

So Far

Few-shot learning via meta-learning

Problem: Given data from $\mathcal{T}_1, \dots, \mathcal{T}_n$, solve new task $\mathcal{T}_{\text{test}}$ more quickly / proficiently / stably

Methods: black-box, optimization-based, non-parametric

What if you don't have a lot of tasks?

What if you *only* have one batch of unlabeled data?

This Lecture

Unsupervised representation learning for few-shot learning

Part I: Contrastive learning

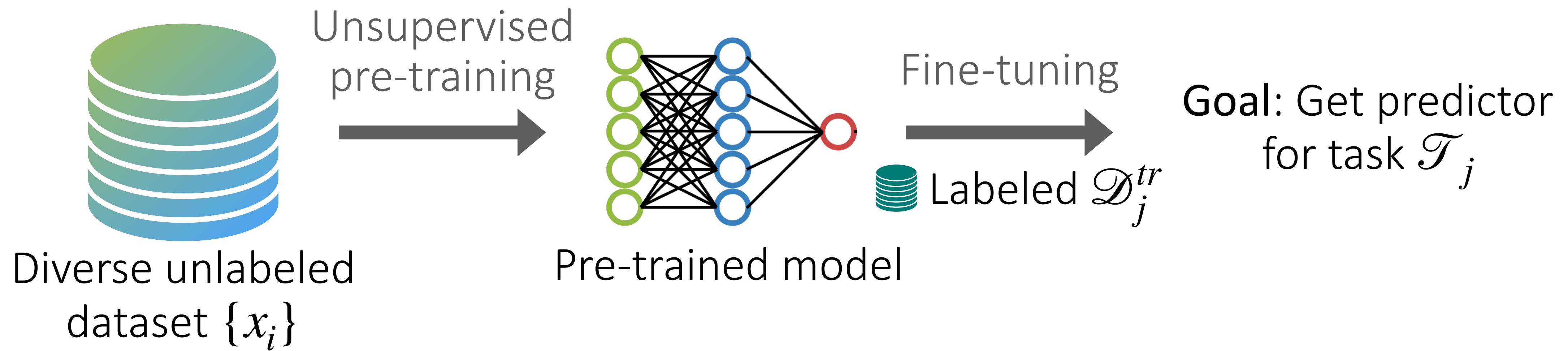
Part II (next time): Reconstruction-based methods

Relation to meta-learning.

Goals for the lecture:

- Understand contrastive learning: intuition, design choices, how to implement
- How contrastive learning relates to meta-learning

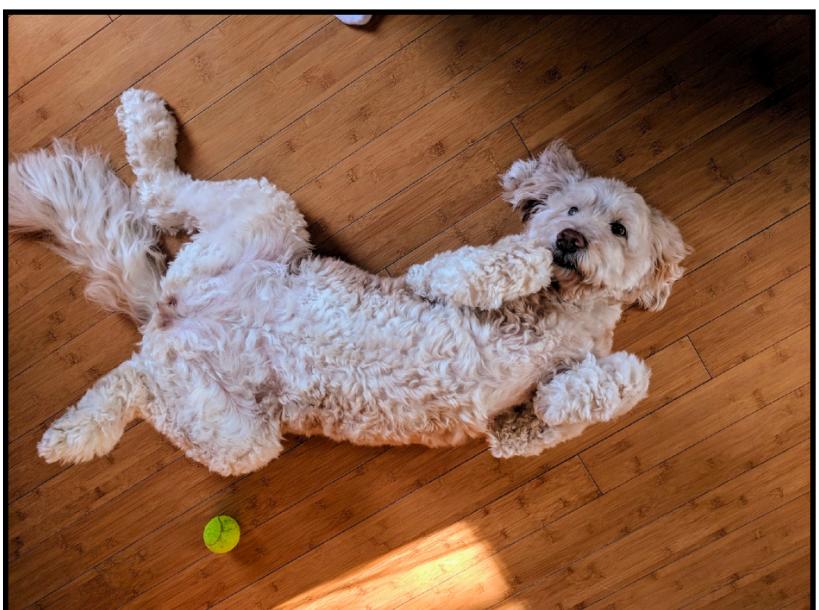
Unsupervised Pre-Training Set-Up



Key Idea of Contrastive Learning

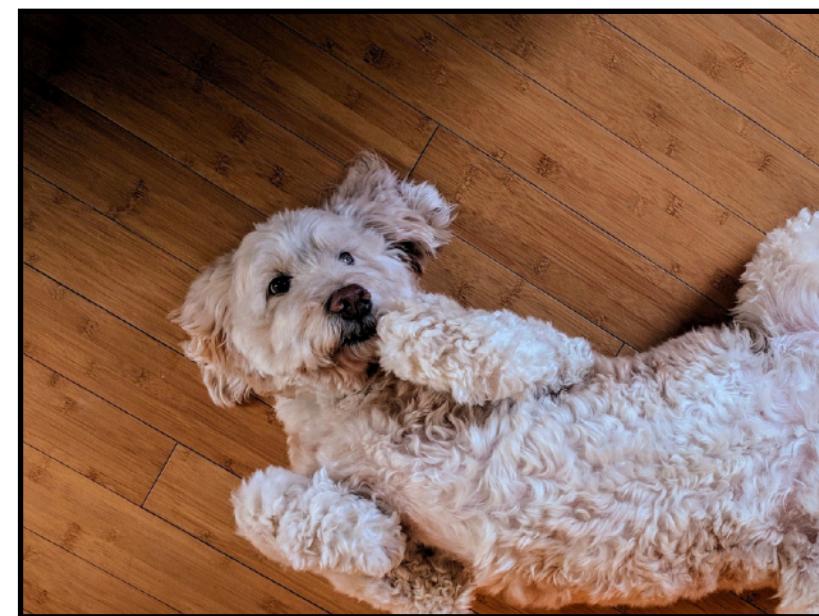
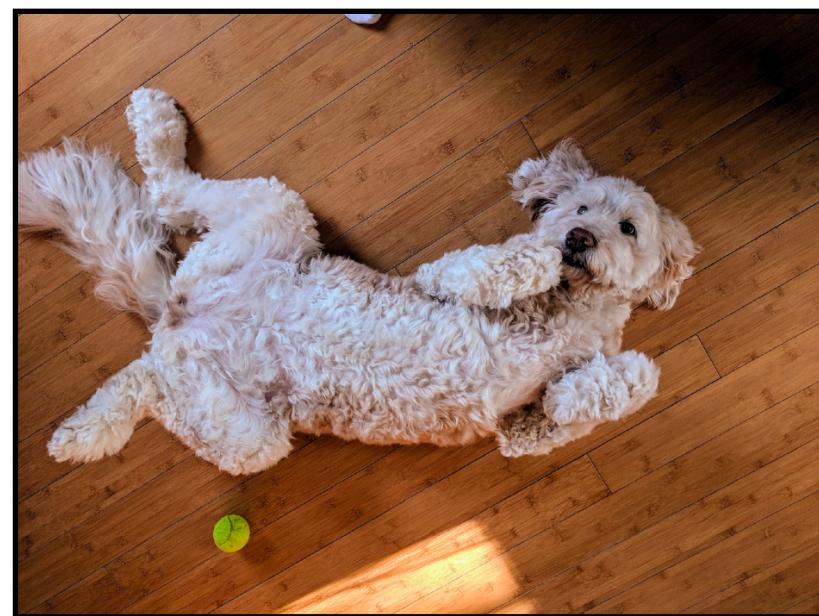
Similar examples should have similar representations

Examples with the same class label



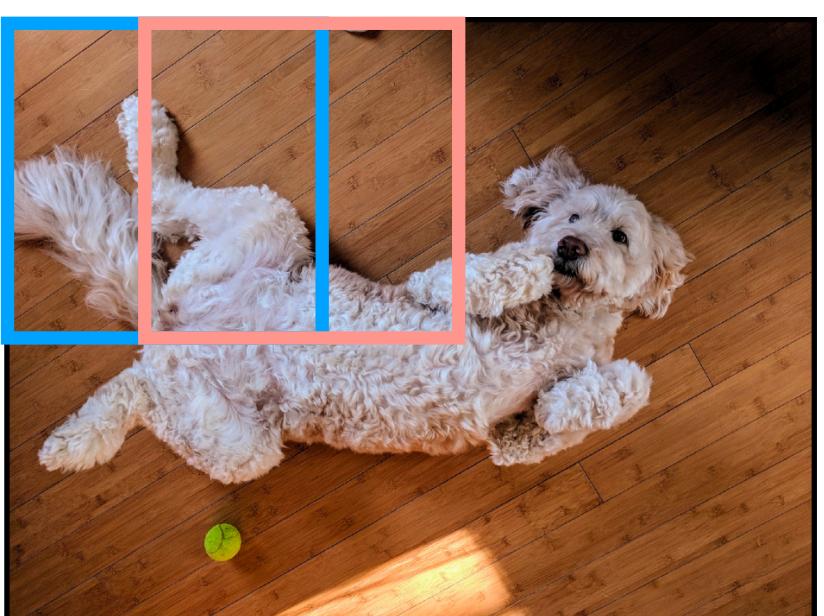
(Requires labels, related to Siamese nets, ProtoNets)

Augmented versions of the example



(flip & crop)

Nearby image patches



Nearby video frames



van den Oord, Li, Vinyals. CPC. 2018

Chen, Kornblith, Norouzi, Hinton. SimCLR. ICML 2020

Key Idea of Contrastive Learning

Similar examples should have similar representations



Similar representations



Similar representations

Question: Why not simply minimize difference between representations?

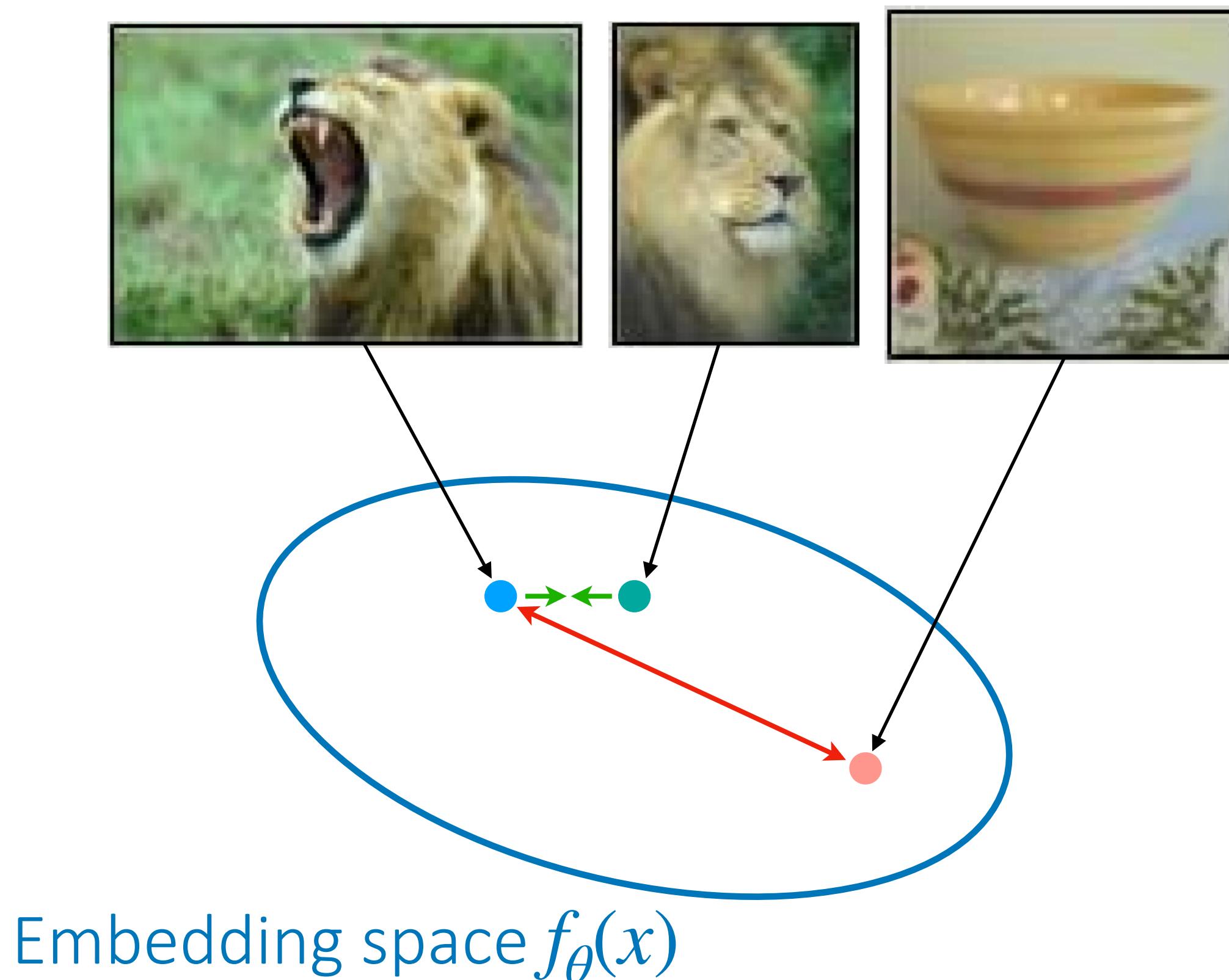
$$\min_{\theta} \sum_{(x_i, x_j)} \|f_{\theta}(x_i) - f_{\theta}(x_j)\|^2$$

Need to both compare & *contrast*!

Key Idea of Contrastive Learning

Similar examples should have **similar representations**

Need to both compare & *contrast*!



Bring together representations of similar examples.

Push apart representations of differing examples.

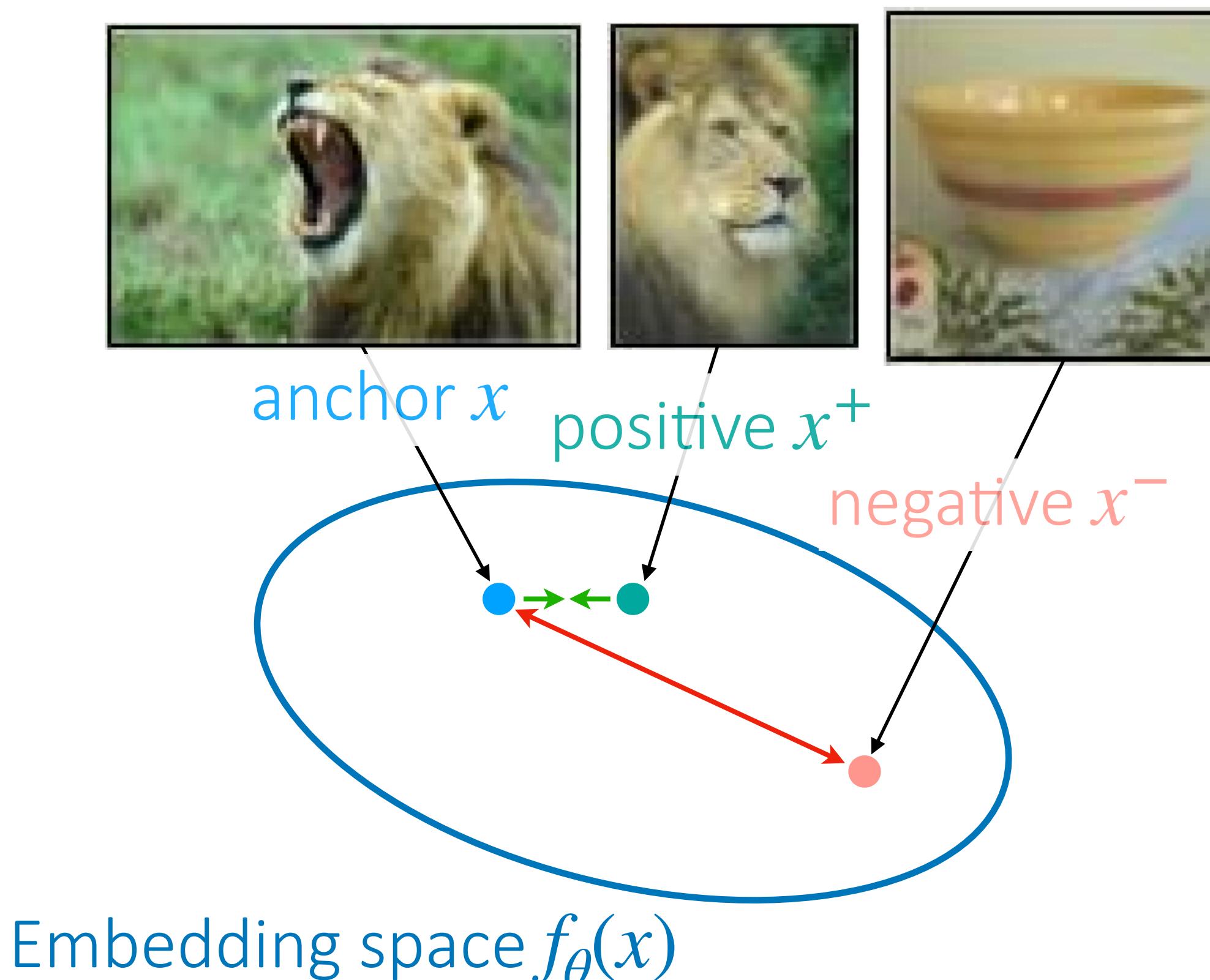
Key design choices:

1. Implementation of contrastive loss
2. Choosing what to compare/contrast

Contrastive Learning Implementation

Similar examples should have **similar representations**

Need to both compare & *contrast*!



V1. Triplet loss:

$$\min_{\theta} \sum_{(x,x^+,x^-)} \max(0, \|f_\theta(\textcolor{blue}{x}) - f_\theta(\textcolor{teal}{x}^+)\|^2 - \|f_\theta(\textcolor{blue}{x}) - f_\theta(\textcolor{red}{x}^-)\|^2 + \epsilon)$$

Compare to Siamese networks:

Classify (x, x') as same class if $\|f(x) - f(x')\|^2$ is small.

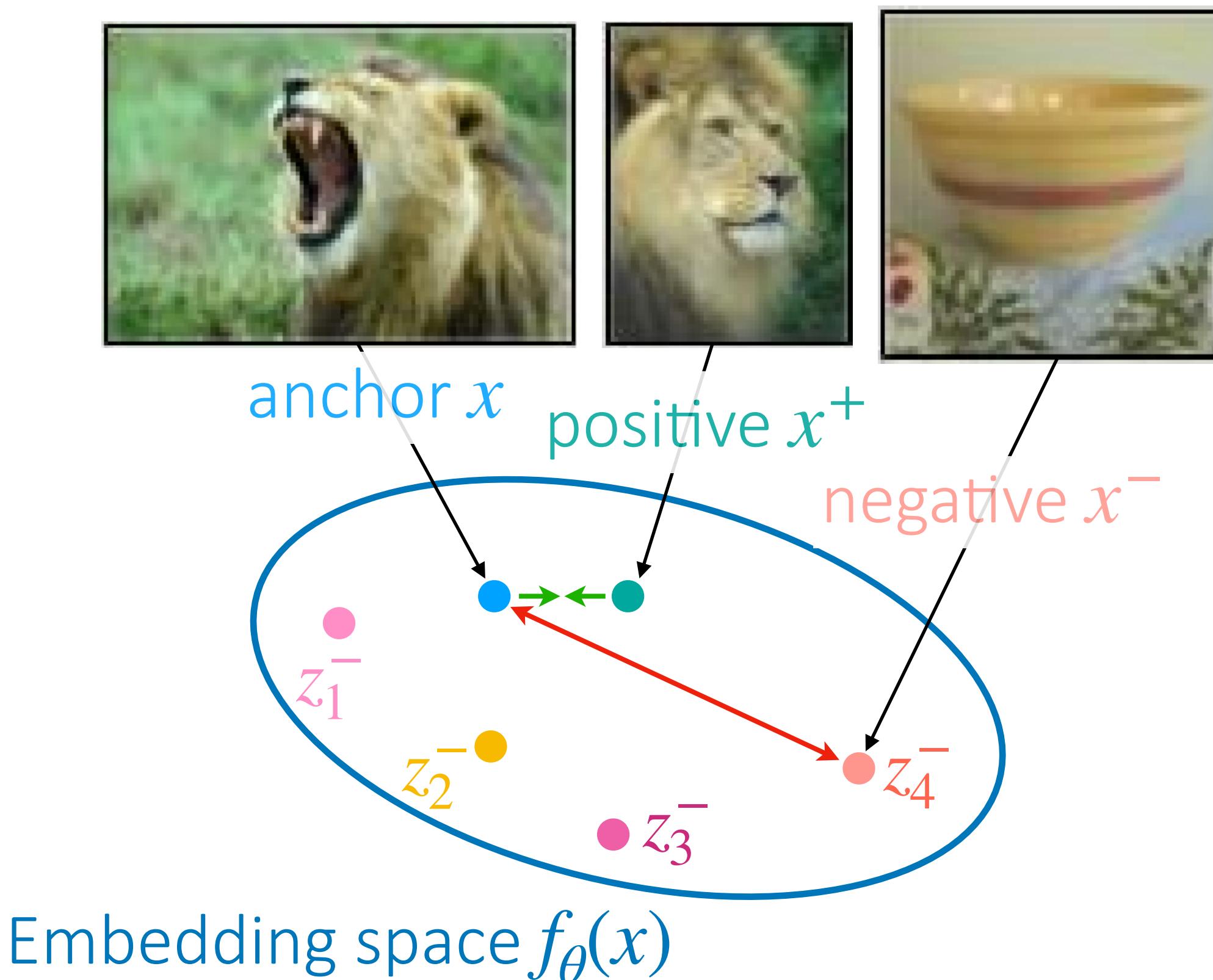
Key difference: learns a metric space, not just a classifier

Challenge: need to find difficult negatives.

Contrastive Learning Implementation

Similar examples should have **similar representations**

Need to both compare & *contrast*!



V2. From binary to N-way classification:

$$\mathcal{L}_{\text{N-way}}(\theta) = - \sum_z \log \frac{\exp(-d(z, z^+))}{\sum_i \exp(-d(z, z_i^-))}$$

- generalization of triplet loss to multiple negatives

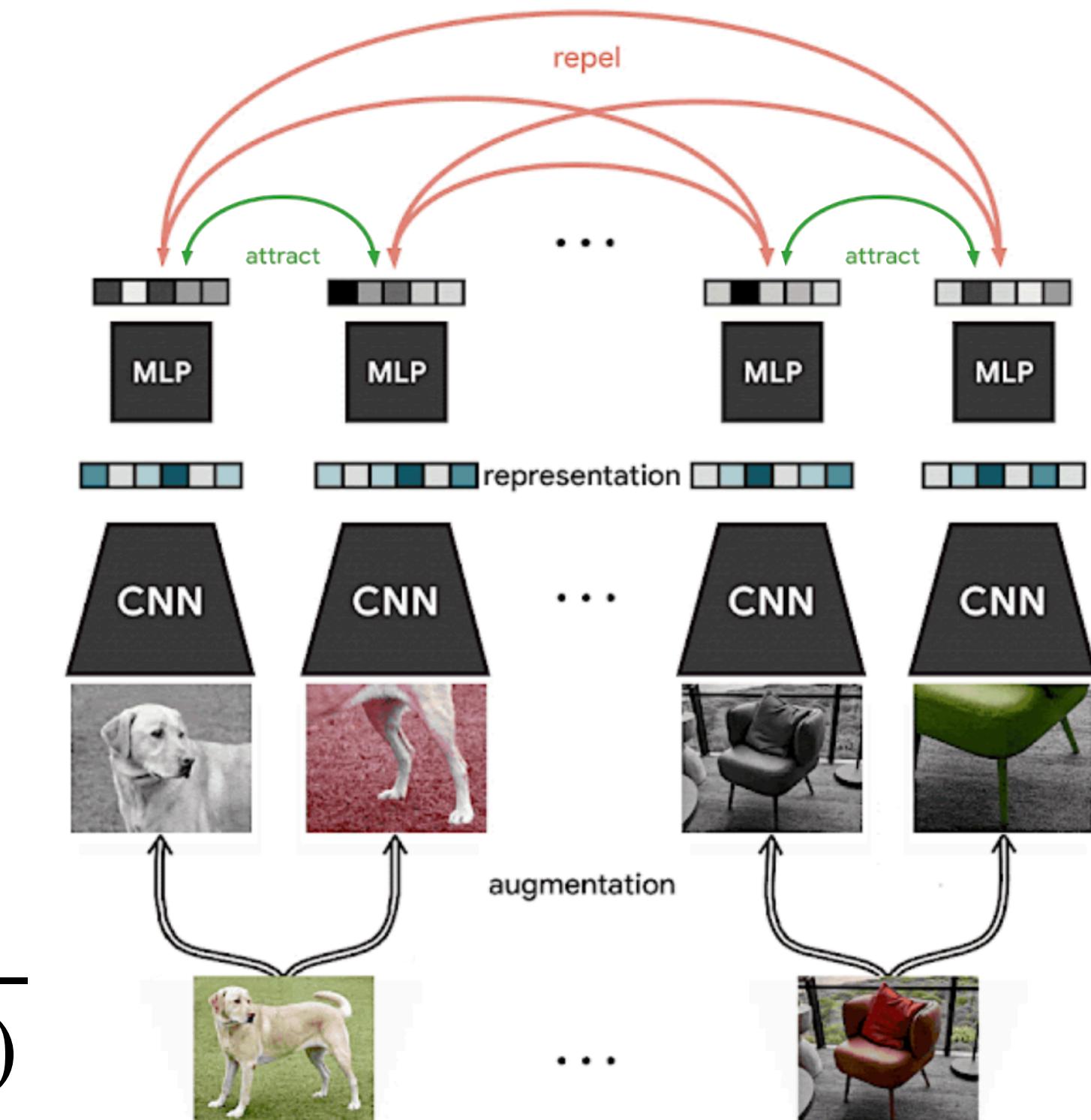
Sohn. N-Pair Loss Objective. NIPS 2016

Contrastive Learning Implementation

SimCLR Algorithm

Unsupervised Pre-Training

1. Sample minibatch of examples x_1, \dots, x_N
2. Augment each example twice to get $\tilde{x}_1, \dots, \tilde{x}_N, \tilde{x}'_1, \dots, \tilde{x}'_N$
3. Embed examples with f_θ to get $\tilde{z}_1, \dots, \tilde{z}_N, \tilde{z}'_1, \dots, \tilde{z}'_N$
4. Compute all pairwise distances $d(z_i, z_j) = -\frac{z_i^T z_j}{\|z_i\| \|z_j\|}$
5. Update θ w.r.t. loss $\mathcal{L}_{N\text{-way}}(\theta) = -\sum_i \log \frac{\exp(-d(\tilde{z}_i, \tilde{z}'_i))}{\sum_{j \neq i} \exp(-d(\tilde{z}_i, \tilde{z}_j))}$



After Pre-Training: train classifier on top of representation or fine-tune entire network.

Performance of Contrastive Learning

ImageNet Classification Results

Method	Architecture	Label fraction		
		1%	10%	Top 5
Supervised baseline	ResNet-50	48.4	80.4	
<i>Methods using other label-propagation:</i>				
Pseudo-label	ResNet-50	51.6	82.4	
VAT+Entropy Min.	ResNet-50	47.0	83.4	
UDA (w. RandAug)	ResNet-50	-	88.5	
FixMatch (w. RandAug)	ResNet-50	-	89.1	
S4L (Rot+VAT+En. M.)	ResNet-50 (4×)	-	91.2	
<i>Methods using representation learning only:</i>				
InstDisc	ResNet-50	39.2	77.4	
BigBiGAN	RevNet-50 (4×)	55.2	78.8	
PIRL	ResNet-50	57.2	83.8	
CPC v2	ResNet-161(*)	77.9	91.2	
SimCLR (ours)	ResNet-50	75.5	87.8	
SimCLR (ours)	ResNet-50 (2×)	83.0	91.2	
SimCLR (ours)	ResNet-50 (4×)	85.8	92.6	

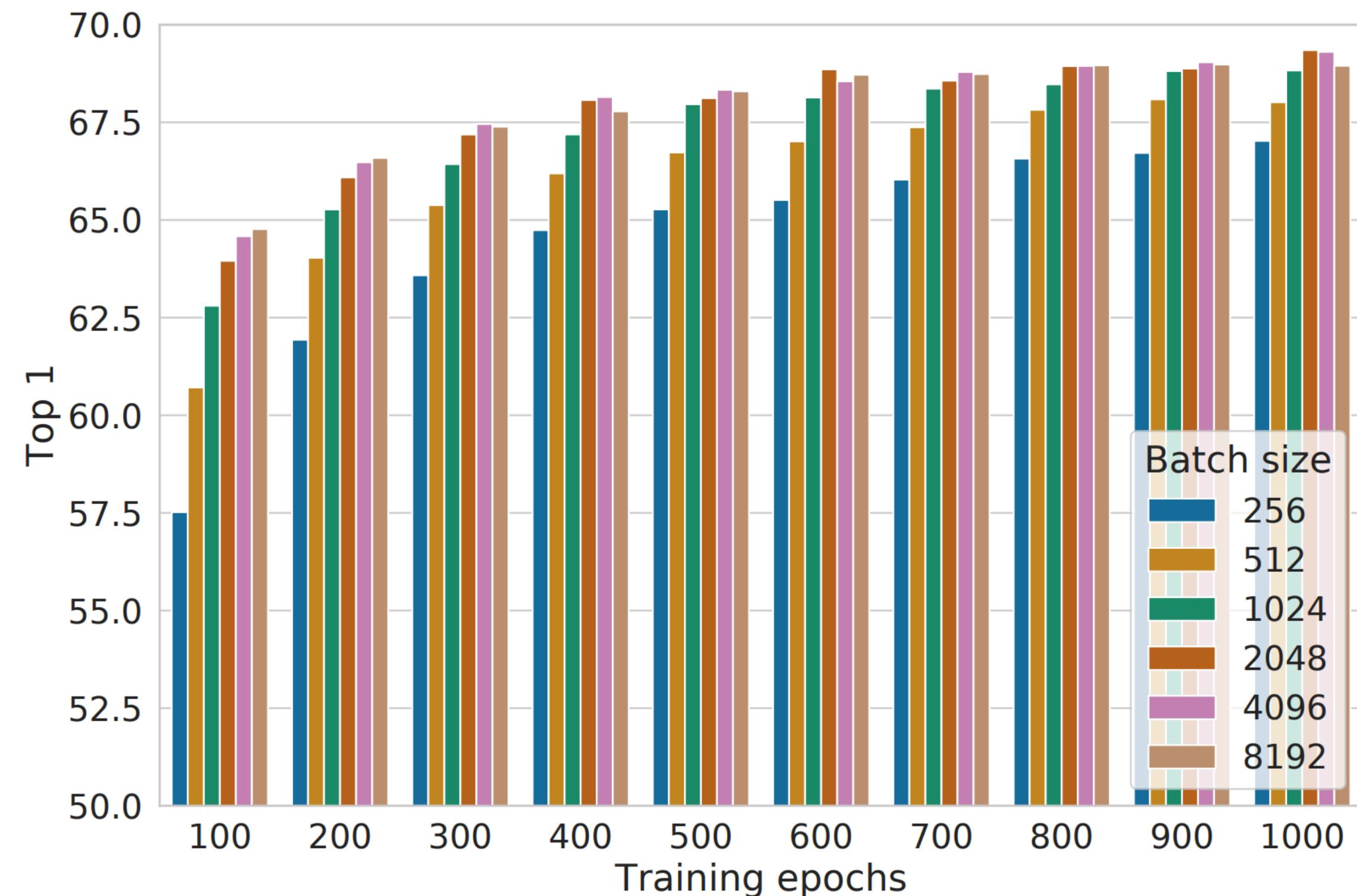
Table 7. ImageNet accuracy of models trained with few labels.

1% labels: ~12.8 images/class

- Substantial improvements over training from scratch
- Improvements over other methods, especially in 1% label setting

Performance of Contrastive Learning

Effect of Batch Size & Number of Training Epochs



- Important to train for longer (~600+ epochs)
- Requires large batch size

Why does contrastive learning need a large batch size?

Interpretation of loss: classifying augmented example from rest of dataset

$$\mathcal{L}_{\text{N-way}}(\theta) = - \sum_i \log \frac{\exp(-d(\tilde{z}_i, \tilde{z}'_i))}{\sum_{j \neq i} \exp(-d(\tilde{z}_i, \tilde{z}_j))} \quad \leftarrow \text{summation over entire dataset}$$

Intuition: Closest z will dominate the denominator, can be missed when subsampling

Mathematically?

Why does contrastive learning need a large batch size?

Interpretation of loss: classifying augmented example from rest of dataset

$$\mathcal{L}_{\text{N-way}}(\theta) = - \sum_i \log \frac{\exp(-d(\tilde{z}_i, \tilde{z}'_i))}{\sum_{j \neq i} \exp(-d(\tilde{z}_i, \tilde{z}_j))} \quad \leftarrow \text{summation over entire dataset}$$

Intuition: Closest z will dominate the denominator, can be missed when subsampling

Mathematically: Minimizing a lower-bound.

Solutions to requiring a large batch size

1. Store representations from previous batches (“momentum contrast”)

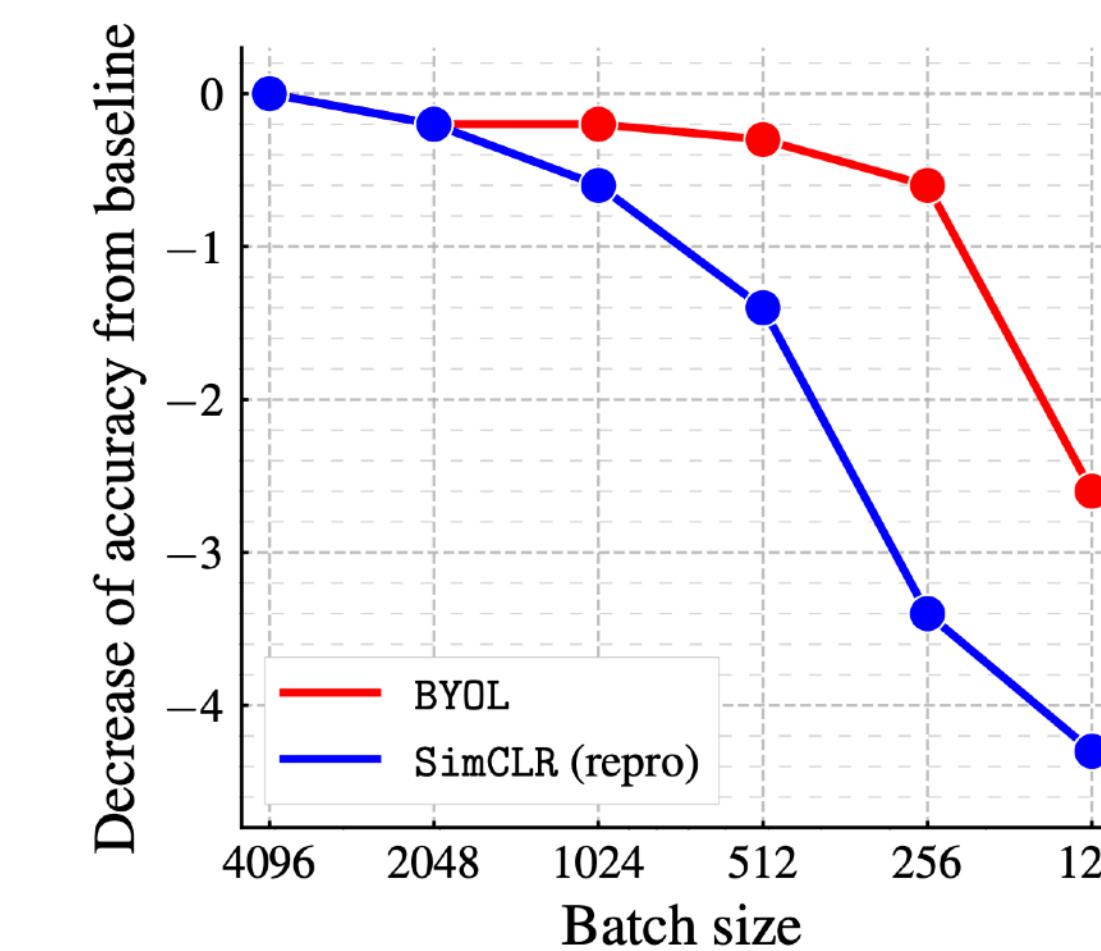
He, Fan, Wu, Xie, Girshick. MoCo. CVPR 2020

- Good results with mini batch size of 256

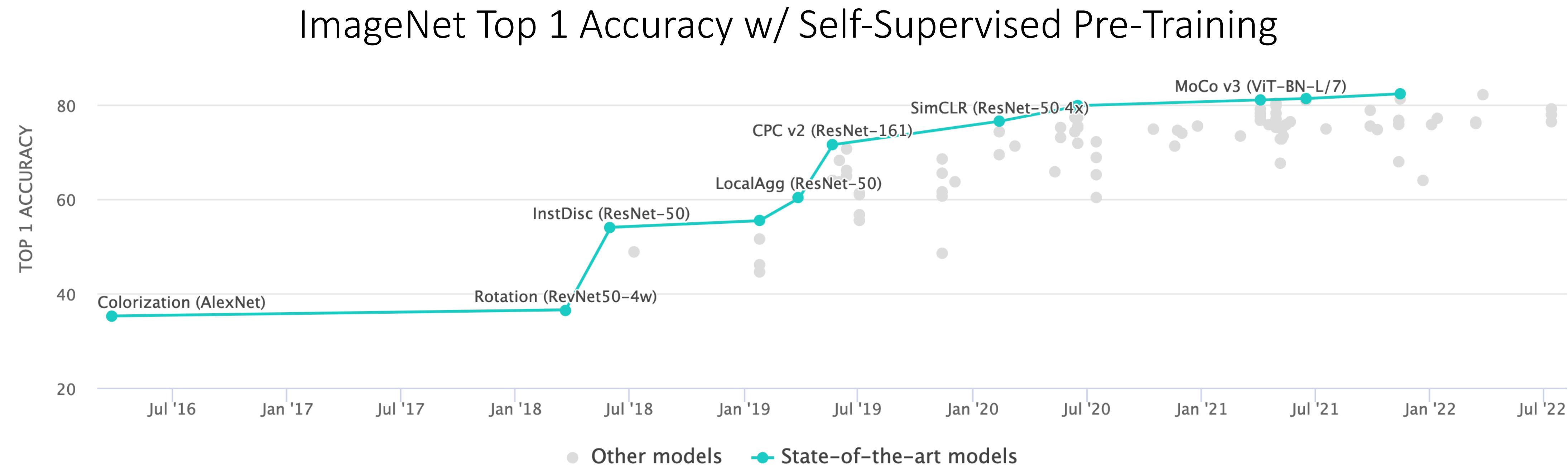
2. Predict representation of same image under different augmentation (“BYOL”)

Grill*, Strub*, Altché*, Tallec* Richemond*, et al. BYOL. NeurIPS 2020

- No negatives required!
- More resilient to batch size



Performance of contrastive learning



Contrastive methods are near state-of-the-art in self-supervised pre-training for visual data.

Contrastive learning beyond augmentations

We don't have good engineered augmentations for many applications!

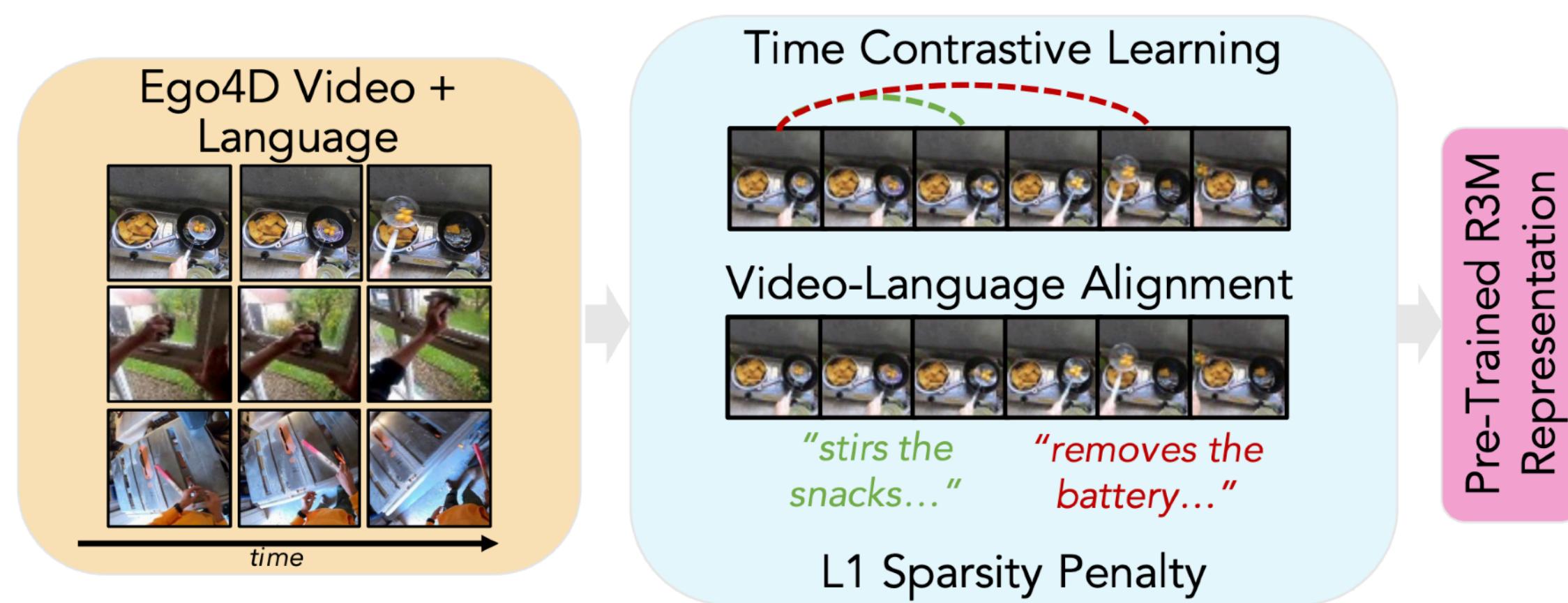
1. *Learn* the augmentations in adversarial manner (but perturbations bounded to ℓ_1 sphere)

Tamkin, Wu, Goodman. Viewmaker Networks. ICLR 2021

- > competitive with SimCLR on image data
- > good results on speech & sensor data

2. *Time-contrastive learning* on *videos* effective for robotics pre-training

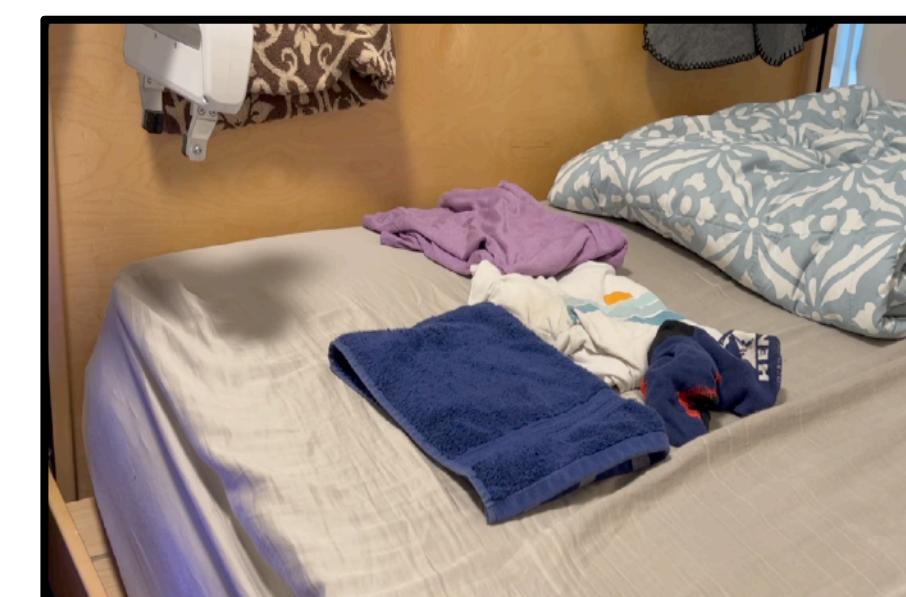
Nair, Rajeswaran, Kumar, Finn, Gupta. R3M. CoRL 2022.



Given 20 demos (<10 min of supervision)



60% success



40% success

Contrastive learning beyond augmentations

We don't have good engineered augmentations for many applications!

1. *Learn* the augmentations in adversarial manner (but perturbations bounded to ℓ_1 sphere)

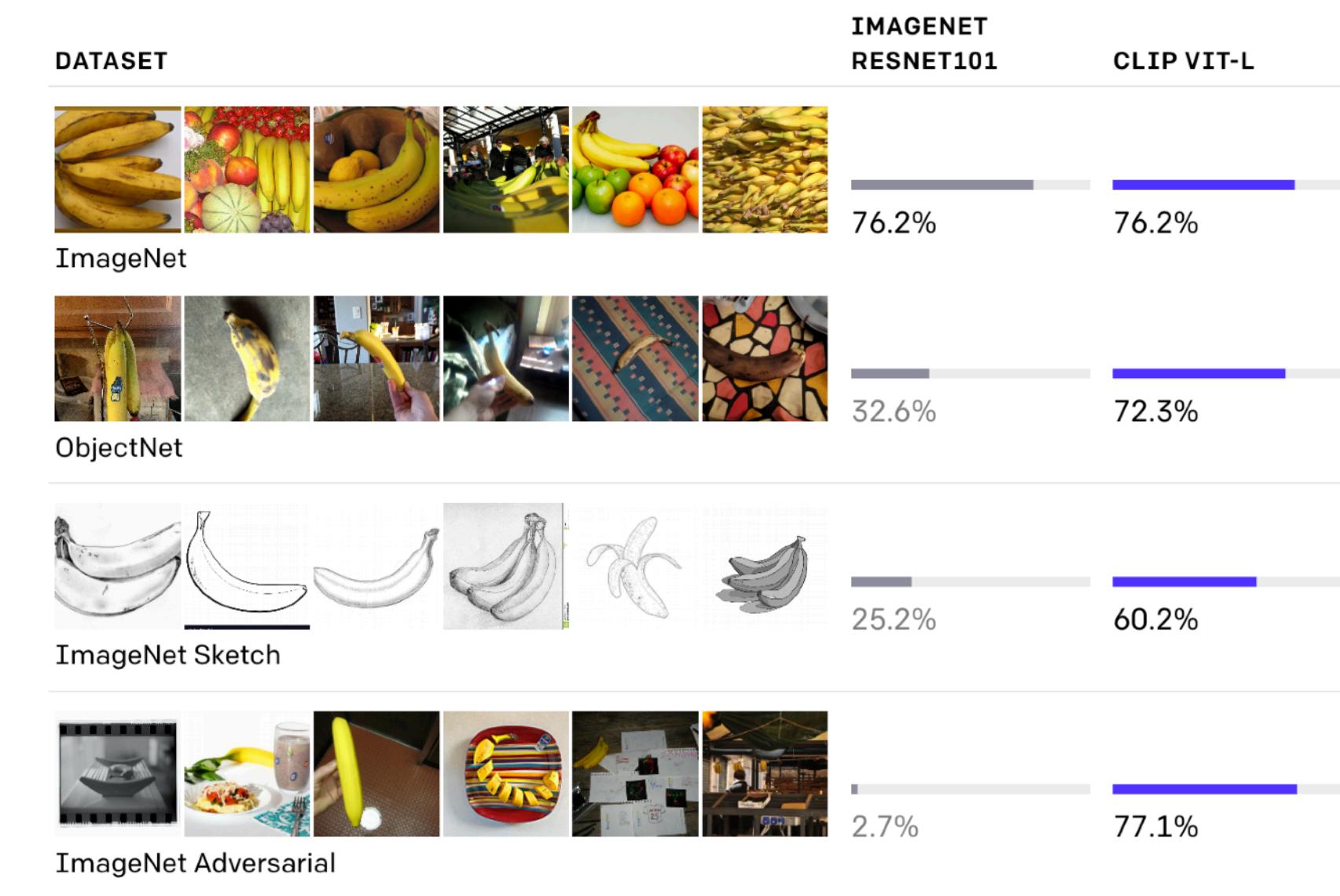
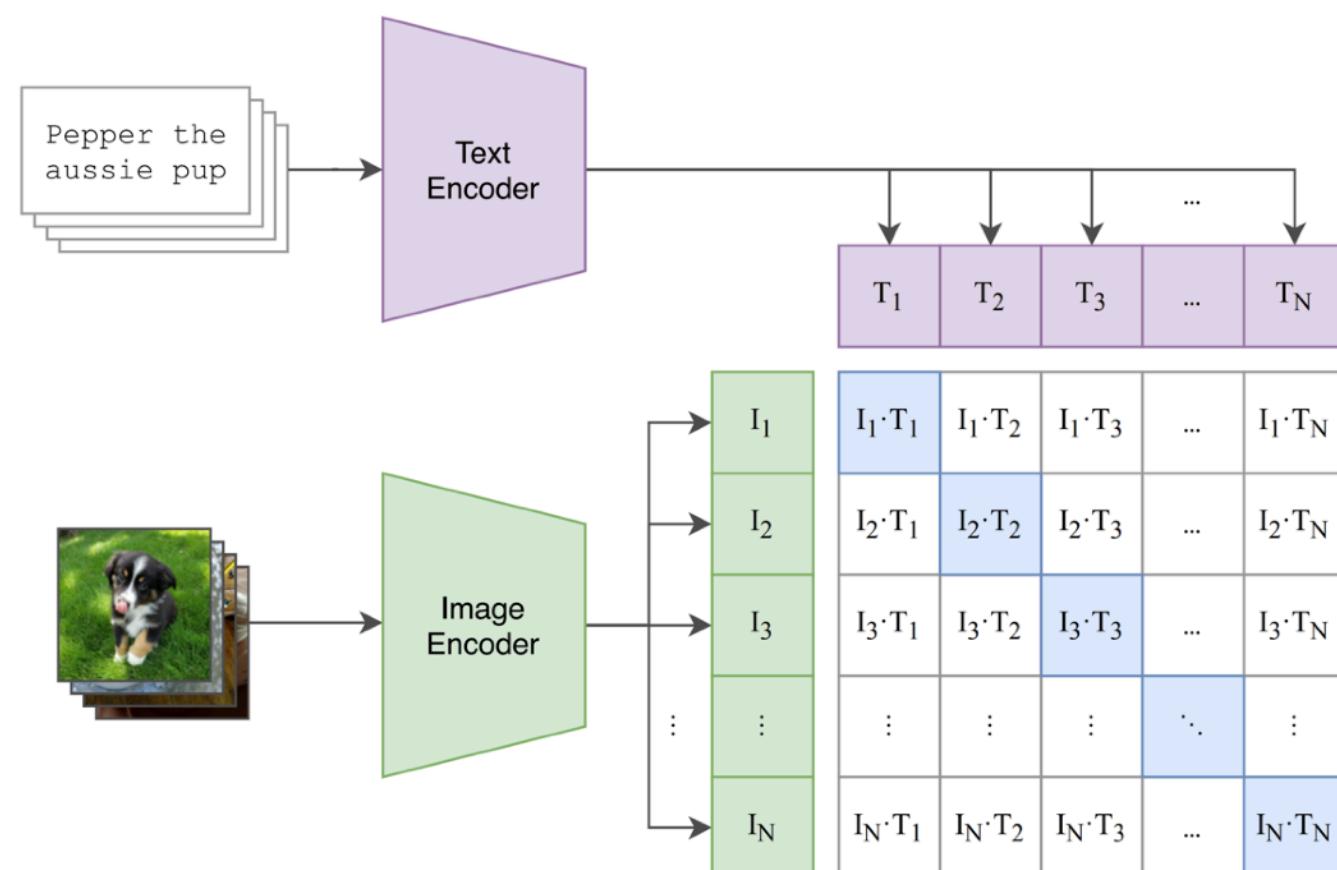
Tamkin, Wu, Goodman. Viewmaker Networks. ICLR 2021

2. *Time-contrastive learning* on *videos* effective for robotics pre-training

Nair, Rajeswaran, Kumar, Finn, Gupta. R3M. CoRL 2022.

3. Image-text contrastive pre-training produces robust zero-shot models

Radford*, Kim*, et al. CLIP. 2021.



Summary of Contrastive Learning

Pros:

- + General, effective framework
- + No generative modeling required
- + Can incorporate domain knowledge through augmentations

Challenges:

- *Negatives* can be hard to select
- Often requires *large batch size*
- Most successful with augmentations

Contrastive Learning as Meta-Learning

Meta-learning algorithm

1. Given unlabeled dataset $\{x_i\}$.
2. Create image class y_i from each datapoint via data augmentation $\mathcal{D}_i := \{\tilde{x}_i, \tilde{x}'_i, \dots\}$
3. Run your favorite meta-learning algorithm.

Differences:

- SimCLR samples **one task** per minibatch; meta-learning usually samples **multiple**
- SimCLR compares **all pairs** of samples; meta-learning compares query examples only to support examples & not to other query examples.

Contrastive Learning as Meta-Learning

Meta-learning algorithm

1. Given unlabeled dataset $\{x_i\}$.
2. Create image class y_i from each datapoint via data augmentation $\mathcal{D}_i := \{\tilde{x}_i, \tilde{x}'_i, \dots\}$
3. Run your favorite meta-learning algorithm.

Contrastive vs. meta-learning representations, transfer from ImageNet

	Flowers102	DTD	VOC2007	Aircraft	Food101	SUN397	CIFAR-10	CIFAR-100
SimCLR	92.4	72.7	66.0	83.7	86.3	57.4	94.8	79.1
ProtoNet	92.7	71.5	64.7	83.9	86.2	56.4	96.0	79.1
R2-D2	94.5	73.8	69.9	86.2	86.9	59.7	96.7	82.8

Representations transfer similarly well.

Plan for Today

Recap

- Problem formulation
- Contrastive learning

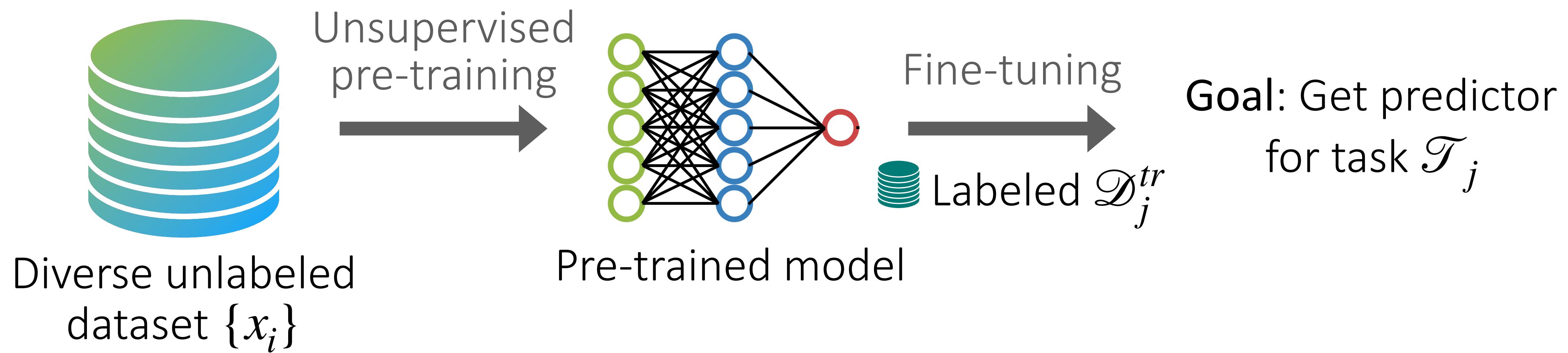
Reconstruction-based unsupervised pre-training

- Why reconstruction?
- Autoencoders
- Masked autoencoders: BERT, MAE
- Autoregressive models: GPT, Flamingo

Goals for by the end of lecture:

- Familiarize you with **widely-used** methods for unsupervised pre-training
- Introduce methods for **efficient fine-tuning** of pre-trained models
- Prepare you for **HW3**

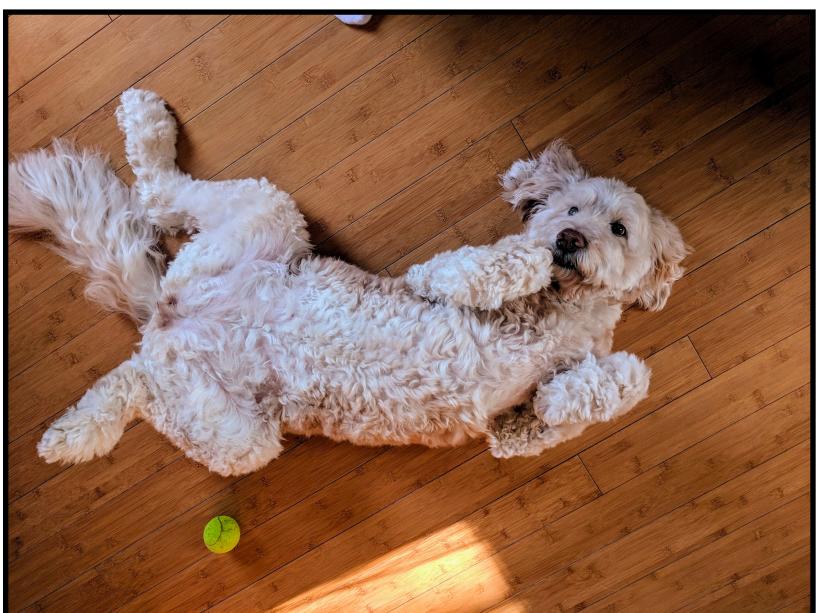
Unsupervised Pre-Training Set-Up



Key Idea of Contrastive Learning

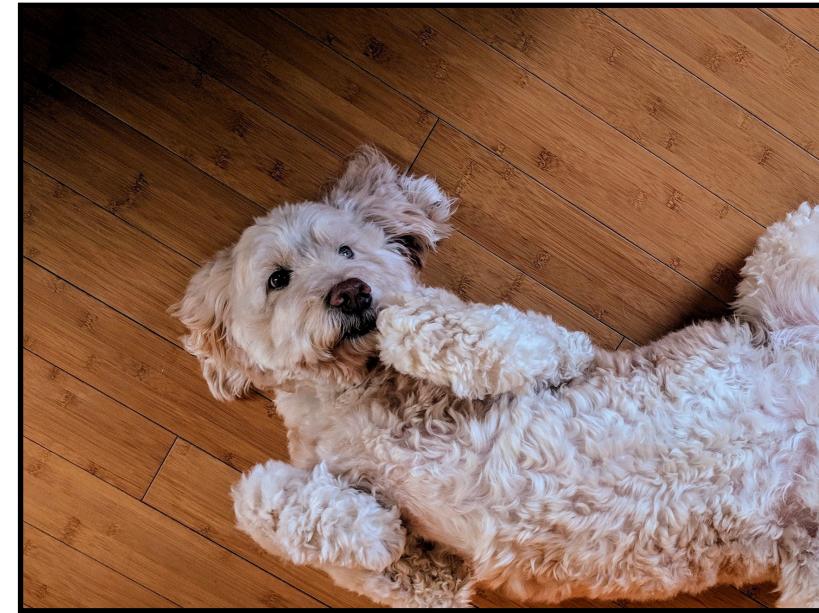
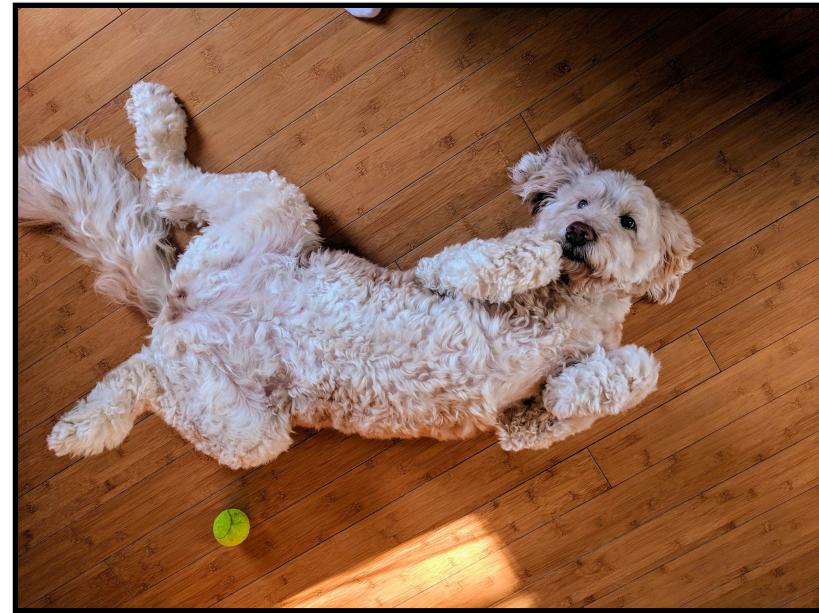
Similar examples should have similar representations

Examples with the same class label



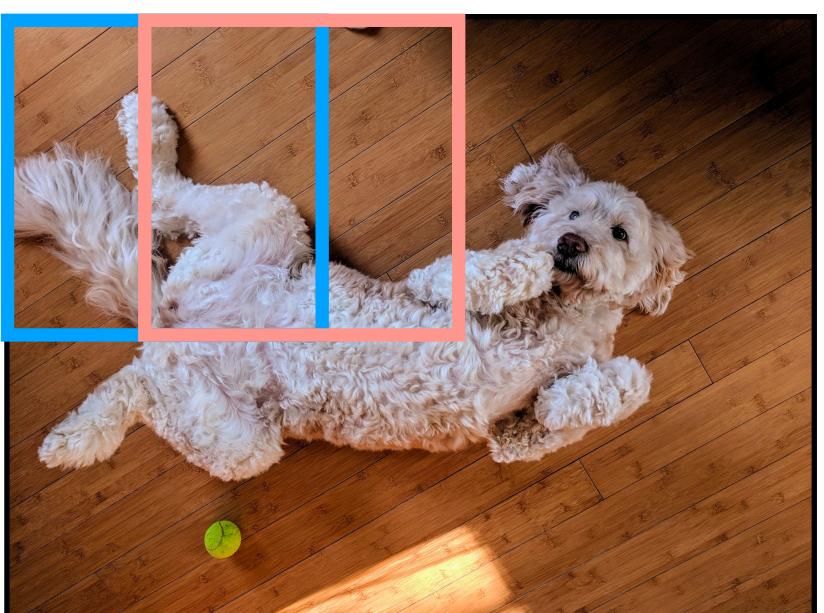
(Requires labels, related to Siamese nets, ProtoNets)

Augmented versions of the example



(flip & crop)

Nearby image patches



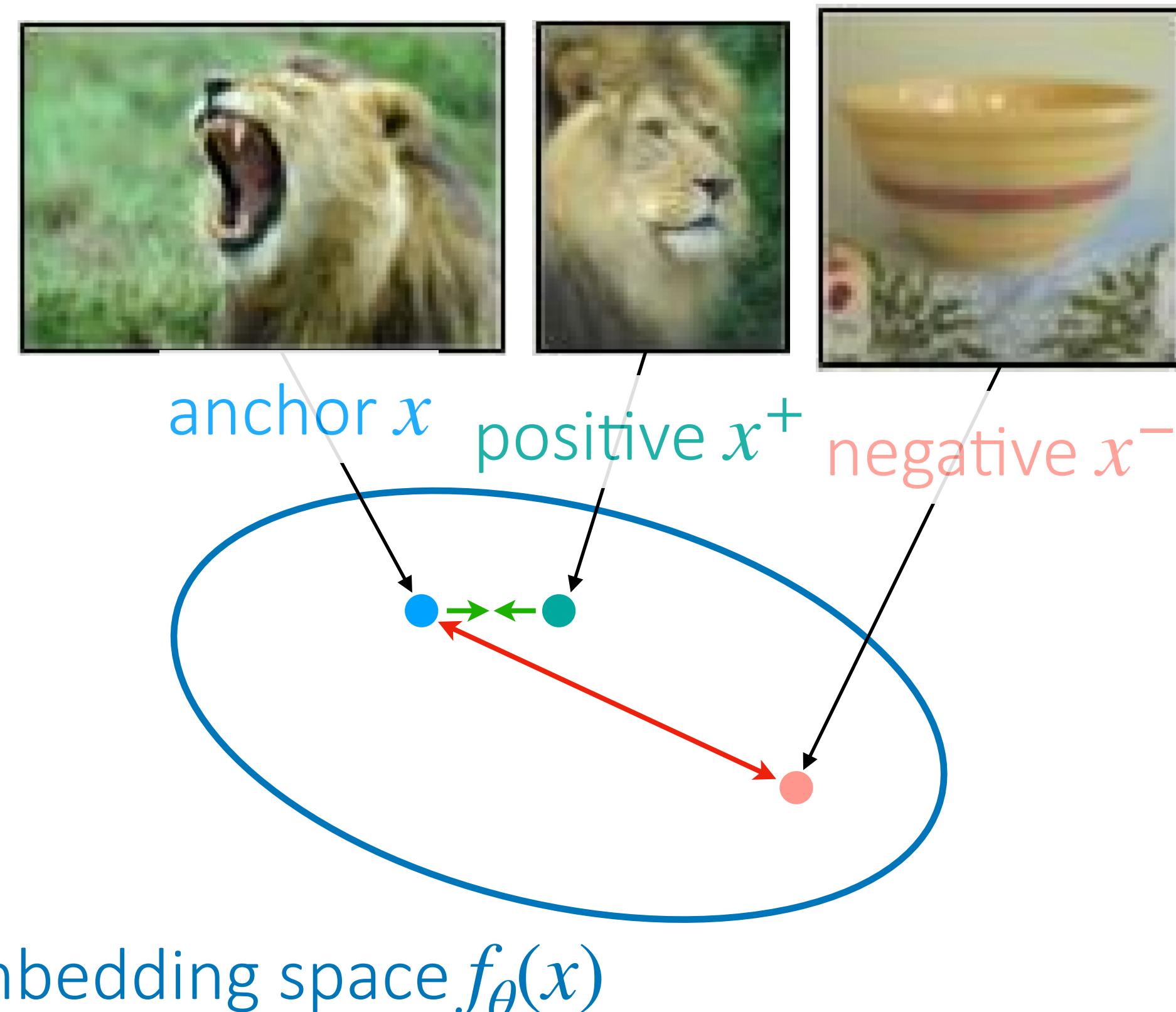
Nearby video frames



Contrastive Learning Implementation

Similar examples should have similar representations

Need to both compare & *contrast*!



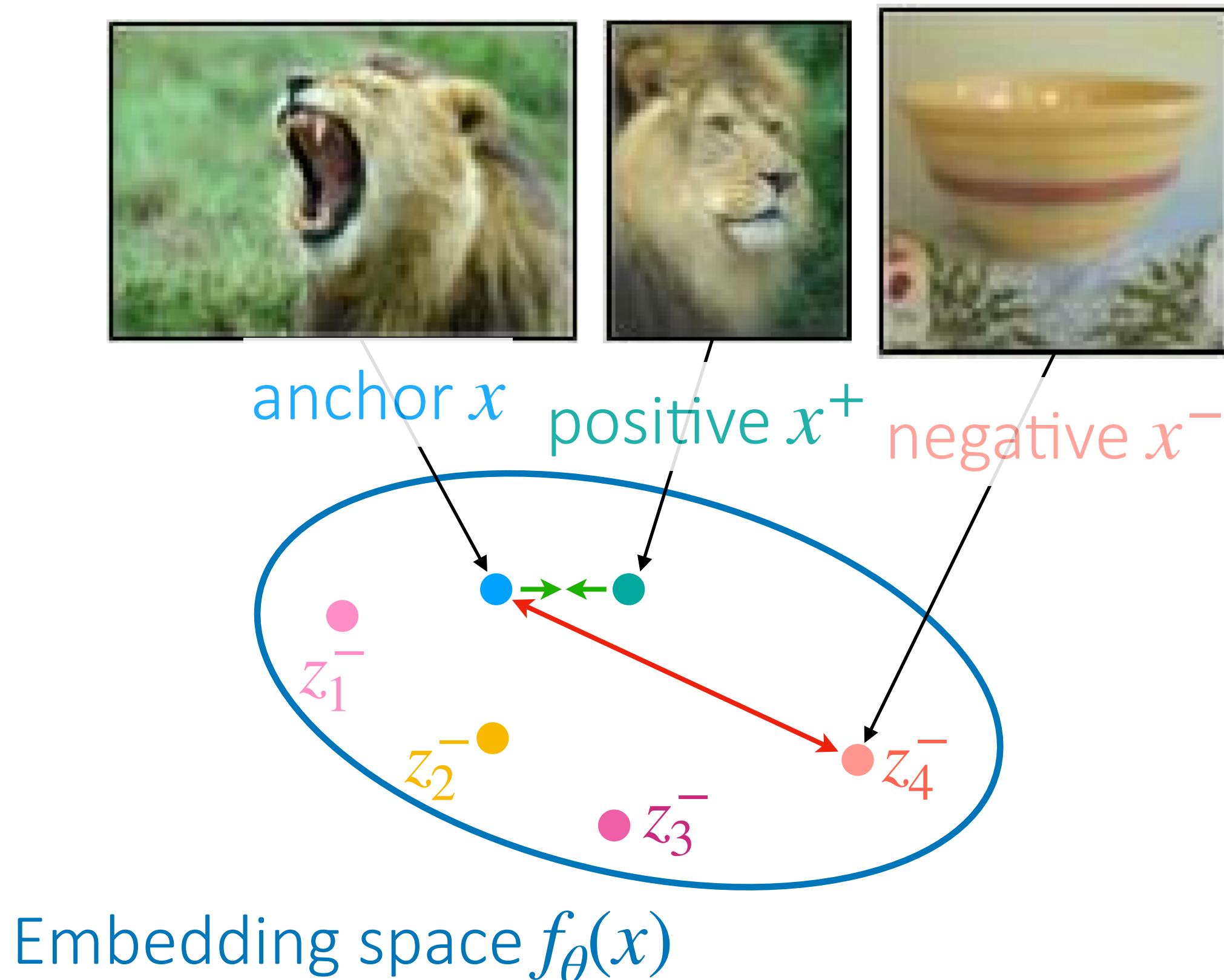
V1. Triplet loss:

$$\min_{\theta} \sum_{(x,x^+,x^-)} \max(0, \|f_\theta(\textcolor{blue}{x}) - f_\theta(\textcolor{teal}{x}^+)\|^2 - \|f_\theta(\textcolor{blue}{x}) - f_\theta(\textcolor{red}{x}^-)\|^2 + \epsilon)$$

Contrastive Learning Implementation

Similar examples should have similar representations

Need to both compare & *contrast*!



V1. Triplet loss:

$$\min_{\theta} \sum_{(x, x^+, x^-)} \max(0, \|f_\theta(x) - f_\theta(x^+)\|^2 - \|f_\theta(x) - f_\theta(x^-)\|^2 + \epsilon)$$

V2. From binary to N-way classification (aka **simCLR***):

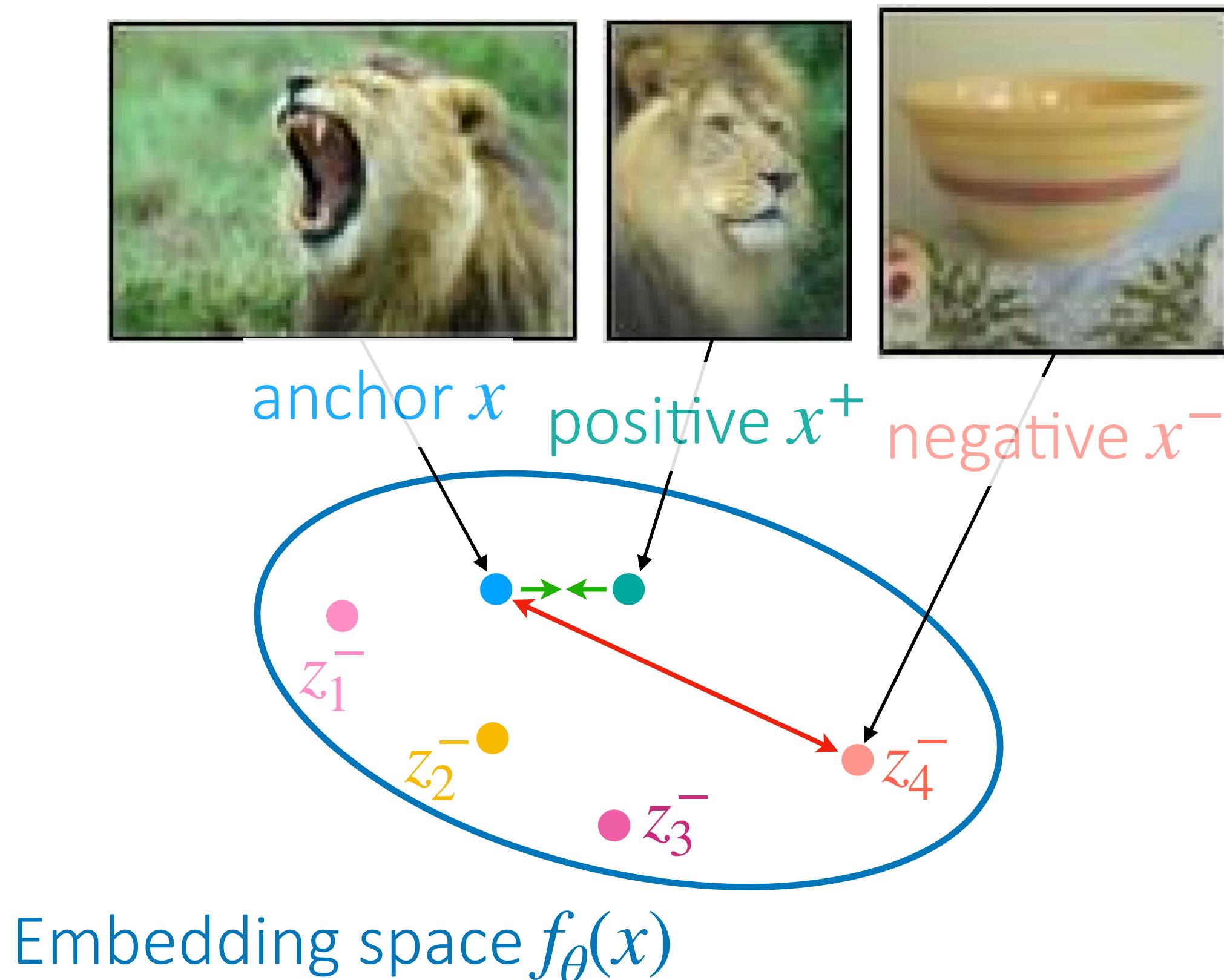
$$\mathcal{L}_{\text{N-way}}(\theta) = - \sum_z \log \frac{\exp(-d(z, z^+))}{\sum_i \exp(-d(z, z_i^-))}$$

*also known as the **NT-Xent** loss, when $d(\cdot, \cdot)$ is **scaled cosine similarity**

Contrastive Learning Implementation

Similar examples should have similar representations

Need to both compare & *contrast*!



V1. Triplet loss:

$$\min_{\theta} \sum_{(x,x^+,x^-)} \max(0, \|f_\theta(\textcolor{blue}{x}) - f_\theta(\textcolor{teal}{x}^+)\|^2 - \|f_\theta(\textcolor{blue}{x}) - f_\theta(\textcolor{red}{x}^-)\|^2 + \epsilon)$$

V2. From binary to N-way classification (aka **simCLR***):

$$\mathcal{L}_{\text{N-way}}(\theta) = - \sum_z \log \frac{\exp(-d(z, z^+))}{\exp(-d(z, z^+)) + \sum_i \exp(-d(z, z_i^-))}$$

Positive score in denominator → loss read as “*classification loss when discriminating positive pair from negatives*”

*also known as the **NT-Xent** loss, when $d(\cdot, \cdot)$ is **scaled cosine similarity**

Plan for Today

Recap

- Problem formulation
- Contrastive learning

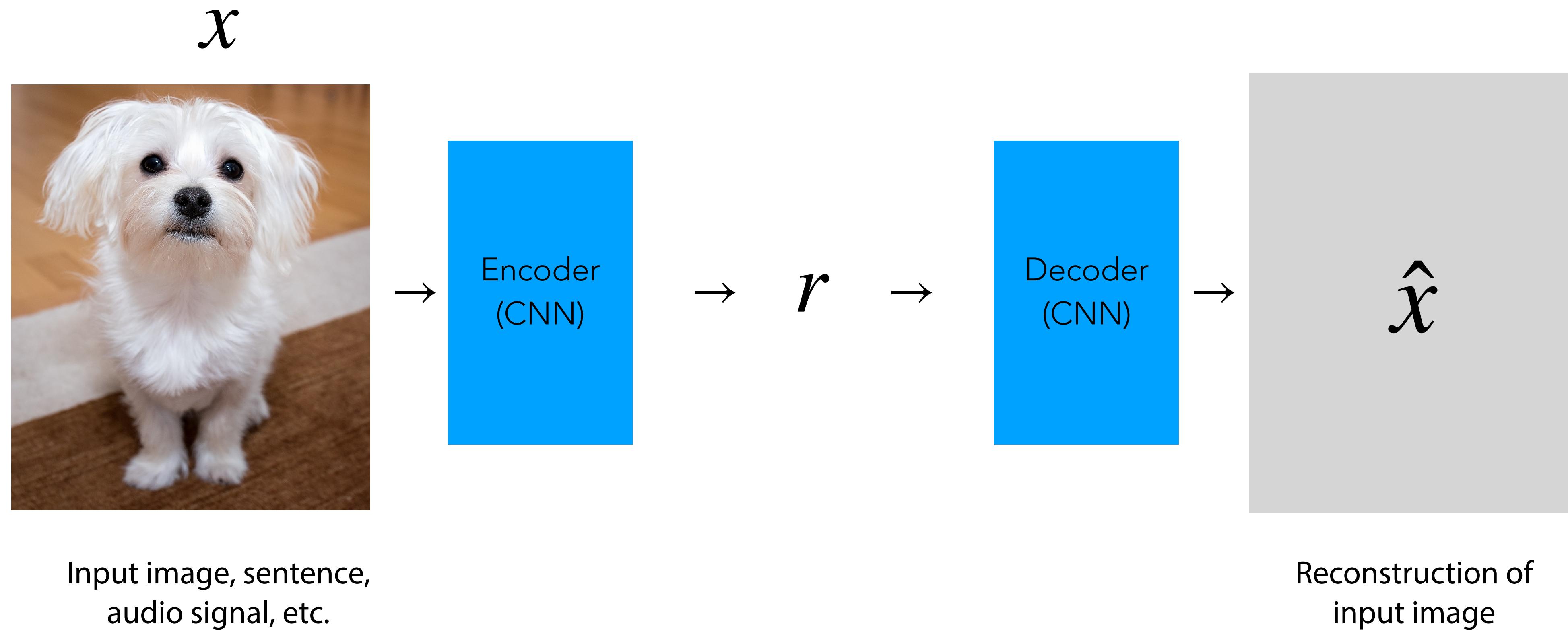
Reconstruction-based unsupervised pre-training

- **Why reconstruction?**
- Autoencoders
- Masked autoencoders: BERT, MAE
- Autoregressive models: GPT, Flamingo
- Emergent behaviors in large models

Why reconstruction?

Simple intuition: a good representation of an input should be sufficient to **reconstruct** it

Bonus: no need to worry about pesky things like **sampling negatives** or **large batch sizes!**

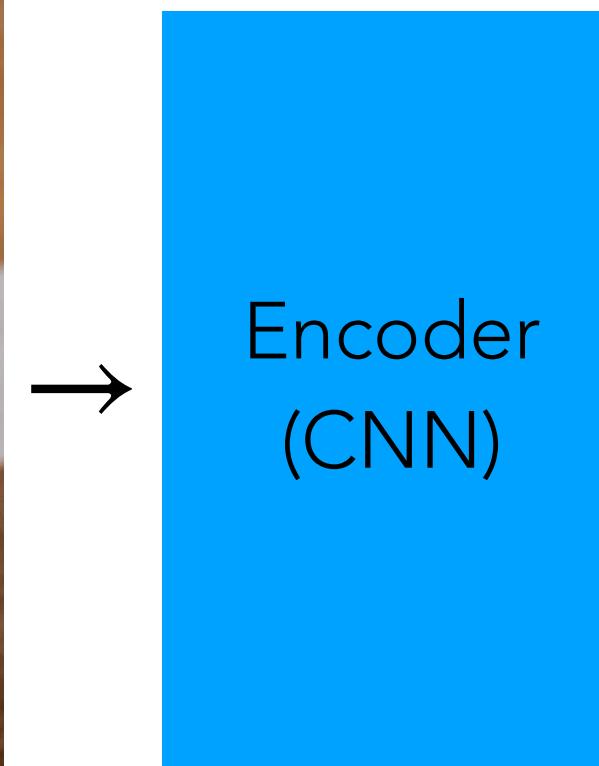


If the encoder is producing a “good” representation, a reasonably-sized decoder should be able to produce **reconstruction** \hat{x} very close to **input** x from **representation** r

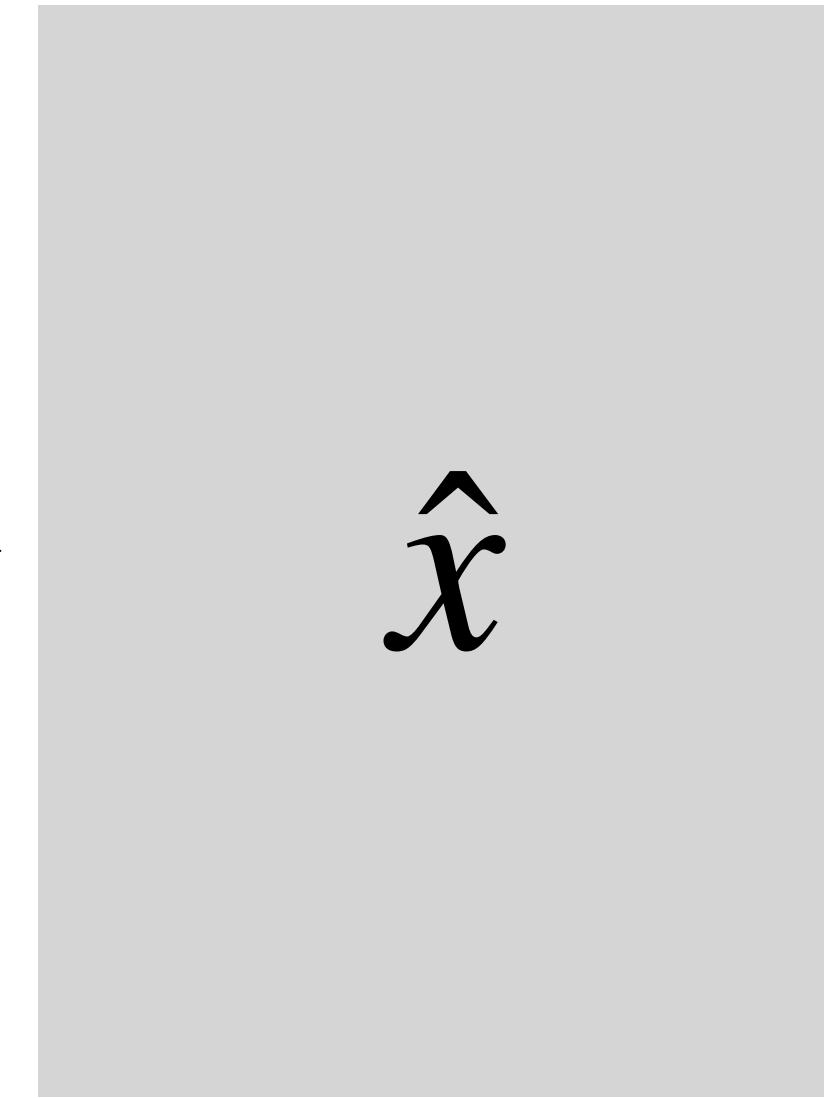
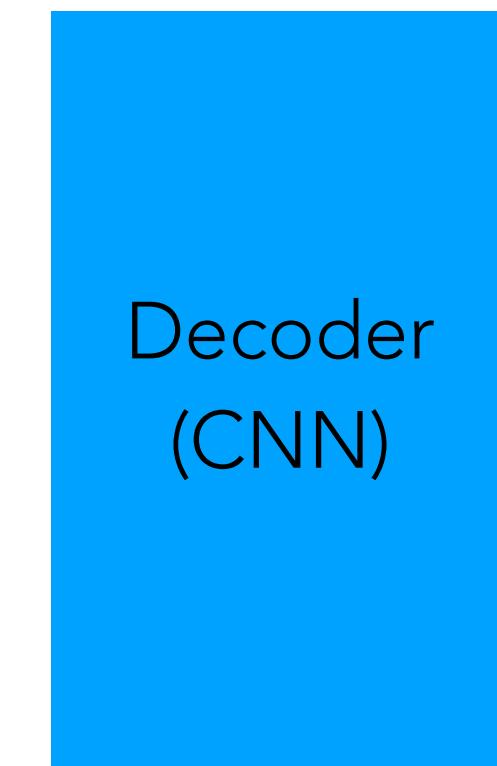
Autoencoders: a first attempt

Simple intuition: a good representation lets us **reconstruct** the input

x



→ r →



Input image, sentence,
audio signal, etc.

Reconstruction of
input image

$$\mathcal{L} = d(x, \hat{x})$$

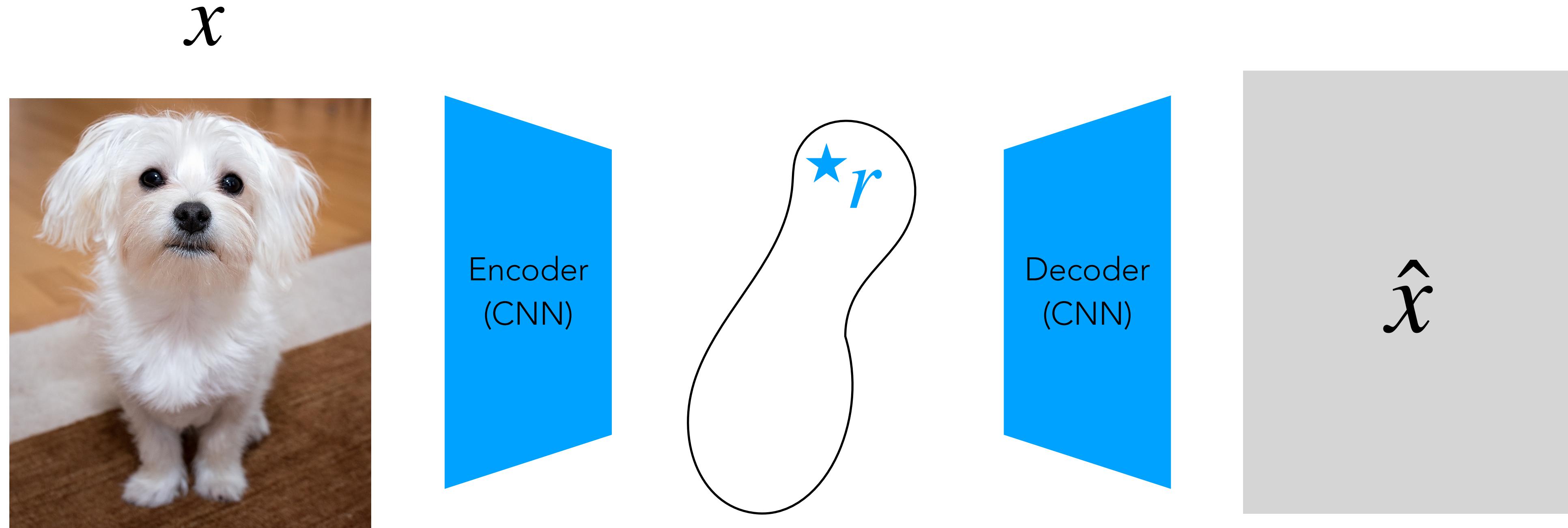
Loss function is reconstruction
error, e.g. L2 distance:

$$d(x, \hat{x}) = \|x - \hat{x}\|^2$$

What can go wrong here?

Is the **identity function** a good encoder/decoder?

Autoencoders: adding a bottleneck

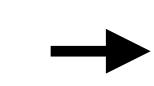


Input image, sentence,
audio signal, etc.

Compact, latent
representation of input image

\hat{x}

Key idea: latent representation is **bottlenecked**,
e.g., **lower-dimensional** than the input



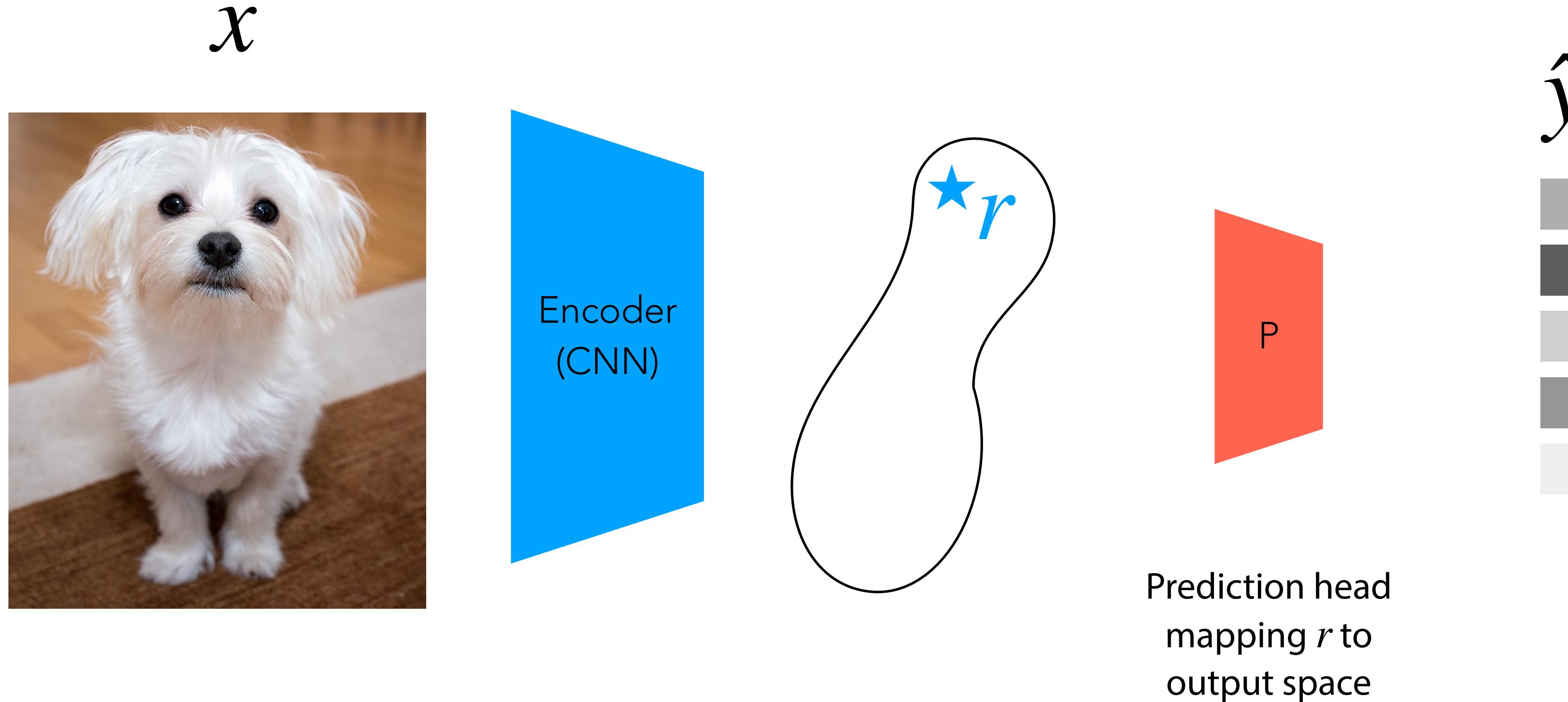
Hope: latent dimensions are forced to represent
high-level concepts that **generalize** to other tasks

$$\mathcal{L} = d(x, \hat{x})$$

Loss function is reconstruction
error, e.g. L2 distance:

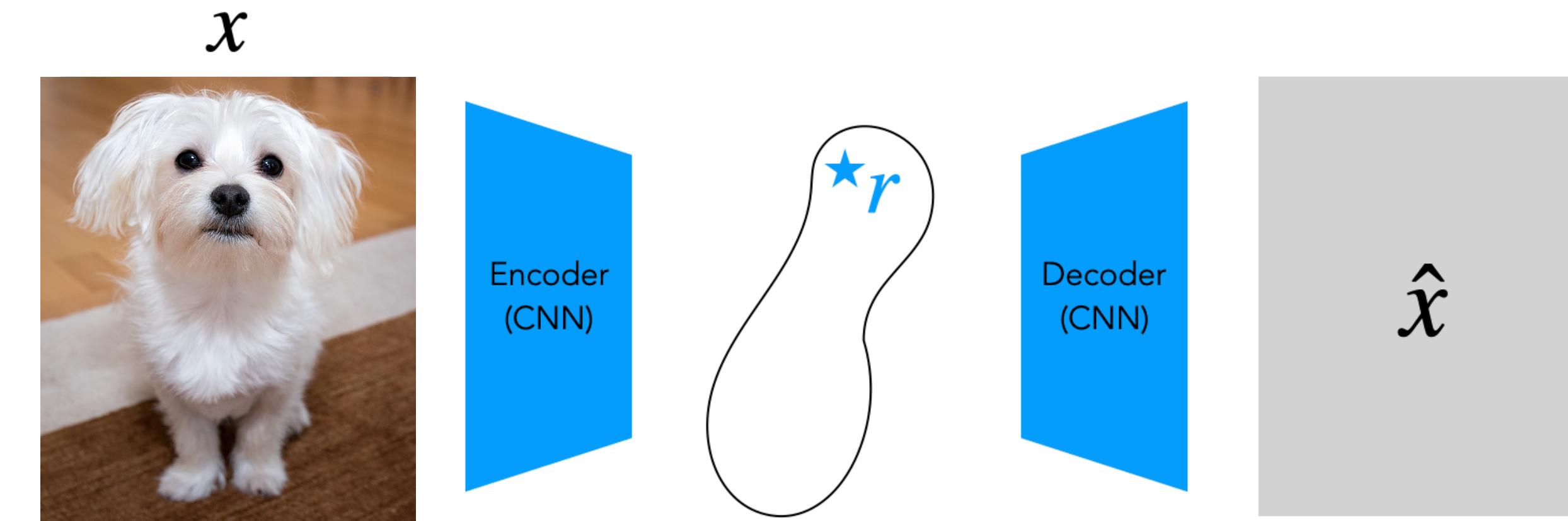
$$d(x, \hat{x}) = \|x - \hat{x}\|^2$$

Autoencoders: few-shot learning



Few-shot learning recipe: freeze **encoder**, fine-tune **prediction head** using our few-shot data
(e.g., a linear layer)

Autoencoders



Pros:

- Simple, general
- Just need to pick $d(x, \hat{x})$
- No need to select positive/negative pairs

Cons:

- Need to design a bottlenecking mechanism
- Relatively poor few-shot performance

Why?

r is just **memorizing** details of x needed to minimize pixel-level reconstruction loss



r is more like a **hash** of x than a **conceptual summary**

How do we encourage the encoder to extract high-level features?

One strategy is **other types of bottlenecks**:

- **information** bottlenecks (adding noise)
- **sparsity** bottlenecks (zero most dimensions)
- **capacity** bottlenecks (weak decoder)

In practice, we'll stop worrying about designing bottlenecks and just make the task a little harder

Plan for Today

Recap

- Problem formulation
- Contrastive learning

Reconstruction-based unsupervised pre-training

- Why reconstruction?
- Autoencoders
- **Masked autoencoders:** BERT, MAE
- Autoregressive models: GPT, Flamingo

Beyond the bottleneck: *masked autoencoders*

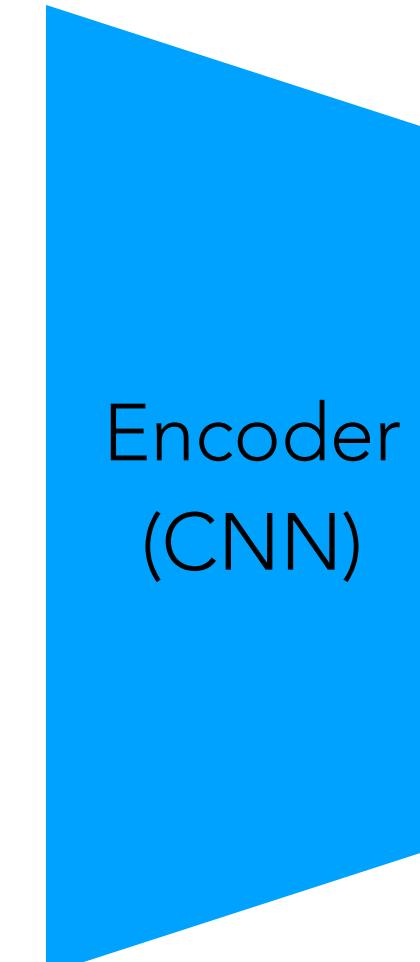
Ultimately, **regular** autoencoders are trying to predict x from... x (through z)

We bottleneck z to avoid **totally degenerate** solutions, but what if the task is just “too easy”, admitting unhelpful solutions?

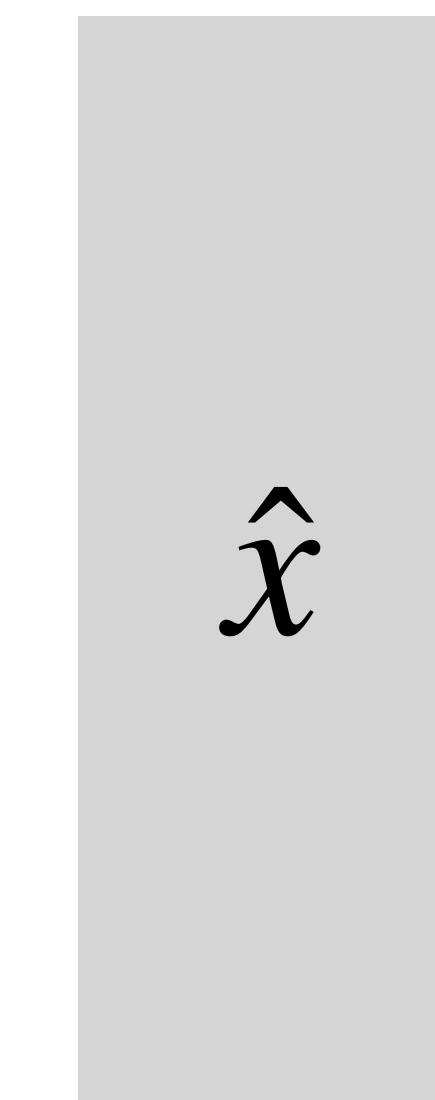
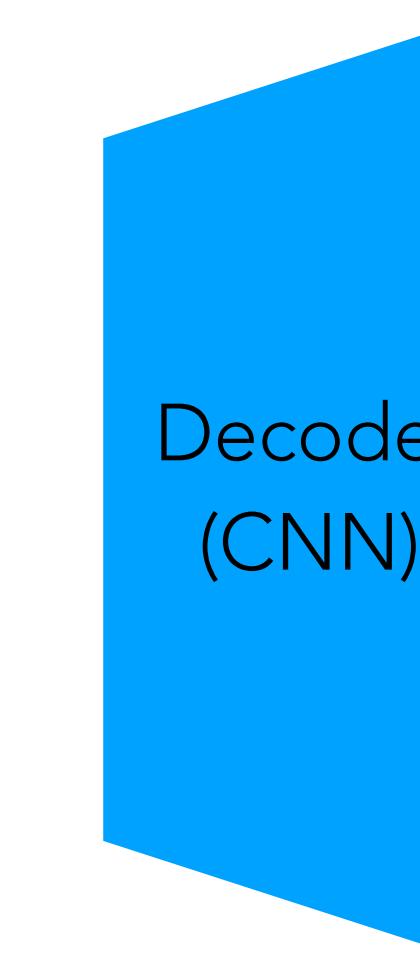
Masked autoencoders use a **more difficult** learning task to encourage the encoder to extract more meaningful features



Masked input image

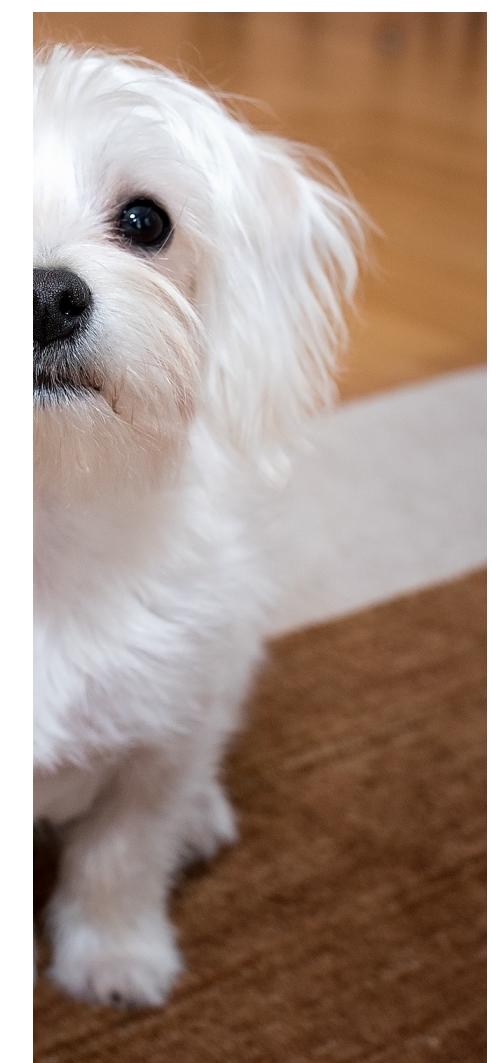


Latent representation



\hat{x}

\approx



Reconstruction of **masked portion** (or **entire**) input image

Beyond the bottleneck: *masked autoencoders*

General recipe for pre-training masked autoencoder f_θ :

1. Choose **distance function** $d(\cdot, \cdot) \rightarrow \mathbb{R}$
2. For **train batch** examples x_i :

These pieces
are our design
choices/control
knobs

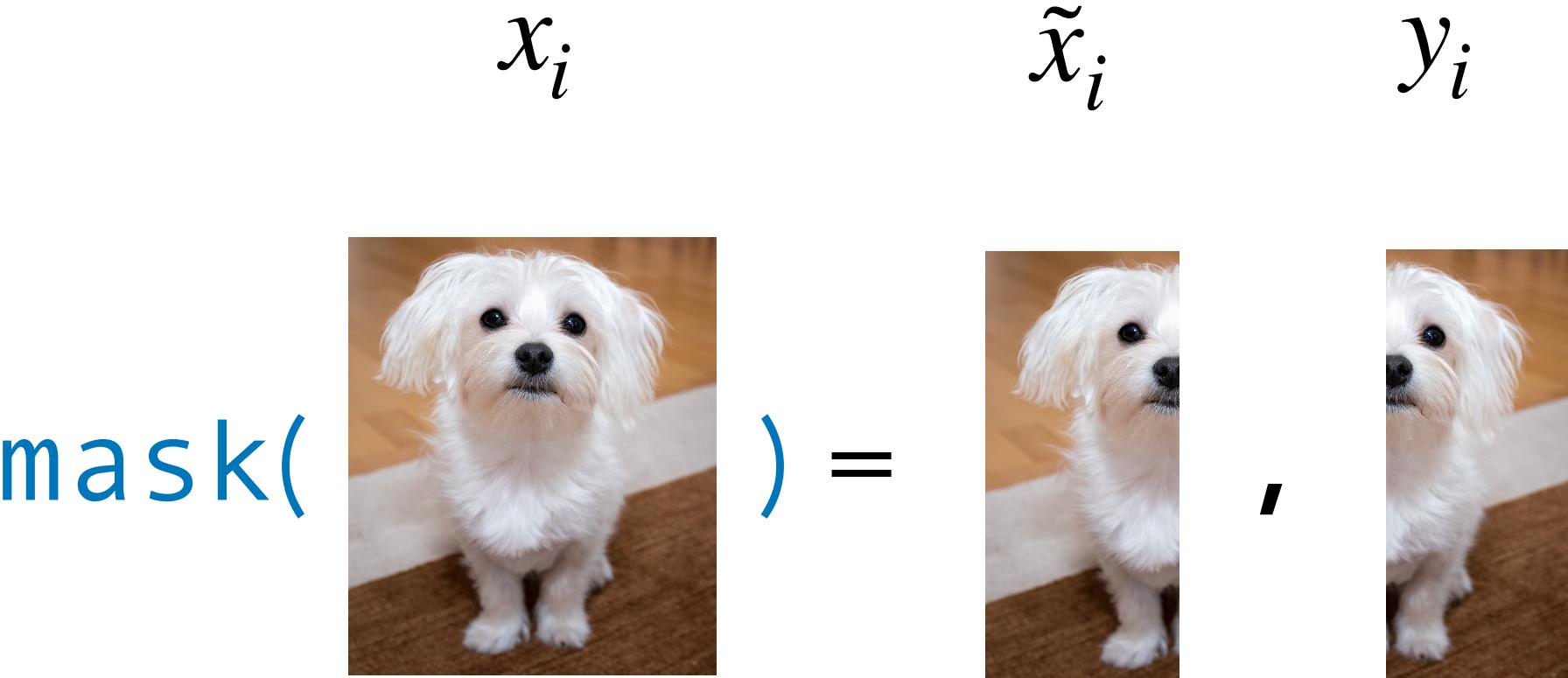
A. Sample $\tilde{x}_i, y_i \sim \text{mask}(x_i)$

\tilde{x}_i, y_i are typically two **disjoint** sub-regions of x_i

B. Make prediction $\hat{y}_i = f_\theta(\tilde{x}_i)$

in some cases, the target y_i may be all of x_i

C. Compute loss $\mathcal{L}_i = d(y_i, \hat{y}_i)$



f_θ : CNN or **Transformer** (stay tuned)

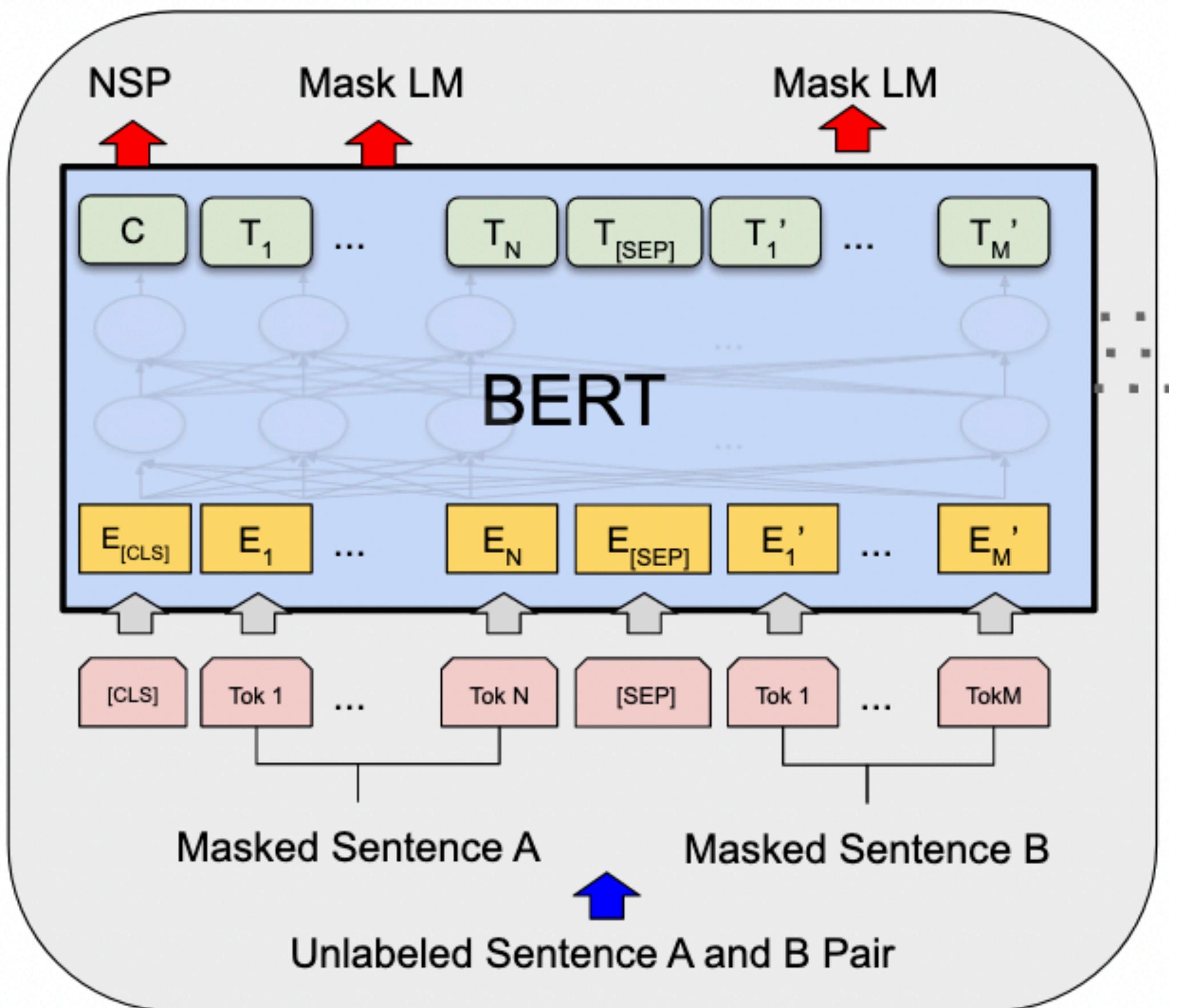
$$d(y, \hat{y}) = \|y - \hat{y}\|^2$$

x_i
 $\text{mask(Joe Biden is the US president)} =$
 \tilde{x}_i
 y_i
 $\text{Joe } <\text{mask}> \text{ is the US } <\text{mask}>, \{ \text{Biden; president} \}$

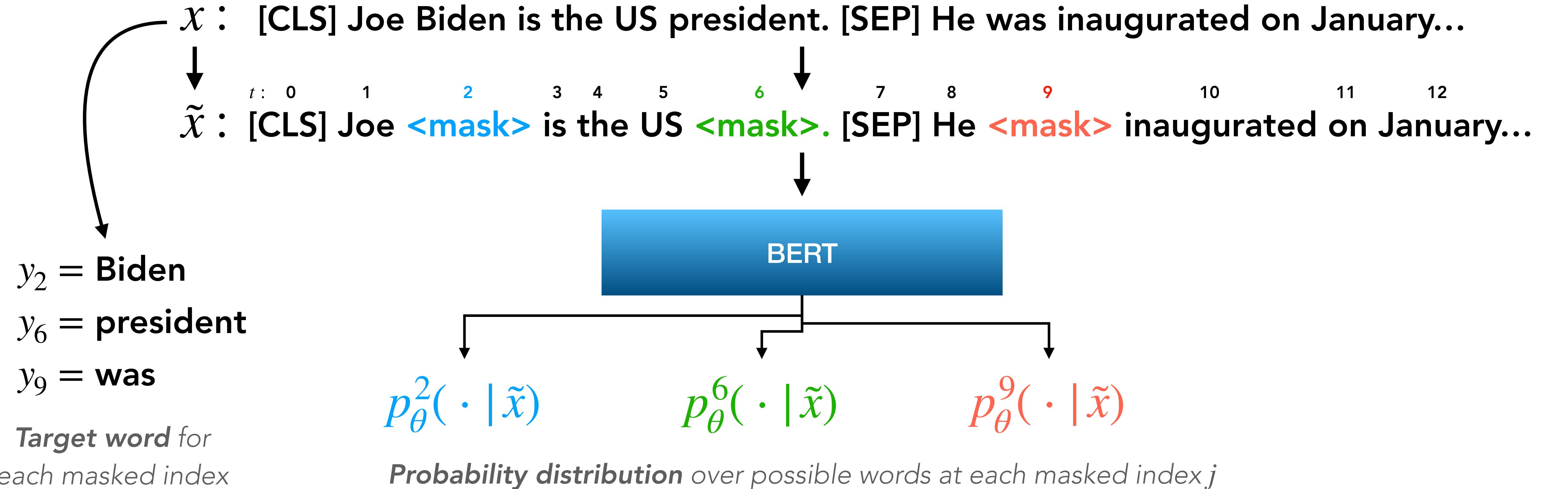
f_θ : **Transformer** (e.g., BERT; stay tuned)

$$d(y, \hat{y}) = \text{KL}(y \parallel \hat{y})$$

Masked autoencoders for language: **BERT** (Devlin et al, 2017)



Case study: BERT as a masked autoencoder



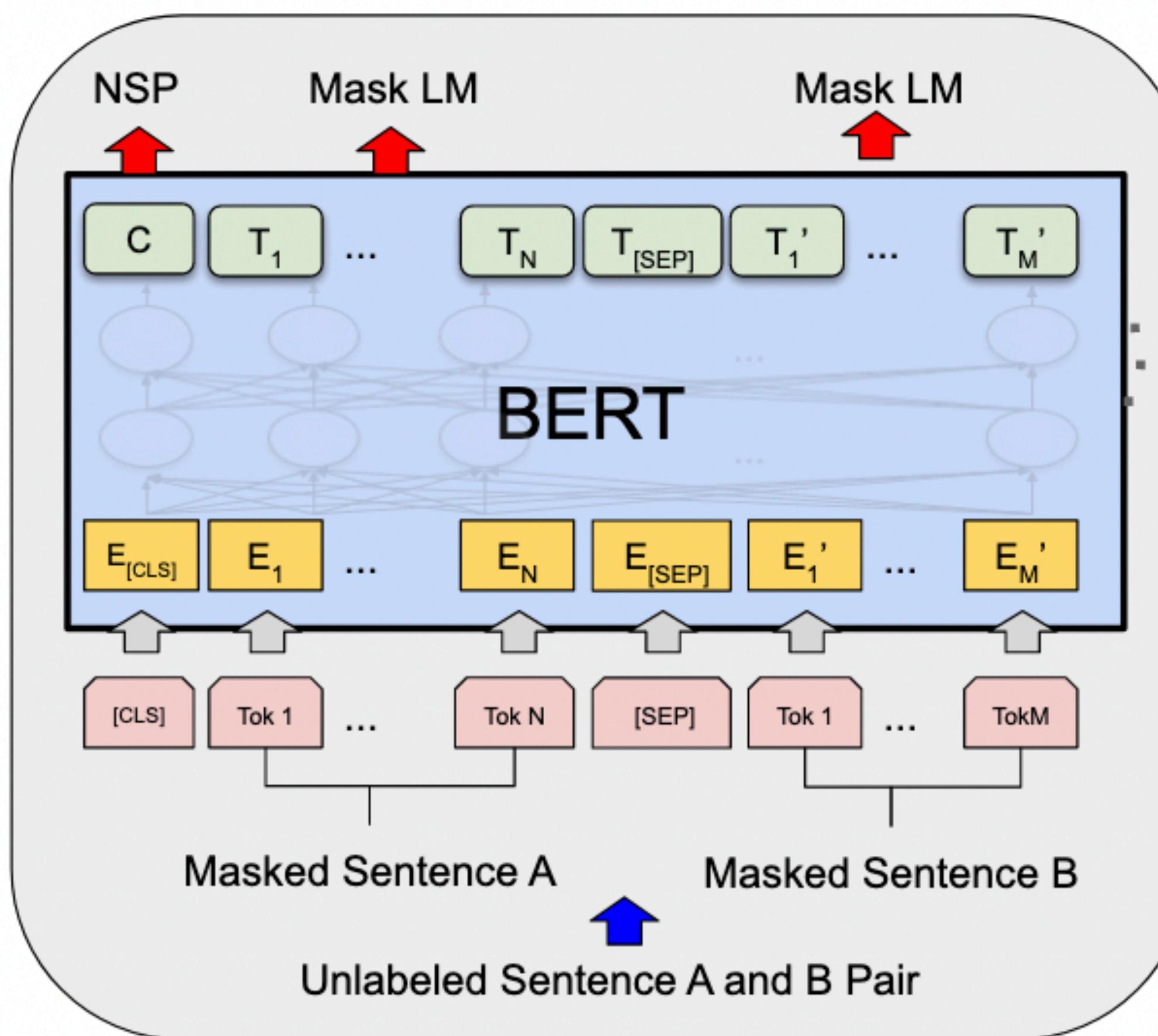
$$d(y, \hat{y}) = \sum_j \text{KL}(y_j \parallel \hat{y}_j) = -\log p_{\theta}^2(\text{Biden} | \tilde{x}) - \log p_{\theta}^6(\text{president} | \tilde{x}) - \log p_{\theta}^9(\text{was} | \tilde{x})$$

Details of BERT masking:

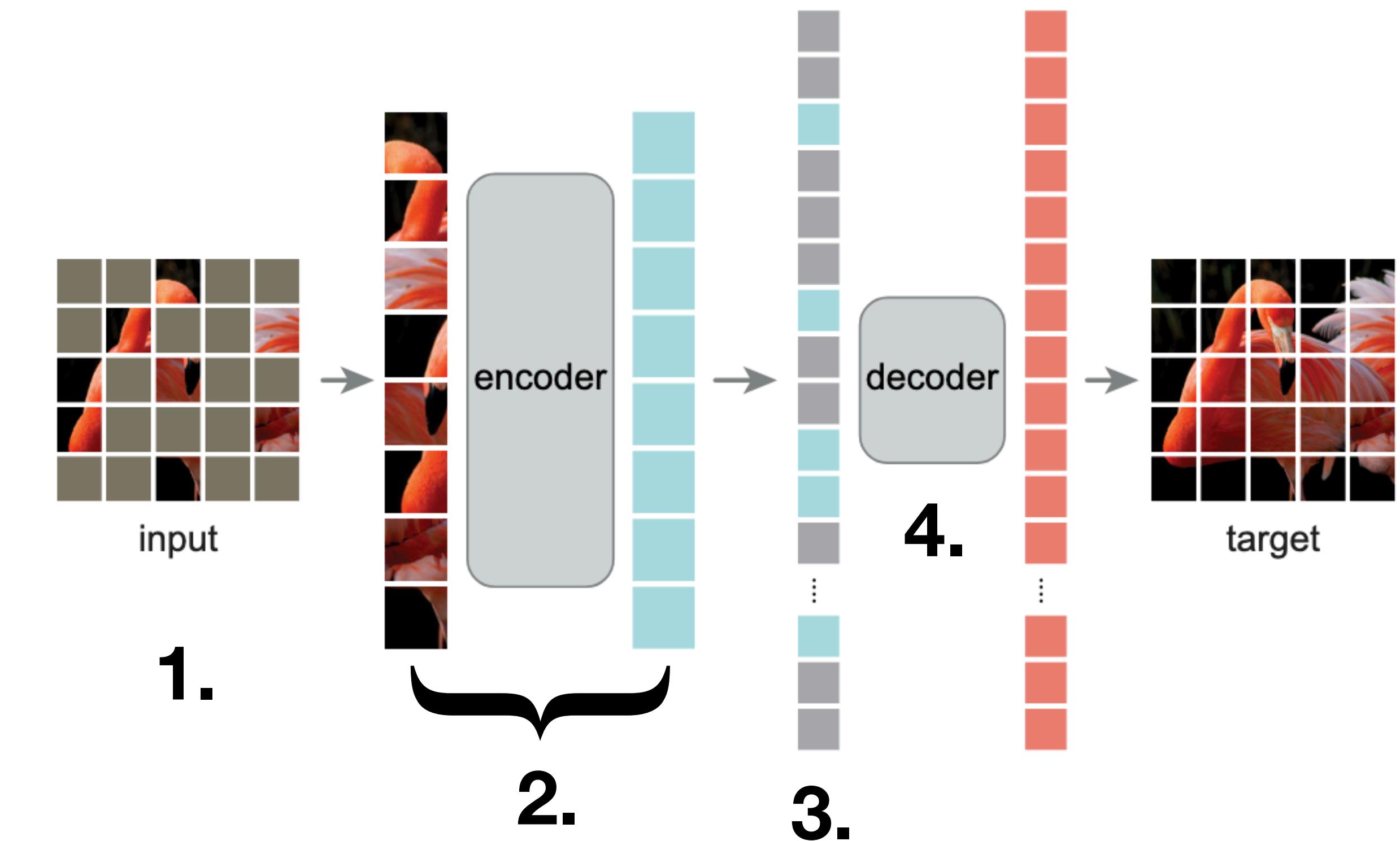
1. Choose **random 15%*** of input timesteps
2. Of these, **80%** are replaced with <mask> token
3. Replace **other 20%** with a **random** token

- *It's possible we can do better than just picking **random** timesteps:
- Mask **longer** spans of text
 - Selecting for **information-dense** spans

Masked autoencoders for language: **BERT** (Devlin et al, 2017)



For **images**: **MAE** (He et al, 2021)



Instead of words, we have a sequence of **image patches**

1. Mask ~75% of image patches
2. Compute representations of **only** unmasked patches
3. Insert **placeholder** patches at masked locations
4. Decode back into original image

Fine-tune on top of the output of **step 2**

More recently: Masked AEs give state-of-the-art **few-shot image classification** performance

The unsupervised masked autoencoding recipe works better than pre-training **with labels** on the **same data!**

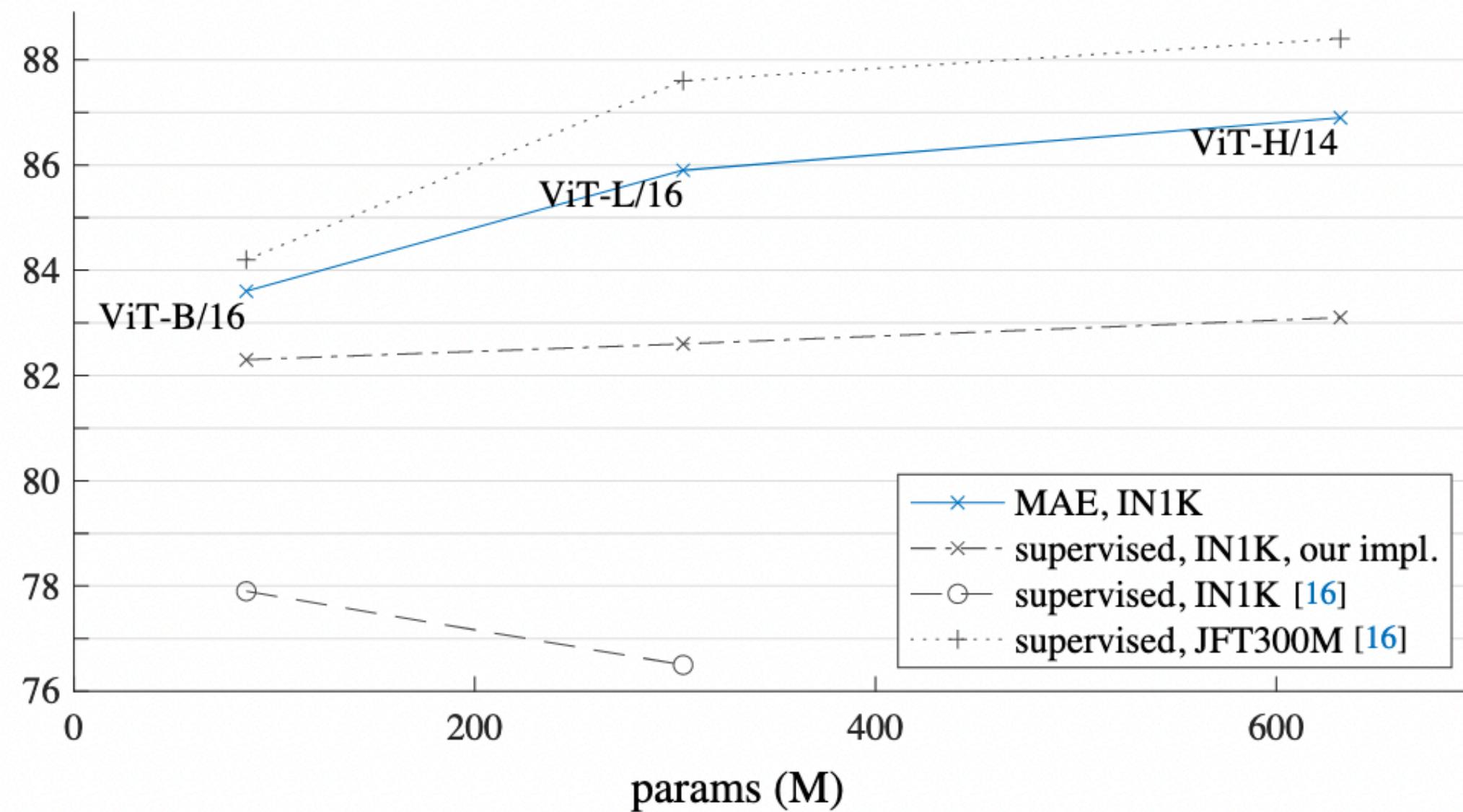


Figure 8. **MAE pre-training vs. supervised pre-training**, evaluated by fine-tuning in ImageNet-1K (224 size). We compare with the original ViT results [16] trained in IN1K or JFT300M.

When **fine-tuning** (not just **linear probing** on frozen pre-trained model), better than **contrastive learning**!

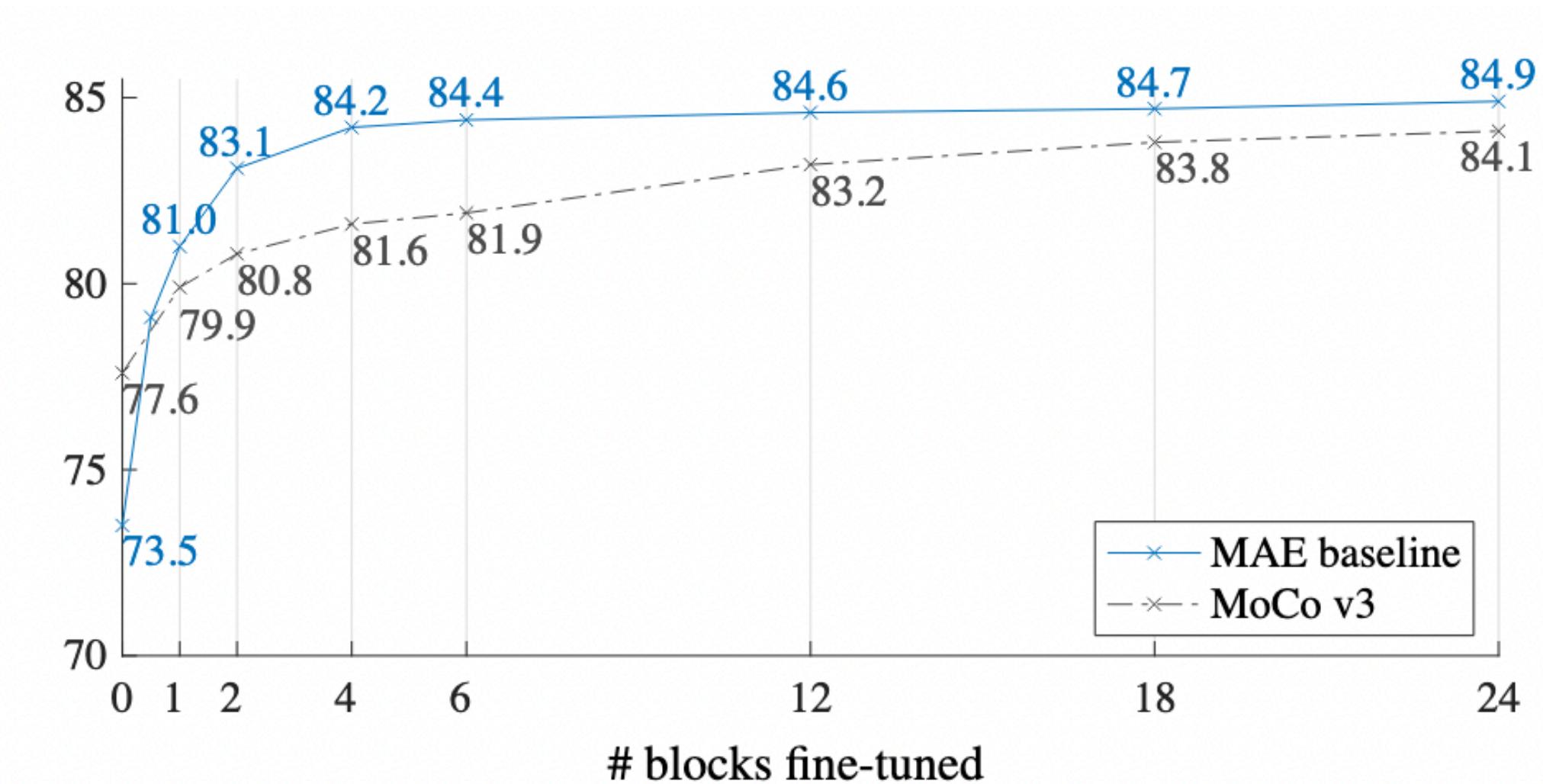
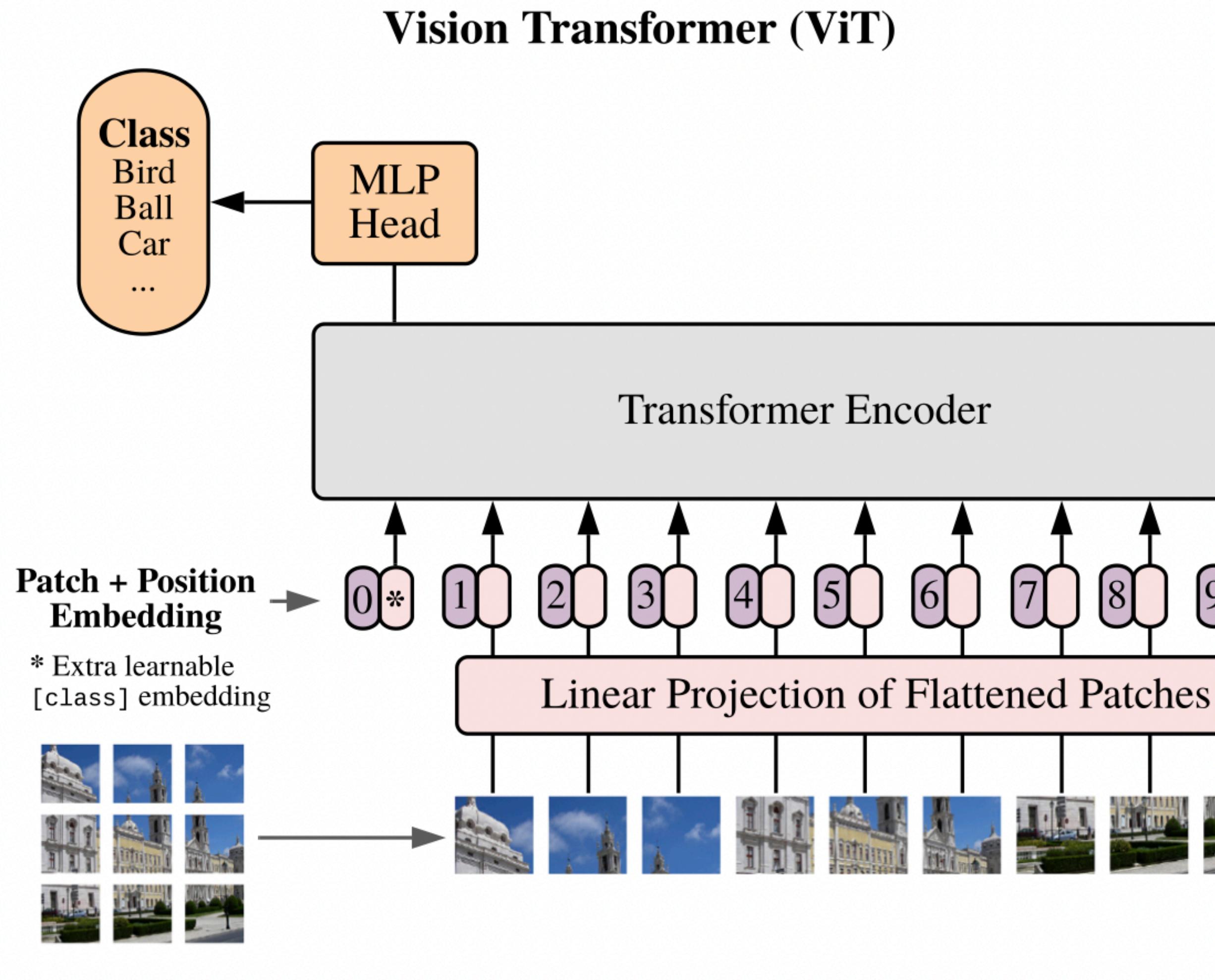


Figure 9. **Partial fine-tuning** results of ViT-L w.r.t. the number of fine-tuned Transformer blocks under the default settings from Table 1. Tuning 0 blocks is linear probing; 24 is full fine-tuning. Our MAE representations are less linearly separable, but are consistently better than MoCo v3 if one or more blocks are tuned.

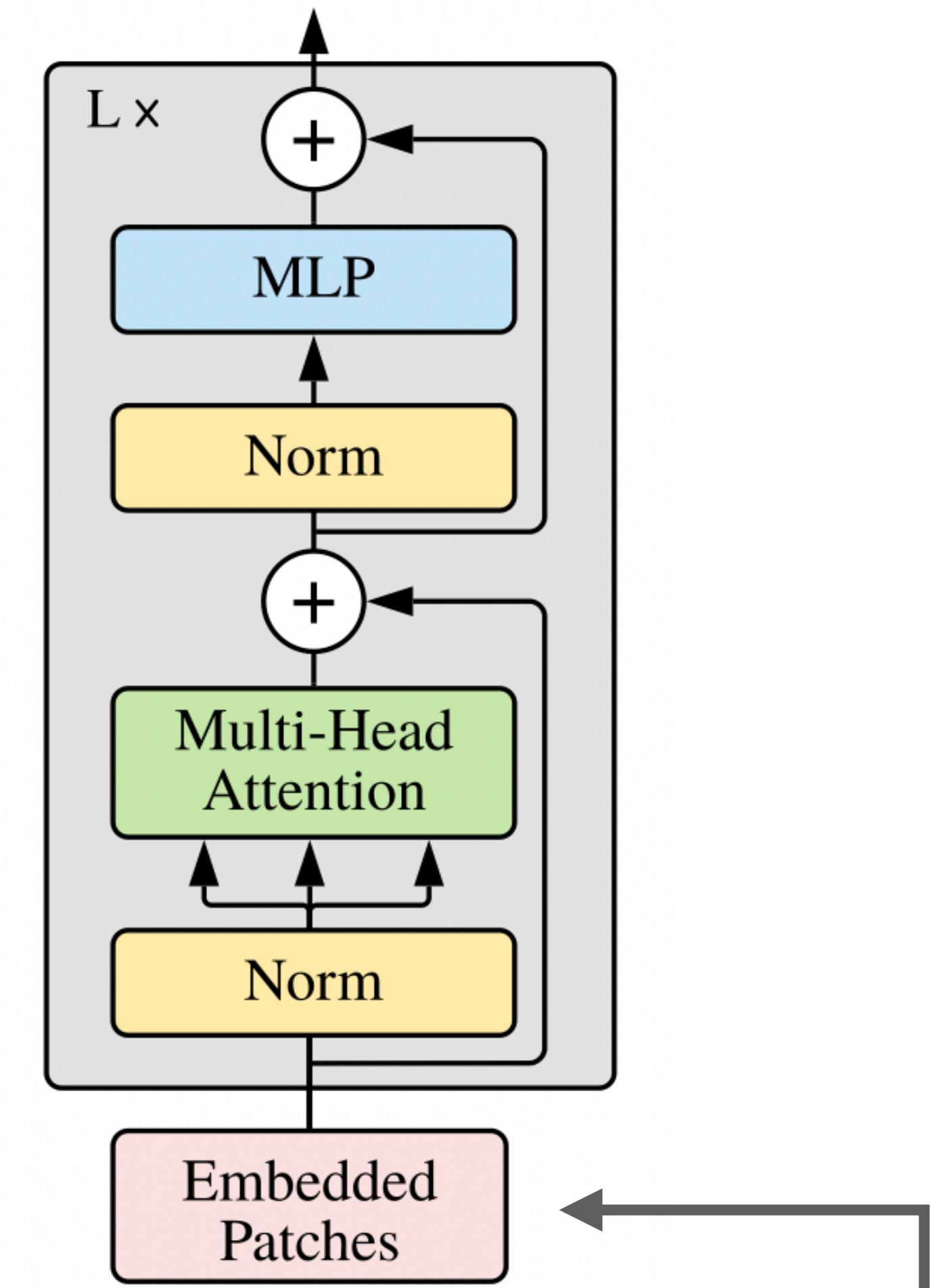
A (very quick) overview of Transformers



A (very quick) overview of Transformers

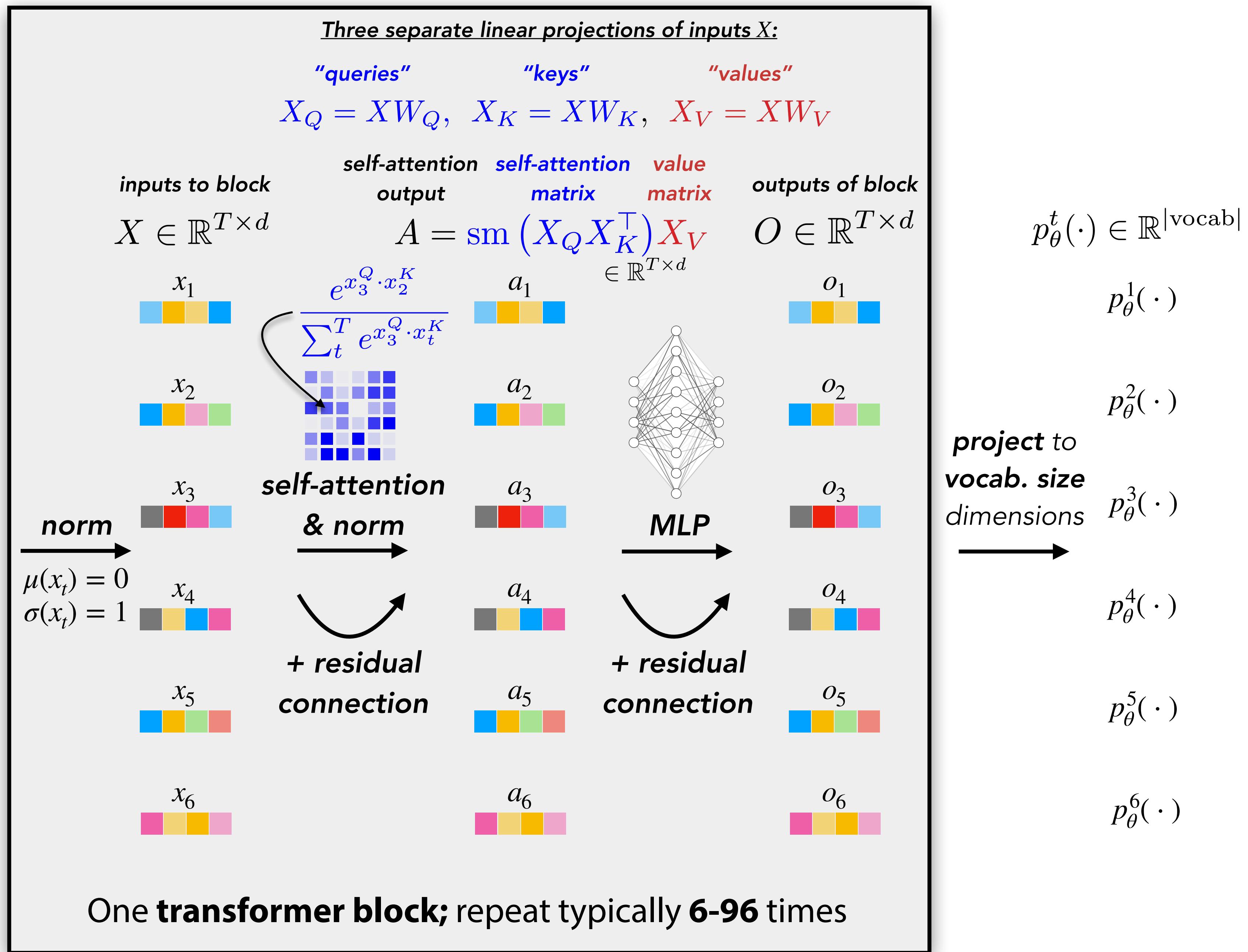
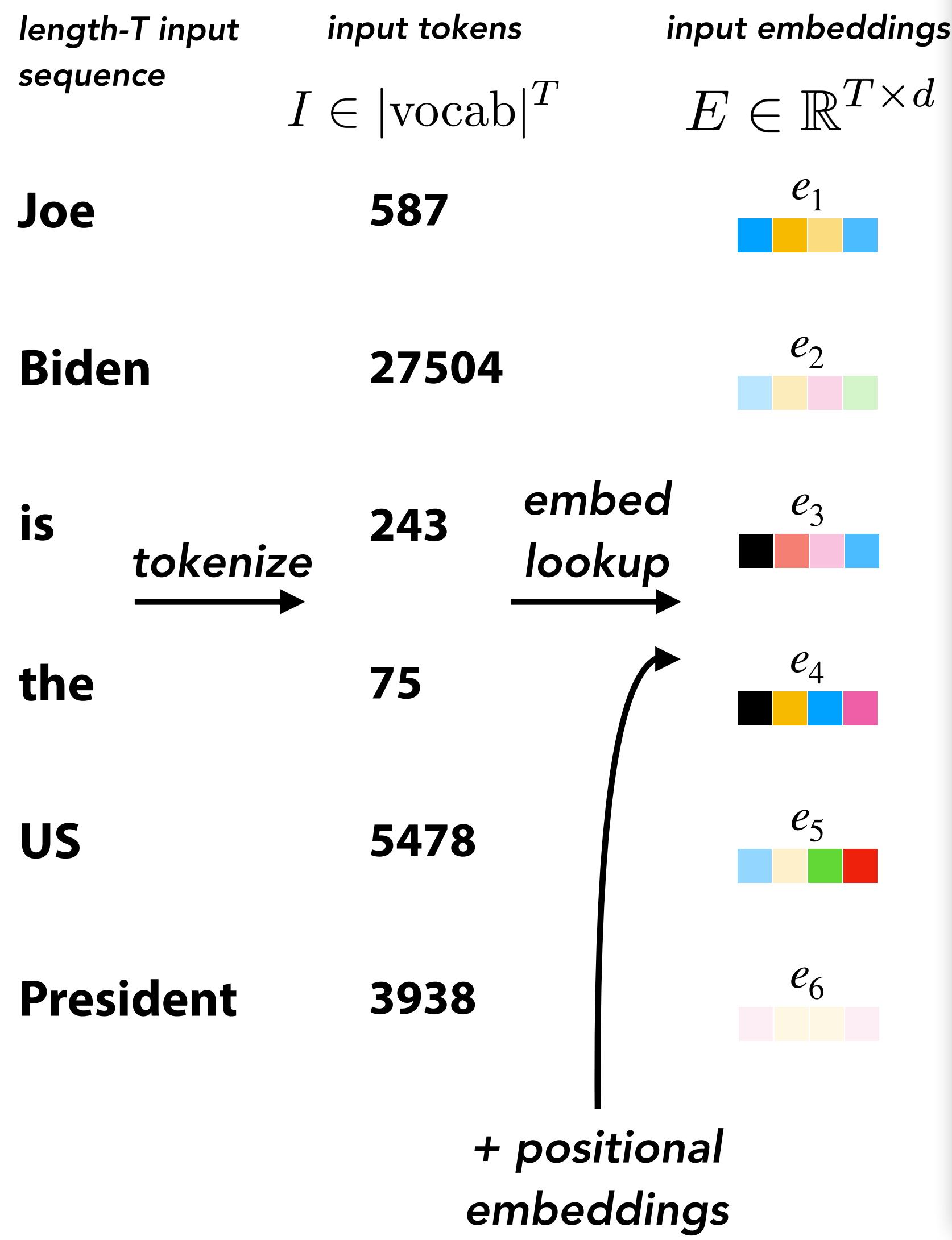


Transformer Encoder

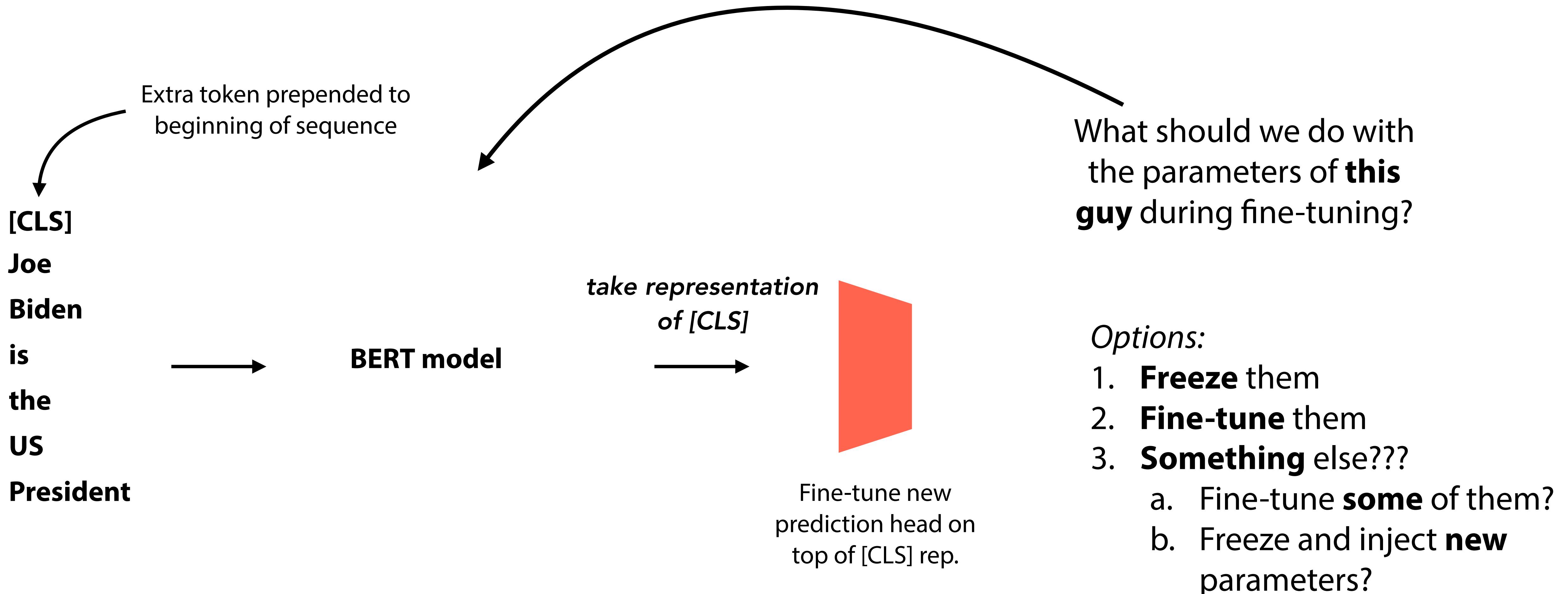


The ~only difference between Transformers for vision/language/RL/molecules/etc. is what we do for this initial **embedding step** —

Transformers in a bit more detail



So... how do we pre-train fine-tune Transformers?



LoRA: Low-rank adaptation of language models (Hu et al., 2021)

What if we just want to fine-tune our model... “a little bit”?

What does “a little bit” even mean? <discuss>

1. Preserve the **knowledge** in the **pre-trained model** (to avoid overfitting)
2. Avoid needing to store a **new version of every single** parameter in the model (to save space)

LoRA: Low-rank adaptation of language models (Hu et al., 2021)

What if we just want to fine-tune our model... "a little bit"?

What does "a little bit" even mean? <discuss>

Associative [key-value] memory view of linear transform (**Kohonen, 1972**)

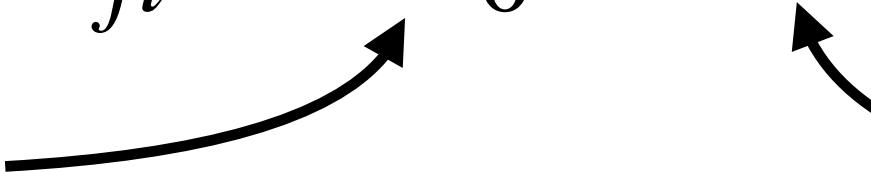
Consider the **linear transform**, the building block of NNs & Transformers

$$W = \sum_r v_r u_r^\top \quad \text{For rank-}r \text{ matrix, we have this decomposition (with orthogonal } u_r \text{ by SVD)}$$

Therefore, $Wx = \left(\sum_r v_r u_r^\top \right) x = \sum_r v_r (u_r^\top x) \rightarrow Wx \text{ produces a sum over the 'memories' } v_r \text{ weighted by the relevance } u_r^\top x \text{ (each } u_r \text{ is a 'key')}$

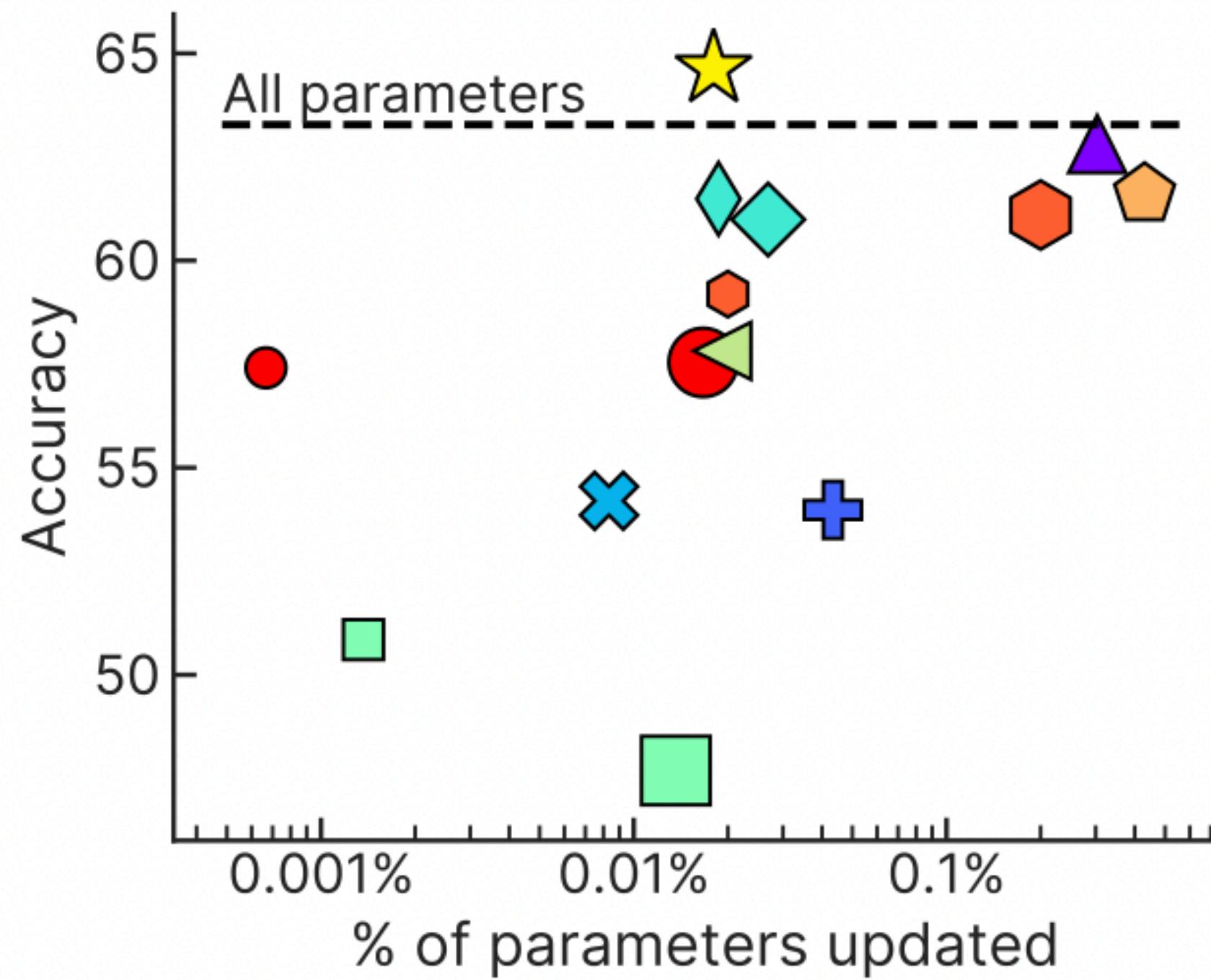
"A little bit" means **only add a few memories** → only make a **low-rank** change to W

LoRA: $W_{ft} = W_0 + AB^\top, \quad A, B \in \mathbb{R}^{d \times p}$

pre-trained weights (frozen) 

new low-rank residual (fine-tuned)
 AB^\top should be zero-initialized (how?)

(Many) other approaches to “lightweight” fine-tuning

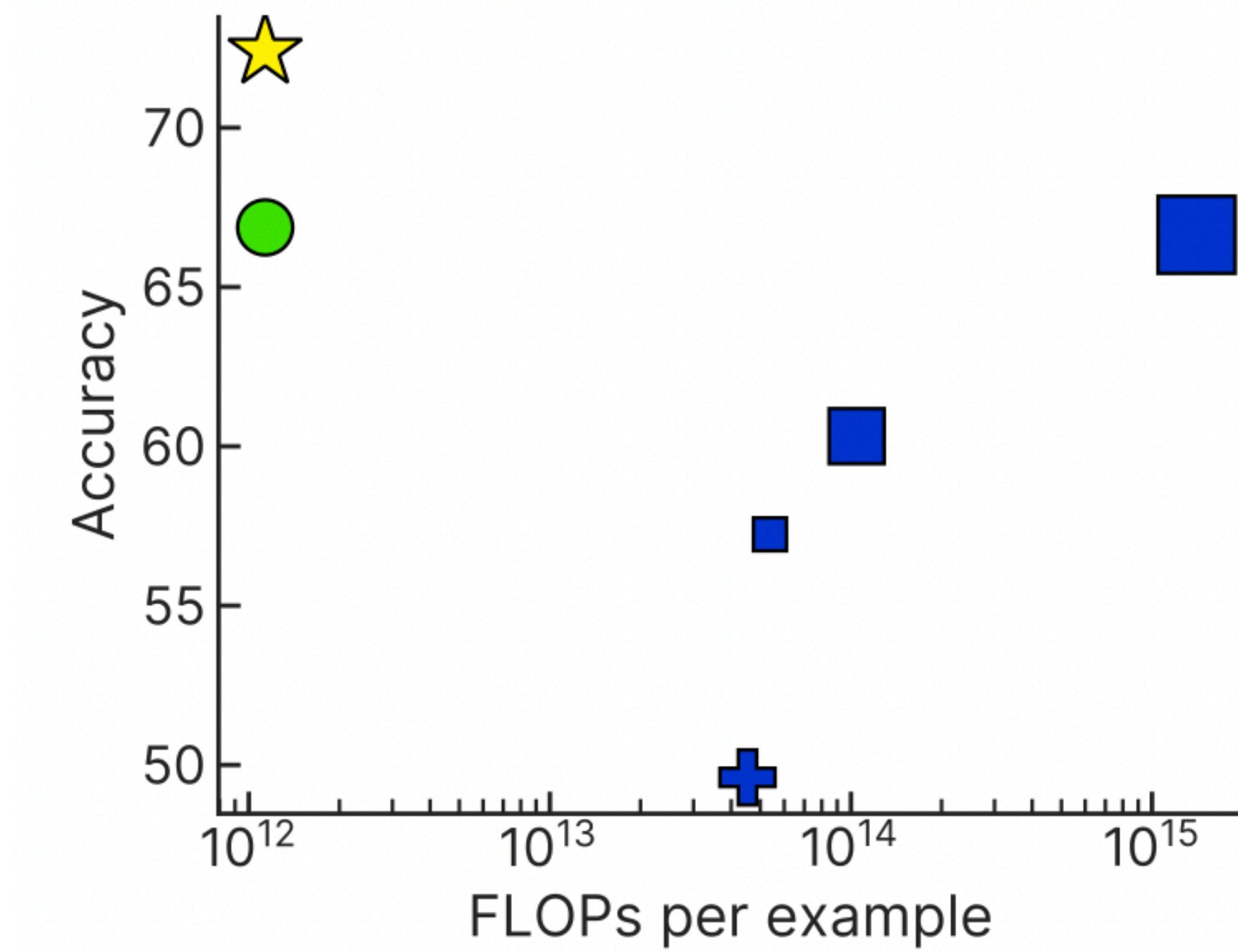


Legend:

- (IA)³
- LoRA
- BitFit
- Layer Norm
- Compacter
- Compacter++

Legend:

- Prompt Tuning
- Prefix Tuning
- Adapter
- FISH Mask
- Intrinsic SAID



Legend:

- T-Few
- T0
- T5+LM

Legend:

- GPT-3 6.7B
- GPT-3 13B
- GPT-3 175B

When “few-shot” means ~20-70, lightweight fine-tuning (**T-Few**) can outperform in-context learning in **much** larger models!

T-Few; Lu, Tam, Muqeeth, et al. (2022)

You will compare fine-tuning and in-context learning in HW3!

Plan for Today

Recap

- Problem formulation
- Contrastive learning

Reconstruction-based unsupervised pre-training

- Why reconstruction?
- Autoencoders
- Masked autoencoders: BERT, MAE
- **Autoregressive models:** GPT, Flamingo

Striving for simplicity: autoregressive models

(recall from the black-box meta-learning lecture!)

What are some **downsides** of masked autoencoders?

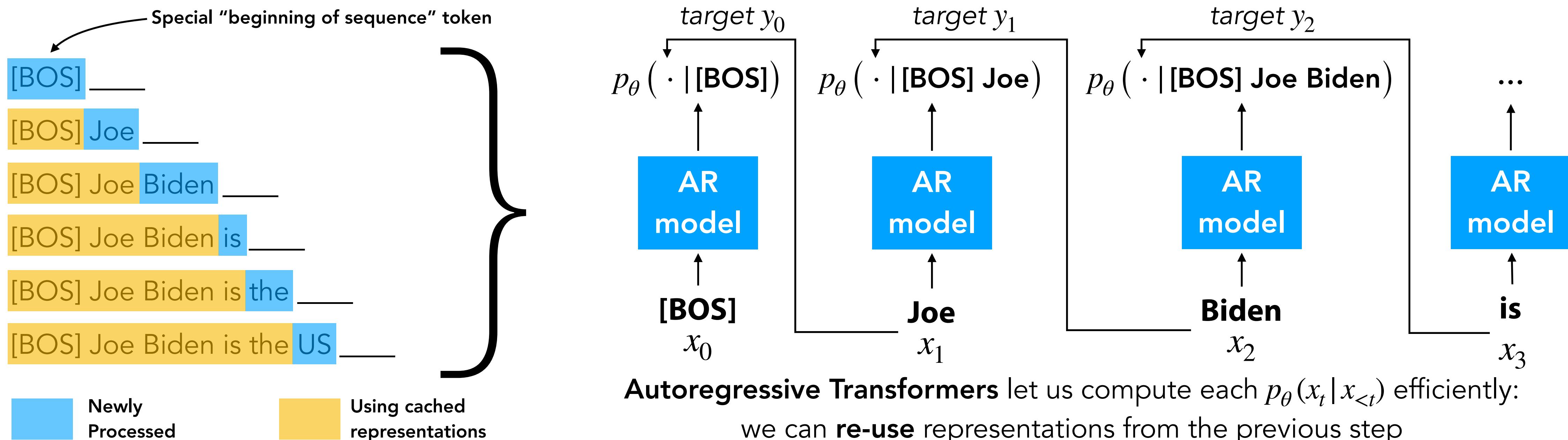
1. Need to pick **mask**
2. Only using ~15% of the example for training
3. Difficult to sample from

Instead of masking a **random subset**, what if we just predict the next word/pixel/token?



No need to pick a masking strategy; **mask every token!**

Simply learn $p_{\theta}(x_t | x_{<t})$, probability of the **next token** given the **previous tokens**



Autoregressive Transformers are *everywhere* these days

...for vision too!
...and RL/decision-making!
...and vision + language!

Improving Language Understanding by Generative Pre-Training

Language Models are Unsupervised Multitask Learners

Language Models are Few-Shot Learners

Megatron-LM: Training Multi-Billion Parameter Language Models Using Model Parallelism

OPT: Open Pre-trained Transformer Language Models

Announcing GPT-NeoX-20B

Announcing GPT-NeoX-20B, a 20 billion parameter model trained in collaboration with CoreWeave.

February 2, 2022 · Connor Leahy

As of February 9, 2022, GPT-NeoX-20B checkpoints are available for [download from The Eye under Apache 2.0](#). More in-depth information on GPT-NeoX-20B can be found in the [associated technical report on arXiv](#).

Looking for a demo? Try GPT-NeoX-20B via CoreWeave and Anllan's inference service, [GooseAI](#)!



28-04-2022

Flamingo: a Visual Language Model for Few-Shot Learning

Jean-Baptiste Alayrac^{*,‡}, Jeff Donahue^{*}, Pauline Luc^{*}, Antoine Miech^{*}, Iain Barr[†], Yana Hasson[†], Karel Lenc[†], Arthur Mensch[†], Katie Millican[†], Malcolm Reynolds[†], Roman Ring[†], Eliza Rutherford[†], Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, Karen Simonyan^{*,‡}

*Equal contributions, ordered alphabetically, †Equal contributions, ordered alphabetically, ‡Equal senior contributions

Case study: Flamingo

Flamingo: a Visual Language Model for Few-Shot Learning

Jean-Baptiste Alayrac ^{*,‡}, Jeff Donahue ^{*}, Pauline Luc ^{*}, Antoine Miech ^{*}, Iain Barr [†], Yana Hasson [†], Karel Lenc [†], Arthur Mensch [†], Katie Millican [†], Malcolm Reynolds [†], Roman Ring [†], Eliza Rutherford [†], Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, Karen Simonyan ^{*,‡}

^{*}Equal contributions, ordered alphabetically, [†]Equal contributions, ordered alphabetically, [‡]Equal senior contributions

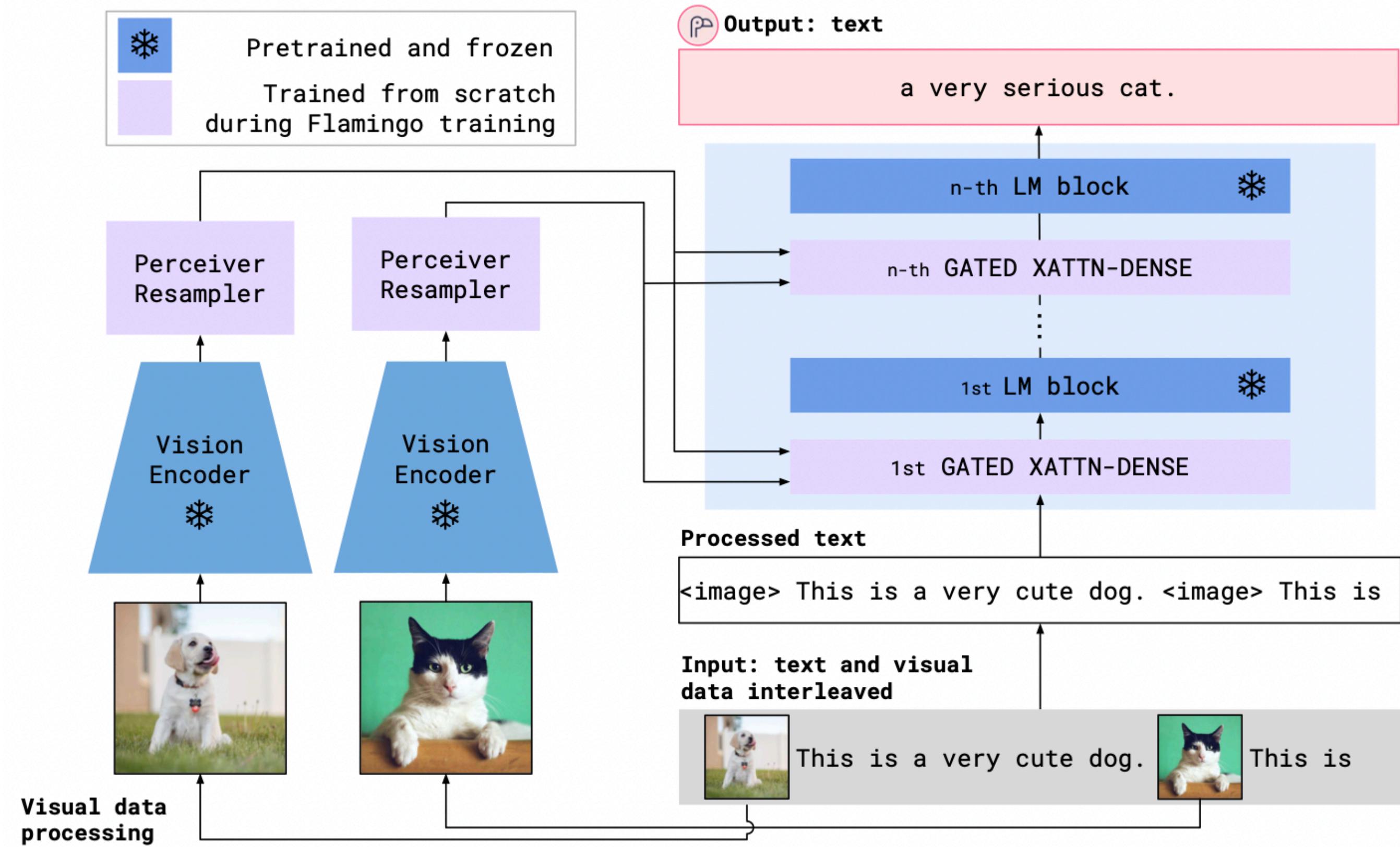
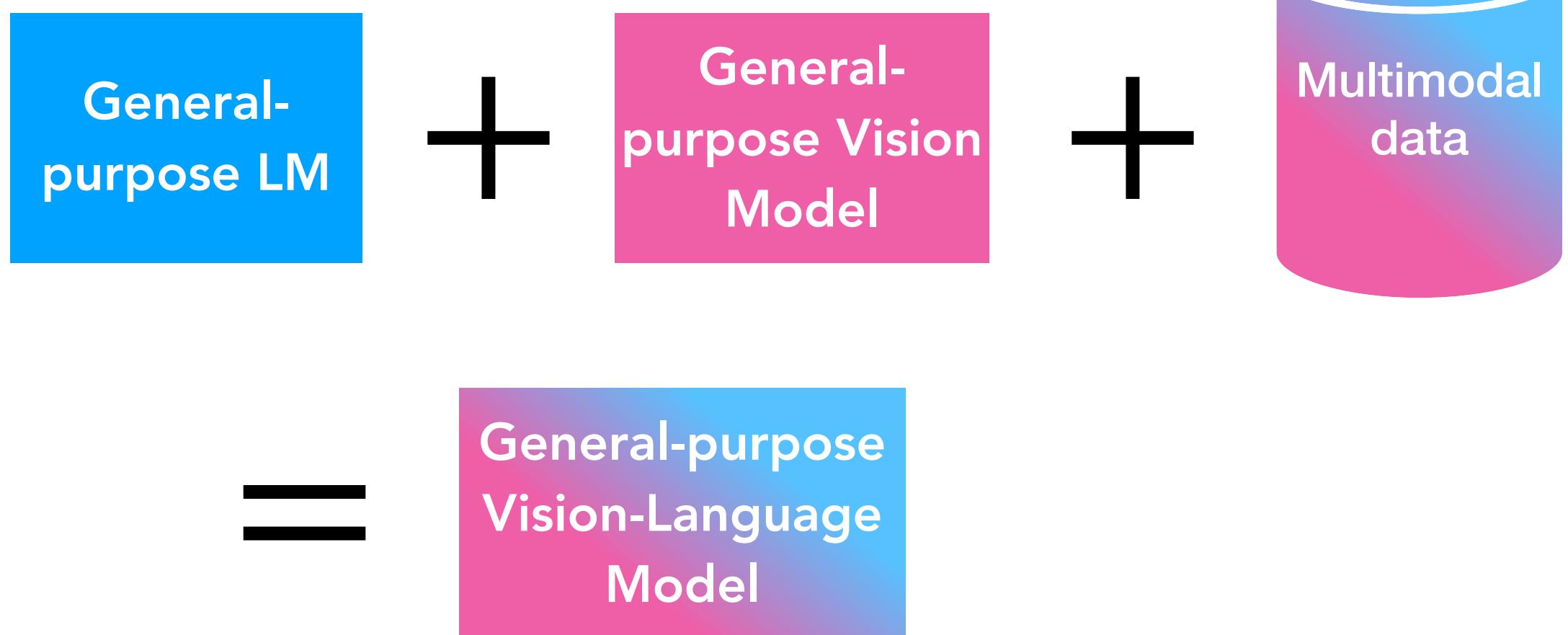
How would you build a multimodal autoregressive model? From scratch? (NO)

[so far] Fine-tuning to **specialize**:



Flamingo

Fine-tuning to **combine models**:



Case study: Flamingo

Flamingo: a Visual Language Model for Few-Shot Learning

Jean-Baptiste Alayrac ^{*,‡}, Jeff Donahue ^{*}, Pauline Luc ^{*}, Antoine Miech ^{*}, Iain Barr [†], Yana Hasson [†], Karel Lenc [†], Arthur Mensch [†], Katie Millican [†], Malcolm Reynolds [†], Roman Ring [†], Eliza Rutherford [†], Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, Karen Simonyan ^{*,‡}

^{*}Equal contributions, ordered alphabetically, [†]Equal contributions, ordered alphabetically, [‡]Equal senior contributions



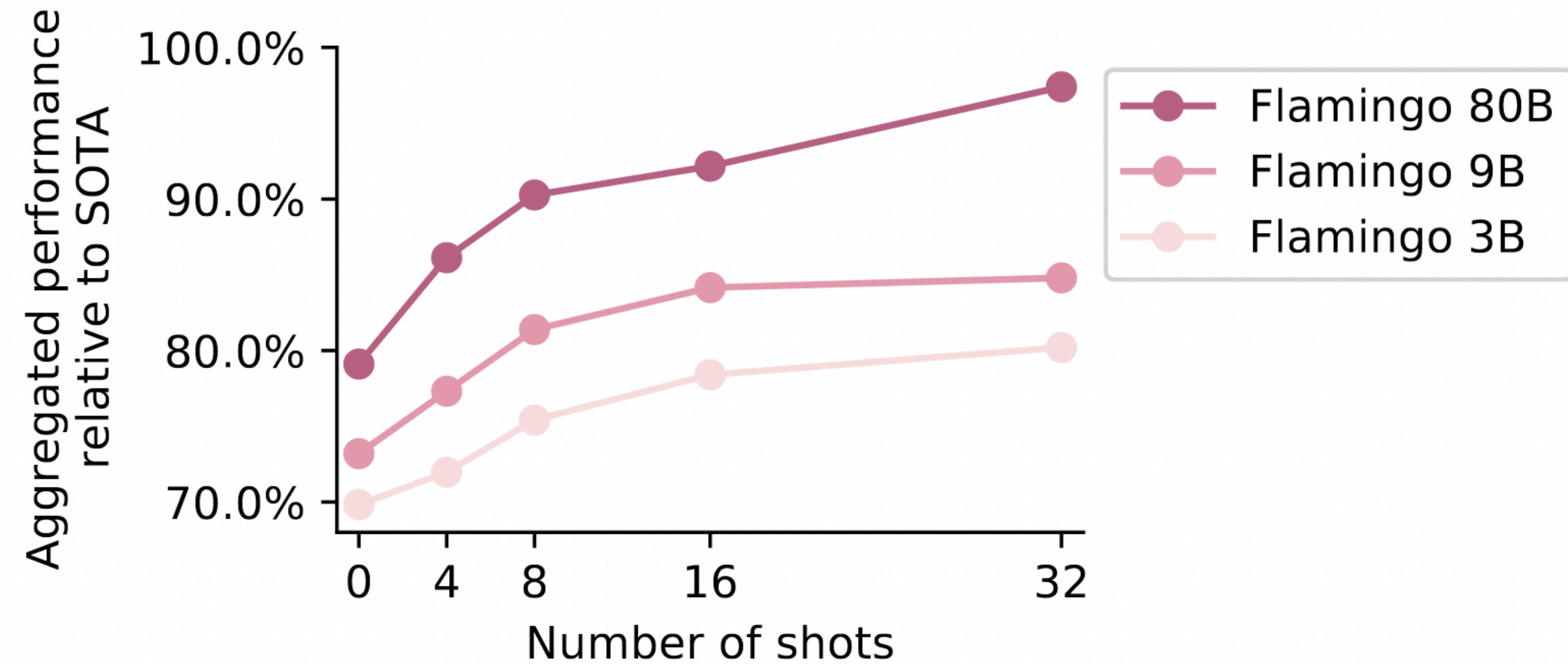
In-context few-shot learning on sequences that freely mix **text** and **images!** Enables few-shot captioning, visual question-answering, etc.

Case study: Flamingo

Flamingo: a Visual Language Model
for Few-Shot Learning

Jean-Baptiste Alayrac ^{*,‡}, Jeff Donahue ^{*}, Pauline Luc ^{*}, Antoine Miech ^{*}, Iain Barr [†], Yana Hasson [†],
Karel Lenc [†], Arthur Mensch [†], Katie Millican [†], Malcolm Reynolds [†], Roman Ring [†], Eliza Rutherford [†],
Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick,
Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski,
Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, Karen Simonyan ^{*,‡}

^{*}Equal contributions, ordered alphabetically, [†]Equal contributions, ordered alphabetically, [‡]Equal senior contributions



Few-shot Flamingo \approx Non-Few-shot state of the art!

Are AR models really **different** from masked autoencoders?

General recipe for training masked autoencoder f_θ :

1. Choose **distance function** $d(\cdot, \cdot) \rightarrow \mathbb{R}$

2. For **train batch** examples x_i :

A. Sample $\tilde{x}_i, y_i \sim \text{mask}(x_i)$

B. Make prediction $\hat{y}_i = f_\theta(\tilde{x}_i)$

C. Compute loss $= d(y_i, \hat{y}_i)$

Masked autoencoder:

$\tilde{x} :$
Joe

$y :$

$\tilde{x} :$
Joe

$y :$

<mask>

Biden

Biden

is

is

the

the

<mask>

US

US

President

President

**AR models are just masked AEs
with a special choice of mask**

Summary of today

1. *Intuition for autoencoders (AEs): “A good **representation** lets us **reconstruct** the input”*
2. ***Masked AEs** learn to restore a **partially-deleted** input & help avoid degeneracies in unmasked AEs*
3. ***State of the art** in pre-training for few-shot learning in **language** & **vision***
4. ***Autoregressive** models (e.g., GPT-3) are **special case** of masked AEs; give a generative model for free at some cost to fine-tuning performance*

Contrastive Learning vs AEs vs Masked AEs

Contrastive learning:

- + Learns very high-quality representations
- + Don't need as large a model
- Need to select negatives carefully*
- Generally needs larger batch size*
- Cross-example dependencies can make implementation more difficult

* new methods are addressing these downsides but are more difficult to interpret/analyze

(Bottlenecked) Autoencoders: Masked autoencoders:

- + Simple to implement
- + No need to select pos/neg pairs; just $d(x, \hat{x})$
- Generally need a larger model
- Need to design a bottleneck
- (Comparatively) poor few-shot performance
- Not generally used in practice
- + **Few-shot** performance as good or better than contrastive
- + **AR special case** gives generative model for free
- **Raw representations** (without fine-tuning) still can be lower quality than contrastive