

# Self-supervised learning in computer vision

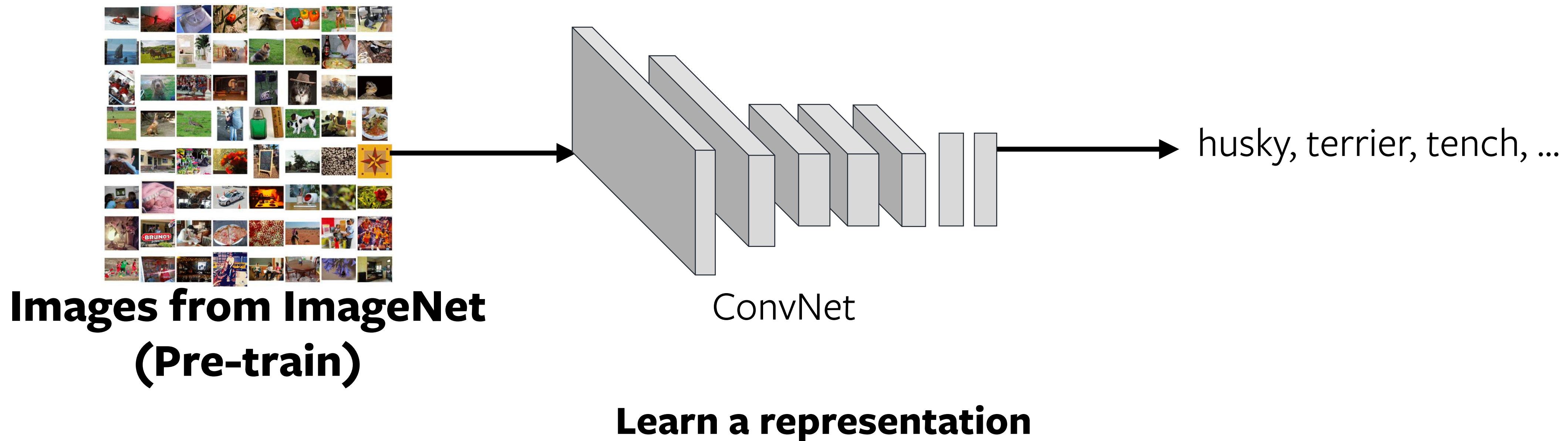
Ishan Misra

Facebook AI Research

With slides from Andrew Zisserman, Carl Doersch

# Success story of supervision: Pre-training

- Features from networks pre-trained on ImageNet can be used for a variety of different downstream tasks



# Success story of supervision: Recipe for good solutions

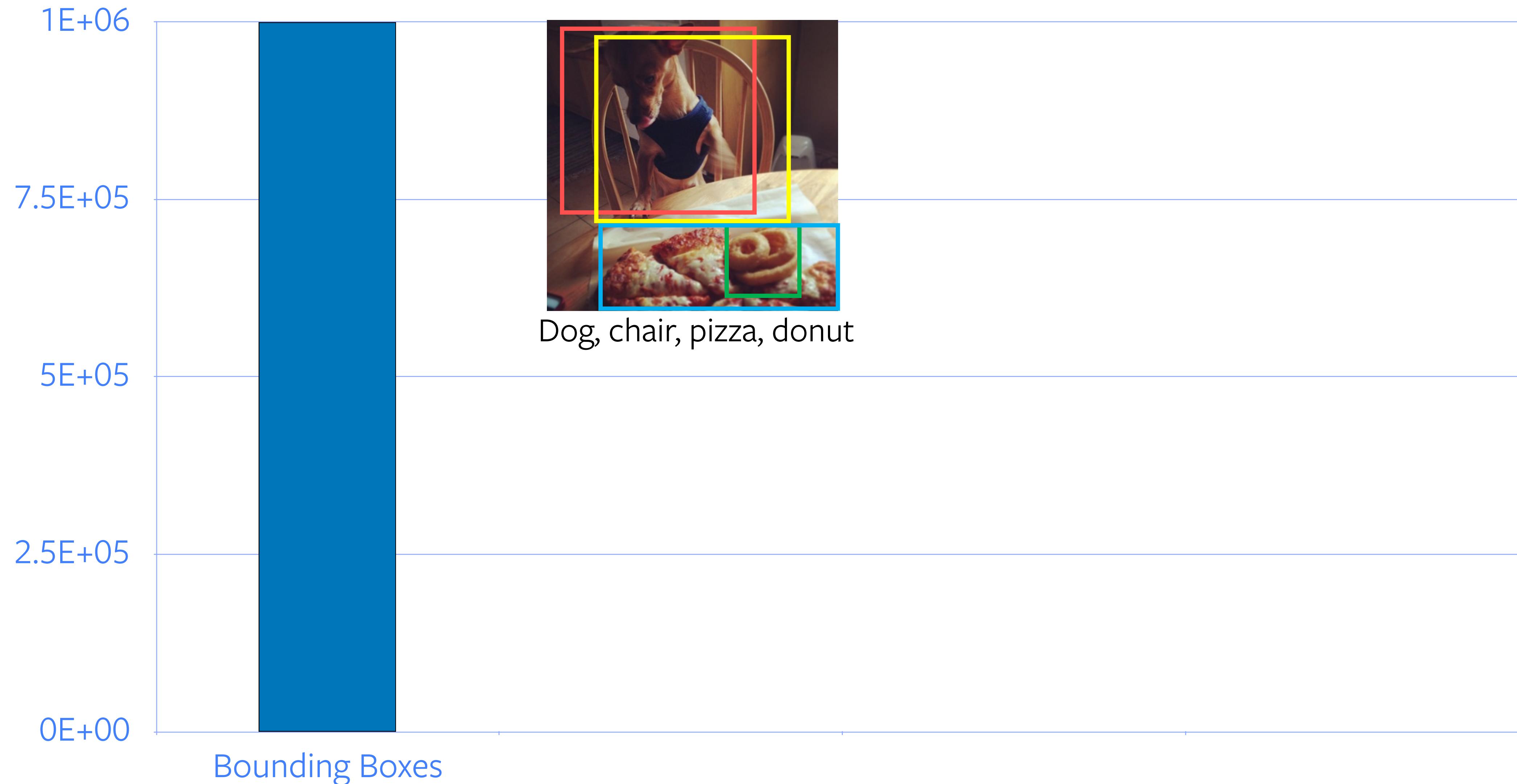
- Pre-train on a large supervised dataset.
- Collect a dataset of “supervised” images
- Train a ConvNet

# The promise of "alternative" supervision

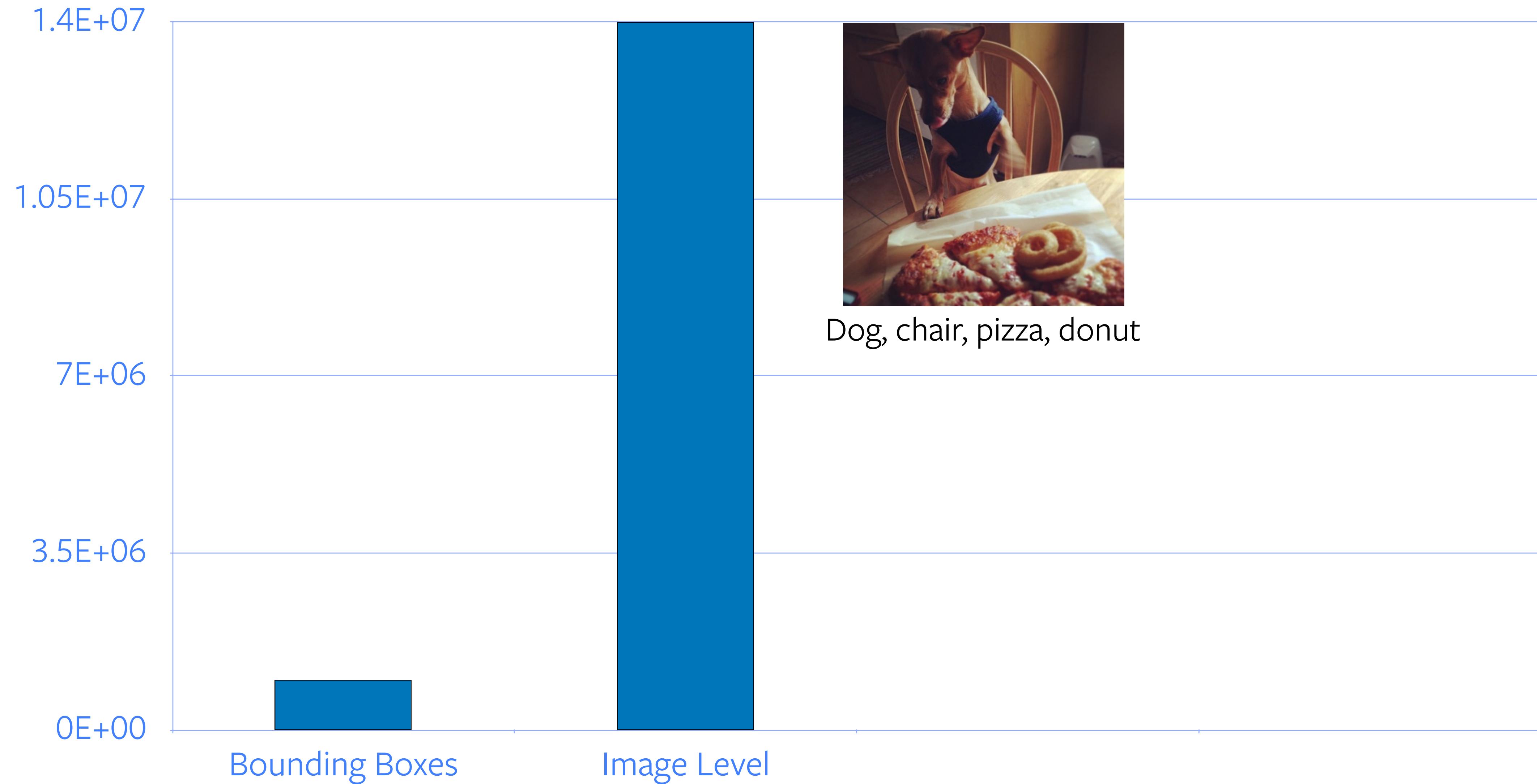
- Getting "real" labels is difficult and expensive
  - ImageNet with 14M images took 22 human years.
- Obtain labels using a "semi-automatic" process
  - Hashtags
  - GPS locations
  - Using the data itself: "self"-supervised

# Can we get labels for all data?

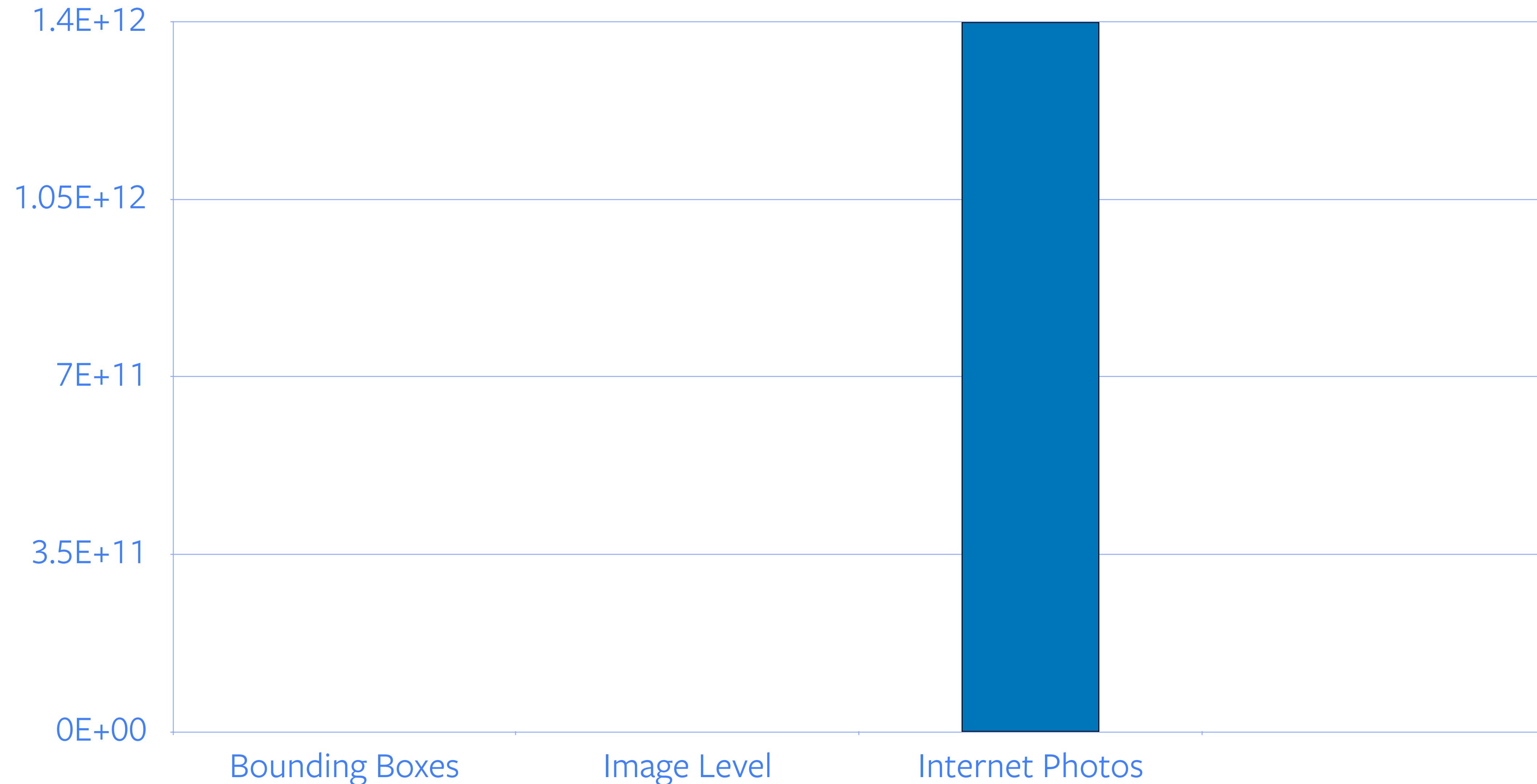
# Can we get labels for all data?



# Can we get labels for all data?



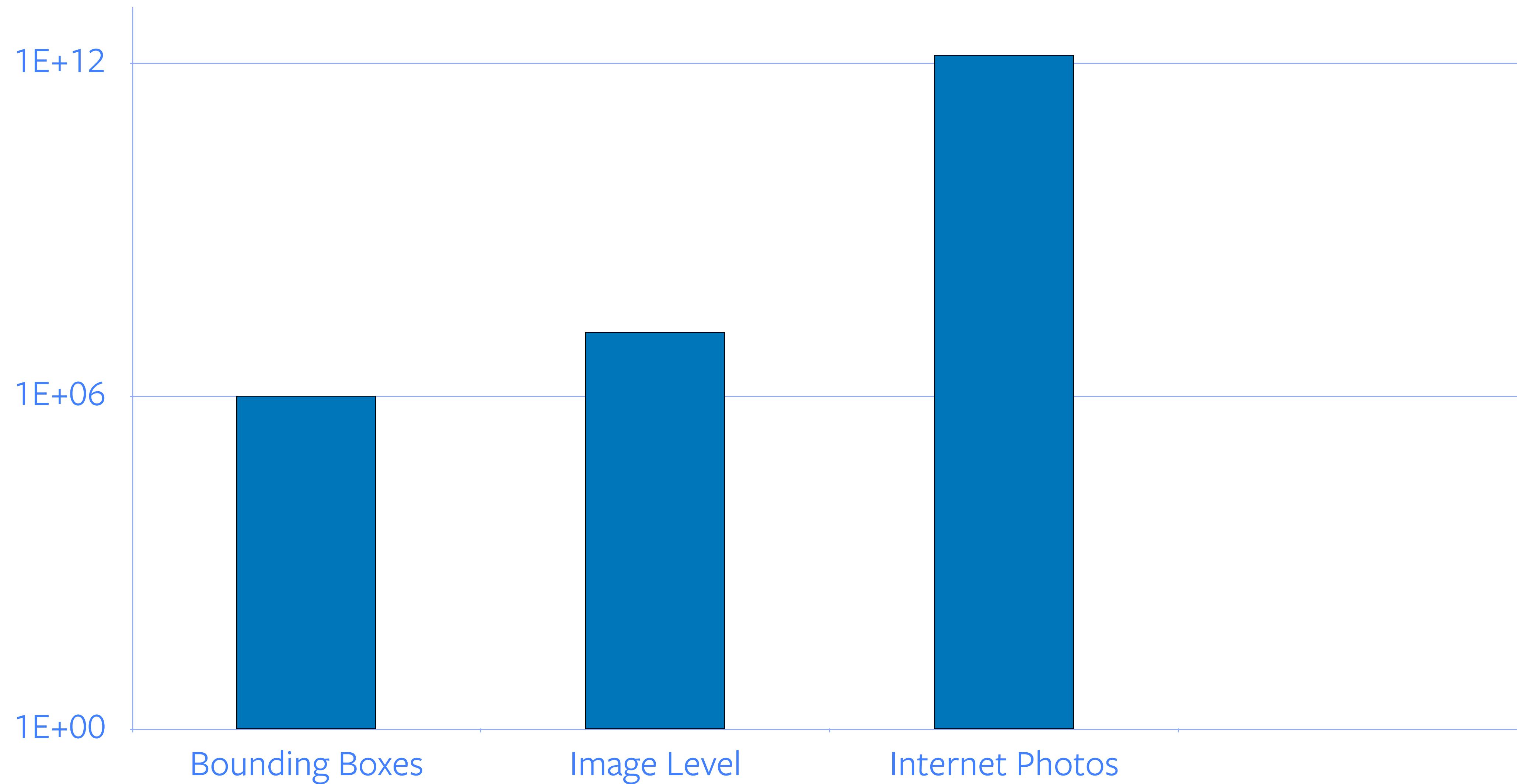
# Can we get labels for all data?



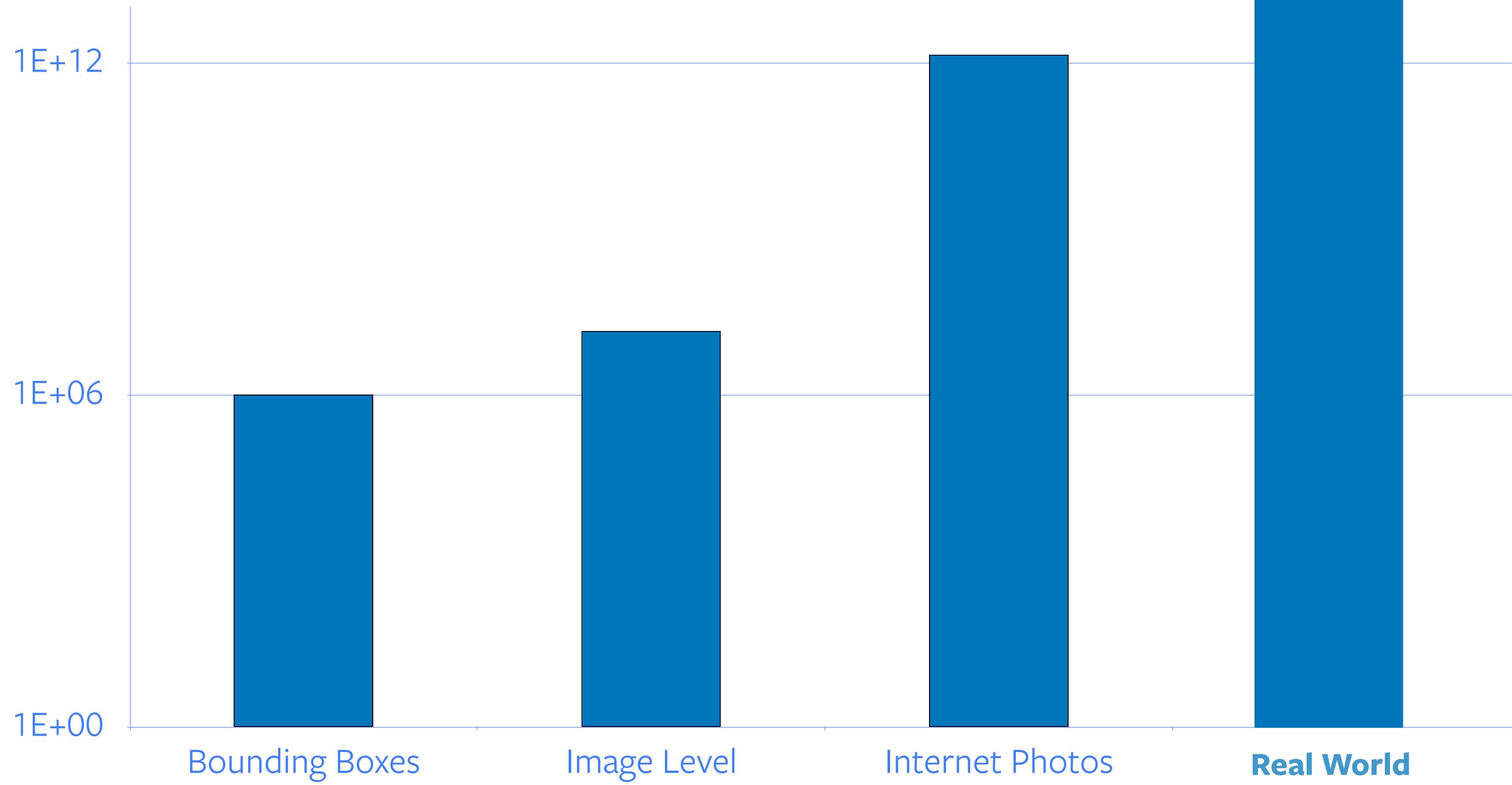
[forbes.com](https://www.forbes.com/sites/bernardmarr/2018/05/21/how-much-data-do-we-create-every-day-the-mind-blowing-stats-everyone-should-read/)

<https://www.forbes.com/sites/bernardmarr/2018/05/21/how-much-data-do-we-create-every-day-the-mind-blowing-stats-everyone-should-read/>

# Can we get labels for all data?



# Can we get labels for all data?

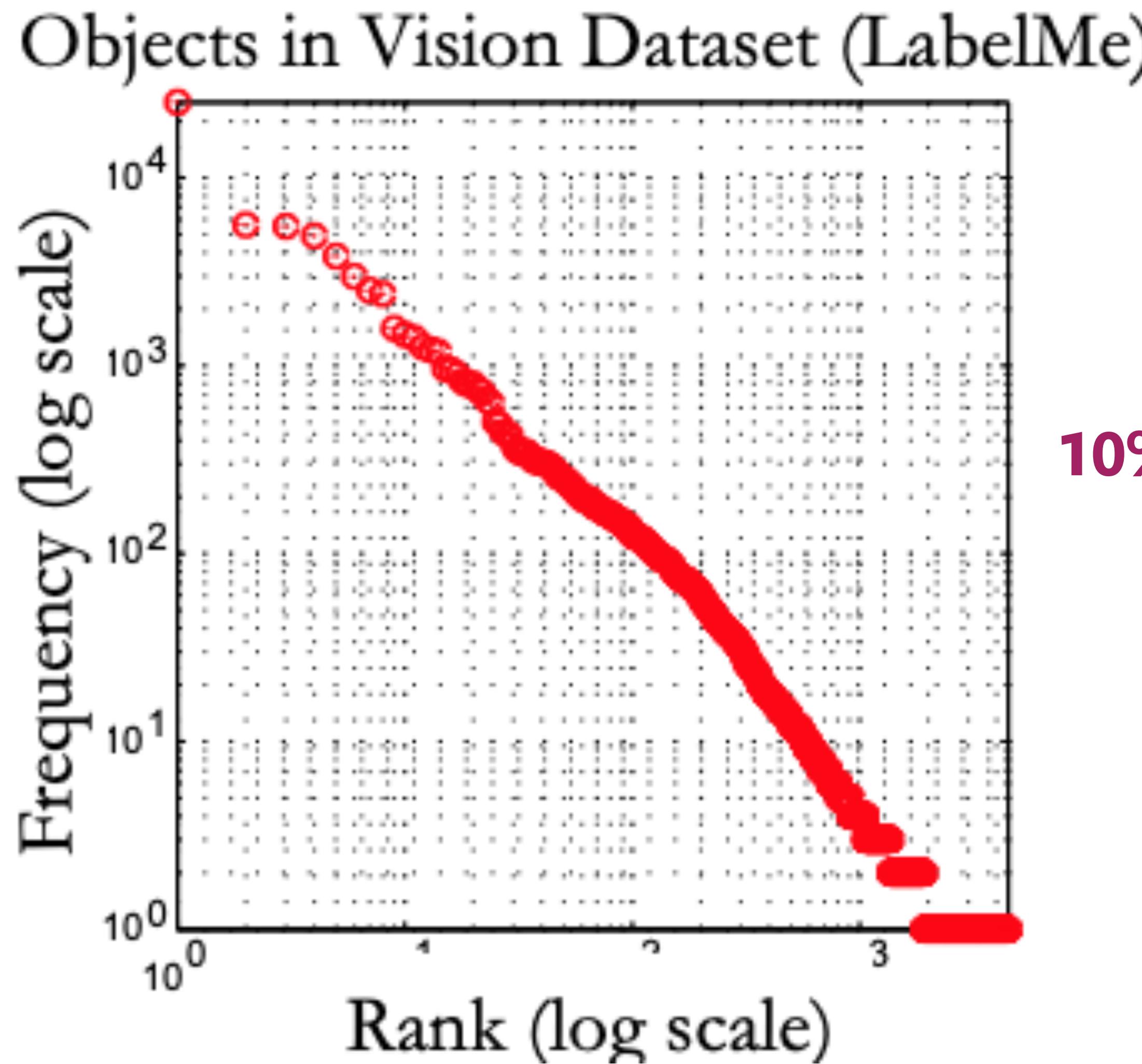


**ImageNet (14 million images) needed 22 human years to label**

# Can we get labels for all data?

- What about complex concepts?
  - Video?
- Labelling cannot scale to the size of the data we generate

# Rare concepts?



**10% of the classes account  
for 93% of the data**

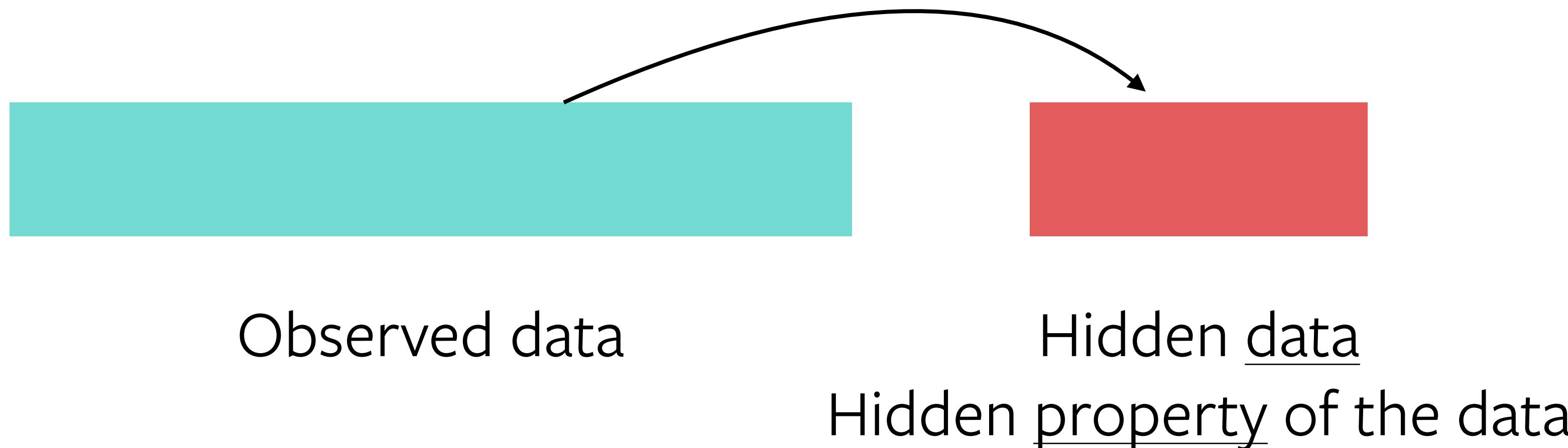
# Different Domains?



ImageNet pre-training may not work

# What is “self” supervision?

- Obtain “labels” from the data itself by using a “semi-automatic” process
- Predict part of the data from other parts

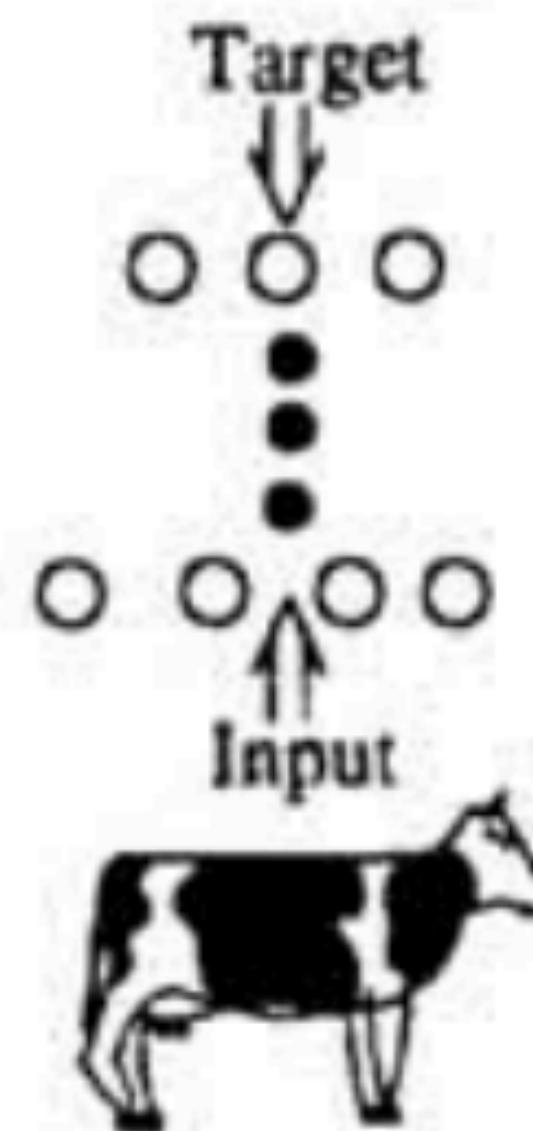


# What is "self" supervision?

## Supervised

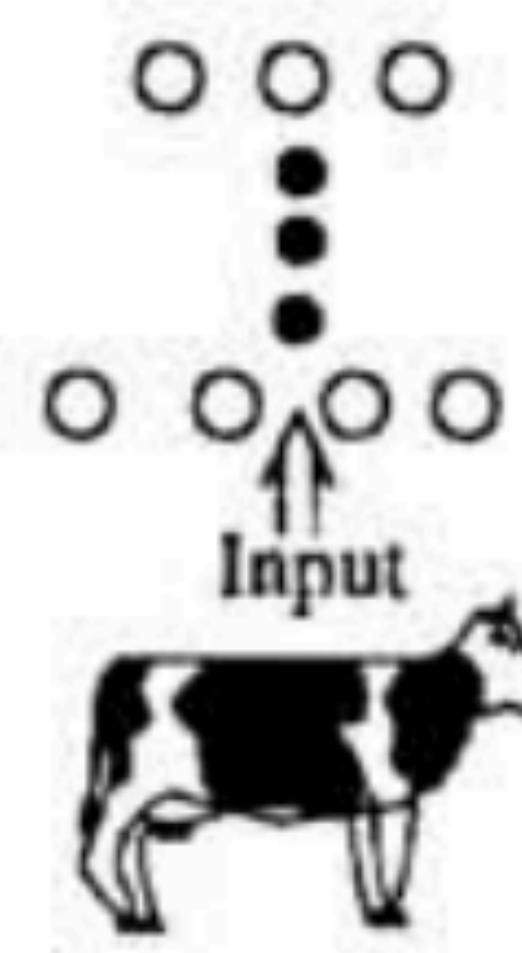
- implausible label

"COW"



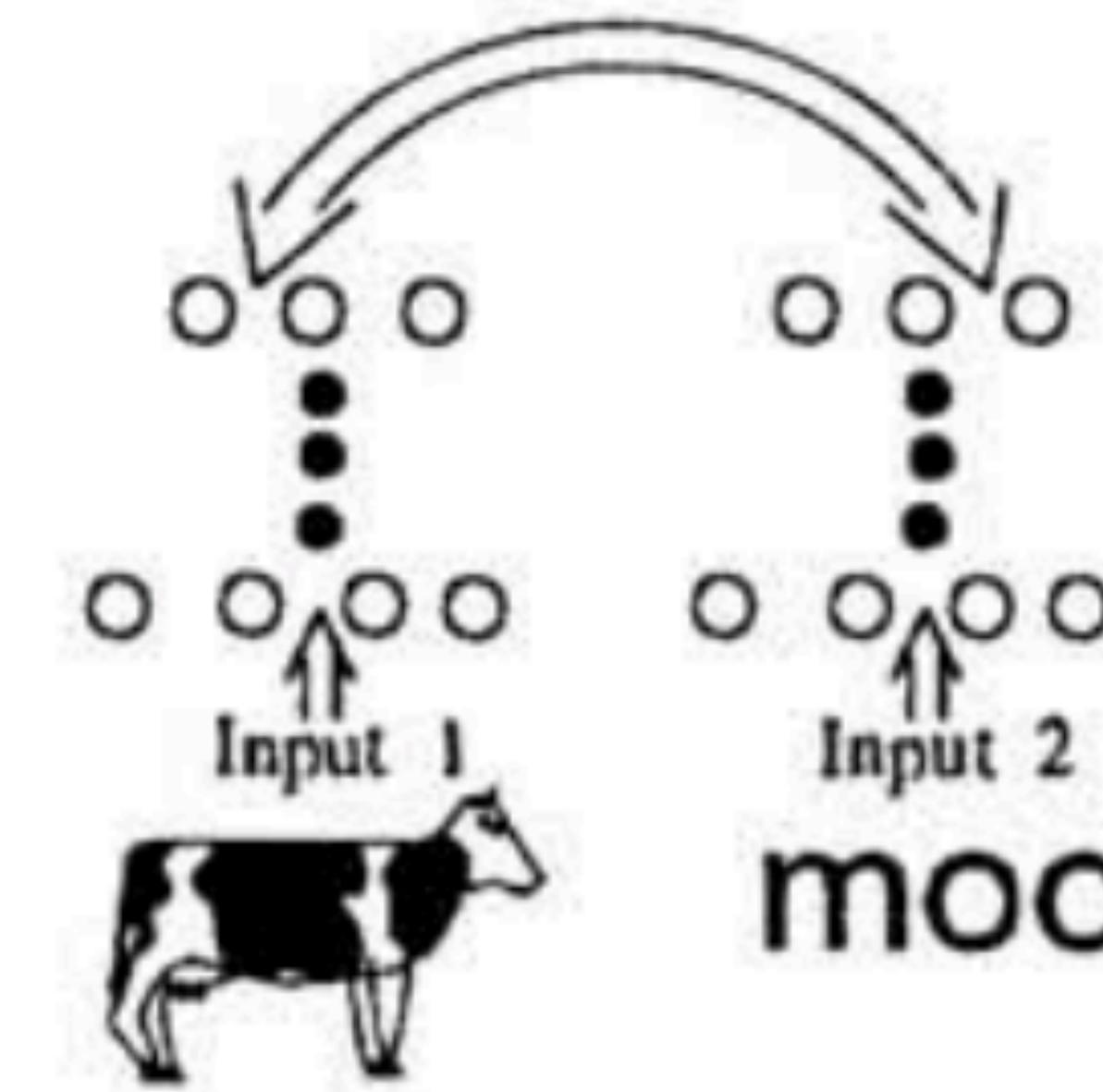
## Unsupervised

- limited power



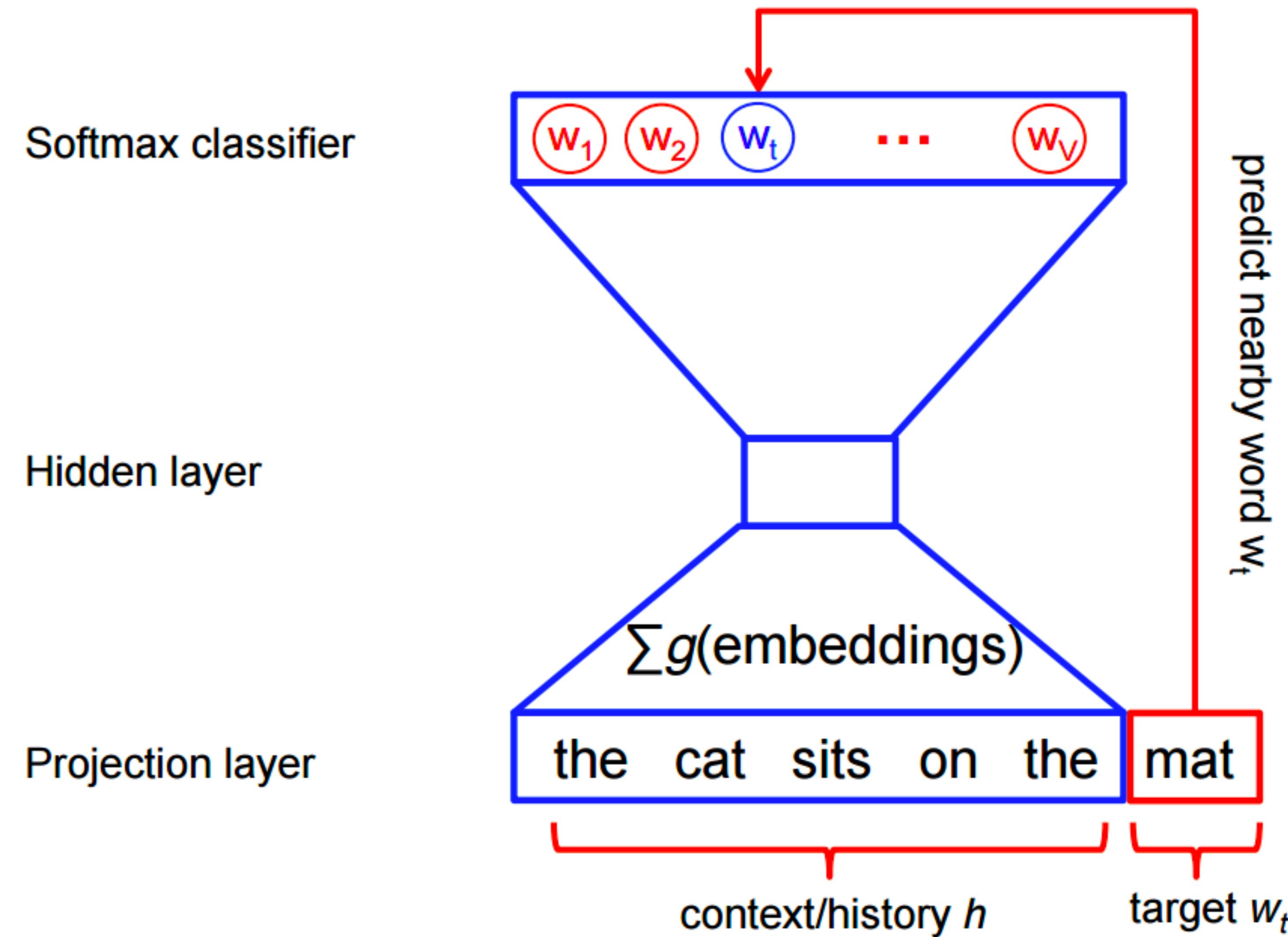
## Self-Supervised

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



# Word2vec

- Fill in the blanks



# Success of self-supervised learning in NLP

- Fill in the blanks is a powerful signal to learn representations
- Sentence/Word representations: BERT - Devlin et al., 2018

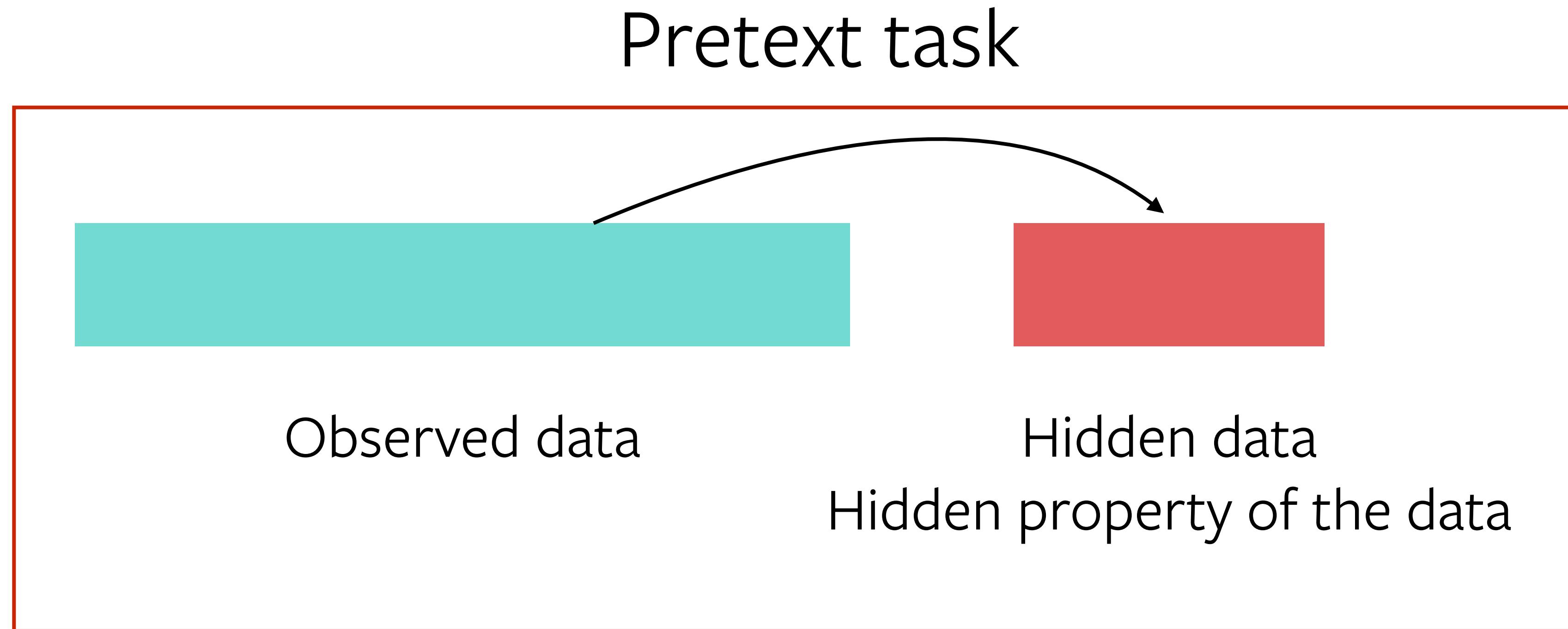
# Why self supervision?

- Helps us learn using observations and interactions
- Does not require exhaustive annotation of concepts
- Leverage multiple modalities or structure in the domain

In the context of  
Computer Vision

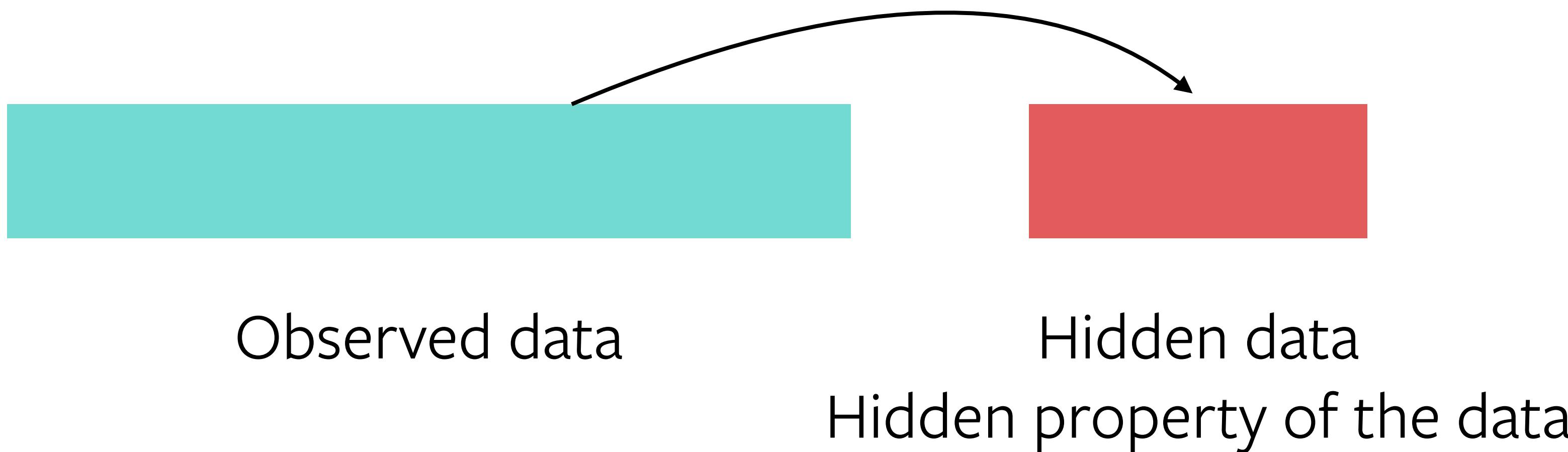
# Pretext task

- Self-supervised task used for learning representations
- Often, not the “real” task (like image classification) we care about



# Pretext task

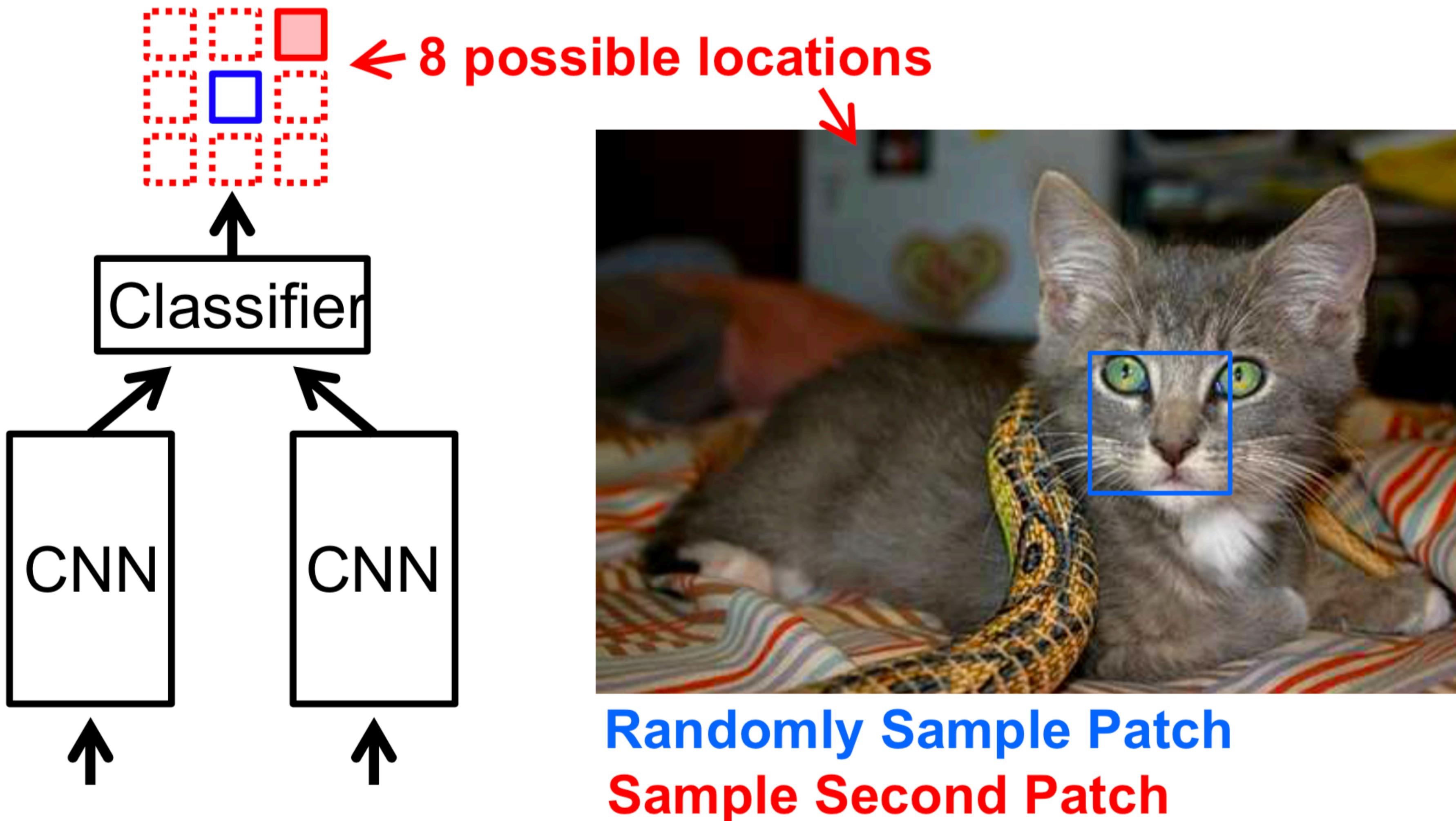
- Using images
- Using video
- Using video and sound



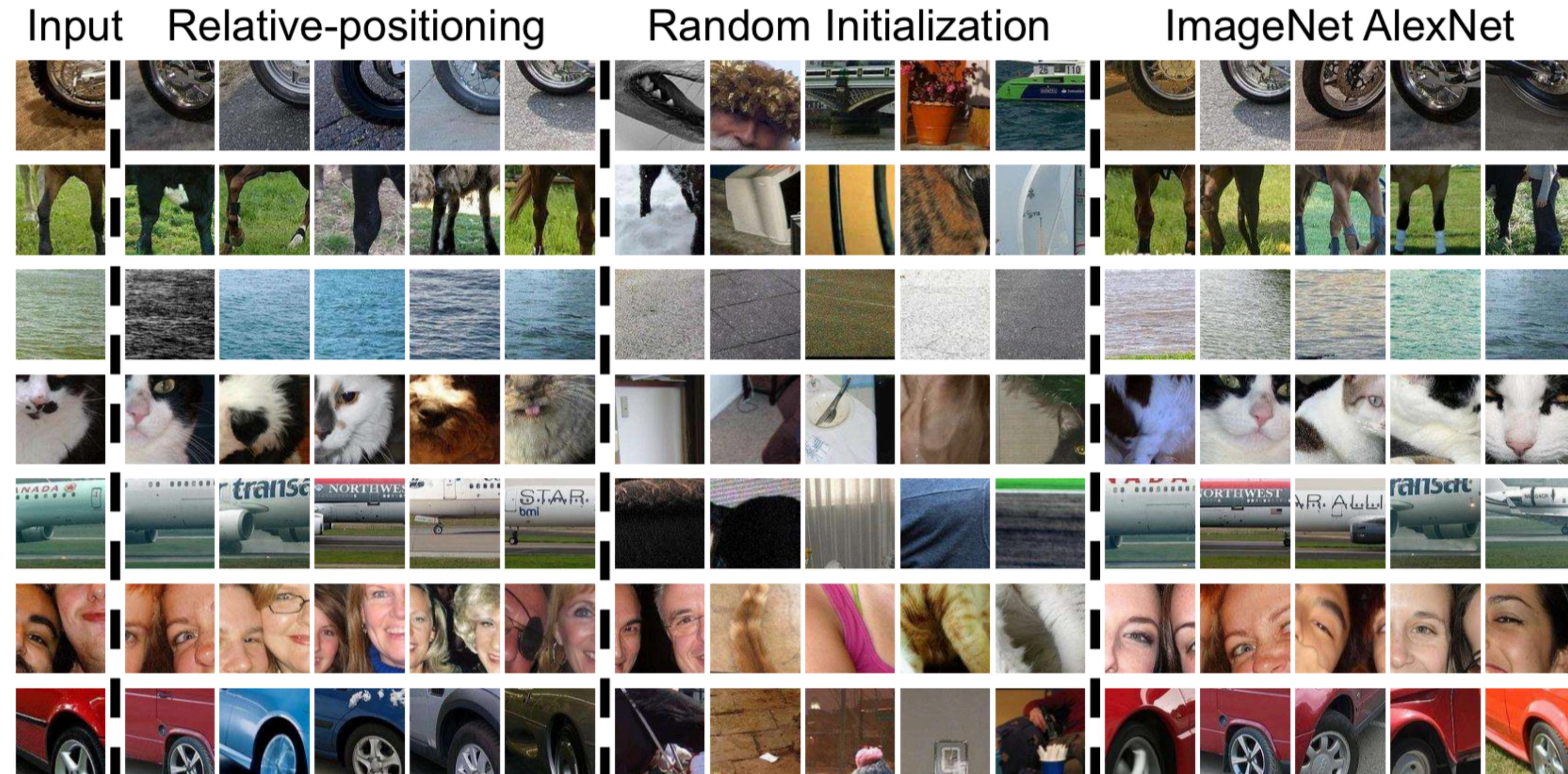
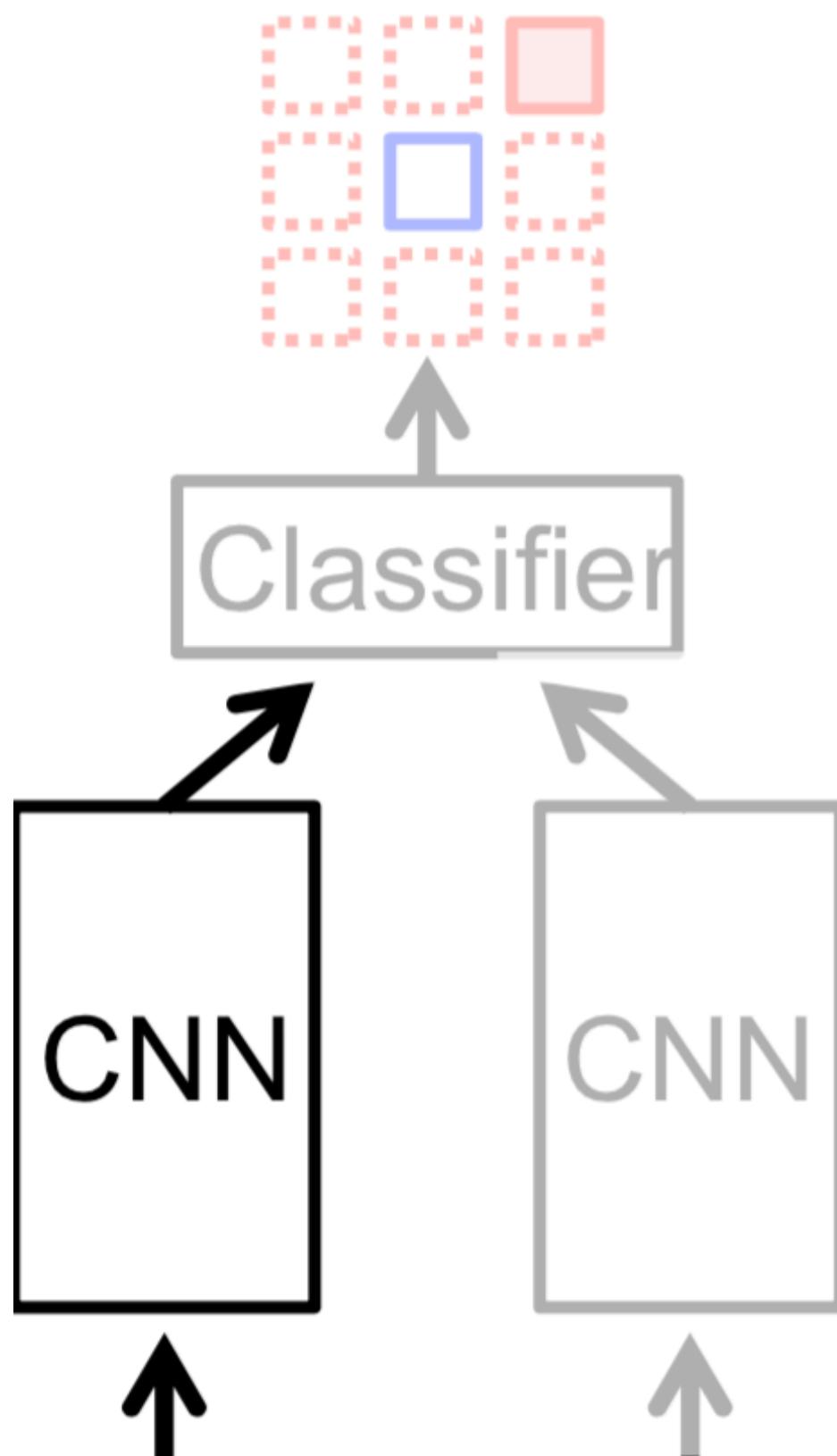
# Pretext task

- Using images
- Using video
- Using video and sound

# Relative Position of patches



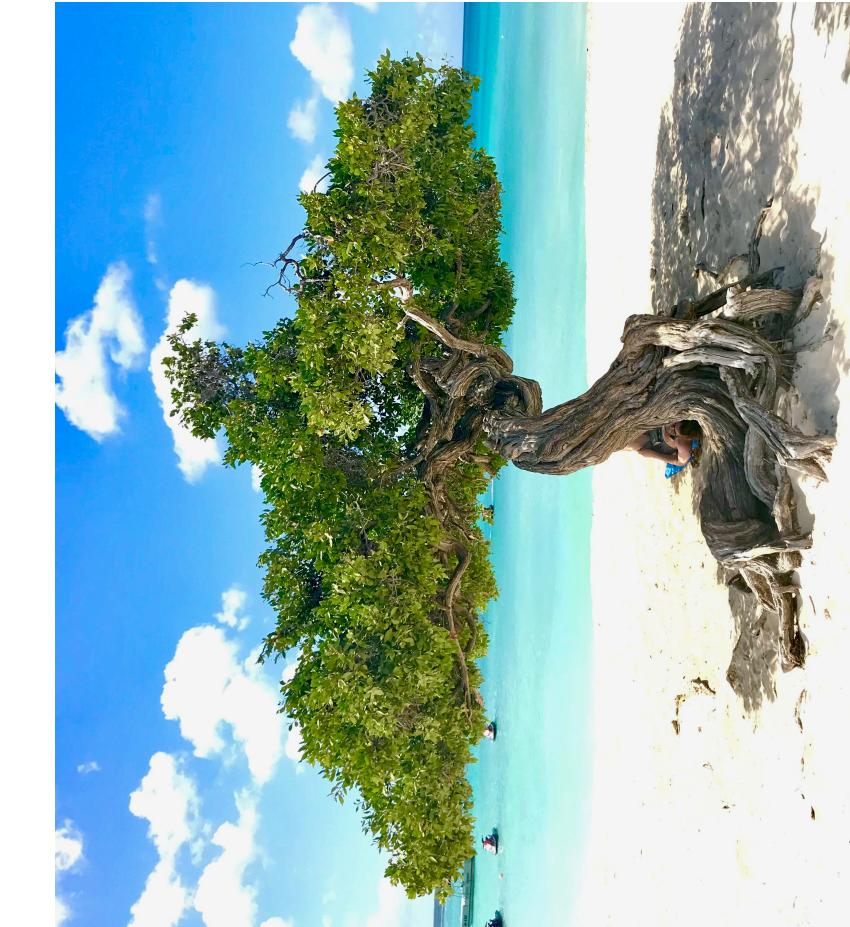
# Relative Position: Nearest Neighbors in features



# Predicting Rotations



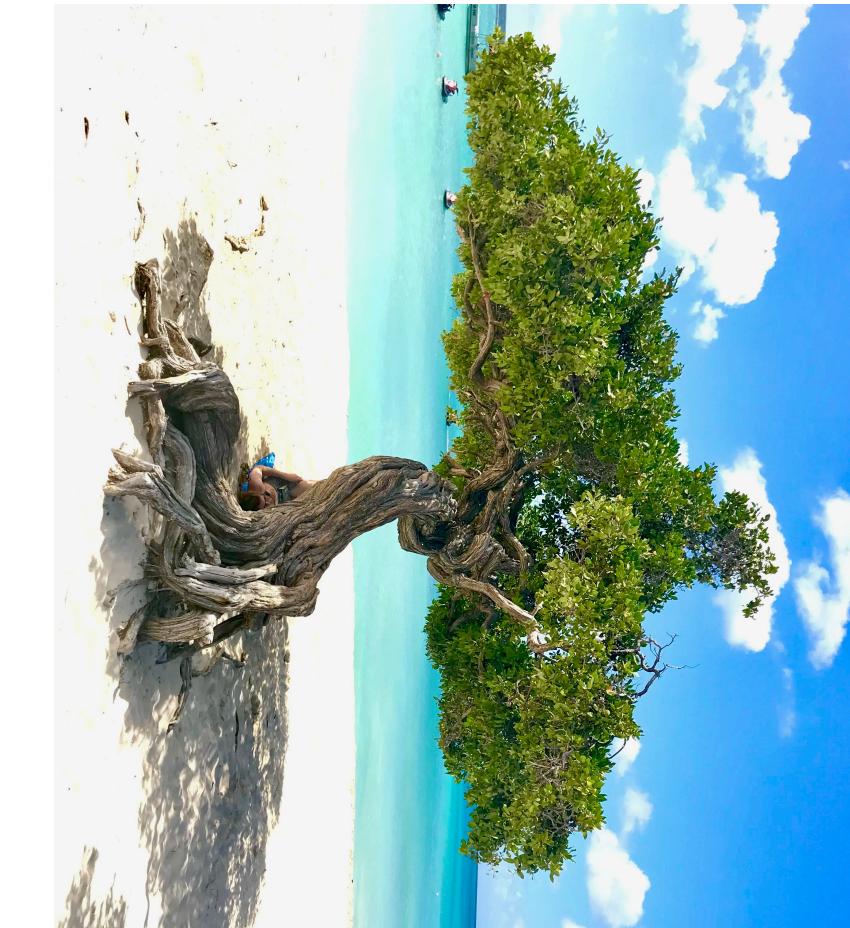
$\rightarrow 0^\circ$



$\rightarrow 90^\circ$

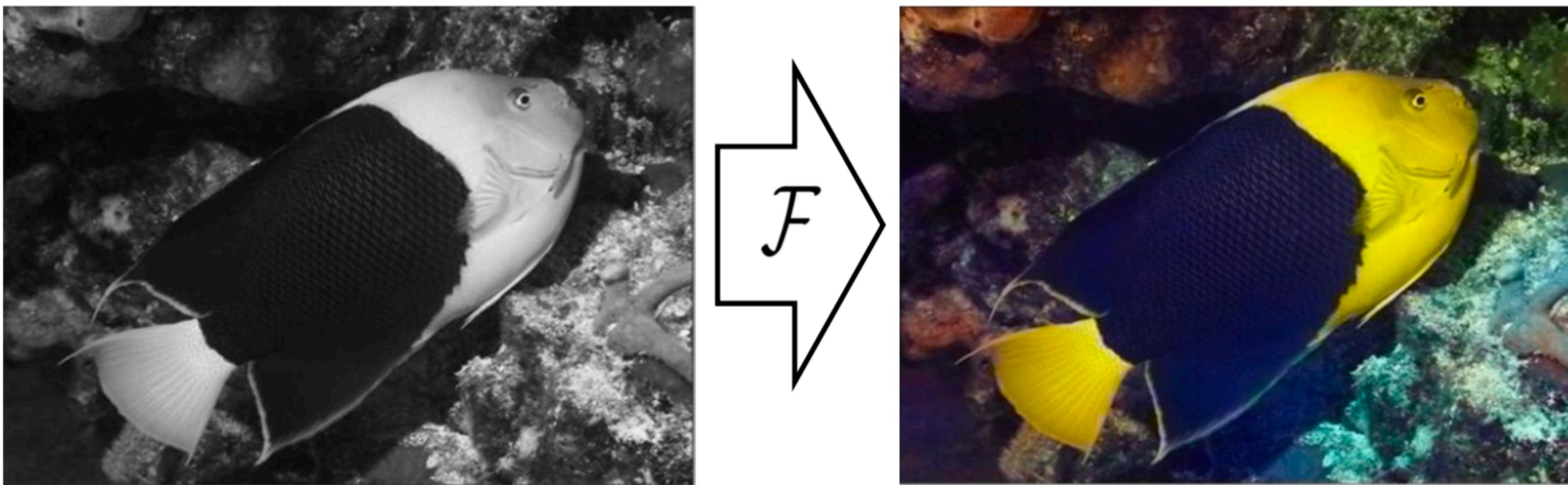


$\rightarrow 180^\circ$



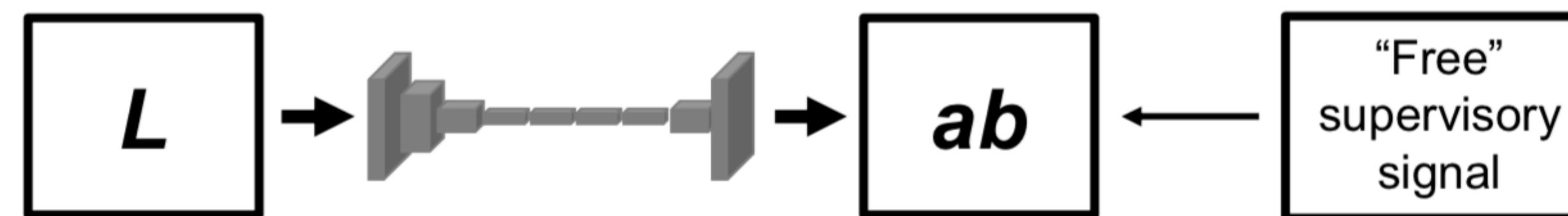
$\rightarrow 270^\circ$

# Colorization



Grayscale image:  $L$  channel

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



Concatenate (L,ab)

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

# Fill in the blanks



Pathak et al., 2016, Context auto encoders

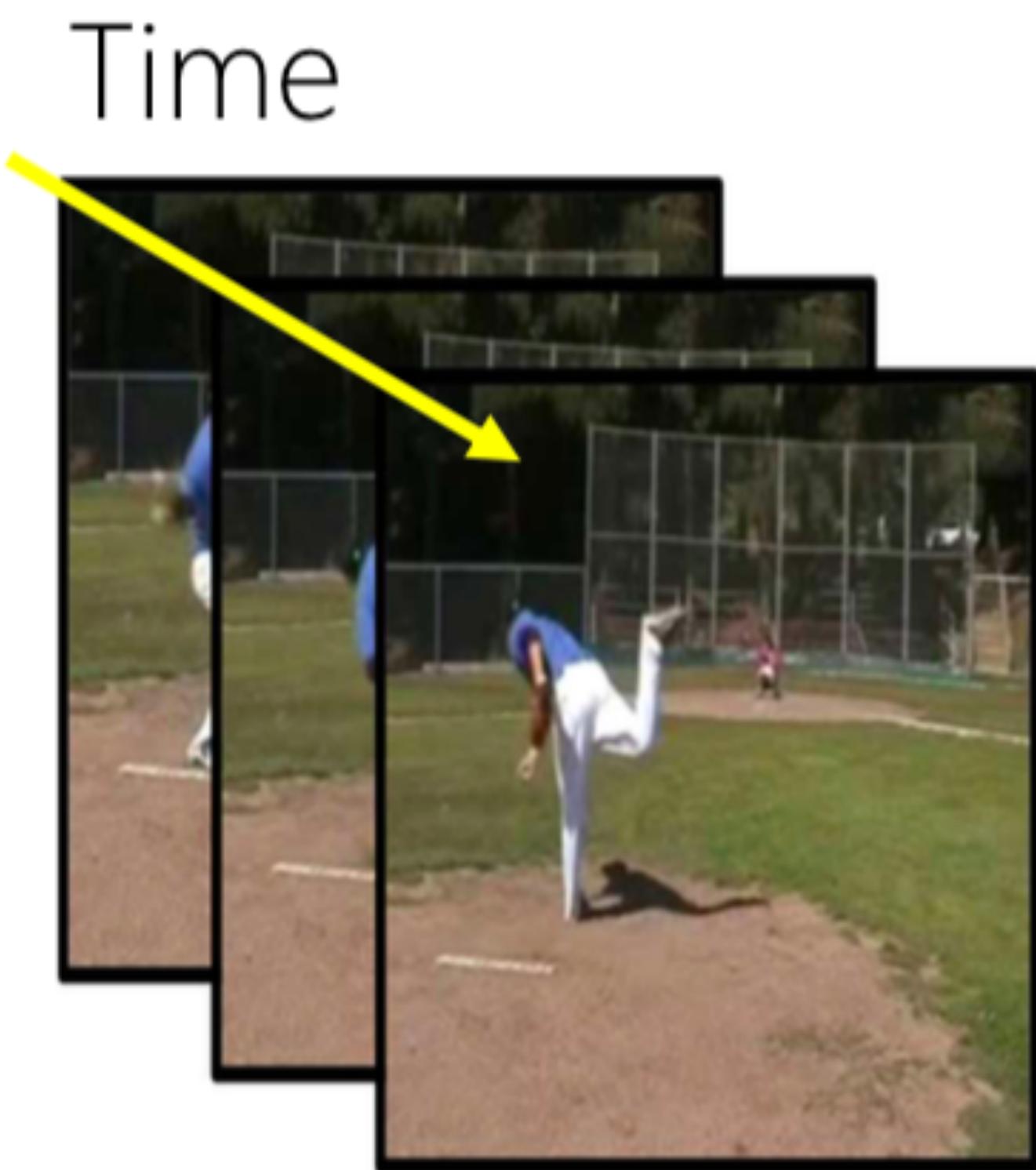
# Self-supervision in computer vision

- Using images
- Using video
- Using video and sound

# Video

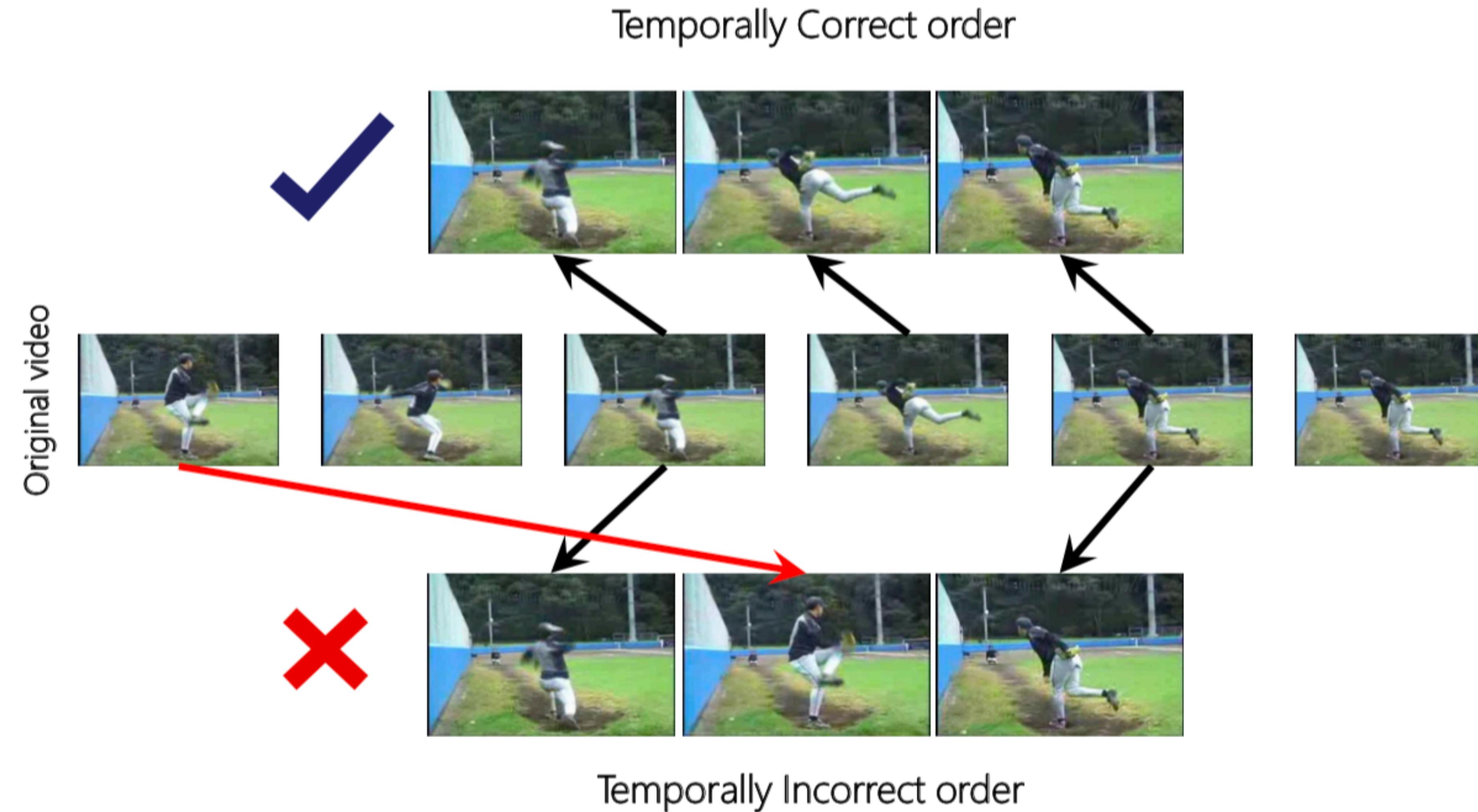
- Video is a “sequence” of frames
- How to get “self-supervision”?

- Predict order of frames
- Fill in the blanks
- Track objects and predict their position

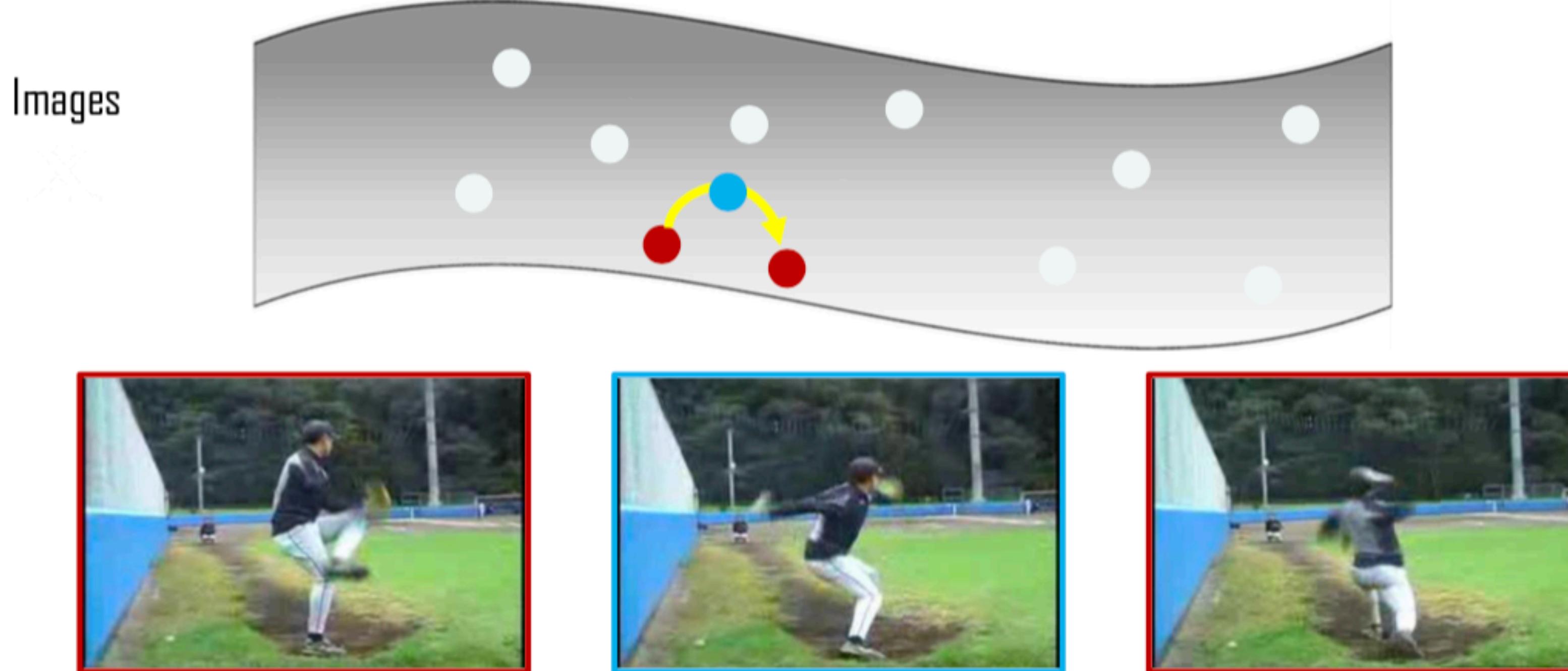


“Sequence” of data

# Shuffle & Learn

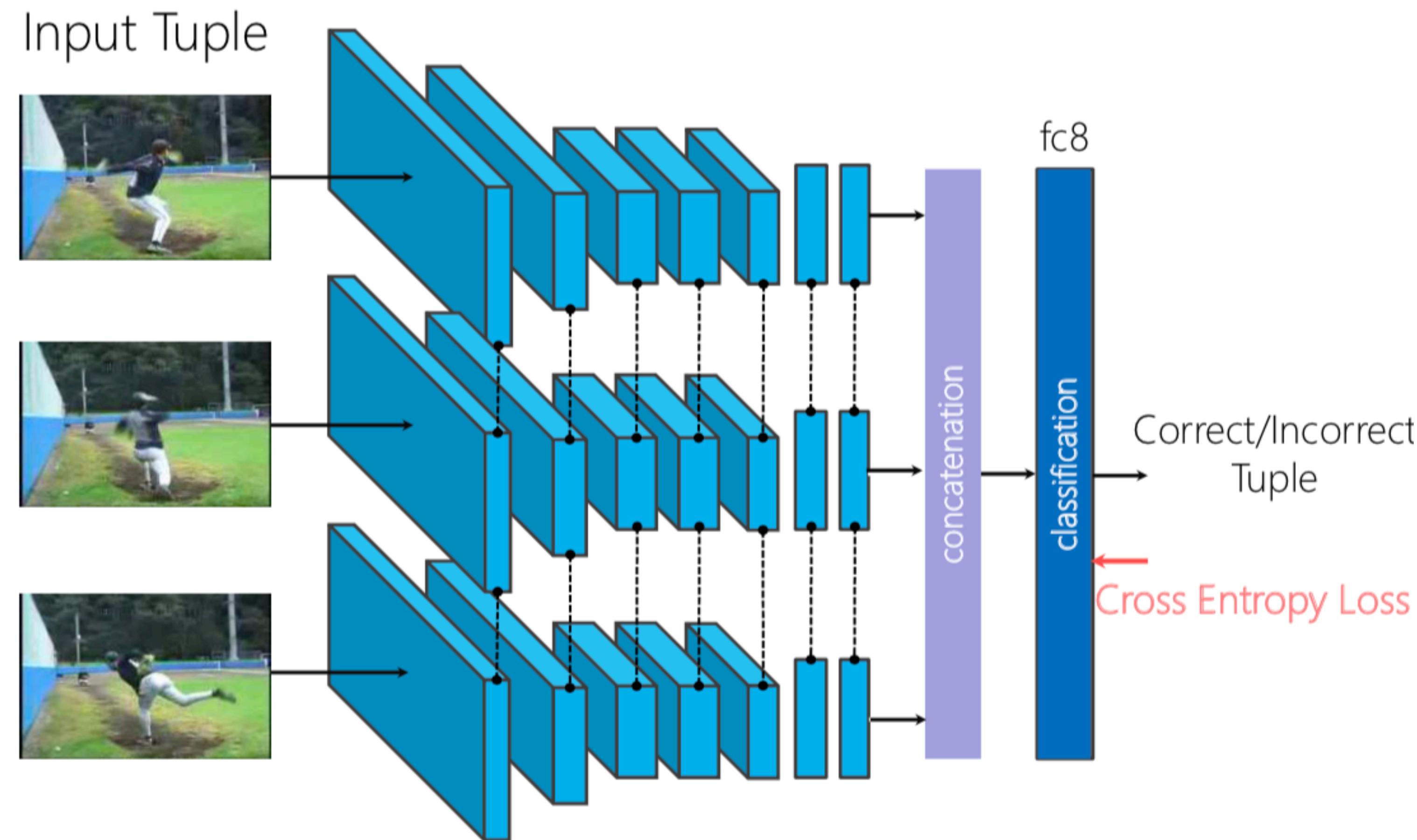


# Shuffle & Learn



Given a start and an end, can this point lie in between?

# Shuffle & Learn



# Nearest Neighbors of Query Frame (fc7 features)

Query



ImageNet



Shuffle & Learn

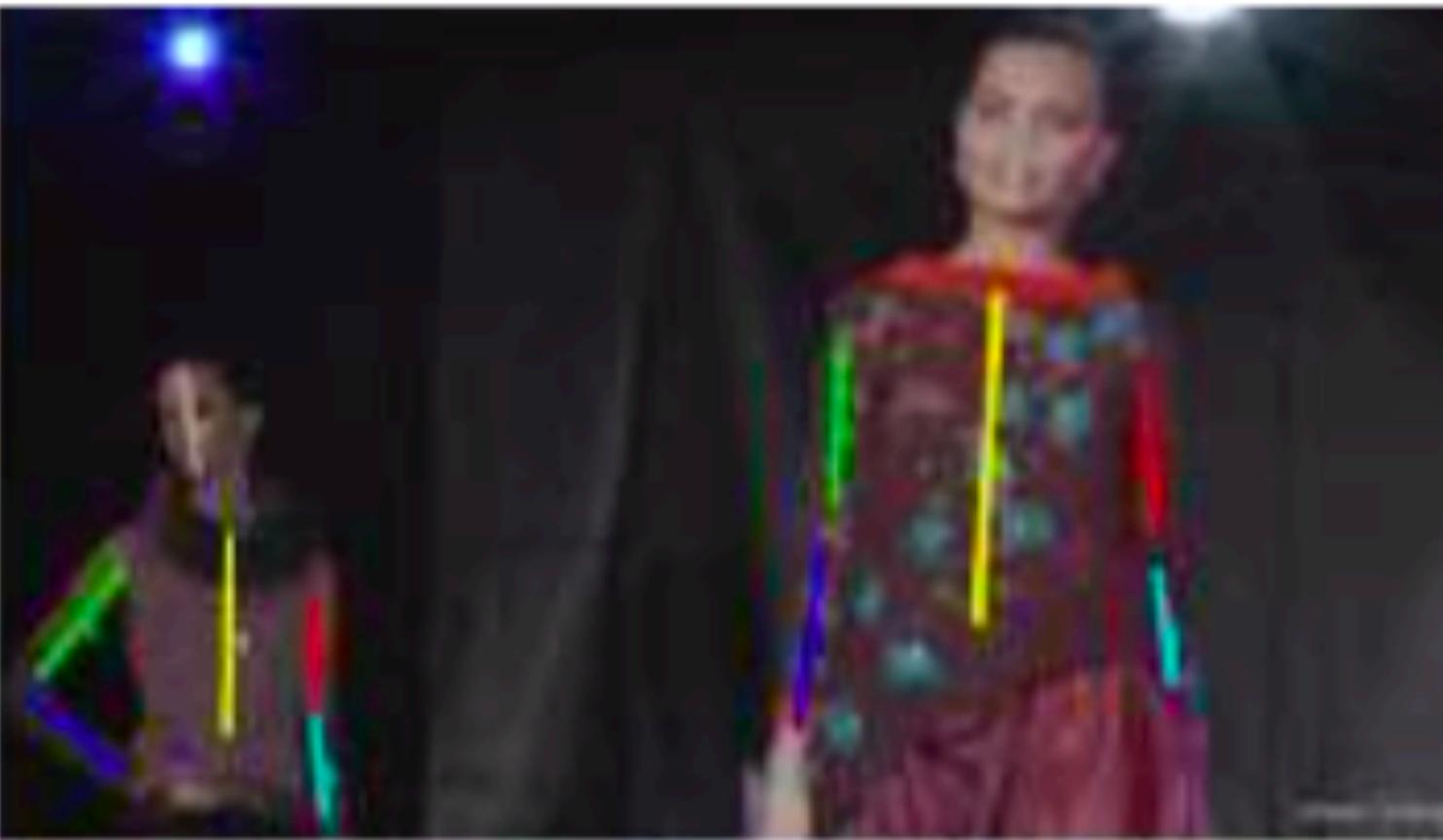


Random



# Shuffle & Learn

## Fine-tune on Human Keypoint Estimation

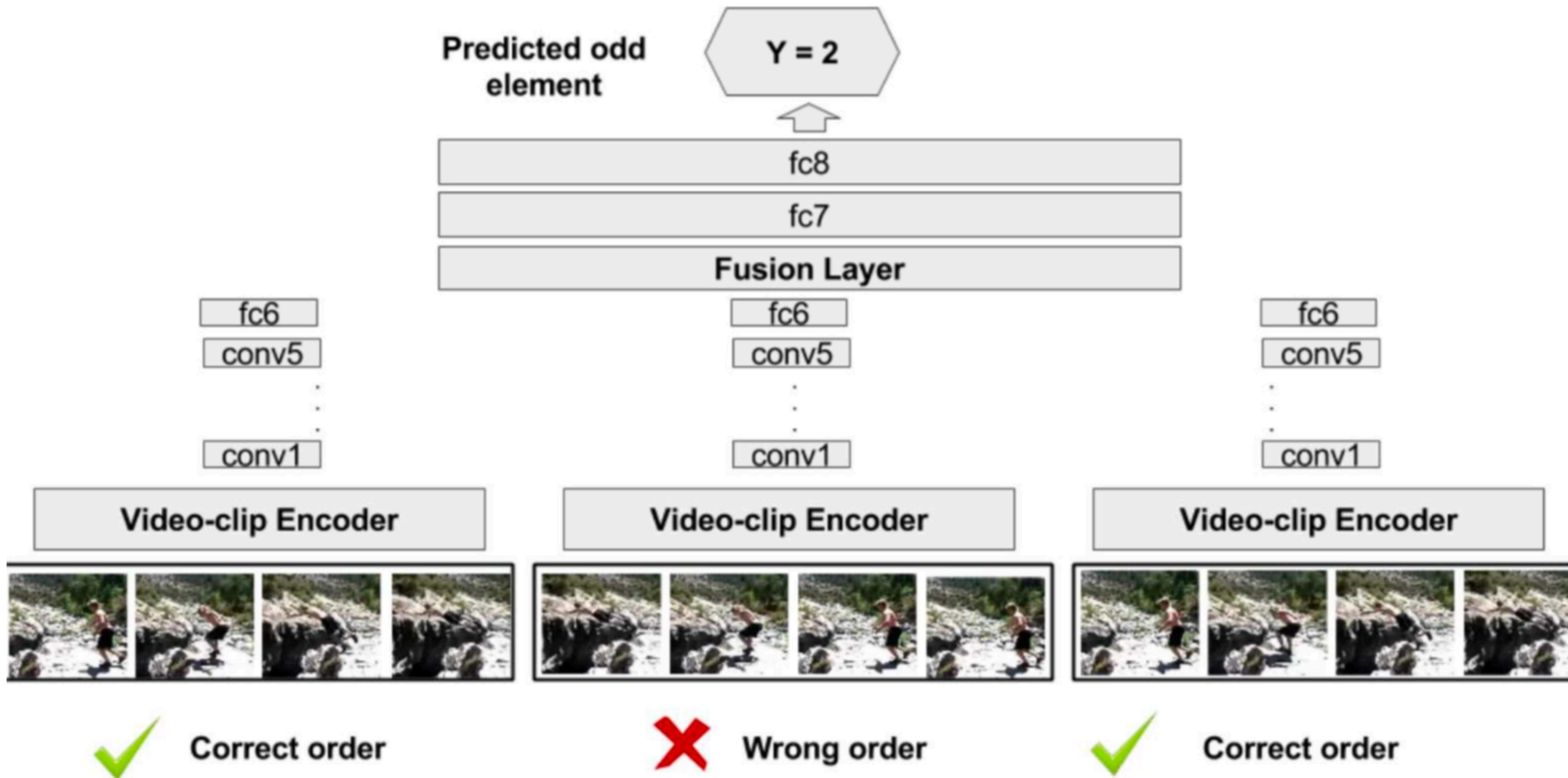


# Shuffle & Learn

Fine-tune on Human Keypoint Estimation

Initialization (AlexNet)	End task	
	FLIC Dataset Keypoints AUC	MPII Dataset Keypoints AUC
ImageNet Supervised	<b>51.3</b>	47.2
Shuffle and Learn (Self-supervised)	49.6	<b>47.6</b>

# Odd-one-out Networks



# Self-supervision in computer vision

- Using images
- Using video
- Using video and sound

# Audio-Visual co-supervision

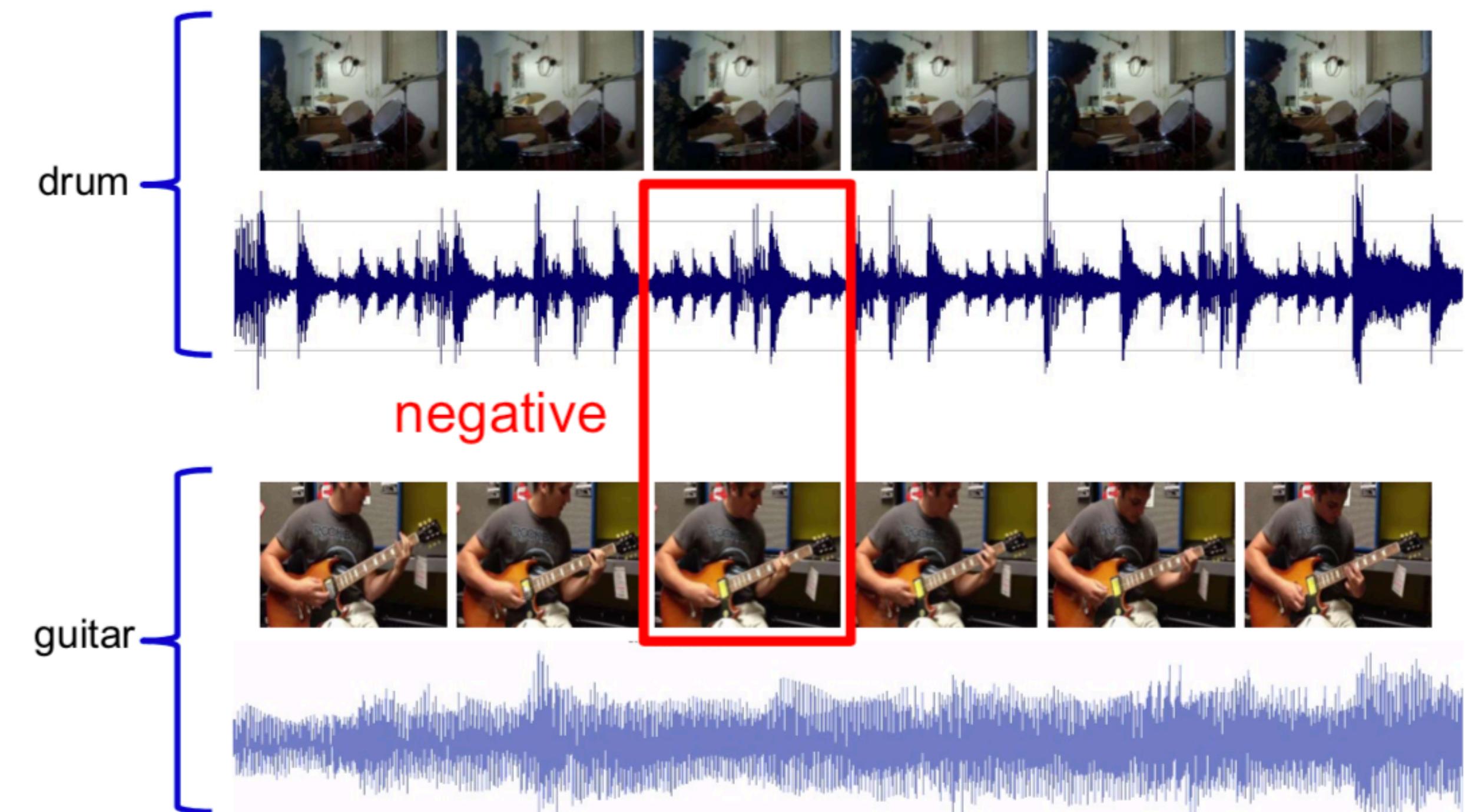
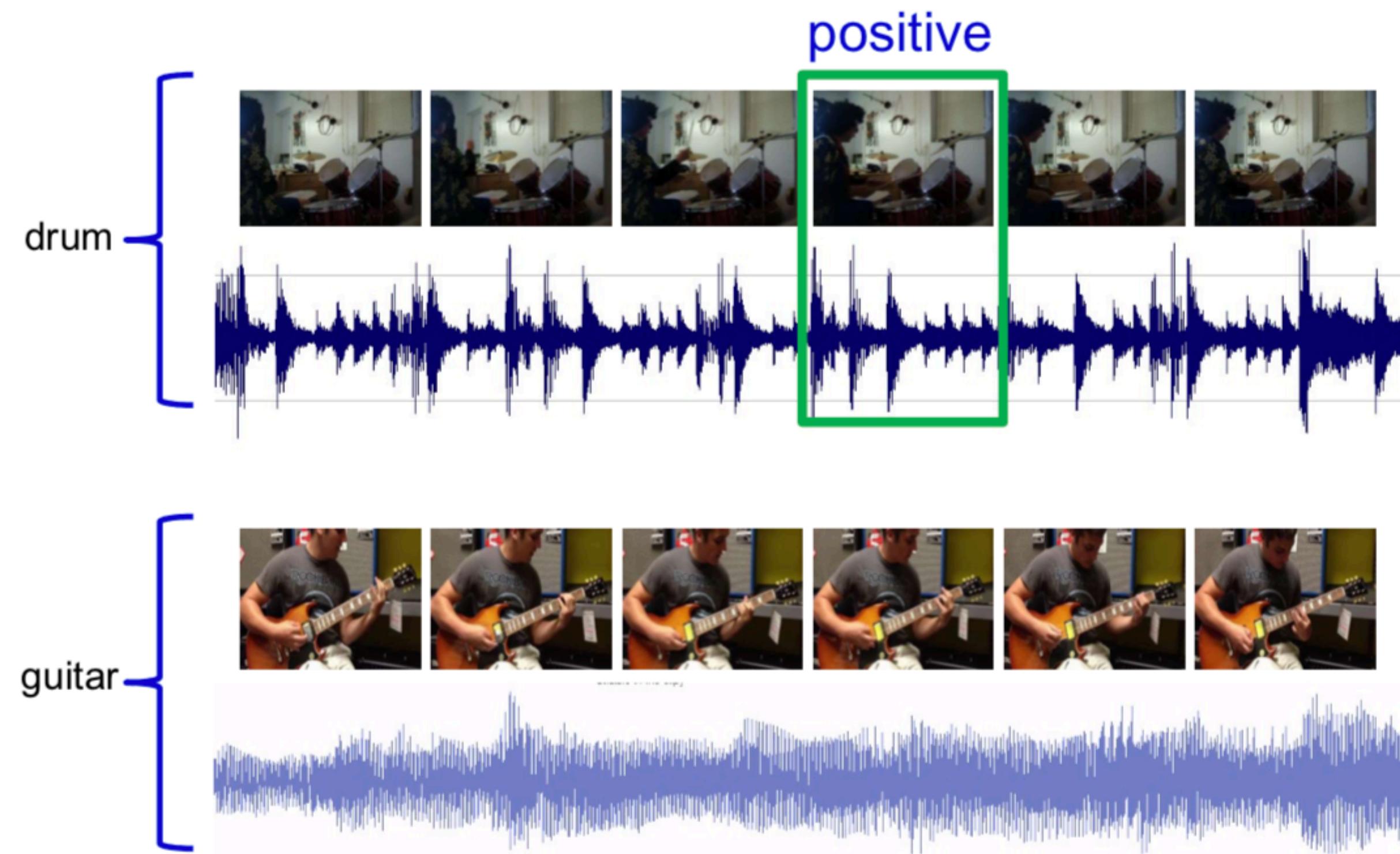
Train a network to predict if **image** and audio clip correspond



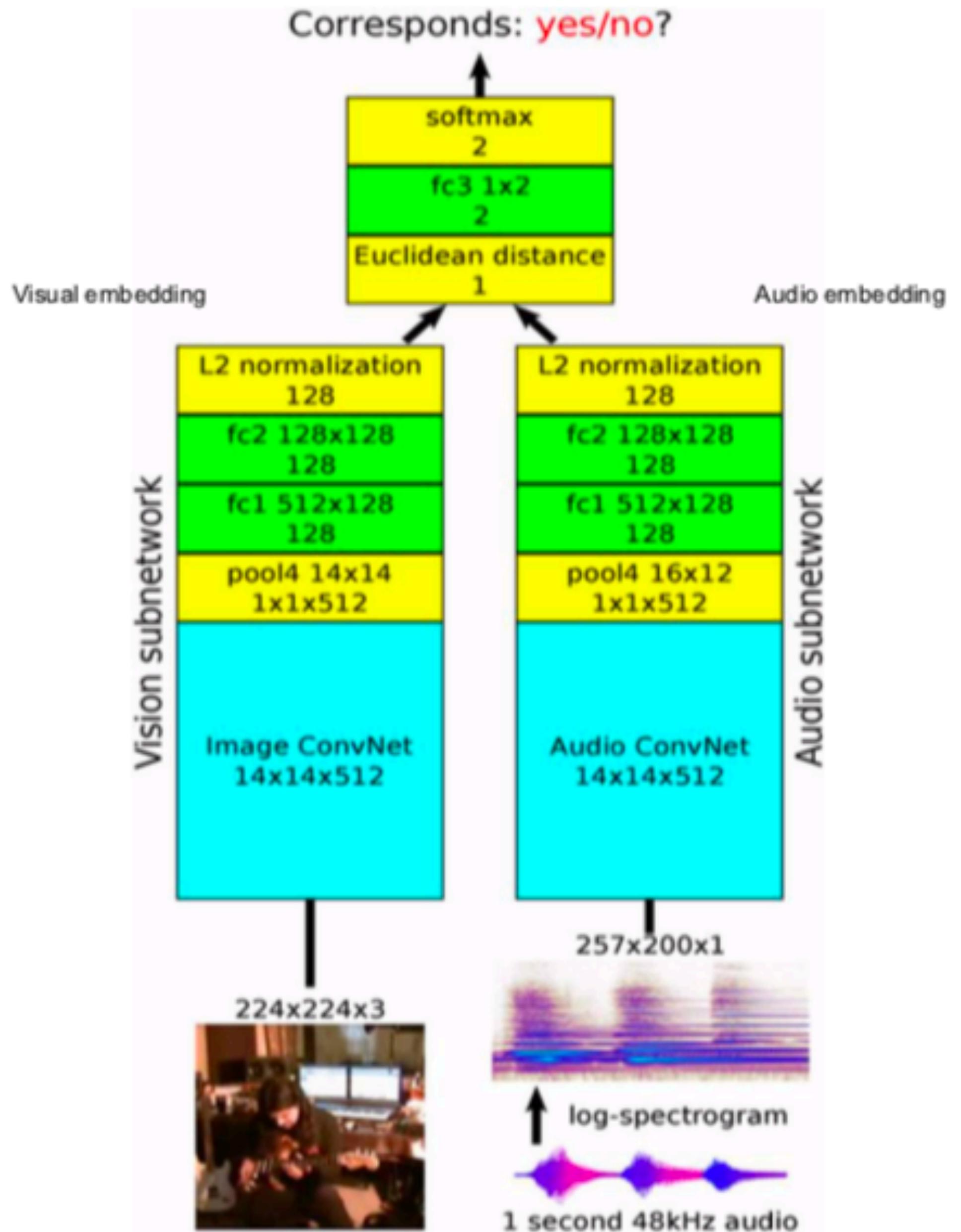
Correspond?



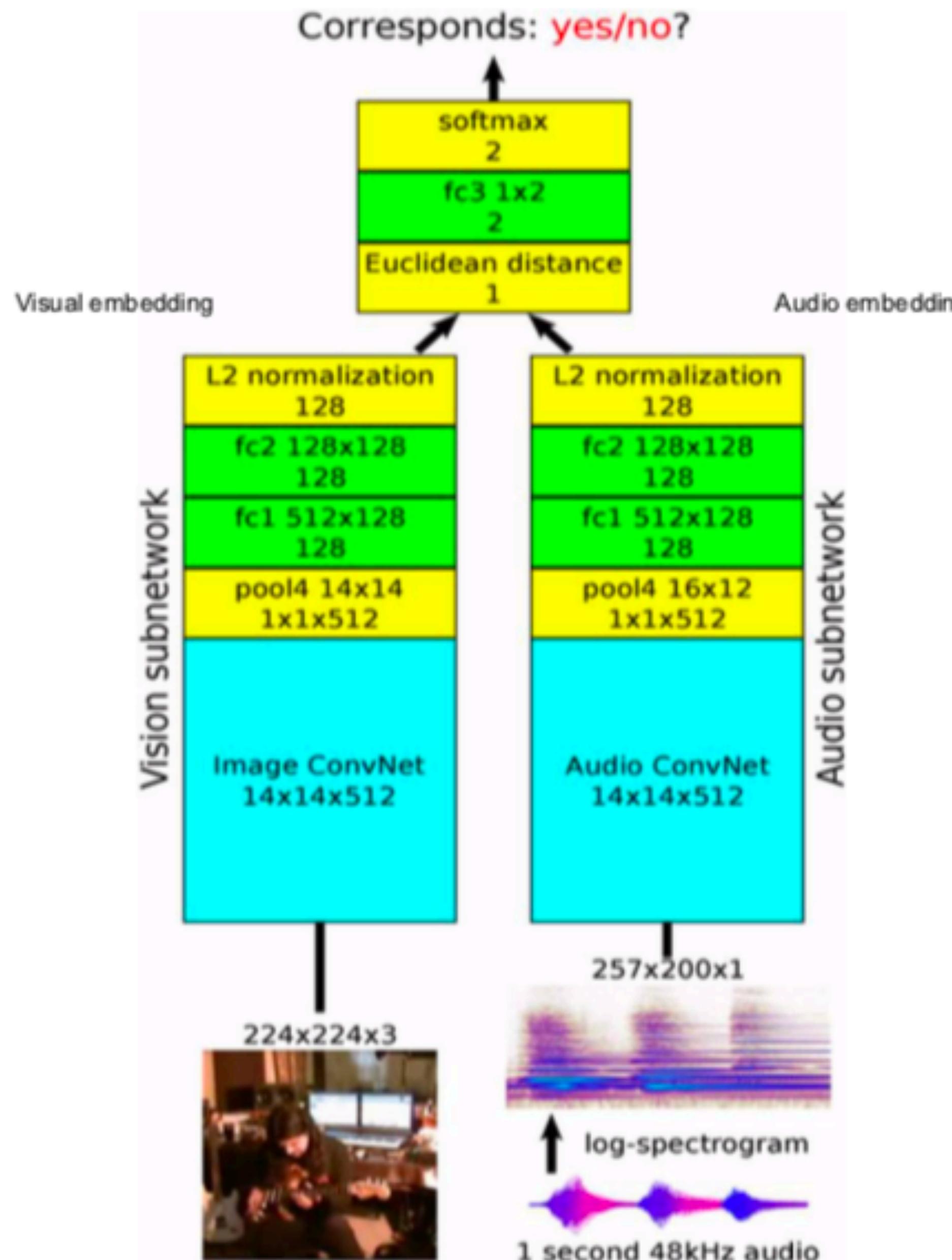
# Objects that Sound



# Objects that Sound



# Objects that Sound



## What can be learnt?

- Good representations – Visual features
  - Audio features
- Intra- and cross-modal retrieval
  - Aligned audio and visual embeddings
- “What is making the sound?”
  - Learn to localize objects that sound

# Objects that Sound

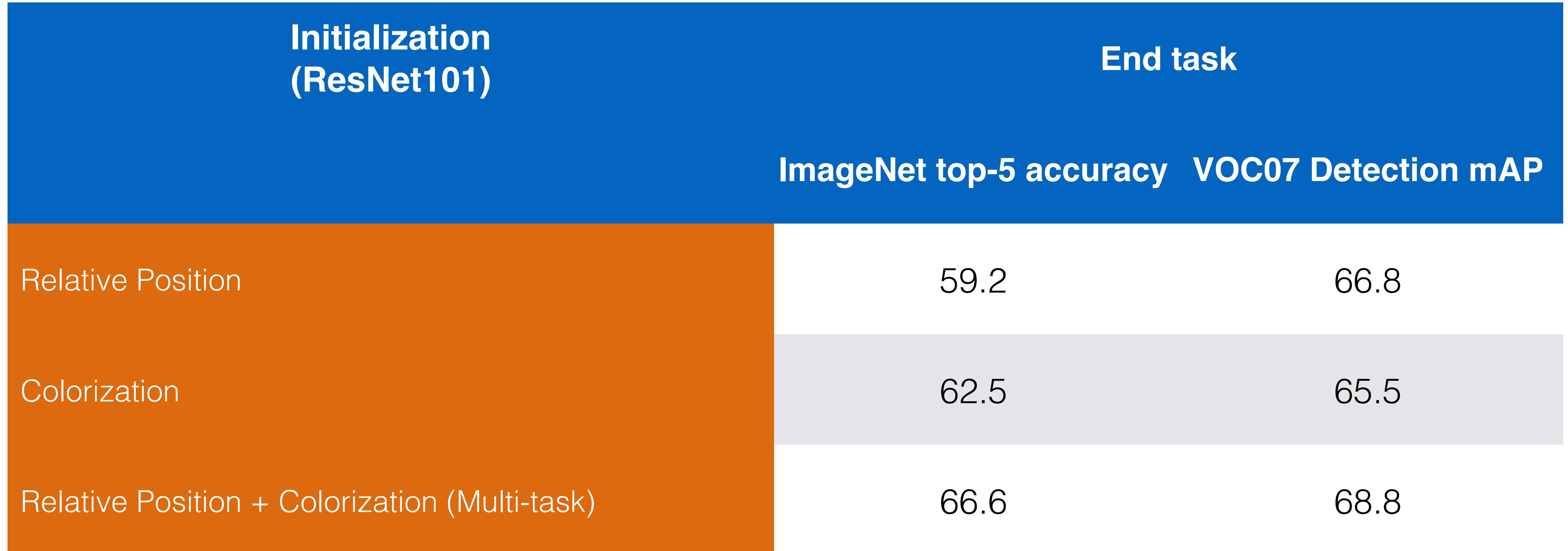
**What would make this sound?**



**Note, no video (motion) information is used**

Understanding what the “pretext” task learns

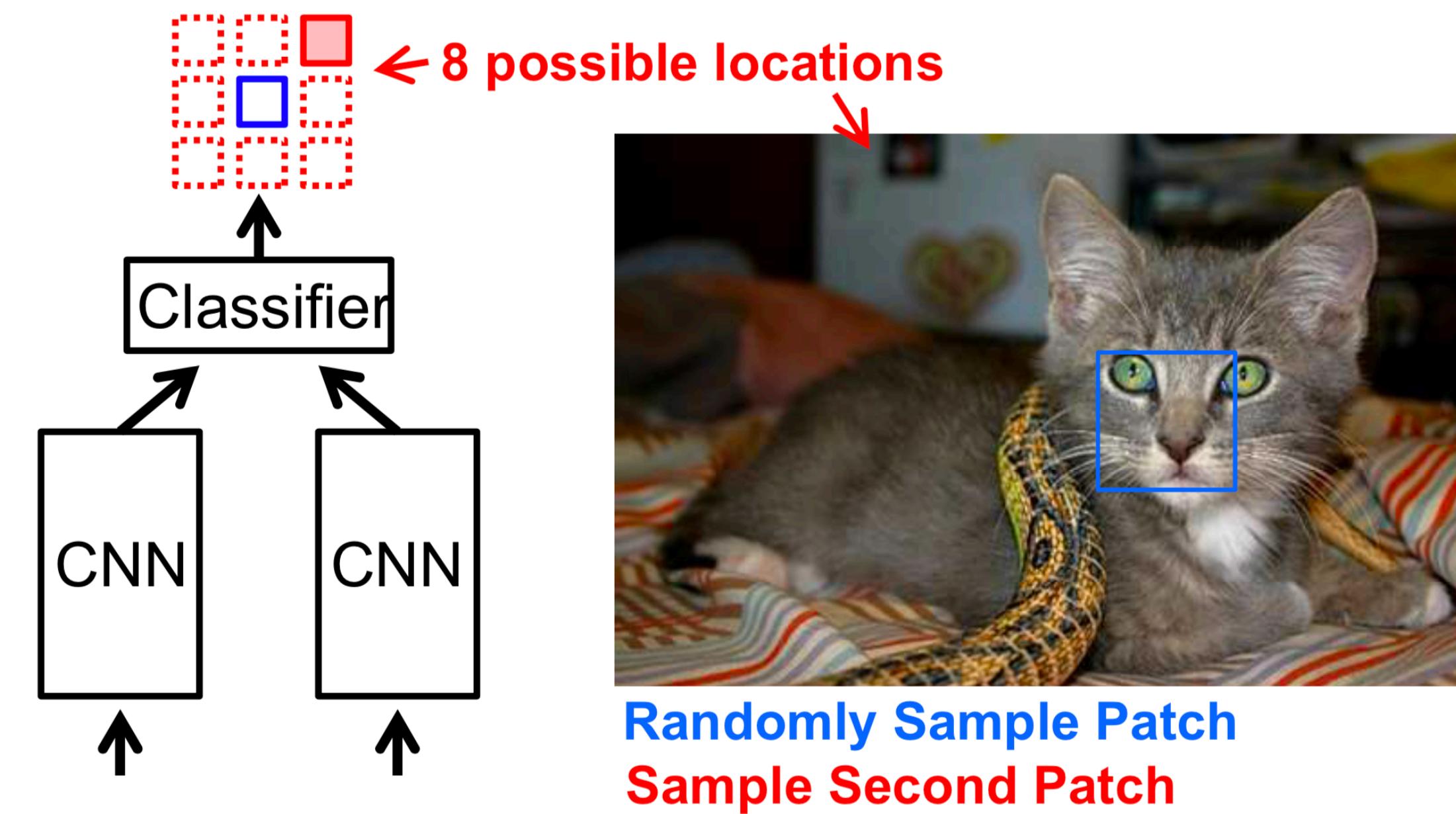
# Are they complementary?



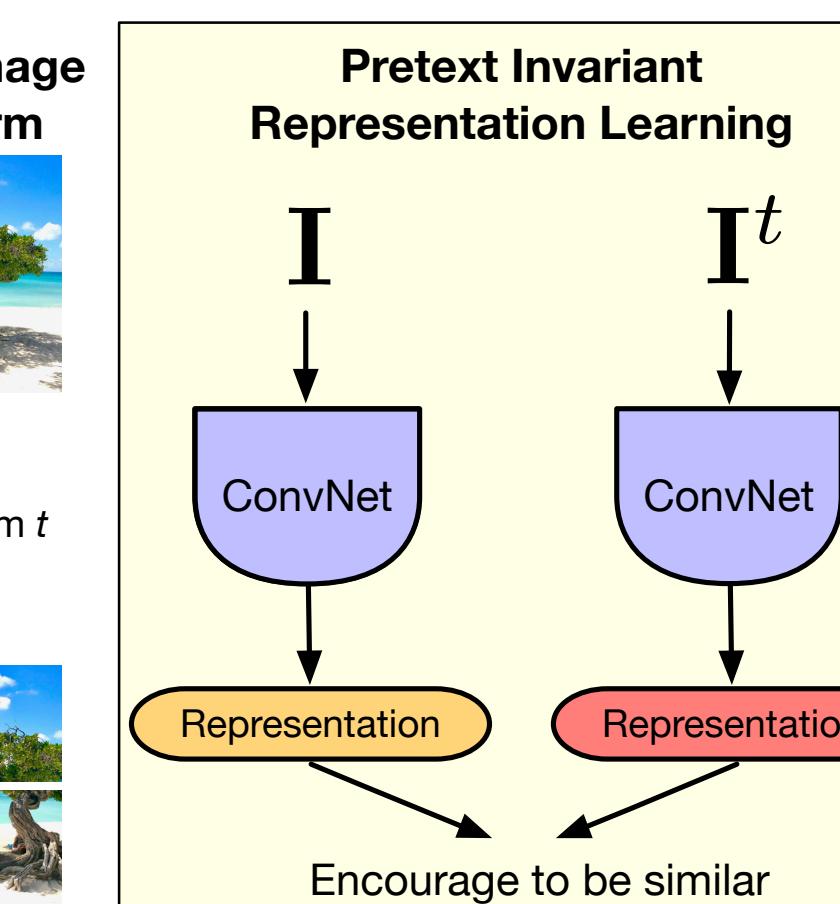
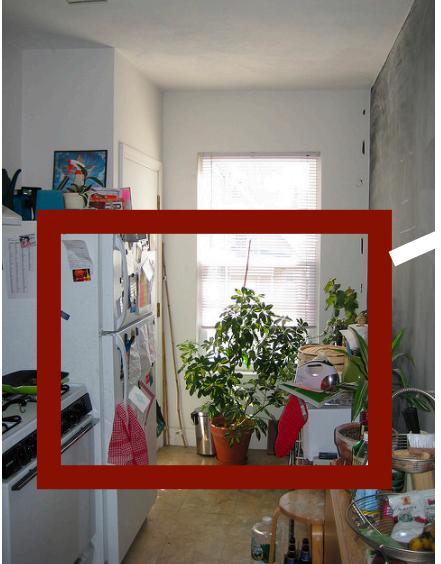
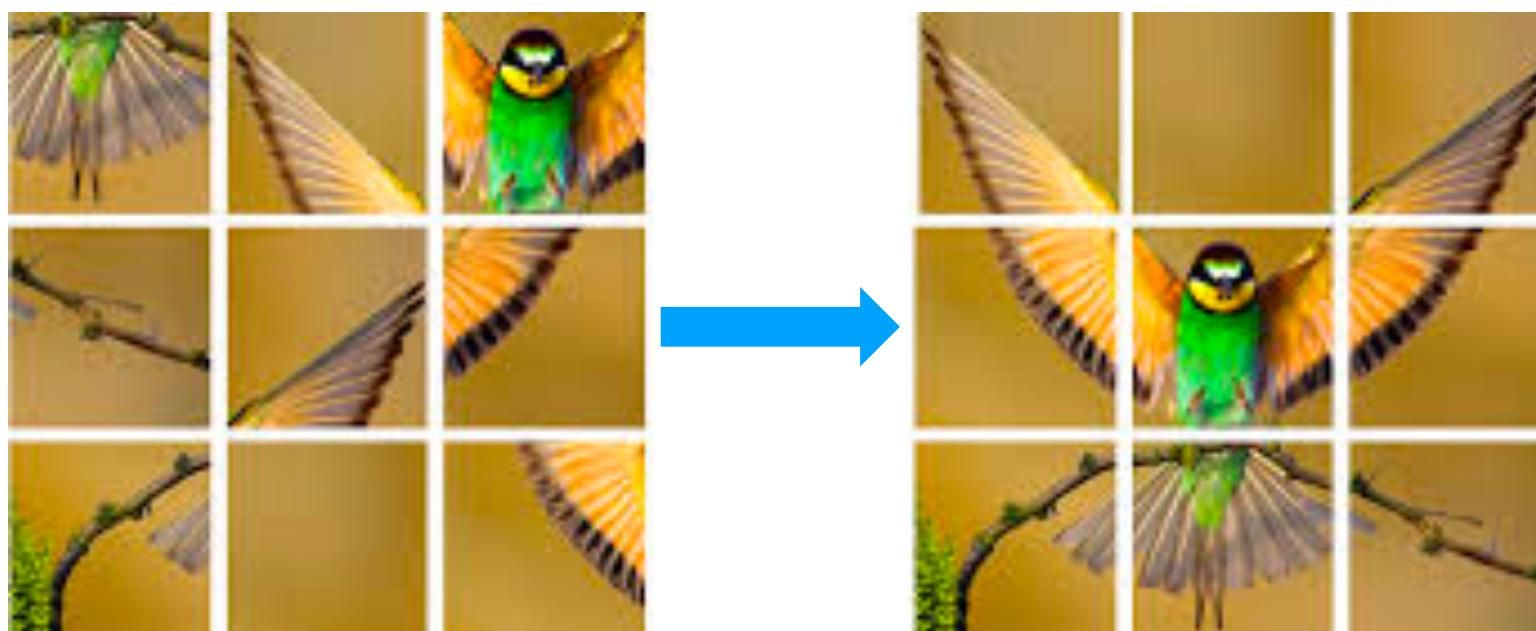
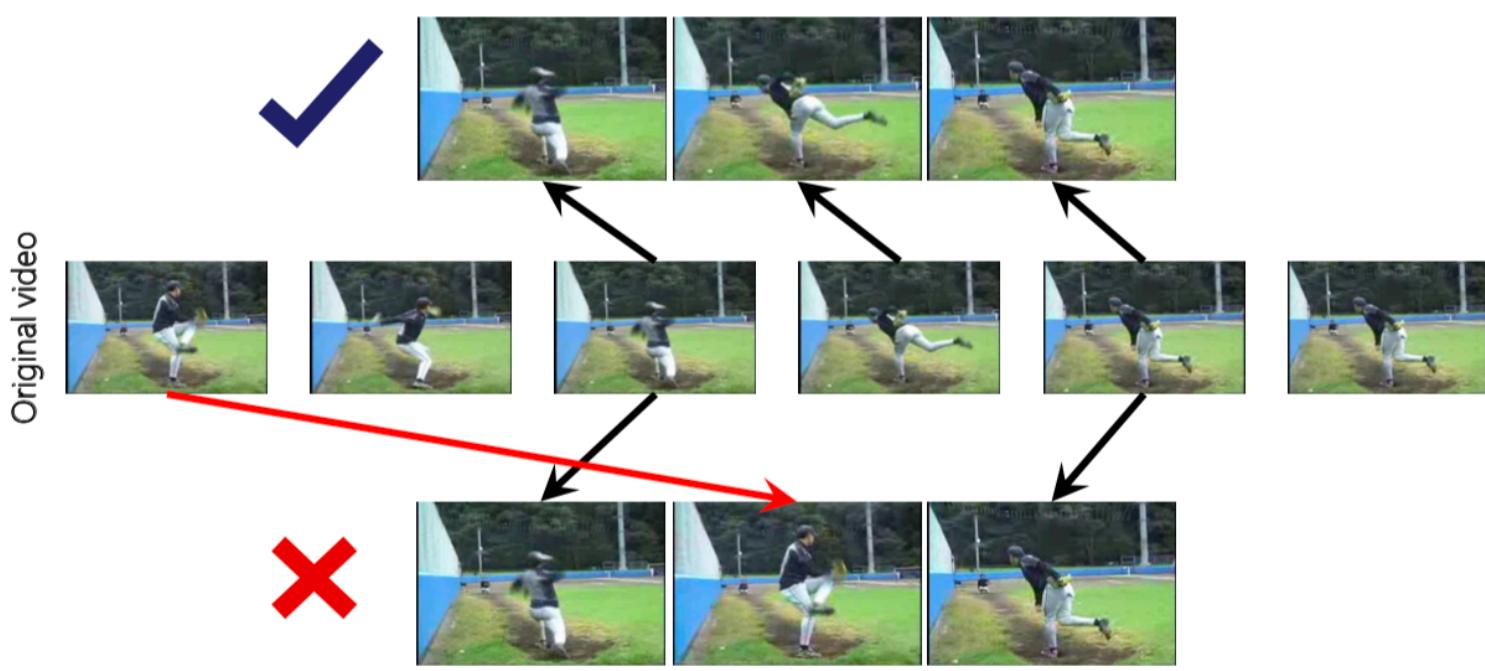
# Information predicted: varies across tasks

# Less

More



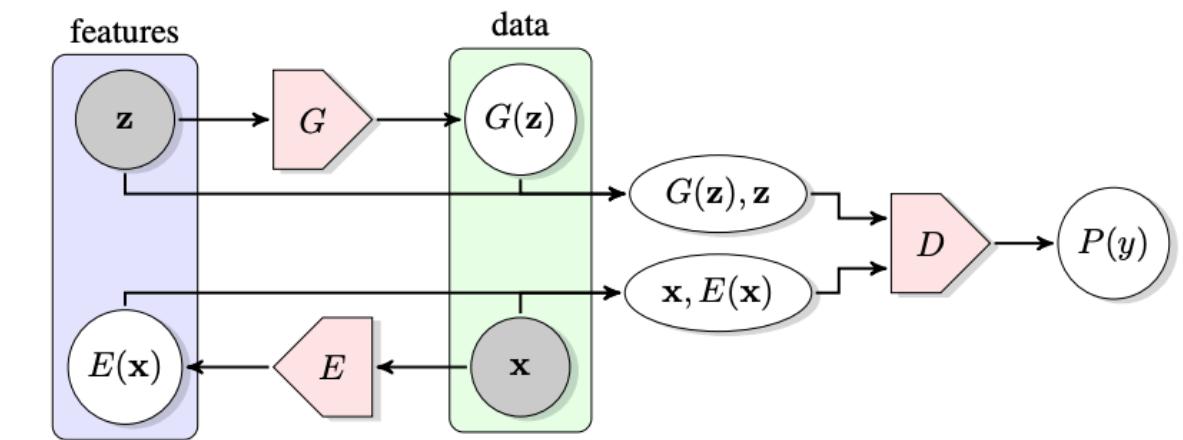
# Pretext tasks



→ Predict more information

# Contrastive

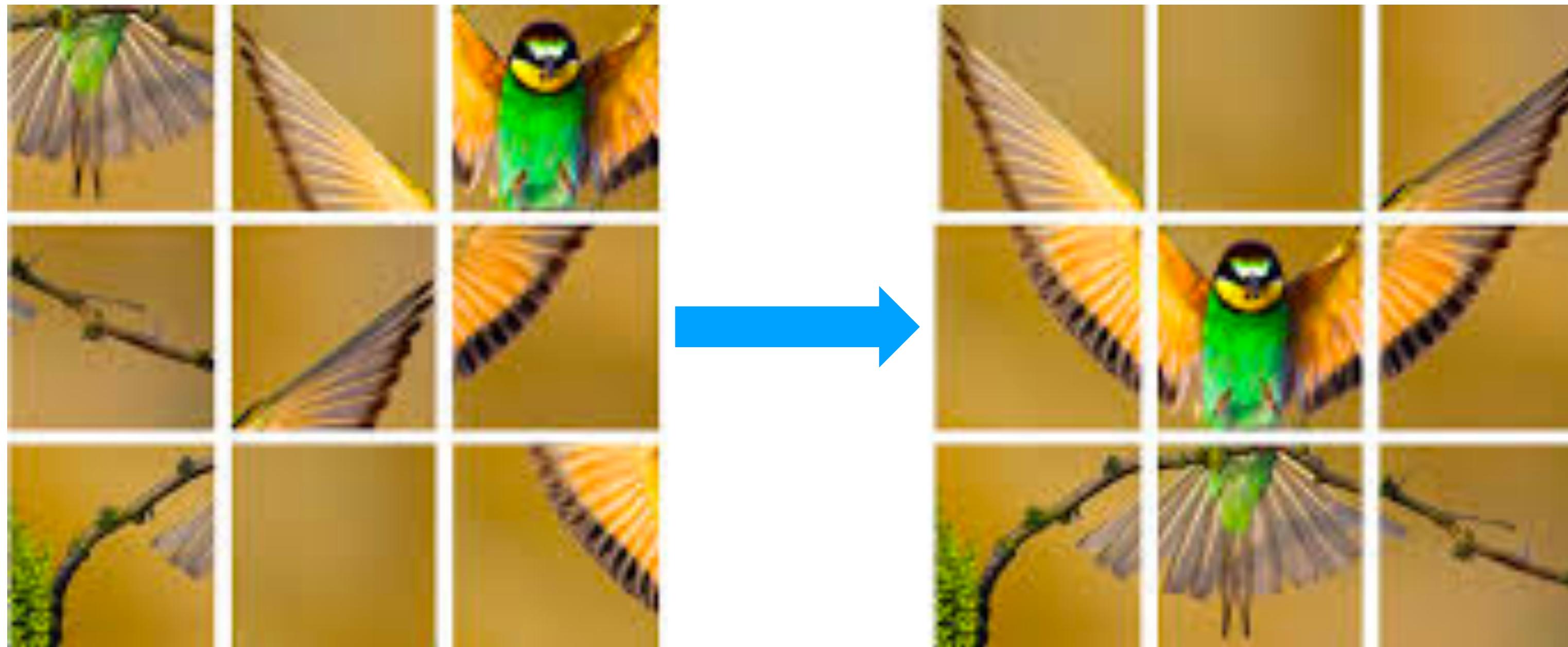
# Generative



AutoEncoder,  
VAE, GAN,  
BiGAN

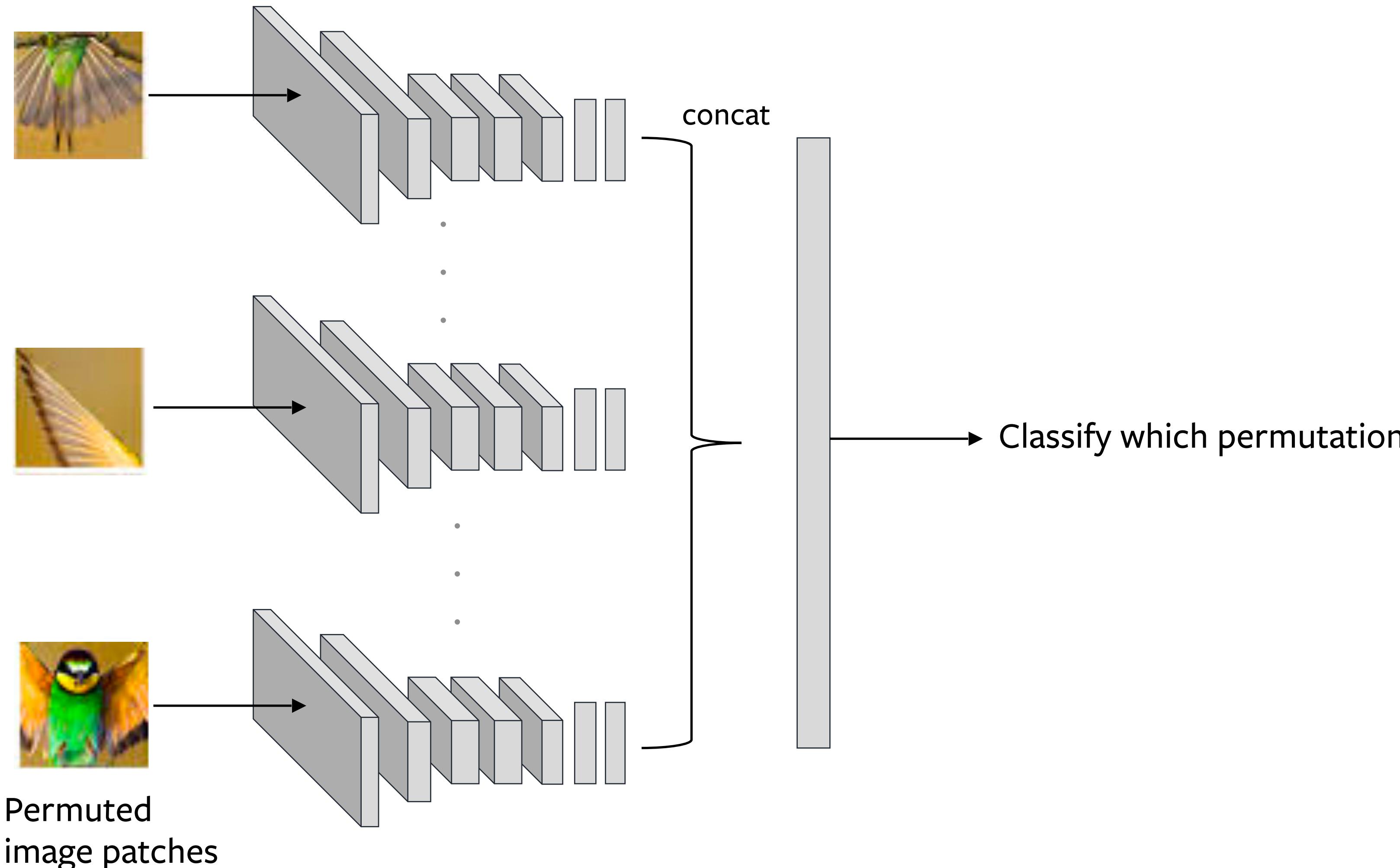
→ Predict more information

# Scaling self-supervised learning



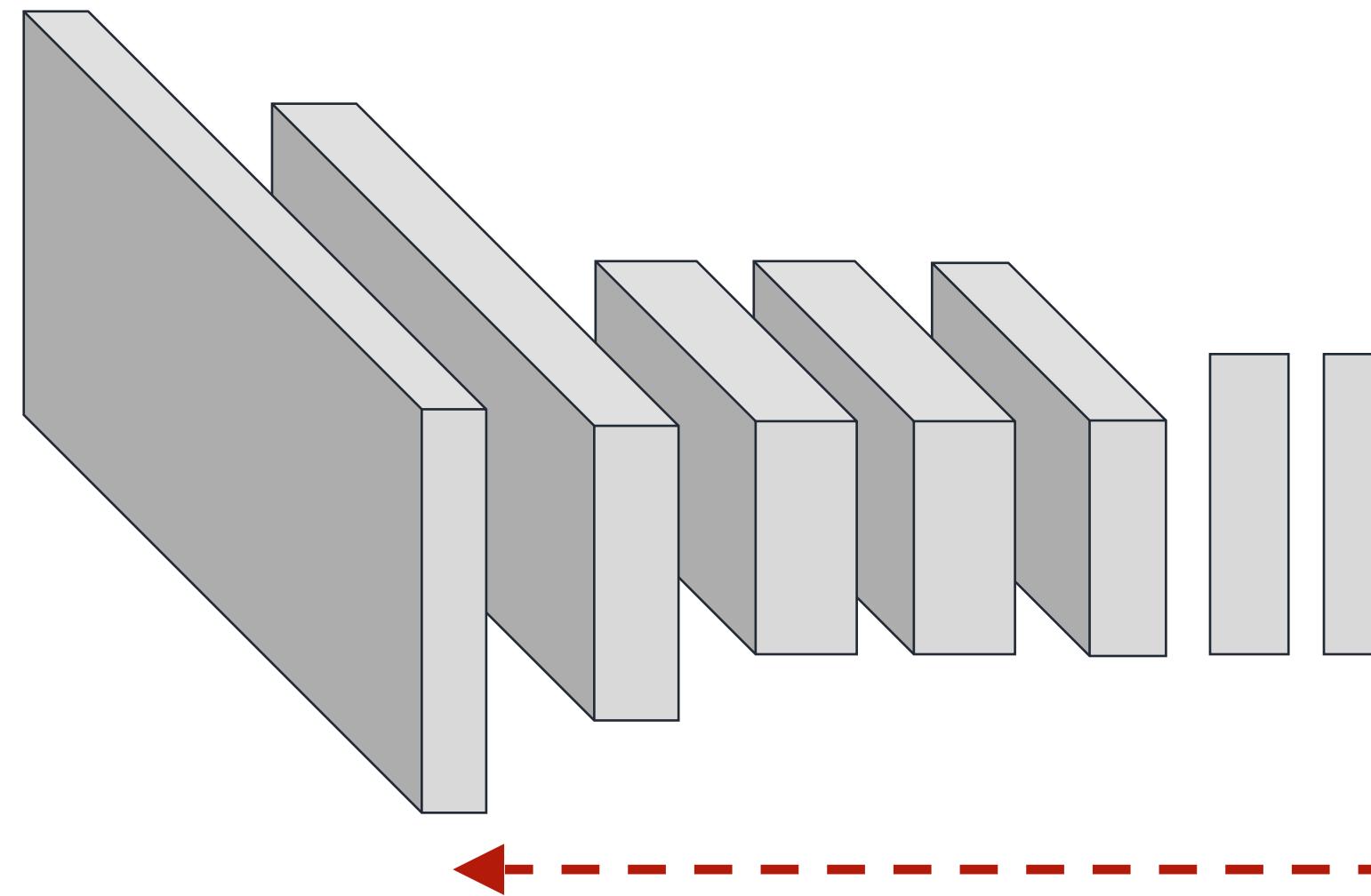
Jigsaw puzzles  
(Noorozi & Favaro, 2016)

# Jigsaw Puzzles

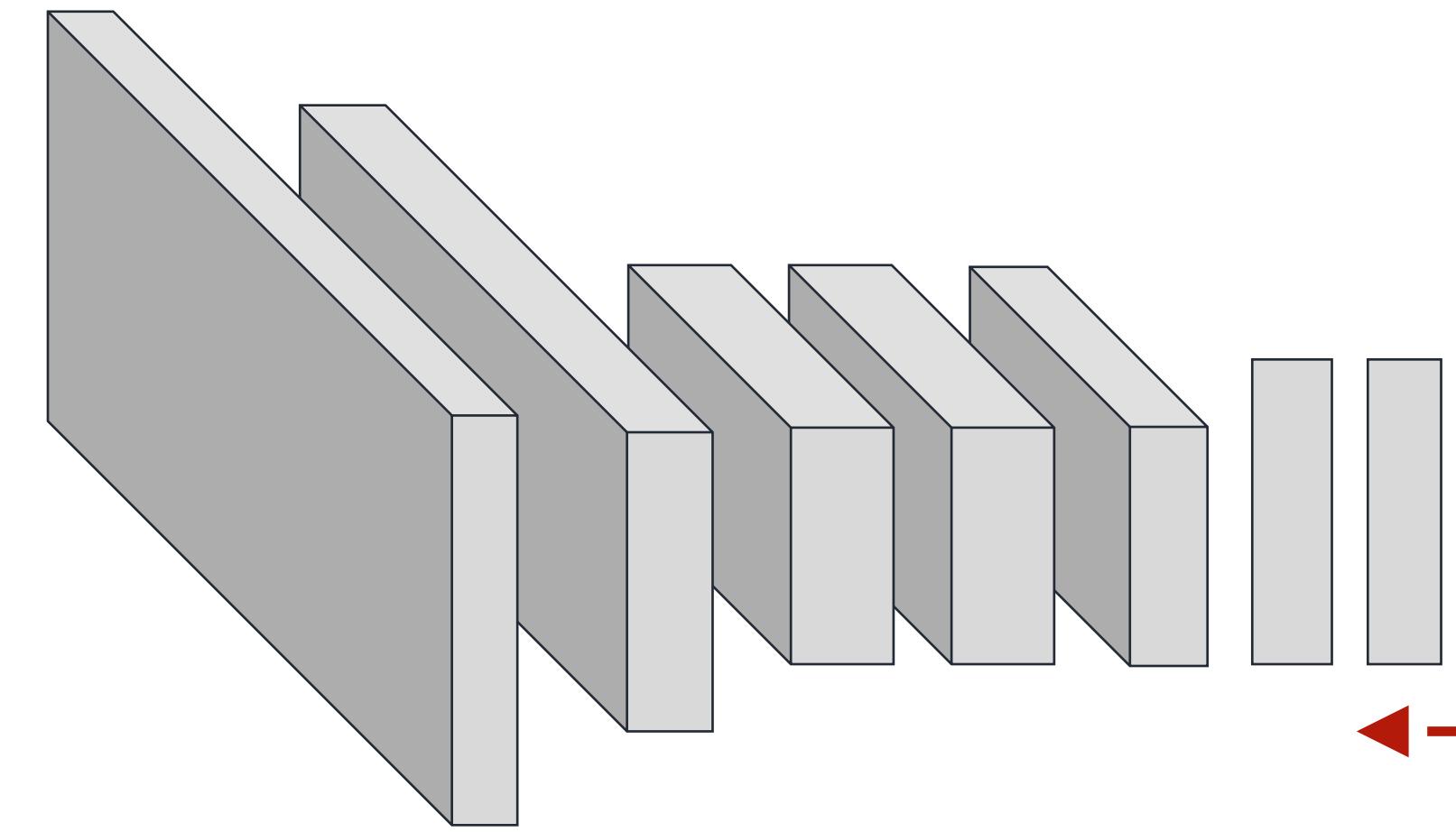


- Use  $N=9$  patches
- In practice, use a subset of permutations
- E.g. 100 from  $9!$
- Each patch is processed independently
- N-way ConvNet (shared params)
- Problem Complexity
  - Size of subset

# Evaluation – fine-tuning vs. linear classifier



Fine-tune all layers



Linear classifier

*A good representation transfers with **little training***

# Evaluation – many tasks

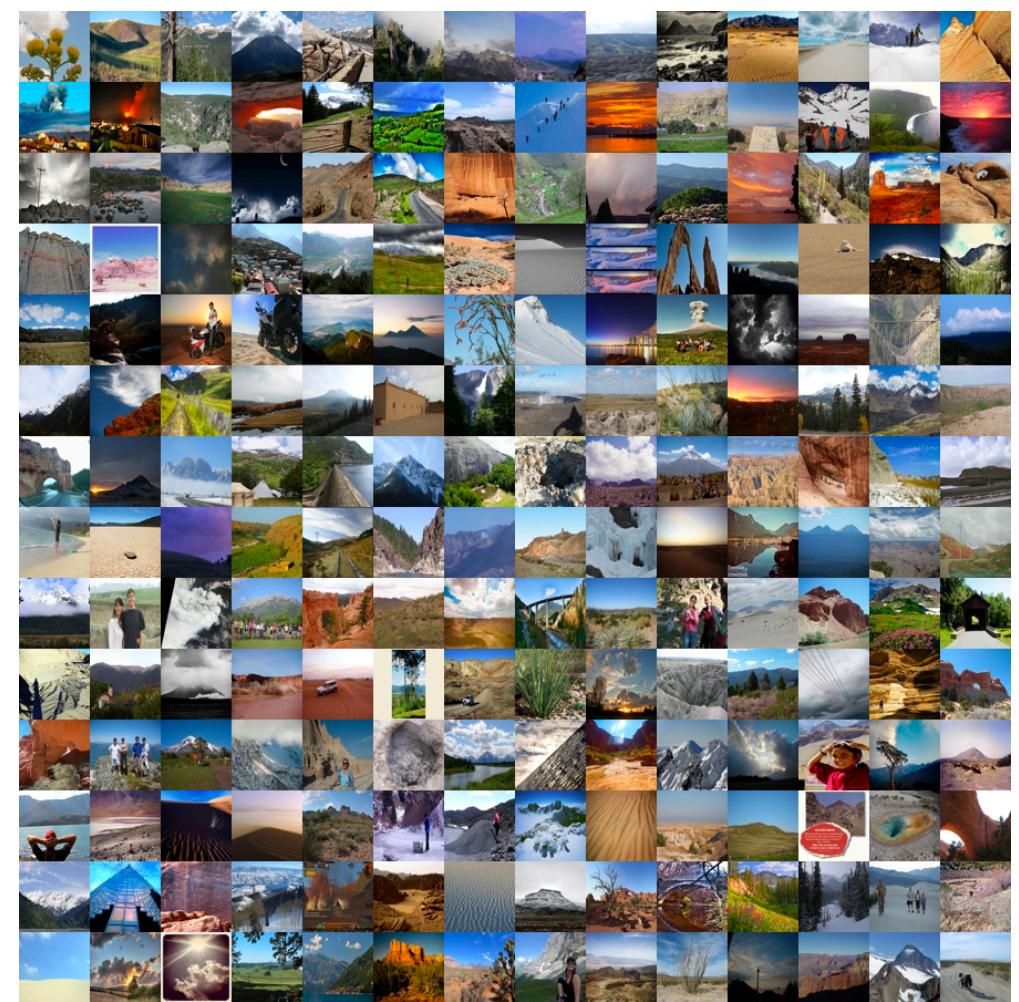
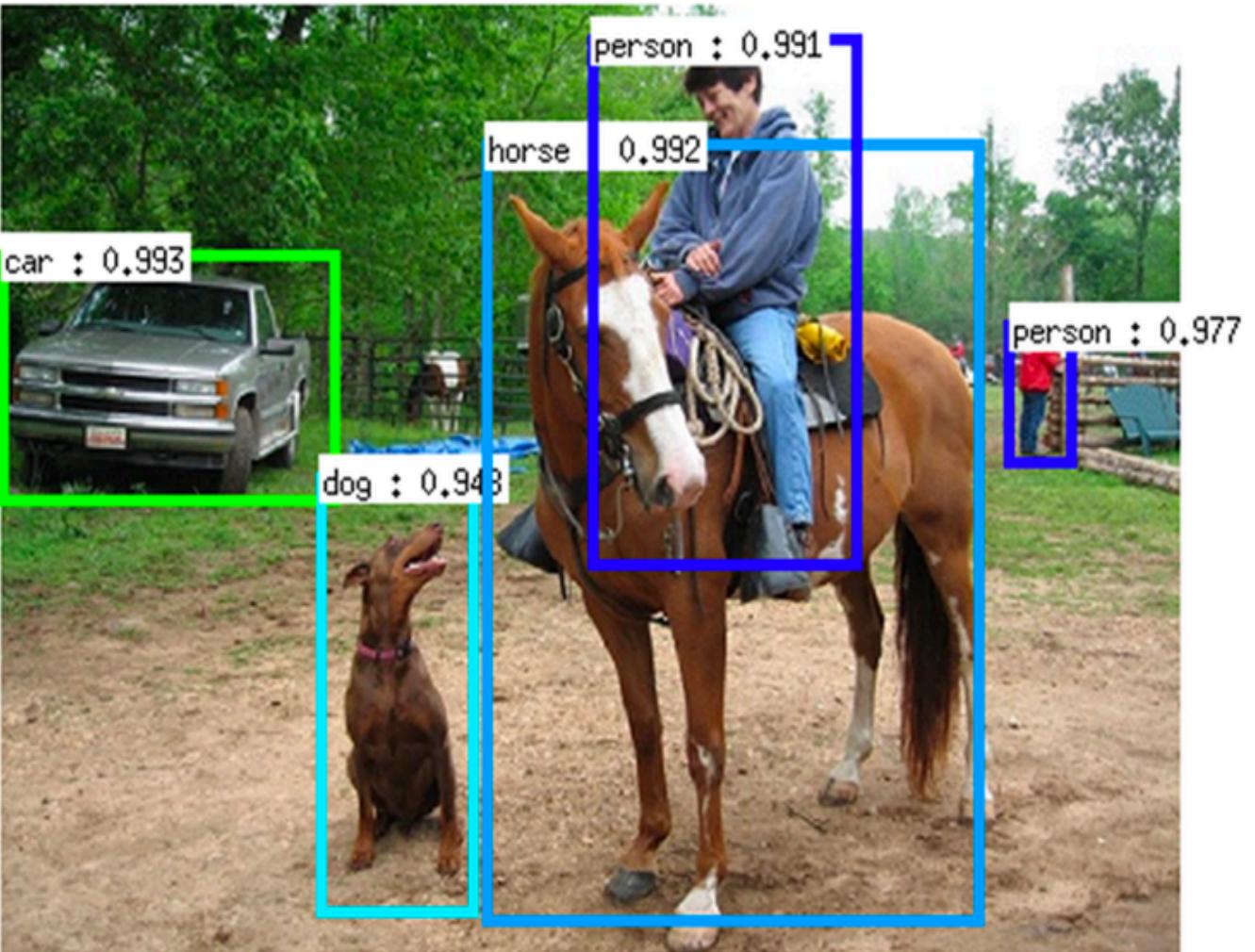


Image classification  
Few-shot learning

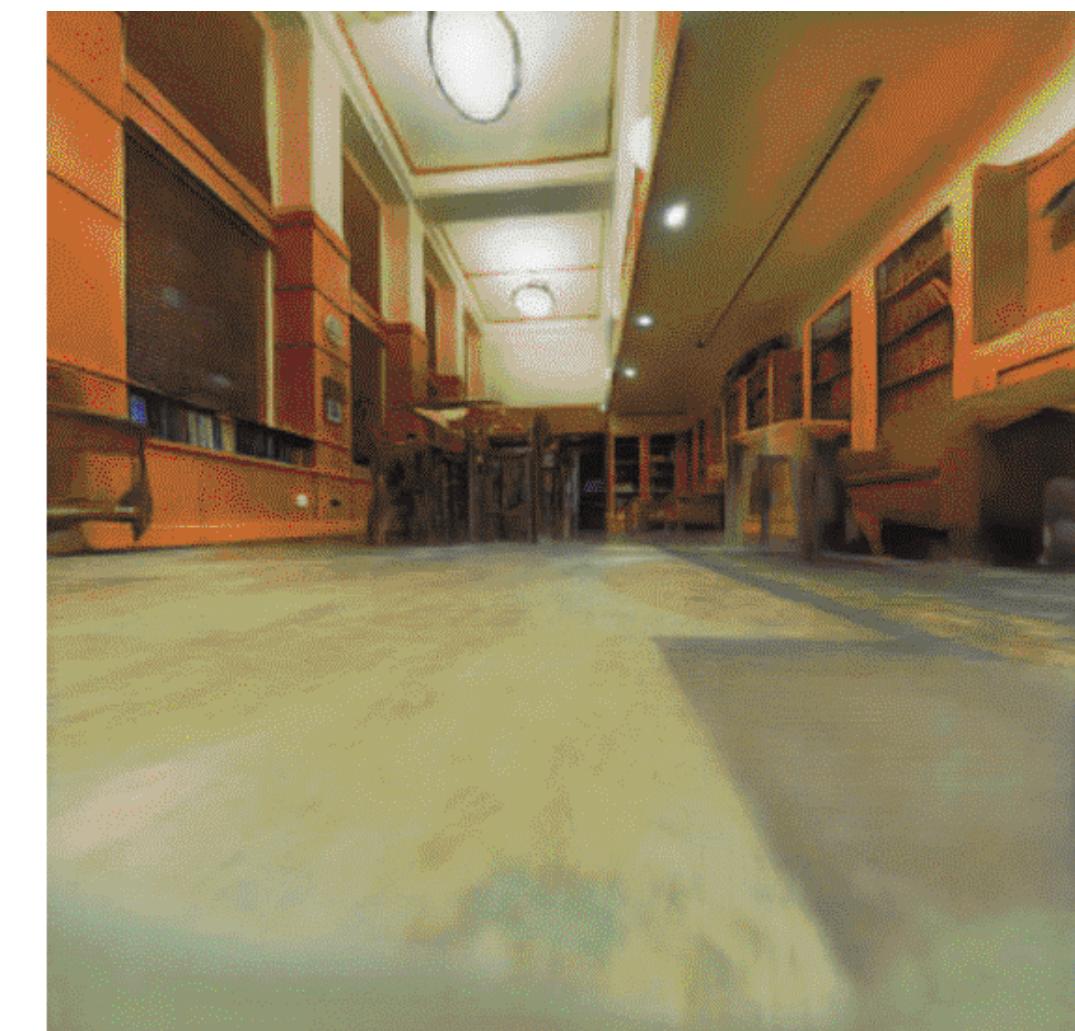
ImageNet, Places-205, VOC'07,  
COCO



Object detection  
VOC'07



3D Understanding  
Surface Normals – NYUv2

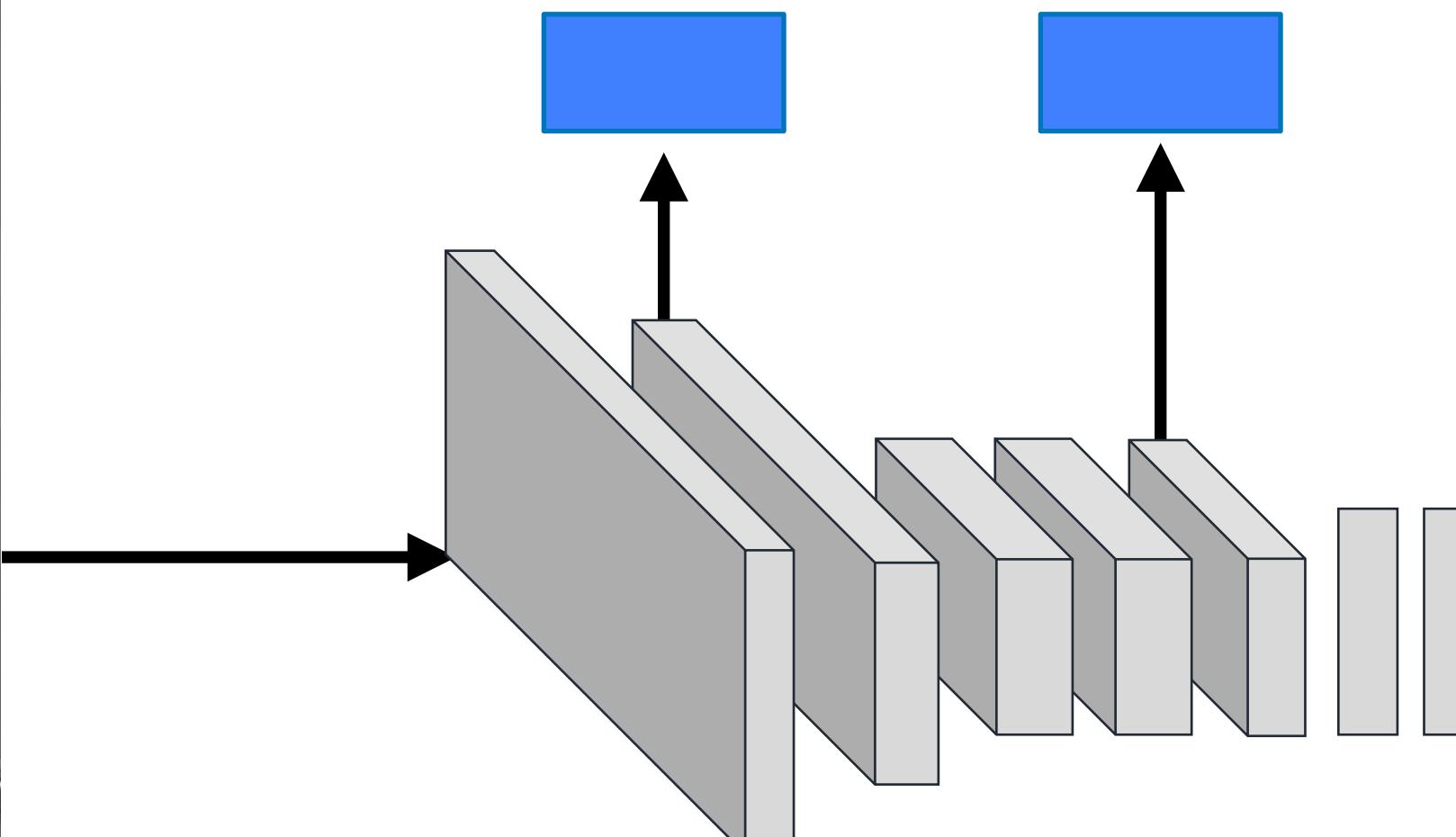


Navigation  
Gibson environment

# Evaluating the representation



Extract "fixed" features



# Evaluating the representation

- Train a Linear SVM on **fixed feature** representations
- Use the VOC07 image classification task



aero

bicycle

bird

boat

bottle



bus

car

cat

chair

cow



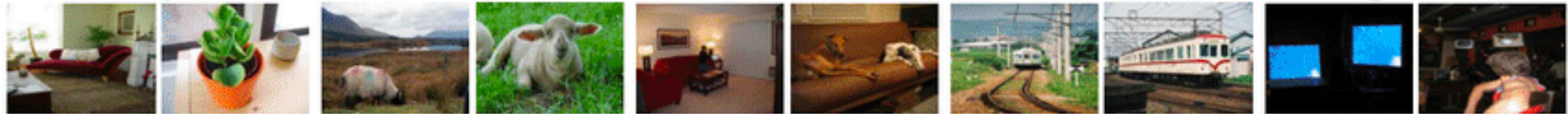
diningtable

dog

horse

mbike

person



plant

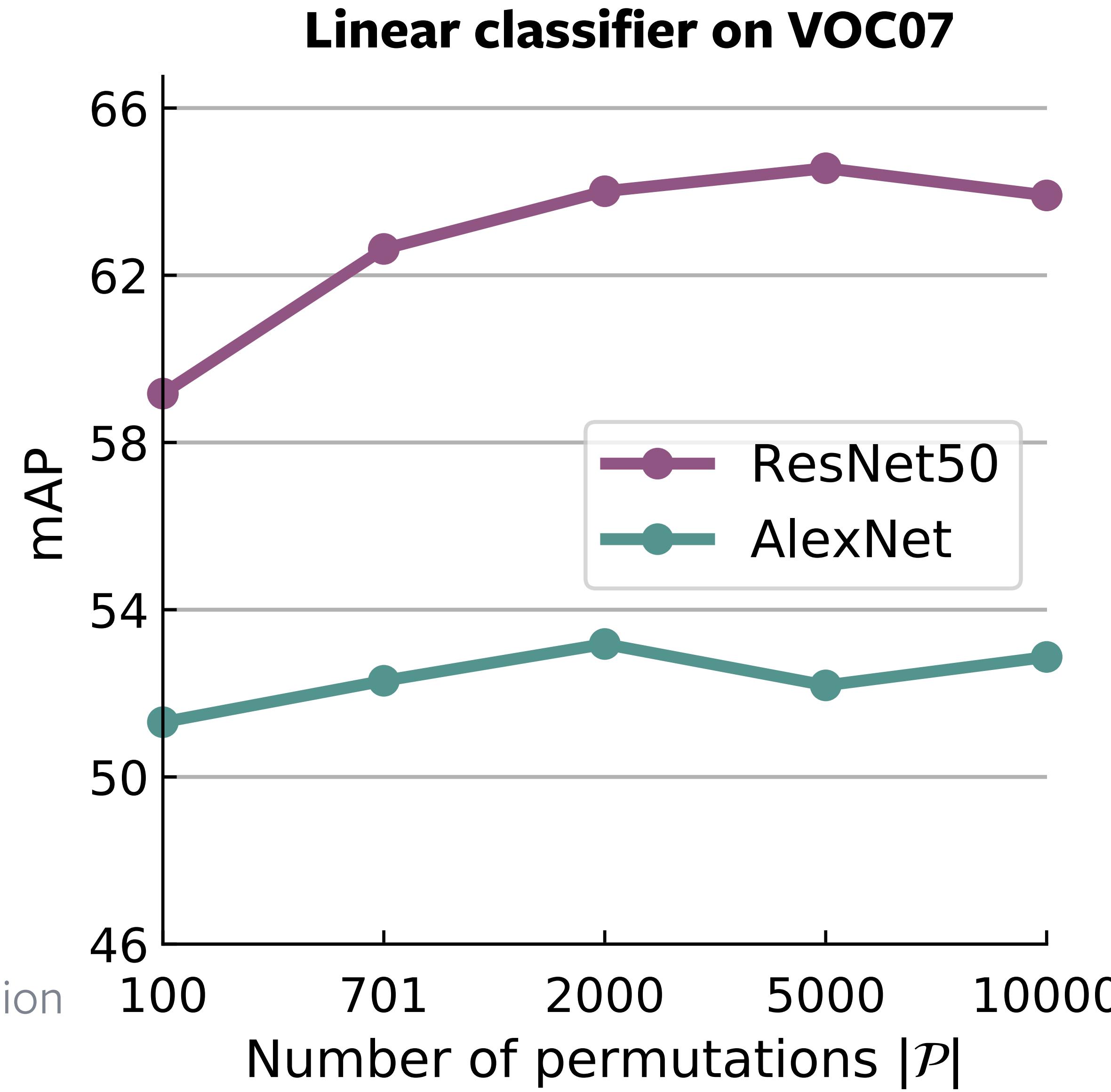
sheep

sofa

train

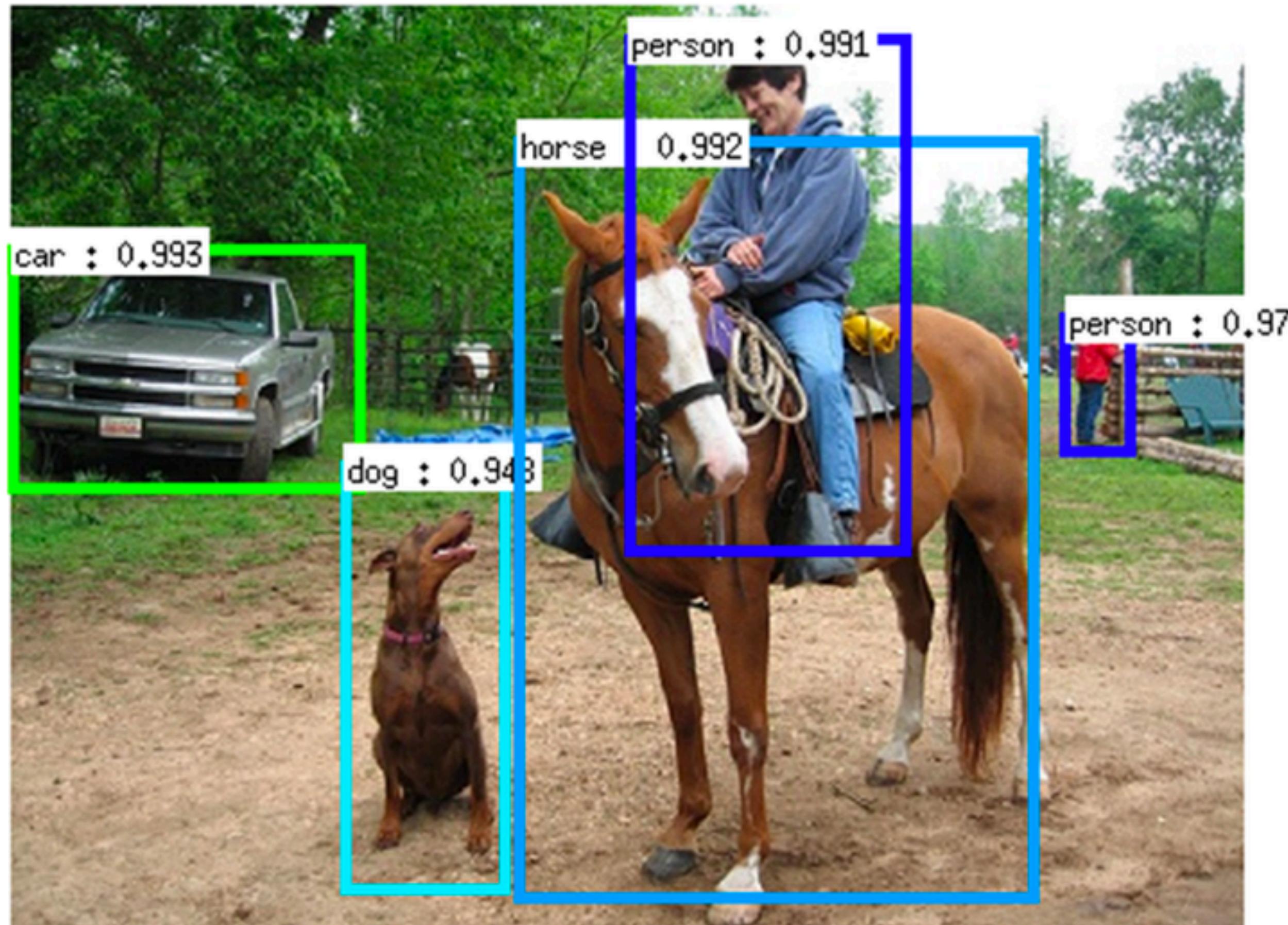
tv

# Increasing amount of information predicted

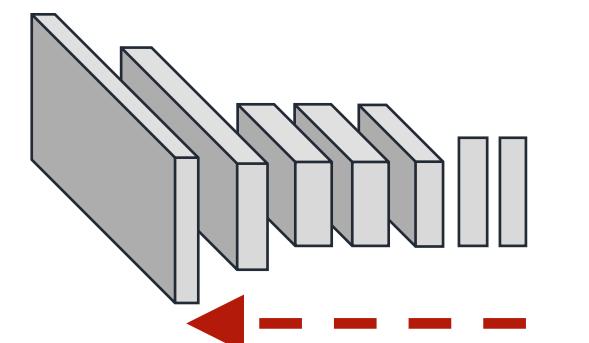


# Object Detection

- Fast R-CNN (Girshick et al., 2015)
  - Same optimization parameters for all methods (including supervised)
  - Use VOC'07



# Object Detection



VOC07 test set.

Fast R-CNN ResNet50

Initialization	Train Set	
	VOC07	VOC07+12
ImageNet Supervised	70.5	76.2
Jigsaw ImageNet 14M	69.2	75.4

**Within error of ImageNet pre-trained network**

# Surface Normal Estimation

- Predict surface normals on NYU-v2
  - Same optimization parameters for all methods (including supervised)
  - PSPNet Architecture
  - Train last few layers only (res5 onwards)



**Input**



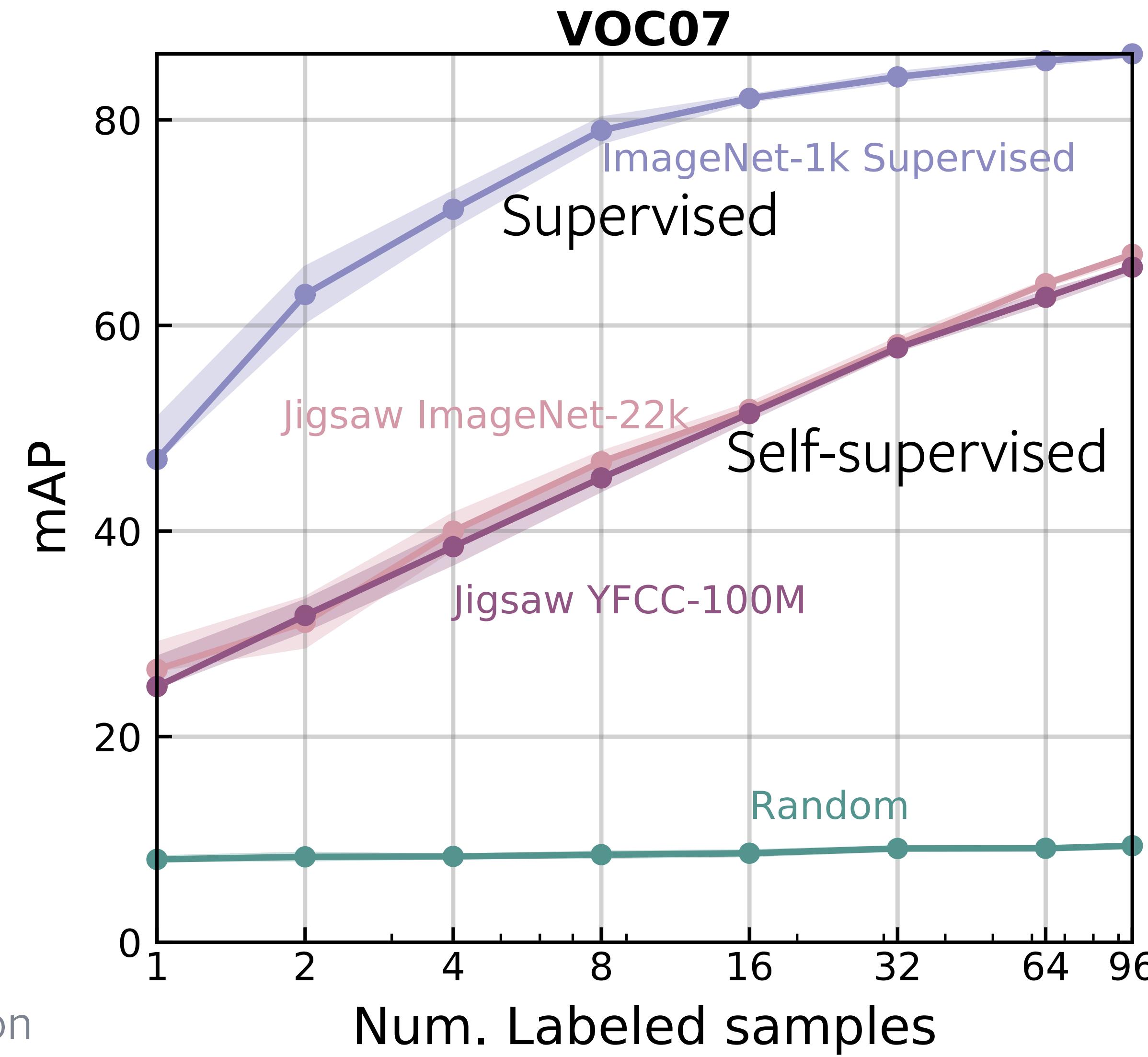
**Output**

# Surface Normal Estimation



**Outperforms ImageNet supervised**

# Few shot learning

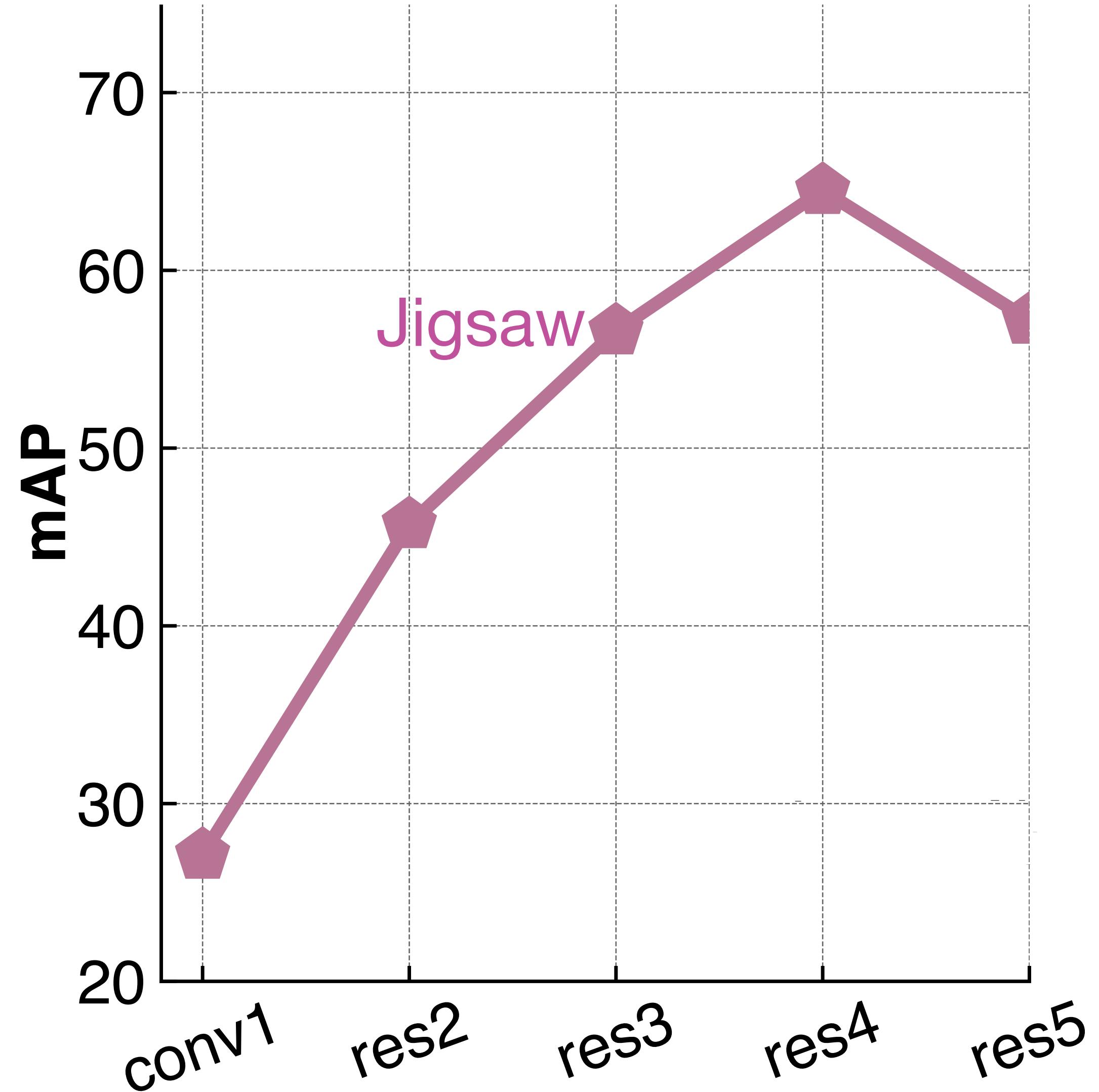


mAP = mean Average Precision  
(Higher is better)

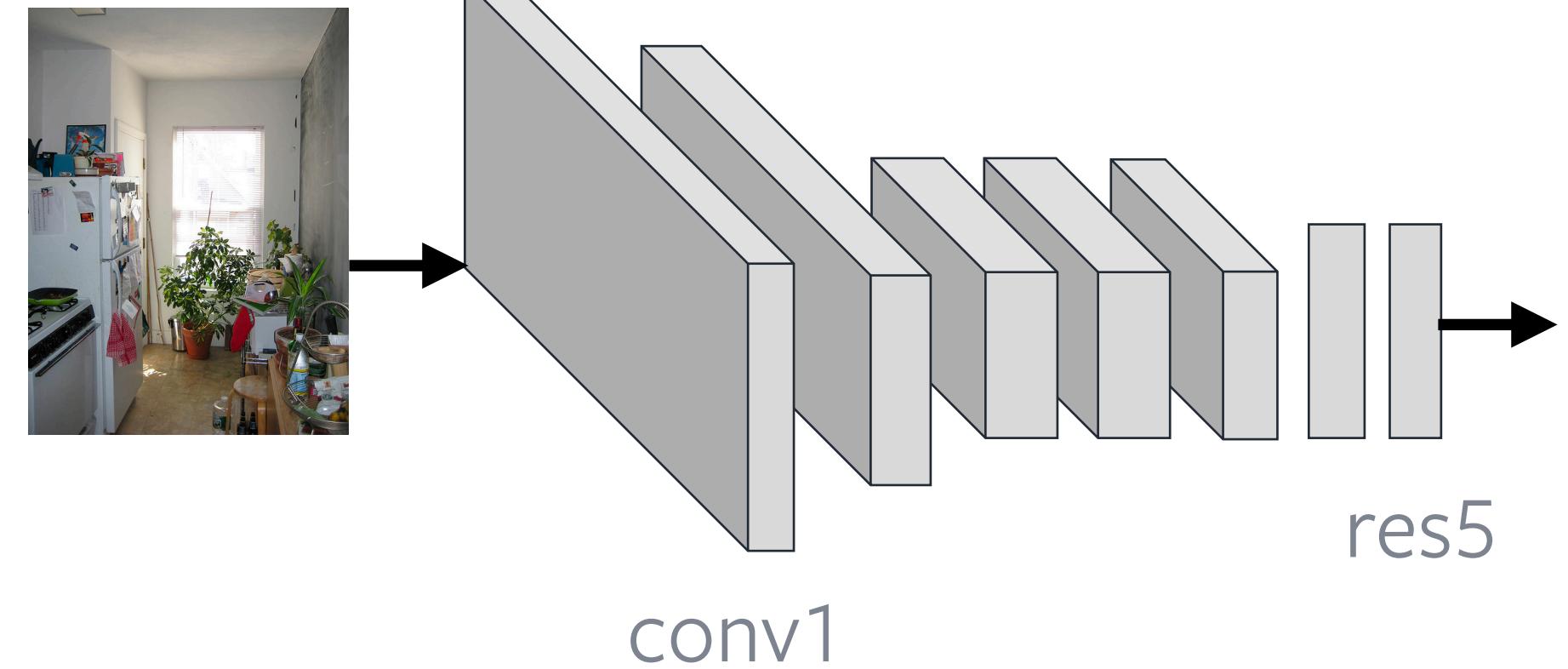
**Self-supervised representations are not as sample efficient**

# What does each layer learn?

**Linear classifier on VOC07**



mAP = mean Average  
Precision  
(Higher is better)

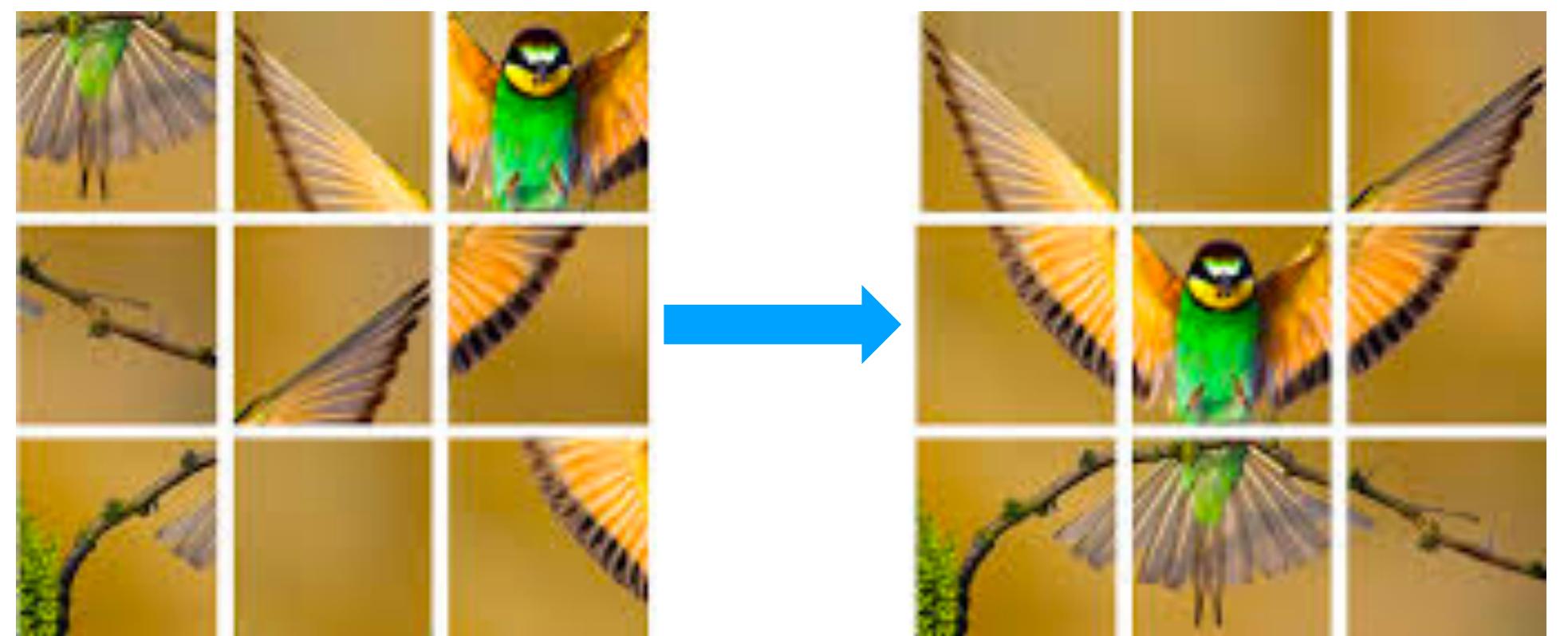


What is missing from “pretext” tasks?  
Or in general “proxy” tasks

# Pretext tasks



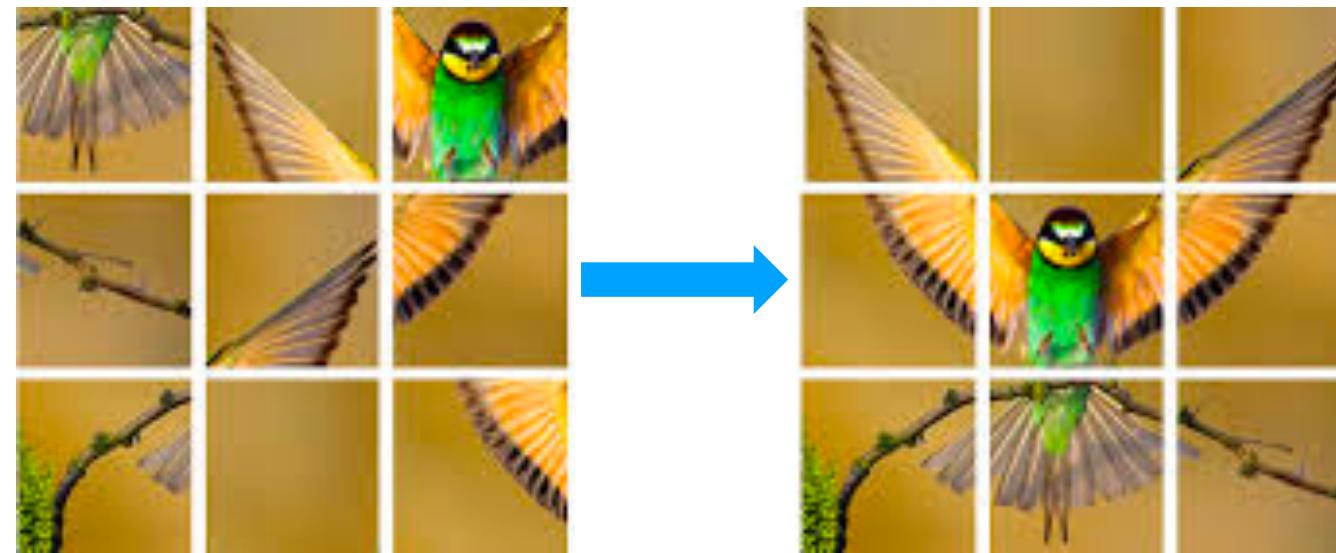
Rotation  
(Gidaris et al., 2018)



Jigsaw puzzles  
(Noroozi et al., 2016)

# The hope of generalization

- We really **hope** that the pre-training task and the transfer task are "aligned"



Pre-training  
Self-supervised



Transfer Tasks



# The hope of generalization

- We really hope that the pre-training task and the transfer task are "aligned"



Pre-training

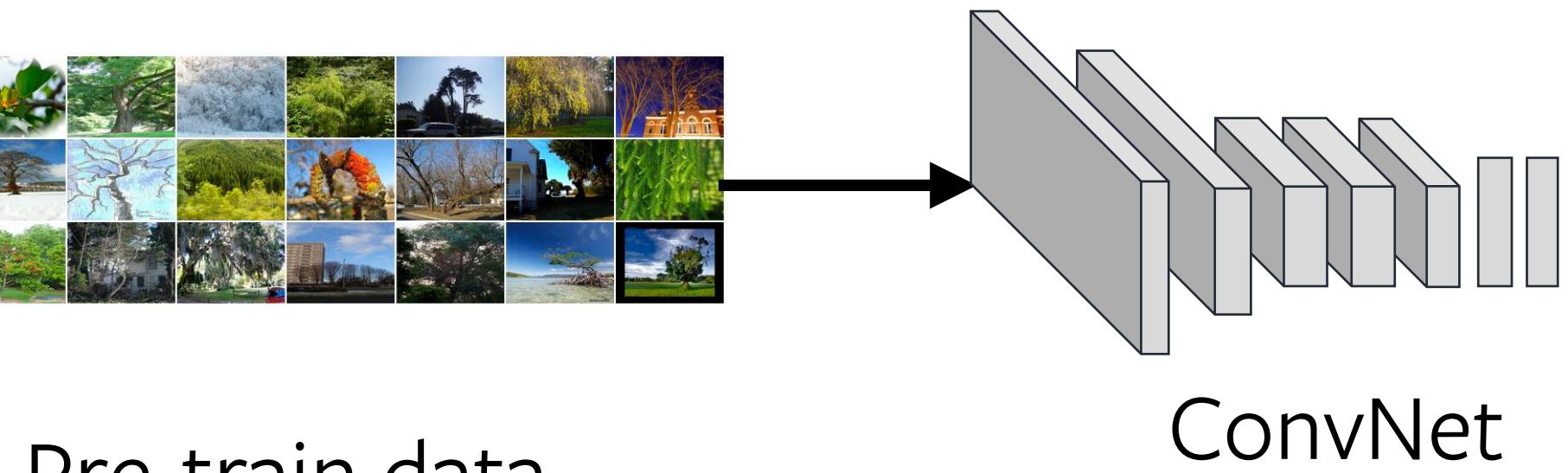
Weak or self-supervised

Transfer Tasks

Why should solving Jigsaw puzzles teach about "semantics"?

Why should performing a non semantic task produce good features?

# The hope of generalization ... ?



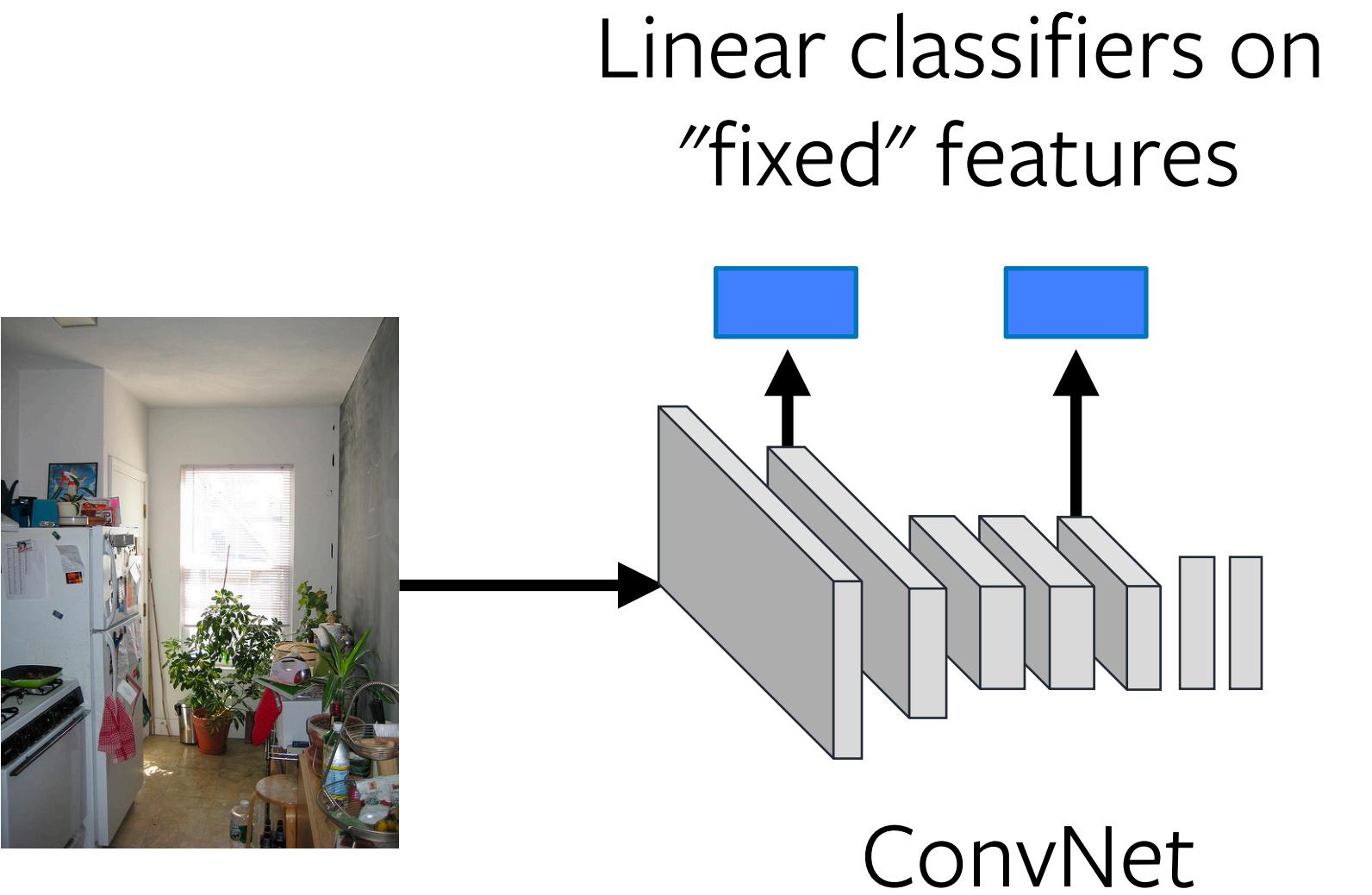
Pre-train data

ConvNet

Jigsaw

Pre-training

Weak or self-supervised

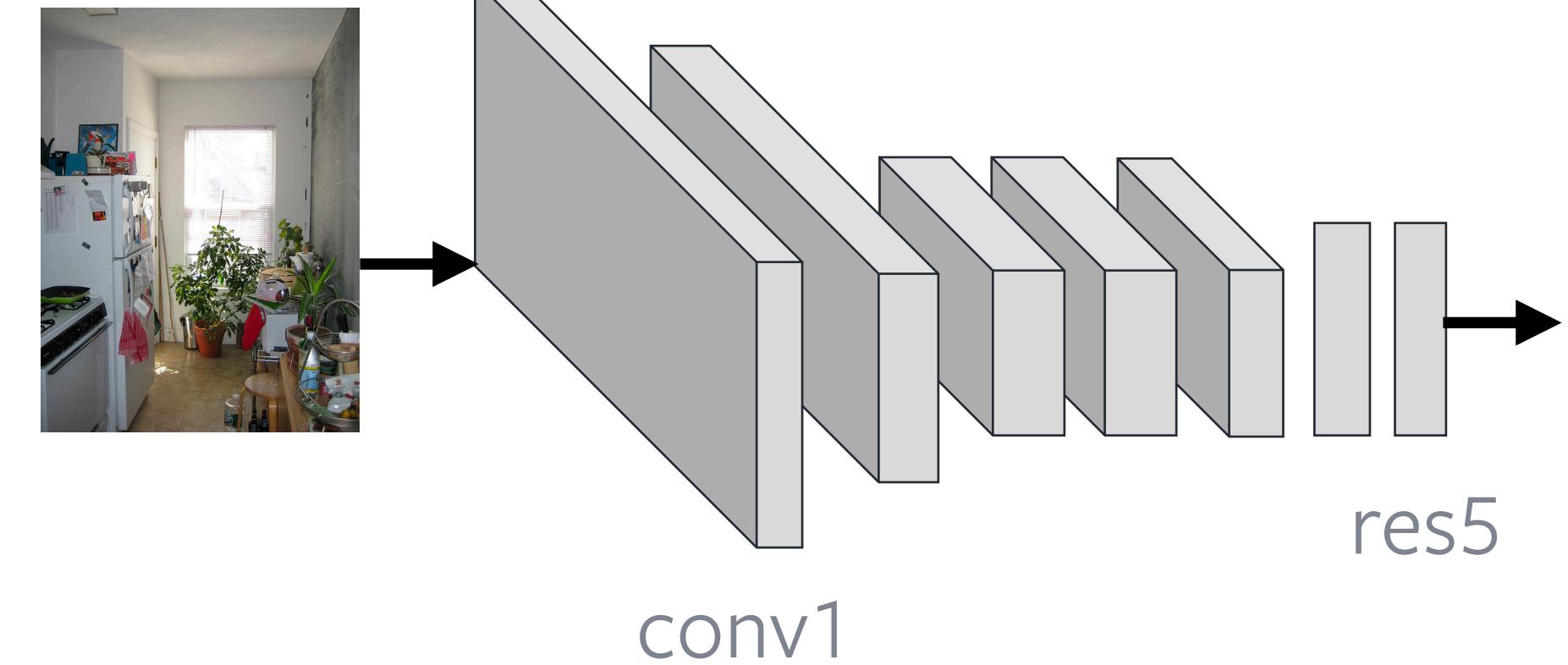
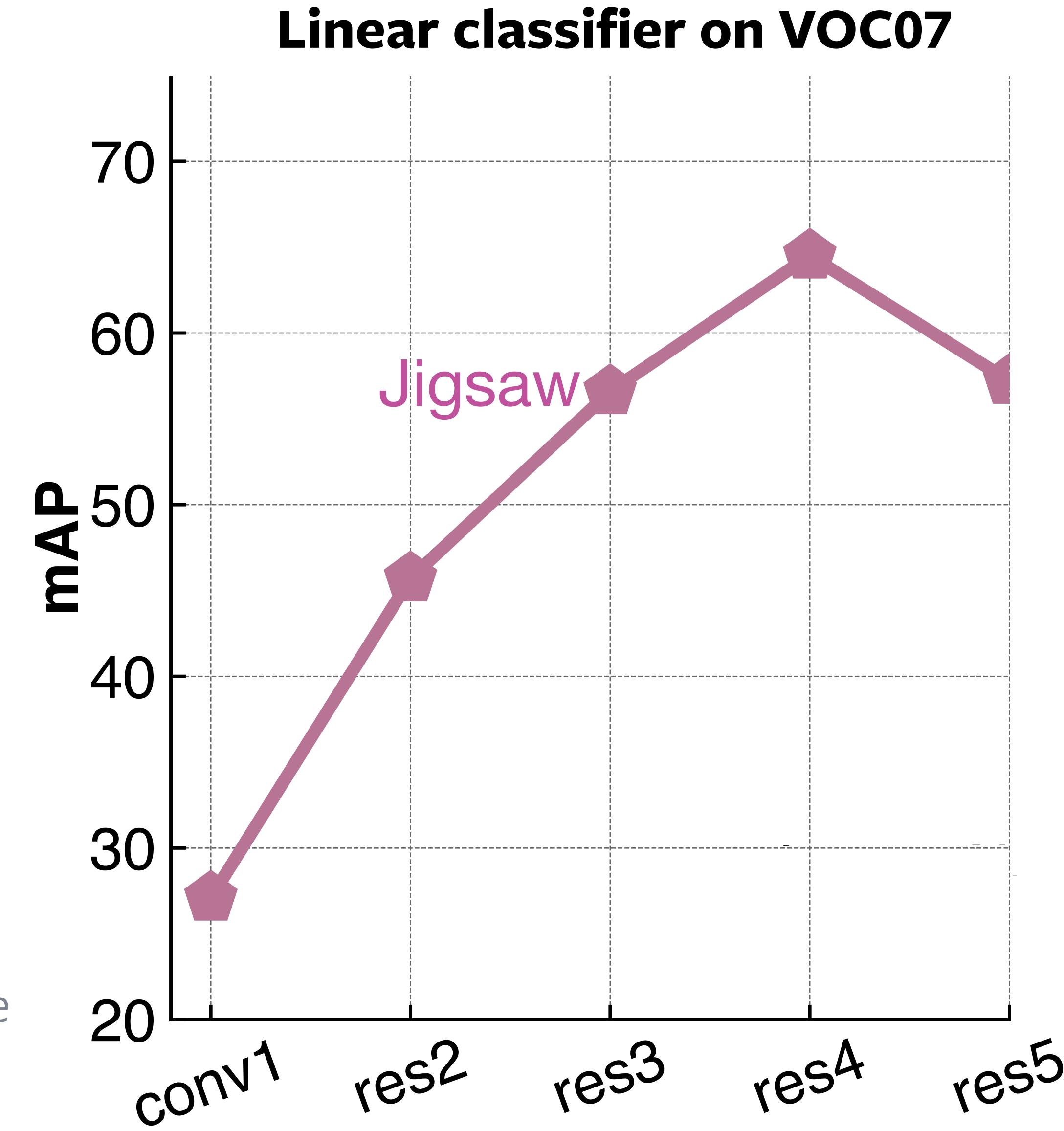


Linear classifiers on  
"fixed" features

ConvNet

Transfer

# Higher layers do not generalize ...



# Pre-trained features should ...

- Represent how images relate to one another
- Be robust to "nuisance factors" -- Invariance
  - e.g., exact location of objects, lighting, exact color

# Pre-trained features should ...

- Represent how images relate to one another
- Be robust to "nuisance factors" -- Invariance
  - e.g., exact location of objects, lighting, exact color

**Clustering and Contrastive Learning  
are two ways to achieve the above**

# Pre-trained features should ...

- Represent how images relate to one another
  - **ClusterFit:** Improving Generalization of Visual Representations
- Be robust to invariances, e.g., data transforms of the input
  - **PIRL:** Self-supervised learning of Pre-text Invariant Representations

See also:

DrLim (Hadsell et al. 2006), Exemplar CNN (Dosovitsky et al, 2014, Bautista et al., 2016), Visual Groups (Isola et al., 2016), Cluster (Xie et al., 2016), NAT (Bojanowski et al., 2017) Transitive (Wang et al., 2017), Boosting Knowledge (Noroozi et al., 2018), DeepCluster (Caron et al., 2018, 2019), CPC (Oord et al., 2018, Henaff et al., 2019), CMC (Tian et al., 2019) MoCo (He et al., 2019)

# ClusterFit: Improving Generalization of Visual Representations

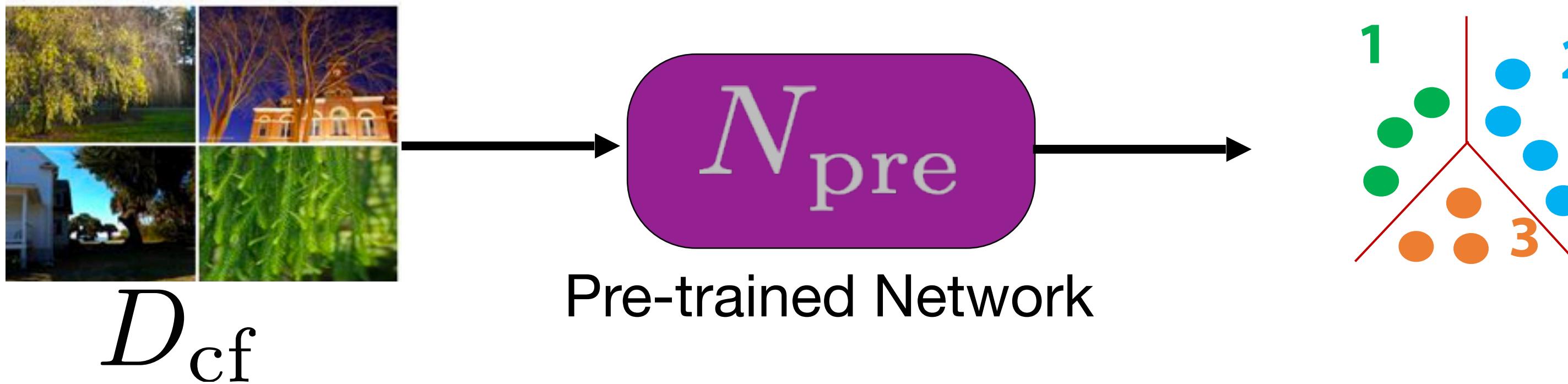
Yan\*, Misra\* et al., CVPR 2020

# Understanding how images relate to each other

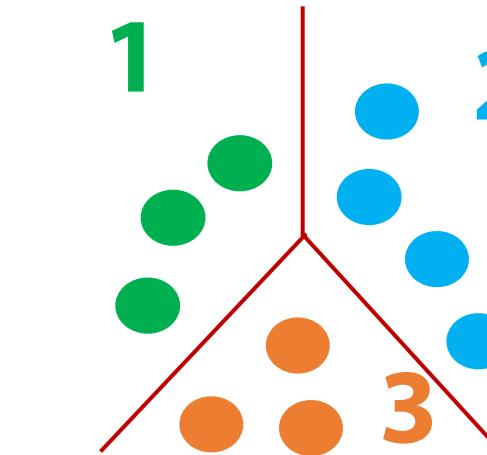
- Clustering the feature space is a way to see what images relate to one another

# Main Idea

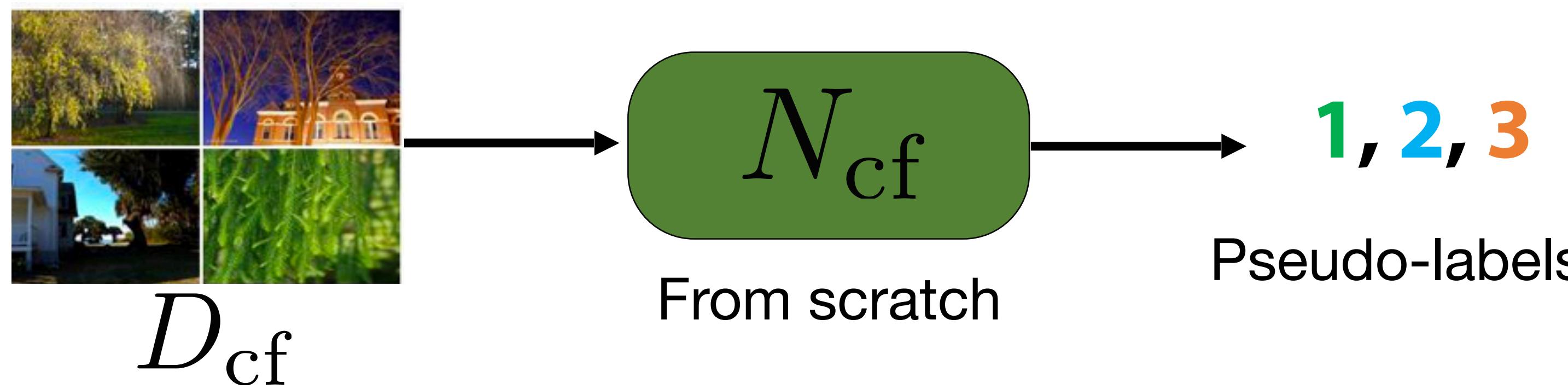
## 1. Cluster: Feature clustering



Pre-trained Network



## 2. Fit: Predict Cluster Assignments



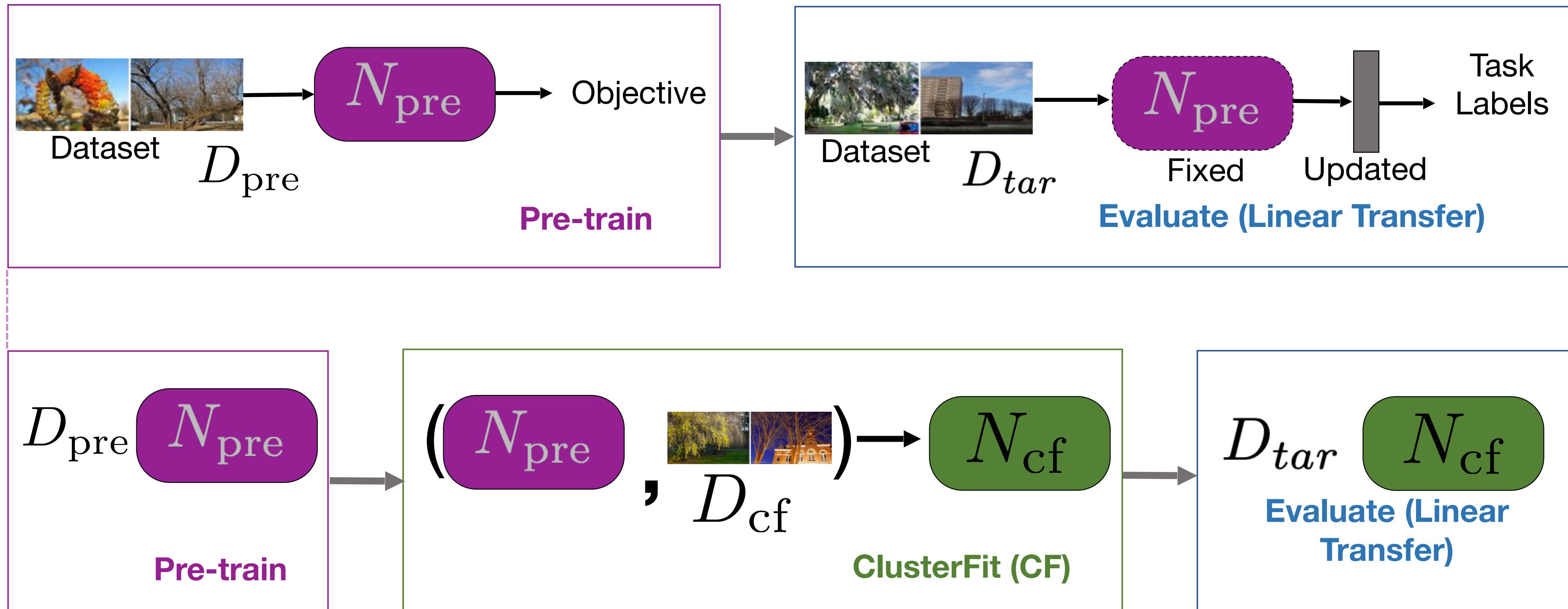
Pseudo-labels

From scratch

See also: Boosting Knowledge (Noroozi et al., 2018); DeepCluster (Caron et al., 2018); DeeperCluster (Caron et al., 2019)

# Applied to any pre-trained network

## “Standard” Pre-train and Transfer



## “Standard” Pre-train + ClusterFit

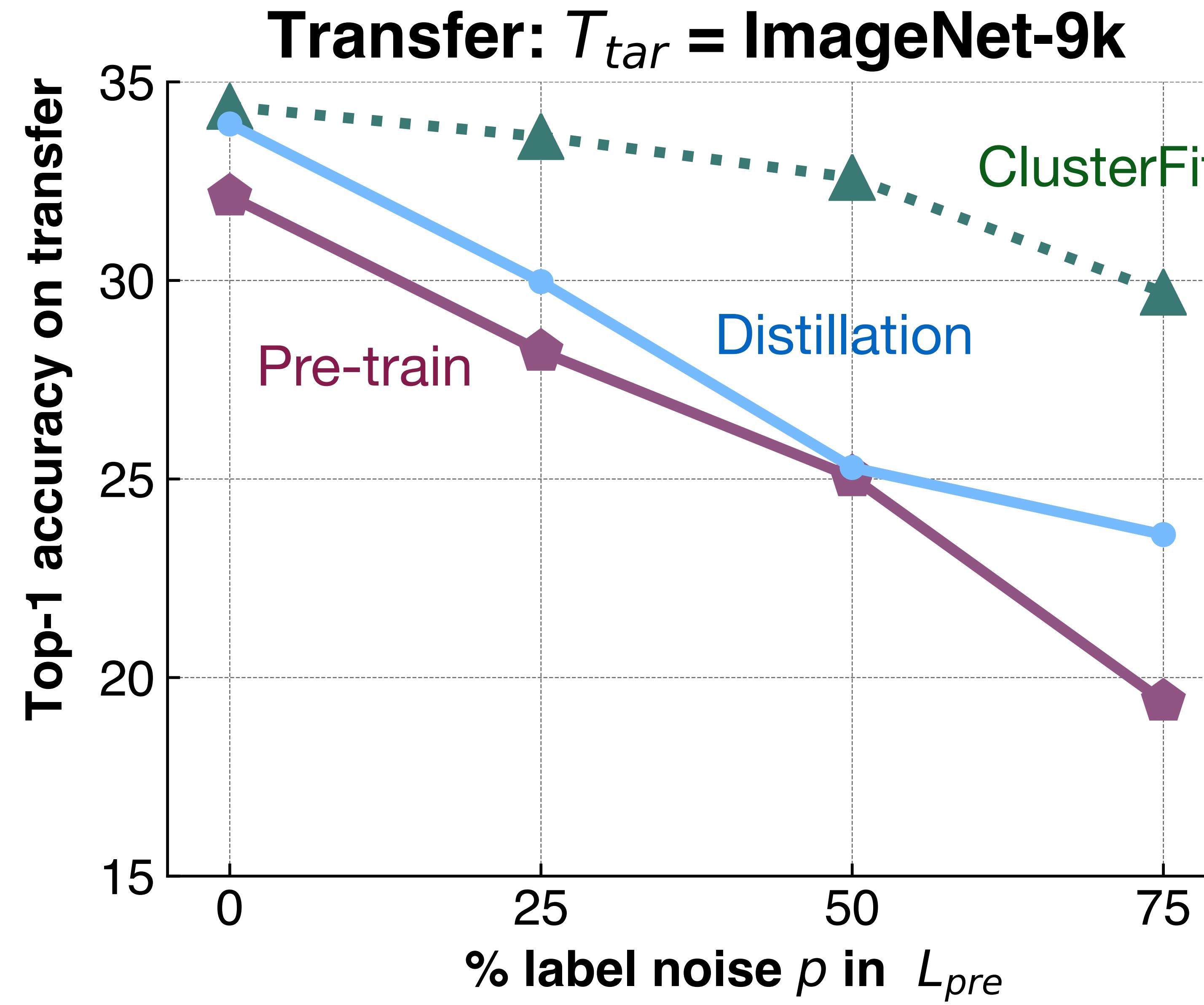
# Why does it work?

- Clustering introduces a "bottleneck" and only the most important information is transferred

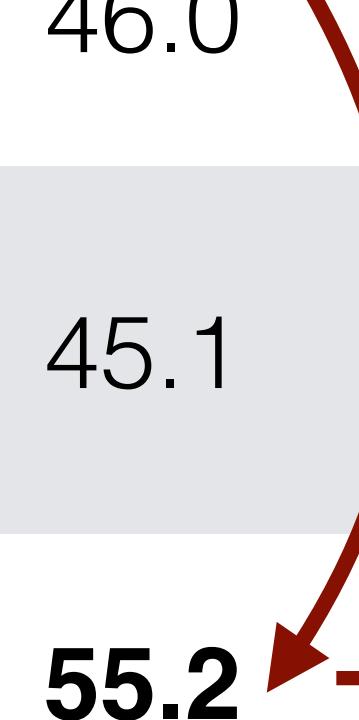
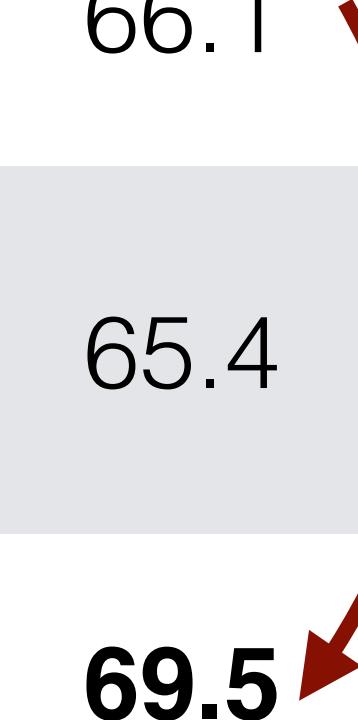
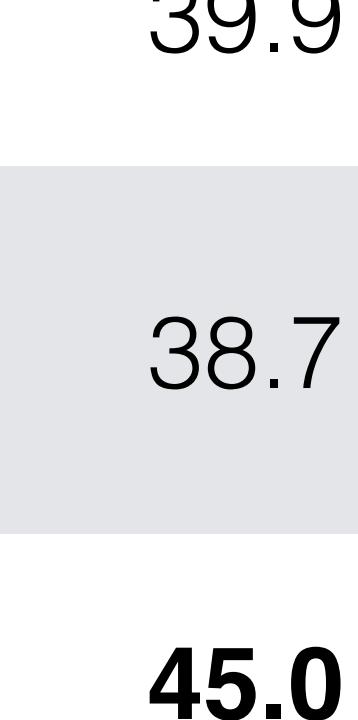
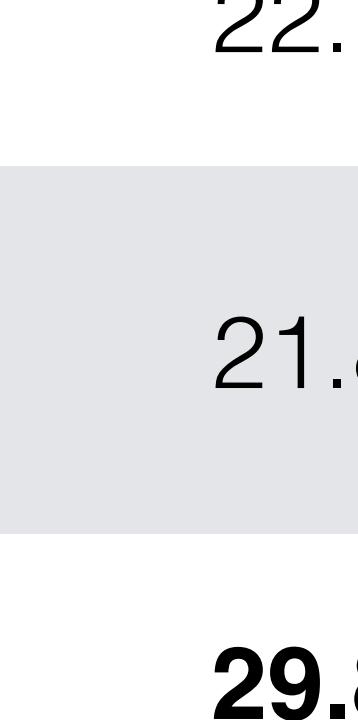
# Synthetic noise

- Add label noise to ImageNet-1K (1000 classes)
- Measure transfer performance on ImageNet-9K (9000 classes)

# Synthetic noise



# Self-supervised Images

Method	Transfer Dataset			
	ImageNet-1M	VOC07	Places205	iNaturalist
Jigsaw Pretrain	46.0	66.1	39.9	22.1
Pretrain 2x	45.1	65.4	38.7	21.8
ClusterFit	<b>55.2</b>  <b>+9</b>	<b>69.5</b>  <b>+3</b>	<b>45.0</b>  <b>+5</b>	<b>29.8</b>  <b>+7</b>

ResNet50, ImageNet-1M  
**Similar gains for RotNet**

# Applied to any pre-trained network

Pre-training Setup	Network	Dataset	Gain of ClusterFit	
			Top-1 Accuracy	
Fully Supervised	ResNet50	ImageNet1K	2.1% on ImageNet-9K	
Weakly Supervised Images	ResNet50	IG 1B	4.6% on ImageNet-9K 5.8% on iNaturalist	
Weakly Supervised Videos	R(2+1)D-34	IG 19M	3.2% on Kinetics 4.3% on Sports1M	
Self-supervised (Jigsaw, RotNet)	ResNet50	ImageNet1K	7-9% on ImageNet-1K 3-5% on Places-205	

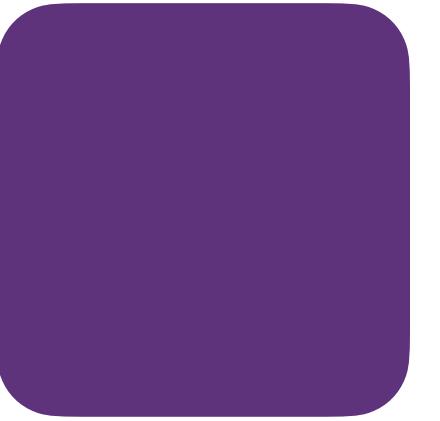
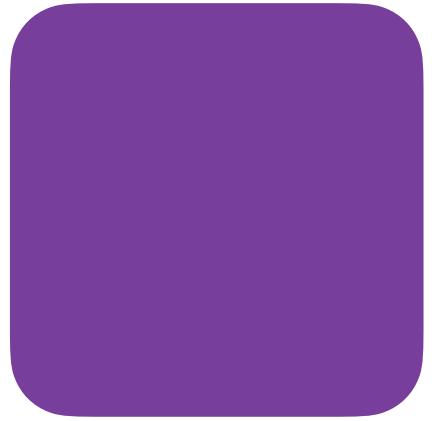
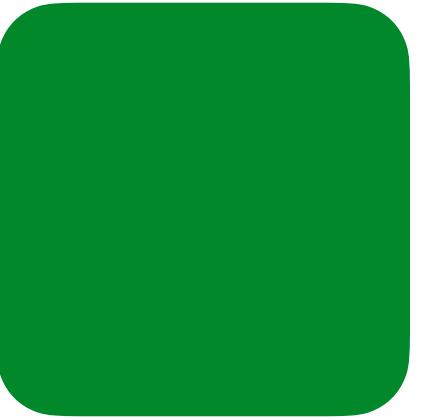
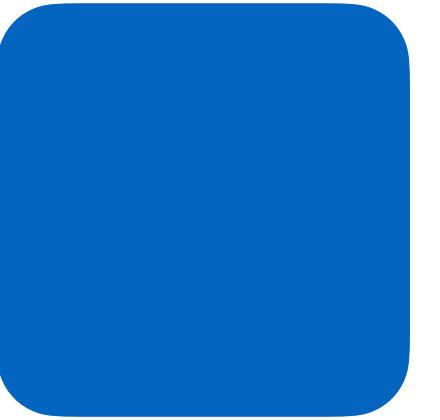
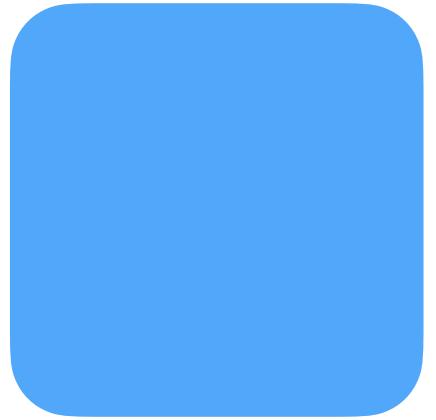
**Gains without extra data, labels or changes in architecture**

# Self-supervised Learning of Pretext Invariant Representations (PIRL)

Misra & van der Maaten, CVPR 2020

# Contrastive Learning

Groups of  
Related and Unrelated  
Images

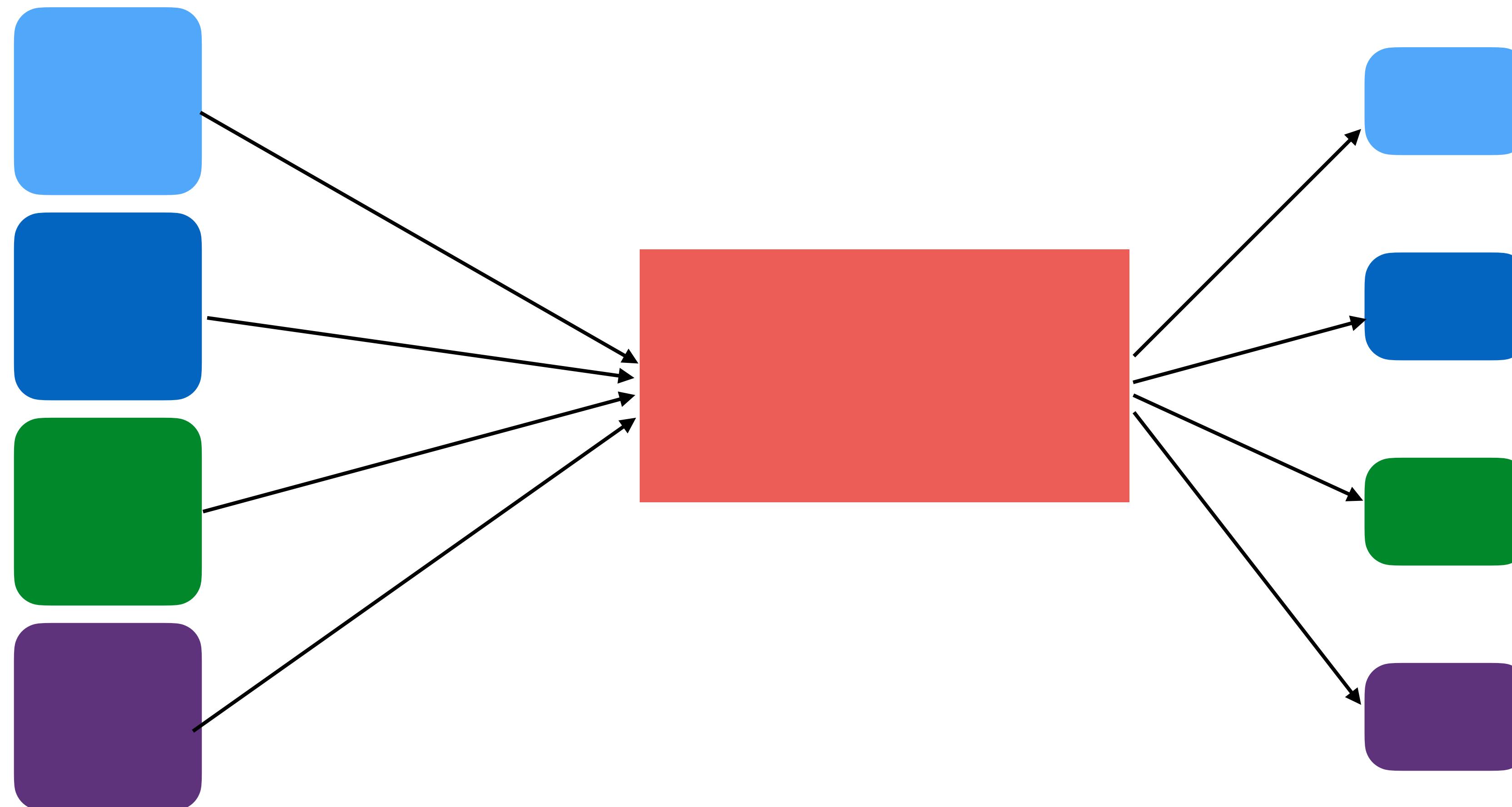


# Contrastive Learning

Groups of  
Related and Unrelated  
Images

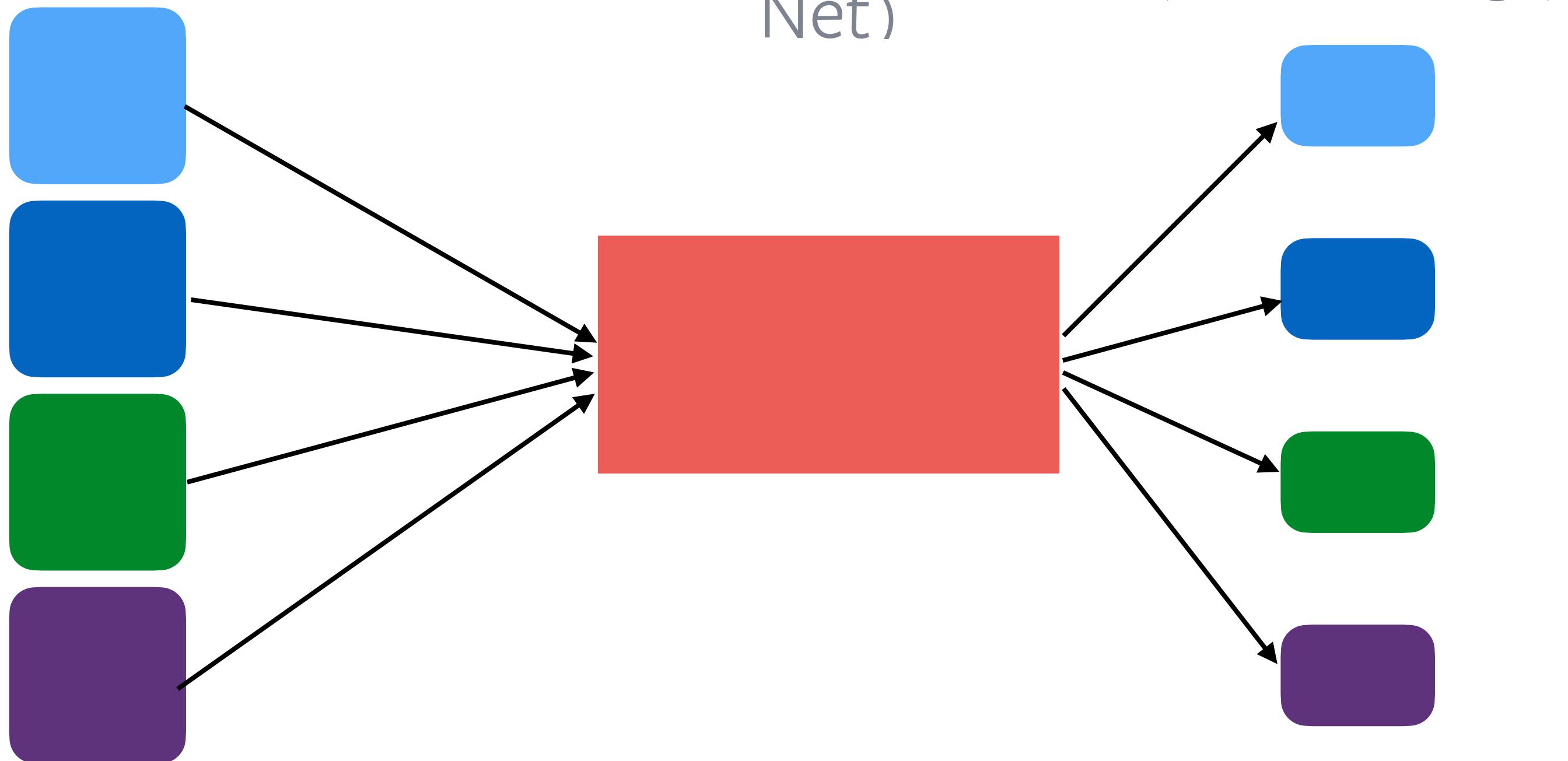
Shared network  
(Siamese Net)

Image Features  
(Embeddings)



# Contrastive Learning

Related and  
Unrelated  
Images



Shared  
network  
(Siamese  
Net)

Image  
Features  
(Embeddings)

## Loss Function

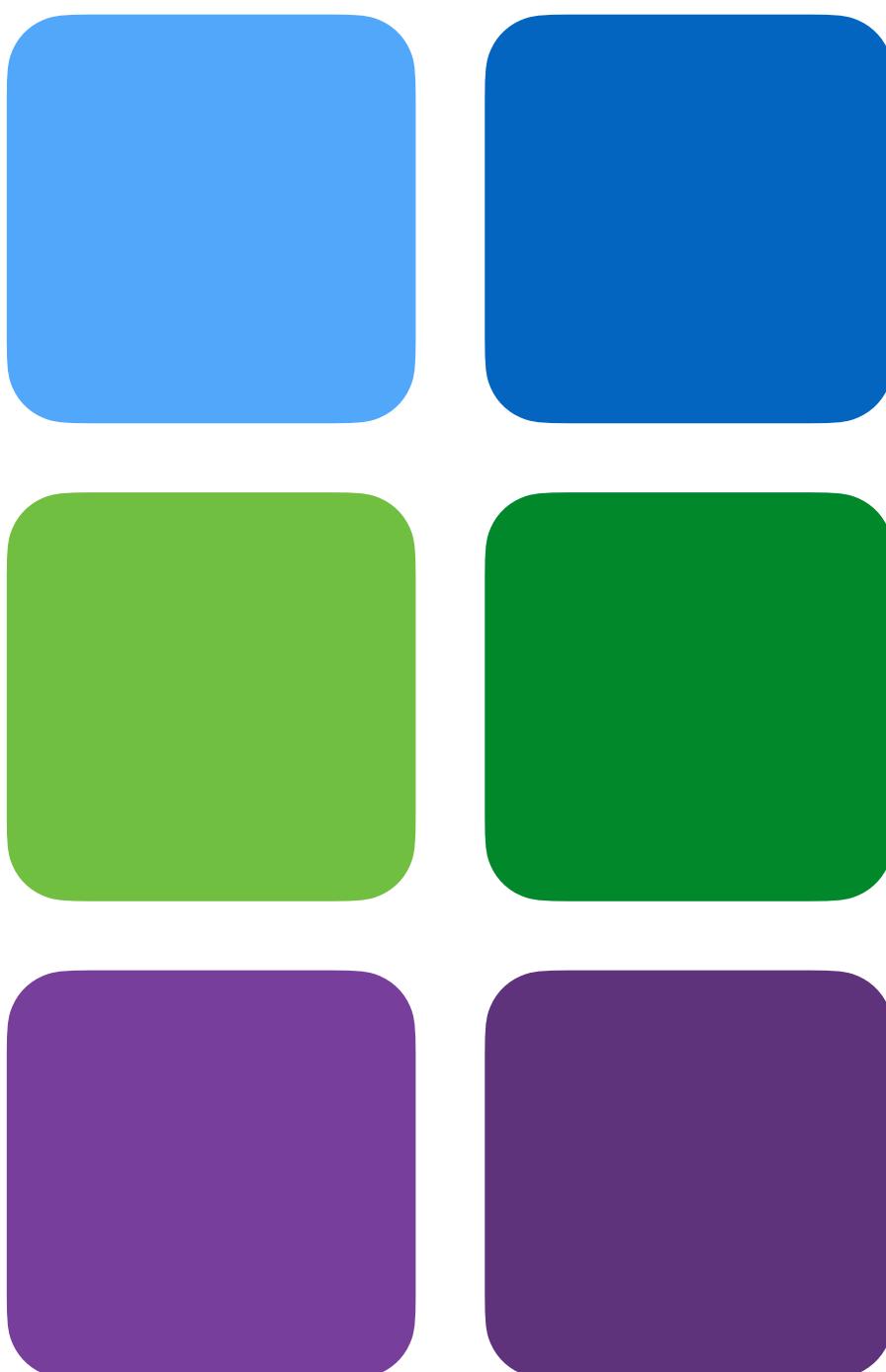
Embeddings from related images should be closer than embeddings from unrelated images

$$d(\text{blue}, \text{blue}) < d(\text{blue}, \text{green})$$
$$d(\text{blue}, \text{blue}) < d(\text{blue}, \text{purple})$$

# Contrastive Learning

- How to define what images are "related" and "unrelated"?

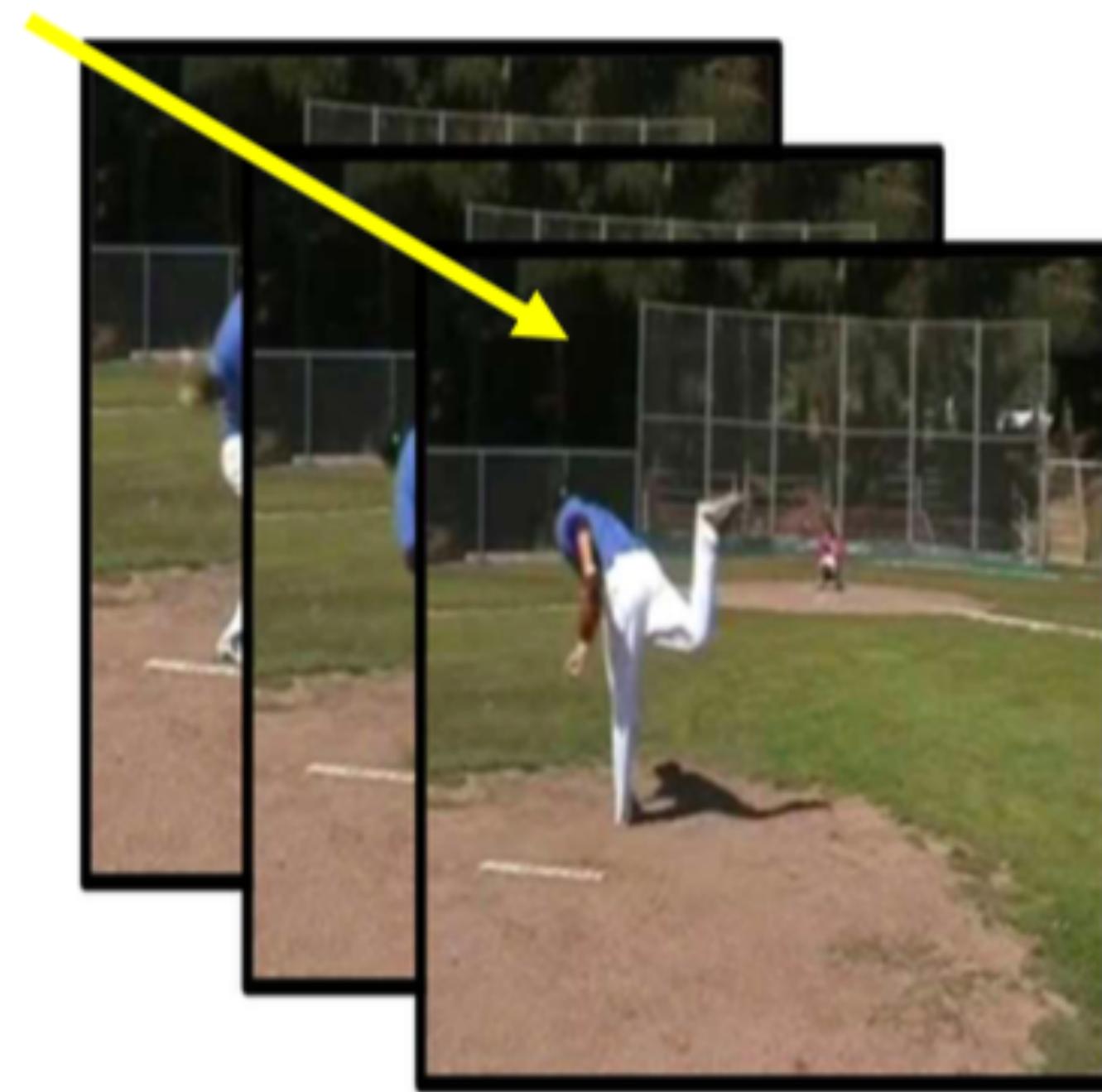
Related and Unrelated  
Images



# Frames of a video

# Video & Audio

Time



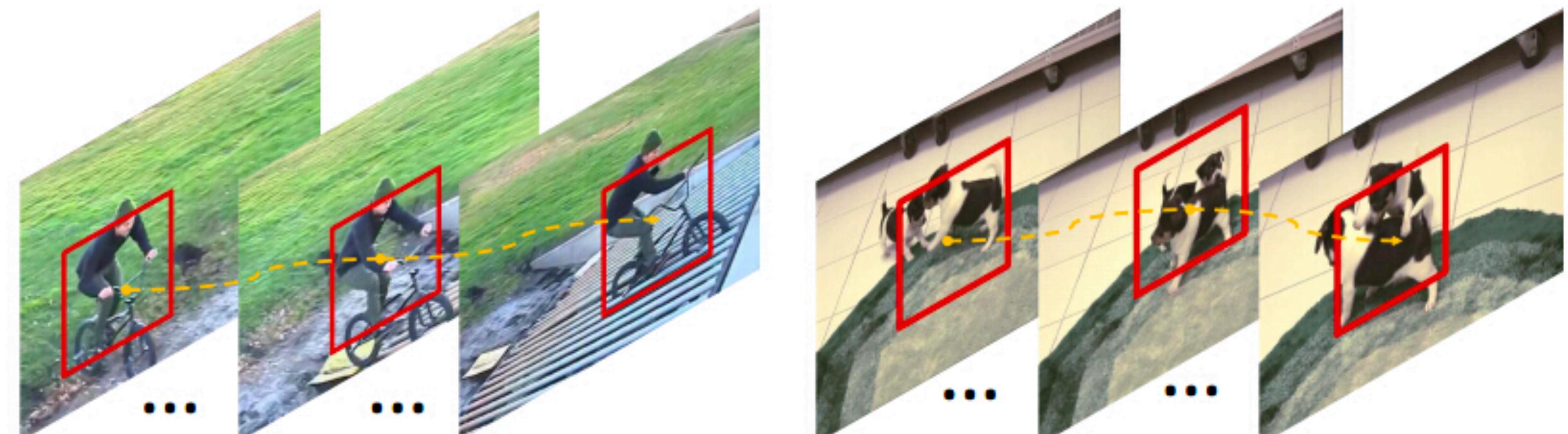
"Sequence" of data

Hadsell et al., 2005, DrLim  
van der Oord et al., 2018, CPC

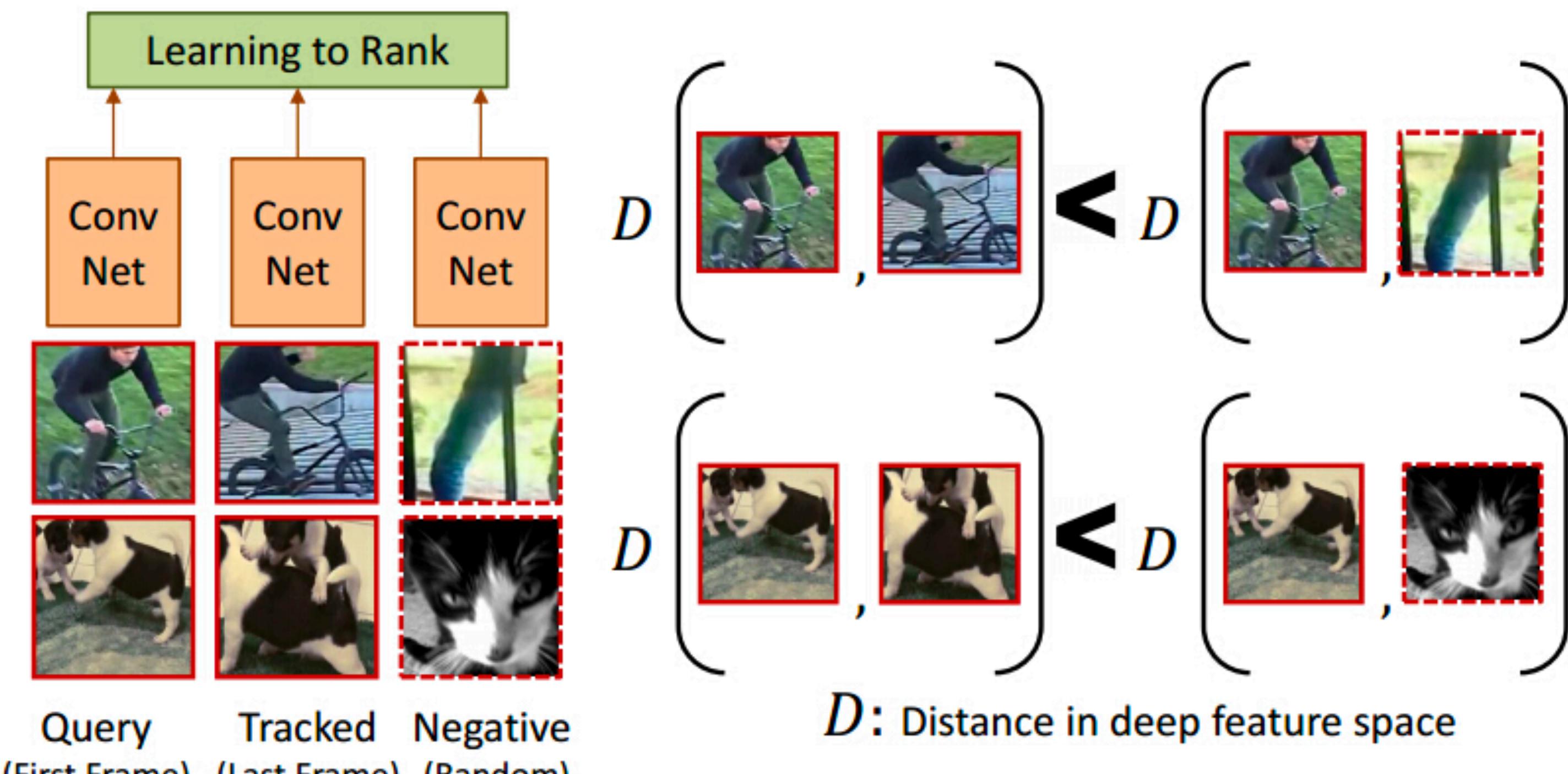


Coming soon

# Tracking Objects



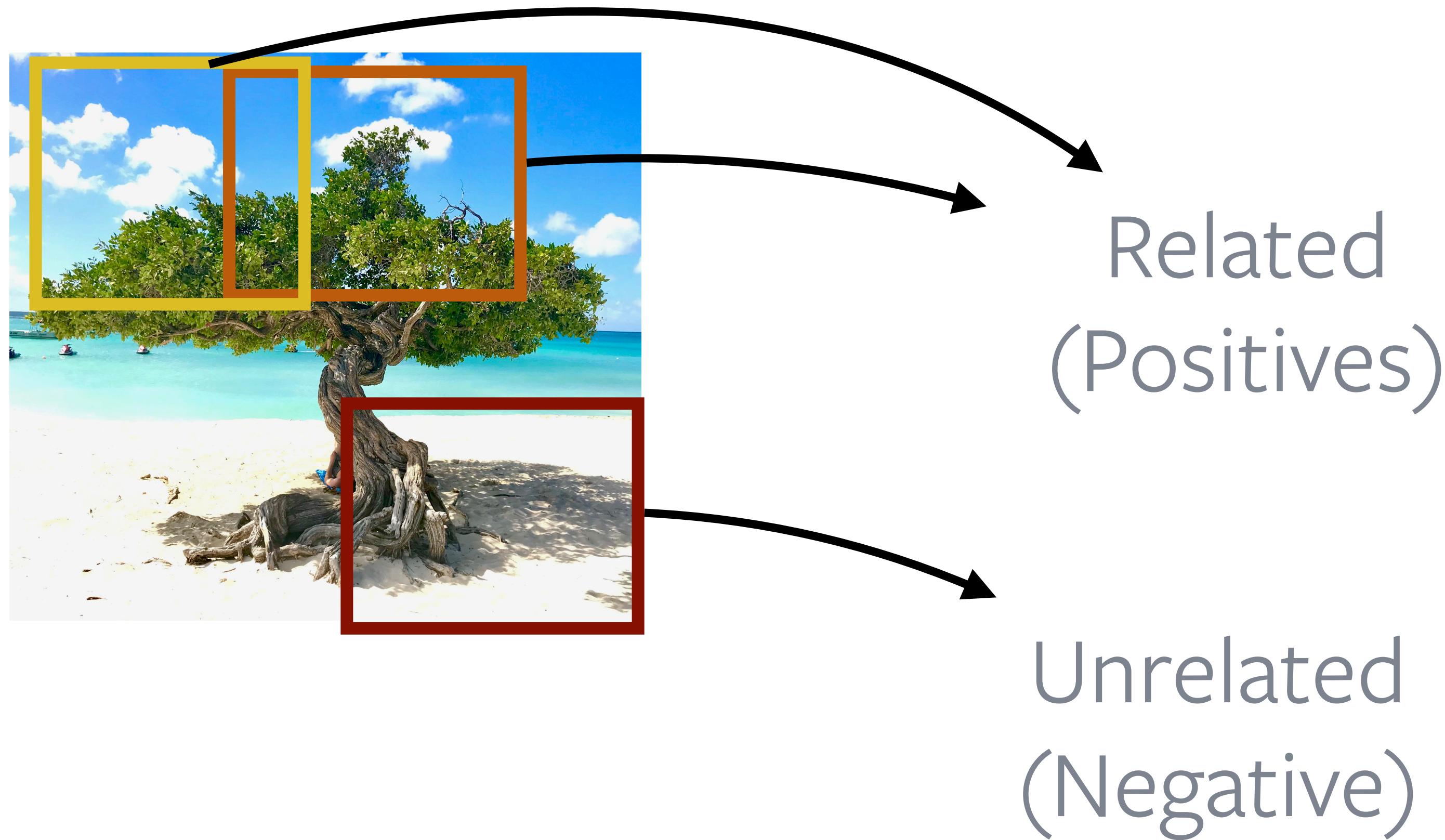
### (a) Unsupervised Tracking in Videos



### (b) Siamese-triplet Network

### (c) Ranking Objective

# Nearby patches vs. distant patches of an Image



van der Oord et al., 2018,  
Henaff et al., 2019  
Contrastive Predictive Coding

# Patches of an image vs. patches of other images



Related  
(Positives)

Wu et al., 2018, Instance Discrimination

He et al., 2019, MoCo

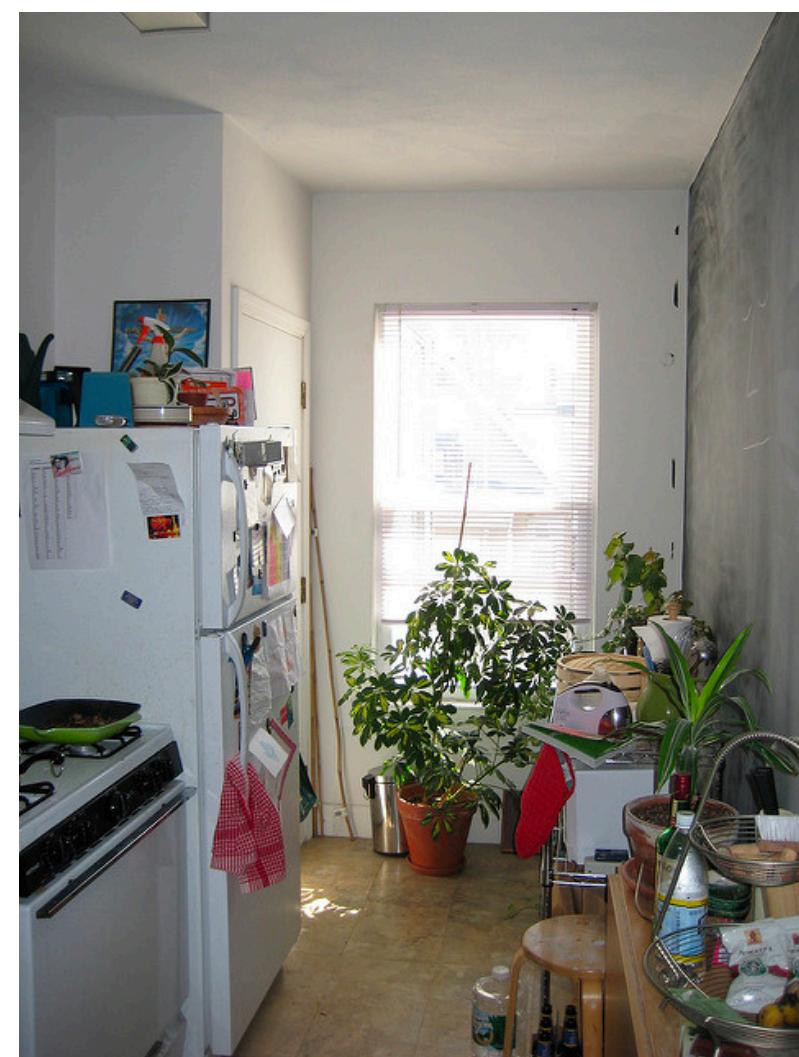
Misra & van der Maaten, 2019, PIRL

Chen et al., 2020, SimCLR



Unrelated  
(Negative)

# Data Augmentations of each patch



Unrelated  
(Negative)

# Underlying Principle for Pretext Tasks

- Apply known image transform  $t$
- Construct task to predict  $t$  from transformed Image ( $I^t$ )
- Final layer representations must carry information about  $t$
- Representations "covary" with  $t$

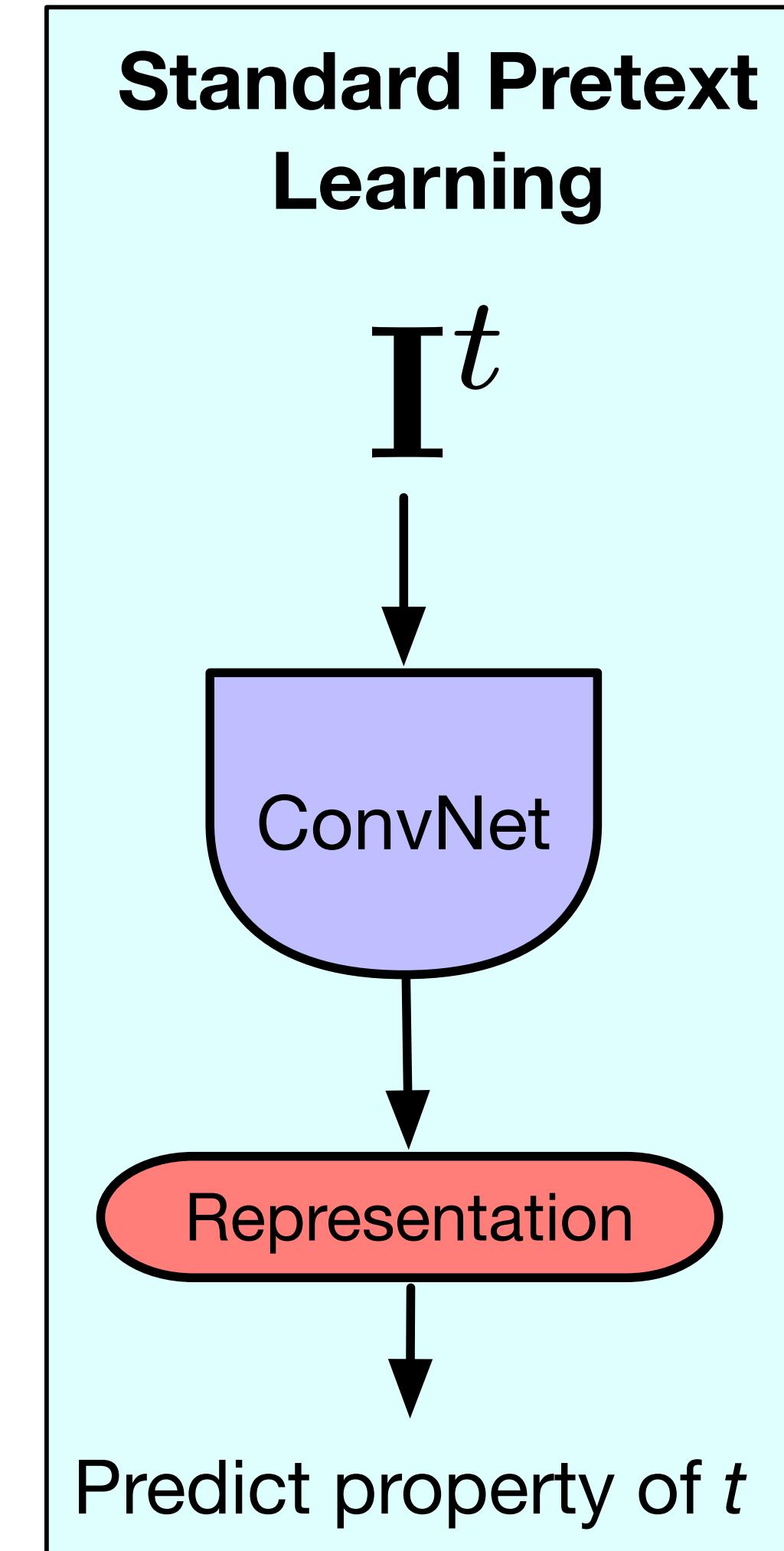
**Pretext Image Transform**



$I$   
Transform  $t$   
 $I^t$



**Standard Pretext Learning**



# Underlying Principle for Pretext Tasks

- Apply known image transform  $t$
- Construct task to predict  $t$  from transformed Image ( $I^t$ )

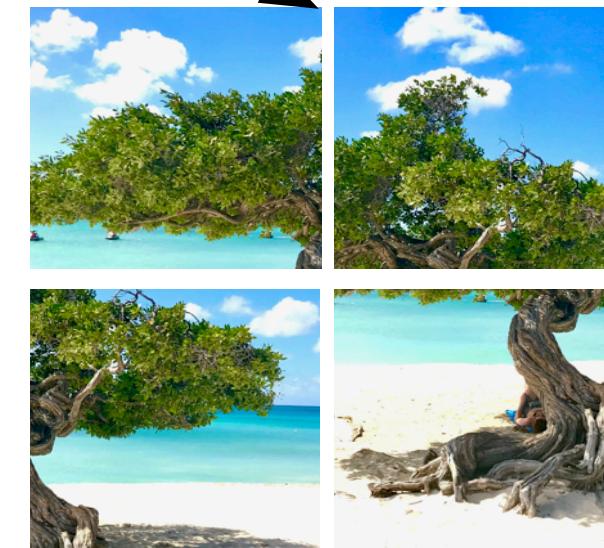
- Final layer representations must carry information about  $t$
- Representations "covary" with  $t$

**But .... shouldn't representations be invariant to low-level image transforms?**

**Pretext Image Transform**

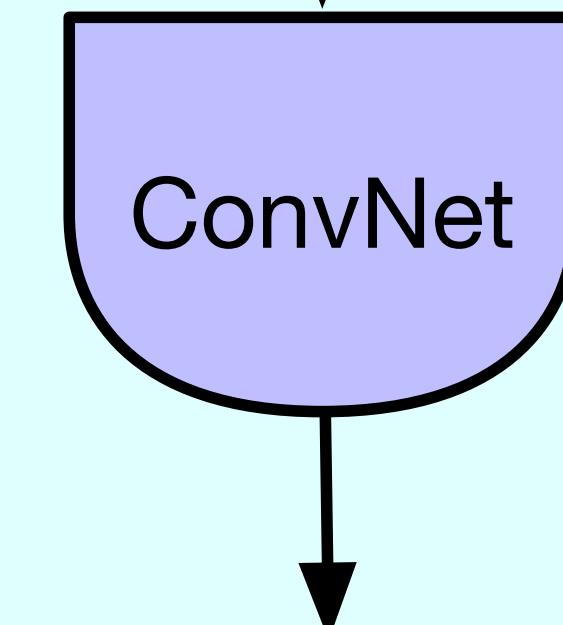


$I$   
Transform  $t$   
 $I^t$



**Standard Pretext Learning**

$I^t$



Representation

Predict property of  $t$

# How important has invariance been?

- Hand-crafted features like SIFT and HOG
- SIFT - Scale **Invariant** Feature Transform
- Supervised systems are trained to be invariant to "data augmentation"



# Pretext-Invariant Representation Learning (PIRL)

- Be invariant to  $t$

**Pretext Image Transform**



$I$   
Transform  $t$



**Standard Pretext Learning**

$I^t$

ConvNet

Representation

Predict property of  $t$

**Pretext Invariant Representation Learning**

$I$

ConvNet

Representation

$I^t$

ConvNet

Representation

Encourage to be similar

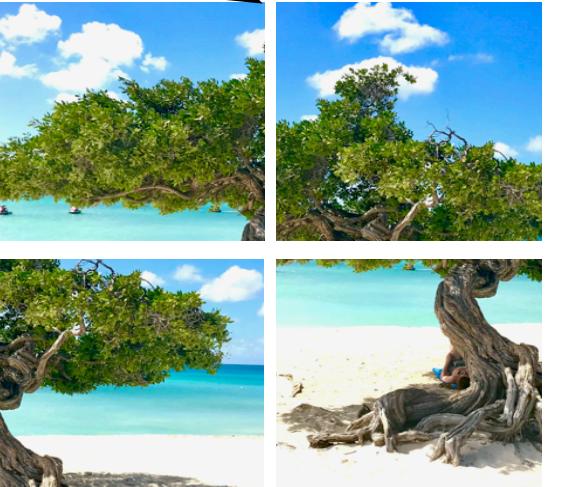
# Pretext-Invariant Representation Learning (PIRL)

- Be invariant to  $t$
- Representation contains no information about  $t$

**Pretext Image Transform**



$I$   
Transform  $t$



$I^t$

**Standard Pretext Learning**

$I^t$

ConvNet

Representation

Predict property of  $t$

**Pretext Invariant Representation Learning**

$I$

ConvNet

Representation

$I^t$

ConvNet

Representation

Encourage to be similar

# PIRL

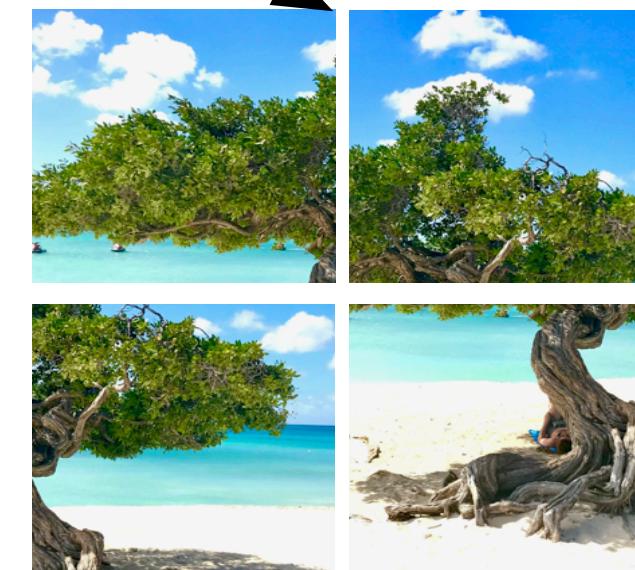
- Representations from  $\mathbf{I}$  and  $\mathbf{I}^t$  should be similar
- $\mathbf{t}$  = Pretext Transforms (Jigsaw/ Rotation, combinations etc.)
- Use a contrastive loss to enforce similarity of features

$$L_{\text{contrastive}}(\mathbf{v}_{\mathbf{I}}, \mathbf{v}_{\mathbf{I}^t})$$

## Pretext Image Transform

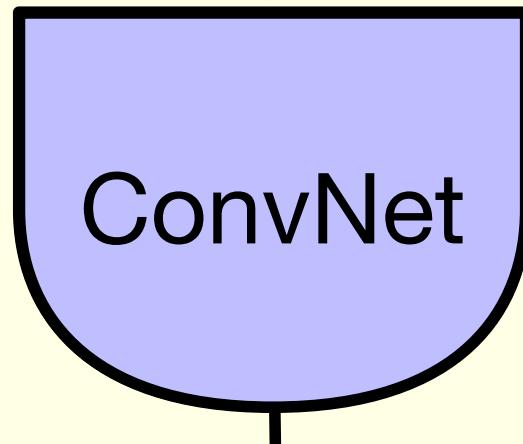


$\mathbf{I}$   
Transform  $t$   
 $\mathbf{I}^t$

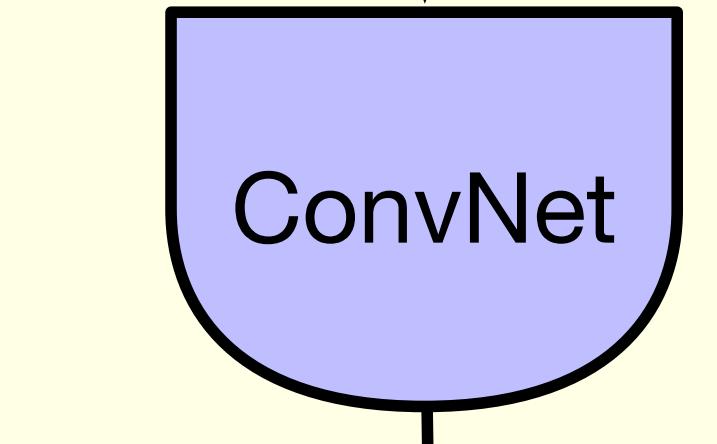


## Pretext Invariant Representation Learning

$\mathbf{I}$



$\mathbf{I}^t$



Representation

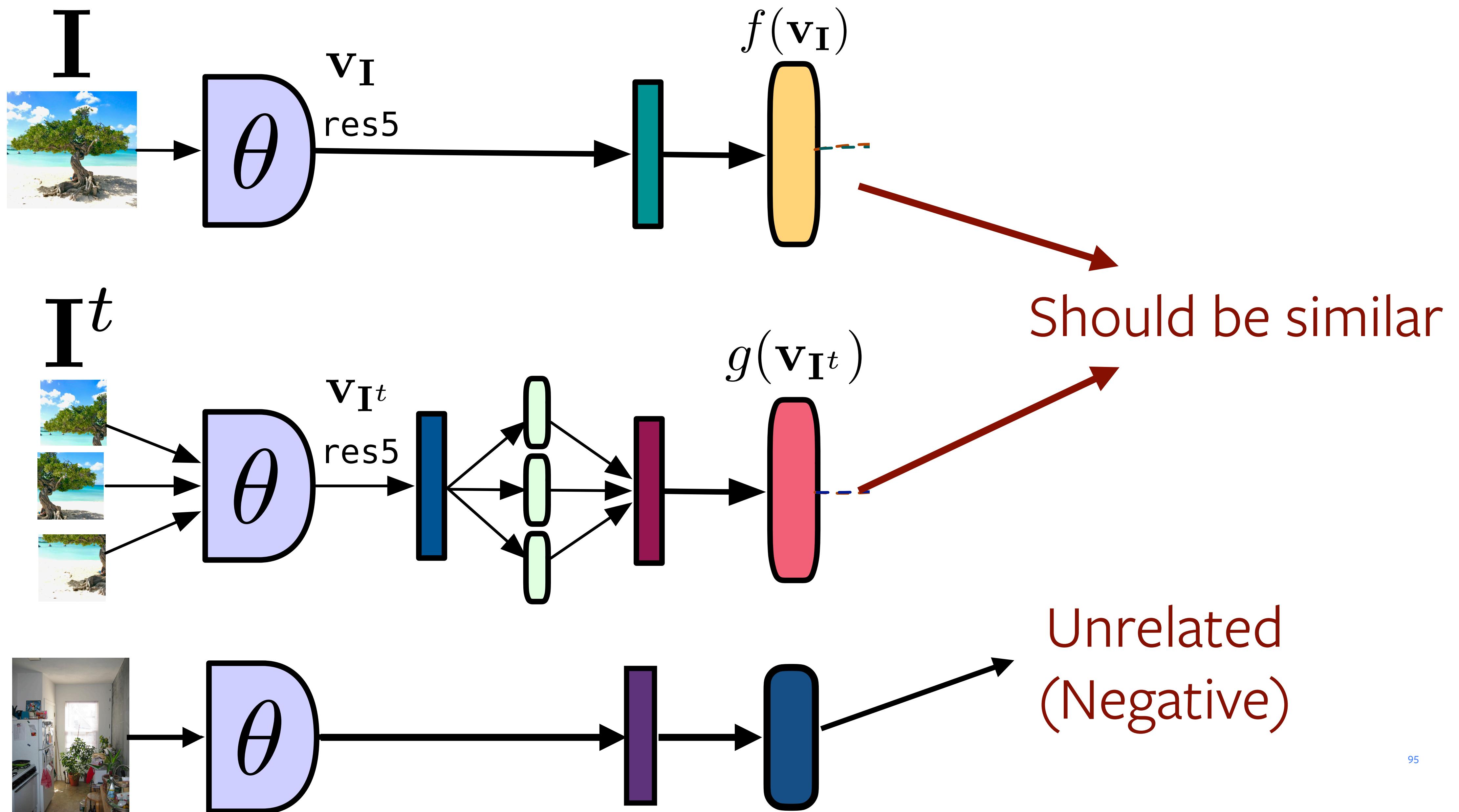
Representation

Encourage to be similar

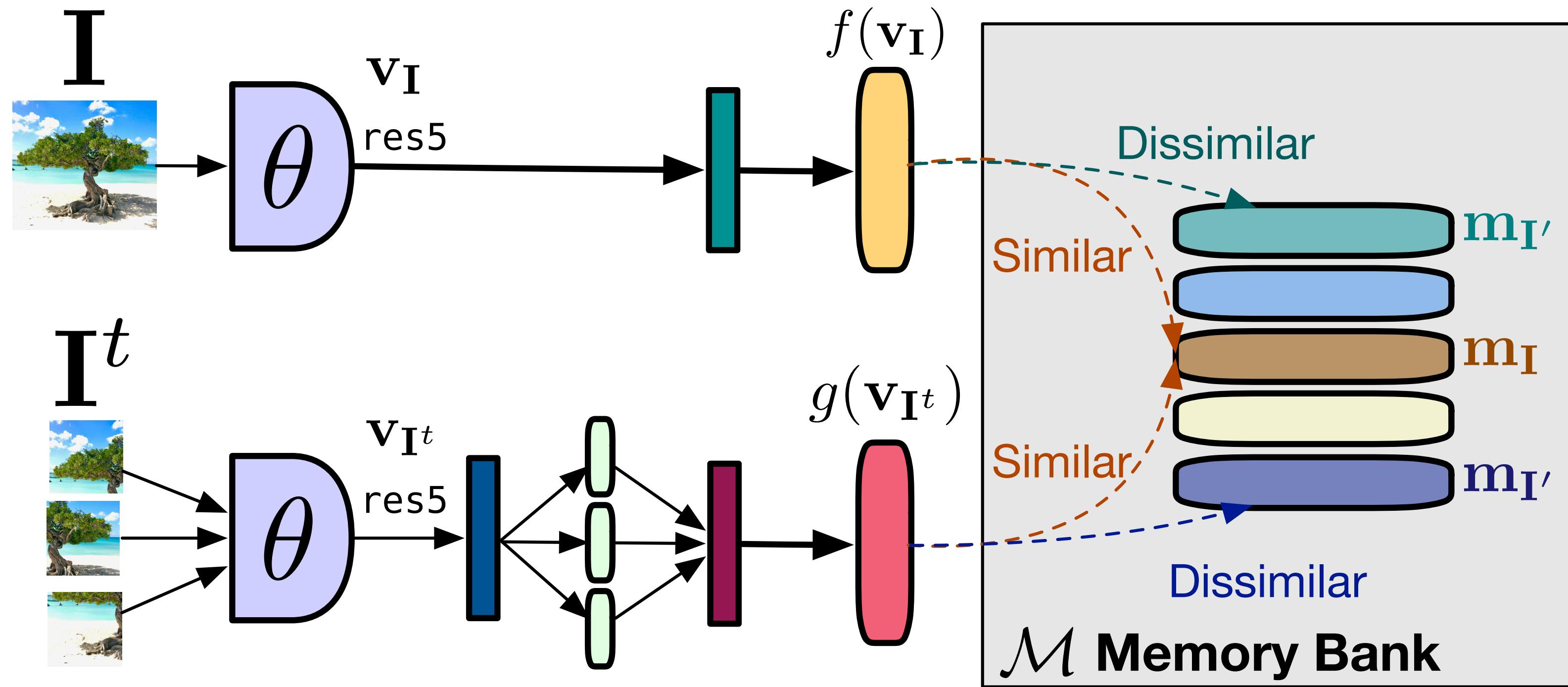
# PIRL - Using Large Number of Negatives

- Large number of negatives is crucial in contrastive learning
- Following Wu et al., PIRL uses a **memory bank** of negatives
- Memory bank stores a feature vector per image in the dataset

# How it works



# How it works



$$L_{\text{NCE}}(g(v_{I^t}), m_I) + L_{\text{NCE}}(f(v_I), m_I)$$

Rep.  $I^t$  and  $m_I$  should be similar

Rep.  $I$  and  $m_I$  should be similar

# PIRL Pre-training

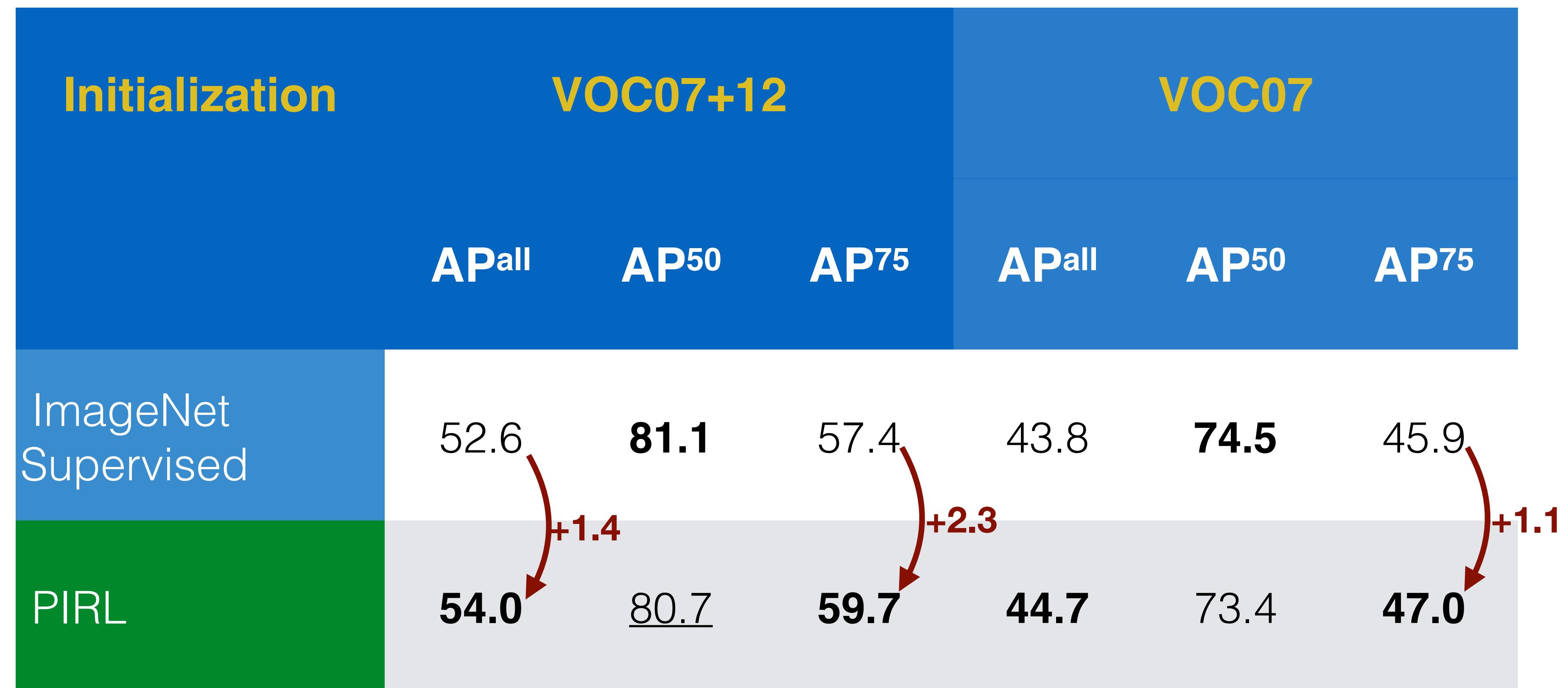
- Pre-trained on images without labels
  - ImageNet-1M
  - YFCC ("in-the-wild"/uncurated images)

# Evaluation

- Evaluate initialization: Full fine-tuning
- Evaluate features: Train linear classifiers on fixed features

# Object Detection

- **Outperforms** ImageNet supervised pre-trained networks
- Full fine-tuning, no bells & whistles
- No extra data, changes in model architecture, fine-tuning schedule



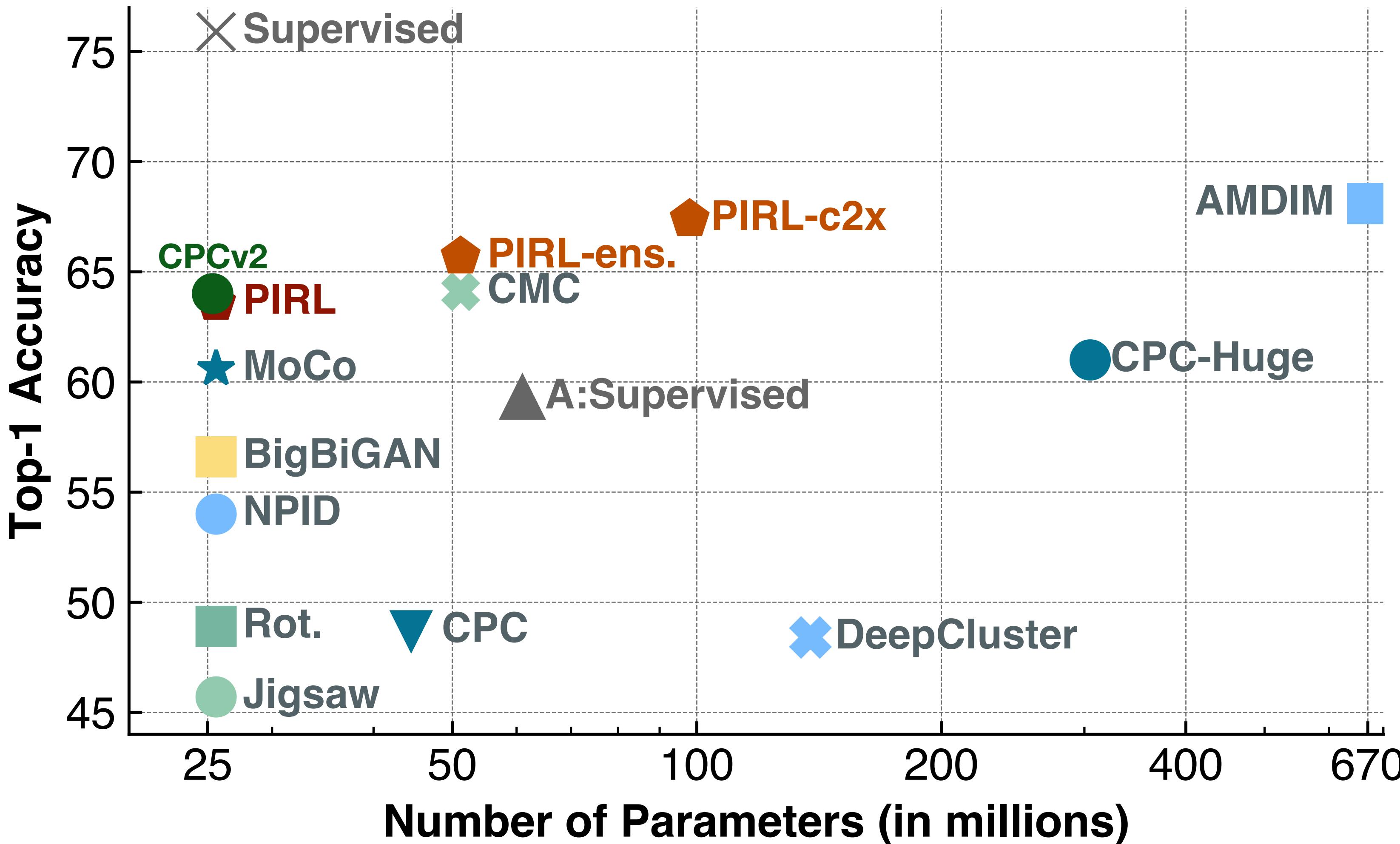
# Semi-supervised Learning

- Fine-tune on fraction of labeled data from **ImageNet-1K**

Method	Top-5 Accuracy	
Fraction of Data	→ 1%	10%
Jigsaw (Goyal et al., 2019)	45.3	79.3
VAT + Ent Min (Grandvalet et al., Miyato et al.)	47.0	83.4
S4L Rotation (Zhai et al., 2019)	53.4	<b>83.8</b>
PIRL	<b>57.2</b>	<b>83.8</b>

# Linear Classification

- Linear classifiers on fixed features. Evaluate on **ImageNet-1K**

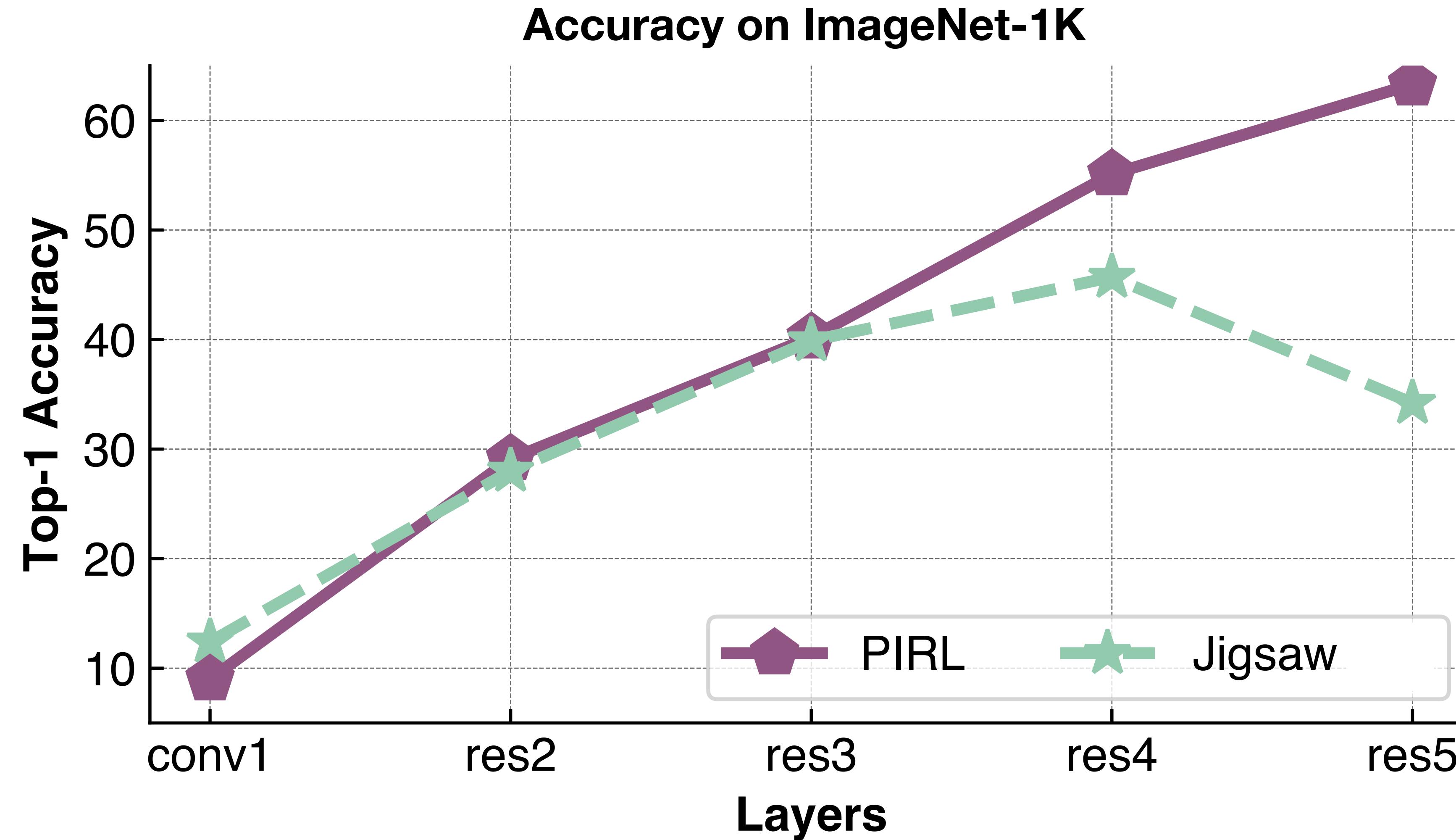


# "In-the-wild" Flickr images

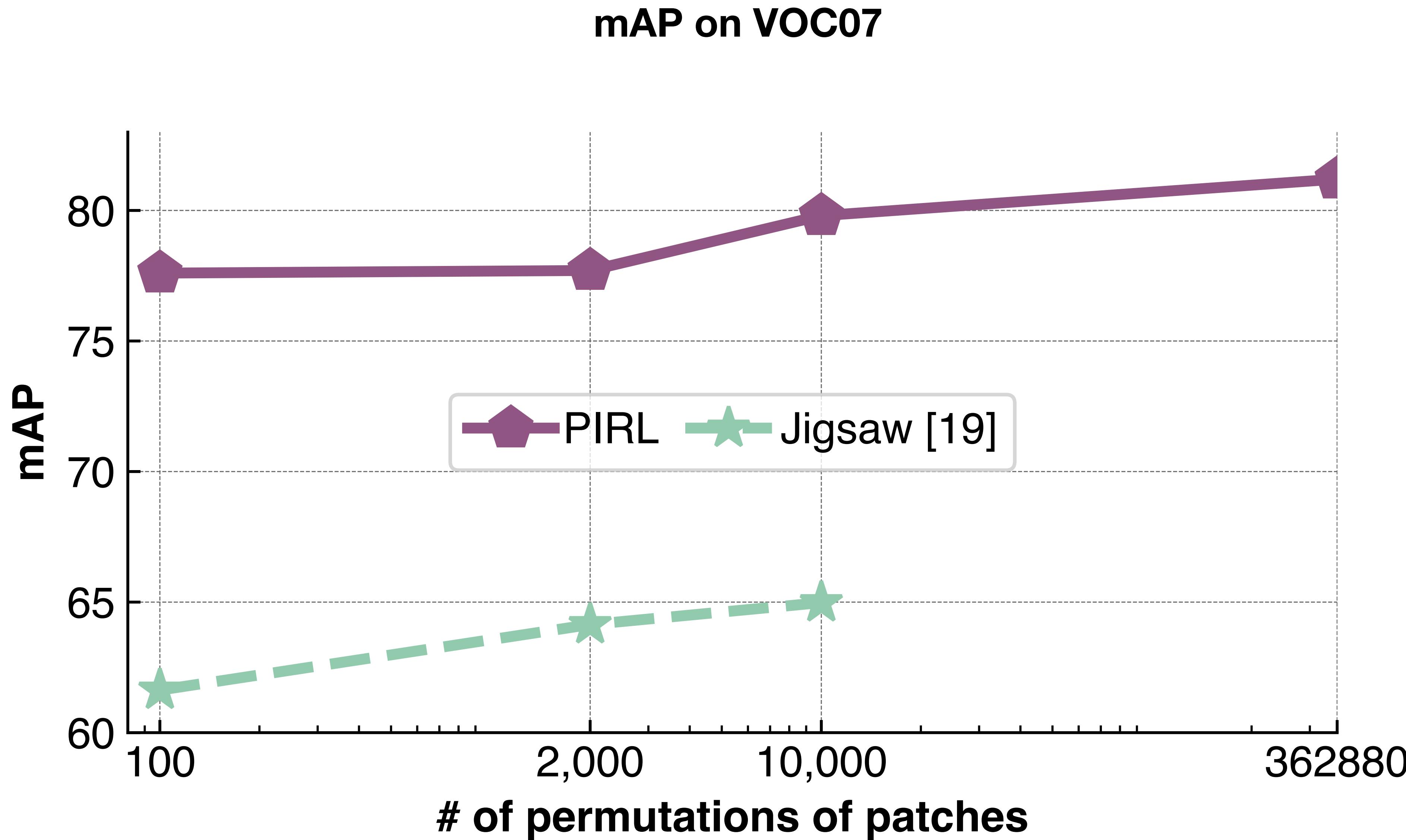
- Yahoo Flickr Creative Commons (YFCC) images. No labels.
- Linear classifiers on fixed features

Method	# Pretrain Images	Evaluation	
		ImageNet	Places-205
DeeperCluster (Caron et al., 2019)	100M	45.6	42.1
Jigsaw (Goyal et al., 2019)	100M	48.3	44.8
PIRL	1M $\frac{1}{100}$	57.8 +9	51.0 +6

# Semantic Features?



# Easily scale to large number of permutations



# Easily Extend to other tasks

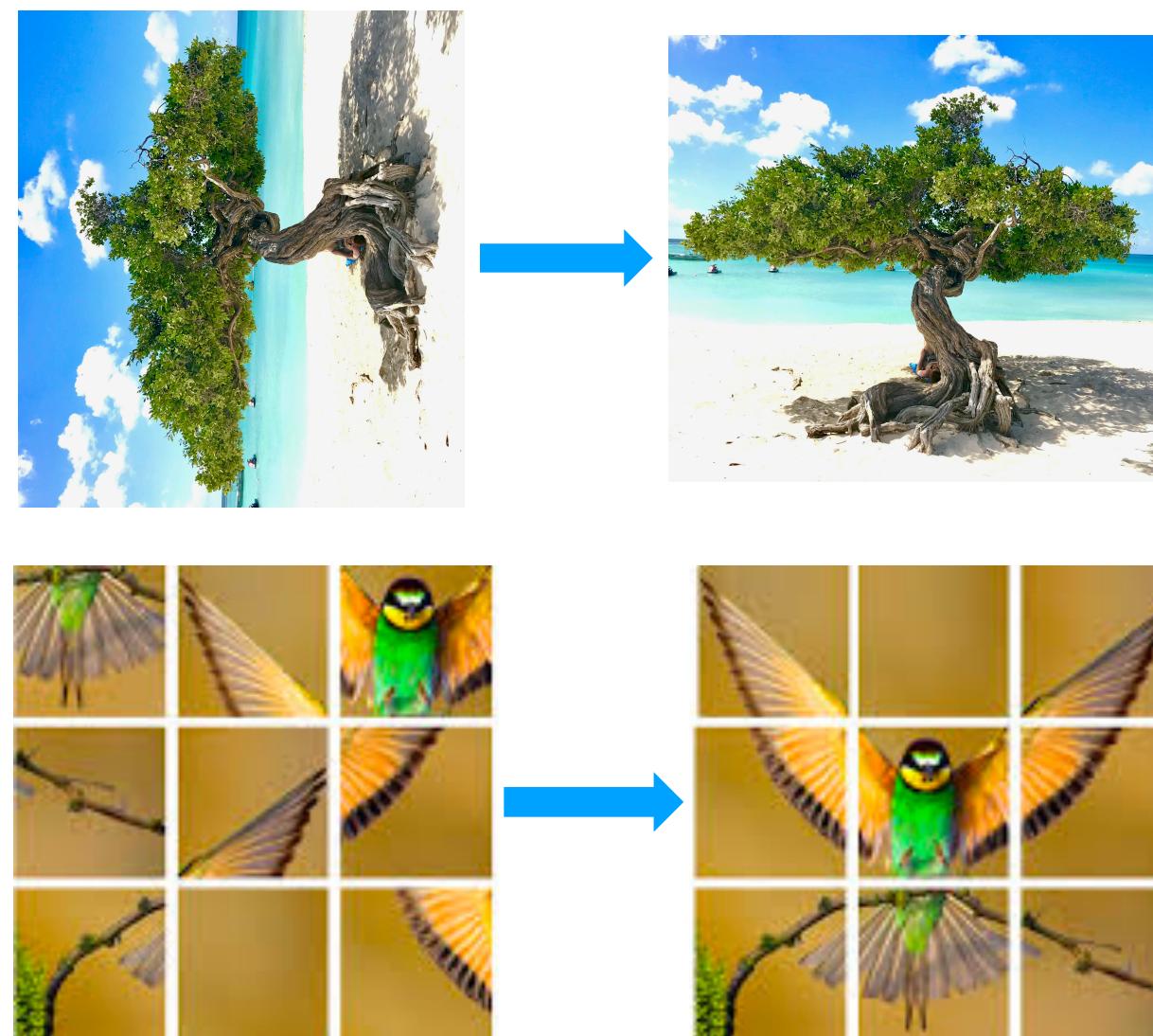
Method	Transfer Dataset			
	ImageNet-1M	VOC07	Places205	iNaturalist
Jigsaw	46.0	66.1	41.4	22.1
Rotation	48.9	63.9	47.6	23
PIRL (Rot)	<b>60.2</b>	<b>77.1</b>	<b>47.6</b>	<b>31.2</b>

# Easily Multi-task

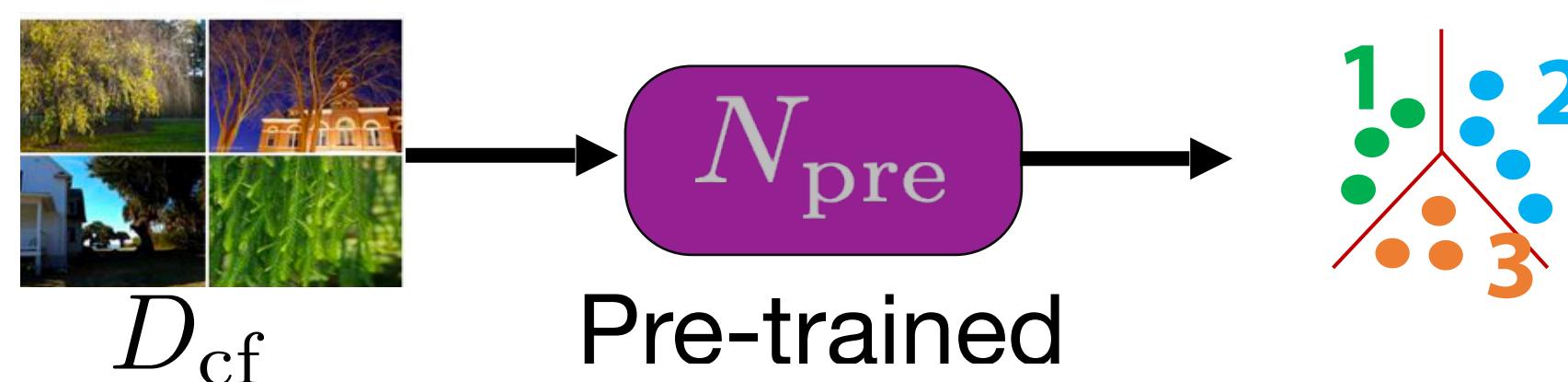
Method	Transfer Dataset			
	ImageNet-1M	VOC07	Places205	iNaturalist
Jigsaw	46.0	66.1	41.4	22.1
Rotation	48.9	63.9	47.6	23
PIRL (Rot)	60.2	77.1	47.6	31.2
PIRL (Jigsaw + Rot)	<b>63.1</b>	<b>80.3</b>	<b>49.7</b>	<b>33.6</b>

# More invariance

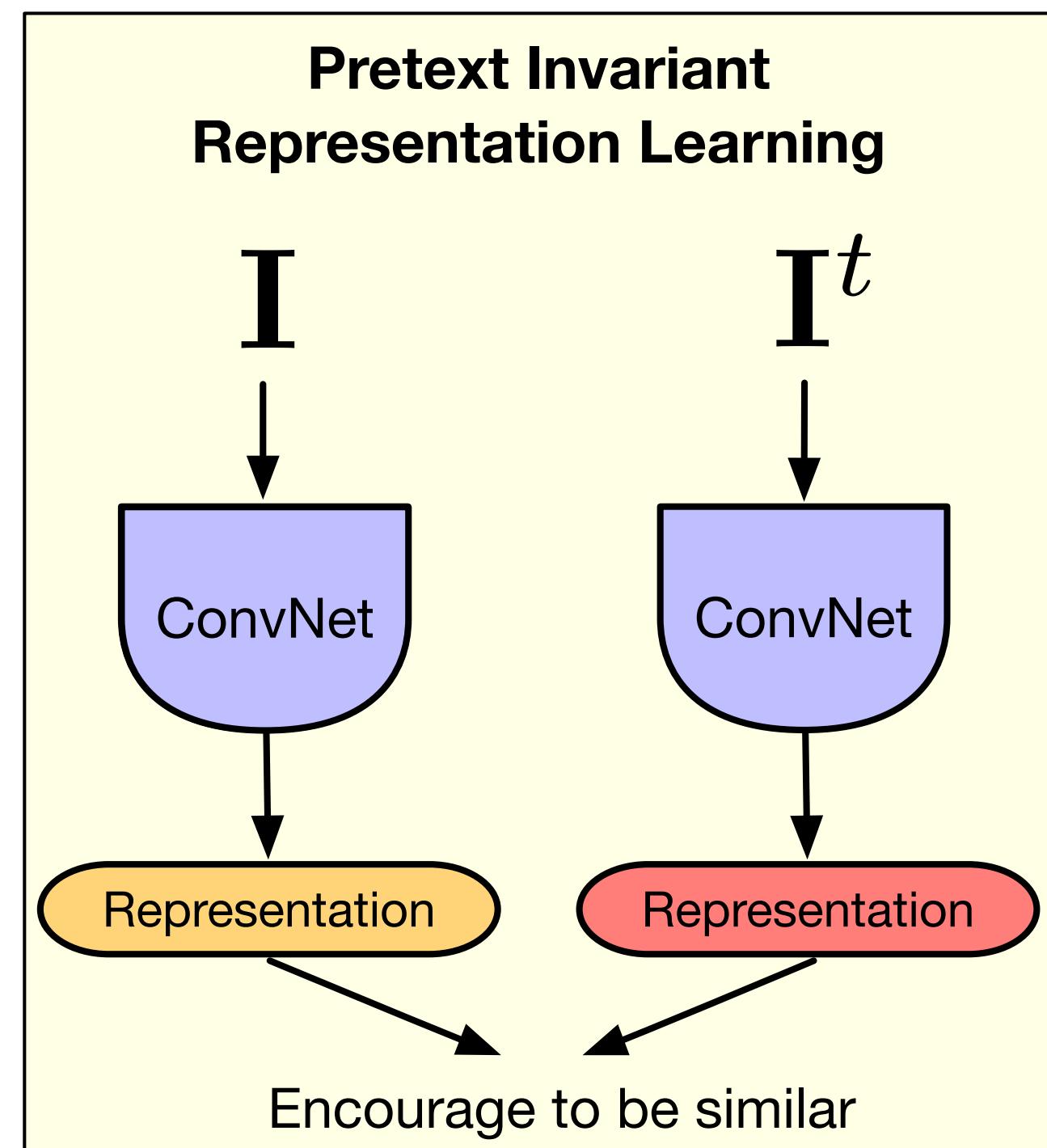
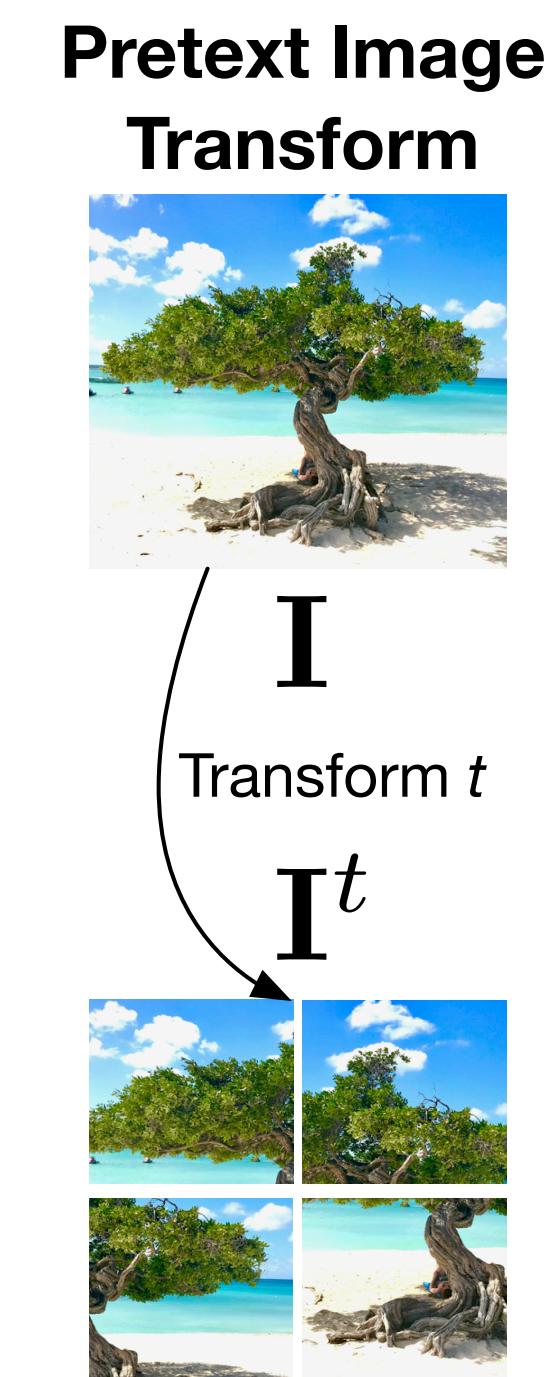
## Pretext tasks



## Clustering



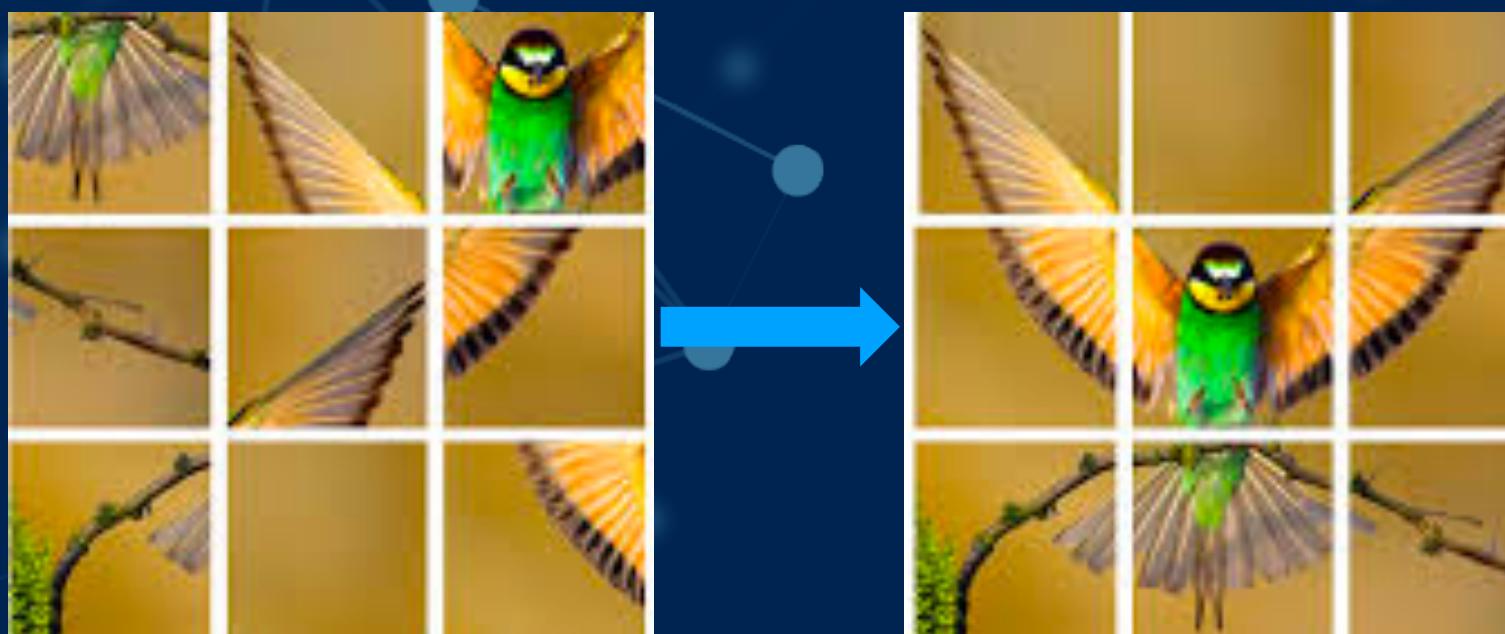
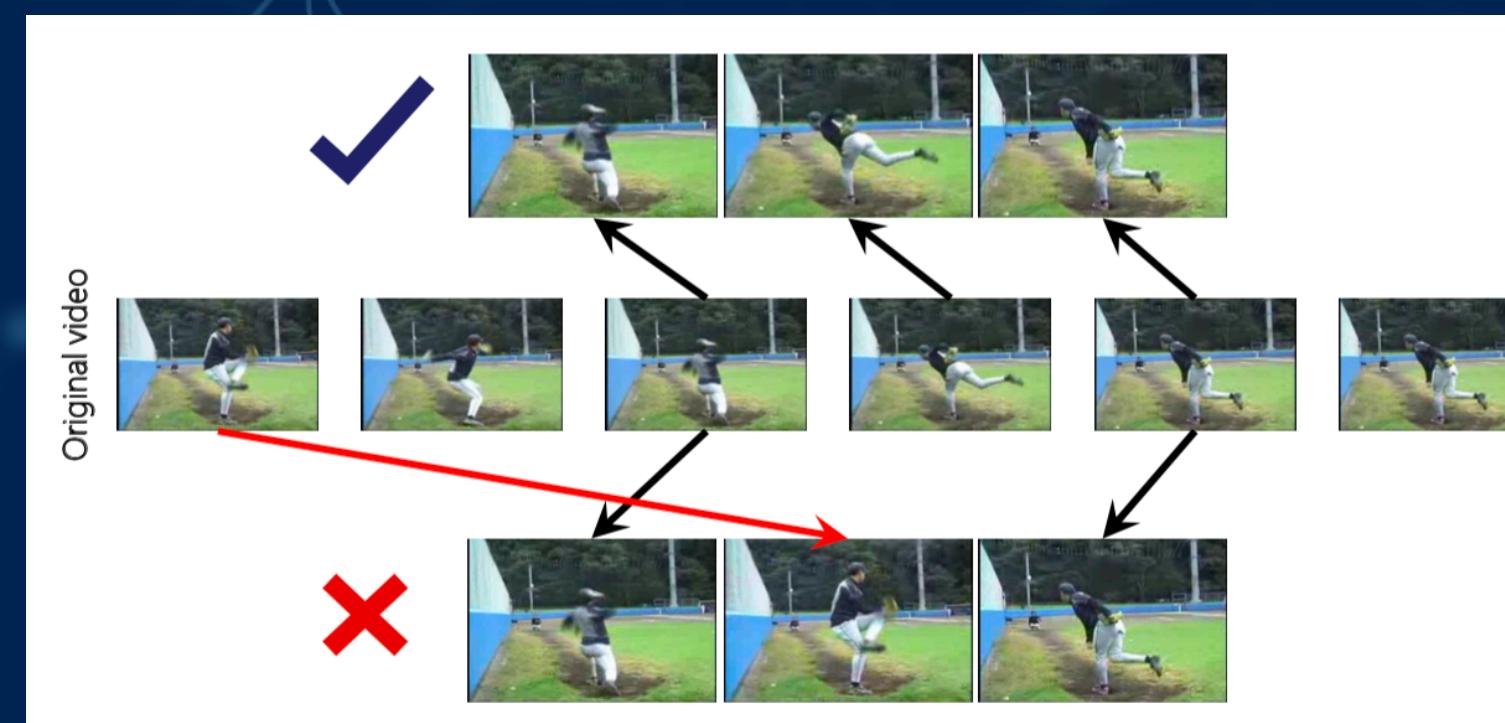
## PIRL



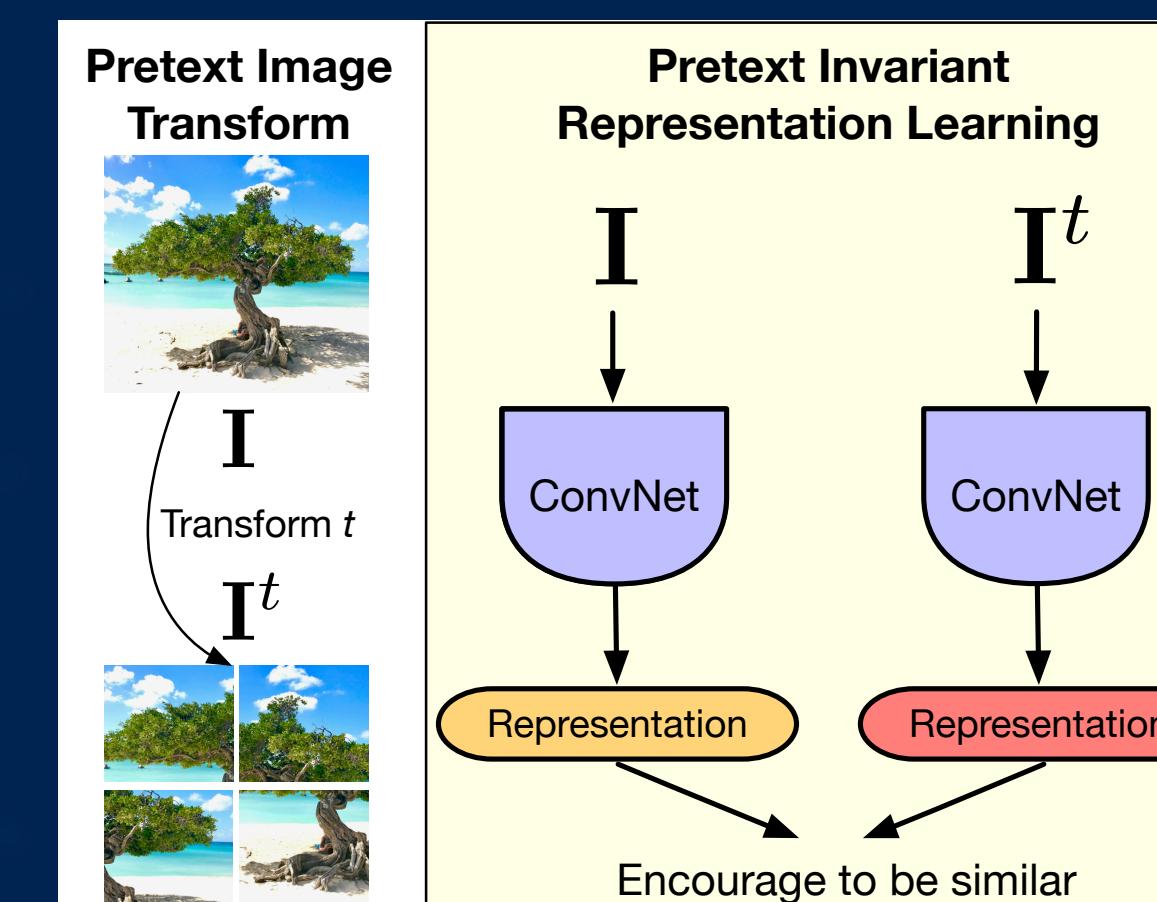
# Shortcomings

- Set of data transforms matters a lot
- Saturation with model size and data size
- What invariances matter?

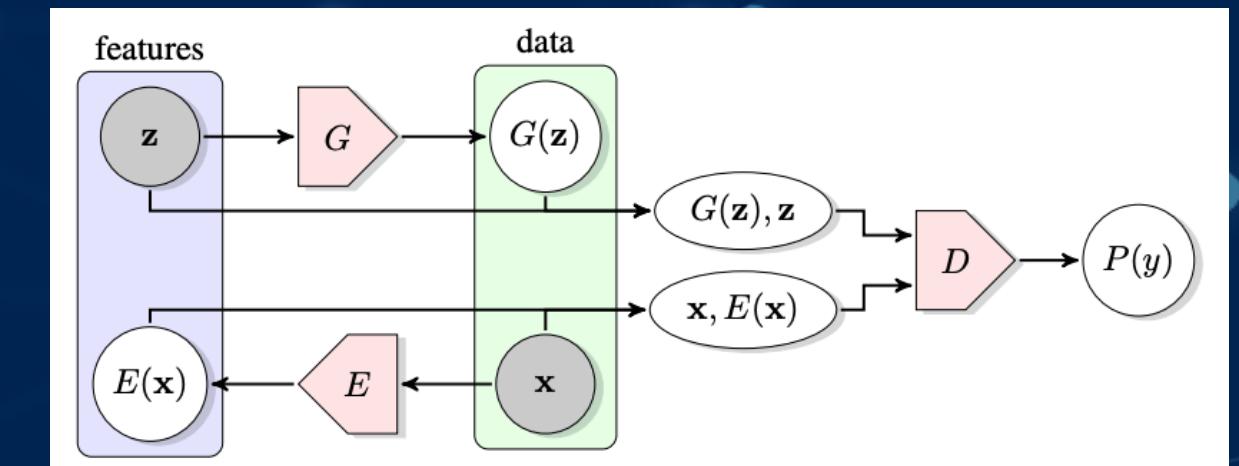
# Pretext tasks



# Contrastive



# Generative



AutoEncoder,  
VAE, GAN,  
BiGAN

Predict more information