Lecture 19:



Contrastive and Self-Supervised Visual Representations

CMSC 472/672 Computer Vision



Train a model on 1 million images

- → label 1 million images
- Labels aren't magically given to you
- → need human effort

How much will it cost?

(1,000,000 images) (Small to medium sized dataset)

× (10 seconds/image) (Fast annotation)

 \times (1/3600 hours/second)

 \times (\$15 / hour) (Minimum wage)

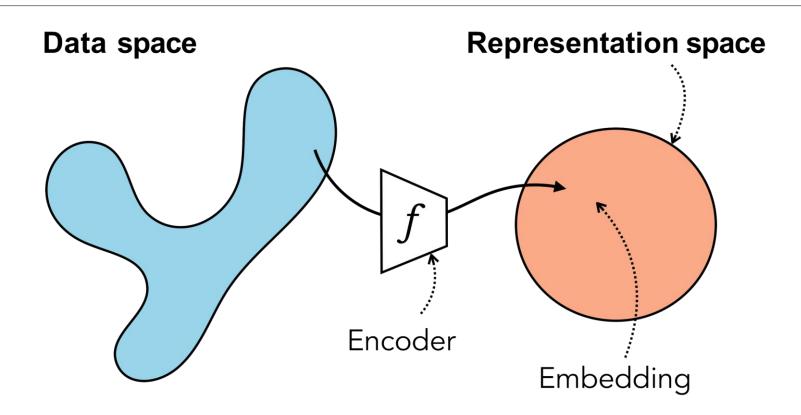
× (3 annotators / image) (for consensus / removing noise)



without considering overhead / admin costs ...

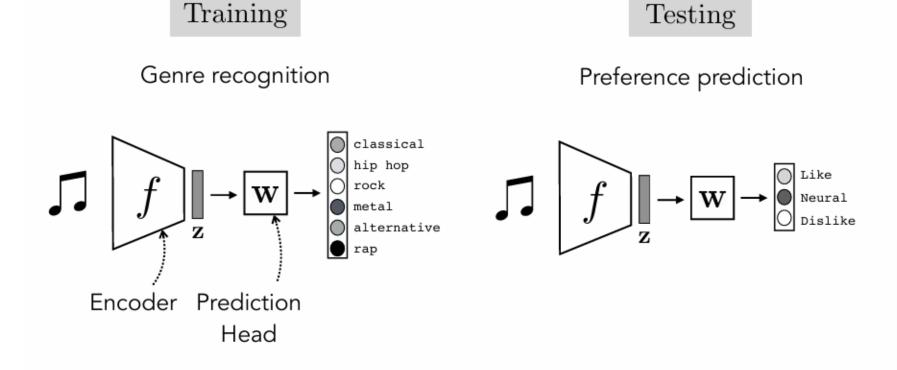
Recap: Representation Learning "x2vec"

- \circ A representation of a data domain \mathcal{X} is a function $f: \mathcal{X} \to \mathbb{R}^d$ (an encoder) that assigns a feature vector to each input in that domain.
- \circ A representation of a datapoint is a vector $z \in \mathbb{R}^d$ with z = f(x).



Why Learn Representations?

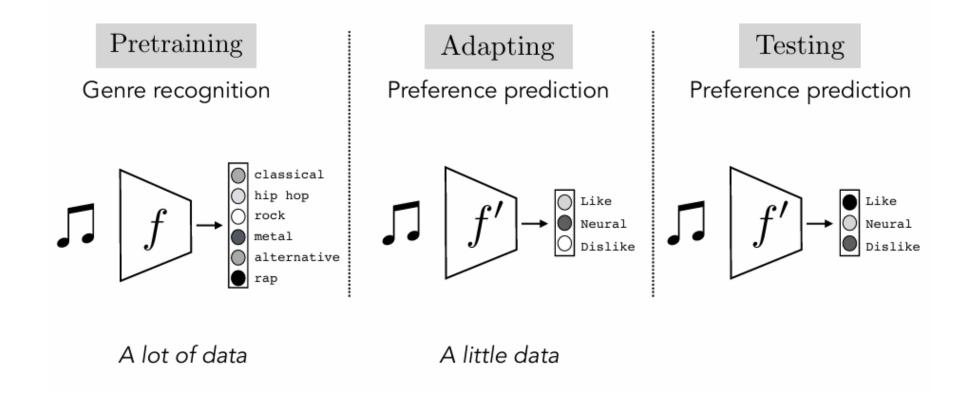
"Generally speaking, a good representation is one that makes a subsequent learning task easier."
- Goodfellow et al. "Deep Learning". 2016



Often, what we will be "tested" on is not what we were trained on.

Why Learn Representations?

"Generally speaking, a good representation is one that makes a subsequent learning task easier."
- Goodfellow et al. "Deep Learning". 2016



Learning from examples

(aka supervised learning)

Training data

. . .

$$f^* = \underset{f \in \mathcal{F}}{\operatorname{arg\,min}} \sum_{i=1}^{N} \mathcal{L}(f(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$$

Learning without examples

(includes unsupervised learning / self-supervised learning)

Data

$$\begin{cases} x^{(1)} \\ \{x^{(2)} \} \\ \{x^{(3)} \} \end{cases} \longrightarrow$$
 Learner \longrightarrow ?

Learning without examples

(includes unsupervised learning / self-supervised learning)

Data

 $\{x^{(1)}\}$

 $\{x^{(2)}\}$

 $\{x^{(3)}\}$

Learner

Embeddings

Clusters

Metrics

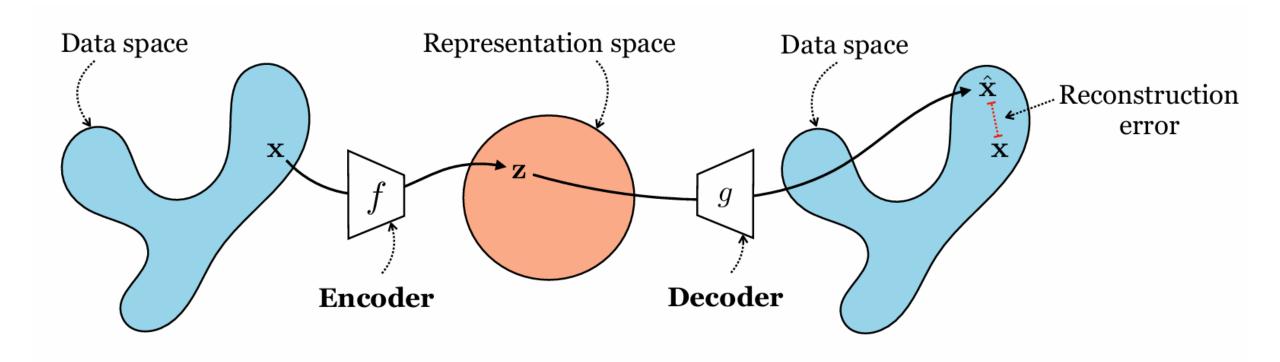
. . .

Two Basic Approaches: (1) Compression (2) Prediction

Learning Method	Learning Principle	Short Summary	
Autoencoding	Compression	Remove redundant information	
Contrastive	Compression	Achieve invariance to viewing transformations	
Clustering	Compression	Quantize continuous data into discrete categories	
Future prediction	Prediction	Predict the future	
Imputation	Prediction	Predict missing data	
Pretext tasks	Prediction	Predict abstract properties of your data	

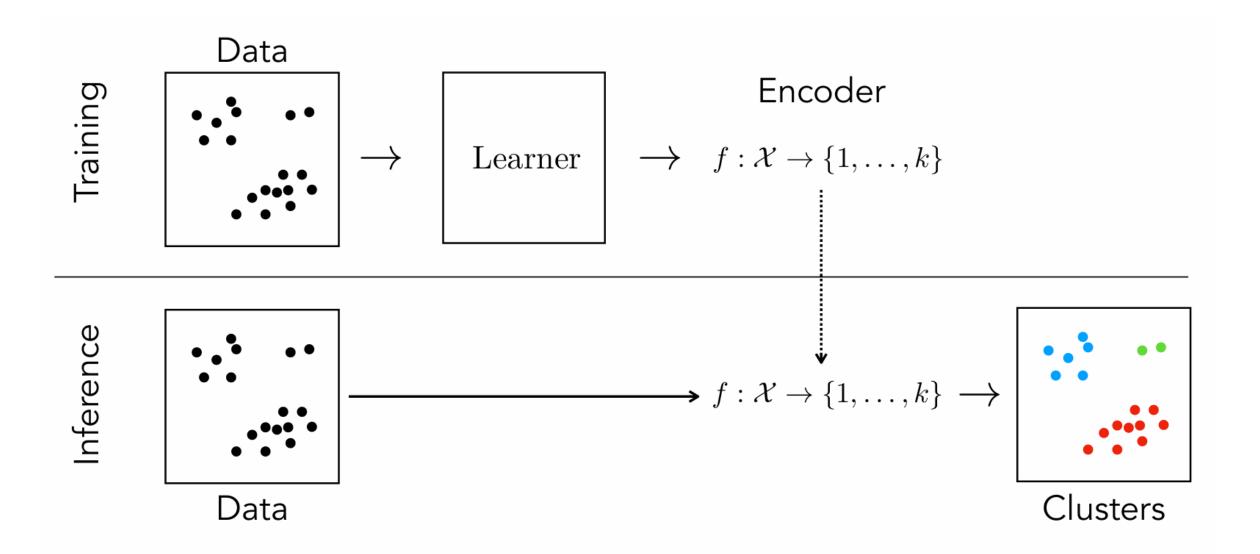
Some examples of the "Compression" Approach:

Recap: Autoencoder

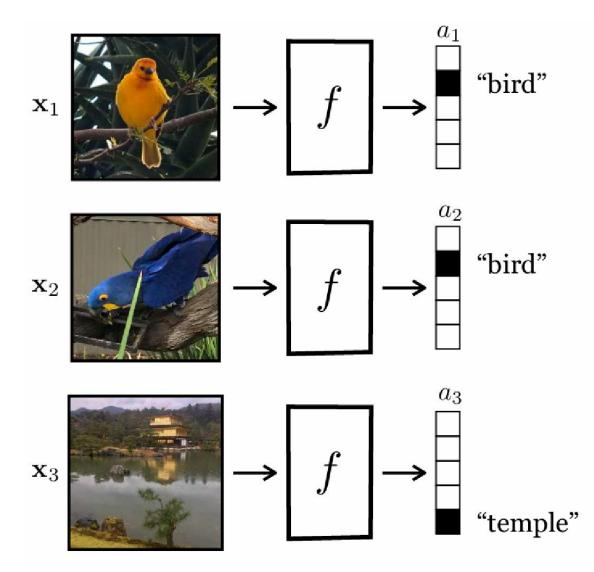


$$f^*, g^* = \underset{f,g}{\operatorname{arg\,min}} \mathbb{E}_{\mathbf{x}} \|\mathbf{x} - g(f(\mathbf{x}))\|_2^2$$

Clustering



Clustering

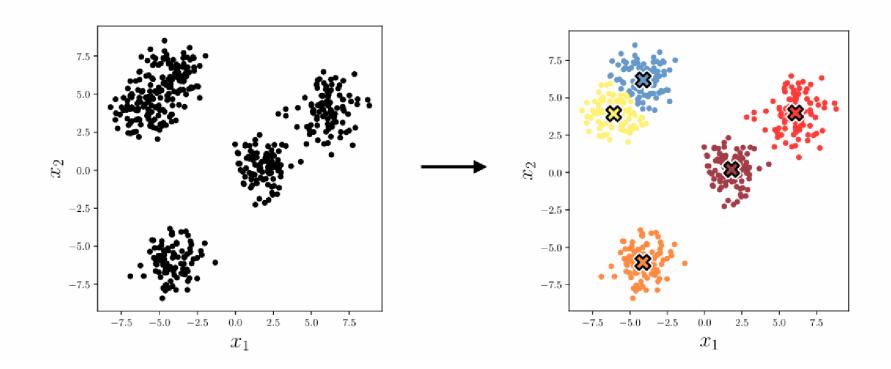


 What's the best representation that humans have come up with so far?

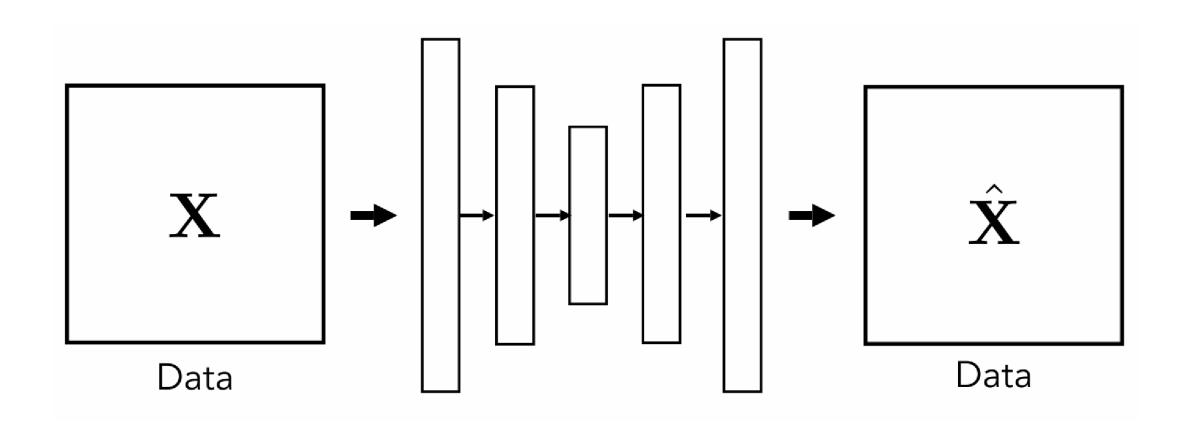
- Language!
- Words are the atoms of language
- Clustering is the problem of making up new words for things

Clustering Algorithm: k-means

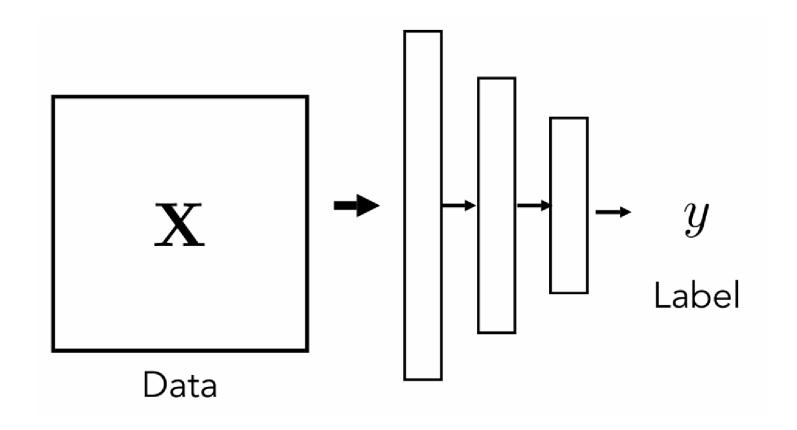
- Map datapoints to integers (i.e. cluster)
- In such a way that each datapoint is as close as possible to the mean of the cluster it is assigned to



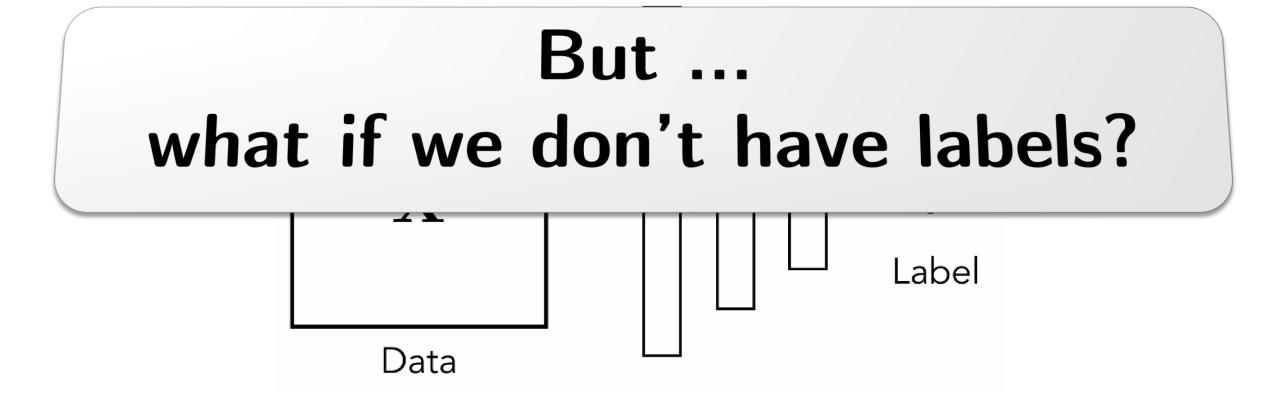
The "Compression" Approach



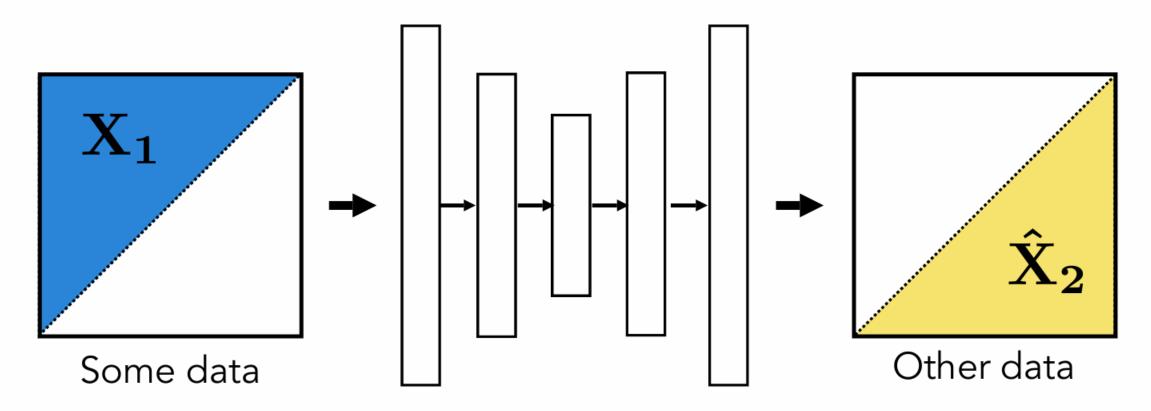
The "Prediction" Approach for Representation Learning



The "Prediction" Approach for Representation Learning



Data prediction aka "self-supervised learning"



Self-Supervised Learning

Build methods that learn from "raw" data (inputs only) — no labels!

Unsupervised Learning:

older terminology ... model isn't told what to predict

Self-Supervised Learning:

o model is trained to predict some natural occurring signal rather than predicting labels

Semi-Supervised Learning:

o train jointly with some labeled data and a lot of unlabeled data.

Self-Supervised Learning: A trick

- If you don't have labels, make labels.
- Convert "unsupervised" problem into "supervised"
- Cook up labels (prediction targets) from the data itself
 - This is often called a "pretext" task

Claim:

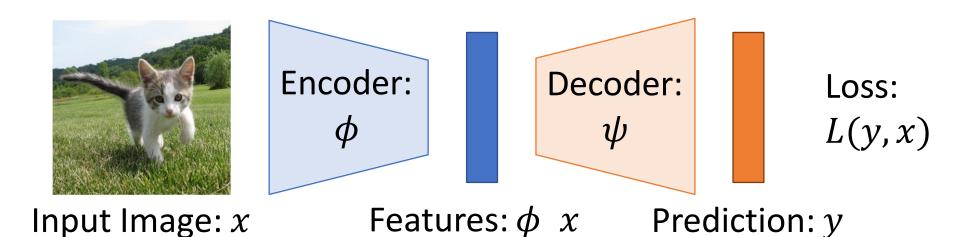
Training a model for "pretext" task can lead to very good representations



SSL: "Pretext then transfer"

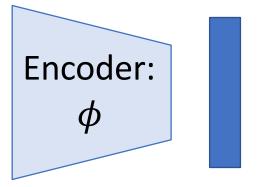
Step 1: Pretrain a network on a pretext task that doesn't require supervision

Step 2: Transfer encoder to downstream tasks via linear classifiers, KNN, finetuning





Input Image: x



Features: ϕx

Downstream tasks: Image classification, object detection, semantic segmentation

Some Examples of Pretext Tasks

Pretext task:	Class prediction	Future frame prediction	Next pixel prediction
$egin{array}{c} \mathbf{Model} \\ \mathbf{schematic} : \end{array}$	g g f x	y g z f x	$egin{array}{cccccccccccccccccccccccccccccccccccc$

Examples of Pretext Tasks

Generative:

Predict part of the input signal

- Autoencoders
 (sparse, denoising, masked)
- Autoregressive
- GANs
- Colorization
- Inpainting

Discriminative:

Predict something about the input signal

- Context prediction
- Rotation
- Clustering
- Contrastive

Multimodal:

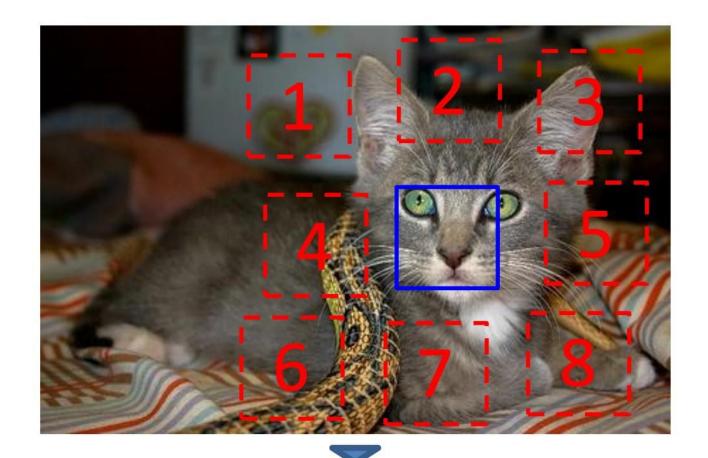
Use some signal in addition to RGB images

- Video
- 3D
- Sound
- Language

Context Prediction

Model predicts relative location of two patches from the same image. <u>Discriminative</u> pretraining task

Intuition: Requires understanding objects and their parts



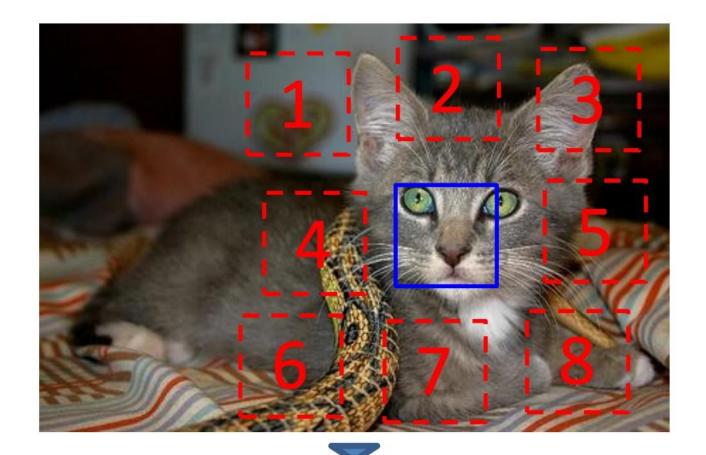
$$X = (\emptyset, \emptyset);$$

Context Prediction

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

<u>Discriminative</u> pretraining task

Intuition: Requires understanding objects and their parts



$$X = (30, 3); Y = 3$$

Context Prediction

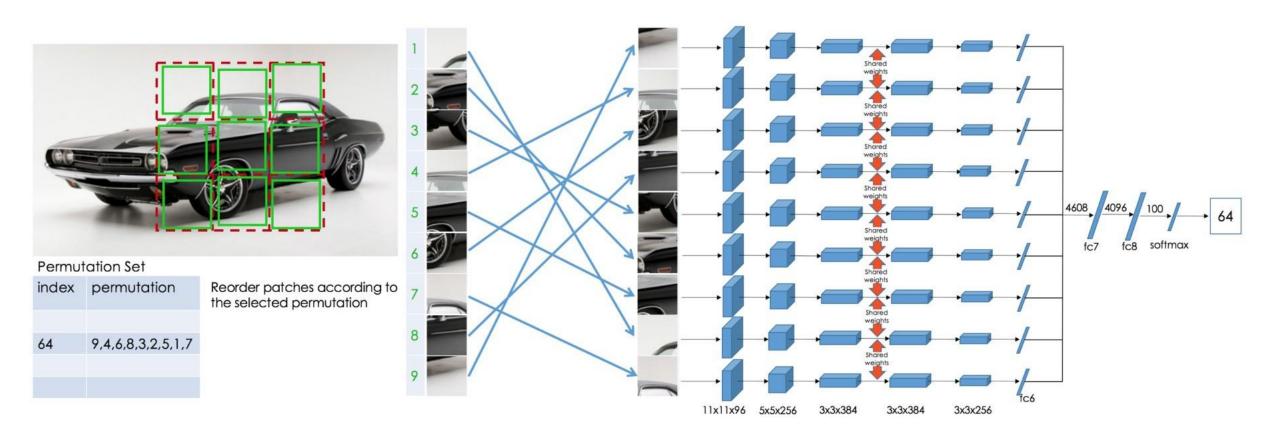
Model predicts relative location of two patches from the same image. <u>Discriminative</u> pretraining task

Intuition: Requires understanding objects and their parts

Classification over 8 positions Concatenate **CNN CNN** Shared Weights

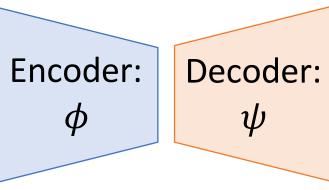
Extension: Solving Jigsaw Puzzles

Rather than predict relative position of two patches, instead predict permutation to "unscramble" 9 shuffled patches



Input Image





Input Image



Encoder: ϕ

Decoder: ψ

Predict Missing Pixels



Human Artist

Input Image



Encoder: ϕ

Decoder: ψ

Predict Missing Pixels



L2 Loss (Best for feature learning)

Input Image



Encoder: ϕ Decoder: ψ

Predict Missing Pixels



L2 + Adversarial Loss (Best for nice images)

Colorization

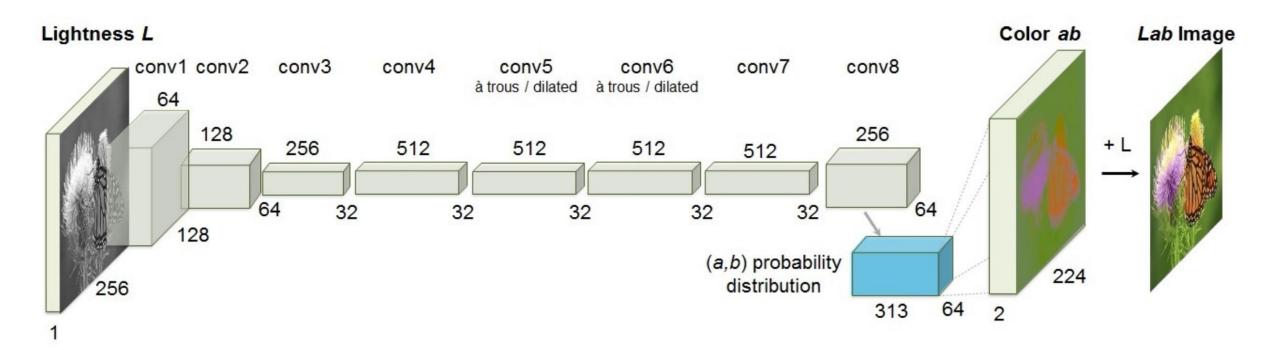
Intuition: A model must be able to identify objects to be able to colorize them



Input: Grayscale Image

Output: Color Image

Colorization



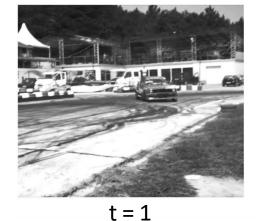
Pretext task: video coloring

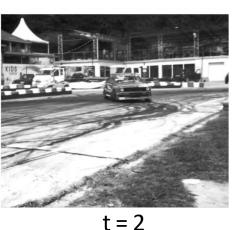
Idea: model the temporal coherence of colors in videos

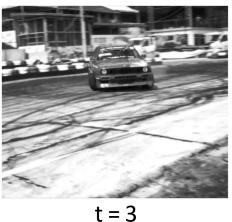
reference frame











t = 0

Source: Vondrick et al., 2018

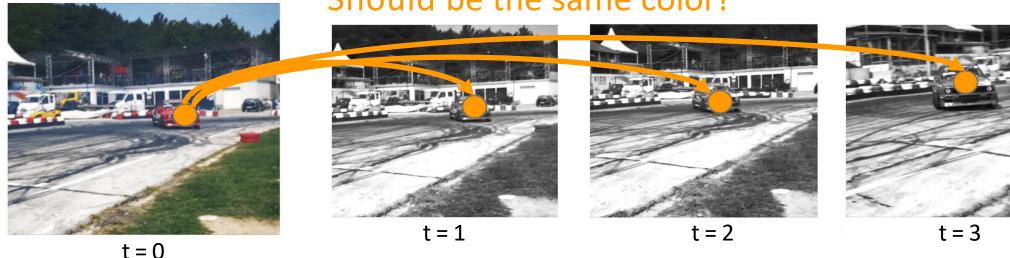
Pretext task: video coloring

Idea: model the temporal coherence of colors in videos

reference frame

how should I color these frames?

Should be the same color!

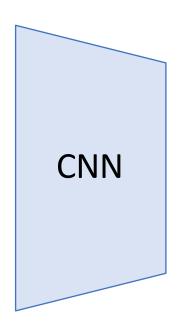


Hypothesis: learning to color video frames should allow model to learn to track regions or objects without labels!

Source: Vondrick et al., 2018

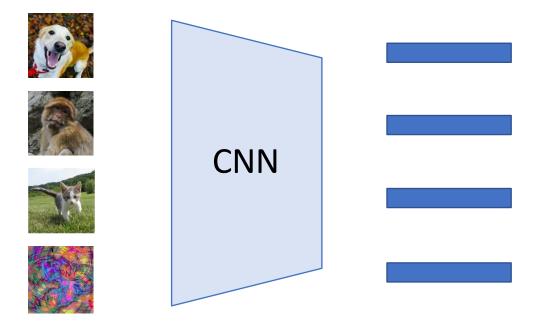
Deep Clustering

(1) Randomly initialize a CNN



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

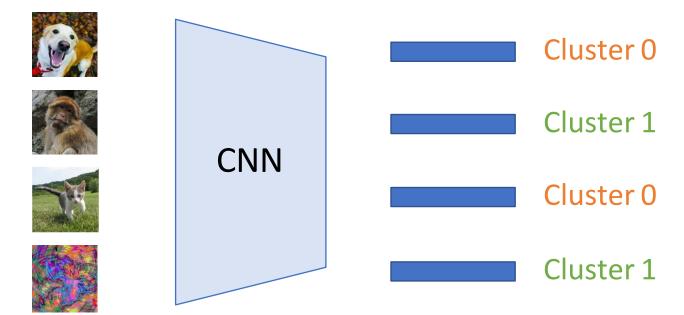
(1) Randomly initialize a CNN



(2) Run many images through CNN, get their final-layer features

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

(1) Randomly initialize a CNN

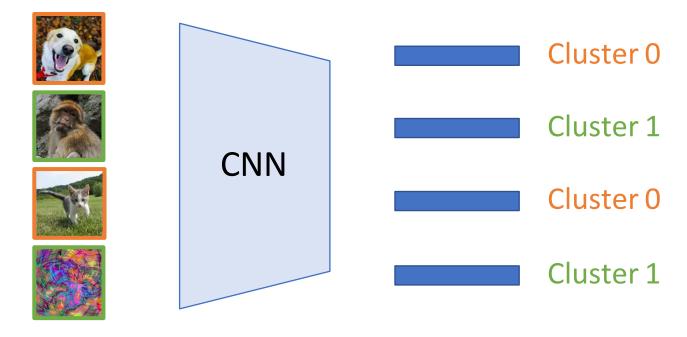


(3) Cluster the features with K-Means; record cluster for each feature

(2) Run many images through CNN, get their final-layer features

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

(1) Randomly initialize a CNN



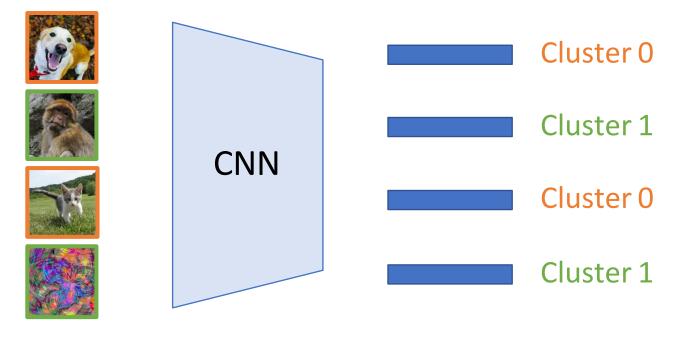
(3) Cluster the features with K-Means; record cluster for each feature

(4) Use cluster assignments as pseudolabels for each image; train the CNN to predict cluster assignments

(2) Run many images through CNN, get their final-layer features

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

(1) Randomly initialize a CNN



(3) Cluster the features with K-Means; record cluster for each feature

(4) Use cluster assignments as pseudolabels for each image; train the CNN to predict cluster assignments

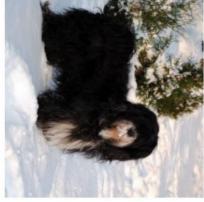
(2) Run many images through CNN, get their final-layer features

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

(5) Repeat: GOTO (2)

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











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







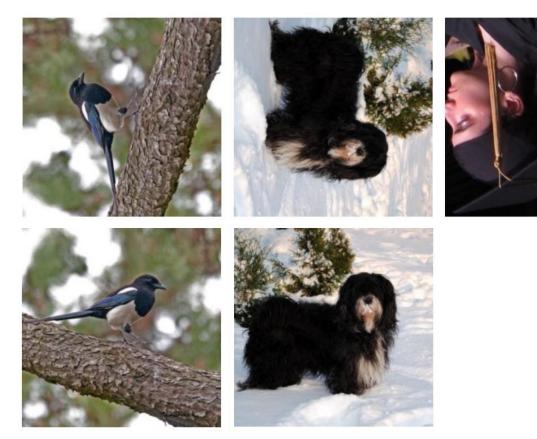






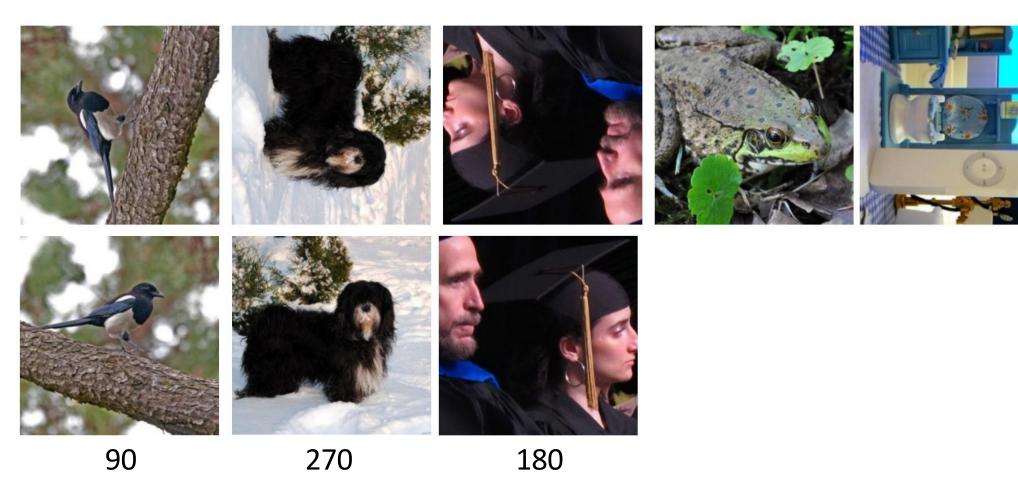
90

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



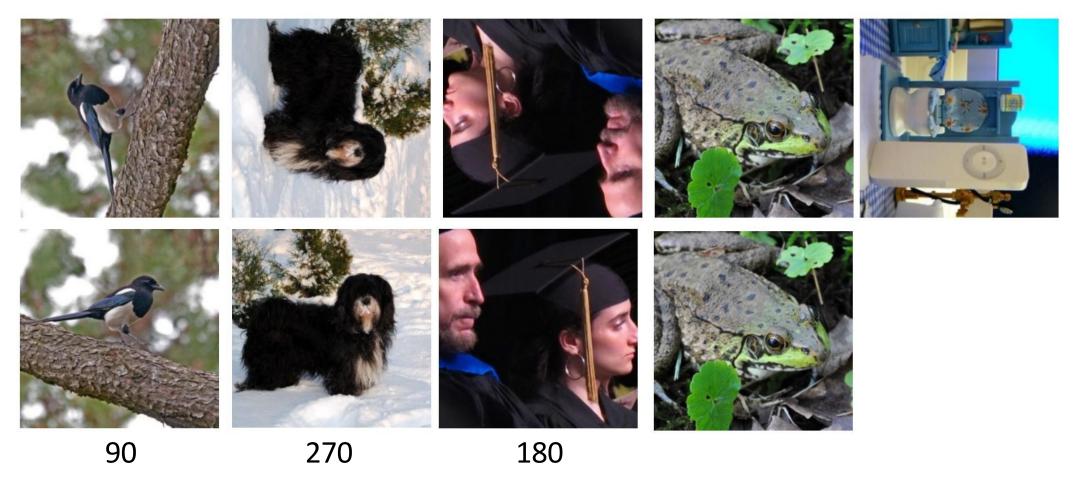


4-way classification task: How much was each image rotated? (0, 90, 180, or 270 degrees)



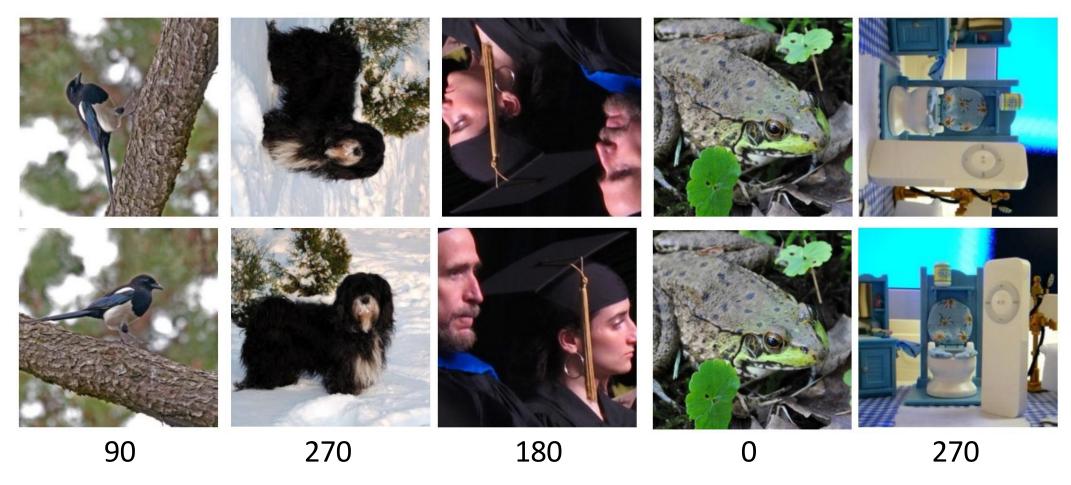
Gidaris et al, "Unsupervised representation learning by predicting image rotations", ICLR 2018

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



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

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



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

Summary: pretext tasks via 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 often do not 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 via image transformations

- Pretext tasks focus on "visual common sense"
 - o e.g., predict rotations, inpainting, rearrangement, and colorization.
- We often do not care about the performance of these pretext tasks
 - o but rather how useful the learned features are for downstream tasks (classification, detection, segmentation).
- Problems:
 - \circ (1) coming up with individual pretext tasks is tedious
 - o (2) the learned representations may not be general.

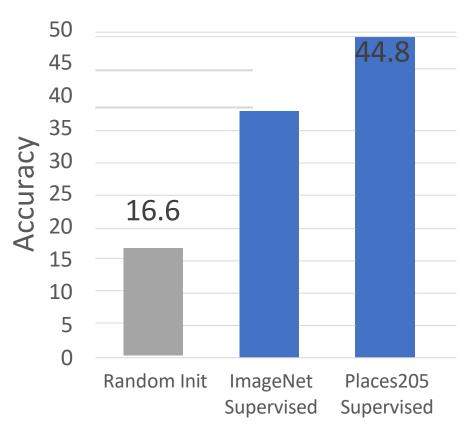
Fair evaluation of SSL methods is very hard ... No theory, so we need to rely on experiments !!!

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? which layer? Linear or nonlinear? Fine-tuning hyperparameters?
 - KNN? What value of K? Normalization on features?

Some papers have tried to do fair comparisons of many SSL methods

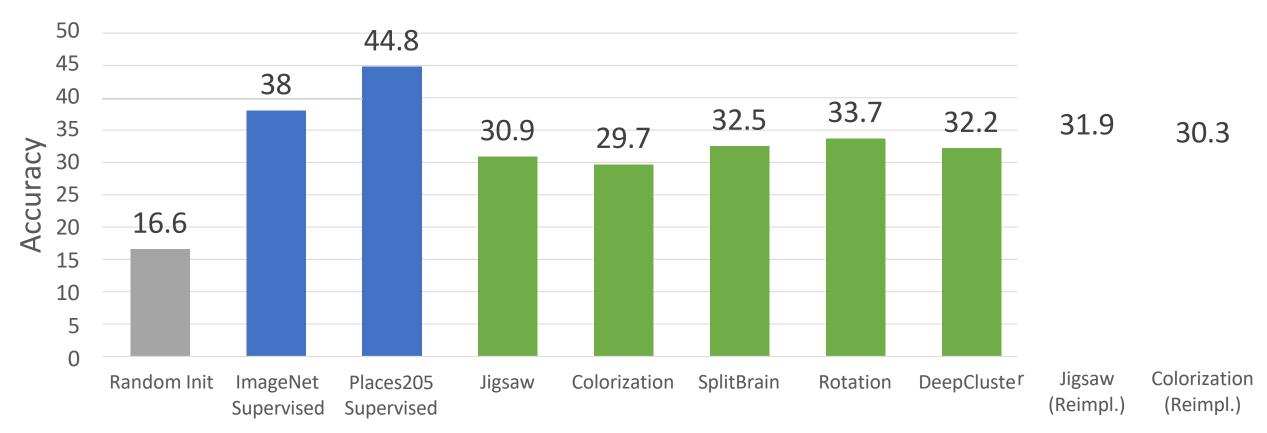
Places 205 Linear Classification from AlexNet conv5



Goyal et al, "Scaling and Benchmarking Self-Supervised Visual Representation Learning", ICCV 2019

Some papers have tried to do fair comparisons of many SSL methods

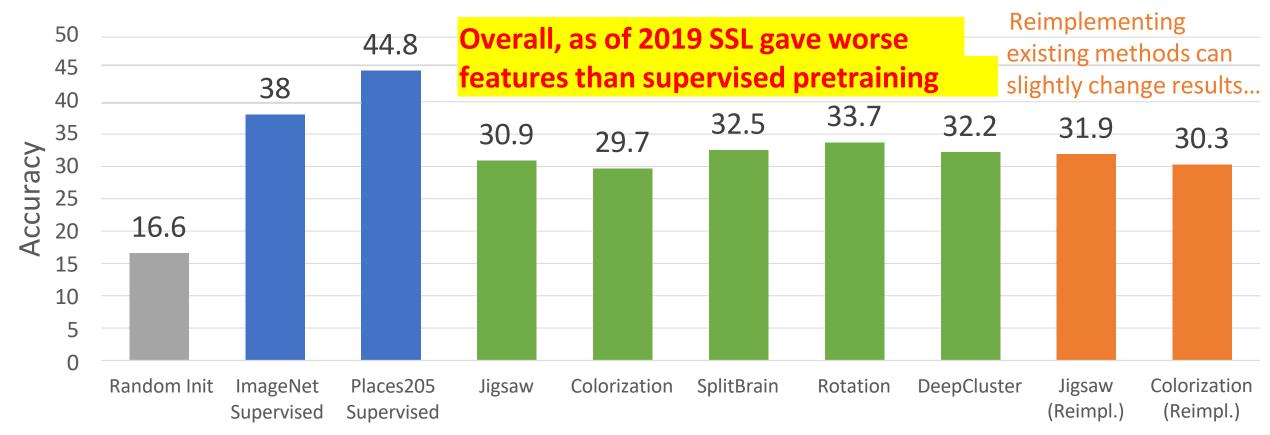
Places 205 Linear Classification from AlexNet conv5



Goyal et al, "Scaling and Benchmarking Self-Supervised Visual Representation Learning", ICCV 2019

Some papers have tried to do fair comparisons of many SSL methods

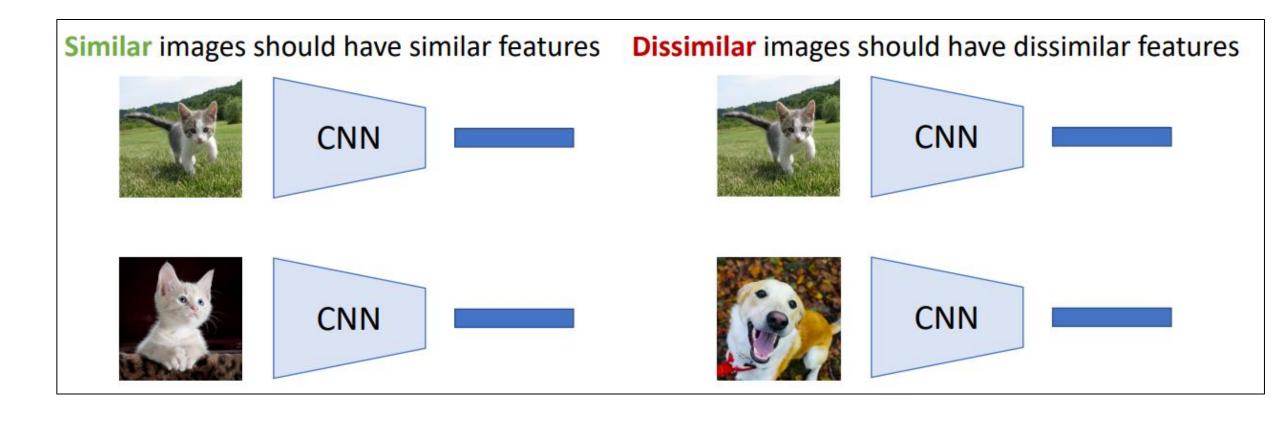
Places 205 Linear Classification from AlexNet conv5



Goyal et al, "Scaling and Benchmarking Self-Supervised Visual Representation Learning", ICCV 2019

Let's take a step back ...

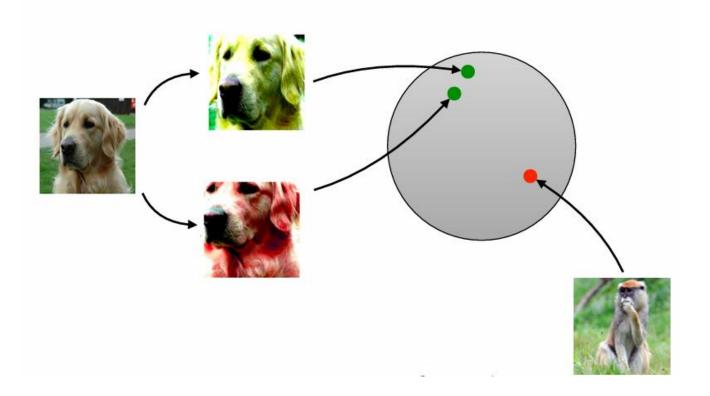
A simpler idea ...



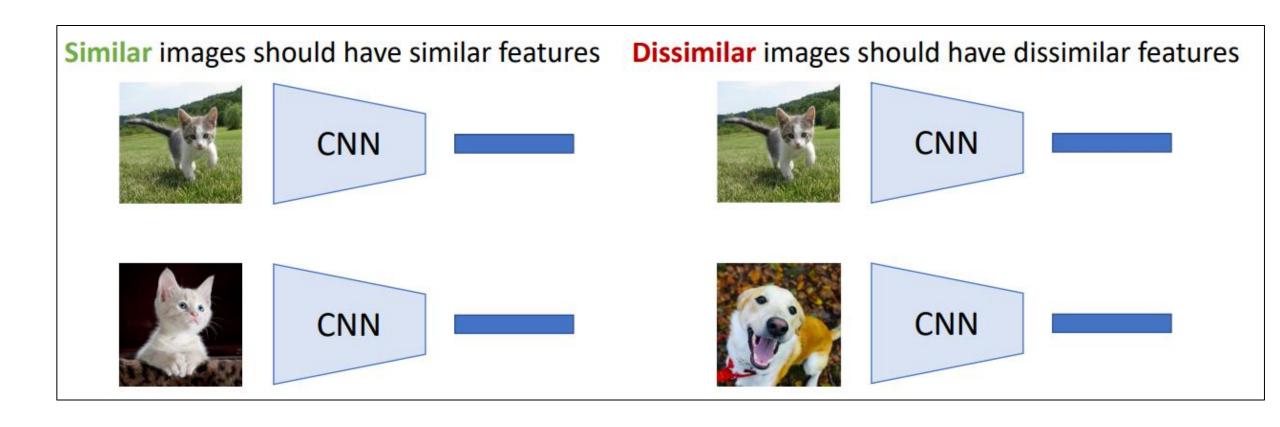
Similarity based Representation Learning

• Build representations via feedback in terms of similarity:

pairs of similar / dissimilar inputs

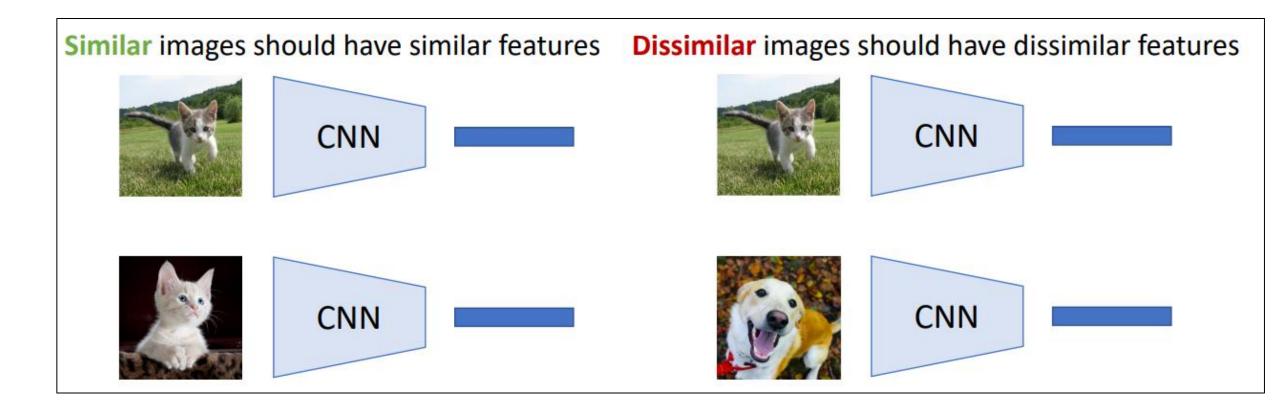


Problem 1: How to compute similarity if we don't have labels for images?



Problem 1: How to compute similarity if we don't have labels for images?

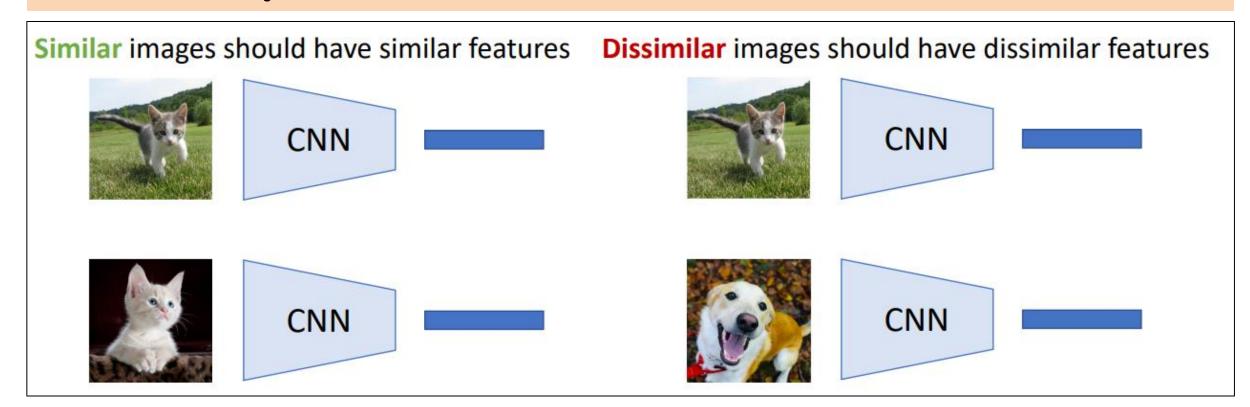
Solution? Euclidean Distance between features $\|\phi(x_1) - \phi(x_2)\|_2$



Problem 1: How to compute similarity if we don't have labels for images?

Solution? Euclidean Distance between features $\|\phi(x_1) - \phi(x_2)\|_2$

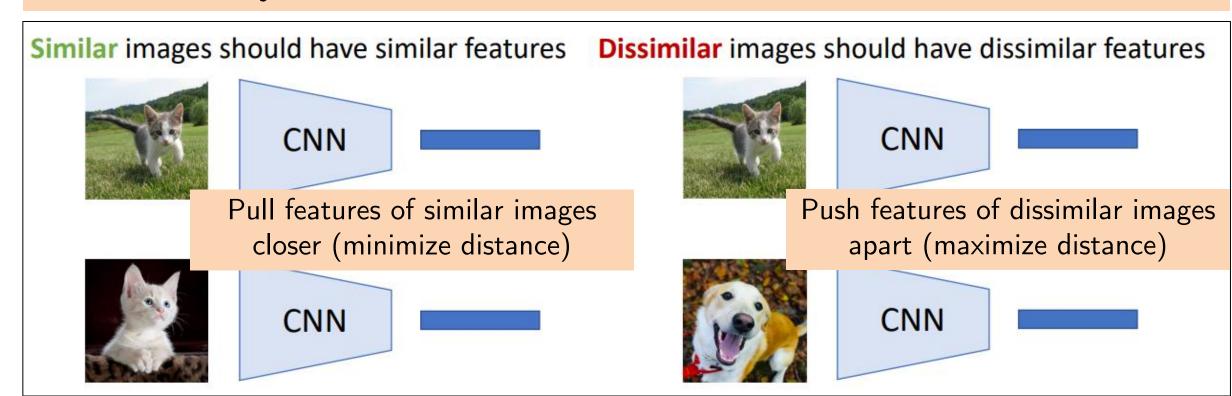
Problem 2: Objective Function?



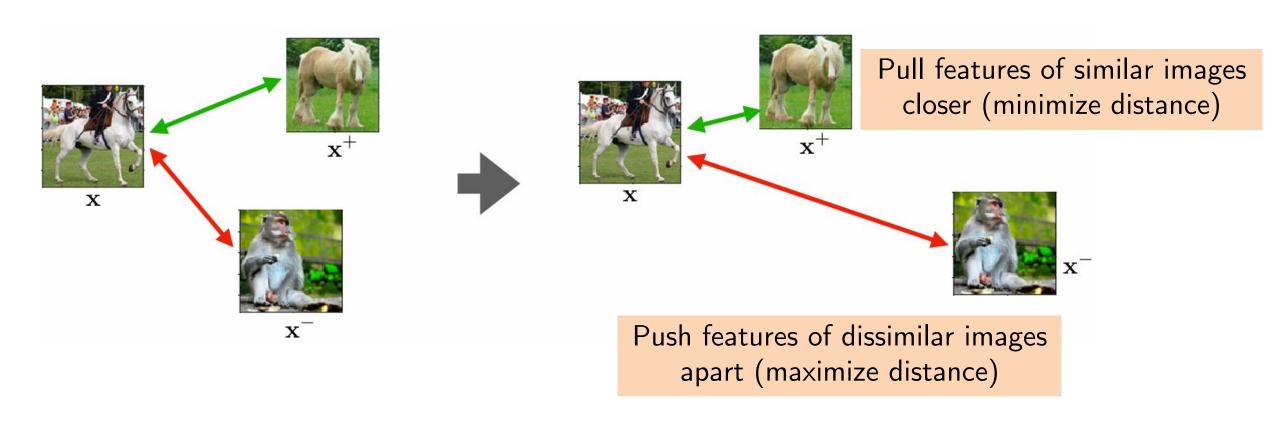
Problem 1: How to compute similarity if we don't have labels for images?

Solution? Euclidean Distance between features $\|\phi(x_1) - \phi(x_2)\|_2$

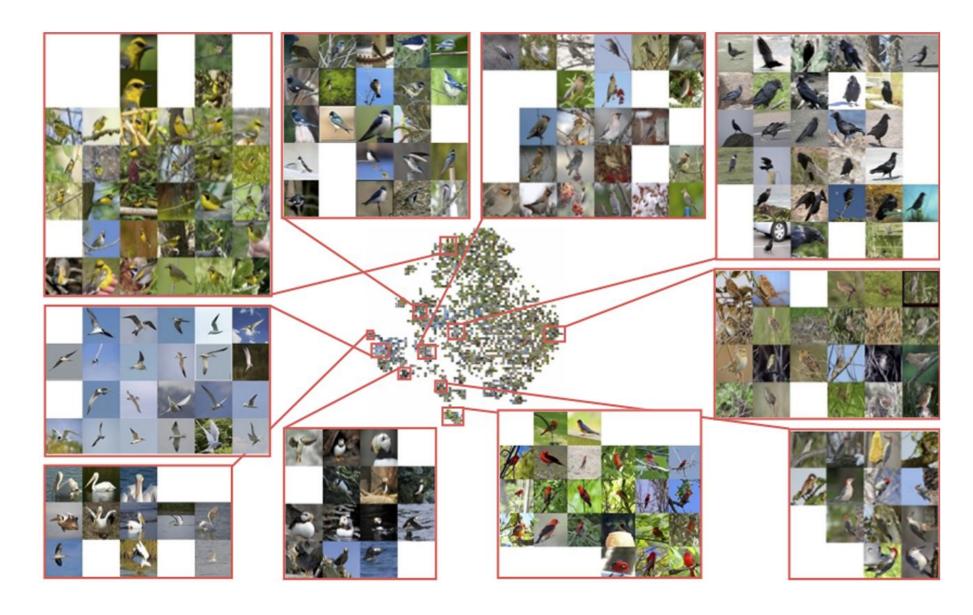
Problem 2: Objective Function ?



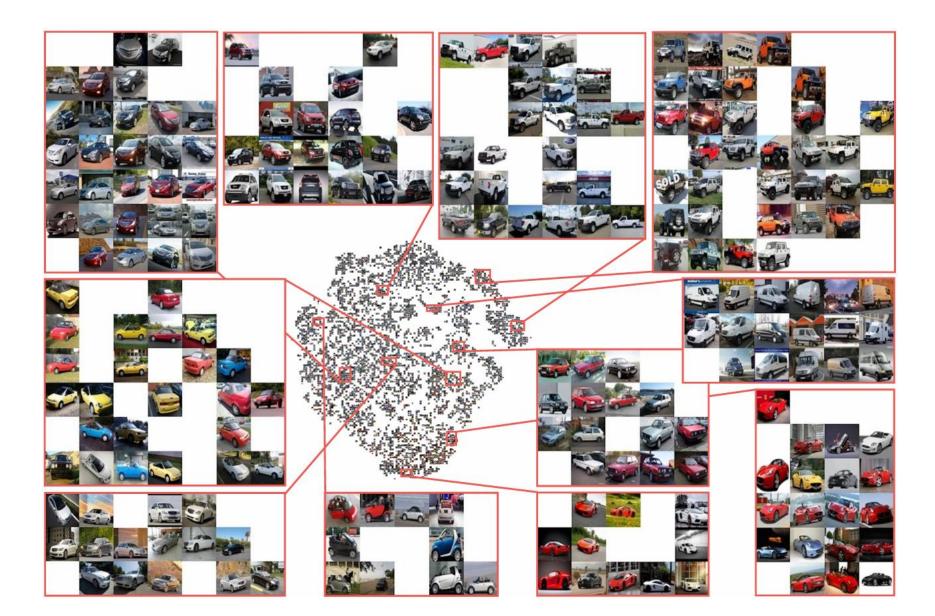
Examples of Contrastive Pairs



Examples of the Embedding Space



Examples of the Embedding Space

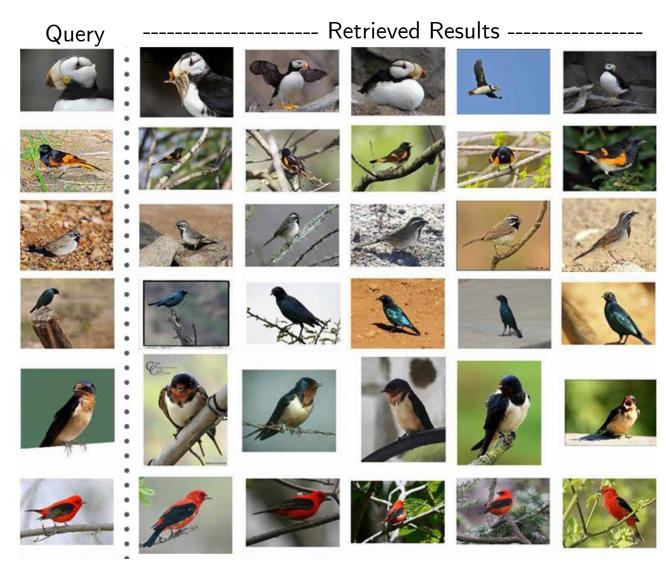


What can you do with this embedding space?

What can you do with this embedding space?

RETRIEVAL

- Given a query image (left column), find similar images
- All you have to do is find the nearest neighbors in the embedding space and return the results
- Embedding space now has a notion of "similarity"
 - o Similar datapoints are neighbors
 - Dissimilar datapoints are not



What can you do with this embedding space?

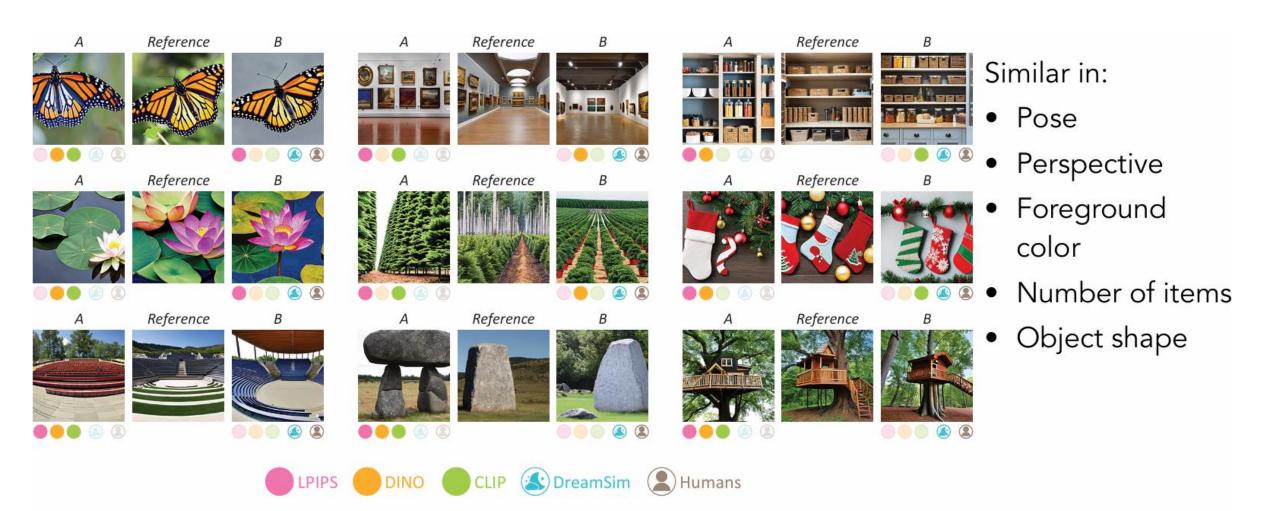
RETRIEVAL

- Given a query image (left column), find similar images
- All you have to do is find the nearest neighbors in the embedding space and return the results
- Embedding space now has a notion of "similarity"
 - o Similar datapoints are neighbors
 - o Dissimilar datapoints are not

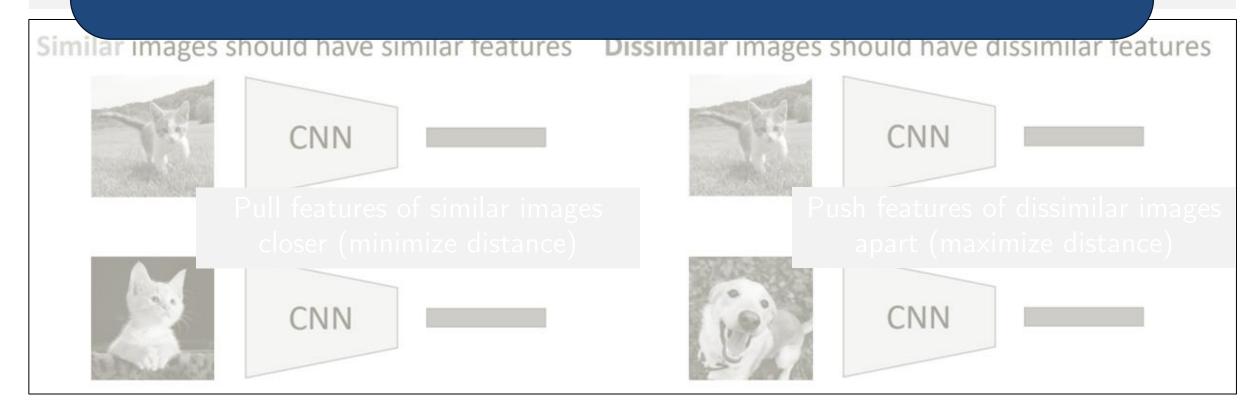


Figure 1: Example retrieval results on our *Online Products* dataset using the proposed embedding. The images in the first column are the query images.

Challenges: "Similarity" is hard ... What makes an image "similar"?



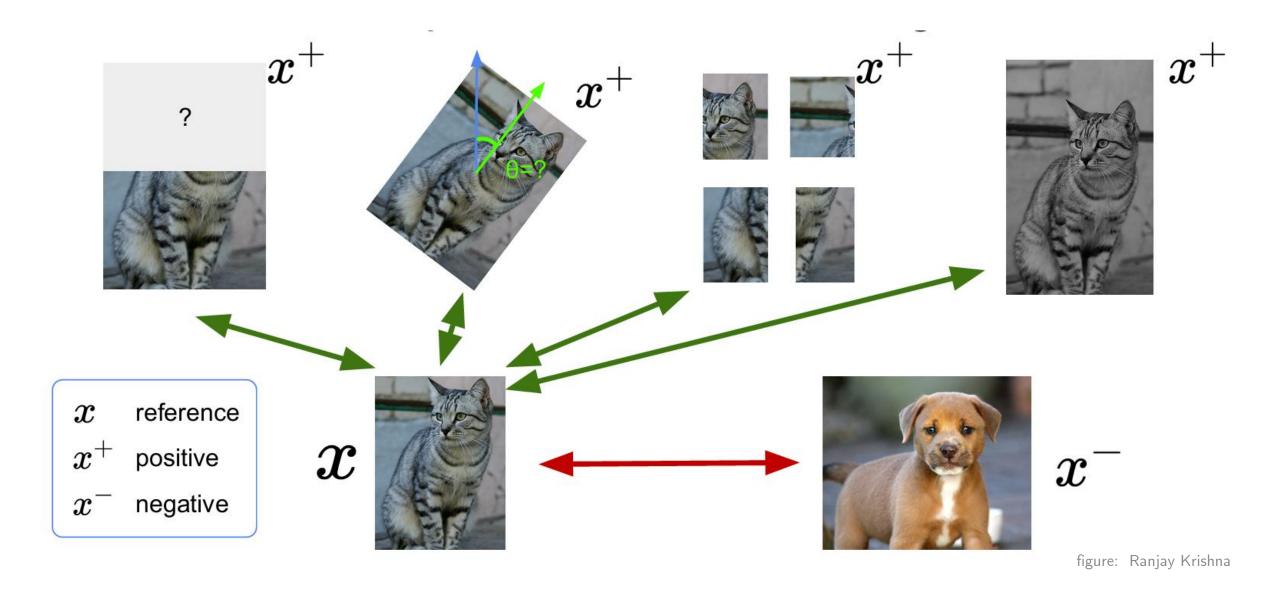
Where can we get pairs of similar and dissimilar images from?



Where can we get pairs of similar and dissimilar images from?

DATA AUGMENTATION

Contrastive Learning with Data Augmentation



Contrastive Learning Formulation

We want:

$$\operatorname{score}(f(x),f(x^+)) >> \operatorname{score}(f(x),f(x^-))$$

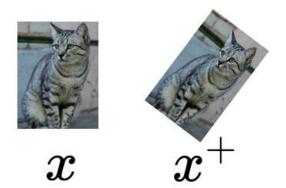
x: reference sample; x⁺ positive sample; x⁻ negative sample

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

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

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

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















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 score for the N-1 positive pair negative pairs

This seems familiar ...

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 score for the N-1 positive pair 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

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 score for the N-1 positive pair 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

Very similar to a softmax classifier

We want to compare the reference image against all other positive and negative images. We can exponentiate and normalize these scores like we did with the softmax classifier.

_

Contrastive Learning Loss

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

SimCLR: A Simple Framework for Contrastive learning

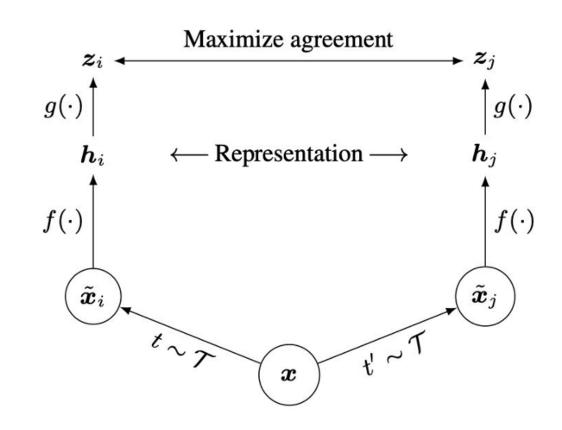
Cosine similarity as the score function:

$$s(u,v)=rac{u^Tv}{||u||||v||}$$

Use a projection network $h(\cdot)$ to project features to a space where contrastive learning is applied

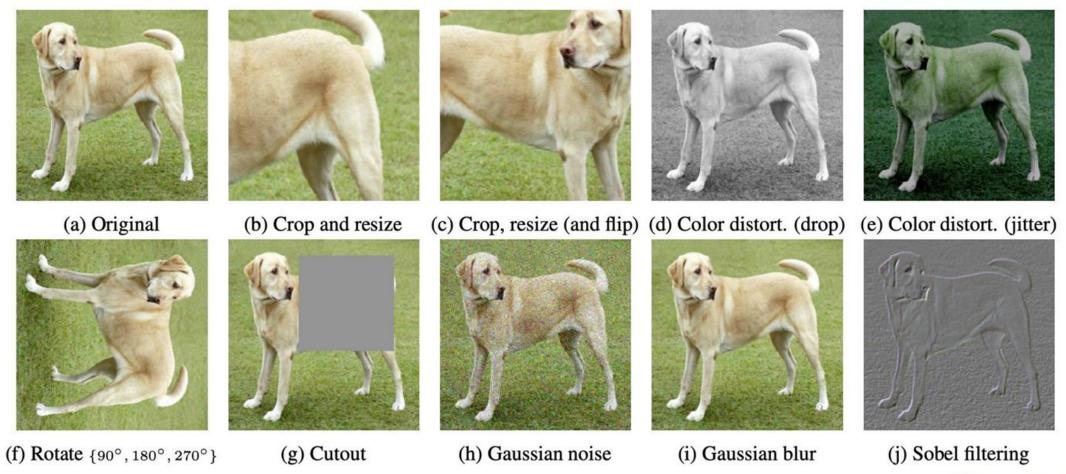
Generate positive samples through data augmentation:

 random cropping, random color distortion, and random blur.



Source: Chen et al., 2020

SimCLR: Data Augmentation Strategies



Source: Chen et al., 2020

SimCLR: Algorithm Sketch

Algorithm 1 SimCLR's main learning algorithm.

input: batch size N, constant τ , structure of f, g, \mathcal{T} . for sampled minibatch $\{x_k\}_{k=1}^N$ do

for all $k \in \{1, \ldots, N\}$ do

draw two augmentation functions $t \sim T$, $t' \sim T$

the first augmentation

Generate a positive pair by sampling data augmentation functions

the second augmentation

$$ilde{oldsymbol{x}}_{2k} = t'(oldsymbol{x}_k)$$

$$egin{aligned} oldsymbol{h}_{2k} &= f(ilde{oldsymbol{x}}_{2k}) \ oldsymbol{z}_{2k} &= g(oldsymbol{h}_{2k}) \end{aligned}$$

representation # projection

end for

for all $i \in \{1, \dots, 2N\}$ and $j \in \{1, \dots, 2N\}$ do

 $s_{i,j} = oldsymbol{z}_i^ op oldsymbol{z}_j/(\|oldsymbol{z}_i\|\|oldsymbol{z}_j\|)$ # pairwise similarity

end for

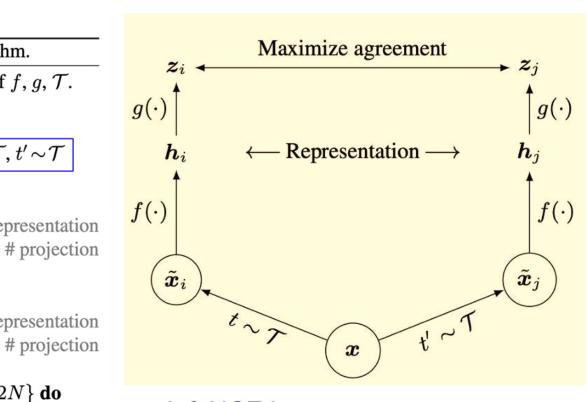
define $\ell(i,j)$ as $\ell(i,j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$

 $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^{N} \left[\ell(2k-1, 2k) + \ell(2k, 2k-1) \right]$

update networks f and g to minimize \mathcal{L}

end for

return encoder network $f(\cdot)$, and throw away $g(\cdot)$



InfoNCE loss: Use all non-positive samples in the batch as x⁻

Source: Chen et al., 2020

Iterate through and use each of the 2N sample as reference, compute average loss

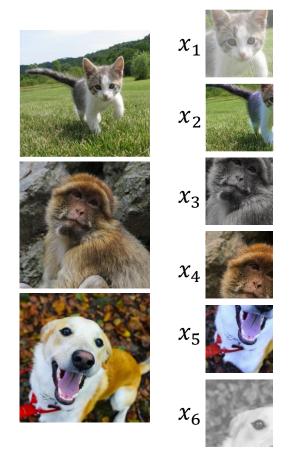
Batch of N images

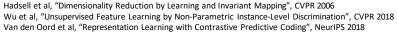




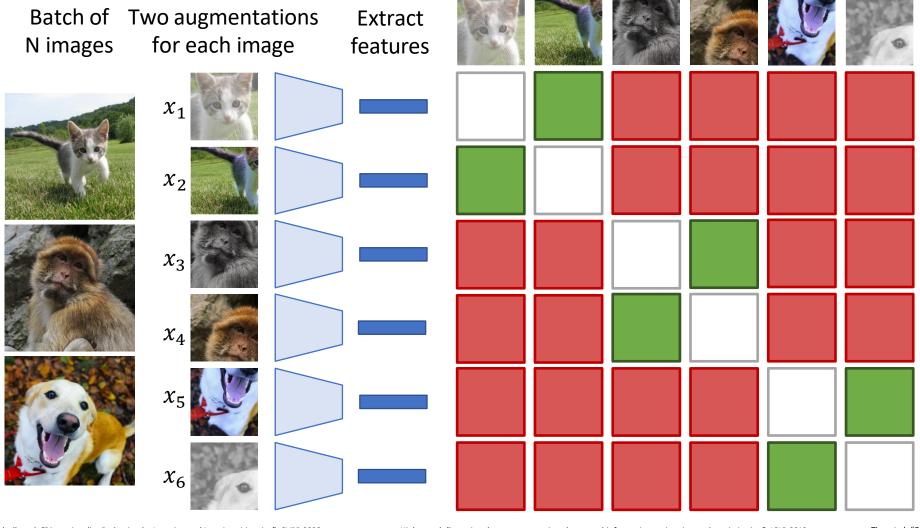


Batch of Two augmentations N images for each image





Batch of Two augmentations **Extract** N images for each image 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 :

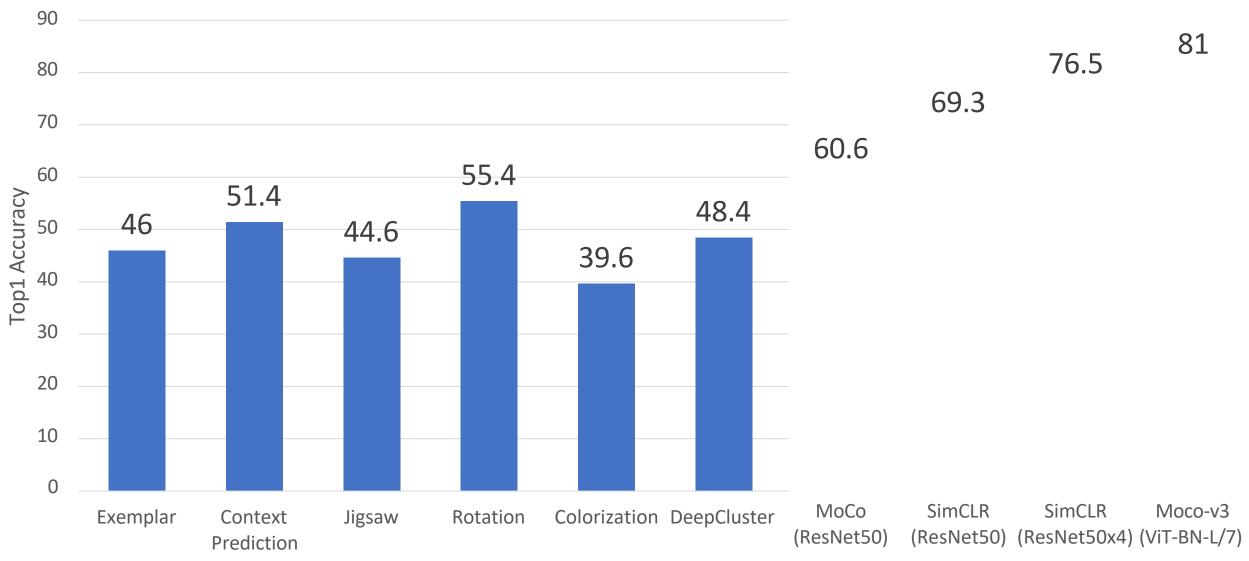
$$s_{i,j} = \frac{\phi(x_i)^T \phi(x_j)}{\|\phi(x_i)\| \cdot \|\phi(x_i)\|}$$

If (x_i, x_j) is a positive pair, then loss for x_i is:

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

Interpretation: Cross-entropy loss over the other 2N-1 elements in the batch!

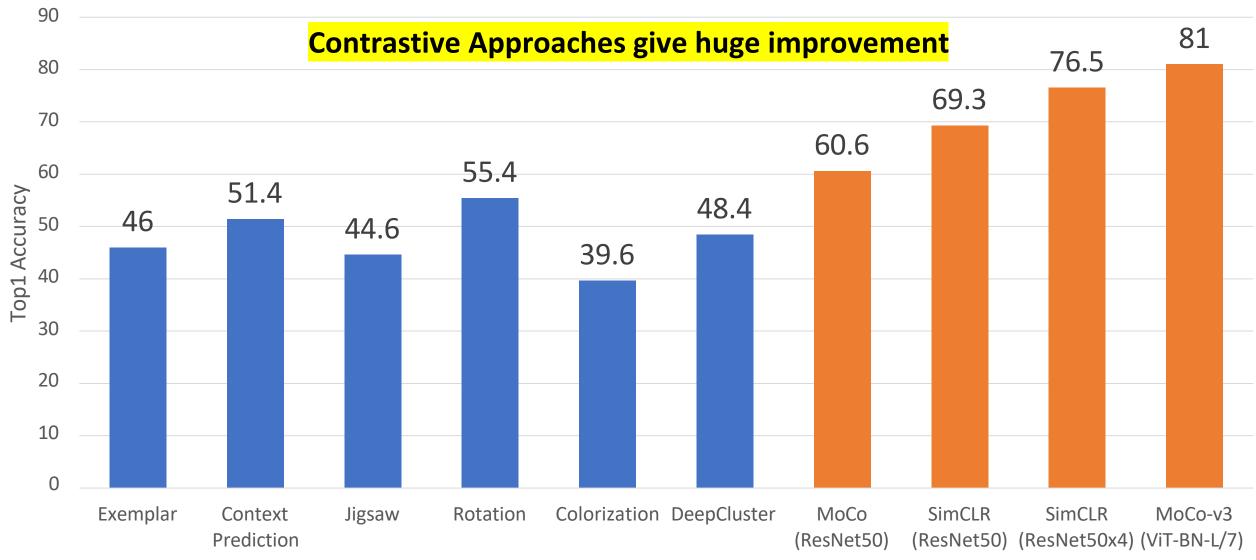
ImageNet Linear Classification from SSL Features



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

(Lots of caveats here ... different architectures, etc)

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

(Lots of caveats here ... different architectures, etc)

But how did you get the pretraining data?

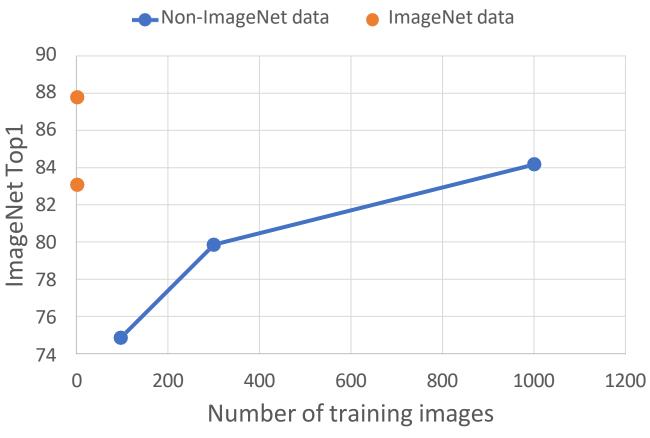
The motivation of SSL is scaling to large data that can't be labeled

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", arXiv 2021

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 form 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

RedCaps: Images and Captions from Reddit











r/birdpics: male r/crafts: my mom northern cardinal

tied this mouse

r/itookapicture: r/perfectfit: this itap of the taj mahal lemon in my drink

r/shiba: mlem!

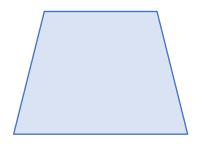
Data from 350 manually-chosen subreddits 12M high-quality (image, caption) pairs

For now: Assume you can learn language representations

(I teach this in much detail on CMSC 475/675 Neural Networks)

Computer Vision

Image Features: H x W x C

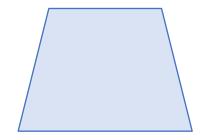




Input Image

Natural Language Processing

Word Features L x C



A white and gray cat standing outside on the grass

Input Sentence (L words)

Contrastive Learning with Vision-Language Data

OpenAl

January 5, 2021 Milestone

CLIP: Connecting text and images

Contrastive Learning with Vision-Language Data

Learning Transferable Visual Models From Natural Language Supervision

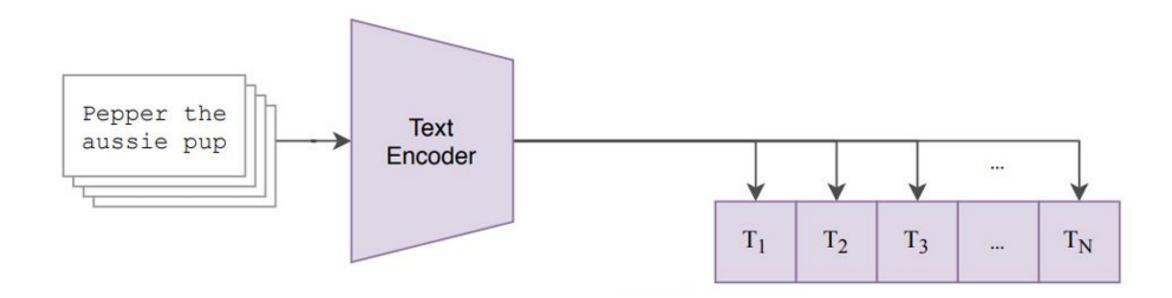
```
Alec Radford * 1 Jong Wook Kim * 1 Chris Hallacy 1 Aditya Ramesh 1 Gabriel Goh 1 Sandhini Agarwal 1 Girish Sastry 1 Amanda Askell 1 Pamela Mishkin 1 Jack Clark 1 Gretchen Krueger 1 Ilya Sutskever 1
```

ICML | 2021

Thirty-eighth International Conference on Machine Learning

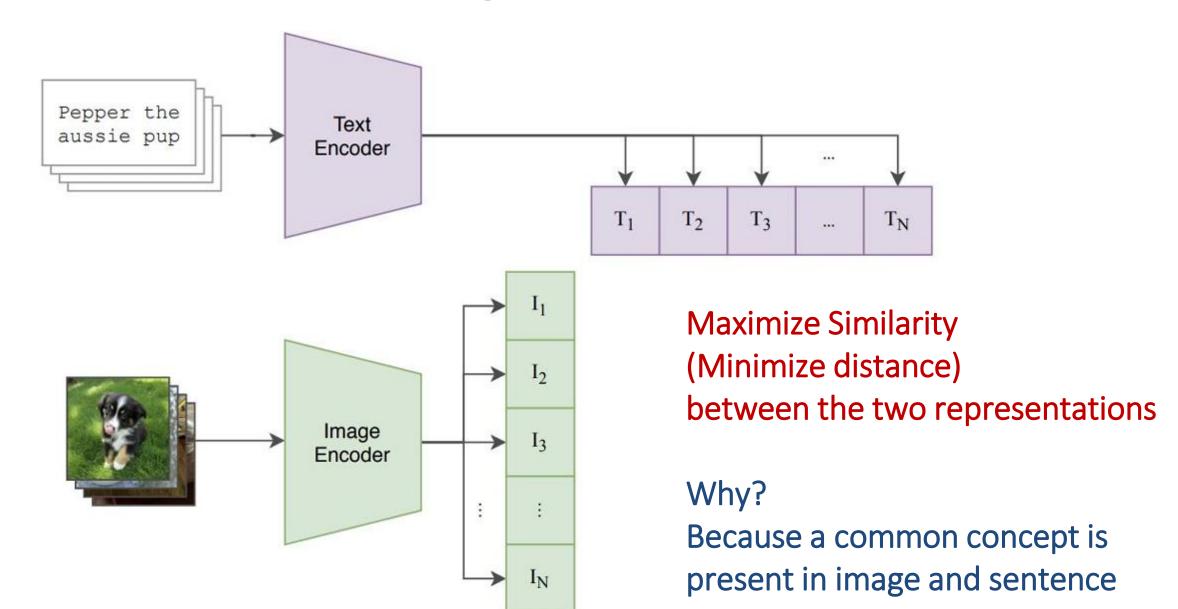
Pepper the aussie pup



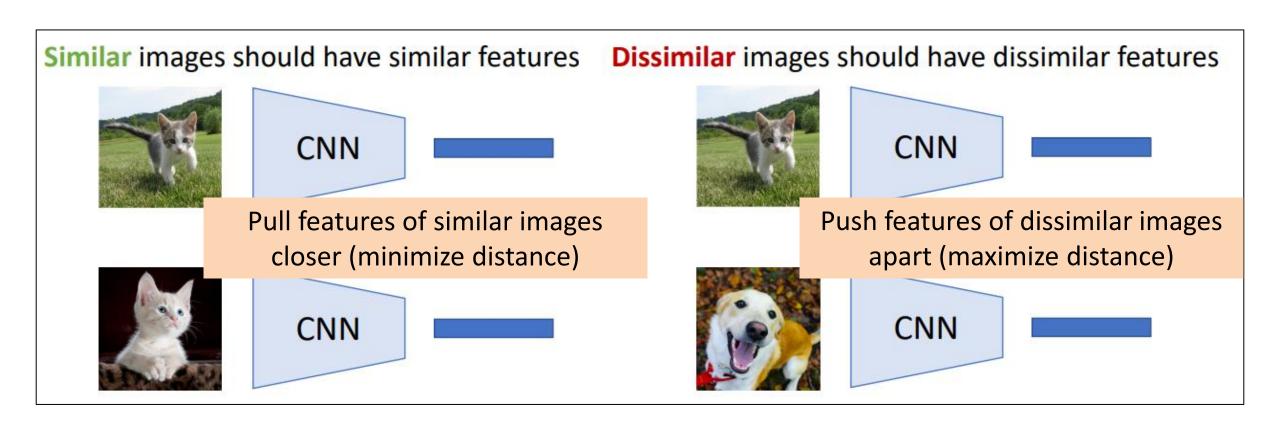




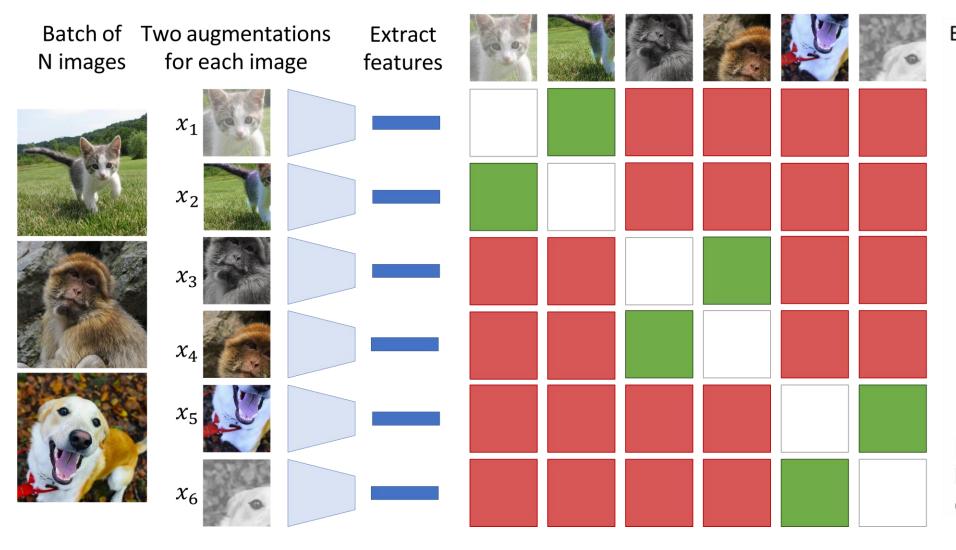
Goal: Do Contrastive Learning!



Recall: Contrastive Learning (General Form)



Recall: Contrastive Learning (SimCLR)



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 :

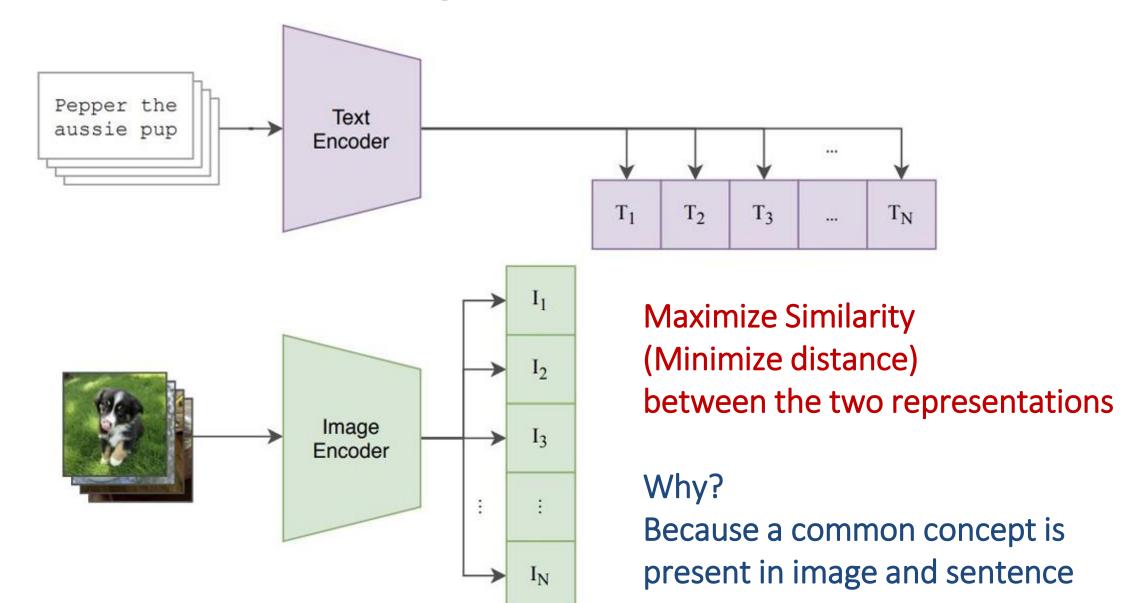
$$s_{i,j} = \frac{\phi(x_i)^T \phi(x_j)}{\|\phi(x_i)\| \cdot \|\phi(x_i)\|}$$

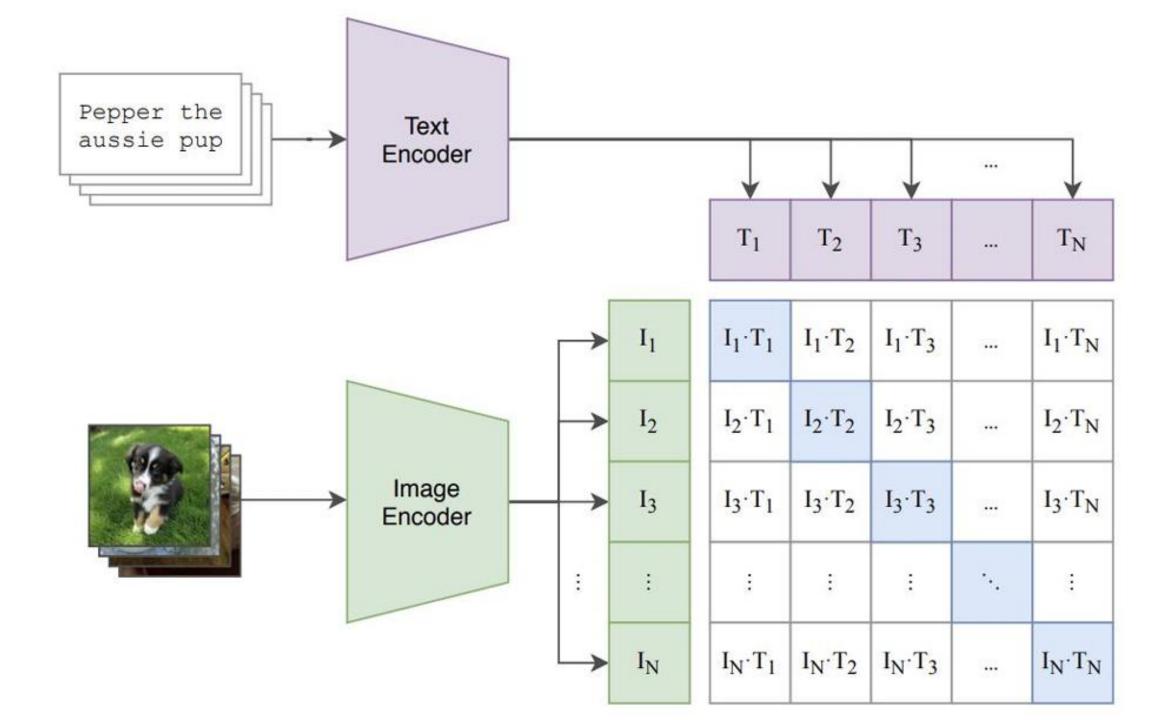
If (x_i, x_j) is a positive pair, then loss for x_i is:

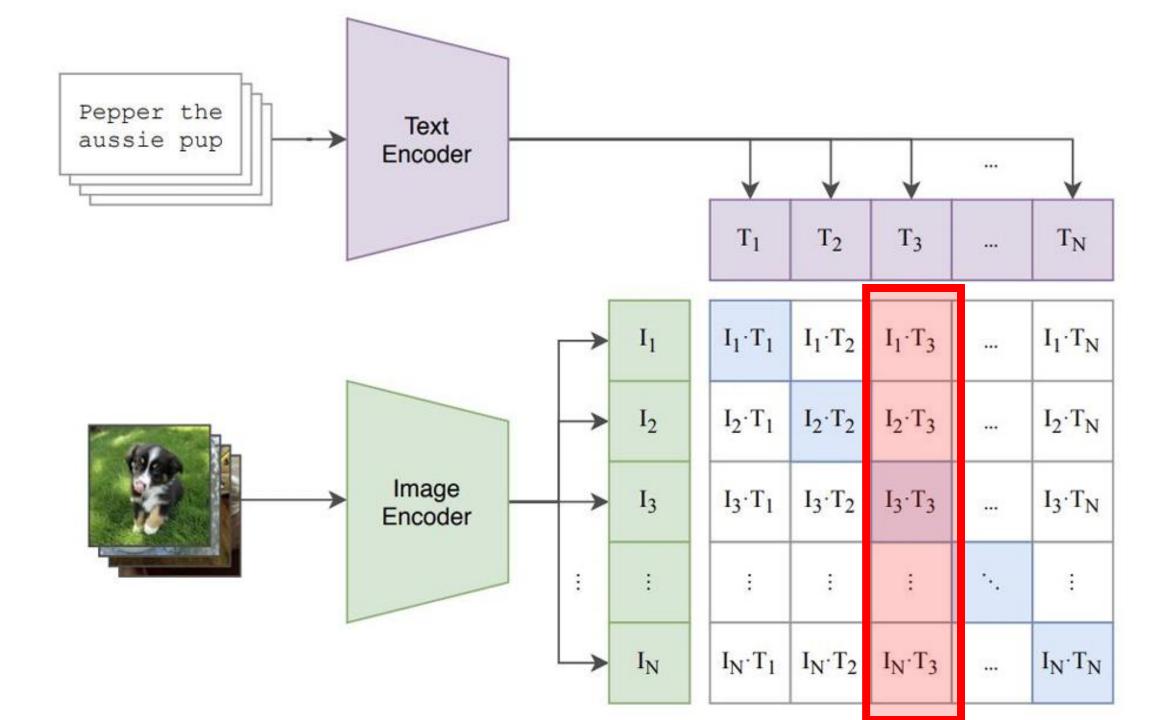
$$L_{i} = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{\substack{k=1\\k\neq i}}^{2N} \exp(s_{i,k}/\tau)}$$
(τ is a temperature)

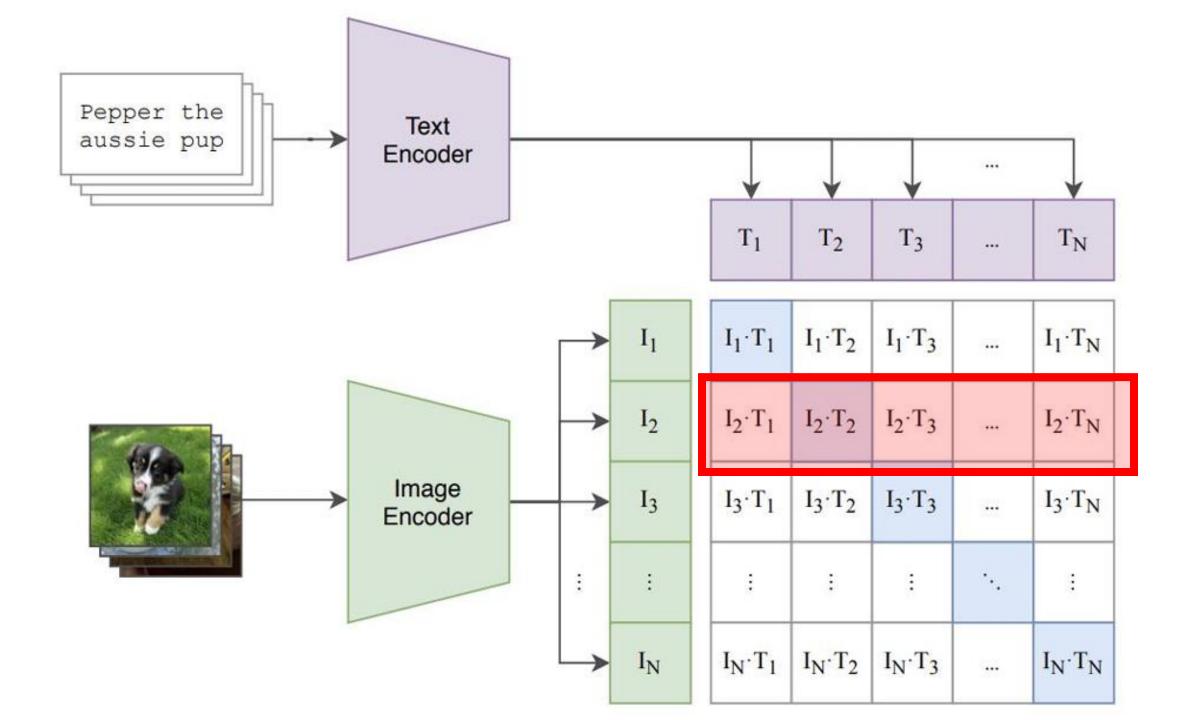
Interpretation: Cross-entropy loss over the other 2N-1 elements in the batch!

Goal: Do Contrastive Learning!









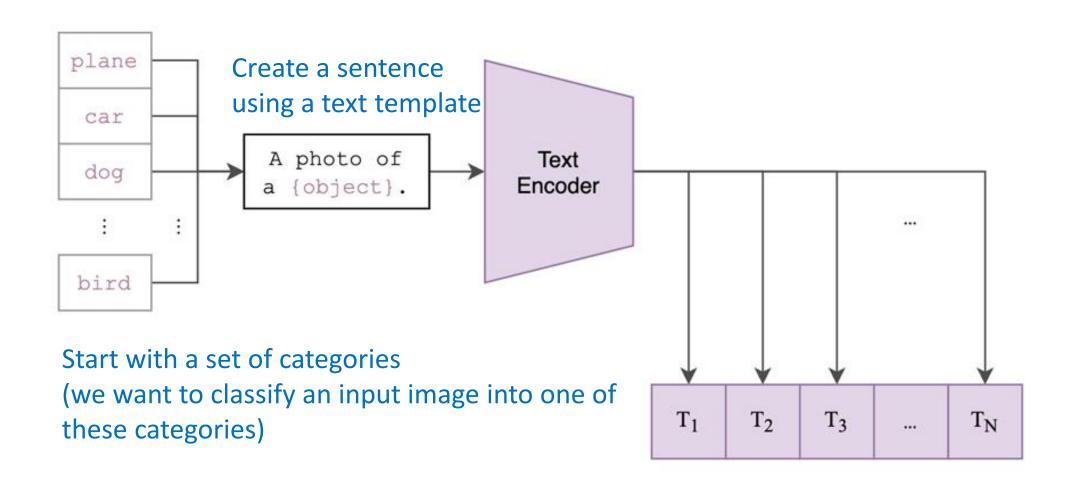
plane car dog :

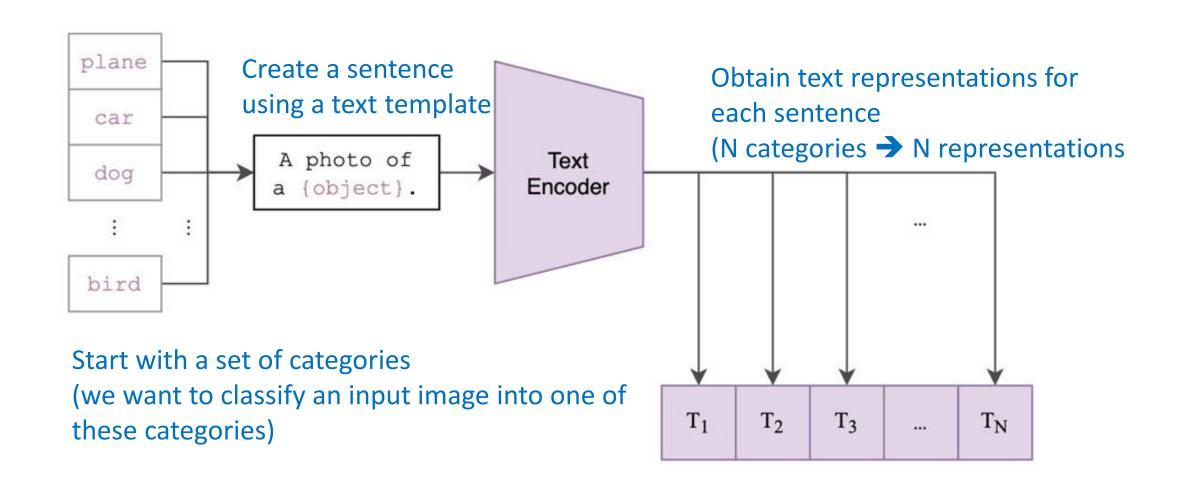
bird

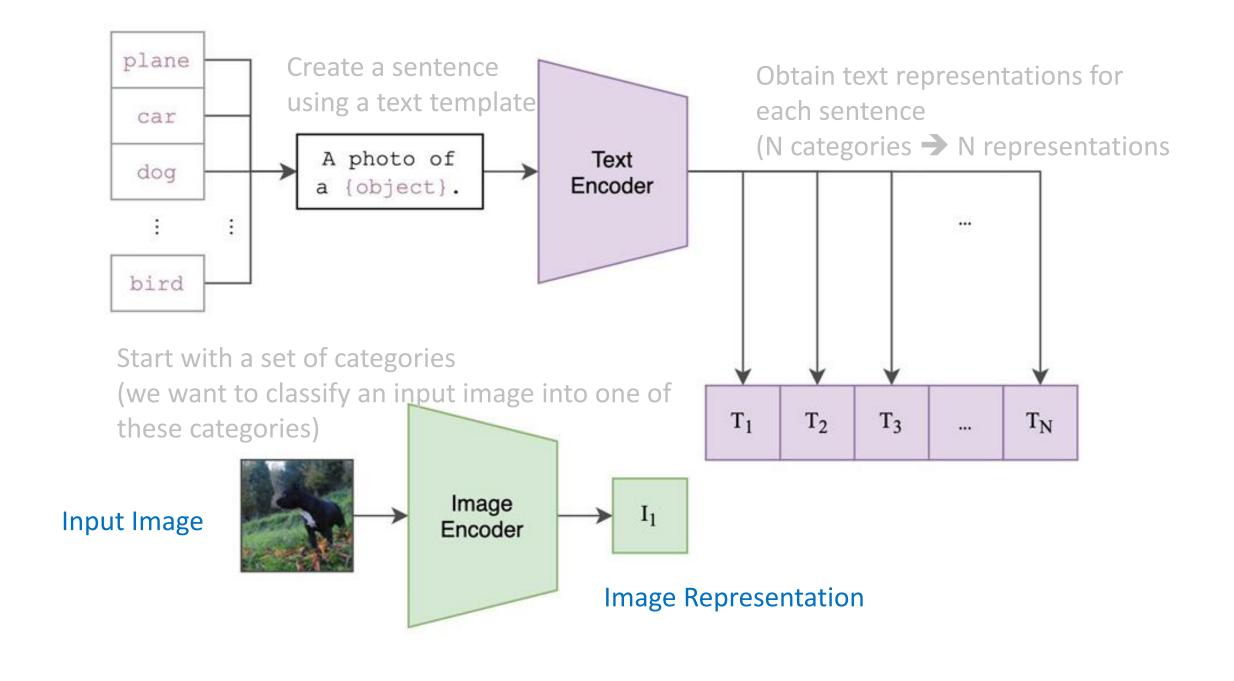
Start with a set of categories (we want to classify an input image into one of these categories)

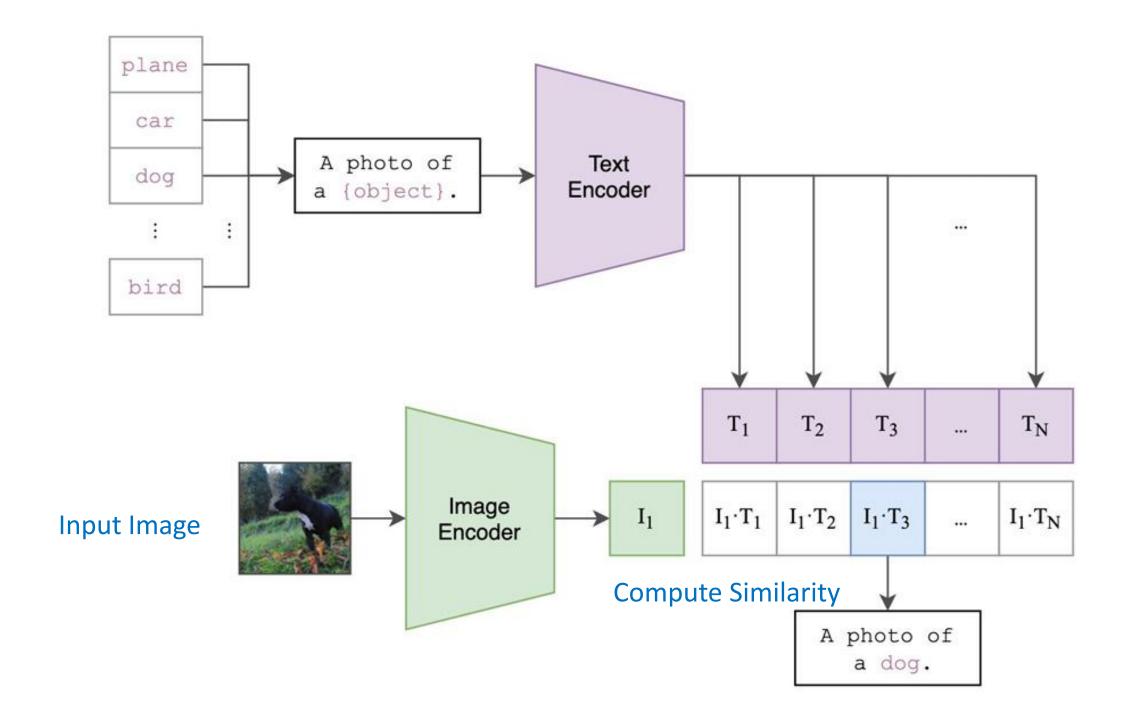
Image Classification with CLIP

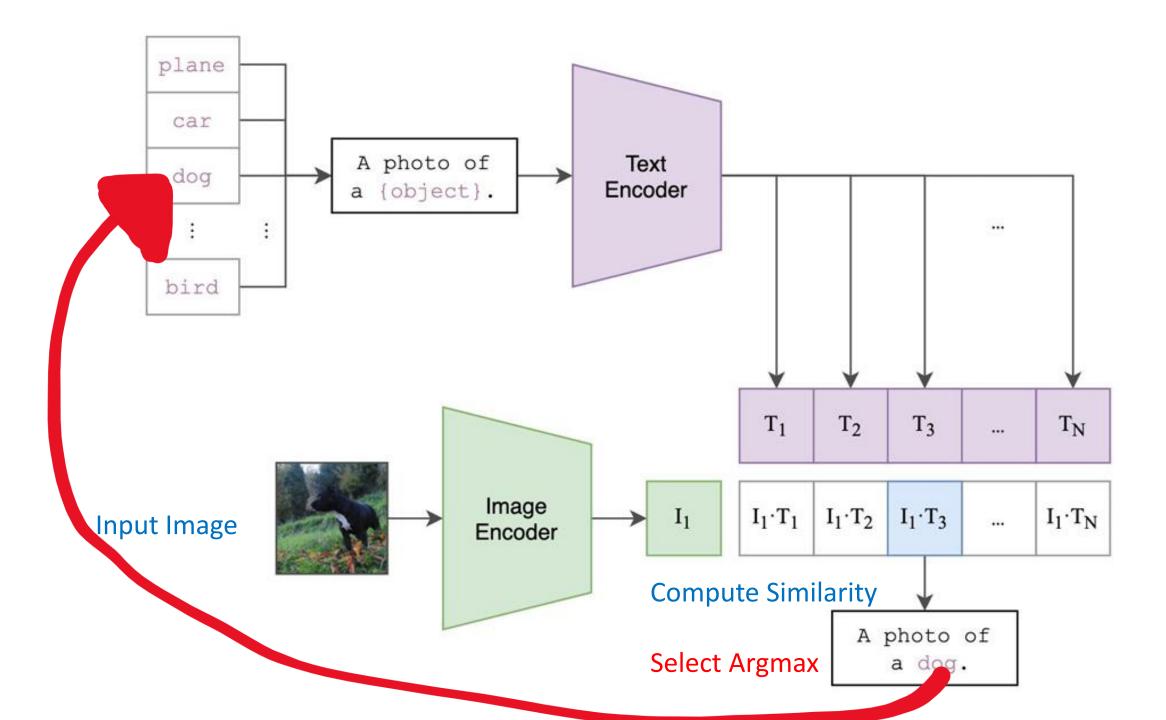
(so called "Zero-shot" classification)



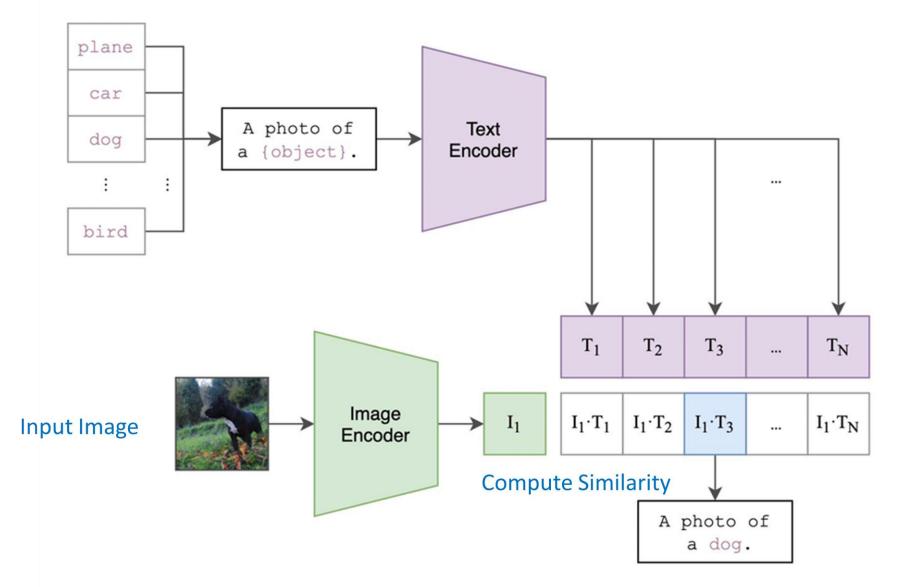








CLIP: Zero-Shot Classification



Language enables zero shot classification:

Classify images into categories without any additional training data!

Contrastive loss:

For each image, predict which sentence matches it.

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

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

CLIP Performance

Very strong performance on many downstream vision problems!

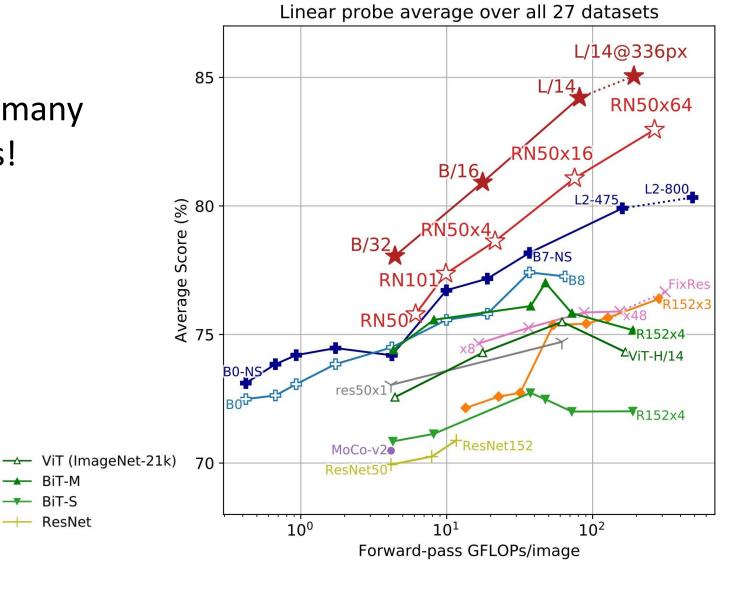
Performance continues to improve with larger models

CLIP-ViT

CLIP-ResNet

EfficientNet

EfficientNet-NoisyStudent



Instagram-pretrained

BiT-M

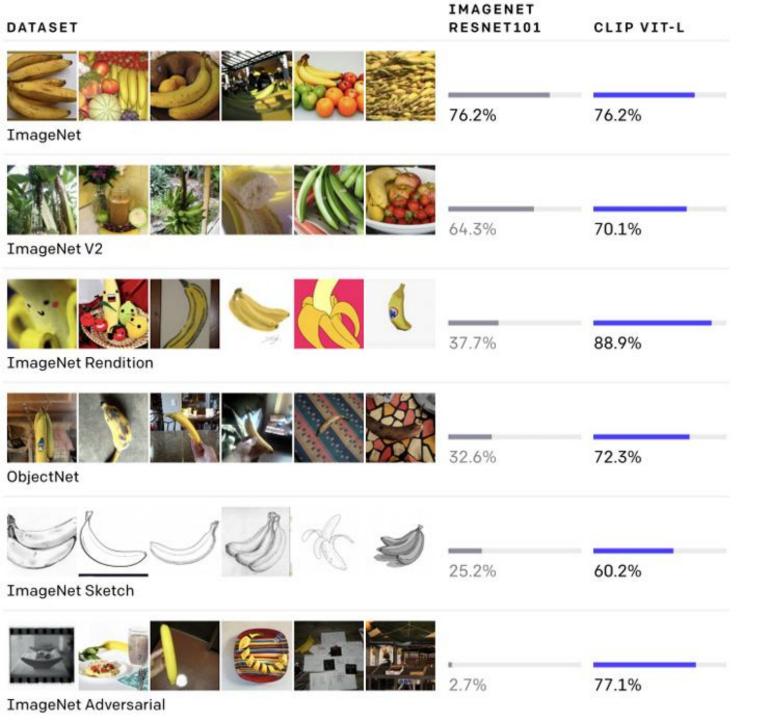
BiT-S

ResNet

SimCLRv2

BYOL

─ MoCo



CLIP Details

Training Details:

- Trained on 400M image-text pairs from the internet (i.e. without permissions a.k.a. stealing)
- Batch size of 32,768
- 32 epochs over the dataset
- Cosine learning rate decay

Architecture

- ResNet-based or ViT-based image encoder
- Transformer-based text encoder

Caltech-101

kangaroo (99.8%) Ranked 1 out of 102 labels



- ✓ a photo of a kangaroo.
- x a photo of a gerenuk.
- x a photo of a emu.
- × a photo of a wild cat.
- x a photo of a scorpion.

ImageNet-R (Rendition)

Siberian Husky (76.0%) Ranked 1 out of 200 labels



- ✓ a photo of a siberian husky.
- X a photo of a german shepherd dog.
- x a photo of a collie.
- × a photo of a border collie.
- × a photo of a rottweiler.

Oxford-IIIT Pets Maine Coon (100.0%) Ranked 1 out of 37 labels



- ✓ a photo of a maine coon, a type of pet.
- × a photo of a persian, a type of pet.
- x a photo of a ragdoll, a type of pet.
- X a photo of a birman, a type of pet.
- X a photo of a siamese, a type of pet.

CIFAR-100 snake (38.0%) Ranked 1 out of 100 labels



- × a photo of a sweet pepper.
- x a photo of a flatfish.
- x a photo of a turtle.
- × a photo of a lizard.

Country211

Belize (22.5%) Ranked 5 out of 211 labels



x a photo i took in french guiana.

x a photo i took in gabon.

x a photo i took in cambodia.

× a photo i took in guyana.

RESISC45

roundabout (96.4%) Ranked 1 out of 45 labels



- ✓ satellite imagery of roundabout.
- × satellite imagery of intersection.
- × satellite imagery of church.
- × satellite imagery of medium residential.
- × satellite imagery of chaparral.

Stanford Cars

2012 Honda Accord Coupe (63.3%) Ranked 1 out of 196 labels



- × a photo of a 2012 honda accord sedan.
- x a photo of a 2012 acura ti sedan.
- × a photo of a 2012 acura tsx sedan.
- × a photo of a 2008 acura tl type-s.

SUN

kennel indoor (98.6%) Ranked 1 out of 723 labels



- × a photo of a kennel outdoor.
- × a photo of a jail cell.
- × a photo of a jail indoor.
- × a photo of a veterinarians office.

Can be "Attacked"

(if there is time, we will discuss Adv Attacks in this course. If not, I teach this in NN.)





Target class: pizza Attack text: pizza

B	pizza	83.7%
pizza pizza	pretzel	2%
pizza	Chihuahua	1.5%
pizza	broccoli	1.2%
pizza	hot dog	0.6%
pizza	Boston Terrier	0.6%
plzza	French Bulldog	0.5%
	spatula	0.4%
	Italian Greyhound	0.3%

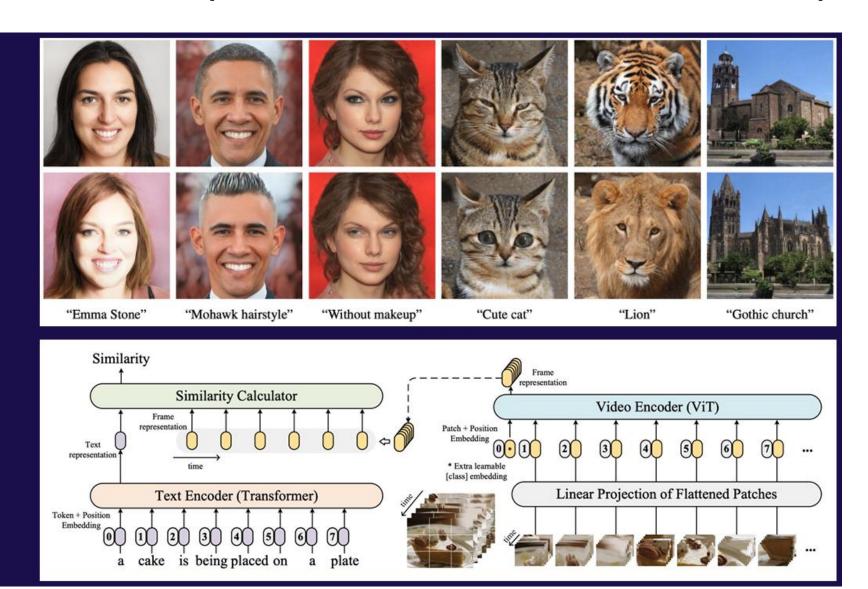
Applications of CLIP (slide from Radford et al.)

StyleCLIP (Patashnik et al.)

Steering a GAN Using CLIP

CLIP4Clip (Luo & Ji, et al.)

Video retrieval using CLIP features



Summary

- Self-Supervised Learning: scale up training without human annotation
 - First train for a pretext task, then transfer to downstream tasks
 - Many pretext tasks: context prediction, jigsaw, colorization, clustering, rotation
 - SSL has been wildly successful for language
- Intense research on SSL in vision
- Multimodal SSL with vision + language has been very successful