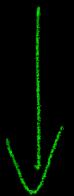


# Transformers & BERT



Bidirectional encoder  
representation from  
Transformers.

Assumption : Studied encoder-decoder models

Attention

LSTMs

CNNs

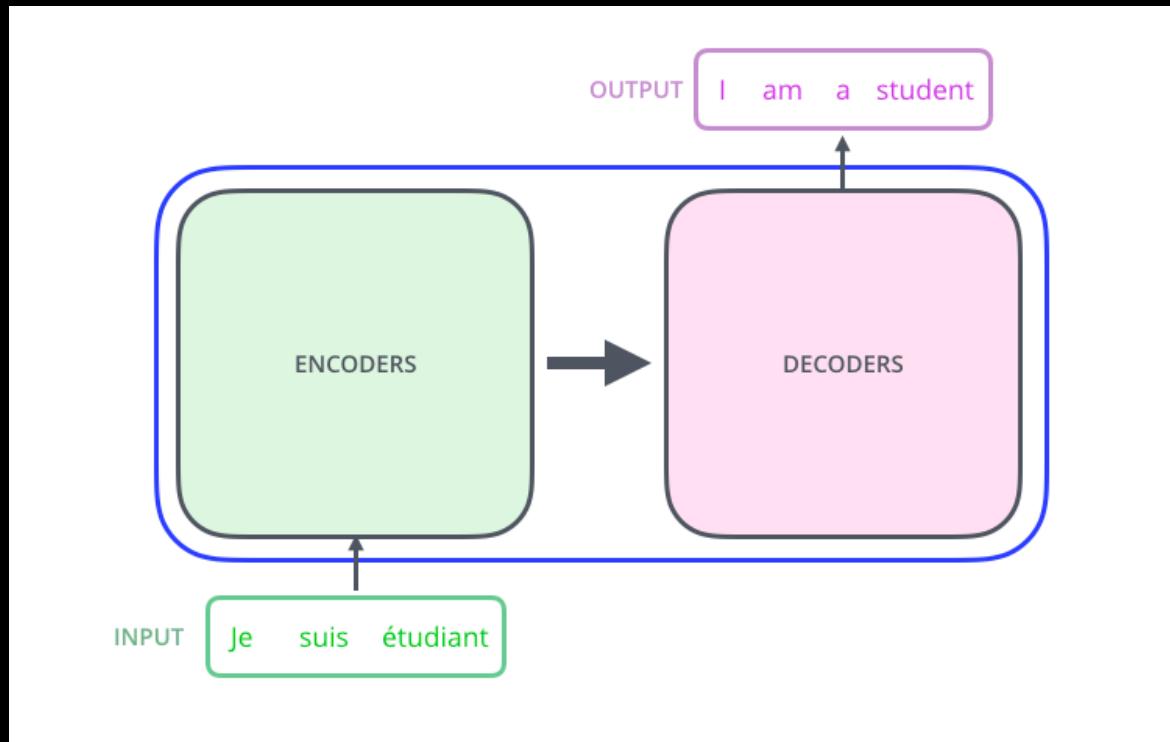
Source: <http://jalammar.github.io/illustrated-transformer/>

Transformers:

top-down approach

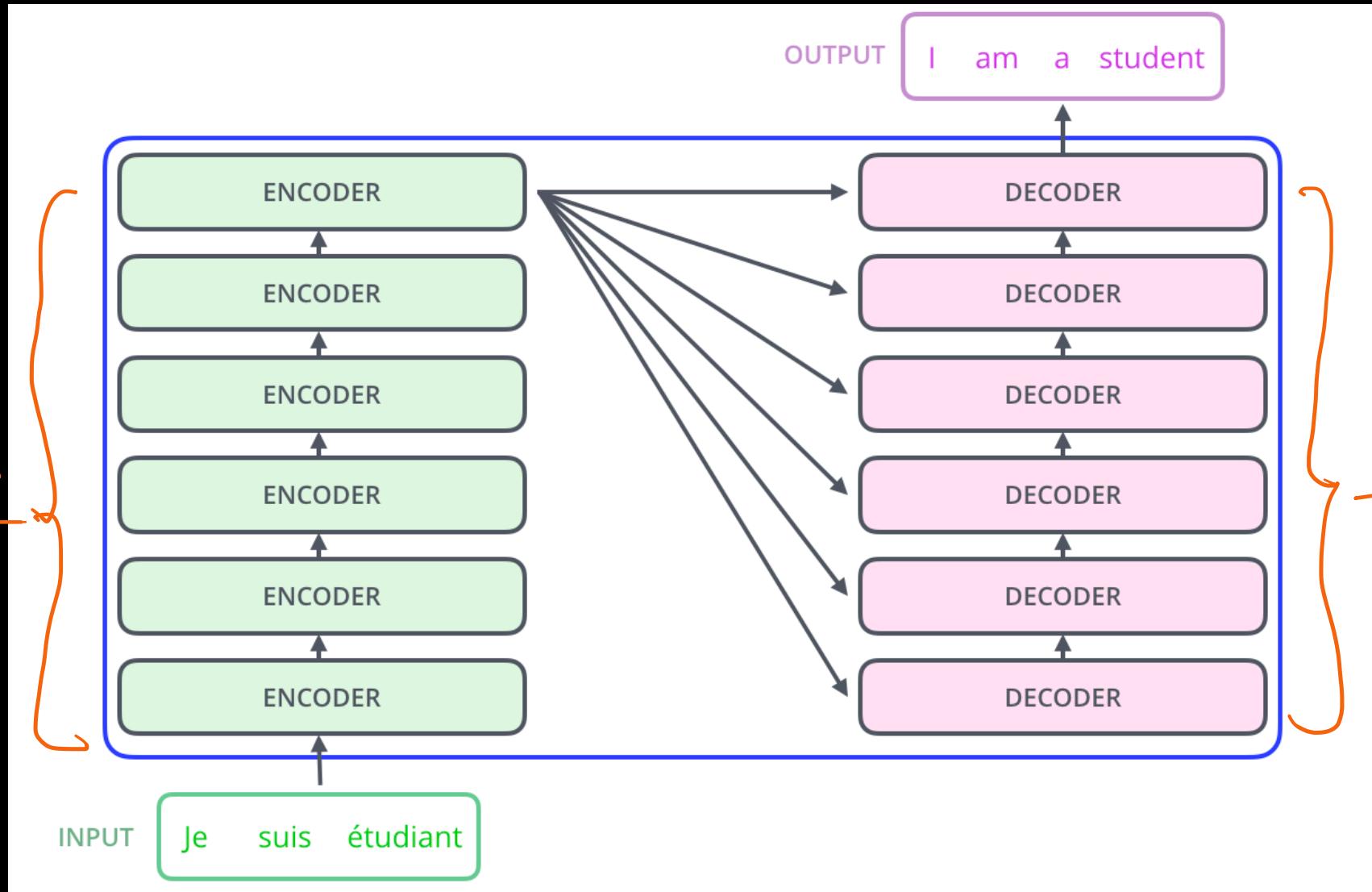


Seq-to-seq model

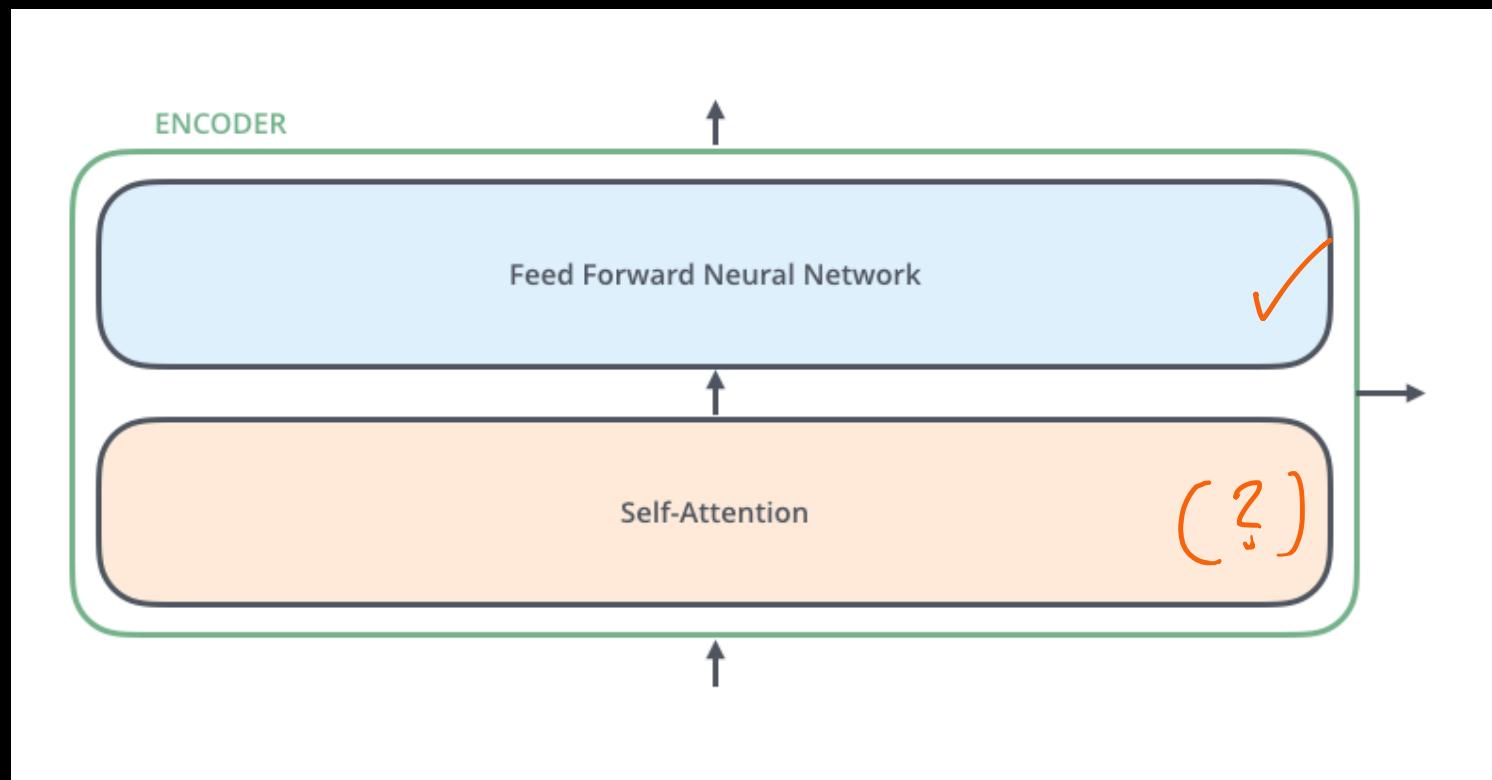


encoder-decoder

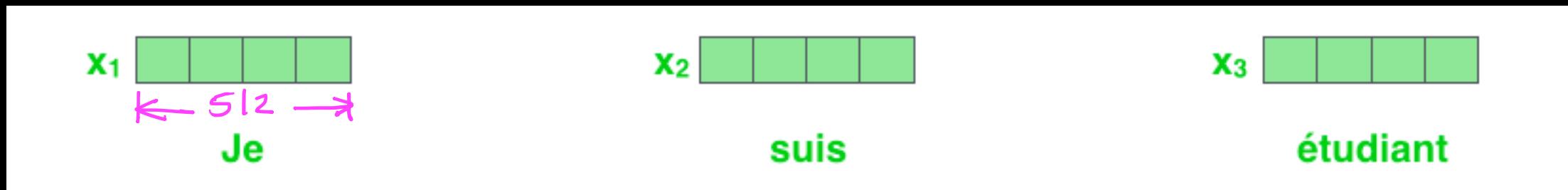
Stack of 6 encoders & decoders



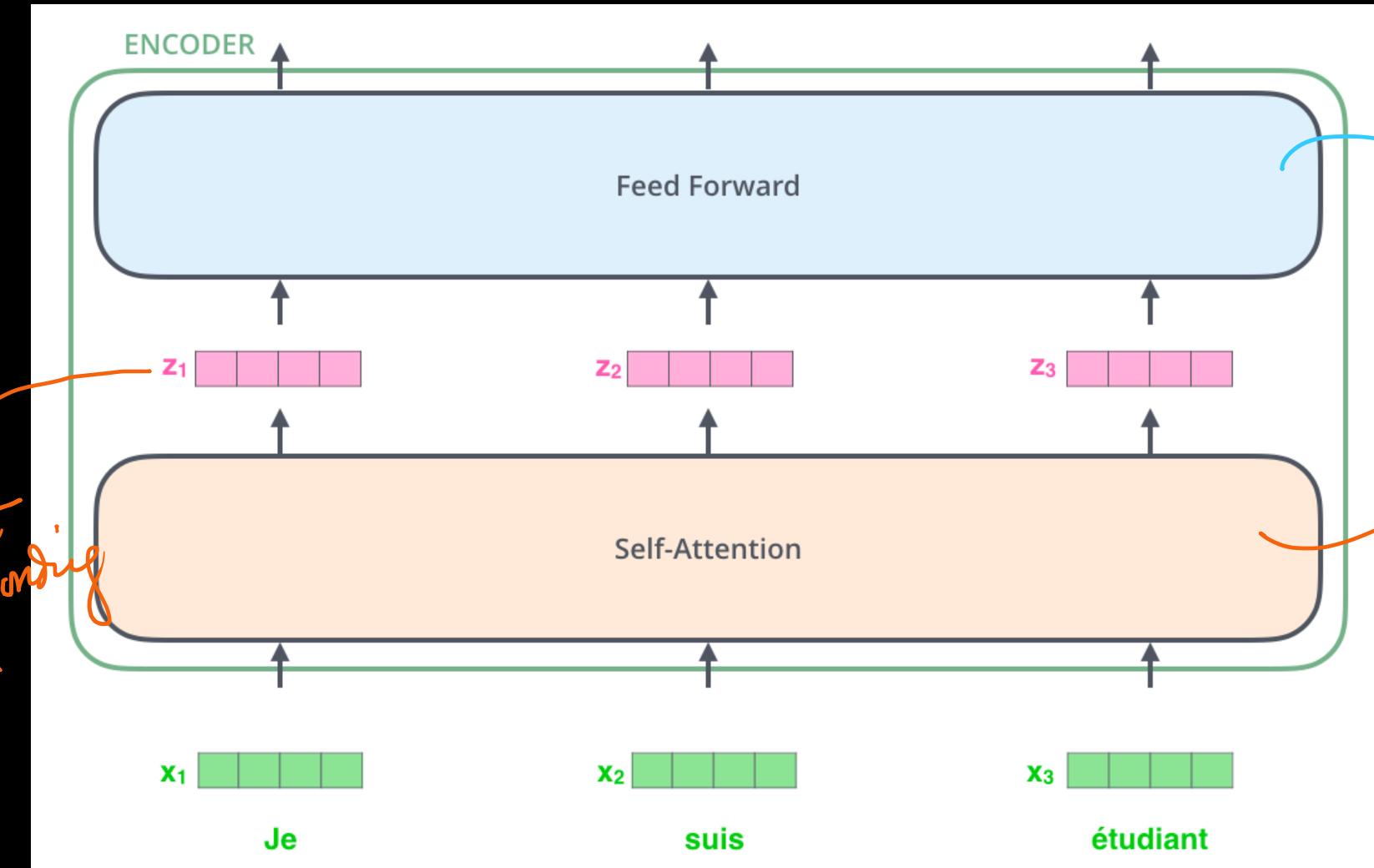
Structure of an encoder:



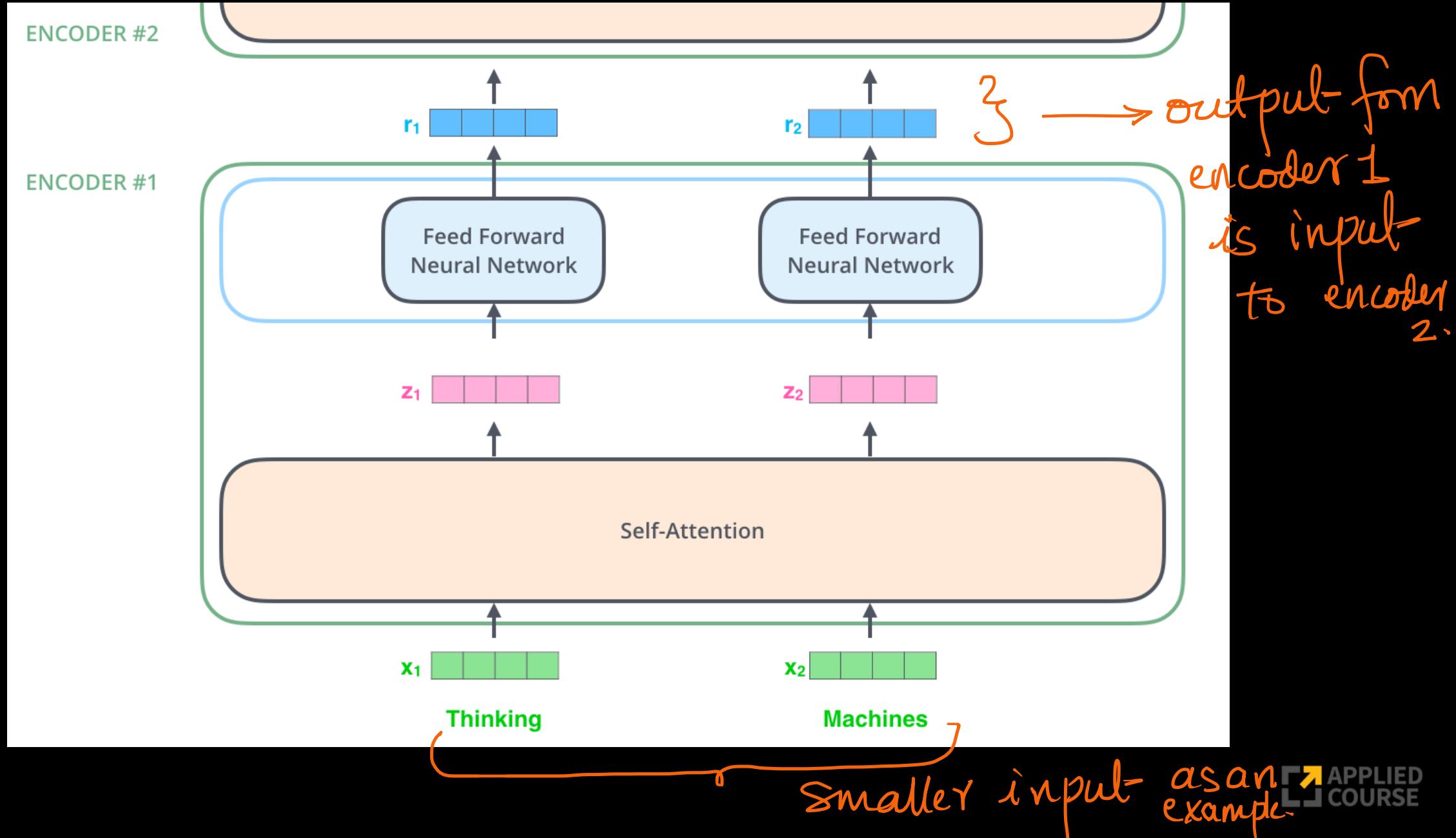
input: word-vectors/tensors



word2vec



we will see  
later how  
this works



## Self-attention:

Input: "The animal didn't cross the street because it was too tired"

↳ animal or street?

1

Input	Thinking		Machines	
Embedding	$x_1$	512	$x_2$	
Queries	$q_1$	64	$q_2$	$W^Q$
Keys	$k_1$	64	$k_2$	$W^K$
Values	$v_1$	64	$v_2$	$W^V$

$$q_i = x_i w^Q$$

$$k_i = x_i w^K$$

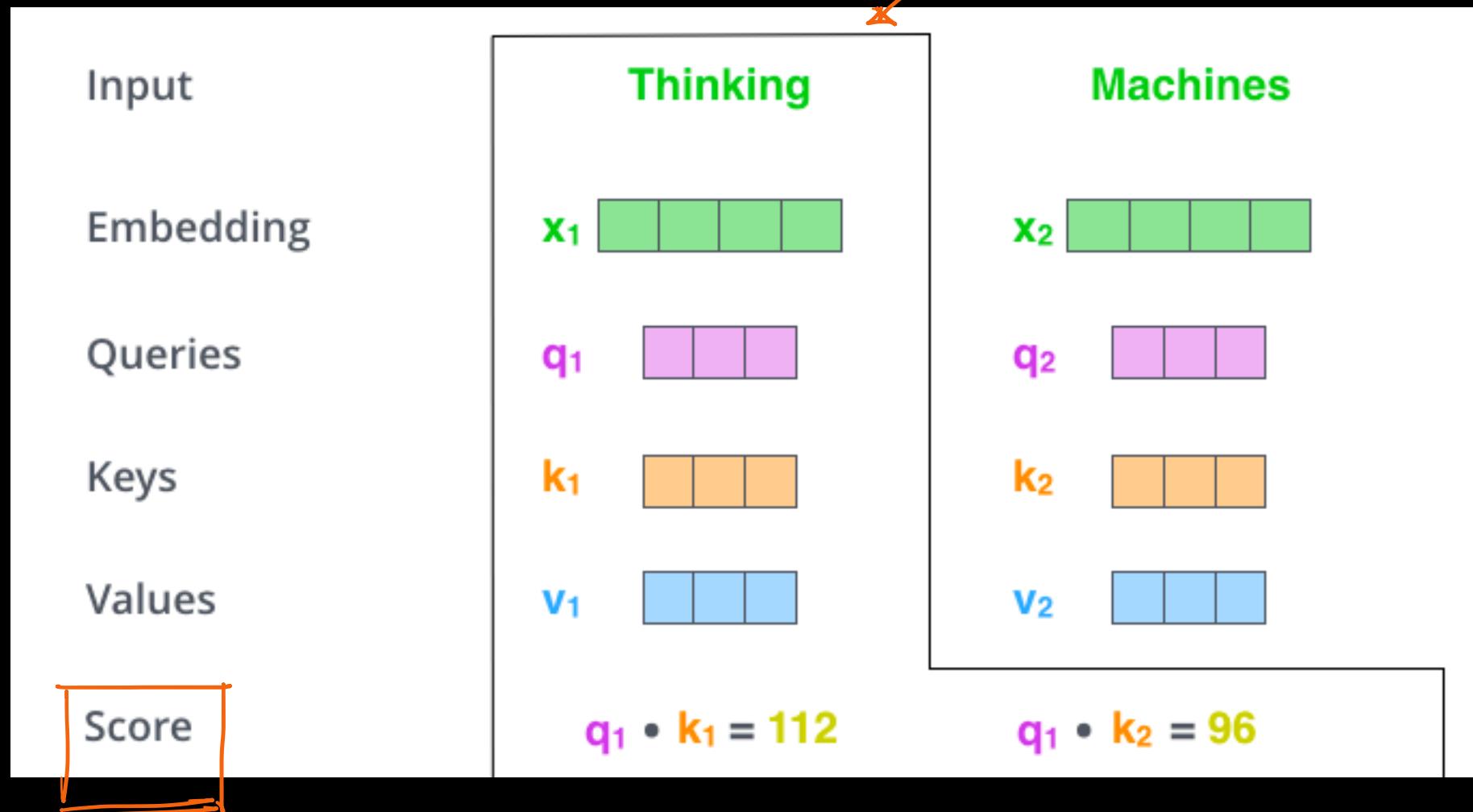
$$v_i = x_i w^V$$

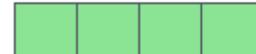
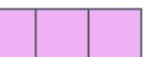
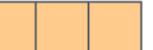
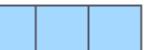
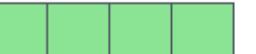
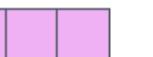
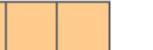
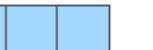


Trainable  
matrices -

②

computing self-attention for 1st word ( $x_1$ ) to generate ( $z_1$ )



Input	<b>Thinking</b>	
Embedding	$x_1$	
Queries	$q_1$	
Keys	$k_1$	
Values	$v_1$	
Score	$q_1 \cdot k_1 = 112$	
Divide by 8 ( $\sqrt{d_k}$ )	$= \frac{112}{\sqrt{64}} = 14$	
Softmax	0.88	
	<b>Machines</b>	
Embedding	$x_2$	
Queries	$q_2$	
Keys	$k_2$	
Values	$v_2$	
Score	$q_1 \cdot k_2 = 96$	
Divide by 8 ( $\sqrt{d_k}$ )	$= \frac{96}{\sqrt{64}} = 12$	
Softmax	0.12	

More stable  
↑ gradients

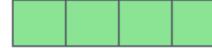
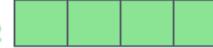
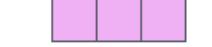
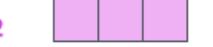
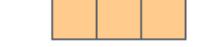
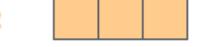
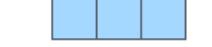
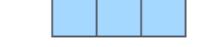
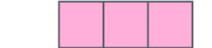
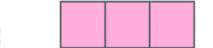
③

④

Sum=1

5

6

Input	<b>Thinking</b>		<b>Machines</b>	
Embedding	$x_1$		$x_2$	
Queries	$q_1$		$q_2$	
Keys	$k_1$		$k_2$	
Values	$v_1$		$v_2$	
Score	$q_1 \cdot k_1 = 112$		$q_1 \cdot k_2 = 96$	
Divide by 8 ( $\sqrt{d_k}$ )	14		12	
Softmax	0.88		0.12	
Softmax X Value	$v_1$		$v_2$	
Sum	$z_1$		$z_2$	

# Matrix-computation :

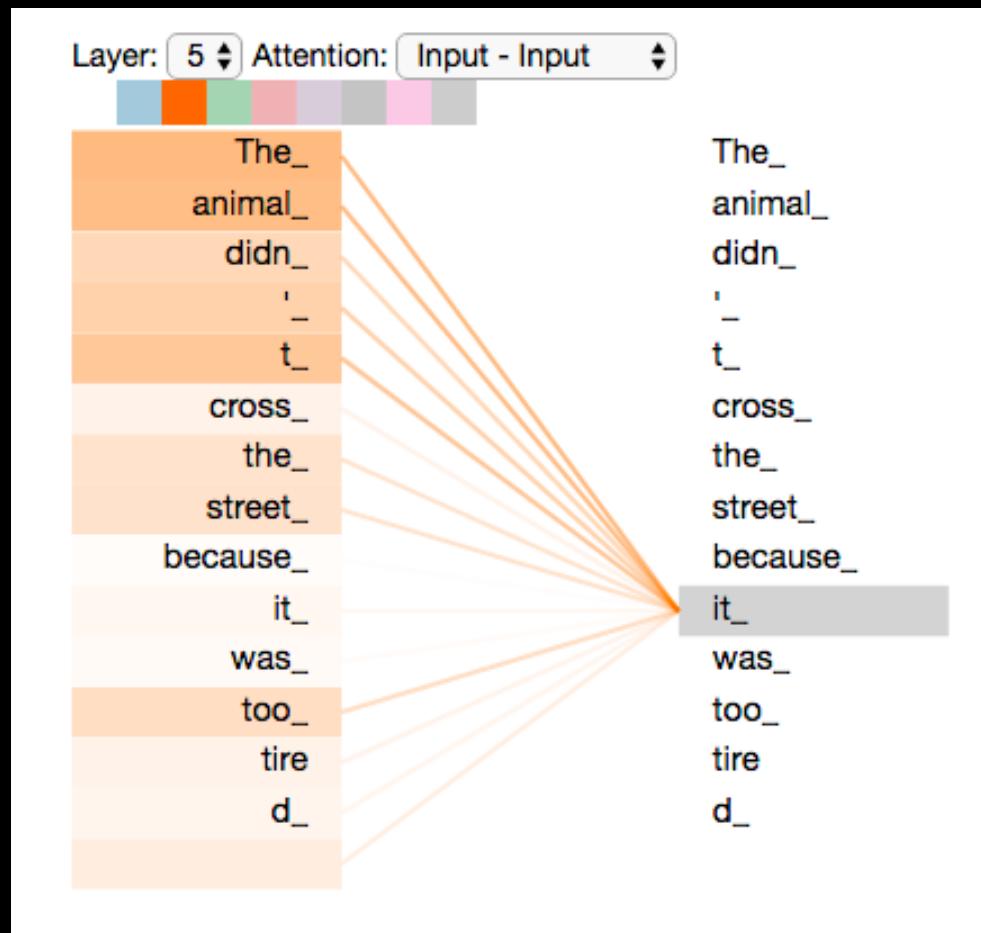
1

$$\begin{array}{ccc} \mathbf{X} & & \\ \text{X}_1: \begin{matrix} \text{green} & \text{green} & \text{green} \\ \text{green} & \text{green} & \text{green} \end{matrix} & \times & \begin{matrix} \text{purple} \\ \text{purple} \\ \text{purple} \\ \text{purple} \end{matrix} \\ & & = \begin{matrix} \text{purple} \\ \text{purple} \\ \text{purple} \\ \text{purple} \end{matrix} \end{array}$$
$$\begin{array}{ccc} \mathbf{X} & & \\ \text{X}_2: \begin{matrix} \text{green} & \text{green} & \text{green} \\ \text{green} & \text{green} & \text{green} \end{matrix} & \times & \begin{matrix} \text{orange} \\ \text{orange} \\ \text{orange} \\ \text{orange} \end{matrix} \\ & & = \begin{matrix} \text{orange} \\ \text{orange} \\ \text{orange} \\ \text{orange} \end{matrix} \end{array}$$
$$\begin{array}{ccc} \mathbf{X} & & \\ \text{X}_3: \begin{matrix} \text{green} & \text{green} & \text{green} \\ \text{green} & \text{green} & \text{green} \end{matrix} & \times & \begin{matrix} \text{blue} \\ \text{blue} \\ \text{blue} \\ \text{blue} \end{matrix} \\ & & = \begin{matrix} \text{blue} \\ \text{blue} \\ \text{blue} \\ \text{blue} \end{matrix} \end{array}$$

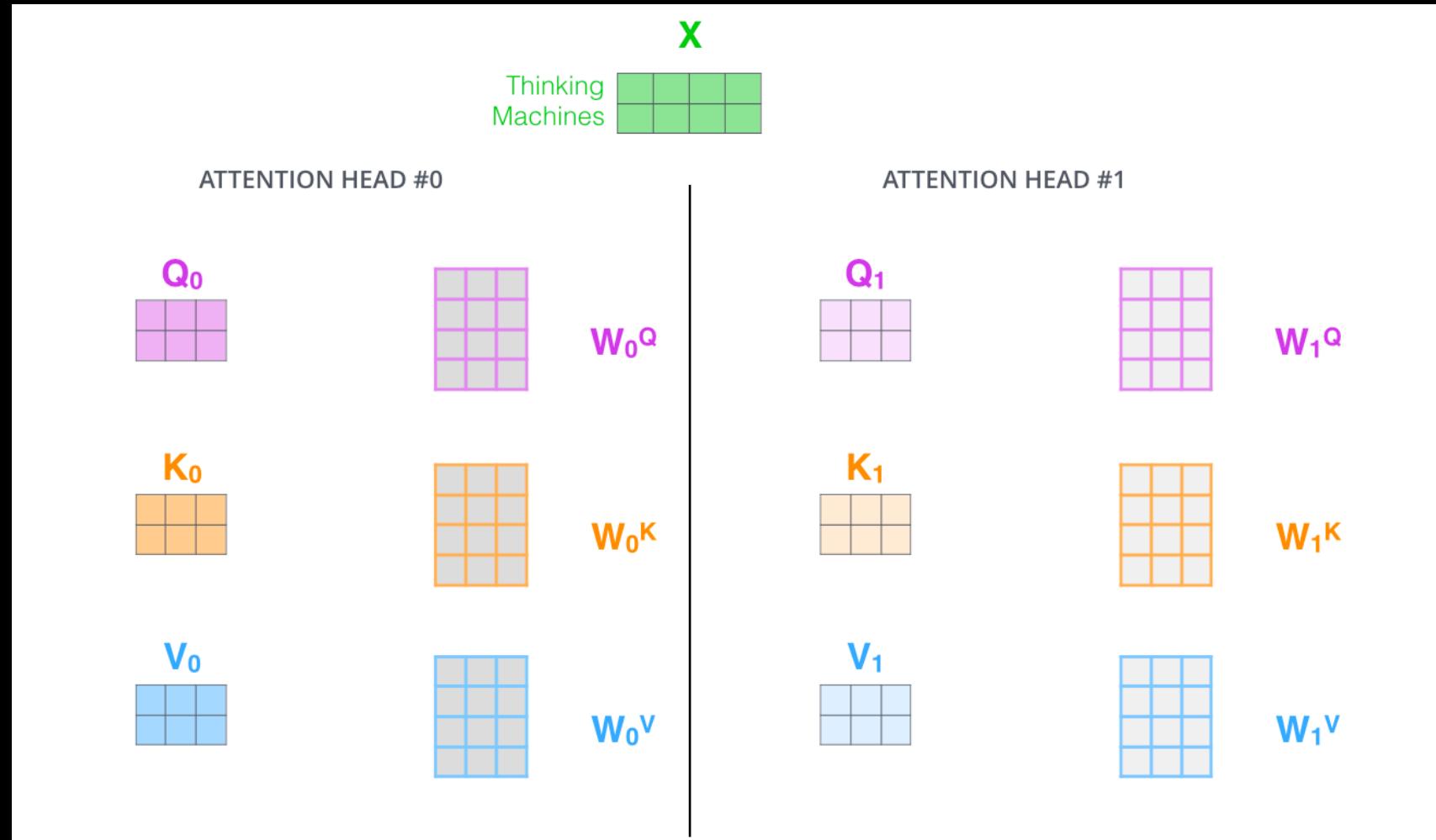
2

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

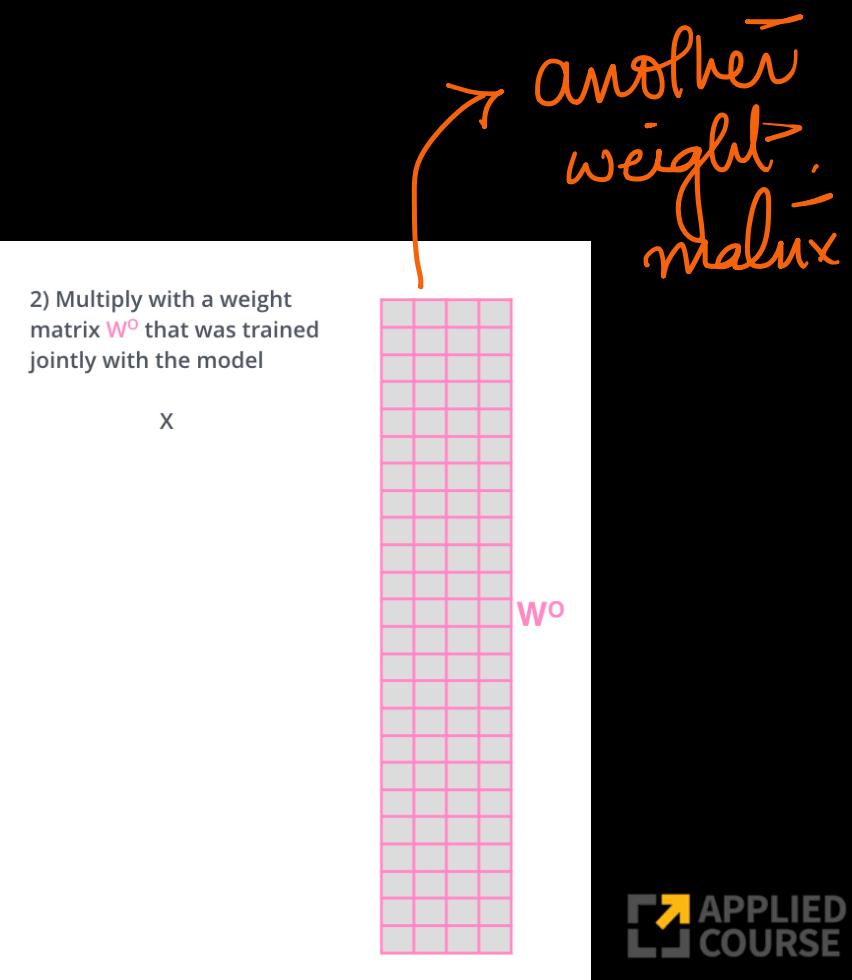
Attention helps us to focus on some of the words in the vicinity which matter



# Multi-headed Attention: 8-headed (increases representational power)

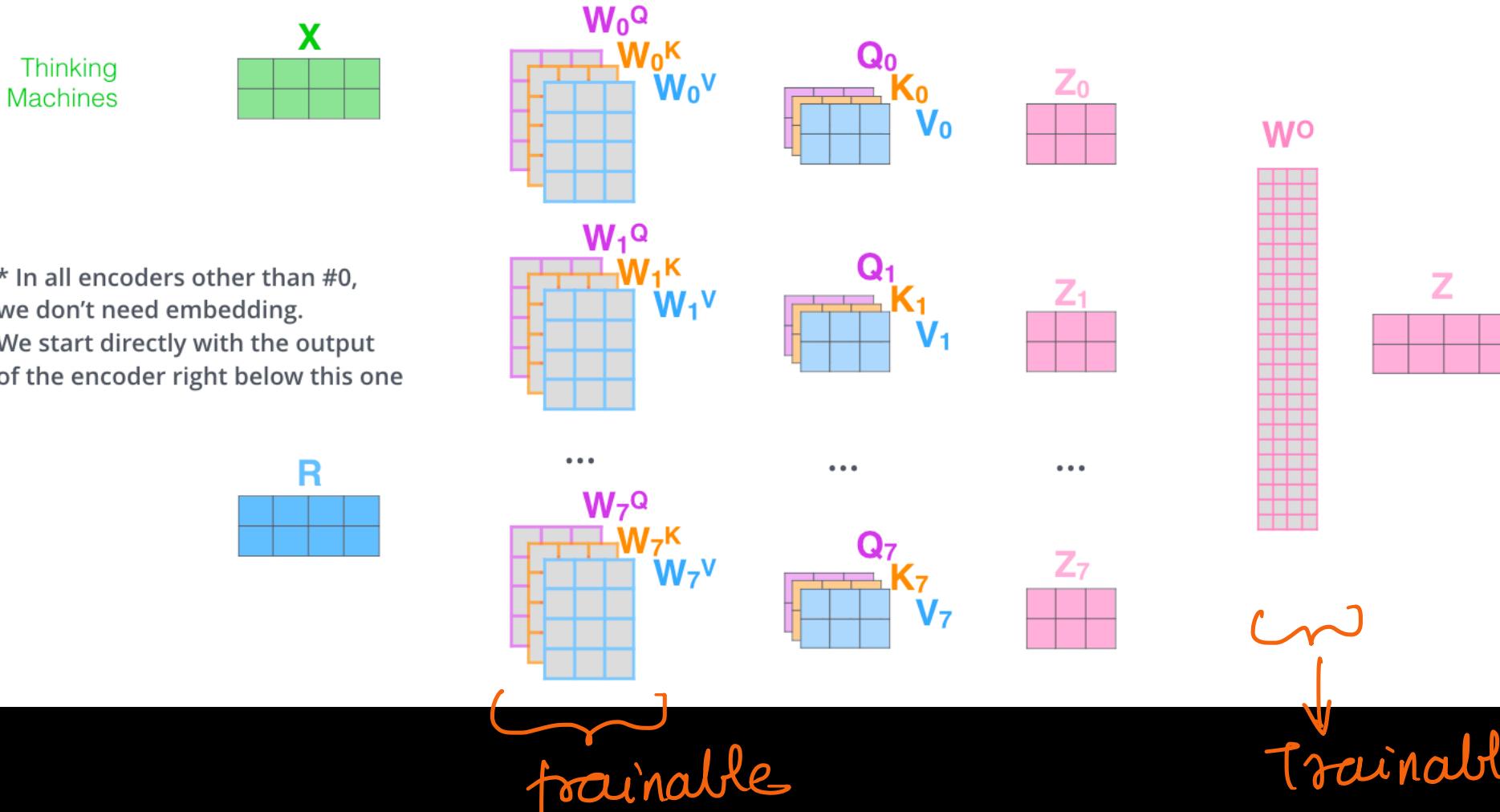


single-headed  
 $z_1$  could focus  
on 1st word  
itself more



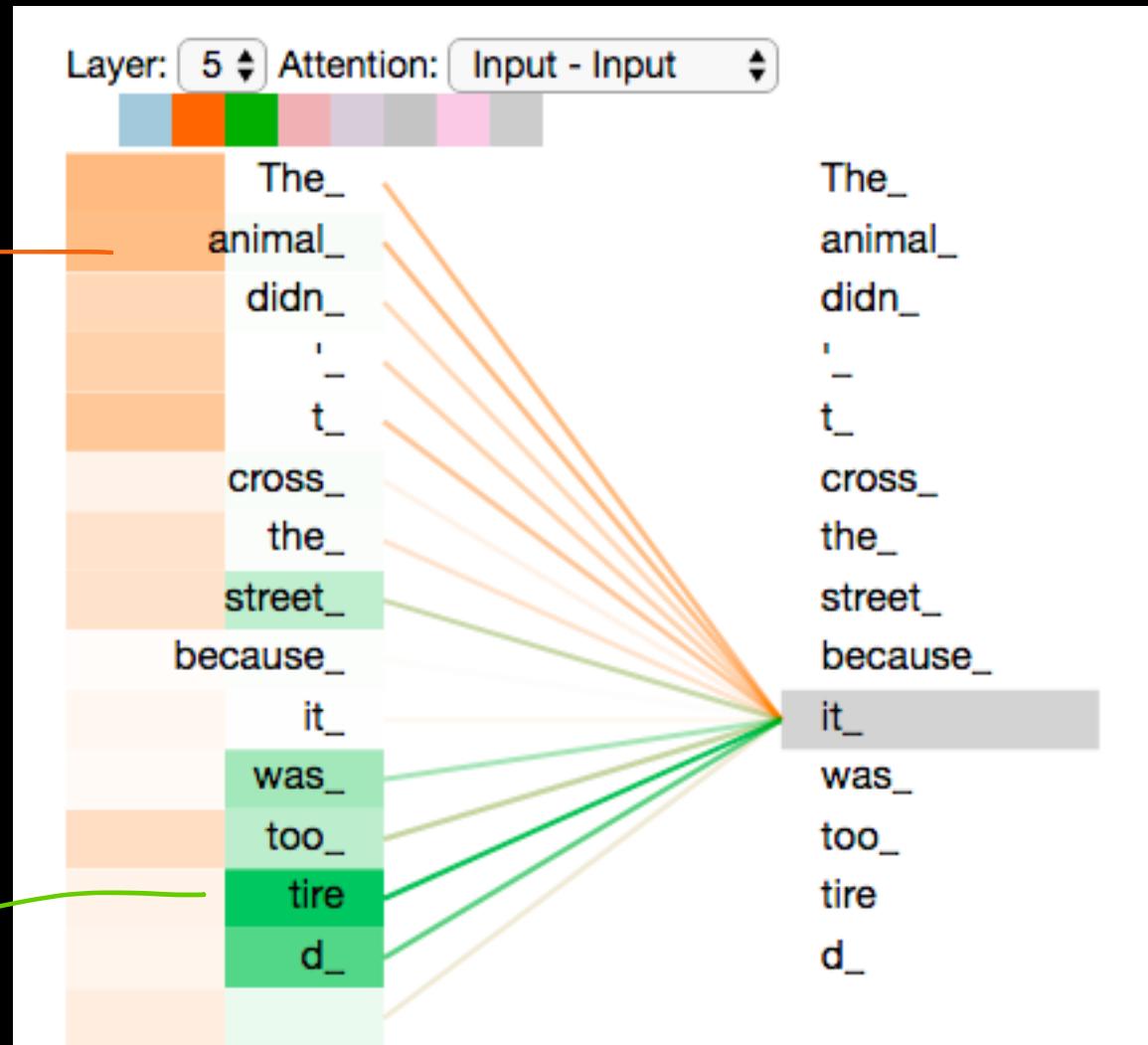
# Summary:

- 1) This is our input sentence\*
- 2) We embed each word\*
- 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

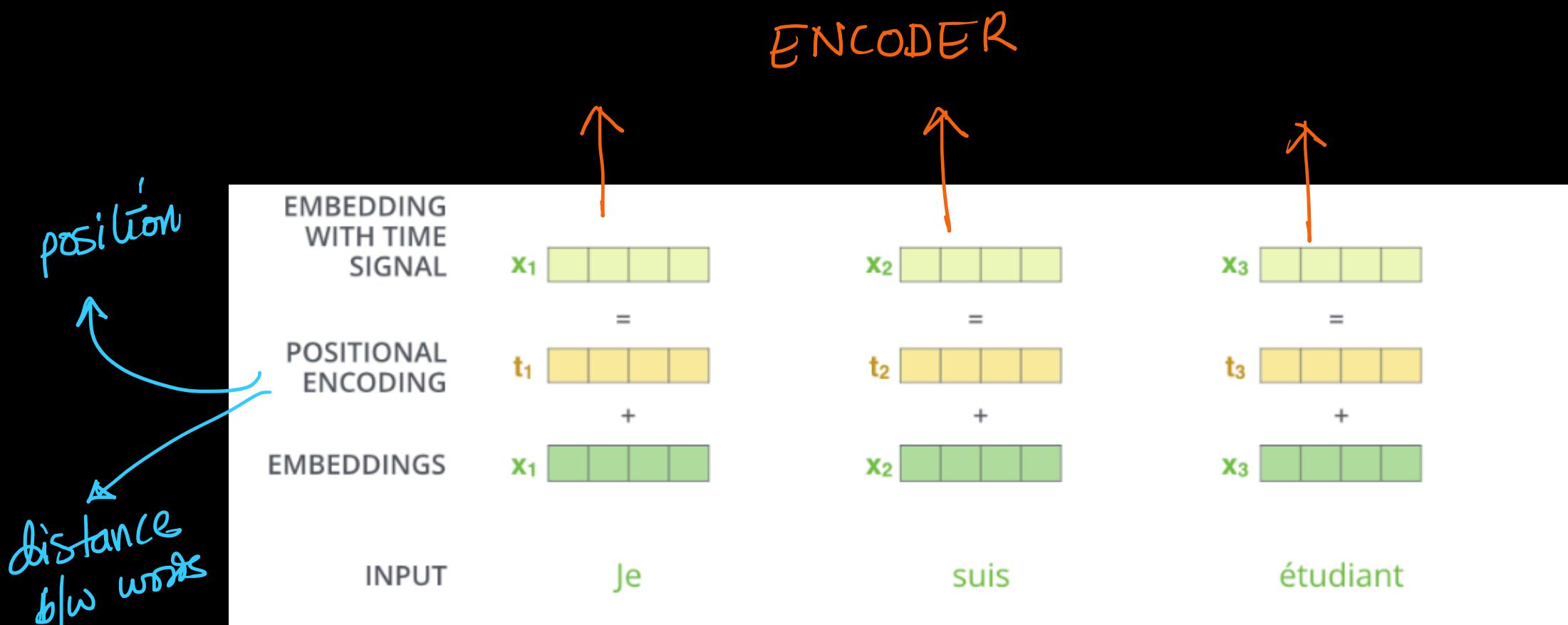


1st head  
focus on 'animal'

2nd head  
focus on 'fire'



Order of words:



## POSITIONAL ENCODING

0	0	1	1
---	---	---	---

+

## EMBEDDINGS

$x_1$	[green squares]
-------	-----------------

0.84	0.0001	0.54	1
------	--------	------	---

+

0.91	0.0002	-0.42	1
------	--------	-------	---

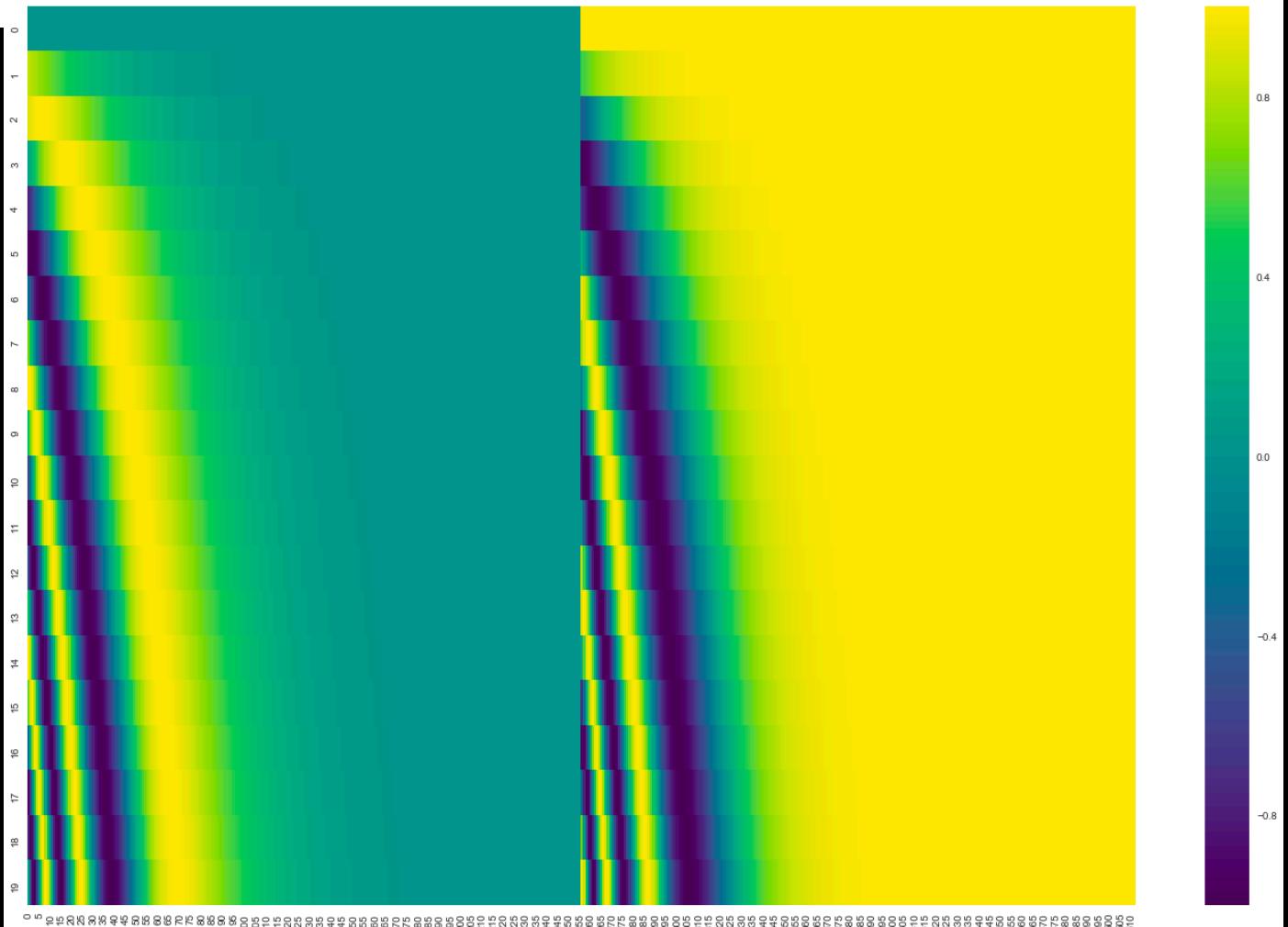
+

## INPUT

Je

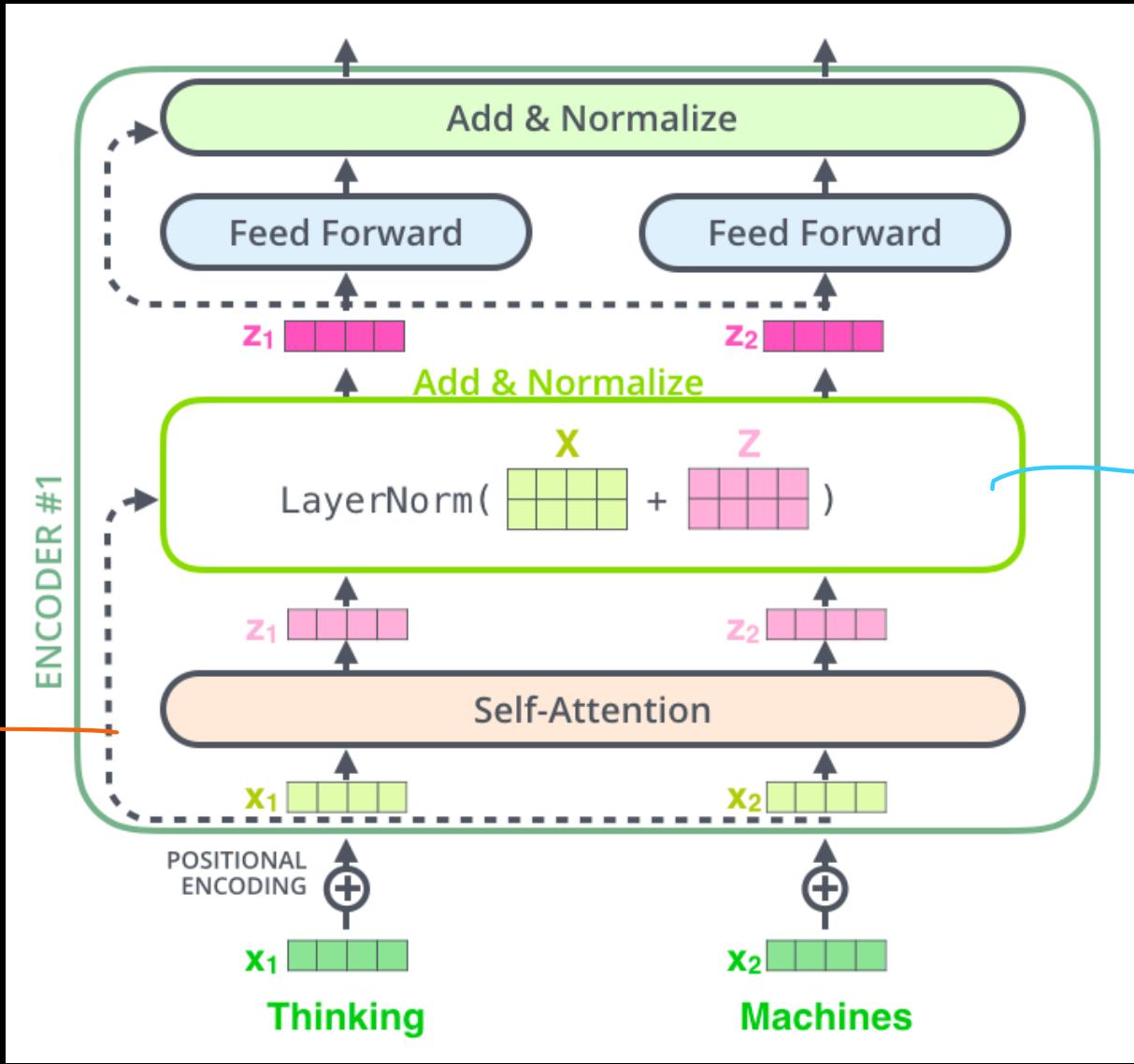
suis

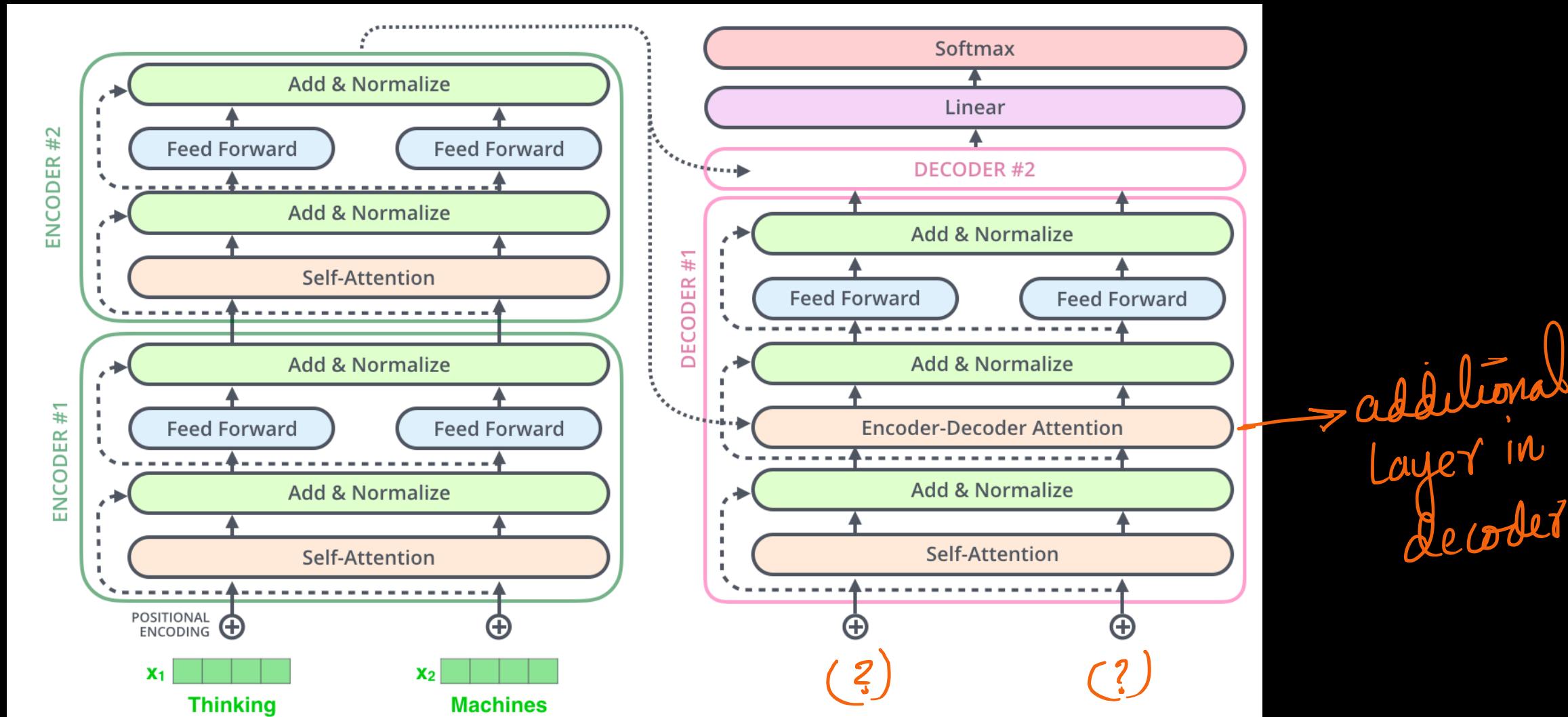
étudiant



Residual  
connection

(like in ResNets)





Decoding time-steps :

①

[http://jalammar.github.io/images/t/transformer\\_decoding\\_1.gif](http://jalammar.github.io/images/t/transformer_decoding_1.gif)

②

[http://jalammar.github.io/images/t/transformer\\_decoding\\_2.gif](http://jalammar.github.io/images/t/transformer_decoding_2.gif)

Which word in our vocabulary  
is associated with this index?

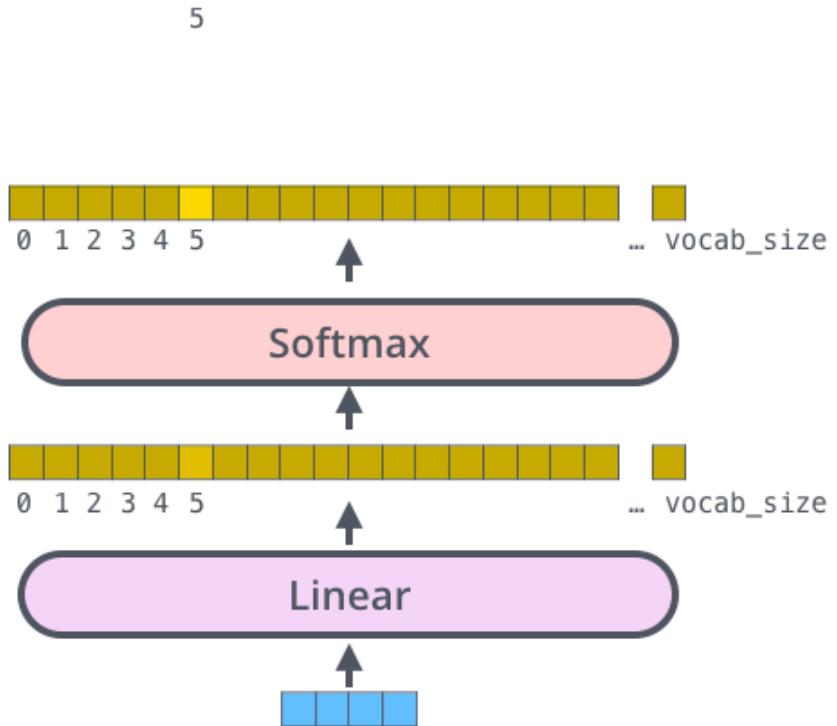
am

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

log\_probs

logits

Decoder stack output



$L$  : Cross-  
entropy

## Key-take aways:

- faster than RNN-based models as all input is ingested once
- inspired from neighborhood-like concept in CNNs
- state of the art.

BERT : bi-directional encoder representation from  
Transformers



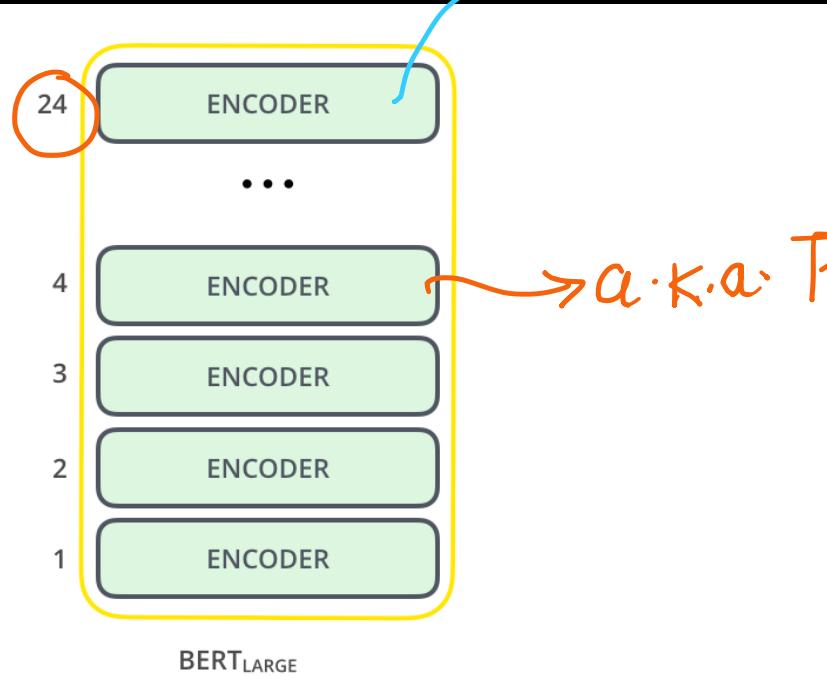
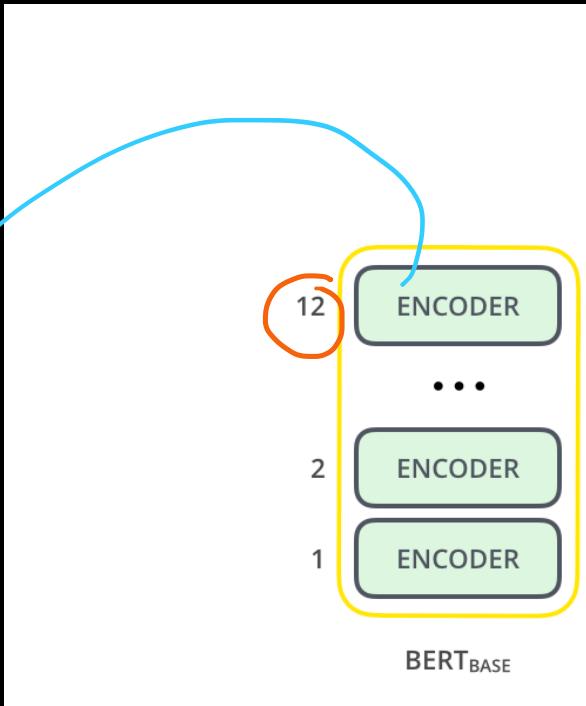
SQuAD1.1 Leaderboard			
Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar et al. '16)	82.304	91.221
1 Oct 05, 2018	BERT (ensemble) Google AI Language <a href="https://arxiv.org/abs/1810.04805">https://arxiv.org/abs/1810.04805</a>	87.433	93.160
2 Sep 09, 2018	nlnet (ensemble) Microsoft Research Asia	85.356	91.202
3 Jul 11, 2018	QANet (ensemble) Google Brain & CMU	84.454	90.490

Transfer Learning in NLP.

→ comparable to humans on multiple NLP tasks.

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

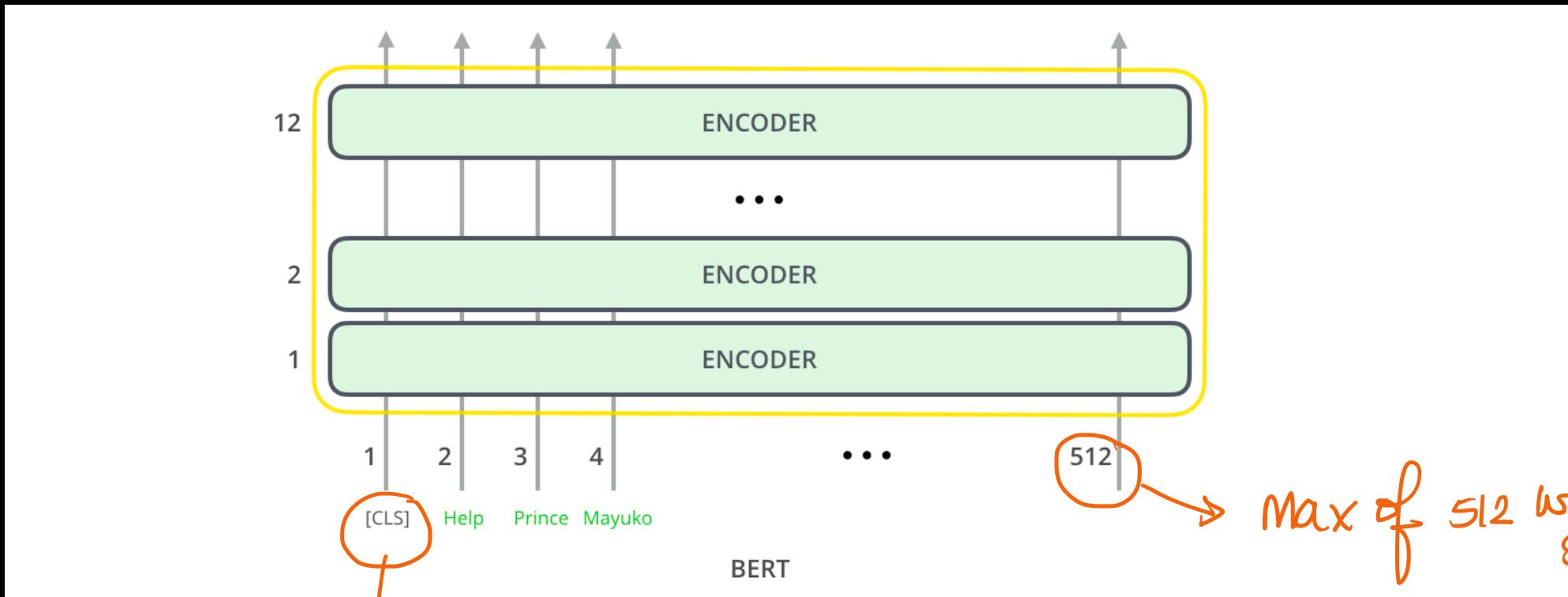
FFN: 768  
12 Attention heads



Encoder-stack

NO-DECODERS

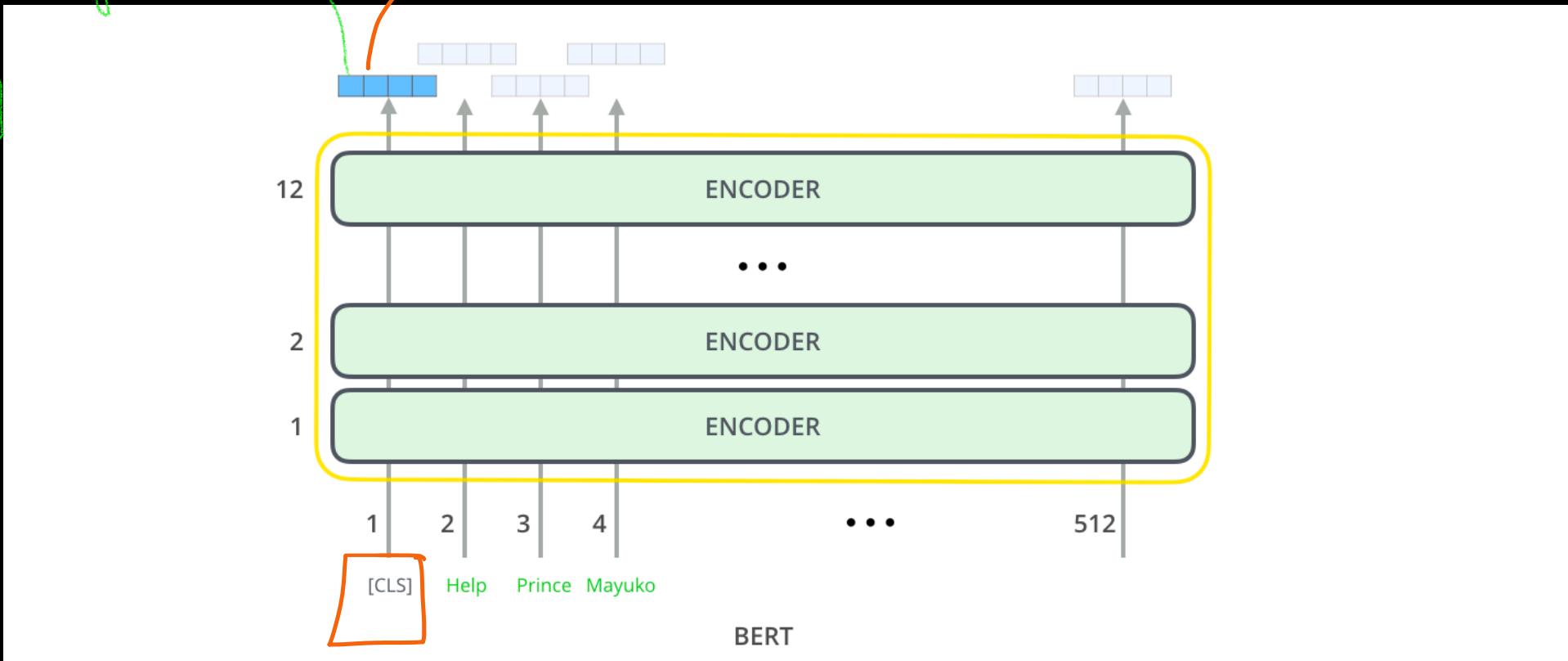
# Spam or NOT (Transfer Learning)

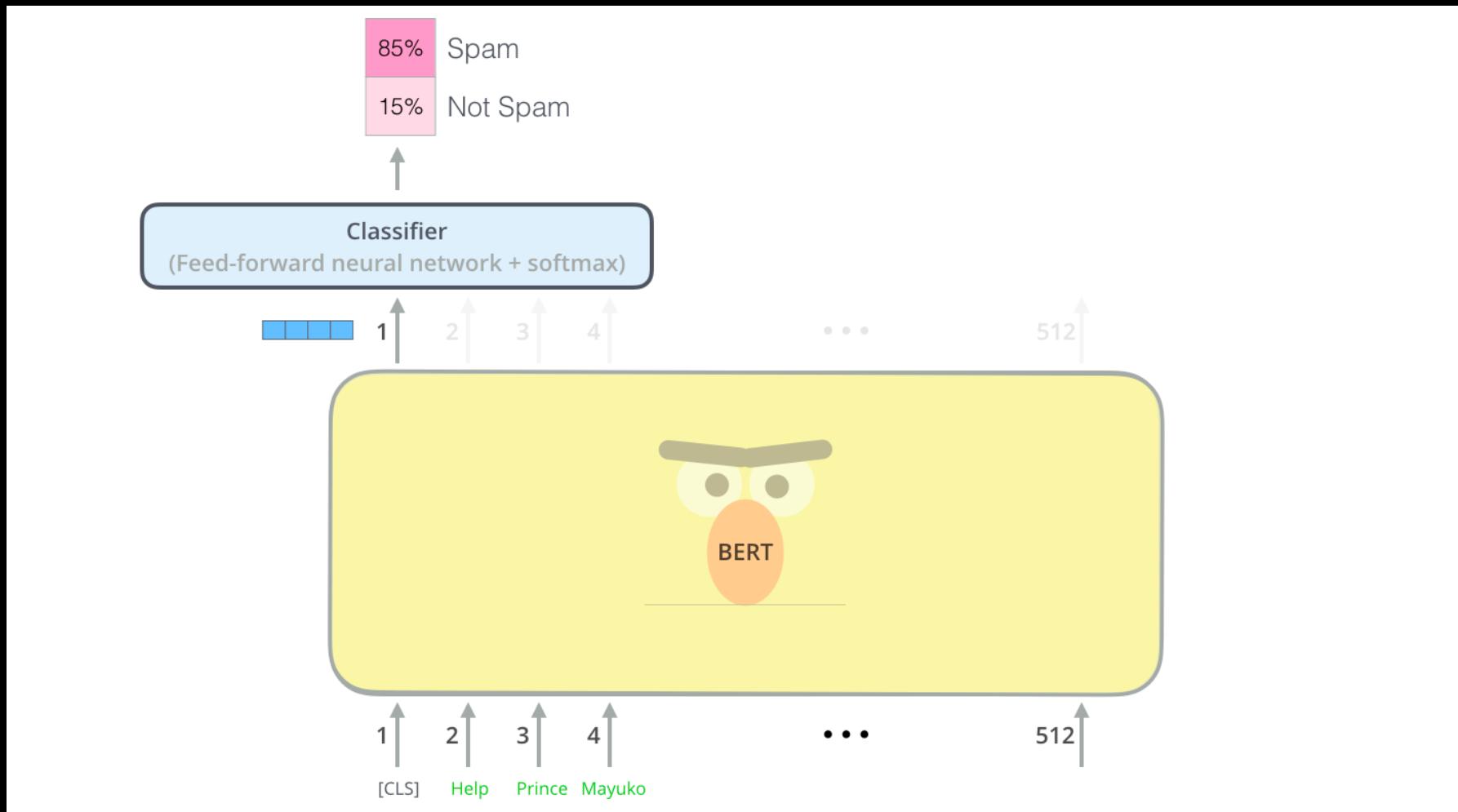


classification

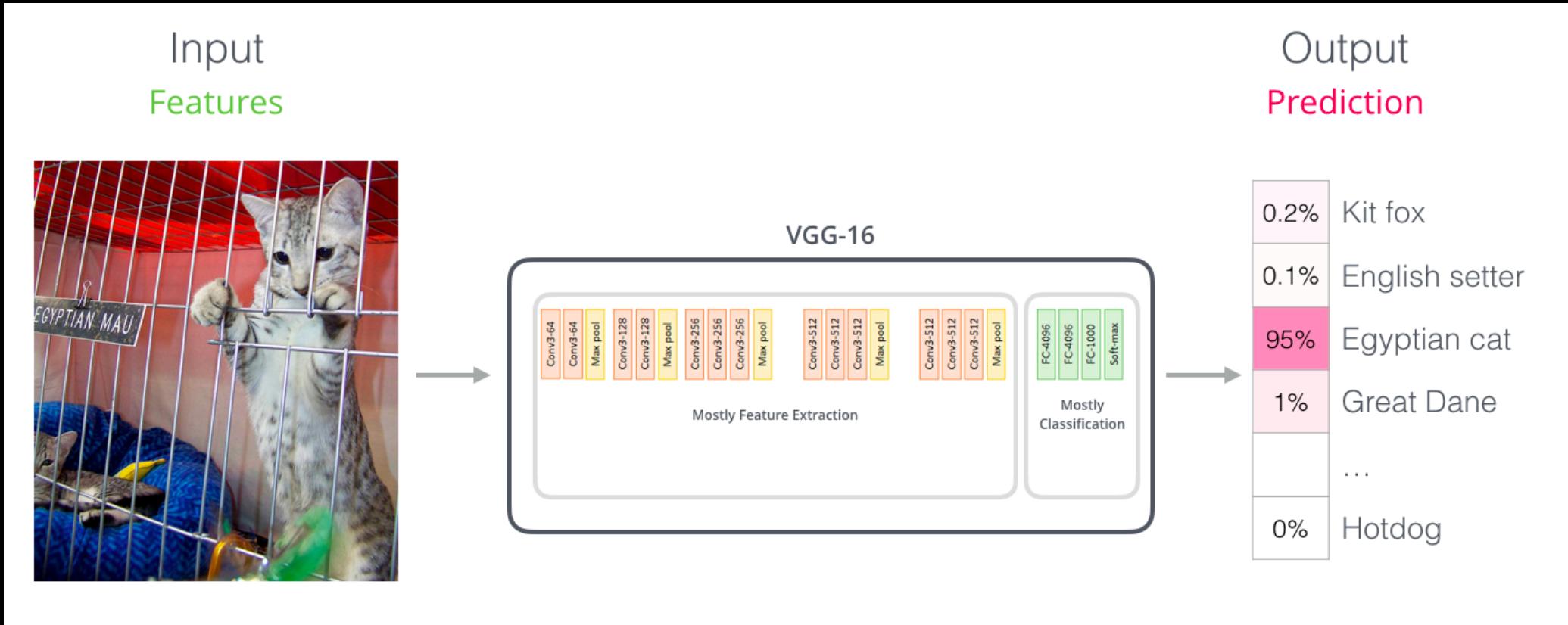
output  
corresponding  
to [CLS]

768-dim vector in BERT-base





Very similar to popular CNN-Models:



BERT training:

→ Bi-directional models are better than forward-only models.



looking @ words before & after matters

# Masked Language model

Similar in idea

to word2Vec



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

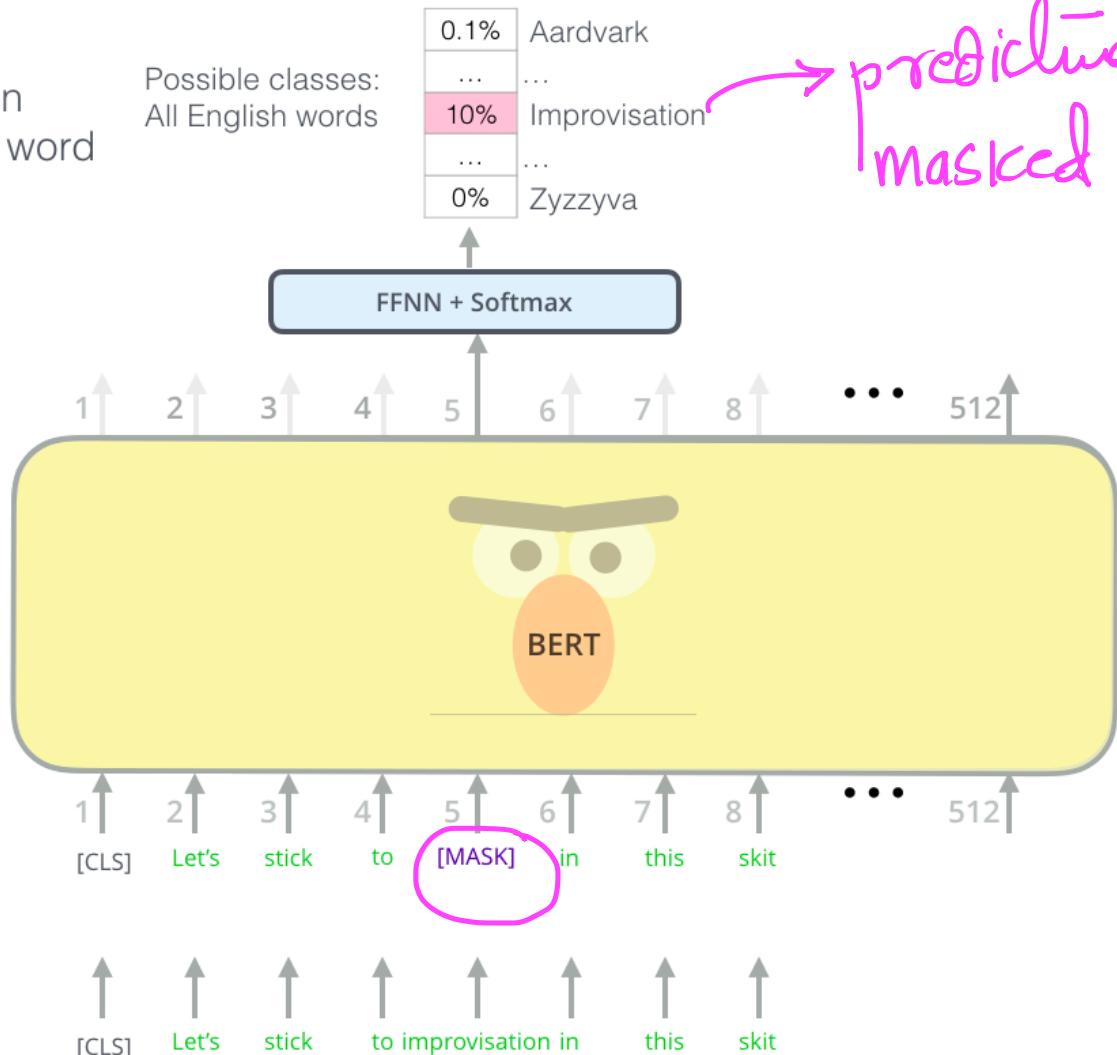
Randomly mask 15% of tokens

Input

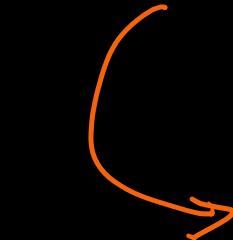
Possible classes:  
All English words



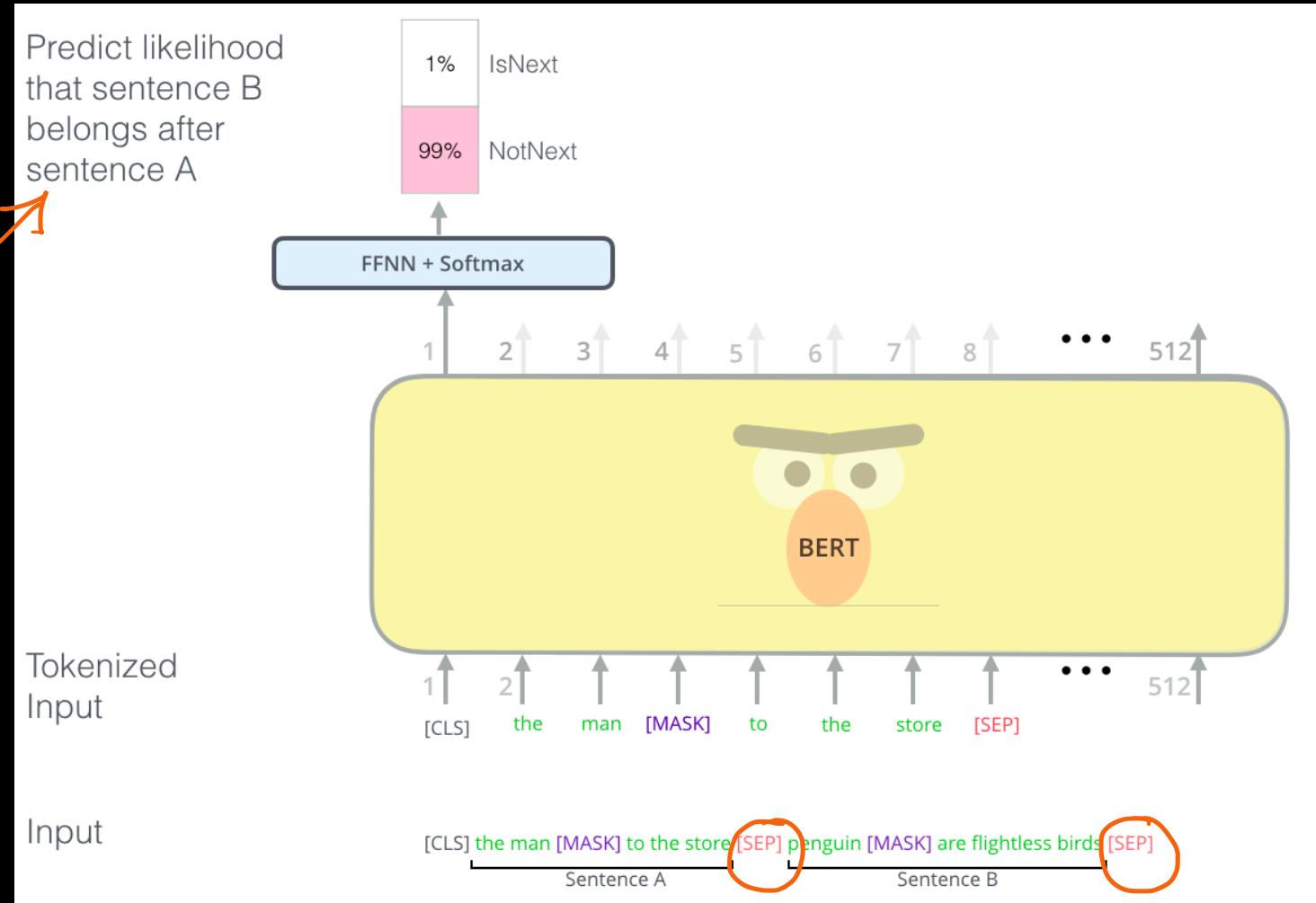
→ predicting the masked word.



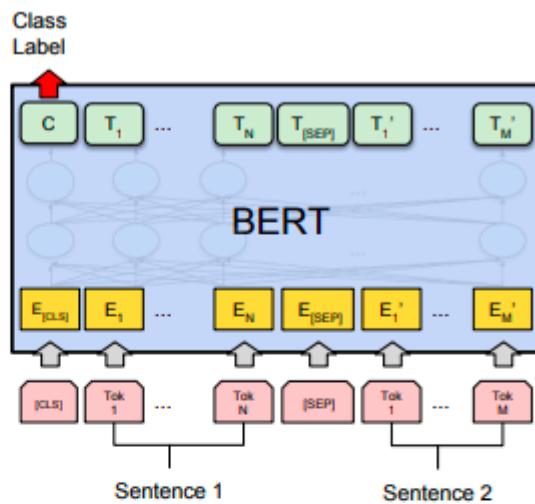
## 2-Sentence-tasks:



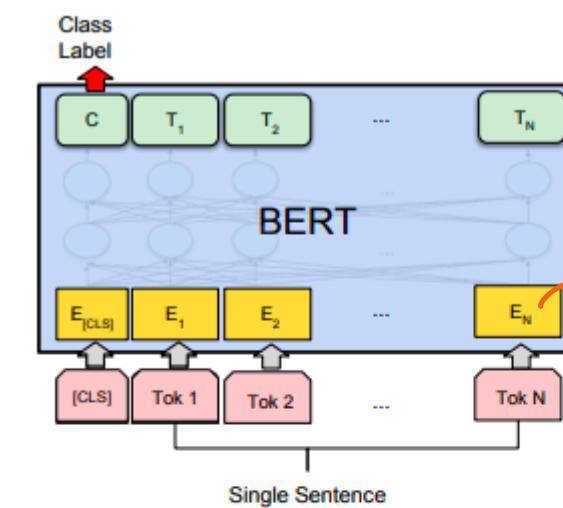
Task:



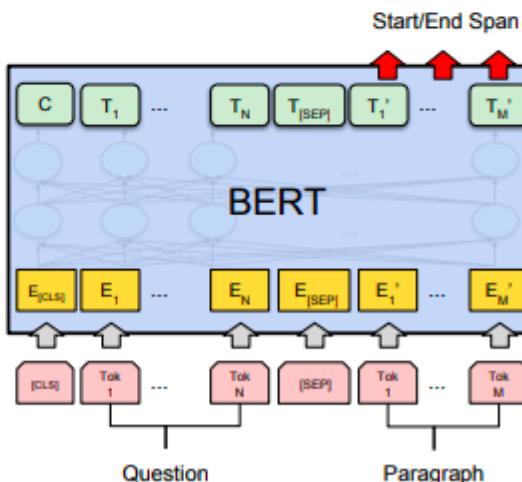
# Task-specific - Models:



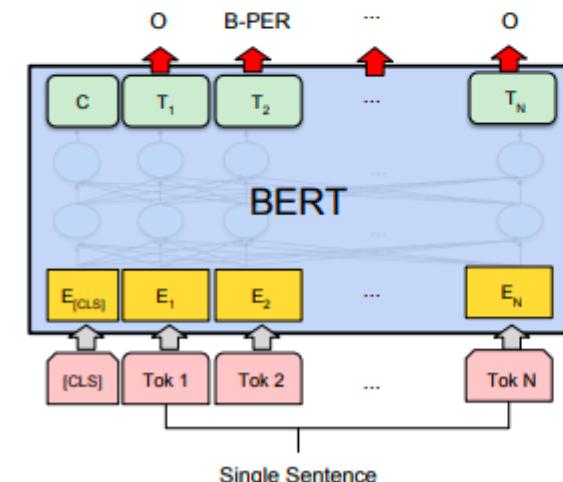
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



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



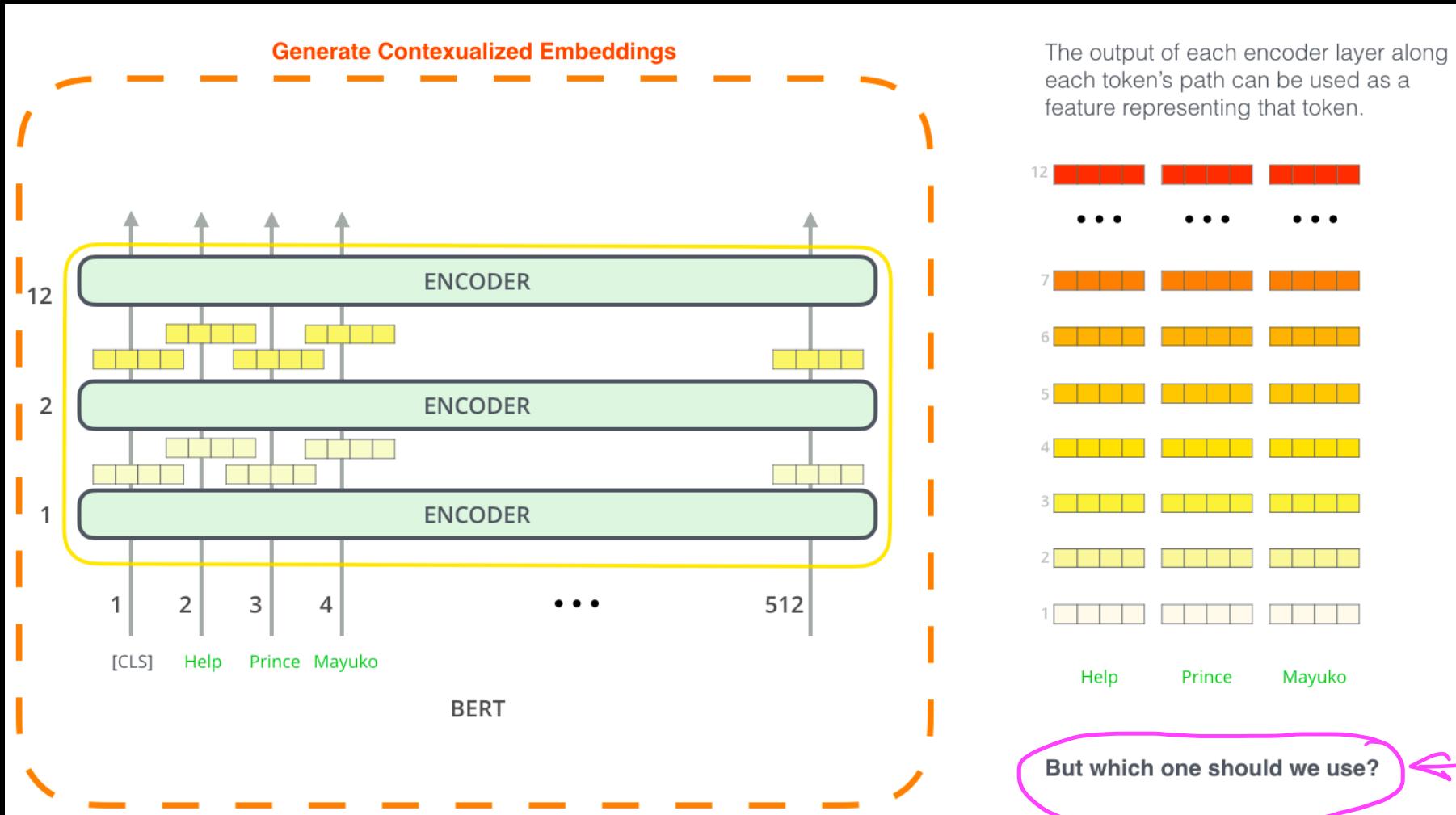
(c) Question Answering Tasks:  
SQuAD v1.1



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

word/token embeddings

# BERT - 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 

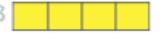
• • •

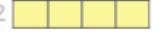
7 

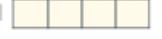
6 

5 

4 

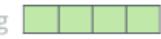
3 

2 

1 



First Layer

Embedding 

91.0

Last Hidden Layer

12 

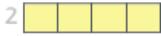
94.9

Sum All 12  
Layers

12 

+

...

2 

+

1 

=

95.5

Second-to-Last  
Hidden Layer

11 

95.6

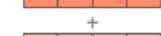
Sum Last Four  
Hidden

12 

+

11 

+

10 

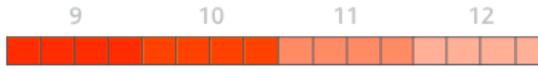
+

9 

=

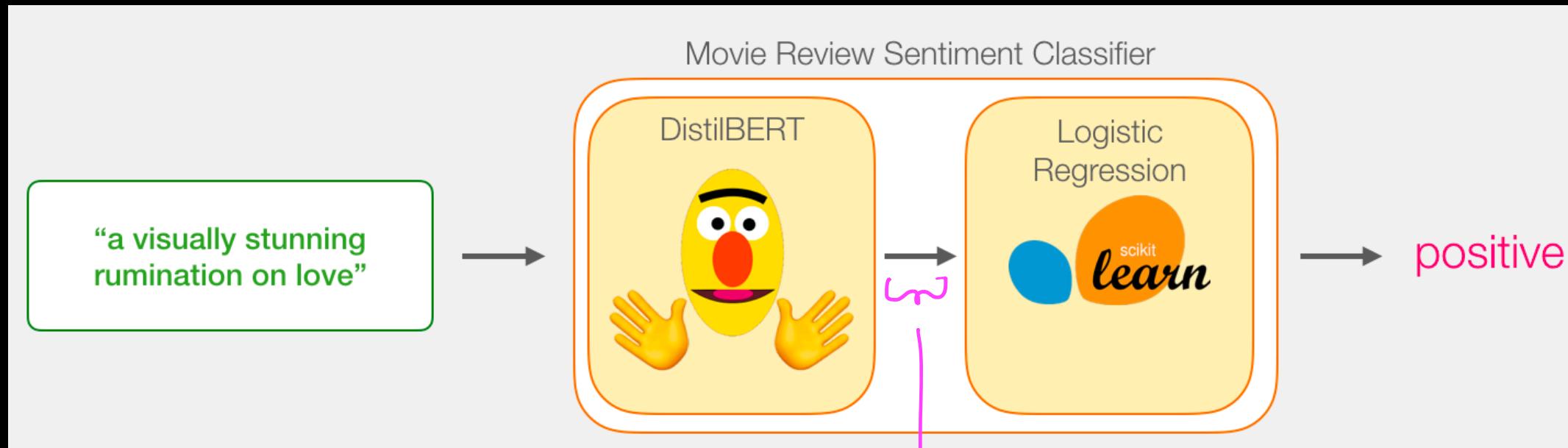
95.9

Concat Last  
Four Hidden

9    10    11    12 

96.1

Example to use BERT pre-trained model for sentence classification:  
<http://jalammar.github.io/a-visual-guide-to-using-bert-for-the-first-time/>



Sentence-vectors