

# CSI2121: Big Data

Ch4. Transformers and Langauge Models

Jinyoung Yeo

Yonsei AI

# Outline

- Transformer
- Language Models

# Transformer

---

## Attention Is All You Need

---

**Ashish Vaswani\***

Google Brain

avaswani@google.com

**Noam Shazeer\***

Google Brain

noam@google.com

**Niki Parmar\***

Google Research

nikip@google.com

**Jakob Uszkoreit\***

Google Research

usz@google.com

**Llion Jones\***

Google Research

llion@google.com

**Aidan N. Gomez\*** †

University of Toronto

aidan@cs.toronto.edu

**Lukasz Kaiser\***

Google Brain

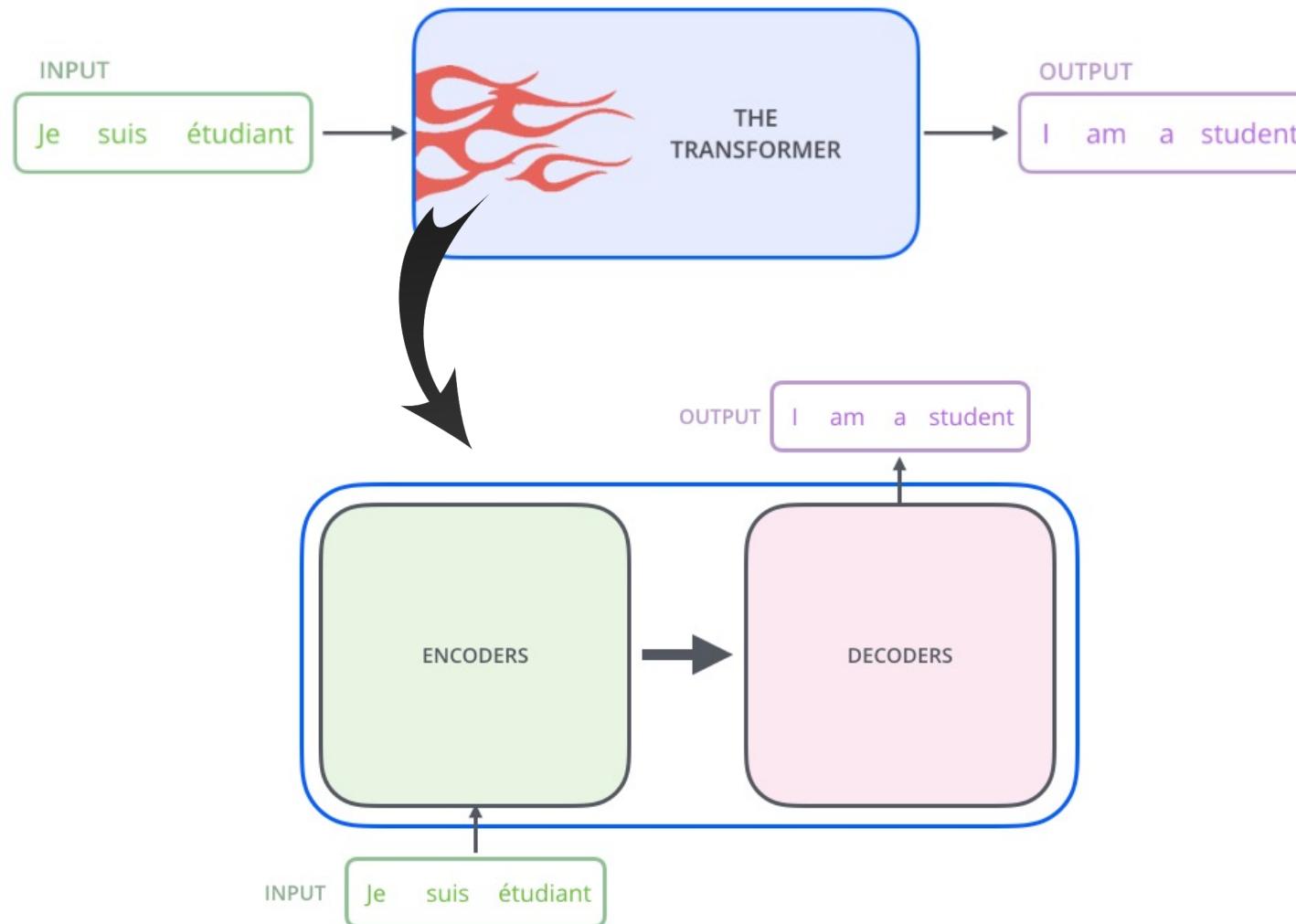
lukaszkaiser@google.com

**Illia Polosukhin\*** ‡

illia.polosukhin@gmail.com

RNN → LSTM → LSTM + attention → Only attention (Transformer)

# A High-level Look



# A High-level Look

```
class EncoderDecoder(nn.Module):
    """
    A standard Encoder-Decoder architecture. Base for this and many
    other models.
    """
    def __init__(self, encoder, decoder, src_embed, tgt_embed, generator):
        super(EncoderDecoder, self).__init__()
        self.encoder = encoder
        self.decoder = decoder
        self.src_embed = src_embed
        self.tgt_embed = tgt_embed
        self.generator = generator

    def forward(self, src, tgt, src_mask, tgt_mask):
        "Take in and process masked src and target sequences."
        return self.decode(self.encode(src, src_mask), src_mask,
                           tgt, tgt_mask)

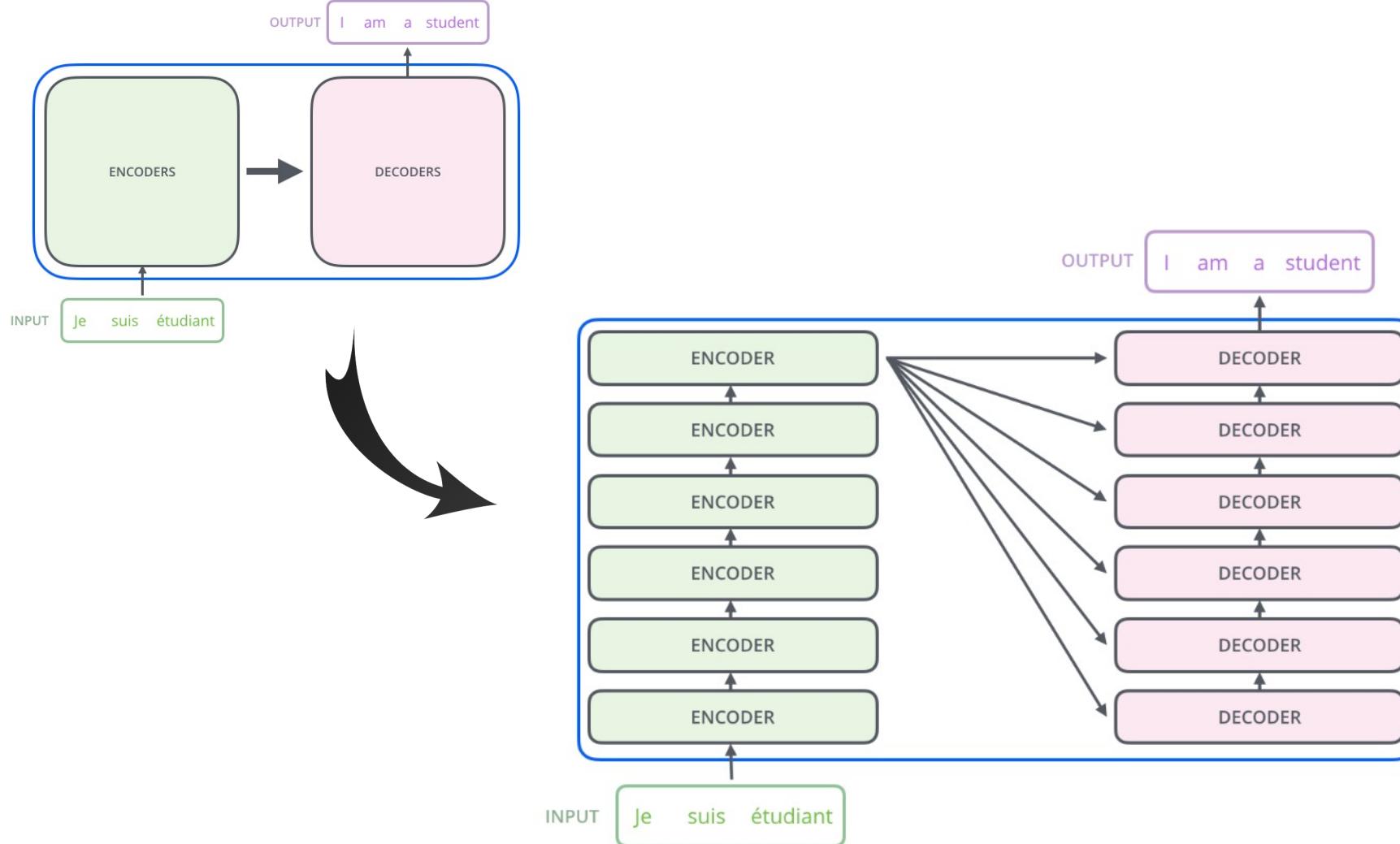
    def encode(self, src, src_mask):
        return self.encoder(self.src_embed(src), src_mask)

    def decode(self, memory, src_mask, tgt, tgt_mask):
        return self.decoder(self.tgt_embed(tgt), memory, src_mask, tgt_mask)
```

```
class Generator(nn.Module):
    "Define standard linear + softmax generation step."
    def __init__(self, d_model, vocab):
        super(Generator, self).__init__()
        self.proj = nn.Linear(d_model, vocab)

    def forward(self, x):
        return F.log_softmax(self.proj(x), dim=-1)
```

# A High Level Look



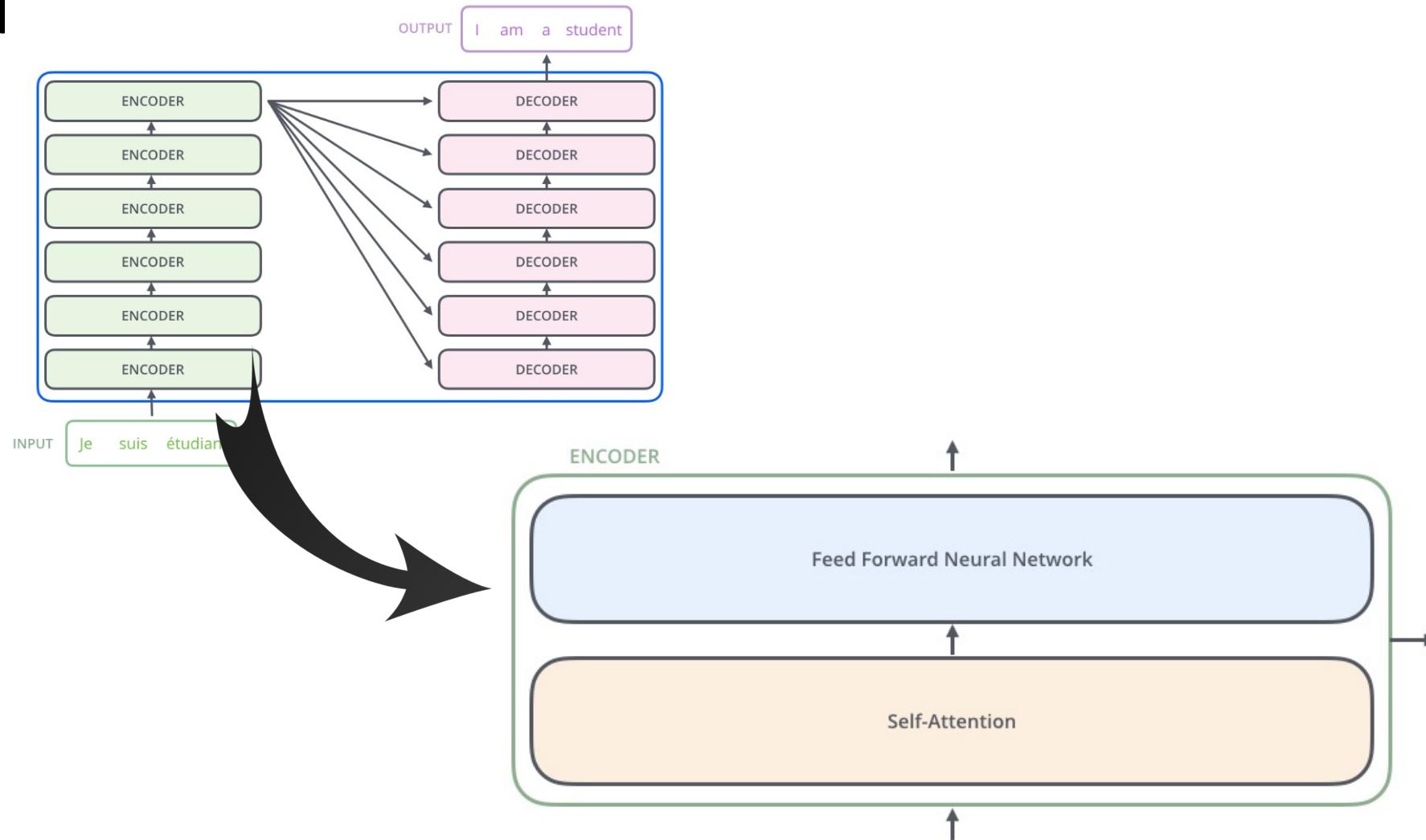
# A High-level Look

```
def clones(module, N):
    "Produce N identical layers."
    return nn.ModuleList([copy.deepcopy(module) for _ in range(N)])
```

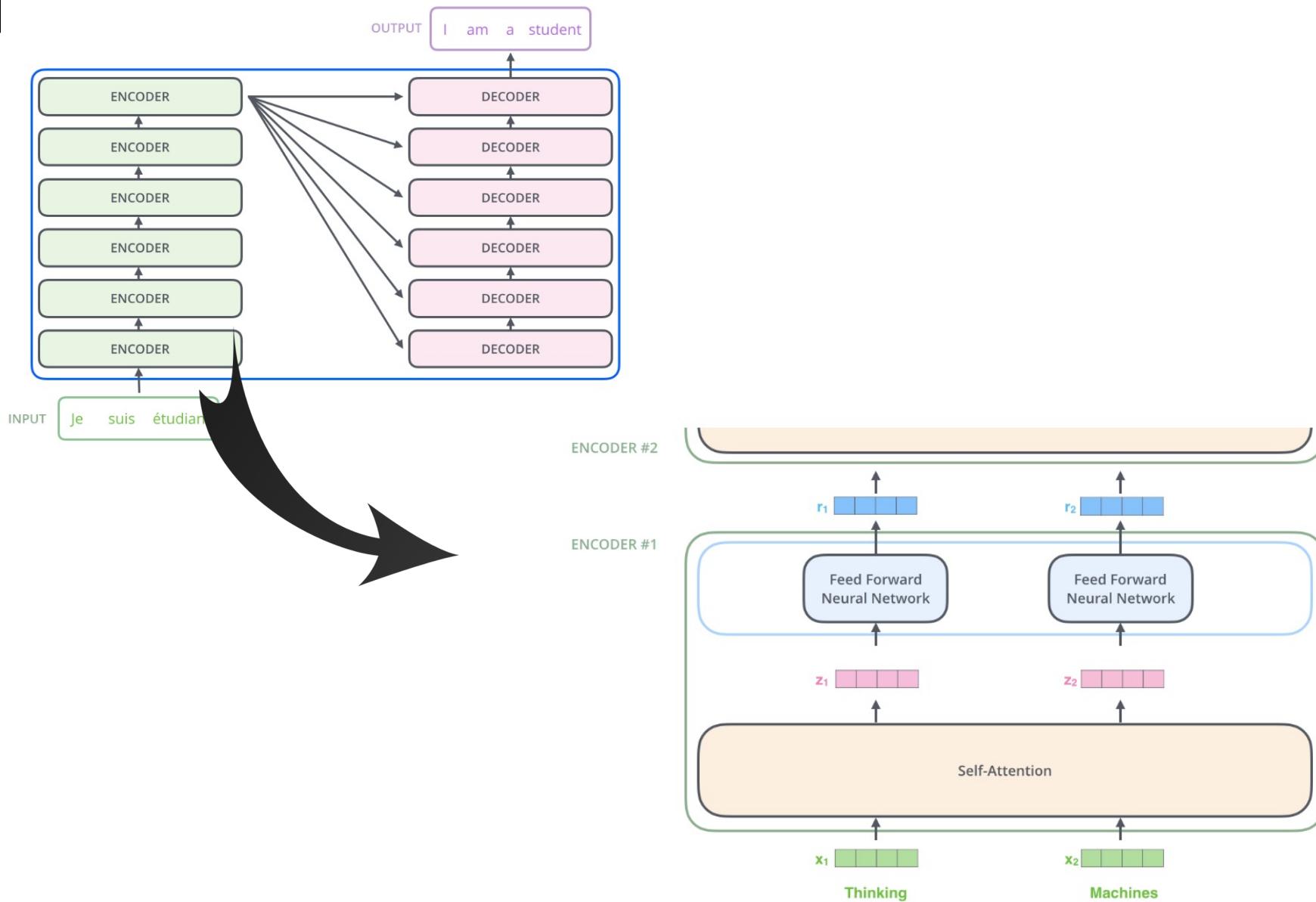
```
class Encoder(nn.Module):
    "Core encoder is a stack of N layers"
    def __init__(self, layer, N):
        super(Encoder, self).__init__()
        self.layers = clones(layer, N)
        self.norm = LayerNorm(layer.size)

    def forward(self, x, mask):
        "Pass the input (and mask) through each layer in turn."
        for layer in self.layers:
            x = layer(x, mask)
        return self.norm(x)
```

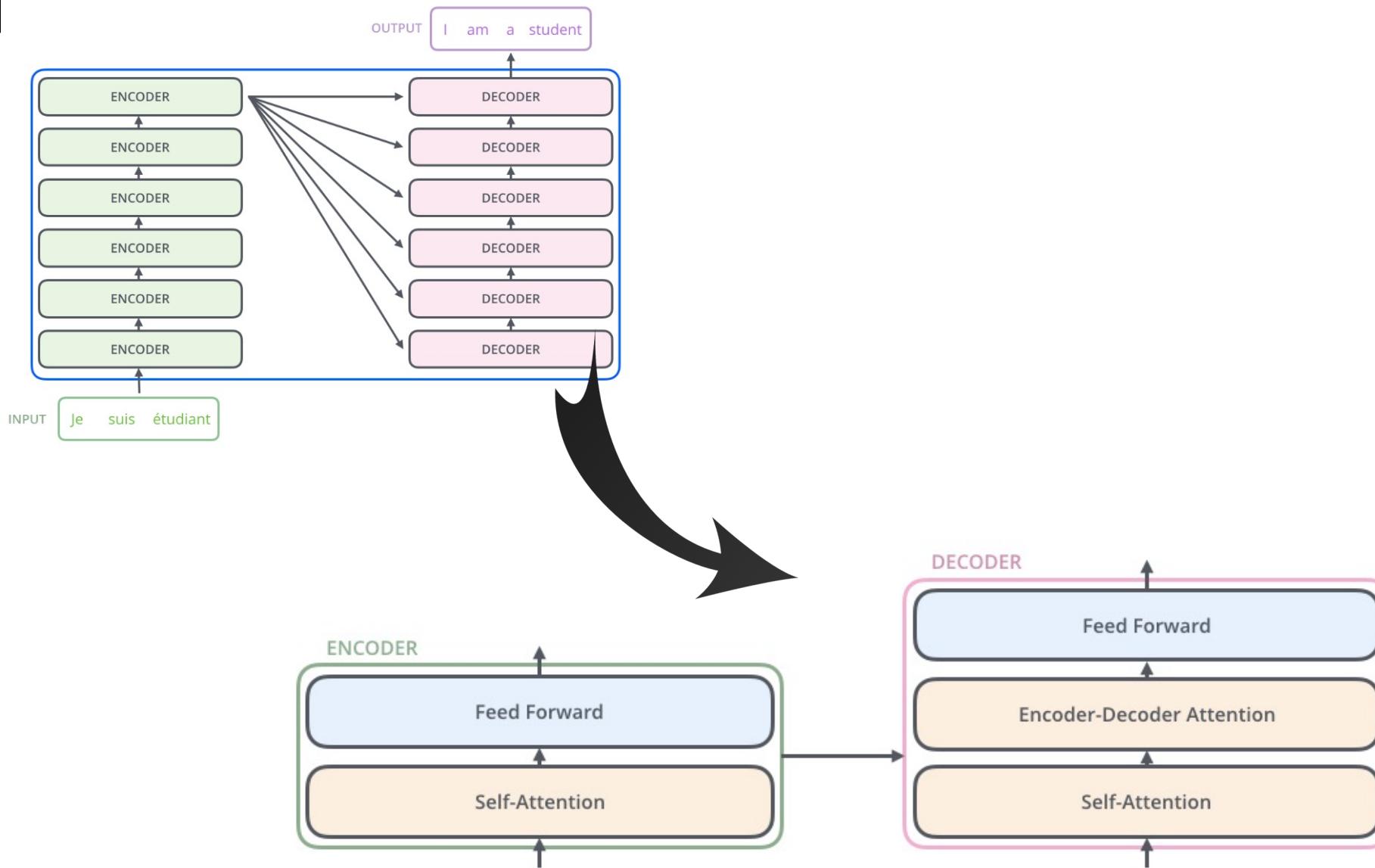
A Hierarchical Encoder-Decoder



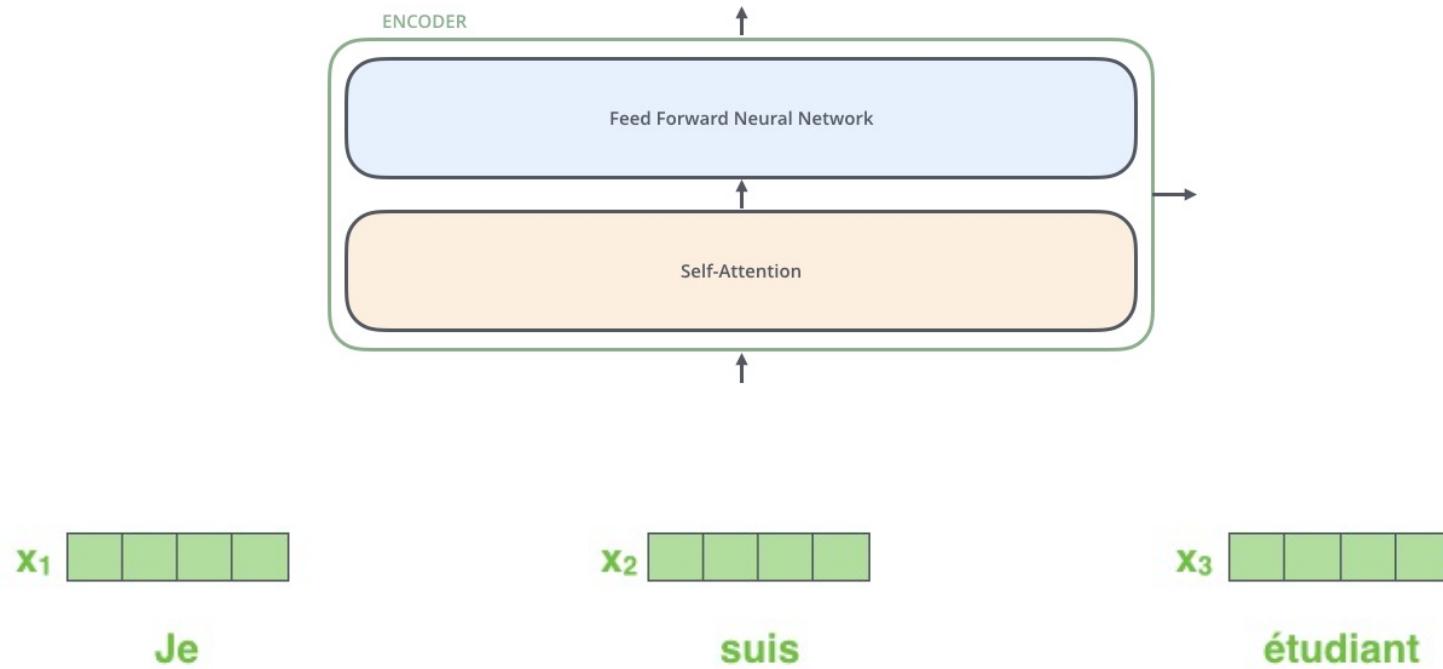
A H : -L- L- L- L- L-



# A Hierarchical Encoder-Decoder

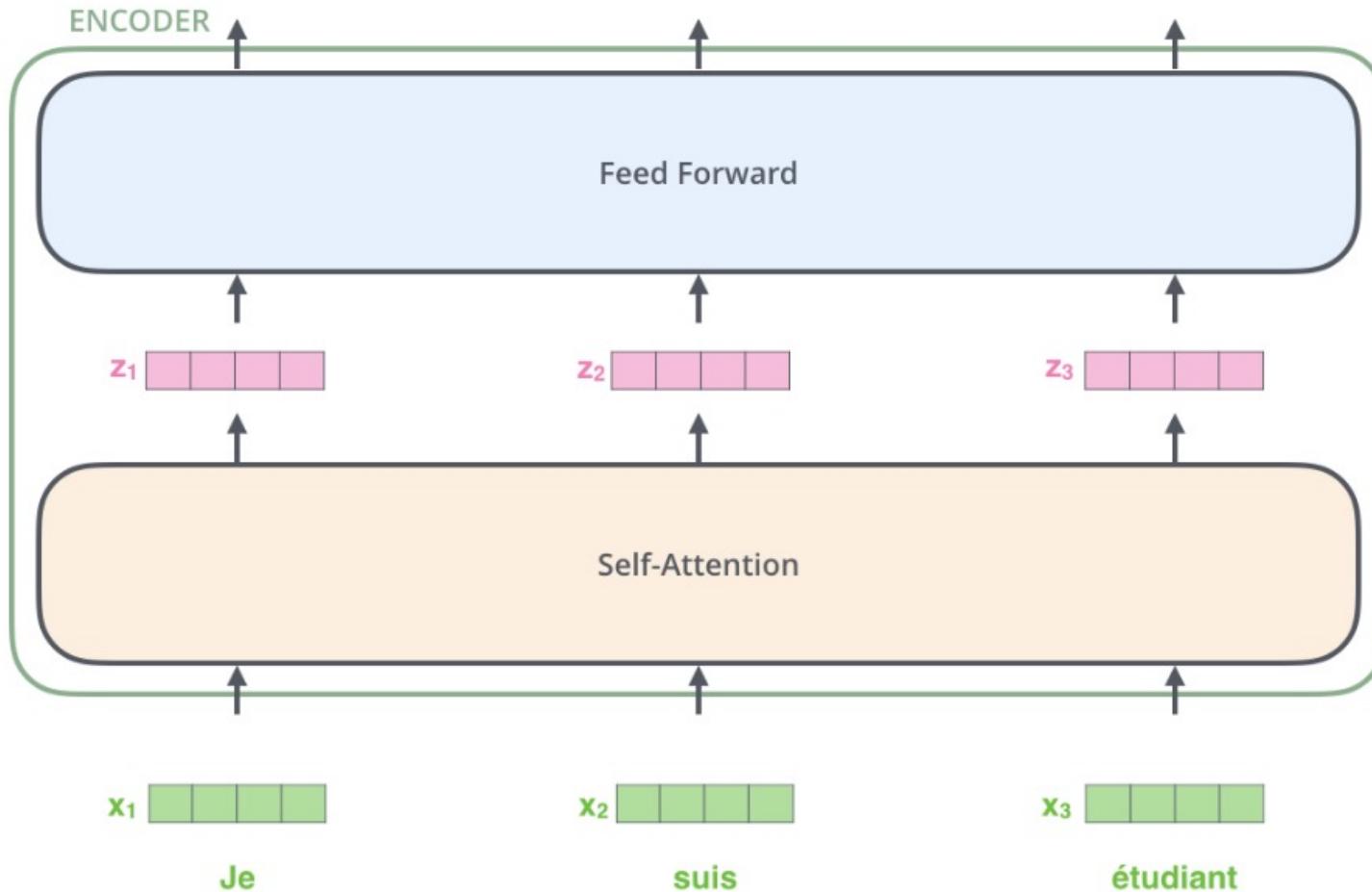


# Bringing The Tensors Into The Picture

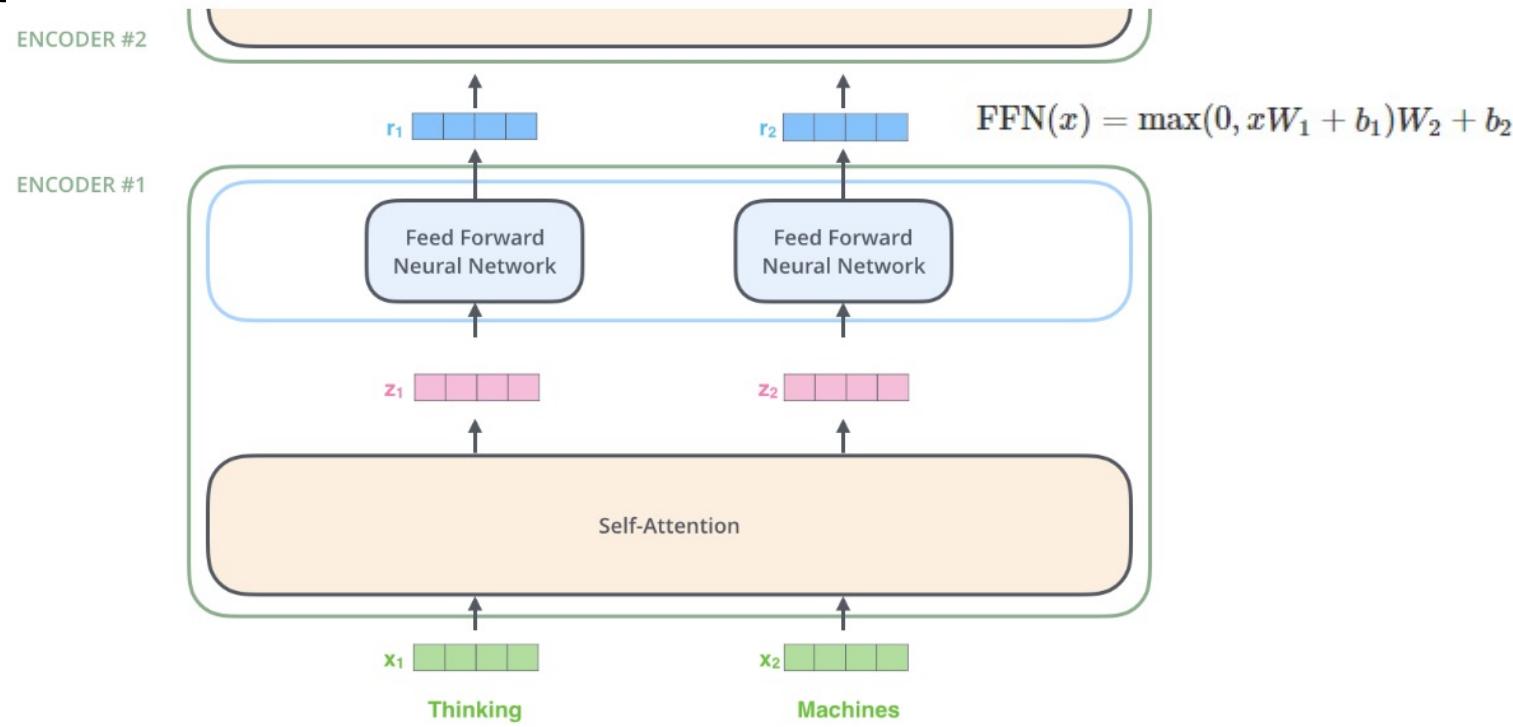


Each word is embedded into a vector of size 512. We'll represent those vectors with these simple boxes.

# Bringing The Tensors Into The Picture



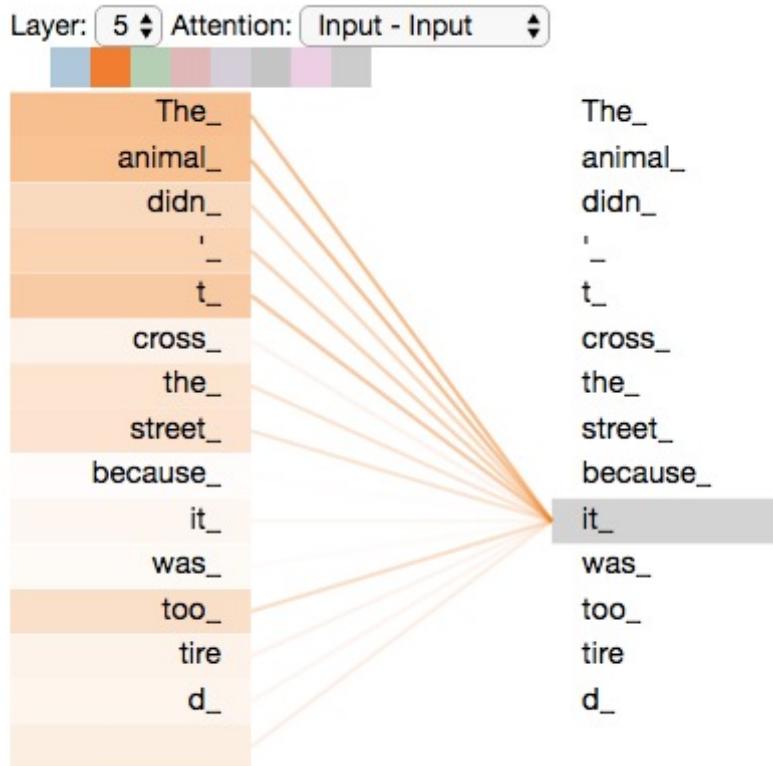
# Now We're Fncoding!



The word at each position passes through a self-attention process. Then, they each pass through a feed-forward neural network -- the exact same network with each vector flowing through it separately.

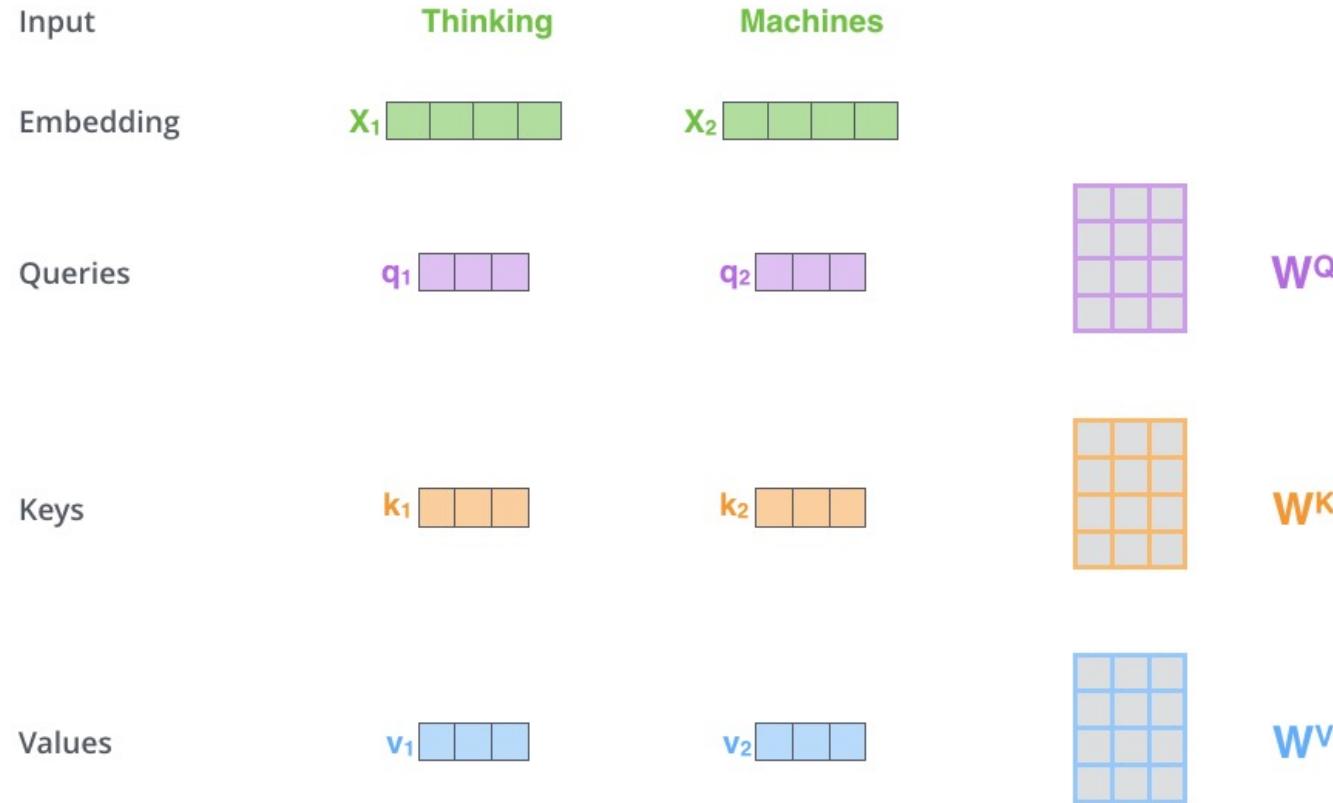
# Self-Attention at a High Level

[Tensor2Tensor notebook](#)



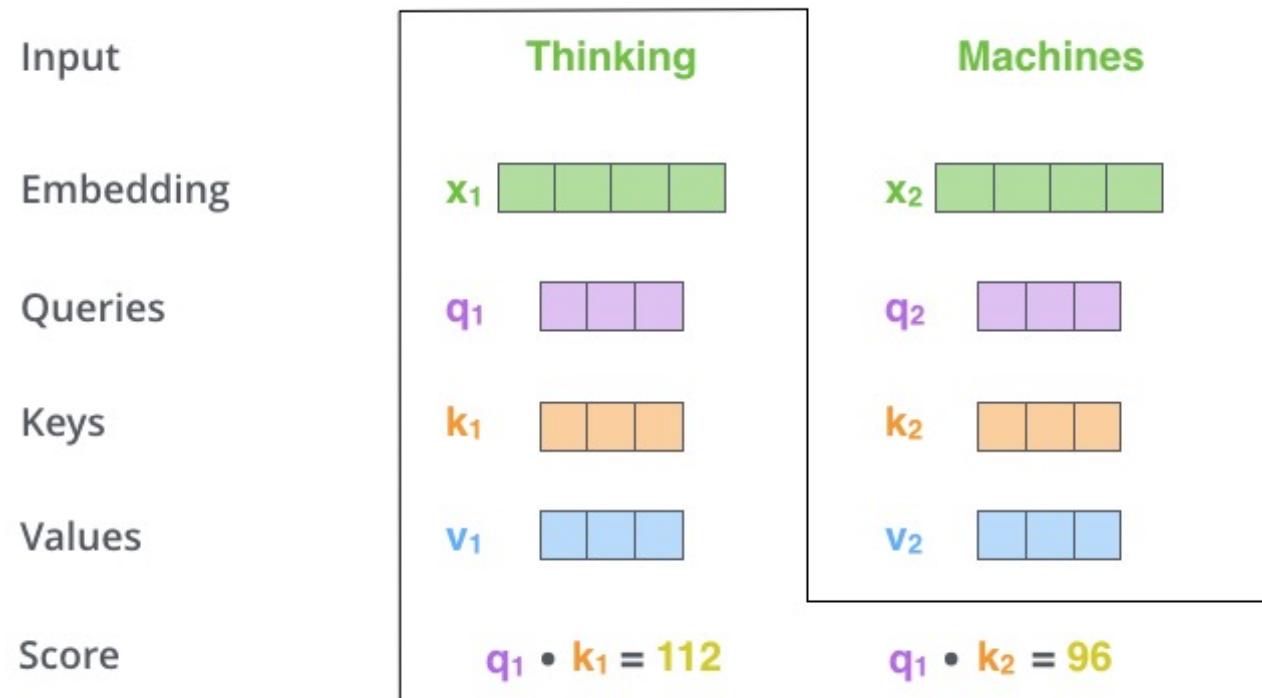
As we are encoding the word "it" in encoder #5 (the top encoder in the stack), part of the attention mechanism was focusing on "The Animal", and baked a part of its representation into the encoding of "it".

# Self-Attention in Detail

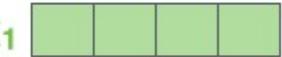
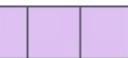
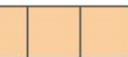
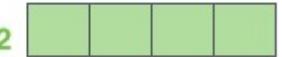
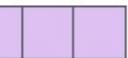
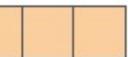


Multiplying  $x_1$  by the  $W^Q$  weight matrix produces  $q_1$ , the "query" vector associated with that word. We end up creating a "query", a "key", and a "value" projection of each word in the input sentence.

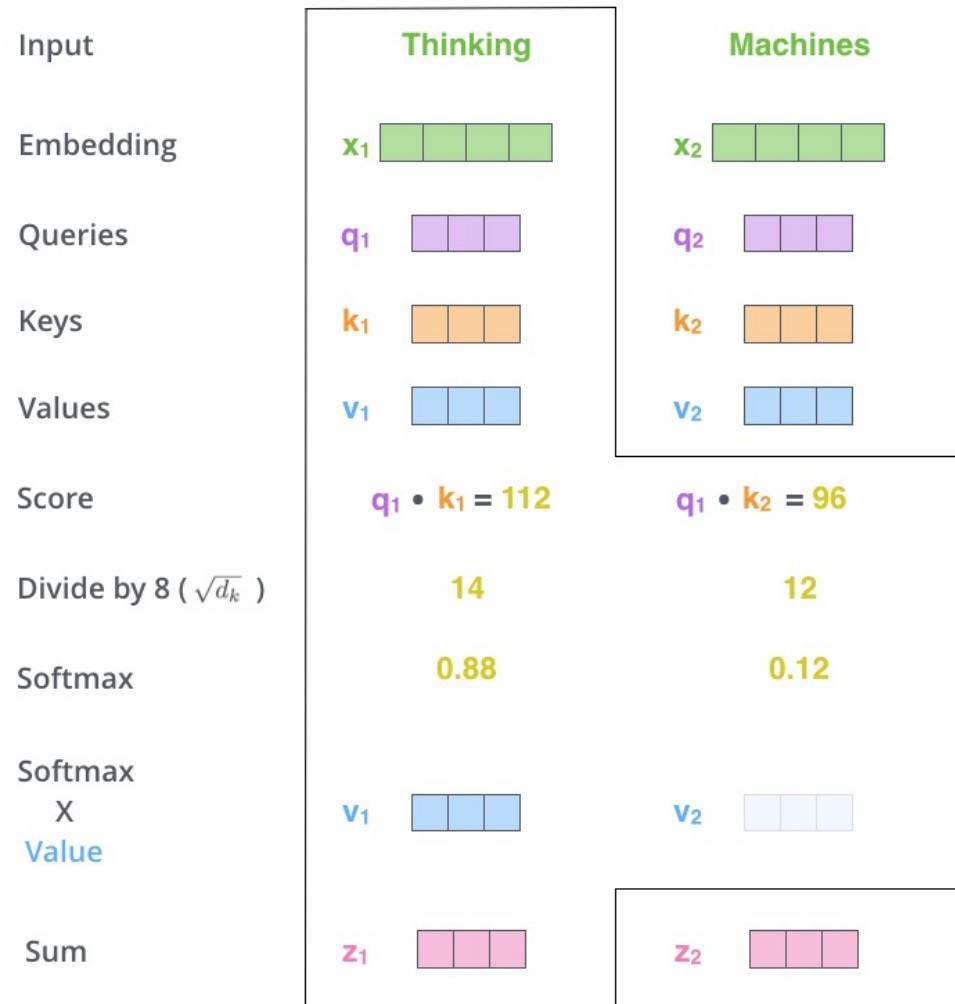
# Self-Attention in Detail



# Self-Attention in Detail

Input	<b>Thinking</b>	
Embedding	$x_1$	
Queries	$q_1$	
Keys	$k_1$	
Values	$v_1$	
Score	$q_1 \bullet k_1 = 112$	
Divide by 8 ( $\sqrt{d_k}$ )	14	
Softmax	0.88	
Input	<b>Machines</b>	
Embedding	$x_2$	
Queries	$q_2$	
Keys	$k_2$	
Values	$v_2$	
Score	$q_1 \bullet k_2 = 96$	
Divide by 8 ( $\sqrt{d_k}$ )	12	
Softmax	0.12	

# Self-Attention in Detail



# Matrix Calculation of Self-Attention

$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|c|}\hline & & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{WQ} \\ \begin{array}{|c|c|c|c|}\hline & & & \\ \hline \end{array} \end{matrix} = \begin{matrix} \text{Q} \\ \begin{array}{|c|c|c|c|}\hline & & & \\ \hline \end{array} \end{matrix}$$

$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|c|}\hline & & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{WK} \\ \begin{array}{|c|c|c|c|}\hline & & & \\ \hline \end{array} \end{matrix} = \begin{matrix} \text{K} \\ \begin{array}{|c|c|c|c|}\hline & & & \\ \hline \end{array} \end{matrix}$$

$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|c|}\hline & & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{WV} \\ \begin{array}{|c|c|c|c|}\hline & & & \\ \hline \end{array} \end{matrix} = \begin{matrix} \text{V} \\ \begin{array}{|c|c|c|c|}\hline & & & \\ \hline \end{array} \end{matrix}$$

Every row in the **X** matrix corresponds to a word in the input sentence. We again see the difference in size of the embedding vector (512, or 4 boxes in the figure), and the q/k/v vectors (64, or 3 boxes in the figure)

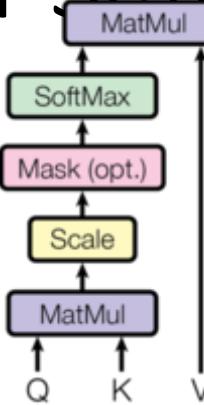
# Matrix Calculation of Self-Attention

$$\text{softmax} \left( \frac{\begin{matrix} \text{Q} & \text{K}^T \\ \times & \end{matrix}}{\sqrt{d_k}} \right) \text{V}$$
$$= \text{Z}$$

The diagram illustrates the matrix calculation of self-attention. It shows the multiplication of matrices  $\text{Q}$  and  $\text{K}^T$ , followed by division by the square root of  $d_k$ , and finally the multiplication by matrix  $\text{V}$ . The result is labeled  $\text{Z}$ .

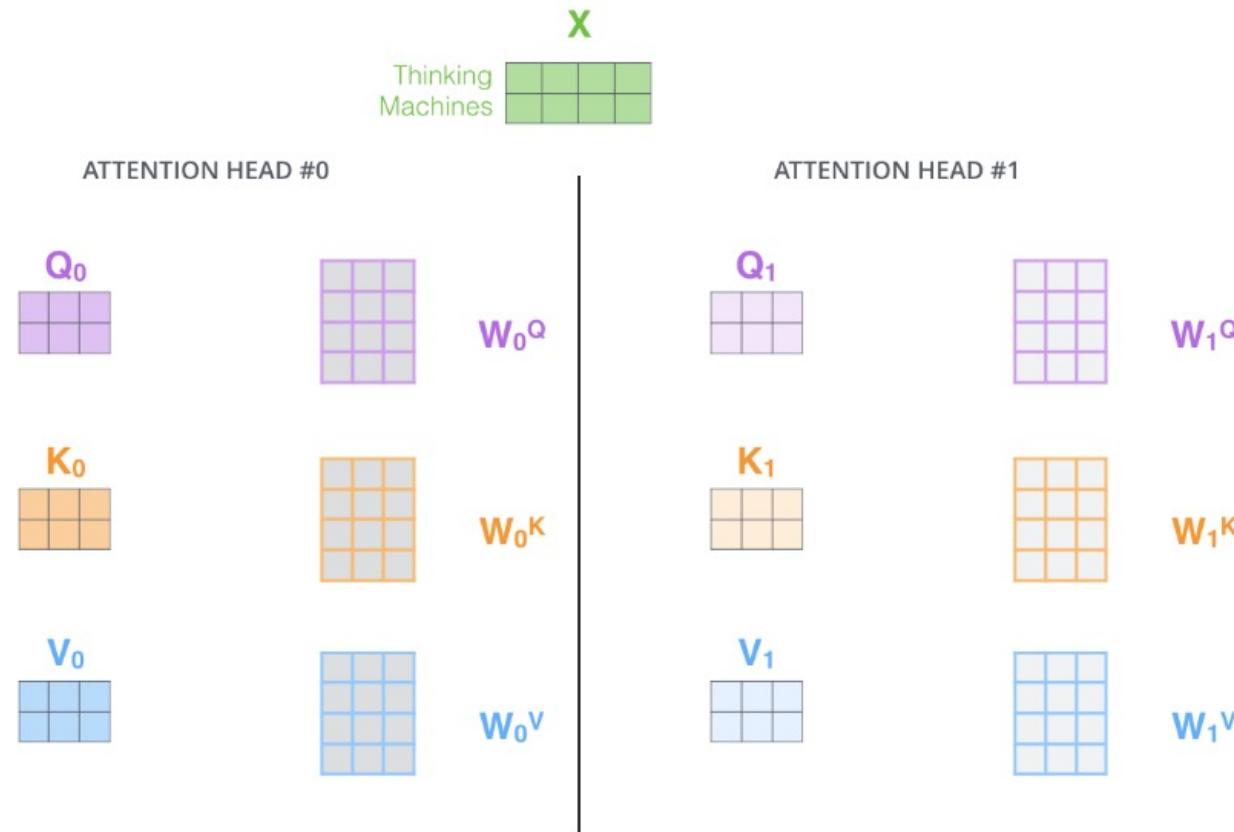
Matrix  $\text{Q}$  is a 2x4 matrix with purple squares. Matrix  $\text{K}^T$  is a 4x2 matrix with orange squares. Matrix  $\text{V}$  is a 2x2 matrix with blue squares. The result matrix  $\text{Z}$  is a 2x2 matrix with pink squares.

# Matrix Calculation of Self-Attention



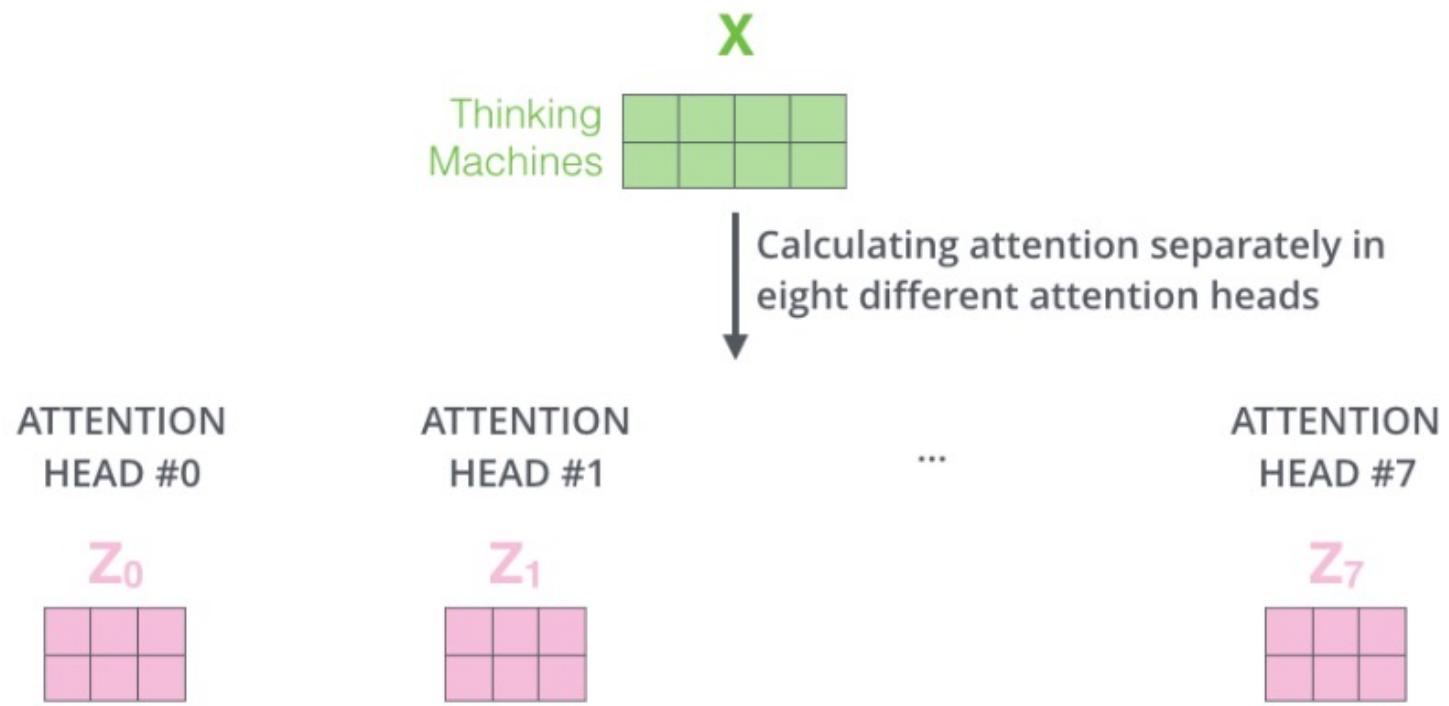
```
def attention(query, key, value, mask=None, dropout=None):
    "Compute 'Scaled Dot Product Attention'"
    d_k = query.size(-1)
    scores = torch.matmul(query, key.transpose(-2, -1)) \
        / math.sqrt(d_k)
    if mask is not None:
        scores = scores.masked_fill(mask == 0, -1e9)
    p_attn = F.softmax(scores, dim = -1)
    if dropout is not None:
        p_attn = dropout(p_attn)
    return torch.matmul(p_attn, value), p_attn
```

# Multi-head Attention



With multi-headed attention, we maintain separate Q/K/V weight matrices for each head resulting in different Q/K/V matrices. As we did before, we multiply  $X$  by the  $WQ/WK/WV$  matrices to produce Q/K/V matrices.

# Multi-head Attention



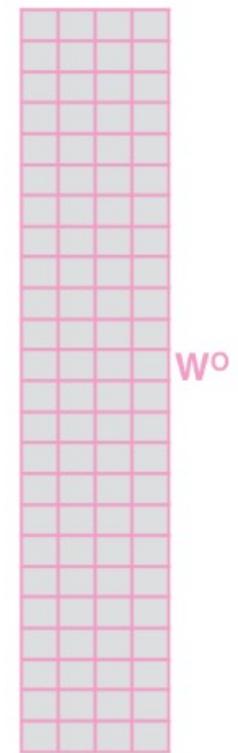
# Multi-head Attention

1) Concatenate all the attention heads



2) Multiply with a weight matrix  $W^o$  that was trained jointly with the model

$\times$



3) The result would be the  $Z$  matrix that captures information from all the attention heads. We can send this forward to the FFNN

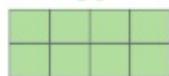
$$= \begin{matrix} Z \\ \begin{matrix} 4 \times 4 & \text{grid} \end{matrix} \end{matrix}$$

# Multi-head Attention

- 1) This is our input sentence\* each word\*
- 2) We embed
- 3) Split into 8 heads. We multiply  $X$  or  $R$  with weight matrices
- 4) Calculate attention using the resulting  $Q/K/V$  matrices
- 5) Concatenate the resulting  $Z$  matrices, then multiply with weight matrix  $W^O$  to produce the output of the layer

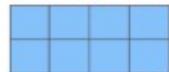
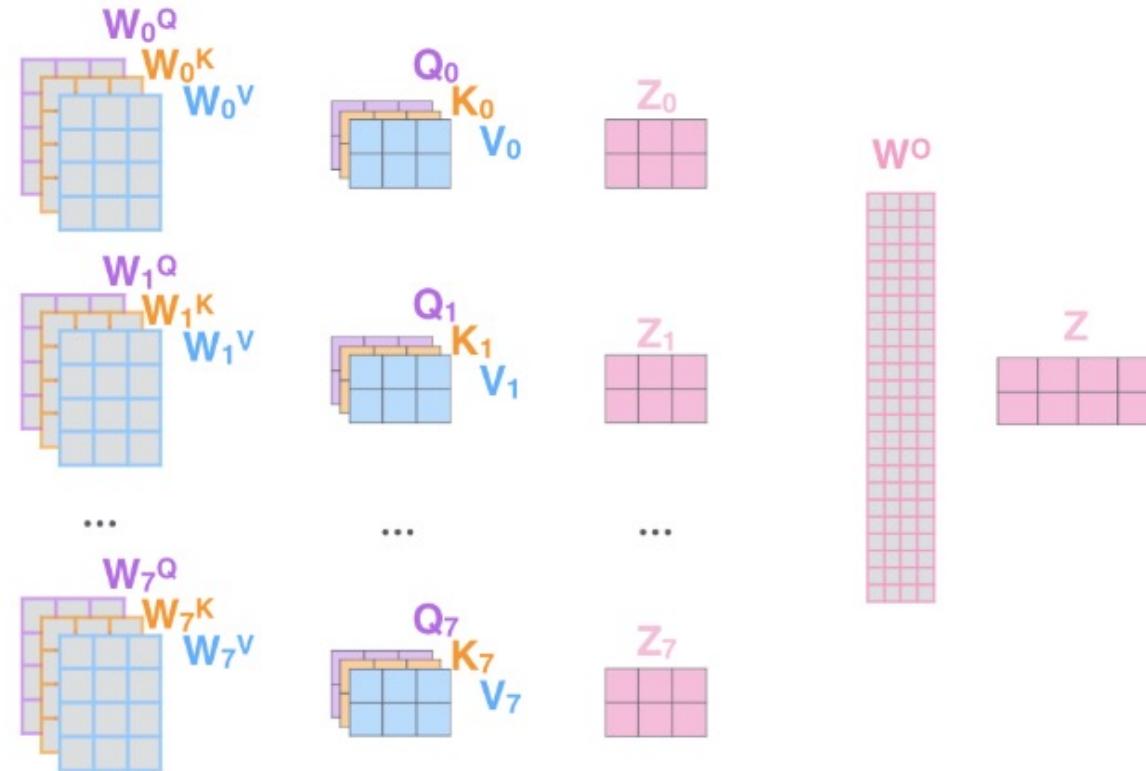
Thinking Machines

$X$

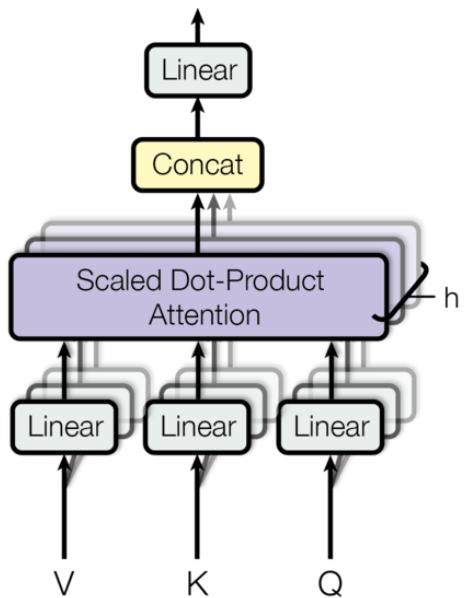
A green 4x4 grid labeled 'X'.

\* In all encoders other than #0, we don't need embedding.  
We start directly with the output of the encoder right below this one

$R$

A blue 4x4 grid labeled 'R'.

# Multi-head Attention



```
class MultiHeadedAttention(nn.Module):
    def __init__(self, h, d_model, dropout=0.1):
        "Take in model size and number of heads."
        super(MultiHeadedAttention, self).__init__()
        assert d_model % h == 0
        # We assume d_v always equals d_k
        self.d_k = d_model // h
        self.h = h
        self.linears = clones(nn.Linear(d_model, d_model), 4)
        self.attn = None
        self.dropout = nn.Dropout(p=dropout)

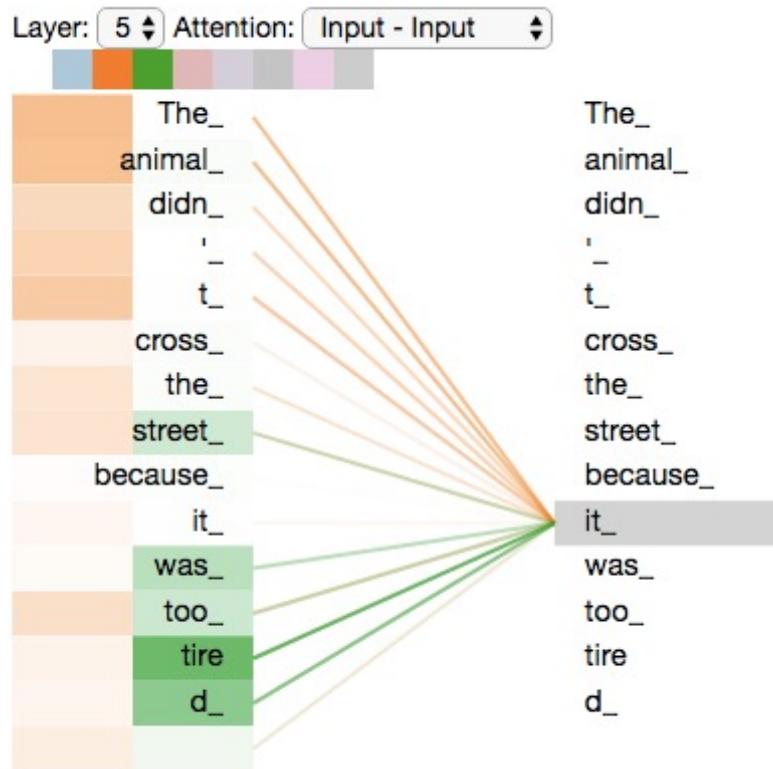
    def forward(self, query, key, value, mask=None):
        "Implements Figure 2"
        if mask is not None:
            # Same mask applied to all h heads.
            mask = mask.unsqueeze(1)
            nbatches = query.size(0)

        # 1) Do all the linear projections in batch from d_model => h x d_k
        query, key, value = \
            [l(x).view(nbatches, -1, self.h, self.d_k).transpose(1, 2)
             for l, x in zip(self.linears, (query, key, value))]

        # 2) Apply attention on all the projected vectors in batch.
        x, self.attn = attention(query, key, value, mask=mask,
                               dropout=self.dropout)

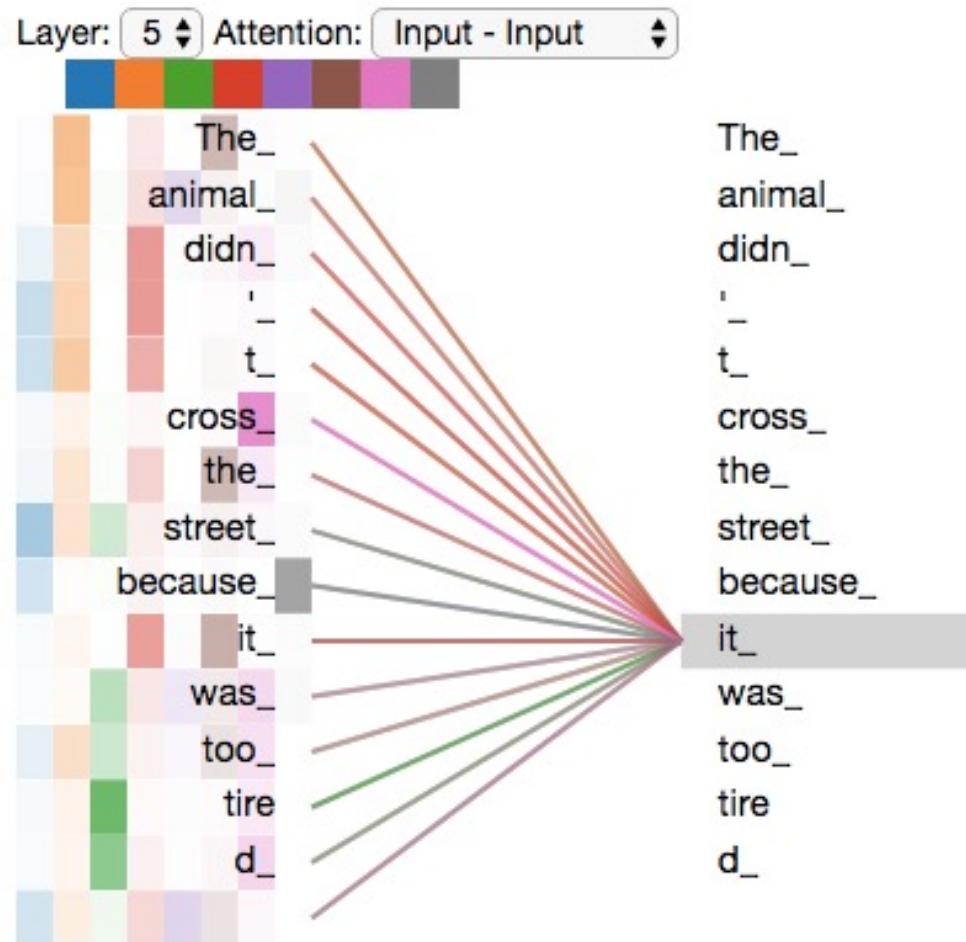
        # 3) "Concat" using a view and apply a final linear.
        x = x.transpose(1, 2).contiguous() \
            .view(nbatches, -1, self.h * self.d_k)
        return self.linears[-1](x)
```

# Multi-head Attention

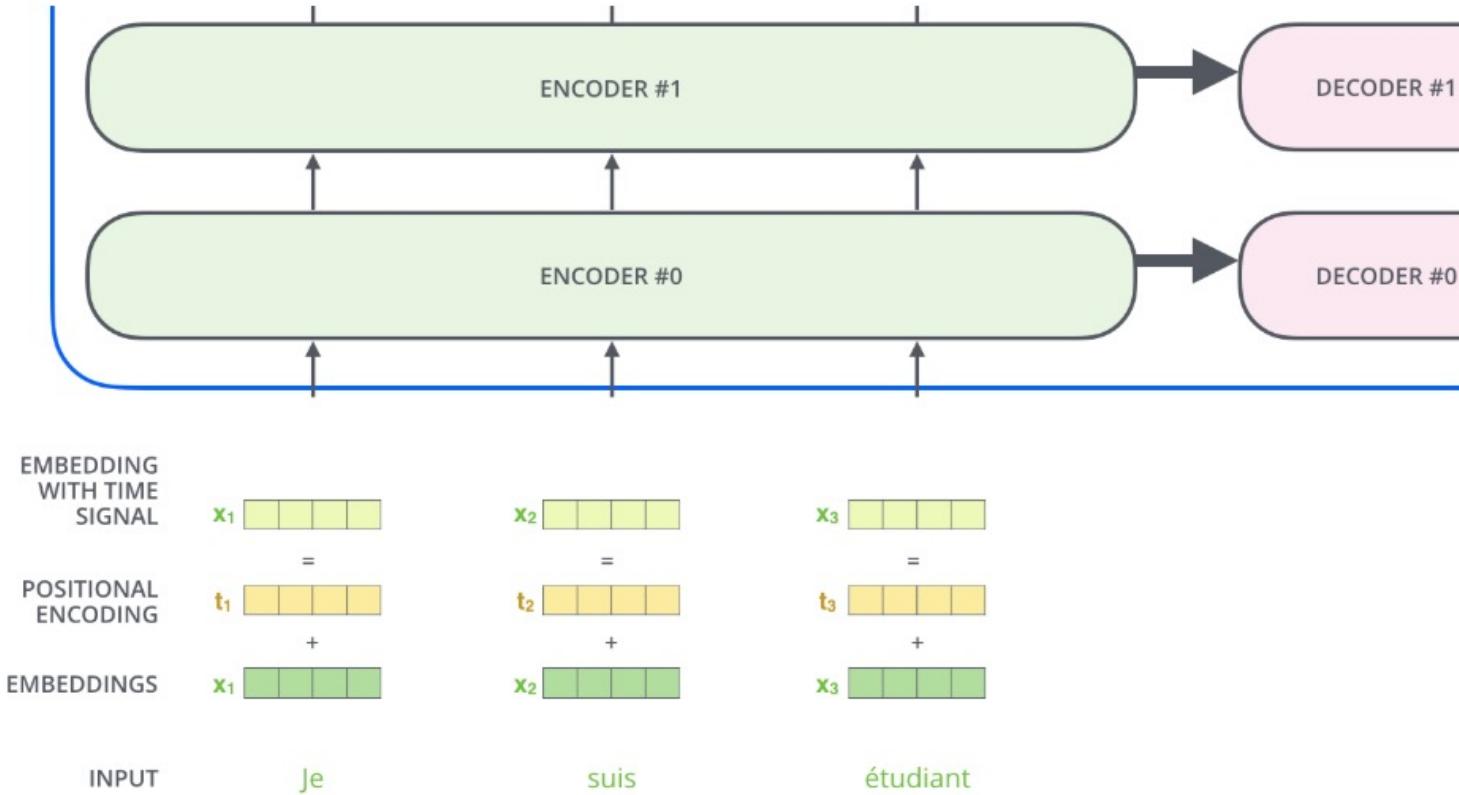


As we encode the word "it", one attention head is focusing most on "the animal", while another is focusing on "tired" -- in a sense, the model's representation of the word "it" bakes in some of the representation of both "animal" and "tired".

# Multi-head Attention

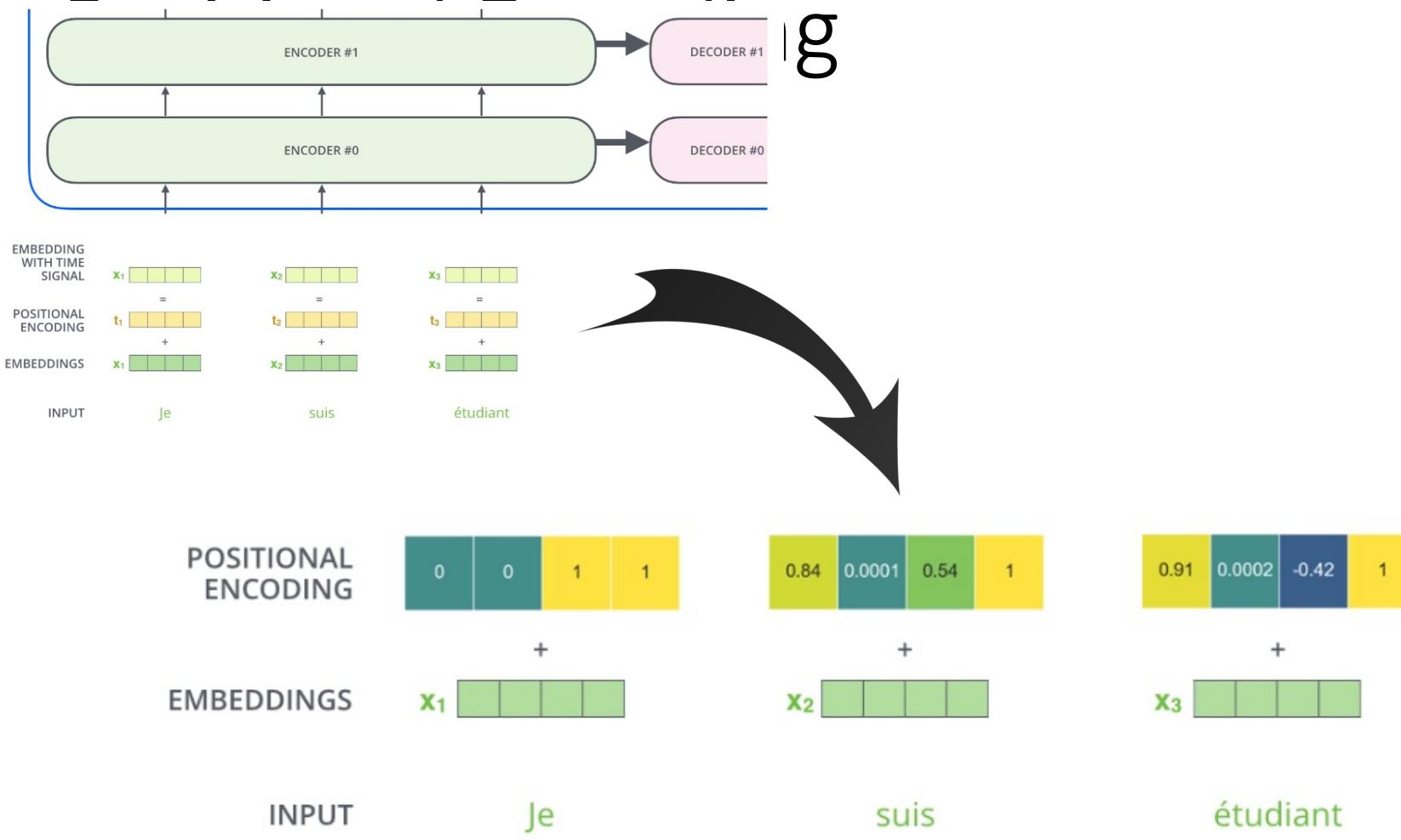


# Representing The Order of The Sequence Using F



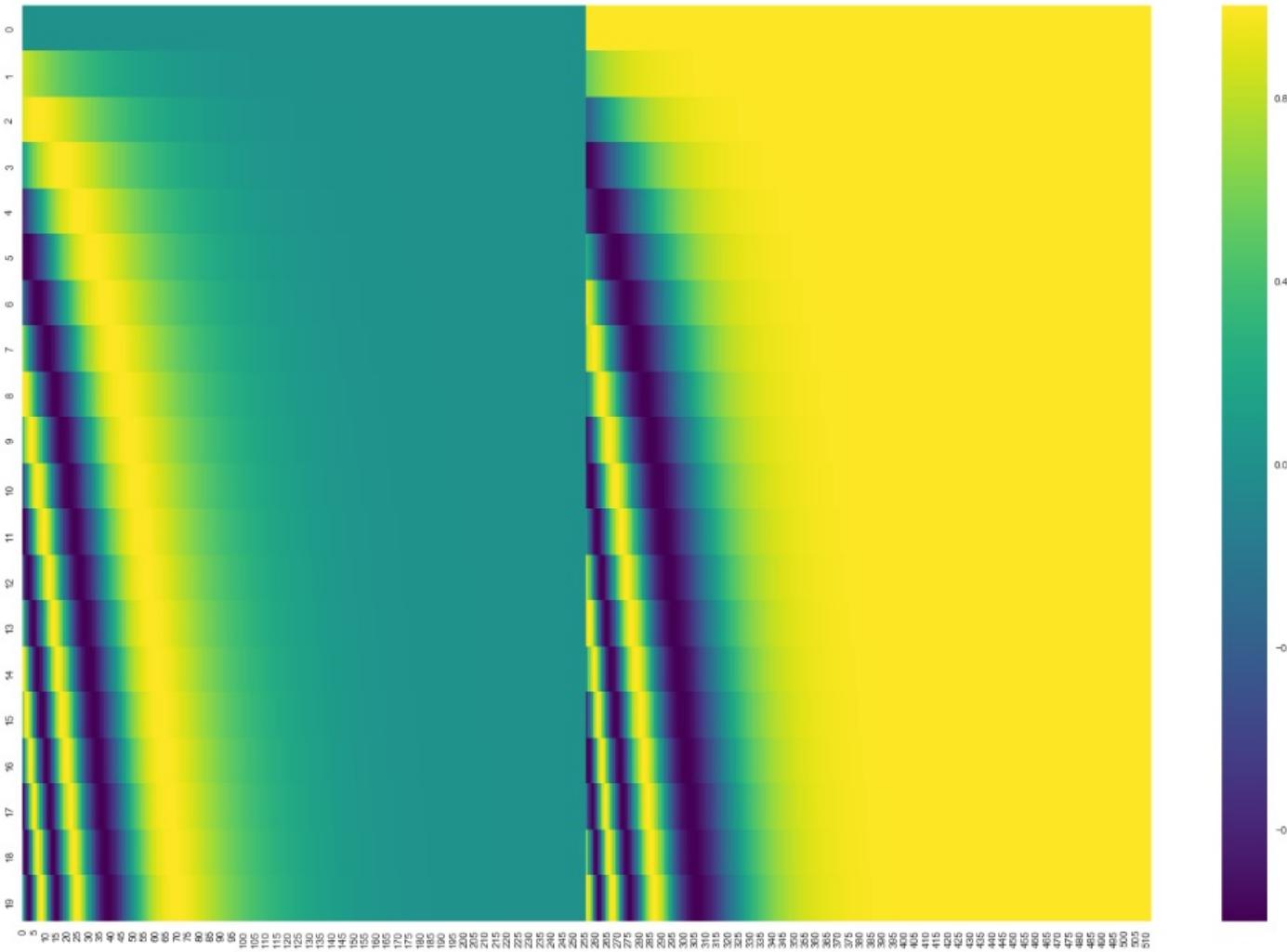
To give the model a sense of the order of the words, we add positional encoding vectors -- the values of which follow a specific pattern.

# Representing The Order of The Sequence Using Positional Encoding

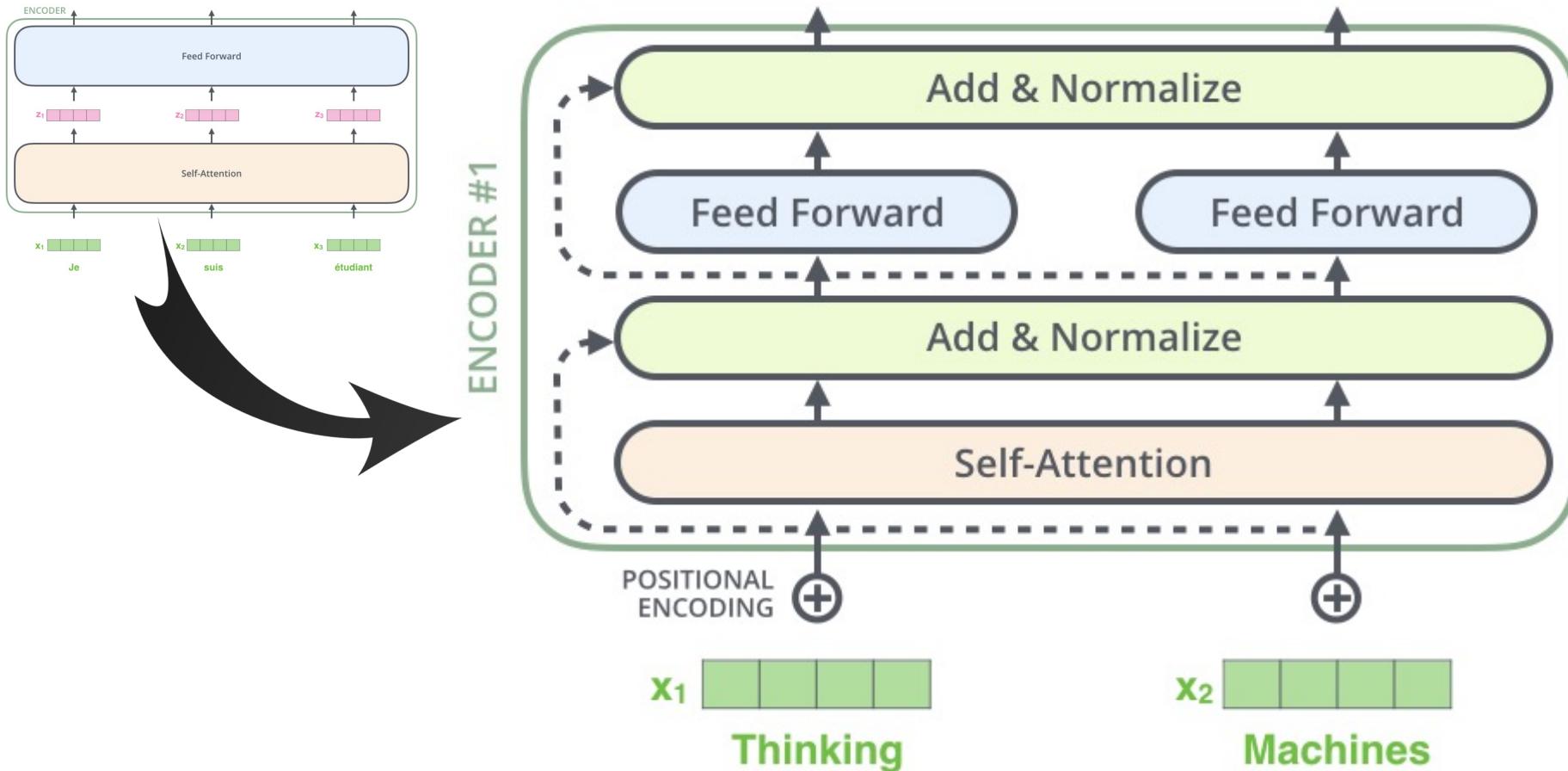


A real example of positional encoding with a toy embedding size of 4

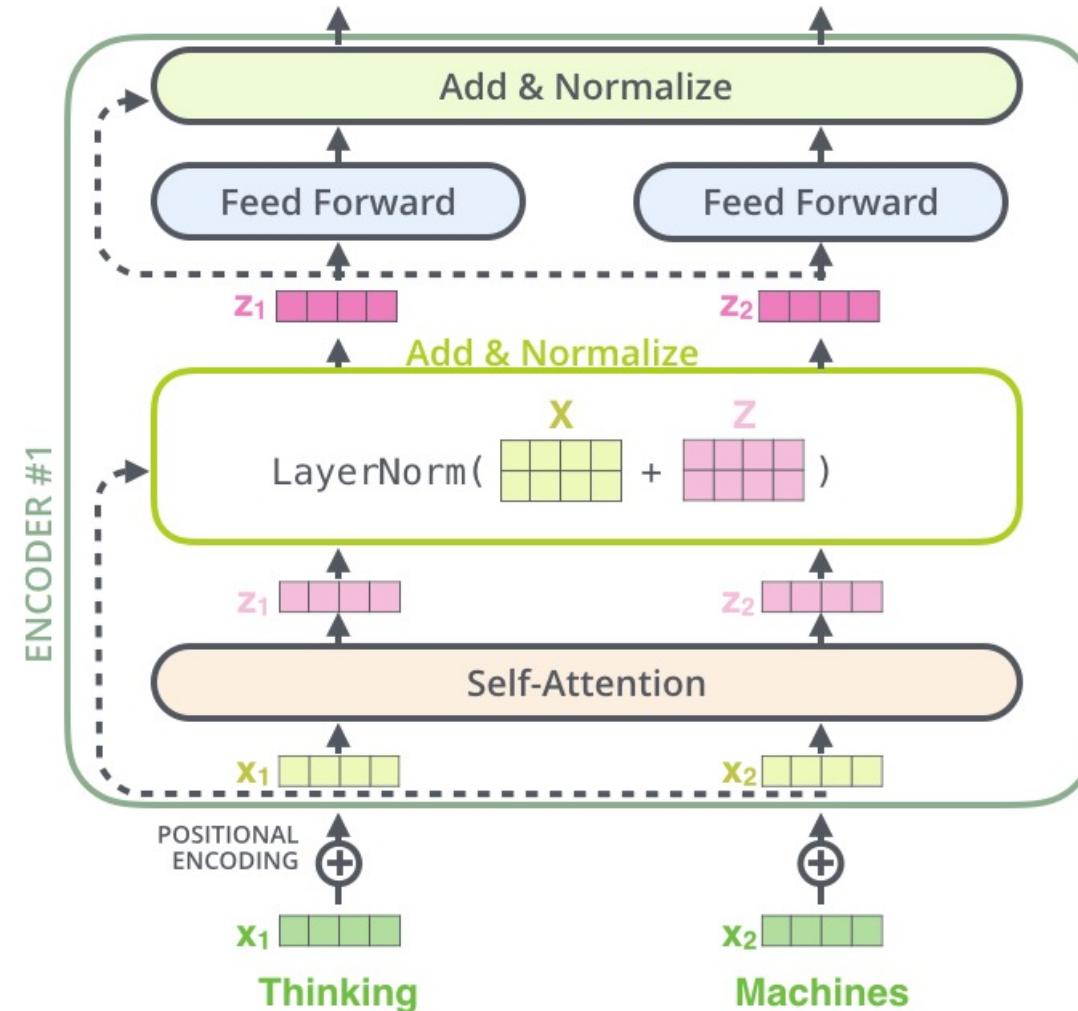
# Representing The Order of The Sequence Using Positional Encoding



# The Residuals



# The Residuals



# The Residuals

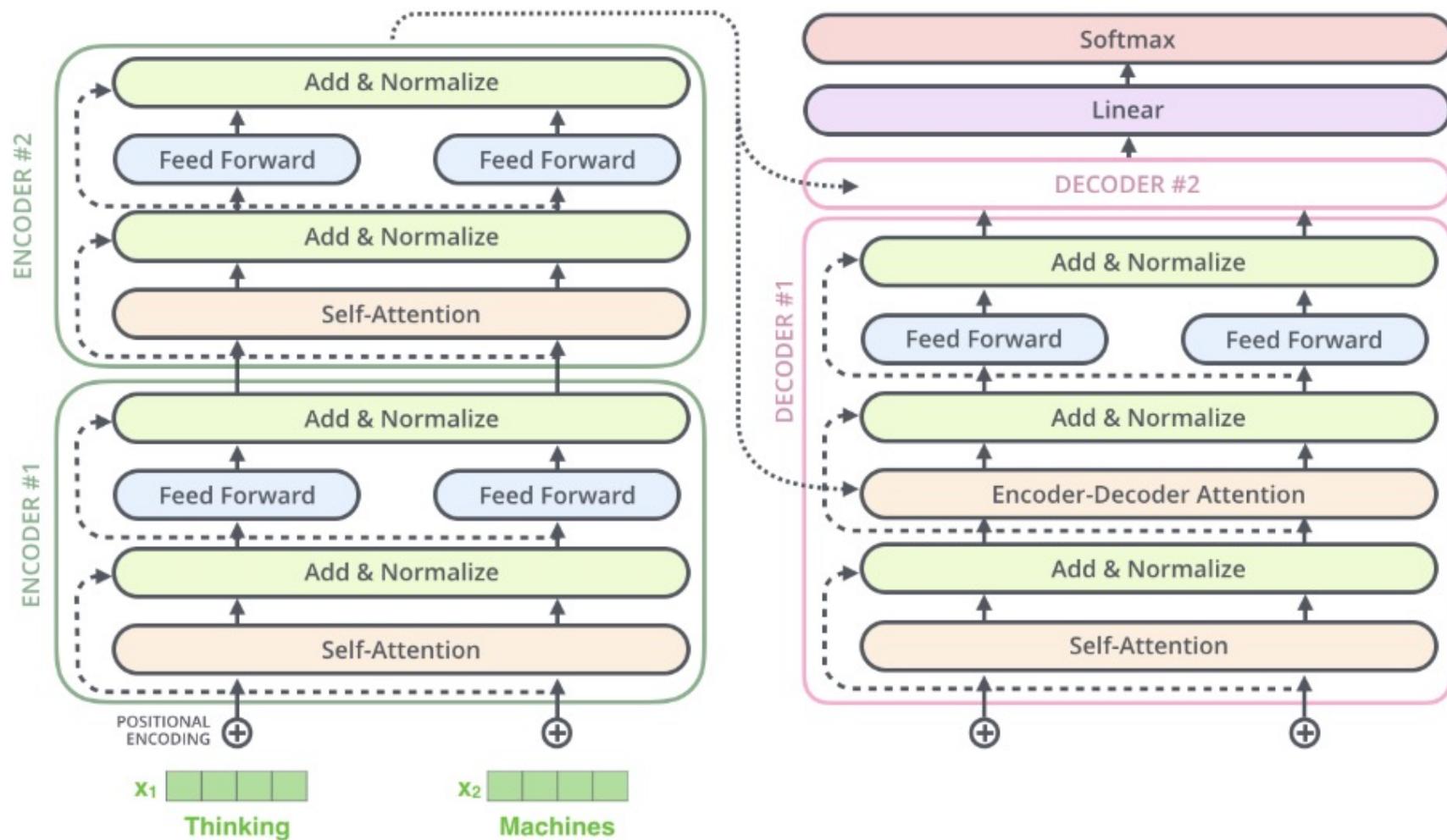
```
class SublayerConnection(nn.Module):
    """
    A residual connection followed by a layer norm.
    Note for code simplicity the norm is first as opposed to last.
    """
    def __init__(self, size, dropout):
        super(SublayerConnection, self).__init__()
        self.norm = LayerNorm(size)
        self.dropout = nn.Dropout(dropout)

    def forward(self, x, sublayer):
        "Apply residual connection to any sublayer with the same size."
        return x + self.dropout(sublayer(self.norm(x)))
```

```
class LayerNorm(nn.Module):
    "Construct a layernorm module (See citation for details)."
    def __init__(self, features, eps=1e-6):
        super(LayerNorm, self).__init__()
        self.a_2 = nn.Parameter(torch.ones(features))
        self.b_2 = nn.Parameter(torch.zeros(features))
        self.eps = eps

    def forward(self, x):
        mean = x.mean(-1, keepdim=True)
        std = x.std(-1, keepdim=True)
        return self.a_2 * (x - mean) / (std + self.eps) + self.b_2
```

# The Residuals

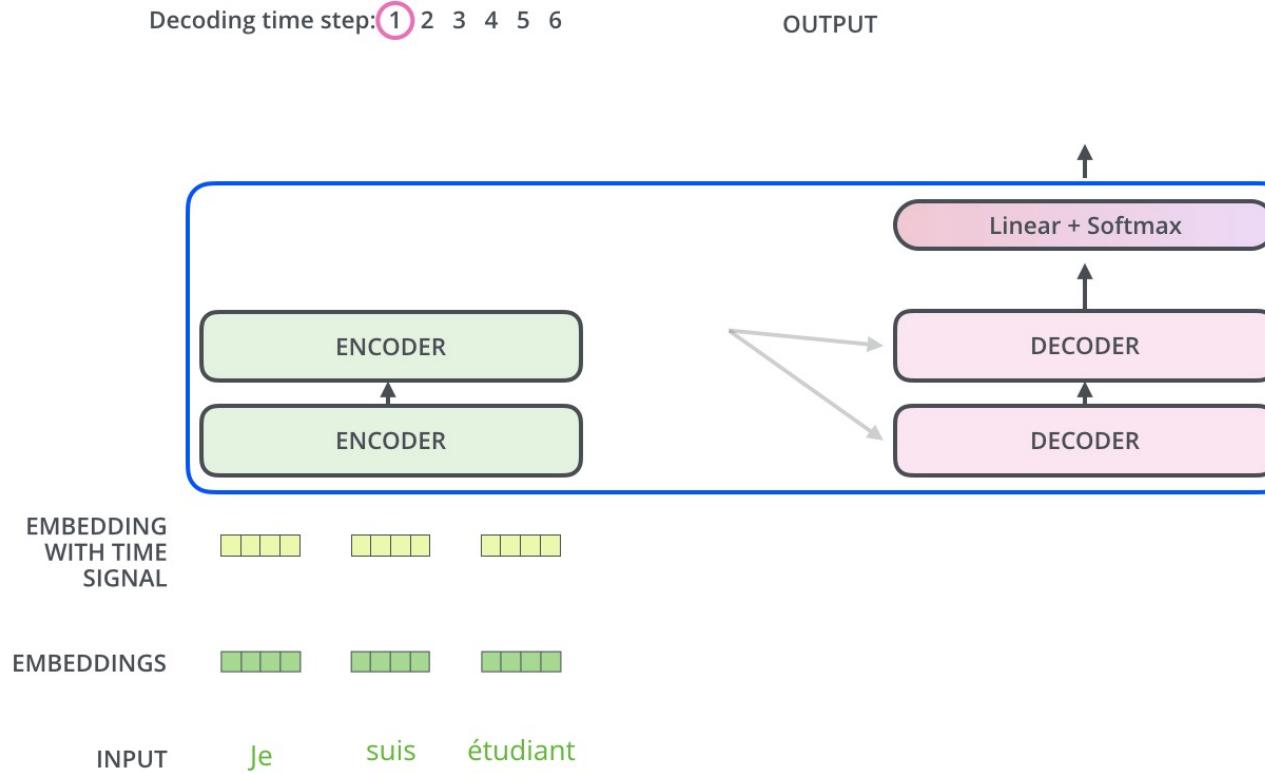


# The Residuals

```
class EncoderLayer(nn.Module):
    "Encoder is made up of self-attn and feed forward (defined below)"
    def __init__(self, size, self_attn, feed_forward, dropout):
        super(EncoderLayer, self).__init__()
        self.self_attn = self_attn
        self.feed_forward = feed_forward
        self.sublayer = clones(SublayerConnection(size, dropout), 2)
        self.size = size

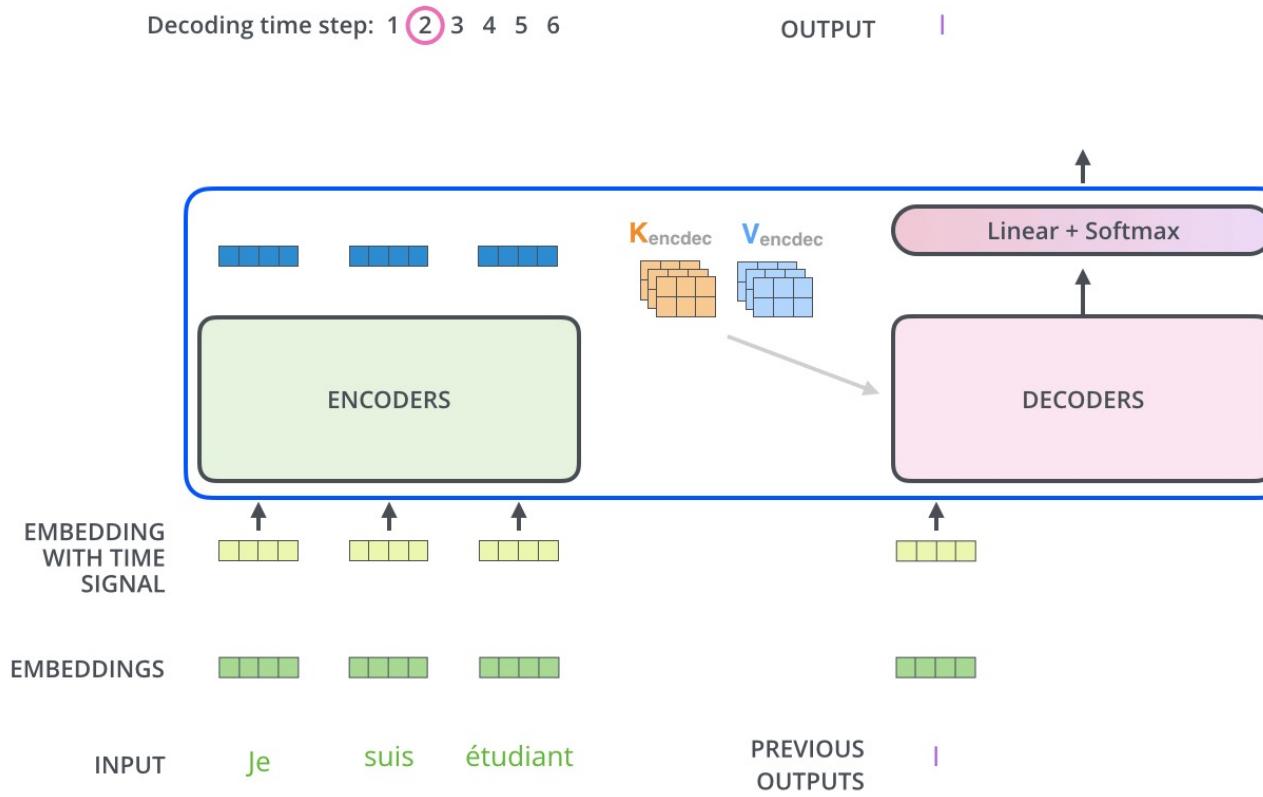
    def forward(self, x, mask):
        "Follow Figure 1 (left) for connections."
        x = self.sublayer[0](x, lambda x: self.self_attn(x, x, x, mask))
        return self.sublayer[1](x, self.feed_forward)
```

# The Decoder Side



After finishing the encoding phase, we begin the decoding phase. Each step in the decoding phase outputs an element from the output sequence (the English translation sentence in this case).

# The Decoder Side



# The Decoder

```
class Decoder(nn.Module):
    "Generic N layer decoder with masking."
    def __init__(self, layer, N):
        super(Decoder, self).__init__()
        self.layers = clones(layer, N)
        self.norm = LayerNorm(layer.size)

    def forward(self, x, memory, src_mask, tgt_mask):
        for layer in self.layers:
            x = layer(x, memory, src_mask, tgt_mask)
        return self.norm(x)
```

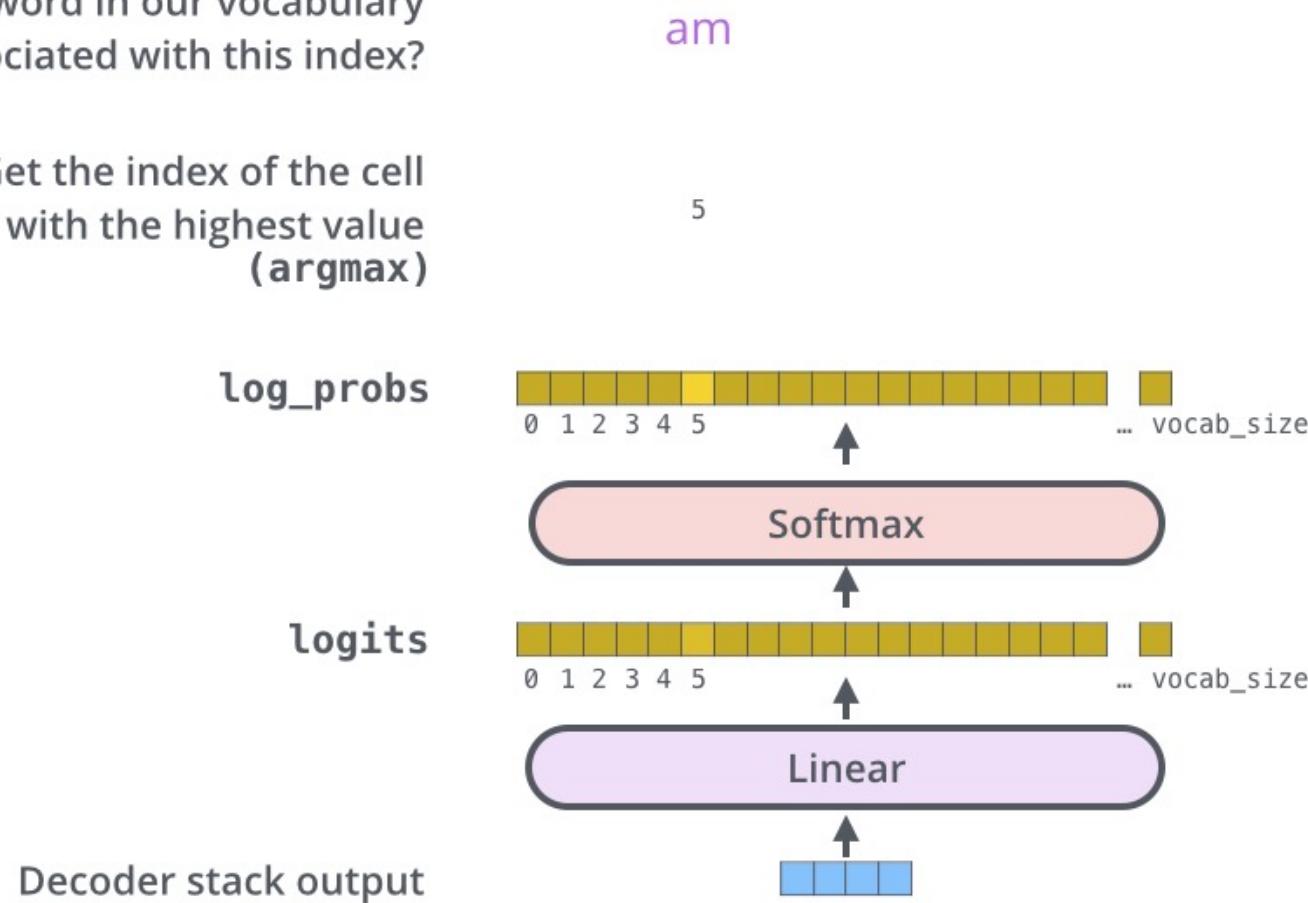
```
class DecoderLayer(nn.Module):
    "Decoder is made of self-attn, src-attn, and feed forward (defined below)"
    def __init__(self, size, self_attn, src_attn, feed_forward, dropout):
        super(DecoderLayer, self).__init__()
        self.size = size
        self.self_attn = self_attn
        self.src_attn = src_attn
        self.feed_forward = feed_forward
        self.sublayer = clones(SublayerConnection(size, dropout), 3)

    def forward(self, x, memory, src_mask, tgt_mask):
        "Follow Figure 1 (right) for connections."
        m = memory
        x = self.sublayer[0](x, lambda x: self.self_attn(x, x, x, tgt_mask))
        x = self.sublayer[1](x, lambda x: self.src_attn(x, m, m, src_mask))
        return self.sublayer[2](x, self.feed_forward)
```

# The First Linear and Softmax Layer

Which word in our vocabulary  
is associated with this index?

Get the index of the cell  
with the highest value  
(`argmax`)



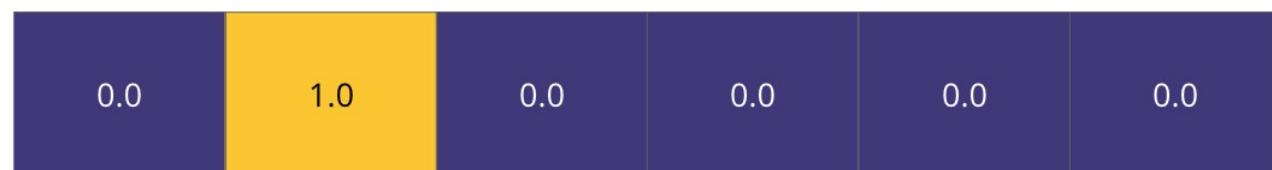
This figure starts from the bottom with the vector produced as the output of the decoder stack. It is then turned into an output word.

# Training Procedure

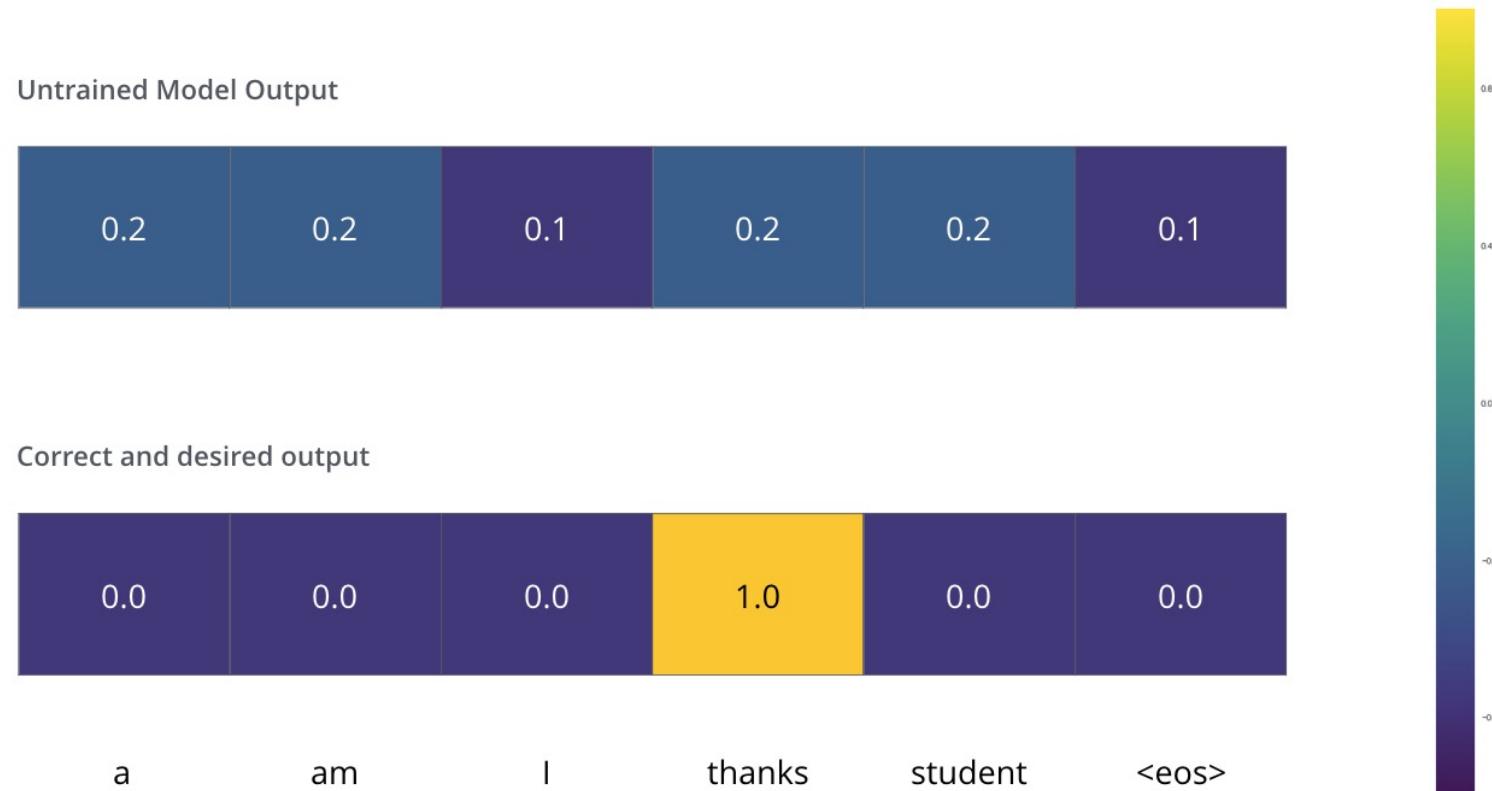
Output Vocabulary

WORD	a	am	I	thanks	student	<eos>
INDEX	0	1	2	3	4	5

One-hot encoding of the word "am"



# The Loss Function



Since the model's parameters (weights) are all initialized randomly, the (untrained) model produces a probability distribution with arbitrary values for each cell/word. We can compare it with the actual output, then tweak all the model's weights using backpropagation to make the output closer to the desired output.

# The Loss Function

## Target Model Outputs

Output Vocabulary: a am I thanks student <eos>



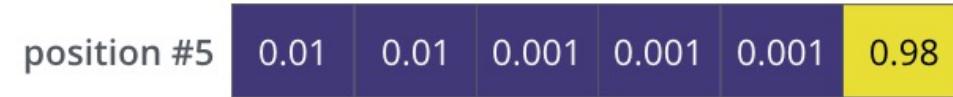
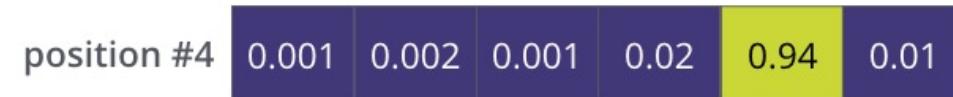
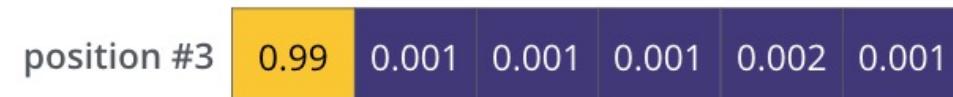
a am I thanks student <eos>



# The Loss Function

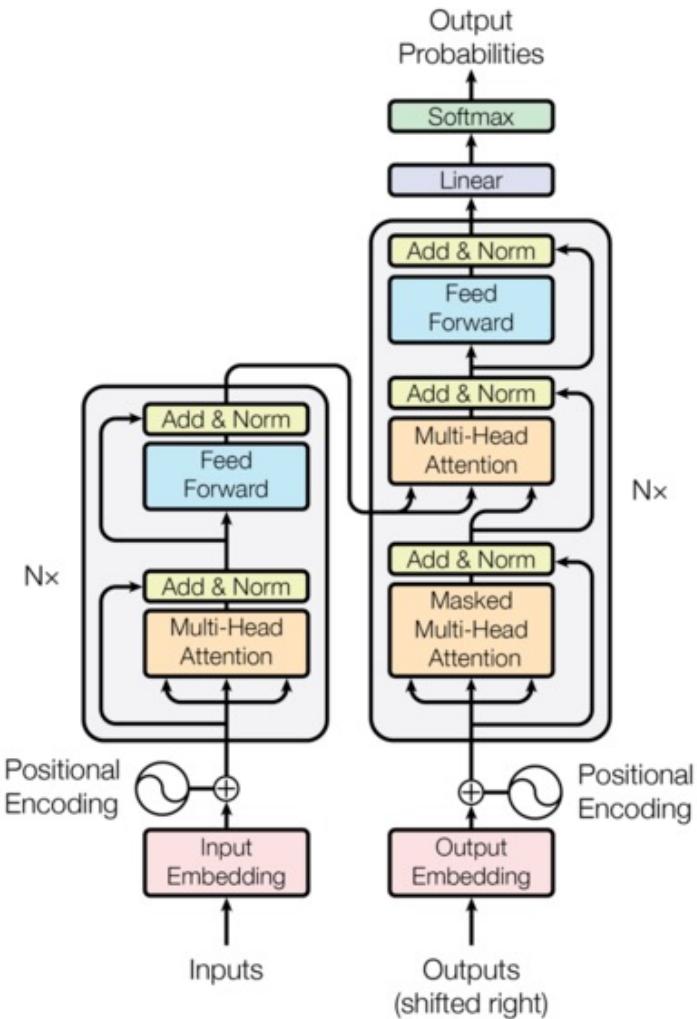
## Trained Model Outputs

Output Vocabulary: a am I thanks student <eos>



a am I thanks student <eos>

# Wrap up!



# References

- The Illustrated Transformer
  - <http://jalammar.github.io/illustrated-transformer/>
- The Annotated Transformer
  - <http://nlp.seas.harvard.edu/2018/04/03/attention.html>
- Transformer Neural Networks - EXPLAINED! (Attention is all you need)
  - <https://www.youtube.com/watch?v=TQQlZhbC5ps>

# ELMo



el Architecture

state-of-the-art neural language models (Jozefowicz et al., 2016; Melis et al., 2017) compute a context-independent token representation  $x_t^{LM}$  (via tokens or a CNN over characters) which passes through  $L$  layers of forward LSTM. At position  $k$ , each LSTM layer outputs a context-dependent representation  $\tilde{x}_t^{LM}$ . The top-layer LSTM output,  $\tilde{x}_{t,k}^{LM}$ , is used to predict the next token  $t_{k+1}$  with a feed-forward layer. ELMo is similar to a forward LM, except it runs over the sequence in reverse, token given the future context.

$$\prod_{k=1}^N p(t_k \mid t_1, t_2, \dots, t_{k-1}) = p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k \mid t_{k+1}, t_{k+2}, \dots, t_N)$$

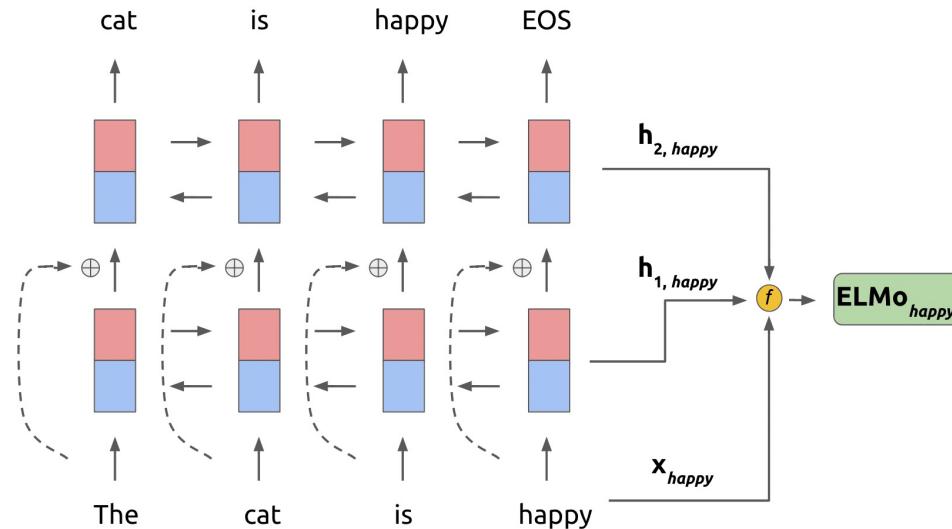
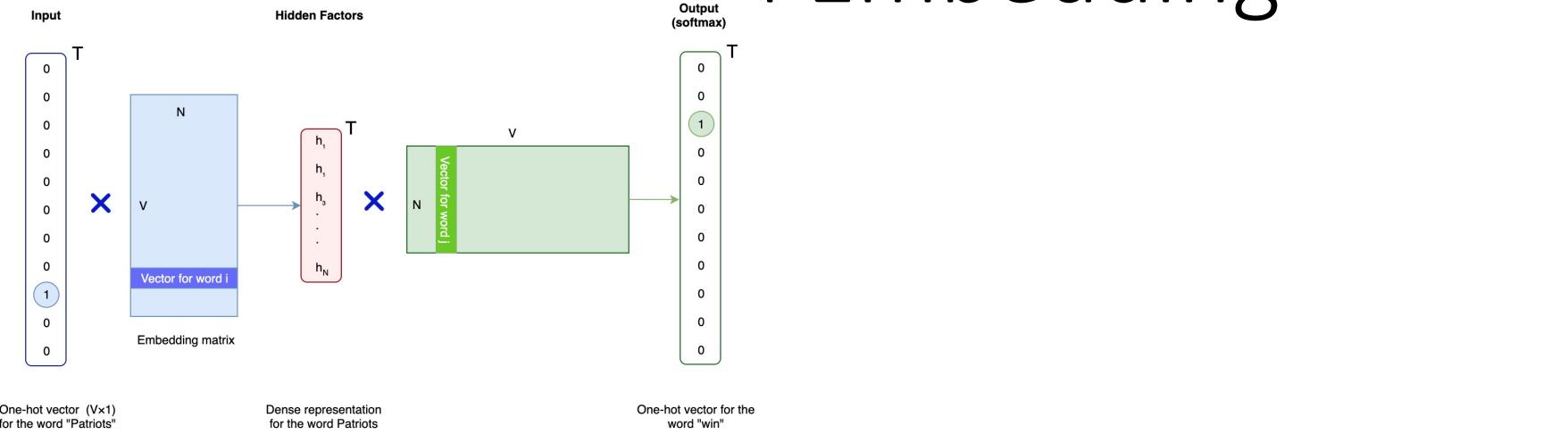
pp\_16.mp4

Deep  
Conte  
Word  
Repre

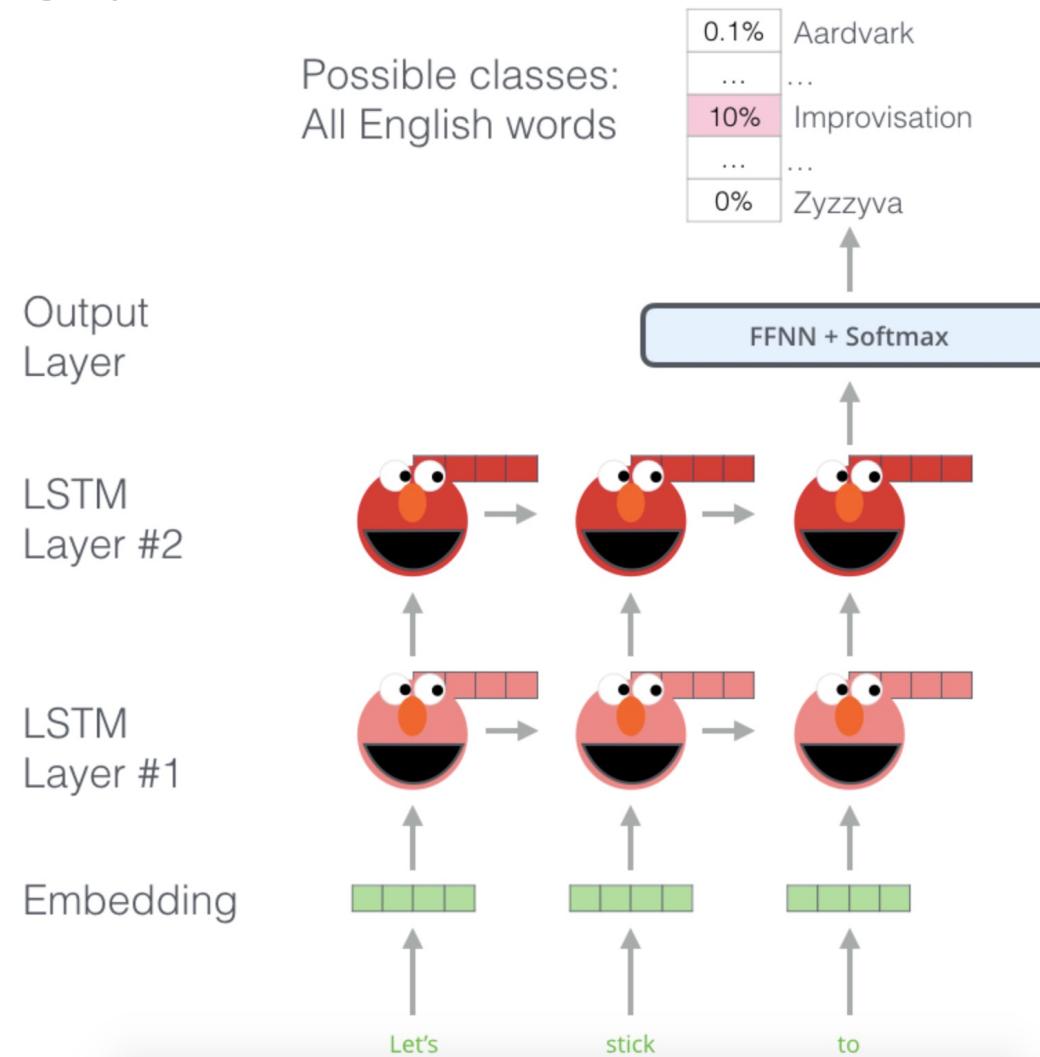
2019521112  
DONGYUP SH

pp\_16.pdf

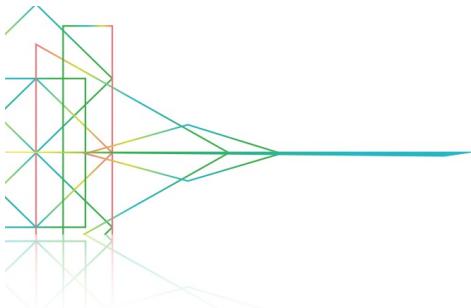
# “Contextualized” Word Embedding



# (Pre-)Training procedure of ELMo



# (Pre-)Training procedure of ELMo



ELMo is a task specific representation. A down-stream task learns weighting parameters

$$\text{ELMo}_k^{\text{task}} = \gamma^{\text{task}} \times \sum \left\{ \begin{array}{l} s_2^{\text{task}} \times h_{k2}^{\text{LM}} \\ s_1^{\text{task}} \times h_{k1}^{\text{LM}} \\ s_0^{\text{task}} \times h_{k0}^{\text{LM}} \end{array} \right. \quad ([x_k; x_k])$$

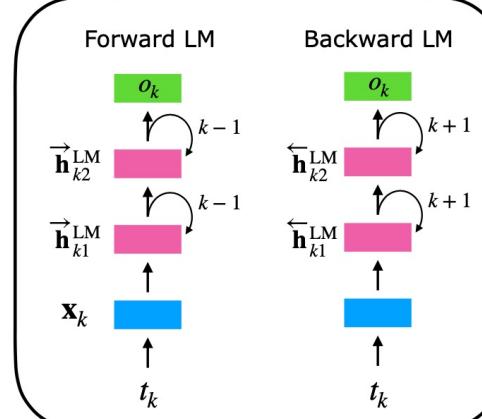
Concatenate hidden layers  
 $[\vec{h}_{kj}^{\text{LM}}; \overleftarrow{h}_{kj}^{\text{LM}}]$

Unlike usual word embeddings, ELMo is assigned to every *token* instead of a *type*

## Method

ELMo represents a word  $t_k$  as a linear combination of corresponding hidden layers (inc. its embedding)

### biLMs



# How to use ELMo when after pre-training

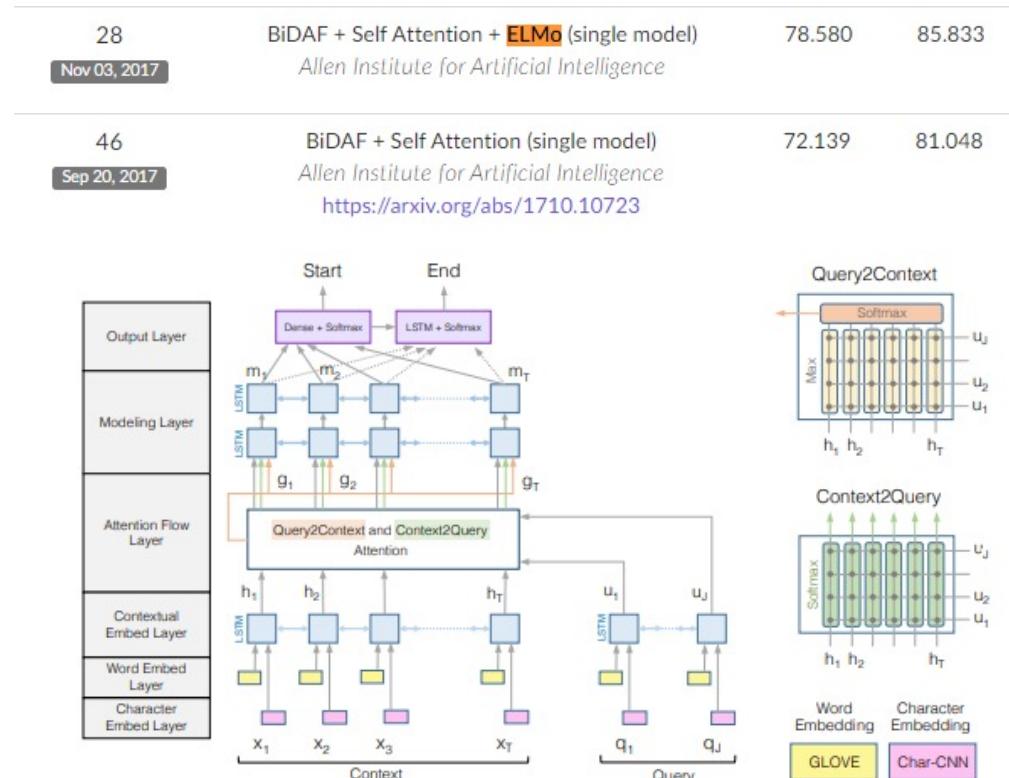


Figure 1: BiDirectional Attention Flow Model (best viewed in color)

# How to use ELMo when after pre-training

Task	Previous SOTA		Our baseline	ELMo + Baseline	Increase (Absolute/Relative)
SQuAD	SAN	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%
Sentiment (5-class)	McCann et al (2017)	53.7	51.4	54.7 +/- 0.5	3.3 / 6.8%

# Deep Pre-trained Language Models

- ~~2017 Jun: Transformer (not LM)~~

## Attention Is All You Need

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, Illia Polosukhin

- 2018 Feb: ELMo

## Deep contextualized word representations

Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, Luke Zettlemoyer

- 2018 Jun: OpenAI GPT

---

## Improving Language Understanding by Generative Pre-Training

---

- 2018 Oct: BERT

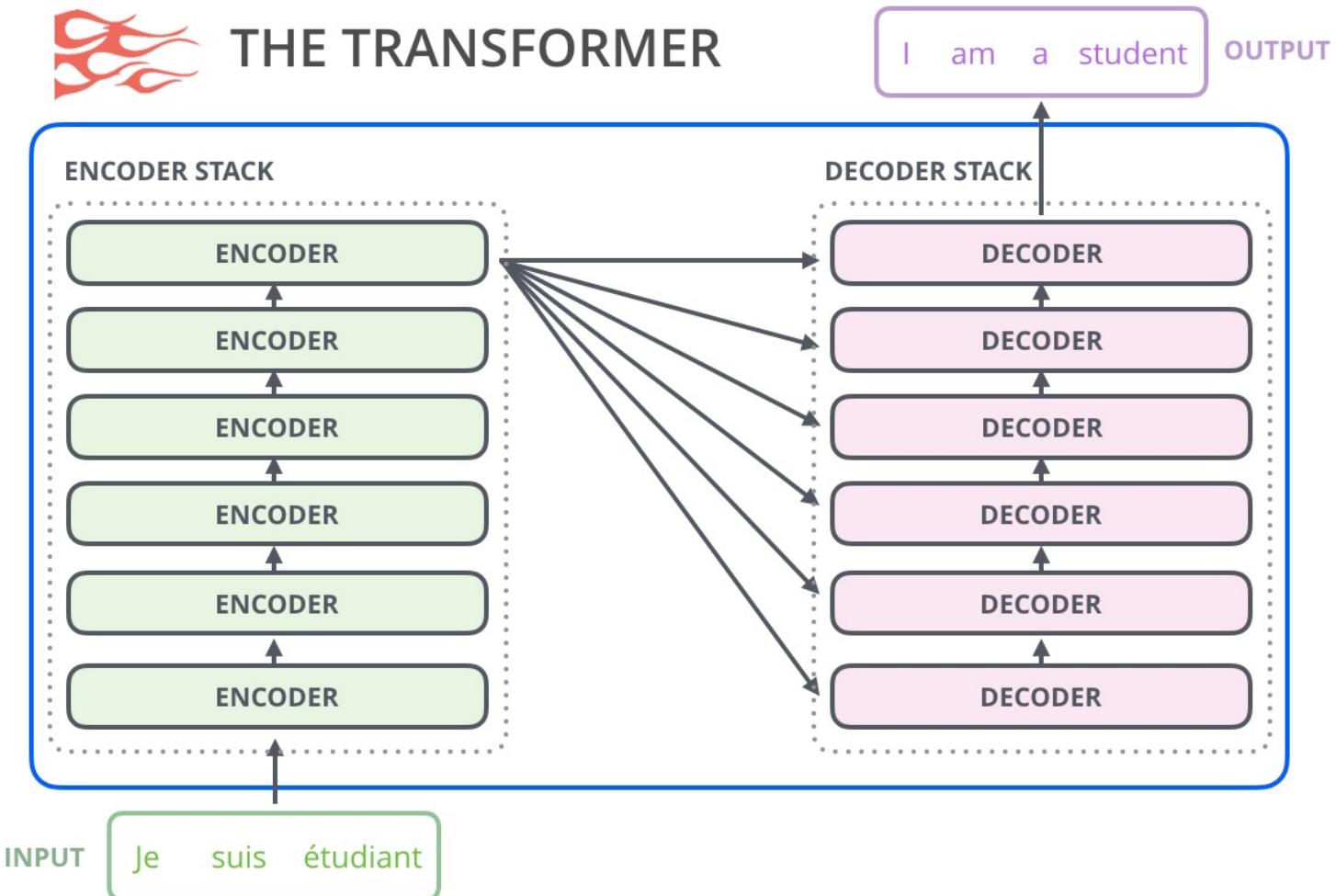
## BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova

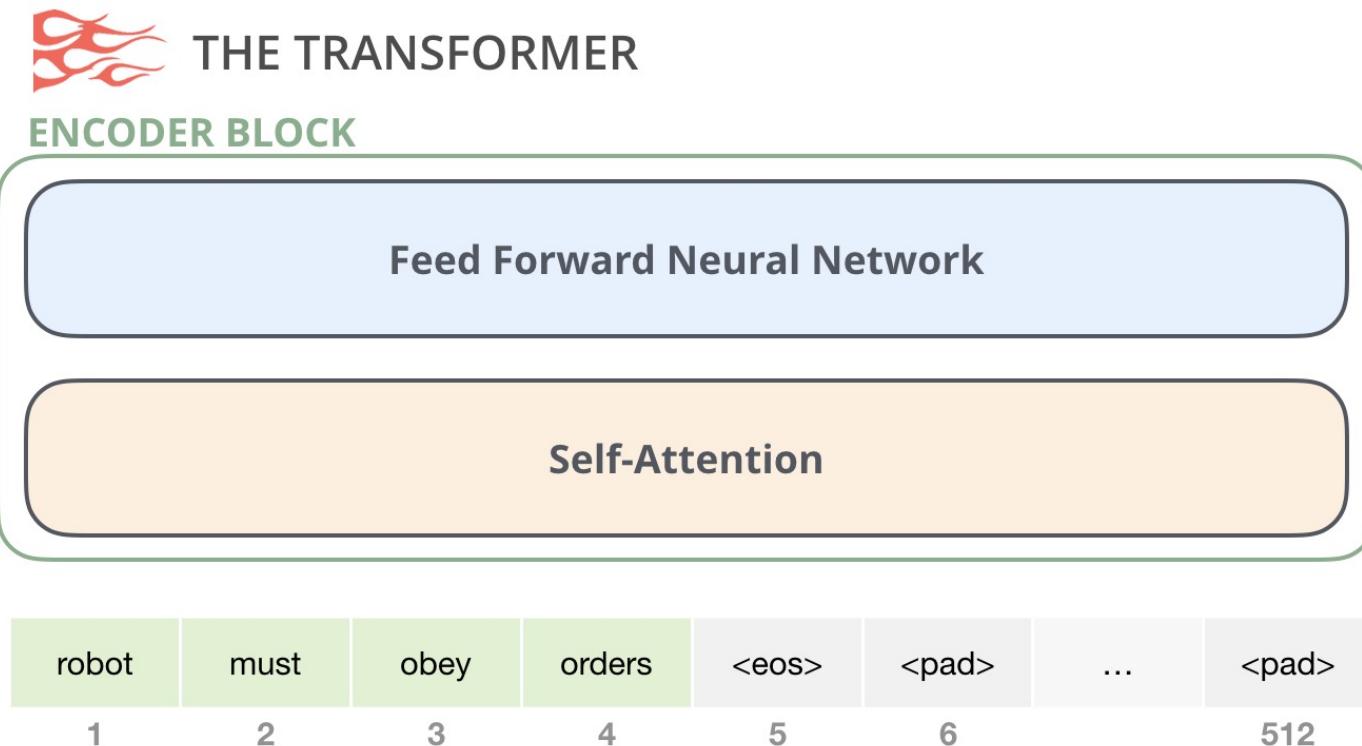
- 2019 Feb: OpenAI GPT-2

...

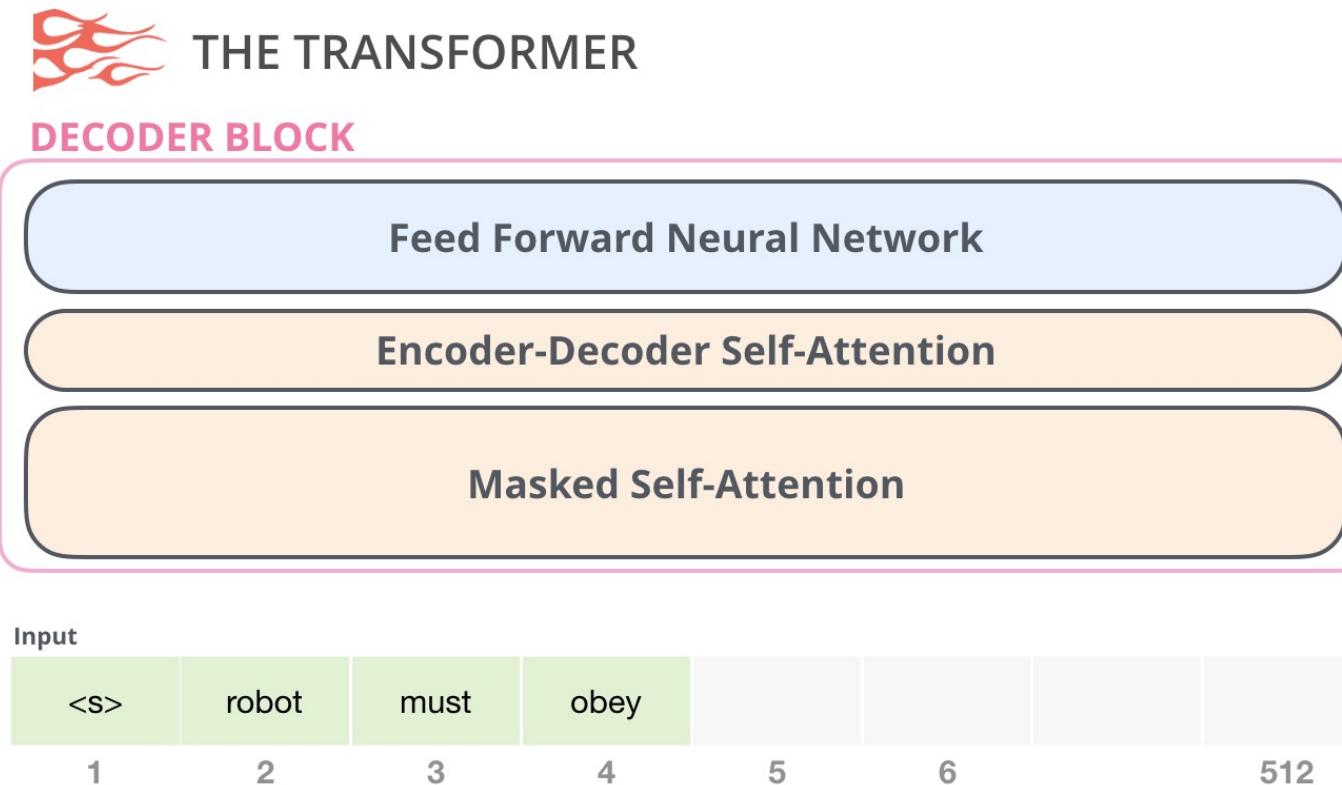
# (Recap) Transformer



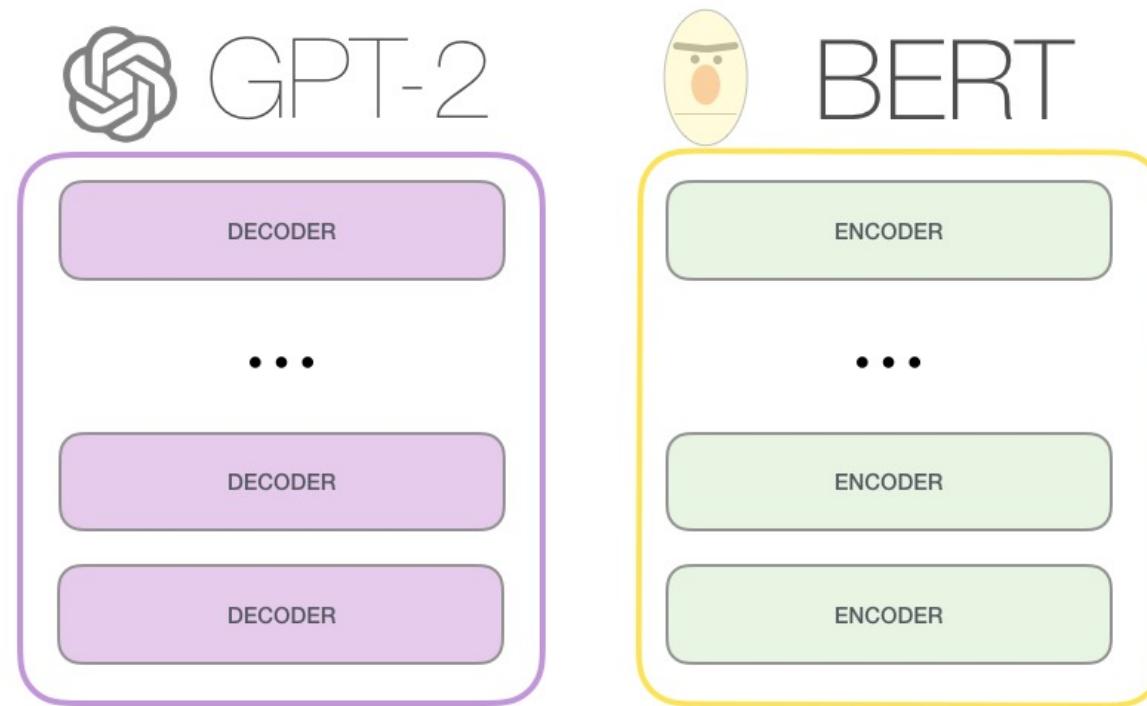
# (Recap) Encoder Block



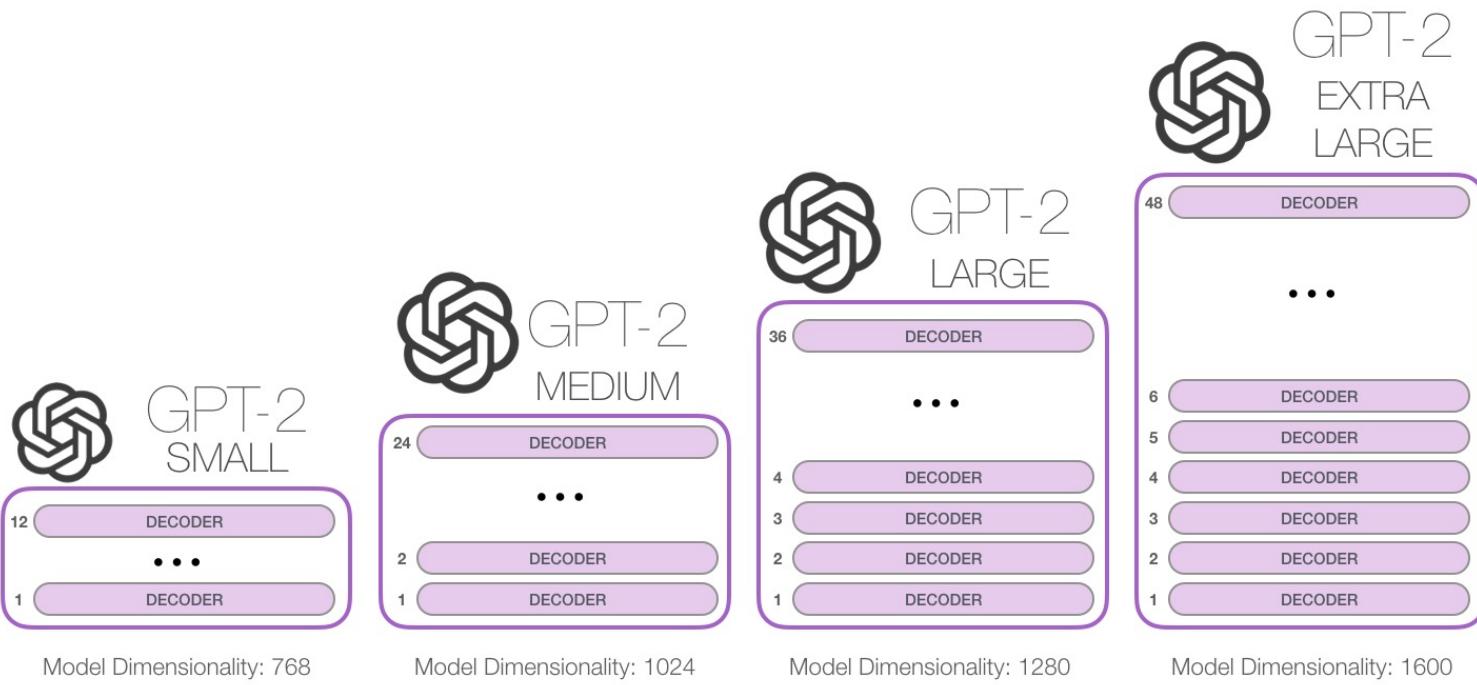
# (Recap) Decoder Block



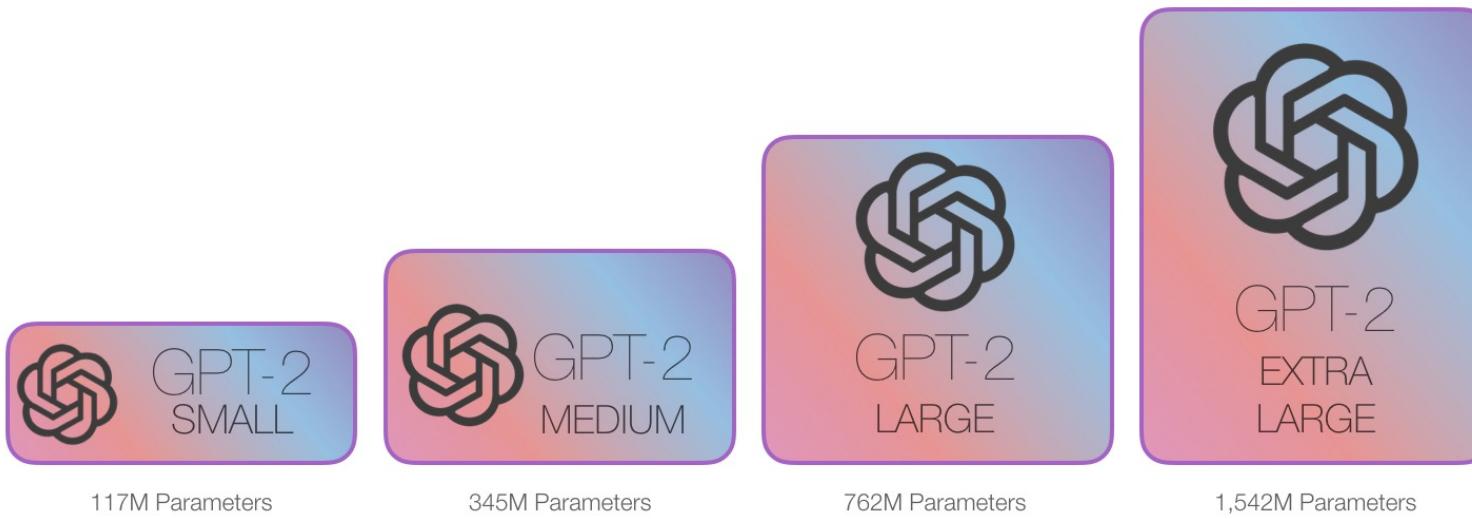
# Transformer-based Language Models



# How high can we stack up these blocks?



# How high can we stack up these blocks?



# GPT: Generative Pre-trained Transformer



---

## Improving Language Understanding by Generative Pre-Training

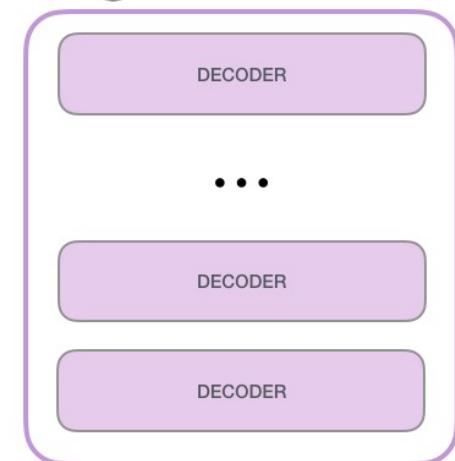
---

Alec Radford  
OpenAI  
[alec@openai.com](mailto:alec@openai.com)

Karthik Narasimhan  
OpenAI  
[karthikn@openai.com](mailto:karthikn@openai.com)

Tim Salimans  
OpenAI  
[tim@openai.com](mailto:tim@openai.com)

Ilya Sutskever  
OpenAI  
[ilyasu@openai.com](mailto:ilyasu@openai.com)

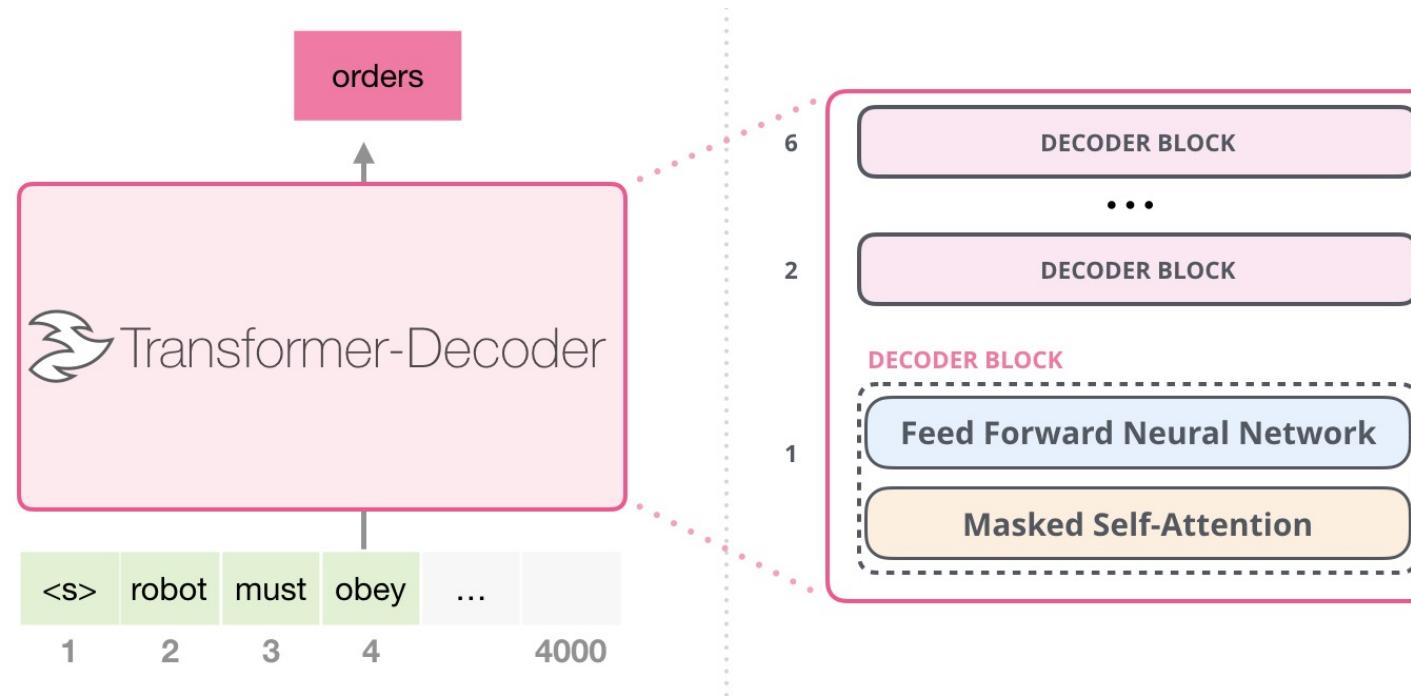


# The Decoder-Only Block

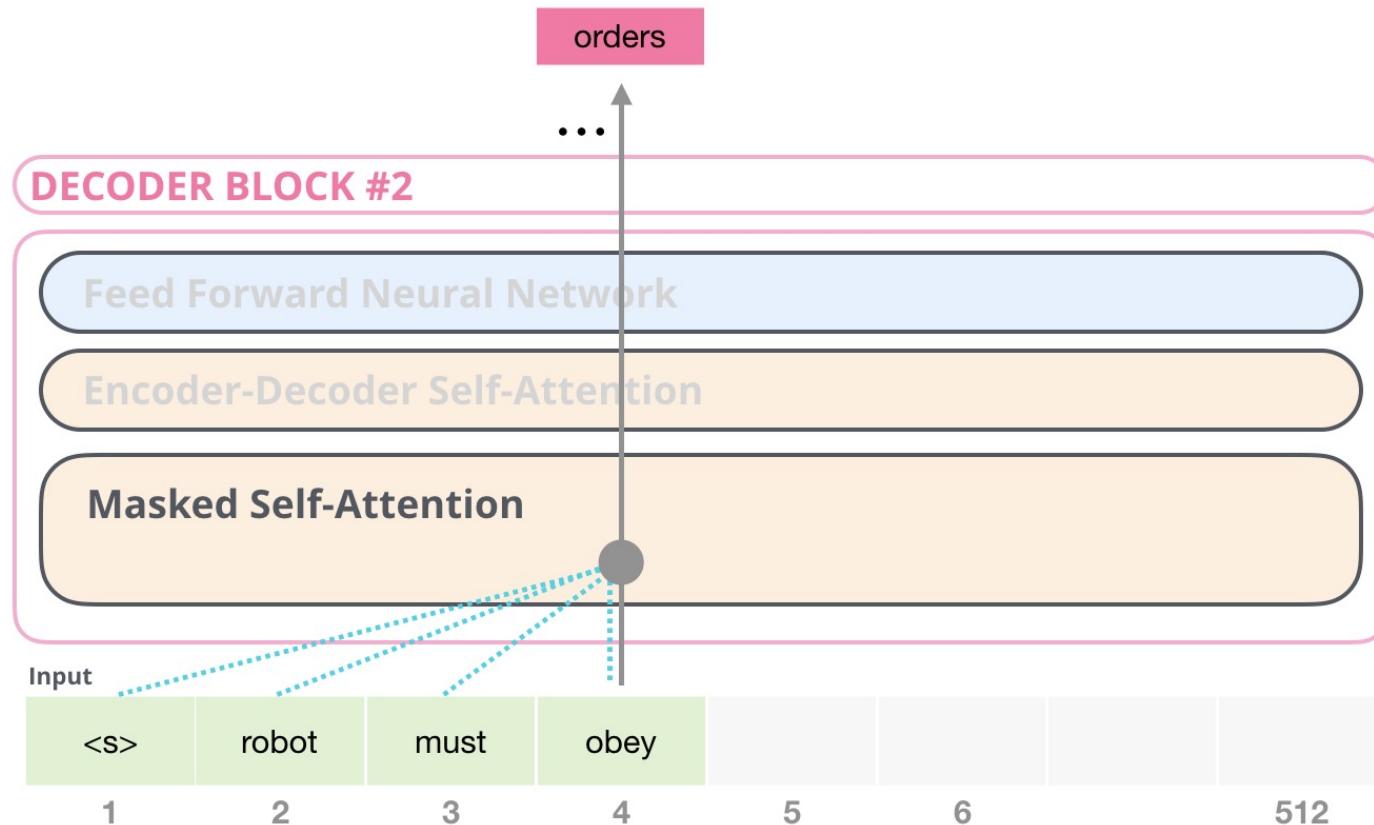
Published as a conference paper at ICLR 2018

## GENERATING WIKIPEDIA BY SUMMARIZING LONG SEQUENCES

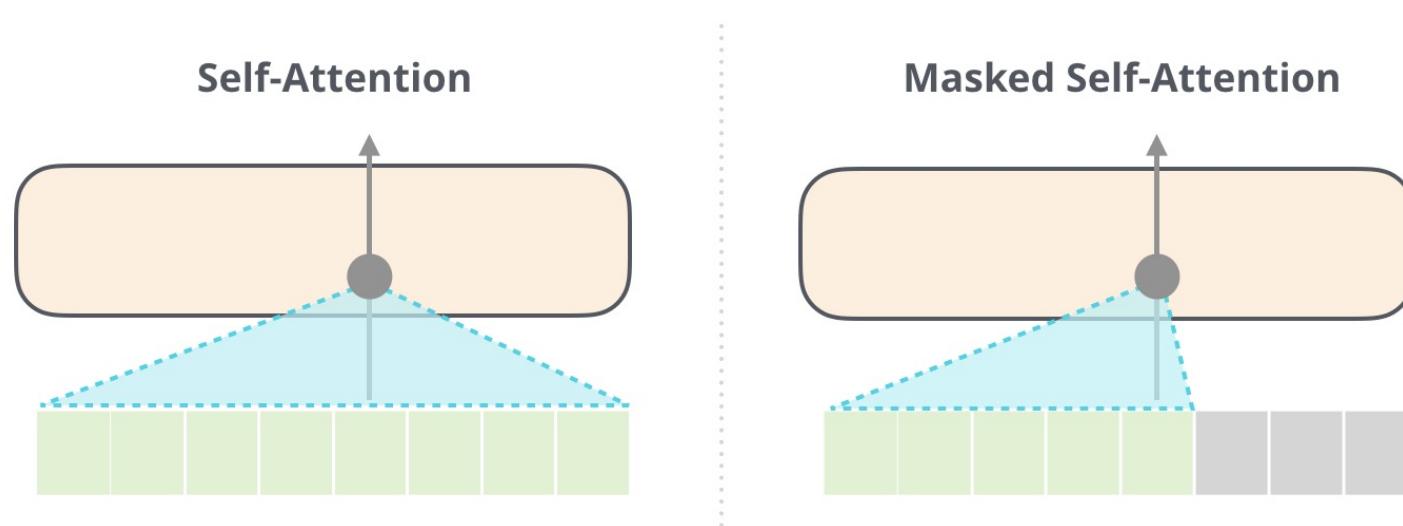
Peter J. Liu\*, Mohammad Saleh\*,  
Etienne Pot†, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, Noam Shazeer  
Google Brain  
Mountain View, CA  
`{peterjliu,msaleh,epot,bgoodrich,rsepassi,lukaszkaiser,noam}@google.com`



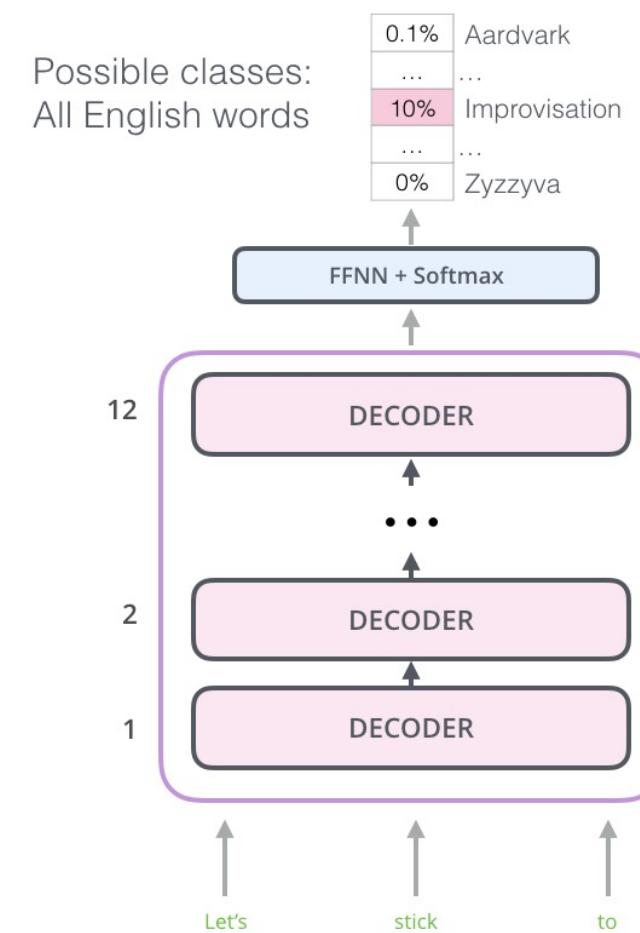
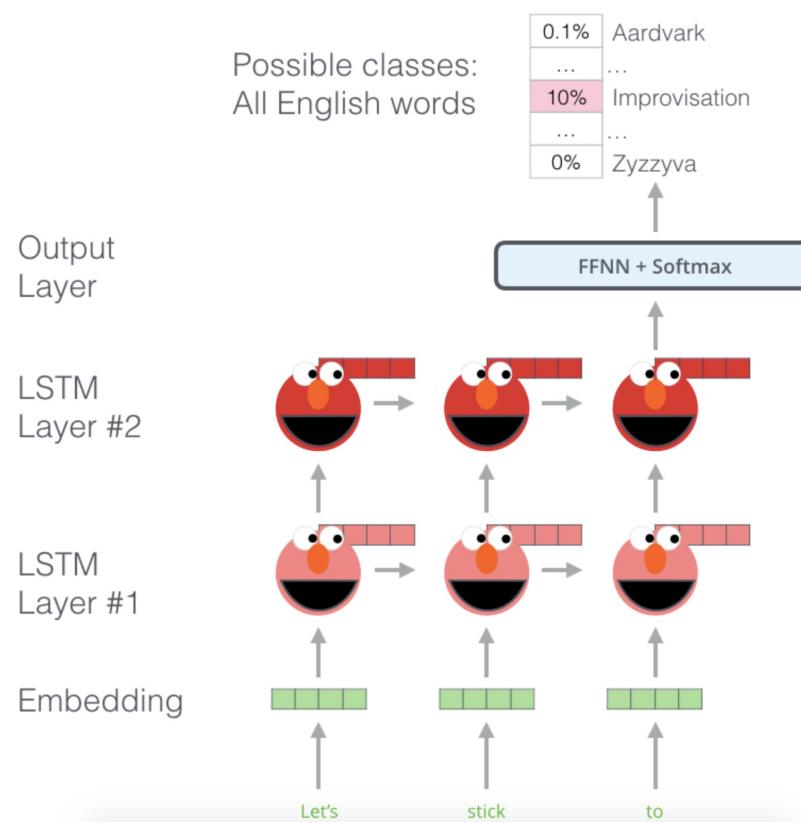
# (Recap) Decoder Block



# (Recap) Self-Attention

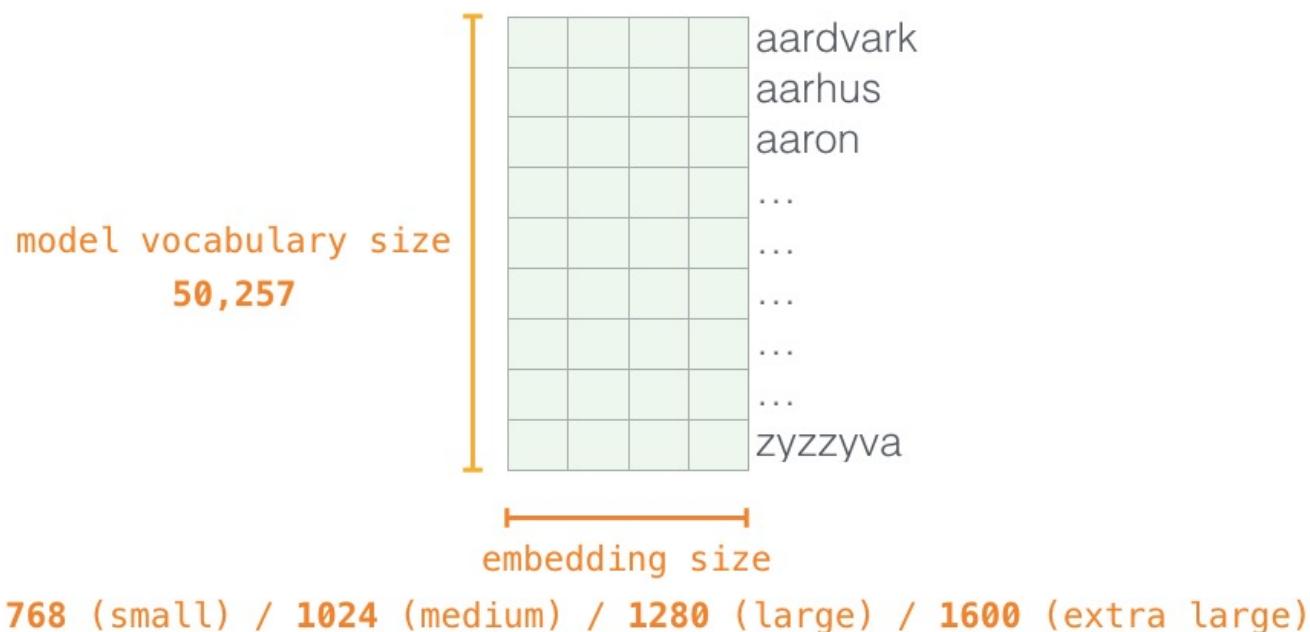


# Pre-training: Language modeling (unsupervised)



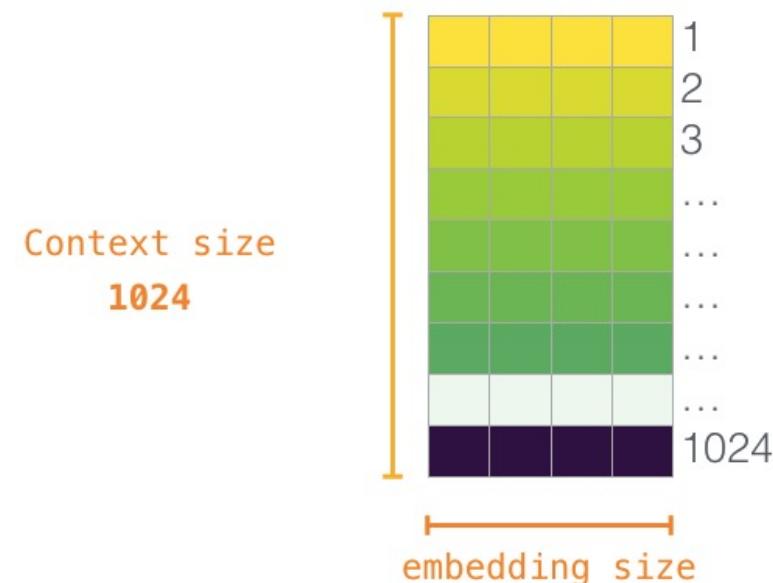
# Input Encoding in GPT

## Token Embeddings (wte)



# Input Encoding in GPT

## Positional Encodings (wpe)

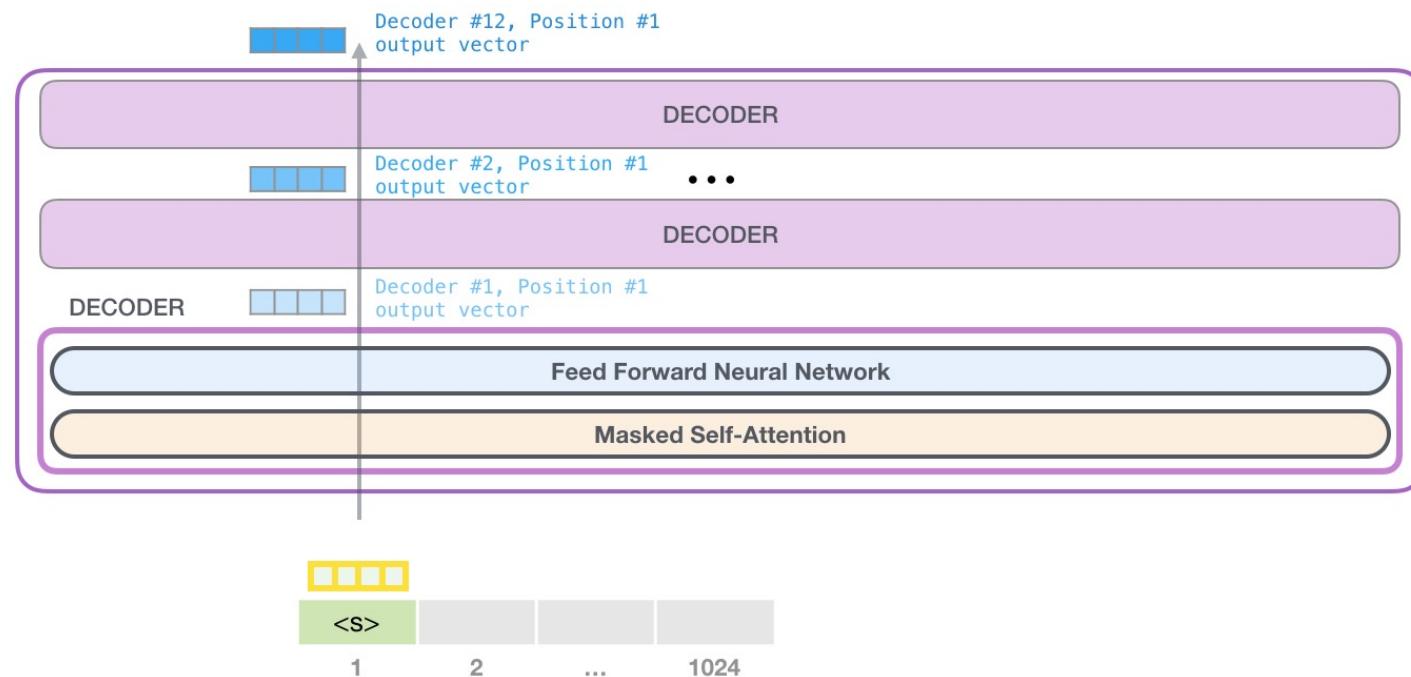


**768** (small) / **1024** (medium) / **1280** (large) / **1600** (extra large)

# Input Encoding in GPT



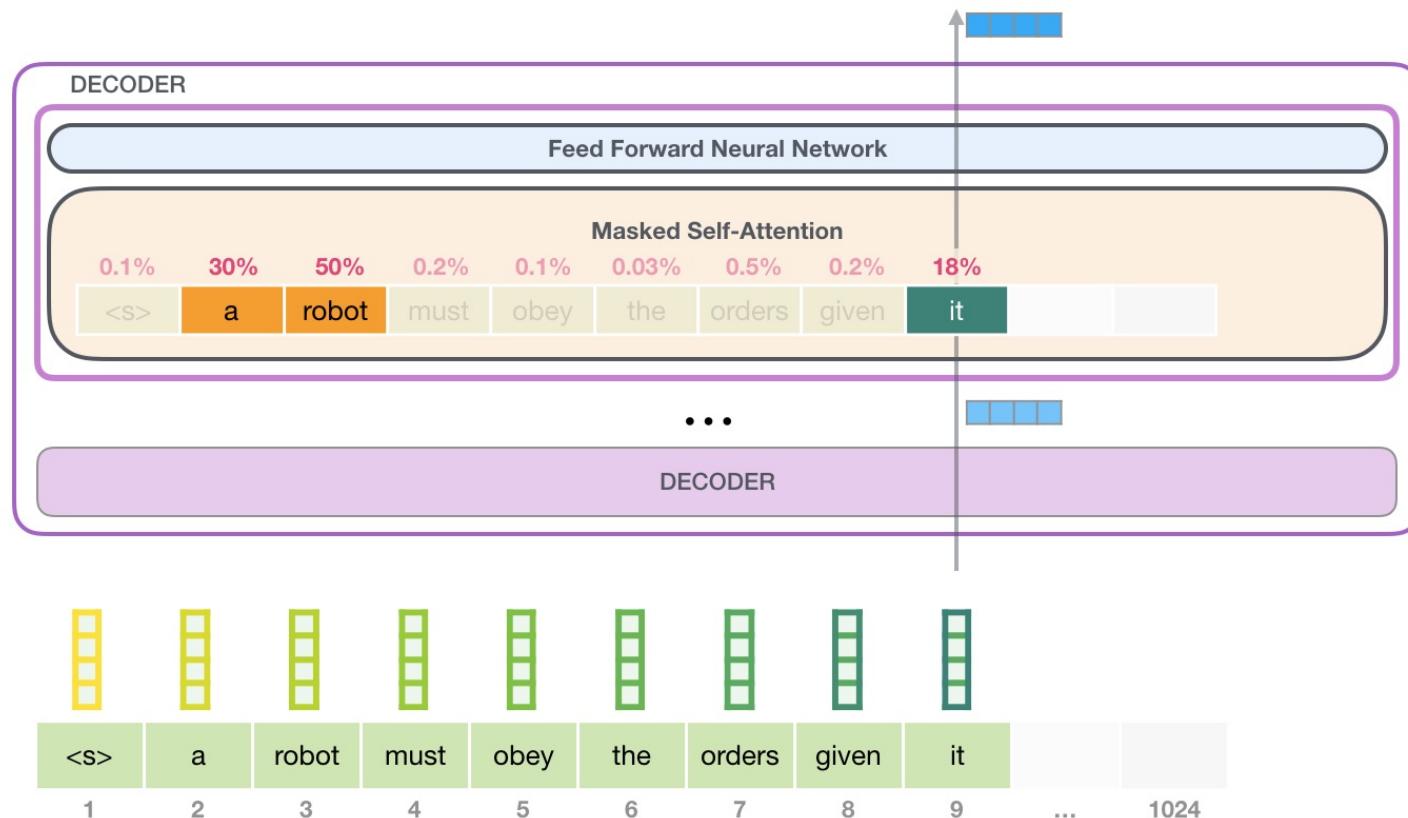
# Input Encoding in GPT



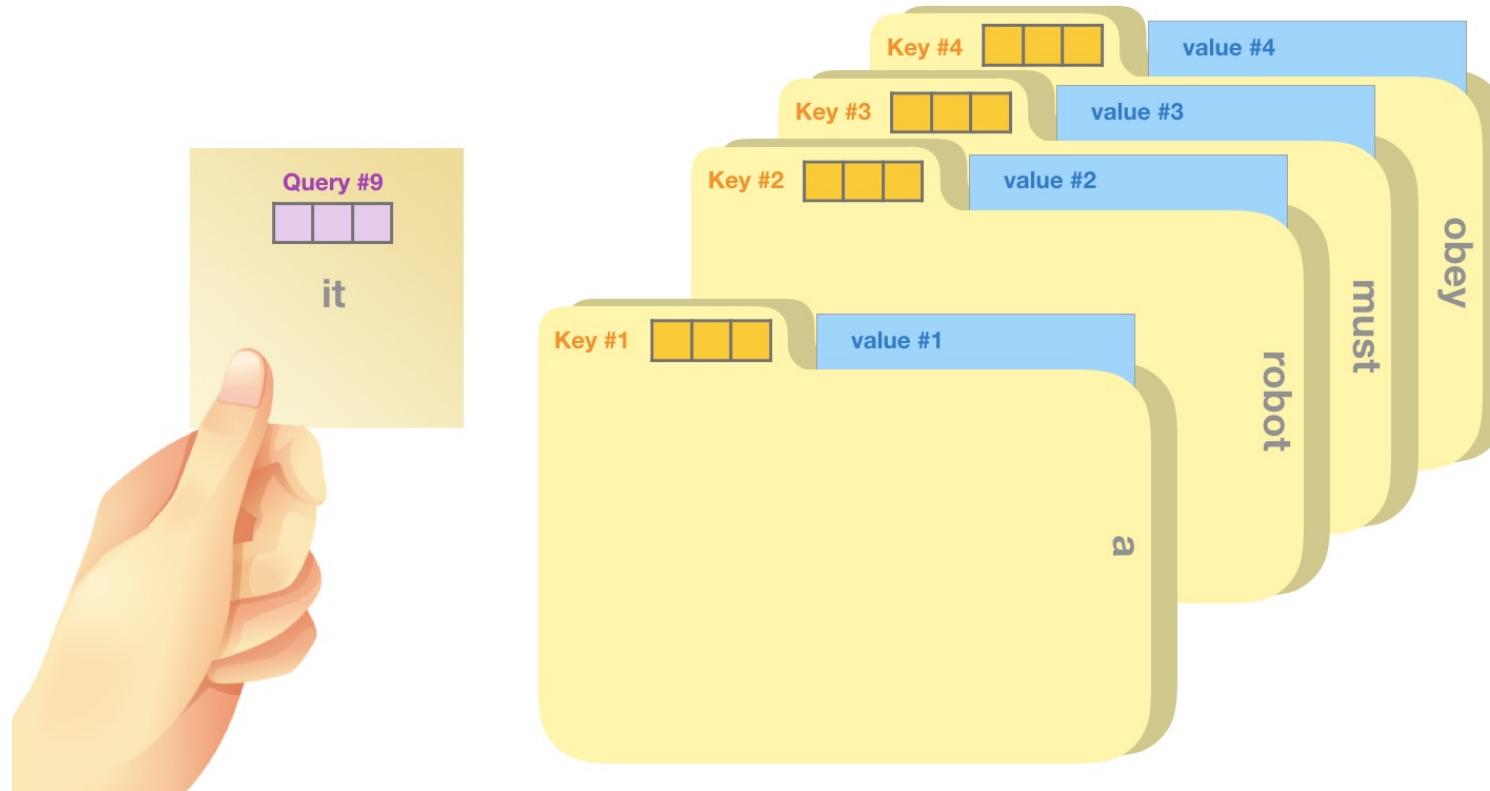
# Self-Attention in GPT

## Second Law of Robotics

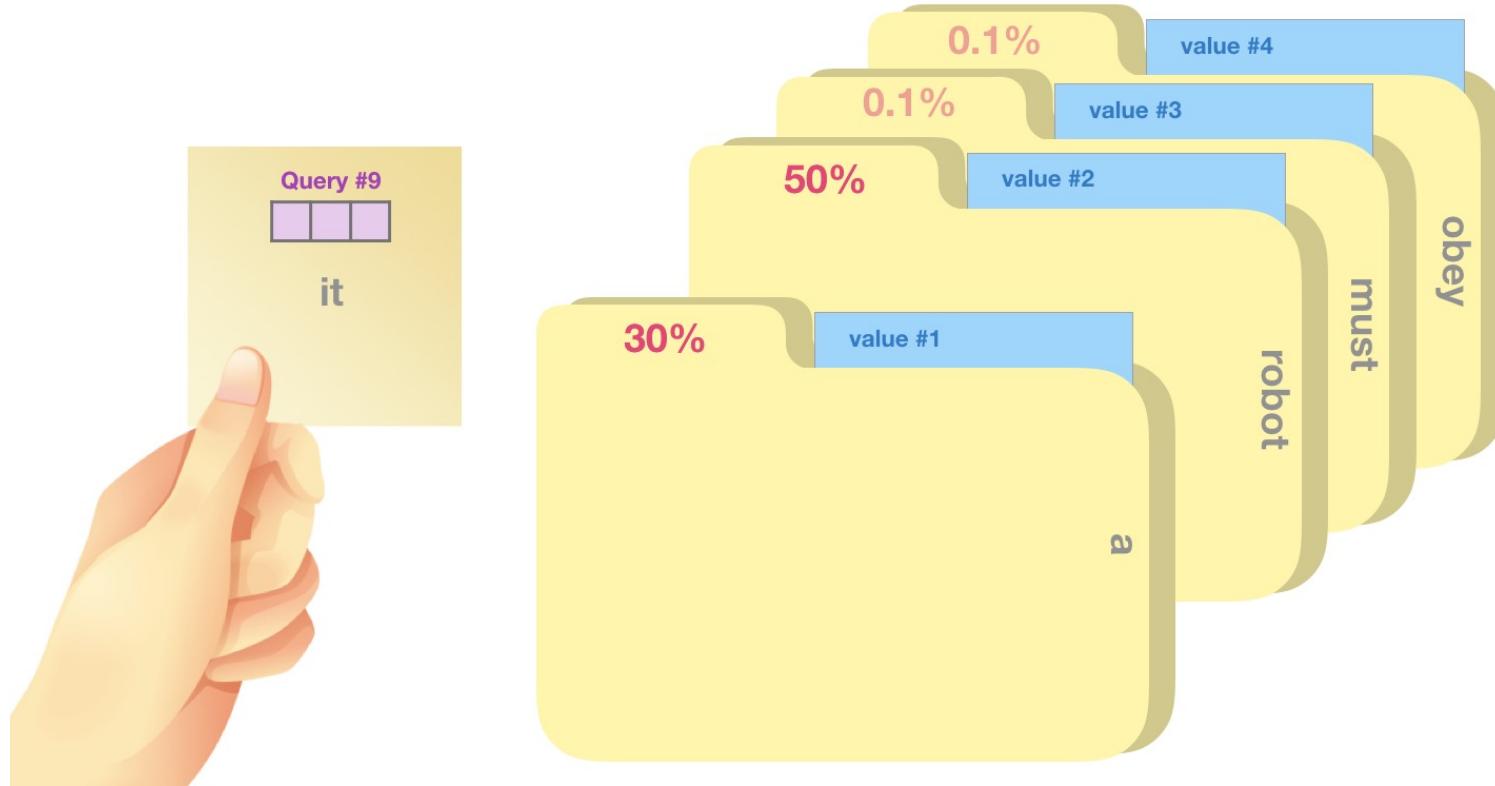
A robot must obey the orders given *it* by human beings except where **such orders** would conflict with the **First Law**.



# Self-Attention in GPT



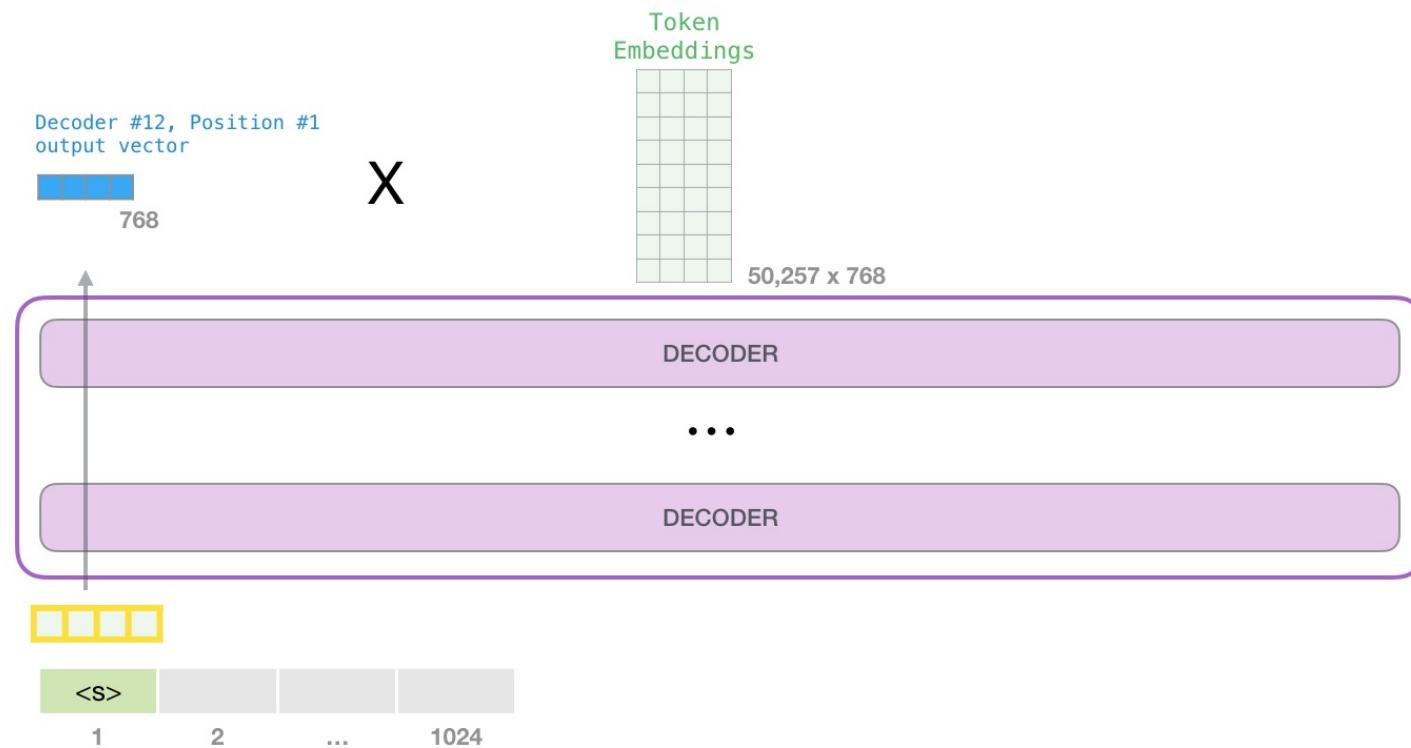
# Self-Attention in GPT



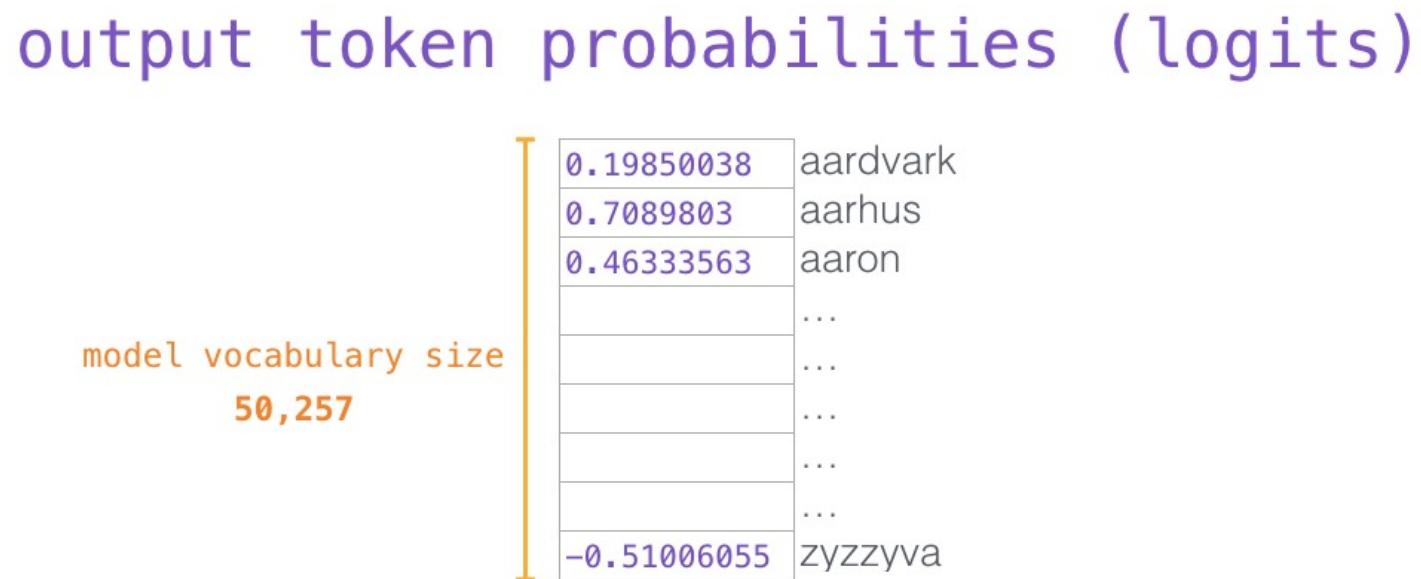
# Self-Attention in GPT

Word	Value vector	Score	Value X Score
<S>		0.001	
a		0.3	
robot		0.5	
must		0.002	
obey		0.001	
the		0.0003	
orders		0.005	
given		0.002	
it		0.19	
		Sum:	

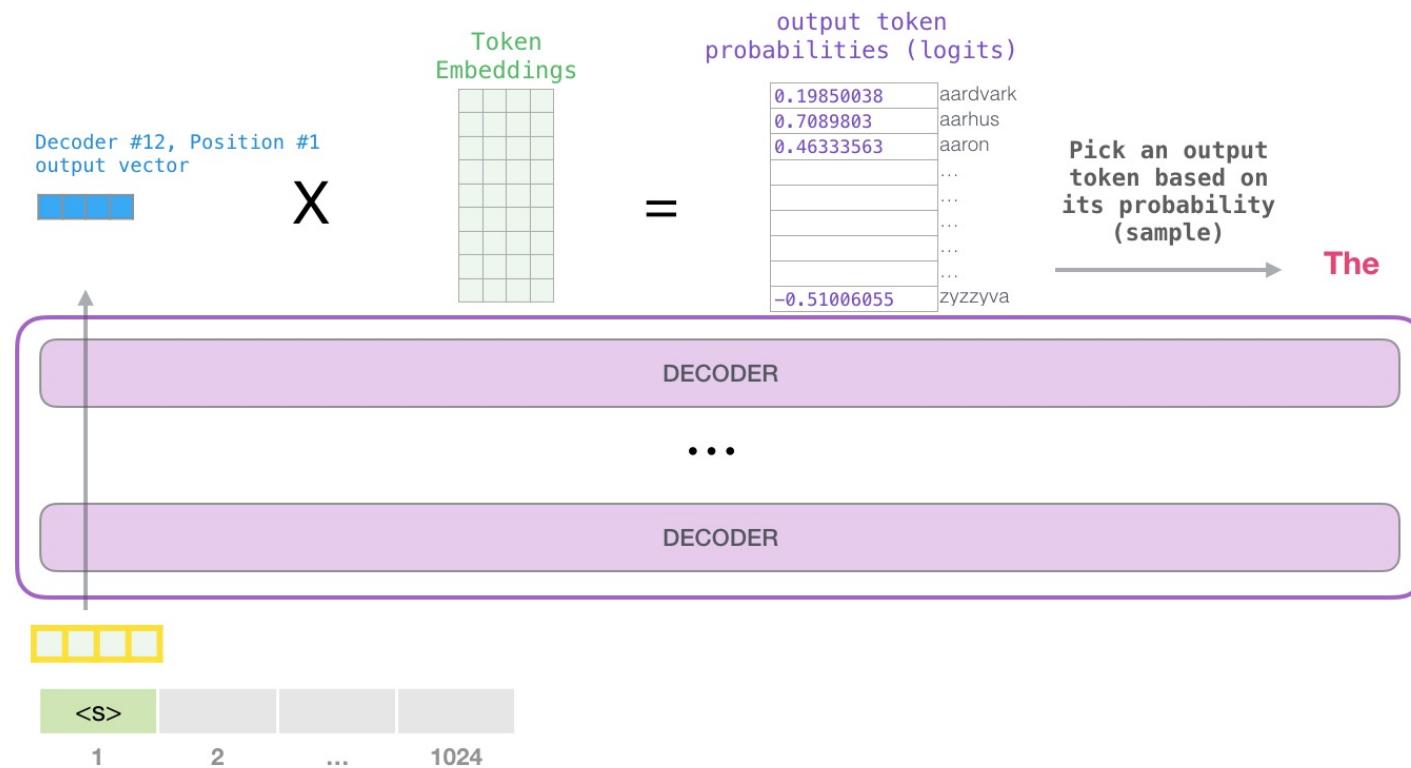
# Output in GPT



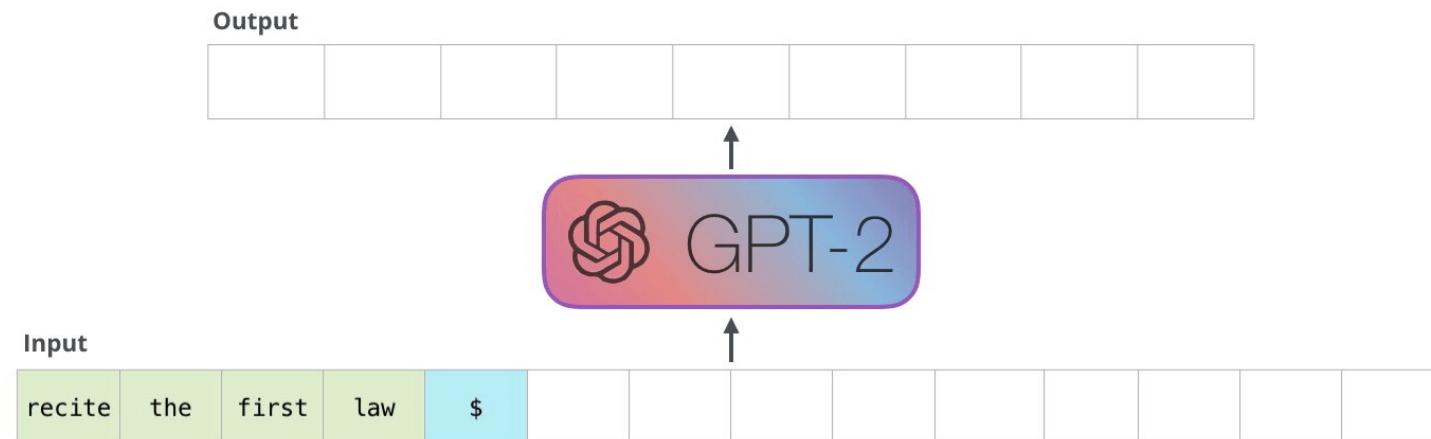
# Output in GPT



# Output in GPT



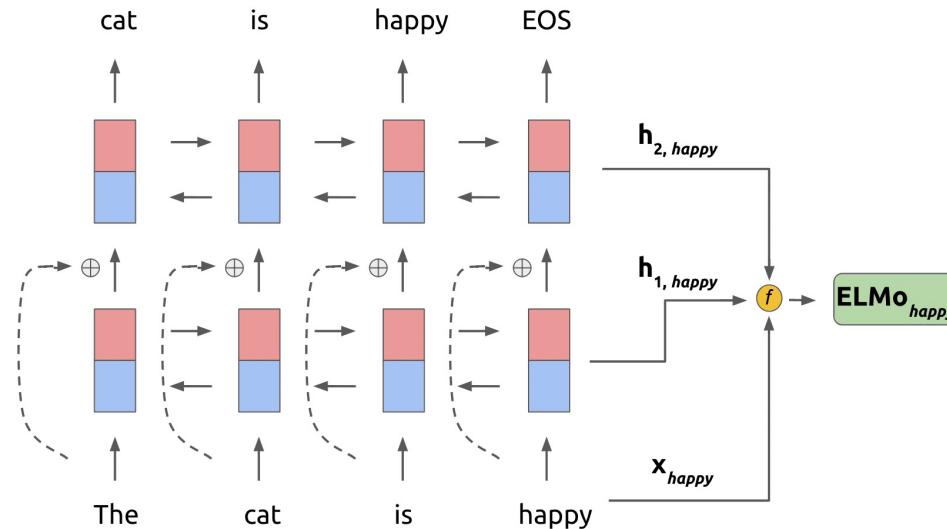
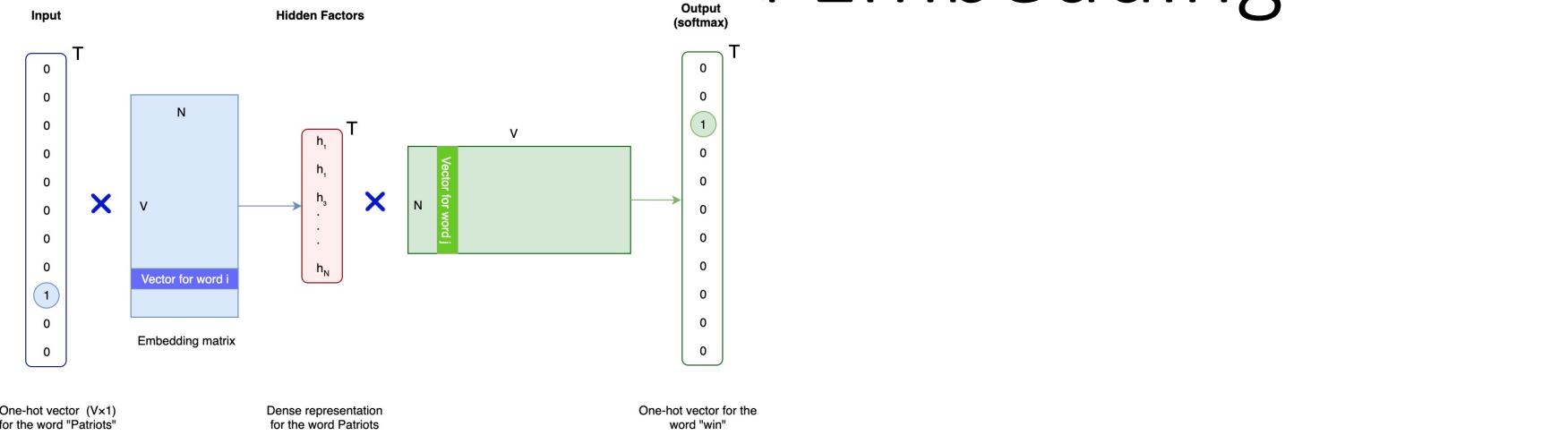
# Auto-regressive Model for Language Modeling



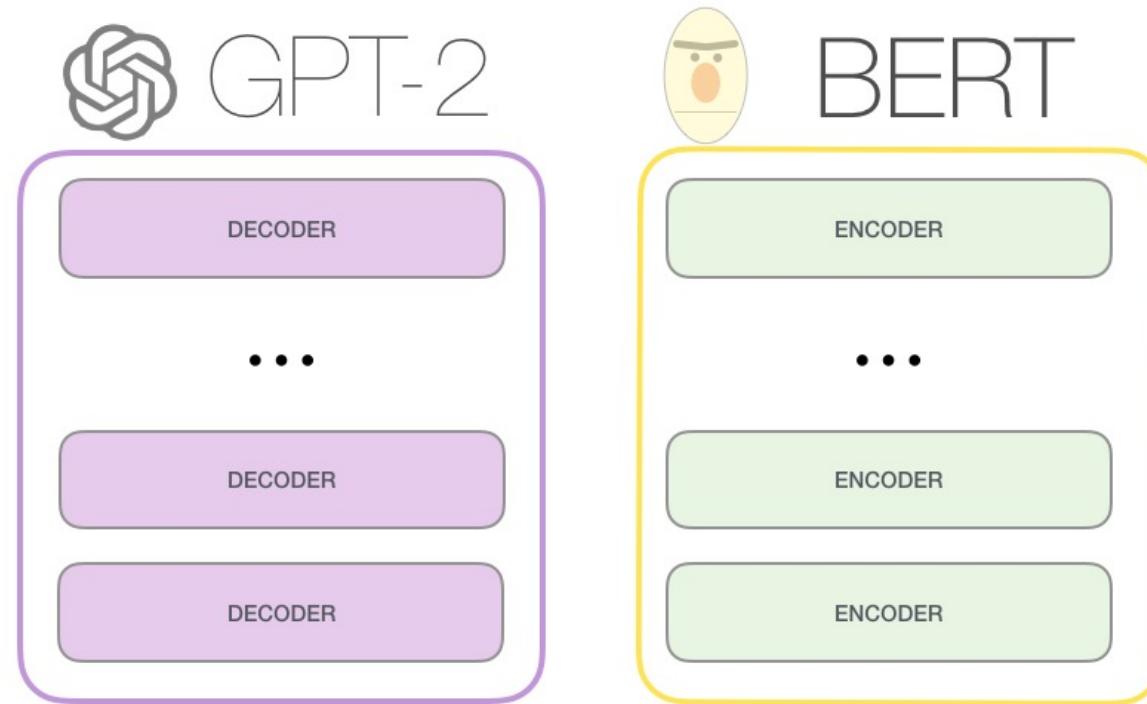
# References

- **The Illustrated BERT, ELMo, and co.**
  - <http://jalammar.github.io/illustrated-bert/>
- **The Illustrated GPT-2**
  - <http://jalammar.github.io/illustrated-gpt2/#part-1-got-and-language-modeling>

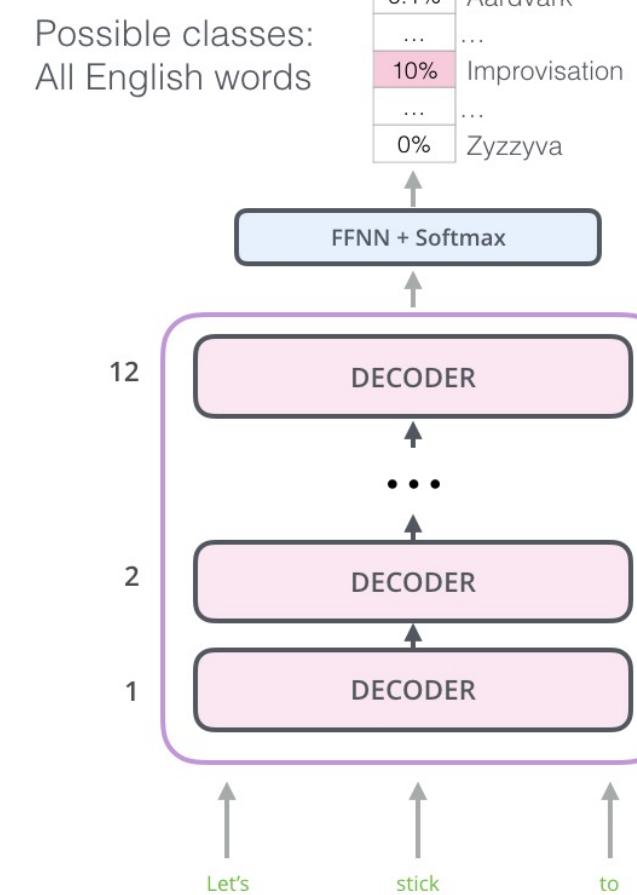
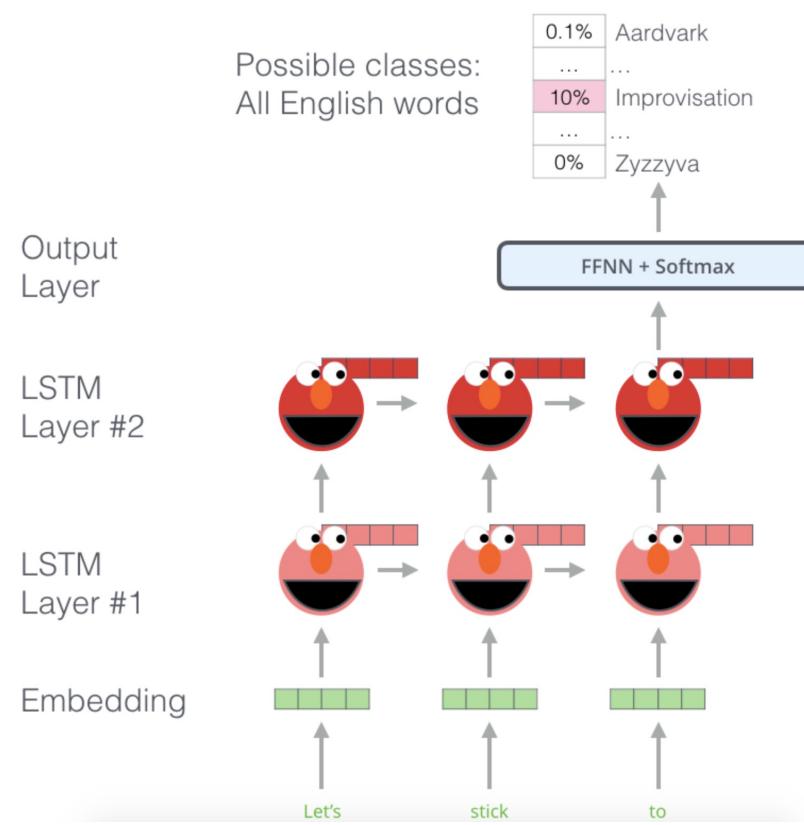
# “Contextualized” Word Embedding



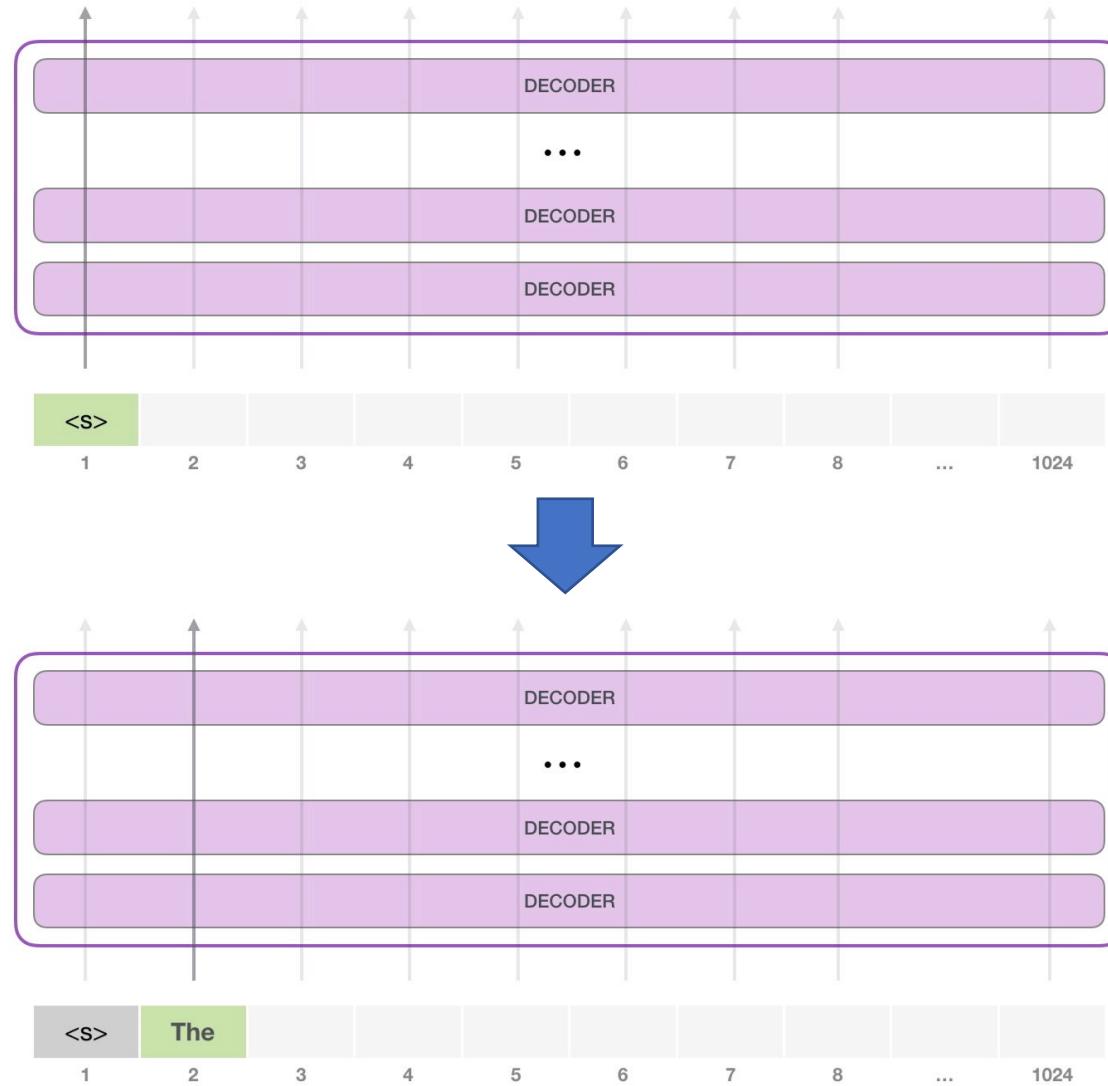
# Transformer-based Language Models



# Pre-training: Language modeling (unsupervised)



# Predicting next words in GPT



# Predicting next words in GPT

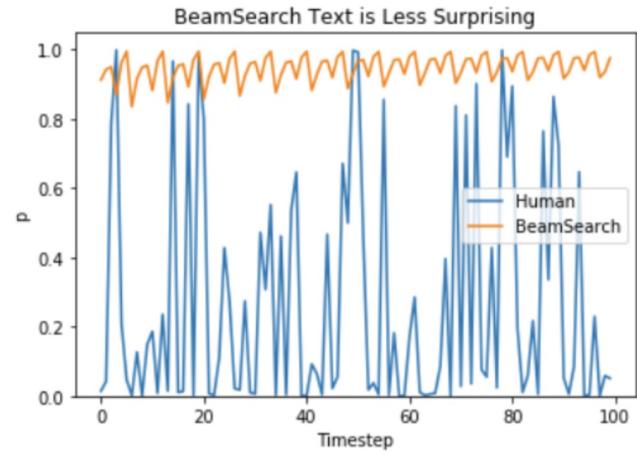
## Context:

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.

### Continuation (BeamSearch, b=10):

Figure 1: Beam search leads to degenerate text, even when generated from GPT-2-117M, in stark contrast with the admirable quality of the text decoded using *top-k* sampling (Radford et al., 2019). The *continuation* is machine generated, conditioned on the *context* provided by a human. Blue text highlights decoded words that have occurred previously in the text.

# Predicting next words in GPT



...get your hopes up. I saw him once and I have no intention of being near him anytime soon. He sat on the edge, the wind tossing around his hair. It was going to be seriously wind-blown later. I sat down next to him and I was trying to forget the dwarfs mangled body. I shook and hugged myself. Are you cold? He asked, his voice full of concern. I just shrugged and squeezed my eyes shut. I saw Kojas glowing eyes and sword, the...

Figure 2: The probability assigned to tokens generated by humans and beam search using GPT-2-117M. Note the increased variance that characterizes the richness of human text.

# Predicting next words in GPT

💡 **Context:** On days when he woke early, the president liked to have coffee in the oval office. There was something about watching the sky grow lighter and lighter as you sat your pajama'd behind in the most powerful chair in the free world, sipping marine strength coffee, that you just couldn't achieve anywhere else.

---

💰 **Greedy:** He'd been in the oval office for over a year now, and he'd never been so happy to be there. He'd been in the oval office for almost a year, and he'd never been so happy to be there.

---

🚦 **BeamSearch (k=4):** The President was a man of his word. He was a man of his word.

---

🔍 **Sampling (t=1.0):** You couldn't be sure if that's what you were really doing, and If you decided to take the day off. The president wanted you to take the day off, but he wanted to maintain a curfew and use his influence wisely.

---

☒ **Top-k Sampling (k=8):** "I'm sorry, sir." "No, it's okay, son," the president said. "I understand." "You're going to have to make a special trip down there to get that kid. He has no idea where he's going."

---

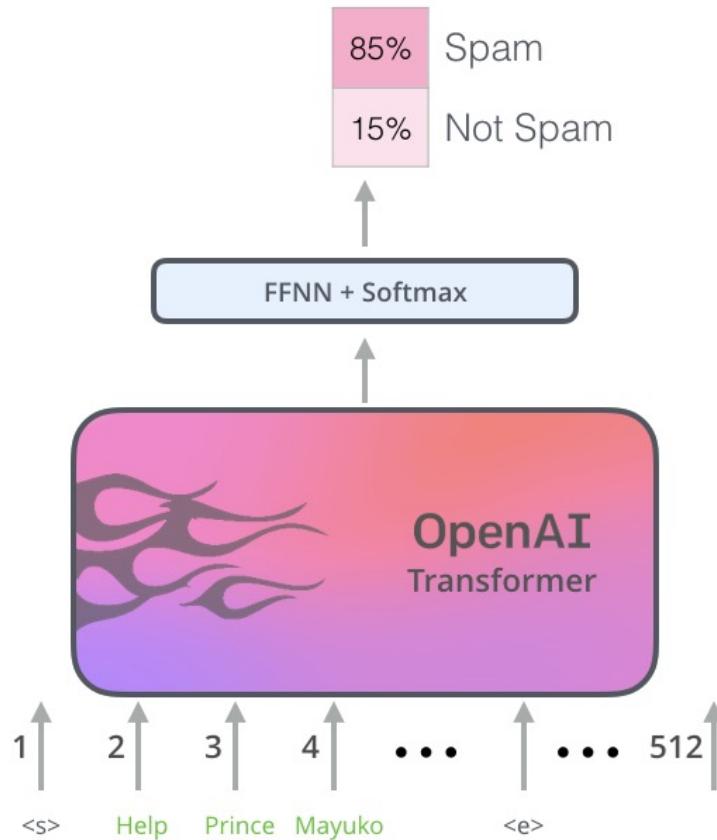
✳️ **Nucleus Sampling ( $p=0.9$ ):** But that wasn't what drew the president's attention. He'd been seated for maybe a minute when he noticed the other man. What was the guy doing here?

---

🗣 **Gold:** He was therefore disagreeably surprised to find a man in an understated grey suit sitting in that selfsame chair sipping tea. The president turned around and went looking for his chief of staff.

Figure 10: Example generations from all discussed decoding strategies, hyperparameters were chosen by experts. All generations for all hyperparameters will be made publicly available.

## Fine-tuning: Transfer learning to downstream tasks

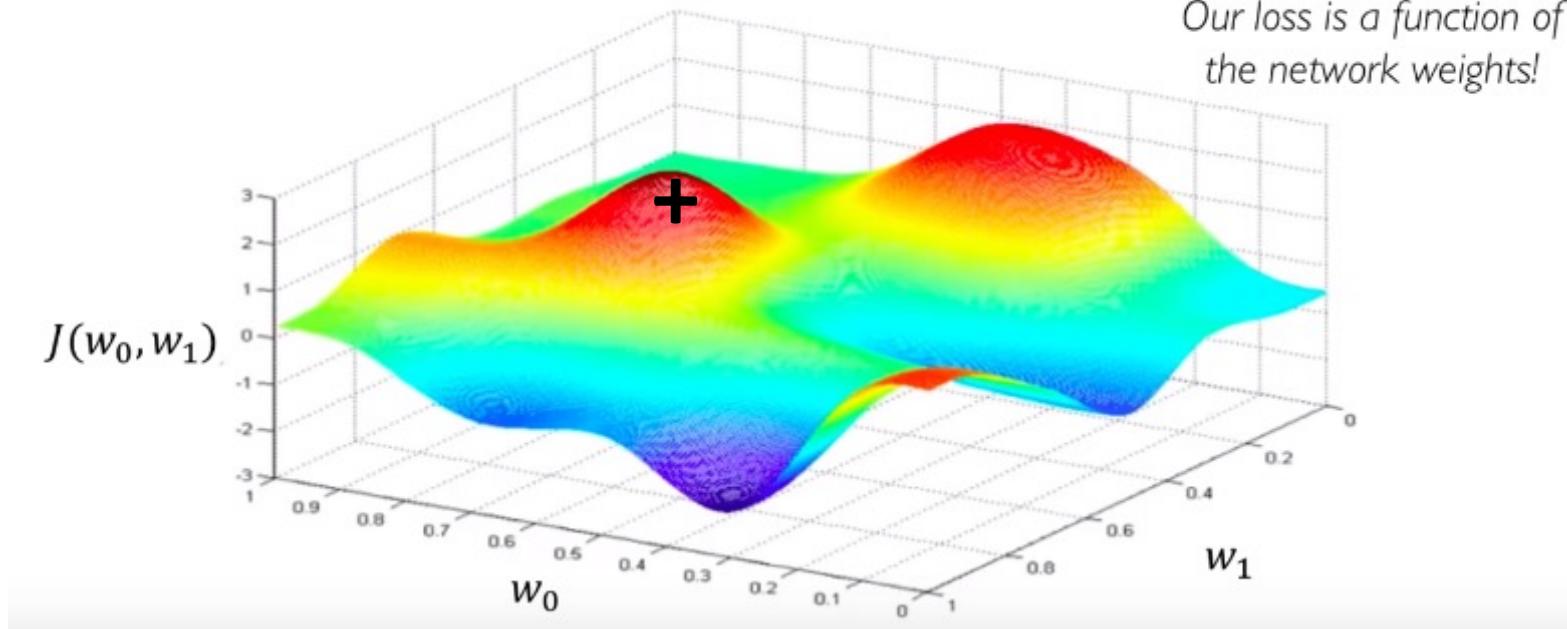


# Why pre-training? Why fine-tuning?

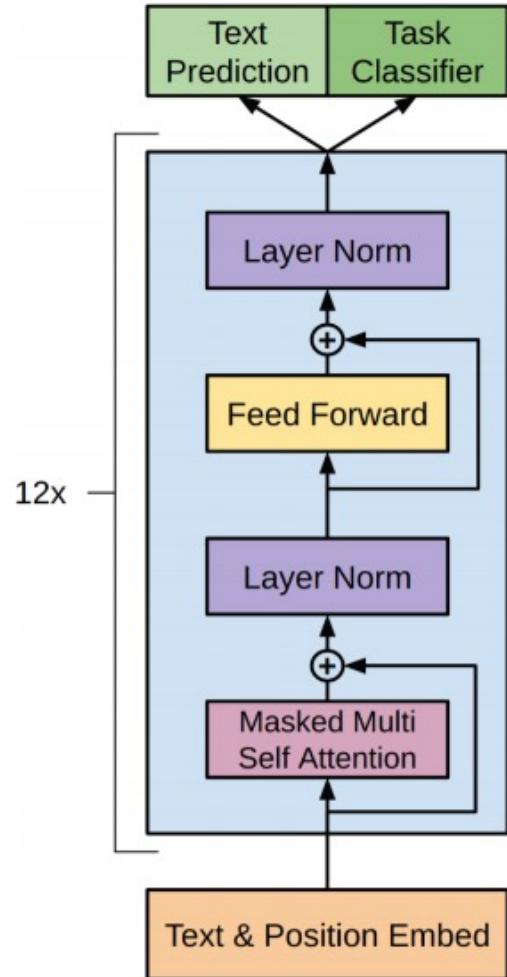
Randomly pick an initial  $(w_0, w_1)$

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$

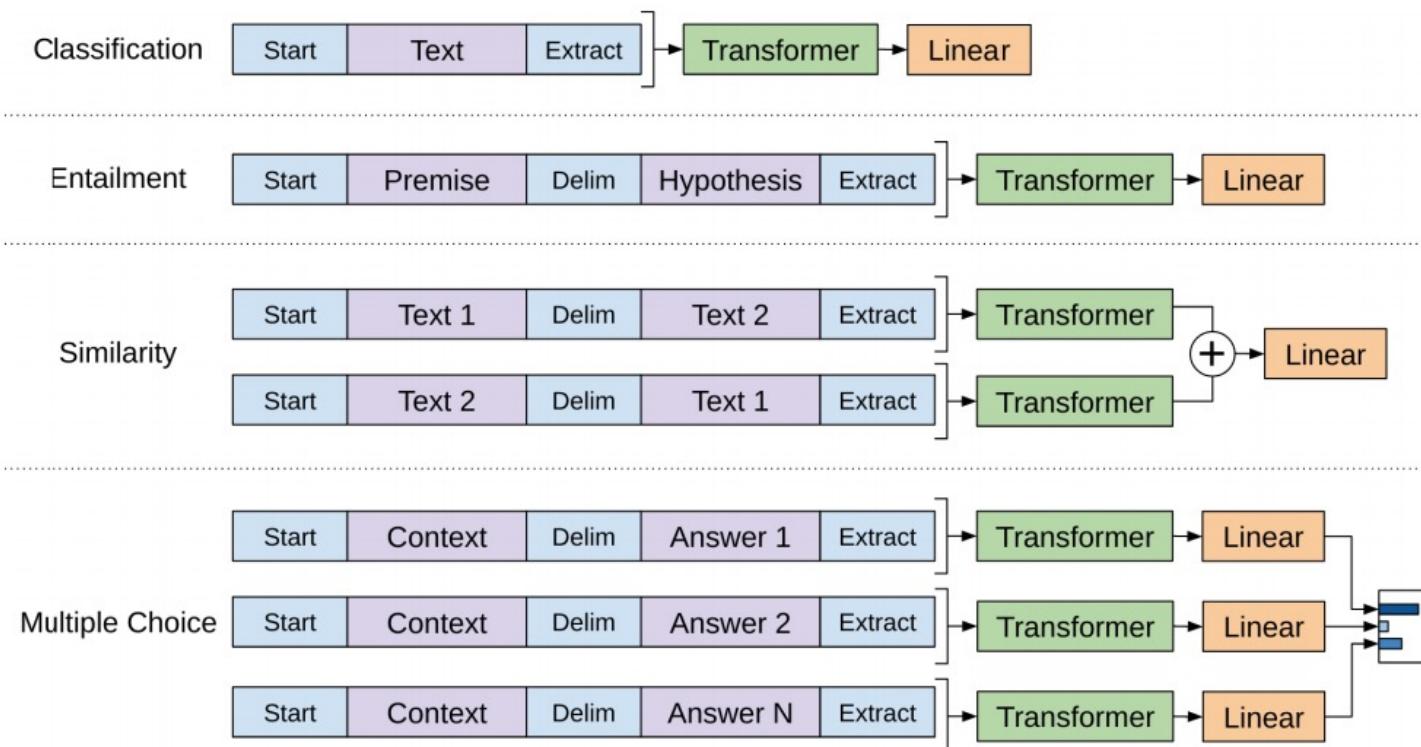
Remember:  
Our loss is a function of  
the network weights!



## Fine-tuning: Transfer learning to downstream tasks



# Fine-tuning: Transfer learning to downstream tasks



# Experiment in GPT

DATASET	TASK	SOTA	OURS
SNLI	Textual Entailment	89.3	<b>89.9</b>
MNLI Matched	Textual Entailment	80.6	<b>82.1</b>
MNLI Mismatched	Textual Entailment	80.1	<b>81.4</b>
SciTail	Textual Entailment	83.3	<b>88.3</b>
QNLI	Textual Entailment	82.3	<b>88.1</b>
RTE	Textual Entailment	<b>61.7</b>	56.0
STS-B	Semantic Similarity	81.0	<b>82.0</b>
QQP	Semantic Similarity	66.1	<b>70.3</b>
MRPC	Semantic Similarity	<b>86.0</b>	82.3
RACE	Reading Comprehension	53.3	<b>59.0</b>
ROCStories	Commonsense Reasoning	77.6	<b>86.5</b>
COPA	Commonsense Reasoning	71.2	<b>78.6</b>
SST-2	Sentiment Analysis	<b>93.2</b>	91.3
CoLA	Linguistic Acceptability	35.0	<b>45.4</b>
GLUE	Multi Task Benchmark	68.9	<b>72.8</b>

# Why Unsupervised Learning?



# BERT: Bidirectional Encoder Representations from Transformers



**BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding**

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova

Google AI Language

{jacobdevlin, mingweichang, kentonl, kristout}@google.com

<https://arxiv.org/pdf/1810.04805.pdf>

google-research / bert

Watch 923

Star 22.6k

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

huggingface / transformers

Watch 536

Star 25.5k

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

# BERT: Bidirectional Encoder Representations from Transformers

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



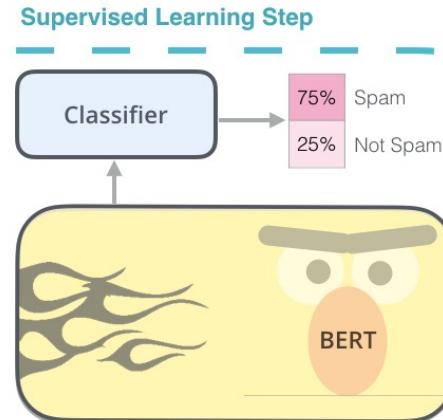
Model:

Dataset:

Objective:

Predict the masked word  
(language modeling)

2 - **Supervised** training on a specific task with a labeled dataset.

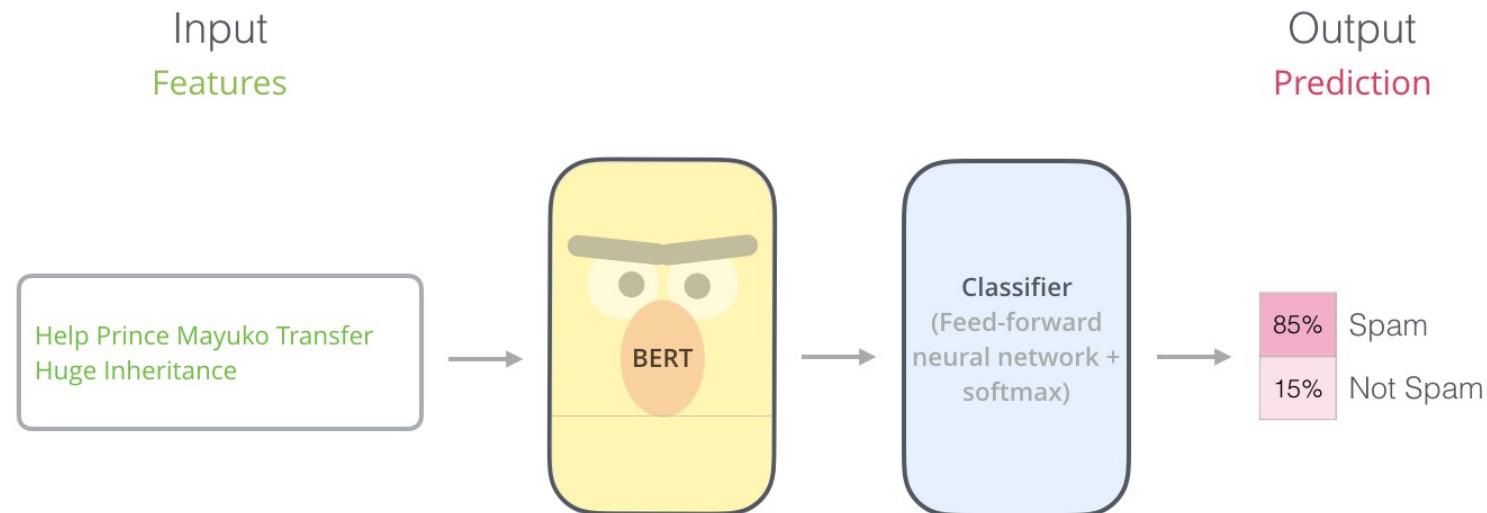


Model:  
(pre-trained  
in step #1)

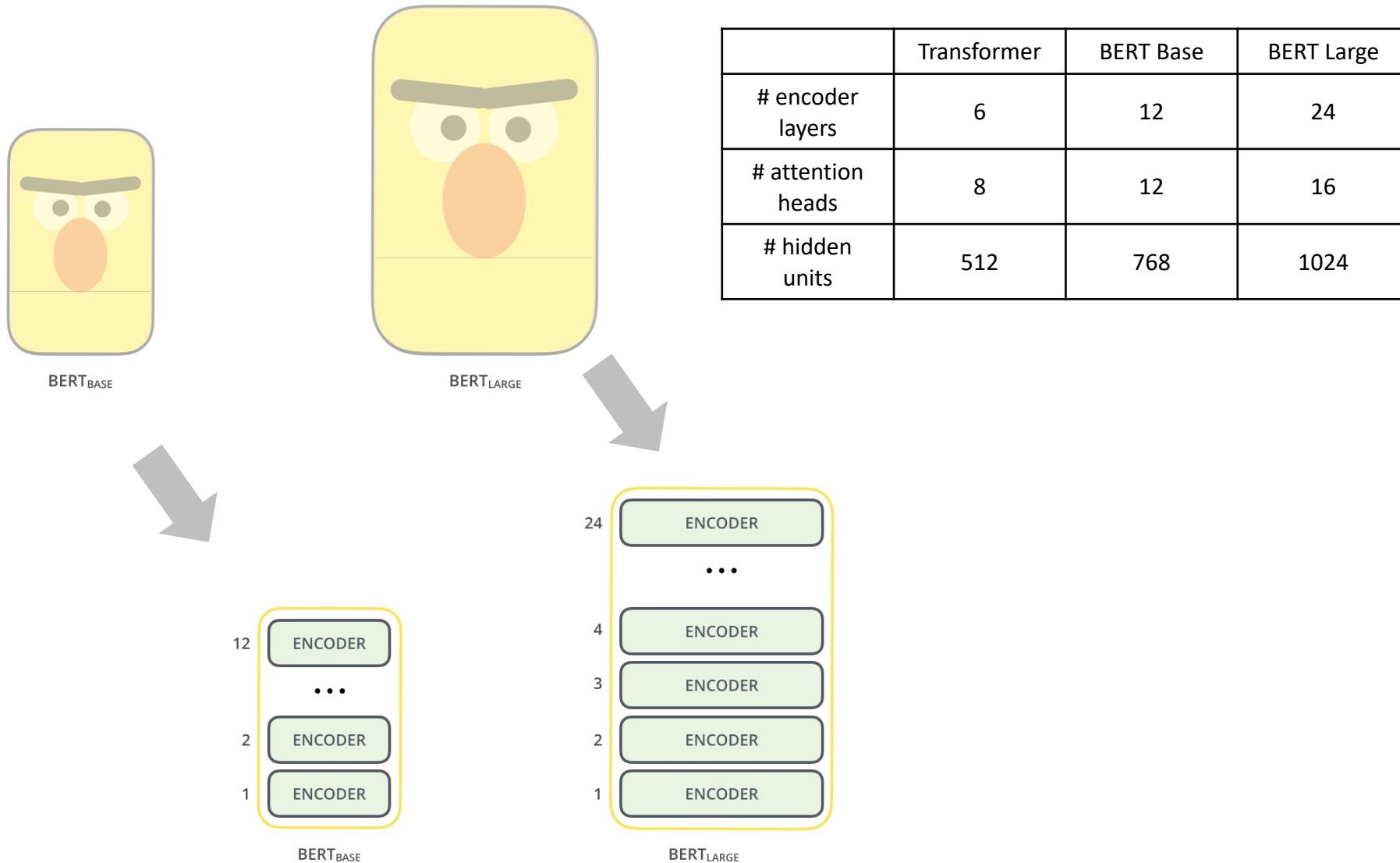
Dataset:

Email message	Class
Buy these pills	Spam
Win cash prizes	Spam
Dear Mr. Atreides, please find attached...	Not Spam

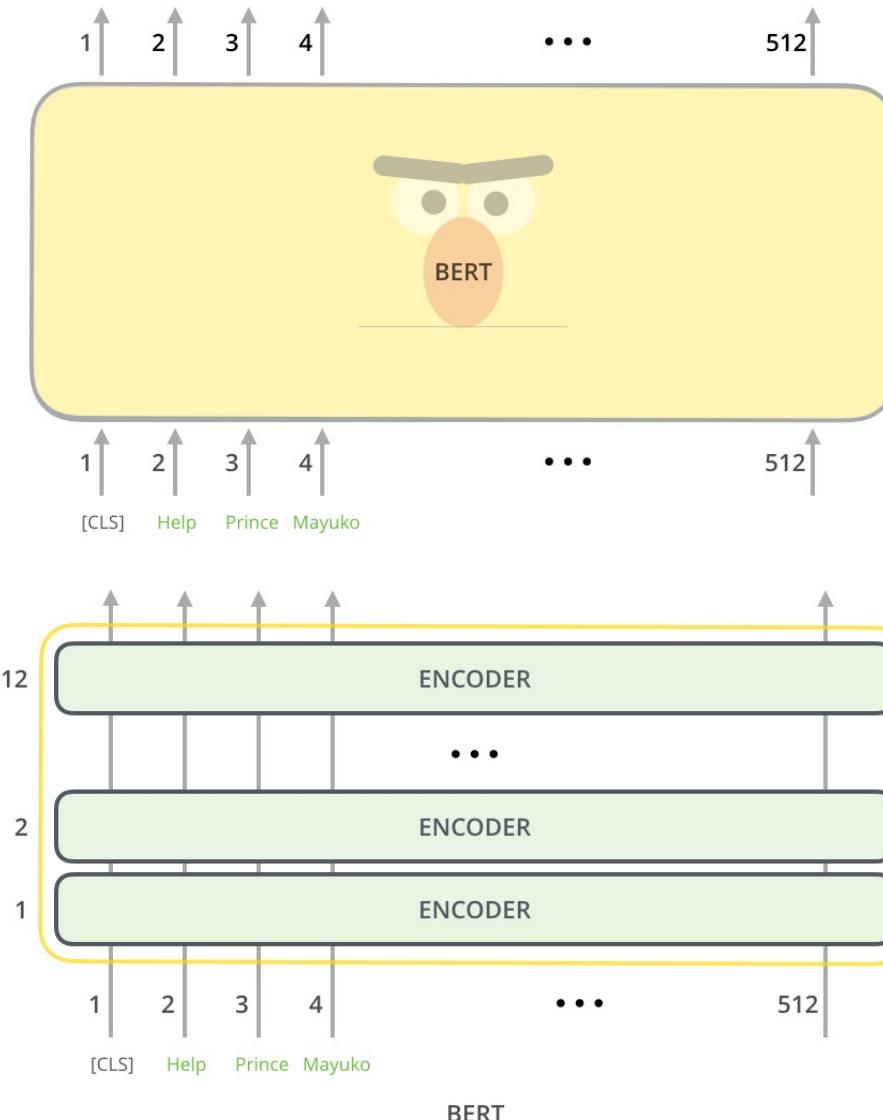
# BERT: Bidirectional Encoder Representations from Transformers



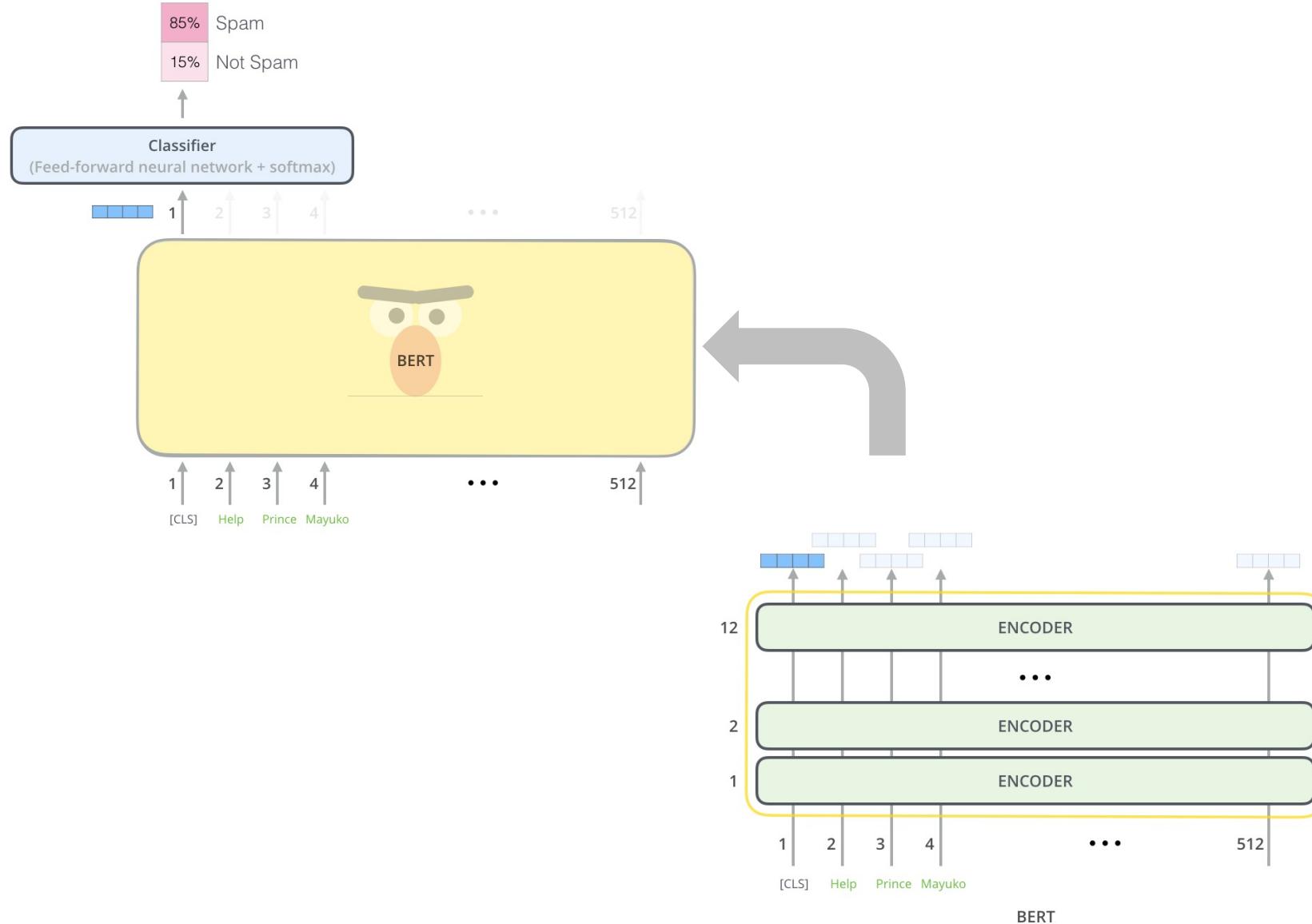
# Architecture of BERT



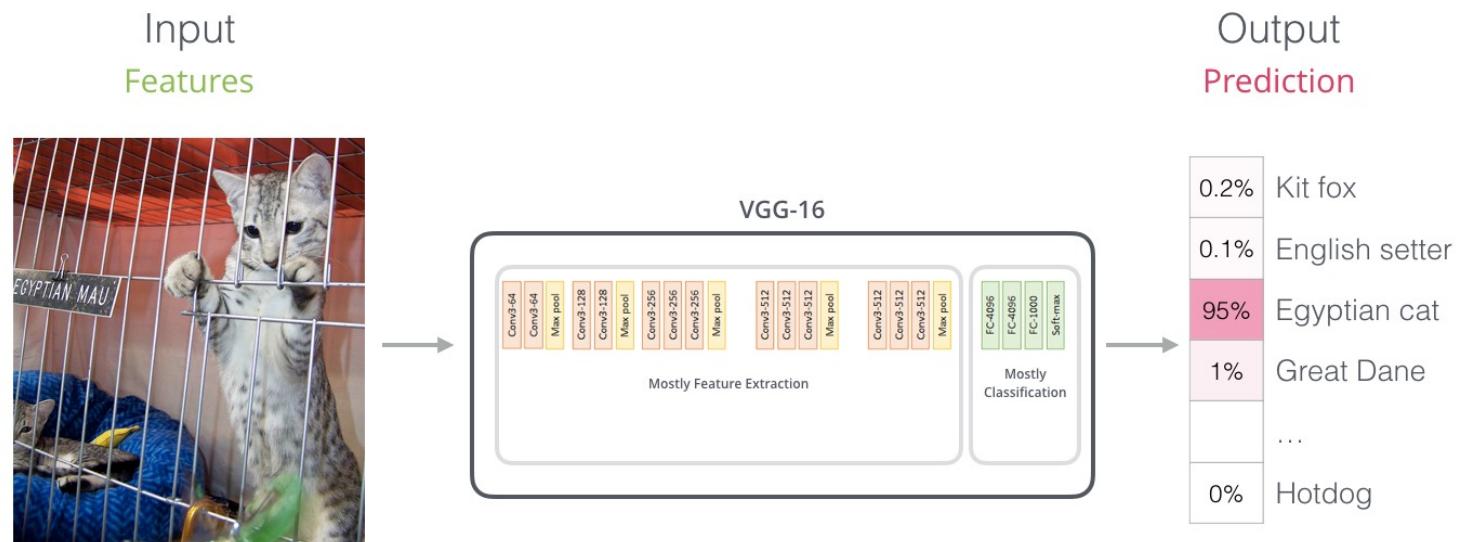
# Architecture of BERT



# Architecture of BERT



# Parallels with Convolutional Nets



# Difference with Vanilla Transformer Encoder

Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	# #ing	[SEP]
<hr/>											
Token Embeddings	$E_{[CLS]}$	$E_{\text{my}}$	$E_{\text{dog}}$	$E_{\text{is}}$	$E_{\text{cute}}$	$E_{[\text{SEP}]}$	$E_{\text{he}}$	$E_{\text{likes}}$	$E_{\text{play}}$	$E_{\#\#\text{ing}}$	$E_{[\text{SEP}]}$
Segment Embeddings	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$
Position Embeddings	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	$E_8$	$E_9$	$E_{10}$

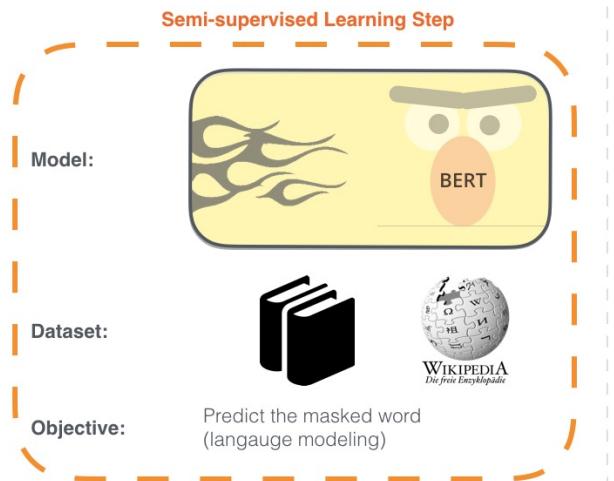
# Difference with Vanilla Transformer Encoder

```
1  """ encoder """
2  class Encoder(nn.Module):
3      def __init__(self, config):
4          super().__init__()
5          self.config = config
6
7          self.enc_emb = nn.Embedding(self.config.n_enc_vocab, self.config.d_hidn)
8          self.pos_emb = nn.Embedding(self.config.n_enc_seq + 1, self.config.d_hidn)
9          self.seg_emb = nn.Embedding(self.config.n_seg_type, self.config.d_hidn)
10
11         self.layers = nn.ModuleList([EncoderLayer(self.config) for _ in range(self.config.n_enc_layers)])
12
13     def forward(self, inputs, segments):
14         positions = torch.arange(inputs.size(1), device=inputs.device, dtype=inputs.dtype)
15         pos_mask = inputs.eq(self.config.i_pad)
16         positions.masked_fill_(pos_mask, 0)
17
18         # (bs, n_enc_seq, d_hidn)
19         outputs = self.enc_emb(inputs) + self.pos_emb(positions) + self.seg_emb(segments)
20
21         # (bs, n_enc_seq, n_enc_seq)
22         attn_mask = get_attn_pad_mask(inputs, inputs, self.config.i_pad)
23
24         attn_probs = []
25         for layer in self.layers:
26             # (bs, n_enc_seq, d_hidn), (bs, n_head, n_enc_seq, n_enc_seq)
27             outputs, attn_prob = layer(outputs, attn_mask)
28             attn_probs.append(attn_prob)
29             # (bs, n_enc_seq, d_hidn), [(bs, n_head, n_enc_seq, n_enc_seq)]
30
31         return outputs, attn_probs
```

# Pre-training

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



**Input** = [CLS] the man went to [MASK] store [SEP]

he bought a gallon [MASK] milk [SEP]

**Label** = IsNext

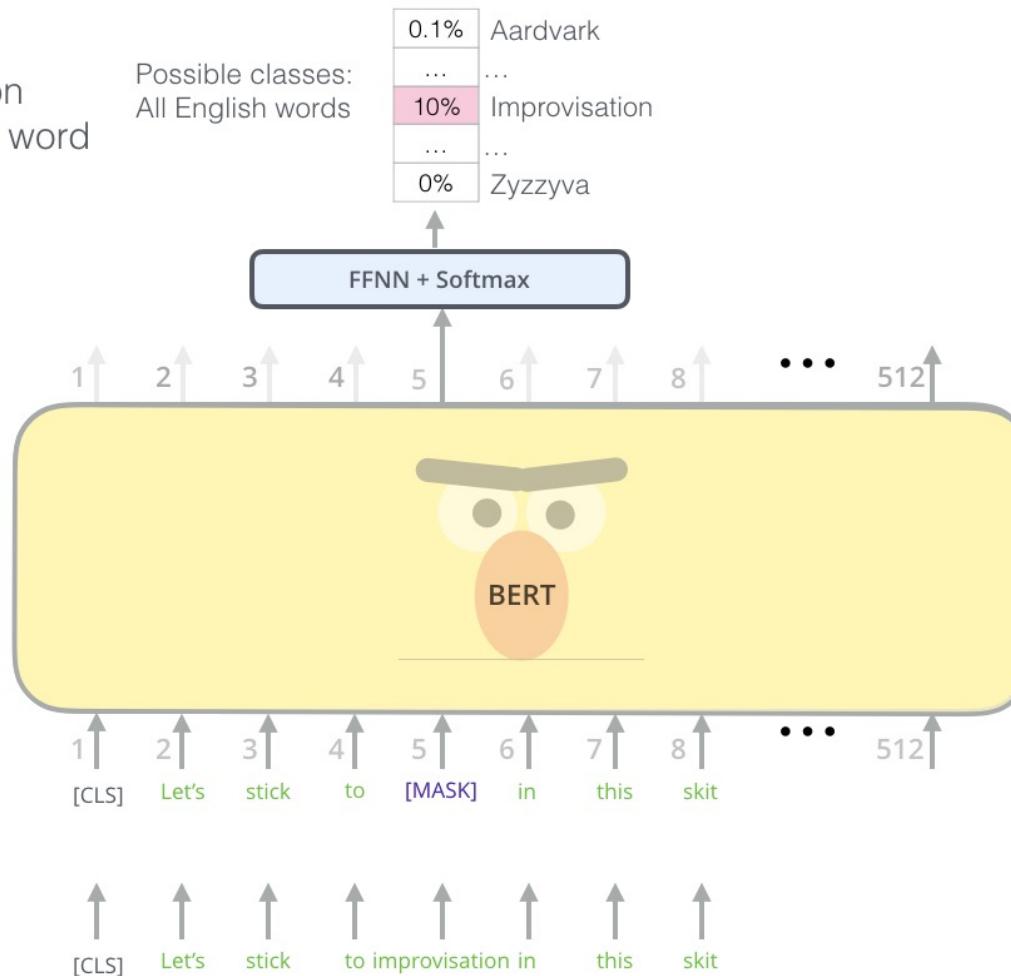
**Input** = [CLS] the man [MASK] to the store [SEP]

penguin [MASK] are flight ##less birds [SEP]

**Label** = NotNext

# Pre-training Task: Mask LM

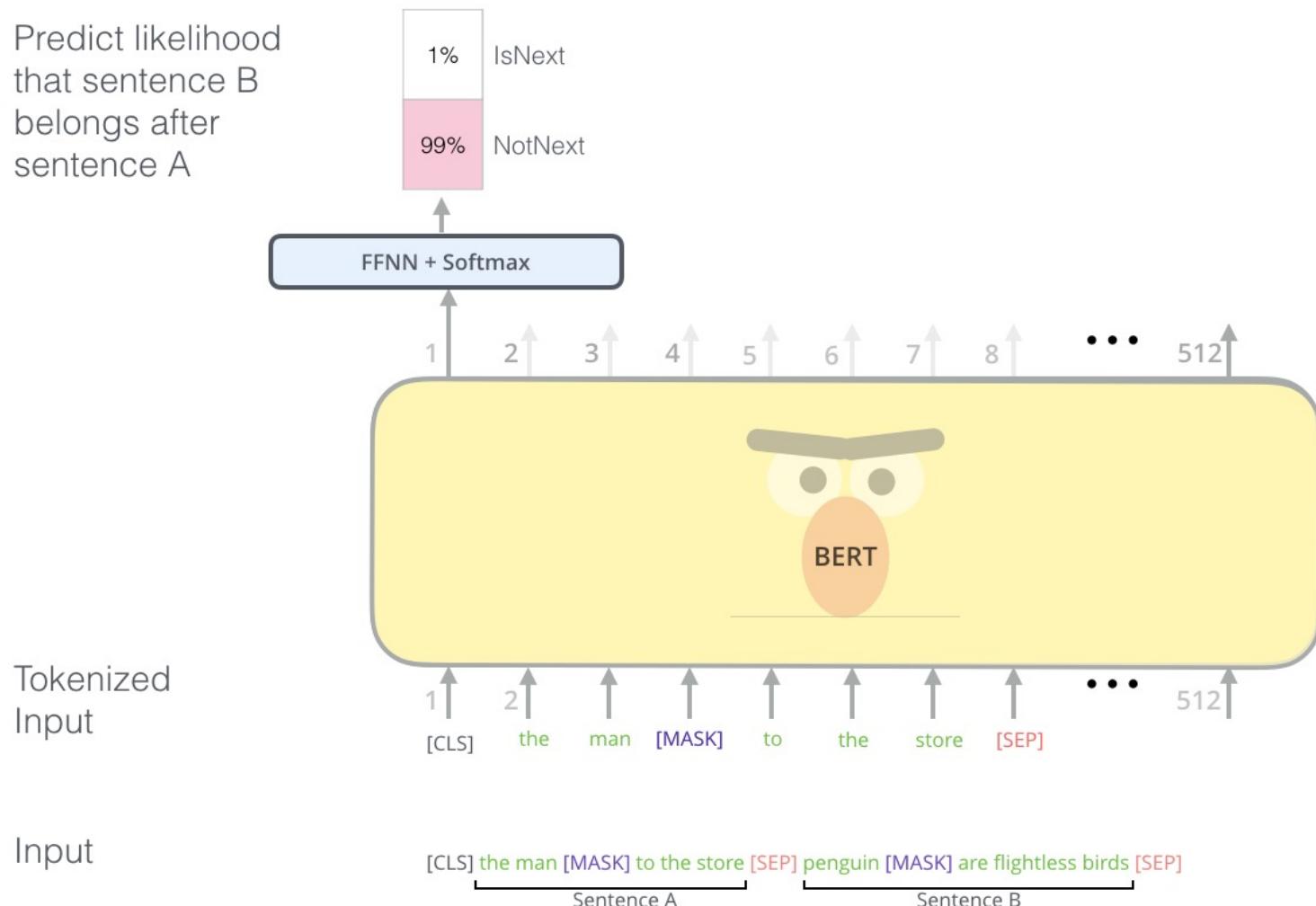
Use the output of the masked word's position to predict the masked word



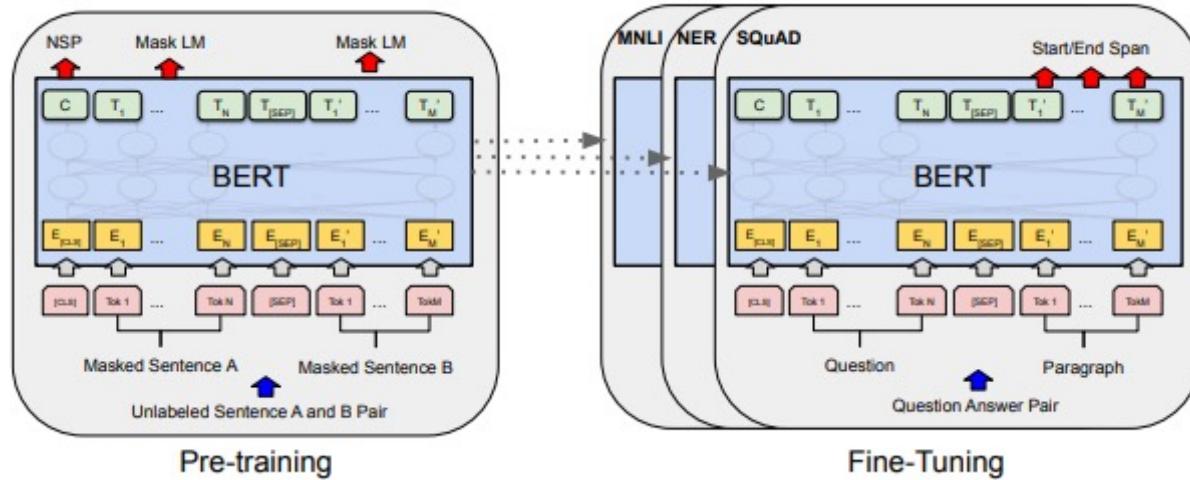
# Pre-training Task: Mask LM

- Rather than *always* replacing the chosen words with [MASK], the data generator will do the following:
- 80% of the time: Replace the word with the [MASK] token, e.g., my dog is hairy → my dog is [MASK]
- 10% of the time: Replace the word with a random word, e.g., my dog is hairy → my dog is apple
- 10% of the time: Keep the word unchanged, e.g., my dog is hairy → my dog is hairy. The purpose of this is to bias the representation towards the actual observed word.

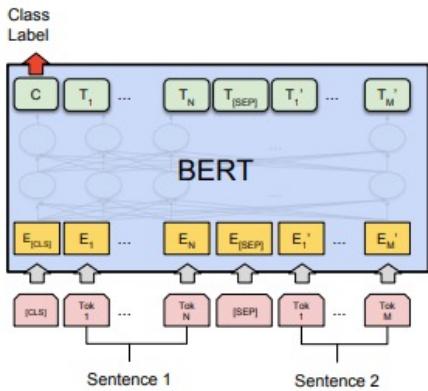
# Pre-training Task: Next Sentence Prediction



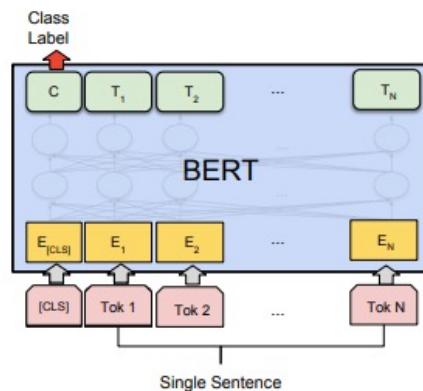
# Fine-tuning Tasks



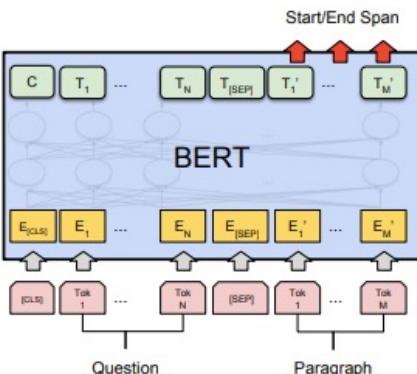
# Fine-tuning Tasks



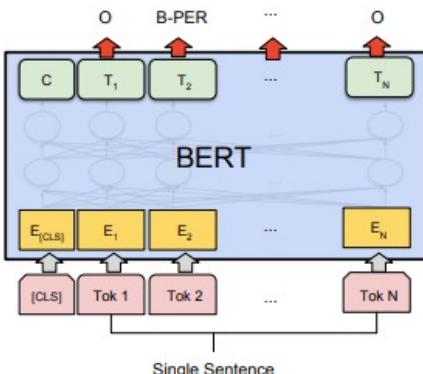
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



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



(c) Question Answering Tasks:  
SQuAD v1.1



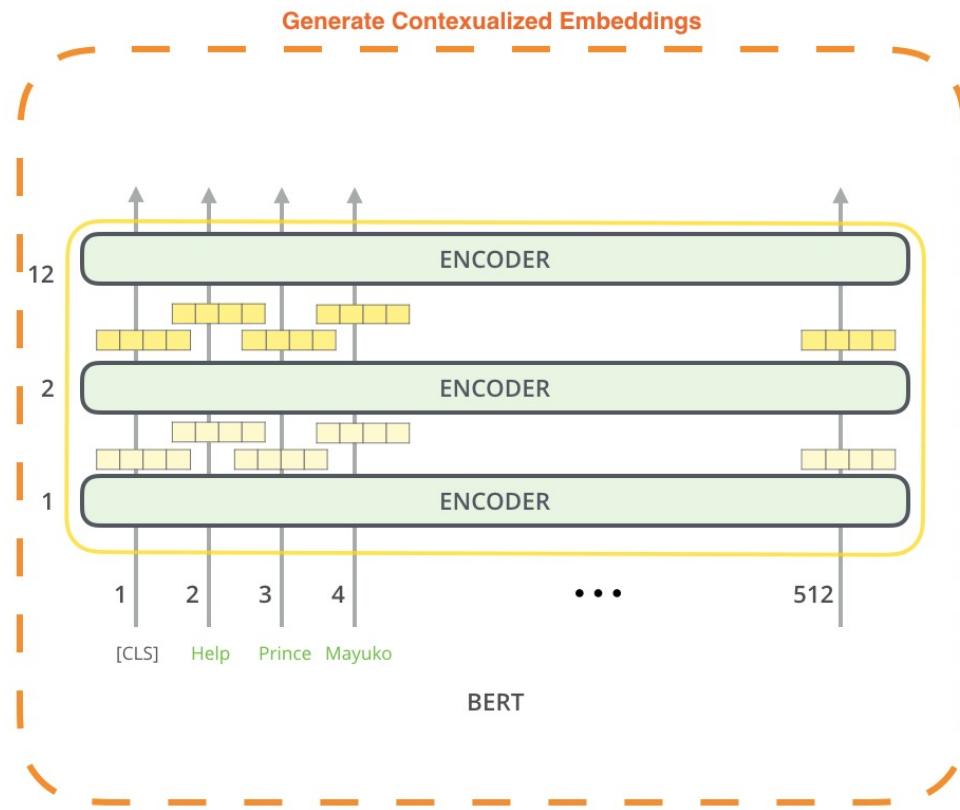
(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

# Fine-tuning Tasks

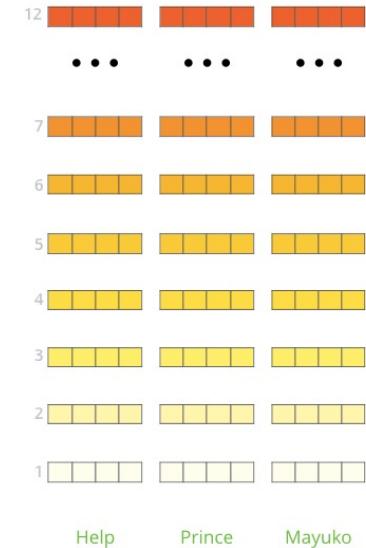
System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT <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>

Table 1: GLUE Test results, scored by the evaluation server (<https://gluebenchmark.com/leaderboard>). The number below each task denotes the number of training examples. The “Average” column is slightly different than the official GLUE score, since we exclude the problematic WNLI set.<sup>8</sup> BERT and OpenAI GPT are single-model, single task. F1 scores are reported for QQP and MRPC, Spearman correlations are reported for STS-B, and accuracy scores are reported for the other tasks. We exclude entries that use BERT as one of their components.

# BERT for feature extraction



The output of each encoder layer along each token's path can be used as a feature representing that token.



But which one should we use?

# BERT for feature extraction

What is the best contextualized embedding for “Help” in that context?

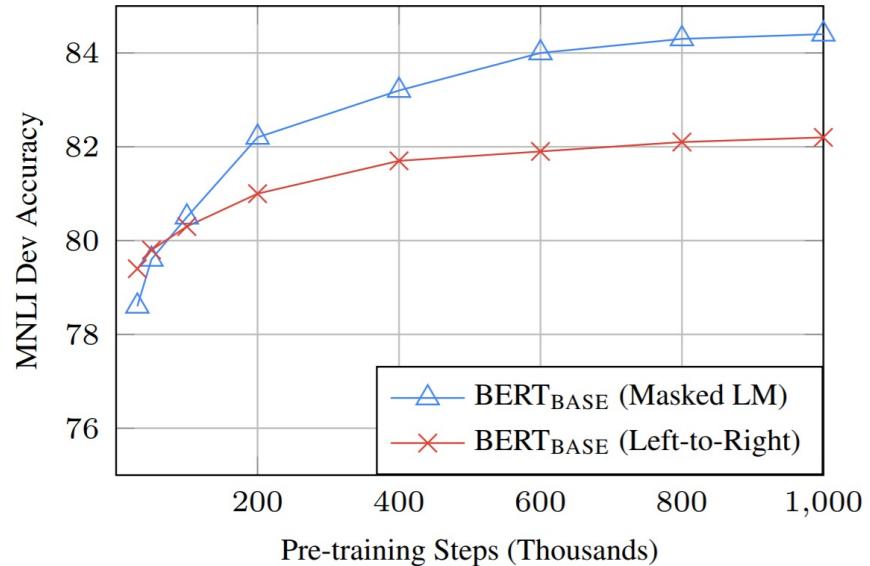
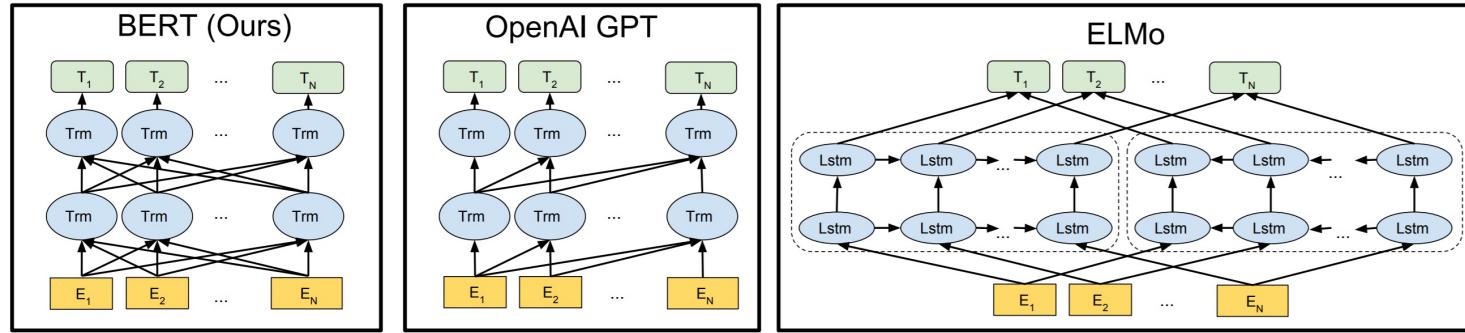
For named-entity recognition task CoNLL-2003 NER

		Dev F1 Score
12		
• • •		
7		
6		
5		
4		
3		
2		
1		
Help		
First Layer	Embedding	91.0
Last Hidden Layer		94.9
Sum All 12 Layers		95.5
Second-to-Last Hidden Layer		95.6
Sum Last Four Hidden		95.9
Concat Last Four Hidden		96.1

# BERT for feature extraction

Layers	Dev F1
Finetune All	96.4
First Layer (Embeddings)	91.0
Second-to-Last Hidden	95.6
Last Hidden	94.9
Sum Last Four Hidden	95.9
Concat Last Four Hidden	96.1
Sum All 12 Layers	95.5

# BERT vs. GPT vs. ELMo



Thank you for listening. Any question?

# References

- **The Illustrated BERT, ELMo, and co.**
  - <http://jalammar.github.io/illustrated-bert/>
- **The Illustrated GPT-2**
  - <http://jalammar.github.io/illustrated-gpt2/#part-1-got-and-language-modeling>