

# Label Propagation for Zero-shot Classification with Vision-Language Models -Supplementary-

Vladan Stojnić<sup>1</sup>

Yannis Kalantidis<sup>2</sup>

Giorgos Tolias<sup>1</sup>

<sup>1</sup> VRG, FEE, Czech Technical University in Prague

<sup>2</sup> NAVER LABS Europe

**Impact of hyper-parameters** We show the effect of  $\gamma$  and  $\alpha$  in Table 1 and observe stability over a wide range of values for both cases.

**Additional backbones** In Table 2, we present transductive and inductive zero-shot classification results on ImageNet with additional CLIP backbones. We present results with two versions of ViT-L-14 from CLIP [4]. Additionally, we present the results with ViT-B-16 and ViT-H-14 from OpenCLIP [1]<sup>1</sup> trained on the LAION-2B [5] dataset. We see from Table 2 that ZLaP improves the results with different backbones in both transductive and inductive setups. This verifies that ZLaP is not backbone dependant and that it is independent of the dataset used for pre-training.

**Per-dataset results** In Tables 3 and 4, we present per dataset results for transductive and inductive setups, respectively.

**Leveraging LLM generated prompts** In the main paper, we present the average results when we leverage LLM generated prompts from CuPL [2]. In Tables 5 and 6, we present per dataset results for transductive and inductive setups, respectively. CuPL prompts, compared to hand-crafted universal class templates, improve CLIP+ZLaP

<sup>1</sup>[https://github.com/mlfoundations/open\\_clip](https://github.com/mlfoundations/open_clip)

$\alpha \backslash \gamma$	0.01	0.05	0.1	0.3	0.5	0.7	0.8	0.9
0.5	33.6	34.4	35.2	37.9	39.9	42.8	43.9	43.6
1	37.5	38.3	39.2	41.6	43.8	46.0	46.8	45.9
5	50.9	51.3	51.5	52.1	52.5	52.3	51.4	49.7
8	52.2	52.7	52.9	53.7	53.6	52.6	51.5	49.6
10	52.8	52.9	53.2	53.8	53.6	52.4	51.4	49.2
12	52.8	52.8	53.1	53.6	53.3	52.4	51.0	49.0
15	1.5	52.0	52.0	53.0	52.8	51.7	50.4	48.6
20	0.5	0.5	0.5	27.0	48.1	47.5	46.2	44.9

Table 1. **Impact of  $\gamma$  and  $\alpha$  hyper-parameters.** Results presented on CUB dataset for transductive inference.

	Transductive	Inductive
<i>Results with ViT-L-14</i>		
CLIP	75.9	75.9
+ ZLaP	<b>77.2</b>	<b>77.3</b>
<i>Results with ViT-L-14@336</i>		
CLIP	77.0	77.0
+ ZLaP	<b>78.0</b>	<b>78.4</b>
<i>Results with ViT-B-16 (LAION-2B)</i>		
CLIP	70.4	70.4
+ ZLaP	<b>72.0</b>	<b>72.1</b>
<i>Results with ViT-H-14 (LAION-2B)</i>		
CLIP	78.0	78.0
+ ZLaP	<b>79.1</b>	<b>79.1</b>

Table 2. **Accuracy on ImageNet using different CLIP backbones.**

from 60.0% to 64.6% and from 58.7% to 64.2% for the transductive and inductive setup, respectively.

**Web-crawled unlabeled images** We construct a new set of unlabeled images with 10,000 images per class that are chosen either randomly, or based on proximity of their image or text features to the class representation. Results are presented in Table 7. Switching to using only the LAION-based unlabeled set, we observe that random selection fails by performing worse than CLIP, but the other two options provide some improvement, with the caption-based neighbors being a bit better. Interestingly, web-crawling is better than the the target distribution images for the Pets dataset, while much worse for Eurosat due to the lack of satellite images on LAION. On the other hand, if the randomly selected set is mixed with that from the target distribution, ZLaP manages to benefit from the relevant images and to deliver an improvement compared to CLIP.

	imagenet	did	eurosat	fgvca	flowers	food	pets	sun	cars	caltech	cifar10	cifar100	cub	ucf	avg
<i>Results with ResNet50</i>															
TPT <sup>†</sup>	60.7	40.8	28.3	17.6	62.7	74.9	84.5	61.5	58.5	87.0	–	–	–	<b>69.8</b>	(58.8)
CLIP-DN	60.2	41.1	28.4	17.3	63.3	77.2	83.1	60.9	54.8	<u>88.3</u>	74.0	44.7	48.9	60.4	57.3
CLIP	60.3	41.1	26.9	16.7	62.9	76.6	83.1	61.2	54.4	87.9	72.3	42.5	47.0	59.9	56.6
+ <b>ZLaP</b>	<u>61.8</u>	41.9	<b>35.5</b>	17.8	<u>65.9</u>	78.8	83.9	63.3	57.8	<b>89.6</b>	78.2	47.9	52.1	65.9	60.0
(vs CLIP)	<b>↑1.5</b>	<b>↑0.8</b>	<b>↑8.6</b>	<b>↑1.1</b>	<b>↑3.0</b>	<b>↑2.2</b>	<b>↑0.8</b>	<b>↑2.1</b>	<b>↑3.4</b>	<b>↑1.7</b>	<b>↑5.9</b>	<b>↑5.4</b>	<b>↑5.1</b>	<b>↑6.0</b>	<b>↑3.4</b>
InMaP	<b>63.8</b>	<u>44.8</u>	33.4	<b>19.0</b>	65.0	<b>79.4</b>	<u>89.0</u>	<u>65.3</u>	<u>61.5</u>	74.5	<u>78.9</u>	<u>49.6</u>	<b>55.5</b>	65.6	<u>60.4</u>
+ <b>ZLaP</b>	<b>63.8</b>	<b>45.9</b>	<u>34.5</u>	<u>18.4</u>	<b>67.1</b>	<u>79.2</u>	<b>89.2</b>	<b>65.9</b>	<b>62.0</b>	80.7	<b>79.2</b>	<b>49.7</b>	<u>55.3</u>	<u>67.8</u>	<b>61.3</b>
(vs CLIP)	<b>↑3.5</b>	<b>↑4.8</b>	<b>↑7.6</b>	<b>↑1.7</b>	<b>↑4.2</b>	<b>↑2.6</b>	<b>↑6.1</b>	<b>↑4.7</b>	<b>↑7.6</b>	<b>↓7.2</b>	<b>↑6.9</b>	<b>↑7.2</b>	<b>↑8.3</b>	<b>↑7.9</b>	<b>↑4.7</b>
(vs InMaP)	<b>↑0.0</b>	<b>↑1.1</b>	<b>↑1.1</b>	<b>↓0.6</b>	<b>↑2.1</b>	<b>↓0.2</b>	<b>↑0.2</b>	<b>↑0.6</b>	<b>↑0.5</b>	<b>↑6.2</b>	<b>↑0.3</b>	<b>↑0.1</b>	<b>↓0.2</b>	<b>↑2.2</b>	<b>↑0.9</b>
<i>Results with ViT-B/16</i>															
TPT <sup>†</sup>	69.0	47.8	42.4	24.8	69.0	84.7	87.8	65.5	66.9	<b>94.2</b>	–	–	–	68.0	(65.5)
CLIP-DN	68.3	45.7	53.3	24.3	68.0	86.0	87.7	66.5	64.0	93.6	91.4	69.6	56.1	68.4	67.3
CLIP	68.8	45.1	50.2	23.0	67.0	85.7	88.3	66.3	63.8	<u>93.9</u>	91.2	68.7	55.2	67.5	66.8
+ <b>ZLaP</b>	69.7	46.0	57.7	26.3	67.9	87.2	87.9	67.8	66.8	91.8	92.6	<u>70.8</u>	58.2	73.8	68.9
(vs CLIP)	<b>↑0.9</b>	<b>↑0.9</b>	<b>↑7.5</b>	<b>↑3.3</b>	<b>↑0.9</b>	<b>↑1.5</b>	<b>↓0.4</b>	<b>↑1.5</b>	<b>↑3.0</b>	<b>↓2.1</b>	<b>↑1.4</b>	<b>↑2.1</b>	<b>↑3.0</b>	<b>↑6.3</b>	<b>↑2.1</b>
InMaP	<u>72.5</u>	50.9	60.1	<u>28.3</u>	<u>70.8</u>	<b>88.0</b>	<b>93.2</b>	<u>71.3</u>	<u>71.7</u>	76.7	<u>93.3</u>	<b>73.3</b>	63.8	75.7	<u>70.7</u>
+ <b>ZLaP</b>	<b>72.7</b>	<b>51.8</b>	<b>60.9</b>	<b>28.4</b>	<b>73.4</b>	<u>87.9</u>	<u>92.8</u>	<b>71.9</b>	<b>72.1</b>	83.7	<b>93.6</b>	<b>73.3</b>	<b>64.1</b>	<b>77.7</b>	<b>71.7</b>
(vs CLIP)	<b>↑3.9</b>	<b>↑6.7</b>	<b>↑10.7</b>	<b>↑5.4</b>	<b>↑6.4</b>	<b>↑2.2</b>	<b>↑4.5</b>	<b>↑5.6</b>	<b>↑8.3</b>	<b>↓10.2</b>	<b>↑2.4</b>	<b>↑4.6</b>	<b>↑8.9</b>	<b>↑10.2</b>	<b>↑4.9</b>
(vs InMaP)	<b>↑0.2</b>	<b>↑0.9</b>	<b>↑0.8</b>	<b>↑0.1</b>	<b>↑2.6</b>	<b>↓0.1</b>	<b>↓0.4</b>	<b>↑0.6</b>	<b>↑0.4</b>	<b>↑7.0</b>	<b>↑0.3</b>	<b>↑0.0</b>	<b>↑0.3</b>	<b>↑2.0</b>	<b>↑1.0</b>

Table 3. **Transductive zero-shot classification accuracy on 14 datasets** for two CLIP backbones. Rows denoted as (vs CLIP) and (vs InMaP) show the absolute accuracy gains of our method over CLIP and InMaP, respectively. <sup>†</sup> denotes numbers taken from InMaP [3].

	imagenet	dtd	eurosat	fgvca	flowers	food	pets	sun	cars	caltech	cifar10	cifar100	cut	ucf	avg
<i>Results with ResNet50</i>															
TPT <sup>†</sup>	60.7	40.8	28.3	17.6	62.7	74.9	<u>84.5</u>	61.5	58.5	87.0	–	–	–	<b>69.8</b>	(58.8)
CLIP-DN	60.2	41.2	28.3	17.2	63.3	77.2	83.3	60.8	54.9	<b>88.3</b>	74.0	44.7	48.9	60.4	57.3
CLIP	60.3	41.1	26.9	16.7	62.9	76.7	83.1	61.2	54.5	<u>87.9</u>	72.3	42.5	47.0	59.9	56.6
+ <b>ZLaP</b>	62.2	42.8	31.9	17.4	<u>69.3</u>	77.9	80.3	61.8	56.4	86.9	76.3	46.0	49.7	62.8	58.7
(vs CLIP)	<b>↑1.9</b>	<b>↑1.7</b>	<b>↑5.0</b>	<b>↑0.7</b>	<b>↑6.4</b>	<b>↑1.2</b>	<b>↓2.8</b>	<b>↑0.6</b>	<b>↑1.9</b>	<b>↓1.0</b>	<b>↑4.0</b>	<b>↑3.5</b>	<b>↑2.7</b>	<b>↑2.9</b>	<b>↑2.1</b>
+ <b>ZLaP*</b>	<u>62.9</u>	43.1	<b>38.8</b>	17.9	68.8	78.3	77.7	61.2	55.8	86.3	78.6	48.0	51.1	64.2	59.5
(vs CLIP)	<b>↑2.6</b>	<b>↑2.0</b>	<b>↑11.9</b>	<b>↑1.2</b>	<b>↑5.9</b>	<b>↑1.6</b>	<b>↓5.4</b>	<b>↑0.0</b>	<b>↑1.3</b>	<b>↓1.6</b>	<b>↑6.3</b>	<b>↑5.5</b>	<b>↑4.1</b>	<b>↑4.3</b>	<b>↑2.9</b>
InMaP	<u>62.9</u>	45.7	33.6	<b>19.2</b>	66.4	<b>79.2</b>	<b>85.7</b>	<u>65.0</u>	<b>62.0</b>	76.2	79.0	49.7	<b>55.4</b>	66.0	60.4
+ <b>ZLaP</b>	<u>62.9</u>	<b>46.6</b>	<u>36.3</u>	18.7	69.1	<u>79.0</u>	83.4	64.9	61.8	80.2	<u>79.1</u>	<b>50.6</b>	54.9	66.5	<b>61.0</b>
(vs CLIP)	<b>↑2.6</b>	<b>↑5.5</b>	<b>↑9.4</b>	<b>↑2.0</b>	<b>↑6.2</b>	<b>↑2.3</b>	<b>↑0.3</b>	<b>↑3.7</b>	<b>↑7.3</b>	<b>↓7.7</b>	<b>↑6.8</b>	<b>↑8.1</b>	<b>↑7.9</b>	<b>↑6.6</b>	<b>↑4.4</b>
(vs InMaP)	<b>↑0.0</b>	<b>↑0.9</b>	<b>↑2.7</b>	<b>↓0.5</b>	<b>↑2.7</b>	<b>↓0.2</b>	<b>↓2.3</b>	<b>↓0.1</b>	<b>↓0.2</b>	<b>↑4.0</b>	<b>↑0.1</b>	<b>↑0.9</b>	<b>↓0.5</b>	<b>↑0.5</b>	<b>↑0.6</b>
+ <b>ZLaP*</b>	<b>63.0</b>	<u>46.3</u>	36.2	<u>18.9</u>	<b>69.4</b>	<b>79.2</b>	81.4	<b>65.1</b>	<u>61.9</u>	79.3	<b>79.2</b>	<u>50.5</u>	<u>55.1</u>	<u>67.0</u>	<u>60.9</u>
(vs CLIP)	<b>↑2.7</b>	<b>↑5.2</b>	<b>↑9.3</b>	<b>↑2.2</b>	<b>↑6.5</b>	<b>↑2.5</b>	<b>↓1.7</b>	<b>↑3.9</b>	<b>↑7.4</b>	<b>↓8.6</b>	<b>↑6.9</b>	<b>↑8.0</b>	<b>↑8.1</b>	<b>↑7.1</b>	<b>↑4.3</b>
(vs InMaP)	<b>↑0.1</b>	<b>↑0.6</b>	<b>↑2.6</b>	<b>↓0.3</b>	<b>↑3.0</b>	<b>↑0.0</b>	<b>↓4.3</b>	<b>↑0.1</b>	<b>↓0.1</b>	<b>↑3.1</b>	<b>↑0.2</b>	<b>↑0.8</b>	<b>↓0.3</b>	<b>↑1.0</b>	<b>↑0.5</b>
<i>Results with ViT-B/16</i>															
TPT <sup>†</sup>	69.0	47.8	42.4	24.8	69.0	84.7	87.8	65.5	66.9	<b>94.2</b>	–	–	–	68.0	(65.5)
CLIP-DN	68.3	45.6	53.3	24.3	67.9	86.0	87.7	66.5	64.1	93.6	91.5	69.6	56.0	68.4	67.3
CLIP	68.8	45.1	50.2	23.0	67.0	85.7	88.3	66.3	63.8	<u>93.9</u>	91.2	68.7	55.2	67.5	66.8
+ <b>ZLaP</b>	70.2	48.6	55.6	25.4	73.5	86.9	87.1	67.4	65.6	93.1	92.2	71.0	59.4	71.5	69.1
(vs CLIP)	<b>↑1.4</b>	<b>↑3.5</b>	<b>↑5.4</b>	<b>↑2.4</b>	<b>↑6.5</b>	<b>↑1.2</b>	<b>↓1.2</b>	<b>↑1.1</b>	<b>↑1.8</b>	<b>↓0.8</b>	<b>↑1.0</b>	<b>↑2.3</b>	<b>↑4.2</b>	<b>↑4.0</b>	<b>↑2.3</b>
+ <b>ZLaP*</b>	71.0	49.1	58.2	25.8	72.6	87.3	86.3	67.2	66.1	92.1	92.7	72.0	59.1	72.2	69.4
(vs CLIP)	<b>↑2.2</b>	<b>↑4.0</b>	<b>↑8.0</b>	<b>↑2.8</b>	<b>↑5.6</b>	<b>↑1.6</b>	<b>↓2.0</b>	<b>↑0.9</b>	<b>↑2.3</b>	<b>↓1.8</b>	<b>↑1.5</b>	<b>↑3.3</b>	<b>↑3.9</b>	<b>↑4.7</b>	<b>↑2.6</b>
InMaP	<u>72.0</u>	49.6	59.4	<u>29.0</u>	71.9	<b>87.9</b>	<b>91.6</b>	<b>71.4</b>	<b>71.9</b>	79.0	<u>93.4</u>	73.7	63.9	<u>75.4</u>	70.7
+ <b>ZLaP</b>	<b>72.1</b>	<b>51.2</b>	<b>63.2</b>	<b>29.1</b>	<b>75.9</b>	<u>87.8</u>	<u>90.0</u>	<u>71.0</u>	71.2	84.0	<u>93.4</u>	<u>74.0</u>	<b>64.3</b>	<b>76.3</b>	<b>71.7</b>
(vs CLIP)	<b>↑3.3</b>	<b>↑6.1</b>	<b>↑13.0</b>	<b>↑6.1</b>	<b>↑8.9</b>	<b>↑2.1</b>	<b>↑1.7</b>	<b>↑4.7</b>	<b>↑7.4</b>	<b>↓9.9</b>	<b>↑2.2</b>	<b>↑5.3</b>	<b>↑9.1</b>	<b>↑8.8</b>	<b>↑4.9</b>
(vs InMaP)	<b>↑0.1</b>	<b>↑1.6</b>	<b>↑3.8</b>	<b>↑0.1</b>	<b>↑4.0</b>	<b>↓0.1</b>	<b>↓1.6</b>	<b>↓0.4</b>	<b>↓0.7</b>	<b>↑5.0</b>	<b>↑0.0</b>	<b>↑0.3</b>	<b>↑0.4</b>	<b>↑0.9</b>	<b>↑1.0</b>
+ <b>ZLaP*</b>	<b>72.1</b>	<u>51.0</u>	<u>62.7</u>	<u>29.0</u>	<u>75.5</u>	<b>87.9</b>	89.0	<b>71.4</b>	<u>71.8</u>	83.1	<b>93.6</b>	<b>74.2</b>	<u>64.2</u>	<b>76.3</b>	<u>71.6</u>
(vs CLIP)	<b>↑3.3</b>	<b>↑5.9</b>	<b>↑12.5</b>	<b>↑6.0</b>	<b>↑8.5</b>	<b>↑2.2</b>	<b>↑0.7</b>	<b>↑5.1</b>	<b>↑8.0</b>	<b>↓10.8</b>	<b>↑2.4</b>	<b>↑5.5</b>	<b>↑9.0</b>	<b>↑8.8</b>	<b>↑4.8</b>
(vs InMaP)	<b>↑0.1</b>	<b>↑1.4</b>	<b>↑3.3</b>	<b>↑0.0</b>	<b>↑3.6</b>	<b>↑0.0</b>	<b>↓2.6</b>	<b>↑0.0</b>	<b>↓0.1</b>	<b>↑4.1</b>	<b>↑0.2</b>	<b>↑0.5</b>	<b>↑0.3</b>	<b>↑0.9</b>	<b>↑0.9</b>

Table 4. **Inductive zero-shot classification accuracy on 14 datasets** for two CLIP backbones. Rows denoted as (vs CLIP) and (vs InMaP) show the absolute accuracy gains of our method over CLIP and InMaP, respectively. \* denotes our method with approximation of  $\hat{Y}$ . <sup>†</sup> denotes numbers taken from InMaP [3].

	imagenet	dtd	fgvca	flowers	food	pets	sun	cars	caltech	cifar10	cifar100	ucf	avg
<i>Results with ResNet50</i>													
CLIP	61.7	49.1	18.5	<u>67.9</u>	77.8	87.5	63.8	55.8	<u>88.7</u>	76.4	45.2	63.5	63.0
+ ZLaP	62.7	51.4	<u>20.2</u>	67.6	78.9	88.1	65.2	58.8	<b>89.8</b>	77.6	47.4	67.8	64.6
(vs CLIP)	$\uparrow 1.0$	$\uparrow 2.3$	$\uparrow 1.7$	$\downarrow 0.3$	$\uparrow 1.1$	$\uparrow 0.6$	$\uparrow 1.4$	$\uparrow 3.0$	$\uparrow 1.1$	$\uparrow 1.2$	$\uparrow 2.2$	$\uparrow 4.3$	$\uparrow 1.6$
InMaP	<b>64.4</b>	<u>54.5</u>	<b>22.2</b>	67.2	<b>79.3</b>	<b>89.9</b>	<u>67.4</u>	<u>62.8</u>	73.7	<u>78.2</u>	<u>50.2</u>	<u>68.2</u>	<u>64.8</u>
+ ZLaP	<u>64.3</u>	<b>55.6</b>	<b>22.2</b>	<b>69.8</b>	<u>79.2</u>	<u>89.5</u>	<b>67.8</b>	<b>63.2</b>	78.9	<b>78.9</b>	<b>50.5</b>	<b>70.2</b>	<b>65.8</b>
(vs CLIP)	$\uparrow 2.6$	$\uparrow 6.5$	$\uparrow 3.7$	$\uparrow 1.9$	$\uparrow 1.4$	$\uparrow 2.0$	$\uparrow 4.0$	$\uparrow 7.4$	$\downarrow 9.8$	$\uparrow 2.5$	$\uparrow 5.3$	$\uparrow 6.7$	$\uparrow 2.8$
(vs InMaP)	$\downarrow 0.1$	$\uparrow 1.1$	$\uparrow 0.0$	$\uparrow 2.6$	$\downarrow 0.1$	$\downarrow 0.4$	$\uparrow 0.4$	$\uparrow 0.4$	$\uparrow 5.2$	$\uparrow 0.7$	$\uparrow 0.3$	$\uparrow 2.0$	$\uparrow 1.0$
<i>Results with ViT-B/16</i>													
CLIP	70.0	53.2	27.9	73.4	86.3	91.7	69.5	66.1	<b>94.4</b>	90.7	69.4	70.5	71.9
+ ZLaP	<u>70.5</u>	54.0	30.1	72.2	86.9	91.8	69.7	<u>67.3</u>	<u>92.7</u>	92.4	69.9	74.0	72.6
(vs CLIP)	$\uparrow 0.5$	$\uparrow 0.8$	$\uparrow 2.2$	$\downarrow 1.2$	$\uparrow 0.6$	$\uparrow 0.1$	$\uparrow 0.2$	$\uparrow 1.2$	$\downarrow 1.7$	$\uparrow 1.7$	$\uparrow 0.5$	$\uparrow 3.5$	$\uparrow 0.7$
InMaP	<b>73.3</b>	<u>57.3</u>	<b>31.9</b>	<u>74.1</u>	<b>88.1</b>	<b>93.7</b>	<u>73.3</u>	<b>72.8</b>	78.0	<u>93.4</u>	<b>73.3</b>	<u>77.1</u>	73.9
+ ZLaP	<b>73.3</b>	<b>57.9</b>	<u>31.7</u>	<b>76.9</b>	<u>88.0</u>	<u>93.3</u>	<b>73.7</b>	<b>72.8</b>	83.3	<b>93.6</b>	<u>73.2</u>	<b>79.5</b>	<b>74.8</b>
(vs CLIP)	$\uparrow 3.3$	$\uparrow 4.7$	$\uparrow 3.8$	$\uparrow 3.5$	$\uparrow 1.7$	$\uparrow 1.6$	$\uparrow 4.2$	$\uparrow 6.7$	$\downarrow 11.1$	$\uparrow 2.9$	$\uparrow 3.8$	$\uparrow 9.0$	$\uparrow 2.9$
(vs InMaP)	$\uparrow 0.0$	$\uparrow 0.6$	$\downarrow 0.2$	$\uparrow 2.8$	$\downarrow 0.1$	$\downarrow 0.4$	$\uparrow 0.4$	$\uparrow 0.0$	$\uparrow 5.3$	$\uparrow 0.2$	$\downarrow 0.1$	$\uparrow 2.4$	$\uparrow 0.9$

Table 5. **Transductive zero-shot classification accuracy on 12 datasets** for two CLIP backbones and prompts generated by a LLM [2]. Rows denoted as (vs CLIP) and (vs InMaP) show the absolute accuracy gains of our method over CLIP and InMaP, respectively.

	imagenet	dtd	fgvca	flowers	food	pets	sun	cars	caltech	cifar10	cifar100	ucf	avg
<i>Results with ResNet50</i>													
CLIP	61.7	49.1	18.5	67.9	77.8	<b>87.5</b>	63.8	55.8	<b>88.7</b>	76.4	45.2	63.5	63.0
+ ZLaP	<u>63.1</u>	51.4	<u>20.0</u>	<b>72.7</b>	<u>78.4</u>	85.4	63.3	57.8	<u>88.3</u>	77.9	48.0	63.6	64.2
(vs CLIP)	$\uparrow 1.4$	$\uparrow 2.3$	$\uparrow 1.5$	$\uparrow 4.8$	$\uparrow 0.6$	$\downarrow 2.1$	$\downarrow 0.5$	$\uparrow 2.0$	$\downarrow 0.4$	$\uparrow 1.5$	$\uparrow 2.8$	$\uparrow 0.1$	$\uparrow 1.2$
InMaP	<b>63.4</b>	<b>54.6</b>	<b>22.6</b>	68.8	<b>79.1</b>	<u>86.4</u>	<b>66.6</b>	<b>62.5</b>	75.7	<u>78.2</u>	<u>50.4</u>	<u>67.5</u>	<u>64.6</u>
+ ZLaP	<b>63.4</b>	<u>54.1</u>	<b>22.6</b>	<u>71.5</u>	<b>79.1</b>	83.4	<u>66.5</u>	<u>62.4</u>	79.4	<b>78.8</b>	<b>51.0</b>	<b>67.7</b>	<b>65.0</b>
(vs CLIP)	$\uparrow 1.7$	$\uparrow 5.0$	$\uparrow 4.1$	$\uparrow 3.6$	$\uparrow 1.3$	$\downarrow 4.1$	$\uparrow 2.7$	$\uparrow 6.6$	$\downarrow 9.3$	$\uparrow 2.4$	$\uparrow 5.8$	$\uparrow 4.2$	$\uparrow 2.0$
(vs InMaP)	$\uparrow 0.0$	$\downarrow 0.5$	$\uparrow 0.0$	$\uparrow 2.7$	$\uparrow 0.0$	$\downarrow 3.0$	$\downarrow 0.1$	$\downarrow 0.1$	$\uparrow 3.7$	$\uparrow 0.6$	$\uparrow 0.6$	$\uparrow 0.2$	$\uparrow 0.4$
<i>Results with ViT-B/16</i>													
CLIP	70.0	53.2	27.9	73.4	86.3	<u>91.7</u>	69.5	66.1	<b>94.4</b>	90.7	69.4	70.5	71.9
+ ZLaP	71.2	55.5	<u>29.8</u>	<u>77.7</u>	87.2	91.1	69.7	<u>67.5</u>	<b>94.4</b>	91.6	71.3	<u>72.6</u>	73.3
(vs CLIP)	$\uparrow 1.2$	$\uparrow 2.3$	$\uparrow 1.9$	$\uparrow 4.3$	$\uparrow 0.9$	$\downarrow 0.6$	$\uparrow 0.2$	$\uparrow 1.4$	$\uparrow 0.0$	$\uparrow 0.9$	$\uparrow 1.9$	$\uparrow 2.1$	$\uparrow 1.4$
InMaP	<u>72.4</u>	<b>57.2</b>	<b>32.8</b>	75.8	<b>88.0</b>	<b>92.3</b>	<b>73.0</b>	<b>72.9</b>	79.8	<u>93.3</u>	<u>73.7</u>	<b>76.6</b>	<u>74.0</u>
+ ZLaP	<b>72.5</b>	<u>56.3</u>	<b>32.8</b>	<b>78.5</b>	<u>87.9</u>	89.6	<u>72.6</u>	<b>72.9</b>	<u>83.9</u>	<b>93.5</b>	<b>73.8</b>	<b>76.6</b>	<b>74.2</b>
(vs CLIP)	$\uparrow 2.5$	$\uparrow 3.1$	$\uparrow 4.9$	$\uparrow 5.1$	$\uparrow 1.6$	$\downarrow 2.1$	$\uparrow 3.1$	$\uparrow 6.8$	$\downarrow 10.5$	$\uparrow 2.8$	$\uparrow 4.4$	$\uparrow 6.1$	$\uparrow 2.3$
(vs InMaP)	$\uparrow 0.1$	$\downarrow 0.9$	$\uparrow 0.0$	$\uparrow 2.7$	$\downarrow 0.1$	$\downarrow 2.7$	$\downarrow 0.4$	$\uparrow 0.0$	$\uparrow 4.1$	$\uparrow 0.2$	$\uparrow 0.1$	$\uparrow 0.0$	$\uparrow 0.2$

Table 6. **Inductive zero-shot classification accuracy on 12 datasets** for two CLIP backbones and prompts generated by a LLM [2]. Rows denoted as (vs CLIP) and (vs InMaP) show the absolute accuracy gains of our method over CLIP and InMaP, respectively.

	imagenet	dtd	eurosat	fgvca	flowers	food	pets	sun	cars	caltech	cifar10	cifar100	cub	ucf	avg
<i>Results with ResNet50</i>															
CLIP	60.3	41.1	26.9	16.7	62.9	76.7	83.1	61.2	54.5	<u>87.9</u>	72.3	42.5	47.0	59.9	56.6
+ <b>ZLaP</b> (target distribution)	<b>62.2</b>	<b>42.8</b>	<b>31.9</b>	<b>17.4</b>	<b>69.3</b>	<b>77.9</b>	80.3	61.8	<b>56.4</b>	86.9	<b>76.3</b>	<b>46.0</b>	<b>49.7</b>	<b>62.8</b>	<b>58.7</b>
+ <b>ZLaP</b> (target distr. + LAION random)	<u>61.4</u>	<u>42.3</u>	<u>30.2</u>	15.7	<u>63.5</u>	<u>77.1</u>	80.3	61.6	53.7	87.8	<u>75.1</u>	<u>42.9</u>	47.5	59.8	<u>57.1</u>
+ <b>ZLaP</b> (LAION random)	59.9	41.4	26.2	14.3	59.1	74.5	79.3	61.1	51.4	87.6	70.6	41.4	43.4	58.8	54.6
+ <b>ZLaP</b> (LAION image neighbors)	60.6	41.1	29.1	16.7	<u>63.5</u>	76.9	<u>83.5</u>	<u>61.9</u>	54.7	<b>88.4</b>	69.5	41.1	48.2	59.6	56.8
+ <b>ZLaP</b> (LAION caption neighbors)	60.7	40.5	26.9	<u>16.9</u>	63.0	76.9	<b>83.6</b>	<b>62.0</b>	<u>55.3</u>	<b>88.4</b>	73.0	41.7	<u>48.5</u>	<u>60.1</u>	57.0
<i>Results with ViT-B/16</i>															
CLIP	68.8	45.1	50.2	23.0	67.0	85.7	<u>88.3</u>	66.3	63.8	93.9	91.2	68.7	55.2	67.5	66.8
+ <b>ZLaP</b> (target distribution)	<b>70.2</b>	<b>48.6</b>	<b>55.6</b>	<b>25.4</b>	<b>73.5</b>	<b>86.9</b>	87.1	<b>67.4</b>	<b>65.6</b>	93.1	<b>92.2</b>	<b>71.0</b>	<b>59.4</b>	<b>71.5</b>	<b>69.1</b>
+ <b>ZLaP</b> (target distr. + LAION random)	<u>69.5</u>	<u>45.9</u>	<u>53.1</u>	21.0	67.3	<u>86.3</u>	86.4	66.9	<u>64.7</u>	93.7	<u>91.9</u>	<u>69.3</u>	55.6	<u>67.6</u>	<u>67.1</u>
+ <b>ZLaP</b> (LAION random)	68.6	44.9	49.4	19.8	65.3	85.3	86.8	66.5	63.0	93.6	90.3	68.6	54.0	66.9	65.9
+ <b>ZLaP</b> (LAION image neighbors)	69.0	45.4	49.2	<u>23.8</u>	<u>68.1</u>	85.8	<b>88.4</b>	66.9	64.5	<b>94.2</b>	90.8	68.1	<u>56.7</u>	<u>67.6</u>	67.0
+ <b>ZLaP</b> (LAION caption neighbors)	69.1	45.0	49.4	23.4	<u>68.1</u>	85.9	<b>88.4</b>	<u>67.0</u>	64.6	<u>94.0</u>	90.8	68.5	<u>56.7</u>	<u>67.6</u>	<u>67.1</u>

Table 7. **Inductive zero-shot classification accuracy on 14 datasets using different sources of unlabeled data.** Compared to the original experiments that use unlabeled images from the target distribution, LAION-400M is used to create a web-crawled unlabeled set.

## References

- [1] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *CVPR*, 2023. [1](#)
- [2] Sarah M. Pratt, Rosanne Liu, and Ali Farhadi. What does a platypus look like? generating customized prompts for zero-shot image classification. In *ICCV*, 2023. [1](#), [4](#)
- [3] Qi Qian, Yuanhong Xu, and Juhua Hu. Intra-modal proxy learning for zero-shot visual categorization with clip. In *NeurIPS*, 2023. [2](#), [3](#)
- [4] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021. [1](#)
- [5] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5B: an open large-scale dataset for training next generation image-text models. In *NeurIPS*, 2022. [1](#)