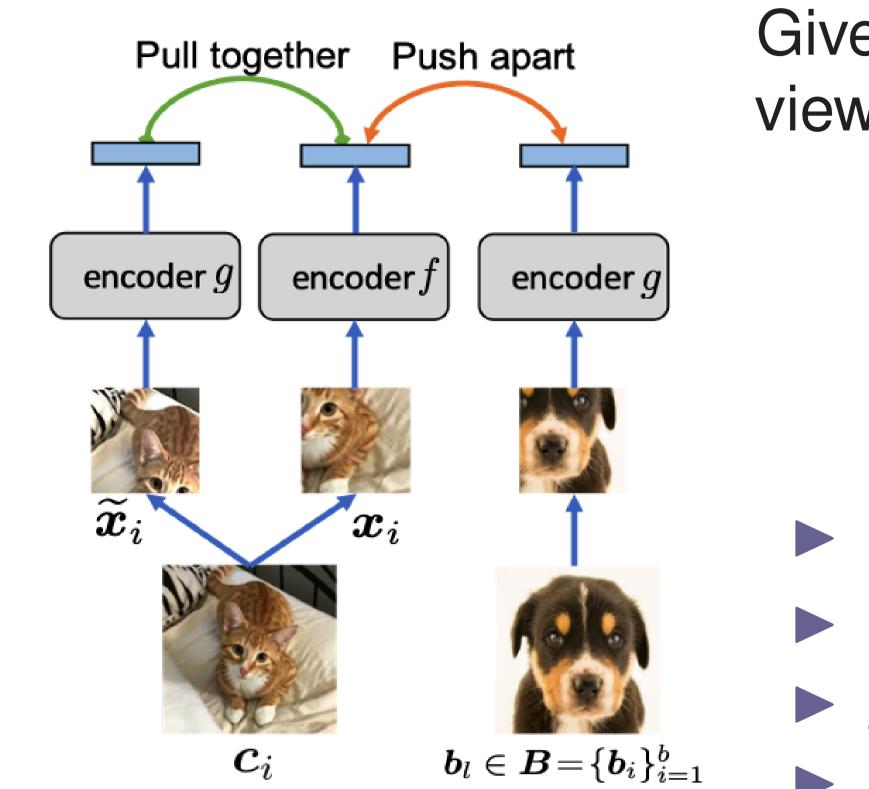
# A Theory-Driven Self-Labeling Refinement Method for Contrastive Representation Learning

## **Problem Setup**

MoCo pulls together crops of same image, pushs away crops of different images



Given a minibatch samples  $\{c_i\}_{i=1}^s$ , it augments each  $c_i$  into two views  $(\mathbf{x}_i, \widetilde{\mathbf{x}}_i)$  and optimizes:

 $\mathcal{L}_{\mathrm{n}}(oldsymbol{w})\!=\!-rac{1}{s}\!\sum_{i=1}^{s}\log\Bigl(rac{\sigma(oldsymbol{x}_{i},\widetilde{oldsymbol{x}}_{i})}{\sigma(oldsymbol{x}_{i},\widetilde{oldsymbol{x}}_{i})+\sum_{l=1}^{b}\sigma(oldsymbol{x}_{i},oldsymbol{b}_{l})}$ query positive kev

- $f_w$ : online network
- $\blacktriangleright$  g<sub>e</sub>: target network
- $\blacktriangleright$   $B = \{b_i\}_{i=1}^{b}$ : negative key buffer (previous minibatch crops)

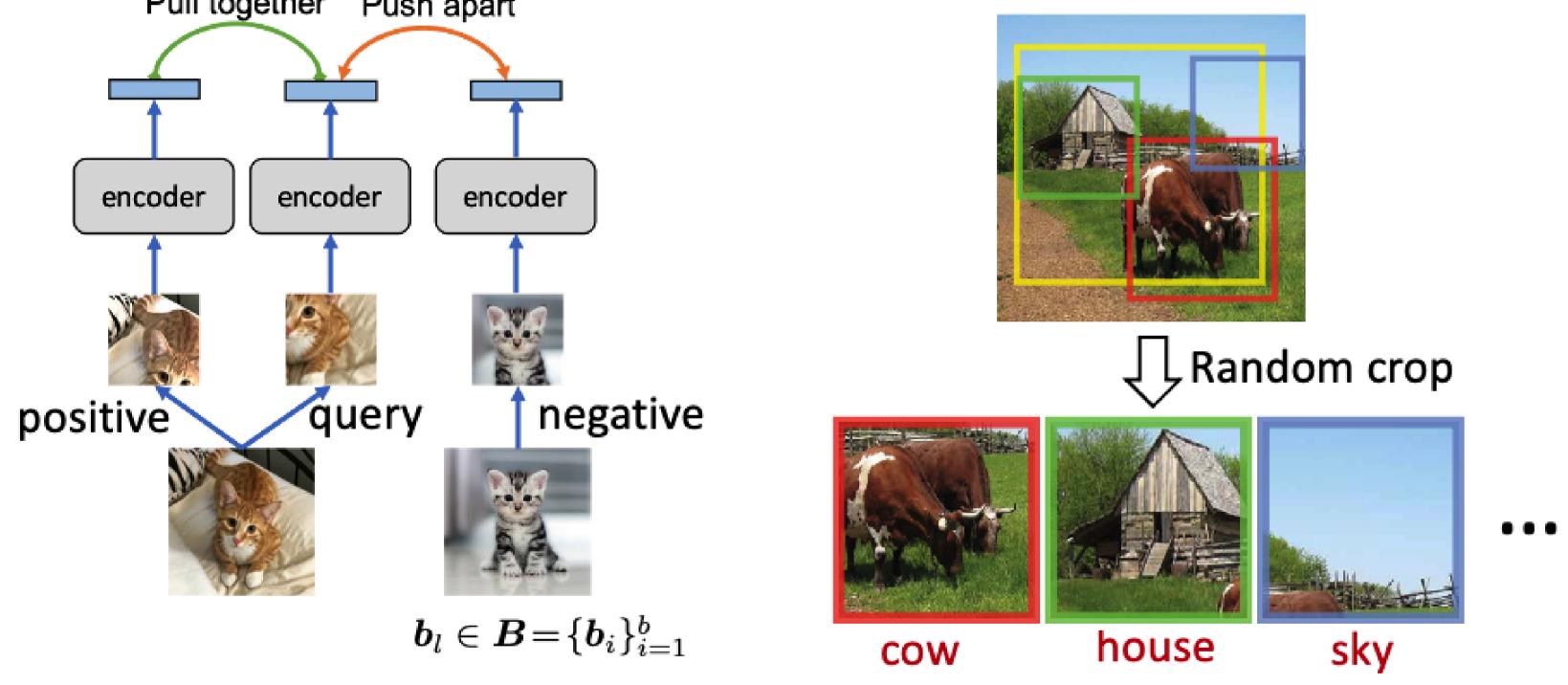
**One-hot label assignment**: query  $\mathbf{x}_i$  has only one positive  $\widetilde{\mathbf{x}}_i$  among  $\widetilde{\mathbf{x}}_i \cup \mathbf{B}$ 

$$\mathcal{L}_{n}(\boldsymbol{w}) = -\frac{1}{s} \sum_{i=1}^{s} \log \left( \frac{\sigma(\boldsymbol{x}_{i}, \widetilde{\boldsymbol{x}}_{i})}{\sigma(\boldsymbol{x}_{i}, \widetilde{\boldsymbol{x}}_{i}) + \sum_{l=1}^{b} \sigma(\boldsymbol{x}_{i}, \widetilde{\boldsymbol{x}}_{l})} \right)$$

## Issues of hot label assignment: imprecise & uninformative:

(1) some negatives & query from same semantic class Pull together Push apart

(2) augmentations gives different semantic crops



**Result:** one-hot label cannot guarantee semantically similar samples to close

### **Theoretical Motivation**

Support that the pair  $(\mathbf{x}_i, \widetilde{\mathbf{x}}_i)$  in the training dataset  $\mathcal{D} = \{(\mathbf{x}_i, \widetilde{\mathbf{x}}_i)\}_{i=1}^n$  sampled from an unknown distribution  $\mathcal{S}$  denotes the positive pair in MoCo. Assume the query  $\mathbf{x}_i$  has ground truth soft label  $\mathbf{y}_i^* \in \mathbb{R}^{b+1}$  over the key set  $\mathbf{B}_i = \{\widetilde{\mathbf{x}}_i \cup \mathbf{B}\}$ where  $y_{it}^*$  measures the semantic similarity between  $x_i$  and the t-th key  $b'_t$  in buffer  $B_i$ 

Theorem 1 (upper bound of generalization error, informal). Under proper assumptions, for MoCo, with probability  $1-\nu$  , the generalization error on instance discrimination task can be upper bounded as :

one-hot label true soft label  $\text{generalization error} \leq \mathcal{O}(\mathbb{E}_{\mathcal{D}\sim} s \left[ \| \boldsymbol{y} - \boldsymbol{y}^* \|_2 \right]) + \mathcal{O}(\sqrt{\frac{V_{\mathcal{D}} \ln(\boldsymbol{y})}{2}})$ 

training & test error gap

where  $V_{\mathcal{D}}$  is the variance of  $f_w$  on data  $\mathcal{D}$ ,  $\mathcal{F}$  is the covering number of encoder  $f_w$ 

**Remark:** the more accurate of the label y, the better the generalization

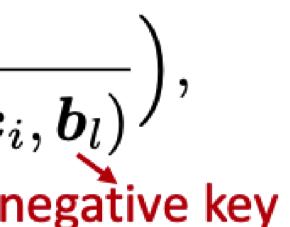
\* Salesforce Research

Pan Zhou\*, Caiming Xiong\*, Xiao-Tong Yuan<sup>†</sup>, Steven HOI\* <sup>†</sup> Nanjing University of Information Science & Technology {cxiong, shoi}@salesforce.com xtyuan@nuist.edu.cn

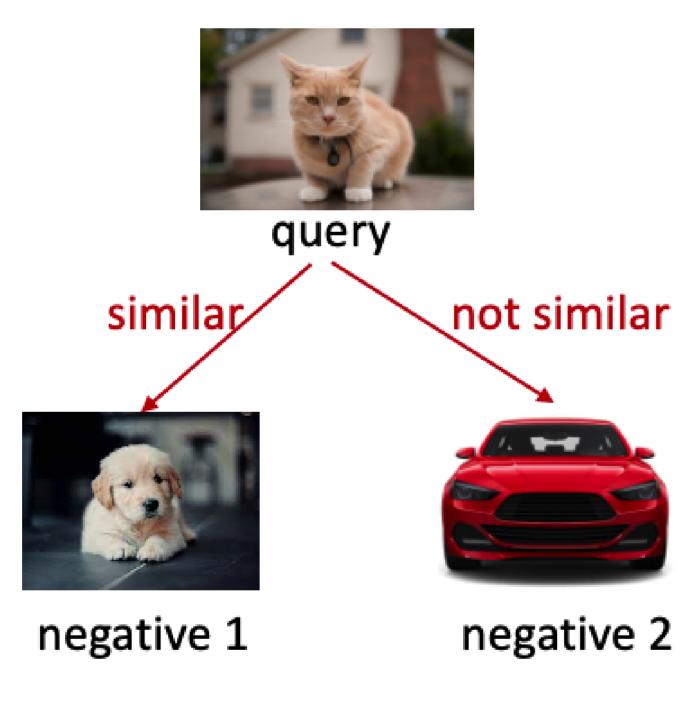
panzhou3@gmail.com

negative key

cosine similarity function



# (3) negatives have different similarity to query



$$\frac{\overline{\left(|\mathcal{F}|/\nu\right)}}{n} + \frac{\ln(|\mathcal{F}|/\nu)}{n}\Big),$$

# **Theoretical Motivation**

Theorem 2 (lower bound of generalization error, informal). Under proper assumptions, for MoCo, there exists a contrastive learning problem such that the generalization error on instance discrimination task is lower bounded as :

Solution: Self-Labeling Refinement

**Self-Labeling Refinement** has two components:

- accuracy
- accuracy

**MoCo Reformulation**: for query  $\mathbf{x}_i$  in minibatch  $\{(\mathbf{x}_i, \widetilde{\mathbf{x}}_i)\}_{i=1}^s$ , we maximize its similarity to its positive  $\mathbf{x}_i$  in key set  $\mathbf{\bar{B}} = \{\mathbf{\tilde{x}}_i\}_{i=1}^s \cup \{\mathbf{b}_i\}_{i=1}^b$  & push it away remaining samples:

$$\mathcal{L}_{c}(\boldsymbol{w}, \{(\boldsymbol{x}_{i}, \boldsymbol{y}_{i})\}) = -\frac{1}{s} \sum_{i=1}^{s} \sum_{k=1}^{s+b} \boldsymbol{y}_{ik} \log\left(\frac{\sigma(\boldsymbol{x}_{i}, \bar{\boldsymbol{b}}_{k})}{\sum_{l=1}^{s+b} \sigma(\boldsymbol{x}_{i}, \bar{\boldsymbol{b}}_{l})}\right)$$

dictionary, and thus can be linearly combined.

Step2. as  $\tilde{x}_i$  is highly similar to itself in  $\bar{B}$ ,  $p_{ii}^t$  is much larger than others, conceals similarity of other semantically similar instances in  $\bar{B}$ . So we remove  $\tilde{x}_i$  from  $\bar{B}$ 

$$\boldsymbol{q}_{ik}^{t} = \sigma^{1/\tau'}(\widetilde{\boldsymbol{x}}_{i}, I)$$

Step3. Linear combination

$$oldsymbol{ar{y}}_i^t = (1$$

Momentum mixup constructs virtual instance as follows:

 $\mathbf{x}'_i = \theta \mathbf{x}_i + (1 - \theta) \widetilde{\mathbf{x}}_k, \quad \mathbf{y}'_i = \theta \overline{\mathbf{y}}_i + (1 - \theta) \overline{\mathbf{y}}_k, \quad \theta \sim \text{Beta distribution}$ where  $\widetilde{\mathbf{x}}_k \overset{random}{\sim} {\{\widetilde{\mathbf{x}}_i\}_{i=1}^s, \overline{\mathbf{y}}_i \text{ is refined label by self-labeling refinery}}$ **Benefits**: the component  $\widetilde{\mathbf{x}}_k$  in  $\mathbf{x}'_i$  increases similarity between query  $\mathbf{x}'_i$  and positive key  $\mathbf{x}'_i$ , which improves the label accuracy



one-hot label true soft label

generalization error  $\geq \mathcal{O}(\mathbb{E}_{\mathcal{D}\sim} \boldsymbol{s}[\|\boldsymbol{y} - \boldsymbol{y}^*\|_2]).$ 

**Remark:** lower & upper bounds show generalization error  $\sim \mathbb{E}_{\mathcal{D}\sim \mathcal{S}}[||y - y^*||_2]$ the more accurate of the label y, the better the generalization

self-labeling refinery: soft label replaces one-hot label to directly improve label

momentum mixup: increase similarity of positive pair to indirectly improve label

where  $\bar{\bm{b}}_k$  is the k-th sample in  $\bar{\bm{B}}$ , the i-th entry  $\bm{y}_{ii}$  of one-hot label  $\bm{y}_i$  of query  $\bm{x}_i$  is 1. Target of Reformulation: labels of different samples are defined on a shared

Self-labeling refinery iteratively uses network to improve labels during training Step1. for query  $\mathbf{x}_i$ , we use its positive  $\widetilde{\mathbf{x}}_i$  to estimate semantic similarity between  $\mathbf{x}_i$ and instances in  $\overline{B} = {\{\widetilde{x}_i\}_{i=1}^s \cup \{b_i\}_{i=1}^b}$ , since  $x_i$  and  $\widetilde{x}_i$  come from the same image:  $\boldsymbol{p}_{ik}^{t} = \sigma^{1/\tau'}(\widetilde{\boldsymbol{x}}_{i}, \overline{\boldsymbol{b}}_{k}) / \sum_{l=1}^{s+b} \sigma^{1/\tau'}(\widetilde{\boldsymbol{x}}_{i}, \overline{\boldsymbol{b}}_{l}),$ 

 $(\bar{\boldsymbol{b}}_{k})/\sum_{l=1}^{s+b}\sigma^{1/\tau'}(\tilde{\boldsymbol{x}}_{l}, \bar{\boldsymbol{b}}_{l}), \ \boldsymbol{q}_{ll}^{t}=0.$ 

$$-\alpha_t - \beta_t \mathbf{y}_i + \alpha_t \mathbf{p}_i^t + \beta_t \mathbf{q}_i^t,$$





# Theoretical Analysis on Self-Labeling Refinary

- ground-truth semantic label  $y_i^* \in \{\gamma_t\}_{t=1}^K$  of  $x_i$  is decided by its corresponding  $c_t$
- the classes are separated:  $|\gamma_i \gamma_k| \ge \delta$ ,  $\|c_i c_k\|_2 \ge 2\varepsilon$ ,  $(\forall i \neq k)$ ,
- for each sample  $m{c}_i$  , at most  $ho n_i$  augmentations are assigned to wrong labels, where  $n_i$ denotes the crop sample number of  $c_i$

# **Remark:** along training, self-label refinery recovers corrupted training labels

that obeys  $\| m{x} - m{c}_k \|_2 \leq \varepsilon$ :

**Remark:** network trained by self-label refinery predicts label well





# Label-corrupted dataset

- Let  $\{(x_i, y_i^*)\}_{i=1}^n$  resp. denote the pairs of crops and ground-truth semantic label
- crop  $x_i$  generated from vanilla sample  $c_t$  obeys  $\|x_i c_t\|_2 \leq \varepsilon$
- Assume online/target networks  $f/g: x \in \mathbb{R}^d \mapsto f(W, x) = v^\top \phi(Wx)$ training loss:  $\mathcal{L}_t(W) = \frac{1}{2} \sum_{i=1}^n (\bar{y}_i^t - f(W, x_i))^2 = \frac{1}{2} ||\bar{y}_i^t - f(W, X)||_2^2$ gradient descent algorithm:  $W_{t+1} = W_t - \eta \nabla \mathcal{L}_t(W_t)$

# Theorem 3 (exact label recovery on training data, informal).

Under proper assumptions, after t iterations, for data  $\{x_i\}_{i=1}^n$ , we have  $\frac{1}{\sqrt{n}} \| \boldsymbol{y}^t - \boldsymbol{y}^* \|_2 \le \frac{1 - \alpha_t}{\sqrt{n}} \| \boldsymbol{y}^0 - \boldsymbol{y}^* \|_2 + \alpha_t (6\rho + \zeta),$ 

- where  $y^t = [y_1^t, \cdots, y_n^t]$ ,  $y^* = [y_1^*, \cdots, y_n^*]$ ,  $\zeta = c \varepsilon K^2 \sqrt{\log K}$  is related to network
- Moreover, if the following two assumption hold, •  $\rho \leq \frac{\delta}{24}$  : label corrupted ratio is small
- $(1 \alpha_0)|y_i^0 y_i^*| + \frac{1}{3}\alpha_0\delta < \frac{1}{2}\delta$ : label noise magnitude is small (  $\delta$  is class label separation) the estimated label  $y_i^t$  predicts true label  $y_i^*$  of any crop  $x_i$ 
  - **Exact label recovery:**  $\gamma_{k^*} = \boldsymbol{y}_i^*$  with  $k^* = \operatorname{argmin}_{1 \le k \le K} |\boldsymbol{y}_i^t \gamma_k|$ . estimated label true label
- Theorem 4 (exact prediction of network trained by self-label refinery, informal) Under proper assumptions, by using the refined label  $m{y}^t$  to train network, the error of network prediction on training data  $X = \{x_i\}_{i=1}^n$  is upper bounded as

$$\frac{1}{\sqrt{n}} \|f(\boldsymbol{W}_t, \boldsymbol{X}) - \boldsymbol{y}^*\|_2 \le 6\rho + \frac{c\zeta}{K\Gamma^2},$$
  
predicted label true label

- Moreover, if the following two assumption hold, •  $\rho \leq \frac{\delta}{24}$  : label corrupted ratio is small
- $(1 \alpha_0)|y_i^0 y_i^*| + \frac{1}{3}\alpha_0\delta < \frac{1}{2}\delta$ : label noise magnitude is small ( $\delta$  is class label separation) for any vanilla sample  $c_k$ , network  $f(W_t, \cdot)$  predicts true label  $y_i^*$  of any test augmentation x
  - **Exact label prediction:**  $\gamma_{k^*} = \gamma_k$  with  $k^* = \operatorname{argmin}_{1 \le i \le \overline{K}} |f(W_t, x) \gamma_i|$ . estimated label true label