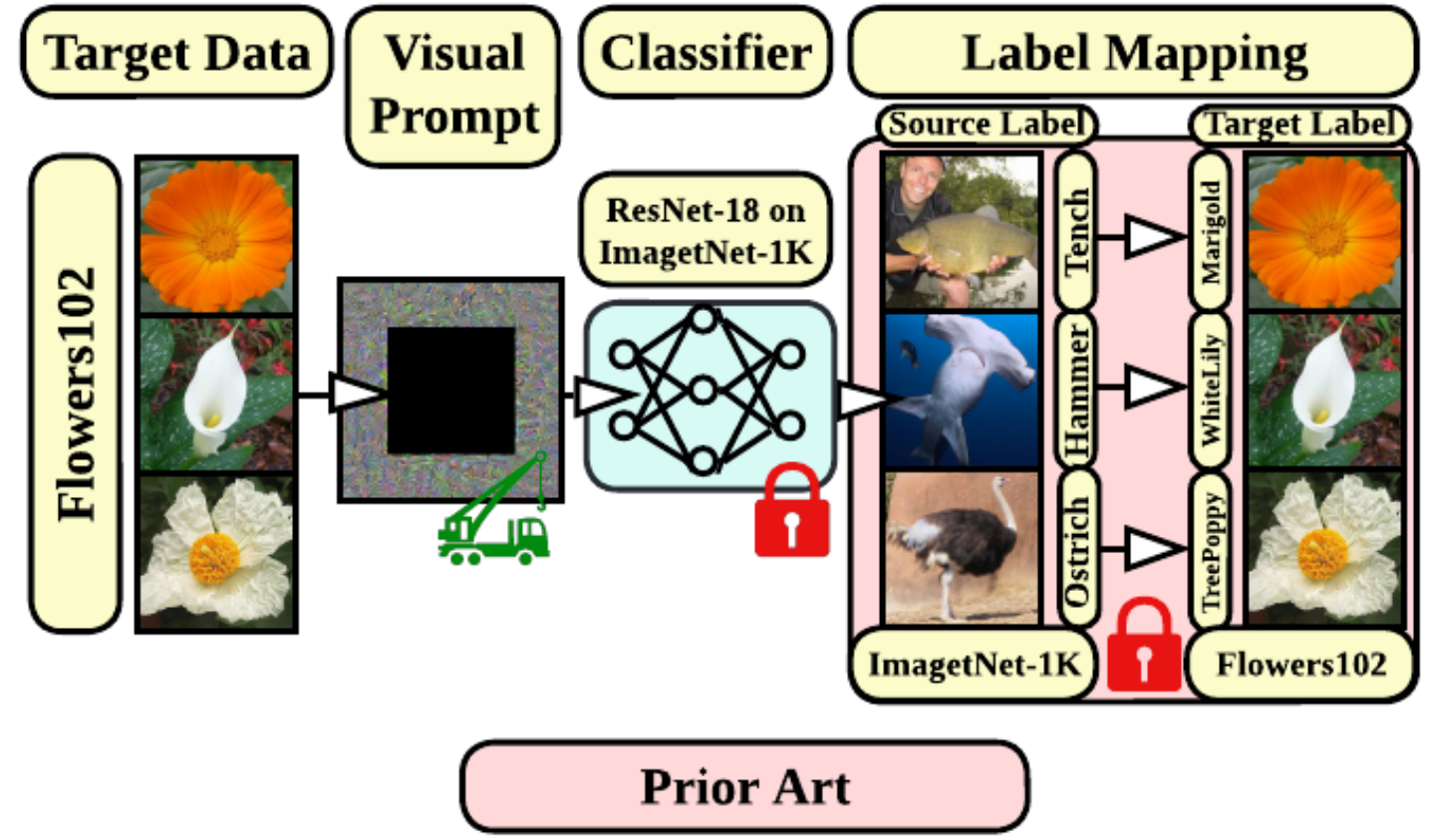


# Understanding and Improving Visual Prompting: A Label- Mapping Perspective

**Aochuan Chen (MSU), Yuguang Yao (MSU), Pin-Yu Chen (IBM  
Research), Yihua Zhang (MSU), Sijia Liu (MSU)**

# Visual Prompting for Transfer Learning



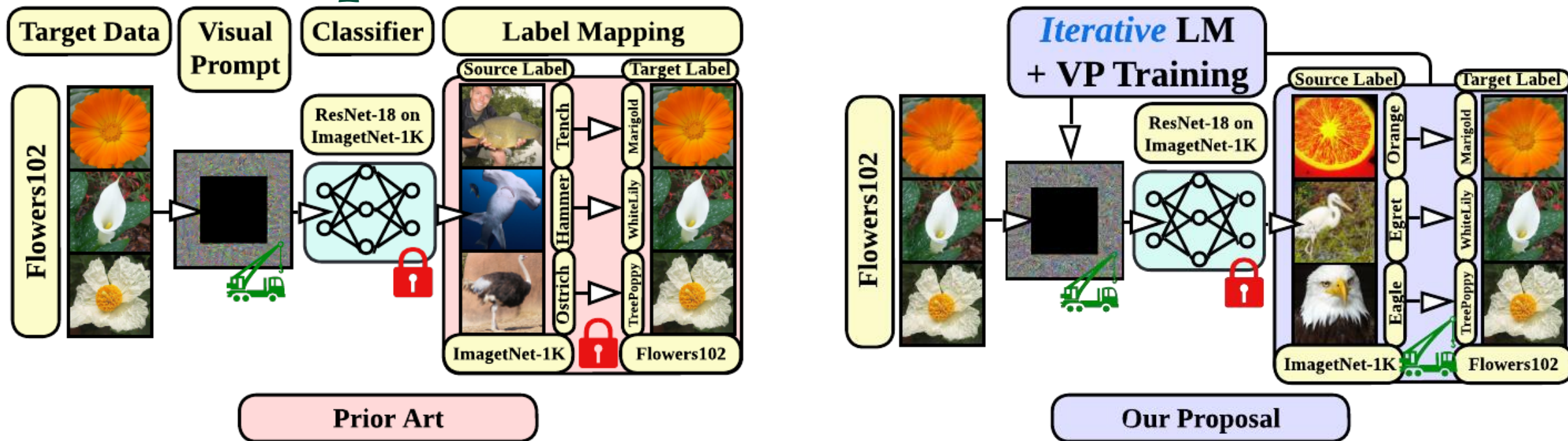
# Label Mapping

- ❖ Mapping from source classes to downstream classes.
- ❖ Existing label mapping methods seem **ruleless**.
- ❖ A key building block for visual prompting.

## Main Research Question

Given the source model, how to build a mapping from the source label space to the target label space so that the model's prediction directs to the correct target label?

# Our Proposal: BLO based ILM-VP



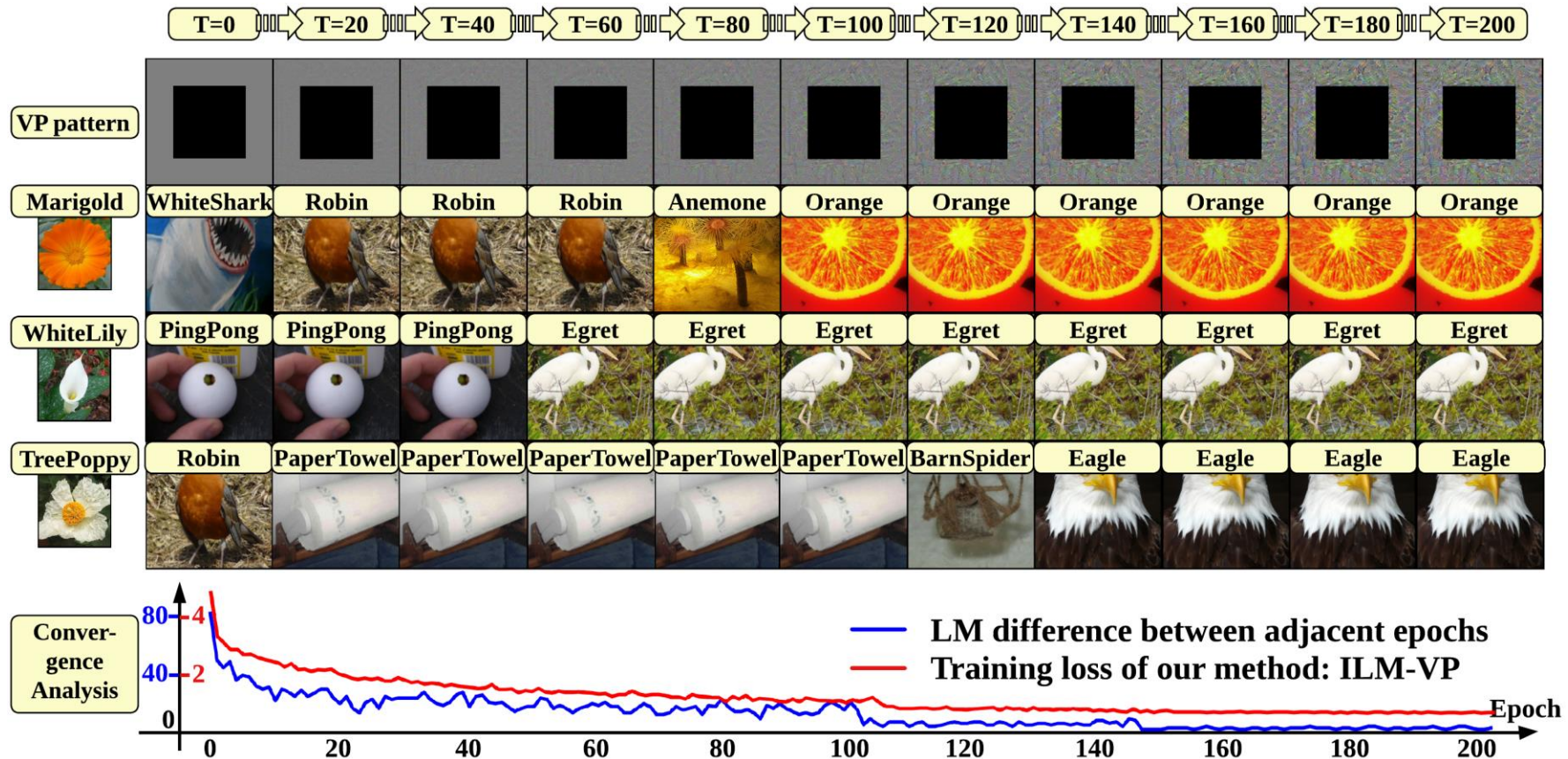

---

## Algorithm 2 Bi-level optimization based ILM-VP algorithm

---

- 1: **Initialize:** Given target training set  $\mathcal{T}_{tr}$ , pre-trained model  $f_{\theta_s}$ , prompt pattern initialization  $\delta_0$ , and upper-level learning rate  $\lambda$  for SGD
  - 2: **for** Epoch  $n = 0, 1, \dots$ , **do**
  - 3:     **Lower-level label mapping:** Given  $\delta_{n-1}$ , call LM for each target class  $y_t$  in  $\mathcal{T}_{tr}$
  - 4:     **Upper-level prompt learning:** Given LM, call SGD to update prompt  $\delta_n \leftarrow \delta_{n-1}$
  - 5: **end for**
-

# Convergence of ILM-VP



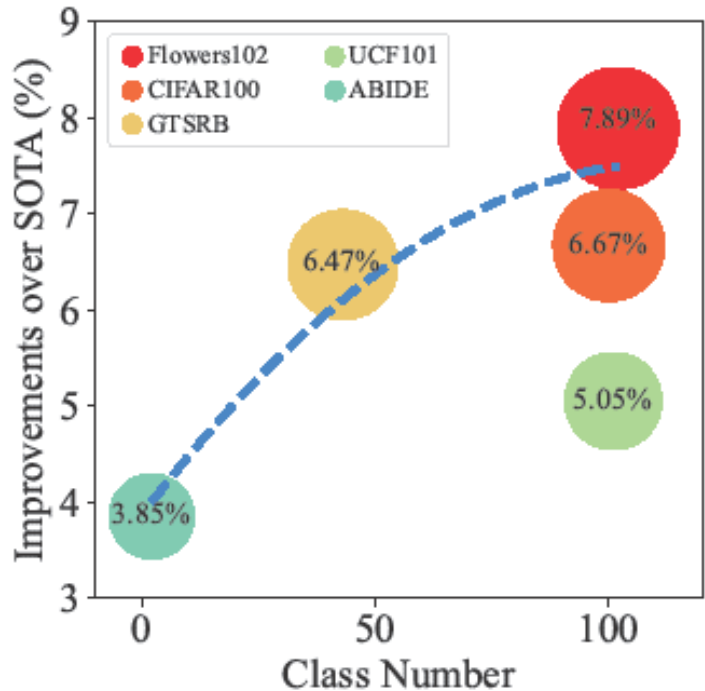
**Figure 1.** ILM-VP training dynamics from epoch 0 to 200. Rows show: (1) VP pattern vs. epoch number; (2-4) Learned source label mapping with respect to target label ‘Marigold’, ‘White Lily’, and ‘Tree Poppy’, together with explanation-by-example-identified source training examples to explain each re-purposed target label; (5) Convergence of training loss and LM difference between adjacent epochs measured by Hamming distance.

# Accuracy Improvements

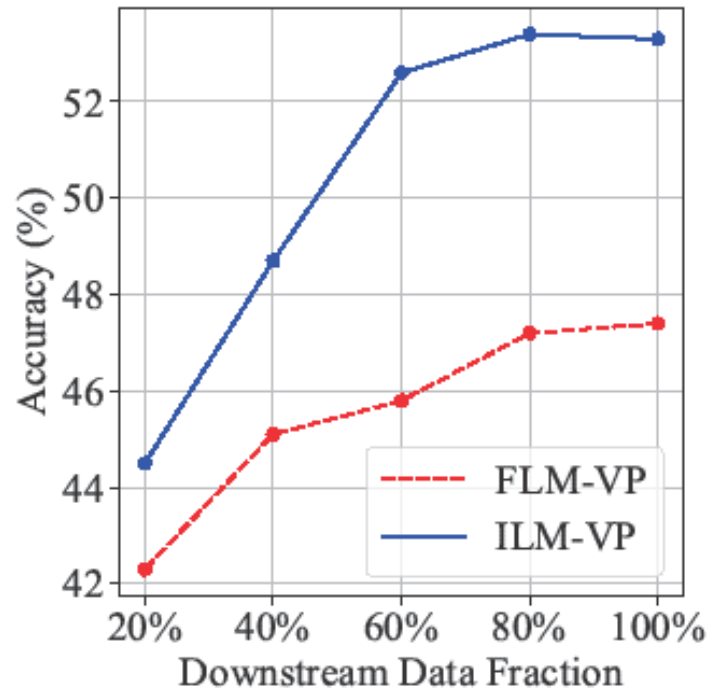
Source Model	ResNet-18 (ImageNet-1K)					ResNet-50 (ImageNet-1K)				ResNeXt-101-32x8d (Instagram)			
Method	<u>Ours</u> ILM-VP	Prompt baseline		Finetuning		<u>Ours</u> ILM-VP	Prompt base.	Finetuning		<u>Ours</u> ILM-VP	Prompt base.	Finetuning	
Parameter Size	0.05M	0.05M	0.05M	0.51M	11.7M	0.05M	0.05M	0.51M	25.6M	0.05M	0.05M	0.51M	88.8M
Flowers102	<b>27.9</b> $\pm$ 0.7	11.0 $\pm$ 0.5	20.0 $\pm$ 0.3	88.0 $\pm$ 0.5	97.1 $\pm$ 0.7	<b>24.6</b> $\pm$ 0.6	20.3 $\pm$ 0.3	90.9 $\pm$ 0.4	97.9 $\pm$ 0.7	<b>27.9</b> $\pm$ 0.3	22.5 $\pm$ 0.5	89.1 $\pm$ 0.2	99.2 $\pm$ 0.5
DTD	<b>35.3</b> $\pm$ 0.9	16.3 $\pm$ 0.7	32.4 $\pm$ 0.5	60.0 $\pm$ 0.6	65.5 $\pm$ 0.9	<b>40.5</b> $\pm$ 0.5	36.9 $\pm$ 0.8	67.6 $\pm$ 0.3	69.7 $\pm$ 0.9	<b>41.4</b> $\pm$ 0.7	40.3 $\pm$ 0.5	69.7 $\pm$ 0.2	69.1 $\pm$ 1.0
UCF101	<b>23.9</b> $\pm$ 0.5	6.6 $\pm$ 0.4	18.9 $\pm$ 0.5	63.2 $\pm$ 0.8	73.0 $\pm$ 0.6	<b>34.6</b> $\pm$ 0.2	33.9 $\pm$ 0.4	70.8 $\pm$ 0.3	78.0 $\pm$ 0.8	<b>43.1</b> $\pm$ 0.8	41.9 $\pm$ 0.6	76.9 $\pm$ 0.5	79.1 $\pm$ 0.7
Food101	<b>14.8</b> $\pm$ 0.2	3.8 $\pm$ 0.3	12.8 $\pm$ 0.1	50.6 $\pm$ 0.3	75.4 $\pm$ 0.8	<b>17.0</b> $\pm$ 0.3	15.3 $\pm$ 0.2	57.6 $\pm$ 0.5	80.3 $\pm$ 0.9	<b>23.0</b> $\pm$ 0.4	20.5 $\pm$ 0.5	76.0 $\pm$ 0.4	82.5 $\pm$ 0.3
GTSRB	<b>52.0</b> $\pm$ 1.2	46.1 $\pm$ 1.3	45.5 $\pm$ 1.0	77.4 $\pm$ 1.2	98.0 $\pm$ 0.3	<b>52.5</b> $\pm$ 1.4	47.6 $\pm$ 1.1	77.8 $\pm$ 0.7	97.6 $\pm$ 1.0	<b>59.9</b> $\pm$ 1.0	56.2 $\pm$ 0.6	73.5 $\pm$ 0.7	97.6 $\pm$ 0.9
EuroSAT	<b>85.2</b> $\pm$ 0.6	82.4 $\pm$ 0.4	83.8 $\pm$ 0.2	93.8 $\pm$ 0.3	98.8 $\pm$ 0.5	83.6 $\pm$ 0.7	<b>84.8</b> $\pm$ 0.3	95.7 $\pm$ 0.2	98.9 $\pm$ 0.6	86.2 $\pm$ 0.8	<b>87.8</b> $\pm$ 0.4	93.4 $\pm$ 0.3	98.9 $\pm$ 0.7
OxfordPets	<b>65.4</b> $\pm$ 0.7	9.3 $\pm$ 0.4	62.9 $\pm$ 0.1	87.2 $\pm$ 0.6	87.8 $\pm$ 0.5	76.2 $\pm$ 0.6	<b>76.4</b> $\pm$ 0.2	90.4 $\pm$ 0.3	91.9 $\pm$ 0.4	<b>78.9</b> $\pm$ 0.8	76.8 $\pm$ 0.6	93.6 $\pm$ 0.4	90.1 $\pm$ 0.9
StanfordCars	<b>4.5</b> $\pm$ 0.1	0.9 $\pm$ 0.1	2.7 $\pm$ 0.1	33.8 $\pm$ 0.2	81.0 $\pm$ 0.1	<b>4.7</b> $\pm$ 0.2	4.2 $\pm$ 0.3	40.6 $\pm$ 0.1	86.4 $\pm$ 0.3	<b>7.0</b> $\pm$ 0.2	4.6 $\pm$ 0.1	64.7 $\pm$ 0.1	92.5 $\pm$ 0.2
SUN397	<b>13.0</b> $\pm$ 0.2	1.0 $\pm$ 0.1	10.4 $\pm$ 0.1	46.1 $\pm$ 0.2	53.2 $\pm$ 0.2	<b>20.3</b> $\pm$ 0.2	19.8 $\pm$ 0.1	53.5 $\pm$ 0.1	59.0 $\pm$ 0.1	<b>23.7</b> $\pm$ 0.2	21.6 $\pm$ 0.3	62.3 $\pm$ 0.1	61.0 $\pm$ 0.2
CIFAR10	65.5 $\pm$ 0.1	63.0 $\pm$ 0.1	<b>65.7</b> $\pm$ 0.6	85.9 $\pm$ 0.5	96.5 $\pm$ 0.4	<b>76.6</b> $\pm$ 0.3	74.8 $\pm$ 0.5	90.1 $\pm$ 0.1	96.6 $\pm$ 0.2	<b>81.7</b> $\pm$ 0.3	80.3 $\pm$ 0.3	94.1 $\pm$ 0.1	97.1 $\pm$ 0.1
CIFAR100	<b>24.8</b> $\pm$ 0.1	12.9 $\pm$ 0.1	18.1 $\pm$ 0.2	63.3 $\pm$ 0.8	82.5 $\pm$ 1.2	<b>38.9</b> $\pm$ 0.3	32.0 $\pm$ 0.4	70.7 $\pm$ 0.7	83.4 $\pm$ 0.9	<b>45.9</b> $\pm$ 0.2	39.7 $\pm$ 0.2	76.2 $\pm$ 0.9	84.6 $\pm$ 1.2
SVHN	<b>75.2</b> $\pm$ 0.2	73.5 $\pm$ 0.3	73.1 $\pm$ 0.2	65.0 $\pm$ 0.2	96.5 $\pm$ 0.3	<b>75.8</b> $\pm$ 0.4	75.6 $\pm$ 0.2	63.5 $\pm$ 0.2	96.9 $\pm$ 0.3	<b>81.4</b> $\pm$ 0.1	79.0 $\pm$ 0.5	51.0 $\pm$ 0.2	97.1 $\pm$ 0.3
ABIDE	<b>76.9</b> $\pm$ 2.1	74.0 $\pm$ 2.2	73.1 $\pm$ 1.6	65.4 $\pm$ 3.8	60.6 $\pm$ 4.2	63.5 $\pm$ 2.2	<b>64.4</b> $\pm$ 3.4	55.8 $\pm$ 2.6	70.2 $\pm$ 2.5	<b>67.3</b> $\pm$ 2.6	65.7 $\pm$ 3.4	54.8 $\pm$ 3.4	73.1 $\pm$ 4.2

**Table 1. (Main Results)** Performance overview of ILM-VP, prompt baseline methods (RLM-VP and FLM-VP), and finetuning methods (LP and FF) over 13 target image classification datasets using 3 pretrained source models.

# Target Dataset Analysis

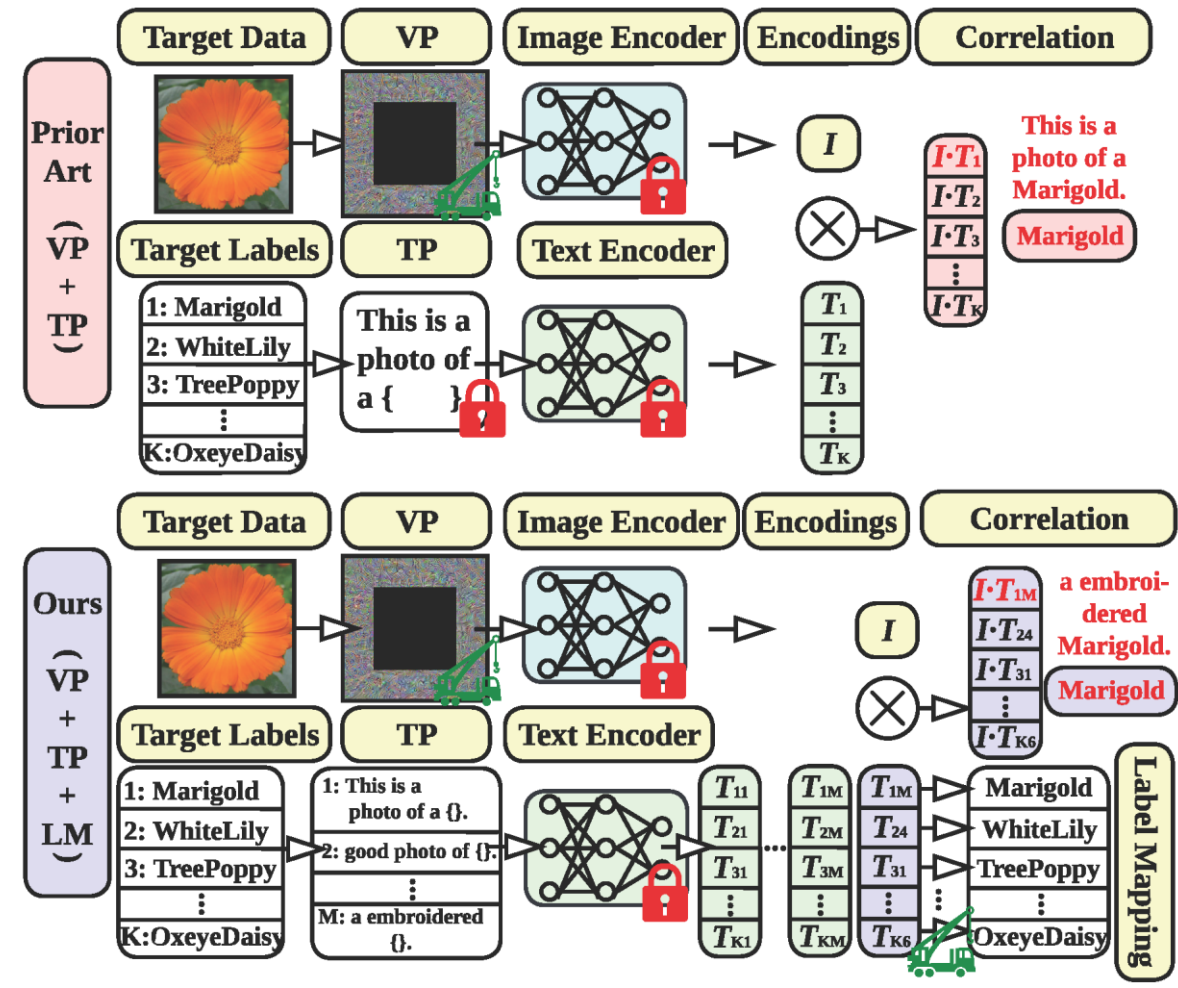


**Figure 2.** ILM-VP's improvements over FLM-VP on representative datasets (datasets with improvements over 3%) using ResNet-18.



**Figure 3.** ILM-VP and FLM-VP performance on different fractions of GTSRB dataset (43 classes and more than 900 training samples per class) using ResNet-18.

# Extension: LM in Text Domain for CLIP





# Extension: LM in Text Domain for CLIP

Methods	VP+TP Acc(%)	Acc(%)	Ours (VP+TP+LM) Examples of context prompt template → target label
Flowers102	70.0	<b>83.7</b>	a close-up photo of a {} → buttercup
DTD	56.8	<b>63.9</b>	graffiti of a {} → blotchy
UCF101	66.0	<b>70.6</b>	a {} in a video game → baseball pitch
Food101	78.9	<b>79.1</b>	a photo of the dirty {} → crab cake
SVHN	89.9	<b>91.2</b>	a photo of a {} → 7
EuroSAT	96.4	<b>96.9</b>	a pixelated photo of a {} → river
StanfordCars	57.2	<b>57.6</b>	the toy {} → 2011 audi s6 sedan
SUN397	60.5	<b>61.2</b>	a photo of a large {} → archive
CIFAR10	93.9	<b>94.4</b>	a pixelated photo of a {} → ship
ImageNet-R	67.5	<b>68.6</b>	a rendition of a {} → gold fish
ImageNet-Sketch	38.5	<b>39.7</b>	a sketch of a {} → eagle

**Table 2.** Results of CLIP-based prompt learning ‘VP+TP+LM’ and the baseline ‘VP+TP’ (restricted to using text prompt template “This is a photo of a {}”) over 11 target datasets.



# More Explanation Examples



Figure 4. Interpretation merit of ILM vs. FLM, visualized by LM results in VP to re-purposed a ResNet-18 to conduct image classification tasks.

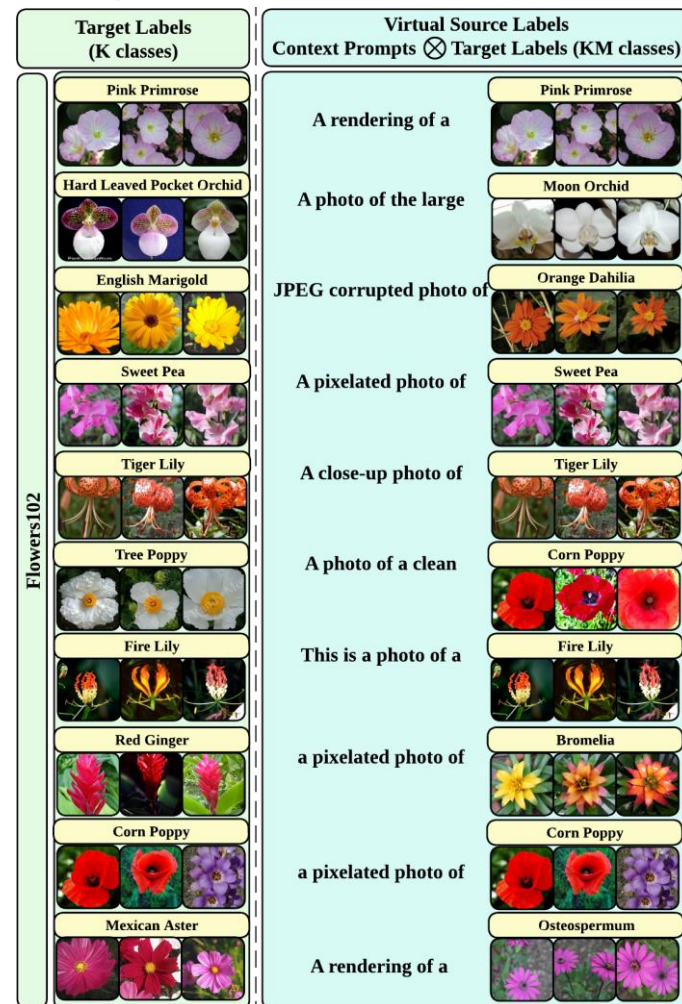


Figure 5. LM results of our proposed 'VP+TP+LM' method for CLIP, which shows significant interpretability.



Thank You

terima kasih  
multumesc  
ありがとうございます  
謝謝  
ngiyabongsa  
uksema  
baie  
dankie  
Met dank  
obrigada  
molte grazie  
Danke schön!  
감사합니다  
obrigado  
謝謝  
gracias  
Благодарность  
شكراً  
Спасиби  
Dziękuje  
dank u  
mahalo  
tusind tak