

Vision Transformer (ViT), Pretrained Transformers

CS60010

Encoding Images? Vision Transformer



An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale: ICLR 2021





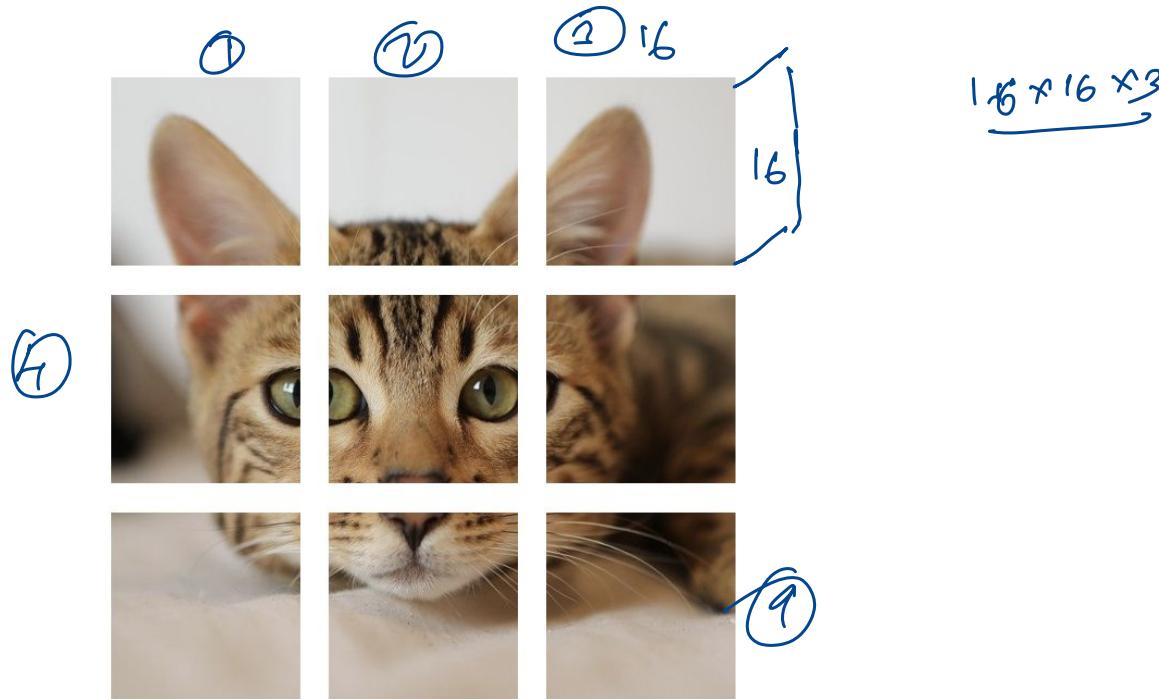
Contextual
Representative

Encoders



Embeddy

Encoding Images? Vision Transformer



An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale: ICLR 2021

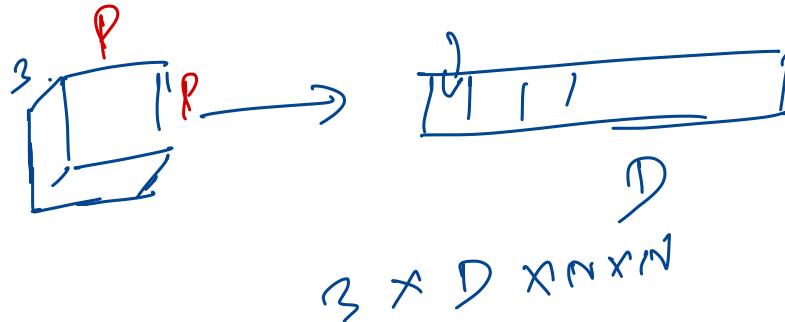
Encoding Images? Vision Transformer

N input patches, each
of shape 3x16x16



An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale: ICLR 2021

Encoding Images? Vision Transformer

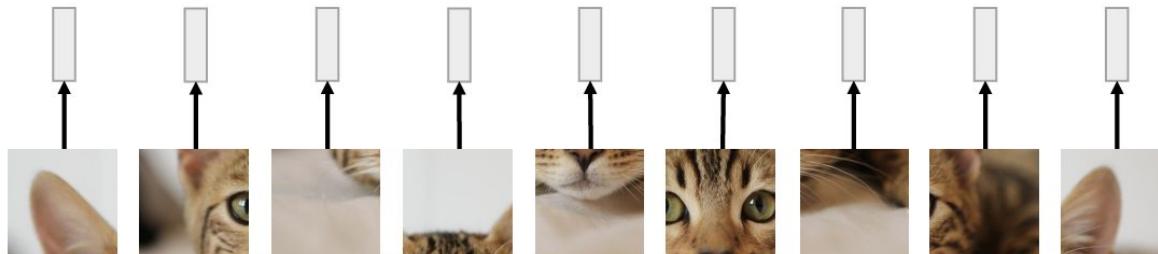


Conv.

$P \times P$ kernel
 P stride

What is this layer?

Linear projection to D-dimensional vector
N input patches, each of shape $3 \times 16 \times 16$

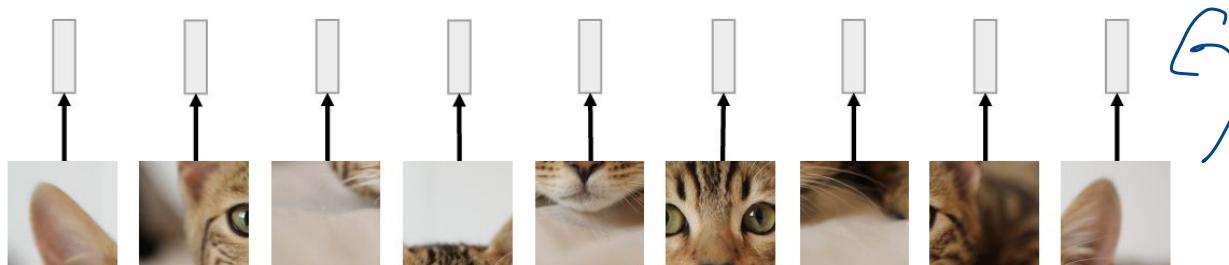


An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale: ICLR 2021

Encoding Images? Vision Transformer

With patch size p , this layer is Conv2D
(filter size: $p \times p$, Input channels: 3,
Output channels: D , stride= p)

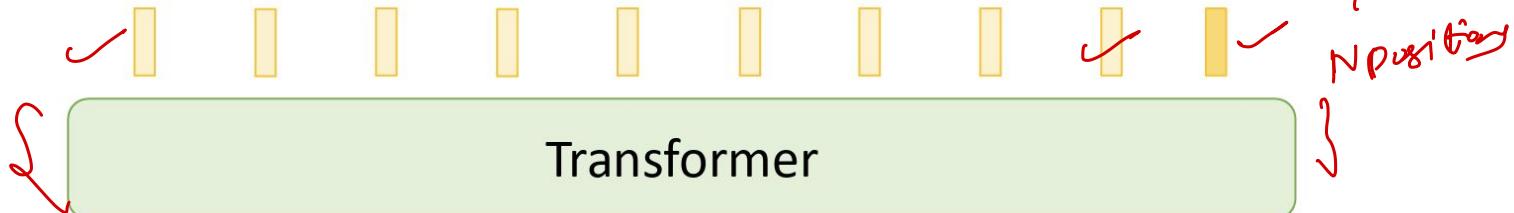
Linear projection to
 D -dimensional vector
N input patches, each
of shape $3 \times 16 \times 16$



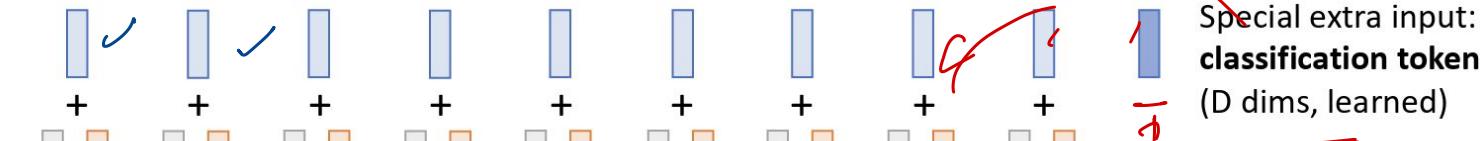
Encoding Images? Vision Transformer

$$\mathbf{z}_0 = [\mathbf{x}_{\text{class}}, \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \dots; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{\text{pos}} \quad \mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{\text{pos}} \in \mathbb{R}^{(N+1) \times D}$$

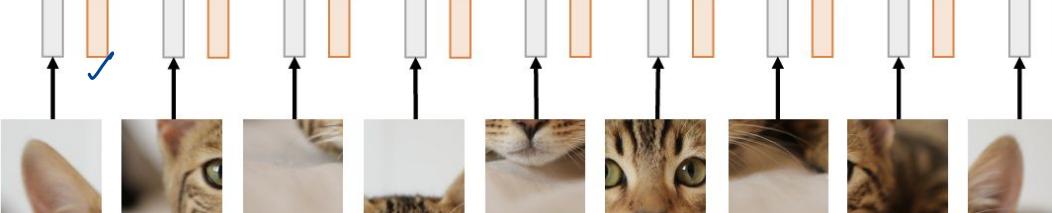
Output vectors



Add positional embedding: learned D-dim vector per position



Linear projection to D-dimensional vector



Encoding Images? Vision Transformer

224×224

What is the sequence length?

$$14 \times 14 = 196$$

What is the size of the attention matrix?

14^4

Output vectors

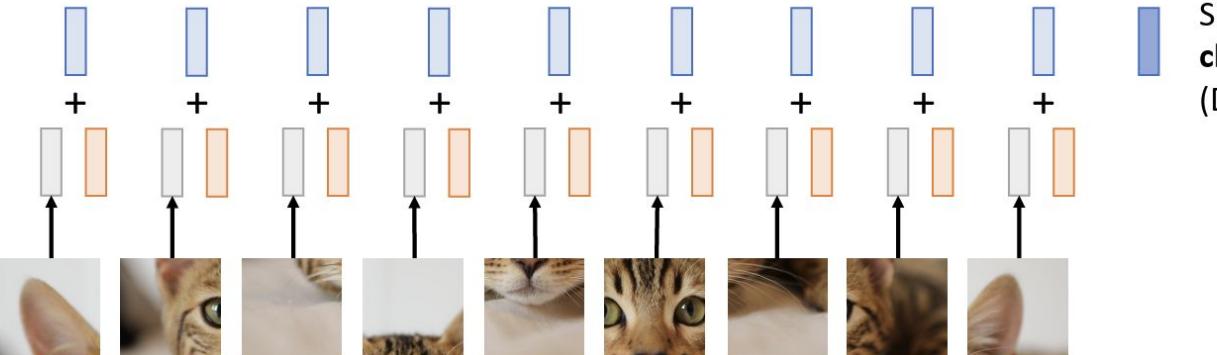


Exact same as
NLP Transformer!

Transformer

Linear projection
to C-dim vector
of predicted
class scores

Add positional
embedding: learned D-
dim vector per position



Linear projection to
D-dimensional vector

N input patches, each
of shape $3 \times 16 \times 16$



Special extra input:
classification token
(D dims, learned)

Encoding Images? Vision Transformer

In practice: take 224x224 input image,
divide into 14x14 grid of 16x16 pixel
patches (or 16x16 grid of 14x14 patches)

Output vectors

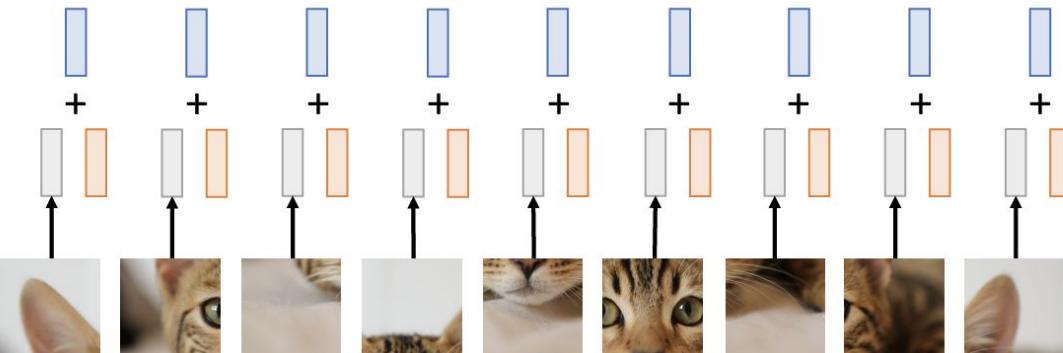


Each attention matrix has $14^4 = 38,416$
entries, takes 150 KB → *1 byte/entry*
(or 65,536 entries, takes 256 KB)

Linear projection to C-dim vector
of predicted class scores

Exact same as
NLP Transformer!

Add positional
embedding: learned D-
dim vector per position



Linear projection to
D-dimensional vector

N input patches, each
of shape 3x16x16

Special extra input:
classification token
(D dims, learned)

Vision Transformer Model Variants

Model	Layers	Hidden size D	MLP size	Heads	Params
ViT-Base	12	768	3072	12	86M
ViT-Large	24	1024	4096	16	307M
ViT-Huge	32	1280	5120	16	632M

Can you make sense of the number of parameters?

↑
model dim
↓
 n

ViT Base

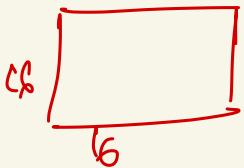
$12 \cdot d^2 \times L$

$$= [12 \times (768)^2 \times 12]$$

$\approx 84-85 M$

I^{P} \rightarrow Embedd

224×224
image



196 set lengths

P.R.

196×768

Conv.

$$3 \times 768 \times 16 \times 16 = 768 \times (196 + 768)$$

$\approx 0.75 \text{ M}$

Op to class

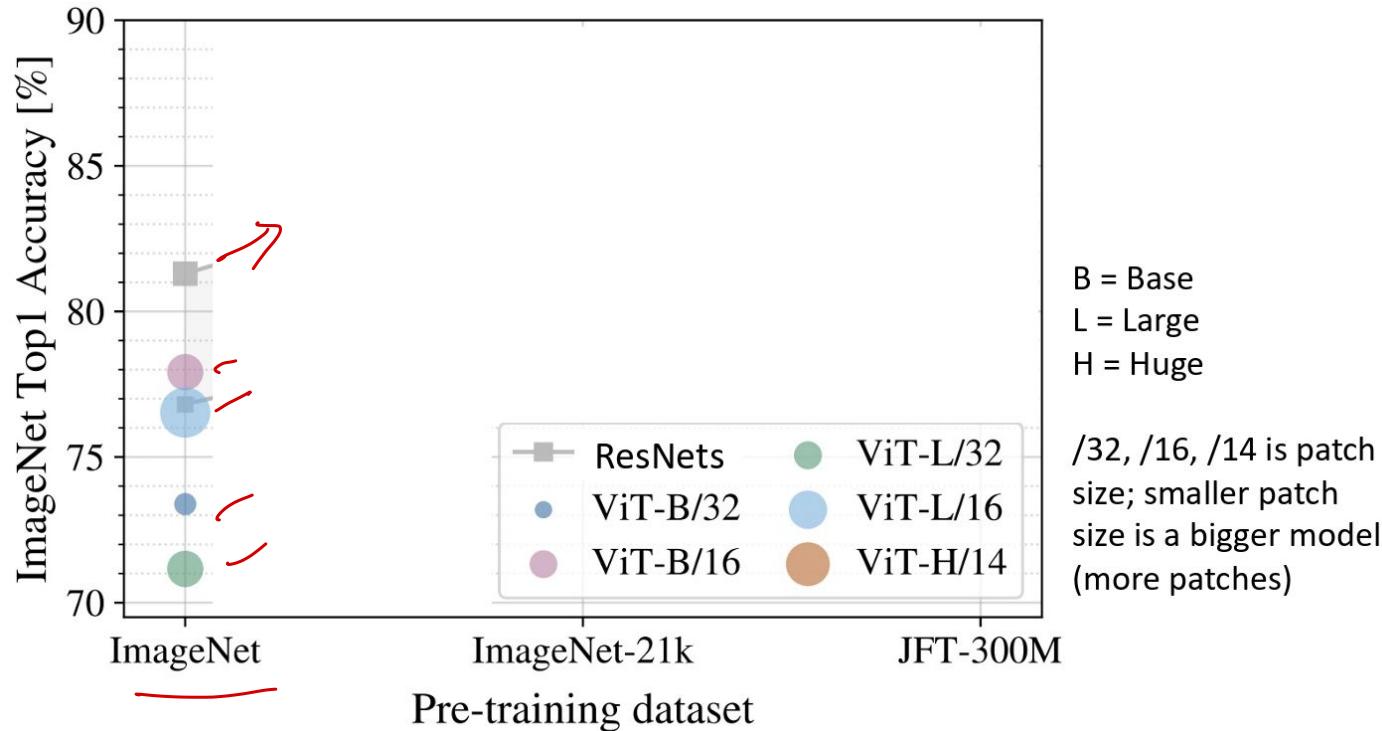
768×1000

$\approx 0.75 \text{ M}$

Performance Comparison: ViTs vs ResNets

Recall: ImageNet dataset has 1k categories, 1.2M images

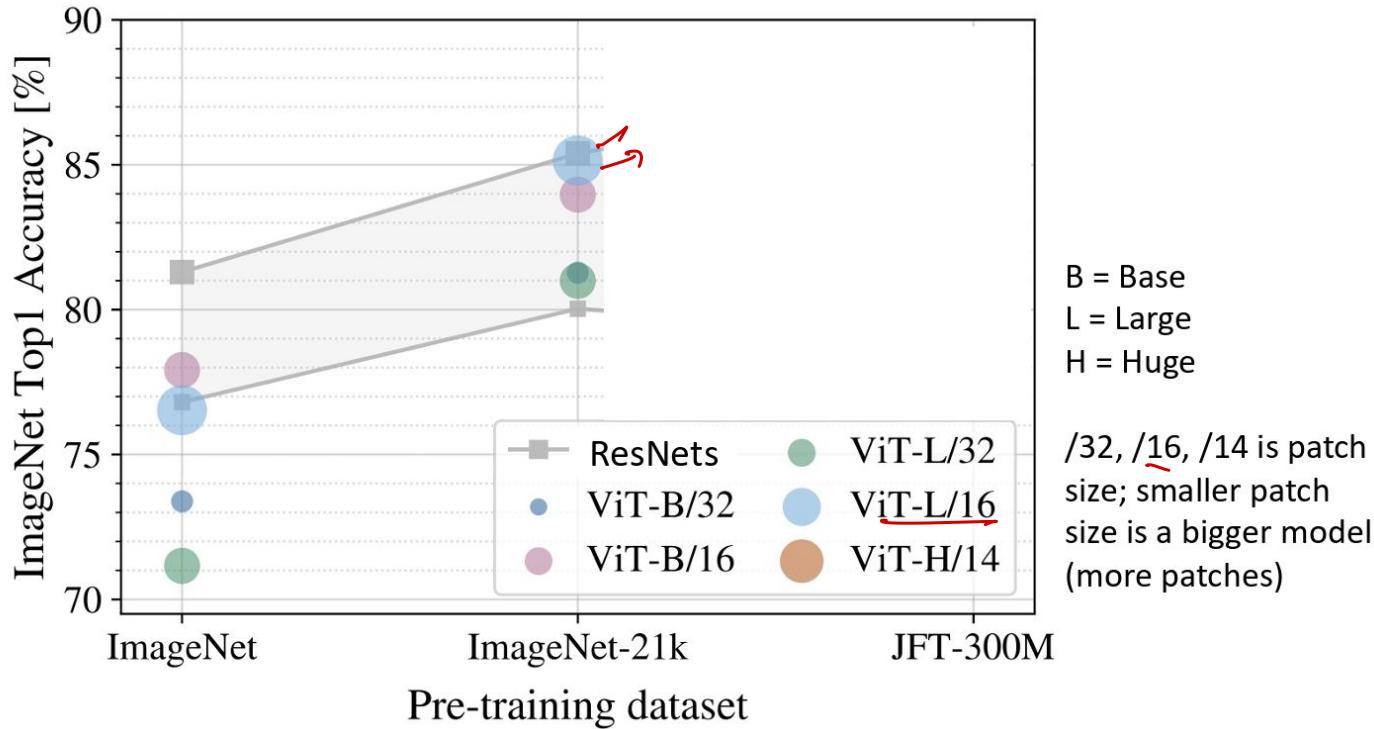
When trained on ImageNet, ViT models perform worse than ResNets



Performance Comparison: ViTs vs ResNets

ImageNet-21k has
14M images with 21k
categories

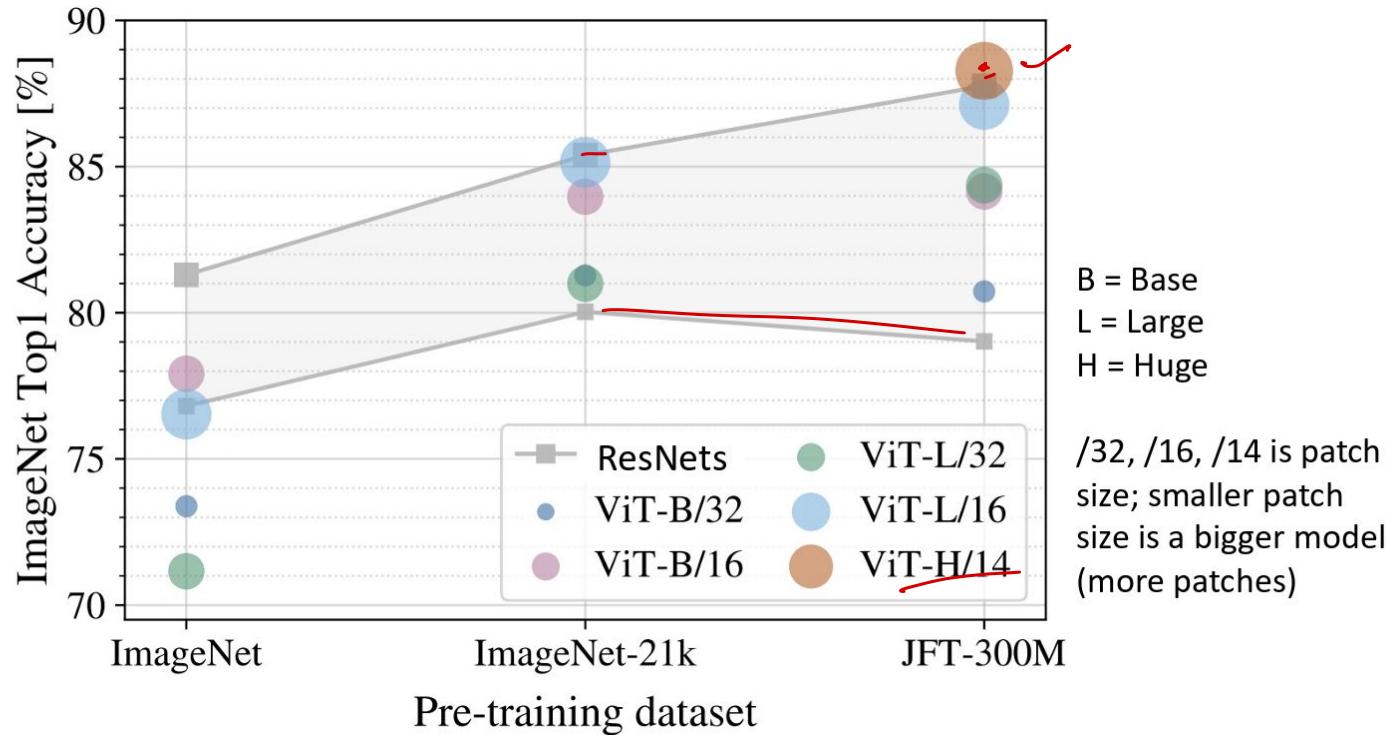
If you pretrain on
ImageNet-21k and
fine-tune on
ImageNet, ViT does
better: big ViTs match
big ResNets



Performance Comparison: ViTs vs ResNets

JFT-300M is an internal Google dataset with 300M labeled images

If you pretrain on JFT and finetune on ImageNet, large ViTs outperform large ResNets



What is pretraining? Let's go back to text

The idea of transfer learning

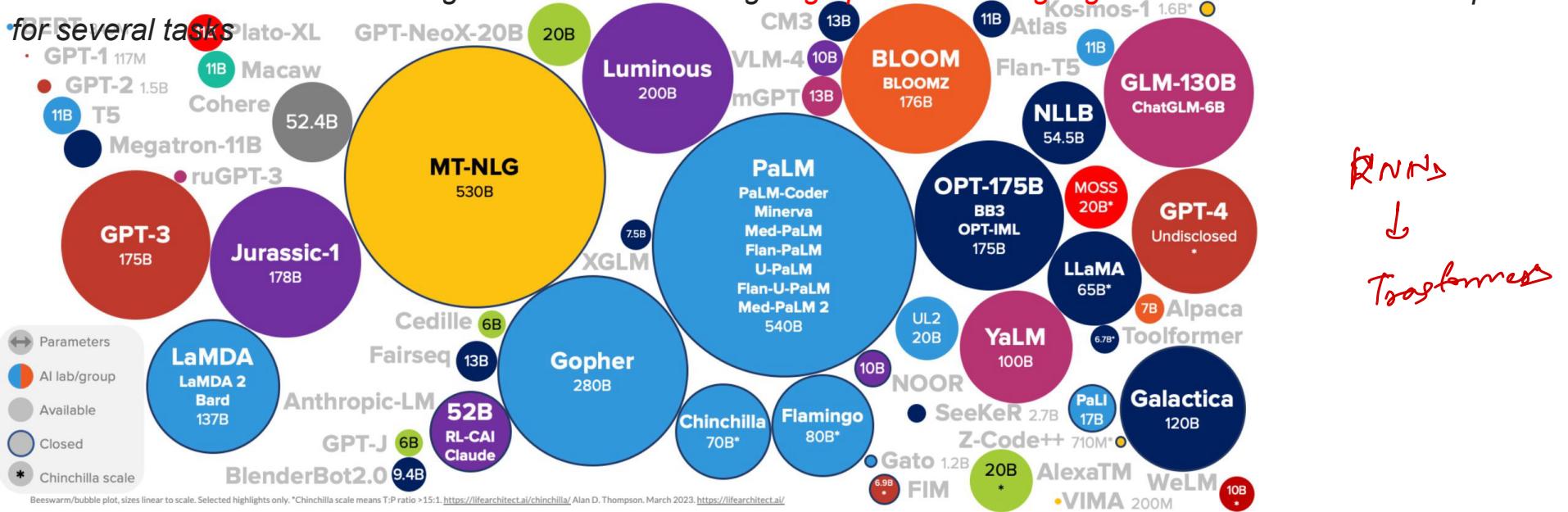
The idea of transfer learning

However, what has made Transformers popular is that they can be combined with the idea of *transfer learning*.

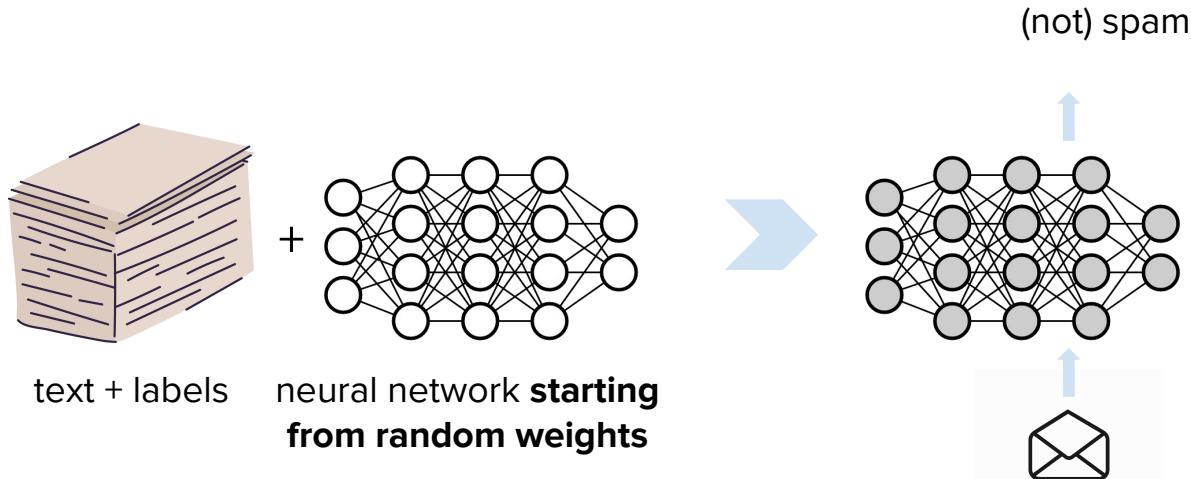
Transformers have become the go-to model for building **large pretrained language models** which can be adapted

• for several tasks

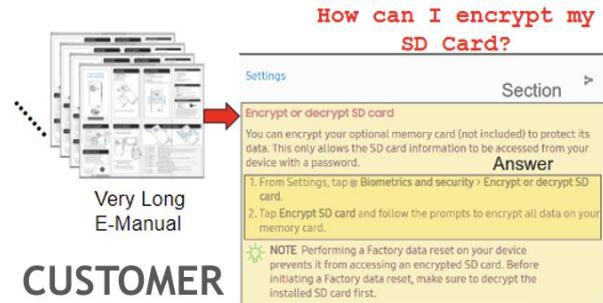
- GPT-1 117M
 - GPT-2 1.5B



Standard Supervised Deep Learning



Abundance in domain-specific NLP Applications



CUSTOMER

SUPPORT

DOMAIN

FINANCIAL

DOMAIN

Cash and Cash Equivalents
Debt Instrument Convertible Conversion Price

As of December 31, 2020, we had cash equivalents of \$24.8 million and a closing stock price of \$18.20 per share.

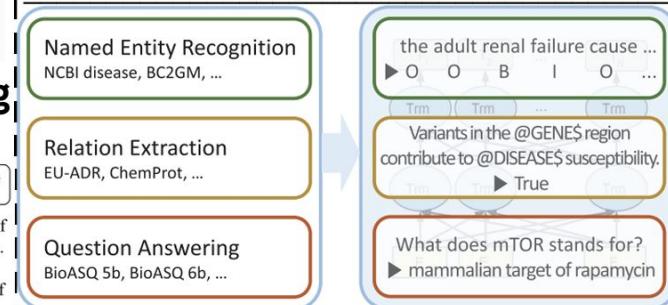
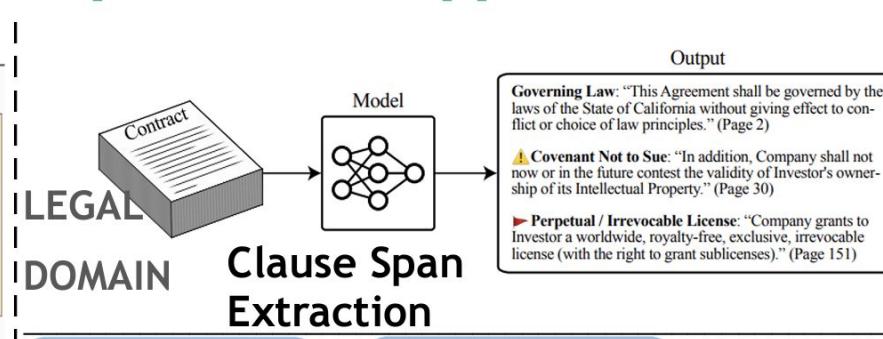
We also acquired a business loan from the U.S. Bancorp of \$60.5 million.

Limit of Credit Facility Maximum Borrowing Capacity
Impairment Loss

Finally, our firm reports no impairment loss for this year.

NER

Very difficult to get sufficient labeled data for each task

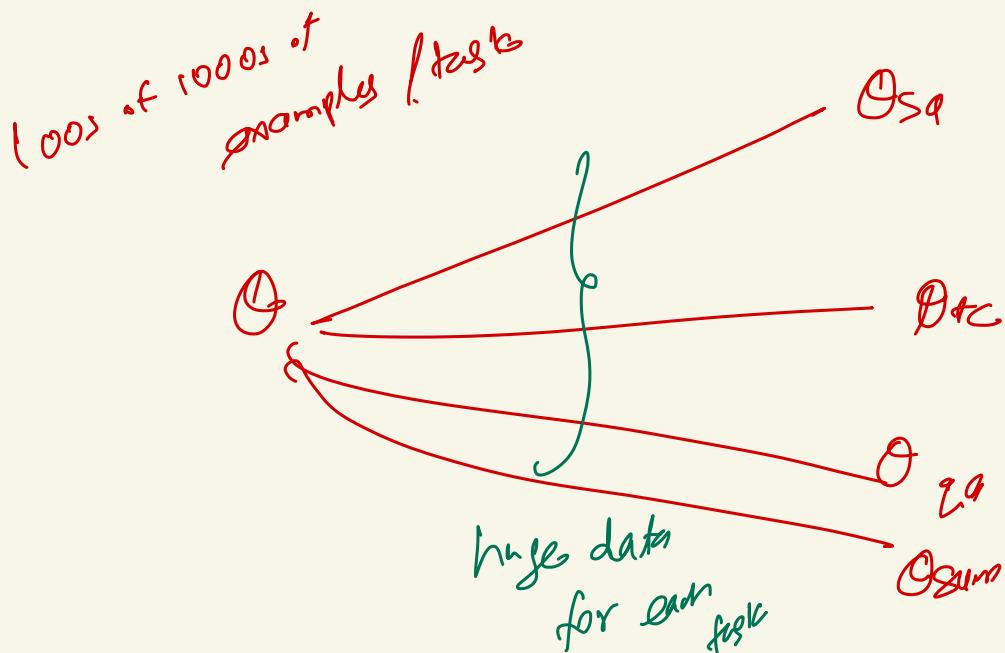


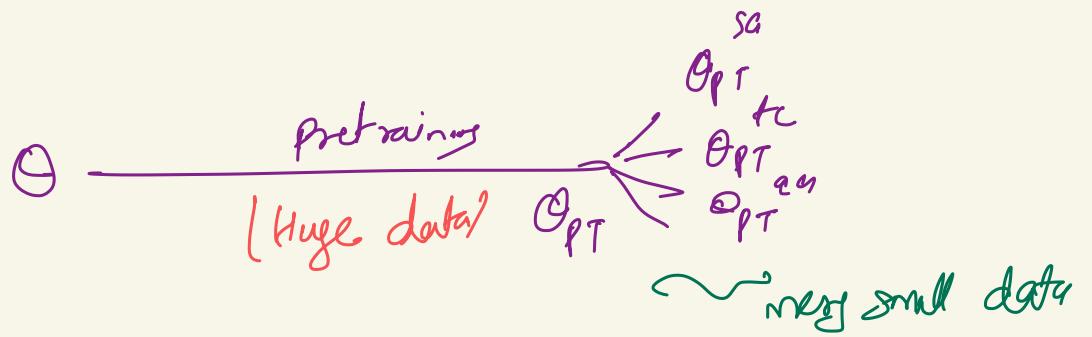
BIOMEDICAL
DOMAIN

topic class → Transformer
sent. analys → Transformer
→ Osa

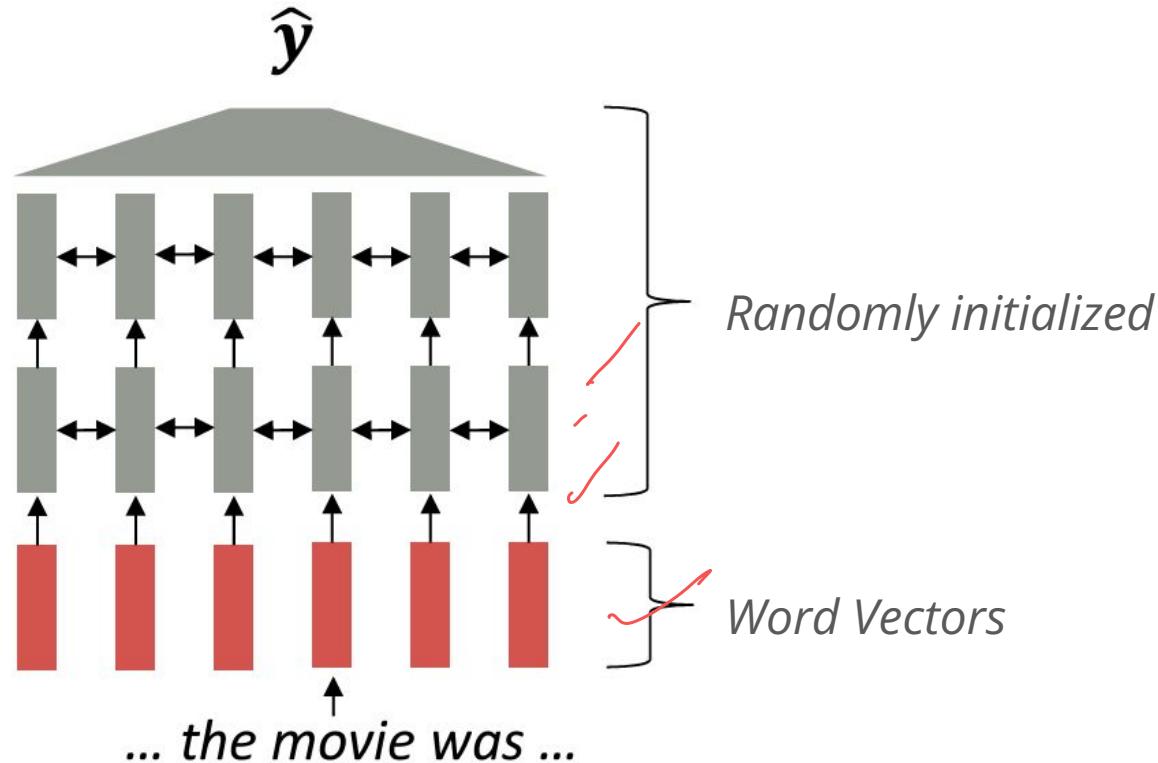
Languages

θ_{sa} , θ_{tc} , θ_{ca} , θ_{sum}





Why is this an issue



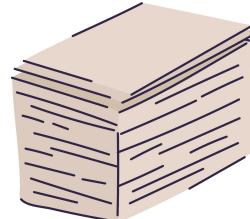
2

The training data we have for our downstream task (like question answering) must be sufficient to teach all contextual aspects of language.

Pretrain-then-Finetune (from 2018)

*Stage 1:
Pretrain a model*

neural network
starting from
random weights



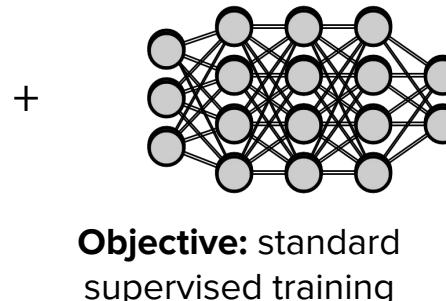
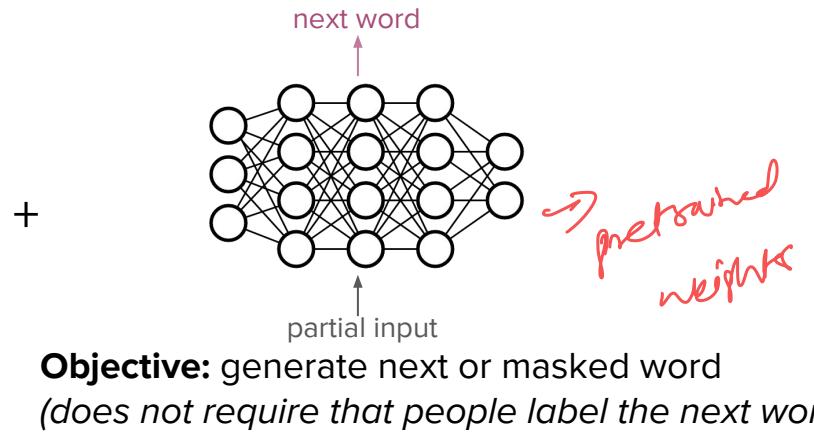
text

*Stage 2:
Finetune the model*

neural network
starting from
pretrained
weights



text + **labels**



Transfer Learning with Language Modeling

Language modelling is the popular choice of pretraining in NLP.

The **key idea** is to **train a neural network language model (on vast text corpora)**. And **then, use it to transfer that knowledge to any target task in NLP** we care about.

This is a **very successful** idea in NLP:

- Supervised systems for various NLP tasks used to require **millions of examples** to learn tasks. NLP systems that use language models as a starting point, need much fewer - **thousands** of examples to do so!

What does a model learn from pretraining?

- There are canines everywhere! One dog in the front room, and two dogs
- It wasn't just big it was enormous
- The author of "A Room of One's Own" is Virginia Woolf
- The doctor told me that he
- The square root of 4 is 2

2

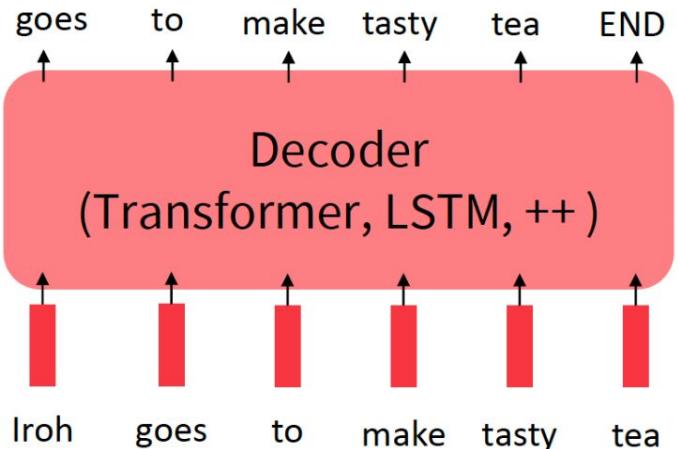
- *Text contains enormous amounts of knowledge*
- *Pretraining on lots of text with all that knowledge is what gives language models their ability to do so much*

The Pretraining / Finetuning paradigm

Pretraining can improve NLP applications by serving as parameter initialization.

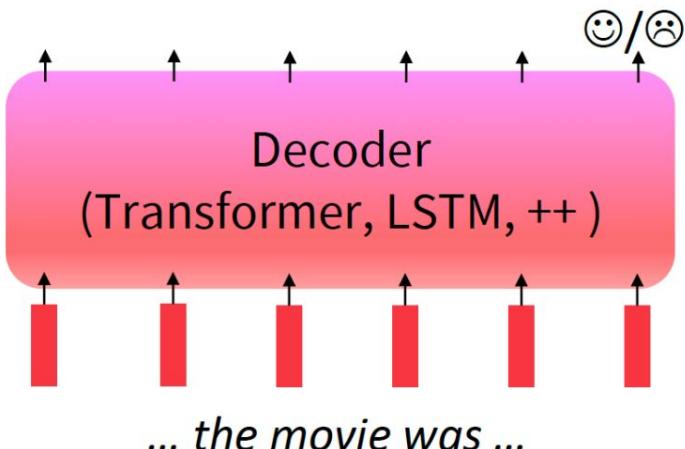
Step 1: Pretrain (on language modeling)

Lots of text; learn general things!

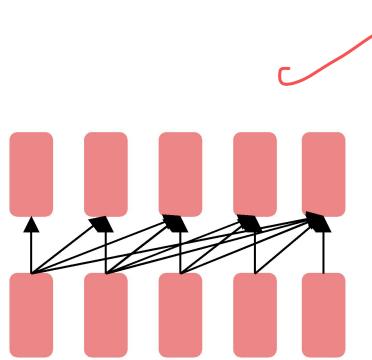


Step 2: Finetune (on your task)

Not many labels; adapt to the task!

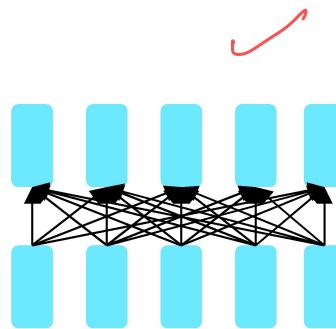


Three architectures for large language models



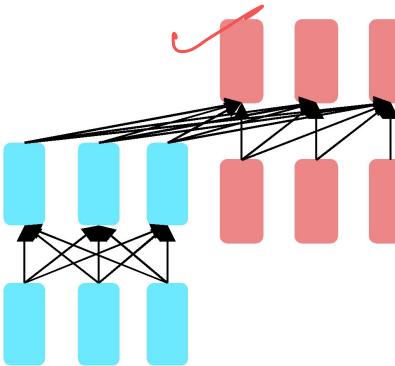
Decoders

GPT, Claude,
Llama, Mixtral



Encoders

BERT family
F



Encoder-decoders

Flan-T5
BART

Pretrain Decoder

Predict next word

Pretrain Encoder

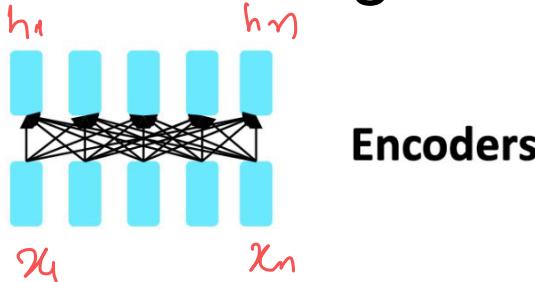


Can't use

next word
prediction

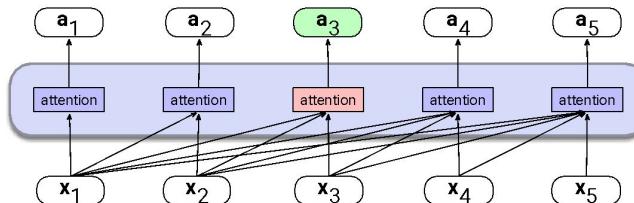
The

Pretraining Encoders

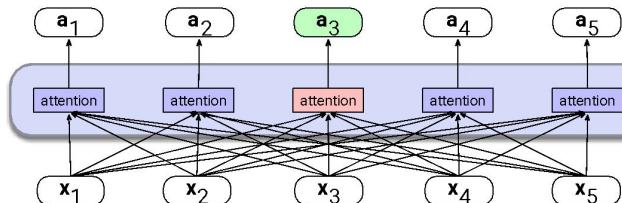


- Gets bidirectional context – can condition on future!
- How do we train them to build strong representations?

- map sequences of input embeddings (x_1, \dots, x_n)
- to sequences of output embeddings of the same length (h_1, \dots, h_n),
- where the output vectors have been contextualized using information from the entire input sequence.



a) A causal self-attention layer



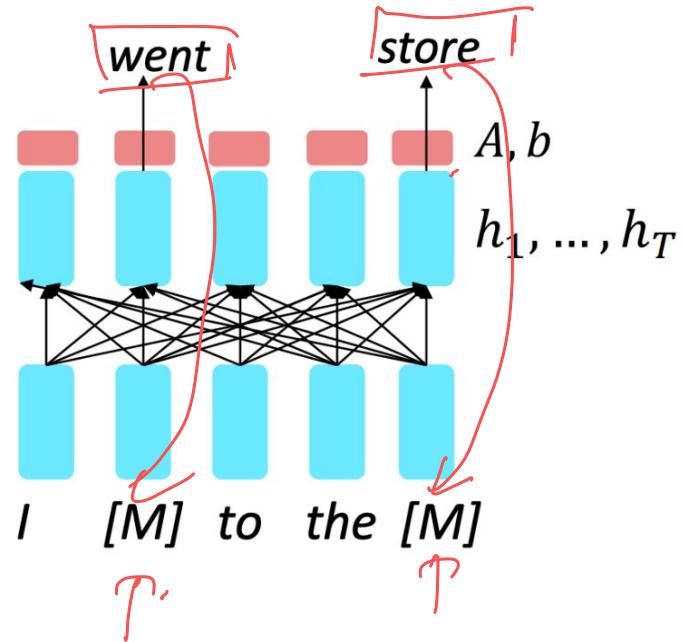
b) A bidirectional self-attention layer

Solution: Use Masks

Idea: replace some fraction of words in the input with a special [MASK] token; predict these words.

$$h_1, \dots, h_T = \text{Encoder}(w_1, \dots, w_T)$$
$$y_i \sim Aw_i + b$$

Only add loss terms from words that are “masked out.” If \tilde{x} is the masked version of x , we’re learning $p_\theta(x|\tilde{x})$. Called **Masked LM**.

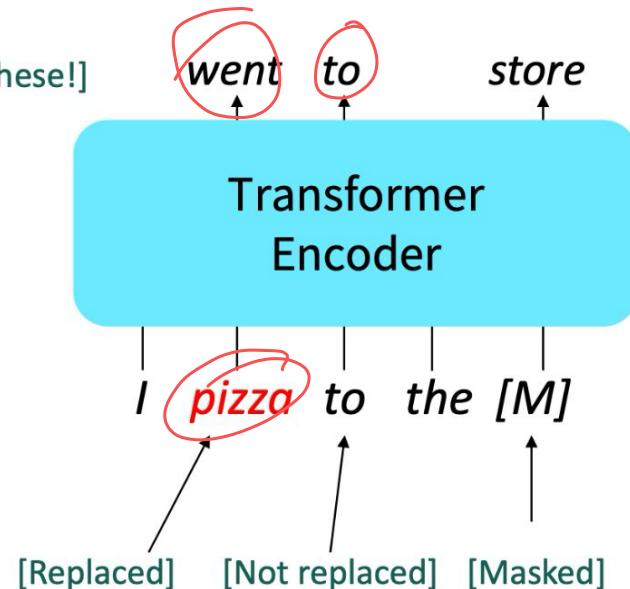


BERT: Bidirectional Encoder Representations from Transformers

Devlin et al., 2018 proposed the “Masked LM” objective and **released the weights of a pretrained Transformer**, a model they labeled BERT.

Some more details about Masked LM for BERT:

- Predict a random 15% of (sub)word tokens.
 - Replace input word with [MASK] 80% of the time (1.2%)
 - Replace input word with a random token 10% of the time (1.5%)
 - Leave input word unchanged 10% of the time (but still predict it!) (1.5%)
- Why? Doesn’t let the model get complacent and not build strong representations of non-masked words.
(No masks are seen at fine-tuning time!)



More Details about BERT

vocab

- Two models were released:
 - BERT-base: 12 layers, 768-dim hidden states, 12 attention heads, 110 million params.
 - BERT-large: 24 layers, 1024-dim hidden states, 16 attention heads, 340 million params.
- Trained on:
 - BooksCorpus (800 million words)
 - English Wikipedia (2,500 million words)
- Pretraining is expensive and impractical on a single GPU.
 - BERT was pretrained with 64 TPU chips for a total of 4 days.
 - (TPUs are special tensor operation acceleration hardware)
- Finetuning is practical and common on a single GPU
 - “Pretrain once, finetune many times.”