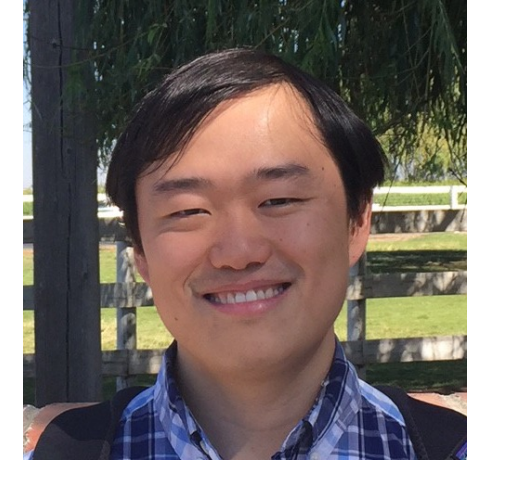


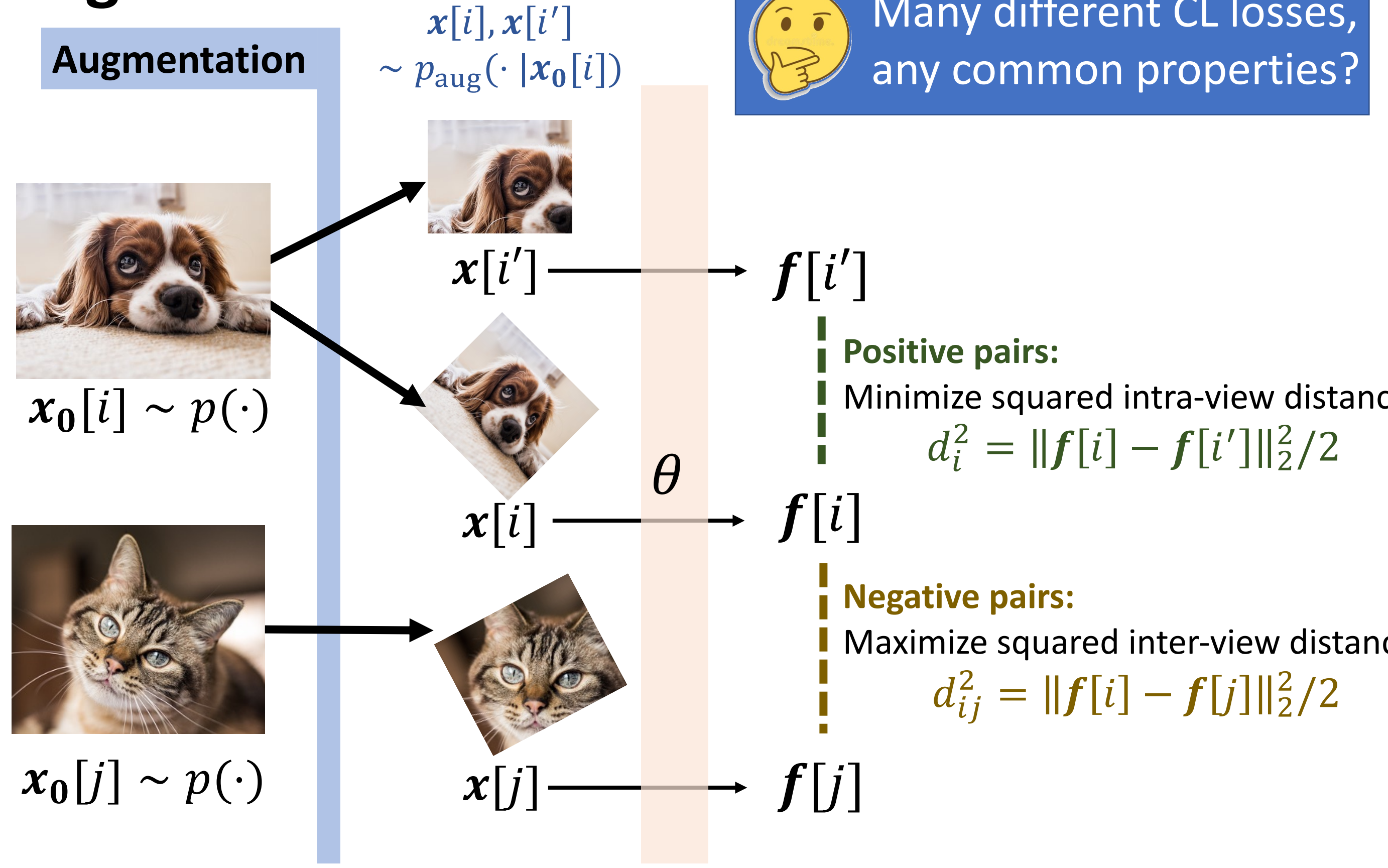
# Understanding Contrastive Learning via Coordinate-wise Optimization

Selected as Oral

Yuandong Tian  
yuandong@meta.com



## Background



## Common piece of various CL loss functions

First we can prove  $\frac{\partial \mathcal{L}_{\phi, \psi}}{\partial \theta} = -\frac{\partial \mathcal{E}_\alpha}{\partial \theta} |_{\alpha = \alpha(\theta)}$

for the energy  $\mathcal{E}_\alpha$  defined as the trace of **contrastive covariance**  $\mathbb{C}_\alpha$ :

$$\mathcal{E}_\alpha(\theta) := \frac{1}{2} \text{tr} \mathbb{C}_\alpha[f_\theta(x)]$$

where the **contrastive covariance** is defined as

$$\mathbb{C}_\alpha[f] := \sum_{i,j} \alpha_{ij} (f[i] - f[j])(f[i] - f[j])^T - (f[i] - f[i'])(f[i] - f[i'])^T$$

Here the **pairwise importance**  $\alpha_{ij} := \phi'(d_i) \psi'(d_{ij}) \geq 0$ , where  $\xi_i := \sum_{j \neq i} \psi(d_i^2 - d_{ij}^2)$

## $\alpha$ as an adversarial player

[Theorem] If  $\psi(x) = e^{x/\tau}$ , then  $\alpha(\theta) = \arg \min_{\alpha \in \mathcal{A}} \mathcal{E}_\alpha(\theta) - \mathcal{R}(\alpha)$

where  $\mathcal{A} := \{\alpha: \forall i, \sum_{j \neq i} \alpha_{ij} = \tau^{-1} \xi_i \phi'(\xi_i), \alpha_{ij} \geq 0\}$

and entropy regularization term  $\mathcal{R}(\alpha) := \tau \sum_{i=1}^N H(\alpha_i)$

Example For infoNCE:

$$\alpha_{ij}(\theta) = \frac{\exp(-d_{ij}^2/\tau)}{\epsilon \exp(-d_i^2/\tau) + \sum_{j \neq i} \exp(-d_{ij}^2/\tau)}$$

Larger  $\alpha_{ij}$  on **small**  $d_{ij} \rightarrow$  distinct samples with similar representations

[Theorem] Minimizing  $\mathcal{L}_{\phi, \psi} \Leftrightarrow$  Coordinate-wise optimization:

$$\alpha_t := \arg \min_{\alpha \in \mathcal{A}} \mathcal{E}_\alpha(\theta_t) - \mathcal{R}(\alpha)$$

$$\theta_{t+1} := \theta_t + \eta \nabla_{\theta} \mathcal{E}_{\alpha_t}(\theta_t)$$

### Max-player $\theta$

Learns the representation to maximize contrastiveness.

### Min-player $\alpha$

Find distinct sample pairs that share similar representation (i.e., **hard negative pairs**)

The pairwise importance  $\alpha$  incorporates the effects of  $\phi$  and  $\psi$ .

| Contrastive Loss                                  | $\phi(x)$                 | $\psi(x)$               |
|---|---------------------------|-------------------------|
| InfoNCE (Oord et al, 2018)                        | $\tau \log(\epsilon + x)$ | $e^{x/\tau}$            |
| MINE (Belghazi et al, 2018)                       | $\log(x)$                 | $e^x$                   |
| Triplet (Schroff et al., 2015)                    | $x$                       | $[x + \epsilon]_+$      |
| Soft Triplet (Tian et al., 2020c)                 | $\tau \log(1 + x)$        | $e^{x/\tau + \epsilon}$ |
| N+1 Tuplet (Sohn, 2016)                           | $\log(1 + x)$             | $e^x$                   |
| Lifted Structured (Oh Song et al., 2016)          | $[\log(x)]_+^2$           | $e^{x+\epsilon}$        |
| Modified Triplet (Eqn. 10 in Coria et al., 2020)) | $x$                       | $\text{sigmoid}(cx)$    |
| Triplet Contrastive (Eqn. 2 in Ji et al., 2021)   | Linear                    | Linear                  |

Different loss functions ( $\phi, \psi$ ) corresponds to the same energy function  $\mathcal{E}$   
**How the min player  $\alpha = \alpha(\theta)$  operates is different.**

## Proposed Unified Framework

General CL loss ( $\phi, \psi$  are monotonous increasing functions)

$$\min_{\theta} \mathcal{L}_{\phi, \psi}(\theta) := \sum_{i=1}^N \phi \left( \sum_{j \neq i} \psi(d_i^2 - d_{ij}^2) \right)$$

For infoNCE:

$$\mathcal{L}_{nce} := -\tau \sum_{i=1}^N \log \frac{e^{-d_i^2/\tau}}{\epsilon e^{-d_i^2/\tau} + \sum_{j \neq i} e^{-d_{ij}^2/\tau}} = \tau \sum_{i=1}^N \log \left( \epsilon + \sum_{j \neq i} \exp \left( \frac{d_i^2 - d_{ij}^2}{\tau} \right) \right)$$

Here  $\phi(x) = \tau \log(\epsilon + x)$  and  $\psi(x) = \exp(x/\tau)$

## Proposed: $\alpha$ -CL



Why we are stuck with coordinate-wise optimization?

Optimize network parameter  $\theta$  using gradient ascent of the energy function  $\mathcal{E}$ :

$$\theta_{t+1} = \theta_t + \eta \nabla_{\theta} \mathcal{E}_{\text{sg}(\alpha_t)}(\theta_t)$$

Pairwise importance  $\alpha_t = \alpha(\theta_t)$

- The pairwise importance  $\alpha$  can be
1. optimized by a separate loss function, or
  2. **directly** specified ( $\alpha$ -CL-direct)

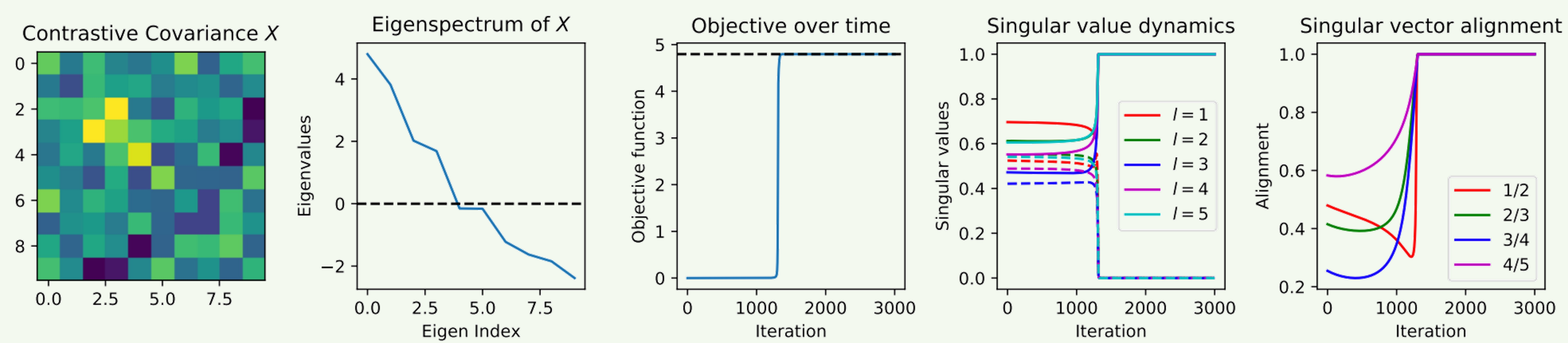
## Theoretical Properties when $\alpha$ is fixed

### Deep linear network

If  $f_\theta(x) = W_L W_{L-1} \dots W_1 x$ , then almost all local optima are global, and CL becomes Principal Component Analysis (PCA).

[Theorem] Let  $X_\alpha := \mathbb{C}_\alpha[x]$ . If  $\lambda_{\max}(X_\alpha) > 0$ , then for any local maximum  $\theta = \{W_L, W_{L-1}, \dots, W_1\}$  whose  $W_{>1}^T W_{>1}$  has distinct maximal eigenvalue, then

- $\theta$  is aligned rank-1 (i.e.,  $W_l = v_l v_{l-1}^T$ ),  $v_0$  is the unit eigenvector for  $\lambda_{\max}(X_\alpha)$ .
- $\theta$  is globally optimal with objective  $2\mathcal{E}^* = \lambda_{\max}(X_\alpha)$ .



### Nonlinear network

Many interesting properties. Detailed in the paper and follow-up works (Please check Workshop on SSL, Theory and Practice on Dec. 3)

## Experimental Results: $\alpha$ -CL

Use ResNet18 backbone, and set different  $\alpha$

- $\alpha$ -CL- $r_H$ : Entropy regularizer
- $\alpha$ -CL- $r_I$ : Inverse regularizer
- $\alpha$ -CL- $r_S$ : Square regularizer
- $\alpha$ -CL-direct: Directly setting  $\alpha$ .

|                           | CIFAR-10            |                     |                     | STL-10              |                     |                     |
|---------------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|
|                           | 100 epochs          | 300 epochs          | 500 epochs          | 100 epochs          | 300 epochs          | 500 epochs          |
| $\mathcal{L}_{quadratic}$ | 63.59 ± 2.53        | 73.02 ± 0.80        | 73.58 ± 0.82        | 55.59 ± 4.00        | 64.97 ± 1.45        | 67.28 ± 1.21        |
| $\mathcal{L}_{nce}$       | 84.06 ± 0.30        | 87.63 ± 0.13        | 87.86 ± 0.12        | 78.46 ± 0.24        | 82.49 ± 0.26        | 83.70 ± 0.12        |
| backprop $\alpha(\theta)$ | 83.42 ± 0.25        | 87.18 ± 0.19        | 87.48 ± 0.21        | 77.88 ± 0.17        | 81.86 ± 0.30        | 83.19 ± 0.16        |
| $\alpha$ -CL- $r_H$       | 84.27 ± 0.24        | <b>87.75 ± 0.25</b> | <b>87.92 ± 0.24</b> | <b>78.53 ± 0.35</b> | <b>82.62 ± 0.15</b> | <b>83.74 ± 0.18</b> |
| $\alpha$ -CL- $r_I$       | 83.72 ± 0.19        | 87.51 ± 0.11        | 87.69 ± 0.09        | 78.22 ± 0.28        | 82.19 ± 0.52        | 83.47 ± 0.34        |
| $\alpha$ -CL- $r_S$       | <b>84.72 ± 0.10</b> | 86.62 ± 0.17        | 86.74 ± 0.15        | 76.95 ± 1.06        | 80.64 ± 0.77        | 81.65 ± 0.59        |
| $\alpha$ -CL-direct       | <b>85.09 ± 0.13</b> | <b>88.00 ± 0.12</b> | <b>88.16 ± 0.12</b> | <b>79.38 ± 0.16</b> | <b>82.99 ± 0.15</b> | <b>84.06 ± 0.24</b> |

More datasets

|                     | CIFAR-100             |                       |                       |
|---------------------|-----------------------|-----------------------|-----------------------|
|                     | 100 epochs            | 300 epochs            | 500 epochs            |
| $\mathcal{L}_{nce}$ | 55.696 ± 0.368        | 59.706 ± 0.360        | 59.892 ± 0.340        |
| $\alpha$ -CL-direct | <b>57.144 ± 0.150</b> | <b>60.110 ± 0.187</b> | <b>60.330 ± 0.194</b> |

$\alpha$ -CL-direct:

$$\alpha_{ij} := \frac{\exp(-\frac{d_{ij}^p}{\tau})}{\sum_{i \neq j} \exp(-\frac{d_{ij}^p}{\tau})}$$

( $p = 4$ )

Backbone = ResNet50

| Dataset   | Method              | 100 epochs            | 300 epochs            | 500 epochs            |
|-----------|---------------------|-----------------------|-----------------------|-----------------------|
| CIFAR-10  | $\mathcal{L}_{nce}$ | 86.388 ± 0.157        | 89.974 ± 0.138        | 90.194 ± 0.232        |
|           | $\alpha$ -CL-direct | <b>87.406 ± 0.227</b> | <b>90.228 ± 0.185</b> | <b>90.366 ± 0.209</b> |
| CIFAR-100 | $\mathcal{L}_{nce}$ | 60.162 ± 0.482        | 65.400 ± 0.310        | 65.532 ± 0.297        |
|           | $\alpha$ -CL-direct | <b>62.650 ± 0.181</b> | <b>65.630 ± 0.263</b> | <b>65.636 ± 0.269</b> |
| STL-10    | $\mathcal{L}_{nce}$ | 81.635 ± 0.244        | 86.570 ± 0.174        | <b>87.900 ± 0.222</b> |
|           | $\alpha$ -CL-direct | <b>82.850 ± 0.171</b> | <b>86.870 ± 0.178</b> | 87.653 ± 0.175        |