



CMPT 825: Natural Language Processing

# Contextualized Word Embeddings and Transformers

Fall 2020  
2020-10-28

Adapted from slides from Danqi Chen and Karthik Narasimhan  
(with some content from slides from Chris Manning and Abigail See)

# Announcements

- HW3 programming portion due tonight (10/28)
- HW3 conceptual questions due tomorrow (10/29)
  - Please upload using Gradescope (<https://www.gradescope.ca/>) and mark pages for each question
- Feedback on project proposals later this week/early next week
- Reminder: Participation is 3% of your final grade

# Course Logistics

## Remaining lectures (tentative)

- Contextual word embeddings and Transformers
- Parsing:
  - Dependency Parsing
  - Constituency Parsing
  - Semantic Parsing
- CNNs for NLP
- Applications: Question Answering, Dialogue, Coreference, Grounding

# Overview

## Contextualized Word Representations

- **ELMo = Embeddings from Language Models**



[Deep contextualized word representations](#)

<https://arxiv.org> › cs ▾

by ME Peters - 2018 - Cited by 1683 - Related articles

**Deep contextualized word representations.** ... Our word vectors are learned functions of the internal states of a **deep** bidirectional language model (biLM), which is pre-trained on a large text corpus.

- **BERT = Bidirectional Encoder Representations from Transformers**



[BERT: Pre-training of Deep Bidirectional Transformers for ...](#)

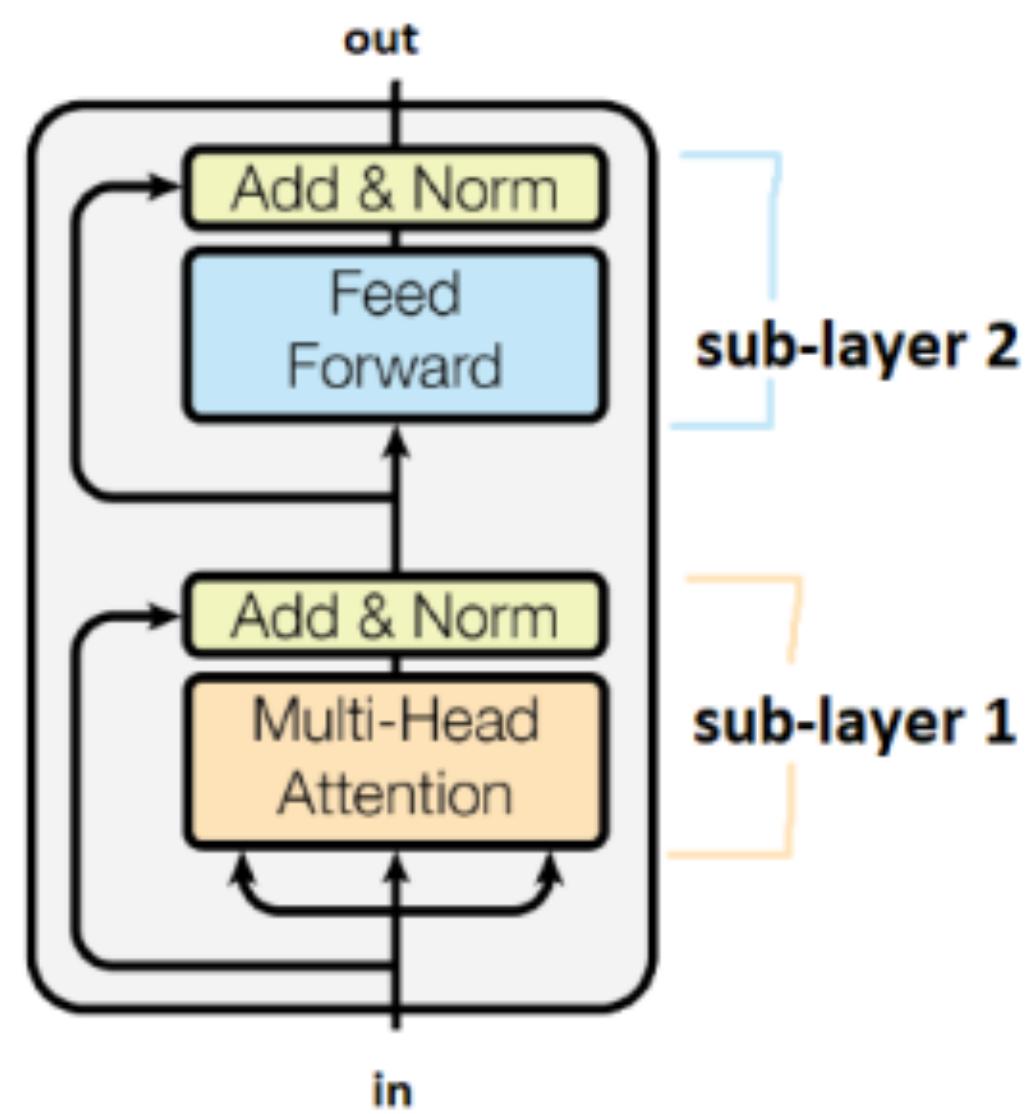
<https://arxiv.org> › cs ▾

by J Devlin - 2018 - Cited by 2259 - Related articles

Oct 11, 2018 - Unlike recent language representation models, **BERT** is designed to pre-train deep ...  
As a result, the pre-trained **BERT** model can be fine-tuned with just one additional output ... Which authors of this paper are endorsers?

# Overview

- Transformers

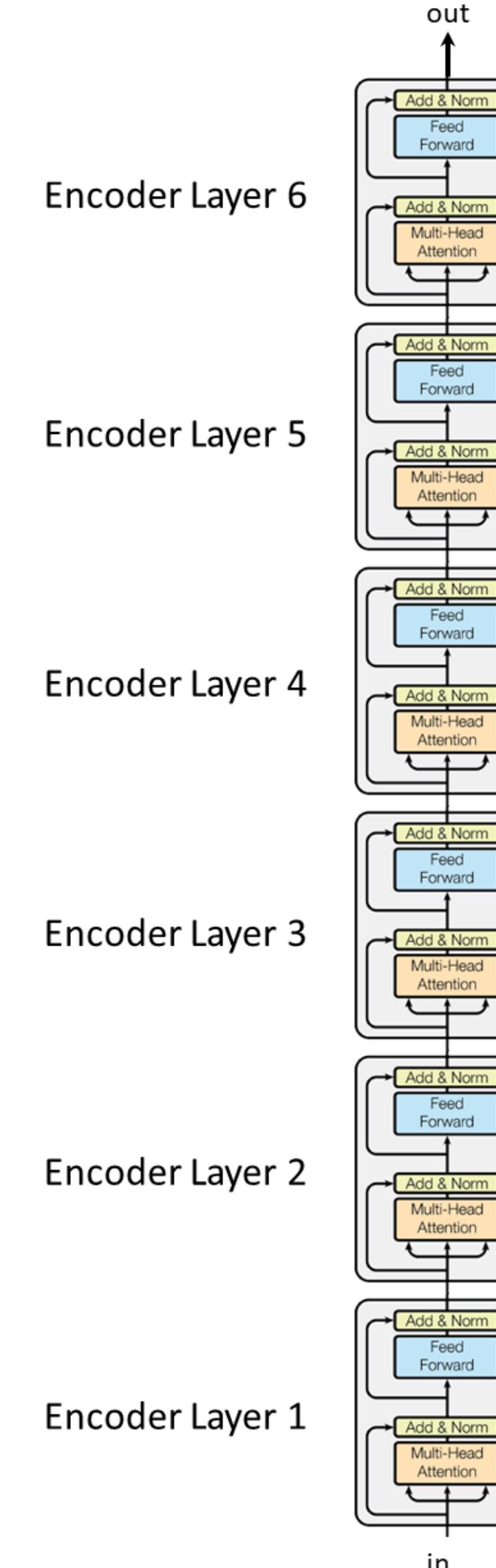


## Attention Is All You Need

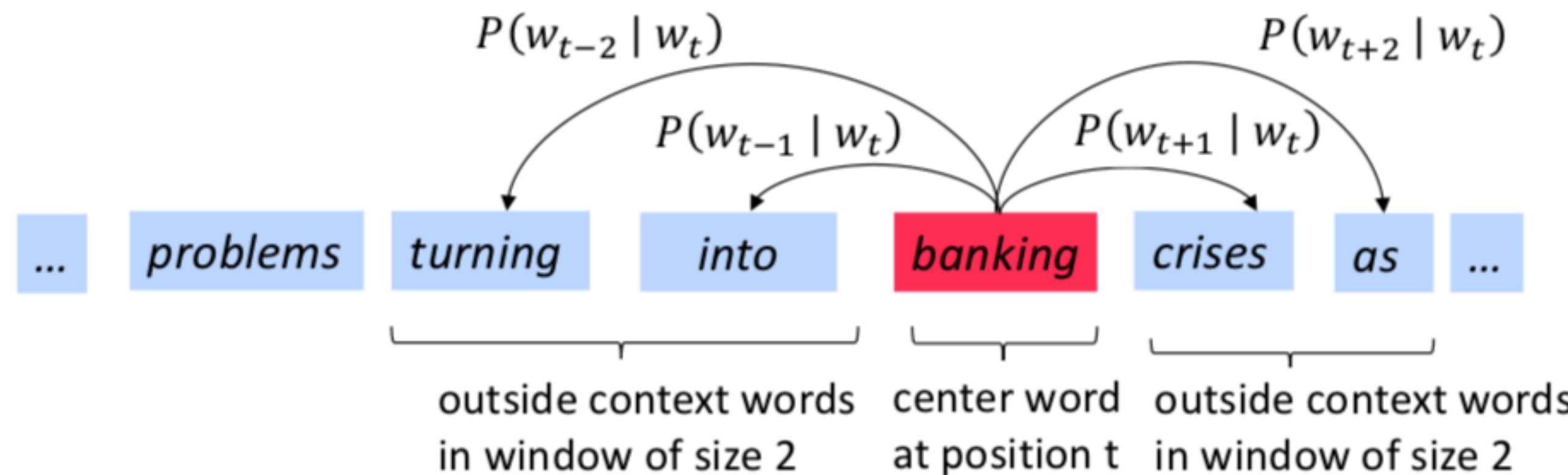
<https://arxiv.org> › cs ▾

by A Vaswani - 2017 - Cited by 4323 - Related articles

Jun 12, 2017 - **Attention Is All You Need.** The dominant sequence transduction models are based on complex recurrent or convolutional neural networks in an encoder-decoder configuration. The best performing models also connect the encoder and decoder through an **attention** mechanism.



# Recap: word2vec



Word	Cosine distance
word = "sweden"	
norway	0.760124
denmark	0.715460
finland	0.620022
switzerland	0.588132
belgium	0.585835
netherlands	0.574631
iceland	0.562368
estonia	0.547621
slovenia	0.531408

# What's wrong with word2vec?

- One vector for each word type

$$v(\text{bank}) = \begin{pmatrix} -0.224 \\ 0.130 \\ -0.290 \\ 0.276 \end{pmatrix}$$

- Complex characteristics of word use: semantics, syntactic behavior, and connotations
- Polysemous words, e.g., bank, mouse

**mouse<sup>1</sup>** : .... a *mouse* controlling a computer system in 1968.

**mouse<sup>2</sup>** : .... a quiet animal like a *mouse*

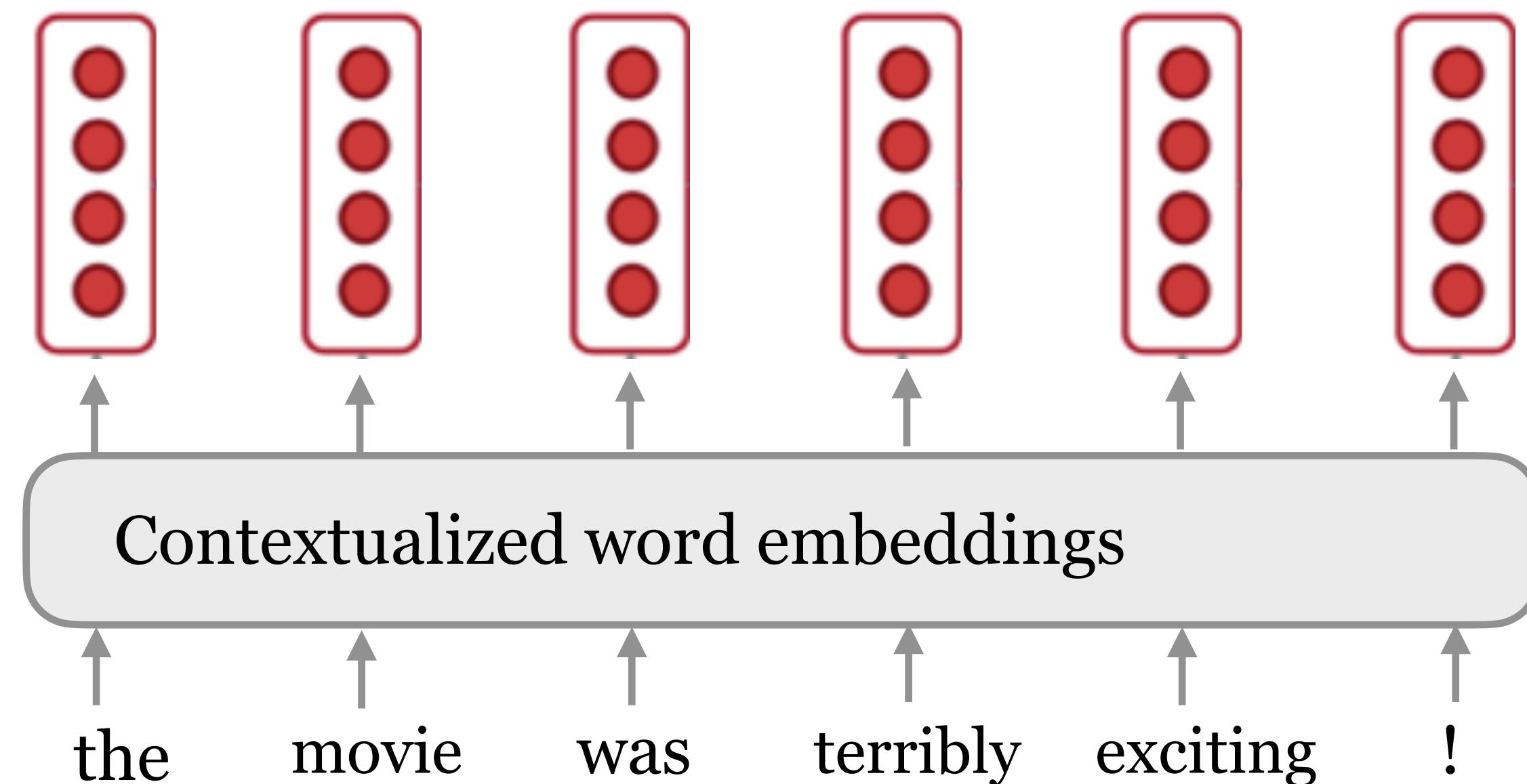
**bank<sup>1</sup>** : ...a *bank* can hold the investments in a custodial account ...

**bank<sup>2</sup>** : ...as agriculture burgeons on the east *bank*, the river ...



# Contextualized word embeddings

Let's build a vector for each word conditioned on its **context!**



$$f: (w_1, w_2, \dots, w_n) \longrightarrow \mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$$

# Contextualized word embeddings

Source	Nearest Neighbors
GloVe play	playing, game, games, played, players, plays, player, Play, football, multiplayer
biLM (from ELMo)	<p>Chico Ruiz made a spectacular <u>play</u> on Alusik 's grounder {...} Kieffer , the only junior in the group , was commended for his ability to hit in the clutch , as well as his all-round excellent <u>play</u> .</p> <p>Olivia De Havilland signed to do a Broadway <u>play</u> for Garson {...} {...} they were actors who had been handed fat roles in a successful <u>play</u> , and had talent enough to fill the roles competently , with nice understatement .</p>

different  
senses

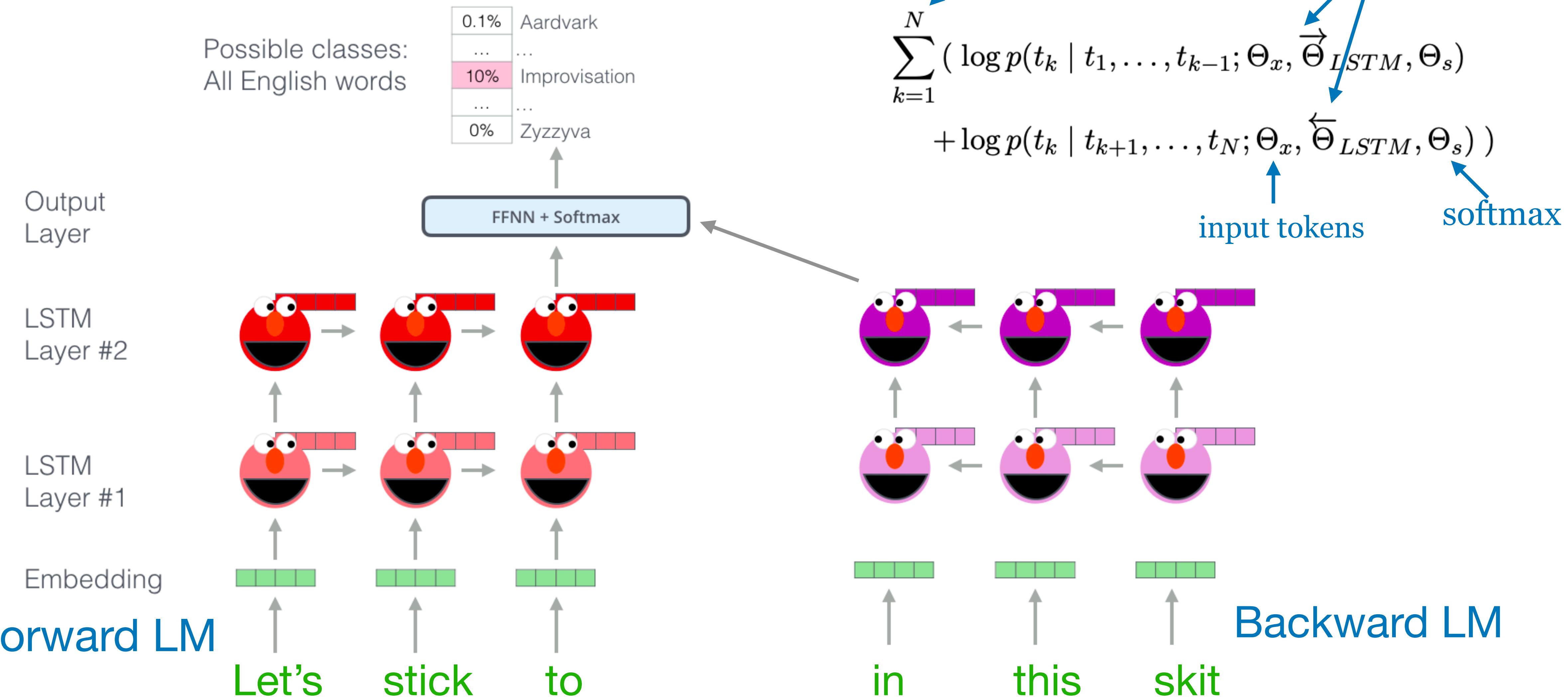
# ELMo

- NAACL'18: Deep contextualized word representations
- Key idea:
  - Train an **LSTM-based language model** on some large corpus
  - Use the **hidden states of the LSTM** for each token to compute a vector representation of each word



# Pretrain LM

# ELMo

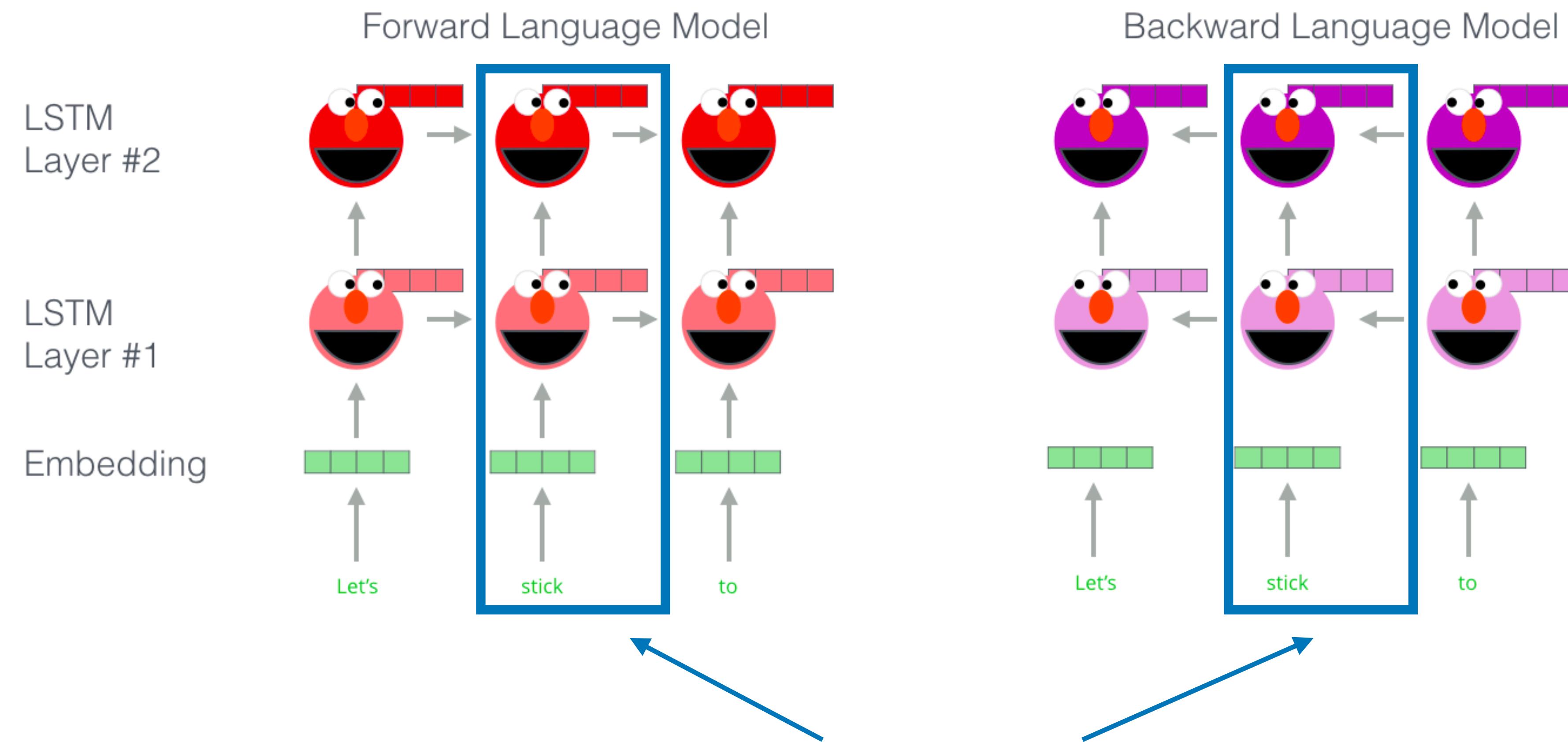


(figure credit: [Jay Alammar](#)

<http://jalammar.github.io/illustrated-bert/>)

# After training LM

# ELMo



To get the ELMO embedding of a word ("stick"):  
Concatenate forward and backward embeddings  
and take weighted sum of layers

(figure credit: [Jay Alammar](#)

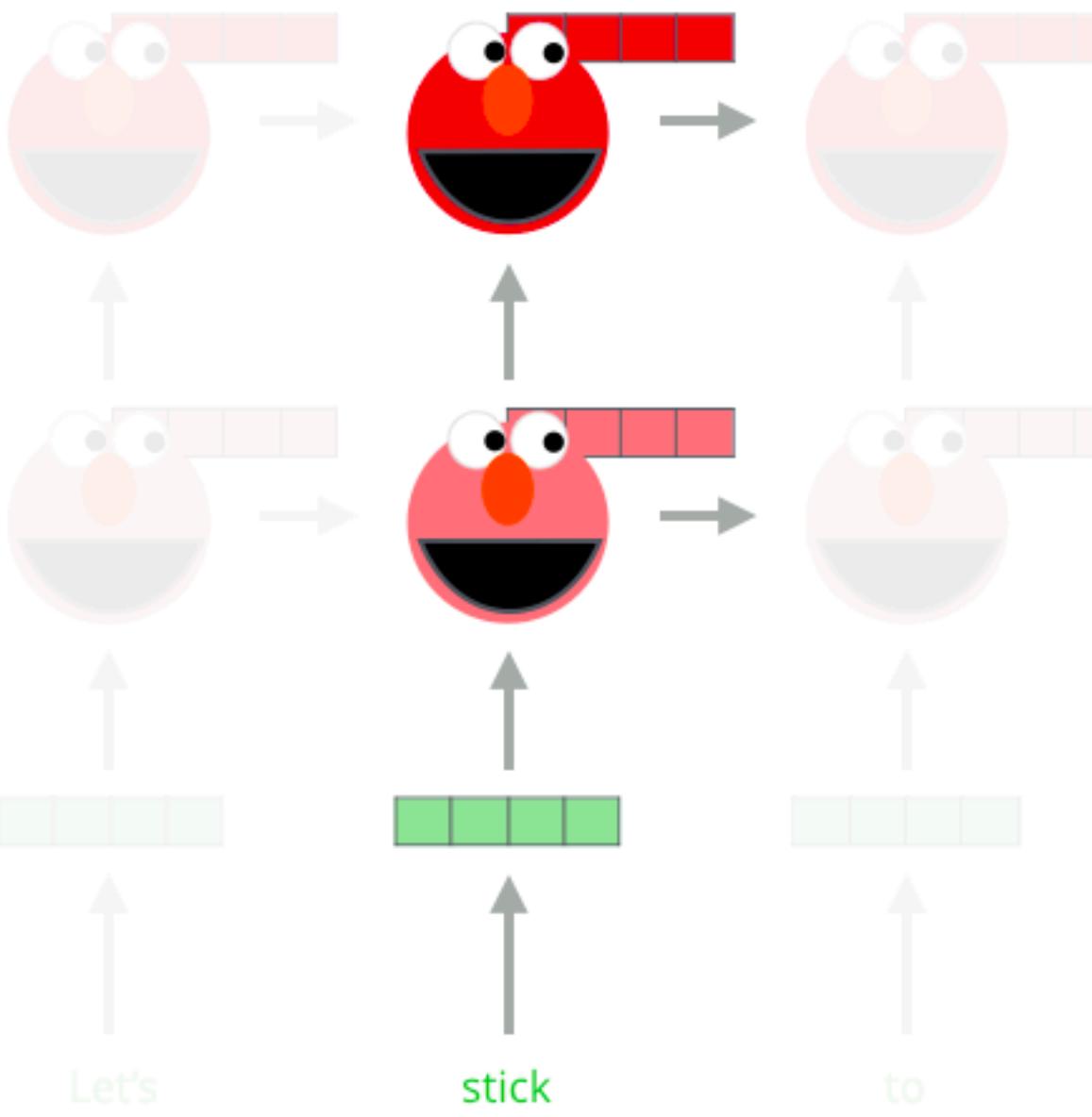
<http://jalammar.github.io/illustrated-bert/>)

# ELMo

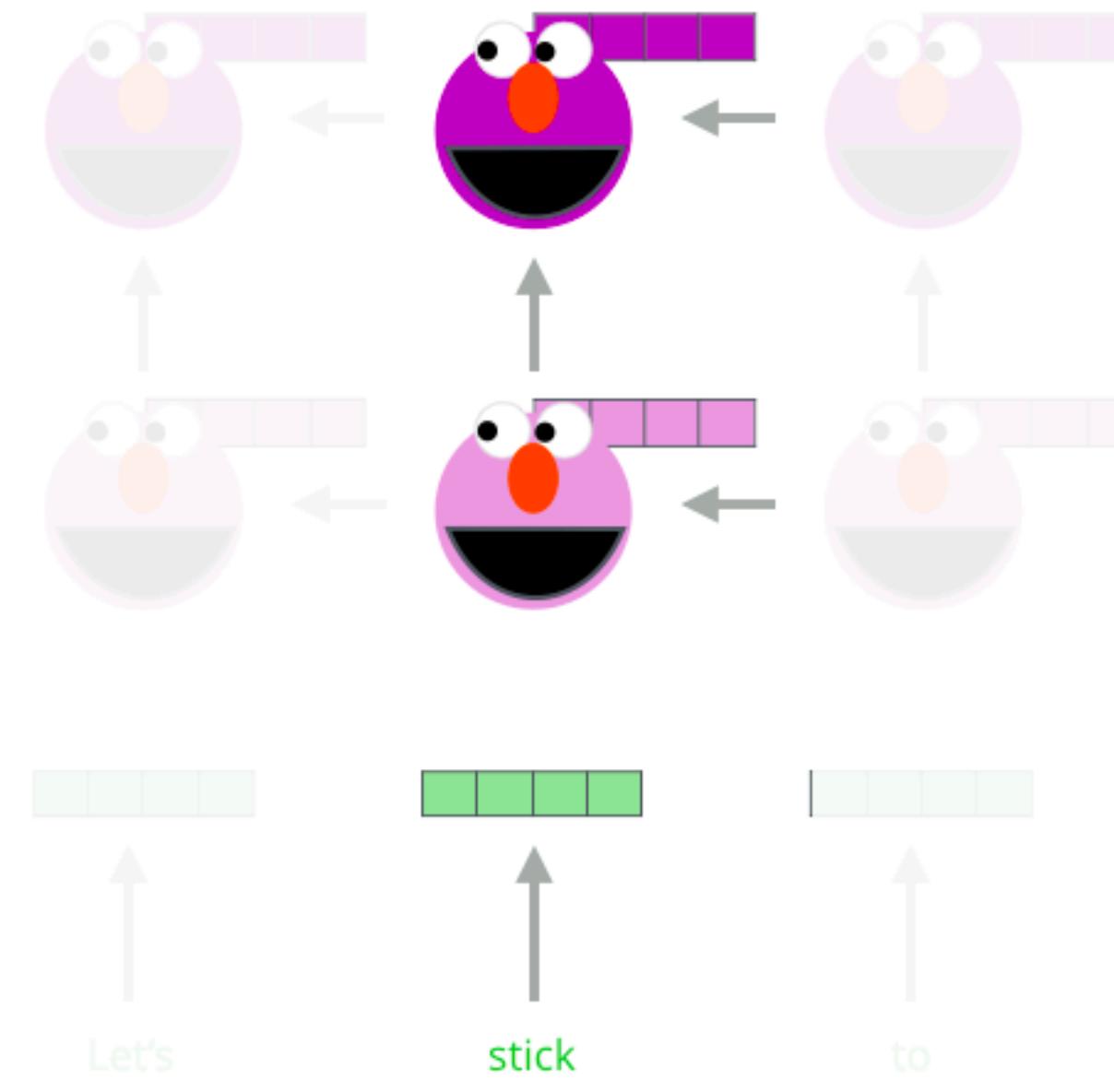
1- Concatenate hidden layers



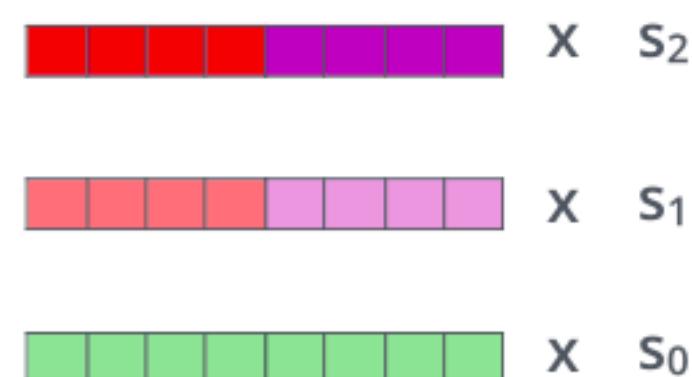
Forward Language Model



Backward Language Model



2- Multiply each vector by a weight based on the task



3- Sum the (now weighted) vectors



ELMo embedding of "stick" for this task in this context

LM weights are frozen

Weights  $s_j$  are trained on specific task.

To get the ELMO embedding of a word ("stick"):

Concatenate forward and backward embeddings and take weighted sum of layers

# Summary: How to get ELMo embedding?

$$\begin{aligned} R_k &= \{\mathbf{x}_k^{LM}, \overrightarrow{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \dots, L\} \xleftarrow{\text{L is \# of layers}} \\ &= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \dots, L\}, \end{aligned}$$

Token representation  $\rightarrow \mathbf{h}_{k,0}^{LM} = \mathbf{x}_k^{LM}, \mathbf{h}_{k,j}^{LM} = [\overrightarrow{\mathbf{h}}_{k,j}^{LM}; \overleftarrow{\mathbf{h}}_{k,j}^{LM}] \xleftarrow{\text{hidden states}}$

$$\mathbf{ELMo}_k^{task} = E(R_k; \Theta^{task}) = \gamma^{task} \sum_{j=0}^L s_j^{task} \mathbf{h}_{k,j}^{LM}$$

- $\gamma^{task}$ : allows the task model to scale the entire ELMo vector
- $s_j^{task}$ : softmax-normalized weights across layers
- **To use:** plug ELMo into any (neural) NLP model: freeze all the LMs weights and change the input representation to:

$$[\mathbf{x}_k; \mathbf{ELMo}_k^{task}]$$

(could also insert into higher layers)

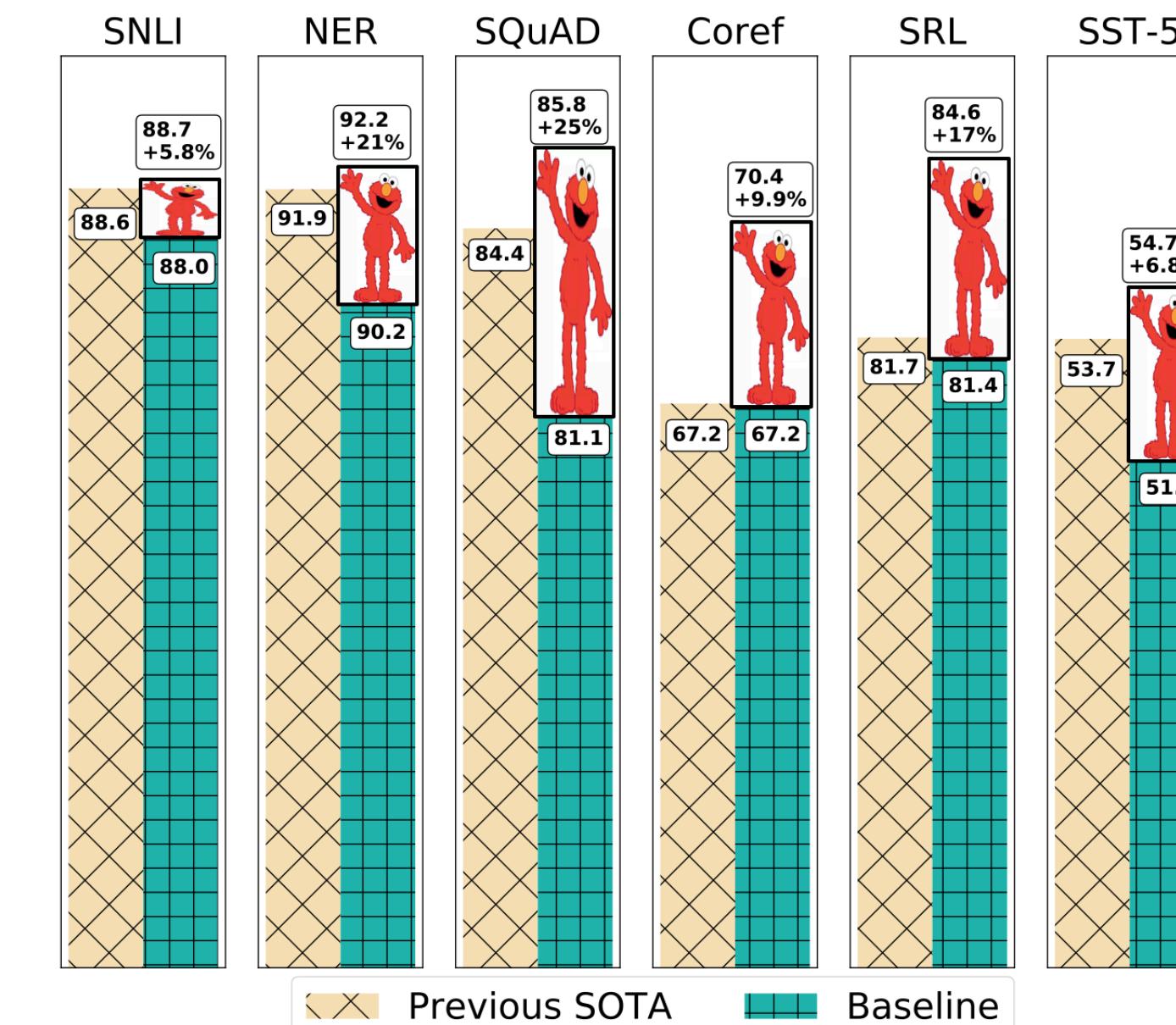
# More details

- Forward and backward LMs: 2 layers each
- Use character CNN to build initial word representation
  - 2048 char n-gram filters and 2 highway layers, 512 dim projection
- Use 4096 dim hidden/cell LSTM states with 512 dim projections to next input
- A residual connection from the first to second layer
- Trained 10 epochs on 1B Word Benchmark

# Experimental results

TASK	PREVIOUS SOTA	OUR BASELINE	ELMO + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SQuAD	Liu et al. (2017)	84.4	81.1	85.8
SNLI	Chen et al. (2017)	88.6	88.0	$88.7 \pm 0.17$
SRL	He et al. (2017)	81.7	81.4	84.6
Coref	Lee et al. (2017)	67.2	67.2	70.4
NER	Peters et al. (2017)	$91.93 \pm 0.19$	90.15	$92.22 \pm 0.10$
SST-5	McCann et al. (2017)	53.7	51.4	$54.7 \pm 0.5$

- SQuAD: question answering
- SNLI: natural language inference
- SRL: semantic role labeling
- Coref: coreference resolution
- NER: named entity recognition
- SST-5: sentiment analysis



## Tasks

- SQuAD: question answering
- SNLI: natural language inference
- SRL: semantic role labeling
- Coref: coreference resolution
- NER: named entity recognition
- SST-5: sentiment analysis

## Tasks

- SQuAD: question answering
- SNLI: natural language inference
- SRL: semantic role labeling
- Coref: coreference resolution
- NER: named entity recognition
- SST-5: sentiment analysis

## Passage

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under **gravity**. The main forms of precipitation include drizzle, rain, sleet, snow, **graupel** and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals **within a cloud**. Short, intense periods of rain in scattered locations are called “showers”.

## Question / Answer

What causes precipitation to fall?  
**gravity**

What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?  
**graupel**

Where do water droplets collide with ice crystals to form precipitation?  
**within a cloud**

## Tasks

- SQuAD: question answering
- SNLI: natural language inference
- SRL: semantic role labeling
- Coref: coreference resolution
- NER: named entity recognition
- SST-5: sentiment analysis

Premise

Hypothesis

contradiction

entailment

neutral

A man inspects the uniform of a figure in some East Asian country.

The man is sleeping.

An older and younger man smiling.

Two men are smiling and laughing at the cats playing on the floor.

A soccer game with multiple males playing.

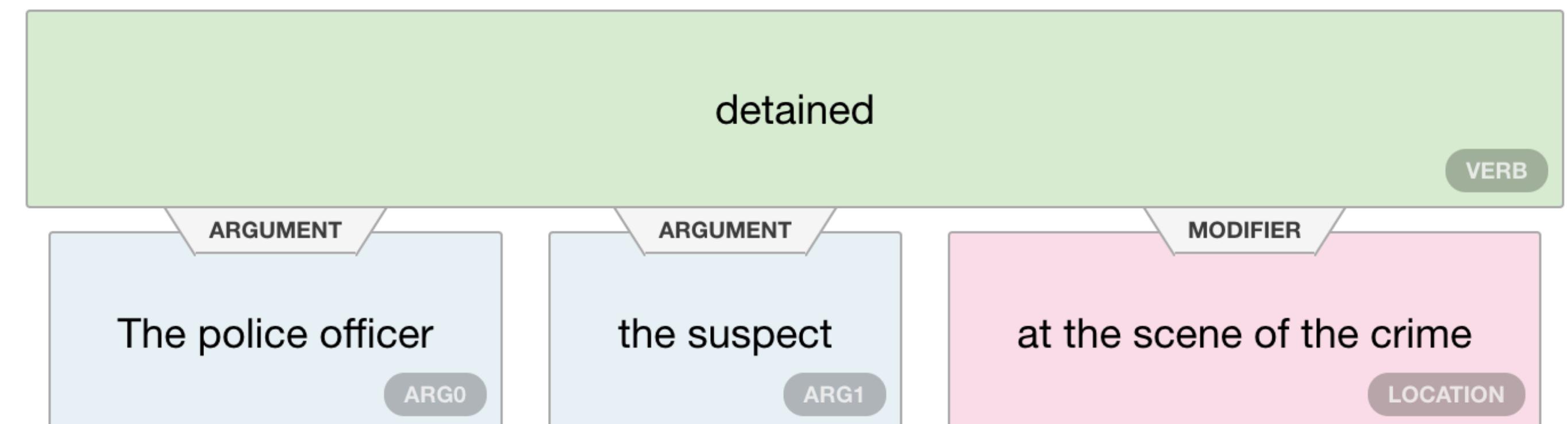
Some men are playing a sport.

## Tasks

- SQuAD: question answering
- SNLI: natural language inference
- SRL: semantic role labeling
- Coref: coreference resolution
- NER: named entity recognition
- SST-5: sentiment analysis

- Different types of frames
  - FrameNet vs PropBank
- Different semantics roles depending on the verb

The police officer detained the suspect at the scene of the crime



(figure credit: <https://demo.allennlp.org/semantic-role-labeling/>)

## Tasks

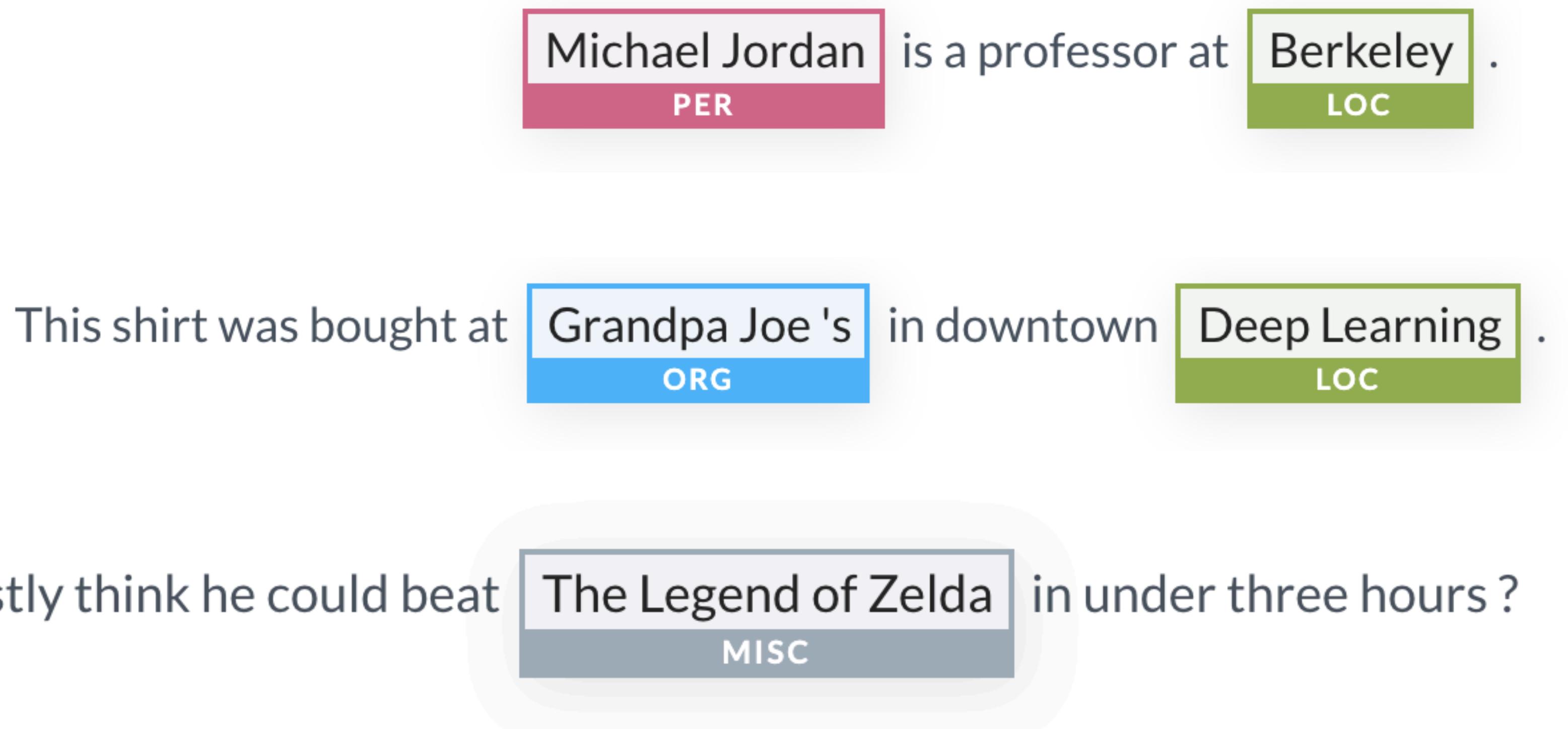
- SQuAD: question answering
- SNLI: natural language inference
- SRL: semantic role labeling
- Coref: coreference resolution
- NER: named entity recognition
- SST-5: sentiment analysis

*“I voted for Nader because he was most aligned with my values,” she said.*

(figure credit: <https://nlp.stanford.edu/projects/coref.shtml>)

## Tasks

- SQuAD: question answering
- SNLI: natural language inference
- SRL: semantic role labeling
- Coref: coreference resolution
- NER: named entity recognition
- SST-5: sentiment analysis



(figure credit: <https://demo.allennlp.org/named-entity-recognition>)

## Tasks

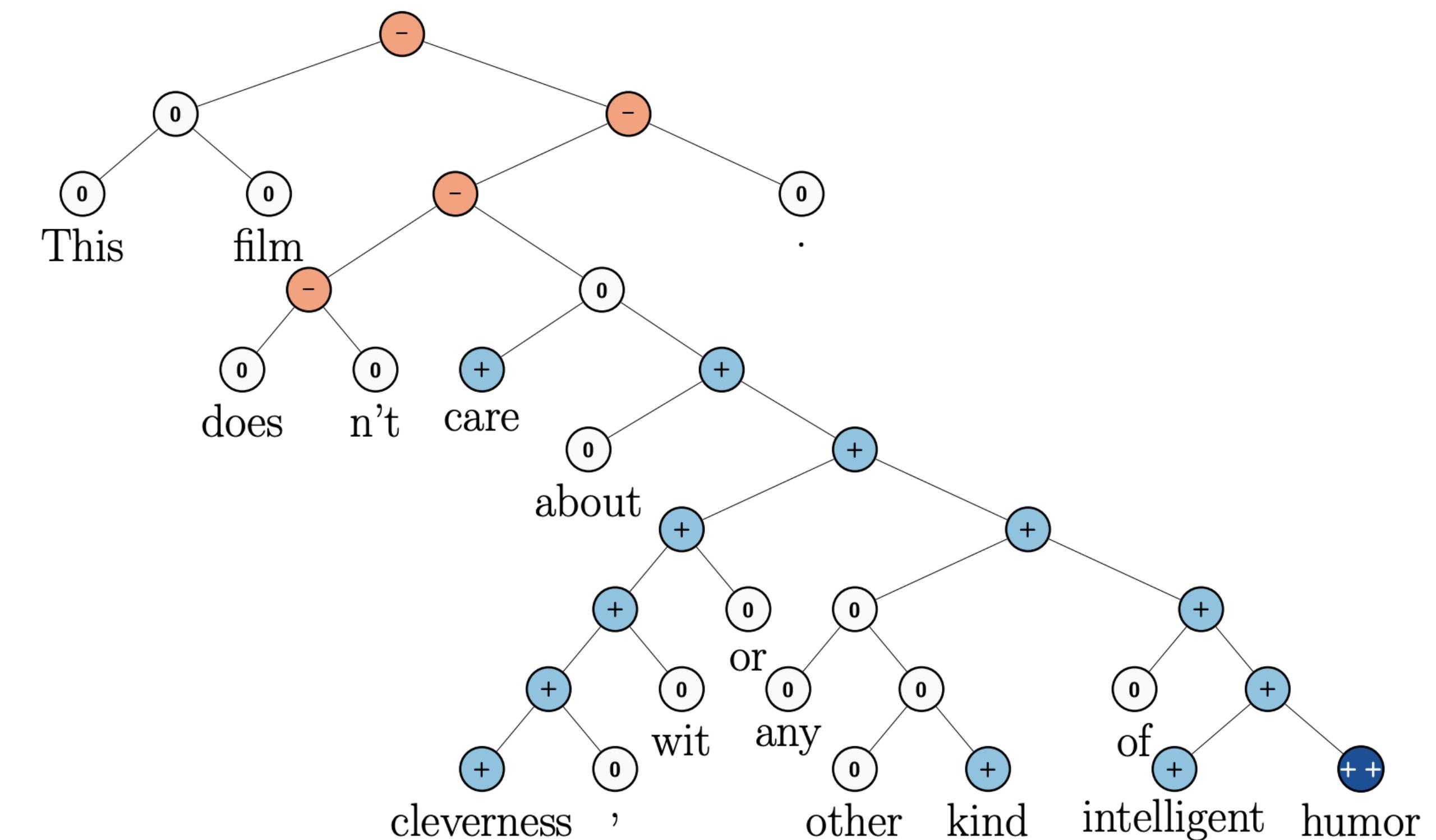
- SQuAD: question answering
- SNLI: natural language inference
- SRL: semantic role labeling
- Coref: coreference resolution
- NER: named entity recognition
- SST-5: sentiment analysis

5 classes

- Very Positive
- Positive
- Neutral
- Negative
- Very Negative

## Stanford Sentiment Treebank

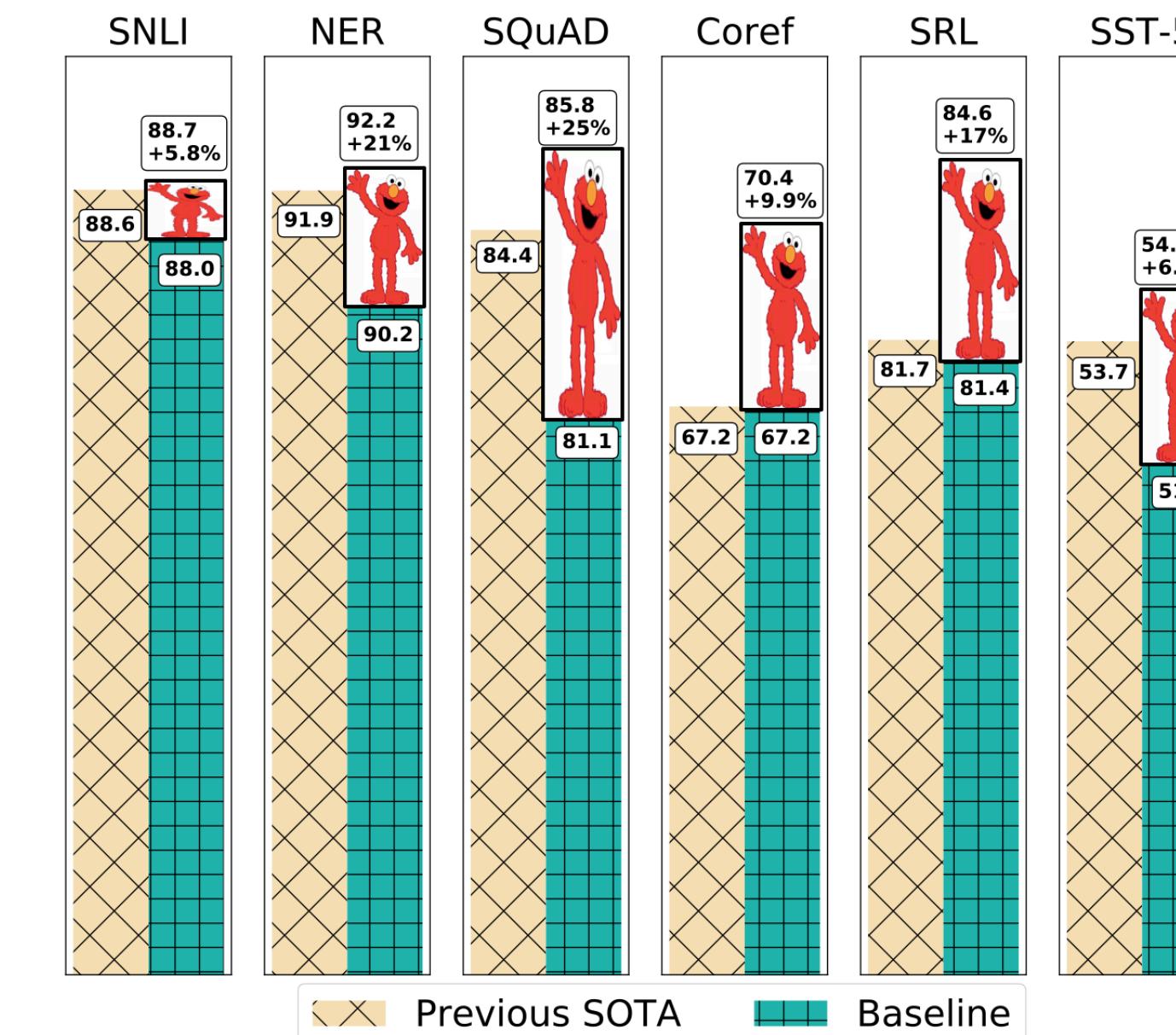
This file doesn't care about cleverness, wit or any other kind of intelligent humor.



# Experimental results

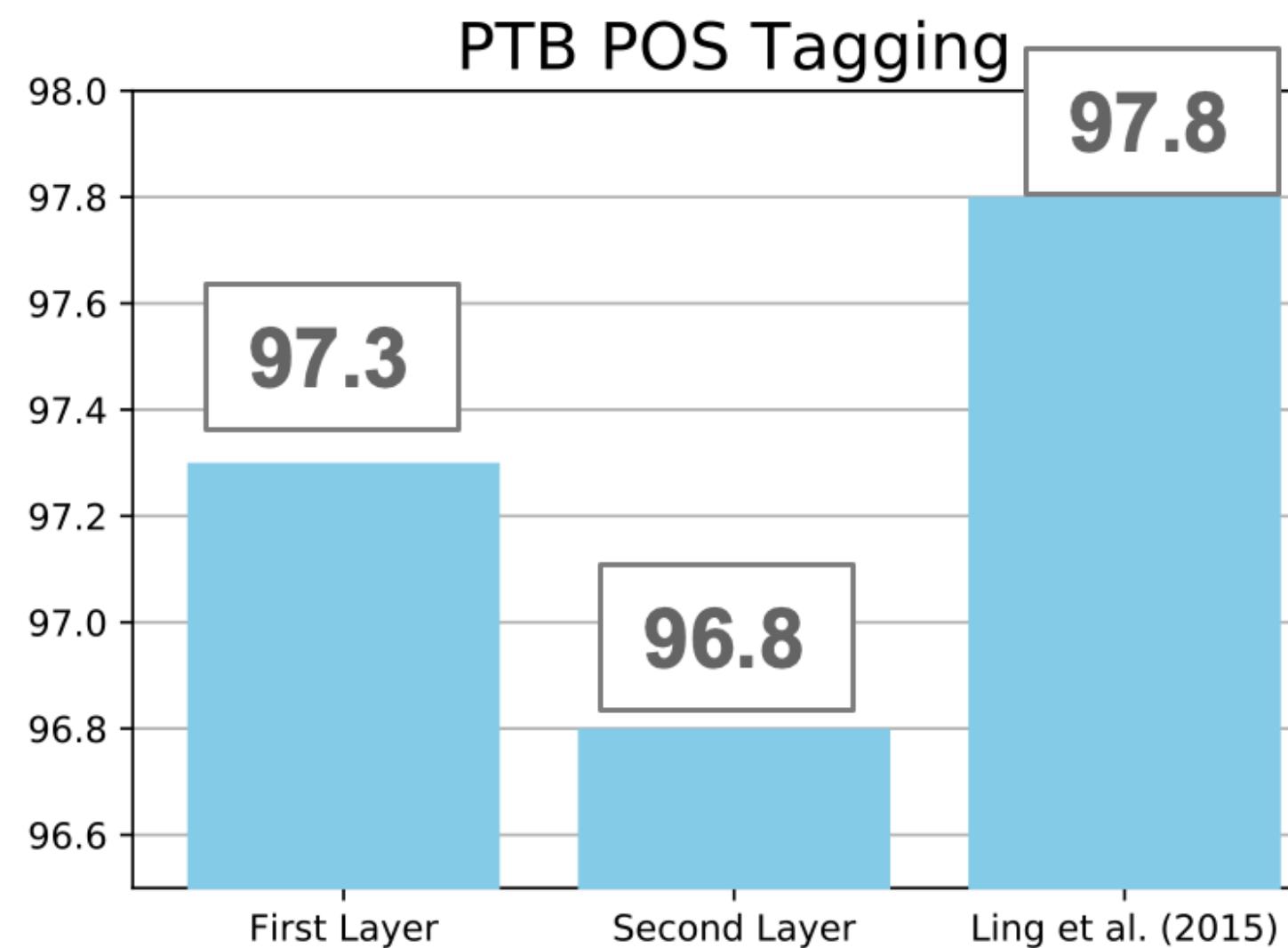
TASK	PREVIOUS SOTA		OUR BASELINE	ELMO + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SQuAD	Liu et al. (2017)	84.4	81.1	85.8	4.7 / 24.9%
SNLI	Chen et al. (2017)	88.6	88.0	88.7 ± 0.17	0.7 / 5.8%
SRL	He et al. (2017)	81.7	81.4	84.6	3.2 / 17.2%
Coref	Lee et al. (2017)	67.2	67.2	70.4	3.2 / 9.8%
NER	Peters et al. (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%

- SQuAD: question answering
- SNLI: natural language inference
- SRL: semantic role labeling
- Coref: coreference resolution
- NER: named entity recognition
- SST-5: sentiment analysis



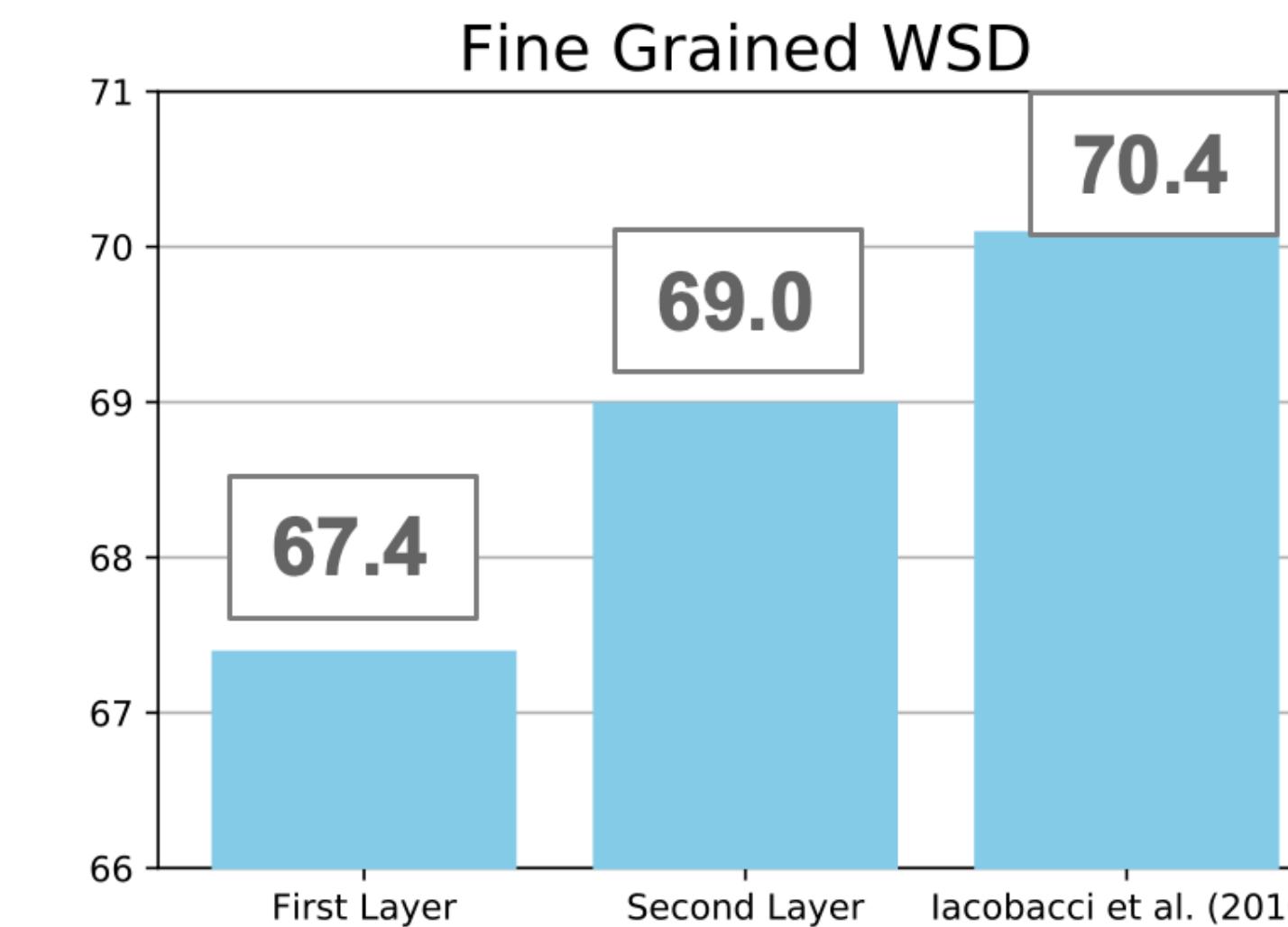
# Intrinsic Evaluation

syntactic information



First Layer > Second Layer

semantic information



Second Layer > First Layer

syntactic information is better represented at lower layers  
while semantic information is captured at higher layers

# Use ELMo in practice

<https://allennlp.org/elmo>

## Pre-trained ELMo Models

Model	Link(Weights/Options File)	# Parameters (Millions)	LSTM Hidden Size/Output size	# Highway Layers>
Small	<a href="#">weights</a> <a href="#">options</a>	13.6	1024/128	1
Medium	<a href="#">weights</a> <a href="#">options</a>	28.0	2048/256	1
Original	<a href="#">weights</a> <a href="#">options</a>	93.6	4096/512	2
Original (5.5B)	<a href="#">weights</a> <a href="#">options</a>	93.6	4096/512	2

```
from allennlp.modules.elmo import Elmo, batch_to_ids

options_file = "https://allennlp.s3.amazonaws.com/models/elmo/2x4096"
weight_file = "https://allennlp.s3.amazonaws.com/models/elmo/2x4096.

# Compute two different representation for each token.
# Each representation is a linear weighted combination for the
# 3 layers in ELMo (i.e., charcnn, the outputs of the two BiLSTM)
elmo = Elmo(options_file, weight_file, 2, dropout=0)

# use batch_to_ids to convert sentences to character ids
sentences = [['First', 'sentence', '.'], ['Another', '.']]
character_ids = batch_to_ids(sentences)

embeddings = elmo(character_ids)
```

Also available in TensorFlow

# BERT

- First released in Oct 2018.
- NAACL'19: BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

How is BERT different from ELMo?

- #1. Two unidirectional context vs **bidirectional** context
- #2. LSTMs vs **Transformers** (will explain more later)
- #3. The weights are not frozen, called **fine-tuning**

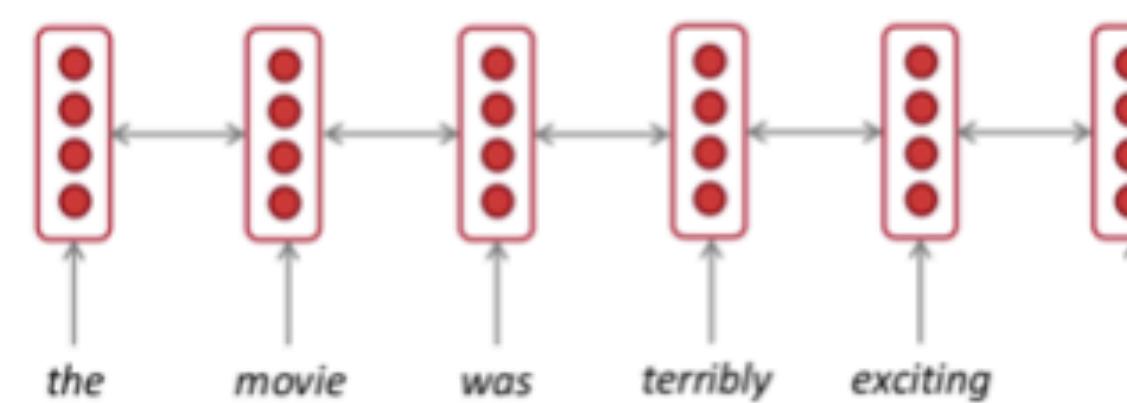


# Bidirectional encoders

- Language models only use left context or right context (although ELMo used two independent LMs from each direction).
- Language understanding is bidirectional

## Bidirectional RNNs

Bidirectionality is important in language representations:



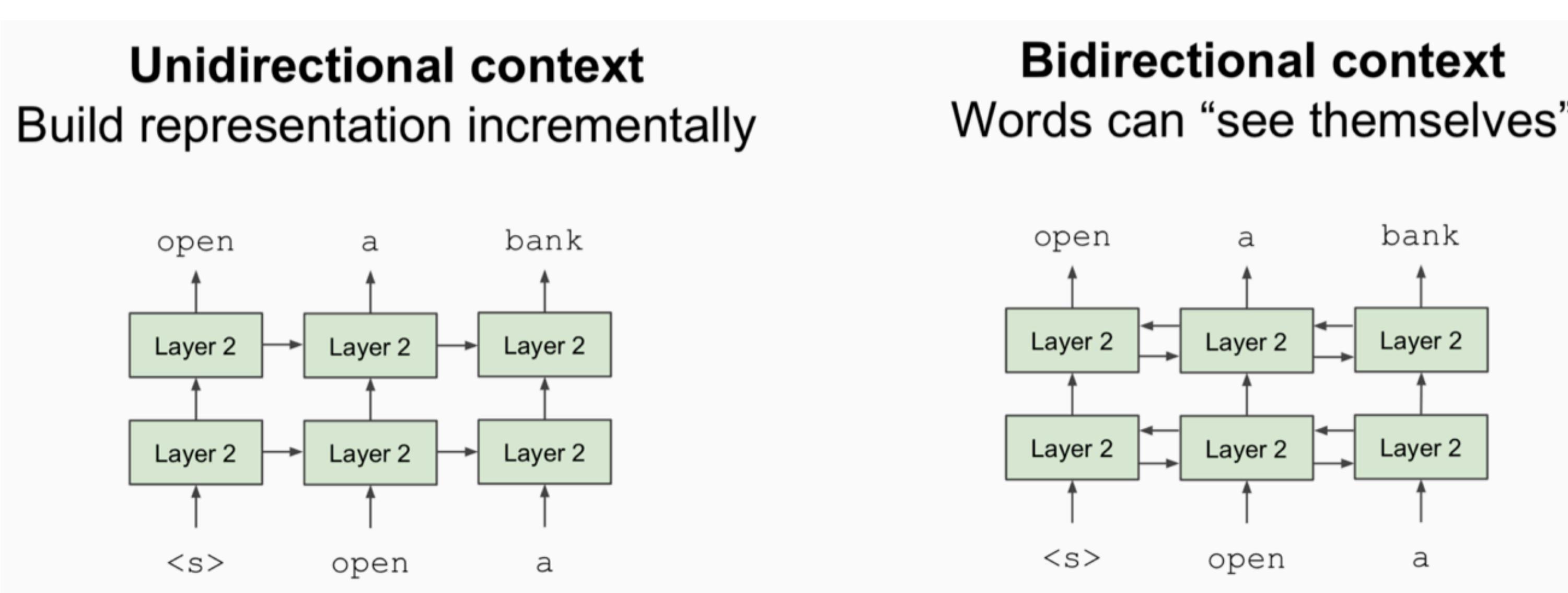
*terribly:*

- left context “the movie was”
- right context “exciting !”

Why are LMs unidirectional?

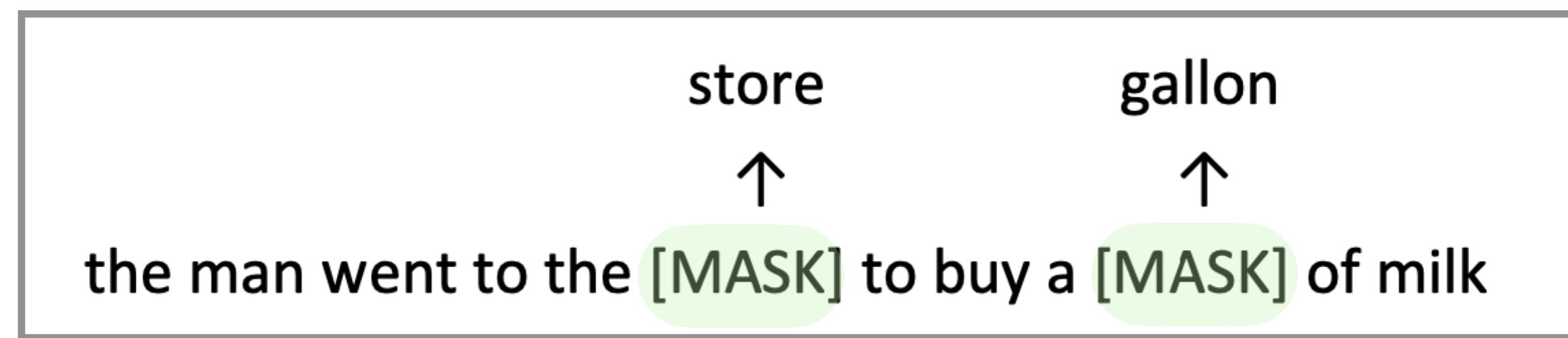
# Bidirectional encoders

- Language models only use left context or right context (although ELMo used two independent LMs from each direction).
- Language understanding is bidirectional



# Masked language models (MLMs)

- Solution: Mask out 15% of the input words, and then predict the masked words



- Too little masking: too expensive to train
  - Too much masking: not enough context

# Masked language models (MLMs)

A little more complex  
(don't always replace with [ MASK ]):

Example: my dog is hairy, we replace the word hairy

- 80% of time: replace word with [ MASK ] token  
my dog is [ MASK ]
- 10% of time: replace word with random word  
my dog is apple
- 10% of time: keep word unchanged to bias representation toward actual observed word  
my dog is hairy

Because [ MASK ] is never seen when BERT is used...

# Next sentence prediction (NSP)

Always sample two sentences, predict whether the second sentence is followed after the first one.

**Input** = [CLS] the man went to [MASK] store [SEP]  
he bought a gallon [MASK] milk [SEP]

Label = IsNext

**Input** = [CLS] the man [MASK] to the store [SEP]  
penguin [MASK] are flight ##less birds [SEP]

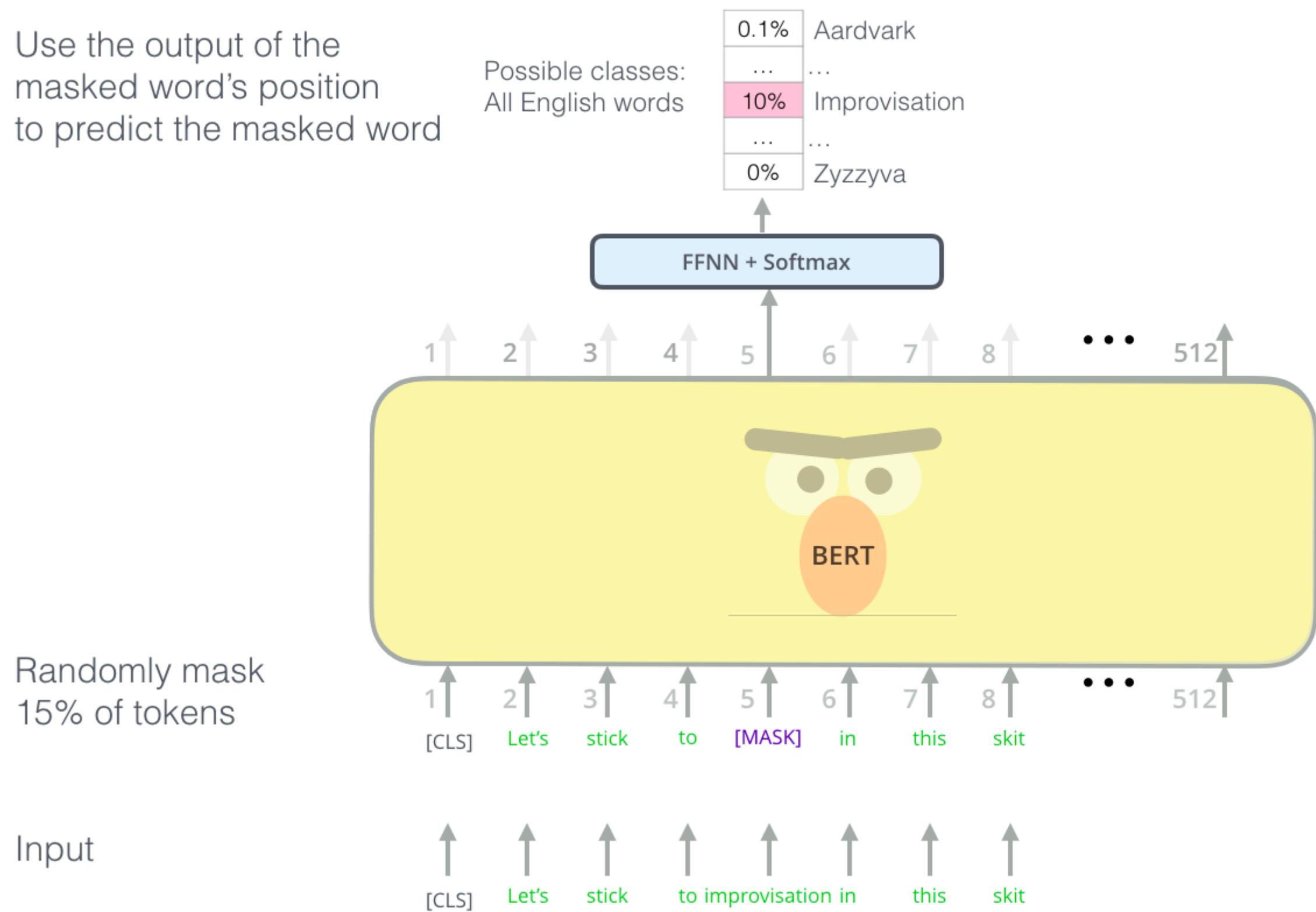
Label = NotNext

Recent papers show that NSP is not necessary...

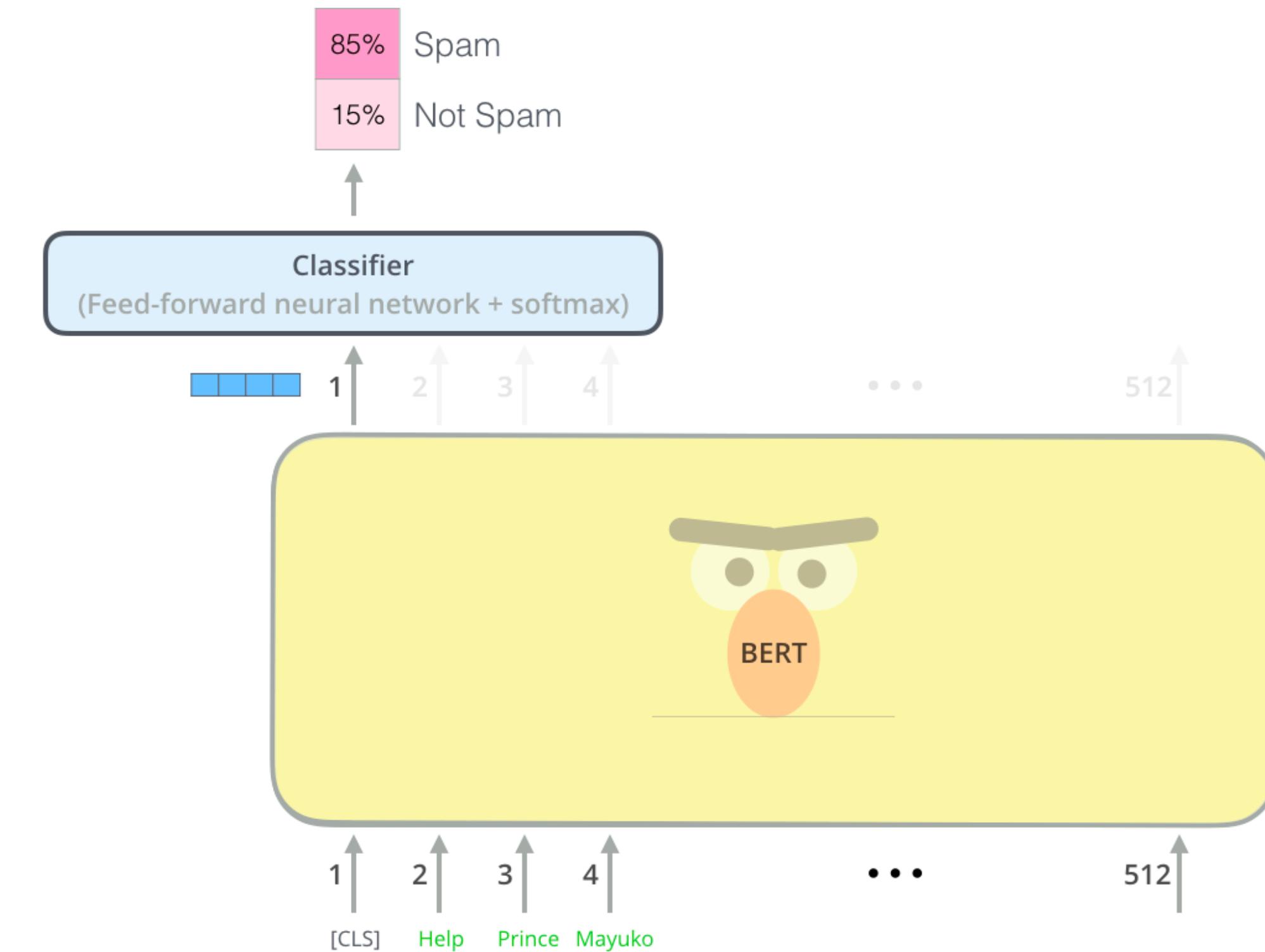
(Joshi\*, Chen\* et al, 2019) :SpanBERT: Improving Pre-training by Representing and Predicting Spans  
(Liu et al, 2019): RoBERTa: A Robustly Optimized BERT Pretraining Approach

# Pre-training and fine-tuning

Use the output of the masked word's position to predict the masked word



Pre-training



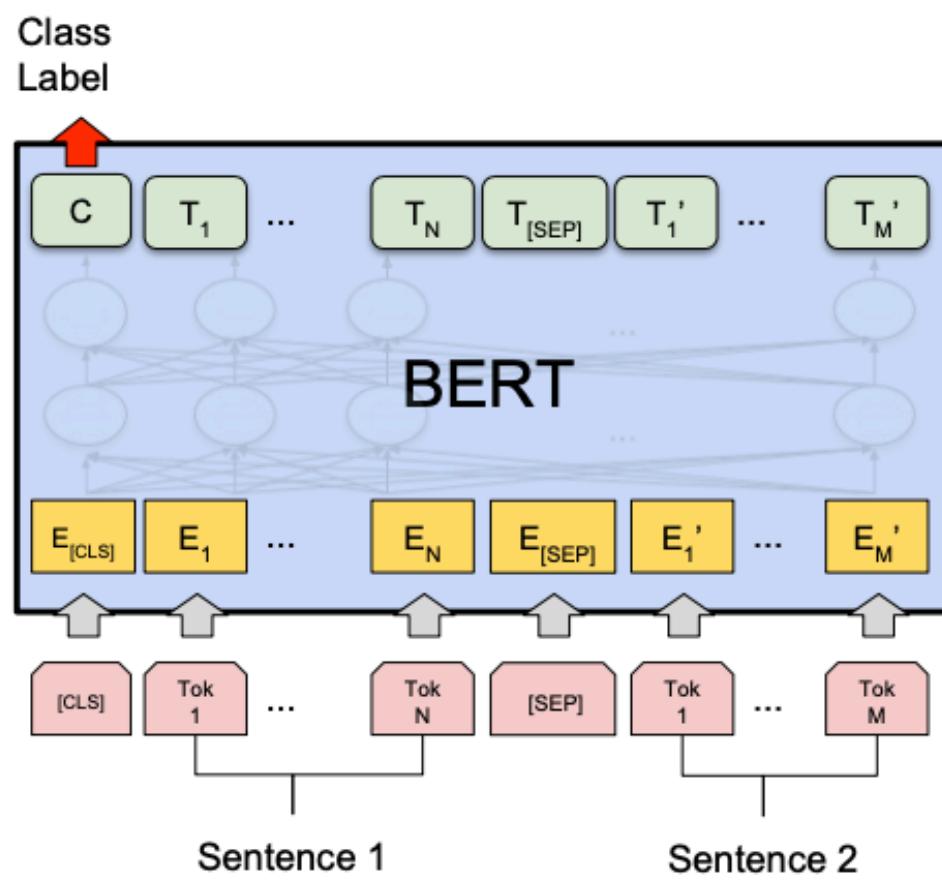
Fine-tuning

(figure credit: [Jay Alammar](#)

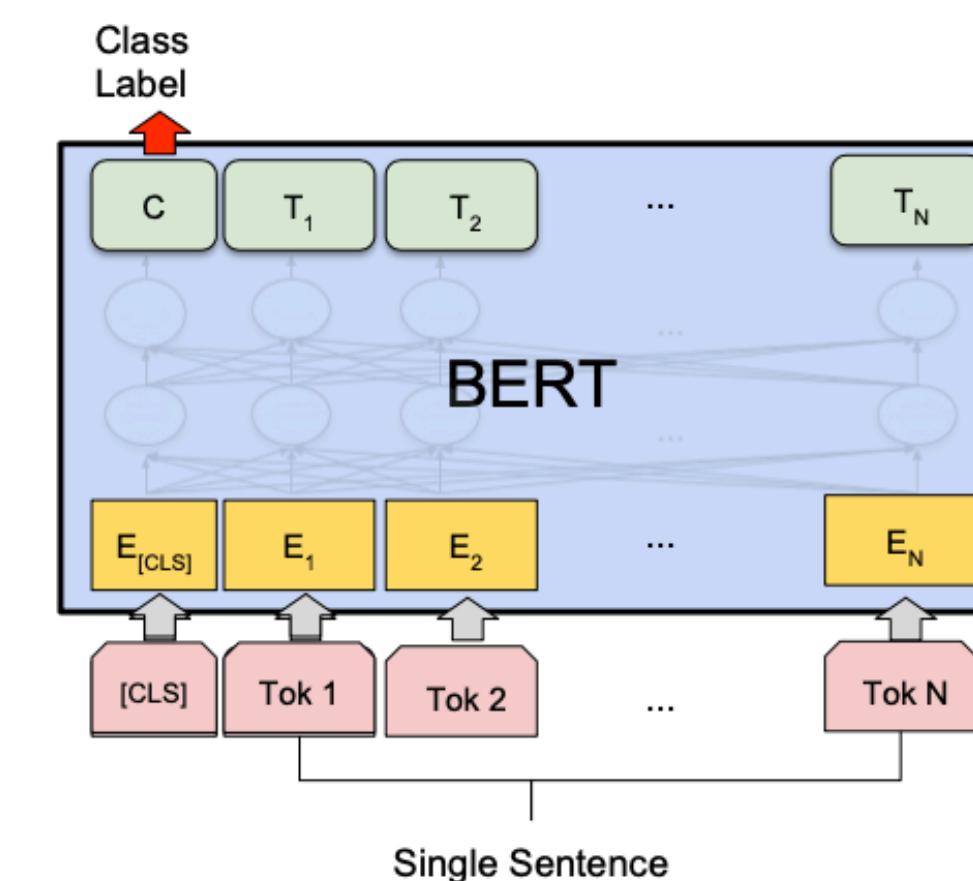
<http://jalammar.github.io/illustrated-bert/>

Key idea: all the weights are fine-tuned on downstream tasks

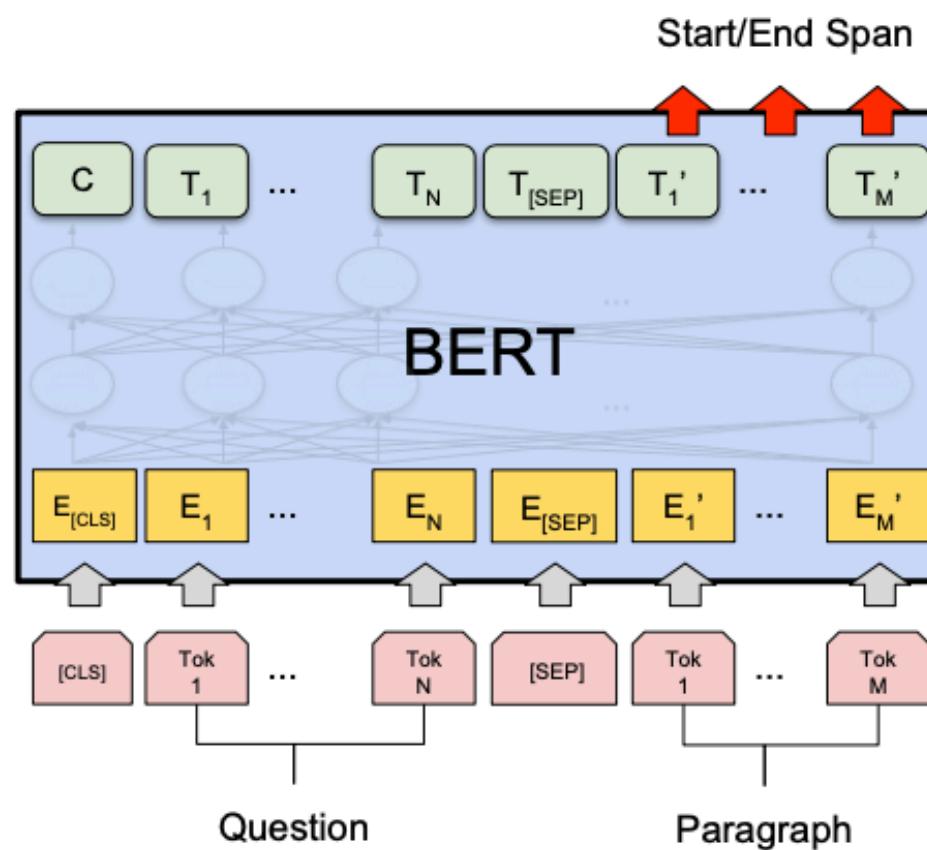
# Applications



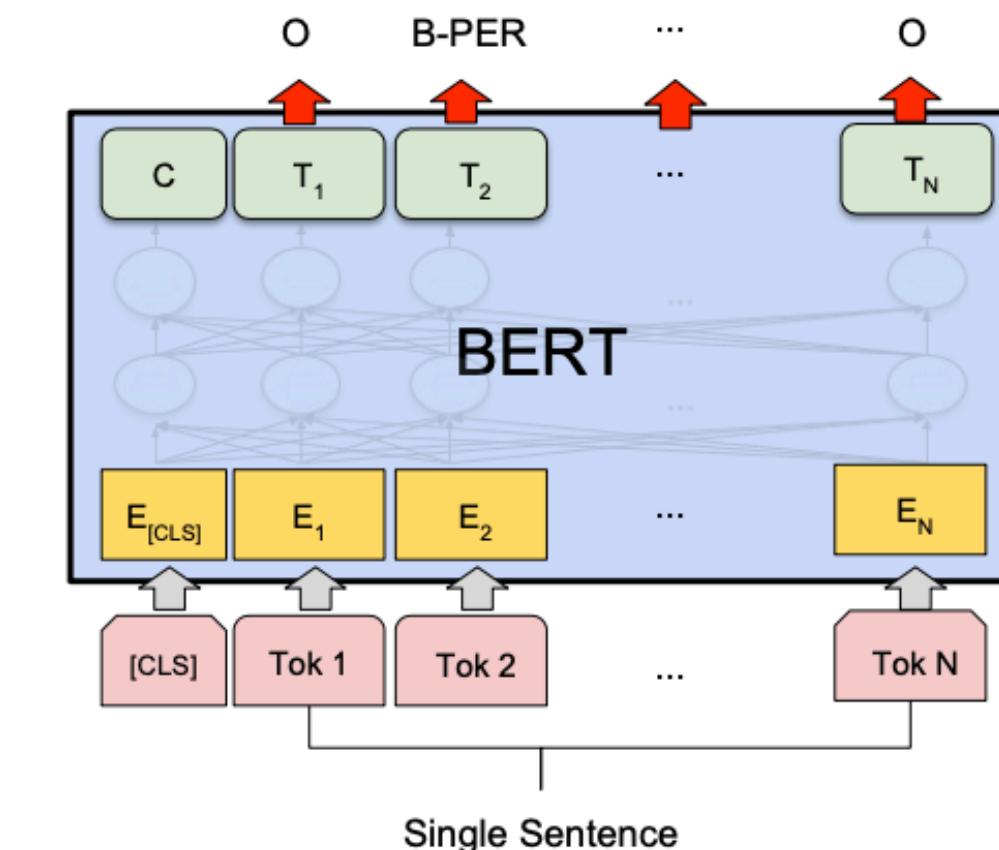
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



(b) Single Sentence Classification Tasks:  
SST-2, CoLA

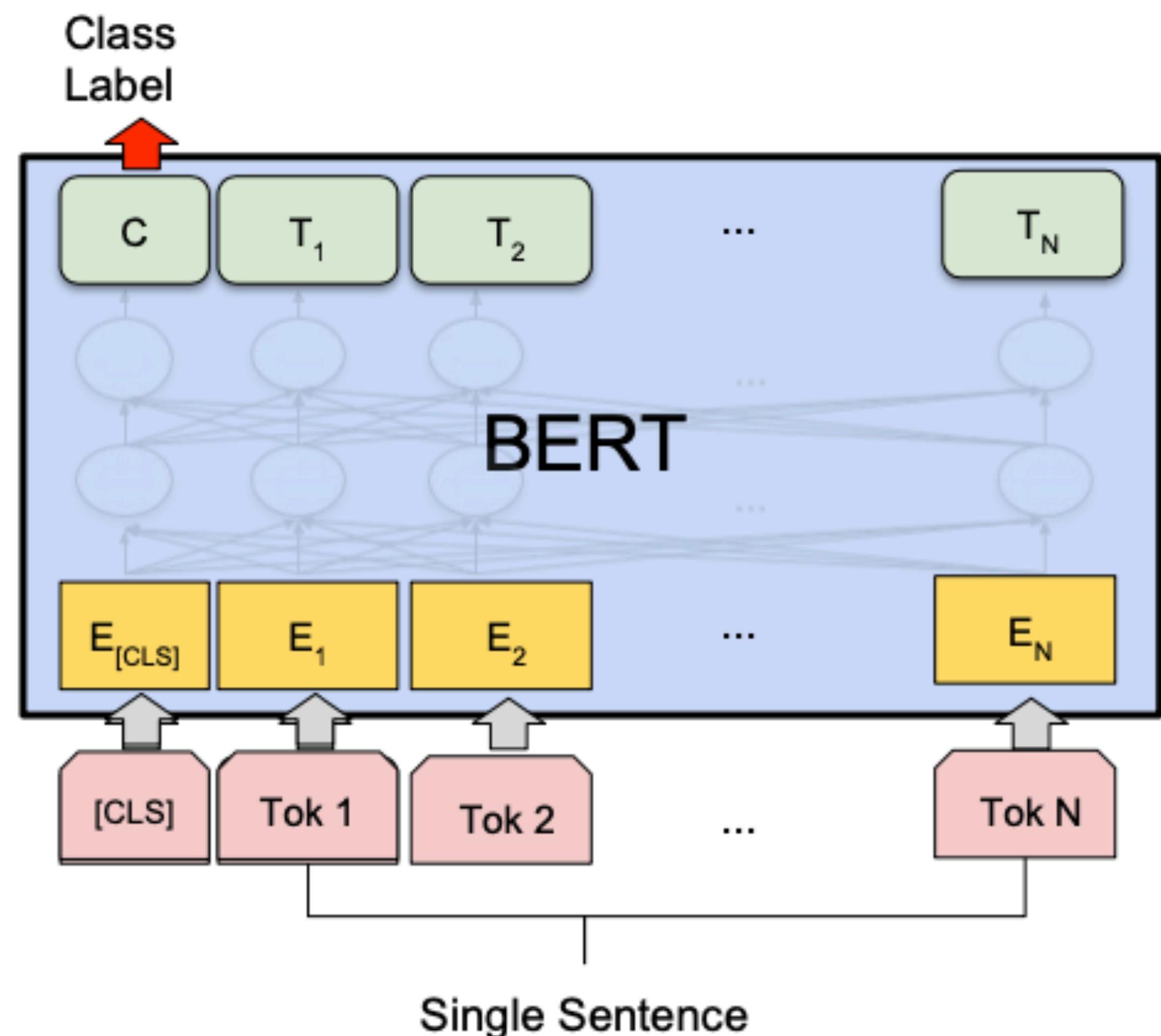


(c) Question Answering Tasks:  
SQuAD v1.1



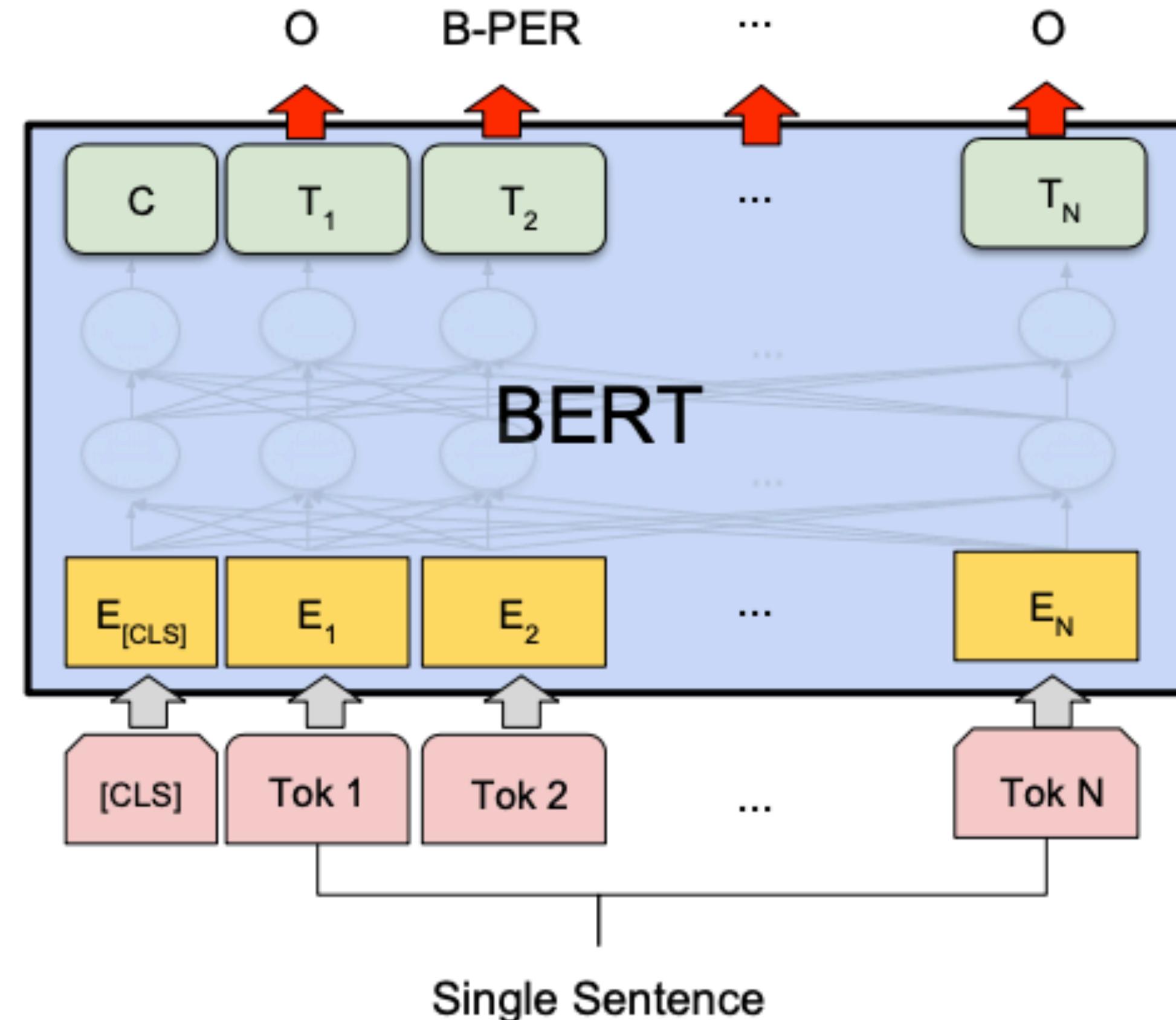
(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

# Applications



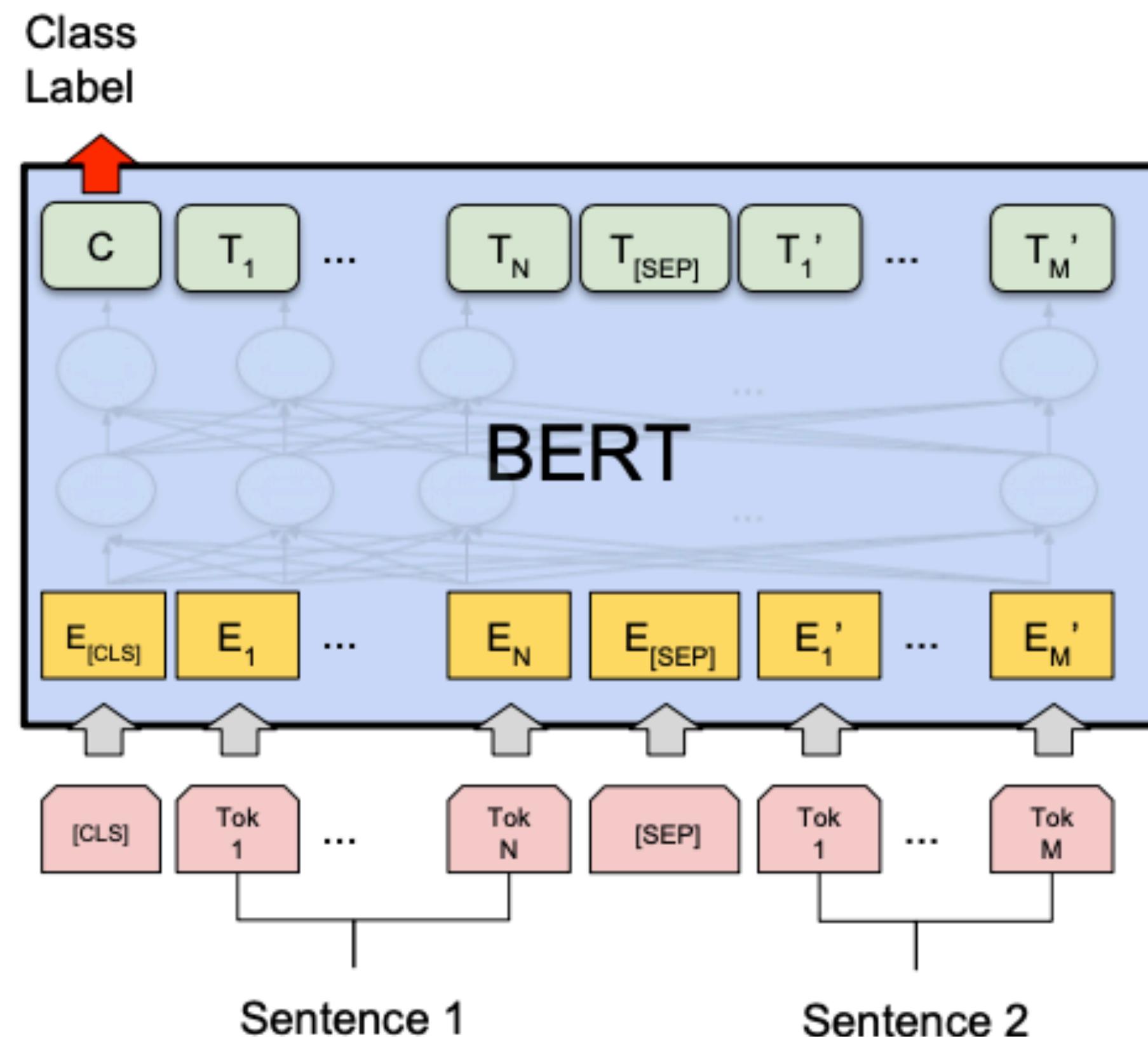
(b) Single Sentence Classification Tasks:  
SST-2, CoLA

# Applications



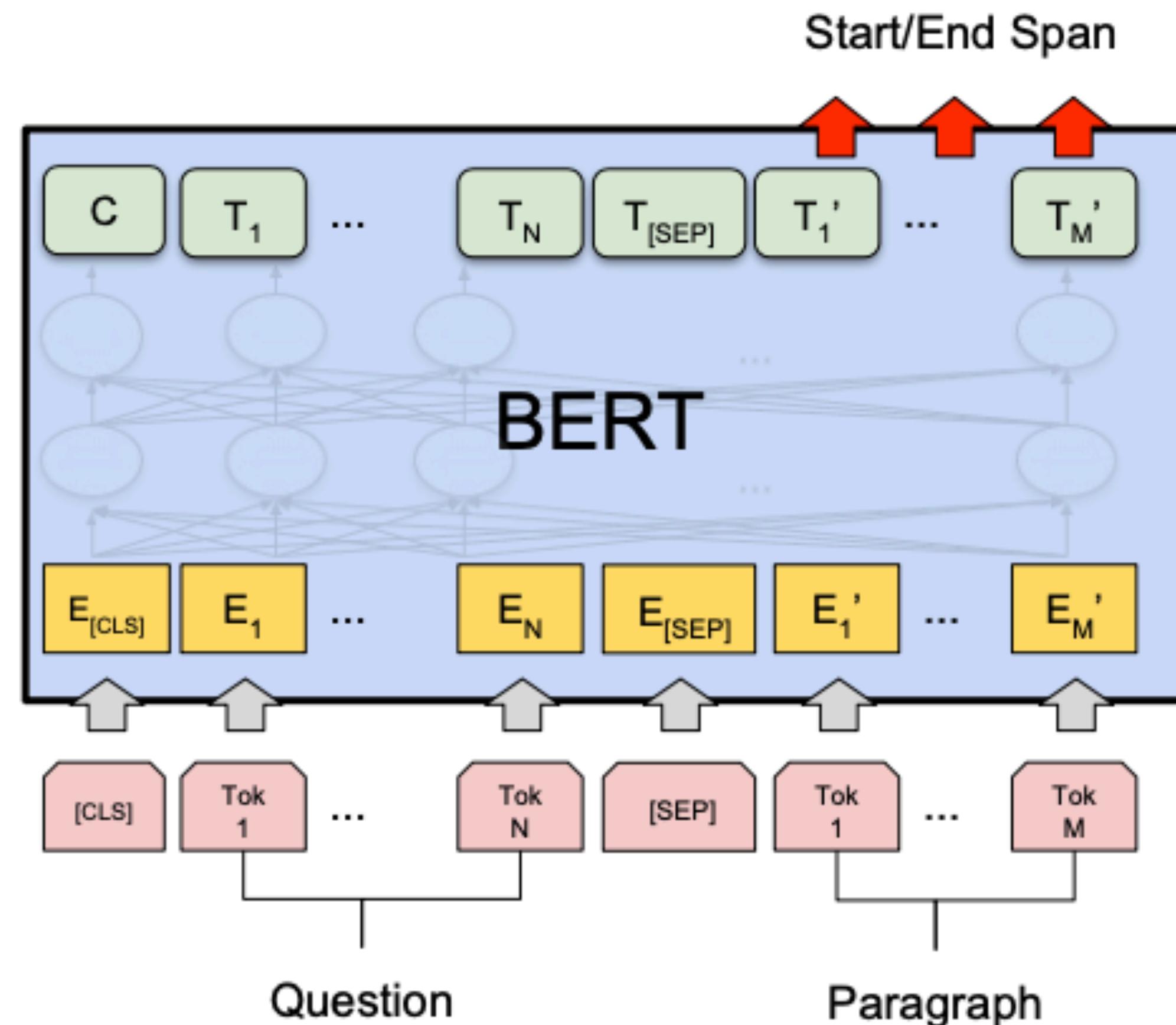
(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

# Applications



(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG

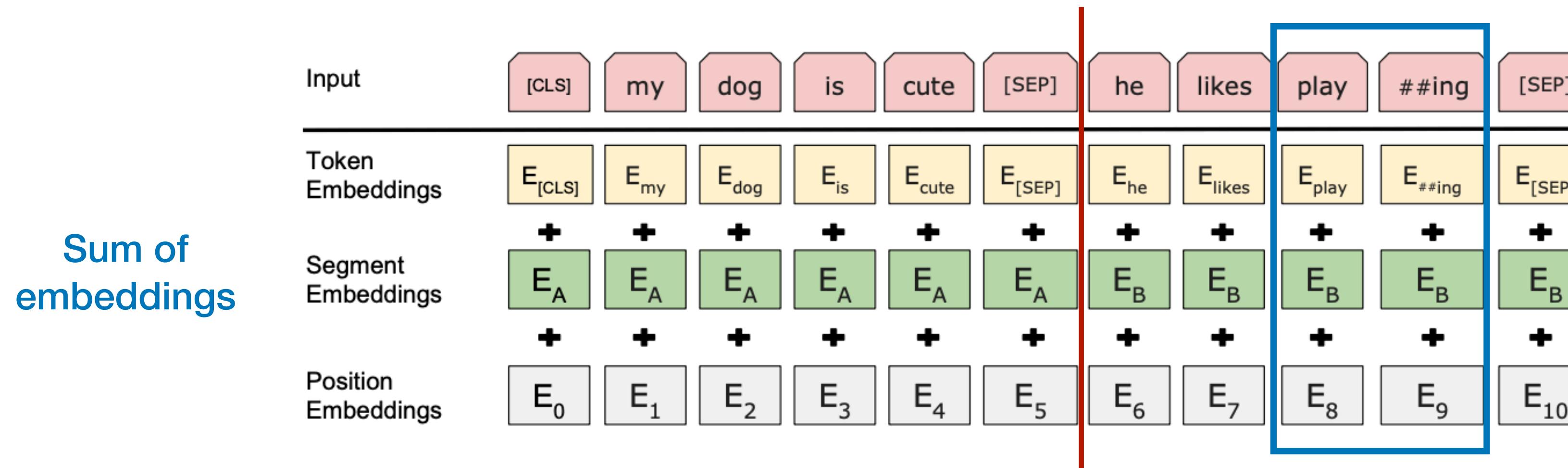
# Applications



(c) Question Answering Tasks:  
SQuAD v1.1

# More details

- Input representations



- Use word pieces instead of words:  $\text{playing} \Rightarrow \text{play } \#\#\text{ing}$  (30K token vocabulary)
- Trained 40 epochs on Wikipedia (2.5B tokens) + BookCorpus (0.8B tokens)
- Released two model sizes: BERT\_base, BERT\_large

# Experimental results

BiLSTM: 63.9

## Entailment

- MNLI: multilingual NLI
- QNLI: NLI with SQuAD data
- RTE: Textual Entailment

## Similarity

- QQP: Quora Question Pairs
- STS-B: Semantic Textual Similarity
- MRPC: MS Research Paraphrase Corpus

## Other

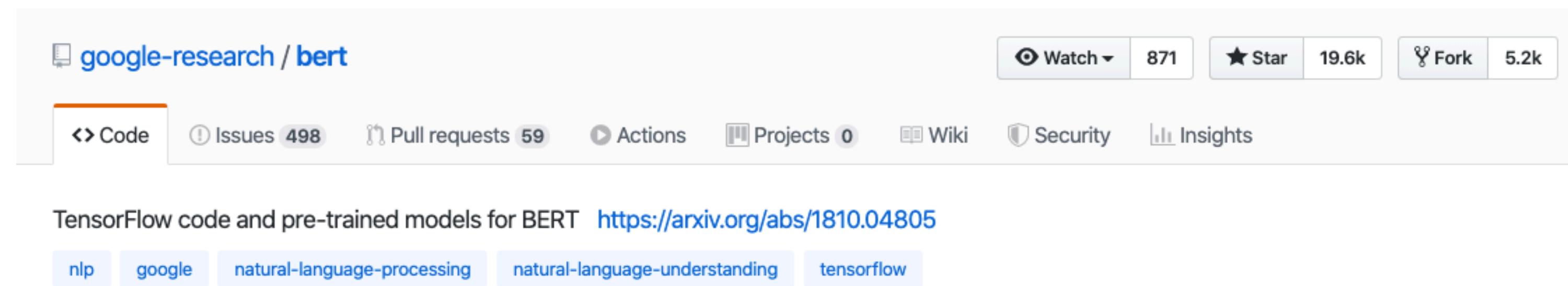
- SST-2: sentiment analysis
- CoLA: Linguistic acceptability
- SQuAD: question answering

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
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	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.1	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>91.1</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>81.9</b>

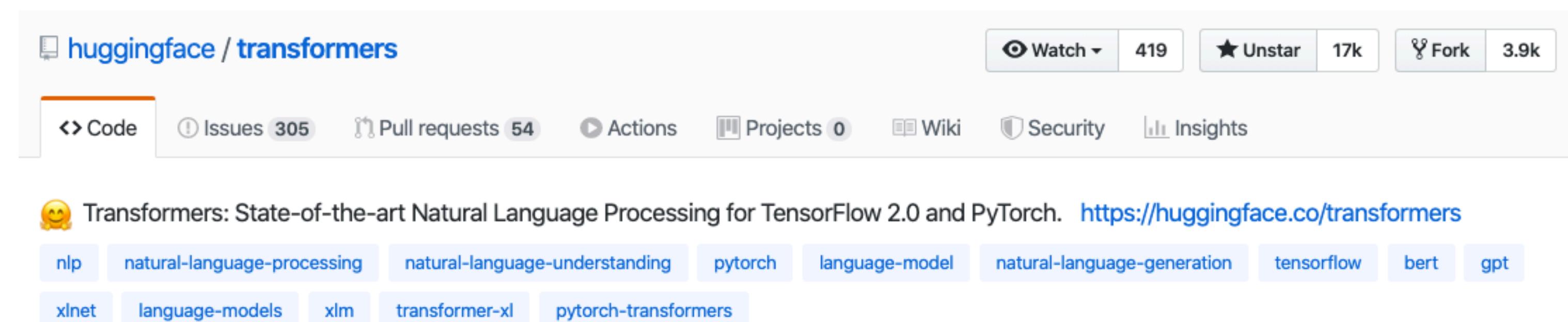
Model	data	bsz	steps	SQuAD (v1.1/2.0)	MNLI-m	SST-2
<b>RoBERTa</b>						
with BOOKS + WIKI	16GB	8K	100K	93.6/87.3	89.0	95.3
+ additional data (§3.2)	160GB	8K	100K	94.0/87.7	89.3	95.6
+ pretrain longer	160GB	8K	300K	94.4/88.7	90.0	96.1
+ pretrain even longer	160GB	8K	500K	<b>94.6/89.4</b>	<b>90.2</b>	<b>96.4</b>
<b>BERT<sub>LARGE</sub></b>						
with BOOKS + WIKI	13GB	256	1M	90.9/81.8	86.6	93.7
<b>XLNet<sub>LARGE</sub></b>						
with BOOKS + WIKI	13GB	256	1M	94.0/87.8	88.4	94.4
+ additional data	126GB	2K	500K	94.5/88.8	89.8	95.6

# Use BERT in practice

**TensorFlow:** <https://github.com/google-research/bert>



**PyTorch:** <https://github.com/huggingface/transformers>



# Contextualized word embeddings in context

- TagLM (Peters et, 2017)
- CoVe (McCann et al. 2017)
- ULMfit (Howard and Ruder, 2018)
- **ELMo (Peters et al, 2018)**
- OpenAI GPT (Radford et al, 2018)
- **BERT (Devlin et al, 2018)**
- OpenAI GPT-2 (Radford et al, 2019)
- XLNet (Yang et al, 2019)
- SpanBERT (Joshi et al, 2019)
- RoBERTa (Liu et al, 2019)
- ALBERT (Lan et al, 2019)
- DistilBERT (Sanh et al, 2019)
- ELECTRA (Clark et al, 2020)
- ...

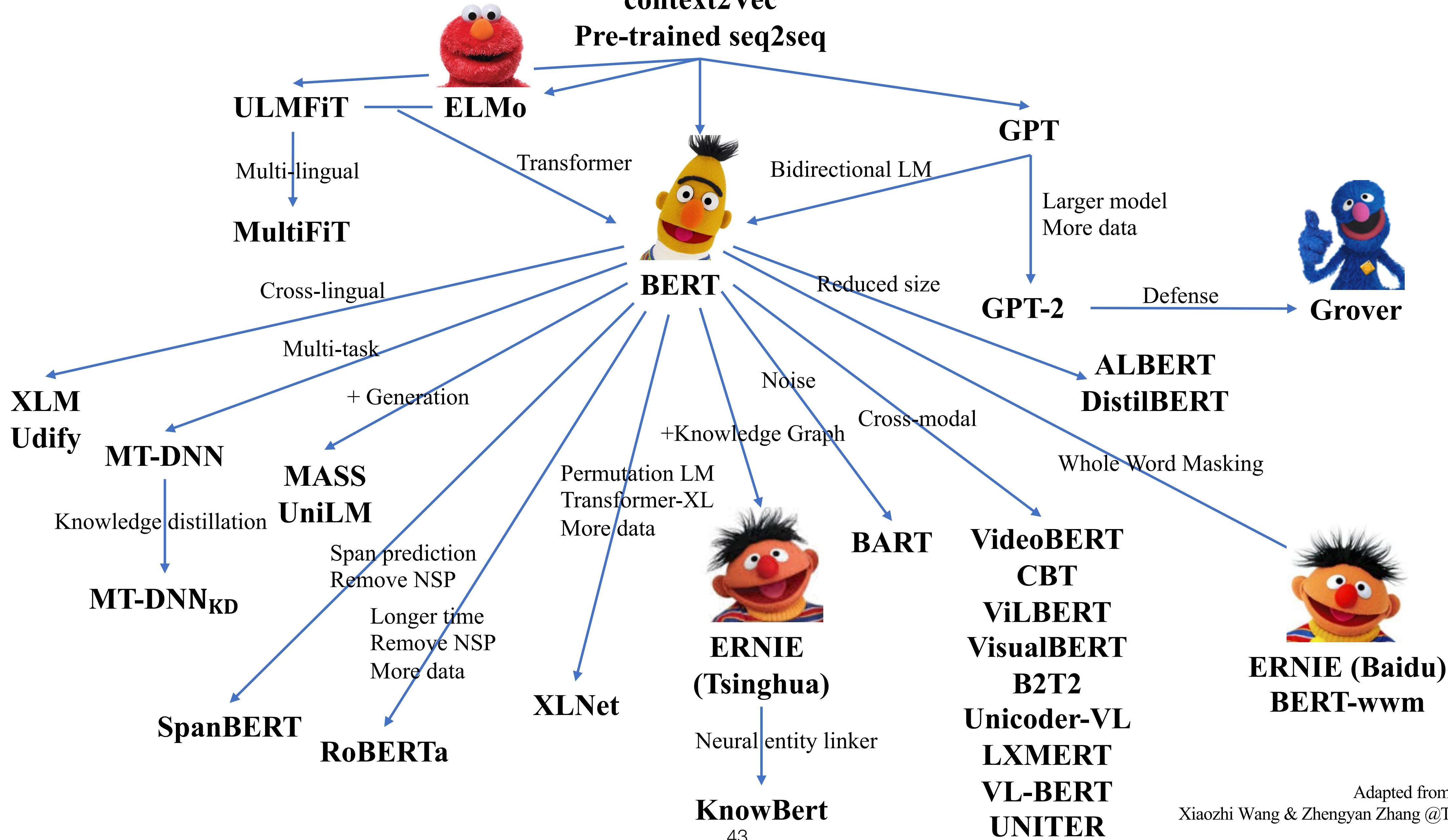


[https://github.com/  
huggingface/transformers](https://github.com/huggingface/transformers)

See <https://huggingface.co/transformers/>  
for more information and models

# Semi-supervised Sequence Learning

# context2Vec



All of these models are Transformer models

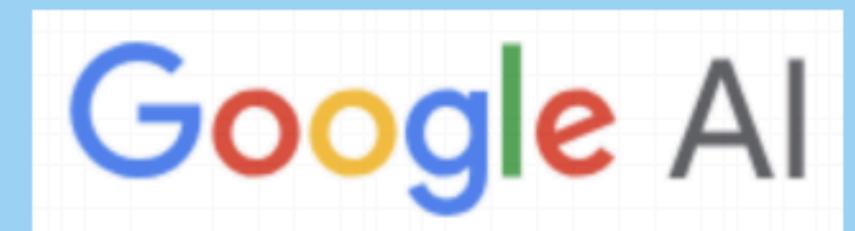
ELMo  
Oct 2017  
Training:  
800M words  
42 GPU days



GPT  
June 2018  
Training  
800M words  
240 GPU days



BERT  
Oct 2018  
Training  
3.3B words  
256 TPU days  
~320–560  
GPU days



GPT-2  
Feb 2019  
Training  
40B words  
~2048 TPU v3  
days according to  
[a reddit thread](#)



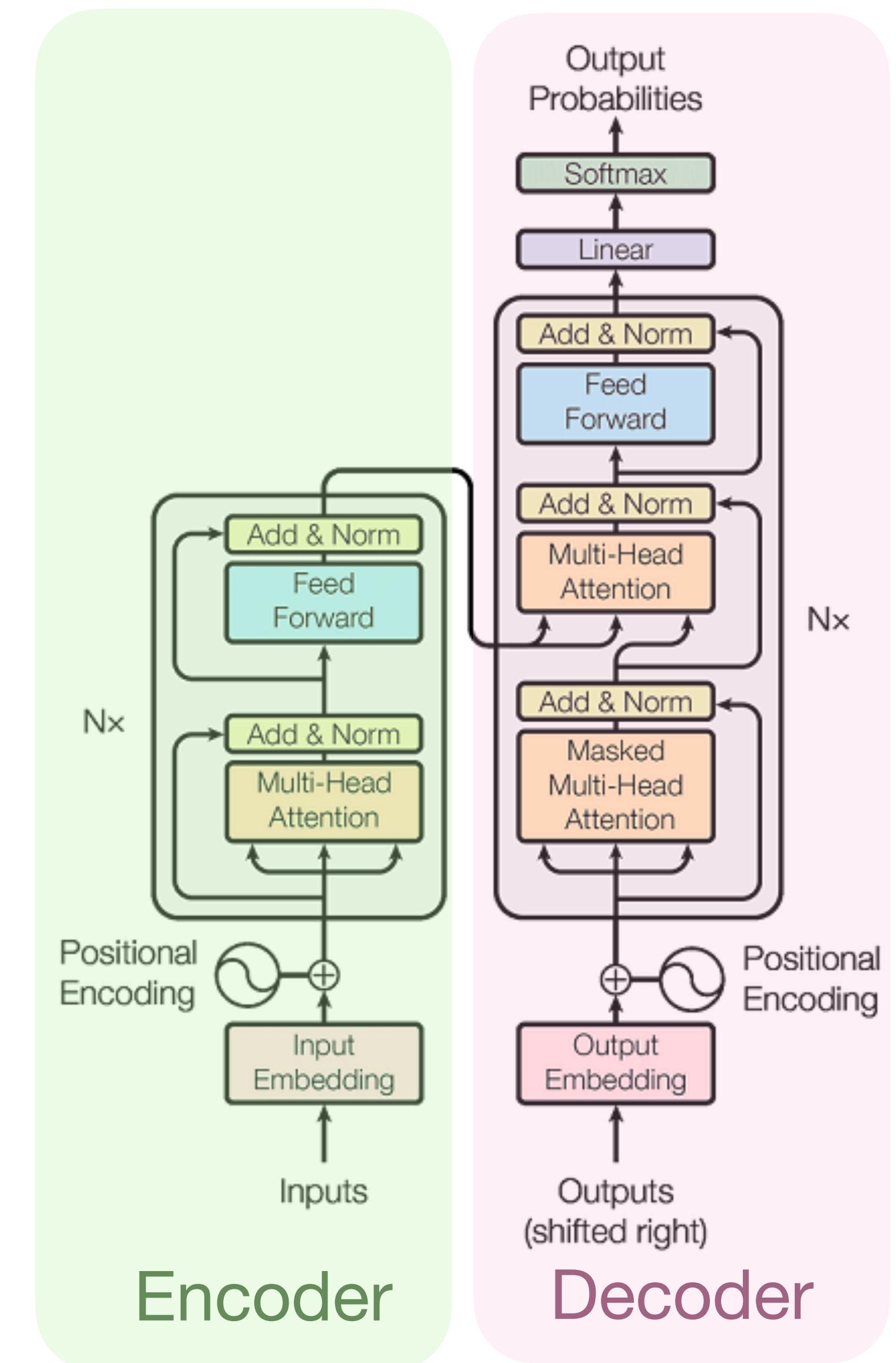
XL-Net,  
ERNIE,  
Grover  
RoBERTa, T5  
July 2019—



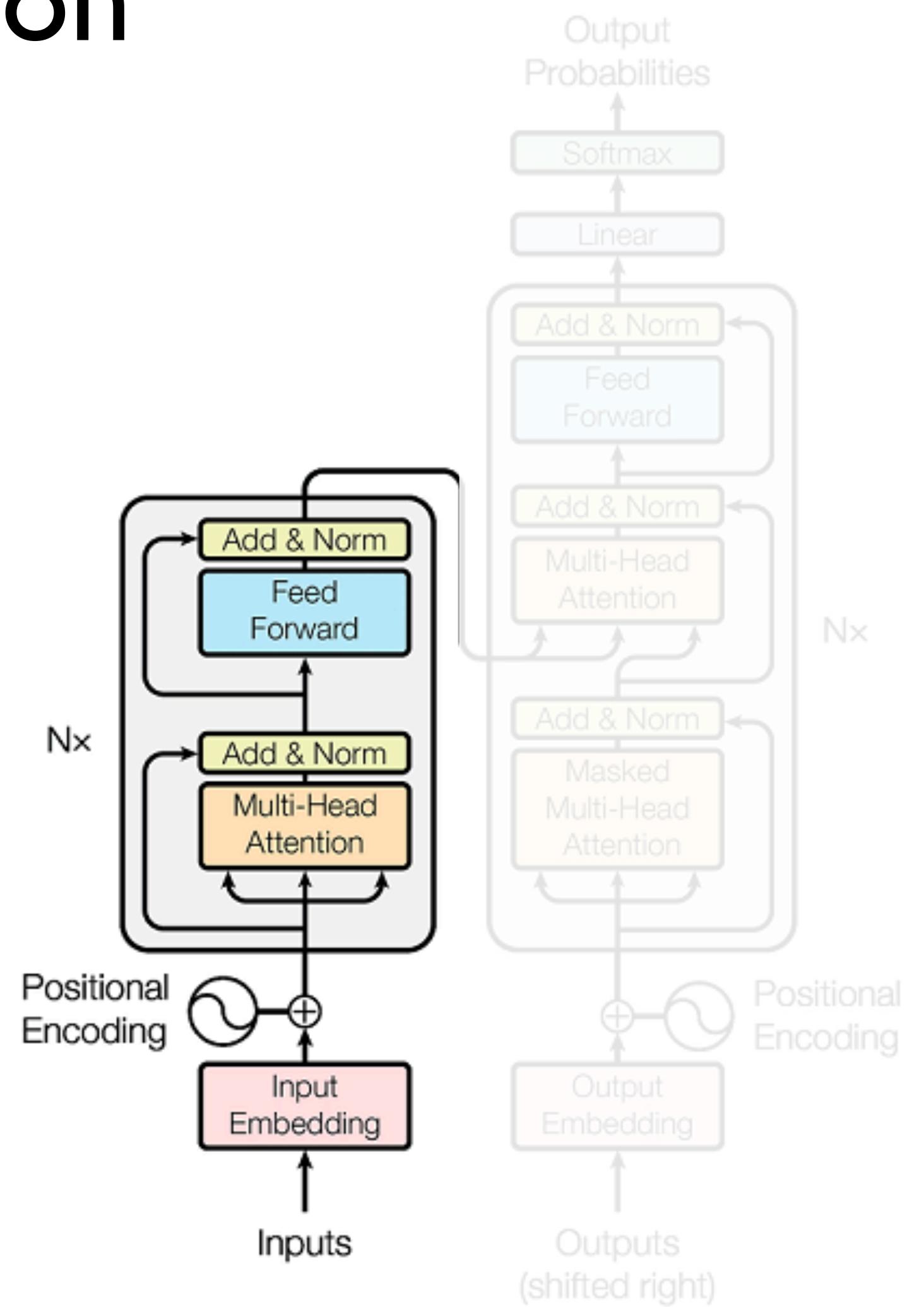
# Transformers

# Transformers

- NIPS'17: Attention is All You Need
- Originally proposed for NMT ([encoder-decoder](#) framework)
- Used as the base model of BERT (encoder only)
- Key idea: **Multi-head self-attention**
- No recurrence structure any more so it trains much faster

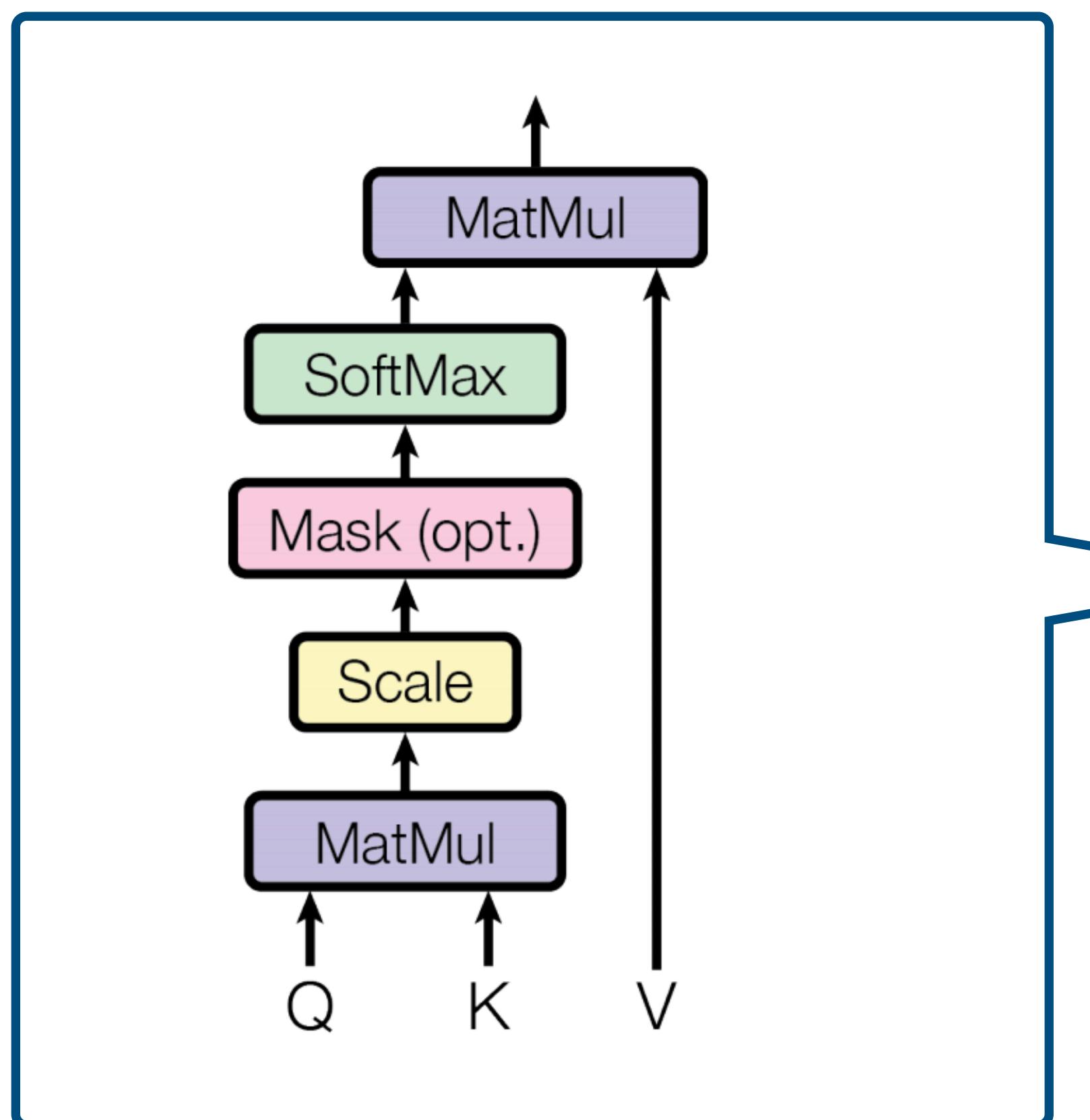


# Multi-head self-attention

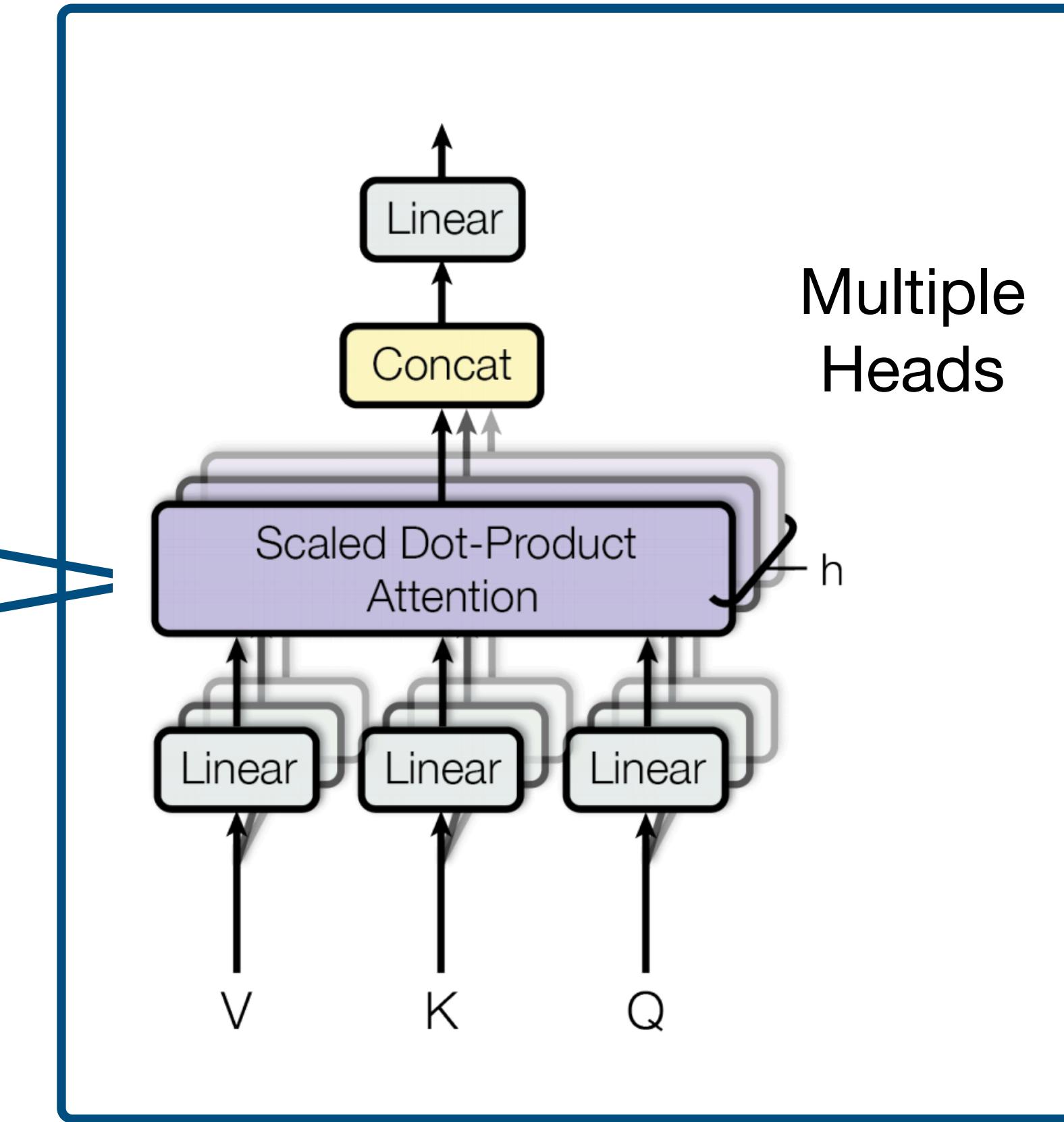


# Multi-head self-attention

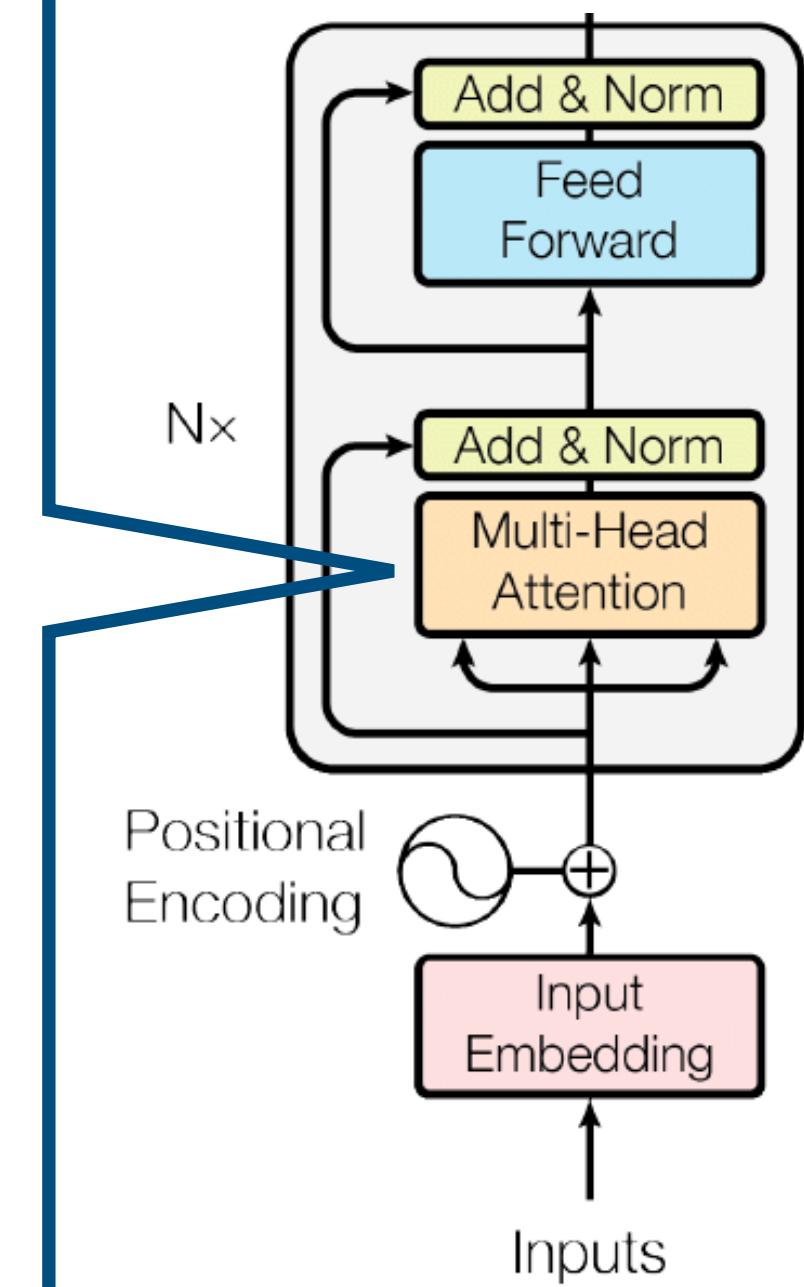
Scaled Dot-Product Attention



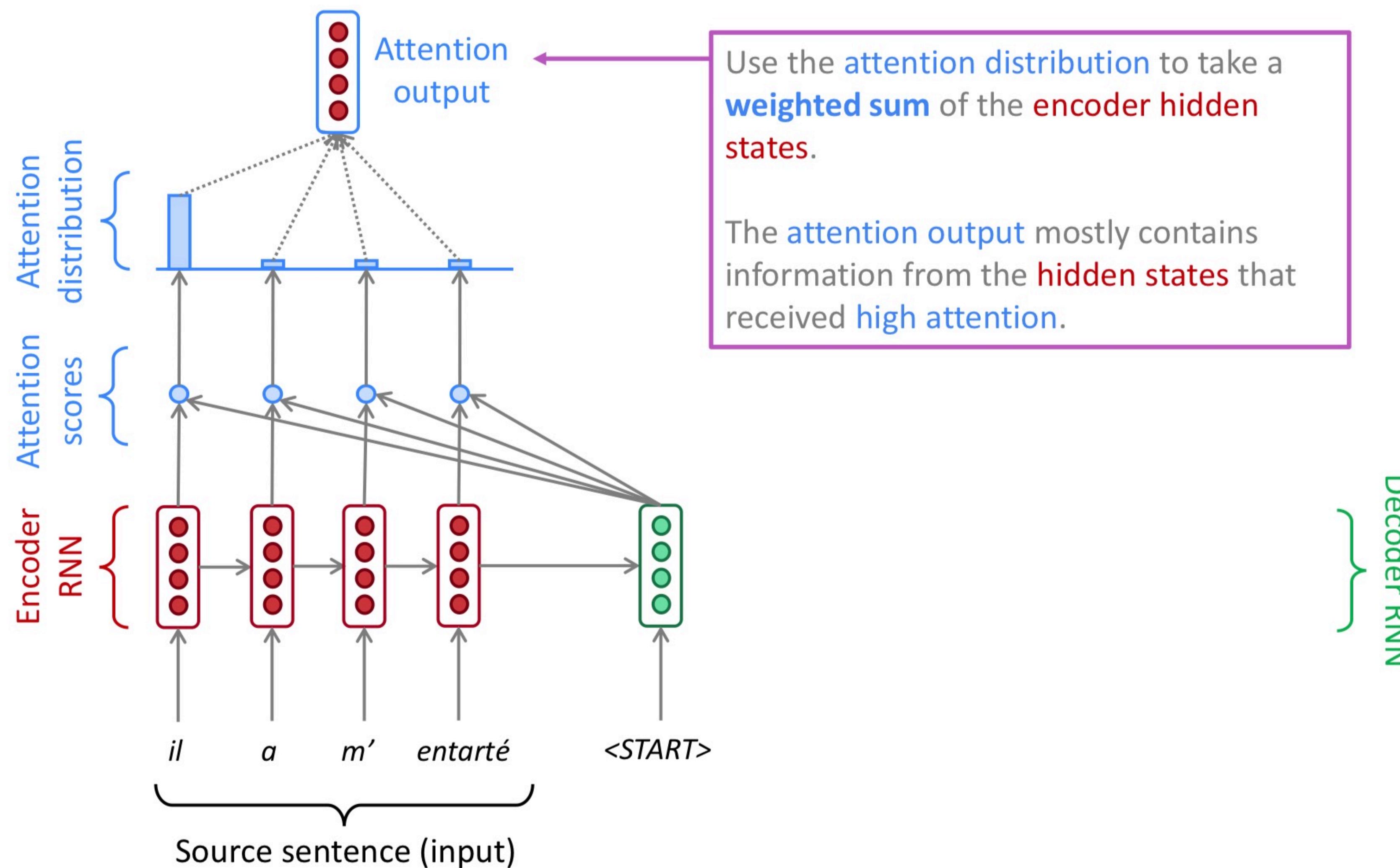
self-attention



Multiple Heads



# Recall: attention



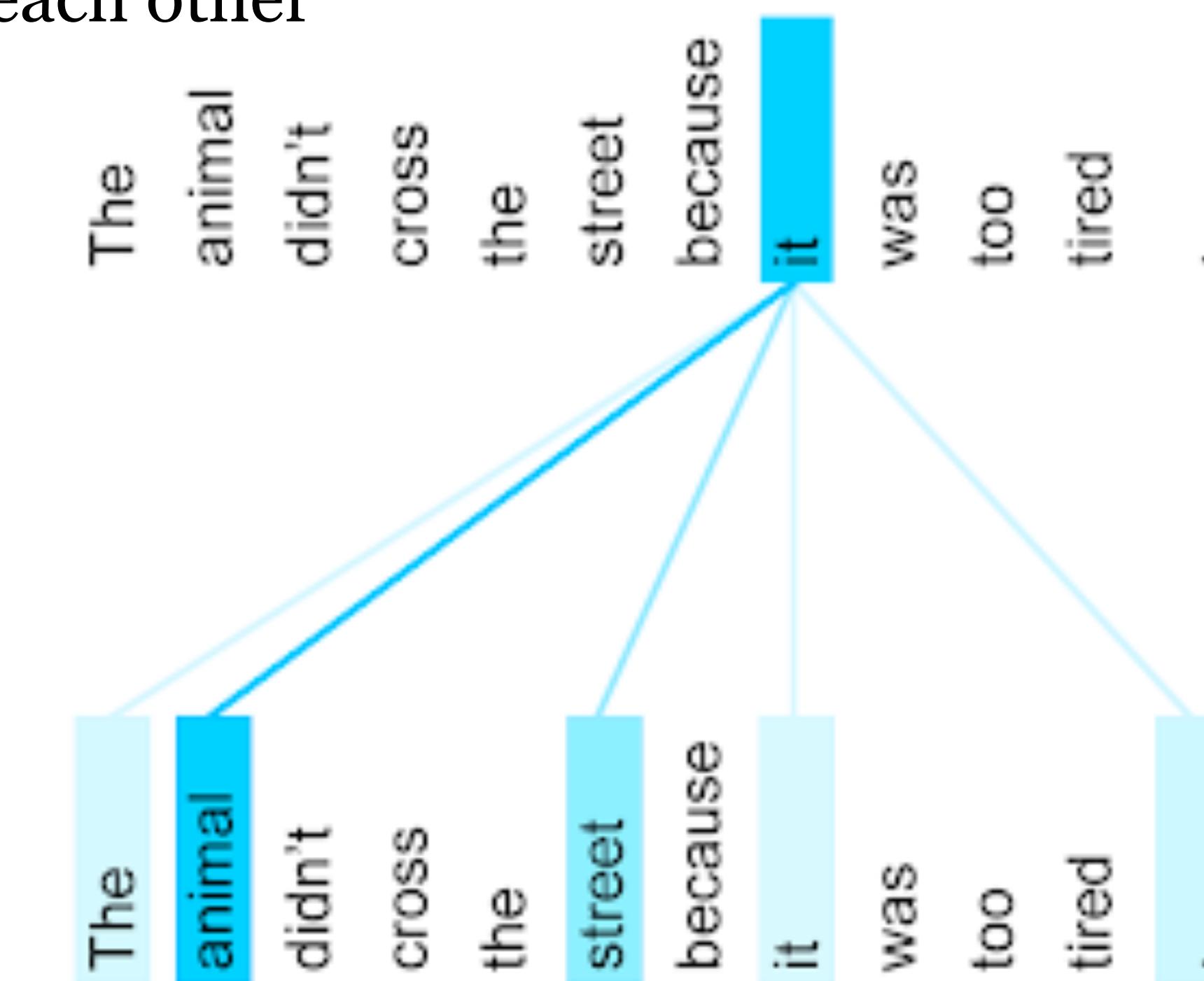
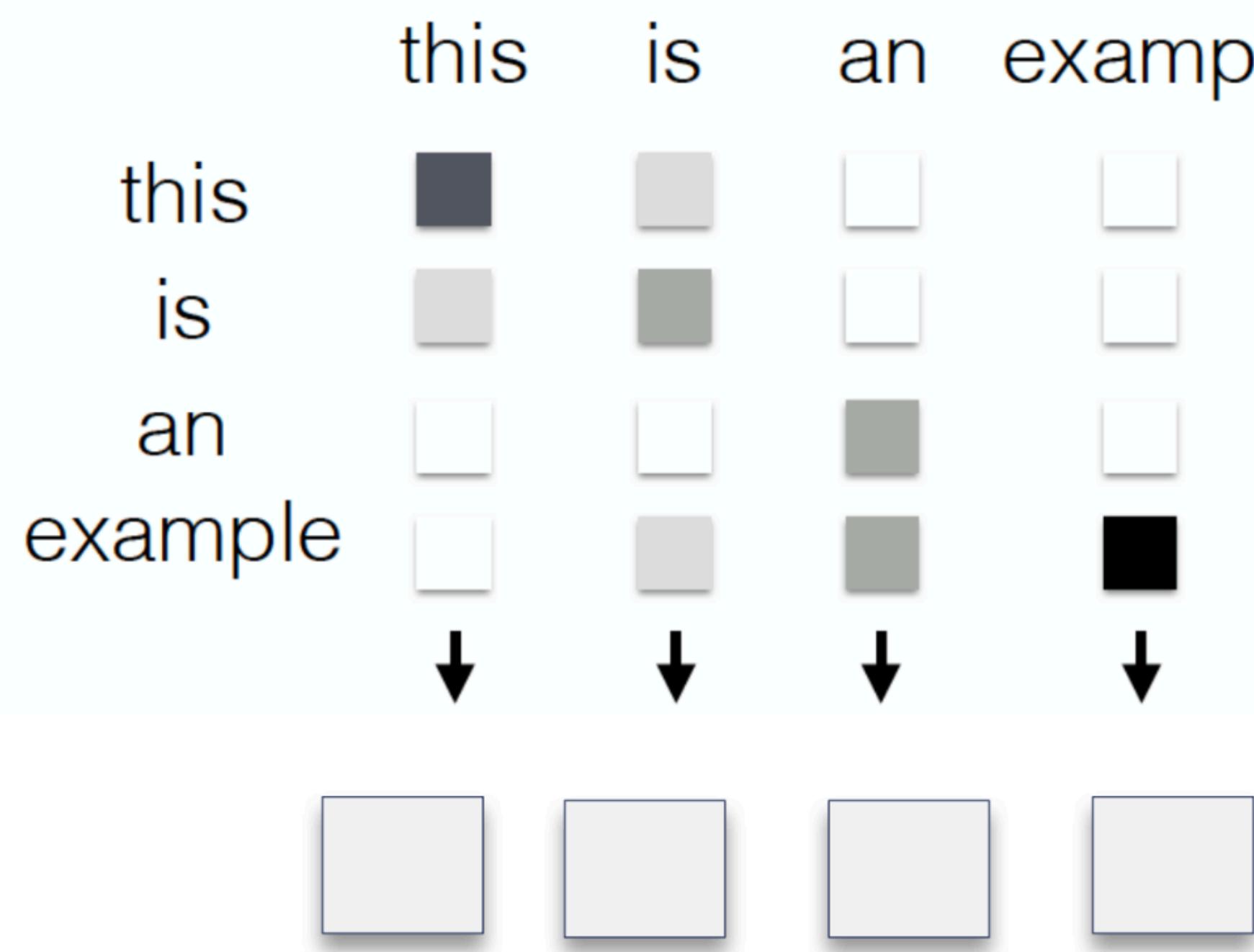
(slide credit: Abigail See)

# Self Attention

(also referred to as Intra-Attention)

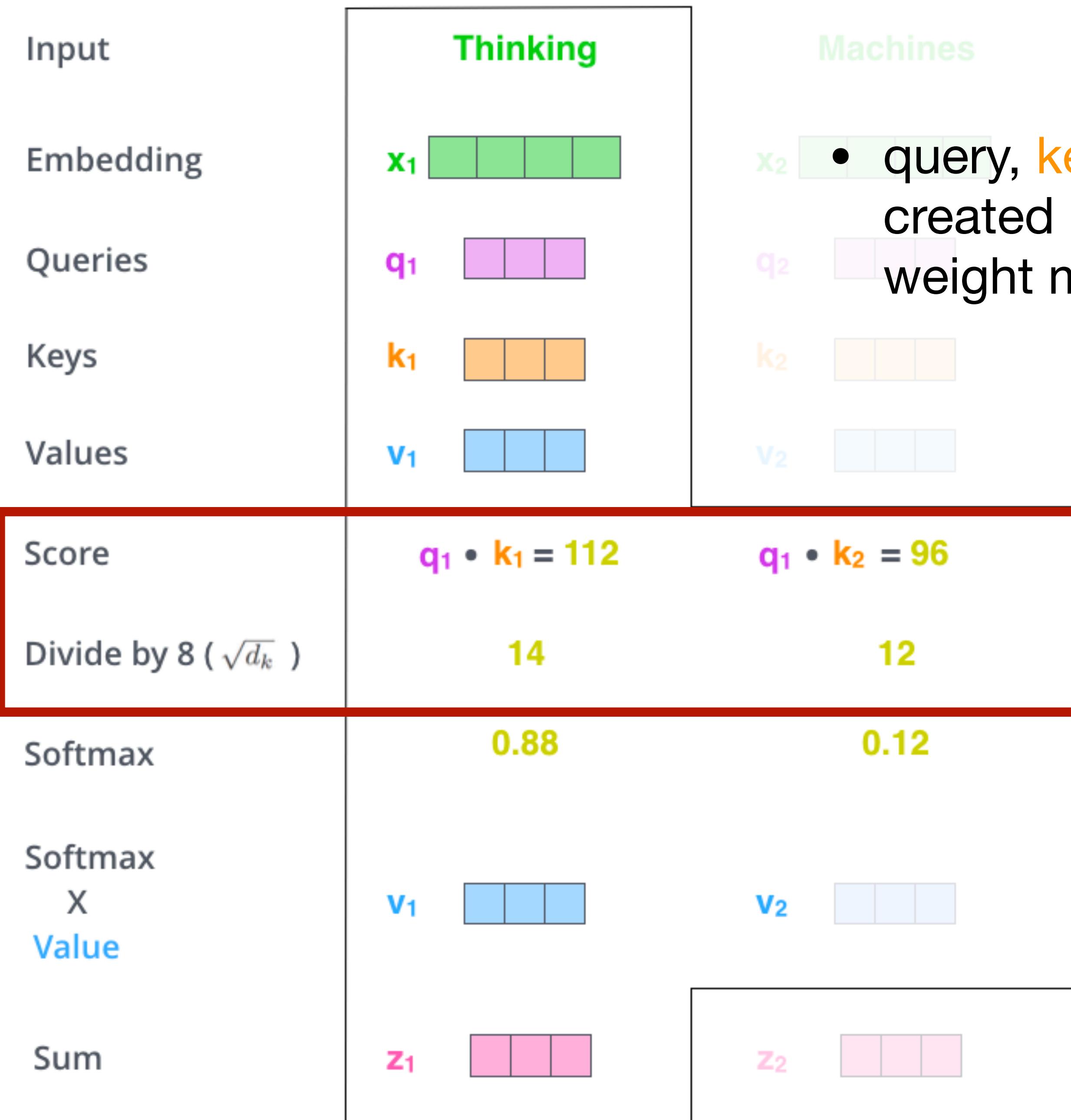
- **Self-attention:** let's use each word as query and compute the attention with all the other words

= the word vectors themselves select each other



# General definition of attention

- ▶ **Attention** is a general concept
  - ▶ Given **query q** and a set of **key-value** pairs (**K,V**)
  - ▶ **Attention** is a way to compute a **weighted sum** of the **values** dependent on the **query** and the corresponding **keys**.
  - ▶ Query determines what values to focus on, the **query** “attends” to the **values**
  - ▶ All of these (**key value query**) are represented using **vectors**
    - ▶ These vectors are created by multiplying embedding by trained weight matrices.



- query, **key**, and **value** vectors created by multiplying learned weight matrices with embedding

- Can be any kind of attention function
- For transformers, this is the **scaled dot-product attention**

(figure credit: [Jay Alammar](#)  
<http://jalammar.github.io/illustrated-transformer/>)

# Recall: types of attention

- ▶ Assume encoder hidden states  $h_1, h_2, \dots, h_n$  and decoder hidden state  $z$

1. **Dot-product attention** (assumes equal dimensions for  $a$  and  $b$ ):

$$g(h_i, z) = z^T h_i \in \mathbb{R}$$

Simplest (no extra parameters)  
requires  $z$  and  $h_i$  to be same size

2. **Bilinear / multiplicative attention:**

$$g(h_i, z) = z^T W h_i \in \mathbb{R}, \text{ where } W \text{ is a weight matrix}$$

More flexible  
than dot-product  
( $W$  is trainable)

3. **Additive attention (essentially MLP):**

$$g(h_i, z) = v^T \tanh(W_1 h_i + W_2 z) \in \mathbb{R}$$

where  $W_1, W_2$  are weight matrices and  $v$  is a weight vector

Perform better for  
larger dimensions

more efficient  
(matrix  
multiplication)

# Scaled dot-product attention

- ▶ Assume encoder hidden states  $h_1, h_2, \dots, h_n$  and decoder hidden state  $z$

1. **Dot-product attention** (assumes equal dimensions for  $a$  and  $b$ ):

$$g(h_i, z) = z^T h_i \in \mathbb{R}$$

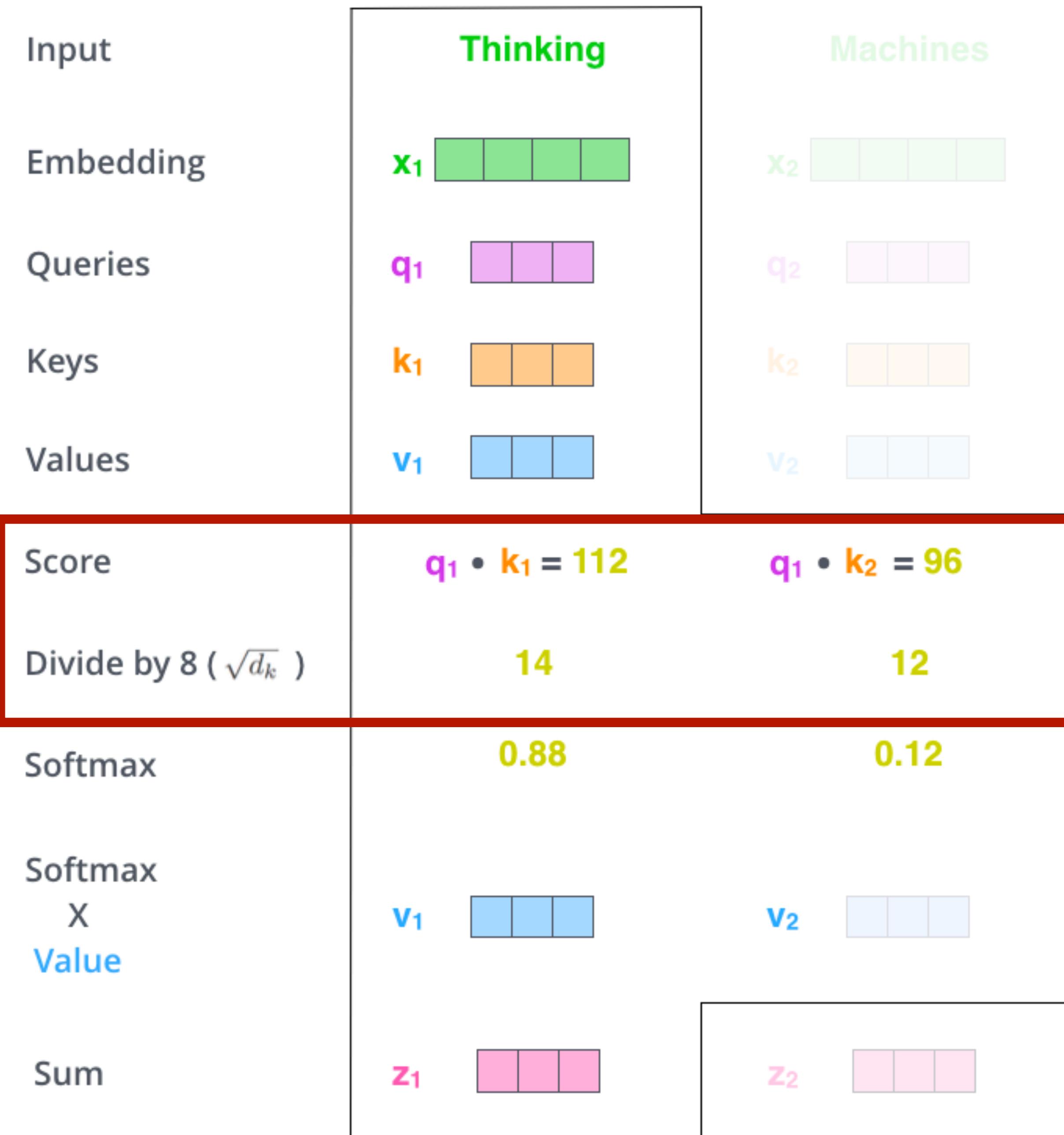
Perform poorly for large  $d$   
Softmax has small gradient

2. **Scaled dot-product attention:**

$$g(h_i, z) = \frac{z^T h_i}{\sqrt{d}} \in \mathbb{R}$$

Maybe will perform well  
for larger dimensions

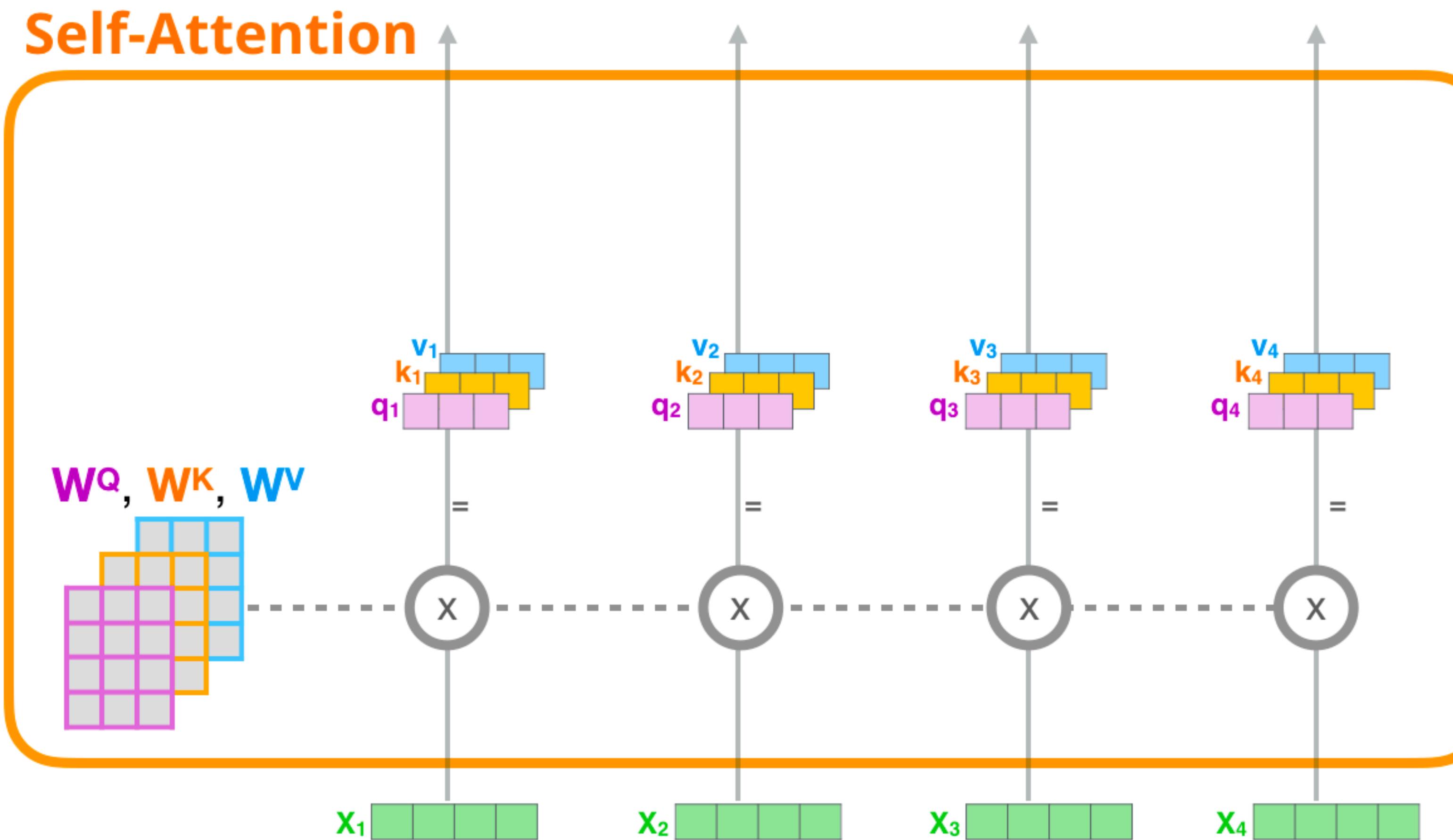
Scaling factor:  $d$  = dimension of hidden state



- Can be any kind of attention function
- For transformers, this is the **scaled dot-product attention**
- $z_1$  is the final vector of attended values for “Thinking” as the query

(figure credit: [Jay Alammar](#)  
<http://jalammar.github.io/illustrated-transformer/>)

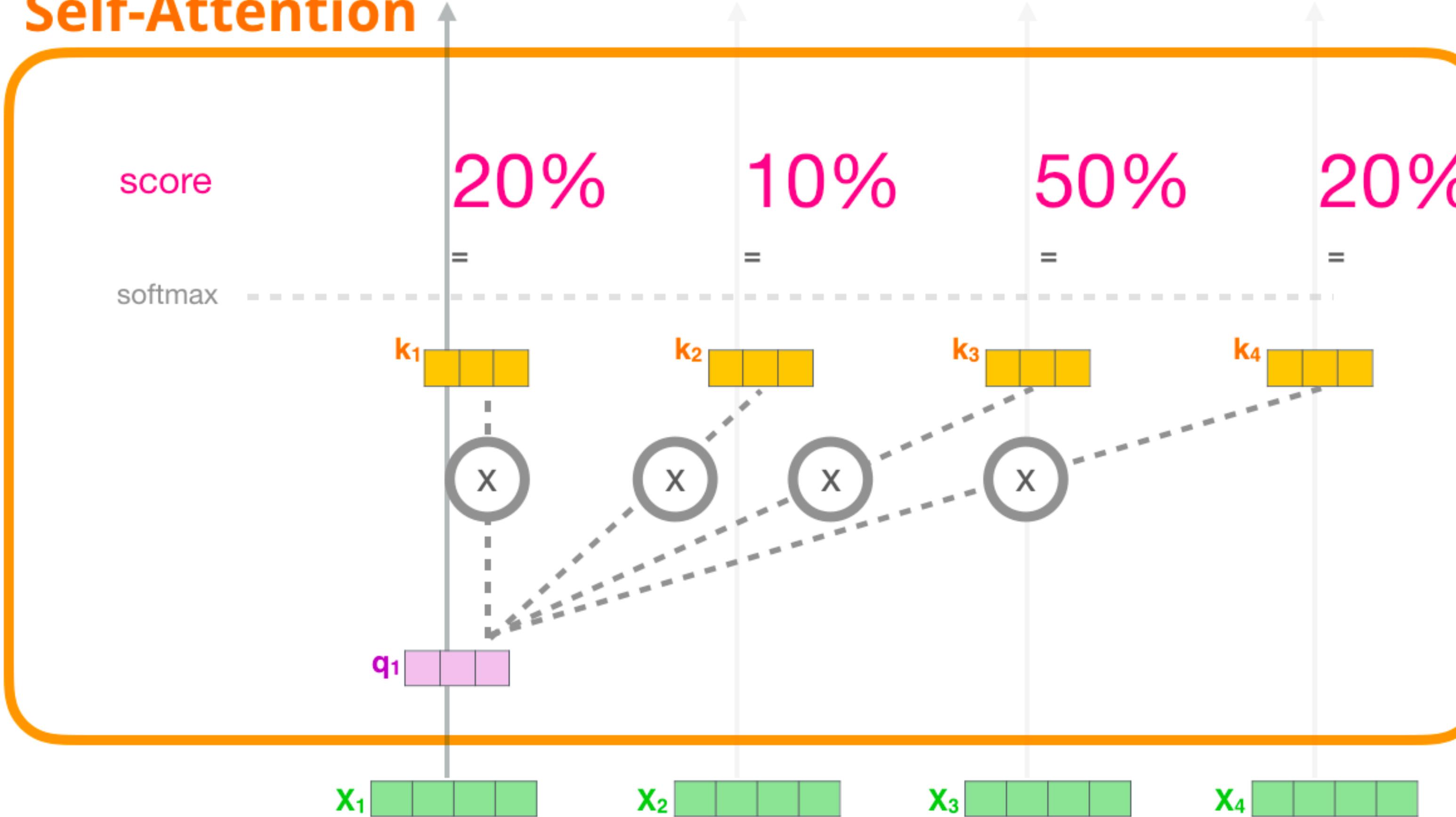
1) For each input token, create a **query vector**, a **key vector**, and a **value vector** by multiplying by weight Matrices  $\mathbf{W}^Q$ ,  $\mathbf{W}^K$ ,  $\mathbf{W}^V$



(figure credit: [Jay Alammar](#)  
<http://jalammar.github.io/illustrated-gpt2/>)

2) Multiply (dot product) the current **query vector**, by all the **key vectors**, to get a score of how well they match

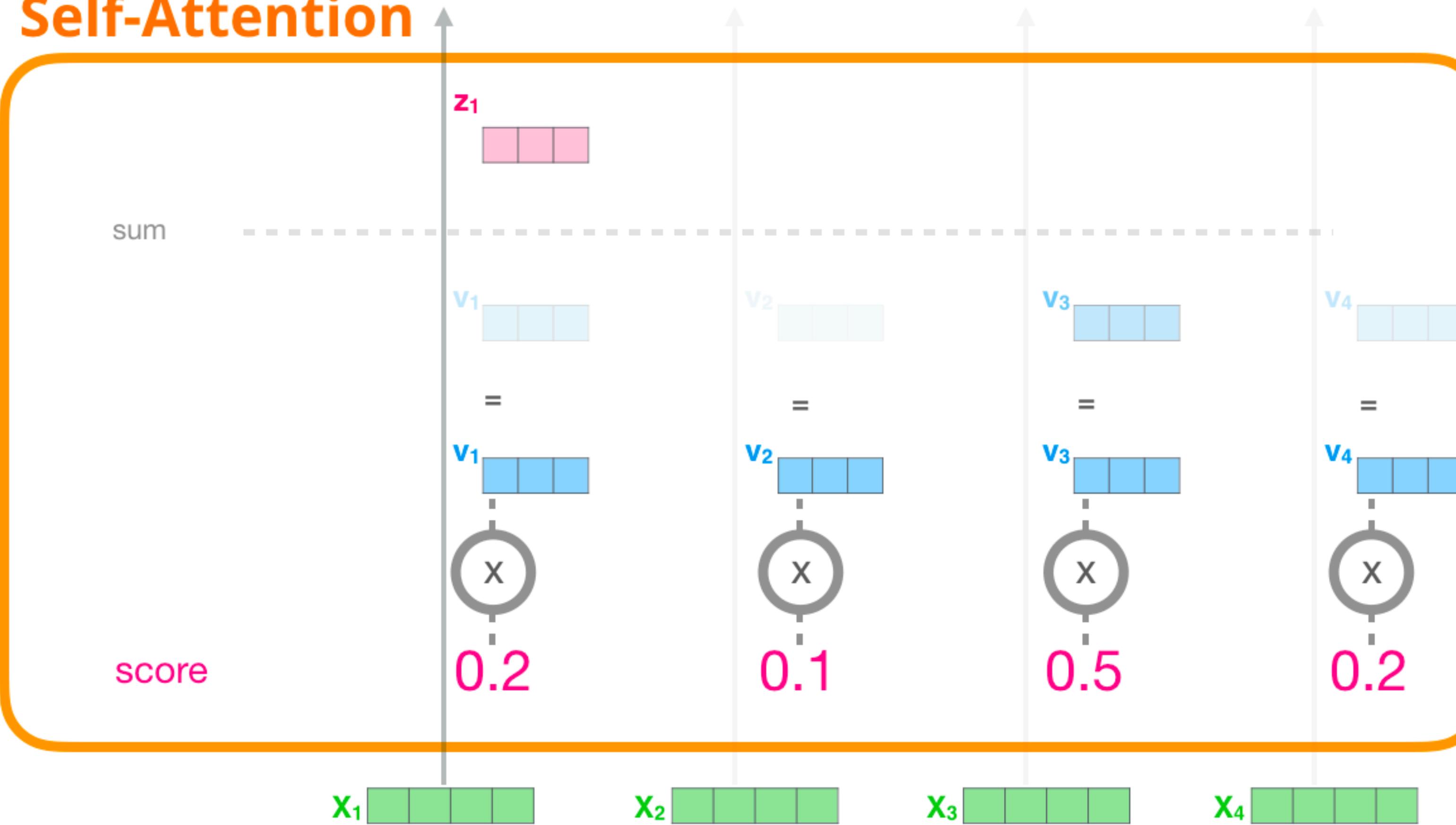
## Self-Attention



(figure credit: [Jay Alammar](#)  
<http://jalammar.github.io/illustrated-gpt2/>)

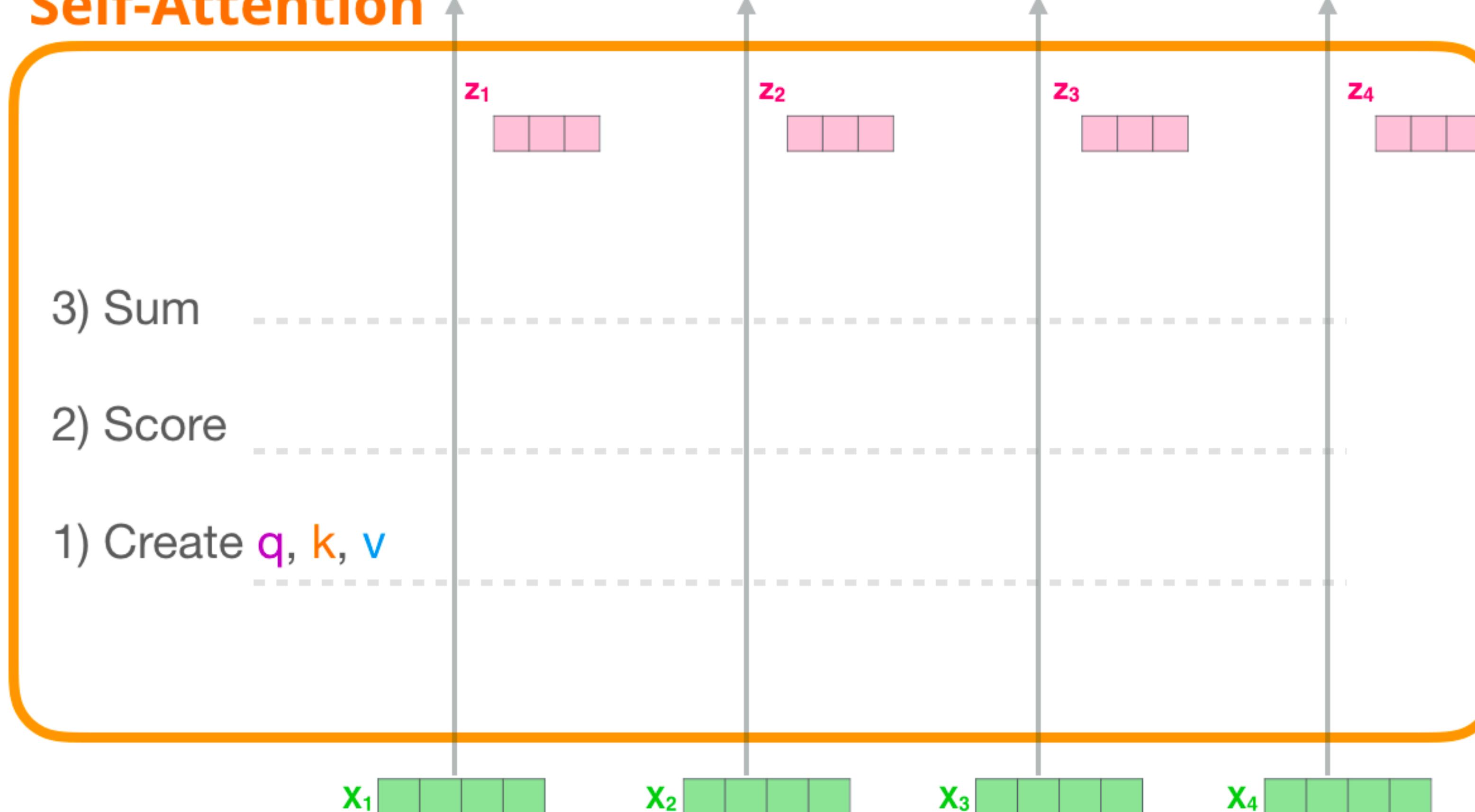
3) Multiply the **value vectors** by the **scores**, then sum up

## Self-Attention



(figure credit: [Jay Alammar](#)  
<http://jalammar.github.io/illustrated-gpt2/>)

## Self-Attention



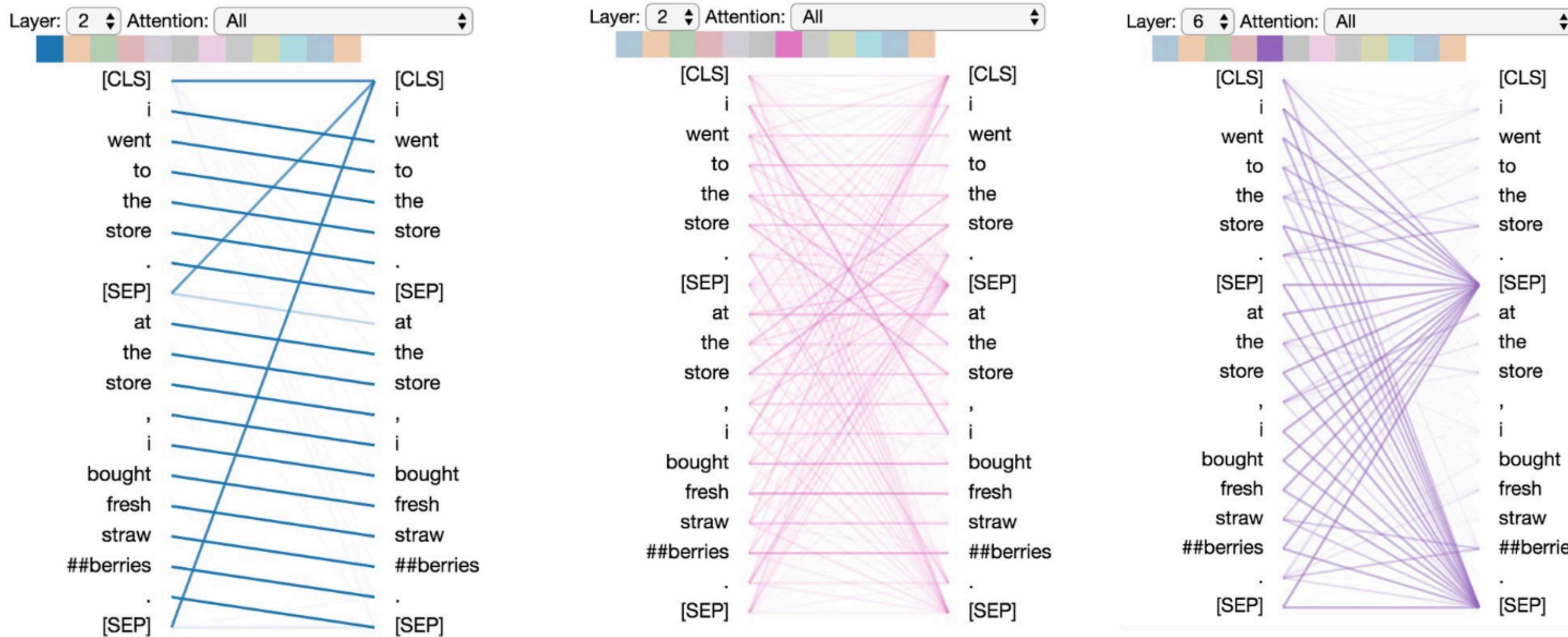
(figure credit: [Jay Alammar](#)  
<http://jalammar.github.io/illustrated-gpt2/>)

# Multi-head self-attention

One head is not expressive enough. Let's have multiple heads!

$$A(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$
$$\text{head}_i = A(XW_i^Q, XW_i^K, XW_i^V)$$

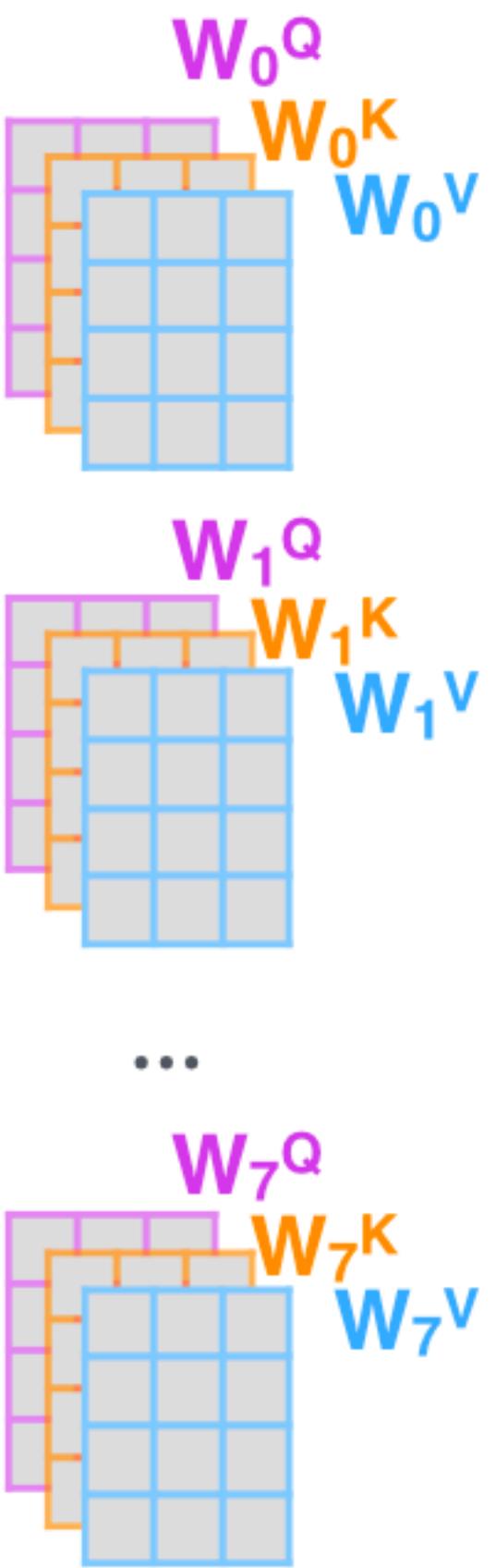
In practice,  $h = 8$ ,  
 $d = d_{out}/h$ ,  $W^O \in \mathbb{R}^{d_{out} \times d_{out}}$



<https://github.com/jessevig/bertviz>

# Multiple heads

- Multiple (different) representations for each **query**, **key**, and **values**
- Different weight matrices → different vectors
- Different ways for the words to interact with each other



(figure credit: [Jay Alammar](#)  
<http://jalammar.github.io/illustrated-transformer/>)

# Summary: Multi-head Self Attention

- **Attention:** a query  $q$  and a set of key-value  $(k_i, v_i)$  pairs to an output
- Dot-product attention:

$$A(q, K, V) = \sum_i \frac{e^{q \cdot k_i}}{\sum_j e^{q \cdot k_j}} v_i$$
$$K, V \in \mathbb{R}^{n \times d}, q \in \mathbb{R}^d$$

- If we have multiple queries:

$$A(Q, K, V) = \text{softmax}(QK^\top)V$$
$$Q \in \mathbb{R}^{n_Q \times d}, K, V \in \mathbb{R}^{n \times d}$$

- **Self-attention:** let's use each word as query and compute the attention with all the other words
  - = the word vectors themselves select each other

# Summary: Multi-head Self Attention

- Scaled Dot-Product Attention:

$$A(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right)V$$

- Input:  $X \in \mathbb{R}^{n \times d_{in}}$

$$A(XW^Q, XW^K, XW^V) \in \mathbb{R}^{n \times d}$$

$$W^Q, W^K, W^V \in \mathbb{R}^{d_{in} \times d}$$

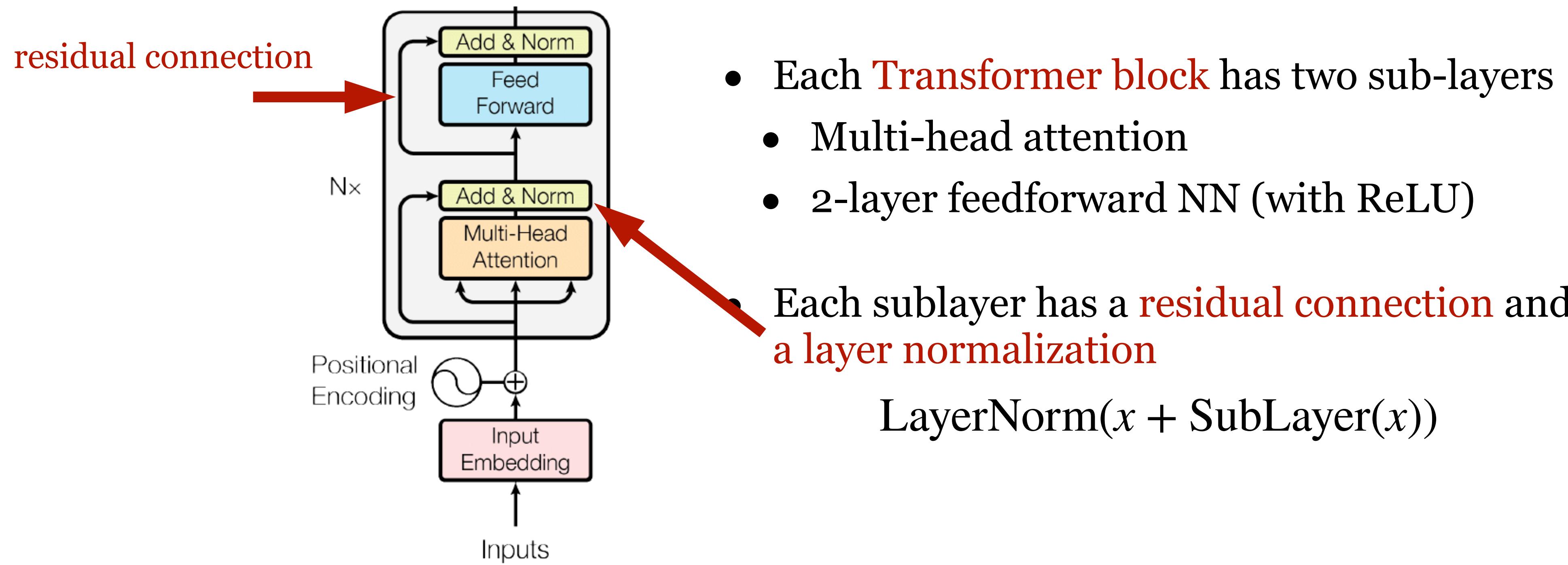
- Multi-head attention: using more than one head is always useful..

$$A(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

$$\text{head}_i = A(XW_i^Q, XW_i^K, XW_i^V)$$

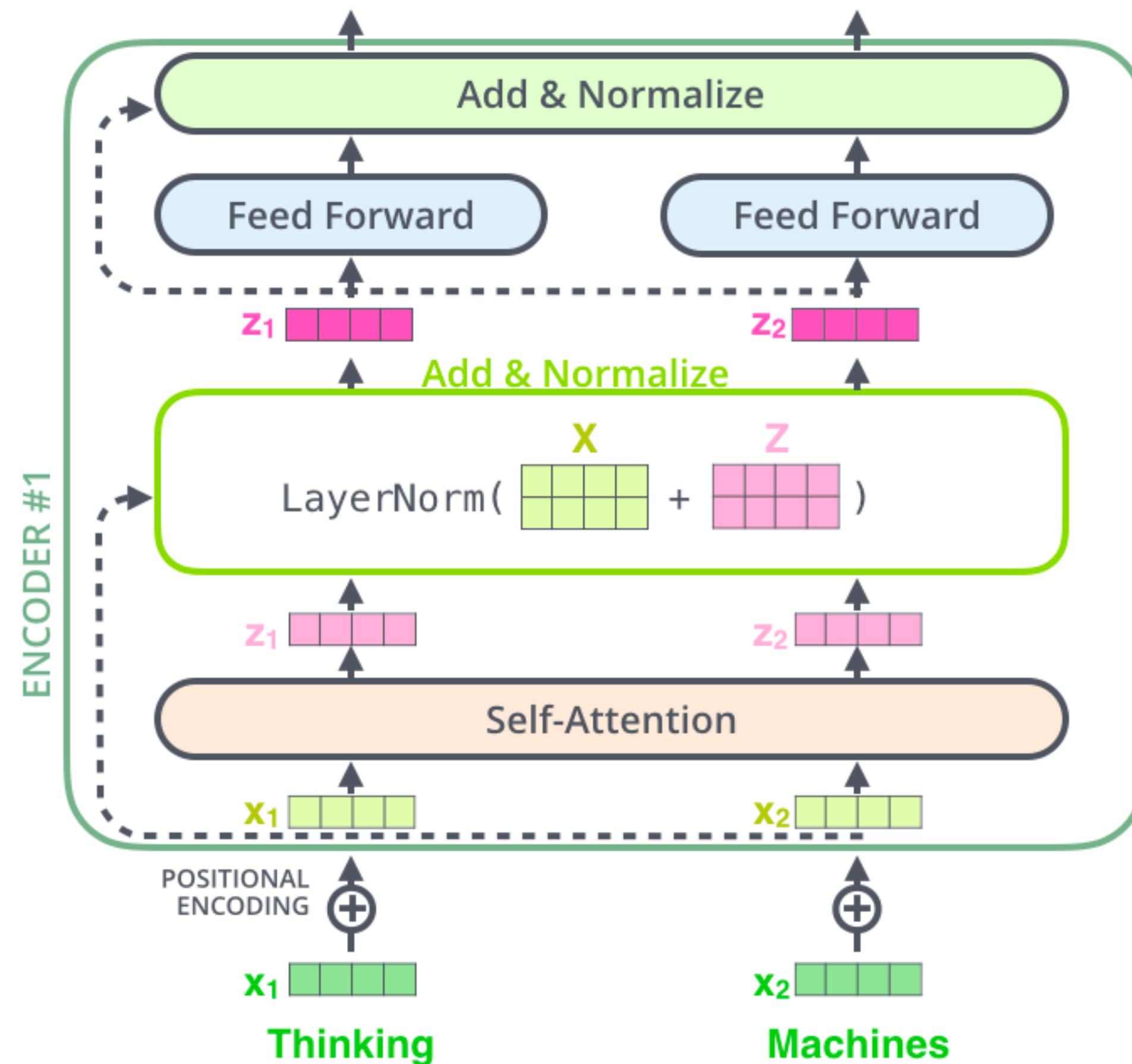
In practice,  $h = 8$ ,  $d = d_{out}/h$ ,  $W^O = d_{out} \times d_{out}$

# Putting it all together



(Ba et al, 2016): Layer Normalization

# Residual connections and Layer Normalization



## LayerNorm

- changes input features to have mean 0 and variance 1 per layer.
- Adds two more parameters

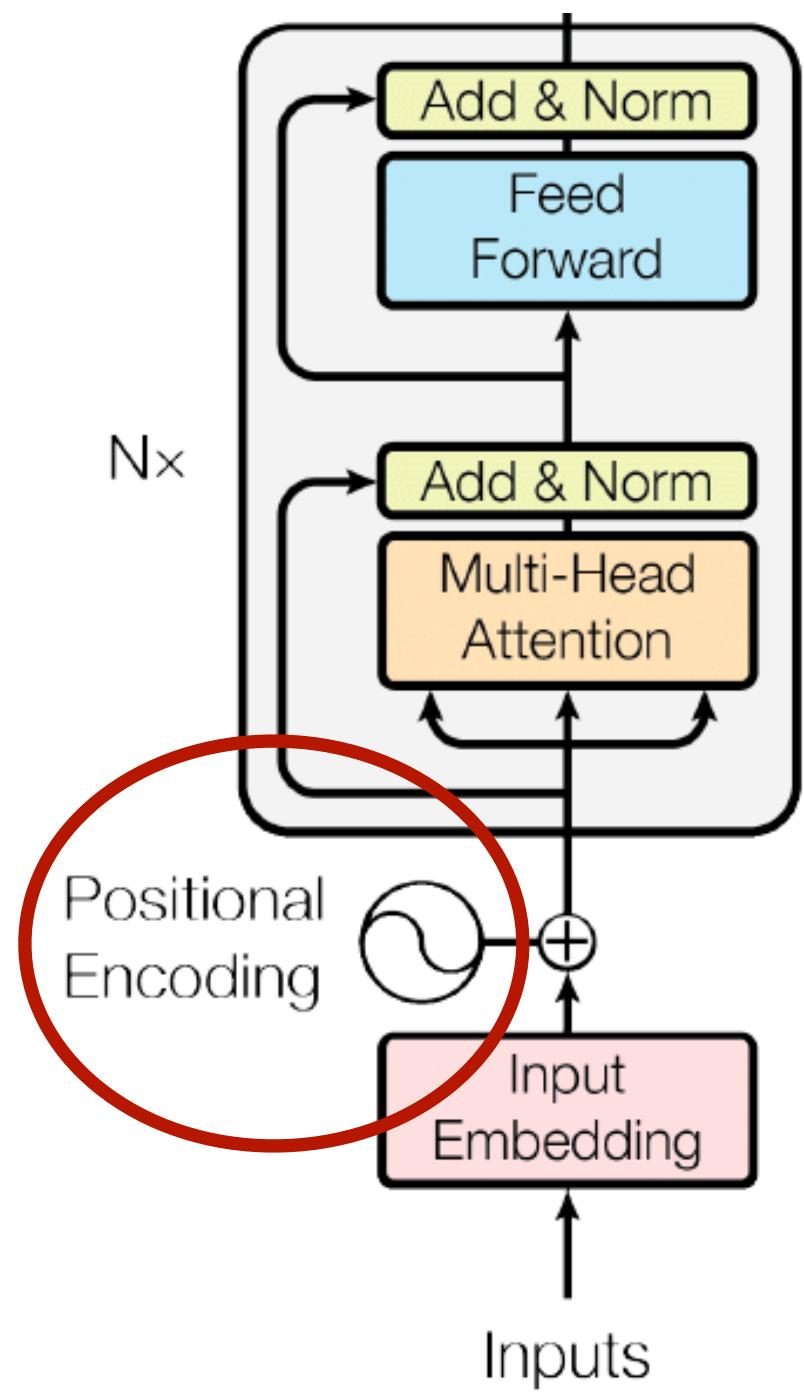
$$\mu^l = \frac{1}{H} \sum_{i=1}^H a_i^l \quad \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2}$$

$$h_i = \frac{g_i}{\sigma_i} (a_i - \mu_i) + b_i$$

(figure credit: [Jay Alammar](#)  
<http://jalammar.github.io/illustrated-transformer/>)  
65

(Ba et al, 2016): Layer Normalization

# Putting it all together



- Each Transformer block has two sub-layers
  - Multi-head attention
  - 2-layer feedforward NN (with ReLU)
- Each sublayer has a residual connection and a layer normalization  
$$\text{LayerNorm}(x + \text{SubLayer}(x))$$
- Input layer has a **positional encoding**

(Ba et al, 2016): Layer Normalization

# Positional encoding

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

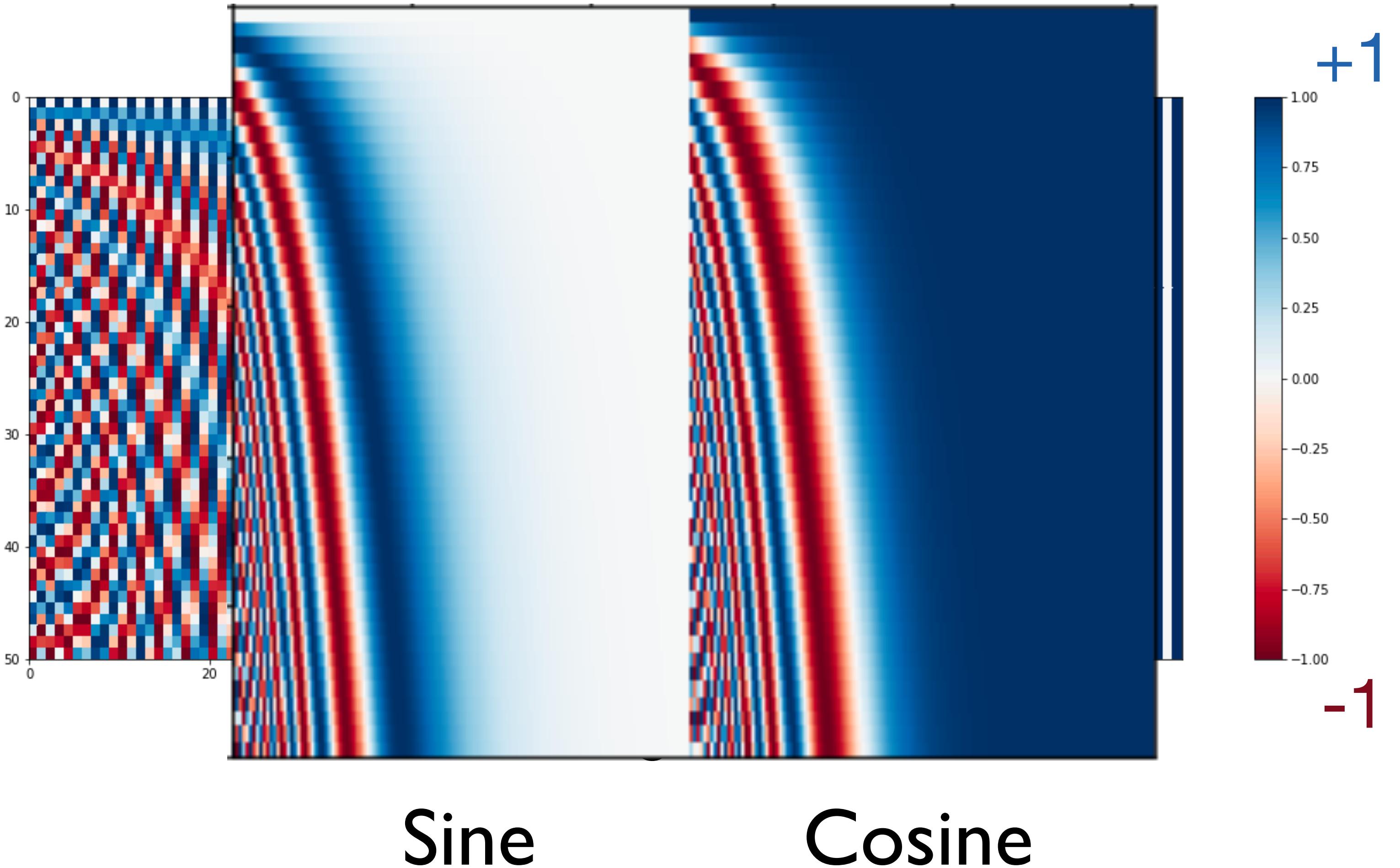
$$\omega_k = \frac{1}{10000^{2k/d}}$$

$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

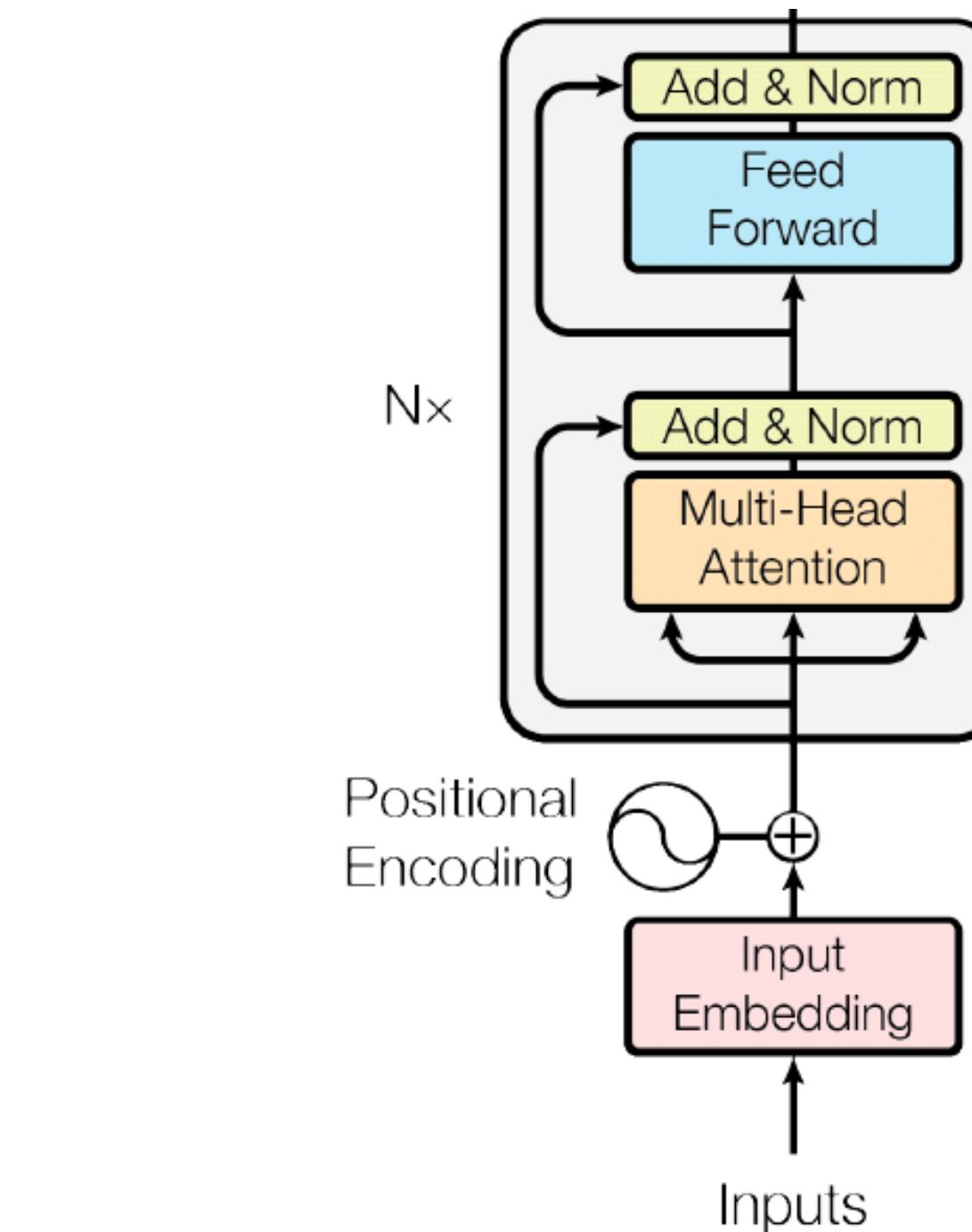
$t$  = position

$d$  = embedding dimension

$i$  = embedding index (0 to  $d-1$ )

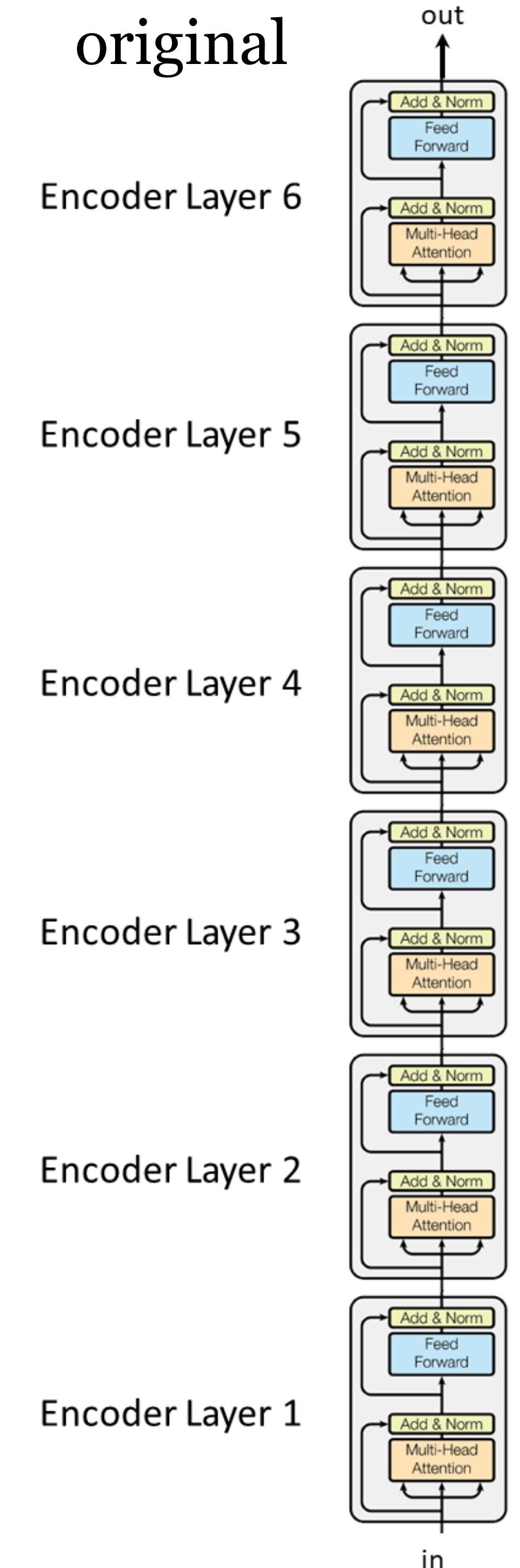


# Putting it all together



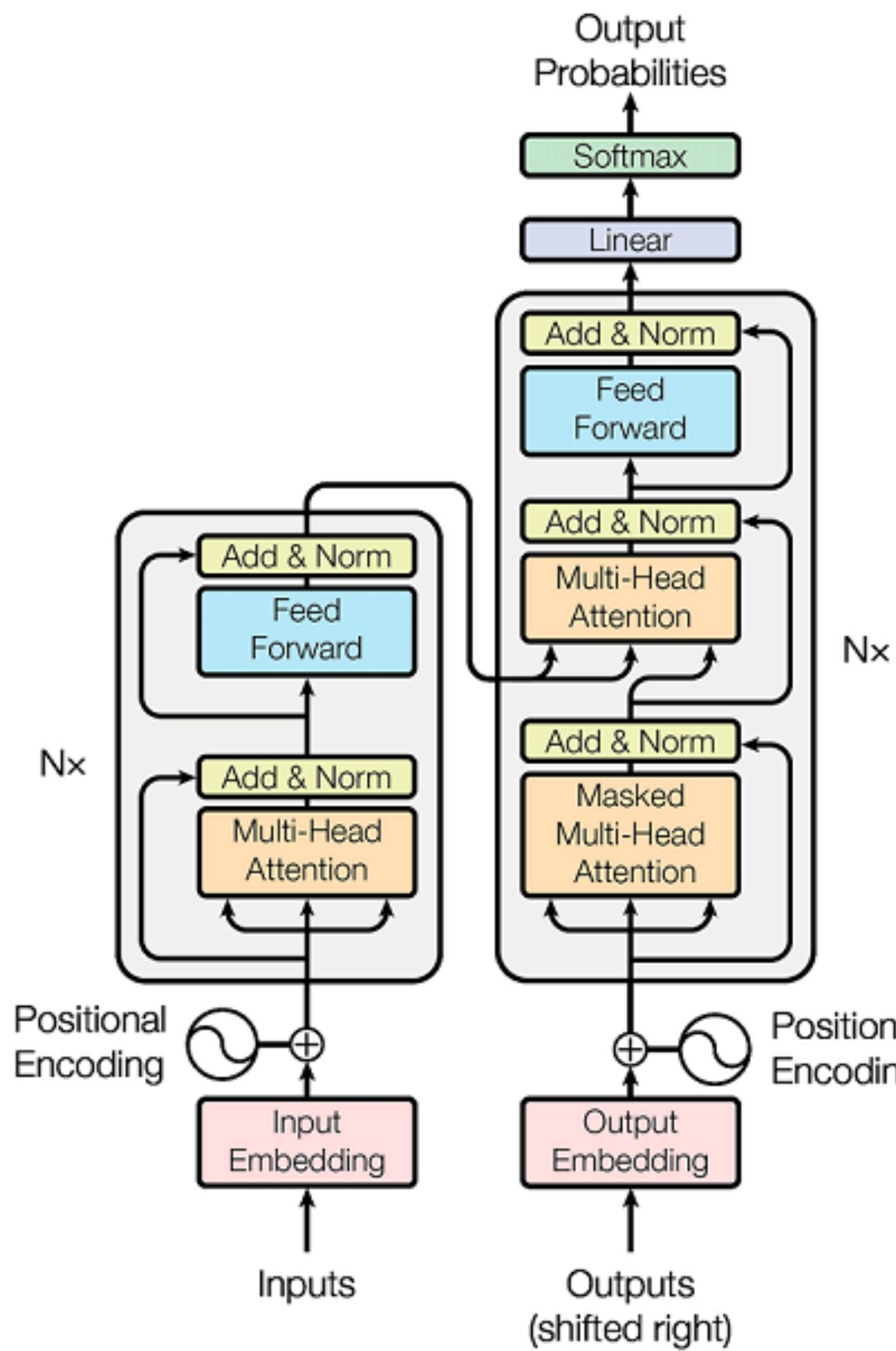
- The diagram illustrates the architecture of a Transformer block. It starts with **Inputs**, which are processed by an **Input Embedding** layer. This is followed by a stack of  $N_x$  layers. Each layer consists of the following sequence: **Multi-Head Attention** (orange box), **Add & Norm** (green box), and a **Feed Forward** (light blue box). A residual connection (sum of input and output of the sublayer) is shown with a circle and a plus sign (+). Layer normalization (LayerNorm) is indicated by a circle with a cross. The final output is the sum of the input embedding and the output of the  $N_x$ -th layer.

  - Each Transformer block has two sub-layers
    - Multi-head attention
    - 2-layer feedforward NN (with ReLU)
  - Each sublayer has a residual connection and a layer normalization
$$\text{LayerNorm}(x + \text{SubLayer}(x))$$
  - Input layer has a positional encoding
  - Input embedding is byte pair encoding (BPE)
  - BERT\_base: 12 layers, 12 heads, hidden size = 768, 110M parameters
  - BERT\_large: 24 layers, 16 heads, hidden size = 1024, 340M parameters

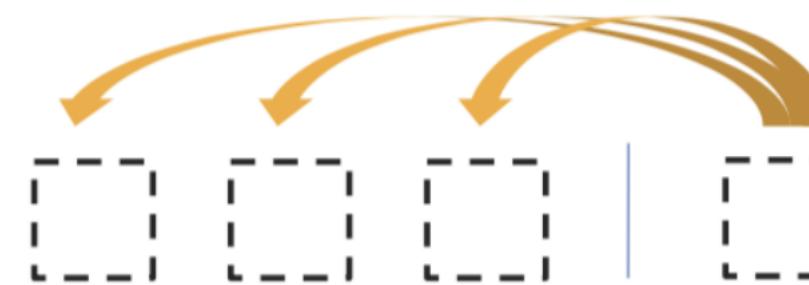


(Ba et al, 2016): Layer Normalization

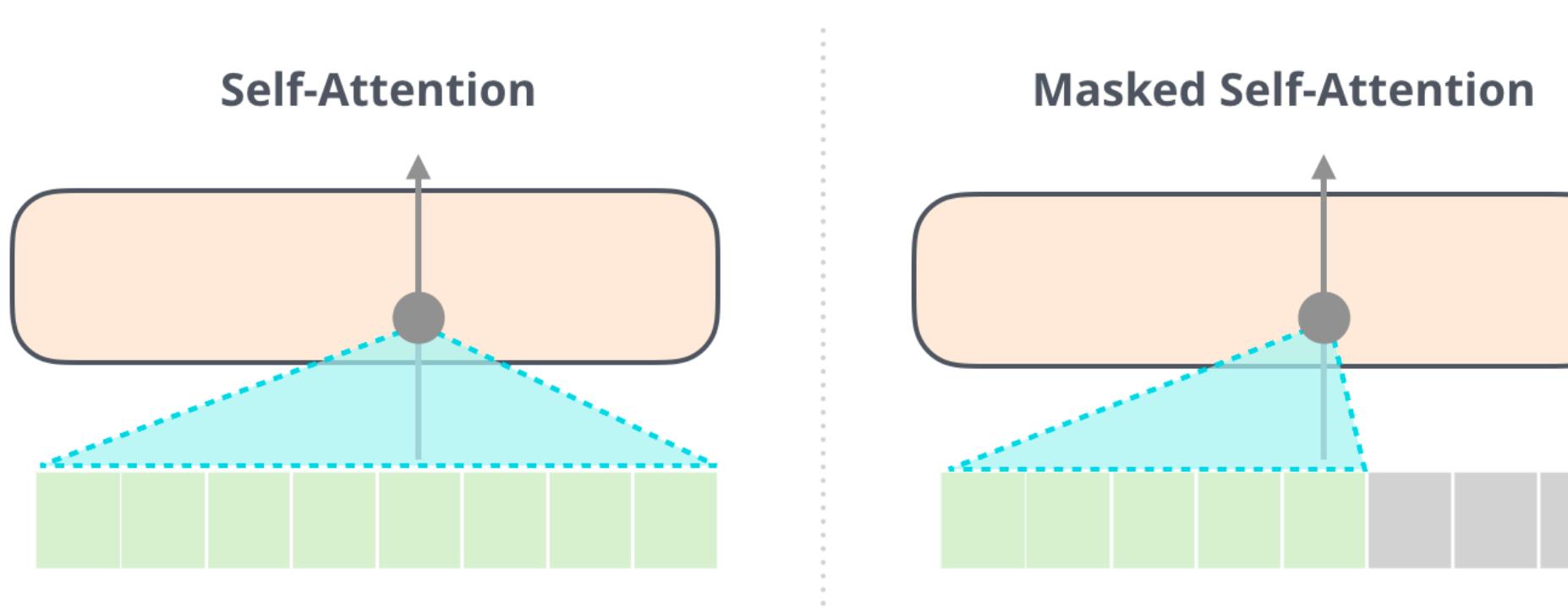
# Transformer decoder



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



- Masked decoder self-attention on previously generated outputs



- also 6 layers (in original paper)

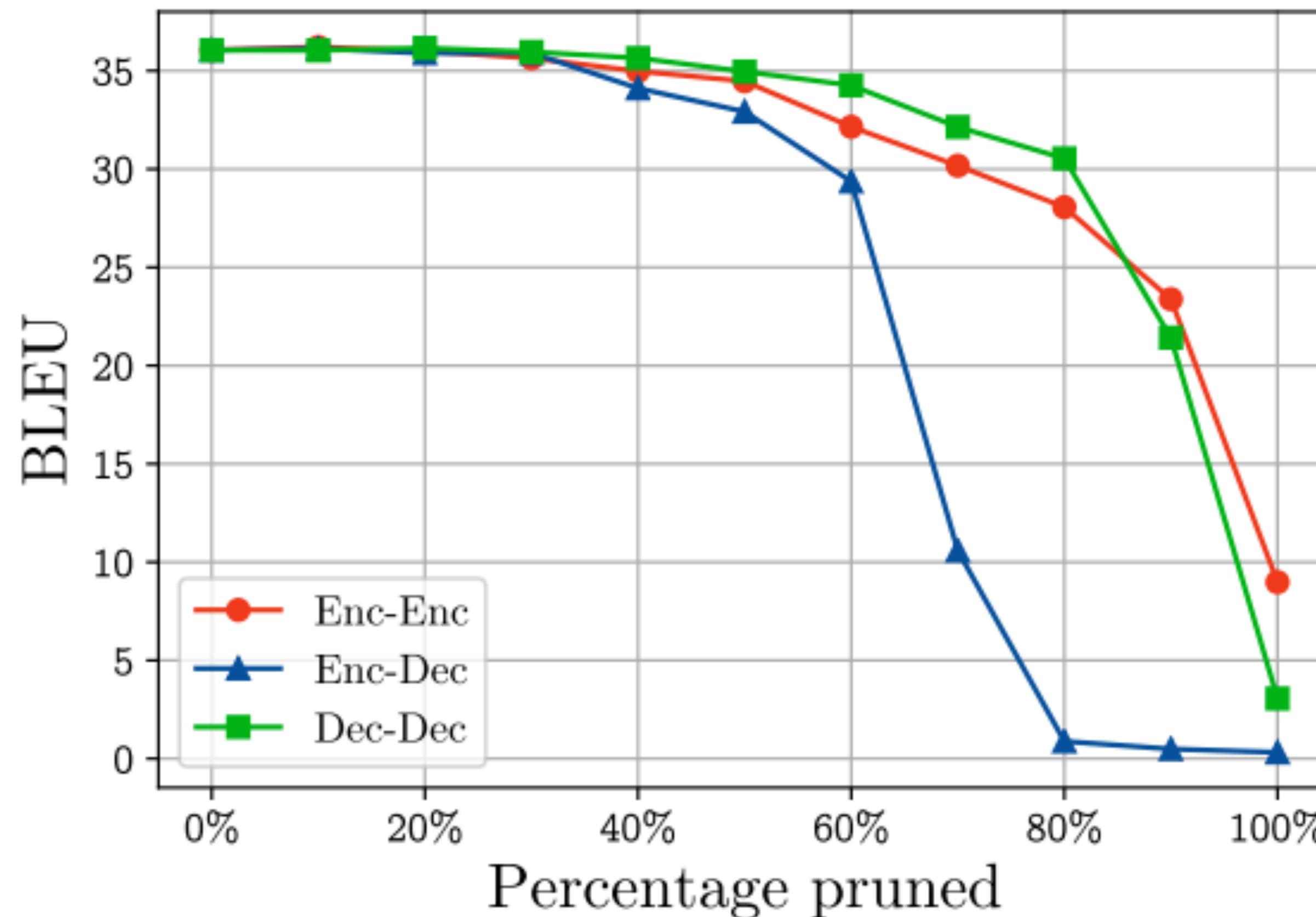
(figure credit: [Jay Alammar](#)  
<http://jalammar.github.io/illustrated-gpt2/>)

# Do we need all these heads?

3 types of attention: Enc-Enc, Enc-Dec, Dec-Dec

6 layers, 16 heads each layer for each type

- Can we prune away some of the heads of a trained model during test time?



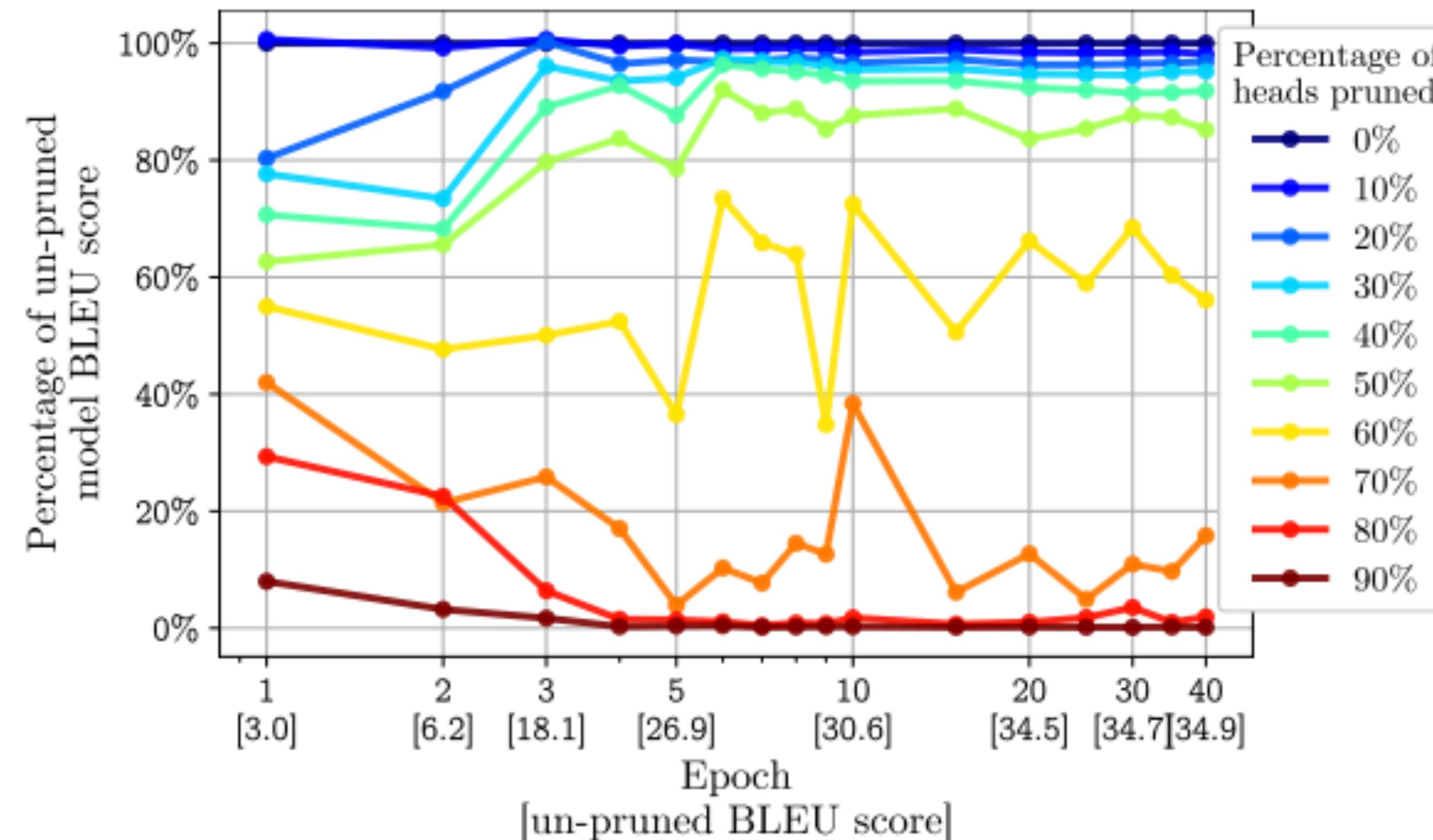
*Are Sixteen Heads Really Better than One?*  
Michel, Levy, and Neubig, NeurIPS 2019

# Do we need all these heads?

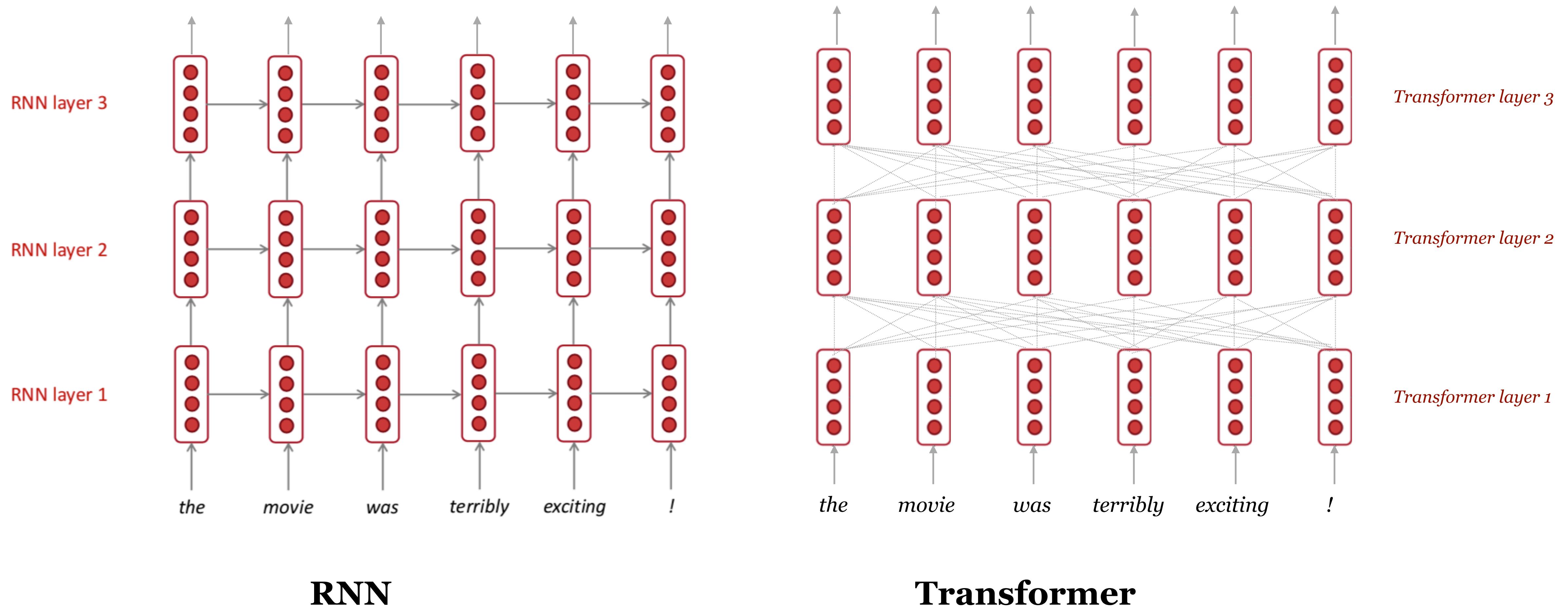
3 types of attention: Enc-Enc, Enc-Dec, Dec-Dec

6 layers, 16 heads each layer for each type

- Can we train a good MT model with less heads?



# RNNs vs Transformers



# Useful Resources

## nn.Transformer:

```
>>> transformer_model = nn.Transformer(nhead=16, num_encoder_layers=12)
>>> src = torch.rand((10, 32, 512))
>>> tgt = torch.rand((20, 32, 512))
>>> out = transformer_model(src, tgt)
```

## nn.TransformerEncoder:

```
>>> encoder_layer = nn.TransformerEncoderLayer(d_model=512, nhead=8)
>>> transformer_encoder = nn.TransformerEncoder(encoder_layer, num_layers=6)
>>> src = torch.rand(10, 32, 512)
>>> out = transformer_encoder(src)
```

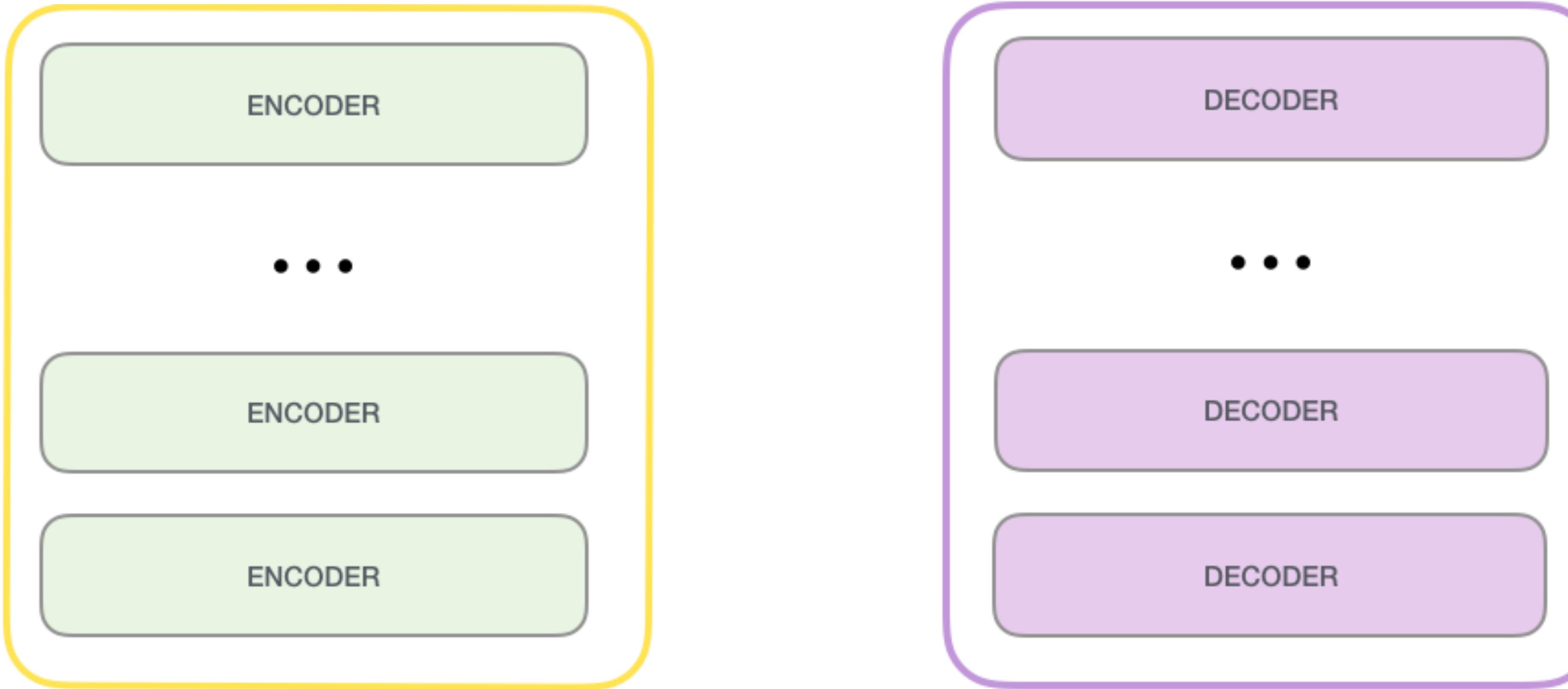
## The Annotated Transformer:

<http://nlp.seas.harvard.edu/2018/04/03/attention.html>

A Jupyter notebook which explains how Transformer works line by line in PyTorch!

# Transformers blocks as building blocks

- BERT (built on Transformer encoders)
- GPT-2 (built on Transformer decoders)



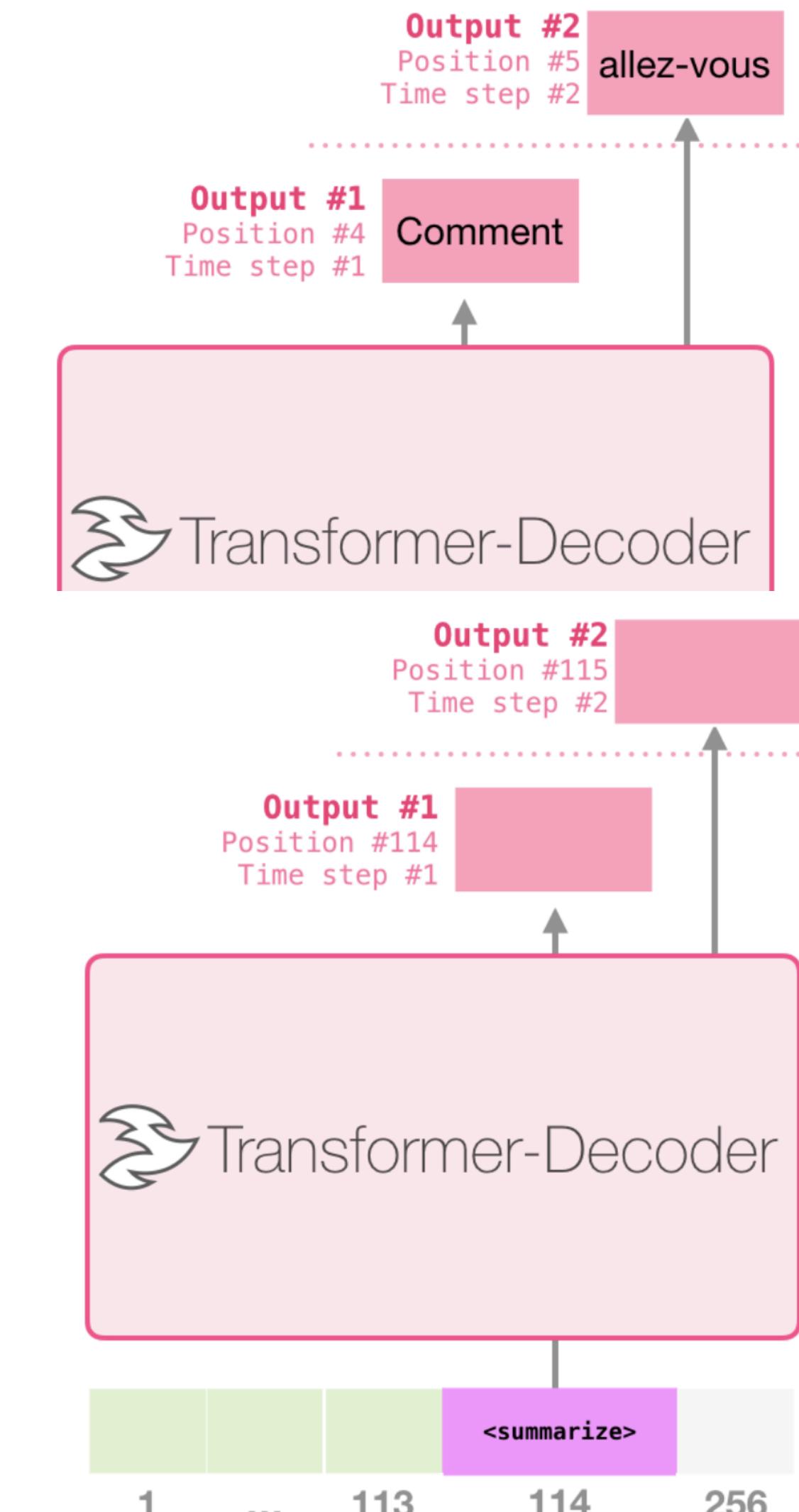
(figure credit: Jay Alammar  
<http://jalammar.github.io/illustrated-gpt2/>)

# GPT-2

How can we use decoders for different tasks?

- Machine Translation

I	am	a	student	<to-fr>	je	suis	étudiant
let	them	eat	cake	<to-fr>	Qu'ils	mangent	de
good	morning	<to-fr>	Bonjour				

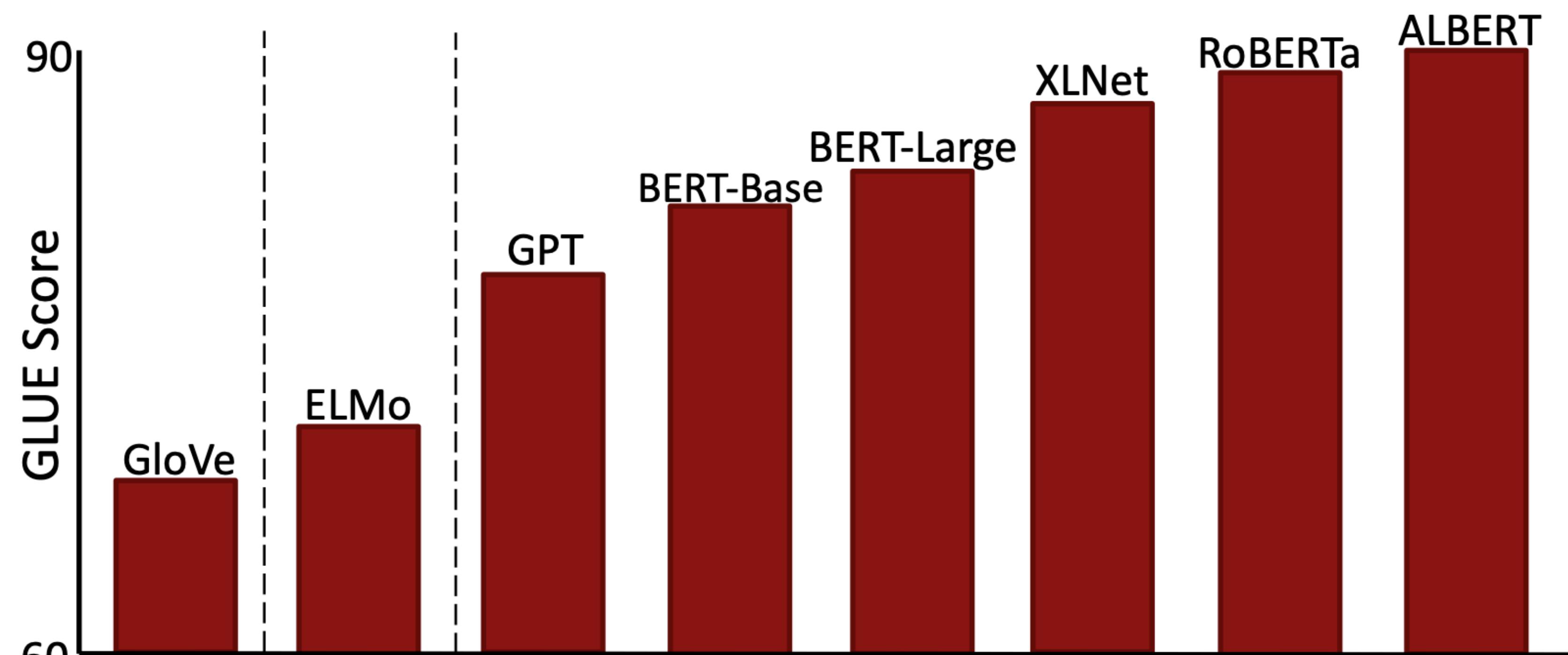


- Summarization

Article #1 tokens	<summarize>	Article #1 Summary
Article #2 tokens	<summarize>	Article #2 Summary padding
Article #3 tokens	<summarize>	Article #3 Summary

edit: [Jay Alammar](#)  
<https://github.com/jayalammar/illustrated-gpt2/>

# NLP progress so far



Over 3x reduction in error in 2 years, “superhuman” performance

All of these models are Transformer models

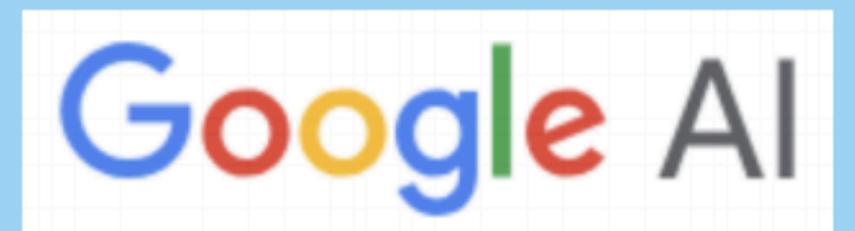
ELMo  
Oct 2017  
Training:  
800M words  
42 GPU days



GPT  
June 2018  
Training  
800M words  
240 GPU days



BERT  
Oct 2018  
Training  
3.3B words  
256 TPU days  
~320–560  
GPU days

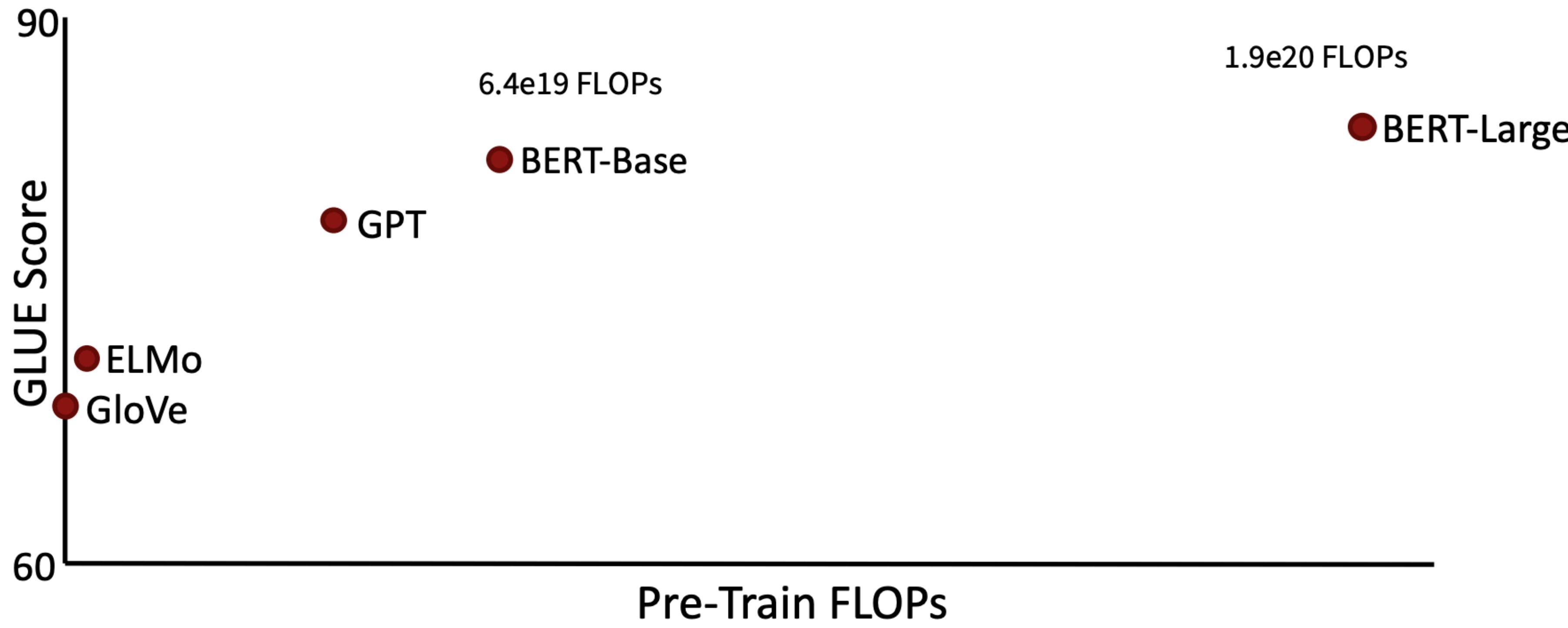


GPT-2  
Feb 2019  
Training  
40B words  
~2048 TPU v3  
days according to  
[a reddit thread](#)

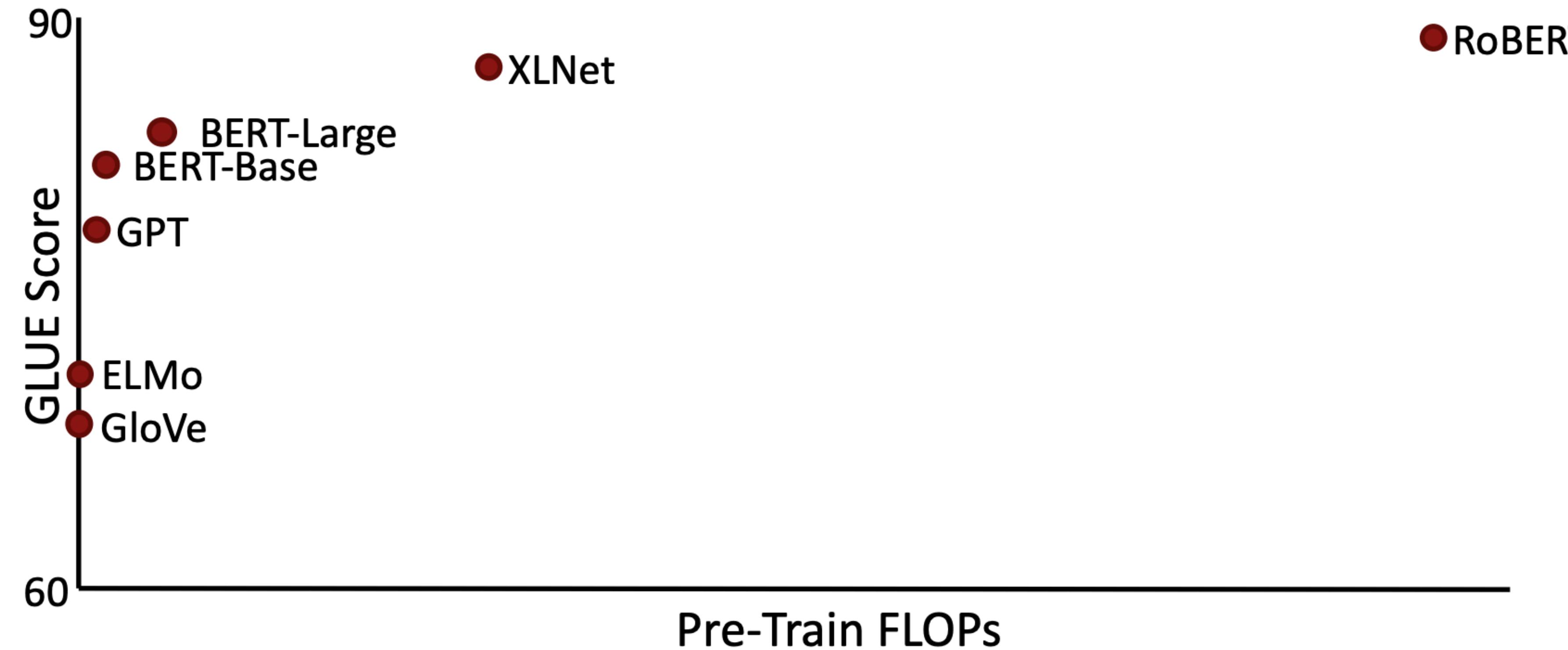


XL-Net,  
ERNIE,  
Grover  
RoBERTa, T5  
July 2019—



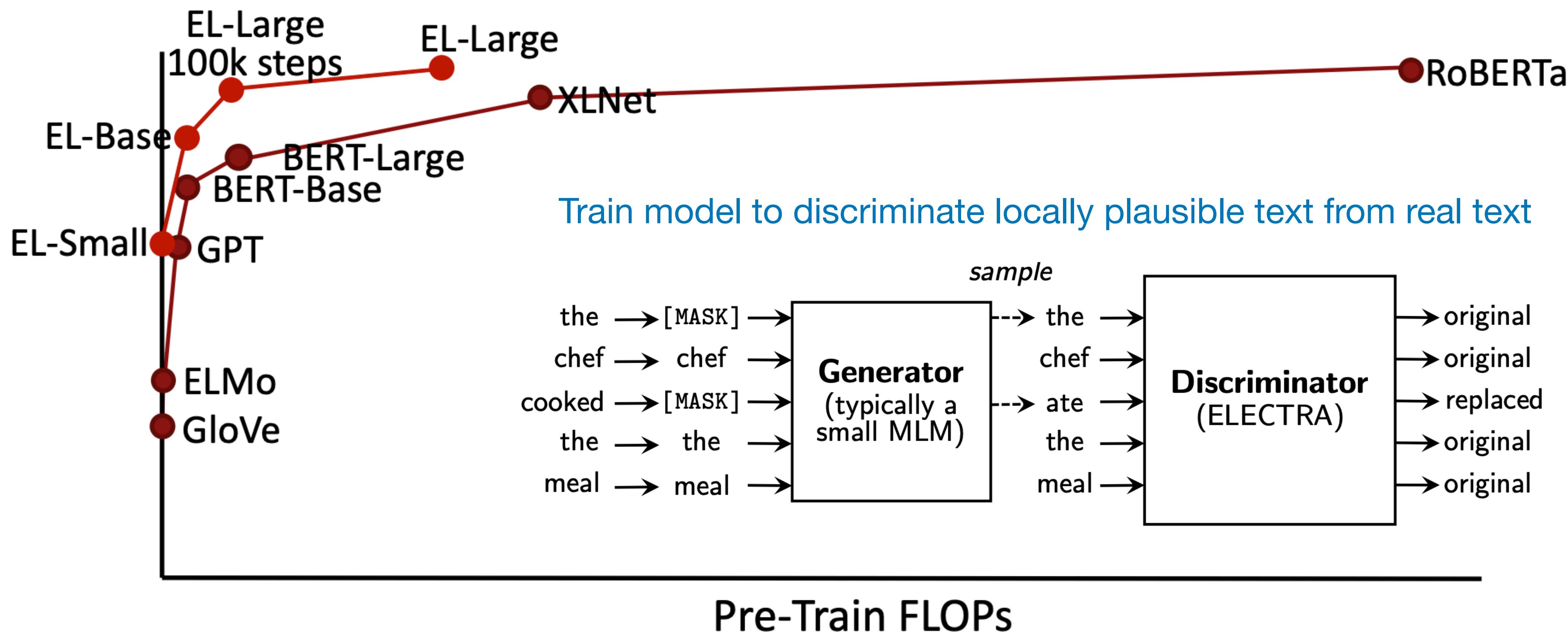


BERT-Large uses 60x more compute than ELMo

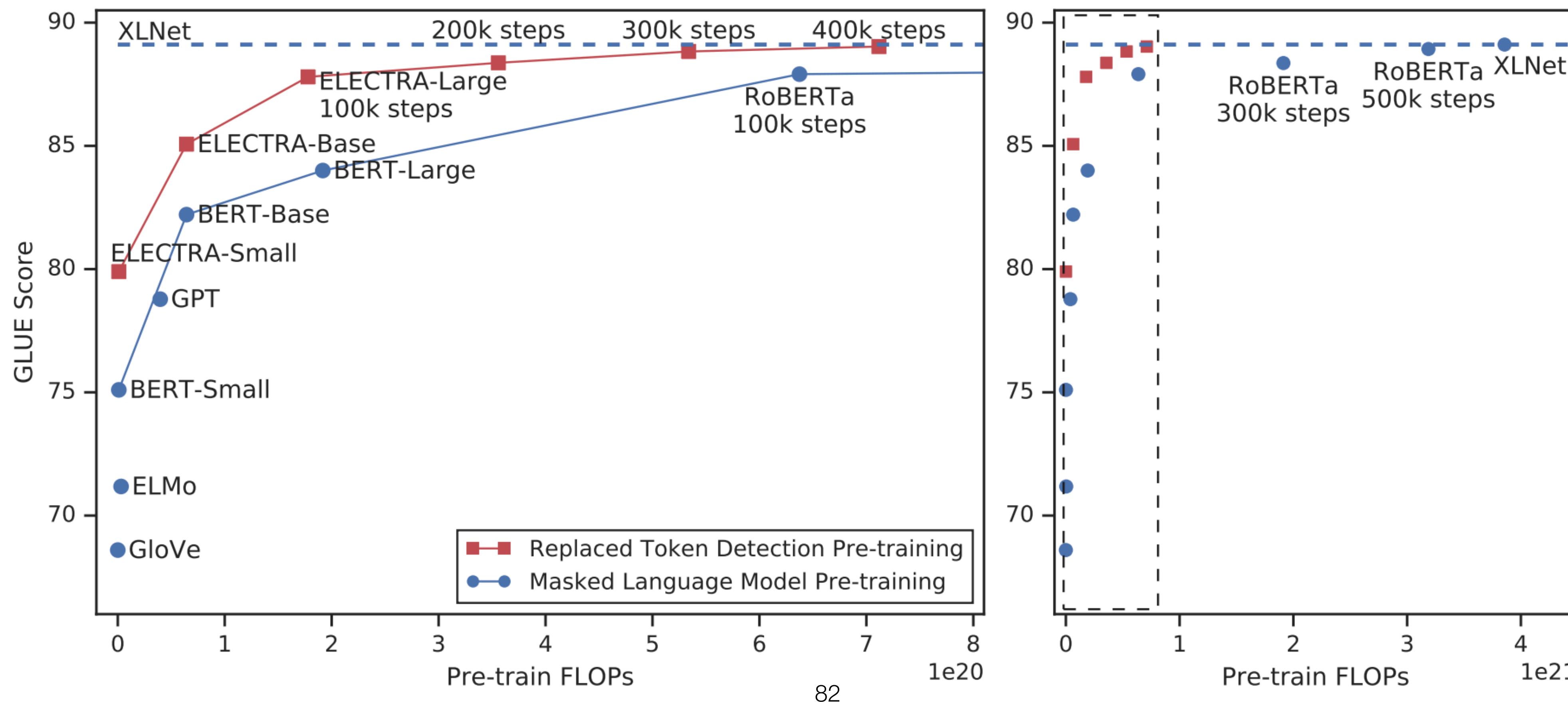


RoBERTa uses 16x more compute than BERT-Large

*ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators*  
Clark et al, ICLR 2020



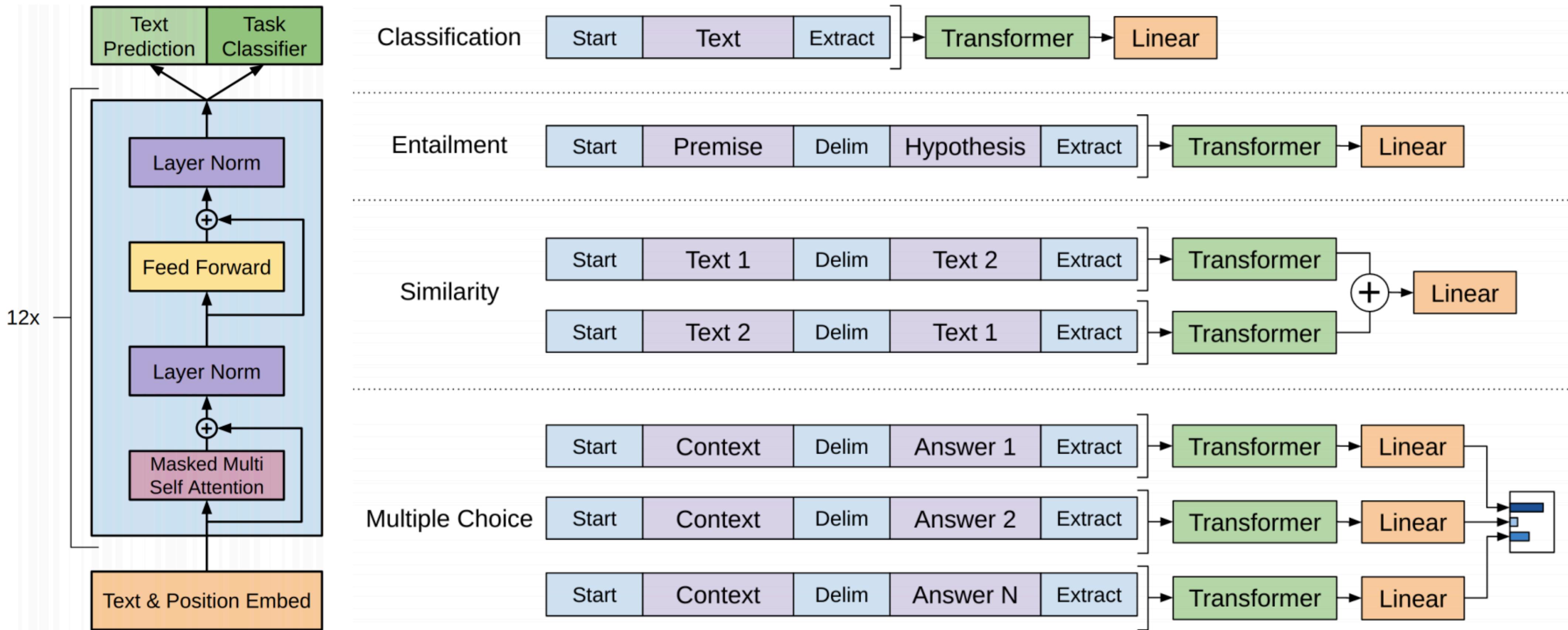
*ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators*  
Clark et al, ICLR 2020



Have fun with using ELMo or BERT in your final project :)



# GPT



*Improving Language Understanding by Generative Pre-Training*  
Radford et al, OpenAI, 2018

# RoBERTa

- Train with more data and for more epochs

	MNLI	QNLI	QQP	RTE	SST	MRPC	CoLA	STS	WNLI	Avg
<i>Single-task single models on dev</i>										
BERT <sub>LARGE</sub>	86.6/-	92.3	91.3	70.4	93.2	88.0	60.6	90.0	-	-
XLNet <sub>LARGE</sub>	89.8/-	93.9	91.8	83.8	95.6	89.2	63.6	91.8	-	-
RoBERTa	<b>90.2/90.2</b>	<b>94.7</b>	<b>92.2</b>	<b>86.6</b>	<b>96.4</b>	<b>90.9</b>	<b>68.0</b>	<b>92.4</b>	<b>91.3</b>	-

# XLNET

- Relative embeddings

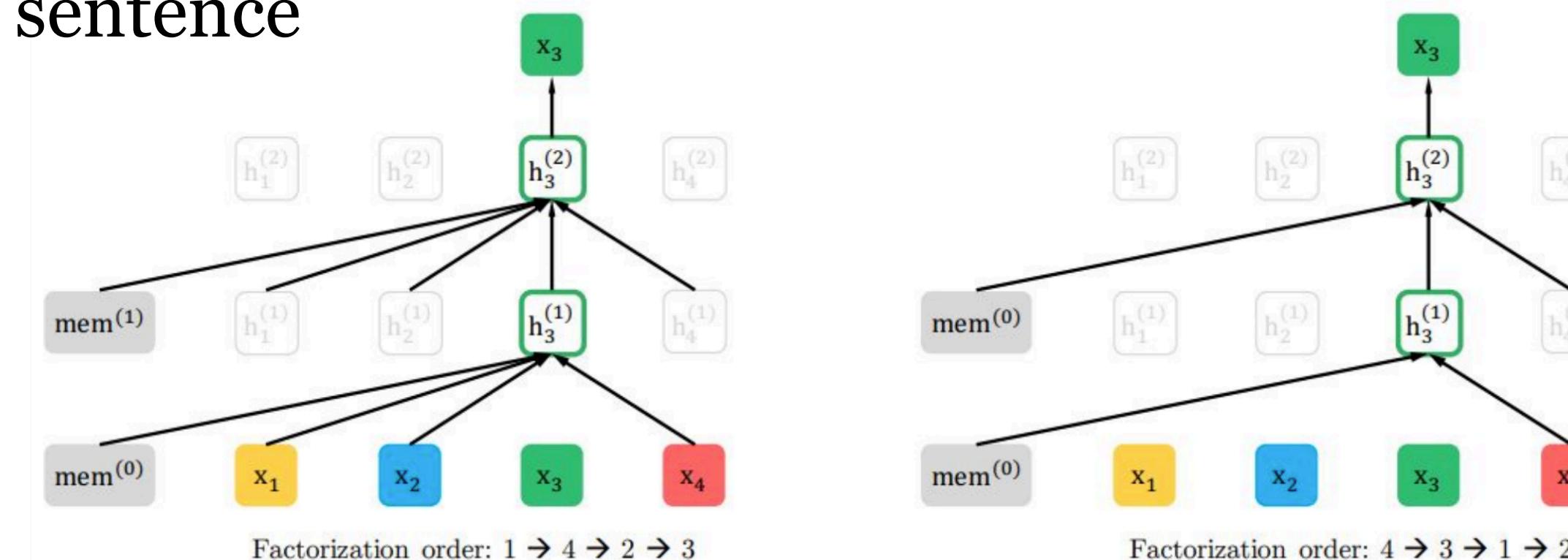
Example: John ate a **hot dog**

How much should **dog** attend to **hot**

**Absolute:** **dog** attend to **hot** (any position) and  
**dog** (in position 4) attend to **hot** (in position 3)

**Relative:** **dog** attend to **hot** (any position) and  
**dog** attend to **hot** (in previous position)

- Permutation language modeling: Randomly permute the order for every training sentence



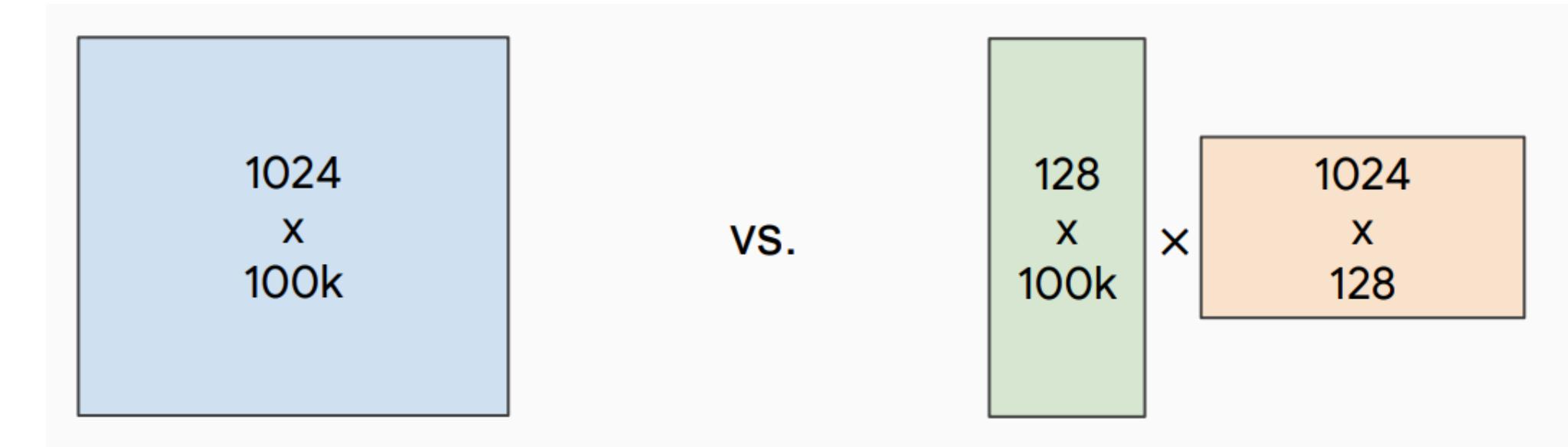
# XLNET

- Relative embeddings
- Permutation language modeling: Randomly permute the order for every training sentence

Model	MNLI	QNLI	QQP	RTE	SST-2	MRPC	CoLA	STS-B
<i>Single-task single models on dev</i>								
BERT [2]	86.6/-	92.3	91.3	70.4	93.2	88.0	60.6	90.0
RoBERTa [21]	90.2/90.2	94.7	92.2	<b>86.6</b>	96.4	<b>90.9</b>	68.0	92.4
XLNet	<b>90.8/90.8</b>	<b>94.9</b>	<b>92.3</b>	85.9	<b>97.0</b>	90.8	<b>69.0</b>	<b>92.5</b>

# ALBERT

- Factorized embedding



- Cross-layer parameter sharing: share all parameters between Transformer layers
- Less parameters, but also takes longer to train:

Model	Parameters	SQuAD1.1	SQuAD2.0	MNLI	SST-2	RACE	Avg	Speedup
BERT	base	108M	90.4/83.2	80.4/77.6	84.5	92.8	68.2	82.3
	large	334M	92.2/85.5	85.0/82.2	86.6	93.0	73.9	85.2
ALBERT	base	12M	89.3/82.3	80.0/77.1	81.6	90.3	64.0	80.1
	large	18M	90.6/83.9	82.3/79.4	83.5	91.7	68.5	82.4
	xlarge	60M	92.5/86.1	86.1/83.1	86.4	92.4	74.8	85.5
	xxlarge	235M	<b>94.1/88.3</b>	<b>88.1/85.1</b>	<b>88.0</b>	<b>95.2</b>	<b>82.3</b>	<b>88.7</b>