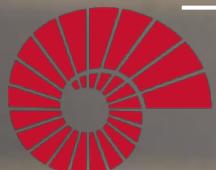


# COMP547

## DEEP UNSUPERVISED LEARNING

Lecture #11 – Self-Supervised Learning



KOÇ  
UNIVERSITY

Aykut Erdem // Koç University // Spring 2025

# Previously on COMP547

- Fundamentals of Diffusion Models
- Video Generation
- Evaluation



# Lecture overview

- Motivation
- Reconstruct from a corrupted (or partial) version
- Proxy tasks in computer vision
- Contrastive learning

**Disclaimer:** Much of the material and slides for this lecture were borrowed from  
—Pieter Abbeel, Peter Chen, Jonathan Ho, Aravind Srinivas' Berkeley CS294-158 class  
—Aaron Courville's Université de Montréal IFT6268 class

# Lecture overview

- Motivation
- Reconstruct from a corrupted (or partial) version
- Proxy tasks in computer vision
- Contrastive learning

# Course so far...

- Density modelling
  - Autoregressive, Normalizing Flows, Variational Inference
- Implicit models
  - Generative Adversarial Networks
- Diffusion models
  - Score-based Models, Denoising Diffusion Models
- Applications of generative modelling

# Today...

- How do we learn rich and useful features from raw unlabeled data that can be useful for several downstream tasks?
- What are the various pretext (proxy) tasks that can be used to learn representations from unlabeled data?
- How can we improve data-efficiency and performance of downstream tasks with a good pre-trained network?

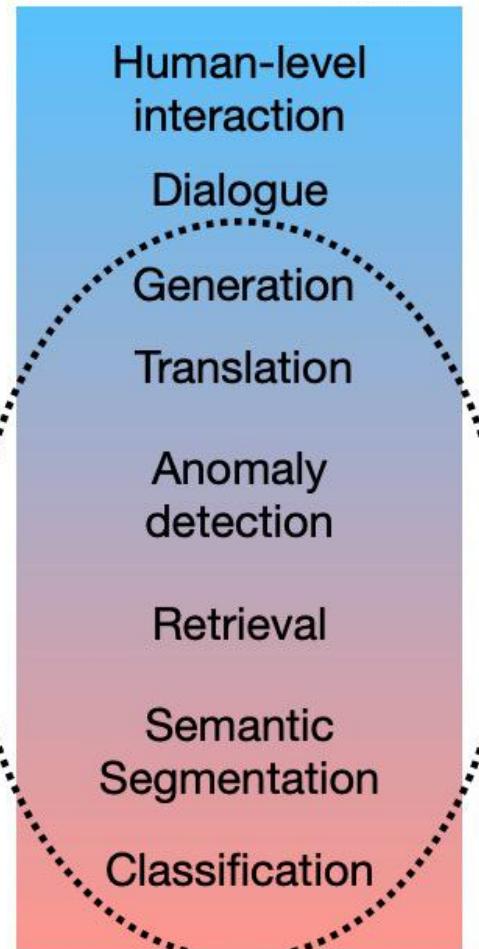
# Learning “really useful” representations

- Longstanding dream of the Deep Learning community:
  - Use unsupervised learning to learn some feature representation that can be used to support effective supervised learning (like classification)

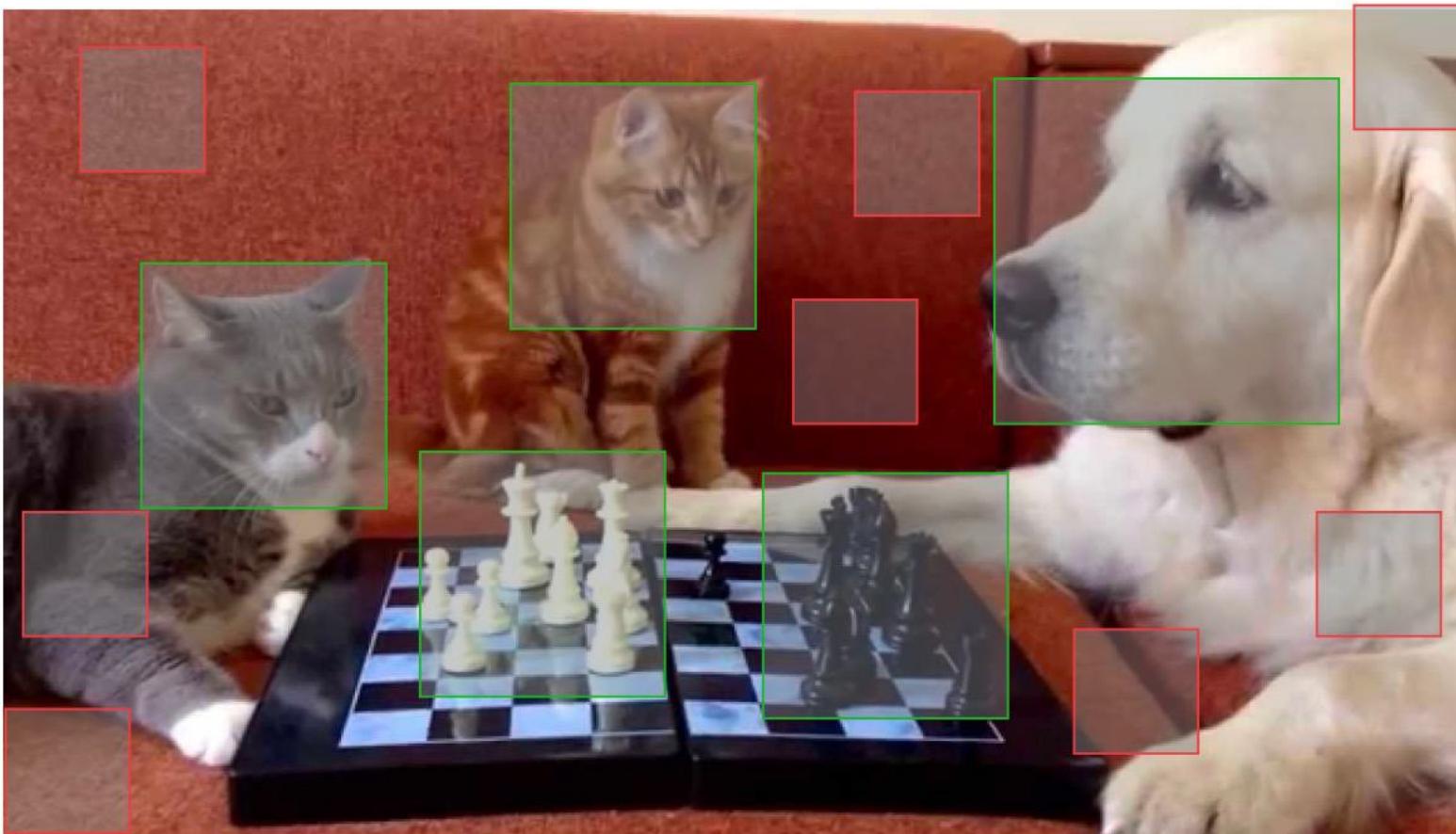
Input



Tasks



# We don't need generation/reconstruction



- Interesting thing
- Not interesting thing

- Generative models (in principle) care about all the pixels

# Representation Learning for Supervised Learning

- Using generative models (AEs, VAEs, etc) have largely been ineffective with two exceptions:
  1. Natural Language Modelling (all SOTA models are build on BERT-like representations)
  2. In the very small dataset regime, unsupervised learning can actually help.
- Gradient-based supervised training with the right model (e.g. CNNs for vision problems) has been very difficult to beat with unsupervised methods.

# It's worth asking ... Why?

A speculative answer:

- Most (essentially all) existing unsupervised methods learn features that are overwhelmingly low-level (non-semantic).
  - The features describe superficial aspects of the data and preserve few of the invariances that one would want from a representation learning scheme.
- Modern supervised learning methods (i.e. with NN) learn layers of representations that learn the relevant axes of variance in the data.
  - Eg. Higher level features of a CNN trained to recognize car makes and models should be relatively invariant to color but very sensitive to subtle differences in shape.

# Self-Supervised Learning

- A version of unsupervised learning where data provides the supervision.
- In general, withhold some part of the data and the task a neural network to predict it from the remaining parts.
- Details decide what proxy loss or pretext task the network tries to solve, and depending on the quality of the task, good semantic features can be obtained without actual labels.

# Motivation

- Supervised learning success story is heavily because of the utility of pre-trained classifier features for commercially useful downstream tasks like segmentation, detection, etc.
- Recipe is clear: Collect a large labeled dataset, train a model, deploy. Good data and sufficient data are what you need.
- Goal of self-supervised learning:
  - Learn equally good (if not better) features without supervision
  - Be able to deploy similar quality systems without relying on too many labels for the downstream tasks
  - Generalize better potentially because you learn more about the world

# How Much Information Does the Machine Need to Predict?

Y LeCun

## ■ "Pure" Reinforcement Learning (cherry)

- ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**

## ■ Supervised Learning (icing)

- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

## ■ Unsupervised/Predictive Learning (cake)

- ▶ The machine predicts any part of its input for any observed part.
- ▶ Predicts future frames in videos
- ▶ **Millions of bits per sample**

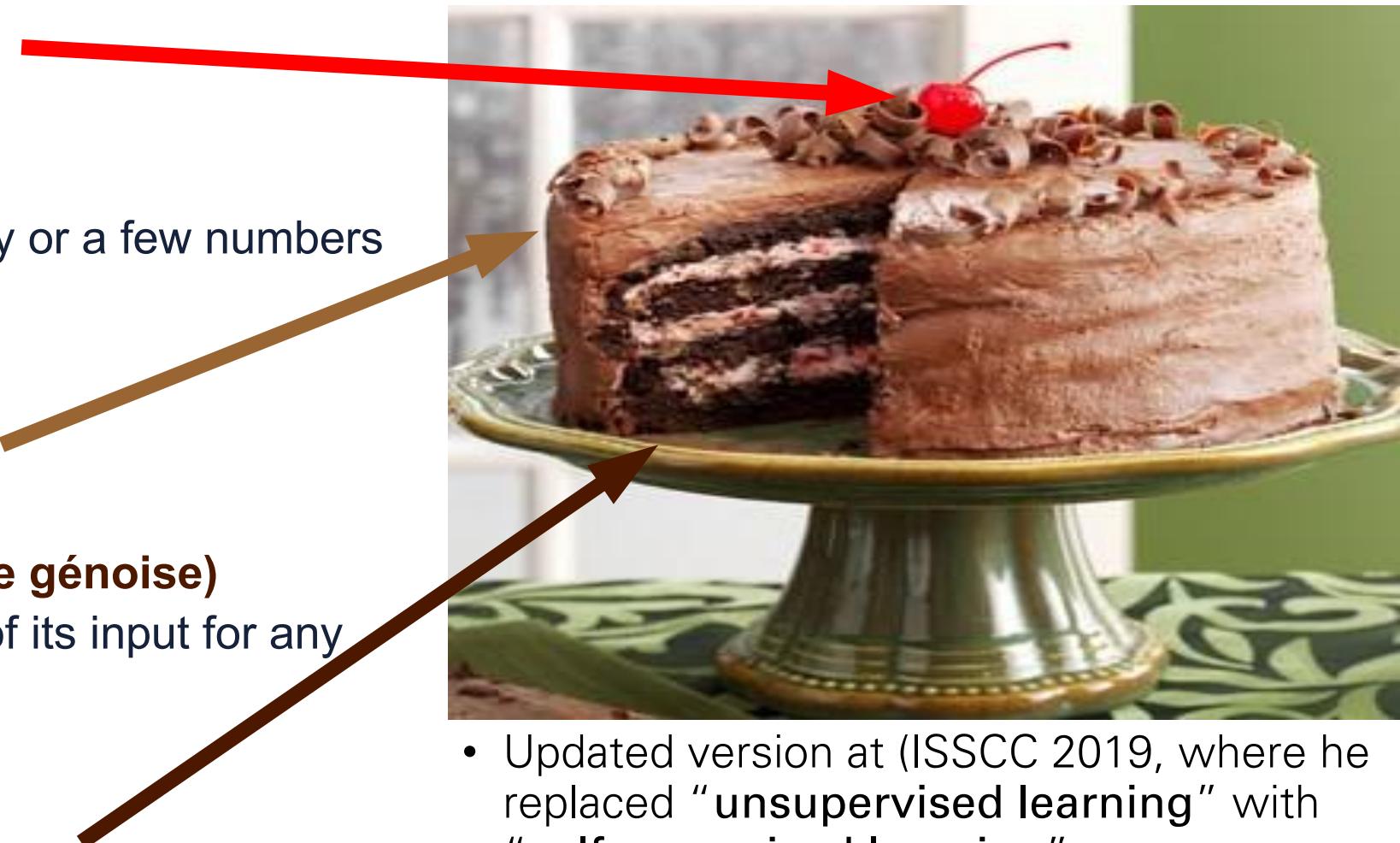


■ (Yes, I know, this picture is slightly offensive to RL folks. But I'll make it up)

- LeCun's original cake analogy slide, presented at his keynote speech in NIPS 2016.

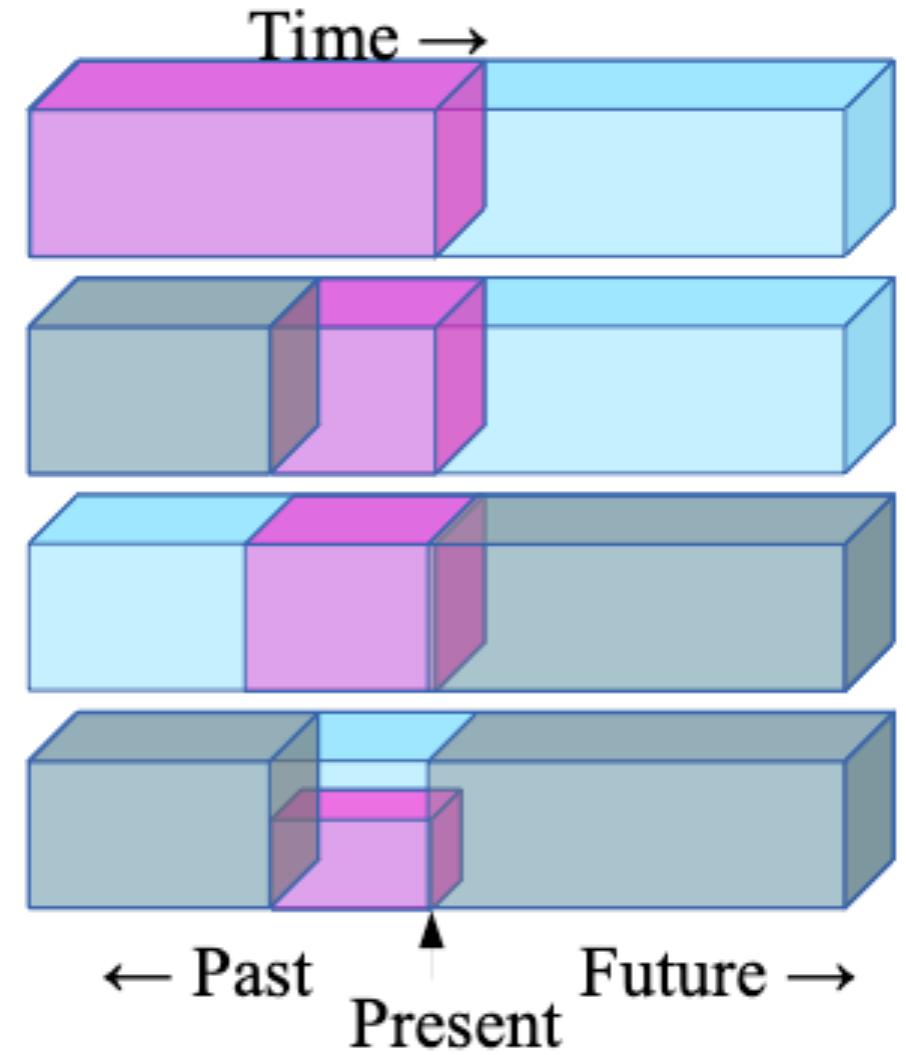
# How Much Information is the Machine Given during Learning?

- ▶ “Pure” Reinforcement Learning (**cherry**)
  - ▶ The machine predicts a scalar reward given once in a while.
  - ▶ **A few bits for some samples**
- ▶ Supervised Learning (**icing**)
  - ▶ The machine predicts a category or a few numbers for each input
  - ▶ Predicting human-supplied data
  - ▶ **10→10,000 bits per sample**
- ▶ Self-Supervised Learning (**cake génoise**)
  - ▶ The machine predicts any part of its input for any observed part.
  - ▶ Predicts future frames in videos
  - ▶ **Millions of bits per sample**



# Self-Supervised/Predictive Learning

- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.
- ▶ Predict the occluded from the visible
- ▶ **Pretend there is a part of the input you don't know and predict that.**



# What/Why Self-Supervision?

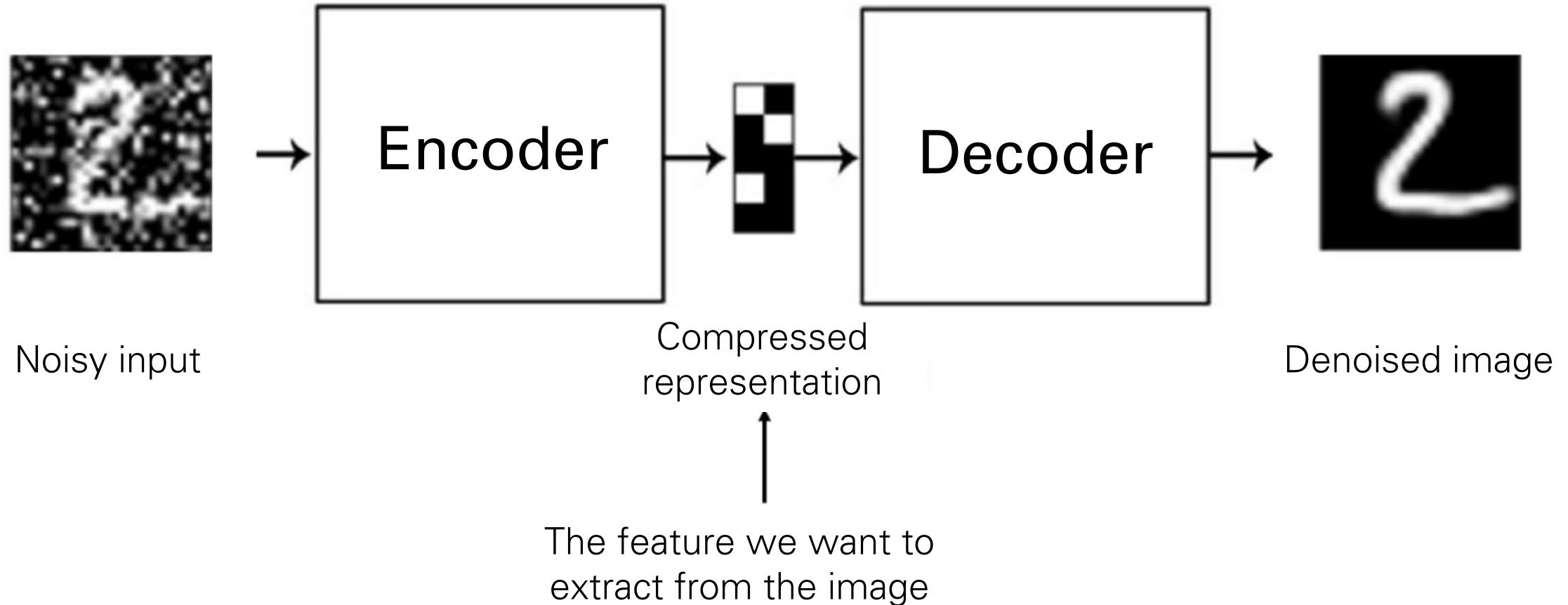
**Self-supervision:** Recover useful/semantic representations by training models to answer specific questions about the data.

- **Good:**
  - Can procedurally generate potentially infinite amounts of annotation.
  - We can borrow tricks from supervised learning without labels.
  - Focus on only the information that you need (e.g., not pixels).
  - Answering these questions requires more fundamental understanding of data.
- **Not so good:** designing good questions also requires some fundamental understanding of the data (e.g., structure).

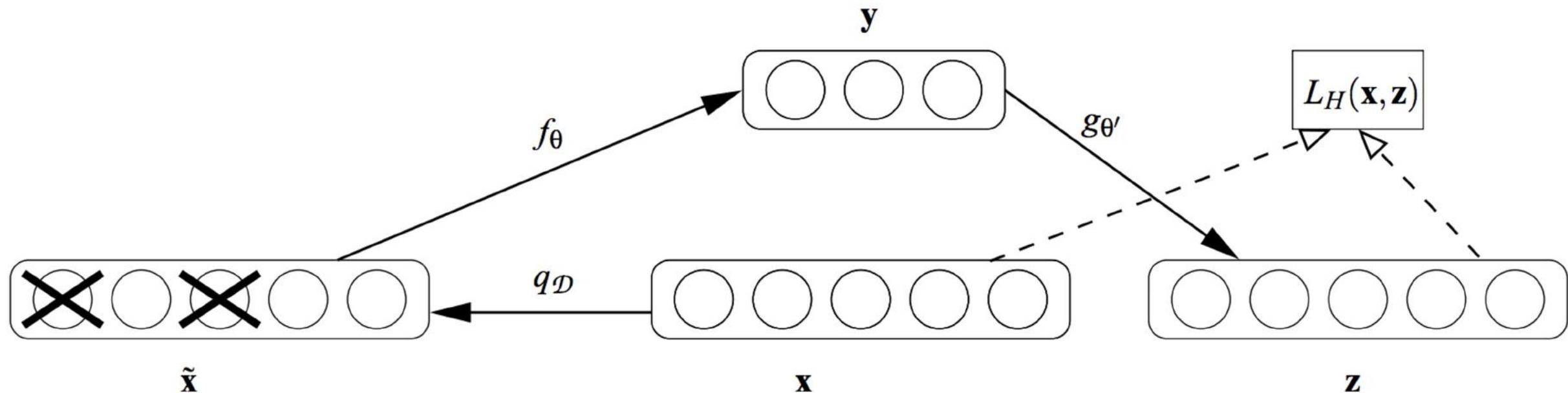
# Lecture overview

- Motivation
- Reconstruct from a corrupted (or partial) version
  - Denoising Autoencoder
  - In-painting / Masked AutoEncoder
  - Colorization, Split-Brain Autoencoder
- Proxy tasks in computer vision
- Contrastive Learning

# Denoising Autoencoder



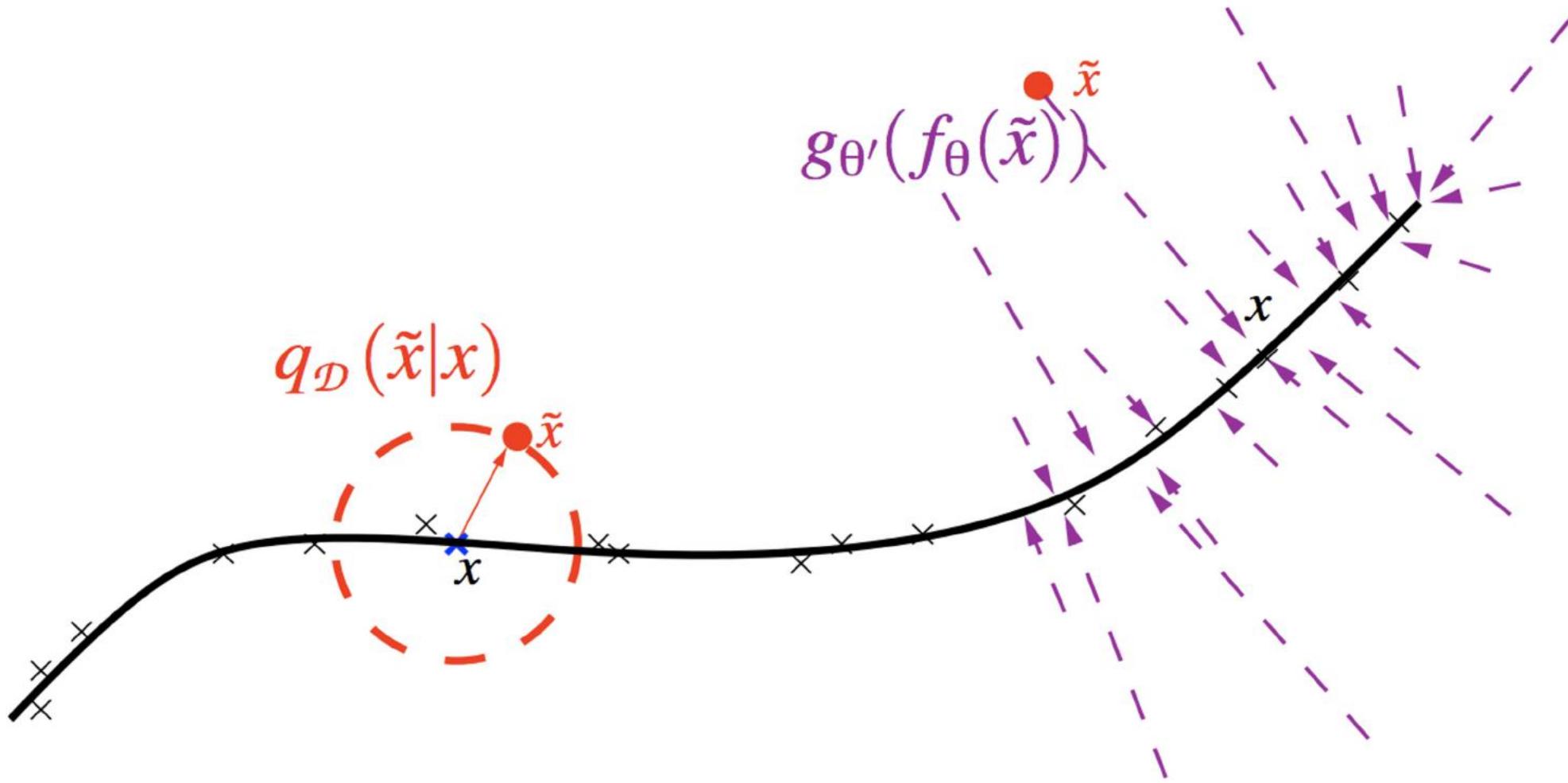
# Denoising Autoencoder



# Denoising Autoencoder

- Additive isotropic *Gaussian noise* (GS):  $\tilde{\mathbf{x}}|\mathbf{x} \sim \mathcal{N}(\mathbf{x}, \sigma^2 I)$ ;
- *Masking noise* (MN): a fraction  $v$  of the elements of  $\mathbf{x}$  (chosen at random for each example) is forced to 0;
- *Salt-and-pepper noise* (SP): a fraction  $v$  of the elements of  $\mathbf{x}$  (chosen at random for each example) is set to their minimum or maximum possible value (typically 0 or 1) according to a fair coin flip.

# Denoising Autoencoder

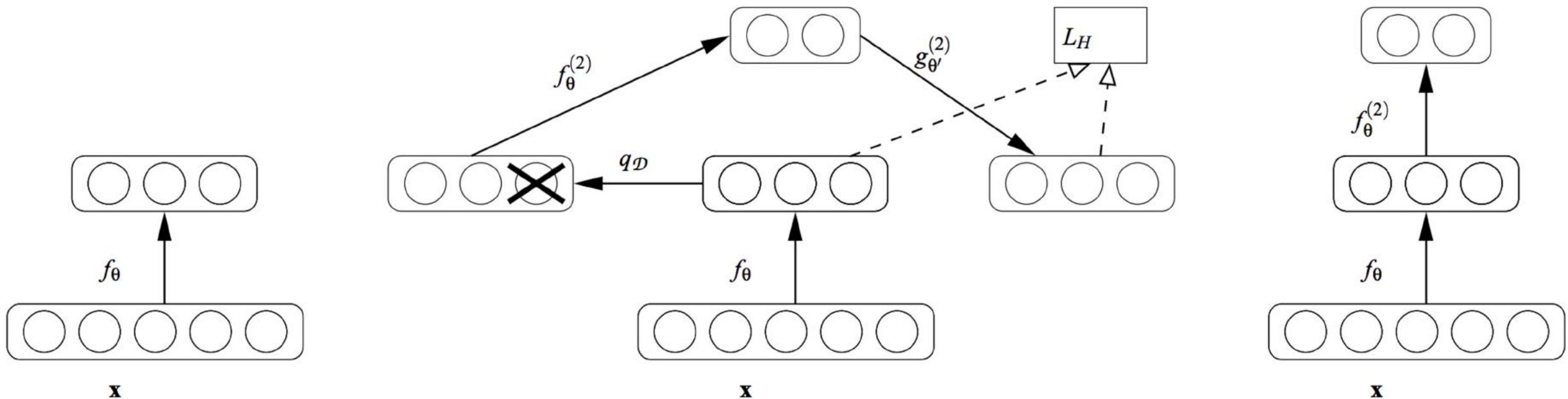


# Emphasizing corrupted dimensions

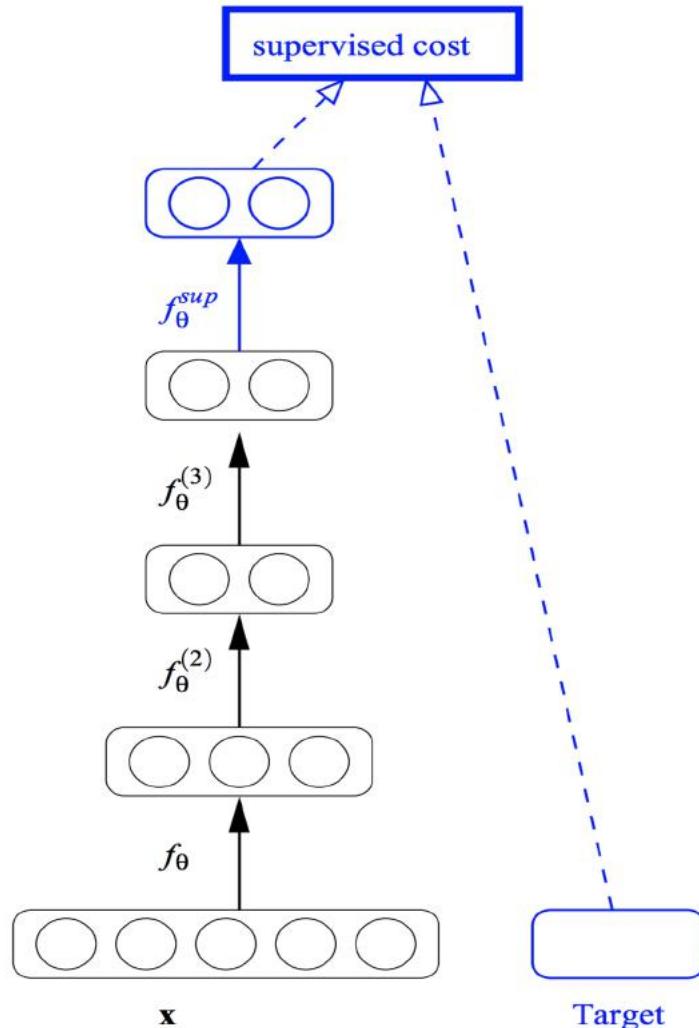
$$L_{2,\alpha}(\mathbf{x}, \mathbf{z}) = \alpha \left( \sum_{j \in \mathcal{J}(\tilde{\mathbf{x}})} (\mathbf{x}_j - \mathbf{z}_j)^2 \right) + \beta \left( \sum_{j \notin \mathcal{J}(\tilde{\mathbf{x}})} (\mathbf{x}_j - \mathbf{z}_j)^2 \right)$$

$$\begin{aligned} L_{\text{IH},\alpha}(\mathbf{x}, \mathbf{z}) &= \alpha \left( - \sum_{j \in \mathcal{J}(\tilde{\mathbf{x}})} [\mathbf{x}_j \log \mathbf{z}_j + (1 - \mathbf{x}_j) \log(1 - \mathbf{z}_j)] \right) \\ &\quad + \beta \left( - \sum_{j \notin \mathcal{J}(\tilde{\mathbf{x}})} [\mathbf{x}_j \log \mathbf{z}_j + (1 - \mathbf{x}_j) \log(1 - \mathbf{z}_j)] \right) \end{aligned}$$

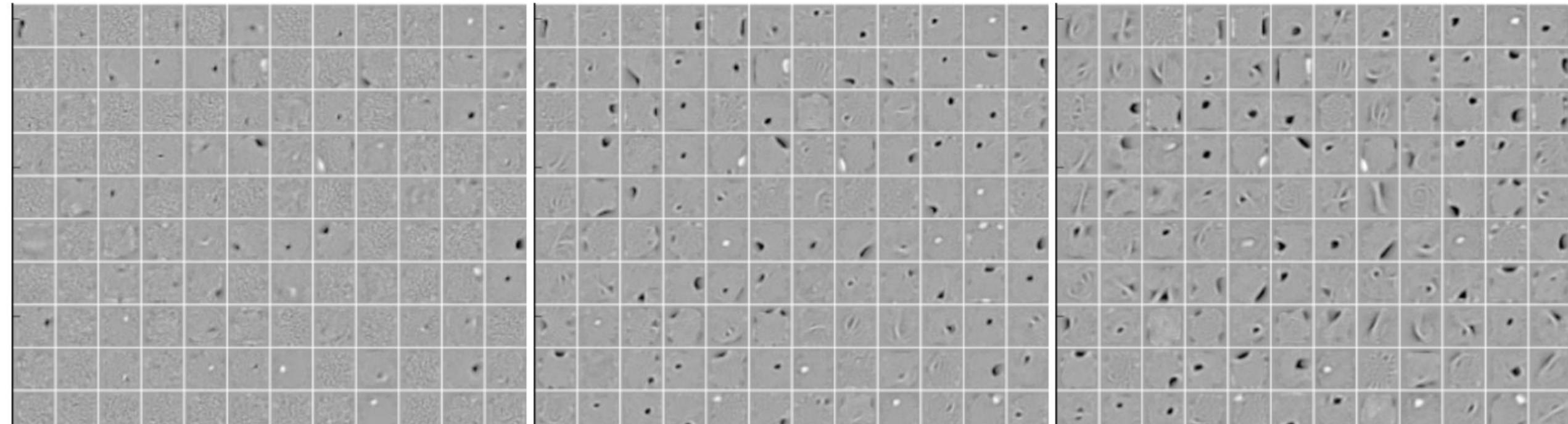
# Stacked Denoising Autoencoder



# Denoising Autoencoder



# Denoising Autoencoder



(a) No destroyed inputs

(b) 25% destruction

(c) 50% destruction

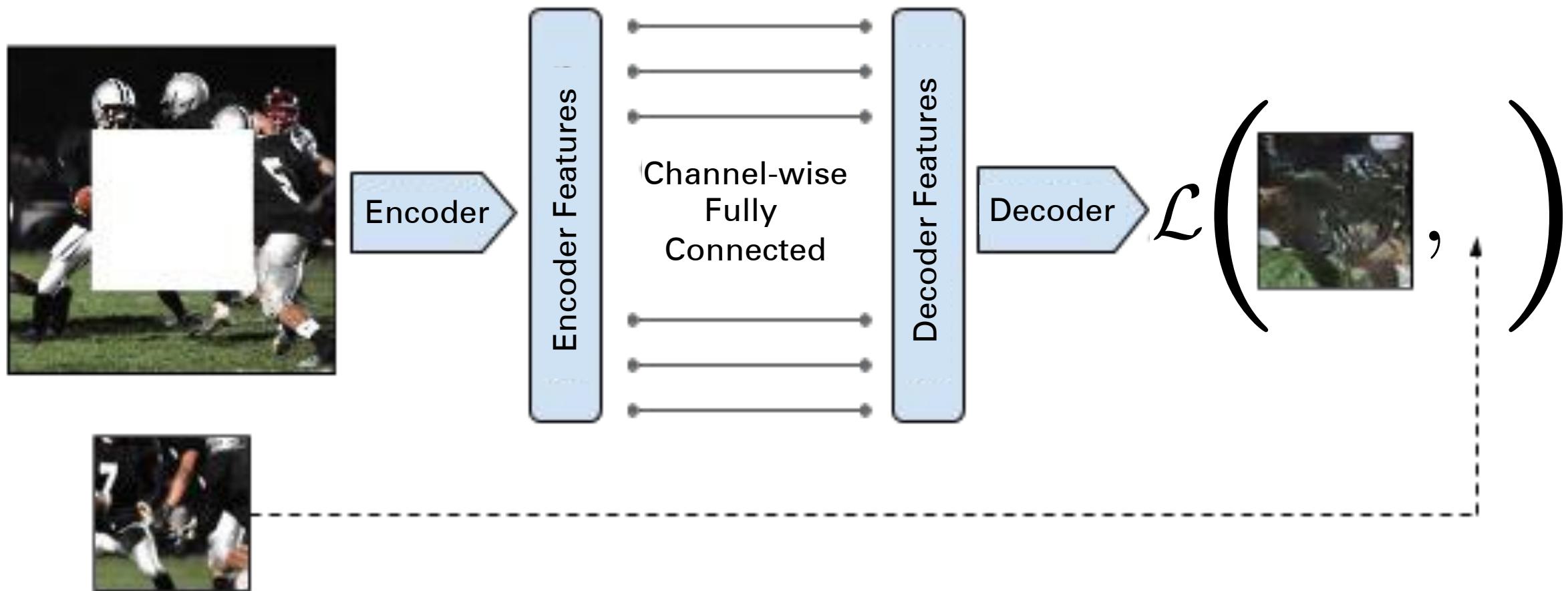
# Denoising Autoencoder

Dataset	SVM <sub>rbf</sub>	SVM <sub>poly</sub>	DBN-1	SAA-3	DBN-3	SdA-3 ( $\nu$ )
basic	<b>3.03±0.15</b>	3.69±0.17	3.94±0.17	3.46±0.16	3.11±0.15	<b>2.80±0.14</b> (10%)
rot	11.11±0.28	15.42±0.32	14.69±0.31	<b>10.30±0.27</b>	<b>10.30±0.27</b>	<b>10.29±0.27</b> (10%)
bg-rand	14.58±0.31	16.62±0.33	9.80±0.26	11.28±0.28	<b>6.73±0.22</b>	10.38±0.27 (40%)
bg-img	22.61±0.37	24.01±0.37	<b>16.15±0.32</b>	23.00±0.37	<b>16.31±0.32</b>	<b>16.68±0.33</b> (25%)
rot-bg-img	55.18±0.44	56.41±0.43	52.21±0.44	51.93±0.44	47.39±0.44	<b>44.49±0.44</b> (25%)
rect	<b>2.15±0.13</b>	<b>2.15±0.13</b>	4.71±0.19	2.41±0.13	2.60±0.14	<b>1.99±0.12</b> (10%)
rect-img	24.04±0.37	24.05±0.37	23.69±0.37	24.05±0.37	22.50±0.37	<b>21.59±0.36</b> (25%)
convex	19.13±0.34	19.82±0.35	19.92±0.35	<b>18.41±0.34</b>	<b>18.63±0.34</b>	<b>19.06±0.34</b> (10%)

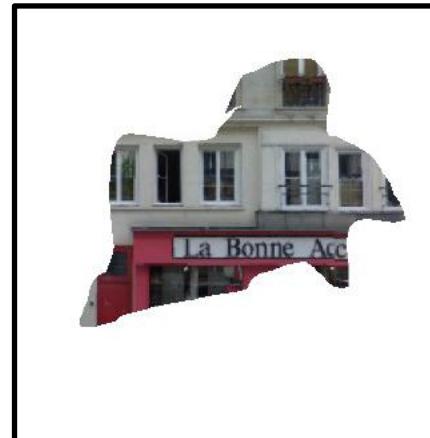
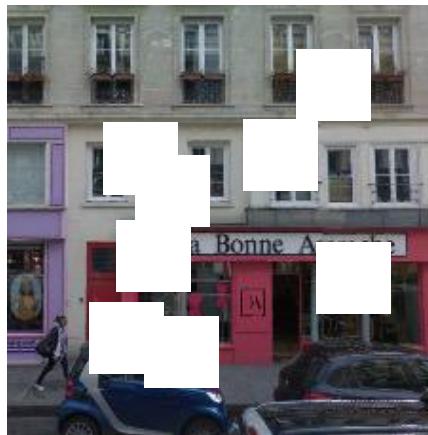
# Predict missing pieces



# Context Encoders



# Context Encoders



(a) Center Region

(b) Random Blocks

(c) Random Shapes

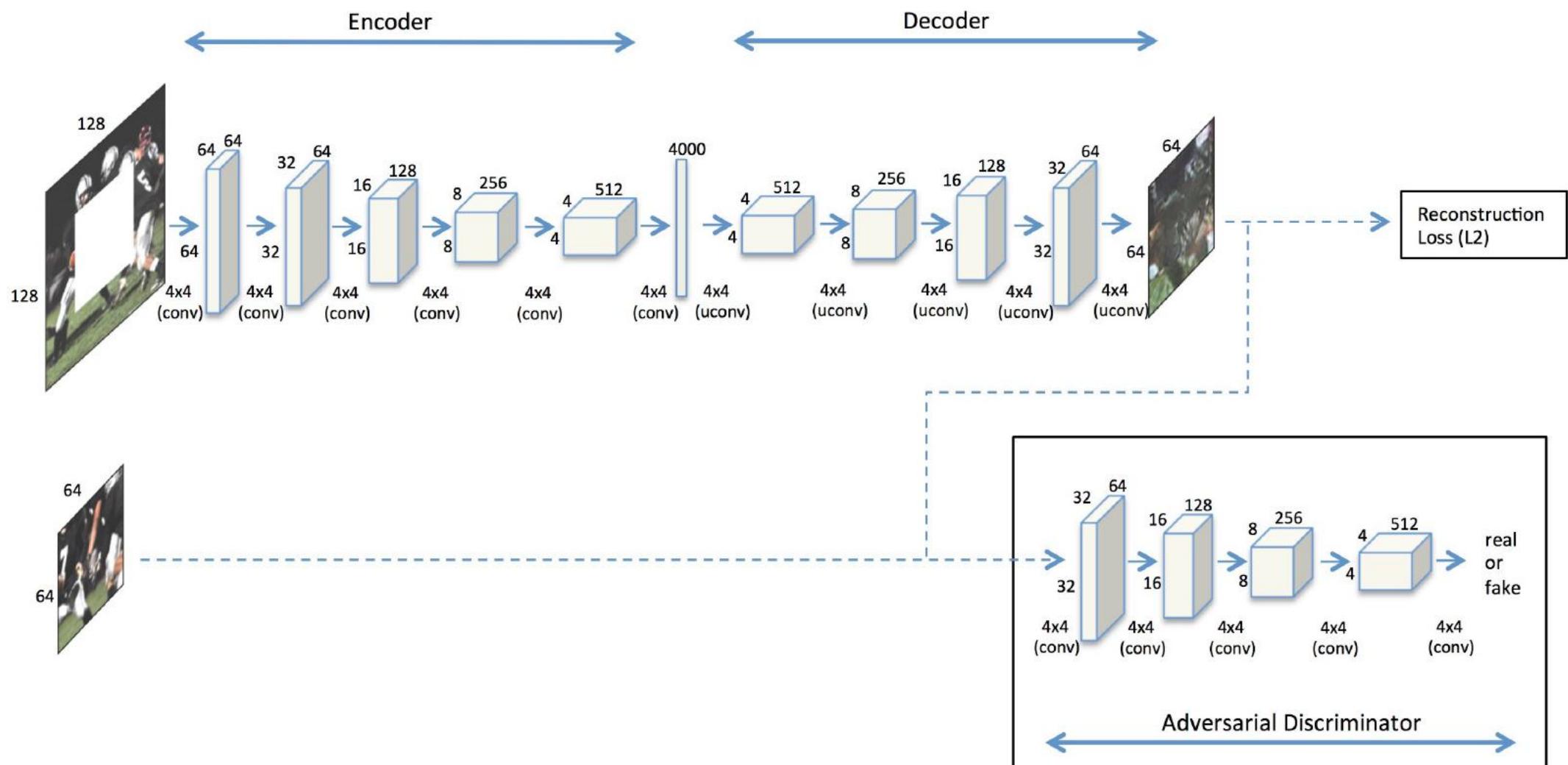
# Context Encoders

$$\mathcal{L}_{rec}(x) = \|\hat{M} \odot (x - F((1 - \hat{M}) \odot x))\|_2^2$$

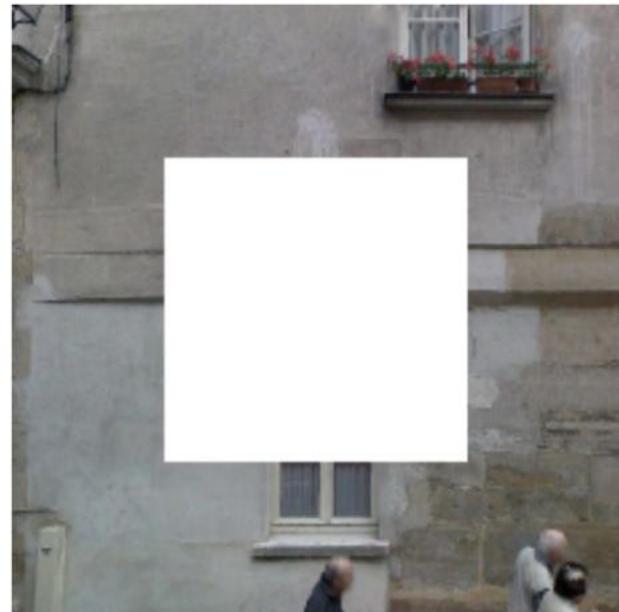
$$\begin{aligned}\mathcal{L}_{adv} = \max_D \mathbb{E}_{x \in \mathcal{X}} & [\log(D(x)) \\ & + \log(1 - D(F((1 - \hat{M}) \odot x)))]\end{aligned}$$

$$\mathcal{L} = \lambda_{rec} \mathcal{L}_{rec} + \lambda_{adv} \mathcal{L}_{adv}$$

# Context Encoders



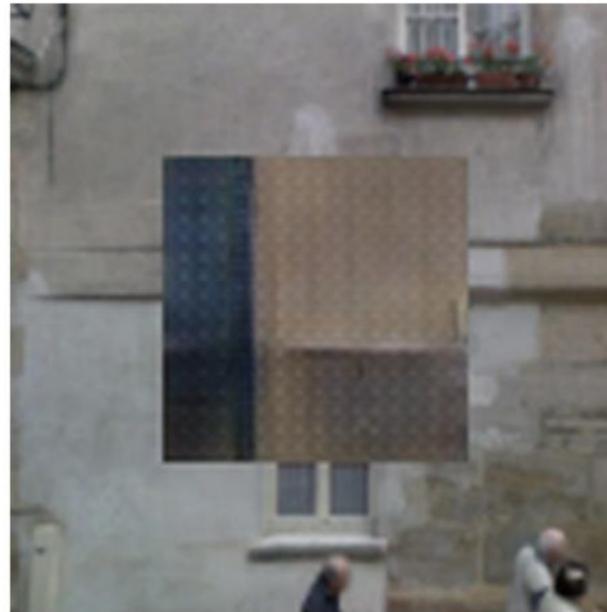
# Context Encoders



Input Image



L2 Loss



Adversarial Loss



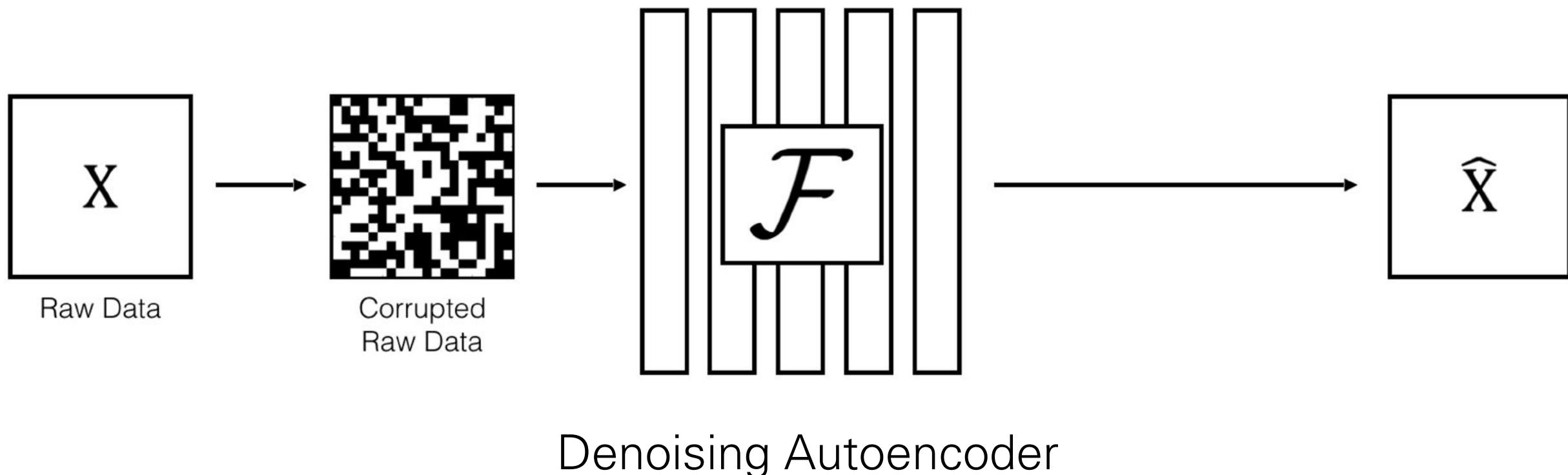
Joint Loss

# Context Encoders

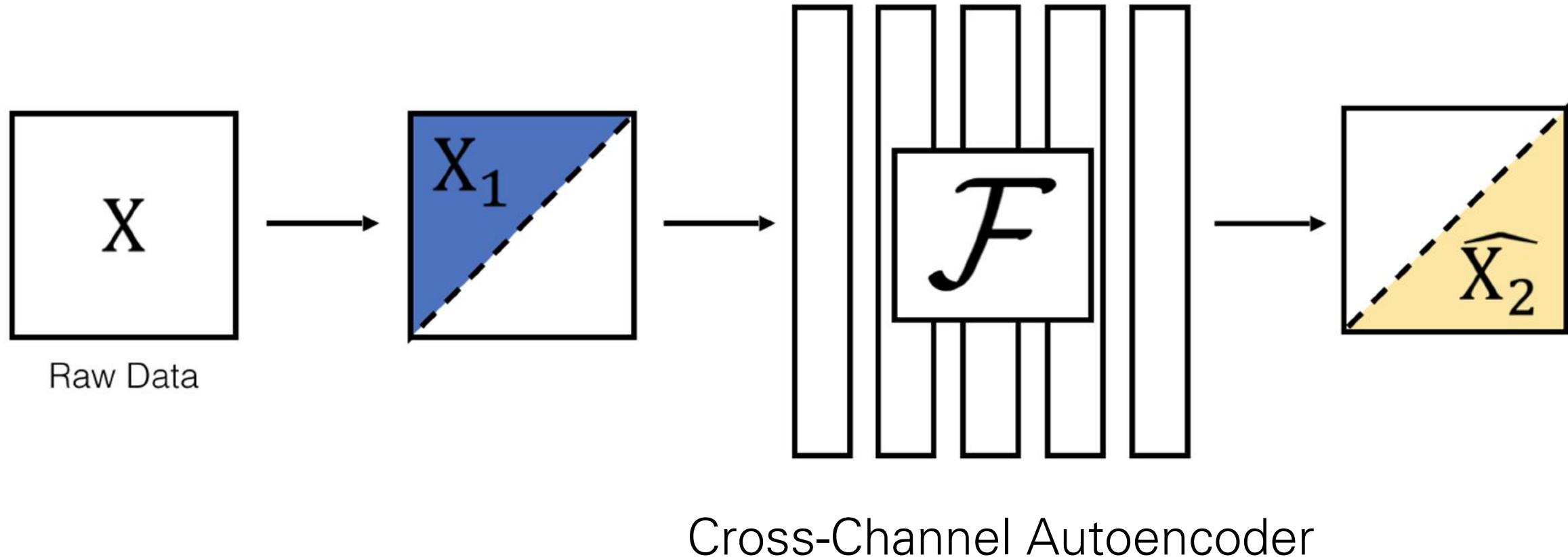
Pretraining Method	Supervision	Pretraining time	Classification	Detection	Segmentation
ImageNet [26]	1000 class labels	3 days	<b>78.2%</b>	<b>56.8%</b>	<b>48.0%</b>
Random Gaussian	initialization	< 1 minute	53.3%	43.4%	19.8%
Autoencoder	-	14 hours	53.8%	41.9%	25.2%
Agrawal <i>et al.</i> [1]	egomotion	10 hours	52.9%	41.8%	-
Doersch <i>et al.</i> [7]	context	4 weeks	55.3%	<b>46.6%</b>	-
Wang <i>et al.</i> [39]	motion	1 week	<b>58.4%</b>	44.0%	-
Ours	context	14 hours	56.5%	44.5%	<b>29.7%</b>

Table 2: Quantitative comparison for classification, detection and semantic segmentation. Classification and Fast-RCNN Detection results are on the PASCAL VOC 2007 test set. Semantic segmentation results are on the PASCAL VOC 2012 validation set from the FCN evaluation described in Section 5.2.3, using the additional training data from [18], and removing overlapping images from the validation set [28].

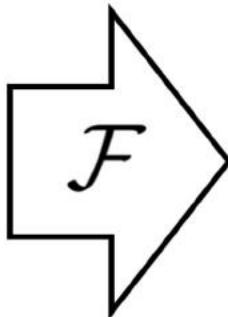
# Predicting one view from another



# Predicting one view from another

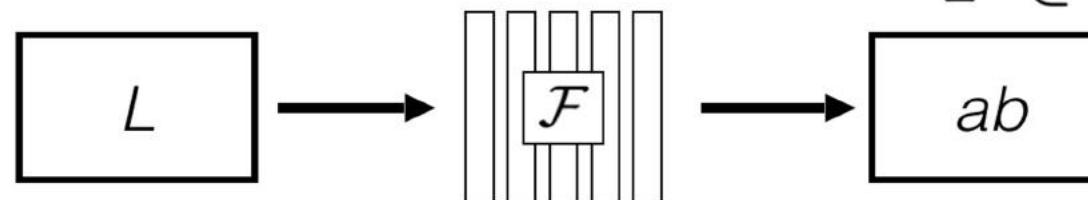


# Predicting one view from another



Grayscale image:  $L$  channel

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



Color information:  $ab$  channels

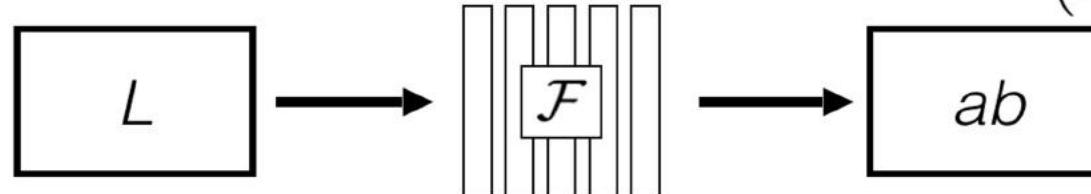
$$\hat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 2}$$

# Predicting one view from another



Grayscale image:  $L$  channel

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



Concatenate  $(L, ab)$  channels  
 $(\mathbf{X}, \hat{\mathbf{Y}})$

# Predicting one view from another



Ground Truth

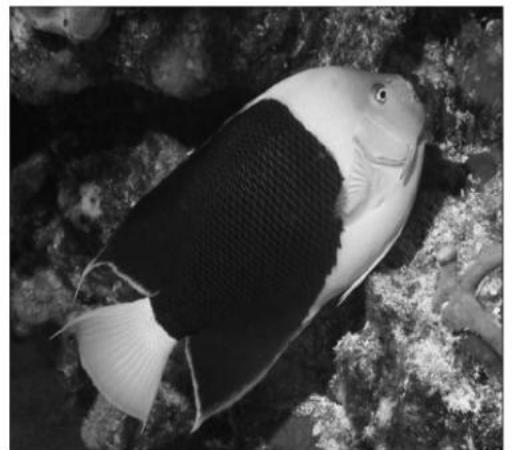


L2 regression



Pixelwise  
classification

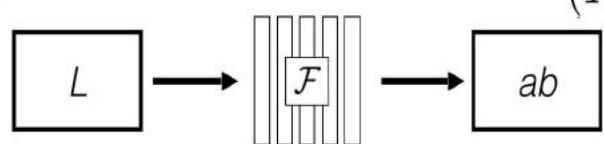
# Predicting one view from another



Grayscale image:  $L$  channel  
 $\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$



Concatenate ( $L, ab$ ) channels  
 $(\mathbf{X}, \hat{\mathbf{Y}})$



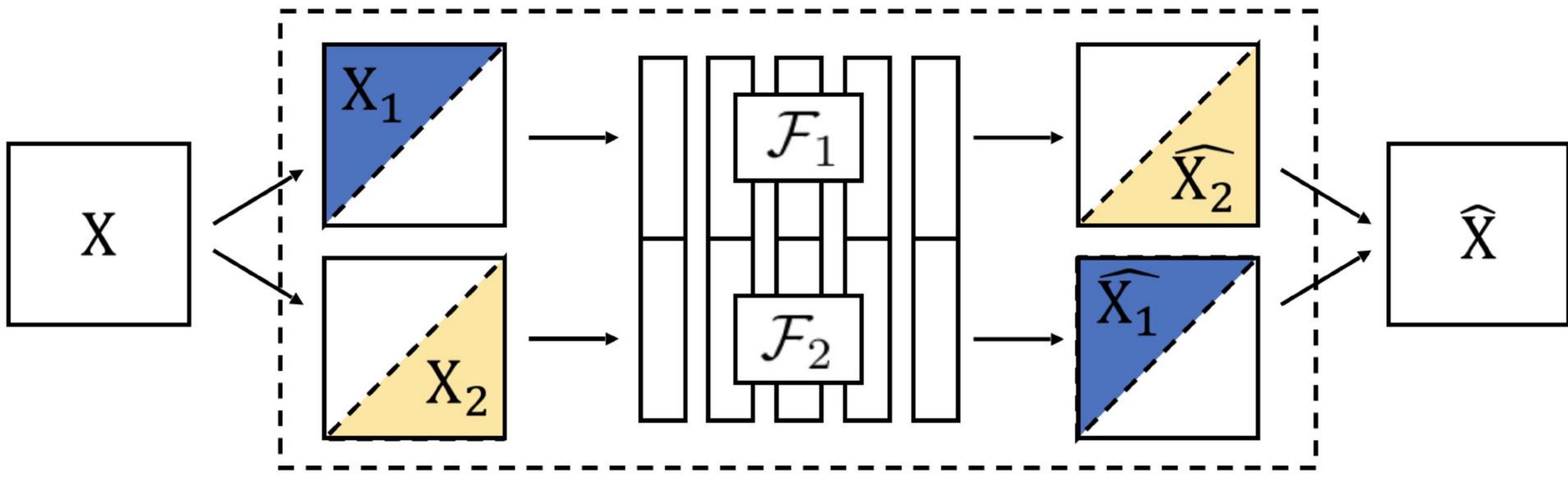
$$L_2(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2} \sum_{h,w} \|\mathbf{Y}_{h,w} - \hat{\mathbf{Y}}_{h,w}\|_2^2$$

6



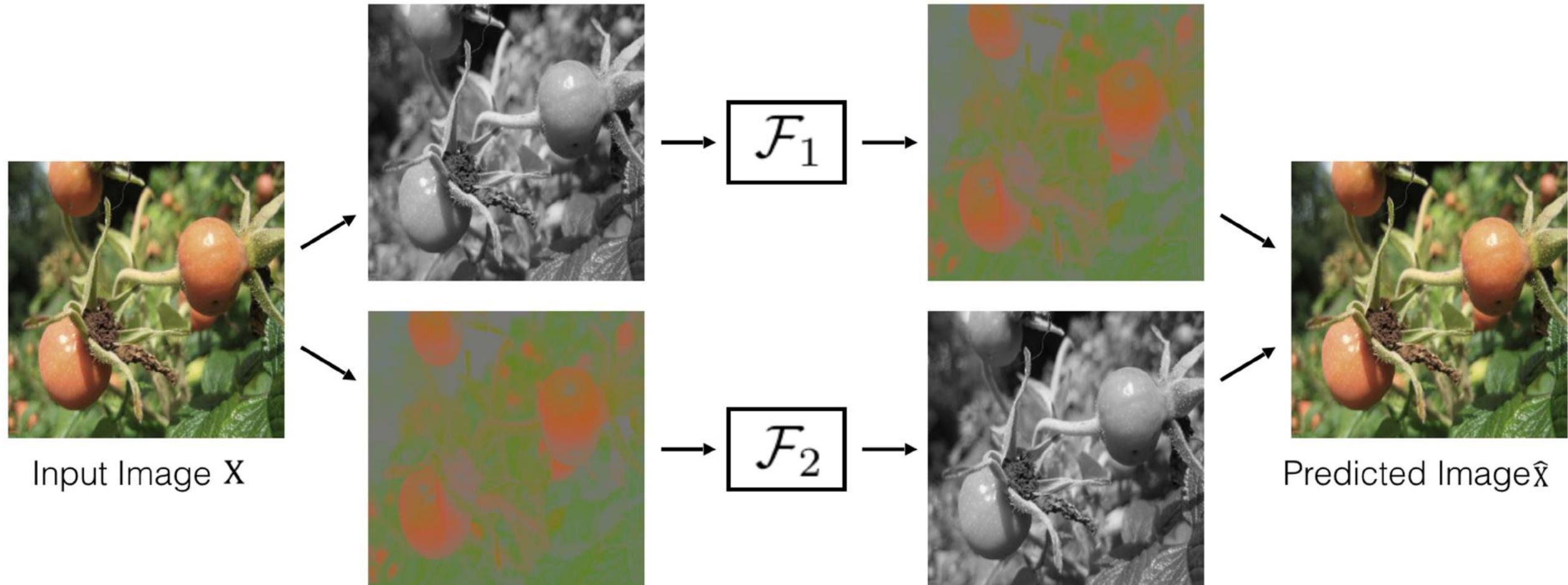
$$L(\hat{\mathbf{Z}}, \mathbf{z}) = -\frac{1}{HW} \sum_{h,w} \sum_q \mathbf{z}_{h,w,q} \log(\hat{\mathbf{z}}_{h,w,q})$$

# Predicting one view from another

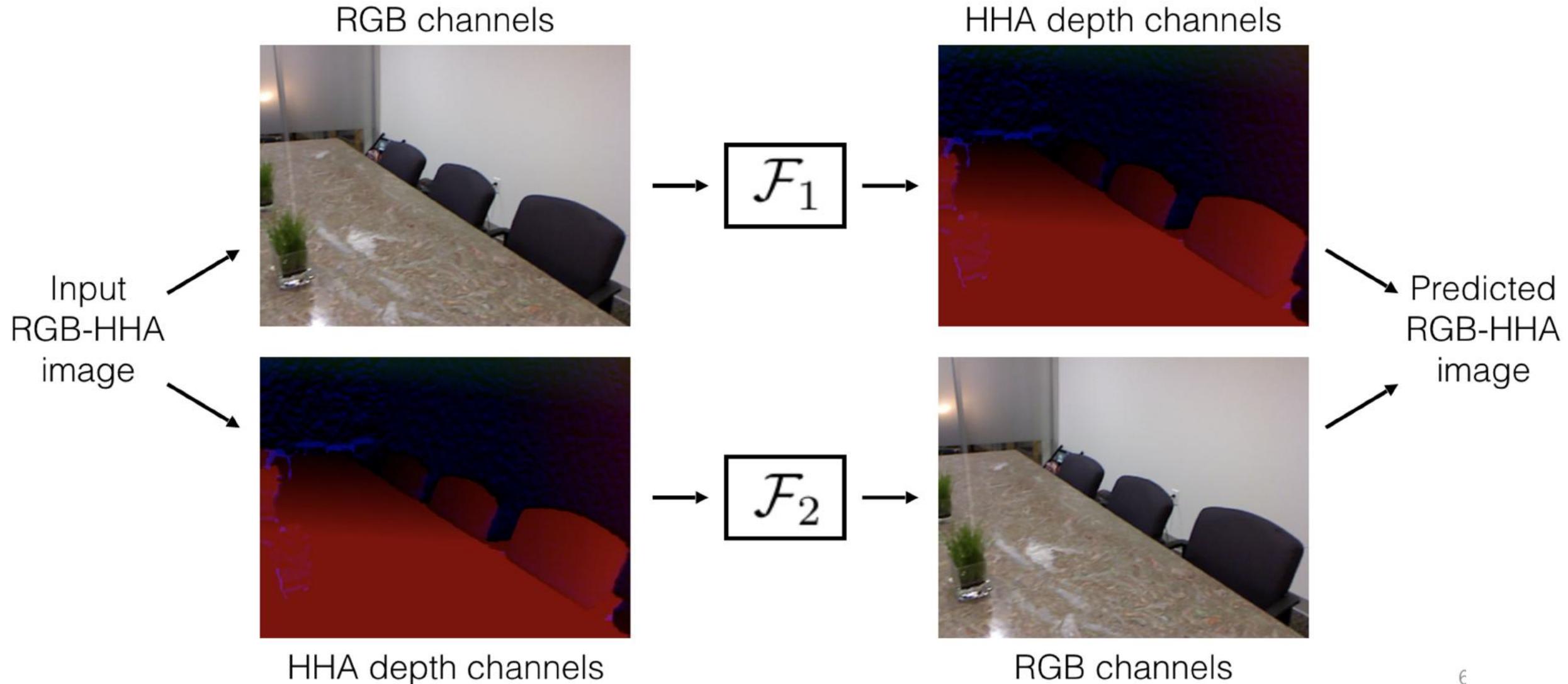


Split-Brain Autoencoder

# Predicting one view from another



# Predicting one view from another



€

# Masked Autoencoders (MAEs)

## **Masked Autoencoders Are Scalable Vision Learners**

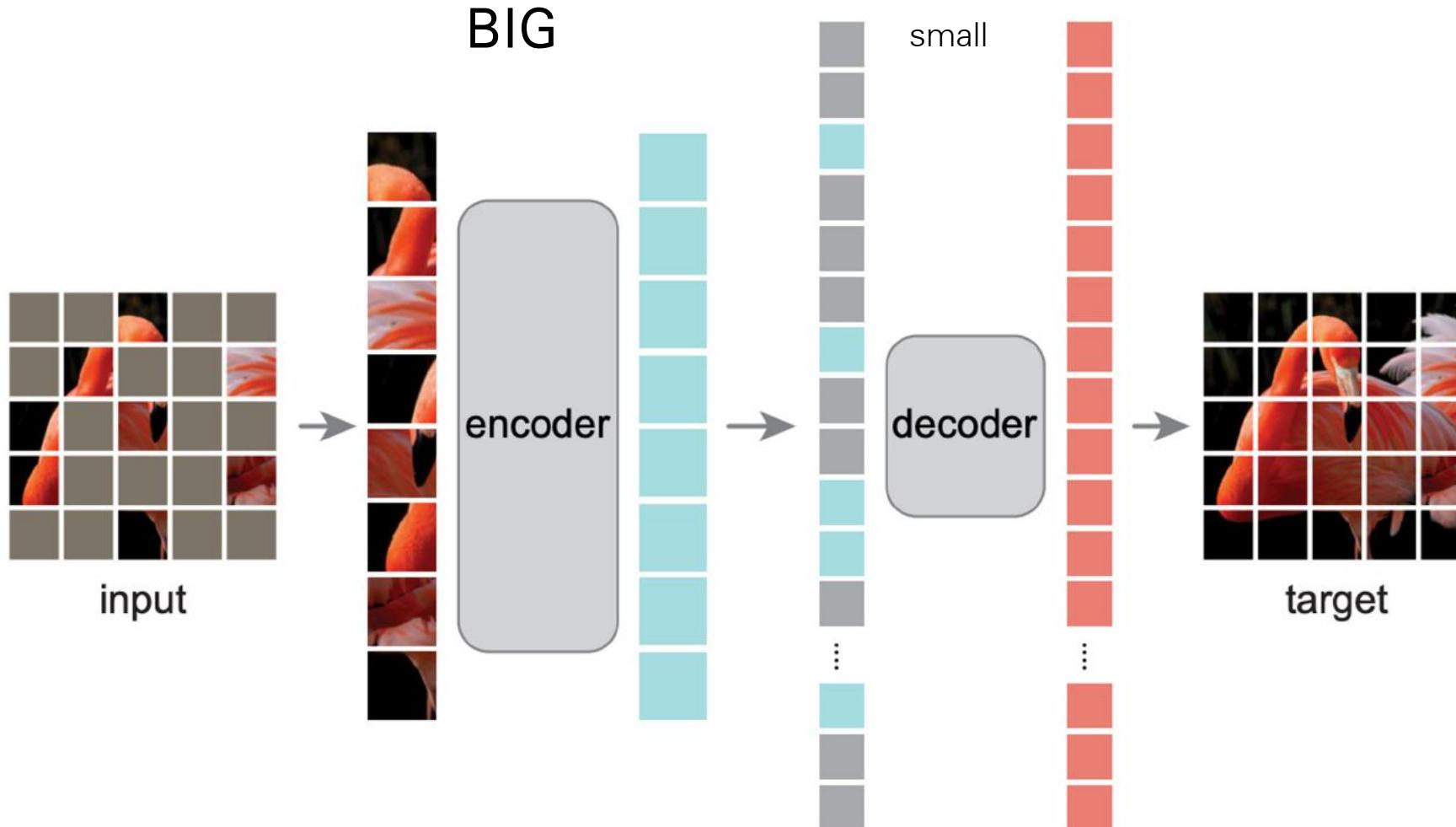
Kaiming He<sup>\*,†</sup> Xinlei Chen<sup>\*</sup> Saining Xie Yanghao Li Piotr Dollár Ross Girshick

<sup>\*</sup>equal technical contribution <sup>†</sup>project lead

Facebook AI Research (FAIR)

Nov, 2021

# Masked Autoencoders (MAEs)



Architecture: Vision Transformer (ViT)

# MAE on ImageNet validation images

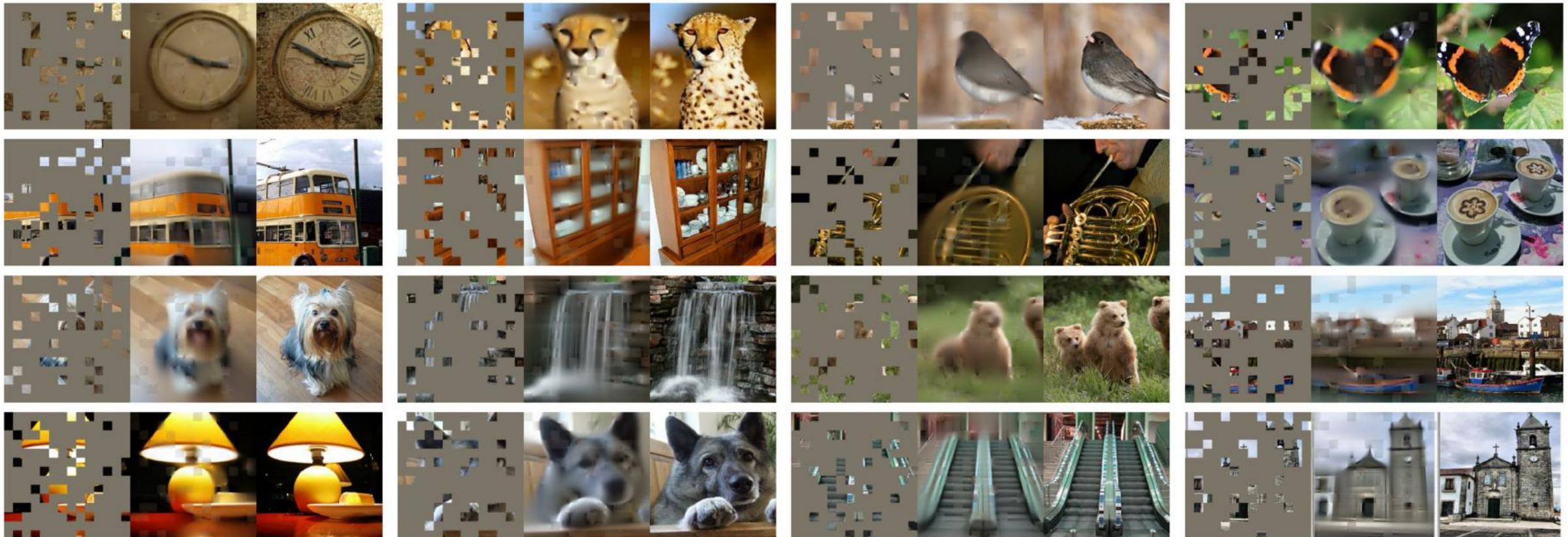


Figure 2. Example results on ImageNet *validation* images. For each triplet, we show the masked image (left), our MAE reconstruction<sup>†</sup> (middle), and the ground-truth (right). The masking ratio is 80%, leaving only 39 out of 196 patches. More examples are in the appendix.

<sup>†</sup>As no loss is computed on visible patches, the model output on visible patches is qualitatively worse. One can simply overlay the output with the visible patches to improve visual quality. We intentionally opt not to do this, so we can more comprehensively demonstrate the method’s behavior.

# MAE on CoCo validation images



Figure 3. Example results on COCO validation images, using an MAE trained on ImageNet (the same model weights as in Figure 2). Observe the reconstructions on the two right-most examples, which, although different from the ground truth, are semantically plausible.

# Masking Ratio

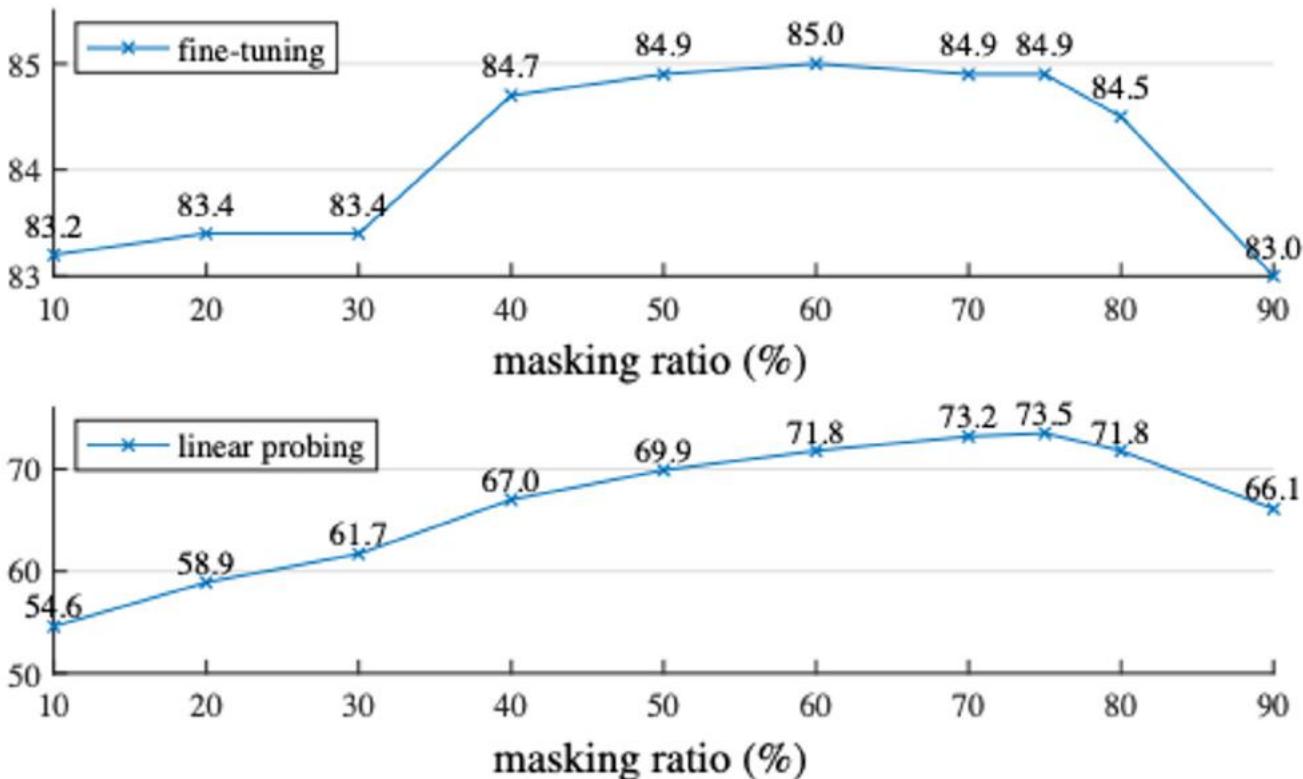


Figure 5. **Masking ratio.** A high masking ratio (75%) works well for both fine-tuning (top) and linear probing (bottom). The y-axes are ImageNet-1K validation accuracy (%) in all plots in this paper.

# Comparison with Prior SOTA

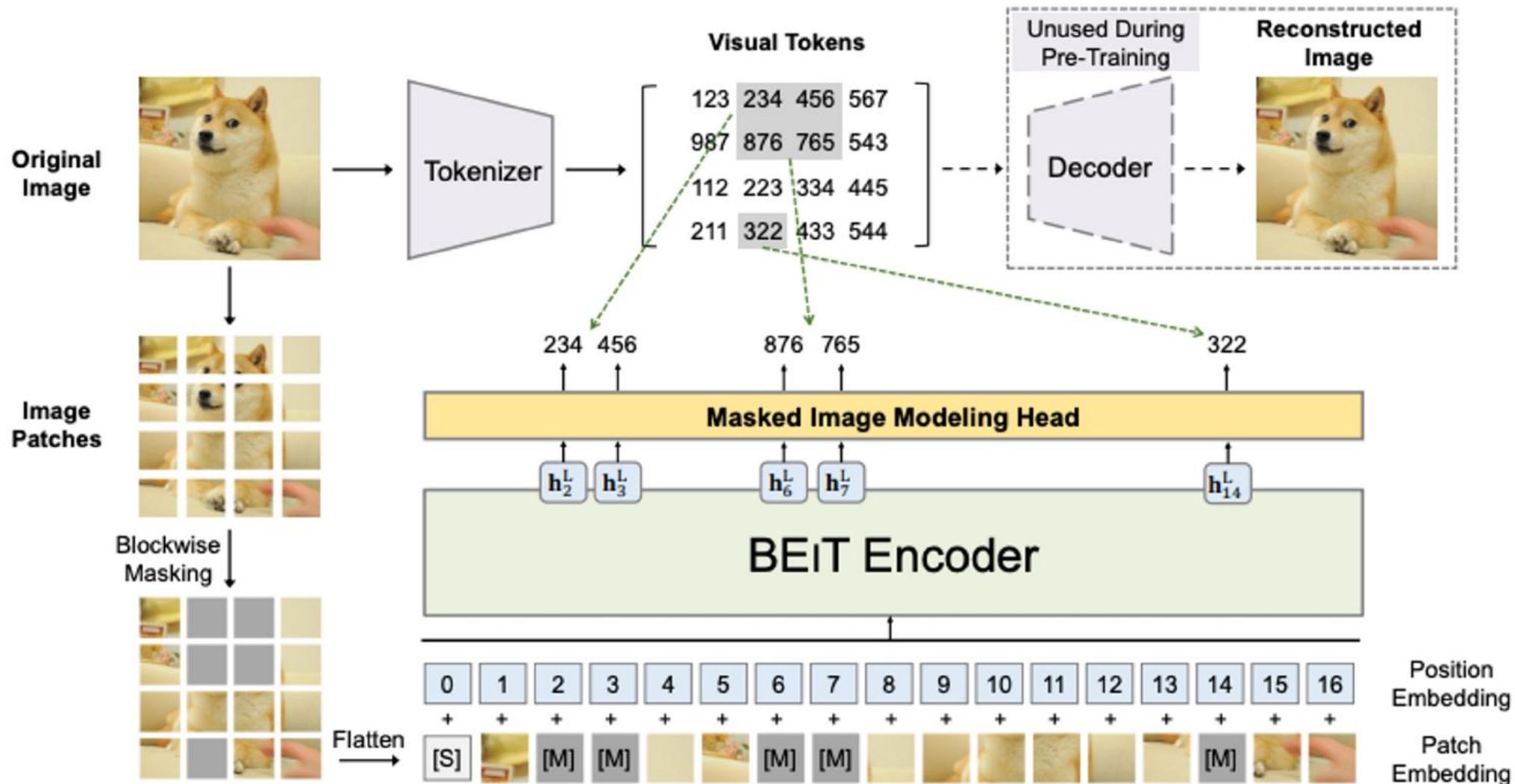
method	pre-train data	ViT-B	ViT-L	ViT-H	ViT-H <sub>448</sub>
scratch, our impl.	-	82.3	82.6	83.1	-
DINO [5]	IN1K	<u>82.8</u>	-	-	-
MoCo v3 [9]	IN1K	83.2	<u>84.1</u>	-	-
BEiT [2]	IN1K+DALLE	83.2	<u>85.2</u>	-	-
MAE	IN1K	<u>83.6</u>	<u>85.9</u>	<u>86.9</u>	<u>87.8</u>

Table 3. **Comparisons with previous results on ImageNet-1K.** The pre-training data is the ImageNet-1K training set (except the tokenizer in BEiT was pre-trained on 250M DALLE data [50]). All self-supervised methods are evaluated by end-to-end fine-tuning. The ViT models are B/16, L/16, H/14 [16]. The best for each column is underlined. All results are on an image size of 224, except for ViT-H with an extra result on 448. Here our MAE reconstructs normalized pixels and is pre-trained for 1600 epochs.

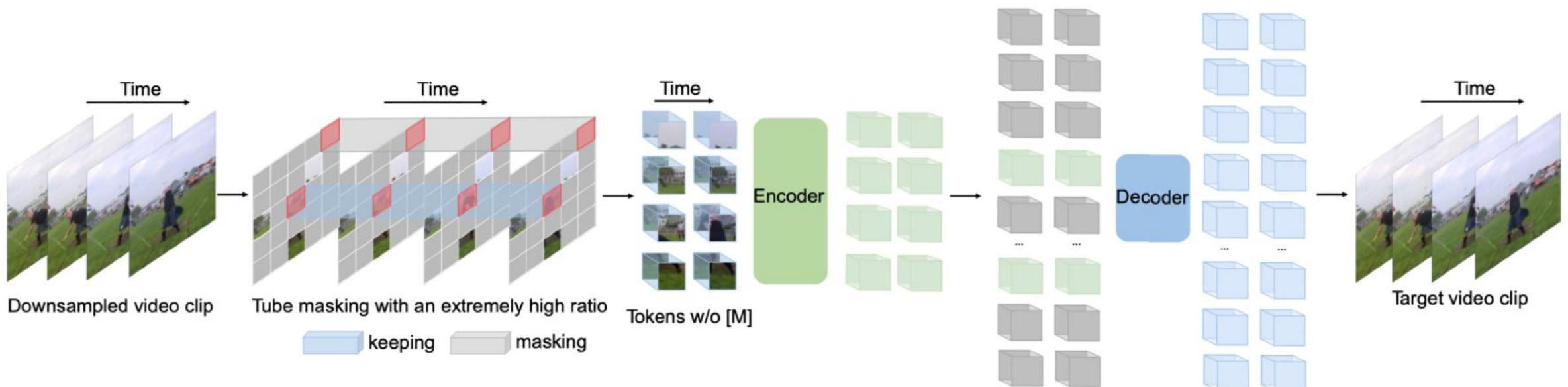
# MAE Cousins / Derivatives

- BeIT
- VideoMAE
- SiamMAE
- Audio-MAE
- M3AE
- MultiMAE
- Multi-View / Masked World Models for Visual Control

# BEiT Architecture



# VideoMAE



# Observations

- High masking ratio: 90% to 95%
- Impressive results even on very small datasets, e.g. 3k videos
- Data quality is more important than data quantity for Self Supervised Video Pretraining. Domain shift between pre-training and target datasets is an important factor.
- VideoMAE with the vanilla ViT backbone can achieve 87.4% on Kinects-400, 75.4% on SomethingSomething V2, 91.3% on UCF101, and 62.6% on HMDB51, without using any extra data.

# Experiments on Something-Something V2

Method	Backbone	Extra data	Ex. labels	Frames	GFLOPs	Param	Top-1	Top-5
TEINet <sub>En</sub> [40]	ResNet50 <sub>×2</sub>		✓	8+16	99×10×3	50	66.5	N/A
TANet <sub>En</sub> [41]	ResNet50 <sub>×2</sub>	ImageNet-1K	✓	8+16	99×2×3	51	66.0	90.1
TDN <sub>En</sub> [75]	ResNet101 <sub>×2</sub>		✓	8+16	198×1×3	88	69.6	92.2
SlowFast [23]	ResNet101		✓	8+32	106×1×3	53	63.1	87.6
MViTv1 [22]	MViTv1-B	Kinetics-400	✓	64	455×1×3	37	67.7	90.9
TimeSformer [6]	ViT-B		✓	8	196×1×3	121	59.5	N/A
TimeSformer [6]	ViT-L	ImageNet-21K	✓	64	5549×1×3	430	62.4	N/A
ViViT FE [3]	ViT-L		✓	32	995×4×3	N/A	65.9	89.9
Motionformer [51]	ViT-B		✓	16	370×1×3	109	66.5	90.1
Motionformer [51]	ViT-L	IN-21K+K400	✓	32	1185×1×3	382	68.1	91.2
Video Swin [39]	Swin-B		✓	32	321×1×3	88	69.6	92.7
VIMPAC [65]	ViT-L	HowTo100M+DALLE	✗	10	N/A×10×3	307	68.1	N/A
BEVT [77]	Swin-B	IN-1K+K400+DALLE	✗	32	321×1×3	88	70.6	N/A
MaskFeat↑312 [80]	MViT-L	Kinetics-600	✓	40	2828×1×3	218	75.0	95.0
VideoMAE	ViT-B	Kinetics-400	✗	16	180×2×3	87	69.7	92.3
VideoMAE	ViT-L	Kinetics-400	✗	16	597×2×3	305	74.0	94.6
VideoMAE	ViT-S		✗	16	57×2×3	22	66.8	90.3
VideoMAE	ViT-B	<i>no external data</i>	✗	16	180×2×3	87	70.8	92.4
VideoMAE	ViT-L		✗	16	597×2×3	305	74.3	94.6
VideoMAE	ViT-L		✗	32	1436×1×3	305	<b>75.4</b>	<b>95.2</b>

Table 6: **Comparison with the state-of-the-art methods on Something-Something V2.** Our Video-MAE reconstructs normalized cube pixels and is pre-trained with a masking ratio of 90% for 2400 epochs. “Ex. labels  $\times$ ” means only *unlabelled* data is used during the pre-training phase. “N/A” indicates the numbers are not available for us.

# Experiments on Kinetics 400

Method	Backbone	Extra data	Ex. labels	Frames	GFLOPs	Param	Top-1	Top-5
NL I3D [78]	ResNet101	ImageNet-1K	✓	128	359×10×3	62	77.3	93.3
TANet [41]	ResNet152		✓	16	242×4×3	59	79.3	94.1
TDN <sub>En</sub> [75]	ResNet101		✓	8+16	198×10×3	88	79.4	94.4
TimeSformer [6]	ViT-L	ImageNet-21K	✓	96	8353×1×3	430	80.7	94.7
ViViT FE [3]	ViT-L		✓	128	3980×1×3	N/A	81.7	93.8
Motionformer [51]	ViT-L		✓	32	1185×10×3	382	80.2	94.8
Video Swin [39]	Swin-L		✓	32	604×4×3	197	83.1	95.9
ViViT FE [3]	ViT-L	JFT-300M	✓	128	3980×1×3	N/A	83.5	94.3
ViViT [3]	ViT-H	JFT-300M	✓	32	3981×4×3	N/A	84.9	95.8
VIMPAC [65]	ViT-L	HowTo100M+DALLE	✗	10	N/A×10×3	307	77.4	N/A
BEVT [77]	Swin-B	IN-1K+DALLE	✗	32	282×4×3	88	80.6	N/A
MaskFeat↑352 [80]	MViT-L	Kinetics-600	✗	40	3790×4×3	218	87.0	97.4
ip-CSN [69]	ResNet152	no external data	✗	32	109×10×3	33	77.8	92.8
SlowFast [23]	R101+NL		✗	16+64	234×10×3	60	79.8	93.9
MViTv1 [22]	MViTv1-B		✗	32	170×5×1	37	80.2	94.4
MaskFeat [80]	MViT-L		✗	16	377×10×1	218	84.3	96.3
VideoMAE	ViT-S	no external data	✗	16	57×5×3	22	79.0	93.8
VideoMAE	ViT-B		✗	16	180×5×3	87	81.5	95.1
VideoMAE	ViT-L		✗	16	597×5×3	305	85.2	96.8
VideoMAE	ViT-H		✗	16	1192×5×3	633	<b>86.6</b>	<b>97.1</b>
VideoMAE↑320	ViT-L	no external data	✗	32	3958×4×3	305	86.1	97.3
VideoMAE↑320	ViT-H		✗	32	7397×4×3	633	<b>87.4</b>	<b>97.6</b>

Table 7: **Comparison with the state-of-the-art methods on Kinetics-400.** Our VideoMAE reconstructs normalized cube pixels. Here models are self-supervised pre-trained with a masking ratio of 90% for 1600 epochs on Kinetics-400. VideoMAE↑320 is initialized from its  $224^2$  resolution counterpart and then fine-tuned for evaluation. “Ex. labels ✗” means only *unlabelled* data is used during the pre-training phase. “N/A” indicates the numbers are not available for us.

# SiamMAE

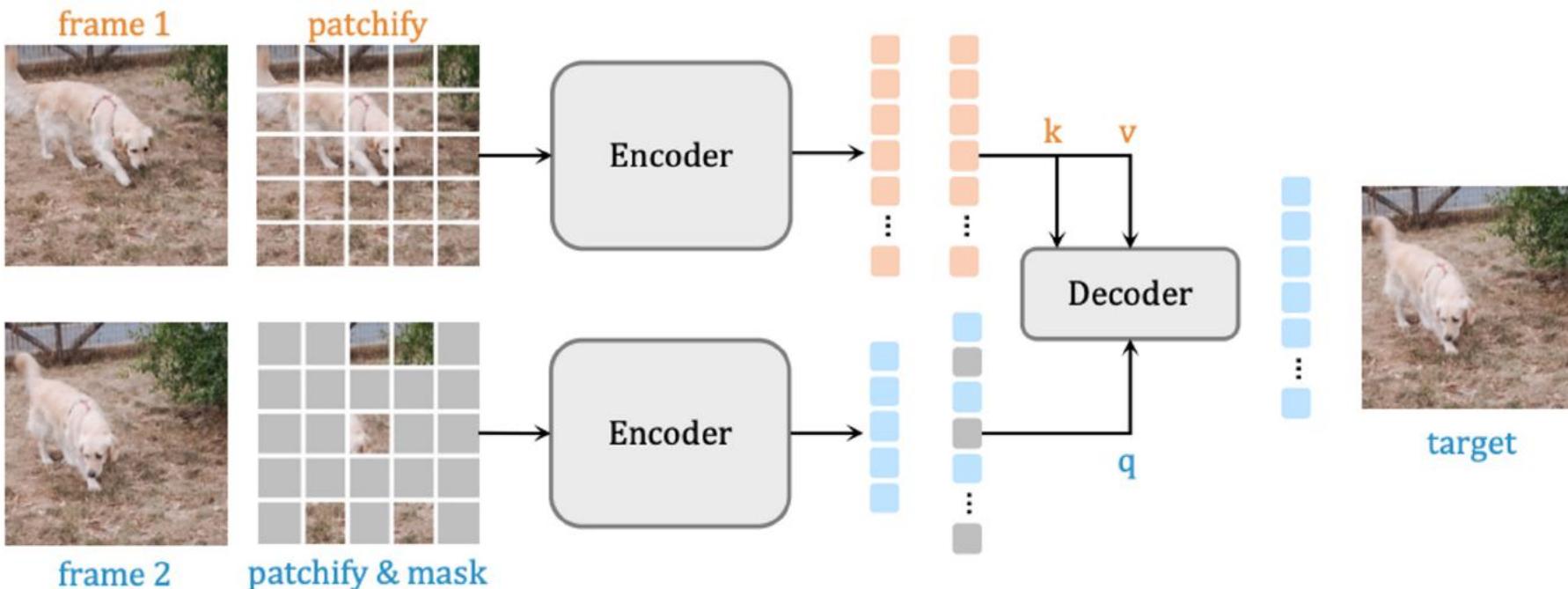
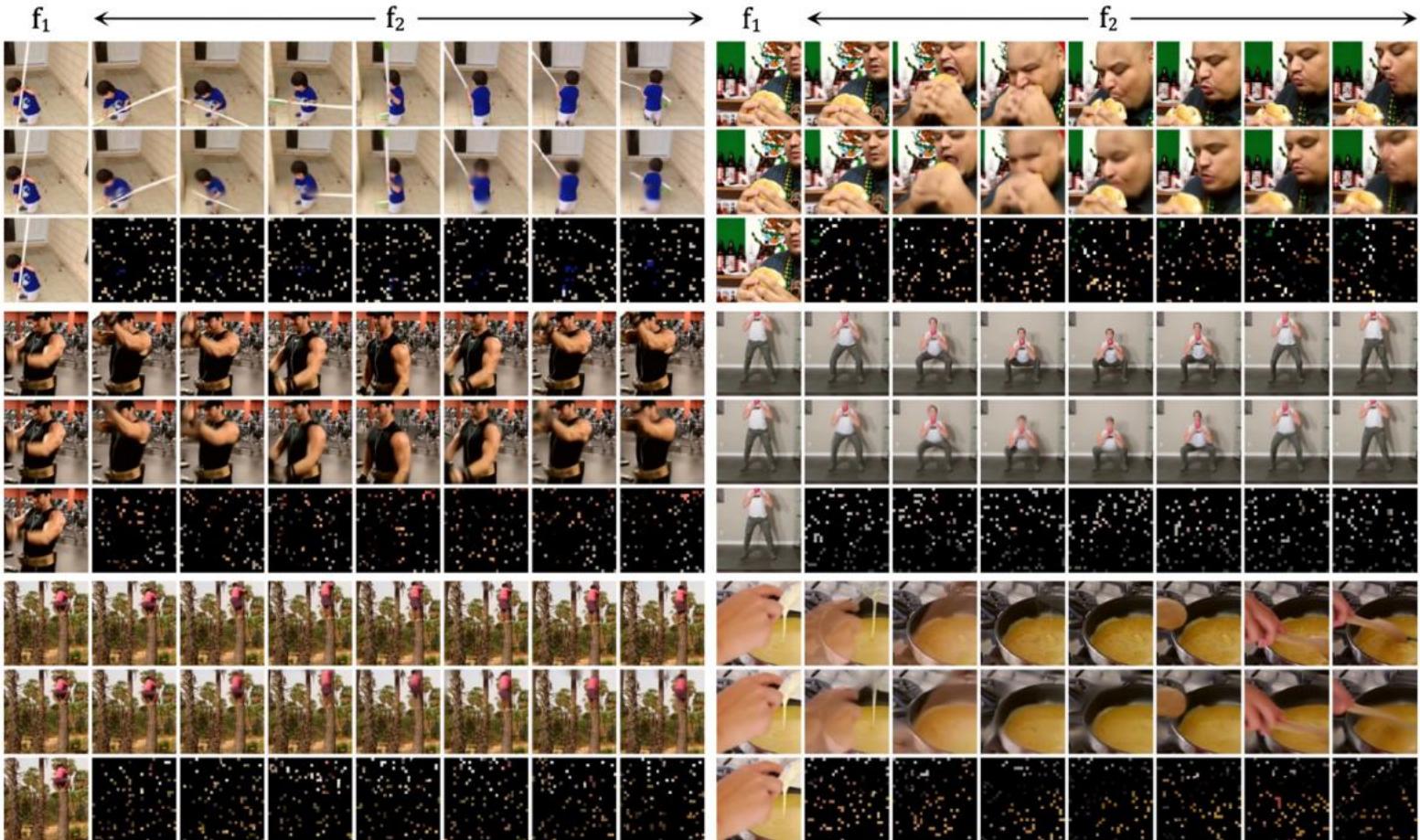


Figure 1: **Siamese Masked Autoencoders.** During pre-training we randomly sample a pair of video frames and randomly mask a huge fraction (95%) of patches of the future frame while leaving the past frame unchanged. The two frames are processed *independently* by a siamese encoder parametrized by a ViT [31]. The decoder consists of a sequence of cross-attention layers and predicts missing patches in the future frame. Videos available at [this project page](#).

# SiamMAE

- By masking a large fraction (95%) of patches in the future frame while leaving the past frame unchanged, SiamMAE encourages the network to focus on object motion and learn object-centric representations.
- SiamMAE outperform state-of-the-art self-supervised methods on video object segmentation, pose keypoint propagation, and semantic part propagation tasks

# SiamMAE



**Figure 2: Visualizations** on the Kinetics-400 [93] validation set (masking ratio 90%). For each video sequence, we sample a clip of 8 frames with a frame gap of 4 and show the original video (top), SiamMAE output (middle), and masked future frames (bottom). Reconstructions are shown with  $f_1$  as the first frame of the video clip and  $f_2$  as the remaining frames, using a SiamMAE pre-trained ViT-S/8 encoder with a masking ratio of 95%.

# SiamMAE Experiments

Method	Backbone	Dataset	DAVIS			VIP mIoU	JHMDB	
			$\mathcal{J}$ & $\mathcal{F}_m$	$\mathcal{J}_m$	$\mathcal{F}_m$		PCK@0.1	PCK@0.2
Supervised [98]	ResNet-50	ImageNet	66.0	63.7	68.4	39.5	59.2	78.3
SimSiam [20]	ResNet-50	ImageNet	66.3	64.5	68.2	35.0	58.4	77.5
MoCo [19]	ResNet-50	ImageNet	65.4	63.2	67.6	36.1	60.4	79.3
TimeCycle [14]	ResNet-50	VLOG	40.7	41.9	39.4	28.9	57.7	78.5
UVC [12]	ResNet-50	Kinetics	56.3	54.5	58.1	34.2	56.0	76.6
VFS [16]	ResNet-50	Kinetics	68.9	66.5	71.3	43.2	60.9	80.7
MAE-ST [27]	ViT-L/16	Kinetics	54.6	55.5	53.6	33.2	44.4	72.5
MAE [24]	VIT-B/16	ImageNet	53.5	52.1	55.0	28.1	44.6	73.4
VideoMAE [28]	ViT-S/16	Kinetics	39.3	39.7	38.9	23.3	41.0	67.9
Dino [17]	ViT-S/16	ImageNet	61.8	60.2	63.4	36.2	45.6	75.0
<b>SiamMAE (ours)</b>	ViT-S/16	Kinetics	<b>62.0</b>	<b>60.3</b>	<b>63.7</b>	<b>37.3</b>	<b>47.0</b>	<b>76.1</b>
Dino [17]	ViT-S/8	ImageNet	69.9	66.6	73.1	39.5	56.5	80.3
<b>SiamMAE (ours)</b>	ViT-S/8	Kinetics	<b>71.4</b>	<b>68.4</b>	<b>74.5</b>	<b>45.9</b>	<b>61.9</b>	<b>83.8</b>

Table 1: **Comparison with prior work** on three downstream tasks: video object segmentation (DAVIS-2017 [95]), human pose propagation (JHMDB [96]) and semantic part propagation (VIP [97])

# Audio-MAE

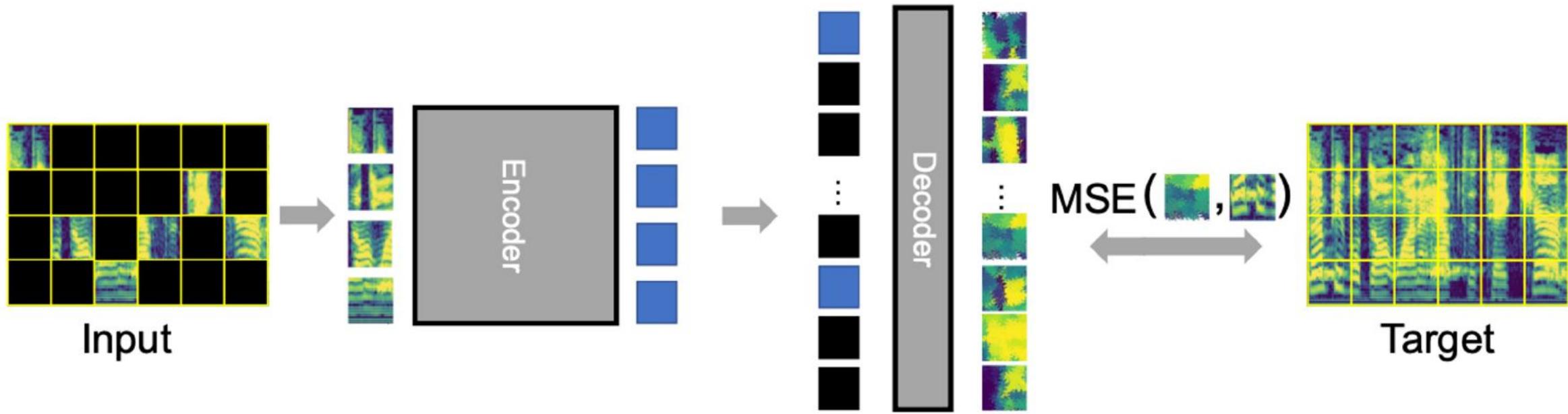


Figure 1: **Audio-MAE for audio self-supervised learning.** An audio recording is first transformed into a spectrogram and split into patches. We embed patches and mask out a large subset (80%). An encoder then operates on the visible (20%) patch embeddings. Finally, a decoder processes the order-restored embeddings and mask tokens to reconstruct the input. Audio-MAE is minimizing the mean square error (MSE) on the masked portion of the reconstruction and the input spectrogram.

# Audio-MAE Experiments

Model	Backbone	PT-Data	AS-20K	AS-2M	ESC-50	SPC-2	SPC-1	SID
<b>No pre-training</b>								
ERANN [58]	CNN	-	-	45.0	89.2	-	-	-
PANN [59]	CNN	-	27.8	43.1	83.3	61.8	-	-
<b>In-domain self-supervised pre-training</b>								
wav2vec 2.0 [33]	Transformer	LS	-	-	-	-	96.2*	75.2*
HuBERT [35]	Transformer	LS	-	-	-	-	96.3*	81.4*
Conformer [37]	Conformer	AS	-	41.1	88.0	-	-	-
SS-AST [18]	ViT-B	AS+LS	31.0	-	88.8	98.0	96.0	64.3
<i>Concurrent MAE-based works</i>								
MaskSpec [43]	ViT-B	AS	32.3	47.1	89.6	97.7	-	-
MAE-AST [38]	ViT-B	AS+LS	30.6	-	90.0	97.9	95.8	63.3
Audio-MAE (global)	ViT-B	AS	$36.6 \pm .11$	$46.8 \pm .06$	$93.6 \pm .11$	<b>98.3 <math>\pm .06</math></b>	<b>97.6 <math>\pm .06</math></b>	$94.1 \pm .06$
Audio-MAE (local)	ViT-B	AS	<b><math>37.0 \pm .11</math></b>	<b><math>47.3 \pm .11</math></b>	<b><math>94.1 \pm .10</math></b>	<b>98.3 <math>\pm .06</math></b>	$96.9 \pm .00$	<b>94.8 <math>\pm .11</math></b>
<b>Out-of-domain supervised pre-training</b>								
PSLA [30]	EffNet [60]	IN	31.9	44.4	-	96.3	-	-
AST [10]	DeiT-B	IN	34.7	45.9	88.7	98.1	95.5	41.1
MBT [11]	ViT-B	IN-21K	31.3	44.3	-	-	-	-
HTS-AT [29]	Swin-B	IN	-	47.1	$97.0^\dagger$	98.0	-	-
PaSST [28]	DeiT-B	IN	-	47.1	$96.8^\dagger$	-	-	-

Table 2: **Comparison with other state-of-the-art models** on audio and speech classification tasks. Metrics are mAP for AS and accuracy (%) for ESC/SPC/SID. For pre-training (PT) dataset, AS:AudioSet, LS:LibriSpeech, and IN:ImageNet.  $^\dagger$ : Fine-tuning results with additional supervised training on AS-2M. We gray-out models pre-trained with external non-audio datasets (e.g., ImageNet). Best single models in AS-2M are compared (no ensembles). \*: linear evaluation results from [53].

# MultiMAE

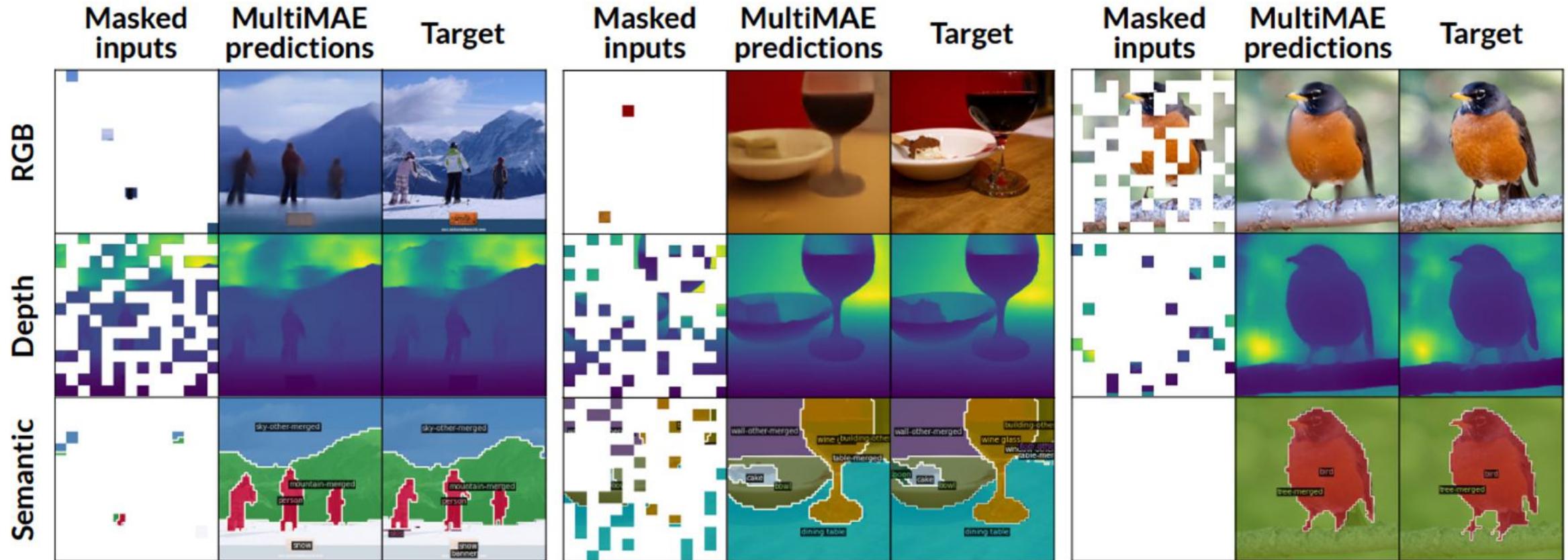


Figure 1. **MultiMAE pre-training objective.** We randomly select 1/6 of all  $16 \times 16$  image patches from multiple modalities and learn to reconstruct the remaining 5/6 masked patches from them. The figure shows validation examples from ImageNet, where masked inputs (left), predictions (middle), and non-masked images (right) for **RGB** (top), **depth** (middle), and **semantic segmentation** (bottom) are provided. Since we do not compute a loss on non-masked patches, we overlay the input patches on the predictions.

# MultimAE

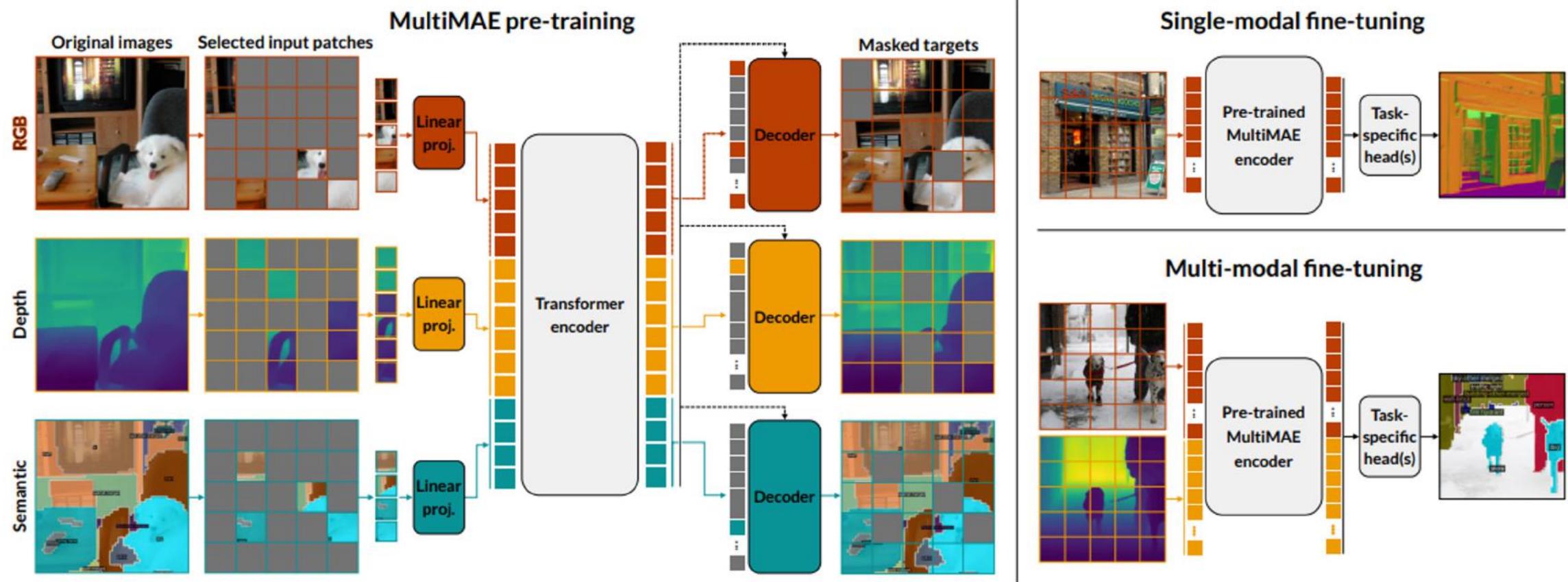


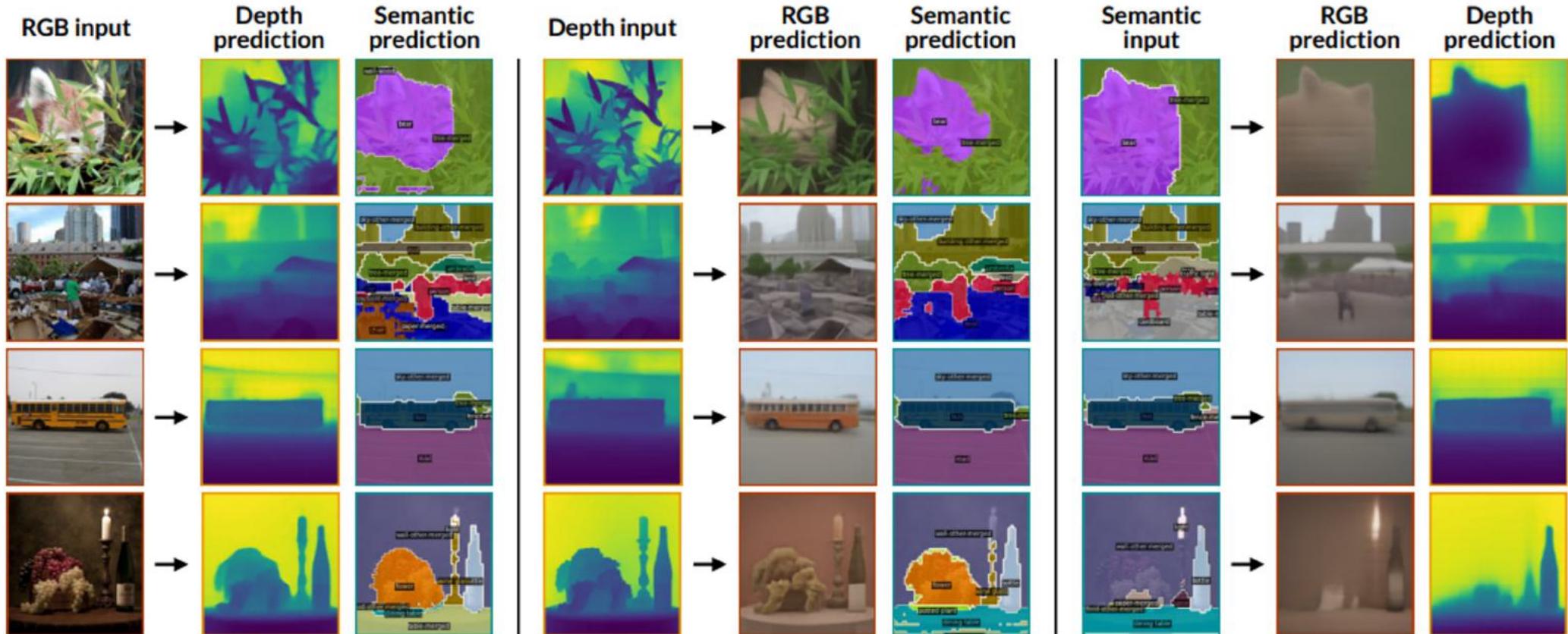
Figure 2. **(Left) MultiMAE pre-training:** A small subset of randomly sampled patches from multiple modalities (e.g., **RGB**, **depth**, and **semantic segmentation**) is linearly projected to tokens with a fixed dimension and encoded using a Transformer. Task-specific decoders reconstruct the masked-out patches by first performing a cross-attention step from queries to the encoded tokens, followed by a shallow Transformer. The queries consist of mask tokens (in gray), with the task-specific encoded tokens added at their respective positions. **(Right) Fine-tuning:** By pre-training on multiple modalities, MultiMAE lends itself to fine-tuning on single-modal and multi-modal downstream tasks. No masking is performed at transfer time.

# MultimAE Experiments

Method	IN-1K (C)	ADE20K (S)	Hypersim (S)	NYUv2 (S)	NYUv2 (D)
Supervised [81]	81.8	45.8	33.9	50.1	80.7
DINO [12]	83.1	44.6	32.5	47.9	81.3
MoCo-v3 [17]	82.8	43.7	31.7	46.6	80.9
MAE [35]	<b>83.3</b>	<b>46.2</b>	<u>36.5</u>	<u>50.8</u>	<u>85.1</u>
MultiMAE	<b>83.3</b>	<b>46.2</b>	<b>37.0</b>	<b>52.0</b>	<b>86.4</b>

Table 1. **Fine-tuning with RGB-only.** We report the top-1 accuracy ( $\uparrow$ ) on ImageNet-1K (IN-1K) [23] classification (C), mIoU ( $\uparrow$ ) on ADE20K [102], Hypersim [68], and NYUv2 [73] semantic segmentation (S), as well as  $\delta_1$  accuracy ( $\uparrow$ ) on NYUv2 depth (D). Text in **bold** and underline indicates the first and second-best results, respectively. All methods are pre-trained on ImageNet-1K (with pseudo labels for MultiMAE).

# MultiMAE Experiments



**Figure 4. Single-modal predictions.** We visualize MultiMAE cross-modal predictions on ImageNet-1K validation images. Only a single, full modality is used as input. The predictions remain plausible despite the absence of input patches from other modalities.

# Lecture overview

- Motivation
- Reconstruct from a corrupted (or partial) version
- **Proxy tasks in computer vision**
  - Relative patch prediction
  - Jigsaw puzzles
  - Rotation
- Contrastive learning

# Relative Position of Image Patches

---

## Unsupervised Visual Representation Learning by Context Prediction

Carl Doersch<sup>1,2</sup> Abhinav Gupta<sup>1</sup> Alexei A. Efros<sup>2</sup>

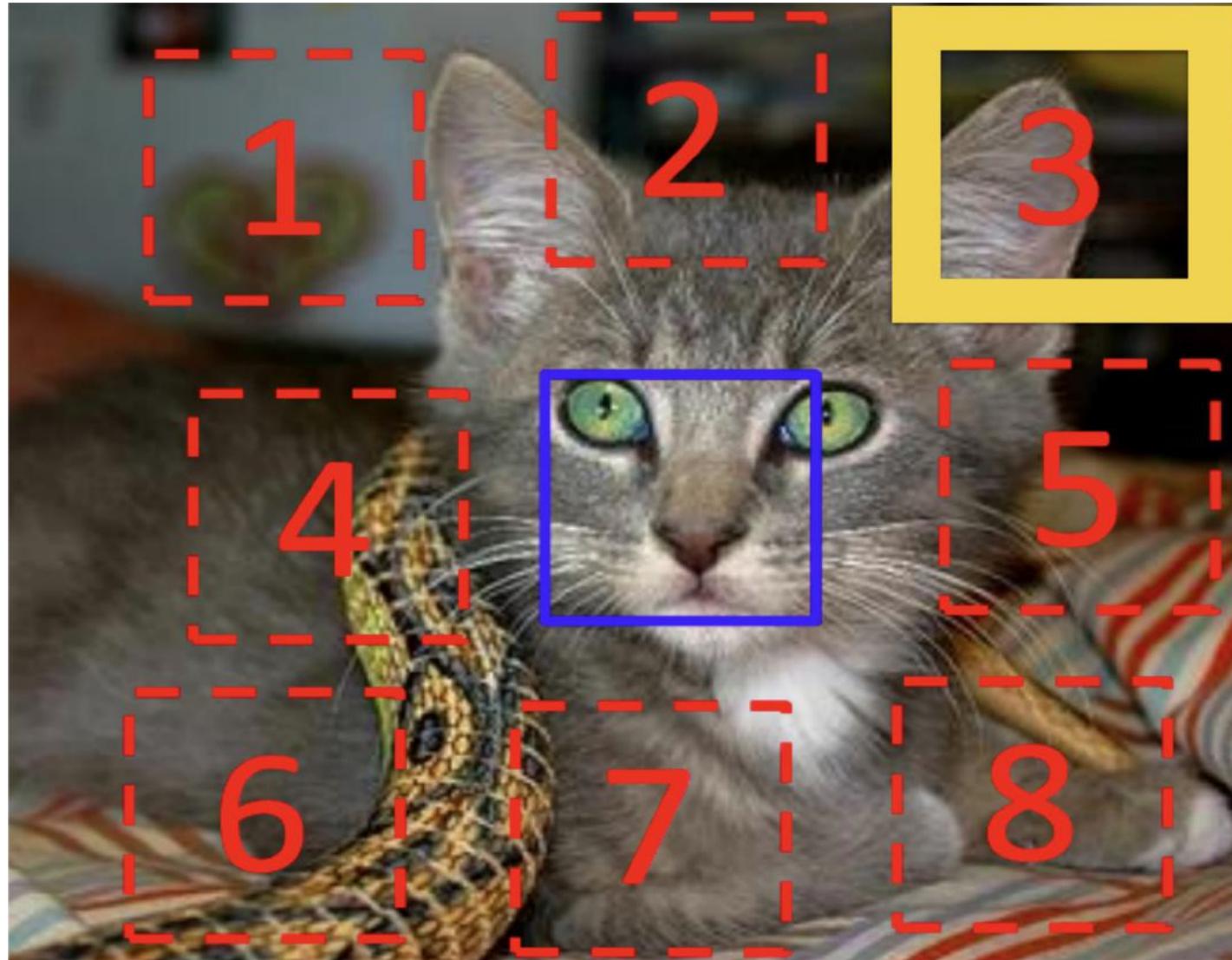
<sup>1</sup> School of Computer Science  
Carnegie Mellon University

<sup>2</sup> Dept. of Electrical Engineering and Computer Science  
University of California, Berkeley



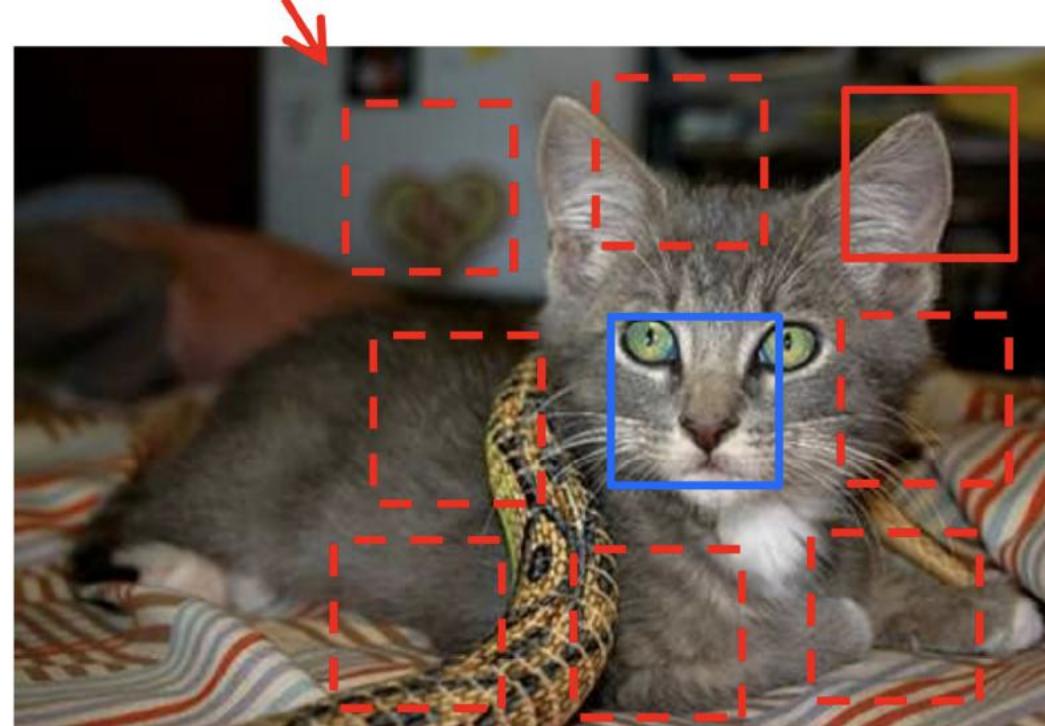
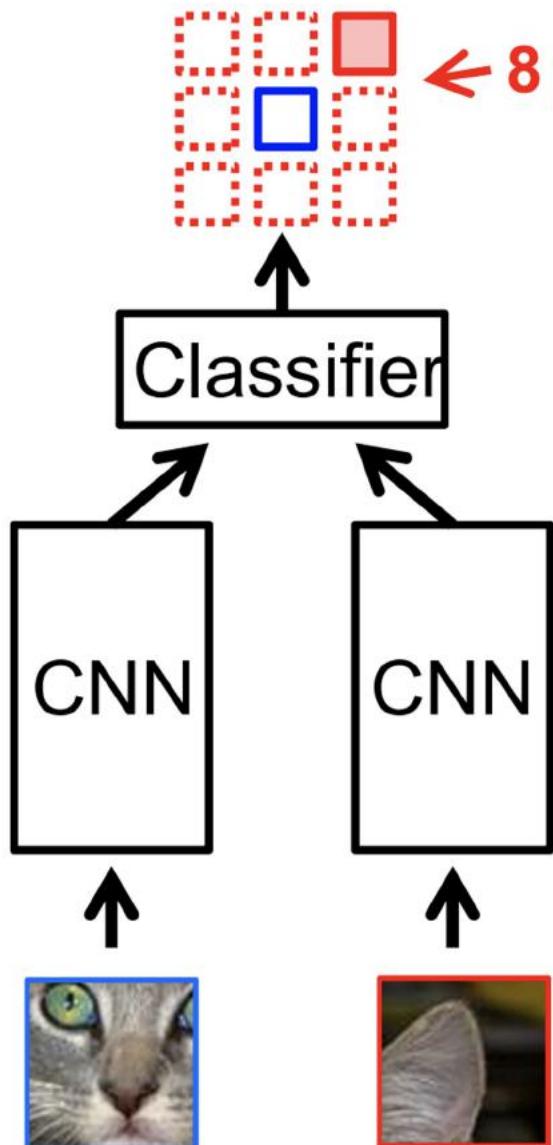
**Task:** Predict the relative position of the second patch with respect to the first

# Relative Position of Image Patches



Doersch,  
Gupta, Efros

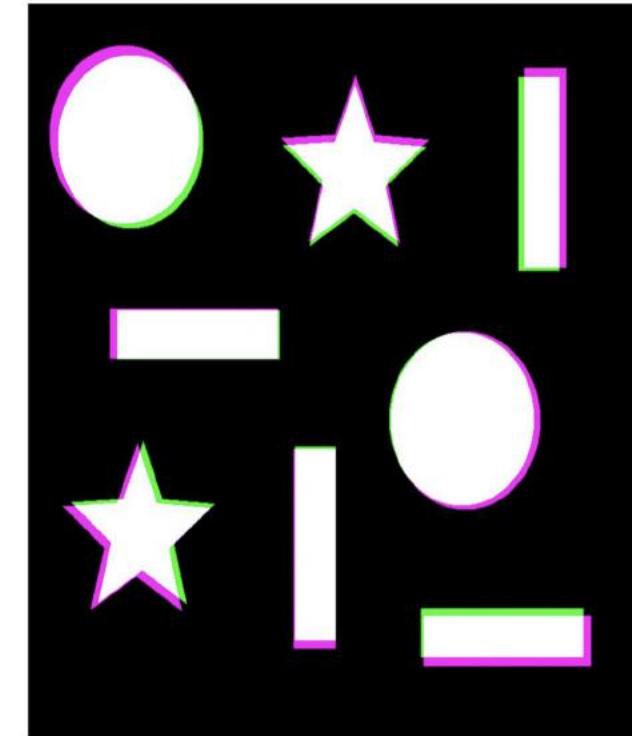
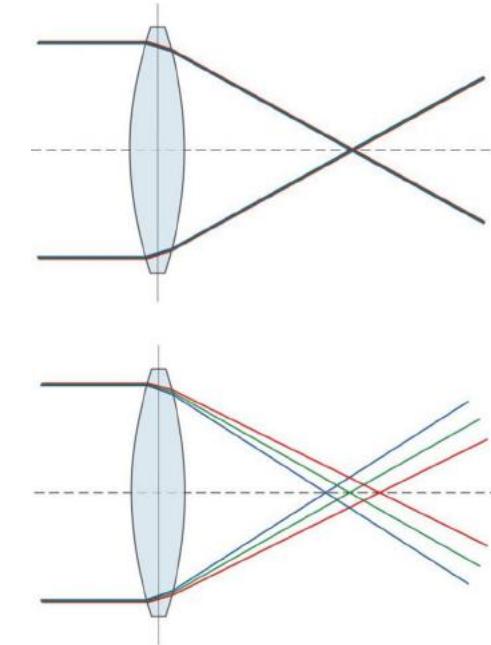
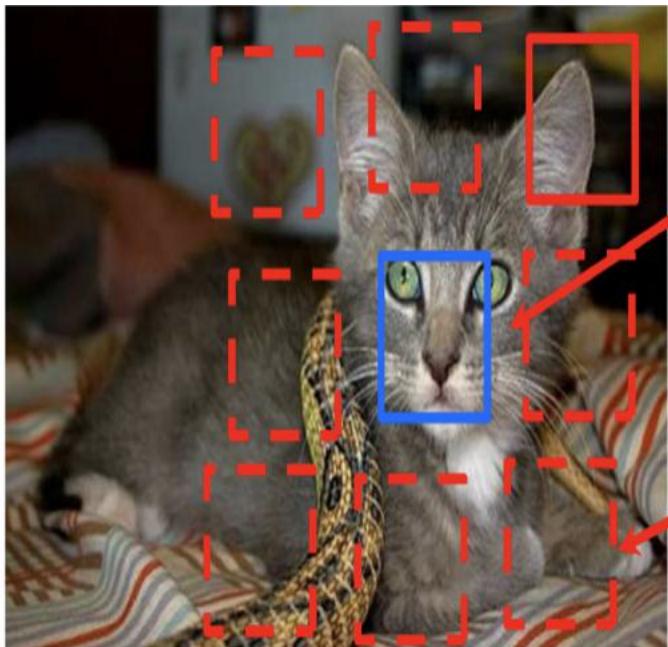
# Relative Position of Image Patches



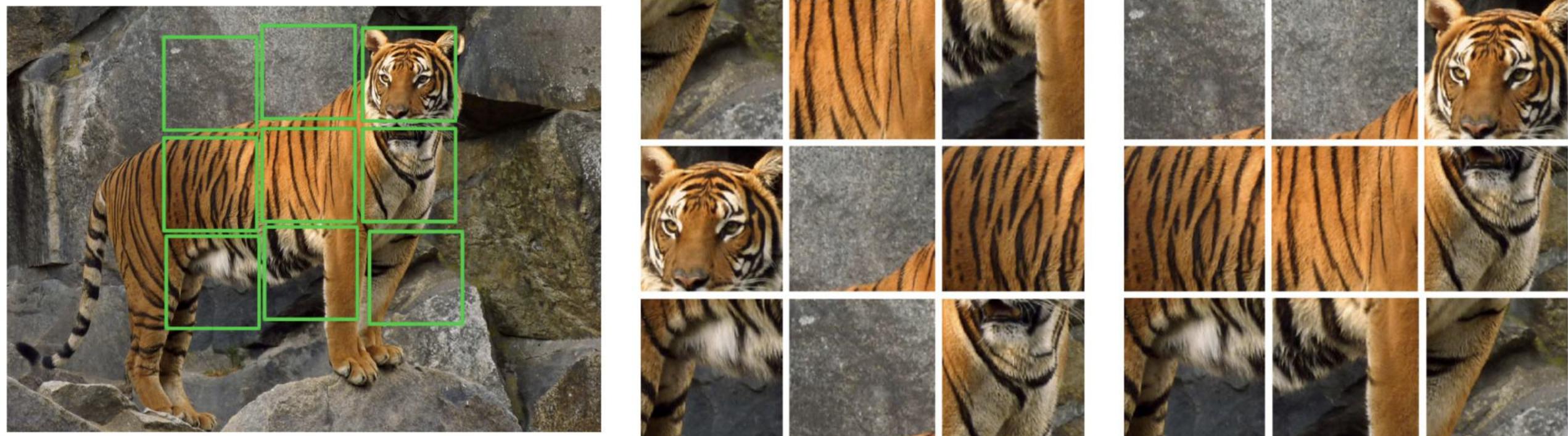
**Randomly Sample Patch  
Sample Second Patch**

Unsupervised visual representation learning by context prediction,  
Carl Doersch, Abhinav Gupta, Alexei A. Efros, ICCV 2015

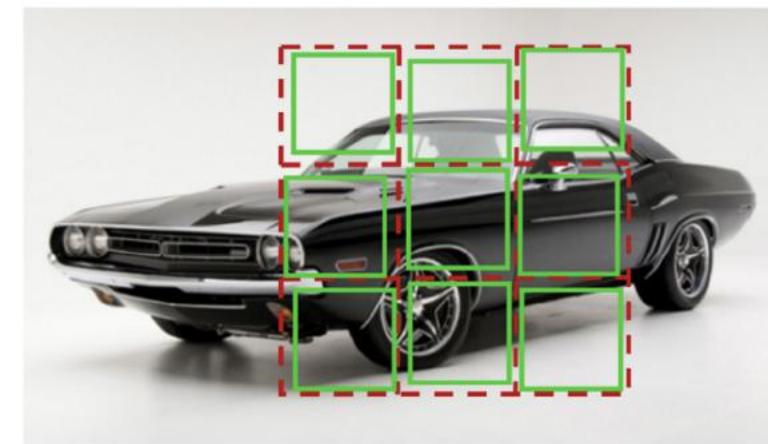
# Relative Position of Image Patches



# Solving Jigsaw Puzzles

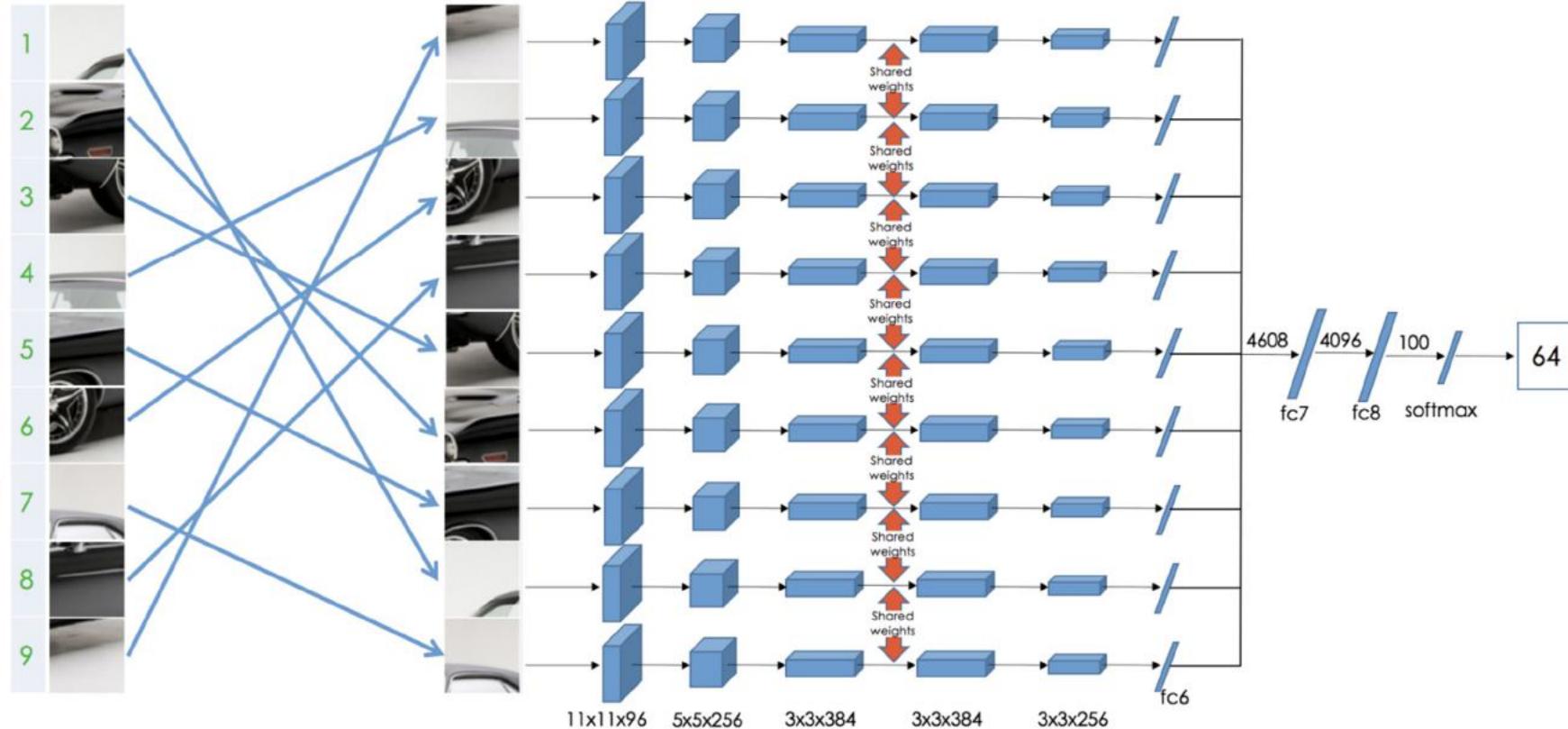


# Solving Jigsaw Puzzles



Permutation Set	
index	permutation
64	9,4,6,8,3,2,5,1,7

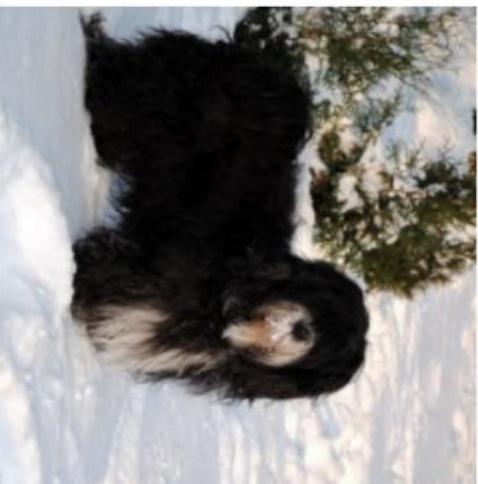
Reorder patches according to the selected permutation



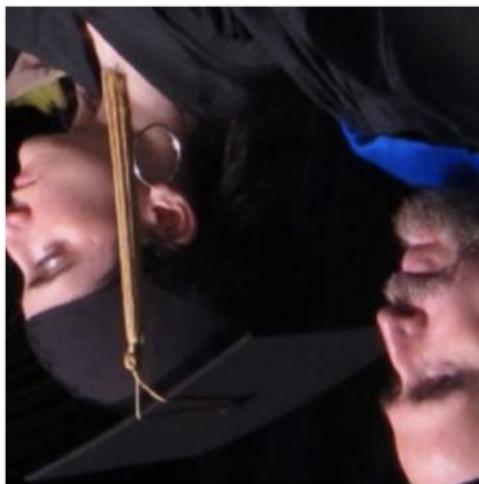
# Rotation



90° rotation



270° rotation



180° rotation

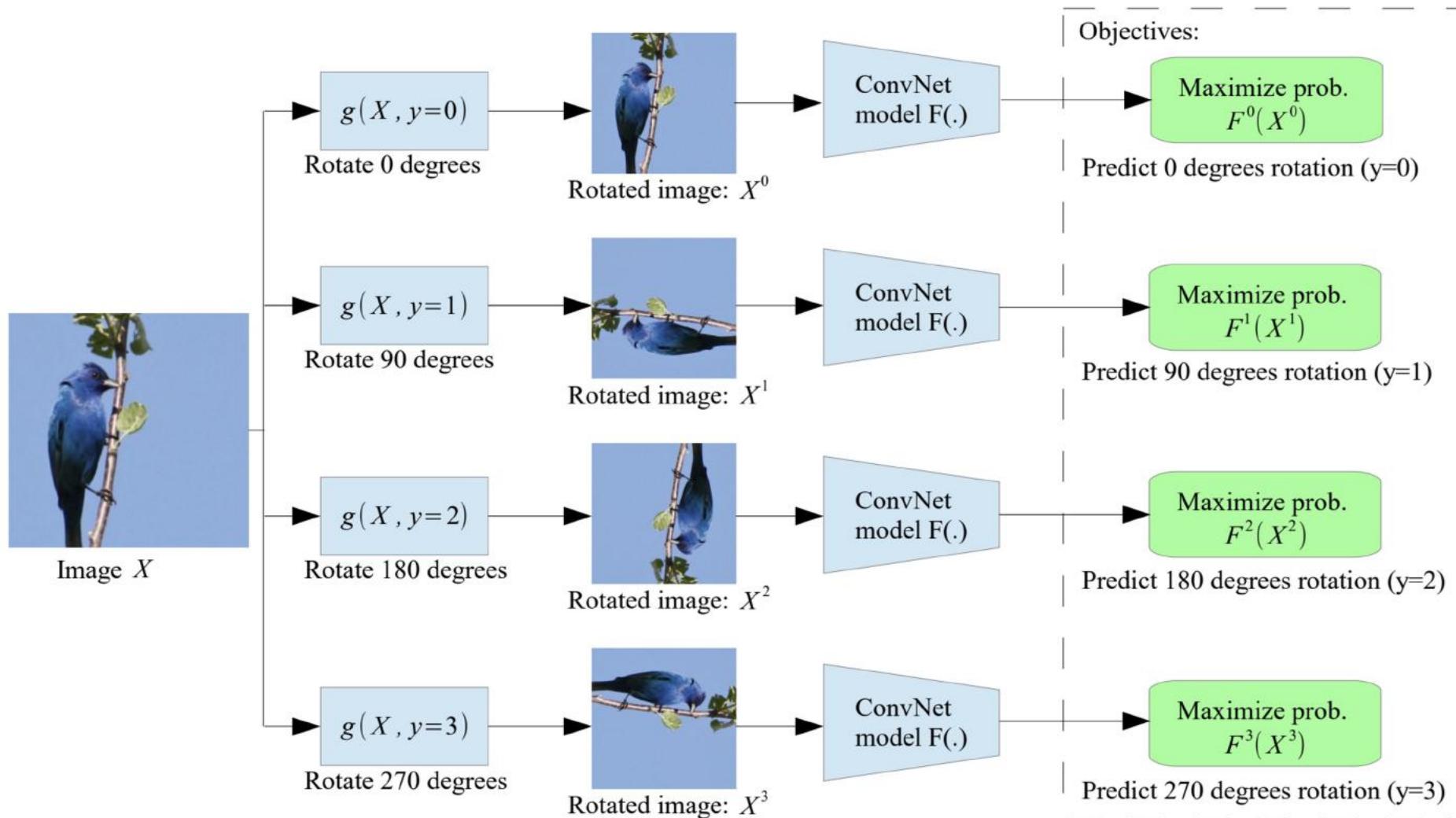


0° rotation



270° rotation

# Rotation



# Rotation

# Rotations	Rotations	CIFAR-10 Classification Accuracy
4	$0^\circ, 90^\circ, 180^\circ, 270^\circ$	<b>89.06</b>
8	$0^\circ, 45^\circ, 90^\circ, 135^\circ, 180^\circ, 225^\circ, 270^\circ, 315^\circ$	88.51
2	$0^\circ, 180^\circ$	87.46
2	$90^\circ, 270^\circ$	85.52

# Rotation

Table 4: **Task Generalization: ImageNet top-1 classification with non-linear layers.** We compare our unsupervised feature learning approach with other unsupervised approaches by training non-linear classifiers on top of the feature maps of each layer to perform the 1000-way ImageNet classification task, as proposed by Noroozi & Favaro (2016). For instance, for the conv5 feature map we train the layers that follow the conv5 layer in the AlexNet architecture (i.e., fc6, fc7, and fc8). Similarly for the conv4 feature maps. We implemented those non-linear classifiers with batch normalization units after each linear layer (fully connected or convolutional) and without employing drop out units. All approaches use AlexNet variants and were pre-trained on ImageNet without labels except the ImageNet labels and Random entries. During testing we use a single crop and do not perform flipping augmentation. We report top-1 classification accuracy.

Method	Conv4	Conv5
ImageNet labels from (Bojanowski & Joulin, 2017)	59.7	59.7
Random from (Noroozi & Favaro, 2016)	27.1	12.0
Tracking Wang & Gupta (2015)	38.8	29.8
Context (Doersch et al., 2015)	45.6	30.4
Colorization (Zhang et al., 2016a)	40.7	35.2
Jigsaw Puzzles (Noroozi & Favaro, 2016)	45.3	34.6
BIGAN (Donahue et al., 2016)	41.9	32.2
NAT (Bojanowski & Joulin, 2017)	-	36.0
(Ours) RotNet	50.0	43.8

# Rotation

Table 5: **Task Generalization: ImageNet top-1 classification with linear layers.** We compare our unsupervised feature learning approach with other unsupervised approaches by training logistic regression classifiers on top of the feature maps of each layer to perform the 1000-way ImageNet classification task, as proposed by Zhang et al. (2016a). All weights are frozen and feature maps are spatially resized (with adaptive max pooling) so as to have around 9000 elements. All approaches use AlexNet variants and were pre-trained on ImageNet without labels except the ImageNet labels and Random entries.

Method	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
Random rescaled Krähenbühl et al. (2015)	17.5	23.0	24.5	23.2	20.6
Context (Doersch et al., 2015)	16.2	23.3	30.2	31.7	29.6
Context Encoders (Pathak et al., 2016b)	14.1	20.7	21.0	19.8	15.5
Colorization (Zhang et al., 2016a)	12.5	24.5	30.4	31.5	30.3
Jigsaw Puzzles (Noroozi & Favaro, 2016)	18.2	28.8	34.0	33.9	27.1
BIGAN (Donahue et al., 2016)	17.7	24.5	31.0	29.9	28.0
Split-Brain (Zhang et al., 2016b)	17.7	29.3	35.4	35.2	32.8
Counting (Noroozi et al., 2017)	18.0	30.6	34.3	32.5	25.7
(Ours) RotNet	<b>18.8</b>	<b>31.7</b>	<b>38.7</b>	<b>38.2</b>	<b>36.5</b>

# Rotation

**Table 7: Task & Dataset Generalization: PASCAL VOC 2007 classification and detection results, and PASCAL VOC 2012 segmentation results.** We used the publicly available testing frameworks of Krahenbuhl et al. (2015) for classification, of Girshick (2015) for detection, and of Long et al. (2015) for segmentation. For classification, we either fix the features before conv5 (column *fc6-8*) or we fine-tune the whole model (column *all*). For detection we use multi-scale training and single scale testing. All approaches use AlexNet variants and were pre-trained on ImageNet without labels except the ImageNet labels and Random entries. After unsupervised training, we absorb the batch normalization units on the linear layers and we use the weight rescaling technique proposed by Krahenbuhl et al. (2015) (which is common among the unsupervised methods). As customary, we report the mean average precision (mAP) on the classification and detection tasks, and the mean intersection over union (mIoU) on the segmentation task.

	Classification (%mAP)	Detection (%mAP)	Segmentation (%mIoU)
Trained layers	fc6-8	all	all
ImageNet labels	78.9	79.9	56.8
Random		53.3	43.4
Random rescaled Krähenbühl et al. (2015)	39.2	56.6	48.0
Egomotion (Agrawal et al., 2015)	31.0	54.2	43.9
Context Encoders (Pathak et al., 2016b)	34.6	56.5	44.5
Tracking (Wang & Gupta, 2015)	55.6	63.1	47.4
Context (Doersch et al., 2015)	55.1	65.3	51.1
Colorization (Zhang et al., 2016a)	61.5	65.6	46.9
BIGAN (Donahue et al., 2016)	52.3	60.1	46.9
Jigsaw Puzzles (Noroozi & Favaro, 2016)	-	67.6	53.2
NAT (Bojanowski & Joulin, 2017)	56.7	65.3	49.4
Split-Brain (Zhang et al., 2016b)	63.0	67.1	46.7
ColorProxy (Larsson et al., 2017)		65.9	38.4
Counting (Noroozi et al., 2017)	-	67.7	51.4
<b>(Ours) RotNet</b>	<b>70.87</b>	<b>72.97</b>	<b>54.4</b>
			<b>39.1</b>

# Temporal coherence of color

Task: given a color video ...

Colorize all frames of a gray scale version using a reference frame



reference frame

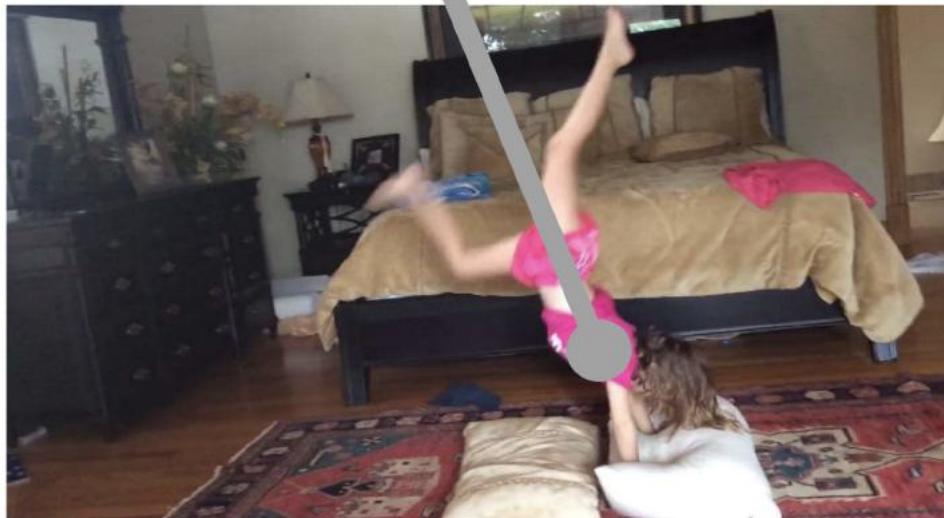


gray-scale video

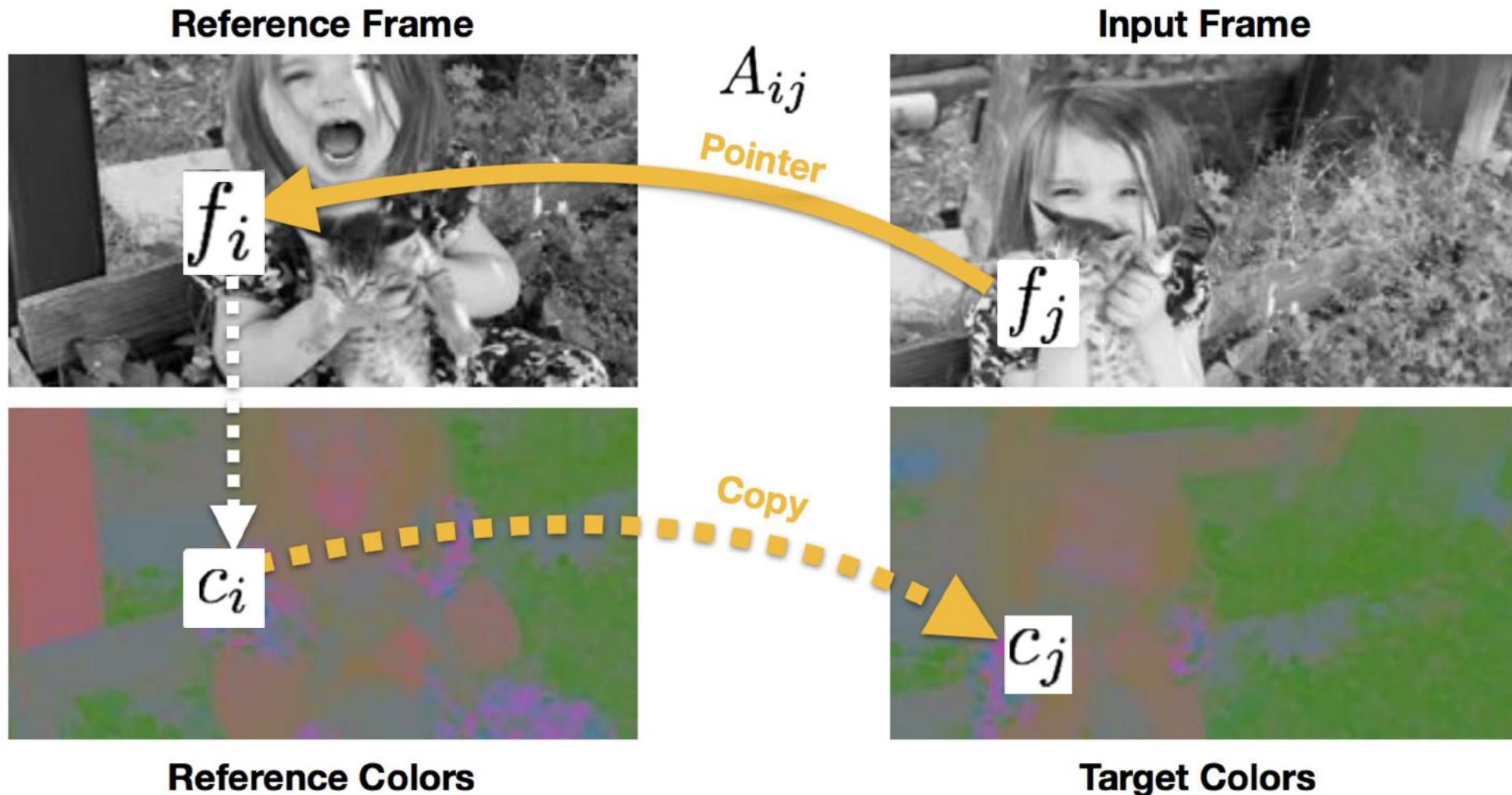
# Temporal coherence of color



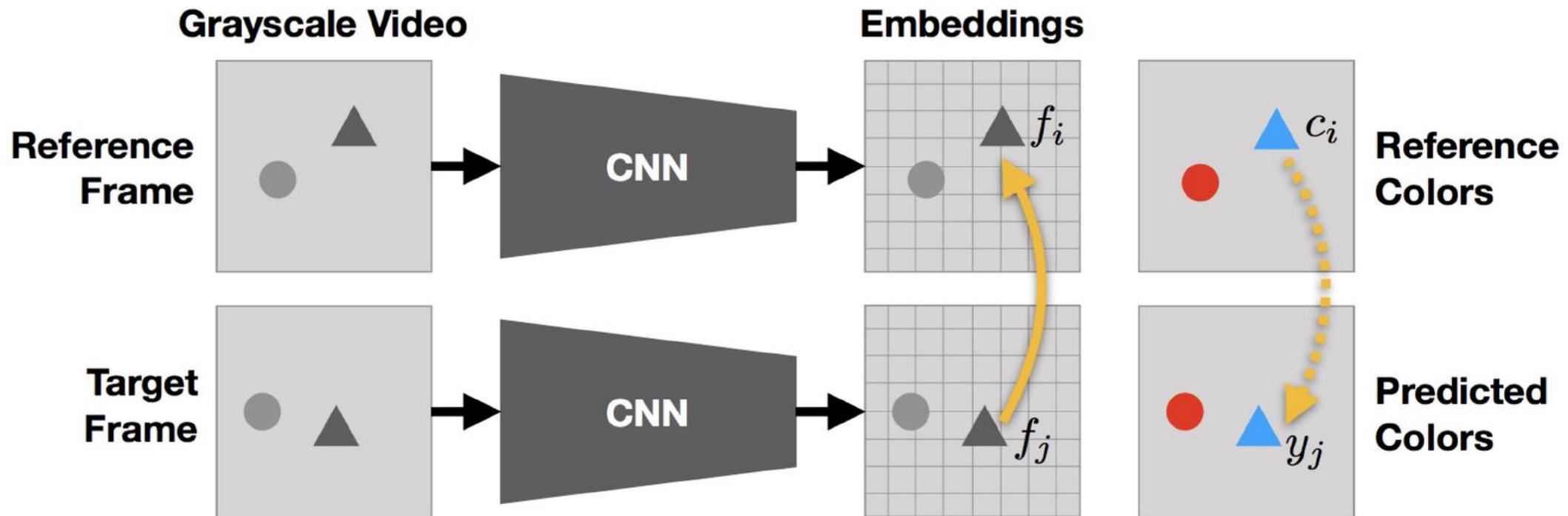
# Temporal coherence of color



# Tracking emerges from colorization



# Tracking emerges from colorization

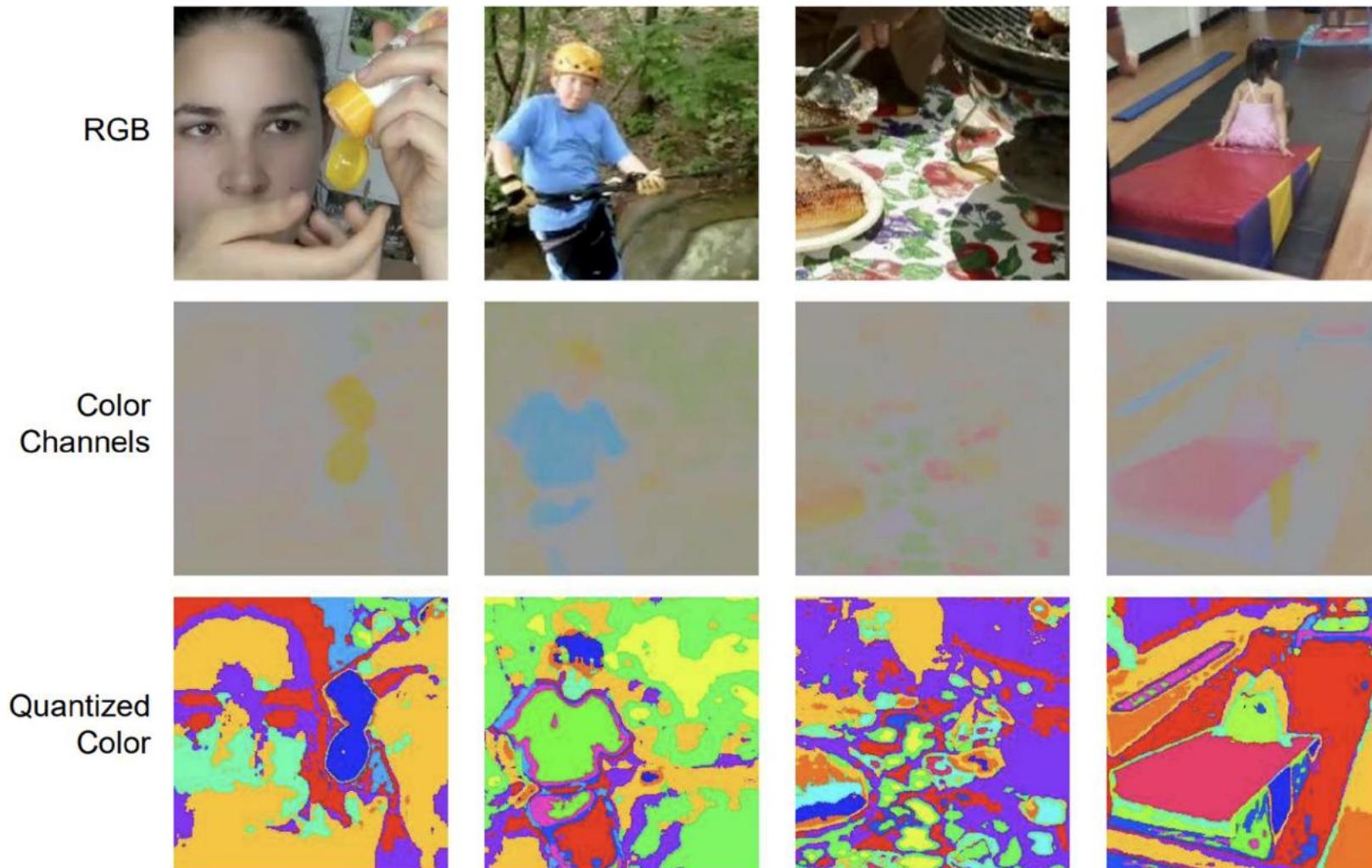


$$A_{ij} = \frac{\exp(f_i^T f_j)}{\sum_k \exp(f_k^T f_j)}$$

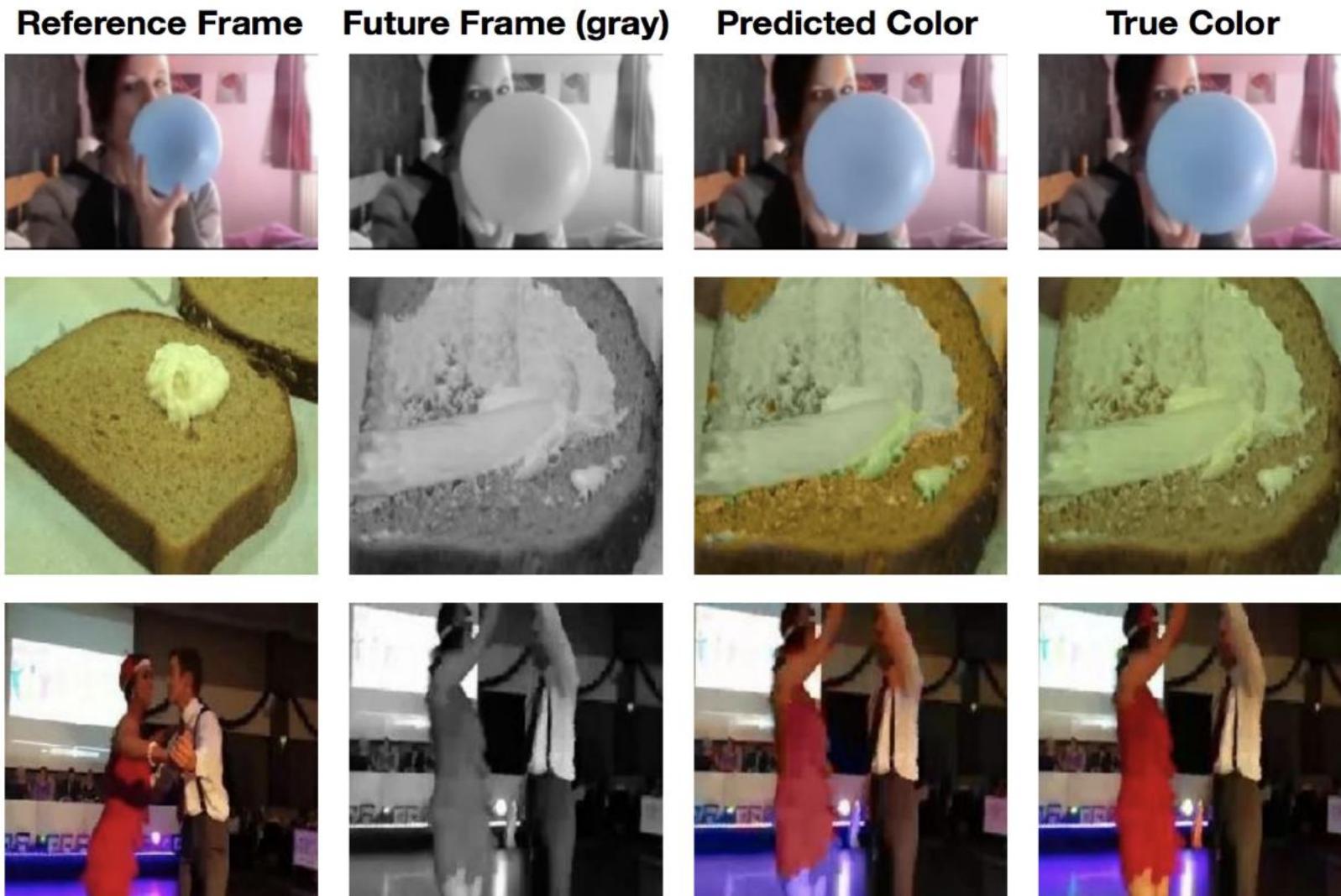
$$\hat{c}_j = \sum_i A_{ij} c_i$$

$$\min_f \mathcal{L} \left( c_j, \sum_i A_{ij} c_i \right)$$

# Tracking emerges from colorization



# Tracking emerges from colorization



# Tracking emerges from colorization



# Last time on COMP547

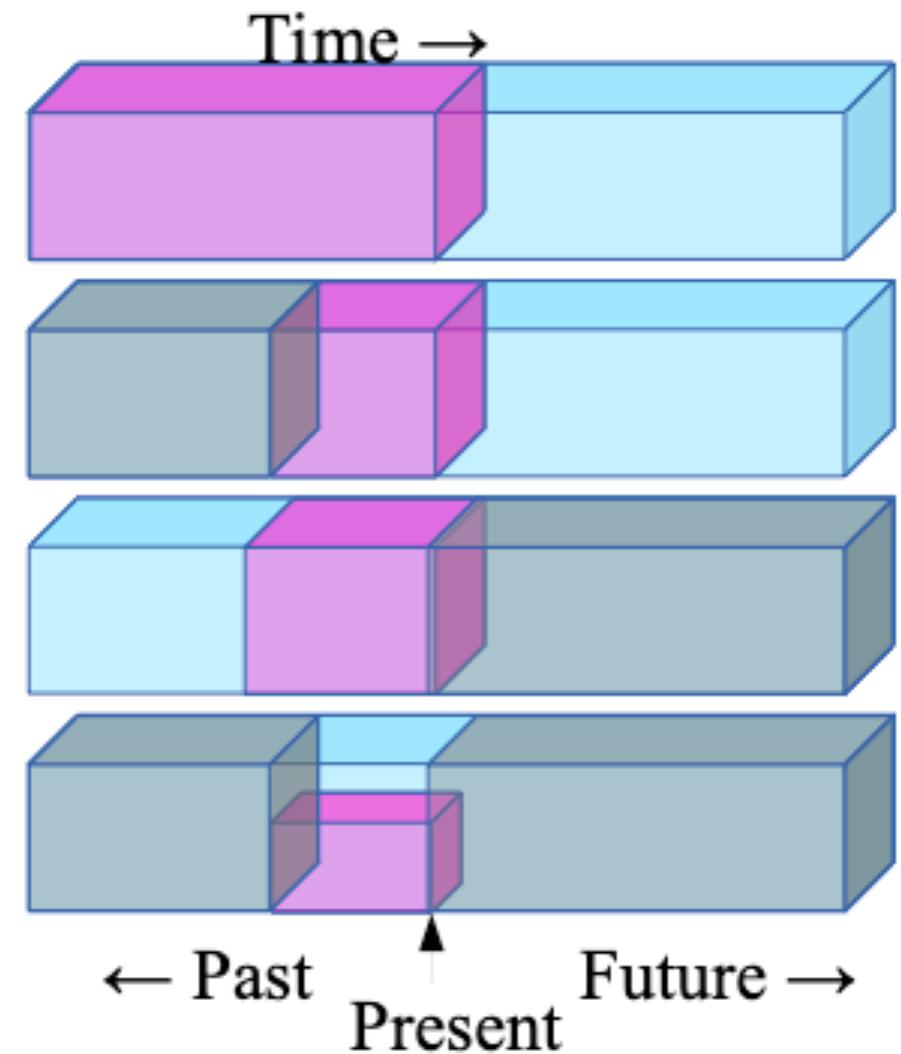
- Motivation
- Reconstruct from a corrupted (or partial) version
- Proxy tasks in computer vision
- Contrastive learning

# Lecture overview

- Motivation
- Reconstruct from a corrupted (or partial) version
- Proxy tasks in computer vision
- Contrastive learning
  - Contrastive Predictive Coding (CPC)
  - Instance Discrimination
  - Recent State-of-the-art progress

# Predicting neighbouring context

- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.
- ▶ Predict the occluded from the visible
- ▶ **Pretend there is a part of the input you don't know and predict that.**



# Deep InfoMax

## LEARNING DEEP REPRESENTATIONS BY MUTUAL INFORMATION ESTIMATION AND MAXIMIZATION

**R Devon Hjelm**

MSR Montreal, MILA, UdeM, IVADO  
devon.hjelm@microsoft.com

**Alex Fedorov**

MRN, UNM

**Samuel Lavoie-Marchildon**

MILA, UdeM

**Karan Grewal**

U Toronto

**Phil Bachman**

MSR Montreal

**Adam Trischler**

MSR Montreal

**Yoshua Bengio**

MILA, UdeM, IVADO, CIFAR

### ABSTRACT

This work investigates unsupervised learning of representations by maximizing mutual information between an input and the output of a deep neural network encoder. Importantly, we show that structure matters: incorporating knowledge about locality in the input into the objective can significantly improve a representation’s suitability for downstream tasks. We further control characteristics of the representation by matching to a prior distribution adversarially. Our method, which we call Deep InfoMax (DIM), outperforms a number of popular unsupervised learning methods and compares favorably with fully-supervised learning on several classification tasks in with some standard architectures. DIM opens new avenues for unsupervised learning of representations and is an important step towards flexible formulations of representation learning objectives for specific end-goals.

# Deep InfoMax

- Network encodes the input
- The discriminator estimates mutual information (batch-wise)
- Estimate is used to maximize the MI between encoder input and output

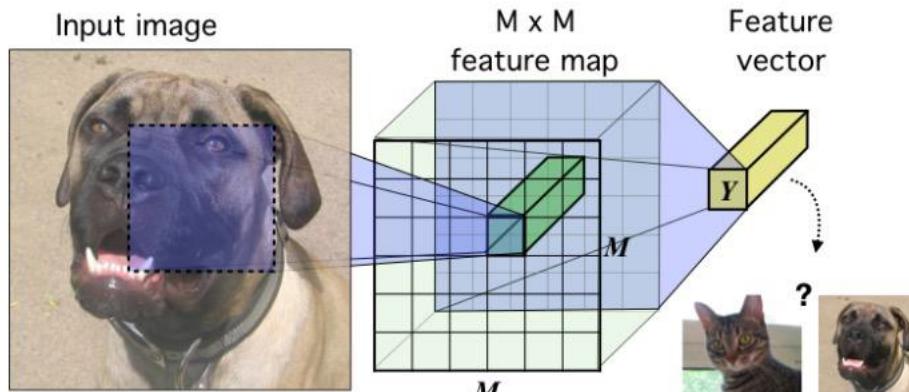


Figure 1: **The base encoder model in the context of image data.** An image (in this case) is encoded using a convnet until reaching a feature map of  $M \times M$  feature vectors corresponding to  $M \times M$  input patches. These vectors are summarized into a single feature vector,  $Y$ . Our goal is to train this network such that useful information about the input is easily extracted from the high-level features.

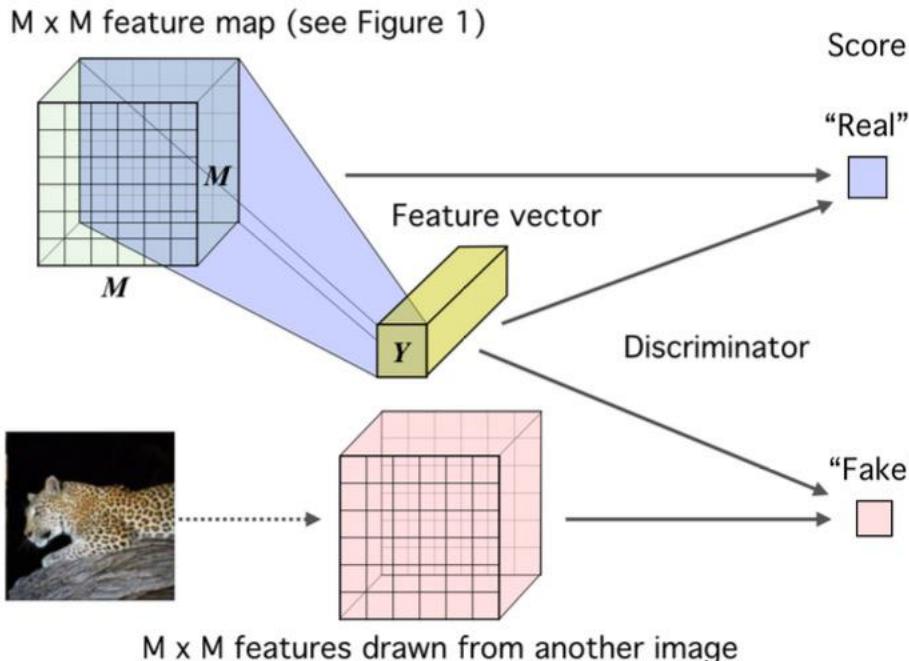
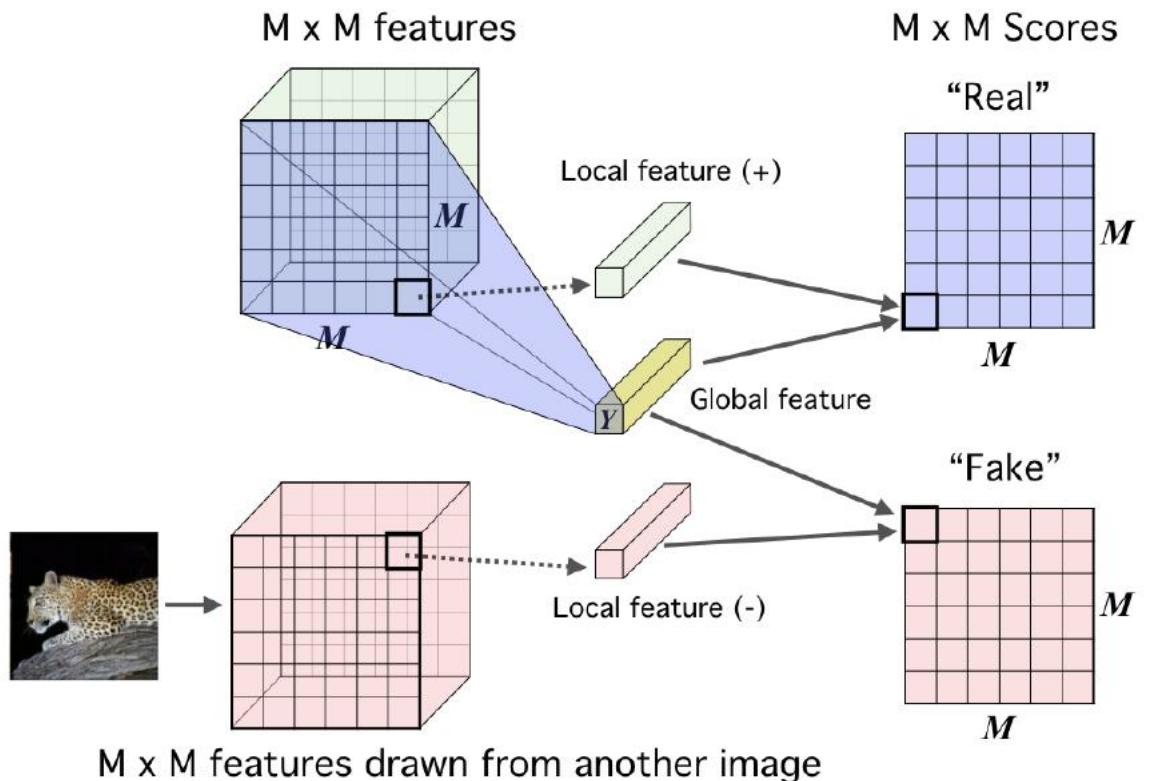


Figure 2: **Deep InfoMax (DIM) with a global  $\text{MI}(X; Y)$  objective.** Here, we pass both the high-level feature vector,  $Y$ , and the lower-level  $M \times M$  feature map (see Figure 1) through a discriminator to get the score. Fake samples are drawn by combining the same feature vector with a  $M \times M$  feature map from another image.

# Deep InfoMax



**Figure 3: Maximizing mutual information between local features and global features.** First we encode the image to a feature map that reflects some structural aspect of the data, e.g. spatial locality, and we further summarize this feature map into a global feature vector (see Figure 1). We then concatenate this feature vector with the lower-level feature map *at every location*. A score is produced for each local-global pair through an additional function (see the Appendix A.2 for details).

# Deep InfoMax

Table 1: Classification accuracy (top 1) results on CIFAR10 and CIFAR100. DIM(L) (i.e., with the local-only objective) outperforms all other unsupervised methods presented by a wide margin. In addition, DIM(L) approaches or even surpasses a fully-supervised classifier with similar architecture. DIM with the global-only objective is competitive with some models across tasks, but falls short when compared to generative models and DIM(L) on CIFAR100. Fully-supervised classification results are provided for comparison.

Model	CIFAR10			CIFAR100		
	conv	fc (1024)	Y(64)	conv	fc (1024)	Y(64)
Fully supervised		75.39			42.27	
VAE	60.71	60.54	54.61	37.21	34.05	24.22
AE	62.19	55.78	54.47	31.50	23.89	27.44
$\beta$ -VAE	62.4	57.89	55.43	32.28	26.89	28.96
AAE	59.44	57.19	52.81	36.22	33.38	23.25
BiGAN	62.57	62.74	52.54	37.59	33.34	21.49
NAT	56.19	51.29	31.16	29.18	24.57	9.72
DIM(G)	52.2	52.84	43.17	27.68	24.35	19.98
DIM(L) (DV)	<b>72.66</b>	<b>70.60</b>	<b>64.71</b>	<b>48.52</b>	<b>44.44</b>	<b>39.27</b>
DIM(L) (JSD)	<b>73.25</b>	<b>73.62</b>	<b>66.96</b>	<b>48.13</b>	<b>45.92</b>	<b>39.60</b>
DIM(L) (infoNCE)	<b>75.21</b>	<b>75.57</b>	<b>69.13</b>	<b>49.74</b>	<b>47.72</b>	<b>41.61</b>

# Deep InfoMax

Table 2: Classification accuracy (top 1) results on Tiny ImageNet and STL-10. For Tiny ImageNet, DIM with the local objective outperforms all other models presented by a large margin, and approaches accuracy of a fully-supervised classifier similar to the Alexnet architecture used here.

	Tiny ImageNet			STL-10 (random crop pretraining)			
	conv	fc (4096)	$Y(64)$	conv	fc (4096)	$Y(64)$	SS
Fully supervised	36.60			68.7			
VAE	18.63	16.88	11.93	58.27	56.72	46.47	68.65
AE	19.07	16.39	11.82	58.19	55.57	46.82	70.29
$\beta$ -VAE	19.29	16.77	12.43	57.15	55.14	46.87	70.53
AAE	18.04	17.27	11.49	59.54	54.47	43.89	64.15
BiGAN	24.38	20.21	13.06	71.53	67.18	58.48	74.77
NAT	13.70	11.62	1.20	64.32	61.43	48.84	70.75
DIM(G)	11.32	6.34	4.95	42.03	30.82	28.09	51.36
DIM(L) (DV)	30.35	29.51	28.18	69.15	63.81	61.92	71.22
DIM(L) (JSD)	<b>33.54</b>	<b>36.88</b>	<b>31.66</b>	<b>72.86</b>	<b>70.85</b>	<b>65.93</b>	<b>76.96</b>
DIM(L) (infoNCE)	<b>34.21</b>	<b>38.09</b>	<b>33.33</b>	<b>72.57</b>	<b>70.00</b>	<b>67.08</b>	<b>76.81</b>

# Contrastive Predictive Coding

---

## Representation Learning with Contrastive Predictive Coding

---

**Aaron van den Oord**

DeepMind

[avdnoord@google.com](mailto:avdnoord@google.com)

**Yazhe Li**

DeepMind

[yazhe@google.com](mailto:yazhe@google.com)

**Oriol Vinyals**

DeepMind

[vinyals@google.com](mailto:vinyals@google.com)

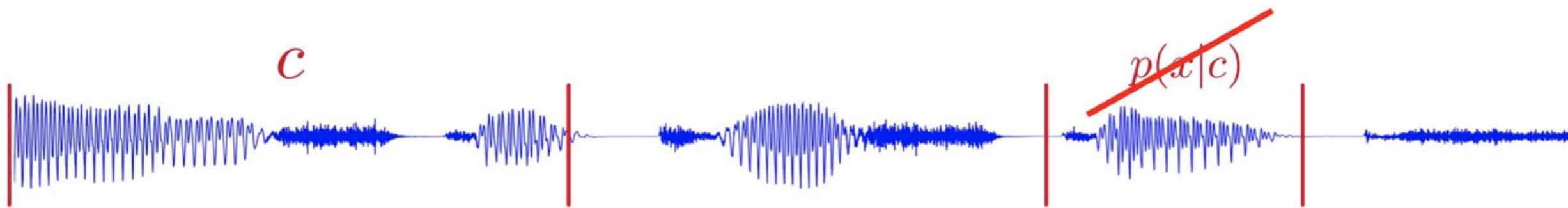
### Abstract

While supervised learning has enabled great progress in many applications, unsupervised learning has not seen such widespread adoption, and remains an important and challenging endeavor for artificial intelligence. In this work, we propose a universal unsupervised learning approach to extract useful representations from high-dimensional data, which we call Contrastive Predictive Coding. The key insight of our model is to learn such representations by predicting the future in *latent* space by using powerful autoregressive models. We use a probabilistic contrastive loss which induces the latent space to capture information that is maximally useful to predict future samples. It also makes the model tractable by using negative sampling. While most prior work has focused on evaluating representations for a particular modality, we demonstrate that our approach is able to learn useful representations achieving strong performance on four distinct domains: speech, images, text and reinforcement learning in 3D environments.

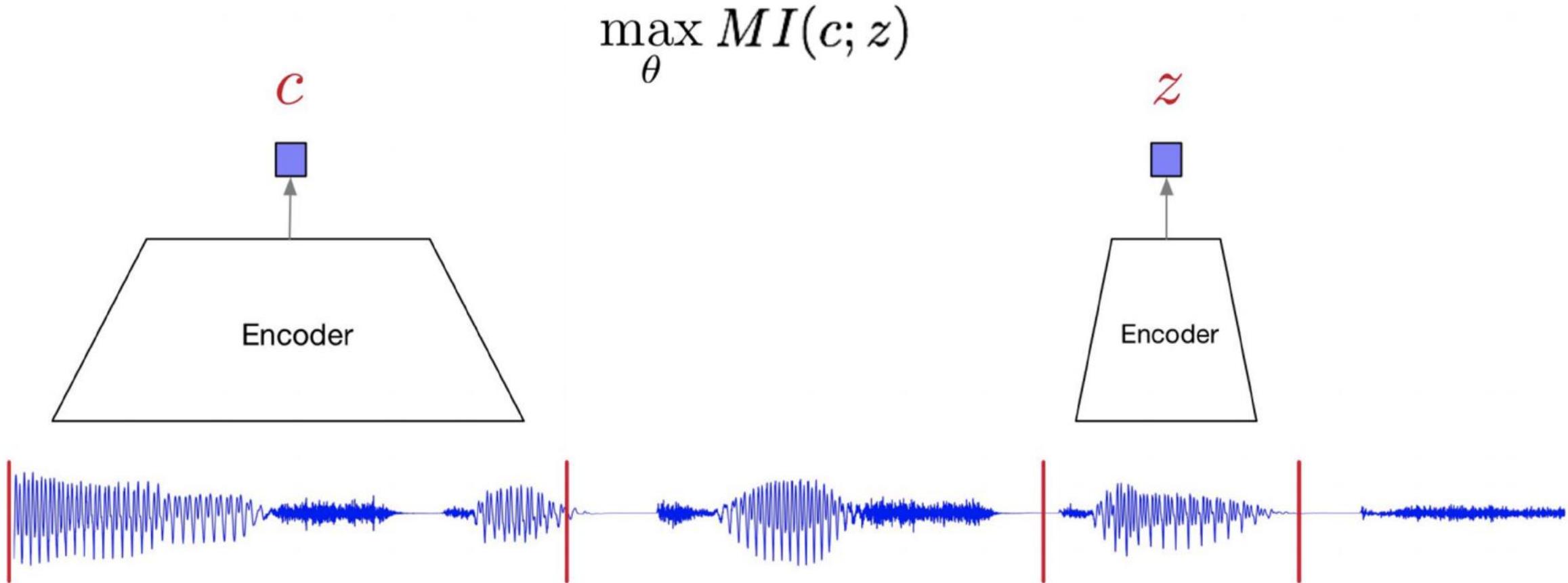
# Contrastive Predictive Coding



# Contrastive Predictive Coding



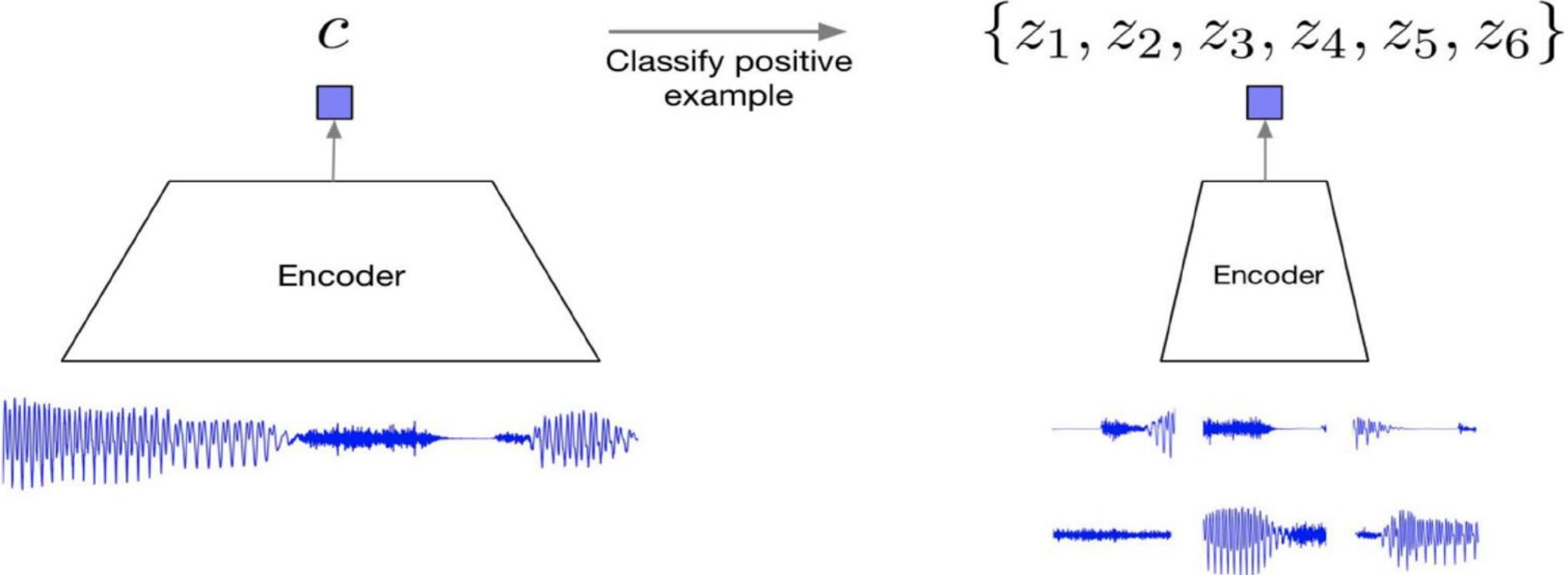
# Contrastive Predictive Coding



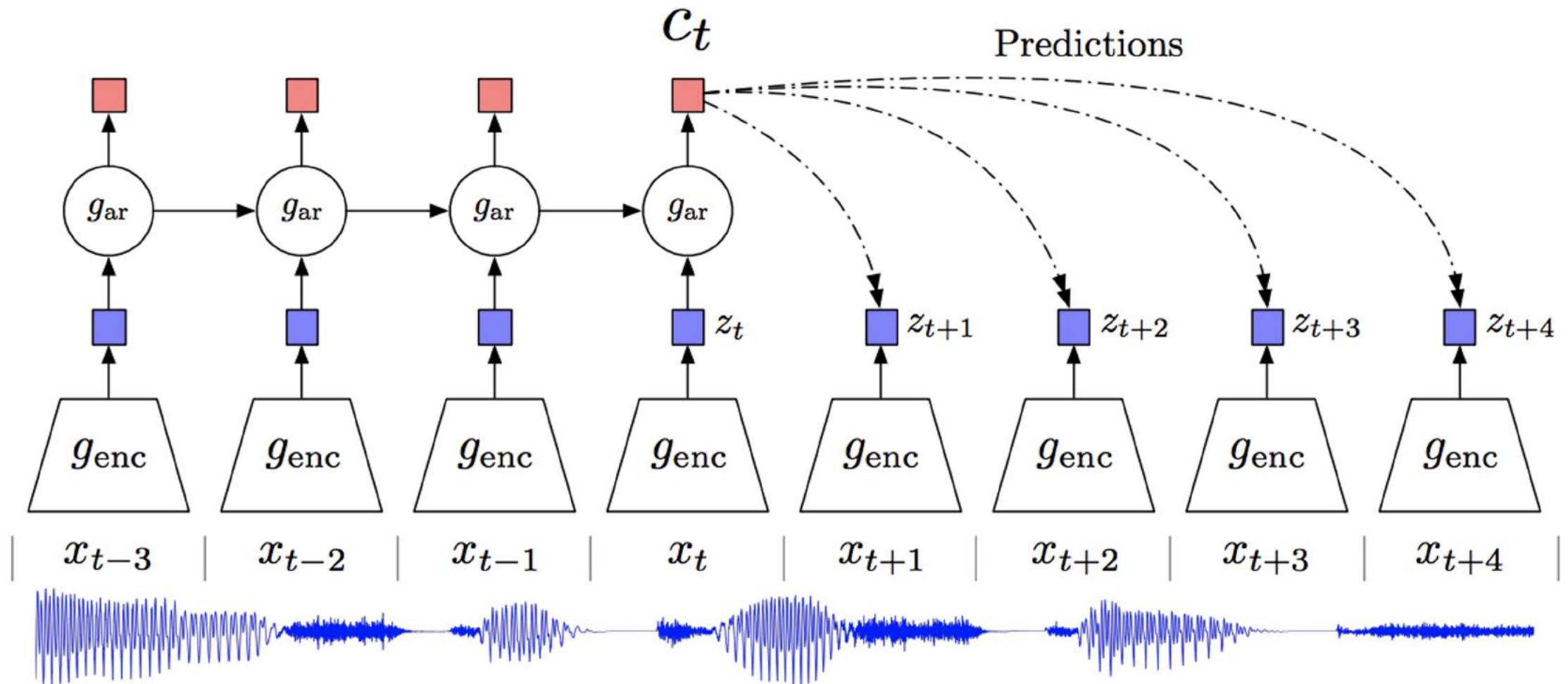
# Contrastive Predictive Coding

$$\frac{\exp f(c, z_i)}{\sum_j \exp f(c, z_j)}$$

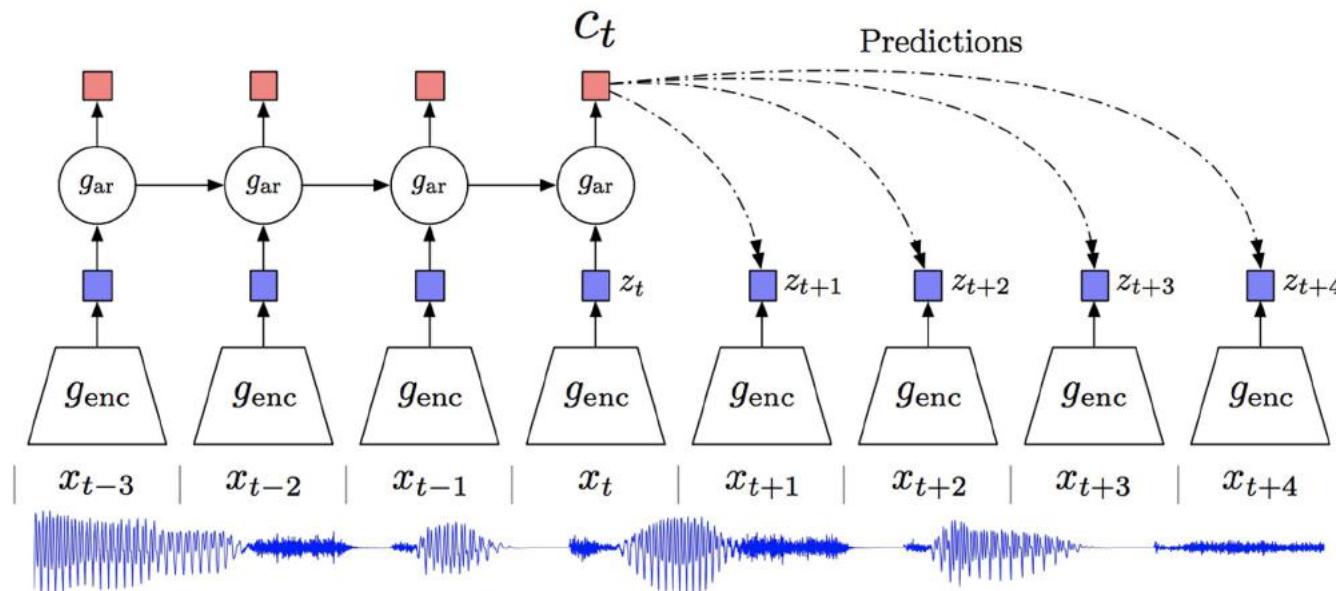
$$f_k(x_{t+k}, c_t) = \exp \left( z_{t+k}^T W_k c_t \right)$$



# Contrastive Predictive Coding



# Contrastive Predictive Coding



$$f_k(x_{t+k}, c_t) = \exp \left( z_{t+k}^T W_k c_t \right)$$

$$\mathcal{L}_{\text{N}} = - \mathbb{E}_X \left[ \log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right]$$

# CPC - Speech

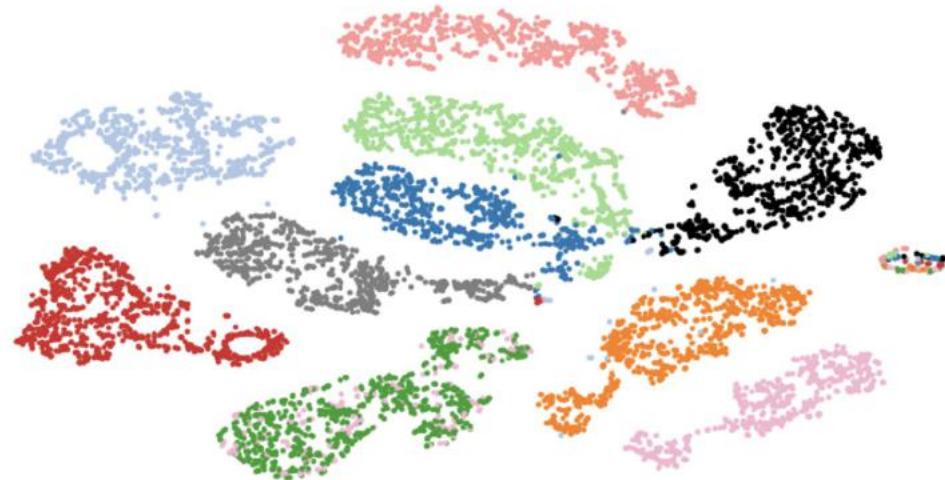


Figure 2: t-SNE visualization of audio (speech) representations for a subset of 10 speakers (out of 251). Every color represents a different speaker.

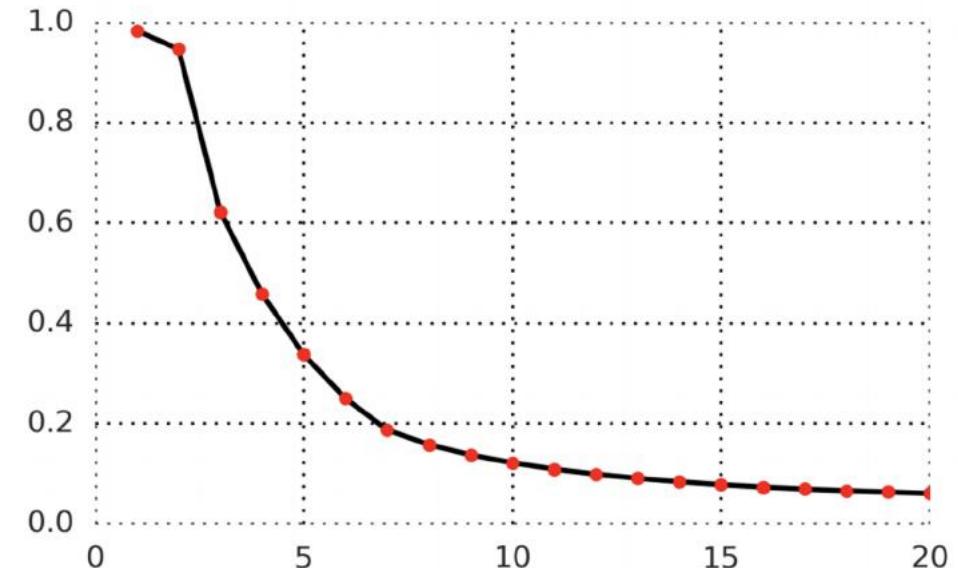


Figure 3: Average accuracy of predicting the positive sample in the contrastive loss for 1 to 20 latent steps in the future of a speech waveform. The model predicts up to 200ms in the future as every step consists of 10ms of audio.

# CPC - Speech

Method	ACC
<b>Phone classification</b>	
Random initialization	27.6
MFCC features	39.7
CPC	64.6
Supervised	74.6
<b>Speaker classification</b>	
Random initialization	1.87
MFCC features	17.6
CPC	97.4
Supervised	98.5

Table 1: LibriSpeech phone and speaker classification results. For phone classification there are 41 possible classes and for speaker classification 251. All models used the same architecture and the same audio input sizes.

Method	ACC
<b>#steps predicted</b>	
2 steps	28.5
4 steps	57.6
8 steps	63.6
12 steps	64.6
16 steps	63.8
<b>Negative samples from</b>	
Mixed speaker	64.6
Same speaker	65.5
Mixed speaker (excl.)	57.3
Same speaker (excl.)	64.6
Current sequence only	65.2

Table 2: LibriSpeech phone classification ablation experiments. More details can be found in Section 3.1.

# CPC - ImageNet

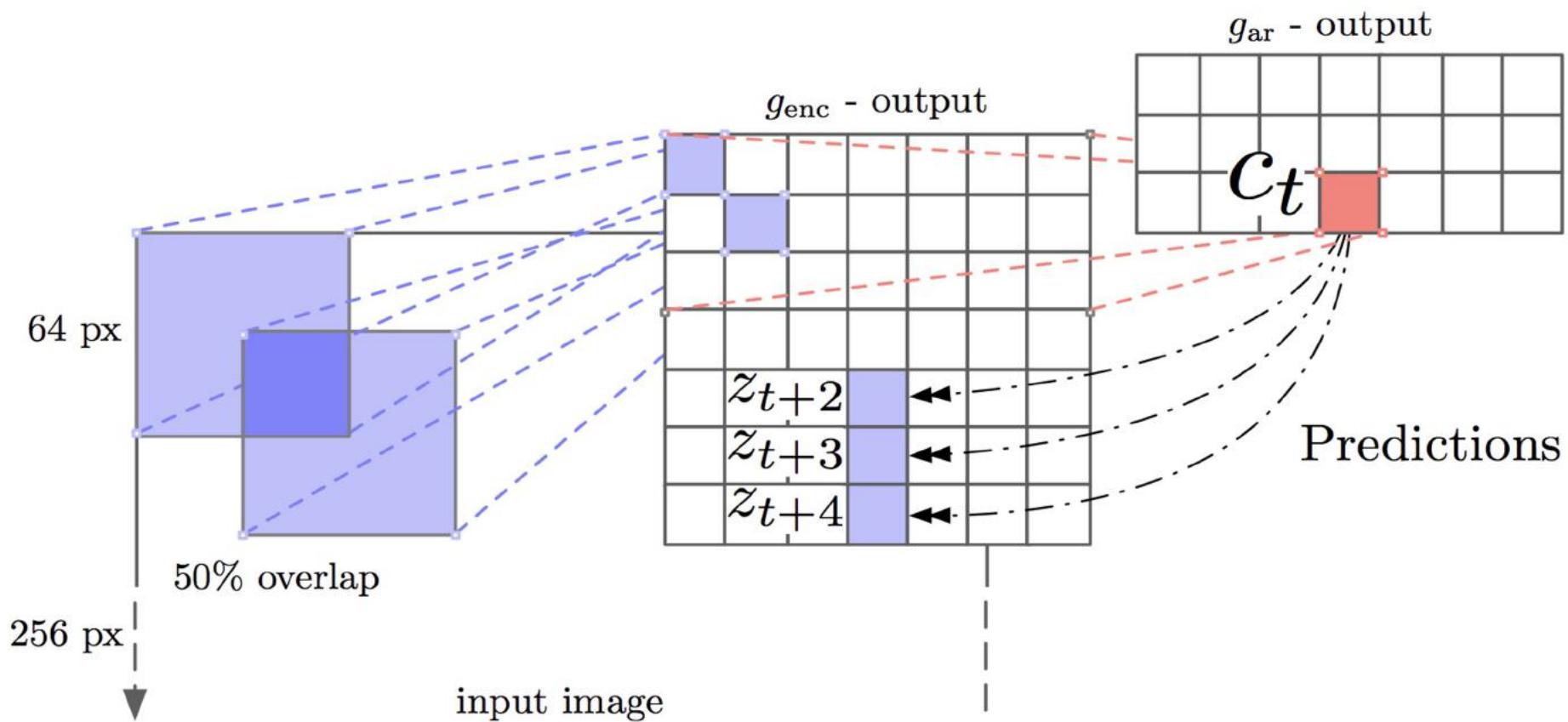


Figure 4: Visualization of Contrastive Predictive Coding for images (2D adaptation of Figure 1).

# CPC - ImageNet



# CPC - ImageNet

Method	Top-1 ACC
<b>Using AlexNet conv5</b>	
Video [28]	29.8
Relative Position [11]	30.4
BiGan [35]	34.8
Colorization [10]	35.2
Jigsaw [29] *	38.1
<b>Using ResNet-V2</b>	
Motion Segmentation [36]	27.6
Exemplar [36]	31.5
Relative Position [36]	36.2
Colorization [36]	39.6
<b>CPC</b>	<b>48.7</b>

Table 3: ImageNet top-1 unsupervised classification results. \*Jigsaw is not directly comparable to the other AlexNet results because of architectural differences.

Method	Top-5 ACC
Motion Segmentation (MS)	48.3
Exemplar (Ex)	53.1
Relative Position (RP)	59.2
Colorization (Col)	62.5
Combination of MS + Ex + RP + Col	69.3
<b>CPC</b>	<b>73.6</b>

Table 4: ImageNet top-5 unsupervised classification results. Previous results with MS, Ex, RP and Col were taken from [36] and are the best reported results on this task.

# CPC - ImageNet

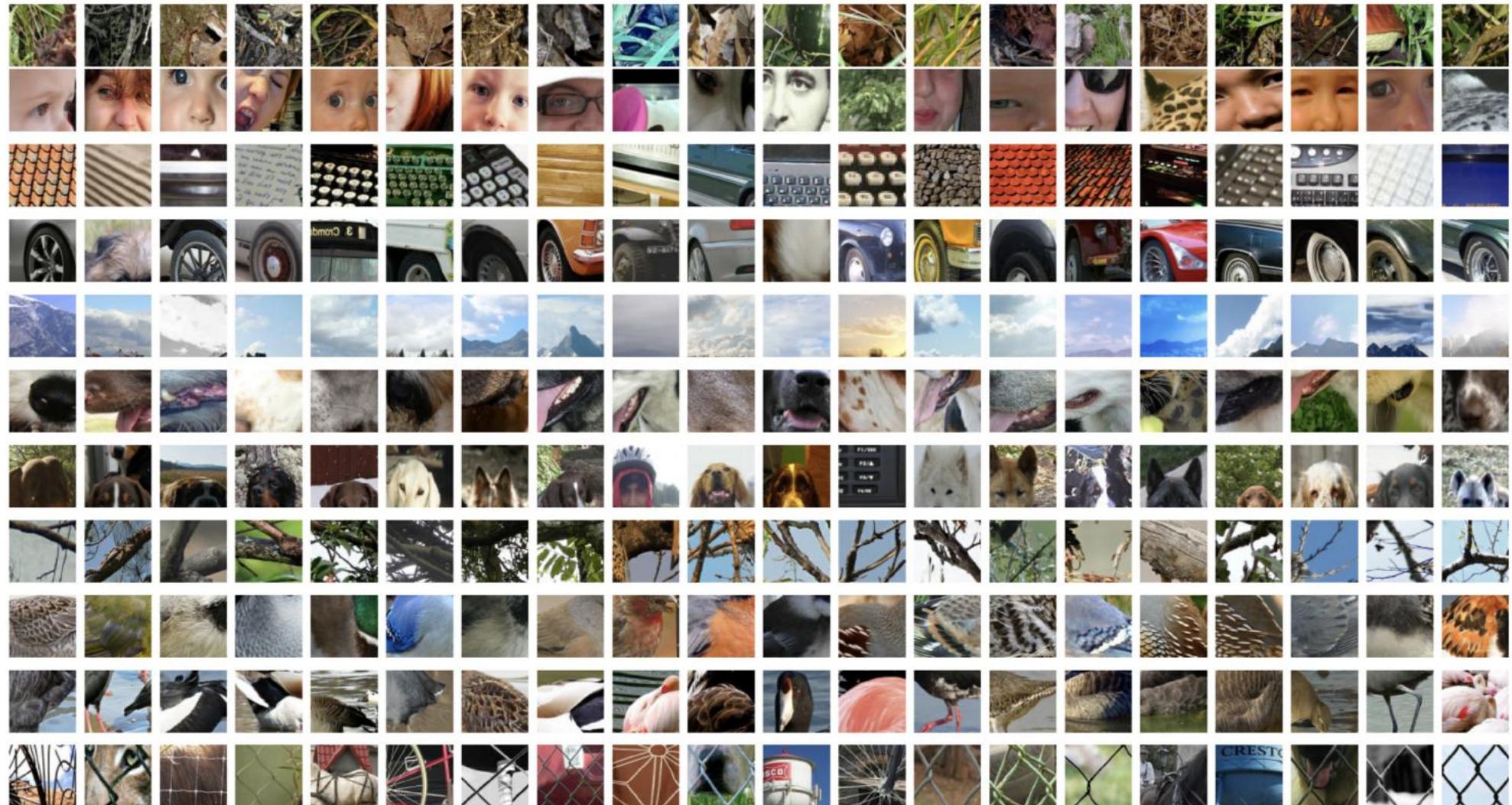


Figure 5: Every row shows image patches that activate a certain neuron in the CPC architecture.

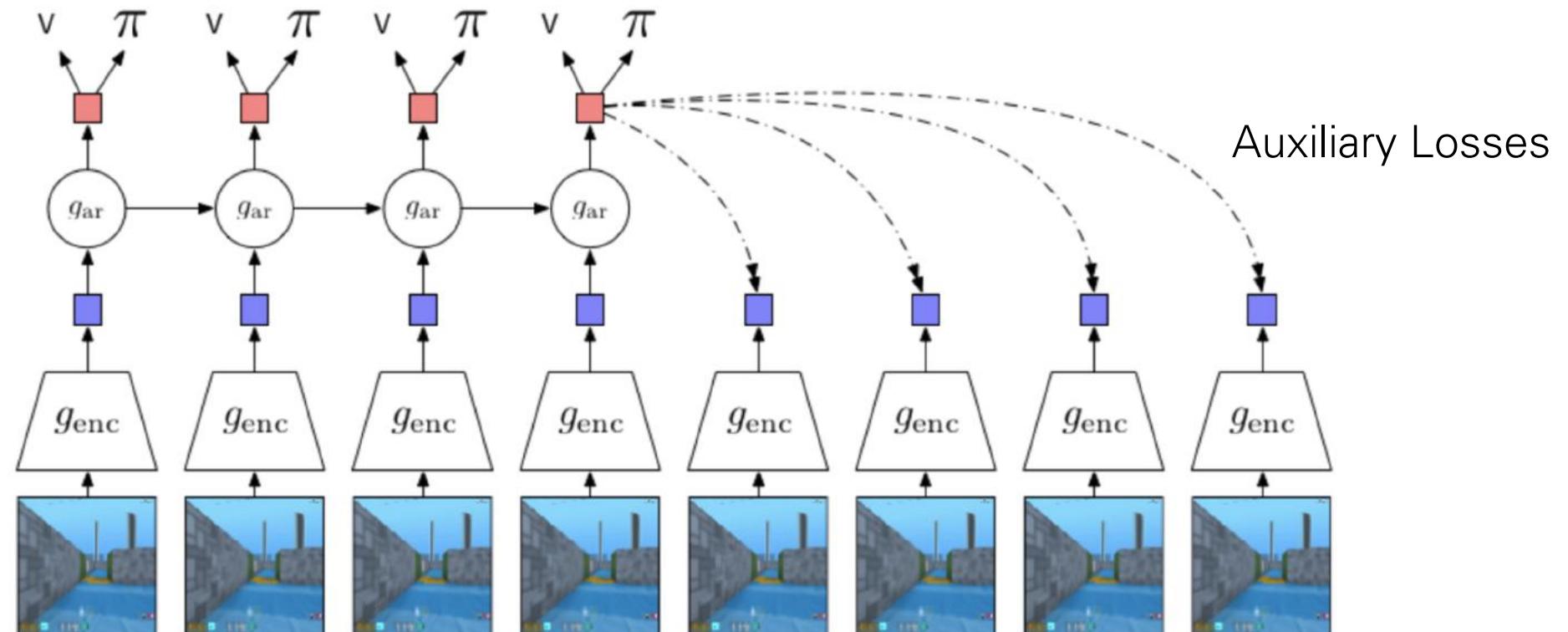
# CPC - Natural Language Processing

<b>Method</b>	<b>MR</b>	<b>CR</b>	<b>Subj</b>	<b>MPQA</b>	<b>TREC</b>
Paragraph-vector [40]	74.8	78.1	90.5	74.2	91.8
Skip-thought vector [26]	75.5	79.3	92.1	86.9	91.4
Skip-thought + LN [41]	79.5	82.6	93.4	89.0	-
CPC	76.9	80.1	91.2	87.7	96.8

Table 5: Classification accuracy on five common NLP benchmarks. We follow the same transfer learning setup from Skip-thought vectors [26] and use the BookCorpus dataset as source. [40] is an unsupervised approach to learning sentence-level representations. [26] is an alternative unsupervised learning approach. [41] is the same skip-thought model with layer normalization trained for 1M iterations.

# CPC - Reinforcement Learning

Auxiliary loss is on policy  
Predict 30 steps in the future



# CPCv2 - Large Scale CPC on ImageNet

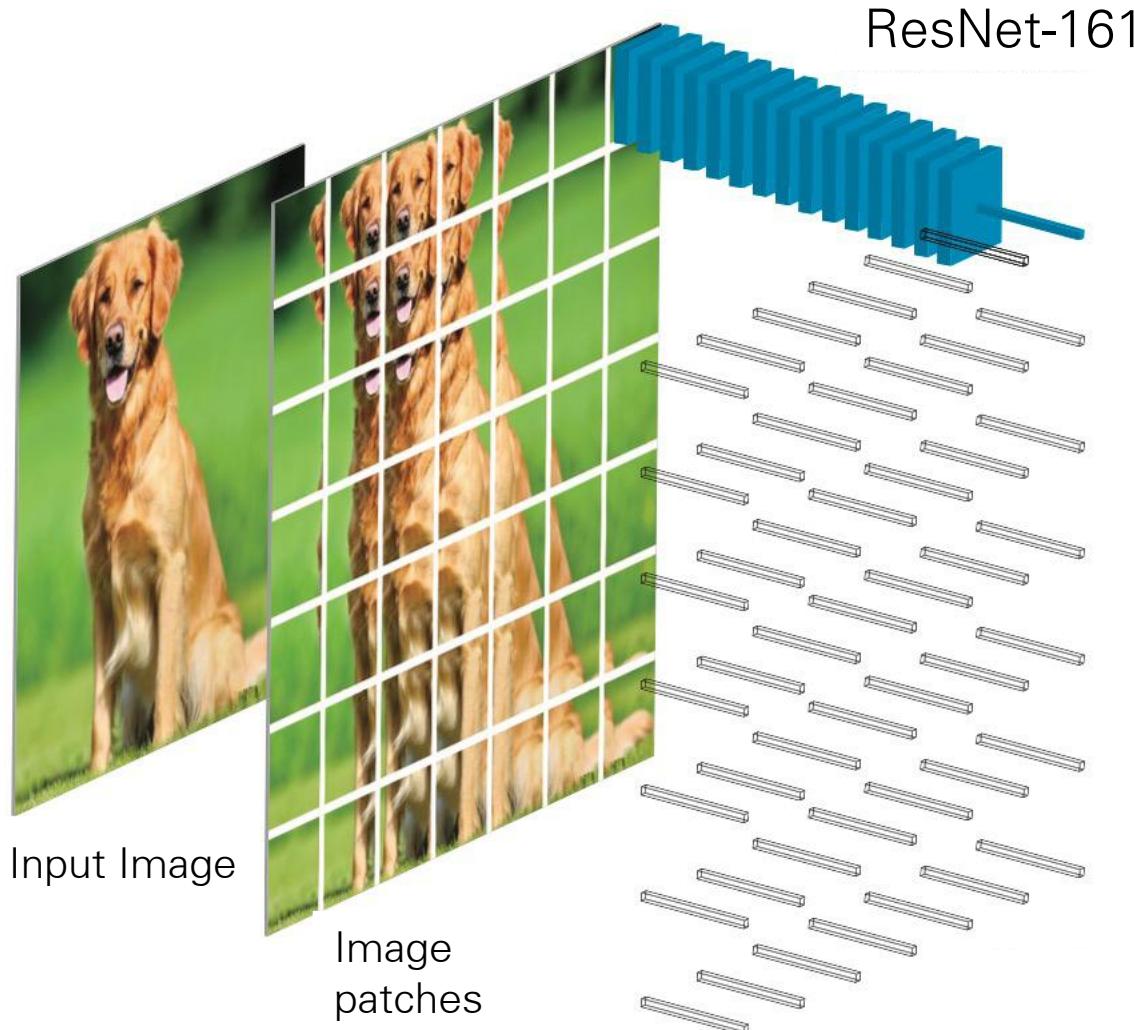
## DATA-EFFICIENT IMAGE RECOGNITION WITH CONTRASTIVE PREDICTIVE CODING

**Olivier J. Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi,  
Carl Doersch, S. M. Ali Eslami, Aaron van den Oord**  
DeepMind  
London, UK

### ABSTRACT

Human observers can learn to recognize new categories of images from a handful of examples, yet doing so with machine perception remains an open challenge. We hypothesize that data-efficient recognition is enabled by representations which make the variability in natural signals more predictable. We therefore revisit and improve Contrastive Predictive Coding, an unsupervised objective for learning such representations. This new implementation produces features which support state-of-the-art linear classification accuracy on the ImageNet dataset. When used as input for non-linear classification with deep neural networks, this representation allows us to use  $2\text{--}5\times$  less labels than classifiers trained directly on image pixels. Finally, this unsupervised representation substantially improves transfer learning to object detection on PASCAL VOC-2007, surpassing fully supervised pre-trained ImageNet classifiers.

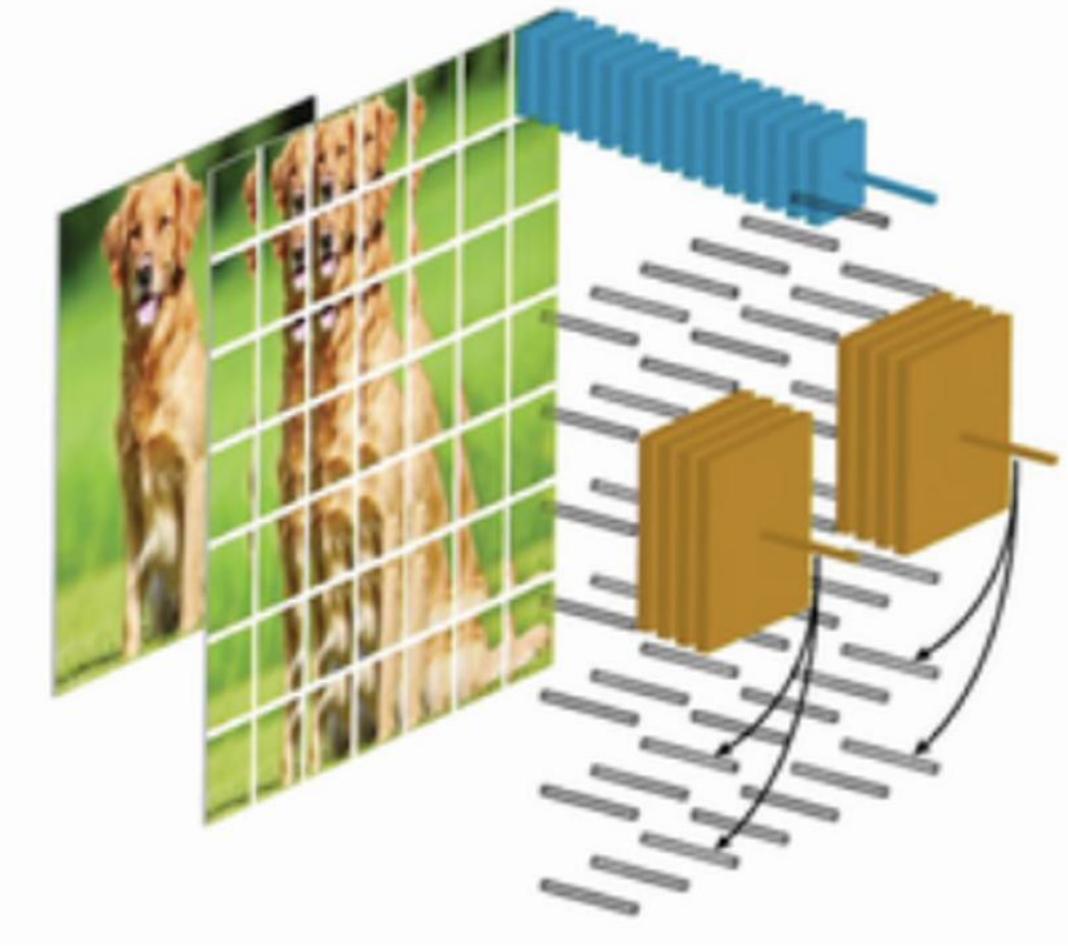
# CPCv2 - Large Scale CPC on ImageNet



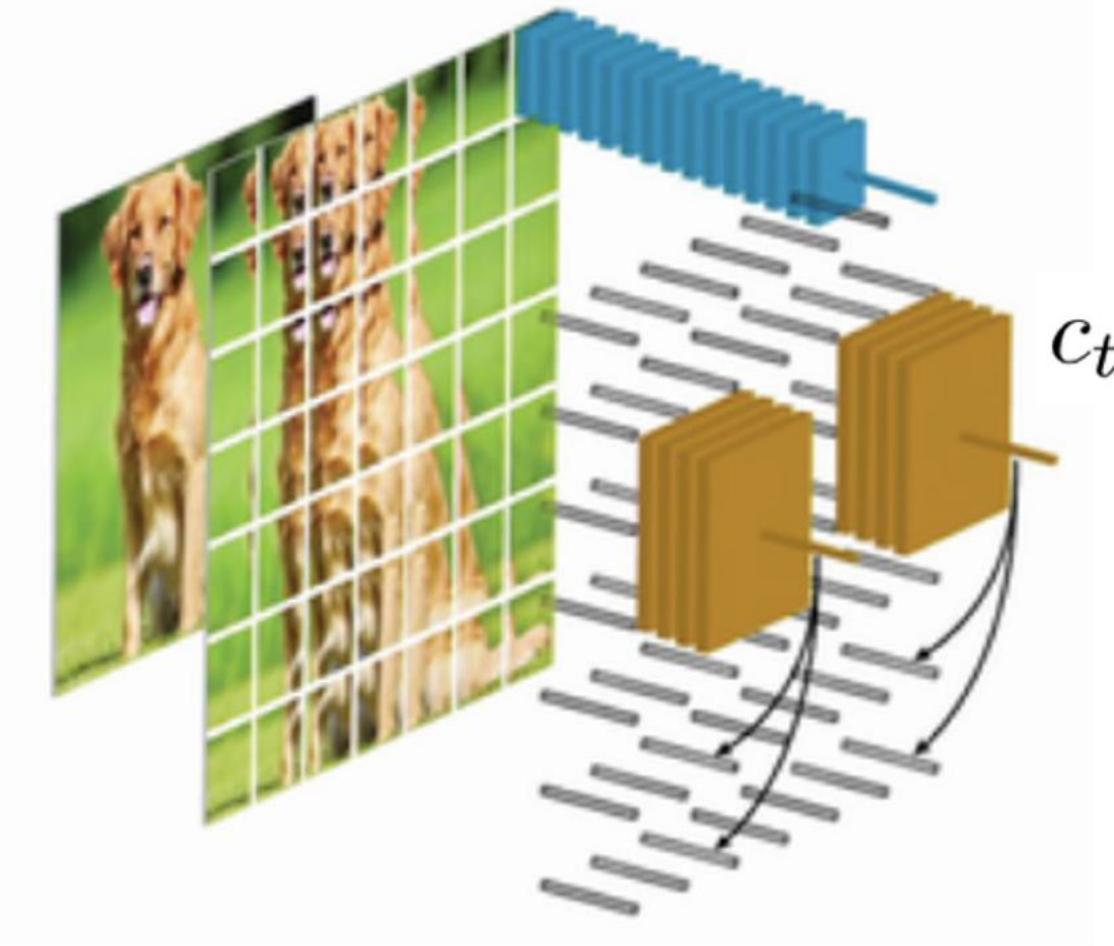
# CPCv2 - Large Scale CPC on ImageNet



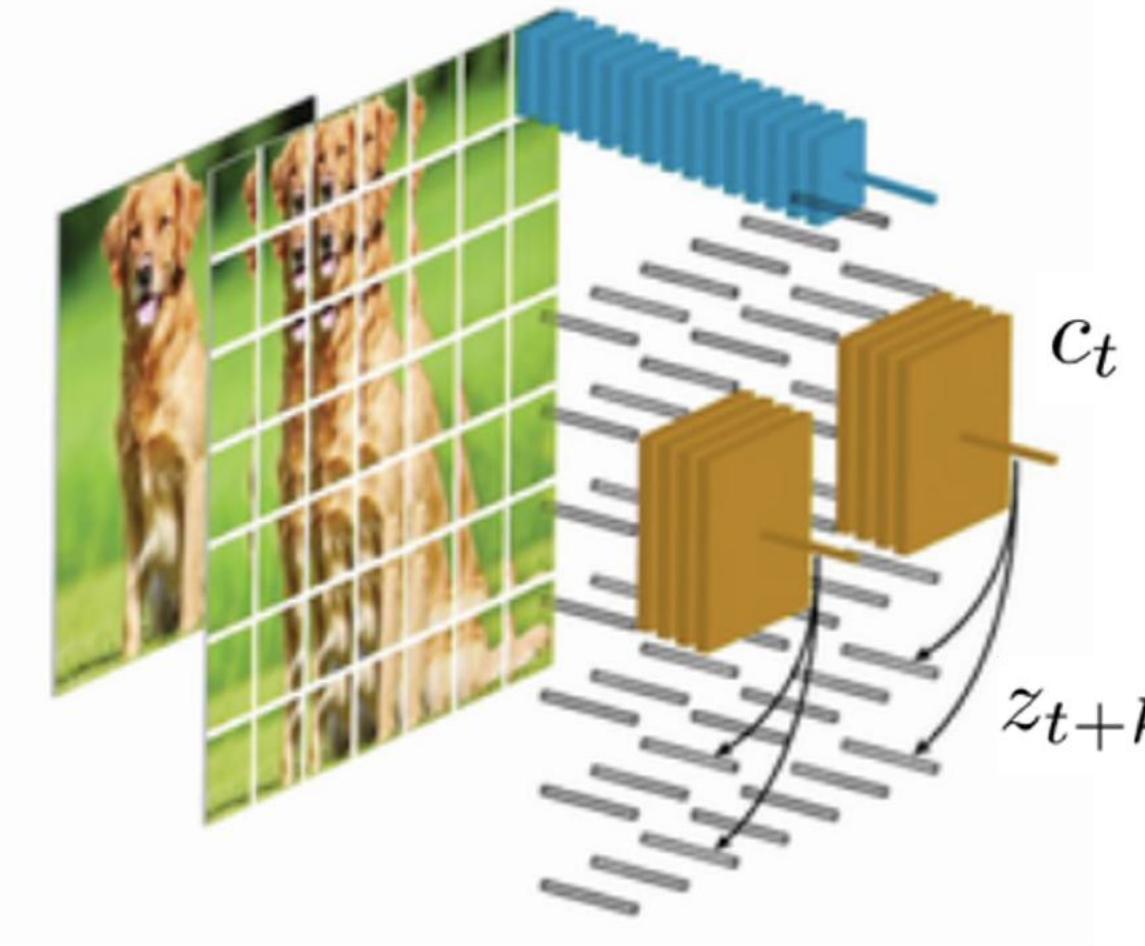
# CPCv2 - Large Scale CPC on ImageNet



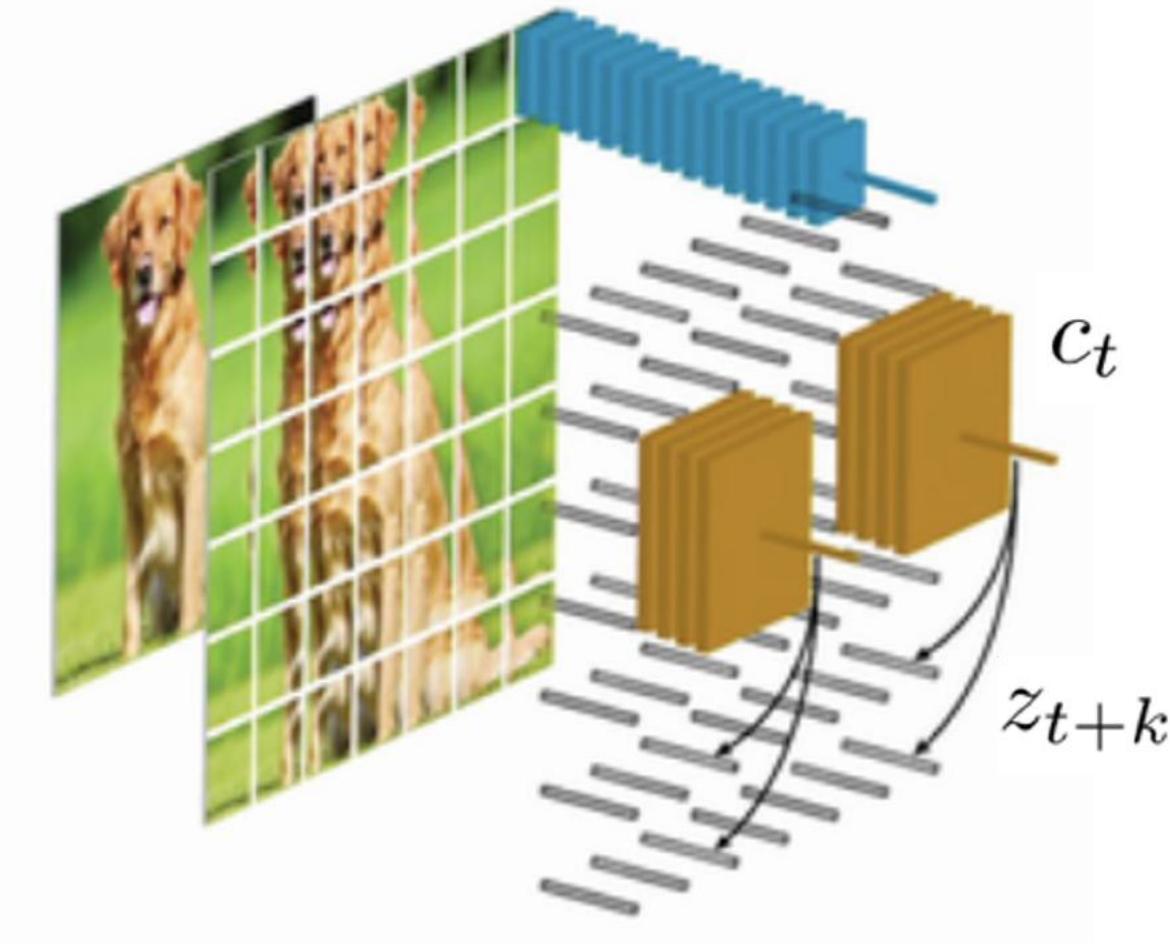
# CPCv2 - Large Scale CPC on ImageNet



# CPCv2 - Large Scale CPC on ImageNet

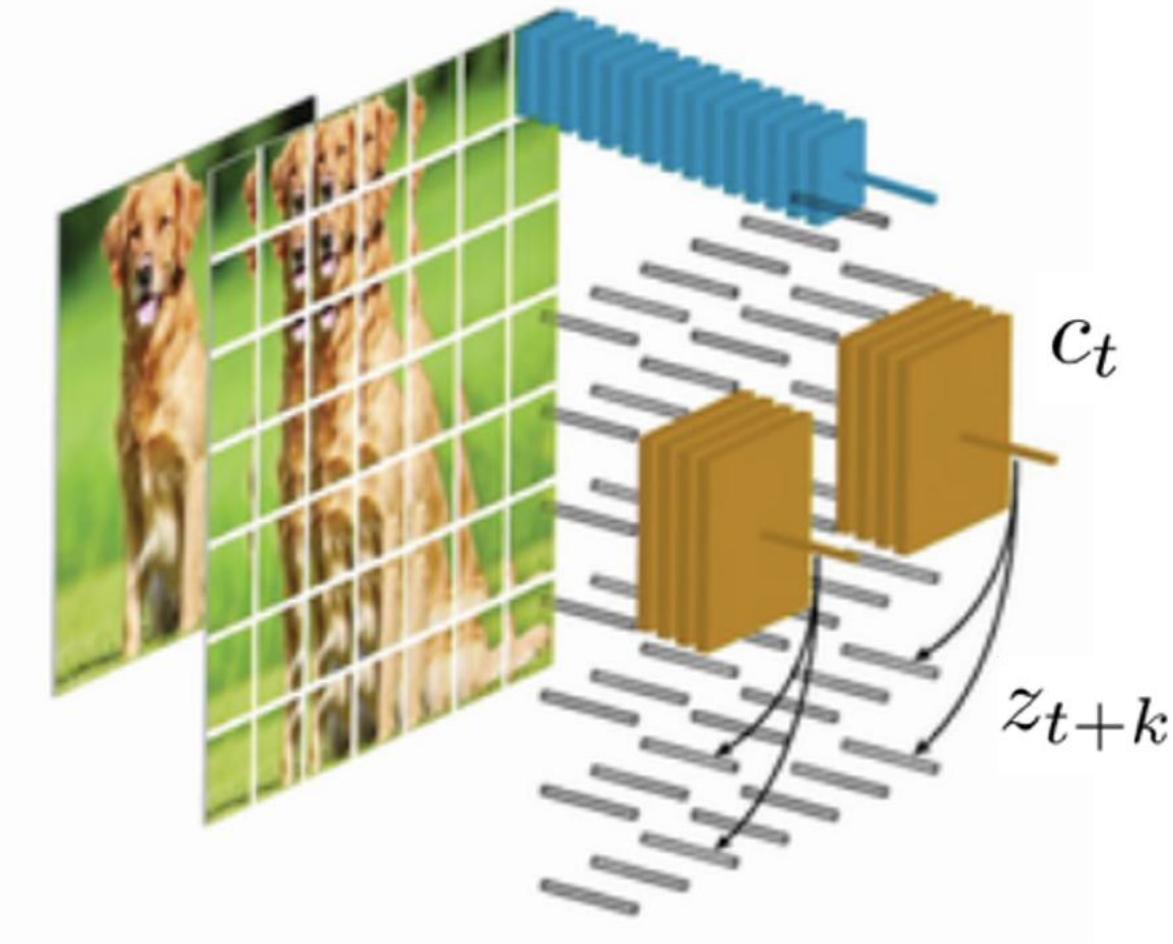


# CPCv2 - Large Scale CPC on ImageNet



$$f_k(x_{t+k}, c_t) = \exp \left( z_{t+k}^T W_k c_t \right)$$

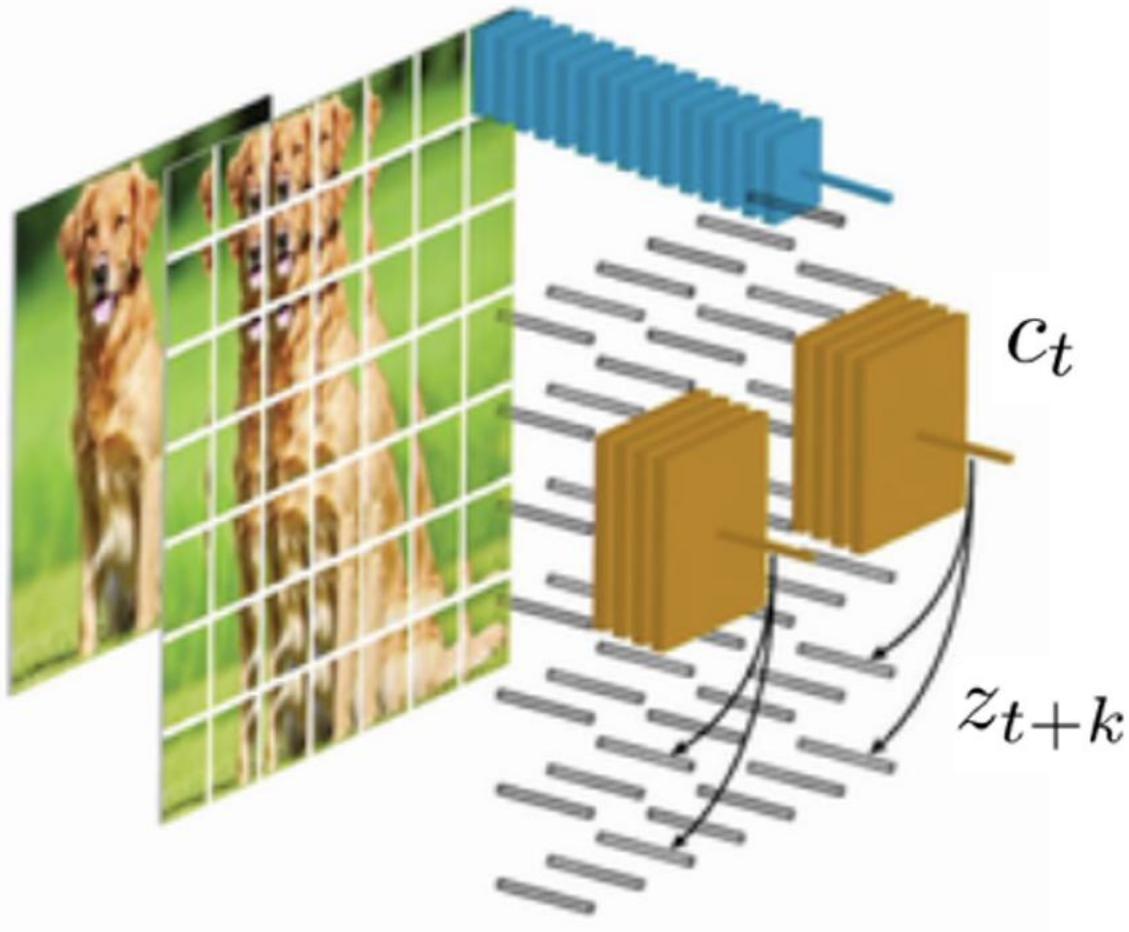
# CPCv2 - Large Scale CPC on ImageNet



$$f_k(x_{t+k}, c_t) = \exp\left(z_{t+k}^T W_k c_t\right)$$

$$\mathcal{L}_N = -\mathbb{E}_X \left[ \log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right]$$

# CPCv2 - Large Scale CPC on ImageNet



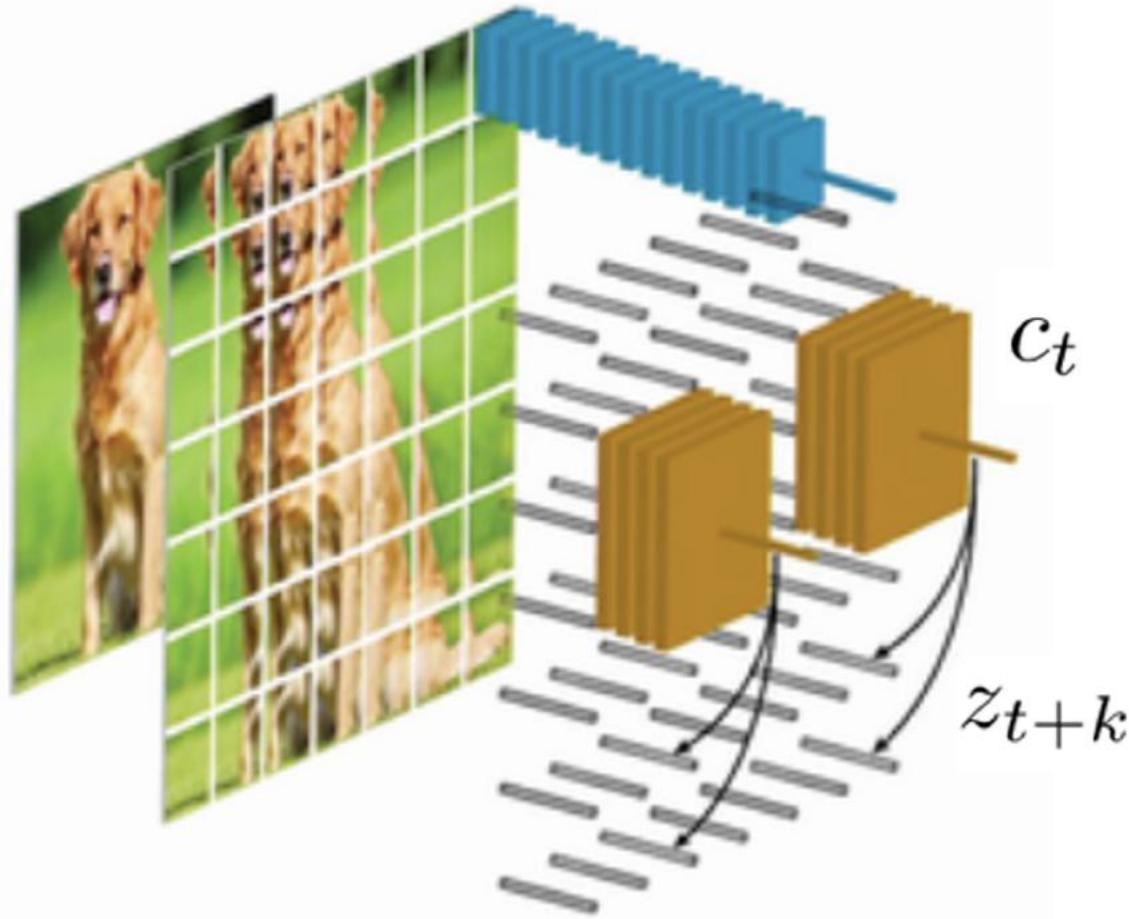
$$f_k(x_{t+k}, c_t) = \exp(z_{t+k}^T W_k c_t)$$

$$\mathcal{L}_N = -\mathbb{E}_X \left[ \log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right]$$

↓ Negatives

1. Other patches within image
2. Patches from other images

# CPCv2 - Large Scale CPC on ImageNet



## InfoNCE Loss

NCE: Noise-Contrastive Estimation

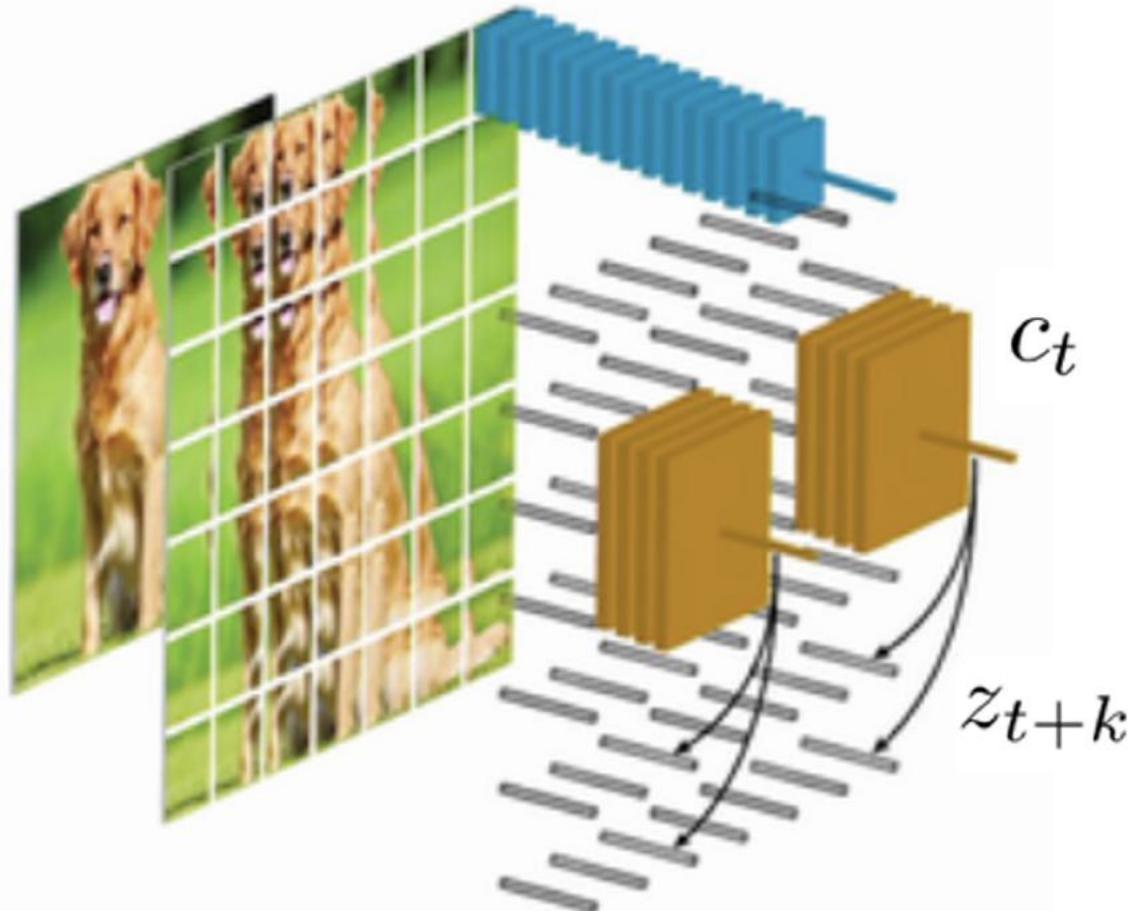
$$f_k(x_{t+k}, c_t) = \exp\left(z_{t+k}^T W_k c_t\right)$$

$$\mathcal{L}_N = -\mathbb{E}_X \left[ \log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right]$$

↓ Negatives

1. Other patches within image
2. Patches from other images

# CPCv2 - Large Scale CPC on ImageNet



**Parallel Implementation  
with PixelCNN (masked conv) and  
1x1 conv**

## InfoNCE Loss

NCE: Noise-Contrastive Estimation

$$f_k(x_{t+k}, c_t) = \exp\left(z_{t+k}^T W_k c_t\right)$$

$$\mathcal{L}_N = -\mathbb{E}_X \left[ \log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right]$$

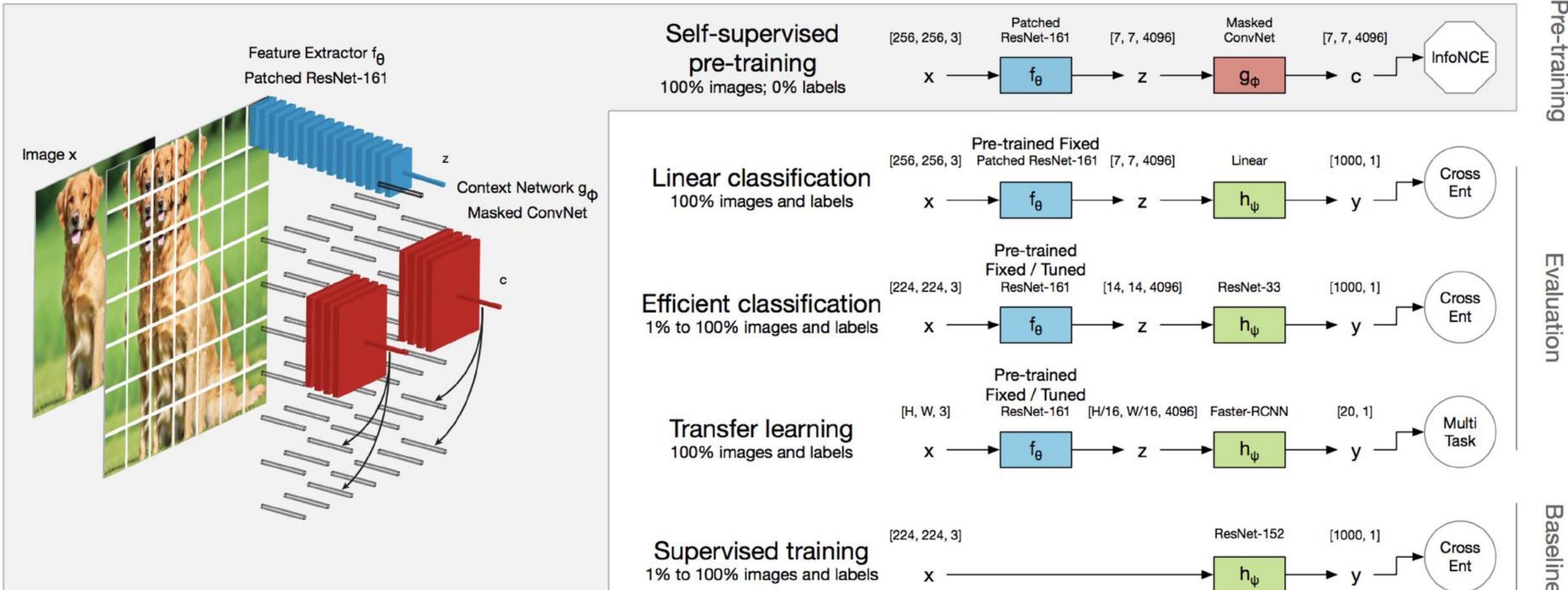
↓  
Negatives

1. Other patches within image
2. Patches from other images

# CPCv2 - Large Scale CPC on ImageNet

- Train CPC on unlabeled ImageNet
- Train as long as possible (500 epochs) – 1 week
- Augment every patch with a lot of spatial and color augmentation  
**[extremely crucial]**
- Effective number of negatives = number of instances \* number of patches per instance =  $16 * 36 = 576$

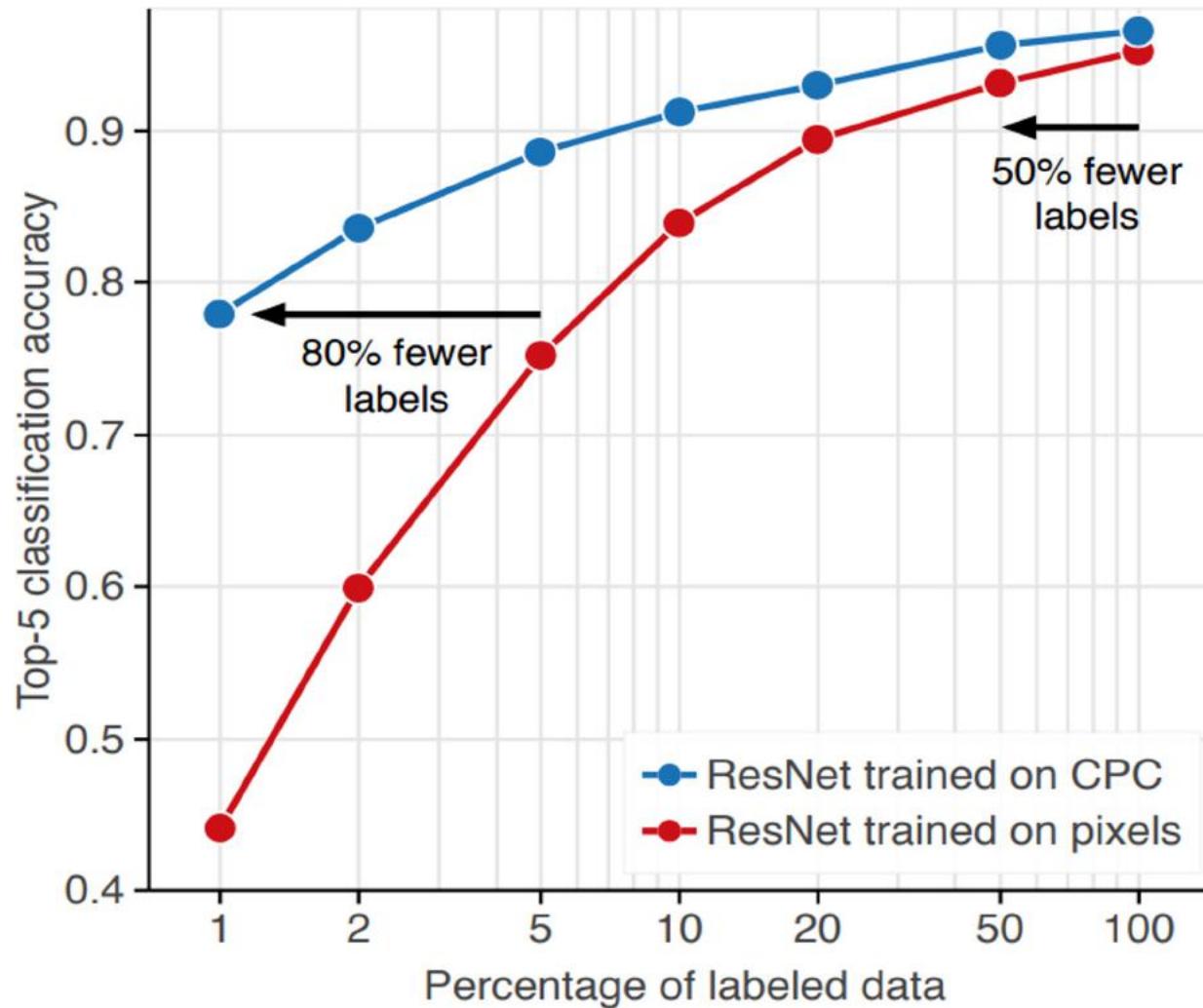
# CPCv2 - Large Scale CPC on ImageNet



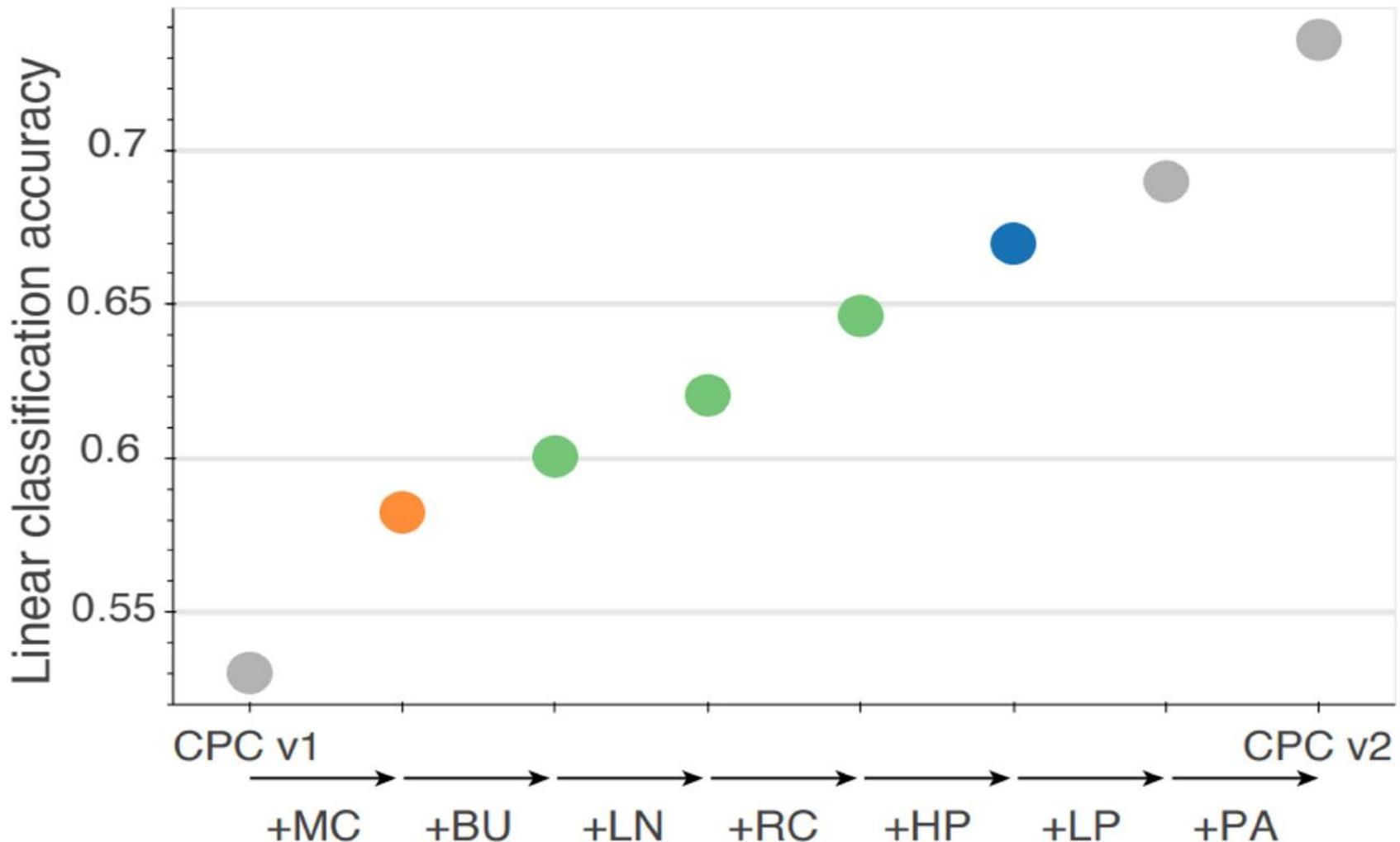
# CPCv2 - Linear Classification on ImageNet

METHOD	PARAMS (M)	TOP-1	TOP-5
<i>Methods using ResNet-50:</i>			
INSTANCE DISCR. [1]	24	54.0	-
LOCAL AGGR. [2]	24	58.8	-
MoCo [3]	24	60.6	-
PIRL [4]	24	63.6	-
CPC v2 - RESNET-50	24	<b>63.8</b>	<b>85.3</b>
<i>Methods using different architectures:</i>			
MULTI-TASK [5]	28	-	69.3
ROTATION [6]	86	55.4	-
CPC v1 [7]	28	48.7	73.6
BIGBiGAN [8]	86	61.3	81.9
AMDIM [9]	626	68.1	-
CMC [10]	188	68.4	88.2
MoCo [2]	375	68.6	-
CPC v2 - RESNET-161	305	<b>71.5</b>	<b>90.1</b>

# CPCv2 - Data-Efficient Image Recognition



# CPCv1 → CPCv2



- MC:** model capacity
- BU:** bottom-up spatial predictions
- LN:** layer normalization
- RC:** random color-dropping
- HP:** horizontal spatial predictions
- LP:** larger patches.
- PA:** further patch-based augmentation.

# CPCv2 - Data-Efficient Supervised Learning

<b>Method</b>	<b>Architecture</b>	<b>Top-5 accuracy</b>				
		1%	5%	10%	50%	100%
Labeled data						
† Supervised baseline						
	ResNet-200	44.1	75.2*	83.9	93.1	95.2#
<b>Methods using label-propagation:</b>						
Pseudolabeling [63]	ResNet-50	51.6	-	82.4	-	-
VAT + Entropy Minimization [63]	ResNet-50	47.0	-	83.4	-	-
Unsup. Data Augmentation [61]	ResNet-50	-	-	88.5	-	-
Rotation + VAT + Ent. Min. [63]	ResNet-50 $\times 4$	-	-	<b>91.2</b>	-	95.0
<b>Methods using representation learning only:</b>						
Instance Discrimination [60]	ResNet-50	39.2	-	77.4	-	-
Rotation [63]	ResNet-152 $\times 2$	57.5	-	86.4	-	-
ResNet on BigBiGAN (fixed)	RevNet-50 $\times 4$	55.2	73.7	78.8	85.5	87.0
ResNet on AMDIM (fixed)	Custom-103	67.4	81.8	85.8	91.0	92.2
ResNet on CPC v2 (fixed)	ResNet-161	77.1	87.5	90.5	95.0	96.2
ResNet on CPC v2 (fine-tuned)	ResNet-161	<b>77.9*</b>	<b>88.6</b>	<b>91.2</b>	<b>95.6#</b>	<b>96.5</b>

# CPCv2 - PASCAL VOC-07 Detection

Method	Architecture	mAP
<b><i>Transfer from labeled data:</i></b>		
Supervised baseline	ResNet-152	74.7
<b><i>Transfer from unlabeled data:</i></b>		
Exemplar [17] by [13]	ResNet-101	60.9
Motion Segmentation [47] by [13]	ResNet-101	61.1
Colorization [64] by [13]	ResNet-101	65.5
Relative Position [14] by [13]	ResNet-101	66.8
Multi-task [13]	ResNet-101	70.5
Instance Discrimination [60]	ResNet-50	65.4
Deep Cluster [7]	VGG-16	65.9
Deeper Cluster [8]	VGG-16	67.8
Local Aggregation [66]	ResNet-50	69.1
Momentum Contrast [25]	ResNet-50	74.9
Faster-RCNN trained on CPC v2	ResNet-161	<b>76.6</b>

# Instance Discrimination



# Instance Discrimination



1. MoCo
2. SimCLR

# Momentum Contrast (MoCo)

## Momentum Contrast for Unsupervised Visual Representation Learning

Kaiming He Haoqi Fan Yuxin Wu Saining Xie Ross Girshick

### Abstract

*We present Momentum Contrast (MoCo) for unsupervised visual representation learning. From a perspective on contrastive learning [29] as dictionary look-up, we build a dynamic dictionary with a queue and a moving-averaged encoder. This enables building a large and consistent dictionary on-the-fly that facilitates contrastive unsupervised learning. MoCo provides competitive results under the common linear protocol on ImageNet classification. More importantly, the representations learned by MoCo transfer well to downstream tasks. MoCo can **outperform** its supervised pre-training counterpart in 7 detection/segmentation tasks on PASCAL VOC, COCO, and other datasets, sometimes surpassing it by large margins. This suggests that the gap between unsupervised and supervised representation learning has been largely closed in many vision tasks.*

# Momentum Contrast (MoCo)

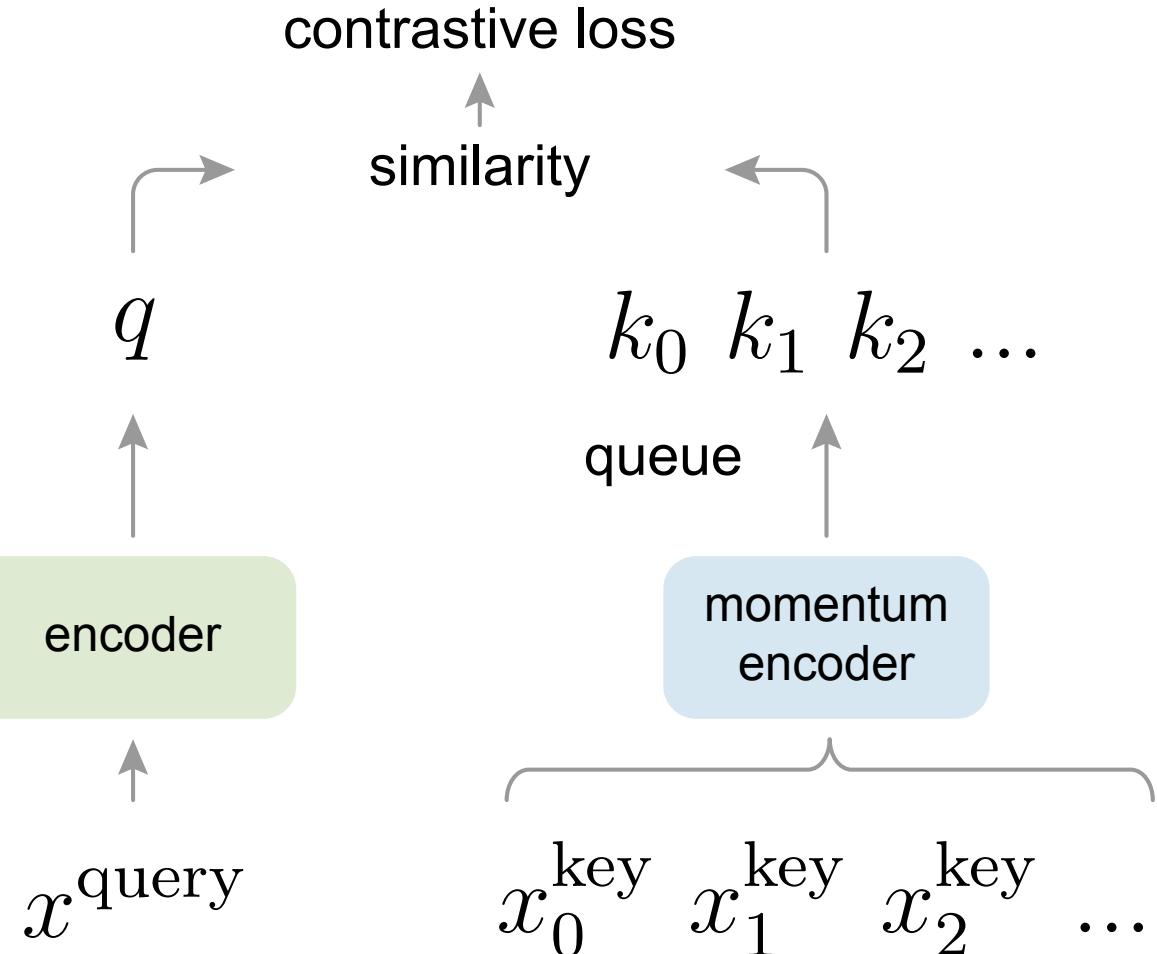
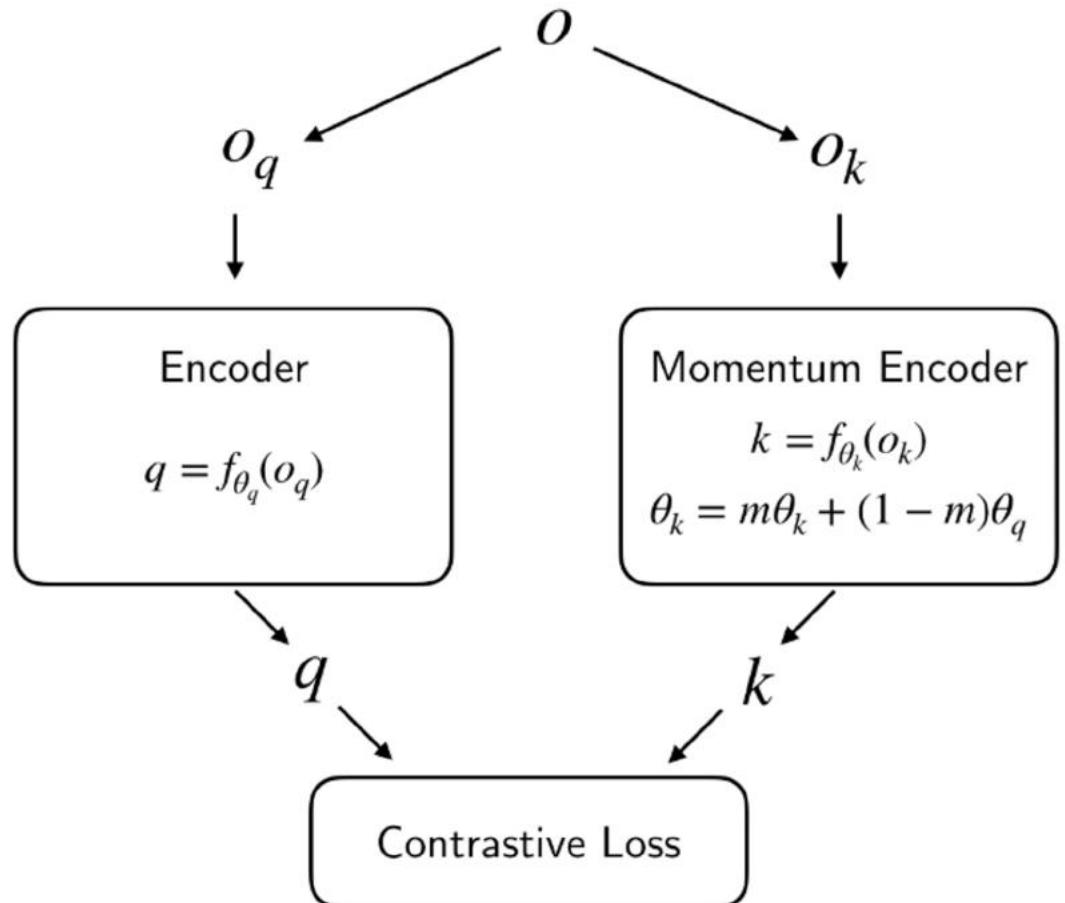


Figure 1. Momentum Contrast (MoCo) trains a visual representation encoder by matching an encoded query  $q$  to a dictionary of encoded keys using a contrastive loss. The dictionary keys  $\{k_0, k_1, k_2, \dots\}$  are defined on-the-fly by a set of data samples. The dictionary is built as a queue, with the current mini-batch enqueued and the oldest mini-batch dequeued, decoupling it from the mini-batch size. The keys are encoded by a slowly progressing encoder, driven by a momentum update with the query encoder. This method enables a large and consistent dictionary for learning visual representations.

# Momentum Contrast (MoCo)



$$\mathcal{L}_q = - \log \frac{\exp(q \cdot k_+ / \tau)}{\sum_{i=0}^K \exp(q \cdot k_i / \tau)}$$

# Momentum Contrast (MoCo)

---

**Algorithm 1** Pseudocode of MoCo in a PyTorch-like style.

---

```
# f_q, f_k: encoder networks for query and key
# queue: dictionary as a queue of K keys (CxK)
# m: momentum
# t: temperature

f_k.params = f_q.params # initialize
for x in loader: # load a minibatch x with N samples
    x_q = aug(x) # a randomly augmented version
    x_k = aug(x) # another randomly augmented version

    q = f_q.forward(x_q) # queries: NxC
    k = f_k.forward(x_k) # keys: NxC
    k = k.detach() # no gradient to keys

    # positive logits: Nx1
    l_pos = bmm(q.view(N,1,C), k.view(N,C,1))

    # negative logits: NxK
    l_neg = mm(q.view(N,C), queue.view(C,K))

    # logits: Nx(1+K)
    logits = cat([l_pos, l_neg], dim=1)

    # contrastive loss, Eqn.(1)
    labels = zeros(N) # positives are the 0-th
    loss = CrossEntropyLoss(logits/t, labels)

    # SGD update: query network
    loss.backward()
    update(f_q.params)

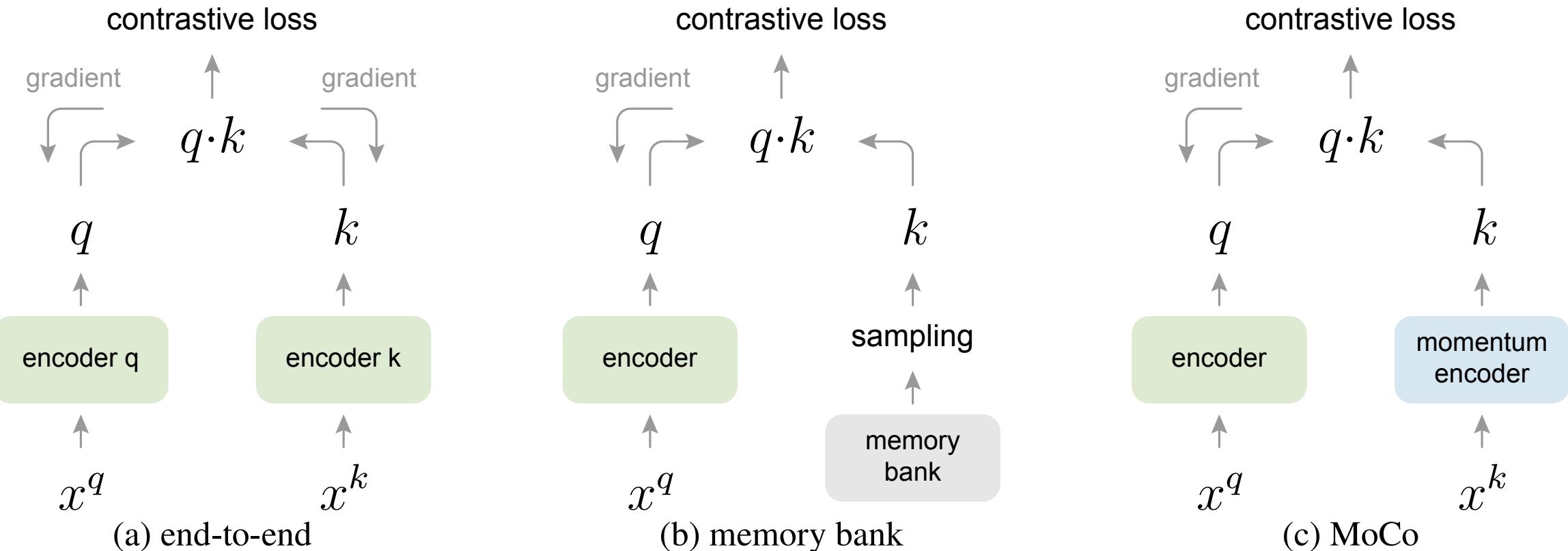
    # momentum update: key network
    f_k.params = m*f_k.params+(1-m)*f_q.params

    # update dictionary
    enqueue(queue, k) # enqueue the current minibatch
    dequeue(queue) # dequeue the earliest minibatch
```

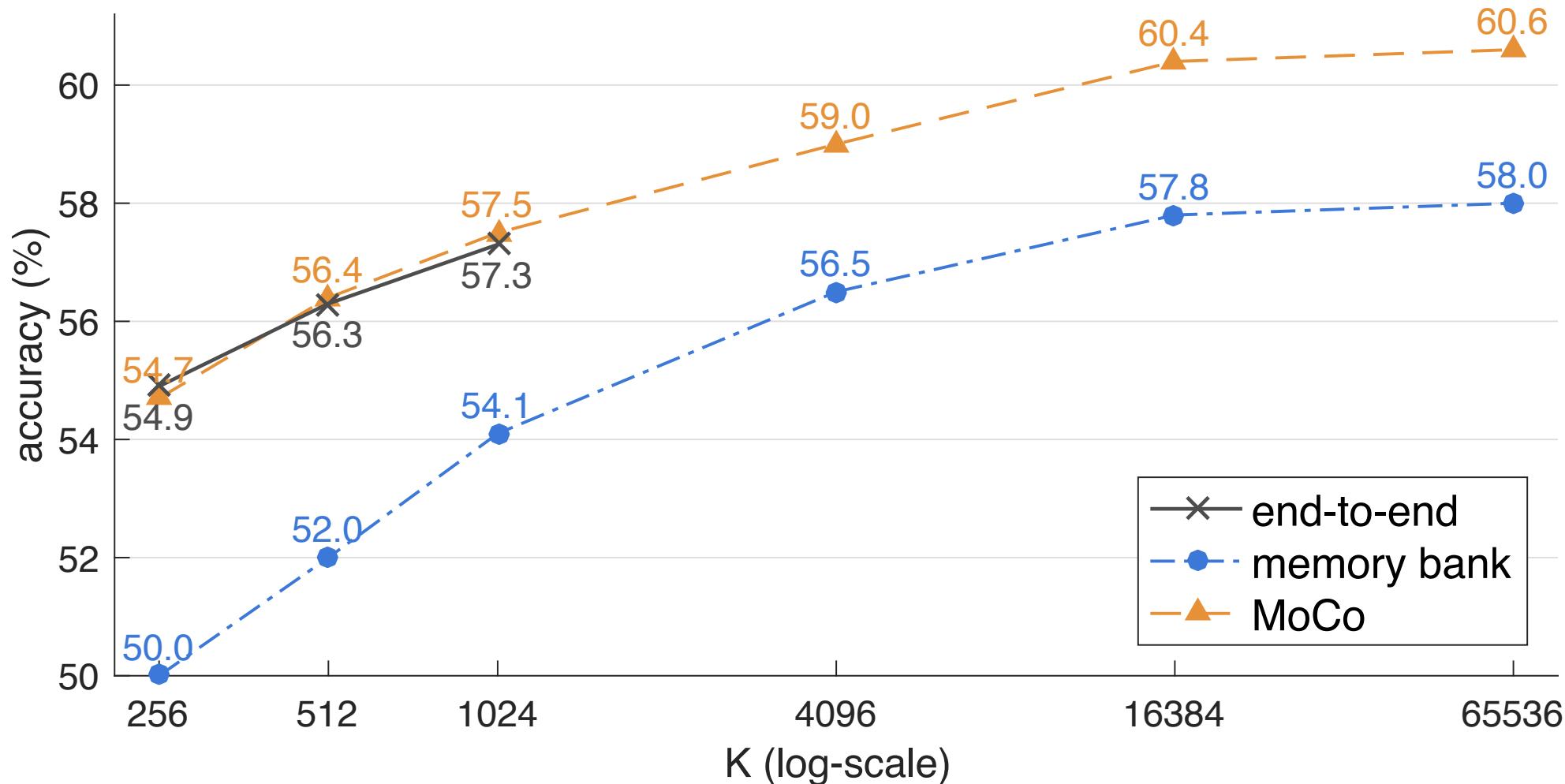
---

bmm: batch matrix multiplication; mm: matrix multiplication; cat: concatenation.

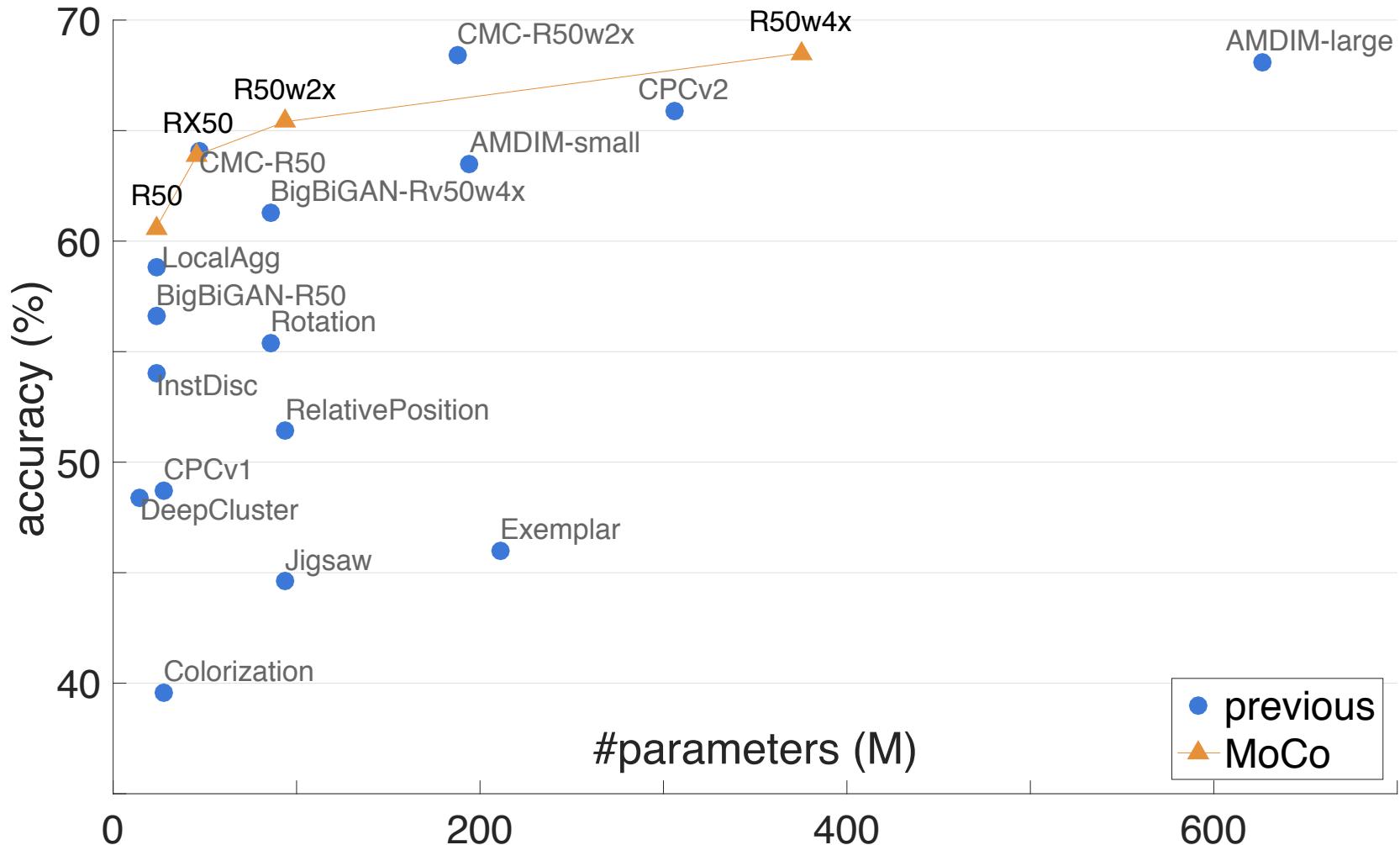
# Momentum Contrast (MoCo)



# Momentum Contrast (MoCo)



# Momentum Contrast (MoCo)



# SimCLR

---

## A Simple Framework for Contrastive Learning of Visual Representations

---

Ting Chen<sup>1</sup> Simon Kornblith<sup>1</sup> Mohammad Norouzi<sup>1</sup> Geoffrey Hinton<sup>1</sup>

### Abstract

This paper presents *SimCLR*: a simple framework for contrastive learning of visual representations. We simplify recently proposed contrastive self-supervised learning algorithms without requiring specialized architectures or a memory bank. In order to understand what enables the contrastive prediction tasks to learn useful representations, we systematically study the major components of our framework. We show that (1) composition of data augmentations plays a critical role in defining effective predictive tasks, (2) introducing a learnable nonlinear transformation between the representation and the contrastive loss substantially improves the quality of the learned representations, and (3) contrastive learning benefits from larger batch sizes and more training steps compared to supervised learning. By combining these findings, we are able to considerably outperform previous methods for self-supervised and semi-supervised learning on ImageNet. A linear classifier trained on self-supervised representations learned by SimCLR achieves 76.5% top-1 accuracy, which is a 7% relative improvement over previous state-of-the-art, matching the performance of a supervised ResNet-50. When fine-tuned on only 1% of the labels, we achieve 85.8% top-5 accuracy, outperforming AlexNet with 100× fewer labels.

# SimCLR

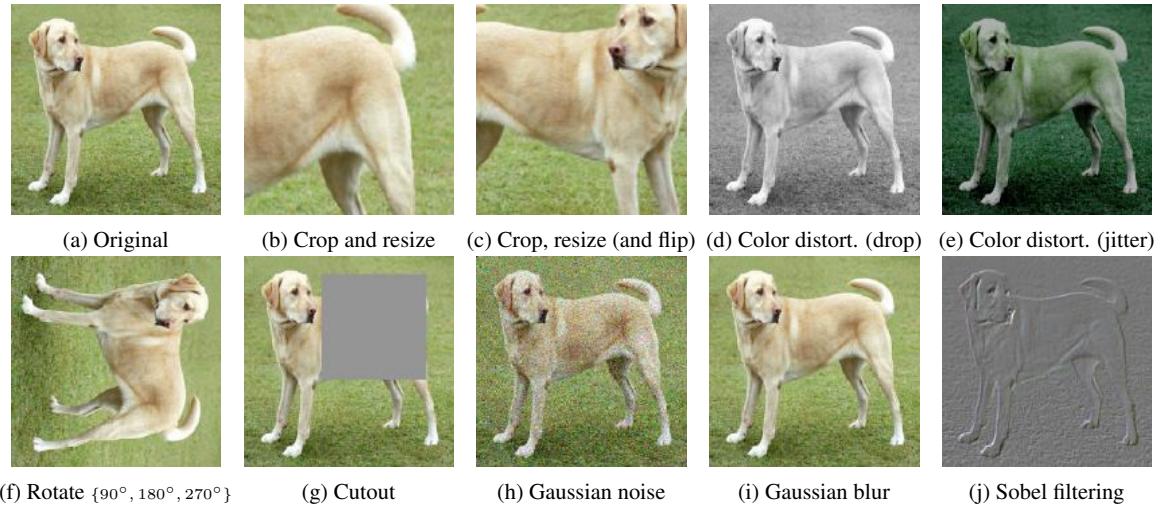
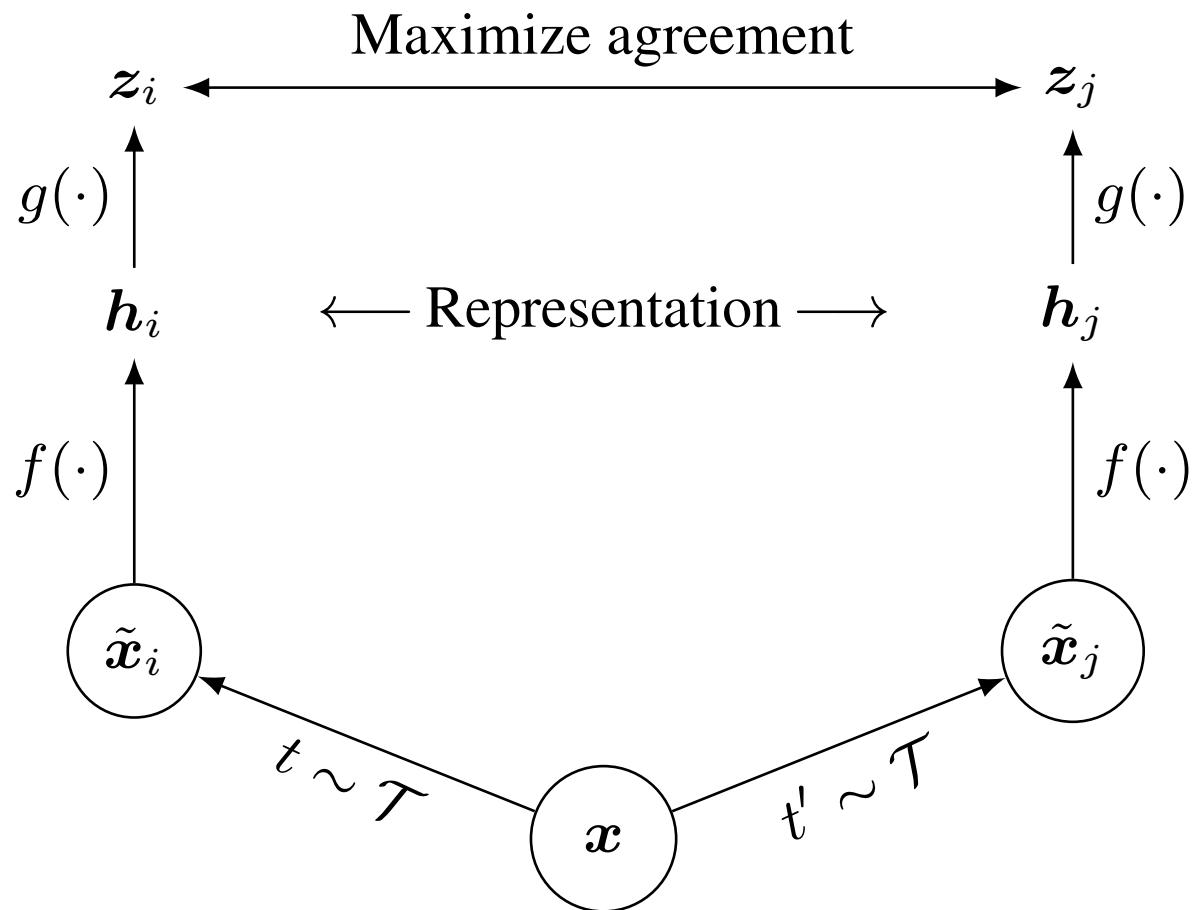


Figure 2. A simple framework for contrastive learning of visual representations. Two separate data augmentation operators are sampled from the same family of augmentations ( $t \sim \mathcal{T}$  and  $t' \sim \mathcal{T}$ ) and applied to each data example to obtain two correlated views. A base encoder network  $f(\cdot)$  and a projection head  $g(\cdot)$  are trained to maximize agreement using a contrastive loss. After training is completed, we throw away the projection head  $g(\cdot)$  and use encoder  $f(\cdot)$  and representation  $h$  for downstream tasks.

# SimCLR

---

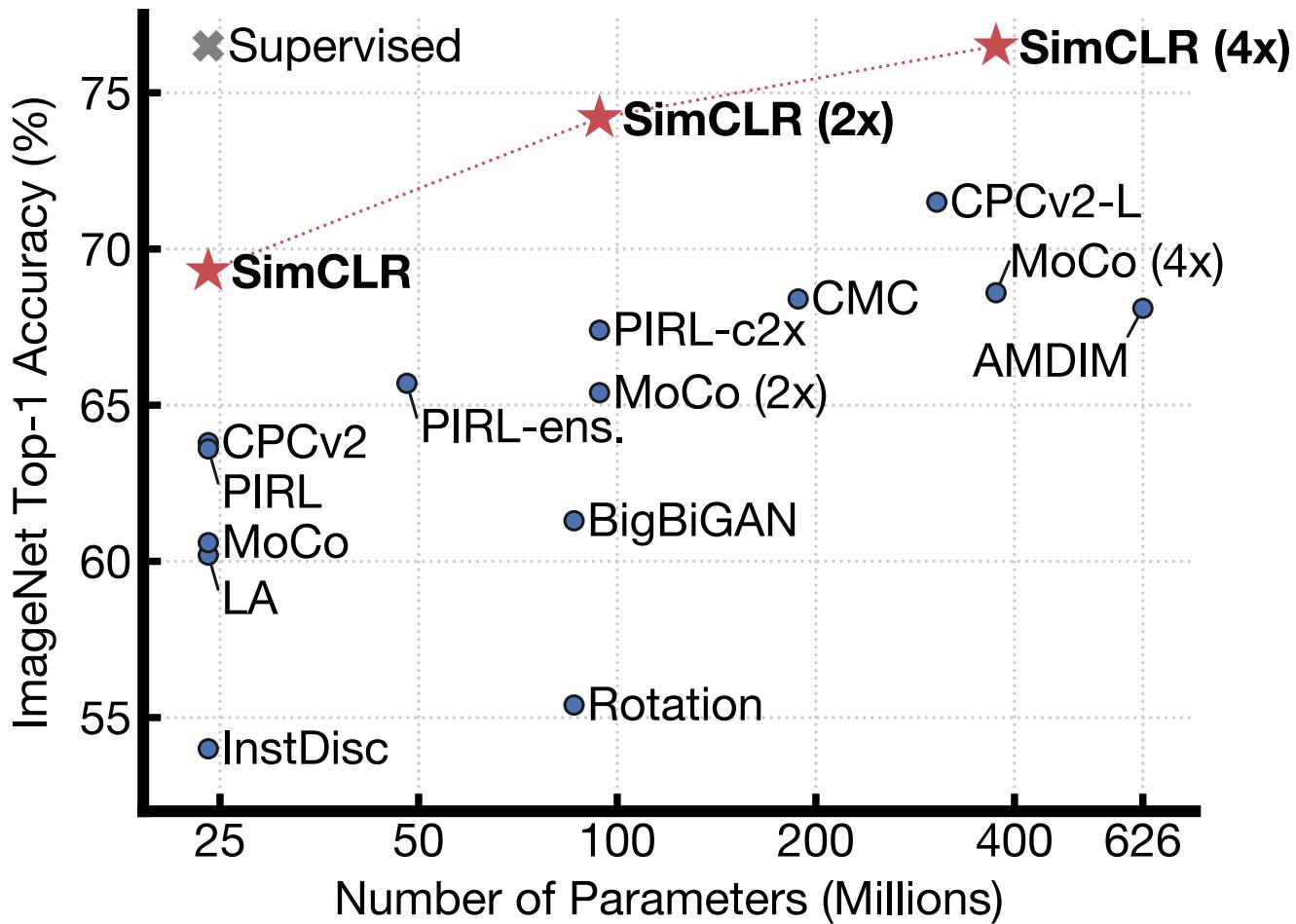
**Algorithm 1** SimCLR's main learning algorithm.

---

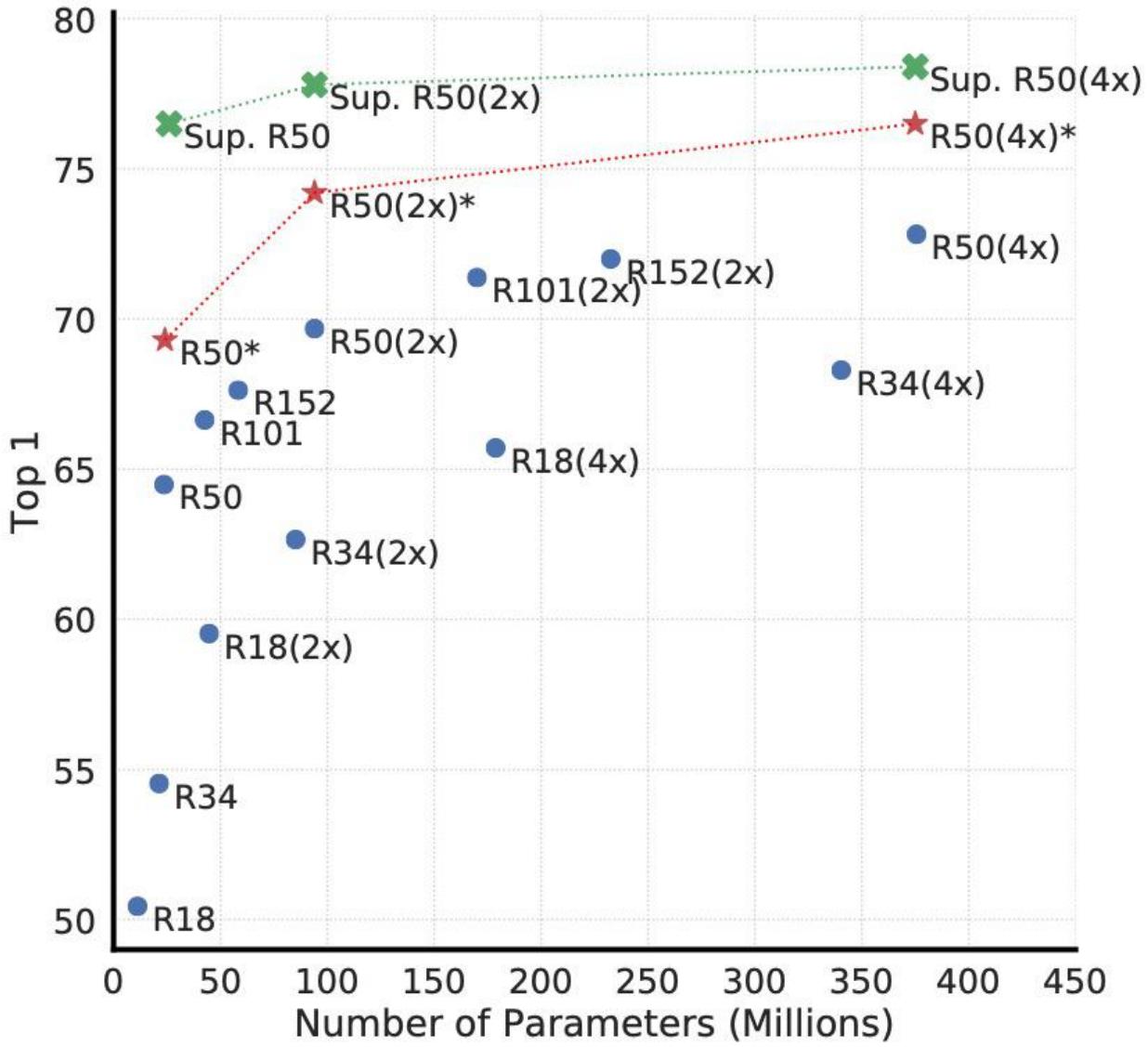
**input:** batch size  $N$ , constant  $\tau$ , structure of  $f, g, \mathcal{T}$ .  
**for** sampled minibatch  $\{\mathbf{x}_k\}_{k=1}^N$  **do**  
    **for all**  $k \in \{1, \dots, N\}$  **do**  
        draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$   
        # the first augmentation  
         $\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$   
         $\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$  # representation  
         $\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$  # projection  
        # the second augmentation  
         $\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$   
         $\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$  # representation  
         $\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$  # projection  
    **end for**  
    **for all**  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  **do**  
         $s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$  # pairwise similarity  
    **end for**  
    **define**  $\ell(i, j)$  **as**  $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$   
     $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$   
    update networks  $f$  and  $g$  to minimize  $\mathcal{L}$   
    **end for**  
    **return** encoder network  $f(\cdot)$ , and throw away  $g(\cdot)$

---

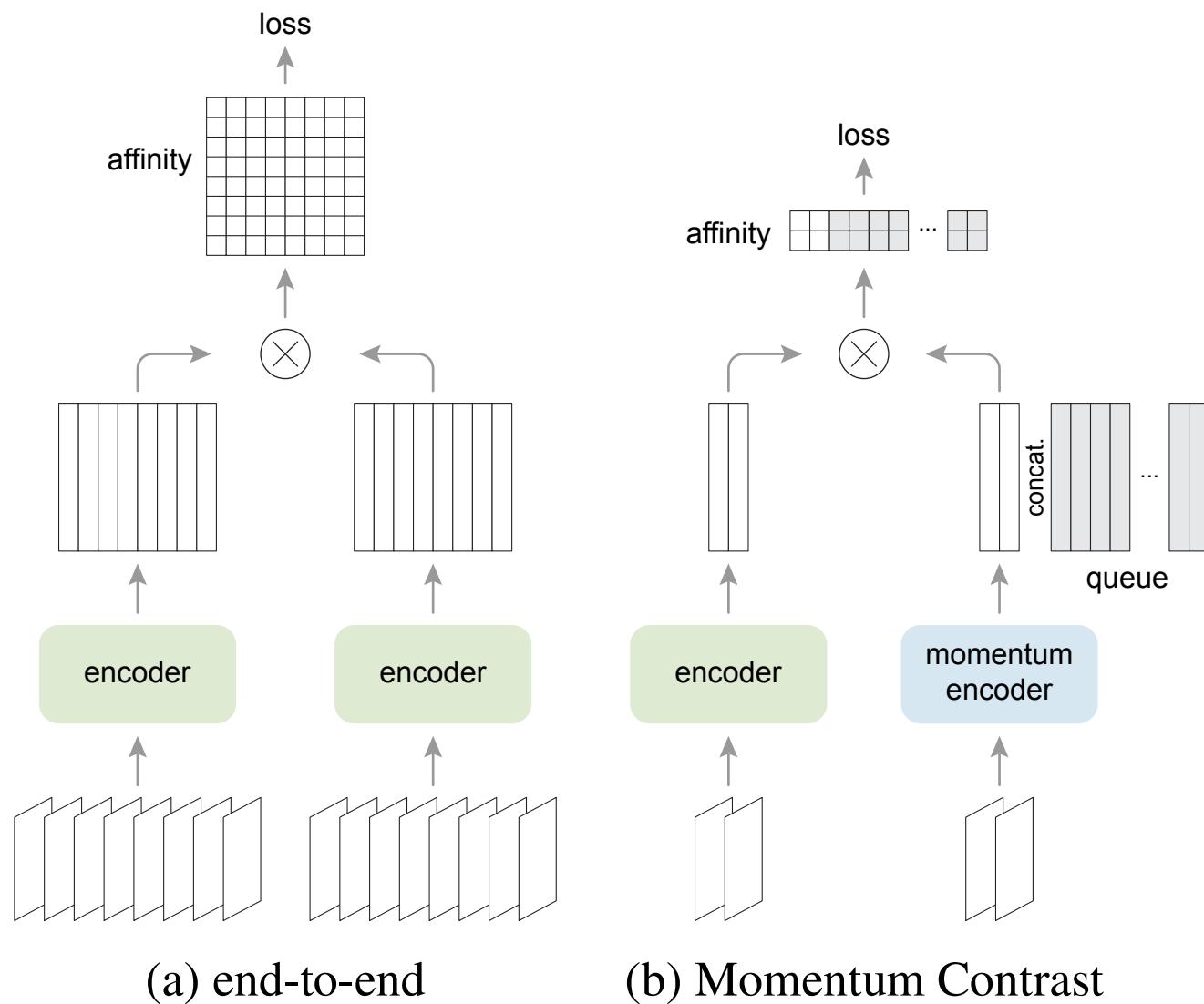
# SimCLR



# SimCLR



# MoCo v2 vs SimCLR



## Bootstrap Your Own Latent A New Approach to Self-Supervised Learning

---

Jean-Bastien Grill<sup>\*1</sup> Florian Strub<sup>\*1</sup> Florent Altché<sup>\*1</sup> Corentin Tallec<sup>\*1</sup> Pierre H. Richemond<sup>\*1,2</sup>  
Elena Buchatskaya<sup>1</sup> Carl Doersch<sup>1</sup> Bernardo Avila Pires<sup>1</sup> Zhaohan Daniel Guo<sup>1</sup>  
Mohammad Gheshlaghi Azar<sup>1</sup> Bilal Piot<sup>1</sup> Koray Kavukcuoglu<sup>1</sup> Rémi Munos<sup>1</sup> Michal Valko<sup>1</sup>

<sup>1</sup>DeepMind      <sup>2</sup>Imperial College

[jbgrill,fstrub,altche,corentint,richemond]@google.com

### Abstract

We introduce **Bootstrap Your Own Latent** (BYOL), a new approach to self-supervised image representation learning. BYOL relies on two neural networks, referred to as *online* and *target* networks, that interact and learn from each other. From an augmented view of an image, we train the online network to predict the target network representation of the same image under a different augmented view. At the same time, we update the target network with a slow-moving average of the online network. While state-of-the art methods intrinsically rely on negative pairs, BYOL achieves a new state of the art *without them*. BYOL reaches 74.3% top-1 classification accuracy on ImageNet using the standard linear evaluation protocol with a ResNet-50 architecture and 79.6% with a larger ResNet. We show that BYOL performs on par or better than the current state of the art on both transfer and semi-supervised benchmarks.

# BYOL

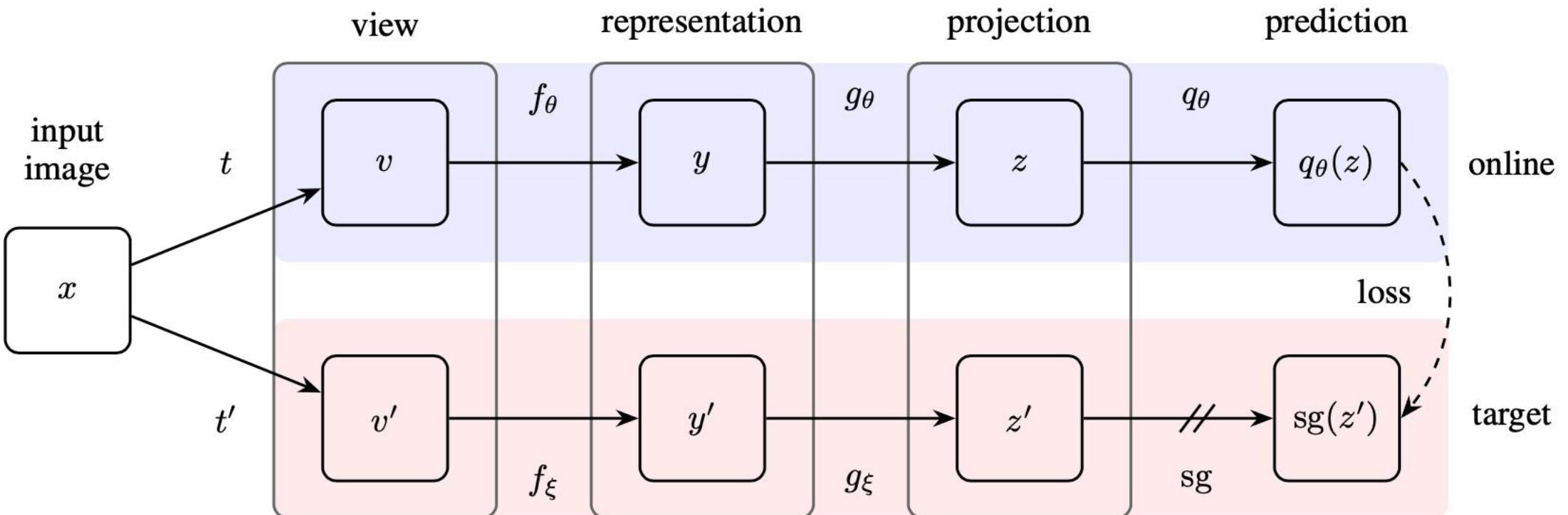
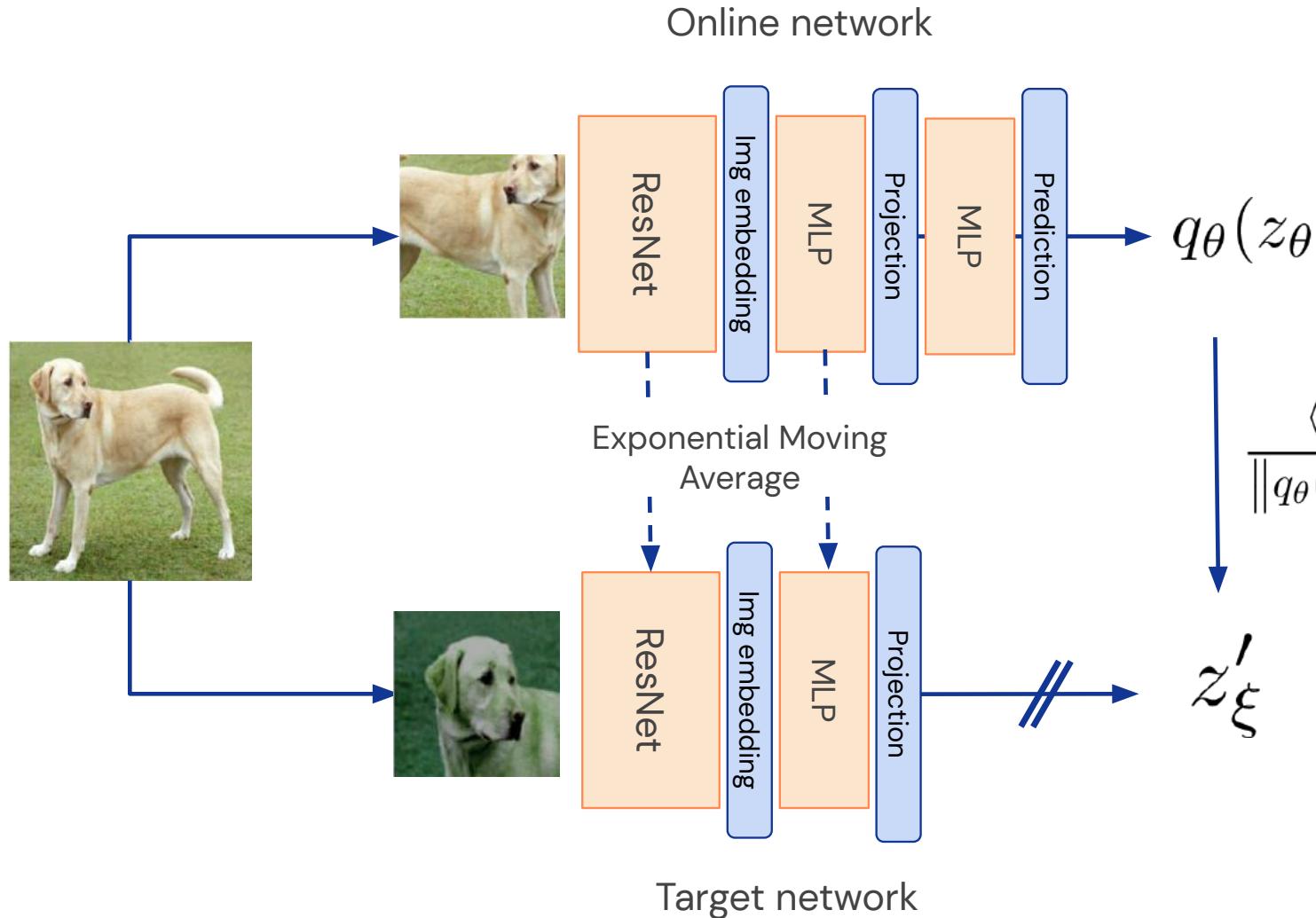


Figure 2: BYOL's architecture. BYOL minimizes a similarity loss between  $q_\theta(z)$  and  $sg(z')$ , where  $\theta$  are the trained weights,  $\xi$  are an exponential moving average of  $\theta$  and sg means stop-gradient. At the end of training, everything but  $f_\theta$  is discarded and  $y$  is used as the image representation.

# BYOL



Normalize features

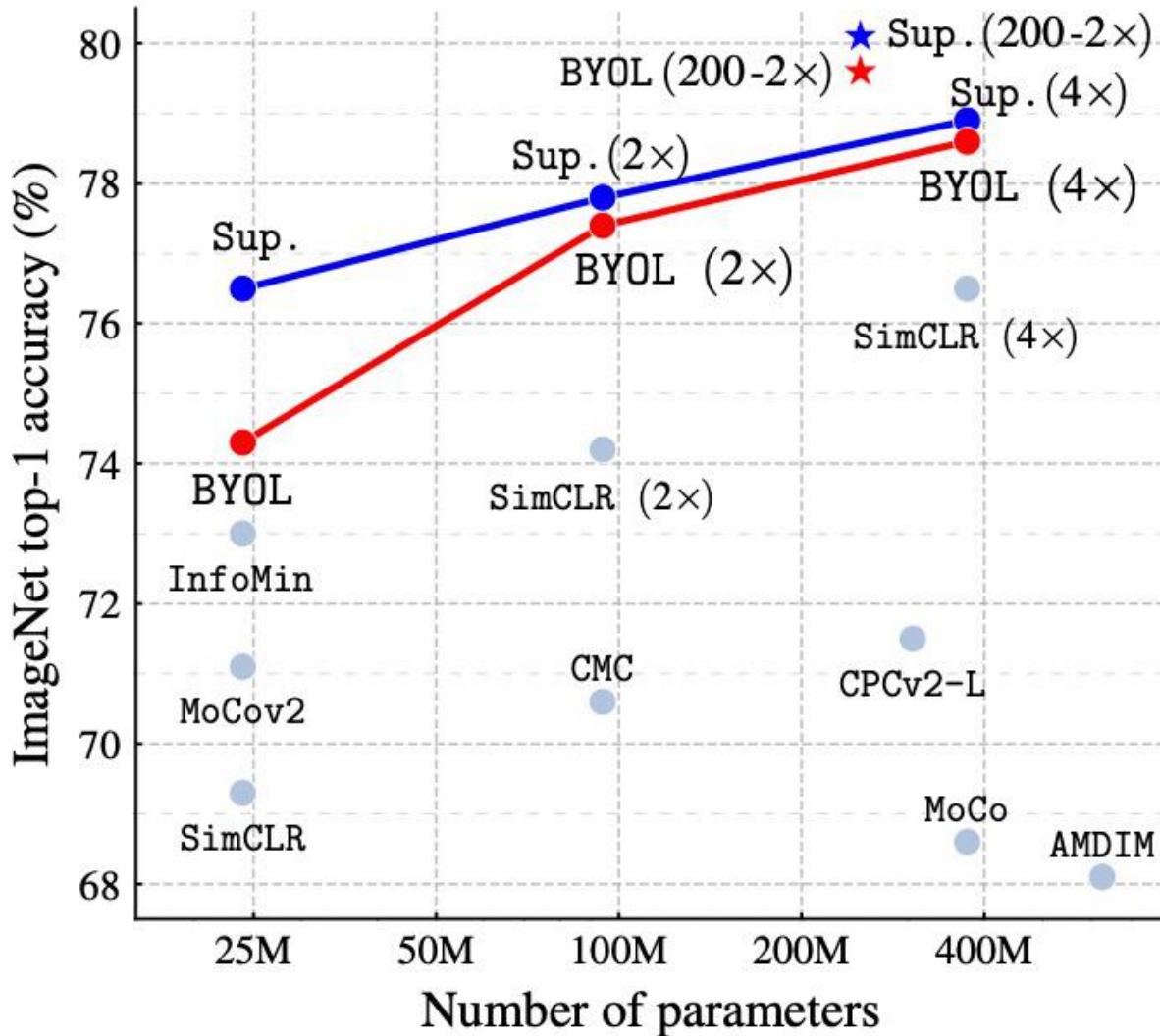
$$\bar{z}'_\xi \triangleq z'_\xi / \|z'_\xi\|_2$$

$$\bar{q}_\theta(z_\theta) \triangleq q_\theta(z_\theta) / \|q_\theta(z_\theta)\|_2$$

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

- Does not use negative examples!

# BYOL



# BYOL

Method	Top-1	Top-5
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.

Table 1: Top-1 and top-5 accuracies (in %) under linear evaluation on ImageNet.

# BYOL

Method	Food101	CIFAR10	CIFAR100	Birdsnap	SUN397	Cars	Aircraft	VOC2007	DTD	Pets	Caltech-101	Flowers
<i>Linear evaluation:</i>												
BYOL (ours)	<b>75.3</b>	91.3	<b>78.4</b>	<b>57.2</b>	<b>62.2</b>	<b>67.8</b>	60.6	82.5	75.5	90.4	94.2	<b>96.1</b>
SimCLR (repro)	72.8	90.5	74.4	42.4	60.6	49.3	49.8	81.4	<b>75.7</b>	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	<b>93.6</b>	78.3	53.7	61.9	66.7	<b>61.0</b>	<b>82.8</b>	74.9	<b>91.5</b>	<b>94.5</b>	94.7
<i>Fine-tuned:</i>												
BYOL (ours)	<b>88.5</b>	<b>97.8</b>	86.1	<b>76.3</b>	63.7	91.6	<b>88.1</b>	<b>85.4</b>	<b>76.2</b>	91.7	<b>93.8</b>	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	<b>86.4</b>	75.8	<b>64.3</b>	<b>92.1</b>	86.0	85.0	74.6	<b>92.1</b>	93.3	<b>97.6</b>
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

Table 3: Transfer learning results from ImageNet (IN) with the standard ResNet-50 architecture.

# DINO

## Emerging Properties in Self-Supervised Vision Transformers

Mathilde Caron<sup>1,2</sup> Hugo Touvron<sup>1,3</sup> Ishan Misra<sup>1</sup> Hervé Jegou<sup>1</sup>  
Julien Mairal<sup>2</sup> Piotr Bojanowski<sup>1</sup> Armand Joulin<sup>1</sup>

<sup>1</sup> Facebook AI Research

<sup>2</sup> Inria\*

<sup>3</sup> Sorbonne University



Figure 1: **Self-attention from a Vision Transformer with  $8 \times 8$  patches trained with no supervision.** We look at the self-attention of the [CLS] token on the heads of the last layer. This token is not attached to any label nor supervision. These maps show that the model automatically learns class-specific features leading to unsupervised object segmentations.

## Abstract

In this paper, we question if self-supervised learning provides new properties to Vision Transformer (ViT) [16] that stand out compared to convolutional networks (convnets). Beyond the fact that adapting self-supervised methods to this architecture works particularly well, we make the following observations: first, self-supervised ViT features contain explicit information about the semantic segmentation of an image, which does not emerge as clearly with supervised ViTs, nor with convnets. Second, these features are also excellent  $k$ -NN classifiers, reaching 78.3% top-1 on ImageNet with a small ViT. Our study also underlines the importance of momentum encoder [26], multi-crop training [9], and the use of small patches with ViTs. We implement our findings into a simple self-supervised method, called DINO, which we interpret as a form of self-distillation with **no labels**. We show the synergy between DINO and ViTs by achieving 80.1% top-1 on ImageNet in linear evaluation with ViT-Base.

# DINO

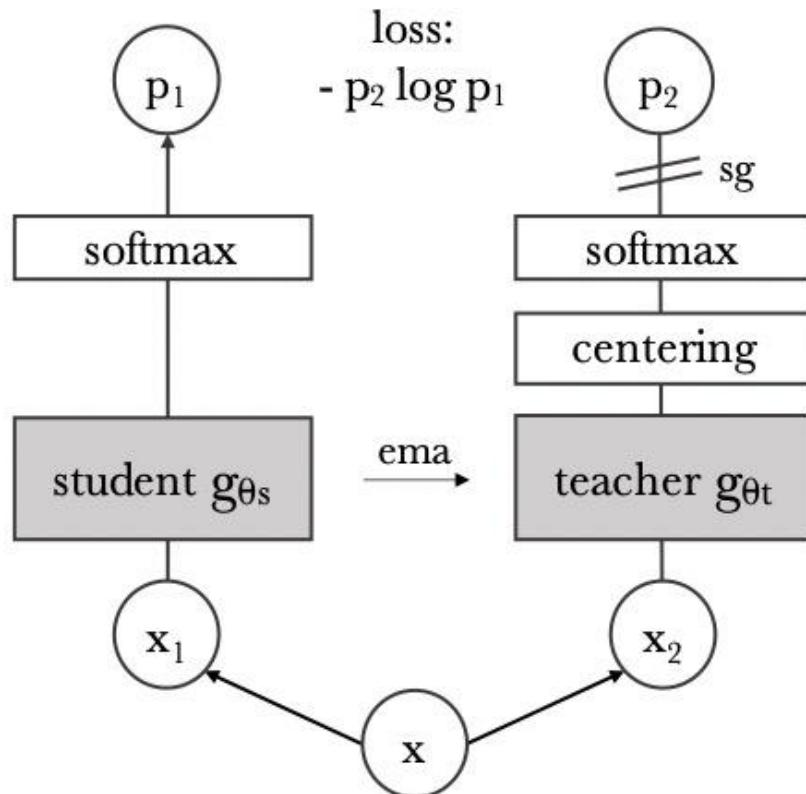


Figure 2: **Self-distillation with no labels.** We illustrate DINO in the case of one single pair of views ( $x_1, x_2$ ) for simplicity. The model passes two different random transformations of an input image to the student and teacher networks. Both networks have the same architecture but different parameters. The output of the teacher network is centered with a mean computed over the batch. Each networks outputs a  $K$  dimensional feature that is normalized with a temperature softmax over the feature dimension. Their similarity is then measured with a cross-entropy loss. We apply a stop-gradient (sg) operator on the teacher to propagate gradients only through the student. The teacher parameters are updated with an exponential moving average (ema) of the student parameters.

# DINO

---

## Algorithm 1 DINO PyTorch pseudocode w/o multi-crop.

---

```
# gs, gt: student and teacher networks
# C: center (K)
# tps, tpt: student and teacher temperatures
# l, m: network and center momentum rates
gt.params = gs.params
for x in loader: # load a minibatch x with n samples
    x1, x2 = augment(x), augment(x) # random views

    s1, s2 = gs(x1), gs(x2) # student output n-by-K
    t1, t2 = gt(x1), gt(x2) # teacher output n-by-K

    loss = H(t1, s2)/2 + H(t2, s1)/2
    loss.backward() # back-propagate

    # student, teacher and center updates
    update(gs) # SGD
    gt.params = l*gt.params + (1-l)*gs.params
    C = m*C + (1-m)*cat([t1, t2]).mean(dim=0)

def H(t, s):
    t = t.detach() # stop gradient
    s = softmax(s / tps, dim=1)
    t = softmax((t - C) / tpt, dim=1) # center + sharpen
    return - (t * log(s)).sum(dim=1).mean()
```

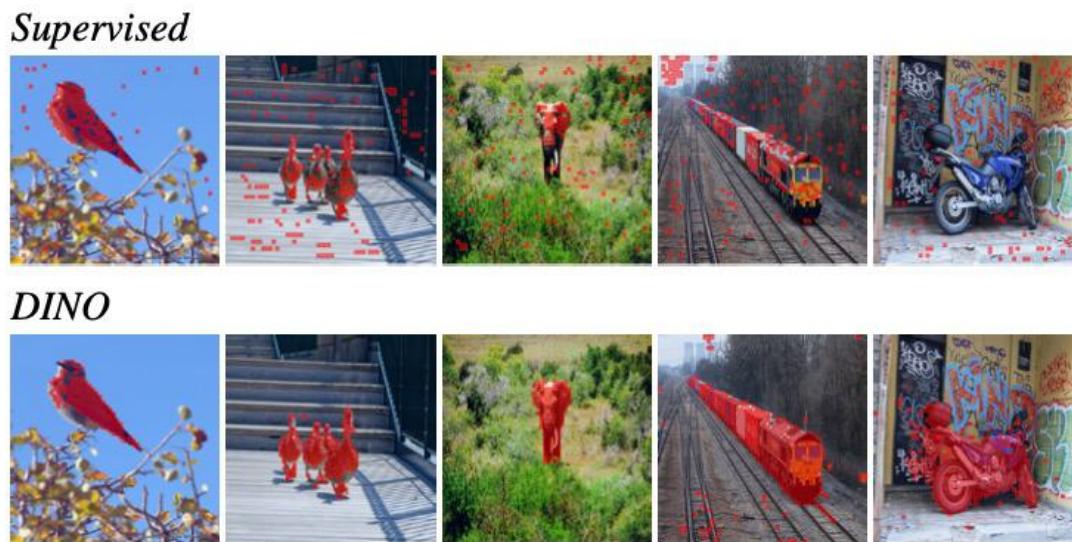
---

# DINO

Table 2: **Linear and  $k$ -NN classification on ImageNet.** We report top-1 accuracy for linear and  $k$ -NN evaluations on the validation set of ImageNet for different self-supervised methods. We focus on ResNet-50 and ViT-small architectures, but also report the best results obtained across architectures. \* are run by us. We run the  $k$ -NN evaluation for models with official released weights. The throughput (im/s) is calculated on a NVIDIA V100 GPU with 128 samples per forward. Parameters (M) are of the feature extractor.

Method	Arch.	Param.	im/s	Linear	$k$ -NN	
Supervised	RN50	23	1237	79.3	79.3	
SCLR [11]	RN50	23	1237	69.1	60.7	
MoCov2 [13]	RN50	23	1237	71.1	61.9	
InfoMin [54]	RN50	23	1237	73.0	65.3	
BarlowT [66]	RN50	23	1237	73.2	66.0	
OBoW [21]	RN50	23	1237	73.8	61.9	
BYOL [23]	RN50	23	1237	74.4	64.8	
DCv2 [9]	RN50	23	1237	75.2	67.1	
SwAV [9]	RN50	23	1237	<b>75.3</b>	65.7	
<b>DINO</b>	RN50	23	1237	<b>75.3</b>	<b>67.5</b>	
Supervised	ViT-S	21	1007	79.8	79.8	
BYOL* [23]	ViT-S	21	1007	71.4	66.6	
MoCov2* [13]	ViT-S	21	1007	72.7	64.4	
SwAV* [9]	ViT-S	21	1007	73.5	66.3	
<b>DINO</b>	ViT-S	21	1007	<b>77.0</b>	<b>74.5</b>	
<i>Comparison across architectures</i>						
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	–
<b>DINO</b>	ViT-B/16		85	312	78.2	<b>76.1</b>
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
<b>DINO</b>	ViT-S/8		21	180	79.7	<b>78.3</b>
SCLRV2 [12]	RN152w3+SK		794	46	79.8	73.1
<b>DINO</b>	ViT-B/8		85	63	<b>80.1</b>	77.4

# DINO



	Random	Supervised	DINO
ViT-S/16	22.0	27.3	45.9
ViT-S/8	21.8	23.7	44.7

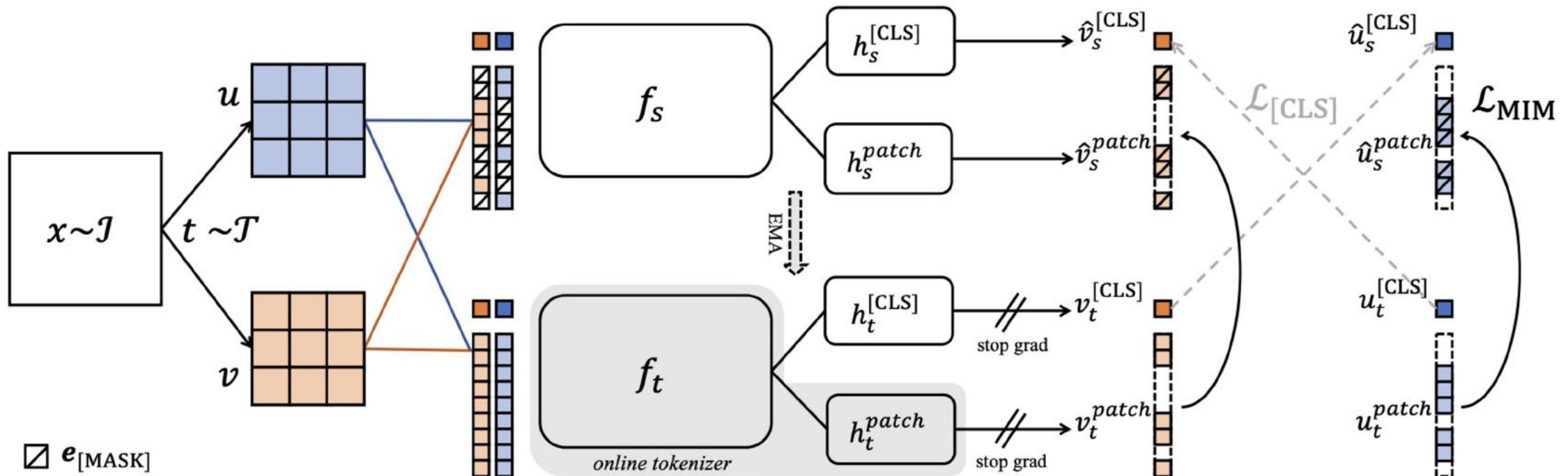
**Figure 4: Segmentations from supervised versus DINO.** We visualize masks obtained by thresholding the self-attention maps to keep 60% of the mass. On top, we show the resulting masks for a ViT-S/8 trained with supervision and DINO. We show the best head for both models. The table at the bottom compares the Jaccard similarity between the ground truth and these masks on the validation images of PASCAL VOC12 dataset.

# iBOT

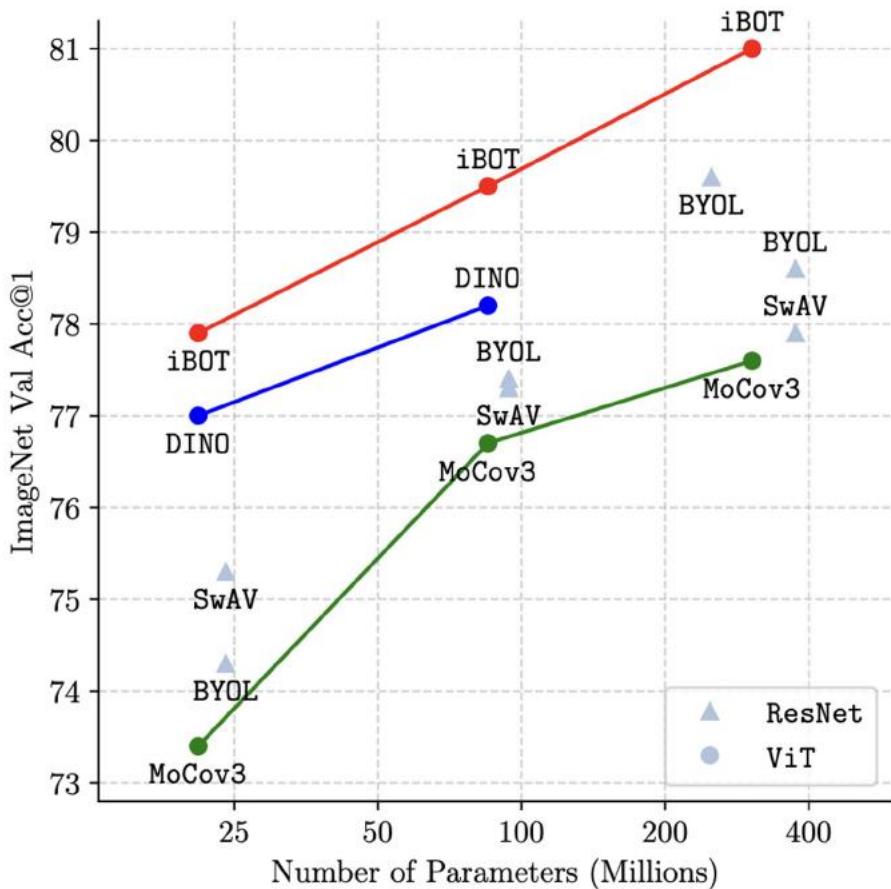
## IBOT 🤖: IMAGE BERT PRE-TRAINING WITH ONLINE TOKENIZER

Jinghao Zhou<sup>1</sup> Chen Wei<sup>2</sup> Huiyu Wang<sup>2</sup> Wei Shen<sup>3</sup> Cihang Xie<sup>4</sup> Alan Yuille<sup>2</sup> Tao Kong<sup>1</sup>

<sup>1</sup>ByteDance <sup>2</sup>Johns Hopkins University <sup>3</sup>Shanghai Jiao Tong University <sup>4</sup>UC Santa Cruz



# iBOT



**Table 6: Object detection (Det.) & instance segmentation (ISeg.) on COCO and Semantic segmentation (Seg.) on ADE20K.** We report the results of ViT-S/16 (left) and ViT-B/16 (right). Seg.<sup>†</sup> denotes using a linear head for semantic segmentation.

Method	Arch.	Param.	Det.			ISeg.		Seg.		Method	Det.			ISeg.		Seg.	
			AP <sup>b</sup>	AP <sup>m</sup>	mIoU	AP <sup>b</sup>	AP <sup>m</sup>	mIoU	AP <sup>b</sup>		AP <sup>b</sup>	AP <sup>m</sup>	mIoU	AP <sup>b</sup>	AP <sup>m</sup>	mIoU	
Sup.	Swin-T	29	48.1	41.7	44.5					Sup.	49.8	43.2	35.4	46.6			
MoBY	Swin-T	29	48.1	41.5	44.1					BEiT	50.1	43.5	27.4	45.8			
Sup.	ViT-S/16	21	46.2	40.1	44.5					DINO	50.1	43.4	34.5	46.8			
iBOT	ViT-S/16	21	<b>49.4</b>	<b>42.6</b>	<b>45.4</b>					iBOT	<b>51.2</b>	<b>44.2</b>	<b>38.3</b>	<b>50.0</b>			

# DINO v2

## DINOv2: Learning Robust Visual Features without Supervision

Maxime Oquab\*\*, Timothée Darcet\*\*, Théo Moutakanni\*\*,  
Huy V. Vo\*, Marc Szafraniec\*, Vasil Khalidov\*, Pierre Fernandez, Daniel Haziza,  
Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba,  
Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat,  
Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal<sup>1</sup>,  
Patrick Labatut\*, Armand Joulin\*, Piotr Bojanowski\*

*Meta AI Research*

<sup>1</sup>*Inria*

# DINO v2

ViT-L

	INet-1k k-NN	INet-1k linear
iBOT	72.9	82.3
+ (our reproduction)	74.5 $\uparrow$ 1.6	83.2 $\uparrow$ 0.9
+ LayerScale, Stochastic Depth	75.4 $\uparrow$ 0.9	82.0 $\downarrow$ 1.2
+ 128k prototypes	76.6 $\uparrow$ 1.2	81.9 $\downarrow$ 0.1
+ KoLeo	78.9 $\uparrow$ 2.3	82.5 $\uparrow$ 0.6
+ SwiGLU FFN	78.7 $\downarrow$ 0.2	83.1 $\uparrow$ 0.6
+ Patch size 14	78.9 $\uparrow$ 0.2	83.5 $\uparrow$ 0.4
+ Teacher momentum 0.994	79.4 $\uparrow$ 0.5	83.6 $\uparrow$ 0.1
+ Tweak warmup schedules	80.5 $\uparrow$ 1.1	83.8 $\uparrow$ 0.2
+ Batch size 3k	81.7 $\uparrow$ 1.2	84.7 $\uparrow$ 0.9
+ Sinkhorn-Knopp	81.7 =	84.7 =
+ Untying heads = DINOV2	82.0 $\uparrow$ 0.3	84.5 $\downarrow$ 0.2

# DINO v2

Method	Arch.	Data	Text sup.	kNN	linear		
				val	val	ReaL	V2
<b>Weakly supervised</b>							
CLIP	ViT-L/14	WIT-400M	✓	79.8	84.3	88.1	75.3
CLIP	ViT-L/14 <sub>336</sub>	WIT-400M	✓	80.5	85.3	88.8	75.8
SWAG	ViT-H/14	IG3.6B	✓	82.6	85.7	88.7	77.6
OpenCLIP	ViT-H/14	LAION-2B	✓	81.7	84.4	88.4	75.5
OpenCLIP	ViT-G/14	LAION-2B	✓	83.2	86.2	89.4	77.2
EVA-CLIP	ViT-g/14	custom*	✓	<b>83.5</b>	86.4	89.3	77.4
<b>Self-supervised</b>							
MAE	ViT-H/14	INet-1k	✗	49.4	76.6	83.3	64.8
DINO	ViT-S/8	INet-1k	✗	78.6	79.2	85.5	68.2
SEERv2	RG10B	IG2B	✗	—	79.8	—	—
MSN	ViT-L/7	INet-1k	✗	79.2	80.7	86.0	69.7
EsViT	Swin-B/W=14	INet-1k	✗	79.4	81.3	87.0	70.4
Mugs	ViT-L/16	INet-1k	✗	80.2	82.1	86.9	70.8
iBOT	ViT-L/16	INet-22k	✗	72.9	82.3	87.5	72.4
DINOv2	ViT-S/14	LVD-142M	✗	79.0	81.1	86.6	70.9
	ViT-B/14	LVD-142M	✗	82.1	84.5	88.3	75.1
	ViT-L/14	LVD-142M	✗	<b>83.5</b>	86.3	89.5	78.0
	ViT-g/14	LVD-142M	✗	<b>83.5</b>	<b>86.5</b>	<b>89.6</b>	<b>78.4</b>

# DINO v2

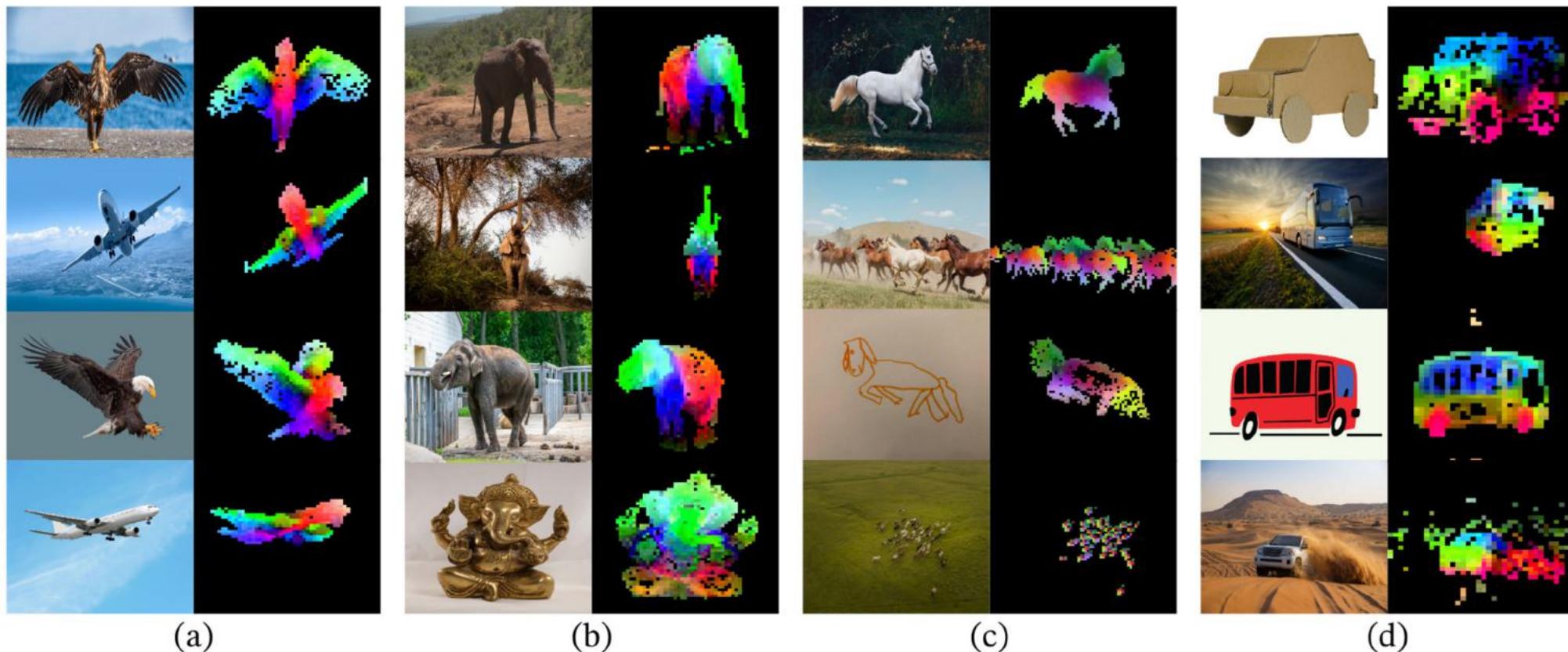
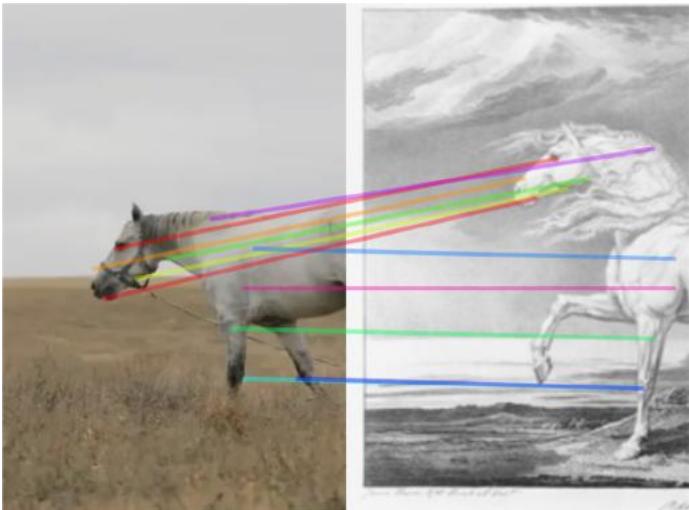


Figure 1: **Visualization of the first PCA components.** We compute a PCA between the patches of the images from the same column (a, b, c and d) and show their first 3 components. Each component is matched to a different color channel. Same parts are matched between related images despite changes of pose, style or even objects. Background is removed by thresholding the first PCA component.

# DINO v2



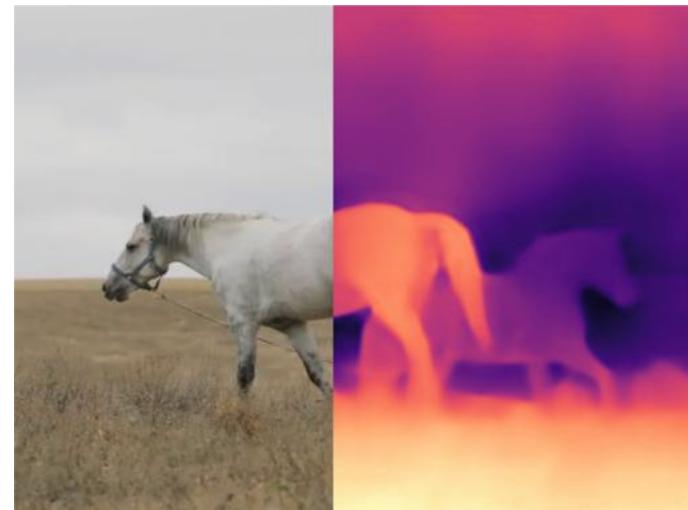
Feature  
matching



Segmentation



Image  
retrieval



Depth  
prediction

# CLIP Contrastive Language-Image Pre-training

---

## Learning Transferable Visual Models From Natural Language Supervision

---

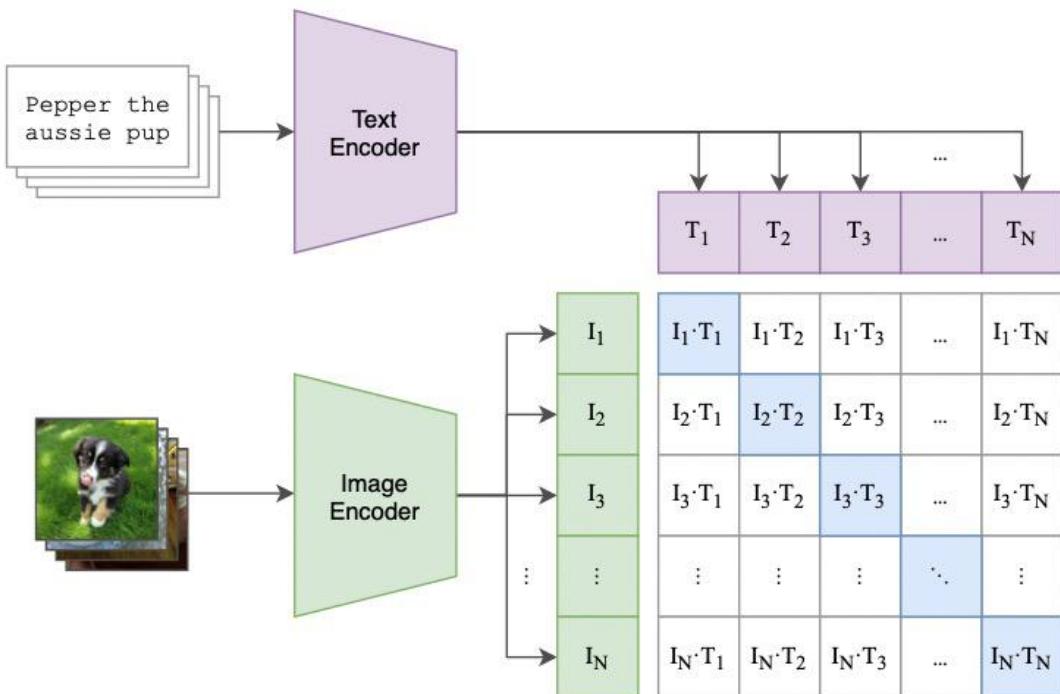
**Alec Radford**<sup>\* 1</sup> **Jong Wook Kim**<sup>\* 1</sup> **Chris Hallacy**<sup>1</sup> **Aditya Ramesh**<sup>1</sup> **Gabriel Goh**<sup>1</sup> **Sandhini Agarwal**<sup>1</sup>  
**Girish Sastry**<sup>1</sup> **Amanda Askell**<sup>1</sup> **Pamela Mishkin**<sup>1</sup> **Jack Clark**<sup>1</sup> **Gretchen Krueger**<sup>1</sup> **Ilya Sutskever**<sup>1</sup>

### Abstract

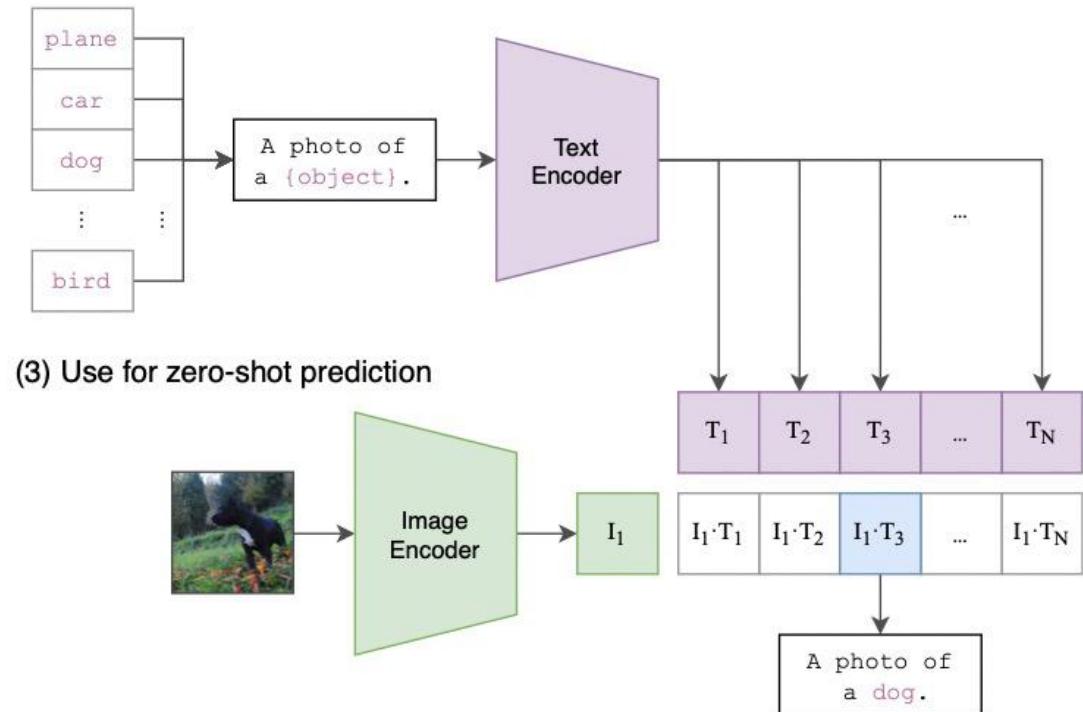
State-of-the-art computer vision systems are trained to predict a fixed set of predetermined object categories. This restricted form of supervision limits their generality and usability since additional labeled data is needed to specify any other visual concept. Learning directly from raw text about images is a promising alternative which leverages a much broader source of supervision. We demonstrate that the simple pre-training task of predicting which caption goes with which image is an efficient and scalable way to learn SOTA image representations from scratch on a dataset of 400 million (image, text) pairs collected from the internet. After pre-training, natural language is used to reference learned visual concepts (or describe new ones) enabling zero-shot transfer of the model to downstream tasks. We study the performance of this approach by benchmarking on over 30 different existing computer vision datasets, spanning tasks such as OCR, action recognition in videos, geo-localization, and many types of fine-grained object classification. The model transfers non-trivially to most tasks and is often competitive with a fully supervised baseline without the need for any dataset specific training. For instance, we match the accuracy of the original ResNet-50 on ImageNet zero-shot without needing to use any of the 1.28 million training examples it was trained on.

# CLIP

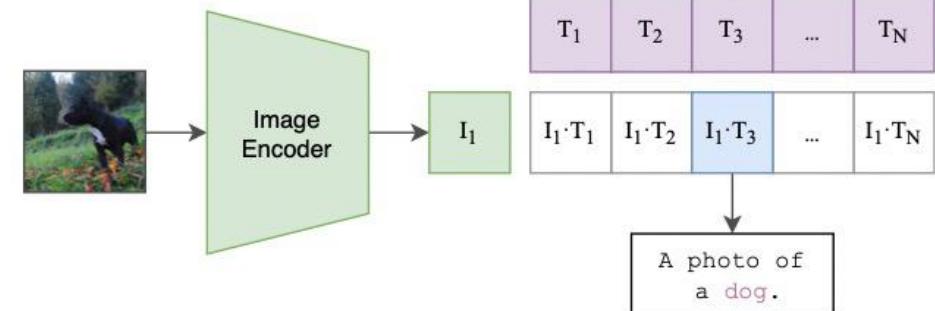
(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction



*Figure 1.* Summary of our approach. While standard image models jointly train an image feature extractor and a linear classifier to predict some label, CLIP jointly trains an image encoder and a text encoder to predict the correct pairings of a batch of (image, text) training examples. At test time the learned text encoder synthesizes a zero-shot linear classifier by embedding the names or descriptions of the target dataset's classes.

# CLIP

```
# image_encoder - ResNet or Vision Transformer
# text_encoder - CBOW or Text Transformer
# I[n, h, w, c] - minibatch of aligned images
# T[n, l]        - minibatch of aligned texts
# W_i[d_i, d_e] - learned proj of image to embed
# W_t[d_t, d_e] - learned proj of text to embed
# t              - learned temperature parameter

# extract feature representations of each modality
I_f = image_encoder(I) #[n, d_i]
T_f = text_encoder(T) #[n, d_t]

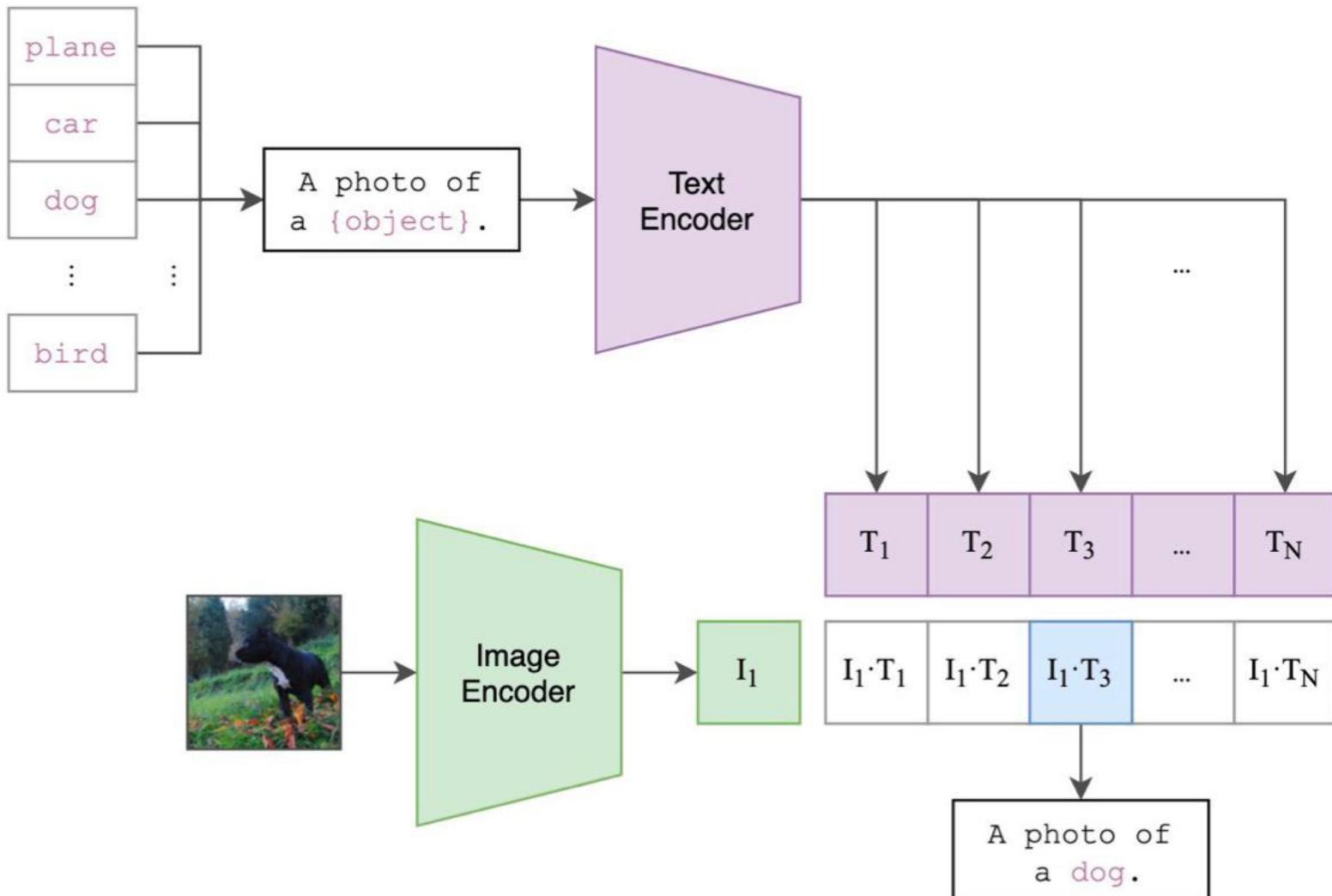
# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)

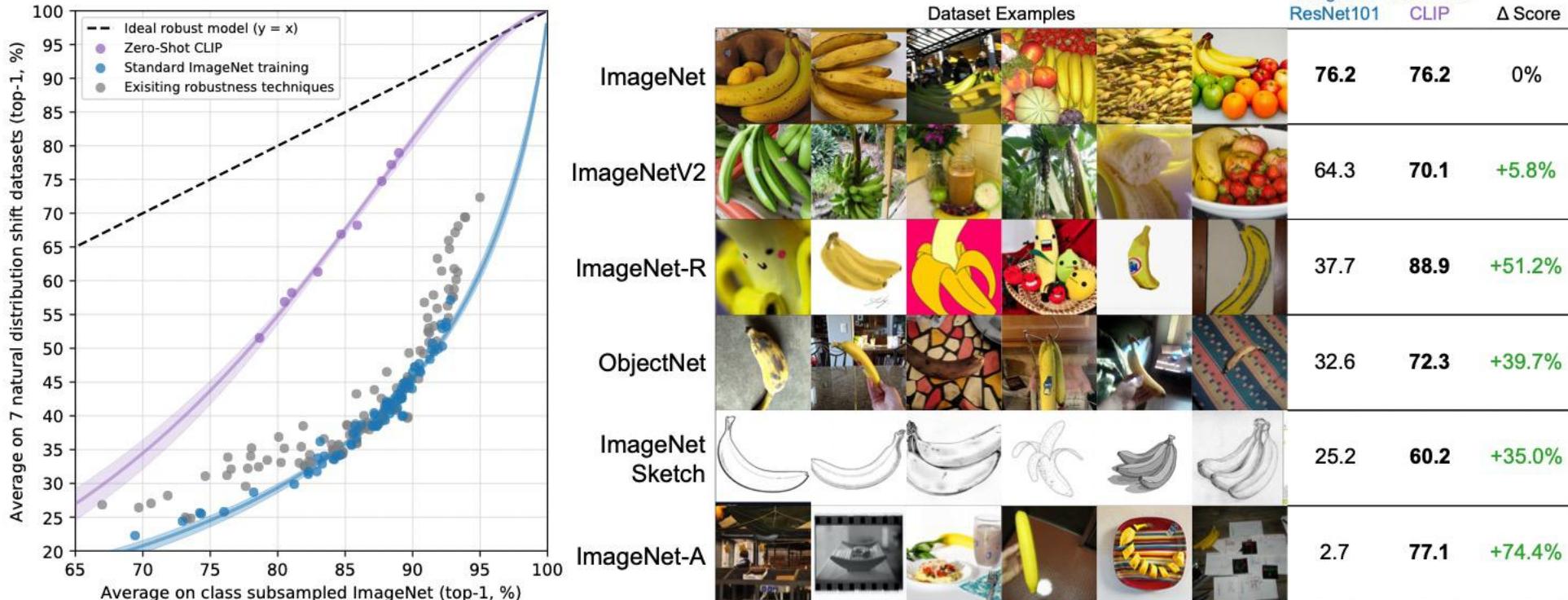
# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss   = (loss_i + loss_t)/2
```

Figure 3. Numpy-like pseudocode for the core of an implementation of CLIP.

# CLIP



# CLIP



**Figure 13. Zero-shot CLIP is much more robust to distribution shift than standard ImageNet models.** (Left) An ideal robust model (dashed line) performs equally well on the ImageNet distribution and on other natural image distributions. Zero-shot CLIP models shrink this “robustness gap” by up to 75%. Linear fits on logit transformed values are shown with bootstrap estimated 95% confidence intervals. (Right) Visualizing distribution shift for bananas, a class shared across 5 of the 7 natural distribution shift datasets. The performance of the best zero-shot CLIP model, ViT-L/14@336px, is compared with a model that has the same performance on the ImageNet validation set, ResNet-101.

# CLIP

