# Understanding Deep Contrastive Learning via Coordinate-wise Optimization

Yuandong Tian

Research Scientist and Senior Manager
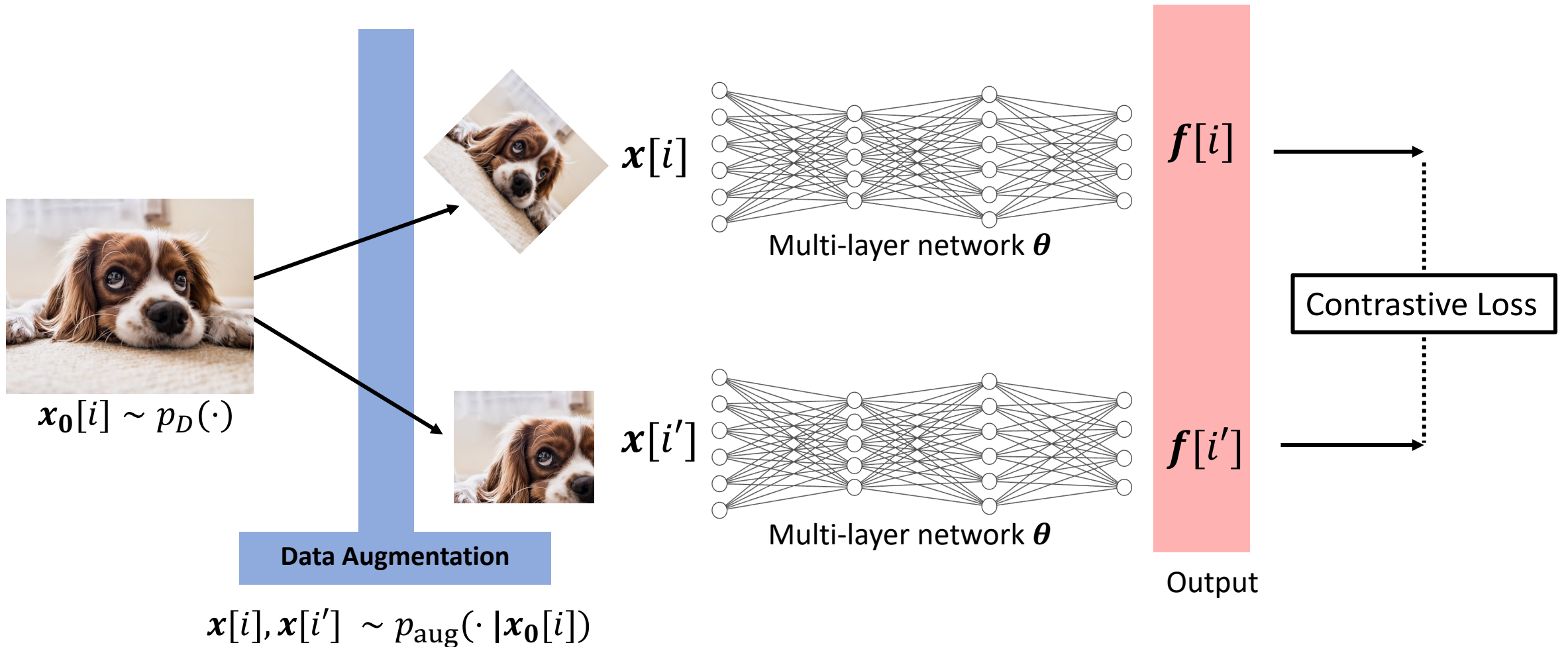
Meta AI (FAIR)
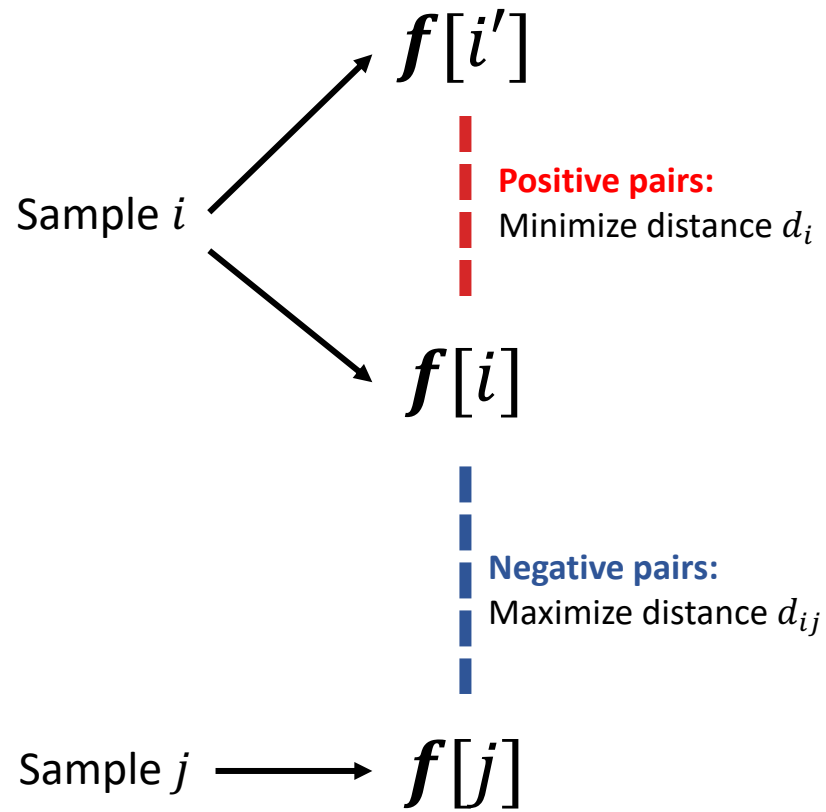
[Published in NeurIPS'22]

# Great Empirical Success of Deep Models

# Contrastive Learning (CL)



$x[i]$

Multi-layer network $\boldsymbol{\theta}$

$f[i]$

Contrastive Loss

$x_0[i] \sim p_D(\cdot)$

$x[i']$

Multi-layer network $\boldsymbol{\theta}$

$f[i']$

**Data Augmentation**

Output

$x[i], x[i'] \sim p_{\mathrm{aug}}(\cdot \,|\, x_0[i])$

# Formulation of Contrastive Learning

$$f[i']$$

Sample $i$

**Positive pairs:**
Minimize distance $d_i$

$$f[i]$$

**Negative pairs:**
Maximize distance $d_{ij}$

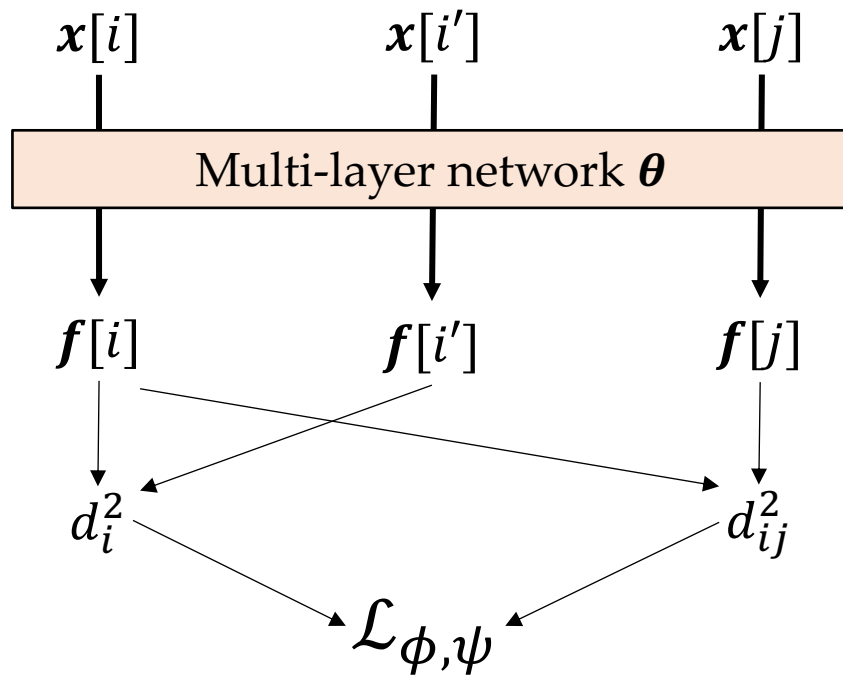Sample $j \longrightarrow$ $$f[j]$$

InfoNCE loss:

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

Intra-view distance $d_i^2 = \|f[i] - f[i']\|_2^2/2$

Inter-view distance $d_{ij}^2 = \|f[i] - f[j]\|_2^2/2$

# A family of contrastive losses

General Loss function we consider ($\phi, \psi$ are monotonous increasing functions)



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

Intra-view distance $d_i^2 = \|\boldsymbol{f}[i] - \boldsymbol{f}[i']\|_2^2/2$

Inter-view distance $d_{ij}^2 = \|\boldsymbol{f}[i] - \boldsymbol{f}[j]\|_2^2/2$

# A general family

| 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) | $[\log(x)]_+^2$ | $e^{x+\epsilon}$ |
| (Coria et al., 2020) | $x$ | $\mathrm{sigmoid}(cx)$ |
| (Ji et al., 2021) | linear | linear |

# Example: InfoNCE

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

$$= \tau \sum_{i=1}^{N} \log \left( \epsilon + \sum_{j \neq i} \exp\left( \frac{d_i^2 - d_{ij}^2}{\tau} \right) \right)$$

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

# Coordinate-wise Optimization

**Claim:** if $\psi(x) = e^{x/\tau}$, minimizing $\mathcal{L}_{\phi,\psi}$ $\Leftrightarrow$ Coordinate-wise optimization:

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

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

**Max-player $\theta$**

Learns the representation to maximize contrastiveness.

**Min-player $\alpha$**

Emphasize distinct sample pairs that share similar representation (**hard negative pairs**)

# Different Losses, Same Energy Function

| 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) | $[\log(x)]_+^2$ | $e^{x + \epsilon}$ |
| (Coria et al., 2020) | $x$ | $\mathrm{sigmoid}(cx)$ |
| (Ji et al., 2021) | linear | linear |

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

# How min player $\boldsymbol{\alpha}$ is determined?

If $\psi(x) = e^{x/\tau}$, then we have $\alpha(\boldsymbol{\theta}) := \arg\min_{\alpha \in \mathcal{A}} \mathcal{E}_\alpha(\boldsymbol{\theta}) - \mathcal{R}(\alpha)$

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

and entropy regularization term $\mathcal{R}(\alpha) := 2\tau \sum_{i=1}^N H(\alpha_i.)$ $\qquad \xi_i := \sum_{j \neq i} \psi(d_i^2 - d_{ij}^2)$

For infoNCE with $\epsilon = 0$, solving the optimization problem yields:

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

We put more weights on **small $\boldsymbol{d_{ij}}$**, i.e., distinct samples with similar representations

# Coordinate-wise Optimization

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

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

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

# Coordinate-wise Optimization

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

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

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

# Proposed: Pair-weighed CL ($\boldsymbol{\alpha}$-CL)

The min player $\alpha$ can be optimized by a loss function, or **_directly_** specified:

**Pairwise importance**

$$\alpha_t = \text{sg}(\alpha(\boldsymbol{\theta}_t))$$

$$\boldsymbol{\theta}_{t+1} := \boldsymbol{\theta}_t + \eta \nabla_{\boldsymbol{\theta}} \mathcal{E}_{\boxed{\alpha_t}}(\boldsymbol{\theta}_t)$$

# Experimental Results

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

- ($\alpha$-CL-$r_H$) Entropy regularizer $r_H(\alpha_{ij}) = -2\tau\alpha_{ij}\log\alpha_{ij}$;

- ($\alpha$-CL-$r_\gamma$) Inverse regularizers $r_\gamma(\alpha_{ij}) = \frac{2\tau}{1-\gamma}\alpha_{ij}^{1-\gamma}$ ($\gamma > 1$).

- ($\alpha$-CL-$r_s$) Square regularizer $r_s(\alpha_{ij}) = -\frac{\tau}{2}\alpha_{ij}^2$.   - ($\alpha$-CL-direct) Directly setting $\alpha$: $\alpha_{ij} = \exp(-d_{ij}^p/\tau)$ ($p > 1$).

# Experimental Results

More datasets

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

Backbone = ResNet50

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

# Roadmap of $\alpha$-CL

$$\mathcal{E}_\alpha(\boldsymbol{\theta}) := \operatorname{tr} \mathbb{C}_\alpha[\boldsymbol{f_\theta}(\boldsymbol{x})]$$

$\alpha$-CL

$$\min_{\boldsymbol{\theta}} \mathcal{L}_{\phi,\psi}(\boldsymbol{\theta})$$

Minimization of various CL losses

■ Applications ➡

Finding the best $\alpha = \alpha(\boldsymbol{\theta})$ for performance gain

Receptive-field specific $\alpha$

More applications (e.g., CL in GNN)

Understanding

Dynamics of $\boldsymbol{\theta}$ with fixed $\alpha$ in the linear setting

Dynamics of $\boldsymbol{\theta}$ in the nonlinear setting

Hierarchical representation learning

# Roadmap of $\alpha$-CL

$$\mathcal{E}_\alpha(\boldsymbol{\theta}) := \operatorname{tr} \mathbb{C}_\alpha[\boldsymbol{f}_{\boldsymbol{\theta}}(\boldsymbol{x})]$$

$\alpha$-CL

$$\min_{\boldsymbol{\theta}} \mathcal{L}_{\phi,\psi}(\boldsymbol{\theta})$$

Minimization of various CL losses

■ Applications ⇒

Finding the best $\alpha = \alpha(\boldsymbol{\theta})$ for performance gain

Receptive-field specific $\alpha$

More applications (e.g., CL in GNN)

Understanding

Dynamics of $\boldsymbol{\theta}$ with fixed $\alpha$ in the linear setting

Dynamics of $\boldsymbol{\theta}$ in the nonlinear setting

Hierarchical representation learning

# Deep linear case with fixed $\alpha$

If $f_\theta(x) = W(\theta)x$, then Contrastive Learning reduces to PCA objective

**Corollary 2** (Representation learning in Deep Linear CL reparameterizes Principal Component Analysis (PCA)). *When $z = W(\theta)x$ with a constraint $WW^\top = I$, $\mathcal{E}_\alpha$ is the objective of Principal Component Analysis (PCA) with reparameterization $W = W(\theta)$:*

$$\max_{\theta} \mathcal{E}_\alpha(\theta) = \mathrm{tr}(W(\theta)X_\alpha W^\top(\theta)) \quad \text{s.t. } WW^\top = I \tag{9}$$

*here $X_\alpha := \mathbb{C}_\alpha[x]$ is the contrastive covariance of input $x$.*

# Deep linear case with fixed $\alpha$

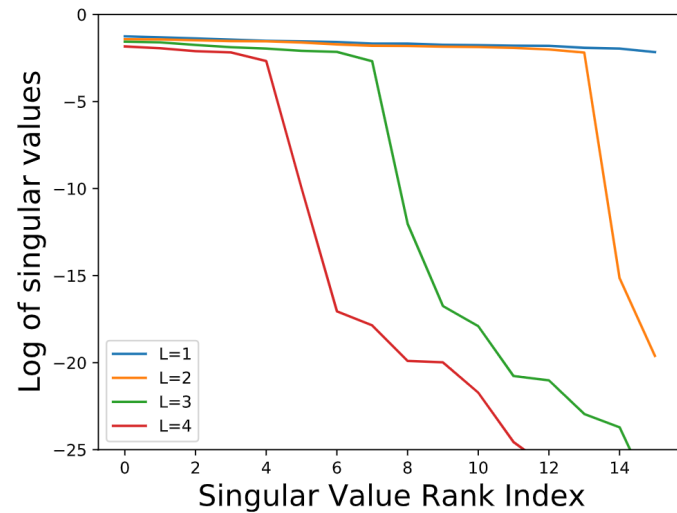If $f_\theta(x) = W_L W_{L-1} \dots W_1 x$, then almost all local optima are global and it is PCA

**Theorem 3** (Representation Learning with DeepLin is PCA). *If $\lambda_{\max}(X_\alpha) > 0$, then for any local maximum $\theta \in \Theta$ of Eqn. 11 whose $W_{>1}^\top W_{>1}$ has distinct maximal eigenvalue:*

- *there exists a set of unit vectors $\{v_l\}_{l=0}^L$ so that $\boxed{W_l = v_l v_{l-1}^\top}$ for $1 \le l \le L$, in particular, $v_0$ is the unit eigenvector corresponding to $\lambda_{\max}(X_\alpha)$,*

  <span style="color:red">1. Nearby weights align</span>

  <span style="color:red">2. All $W_l$ has rank-1 structure</span>

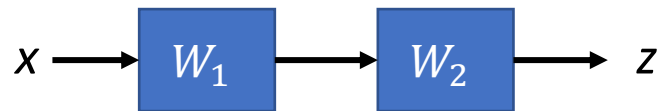- *$\theta$ is global optimal with objective $\mathcal{E}^* = \lambda_{\max}(X_\alpha)$.*

**Corollary 3.** *If we additionally use per-filter normalization (i.e., $\|w_{lk}\|_2 = 1/\sqrt{n_l}$), then Thm. 3 holds and $v_l$ is more constrained: $[v_l]_k = \pm 1/\sqrt{n_l}$ for $1 \le l \le L-1$.*

# Dimensional Collapsing in CL

Shouldn't contrastive SSL make full use of all dimensions? The answer is **No...**



(a) multiple layers

$$X \longrightarrow \boxed{W_1} \longrightarrow \boxed{W_2} \longrightarrow z$$

$W_1$ and $W_2$ will align with each other

If things are aligned, why not let them align directly?

| Loss function | Projector | Top-1 Accuracy |
|---|---|---|
| SimCLR | 2-layer nonlinear projector | 66.5 |
| SimCLR | 1-layer linear projector | 61.1 |
| SimCLR | no projector | 51.5 |
| *DirectCLR* | no projector | 62.7 |

**DirectCLR** [*L. Jing, P. Vincent, Y. LeCun, **Y. Tian**, Understanding Dimensional Collapse in Contrastive Self-supervised Learning, ICLR'22*]

# Roadmap of $\alpha$-CL

$$\mathcal{E}_\alpha(\boldsymbol{\theta}) := \operatorname{tr} \mathbb{C}_\alpha[f_{\boldsymbol{\theta}}(x)]$$

$\alpha$-CL

$$\min_{\boldsymbol{\theta}} \mathcal{L}_{\phi,\psi}(\boldsymbol{\theta})$$

Minimization of various CL losses

■ Applications ➡

Finding the best $\alpha = \alpha(\boldsymbol{\theta})$ for performance gain

Receptive-field specific $\alpha$

More applications (e.g., CL in GNN)

Understanding

Dynamics of $\boldsymbol{\theta}$ with fixed $\alpha$ in the linear setting

Dynamics of $\boldsymbol{\theta}$ in the nonlinear setting

**NeurIPS 2022 Workshop:
Self-Supervised Learning - Theory and Practice**

Hierarchical representation learning

# Thanks!