

Table 1: Comparison with existing methods in the base-to-novel generalization setting on 11 datasets with the 16-shot samples from the base classes. HM: Harmonic mean.

	Base	Novel	HM
ProGrad <sup>1</sup> (ICCV’23)	82.48	70.75	76.16
CoPrompt <sup>2</sup> (ICLR’24)	94.00	77.23	80.48
DeKgTCP <sup>3</sup> (ICLR’25)	84.96	76.38	80.44
TAP <sup>4</sup> (ICLR’25)	84.75	77.63	81.04
MMRL <sup>5</sup> (CVPR’25)	85.68	77.16	81.20
<b>CoLD (Ours)</b>	<b>86.19</b>	<b>78.30</b>	<b>82.06</b>

<sup>1</sup> Zhu et al. “ProGrad: Prompt-aligned Gradient for Prompt Tuning.”

<sup>2</sup> Roy et al. “CoPrompt: Consistency-guided prompt learning for vision-language models.”

<sup>3</sup> Li et al. “DeKgTCP: Divergence-enhanced knowledge-guided context optimization for visual-language prompt tuning.”

<sup>4</sup> Ding et al. “TAP: Tree of Attributes Prompt Learning for Vision-Language Models.”

<sup>5</sup> Guo et al. “MMRL: Multi-Modal Representation Learning for Vision-Language Models.”

Table 2: Comparison of CoLD with existing methods in the domain generalization setting.

	Source	Target				
	ImageNet	-V	-S	-A	-S	Avg.
ProGrad (ICCV’23)	<b>72.24</b>	64.73	47.61	49.39	74.58	59.08
CoPrompt (ICLR’24)	70.80	64.25	49.43	50.50	77.51	60.42
MMA <sup>1</sup> (CVPR’24)	71.00	64.33	49.13	51.12	77.32	60.48
MMRL (CVPR’25)	72.03	64.47	49.13	<b>51.20</b>	77.53	60.58
<b>CoLD (Ours)</b>	71.82	<b>65.31</b>	<b>50.64</b>	51.07	<b>77.72</b>	<b>61.19</b>

<sup>1</sup> Yang et al. “Mma: Multi-modal adapter for vision-language models.”

Table 3: Comparison of CoLD in few shot classification results with 16 shots.

	16-Shot Classification											
	Average	ImageNet	Caltech	Pets	Cars	Flowers	Food	FGVC	SUN397	DTD	EuroSAT	UCF101
CoOp (IJCV’22)	79.89	71.87	95.57	91.87	83.07	97.07	84.20	43.40	74.67	69.87	84.93	82.23
CoCoOp (CVPR’22)	74.90	70.83	95.16	93.34	71.57	87.84	87.25	31.21	72.15	63.04	73.32	78.14
MaPLe (CVPR’23)	81.79	72.33	96.00	92.83	83.57	97.00	85.33	48.40	75.53	71.33	92.33	85.03
PSRC (CVPR’23)	82.87	73.17	96.07	93.67	83.83	97.60	87.50	50.83	77.23	72.73	92.43	86.47
LLaMP (CVPR’24)	83.81	73.49	97.08	94.21	86.07	98.06	<b>87.62</b>	56.07	77.02	74.17	91.31	86.84
TAP (ICLR’25)	83.37	<b>73.76</b>	96.73	93.90	85.37	98.10	87.53	50.43	77.30	74.90	91.90	87.17
MMRL (CVPR’25)	84.34	73.40	97.13	93.83	86.43	98.40	87.03	57.60	<b>77.70</b>	75.30	<b>93.37</b>	87.60
CoLD (Ours)	<b>84.62</b>	73.68	<b>97.17</b>	<b>94.65</b>	<b>87.39</b>	<b>98.72</b>	87.58	<b>58.79</b>	77.11	<b>76.56</b>	91.29	<b>87.91</b>

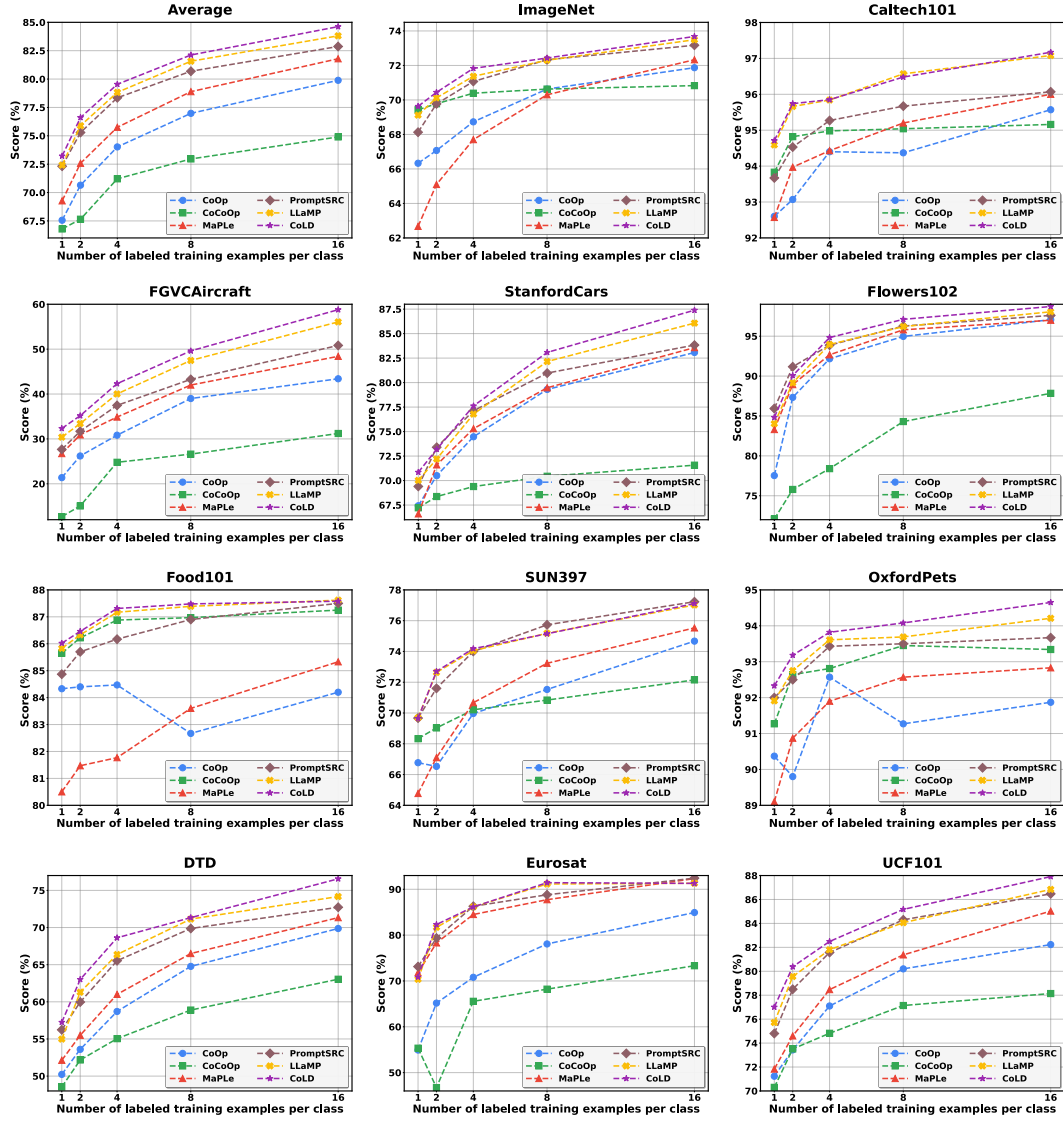


Figure 1: Comparison of CoLD with previous methods on few-shot learning across 11 datasets.

Table 4: Comparison in the cross-dataset evaluation. The model is trained on the entire class of ImageNet (16 shots) and evaluated on the other 10 datasets.

	Source	Target										
	<i>ImageNet</i>	<i>Caltech</i>	<i>Pets</i>	<i>Cars</i>	<i>Flowers</i>	<i>Food</i>	<i>FGVC</i>	<i>SUN397</i>	<i>DTD</i>	<i>EuroSAT</i>	<i>UCF101</i>	<i>Average</i>
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
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
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
ProGrad	72.24	91.52	89.64	62.39	67.87	85.40	20.16	62.47	36.42	43.46	64.29	62.36
CoPrompt	70.80	94.50	90.73	65.67	72.30	86.43	24.00	67.57	47.07	51.90	<b>69.73</b>	67.00
DeKgTCP	<b>72.33</b>	94.73	90.02	65.49	72.39	<b>86.59</b>	25.05	67.19	44.47	51.37	68.78	66.61
TAP	72.30	94.30	90.70	65.60	70.93	86.10	24.57	<b>68.30</b>	<b>50.20</b>	46.00	68.90	66.56
MMRL	72.03	94.67	<b>91.43</b>	66.10	72.77	86.40	26.30	67.57	45.90	53.10	68.27	67.25
CoLD	71.82	<b>94.75</b>	91.13	<b>66.25</b>	<b>72.86</b>	86.21	<b>26.51</b>	67.64	47.23	<b>53.14</b>	68.96	<b>67.47</b>