



# Roll With the Punches: Expansion and Shrinkage of Soft Label Selection for Semi-supervised Fine-Grained Learning

**Yue Duan<sup>1</sup>, Zhen Zhao<sup>2</sup>, Lei Qi<sup>3</sup>, Luping Zhou<sup>2</sup>, Lei Wang<sup>4</sup>, and Yinghuan Shi<sup>1</sup>**

<sup>1</sup> Nanjing University, China <sup>2</sup> University of Sydney, Australia <sup>3</sup> Southeast University, China <sup>4</sup> University of Wollongong, Australia

Arxiv: <https://arxiv.org/abs/2312.12237>

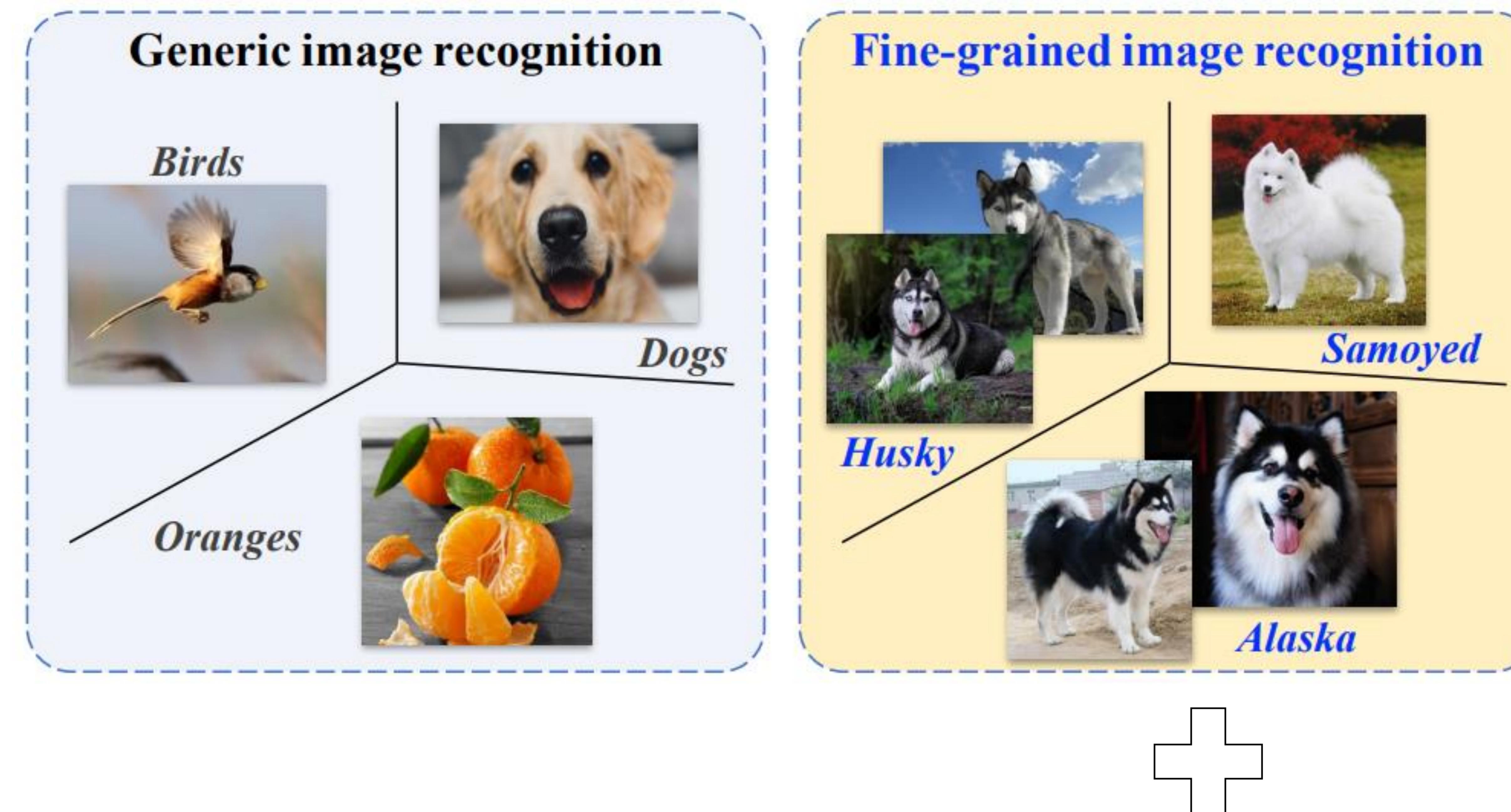
Github: <https://github.com/NJUYued/SoC4SS-FGVC>

Homepage: <https://njuyued.github.io/>

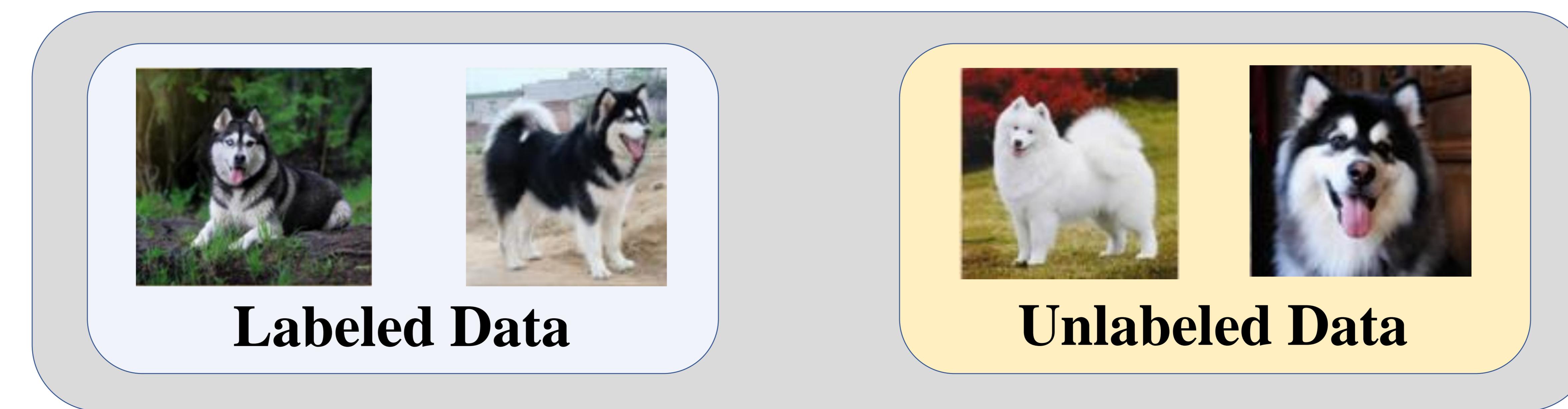


## Introduction

Generic image recognition V.S. Fine-grained image recognition (FGVC)



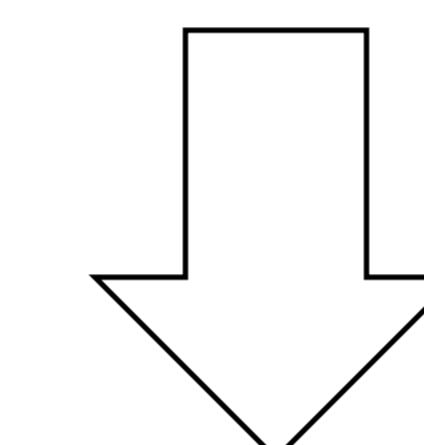
## Semi-supervised Learning



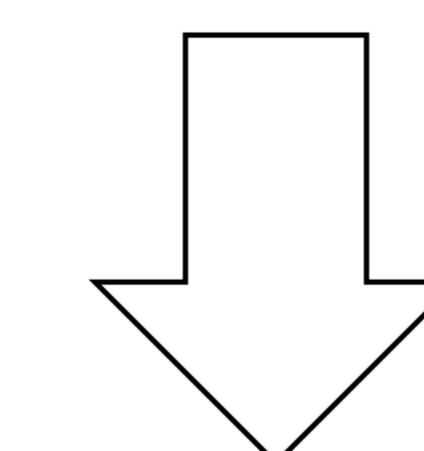


## Introduction

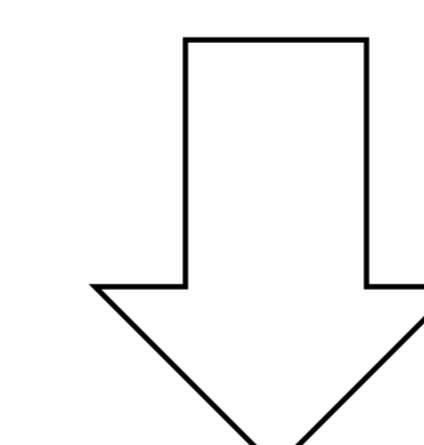
Semi-supervised learning (SSL) aims to leverage a pool of unlabeled data to alleviate the dependence of deep models on labeled data.



However, current SSL approaches achieve promising performance with clean and ordinary data, *but become unstuck against the indiscernible unlabeled data.*

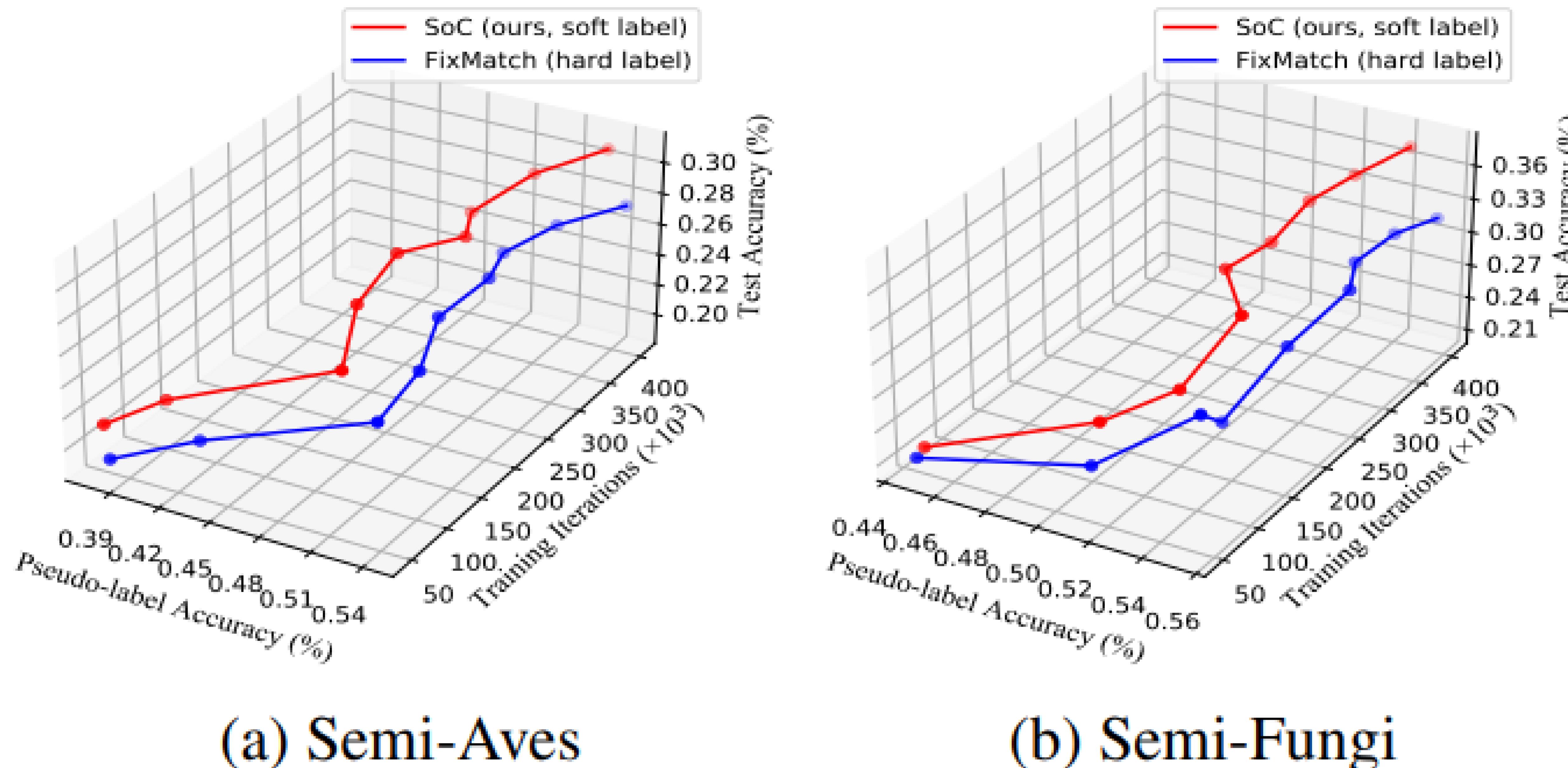


A typical and worth discussing example is the *semi-supervised fine-grained visual classification (SS-FGVC)* [1].



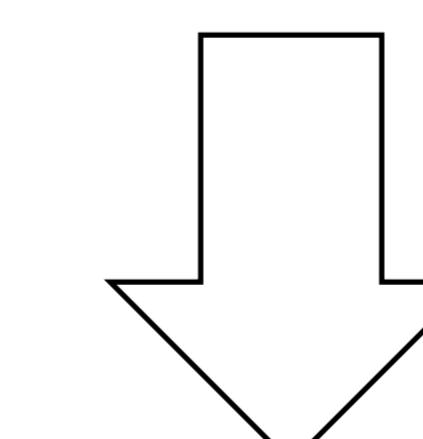
**Fine-grained data may severely affect the quality of pseudo-labels and consequently pull down the model performance.**

# Roll With the Punches: Expansion and Shrinkage of Soft Label Selection for Semi-supervised Fine-Grained Learning



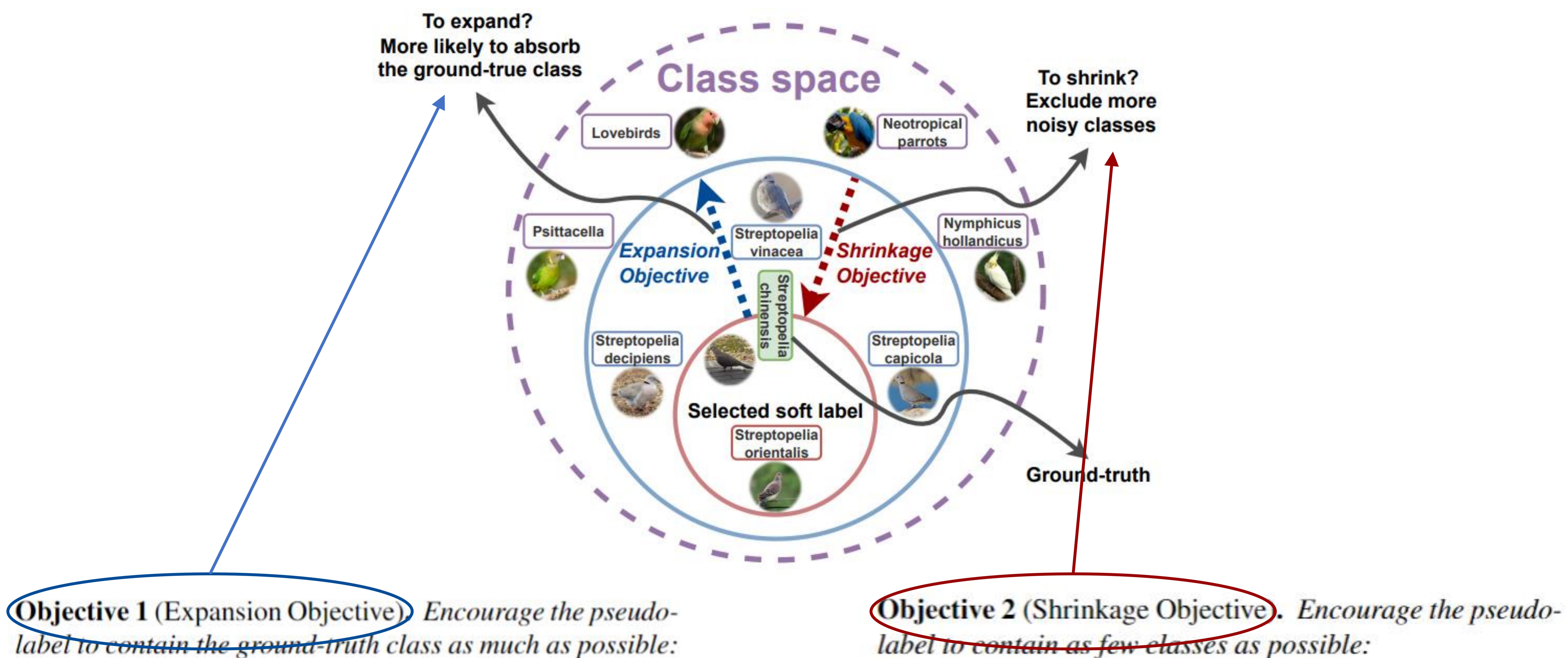
## Motivation

The bad effects of incorrect hard label on the model have overridden its low entropy advantage while *soft label can benefit the model because it could still provide useful information although it is wrong.*



**Reasonably utilizing soft labels will enable the model to be more composed when dealing with fine-grained data that is difficult to distinguish.**

## A Coupled Optimization Goal for SS-FGVC



$$\max_{\theta} \mathbb{E}_{\mathbf{x}_i^{\text{ulb}} \in \mathcal{D}^{\text{ulb}}} [\mathbb{1}(y_i^* \in \{c \mid g_{i,(c)} = 1\}) p_{i,(y_i^*)}], \quad (5)$$

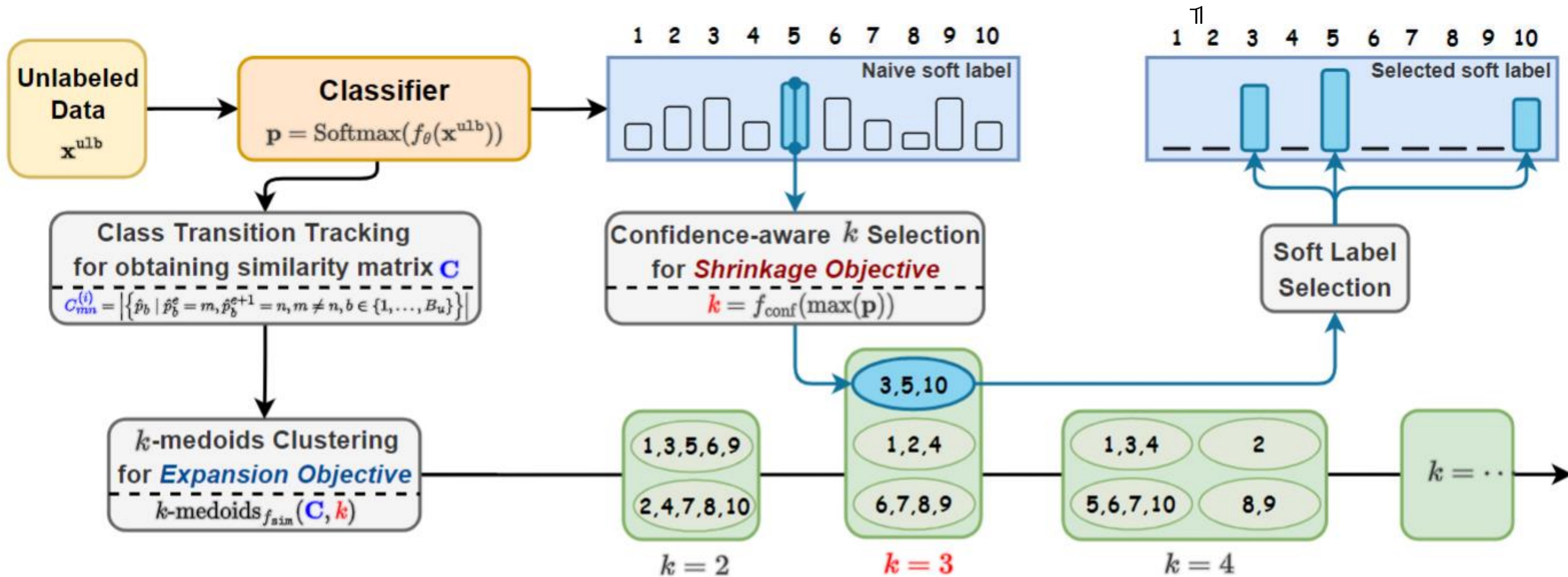
where  $y_i^* \in \mathcal{Y}$  is the ground-truth label of  $\mathbf{x}_i^{\text{ulb}}$ . For simplicity, we denote  $p_{i,(y_i^*)} \mathbb{1}(y_i^* \in \{c \mid g_{i,(c)} = 1\})$  as  $z_i^{\text{obj1}}$ .

$$\min_{\theta} \mathbb{E}_{\mathbf{x}_i^{\text{ulb}} \in \mathcal{D}^{\text{ulb}}} \sum_{c=1}^K g_{i,(c)}, \quad (6)$$

where  $\sum_{c=1}^K g_{i,(c)}$  is denoted as  $z_i^{\text{obj2}}$  for simplicity.



## Overview of Method



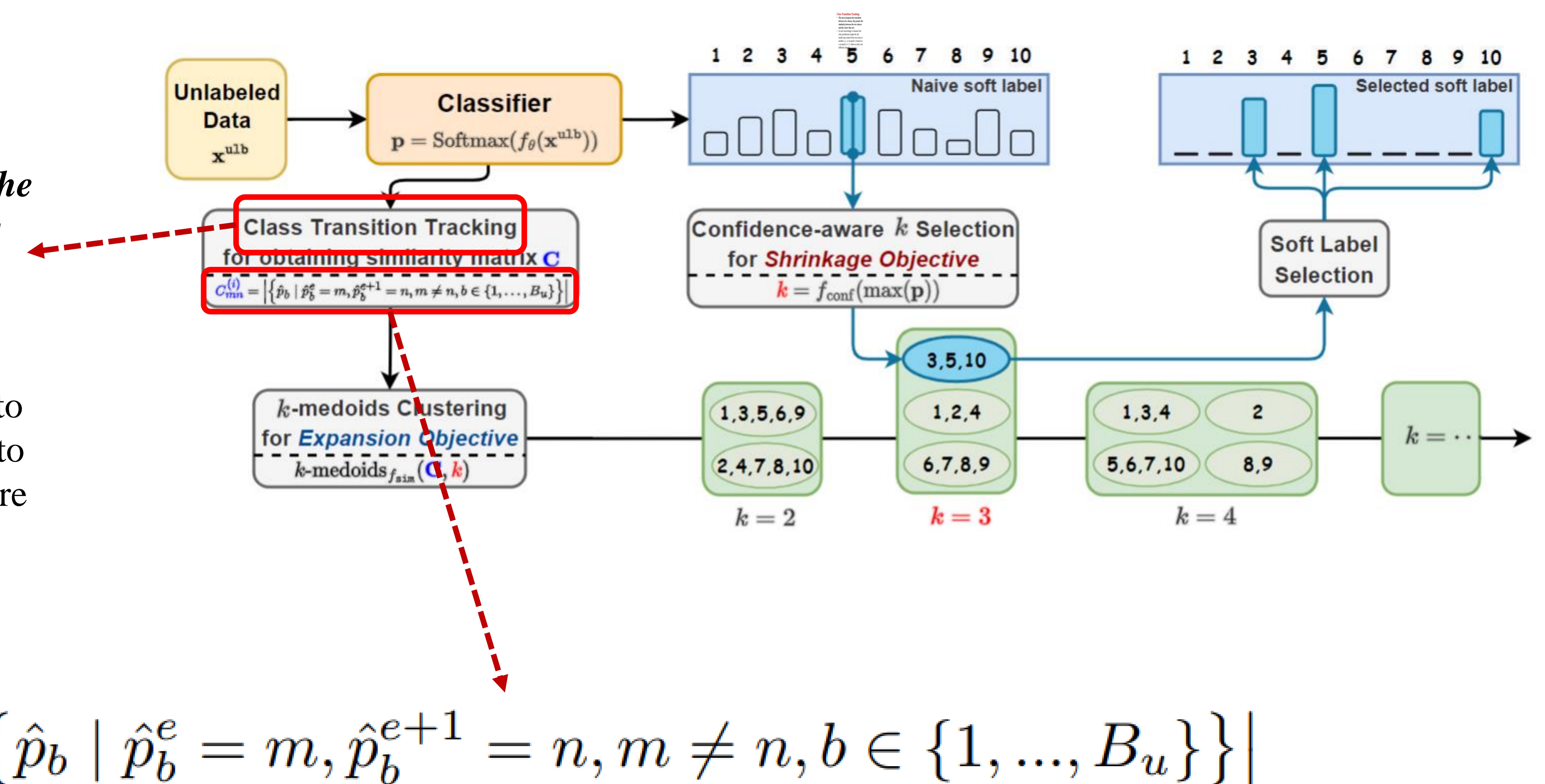
- For **Expansion Objective**, use class-level similarity obtained through class transfer tracking to acquire **candidate class clusters**.
- For **Shrinkage Objective**, select a more appropriate **clustering granularity** based on prediction confidence.

## Method

- For **Expansion Objective**, use class-level similarity obtained through class transfer tracking to acquire candidate class clusters.

### Class Transition Tracking:

- The more frequent the transition between two classes, the greater the similarity between the two classes and the closer they are.*
- As new knowledge is learned, the class predictions output by the model may transit from one class to another, i.e.,  $m$  at epoch  $e$  transits to  $n$  at epoch  $e + 1$ , where  $m$  and  $n$  are different classes.



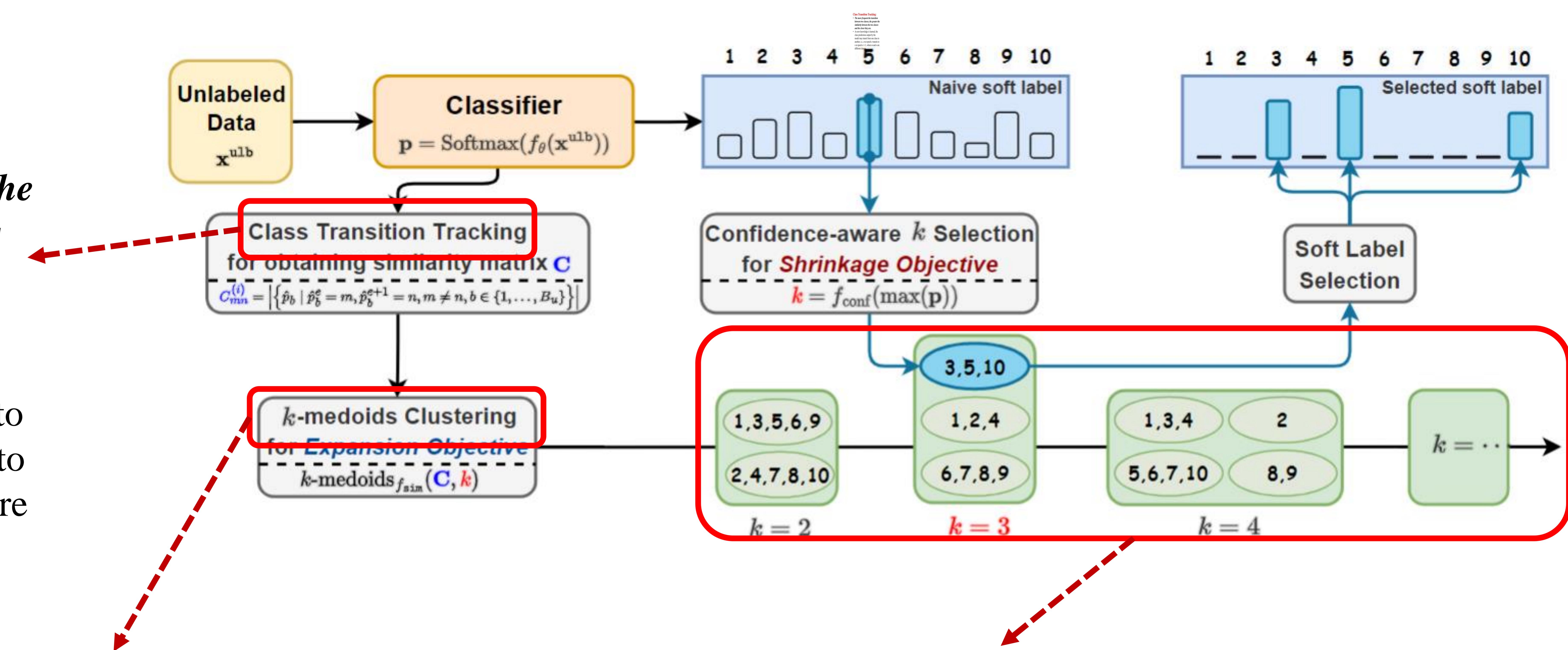
$$C_{mn}^{(i)} = \left| \left\{ \hat{p}_b \mid \hat{p}_b^e = m, \hat{p}_b^{e+1} = n, m \neq n, b \in \{1, \dots, B_u\} \right\} \right|$$

## Method

- For **Expansion Objective**, use class-level similarity obtained through class transfer tracking to acquire candidate class clusters.

### Class Transition Tracking:

- The more frequent the transition between two classes, the greater the similarity between the two classes and the closer they are.*
- As new knowledge is learned, the class predictions output by the model may transit from one class to another, i.e.,  $m$  at epoch  $e$  transits to  $n$  at epoch  $e + 1$ , where  $m$  and  $n$  are different classes.

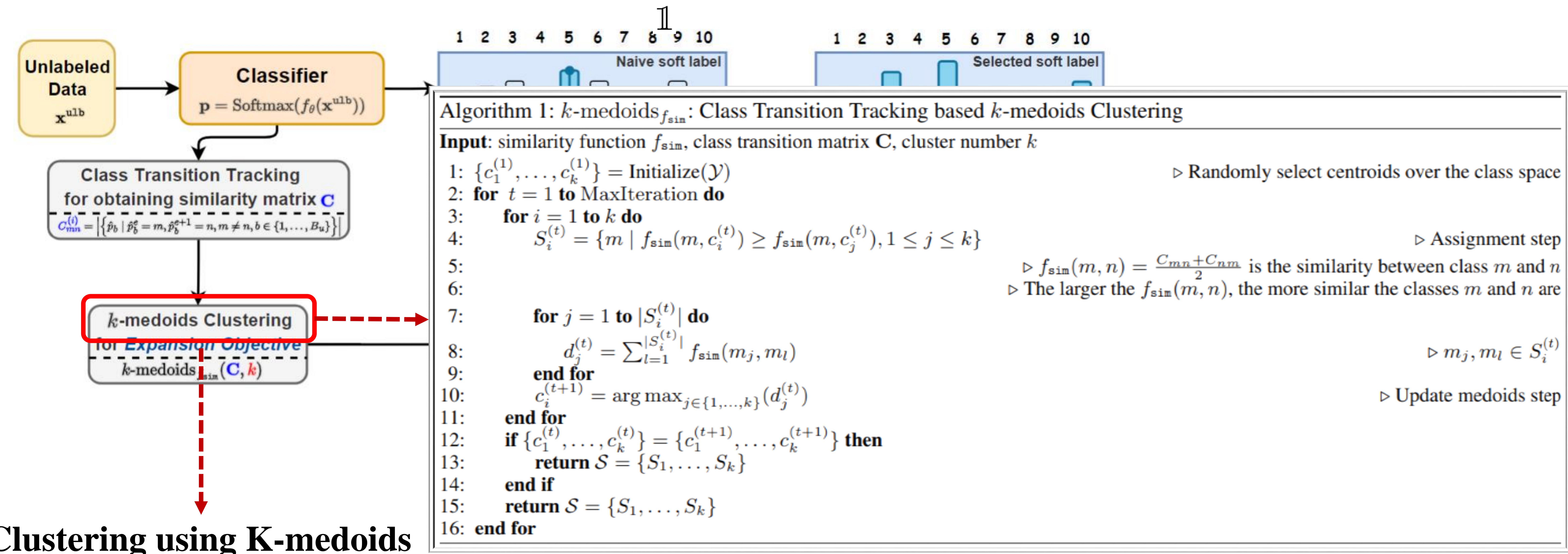


$$\mathcal{S} = k\text{-medoids}_{f_{sim}}(\mathbf{C}, k)$$

$$\mathcal{C}_i = S_s, s = \arg \max_{c \in \{1, \dots, k\}} \mathbb{1}(\hat{p}_i \in S_c)$$

## Method

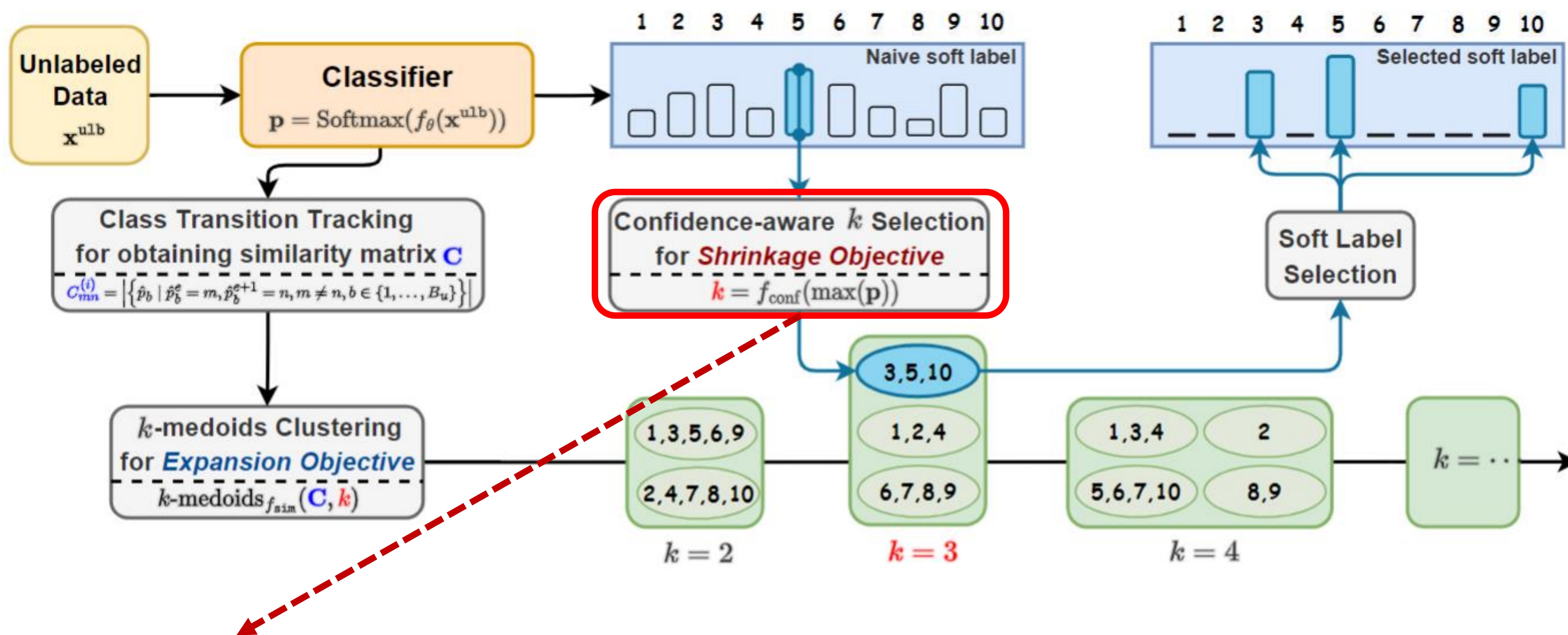
- For **Expansion Objective**, use class-level similarity obtained through class transfer tracking to acquire candidate class clusters.



$$\mathcal{S} = k\text{-medoids}_{f_{sim}}(\mathbf{C}, k)$$

## Method

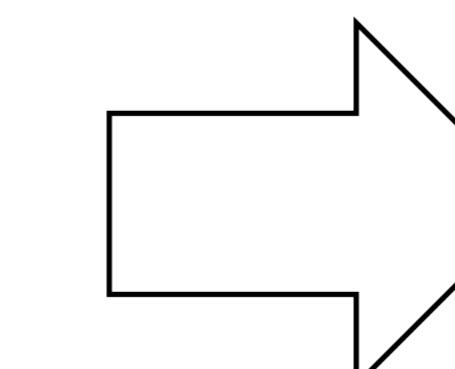
- For **Shrinkage Objective**, select a more appropriate **clustering granularity** based on prediction confidence.



The number of clusters is selected based on confidence, i.e., the granularity of the clustering.

In this implementation, a simple and effective linear function is used.

$$k_i = f_{\text{conf}}(\max(p_i)),$$



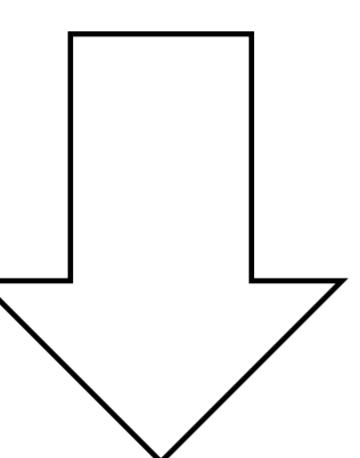
$$k_i = \lceil \left( \frac{\max(p_i)}{\alpha} + \frac{2}{K} \right) \times K - \frac{1}{2} \rceil,$$

# Roll With the Punches: Expansion and Shrinkage of Soft Label Selection for Semi-supervised Fine-Grained Learning

## Method

- Putting It All Together

$$\mathcal{L}_{\text{sup}} = \frac{1}{B} \sum_{n=1}^B H(\mathbf{y}_i, f_\theta(\text{Aug}_w(\mathbf{x}_i^{\text{ulb}}))),$$

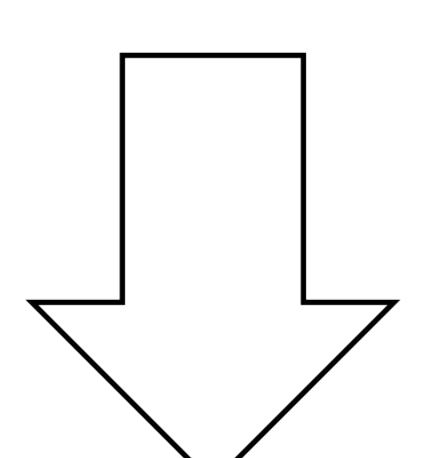


**Step 1.** The soft pseudo-label for weakly-augmented  $\mathbf{x}_i^{\text{ulb}}$  is computed as  $\mathbf{p}_i^w = \text{Softmax}(f_\theta(\text{Aug}_w(\mathbf{x}_i^{\text{ulb}})))$ . With obtained  $\mathbf{p}_i^w$ , the  $k_i^w$  for clustering is selected by Eq. (13).

**Step 2.** Class transition matrix  $\mathbf{C}$  is calculated by Eq. (7) in the current iteration. By Eq. (8), the CTT based  $k$ -medoids clustering is performed with  $k_i^w$  and  $\mathbf{C}$  to obtain  $\mathcal{C}_i^w$ .

**Step 3.**  $\mathbf{g}_i^w$  is computed by Eq. (3) with  $\mathcal{C}_i^w$  and then the selected pseudo-label  $\tilde{\mathbf{p}}_i^w$  is computed by Eq. (4) with  $\mathbf{g}_i^w$ .

$$\mathcal{L}_{\text{cos}} = \frac{1}{\mu B} \sum_{i=1}^{\mu B} H(\tilde{\mathbf{p}}_i^w, f_\theta(\text{Aug}_s(\mathbf{x}_i^{\text{ulb}}))).$$



$$\mathcal{L} = \mathcal{L}_{\text{sup}} + \lambda_{\text{cos}} \mathcal{L}_{\text{cos}},$$

### Algorithm 2: SoC: Soft Label Selection with Confidence-Aware Clustering based on Class Transition Tracking

**Input:** class number:  $K$ , labeled data set:  $\mathcal{D}^{\text{lb}} = \{(\mathbf{x}_i^{\text{lb}}, \mathbf{y}_i)\}_{i=1}^N$ , unlabeled data set:  $\mathcal{D}^{\text{ulb}} = \{\mathbf{x}_i^{\text{ulb}}\}_{i=1}^M$ , model:  $f_\theta$ , weak and strong augmentation:  $\text{Aug}_w$  and  $\text{Aug}_s$ , class prediction bank:  $\{l_i\}_{i=1}^M$ , class tracking matrices:  $\{\mathbf{C}^{(i)}\}_{i=1}^{N_b}$ , CTT based  $k$ -medoids Clustering:  $k$ -medoids  $f_{\text{sim}}$ , mapping function from confidence to  $k$ :  $f_{\text{conf}}$

```

1: for  $t = 1$  to MaxIteration do
2:   Sample labeled data batch  $\{(\mathbf{x}_i^{\text{lb}}, \mathbf{y}_i)\}_{i=1}^B \subset \mathcal{D}^{\text{lb}}$ 
3:   Sample unlabeled data batch  $\{\mathbf{x}_i^{\text{ulb}}\}_{i=1}^{\mu B} \subset \mathcal{D}^{\text{ulb}}$ 
4:    $\mathcal{L}_{\text{sup}} = \frac{1}{B} \sum_{i=1}^B H(\mathbf{y}_i, f_\theta(\mathbf{x}_i^{\text{lb}}))$  ▷ Compute the supervised loss
5:   for  $i = 1$  to  $\mu B$  do
6:      $d = \text{Index}(\mathbf{x}_i^{\text{ulb}})$  ▷ Obtain the index of  $\mathbf{x}_i^{\text{ulb}}$ 
7:      $\mathbf{p}_i^w = \text{Softmax}(f_\theta(\text{Aug}_w(\mathbf{x}_i^{\text{ulb}})))$  ▷ Compute soft pseudo-label for weakly-augmented  $\mathbf{x}_i^{\text{ulb}}$ 
8:      $\hat{p}_i = \arg \max(\mathbf{p}_i^w)$  ▷ Compute class prediction
9:     if  $l_d \neq \hat{p}_i$  then
10:       $C_{l_d \hat{p}_i}^{(n)} = C_{l_d \hat{p}_i}^{(n)} + 1$  ▷ Perform class transition tracking  $l_d = \hat{p}_i$ 
11:    end if
12:     $k_i = f_{\text{conf}}(\max(\mathbf{p}_i^w))$ 
13:     $\{S_1, \dots, S_k\} = k\text{-medoids}_{f_{\text{sim}}} \left( \text{Average}(\{\mathbf{C}^{(i)}\}_{i=1}^{N_b}), k_i \right)$ 
14:     $\mathcal{C}_i = S_s, s = \arg \max_{c \in \{1, \dots, k\}} \mathbb{1}(\hat{p}_i \in S_c)$ 
15:     $\mathbf{g}_i = (\mathbb{1}(1 \in \mathcal{C}_i), \dots, \mathbb{1}(K \in \mathcal{C}_i))$ 
16:     $\tilde{\mathbf{p}}_i = \text{Normalize}(\mathbf{g}_i \circ \mathbf{p}_i)$  ▷ See Algorithm 1
17:  end for
18:   $\mathcal{L}_{\text{cos}} = \frac{1}{\mu B} \sum_{i=1}^{\mu B} H(\tilde{\mathbf{p}}_i^w, f_\theta(\text{Aug}_s(\mathbf{x}_i^{\text{ulb}})))$  ▷ Compute label selection indicator  $\mathbf{g}_i$ 
19:  return  $\mathcal{L} = \mathcal{L}_{\text{sup}} + \lambda_{\text{cos}} \mathcal{L}_{\text{cos}}$  ▷ Compute selected soft label
20: end for ▷ Compute the consistency loss
▷ Optimize total loss

```

## Theoretical Support

The **Shrinkage Objective** is equivalent to the **entropy minimization** objective (see the full version of this paper [2]).

In SoC, for  $\mathbf{x}_i^{\text{ulb}}$ , we obtain the set of candidate classes  $\mathcal{C}_i$ , the prediction probabilities  $\mathbf{p}_i$  and the selected soft label  $\tilde{\mathbf{p}}_i$ . First, we resort the vector sequence  $\mathbf{p}_i = (p_{i,(1)}, \dots, p_{i,(K)})$  to  $(p_{i,(a_i)}), \dots, p_{i,(|C_i|)}, p_{i,(b_1)}, \dots, p_{i,(b_{|\mathcal{Y} \setminus C_i|})}$ , where  $K$  is the number of classes,  $a_{j_a} \in \mathcal{C}_i$  and  $b_{j_b} \in \mathcal{Y} \setminus \mathcal{C}_i$ , i.e., we put the probabilities of unselected classes at the back end of the vector sequence. Meanwhile, the subsequences  $(p_{i,(a_1)}, \dots, p_{i,(a_{|C_i|})})$  and  $(p_{i,(b_1)}, \dots, p_{i,(b_{|\mathcal{Y} \setminus C_i|})})$  in  $\mathbf{p}_i$  are sorted by value in descending order, i.e.,  $\max(p_{i,(a_{j_a})}) = p_{i,(a_1)} = p_{i,(1)}$  and  $\max(p_{i,(b_{j_b})}) = p_{i,(b_1)} = p_{i,(|C_i|+1)}$ . Then, denoting  $d = \sum_{c=|C_i|+1}^K p_{i,(c)}$ , we prove a equivalent form of Lemma 1:

$$\begin{aligned} -\sum_{c=1}^{|C_i|} p_{i,(c)} \log p_{i,(c)} &\geq -\sum_{c=1}^{|C_i|} \frac{p_{i,(c)}}{1-d} \log \frac{p_{i,(c)}}{1-d} \\ &= -\sum_{c=1}^{|C_i|} \frac{p_{i,(c)}}{1-d} \log p_{i,(c)} + \sum_{c=1}^{|C_i|} \frac{p_{i,(c)}}{1-d} \log(1-d) \\ &= -\sum_{c=1}^{|C_i|} \frac{p_{i,(c)}}{1-d} \log p_{i,(c)} + \log(1-d). \end{aligned} \quad (18)$$

Rewriting Eq. (18) we obtain its equivalent:

$$\begin{aligned} \sum_{c=1}^{|C_i|} \frac{p_{i,(c)}}{1-d} \log p_{i,(c)} - \sum_{c=1}^{|C_i|} p_{i,(c)} \log p_{i,(c)} &\geq \sum_{c=|C_i|+1}^K p_{i,(c)} \log p_{i,(c)} + \log(1-d) \\ \left(\frac{d}{1-d}\right) \sum_{c=1}^{|C_i|} p_{i,(c)} \log p_{i,(c)} &\geq \sum_{c=|C_i|+1}^K p_{i,(c)} \log p_{i,(c)} + \log(1-d). \end{aligned} \quad (19)$$

When  $d = 0$ , Eq. (19) obviously holds. Given  $d = \sum_{c=|C_i|+1}^K p_{i,(c)} > 0$ , according to the property of  $h(x) = x \log x$ , when only one entry  $p_{i,(z)}$  of  $(p_{i,(|C_i|+1)}, \dots, p_{i,(K)})$  has a non-zero value,  $\sum_{c=|C_i|+1}^K p_{i,(c)} \log p_{i,(c)}$  is maximized, i.e., the term on the right side of the inequality Eq. (19) obtains the maximum. Given  $\max(p_{i,(|C_i|+1)}, \dots, p_{i,(K)}) = p_{i,(|C_i|+1)}$ , the non-zero entry  $p_{i,(z)}$  is obviously  $p_{i,(|C_i|+1)}$ . In this case, we have  $d = p_{i,(|C_i|+1)}$  and  $\sum_{c=|C_i|+1}^K p_{i,(c)} \log p_{i,(c)} = p_{i,(|C_i|+1)} \log p_{i,(|C_i|+1)}$ . Thus, we rewrite Eq. (19) as

$$\begin{aligned} \left(\frac{p_{i,(|C_i|+1)}}{1-p_{i,(|C_i|+1)}}\right) \sum_{c=1}^{|C_i|} p_{i,(c)} \log p_{i,(c)} &\geq p_{i,(|C_i|+1)} \log p_{i,(|C_i|+1)} + \log(1-p_{i,(|C_i|+1)}) \\ \left(\frac{p_{i,(|C_i|+1)}}{1-p_{i,(|C_i|+1)}}\right) (p_{i,(1)} \log p_{i,(1)} + \sum_{c=2}^{|C_i|} p_{i,(c)} \log p_{i,(c)}) &\geq p_{i,(|C_i|+1)} \log p_{i,(|C_i|+1)} + \log(1-p_{i,(|C_i|+1)}). \end{aligned} \quad (20)$$

It's worth noting that we have  $\arg \max(\mathbf{p}_i) \in \mathcal{C}_i$  by Eq. (9), therefore we obtain  $\max(\mathbf{p}_i) \in \{p_{i,(1)}, \dots, p_{i,(|C_i|)}\}$ , i.e.,  $p_{i,(1)} \geq p_{i,(|C_i|+1)}$ . In other words, there is no constraint on the value of  $\sum_{c=2}^{|C_i|} p_{i,(c)} \log p_{i,(c)}$  when  $p_{i,(1)}$  and  $p_{i,(|C_i|+1)}$  are given. Thus, by Jensen Inequality we know that when  $p_{i,(2)} = p_{i,(3)} = \dots = p_{i,(|C_i|)} = \frac{1-p_{i,(1)}-p_{i,(|C_i|+1)}}{|C_i|-1}$ ,  $\sum_{c=2}^{|C_i|} p_{i,(c)} \log p_{i,(c)}$  obtains the minimum, i.e., the left side of the inequality Eq. (20) obtains the minimum with the given  $p_{i,(1)}$  and  $p_{i,(|C_i|+1)}$  where  $p_{i,(1)} \geq p_{i,(|C_i|+1)}$ . In this case, we can rewrite Eq. (20) as

$$\begin{aligned} \frac{p_{i,(|C_i|+1)}}{1-p_{i,(|C_i|+1)}} p_{i,(1)} \log p_{i,(1)} + \frac{p_{i,(|C_i|+1)}(|C_i|-1)}{1-p_{i,(|C_i|+1)}} \times \frac{1-p_{i,(1)}-p_{i,(|C_i|+1)}}{|C_i|-1} \log \frac{1-p_{i,(1)}-p_{i,(|C_i|+1)}}{|C_i|-1} \\ \geq p_{i,(|C_i|+1)} \log p_{i,(|C_i|+1)} + \log(1-p_{i,(|C_i|+1)}). \end{aligned} \quad (21)$$

By Eq. (21), we obtain

$$\begin{aligned} \frac{p_{i,(|C_i|+1)}}{1-p_{i,(|C_i|+1)}} p_{i,(1)} \log p_{i,(1)} + \frac{p_{i,(|C_i|+1)}(1-p_{i,(1)}-p_{i,(|C_i|+1)})}{1-p_{i,(|C_i|+1)}} [\log(1-p_{i,(1)}-p_{i,(|C_i|+1)}) - \log(|C_i|-1)] \\ \geq p_{i,(|C_i|+1)} \log p_{i,(|C_i|+1)} + \log(1-p_{i,(|C_i|+1)}) - \\ \frac{p_{i,(|C_i|+1)}}{1-p_{i,(|C_i|+1)}} p_{i,(1)} \log p_{i,(1)} + \frac{p_{i,(|C_i|+1)}(1-p_{i,(1)}-p_{i,(|C_i|+1)})}{1-p_{i,(|C_i|+1)}} \log(1-p_{i,(1)}-p_{i,(|C_i|+1)}) - \\ p_{i,(|C_i|+1)} \log p_{i,(|C_i|+1)} - \log(1-p_{i,(|C_i|+1)}) \geq \frac{p_{i,(|C_i|+1)}(1-p_{i,(1)}-p_{i,(|C_i|+1)})}{1-p_{i,(|C_i|+1)}} \log(|C_i|-1) \\ \frac{p_{i,(1)}}{(1-p_{i,(1)}-p_{i,(|C_i|+1)})} \log p_{i,(1)} + \log(1-p_{i,(1)}-p_{i,(|C_i|+1)}) - \\ \frac{(1-p_{i,(|C_i|+1)})(p_{i,(|C_i|+1)} \log p_{i,(|C_i|+1)} + \log(1-p_{i,(|C_i|+1)}))}{p_{i,(|C_i|+1)}(1-p_{i,(1)}-p_{i,(|C_i|+1)})} \geq \log(|C_i|-1). \end{aligned} \quad (22)$$

Letting  $x = p_{i,(1)}$  and  $y = p_{i,(|C_i|+1)}$ , we define the function

$$f(x, y) = \frac{x}{1-x-y} \log x + \log(1-x-y) - \frac{(1-y)(y \log y + \log(1-y))}{y(1-x-y)}. \quad (23)$$

Then we can obtain

$$\frac{\partial f(x, y)}{\partial x} = \frac{(y-1)(y \log(x) - \log(1-y) - y \log(y))}{y(1-x-y)^2} \geq 0. \quad (24)$$

Given  $x \geq y$  (due to  $p_{i,(1)} \geq p_{i,(|C_i|+1)}$ ),  $f(x, y)$  is minimized when  $x = y$ . Plugging  $x = y$  into  $f(x, y)$ , we obtain

$$f(y) = \frac{y}{1-2y} \log y + \log(1-2y) - \frac{(1-y)(y \log y + \log(1-y))}{y(1-2y)} \quad (25)$$

and the first derivative of  $f(y)$ :

$$f'(y) = \frac{(2y^2 - 4y + 1) \log(1-y)}{(1-2y)^2 y^2}. \quad (26)$$

Solving  $f'(y) = 0$  we obtain  $y = 1 - \frac{\sqrt{2}}{2}$ . It is easy to verify that  $f'(y)$  exists for all  $y$  such that  $0 < y < \frac{1}{2}$  and the sign of  $f'(y)$  changes from negative to positive. Thus,  $f(y)$  has a minimum at  $y = 1 - \frac{\sqrt{2}}{2}$ , i.e.,

$$\min f(y) = f\left(1 - \frac{\sqrt{2}}{2}\right) - \frac{1}{4} \left( \sqrt{2} \log(16) + \log(64) - 4 \log\left(1 - \frac{\sqrt{2}}{2}\right) + 4 \log(\sqrt{2}-1) \right) \approx 2.36655. \quad (27)$$

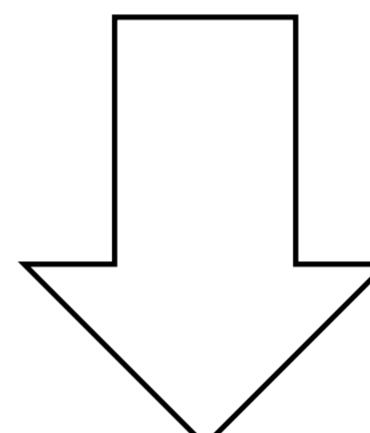
Obviously, we have  $\log 10 < f\left(1 - \frac{\sqrt{2}}{2}\right) < \log 11$ . By Eq. (27), we know that Eq. (22) holds when  $|\mathcal{C}_i| \leq 11$ , i.e.,  $z_i^{\text{obj2}} \leq 11$  (by  $z_i^{\text{obj2}} = \sum_{c=1}^K g_{i,(c)}$  and Eq. (3)). It is worth noting that the strict requirement for  $z_i^{\text{obj2}}$  is to ensure that Eq. (22) holds in the worst case, where the class distribution of *selected classes* excluding  $\arg \max(\mathbf{p}_i)$  is “uniform distribution” (i.e., for any  $a, b \in \mathcal{C}_i \wedge a, b \neq \arg \max(\mathbf{p}_i)$ ,  $p_{i,(a)} = p_{i,(b)}$ ) and the class distribution of *unselected classes* is “one-hot distribution” (i.e.,  $|\{p_{i,(z)} \mid p_{i,(z)} \neq 0 \wedge z \in \mathcal{Y} \setminus \mathcal{C}_i\}| = 1$ ). In the vast majority of cases, any  $z_i^{\text{obj2}}$  will make Eq. (22). Meanwhile, given an opposite  $\alpha$ , thanks to Confidence-Aware  $k$  Selection in Sec. 3.3,  $z_i^{\text{obj2}} \leq 11$  for the worst case also can be guaranteed after a few training iterations because the model becomes more confident on  $\mathbf{x}_i^{\text{ulb}}$  (implying a smaller  $z_i^{\text{obj2}}$ ).

In fact, the mentioned extreme situation is almost impossible to occur in reality. For specific, with training, the prediction distribution outputted by the model will tend to be non-uniform unless the model learns nothing. Thus, the requirement for  $z_i^{\text{obj2}}$  will be greatly relaxed. To verify this point, we show an intuitive example: even if the class distribution of selected classes is uniform distribution (i.e., for any  $a, b \in \mathcal{C}_i$ ,  $p_{i,(a)} = p_{i,(b)}$ ), Lemma 1 holds no matter what value  $z_i^{\text{obj2}}$  takes.

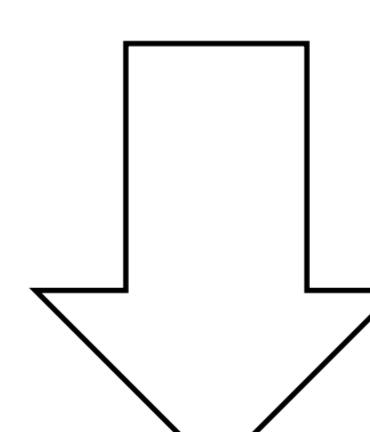
**Lemma 1.** Given  $\tilde{\mathbf{p}}_i$  in SoC (implying  $\mathcal{C}_i \subsetneq \mathcal{Y}$  and  $\hat{p}_i \in \mathcal{C}_i$ ), we show that the entropy of  $\tilde{\mathbf{p}}_i$  is smaller than that of  $\mathbf{p}_i$ :

$$\mathcal{H}(\tilde{\mathbf{p}}_i) \leq \mathcal{H}(\mathbf{p}_i), \quad (11)$$

where  $\mathcal{H}(\cdot)$  refers to the entropy.



**Proof for Theorem 1..** Given  $\mathcal{C}_i^{(2)} \subsetneq \mathcal{C}_i^{(1)}$ , we can treat  $\tilde{\mathbf{p}}_i^{(1)}$  as  $\mathbf{p}_i$  and  $\tilde{\mathbf{p}}_i^{(2)}$  as  $\tilde{\mathbf{p}}_i$  in Lemma 1, i.e., treat the previously obtained  $\tilde{\mathbf{p}}_i$  as the naive soft pseudo-label. Thus,  $\mathcal{H}(\tilde{\mathbf{p}}_i^{(2)}) \leq \mathcal{H}(\tilde{\mathbf{p}}_i^{(1)})$  holds. For any  $\mathcal{C}_i^{(j)} \subsetneq \mathcal{C}_i^{(j-1)}$ , we can repeat the above proof to obtain  $\mathcal{H}(\tilde{\mathbf{p}}_i^{(j)}) \leq \mathcal{H}(\tilde{\mathbf{p}}_i^{(j-1)})$ . Thus, we can obtain  $\mathcal{H}(\tilde{\mathbf{p}}_i^{(m)}) \leq \dots \leq \mathcal{H}(\tilde{\mathbf{p}}_i^{(1)})$ .  $\square$



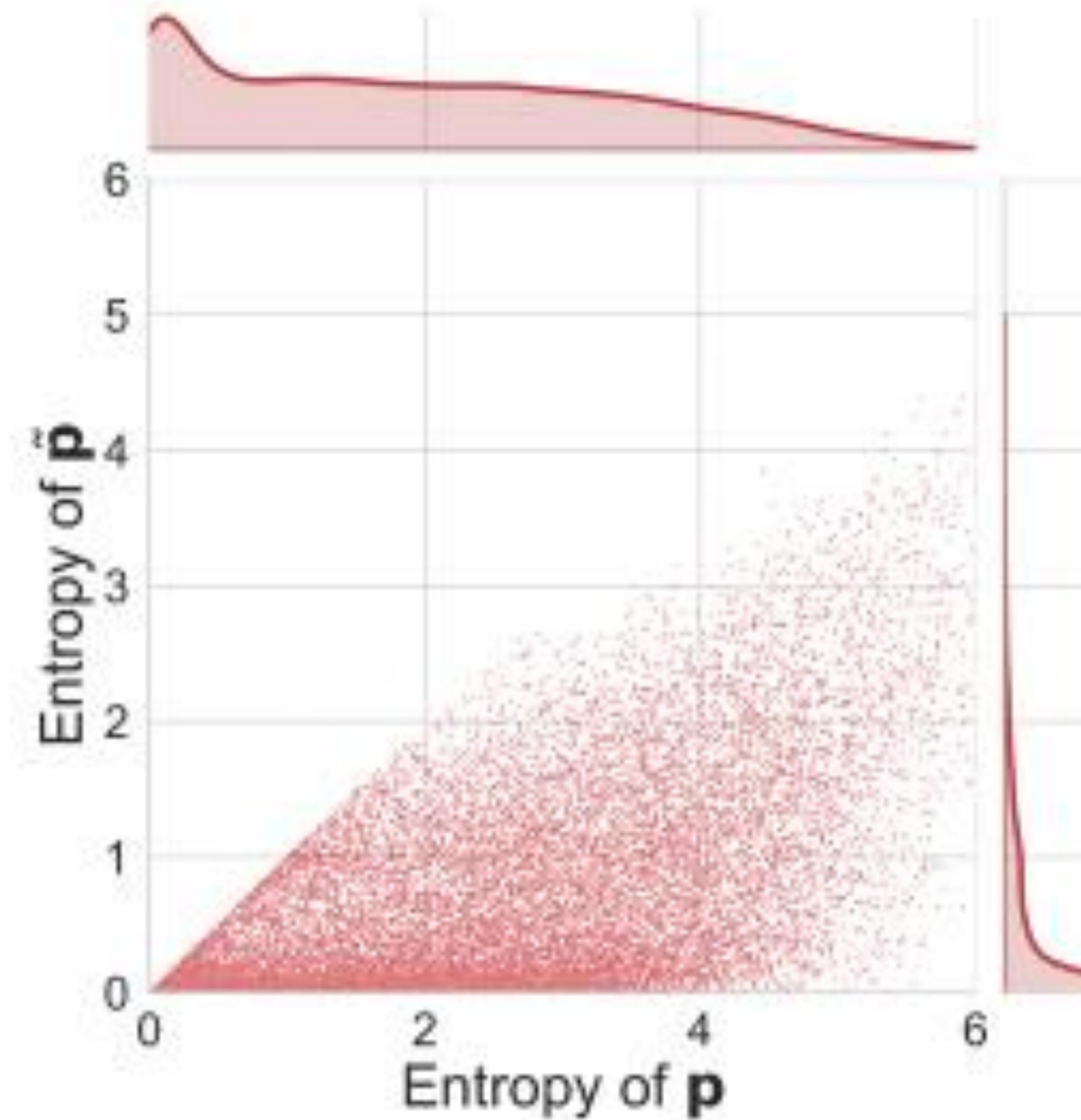
**Theorem 1.** In SoC, minimizing  $|\mathcal{C}_i^{(1)}|$  to  $|\mathcal{C}_i^{(m)}|$ :  $\mathcal{C}_i^{(m)} \subsetneq \dots \mathcal{C}_i^{(2)} \subsetneq \mathcal{C}_i^{(1)}$  (i.e.,  $\sum_{c=1}^K g_{i,(c)}^{(m)} < \dots < \sum_{c=1}^K g_{i,(c)}^{(2)} < \sum_{c=1}^K g_{i,(c)}^{(1)}$ ), we show that the entropy of  $\mathbf{p}_i$  is minimizing:

$$\mathcal{H}(\tilde{\mathbf{p}}_i^{(m)}) \leq \mathcal{H}(\tilde{\mathbf{p}}_i^{(m-1)}) \leq \dots \mathcal{H}(\tilde{\mathbf{p}}_i^{(2)}) \leq \mathcal{H}(\tilde{\mathbf{p}}_i^{(1)}), \quad (10)$$

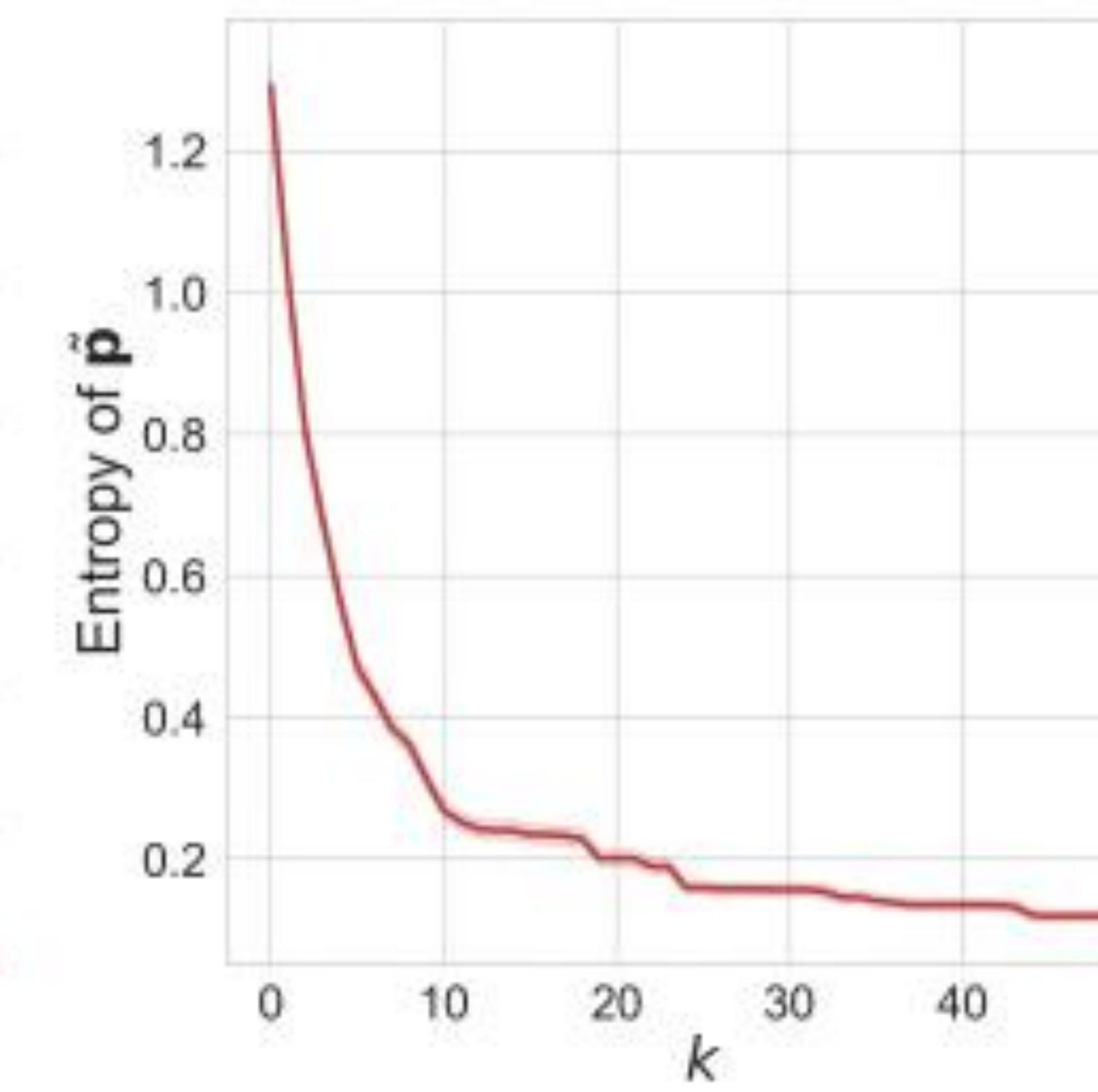
where  $\mathcal{H}(\cdot)$  refers to the entropy.

## Experimental Support

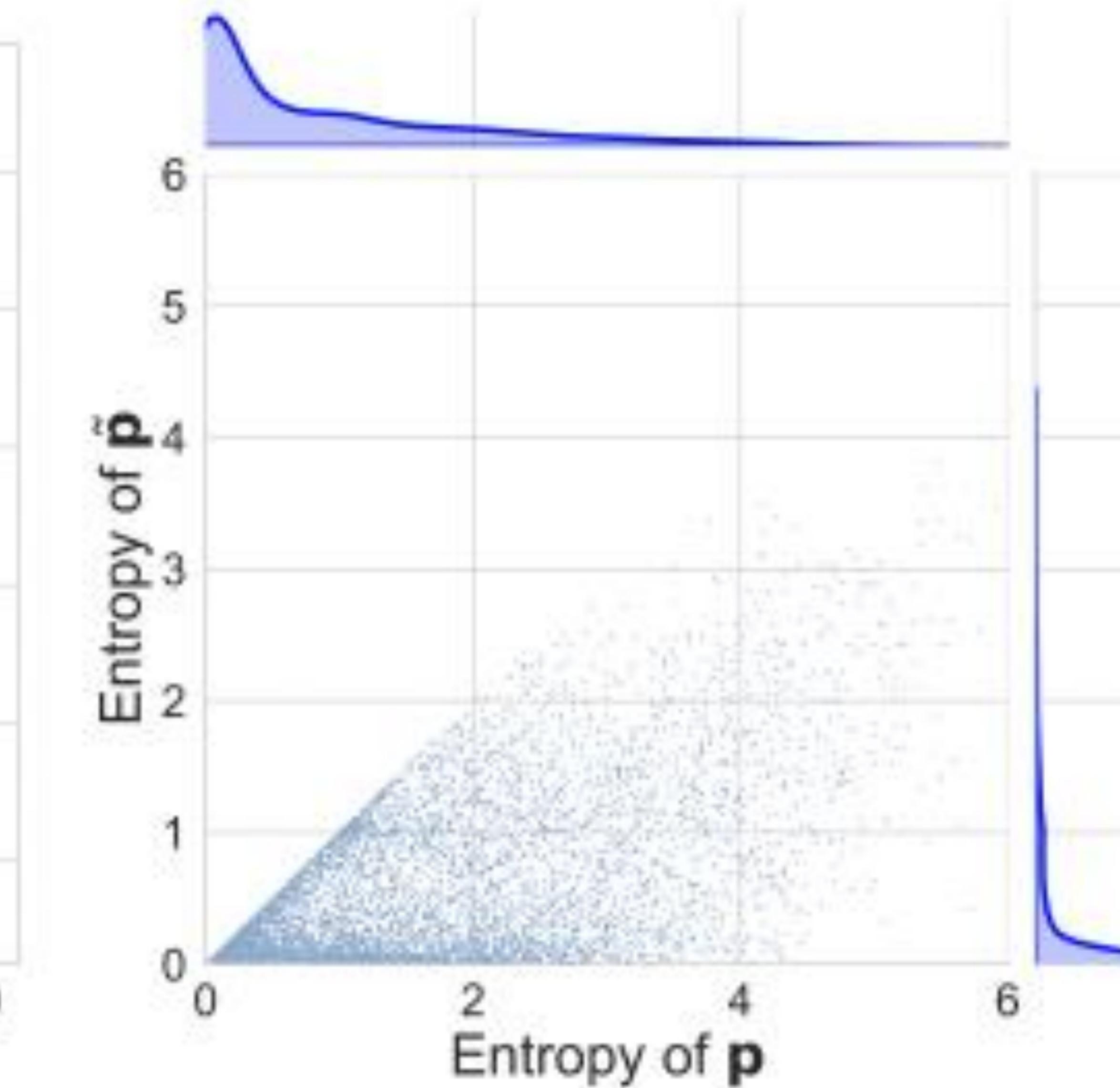
The **Shrinkage Objective** is equivalent to the **entropy minimization** objective (see the full version of this paper [2]).



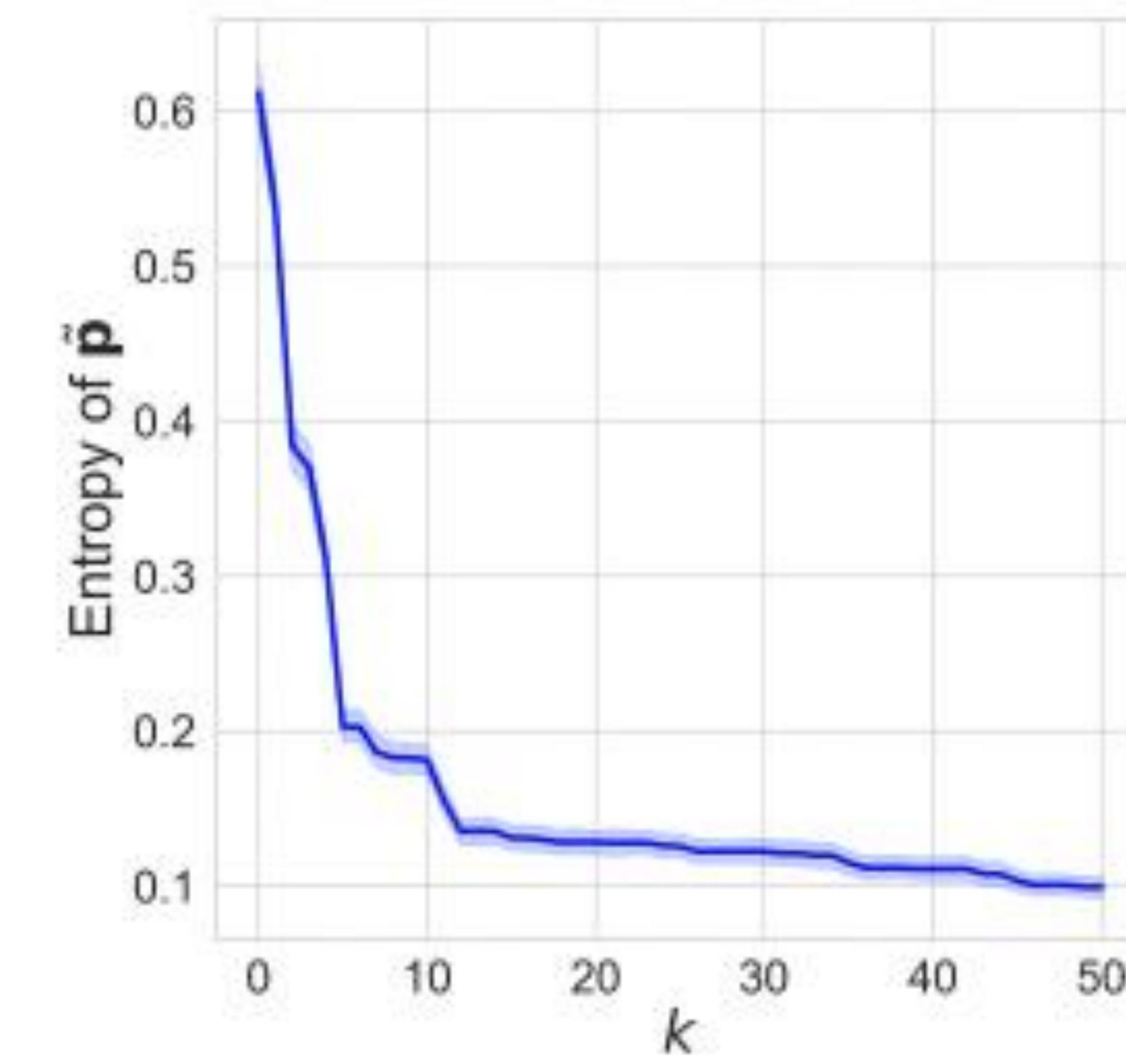
(a) Semi-Aves



(b) Semi-Aves



(c) Semi-Fungi

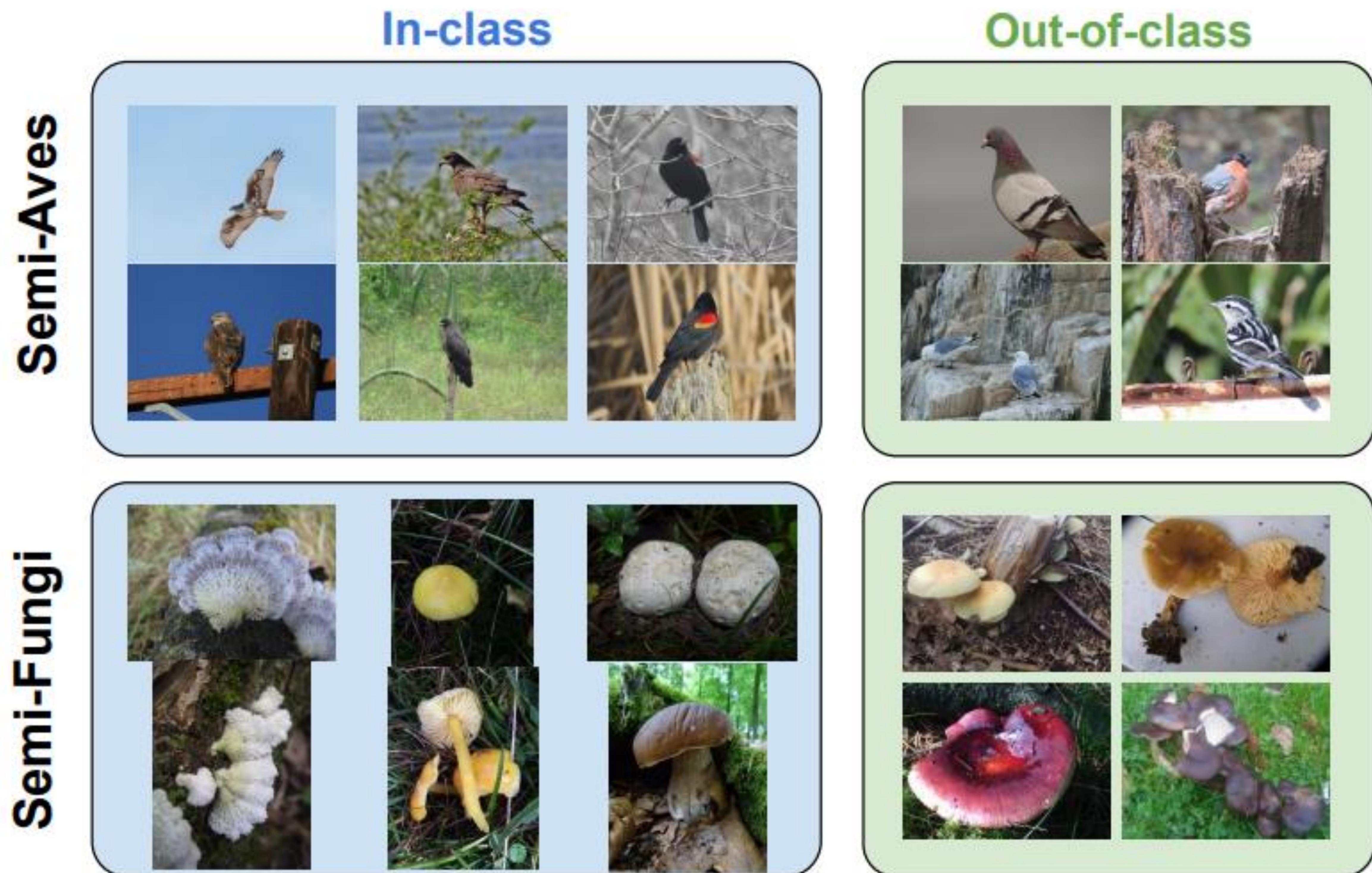
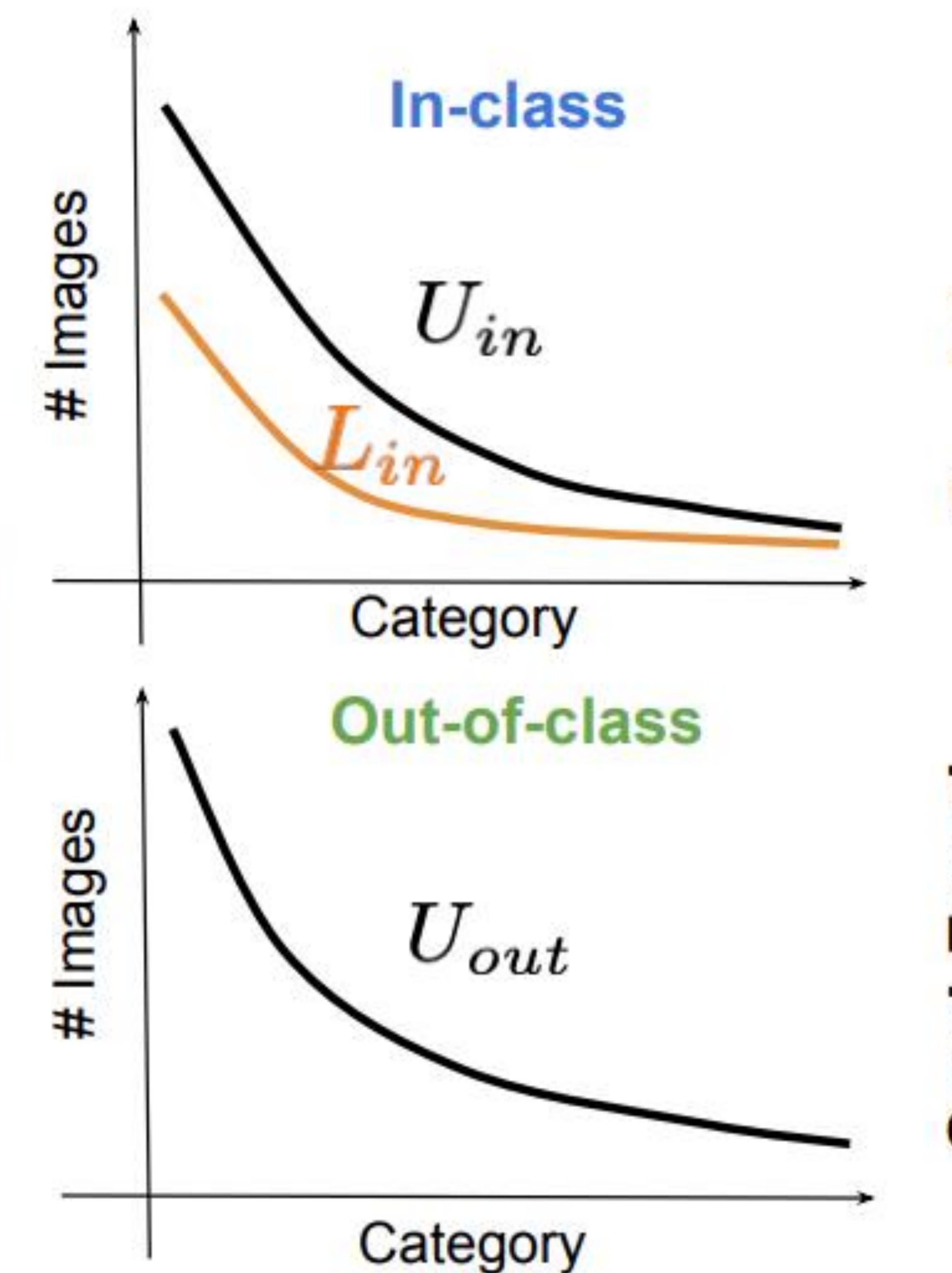
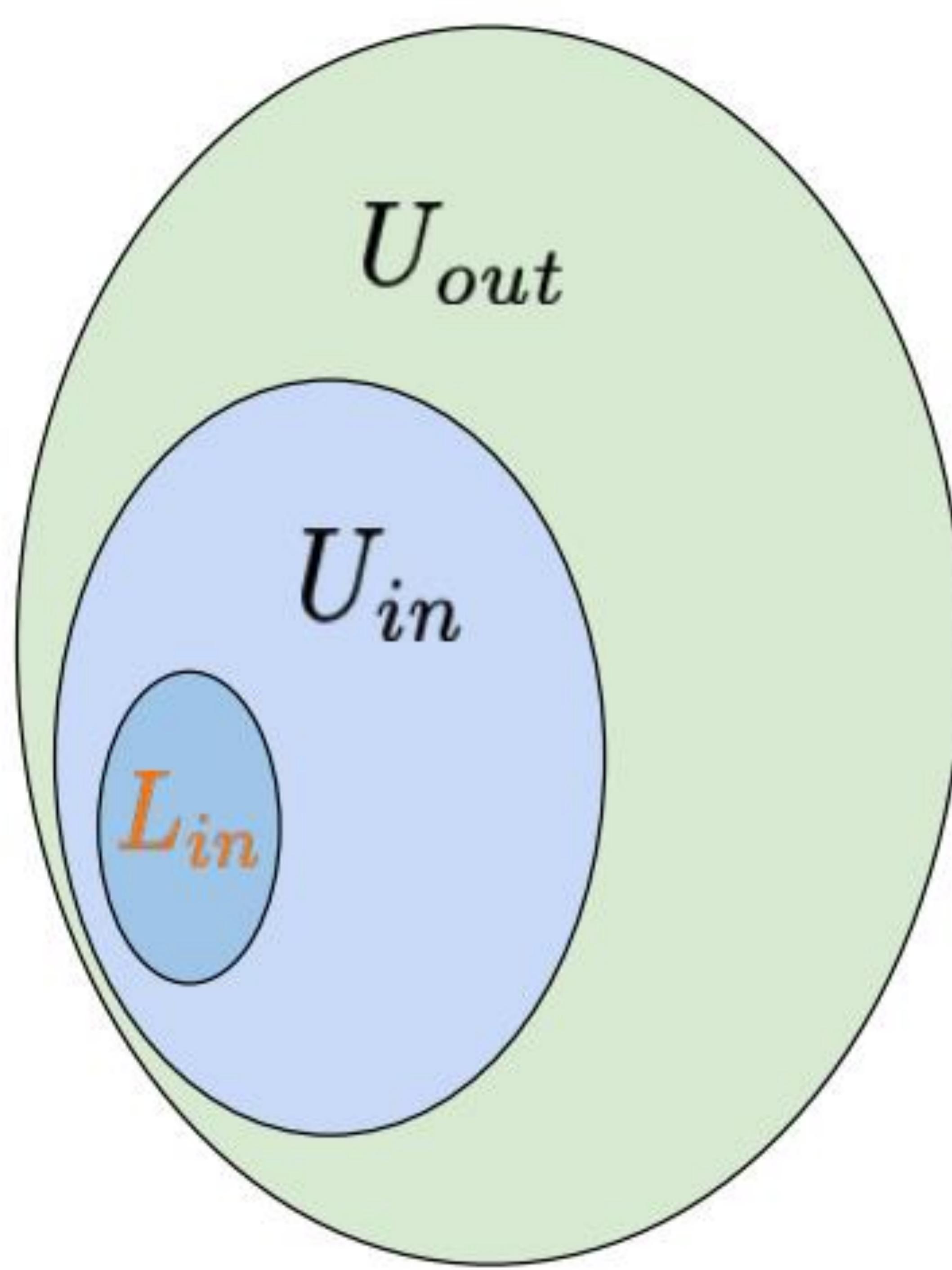


(d) Semi-Fungi

- The entropy of the pseudo-labels produced by our method is smaller than that of the original pseudo-labels.
- **The larger  $k$  is, the smaller the entropy is.**

## Experiments

- Datasets: Semi-Aves (200 classes) and Semi-Fungi (200 classes) [1]



- Semi-Aves** is divided into the training set and validation set with a total of 5,959 labeled images and 26,640 unlabeled images, and the test set with 8,000 images.
- Semi-Fungi** is divided into the training set and validation set with a total of 4,141 labeled images and 13,166 unlabeled images, and the test set with 4,000 images.

[1] A realistic evaluation of semi-supervised learning for fine-grained classification. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*.



## Experiments

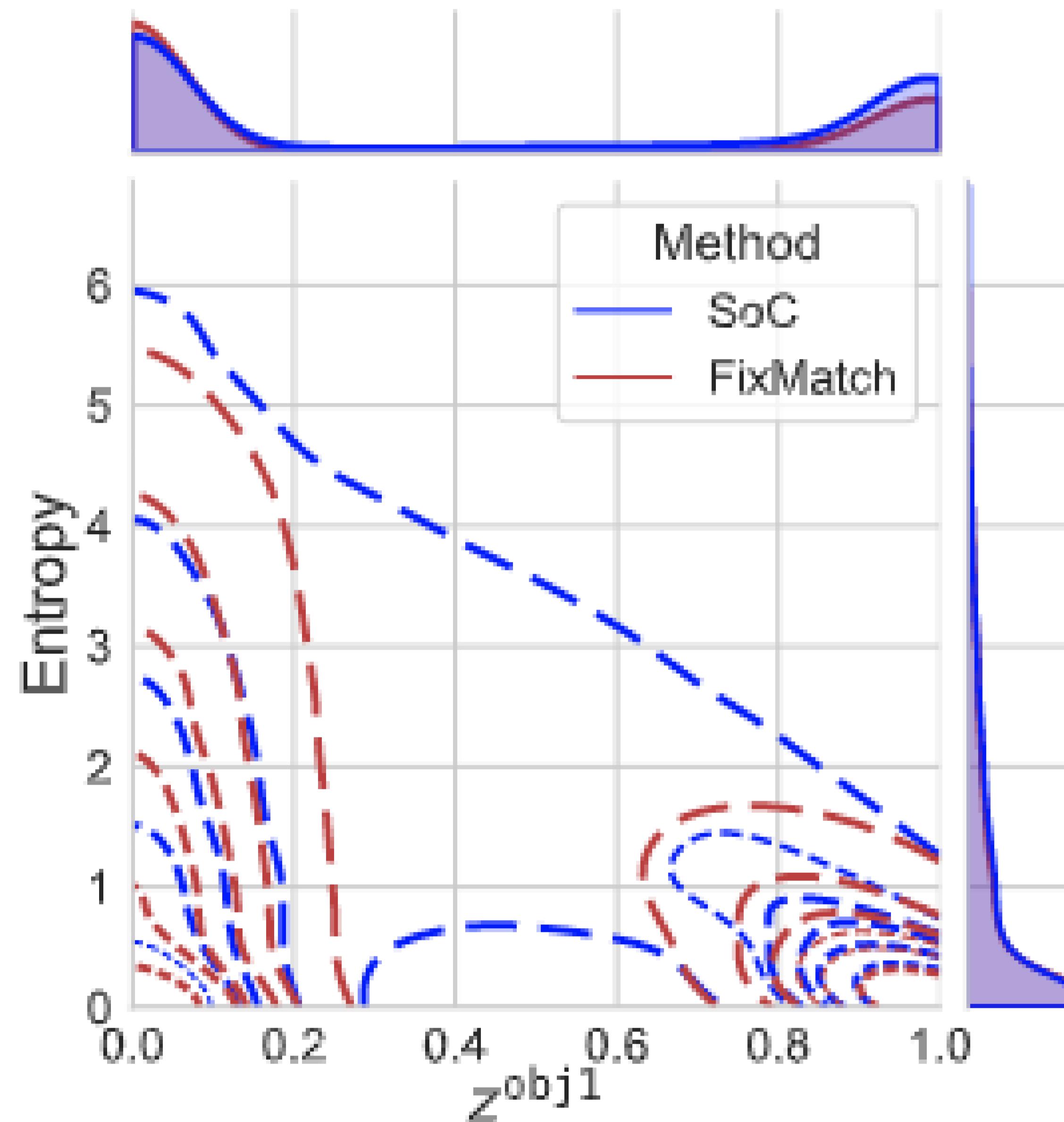
- SS-FGVC results on Semi-Aves and Semi-Fungi

Dataset	Pseudo-label	Method	Year	from scratch		from ImageNet		from iNat	
				Top1	Top5	Top1	Top5	Top1	Top5
Semi-Aves	Hard label	Supervised oracle	—	57.4±0.3	79.2±0.1	68.5±1.4	88.5±0.4	69.9±0.5	89.8±0.7
		MoCo (He et al. 2020)	CVPR' 20	28.2±0.3	53.0±0.1	52.7±0.1	78.7±0.2	68.6±0.1	87.7±0.1
		Pseudo-Label (Lee et al. 2013)	ICML' 13	16.7±0.2	36.5±0.8	54.4±0.3	78.8±0.3	65.8±0.2	86.5±0.2
		Curriculum Pseudo-Label (Cascante-Bonilla et al. 2021)	AAAI' 21	20.5±0.5	41.7±0.5	53.4±0.8	78.3±0.5	69.1±0.3	87.8±0.1
		FixMatch (Sohn et al. 2020)	NIPS' 20	28.1±0.1	51.8±0.6	57.4±0.8	78.5±0.5	70.2±0.6	87.0±0.1
	Soft label	FlexMatch (Zhang et al. 2021)*	NIPS' 21	27.3±0.5	49.7±0.8	53.4±0.2	77.9±0.3	67.6±0.5	87.0±0.2
		MoCo + FlexMatch*	NIPS' 21	35.0±1.2	58.5±1.0	53.4±0.4	77.0±0.2	68.9±0.3	87.7±0.2
		KD-Self-Training (Su, Cheng, and Maji 2021)	CVPR' 21	22.4±0.4	44.1±0.1	55.5±0.1	79.8±0.1	67.7±0.2	87.5±0.2
		MoCo + KD-Self-Training (Su, Cheng, and Maji 2021)	CVPR' 21	31.9±0.1	56.8±0.1	55.9±0.2	80.3±0.1	70.1±0.2	88.1±0.1
		SimMatch (Zheng et al. 2022)*	CVPR' 22	24.8±0.5	48.1±0.6	53.3±0.5	77.9±0.8	65.4±0.2	86.9±0.3
Semi-Fungi	Hard label	MoCo + SimMatch*	CVPR' 22	32.9±0.4	57.9±0.3	53.7±0.2	78.8±0.5	65.7±0.3	87.1±0.2
		SoC	Ours	31.3±0.8 ( $\uparrow$ 11.4%)	55.3±0.7 ( $\uparrow$ 16.8%)	57.8±0.5 ( $\uparrow$ 0.7%)	80.8±0.5 ( $\uparrow$ 1.3%)	71.3±0.3 ( $\uparrow$ 1.6%)	88.8±0.2 ( $\uparrow$ 1.1%)
		MoCo + SoC	Ours	39.3±0.2 ( $\uparrow$ 12.3%)	62.4±0.4 ( $\uparrow$ 6.7%)	58.0±0.4 ( $\uparrow$ 3.8%)	81.7±0.4 ( $\uparrow$ 1.7%)	70.8±0.4 ( $\uparrow$ 1.0%)	88.9±0.5 ( $\uparrow$ 0.9%)
		Supervised oracle	—	60.2±0.8	83.3±0.9	73.3±0.1	92.5±0.3	73.8±0.3	92.4±0.3
		MoCo (He et al. 2020)	CVPR' 20	33.6±0.2	59.4±0.3	55.2±0.2	82.9±0.2	52.5±0.4	79.5±0.2
	Soft label	Pseudo-Label (Lee et al. 2013)	ICML' 13	19.4±0.4	43.2±1.5	51.5±1.2	81.2±0.2	49.5±0.4	78.5±0.2
		Curriculum Pseudo-Label (Cascante-Bonilla et al. 2021)	AAAI' 21	31.4±0.6	55.0±0.6	53.7±0.2	80.2±0.1	53.3±0.5	80.0±0.5
		FixMatch (Sohn et al. 2020)	NIPS' 20	32.2±1.0	57.0±1.2	56.3±0.5	80.4±0.5	58.7±0.7	81.7±0.2
		FlexMatch (Zhang et al. 2021)*	NIPS' 21	36.0±0.9	59.9±1.1	59.6±0.5	82.4±0.5	60.1±0.6	82.2±0.5
		MoCo + FlexMatch*	NIPS' 21	44.2±0.6	67.0±0.8	59.9±0.8	82.8±0.7	61.4±0.6	83.2±0.4
[3] He, K.; Fan, H.; Wu, Y.; Xie, S.; and Girshick, R. 2020. Momentum contrast for unsupervised visual representation learning. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> .	[3] He, K.; Fan, H.; Wu, Y.; Xie, S.; and Girshick, R. 2020. Momentum contrast for unsupervised visual representation learning. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> .	KD-Self-Training (Su, Cheng, and Maji 2021)	CVPR' 21	32.7±0.2	56.9±0.2	56.9±0.3	81.7±0.2	55.7±0.3	82.3±0.2
		MoCo + KD-Self-Training (Su, Cheng, and Maji 2021)	CVPR' 21	39.4±0.3	64.4±0.5	58.2±0.5	84.4±0.2	55.2±0.5	82.9±0.2
		SimMatch (Zheng et al. 2022)*	CVPR' 22	36.5±0.9	61.7±1.0	56.6±0.4	81.8±0.6	56.7±0.3	80.9±0.4
		MoCo + SimMatch*	CVPR' 22	42.2±0.5	67.0±0.4	56.5±0.2	82.5±0.3	57.4±0.2	81.3±0.4
		SoC	Ours	39.4±2.3 ( $\uparrow$ 7.9%)	62.5±1.1 ( $\uparrow$ 1.3%)	61.4±0.4 ( $\uparrow$ 3.0%)	83.9±0.6 ( $\uparrow$ 1.8%)	62.4±0.2 ( $\uparrow$ 3.8%)	85.1±0.2 ( $\uparrow$ 3.5%)
		MoCo + SoC	Ours	47.2±0.5 ( $\uparrow$ 6.8%)	71.3±0.2 ( $\uparrow$ 6.4%)	61.9±0.3 ( $\uparrow$ 3.3%)	85.8±0.2 ( $\uparrow$ 3.6%)	62.5±0.4 ( $\uparrow$ 1.8%)	84.7±0.2 ( $\uparrow$ 1.8%)

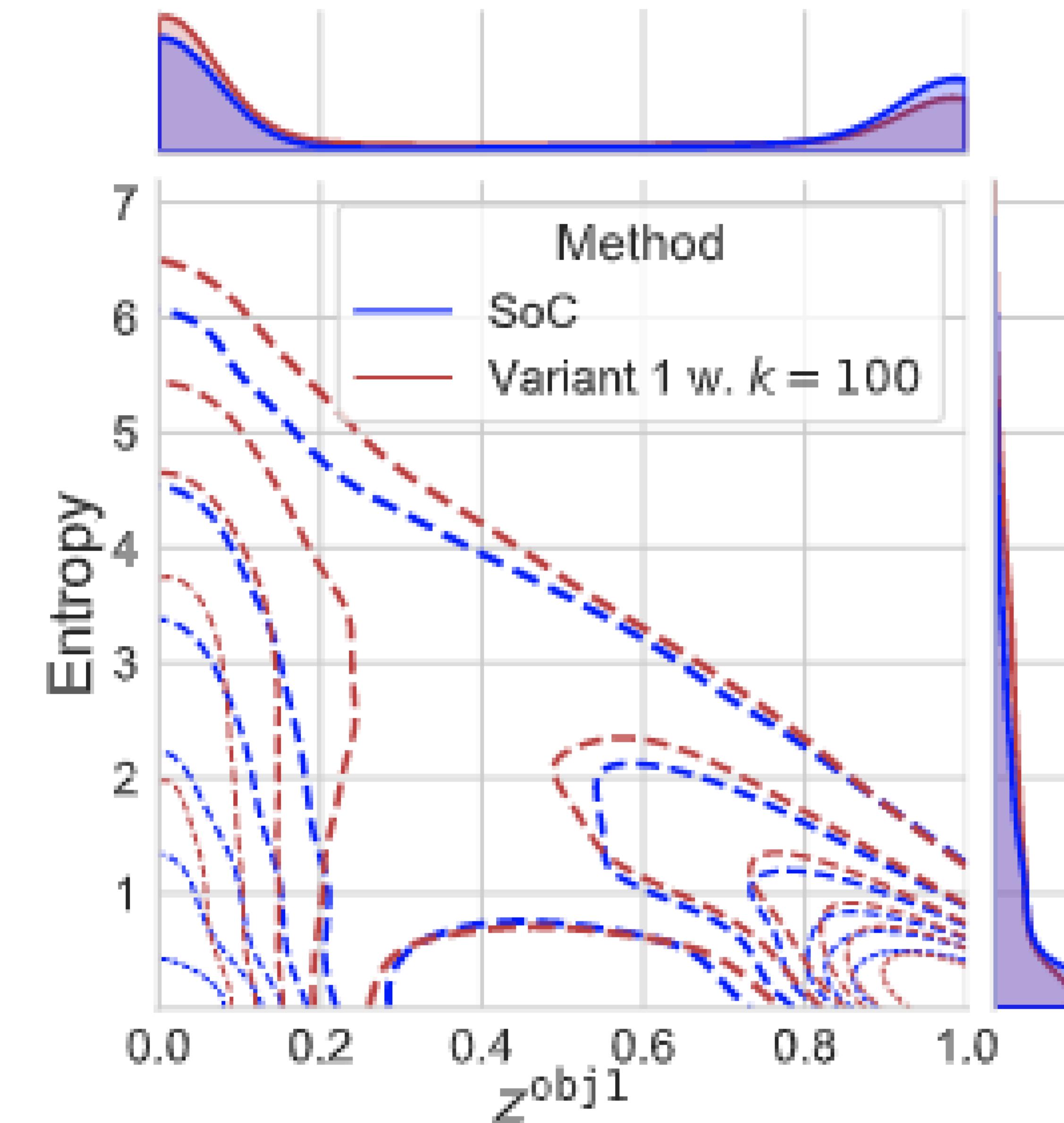
- The models are trained from scratch, or ImageNet/iNat pre-trained or the model initialized with MoCo [3] learning on the unlabeled data.

## Experiments

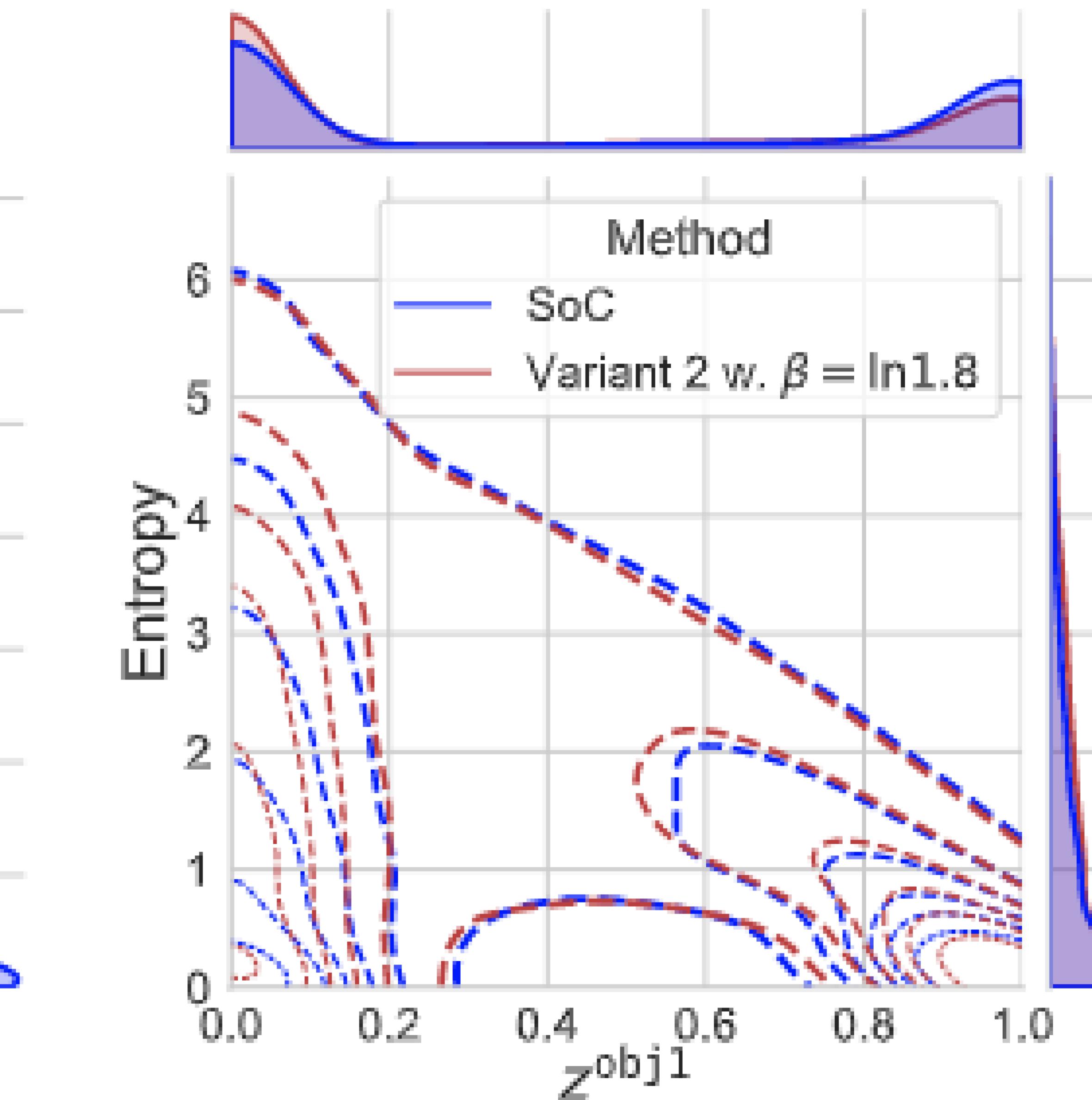
- Ablation studies: comparison with different SoC variants**



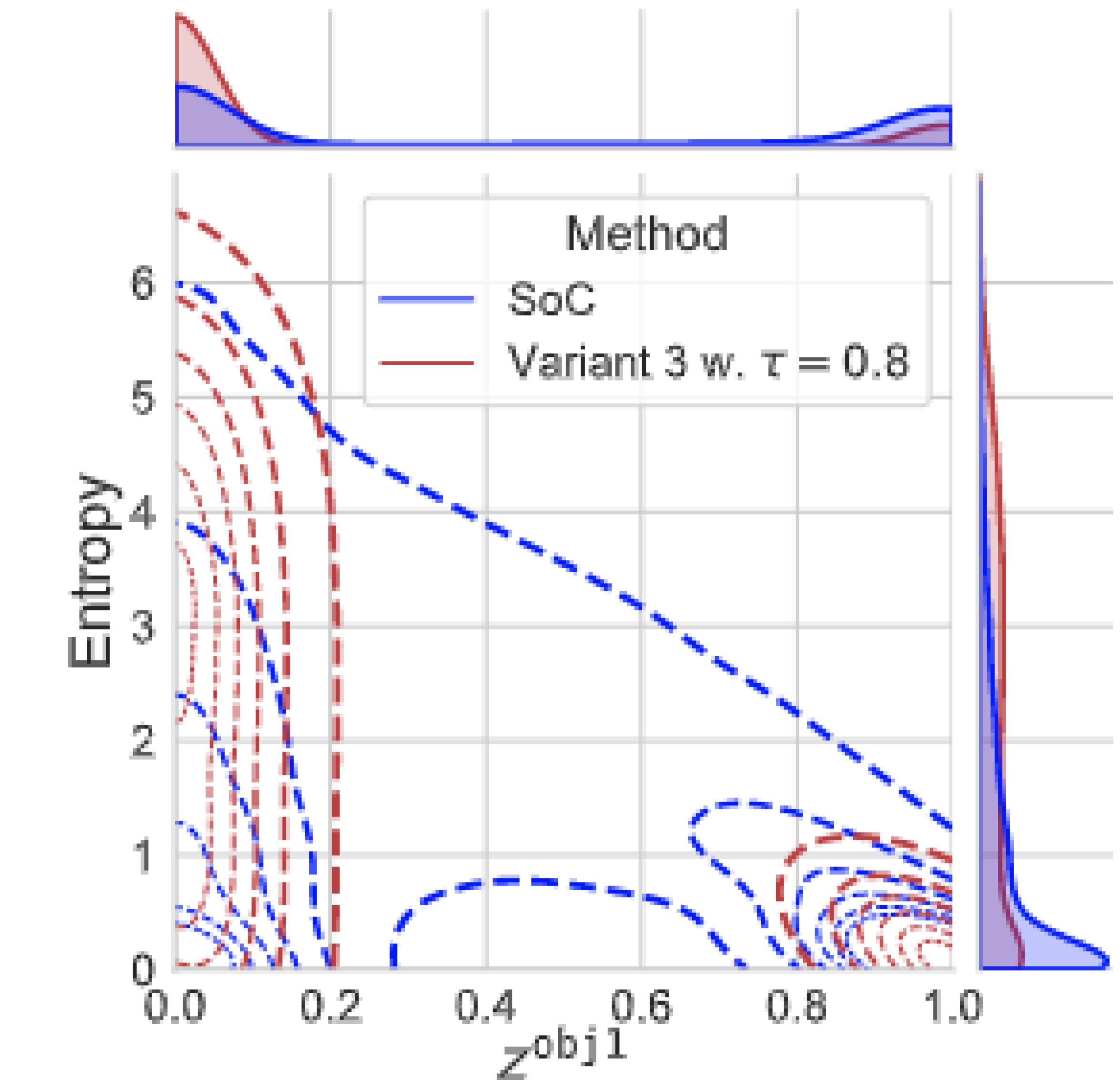
(a) SoC vs. FixMatch



(b) SoC vs. Variant 1



(c) SoC vs. Variant 2



(d) SoC vs. Variant 3

- SoC vs. FixMatch:** Highlighting the performance advantage of soft label.
- Variant 1:** Fixing the number of clusters.
- Variant 2:** Adjusting the number of clusters non-linearly based on confidence.
- Variant 3:** Filtering pseudo-labels with a fixed threshold.

# Roll With the Punches: Expansion and Shrinkage of Soft Label Selection for Semi-supervised Fine-Grained Learning



# Thanks!

Arxiv: <https://arxiv.org/abs/2312.12237>

Github: <https://github.com/NJUyued/SoC4SS-FGVC>

Homepage: <https://njuyued.github.io/>