	43.50	60.60	<u>64.80</u>	62.30	42.60	<u>62.73</u>	Base	53.70	75.12	74.73	73.80	51.35	<b>75.97</b>
	gNovel	<u>41.29</u>	27.05	40.45	27.05	18.30	<b>44.27</b>	Novel	55.60	37.08	<u>48.39</u>	46.38	29.85	<b>57.90</b>
	gHM	42.37	37.40	<u>49.81</u>	37.72	25.60	<b>51.91</b>	HM	54.63	49.65	<b>58.74</b>	56.96	37.75	<b>65.72</b>
Oxford Pets	gBase	85.90	84.70	85.75	<u>85.95</u>	67.65	<b>87.00</b>	Base	91.20	90.15	<b>92.57</b>	<u>92.36</u>	77.60	91.90
	gNovel	85.59	85.25	<b>90.45</b>	87.10	65.65	<u>89.27</u>	Novel	93.90	90.70	<b>94.61</b>	94.48	79.40	<u>94.57</u>
	gHM	85.74	84.97	<u>88.04</u>	86.52	66.63	<b>88.12</b>	HM	92.53	90.42	<u>93.58</u>	93.41	78.49	93.22
Stanford Cars	gBase	48.29	<u>64.70</u>	62.25	64.30	17.35	<b>68.20</b>	Base	55.50	68.89	63.28	<b>71.79</b>	26.35	<u>71.43</u>
	gNovel	<b>64.09</b>	48.05	<u>59.20</u>	53.45	21.65	57.60	Novel	66.50	57.13	<b>66.92</b>	59.36	25.50	<u>64.77</u>
	gHM	55.08	55.15	<u>60.69</u>	58.38	19.26	<b>62.45</b>	HM	60.50	62.46	65.05	<u>64.99</u>	25.92	<b>67.94</b>
Flowers102	gBase	62.59	<u>89.40</u>	85.70	88.80	65.00	<b>92.52</b>	Base	69.70	<u>95.22</u>	91.45	94.71	73.75	<b>96.40</b>
	gNovel	<b>68.30</b>	50.70	<u>63.85</u>	52.75	10.85	61.62	Novel	73.90	59.53	<b>71.75</b>	68.86	19.75	<u>70.30</u>
	gHM	65.32	64.70	<u>73.18</u>	66.18	18.60	<b>73.97</b>	HM	71.74	73.26	<u>80.41</u>	79.74	31.16	<b>81.31</b>
Food101	gBase	<b>75.80</b>	73.80	<b>78.30</b>	76.30	32.65	<u>77.97</u>	Base	83.10	81.70	<u>83.90</u>	83.77	37.85	83.57
	gNovel	<b>78.90</b>	68.50	<u>78.25</u>	72.90	17.60	76.60	Novel	84.50	78.13	<u>85.23</u>	83.74	27.15	84.13
	gHM	<u>77.32</u>	71.05	<b>78.27</b>	74.56	22.87	77.28	HM	83.79	79.88	<u>84.56</u>	83.75	31.62	83.85
FGVC Aircraft	gBase	12.69	<b>24.15</b>	20.20	21.60	8.65	<u>23.17</u>	Base	18.80	28.39	24.91	<b>30.17</b>	14.20	<u>28.97</u>
	gNovel	22.10	14.75	<b>18.20</b>	14.25	6.95	<u>17.45</u>	Novel	26.00	20.02	<b>25.69</b>	19.70	9.05	<u>25.33</u>
	gHM	16.12	18.31	<u>19.15</u>	17.17	7.71	<b>19.91</b>	HM	21.82	23.48	<u>25.29</u>	23.84	11.05	<b>27.03</b>
SUN397	gBase	56.70	66.65	67.05	<u>67.15</u>	54.25	<b>70.40</b>	Base	66.40	76.33	75.33	<u>76.90</u>	63.25	<b>78.20</b>
	gNovel	60.50	53.30	<u>61.80</u>	56.50	45.85	<b>62.20</b>	Novel	70.10	62.89	72.25	68.09	57.50	<b>73.27</b>
	gHM	58.54	59.23	<u>64.32</u>	61.37	49.70	<b>66.05</b>	HM	70.10	68.96	<u>73.76</u>	72.23	60.24	<b>75.65</b>
Caltech101	gBase	88.59	91.35	<u>91.65</u>	91.50	79.35	<b>92.13</b>	Base	91.00	95.20	95.35	<b>95.72</b>	84.80	<u>95.67</u>
	gNovel	81.69	82.15	<b>88.05</b>	86.30	58.65	<u>87.50</u>	Novel	90.60	87.55	<b>91.92</b>	89.92	65.65	<u>91.13</u>
	gHM	85.00	86.51	<b>89.81</b>	88.82	67.45	<u>89.76</u>	HM						

Class (Dataset)	Descriptions	Class (Dataset)	Descriptions
Breast stroke (UCF101)	<ol style="list-style-type: none"> <li>1. Arms moving in a circular motion</li> <li>2. Kicking legs in a frog-like motion</li> <li>3. Head above water during stroke</li> <li>4. Positioned horizontally in the water</li> <li>5. Pushing water forward and outwards</li> </ol>	Diving (UCF101)	<ol style="list-style-type: none"> <li>1. Person in mid-air or jumping</li> <li>2. Person wearing diving gear</li> <li>3. water splashing or ripples</li> <li>4. Person wearing goggles</li> <li>5. Person wearing swim cap</li> </ol>
Highway or road (EuroSAT)	<ol style="list-style-type: none"> <li>1. Long and straight path</li> <li>2. Multiple lanes for traffic</li> <li>3. Traffic signs</li> <li>4. Smooth and paved surface</li> <li>5. Guardrails or barriers</li> </ol>	Permanent cropland (EuroSAT)	<ol style="list-style-type: none"> <li>1. Uniform vegetation or crops</li> <li>2. Irrigation systems or canals</li> <li>3. Organized rows or patterns</li> <li>4. Fences or boundaries</li> <li>5. Distinct crop types or varieties</li> </ol>
Striped (DTD)	<ol style="list-style-type: none"> <li>1. Alternating bands or lines</li> <li>2. Regular pattern of stripes</li> <li>3. Varying widths of stripes</li> <li>4. Contrasting colors between stripes</li> <li>5. Horizontal, vertical, diagonal stripes</li> </ol>	Wrinkled (DTD)	<ol style="list-style-type: none"> <li>1. Irregular and uneven surface</li> <li>2. Creases or folds</li> <li>3. Shadows indicating unevenness</li> <li>4. Lack of smoothness</li> <li>5. Distorted or crumpled appearance</li> </ol>
Maine coon (Oxford Pets)	<ol style="list-style-type: none"> <li>1. Large domestic cat</li> <li>2. Long, bushy tail</li> <li>3. Tufted ears with lynx-like tips</li> <li>4. Rectangular body shape</li> <li>5. Tufted paws</li> </ol>	Chihuahua (Oxford Pets)	<ol style="list-style-type: none"> <li>1. Small breed of dog</li> <li>2. Rounded apple-shaped head</li> <li>3. Erect, pointy ears</li> <li>4. Short snout</li> <li>5. Short legs and long tail</li> </ol>
2008 chrysler pt cruiser convertible (Stanford Cars)	<ol style="list-style-type: none"> <li>1. Convertible top</li> <li>2. Chrome grille</li> <li>3. PT cruiser badge</li> <li>4. Alloy wheels</li> <li>5. Boxy shape</li> </ol>	2012 ferrari ff coupe (Stanford Cars)	<ol style="list-style-type: none"> <li>1. Sleek and sporty design</li> <li>2. Large and stylish alloy wheels</li> <li>3. Low and wide stance</li> <li>4. Ferrari logo on the front and rear</li> <li>5. Dual exhaust pipes</li> </ol>
Watercress (Flowers102)	<ol style="list-style-type: none"> <li>1. Small, round-shaped leaves</li> <li>2. Vibrant green color</li> <li>3. Thin, delicate stems</li> <li>4. Water or moist environments</li> <li>5. Clusters of small white flowers</li> </ol>	Trumpet creeper (Flowers102)	<ol style="list-style-type: none"> <li>1. Bright orange or red flowers</li> <li>2. Trumpet-shaped blossoms</li> <li>3. Long, tubular petals</li> <li>4. Green leaves with serrated edges</li> <li>5. Hummingbirds and bees</li> </ol>
Hot dog (Food101)	<ol style="list-style-type: none"> <li>1. Cylindrical-shaped food</li> <li>2. Bun or bread</li> <li>3. Sausage or frankfurter</li> <li>4. Visible grill marks</li> <li>5. Toppings like onions or relish</li> </ol>	Sushi (Food101)	<ol style="list-style-type: none"> <li>1. Bite-sized and compact</li> <li>2. Rice as a base</li> <li>3. Raw or cooked fish</li> <li>4. Seaweed wrapping (nori)</li> <li>5. Served with soy sauce</li> </ol>
737-200 (FGVC Aircraft)	<ol style="list-style-type: none"> <li>1. Two engines on the wings</li> <li>2. Low wing configuration</li> <li>3. Narrow body</li> <li>4. Distinctive short fuselage</li> <li>5. Swept-back wings</li> </ol>	Industrial area (SUN397)	<ol style="list-style-type: none"> <li>1. Factories or warehouses</li> <li>2. Smokestacks or chimneys</li> <li>3. Cranes or heavy machinery</li> <li>4. Conveyor belts or assembly lines</li> <li>5. Trucks or shipping containers</li> </ol>
Gramophone (Caltech101)	<ol style="list-style-type: none"> <li>1. Phonograph Cylinder or Disc</li> <li>2. Horn Speaker</li> <li>3. Hand-Cranked Operation</li> <li>4. Nostalgic and Vintage Appeal</li> <li>5. Vinyl or Shellac Records</li> </ol>	Buckle (Imagenet)	<ol style="list-style-type: none"> <li>1. Metal or plastic object</li> <li>2. Rectangular or circular shape</li> <li>3. Fastening or securing</li> <li>4. Opened and closed</li> <li>5. Found on belts or straps</li> </ol>

Table 23: Sample classes from various datasets and the corresponding descriptions provided by GPT-3.5.