

# 11. Representation Learning

## GEV6135 Deep Learning for Visual Recognition and Applications

**Kibok Lee**  
Assistant Professor of  
Applied Statistics / Statistics and Data Science  
Nov 24, 2022



# Assignment 6

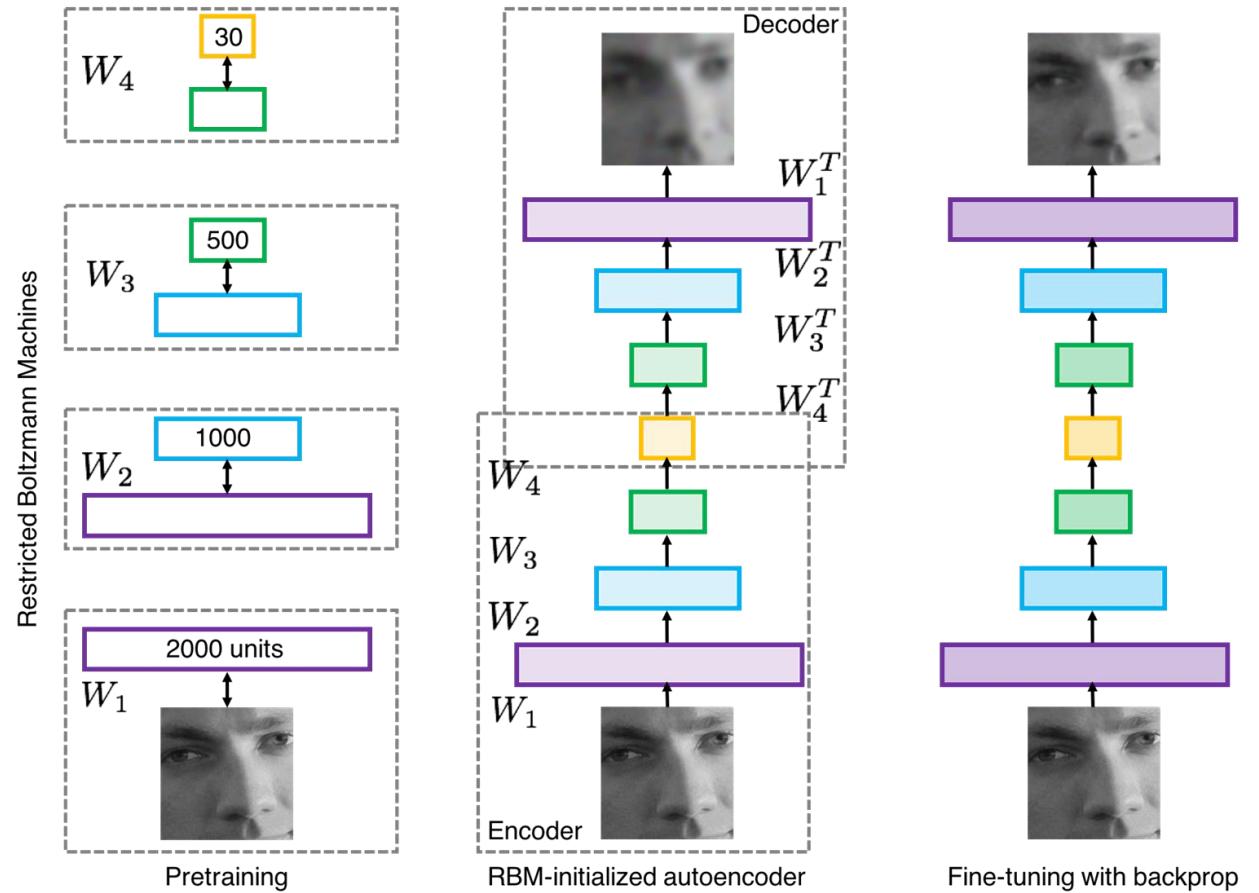
- Due **Friday 12/2, 11:59pm KST**
- Convolutional networks
  - Modularized implementation (loss is given!)
  - BatchNorm is given, but should be plugged into DeepConvNet
- Before submitting your work, we recommend you
  - Re-download clean files
  - Copy-paste your solution to clean py
  - Re-run clean ipynb only once
- If you feel difficult, consider to take **option 2.**

# Outline

- Trends of Representation Learning by Decade
- Supervised Learning for Transfer Learning
- Self-Supervised Learning (Roughly, ~ 2019)
- Self-Supervised Learning (Roughly, 2020 ~)

# Representation Learning in the 2000s

- To make DNNs working
  - Training DNNs was challenging
  - Not a mainstream research topic at the time
- Layer-wise RBM
- Sparse coding
- Autoencoding



Hinton and Salakhutdinov, "Reducing the Dimensionality of Data with Neural Networks", Science 2006

Bengio et al, "Greedy Layer-Wise Training of Deep Networks", NeurIPS 2007

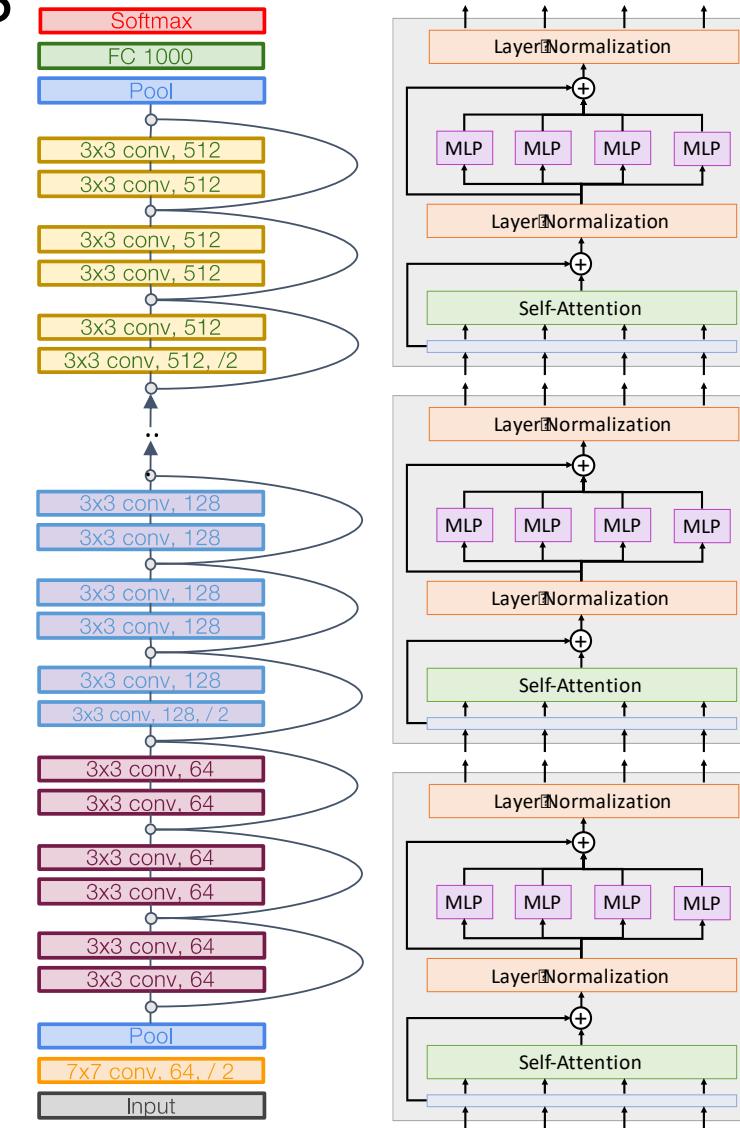
Lee et al, "Sparse deep belief net model for visual area V2", NeurIPS 2007

Lee et al, "Convolutional Deep Belief Networks for Scalable Unsupervised Learning of Hierarchical Representations", ICML 2009

Vincent et al, "Stacked Denoising Autoencoders", JMLR 2010

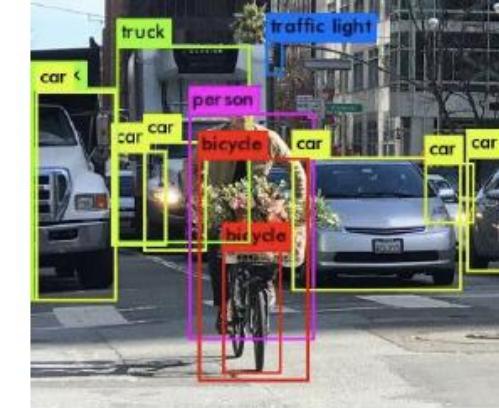
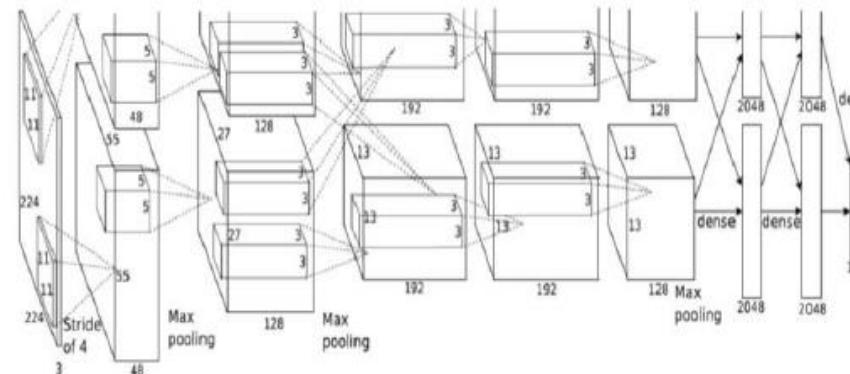
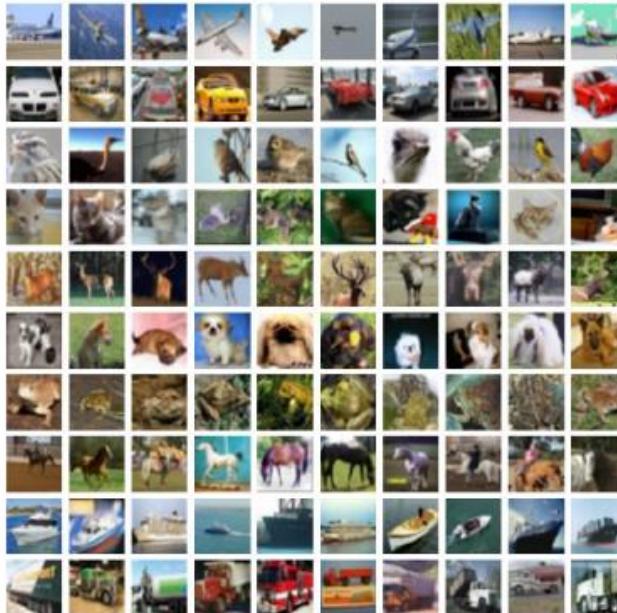
# Representation Learning in the 2020s

- To train DNNs efficiently
  - Also better few-shot performance
- Supervised Learning for Transfer Learning
  - State-of-the-art representation learning in the 2010s
  - Still competitive if no unlabeled data in the 2020s
- Self-Supervised Learning
  - Taking advantage of huge unlabeled datasets
  - Outperforms supervised learning in the 2020s



# How Deep Learning Works

airplane  
automobile  
bird  
cat  
deer  
dog  
frog  
horse  
ship  
truck



- Define the set of (visual) concepts to be learned
- Collect diverse and large number of examples for each of them
- Train a deep model for several GPU hours (or days)

# Supervised Learning Is Expensive

- Supervised learning requires to curate large human-annotated datasets
  - Annotating + cleaning raw data
- Time consuming and expensive



- Annotating such image: ~ 1.5h
- Error prone (human mistakes)

# Supervised Learning Is Expensive

- Assume you want to label **1M** images. How much will it cost?

(1,000,000 images) (Small to medium sized dataset)  
× (10 seconds/image) (Fast annotation)  
× (1/3600 hours/second)  
× (10k won/hour) (Low wage paid to annotator)  
**= 2777,7778 won**

(Other assumptions: one annotator per image, no benefits / payroll tax / crowdsourcing fee for annotators; not accounting for time to set up tasks for annotators, etc. Real costs could easily be 3x this or more)

# Supervised Learning Is Expensive

- Assume you want to label **1B** images. How much will it cost?

(1,000,000,000 images)

(Large dataset)

× (10 seconds/image)

(Fast annotation)

× (1/3600 hours/second)

× (10k won/hour)

(Low wage paid to annotator)

= **277,777,778 won**

(Other assumptions: one annotator per image, no benefits / payroll tax / crowdsourcing fee for annotators; not accounting for time to set up tasks for annotators, etc. Real costs could easily be 3x this or more)

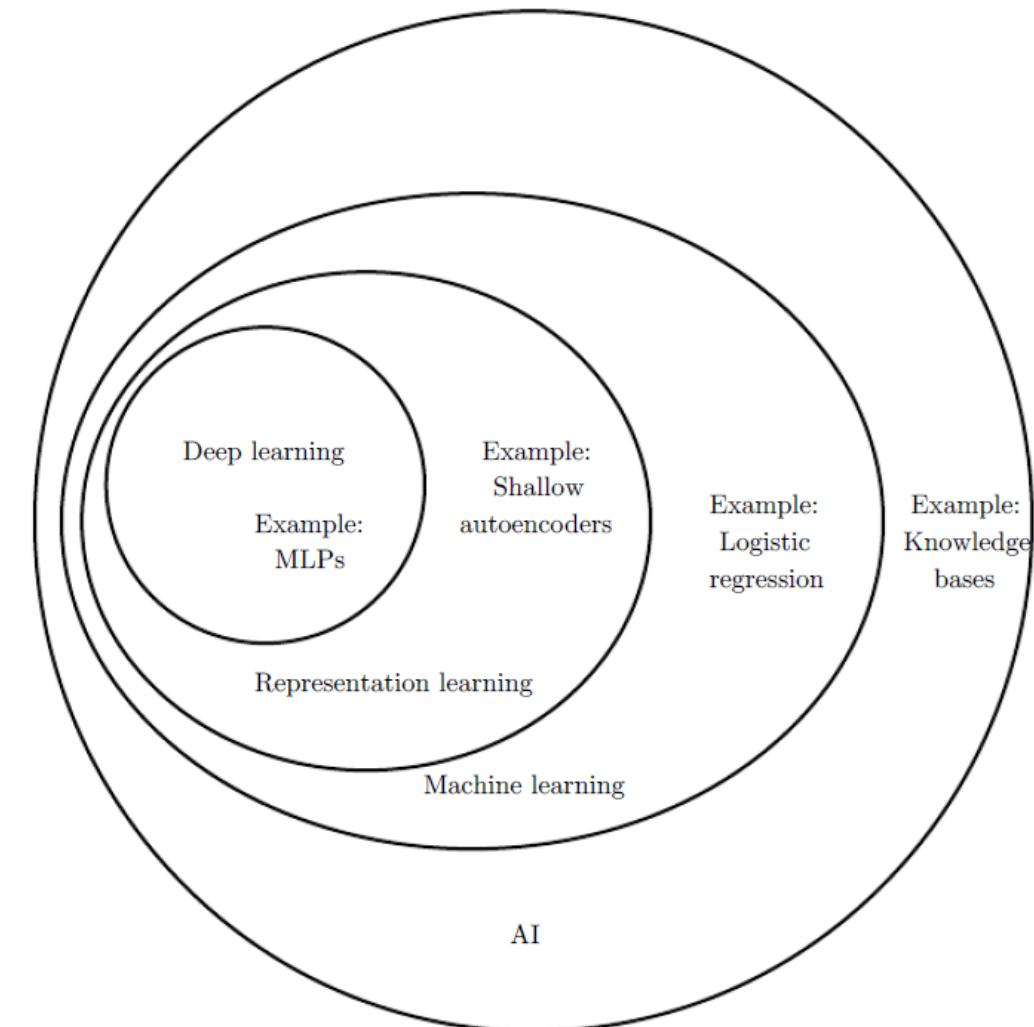
# Supervised Learning Is Not How We Learn

Babies don't get supervision  
for everything they see!



# Representation Learning

- “Good representations” are
  - reusable, smooth, spatially coherent, disentangled, hierarchical, semantically meaningful
- (Most) Deep learning models consist of backbone and task-specific layers
  - Backbone = a good feature extractor
    - C.f. before DL, hand-designed features
  - Task-specific layers = classifier/regressor
- **Q. Can we pretrain a good backbone for efficient deep learning?**



<https://dmitry.ai/t/topic/175>

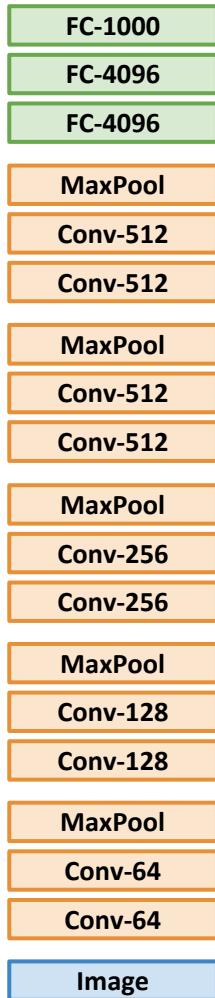
# Supervised Learning for Transfer Learning

# Supervised Learning for Transfer Learning

1. Pretrain on a supervised learning task
  - Usually ImageNet-1k classification task
2. Remove the task-specific layer
  - Usually keep all conv layers and remove the linear layer/all MLPs
  - This simple idea surprisingly works well!

# Transfer Learning with CNNs: Feature Extraction

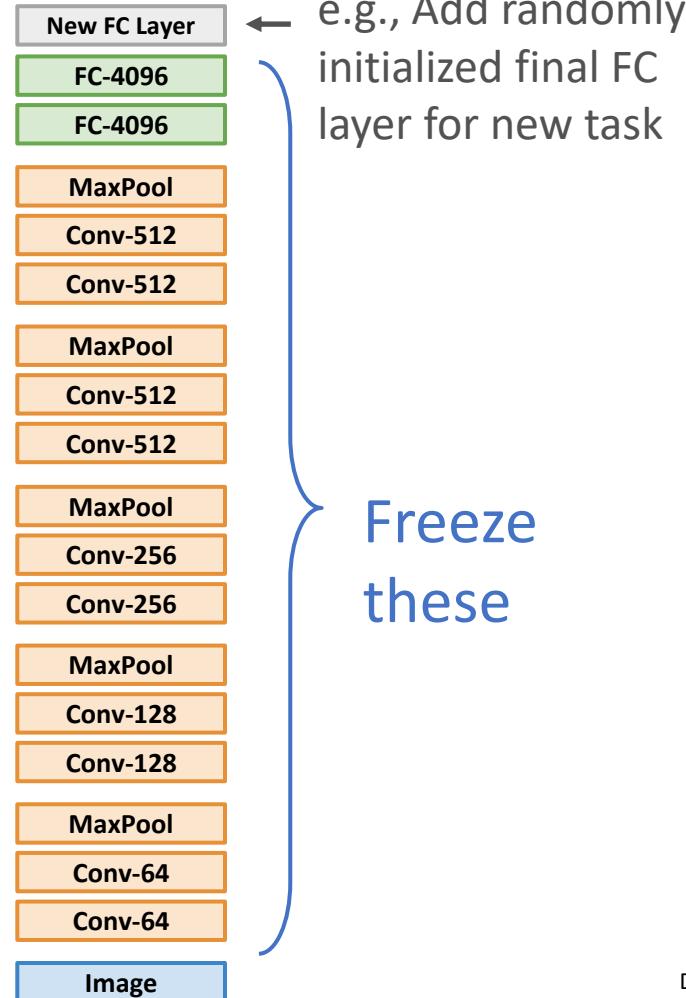
## 1. Train on ImageNet



Remove last layer

Initialize from ImageNet model

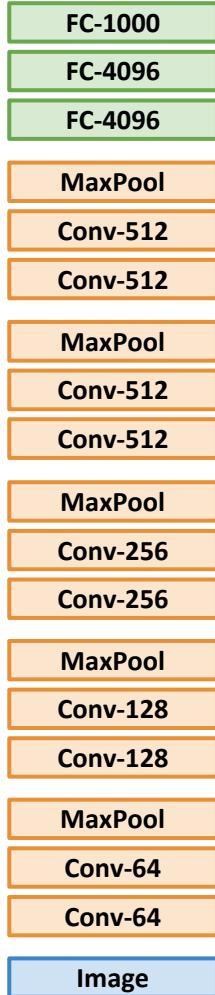
## 2. Use CNN as a feature extractor



Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014

# Transfer Learning with CNNs: Fine-Tuning

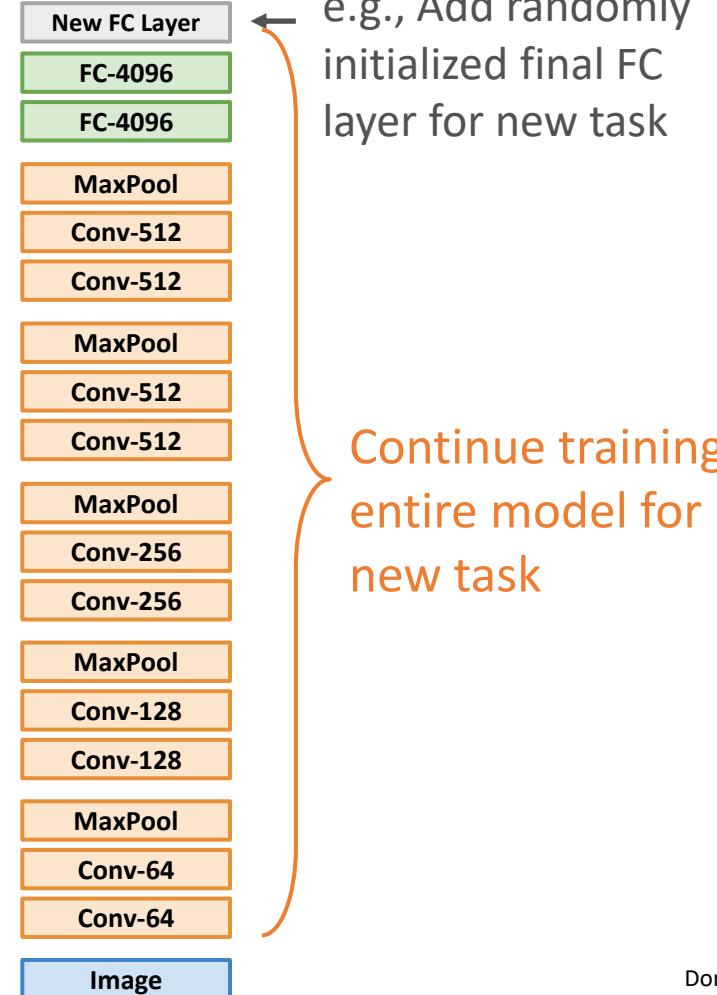
## 1. Train on ImageNet



Remove  
last layer

Initialize from  
ImageNet model

## 2. Fine-tune CNN



e.g., Add randomly initialized final FC layer for new task

Continue training entire model for new task

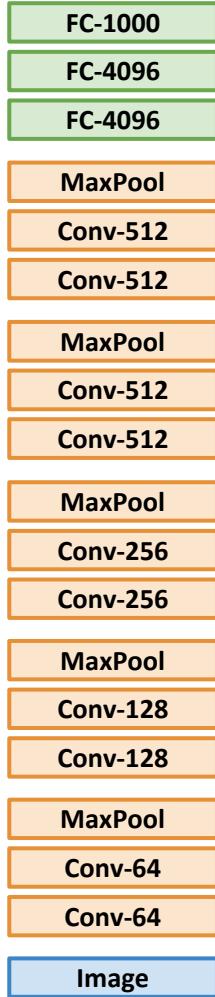
Some tricks:

- Train with frozen feature extraction first before fine-tuning
- Lower the learning rate: use  $\sim 1/10$  of LR used in original training
- Sometimes freeze lower layers to save computation
- Train with BatchNorm in “test” mode

Donahue et al, “DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition”, ICML 2014

# Transfer Learning with CNNs: Fine-Tuning

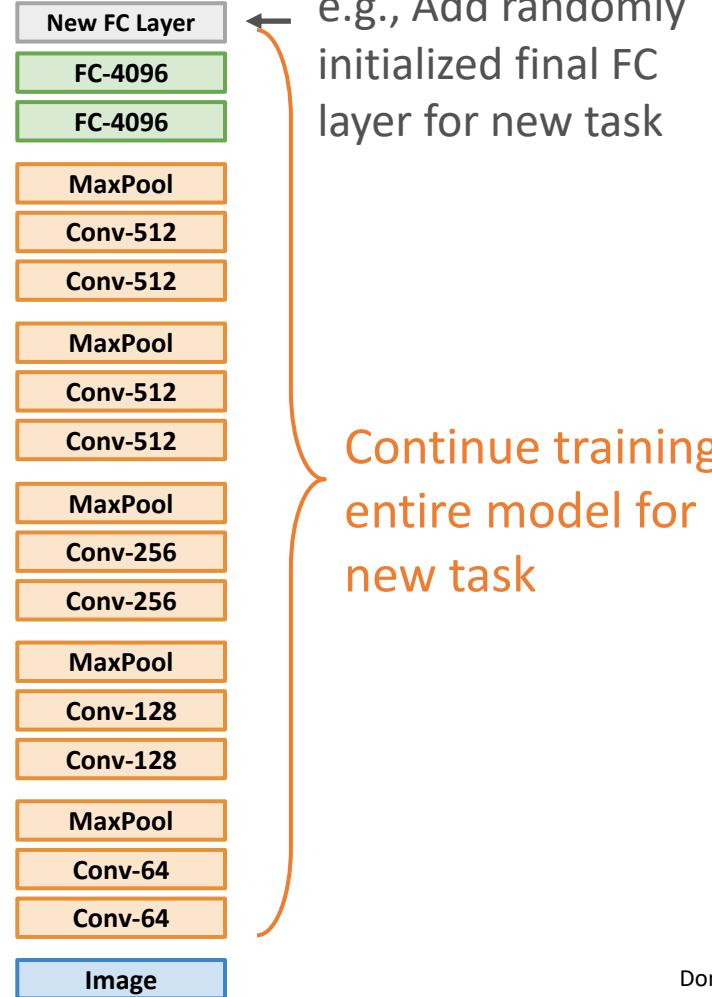
## 1. Train on ImageNet



Remove last layer

Initialize from ImageNet model

## 2. Fine-tune CNN



e.g., Add randomly initialized final FC layer for new task

Continue training entire model for new task

Compared with Feature Extraction, Fine-Tuning:

- Requires more data
- Is more computationally expensive
- Can give higher accuracies

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014

# Transfer Learning with CNNs



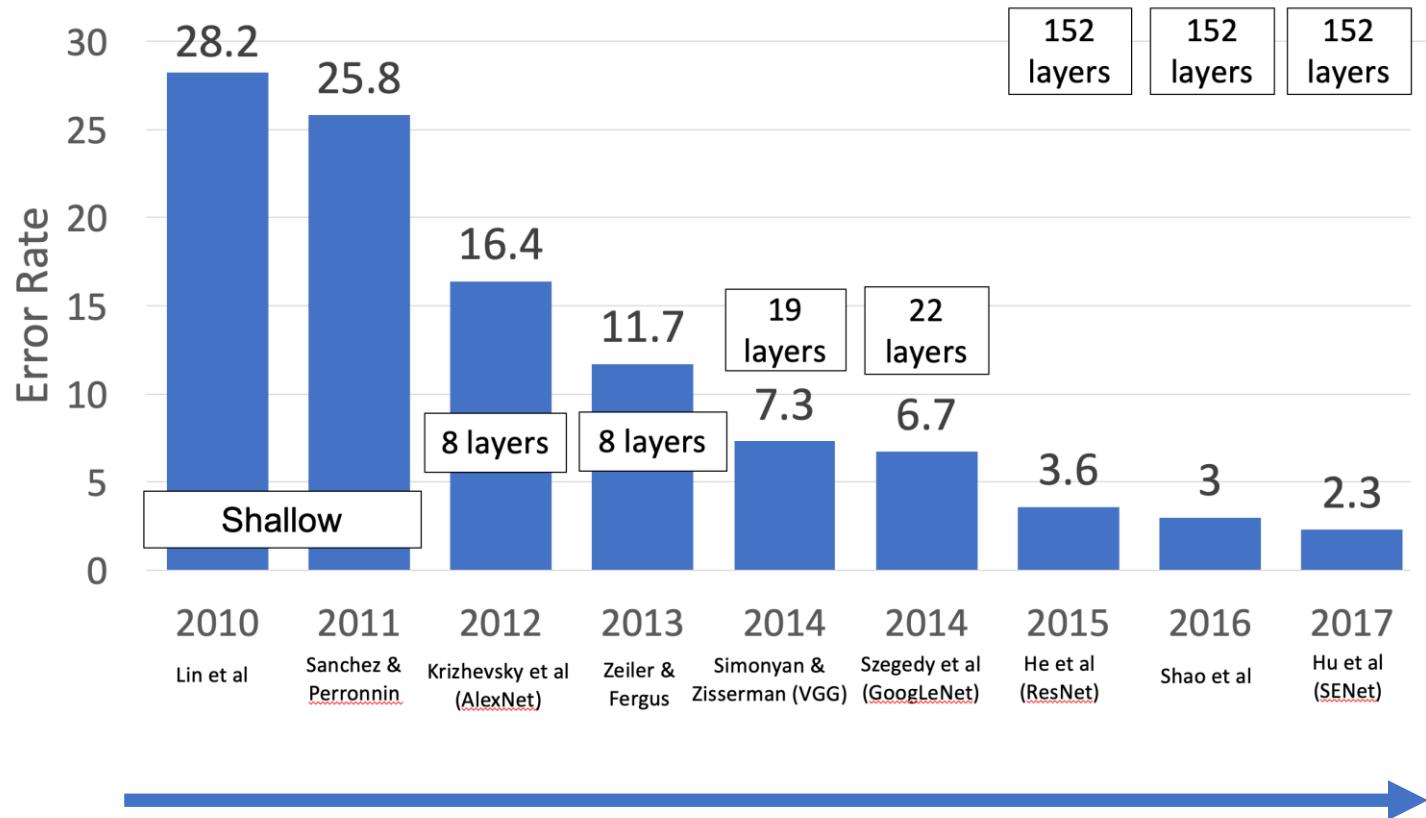
More specific

More generic

	<b>Dataset similar to ImageNet</b>	<b>Dataset very different from ImageNet</b>
<b>very little data (10s to 100s)</b>	Use Linear Classifier on top layer	You're in trouble... Try linear classifier from different stages
<b>quite a lot of data (100s to 1000s)</b>	Finetune a few layers	Finetune a larger number of layers

# Transfer Learning with CNNs: Architecture Matters!

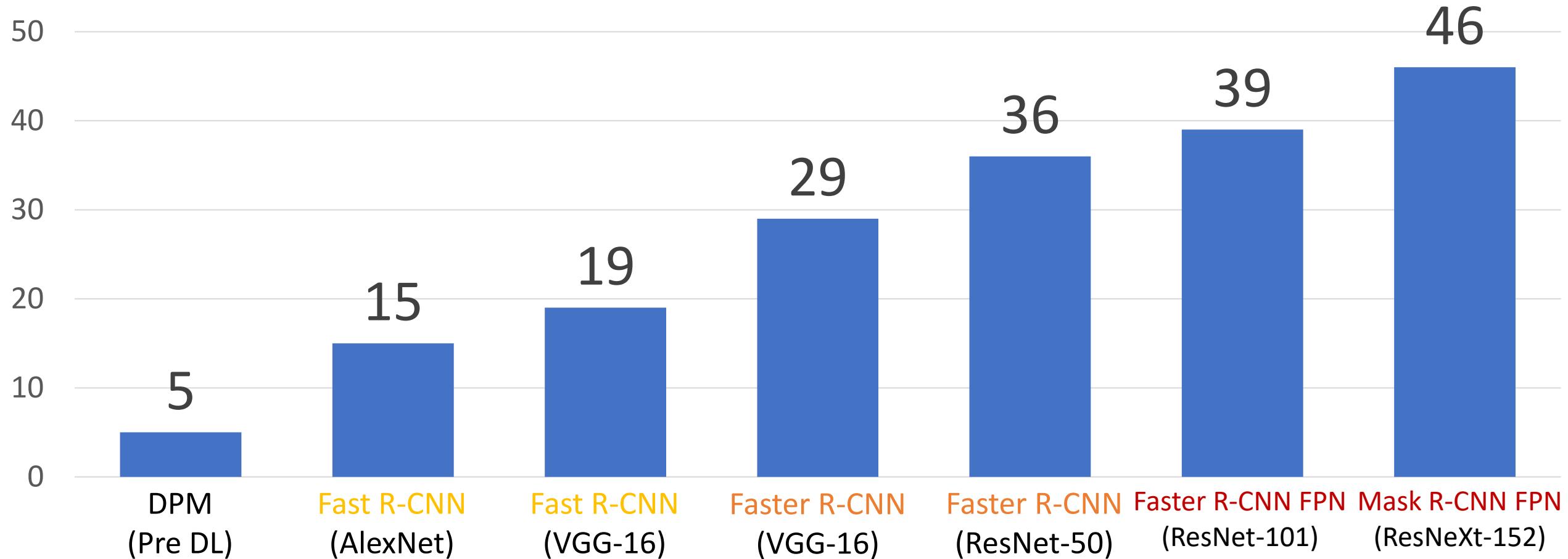
## ImageNet Classification Challenge



Improvements in CNN architectures lead to improvements in many downstream tasks thanks to transfer learning!

# Transfer Learning with CNNs: Architecture Matters!

## Object Detection on COCO



Ross Girshick, "The Generalized R-CNN Framework for Object Detection", ICCV 2017 Tutorial on Instance-Level Visual Recognition

# Self-Supervised Learning (Roughly, ~2019)

# Self-Supervised Representation Learning

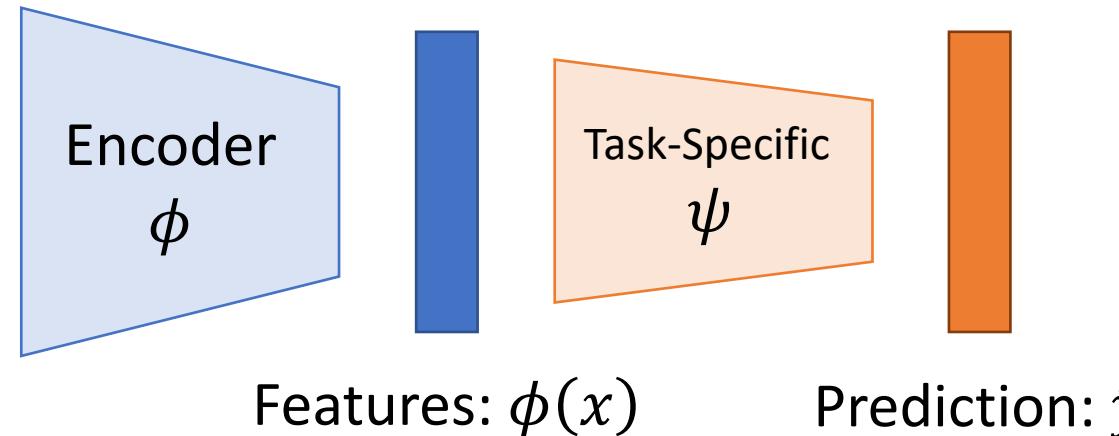
- Building methods that learn representations from “raw” data – no annotations required
- **Unsupervised Learning:** Model isn’t told what to predict. This is somewhat old terminology, not used as much these days.
- **Self-Supervised Learning:** Model is trained to predict some naturally-occurring signal in the raw data rather than human annotations.
- **Semi-Supervised Learning:** Train jointly with some labeled data and (a lot of) unlabeled data.

# Self-Supervised Learning Pipeline

**Stage 1: Pretraining  
(Representation learning)**  
Pretrain a network on a  
pretext task that doesn't  
require supervision;



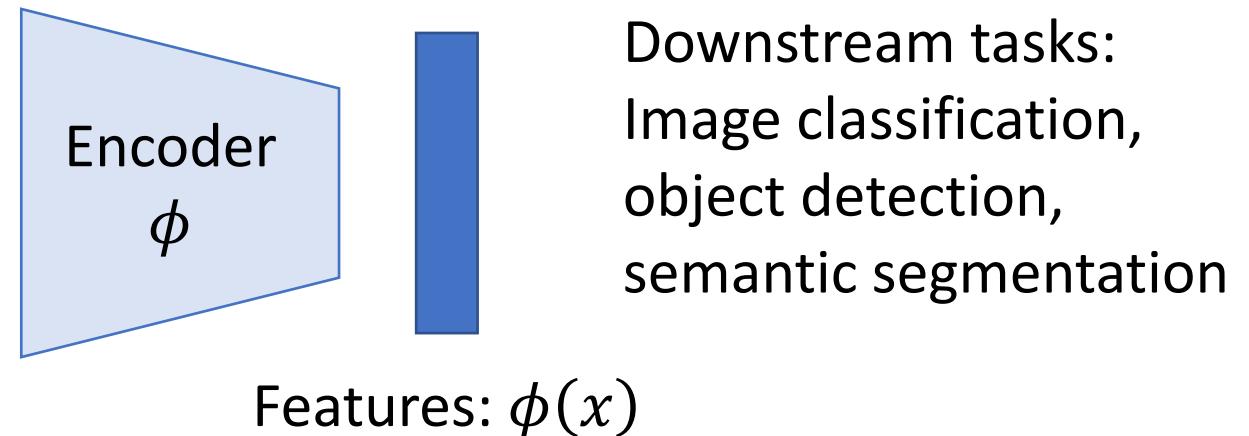
Input Image:  $x$



**Stage 2: Transfer learning**  
Transfer encoder to  
downstream tasks via  
linear classifier,  
K-NN, fine-tuning



Input Image:  $x$

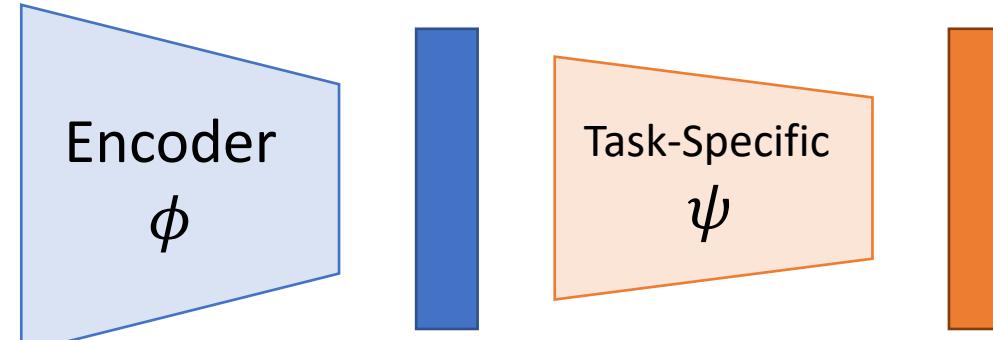


# Self-Supervised Learning Pipeline

**Stage 1: Pretraining  
(Representation learning)**  
Pretrain a network on a  
pretext task that doesn't  
require supervision;



Input Image:  $x$

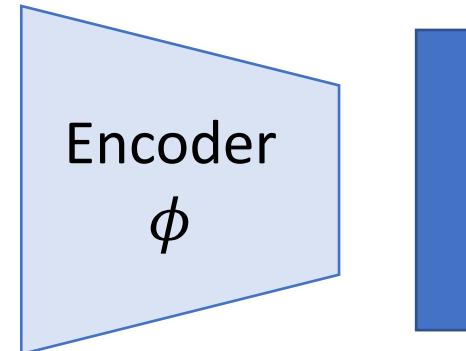


$$\text{Loss: } L(\hat{y}, y)$$

**Stage 2: Transfer learning**  
Transfer encoder to  
downstream tasks via  
linear classifier,  
K-NN, fine-tuning



Input Image:  $x$



Features:  $\phi(x)$

**Goal:** Pretrain + Transfer  
does better than  
supervised pretraining,  
and better than directly  
training on downstream

# Pretext Tasks

**Generative:** Predict a part of the input signal

- Autoencoding
- Autoregression
- GANs
- Inpainting
- Colorization

**Discriminative:** Predict something about the input signal

- Context prediction
- Clustering
- Rotation prediction
- Contrastive learning

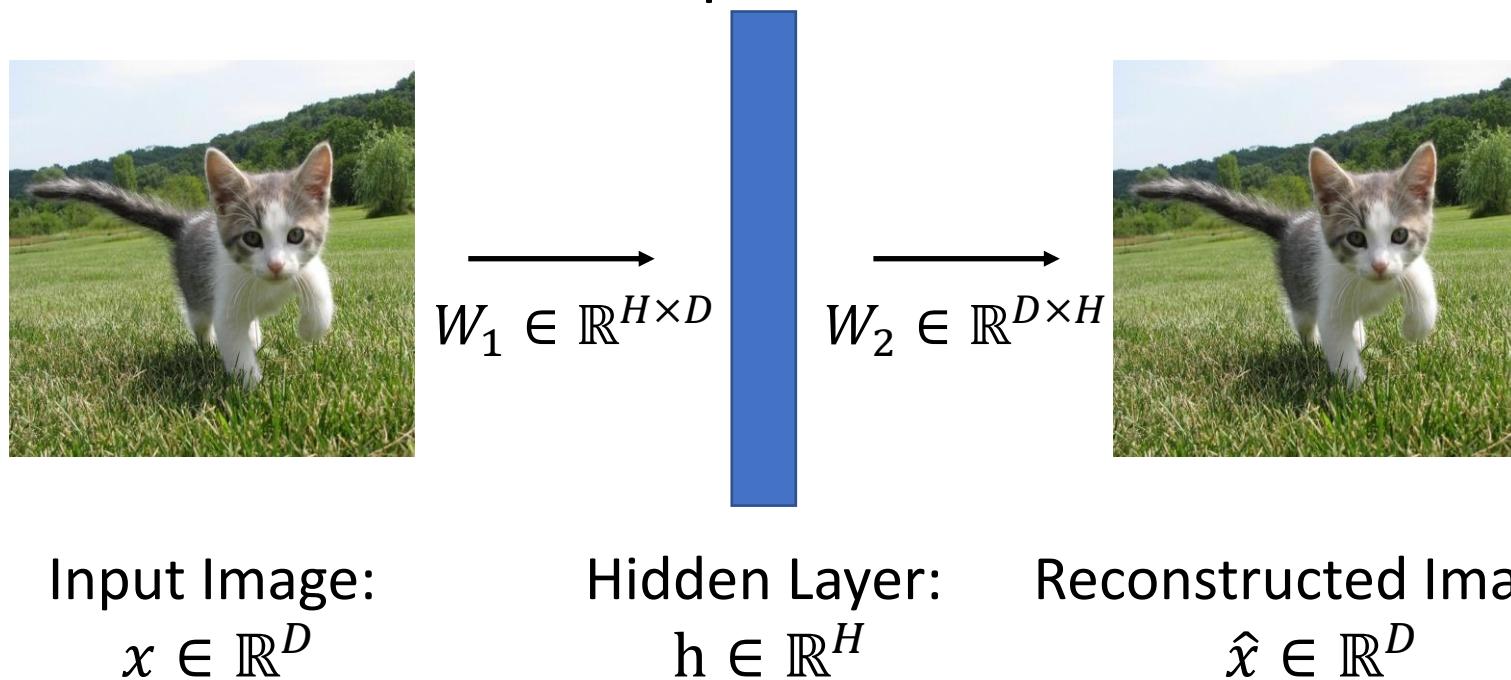
**Multimodal:** Use some additional signal in addition to RGB images

- Video
- 3D
- Sound
- Language

# Pretext Task: Autoencoding

- Autoencoder tries to reconstruct inputs.
- Hidden layer (hopefully) learns good representations.
- $H < D$  to get non-trivial hidden representations

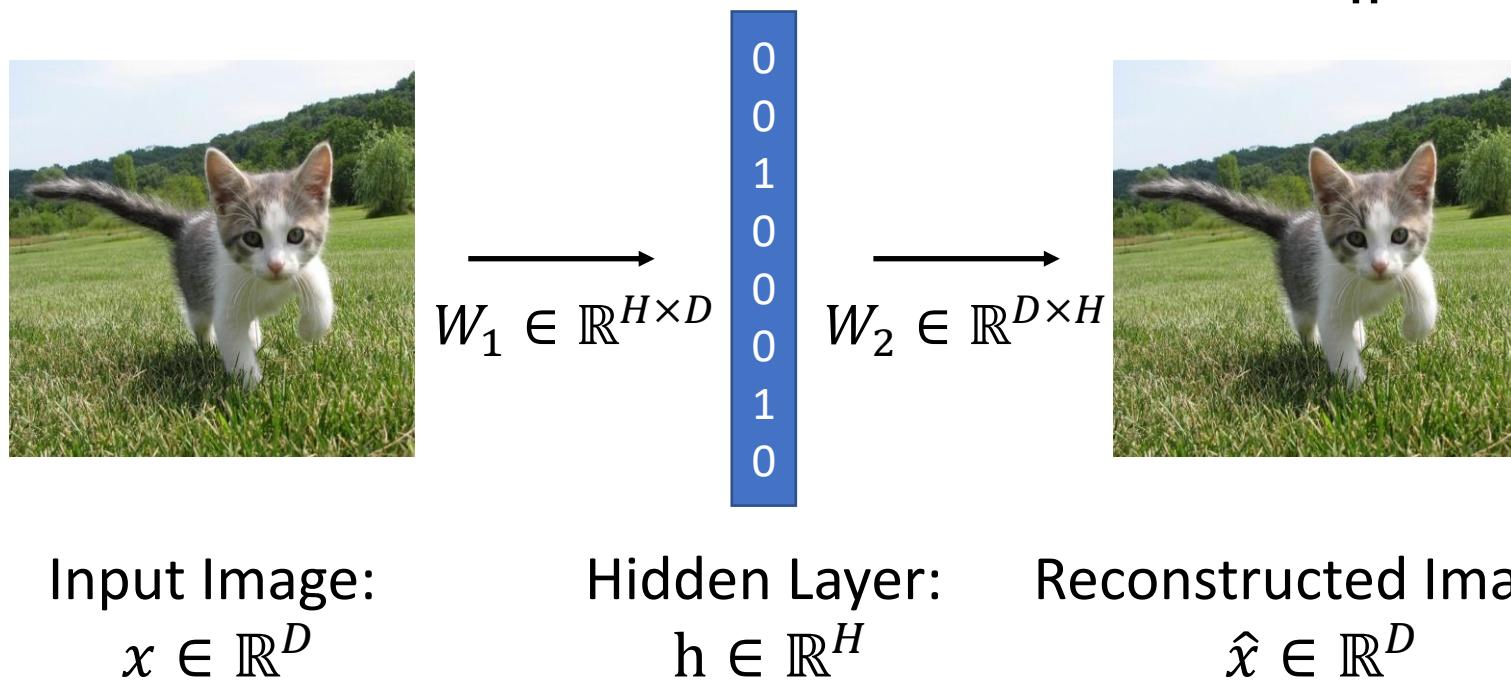
$$\begin{aligned} L(x) &= R(x, \hat{x}) \\ &= \|x - \hat{x}\|_2^2 \end{aligned}$$



Lee et al, "Efficient Sparse Coding Algorithms", NeurIPS 2006; Ranzato et al, "Efficient Learning of Sparse Representations with an Energy-Based Model", NeurIPS 2006;  
Lee et al, "Sparse deep belief net models for visual area V2", NeurIPS 2007; Ng, "Sparse Autoencoder", CS294A Lecture Notes

# Pretext Task: Sparse Autoencoder

- Reconstruct inputs with **sparse activations** (mostly 0)
- Many ways to implement sparsity penalties  $L(x) = R(x, \hat{x}) + \lambda S(h)$   
 $= \|x - \hat{x}\|_2^2 + \lambda \|h\|_1$

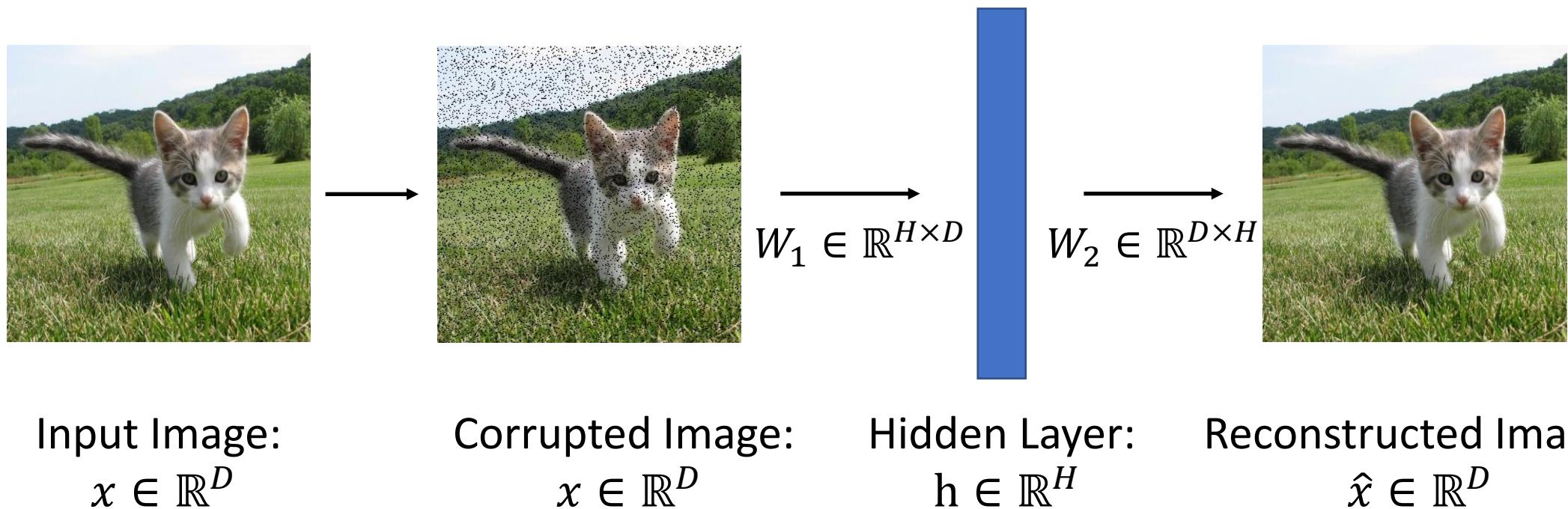


Lee et al, "Efficient Sparse Coding Algorithms", NeurIPS 2006; Ranzato et al, "Efficient Learning of Sparse Representations with an Energy-Based Model", NeurIPS 2006;  
Lee et al, "Sparse deep belief net models for visual area V2", NeurIPS 2007; Ng, "Sparse Autoencoder", CS294A Lecture Notes

# Pretext Task: Denoising Autoencoder

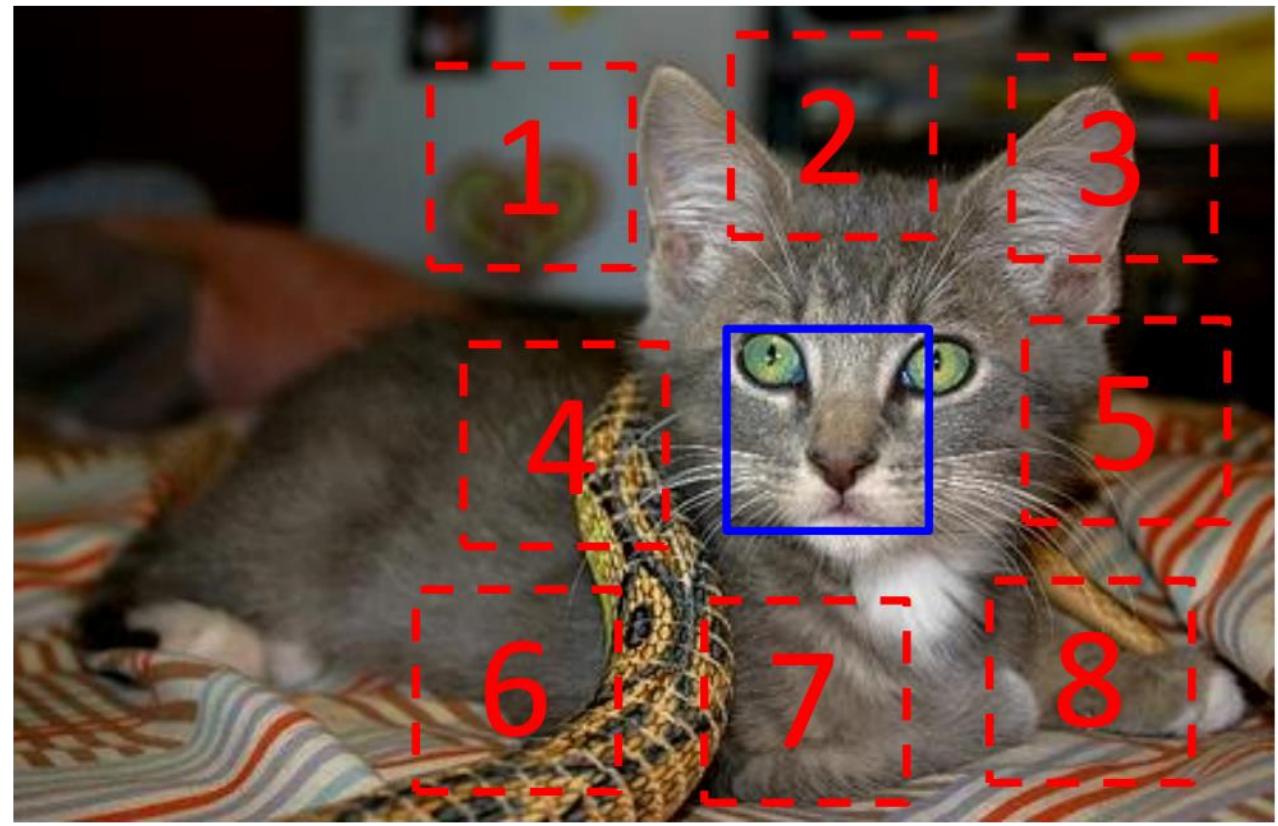
- Reconstruct noisy inputs (pixels randomly set to zero)

$$\begin{aligned} L(x) &= R(x, \hat{x}) \\ &= \|x - \hat{x}\|_2^2 \end{aligned}$$



# Pretext Task: Context Prediction

- Model predicts relative location of two patches from the same image.
- Discriminative pretraining task
- Intuition: Requires understanding objects and their parts



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

Doersch et al, "Unsupervised Visual Representation Learning by Context Prediction", ICCV 2015

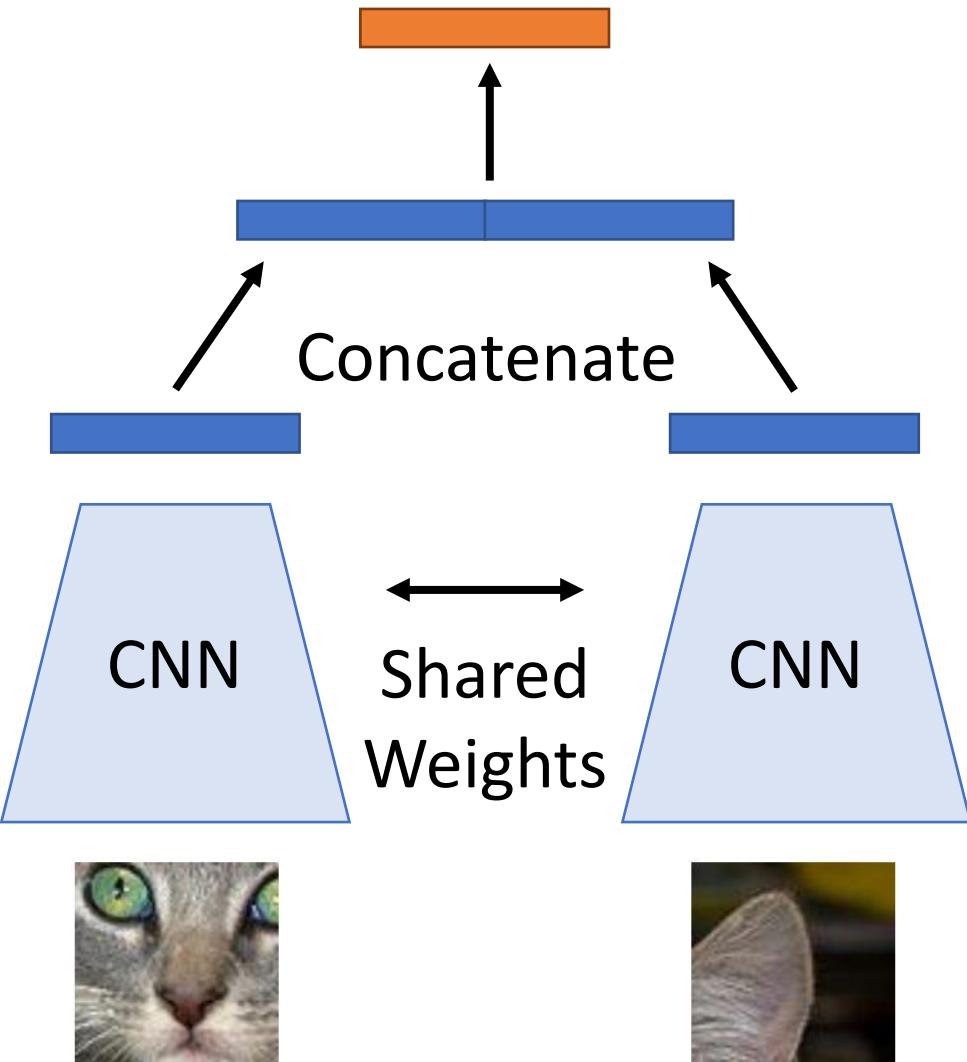
# Pretext Task: Context Prediction

Classification over 8 positions

Model predicts relative location of two patches from the same image.  
Discriminative pretraining task

Intuition: Requires understanding objects and their parts

Two networks with shared weights sometimes called a "Siamese network"



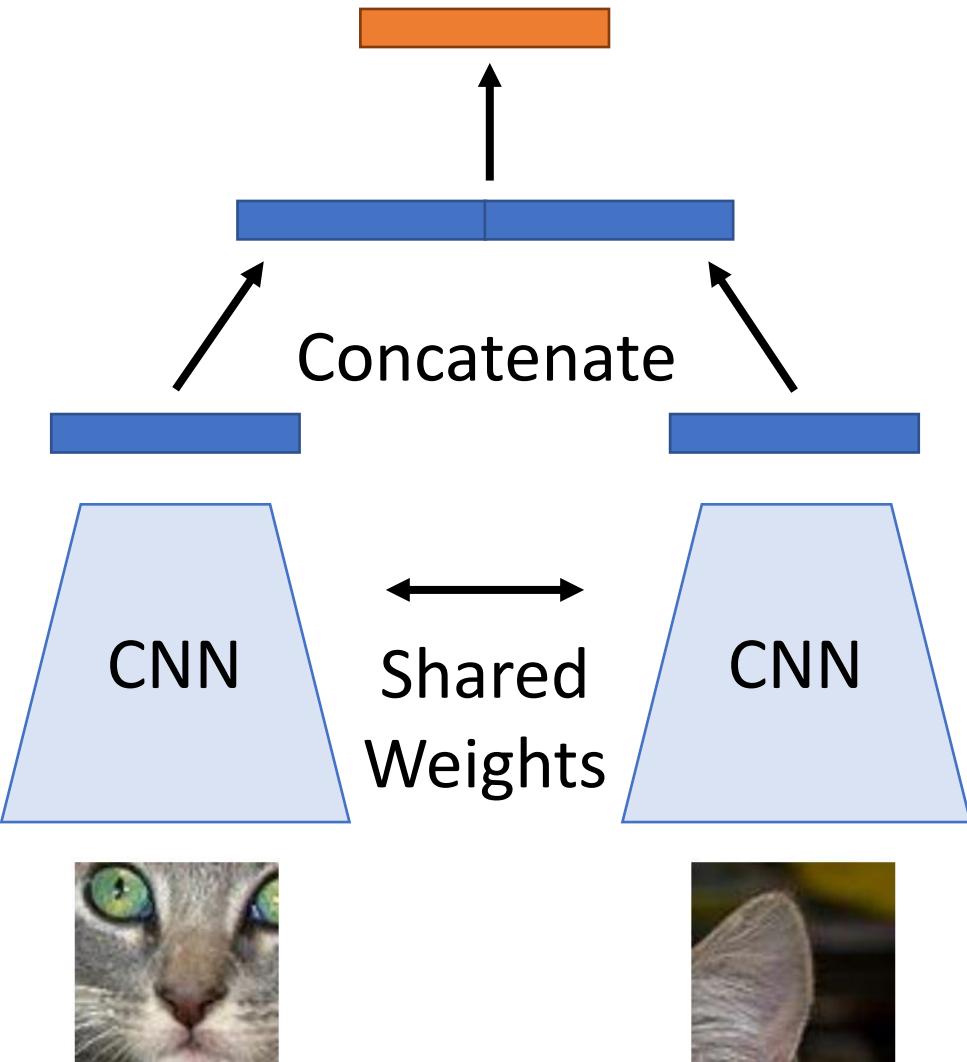
# Pretext Task: Context Prediction

Classification over 8 positions

Model predicts relative location of two patches from the same image.  
Discriminative pretraining task

Intuition: Requires understanding objects and their parts

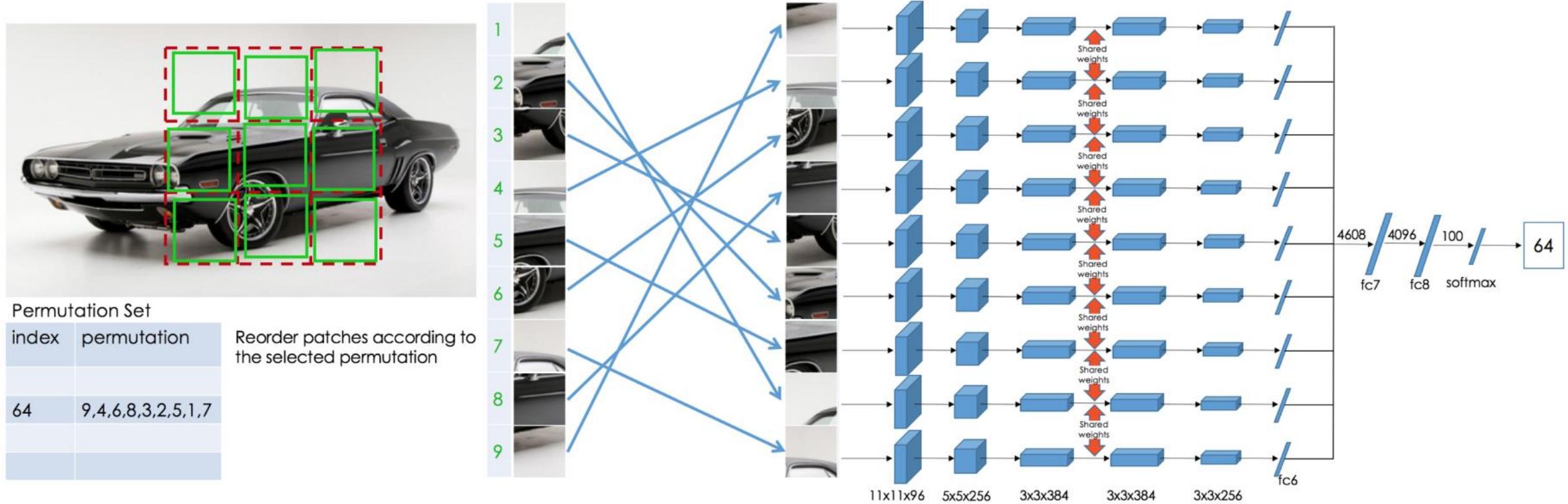
*“For experiments, we use a ConvNet trained on a K40 GPU for approximately four weeks.”*



# Pretext Task: Solving Jigsaw Puzzles

Problem: These methods only work on patches, not whole images!

- Rather than predict relative position of two patches, instead predict permutation to “unscramble” 9 shuffled patches



# Pretext Task: Inpainting

Input Image



Encoder:  
 $\phi$

Decoder:  
 $\psi$

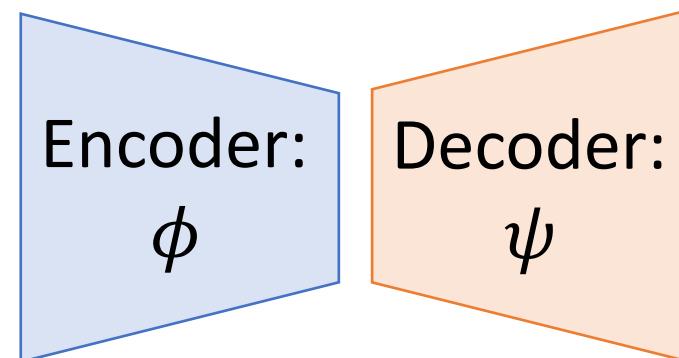
Predict Missing Pixels



Human Artist

# Pretext Task: Inpainting

Input Image



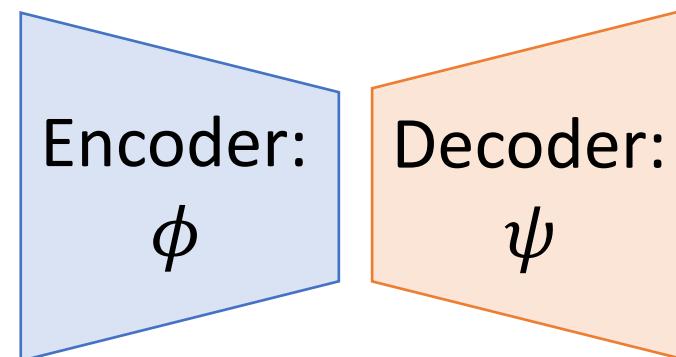
Predict Missing Pixels



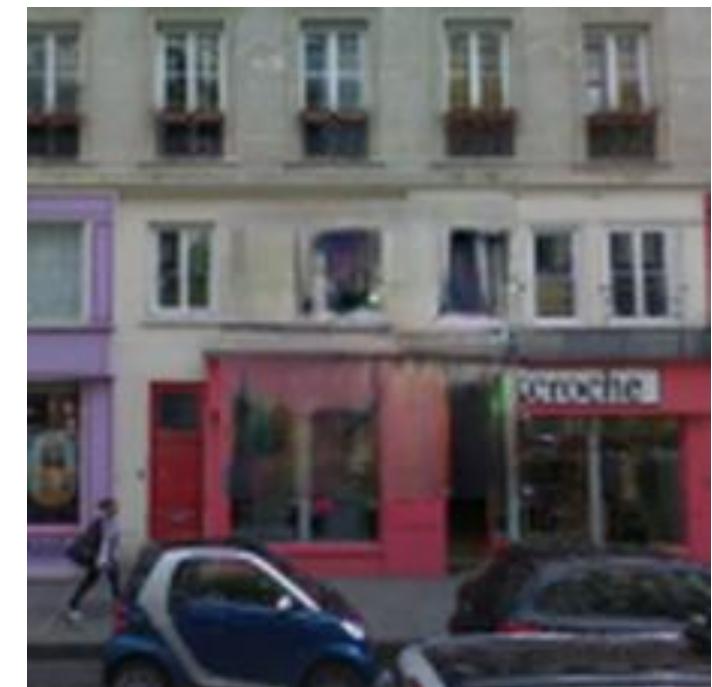
L2 Loss  
(Best for feature learning)

# Pretext Task: Inpainting

Input Image



Predict Missing Pixels



L2 + Adversarial Loss  
(Best for nice images)

# Pretext Task: Colorization

- Colorize grayscale images

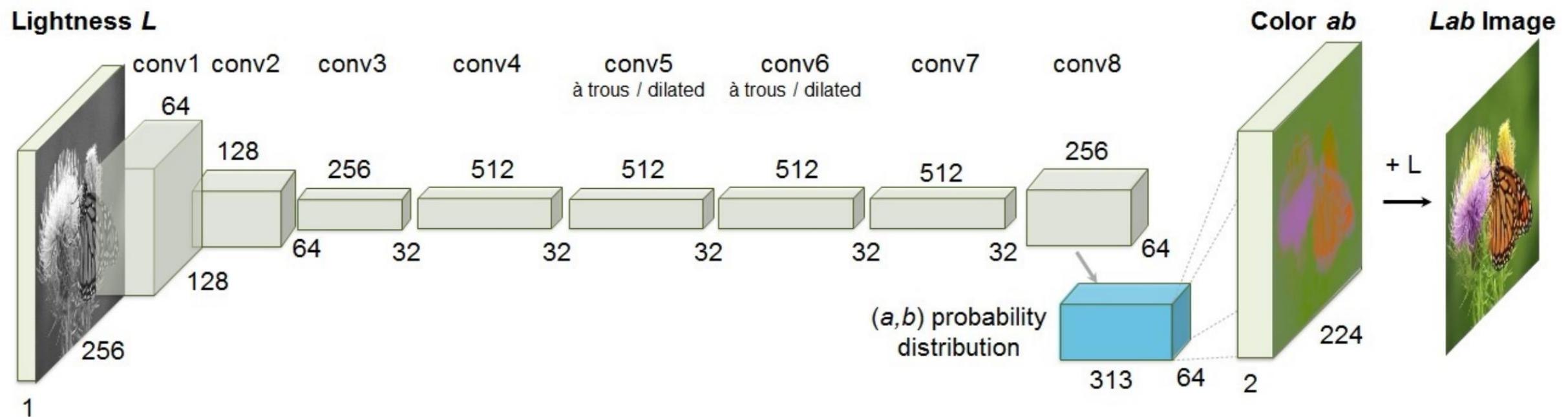


Input: Grayscale Image

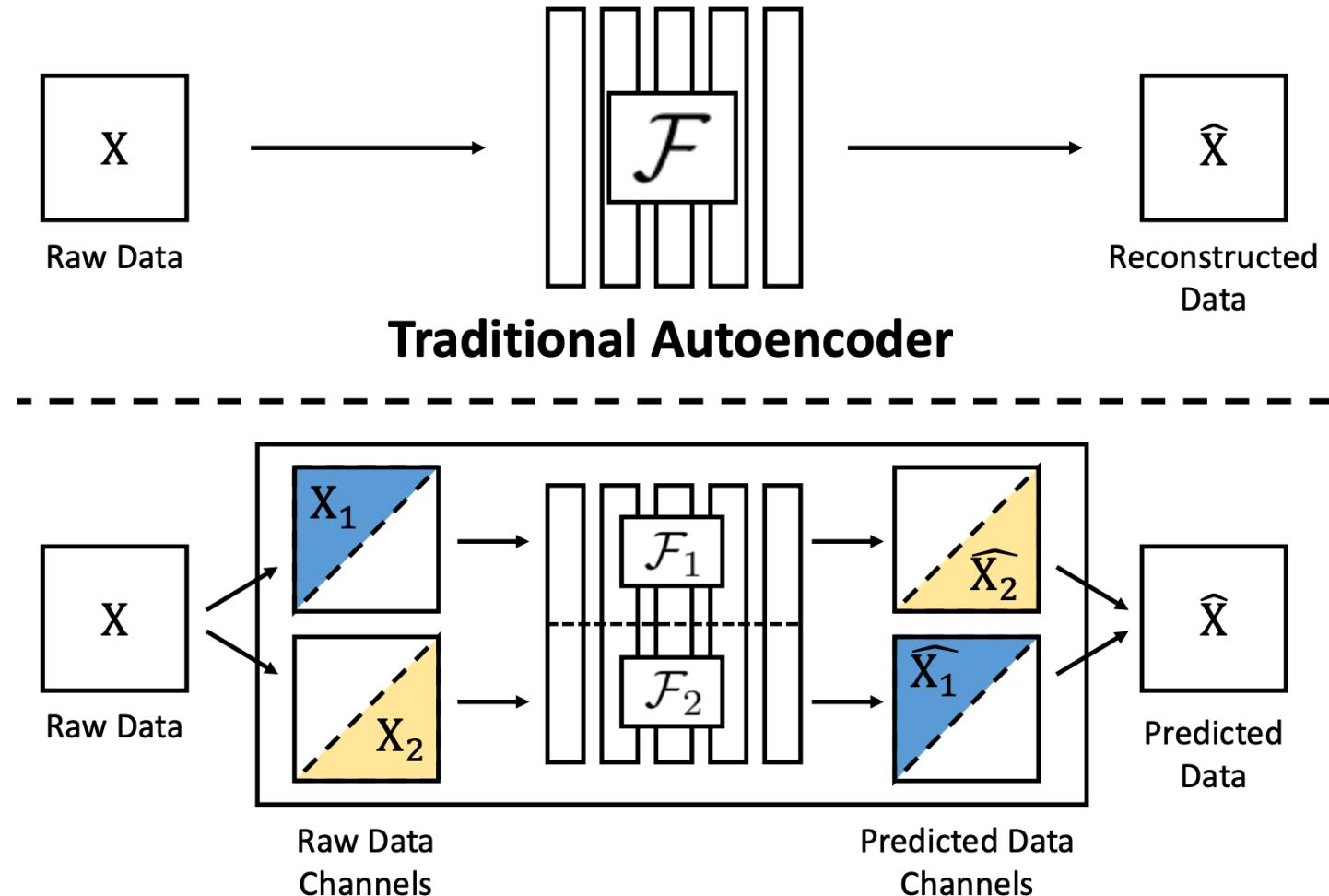


Output: Color Image

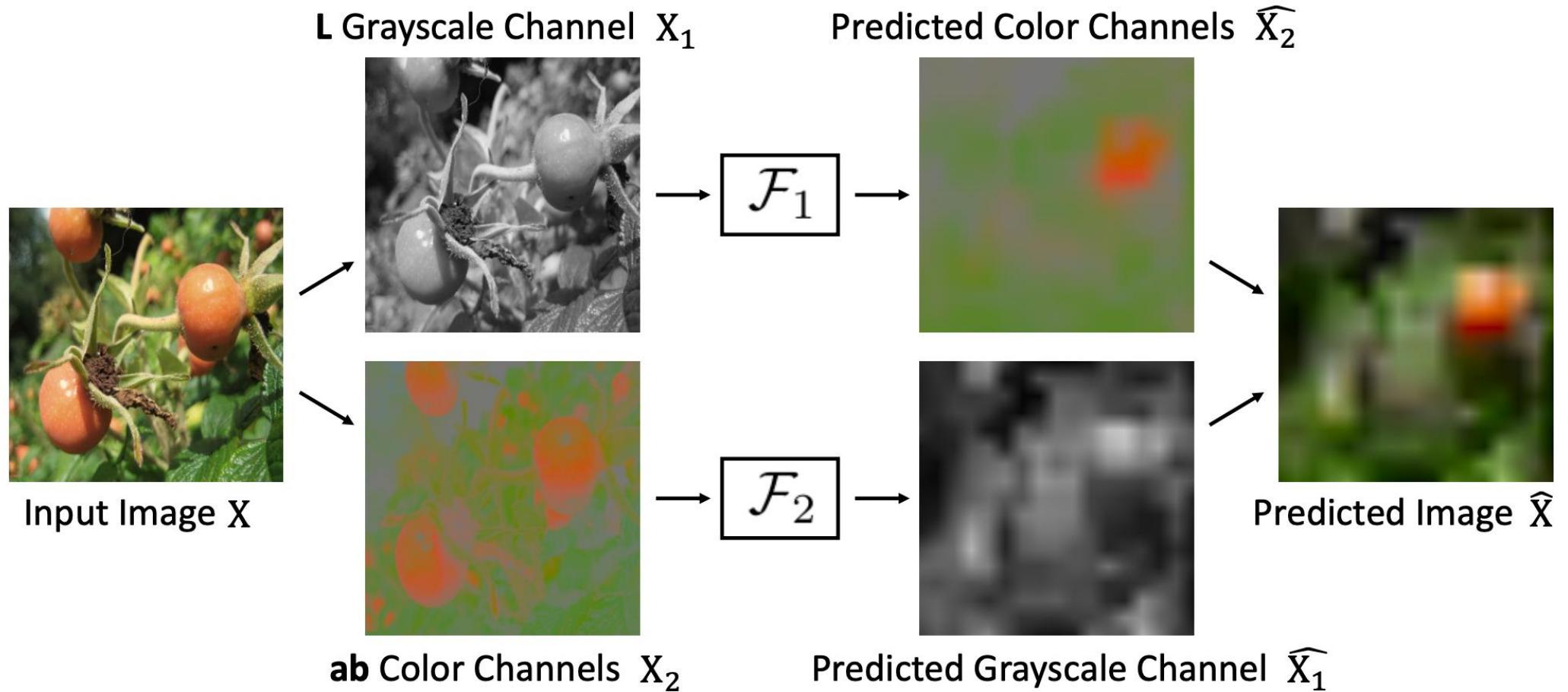
# Pretext Task: Colorization



# Colorization Extension: Split-Brain Autoencoder



# Colorization Extension: Split-Brain Autoencoder



# Colorization Extension: Split-Brain Autoencoder

**Concern:** Generative pretexts encourage spending model capacity on details unimportant for downstream tasks (e.g., regressing exact right shade of orange)



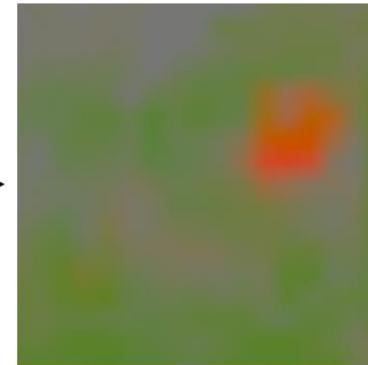
Input Image  $X$

L Grayscale Channel  $X_1$



$$\rightarrow \mathcal{F}_1 \rightarrow$$

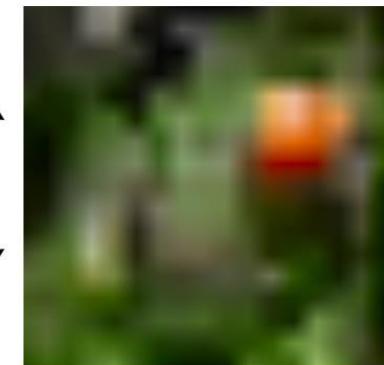
Predicted Color Channels  $\widehat{X}_2$



ab Color Channels  $X_2$

$$\rightarrow \mathcal{F}_2 \rightarrow$$

Predicted Grayscale Channel  $\widehat{X}_1$

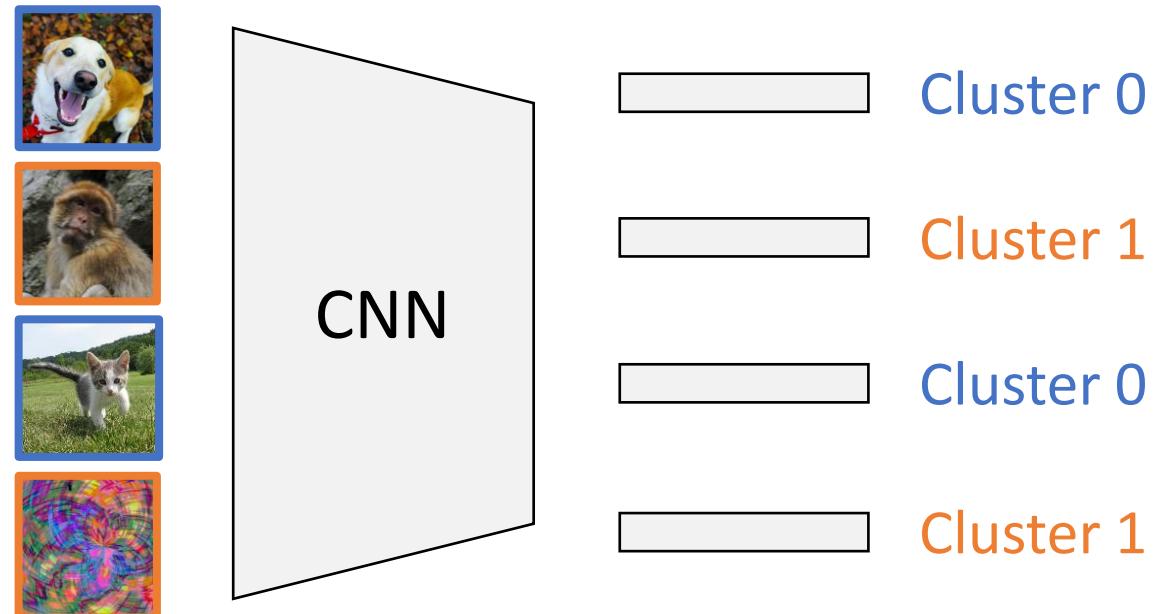


Predicted Image  $\widehat{X}$

**Solution:** Discriminative pretext tasks that require classification, not generation

# Pretext Task: Deep Clustering

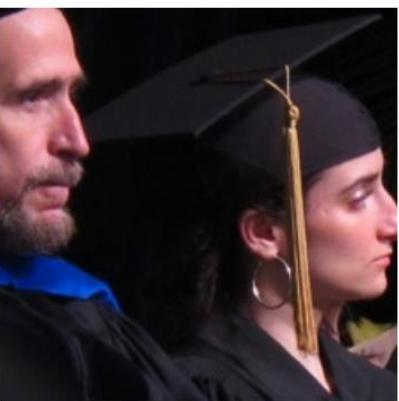
- A naïve example of deep clustering:
  1. Randomly initialize a CNN
  2. Run many images through CNN, get their final-layer features
  3. Cluster the features with K-Means; record cluster assignments
  4. Use cluster assignments as pseudo-labels for each image; train the CNN to predict cluster assignments
  5. GOTO 2.



Caron et al, "Deep Clustering for Unsupervised Learning of Visual Features", ECCV 2018  
Caron et al, "Unsupervised Pre-Training of Image Features on Non-Curated Data", ICCV 2019  
Yan et al, "ClusterFit: Improving Generalization of Visual Representations", CVPR 2020  
Caron et al, "Unsupervised Learning of Visual Features by Contrasting Cluster Assignments", NeurIPS 2020

# Pretext Task: Rotation Prediction

- 4-way classification task: How much was each image rotated?  
(0, 90, 180, or 270 degrees)



90

270

180

0

270

Gidaris et al, "Unsupervised representation learning by predicting image rotations", ICLR 2018

# Which SSL Method is best?

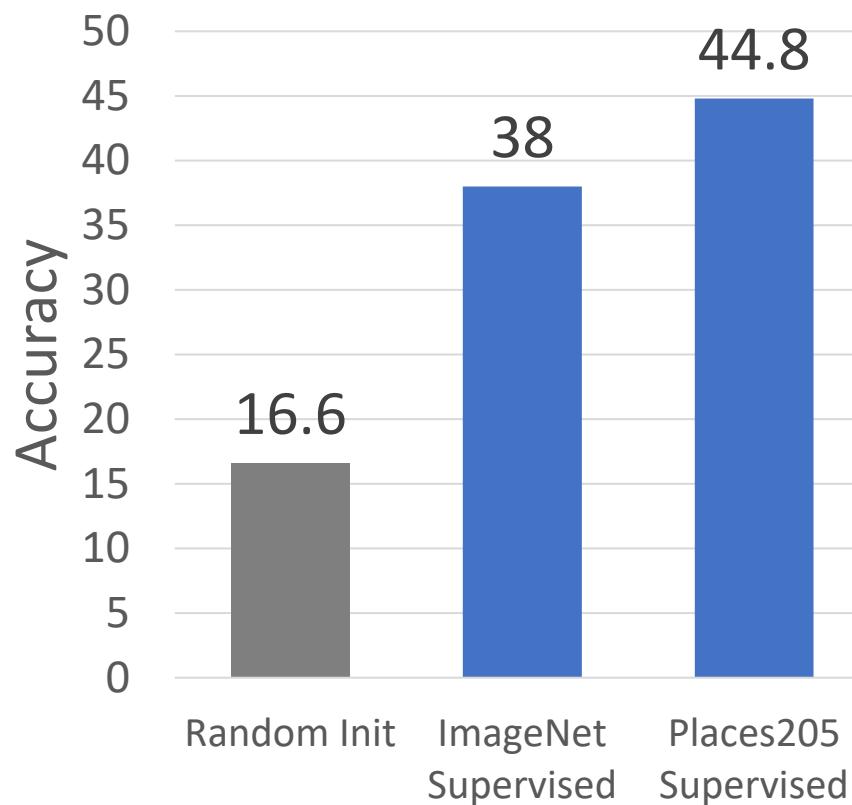
Fair evaluation of SSL methods is very hard! No theory, so we need to rely on experiment

Many choices in experimental setup, huge variations from paper to paper:

- CNN architecture? AlexNet, ResNet50, something else?
- Pretraining dataset? ImageNet, or something else?
- Downstream task? ImageNet classification, detection, something else?
- Pretraining hyperparameters? Learning rates, training iterations, data augmentation?
- Transfer learning protocol?
  - Linear probe? From which layer? How to train linear models? SGD, something else?
  - Transfer learning hyperparameters? Data augmentation or BatchNorm during transfer learning?
  - Fine-tune? From which layer? Architecture of “head” you attach? Linear or nonlinear? Fine-tuning hyperparameters?
  - K-NN? What value of K? Normalization on features?

# Supervised vs. Self-Supervised Pretraining

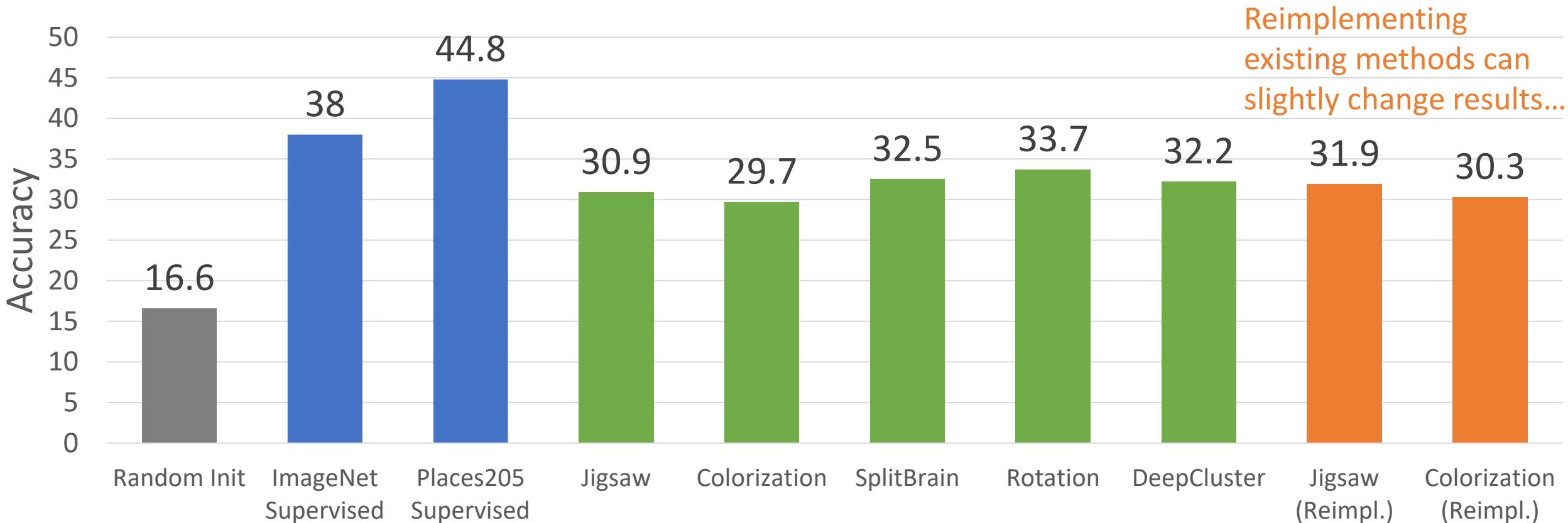
Places205 Linear Classification from AlexNet conv5



# Supervised vs. Self-Supervised Pretraining

- As of 2019, SSL gave worse features than supervised pretraining

Places205 Linear Classification from AlexNet conv5

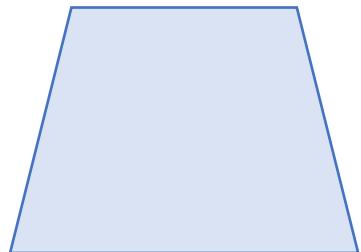


# Self-Supervised Learning for Natural Language

Computer Vision

Image Features:

$H \times W \times C$

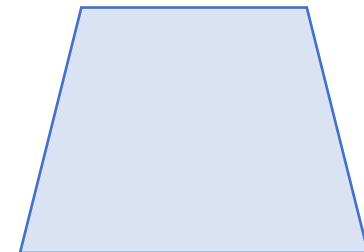


Input Image

Natural Language Processing

Word Features

$L \times C$



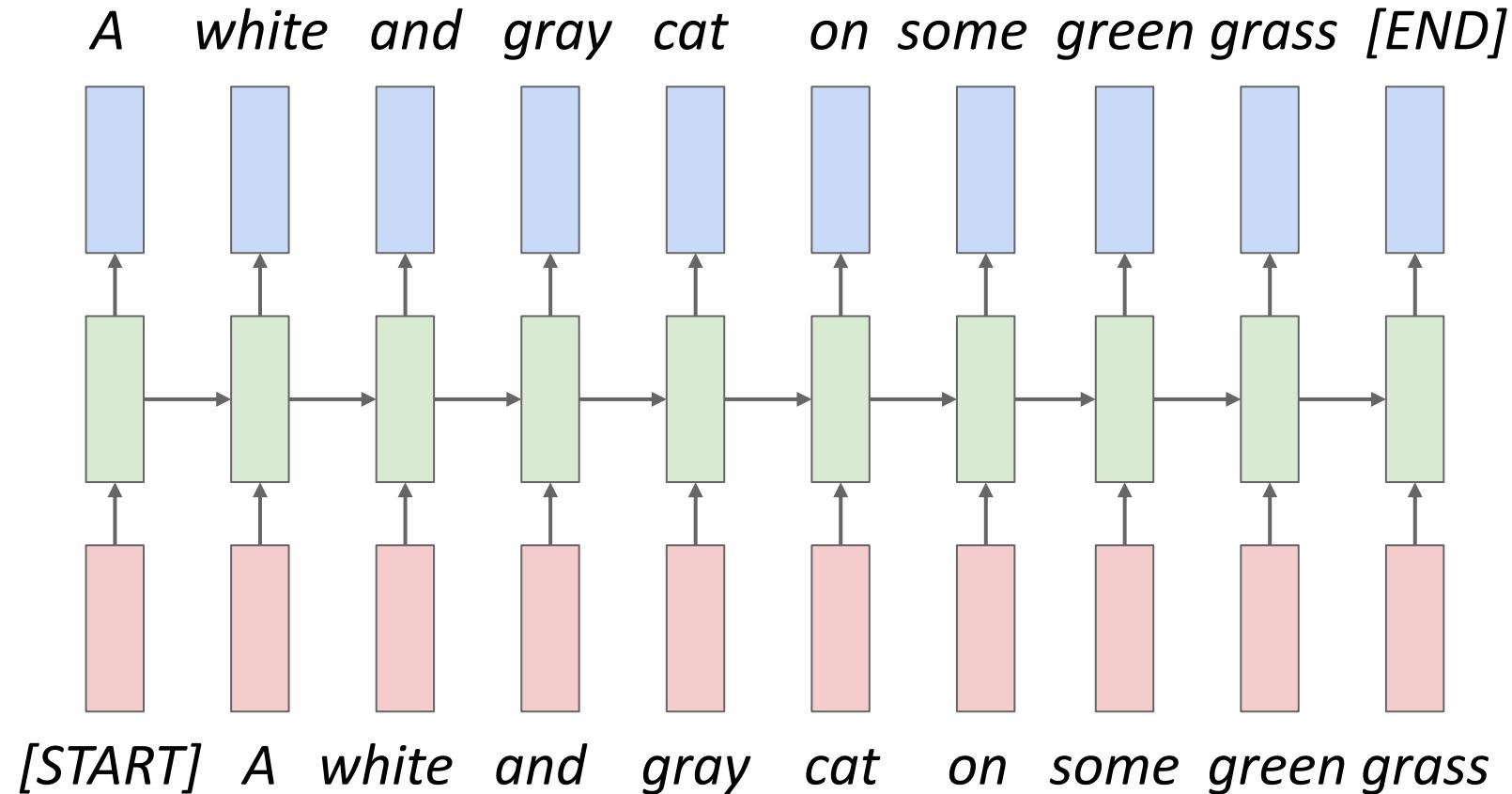
*A white and gray  
cat standing outside  
on the grass*

Input Sentence (L words)

# Self-Supervised Learning for Natural Language

RNN language models train on raw text – no human labels required!

Their hidden states give features that transfer to many downstream tasks!

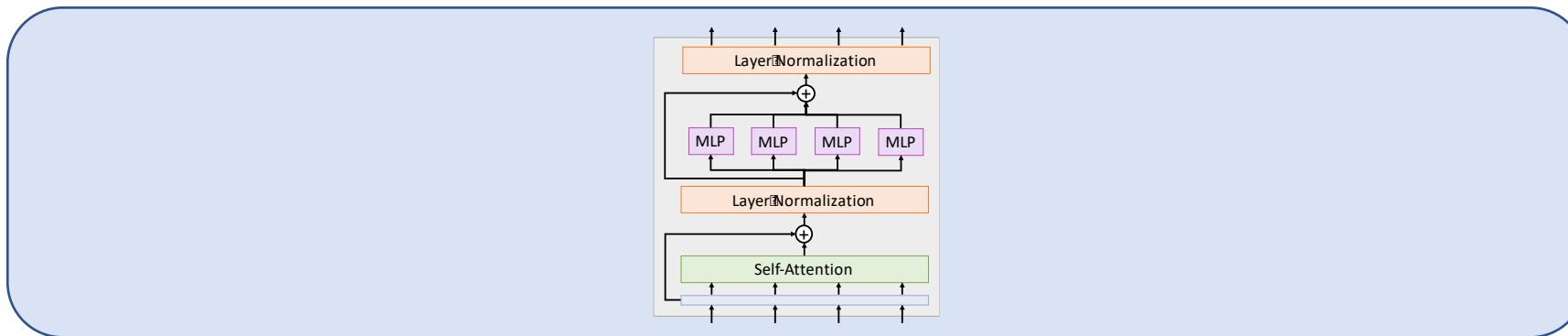


# Self-Supervised Learning for Natural Language

Transformer-based language models work even better! Can scale up to very large datasets, and give extremely powerful features that transfer to downstream tasks

Wildly successful: larger models, larger datasets give better features that improve performance on many downstream NLP tasks. The dream of SSL made real!

*A white and gray cat on some green grass [END]*



*[START] A white and gray cat on some green grass*

Radford et al, "Language models are unsupervised multitask learners", 2019

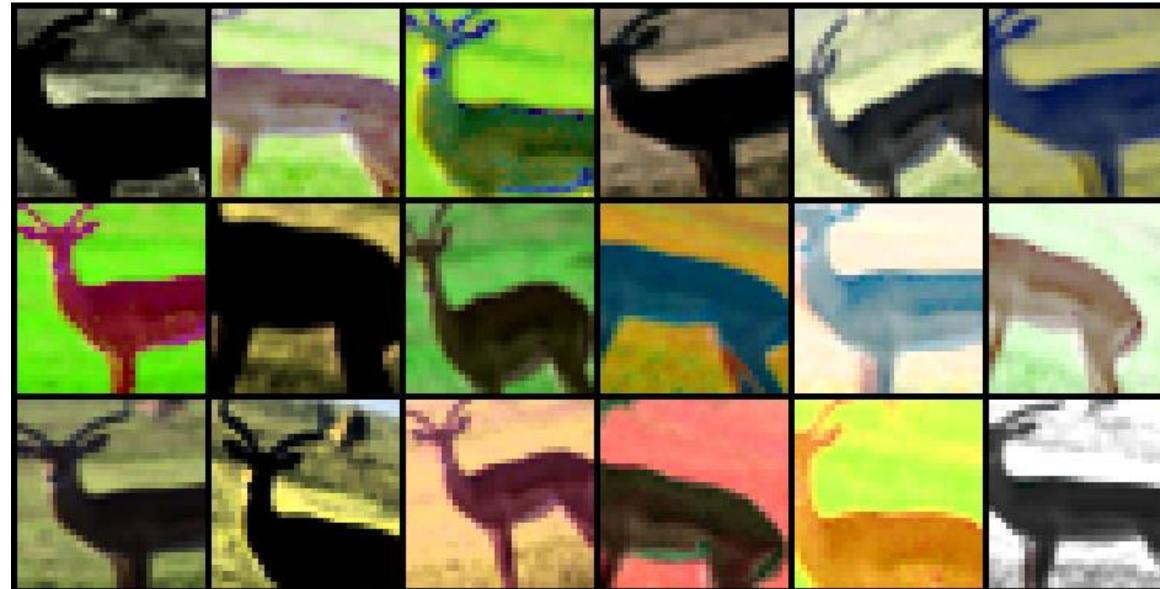
Brown et al, "Language Models are Few-Shot Learners", arXiv 2020

Rae et al, "Scaling Language Models: Methods, Analysis, & Insights from Training Gopher", arXiv 2021

# Self-Supervised Learning (Roughly, 2020 ~)

# Pretext Task: Exemplar CNN

Quiz: What is this?



Different data augmentations (scale, shift, color jitter) of the same initial image patch

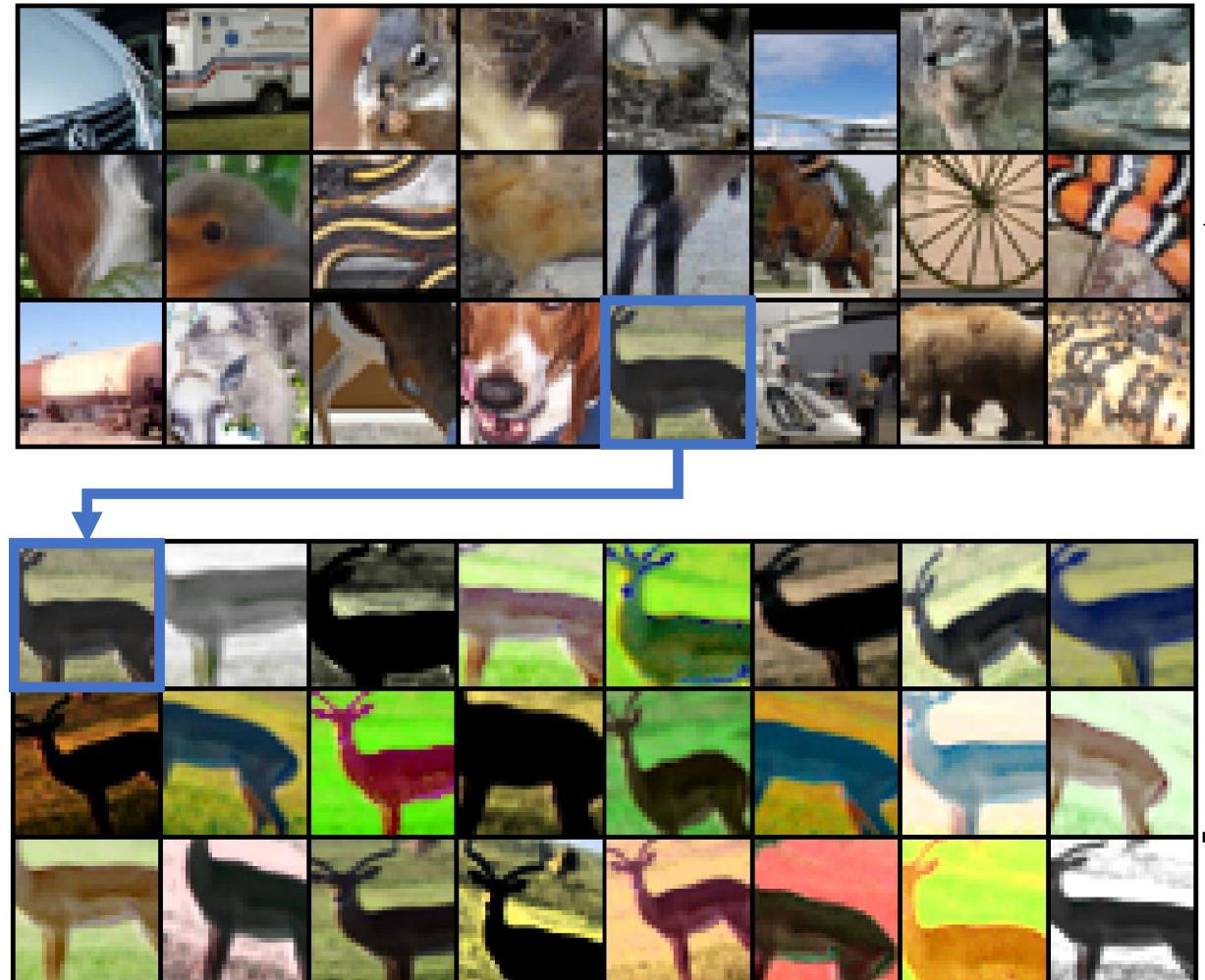
Answer: Deer!

# Pretext Task: Exemplar CNN

Given an initial dataset of  $N$  images

**Problem:** number of parameters in the final layer depends on  $N$ ; hard to scale

Sample  $K$  different augmentations for each; now have  $K \cdot N$  total patches



Predicts which of the  $N$  original images it came from (N-way classification)

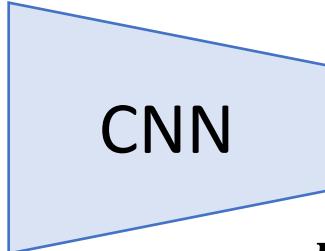
CNN

CNN inputs an augmented patch

# Pretext Task: Contrastive Learning

- Assume we don't have labels for images, but we know whether some pairs of images are **similar** or **dissimilar**
- Let  $d = \|\phi(x_1) - \phi(x_2)\|_2$  be the Euclidean distance btw features
- **Problem: where to get positive and negative pairs?**

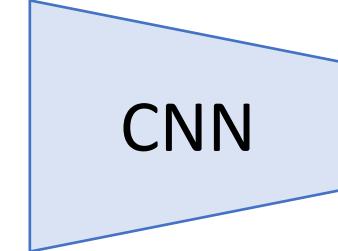
**Similar** images should have similar features



$$L_S(x_1, x_2) = d^2$$

Pull features together

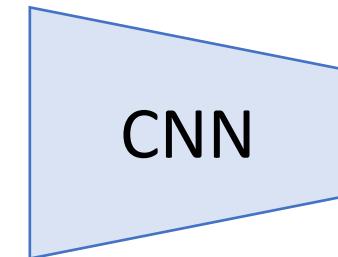
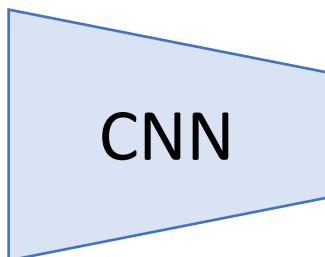
**Dissimilar** images should have dissimilar features



$$L_D(x_1, x_2)$$

$$= \max(0, m - d)^2$$

Push features apart  
(up to margin  $m$ )



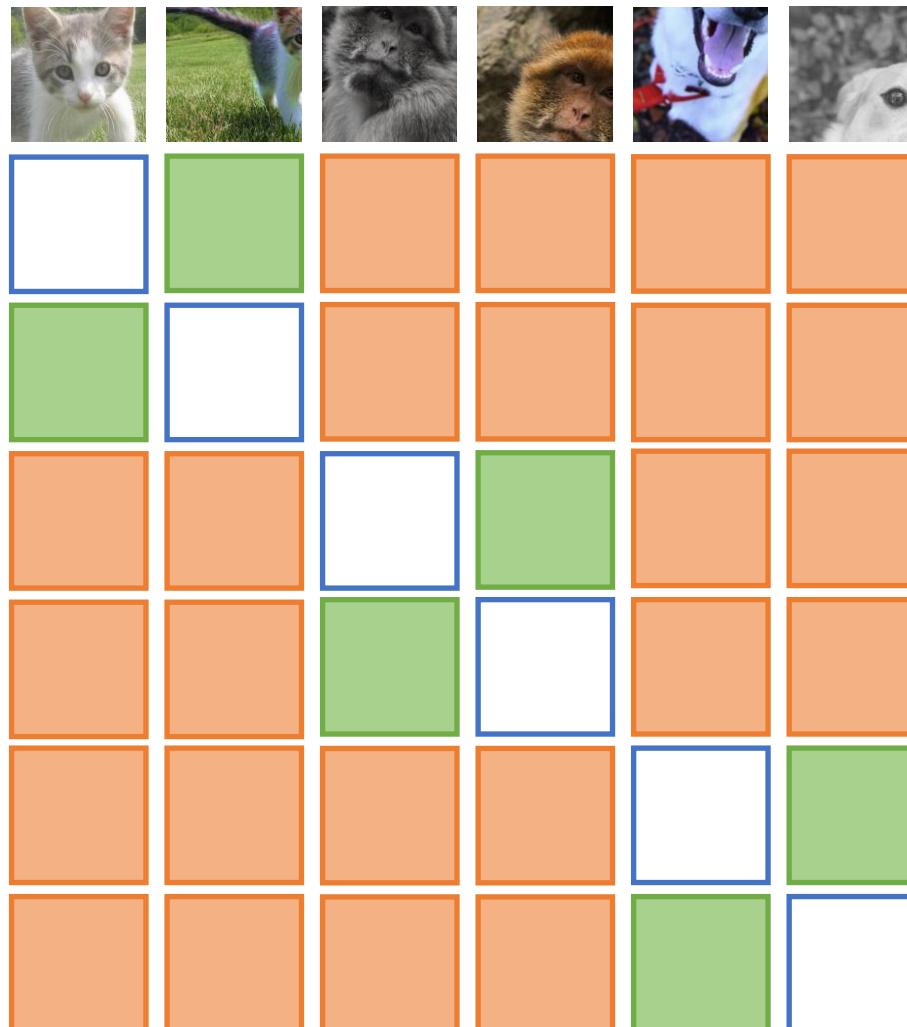
Hadsell et al, "Dimensionality Reduction by Learning and Invariant Mapping", CVPR 2006

# Pretext Task: Contrastive Learning

Batch of N images    Two augmentations for each image



Extract features



- Each image tries to predict which of the other  $2N-1$  images came from the same original image.
- Similarity between  $x_i$  and  $x_j$ :  

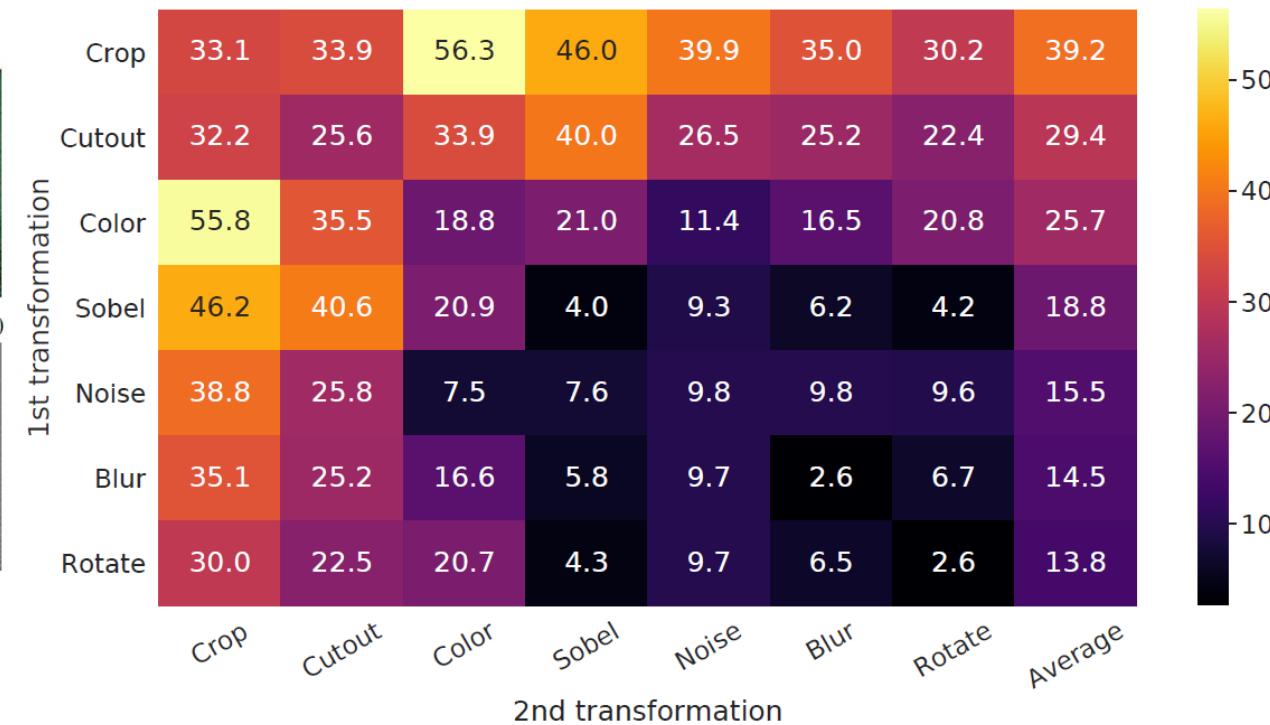
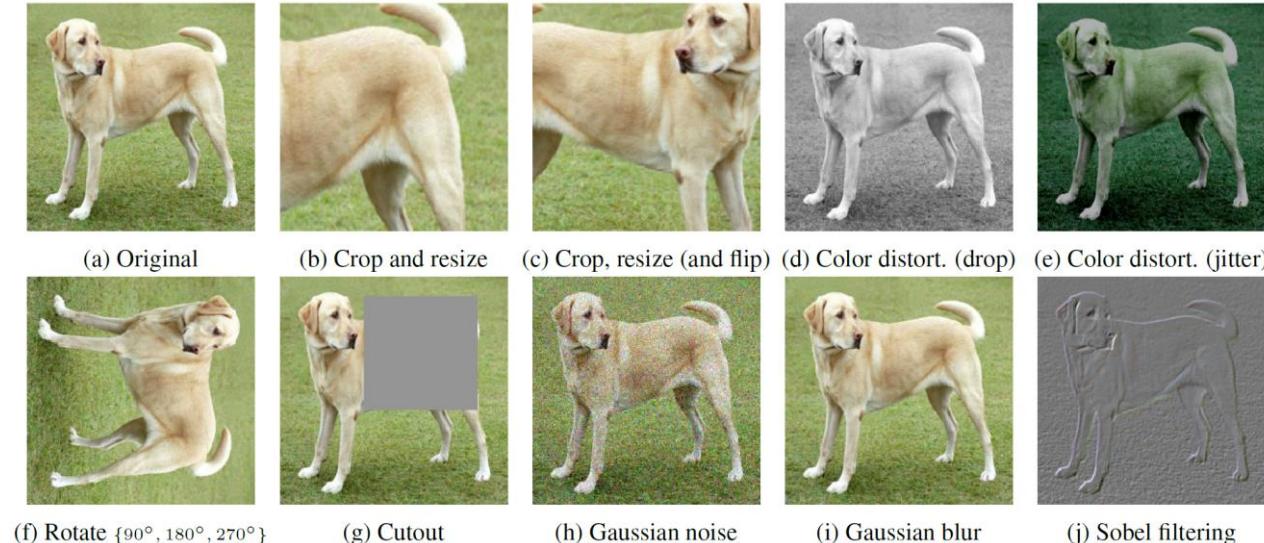
$$s_{i,j} = \frac{\phi(x_i)^T \phi(x_j)}{\|\phi(x_i)\| \cdot \|\phi(x_j)\|}$$
- If  $(x_i, x_j)$  is a positive pair, then loss for  $x_i$  is:  

$$L_i = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1, k \neq i}^{2N} \exp(s_{i,k}/\tau)}$$

( $\tau$  is a temperature)
- Interpretation: Cross-entropy loss over the other  $2N-1$  elements in the batch.

# Impact of Data Augmentation

- SimCLR: Composition of data augmentation operations is crucial for learning good representations.
  - Random resized crop, color jitter, Gaussian blurring, ...



Chen et al, "A Simple Framework for Contrastive Learning of Visual Representations", ICML 2020

# Contrastive vs. Feature Reconstruction

- SimCLR

$$L_i = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{\substack{k=1 \\ k \neq i}}^{2N} \exp(s_{i,k}/\tau)}$$

- Contrastive learning
- $2N-2$  negatives

- MoCo

$$L_i = -\log \frac{\exp(s_{i,j}/\tau)}{\exp(s_{i,j}/\tau) + \sum_{k=1}^K \exp(s_{i,k}/\tau)}$$

- Contrastive learning
- Momentum encoder
- Negatives from queue (size K)

- BYOL

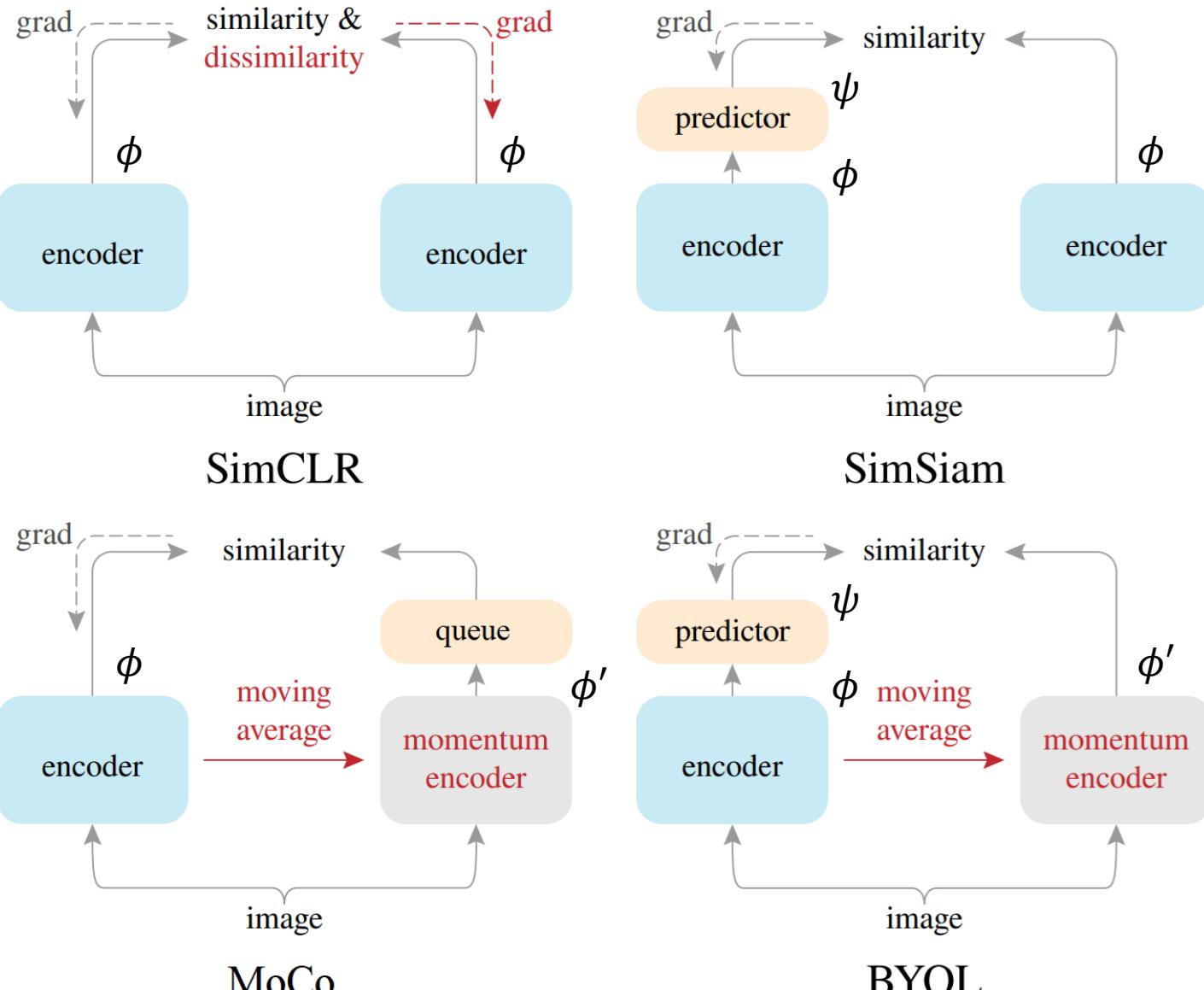
$$L_i = -\frac{\psi(\phi(x_i))^T \phi'(x_j)}{\|\psi(\phi(x_i))\| \cdot \|\phi'(x_j)\|}$$

- Feature reconstruction
- Momentum encoder

- SimSiam

$$L_i = -\frac{\psi(\phi(x_i))^T \phi(x_j)}{\|\psi(\phi(x_i))\| \cdot \|\phi(x_j)\|}$$

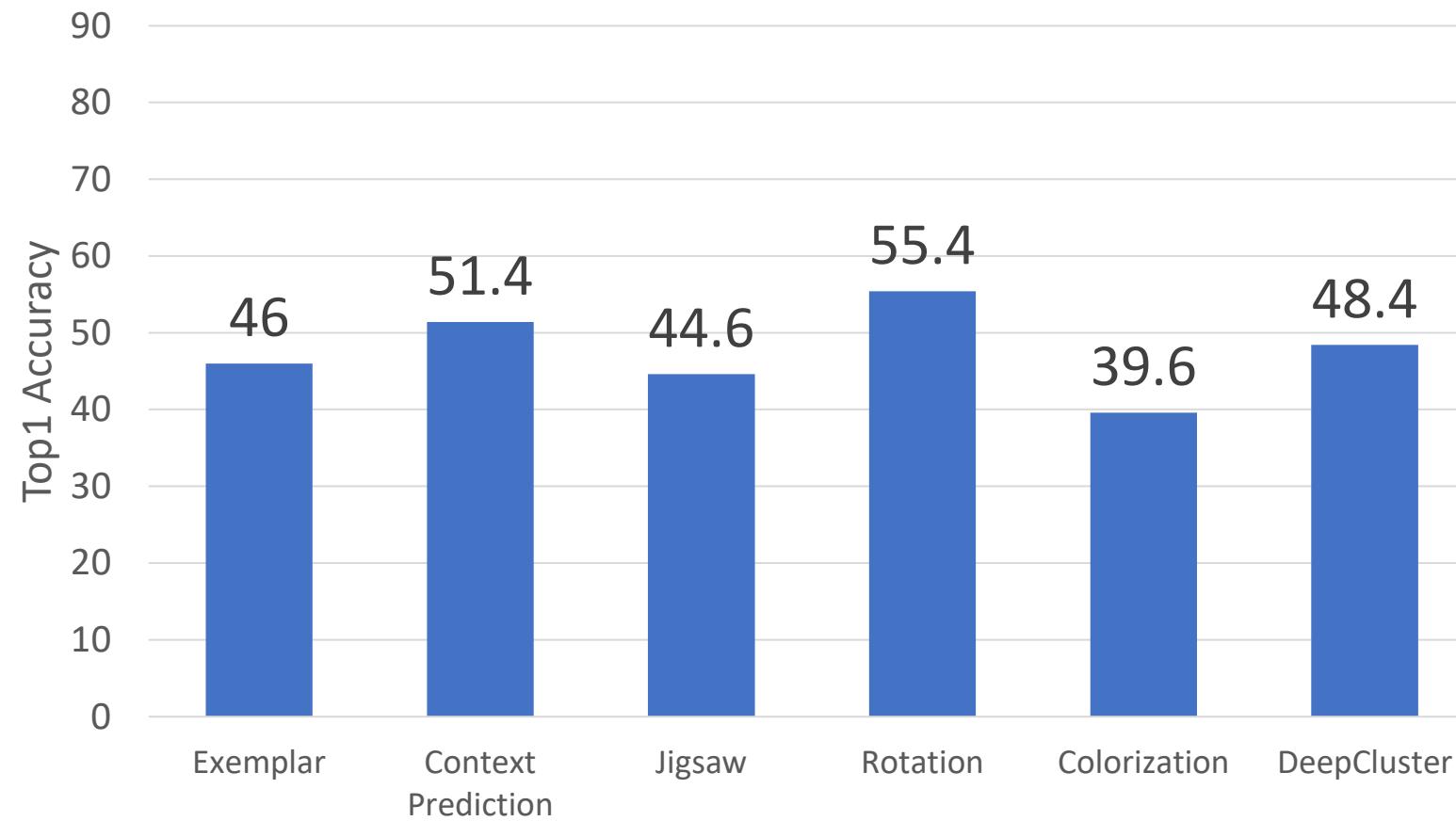
- Feature reconstruction



He et al, "Momentum Contrast for Unsupervised Visual Representation Learning", CVPR 2020  
 Chen et al, "A Simple Framework for Contrastive Learning of Visual Representations", ICML 2020  
 Grill et al, "Bootstrap Your Own Latent: A New Approach to Self-Supervised Learning", NeurIPS 2020  
 Chen and He, "Exploring simple Siamese representation learning", CVPR 2021

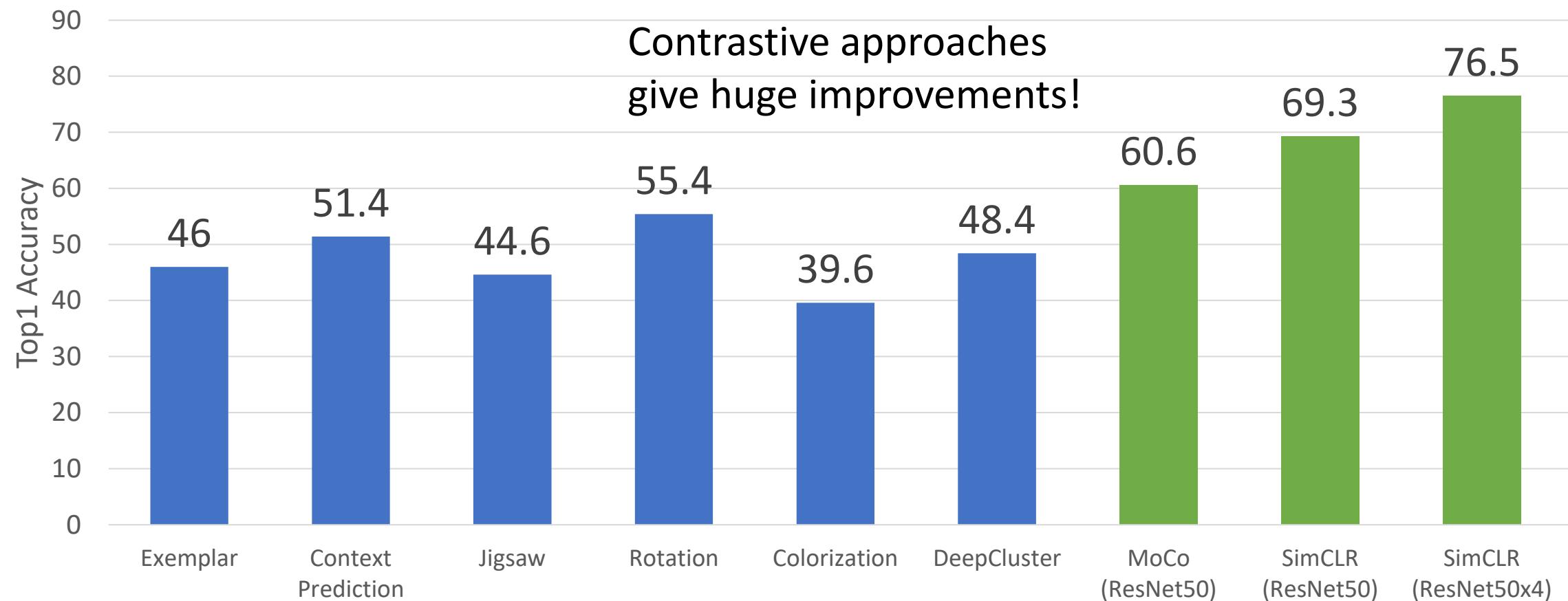
# Comparison of SSL Methods

ImageNet Linear Classification from SSL Features



# Comparison of SSL Methods

## ImageNet Linear Classification from SSL Features



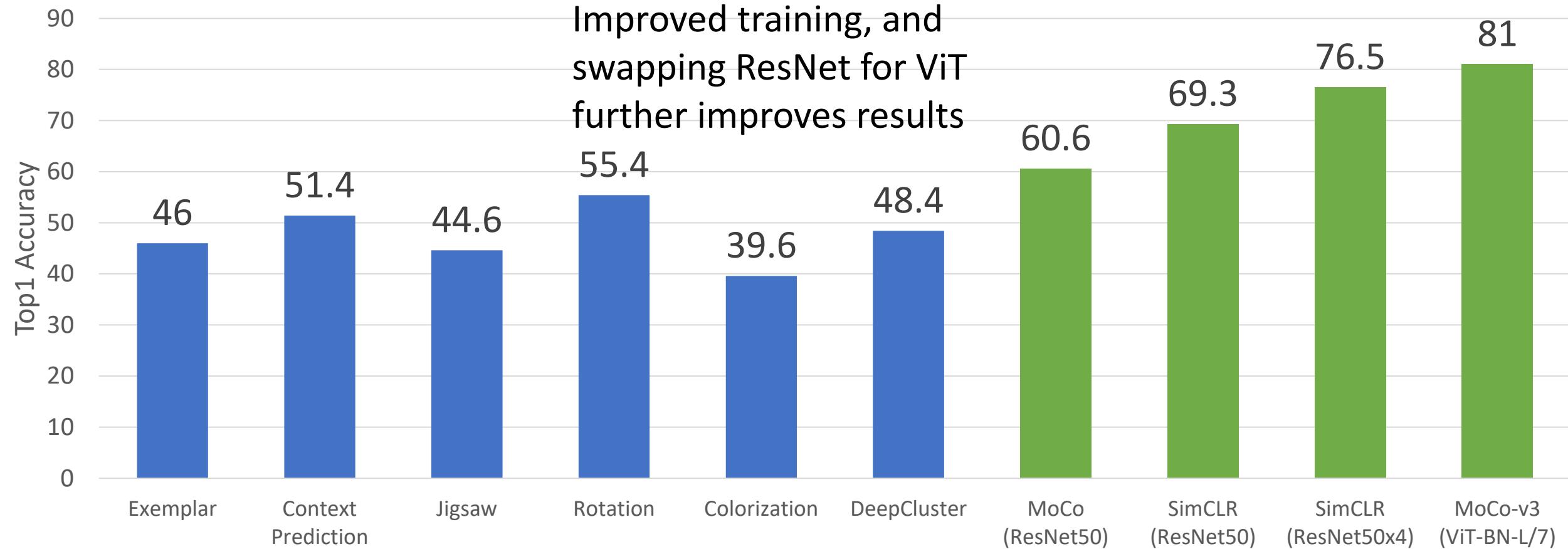
He et al, "Momentum Contrast for Unsupervised Visual Representation Learning", CVPR 2020

Chen et al, "A Simple Framework for Contrastive Learning of Visual Representations", ICML 2020

Chen et al, "An Empirical Study of Training Self-Supervised Vision Transformers", ICCV 2021

# Comparison of SSL Methods

ImageNet Linear Classification from SSL Features



He et al, "Momentum Contrast for Unsupervised Visual Representation Learning", CVPR 2020

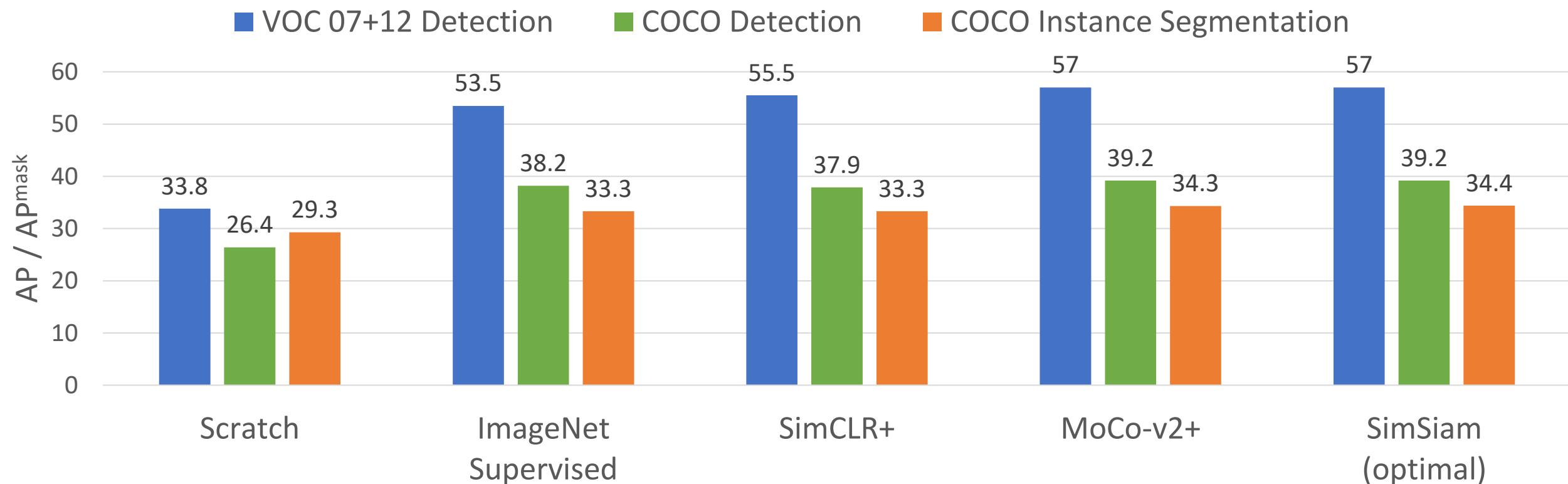
Chen et al, "A Simple Framework for Contrastive Learning of Visual Representations", ICML 2020

Chen et al, "An Empirical Study of Training Self-Supervised Vision Transformers", ICCV 2021

# Supervised vs. Self-Supervised Learning

- SSL pretraining on ImageNet matches its supervised counterpart

Contrastive SSL Pretraining then Finetuning on Detection



# Masked Autoencoders (MAE)

A ~~new~~ old method dethrones contrastive learning? Denoising Autoencoder with Vision Transformer

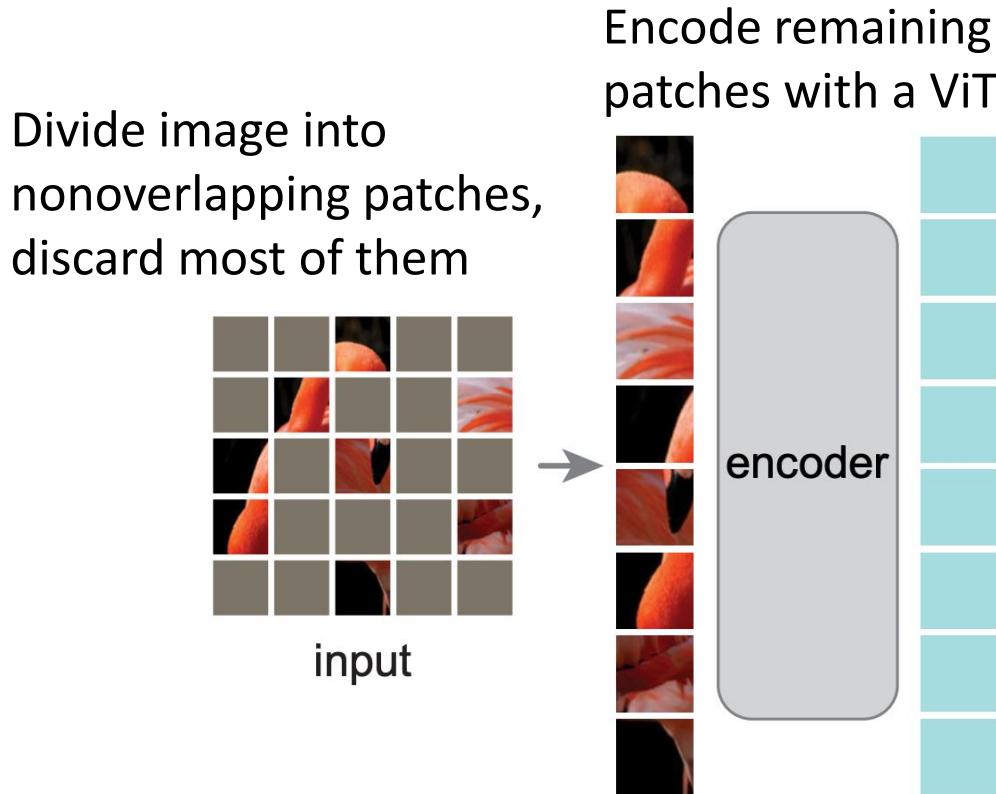
Divide image into  
nonoverlapping patches,  
discard most of them



input

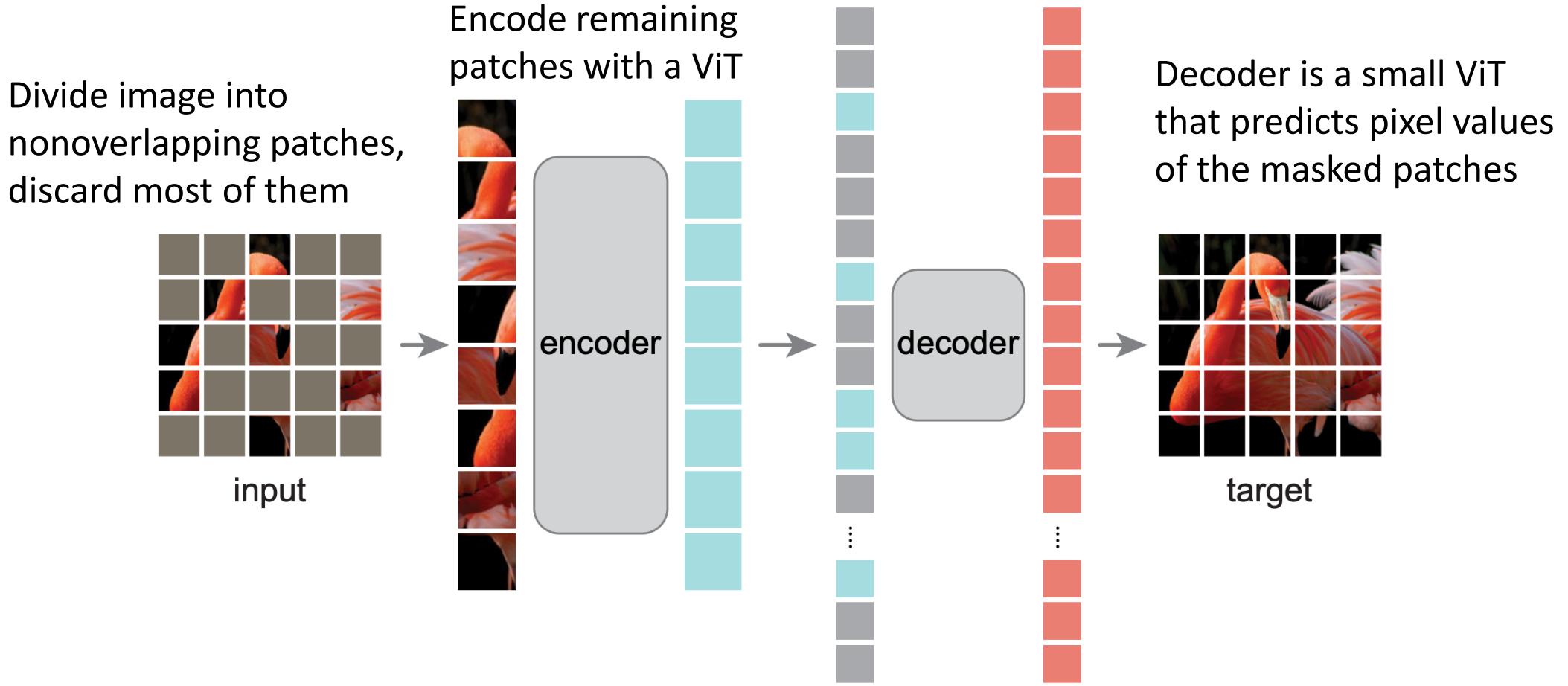
# Masked Autoencoders (MAE)

A new old method dethrones contrastive learning? Denoising Autoencoder with Vision Transformer



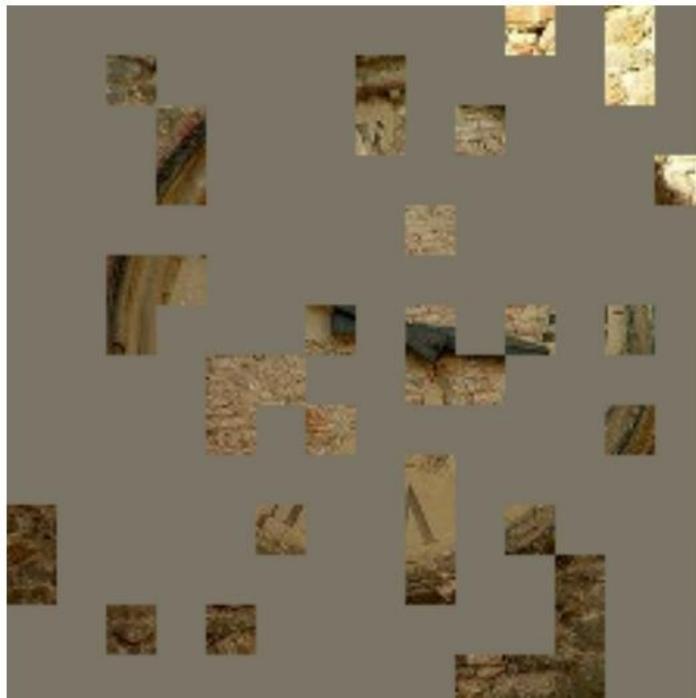
# Masked Autoencoders (MAE)

A new old method dethrones contrastive learning? Denoising Autoencoder with Vision Transformer



# Masked Autoencoders (MAE): Reconstructions

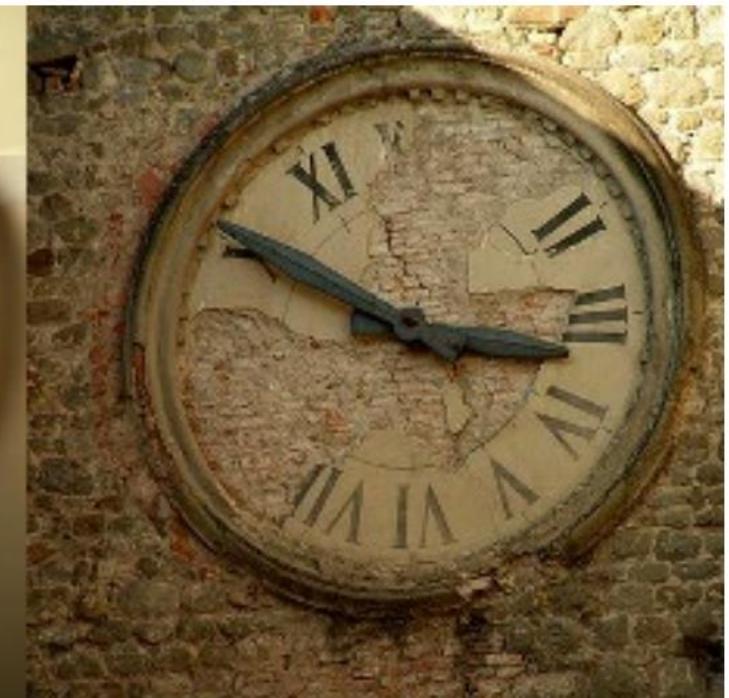
Input Patches



Prediction



Actual Image

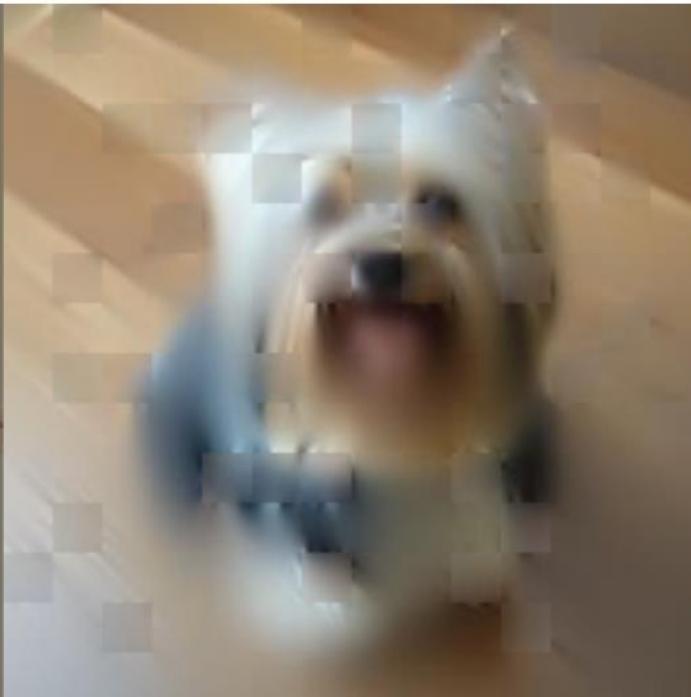


# Masked Autoencoders (MAE): Reconstructions

Input Patches



Prediction



Actual Image



# Masked Autoencoders (MAE): Reconstructions

Input Patches



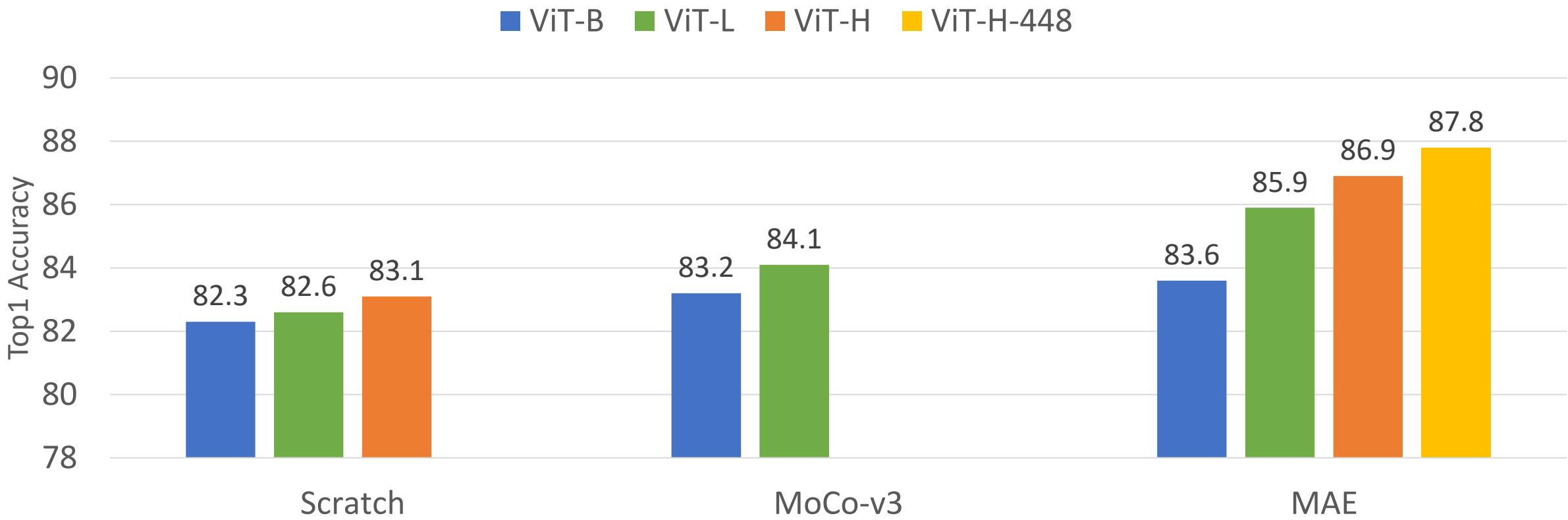
Prediction



Actual Image



## SSL Pretraining, then finetuning for ImageNet Classification



MAE Pretraining outperforms training from scratch, and allows scaling to larger ViT models

# What Is Going on These Days

- Active ongoing researches on self-supervised learning
  - SimCLR is cited > 7k times in < 3 years (Feb 13, 2020 ~ Nov 24, 2022)  
[A Simple Framework for Contrastive Learning of Visual ... - arXiv](#)  
T Chen 저술 · 2020 · 7103회 인용 — Abstract: This paper presents SimCLR: a simple framework for contrastive learning of visual representations. We simplify recently...
- Further improvements
  - [SwAV] Caron et al, “Unsupervised Learning of Visual Features by Contrasting Cluster Assignments”, NeurIPS 2020
  - [OBoW] Gidaris et al, “OBoW: Online Bag-of-Visual-Words Generation for Self-Supervised Learning”, CVPR 2021
  - [DeiT] Brigato et al, “Training data-efficient image transformers & distillation through attention”, ICML 2021
  - [Barlow Twins] Zbontar et al, “Barlow Twins: Self-Supervised Learning via Redundancy Reduction”, ICML 2021
  - [MoCo v3] Chen et al., “An empirical study of training self-supervised Vision Transformers”, ICCV 2021
  - [DINO] Caron et al, “Emerging Properties in Self-Supervised Vision Transformers”, ICCV 2021
  - [VICReg] Bardes et al, “VICReg: Variance-Invariance-Covariance Regularization for Self-Supervised Learning”, ICLR 2022
  - [BEiT] Bao et al, “BEiT: BERT Pre-Training of Image Transformers”, ICLR 2022
  - [SimMIM] Xie et al, “SimMIM: A Simple Framework for Masked Image Modeling”, CVPR 2022
  - [MSN] Assran et al, “Masked Siamese Networks for Label-Efficient Learning”, arXiv 2022
  - [MCMAE] Gao et al, “MCMAE: Masked Convolution Meets Masked Autoencoders”, NeurIPS 2022
  - [CAE] X Chen et al, “Context Autoencoder for Self-Supervised Representation Learning”, arXiv 2022

# An Issue on Self-Supervised Learning in Vision

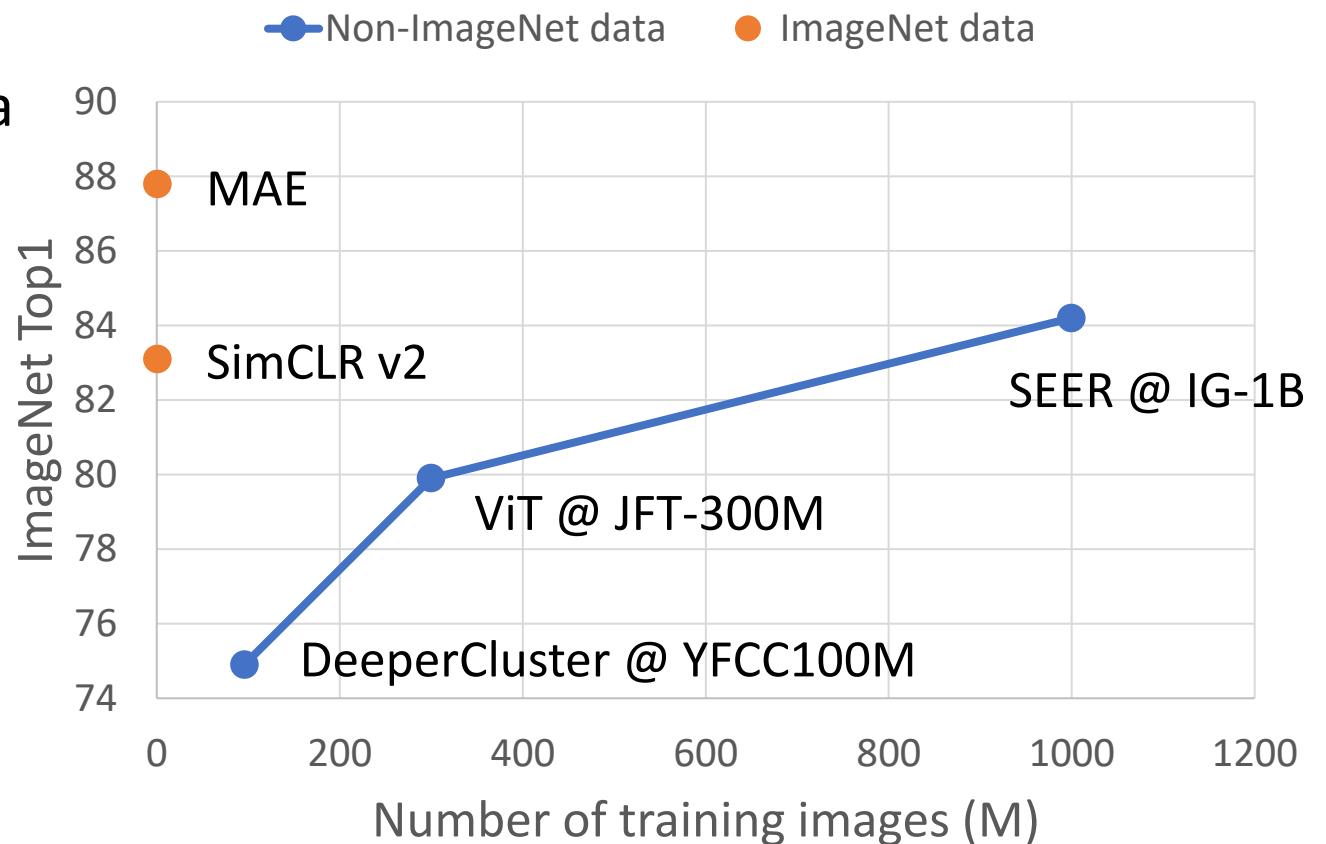
The motivation of SSL is scaling to large data that can't be labeled.

Most papers pretrain on (unlabeled) ImageNet, then evaluate on ImageNet!

Unlabeled ImageNet is still curated: single object per image, balanced classes.

Self-Supervised Learning on larger datasets hasn't been as successful as NLP

Idea: What if we go beyond isolated images?



Caron et al, "Unsupervised pre-training of images features on non-curated data", ICCV 2019  
Chen et al, "Big self-supervised models are strong semi-supervised learners", NeurIPS 2020  
Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021  
Goyal et al, "Self-supervised Pretraining of Visual Features in the Wild", arXiv 2021  
He et al, "Masked Autoencoders are Scalable Vision Learners", CVPR 2022

# Multimodal Self-Supervised Learning

Don't learn from isolated images -- take images together with some **context**

## **Video:** Image together with adjacent video frames

Agrawal et al, "Learning to See by Moving", ICCV 2015

Wang et al, "Unsupervised Learning of Visual Representations using Videos", ICCV 2015

Pathak et al, "Learning Features by Watching Objects Move", CVPR 2017

## **Sound:** Image with audio track from video

Owens et al, "Ambient Sound Provides Supervision for Visual Learning", ECCV 2016

Arandjelovic and Zisserman, "Look, Listen and Learn", ICCV 2017

## **3D:** Image with depth map or point cloud

Xie et al, "PointContrast: Unsupervised Pre-training for 3D Point Cloud Understanding", ECCV 2020

Zhang et al, "Self-supervised pretraining of 3D features on any point-cloud", CVPR 2021

## **Language:** Image with natural-language text

Sariyildiz et al, "Learning Visual Representations with Caption Annotations", ECCV 2020

Desai and Johnson, "VirTex: Learning Visual Representations from Textual Annotations", CVPR 2021

Radford et al, "Learning Transferable Visual Models from Natural Language Supervision", ICML 2021

Jia et al, "Scaling up Visual and Vision-Language Representation Learning with Noisy Text Supervision", ICLR 2021

Desai et al, "RedCaps: Web-curated Image-Text data created by the people, for the people", NeurIPS 2021

# Why Language?

Large dataset of  
(image, caption)



a dog with his  
head out the  
window of the car



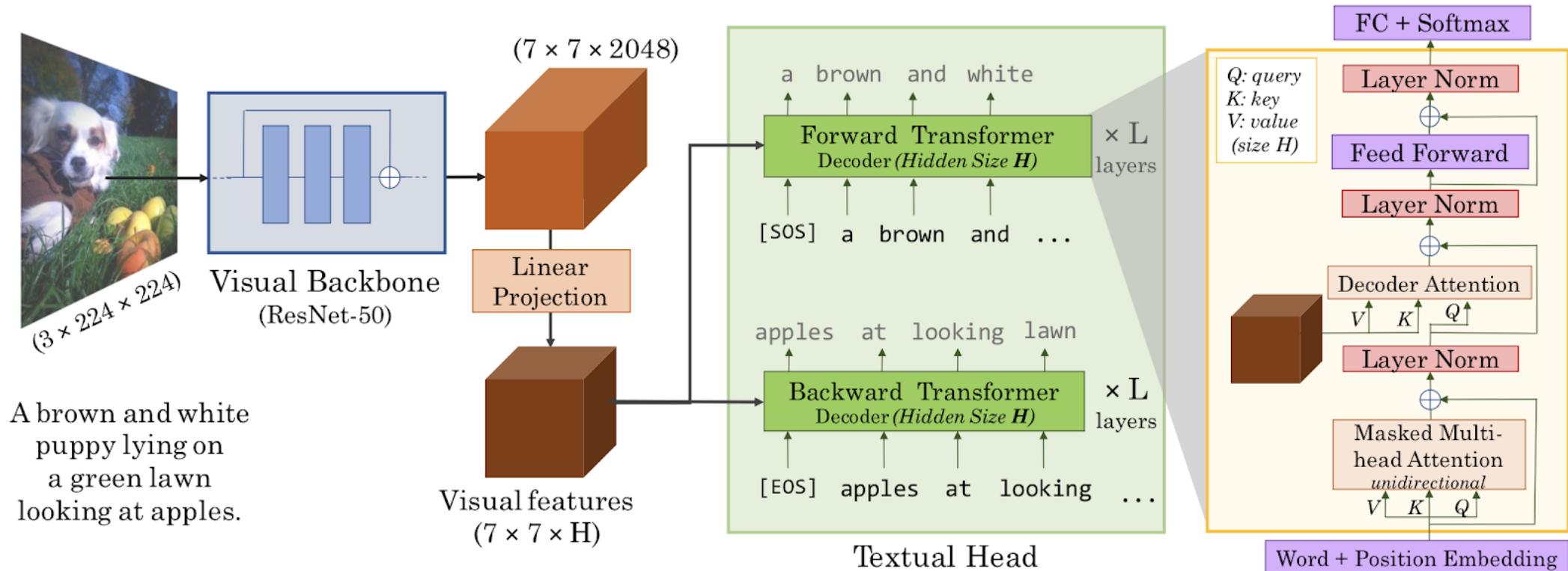
a black and orange  
cat is resting on a  
keyboard and yellow  
back scratcher

**1. Semantic density:** Just a few words give rich information

**2. Universality:** Language can describe any concept

**3. Scalability:** Non-experts can easily caption images; data can also be collected from the web at scale

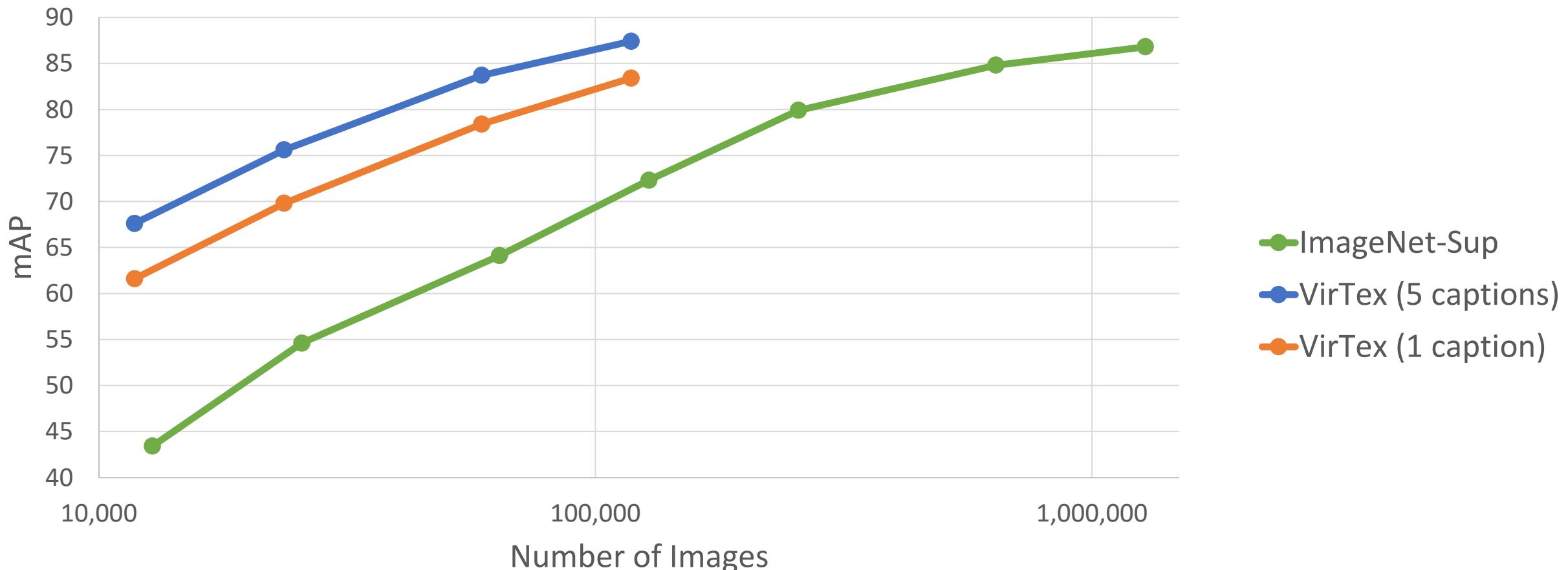
# Generating Captions



Desai and Johnson, "Desai and Johnson, "VirTex: Learning Visual Representations from Textual Annotations", CVPR 2021

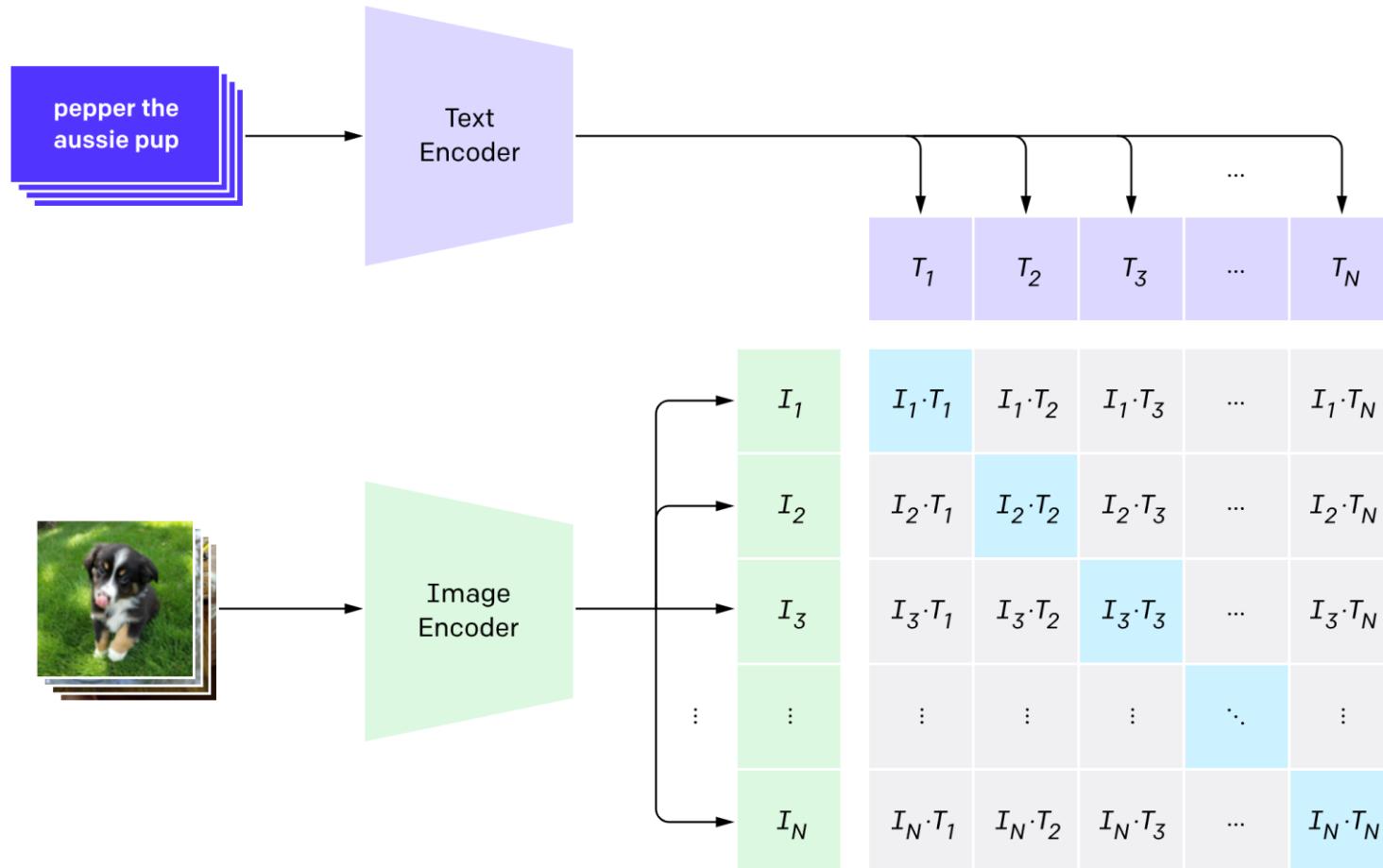
# Generating Captions

## PASCAL VOC Linear Classification



Desai and Johnson, "VirTex: Learning Visual Representations from Textual Annotations", CVPR 2021

# Matching Images and Text: CLIP



Contrastive loss: Each image predicts which caption matches

Large-scale training on 400M (image, text) pairs from the internet

Radford et al, "Learning Transferable Visual Models from Natural Language Supervision", ICML 2021

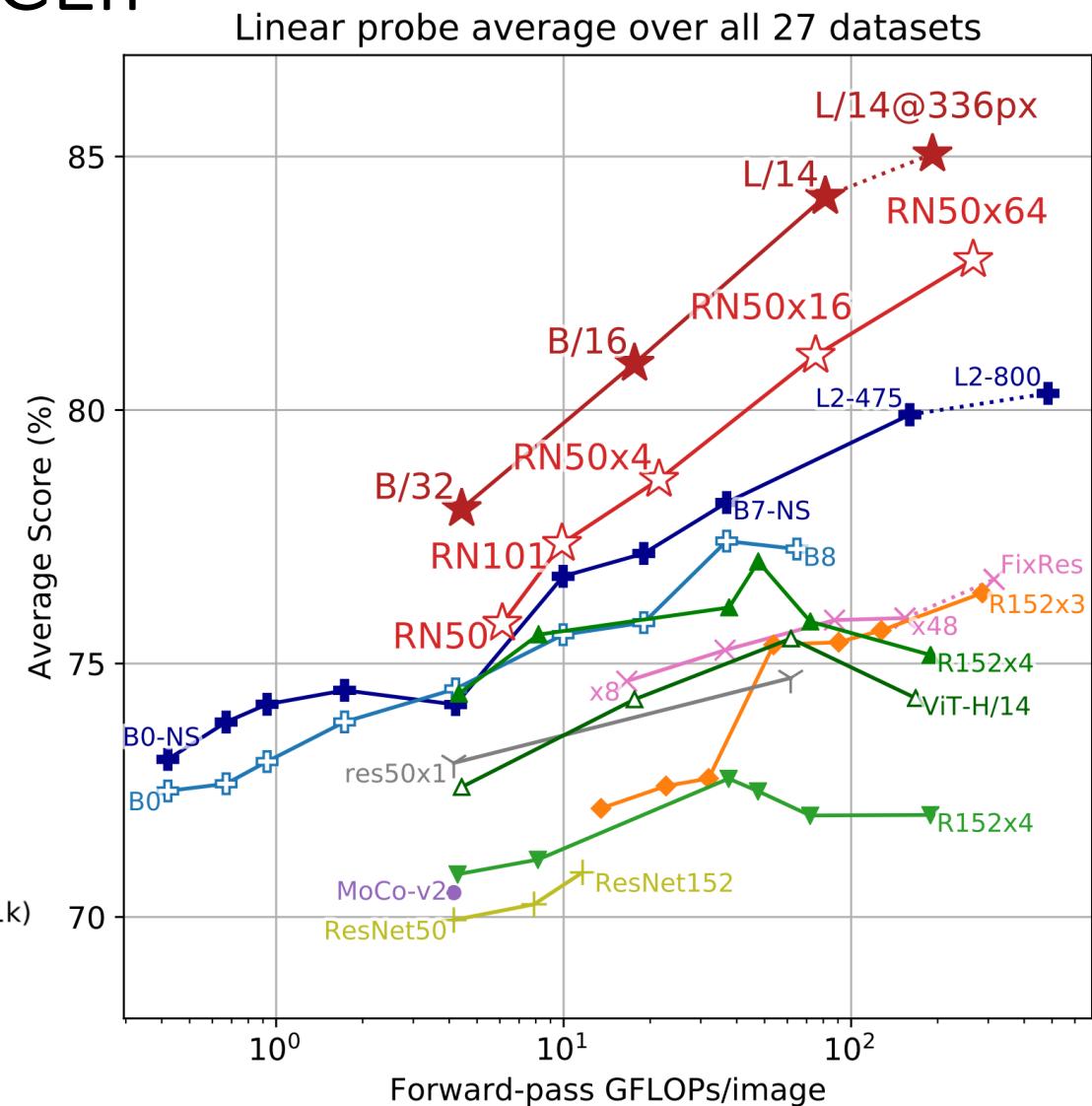
Jia et al, "Scaling up Visual and Vision-Language Representation Learning with Noisy Text Supervision", ICLR 2021

# Matching Images and Text: CLIP

Very strong performance on many downstream vision problems!

Performance continues to improve with larger models

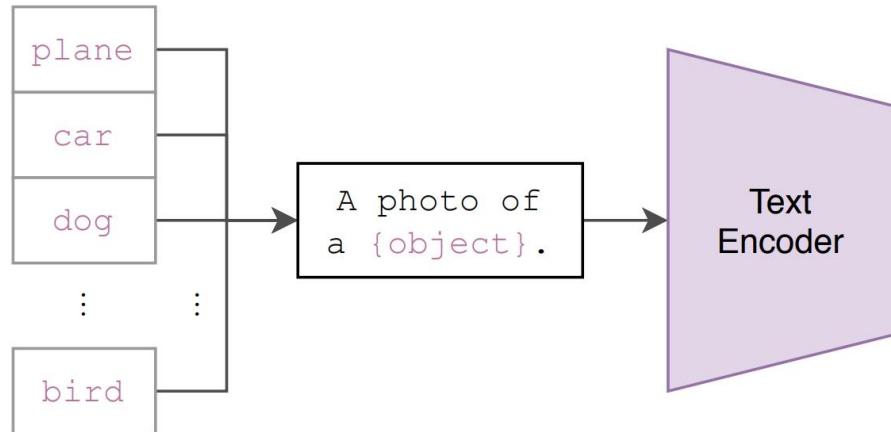
- ★ CLIP-ViT
- ★ CLIP-ResNet
- EfficientNet-NoisyStudent
- + EfficientNet
- Instagram-pretrained
- SimCLRv2
- BYOL
- MoCo
- ViT (ImageNet-21k)
- BiT-M
- BiT-S
- ResNet



Radford et al, "Learning Transferable Visual Models from Natural Language Supervision", ICML 2021

# CLIP: Zero-Shot Classification

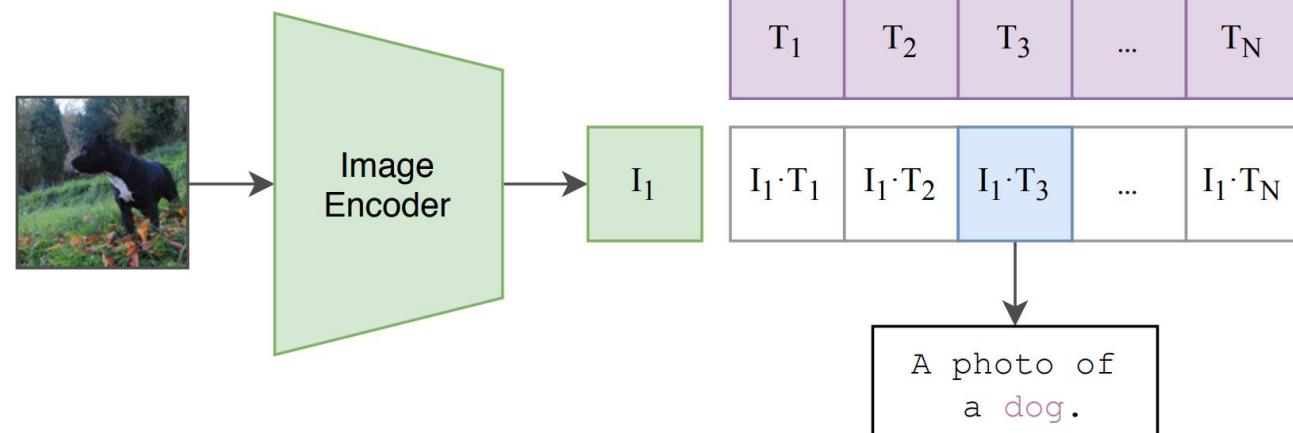
(2) Create dataset classifier from label text



**Caution:** CLIP training dataset is private; can't reproduce results

(3) Use for zero-shot prediction

Language enables **zero-shot classification**:  
Classify images into categories without any additional training data!



Radford et al, "Learning Transferable Visual Models from Natural Language Supervision", ICML 2021

# DALL-E 2: Text-to-Image Generation

A rabbit detective sitting on a park bench and reading a newspaper in a victorian setting



A shark and a dolphin cruise hand-in-hand with an undersea city in the background



Robot dinosaurs versus monster trucks in the colosseum



Ramesh et al, "DALL-E 2", 2022. <https://openai.com/dall-e-2/>  
Source: <https://twitter.com/sama/status/1511724264629678084>

# Phenaki: Text-to-Video Generation

A photorealistic teddy bear is swimming in  
the ocean at San Francisco

The teddy bear goes under water

The teddy bear keeps swimming under the  
water with colorful fishes

A panda bear is swimming under water

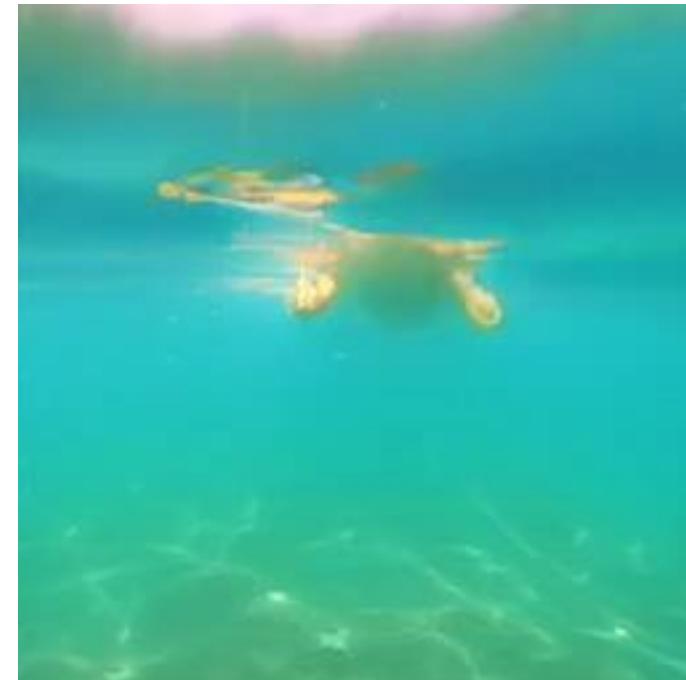


A teddy bear diving in the ocean

A teddy bear emerges from the water

A teddy bear walks on the beach

Camera zooms out to the teddy bear in the  
campfire by the beach



Side view of an astronaut is walking through  
a puddle on mars

The astronaut is dancing on mars

The astronaut walks his dog on mars

The astronaut and his dog watch fireworks



# Representation Learning Summary

- Transferring Supervised Learning model was pervasive in the 2010s
- Self-Supervised Learning (SSL) aims to scale up to larger datasets without human annotation; promising strategy in the 2020s
- First train for a **pretext** task, then **transfer** to **downstream** tasks
- Many pretext tasks: autoencoding, context prediction, jigsaw, inpainting, colorization, clustering, rotation prediction, ...
- SSL has been successful in natural language processing
- Intense research on SSL in vision; current best are
  - Contrastive learning (and variants like feature reconstruction)
  - Masked autoencoding
- Multimodal SSL has been very successful; it seems promising!