

# SPICE: Semantic Pseudo-Labeling for Image Clustering

{ instance-level similarity  
 cluster-level discrepancy

$X = \{x_i\}_{i=1}^N$   $N$  images  $\xrightarrow{\text{cluster}}$   $K$  classes

two parts {
 

- a feature model: images  $\rightarrow$  feature vectors  
 $f_i = F(x_i; \theta_F)$  *measure the similarity among samples*
- a clustering head: feature vectors  $\rightarrow$  probabilities  
 $P_i = C(f_i; \theta_C)$  *identify the discrepancy between clusters*

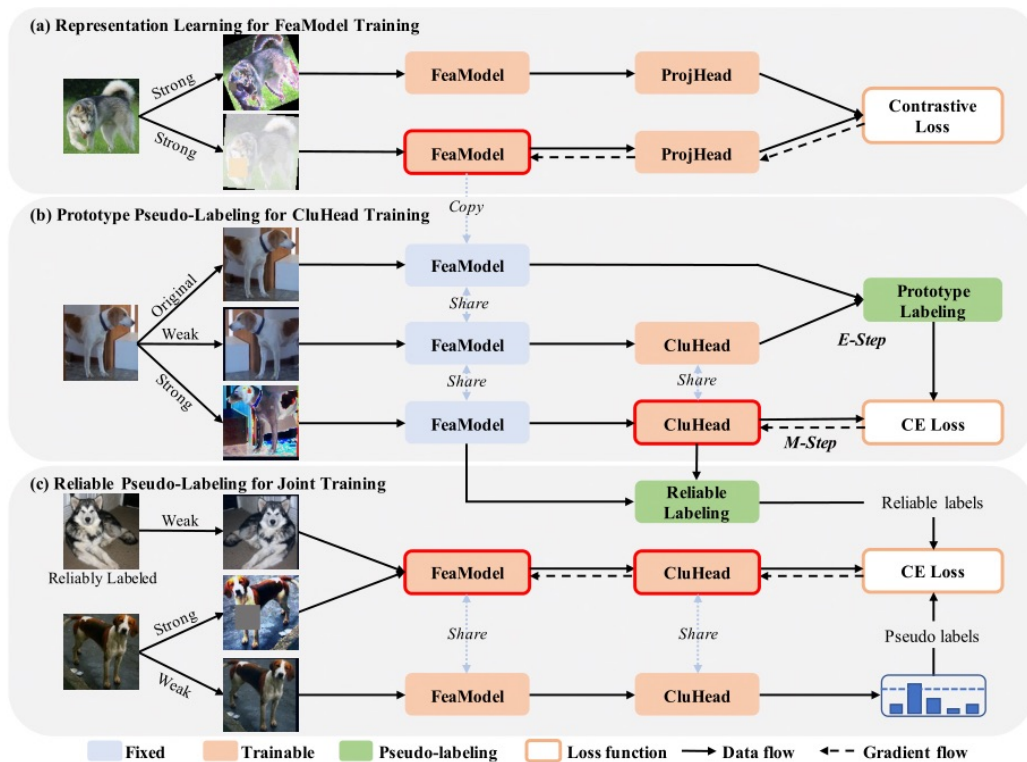
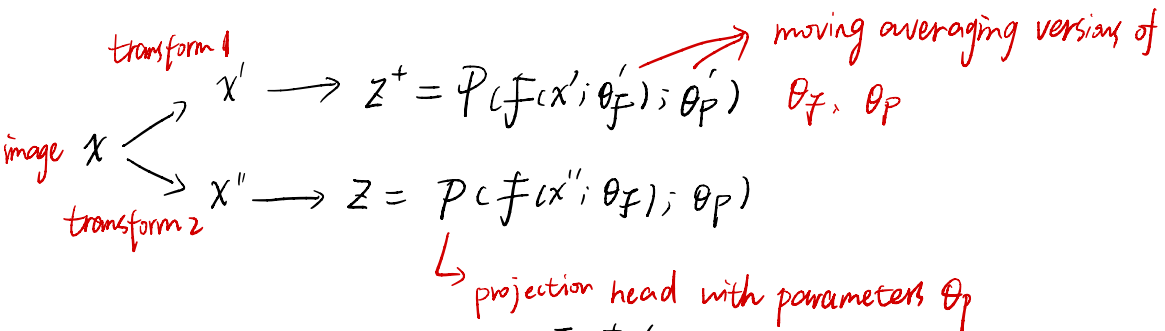


Fig. 3. Illustration of the SPICE framework. (a) Train the feature model with the contrastive learning based unsupervised representation learning. (b) Train the clustering head via the prototype pseudo-labeling algorithm in an EM framework. (c) Jointly train the feature model and the clustering head through the reliable pseudo-labeling algorithm.

three stages:

### A. Feature Model Training with contrastive learning



$$\mathcal{L}_{\text{fea}} = -\log \left( \frac{\exp(z^T z^+ / \tau)}{\sum_{i=1}^{N_q} \exp(z^T z_i^- / \tau) + \exp(z^T z^+ / \tau)} \right) \quad N_q: \text{queue size}$$

Finally, the optimized feature model parameters are denoted as  $\theta_f^s$ .

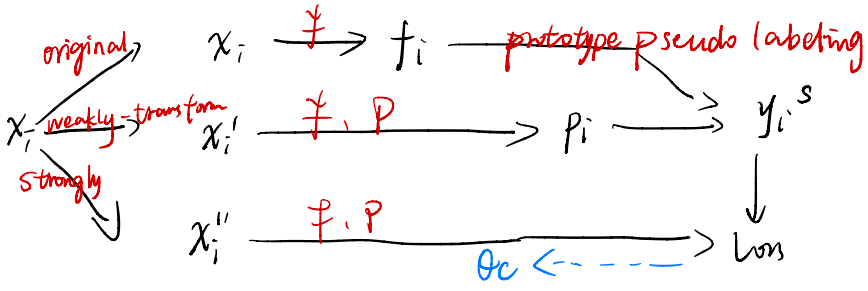
### B. Clustering Head Training with Prototype Pseudo-Labeling

- the parameters of clustering head  $C$ ,  $\theta_c$   $P_i = C(f_i; \theta_c)$
- the cluster labels  $\{y_i^s\}$  of  $X$  over  $K$  clusters, where  $f_i = f(x_i; \theta_f^s)$

EM framework:

Expectation (E) step  $\rightarrow$  solve  $\{y_i^s\}$  given  $\theta_c$

Maximization (M) step  $\rightarrow$  solve  $\theta_c$  given  $\{y_i^s\}$



B-1: Prototype Pseudo-Labeling (E-step)

top branch:  $F = [f_1, f_2, \dots, f_m]^T \in \mathbb{R}^{m \times D}$  for  $\mathcal{X}_b$  Given

middle branch:  $P = [p_1, p_2, \dots, p_m]^T \in \mathbb{R}^{m \times K}$  for  $\alpha(\mathcal{X}_b)$

Select the top confident samples for  $k$ -th cluster:

$$\bar{\mathcal{X}}_k = \{f_i \mid i \in \text{argtop}_k(P_{:,k}, \frac{m}{K}), \forall i=1, 2, \dots, m\} \quad (2)$$

$k$ -th column of  $P$

$\text{argtop}_k(P_{:,k}, \frac{m}{K})$  returns the top  $\frac{m}{K}$  confident sample indices

eg.  $P = \begin{bmatrix} 0.7 & 0.6 & 0.2 & 0.9 \\ \hline 0.3 & 0.4 & 0.8 & 0.1 \end{bmatrix}^T$   $m=4, K=2$   $\bar{\mathcal{X}}_1 = \{f_1, f_4\}$   
 $\bar{\mathcal{X}}_2 = \{f_2, f_3\}$

The cluster centers  $\{\bar{\mathcal{X}}_k\}_{k=1}^K$ :  $\bar{\mathcal{X}}_k = \frac{K}{m} \sum_{f_i \in \bar{\mathcal{X}}_k} f_i, \forall k=1, 2, \dots, K$  (3)

balanced?

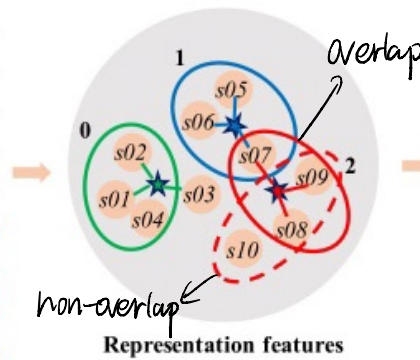
Compute the cosine similarity between  $f_i$  and  $\bar{\mathcal{X}}_k$ , select  $\frac{m}{K}$  nearest samples to  $\bar{\mathcal{X}}_k$ , denoted by  $\mathcal{X}^k \rightarrow y_i^s = k$

$$\mathcal{X}^s = \{(x_i, y_i^s) \mid \forall x_i \in \mathcal{X}^k, k=1, 2, \dots, K\} \quad (4)$$

Mini-batch of images with semantic pseudo-labels

S01	0.99	0.01	0.00
S02	0.80	0.10	0.10
S03	0.55	0.20	0.25
S04	0.40	0.15	0.45
S05	0.01	0.98	0.01
S06	0.15	0.80	0.05
S07	0.01	0.50	0.49
S08	0.01	0.02	0.97
S09	0.10	0.08	0.82
S10	0.20	0.35	0.45

Probabilities over  $K$  clusters



Representation features

0
0
-1
0
1
1
1 2
2
2
-1

Pseudo labels

## B-2: Training Clustering Head (M-step)

Given  $\mathcal{X}^S \rightarrow$  supervised learning

bottom batch: compute  $\beta(\mathcal{X}^S)$ 's probabilities

$\mathcal{C}$  can be optimized by minimizing CE loss:

$$L_{clu} = \frac{1}{M} \sum_{i=1}^M L_{ce}(y_i^S, p_i') \quad (5)$$

where  $p_i' = \text{softmax}(p_i)$ ,  $p_i = \mathcal{A}(\mathcal{F}(\beta(x_i); \theta_{\mathcal{F}}^S); \theta_{\mathcal{A}})$

double softmax ( $p_i$  already softmax)

---

### Algorithm 1: Training Clustering Head.

---

**Input:** Dataset  $\mathcal{X} = \{x_i\}_{i=1}^N$ ,  $\theta_{\mathcal{F}}^s$ ,  $K$ ,  $M$ ,  $m$ ,  $T$ ,  $\alpha$ ,  $\beta$   
**Output:** Cluster label  $y_i^s$  of  $x_i \in \mathcal{X}$

- 1 Set feature model parameters to  $\theta_{\mathcal{F}}^s$ ,  $t = 0$ , and initialize  $\theta_{\mathcal{C}}$ ;
- 2 **while**  $t < T$  **do**
- 3     **for**  $b = 1, 2, \dots, \lfloor \frac{N}{M} \rfloor$  **do**
- 4         **E-step:**
- 5         Select  $M$  samples from  $\mathcal{X}$  as  $\mathcal{X}_b$ ;
- 6         Compute embedding features  $\mathbf{F} = \mathcal{F}(\mathcal{X}_b; \theta_{\mathcal{F}}^s)$ ;
- 7         Predict probabilities  $\mathbf{P} = \mathcal{C}(\mathcal{F}(\alpha(\mathcal{X}_b); \theta_{\mathcal{F}}^s); \theta_{\mathcal{C}})$ ;
- 8         Construct labeled image set  $\mathcal{X}^s$  with Eqs. (2), (3), and (4) ;
- 9         **M-step:**
- 10         Compute probabilities  $\mathbf{P} = \mathcal{C}(\mathcal{F}(\beta(\mathcal{X}_b); \theta_{\mathcal{F}}^s); \theta_{\mathcal{C}})$  ;
- 11         Optimize  $\theta_{\mathcal{C}}$  by minimizing Eq. (5) ;
- 12     **end**
- 13      $t \leftarrow t + 1$
- 14 **end**
- 15 Select the best clustering head with the minimum loss as  $\theta_{\mathcal{C}}^s$  ;
- 16 **foreach**  $x_i \in \mathcal{X}$  **do**
- 17      $p_i = \mathcal{C}(\mathcal{F}(x_i; \theta_{\mathcal{F}}^s); \theta_{\mathcal{C}}^s)$  ;
- 18      $y_i^s = \arg \max_k (p_i)_k$  ;
- 19 **end**

## C. Joint Training with Reliable Pseudo-Labeling

Reliable Pseudo-Labeling:

Given  $\{(x_i, f_i, y_i^s)\}_{i=1}^N$  from "B".

For each  $x_i$ , select  $N_s$  nearest samples, corresponding labels  $\mathcal{P}_i$ .

The semantically consistent fraction  $r_i$  of  $x_i$ :

$$r_i = \frac{1}{N_s} \sum_{y \in \mathcal{P}_i} \mathbb{1}(y = y_i^s)$$

Given a threshold  $\lambda$ , if  $r_i > \lambda$ ,  $(x_i, y_i^s)$  is reliably labeled  
otherwise, is ignored.

partially labeled

$$\mathcal{X}^r = \{(x_i, y_i^s) \mid r_i > \lambda, \forall i = 1, 2, \dots, N\}$$

Joint Training: Semi-supervised learning.

We adapt a simple semi-supervised learning method [54]. During training, the subset of reliably labeled samples keep fixed. On the other hand, all training samples should be consistently clustered, i.e., different transformations of the same image are constrained to have the consistent prediction. To this end, as shown in Fig. 3(c), the confidently predicted label of weak transformations is used as the pseudo-label for strong transformations of the same image. Formally, the consistency pseudo label  $y_j^u$  of the sample  $x_j$  is calculated as in Eq. (8):

$$y_j^u = \begin{cases} \arg \max(\mathbf{p}_j) & \text{if } \max(\mathbf{p}_j) \geq \eta, \\ -1 & \text{otherwise} \end{cases} \quad (8)$$

where  $\mathbf{p}_j = \mathcal{C}(\mathcal{F}(\alpha(x_j); \theta_{\mathcal{F}}); \theta_{\mathcal{C}})$ , and  $\eta$  is the confidence threshold.

Then, the whole network parameters  $\theta_{\mathcal{F}}$  and  $\theta_{\mathcal{C}}$  are optimized with the following loss function:

$$\mathcal{L}_{\text{joint}} = \frac{1}{L} \sum_{i=1}^L \underbrace{\mathcal{L}_{ce}(y_i^s, \mathcal{C}(\mathcal{F}(\alpha(x_i); \theta_{\mathcal{F}}); \theta_{\mathcal{C}}))}_{\text{partial samples with reliable pseudo-labels}} \quad (9) \\ + \frac{1}{U} \sum_{j=1}^U \underbrace{\mathbb{1}(y_j^u \geq 0) \mathcal{L}_{ce}(y_j^u, \mathcal{C}(\mathcal{F}(\beta(x_j); \theta_{\mathcal{F}}); \theta_{\mathcal{C}}))}_{\text{all samples with consistency pseudo-labels}},$$

where the first term is computed with reliably labeled samples  $(x_i, y_i^s)$  drawn from  $\mathcal{X}^r$ , and the second term is computed with pseudo-labeled samples  $(x_j, y_j^u)$  drawn from the whole dataset  $\mathcal{X}$ , which is dynamically labeled by thresholding the confident predictions as in Eq. (8).  $L$  and  $U$  denote the numbers of labeled and unlabeled images in a mini-batch.