



SINGAPORE UNIVERSITY OF  
TECHNOLOGY AND DESIGN

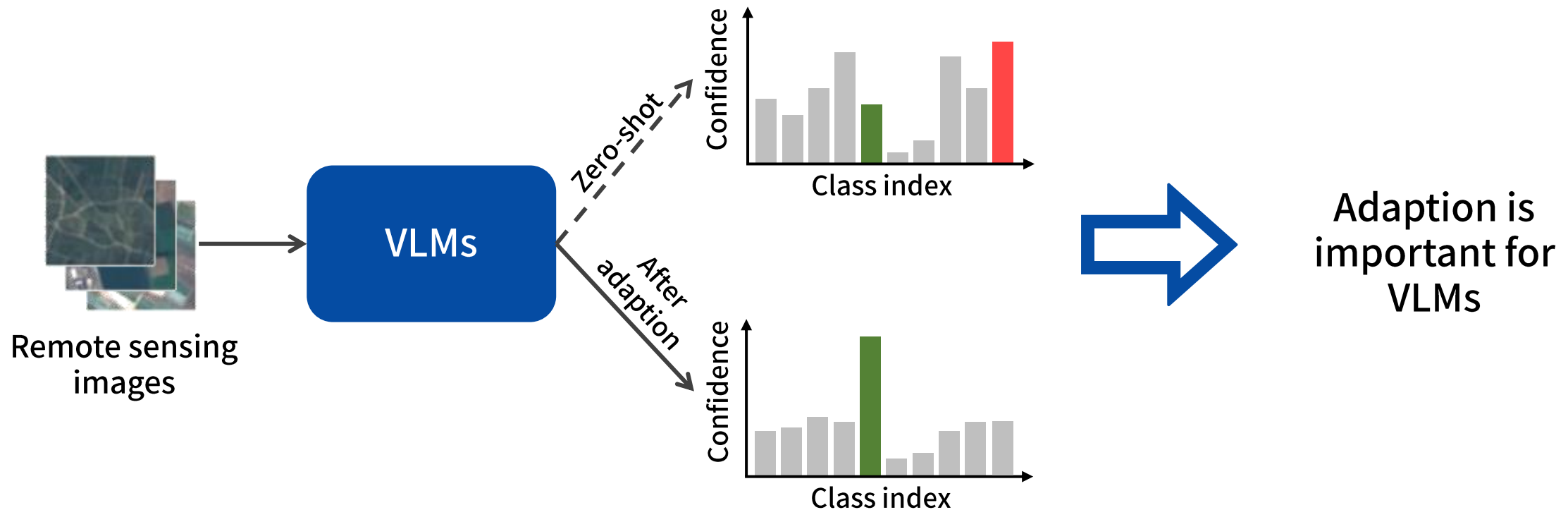


# **Candidate Pseudolabel Learning: Enhancing Vision-Language Models by Prompt Tuning with Unlabeled Data**

**Authors:** Jiahan Zhang, Qi Wei, Feng Liu, Lei Feng

**Date:** 07/24/2024

### ❖ Image classification task for VLMs with specialized domain



### ❖ Image classification task for VLMs with specialized domain

Adaption is important for VLMs



But it still needs a lot of well-labeled data for adaption



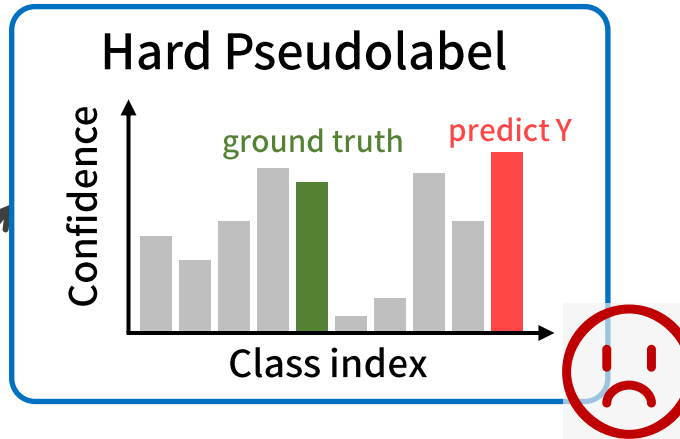
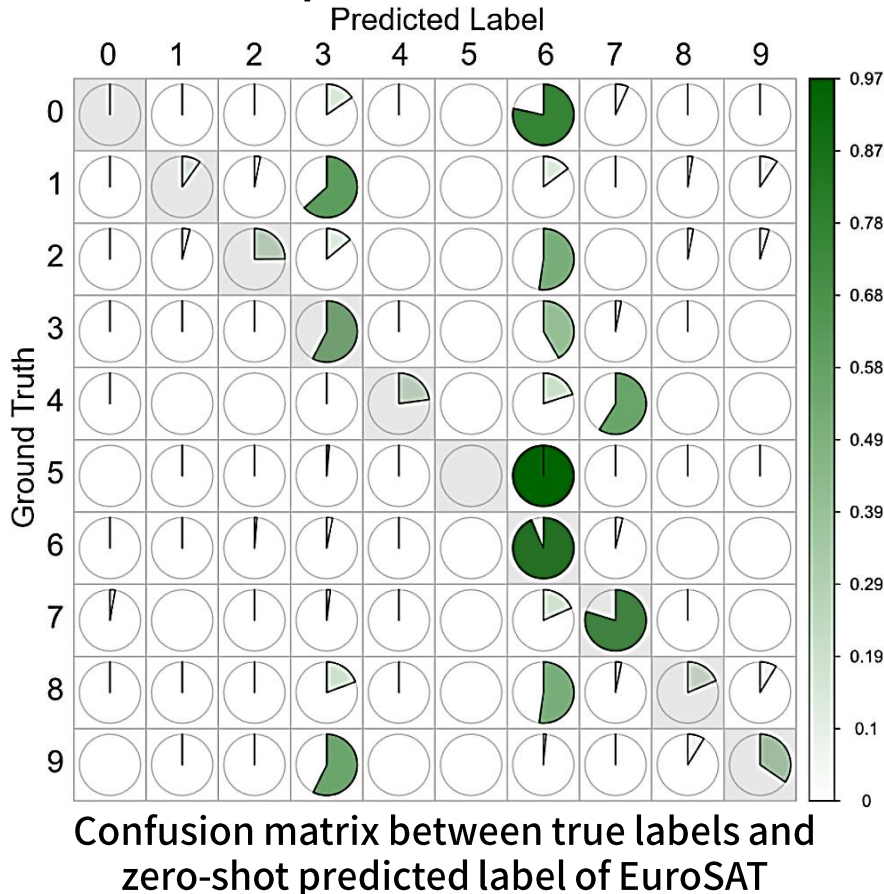
When a large amount of unlabeled data is available

Can we use the zero-shot capabilities of VLMs to exploit unlabeled data for adaption?



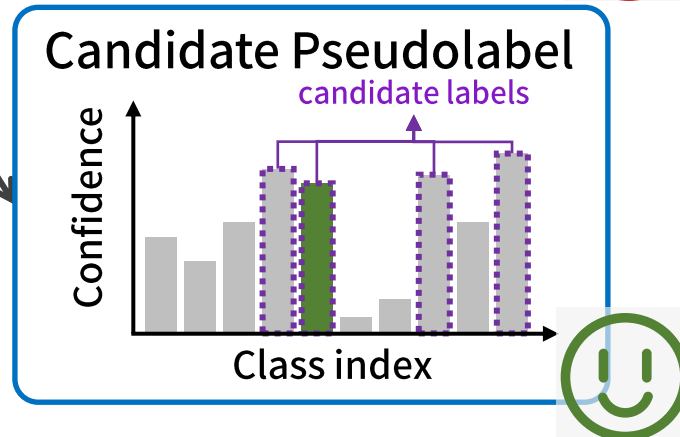
## Comparison and Motivation

- ❖ **How** can we better use the zero-shot capabilities of VLMs to exploit unlabeled data?



We observed that lots of zero-shot predicted labels are **incorrect and imbalanced**:

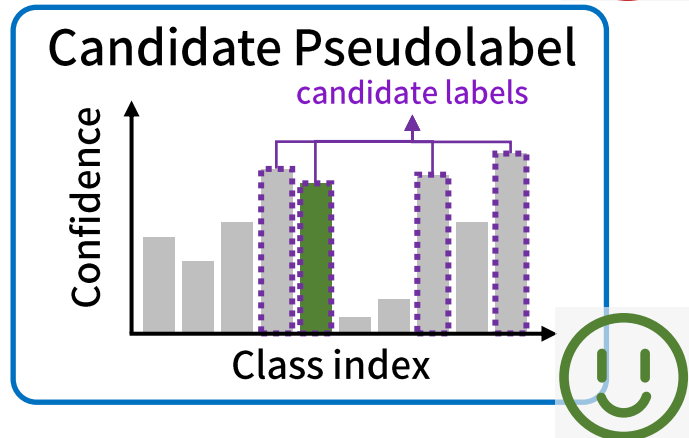
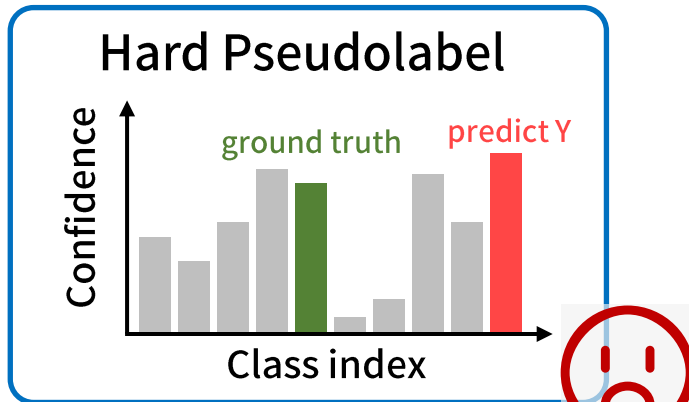
→ Directly applying them as hard pseudolabels will inherit these shortcomings



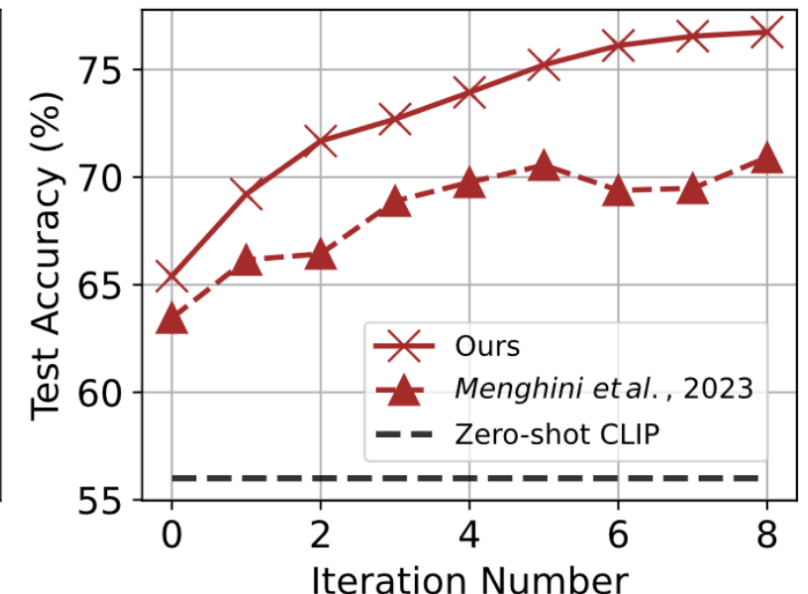
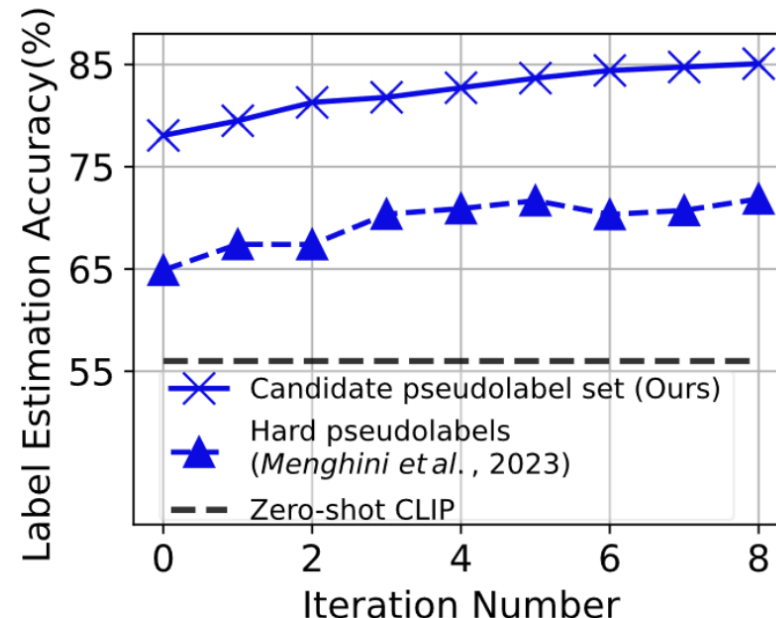
→ Generating a set of potential pseudolabels from the predictions will alleviate these issues

## Comparison and Motivation

### Advantages of candidate pseudolabels over hard pseudolabels

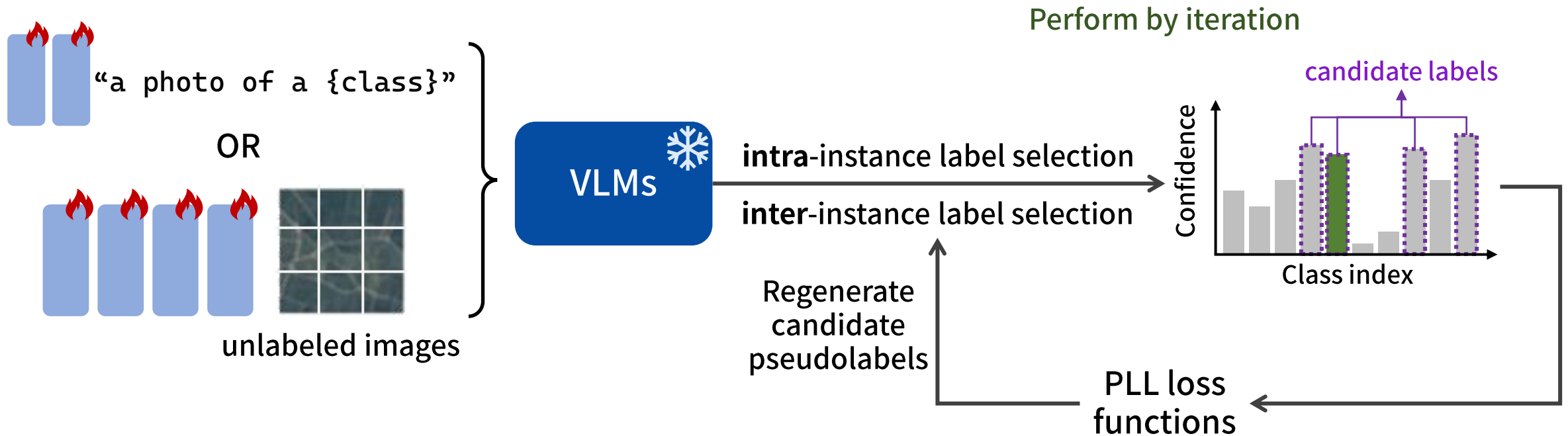


- Higher label estimation accuracy (the rate at which the true label is included in the pseudolabels)
- Fully leverage the zero-shot capabilities/information from VLMs
- Result in improved performance on test accuracy



## Overview of Workflow

- ❖ Overall workflow for prompt tuning with unlabeled data in **CPL** (Candidate Pseudolabel Learning)

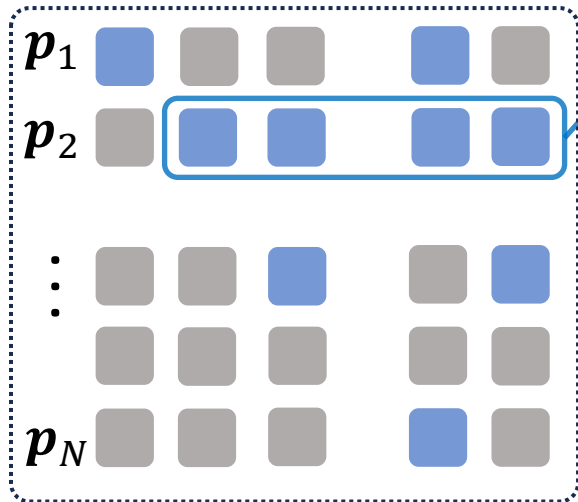
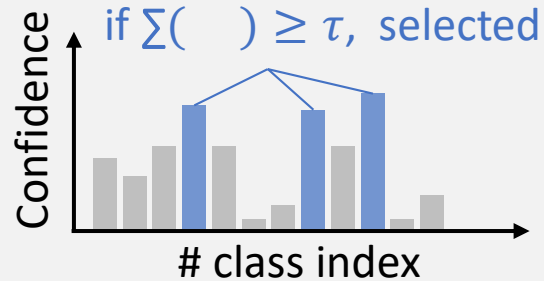


## Scheme for Generating Candidate Pseudolabels

### 1. The strategy of intra-instance label selection

#### 1) Intra-inst. Label Selection

Conf. distribution per instance



Confidence Score Matrix

$$(\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_N)^\top$$

- Similar to top-K confident labels selection
- And further consider varying levels of identification difficulty to determine different K values for each instance

- The vector of conf. scores for instance  $i$ :

$$\mathbf{p}_i = \mathbf{g}(\mathbf{f}_\theta(\mathbf{x})) = (p_{i1}, p_{i2}, \dots, p_{iC})^\top$$

- The intra-instance candidate set:

$$S_i^{\text{intra}} = \text{MinSize}(\{c \mid \sum_{c=1}^C p_{ic} \geq \tau\})$$

to get the set with the minimal size that the cumulative conf. score just surpasses a threshold

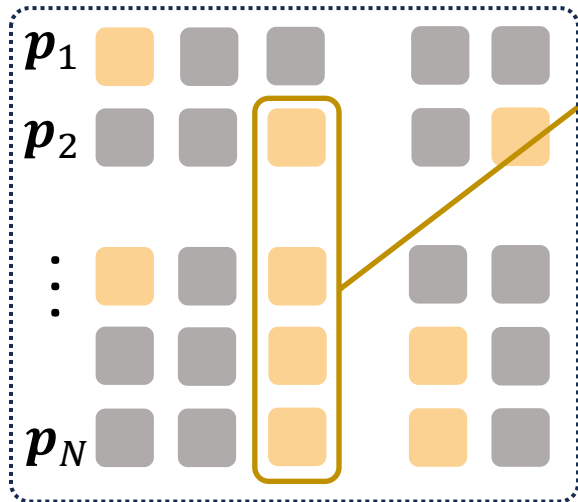
- Determine the threshold:

$$\tau = \text{Quantile}(\text{Sort}(\hat{\mathbf{p}}), \alpha)$$

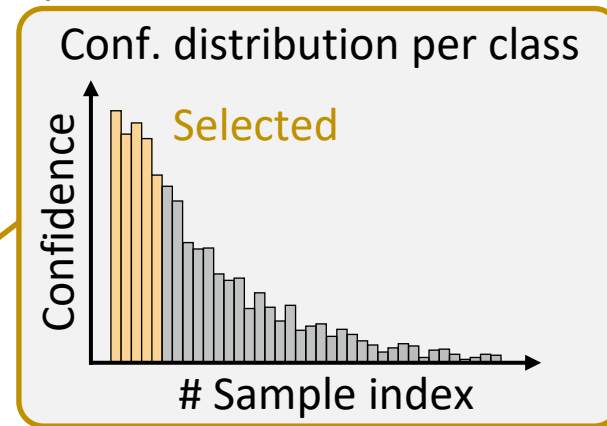
a function that returns the value at the given quantile  $\alpha$

### 2. The strategy of inter-instance label selection

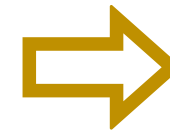
#### 2) Inter-inst. Label Selection



Confidence Score Matrix  
 $(p_1, p_2, \dots, p_N)^T$



- To balance the ratio of each category in the candidate pseudolabel set
- To further refine the candidate pseudolabels



- Conf. scores across all unlabeled instances for class  $c$ :

$$q_c = (p_{1c}, p_{2c}, \dots, p_{Nc})$$

- The inter-instance candidate set:

$$S_i^{\text{inter}} = \{c \mid p_{ic} > \text{Quantile}(\text{Sort}(q_c), \beta)\}_{c=1}^C$$

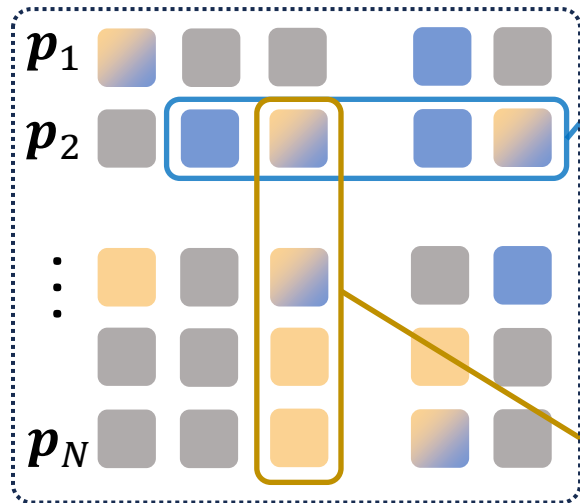
to select the candidate labels of relatively higher confidence levels within the class



## Scheme for Generating Candidate Pseudolabels

### Final candidate sets and training set for unlabeled instances

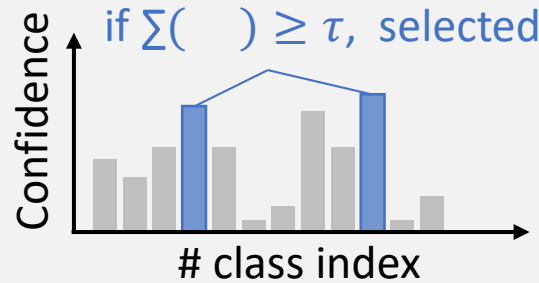
- Intra-inst. selected candidate labels
- Inter-inst. selected candidate labels
- The candidate labels selected by both



Confidence Score Matrix  
 $(\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_N)^\top$

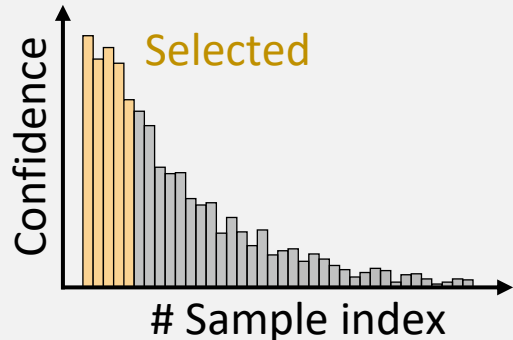
#### 1) Intra-inst. Label Selection

Conf. distribution per instance



#### 2) Inter-inst. Label Selection

Conf. distribution per class



$$S_i = S_i^{\text{intra}} \cap S_i^{\text{inter}}$$

- Ensures a more balanced and accurate candidate set
- Enhancing the model's ability to learn from a diverse and equitable distribution of pseudolabels

Finally, the construction of training set for unlabeled data:

$$D_T = \{(\mathbf{x}_i, S_i) \mid |S_i| > 0\}_{i=1}^N$$

If an instance has an empty candidate pseudolabel set, then filter it out

## Learning with Candidate Pseudolabels

- ❖ The construction of training target for unlabeled data

The training set for unlabeled data:

$$D_T = \{(\mathbf{x}_i, S_i) \mid |S_i| > 0\}_{i=1}^N$$

Transform the candidate set into training target:

$$\begin{cases} s_{ic} = 1 & \text{if } c \in S_i \\ s_{ic} = 0 & \text{if } c \notin S_i \text{ for } c \in [C] \end{cases} \longrightarrow (\mathbf{x}_i, \mathbf{s}_i)$$

- ❖ The training objectives:

If a labeled set  $D_L$  and an unlabeled set  $D_{UL}$  are provided

$$D_L = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^O$$

$$\mathcal{L} = \mathcal{L}_L + \lambda \mathcal{L}_{UL}$$

$$= \frac{1}{b_1} \sum_{i=1}^{b_1} L_{ce}(\mathbf{x}_i, \mathbf{y}_i) + \lambda \frac{1}{b_2} \sum_{i=1}^{b_2} L_{pl}(\mathbf{x}_i, \mathbf{s}_i)$$

partial-label  
learning loss func.

If only an unlabeled data  $D_{UL}$  is provided

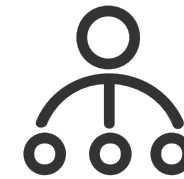
$$\mathcal{L} = \mathcal{L}_{UL} = \frac{1}{b_2} \sum_{i=1}^{b_2} L_{pl}(\mathbf{x}_i, \mathbf{s}_i)$$

## Experimental Settings

Design extensive experiments across 3 dimensions

❖ Learning paradigm variety:

semi-supervised learning (SSL), unsupervised learning (UL), and transductive zero-shot learning (TRZSL). All involving the access of unlabeled data



❖ Prompt tuning variety:

visual prompt tuning (VPT) or text prompt tuning (TPT)



❖ Task variety:

FGVC-Aircraft, EuroSAT, CUB, Flowers102, RESISC45, DTD, CALTECH-101, UCF-101, and CIFAR-100 (9 classification tasks)



## Experimental Results

### ❖ The pros of CPL compared with previous hard pseudolabel on prompt tuning:

Table 1: Comparison results of top-1 test accuracy (%) on six benchmarks when applying **Textual prompts** as tuning strategy. Note that “✓” and “✗” denote whether full unlabeled data are utilized for fine-tuning or not, respectively.

Methods	Flowers102			RESISC45			DTD		
	SSL	UL	TRZSL	SSL	UL	TRZSL	SSL	UL	TRZSL
Zero-shot CLIP	63.67 <sub>0.00</sub>			54.48 <sub>0.00</sub>			43.24 <sub>0.00</sub>		
FPL ✗	75.96 <sub>0.74</sub>	65.67 <sub>0.23</sub>	80.97 <sub>0.00</sub>	68.13 <sub>0.55</sub>	63.07 <sub>0.38</sub>	72.11 <sub>0.00</sub>	37.10 <sub>5.45</sub>	44.96 <sub>0.55</sub>	46.30 <sub>0.03</sub>
CPL (Ours) ✗	<b>77.36</b> <sub>0.24</sub>	<b>70.01</b> <sub>0.21</sub>	<b>84.60</b> <sub>0.10</sub>	<b>71.73</b> <sub>0.57</sub>	<b>68.47</b> <sub>0.34</sub>	<b>72.16</b> <sub>0.26</sub>	<b>54.63</b> <sub>0.79</sub>	<b>48.92</b> <sub>0.17</sub>	<b>59.79</b> <sub>1.32</sub>
GRIP ✓	83.60 <sub>0.48</sub>	69.84 <sub>1.06</sub>	86.26 <sub>0.00</sub>	74.11 <sub>0.68</sub>	70.55 <sub>0.88</sub>	81.07 <sub>0.00</sub>	56.07 <sub>0.85</sub>	46.09 <sub>1.06</sub>	65.30 <sub>0.01</sub>
CPL (Ours) ✓	<b>89.66</b> <sub>0.36</sub>	<b>72.90</b> <sub>0.78</sub>	<b>87.35</b> <sub>0.76</sub>	<b>80.98</b> <sub>0.11</sub>	<b>77.39</b> <sub>0.44</sub>	<b>85.85</b> <sub>0.49</sub>	<b>61.21</b> <sub>0.56</sub>	<b>51.91</b> <sub>0.71</sub>	<b>68.00</b> <sub>0.34</sub>

Methods	CUB			EuroSAT			FGVCAircraft		
	SSL	UL	TRZSL	SSL	UL	TRZSL	SSL	UL	TRZSL
Zero-shot CLIP	51.82 <sub>0.00</sub>			32.88 <sub>0.00</sub>			17.58 <sub>0.00</sub>		
FPL ✗	55.29 <sub>0.59</sub>	53.04 <sub>0.53</sub>	55.44 <sub>0.20</sub>	62.05 <sub>1.64</sub>	48.96 <sub>1.49</sub>	53.70 <sub>26.87</sub>	20.02 <sub>0.77</sub>	16.62 <sub>0.67</sub>	17.55 <sub>0.37</sub>
CPL (Ours) ✗	<b>56.37</b> <sub>0.45</sub>	<b>54.18</b> <sub>0.05</sub>	<b>64.01</b> <sub>0.17</sub>	<b>64.84</b> <sub>2.15</sub>	<b>51.45</b> <sub>1.97</sub>	<b>54.03</b> <sub>2.27</sub>	<b>22.37</b> <sub>0.66</sub>	<b>18.90</b> <sub>0.20</sub>	<b>28.47</b> <sub>0.43</sub>
GRIP ✓	56.65 <sub>0.33</sub>	51.42 <sub>0.21</sub>	59.48 <sub>0.38</sub>	58.66 <sub>2.64</sub>	57.21 <sub>1.77</sub>	92.33 <sub>0.69</sub>	16.98 <sub>0.82</sub>	15.22 <sub>0.71</sub>	26.08 <sub>0.25</sub>
CPL (Ours) ✓	<b>58.53</b> <sub>0.24</sub>	<b>53.47</b> <sub>0.36</sub>	<b>66.20</b> <sub>0.50</sub>	<b>77.51</b> <sub>0.80</sub>	<b>67.26</b> <sub>0.47</sub>	<b>93.78</b> <sub>0.12</sub>	<b>22.48</b> <sub>0.63</sub>	<b>18.35</b> <sub>0.27</sub>	<b>30.86</b> <sub>0.70</sub>

Table 2: Comparison results of top-1 test accuracy (%) on six benchmarks when applying **Visual prompts** as tuning strategy. Note that “✓” and “✗” denote whether full unlabeled data are utilized for fine-tuning or not, respectively.

Methods	Flowers102			RESISC45			DTD		
	SSL	UL	TRZSL	SSL	UL	TRZSL	SSL	UL	TRZSL
Zero-shot CLIP	63.67 <sub>0.00</sub>			54.48 <sub>0.00</sub>			43.24 <sub>0.00</sub>		
FPL ✗	67.03 <sub>0.65</sub>	65.50 <sub>0.41</sub>	71.94 <sub>0.00</sub>	65.14 <sub>0.25</sub>	62.24 <sub>0.22</sub>	67.85 <sub>0.00</sub>	47.60 <sub>1.09</sub>	47.69 <sub>0.48</sub>	52.43 <sub>0.00</sub>
CPL (Ours) ✗	<b>70.58</b> <sub>0.13</sub>	<b>68.94</b> <sub>0.16</sub>	<b>78.13</b> <sub>0.31</sub>	<b>68.85</b> <sub>0.13</sub>	<b>67.97</b> <sub>0.17</sub>	<b>72.18</b> <sub>0.27</sub>	<b>52.64</b> <sub>0.68</sub>	<b>50.37</b> <sub>0.46</sub>	<b>55.90</b> <sub>0.69</sub>
GRIP ✓	67.95 <sub>1.2</sub>	63.09 <sub>0.56</sub>	77.18 <sub>0.00</sub>	71.22 <sub>0.77</sub>	68.43 <sub>0.61</sub>	82.19 <sub>0.00</sub>	54.57 <sub>4.86</sub>	50.51 <sub>0.99</sub>	62.78 <sub>0.00</sub>
CPL (Ours) ✓	<b>73.52</b> <sub>0.37</sub>	<b>67.25</b> <sub>0.41</sub>	<b>80.14</b> <sub>0.73</sub>	<b>78.46</b> <sub>0.74</sub>	<b>72.97</b> <sub>0.58</sub>	<b>86.67</b> <sub>0.33</sub>	<b>58.74</b> <sub>0.81</sub>	<b>53.42</b> <sub>0.56</sub>	<b>65.31</b> <sub>0.78</sub>

Methods	CUB			EuroSAT			FGVCAircraft		
	SSL	UL	TRZSL	SSL	UL	TRZSL	SSL	UL	TRZSL
Zero-shot CLIP	51.82 <sub>0.00</sub>			32.88 <sub>0.00</sub>			17.58 <sub>0.00</sub>		
FPL ✗	52.86 <sub>0.42</sub>	53.17 <sub>0.06</sub>	54.17 <sub>0.16</sub>	52.47 <sub>2.53</sub>	48.79 <sub>3.69</sub>	68.68 <sub>14.74</sub>	20.14 <sub>0.26</sub>	18.28 <sub>0.33</sub>	16.28 <sub>0.45</sub>
CPL (Ours) ✗	<b>53.37</b> <sub>0.55</sub>	<b>53.28</b> <sub>0.31</sub>	<b>56.43</b> <sub>0.21</sub>	<b>66.37</b> <sub>2.10</sub>	<b>52.83</b> <sub>2.10</sub>	<b>74.02</b> <sub>1.34</sub>	<b>21.52</b> <sub>0.68</sub>	<b>20.10</b> <sub>0.51</sub>	<b>26.73</b> <sub>0.08</sub>
GRIP ✓	<b>53.83</b> <sub>0.11</sub>	<b>52.91</b> <sub>0.26</sub>	54.92 <sub>0.17</sub>	63.48 <sub>3.09</sub>	63.68 <sub>3.42</sub>	96.97 <sub>0.77</sub>	19.43 <sub>0.50</sub>	17.51 <sub>0.61</sub>	26.42 <sub>0.30</sub>
CPL (Ours) ✓	49.50 <sub>0.42</sub>	52.11 <sub>0.24</sub>	<b>56.37</b> <sub>0.06</sub>	<b>72.03</b> <sub>1.24</sub>	<b>68.93</b> <sub>1.15</sub>	<b>98.31</b> <sub>0.18</sub>	<b>20.51</b> <sub>0.68</sub>	<b>18.26</b> <sub>0.38</sub>	<b>30.26</b> <sub>0.46</sub>

### ❖ The pros of CPL when combined with existing label-free CLIP fine-tuning pipeline:

	Flowers-102	UCF-101	CIFAR-100	EuroSAT	DTD	CALTECH-101
CLIP	66.6	61.0	64.2	45.1	42.9	90.5
CLIP-PR	57.7	57.9	63.2	44.2	40.1	84.8
UPL	<u>71.5</u>	63.9	65.8	62.2	<u>48.0</u>	90.6
LaFTer	71.0	<u>68.2</u>	<u>74.6</u>	<u>73.9</u>	46.1	<u>93.3</u>
LaFTer + Ours	<b>76.7</b>	<b>71.0</b>	<b>77.3</b>	<b>82.2</b>	<b>56.3</b>	<b>93.4</b>

## More Experimental Results

- ❖ CPL performance when employing various PLL loss functions:
- ❖ CPL performance on imbalance dataset:

Methods		SSL	UL	TRZSL	
Zero-shot CLIP		43.24 <sub>0.00</sub>		43.45 <sub>0.00</sub>	
Use a fixed amount of unlabeled data	FPL	✗	37.10 <sub>5.45</sub>	44.96 <sub>0.55</sub>	46.30 <sub>0.03</sub>
	CPL <sub>Soft CE</sub>	✗	51.83 <sub>0.62</sub>	47.02 <sub>0.37</sub>	59.69 <sub>0.59</sub>
	CPL <sub>CC</sub>	✗	54.63 <sub>0.79</sub>	48.92 <sub>0.17</sub>	59.79 <sub>1.32</sub>
	CPL <sub>RC</sub>	✗	54.98 <sub>0.49</sub>	49.96 <sub>0.15</sub>	59.42 <sub>0.44</sub>
	CPL <sub>CAV</sub>	✗	55.50 <sub>0.29</sub>	48.69 <sub>0.66</sub>	59.44 <sub>0.13</sub>
	CPL <sub>LW</sub>	✗	55.21 <sub>0.74</sub>	49.82 <sub>0.91</sub>	59.24 <sub>0.72</sub>
Use complete unlabeled data	GRIP	✓	56.07 <sub>0.85</sub>	46.09 <sub>1.06</sub>	65.30 <sub>0.01</sub>
	CPL <sub>Soft CE</sub>	✓	60.83 <sub>0.66</sub>	49.13 <sub>0.10</sub>	66.26 <sub>0.77</sub>
	CPL <sub>CC</sub>	✓	61.21 <sub>0.56</sub>	51.91 <sub>0.71</sub>	68.00 <sub>0.34</sub>
	CPL <sub>RC</sub>	✓	60.21 <sub>0.46</sub>	51.58 <sub>0.11</sub>	67.95 <sub>0.31</sub>
	CPL <sub>CAV</sub>	✓	61.06 <sub>0.50</sub>	49.31 <sub>0.19</sub>	67.76 <sub>0.53</sub>
	CPL <sub>LW</sub>	✓	60.20 <sub>0.69</sub>	52.23 <sub>0.84</sub>	68.29 <sub>0.99</sub>

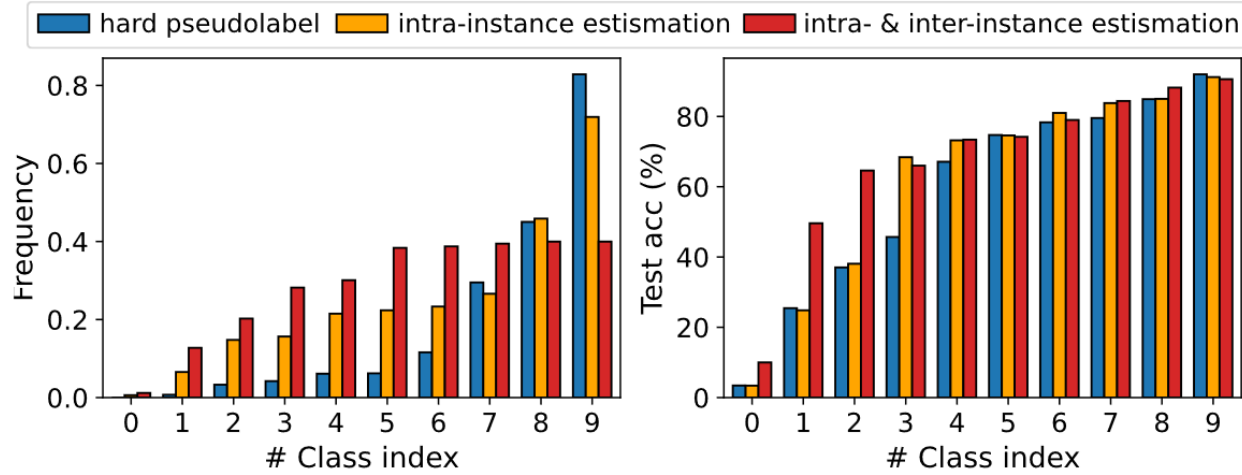
Methods	Balanced	Imbalanced $\delta=100$	Imbalanced $\delta=50$
LaFTer	74.64	65.63	66.59
CPL (w/o inter)	76.07	66.68	67.85
CPL	<b>77.32</b>	<b>67.70</b>	<b>69.65</b>

- ❖ CPL performance with different VLM encoder:

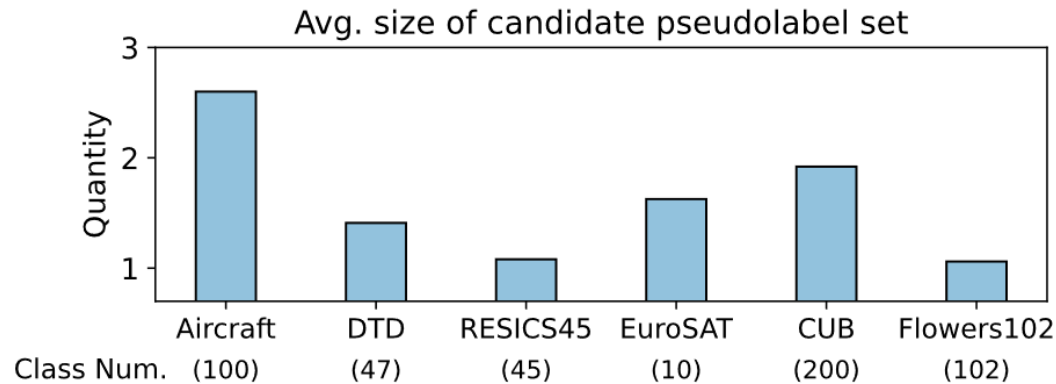
Methods		SSL	UL	TRZSL	
DTD	Zero-shot CLIP		52.45 <sub>0.00</sub>	51.61 <sub>0.00</sub>	
	FPL	✗	60.61 <sub>1.56</sub>	52.99 <sub>0.43</sub>	60.77 <sub>0.54</sub>
	CPL	✗	<b>62.78</b> <sub>0.17</sub>	<b>57.23</b> <sub>0.19</sub>	<b>62.52</b> <sub>1.38</sub>
	GRIP	✓	60.91 <sub>0.00</sub>	54.40 <sub>0.00</sub>	64.92 <sub>0.00</sub>
	CPL	✓	<b>69.82</b> <sub>0.32</sub>	<b>57.20</b> <sub>0.45</sub>	<b>71.97</b> <sub>0.46</sub>
RESISC45	Zero-shot CLIP		62.67 <sub>0.00</sub>	62.13 <sub>0.00</sub>	
	FPL	✗	79.01 <sub>0.55</sub>	70.85 <sub>0.66</sub>	77.69 <sub>0.83</sub>
	CPL	✗	<b>80.38</b> <sub>0.37</sub>	<b>76.01</b> <sub>0.19</sub>	<b>79.97</b> <sub>0.77</sub>
	GRIP	✓	81.53 <sub>0.00</sub>	76.86 <sub>0.00</sub>	86.88 <sub>0.00</sub>
	CPL	✓	<b>87.75</b> <sub>0.29</sub>	<b>80.88</b> <sub>0.86</sub>	<b>89.73</b> <sub>1.73</sub>
Flowers102	Zero-shot CLIP		73.98 <sub>0.00</sub>	73.05 <sub>0.00</sub>	
	FPL	✗	<b>89.07</b> <sub>0.94</sub>	77.81 <sub>0.30</sub>	91.84 <sub>0.73</sub>
	CPL	✗	88.37 <sub>0.39</sub>	<b>82.98</b> <sub>0.14</sub>	<b>96.65</b> <sub>0.08</sub>
	GRIP	✓	94.21 <sub>0.00</sub>	82.33 <sub>0.00</sub>	96.18 <sub>0.00</sub>
	CPL	✓	<b>96.80</b> <sub>0.63</sub>	<b>83.94</b> <sub>0.69</sub>	<b>97.34</b> <sub>0.74</sub>

## Ablation and Visualization

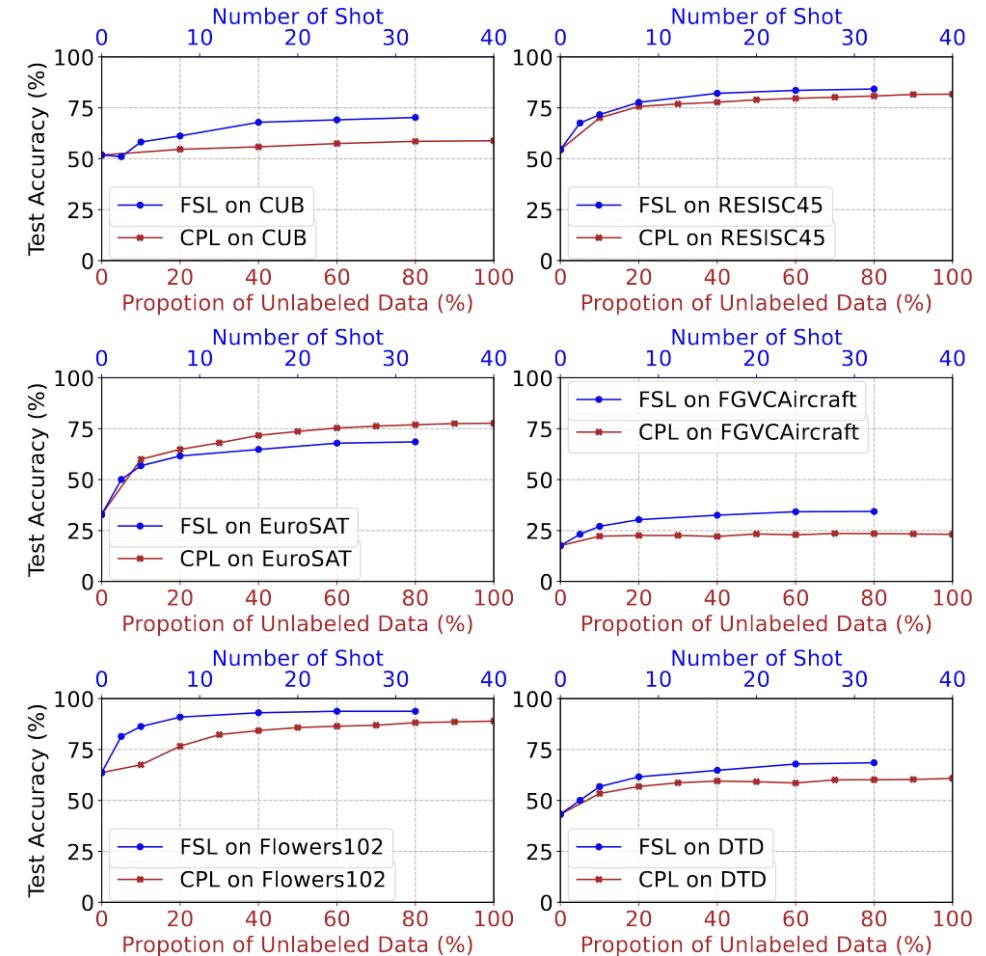
- ❖ The class-level distribution of the generated pseudolabels:



- ❖ The average size of the candidate set before the last iteration:



- ❖ Compare fully supervised few-shot learning and CPL







**Thanks For Your Listening!**

Project Page:

