

# Unsupervised pre-training for few-shot learning, vol. 2: reconstruction-based methods

CS 330

# Logistics

Project proposal due TODAY!

Homework 2 due Monday, October 24

**Kyle's** office hours are hybrid going forward (see Ed for details)

Azure invites have been re-sent - **you have one week to accept!**

**You will need Azure for HW3, so do this today!**

# Plan for Today

## *Recap*

- Problem formulation
- Contrastive learning

## *Reconstruction-based unsupervised pre-training*

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

}

Topic of Homework 3!

## 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**

# Plan for Today

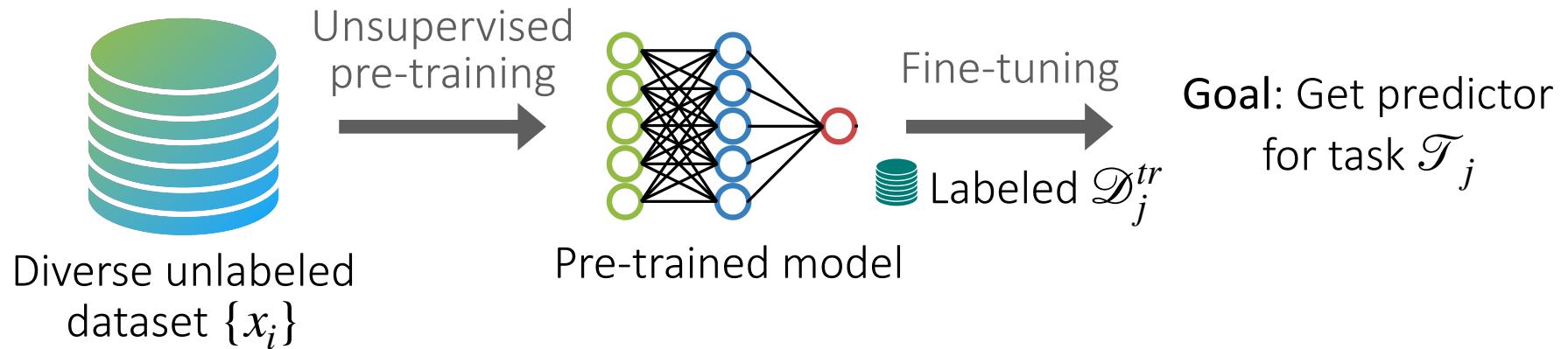
## *Recap*

- **Problem formulation**
- Contrastive learning

## *Reconstruction-based unsupervised pre-training*

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

# Unsupervised Pre-Training Set-Up



# Plan for Today

## *Recap*

- Problem formulation
- **Contrastive learning**

## *Reconstruction-based unsupervised pre-training*

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

# Key Idea of Contrastive Learning

**Similar examples should have similar representations**

Examples with the same class label



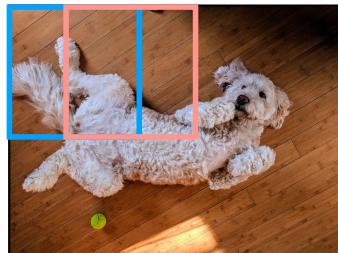
(Requires labels, related to Siamese nets, ProtoNets)

Augmented versions of the example



(flip & crop)

Nearby image patches



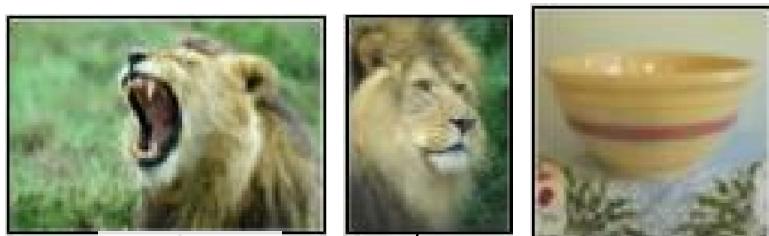
Nearby video frames



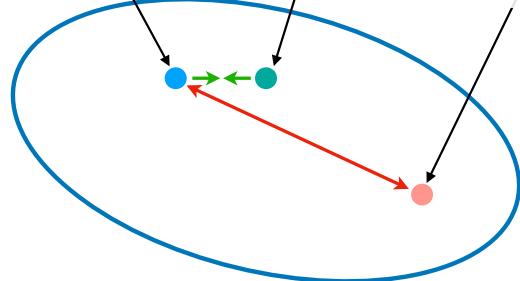
# Contrastive Learning Implementation

**Similar examples** should have **similar representations**

Need to both compare & contrast!



anchor  $x$  positive  $x^+$  negative  $x^-$



V1. Triplet loss:

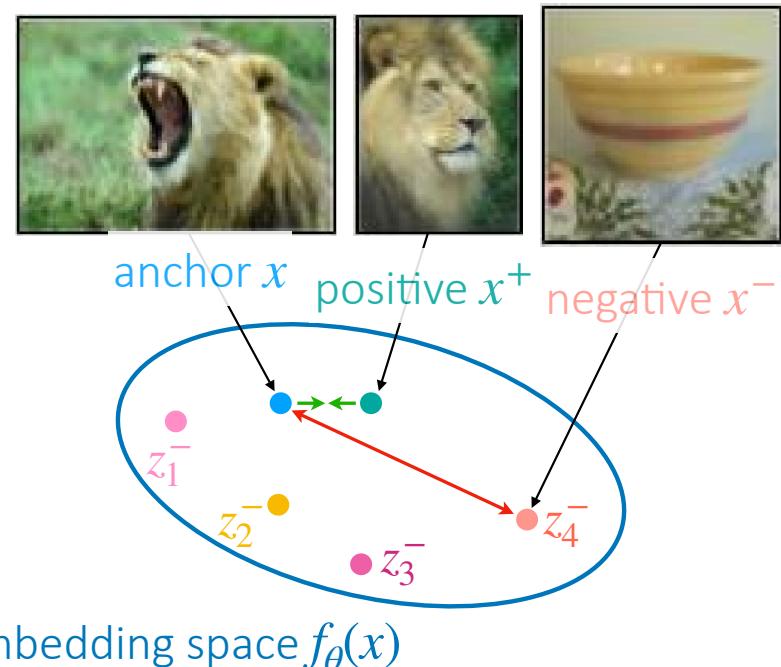
$$\min_{\theta} \sum_{(x,x^+,x^-)} \max (0, \|f_\theta(x) - f_\theta(x^+)\|^2 - \|f_\theta(x) - f_\theta(x^-)\|^2 + \epsilon)$$

Embedding space  $f_\theta(x)$

# 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\***):

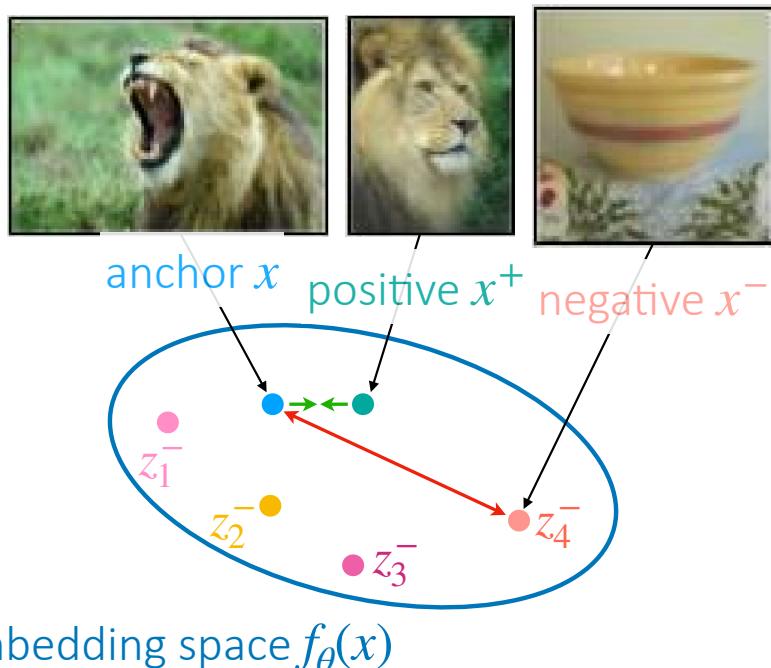
$$\mathcal{L}_{\text{N-way}}(\theta) = - \sum_z \log \frac{\exp(-d(z, z^+))}{\sum_i \exp(-d(z, z_i^-))}$$

\*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:

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

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

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

↳ Useful bcs. the  $\mathcal{L}$  is guaranteed to be  $0 \leq \mathcal{L} \leq 1$ .

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**  
Ⓐ Condition!

# Plan for Today

## *Recap*

- Problem formulation
- Contrastive learning

## *Reconstruction-based unsupervised pre-training*

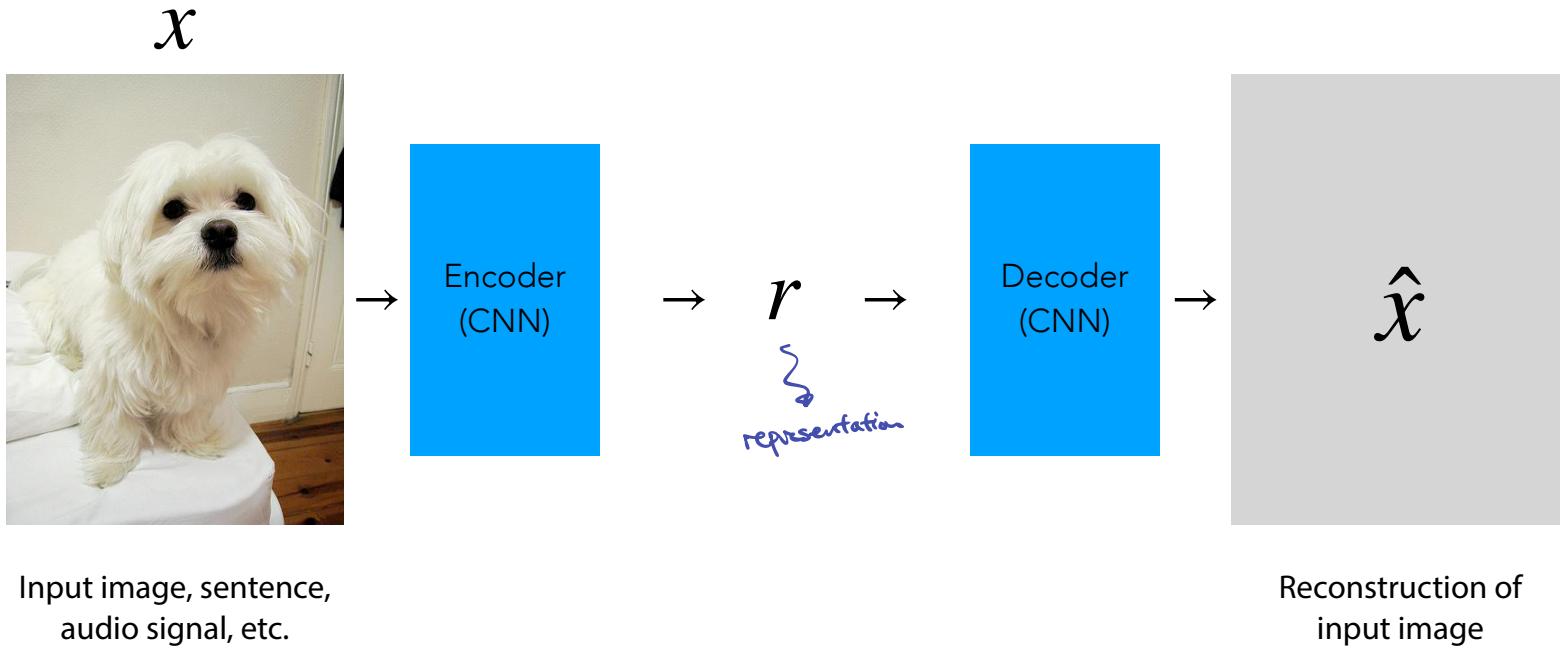
- **Why reconstruction?**
- Autoencoders
- Masked autoencoders: BERT, MAE
- Autoregressive models: GPT, Flamingo
- Emergent behaviors in large models

# 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**  $r$  →  $x \approx \hat{x}$

# Plan for Today

## *Recap*

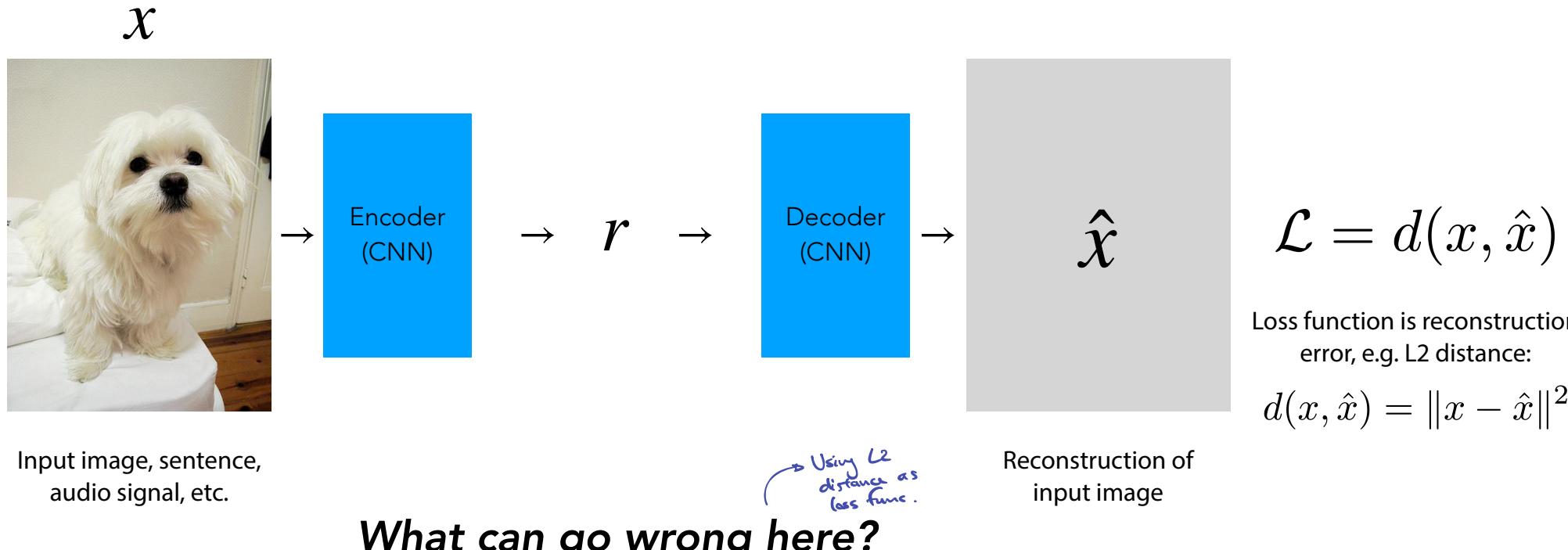
- Problem formulation
- Contrastive learning

## *Reconstruction-based unsupervised pre-training*

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

# Autoencoders: a first attempt

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

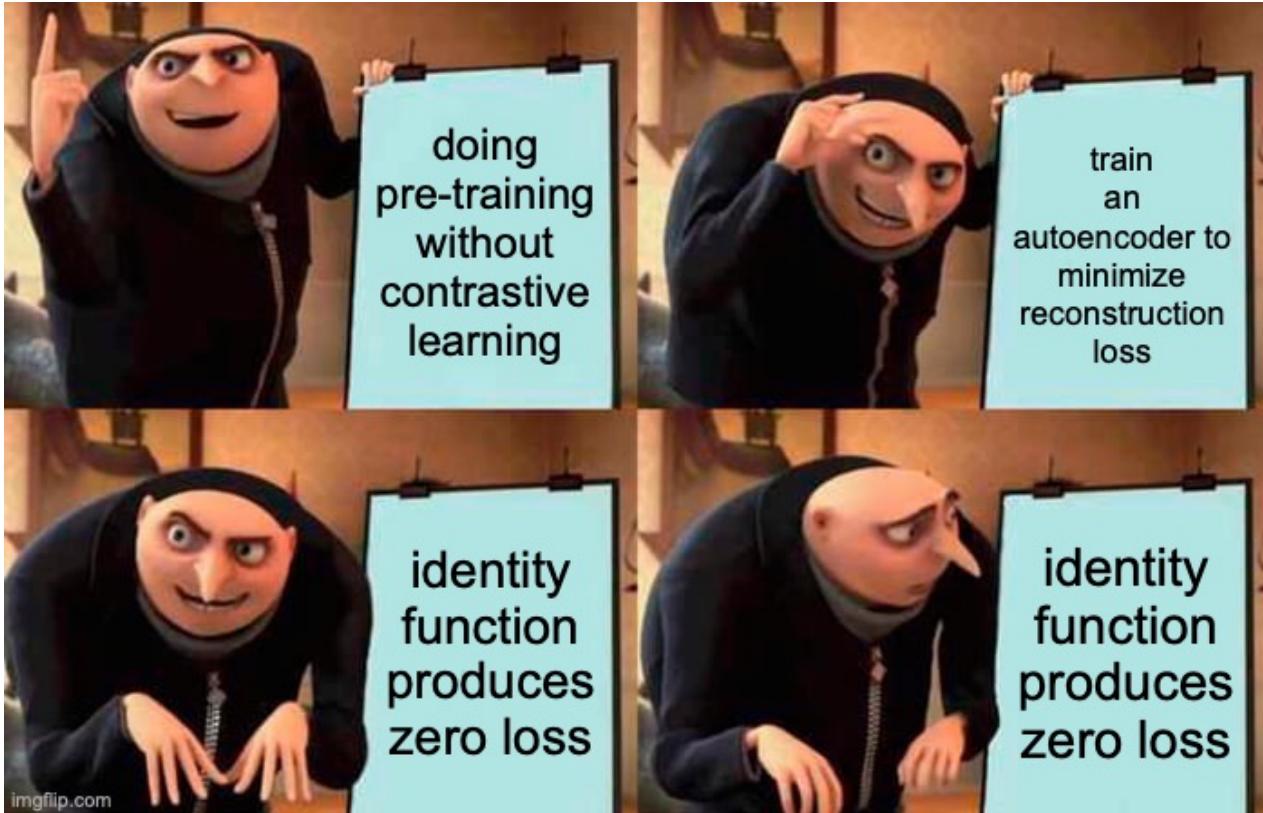


Is the **identity function** a good encoder/decoder?

↳ It'd actually be a perfect encoder/decoder for this loss (L2 dist)

→ How do we fix this? → Constrain the dimensionality on  $r$ ! (Bottleneck)

# Autoencoders: a first attempt



***How to fix???***

# Autoencoders: adding a bottleneck

$x$



Input image, sentence,  
audio signal, etc.



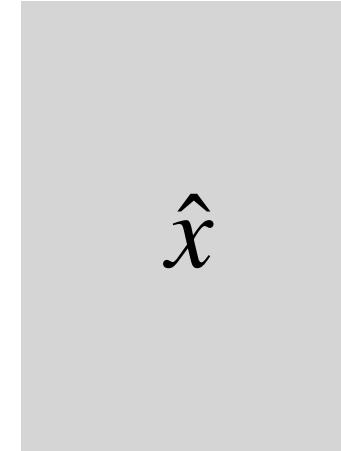
Compact, latent  
representation of input image

$r$



Decoder (CNN)

$\hat{x}$



Reconstruction of  
input image

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

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

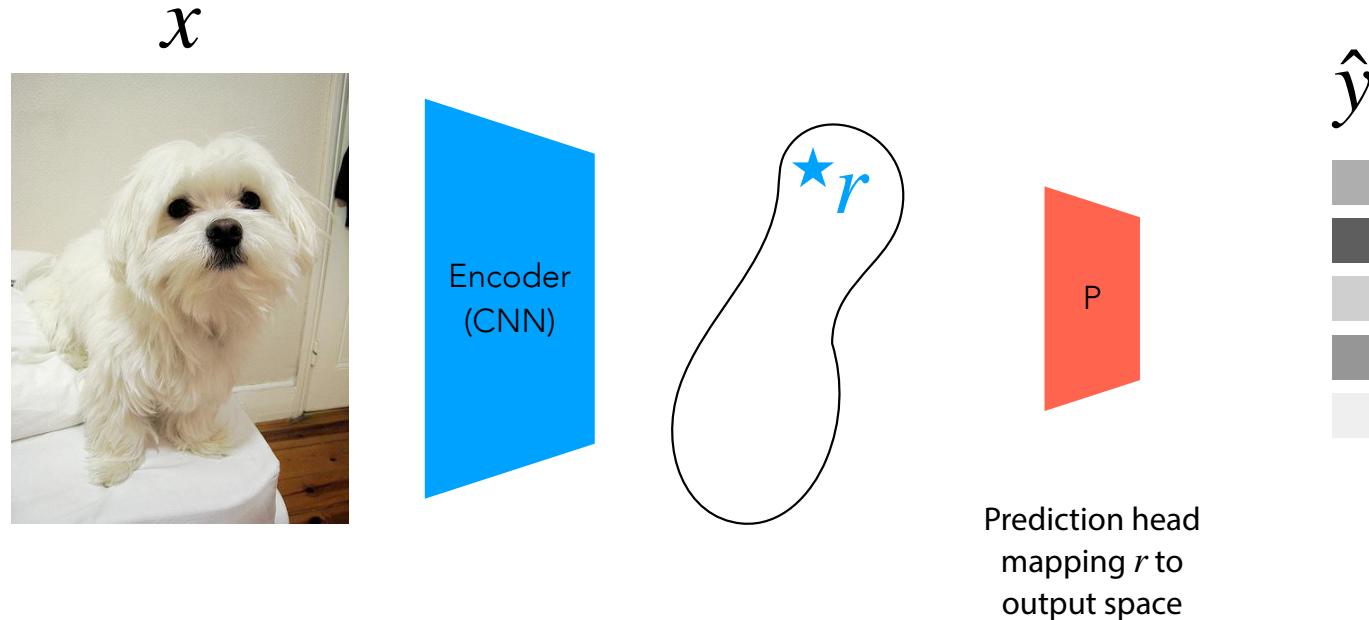
$$d(x, \hat{x}) = \|x - \hat{x}\|^2$$

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



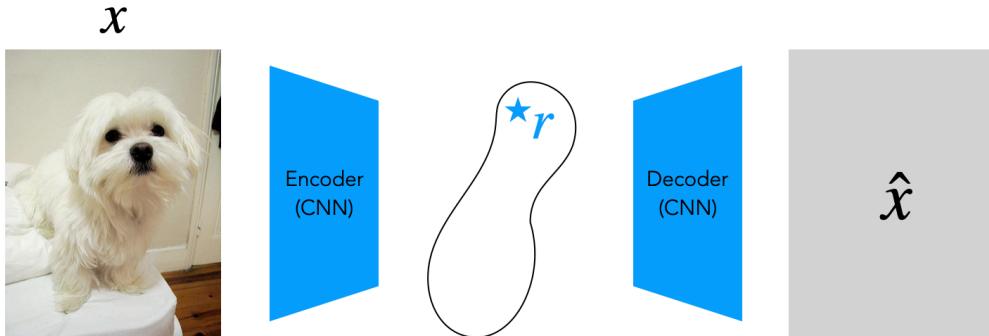
*Hope:* latent dimensions are forced to represent  
**high-level** concepts that **generalize** to other tasks

# 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**

In practice it fails to give you good low-dimensional representation

generally difficult to invert

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

Constrain the no. of dimensions that can be non-zero

Force the decoder to be relatively weak (so that the repr. is forced to be easier to decode)

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

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

→ masked autoencoders

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

→ to force the encoder to learn sth. less trivial.

# 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**

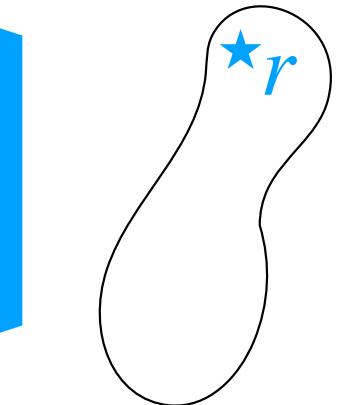
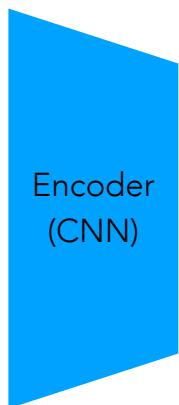
Ultimately, **regular** autoencoders are trying to predict  $x$  from...  $x$  (through  $r$ )

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

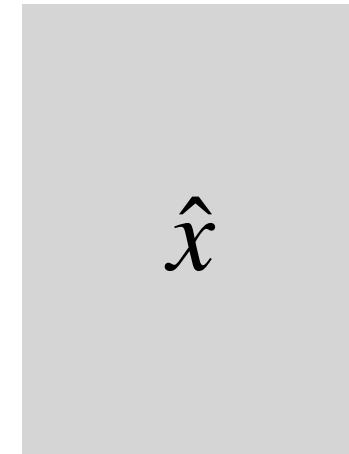
Masked autoencoders use a **more difficult** learning task to encourage  
 $x$  the encoder to extract more meaningful features



Input image, sentence,  
audio signal, etc.



Compact, latent  
representation of input image



Reconstruction of  
input image

## Beyond the bottleneck: **masked autoencoders**

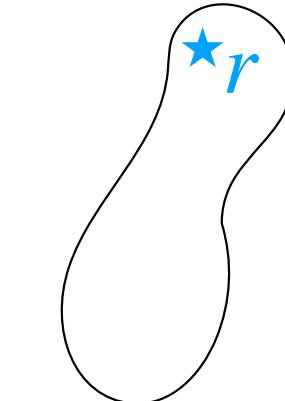
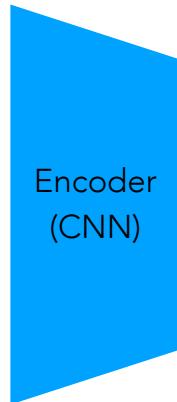
Ultimately, regular autoencoders are trying to predict  $x$  from...  $x$  (through  $z$ )

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

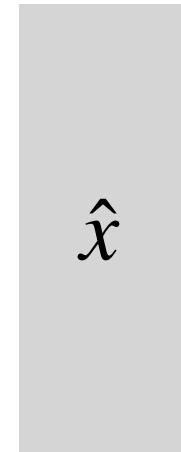
Masked autoencoders use a **more difficult** learning task to encourage  
the encoder to extract more meaningful features



Masked input image



Latent representation



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



$\approx$

Predict the missing part

## Beyond the bottleneck: **masked autoencoders**

General recipe for **pre-training** masked autoencoder  $f_\theta$ :

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

→ select masking function

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

See example below.

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

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

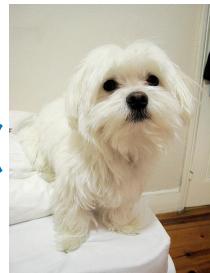
in some cases, the target  $y_i$  may be all of  $x_i$

→ doesn't make a  
huge difference  
in practice

$x_i$

$\tilde{x}_i$

$y_i$



$\text{mask}(\quad)$  =

$x_i$

$\text{mask}(\text{ Joe Biden is the US president }) =$

$\tilde{x}_i$

$y_i$

$\text{Joe } <\text{mask}> \text{ is the US } <\text{mask}>, \{ \text{Biden}; \text{president} \}$

$f_\theta$ : CNN or **Transformer** (stay tuned)

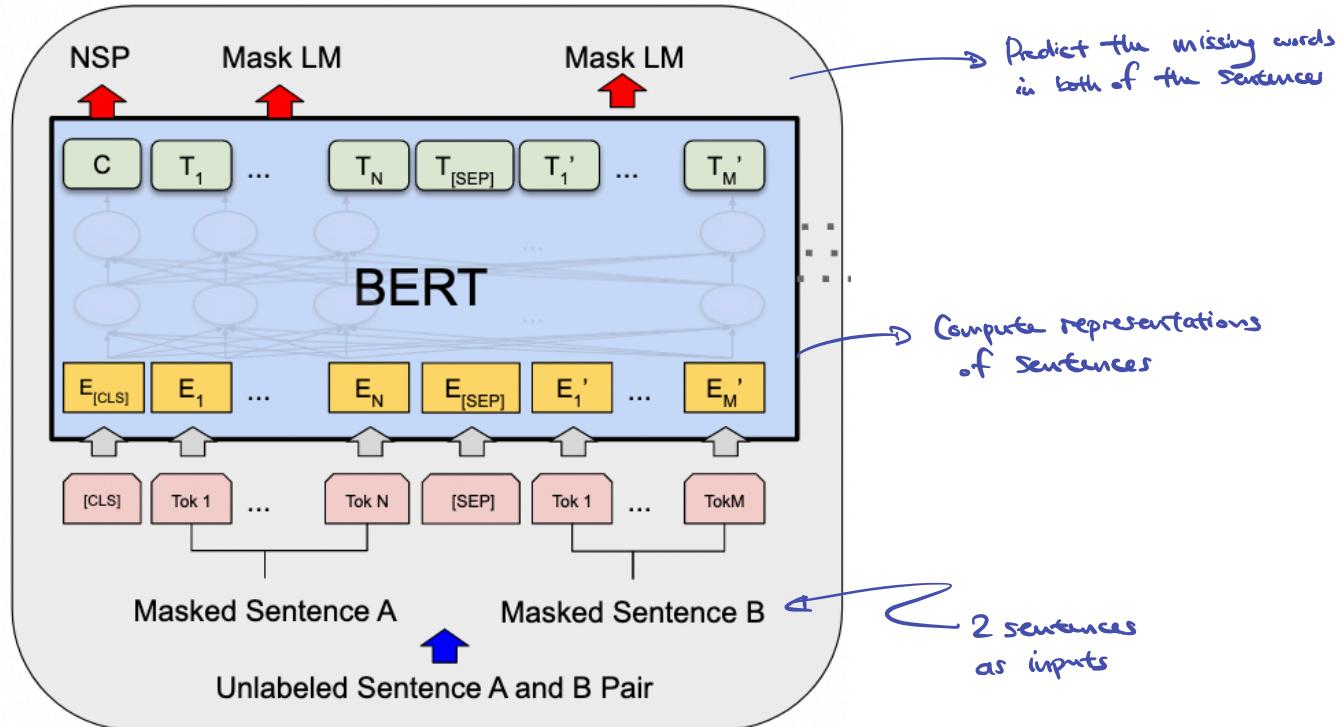
$d(y, \hat{y}) = \|y - \hat{y}\|^2$  → e.g.) L2 distance

$f_\theta$ : **Transformer** (e.g., **BERT**; stay tuned)

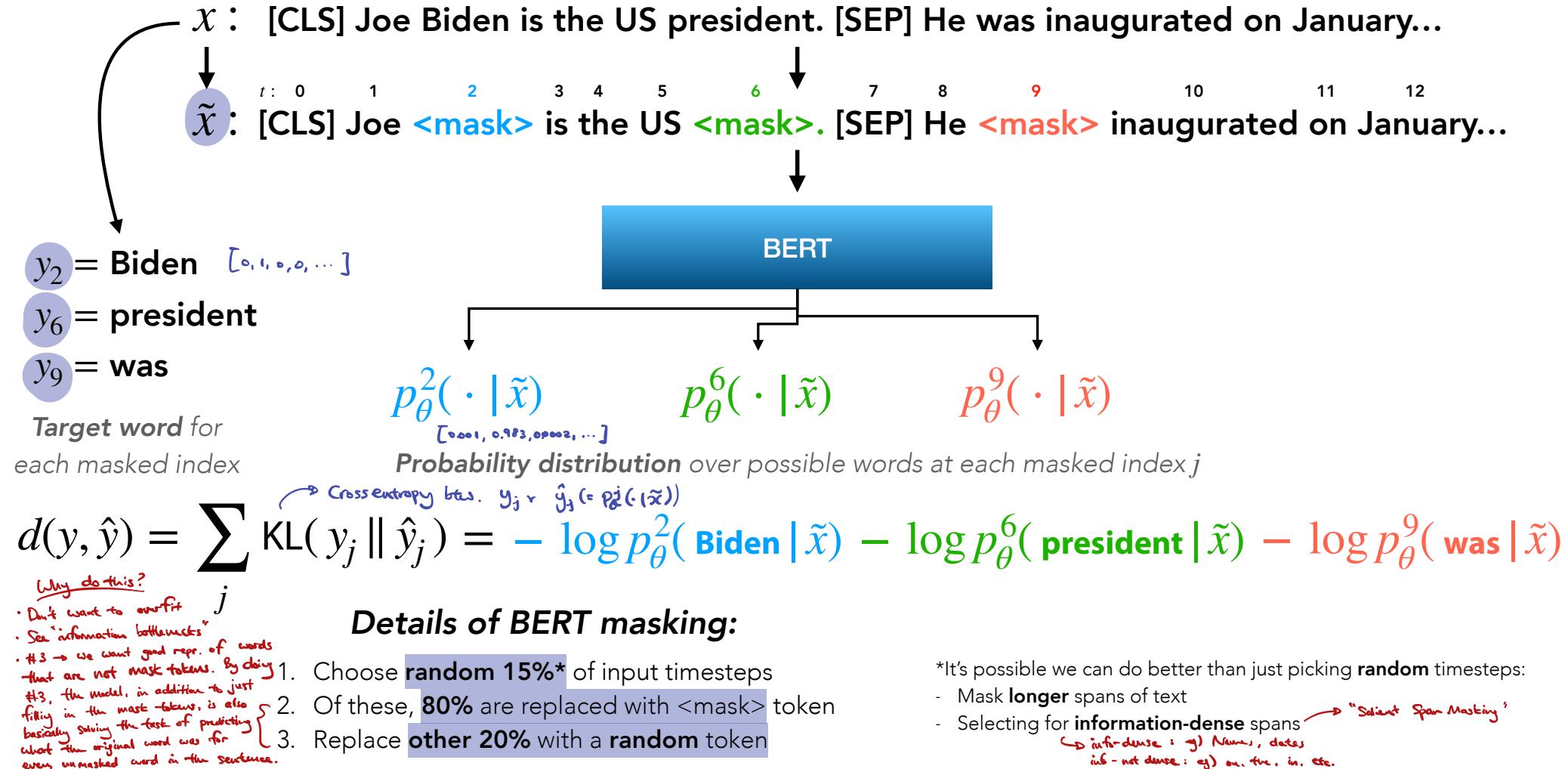
$d(y, \hat{y}) = \text{KL}(y \parallel \hat{y})$  → e.g.) KL-divergence

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

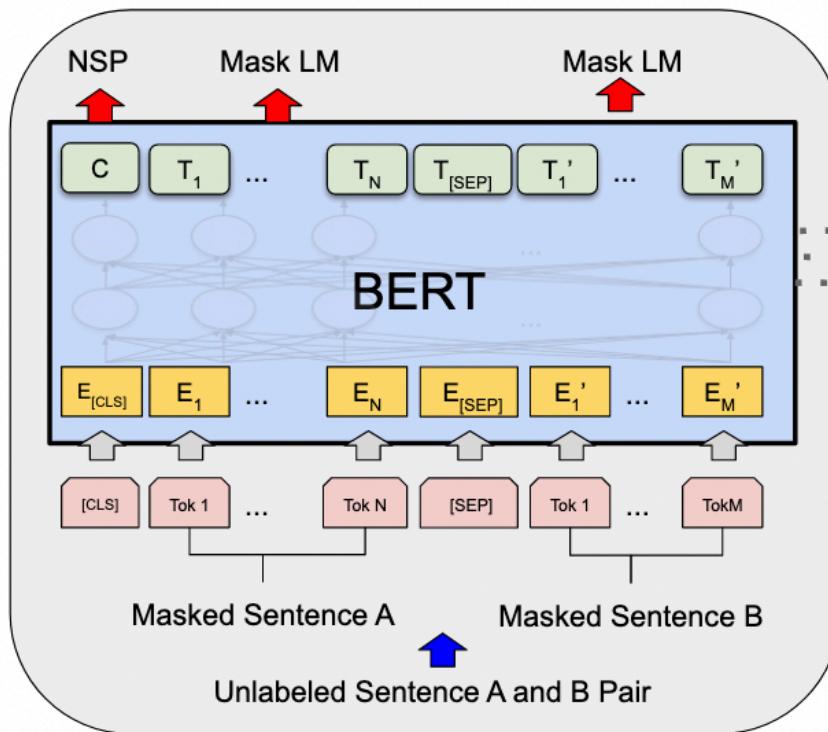
↳ prob most famous



## Case study: BERT as a masked autoencoder

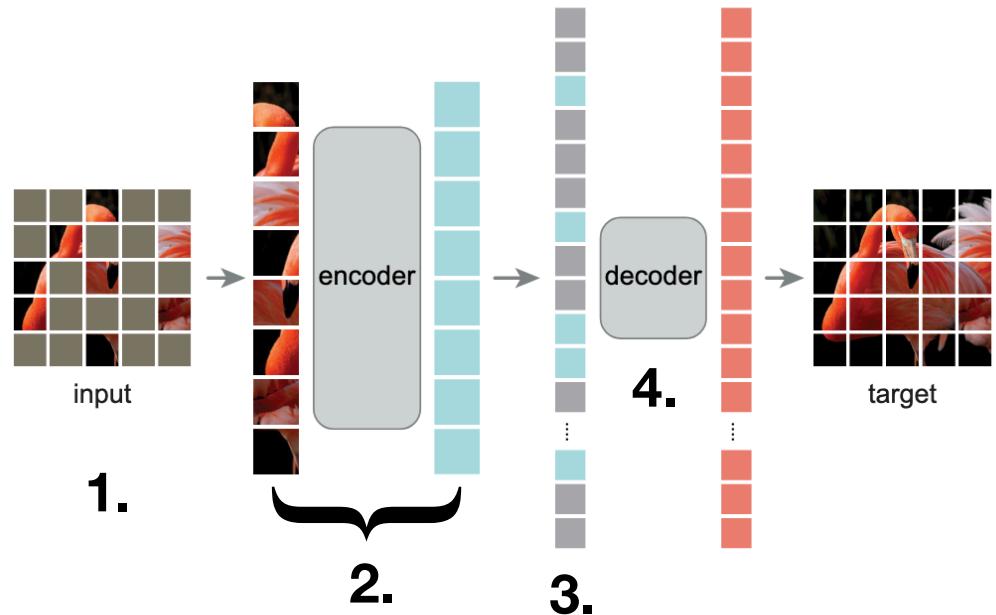


# Masked autoencoders for language: **BERT** (Devlin et al, 2017)



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

Check out  
this paper



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!**

→ get a better feature extractor this way!  
→ Extract even more info from the dataset

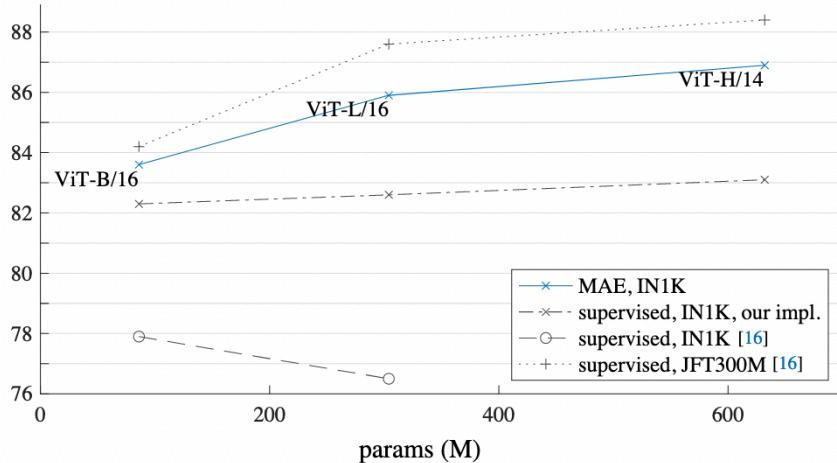


Figure 8. **MAE pre-training vs. supervised pre-training**, evaluated by fine-tuning in ImageNet-1K (224 size). We compare with the original ViT results [16] trained in IN1K or JFT300M.

MAE vs. Contrastive Learning (MoCo v3)

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

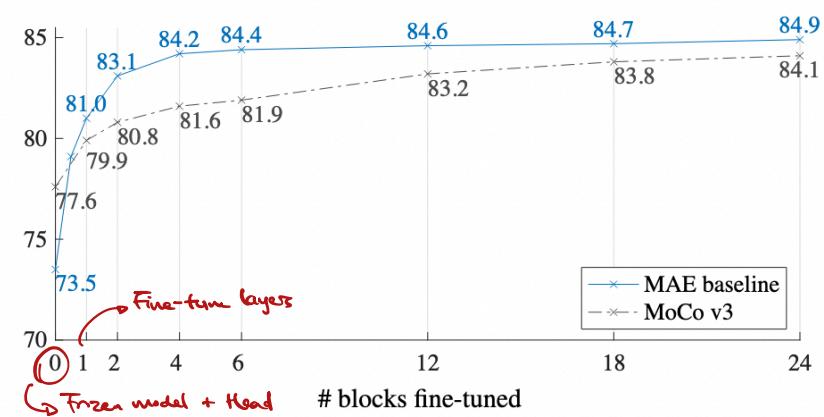


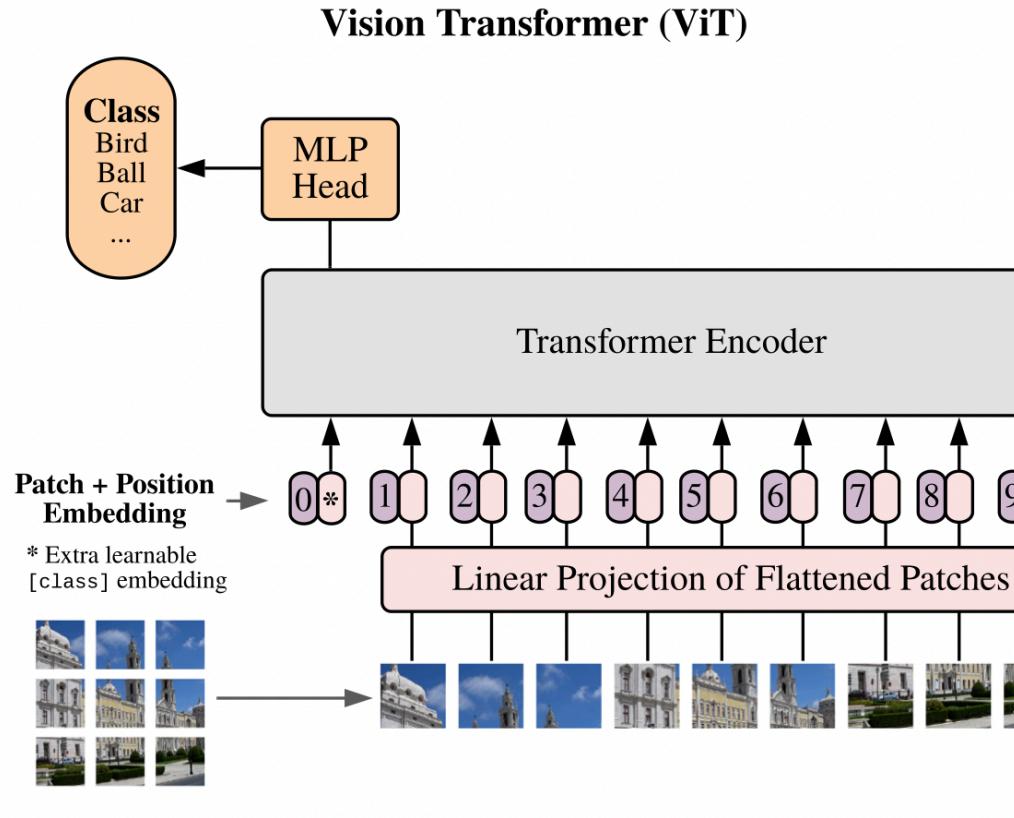
Figure 9. **Partial fine-tuning** results of ViT-L w.r.t. the number of fine-tuned Transformer blocks under the default settings from Table 1. Tuning 0 blocks is linear probing; 24 is full fine-tuning. Our MAE representations are less linearly separable, but are consistently better than MoCo v3 if one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

# A (very quick) overview of Transformers

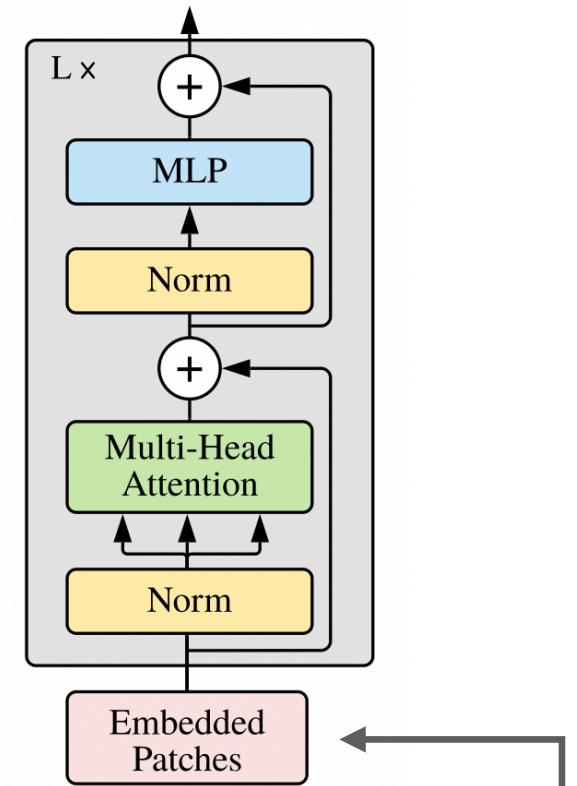
↳ Very general-purpose, modality-independent.



# A (very quick) overview of Transformers

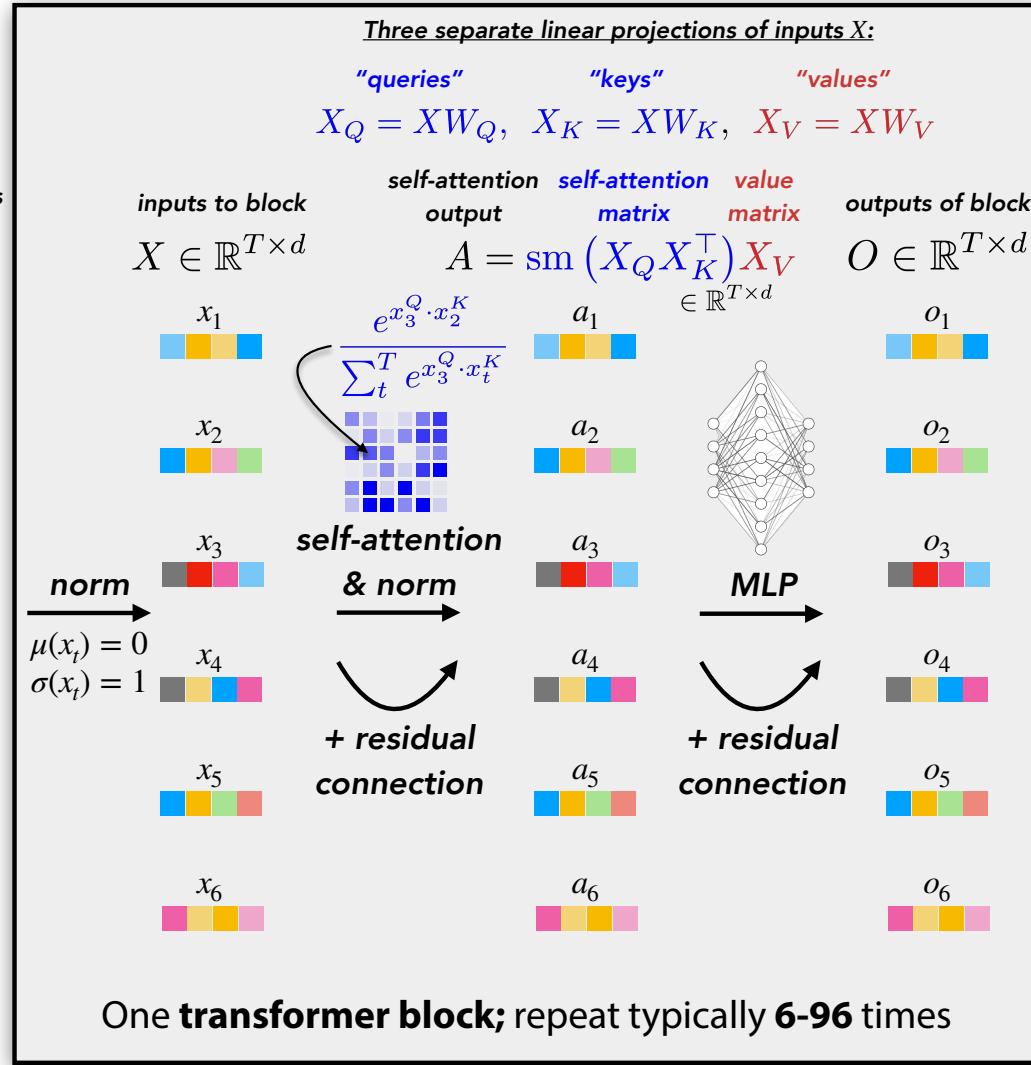
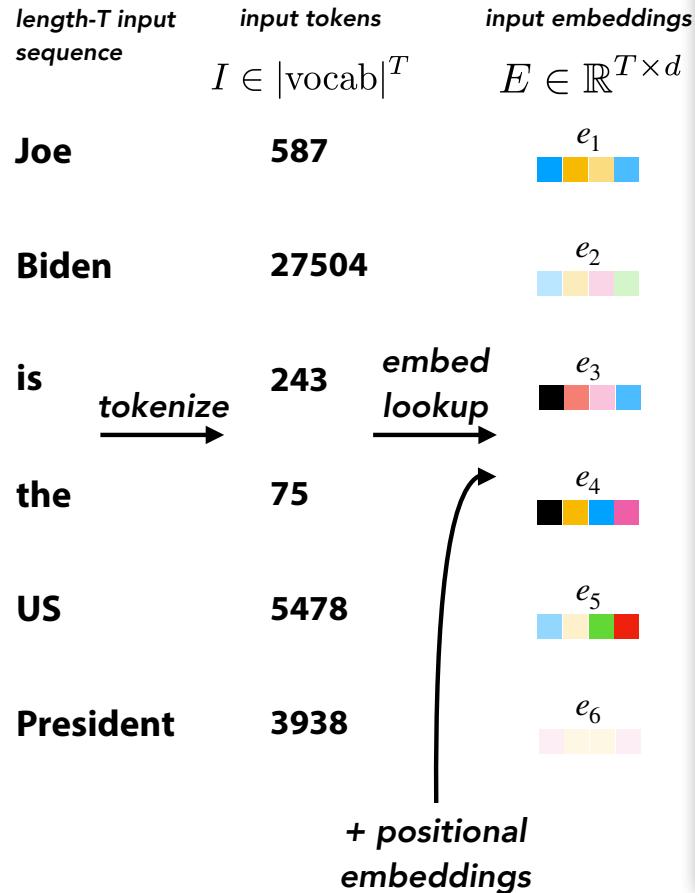


## Transformer Encoder



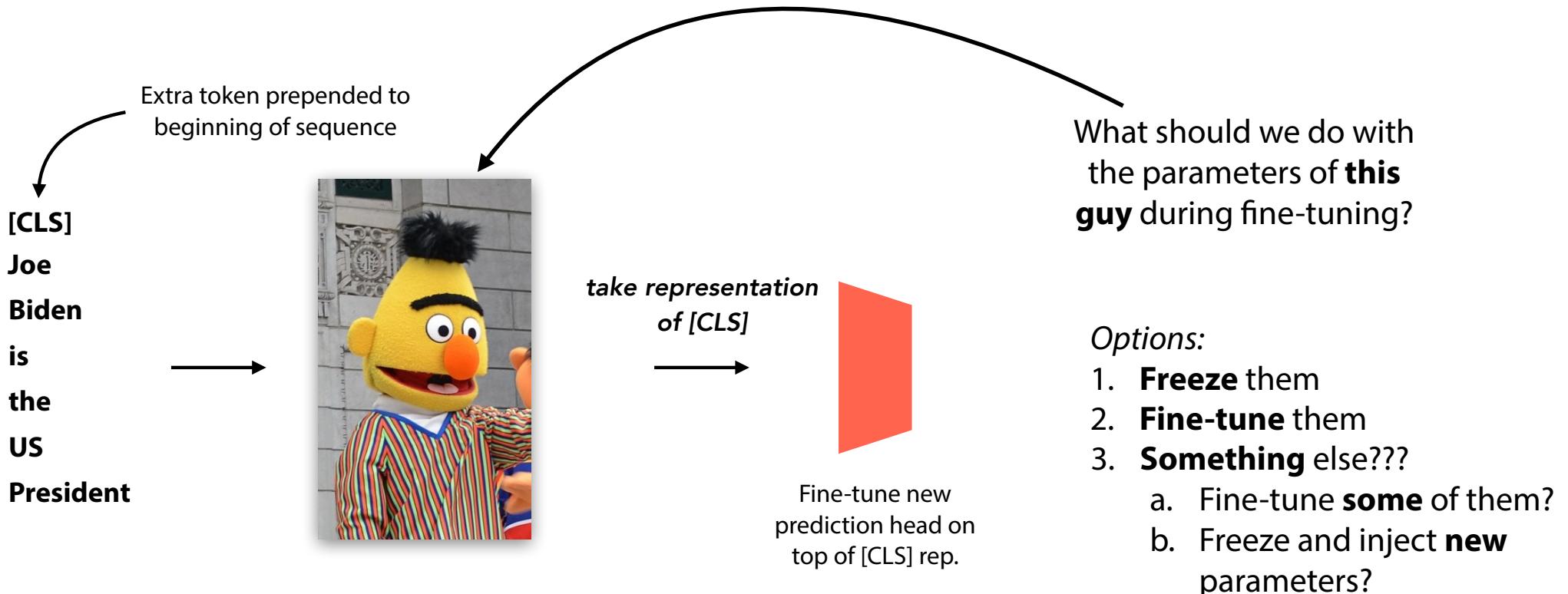
The **~only difference** between Transformers for vision/language/RL/molecules/etc. is what we do for this initial **embedding step**

# Transformers in a bit more detail



$p_\theta^t(\cdot) \in \mathbb{R}^{|\text{vocab}|}$   
 $p_\theta^1(\cdot)$   
 $p_\theta^2(\cdot)$   
**project to vocab. size dimensions**  
 $p_\theta^3(\cdot)$   
 $p_\theta^4(\cdot)$   
 $p_\theta^5(\cdot)$   
 $p_\theta^6(\cdot)$

# 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

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)

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

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

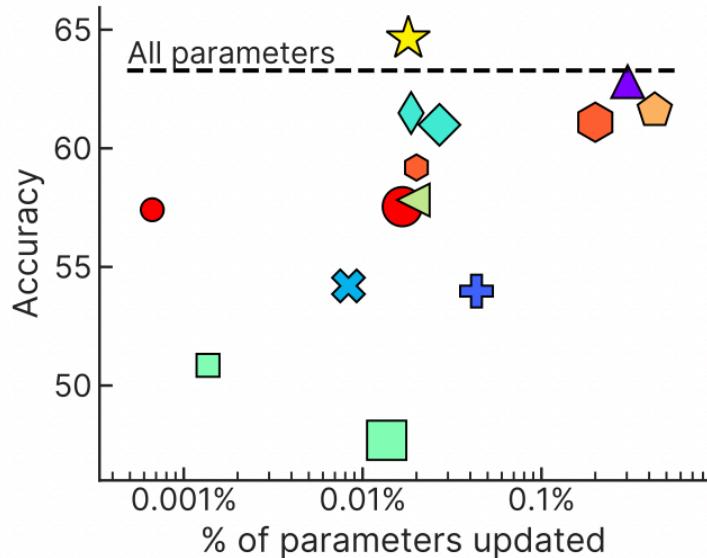
"A little bit" means **only add a few memories** → only make a **low-rank** change to  $W$

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

pre-trained weights (frozen)

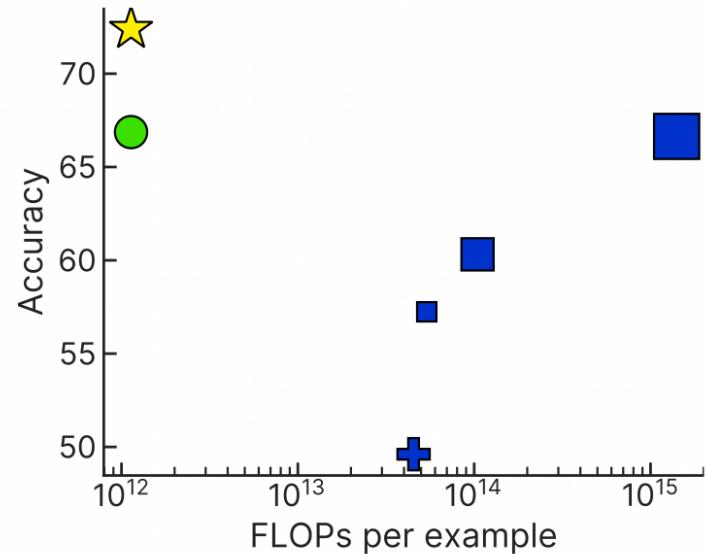
new low-rank residual (fine-tuned)  
 $AB^\top$  should be zero-initialized (how?)

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



Legend:

- (IA)<sup>3</sup>
- LoRA
- BitFit
- Layer Norm
- Compacter
- Compacter++
- Prompt Tuning
- Prefix Tuning
- Adapter
- FISH Mask
- Intrinsic SAID



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!)

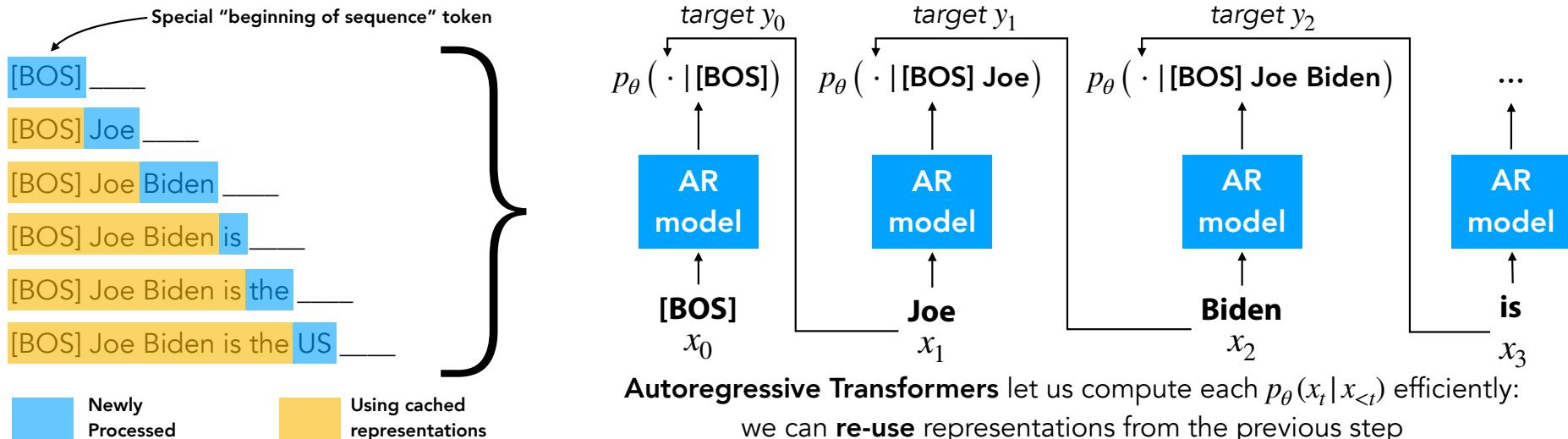
What are some **downsides** of masked autoencoders?

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

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](#).

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

Generative Pretraining from Pixels

Mark Chen<sup>1</sup> Alec Radford<sup>1</sup> Rewon Child<sup>1</sup> Jeff Wu<sup>1</sup> Heewoo Jun<sup>1</sup> Prafulla Dhariwal<sup>1</sup> David Luan<sup>1</sup>

Decision Transformer: Reinforcement  
Learning via Sequence Modeling

WebGPT: Browser-assisted question-answering with  
human feedback



28-04-2022

## Flamingo: a Visual Language Model for Few-Shot Learning

Jean-Baptiste Alayrac<sup>\*,‡</sup>, Jeff Donahue<sup>\*</sup>, Pauline Luc<sup>\*</sup>, Antoine Miech<sup>\*</sup>, Iain Barr<sup>†</sup>, Yana Hasson<sup>†</sup>, Karel Lenc<sup>†</sup>, Arthur Mensch<sup>†</sup>, Katie Millican<sup>†</sup>, Malcolm Reynolds<sup>†</sup>, Roman Ring<sup>†</sup>, Eliza Rutherford<sup>†</sup>, 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<sup>\*,‡</sup>

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

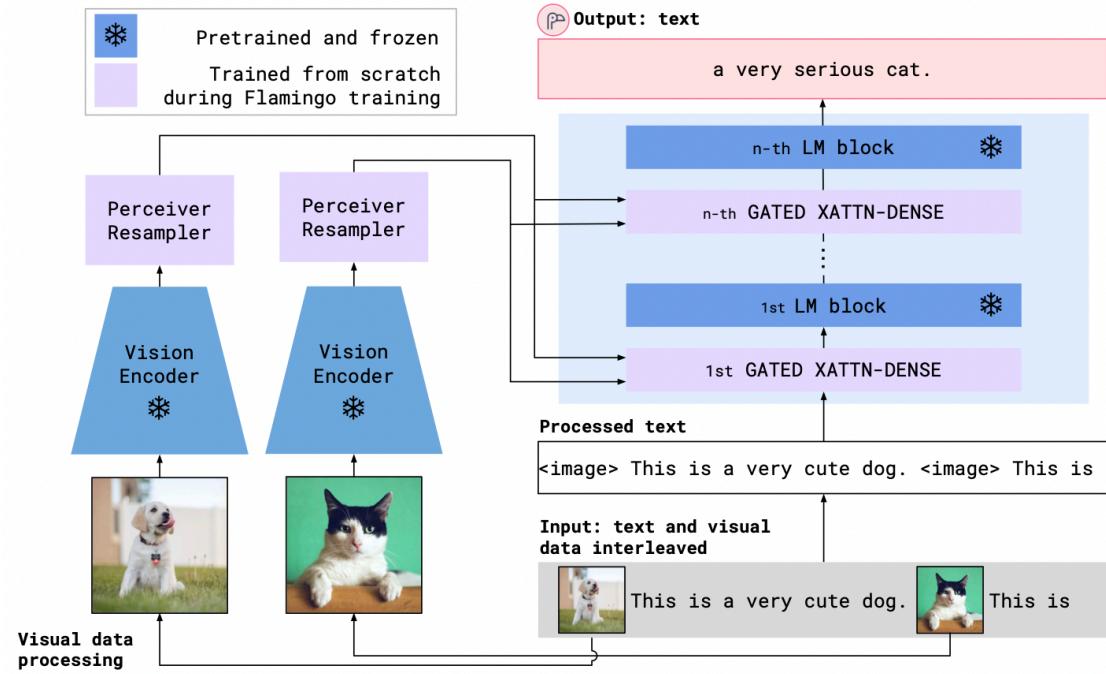
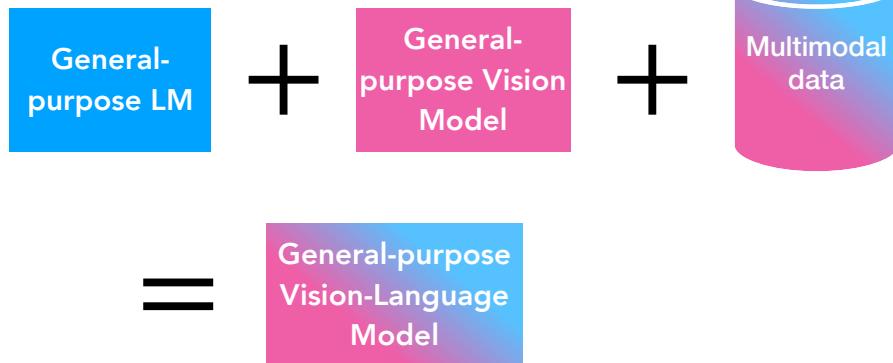
# Case study: Flamingo

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



## Flamingo

Fine-tuning to **combine models**:



## Flamingo: a Visual Language Model for Few-Shot Learning

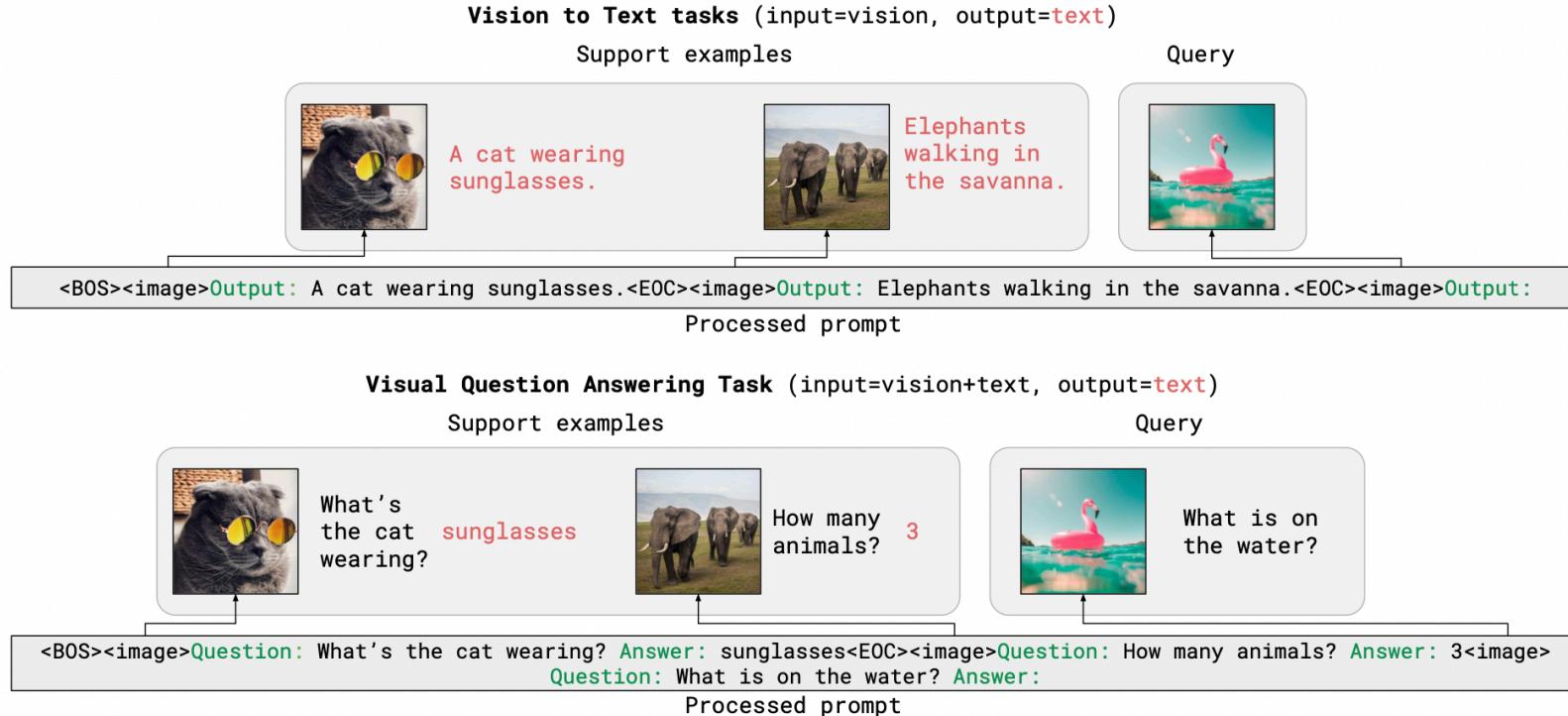
Jean-Baptiste Alayrac<sup>\*‡</sup>, Jeff Donahue<sup>\*</sup>, Pauline Luc<sup>\*</sup>, Antoine Miech<sup>\*</sup>, Iain Barr<sup>†</sup>, Yana Hasson<sup>†</sup>, Karel Lenc<sup>†</sup>, Arthur Mensch<sup>†</sup>, Katie Millican<sup>†</sup>, Malcolm Reynolds<sup>†</sup>, Roman Ring<sup>†</sup>, Eliza Rutherford<sup>†</sup>, 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<sup>\*‡</sup>

<sup>\*</sup>Equal contributions, ordered alphabetically. <sup>†</sup>Equal contributions, ordered alphabetically. <sup>‡</sup>Equal senior contributions

# Case study: Flamingo

## Flamingo: a Visual Language Model for Few-Shot Learning

Jean-Baptiste Alayrac<sup>\*‡</sup>, Jeff Donahue<sup>\*</sup>, Pauline Luc<sup>\*</sup>, Antoine Miech<sup>\*</sup>, Iain Barr<sup>†</sup>, Yana Hasson<sup>‡</sup>, Karel Lenc<sup>‡</sup>, Arthur Mensch<sup>‡</sup>, Katie Millican<sup>†</sup>, Malcolm Reynolds<sup>†</sup>, Roman Ring<sup>†</sup>, Eliza Rutherford<sup>†</sup>, 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<sup>\*‡</sup>  
<sup>\*</sup>Equal contributions, ordered alphabetically. <sup>†</sup>Equal contributions, ordered alphabetically. <sup>‡</sup>Equal senior contributions



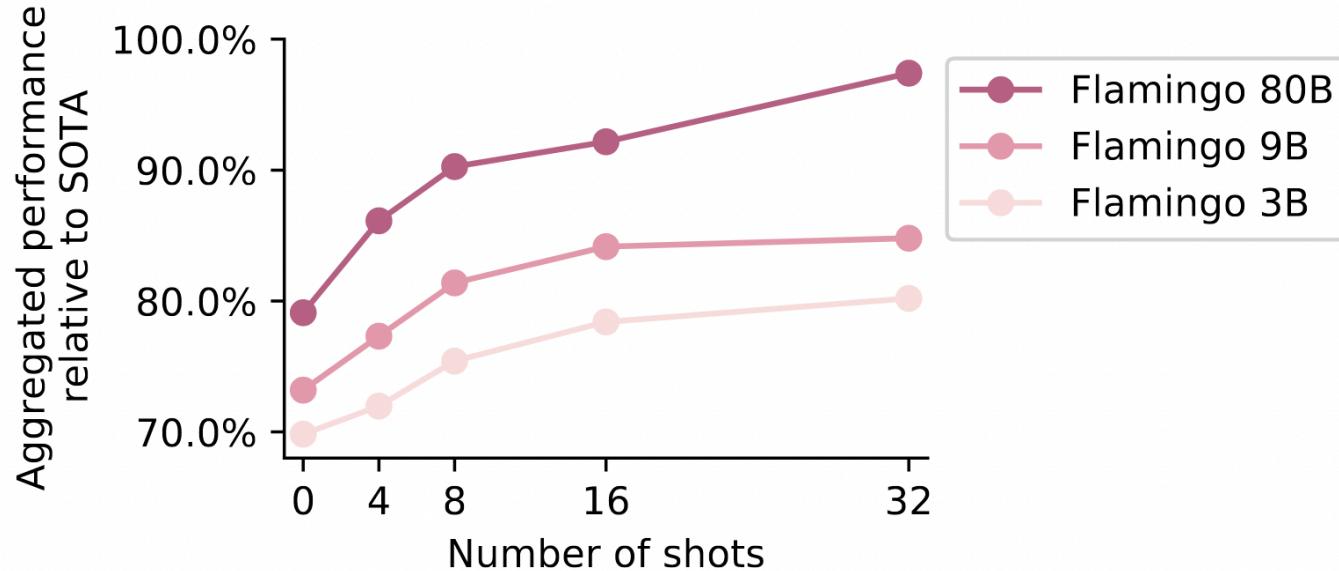
In-context few-shot learning on sequences that freely mix **text** and **images!** Enables few-shot captioning, visual question-answering, etc.

# Case study: Flamingo

## Flamingo: a Visual Language Model for Few-Shot Learning

Jean-Baptiste Alayrac<sup>\*‡</sup>, Jeff Donahue<sup>\*</sup>, Pauline Luc<sup>\*</sup>, Antoine Miech<sup>\*</sup>, Iain Barr<sup>†</sup>, Yana Hasson<sup>‡</sup>, Karel Lenc<sup>‡</sup>, Arthur Mensch<sup>\*</sup>, Katie Millican<sup>†</sup>, Malcolm Reynolds<sup>†</sup>, Roman Ring<sup>†</sup>, Eliza Rutherford<sup>†</sup>, 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<sup>\*‡</sup>

<sup>\*</sup>Equal contributions, ordered alphabetically. <sup>†</sup>Equal contributions, ordered alphabetically. <sup>‡</sup>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  $f_\theta$ :

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**

$y :$

AR model:

$\tilde{x} :$   
**Joe**

**<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

# Reminders

Project proposal due TODAY!

Homework 2 due Monday, October 24

**Kyle's** office hours are hybrid going forward (see Ed for details)

Azure invites have been re-sent - **you have one week to accept!**

**You will need Azure for HW3, so do this today!**