

# Theoretical Insights into Self-Supervised Contrastive Learning

**Weiran Huang**

Associate professor

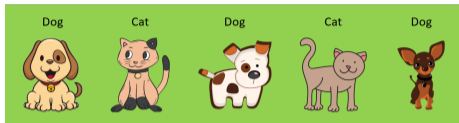
Qing Yuan Research Institute  
Shanghai Jiao Tong University



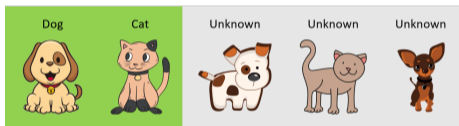
*December, 2023*

# Introduction to Self-Supervised Contrastive Learning

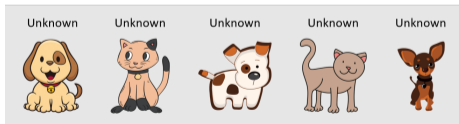
# Representation Learning Paradigm Evolution



Supervised Learning



Semi-Supervised Learning



Self-Supervised Learning



# Self-Supervised Learning

Self-Supervised Learning learns data representations through manually designed supervision signals, and then uses the learned representations for downstream tasks.

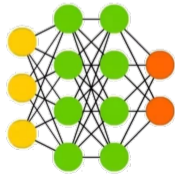
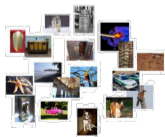
Training Dataset (no labels)



Self-Supervised Task

Pretext Task  
Contrastive Learning

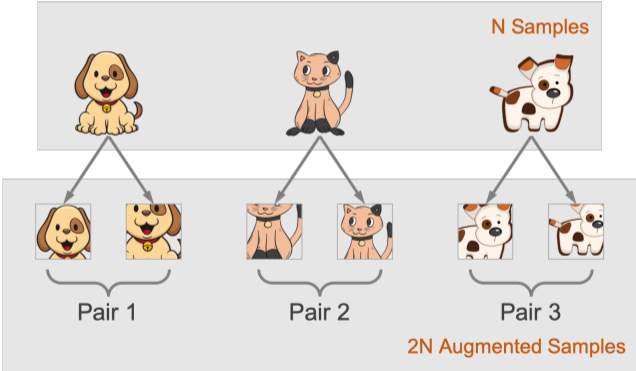
Test Dataset (with labels)



Downstream Task

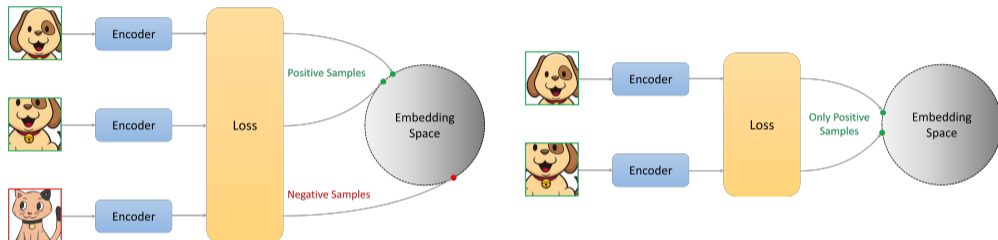
# Self-Supervised Contrastive Learning

Step 1 of 2: Construct similar sample pairs by data augmentation.



# Self-Supervised Contrastive Learning

Step 2 of 2: Pull the similar sample pairs close to each other in the embedding space.



The objectives of most contrastive learning algorithms (including SimCLR, MoCo, Barlow Twins, etc.) can be re-formulated as

$$\min \mathcal{L}(f) = \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in A(\mathbf{x})} \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|^2 + \mathcal{L}_{\text{regularization}}(f).$$

# Observations in Self-Supervised Contrastive Learning

1. Aligning positive samples (augmented from the “same data point”) is able to gather the samples from the “same latent class” into a cluster.

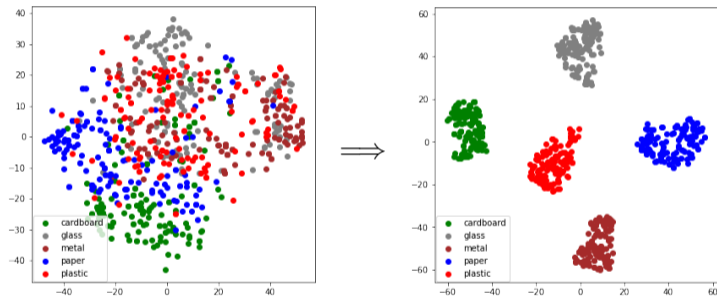
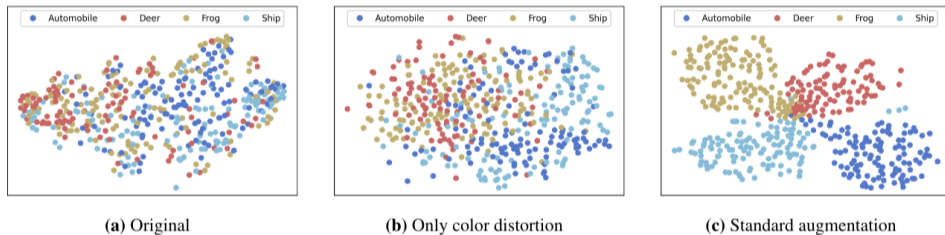


Figure: Embedding Space  
(<https://github.com/mwdhont/SimCLRv1-keras-tensorflow>).

# Observations in Self-Supervised Contrastive Learning

2. Richer data augmentation leads to a more clustered structure in the embedding space.



**Figure:** SimCLR's embedding space with different richnesses of data augmentations.



# Observations in Self-Supervised Contrastive Learning

3. The best composition of augmentations: random cropping and random color distortion.

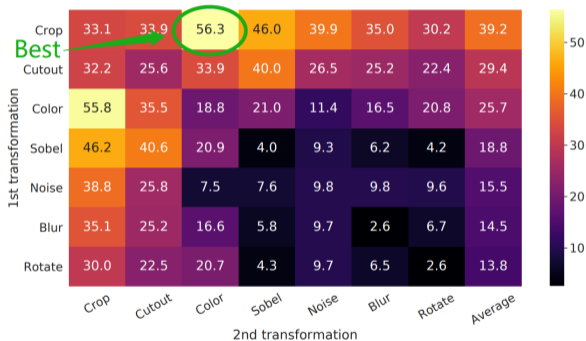


Figure: Experimental results reported in SimCLR paper.

# Observations in Self-Supervised Contrastive Learning

4. Barlow Twins decorrelates components of representation instead of directly optimizing the geometry of embedding space, but it still results in the clustered structure.

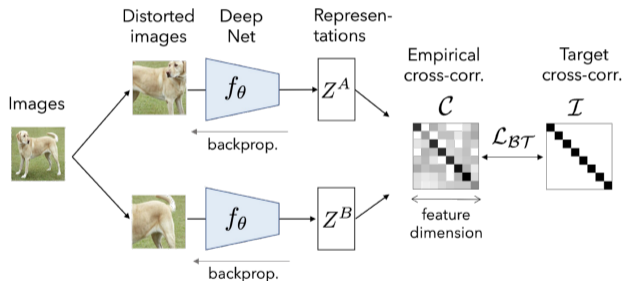


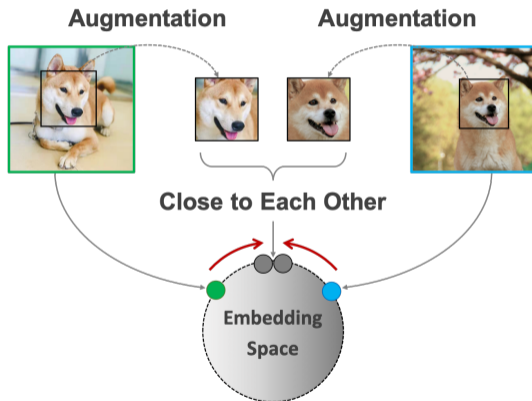
Figure: Barlow Twins aims to decorrelate the components of representation.

# Theoretical Analysis of Self-Supervised Contrastive Learning

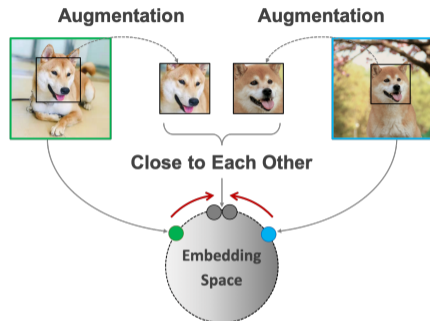
“Towards the Generalization of Contrastive Self-Supervised Learning.”  
**Huang**<sup>\*,†</sup>, Yi<sup>\*</sup>, Zhao<sup>\*</sup>, Jiang. [ICLR 2023](#).

# Intuition

Why does contrastive learning work?



# Intuition



For a given data augmentation set  $A$ , we define the **augmented distance** between two different samples as

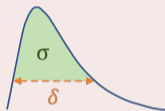
$$d_A(\mathbf{x}_1, \mathbf{x}_2) = \min_{\mathbf{x}'_1 \in A(\mathbf{x}_1), \mathbf{x}'_2 \in A(\mathbf{x}_2)} \|\mathbf{x}'_1 - \mathbf{x}'_2\|.$$

# Data Augmentation Modeling

## Definition 1 (( $\sigma, \delta$ )-Augmentation)

The data augmentation set  $A$  is called a  $(\sigma, \delta)$ -augmentation, if for each class  $C_k$ , there exists a subset  $C_k^0 \subseteq C_k$  (called the main part of  $C_k$ ) such that

- $\mathbb{P}[\mathbf{x} \in C_k^0] \geq \sigma \mathbb{P}[\mathbf{x} \in C_k]$  where  $\sigma \in (0, 1]$ ,
- $\sup_{\mathbf{x}_1, \mathbf{x}_2 \in C_k^0} d_A(\mathbf{x}_1, \mathbf{x}_2) \leq \delta$ .



The sharpness of concentration:

- Larger  $\sigma$  and smaller  $\delta$  indicate the sharper concentration of augmented data.
- Richer data augmentation leads to sharper concentration.

# Performance Guarantee of Self-Supervised Contrastive Learning

## Theorem 1

*Under mild assumptions, if the augmentation used in contrastive learning is  $(\sigma, \delta)$ -augmented, and*

$$\mu_k^\top \mu_\ell < r^2 \left( 1 - \rho_{\max}(\sigma, \delta, \varepsilon) - \sqrt{2\rho_{\max}(\sigma, \delta, \varepsilon)} - \frac{\Delta_\mu}{2} \right)$$

*holds for any pair of  $(\ell, k)$  with  $\ell \neq k$ , then the error rate of downstream classification*

$$\text{Err}(G_f) \leq (1 - \sigma) + R_\varepsilon,$$

*where  $\rho_{\max}(\sigma, \delta, \varepsilon) = 2(1 - \sigma) + \frac{R_\varepsilon}{\min_\ell \rho_\ell} + \sigma \left( \frac{L\delta}{r} + \frac{2\varepsilon}{r} \right)$  and  $\Delta_\mu = 1 - \min_{k \in [K]} \frac{\|\mu_k\|^2}{r^2}$ .*

## Messages From Theorem 1

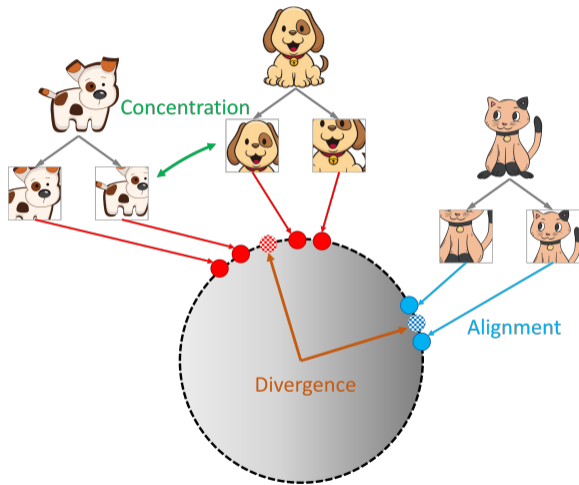
The generalization ability depends on three key factors:

- ① (**Alignment of positive samples**) How close positive samples are located to each other in the embedding space;
- ② (**Divergence of class centers**) How far apart class centers are located from each other in the embedding space;
- ③ (**Concentration of augmented data**) How sharp the concentration of augmented data is.

Only the first two factors can be optimized during the learning process. In contrast, the third factor is priorly decided by the pre-defined data augmentation and is independent of the learning process.



# Messages From Theorem 1



## Loss Functions

- InfoNCE (e.g., SimCLR): pull close positive pairs and push away negative pairs.

$$\mathcal{L}_{\text{InfoNCE}} = - \mathbb{E}_{\substack{\mathbf{x}, \mathbf{x}' \\ \mathbf{x}_1, \mathbf{x}_2 \in A(\mathbf{x}) \\ \mathbf{x}^- \in A(\mathbf{x}')}} \log \frac{e^{f(\mathbf{x}_1)^\top f(\mathbf{x}_2)}}{e^{f(\mathbf{x}_1)^\top f(\mathbf{x}_2)} + e^{f(\mathbf{x}_1)^\top f(\mathbf{x}^-)}},$$

where  $\mathbf{x}, \mathbf{x}'$  are two random samples and  $A$  is the data augmentation set.

- Cross-Correlation (e.g., Barlow Twins): decorrelate feature components.

$$\mathcal{L}_{\text{Cross-Corr}} = \sum_{i=1}^d (1 - C_{ii})^2 + \lambda \sum_{i=1}^d \sum_{i \neq j} C_{ij}^2, \quad \left( \mathbb{E} \left[ f(\mathbf{x}_1) f(\mathbf{x}_2)^\top \right] \rightarrow I_{d \times d} \right)$$

where  $C_{ij} = \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in A(\mathbf{x})} [f_i(\mathbf{x}_1) f_j(\mathbf{x}_2)]$ ,  $d$  is the dimension of encoder  $f$ , and  $f$  is normalized as  $\mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}' \in A(\mathbf{x})} [f_i(\mathbf{x}')^2] = 1$  for each dimension.

# Loss Functions

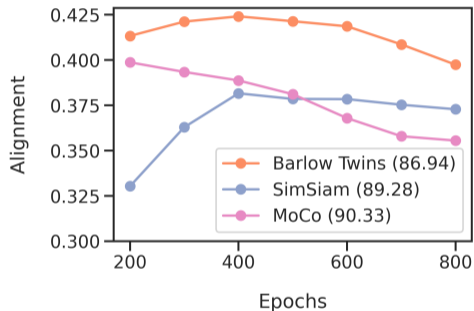
The above two losses can be split into two parts:

$$\mathcal{L}(f) = \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in A(\mathbf{x})} \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|^2 + \mathcal{L}_{\text{regularization}}(f).$$

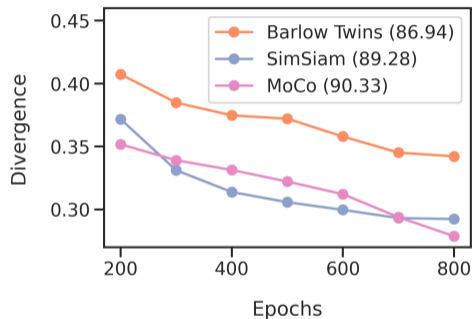
- For InfoNCE, we prove that  $\mu_k^\top \mu_\ell \lesssim \mathcal{L}_{\text{regularization}}(f)$ ;
- For Cross-Correlation, we prove that  $\mu_k^\top \mu_\ell \lesssim \sqrt{\mathcal{L}_{\text{regularization}}(f)}$ .

Therefore, minimizing these two losses can achieve **good alignment** and **large divergence**.

# How Alignment and Divergence Change During Training Process



(a) Alignment change of different algorithms



(b) Divergence change of different algorithms

## A Cookbook of Self-Supervised Learning

Randall Balestriero\*, Mark Ibrahim\*, Vlad Sobal\*, Ari Morcos\*, Shashank Shekhar\*, Tom Goldstein<sup>1</sup>, Florian Bordes<sup>2</sup>, Adrien Bardes\*, Gregoire Mialon\*, Yuandong Tian<sup>3</sup>, Avi Schwarzschild<sup>1</sup>, Andrew Gordon Wilson<sup>4</sup>, Jonas Geiping<sup>5</sup>, Quentin Garrido<sup>3</sup>, Pierre Fernandez<sup>6</sup>, Amir Bar<sup>7</sup>, Hamed Pirsiavash<sup>4</sup>, Yann LeCun<sup>8</sup> and Micah Goldblum<sup>8</sup>

\*Meta AI, FAIR

<sup>2</sup>New York University

<sup>1</sup>University of Maryland

<sup>4</sup>University of California, Davis

<sup>1</sup>Universite de Montreal, Mila

<sup>5</sup>Univ Gustave Eiffel, CNRS, LIGM

<sup>6</sup>Univ. Rennes, Inria, CNRS, IRISA

<sup>8</sup>Equal contributions, randomized ordering

## Contents

<b>1 What is Self-Supervised Learning and Why Bother?</b>	<b>3</b>
1.1 Why a Cookbook for Self-Supervised Learning?	3
<b>2 The Families and Origins of SSL</b>	<b>4</b>
2.1 Origins of SSL	5
2.2 The Deep Metric Learning Family: SimCLR/NNCLR/MeanSHIFT/SCL	7
2.3 The Self-Distillation Family: BYOL/SimSIAM/DINO	8
2.4 The Canonical Correlation Analysis Family: VICReg/BarlowTwins/SWAV/W-MSE	13
2.5 Masked Image Modeling	14
2.6 A Theoretical Unification Of Self-Supervised Learning	16
2.6.1 Theoretical Study of SSL	16
2.6.2 Dimensional Collapse of Representations	18
2.7 Pretraining Data	19

### 2.6.1 Theoretical Study of SSL

Numerous works have attempted to unify various SSL methods. In Huang et al. [2021], Barlow Twins' criterion is shown to be linked to an upper bound of a contrastive loss. This suggests a link exists between contrastive and covariance-based methods. This direction was further pursued in Garrido et al. [2022b], where a covariance-based and contrastive criterion are shown to be equivalent up to normalization by deriving the precise gap between the two approaches. These results were further validated empirically as methods were shown to exhibit similar performance and representation properties at ImageNet's scale (1.2 million samples). The similarities among methods was also studied in Tao et al. [2021] where this unification was tackled from a study of the losses' gradients.

# The Effect of Concentration Factor

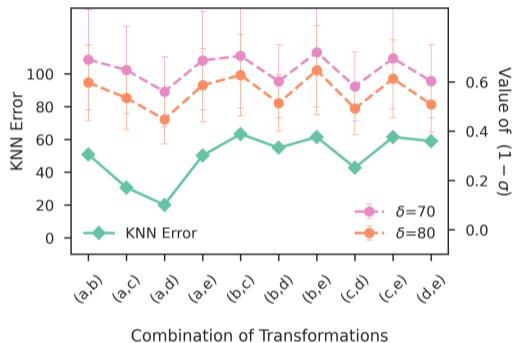
Dataset	Transformations					Accuracy			
	(a)	(b)	(c)	(d)	(e)	SimCLR	Barlow Twins	MoCo	SimSiam
CIFAR-10	✓	✓	✓	✓	✓	<b>89.76 ± 0.12</b>	<b>86.91 ± 0.09</b>	<b>90.12 ± 0.12</b>	<b>90.59 ± 0.11</b>
	✓	✓	✓	✓		88.48 ± 0.22	85.38 ± 0.37	89.69 ± 0.11	89.34 ± 0.09
	✓	✓	✓			83.50 ± 0.14	82.00 ± 0.59	86.78 ± 0.07	85.38 ± 0.09
	✓	✓				63.23 ± 0.05	67.83 ± 0.94	75.12 ± 0.28	63.27 ± 0.30
	✓					62.74 ± 0.18	67.77 ± 0.69	74.94 ± 0.22	61.47 ± 0.74
CIFAR-100	✓	✓	✓	✓	✓	<b>57.74 ± 0.12</b>	<b>57.99 ± 0.29</b>	<b>64.19 ± 0.14</b>	<b>63.48 ± 0.16</b>
	✓	✓	✓	✓		55.43 ± 0.10	55.22 ± 0.25	62.50 ± 0.28	60.31 ± 0.41
	✓	✓	✓			45.10 ± 0.25	50.40 ± 0.64	57.04 ± 0.21	51.42 ± 0.14
	✓	✓				28.01 ± 0.18	34.11 ± 0.59	40.18 ± 0.04	26.26 ± 0.30
	✓					27.95 ± 0.09	34.05 ± 1.13	39.63 ± 0.31	25.90 ± 0.83

(a) random cropping; (b) random Gaussian blur;  
(c) color dropping; (d) color distortion;  
(e) random horizontal flipping.

## The Effect of Concentration Factor

Dataset	Color Distortion Strength	Accuracy			
		SimCLR	Barlow Twins	MoCo	SimSiam
CIFAR-10	1	<b>82.75 ± 0.24</b>	<b>82.58 ± 0.25</b>	<b>86.68 ± 0.05</b>	<b>82.50 ± 1.05</b>
	1/2	78.76 ± 0.18	81.88 ± 0.25	84.30 ± 0.14	81.80 ± 0.15
	1/4	76.37 ± 0.11	79.64 ± 0.34	82.76 ± 0.09	78.80 ± 0.17
	1/8	74.23 ± 0.16	77.96 ± 0.16	81.20 ± 0.12	76.09 ± 0.50
CIFAR-100	1	<b>46.67 ± 0.42</b>	<b>50.39 ± 1.09</b>	<b>58.50 ± 0.51</b>	<b>49.94 ± 2.01</b>
	1/2	40.21 ± 0.05	48.76 ± 0.25	55.08 ± 0.09	46.27 ± 0.46
	1/4	36.67 ± 0.08	46.22 ± 0.71	52.09 ± 0.18	42.02 ± 0.34
	1/8	34.75 ± 0.20	44.72 ± 0.26	49.43 ± 0.16	36.26 ± 0.34

# The Effect of Concentration Factor



- (a) random cropping;
- (b) random Gaussian blur;
- (c) color dropping;
- (d) color distortion;
- (e) random horizontal flipping.

- Fix one transformation as (a), we observe that  $(a, d) < (a, c) < (a, e) \approx (a, b)$ ;
- Composition (a, d) has the **sharpest concentration** and **best performance**.



## Short Summary

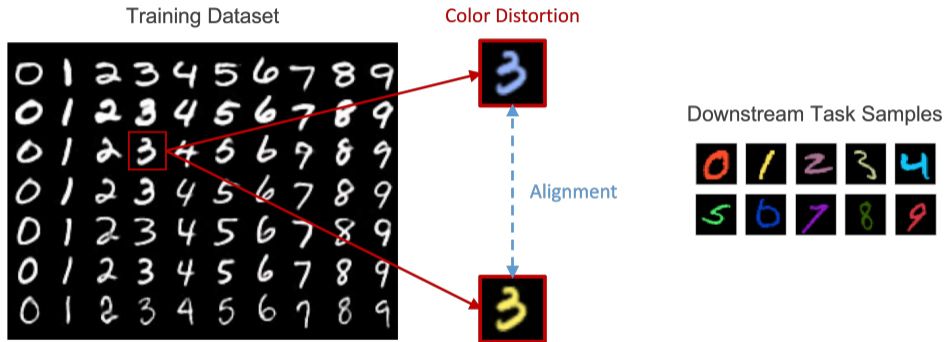
- We provide a mathematical formulation to model the data augmentation.
- We show that alignment of positive samples, divergence of class centers and concentration of augmented data are three key factors of self-supervised contrastive learning.
- We prove that SimCLR and Barlow Twins implicitly optimize the first two factors.
- We empirically verify that sharper concentration results in better generalization.

PS: Can Masked Auto-Encoder (MAE) be analyzed by the proposed framework?

# Transferability of Self-Supervised Contrastive Learning

“ArCL: Enhancing Contrastive Learning with Augmentation-Robust Representations.”  
Zhao<sup>\*</sup>, Du<sup>\*</sup>, Wang, Yao, **Huang**<sup>†</sup>. [ICLR 2023](#).

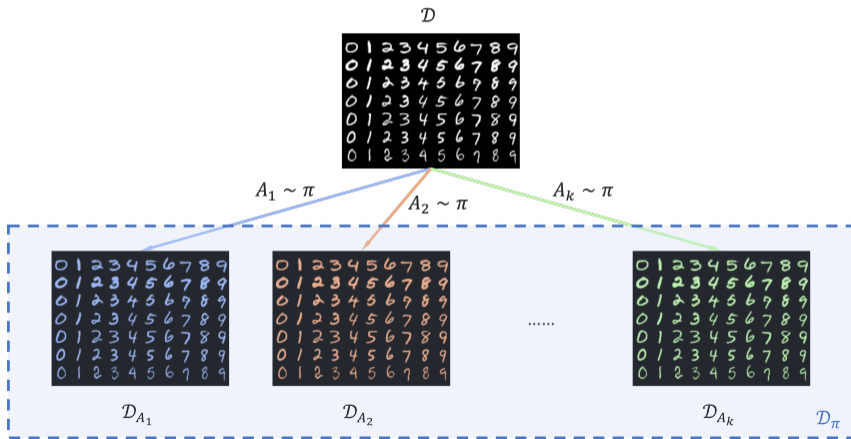
# Paradox



**Question:** Can contrastive learning extract augmentation-invariant features?

## Notations for Distributions

Let  $\mathcal{D}$  denote the original upstream distribution, and  $\mathcal{D}_A$  denote the augmented distribution after applying some random transformation  $A$  which follows distribution  $\pi$ .



## Transferability Evaluation

We first define “risk” of representation  $f$  over downstream distribution  $\mathcal{D}^{\text{tar}}$  as

$$\mathcal{R}(f; \mathcal{D}^{\text{tar}}) := \min_h \mathbb{E}_{(X, Y) \sim \mathcal{D}^{\text{tar}}} \ell(h \cdot f(X), Y),$$

where  $h$  is a linear classifier and  $\ell$  is the loss function.

For an augmentation-invariant representation, the risk should be at about the same level on two different downstream datasets  $\mathcal{D}_{A_1}, \mathcal{D}_{A_2}$  augmented from  $\mathcal{D}$ .

In other words,  $|\mathcal{R}(f; \mathcal{D}_{A_1}) - \mathcal{R}(f; \mathcal{D}_{A_2})|$  should be small.

## A Counter-Example

### Example 1

Data distribution  $(X_1, X_2) \sim \mathcal{N}(0, I)$ ,  $Y = \mathbb{1}(X_1 \geq 0)$ .

Augmentation distribution  $A_\theta(X_1, X_2) = (X_1, \theta \cdot X_2)$  where  $\theta \sim \mathcal{N}(0, 1)$ .

In this case,  $X_1$  is the augmentation-invariant feature.

In fact, we can prove that

$$\forall \varepsilon > 0, \exists f, \mathcal{D}_{A_1}, \mathcal{D}_{A_2}, \text{ s.t. } \mathcal{L}_{\text{align}}(f; \mathcal{D}) < \varepsilon \text{ and } |\mathcal{R}(f; \mathcal{D}_{A_1}) - \mathcal{R}(f; \mathcal{D}_{A_2})| > \text{const.}$$

## Proof (1/2)

For any  $\varepsilon > 0$ , let  $f(x_1, x_2) = x_1 + \frac{\sqrt{\varepsilon}}{2} \cdot x_2$ .

$$\begin{aligned}\mathcal{L}_{\text{align}}(f; \mathcal{D}, \pi) &= \mathbb{E}_{X \sim \mathcal{D}} \mathbb{E}_{(A_1, A_2) \sim \pi^2} \|f(A_1(X)) - f(A_2(X))\|^2 \\ &= \mathbb{E}_{(X_1, X_2) \sim \mathcal{N}(0, I)} \mathbb{E}_{(\theta_1, \theta_2) \sim \mathcal{N}(0, I)} \left| \left( X_1 + \frac{\sqrt{\varepsilon}}{2} \theta_1 X_2 \right) - \left( X_1 + \frac{\sqrt{\varepsilon}}{2} \theta_2 X_2 \right) \right|^2 \\ &= \frac{\varepsilon}{4} \mathbb{E}_{X_2 \sim \mathcal{N}(0, 1)} X_2^2 \mathbb{E}_{(\theta_1, \theta_2) \sim \mathcal{N}(0, I)} (\theta_1 - \theta_2)^2 = 2 \left( \frac{\sqrt{\varepsilon}}{2} \right)^2 < \varepsilon.\end{aligned}$$

## Proof (2/2)

Let  $c = 0$  and  $c' = 2/\sqrt{\varepsilon}$ . Then we have two domains

$$\mathcal{D}_c = \{(X_1, 0) : X_1 \sim \mathcal{N}(0, 1)\}$$

$$\mathcal{D}_{c'} = \{(X_1, 2X_2/\sqrt{\varepsilon}) : X_1 \sim \mathcal{N}(0, 1), X_2 \sim \mathcal{N}(0, 1)\}$$

Therefore, we can get  $\mathcal{R}(f; \mathcal{D}_c) = 0$ , but

$$\begin{aligned}\mathcal{R}(f; \mathcal{D}_{c'}) &= P(Y = 0, hf(X) \geq 0) + P(Y = 1, hf(X) < 0) \\ &\quad \text{(suppose } h \in \mathbb{R}^+ \text{ without loss of generality)} \\ &= P(X_1 < 0, f(X) \geq 0) + P(X_1 \geq 0, f(X) < 0) \\ &= P(X_1 < 0, X_1 + X_2 \geq 0) + P(X_1 \geq 0, X_1 + X_2 \leq 0) \\ &= \frac{1}{8} + \frac{1}{8} = \frac{1}{4}.\end{aligned}$$



## Augmentation-Robust Loss

We define the Augmentation-Robust loss as

$$\mathcal{L}_{AR}(f; \mathcal{D}) := \mathbb{E}_{X \in \mathcal{D}} \sup_{A_1, A_2} \|f(A_1(X)) - f(A_2(X))\|^2 \geq \mathcal{L}_{align}(f; \mathcal{D}).$$

### Theorem 2

For any  $A$ , let  $h_A \in \arg \min_h \mathcal{R}(h \circ f, \mathcal{D}_A)$ , we have

$$0 \leq \mathcal{R}(h_{A'} \circ f; \mathcal{D}_A) - \mathcal{R}(h_A \circ f; \mathcal{D}_A) \leq c \cdot (\|h_A\| + \|h_{A'}\|) \mathcal{L}_{AR}(f, \mathcal{D}).$$

Note that for any augmentation-invariant feature  $f$ ,  $\mathcal{R}(h_{A'} \circ f; \mathcal{D}_A) - \mathcal{R}(h_A \circ f; \mathcal{D}_A) = 0$ .

The empirical version of AR loss is

$$\hat{\mathcal{L}}_{AR}(f) := \frac{1}{n} \sum_{k=1}^n \max_{A_i, A_j} \|f(A_i(X_k)) - f(A_j(X_k))\|^2.$$

# Plug-And-Play ArCL

---

**Algorithm 1:** SimCLR + ArCL

---

**input** : Batch size  $N$ , temperature  $\tau$ , augmentation  $\pi$ , number of views  $m$ , epoch  $T$ ,  
encoder  $f$ , projector  $g$ .

- 1 **for**  $t = 1, \dots, T$  **do**
- 2     sample minibatch  $\{X_i\}_{i=1}^N$ ;
- 3     **for**  $i = 1, \dots, N$  **do**
- 4         draw  $m$  augmentations  $\hat{A} = \{A_1, \dots, A_m\} \sim \pi$ ;
- 5          $z_{i,j} = g(f(A_j X_i))$  for  $j \in [m]$ ;
- 6         *# select the worst positive samples;*
- 7          $s_i^+ = \min_{j,k \in [m]} \{z_{i,j}^\top z_{i,k} / (\|z_{i,j}\| \|z_{i,k}\|)\}$ ;
- 8         *# select the negative samples;*
- 9         **for**  $j = 1, \dots, N$  **do**
- 10              $s_{i,j}^- = z_{i,1}^\top z_{j,1} / (\|z_{i,1}\| \|z_{j,1}\|)$ ;
- 11              $s_{i,j+N}^- = z_{i,1}^\top z_{j,2} / (\|z_{i,1}\| \|z_{j,2}\|)$ ;
- 12     compute  $L = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(s_i^+ / \tau)}{\sum_{j=1, j \neq i}^{2N} \exp(s_{i,j}^- / \tau)}$ ;
- 13     update  $f$  and  $g$  to minimize  $L$ ;
- 14 **return**  $f$

---

## Performance on Augmented Datasets

	Method	Batch Size	Aug 1	Aug 2	Aug 3	Aug 4	Aug 5	Original
CIFAR10	SimCLR	256	86.36	83.21	86.93	86.42	86.13	86.76
	SimCLR + ArCL (views=4)	256	88.68	86.77	89.01	88.70	88.31	88.95
	SimCLR + ArCL (views=6)	256	<b>88.95</b>	<b>87.18</b>	<b>89.54</b>	<b>88.92</b>	<b>88.61</b>	<b>89.11</b>
	SimCLR	512	88.62	86.27	88.96	88.56	88.37	88.81
	SimCLR + ArCL (views=4)	512	89.97	88.06	90.48	89.91	89.59	90.20
	SimCLR + ArCL (views=6)	512	<b>90.24</b>	<b>89.54</b>	<b>90.69</b>	<b>90.43</b>	<b>90.07</b>	<b>90.69</b>
CIFAR100	SimCLR	256	51.65	47.55	53.17	52.05	51.36	52.75
	SimCLR+ArCL(views=4)	256	53.76	49.80	55.68	54.19	52.96	54.83
	SimCLR+ArCL(views=6)	256	<b>54.13</b>	<b>50.74</b>	<b>55.74</b>	<b>54.75</b>	<b>53.46</b>	<b>55.29</b>
	SimCLR	512	52.28	48.09	53.45	52.58	51.53	53.12
	SimCLR+ArCL(views=4)	512	53.40	50.16	54.92	53.77	52.61	54.20
	SimCLR+ArCL(views=6)	512	<b>54.00</b>	<b>50.57</b>	<b>56.24</b>	<b>55.04</b>	<b>53.77</b>	<b>55.60</b>

Aug 1: Grayscale; Aug 2: RandomCrop; Aug 3: HorizontalFlip; Aug 4: ColorJitter;  
Aug 5: Aug 1 + Aug 4.

# Performance on OOD Datasets

	Epochs	Aircraft	Caltech101	Cars	CIFAR10	CIFAR100	DTD	Flowers	Food	Pets	<b>Avg</b>
Linear	MoCo	41.79	87.92	39.31	92.28	74.90	73.88	90.07	68.95	83.30	72.49
	MoCo + AAL (views=2)	40.53	87.80	38.64	92.23	75.14	<b>74.95</b>	88.64	69.24	83.17	72.26
	MoCo + ArCL (views=2)	<b>44.29</b>	<b>89.79</b>	<b>42.15</b>	<b>93.07</b>	<b>76.70</b>	74.20	<b>90.40</b>	<b>70.94</b>	<b>83.68</b>	<b>73.91</b>
	MoCo + AAL (views=3)	40.41	87.79	42.09	92.64	75.31	74.89	89.23	69.37	83.79	72.84
	MoCo + ArCL (views=3)	<b>44.57</b>	<b>89.48</b>	<b>42.11</b>	<b>93.29</b>	<b>77.33</b>	<b>74.63</b>	<b>91.13</b>	<b>71.16</b>	<b>84.23</b>	<b>74.21</b>
	Finetune	MoCo	83.56	82.54	85.09	95.89	71.81	69.95	95.26	76.81	88.83
MoCo + AAL (views=2)		83.87	82.76	85.90	<b>96.38</b>	71.43	<b>72.71</b>	95.50	76.95	<b>89.05</b>	83.84
MoCo + ArCL (views=2)		<b>86.05</b>	<b>87.38</b>	<b>87.28</b>	96.33	<b>79.39</b>	72.18	<b>95.89</b>	<b>81.36</b>	89.03	<b>86.10</b>
MoCo + AAL (views=3)		83.07	83.21	85.19	96.37	72.02	72.55	95.74	79.62	88.83	84.07
MoCo + ArCL (views=3)		<b>84.03</b>	<b>87.64</b>	<b>86.34</b>	<b>96.88</b>	<b>80.98</b>	<b>72.87</b>	<b>96.14</b>	<b>81.90</b>	<b>89.20</b>	<b>86.22</b>

AAL: Average Alignment Loss.

ArCL: Augmentation-robust Contrastive Loss.

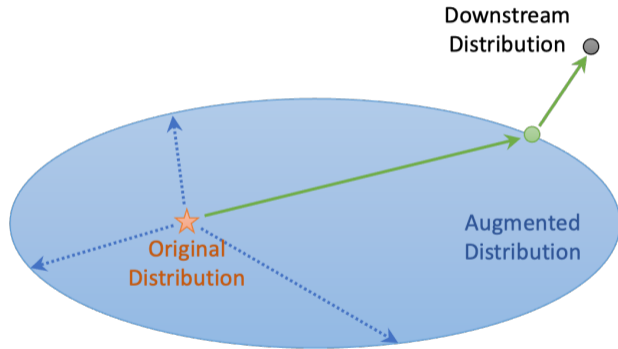
## Short Summary

- We show that contrastive learning fails to learn augmentation-invariant features, which limits its transferability.
- We propose a theory-inspired loss ArCL which can be easily integrated with existing contrastive learning algorithms.
- We empirically verify that ArCL significantly improves the transferability of contrastive learning.

PS: In another ICLR'23 paper, we improve the transferability from the SNE perspective (see “Your Contrastive Learning Is Secretly Doing Stochastic Neighbor Embedding”).

## Conclusion

Training data distribution  $\rightarrow$  Any downstream data distribution?



The Capability Boundary of Self-Supervised Contrastive Learning.

# Thank you!



Interns and visitors are welcome!

Let's explore the most cutting-edge and innovative research together!