

Self-supervised learning in computer vision

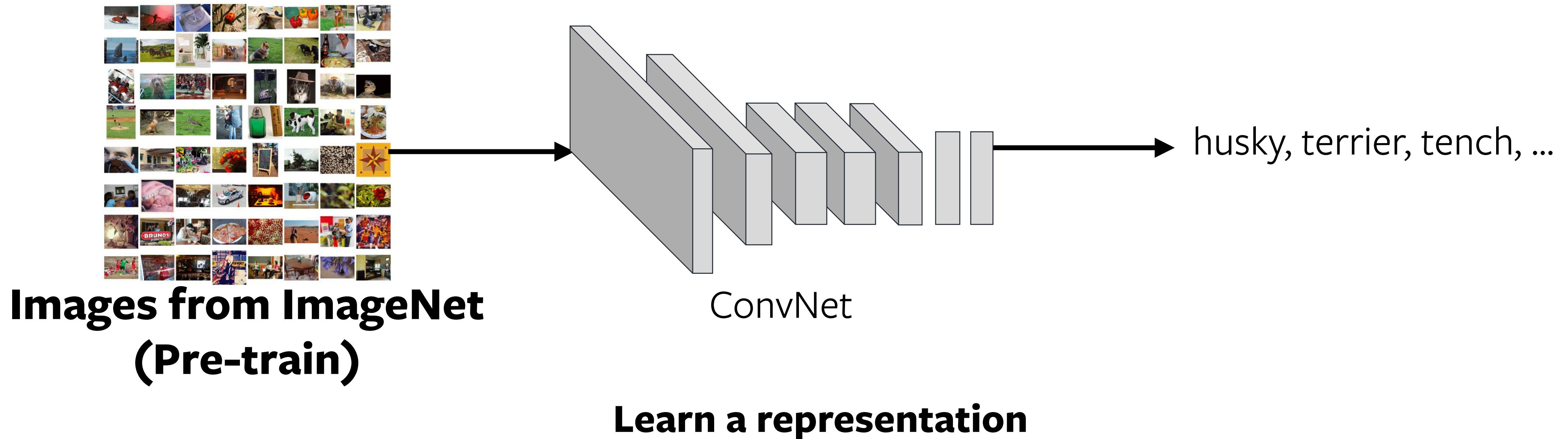
Ishan Misra

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

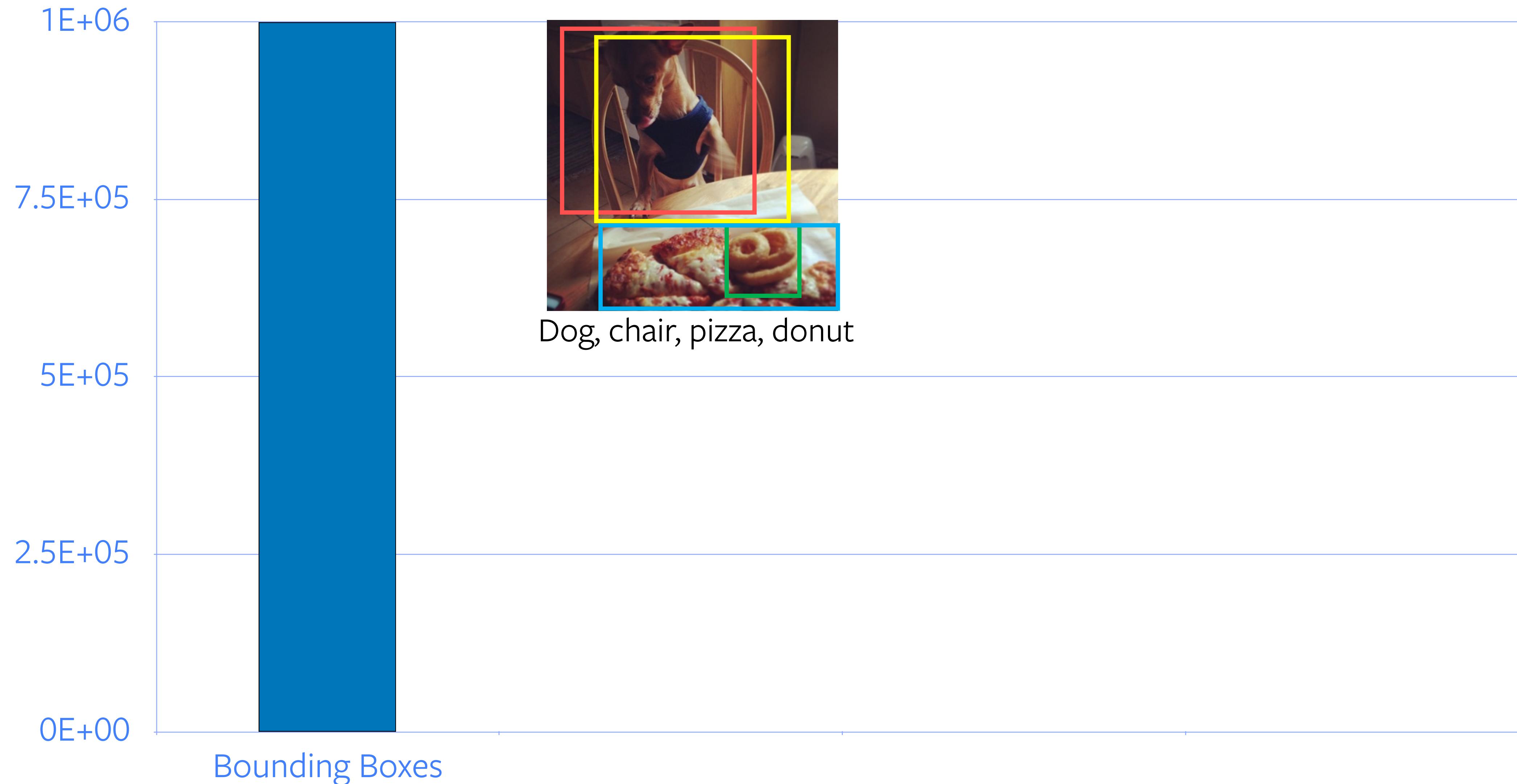
- Pre-train on a large supervised dataset.
- Collect a dataset of “supervised” images
- Train a ConvNet

The promise of "alternative" supervision

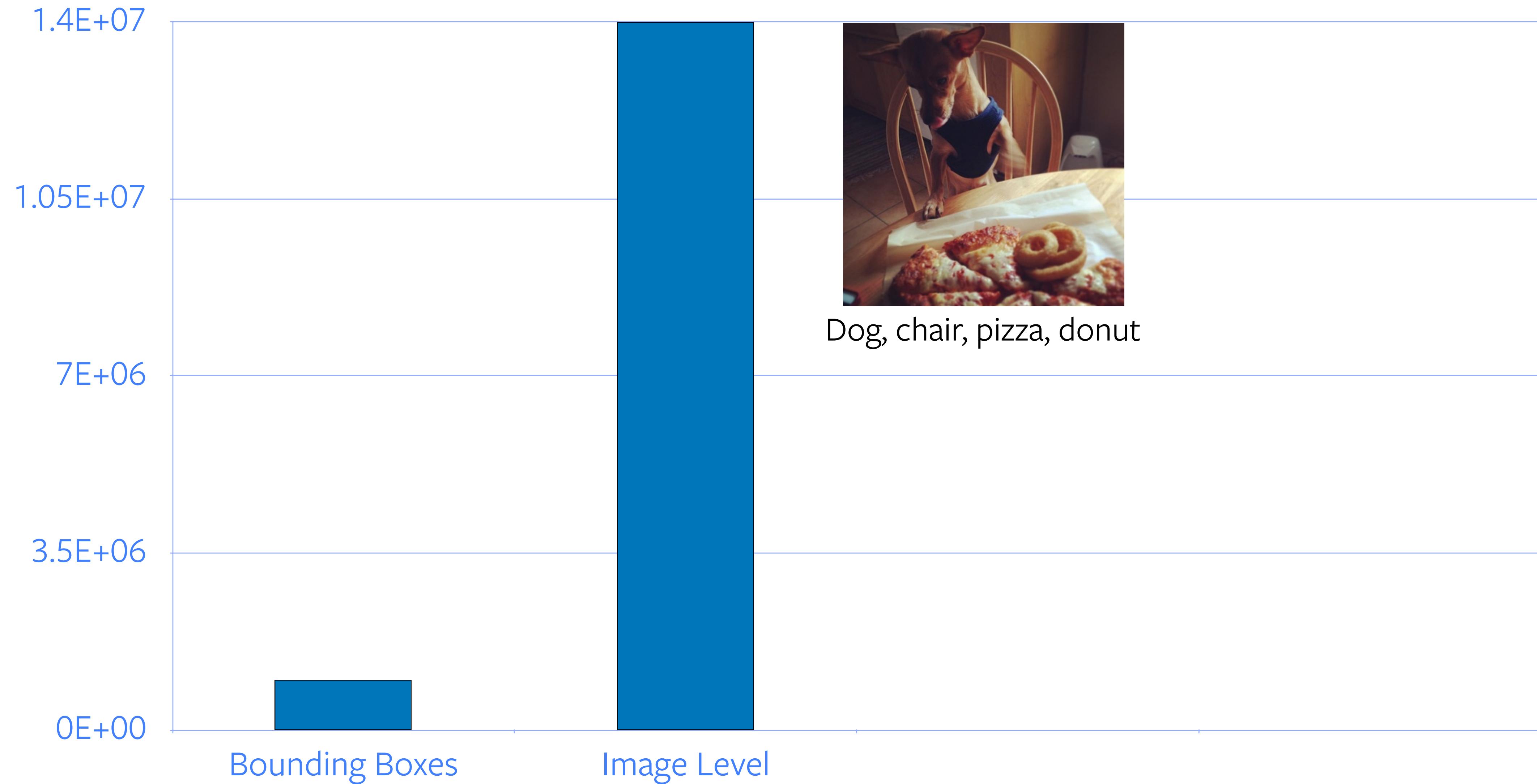
- Getting "real" labels is difficult and expensive
 - ImageNet with 14M images took 22 human years.
- Obtain labels using a "semi-automatic" process
 - Hashtags
 - GPS locations
 - Using the data itself: "self"-supervised

Can we get labels for all data?

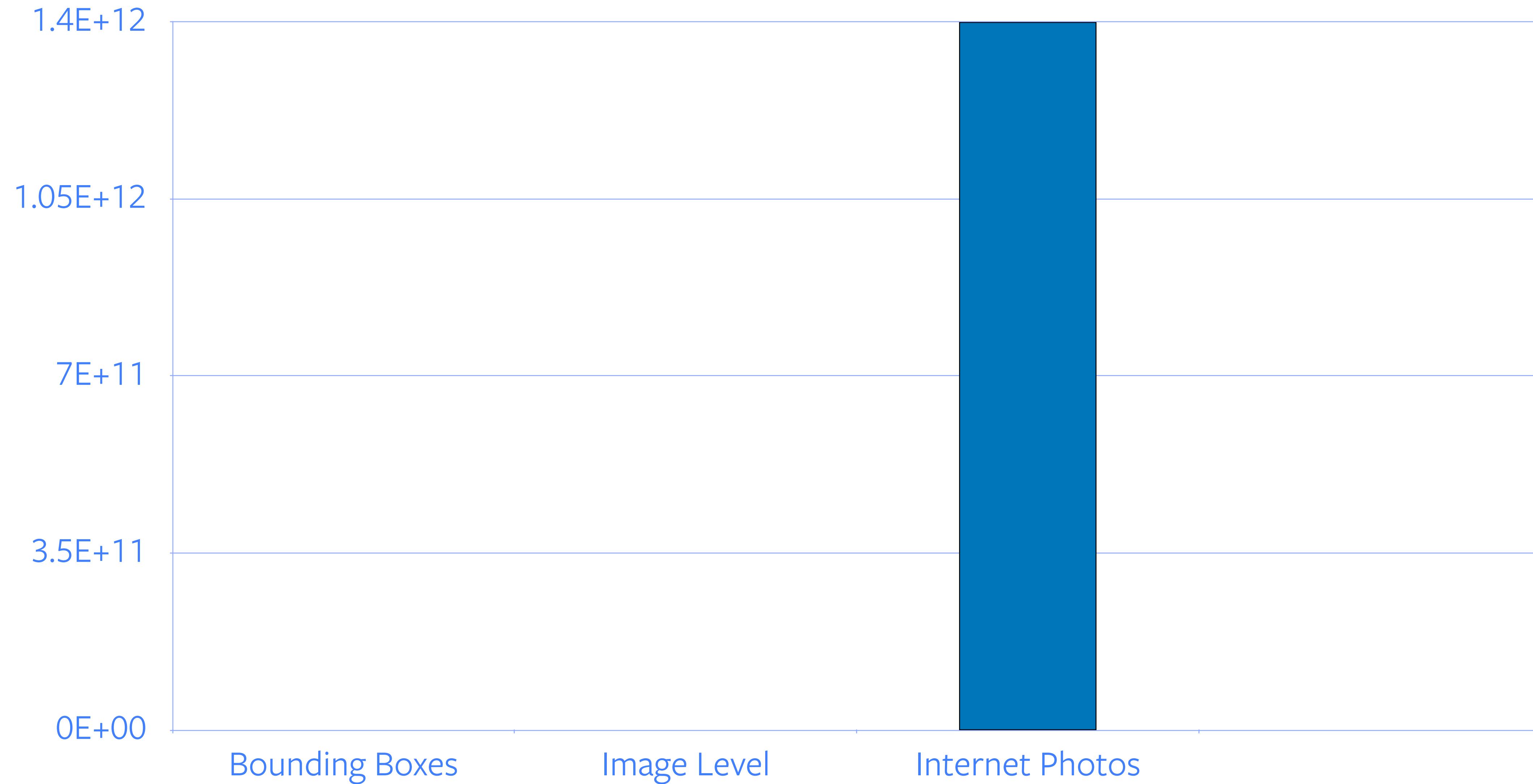
Can we get labels for all data?



Can we get labels for all data?



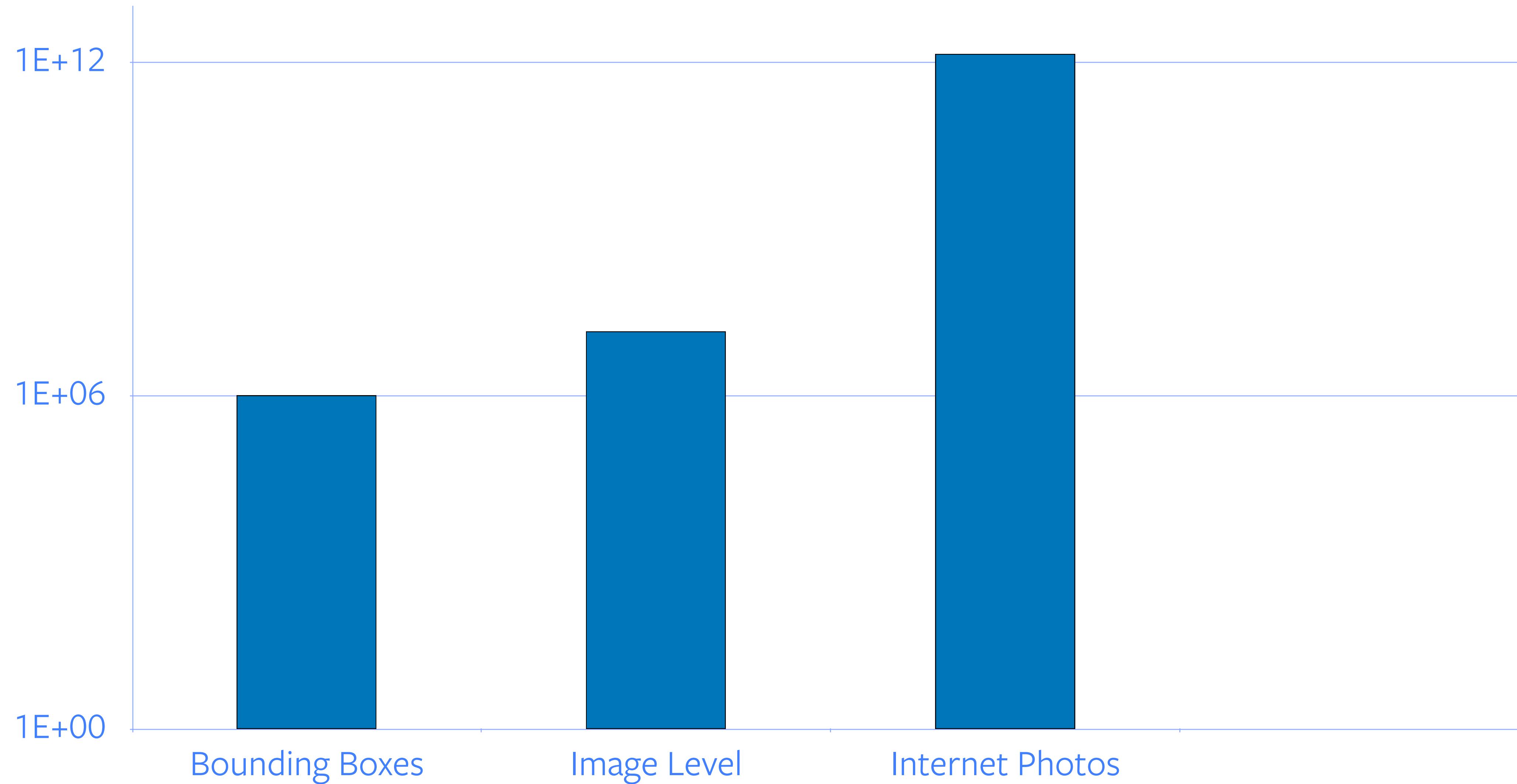
Can we get labels for all data?



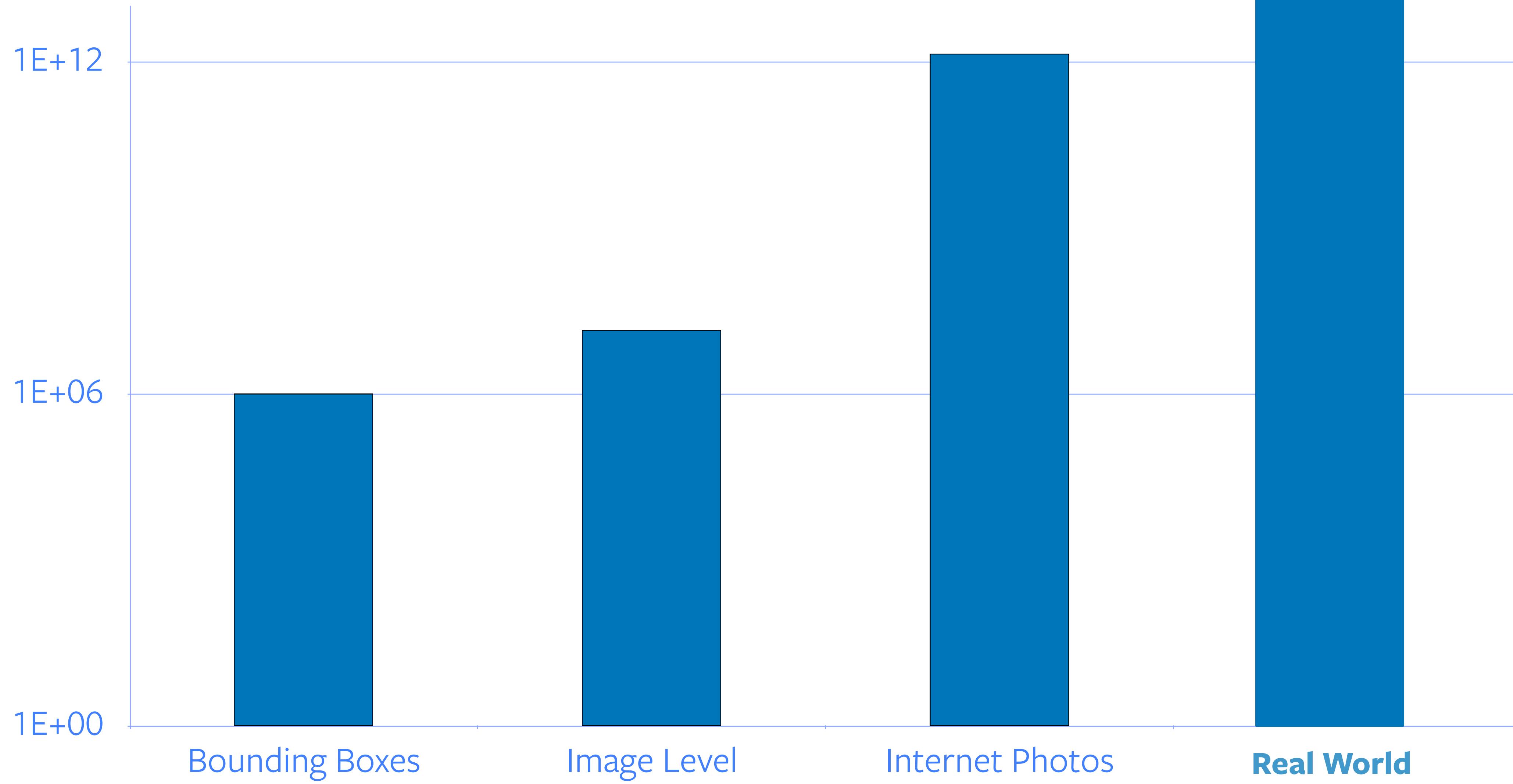
[forbes.com](https://www.forbes.com/sites/bernardmarr/2018/05/21/how-much-data-do-we-create-every-day-the-mind-blowing-stats-everyone-should-read/)

<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?



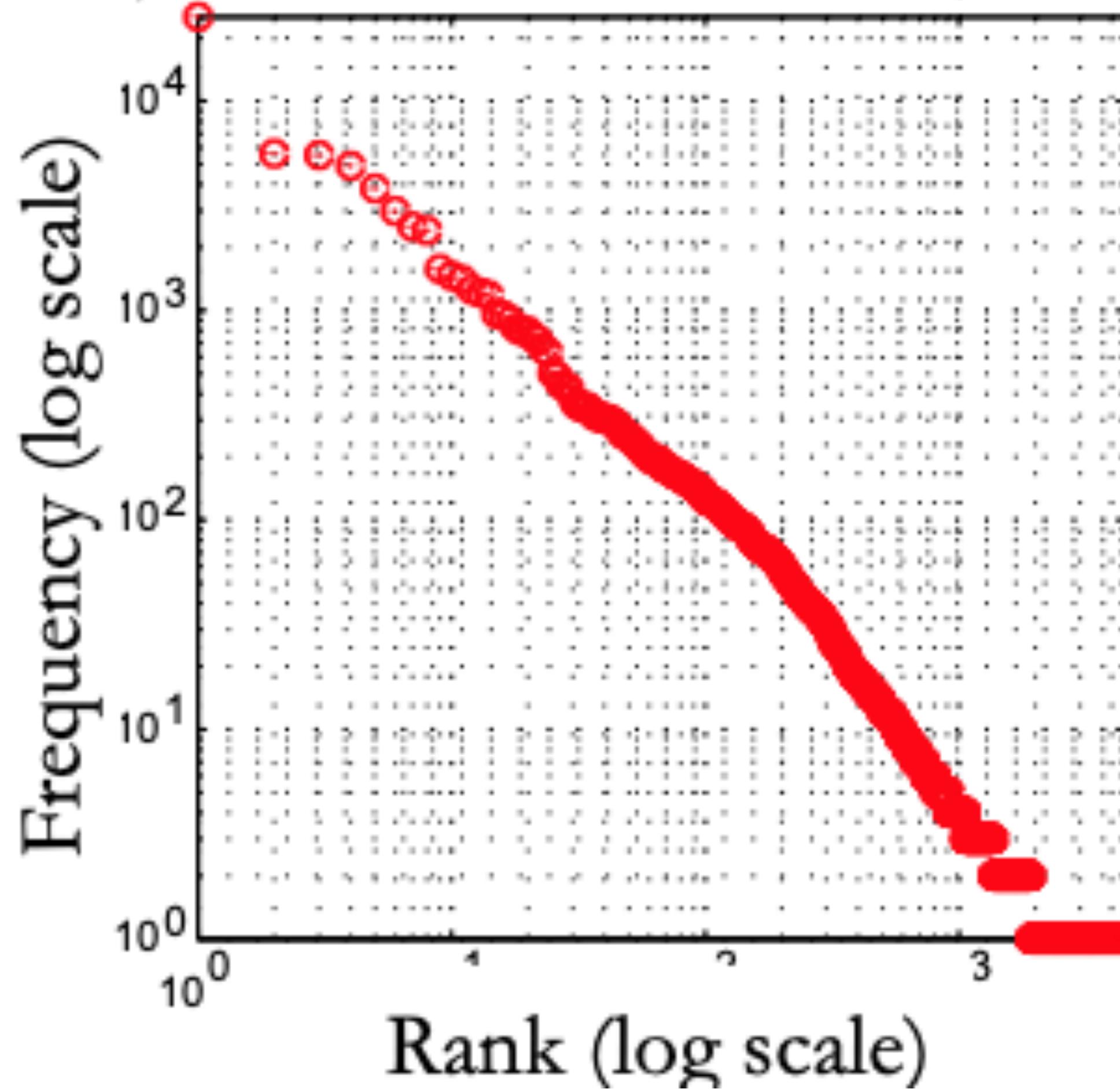
Can we get labels for all data?



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

Rare concepts?

Objects in Vision Dataset (LabelMe)



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

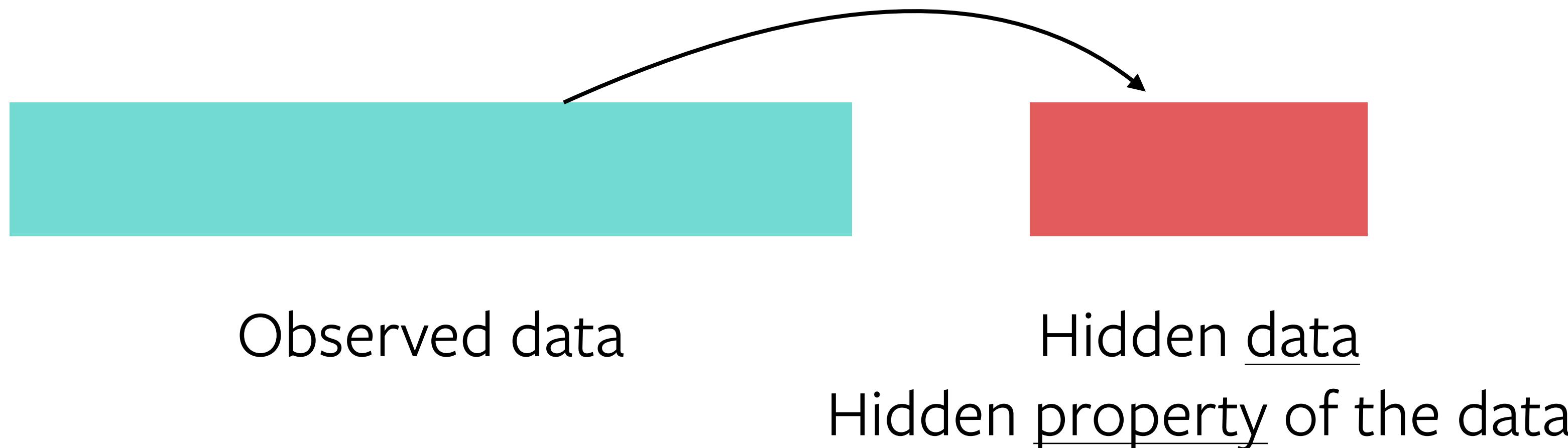
Different Domains?



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

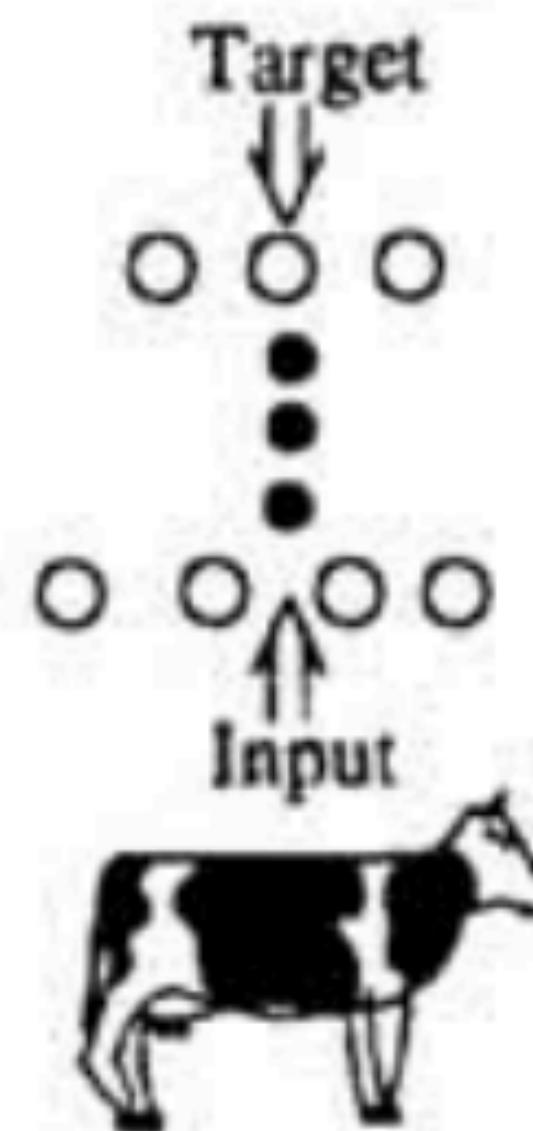


What is "self" supervision?

Supervised

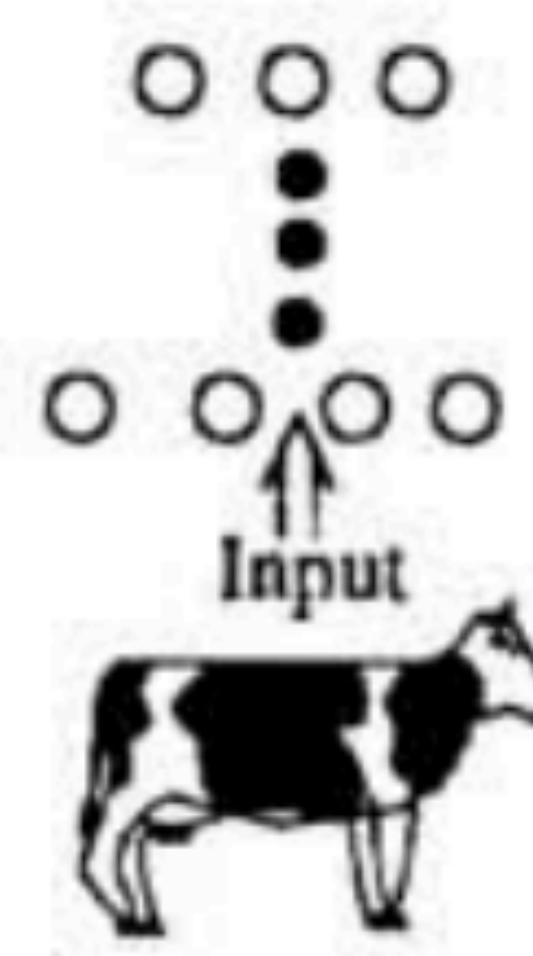
- implausible label

"COW"



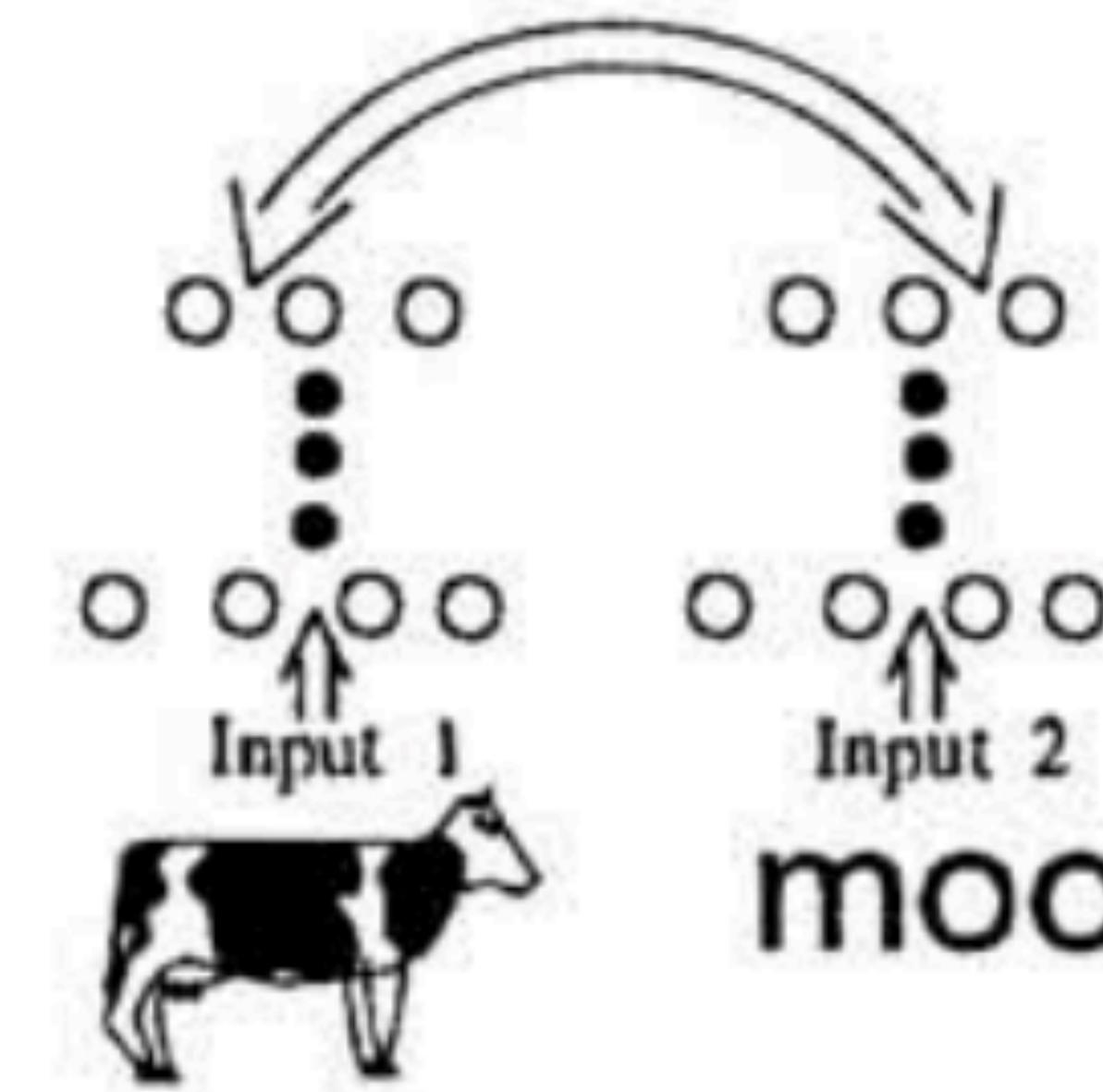
Unsupervised

- limited power



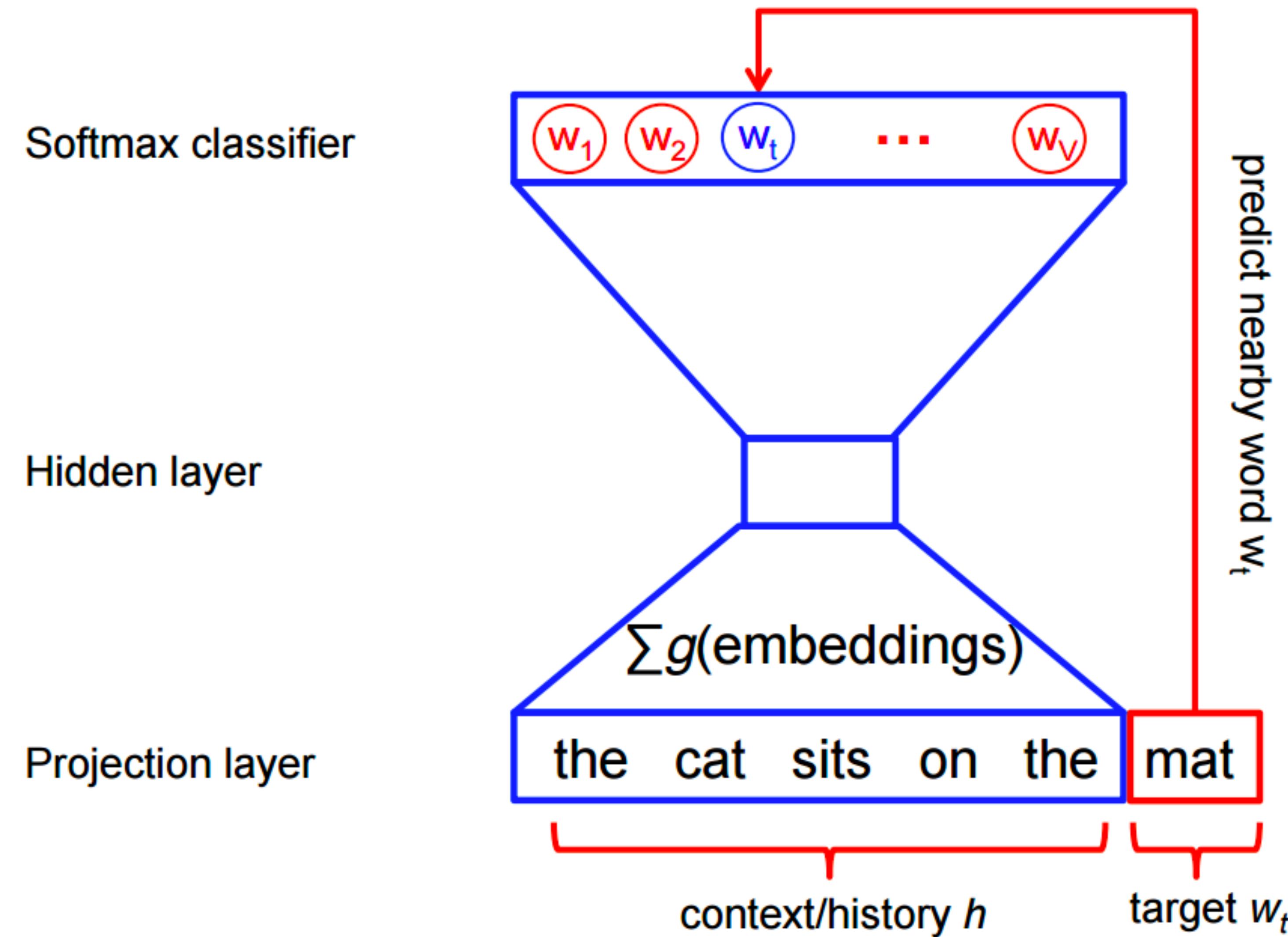
Self-Supervised

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



Word2vec

- Fill in the blanks

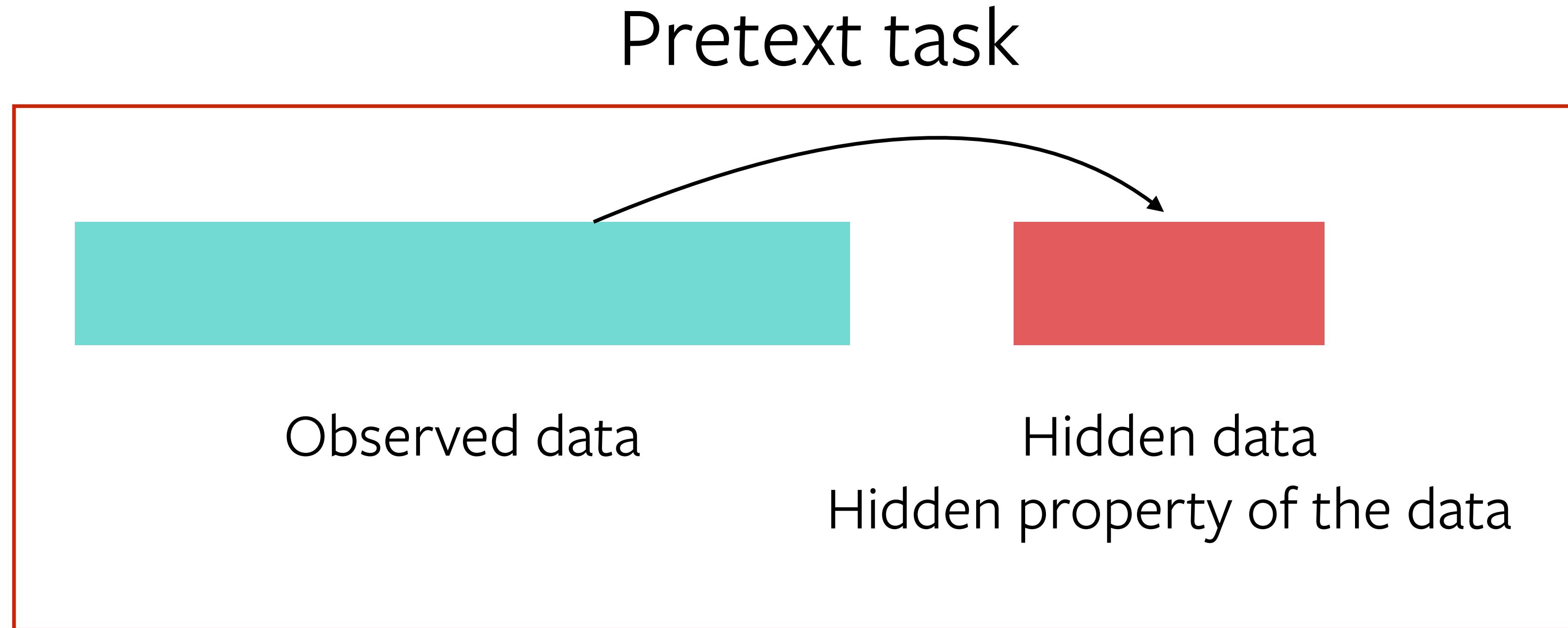


Success of self-supervised learning in NLP

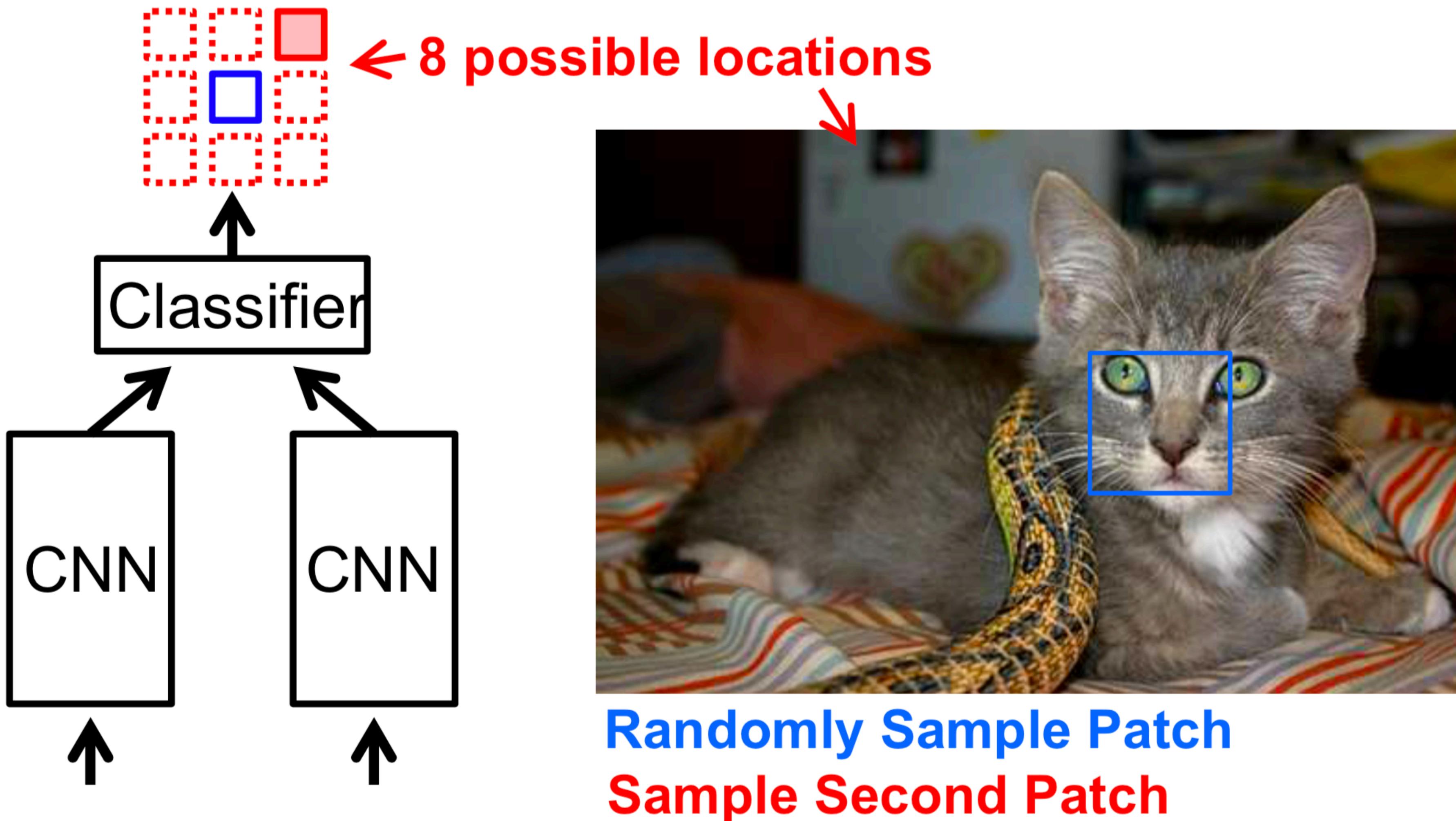
- Fill in the blanks is a powerful signal to learn representations
- Sentence/Word representations: BERT - Devlin et al., 2018

Pretext task

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



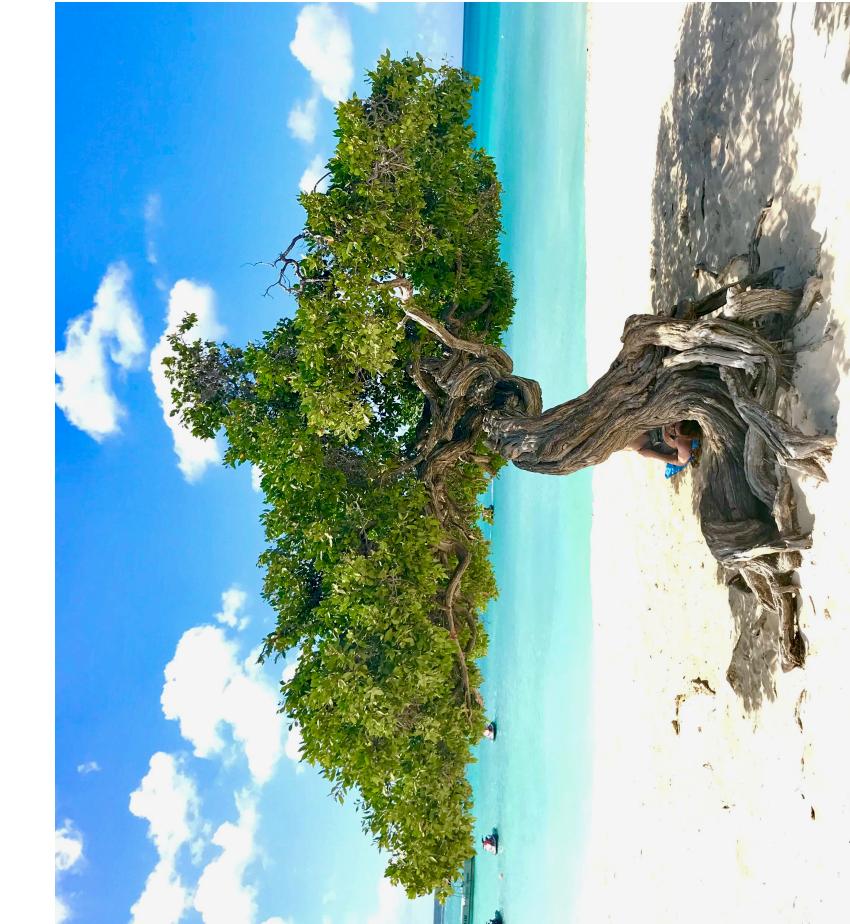
Relative Position of patches



Predicting Rotations



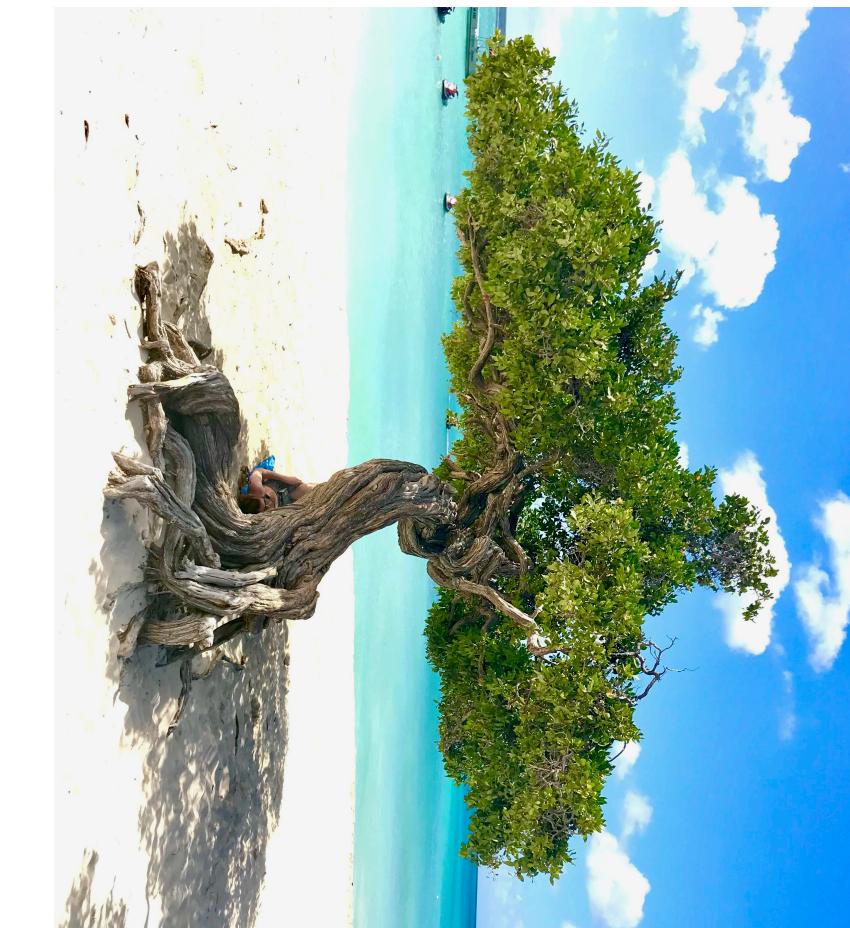
$\rightarrow 0^\circ$



$\rightarrow 90^\circ$

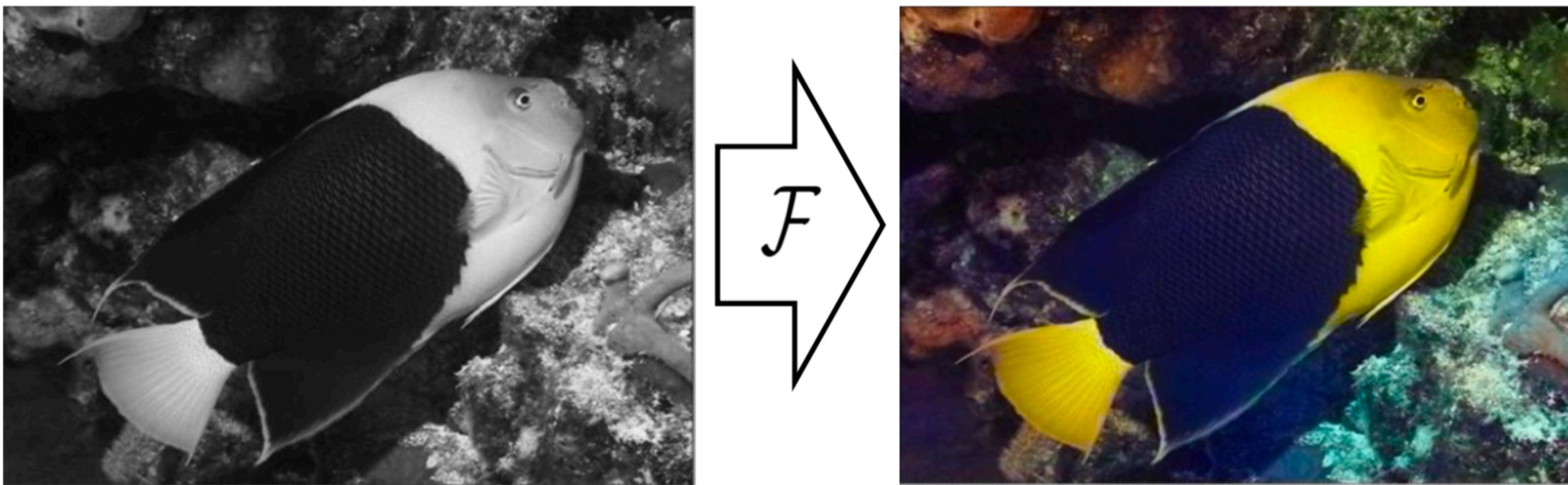


$\rightarrow 180^\circ$



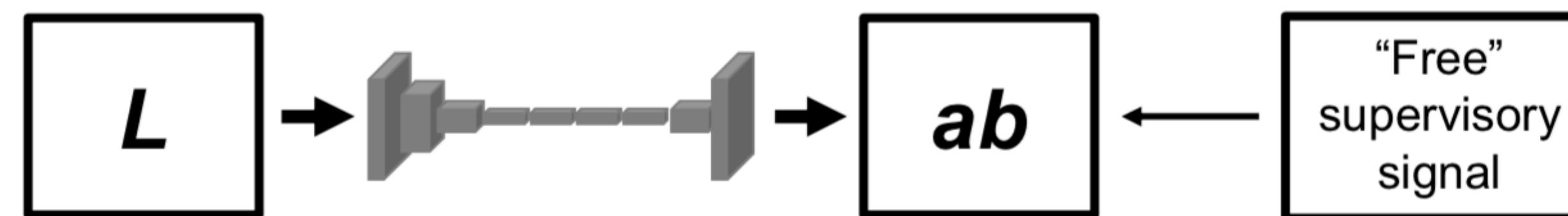
$\rightarrow 270^\circ$

Colorization



Grayscale image: L channel

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



Concatenate (L,ab)

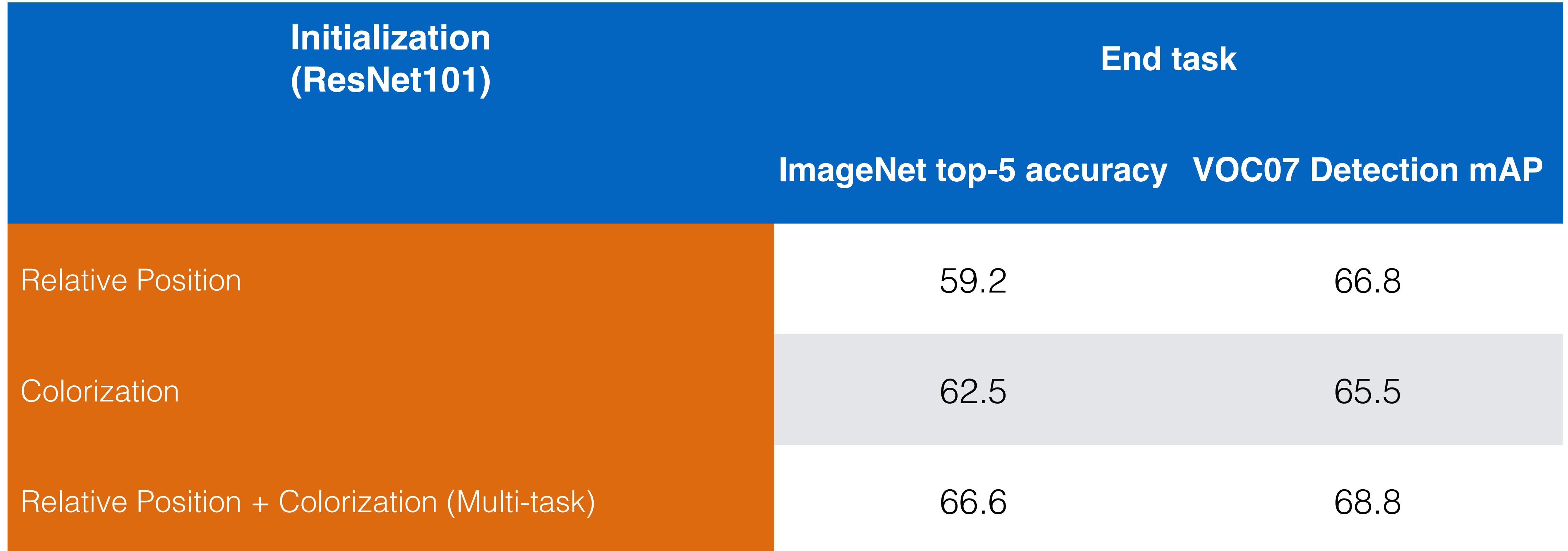
$$(\mathbf{X}, \hat{\mathbf{Y}})$$

Fill in the blanks



Understanding what the “pretext” task learns

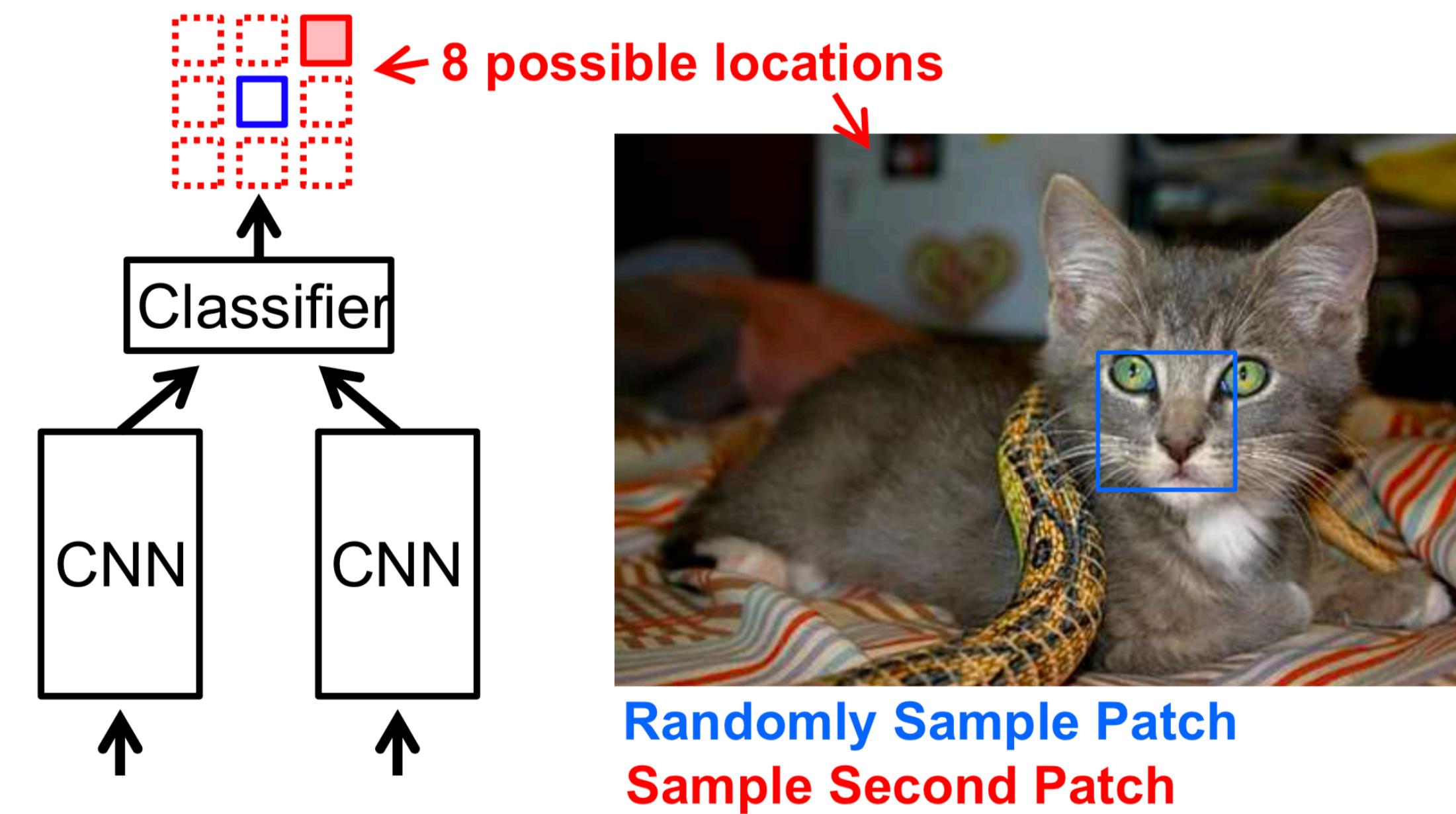
Are they complementary?



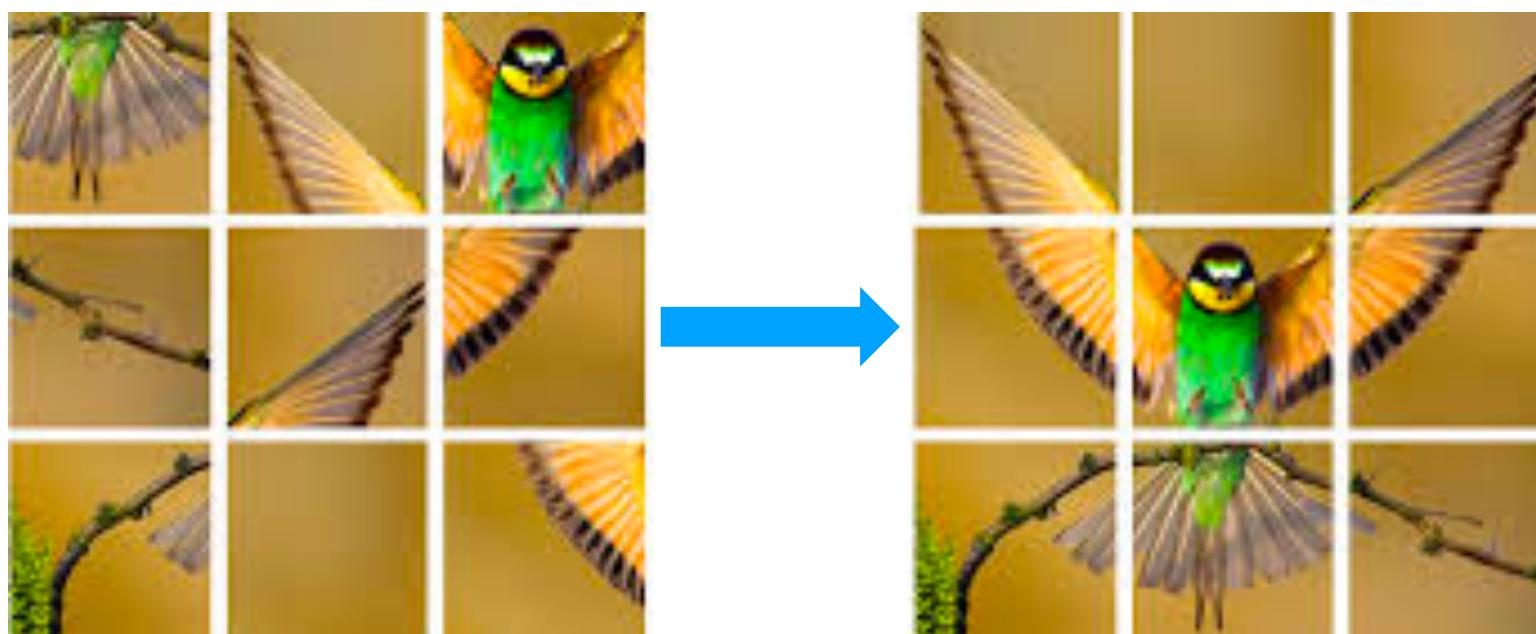
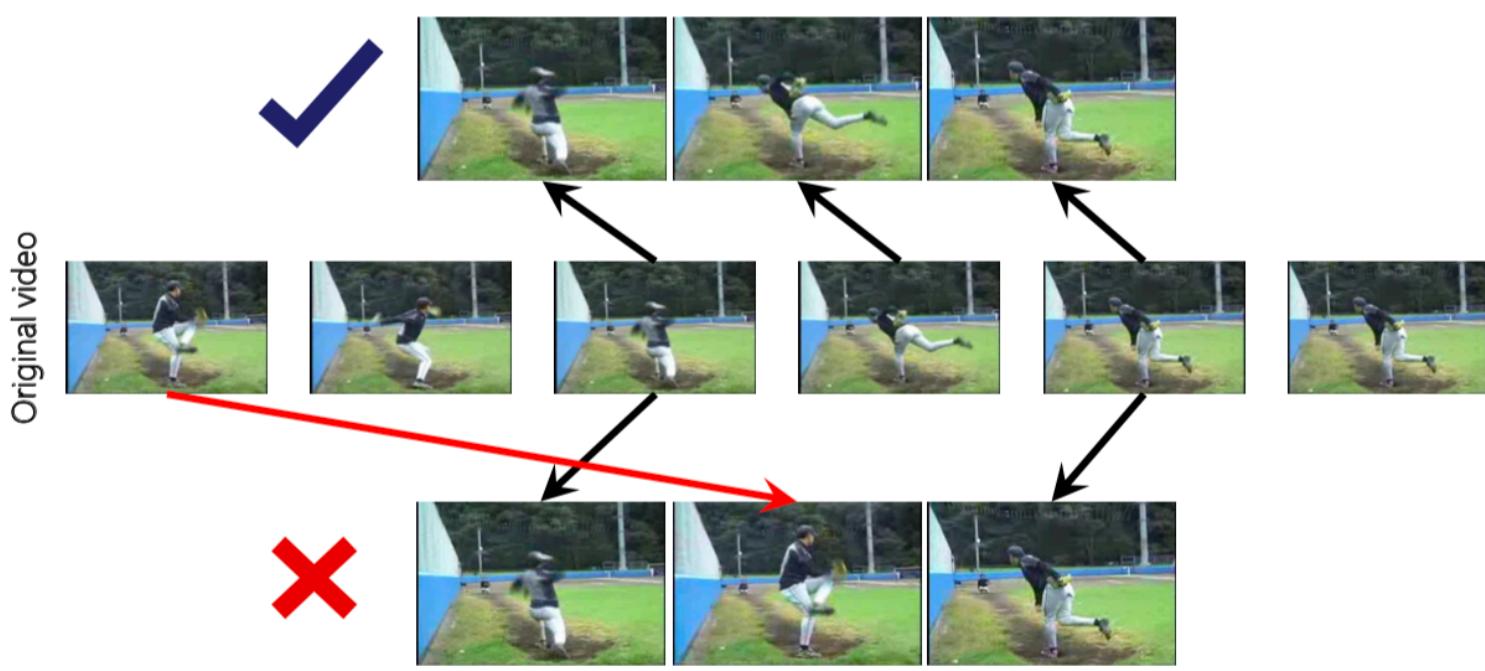
Information predicted: varies across tasks

Less

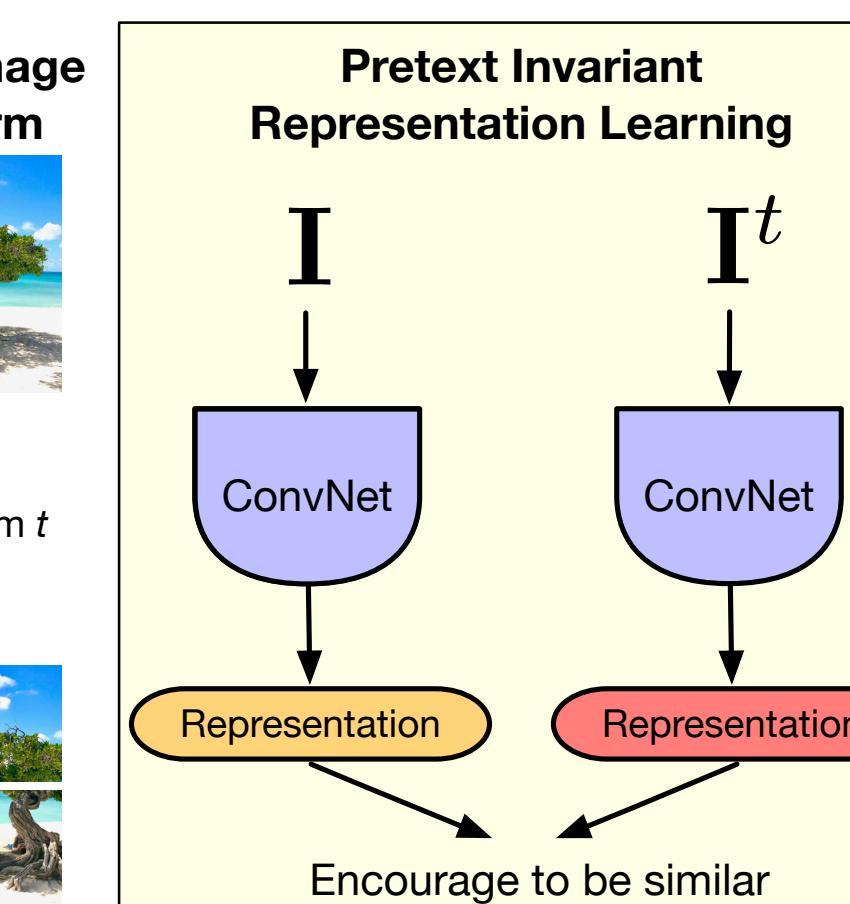
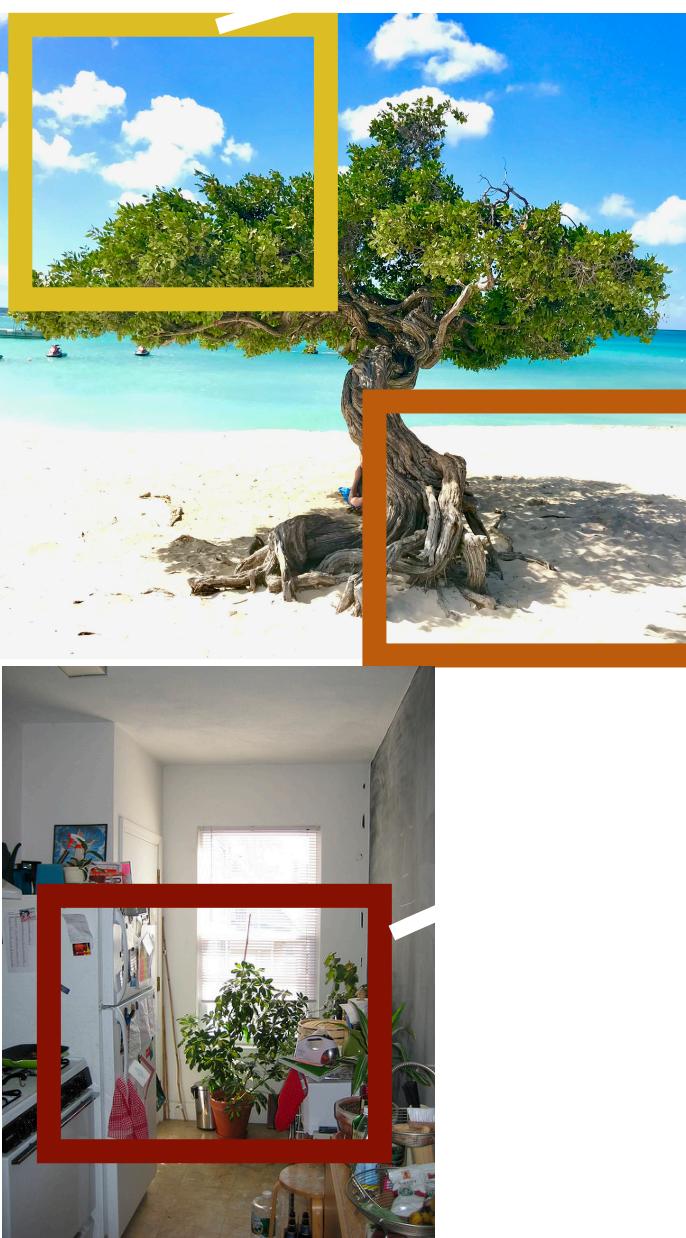
More



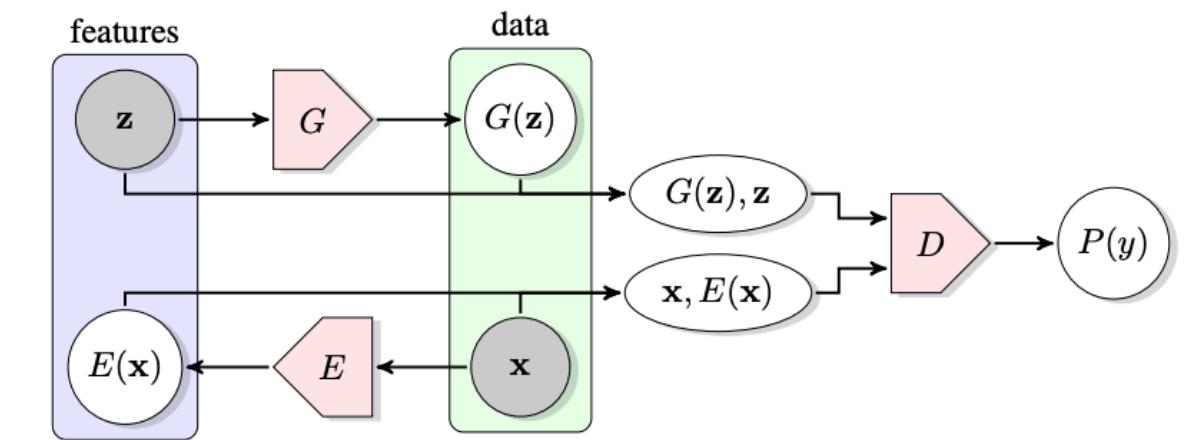
Pretext tasks



Contrastive



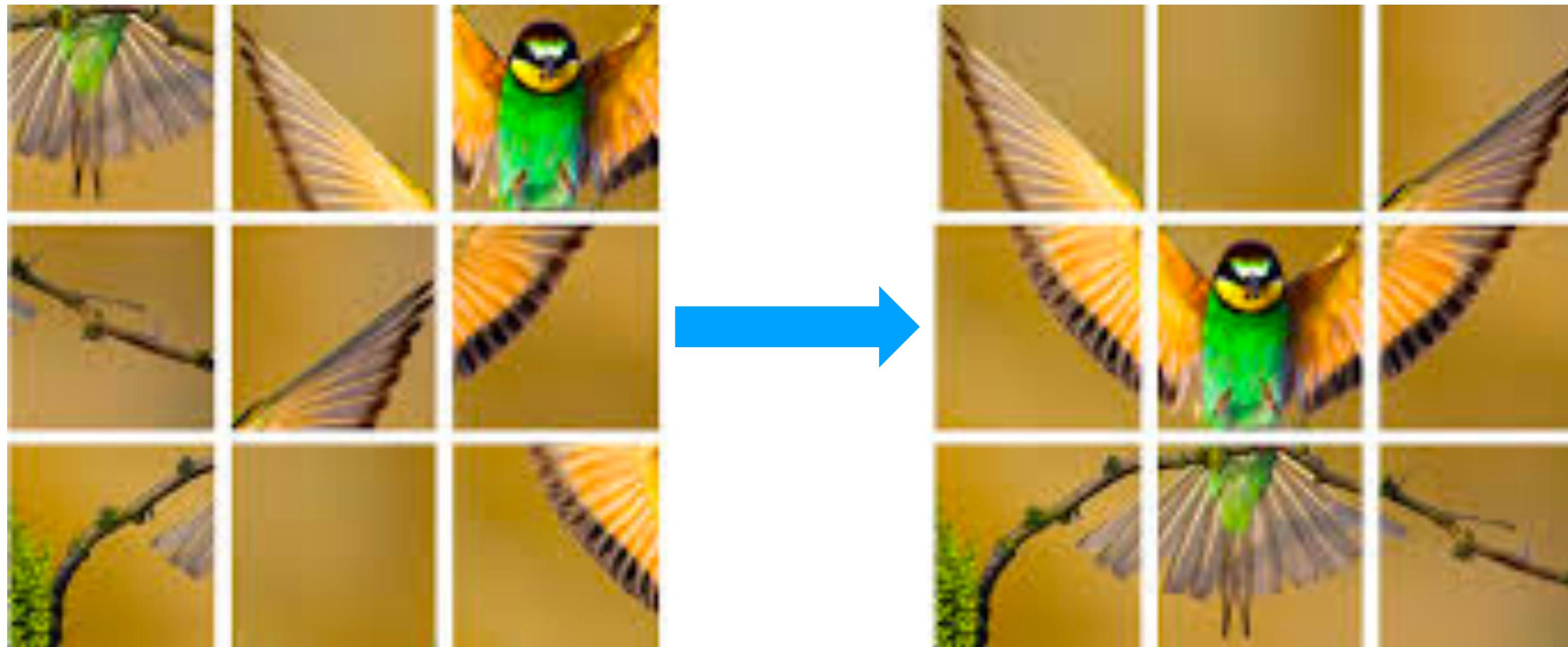
Generative



AutoEncoder,
VAE, GAN,
BiGAN

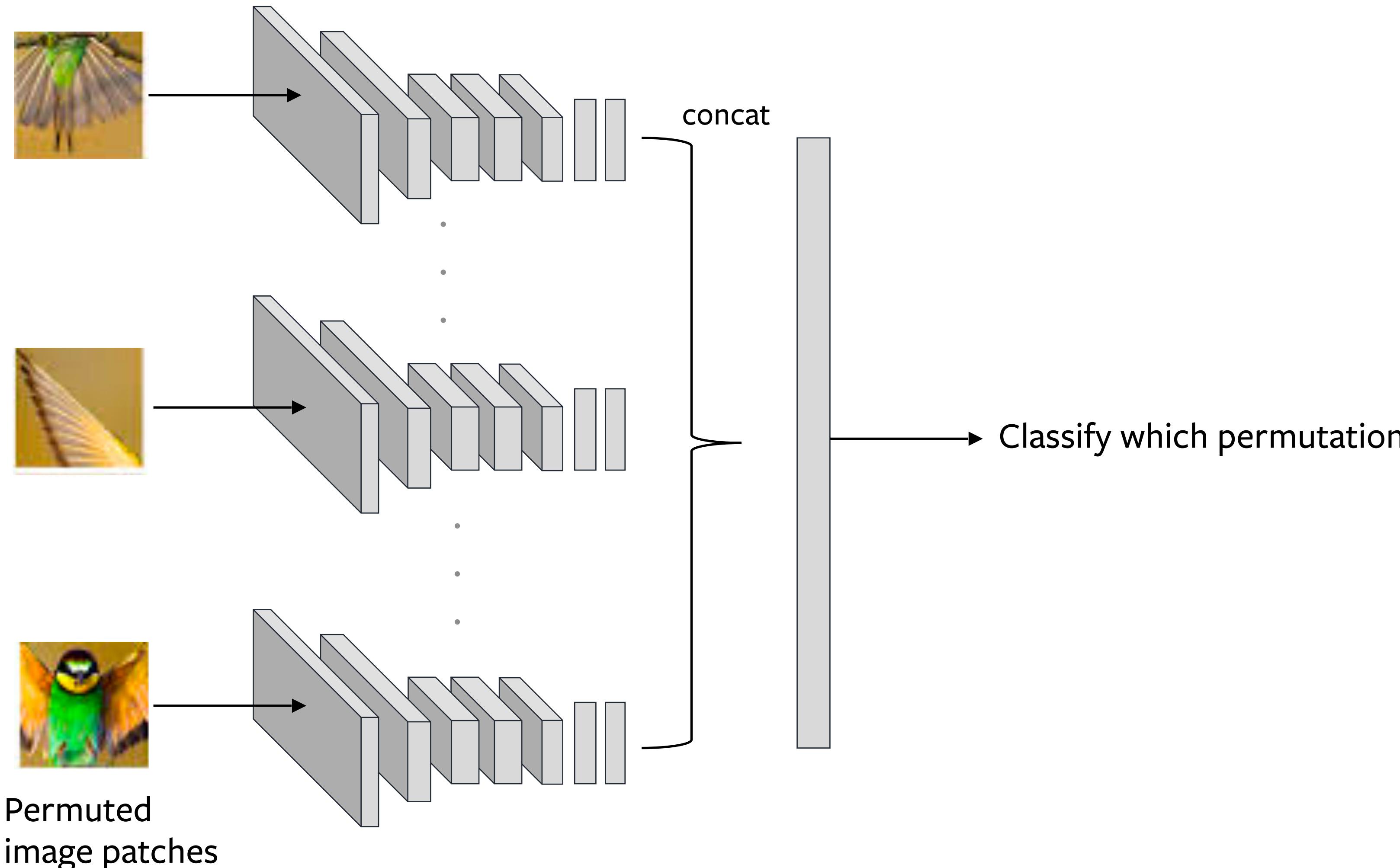
→ Predict more information

Scaling self-supervised learning



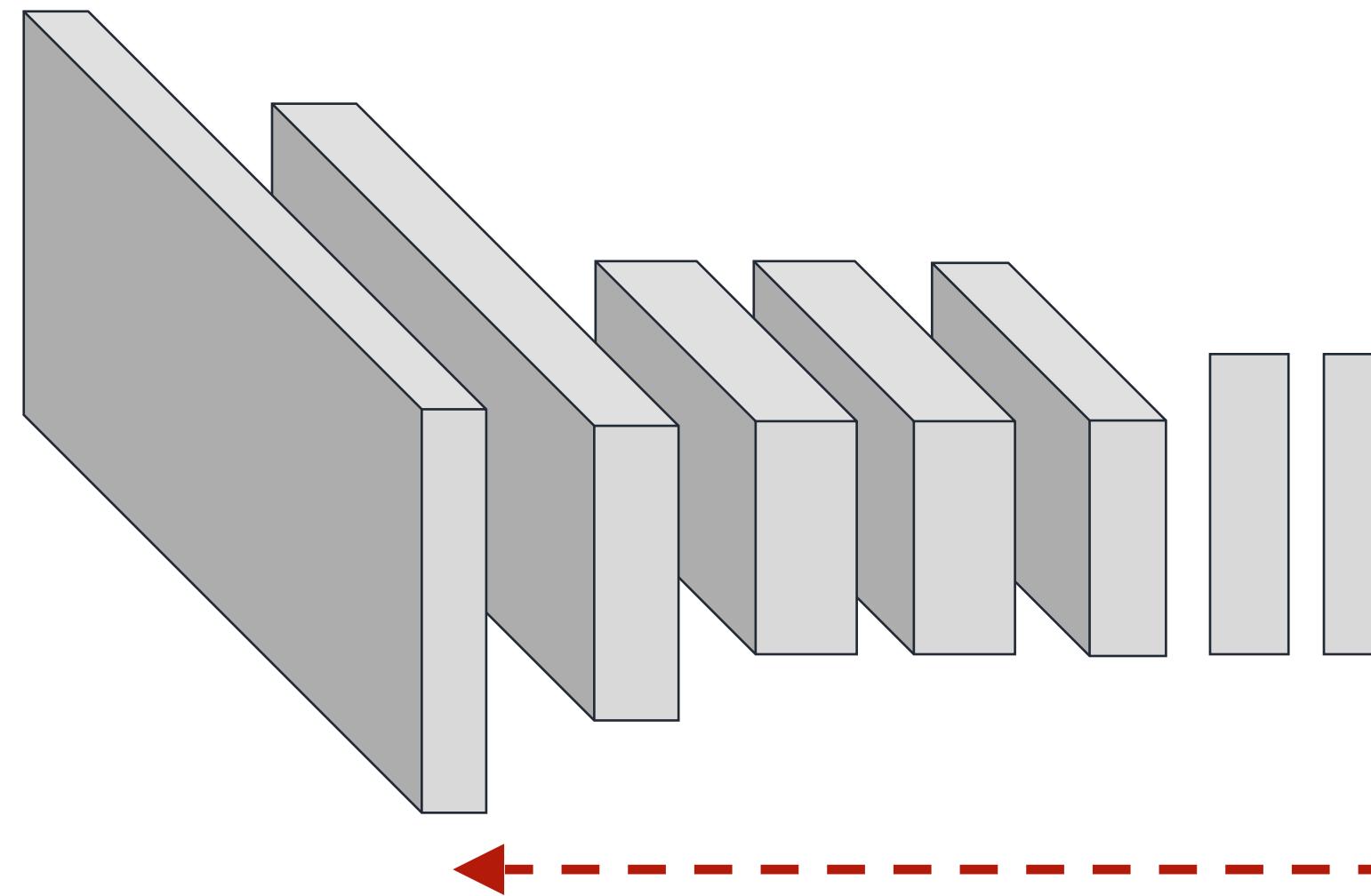
Jigsaw puzzles
(Noorozi & Favaro, 2016)

Jigsaw Puzzles

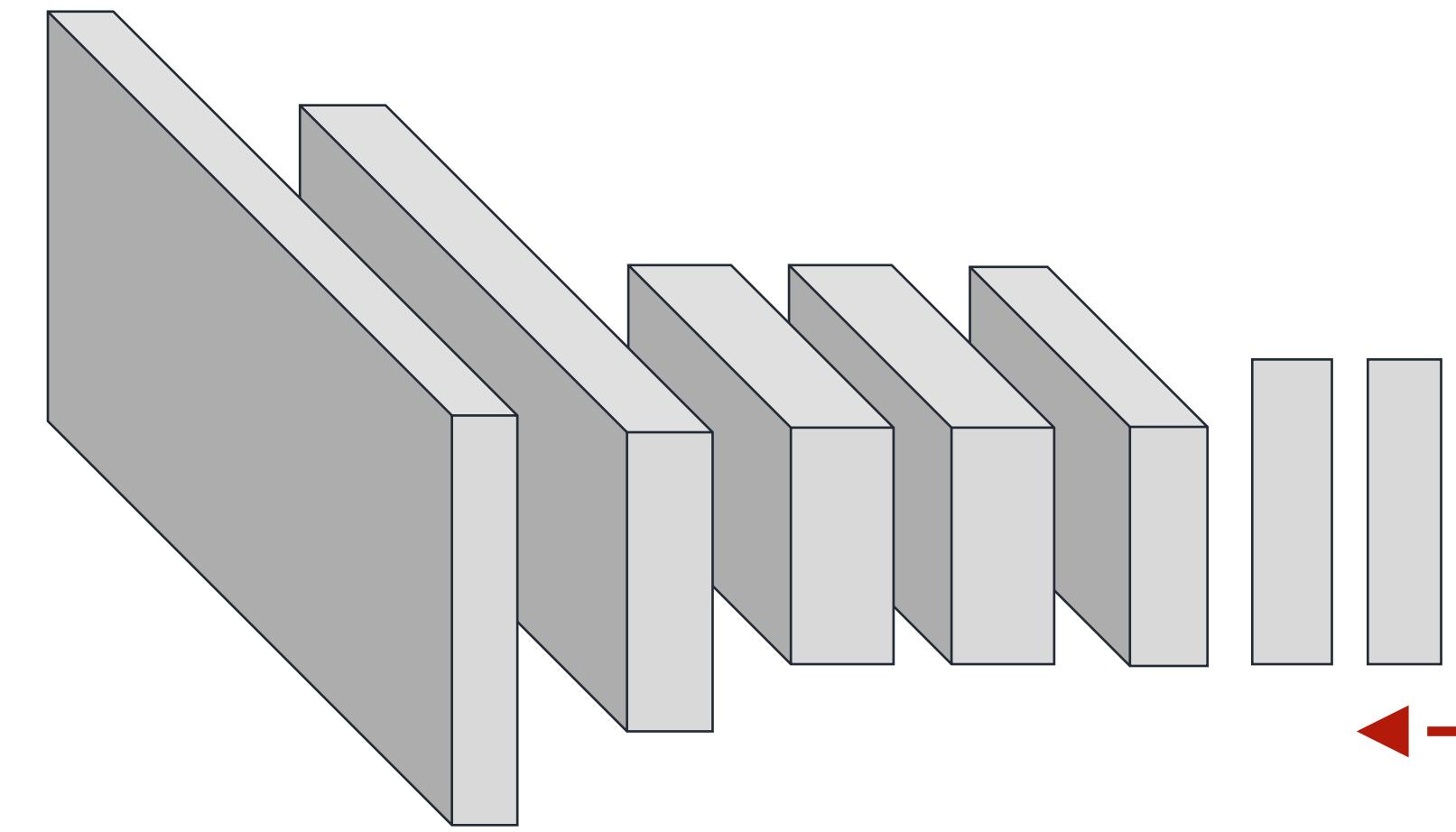


- Use $N=9$ patches
- In practice, use a subset of permutations
- E.g. 100 from $9!$
- Each patch is processed independently
- N-way ConvNet (shared params)
- Problem Complexity
 - Size of subset

Evaluation – fine-tuning vs. linear classifier



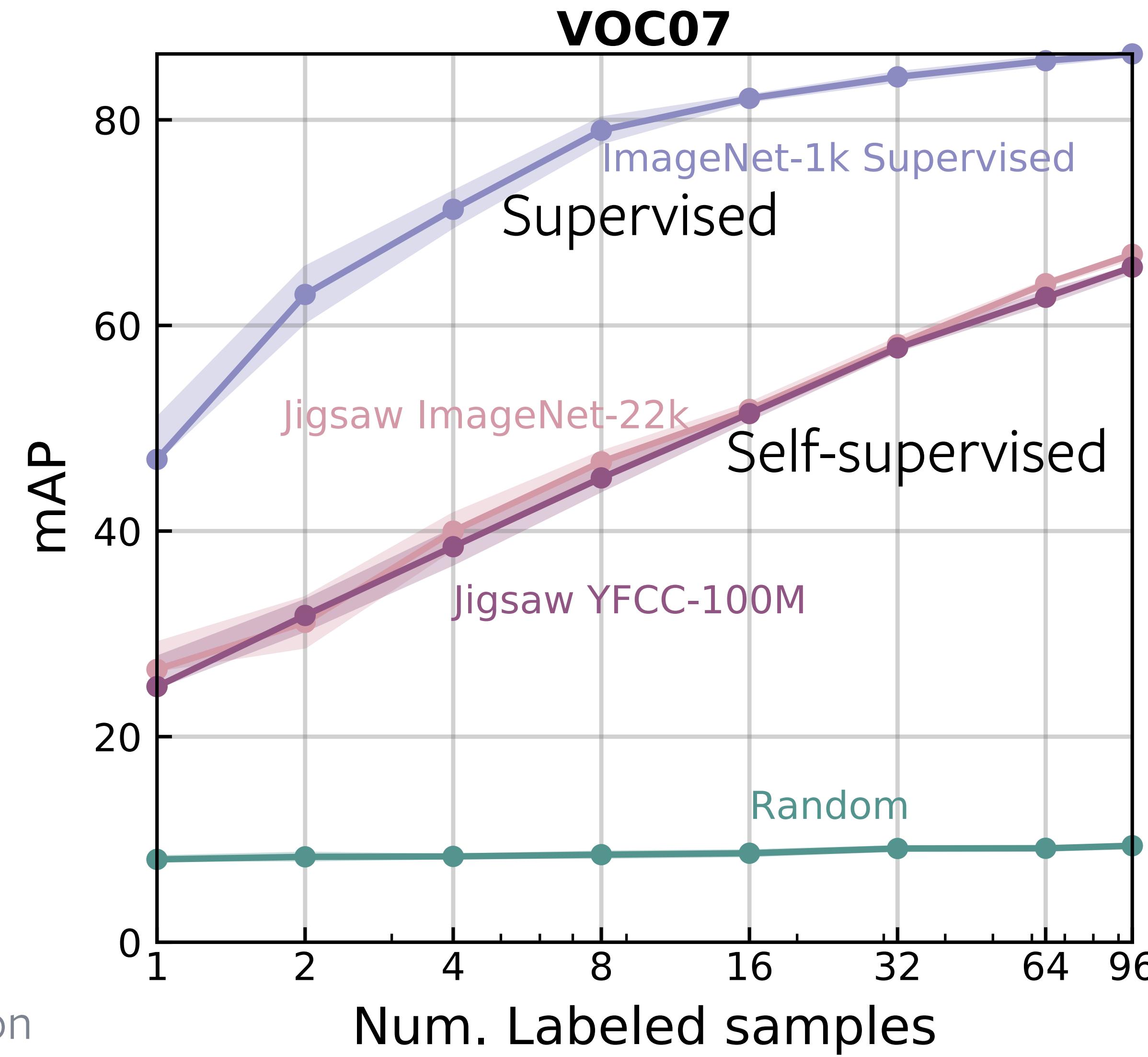
Fine-tune all layers



Linear classifier

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

Few shot learning



mAP = mean Average Precision
(Higher is better)

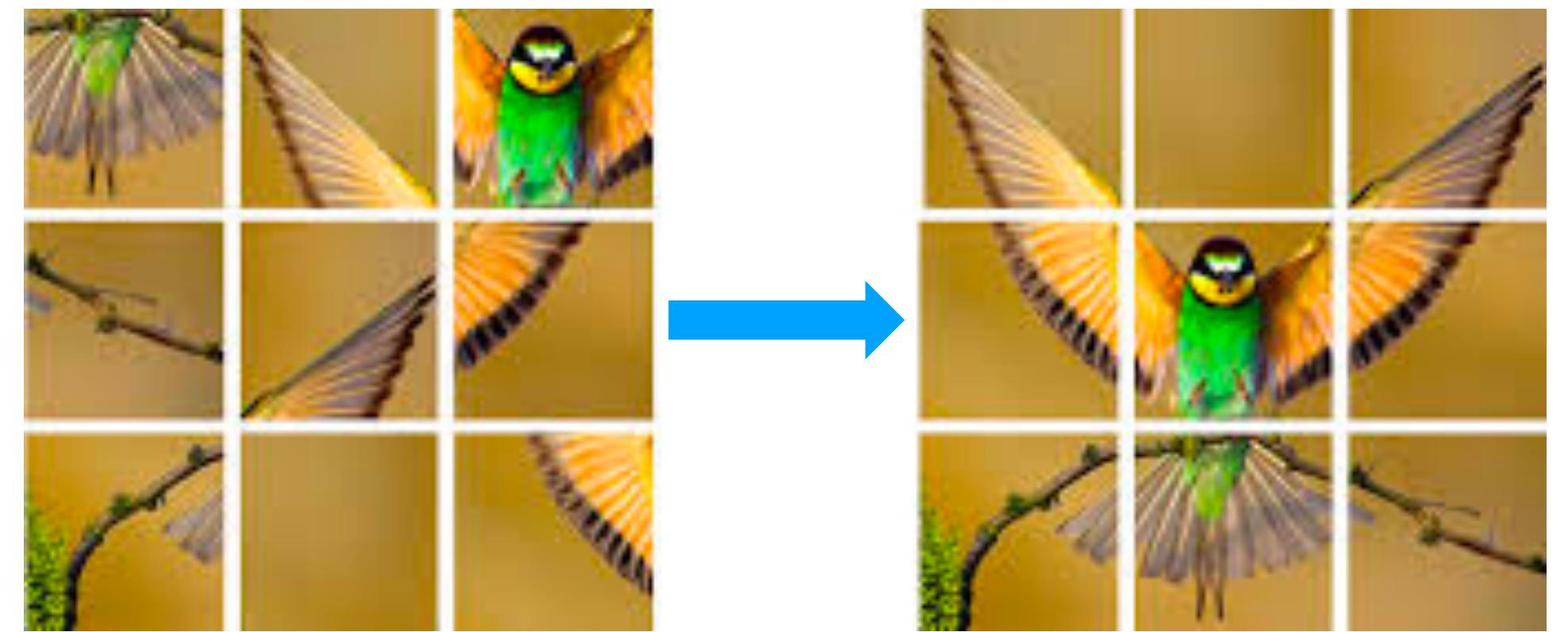
Self-supervised representations are not as sample efficient

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

Pretext tasks



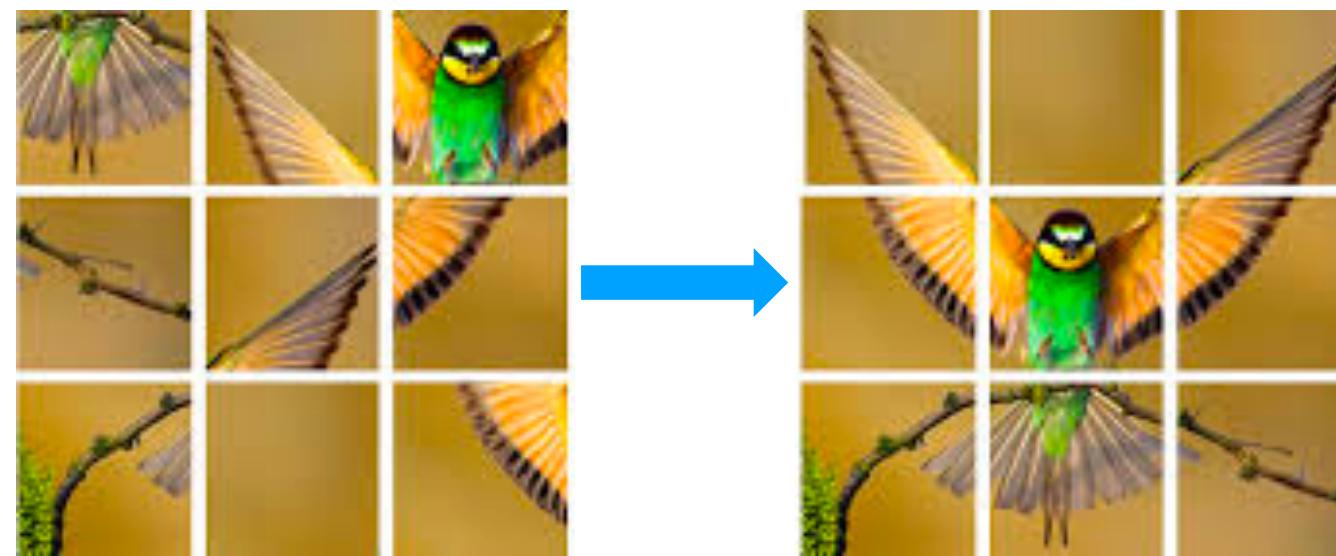
Rotation
(Gidaris et al., 2018)



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

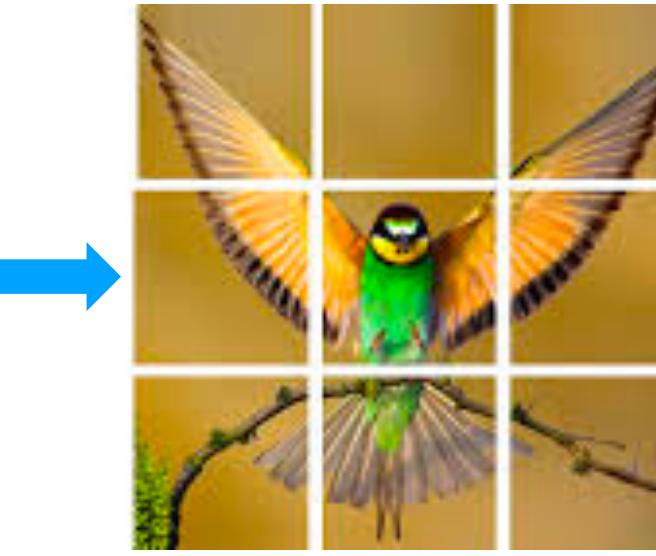


Transfer Tasks



The hope of generalization

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



#sun #nofilter #fun
#tree #aruba

Pre-training

Weak or self-supervised

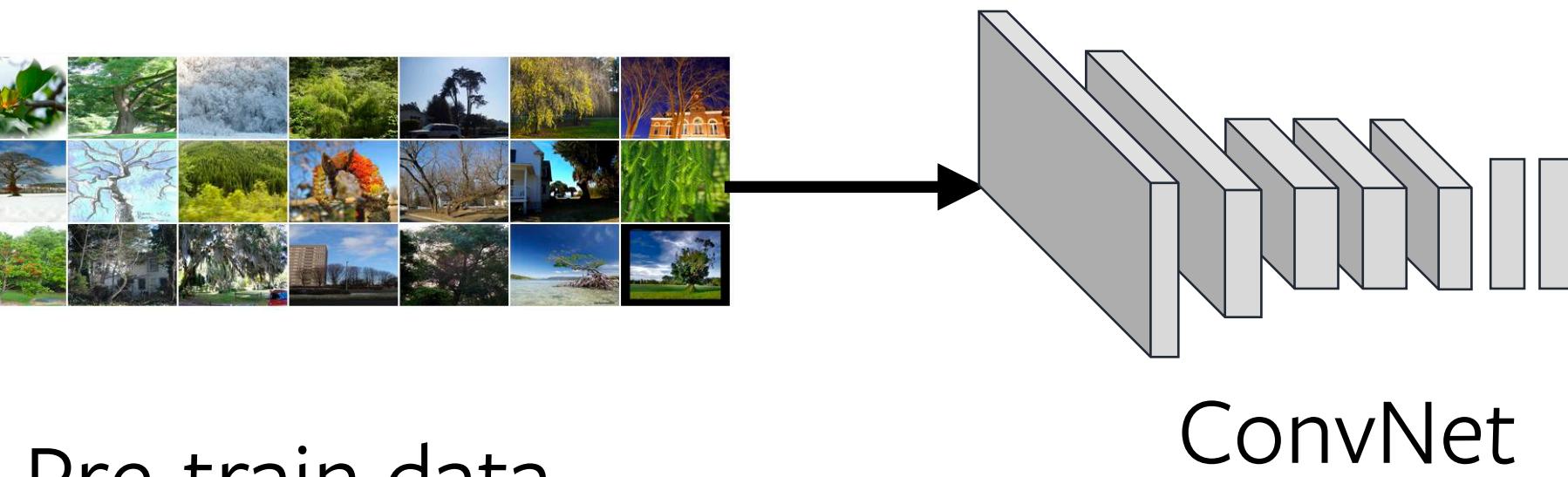


Transfer Tasks

Why should solving Jigsaw puzzles teach about "semantics"?

Why should performing a non semantic task produce good features?

The hope of generalization ... ?

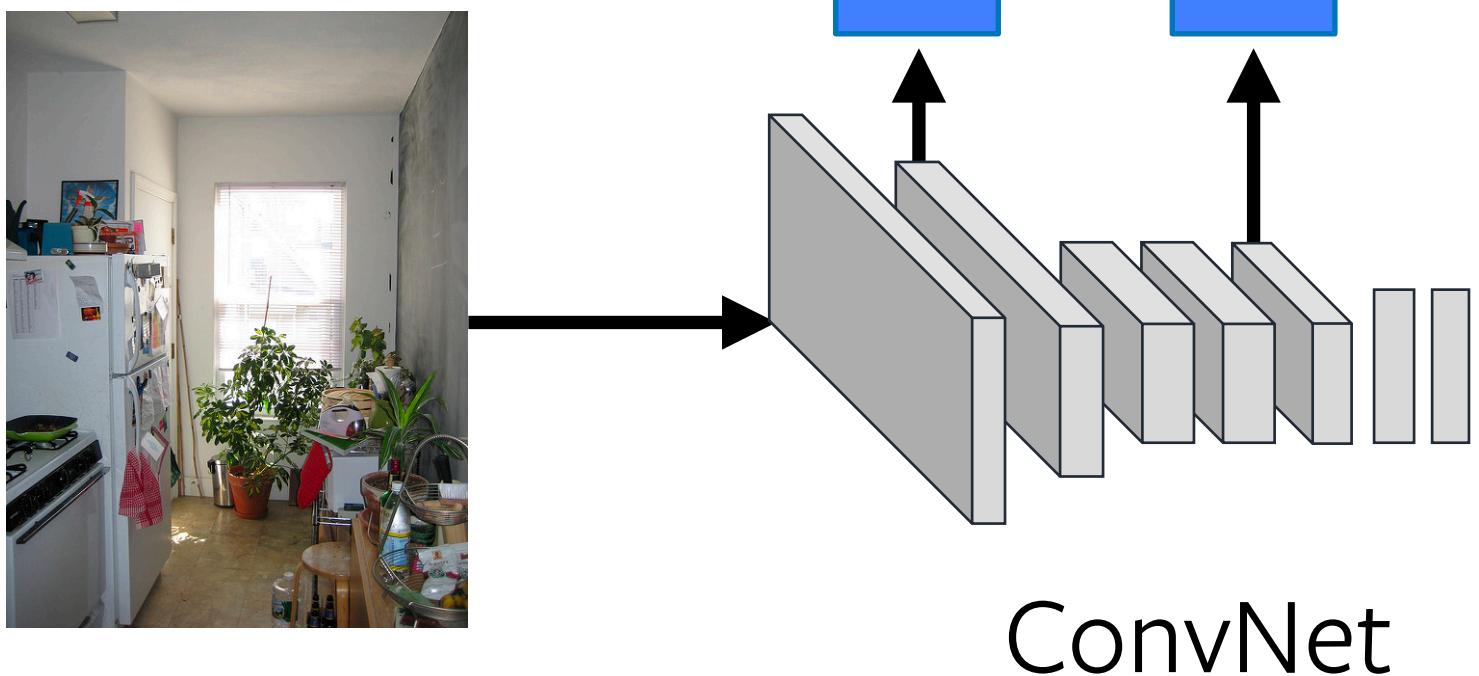


Pre-train data

ConvNet

Jigsaw

Linear classifiers on
"fixed" features

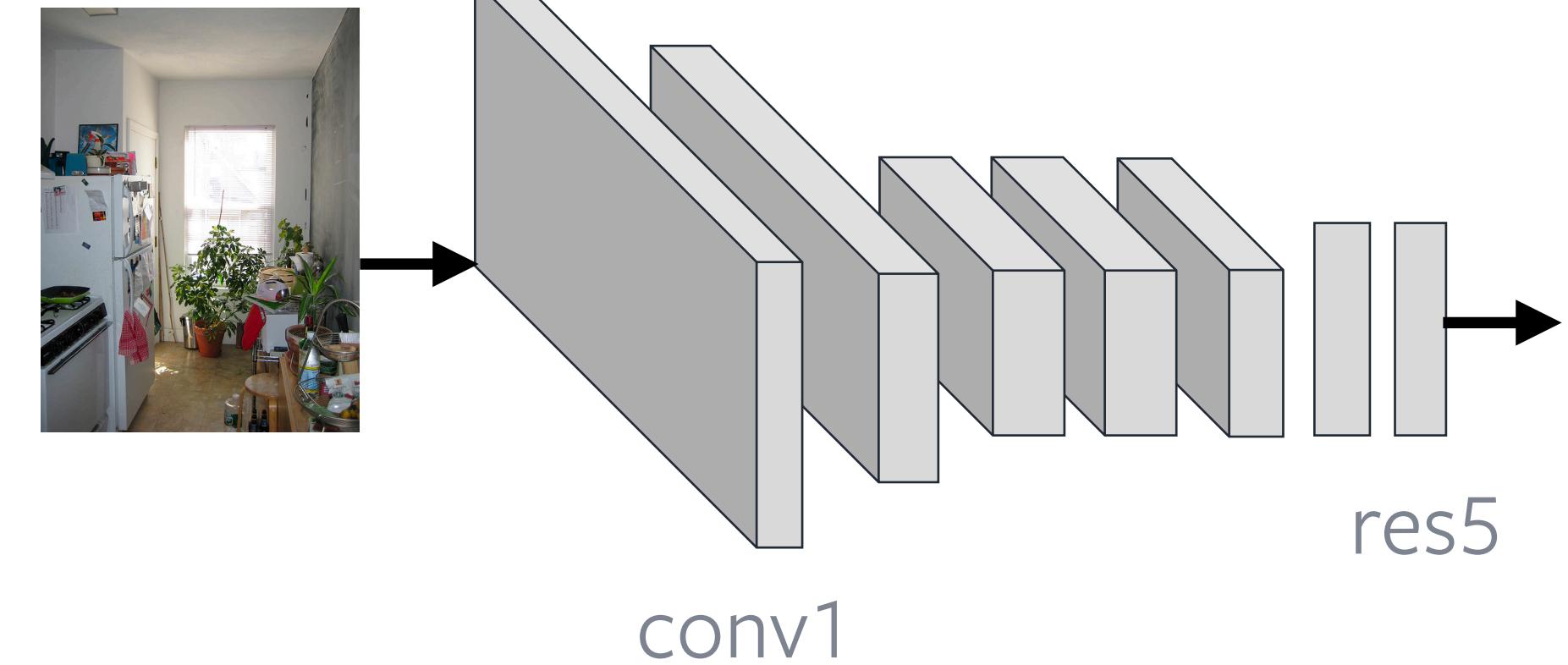
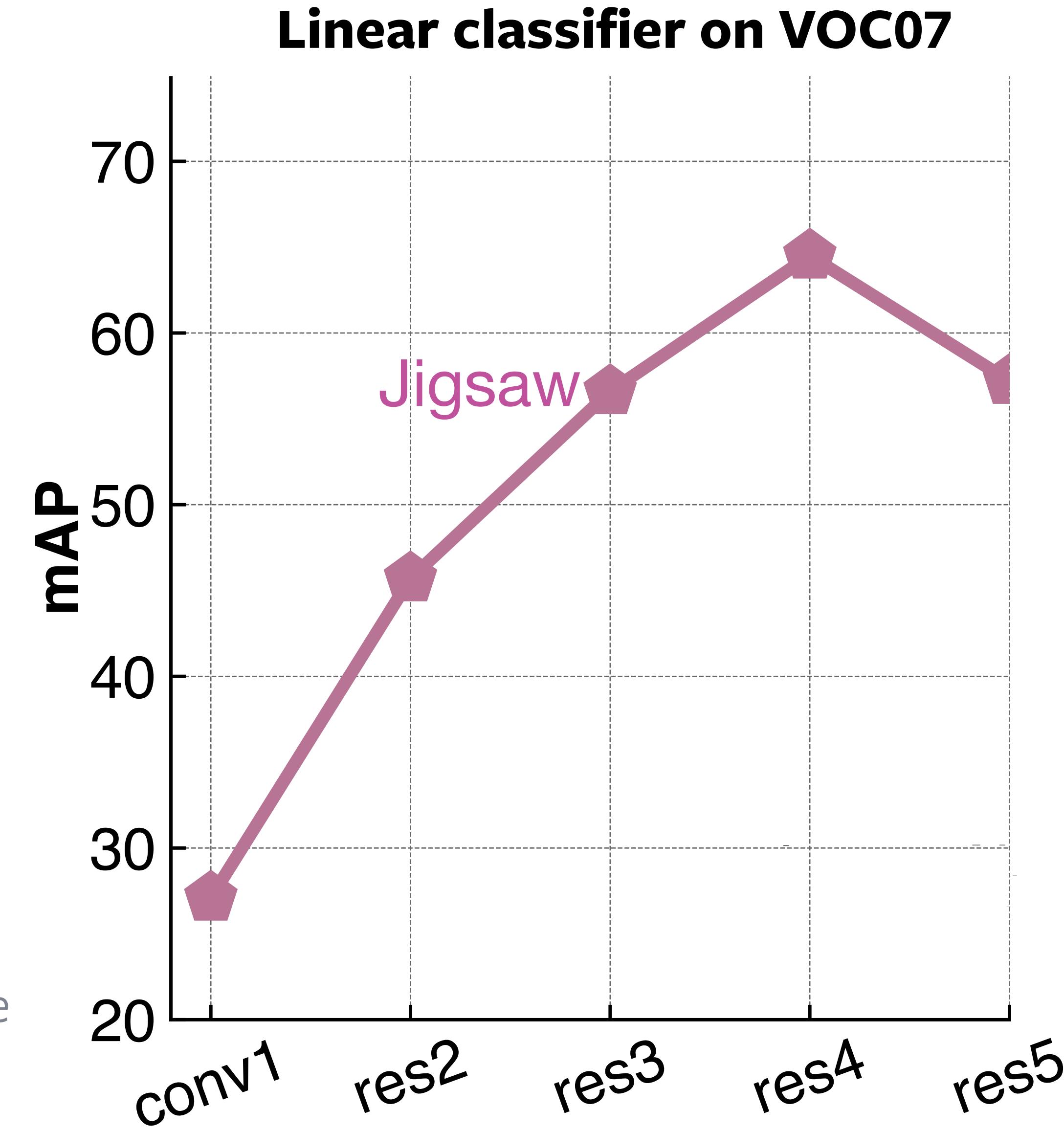


Pre-training

Weak or self-supervised

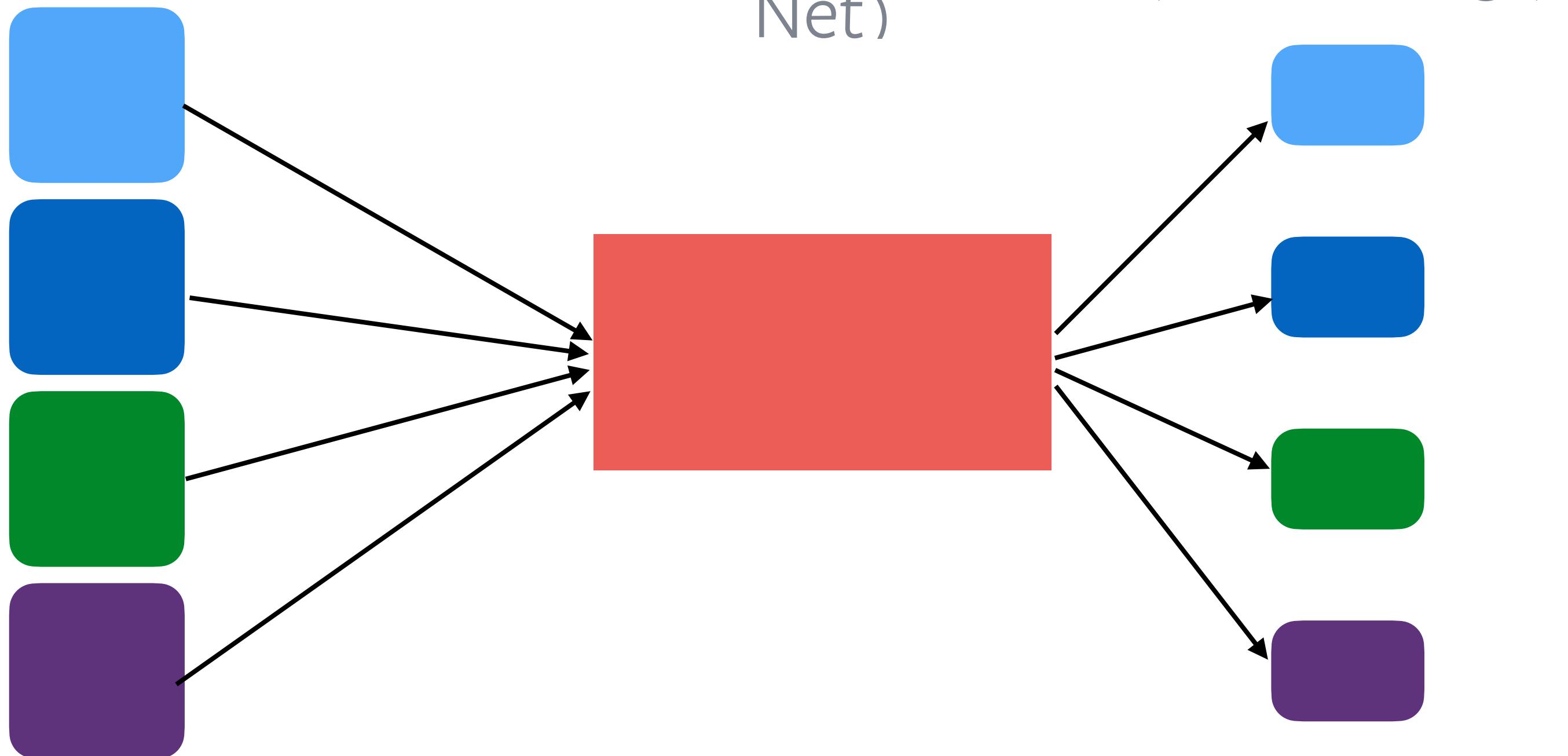
Transfer

Higher layers do not generalize ...



Contrastive Learning

Related and
Unrelated
Images



Shared
network
(Siamese
Net)

Image
Features
(Embeddings)

Loss Function

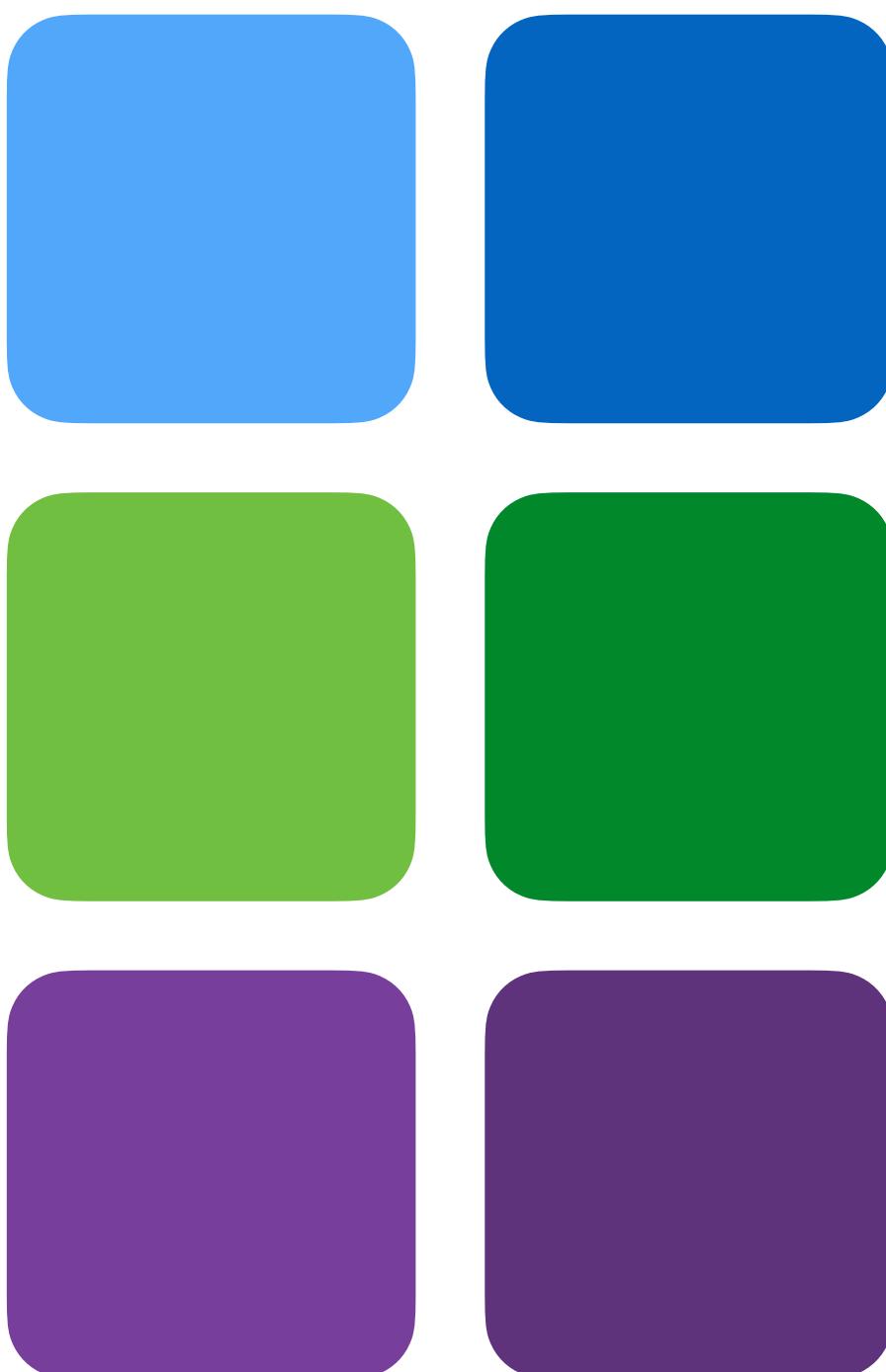
Embeddings from related images should be closer than embeddings from unrelated images

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

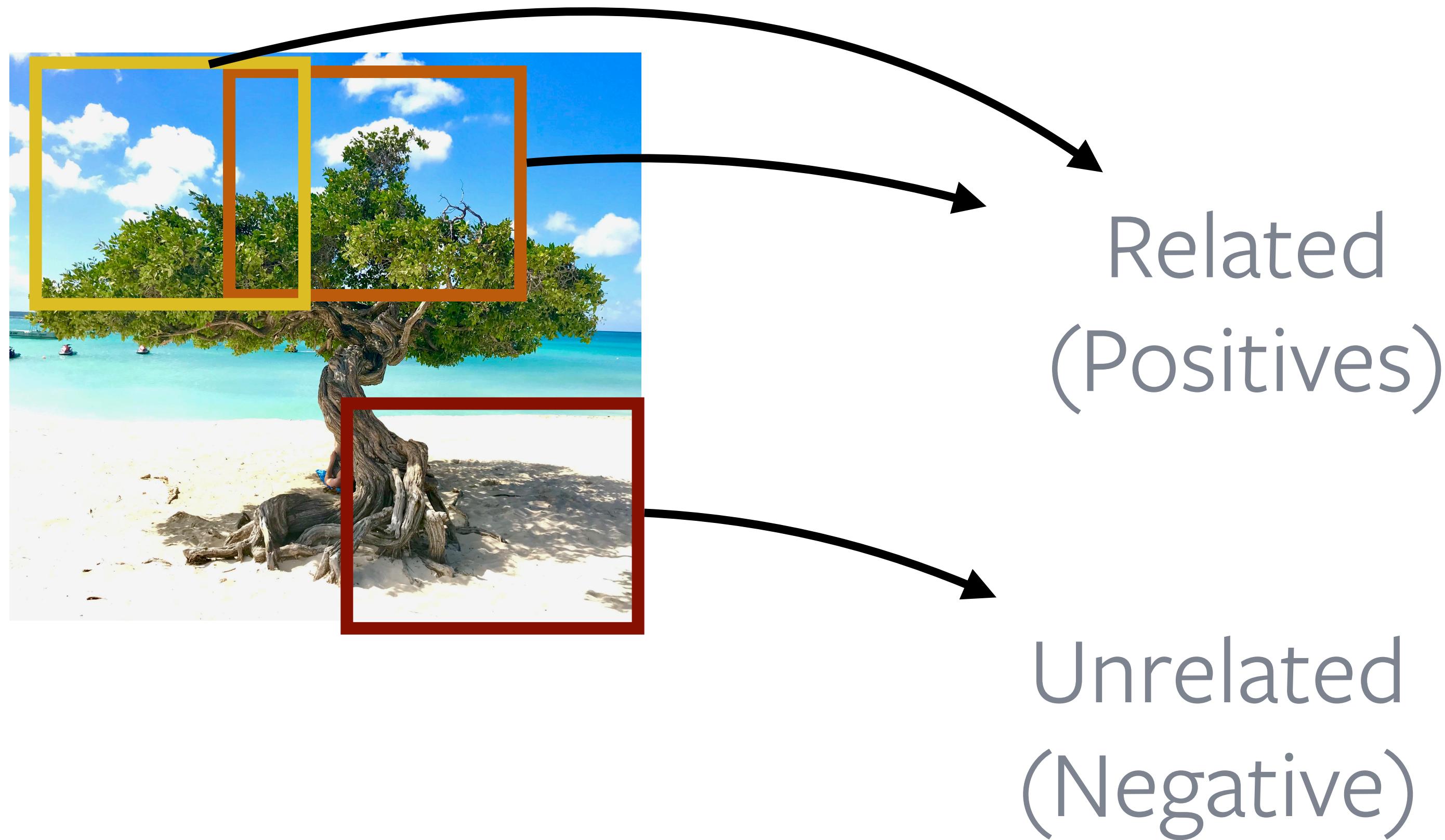
Contrastive Learning

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

Related and Unrelated
Images



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

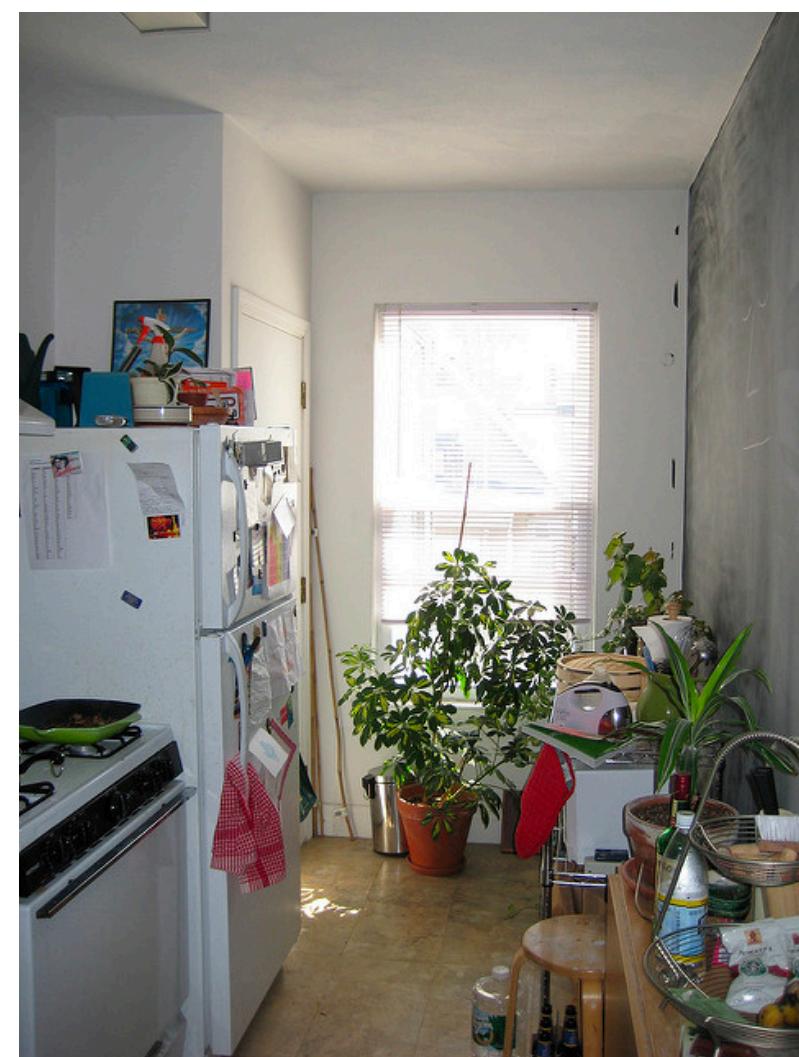
Misra & van der Maaten, 2019, PIRL

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

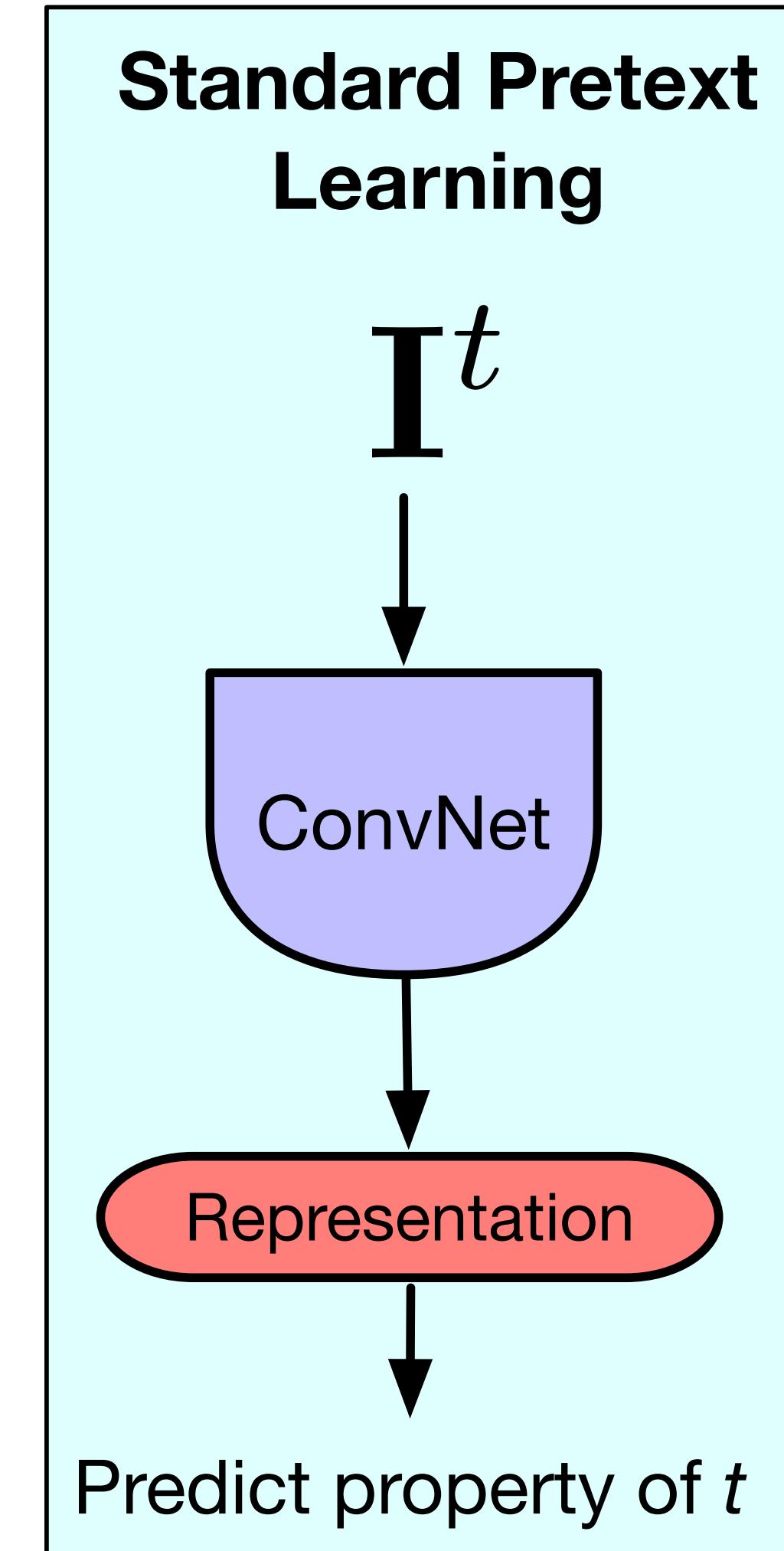
Pretext Image Transform



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)

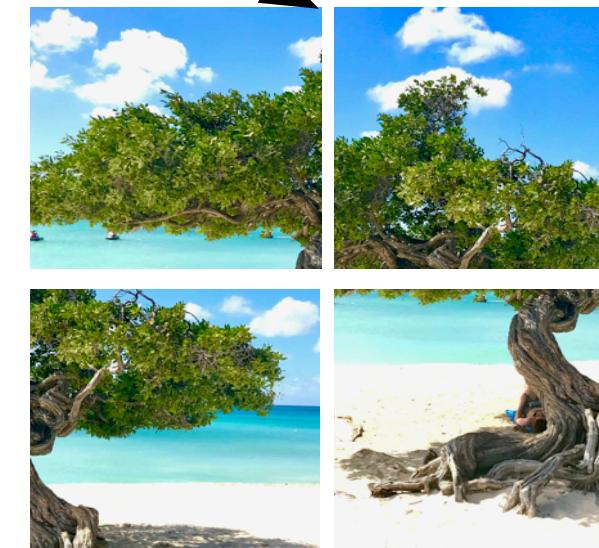
- Final layer representations must carry information about t
- Representations "covary" with t

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

Pretext Image Transform

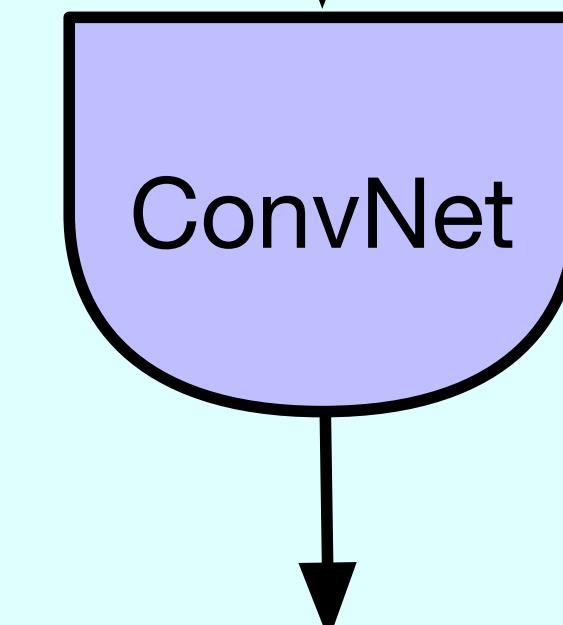


I
Transform t
 I^t



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"



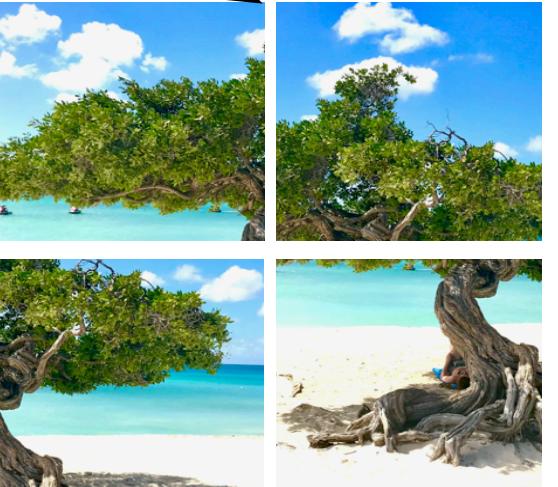
Pretext-Invariant Representation Learning (PIRL)

- Be invariant to t
- 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

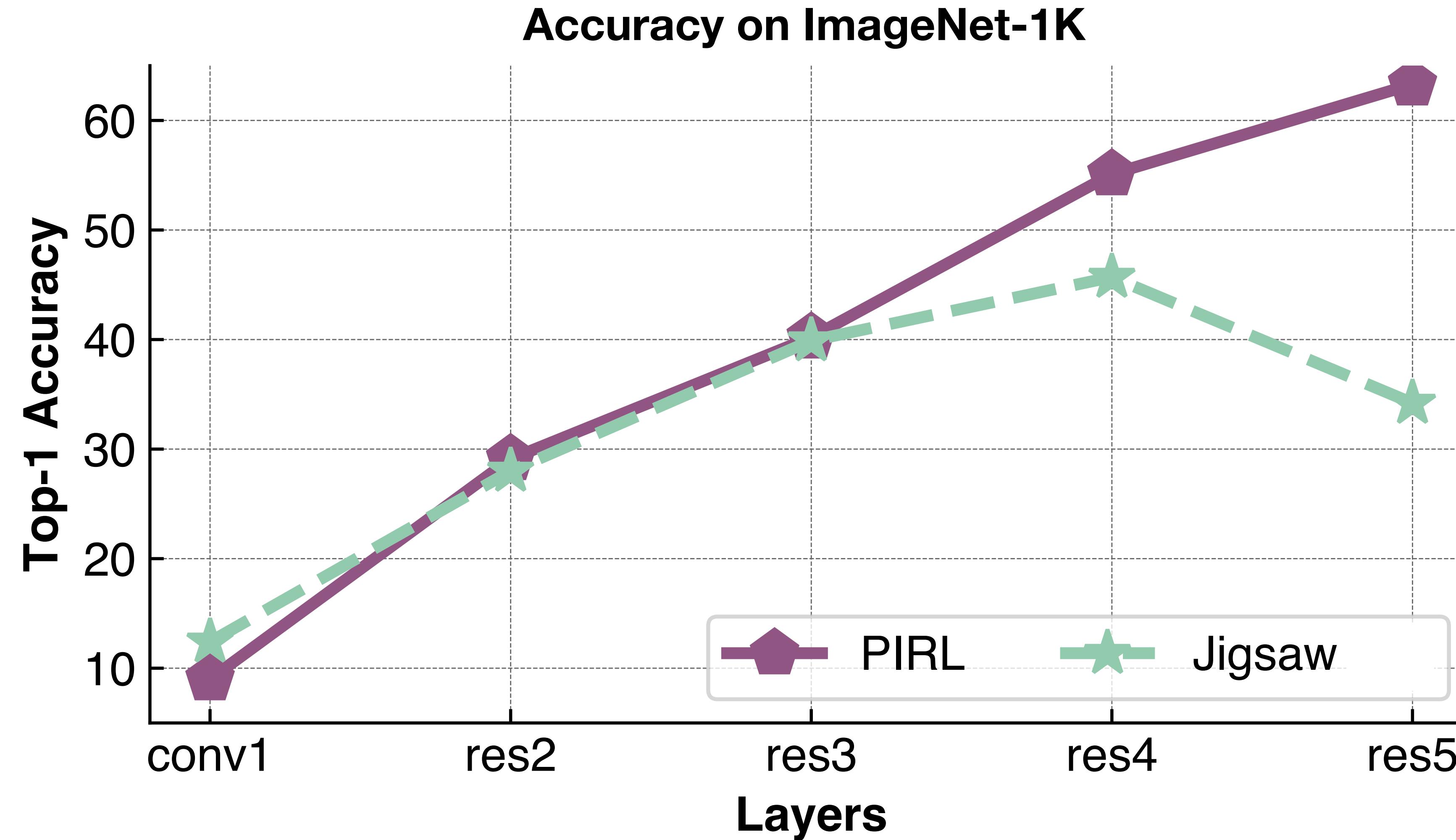
Encourage to be similar

Semi-supervised Learning

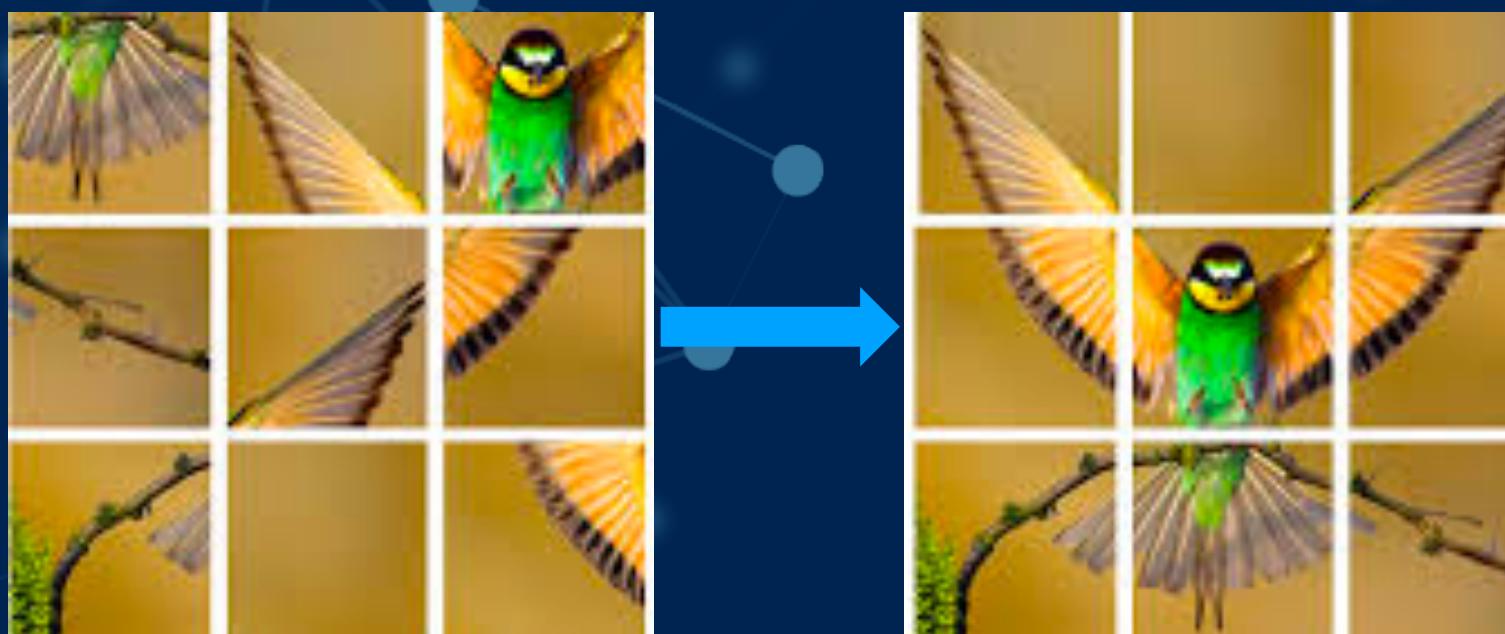
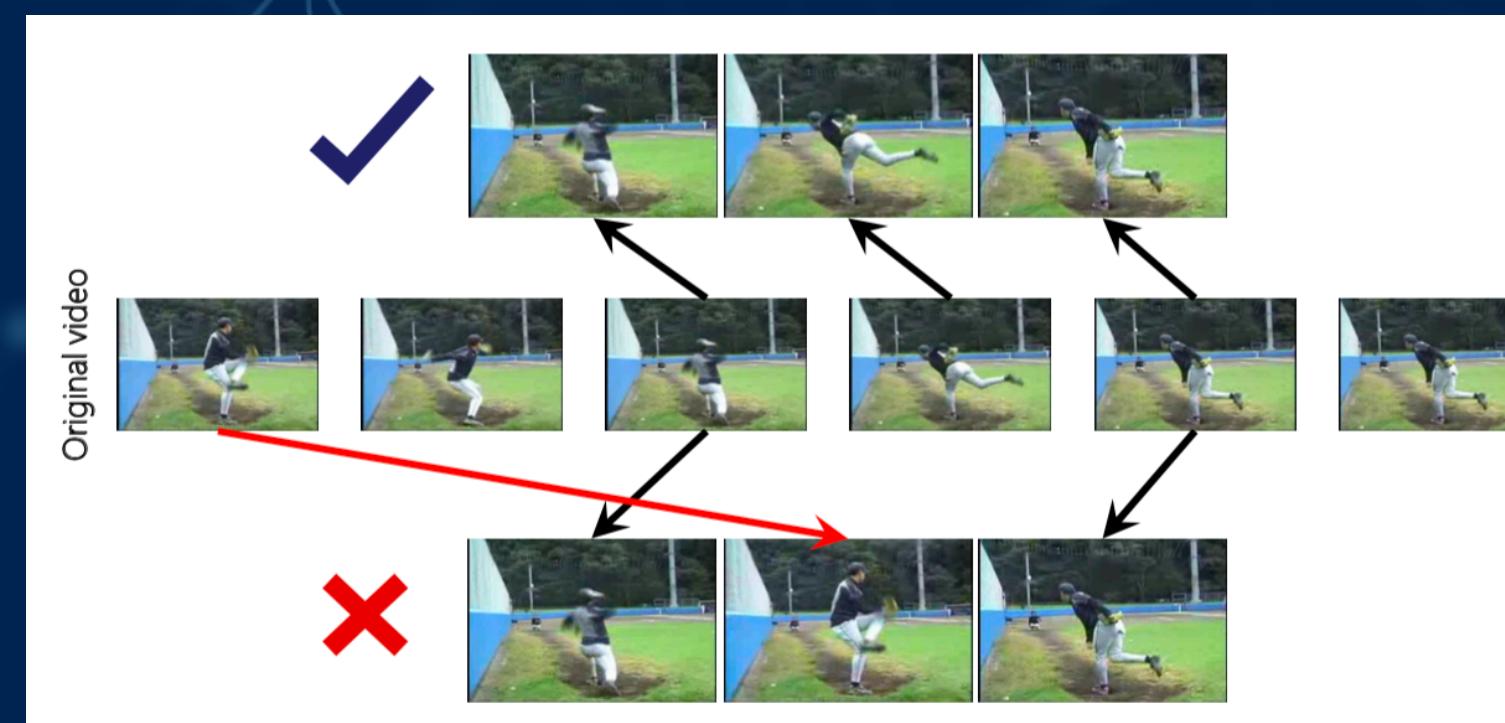
- Fine-tune on fraction of labeled data from **ImageNet-1K**

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	83.8
PIRL	57.2	83.8

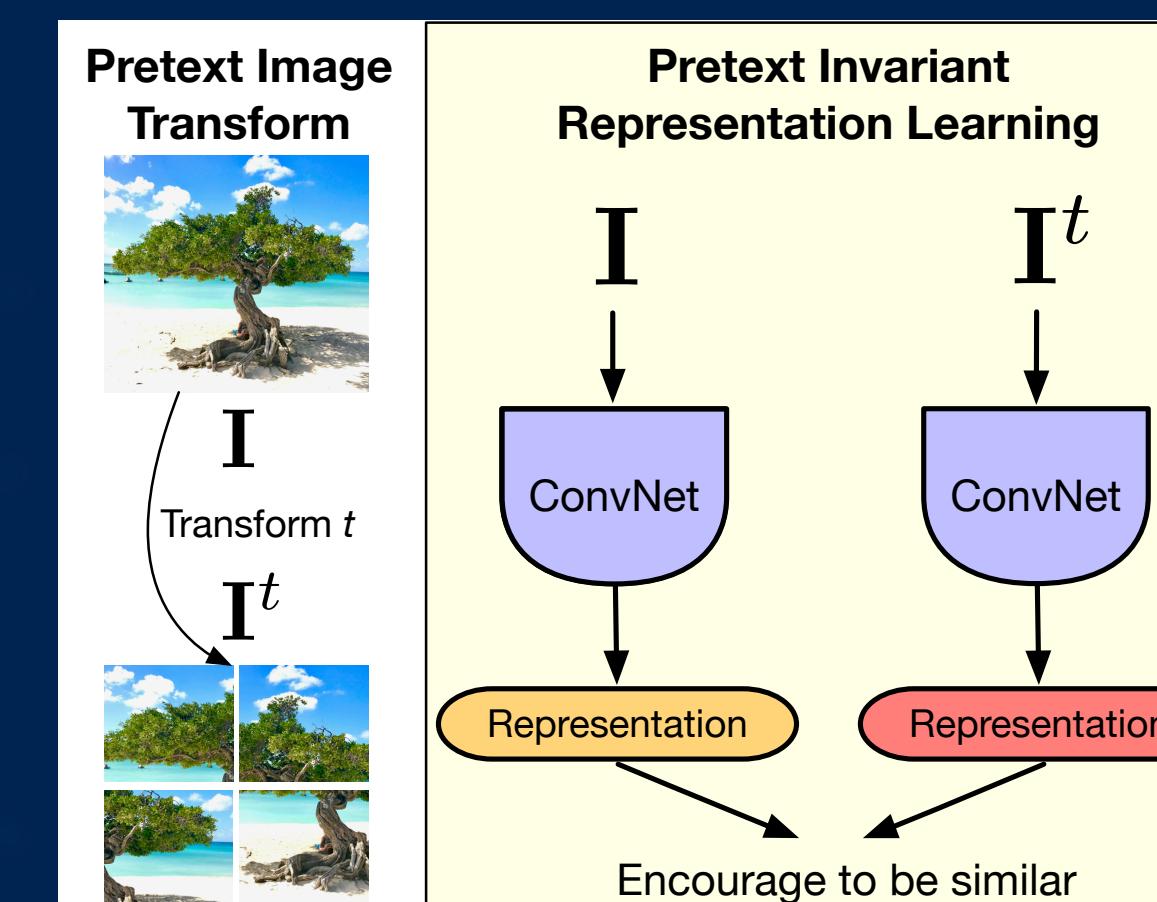
Semantic Features?



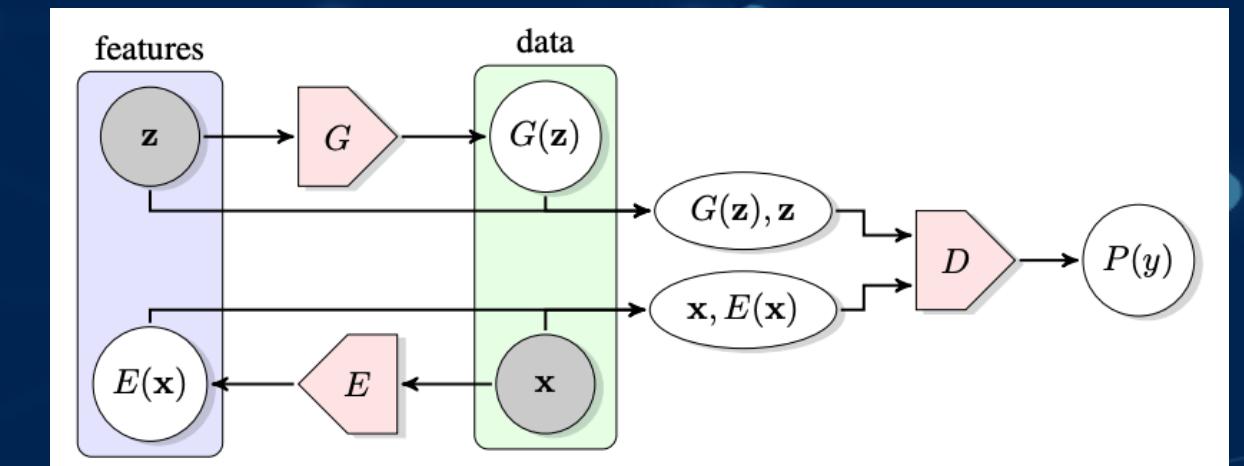
Pretext tasks



Contrastive



Generative



AutoEncoder,
VAE, GAN,
BiGAN

Predict more information