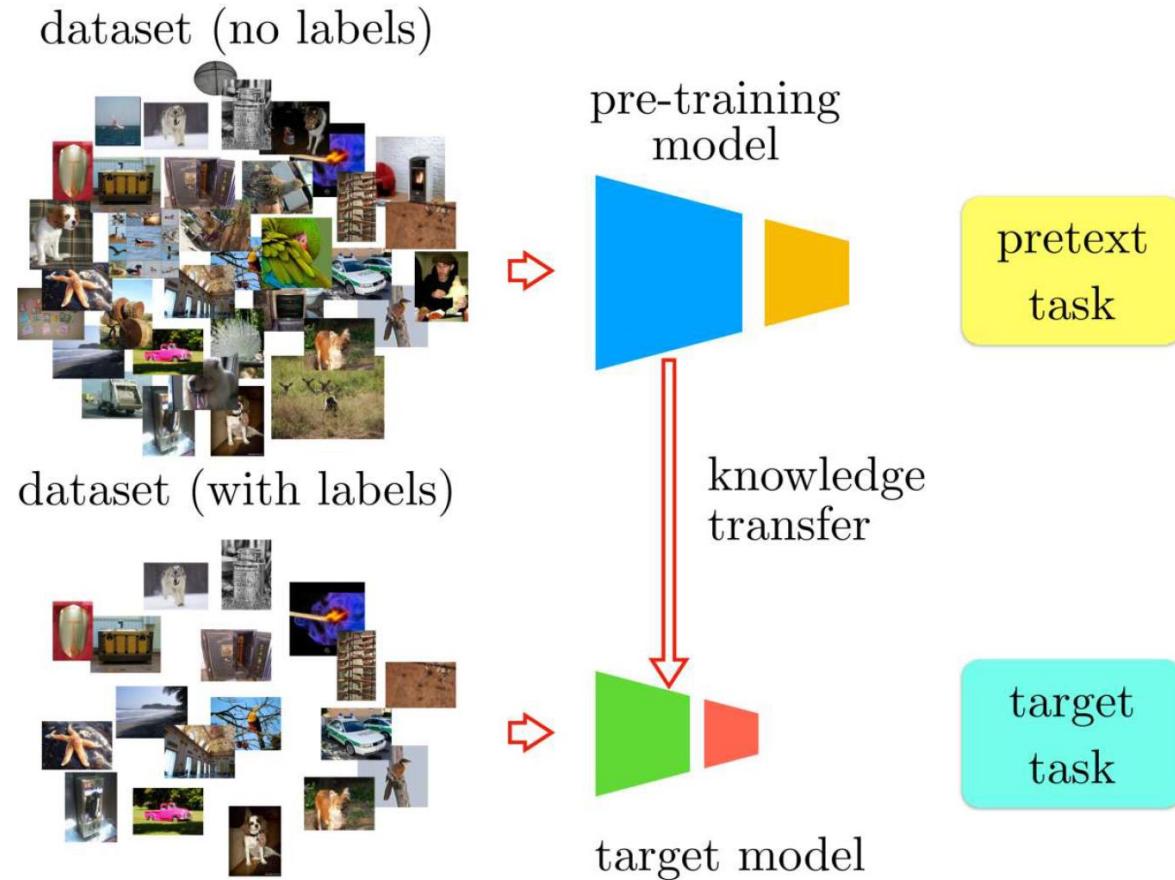


# 自监督视觉学习

Wangmeng Zuo

Centre on Machine Learning Research  
Harbin Institute of Technology

# Self-Supervised Learning

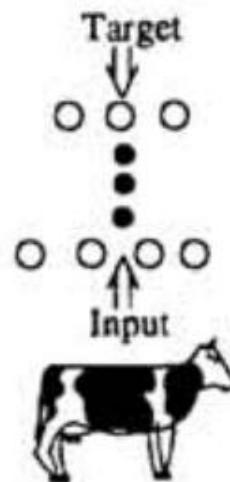


# Self-Supervised Learning in Vision

## Supervised

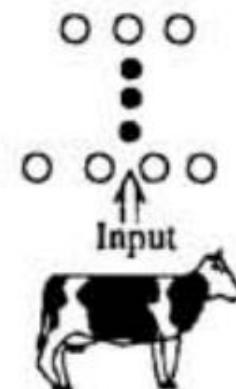
- implausible label

"COW"



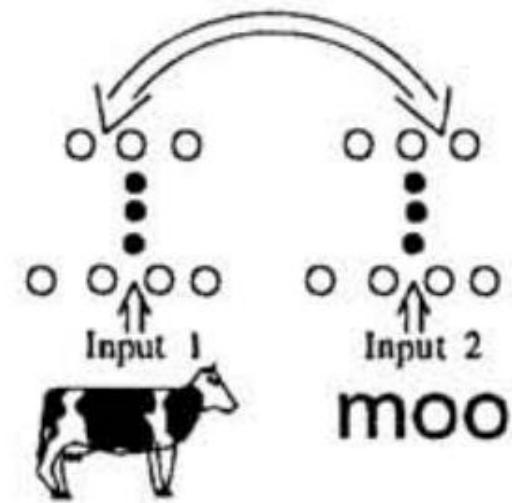
## Unsupervised

- limited power



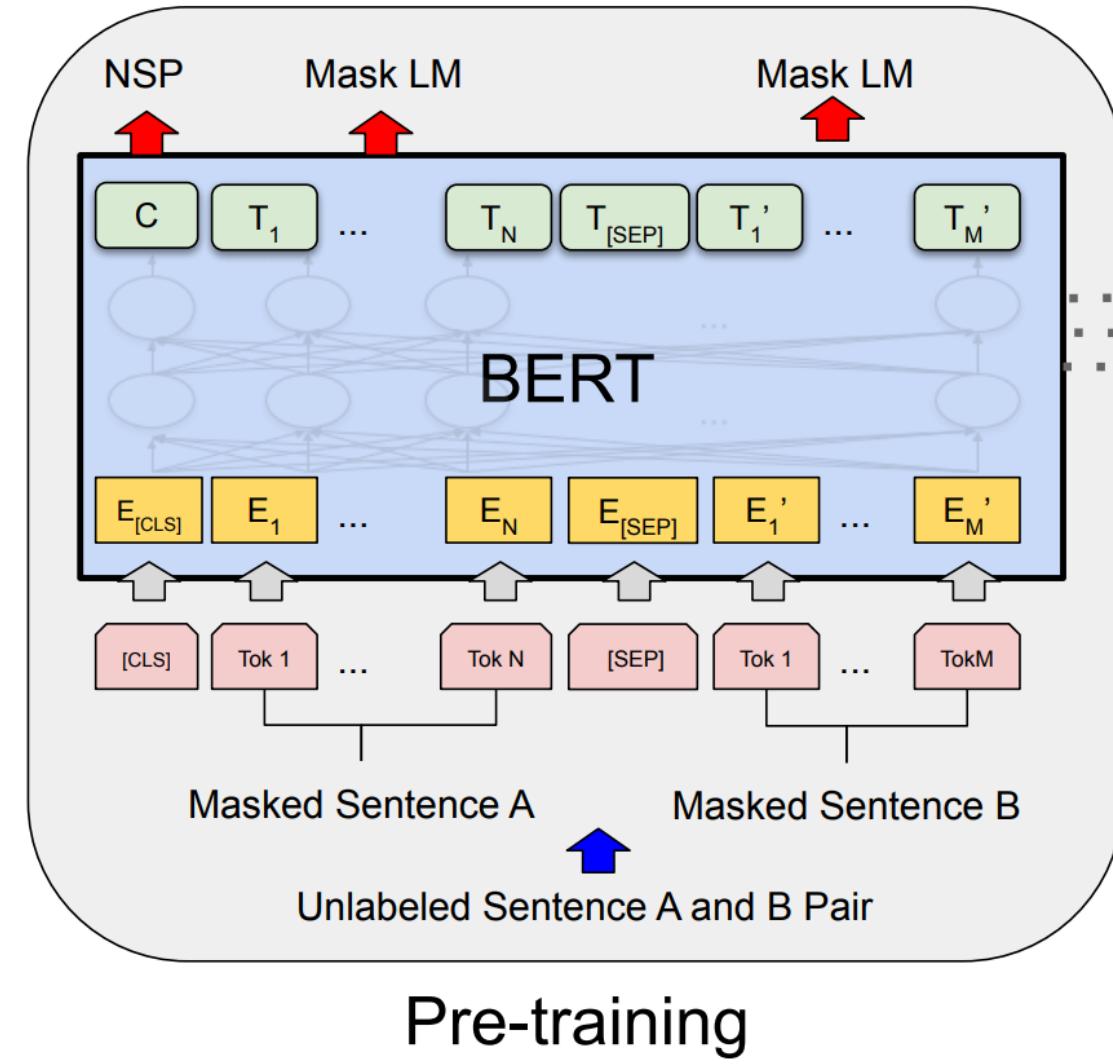
## Self-Supervised

- derives label from a co-occurring input to another modality



# BERT for Language Understanding (Devlin et al., Arxiv2018)

- Citations: 61701 (2018.10 —)
- Randomly masks some of the tokens from the input, and the objective is to predict the original vocabulary id of the masked word based only on its context.



# Content

- Self-supervised Vision Learning
  - 1. Self-supervised Contextual Modeling
  - 2. Contrastive learning-based self-supervised learning
  - 3. Returning of Self-supervised Contextual Modeling

# Deep Context Modeling

- Jigsaw Puzzles

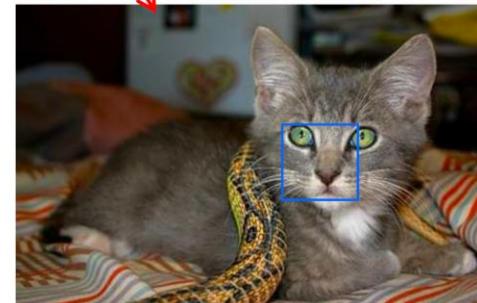


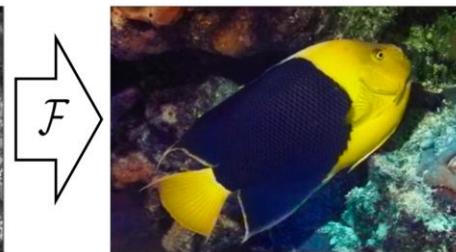
Image Rotation



- Image Rotation



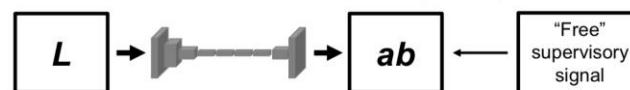
Grayscale image:  $L$  channel



Concatenate ( $L, ab$ )

$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

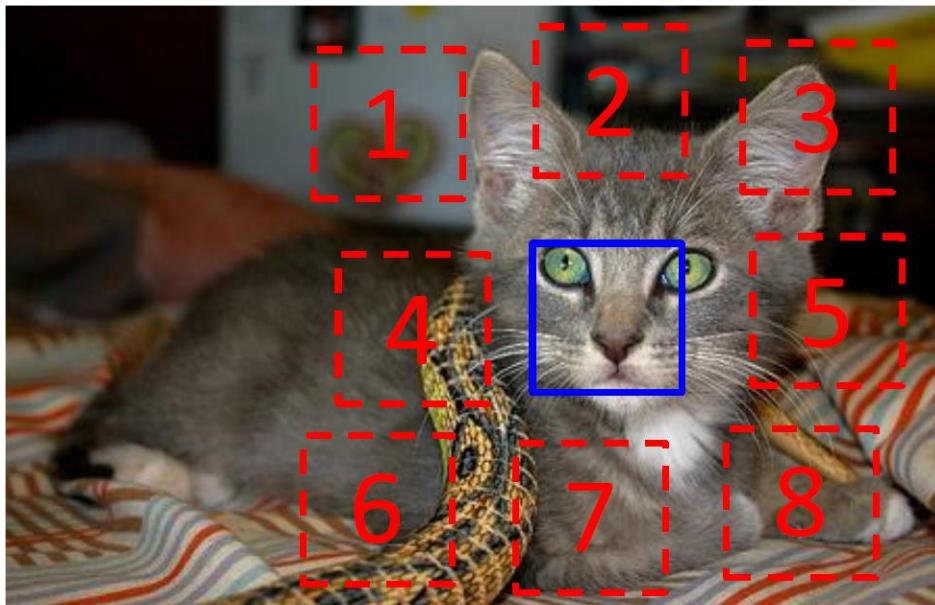
$$(\mathbf{X}, \hat{\mathbf{Y}})$$



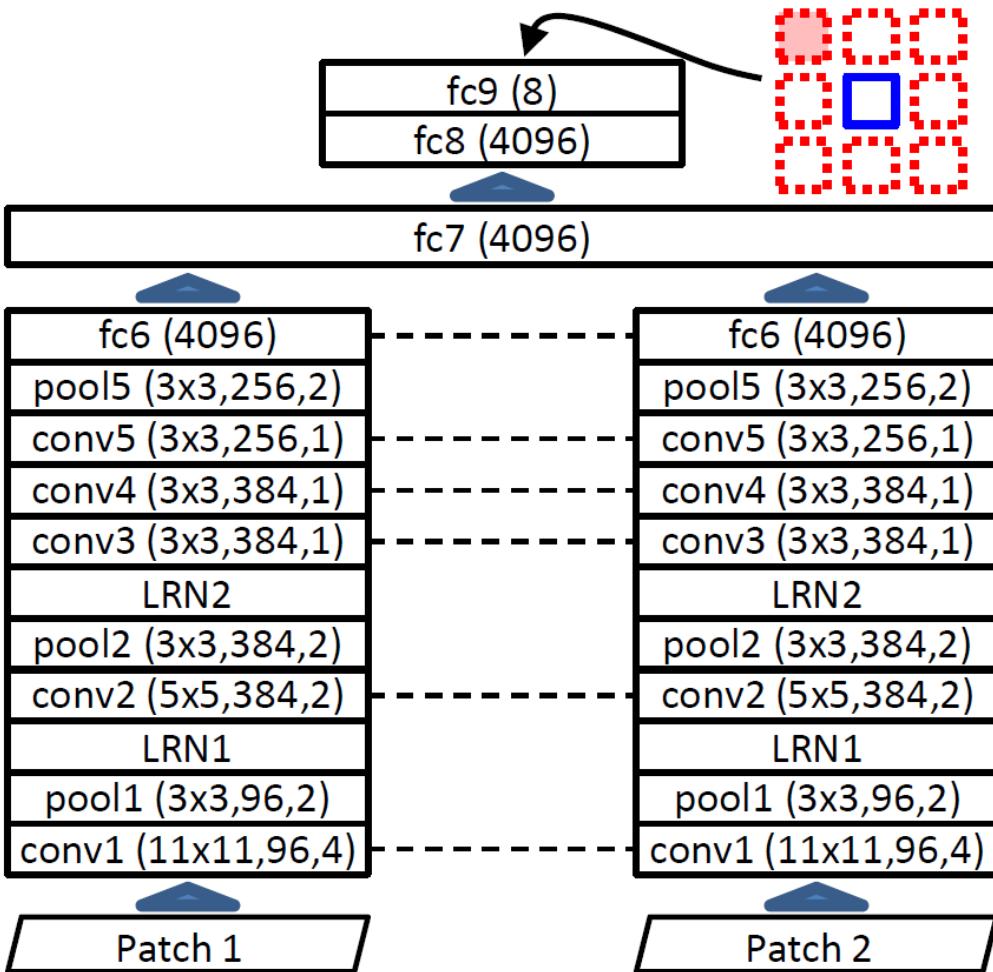
- Image Colorization

- Image Inpainting

# Jigsaw Puzzles (1)



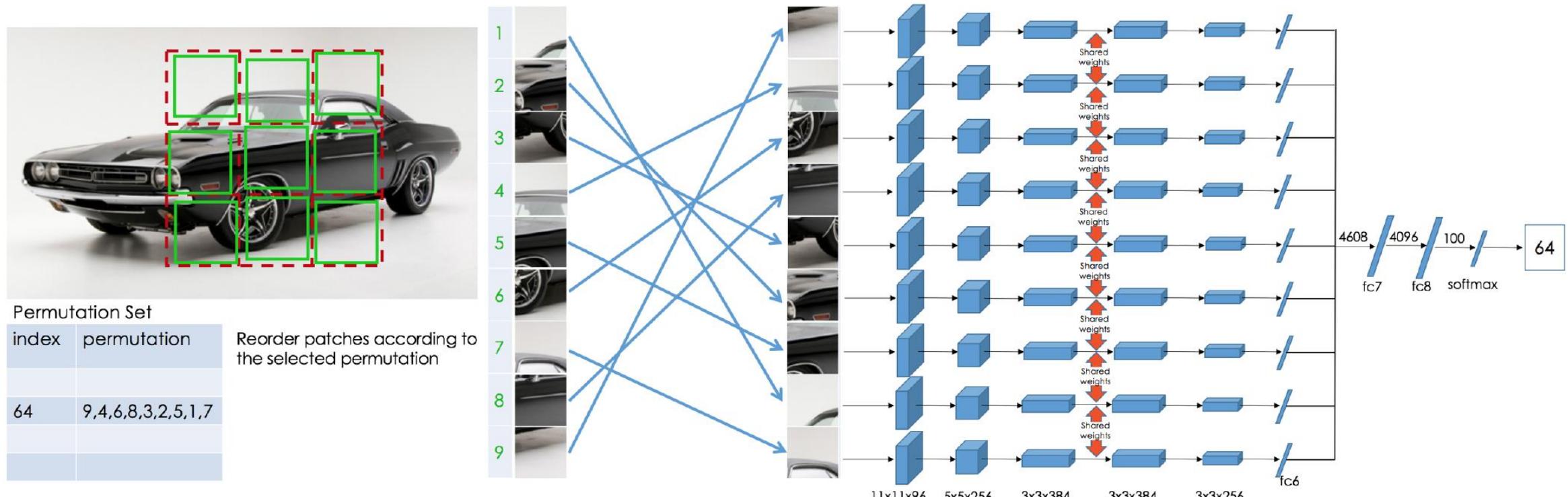
$$X = (\text{Patch 1}, \text{Patch 2}); Y = 3$$



C. Doersch, A. Gupta, and A. A. Efros. Unsupervised Visual Representation Learning by Context Prediction. In ICCV 2015.

# Jigsaw Puzzles (2)

- Permutation set: 1000, Average Hamming distance: 8.00



M. Noroozi and P. Favaro, Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles, ECCV 2016.

# Jigsaw++



(a)



(b)



(c)



(d)

Method	Ref	Class.	Det.		Segm.
			SS	MS	
Supervised [17]		79.9	59.1	59.8	48.0
CC+HOG [6]		70.2	53.2	53.5	39.2
Random	[27]	53.3	43.4	-	19.8
ego-motion [1]	[1]	54.2	43.9	-	-
BiGAN [9]	[9]	58.6	46.2	-	34.9
ContextEncoder [27]	[27]	56.5	44.5	-	29.7
Video [32]	[16]	63.1	47.2	-	-
Colorization [39]	[39]	65.9	46.9	-	35.6
Split-Brain [40]	[40]	67.1	46.7	-	36.0
Context [7]	[16]	55.3	46.6	-	-
Context [7]*	[16]	65.3	51.1	-	-
Counting [23]	[23]	67.7	51.4	-	36.6
WatchingObjectsMove [26]	[26]	61.0	-	52.2	-
Jigsaw [21]	[21]	67.7	53.2	-	-
Jigsaw++		69.8	55.5	55.7	38.1
CC+Context-ColorDrop [7]		67.9	52.8	53.4	-
CC+Context-ColorProjection [7]		66.7	51.5	51.8	-
CC+Jigsaw++		69.9	55.0	55.8	40.0
[37]+vgg-Jigsaw++		70.6	54.8	55.2	38.0
CC+vgg-Context [7]		68.0	53.0	53.5	-
CC+vgg-Jigsaw++		<b>72.5</b>	<b>56.5</b>	<b>57.2</b>	<b>42.6</b>

VOC 2007

# Deep Context Modeling

- Jigsaw Puzzles

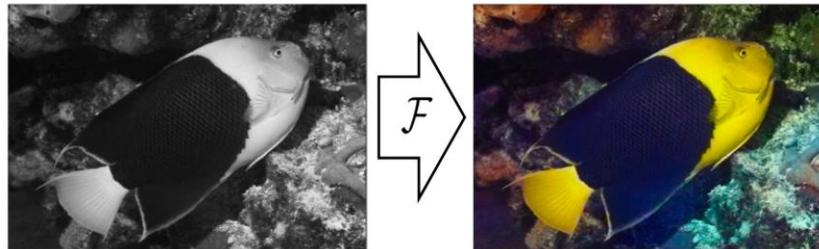


- Image Rotation



Image Rotation

- Image Colorization



Grayscale image:  $L$  channel

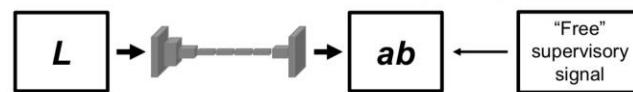
$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

Concatenate ( $L, ab$ )

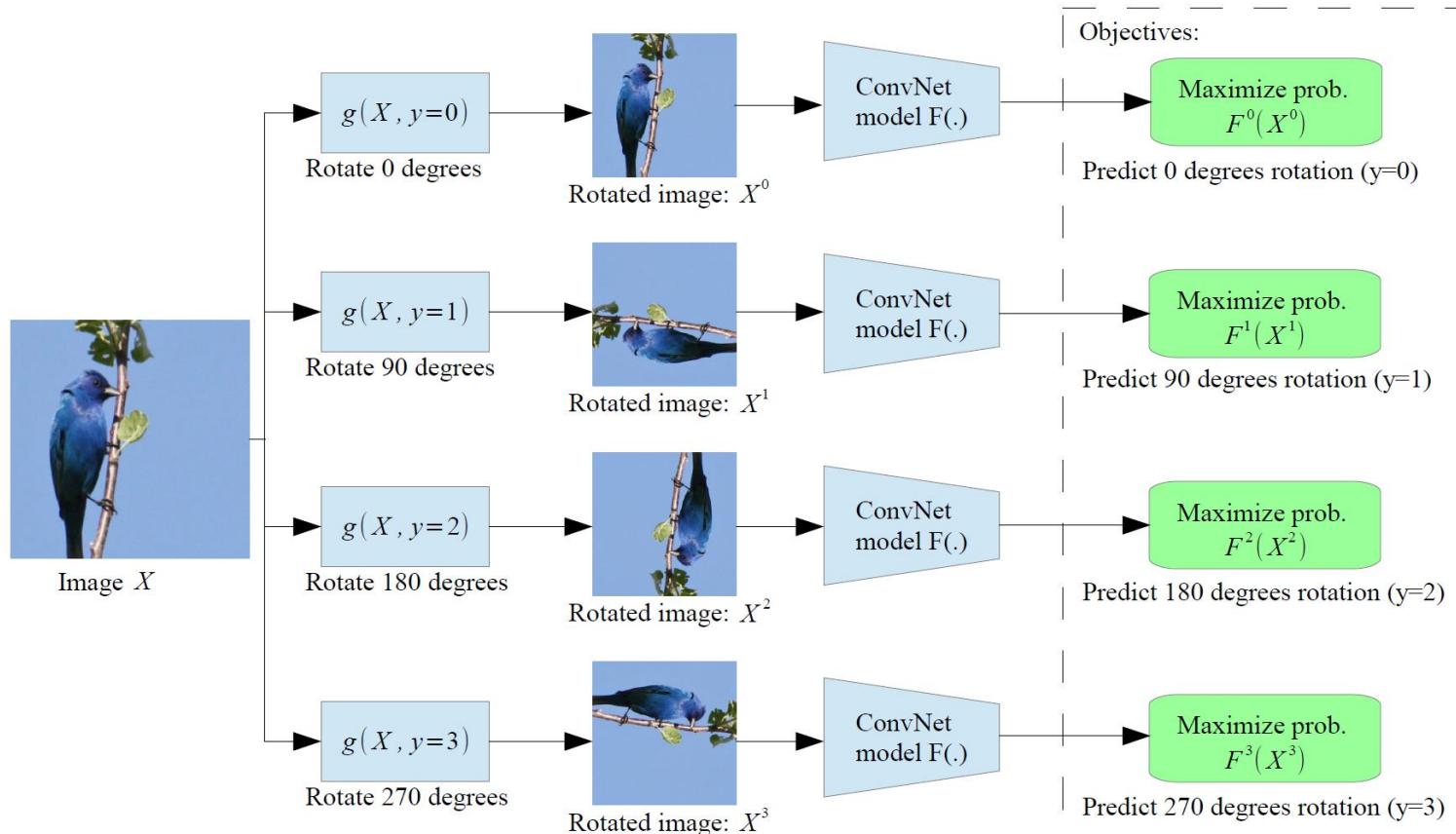
$$(\mathbf{X}, \hat{\mathbf{Y}})$$



- Image Inpainting



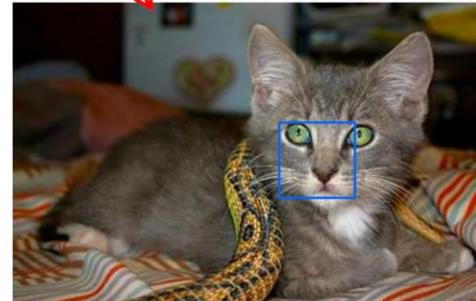
# Predicting Image Rotations



Spyros Gidaris, Praveer Singh, Nikos Komodakis, Unsupervised Representation Learning by Predicting Image Rotations, ICLR 2018.

# Deep Context Modeling

- Jigsaw Puzzles



Randomly Sample Patch  
Sample Second Patch

- Image Rotation

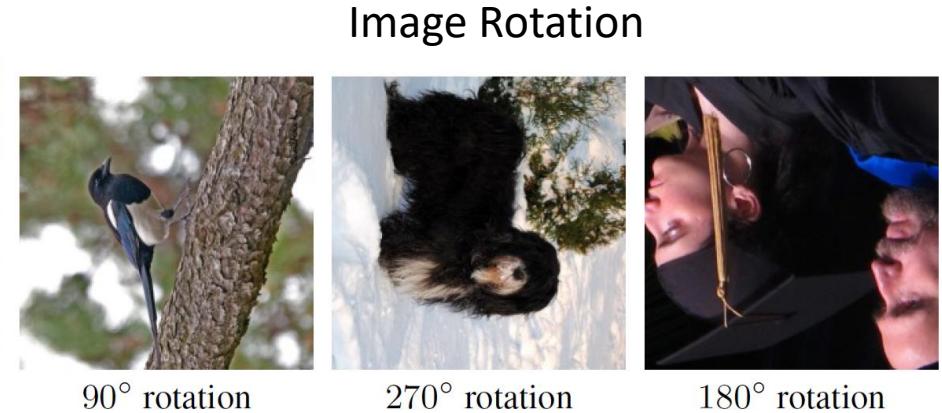
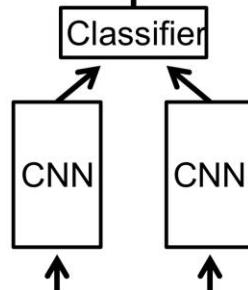
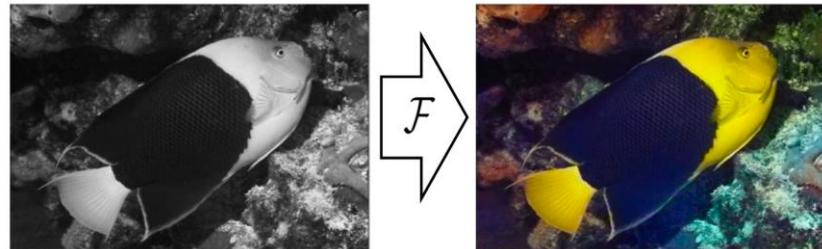


Image Rotation

- Image Colorization

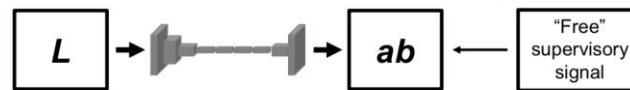


Grayscale image:  $L$  channel

$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

Concatenate ( $L, ab$ )

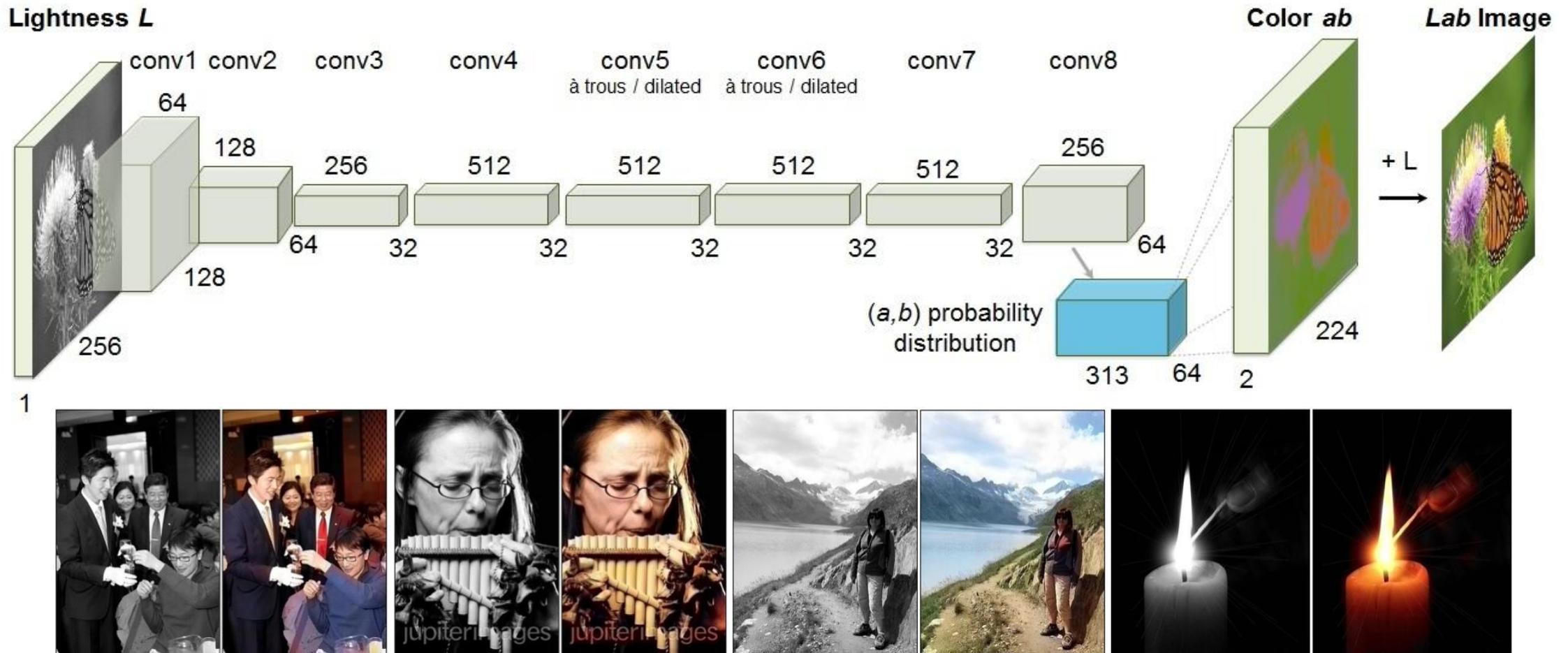
$$(\mathbf{X}, \hat{\mathbf{Y}})$$



- Image Inpainting

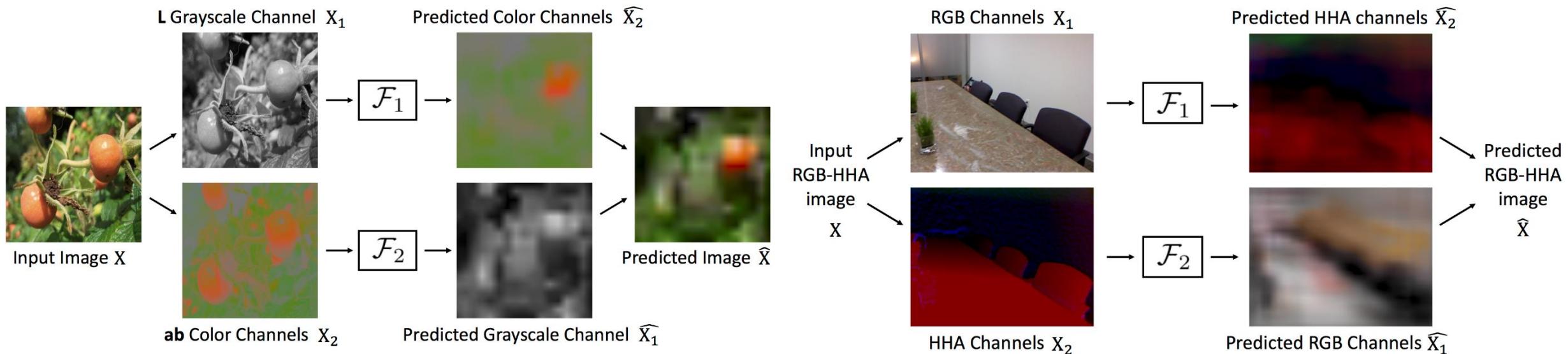


# Image Colorization



R. Zhang, P. Isola, and A.A. Efros. Colorful image colorization, ECCV 2016

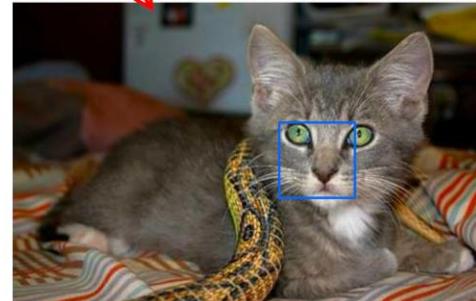
# Split-Brain Autoencoders



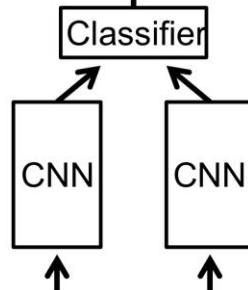
R. Zhang, P. Isola, A. A. Efros, Split-Brain Autoencoders: Unsupervised Learning by Cross-Channel Prediction, CVPR 2017.

# Deep Context Modeling

- Jigsaw Puzzles



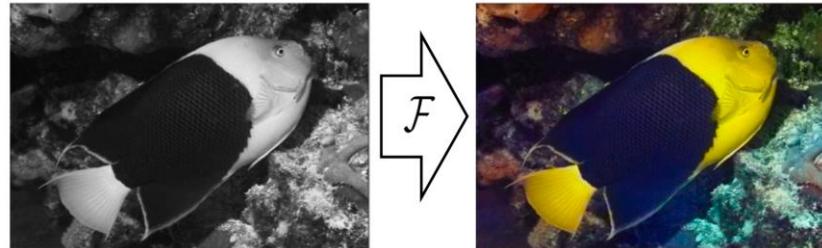
- Image Rotation



Randomly Sample Patch  
Sample Second Patch



- Image Colorization

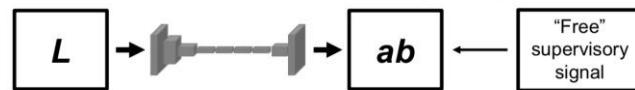


Grayscale image:  $L$  channel

$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

Concatenate ( $L, ab$ )

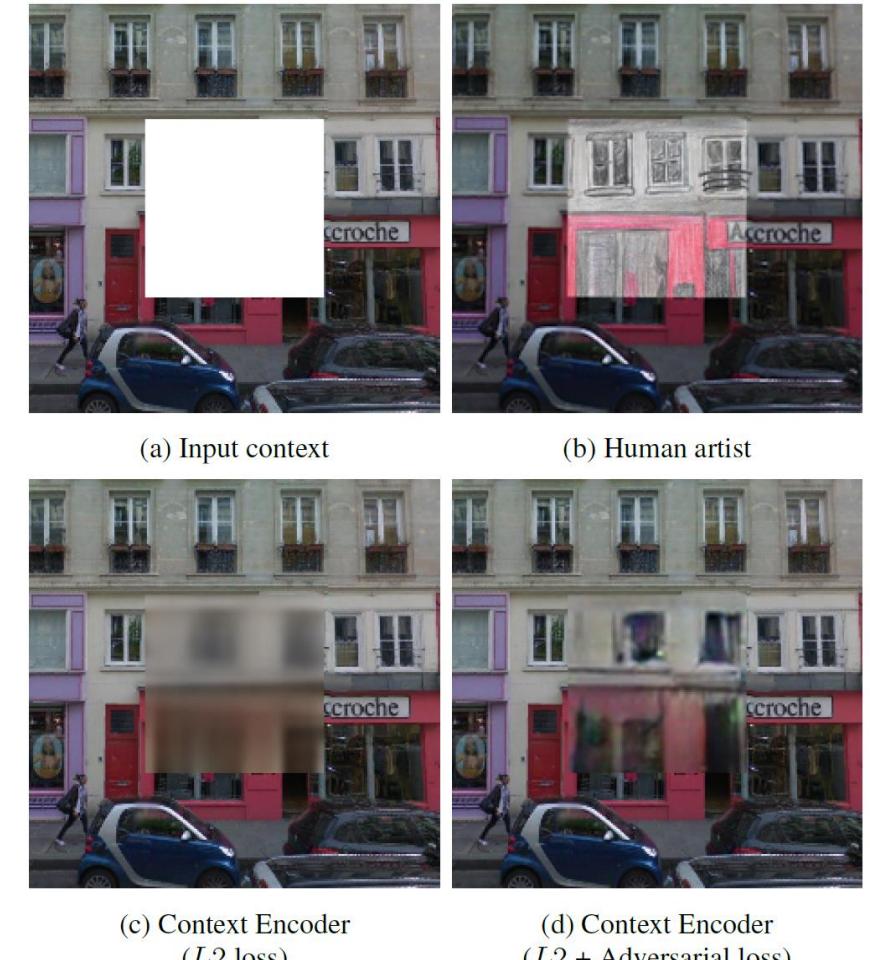
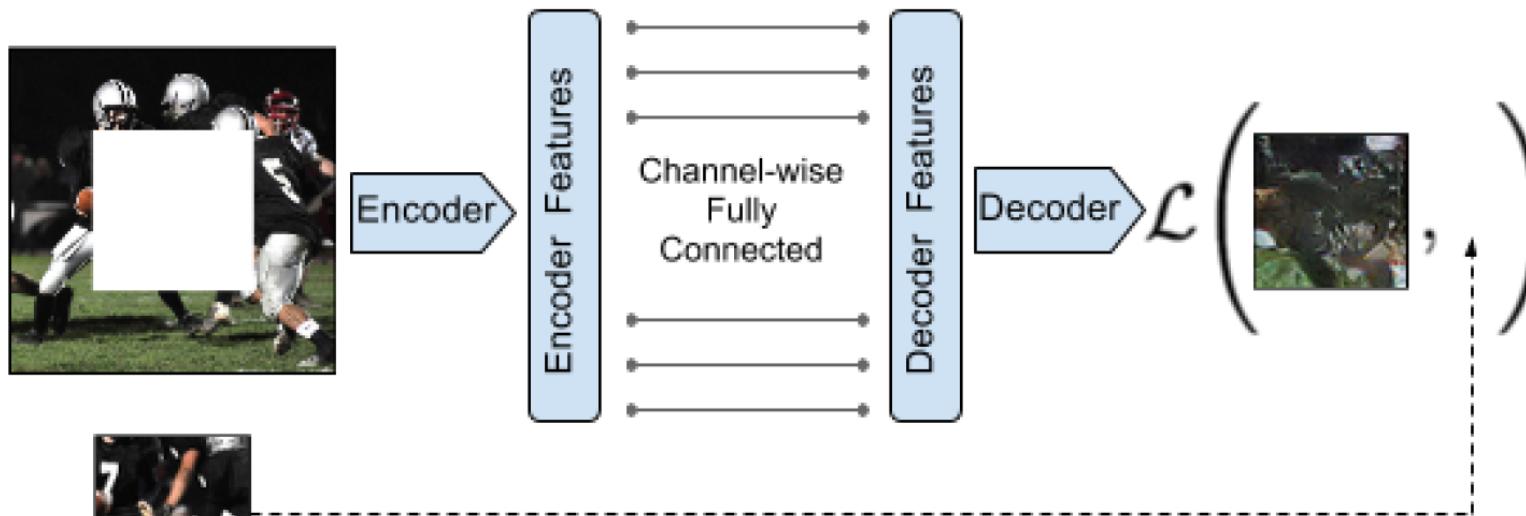
$$(\mathbf{X}, \hat{\mathbf{Y}})$$



- Image Inpainting

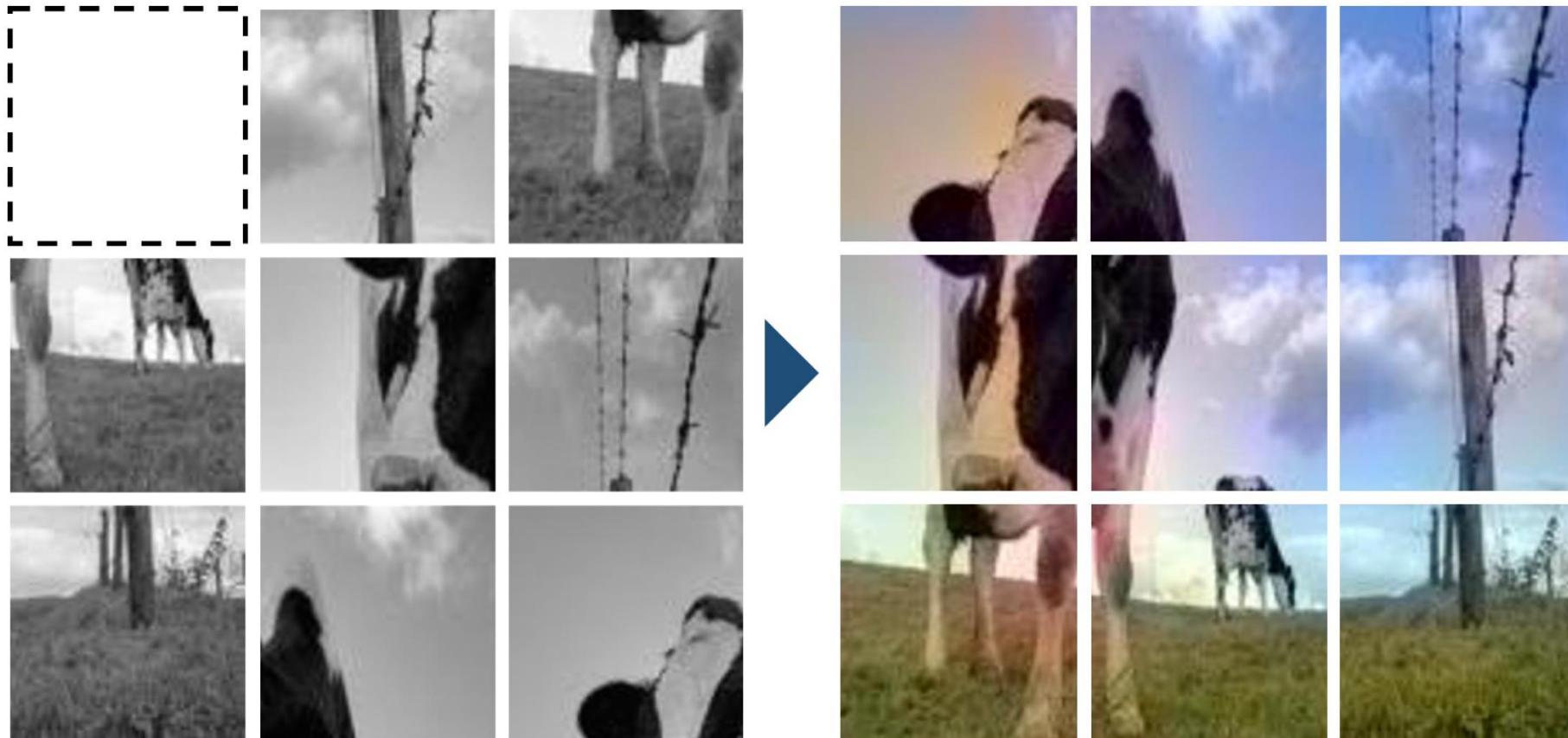


# Context Encoders



D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, A. A. Efros, Context Encoders: Feature Learning by Inpainting, CVPR 2016

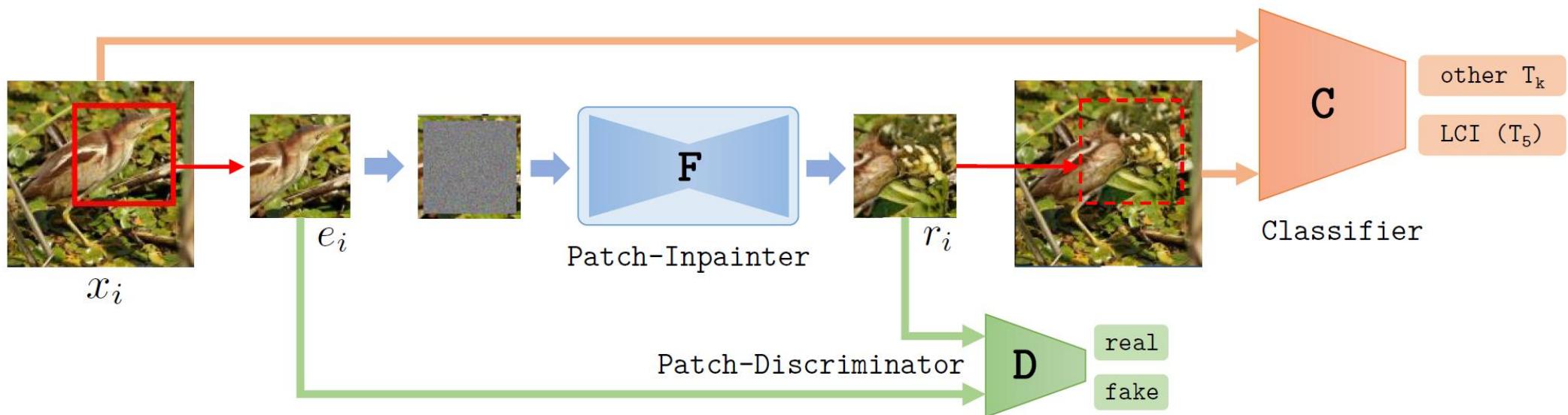
# Jigsaw+Colorization+Inpainting



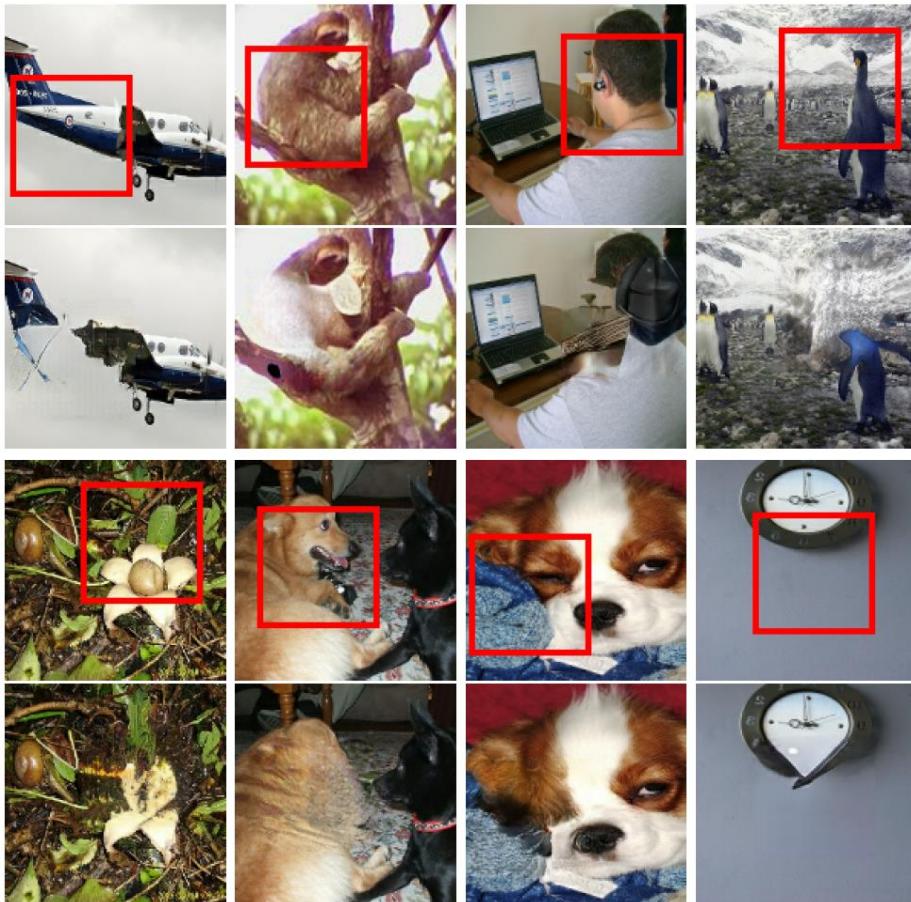
Learning Image Representations by Completing Damaged Jigsaw Puzzles, WACV 2018.

# Inpainting+Rotation

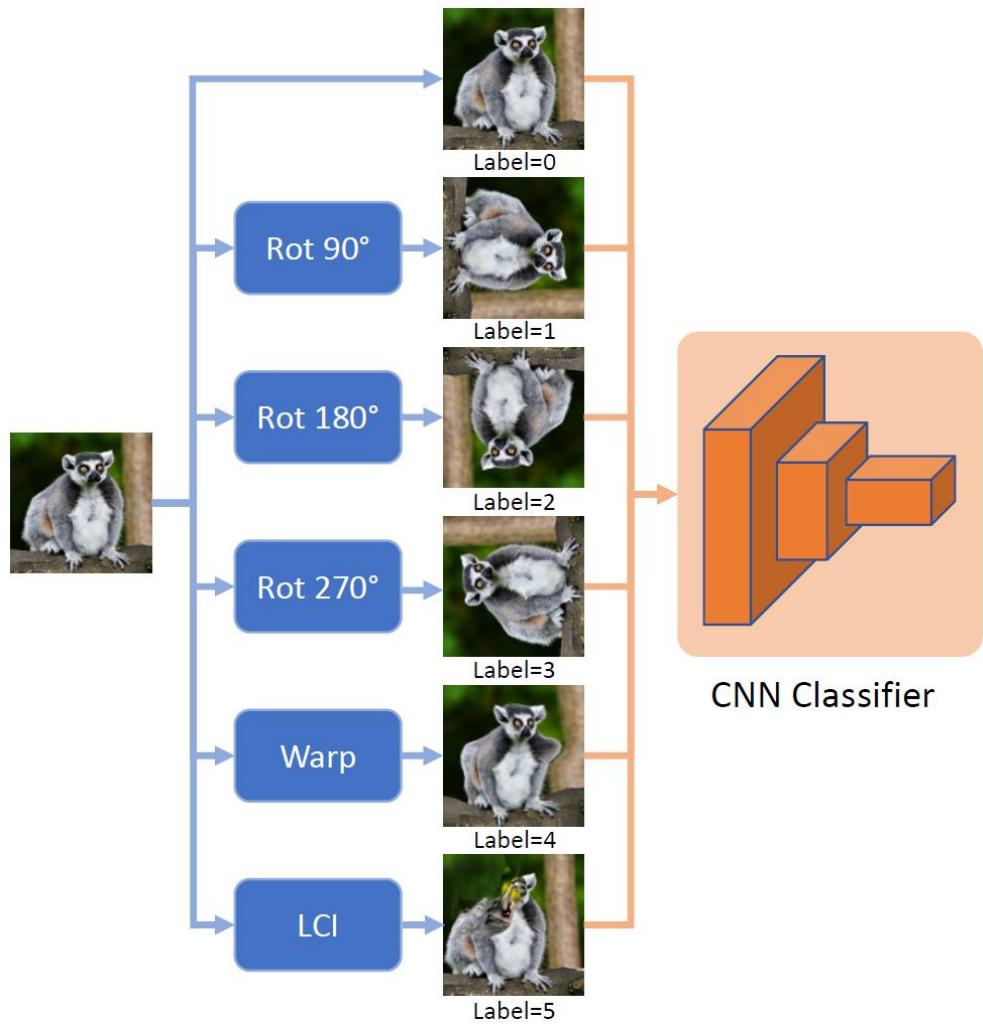
- Limited context inpainting



Simon Jenni, Hailin Jin, Paolo Favaro , Steering Self-Supervised Feature Learning Beyond Local Pixel Statistics, CVPR 2020.



# Learning Global Statistics



Model\Layer	conv1	conv2	conv3	conv4	conv5
ImageNet Labels	19.3%	36.3%	44.2%	48.3%	50.5%
Random	11.6%	17.1%	16.9%	16.3%	14.1%
Donahue <i>et al.</i> [10]	17.7%	24.5%	31.0%	29.9%	28.0%
Feng <i>et al.</i> [15]	19.3%	<u>33.3%</u>	<b>40.8%</b>	<u>41.8%</u>	<b>44.3%</b>
Gidaris <i>et al.</i> [18]	18.8%	31.7%	38.7%	38.2%	36.5%
Huang <i>et al.</i> [25]	15.6%	27.0%	35.9%	39.7%	37.9%
Jenni & Favaro [29]	<u>19.5%</u>	33.3%	37.9%	38.9%	34.9%
Noroozi & Favaro [46]	18.2%	28.8%	34.0%	33.9%	27.1%
Noroozi <i>et al.</i> [47]	18.0%	30.6%	34.3%	32.5%	25.7%
Noroozi <i>et al.</i> [48]	19.2%	32.0%	37.3%	37.1%	34.6%
Tian <i>et al.</i> [57]	18.4%	33.5%	38.1%	40.4%	<u>42.6%</u>
Wu <i>et al.</i> [61]	16.8%	26.5%	31.8%	34.1%	35.6%
Zhang <i>et al.</i> [65]	13.1%	24.8%	31.0%	32.6%	31.8%
Zhang <i>et al.</i> [66]	17.7%	29.3%	35.4%	35.2%	32.8%
Zhang <i>et al.</i> [64]	19.2%	32.8%	<u>40.6%</u>	39.7%	37.7%
Doersch <i>et al.</i> [8]*	16.2%	23.3%	30.2%	31.7%	29.6%
Caron <i>et al.</i> [4]*	12.9%	29.2%	38.2%	39.8%	36.1%
Zhuang <i>et al.</i> [71]*†	18.7%	32.7%	38.1%	42.3%	42.4%
Ours	<b>20.8%</b>	<b>34.5%</b>	40.2%	<b>43.1%</b>	41.4%
Ours†	22.0%	36.4%	42.4%	45.4%	44.4%

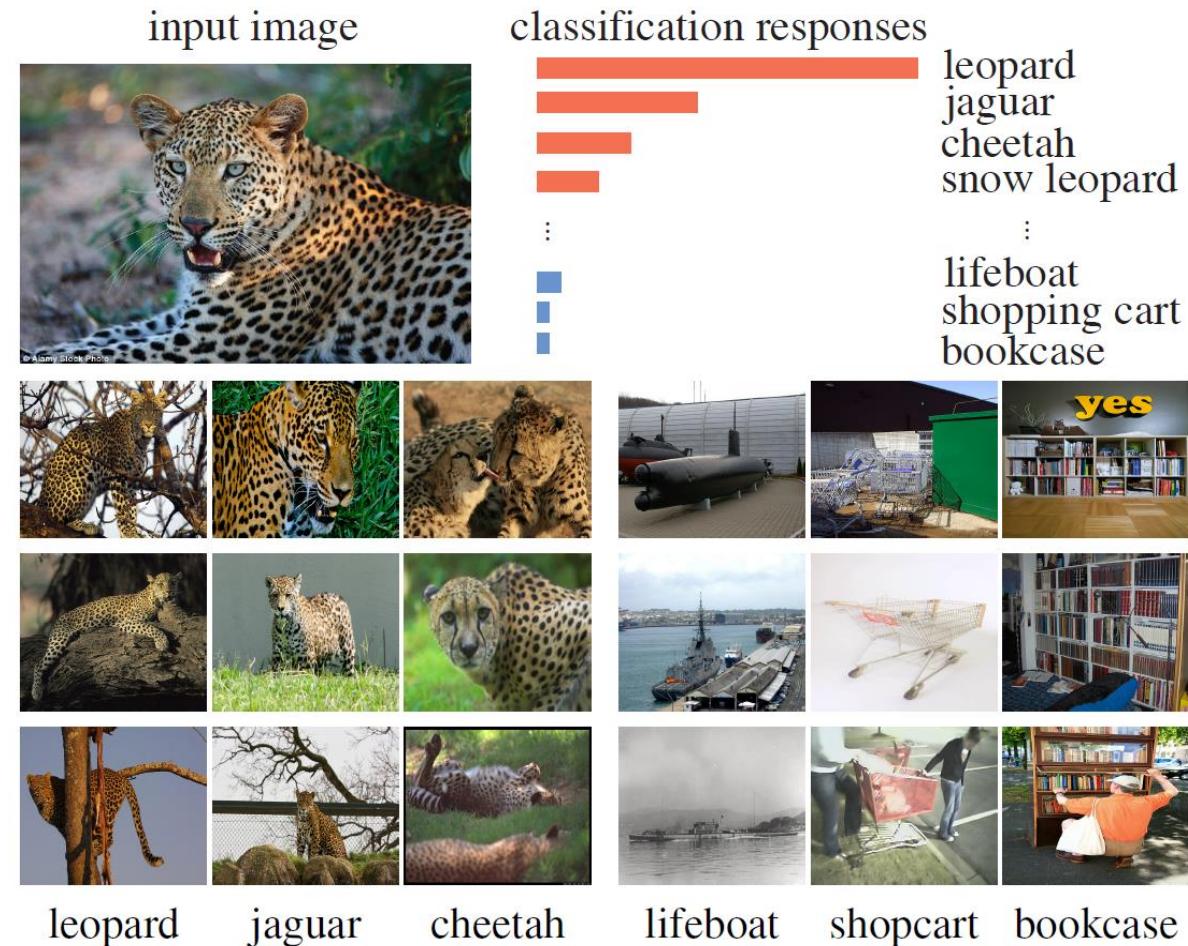
ImageNet (AlexNet)

# Content

- Self-supervised Vision Learning
  - 1. Self-supervised Contextual Modeling
  - 2. Contrastive learning-based self-supervised learning
  - 3. Returning of Self-supervised Contextual Modeling

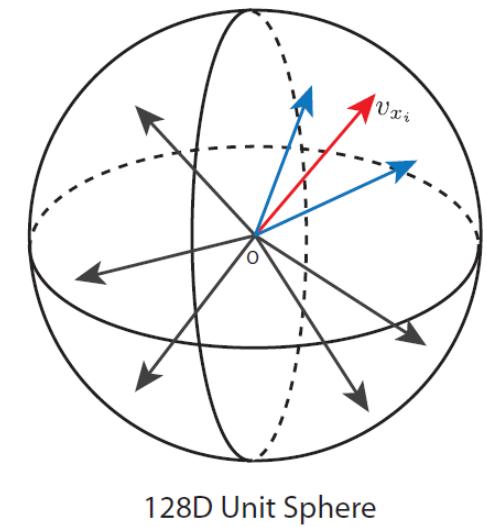
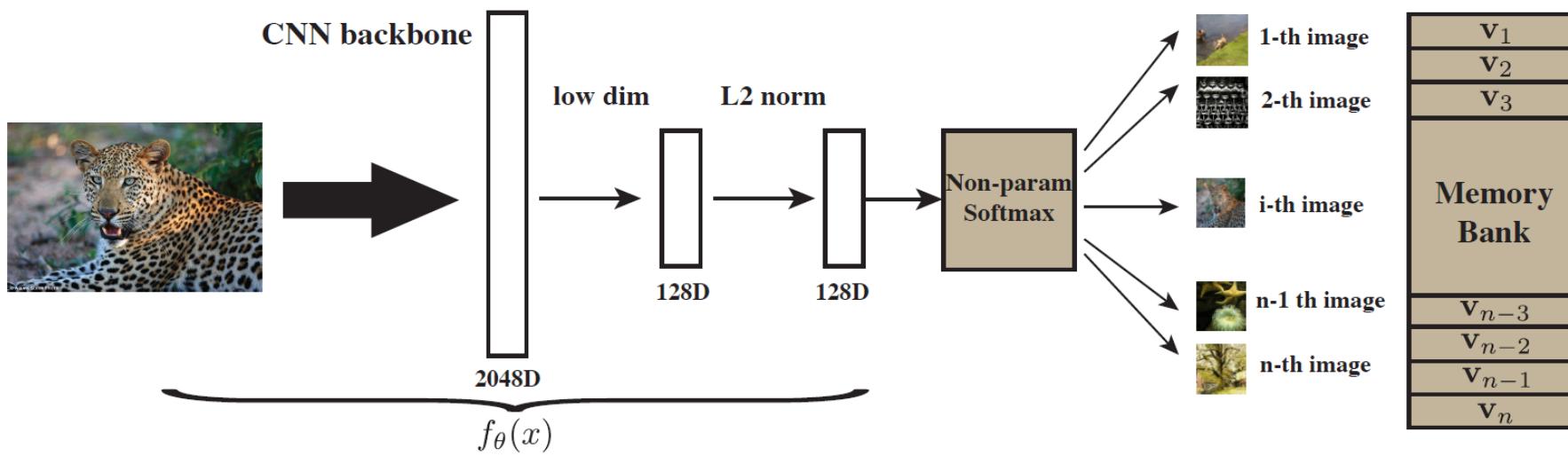
# Instance Discrimination

- Takes the class-wise supervision to the extreme
- Learns a feature representation that discriminates among individual instances



Zhirong Wu, Yuanjun Xiong, Stella Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In CVPR, 2018.

# Instance Discrimination



- Memory bank  $V = \{\mathbf{v}_j\}$
- Forward feature  $\mathbf{f}_i = f_{\theta}(x_i)$
- $\mathbf{f}_i \rightarrow \mathbf{v}_i$

$$\text{Memory Bank} \quad \mathbf{v}_j^T \mathbf{v} \quad \|\mathbf{v}\| = 1$$

# Noise Contrastive Learning

- Softmax classifier

$$P(i|\mathbf{v}) = \frac{\exp(\mathbf{v}^T \mathbf{f}_i / \tau)}{Z_i}$$

$$Z_i = \sum_{j=1}^n \exp(\mathbf{v}_j^T \mathbf{f}_i / \tau)$$

- Posterior probability

$$P_n = 1/n \quad h(i, \mathbf{v}) := P(D = 1 | i, \mathbf{v}) = \frac{P(i|\mathbf{v})}{P(i|\mathbf{v}) + mP_n(i)}$$

- Monte Carlo approximation

$$Z \simeq Z_i \simeq n E_j [\exp(\mathbf{v}_j^T \mathbf{f}_i / \tau)] = \frac{n}{m} \sum_{k=1}^m \exp(\mathbf{v}_{j_k}^T \mathbf{f}_i / \tau)$$

- Negative log-posterior

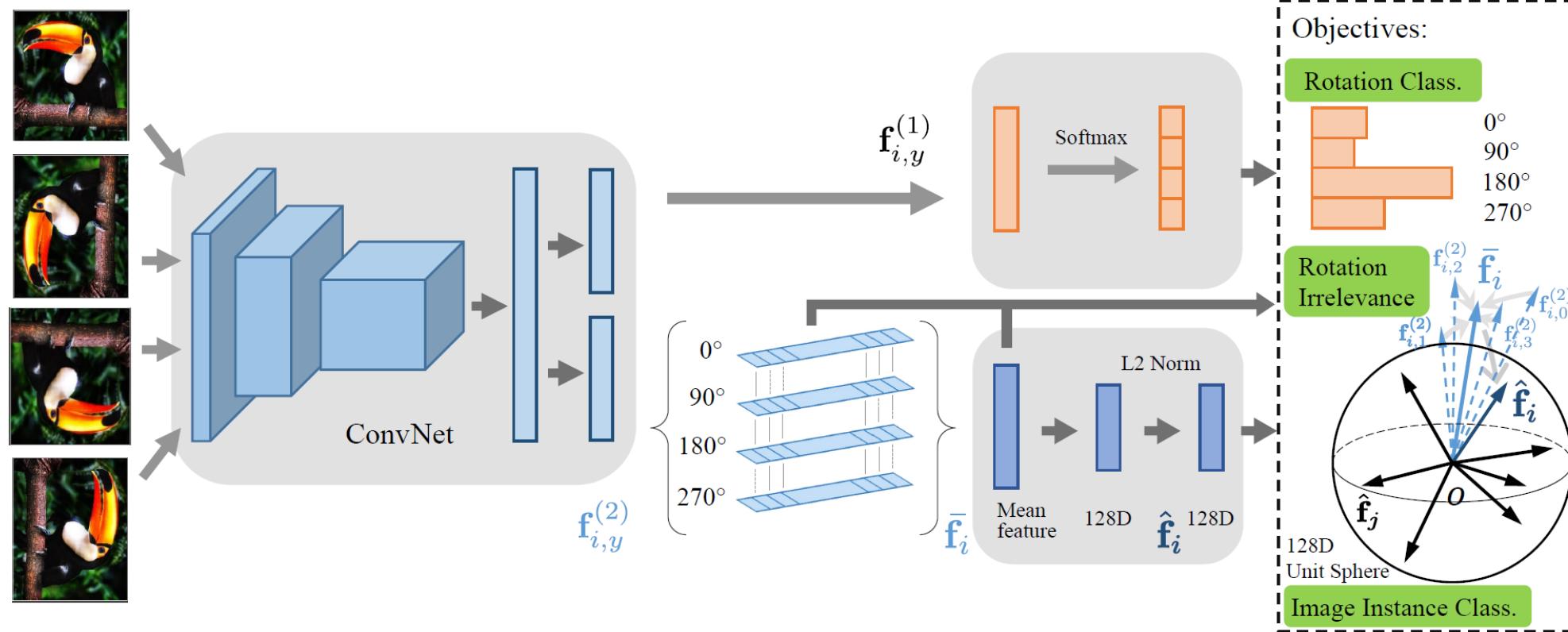
$$\begin{aligned} J_{NCE}(\boldsymbol{\theta}) &= -E_{P_d} \left[ \log h(i, \mathbf{v}_i^{(t-1)}) - \lambda \|\mathbf{v}_i^{(t)} - \mathbf{v}_i^{(t-1)}\|_2^2 \right] \\ &\quad - m \cdot E_{P_n} \left[ \log(1 - h(i, \mathbf{v}'^{(t-1)})) \right]. \end{aligned} \tag{10}$$

Training / Testing	Linear SVM	Nearest Neighbor
Param Softmax	60.3	63.0
Non-Param Softmax	75.4	<b>80.8</b>
NCE $m = 1$	44.3	42.5
NCE $m = 10$	60.2	63.4
NCE $m = 512$	64.3	78.4
NCE $m = 4096$	70.2	<b>80.4</b>

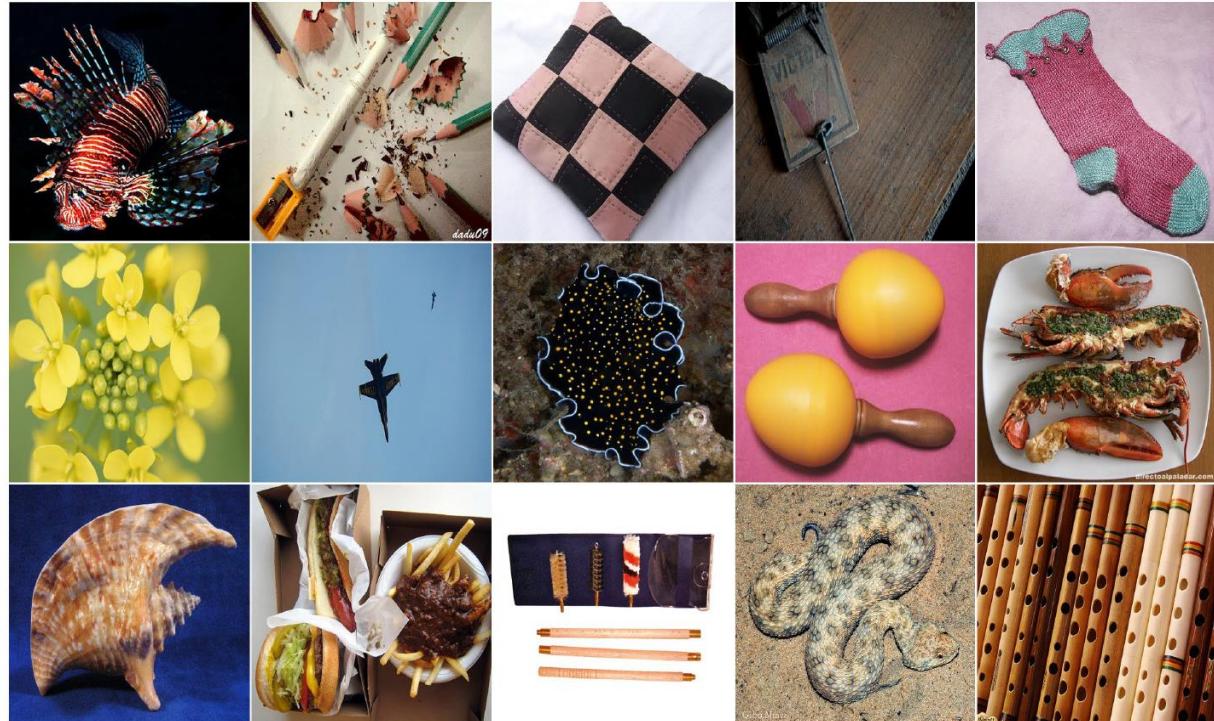
Image Classification Accuracy on ImageNet								
method	conv1	conv2	conv3	conv4	conv5	kNN	#dim	
Random	11.6	17.1	16.9	16.3	14.1	3.5	10K	
Data-Init [16]	17.5	23.0	24.5	23.2	20.6	-	10K	
Context [2]	16.2	23.3	30.2	31.7	29.6	-	10K	
Adversarial [4]	17.7	24.5	31.0	29.9	28.0	-	10K	
Color [47]	13.1	24.8	31.0	32.6	31.8	-	10K	
Jigsaw [27]	19.2	30.1	34.7	33.9	28.3	-	10K	
Count [28]	18.0	30.6	34.3	32.5	25.7	-	10K	
SplitBrain [48]	17.7	29.3	35.4	35.2	32.8	11.8	10K	
Exemplar[3]			31.5			-	4.5K	
Ours Alexnet	16.8	26.5	31.8	34.1	<b>35.6</b>	31.3	128	
Ours VGG16	16.5	21.4	27.6	35.1	<b>39.2</b>	33.9	128	
Ours Resnet18	16.0	19.9	29.8	39.0	<b>44.5</b>	<b>41.0</b>	128	
Ours Resnet50	15.3	18.8	24.9	40.6	<b>54.0</b>	<b>46.5</b>	128	

Table 2: Top-1 classification accuracy on ImageNet.

# Rotation + Instance Discrimination



# Rotation + Instance Discrimination

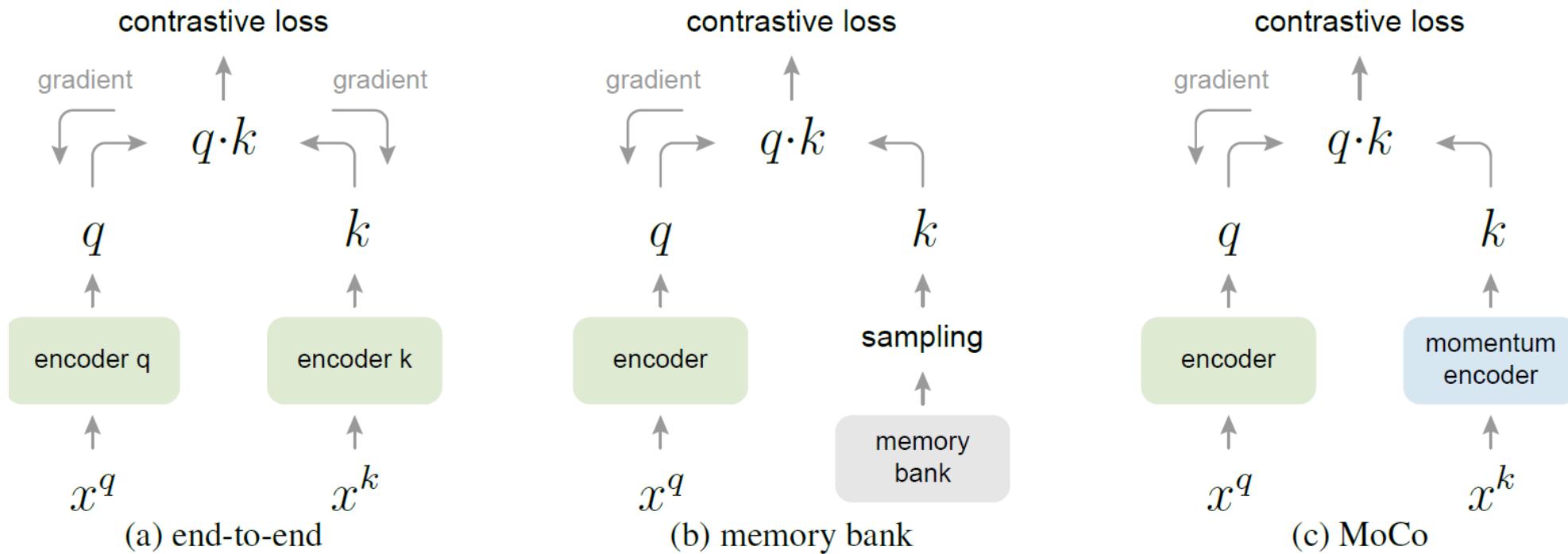


Rotation agnostic

Method \ Layer	conv1	conv2	conv3	conv4	conv5
ImageNet-labels [25, 52]	19.3	36.3	44.2	48.3	50.5
Random [53]	11.6	17.1	16.9	16.3	14.1
Krähenbühl <i>et al.</i> [22]	17.5	23.0	24.5	23.2	20.6
Pathak <i>et al.</i> (Inpainting) [43]	14.1	20.7	21.0	19.8	15.5
Noroozi & Favaro (Jigsaw) [37]	18.2	28.8	34.0	33.9	27.1
Zhang <i>et al.</i> (Colorization) [52]	13.1	24.8	31.0	32.6	31.8
Donahue <i>et al.</i> (BiGANs) [12]	17.7	24.5	31.0	29.9	28.0
Zhang <i>et al.</i> (Split-Brain) [53]	17.7	29.3	35.4	35.2	32.8
Noroozi <i>et al.</i> (Counting) [38]	18.0	30.6	34.3	32.5	25.7
Gidaris <i>et al.</i> (RotNet) [17]	18.8	31.7	<u>38.7</u>	38.2	<u>36.5</u>
Jenni & Favaro [21]	<u>19.5</u>	<b>33.3</b>	37.9	<u>38.9</u>	34.9
Mundhenk <i>et al.</i> [36]	<b>19.6</b>	31.8	37.6	37.8	33.7
Noroozi <i>et al.</i> (CC+) [39]	18.9	30.5	35.7	35.4	32.2
Noroozi <i>et al.</i> (CC+vgg-) [39]	19.2	<u>32.0</u>	37.3	37.1	34.6
Wu <i>et al.</i> [48]	16.8	26.5	31.8	34.1	35.6
Doersch <i>et al.</i> (Context) [10] <sup>*</sup>	16.2	23.3	30.2	31.7	29.6
Ren & Lee [44] <sup>*</sup>	16.5	27.0	30.5	30.1	26.5
Caron <i>et al.</i> (DeepCluster) [6] <sup>*†</sup>	13.4	32.3	41.0	39.6	38.2
Ours	19.3	<b>33.3</b>	<b>40.8</b>	<b>41.8</b>	<b>44.3</b>
Ours ( <i>bigger</i> AlexNet) <sup>*</sup>	20.8	35.2	41.8	44.3	44.4
Ours ( <i>bigger</i> AlexNet) <sup>*†</sup>	22.2	38.2	45.7	48.7	48.3

ImageNet (AlexNet)

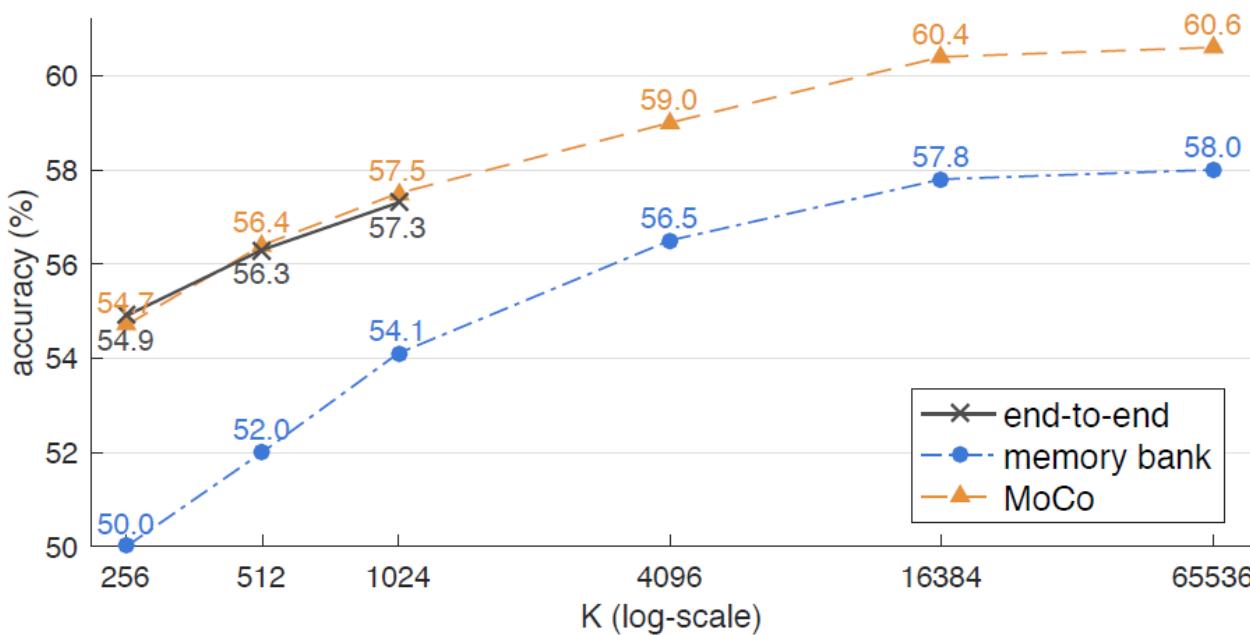
# MoCo: Momentum Contrast



- Dynamic dictionary as a queue of data samples
- Momentum update  $\theta_k \leftarrow m\theta_k + (1 - m)\theta_q \quad m = 0.999$ .

K. He, H. Fan, Y. Wu, S. Xie R. Girshick, Momentum Contrast for Unsupervised Visual Representation Learning, CVPR 2020

# Ablation studies



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

pre-train	R50-dilated-C5			R50-C4		
	AP <sub>50</sub>	AP	AP <sub>75</sub>	AP <sub>50</sub>	AP	AP <sub>75</sub>
end-to-end	77.8	50.1	53.8	79.7	53.0	57.9
memory bank	79.6	51.9	56.3	80.3	53.9	58.9
MoCo	<b>81.1</b>	<b>53.8</b>	<b>58.6</b>	<b>81.4</b>	<b>55.2</b>	<b>61.2</b>

method	architecture	#params (M)	accuracy (%)
Exemplar [15]	R50w3×	211	46.0 [36]
RelativePosition [11]	R50w2×	94	51.4 [36]
Jigsaw [43]	R50w2×	94	44.6 [36]
Rotation [17]	Rv50w4×	86	55.4 [36]
Colorization [62]	R101*	28	39.6 [12]
DeepCluster [3]	VGG [51]	15	48.4 [4]
BigBiGAN [14]	R50	24	56.6
	Rv50w4×	86	61.3

methods based on contrastive learning follow:

InstDisc [59]	R50	24	54.0
LocalAgg [64]	R50	24	58.8
CPC v1 [44]	R101*	28	48.7
CPC v2 [33]	R170* <sub>wider</sub>	303	65.9
CMC [54]	R50 <sub>L+ab</sub>	47	64.1 <sup>†</sup>
	R50w2× <sub>L+ab</sub>	188	68.4 <sup>†</sup>
AMDIM [2]	AMDIM <sub>small</sub>	194	63.5 <sup>†</sup>
	AMDIM <sub>large</sub>	626	68.1 <sup>†</sup>
<b>MoCo</b>	R50	24	60.6
	RX50	46	63.9
	R50w2×	94	65.4
	R50w4×	375	<b>68.6</b>

pre-train	AP <sub>50</sub>	AP	AP <sub>75</sub>
random init.	58.0	32.8	32.5
super. IN-1M	81.5	53.6	58.9
<b>MoCo</b> IN-1M	81.1 (-0.4)	53.8 (+0.2)	58.6 (-0.3)
<b>MoCo</b> IG-1B	81.6 (+0.1)	54.8 ( <b>+1.2</b> )	60.3 ( <b>+1.4</b> )

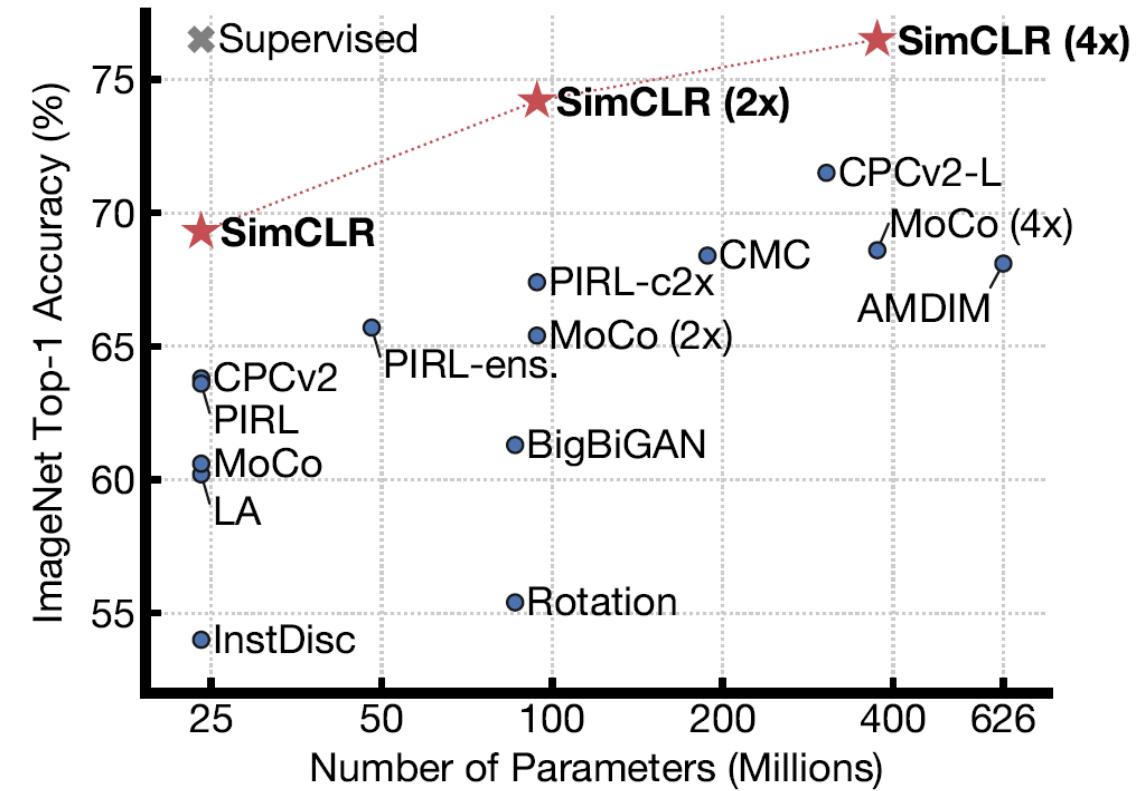
(a) Faster R-CNN, R50-dilated-C5

pre-train	AP <sub>50</sub>	AP	AP <sub>75</sub>
random init.	52.5	28.1	26.2
super. IN-1M	80.8	52.0	56.5
<b>MoCo</b> IN-1M	81.4 ( <b>+0.6</b> )	55.2 ( <b>+3.2</b> )	61.2 ( <b>+4.7</b> )
<b>MoCo</b> IG-1B	82.1 ( <b>+1.3</b> )	56.2 ( <b>+4.2</b> )	62.3 ( <b>+5.8</b> )

(b) Faster R-CNN, R50-C4

# SimCLR

- Data augmentations
- Adding nonlinear transformation
  - Projection head
- Larger batch sizes and more training steps



Ting Chen, Simon Kornblith, Mohammad Norouzi, Geoffrey Hinton, A Simple Framework for Contrastive Learning of Visual Representations, ICML 2020.

# SimCLR

- Loss function

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}$$

- Data Augmentation



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



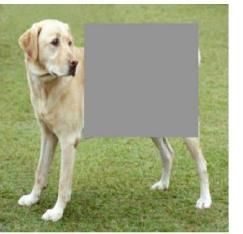
(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate {90°, 180°, 270°}



(g) Cutout



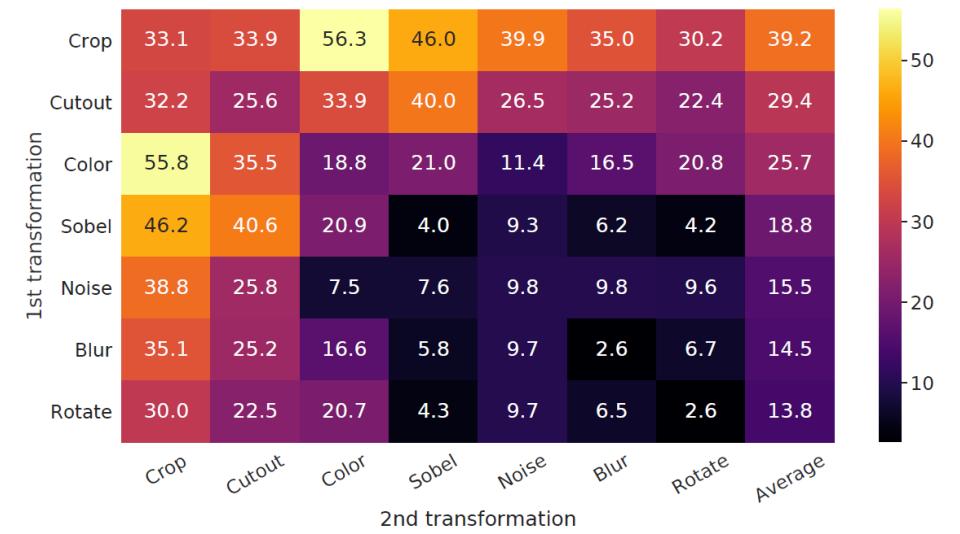
(h) Gaussian noise



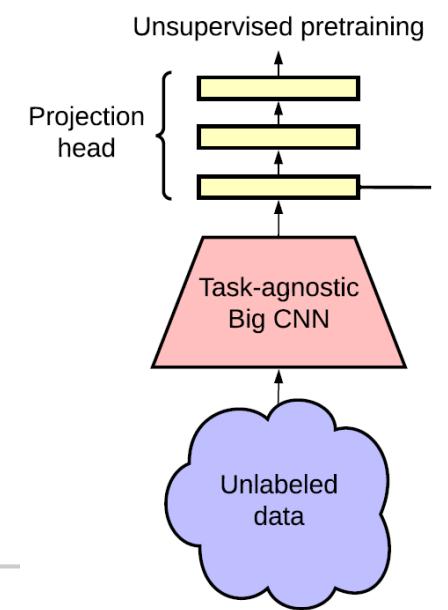
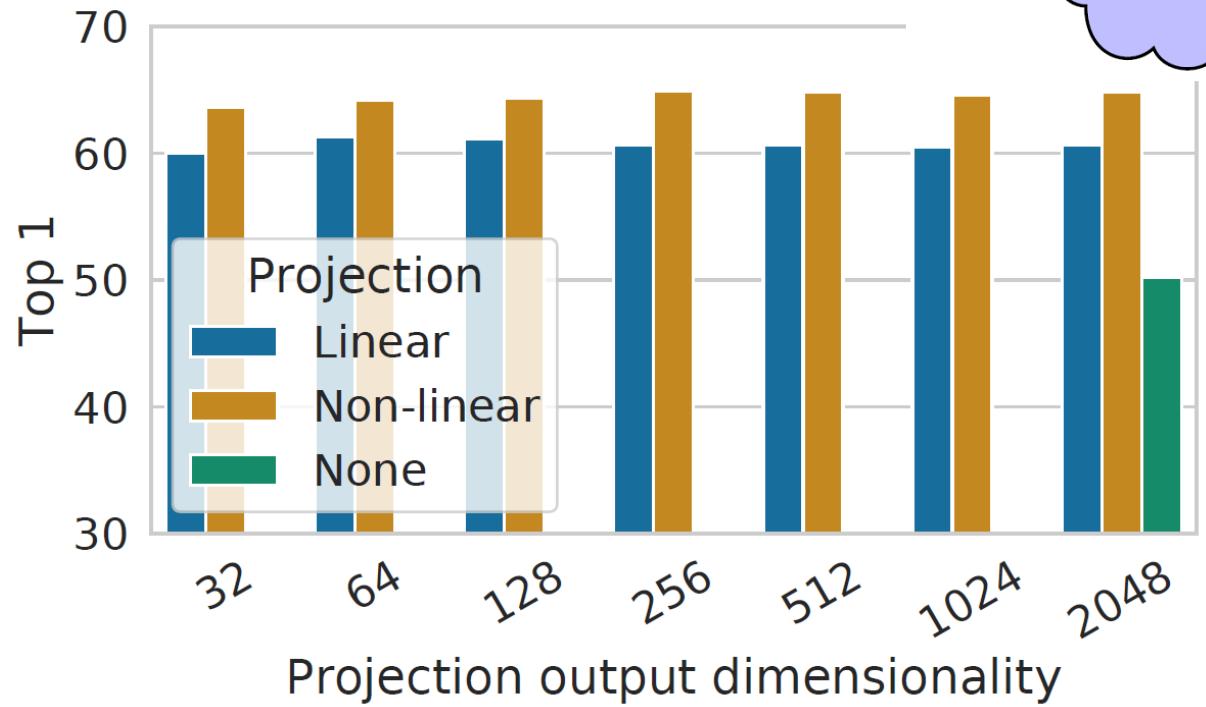
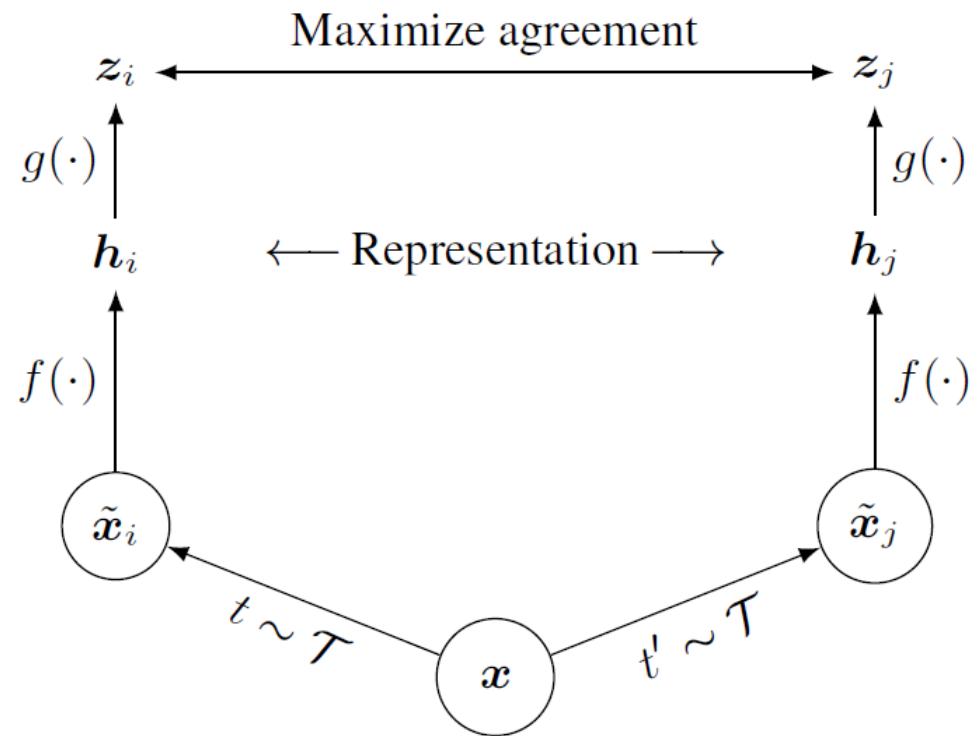
(i) Gaussian blur



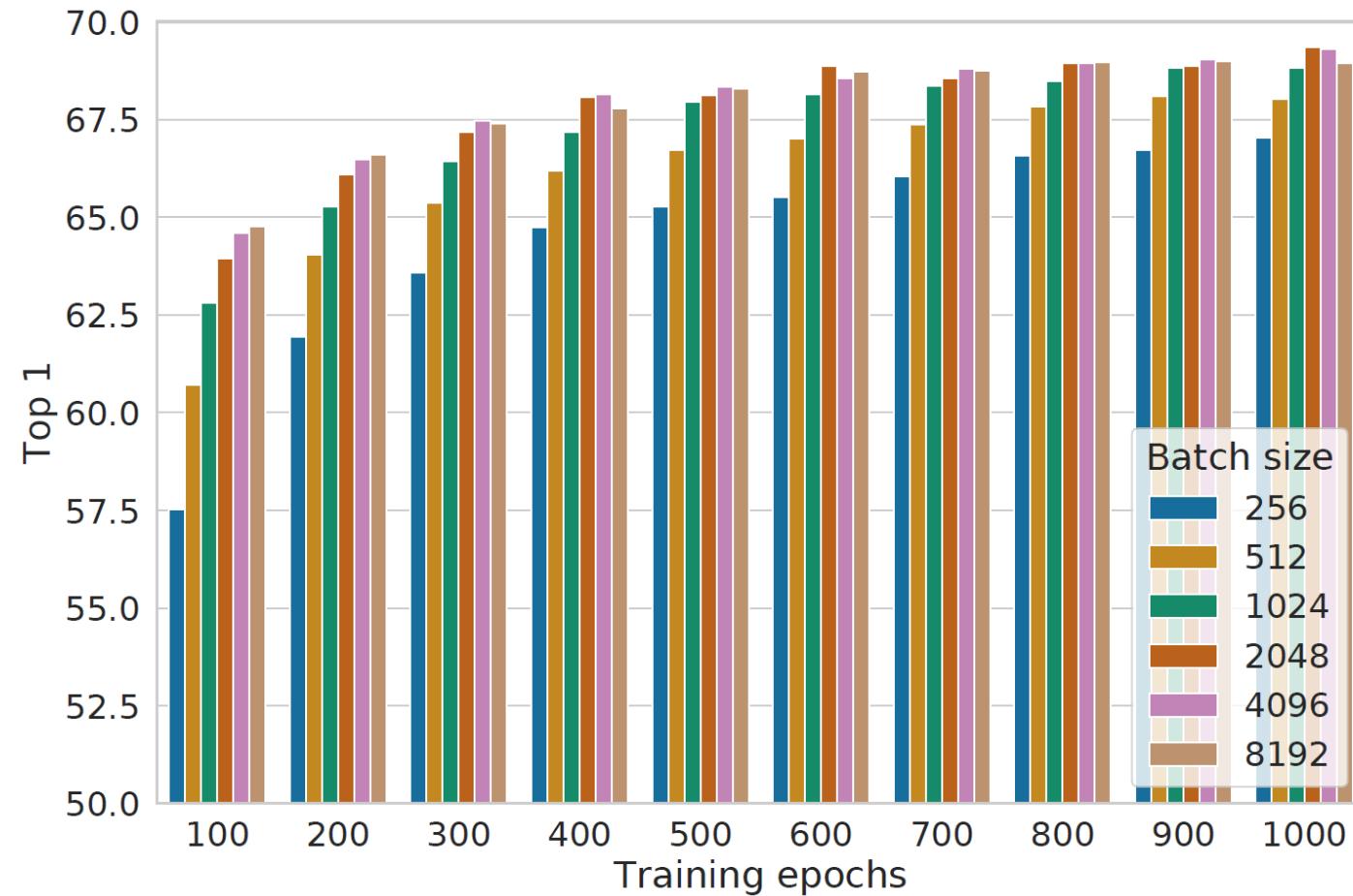
(j) Sobel filtering



# Adding nonlinear transformation



# Larger batch sizes and longer training



Method	Architecture	Param.	Top 1	Top 5
<i>Methods using ResNet-50:</i>				
Local Agg.	ResNet-50	24	60.2	-
MoCo	ResNet-50	24	60.6	-
PIRL	ResNet-50	24	63.6	-
CPC v2	ResNet-50	24	63.8	85.3
SimCLR (ours)	ResNet-50	24	<b>69.3</b>	<b>89.0</b>
<i>Methods using other architectures:</i>				
Rotation	RevNet-50 (4×)	86	55.4	-
BigBiGAN	RevNet-50 (4×)	86	61.3	81.9
AMDIM	Custom-ResNet	626	68.1	-
CMC	ResNet-50 (2×)	188	68.4	88.2
MoCo	ResNet-50 (4×)	375	68.6	-
CPC v2	ResNet-161 (*)	305	71.5	90.1
SimCLR (ours)	ResNet-50 (2×)	94	74.2	92.0
SimCLR (ours)	ResNet-50 (4×)	375	<b>76.5</b>	<b>93.2</b>

Method	Architecture	Label fraction		
		1%	10%	Top 5
<i>Methods using other label-propagation:</i>				
Pseudo-label	ResNet50	51.6	82.4	
VAT+Entropy Min.	ResNet50	47.0	83.4	
UDA (w. RandAug)	ResNet50	-	88.5	
FixMatch (w. RandAug)	ResNet50	-	89.1	
S4L (Rot+VAT+En. M.)	ResNet50 (4×)	-	91.2	
<i>Methods using representation learning only:</i>				
InstDisc	ResNet50	39.2	77.4	
BigBiGAN	RevNet-50 (4×)	55.2	78.8	
PIRL	ResNet-50	57.2	83.8	
CPC v2	ResNet-161(*)	77.9	91.2	
SimCLR (ours)	ResNet-50	75.5	87.8	
SimCLR (ours)	ResNet-50 (2×)	83.0	91.2	
SimCLR (ours)	ResNet-50 (4×)	<b>85.8</b>	<b>92.6</b>	

	Food	CIFAR10	CIFAR100	Birdsnap	SUN397	Cars	Aircraft	VOC2007	DTD	Pets	Caltech-101	Flowers
<i>Linear evaluation:</i>												
SimCLR (ours)	<b>76.9</b>	<b>95.3</b>	80.2	48.4	<b>65.9</b>	60.0	61.2	<b>84.2</b>	<b>78.9</b>	89.2	<b>93.9</b>	<b>95.0</b>
Supervised	75.2	<b>95.7</b>	<b>81.2</b>	<b>56.4</b>	64.9	<b>68.8</b>	<b>63.8</b>	83.8	<b>78.7</b>	<b>92.3</b>	<b>94.1</b>	94.2
<i>Fine-tuned:</i>												
SimCLR (ours)	<b>89.4</b>	<b>98.6</b>	<b>89.0</b>	<b>78.2</b>	<b>68.1</b>	<b>92.1</b>	<b>87.0</b>	<b>86.6</b>	<b>77.8</b>	92.1	<b>94.1</b>	97.6
Supervised	88.7	98.3	<b>88.7</b>	<b>77.8</b>	67.0	91.4	<b>88.0</b>	86.5	<b>78.8</b>	<b>93.2</b>	<b>94.2</b>	<b>98.0</b>
Random init	88.3	96.0	81.9	<b>77.0</b>	53.7	91.3	84.8	69.4	64.1	82.7	72.5	92.5

# MoCo v2

- MLP projection head (in SimCLR)
- More data augmentation (in SimCLR): blur augmentation

case	unsup. pre-train				ImageNet acc.	VOC detection			$\tau$	0.07	0.1	0.2	0.3	0.4	0.5
	MLP	aug+	cos	epochs		AP <sub>50</sub>	AP	AP <sub>75</sub>							
supervised					76.5	81.3	53.5	58.8							
MoCo v1				200	60.6	81.5	55.9	62.6							
(a)	✓			200	66.2	82.0	56.4	62.6	w/o MLP	60.6	<b>60.7</b>	59.0	58.2	57.2	56.4
(b)		✓		200	63.4	82.2	56.8	63.2	w/ MLP	62.9	64.9	<b>66.2</b>	65.7	65.0	64.3
(c)	✓	✓		200	67.3	<b>82.5</b>	57.2	63.9							
(d)	✓	✓	✓	200	67.5	82.4	57.0	63.6							
(e)	✓	✓	✓	800	71.1	<b>82.5</b>	<b>57.4</b>	<b>64.0</b>							

Improved Baselines with Momentum Contrastive Learning, Arxiv 2020.

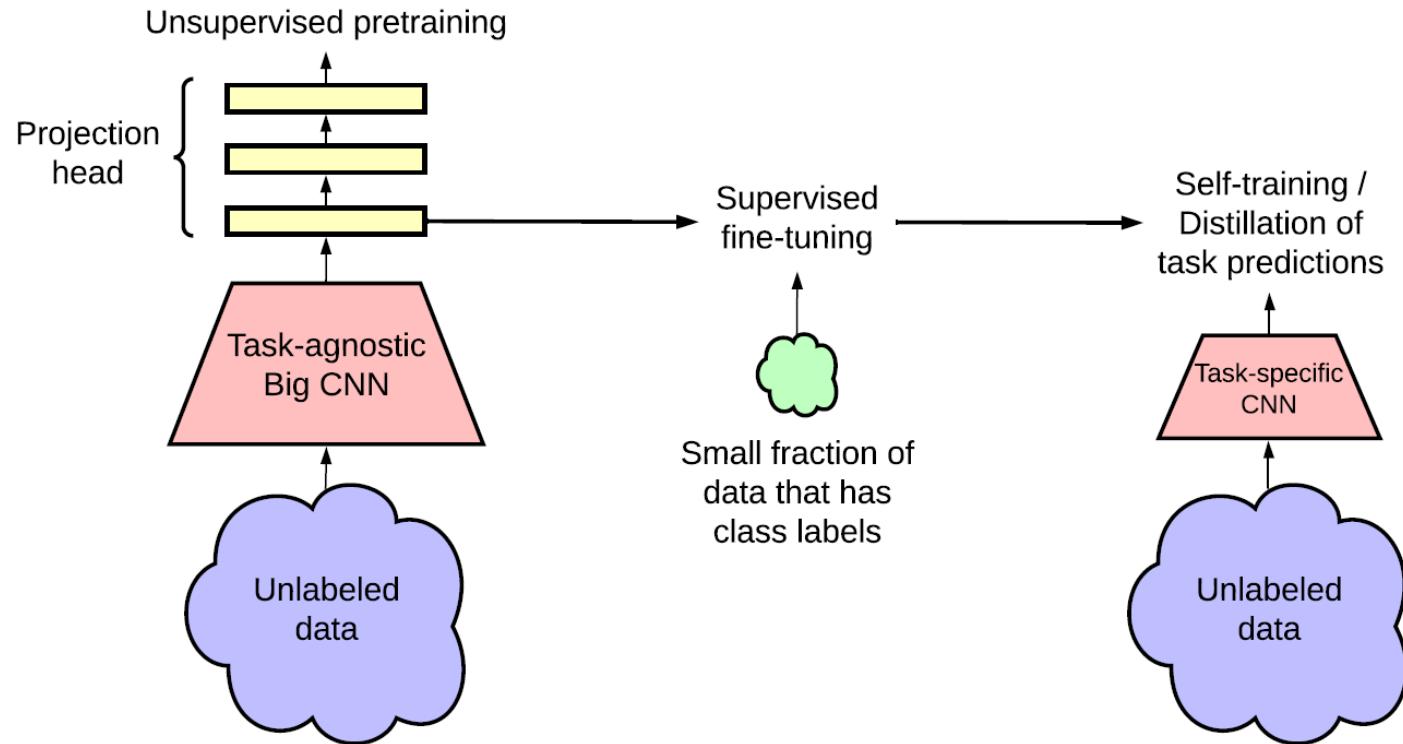
# MoCo v2

case	MLP	unsup. pre-train				ImageNet acc.
		aug+	cos	epochs	batch	
MoCo v1 [6]				200	256	60.6
SimCLR [2]	✓	✓	✓	200	256	61.9
SimCLR [2]	✓	✓	✓	200	8192	66.6
<b>MoCo v2</b>	✓	✓	✓	200	256	<b>67.5</b>

*results of longer unsupervised training follow:*

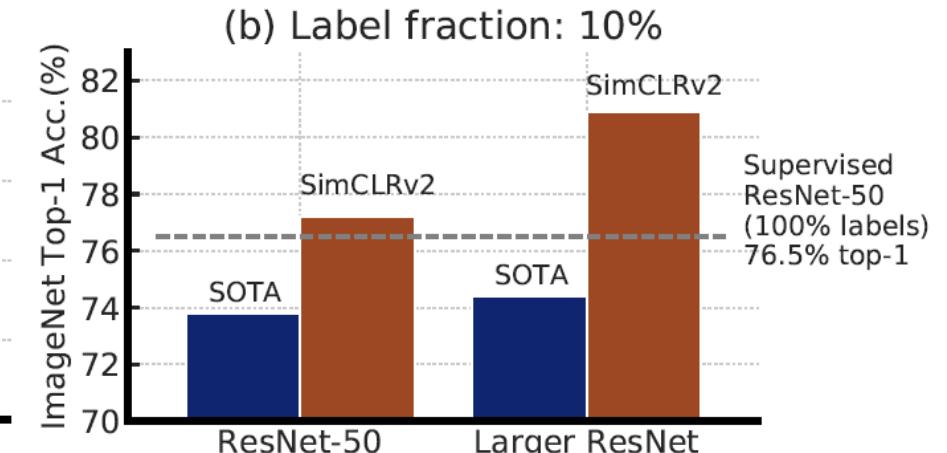
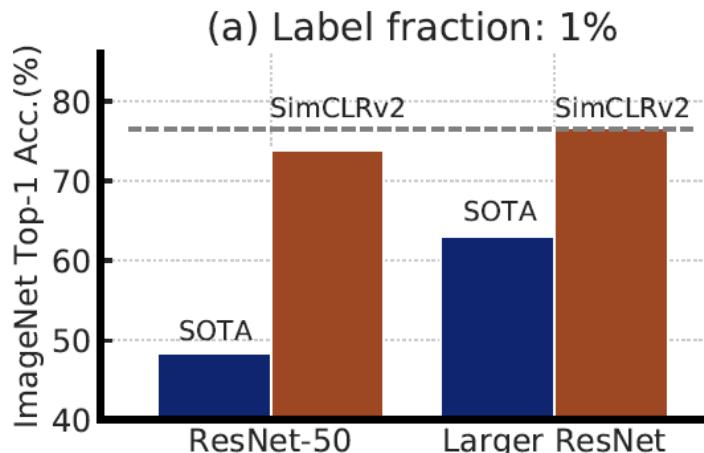
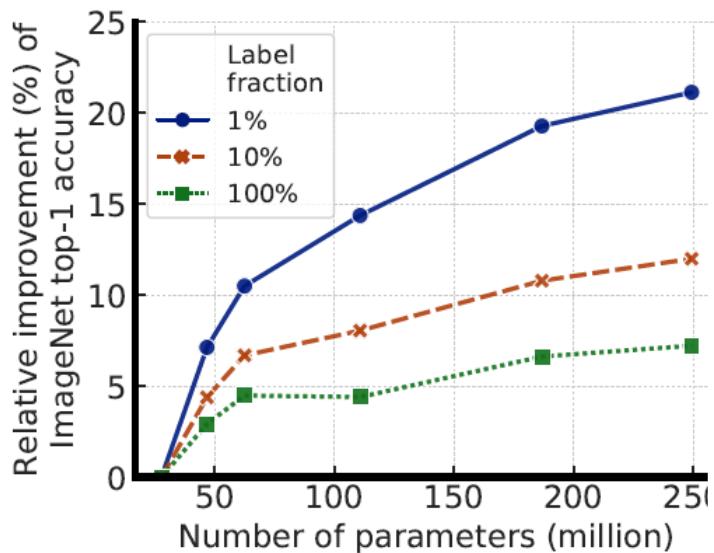
SimCLR [2]	✓	✓	✓	1000	4096	69.3
<b>MoCo v2</b>	✓	✓	✓	800	256	<b>71.1</b>

# SimCLR v2 (NeurIPS 2020)

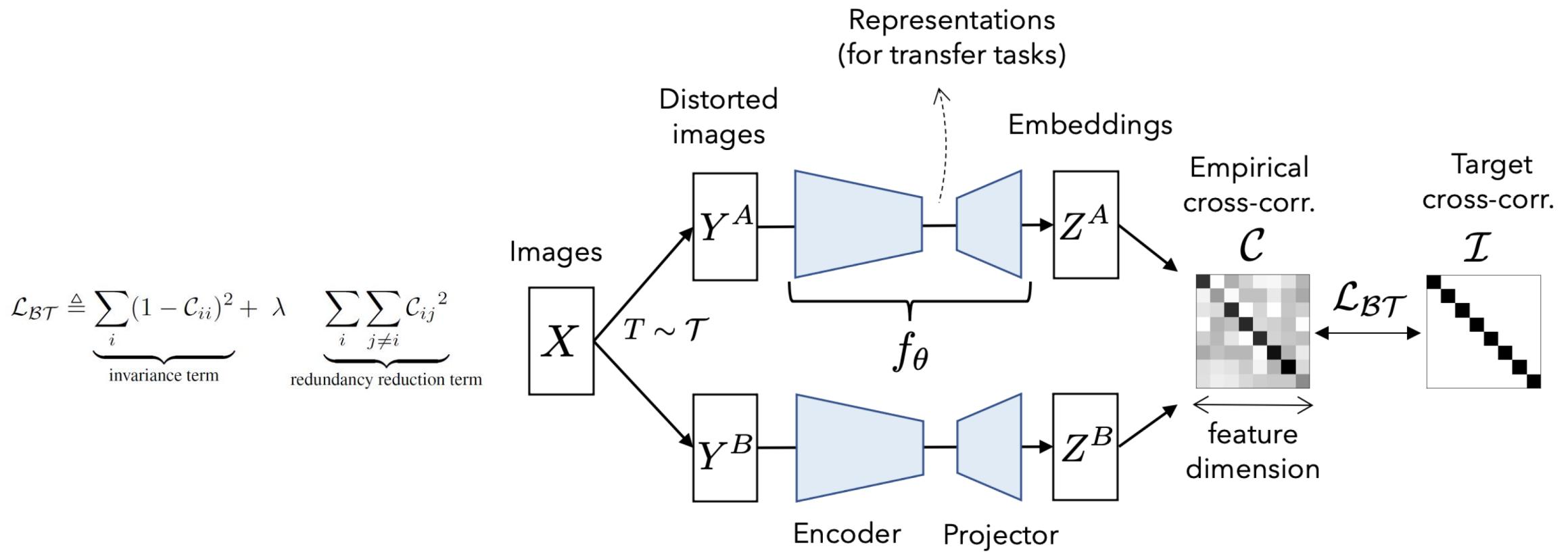


$$\mathcal{L}^{\text{distill}} = - \sum_{\mathbf{x}_i \in \mathcal{D}} \left[ \sum_y P^T(y|\mathbf{x}_i; \tau) \log P^S(y|\mathbf{x}_i; \tau) \right]$$

# SimCLRv2

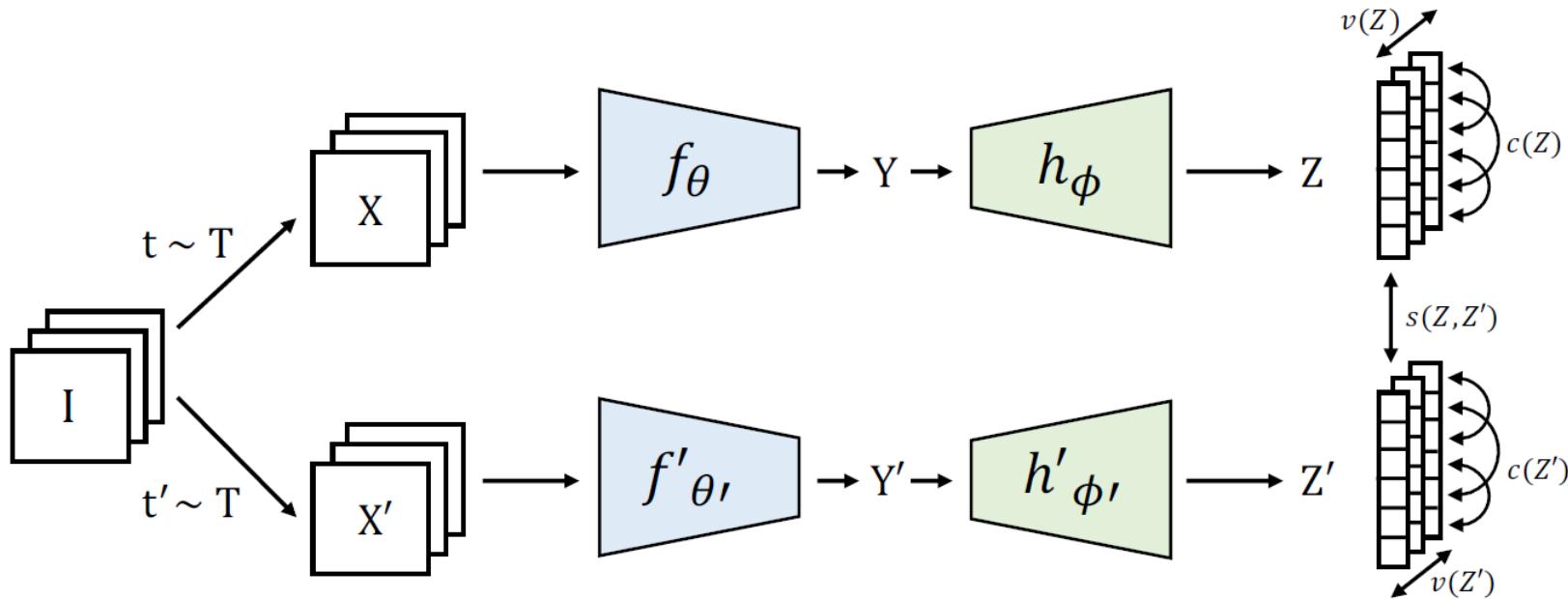


# Barlow Twins (ICML 2021)



Barlow Twins: Self-Supervised Learning via Redundancy Reduction, ICML 2021.

# VICReg (ICLR 2022)



- Variance
- Invariance
- Covariance

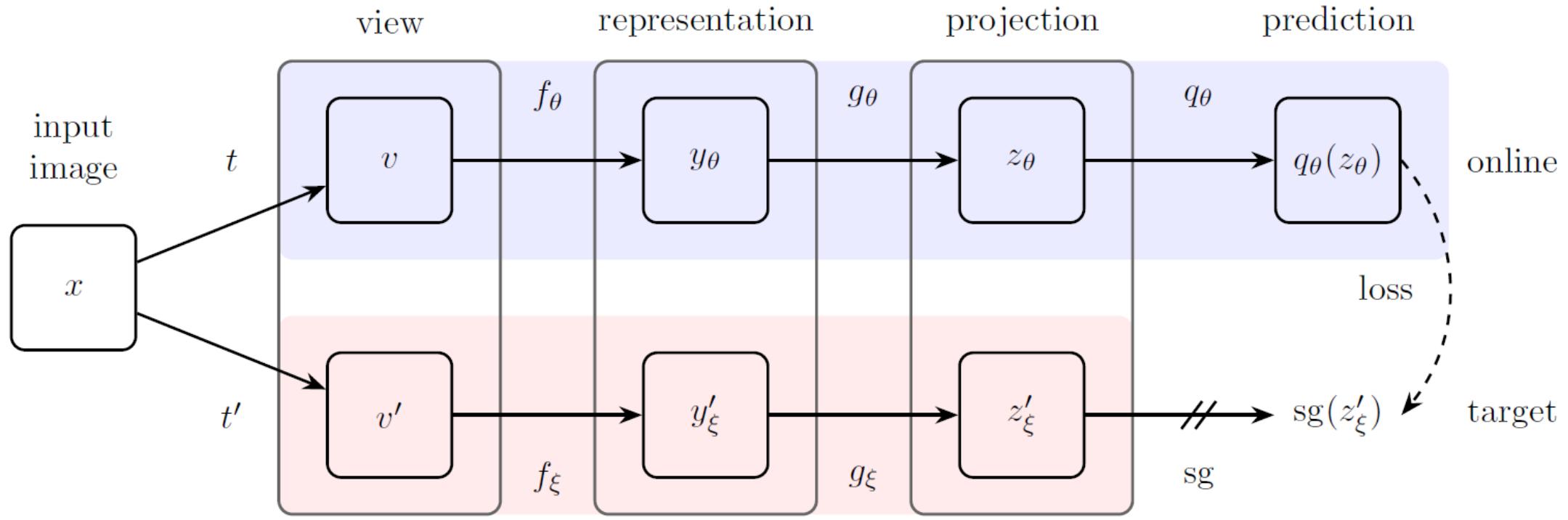
$$v(Z) = \frac{1}{d} \sum_{j=1}^d \max(0, \gamma - S(z^j, \epsilon))$$

$$s(Z, Z') = \frac{1}{n} \sum_i \|z_i - z'_i\|_2^2$$

$$c(Z) = \frac{1}{d} \sum_{i \neq j} [C(Z)]_{i,j}^2$$

$v$	: maintain variance
$c$	: bring covariance to zero
$s$	: minimize distance
$T$	: distribution of transformations
$t, t'$	: random transformations
$f_\theta, f'_\theta$	: encoders
$h_\phi, h'_\phi$	: expanders
$I$	: batch of images
$X, X'$	: batches of views
$Y, Y'$	: batches of representations
$Z, Z'$	: batches of embeddings

# BYOL (NeurIPS 2020)



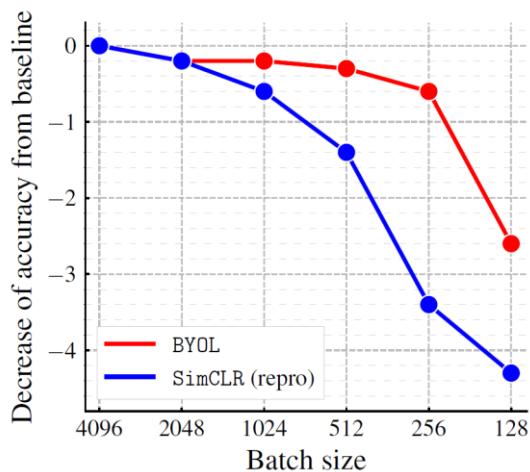
Exponential moving average

$$\xi \leftarrow \tau\xi + (1 - \tau)\theta$$

$$\mathcal{L}_{\theta,\xi} \triangleq \|\overline{q}_\theta(z_\theta) - \overline{z}'_\xi\|_2^2 = 2 - 2 \cdot \frac{\langle q_\theta(z_\theta), z'_\xi \rangle}{\|q_\theta(z_\theta)\|_2 \cdot \|z'_\xi\|_2}$$

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

# BYOL (NeurIPS 2020)



Local Agg.	60.2	-
PIRL [35]	63.6	-
CPC v2 [32]	63.8	85.3
CMC [11]	66.2	87.0
SimCLR [8]	69.3	89.0
MoCo v2 [37]	71.1	-
InfoMin Aug. [12]	73.0	91.1
BYOL (ours)	<b>74.3</b>	<b>91.6</b>

(a) ResNet-50 encoder.

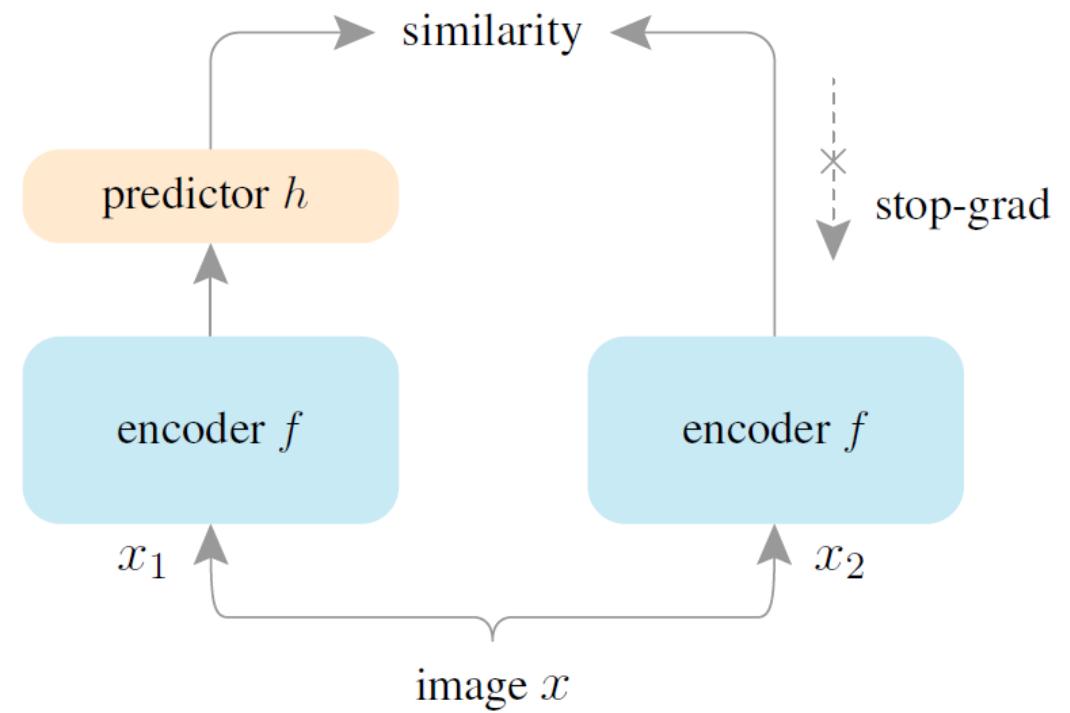
Method	Architecture	Param.	Top-1	Top-5
SimCLR [8]	ResNet-50 (2×)	94M	74.2	92.0
CMC [11]	ResNet-50 (2×)	94M	70.6	89.7
BYOL (ours)	ResNet-50 (2×)	94M	<b>77.4</b>	<b>93.6</b>
CPC v2 [32]	ResNet-161	305M	71.5	90.1
MoCo [9]	ResNet-50 (4×)	375M	68.6	-
SimCLR [8]	ResNet-50 (4×)	375M	76.5	93.2
BYOL (ours)	ResNet-50 (4×)	375M	<b>78.6</b>	<b>94.2</b>
BYOL (ours)	ResNet-200 (2×)	250M	<b>79.6</b>	<b>94.8</b>

(b) Other ResNet encoder architectures.

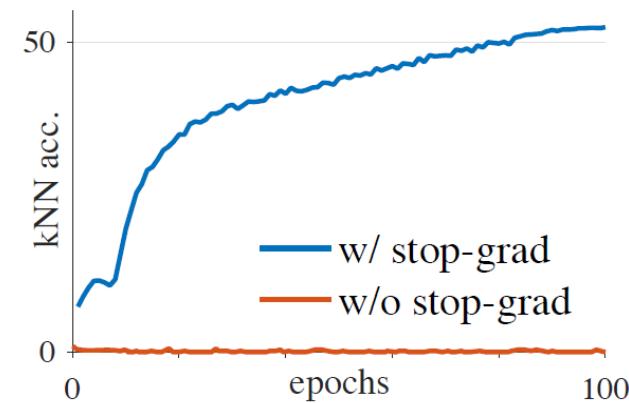
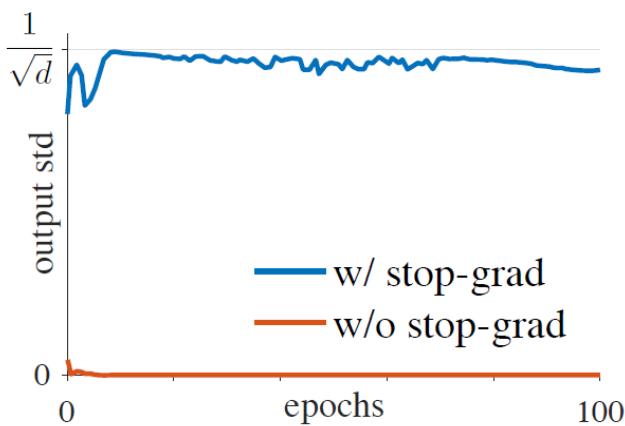
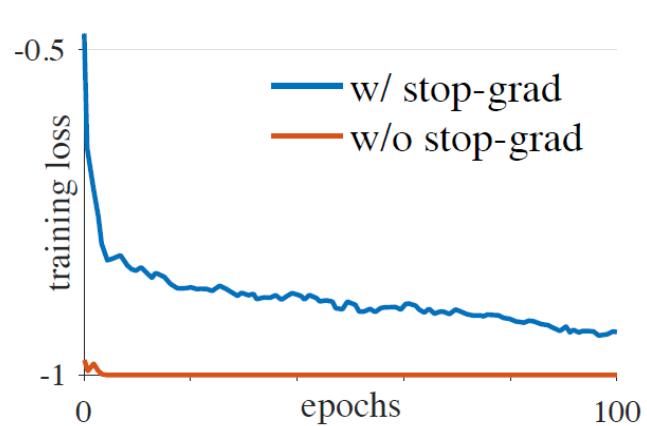
- Large batch size is still required.

# SimSiam (CVPR 2021)

- Can learn meaningful representations even using none of the following:
  - Negative sample pairs
  - Large batches
  - Momentum encoders
- A stop-gradient operation plays an essential role in preventing collapsing



# SimSiam (CVPR 2021)



	acc. (%)
w/ stop-grad	67.7±0.1
w/o stop-grad	0.1

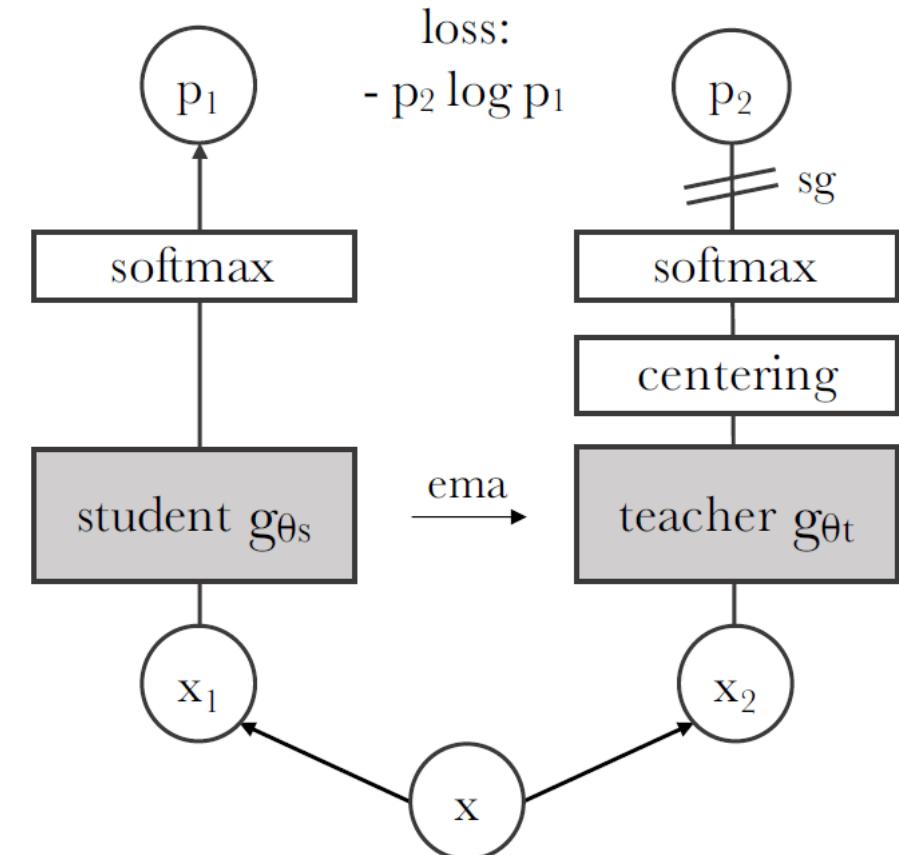
# MoCo v3

- Self-supervised learning for Vision Transformers
  - Batch size: 1k and 2k
  - Learning rate: lr x BatchSize / 256
  - Optimizer: LAMB, AdamW-counterpart of LARS.
  - Freezing the patch projection layer with a fixed random patch projection

framework	model	params	acc. (%)
<i>linear probing:</i>			
iGPT [9]	iGPT-L	1362M	69.0
iGPT [9]	iGPT-XL	6801M	72.0
MoCo v3	ViT-B	86M	76.7
MoCo v3	ViT-L	304M	77.6
MoCo v3	ViT-H	632M	78.1
MoCo v3	ViT-BN-H	632M	79.1
MoCo v3	ViT-BN-L/7	304M	<b>81.0</b>
<i>end-to-end fine-tuning:</i>			
masked patch pred. [16]	ViT-B	86M	79.9 <sup>†</sup>
MoCo v3	ViT-B	86M	83.2
MoCo v3	ViT-L	304M	<b>84.1</b>

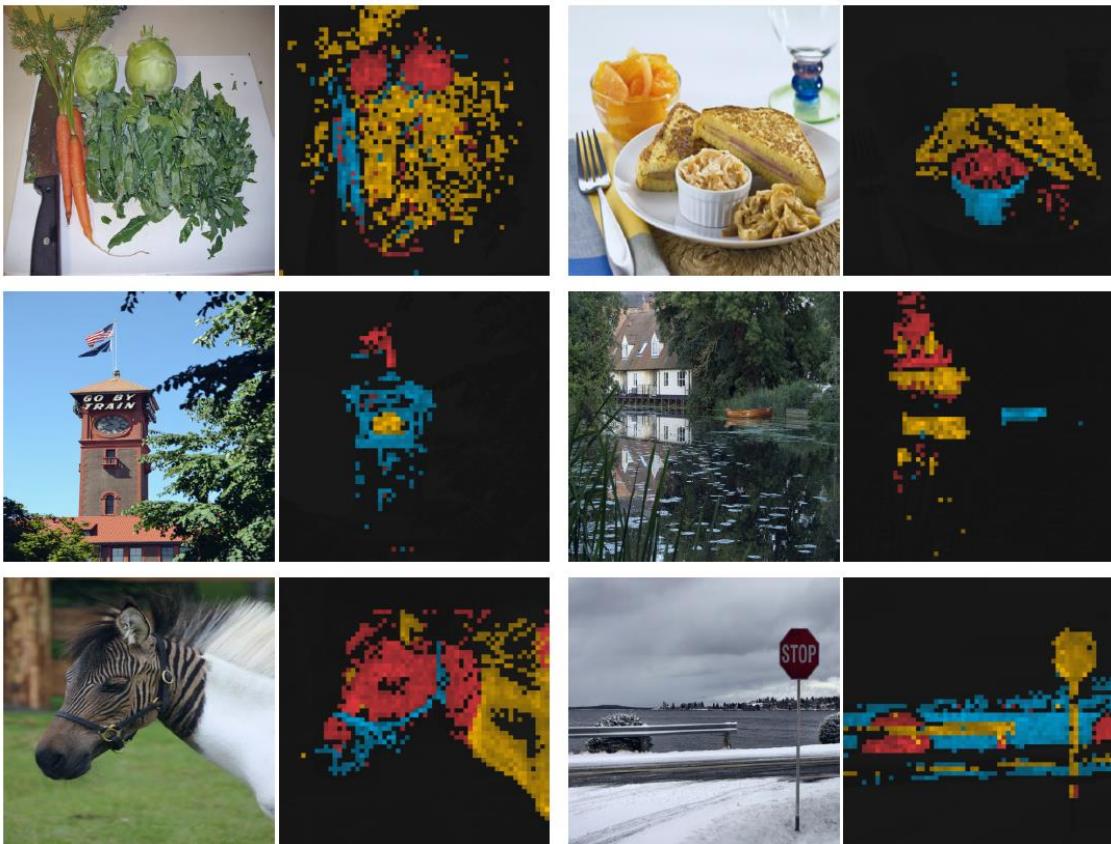
# DINO (ICCV 2021)

- Self-distillation with **no** labels
- Self-supervised ViT
- Containing explicit information about the semantic segmentation of an image
- 78.3% top-1 on ImageNet with a small ViT



# DINO (ICCV 2021)

- Attention map



*Comparison across architectures*

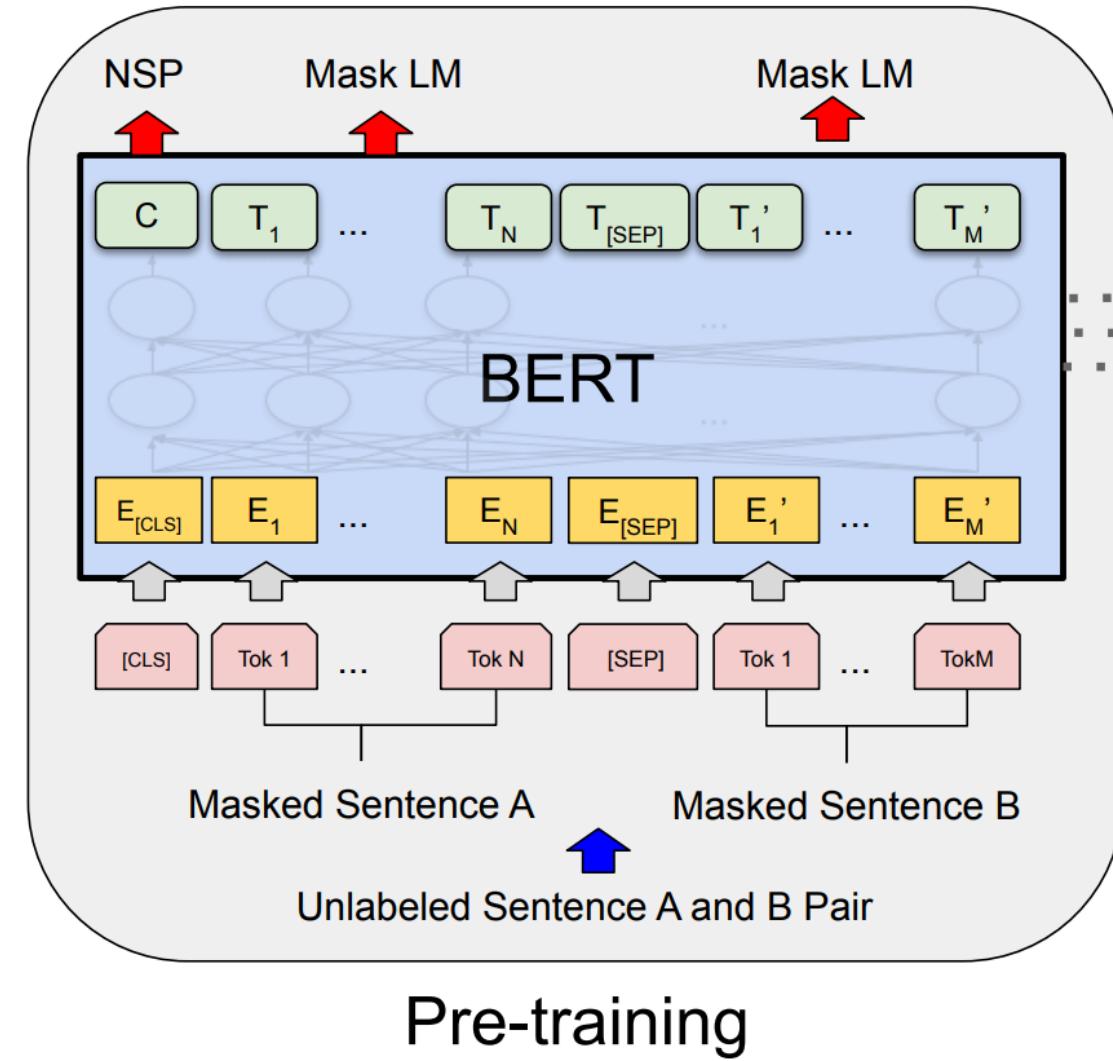
SCLR [11]	RN50w4	375	117	76.8	69.3
SwAV [9]	RN50w2	93	384	77.3	67.3
BYOL [23]	RN50w2	93	384	77.4	–
DINO	ViT-B/16	85	312	78.2	76.1
SwAV [9]	RN50w5	586	76	78.5	67.1
BYOL [23]	RN50w4	375	117	78.6	–
BYOL [23]	RN200w2	250	123	79.6	73.9
DINO	ViT-S/8	21	180	79.7	78.3
SCLRV2 [12]	RN152w3+SK	794	46	79.8	73.1
DINO	ViT-B/8	85	63	<b>80.1</b>	77.4

# Content

- Self-supervised Vision Learning
  - 1. Self-supervised Contextual Modeling
  - 2. Contrastive learning-based self-supervised learning
  - 3. Returning of Self-supervised Contextual Modeling

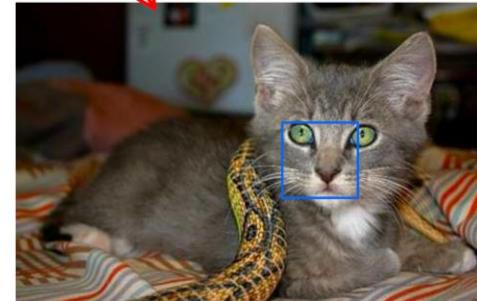
# BERT for Language Understanding (Devlin et al., Arxiv2018)

- Citations: 36172 (2018.10 —)
- Randomly masks some of the tokens from the input, and the objective is to predict the original vocabulary id of the masked word based only on its context.



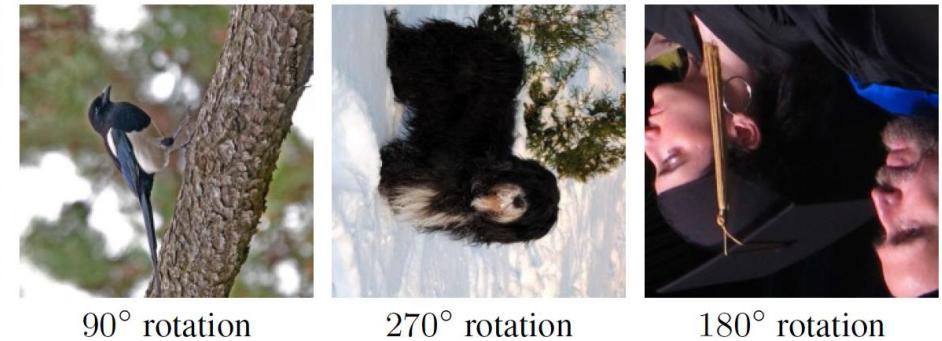
# Deep Context Modeling

- Jigsaw Puzzles

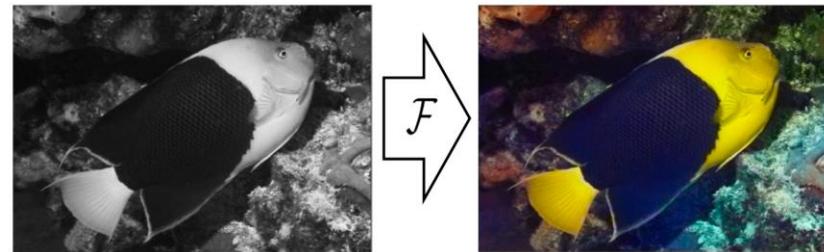


- Image Rotation

Image Rotation



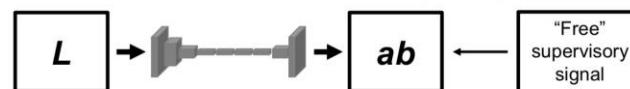
- Image Colorization



- Image Inpainting

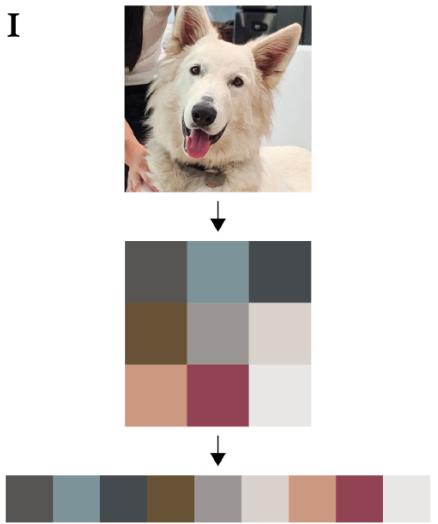
$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

$$(\mathbf{X}, \hat{\mathbf{Y}})$$

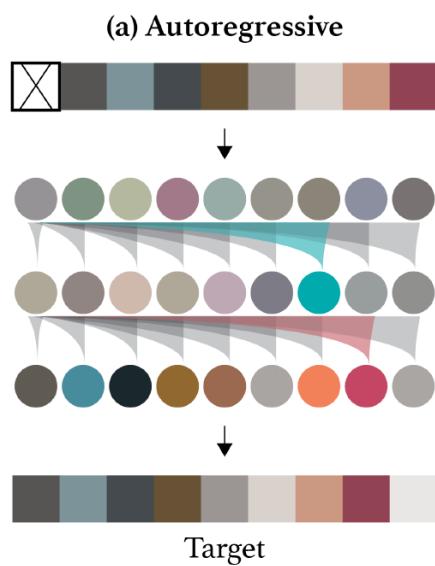


# iGPT (ICML 2021)

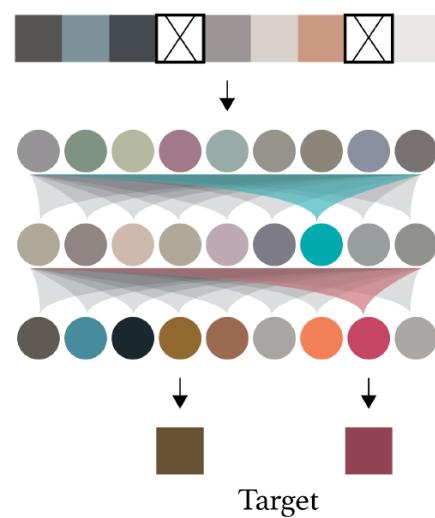
I



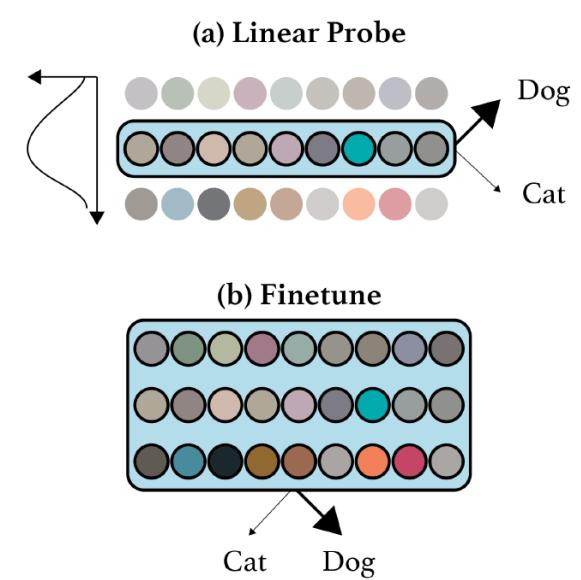
2



(b) BERT

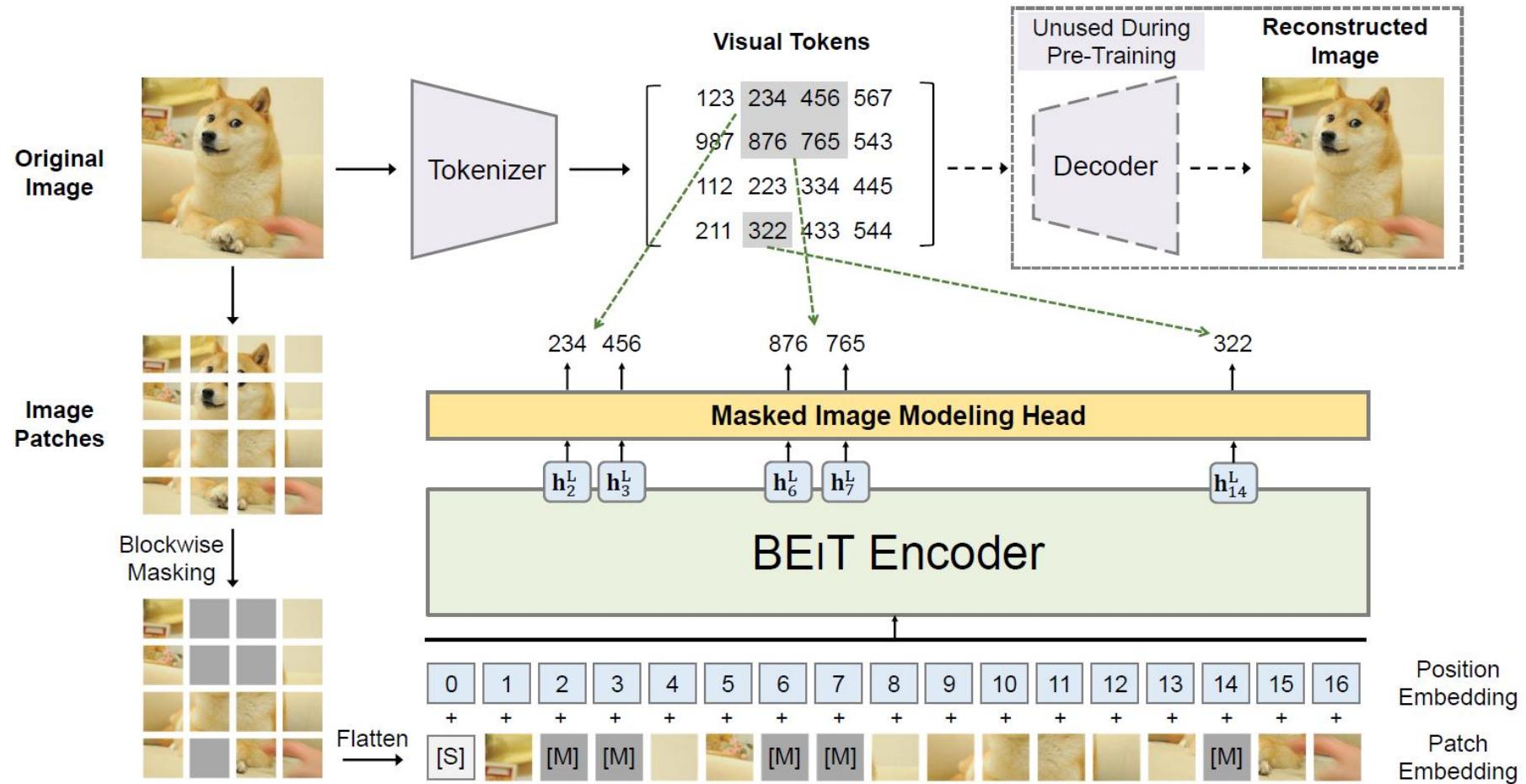


3



$$L_{AR} = \mathbb{E}_{x \sim X} [-\log p(x)] \quad L_{BERT} = \mathbb{E}_{x \sim X} \mathbb{E}_M \sum_{i \in M} [-\log p(x_i | x_{[1,n] \setminus M})]$$

# Masked Image Modeling: BEiT (ICLR 2022)



BEiT: BERT Pre-Training of Image Transformers, ICLR 2022

# BEiT (ICLR 2022)

$$\sum_{(x_i, \tilde{x}_i) \in \mathcal{D}} \left( \underbrace{\mathbb{E}_{z_i \sim q_\phi(z|x_i)} [\log p_\psi(x_i|z_i)]}_{\text{Stage 1: Visual Token Reconstruction}} + \underbrace{\log p_\theta(\hat{z}_i|\tilde{x}_i)}_{\text{Stage 2: Masked Image Modeling}} \right)$$

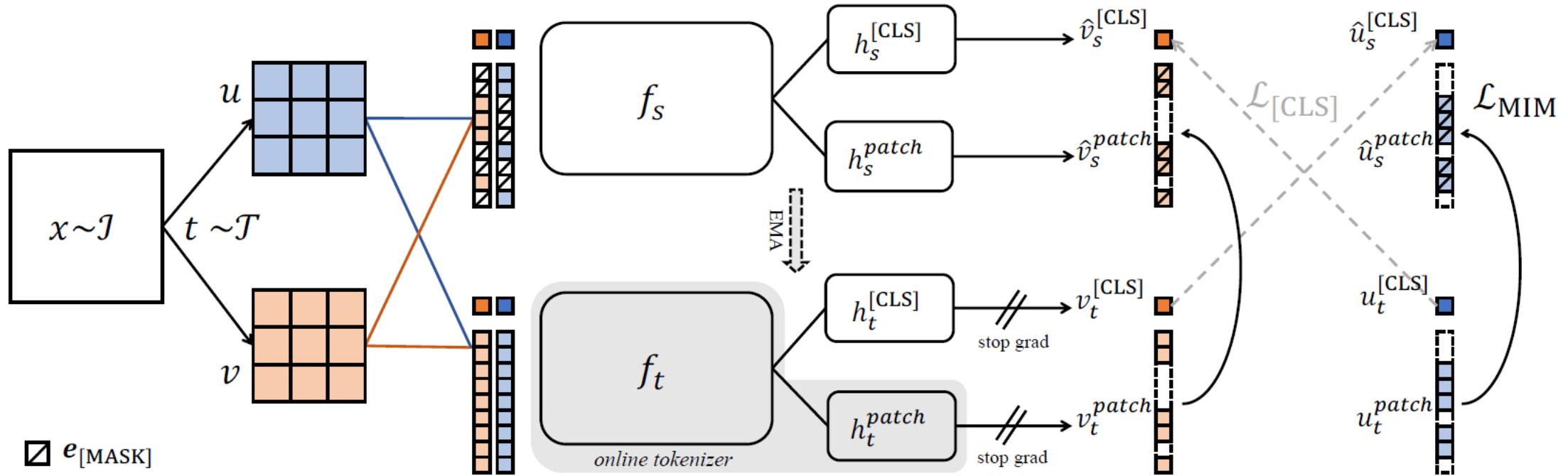
---

*Self-Supervised Pre-Training on ImageNet-1K (without labeled data)*

iGPT-1.36B <sup>†</sup> (Chen et al., 2020a)	1.36B	224 <sup>2</sup>	66.5
ViT <sub>384</sub> -B-JFT300M <sup>‡</sup> (Dosovitskiy et al., 2020)	86M	384 <sup>2</sup>	79.9
MoCo v3-B (Chen et al., 2021)	86M	224 <sup>2</sup>	83.2
MoCo v3-L (Chen et al., 2021)	307M	224 <sup>2</sup>	84.1
DINO-B (Caron et al., 2021)	86M	224 <sup>2</sup>	82.8
BEiT-B (ours)	86M	224 <sup>2</sup>	83.2
BEiT <sub>384</sub> -B (ours)	86M	384 <sup>2</sup>	84.6
BEiT-L (ours)	307M	224 <sup>2</sup>	85.2
BEiT <sub>384</sub> -L (ours)	307M	384 <sup>2</sup>	<b>86.3</b>

---

# iBOT (ICLR 2022)



$$\mathcal{L}_{[\text{CLS}]} = -P_{\theta'}^{[\text{CLS}]}(\hat{v})^T \log P_{\theta}^{[\text{CLS}]}(u) \quad \mathcal{L}_{\text{MIM}} = -\sum_{i=1}^N m_i \cdot P_{\theta'}^{\text{patch}}(u_i)^T \log P_{\theta}^{\text{patch}}(\hat{u}_i)$$

iBOT : IMAGE BERT PRE-TRAINING WITH ONLINE TOKENIZER, ICLR 2022.

# iBOT (ICLR 2022)

Table 1:  **$k$ -NN and linear probing on ImageNet-1K.**  $\dagger$  denotes using selective kernel.  $\ddagger$  denotes pre-training on ImageNet-22K.

Method	Arch.	Par.	im/s	Epo. <sup>†</sup>	$k$ -NN	Lin.
<i>SSL big ResNets</i>						
MoCov3	RN50	23	1237	1600	-	74.6
SwAV	RN50	23	1237	2400	65.7	75.3
DINO	RN50	23	1237	3200	67.5	75.3
BYOL	RN200w2	250	123	2000	73.9	79.6
SCLRv2	RN152w3 $\dagger$	794	46	2000	73.1	79.8
<i>SSL Transformers</i>						
MoCov3	ViT-S/16	21	1007	1200	-	73.4
MoCov3	ViT-B/16	85	312	1200	-	76.7
SwAV	ViT-S/16	21	1007	2400	66.3	73.5
DINO	ViT-S/16	21	1007	3200	74.5	77.0
DINO	ViT-B/16	85	312	1600	76.1	78.2
EsViT	Swin-T/7	28	726	1200	75.7	78.1
EsViT	Swin-T/14	28	593	1200	77.0	78.7
iBOT	ViT-S/16	21	1007	3200	75.2	77.9
iBOT	Swin-T/7	28	726	1200	75.3	78.6
iBOT	Swin-T/14	28	593	1200	76.2	79.3
iBOT	ViT-B/16	85	312	1600	77.1	79.5
iBOT	ViT-L/16	307	102	1200	<b>78.0</b>	81.0
iBOT $\ddagger$	ViT-L/16	307	102	200	72.9	<b>82.3</b>

Table 2: **Fine-tuning on ImageNet-1K.**

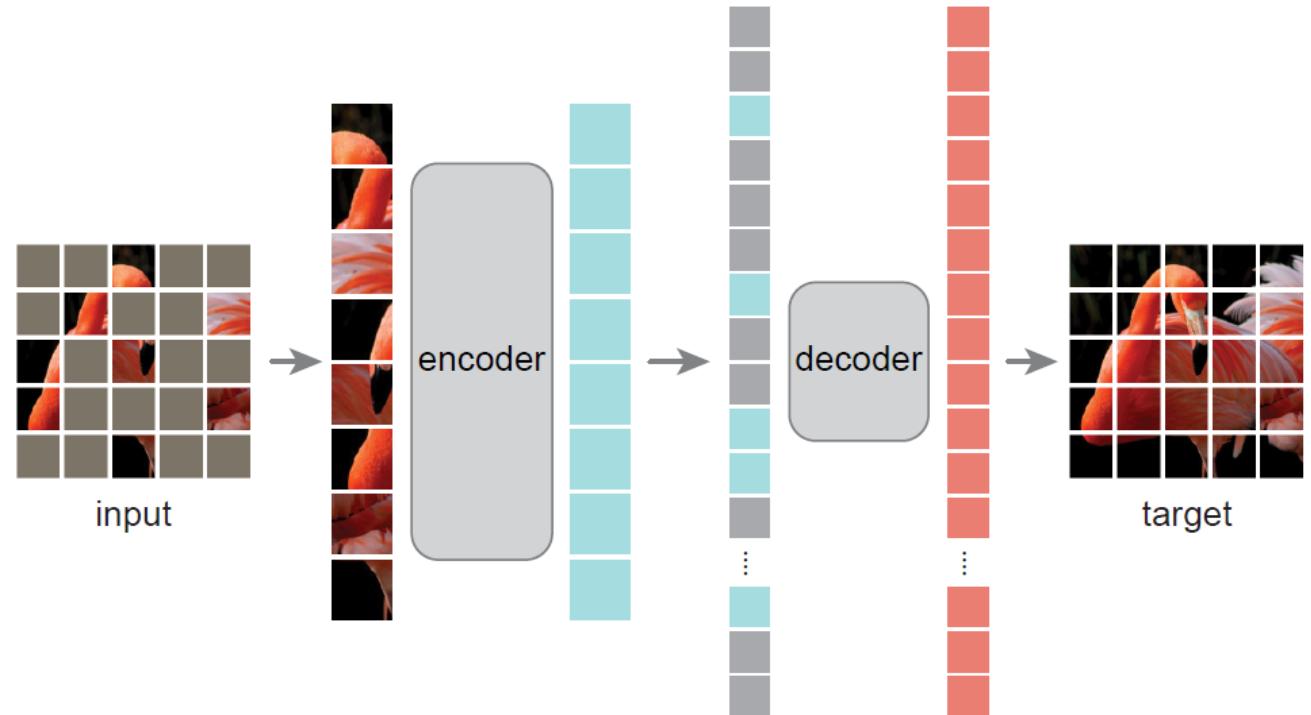
Method	Arch.	Epo. <sup>†</sup>	Acc.
Rand.	ViT-S/16	-	79.9
MoCov3	ViT-S/16	600	81.4
DINO	ViT-S/16	3200	82.0
iBOT	ViT-S/16	3200	<b>82.3</b>
Rand.	ViT-B/16	-	81.8
MoCov3	ViT-B/16	600	83.2
BEiT	ViT-B/16	800	83.4
DINO	ViT-B/16	1600	83.6
iBOT	ViT-B/16	1600	<b>84.0</b>
MoCov3	ViT-L/16	600	84.1
iBOT	ViT-L/16	1000	<b>84.8</b>
BEiT	ViT-L/16	800	<b>85.2</b>

Table 3: **Fine-tuning on ImageNet-1K.**  
Pre-training on ImageNet-22K.

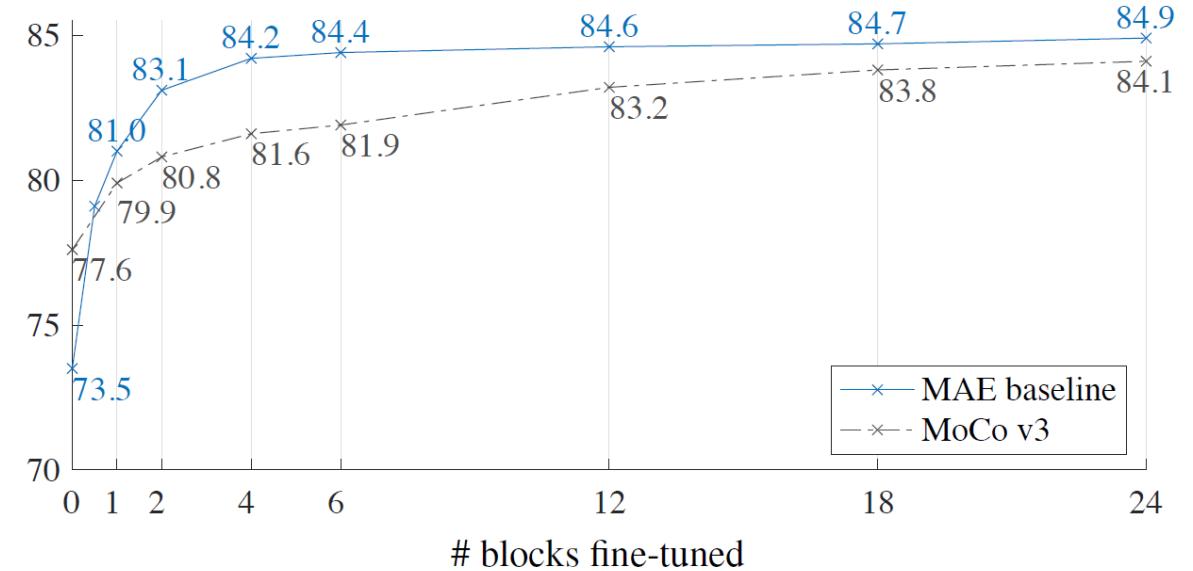
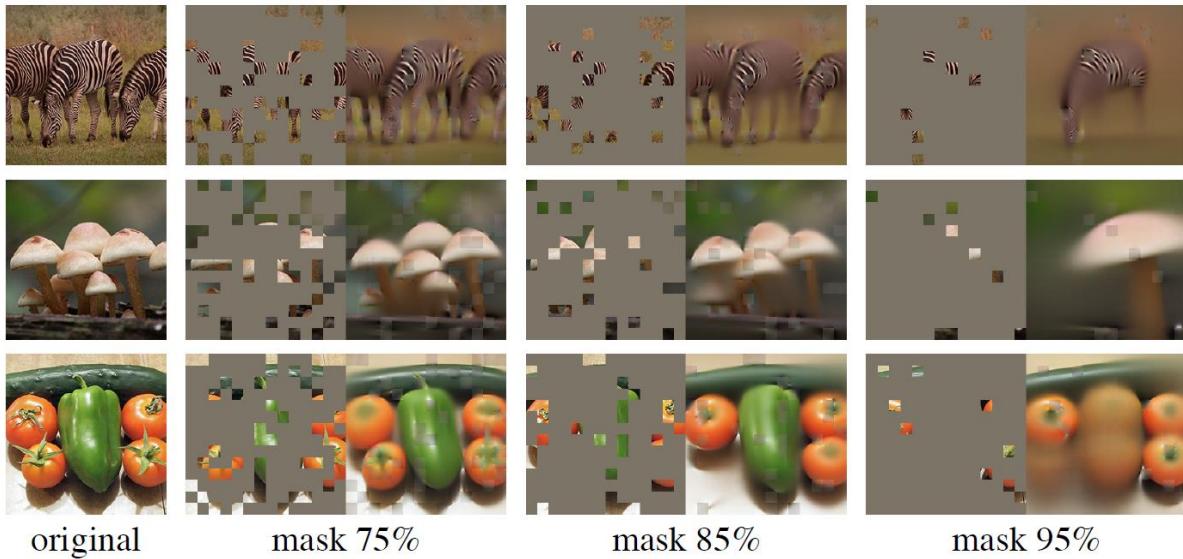
Method	Arch.	Epo. <sup>†</sup>	Acc.
BEiT	ViT-B/16	150	83.7
iBOT	ViT-B/16	320	<b>84.4</b>
BEiT	ViT-L/16	150	86.0
iBOT	ViT-L/16	200	86.6
iBOT	ViT <sub>512</sub> -L/16	200	<b>87.8</b>

# Masked Autoencoders (MAE)(CVPR'22)

- Asymmetric encoder-decoder
- Masking a high proportion of the input image (e.g., 75%)

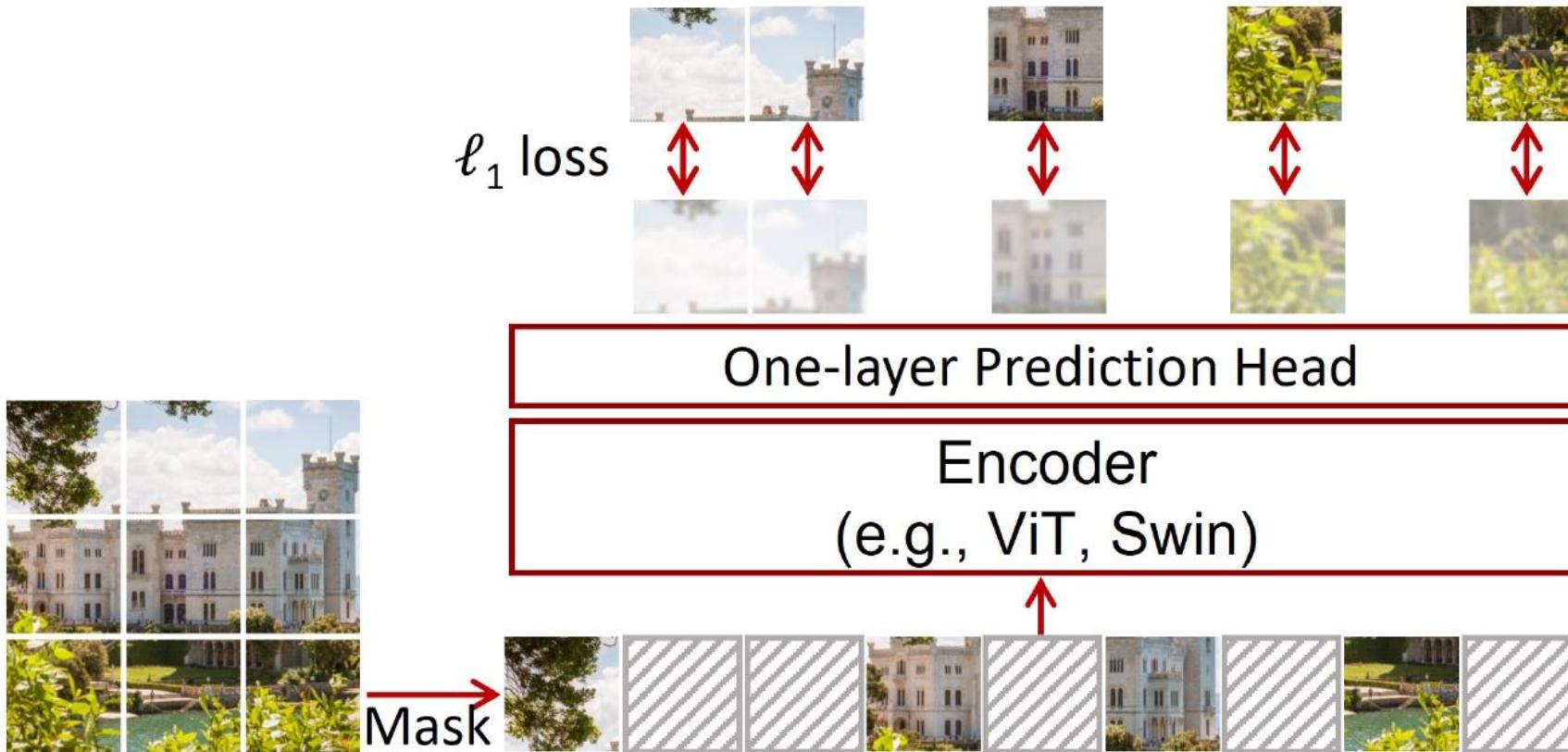


# MAE



method	pre-train data	ViT-B	ViT-L	ViT-H	ViT-H <sub>448</sub>
scratch, our impl.	-	82.3	82.6	83.1	-
DINO [5]	IN1K	82.8	-	-	-
MoCo v3 [9]	IN1K	83.2	84.1	-	-
BEiT [2]	IN1K+DALLE	83.2	85.2	-	-
MAE	IN1K	<u>83.6</u>	<u>85.9</u>	<u>86.9</u>	<b>87.8</b>

# SimMIM (CVPR'22)



SimMIM: A Simple Framework for Masked Image Modeling, Arxiv 2022.

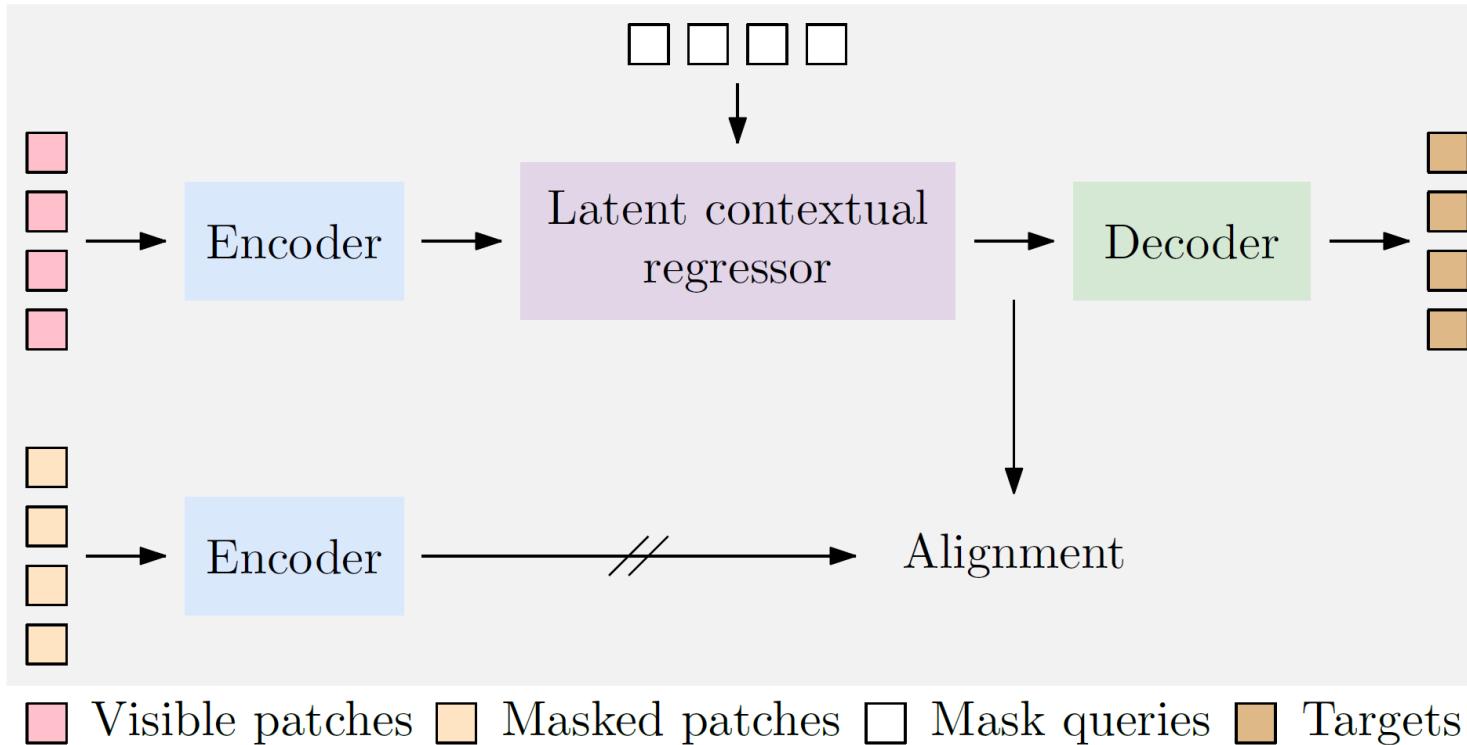
# SimMIM

- Moderately large masked patch size (e.g. 32)
- Predicting raw pixels of RGB values
- Light prediction head (e.g., a linear layer)



Methods	Input Size	Fine-tuning	Linear eval	Pre-training
		Top-1 acc (%)	Top-1 acc (%)	costs
Sup. baseline [44]	$224^2$	81.8	-	-
DINO [5]	$224^2$	82.8	78.2	$2.0\times$
MoCo v3 [9]	$224^2$	83.2	76.7	$1.8\times$
ViT [15]	$384^2$	79.9	-	$\sim 4.0\times$
BEiT [1]	$224^2$	83.2	56.7	$1.5\times^\dagger$
Ours	$224^2$	<b>83.8</b>	56.7	$1.0\times$

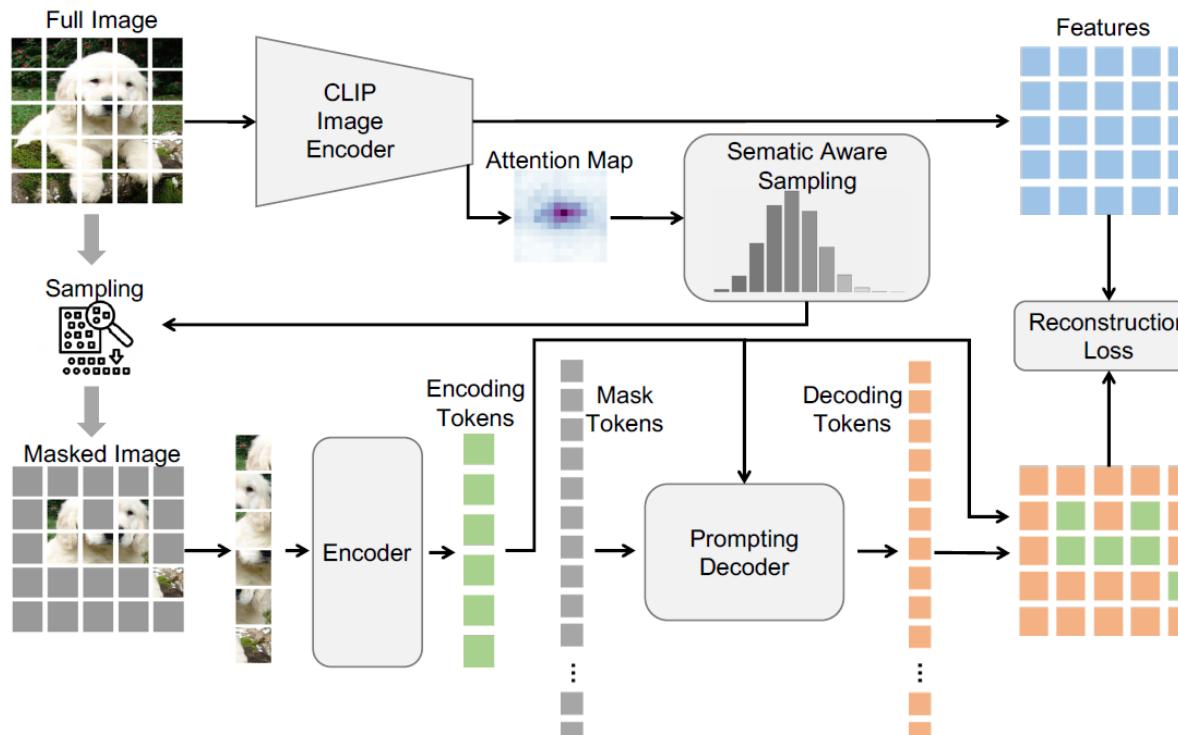
# Context Autoencoder (CAE)



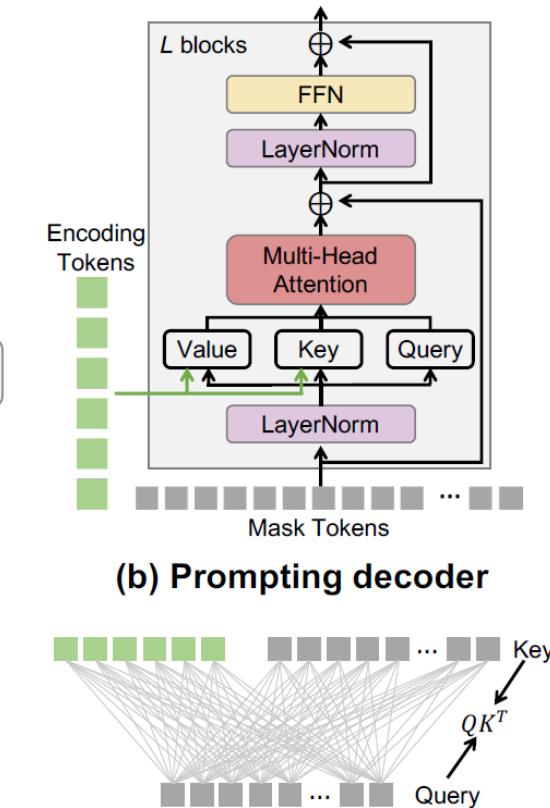
$$\ell_y(\mathbf{Y}_m, \bar{\mathbf{Y}}_m) + \lambda \ell_z(\mathbf{Z}_m, \bar{\mathbf{Z}}_m)$$

Context Autoencoder for Self-Supervised Representation Learning, Arxiv 2022.

# MILAN (NeurIPS'22)



(a) Overall flow of the MILAN framework



(b) Prompting decoder

(c) Attention in the prompting decoder

# MILAN (NeurIPS'22)

Method	Training data	Resolution	ViT-B/16		ViT-L/16	
			Epochs	Top-1 (%)	Epochs	Top-1 (%)
Supervised [50]	IN1K	224	-	83.8 (+1.6)	-	84.9 (+1.8)
<i>contrastive or clustering based</i>						
MoCov3 [11]	IN1K	224	300	83.2 (+2.2)	300	84.1 (+2.6)
DINO [6]	IN1K	224	400	82.8 (+2.6)	-	-
iBOT [69]	IN22K+IN1K	224	320	84.4 (+1.0)	200	86.3 (+0.4)
<i>reconstruction based</i>						
BEiT [3]	DALLE250M+IN22K+IN1K	224	150	83.7 (+1.7)	150	86.0 (+0.7)
mc-BEiT [33]	OpenImages9M+IN1K	224	800	84.1 (+1.3)	800	85.6 (+1.1)
PeCo [18]	IN1K	224	800	84.5 (+0.9)	800	86.5 (+0.2)
SimMIM [61]	IN1K	224	800	83.8 (+1.6)	-	-
MaskFeat [56]	IN1K	224	1600	84.0 (+1.4)	1600	85.7 (+1.0)
data2vec [2]	IN1K	224	800	84.2 (+1.2)	1600	86.6 (+0.1)
CAE [9]	IN1K	224	800	83.6 (+1.8)	-	-
MAE [25]	IN1K	224	1600	83.6 (+1.8)	1600	85.9 (+0.8)
<i>language-image pretraining based</i>						
CLIP [43]	OpenAI400M+IN1K	224	-	82.1 (+3.3)	-	85.3 (+1.4)
MVP [57]	OpenAI400M+IN1K	224	300	84.4 (+1.0)	300	86.3 (+0.4)
<b>MILAN</b>	OpenAI400M+IN1K	224	400	<b>85.4</b>	400	<b>86.7</b>
Supervised [19]	JFT300M+IN1K	384	90	84.2 (+2.2)	90	87.1 (+0.2)
BEiT [3]	DALLE250M+IN1K	384	800	84.6 (+1.8)	800	86.3 (+1.0)
SWAG [47]	IG3.6B+IN1K	384	2	85.3 (+1.1)	-	-
<b>MILAN</b>	OpenAI400M+IN1K	384	400	<b>86.4</b>	400	<b>87.3</b>

# Summary

- Before 2020, Self-supervised Contextual Modeling (百家争鸣)
- Contrastive learning-based self-supervised learning
  - Before 2021, CNN-based
  - 2021, transformer-based
- Masked Image Modeling
  - 2021, transformer-based