

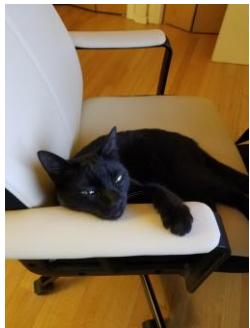


# Cours 5 : Image representations Self-Supervised Learning

Support de cours adapté de  
« CS231n: Deep Learning for Computer Vision, **Stanford - Spring 2022** »,  
*<http://cs231n.stanford.edu/>*

# Representations

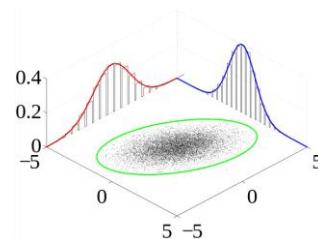
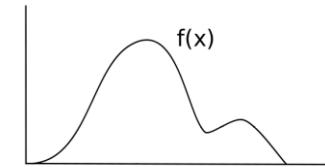
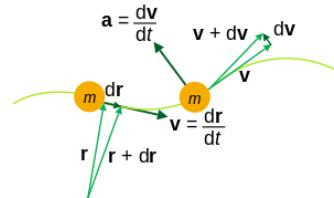
Things...



My heart beats as if the world is dropping,  
you may not feel the love but i do its a heart  
breaking moment of your life. enjoy the times  
that we have, it might not sound good but  
one thing it rhymes it might not be romantic  
but i think it is great,the best rhyme i've ever  
heard.



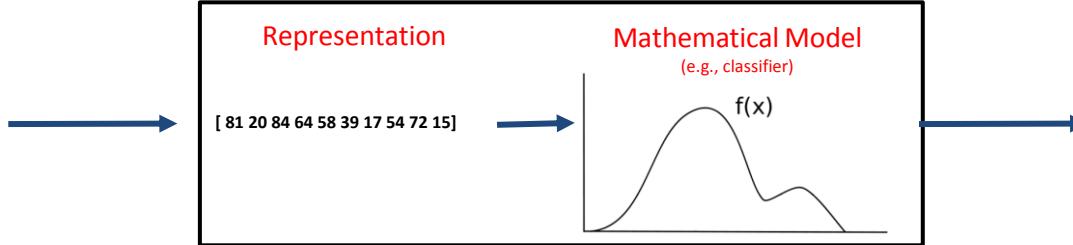
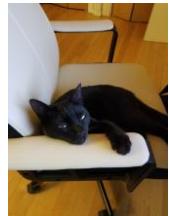
Our Knowledge...



A right-angled triangle with vertices at the corners. The vertical leg is labeled  $a$ , the horizontal leg is labeled  $b$ , and the hypotenuse is labeled  $c$ . The angle at vertex  $C$  is labeled  $\alpha$ . The angle at vertex  $B$  is labeled  $\beta$ . The angle at vertex  $A$  is labeled  $\gamma$ . The triangle is inscribed in a circle with center  $H$ .

$$a^2 + b^2 = c^2, \quad c = \sqrt{a^2 + b^2}$$
$$c^2 - a^2 = b^2, \quad c^2 - b^2 = a^2$$
$$\frac{a}{c} = \frac{HB}{a} \text{ and } \frac{b}{c} = \frac{AH}{b}$$
$$\sin \alpha = \frac{a}{c}, \quad \cos \alpha = \frac{b}{c}, \quad \tan \alpha = \frac{a}{b}$$
$$\cot \alpha = \frac{b}{a}, \quad \sec \alpha = \frac{c}{b}, \quad \csc \alpha = \frac{c}{a}$$

# Representations



I dare not speak of what I have done. Such twisted thoughts overtook my mind and now I am sorry to say that I have done the deed. I have murdered the King of Scotland, King Duncan.

After naming me the *Worthy Thane of Cawdor*, this is how I repay him. I have betrayed him in the most unimaginable way a person possibly could, and I've been forced to do it because I have to Banquo, whom I have loved as a friend. I do feel that I have never done such a treacherous thing, as I am afraid that *I shall sleep no more*.

I was waiting anxiously for my Lady to sound the bell that called me to do the deed. But before she did, a symbol of the supernatural appeared before my eyes. *The dagger of the mind* captured me and the handle was in my hand yet I couldn't grasp it, but *I could see them still*. I knew not whether to follow or to discard it from my eyes, but the *false creation* remained.

As I stepped closer to Duncan's room, I thought that I would panic and freeze, but when I got nearer, a sickening thought made me feel like I was doing the right thing! As soon as I heard the bell I knew that it was the bell that called me.

I heard him pleading as the dagger pierced through his skin.

I dare not speak of what I have done. Such twisted thoughts overtook my mind and now I am sorry to say that I have done the deed. I have murdered the King of Scotland, King Duncan.

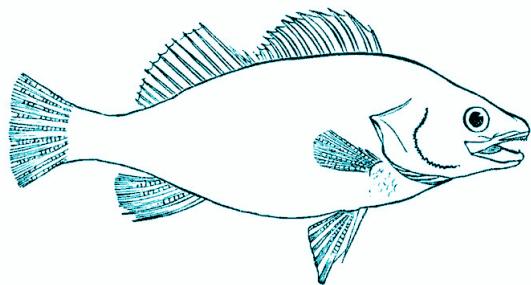
After naming me the *Worthy Thane of Cawdor*, this is how I repay him. I have betrayed him in the most unimaginable way a person possibly could, and I've been disloyal to him, just like I have to Banquo, whom I have loved as a friend. I do feel that I have never done such a treacherous thing, as I am afraid that *I shall sleep no more*.

I was waiting anxiously for my Lady to sound the bell that called me to do the deed. But before she did, a symbol of the supernatural appeared before my eyes. *The dagger of the mind* captured me and the handle was in my hand yet I couldn't grasp it, but *I could see them still*. I knew not whether to follow or to discard it from my eyes, but the *false creation* remained.

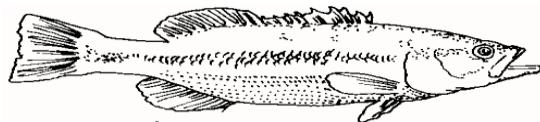
As I stepped closer to Duncan's room, I thought that I would panic and freeze, but when I got nearer, a sickening thought made me feel like I was doing the right thing! As soon as I heard the bell I knew that it was the bell that called me.

I heard him pleading as the dagger pierced through his skin.

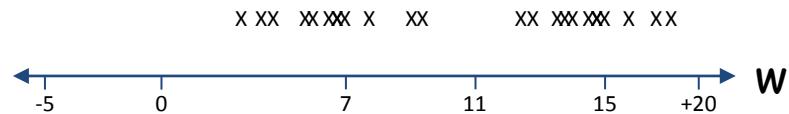
# Representations



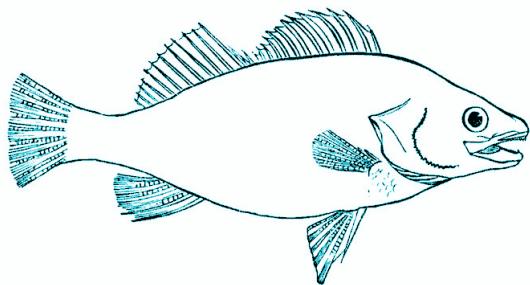
~12 lbs



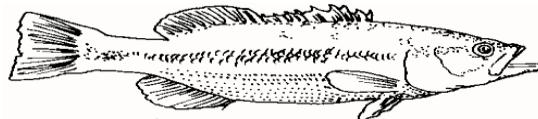
~8 lbs



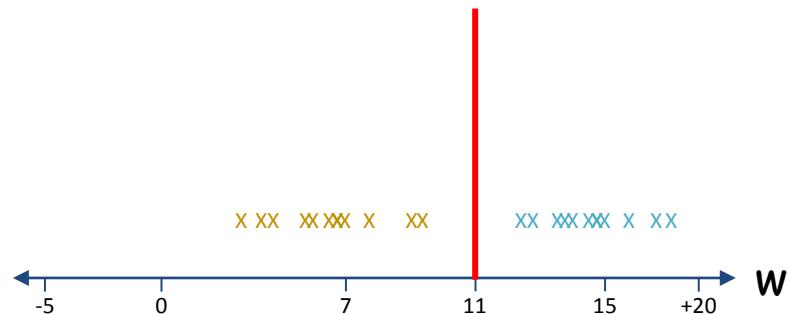
# Representations

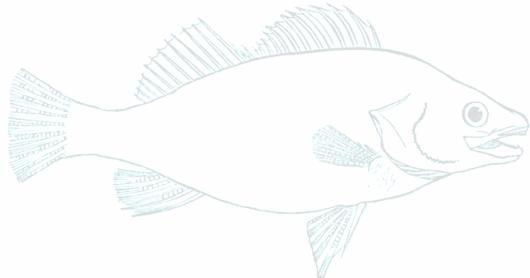


~12 lbs



~8 lbs

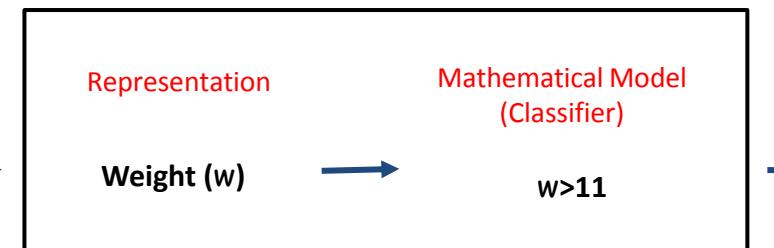
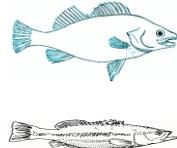
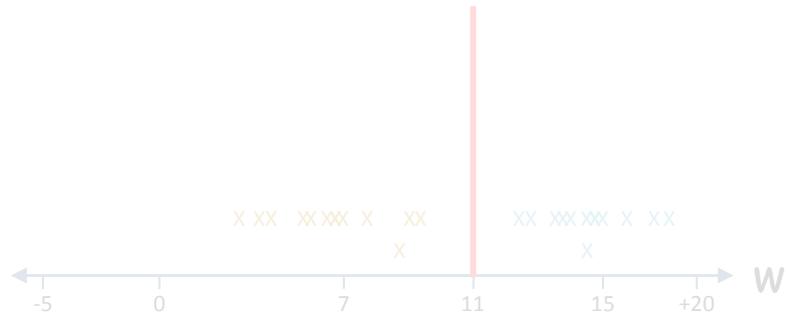




~12 lbs

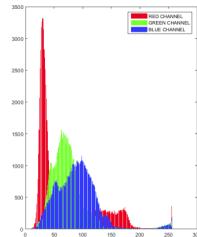


~8 lbs



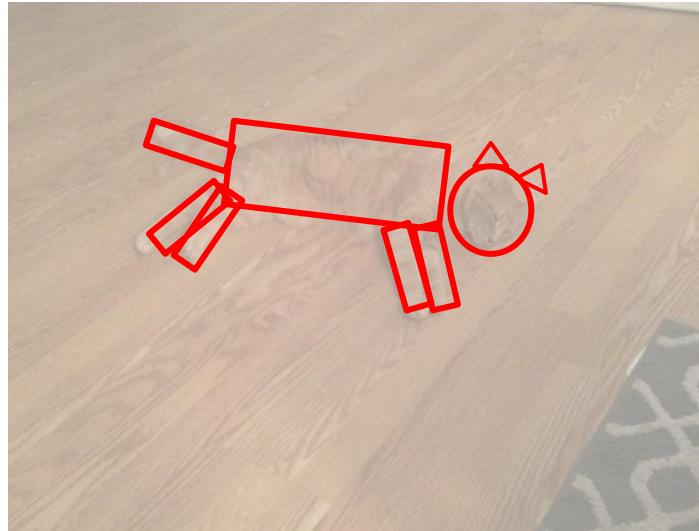
# Handcrafted image representation

Represent these cats for a cat detector!



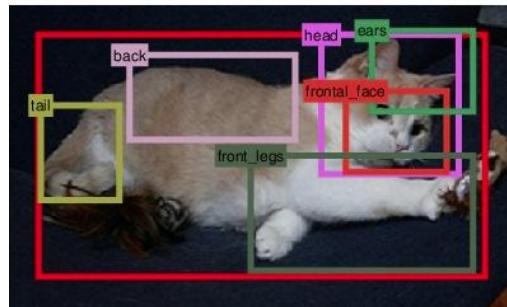
# Handcrafted image representation

Represent these cats for a cat detector!



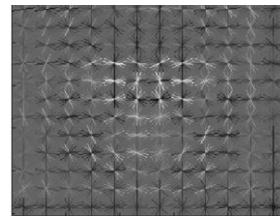
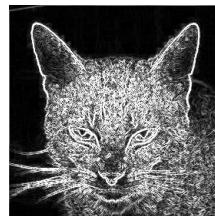
# Handcrafted image representation

Represent these cats for a cat detector!



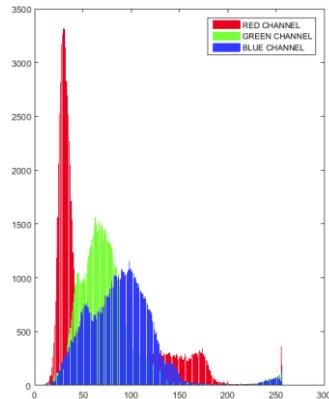
# Handcrafted image representation

Represent these cats for a cat detector!

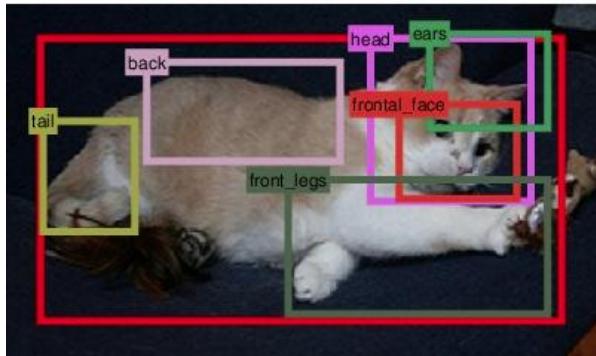


# Handcrafted image representation

Color  
Histograms



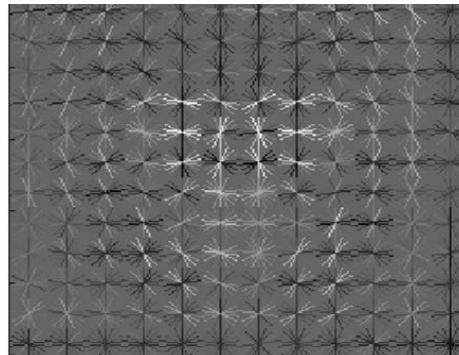
Deformable  
Part based  
Models (DPM)



Models based  
Shapes



Histogram of  
Gradients  
(HOG)



# Handcrafted image representation

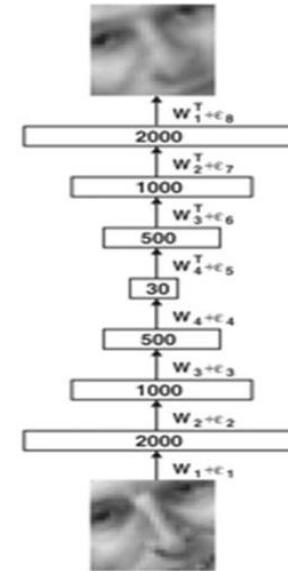
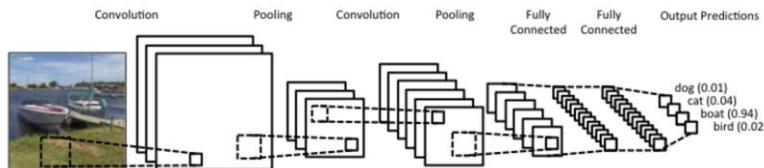
- Was the only way for a long time.
- (almost) Worked for many great applications:
  - Image Retrieval, Structure-from-motion, Face detection, Identification, etc
- Why alternatives?
  - Can't quite find the discriminative signature for a problem.
  - Discriminative signature can be found, but hard to approach programmatically.
  - Too many contributing factors to the problem.
    - Fusion non-trivial. Rule-based fusion outruled.
    - Fusion of contributing factors itself a comparably complex representation problem.

Why a sad  
image?

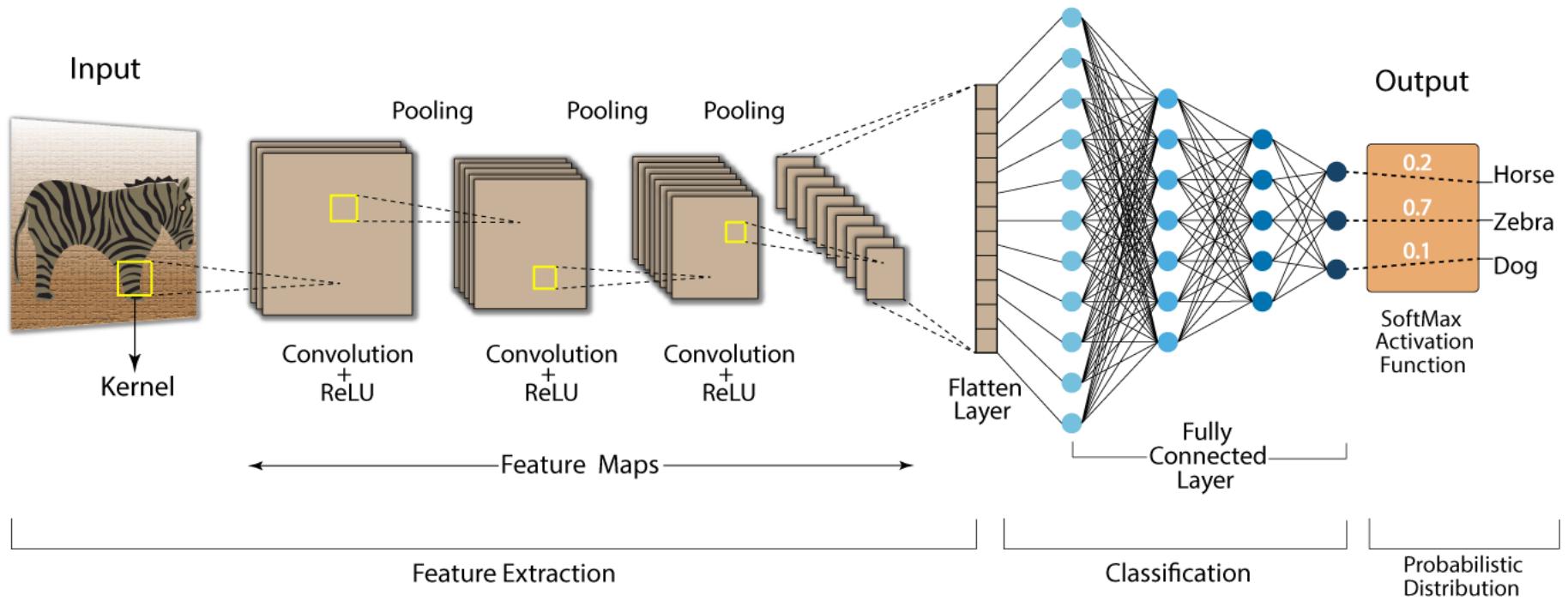


# Representation learning: 2 main approaches

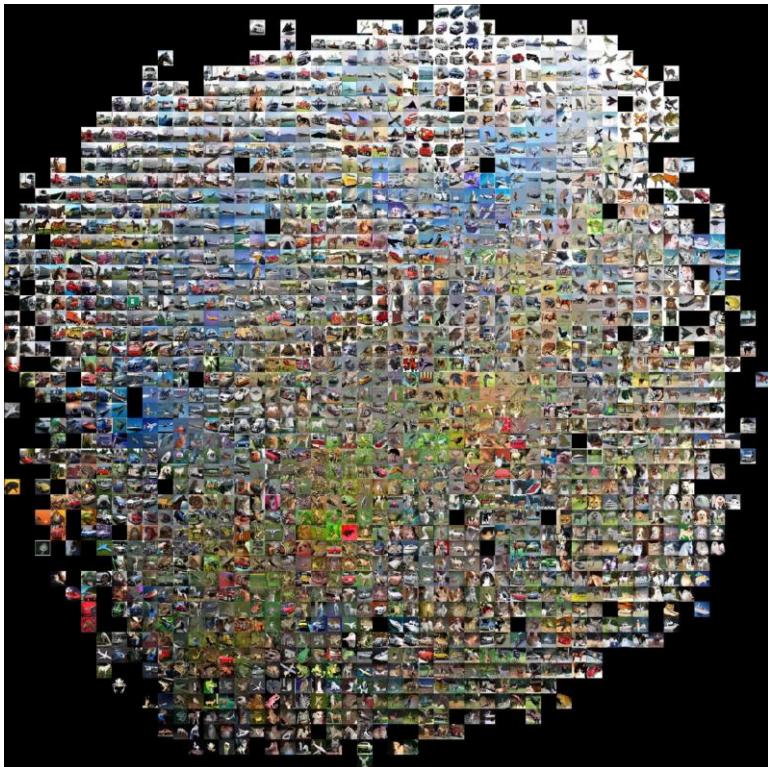
- Supervised
  - Representation constrained on task(s).
- Unsupervised
  - Representation constrained on reconstruction.



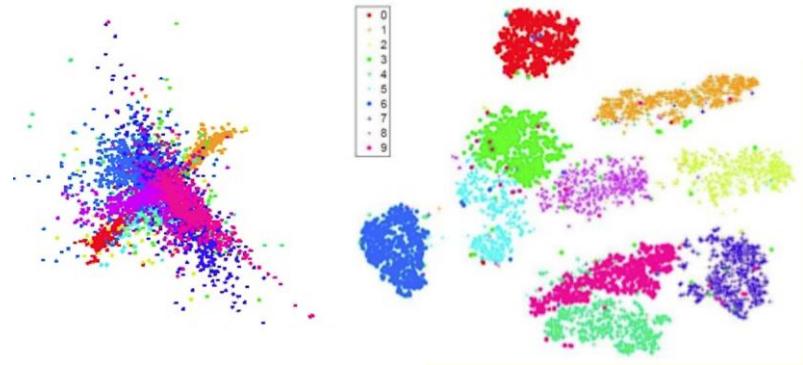
# Supervised approaches



# Visualization



Embedding



# Self-supervised Learning

- Aim to learn from data without manual label annotation.
- Self-supervised learning methods solve “pretext” tasks that produce **good features** for downstream tasks.
  - Learn with supervised learning objectives, e.g., classification, regression.
  - Labels of these pretext tasks are generated *automatically*

# Self-supervised pretext tasks

Example: learn to predict image transformations / complete corrupted images

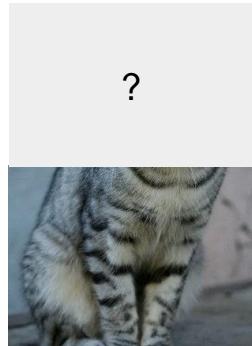
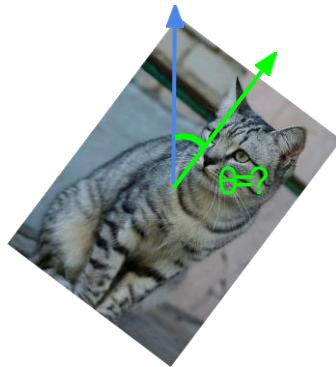


image completion



rotation prediction



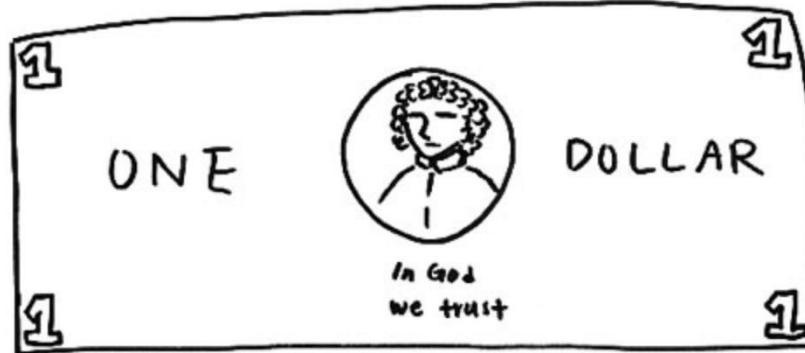
“jigsaw puzzle”



colorization

1. Solving the pretext tasks allow the model to learn good features.
2. We can automatically generate labels for the pretext tasks.

# Generative vs. Self-supervised Learning



Left: Drawing of a dollar bill from memory. Right: Drawing subsequently made with a dollar bill present. Image source: [Epstein, 2016](#)

Learning to generate pixel-level details is often unnecessary; learn high-level semantic features with pretext tasks instead

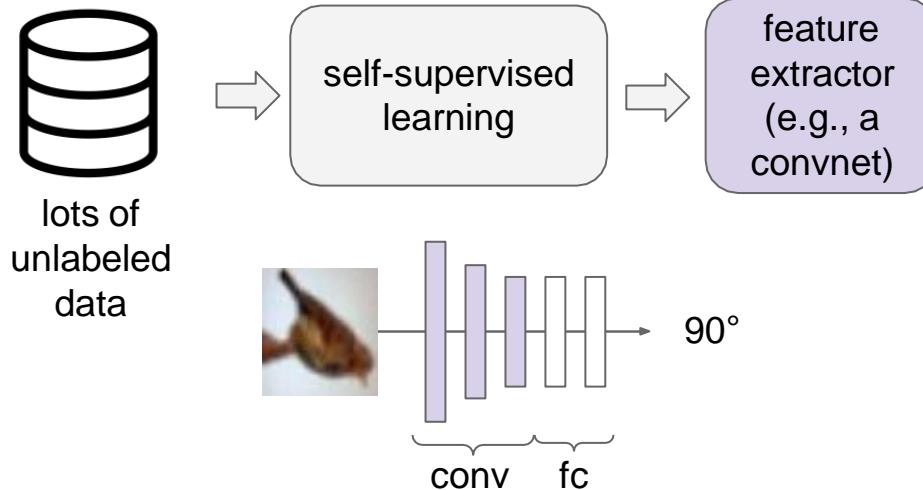
[Source: Anand, 2020](#)

# How to evaluate a self-supervised learning method?

We usually don't care about the performance of the self-supervised learning task, e.g., we don't care if the model learns to predict image rotation perfectly.

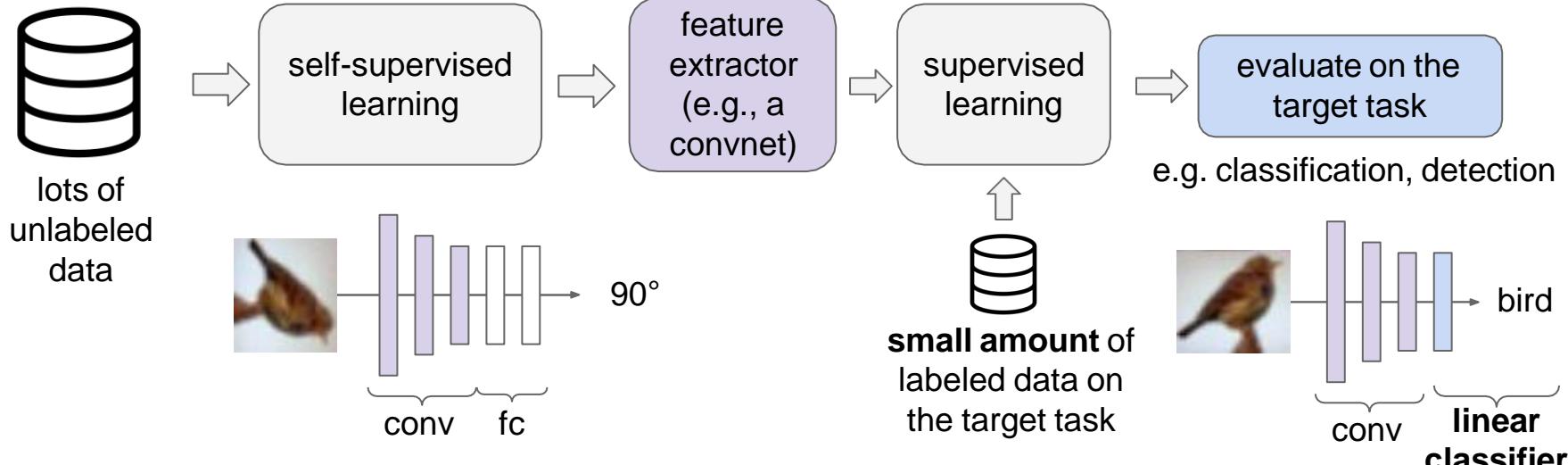
Evaluate the learned feature encoders on downstream *target tasks*

# How to evaluate a self-supervised learning method?



1. Learn good feature extractors from self-supervised pretext tasks, e.g., predicting image rotations

# How to evaluate a self-supervised learning method?



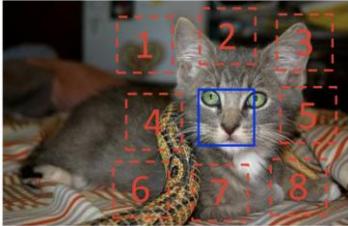
1. Learn good feature extractors from self-supervised pretext tasks, e.g., predicting image rotations

2. Attach a shallow network on the feature extractor; train the shallow network on the target task with small amount of labeled data

# Broader picture

## Today's lecture

### computer vision



Doersch et al., 2015

### robot / reinforcement learning



Dense Object Net (Florence and Manuelli et al., 2018)

### language modeling

#### Language Models are Few-Shot Learners

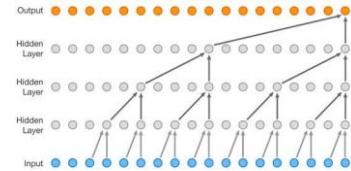
Tom B. Brown\* Benjamin Mann\* Nick Ryder\* Melanie Subbiah\*  
Jared Kaplan† Prafulla Dhariwal Arvind Neelakantan Pranav Shyam Girish Sastry  
Amanda Askell Sandhini Agarwal Ariel Herbert-Voss Gretchen Krueger Tom Henighan  
Rewon Child Aditya Ramesh Daniel M. Ziegler Jeffrey Wu Clemens Winter  
Christopher Hesse Mark Chen Eric Sigler Mateusz Litwin Scott Gray  
Benjamin Chess Jack Clark Christopher Berner  
Sam McCandlish Alec Radford Ilya Sutskever Dario Amodei  
OpenAI

#### Abstract

Recent work has demonstrated substantial gains on many NLP tasks and benchmarks by pre-training on a large corpus of text followed by fine-tuning on a specific task. While typically task-agnostic in architecture, this method still requires task-specific fine-tuning datasets of thousands or tens of thousands of examples. By contrast, humans can generally perform a new language task from only a few examples or from simple context – so what is it about current NLP systems that largely struggle to do this? We show that scaling up language models greatly improves task-agnostic few-shot performance, sometimes even reaching competitiveness with prior state-of-the-art fine-tuning approaches. Specifically, we train GPT-3, an autoregressive language model with 175 billion parameters, 10x more than any previous non-sparse language model, and test its performance in the few-shot regime. For example, GPT-3 is capable without any task-specific fine-tuning, with tasks and few-shot demonstrations specified purely via text interaction with the model. GPT-3 achieves strong performance on many NLP datasets, including translation, question-answering, and cloze tasks, as well as several tasks that require on-the-fly reasoning or domain adaptation, such as unscrambling words, using a novel word in a sentence, or performing 3-digit arithmetic. At the same time, GPT-3 exhibits some interesting quirks, such as generating nonsensical text on some datasets where GPT-3 faces methodological issues related to training on large web corpora. Finally, we find that GPT-3 can generate samples of news articles which human evaluators have difficulty distinguishing from articles written by humans. We discuss broader societal impacts of this finding and of GPT-3 in general.

GPT3 (Brown, Mann, Ryder, Subbiah et al., 2020)

### speech synthesis



Wavenet (van den Oord et al., 2016)

# Today's Agenda

## **Pretext tasks from image transformations**

- Rotation, inpainting, rearrangement, coloring

## **Contrastive representation learning**

- Intuition and formulation
- Instance contrastive learning: SimCLR and MOCO
- Sequence contrastive learning: CPC

# Today's Agenda

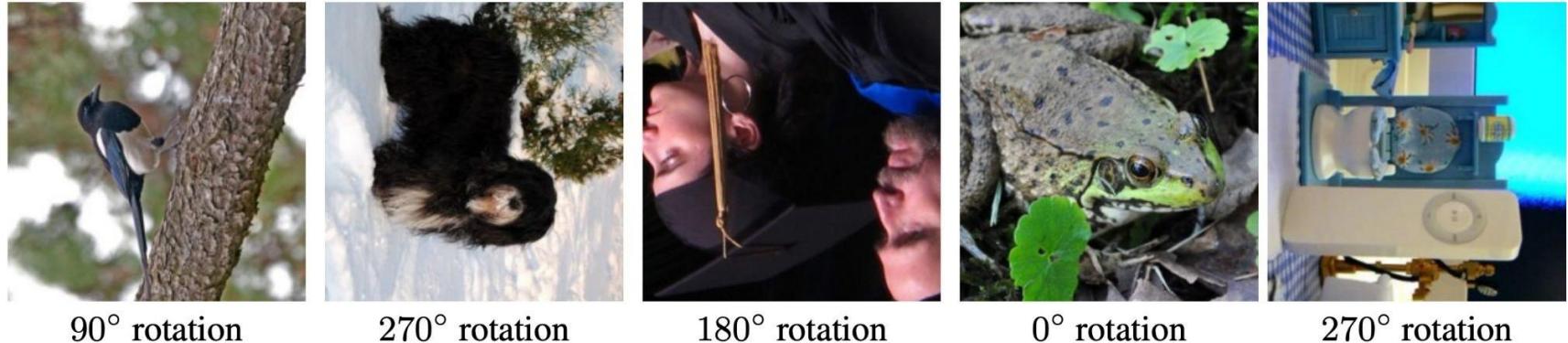
## **Pretext tasks from image transformations**

- Rotation, inpainting, rearrangement, coloring

## **Contrastive representation learning**

- Intuition and formulation
- Instance contrastive learning: SimCLR and MOCO
- Sequence contrastive learning: CPC

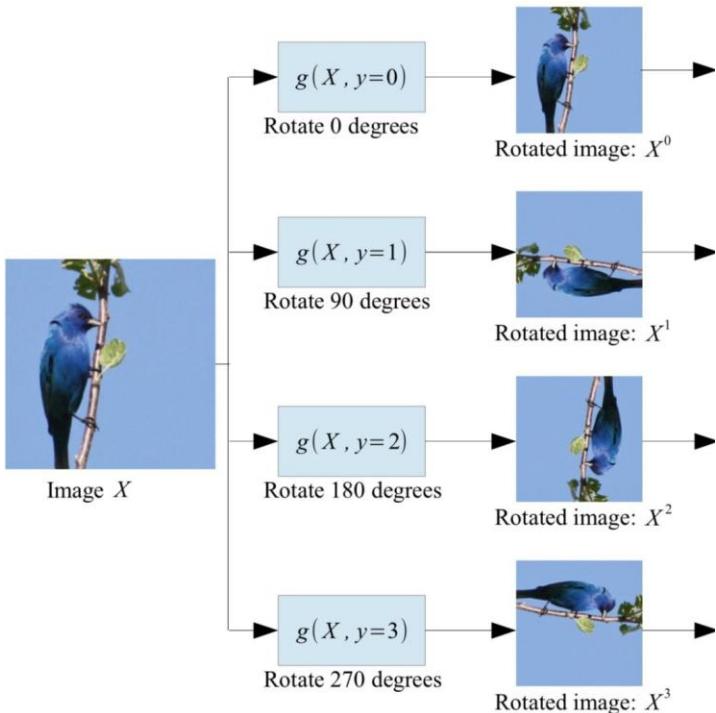
# Pretext task: predict rotations



**Hypothesis:** a model could recognize the correct rotation of an object only if it has the “visual commonsense” of what the object should look like unperturbed.

(Image source: [Gidaris et al. 2018](#))

# Pretext task: predict rotations

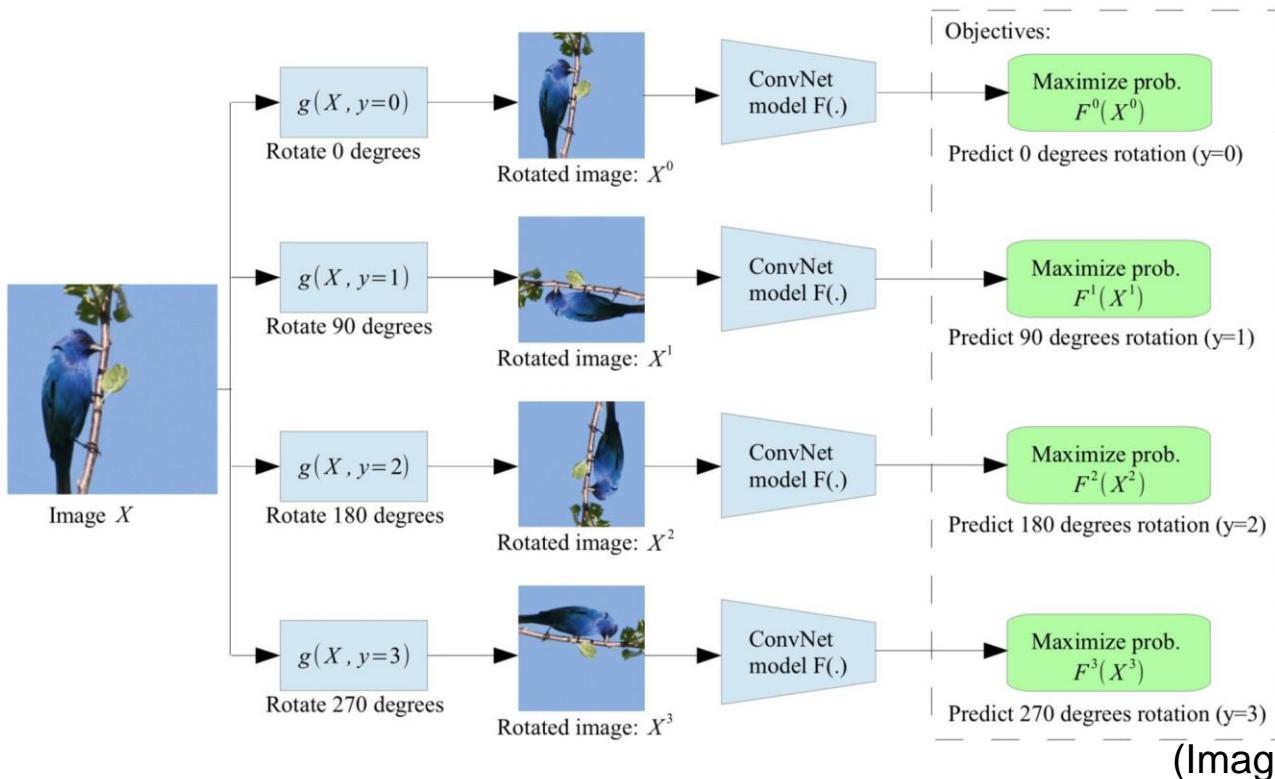


Self-supervised learning by rotating the entire input images.

The model learns to predict which rotation is applied (4-way classification)

(Image source: [Gidaris et al. 2018](#))

# Pretext task: predict rotations

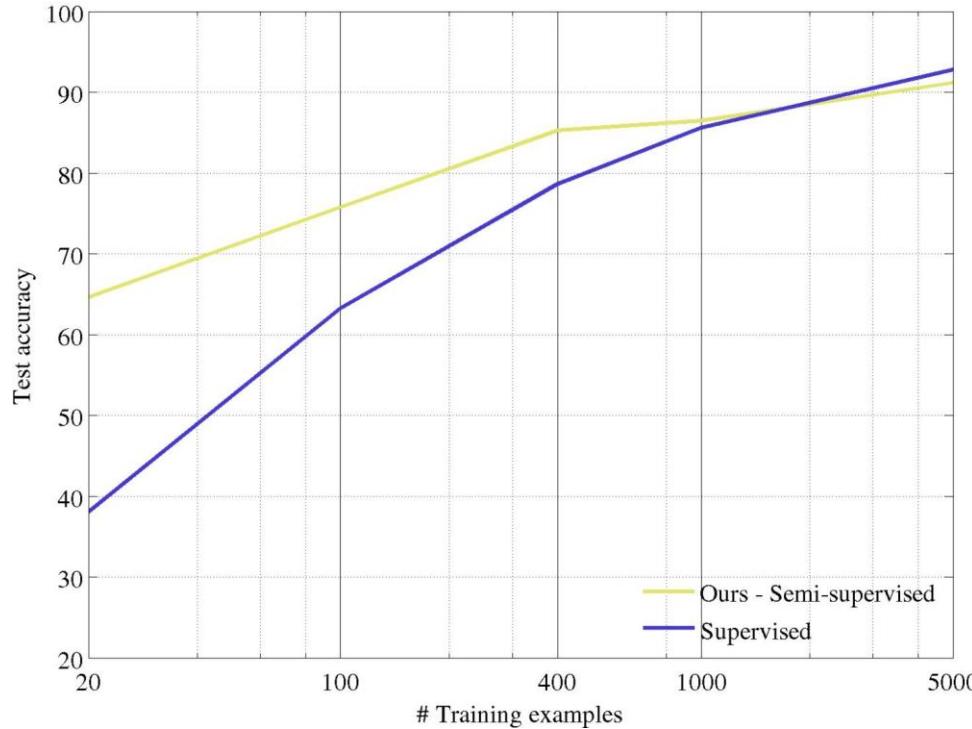


Self-supervised learning by rotating the entire input images.

The model learns to predict which rotation is applied (4-way classification)

(Image source: [Gidaris et al. 2018](#))

# Evaluation on semi-supervised learning



Self-supervised learning on  
**CIFAR10** (entire training set).

Freeze conv1 + conv2  
Learn **conv3 + linear** layers  
with subset of labeled  
CIFAR10 data (classification).

(Image source: [Gidaris et al. 2018](#))

# Transfer learned features to supervised learning

	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	45.6
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
(Ours) RotNet	<b>70.87</b>	<b>72.97</b>	<b>54.4</b>
			<b>39.1</b>

Self-supervised learning with rotation prediction

Pretrained with full  
ImageNet supervision

No pretraining

Self-supervised learning on  
**ImageNet** (entire training  
set) with AlexNet.

Finetune on labeled data  
from **Pascal VOC 2007**.

source: [Gidaris et al. 2018](#)

# Visualize learned visual attentions



Conv1  $27 \times 27$    Conv3  $13 \times 13$    Conv5  $6 \times 6$

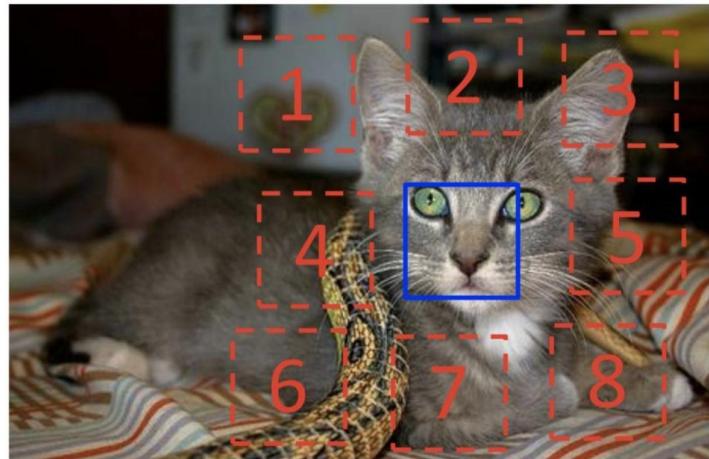
(a) Attention maps of supervised model

Conv1  $27 \times 27$    Conv3  $13 \times 13$    Conv5  $6 \times 6$

(b) Attention maps of our self-supervised model

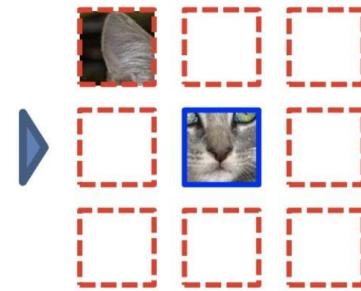
(Image source: [Gidaris et al. 2018](#))

# Pretext task: predict relative patch locations



$$X = (\text{[cat eye]}, \text{[cat ear]}); Y = 3$$

Example:



Question 1:

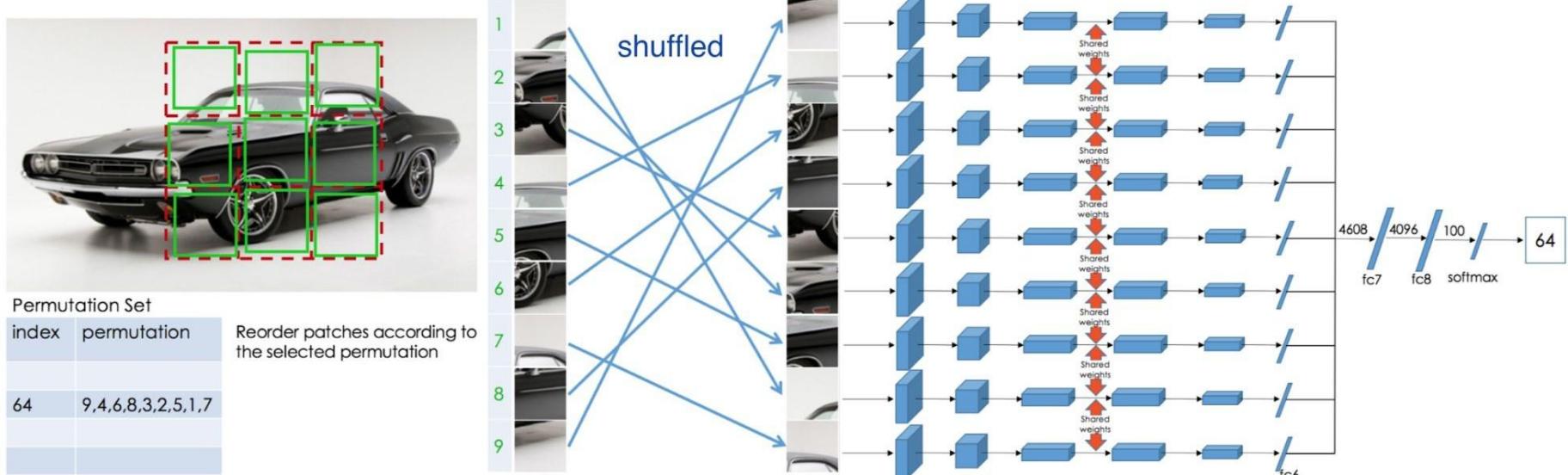


Question 2:



(Image source: [Doersch et al., 2015](#))

# Pretext task: solving “jigsaw puzzles”



(Image source: [Noroozi & Favaro, 2016](#))

# Transfer learned features to supervised learning

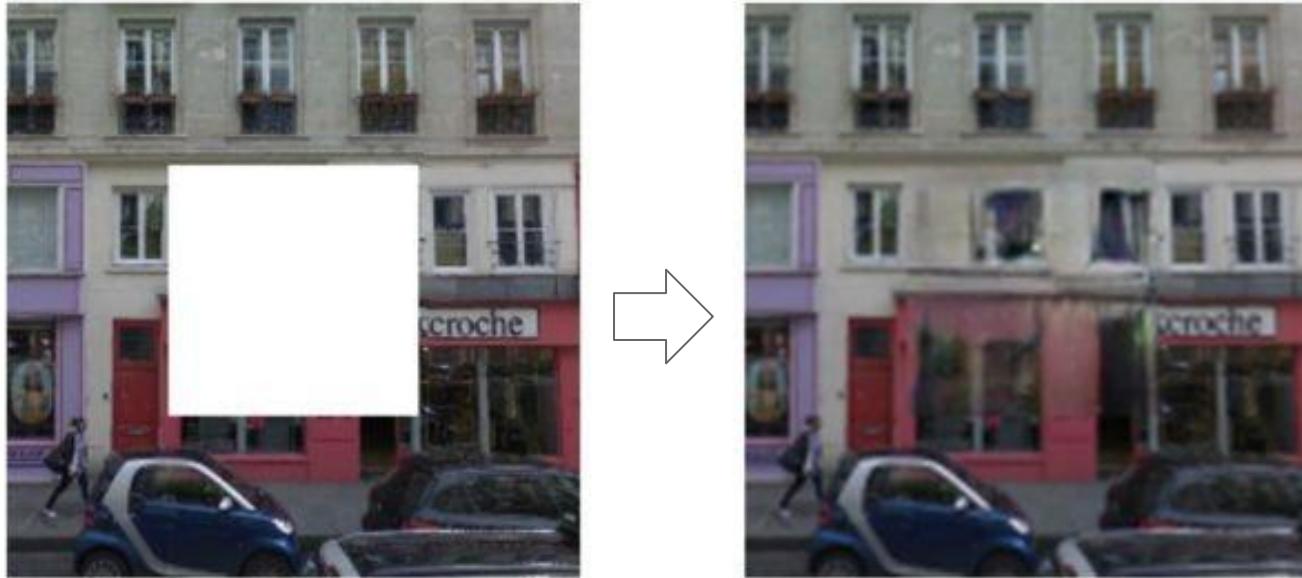
Table 1: Results on PASCAL VOC 2007 Detection and Classification. The results of the other methods are taken from Pathak *et al.* [30].

Method	Pretraining time	Supervision	Classification	Detection	Segmentation
Krizhevsky <i>et al.</i> [25]	3 days	1000 class labels	<b>78.2%</b>	<b>56.8%</b>	<b>48.0%</b>
Wang and Gupta[39]	1 week	motion	58.4%	44.0%	-
Doersch <i>et al.</i> [10]	4 weeks	context	55.3%	46.6%	-
Pathak <i>et al.</i> [30]	14 hours	context	56.5%	44.5%	29.7%
Ours	2.5 days	context	<b>67.6%</b>	<b>53.2%</b>	<b>37.6%</b>

“Ours” is feature learned from solving image Jigsaw puzzles (Noroozi & Favaro, 2016). Doersch et al. is the method with relative patch location

(source: [Noroozi & Favaro, 2016](#))

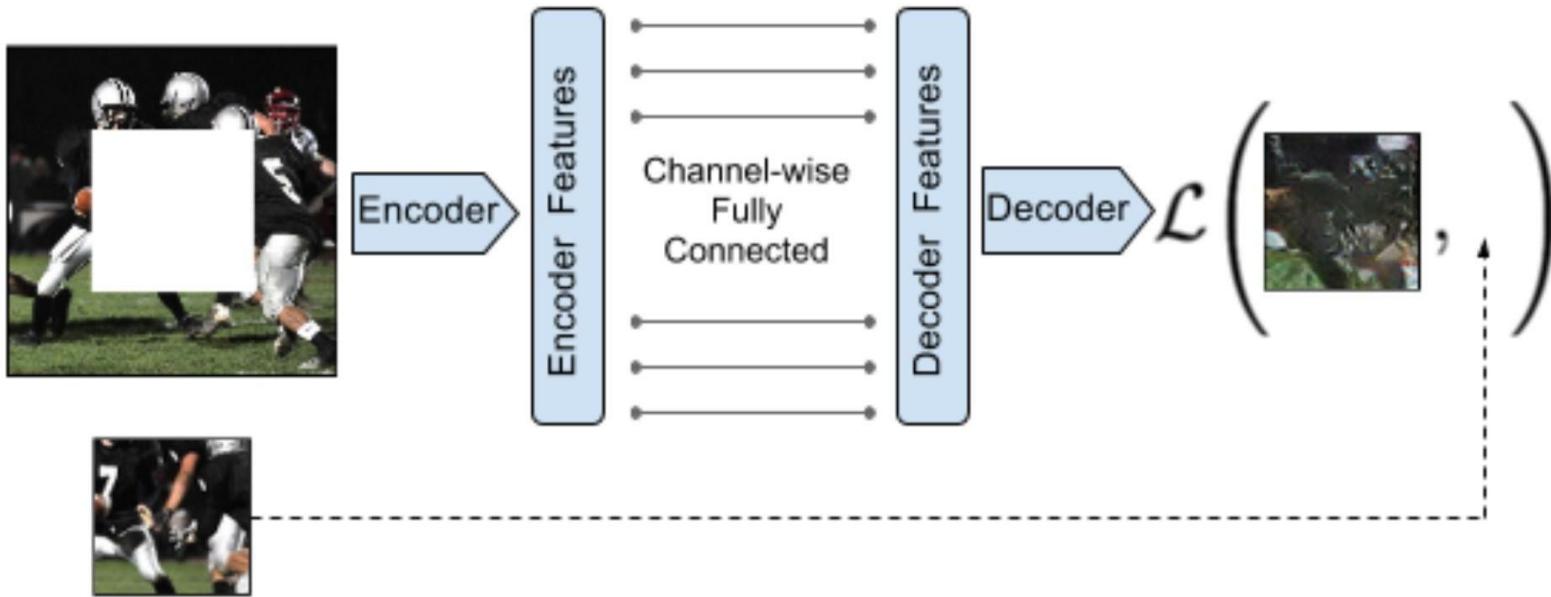
# Pretext task: predict missing pixels (inpainting)



*Context Encoders: Feature Learning by Inpainting* (Pathak et al., 2016)

Source: [Pathak et al., 2016](#)

# Learning to inpaint by reconstruction



Learning to reconstruct the missing pixels

Source: [Pathak et al., 2016](#)

# Inpainting evaluation



Input (context)



reconstruction

Source: [Pathak et al., 2016](#)

# Learning to inpaint by reconstruction

Loss = reconstruction + adversarial learning

$$L(x) = L_{recon}(x) + L_{adv}(x)$$

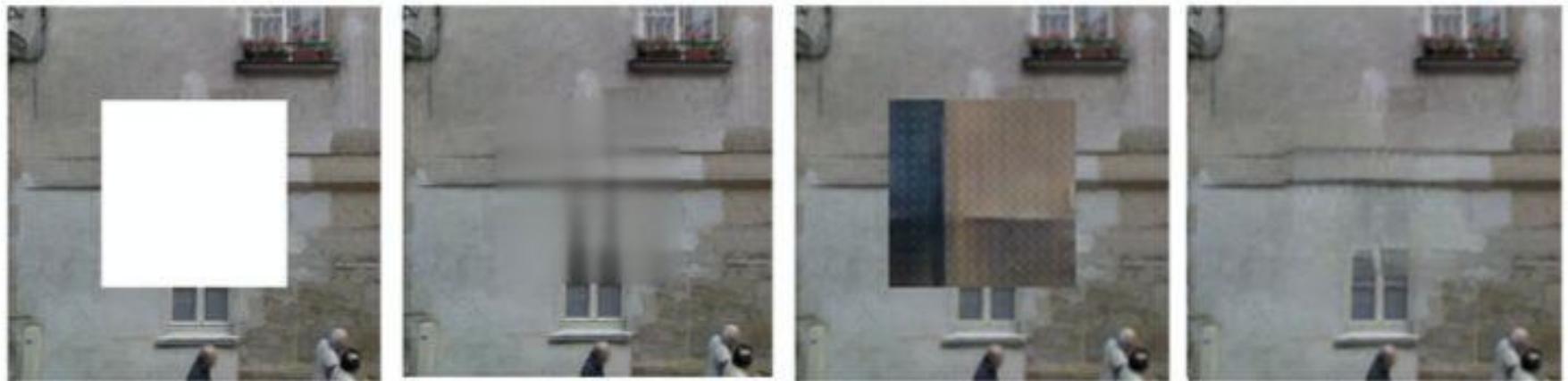
$$L_{recon}(x) = \left\| M * (x - F_\theta((1 - M) * x)) \right\|_2^2$$

$$L_{adv} = \max_D \mathbb{E}[\log(D(x))] + \log(1 - D(F((1 - M) * x)))]$$

Adversarial loss between “real” images and *inpainted images*

Source: [Pathak et al., 2016](#)

# Inpainting evaluation



Input (context)

reconstruction

adversarial

recon + adv

Source: [Pathak et al., 2016](#)

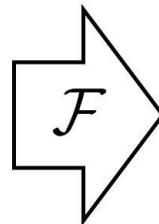
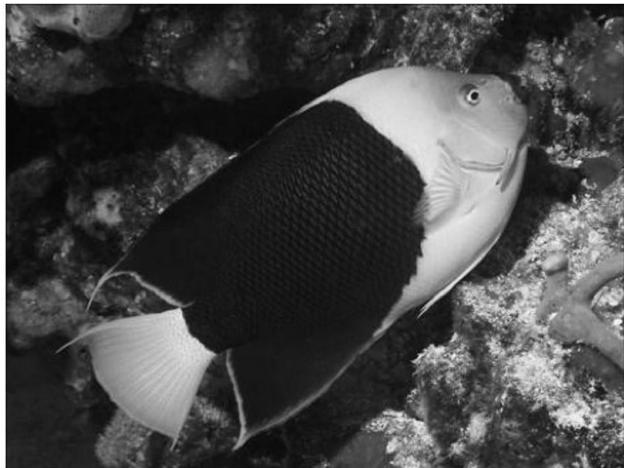
# Transfer learned features to supervised learning

Pretraining Method	Supervision	Pretraining time	Classification	Detection	Segmentation
ImageNet [26]	1000 class labels	3 days	78.2%	56.8%	48.0%
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%	-
Wang <i>et al.</i> [39]	motion	1 week	58.7%	47.4%	-
Doersch <i>et al.</i> [7]	relative context	4 weeks	55.3%	46.6%	-
Ours	context	14 hours	56.5%	44.5%	30.0%

Self-supervised learning on ImageNet training set, transfer to classification (Pascal VOC 2007), detection (Pascal VOC 2007), and semantic segmentation (Pascal VOC 2012)

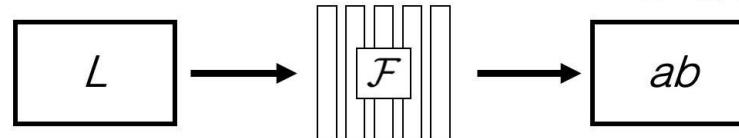
Source: [Pathak et al., 2016](#)

# Pretext task: image coloring



Grayscale image:  $L$  channel

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

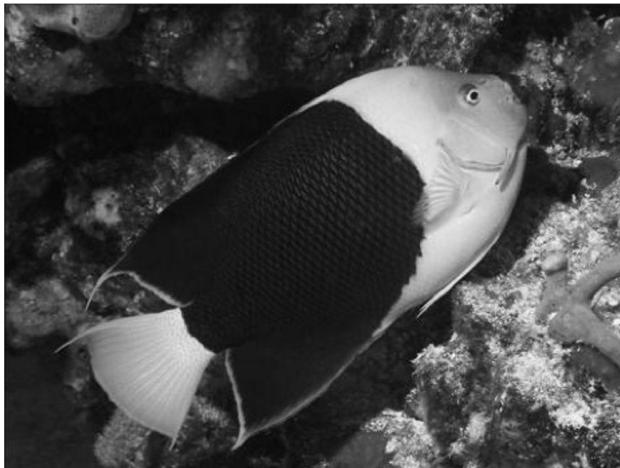


Color information:  $ab$  channels

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

Source: Richard Zhang / Phillip Isola<sup>5</sup>

# Pretext task: image coloring

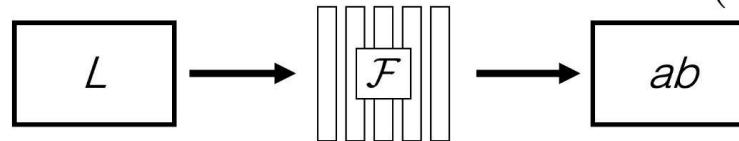


Grayscale image:  $L$  channel

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



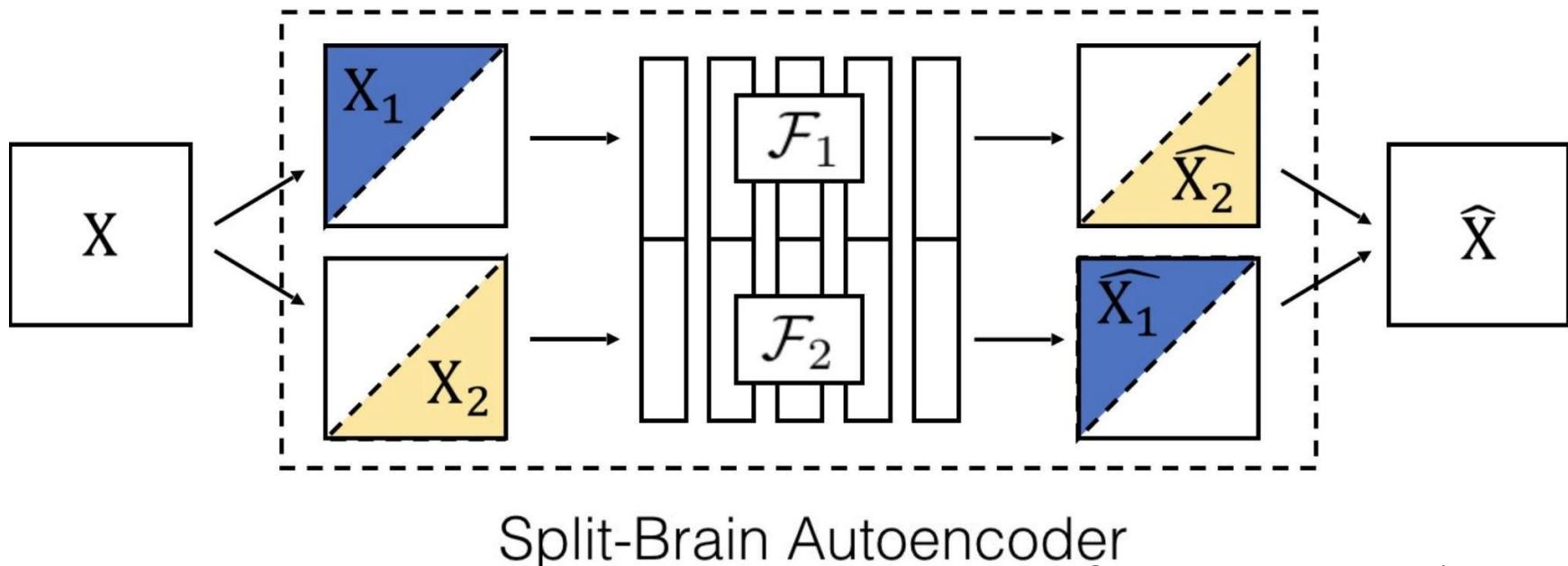
Concatenate  $(L, ab)$  channels  
 $(\mathbf{X}, \hat{\mathbf{Y}})$



Source: Richard Zhang / Phillip Isola

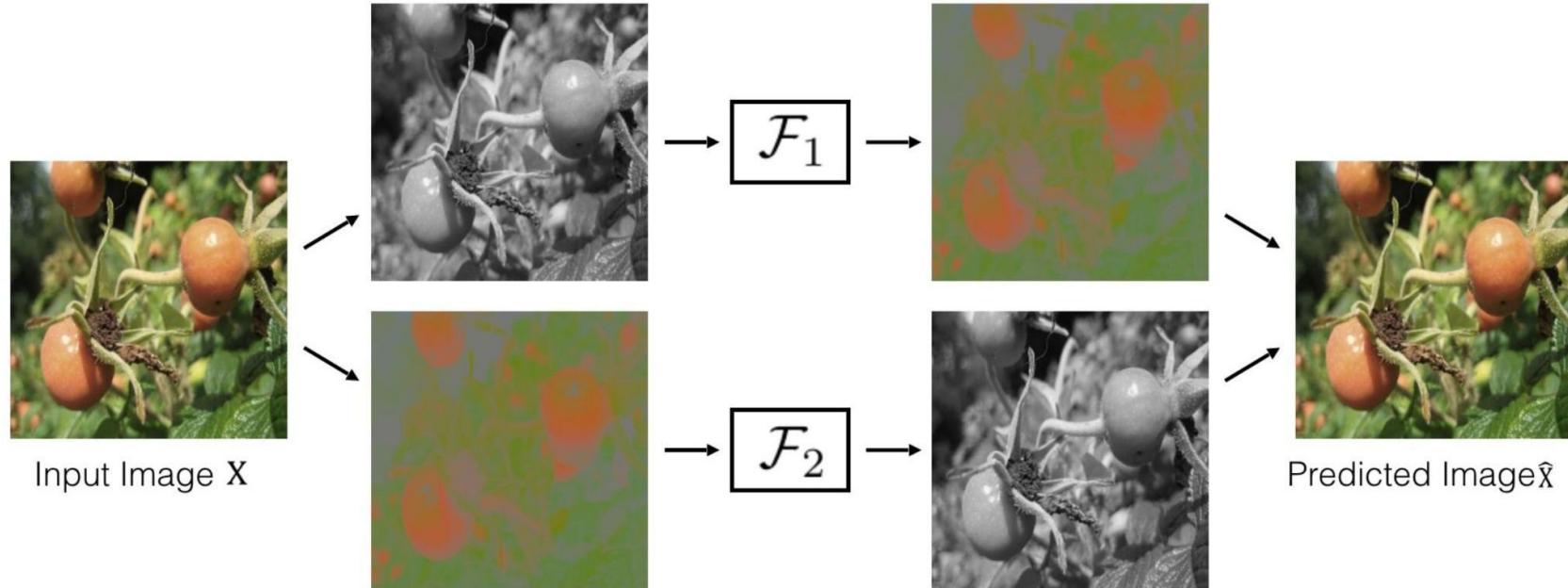
# Learning features from colorization: Split-brain Autoencoder

**Idea:** cross-channel predictions



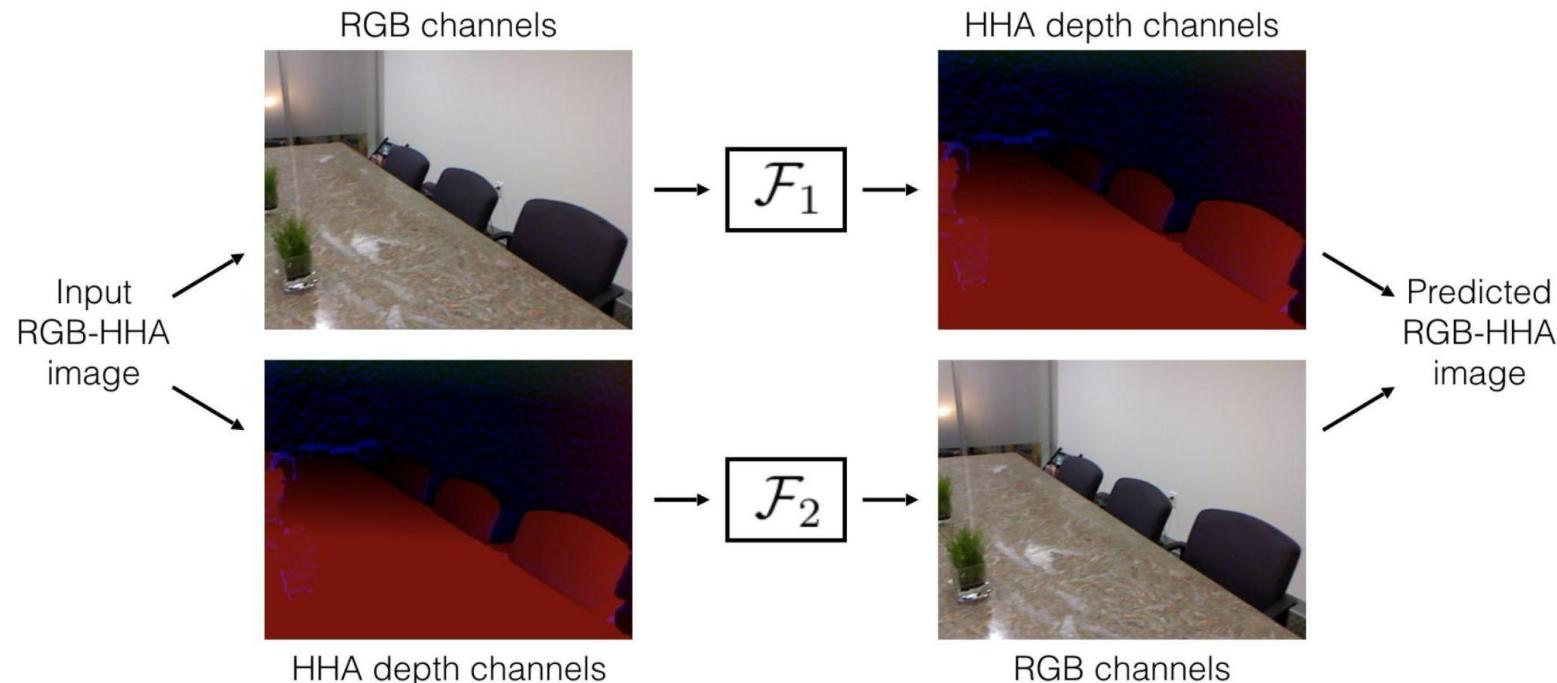
Source: Richard Zhang / Phillip Isola

# Learning features from colorization: Split-brain Autoencoder



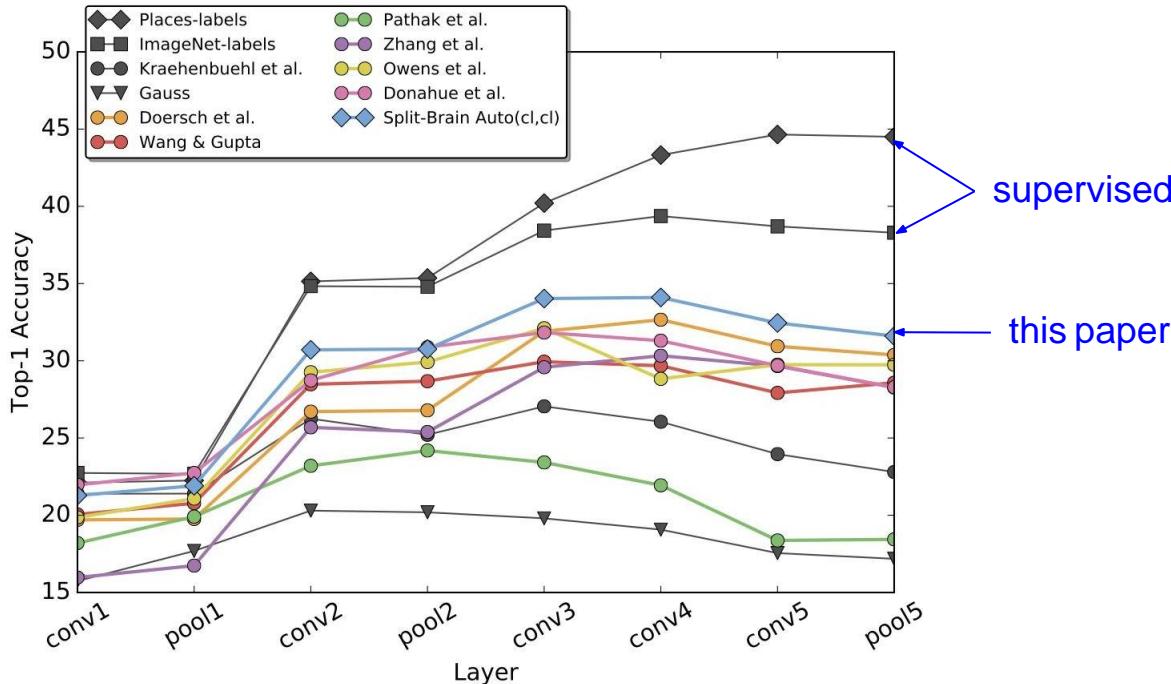
Source: Richard Zhang / Phillip Isola

# Learning features from colorization: Split-brain Autoencoder



Source: Richard Zhang / Phillip Isola

# Transfer learned features to supervised learning



Self-supervised learning on **ImageNet** (entire training set).

Use concatenated features from  $F_1$  and  $F_2$

Labeled data is from the **Places** (Zhou 2016).

Source: [Zhang et al., 2017](#)

# Pretext task: image coloring



Source: Richard Zhang / Phillip Isola

# Pretext task: image coloring



Source: Richard Zhang / Phillip Isola

# Pretext task: video coloring

**Idea:** model the *temporal coherence* of colors in videos

reference frame



$t = 0$

how should I color these frames?



$t = 1$



$t = 2$



$t = 3$

...

Source: [Vondrick et al., 2018](#)

# Pretext task: video coloring

**Idea:** model the *temporal coherence* of colors in videos

reference frame



$t = 0$

how should I color these frames?

**Should be the same color!**



$t = 1$



$t = 2$



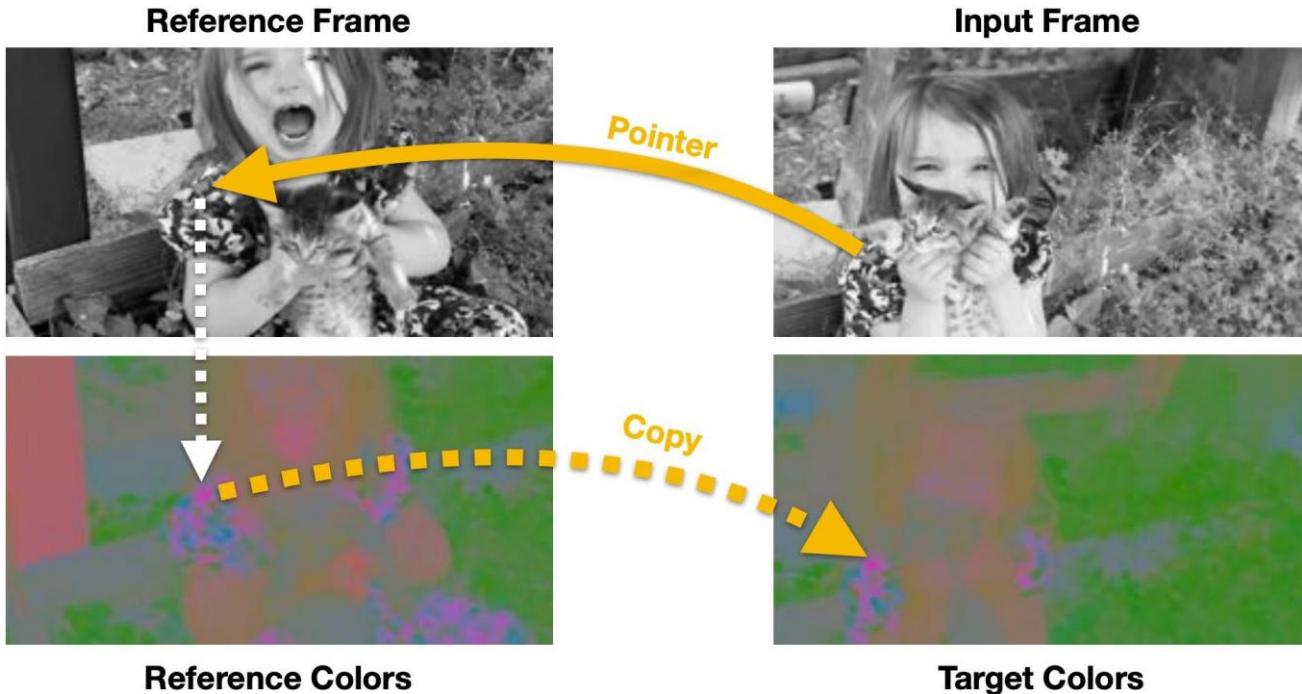
$t = 3$

...

**Hypothesis:** learning to color video frames should allow model to learn to track regions or objects without labels!

Source: [Vondrick et al., 2018](#)

# Learning to color videos



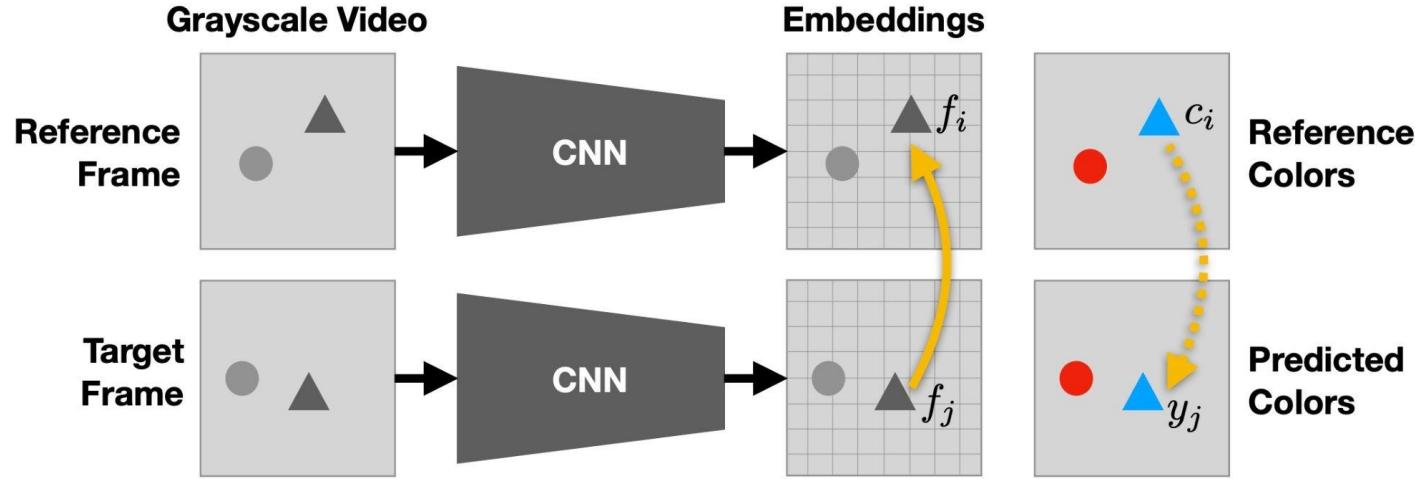
**Learning objective:**

Establish mappings between reference and target frames in a learned feature space.

Use the mapping as “pointers” to copy the correct color (LAB).

Source: [Vondrick et al., 2018](#)

# Learning to color videos

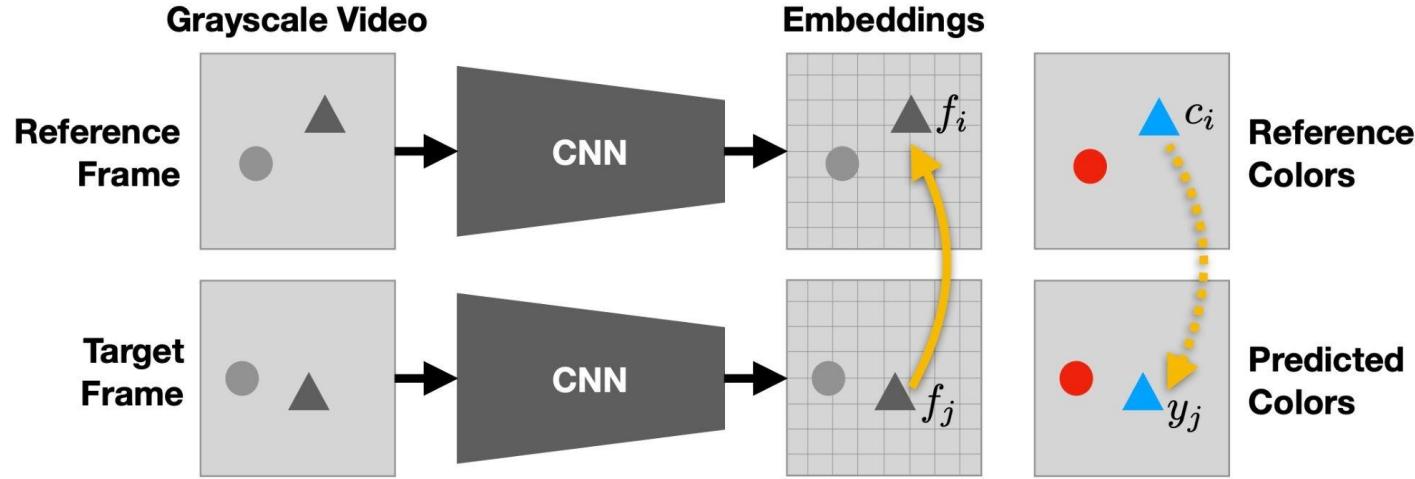


attention map on the  
reference frame

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

Source: [Vondrick et al., 2018](#)

# Learning to color videos



attention map on the  
reference frame

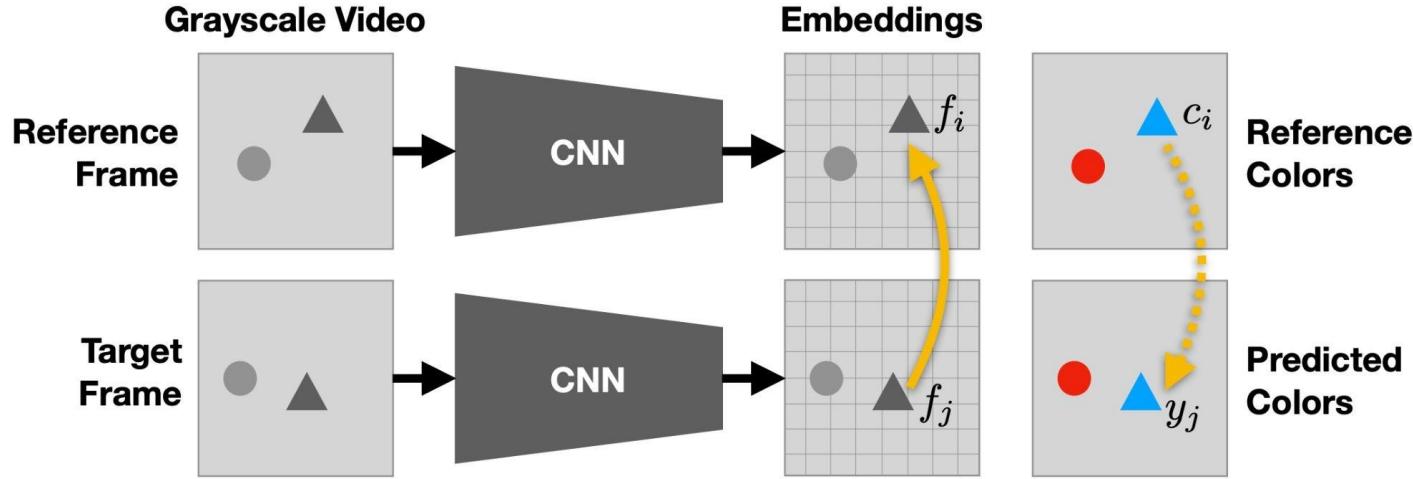
predicted color = weighted  
sum of the reference color

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

$$y_j = \sum_i A_{ij} c_i$$

Source: [Vondrick et al., 2018](#)

# Learning to color videos



attention map on the  
reference frame

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

predicted color = weighted  
sum of the reference color

$$y_j = \sum_i A_{ij} c_i$$

loss between predicted color  
and ground truth color

$$\min_{\theta} \sum_j \mathcal{L}(y_j, c_j)$$

Source: [Vondrick et al., 2018](#)

# Colorizing videos (qualitative)

reference frame



target frames (gray)



predicted color



Source: [Google AI blog post](#)

# Colorizing videos (qualitative)

reference frame



target frames (gray)



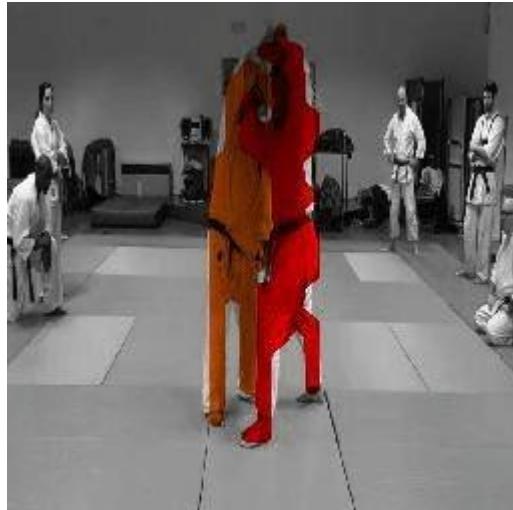
predicted color



Source: [Google AI blog post](#)

# Tracking emerges from colorization

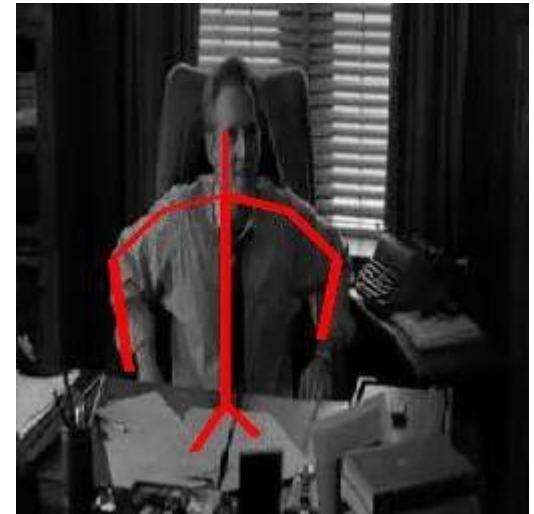
Propagate segmentation masks using learned attention



Source: [Google AI blog post](#)

# Tracking emerges from colorization

Propagate pose keypoints using learned attention



Source: [Google AI blog post](#)

# Summary: pretext tasks from image transformations

- Pretext tasks focus on “visual common sense”, e.g., predict rotations, inpainting, rearrangement, and colorization.
- The models are forced learn good features about natural images, e.g., semantic representation of an object category, in order to solve the pretext tasks.
- We don’t care about the performance of these pretext tasks, but rather how useful the learned features are for downstream tasks (classification, detection, segmentation).

# Summary: pretext tasks from image transformations

- Pretext tasks focus on “visual common sense”, e.g., predict rotations, inpainting, rearrangement, and colorization.
- The models are forced learn good features about natural images, e.g., semantic representation of an object category, in order to solve the pretext tasks.
- We don’t care about the performance of these pretext tasks, but rather how useful the learned features are for downstream tasks (classification, detection, segmentation).
- Problems: 1) coming up with individual pretext tasks is tedious, and 2) the learned representations may not be general.

# Pretext tasks from image transformations

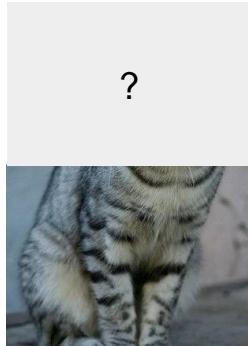
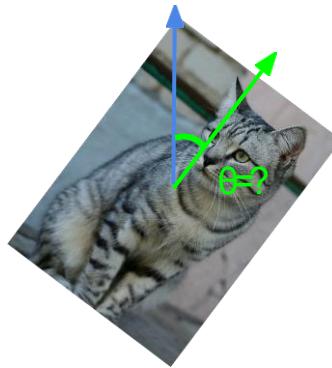


image completion



rotation prediction



“jigsaw puzzle”

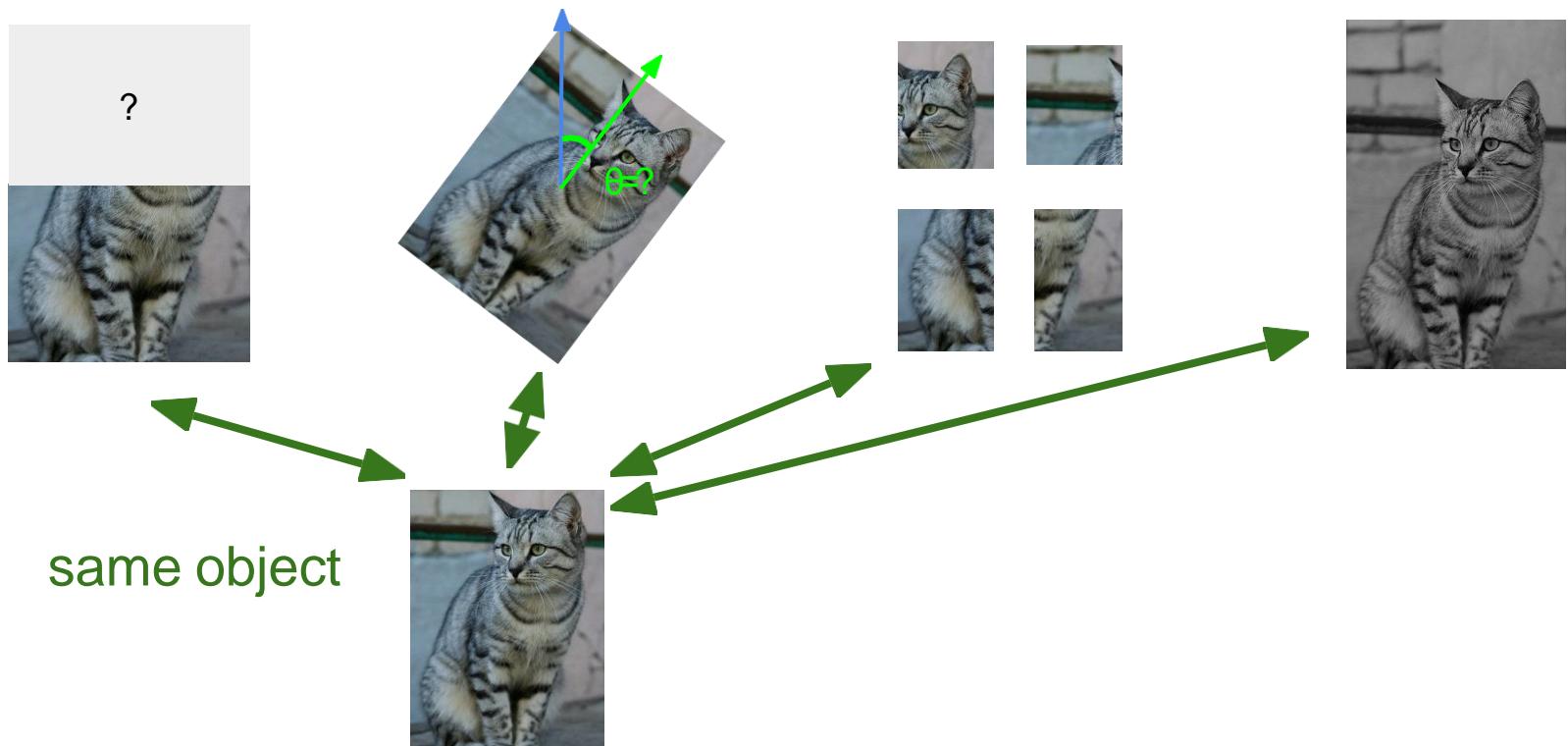


colorization

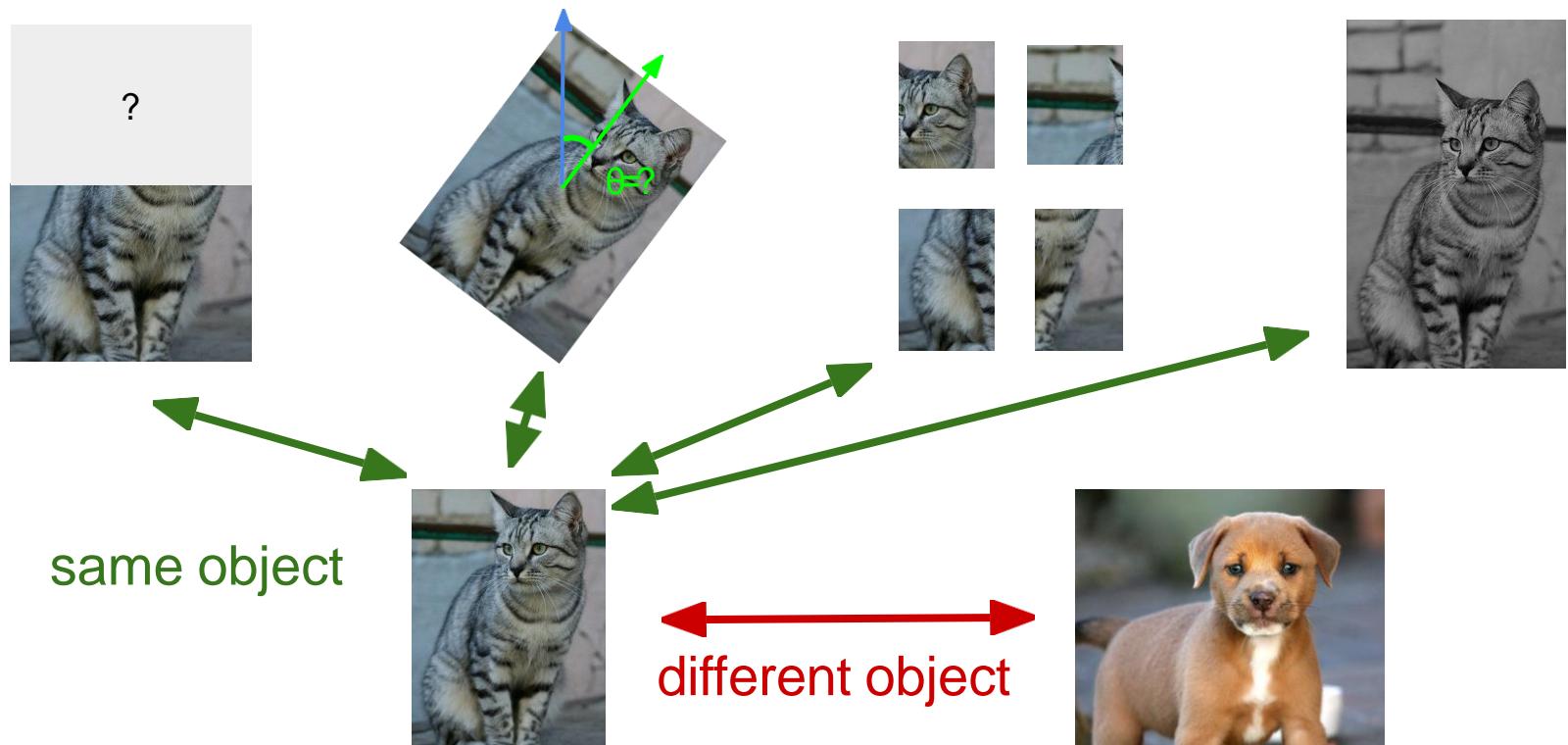
Learned representations may be tied to a specific pretext task!

Can we come up with a more general pretext task?

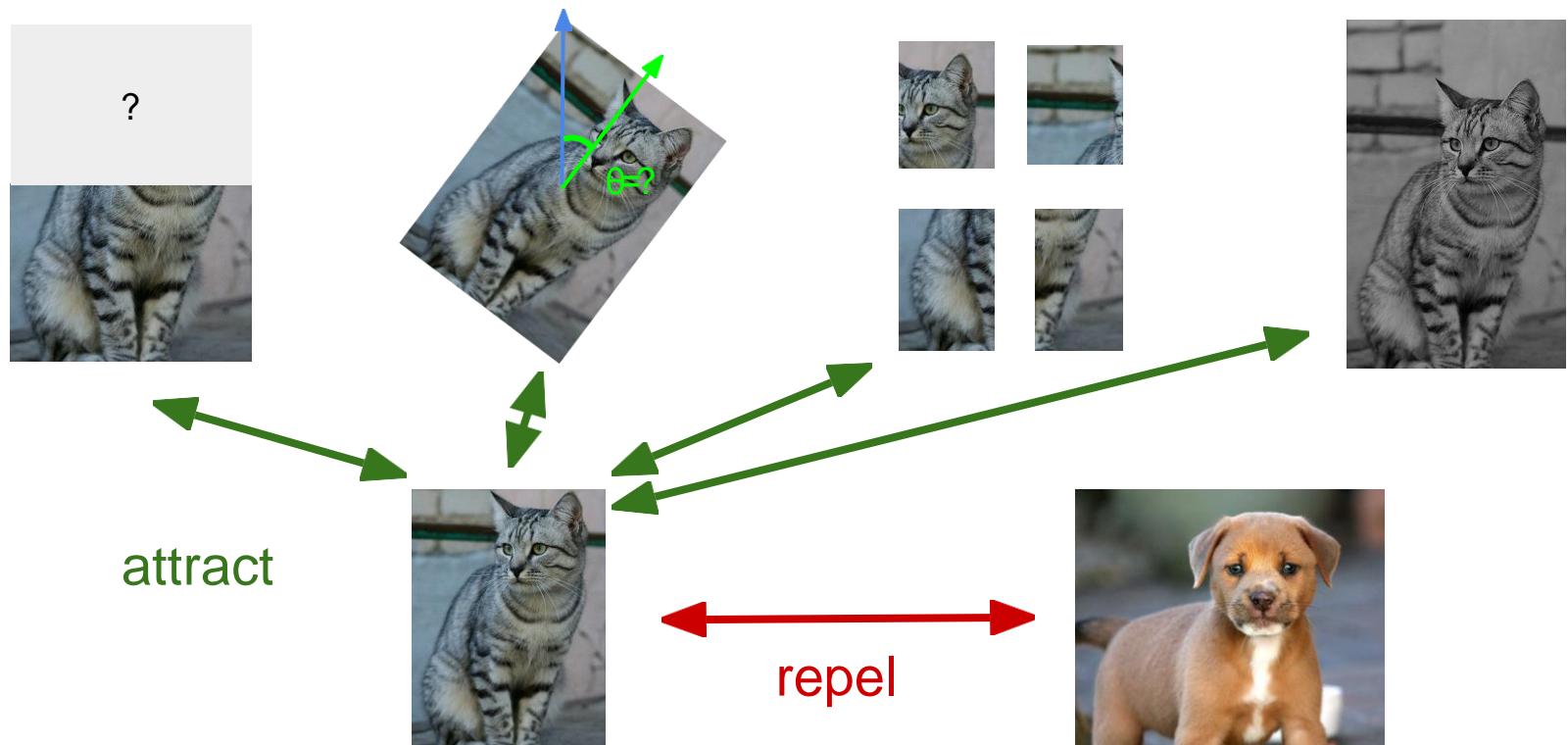
# A more general pretext task?



# A more general pretext task?



# Contrastive Representation Learning



# Today's Agenda

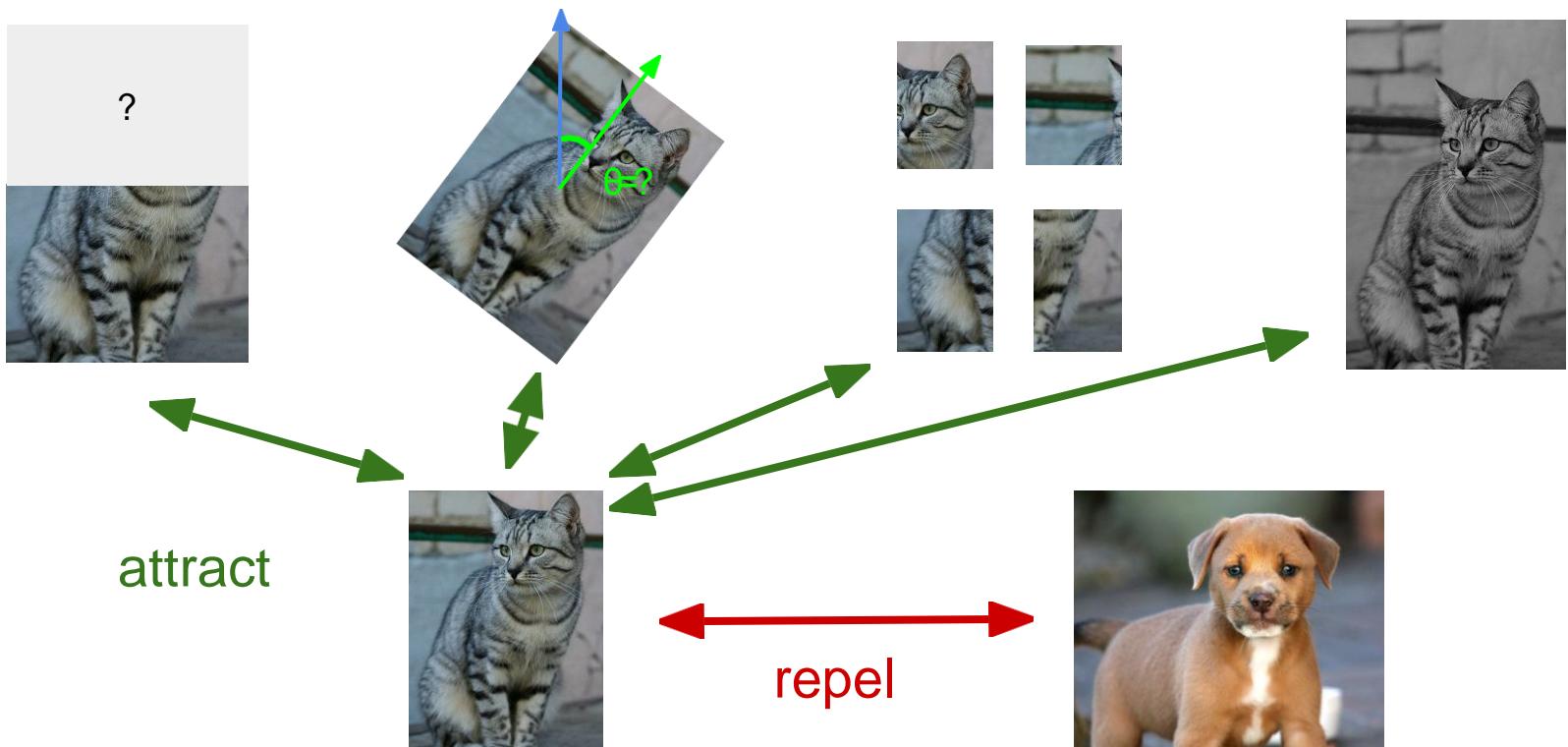
## Pretext tasks from image transformations

- Rotation, inpainting, rearrangement, coloring

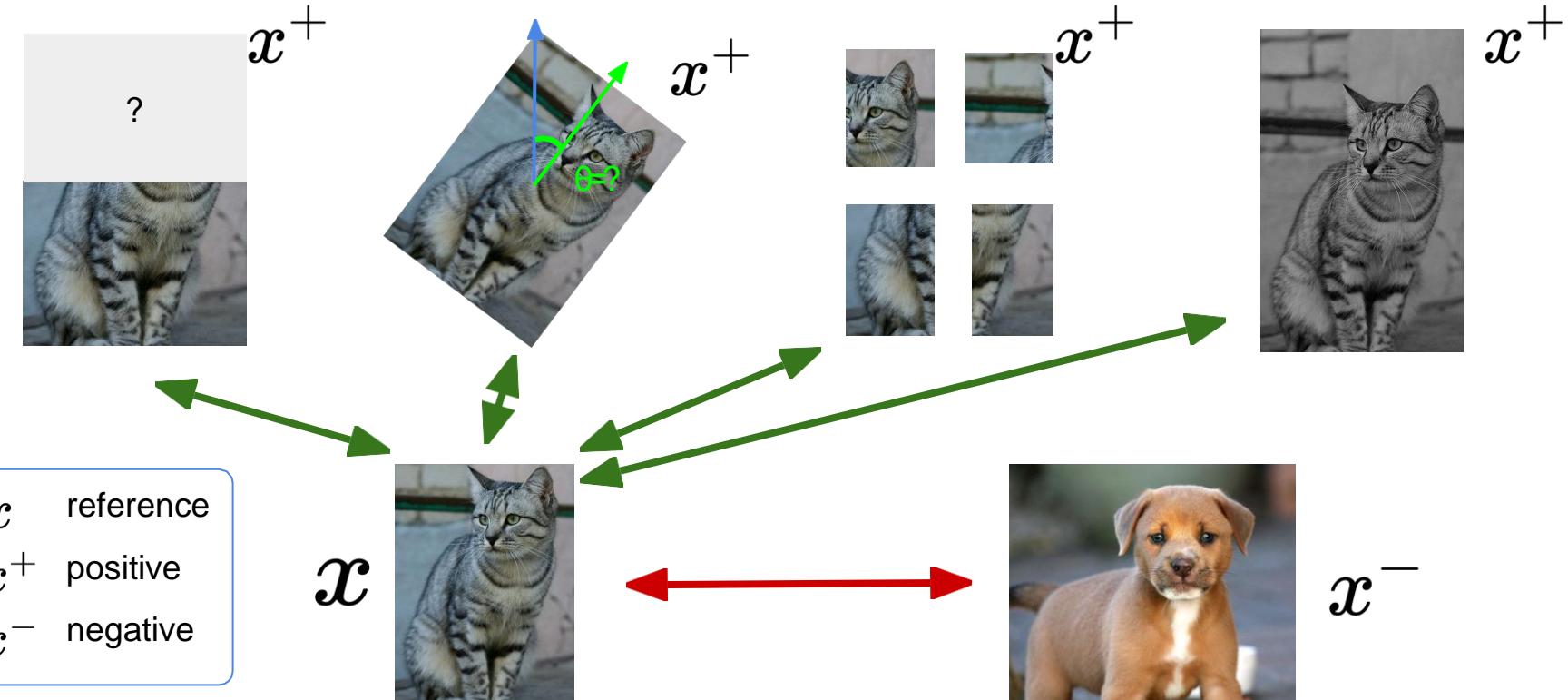
## Contrastive representation learning

- Intuition and formulation
- Instance contrastive learning: SimCLR and MOCO
- Sequence contrastive learning: CPC

# Contrastive Representation Learning



# Contrastive Representation Learning



# A formulation of contrastive learning

What we want:

$$\text{score}(f(x), f(x^+)) \gg \text{score}(f(x), f(x^-))$$

$x$ : reference sample;  $x^+$  positive sample;  $x^-$  negative sample

Given a chosen score function, we aim to learn an **encoder function**  $f$  that yields high score for positive pairs  $(x, x^+)$  and low scores for negative pairs  $(x, x^-)$ .

# A formulation of contrastive learning

Loss function given 1 positive sample and N - 1 negative samples:

$$L = -\mathbb{E}_X \left[ \log \frac{\exp(s(f(x), f(x^+)))}{\exp(s(f(x), f(x^+))) + \sum_{j=1}^{N-1} \exp(s(f(x), f(x_j^-)))} \right]$$

# A formulation of contrastive learning

Loss function given 1 positive sample and  $N - 1$  negative samples:

$$L = -\mathbb{E}_X \left[ \log \frac{\exp(s(f(x), f(x^+)))}{\exp(s(f(x), f(x^+))) + \sum_{j=1}^{N-1} \exp(s(f(x), f(x_j^-)))} \right]$$



$x$



$x^+$



$x$



$x_1^-$



$x_2^-$



$x_3^-$

...

# A formulation of contrastive learning

Loss function given 1 positive sample and  $N - 1$  negative samples:

$$L = -\mathbb{E}_X \left[ \log \frac{\exp(s(f(x), f(x^+)))}{\exp(s(f(x), f(x^+))) + \sum_{j=1}^{N-1} \exp(s(f(x), f(x_j^-)))} \right]$$

score for the positive pair      score for the N-1 negative pairs

## This seems familiar ...

# A formulation of contrastive learning

Loss function given 1 positive sample and  $N - 1$  negative samples:

$$L = -\mathbb{E}_X \left[ \log \frac{\exp(s(f(x), f(x^+)))}{\exp(s(f(x), f(x^+))) + \sum_{j=1}^{N-1} \exp(s(f(x), f(x_j^-)))} \right]$$

score for the positive pair      score for the N-1 negative pairs

This seems familiar ...

# Cross entropy loss for a N-way softmax classifier!

I.e., learn to find the positive sample from the N samples

# A formulation of contrastive learning

Loss function given 1 positive sample and  $N - 1$  negative samples:

$$L = -\mathbb{E}_X \left[ \log \frac{\exp(s(f(x), f(x^+)))}{\exp(s(f(x), f(x^+))) + \sum_{j=1}^{N-1} \exp(s(f(x), f(x_j^-)))} \right]$$

Commonly known as the InfoNCE loss ([van den Oord et al., 2018](#))

A *lower bound* on the mutual information between  $f(x)$  and  $f(x^+)$

$$MI[f(x), f(x^+)] - \log(N) \geq -L$$

The larger the negative sample size ( $N$ ), the tighter the bound

Detailed derivation: [Poole et al., 2019](#)

# SimCLR: A Simple Framework for Contrastive Learning

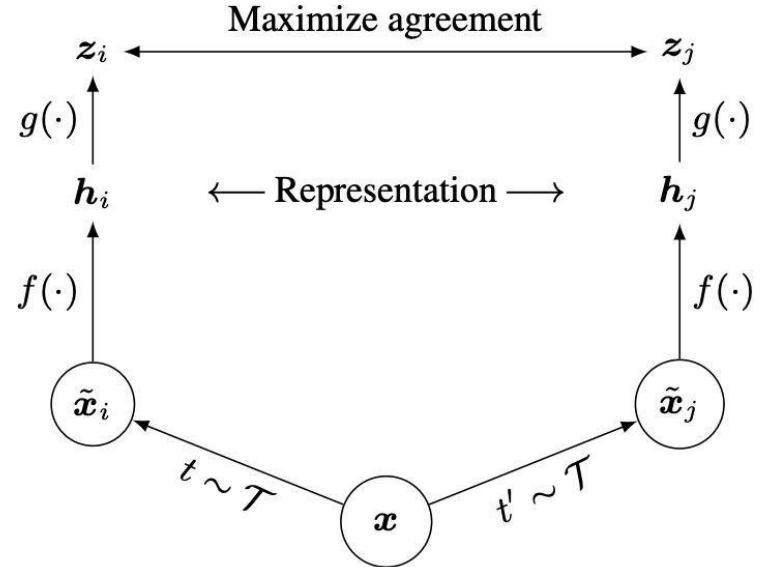
Cosine similarity as the score function:

$$s(u, v) = \frac{u^T v}{\|u\| \|v\|}$$

Use a projection network  $\mathbf{g}(\cdot)$  to project features to a space where contrastive learning is applied

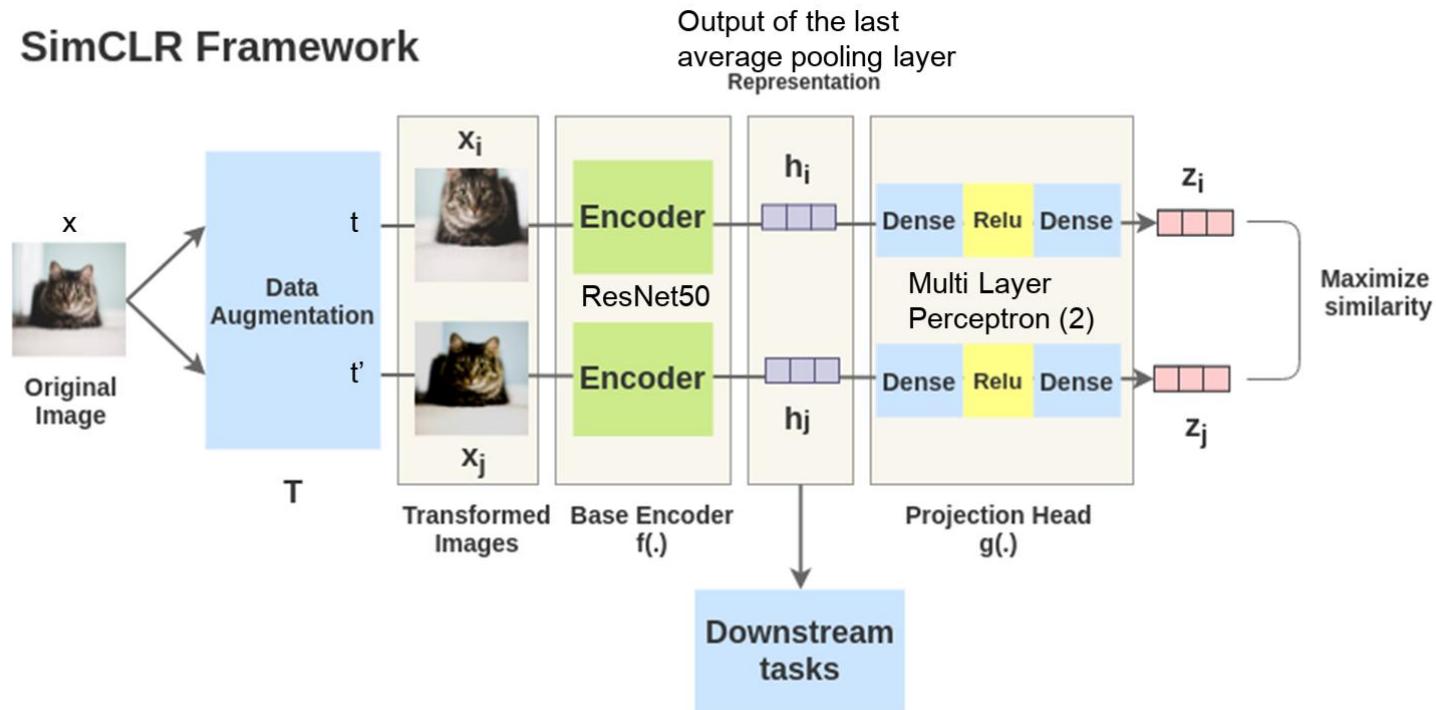
Generate positive samples through data augmentation:

- random cropping, random color distortion, and random blur.



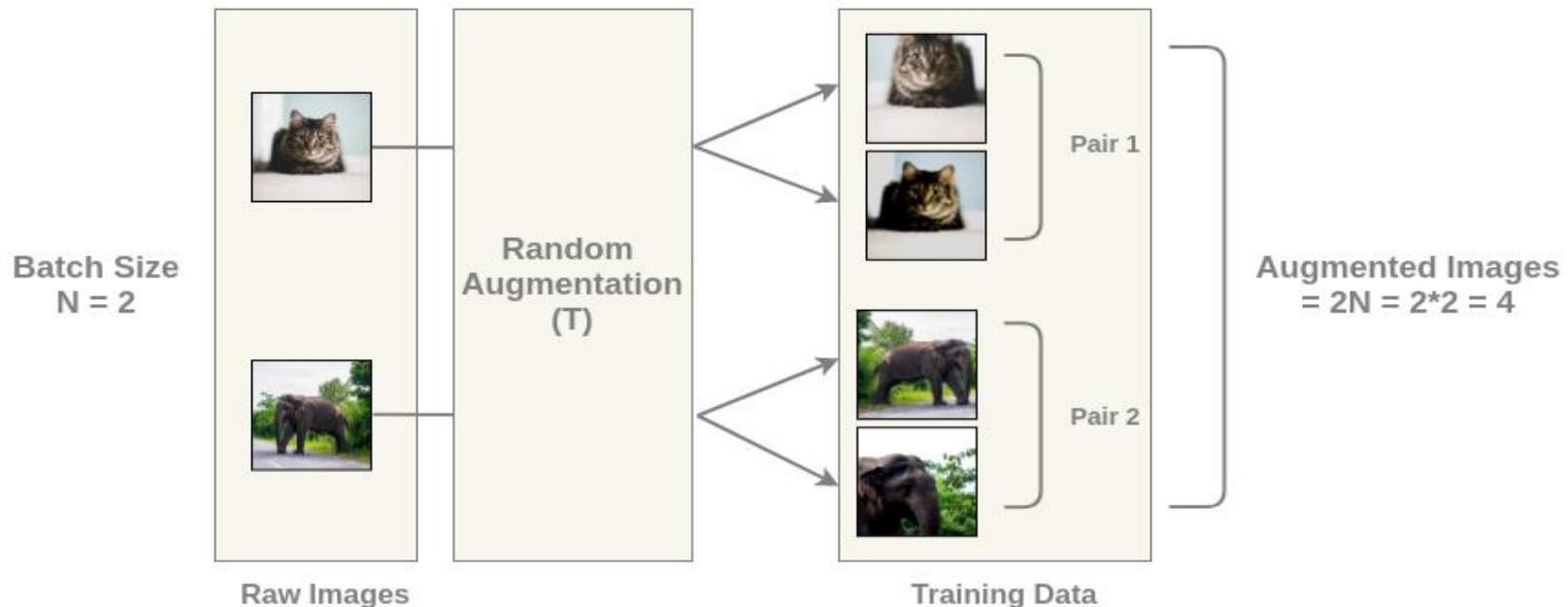
Source: [Chen et al., 2020](#)

# SimCLR: A Simple Framework for Contrastive Learning



# SimCLR: A Simple Framework for Contrastive Learning

Preparing similar pairs in a batch



# SimCLR: generating positive samples from data augmentation



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate  $\{90^\circ, 180^\circ, 270^\circ\}$



(g) Cutout



(h) Gaussian noise



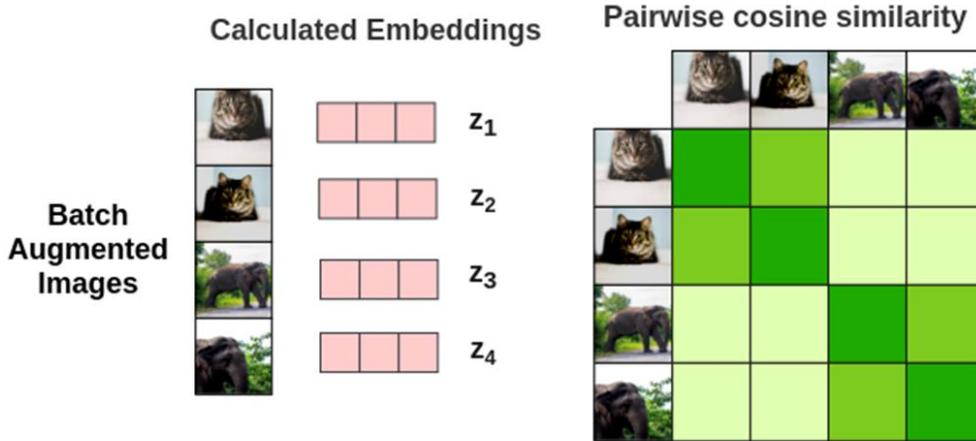
(i) Gaussian blur



(j) Sobel filtering

Source: [Chen et al., 2020](#)

# SimCLR: mini-batch training



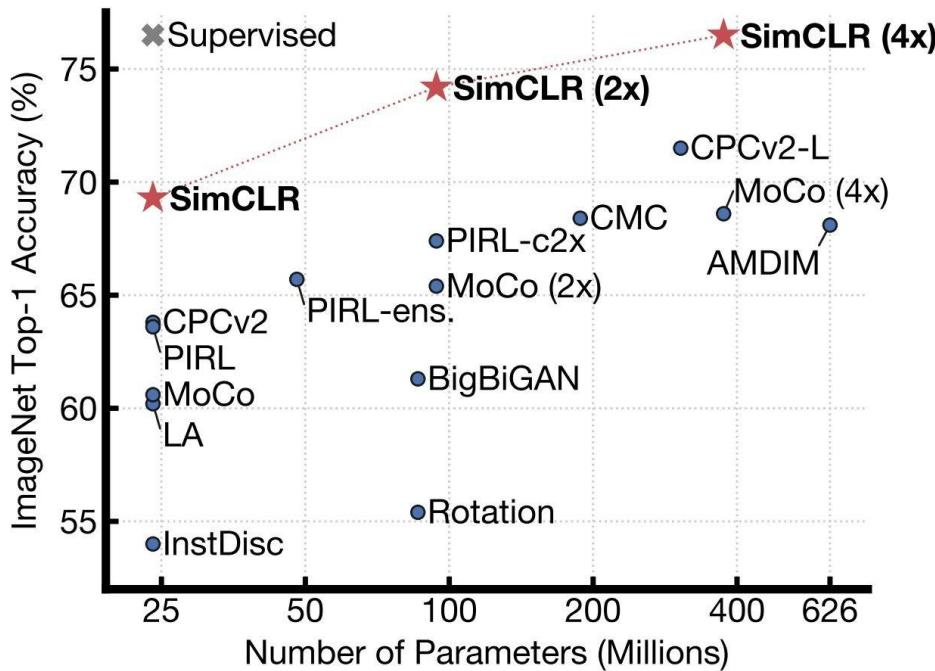
Similarity Calculation of Augmented Images

$$\text{similarity} \left( \begin{matrix} x_i \\ \text{image} \end{matrix}, \begin{matrix} x_j \\ \text{image} \end{matrix} \right) = \text{cosine similarity} \left( \begin{matrix} z_i \\ \text{embedding} \end{matrix}, \begin{matrix} z_j \\ \text{embedding} \end{matrix} \right)$$

$$s_{i,j} = \frac{z_i^T z_j}{(\tau ||z_i|| ||z_j||)}$$

$\tau$  = temperature hyperparameter. It can scale the input and widen the range [-1, 1] of cosine similarity  
 $||z||$  = vector norm

# Training linear classifier on SimCLR features



Train feature encoder on **ImageNet** (entire training set) using SimCLR.

Freeze feature encoder, train a linear classifier on top with labeled data.

Source: [Chen et al., 2020](#)

# Semi-supervised learning on SimCLR features

Method	Architecture	Label fraction		
		1%	10%	Top 5
Supervised baseline	ResNet-50	48.4	80.4	
<i>Methods using other label-propagation:</i>				
Pseudo-label	ResNet-50	51.6	82.4	
VAT+Entropy Min.	ResNet-50	47.0	83.4	
UDA (w. RandAug)	ResNet-50	-	88.5	
FixMatch (w. RandAug)	ResNet-50	-	89.1	
S4L (Rot+VAT+En. M.)	ResNet-50 (4×)	-	91.2	
<i>Methods using representation learning only:</i>				
InstDisc	ResNet-50	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>	

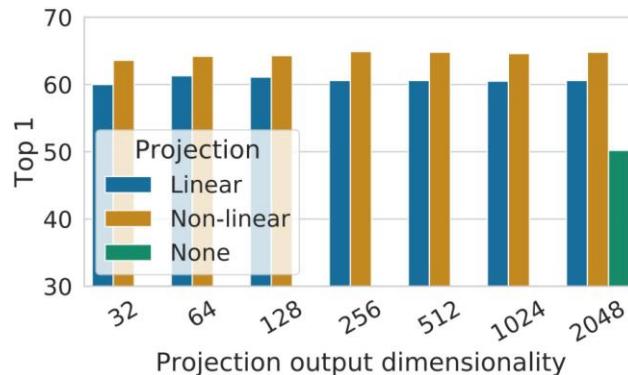
Table 7. ImageNet accuracy of models trained with few labels.

Train feature encoder on **ImageNet** (entire training set) using SimCLR.

**Finetune** the encoder with 1% / 10% of labeled data on ImageNet.

Source: [Chen et al., 2020](#)

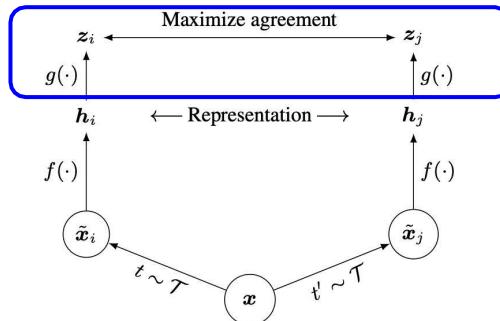
# SimCLR design choices: projection head



Linear / non-linear projection heads improve representation learning.

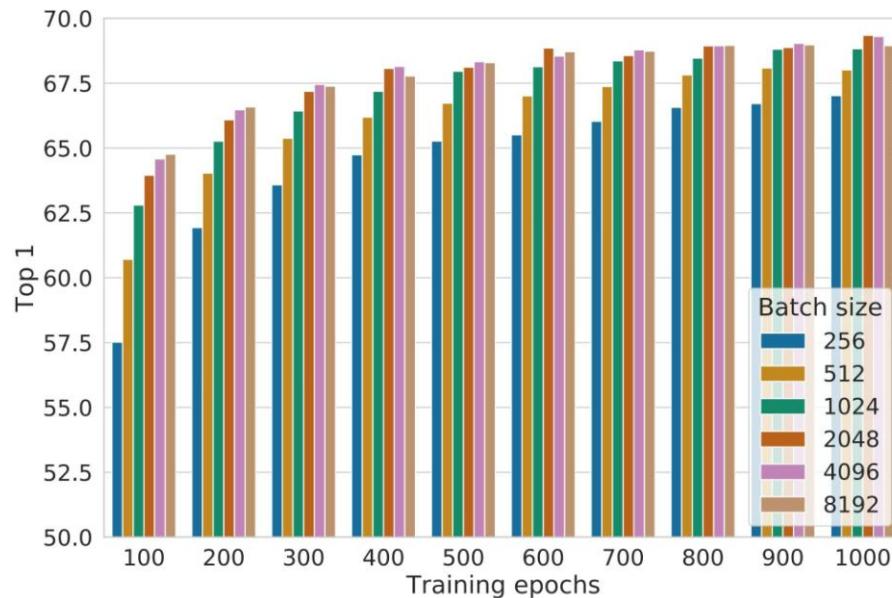
A possible explanation:

- contrastive learning objective may discard useful information for downstream tasks
- representation space  $\mathbf{z}$  is trained to be invariant to data transformation.
- by leveraging the projection head  $\mathbf{g}(\square)$ , more information can be preserved in the  $\mathbf{h}$  representation space



Source: [Chen et al., 2020](#)

# SimCLR design choices: large batch size



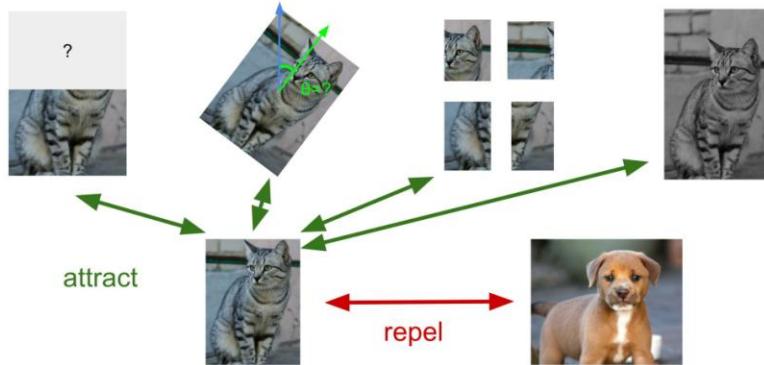
Large training batch size is crucial for SimCLR!

Large batch size causes large memory footprint during backpropagation:  
requires distributed training on TPUs  
(ImageNet experiments)

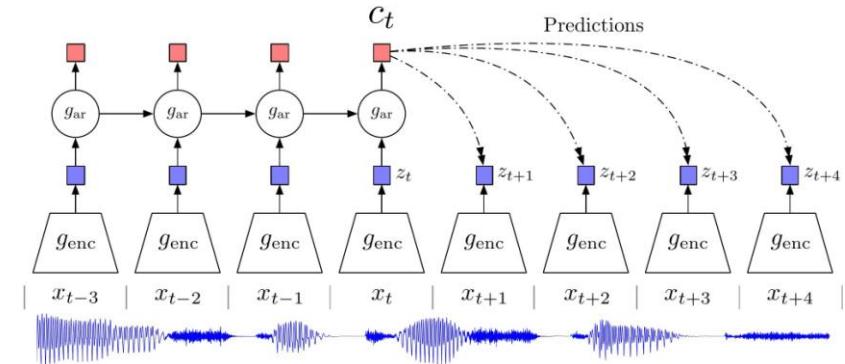
Figure 9. Linear evaluation models (ResNet-50) trained with different batch size and epochs. Each bar is a single run from scratch.<sup>10</sup>

Source: [Chen et al., 2020](#)

# Instance vs. Sequence Contrastive Learning



**Instance-level contrastive learning:**  
contrastive learning based on  
**positive & negative instances.**  
Examples: SimCLR, MoCo



Source: [van den Oord et al., 2018](#)

**Sequence-level contrastive learning:**  
contrastive learning based on  
**sequential / temporal orders.**  
Example: **Contrastive Predictive Coding (CPC)**

# Contrastive Predictive Coding (CPC)

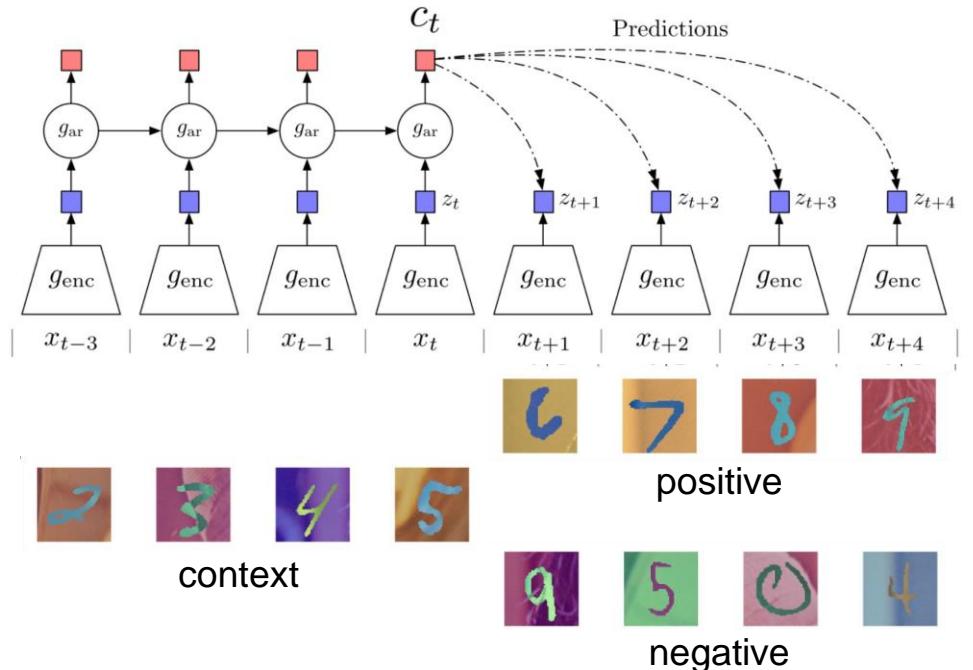
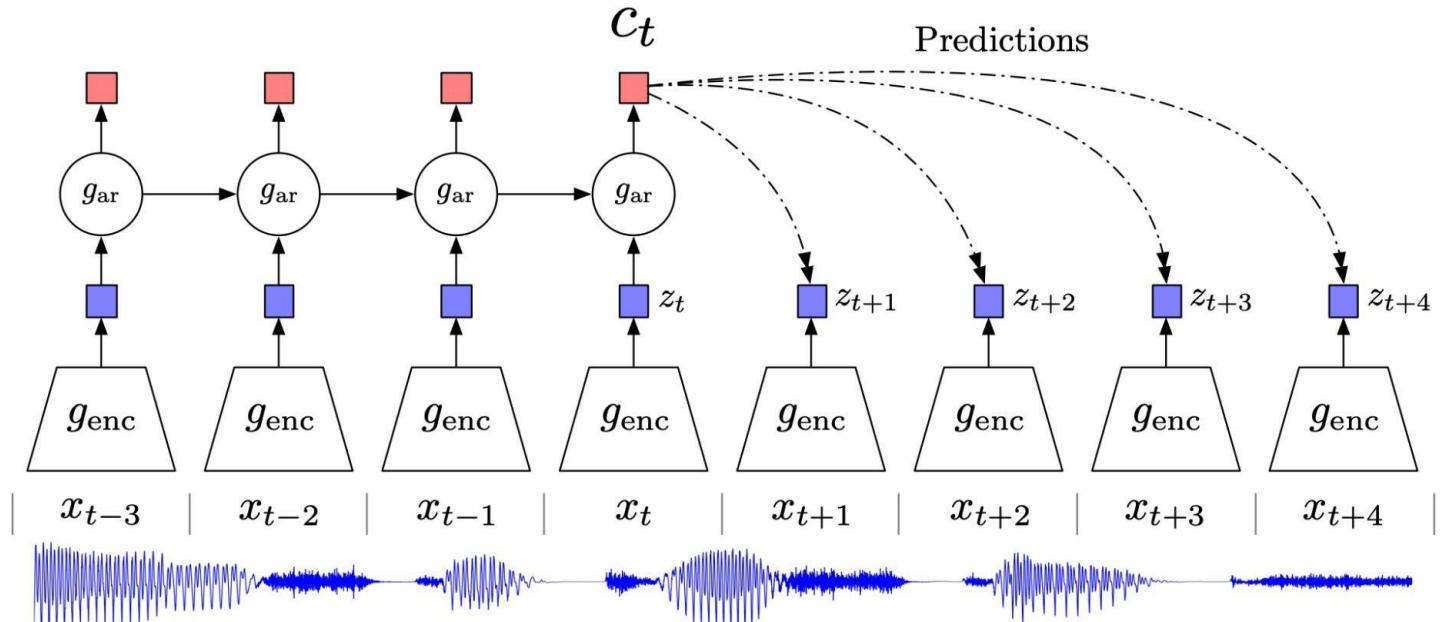


Figure [source](#)

# CPC example: modeling audio sequences



Source: [van den Oord et al., 2018](#),

# CPC example: modeling audio sequences

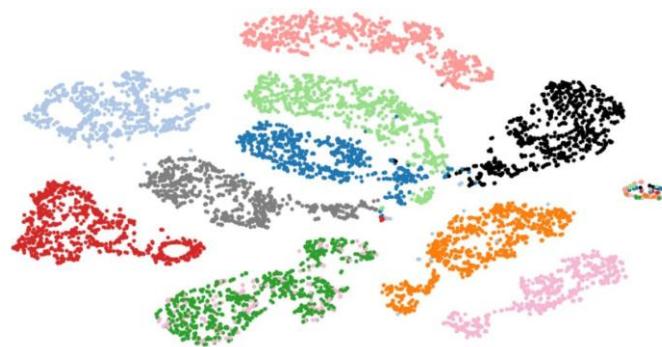


Figure 2: t-SNE visualization of audio (speech) representations for a subset of 10 speakers (out of 251). Every color represents a different speaker.

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

Linear classification on trained representations (LibriSpeech dataset)

Source: [van den Oord et al., 2018](#),

# Summary: Contrastive Representation Learning

A general formulation for contrastive learning:

$$\text{score}(f(x), f(x^+)) \gg \text{score}(f(x), f(x^-))$$

InfoNCE loss: N-way classification among positive and negative samples

$$L = -\mathbb{E}_X \left[ \log \frac{\exp(s(f(x), f(x^+)))}{\exp(s(f(x), f(x^+))) + \sum_{j=1}^{N-1} \exp(s(f(x), f(x_j^-)))} \right]$$

Commonly known as the InfoNCE loss ([van den Oord et al., 2018](#))

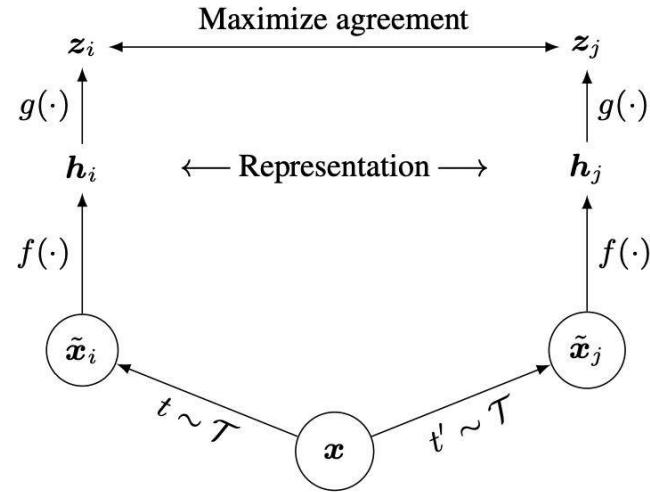
A *lower bound* on the mutual information between  $f(x)$  and  $f(x^+)$

$$MI[f(x), f(x^+)] - \log(N) \geq -L$$

# Summary: Contrastive Representation Learning

**SimCLR**: a simple framework for contrastive representation learning

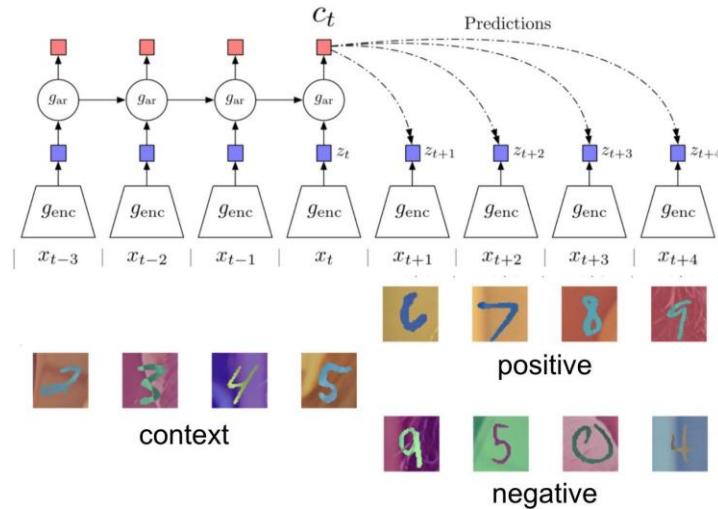
- **Key ideas**: non-linear projection head to allow flexible representation learning
- Simple to implement, effective in learning visual representation
- Requires large training batch size to be effective; large memory footprint



# Summary: Contrastive Representation Learning

**CPC:** sequence-level contrastive learning

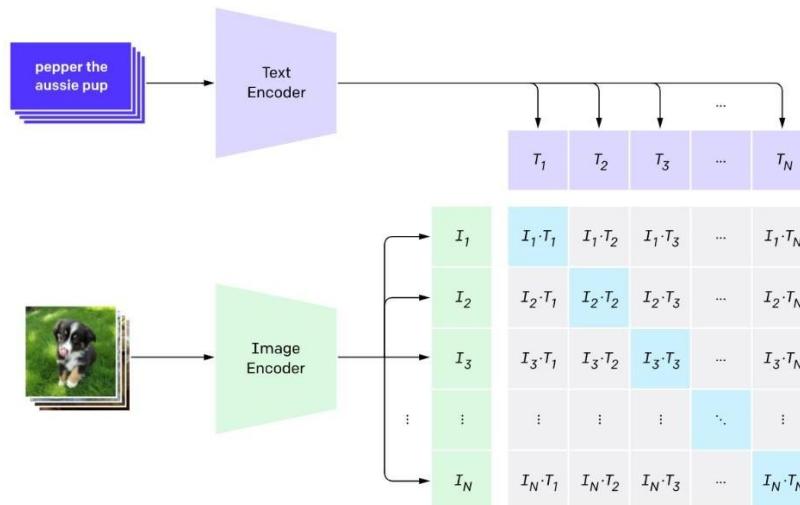
- Contrast “right” sequence with “wrong” sequence.
- InfoNCE loss with a time-dependent score function.
- Can be applied to a variety of learning problems, but not as effective in learning image representations compared to instance-level methods.



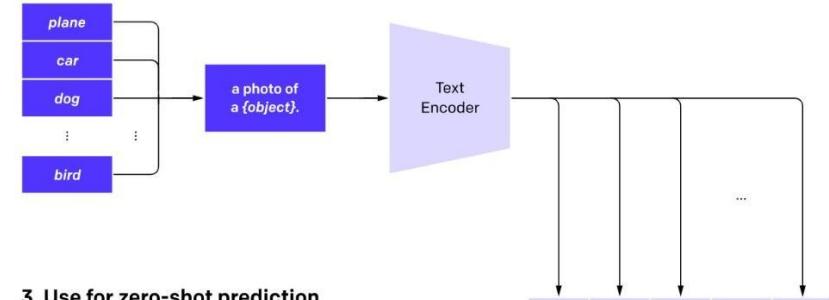
# Other examples

Contrastive learning between image and natural language sentences

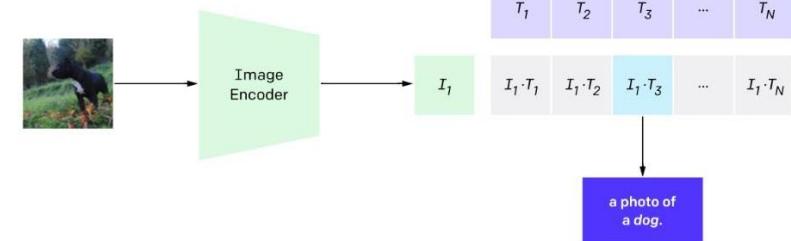
## 1. Contrastive pre-training



## 2. Create dataset classifier from label text



## 3. Use for zero-shot prediction

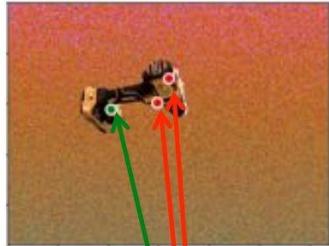


CLIP (*Contrastive Language–Image Pre-training*) Radford *et al.*, 2021

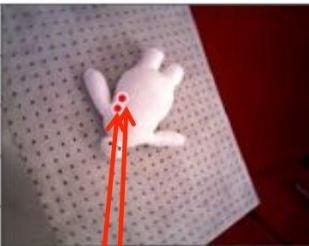
# Other examples

Contrastive learning on pixel-wise feature descriptors

(c) Background Randomization



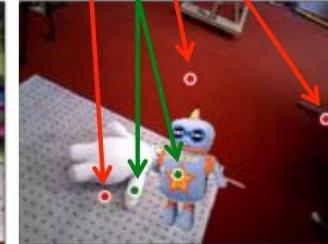
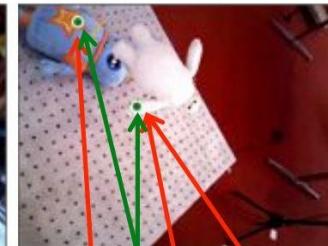
(d) Cross Object Loss



(e) Direct Multi Object

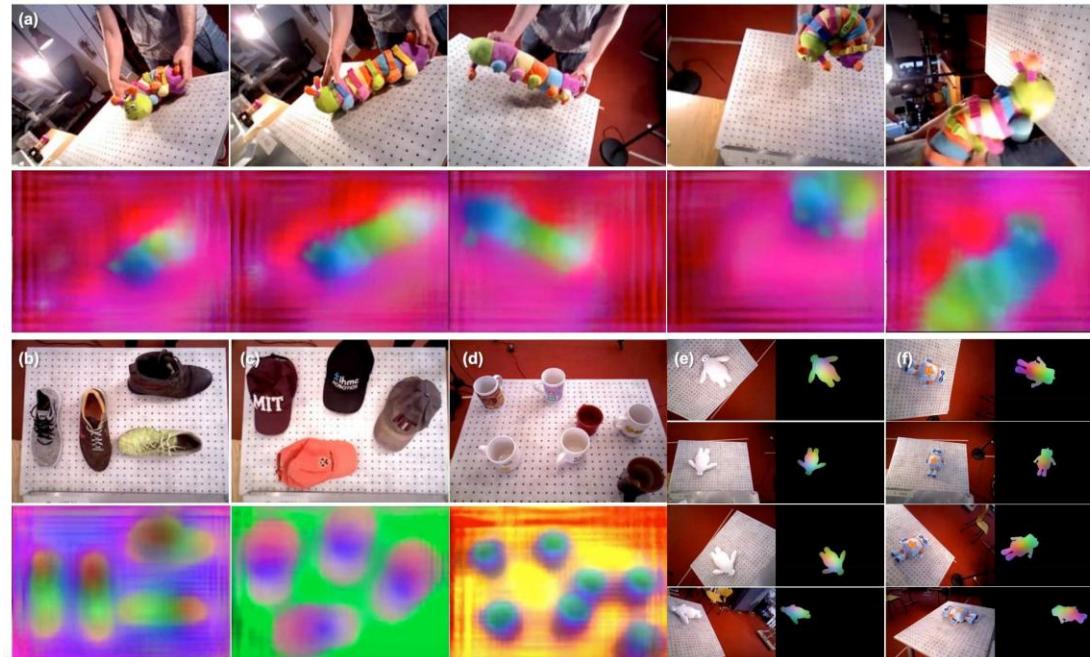


(f) Synthetic Multi Object



Dense Object Net, Florence et al., 2018

# Other examples



Dense Object Net, Florence et al., 2018

# Self-Supervised Learning

General idea: pretend there is a part of the data you don't know and train the neural network to predict that.

Y. LeCun

## Self-Supervised 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.**

Time →

← Past      Present      Future →