

## Supplementary Materials Organization:

---

<b>A</b>	<b>More Experimental Analysis</b>	<b>13</b>
A.1	Base-to-new Generalization . . . . .	13
A.2	Robustness to Different Architectures . . . . .	13
<b>B</b>	<b>Experimental Details</b>	<b>13</b>
B.1	Statistic of Datasets . . . . .	13
B.2	Prompt Templates for Each Dataset . . . . .	15
B.3	Pseudocode . . . . .	15

---

## A MORE EXPERIMENTAL ANALYSIS

## A.1 BASE-TO-NEW GENERALIZATION

**Results.** Our method can be extended to the base-to-new generalization scenario by incorporating the KNN algorithm. To accomplish this, we utilize the text embeddings of the new classes to query the training set and select the  $k$  nearest neighbors as the training data for the new class. Subsequently, we apply our proposed method to generate the classifier for the new classes using the synthesized dataset. In order to compare our approach, we select CLIP (Radford et al., 2021), CoOp (Zhou et al., 2022b), and CoCoOp (Zhou et al., 2022a).

Table 8 presents the results, which demonstrate that our approach outperforms the other methods in terms of base accuracy, new accuracy, and their harmonic mean. On average across 11 datasets, our method surpasses CLIP, CoOp, and CoCoOp by 14.62%, 1.27%, and 3.49% in terms of base accuracy. It also outperforms them by 0.31%, 11.31%, and 2.84% in terms of new accuracy, and by 7.02%, 7.06%, and 2.89% in terms of the harmonic mean. Moreover, our approach achieves the highest harmonic mean in 8 out of 11 datasets. These results clearly indicate the effectiveness of our approach in generalizing to new classes.

## A.2 ROBUSTNESS TO DIFFERENT ARCHITECTURES

We further evaluate the efficacy of our proposed method across 11 datasets with varying visual architectures of CLIP. We selected two approaches for comparison: a training-required method, CoOp (Zhou et al., 2022b), and a training-free method, Tip-Adapter (Zhang et al., 2022). And these methods are trained on the 16-shot dataset. As shown in Table 9, our method yielded a substantial improvement of 17.28%, 18.20%, 16.18%, and 16.62% on average, compared to the zero-shot CLIP (Radford et al., 2021) approach, for ResNet-50, ResNet-101, ViT-B/32, and ViT-B/16 CLIP, respectively, across all 11 datasets. The results demonstrate the effectiveness of our method across different CLIP architectures.

## B EXPERIMENTAL DETAILS

## B.1 STATISTIC OF DATASETS

Following previous work (Zhou et al., 2022a;b; Wang et al., 2023c; Huang et al., 2022; Wang et al., 2023a), we conduct experiments on 17 publicly available image classification datasets. The datasets include ImageNet (Deng et al., 2009), Caltech101 (Li et al., 2004), OxfordPets (Parkhi et al., 2012), StanfordCars (Krause et al., 2013), Flowers102 (Nilsback & Zisserman, 2008), Food101 (Bossard et al., 2014), FGVCAircraft (Maji et al., 2013), EuroSAT (Helber et al., 2019), UCF101 (Soomro et al., 2012), DTD (Cimpoi et al., 2014), SUN397 (Xiao et al., 2010), ImageNetV2 (Recht et al.,

(a) Average over 11 datasets				(b) ImageNet				(c) Caltech101			
	base	new	H		base	new	H		base	new	H
CLIP	69.34	74.22	71.70	CLIP	72.43	68.14	70.22	CLIP	96.84	94.00	95.40
CoOp	82.69	63.22	71.66	CoOp	<b>76.47</b>	67.88	71.92	CoOp	98.00	89.81	93.73
CoCoOp	80.47	71.69	75.83	CoCoOp	75.98	<b>70.43</b>	<b>73.10</b>	CoCoOp	97.96	93.81	95.84
Ours	<b>83.96</b>	<b>74.53</b>	<b>78.72</b>	Ours	75.95	69.79	72.74	Ours	<b>98.04</b>	<b>94.51</b>	<b>96.24</b>
(d) OxfordPets				(e) StanfordCars				(f) Flowers102			
	base	new	H		base	new	H		base	new	H
CLIP	91.17	97.26	94.12	CLIP	63.37	<b>74.89</b>	68.65	CLIP	72.08	<b>77.80</b>	74.83
CoOp	93.67	95.29	94.47	CoOp	78.12	60.40	68.13	CoOp	97.60	59.67	74.06
CoCoOp	<b>95.20</b>	<b>97.69</b>	<b>96.43</b>	CoCoOp	70.49	73.59	72.01	CoCoOp	94.87	71.75	81.71
Ours	94.10	97.15	95.60	Ours	<b>78.71</b>	66.92	<b>72.34</b>	Ours	<b>97.78</b>	72.46	<b>83.24</b>
(g) Food101				(h) FGVCAircraft				(i) SUN397			
	base	new	H		base	new	H		base	new	H
CLIP	90.10	91.22	90.66	CLIP	27.19	<b>36.29</b>	31.09	CLIP	69.36	75.35	72.23
CoOp	88.33	82.26	85.19	CoOp	40.44	22.30	28.75	CoOp	80.60	65.89	72.51
CoCoOp	<b>90.70</b>	<b>91.29</b>	<b>90.99</b>	CoCoOp	33.41	23.71	27.74	CoCoOp	79.74	<b>76.86</b>	78.27
Ours	90.63	91.21	90.92	Ours	<b>45.88</b>	34.09	<b>39.12</b>	Ours	<b>81.95</b>	75.62	<b>78.65</b>
(j) DTD				(k) EuroSAT				(l) UCF101			
	base	new	H		base	new	H		base	new	H
CLIP	53.24	<b>59.90</b>	56.37	CLIP	56.48	64.05	60.03	CLIP	70.53	77.50	73.85
CoOp	79.44	41.18	54.24	CoOp	92.19	54.74	68.69	CoOp	84.69	56.05	67.46
CoCoOp	77.01	56.00	64.85	CoCoOp	87.49	60.04	71.21	CoCoOp	82.33	73.45	77.64
Ours	<b>80.63</b>	59.82	<b>68.69</b>	Ours	<b>93.28</b>	<b>79.21</b>	<b>85.67</b>	Ours	<b>86.63</b>	<b>79.09</b>	<b>82.69</b>

Table 8: **Base-to-new generalization.** Comparison of CLIP, CoOp, CoCoOp, and our method. CoOp and CoCoOp are training-required methods, while our method is a training-free method. base and new denotes the average accuracy of base and new classes, and H denotes their harmonic mean.

Table 9: **Robustness of different architectures on 11 datasets.** The models are trained under the 16-shot setting with different visual architectures of CLIP. **Bold** denotes the highest results.

Method	Pets	Flowers	FGVC	DTD	EuroSAT	Cars	Food	SUN	Cal	UCF	IN	Avg.
<b>ResNet-50</b>												
zero-shot CLIP	85.77	66.14	17.28	42.32	37.56	55.61	77.31	58.52	86.29	61.46	58.18	58.77
CoOp	87.01	94.51	31.26	63.58	83.53	73.36	74.67	69.26	91.83	75.71	62.95	73.42
Tip-Adapter	88.14	89.89	29.76	60.93	70.54	66.77	77.83	66.85	90.18	70.58	62.01	70.32
Ours	<b>88.81</b>	<b>95.72</b>	<b>40.61</b>	<b>66.51</b>	<b>86.12</b>	<b>75.12</b>	<b>79.05</b>	<b>70.70</b>	<b>92.55</b>	<b>77.53</b>	<b>63.82</b>	<b>76.05</b>
<b>ResNet-101</b>												
zero-shot CLIP	86.75	64.03	18.42	38.59	32.59	66.23	80.53	58.96	89.78	60.96	61.62	59.86
CoOp	88.57	95.19	34.76	65.47	83.54	79.74	79.08	71.19	93.42	77.95	<b>66.60</b>	75.96
Tip-Adapter	87.23	90.77	31.51	62.37	66.45	72.96	81.31	67.96	93.01	73.53	64.41	71.96
Ours	<b>91.43</b>	<b>96.17</b>	<b>42.58</b>	<b>68.62</b>	<b>86.32</b>	<b>79.99</b>	<b>82.15</b>	<b>72.07</b>	<b>93.63</b>	<b>79.31</b>	66.33	<b>78.06</b>
<b>ViT-B/32</b>												
zero-shot CLIP	87.49	66.95	19.23	43.97	45.19	60.55	80.50	61.91	90.87	62.01	62.05	61.88
CoOp	88.68	94.97	33.22	65.37	83.43	76.08	78.45	72.38	<b>94.62</b>	78.66	66.85	75.70
Tip-Adapter	88.34	91.61	30.92	61.90	69.53	69.59	80.94	70.27	93.85	73.74	65.41	72.37
Ours	<b>91.21</b>	<b>96.16</b>	<b>41.74</b>	<b>67.63</b>	<b>87.30</b>	<b>77.55</b>	<b>81.84</b>	<b>73.60</b>	94.42	<b>80.17</b>	<b>67.00</b>	<b>78.06</b>
<b>ViT-B/16</b>												
zero-shot CLIP	89.21	71.34	24.72	44.39	47.60	65.32	86.06	62.50	92.94	66.75	66.73	65.23
CoOp	92.53	96.47	42.91	68.50	80.87	<b>83.09</b>	87.21	75.29	95.77	82.24	71.92	79.71
Tip-Adapter	91.54	94.41	39.48	65.68	76.58	75.44	86.47	71.85	95.10	77.94	70.46	76.81
Ours	<b>93.73</b>	<b>97.92</b>	<b>50.33</b>	<b>71.26</b>	<b>89.19</b>	82.63	<b>87.27</b>	<b>75.87</b>	<b>95.79</b>	<b>84.09</b>	<b>72.24</b>	<b>81.85</b>

2019), ImageNet-Sketch (Wang et al., 2019), ImageNet-A (Hendrycks et al., 2021b), ImageNet-R (Hendrycks et al., 2021a), ImageNet-LT (Liu et al., 2019), and Places-LT (Zhou et al., 2017).

Table 10: Detailed statistics of datasets used in experiments.

Dataset	# Classes	# Training	# Test	Task
<b>OxfordPets</b>	37	2,944	3,669	fine-grained pets recognition
<b>Flowers102</b>	102	4,093	2,463	fine-grained flowers recognition
<b>FGVCAircraft</b>	100	3,334	3,333	fine-grained aircraft recognition
<b>DTD</b>	47	2,820	1,692	Textural recognition
<b>EuroSAT</b>	10	13,500	8,100	Satellite image recognition
<b>StanfordCars</b>	196	6,509	8,041	Fine-grained car recognition
<b>Food101</b>	101	50,500	30,300	Fine-grained food recognition
<b>Sun397</b>	397	15,880	19,850	Scene recognition
<b>Caltech101</b>	100	4,128	2,465	Object recognition
<b>UCF101</b>	101	7,639	3,783	Action recognition
<b>ImageNet</b>	1,000	1.28M	50,000	Object recognition
<b>ImageNetV2</b>	1,000	-	10,000	Robustness of collocation
<b>ImageNet-Sketch</b>	1,000	-	50,889	Robustness of sketch domain
<b>ImageNet-A</b>	200	-	7,500	Robustness of adversarial
<b>ImageNet-R</b>	200	-	30,000	Robustness of rendition styles
<b>ImageNet-LT</b>	1,000	115,846	50,000	long-tail object recognition
<b>Places-LT</b>	365	62,500	7300	long-tail place recognition

## B.2 PROMPT TEMPLATES FOR EACH DATASET

For the zero-shot classifier, we employ handcrafted prompts to generate the classifier weight, as proposed in CLIP (Radford et al., 2021). By default, we utilize the prompt template “a photo of {class}.” for class labels, where {class} represents the name of the classes. However, for fine-grained classification datasets such as FGVCAircraft (Maji et al., 2013), we incorporate the name of the superclass or a description into the template. The prompt templates for each dataset are shown as follows.

## B.3 PSEUDOCODE

Table 11: Prompt templates for each class.

Dataset	Prompt template
Caltech101 (Li et al., 2004)	“a photo of a {class}.”
OxfordPets (Parkhi et al., 2012)	“a photo of a {class}, a type of pet.”
StanfordCars (Krause et al., 2013)	“a photo of a {class}.”
Flowers102 (Nilsback & Zisserman, 2008)	“a photo of a {class}, a type of flower.”
Food101 (Bossard et al., 2014)	“a photo of {class}, a type of food.”
FGVCAircraft (Maji et al., 2013)	“a photo of a {class}, a type of aircraft.”
SUN397 (Xiao et al., 2010)	“a photo of a {class}.”
DTD (Cimpoi et al., 2014)	“{class} texture.”
EuroSAT (Helber et al., 2019)	“a centered satellite photo of {class}.”
UCF101 (Soomro et al., 2012)	“a photo of a person doing {class}.”
ImageNet (Deng et al., 2009)	“a bad photo of the {class}.” “a origami {class}.” “a photo of the large {class}.” “a {class} in a video game.” “art of the {class}.” “a photo of the small {class}.”
ImageNet-LT (Liu et al., 2019)	“a photo of a {class}.”
Places-LT (Zhou et al., 2017)	“a photo of a {class}.”

**Algorithm 1** Pytorch-like pseudocode for our method.

---

```

1  # Input:
2  # - X: (N, D) visual features from CLIP visual encoder.
3  # - Y: (N, ) ground-truth label for the features.
4  # - X_test: (M, D) test visual features from CLIP visual encoder.
5  # - Y_test: (M, ) ground-truth label for test features.
6  # - W_c: (K, D) zero-shot classifier generated by prompting.
7  # Output:
8  # - acc: test accuracy.
9
10 def hard_to_beat(X, Y, X_test, Y_test, W_c):
11     # 1. Compute mean vectors for each class.
12     mus = []
13     for i in range(K):
14         idx = torch.where(Y == i)
15         mus.append(X[idx].mean(dim=0))
16     mus = torch.cat(mus)
17
18     # 2. Estimate the precision matrix using Equation (4).
19     # centered features
20     centered_X = torch.cat([(X[torch.where(Y == i)] - mus[i]) for i in range(K)])
21     cov = torch.cov(centered_X)
22     # compute the precision matrix (inverse covariance)
23     inv_cov = D * torch.inv((N - 1) * cov + trace(cov) * eye(D))
24
25     # 3. Compute weight and bias using Equation (3).
26     W = mus @ inv_cov
27     b = log(1 / K) - 0.5 * einsum('nd, dc, nc -> n', mus, inv_cov, mus)
28
29     # 4. Search the hyperparameter using the validation set.
30     alpha = search_hyperparam(W_c, W, b)
31
32     # 5. Test.
33     test_logits = X_test @ W_c.T + alpha * (X_test @ W.T + b)
34     acc = compute_acc(test_logits, Y_test)
35 return acc

```

---