

Contrastive learning for images

Ildus Sadrdinov, 31.01.23

Self-supervised pre-training

generative

generative
pre-text tasks

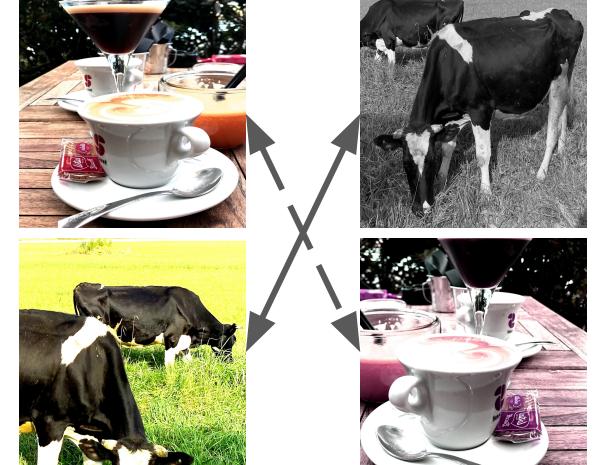


discriminative

discriminative
pre-text tasks



contrastive
tasks



Plan

- **Introduction to information theory**
 - Mutual information & InfoNCE
- **Contrastive learning with negative examples**
 - SimCLR, MoCo
- **Contrastive learning without negative examples**
 - BYOL, SimSiam
- **Contrastive learning inspired by clustering**
 - DeepCluster, SwAV
- **Bonus**
 - Dense Contrastive Learning
 - Supervised Contrastive Learning

Plan

- **Introduction to information theory**
 - Mutual information & InfoNCE
- **Contrastive learning with negative examples**
 - SimCLR, MoCo
- **Contrastive learning without negative examples**
 - BYOL, SimSiam
- **Contrastive learning inspired by clustering**
 - DeepCluster, SwAV
- **Bonus**
 - Dense Contrastive Learning
 - Supervised Contrastive Learning

Information theory

Is it possible to quantify information stored in a random variable?

Information theory

Is it possible to quantify information stored in a random variable?

- (Differential) Entropy: $H(X) = -\mathbb{E}_{p(X)}[\log p(X)] = - \int p(x) \log p(x) dx$

Information theory

Is it possible to quantify information stored in a random variable?

- (Differential) Entropy: $H(X) = -\mathbb{E}_{p(X)}[\log p(X)] = - \int p(x) \log p(x) dx$
- Joint Entropy: $H(X, Y) = -\mathbb{E}_{p(X, Y)}[\log p(X, Y)] = - \int p(x, y) \log p(x, y) dx dy$

Information theory

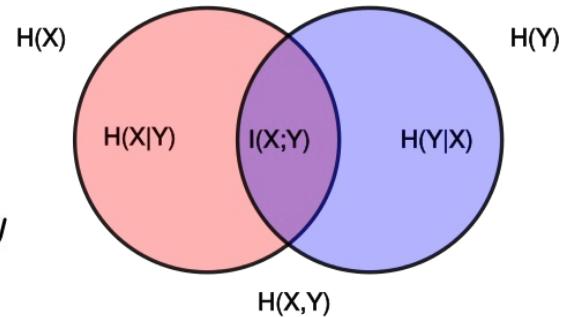
Is it possible to quantify information stored in a random variable?

- (Differential) Entropy: $H(X) = -\mathbb{E}_{p(X)} [\log p(X)] = - \int p(x) \log p(x) dx$
- Joint Entropy: $H(X, Y) = -\mathbb{E}_{p(X, Y)} [\log p(X, Y)] = - \int p(x, y) \log p(x, y) dx dy$
- Conditional Entropy: $H(Y|X) = -\mathbb{E}_{p(X, Y)} [p(Y|X)] = - \int p(x, y) \log \frac{p(x, y)}{p(x)} dx dy$

Information is measured in bits (if $\log = \log_2$) or in nats (if $\log = \ln$)

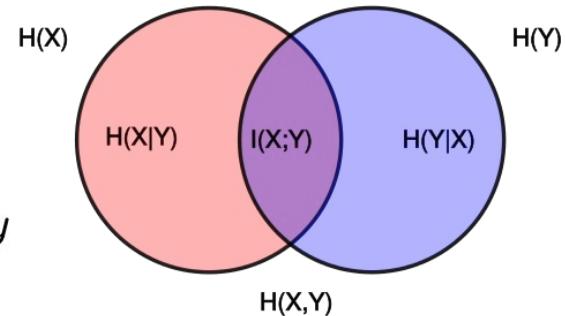
Mutual information

$$I(X;Y) = KL\left(P(X,Y) \middle\| P(X)P(Y)\right) = \int p(x,y) \log \frac{p(x,y)}{p(x)p(y)} dx dy$$



Mutual information

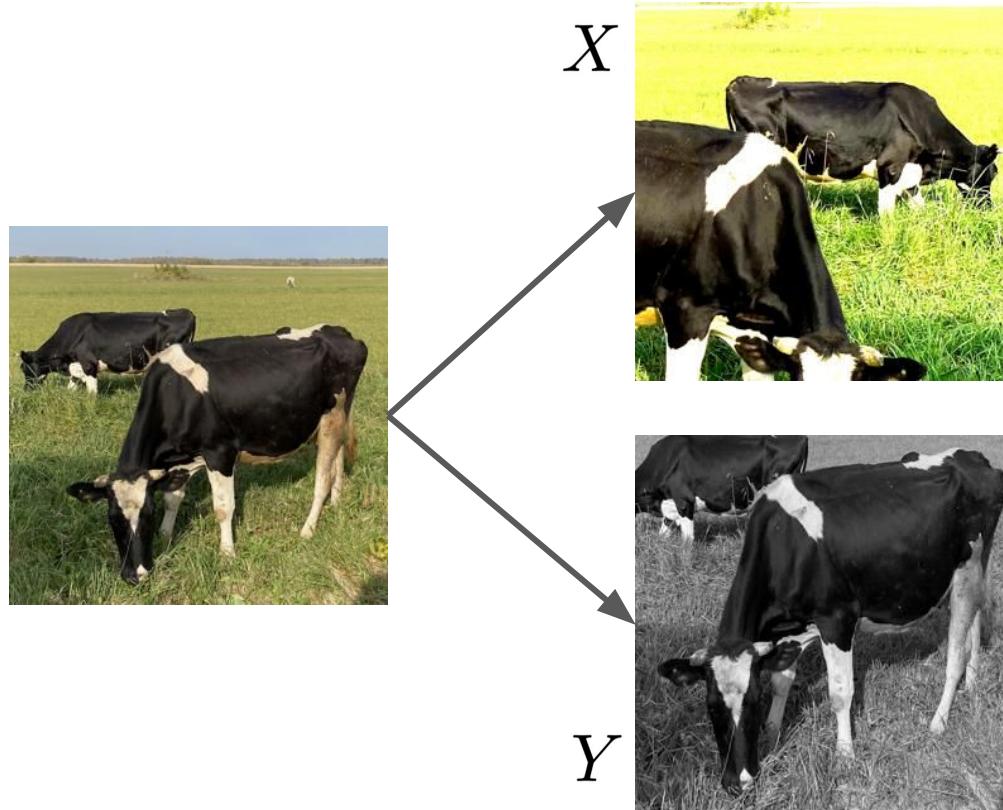
$$I(X;Y) = KL\left(P(X,Y) \middle\| P(X)P(Y)\right) = \int p(x,y) \log \frac{p(x,y)}{p(x)p(y)} dx dy$$



Properties:

- $I(X;Y) \geq 0$
- $I(X;Y) = I(Y;X)$
- $I(X;Y) \equiv H(X) - H(X | Y)$
 $\equiv H(Y) - H(Y | X)$
 $\equiv H(X) + H(Y) - H(X, Y)$
 $\equiv H(X, Y) - H(X | Y) - H(Y | X)$
- When $I(X;Y) = 0$?
- When $I(X;Y) \rightarrow \max$?

So what? (or the main idea of contrastive learning)



$$I\left(f_{\theta}(X), f_{\theta}(Y)\right) \rightarrow \max_{\theta}$$

f_{θ} – our neural network with weights θ

Let's go training!

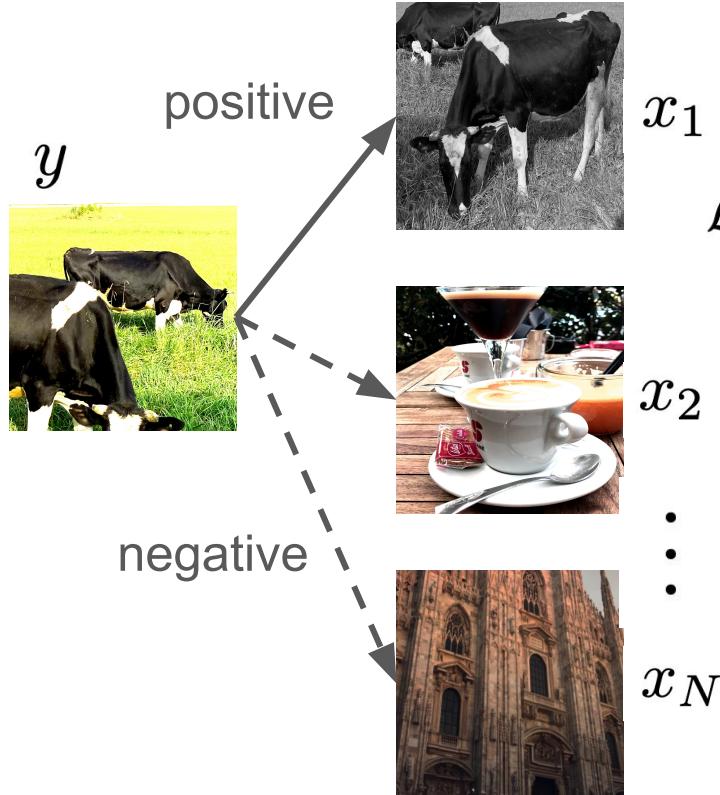
$$I(X, Y) = \int p(x, y) \log \left(\frac{p(x, y)}{p(x)p(y)} \right) dx dy = \mathbb{E}_{p(x, y)} \left[\log \frac{p(x, y)}{p(x)p(y)} \right]$$

Let's go training! Oh, wait...

$$I(X, Y) = \int p(x, y) \log \left(\frac{p(x, y)}{p(x)p(y)} \right) dx dy = \mathbb{E}_{p(x, y)} \left[\log \frac{p(x, y)}{p(x)p(y)} \right]$$

We do not know any of these PDFs! :(

InfoNCE loss and negative examples



$$\mathcal{L}_{NCE}(\theta) = \mathbb{E}_{p(x_{1:N}, y)} \left[-\log \frac{e^{f_\theta(x_1, y)}}{\sum_{n=1}^N e^{f_\theta(x_n, y)}} \right] \rightarrow \min_{\theta}$$

Noise Contrastive Estimation

InfoNCE is a lower bound for mutual information

$$I(X_1; Y) \geq \log N - \mathcal{L}_{NCE}$$

$$\mathcal{L}_{NCE}(\theta) = \mathbb{E}_{p(x_{1:N}, y)} \left[-\log \frac{e^{f_\theta(x_1, y)}}{\sum_{n=1}^N e^{f_\theta(x_n, y)}} \right]$$

Plan

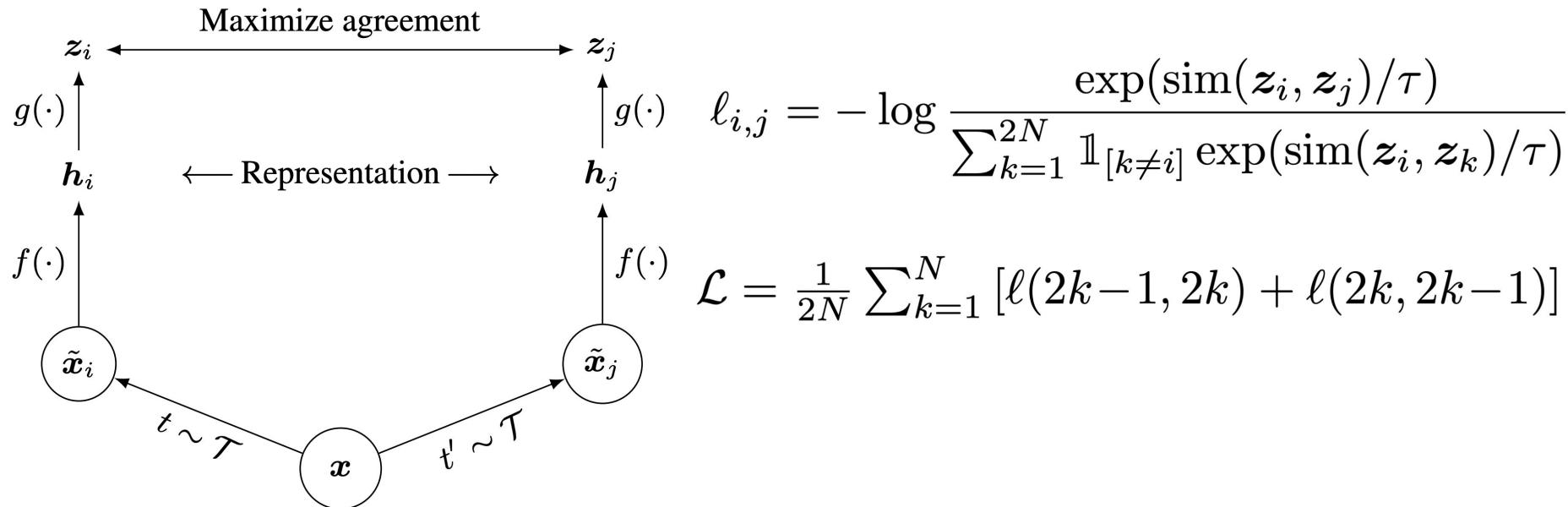
- **Introduction to information theory**
 - Mutual information & InfoNCE
- **Contrastive learning with negative examples**
 - SimCLR, MoCo
- **Contrastive learning without negative examples**
 - BYOL, SimSiam
- **Contrastive learning inspired by clustering**
 - DeepCluster, SwAV
- **Bonus**
 - Dense Contrastive Learning
 - Supervised Contrastive Learning

SimCLR

A Simple Framework for Contrastive Learning of Visual Representations
Google Research

- Make a double batch of augmented images (each image has 2 view)
- Use each of two views both as an anchor (y) and a positive example (x_1)
- Very intense augmentations to make the task meaningful

SimCLR: scheme



SimCLR: augmentations

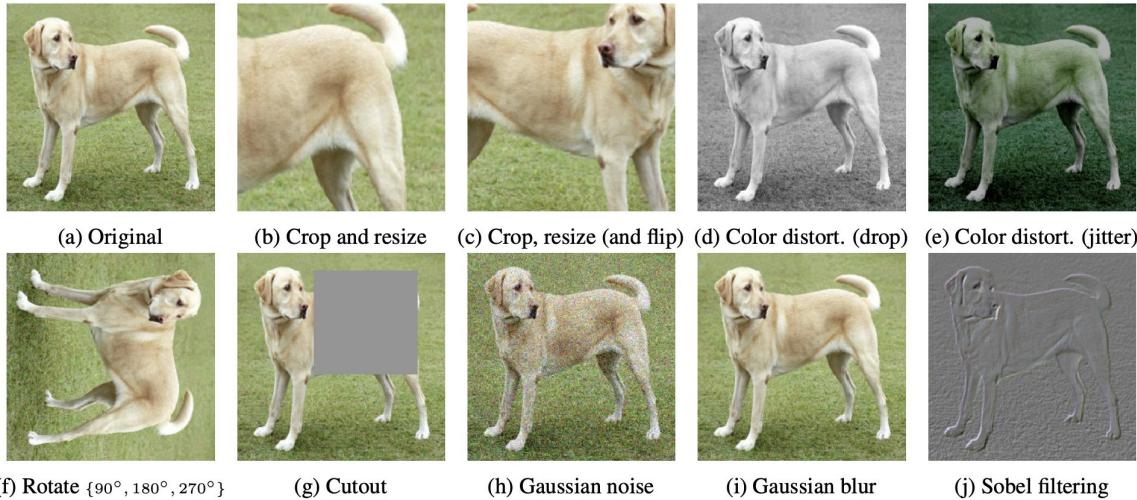


Figure 4. Illustrations of the studied data augmentation operators. Each augmentation can transform data stochastically with some internal parameters (e.g. rotation degree, noise level). Note that we *only* test these operators in ablation, the *augmentation policy used to train our models* only includes *random crop (with flip and resize), color distortion, and Gaussian blur*. (Original image cc-by: Von.grzanka)

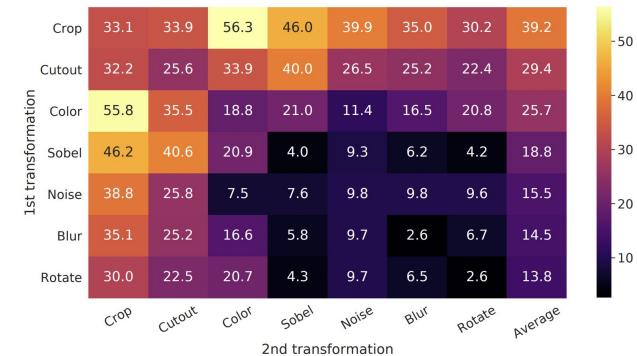
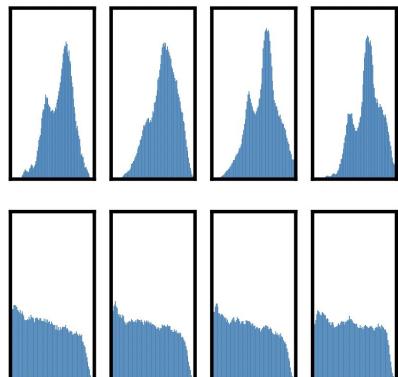
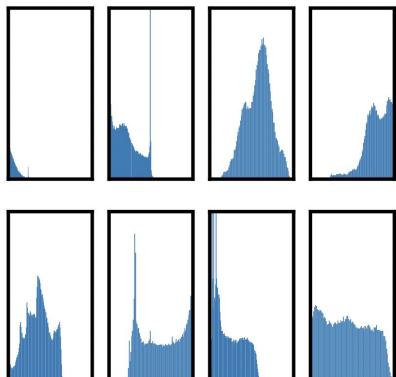


Figure 5. Linear evaluation (ImageNet top-1 accuracy) under individual or composition of data augmentations, applied only to one branch. For all columns but the last, diagonal entries correspond to single transformation, and off-diagonals correspond to composition of two transformations (applied sequentially). The last column reflects the average over the row.

SimCLR: augmentations



(a) Without color distortion.



(b) With color distortion.

Figure 6. Histograms of pixel intensities (over all channels) for different crops of two different images (i.e. two rows). The image for the first row is from Figure 4. All axes have the same range.

Methods	Color distortion strength					AutoAug
	1/8	1/4	1/2	1	1 (+Blur)	
SimCLR	59.6	61.0	62.6	63.2	64.5	61.1
Supervised	77.0	76.7	76.5	75.7	75.4	77.1

Table 1. Top-1 accuracy of unsupervised ResNet-50 using linear evaluation and supervised ResNet-50⁵, under varied color distortion strength (see Appendix A) and other data transformations. Strength 1 (+Blur) is our default data augmentation policy.

SimCLR: projection head

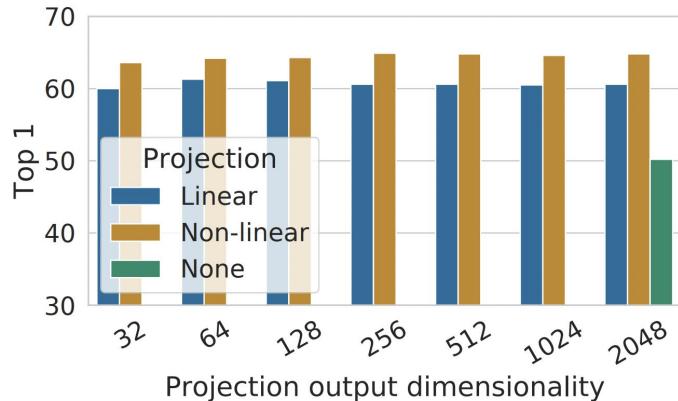


Figure 8. Linear evaluation of representations with different projection heads $g(\cdot)$ and various dimensions of $\mathbf{z} = g(\mathbf{h})$. The representation \mathbf{h} (before projection) is 2048-dimensional here.

What to predict?	Random guess	Representation \mathbf{h}	Representation $g(\mathbf{h})$
Color vs grayscale	80	99.3	97.4
Rotation	25	67.6	25.6
Orig. vs corrupted	50	99.5	59.6
Orig. vs Sobel filtered	50	96.6	56.3

Table 3. Accuracy of training additional MLPs on different representations to predict the transformation applied. Other than crop and color augmentation, we additionally and independently add rotation (one of $\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}$), Gaussian noise, and Sobel filtering transformation during the pretraining for the last three rows. Both \mathbf{h} and $g(\mathbf{h})$ are of the same dimensionality, i.e. 2048.

SimCLR: model size and training time

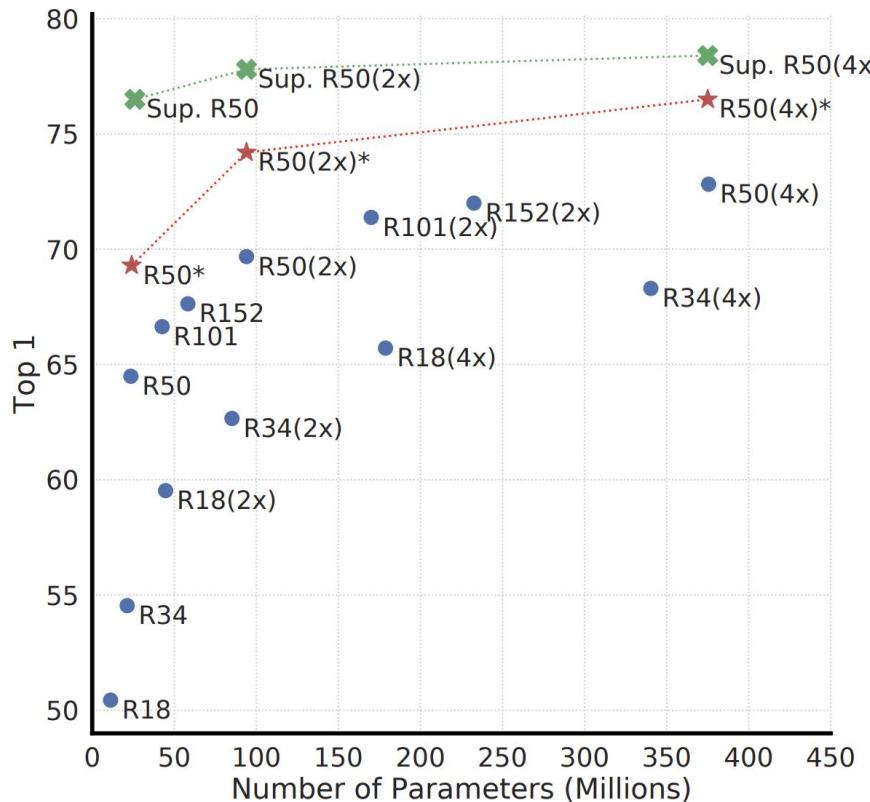


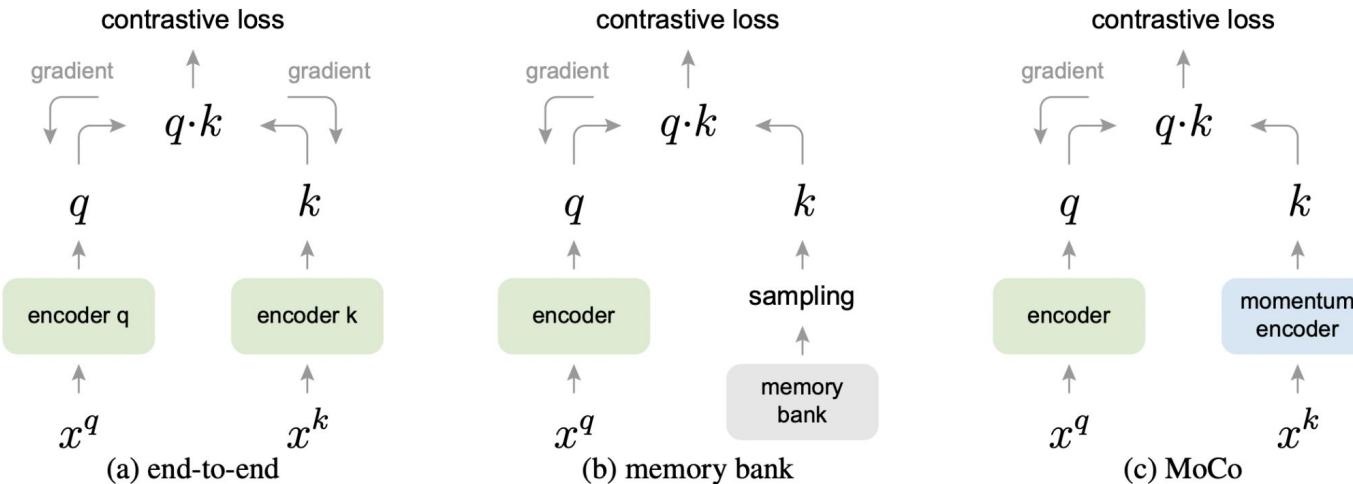
Figure 7. Linear evaluation of models with varied depth and width. Models in blue dots are ours trained for 100 epochs, models in red stars are ours trained for 1000 epochs, and models in green crosses are supervised ResNets trained for 90 epochs⁷ (He et al., 2016).

MoCo

Momentum Contrast, FAIR

- Momentum encoder makes contrastive learning asymmetric
- “Memory bank” to increase the number of negatives

$$\theta_k \leftarrow m\theta_k + (1 - m)\theta_q$$



[He et al., 2019](#)

MoCo: ablations

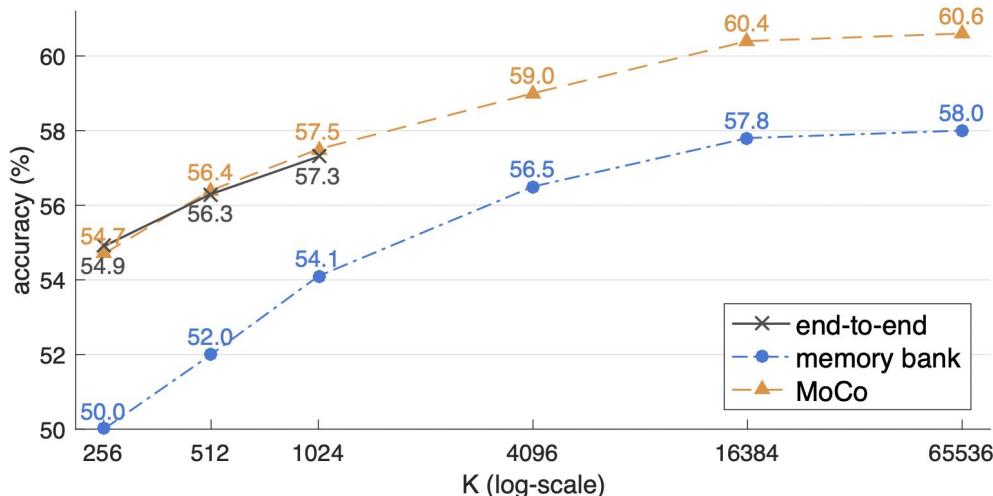


Figure 3. **Comparison of three contrastive loss mechanisms** under the ImageNet linear classification protocol. We adopt the same pretext task (Sec. 3.3) and only vary the contrastive loss mechanism (Figure 2). The number of negatives is K in memory bank and MoCo, and is $K-1$ in end-to-end (offset by one because the positive key is in the same mini-batch). The network is ResNet-50.

Ablation: momentum. The table below shows ResNet-50 accuracy with different MoCo momentum values (m in Eqn.(2)) used in pre-training ($K = 4096$ here) :

momentum m	0	0.9	0.99	0.999	0.9999
accuracy (%)	fail	55.2	57.8	59.0	58.9

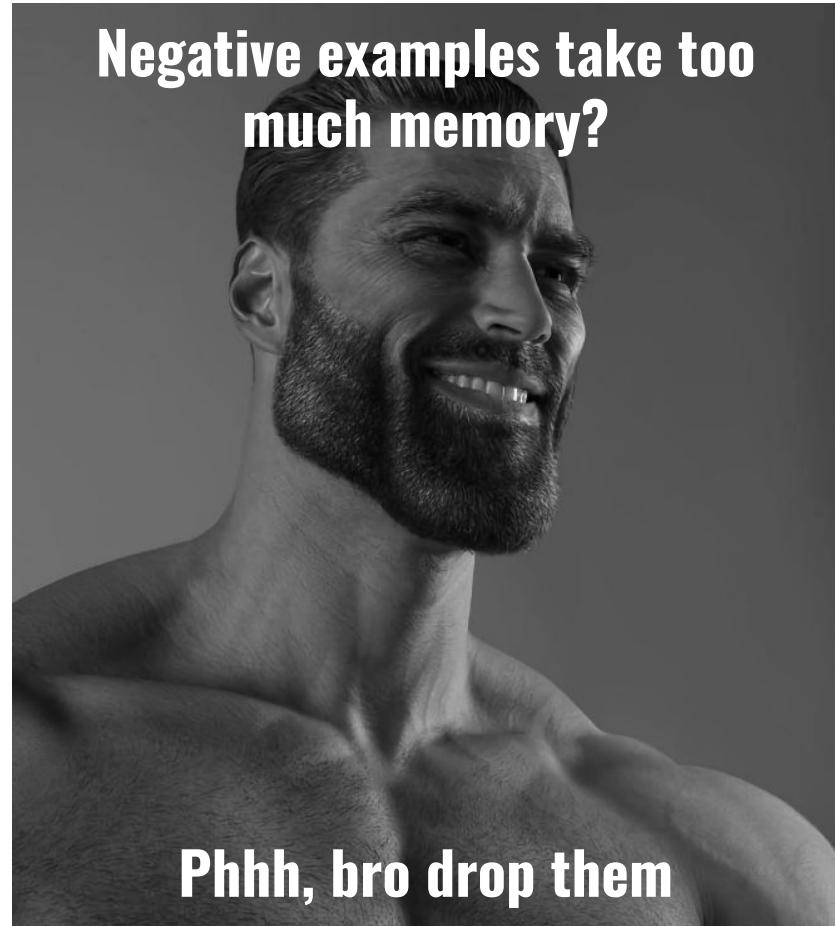
Plan

- **Introduction to information theory**
 - Mutual information & InfoNCE
- **Contrastive learning with negative examples**
 - SimCLR, MoCo
- **Contrastive learning without negative examples**
 - BYOL, SimSiam
- **Contrastive learning inspired by clustering**
 - DeepCluster, SwAV
- **Bonus**
 - Dense Contrastive Learning
 - Supervised Contrastive Learning

BYOL

Bootstrap Your Own Latent, DeepMind

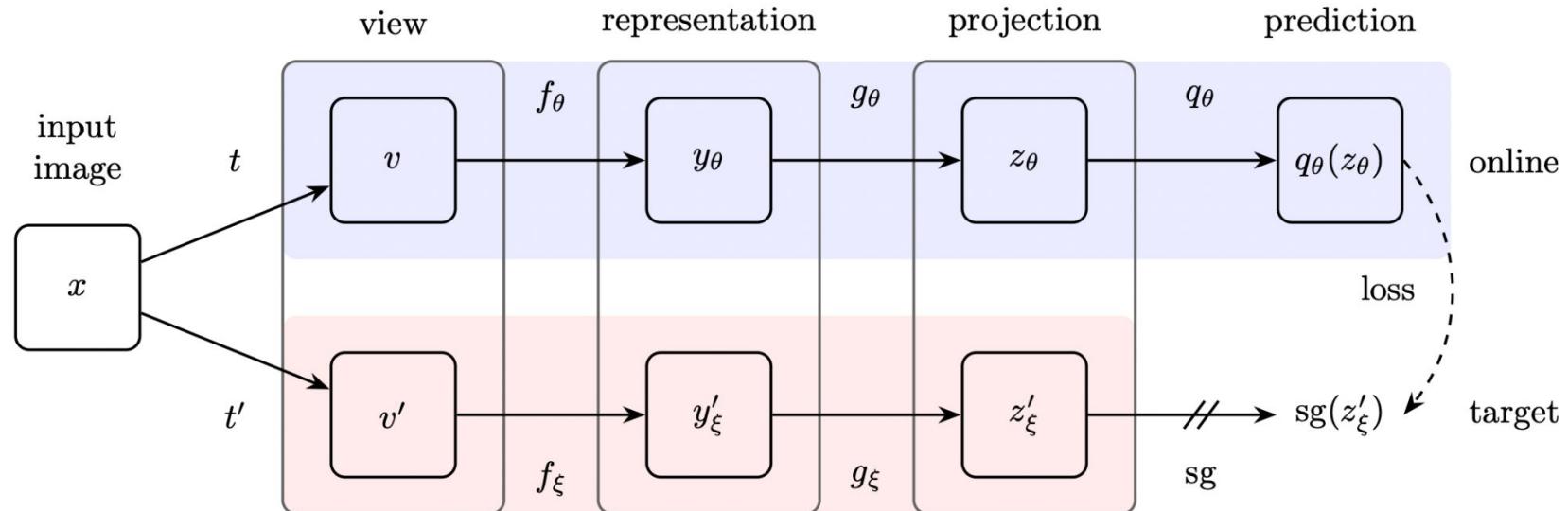
- Further development of momentum encoder from MoCo
- Contrastive learning **without negative examples**
- SOTA in pure contrastive methods



Negative examples take too
much memory?

Phhh, bro drop them

BYOL



$$\mathcal{L}_{\theta,\xi} \triangleq \|\overline{q_\theta}(z_\theta) - \overline{z}'_\xi\|_2^2 = 2 - 2 \cdot \frac{\langle q_\theta(z_\theta), z'_\xi \rangle}{\|q_\theta(z_\theta)\|_2 \cdot \|z'_\xi\|_2} \quad \begin{aligned} \theta &\leftarrow \text{optimizer}(\theta, \nabla_\theta \mathcal{L}_{\theta,\xi}^{\text{BYOL}}, \eta), \\ \xi &\leftarrow \tau \xi + (1 - \tau) \theta, \end{aligned}$$

BYOL

Why do we even need negative examples in contrastive learning?

BYOL

Why do we even need negative examples in contrastive learning?

$$\mathcal{L}_{\theta,\xi} \triangleq \left\| \overline{q_\theta}(z_\theta) - \overline{z}'_\xi \right\|_2^2$$

BYOL: representation collapse

Why do we even need negative examples in contrastive learning?

$$\mathcal{L}_{\theta,\xi} \triangleq \left\| \overline{q_\theta}(z_\theta) - \overline{z}'_\xi \right\|_2^2$$

What if we take constant representations? We end up in a **global minimum** of the loss function with a **representation collapse!**

BYOL: representation collapse

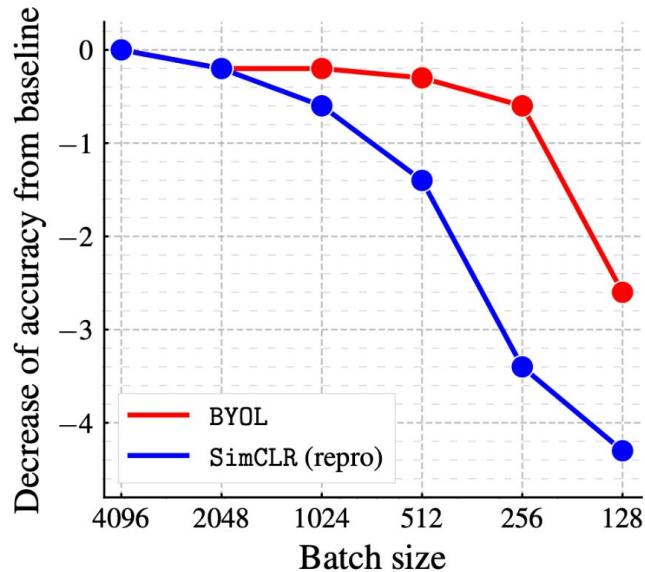
Why does BYOL not collapse in practice?

- Optimization steps are not in the direction of joint gradient $\nabla_{\theta,\xi} \mathcal{L}$

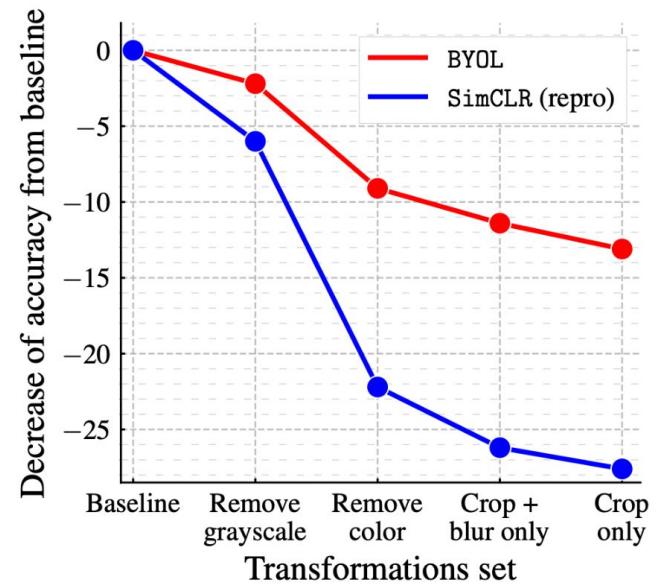
$$\begin{aligned}\theta &\leftarrow \text{optimizer}(\theta, \nabla_{\theta} \mathcal{L}_{\theta,\xi}^{\text{BYOL}}, \eta), \\ \xi &\leftarrow \tau \xi + (1 - \tau) \theta,\end{aligned}$$

- Collapsed global minimum is unstable
- Hyperparameters (EMA, weight decay) are important

BYOL: ablations



(a) Impact of batch size



(b) Impact of progressively removing transformations

Figure 3: Decrease in top-1 accuracy (in % points) of BYOL and our own reproduction of SimCLR at 300 epochs, under linear evaluation on ImageNet.

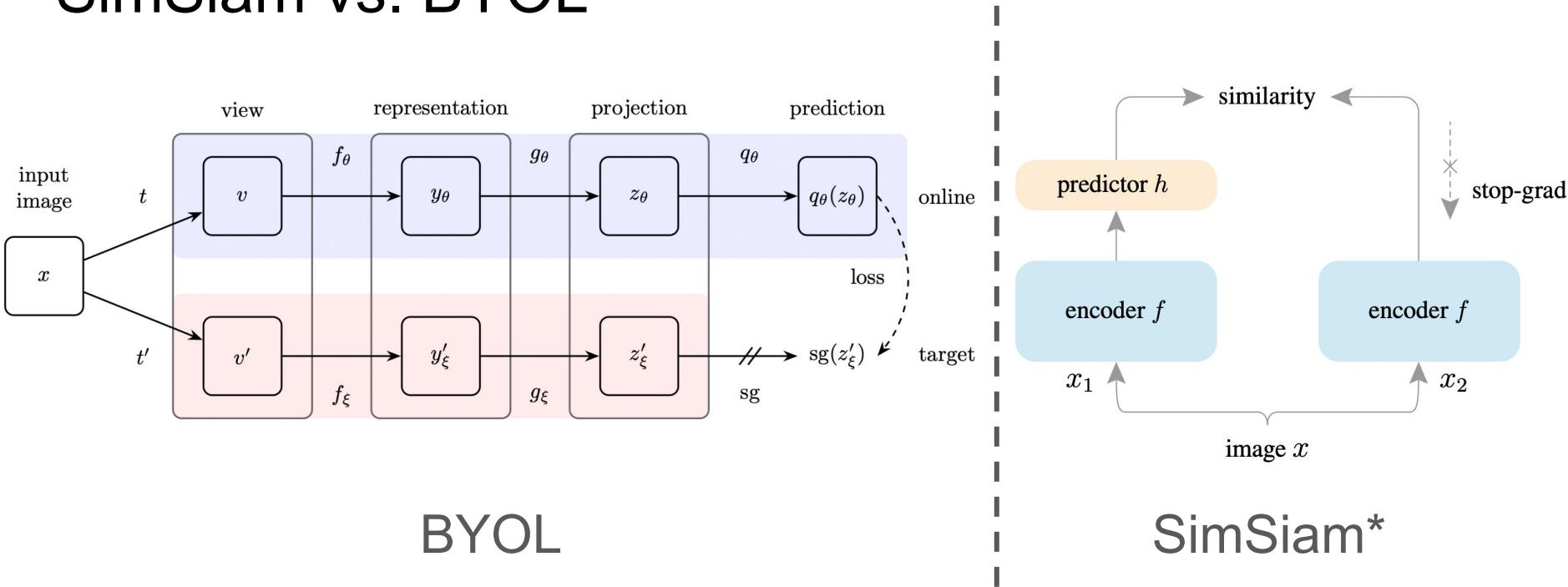
SimSiam

Simple Siamese networks, FAIR

- Momentum encoder is **unnecessary** for contrastive learning
- Overall performance worse than BYOL but still comparable to other methods



SimSiam vs. BYOL



*SimSiam encoder f includes both the convolutional part (BYOL's f_θ) and the MLP projection (BYOL's g_θ)

[Chen and He, 2020](#)

SimSiam: ablations

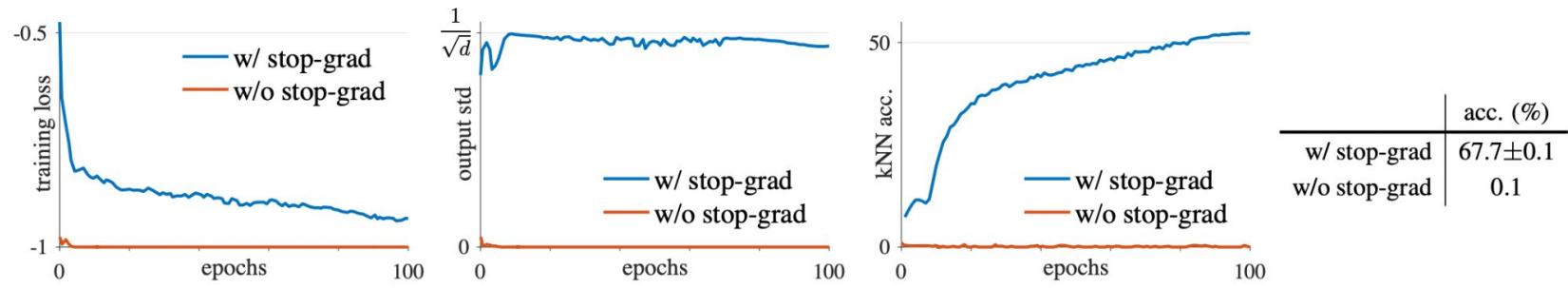


Figure 2. **SimSiam with vs. without stop-gradient.** **Left plot:** training loss. Without stop-gradient it degenerates immediately. **Middle plot:** the per-channel std of the ℓ_2 -normalized output, plotted as the averaged std over all channels. **Right plot:** validation accuracy of a kNN classifier [36] as a monitor of progress. **Table:** ImageNet linear evaluation (“w/ stop-grad” is mean \pm std over 5 trials).

	pred. MLP h	acc. (%)
baseline	lr with cosine decay	67.7
(a)	no pred. MLP	0.1
(b)	fixed random init.	1.5
(c)	lr not decayed	68.1

Table 1. **Effect of prediction MLP** (ImageNet linear evaluation accuracy with 100-epoch pre-training). In all these variants, we use the same schedule for the encoder f (lr with cosine decay).

[Chen and He, 2020](#)

SimSiam is EM algorithm???

$$\mathcal{L}(\theta, \eta) = \mathbb{E}_{x, \mathcal{T}} \left[\left\| \mathcal{F}_\theta(\mathcal{T}(x)) - \eta_x \right\|_2^2 \right]$$

$$\theta^t \quad \leftarrow \quad \arg \min_{\theta} \quad \mathcal{L}(\theta, \eta^{t-1})$$

$$\eta^t \quad \leftarrow \quad \arg \min_{\eta} \quad \mathcal{L}(\theta^t, \eta)$$

x – image

$\mathcal{T}(\cdot)$ – augmentation

$\mathcal{F}_\theta(\cdot)$ – encoder

η_x – representation of

(not augmented) image x

SimSiam is EM algorithm???

$$\mathcal{L}(\theta, \eta) = \mathbb{E}_{x, \mathcal{T}} \left[\left\| \mathcal{F}_\theta(\mathcal{T}(x)) - \eta_x \right\|_2^2 \right]$$

$$\theta^t \quad \leftarrow \quad \arg \min_{\theta} \mathcal{L}(\theta, \eta^{t-1})$$

x – image

$\mathcal{T}(\cdot)$ – augmentation

$\mathcal{F}_\theta(\cdot)$ – encoder

η_x – representation of

(not augmented) image x

$$\eta_x^t \leftarrow \mathbb{E}_{\mathcal{T}} \left[\mathcal{F}_{\theta^t}(\mathcal{T}(x)) \right]$$

SimSiam is EM algorithm???

$$\mathcal{L}(\theta, \eta) = \mathbb{E}_{x, \mathcal{T}} \left[\left\| \mathcal{F}_\theta(\mathcal{T}(x)) - \eta_x \right\|_2^2 \right]$$

$$\eta_x^t \leftarrow \mathbb{E}_{\mathcal{T}} \left[\mathcal{F}_{\theta^t}(\mathcal{T}(x)) \right] \rightarrow \eta_x^t \leftarrow \mathcal{F}_{\theta^t}(\mathcal{T}'(x))$$

x – image

$\mathcal{T}(\cdot)$ – augmentation

$\mathcal{F}_\theta(\cdot)$ – encoder

η_x – representation of

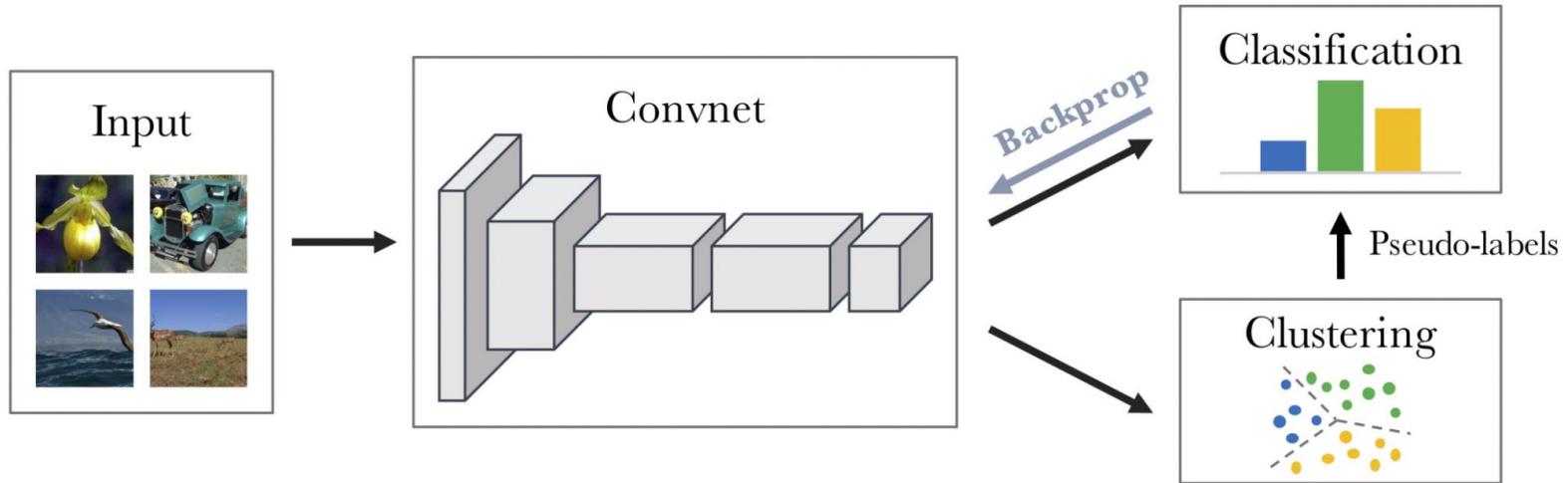
(not augmented) image x

$$\theta^{t+1} \leftarrow \arg \min_{\theta} \mathbb{E}_{x, \mathcal{T}} \left[\left\| \mathcal{F}_\theta(\mathcal{T}(x)) - \mathcal{F}_{\theta^t}(\mathcal{T}'(x)) \right\|_2^2 \right]$$

Plan

- **Introduction to information theory**
 - Mutual information & InfoNCE
- **Contrastive learning with negative examples**
 - SimCLR, MoCo
- **Contrastive learning without negative examples**
 - BYOL, SimSiam
- **Contrastive learning inspired by clustering**
 - DeepCluster, SwAV
- **Bonus**
 - Dense Contrastive Learning
 - Supervised Contrastive Learning

DeepCluster

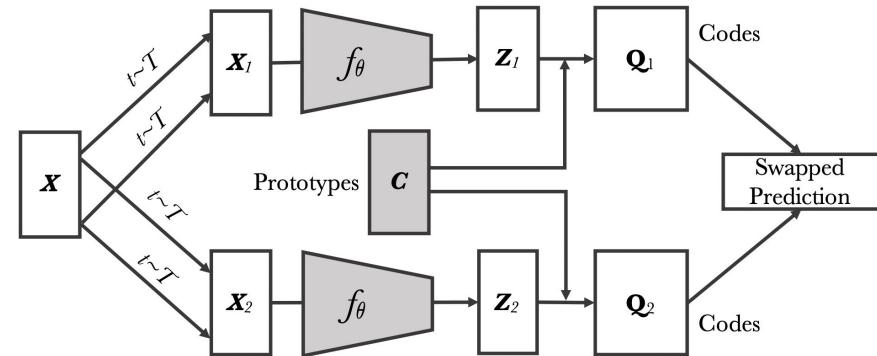


$$\min_{C \in \mathbb{R}^{d \times k}} \frac{1}{N} \sum_{n=1}^N \min_{y_n \in \{0,1\}^k} \|f_\theta(x_n) - Cy_n\|_2^2 \quad \text{such that} \quad y_n^\top \mathbf{1}_k = 1$$

SwAV

Swapping Assignments between Views, FAIR

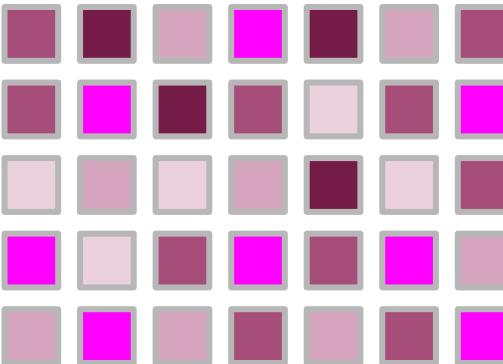
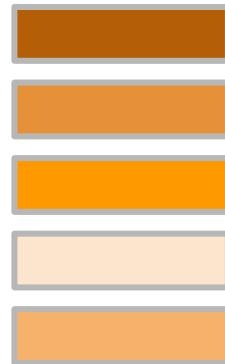
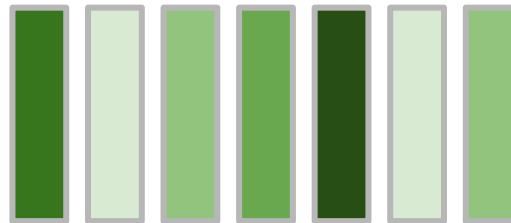
- Further development of DeepCluster
- Make cluster centroids learnable
- Do online clustering
- Swap cluster assignments between two augmented views of the image (i.e. bring contrastive component to the clustering setup)



SwAV: scheme

Image representations

$$\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_B]$$



$$\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_K]$$

Cluster prototypes (“centroids”)

$$\mathbf{Q} \in \mathbb{R}_+^{K \times B}$$

Cluster
assignments

[Caron et al., 2020](#)

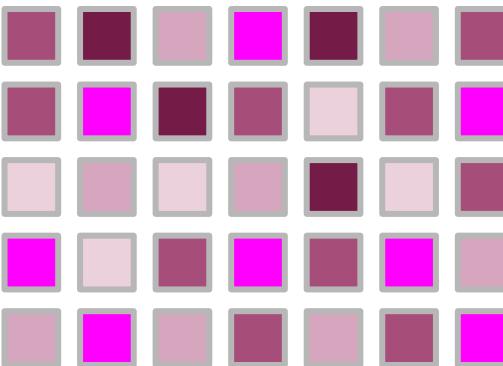
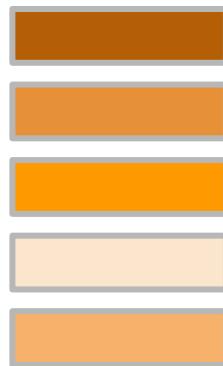
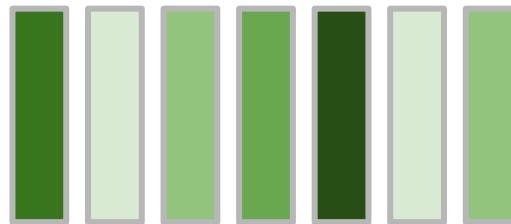
SwAV: scheme

$$\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_K]$$

Cluster prototypes (“centroids”)

Image representations

$$\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_B]$$



$$\max_{\mathbf{Q} \in \mathcal{Q}} \text{Tr} (\mathbf{Q}^\top \mathbf{C}^\top \mathbf{Z}) + \varepsilon H(\mathbf{Q})$$

$$H(\mathbf{Q}) = - \sum_{ij} \mathbf{Q}_{ij} \log \mathbf{Q}_{ij}$$

$$\mathcal{Q} = \left\{ \mathbf{Q} \in \mathbb{R}_+^{K \times B} \mid \mathbf{Q}\mathbf{1}_B = \frac{1}{K}\mathbf{1}_K, \mathbf{Q}^\top \mathbf{1}_K = \frac{1}{B}\mathbf{1}_B \right\}$$

$$\mathbf{Q} \in \mathbb{R}_+^{K \times B}$$

Cluster
assignments

[Caron et al., 2020](#)

SwAV: scheme

$$\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_K]$$

Cluster prototypes (“centroids”)

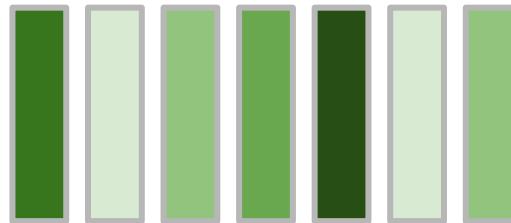
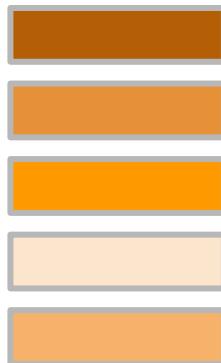
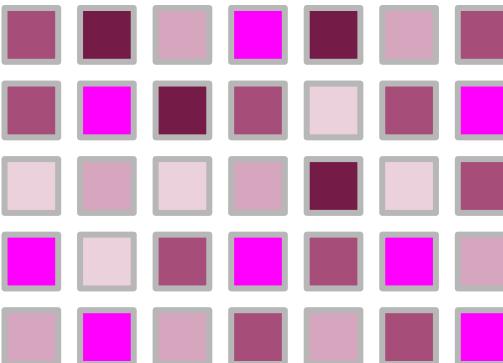


Image representations

$$\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_B]$$



$$\mathbf{Q} \in \mathbb{R}_+^{K \times B}$$

Cluster
assignments

$$\max_{\mathbf{Q} \in \mathcal{Q}} \text{Tr} (\mathbf{Q}^\top \mathbf{C}^\top \mathbf{Z}) + \varepsilon H(\mathbf{Q})$$

$$H(\mathbf{Q}) = - \sum_{ij} \mathbf{Q}_{ij} \log \mathbf{Q}_{ij}$$

$$\mathcal{Q} = \left\{ \mathbf{Q} \in \mathbb{R}_+^{K \times B} \mid \mathbf{Q}\mathbf{1}_B = \frac{1}{K}\mathbf{1}_K, \mathbf{Q}^\top \mathbf{1}_K = \frac{1}{B}\mathbf{1}_B \right\}$$

$$\mathbf{Q}^* = \text{Diag}(\mathbf{u}) \exp \left(\frac{\mathbf{C}^\top \mathbf{Z}}{\varepsilon} \right) \text{Diag}(\mathbf{v})$$

Sinkhorn–Knopp
algorithm

[Caron et al., 2020](#)

SwAV: scheme

$$\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_K]$$

Cluster prototypes (“centroids”)

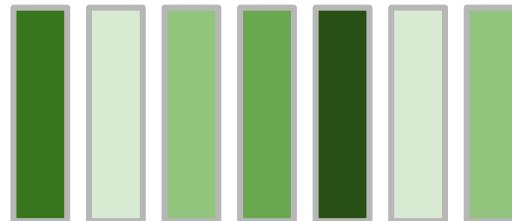
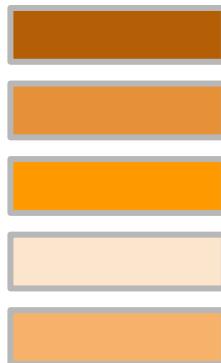
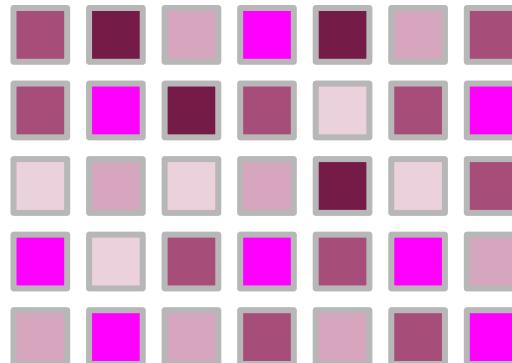


Image representations

$$\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_B]$$

$$\mathbf{Q} \in \mathbb{R}_+^{K \times B}$$

Cluster
assignments



$$\mathbf{Q}^* = \text{Diag}(\mathbf{u}) \exp\left(\frac{\mathbf{C}^\top \mathbf{Z}}{\varepsilon}\right) \text{Diag}(\mathbf{v})$$

$$\mathbf{p}_t^{(k)} = \frac{\exp\left(\frac{1}{\tau} \mathbf{z}_t^\top \mathbf{c}_k\right)}{\sum_{k'} \exp\left(\frac{1}{\tau} \mathbf{z}_t^\top \mathbf{c}_{k'}\right)}$$

$$\ell(\mathbf{z}_t, \mathbf{q}_s) = - \sum_k \mathbf{q}_s^{(k)} \log \mathbf{p}_t^{(k)}$$

SwAV: scheme

$$\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_K]$$

Cluster prototypes (“centroids”)

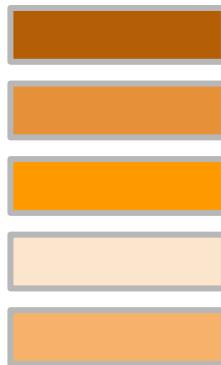
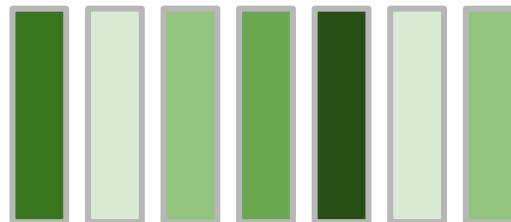


Image representations

$$\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_B]$$



$$\mathbf{Q} \in \mathbb{R}_+^{K \times B}$$

Cluster
assignments

$$\mathbf{Q}^* = \text{Diag}(\mathbf{u}) \exp\left(\frac{\mathbf{C}^\top \mathbf{Z}}{\varepsilon}\right) \text{Diag}(\mathbf{v})$$

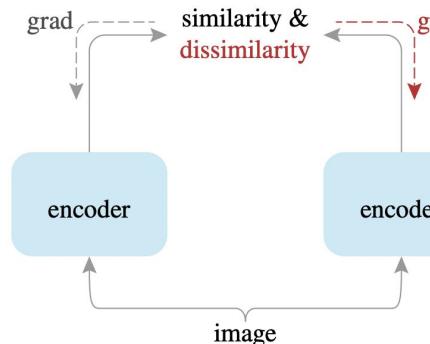
$$\mathbf{p}_t^{(k)} = \frac{\exp\left(\frac{1}{\tau} \mathbf{z}_t^\top \mathbf{c}_k\right)}{\sum_{k'} \exp\left(\frac{1}{\tau} \mathbf{z}_t^\top \mathbf{c}_{k'}\right)}$$

$$\ell(\mathbf{z}_t, \mathbf{q}_s) = - \sum_k \mathbf{q}_s^{(k)} \log \mathbf{p}_t^{(k)}$$

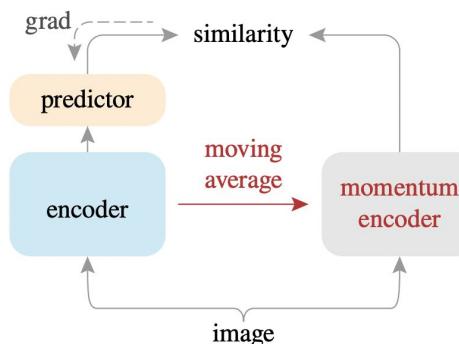
$$L(\mathbf{z}_t, \mathbf{z}_s) = \ell(\mathbf{z}_t, \mathbf{q}_s) + \ell(\mathbf{z}_s, \mathbf{q}_t)$$

[Caron et al., 2020](#)

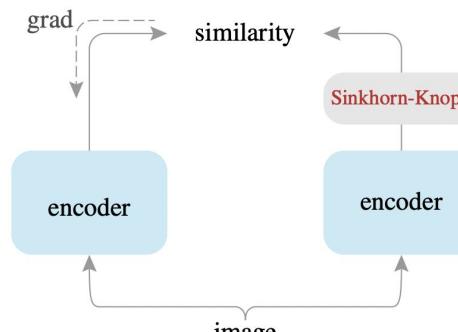
Comparison



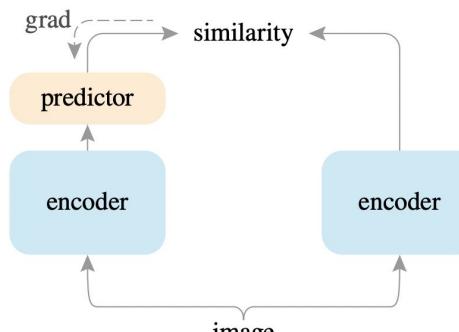
SimCLR



BYOL



SwAV



SimSiam

[Chen and He, 2020](#)

Results: ImageNet LP, VOC & COCO

method	batch size	negative pairs	momentum encoder	100 ep	200 ep	400 ep	800 ep
SimCLR (repro.+)	4096	✓		66.5	68.3	69.8	70.4
MoCo v2 (repro.+)	256	✓	✓	67.4	69.9	71.0	72.2
BYOL (repro.)	4096		✓	66.5	70.6	73.2	74.3
SwAV (repro.+)	4096			66.5	69.1	70.7	71.8
SimSiam	256			68.1	70.0	70.8	71.3

pre-train	VOC 07 detection			VOC 07+12 detection			COCO detection			COCO instance seg.		
	AP ₅₀	AP	AP ₇₅	AP ₅₀	AP	AP ₇₅	AP ₅₀	AP	AP ₇₅	AP ₅₀ ^{mask}	AP ₇₅ ^{mask}	
scratch	35.9	16.8	13.0	60.2	33.8	33.1	44.0	26.4	27.8	46.9	29.3	30.8
ImageNet supervised	74.4	42.4	42.7	81.3	53.5	58.8	58.2	38.2	41.2	54.7	33.3	35.2
SimCLR (repro.+)	75.9	46.8	50.1	81.8	55.5	61.4	57.7	37.9	40.9	54.6	33.3	35.3
MoCo v2 (repro.+)	77.1	48.5	52.5	82.3	57.0	63.3	58.8	39.2	42.5	55.5	34.3	36.6
BYOL (repro.)	77.1	47.0	49.9	81.4	55.3	61.1	57.8	37.9	40.9	54.3	33.2	35.0
SwAV (repro.+)	75.5	46.5	49.6	81.5	55.4	61.4	57.6	37.6	40.3	54.2	33.1	35.1
SimSiam, base	75.5	47.0	50.2	82.0	56.4	62.8	57.5	37.9	40.9	54.2	33.2	35.2
SimSiam, optimal	77.3	48.5	52.5	82.4	57.0	63.7	59.3	39.2	42.1	56.0	34.4	36.7

Results: downstream classification

Method	Food101	CIFAR10	CIFAR100	Birdsnap	SUN397	Cars	Aircraft	VOC2007	DTD	Pets	Caltech-101	Flowers
<i>Linear evaluation:</i>												
BYOL (ours)	75.3	91.3	78.4	57.2	62.2	67.8	60.6	82.5	75.5	90.4	94.2	96.1
SimCLR (repro)	72.8	90.5	74.4	42.4	60.6	49.3	49.8	81.4	75.7	84.6	89.3	92.6
SimCLR [8]	68.4	90.6	71.6	37.4	58.8	50.3	50.3	80.5	74.5	83.6	90.3	91.2
Supervised-IN [8]	72.3	93.6	78.3	53.7	61.9	66.7	61.0	82.8	74.9	91.5	94.5	94.7
<i>Fine-tuned:</i>												
BYOL (ours)	88.5	97.8	86.1	76.3	63.7	91.6	88.1	85.4	76.2	91.7	93.8	97.0
SimCLR (repro)	87.5	97.4	85.3	75.0	63.9	91.4	87.6	84.5	75.4	89.4	91.7	96.6
SimCLR [8]	88.2	97.7	85.9	75.9	63.5	91.3	88.1	84.1	73.2	89.2	92.1	97.0
Supervised-IN [8]	88.3	97.5	86.4	75.8	64.3	92.1	86.0	85.0	74.6	92.1	93.3	97.6
Random init [8]	86.9	95.9	80.2	76.1	53.6	91.4	85.9	67.3	64.8	81.5	72.6	92.0

Results: semi-supervised

Method	Top-1		Top-5		Method	Architecture	Param.	Top-1		Top-5	
	1%	10%	1%	10%				1%	10%	1%	10%
Supervised [77]	25.4	56.4	48.4	80.4	CPC v2 [32]	ResNet-161	305M	-	-	77.9	91.2
InstDisc	-	-	39.2	77.4	SimCLR [8]	ResNet-50 (2×)	94M	58.5	71.7	83.0	91.2
PIRL [35]	-	-	57.2	83.8	BYOL (ours)	ResNet-50 (2×)	94M	62.2	73.5	84.1	91.7
SimCLR [8]	48.3	65.6	75.5	87.8	SimCLR [8]	ResNet-50 (4×)	375M	63.0	74.4	85.8	92.6
BYOL (ours)	53.2	68.8	78.4	89.0	BYOL (ours)	ResNet-50 (4×)	375M	69.1	75.7	87.9	92.5
					BYOL (ours)	ResNet-200 (2×)	250M	71.2	77.7	89.5	93.7

(a) ResNet-50 encoder.

(b) Other ResNet encoder architectures.

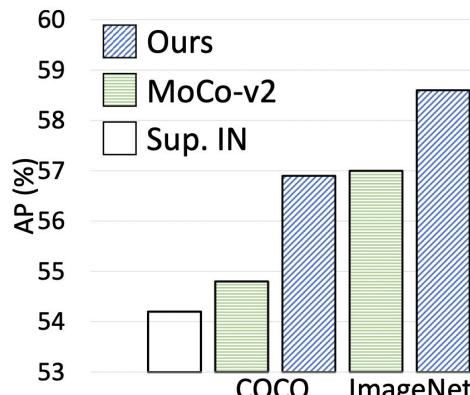
Table 2: Semi-supervised training with a fraction of ImageNet labels.

Plan

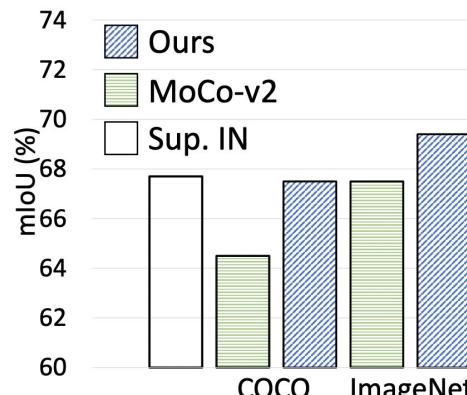
- **Introduction to information theory**
 - Mutual information & InfoNCE
- **Contrastive learning with negative examples**
 - SimCLR, MoCo
- **Contrastive learning without negative examples**
 - BYOL, SimSiam
- **Contrastive learning inspired by clustering**
 - DeepCluster, SwAV
- **Bonus**
 - Dense Contrastive Learning
 - Supervised Contrastive Learning

Dense Contrastive Learning

- **Global contrastive learning** works good for classification tasks
- Do **local contrastive learning** for dense tasks (i.e. detection & segmentation)



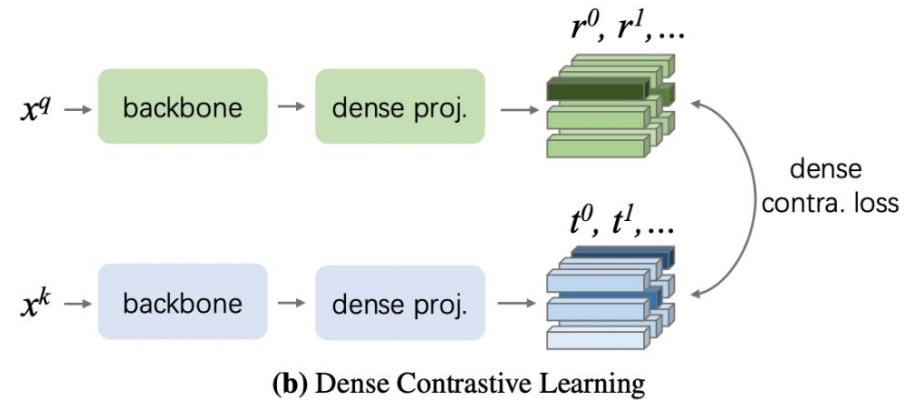
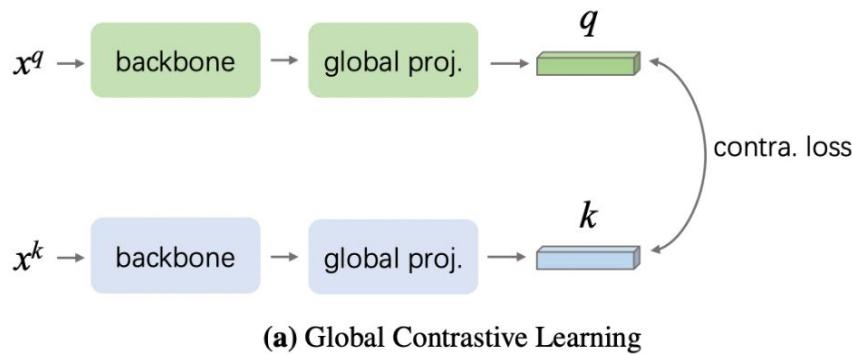
(a) Object Detection



(b) Semantic Segmentation

[Want et al., 2020](#)

Dense Contrastive Learning: scheme



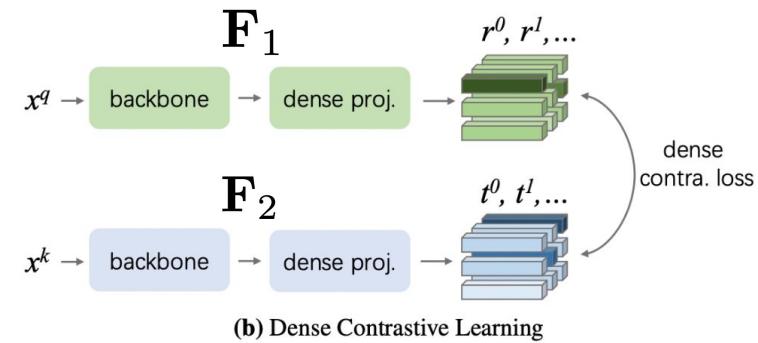
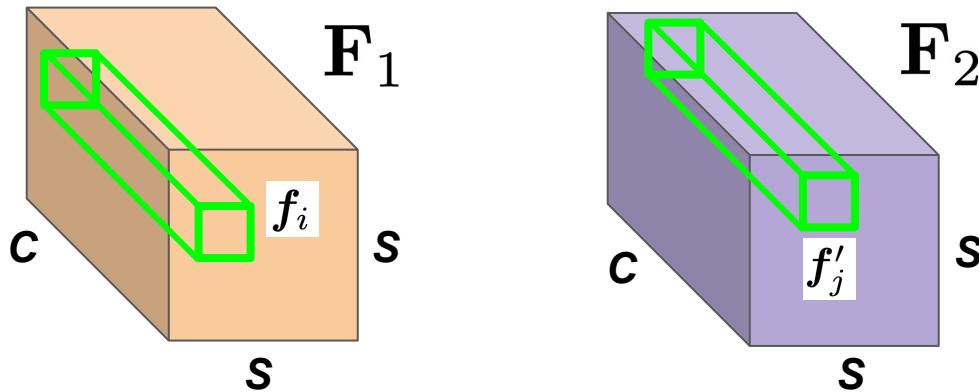
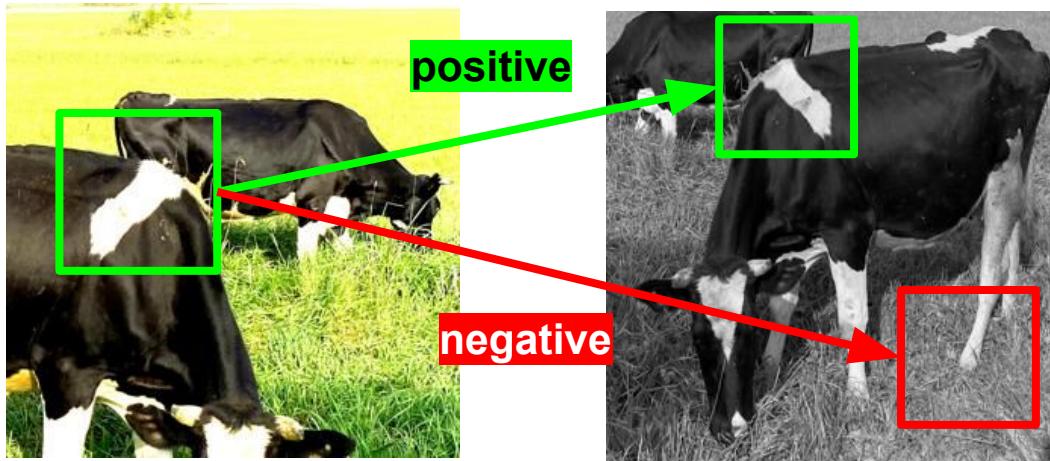
$$\mathcal{L}_q = -\log \frac{\exp(q \cdot k_+ / \tau)}{\exp(q \cdot k_+) + \sum_{k_-} \exp(q \cdot k_- / \tau)}$$

$$\mathcal{L}_r = \frac{1}{S^2} \sum_s -\log \frac{\exp(r^s \cdot t^s_+ / \tau)}{\exp(r^s \cdot t^s_+) + \sum_{t^s_-} \exp(r^s \cdot t^s_- / \tau)}$$

$$\mathcal{L} = (1 - \lambda)\mathcal{L}_q + \lambda\mathcal{L}_r$$

[Want et al., 2020](#)

Dense Contrastive Learning: scheme



$$c_i = \arg \max_j sim(f_i, f'_j)$$

[Want et al., 2020](#)

Dense Contrastive Learning: scheme



Dense Contrastive Learning: results

pre-train	AP ^b	AP ^b ₅₀	AP ^b ₇₅	AP ^m	AP ^m ₅₀	AP ^m ₇₅
random init.	32.8	50.9	35.3	29.9	47.9	32.0
super. IN	39.7	59.5	43.3	35.9	56.6	38.6
MoCo-v2 CC	38.5	58.1	42.1	34.8	55.3	37.3
DenseCL CC	39.6	59.3	43.3	35.7	56.5	38.4
SimCLR IN	38.5	58.0	42.0	34.8	55.2	37.2
BYOL IN	38.4	57.9	41.9	34.9	55.3	37.5
MoCo-v2 IN	39.8	59.8	43.6	36.1	56.9	38.7
DenseCL IN	40.3	59.9	44.3	36.4	57.0	39.2

Table 2 – Object detection and instance segmentation fine-tuned on COCO. ‘CC’ and ‘IN’ indicate the pre-training mod-

pre-train	AP ^b	AP ^b ₅₀	AP ^b ₇₅	AP ^m	AP ^m ₅₀	AP ^m ₇₅
random init.	20.6	34.0	21.5	18.9	31.7	19.8
super. IN	23.6	37.7	25.4	21.8	35.4	23.2
MoCo-v2 CC	22.8	36.4	24.2	20.9	34.6	21.9
DenseCL CC	24.1	38.1	25.6	21.9	36.0	23.0
MoCo-v2 IN	23.8	37.5	25.6	21.8	35.4	23.2
DenseCL IN	24.8	38.8	26.8	22.6	36.8	23.9

Table 3 – Semi-supervised object detection and instance segmentation fine-tuned on COCO. During the fine-tuning,

time/epoch	COCO	ImageNet
MoCo-v2	1'45"	16'48"
DenseCL	1'46"	16'54"

Table 9 – Pre-training time comparison. The training time per epoch is reported. We measure the results on the same 8-GPU machine. The training time overhead introduced by DenseCL is less than 1%.

Dense Contrastive Learning: ablations

λ	Detection			Classification mAP
	AP	AP ₅₀	AP ₇₅	
0.0	54.7	81.0	60.6	82.6
0.1	55.2	81.4	61.4	82.9
0.3	56.2	81.5	62.6	83.3
0.5	56.7	81.7	63.0	82.9
0.7	56.8	81.9	63.1	81.0
0.9	55.5	80.9	61.3	77.8
1.0*	53.5	79.5	58.8	68.9

Table 5 – Ablation study of weight λ . $\lambda = 0$ is the MoCo-v2 baseline. $\lambda = 0.5$ shows the best trade-off between detection and classification. '*' indicates training with warm-up, as discussed in Section 3.4.

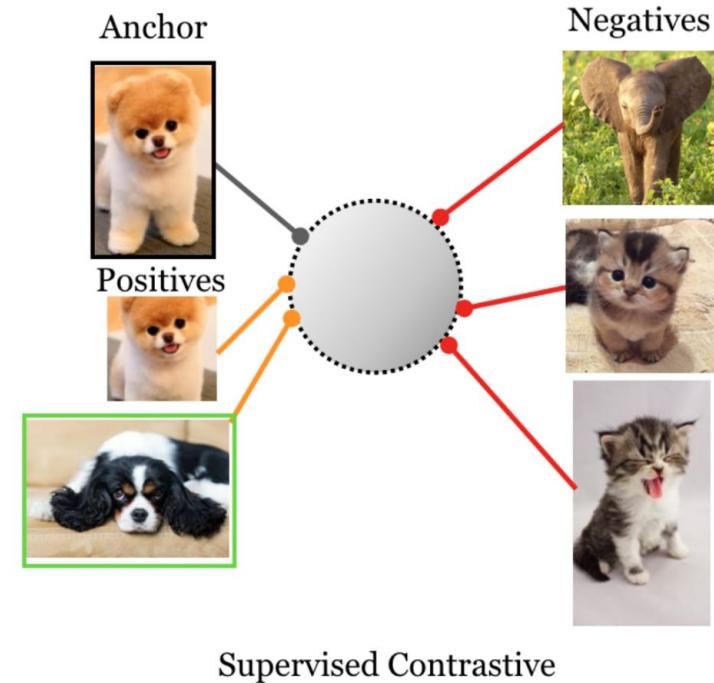
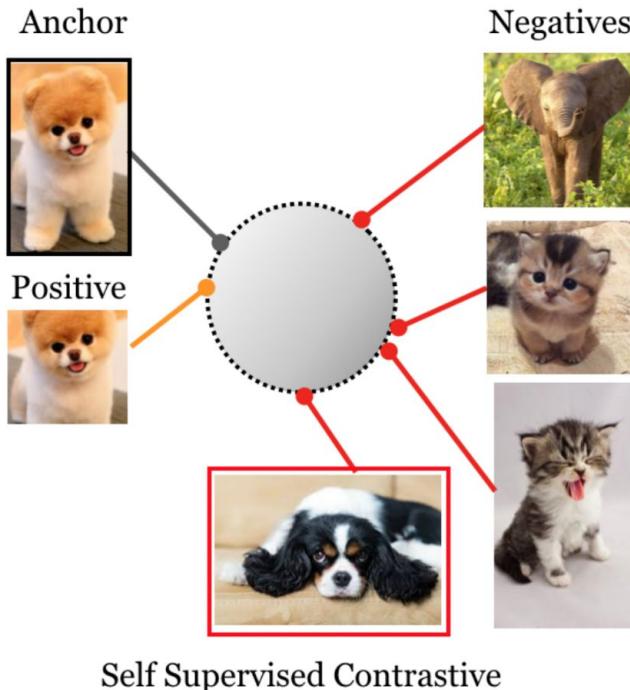
strategy	Detection			Classification mAP
	AP	AP ₅₀	AP ₇₅	
random	56.0	81.3	62.0	81.7
max-sim Θ	56.0	81.5	62.1	81.8
max-sim F	56.7	81.7	63.0	82.9

Table 6 – Ablation study of matching strategy. To extract the dense correspondence according to the backbone features F_1 and F_2 shows the best results.

grid size	Detection			Classification mAP
	AP	AP ₅₀	AP ₇₅	
1	54.6	80.8	60.5	82.2
3	55.6	81.3	61.5	81.6
5	56.1	81.4	62.2	82.6
7	56.7	81.7	63.0	82.9
9	56.7	82.1	63.2	82.9

Table 7 – Ablation study of grid size S . The results increase as the S gets larger. We use grid size being 7 in other experiments, as the performance becomes stable when the S grows beyond 7.

Supervised Contrastive Learning



Supervised Contrastive Learning

$$\mathcal{L}^{self} = \sum_{i \in I} \mathcal{L}_i^{self} = - \sum_{i \in I} \log \frac{\exp(z_i \cdot z_{j(i)}/\tau)}{\sum_{a \in A(i)} \exp(z_i \cdot z_a/\tau)}$$

$$\mathcal{L}_{out}^{sup} = \sum_{i \in I} \mathcal{L}_{out,i}^{sup} = \sum_{i \in I} \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(z_i \cdot z_p/\tau)}{\sum_{a \in A(i)} \exp(z_i \cdot z_a/\tau)}$$

$$\mathcal{L}_{in}^{sup} = \sum_{i \in I} \mathcal{L}_{in,i}^{sup} = \sum_{i \in I} -\log \left\{ \frac{1}{|P(i)|} \sum_{p \in P(i)} \frac{\exp(z_i \cdot z_p/\tau)}{\sum_{a \in A(i)} \exp(z_i \cdot z_a/\tau)} \right\}$$

Supervised Contrastive Learning

$$\mathcal{L}^{self} = \sum_{i \in I} \mathcal{L}_i^{self} = - \sum_{i \in I} \log \frac{\exp(z_i \cdot z_{j(i)}/\tau)}{\sum_{a \in A(i)} \exp(z_i \cdot z_a/\tau)}$$

$$\mathcal{L}_{in}^{sup} \leq \mathcal{L}_{out}^{sup}$$

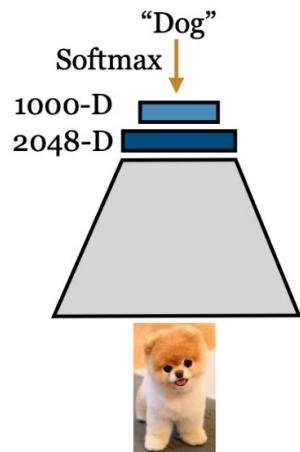
$$\mathcal{L}_{out}^{sup} = \sum_{i \in I} \mathcal{L}_{out,i}^{sup} = \sum_{i \in I} \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(z_i \cdot z_p/\tau)}{\sum_{a \in A(i)} \exp(z_i \cdot z_a/\tau)}$$

$$\mathcal{L}_{in}^{sup} = \sum_{i \in I} \mathcal{L}_{in,i}^{sup} = \sum_{i \in I} -\log \left\{ \frac{1}{|P(i)|} \sum_{p \in P(i)} \frac{\exp(z_i \cdot z_p/\tau)}{\sum_{a \in A(i)} \exp(z_i \cdot z_a/\tau)} \right\}$$

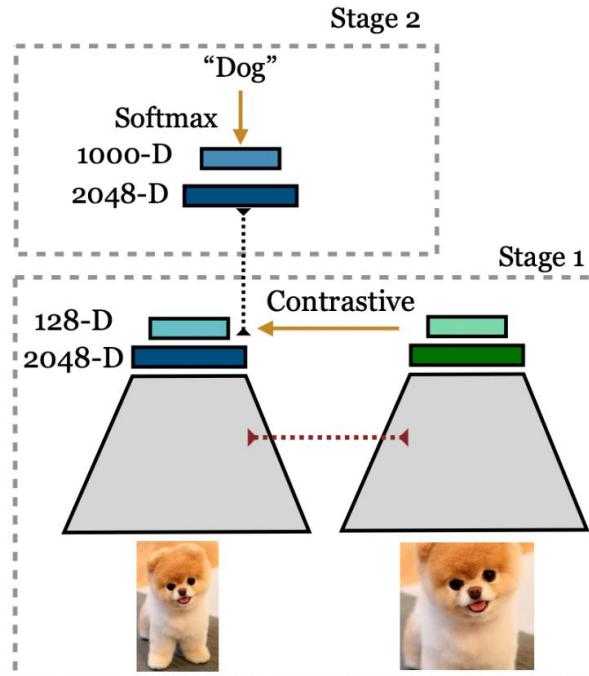
Loss	Top-1
\mathcal{L}_{out}^{sup}	78.7%
\mathcal{L}_{in}^{sup}	67.4%

Supervised Contrastive Learning

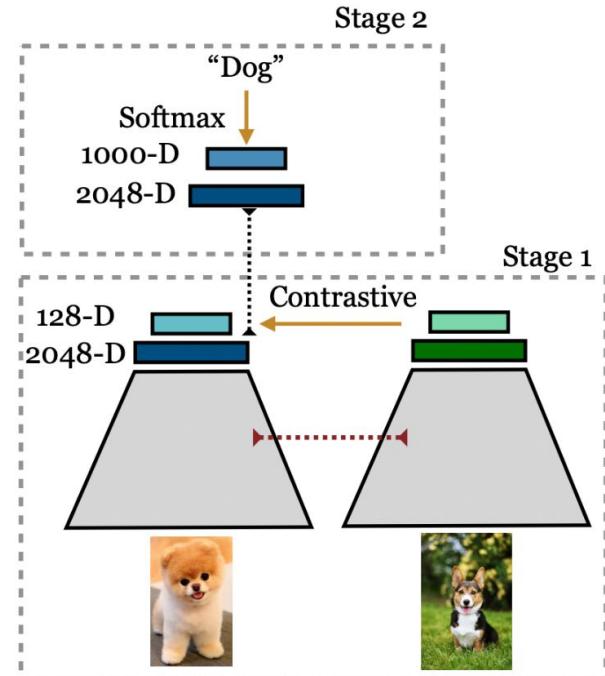
Shared
Weights/Activations
Loss Function



(a) Supervised Cross Entropy



(b) Self Supervised Contrastive



(c) Supervised Contrastive

Supervised Contrastive Learning: results

Loss	Architecture	Augmentation	Top-1	Top-5
Cross-Entropy (baseline)	ResNet-50	MixUp [61]	77.4	93.6
Cross-Entropy (baseline)	ResNet-50	CutMix [60]	78.6	94.1
Cross-Entropy (baseline)	ResNet-50	AutoAugment [5]	78.2	92.9
Cross-Entropy (our impl.)	ResNet-50	AutoAugment [30]	77.6	95.3
SupCon	ResNet-50	AutoAugment [5]	78.7	94.3
Cross-Entropy (baseline)	ResNet-200	AutoAugment [5]	80.6	95.3
Cross-Entropy (our impl.)	ResNet-200	Stacked RandAugment [49]	80.9	95.2
SupCon	ResNet-200	Stacked RandAugment [49]	81.4	95.9
SupCon	ResNet-101	Stacked RandAugment [49]	80.2	94.7

Supervised Contrastive Learning: ablations

Loss	Architecture	rel. mCE	mCE
Cross-Entropy (baselines)	AlexNet [28]	100.0	100.0
	VGG-19+BN [44]	122.9	81.6
	ResNet-18 [17]	103.9	84.7
Cross-Entropy (our implementation)	ResNet-50	96.2	68.6
	ResNet-200	69.1	52.4
Supervised Contrastive	ResNet-50	94.6	67.2
	ResNet-200	66.5	50.6

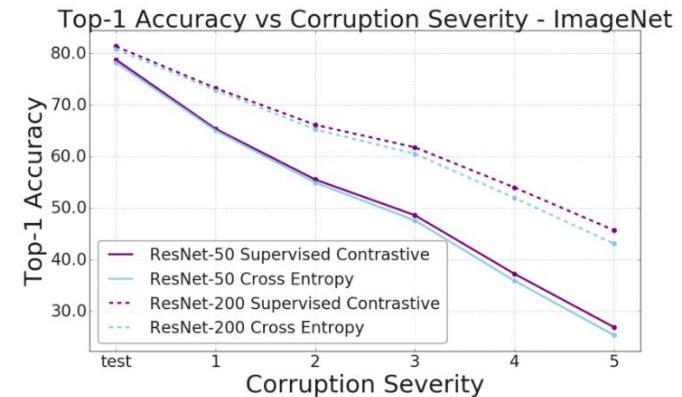
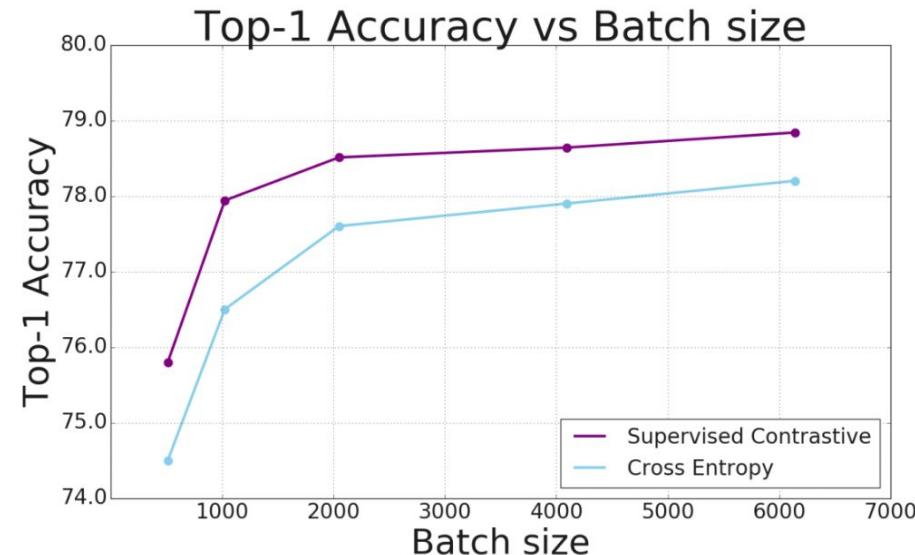
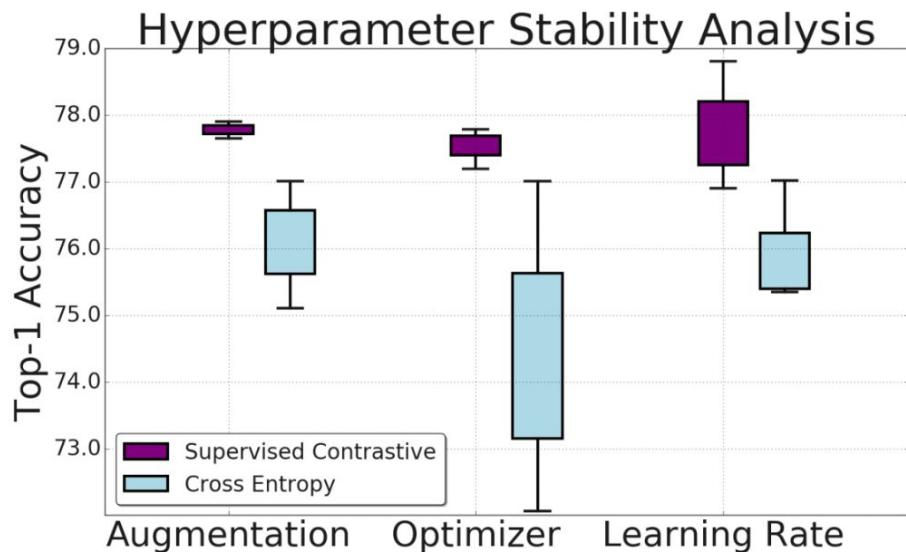


Figure 3: Training with supervised contrastive loss makes models more robust to corruptions in images. **Left:** Robustness as measured by Mean Corruption Error (mCE) and relative mCE over the ImageNet-C dataset [19] (lower is better). **Right:** Mean Accuracy as a function of corruption severity averaged over all various corruptions. (higher is better).

Supervised Contrastive Learning: ablations



Solo-learn library

- Many SSL methods available ✓
- Supports W&B, Nvidia Dali ✓
- Checkpoints available ✓
- MultiGPU training ✓
- PyTorch Lightning ✗



<https://github.com/vturrisi/solo-learn>