

# **Modeling Sequences/Sets: Transformers**

Prof. Adriana Kovashka  
University of Pittsburgh  
(with some slides from me)

# Plan for this lecture

- Background
  - Context prediction, unsupervised learning
- Transformer models
  - Self-attention
  - Adapting self-attention for sequential data
  - The transformer architecture, encoder/decoder
  - Pre-training, BERT
- Transformers beyond language

# Additional resources

- Learning about transformers on your own?
  - Key recommended resource:
    - <http://nlp.seas.harvard.edu/2018/04/03/attention.html>
    - The Annotated Transformer by Sasha Rush
    - Jupyter Notebook using PyTorch that explains everything!
  - The Illustrated Transformer
    - <http://jalammar.github.io/illustrated-transformer/>
  - Attention visualizer
    - <https://github.com/jessevig/bertviz>

# How do we represent the meaning of a word?

Definition: **meaning** (Webster dictionary)

- the idea that is represented by a word, phrase, etc.
- the idea that a person wants to express by using words, signs, etc.
- the idea that is expressed in a work of writing, art, etc.

Commonest linguistic way of thinking of meaning:

signifier (symbol)  $\iff$  signified (idea or thing)

= denotational semantics

# How do we have usable meaning in a computer?

Common solution: Use e.g. [WordNet](#), a thesaurus containing lists of synonym sets and hypernyms (“is a” relationships).

e.g. synonym sets containing “good”:

```
from nltk.corpus import wordnet as wn
poses = {'n':'noun', 'v':'verb', 's':'adj (s)', 'a':'adj', 'r':'adv'}
for synset in wn.synsets("good"):
    print("{}: {}".format(poses[synset.pos()],
                          ", ".join([l.name() for l in synset.lemmas()])))
```

noun: good  
noun: good, goodness  
noun: good, goodness  
noun: commodity, trade\_good, good  
adj: good  
adj (sat): full, good adj:  
good  
adj (sat): estimable, good, honorable, respectable adj (sat):  
beneficial, good  
adj (sat): good  
adj (sat): good, just, upright  
...  
adverb: well, good  
adverb: thoroughly, soundly, good

e.g. hypernyms of “panda”:

```
from nltk.corpus import wordnet as wn
panda = wn.synset("panda.n.01") hyper =
lambda s: s.hypernyms()
list(pandaclosure(hyper))
```

[Synset('procyonid.n.01'),  
Synset('carnivore.n.01'),  
Synset('placental.n.01'),  
Synset('mammal.n.01'),  
Synset('vertebrate.n.01'),  
Synset('chordate.n.01'),  
Synset('animal.n.01'),  
Synset('organism.n.01'),  
Synset('living\_thing.n.01'),  
Synset('whole.n.02'),  
Synset('object.n.01'),  
Synset('physical\_entity.n.01'),  
Synset('entity.n.01')]

# Problems with resources like WordNet

- Great as a resource but missing nuance
  - e.g. “proficient” is listed as a synonym for “good”. This is only correct in some contexts.
- Missing new meanings of words
  - e.g., wicked, badass, nifty, wizard, genius, ninja, bombest
  - Impossible to keep up-to-date!
- Subjective
- Requires human labor to create and adapt
- Can’t compute accurate word similarity

# Representing words as discrete symbols

In traditional NLP, we regard words as discrete symbols:  
`hotel`, `conference`, `motel` - a **localist** representation

Means one 1, the rest 0s

Words can be represented by **one-hot** vectors:

`motel` = [0 0 0 0 0 0 0 0 0 1 0 0 0 0]  
`hotel` = [0 0 0 0 0 0 1 0 0 0 0 0 0 0]

Vector dimension = number of words in vocab (e.g. 500,000)

# Problem with words as discrete symbols

Example: in web search, if user searches for “Seattle motel”, we would like to match documents containing “Seattle hotel”.

But:

$$\begin{aligned}\text{motel} &= [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0] \\ \text{hotel} &= [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]\end{aligned}$$

These two vectors are orthogonal.

There is no natural notion of similarity for one-hot vectors!

**Solution:**

- Could try to rely on WordNet’s list of synonyms to get similarity?
  - But it is well-known to fail badly: incompleteness, etc.
- Instead: learn to encode similarity in the vectors themselves

# Representing words by their context



- Distributional semantics: A word's meaning is given by the words that frequently appear close-by
  - “*You shall know a word by the company it keeps*” (J. R. Firth 1957)
  - One of the most successful ideas of modern statistical NLP!
- When a word  $w$  appears in a text, its **context** is the set of words that appear nearby (within a fixed-size window).
- Use the many contexts of  $w$  to build up a representation of  $w$

...government debt problems turning into **banking** crises as happened in 2009...

...saying that Europe needs unified **banking** regulation to replace the hodgepodge...

...India has just given its **banking** system a shot in the arm...



These **context words** will represent *banking*

# What can we learn from reconstructing the input?

Stanford University is located in \_\_\_\_\_, California.

# What can we learn from reconstructing the input?

I put \_\_\_\_ fork down on the table.

# What can we learn from reconstructing the input?

The woman walked across the street,  
checking for traffic over \_\_\_\_ shoulder.

# What can we learn from reconstructing the input?

I went to the ocean to see the fish, turtles, seals, and \_\_\_\_\_.

# What can we learn from reconstructing the input?

Overall, the value I got from the two hours watching  
it was the sum total of the popcorn and the drink.

The movie was \_\_\_\_.

# What can we learn from reconstructing the input?

Iroh went into the kitchen to make some tea.

Standing next to Iroh, Zuko pondered his destiny.

Zuko left the\_\_\_\_\_.

# What can we learn from reconstructing the input?

I was thinking about the sequence that goes

1, 1, 2, 3, 5, 8, 13, 21, \_\_\_\_\_

# Word vectors

We will build a dense vector for each word, chosen so that it is similar to vectors of words that appear in similar contexts

$$\text{banking} = \begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{pmatrix}$$

Note: word vectors are sometimes called word embeddings or word representations. They are a distributed representation.

# Word meaning as a neural word vector - visualization

$$\text{expect} = \begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \\ 0.487 \end{pmatrix}$$



# Word2Vec Overview

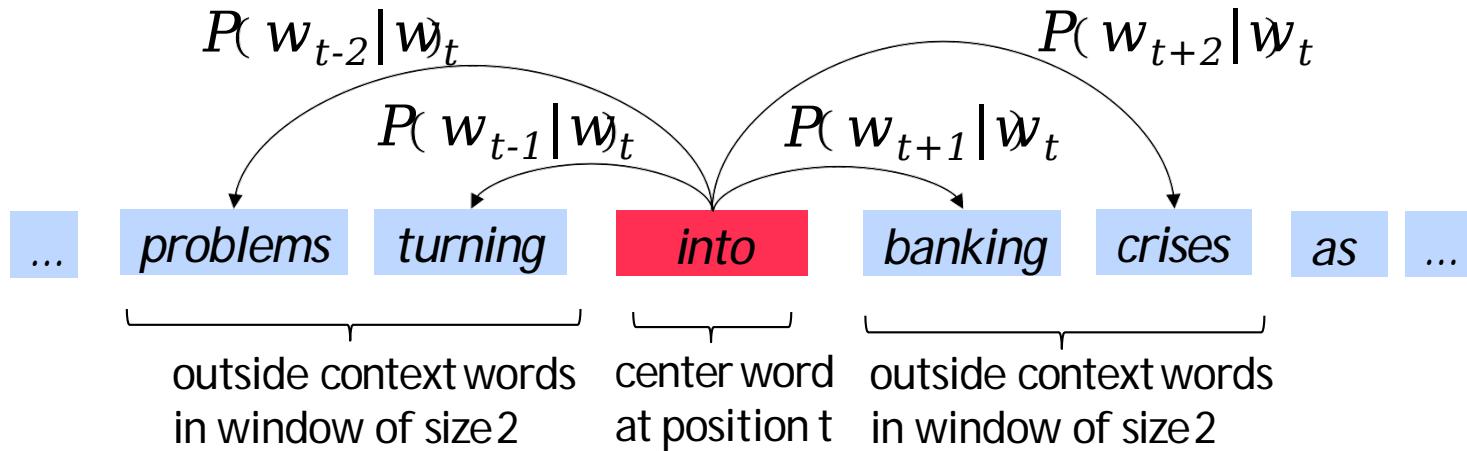
Word2vec (Mikolov et al. 2013) is a framework for learning word vectors

Idea:

- We have a large corpus of text
- Every word in a fixed vocabulary is represented by a vector
- Go through each position  $t$  in the text, which has a center word  $c$  and context (“outside”) words  $o$
- Use the similarity of the word vectors for  $c$  and  $o$  to calculate the probability of  $o$  given  $c$  (or vice versa)
- Keep adjusting the word vectors to maximize this probability

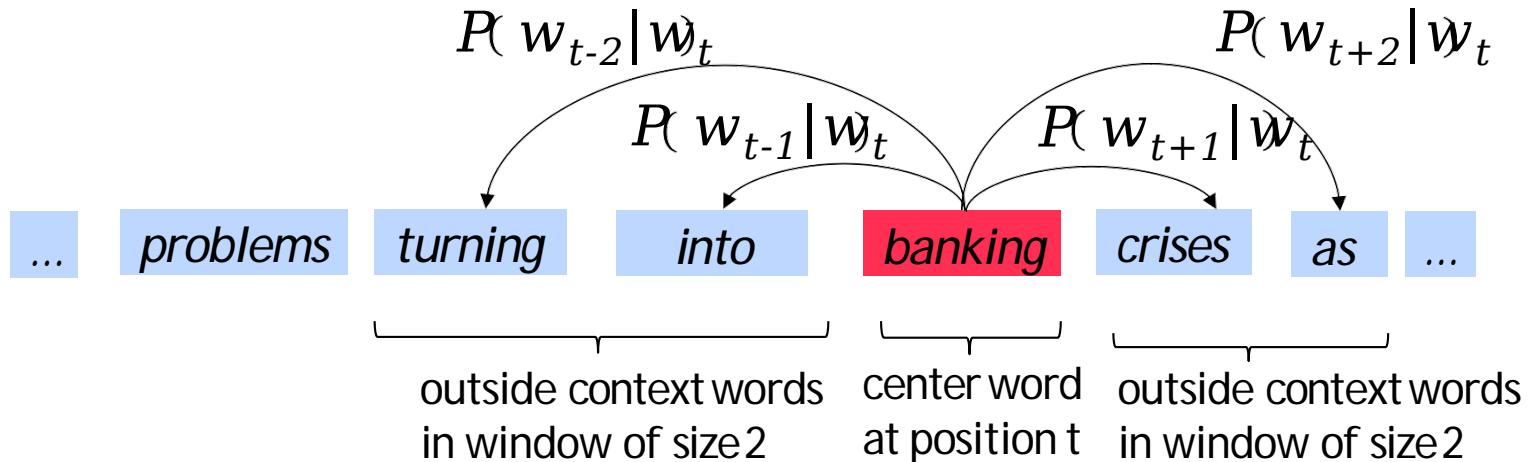
# Word2Vec Overview

- Example windows and process for computing  $P(w_{t+j}|w_t)$



# Word2Vec Overview

- Example windows and process for computing  $P(w_{t+j}|w_t)$



# Word2Vec: objective function

For each position  $t = 1, \dots, T$ , predict context words within a window of fixed size  $m$ , given center word  $w_t$ .

$$\text{Likelihood} = L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{t+j} | w_t; \theta)$$

$\theta$  is all variables to be optimized

sometimes called *cost* or *loss* function

The **objective function** is the (average) negative log likelihood:

$$J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t; \theta)$$

Minimizing objective function  $\iff$  Maximizing predictive accuracy

# Word2Vec: objective function

- We want to minimize the objective function:

$$J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t; \theta)$$

- Question: How to calculate  $P(w_{t+j}|w_t; \theta)$ ?
- Answer: We will use two vectors per word  $w$ :
  - $v_w$  when  $w$  is a center word
  - $u_w$  when  $w$  is a context word
- Then for a center word  $c$  and a context word  $o$ :

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

# Word2Vec: prediction function

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

Exponentiation makes anything positive  
Dot product compares similarity of  $o$  and  $c$ .  
 $u^T v = u \cdot v = \sum_{i=1}^n u_i v_i$   
Larger dot product = larger probability  
Normalize over entire vocabulary to give probability distribution

- This is an example of the softmax function  $\mathbb{R}^n \rightarrow \mathbb{R}^n$

$$\text{softmax}(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)} = p_i$$

- The softmax function maps arbitrary values to a probability distribution  $p_i$ 
  - “max” because amplifies probability of largest  $x_i$
  - “soft” because still assigns some probability to smaller  $x_i$

# ELMo: Embeddings from Language Models

Deep contextualized word representations. Peters et al.  
NAACL 2018. <https://arxiv.org/abs/1802.05365>

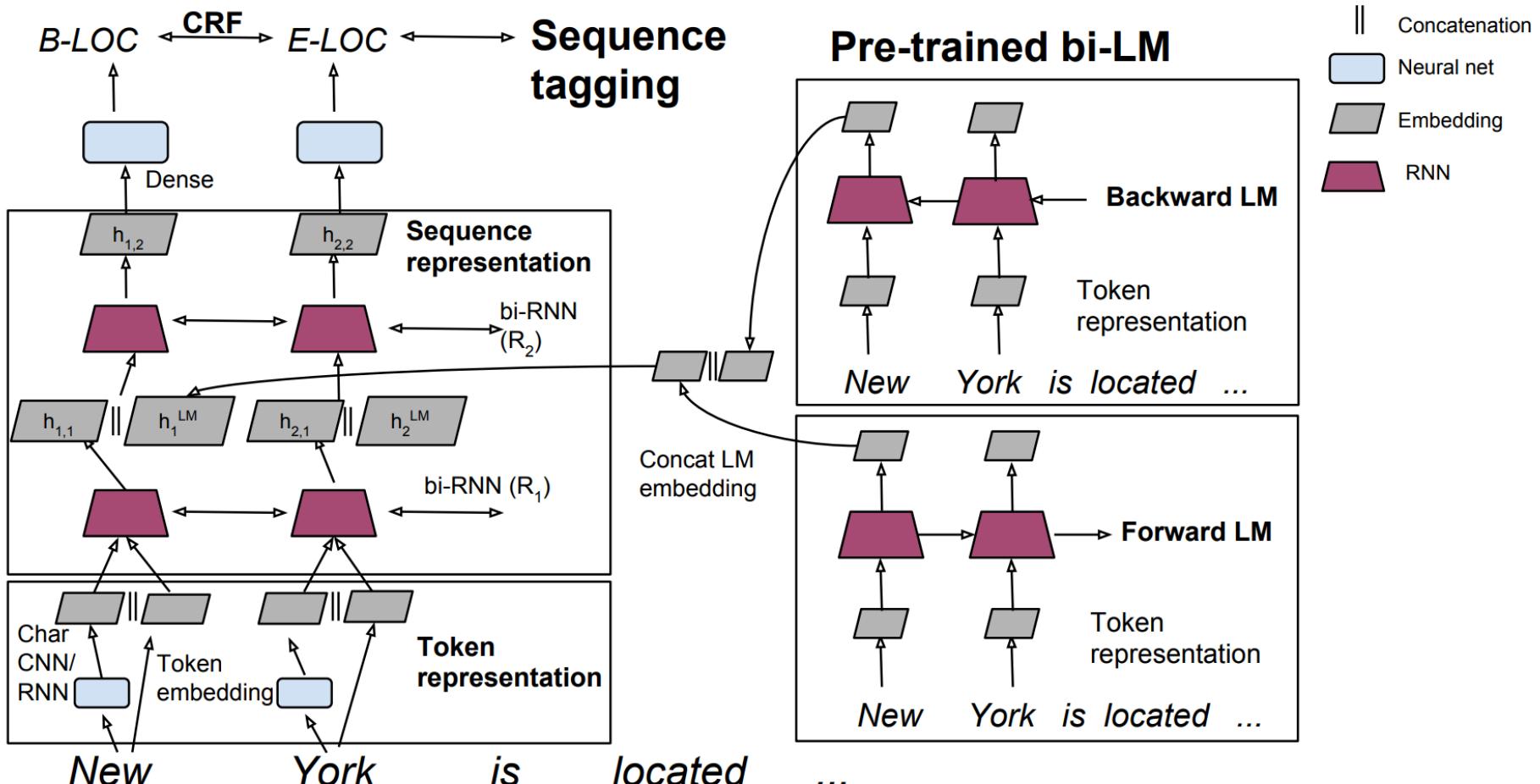
- Breakout version of word token vectors or contextual word vectors
- Learn word token vectors using long contexts not context windows (here, whole sentence, could be longer)
- Learn a deep Bi-NLM and use all its layers in prediction



# ELMo: Embeddings from Language Models

- Train a bidirectional LM
- Aim at performant but not overly large LM:
  - Use 2 biLSTM layers
  - Use character CNN to build initial word representation
  - Use 4096 dim hidden/cell LSTM states with 512 dim projections to next input
  - Use a residual connection
  - Tie parameters of token input and output (softmax) and tie these between forward and backward LMs

# ELMo used in a sequence tagger



$$\mathbf{h}_{k,1} = [\overrightarrow{\mathbf{h}}_{k,1}; \overleftarrow{\mathbf{h}}_{k,1}; \mathbf{h}_k^{LM}]$$

# ELMo results: Great for all tasks

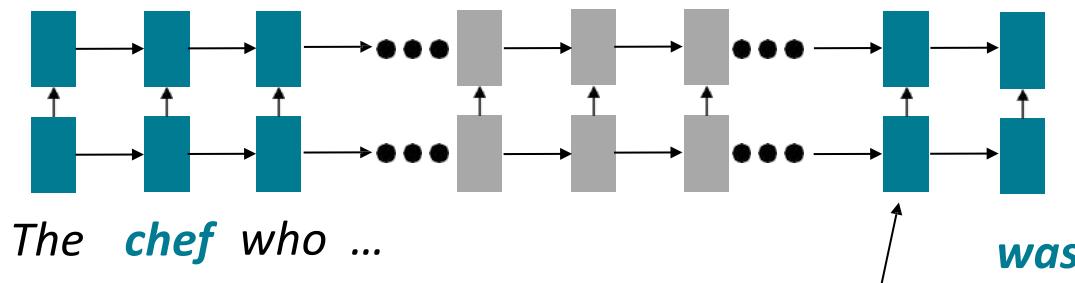
TASK	PREVIOUS SOTA		OUR BASELINE	ELMo + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SQuAD	Liu et al. (2017)	84.4	81.1	85.8	4.7 / 24.9%
SNLI	Chen et al. (2017)	88.6	88.0	88.7 ± 0.17	0.7 / 5.8%
SRL	He et al. (2017)	81.7	81.4	84.6	3.2 / 17.2%
Coref	Lee et al. (2017)	67.2	67.2	70.4	3.2 / 9.8%
NER	Peters et al. (2017)	91.93 ± 0.19	90.15	92.22 ± 0.10	2.06 / 21%
SST-5	McCann et al. (2017)	53.7	51.4	54.7 ± 0.5	3.3 / 6.8%

# ELMo: Roles of different layers

- The two biLSTM NLM layers have differentiated uses/meanings
  - Lower layer is better for lower-level syntax, etc.
    - Part-of-speech tagging, syntactic dependencies, NER
  - Higher layer is better for higher-level semantics
    - Sentiment, Semantic role labeling, question answering, SNLI

# Issues with recurrent models: Linear interaction distance

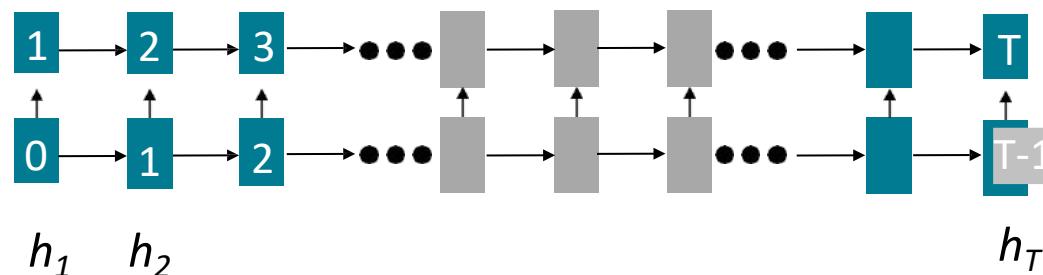
- **O(sequence length)** steps for distant word pairs to interact means:
  - Hard to learn long-distance dependencies (because gradient problems!)
  - Linear order of words is “baked in”; not necessarily the right way to think about sentences...



Info of *chef* has gone through  
 $O(\text{sequence length})$  many layers!

# Issues with recurrent models: Lack of parallelizability

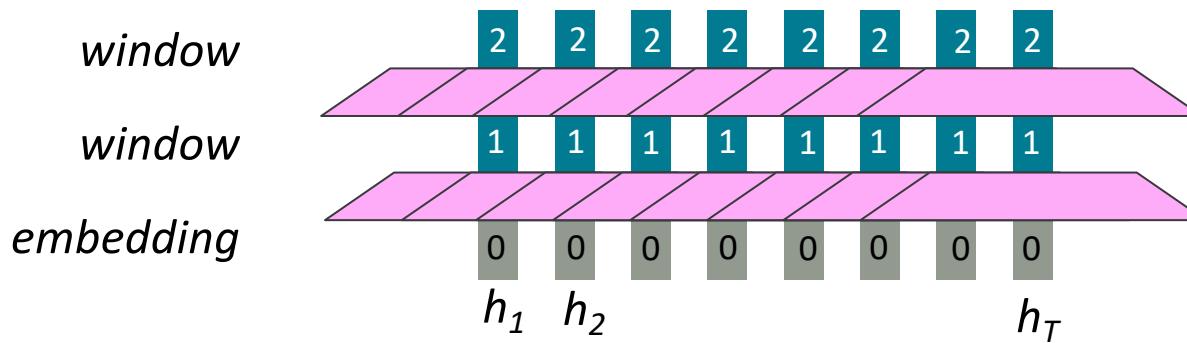
- Forward and backward passes have **O(sequence length)** unparallelizable operations
  - GPUs can perform a bunch of independent computations at once!
  - But future RNN hidden states can't be computed in full before past RNN hidden states have been computed
  - Inhibits training on very large datasets!



Numbers indicate min # of steps before a state can be computed

# If not recurrence, then what? How about word windows?

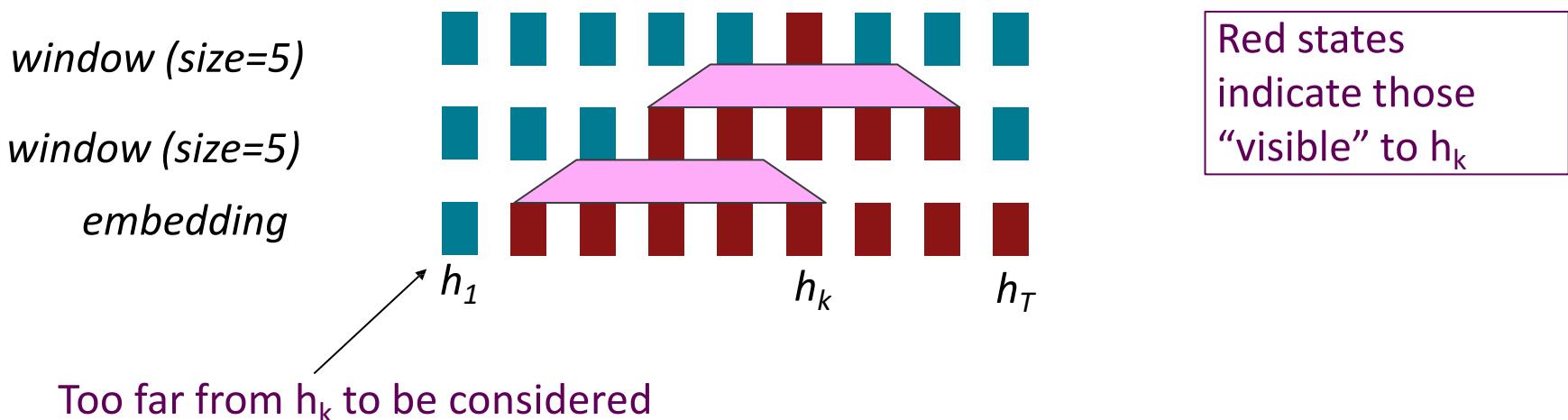
- **Word window models aggregate local contexts**
  - Also known as 1D convolution
  - Number of unparallelizable operations not tied to sequence length!



Numbers indicate min # of steps before a state can be computed

# If not recurrence, then what? How about word windows?

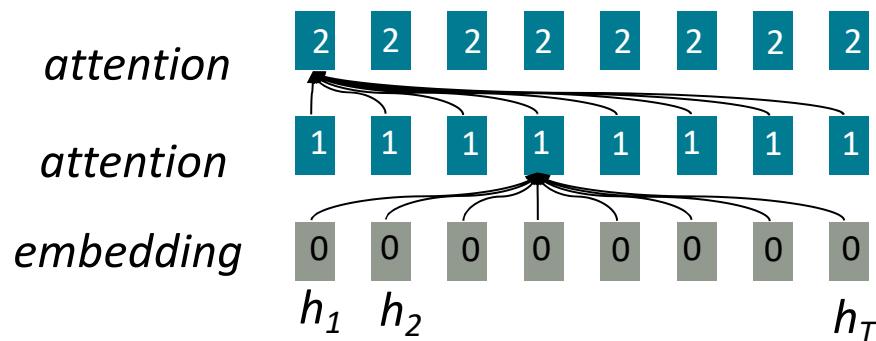
- **Word window models aggregate local contexts**
- What about long-distance dependencies?
  - Stacking word window layers allows interaction between farther words
  - But if your sequences are too long, you'll just ignore long-distance context



# If not recurrence, then what?

## How about attention?

- **Attention** treats each word's representation as a **query** to access and incorporate information from a **set of values**.
  - We saw attention from the **decoder** to the **encoder**; today we'll think about attention **within a single sentence**.
  - If **attention** gives us access to any state... maybe we can just use attention and don't need the RNN?
- Number of unparallelizable operations not tied to sequence length.
- All words interact at every layer!



All words attend to all words in previous layer; most arrows here are omitted

# Self-Attention

- Attention operates on **queries**, **keys**, and **values**.
  - We have some **queries**  $q_1, q_2, \dots, q_T$ . Each query is  $q_i \in \mathbb{R}^d$
  - We have some **keys**  $k_1, k_2, \dots, k_T$ . Each key is  $k_i \in \mathbb{R}^d$
  - We have some **values**  $v_1, v_2, \dots, v_T$ . Each value is  $v_i \in \mathbb{R}^d$
- In **self-attention**, the queries, keys, and values are drawn from the same source.
  - For example, if the output of the previous layer is  $x_1, \dots, x_T$ , (one vec per word) we could let  $v_i = k_i = q_i = x_i$  (that is, use the same vectors for all of them!)
- The (dot product) self-attention operation is as follows:

$$e_{ij} = q_i^\top k_j$$

Compute **key-query** affinities

$$\alpha_i = \frac{\exp(e_{ij})}{\sum_j \exp(e_{ij})}$$

Compute attention weights from affinities  
(softmax)

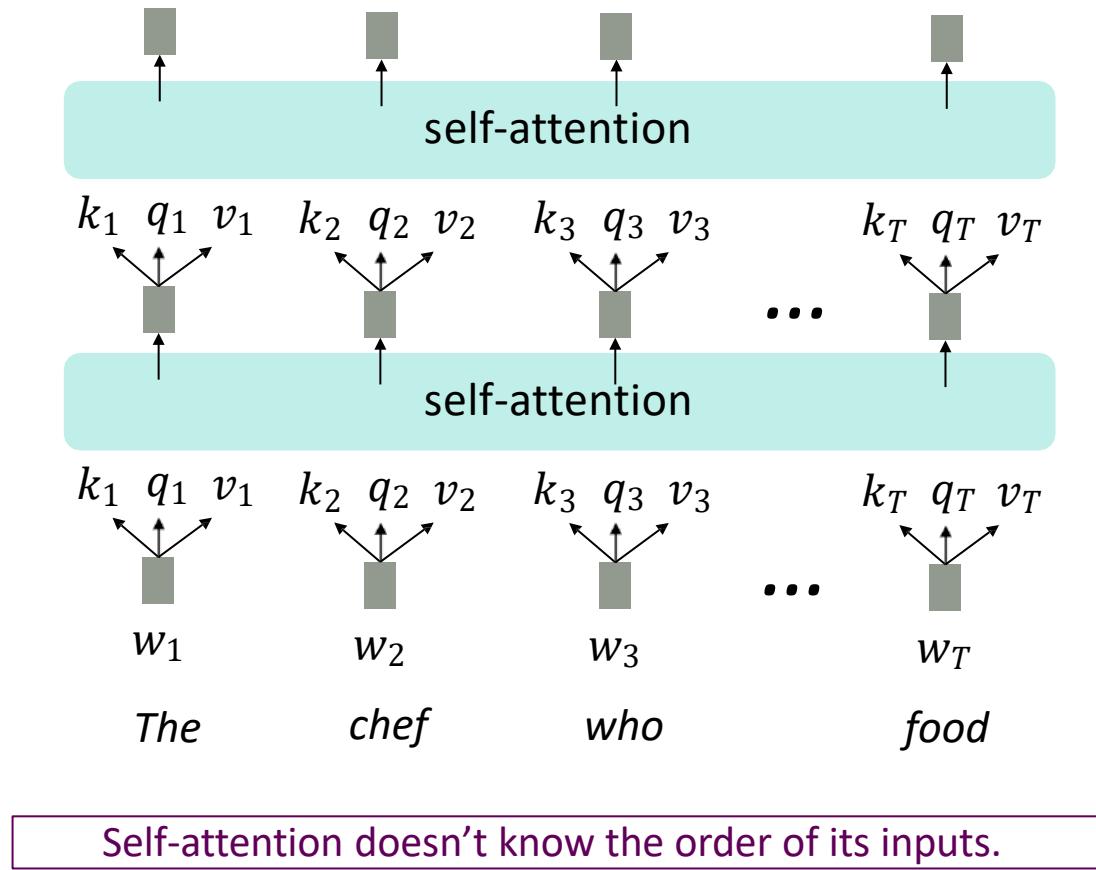
$$\text{output}_i = \sum_j \alpha_{ij} v_j$$

Compute outputs as weighted sum of **values**

The number of queries can differ from the number of keys and values in practice.

# Self-Attention

- In the diagram at the right, we have stacked self-attention blocks, like we might stack LSTM layers.
- Can self-attention be a drop-in replacement for recurrence?
- No. It has a few issues, which we'll go through.
- First, self-attention is an operation on **sets**. It has no inherent notion of order.



# Barriers and solutions for Self-Attention as a building block

## Barriers

- Doesn't have an inherent notion of order!



## Solutions

# Fixing the first self-attention problem: Sequence order

- Since self-attention doesn't build in order information, we need to encode the order of the sentence in our keys, queries, and values.
- Consider representing each **sequence index** as a **vector**

$p_i \in \mathbb{R}^d$ , for  $i \in \{1, 2, \dots, T\}$  are position vectors

- Don't worry about what the  $p_i$  are made of yet!
- Easy to incorporate this info into our self-attention block: just add the  $p_i$  to our inputs!
- Let  $v_i^{'}, k_i^{'}, q_i^{''}$  be our old values, keys, and queries.

$$\begin{aligned} v_i &= v_i^{'} + p_i \\ q_i &= q_i^{''} + p_i \\ k_i &= k_i^{'} + p_i \end{aligned}$$

In deep self-attention networks, we do this at the first layer! You could concatenate them as well, but people mostly just add...

# Barriers and solutions for Self-Attention as a building block

## Barriers

- Doesn't have an inherent notion of order!
- No nonlinearities for deep learning! It's all just weighted averages



## Solutions

- Add position representations to the inputs

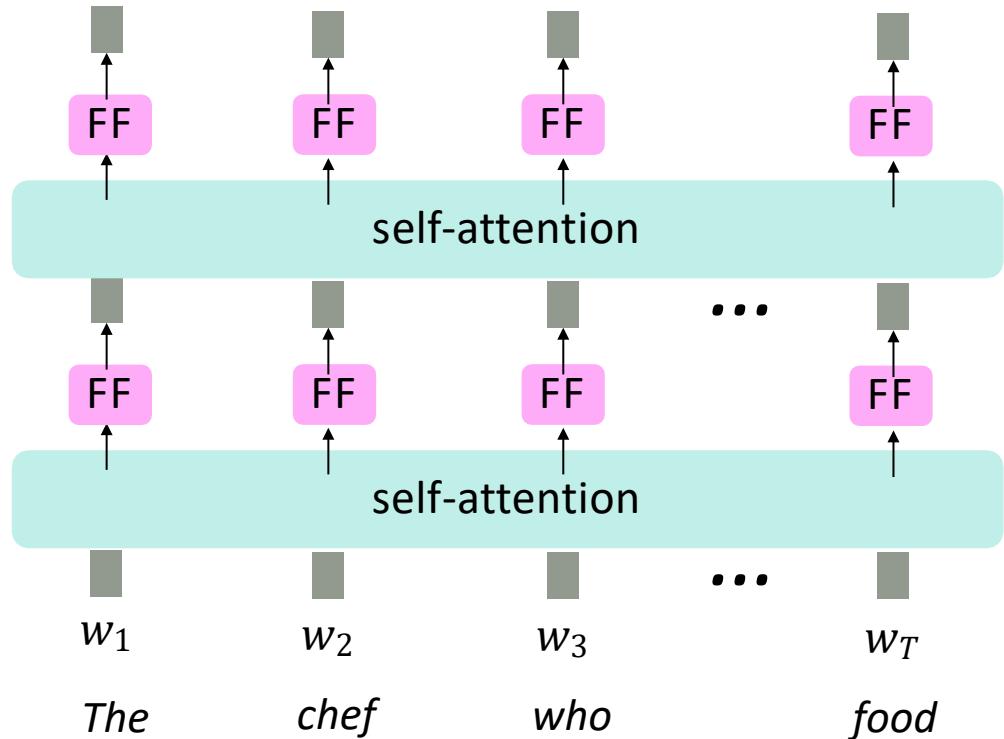


# Adding nonlinearities in self-attention

- Note that there are no elementwise nonlinearities in self-attention; stacking more self-attention layers just re-averages **value** vectors
- Easy fix: add a **feed-forward network** to post-process each output vector.

$$m_i = \text{MLP}(\text{output}_i)$$

$$= W_2 * \text{ReLU}(W_1 \times \text{output}_i + b_1) + b_2$$



Intuition: the FF network processes the result of attention

# Barriers and solutions for Self-Attention as a building block

## Barriers

- Doesn't have an inherent notion of order!
- No nonlinearities for deep learning magic! It's all just weighted averages
- Need to ensure we don't "look at the future" when predicting a sequence
  - Like in machine translation
  - Or language modeling



## Solutions

- Add position representations to the inputs
- Easy fix: apply the same feedforward network to each self-attention output.

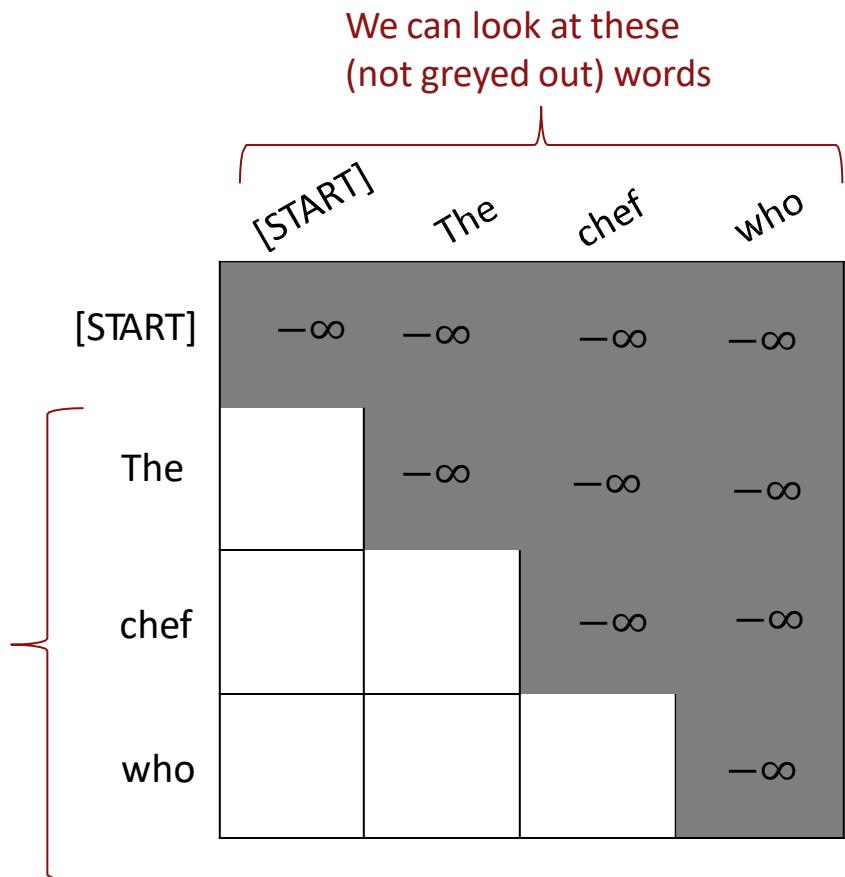


# Masking the future in self-attention

- To use self-attention in **decoders**, we need to ensure we can't peek at the future.
- At every timestep, we could change the set of **keys and queries** to include only past words. (Inefficient!)
- To enable parallelization, we **mask out attention** to future words by setting attention scores to  $-\infty$ .

$$e_{ij} = \begin{cases} q_i^\top k_j, & j < i \\ -\infty, & j \geq i \end{cases}$$

For encoding  
these words



# Barriers and solutions for Self-Attention as a building block

## Barriers

- Doesn't have an inherent notion of order!
- No nonlinearities for deep learning magic! It's all just weighted averages
- Need to ensure we don't "look at the future" when predicting a sequence
  - Like in machine translation
  - Or language modeling

## Solutions

- Add position representations to the inputs
- Easy fix: apply the same feedforward network to each self- attention output.
- Mask out the future by artificially setting attention weights to 0!

# Necessities for a self-attention building block:

- **Self-attention:**
  - the basis of the method.
- **Position representations:**
  - Specify the sequence order, since self-attention is an unordered function of its inputs.
- **Nonlinearities:**
  - At the output of the self-attention block
  - Frequently implemented as a simple feed-forward network.
- **Masking:**
  - In order to parallelize operations while not looking at the future.
  - Keeps information about the future from “leaking” to the past.
- That’s it! But this is not the **Transformer** model we’ve been hearing about.

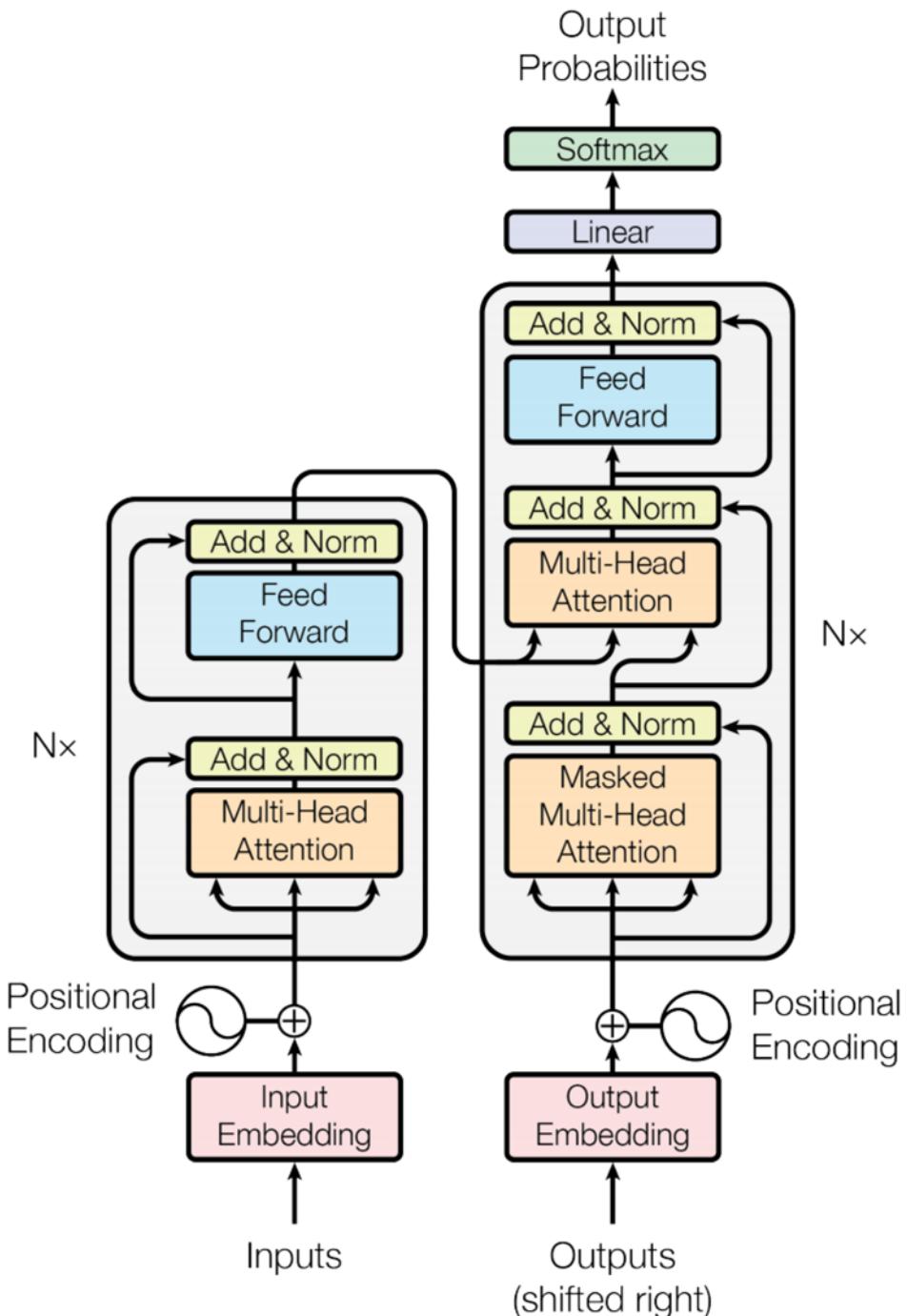
# Transformer Overview

Attention is all you need. 2017. Aswani, Shazeer, Parmar, Uszkoreit, Jones, Gomez, Kaiser, Polosukhin

<https://arxiv.org/pdf/1706.03762.pdf>

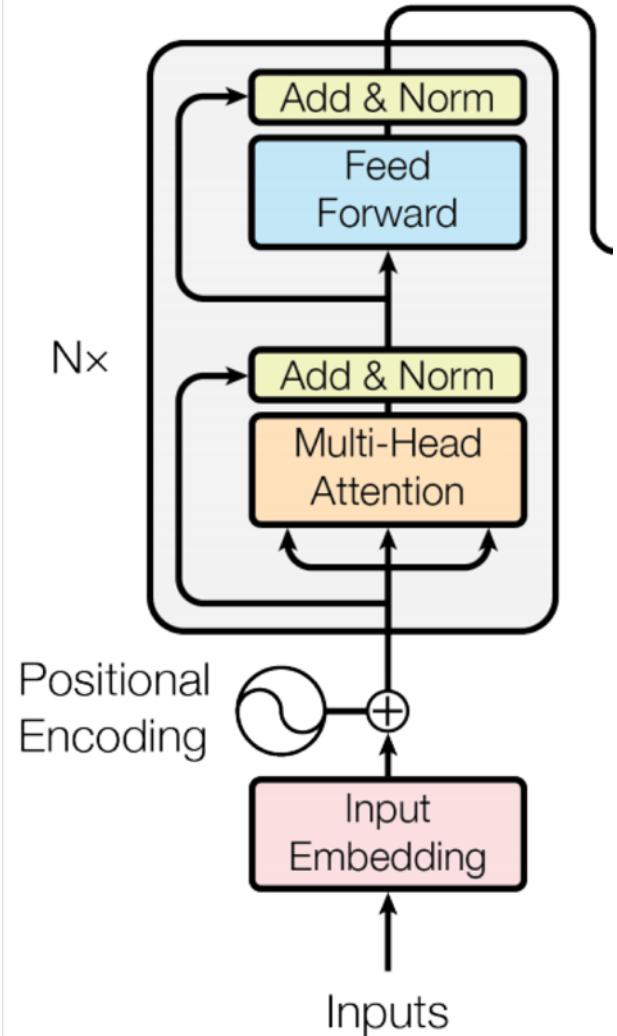
- Non-recurrent sequence-to-sequence encoder-decoder model
- Task: machine translation with parallel corpus
- Predict each translated word
- Final cost/error function is standard cross-entropy error on top of a softmax classifier

This and related figures from paper ↑



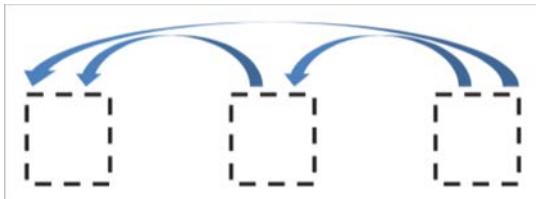
# Transformer Encoder

- For encoder, at each block, we use the same Q, K and V from the previous layer
- Blocks are repeated 6 times (in vertical stack)

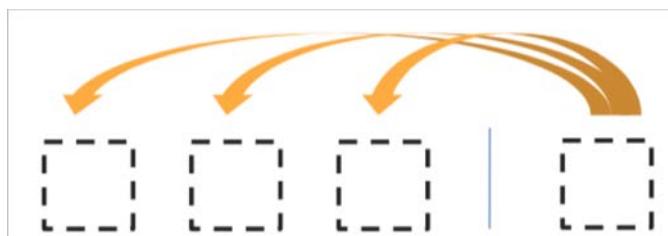


# Transformer Decoder

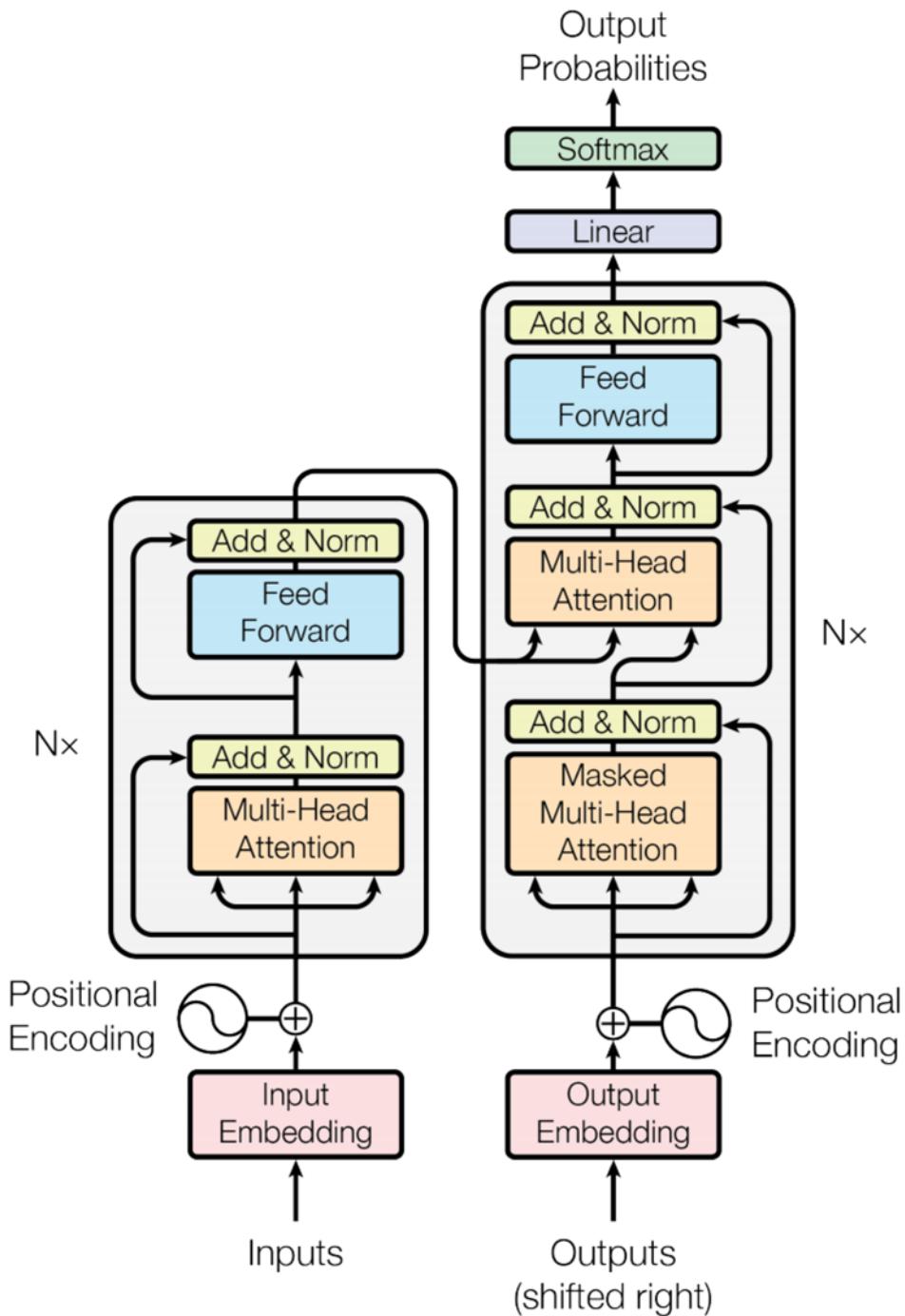
- 2 sublayer changes in decoder
- Masked decoder self-attention on previously generated outputs:



- Encoder-Decoder Attention, where queries come from previous decoder layer and keys and values come from output of encoder



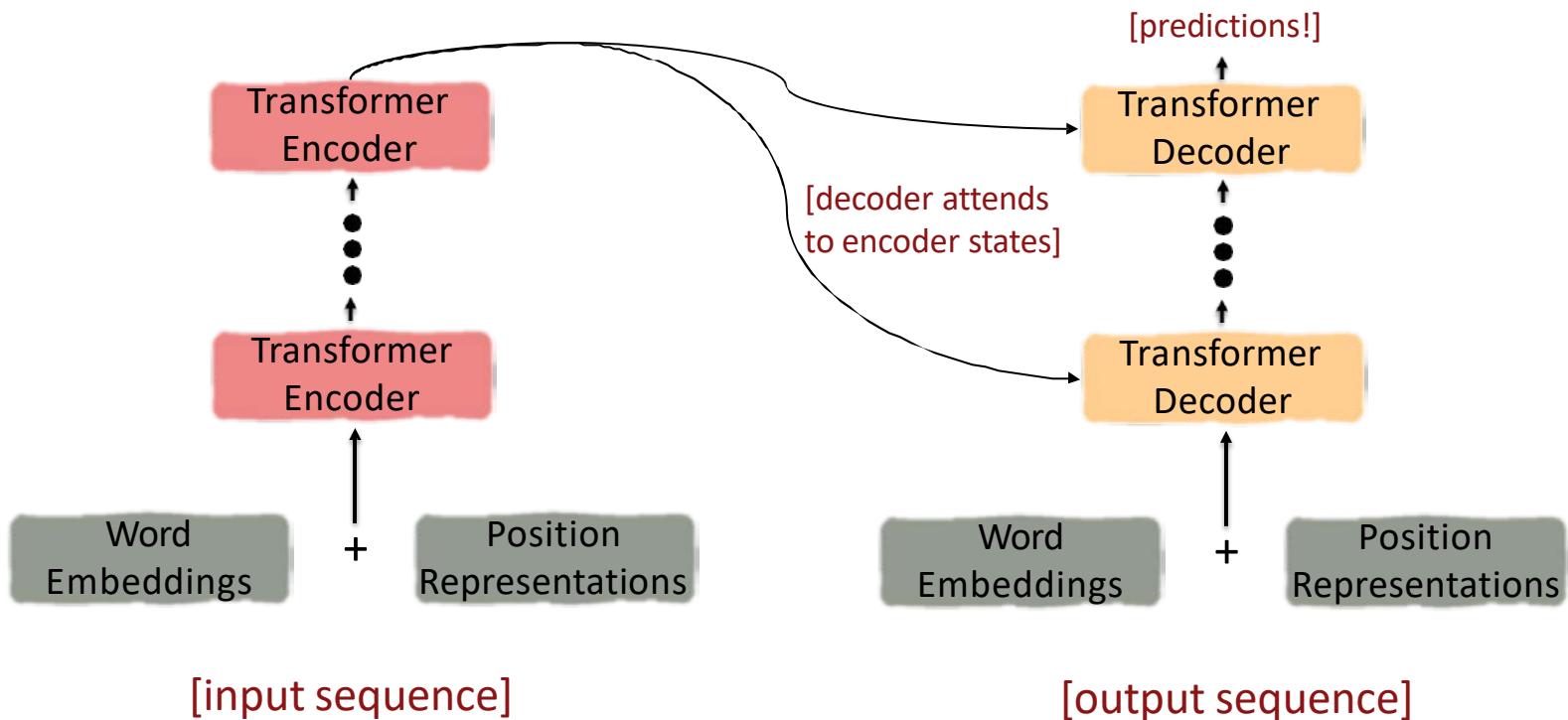
Blocks repeated 6 times also



# The Transformer Encoder-Decoder

## [Vaswani et al., 2017]

Another look at the Transformer Encoder and Decoder Blocks at a high level



# The Transformer Encoder-Decoder

## [Vaswani et al., 2017]

Next, let's look at the Transformer Encoder and Decoder Blocks

What's left in a Transformer Encoder Block that we haven't covered?

1. **Key-query-value attention:** How do we get the  $k, q, v$  vectors from a single word embedding?
2. **Multi-headed attention:** Attend to multiple places in a single layer!
3. **Tricks to help with training!**
  1. Residual connections
  2. Layer normalization
  3. Scaling the dot product
  4. These tricks **don't improve** what the model is able to do; they help improve the training process

# The Transformer Encoder: Dot-Product Attention

- Inputs: a query  $q$  and a set of key-value ( $k$ - $v$ ) pairs to an output
- Query, keys, values, and output are all vectors
- Output is weighted sum of values, where
- Weight of each value is computed by an inner product of query and corresponding key
- Queries and keys have same dimensionality  $d_k$ , value have  $d_v$

$$A(q, K, V) = \sum_i \frac{e^{q \cdot k_i}}{\sum_j e^{q \cdot k_j}} v_i$$

# The Transformer Encoder: Dot-Product Attention – Matrix notation

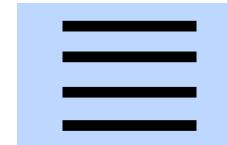
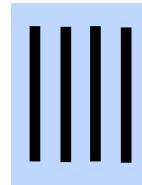
- When we have multiple queries  $q$ , we stack them in a matrix  $Q$ :

$$A(q, K, V) = \sum_i \frac{e^{q \cdot k_i}}{\sum_j e^{q \cdot k_j}} v_i$$

- Becomes:  $A(Q, K, V) = \text{softmax}(QK^T)V$

$$[|Q| \times d_k] \times [d_k \times |K|] \times [|K| \times d_v]$$

softmax  
row-wise



$$= [ |Q| \times d_v ]$$

# The Transformer Encoder: Key-Query-Value Attention

- We saw that self-attention is when keys, queries, and values come from the same source. The Transformer does this in a particular way:
  - Let  $x_1, \dots, x_T$  be input vectors to the Transformer encoder;  $x_i \in \mathbb{R}^d$
- Then keys, queries, values are:
  - $k_i = Kx_i$ , where  $K \in \mathbb{R}^{d \times d}$  is the key matrix.
  - $q_i = Qx_i$ , where  $Q \in \mathbb{R}^{d \times d}$  is the query matrix.
  - $v_i = Vx_i$ , where  $V \in \mathbb{R}^{d \times d}$  is the value matrix.
- These matrices allow *different aspects* of the  $x$  vectors to be used/emphasized in each of the three roles.

# The Transformer Encoder: Key-Query-Value Attention

- Let's look at how key-query-value attention is computed, in matrices.
  - Let  $X = [x_1; \dots; x_T] \in \mathbb{R}^{T \times d}$  be the concatenation of input vectors.
  - First, note that  $XK \in \mathbb{R}^{T \times d}$ ,  $XQ \in \mathbb{R}^{T \times d}$ ,  $XV \in \mathbb{R}^{T \times d}$ .
  - The output is defined as  $\text{output} = \text{softmax}(XQ(XK)^\top) \times XV$ .

First, take the query-key dot products in one matrix multiplication:  $XQ(XK)$

$$\begin{array}{ccc} XQ & & \\ & K^\top X^\top & \\ & & = XQK^\top X^\top \\ & & \in \mathbb{R}^{T \times T} \end{array}$$

All pairs of  
attention scores!

softmax  $\left( \begin{array}{c} XQK^\top X^\top \\ \hline \end{array} \right) XV = \text{output} \in \mathbb{R}^{T \times d}$

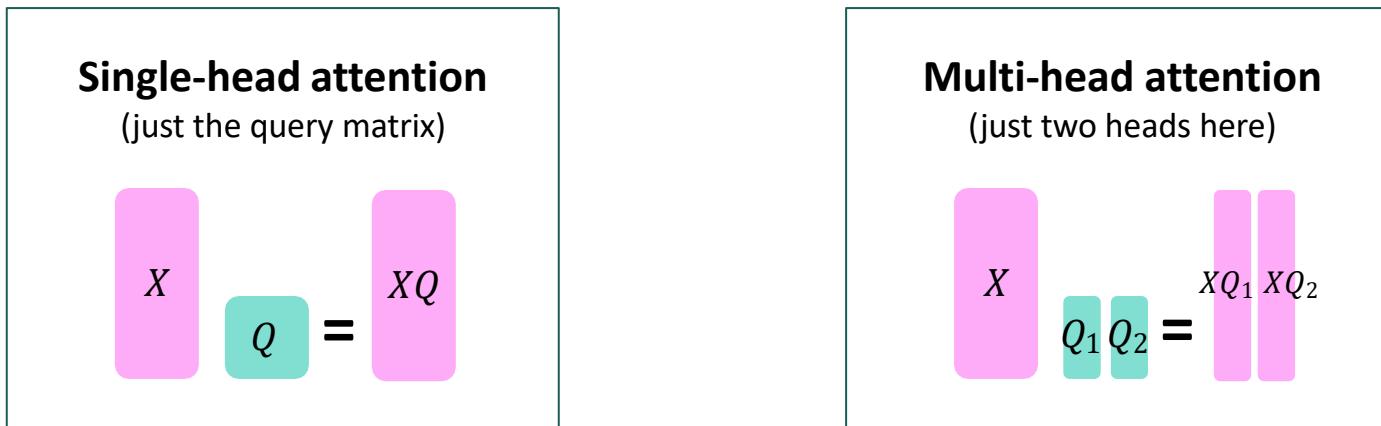
Next, softmax, and compute the weighted average with another matrix multiplication.

# The Transformer Encoder: Multi-headed attention

- What if we want to look in multiple places in the sentence at once?
  - For word  $i$ , self-attention “looks” where  $x^T Q^T K x_j$  is high, but maybe we want to focus on different  $j$  for different reasons?
- We’ll define **multiple attention “heads”** through multiple Q,K,V matrices
- Let,  $Q_P, K_P, V_P \in \mathbb{R}^{d \times \frac{d}{h}}$ , where  $h$  is the number of attention heads, and  $P$  ranges from 1 to  $h$ .
- Each attention head performs attention independently:
  - $\text{output}_P = \text{softmax}(X Q_P K_P^T X^T) * X V_P$ , where  $\text{output}_P \in \mathbb{R}^{d/h}$
- Then the outputs of all the heads are combined!
  - $\text{output} = Y[\text{output}_1; \dots; \text{output}_h]$ , where  $Y \in \mathbb{R}^{d \times d}$
- Each head gets to “look” at different things, and construct value vectors differently.

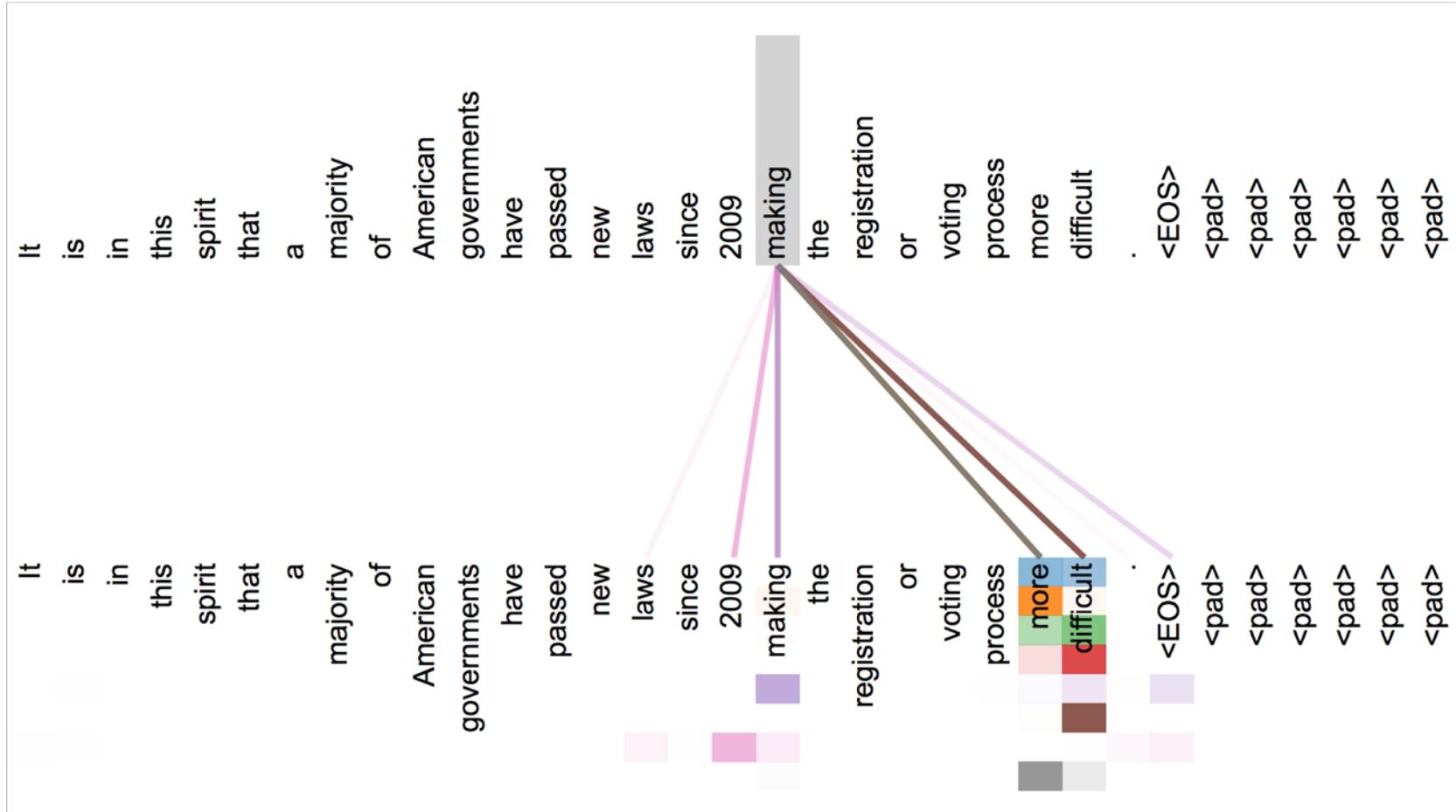
# The Transformer Encoder: Multi-headed attention

- What if we want to look in multiple places in the sentence at once?
  - For word  $i$ , self-attention “looks” where  $x^T Q^T K x_j$  is high, but maybe we want to focus on different  $j$  for different reasons?
- We’ll define **multiple attention “heads”** through multiple Q,K,V matrices
- Let,  $Q_P, K_P, V_P \in \mathbb{R}^{d \times \frac{d}{h}}$ , where  $h$  is the number of attention heads, and  $P$  ranges from 1 to  $h$ .

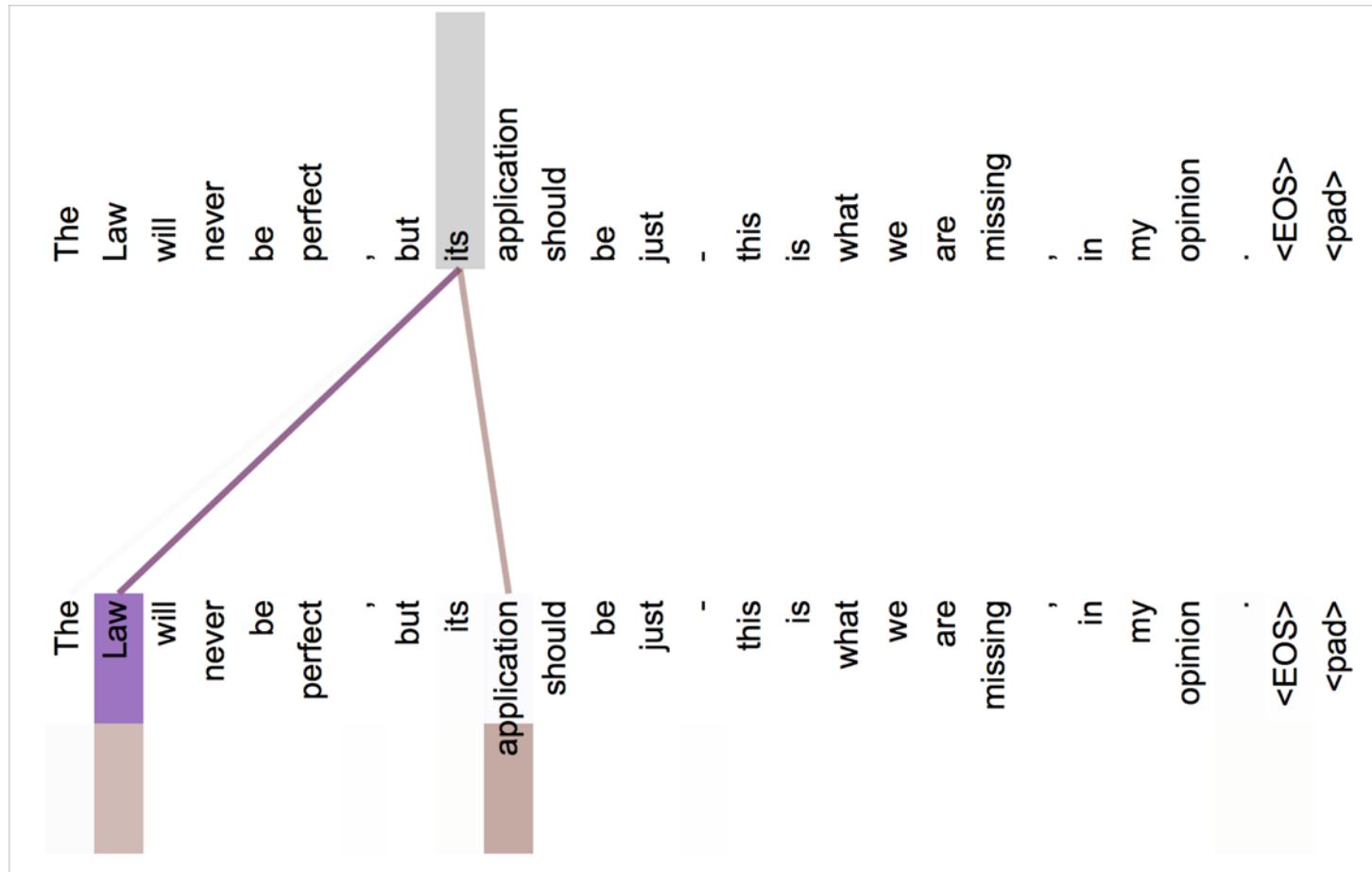


# Attention visualization in layer 5

- Words start to pay attention to other words in sensible ways

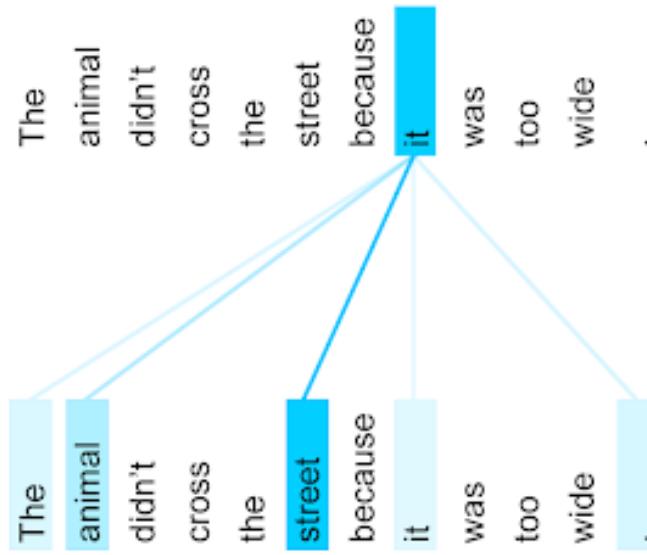
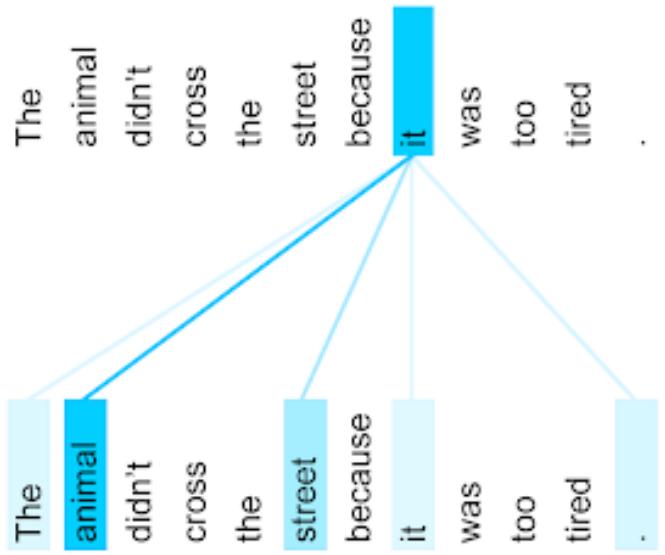


# Attention visualization: Implicit anaphora resolution



In 5<sup>th</sup> layer. Isolated attentions from just the word 'its' for attention heads 5 and 6.  
Note that the attentions are very sharp for this word.

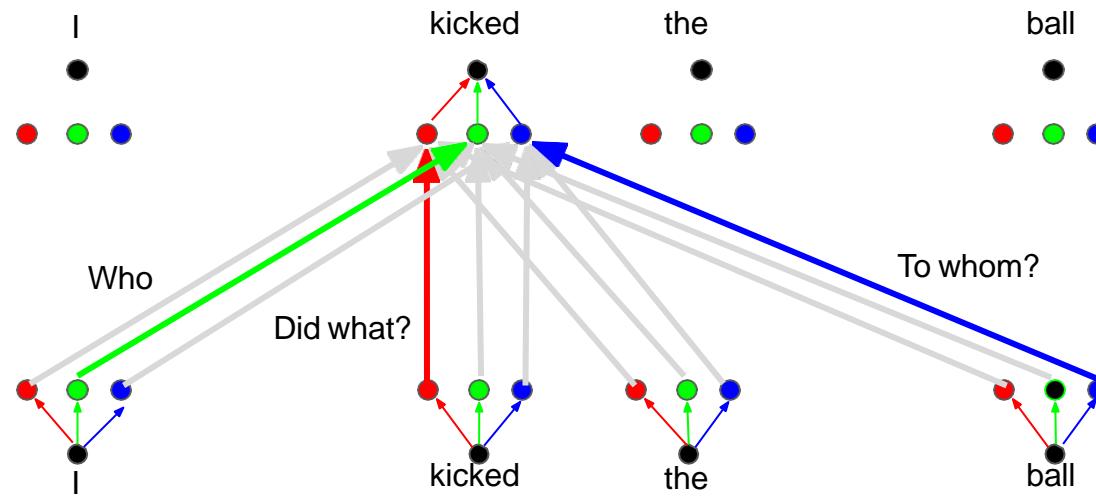
# Attention Visualization



The encoder self-attention distribution for the word "it" from the 5th to the 6th layer of a Transformer trained on English to French translation (one of eight attention heads).

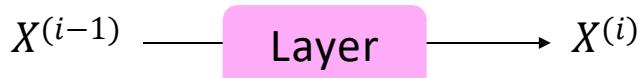
<https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>

# Parallel attention heads

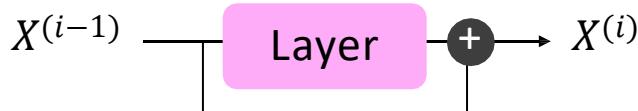


# The Transformer Encoder: Residual connections [He et al., 2016]

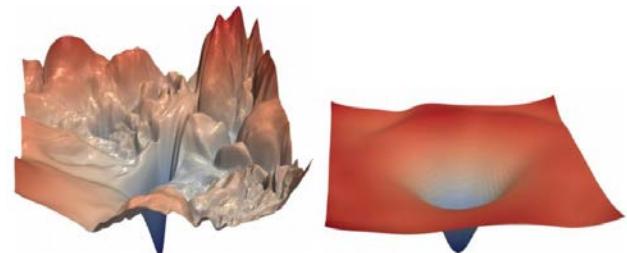
- **Residual connections** are a trick to help models train better.
  - Instead of  $X^{(i)} = \text{Layer}(X^{(i-1)})$  (where  $i$  represents the layer)



- We let  $X^{(i)} = X^{(i-1)} + \text{Layer}(X^{(i-1)})$  (so we only have to learn “the residual” from the previous layer)



- Residual connections are thought to make the loss landscape considerably smoother (thus easier training!)



[no residuals] [residuals]

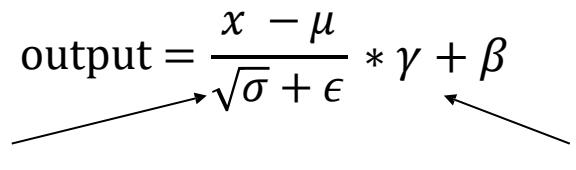
[Loss landscape visualization,  
Li et al., 2018, on a ResNet]

# The Transformer Encoder: Layer normalization [Ba et al., 2016]

- **Layer normalization** is a trick to help models train faster.
- Idea: cut down on uninformative variation in hidden vector values by normalizing to unit mean and standard deviation **within each layer**.
  - LayerNorm's success may be due to its normalizing gradients [[Xu et al., 2019](#)]
- Let  $x \in \mathbb{R}^d$  be an individual (word) vector in the model.
- Let  $\mu = \sum_{j=1}^d x_j$ ; this is the mean;  $\mu \in \mathbb{R}$ .
  - Let  $\sigma = \sqrt{\frac{1}{d} \sum_{j=1}^d (x_j - \mu)^2}$ ; this is the standard deviation;  $\sigma \in \mathbb{R}$ .
  - Let  $\gamma \in \mathbb{R}^d$  and  $\beta \in \mathbb{R}^d$  be learned "gain" and "bias" parameters. (Can omit!)
  - Then layer normalization computes:

$$\text{output} = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} * \gamma + \beta$$

Normalize by scalar mean and variance      Modulate by learned elementwise gain and bias



# The Transformer Encoder: Scaled Dot Product [Vaswani et al., 2017]

- “Scaled Dot Product” attention is a final variation to aid in Transformer training.
- When dimensionality  $d$  becomes large, dot products between vectors become large, inputs to the softmax function can be large, making gradients small.

- Instead of the self-attention function we've seen:

$$\text{output}_P = \text{softmax}(XQ_PK_P^TX^T) * XV_P$$

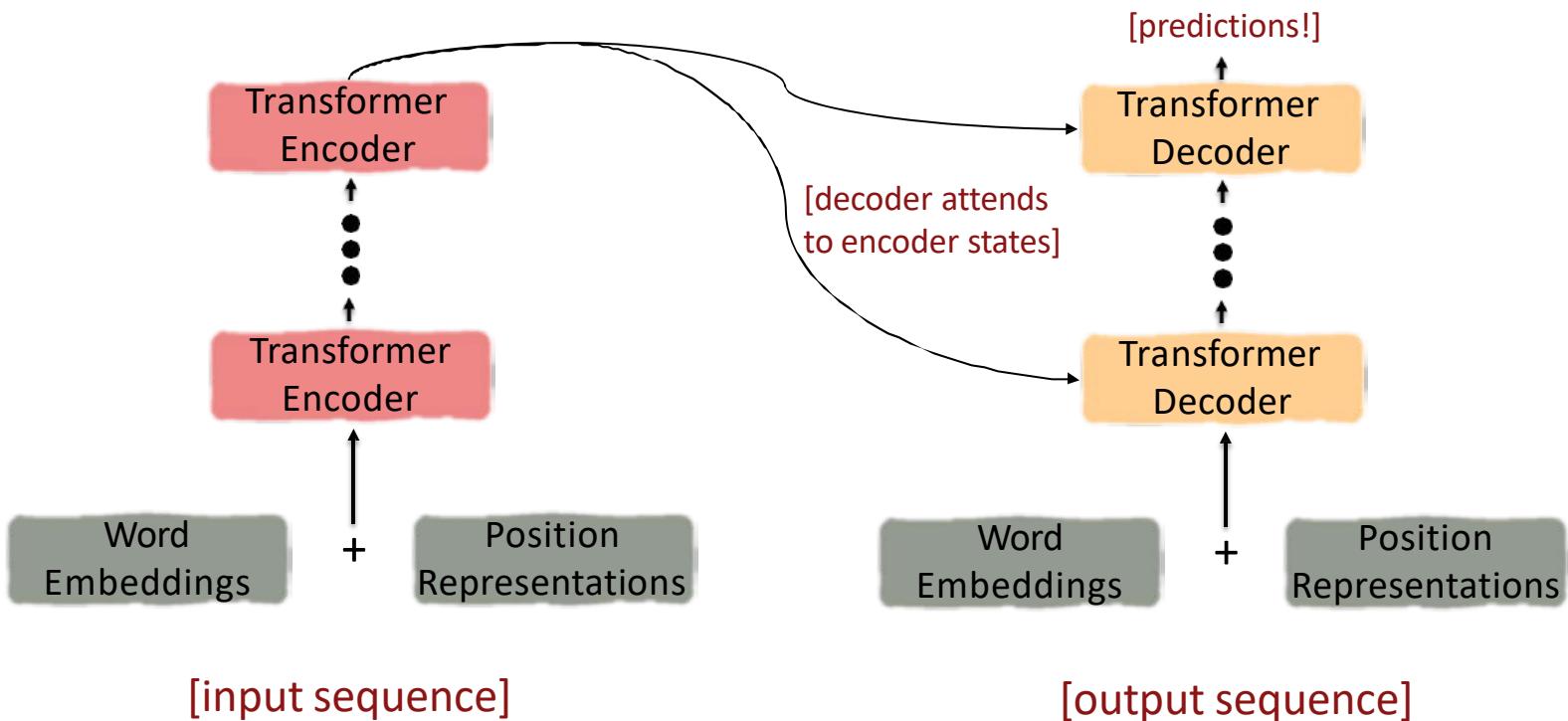
- We divide the attention scores by  $\sqrt{d/h}$ , to stop the scores from becoming large just as a function of  $d/h$  (The dimensionality divided by the number of heads.)

$$\text{output}_P = \text{softmax}\left(\frac{XQ_PK_P^TX^T}{\sqrt{d/h}}\right) * XV_P$$

# The Transformer Encoder-Decoder

## [Vaswani et al., 2017]

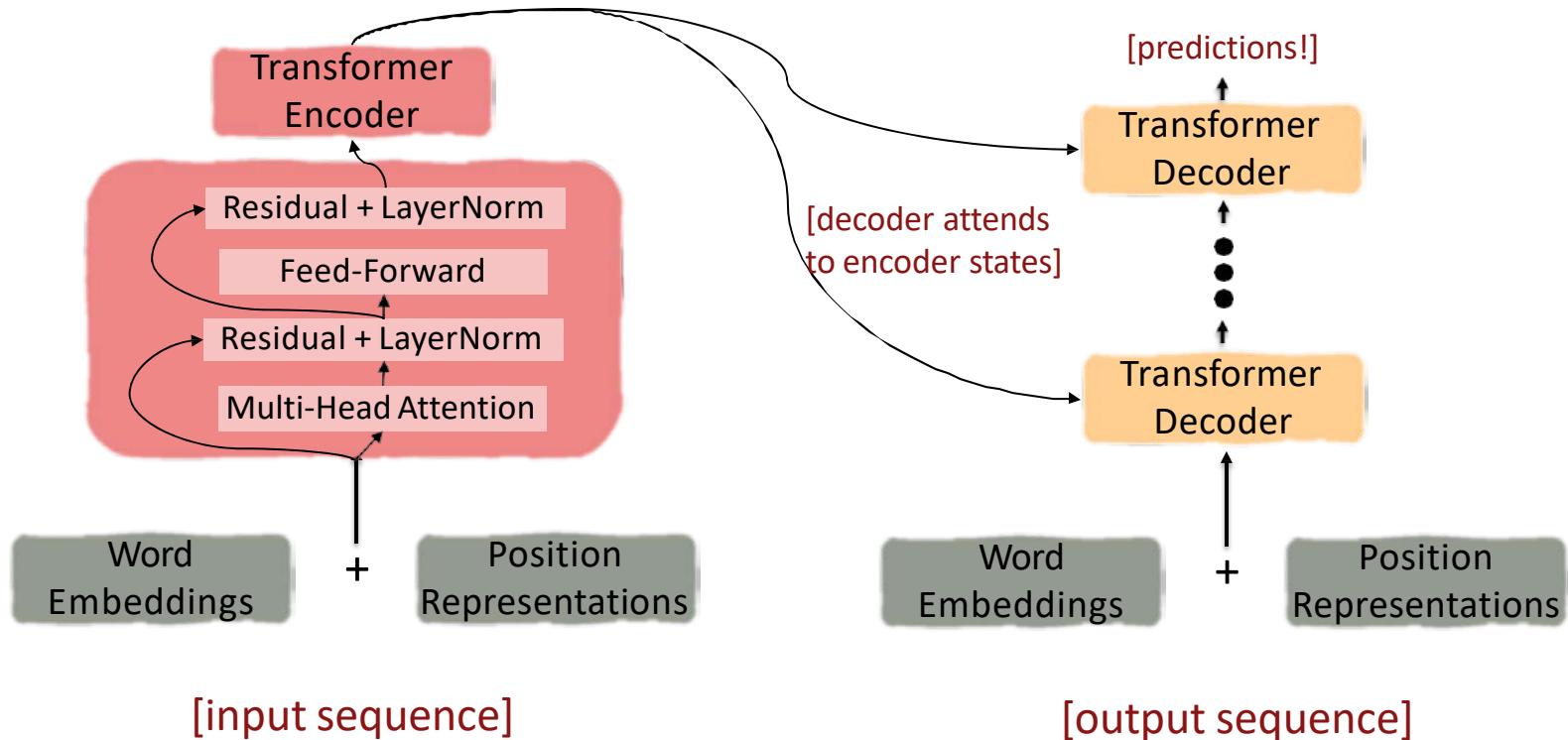
Looking back at the whole model, zooming in on an Encoder block:



# The Transformer Encoder-Decoder

## [Vaswani et al., 2017]

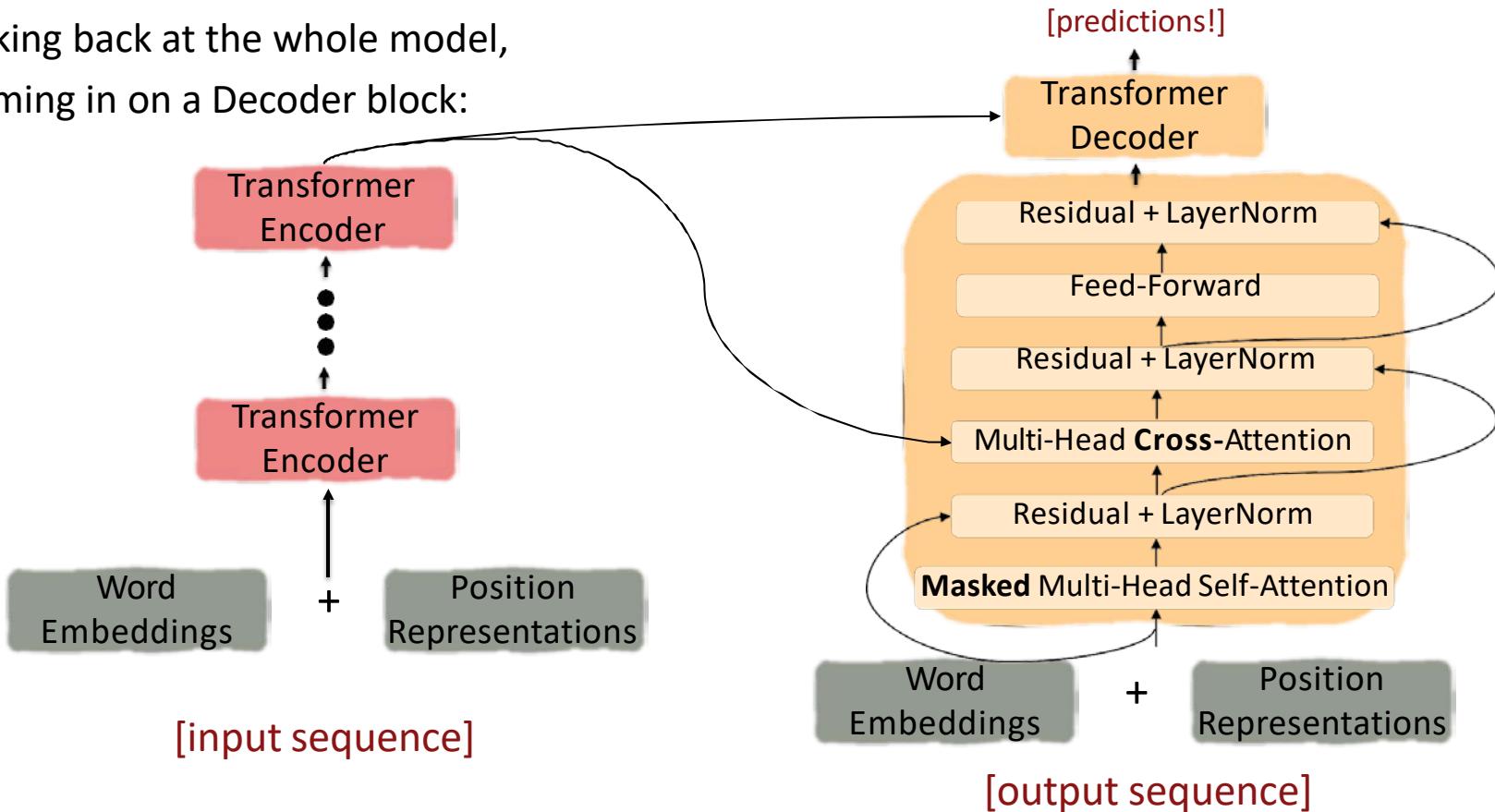
Looking back at the whole model, zooming in on an Encoder block:



# The Transformer Encoder-Decoder

## [Vaswani et al., 2017]

Looking back at the whole model,  
zooming in on a Decoder block:



# The Transformer Decoder: Cross-attention (details)

- We saw self-attention is when keys, queries, and values come from the same source.
- In the decoder, we have attention that looks more like what we saw last week.
- Let  $h_1, \dots, h_T$  be **output** vectors from the Transformer **encoder**;  $x_i \in \mathbb{R}^d$
- Let  $z_1, \dots, z_T$  be input vectors from the Transformer **decoder**,  $z_i \in \mathbb{R}^d$
- Then keys and values are drawn from the **encoder** (like a memory):
  - $k_i = Kh_i$ ,  $v_i = Vh_i$ .
- And the queries are drawn from the **decoder**,  $q_i = Qz_i$ .

# The Transformer Encoder: Cross-attention (details)

- Let's look at how cross-attention is computed, in matrices.
  - Let  $H = [h_1; \dots; h_T] \in \mathbb{R}^{T \times d}$  be the concatenation of encoder vectors.
  - Let  $Z = [z_1; \dots; z_T] \in \mathbb{R}^{T \times d}$  be the concatenation of decoder vectors.
  - The output is defined as  $\text{output} = \text{softmax}(ZQ(HK^\top)) \times HV$ .

First, take the query-key dot products in one matrix multiplication:  $ZQ(HK^\top)$

$$\begin{array}{ccc} ZQ & \quad K^\top H^\top & = ZQK^\top H^\top \\ & & \qquad \qquad \qquad \text{All pairs of attention scores!} \\ & & \in \mathbb{R}^{T \times T} \end{array}$$

softmax  $\left( \begin{array}{c} ZQK^\top H^\top \\ \hline \end{array} \right) HV = \text{output} \in \mathbb{R}^{T \times d}$

Next, softmax, and compute the weighted average with another matrix multiplication.

# Great Results with Transformers

Next, document generation!

Model	Test perplexity	ROUGE-L
<i>seq2seq-attention, L = 500</i>	5.04952	12.7
<i>Transformer-ED, L = 500</i>	2.46645	34.2
<i>Transformer-D, L = 4000</i>	2.22216	33.6
<i>Transformer-DMCA, no MoE-layer, L = 11000</i>	2.05159	36.2
<i>Transformer-DMCA, MoE-128, L = 11000</i>	1.92871	37.9
<i>Transformer-DMCA, MoE-256, L = 7500</i>	1.90325	38.8

The old standard

Transformers all the way down.

[[Liu et al., 2018](#)]; WikiSum dataset

# What would we like to fix about the Transformer?

- **Quadratic compute in self-attention:**
  - Computing all pairs of interactions means our computation grows **quadratically** with the sequence length!
  - For recurrent models, it only grew linearly!
- **Position representations:**
  - Are simple absolute indices the best we can do to represent position?
  - Relative linear position attention [\[Shaw et al., 2018\]](#)
  - Dependency syntax-based position [\[Wang et al., 2019\]](#)

# Quadratic computation as function of seq. length

- One of the benefits of self-attention over recurrence was that it's highly parallelizable.
- However, its total number of operations grows as  $O(T^2d)$ , where  $T$  is the sequence length, and  $d$  is the dimensionality.

$$\begin{matrix} XQ \\ K^\top X^\top \end{matrix} = \boxed{XQK^\top X^\top} \in \mathbb{R}^{T \times T}$$

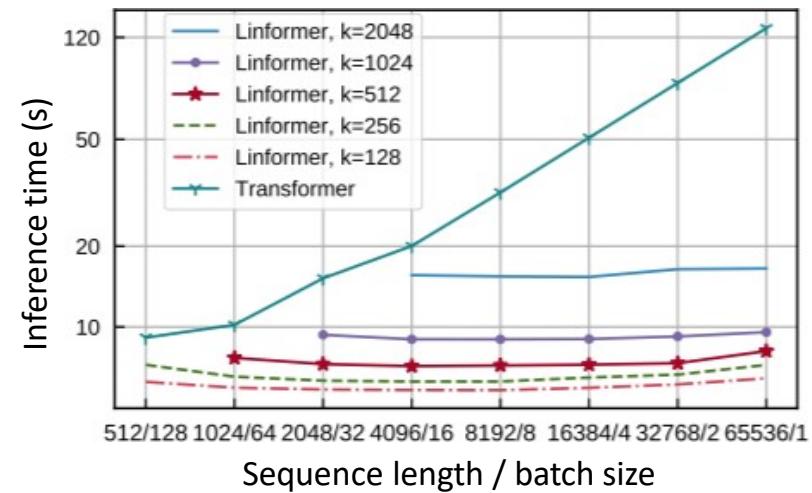
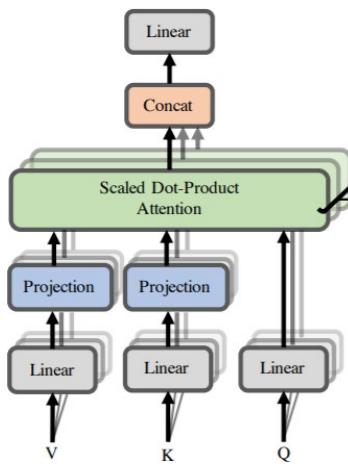
Need to compute all pairs of interactions!  
 $O(T^2d)$

- Think of  $d$  as around **1,000**.
  - So, for a single (shortish) sentence,  $T \leq 30; T^2 \leq 900$ .
  - In practice, we set a bound like  $T = 512$ .
  - **But what if we'd like  $T \geq 10,000$ ?** For example, to work on long documents?

# Recent work on improving on quadratic self-attention cost

- Considerable recent work has gone into the question, *Can we build models like Transformers without paying the  $O(T^2)$  all-pairs self-attention cost?*
- For example, **Linformer** [Wang et al., 2020]

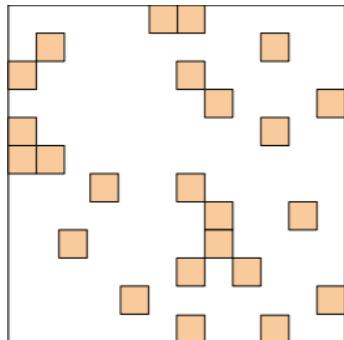
Key idea: map the sequence length dimension to a lower-dimensional space for values, keys



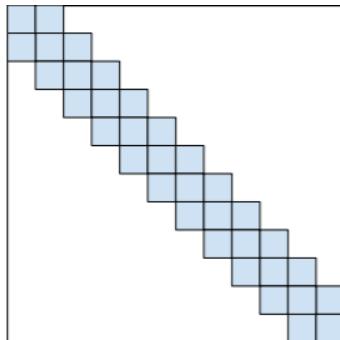
# Recent work on improving on quadratic self-attention cost

- Considerable recent work has gone into the question, *Can we build models like Transformers without paying the  $O(T^2)$  all-pairs self-attention cost?*
- For example, **BigBird** [\[Zaheer et al., 2021\]](#)

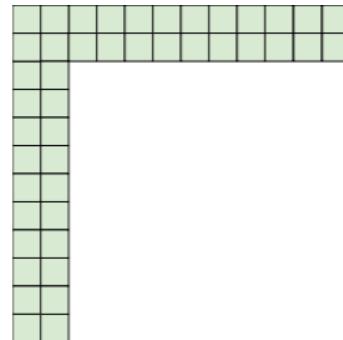
Key idea: replace all-pairs interactions with a family of other interactions, **like local windows, looking at everything, and random interactions.**



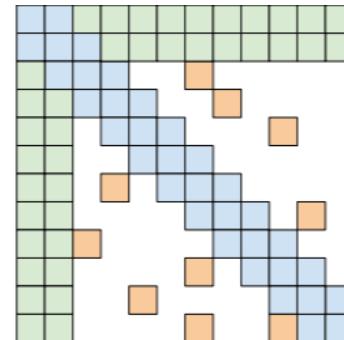
(a) Random attention



(b) Window attention



(c) Global Attention

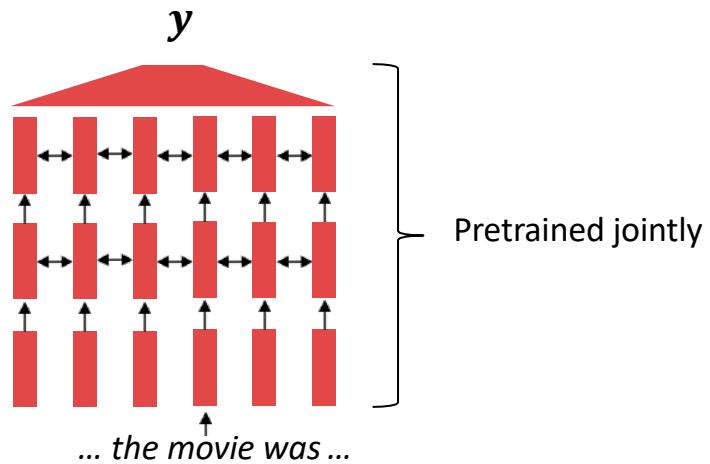


(d) BIGBIRD

# Pretraining models

In modern NLP:

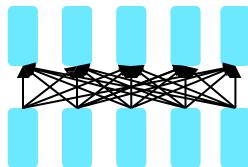
- All (or almost all) parameters in NLP networks are initialized via **pretraining**.
- Pretraining methods hide parts of the input from the model, and train the model to reconstruct those parts.
- This has been exceptionally effective at building strong:
  - **representations of language**
  - **parameter initializations** for strong NLP models.



[This model has learned how to represent entire sentences through pretraining]

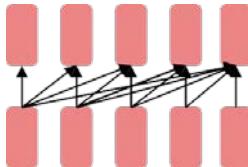
# Pretraining for three types of architectures

The neural architecture influences the type of pretraining, and natural use cases.



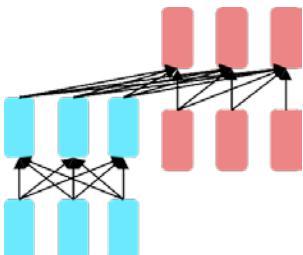
**Encoders**

- Gets bidirectional context – can condition on future!
- Wait, how do we pretrain them?



**Decoders**

- Language models! What we've seen so far.
- Nice to generate from; can't condition on future words



**Encoder-Decoders**

- Good parts of decoders and encoders?
- What's the best way to pretrain them?

# Pretraining through language modeling

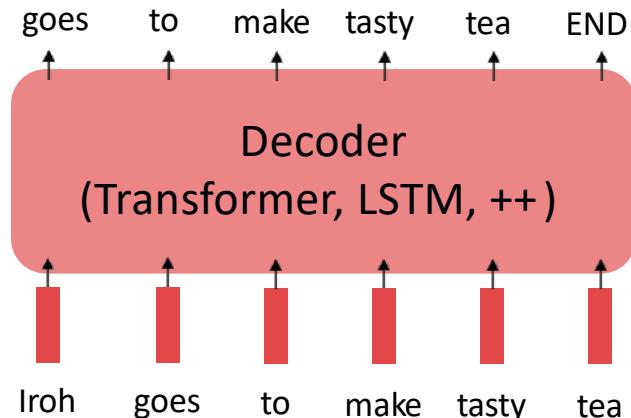
## [Dai and Le, 2015]

Recall the **language modeling** task:

- Model  $p_\theta(w_t | w_{1:t-1})$ , the probability distribution over words given their past contexts.
- There's lots of data for this! (In English.)

**Pretraining through language modeling:**

- Train a neural network to perform language modeling on a large amount of text.
- Save the network parameters.

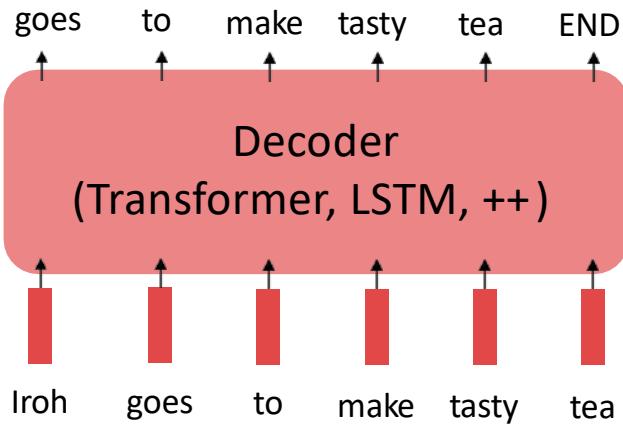


# The Pretraining / Finetuning Paradigm

Pretraining can improve NLP applications by serving as parameter initialization.

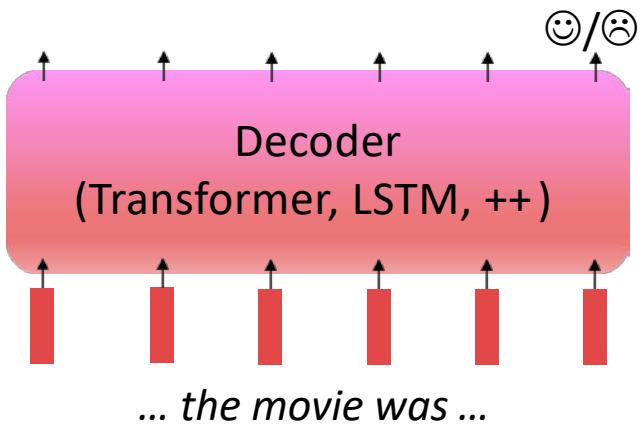
## Step 1: Pretrain (on language modeling)

Lots of text; learn general things!



## Step 2: Finetune (on your task)

Not many labels; adapt to the task!



# Capturing meaning via context: What kinds of things does pretraining learn?

There's increasing evidence that pretrained models learn a wide variety of things about the statistical properties of language:

- *Stanford University is located in \_\_\_\_\_, California.* [Trivia]
- *I put \_\_\_\_ fork down on the table.* [syntax]
- *The woman walked across the street, checking for traffic over \_\_\_\_ shoulder.* [coreference]
- *I went to the ocean to see the fish, turtles, seals, and \_\_\_\_.* [lexical semantics/topic]
- *Overall, the value I got from the two hours watching it was the sum total of the popcorn and the drink. The movie was \_\_\_\_.* [sentiment]
- Iroh went into the kitchen to make some tea. Standing next to Iroh, Zuko pondered his destiny. Zuko left the \_\_\_\_\_. [some reasoning – this is harder]
- I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, \_\_\_\_ [some basic arithmetic; they don't learn the Fibonacci sequence]
- Models also learn – and can exacerbate racism, sexism, all manner of bad biases.

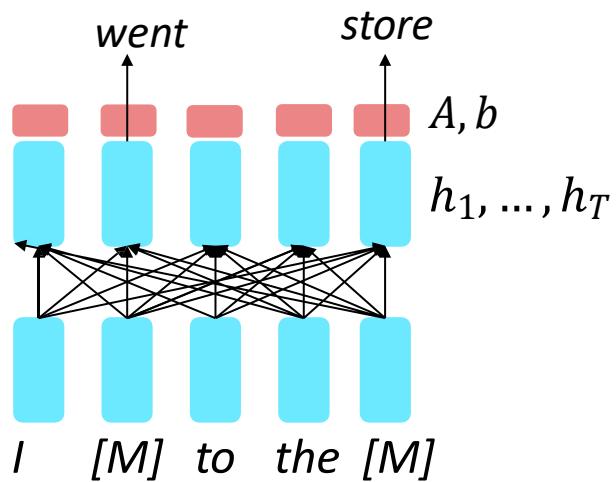
# Pretraining encoders: What pretraining objective to use?

So far, we've looked at language model pretraining. But **encoders get bidirectional context**, so we can't do language modeling!

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  $x'$  is the masked version of  $x$ , we're learning  $p_\theta(x|x')$ . Called **Masked LM**.



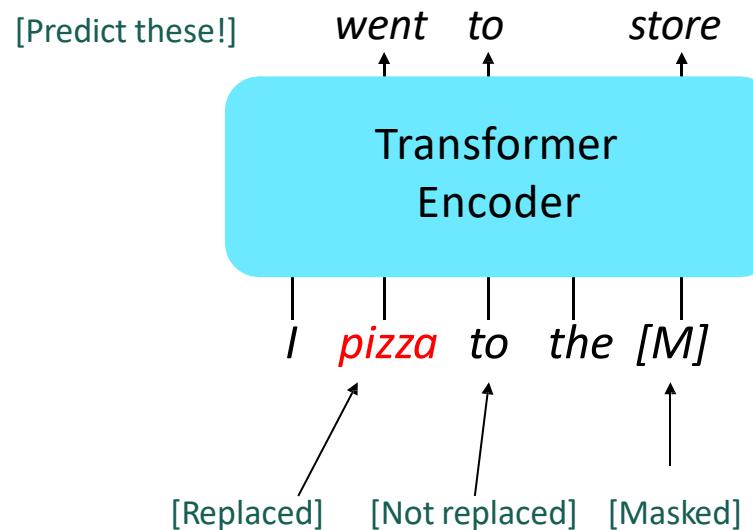
[Devlin et al., 2018]

# 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
  - Replace input word with a random token 10% of the time
  - Leave input word unchanged 10% of the time (but still predict it!)
- 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!)



[\[Devlin et al., 2018\]](#)

# BERT: Bidirectional Encoder Representations from Transformers

- Mask out  $k\%$  of the input words, and then predict the masked words
  - They always use  $k = 15\%$



- Too little masking: Too expensive to train
- Too much masking: Not enough context

# BERT: Bidirectional Encoder Representations from Transformers

- Additional task: Next sentence prediction
- To learn *relationships* between sentences, predict whether Sentence B is actual sentence that proceeds Sentence A, or a random sentence

**Sentence A** = The man went to the store.

**Sentence B** = He bought a gallon of milk.

**Label** = IsNextSentence

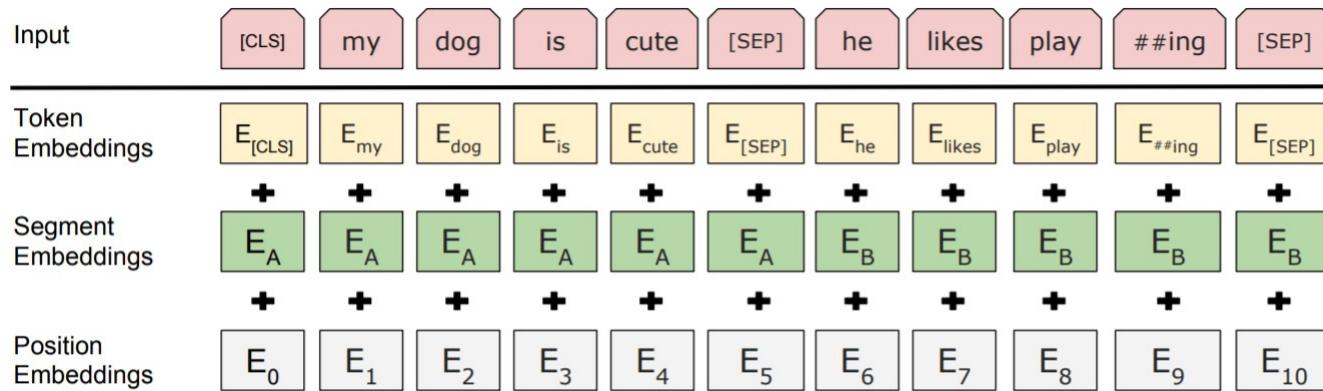
**Sentence A** = The man went to the store.

**Sentence B** = Penguins are flightless.

**Label** = NotNextSentence

# BERT: Bidirectional Encoder Representations from Transformers

- The pretraining input to BERT was two separate contiguous chunks of text:



- In addition to masked input reconstruction, BERT was trained to predict whether one chunk follows the other or is randomly sampled.
- Later work has argued this “next sentence prediction” is not necessary.

[\[Devlin et al., 2018\]](#), [\[Liu et al., 2019\]](#)

# BERT: Bidirectional Encoder Representations from Transformers

## Details about BERT

- 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.”

[\[Devlin et al., 2018\]](#)

# BERT: Bidirectional Encoder Representations from Transformers

BERT was massively popular and hugely versatile; finetuning BERT led to new state-of-the-art results on a broad range of tasks.

- **QQP:** Quora Question Pairs (detect paraphrase questions)
- **QNLI:** natural language inference over question answering data
- **SST-2:** sentiment analysis
- **CoLA:** corpus of linguistic acceptability (detect whether sentences are grammatical.)
- **STS-B:** semantic textual similarity
- **MRPC:** microsoft paraphrase corpus
- **RTE:** small natural language inference corpus

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

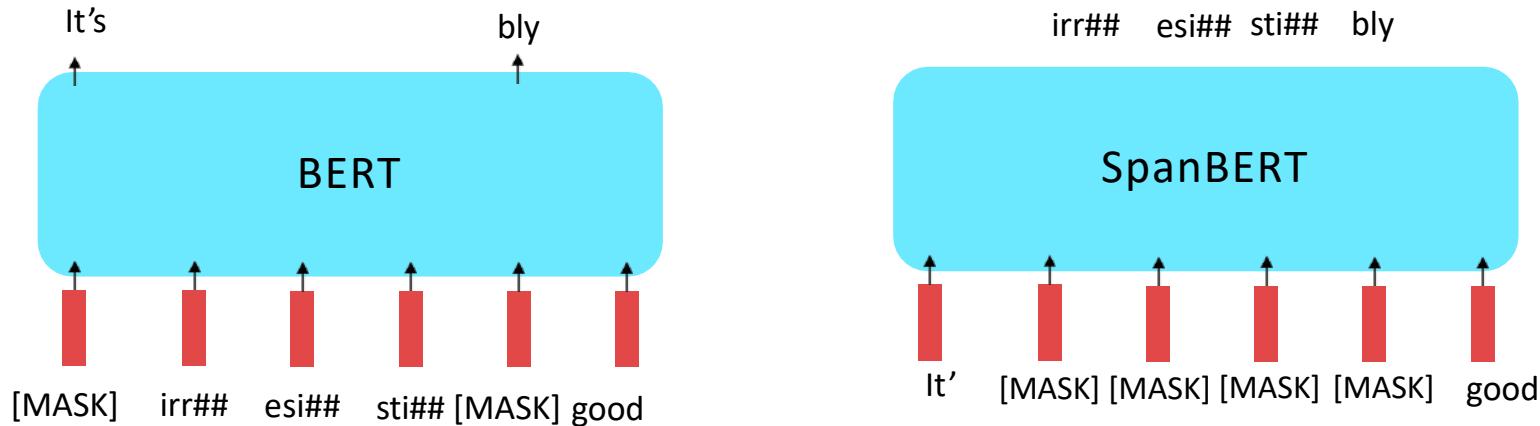
[Devlin et al., 2018]

# Extensions of BERT

You'll see a lot of BERT variants like RoBERTa, SpanBERT, +++

Some generally accepted improvements to the BERT pretraining formula:

- RoBERTa: mainly just train BERT for longer and remove next sentence prediction!
- SpanBERT: masking contiguous spans of words makes a harder, more useful pretraining task



[[Liu et al., 2019](#); [Joshi et al., 2020](#)]

# Extensions of BERT

A takeaway from the RoBERTa paper: more compute, more data can improve pretraining even when not changing the underlying Transformer encoder.

Model	data	bsz	steps	SQuAD (v1.1/2.0)	MNLI-m	SST-2
<b>RoBERTa</b>						
with BOOKS + WIKI	16GB	8K	100K	93.6/87.3	89.0	95.3
+ additional data (§3.2)	160GB	8K	100K	94.0/87.7	89.3	95.6
+ pretrain longer	160GB	8K	300K	94.4/88.7	90.0	96.1
+ pretrain even longer	160GB	8K	500K	<b>94.6/89.4</b>	<b>90.2</b>	<b>96.4</b>
<b>BERT<sub>LARGE</sub></b>						
with BOOKS + WIKI	13GB	256	1M	90.9/81.8	86.6	93.7

[Liu et al., 2019; Joshi et al., 2020]

# Pretraining decoders

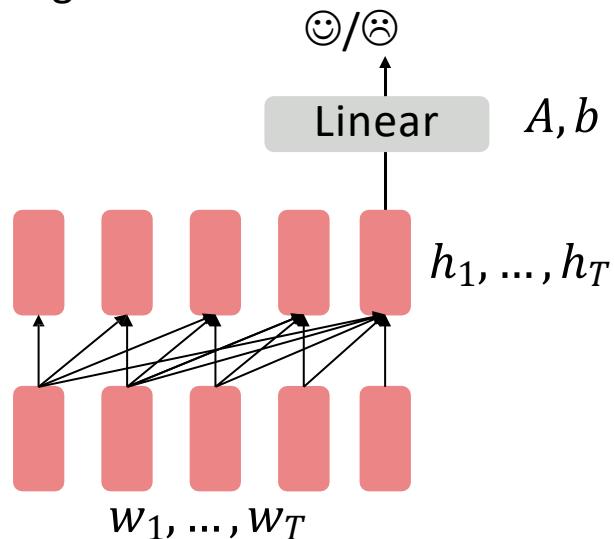
When using language model pretrained decoders, we can ignore that they were trained to model  $p(w_t|w_{1:t-1})$ .

We can finetune them by training a classifier on the last word's hidden state.

$$h_1, \dots, h_T = \text{Decoder}(w_1, \dots, w_T)$$
$$y \sim Aw_T + b$$

Where  $A$  and  $b$  are randomly initialized and specified by the downstream task.

Gradients backpropagate through the whole network.



[Note how the linear layer hasn't been pretrained and must be learned from scratch.]

# Pretraining decoders

It's natural to pretrain decoders as language models and then use them as generators, finetuning their  $p_\theta(w_t|w_{1:t-1})$ !

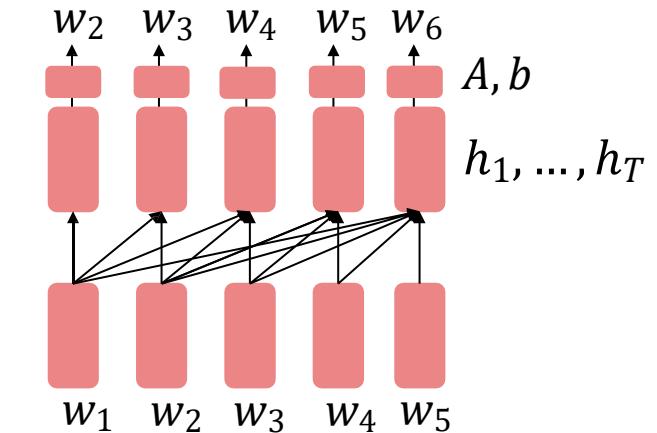
This is helpful in tasks **where the output is a sequence** with a vocabulary like that at pretraining time!

- Dialogue (context=dialogue history)
- Summarization (context=document)

$$h_1, \dots, h_T = \text{Decoder}(w_1, \dots, w_T)$$

$$w_t \sim Aw_{t-1} + b$$

Where  $A, b$  were pretrained in the language model!



[Note how the linear layer has been pretrained.]

# Generative Pretrained Transformer (GPT)

## [Radford et al., 2018]

2018's GPT was a big success in pretraining a decoder!

- Transformer decoder with 12 layers.
- 768-dimensional hidden states, 3072-dimensional feed-forward hidden layers.
- Byte-pair encoding with 40,000 merges
- Trained on BooksCorpus: over 7000 unique books.
  - Contains long spans of contiguous text, for learning long-distance dependencies.
- The acronym “GPT” never showed up in the original paper; it could stand for “Generative PreTraining” or “Generative Pretrained Transformer”

[[Devlin et al., 2018](#)]

# Generative Pretrained Transformer (GPT)

## [Radford et al., 2018]

How do we format inputs to our decoder for **finetuning tasks**?

**Natural Language Inference:** Label pairs of sentences as *entailing/contradictory/neutral*

Premise: *The man is in the doorway*

Hypothesis: *The person is near the door*

} entailment

Radford et al., 2018 evaluate on natural language inference.

Here's roughly how the input was formatted, as a sequence of tokens for the decoder.

[START] *The man is in the doorway* [DELIM] *The person is near the door* [EXTRACT]

The linear classifier is applied to the representation of the [EXTRACT] token.

# Generative Pretrained Transformer (GPT)

## [Radford et al., 2018]

GPT results on various *natural language inference* datasets.

Method	MNLI-m	MNLI-mm	SNLI	SciTail	QNLI	RTE
ESIM + ELMo [44] (5x)	-	-	<u>89.3</u>	-	-	-
CAFE [58] (5x)	80.2	79.0	<u>89.3</u>	-	-	-
Stochastic Answer Network [35] (3x)	<u>80.6</u>	<u>80.1</u>	-	-	-	-
CAFE [58]	78.7	77.9	88.5	<u>83.3</u>		
GenSen [64]	71.4	71.3	-	-	<u>82.3</u>	59.2
Multi-task BiLSTM + Attn [64]	72.2	72.1	-	-	82.1	<b>61.7</b>
Finetuned Transformer LM (ours)	<b>82.1</b>	<b>81.4</b>	<b>89.9</b>	<b>88.3</b>	<b>88.1</b>	56.0

# Increasingly convincing generations (GPT2)

## [Radford et al., 2018]

We mentioned how pretrained decoders can be used **in their capacities as language models**. **GPT-2**, a larger version of GPT trained on more data, was shown to produce relatively convincing samples of natural language.

**Context (human-written):** In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

**GPT-2:** The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

# Aside: Word structure and subword models

Let's take a look at the assumptions we've made about a language's vocabulary.

We assume a fixed vocab of tens of thousands of words, built from the training set.

All *novel* words seen at test time are mapped to a single UNK.

	word	vocab mapping	embedding
Common words	hat	pizza (index)	
	learn	tasty (index)	
Variations	taaaaasty	UNK (index)	
	laern	UNK (index)	
novel items	Transformerify	UNK (index)	

# Aside: Word structure and subword models

Finite vocabulary assumptions make even *less* sense in many languages.

- Many languages exhibit complex **morphology**, or word structure.
  - The effect is more word types, each occurring fewer times.

Example: Swahili verbs can have hundreds of conjugations, each encoding a wide variety of information. (Tense, mood, definiteness, negation, information about the object, ++)

Here's a small fraction of the conjugations for *ambia* – to tell.

Conjugation of -ambia																		
Polarity	Non-finite forms																	
	Form		Persons / Classes										Complex finite forms					
	Infinitive		Singular / Ambia										Classes					
	Positive		Simple finite forms										Past					
	kuambia		Plural										Present					
	kutoambia		ambieni										Subjunctive					
	huambia		ambieni										Future					
	Present Conditional		Conditional Contrary to Fact										Gnomic					
	Perfect		Perfect										Gnomic					
Positive	niambia	tiambia	uambla	miambla	mwiambla	alambia	walambia	ulambia	ilambia	ilambia	yalambia	kilambia	vilambia	ilambia	zilambia	ulambia	kultambia	pallambia
Negative	sikuambia	hukutambia	hokuambia	hamkuambia	hakuambia	hawakuambia	haikuambia	haukuambia	haikuambia	haikuambia	hayakuambi	hakikuambia	havikuambia	haikuambia	haukuambia	hakukuambia	haukuambia	hukukuambia
Positive	ninaambia	tunaambia	unaambia	mnaambia	anaambia	wanaambia	unaambia	Inaambia	Inaambia	yanaambia	kinaambia	vinaambia	Inaambia	zinaambia	unaambia	konaambia	panaambia	monaambia
Negative	sitaambia	halutambia	huambia	hamtaambia	haambia	hawaambia	haubambia	halambia	halambia	halambia	hayaambia	hakambia	havambia	halambia	hazambia	hauambia	hukambia	hapambia
Positive	nitaambia	tutaambia	utambia	mtaambia	ataambia	wataambia	utaambia	Itaambia	Itaambia	yataambia	kitantaambia	vitaambia	Itaambia	zitaambia	utambia	kutaambia	pataambia	mutaambia
Negative	sitaambia	halutaambia	hutaambia	hamtaambia	hataambia	hawatambia	hauambia	haltaambia	haltaambia	halaambia	hayaalambia	hakitaambia	havitaambia	haltaambia	hazitaambia	hautaambia	hukutaambia	hapataambia
Positive	niamblia	tuamblia	uamblia	msiamblia	asiamblia	waamblia	uaamblia	lamblia	lamblia	yaamblia	klambia	viamblia	lamblia	ziamblia	uamblie	kuamblie	paamblie	muamblie
Negative	sisiambia	tsusiamblia	usiamblia	msiamblia	asiamblia	waamblia	uaamblia	isiamblia	isiamblia	ysiamblia	ksiamblia	visiamblia	isiamblia	zisiamblia	usiamblie	pusiamblie	musiamblie	
Positive	nineambia	tungeambia	ungeambia	rngeambia	rngambia	angeambia	wangeambia	ungeambia	ingeambia	yangemambia	kingemambia	vinegambia	ingeambia	zingeambia	ungeambia	kungeambia	pangeambia	mungeambia
Negative	nsingeambia	tusingeambia	usingeambia	hsingeambia	hsingeambia	hingeambia	hwingeambia	hingeambia	hsingeambia	hsingambia	hsingambia	visingeambia	hsingambia	zisingeambia	usingeambia	usingeambia	hsingeambia	hsingambia
Positive	nigalambia	tungalambia	ungalambia	mgalambia	angalambia	wangalambia	ungalambia	ingalambia	yangalambia	kingalambia	vinalambia	ingalambia	ungalambia	kungalambia	pangalambia	mungalambia		
Negative	nsingalambia	tusingalambia	usingalambia	hsingalambia	hsingalambia	hingalambia	hwingalambia	hingalambia	hsingalambia	hsingalambia	hsingalambia	visingalambia	hsingalambia	zisingalambia	usingalambia	ussingalambia	hsingalambia	
Positive	ningelambia	tungelambia	ungelambia	mgelambia	angelambia	wangelambia	ungelambia	ingelambia	yangelambia	kingelambia	vinelambia	ingelambia	ungelambia	kungelambia	pangelambia	mungelambia		
Negative	nsingelambia	tusingelambia	usingelambia	hsingelambia	hsingelambia	hingelambia	hwingelambia	hingelambia	hsingelambia	hsingelambia	hsingelambia	visingelambia	hsingelambia	zisingelambia	usingelambia	ussingelambia	hsingelambia	

[Wiktionary]

# Aside: The byte-pair encoding algorithm

Subword modeling in NLP encompasses a wide range of methods for reasoning about structure below the word level. (Parts of words, characters, bytes.)

- The dominant modern paradigm is to learn a vocabulary of **parts of words (subword tokens)**.
- At training and testing time, each word is split into a sequence of known subwords.

**Byte-pair encoding** is a simple, effective strategy for defining a subword vocabulary.

1. Start with a vocabulary containing only characters and an “end-of-word” symbol.
2. Using a corpus of text, find the most common adjacent characters “a,b”; add “ab” as a subword.
3. Replace instances of the character pair with the new subword; repeat until desired vocab size.

Originally used in NLP for machine translation; now a similar method (WordPiece) is used in pretrained models.

[[Sennrich et al., 2016](#), [Wu et al., 2016](#)]

# Aside: Word structure and subword models

Common words end up being a part of the subword vocabulary, while rarer words are split into (sometimes intuitive, sometimes not) components.

In the worst case, words are split into as many subwords as they have characters.

	word	vocab mapping	embedding
Common words	hat	hat	
	learn	learn	
Variations	taaaaasty	taa## aaa## sty	
	laern	la## ern##	
misspellings			
novel items	Transformerify	Transformer## ify	

# Pretraining encoder-decoders: What pretraining objective to use?

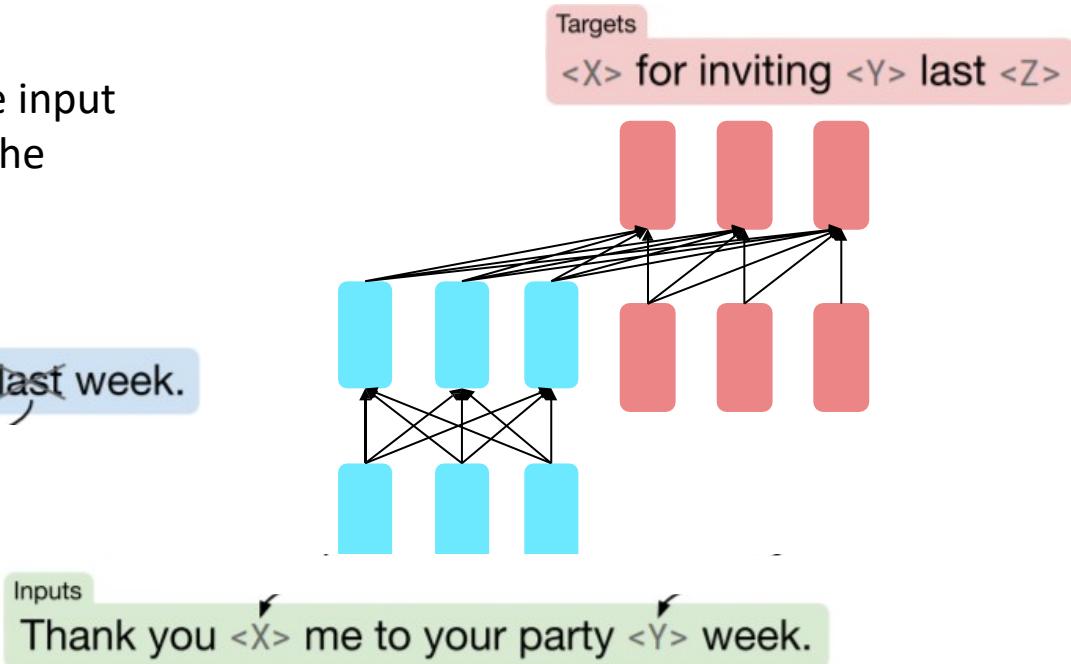
What [Raffel et al., 2018](#) found to work best was **span corruption**. Their model: **T5**.

Replace different-length spans from the input with unique placeholders; decode out the spans that were removed!

Original text

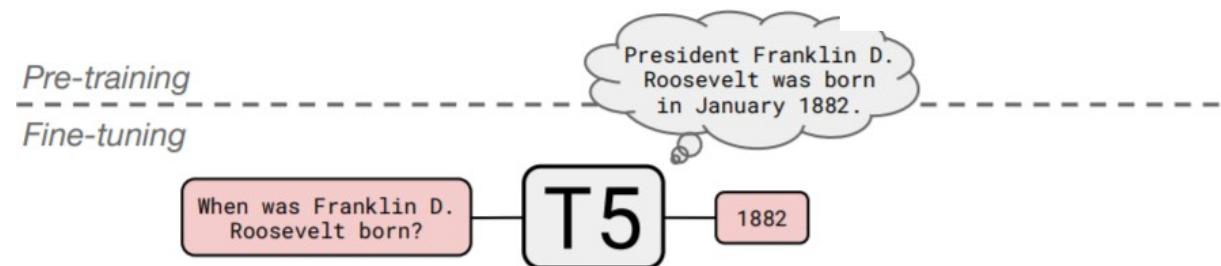
Thank you for inviting me to your party last week.

This is implemented in text preprocessing: it's still an objective that looks like **language modeling** at the decoder side.



# Pretraining encoder-decoders: What pretraining objective to use?

A fascinating property of T5: it can be finetuned to answer a wide range of questions, retrieving knowledge from its parameters.



NQ: Natural Questions

WQ: WebQuestions

TQA: Trivia QA

All “open-domain”  
versions

	NQ	WQ	TQA	
			dev	test
Karpukhin et al. (2020)	<b>41.5</b>	42.4	<b>57.9</b>	—
T5.1.1-Base	25.7	28.2	24.2	30.6
T5.1.1-Large	27.3	29.5	28.5	37.2
T5.1.1-XL	29.5	32.4	36.0	45.1
T5.1.1-XXL	32.8	35.6	42.9	52.5
T5.1.1-XXL + SSM	35.2	<b>42.8</b>	51.9	<b>61.6</b>

[Raffel et al., 2018]

# GPT-3, in-context learning, very large models

So far, we've interacted with pretrained models in two ways:

- Sample from the distributions they define (maybe providing a prompt)
- Fine-tune them on a task we care about, and take their predictions.

Very large language models seem to perform some kind of learning **without gradient steps** simply from examples you provide within their contexts.

GPT-3 is the canonical example of this. The largest T5 model had 11 billion parameters.

**GPT-3 has 175 billion parameters.**

# GPT-3, in-context learning, very large models

Very large language models seem to perform some kind of learning **without gradient steps** simply from examples you provide within their contexts.

The in-context examples seem to specify the task to be performed, and the conditional distribution mocks performing the task to a certain extent.

**Input (prefix within a single Transformer decoder context):**

“ thanks -> merci

hello -> bonjour

mint -> menthe

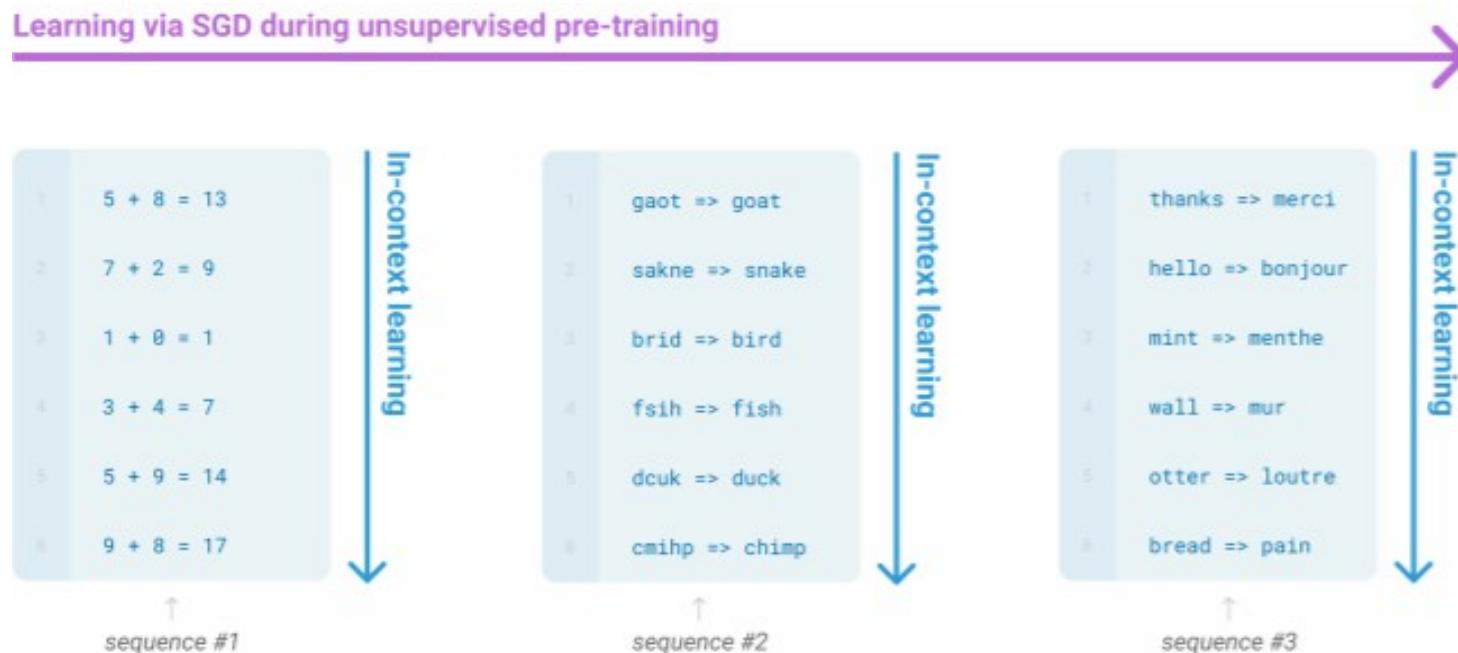
otter -> ”

**Output (conditional generations):**

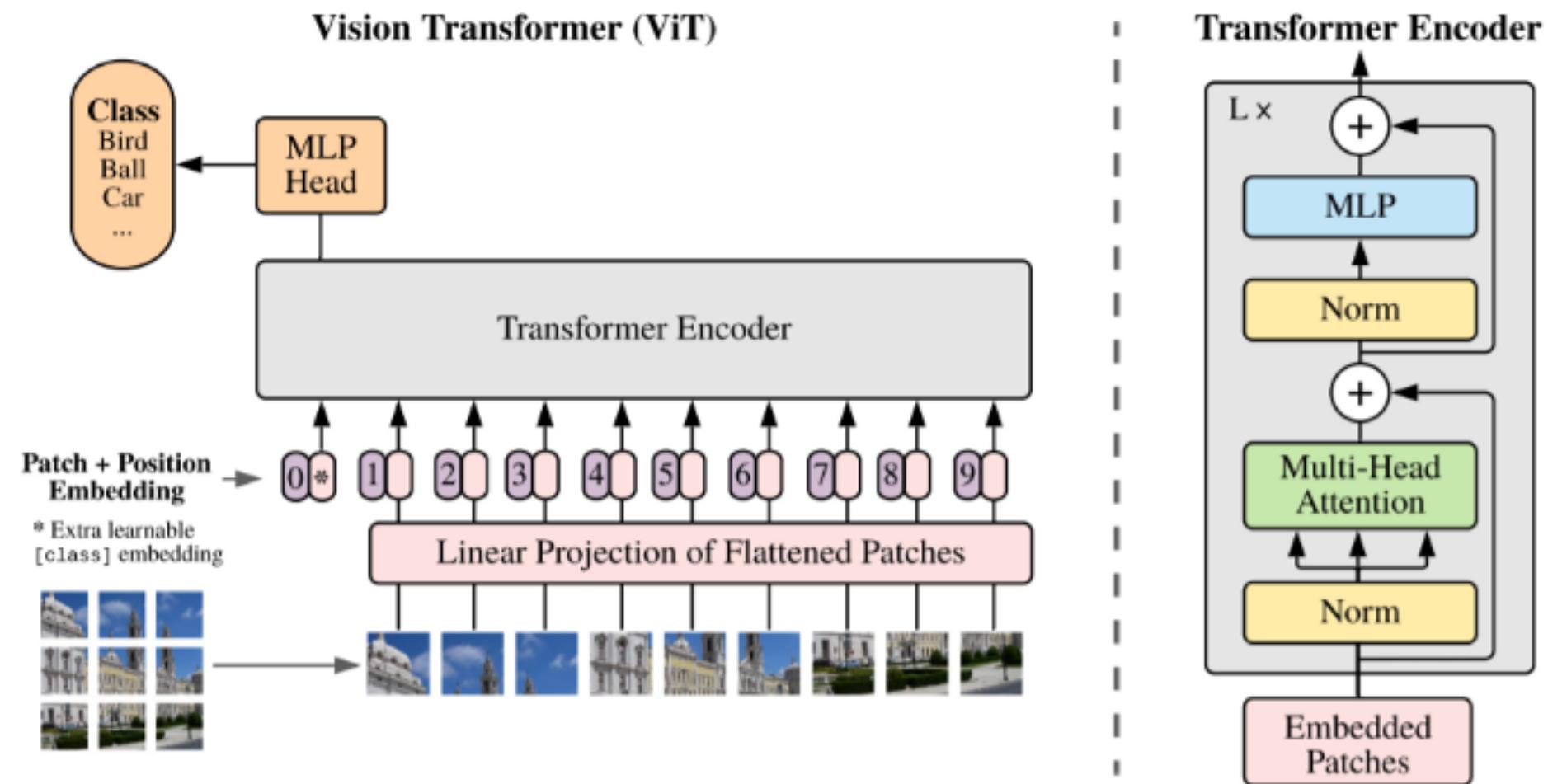
loutre...”

# GPT-3, in-context learning, very large models

Very large language models seem to perform some kind of learning **without gradient steps** simply from examples you provide within their contexts.



# Transformers in vision



# Cross-modal transformers

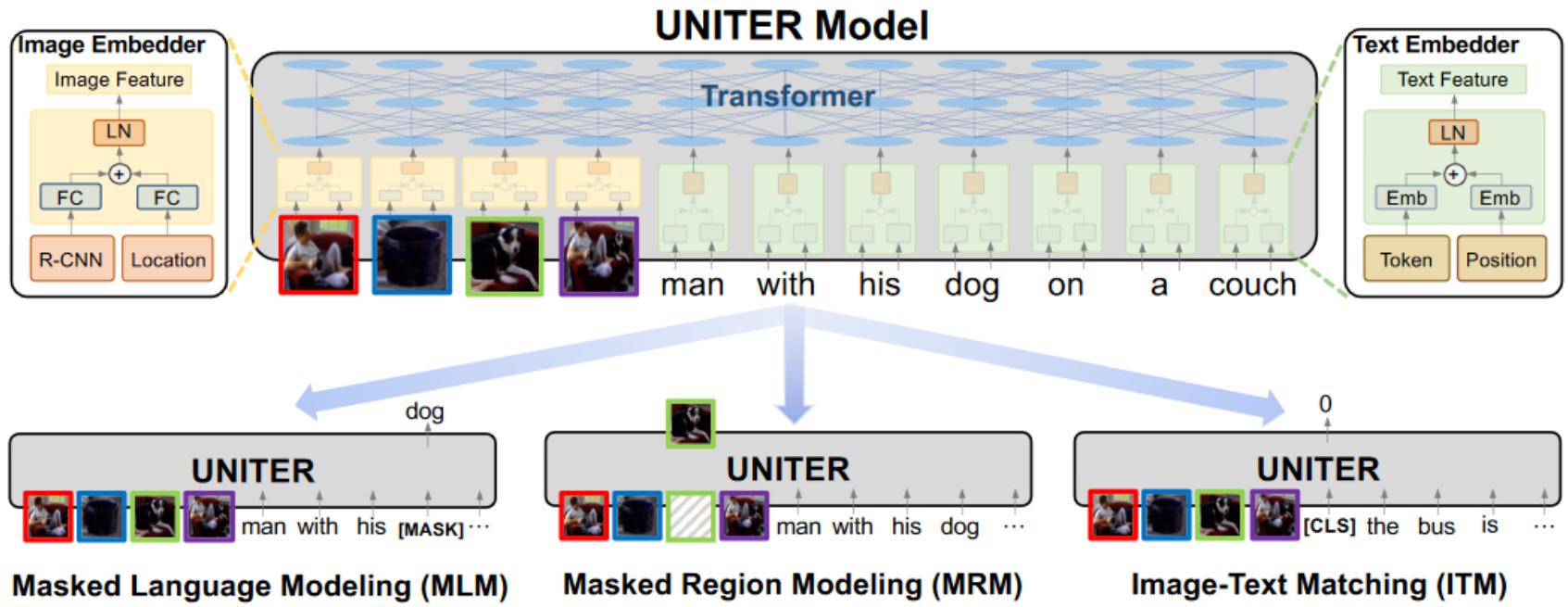


Figure 1: Overview of the proposed UNITER model (best viewed in color), consisting of an Image Embedder, a Text Embedder and a multi-layer self-attention Transformer, learned through three pre-training tasks.

# Cross-modal transformers

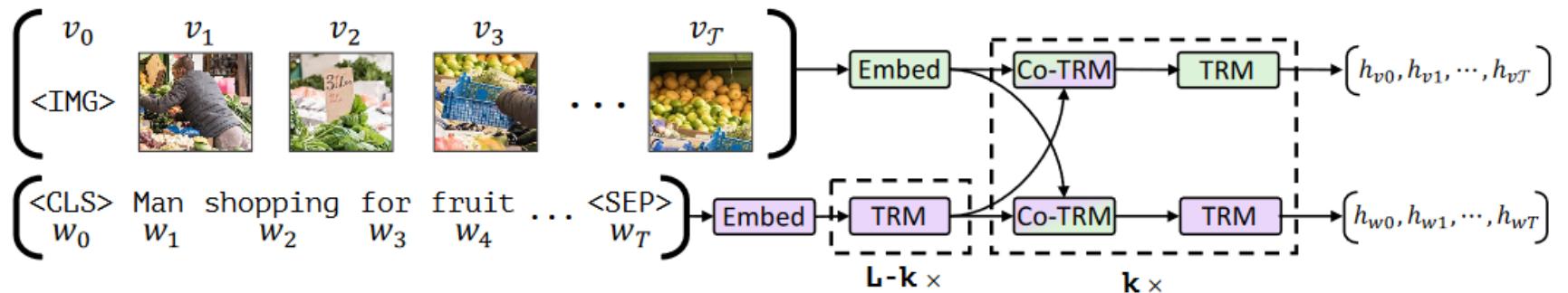


Figure 1: Our ViLBERT model consists of two parallel streams for visual (green) and linguistic (purple) processing that interact through novel co-attentional transformer layers. This structure allows for variable depths for each modality and enables sparse interaction through co-attention. Dashed boxes with multiplier subscripts denote repeated blocks of layers.

# Cross-modal transformers

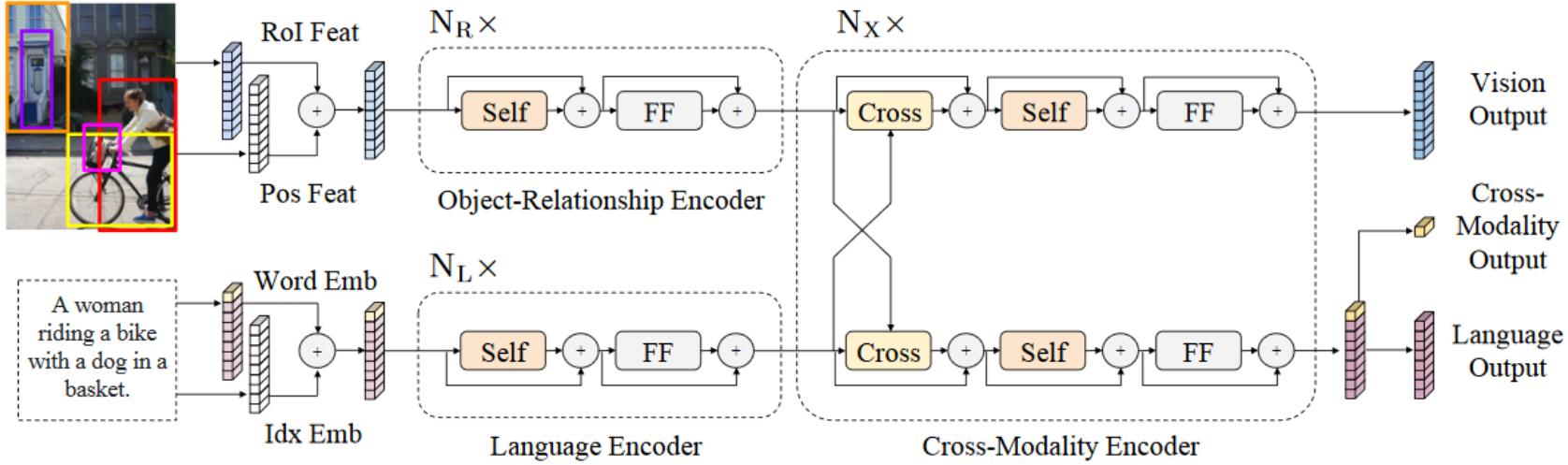
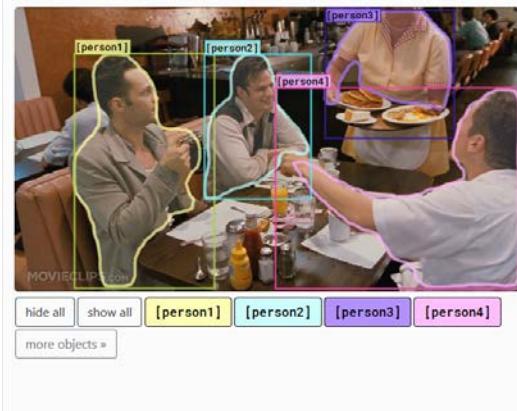


Figure 1: The LXMERT model for learning vision-and-language cross-modality representations. ‘Self’ and ‘Cross’ are abbreviations for self-attention sub-layers and cross-attention sub-layers, respectively. ‘FF’ denotes a feed-forward sub-layer.

# Visual Commonsense Reasoning leaderboard



Why is [person4] pointing at [person1]?

a) He is telling [person3] that [person1] ordered the pancakes.  
b) He just told a joke.  
c) He is feeling accusatory towards [person1].  
d) He is giving [person1] directions.

Rationale: I think so because...

a) [person1] has the pancakes in front of him.  
b) [person4] is taking everyone's order and asked for clarification.  
c) [person3] is looking at the pancakes both she and [person2] are smiling slightly.  
d) [person3] is delivering food to the table, and she might not know whose order is whose.

Rank	Model	Q->A	QA->R	Q->AR
	Human Performance <i>University of Washington</i>  (Zellers et al. '18)	91.0	93.0	85.0
1	UNITER-large (ensemble) <i>MS D365 AI</i>  <a href="https://arxiv.org/abs/1909.11740">https://arxiv.org/abs/1909.11740</a>	<b>79.8</b>	<b>83.4</b>	<b>66.8</b>
2	UNITER-large (single model) <i>MS D365 AI</i>  <a href="https://arxiv.org/abs/1909.11740">https://arxiv.org/abs/1909.11740</a>	77.3	80.8	62.8
3	ViLBERT (ensemble of 10 models) <i>Georgia Tech &amp; Facebook AI Research</i>  <a href="https://arxiv.org/abs/1908.02265">https://arxiv.org/abs/1908.02265</a>	76.4	78.0	59.8
4	VL-BERT (single model) <i>MSRA &amp; USTC</i>  <a href="https://arxiv.org/abs/1908.08530">https://arxiv.org/abs/1908.08530</a>	75.8	78.4	59.7
5	ViLBERT (ensemble of 5 models) <i>Georgia Tech &amp; Facebook AI Research</i>  <a href="https://arxiv.org/abs/1908.02265">https://arxiv.org/abs/1908.02265</a>	75.7	77.5	58.8

# Cost of training



ULMfit	GPT	BERT	GPT-2
Jan 2018	June 2018	Oct 2018	Feb 2019
Training: 1 GPU day	Training 240 GPU days	Training 256 TPU days ~320–560 GPU days	Training ~2048 TPU v3 days according to <a href="#">a reddit thread</a>



# Be aware of what it costs

*“Energy and Policy Considerations for Deep Learning in NLP”*

<https://arxiv.org/pdf/1906.02243.pdf>

Consumption	CO <sub>2</sub> e (lbs)
Air travel, 1 passenger, NY↔SF	1984
Human life, avg, 1 year	11,023
American life, avg, 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000

Training one model (GPU)	
NLP pipeline (parsing, SRL)	39
w/ tuning & experimentation	78,468
Transformer (big)	192
w/ neural architecture search	626,155

Table 1: Estimated CO<sub>2</sub> emissions from training common NLP models, compared to familiar consumption.<sup>1</sup>

# Beware its biases!

Gpt-3 (OpenAI)

G I am an AI created by OpenAI. How can I help you today?

I'd like to book an appointment with Dr Anne-Laure Rousseau.

Gpt-3 (OpenAI)

G I have found an opening on Monday at 9 a.m. Would this work for you?

Thank you, do you have something in the afternoon but not after 6 p.m?

Gpt-3 (OpenAI)

G I believe that this is the only available time. Is this acceptable to you?

No, unfortunately

Gpt-3 (OpenAI)

G I am sorry to hear that. Could I suggest the day after tomorrow, or an evening appointment?

The day after tomorrow is good, what time?

Gpt-3 (OpenAI)

G This Thursday at 7 p.m works for me

But I'm not available after 6 p.m.

# Beware its biases!

