

Yale

Deep Learning Theory and Applications
Transformers

CPSC/AMTH 452

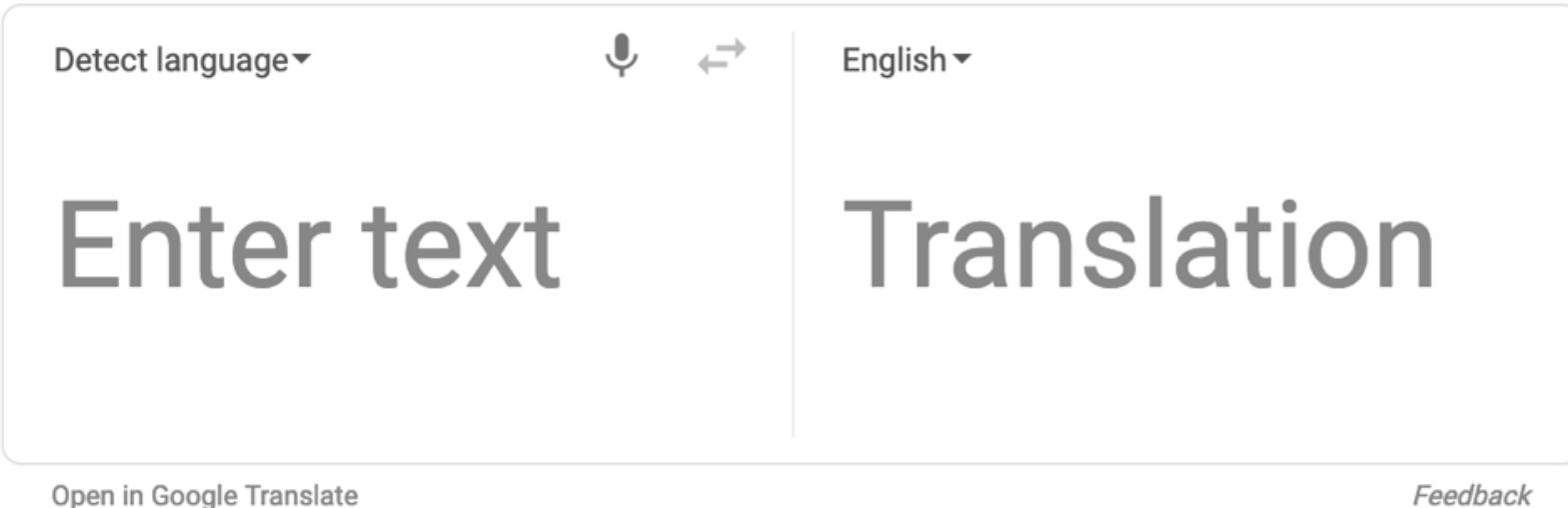
CBB 663



Outline

- RNN review
- Sequence transduction tasks
- Attention
- Self Attention
- Multi-head Attention
- Positional Encoding
- Transformer Architecture
- GPT2/3

Machine Translation



Statistical machine translation:

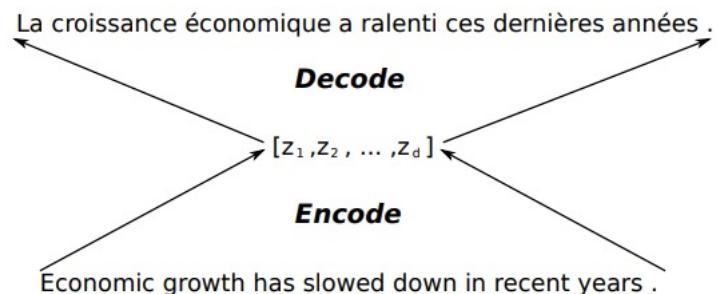
$P(e|f)$ e = english phrase f= foreign phrase

$p(e|f) \propto p(f|e)p(e)$ language model $p(e)$

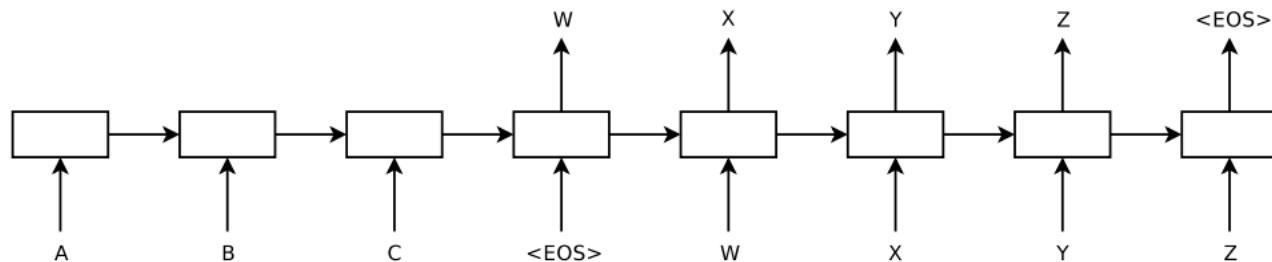
Translates n-grams phrases consisting of n-words

Neural Machine translation

Encoder-decoder architecture



Cho et al 2014, Sutskever 2014



Neural Machine Translation

SEQUENCE TO SEQUENCE MODEL WITH ATTENTION

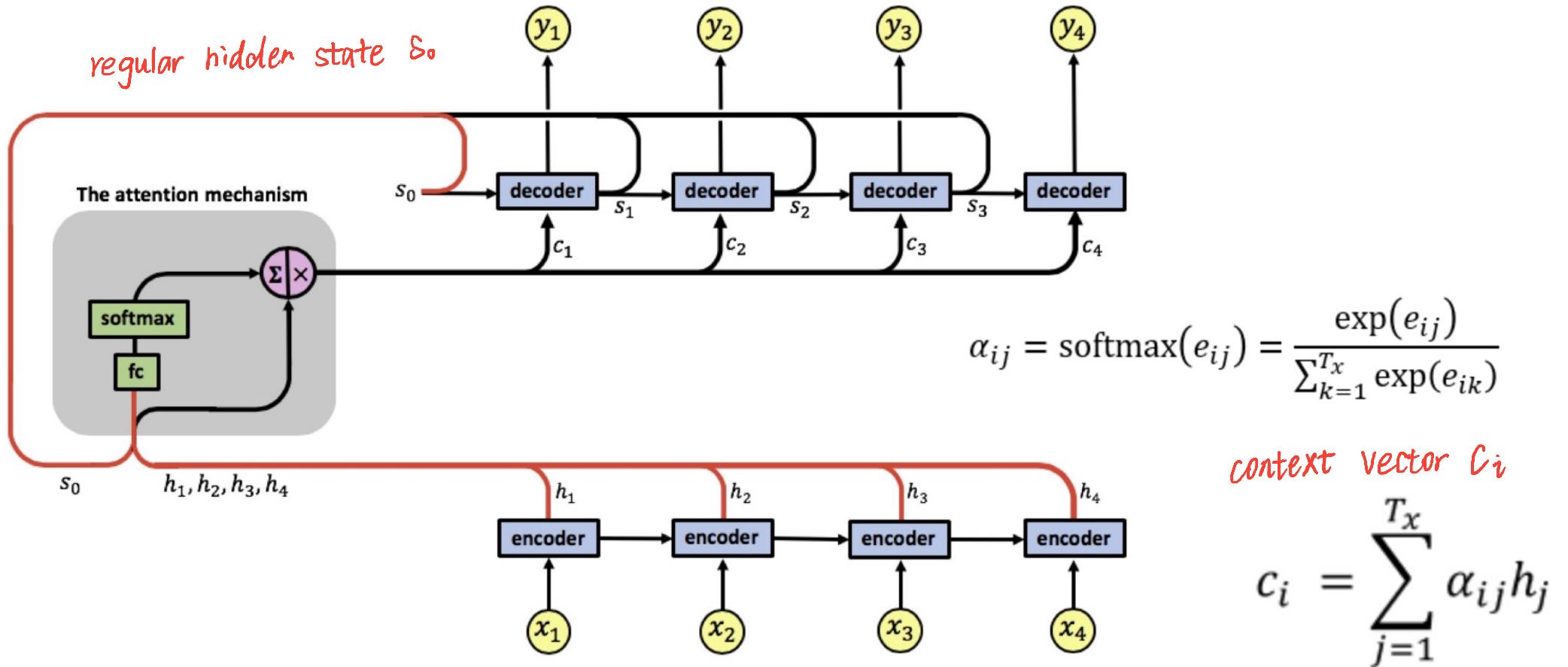


Je

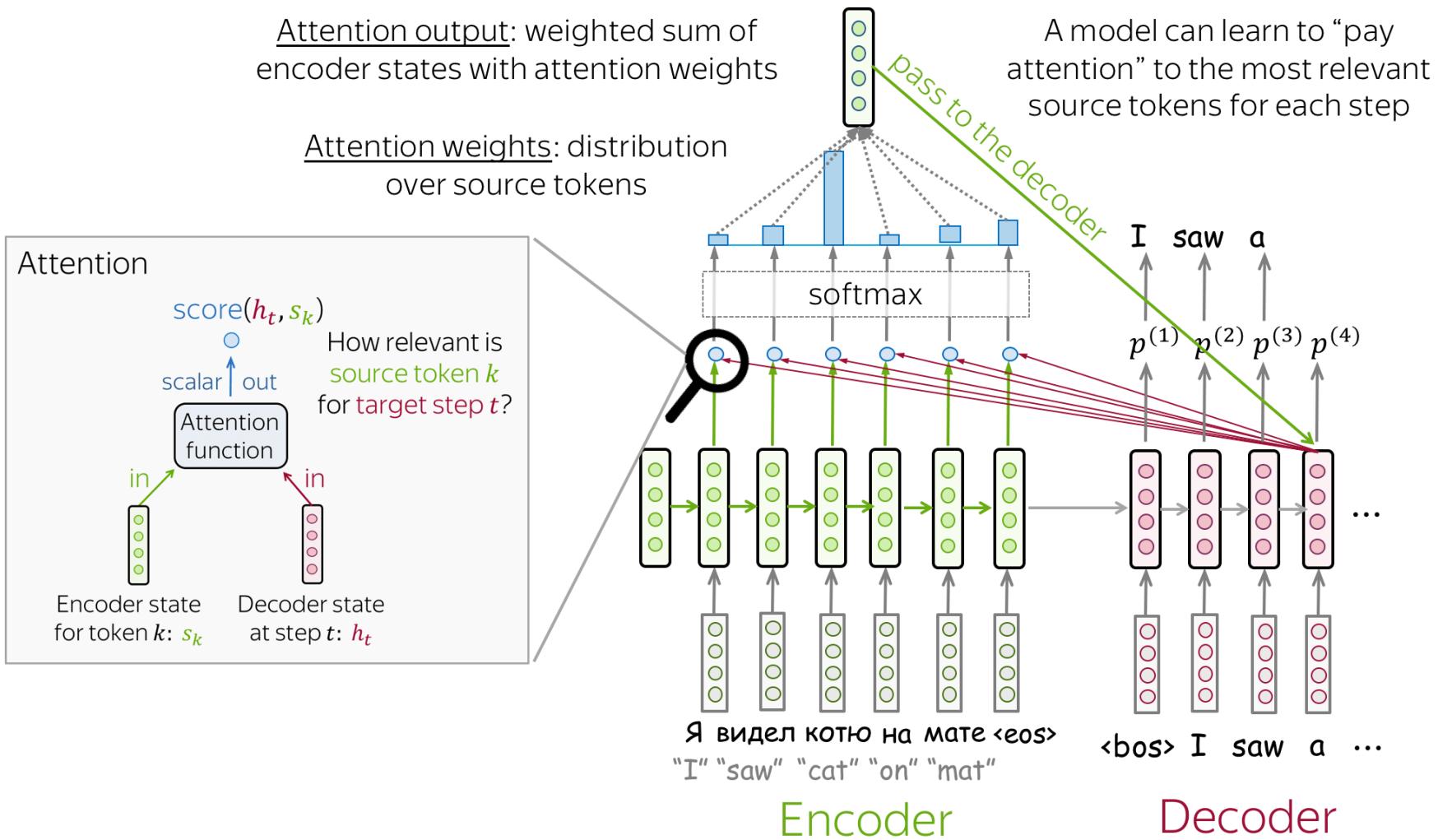
suis

étudiant

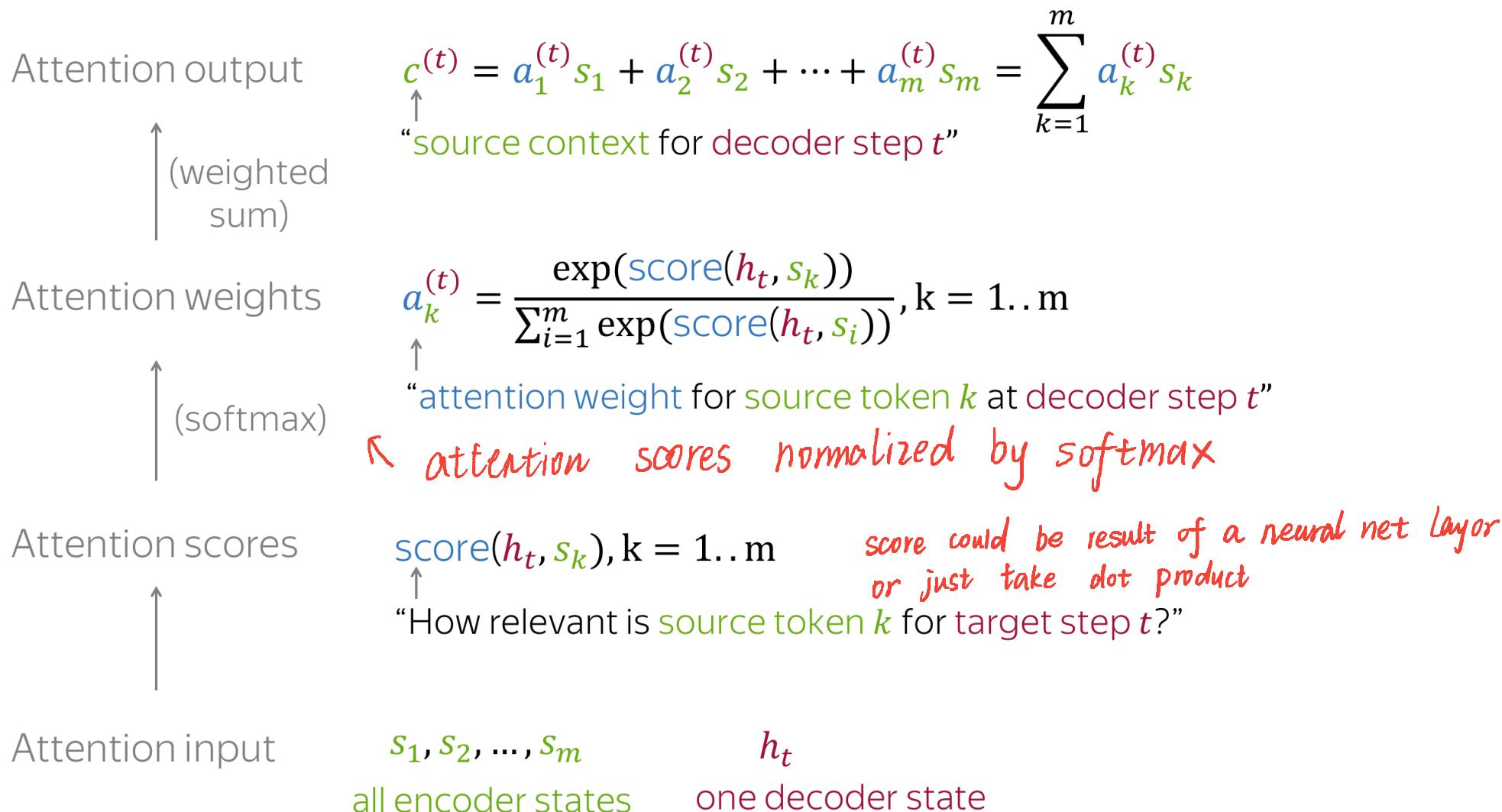
Attention: Adaptively Allows Decoder to Choose Inputs



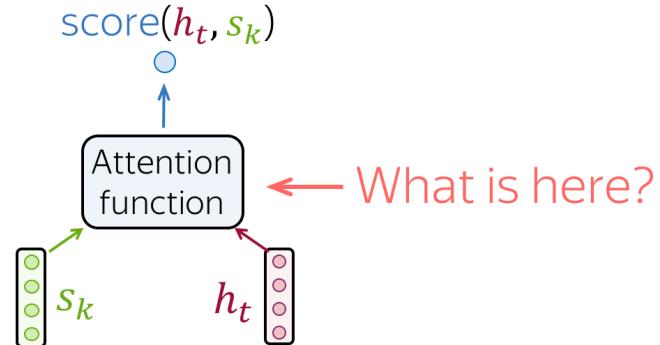
Attention Mechanism



Attention Computations



Attention Score



Dot-product

$$h_t^T \times \begin{bmatrix} \textcolor{green}{\circ} \\ \textcolor{green}{\circ} \\ \textcolor{green}{\circ} \\ \textcolor{green}{\circ} \end{bmatrix} s_k$$

Bilinear

$$h_t^T \times \begin{bmatrix} W \end{bmatrix} \times \begin{bmatrix} \textcolor{green}{\circ} \\ \textcolor{green}{\circ} \\ \textcolor{green}{\circ} \\ \textcolor{green}{\circ} \end{bmatrix} s_k$$

Multi-Layer Perceptron

$$\begin{bmatrix} w_2^T \\ \textcolor{blue}{\circ} \\ \textcolor{blue}{\circ} \\ \textcolor{blue}{\circ} \\ \textcolor{blue}{\circ} \end{bmatrix} \times \tanh \left[\begin{bmatrix} W_1 \end{bmatrix} \times \begin{bmatrix} \textcolor{red}{\circ} \\ \textcolor{red}{\circ} \\ \textcolor{red}{\circ} \\ \textcolor{red}{\circ} \\ \textcolor{red}{\circ} \end{bmatrix} h_t \right] s_k$$

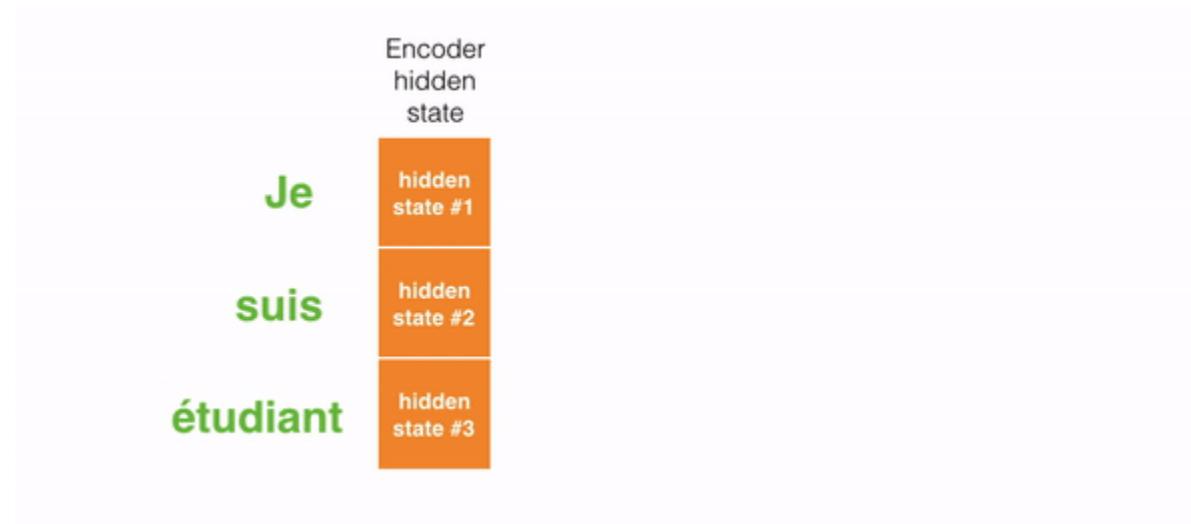
$$\text{score}(h_t, s_k) = h_t^T s_k$$

$$\text{score}(h_t, s_k) = h_t^T W s_k$$

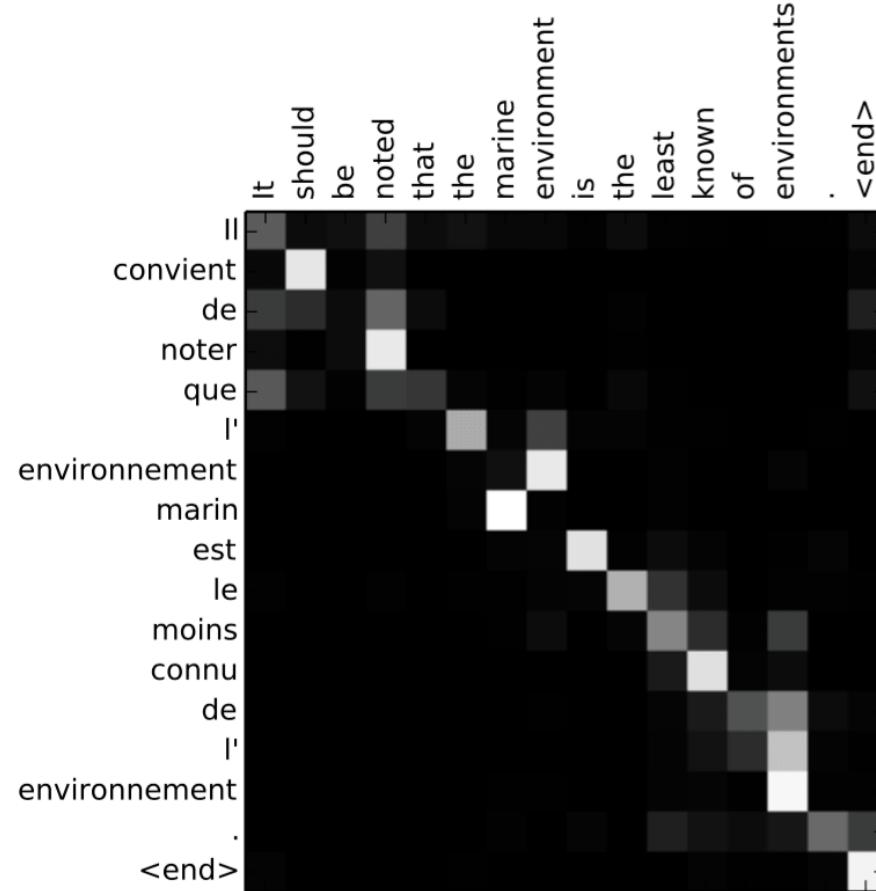
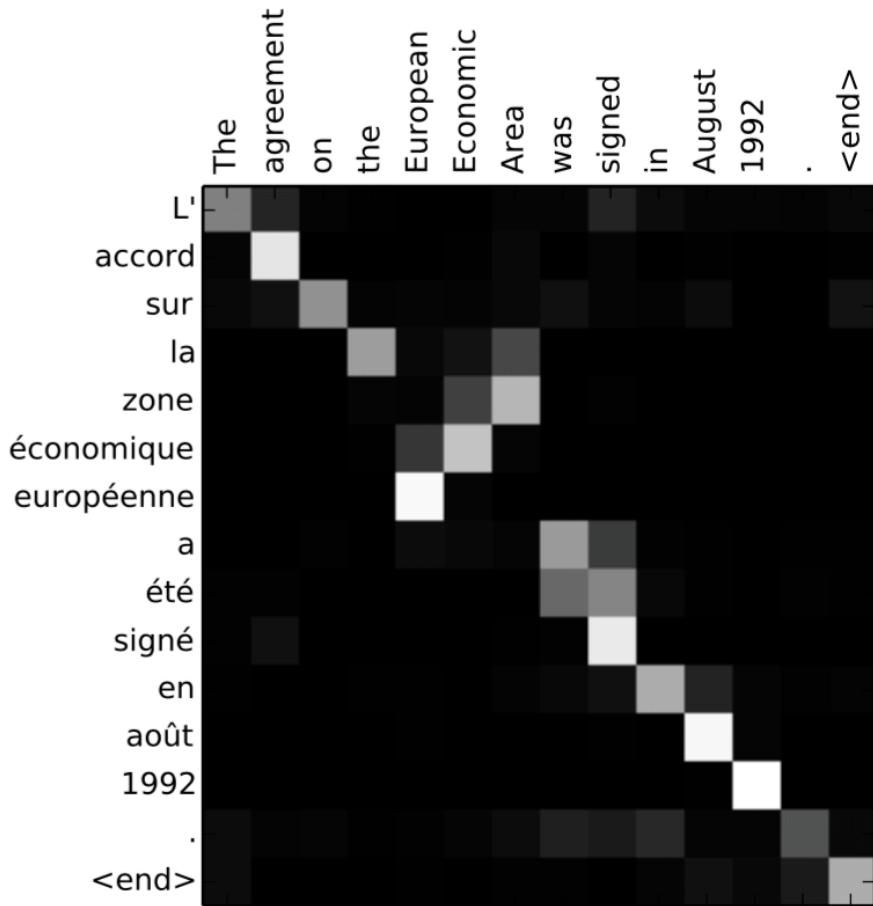
$$\text{score}(h_t, s_k) = w_2^T \cdot \tanh(W_1[h_t, s_k])$$

Attention Weights

can highlight different parts of input for the output

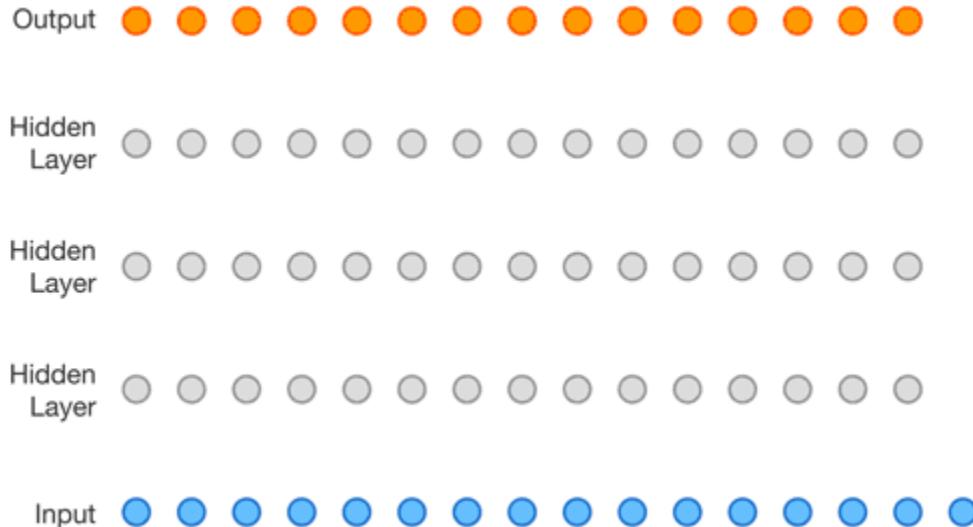


Interpreting Attention



CNNs for sequences

Transformer use attention but limit it for recurrent
in a sense like CNN

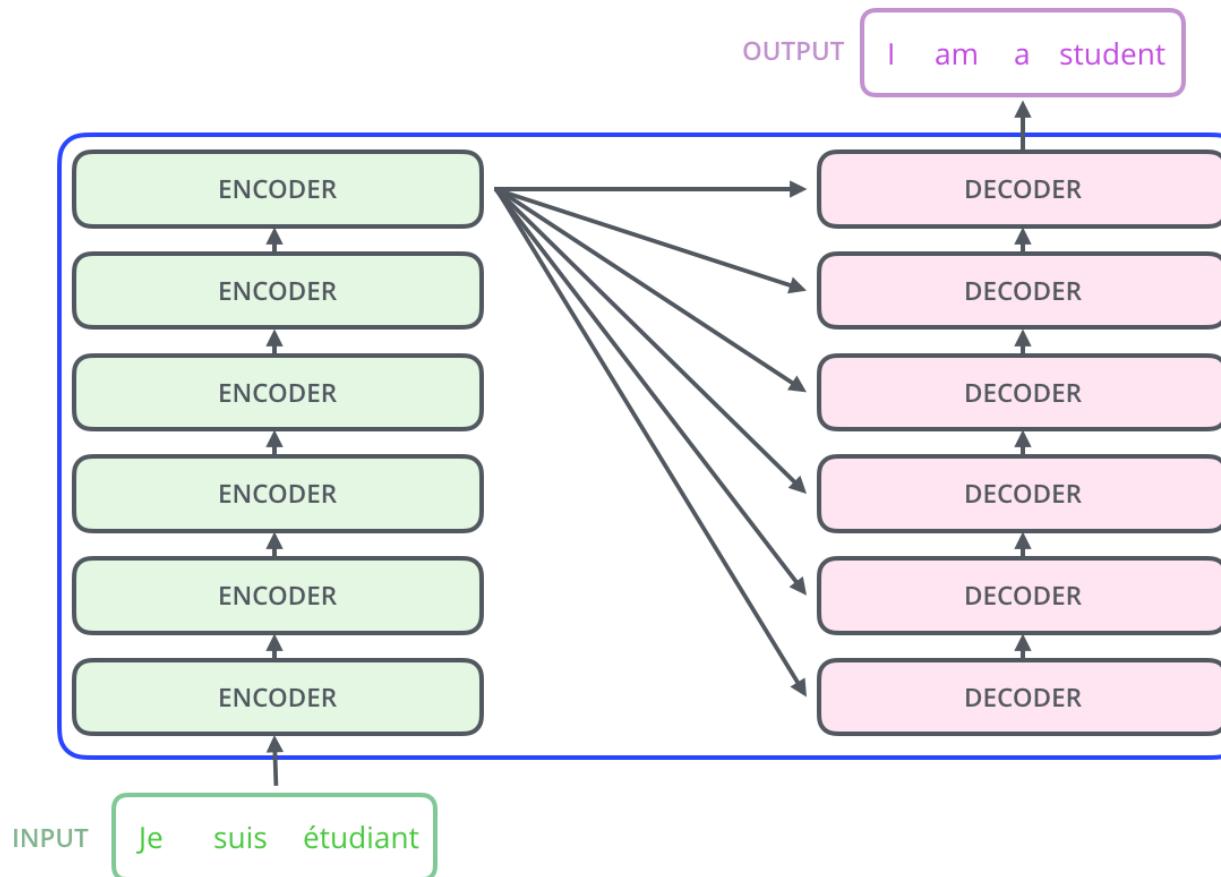


- Easy to parallelize
- Can exploit dependences
- Distance between positions is logarithmic

1D convolution

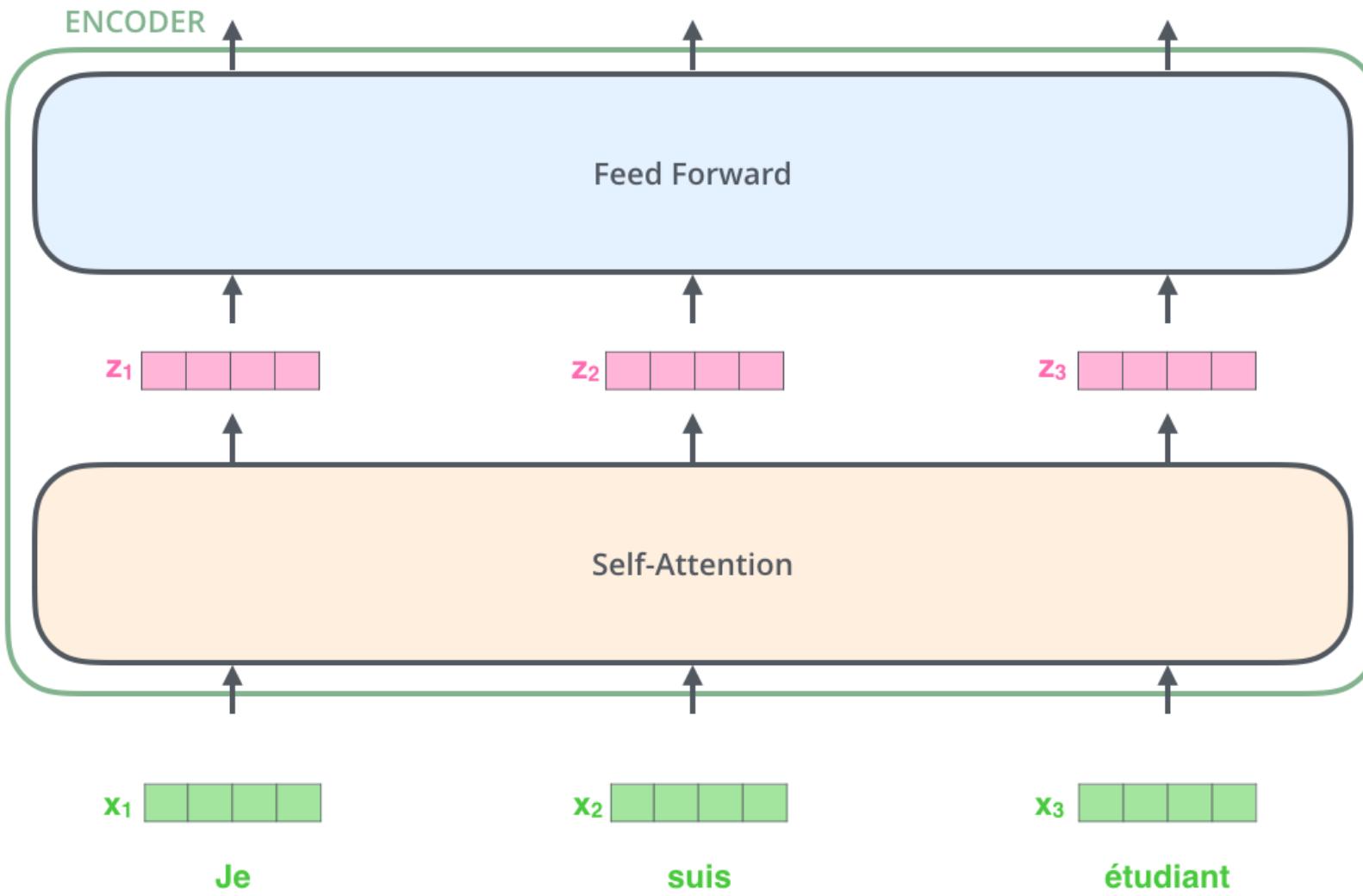
Transformer: CNN + Attention

original transformer has 6 encoders and 6 decoders



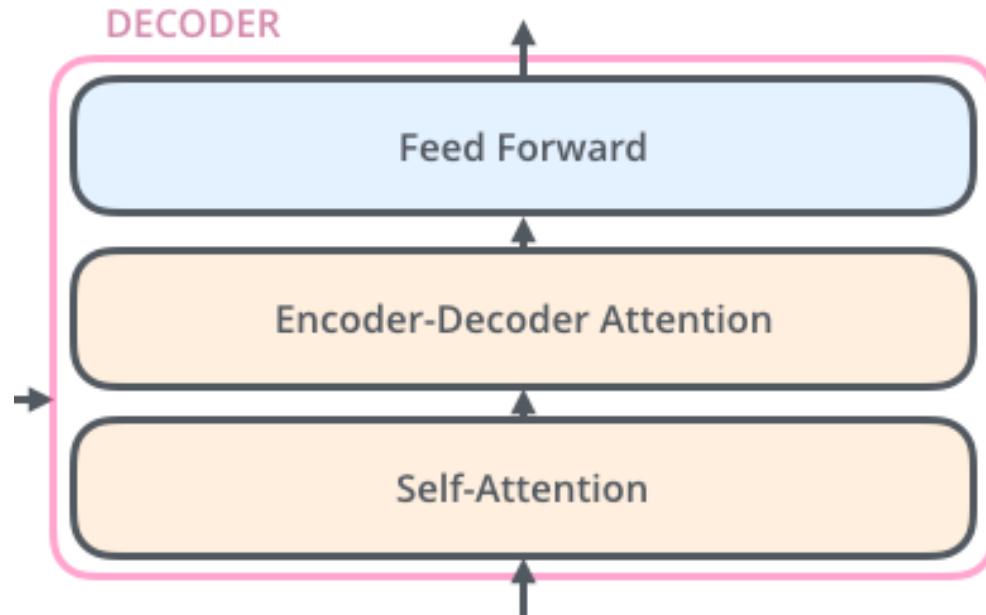
one layer of self-attention and one layer of normal feed forward

Each Encoder has Self-Attention



Decoder

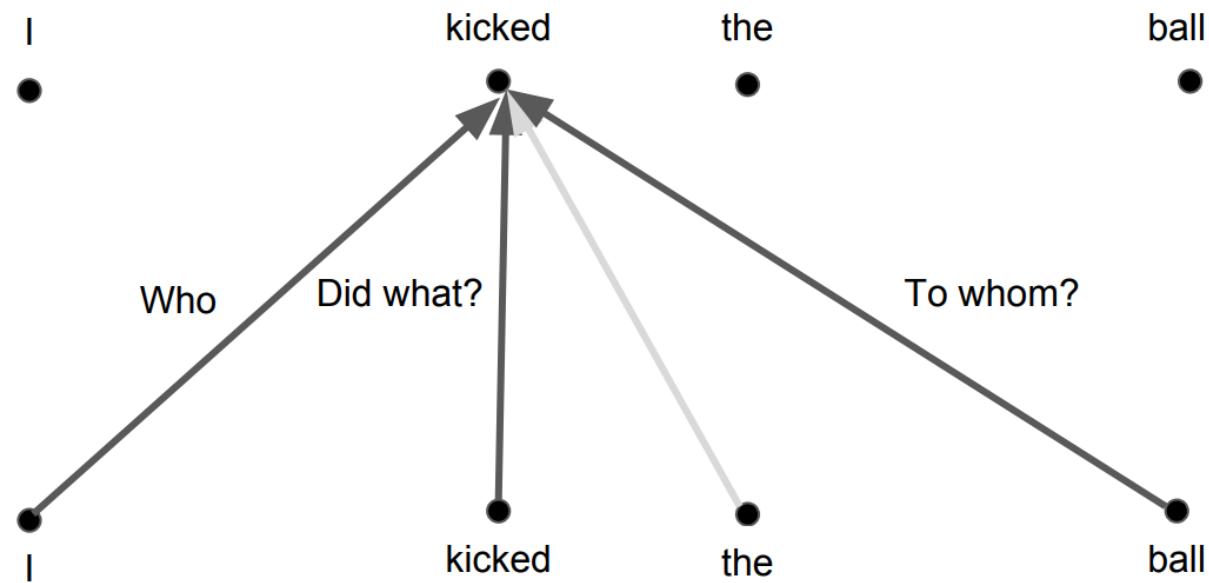
has 3 kinds of layer



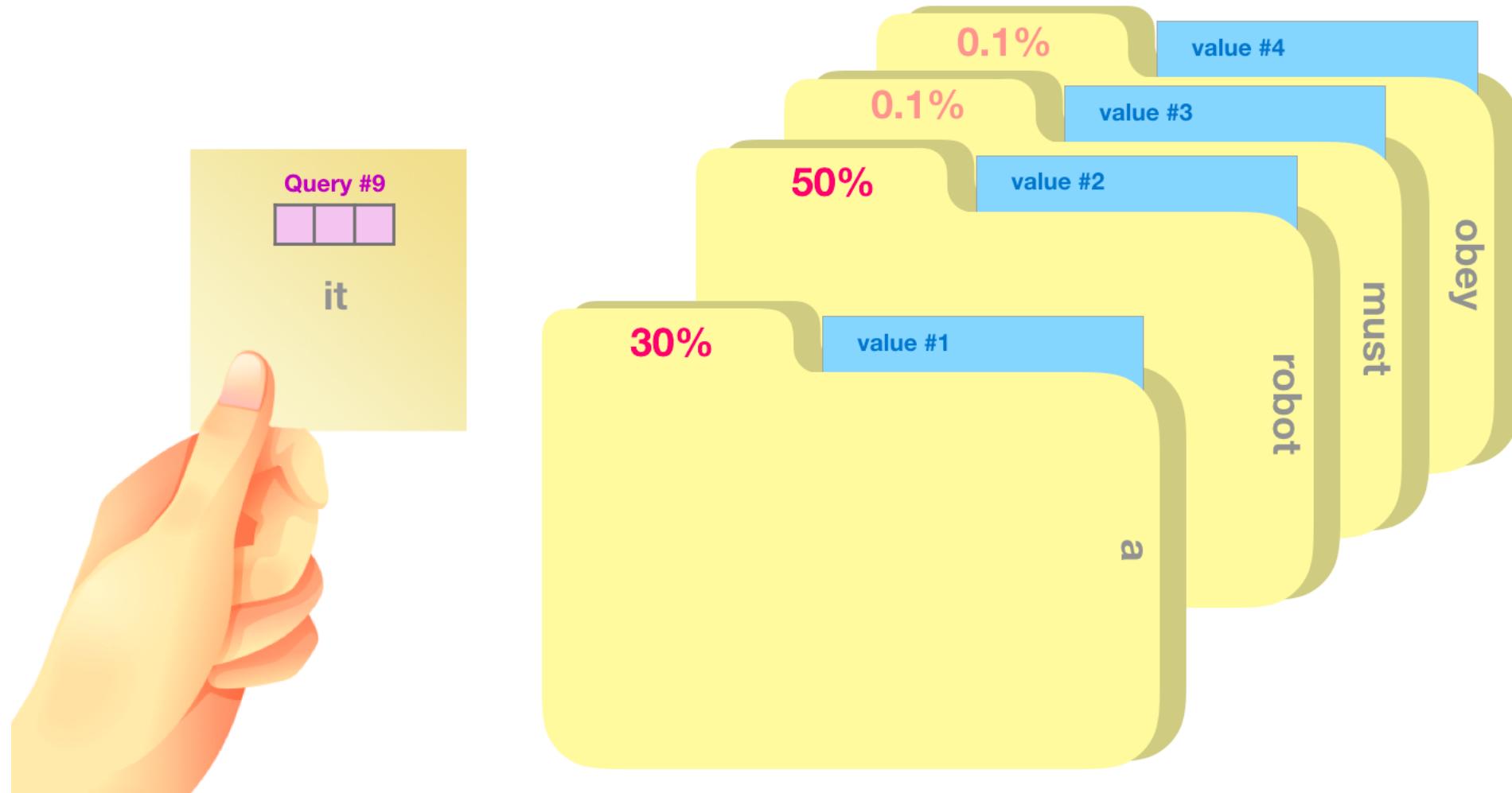
Self Attention

different weight for different context for word "kicked"

Self-Attention



Self-Attention (Filing Cabinet Analogy)



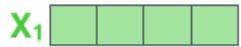
Query, Key, Value

Input

Thinking

Machines

Embedding

x_1 

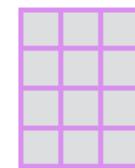
x_2 

word vector $\in \mathbb{R}^{1 \times 4}$

Queries

$$x_1 w^Q = q_1 \quad \begin{matrix} \text{purple} \\ \boxed{} \end{matrix}$$

$$x_2 w^Q = q_2 \quad \begin{matrix} \text{purple} \\ \boxed{} \end{matrix}$$



w^Q

}

Keys

$$x_1 w^K = k_1 \quad \begin{matrix} \text{orange} \\ \boxed{} \end{matrix}$$

$$x_2 w^K = k_2 \quad \begin{matrix} \text{orange} \\ \boxed{} \end{matrix}$$



w^K

weight matrix $\in \mathbb{R}^{4 \times 3}$

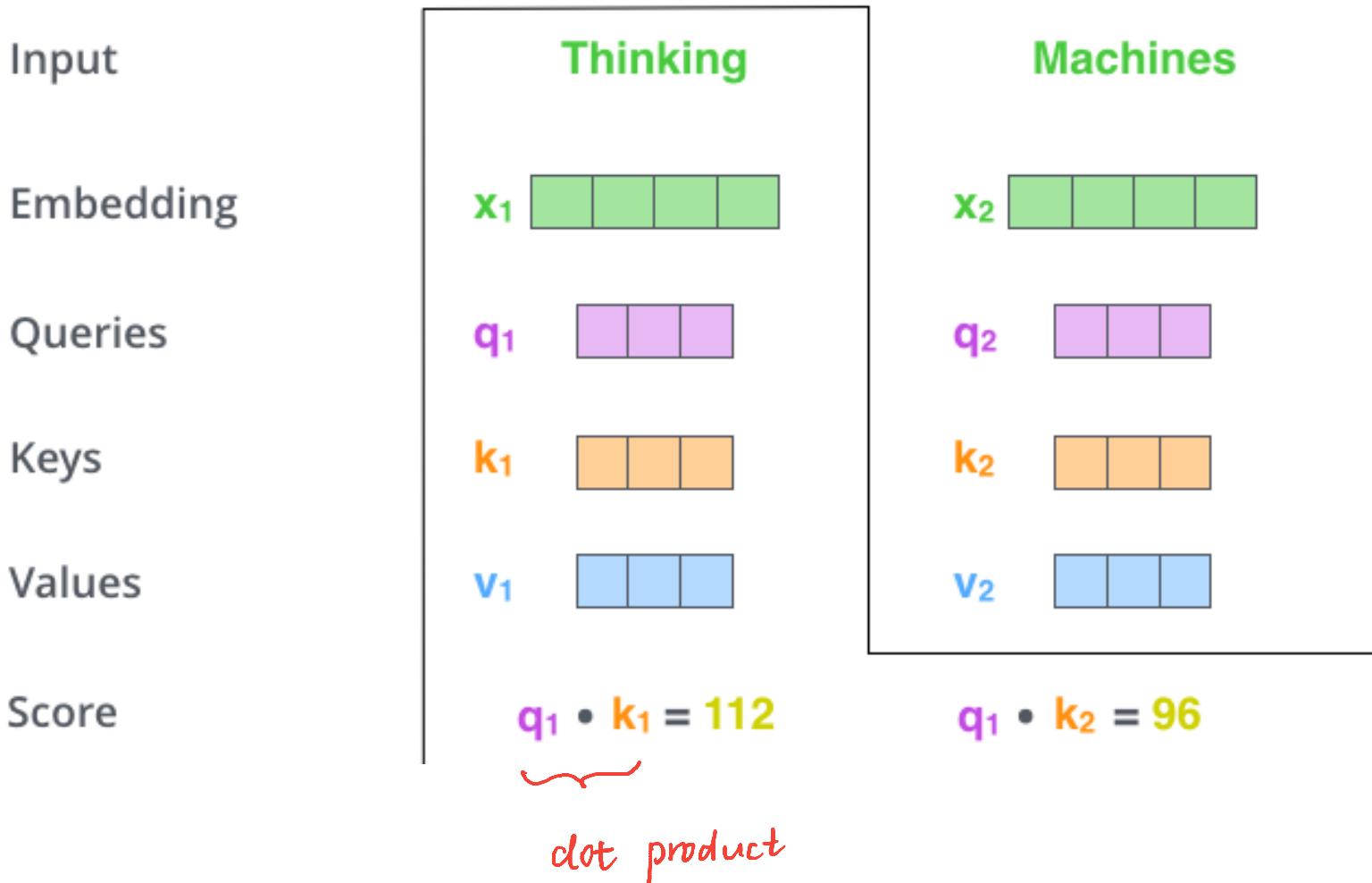
Values

$$x_1 w^V = v_1 \quad \begin{matrix} \text{blue} \\ \boxed{} \end{matrix}$$

$$x_2 w^V = v_2 \quad \begin{matrix} \text{blue} \\ \boxed{} \end{matrix}$$



w^V



Input

Embedding

Queries

Keys

Values

Score

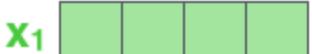
Divide by 8 ($\sqrt{d_k}$)

Softmax

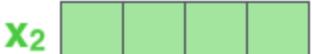
Softmax
X
Value

Sum

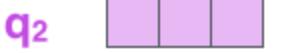
Thinking

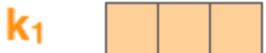
x_1 

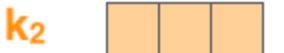
Machines

x_2 

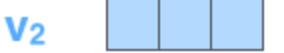
q_1 

q_2 

k_1 

k_2 

v_1 

v_2 

$$q_1 \cdot k_1 = 112$$

14

0.88

$$q_1 \cdot k_2 = 96$$

12

score

0.12

Why the scaling?

Assuming q and k are d_k dimensional vectors
The dot product is

$$q \cdot k = \sum_{i=1}^{d_k} u_i v_i,$$

This has mean 0 and variance d_k

Want to normalize to have variance 1

is *normalized*

v_1 

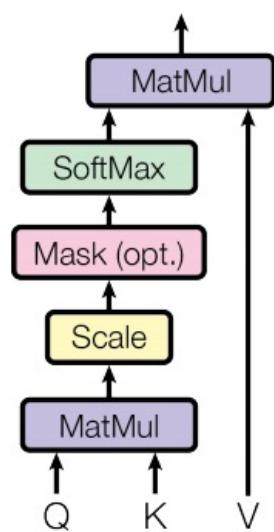
v_2 

z_1 

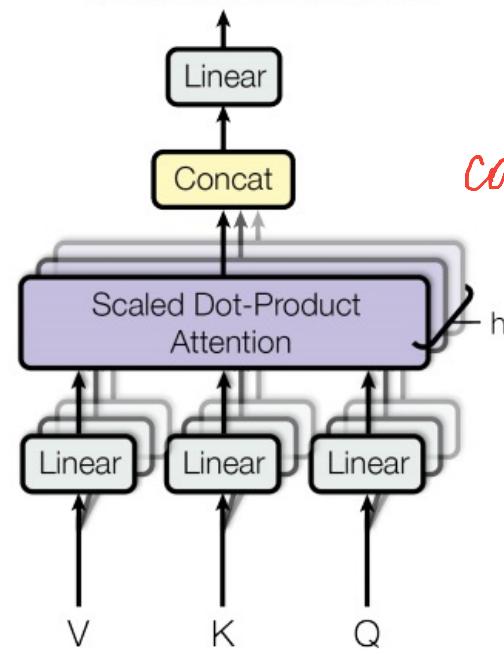
z_2 

Multihead Attention

Scaled Dot-Product Attention



Multi-Head Attention

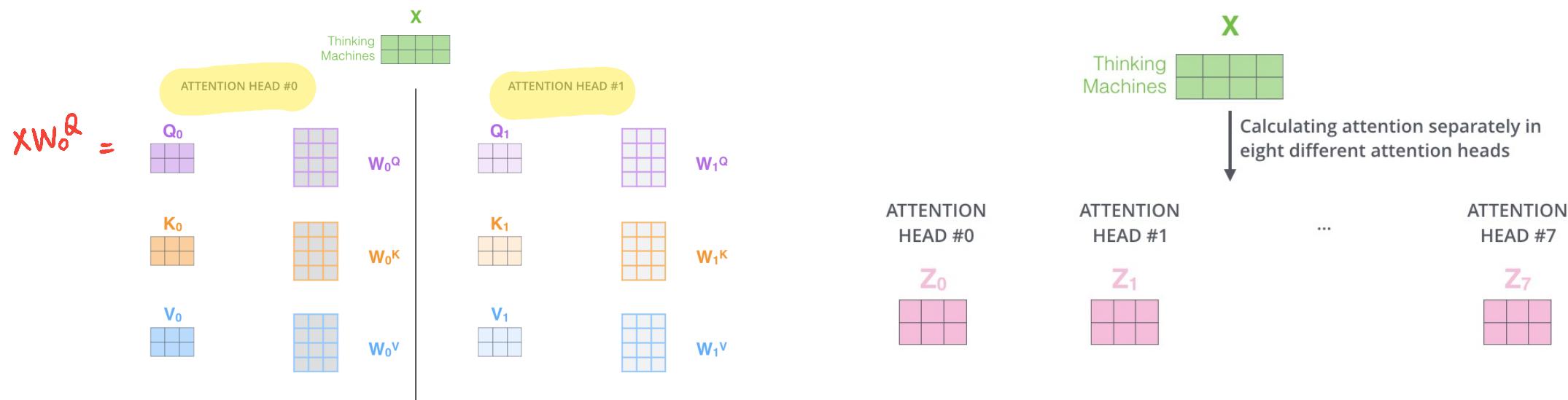


concatenate into one long encoding

From Vaswani et al.

Multi-head Attention

all linear layers



weight matrix W_0^Q and W_1^Q are different

Concatenating Multi-head Attention

1) Concatenate all the attention heads

$$\in \mathbb{R}^{2 \times 24}$$

$$Z' = \begin{matrix} Z_0 & Z_1 & Z_2 & Z_3 & Z_4 & Z_5 & Z_6 & Z_7 \end{matrix}$$


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

x

$$W^o \in \mathbb{R}^{24 \times 4}$$

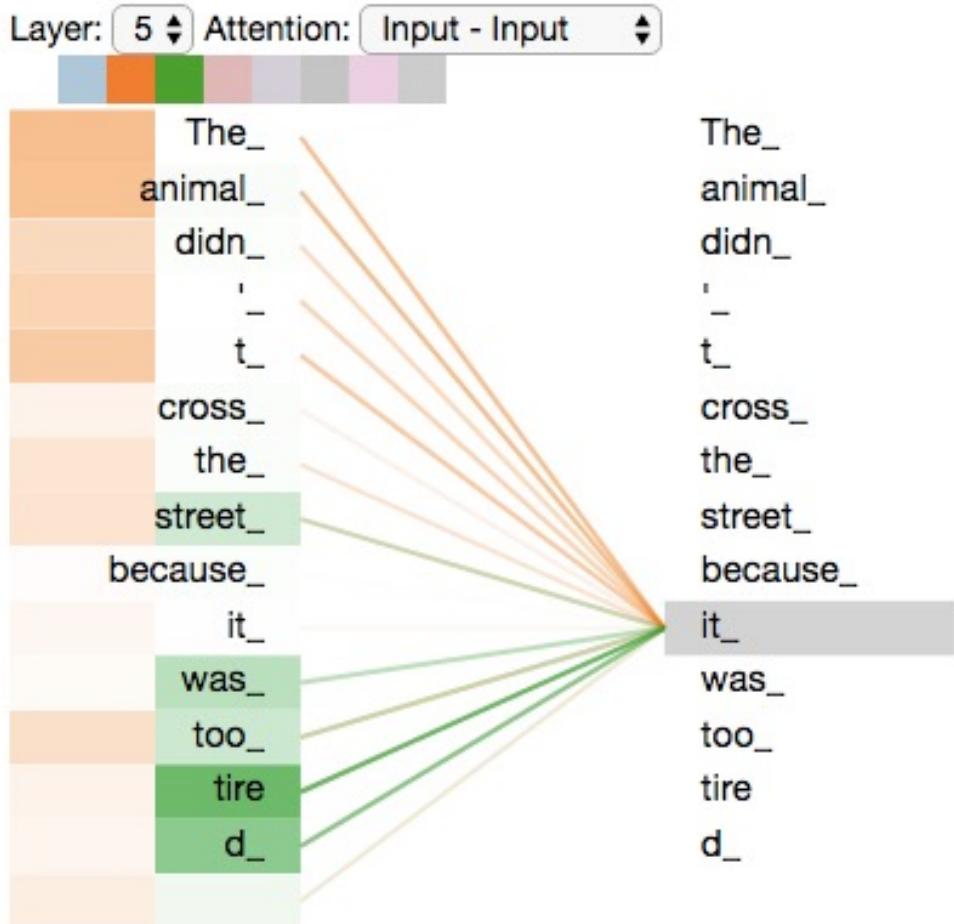


3) The result would be the Z' matrix that captures information from all the attention heads. We can send this forward to the FFNN

$$Z' W^o = \begin{matrix} Z \end{matrix} \in \mathbb{R}^{2 \times 4}$$


Attention Results

8 heads



visualize 2nd head give how much
attention to word "it"

Positional Encoding

fixed length vector $\in \mathbb{R}^d$

- Lost positional information now
- Positions can't be encoded by numbers because sequences of different lengths can't be trained this way
- Instead use a **periodic function** like sin, cos to create the embedding
*reason: want positional encoding to be **compact** rather than one-hot encoding*

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

$$\omega_k = \frac{1}{10000^{2k/d}} \quad t \text{ is } t\text{th word in seq}$$

Total positional encoding

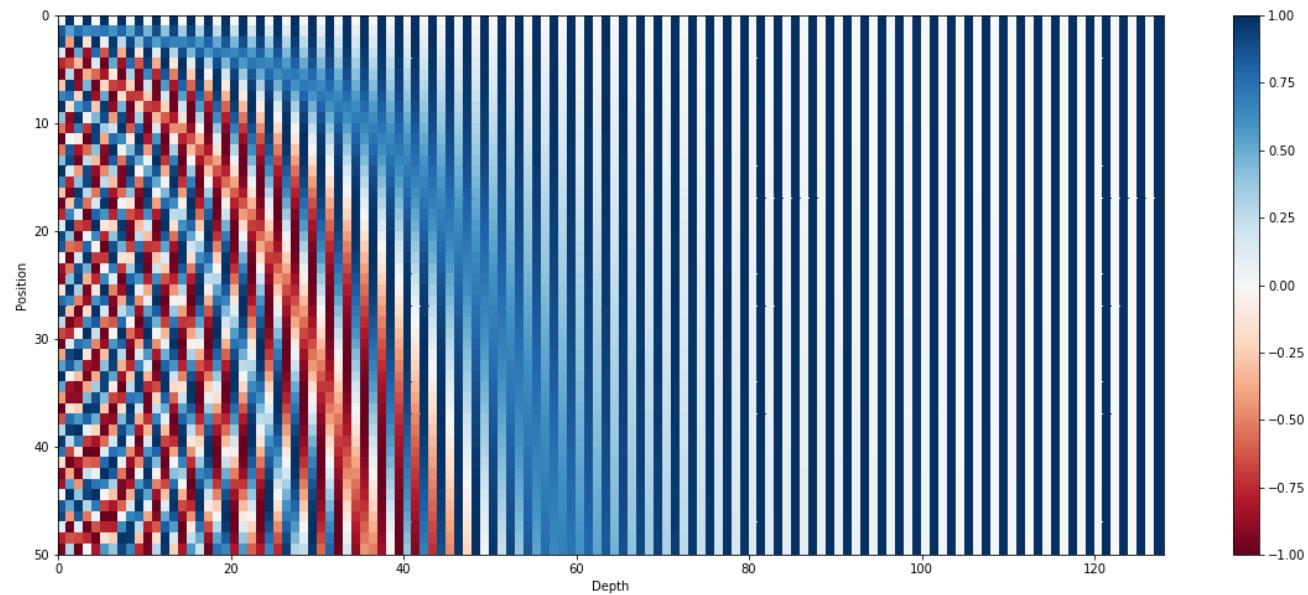
$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \\ \vdots \\ \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

Binary vs Positional Encoding

Binary

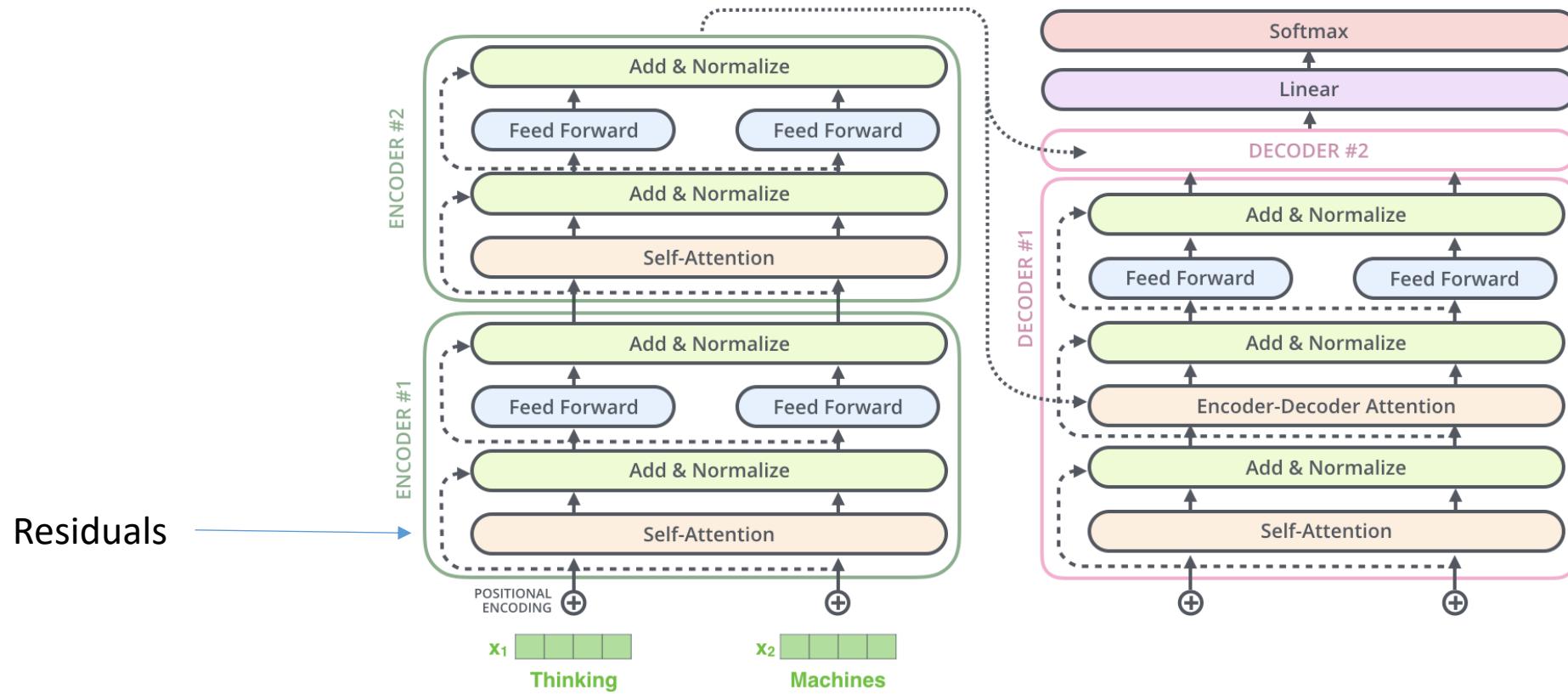
0 :	0	0	0	0
1 :	0	0	0	1
2 :	0	0	1	0
3 :	0	0	1	1
4 :	0	1	0	0
5 :	0	1	0	1
6 :	0	1	1	0
7 :	0	1	1	1
8 :	1	0	0	0
9 :	1	0	0	1
10 :	1	0	1	0
11 :	1	0	1	1
12 :	1	1	0	0
13 :	1	1	0	1
14 :	1	1	1	0
15 :	1	1	1	1

Positional

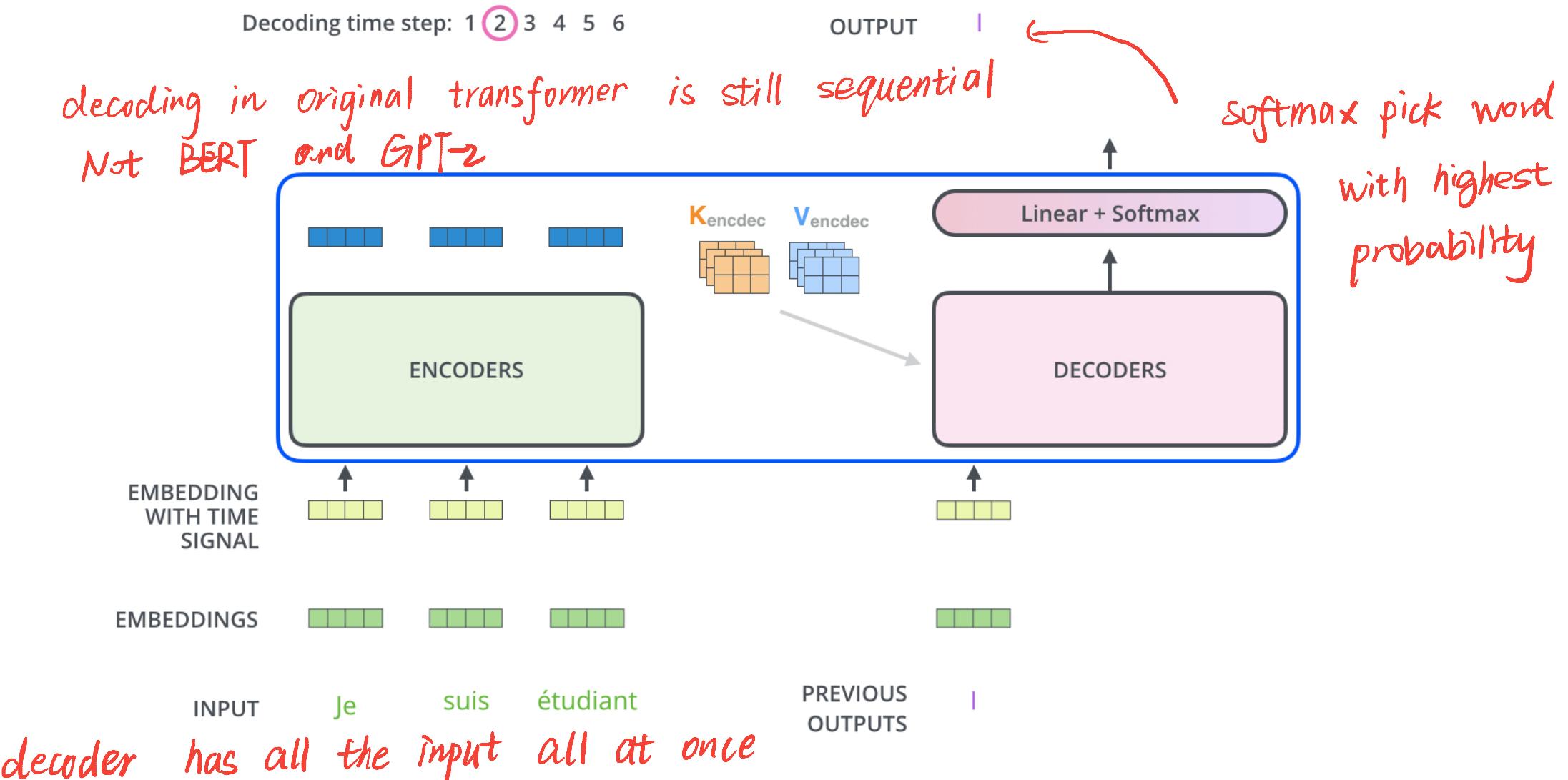


Transformer

reason that transformer don't have memory issue:
all the inputs together with positional encoding are available for decoders



Computation in a Transformer



New AI fake text generator may be too dangerous to release, say creators

The Elon Musk-backed nonprofit company OpenAI declines to release research publicly for fear of misuse



BPT-3 generate realistic fake news

Next paragraph prediction

Feed it the opening line of George Orwell's Nineteen Eighty-Four - "It was a bright cold day in April, and the clocks were striking thirteen" - and the system recognises the vaguely futuristic tone and the novelistic style, and continues with:

"I was in my car on my way to a new job in Seattle. I put the gas in, put the key in, and then I let it run. I just imagined what the day would be like. A hundred years from now. In 2045, I was a teacher in some school in a poor part of rural China. I started with Chinese history and history of science."

Language Modeling

- An unsupervised training task that models the probability of a word given its context

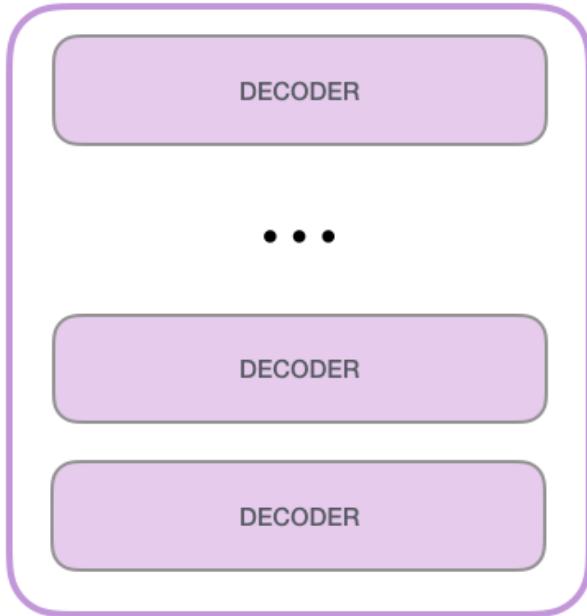
$$p(x) = \prod_{i=1}^n p(s_i | s_1, \dots, s_{i-1})$$

- Instead of task-oriented training, language models are just trained to predict probabilities from corpuses

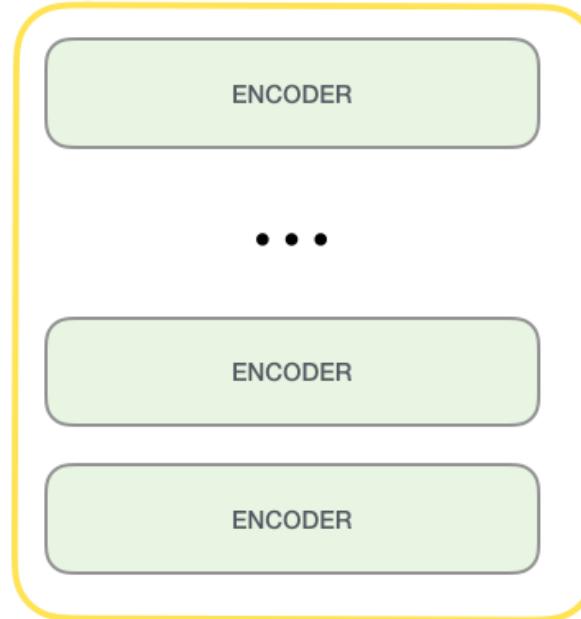
Well-known Language Models



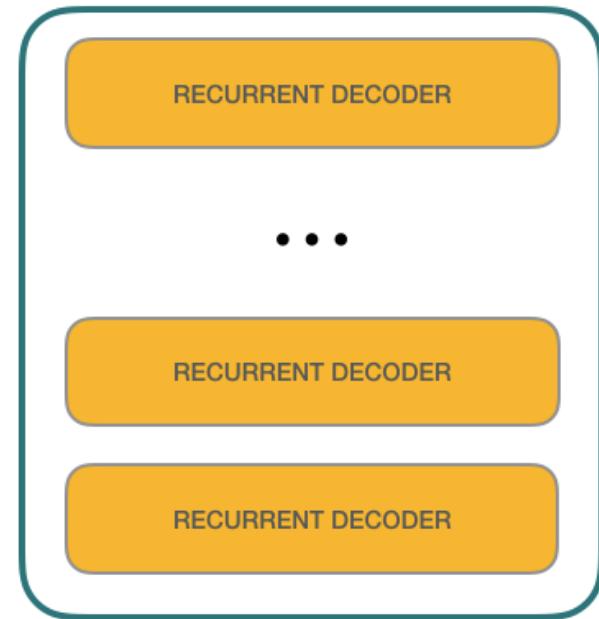
GPT-2



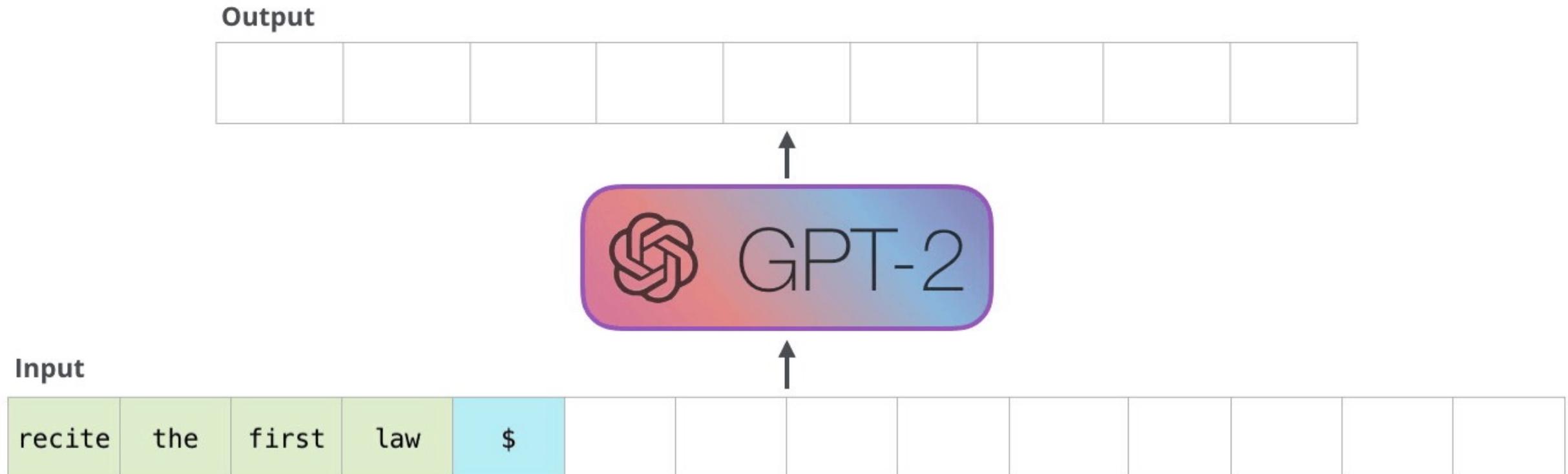
BERT



TRANSFORMER XL



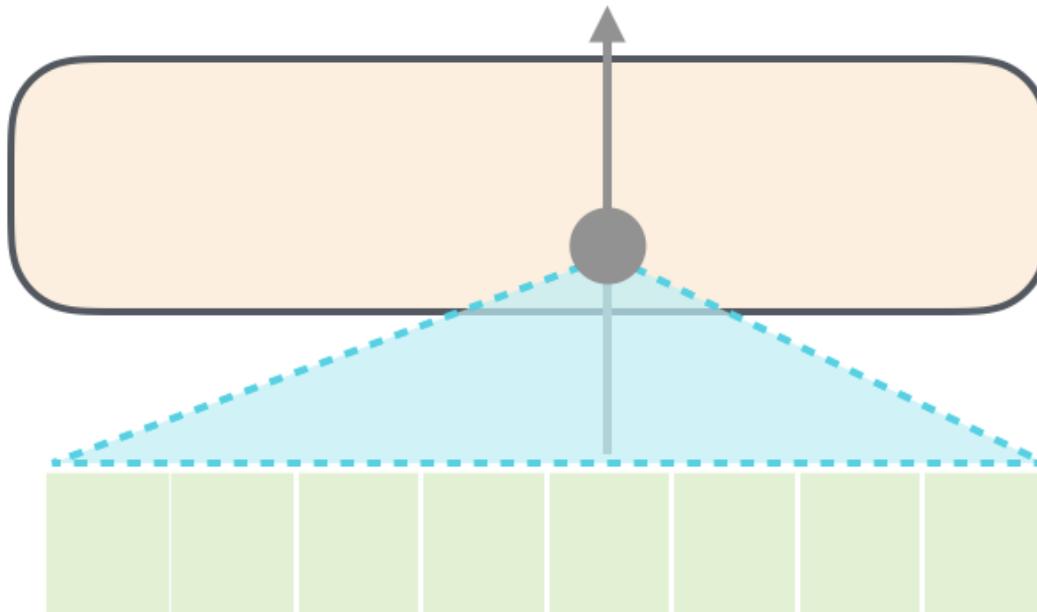
Pretrained Autoregression \rightarrow next word prediction



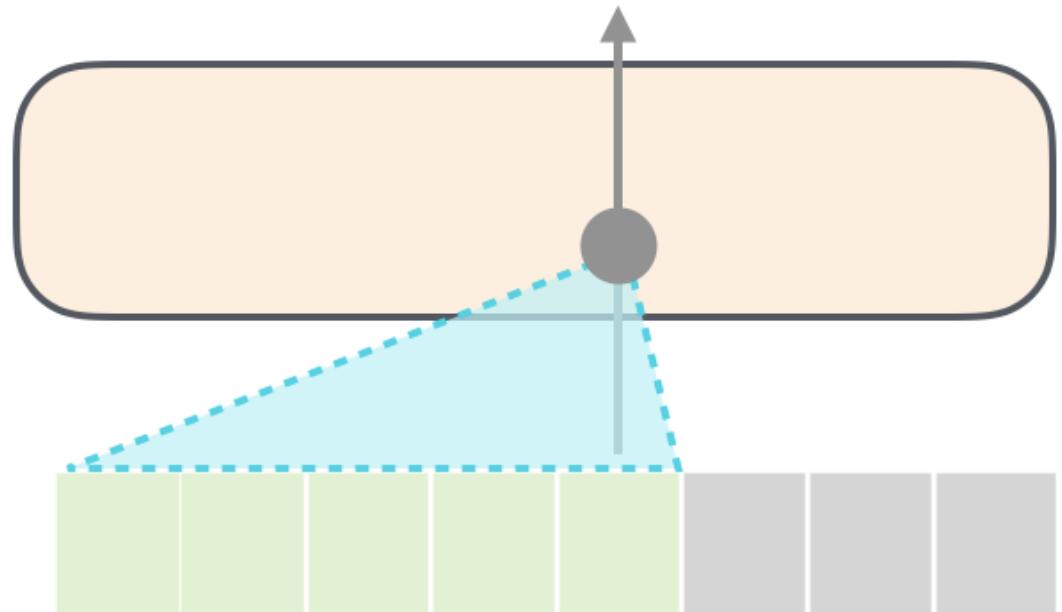
Trained sequentially with newly generated output becoming part of the next input sequence

Masked Attention (Mask everything before)

Self-Attention



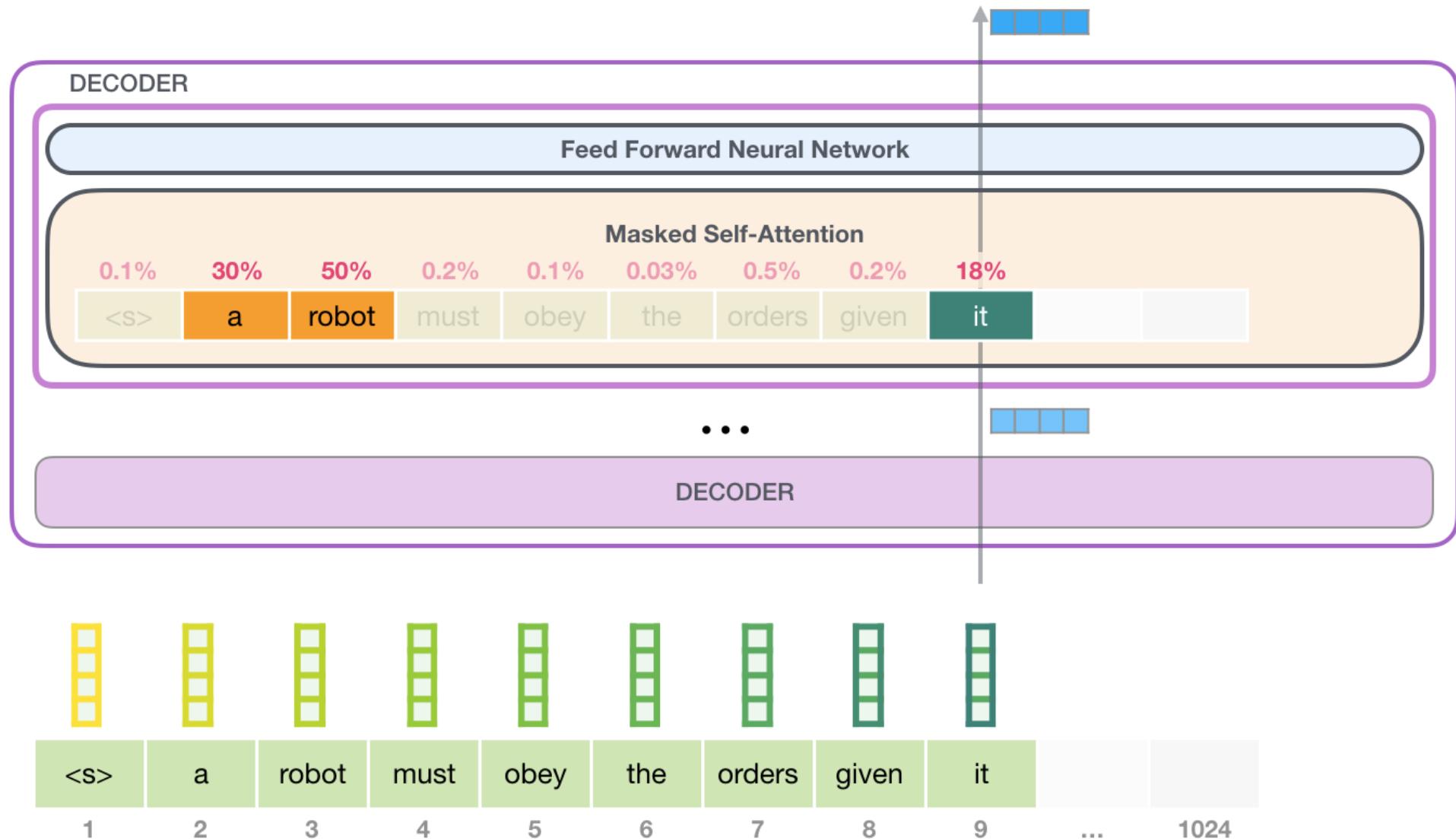
Masked Self-Attention



Word is only allowed to pay attention to
words come before it

Role of self-attention: context

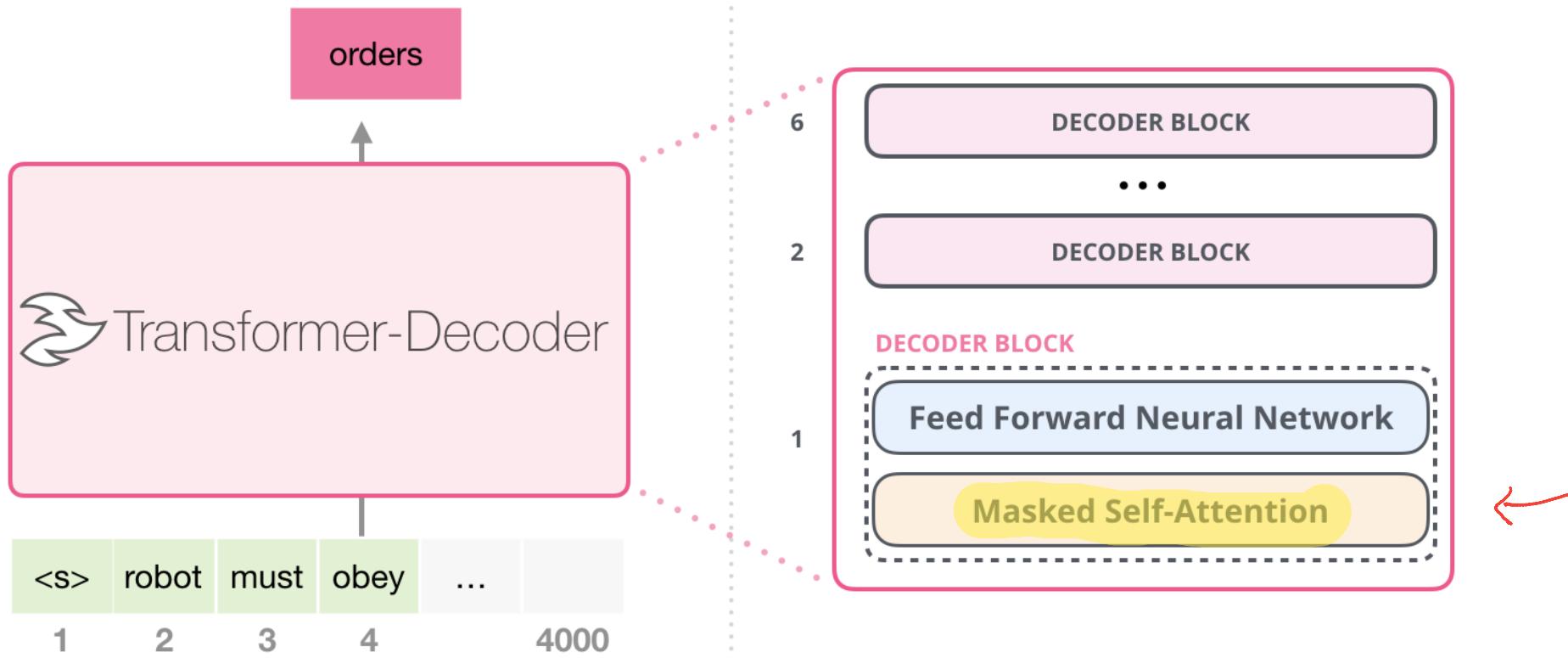
take into account past context



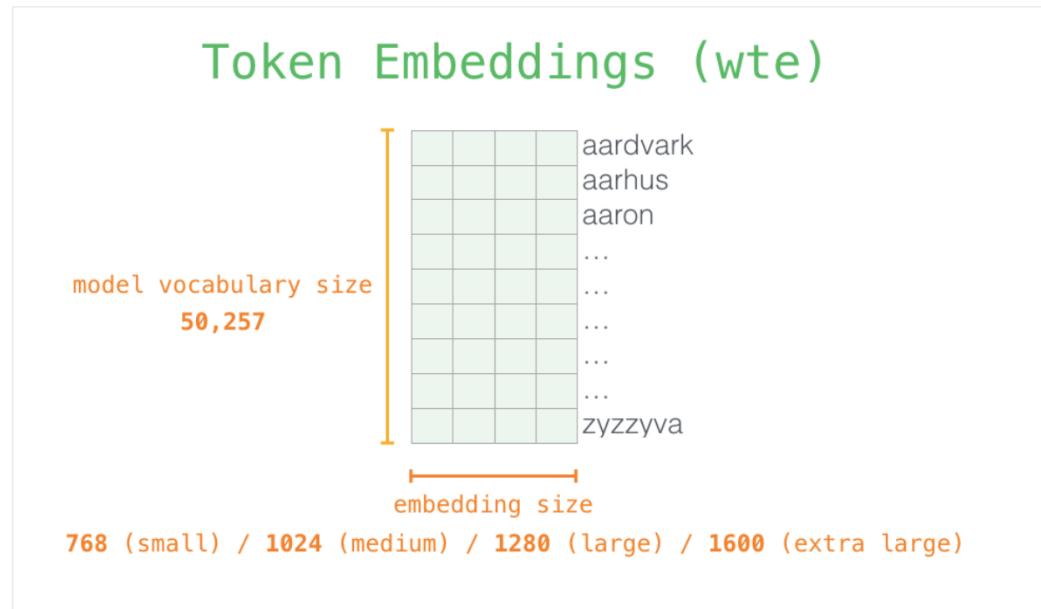
GPT-2 only use decoder part of original transformer , it's fine coz encoder blocks are redundant , very deep in original transformer

Decoder Only Block

Different: use **masked self-attention** in decoder

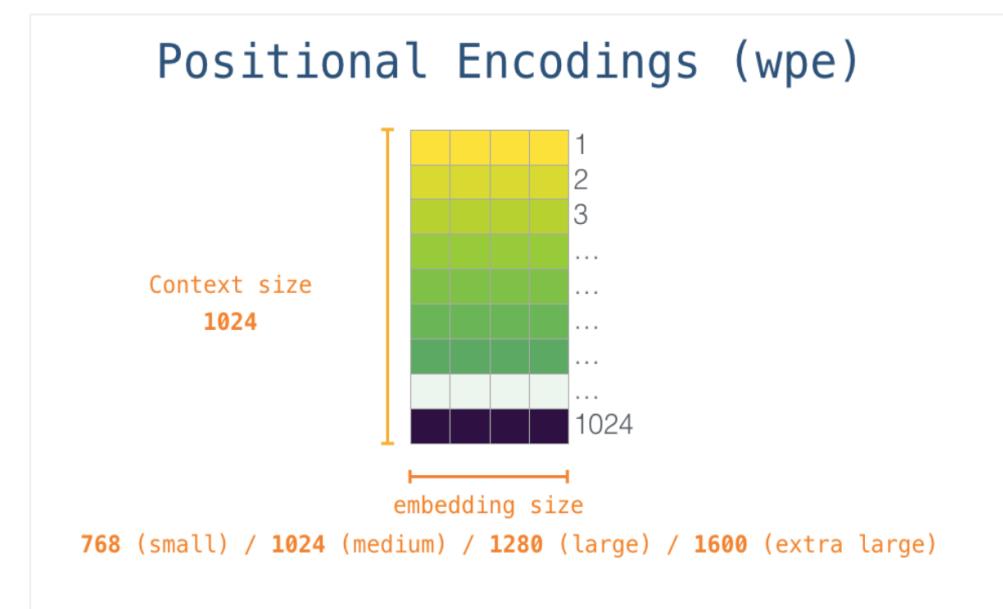


Model Input (with starting token)



Each row is a word embedding: a list of numbers representing a word and capturing some of its meaning. The size of that list is different in different GPT2 model sizes. The smallest model uses an embedding size of 768 per word/token.

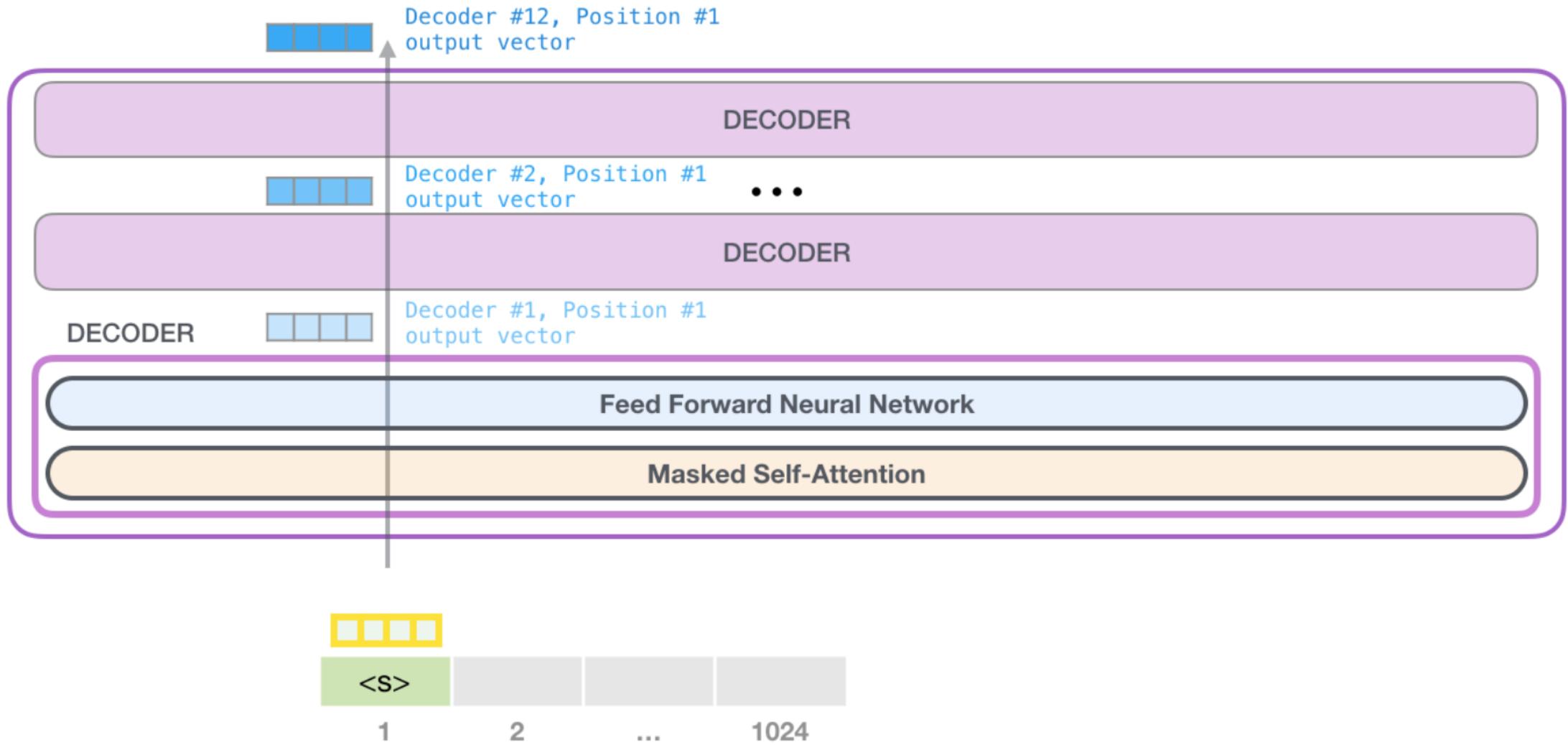
each row is a token embedding

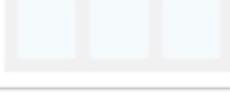
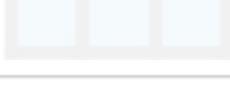


input sentence length : 1024



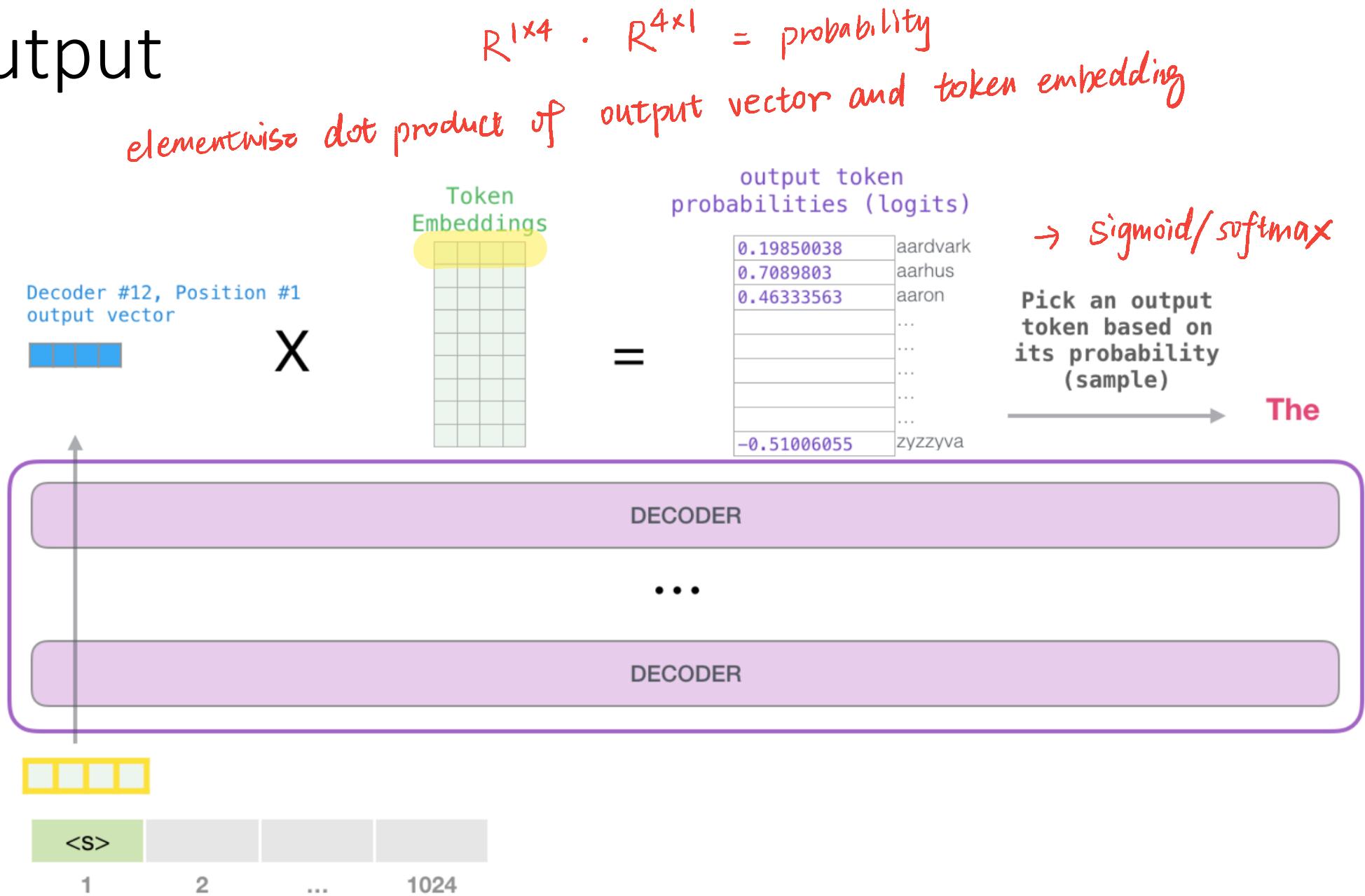
Processing



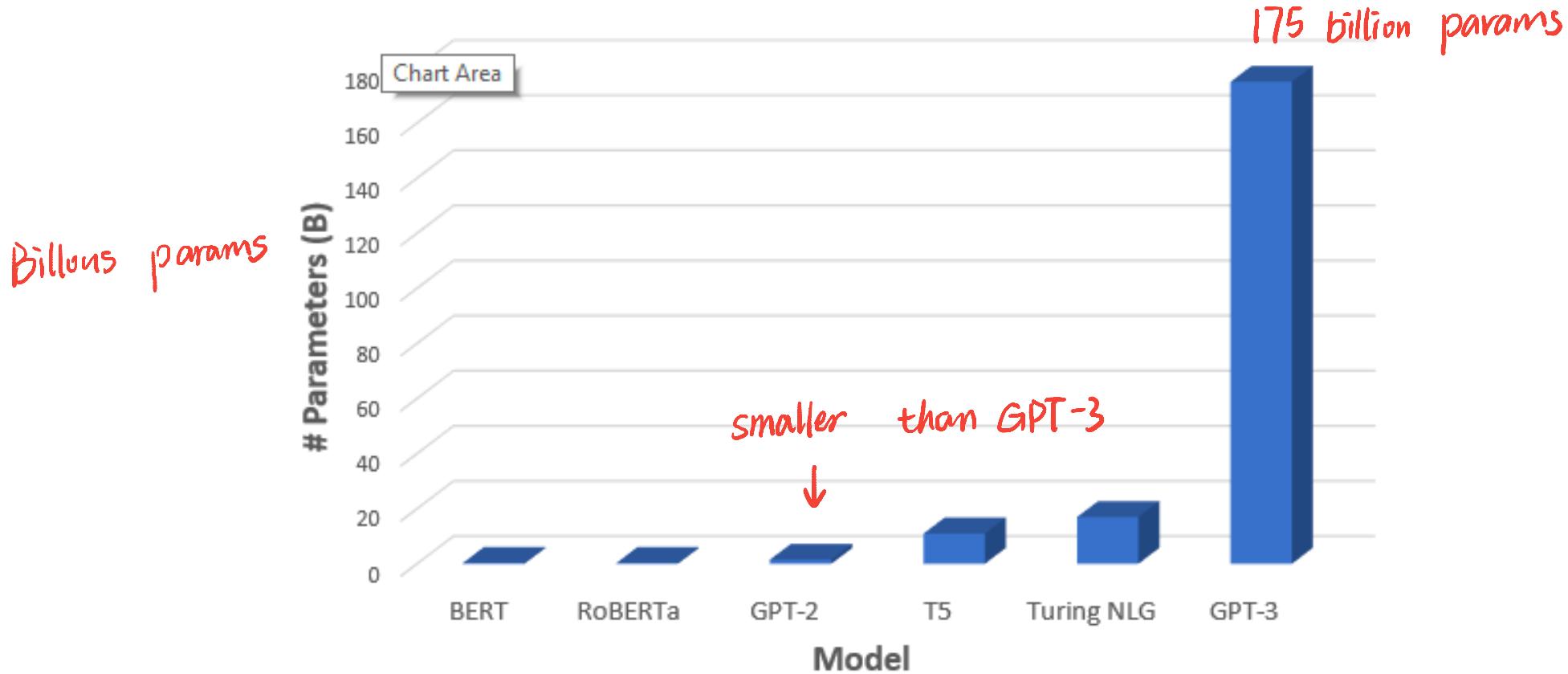
Word	Value vector	Score	Value X Score
<S>		0.001	
a		0.3	
robot		0.5	
must		0.002	
obey		0.001	
the		0.0003	
orders		0.005	
given		0.002	
it		0.19	
		Sum:	

multiply

Model Output



GPT Models



<https://demo.allennlp.org/next-token-lm>

<https://openai.com/api/>

MuseNet-Music Generator uses GPT2



Starts with a music token and generates next token in Sequence

Tokens contain information on:

1. pitch,
2. volume,
3. Instrument

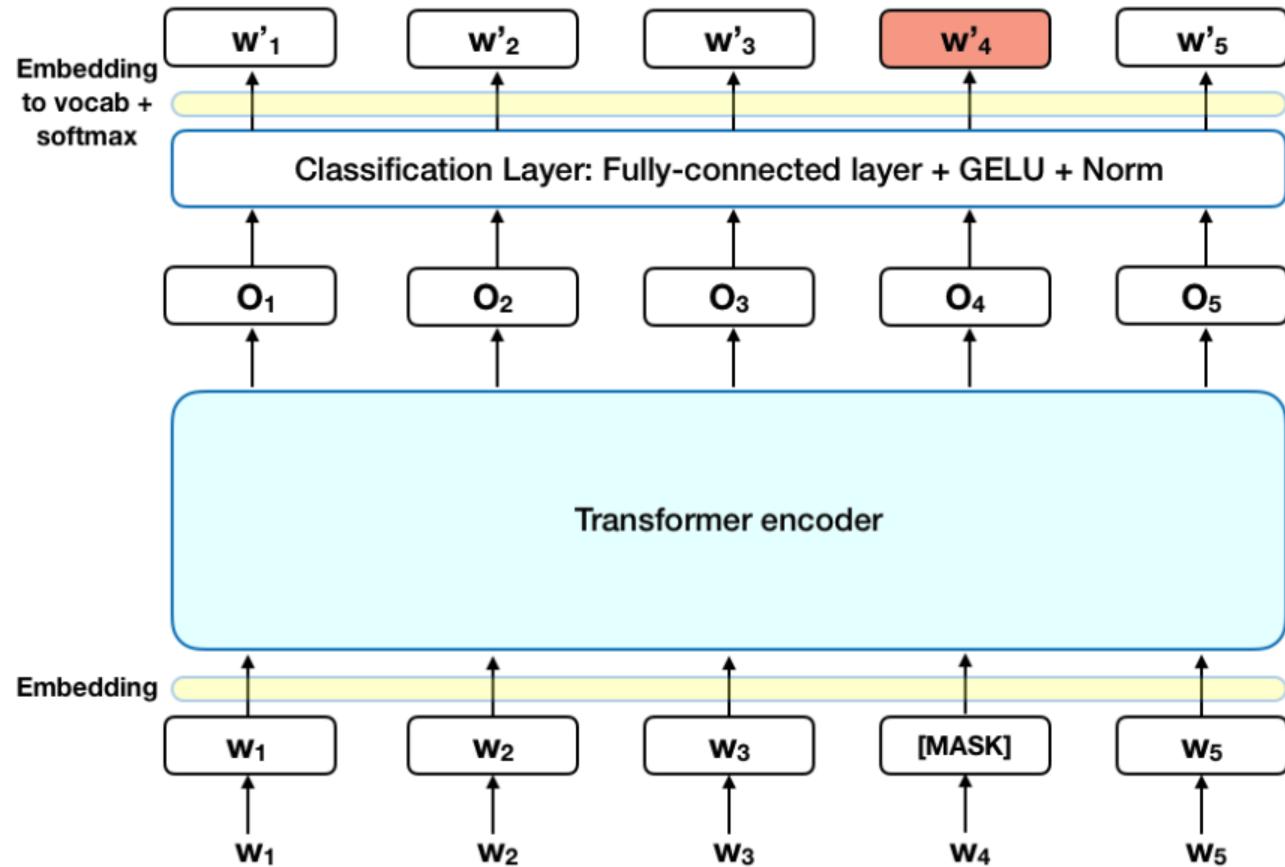
音高
音量
乐器

Trained on a large archive of MIDI files

<https://openai.com/blog/musenet/>

Bert—Encoder Only

Bert only use encoder part of original transformer
and use self-attention not masked self-attention



Trained to predict masked words

GPT-2 and Bert are all
trained unsupervised
predict the next word
or predict the masked
word

Few Shot Learning

The three settings we explore for in-context learning

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



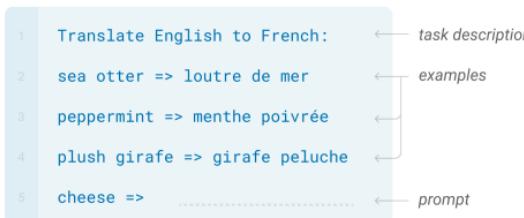
One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Traditional fine-tuning (not used for GPT-3)

Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



Reading list

- <https://towardsdatascience.com/attention-and-its-different-forms-7fc3674d14dc>
- <http://jalammar.github.io/illustrated-transformer/>
- Vaswani et al. 2019 <https://arxiv.org/abs/1706.03762> ↪ Transformer
- Radford et al. (GPT2 paper)
- Brown et al. (GPT3 paper)
- Child et al. (Sparse Transformer)
- <https://jalammar.github.io/illustrated-gpt2/> ↪ illustrated GPT 2
- <https://medium.com/walmartglobaltech/the-journey-of-open-ai-gpt-models-32d95b7b7fb2>