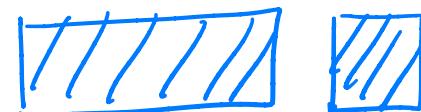


Hilary Mogan  
15 May 21



- \* For language modeling data  
need to
- a) tokenized  $\Rightarrow$  into word
  - b) Converted into word embedding
  - c) Positional embedding
  - d) Segment/Sentence embedding.

\* GPT-2  $\Rightarrow$  Language model that  
Can be used for next word  
Prediction.



given the  
Sequence  
so far      Prediction

500 MB  
storage for its  
parameters



117M Parameters

345M Parameters

762M Parameters

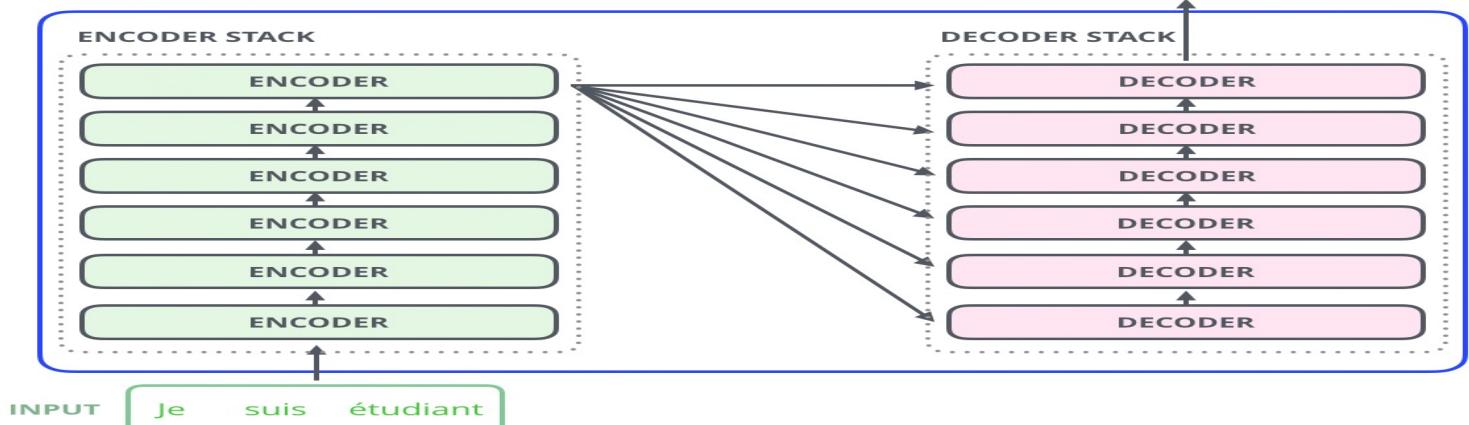
1,542M Parameters



# THE TRANSFORMER

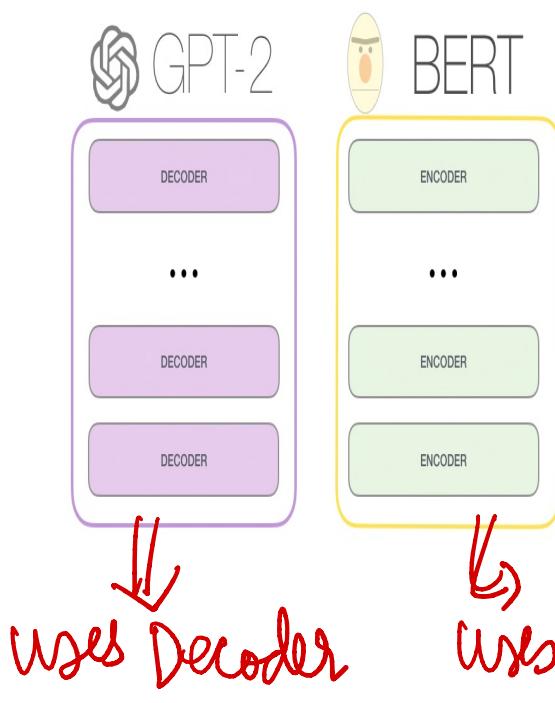
I am a student

OUTPUT



Transformer have an encoder & decoder architecture to tackle Machine translation.

✳ Subsequent models utilize either encoder/decoder or their variants.



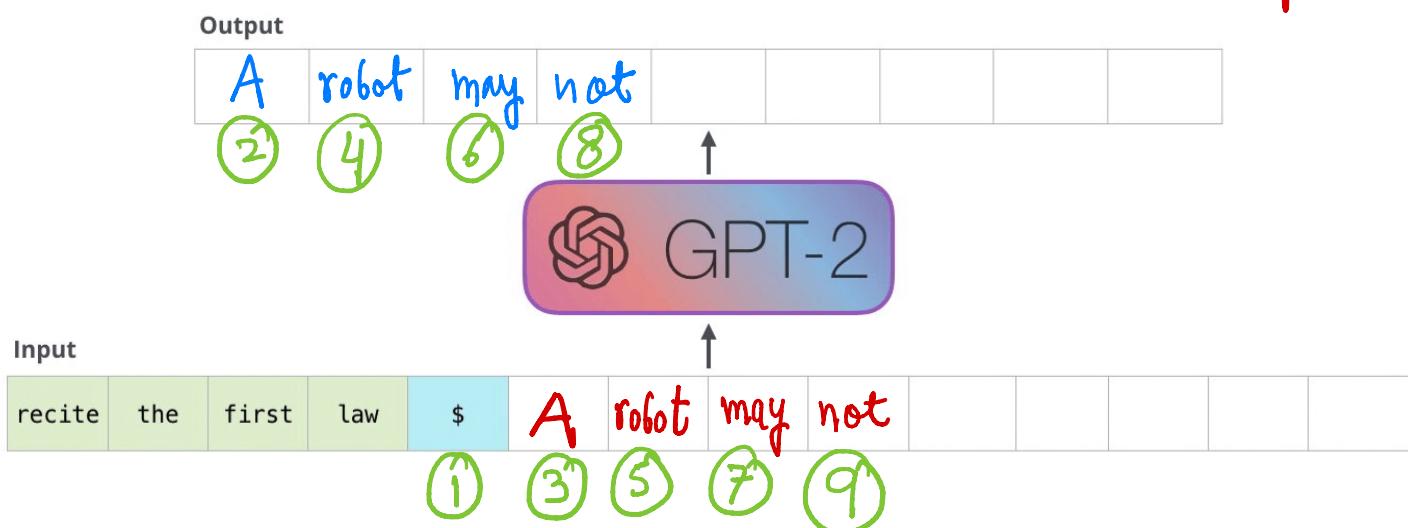
✳ Bert Can be used for better features for any downstream task.

✳ GPT-2 Can be used for language generation.

# GPT-2 O/P one token at a time.



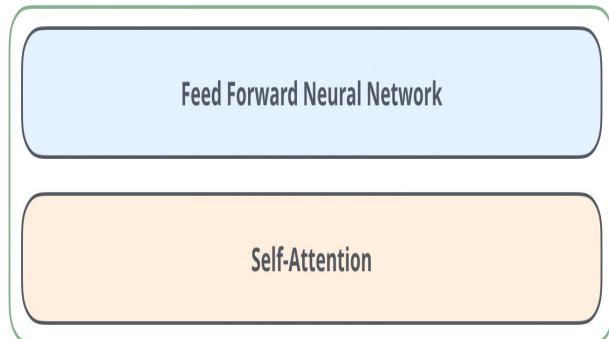
Bert looks/produce all at once  $\Rightarrow$  Never Bidirectional.  
Token after production got added to i/p: Auto-regression



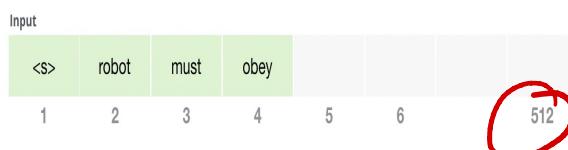
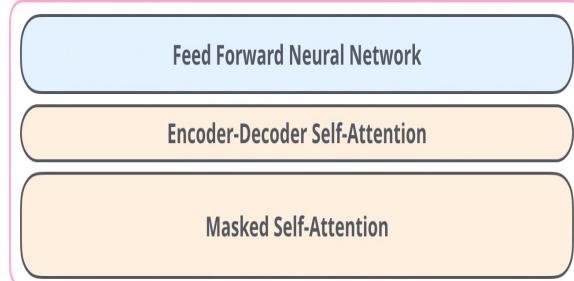
## Transformers En/Decoder



ENCODER BLOCK



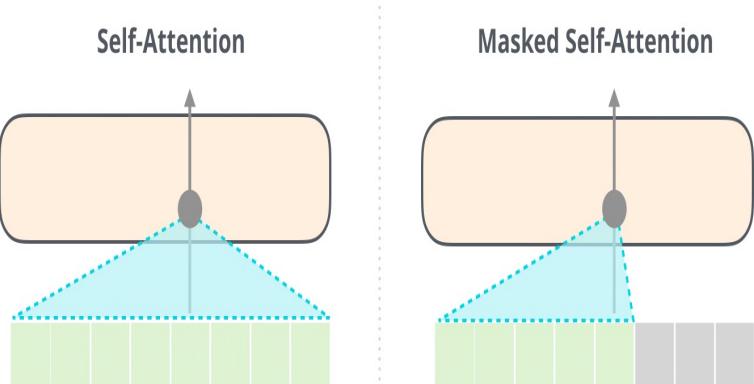
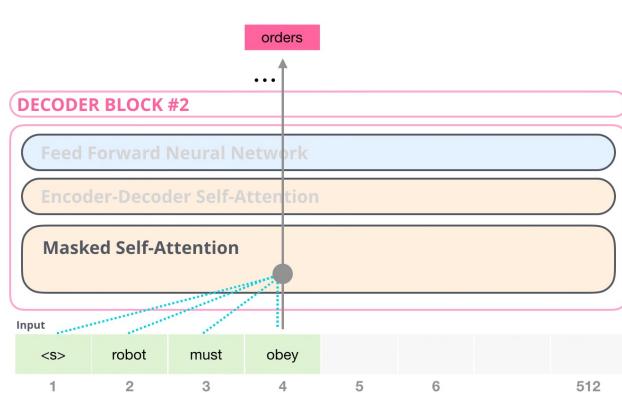
DECODER BLOCK



512 (Size) for original transformer

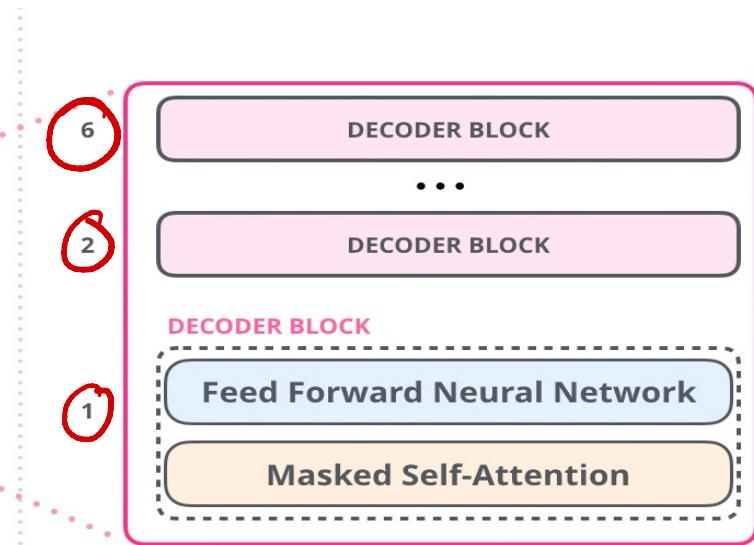
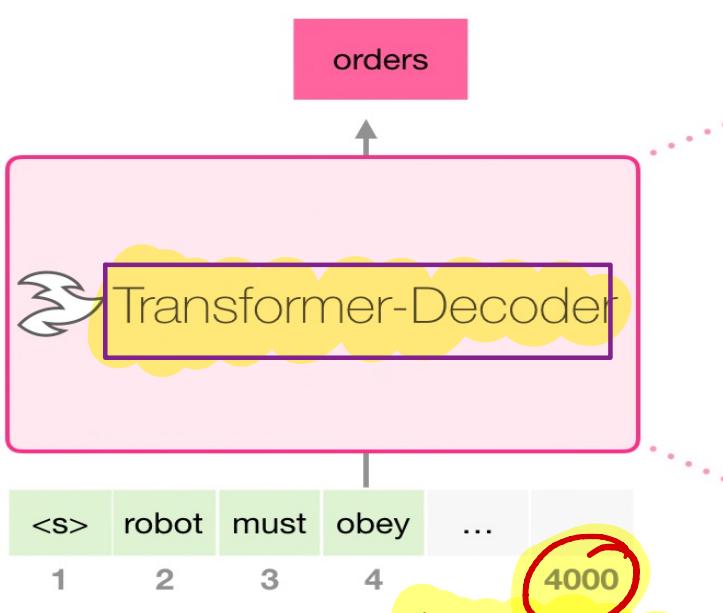
Decoder: En/De SA layer allows to learn attention to specific segments from the encoder.

\* Also decoder's (SA) layer is implemented as Masked SA.



BERT  
(Encoder)

GPT-2  
(Decoder)

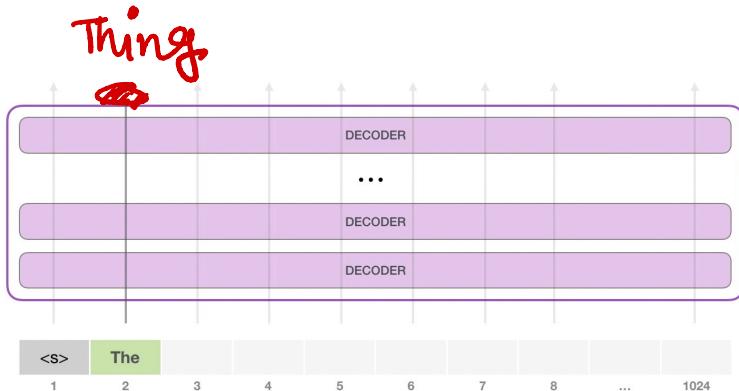


512 → 4000

(Also 2<sup>nd</sup> Cross Attention is dropped)

Subsequently Transformer Decoders are used in GPT-2 to create a language model that predicts Character by character (word by word)

## \* GPT-2 (Working) {Trained model}



\* Give it <S> and it will start spelling next word one by one.

\* final o/p vector

is scored against model vocabulary (50k)

\* Instead of choosing (Max-Prob) have a (top-k) & sample from it.



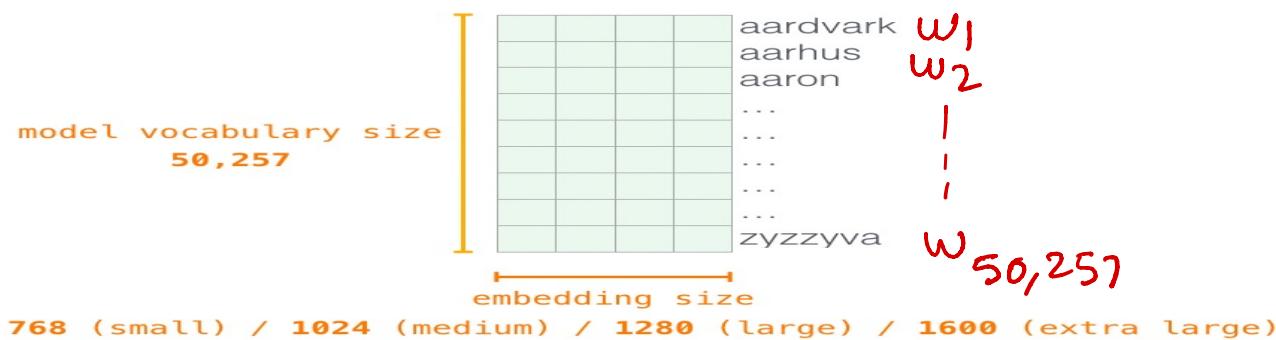
This will help to generate new data (variety)

Each layer of GPT-2 has retained its own interpretation of the first token and will use it in processing the second token (we'll get into more detail about this in the following section about self-attention). GPT-2 does not re-interpret the first token in light of the second token

# I Embeddings/Encoding

## \* Input Encoding

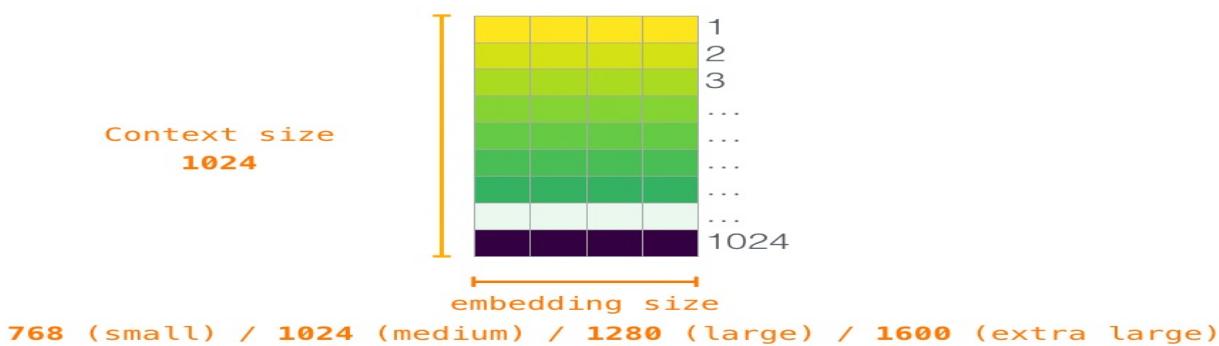
### Token Embeddings (wte)



Embedding Size for different models.

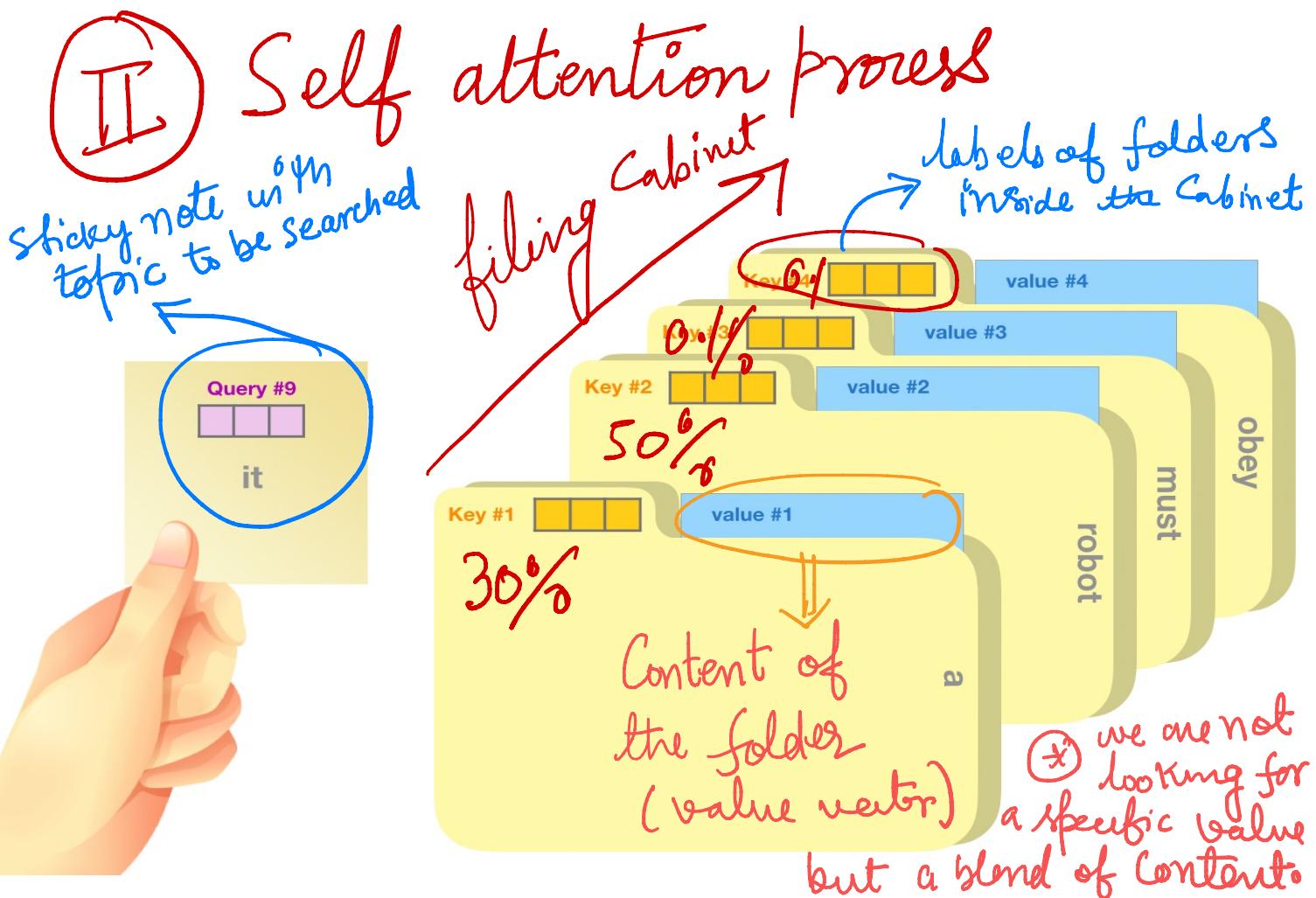
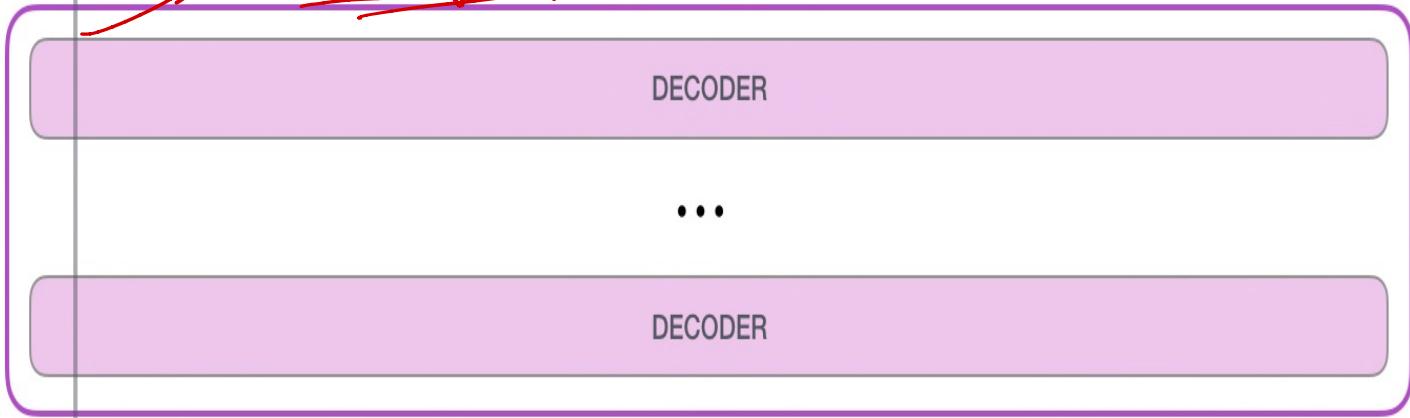
## \* Positional Encoding

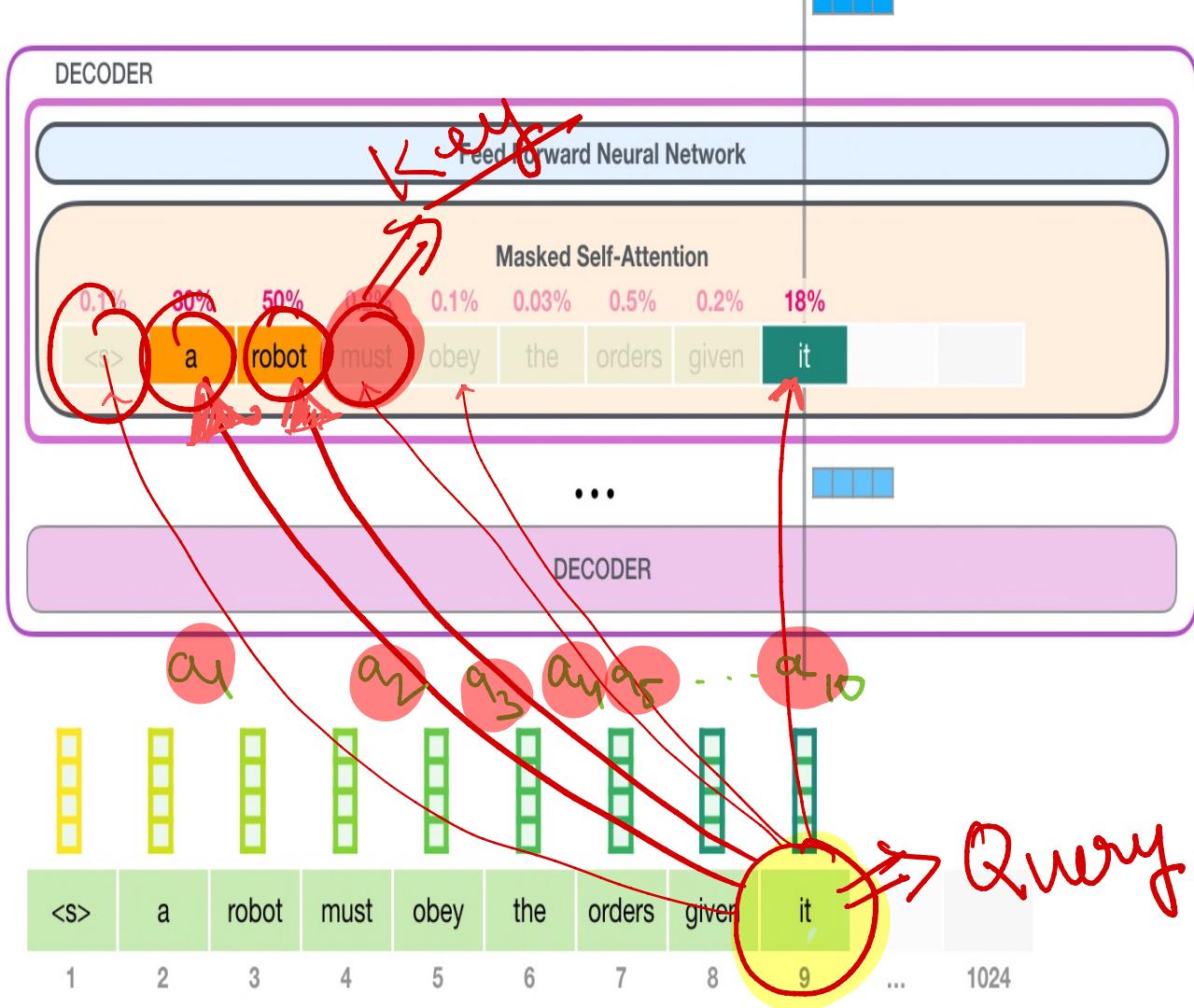
### Positional Encodings (wpe)



Both are added to get final i/p encoding.

# Journey up the Stack





for some ( $i^{th}$ )  
word ( $i^{th}$ ),  
for all ( $j^{th}$ )  
word

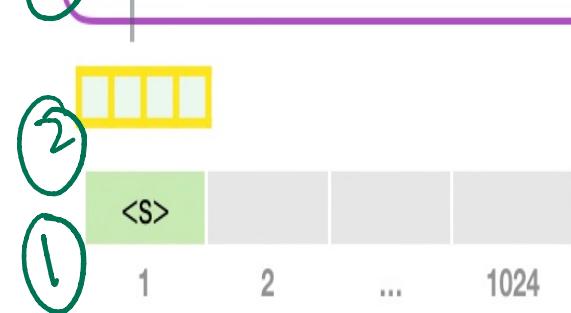
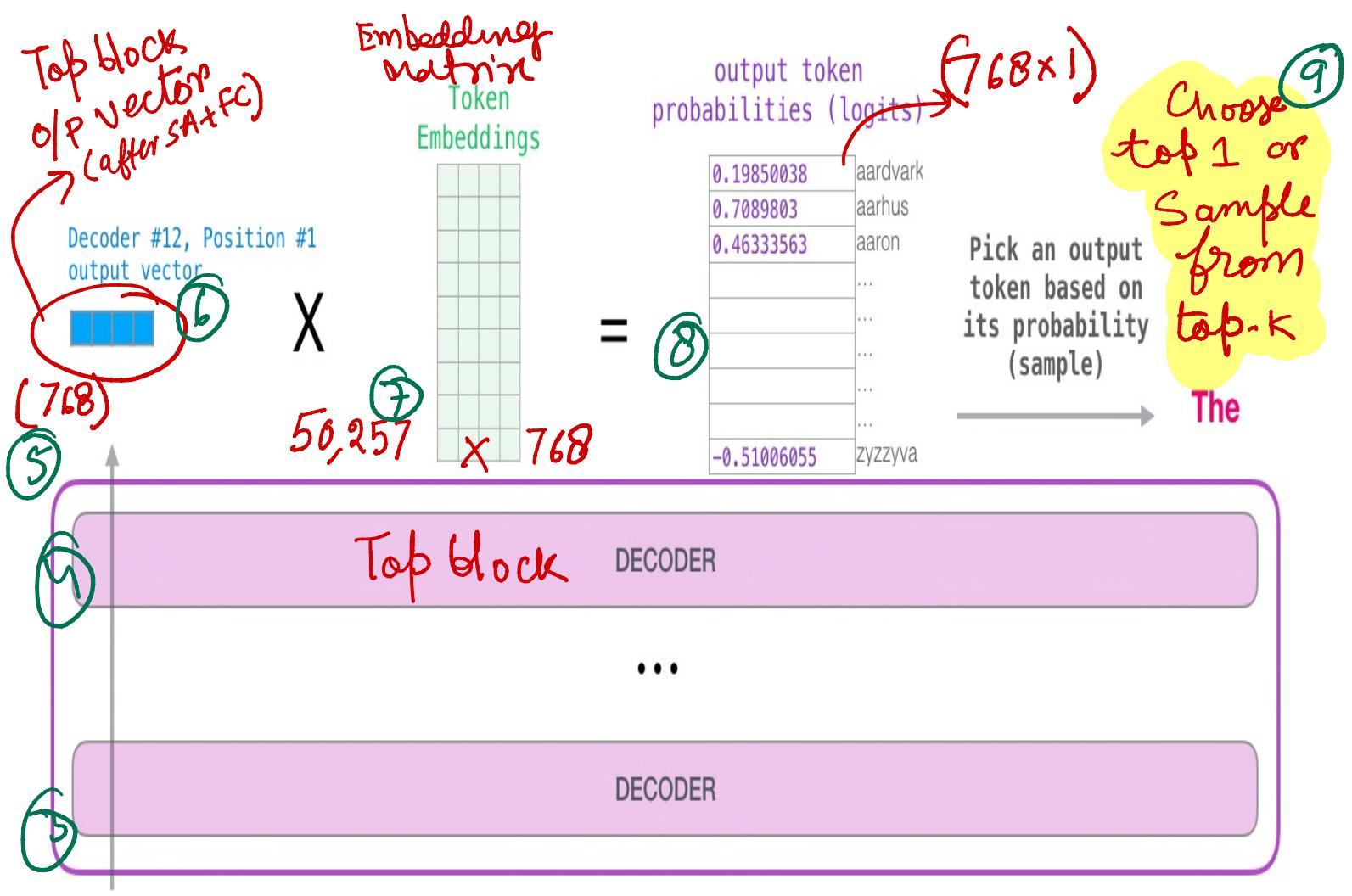
$$a_{ij} = \text{Query}_i * \text{Key}_j$$

$$O/P_i = a_{ij} * \text{Value}_j$$

Word	Value vector	Score	Value X Score
$< s >$	[blue]	0.001	[light blue]
a	[blue]	0.3	[blue]
robot	[blue]	0.5	[blue]
must	[blue]	0.002	[light blue]
obey	[blue]	0.001	[light blue]
the	[blue]	0.0003	[light blue]
orders	[blue]	0.005	[light blue]
given	[blue]	0.002	[light blue]
it	[blue]	0.19	[blue]

(They got most of  
the weight)

Why  
it is not  
very high.  
 $O/P_{ith}$  Word

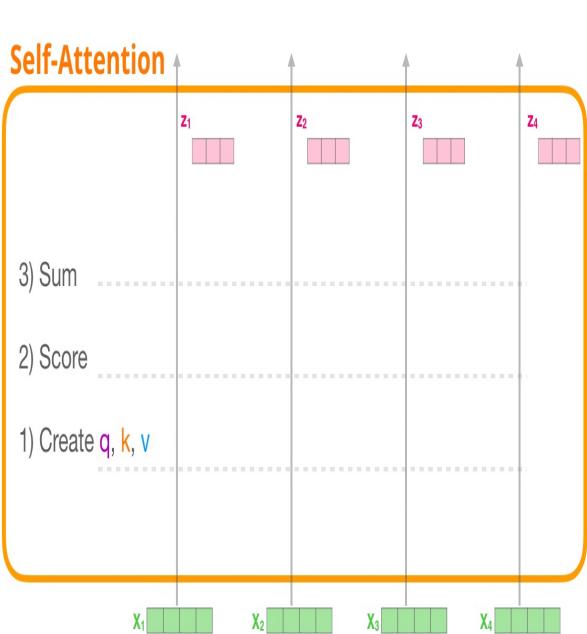


With these ⑨ steps, model completed an iteration resulting in outputting a single word.

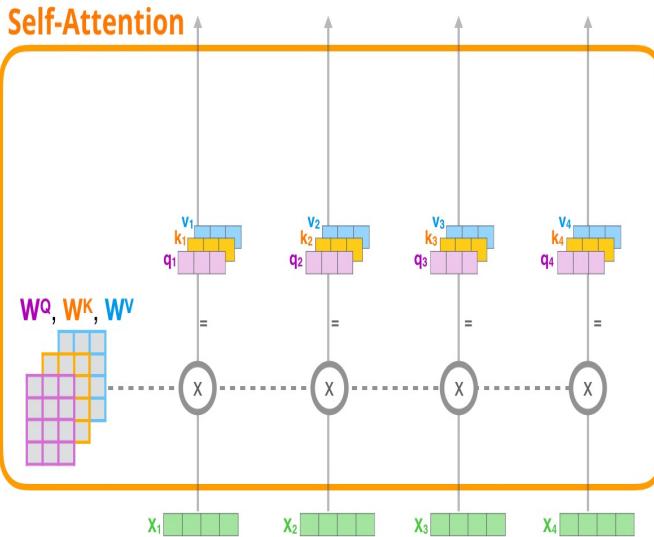
Model continues to iterate till it generates 1024 entire content or EOS.

\* GPT-2 uses Byte pair encoding (BPE)  $\Rightarrow$  to create tokens in its vocabulary.

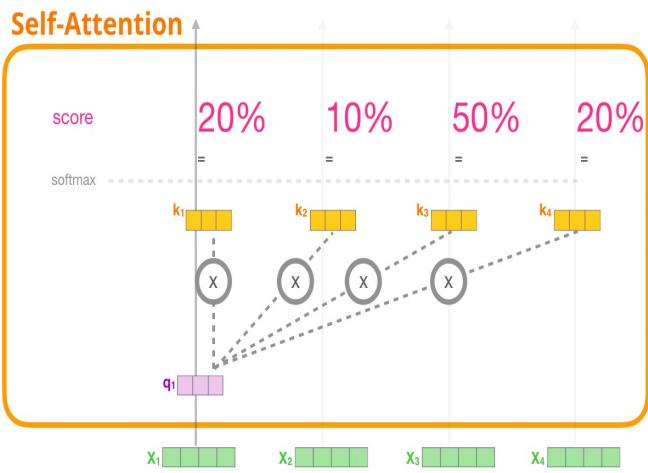
## Visual Self Attention (Recap)



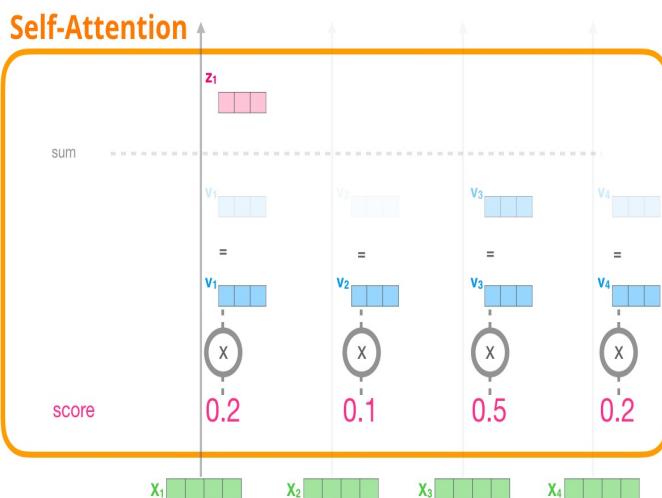
1) For each input token, create a **query vector**, a **key vector**, and a **value vector** by multiplying by weight Matrices  $W_Q, W_K, W_V$



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



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



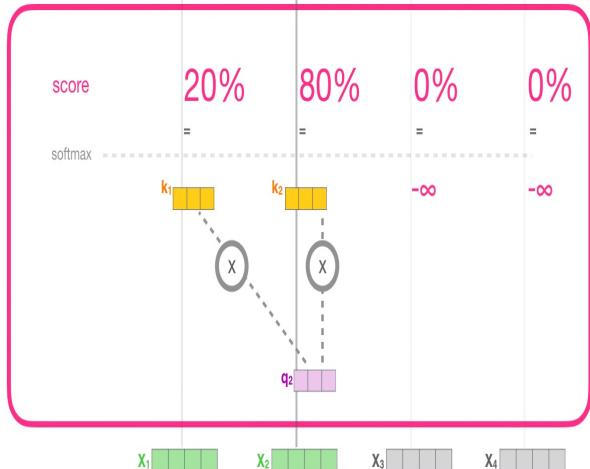
# Masked Self Attention

Masked Self-Attention

A



robot must obey orders.



Model has (2) tokens  
and we observe 2<sup>nd</sup> one.

B

Features

Labels

Example:	position: 1	2	3	4
1	robot	must	obey	orders
2	robot	must	obey	orders
3	robot	must	obey	orders
4	robot	must	obey	orders

must
obey
orders
<eos>

\* This sequence will be absorbed in (4) steps. BatchSize = 4 (for this toy example)

Keys

Scores  
(before softmax)

Queries	robot	must	obey	orders
robot	0.11	0.00	0.81	0.79
must	0.19	0.50	0.30	0.48
obey	0.53	0.98	0.95	0.14
orders	0.81	0.86	0.38	0.90

Not Matrix multiplication

C

robot must obey orders.

$q_1 K_1$	$q_1 K_2$	$q_1 K_3$	$q_1 K_4$
$q_2 K_1$	$q_2 K_2$	$q_2 K_3$	$q_2 K_4$
$q_3 K_1$	$q_3 K_2$	$q_3 K_3$	$q_3 K_4$
$q_4 K_1$	$q_4 K_2$	$q_4 K_3$	$q_4 K_4$

**Scores**  
(before softmax)

0.11	0.00	0.81	0.79
0.19	0.50	0.30	0.48
0.53	0.98	0.95	0.14
0.81	0.86	0.38	0.90

D

**Masked Scores**  
(before softmax)

0.11	-inf	-inf	-inf
0.19	0.50	-inf	-inf
0.53	0.98	0.95	-inf
0.81	0.86	0.38	0.90

Apply Attention Mask



**Masked Scores**  
(before softmax)

0.11	-inf	-inf	-inf
0.19	0.50	-inf	-inf
0.53	0.98	0.95	-inf
0.81	0.86	0.38	0.90

E

**Scores**

1	0	0	0
0.48	0.52	0	0
0.31	0.35	0.34	0
0.25	0.26	0.23	0.26

Softmax  
(along rows)  
(Why?)

When model processes the word "must" (input is "robot must") while processing first word ("robot") full attention was 48% attention is on "robot" & 52% attention is on "must"

