



# Attention, Transformer and BERT

Prof. Kuan-Ting Lai

2021/4/17

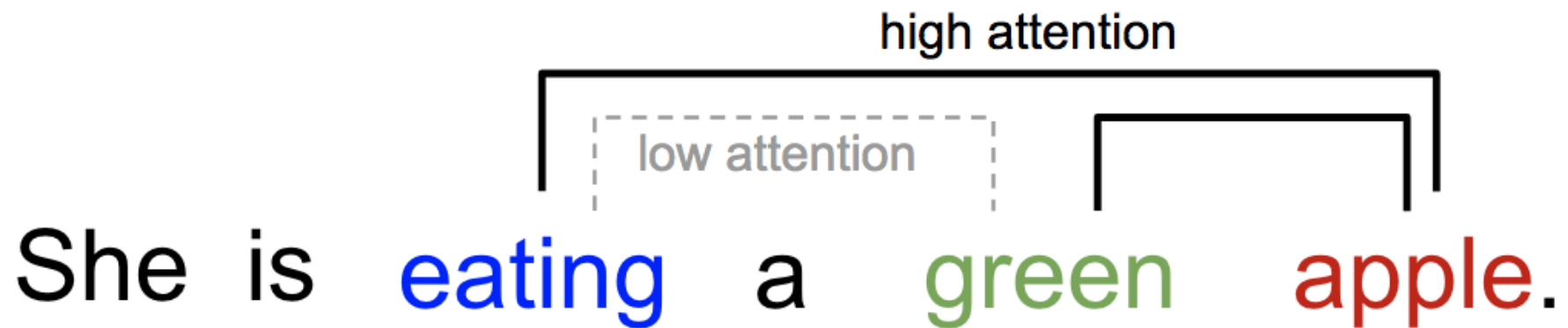
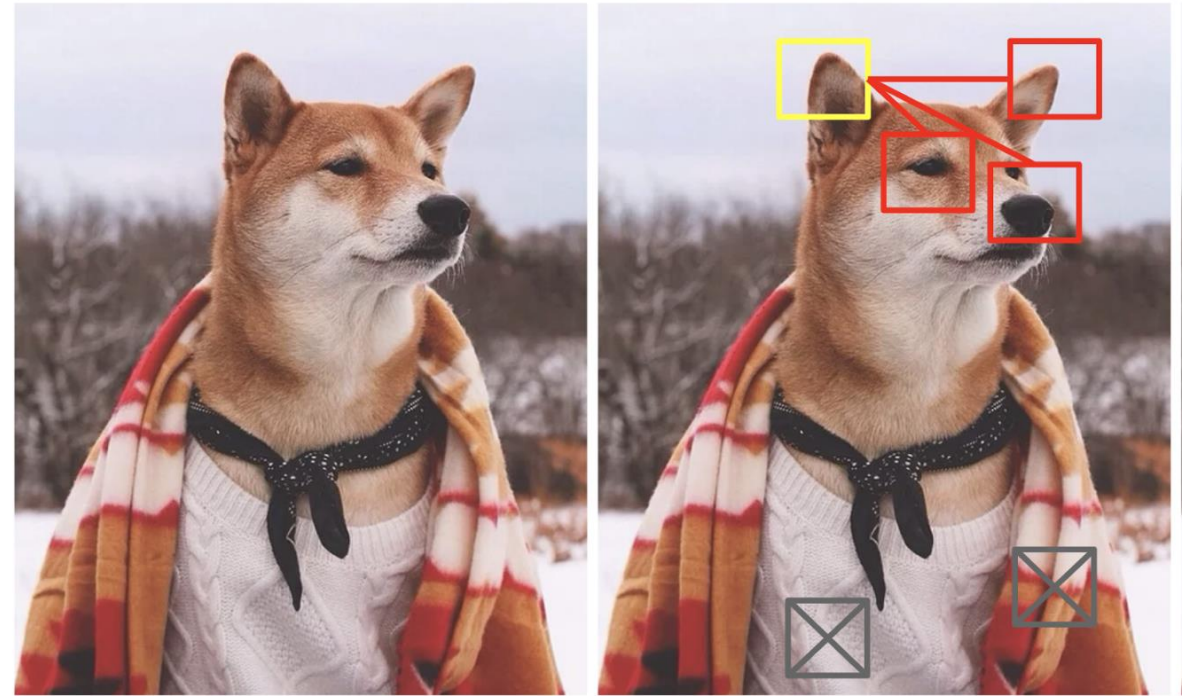


# Attention is All You Need!

A. Waswani et al., *NIPS*, 2017  
Google Brain & University of Toronto

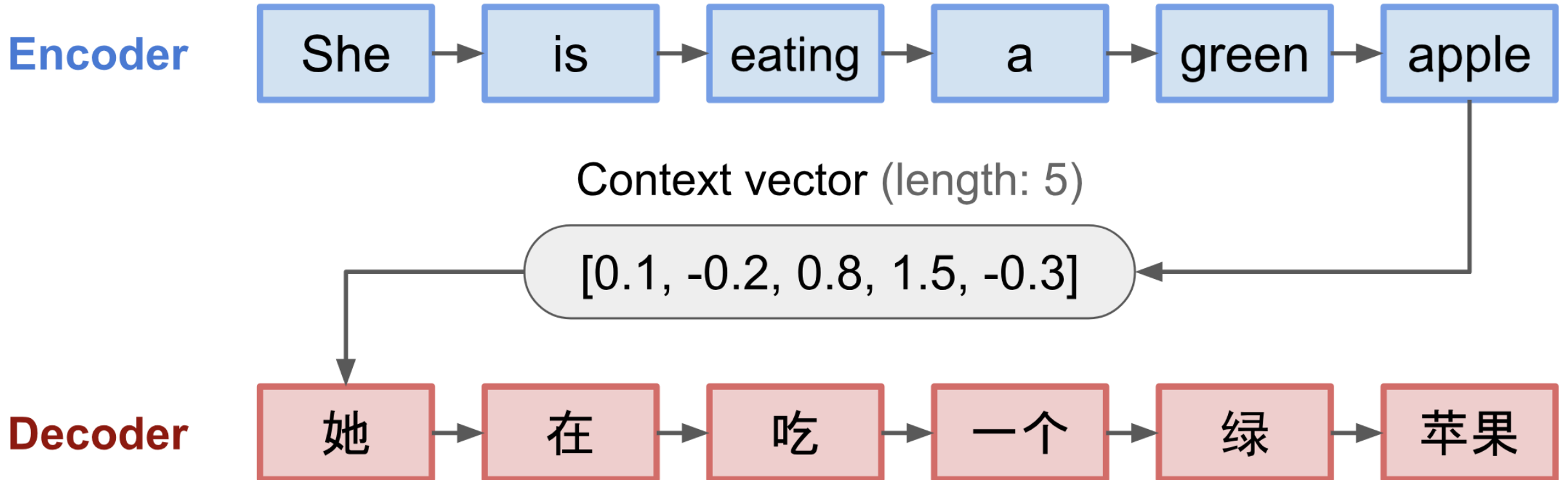
# Attention

- Visual attention and textual attention

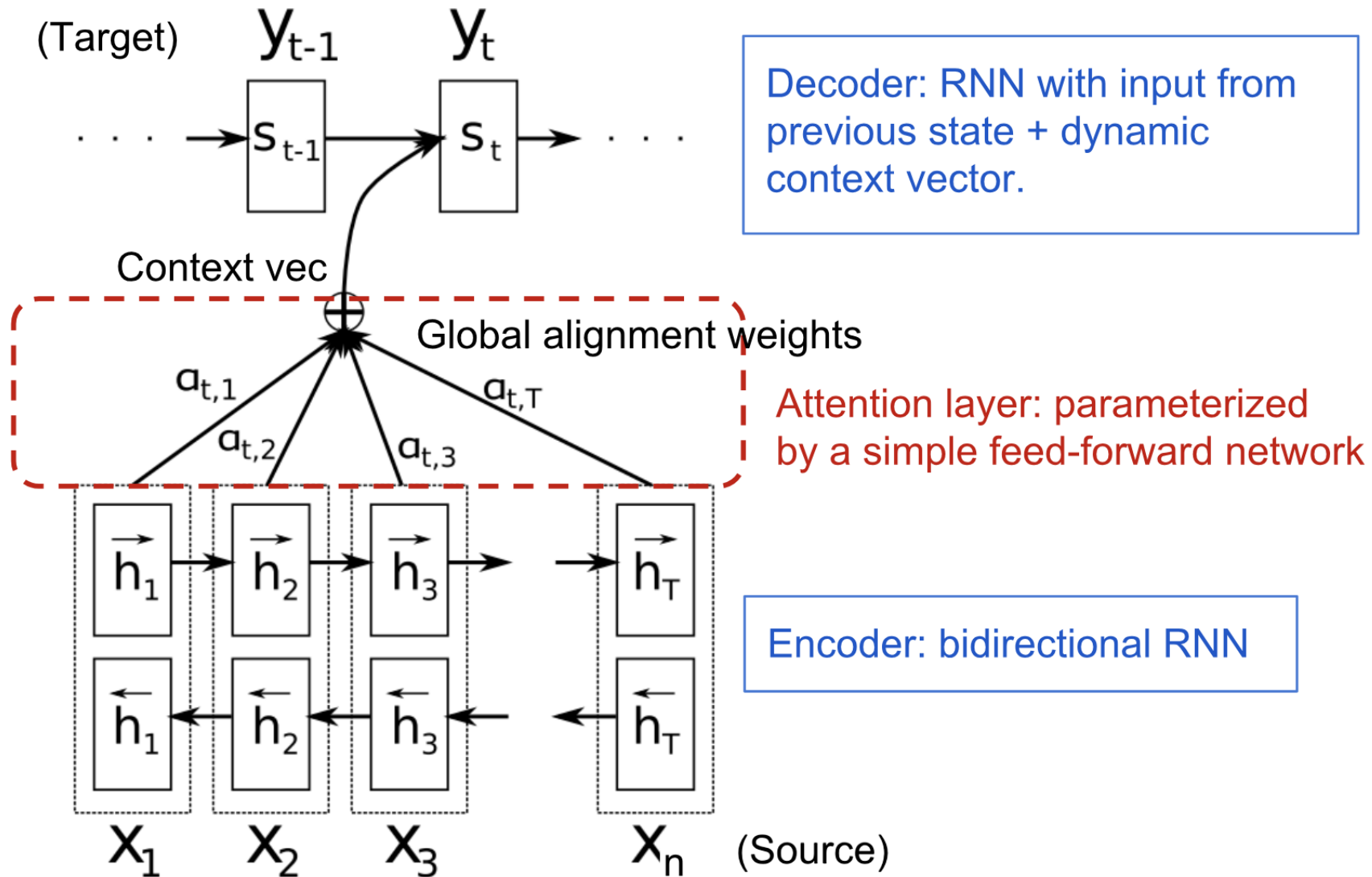


# Seq2seq model

- Language translation



# Attention = Vector of Importance Weights



**Additive Attention**

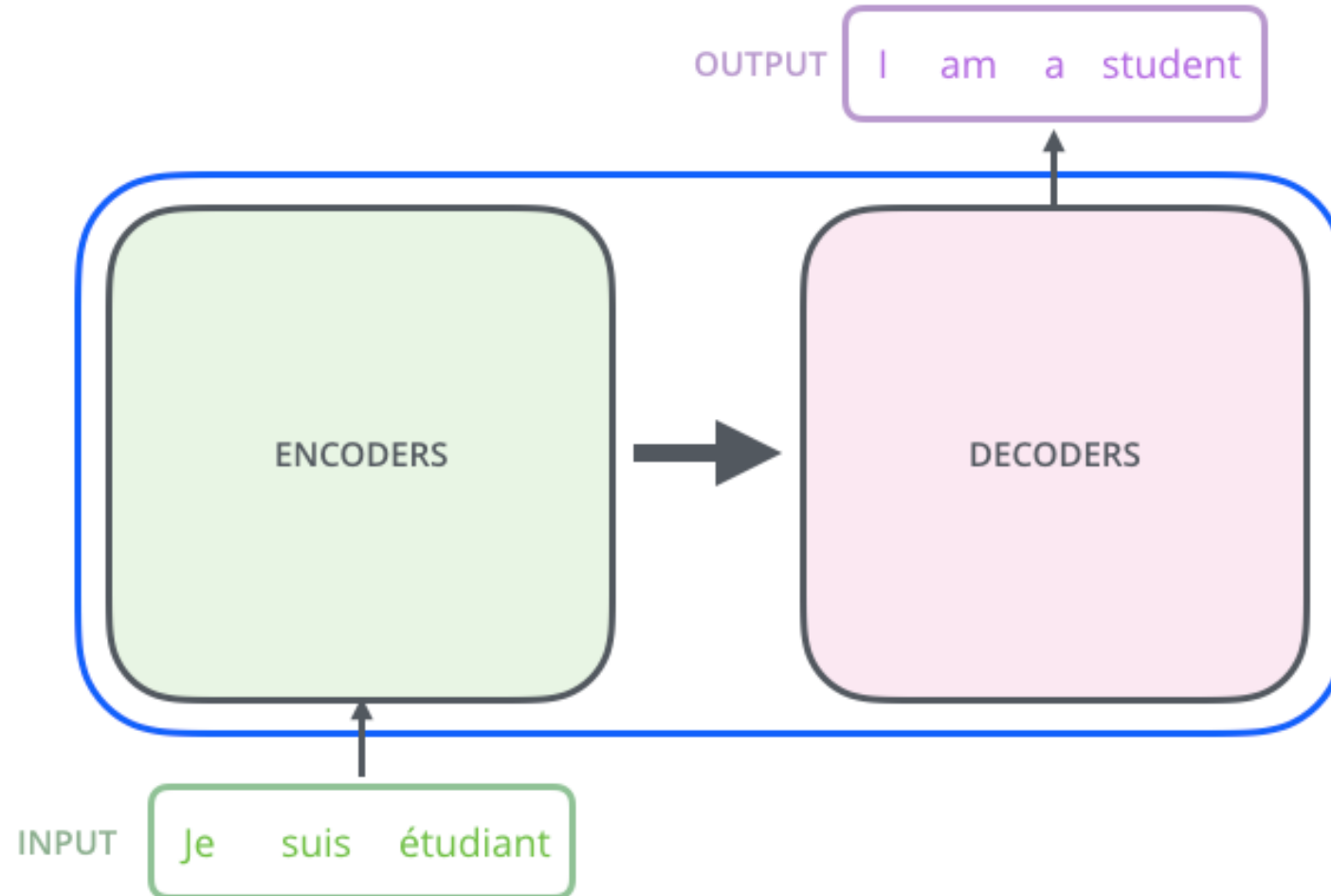
# Transformer

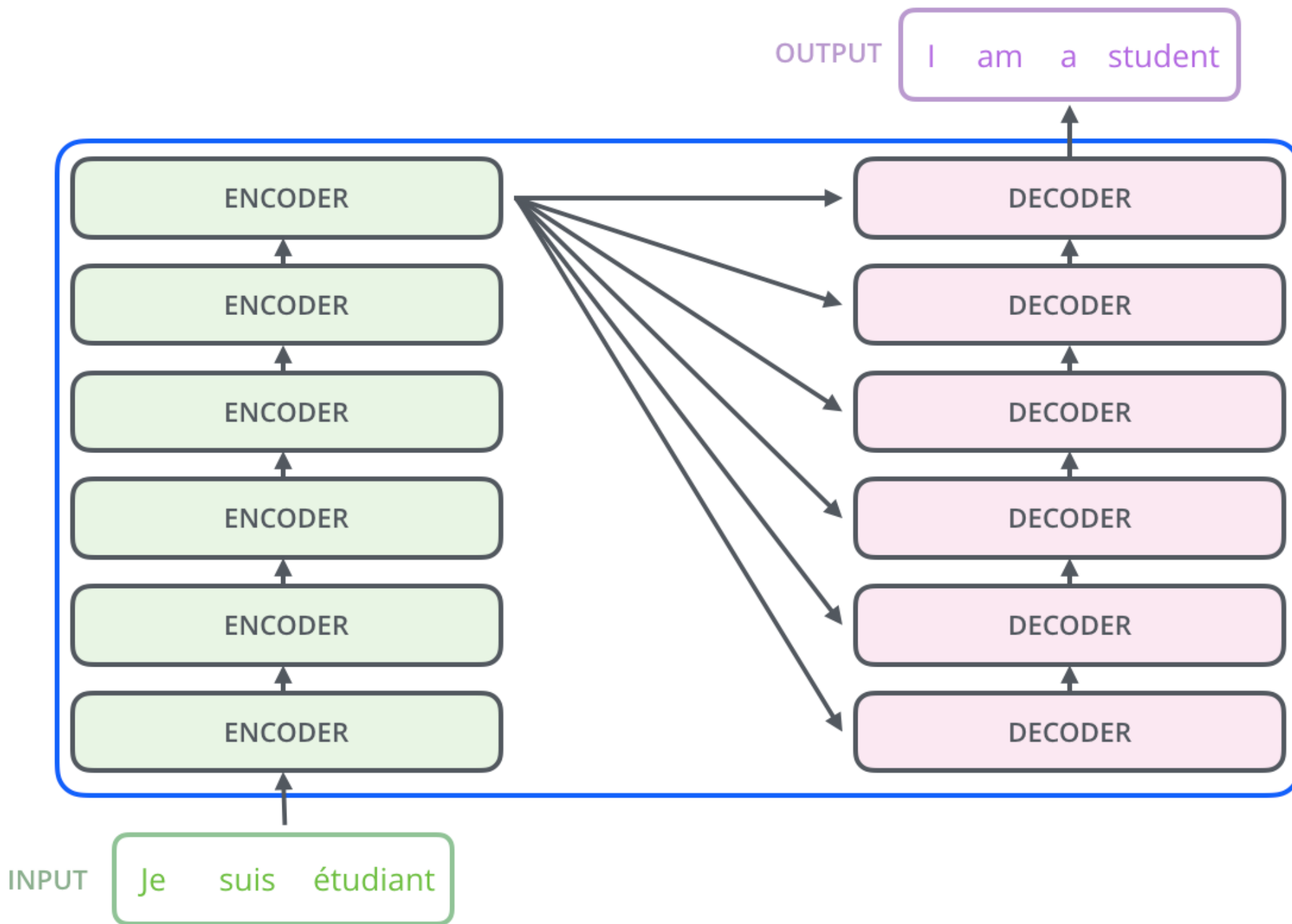


- <http://jalammar.github.io/illustrated-transformer/>



# Encoder and Decoder

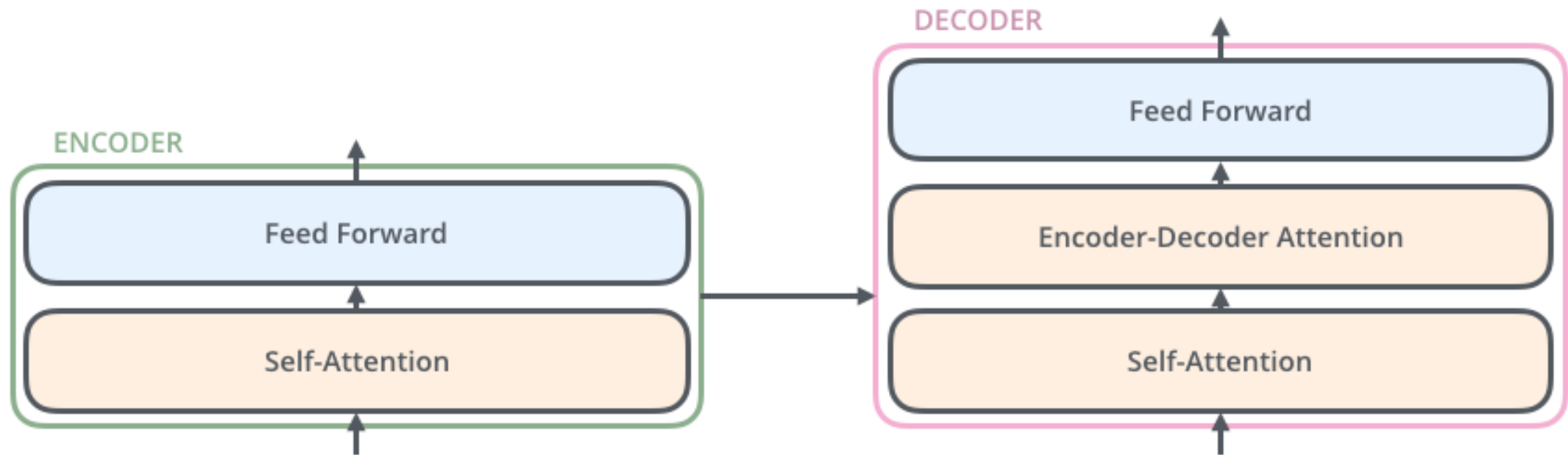






# Structure of the Encoder and Decoder

- Self-attention
- Encoder-decoder attention



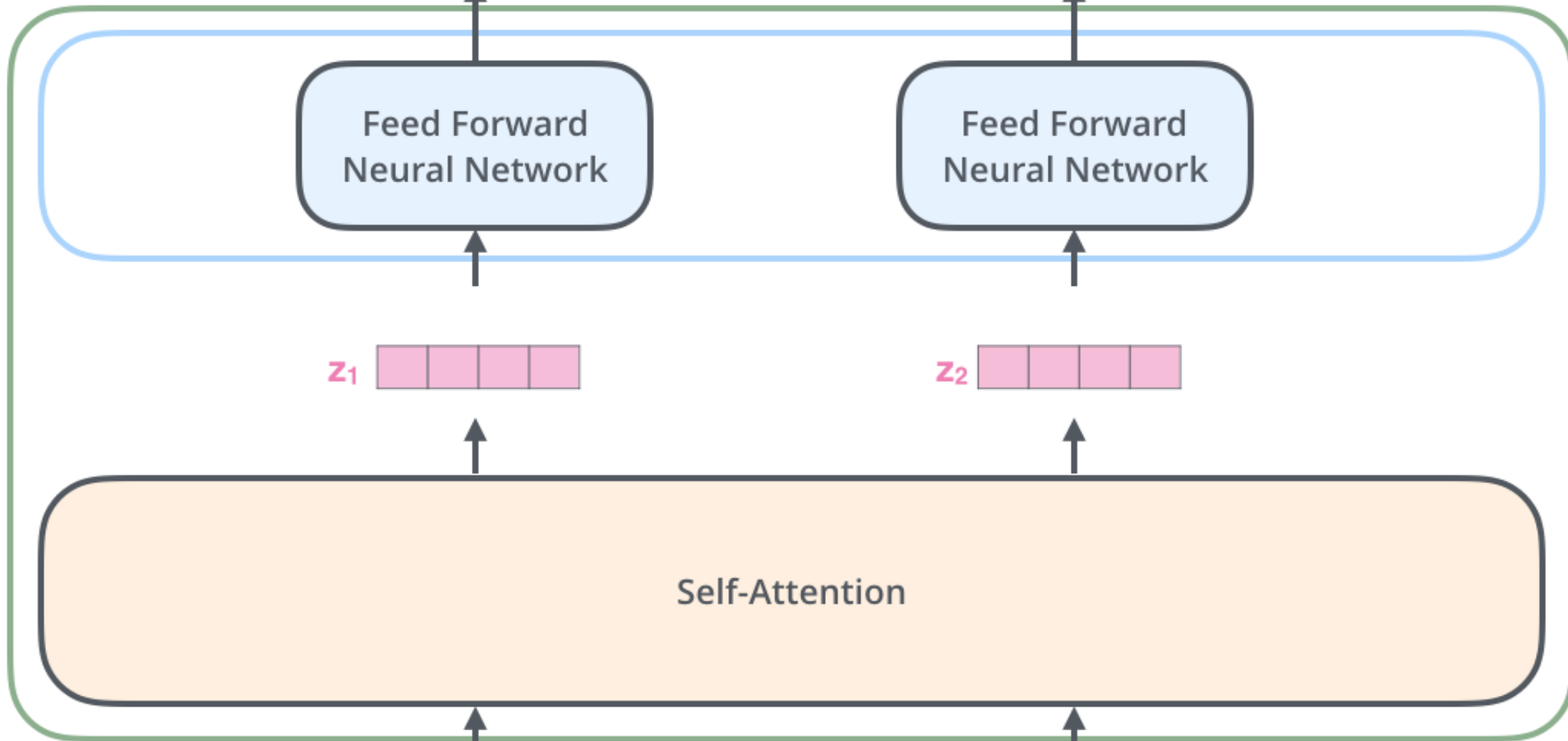
ENCODER #2



$r_1$  [ ] [ ] [ ] [ ]

$r_2$  [ ] [ ] [ ] [ ]

ENCODER #1



$z_1$  [ ] [ ] [ ] [ ]

$z_2$  [ ] [ ] [ ] [ ]

$x_1$  [ ] [ ] [ ] [ ]

$x_2$  [ ] [ ] [ ] [ ]

Thinking

Machines

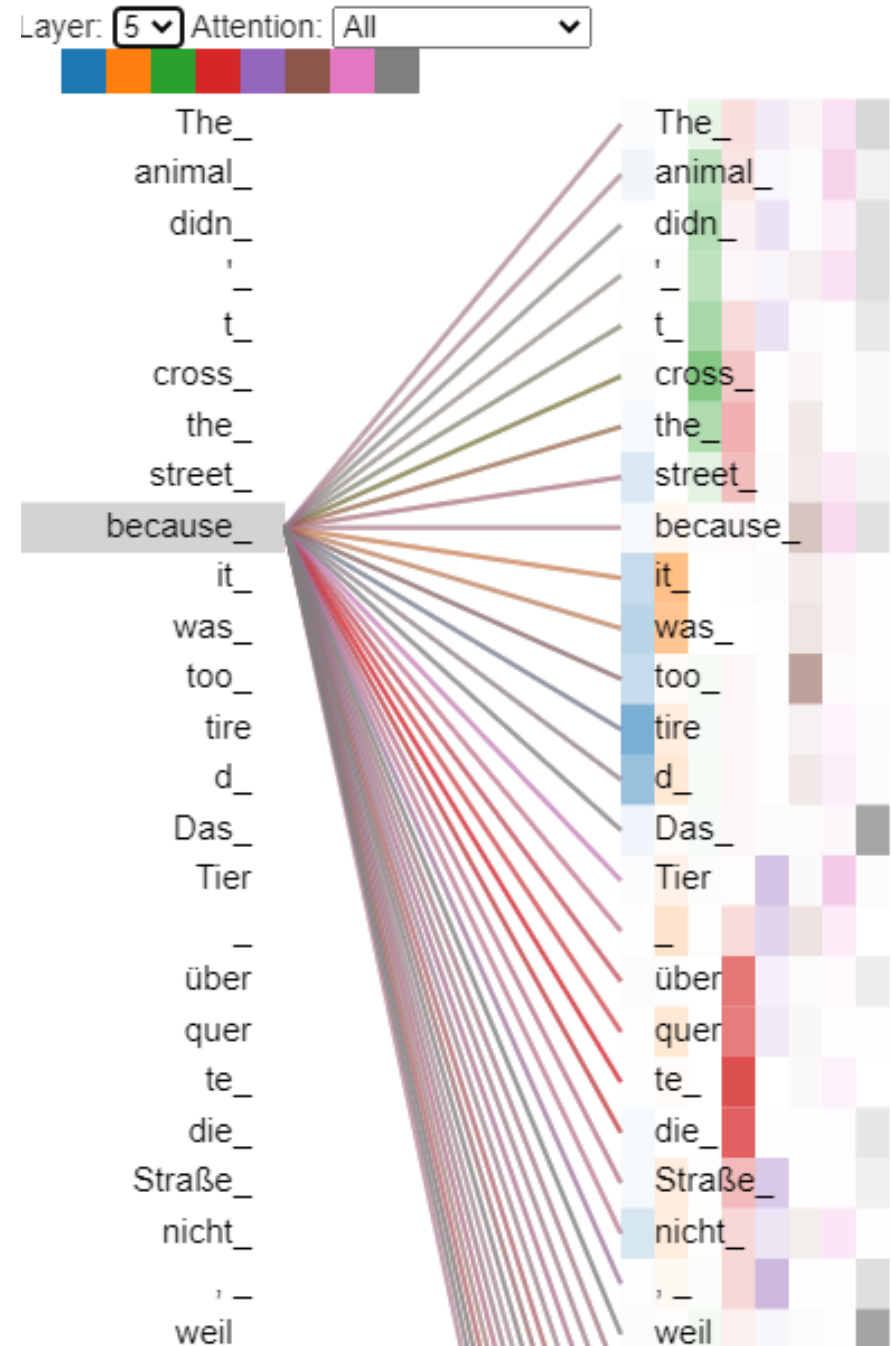
Self-Attention

Feed Forward  
Neural Network

Feed Forward  
Neural Network

# Tensor2Tensor Notebook

- [https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello\\_t2t.ipynb](https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb)

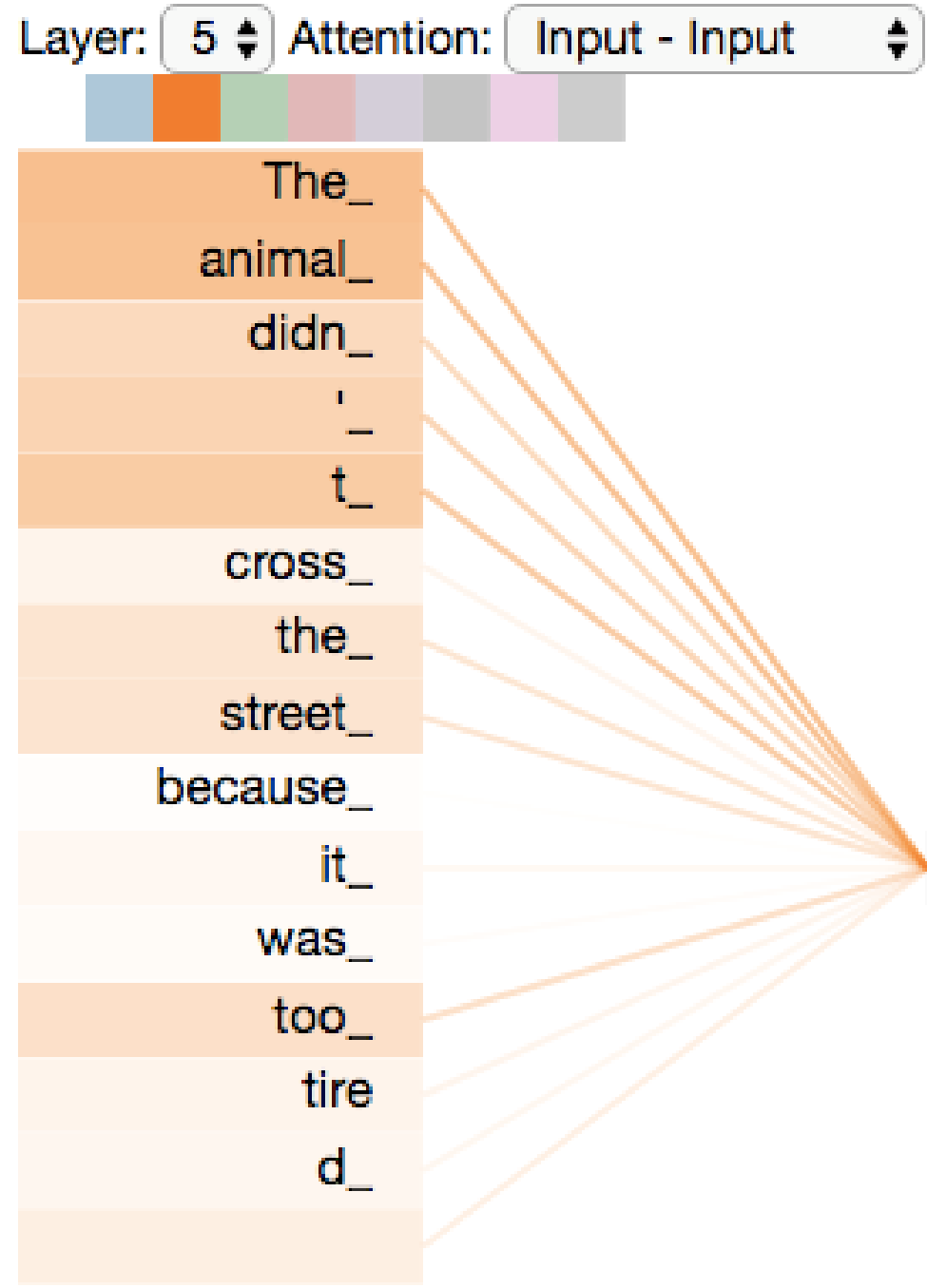




# Self Attention



# Self-attention Example



# Self-attention (query, key, value)

$q$ : query (to match others)

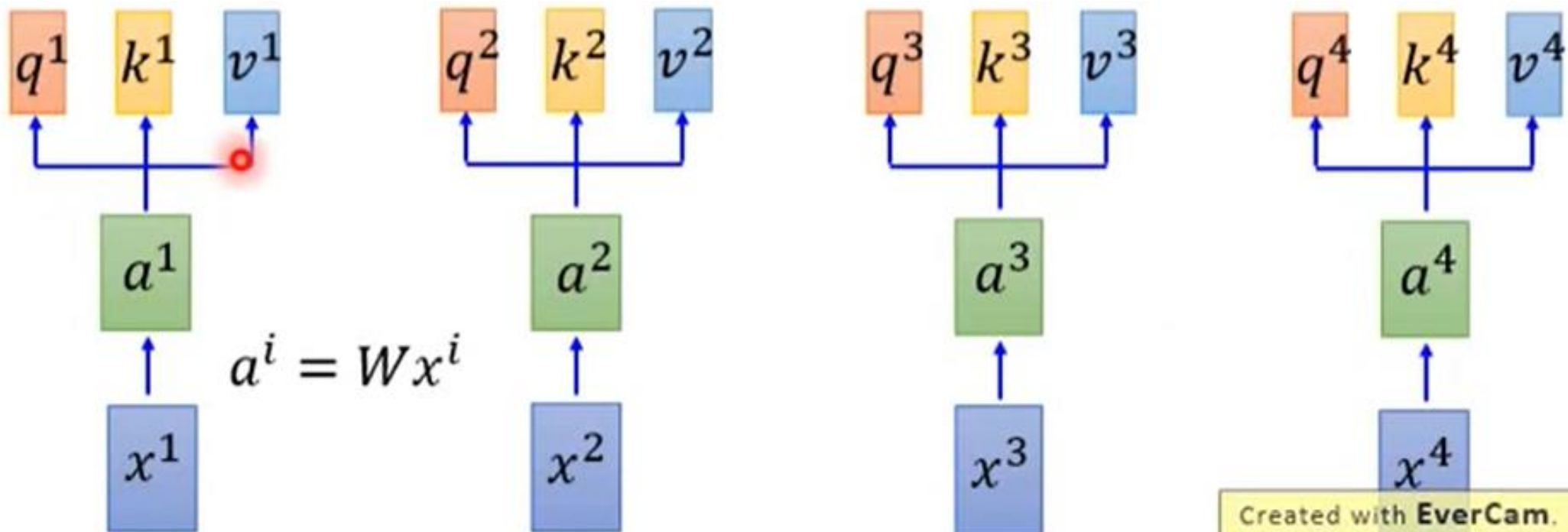
$$q^i = W^q a^i$$

$k$ : key (to be matched)

$$k^i = W^k a^i$$

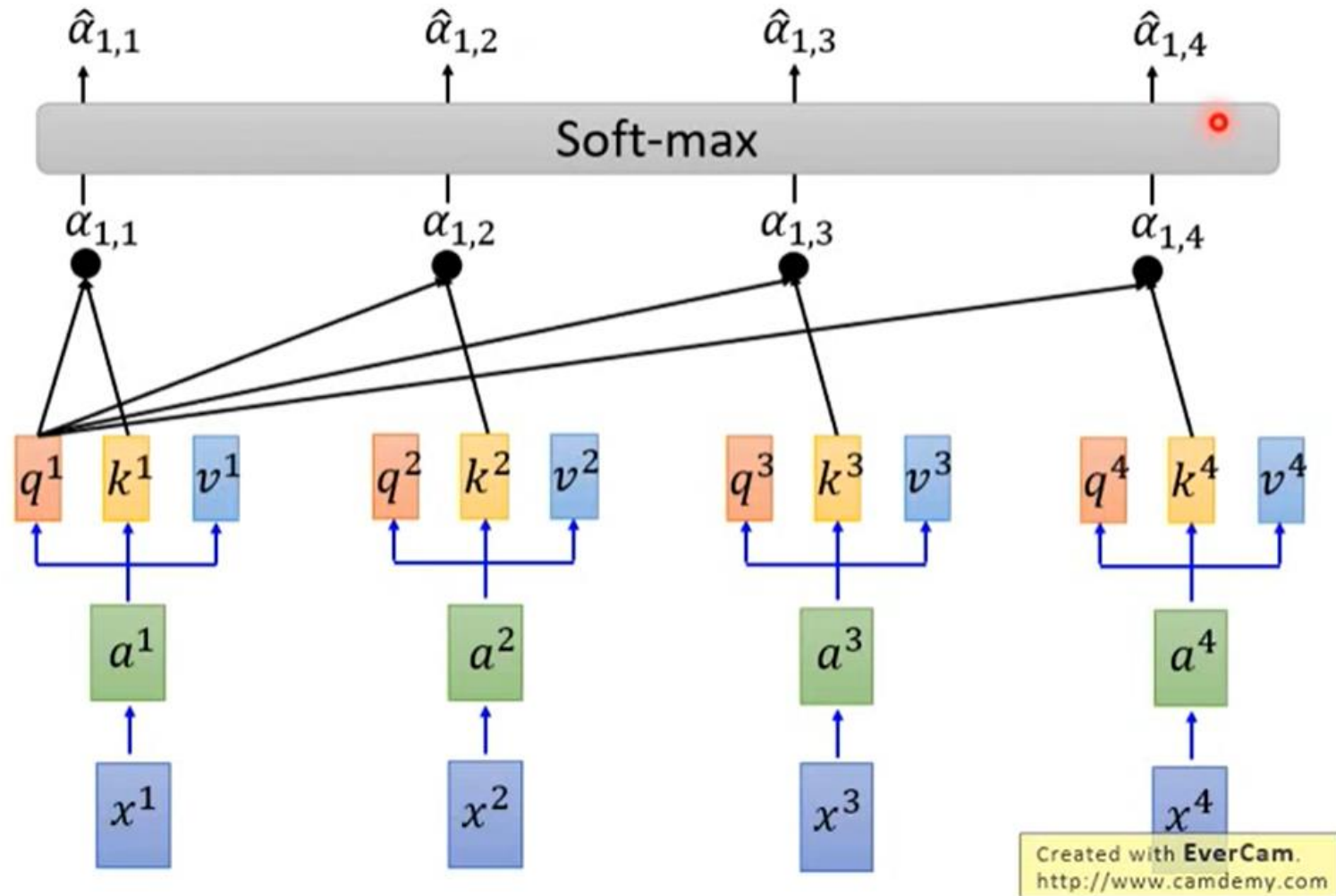
$v$ : information to be extracted

$$v^i = W^v a^i$$



Created with **EverCam**.  
<http://www.camdemy.com>

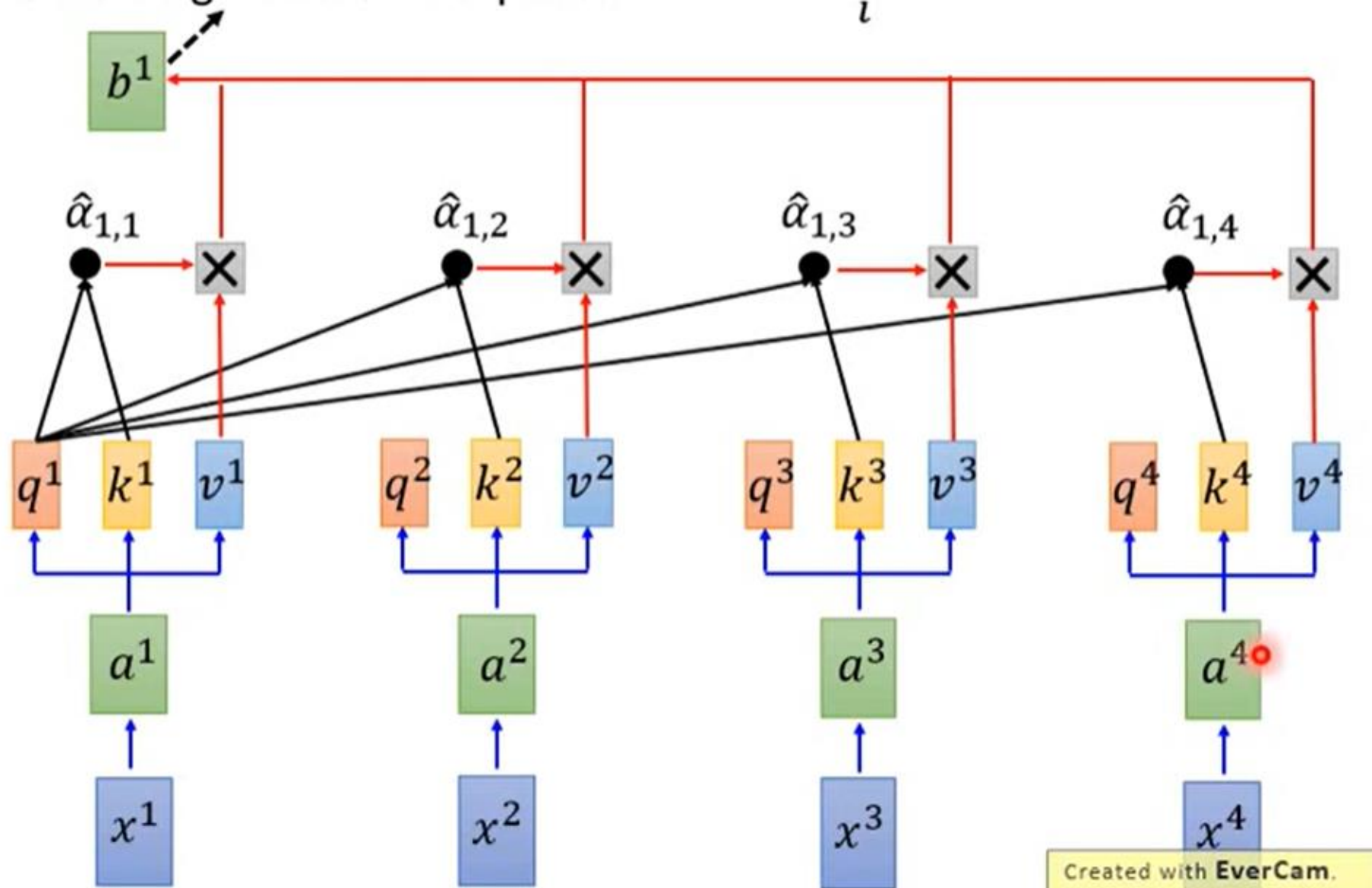
# Self-attention



# Self-attention

Considering the whole sequence

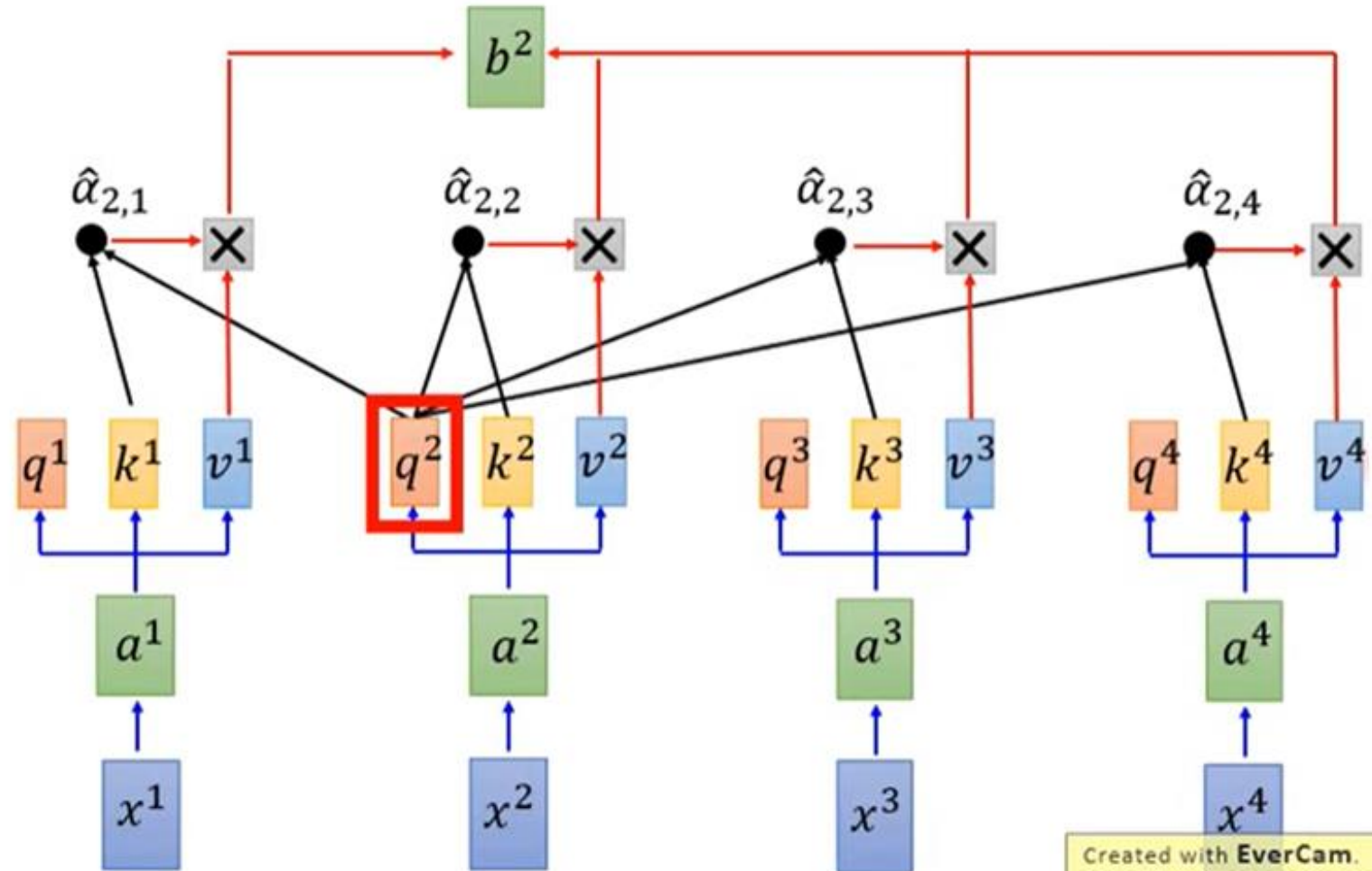
$$b^1 = \sum_i \hat{\alpha}_{1,i} v^i$$





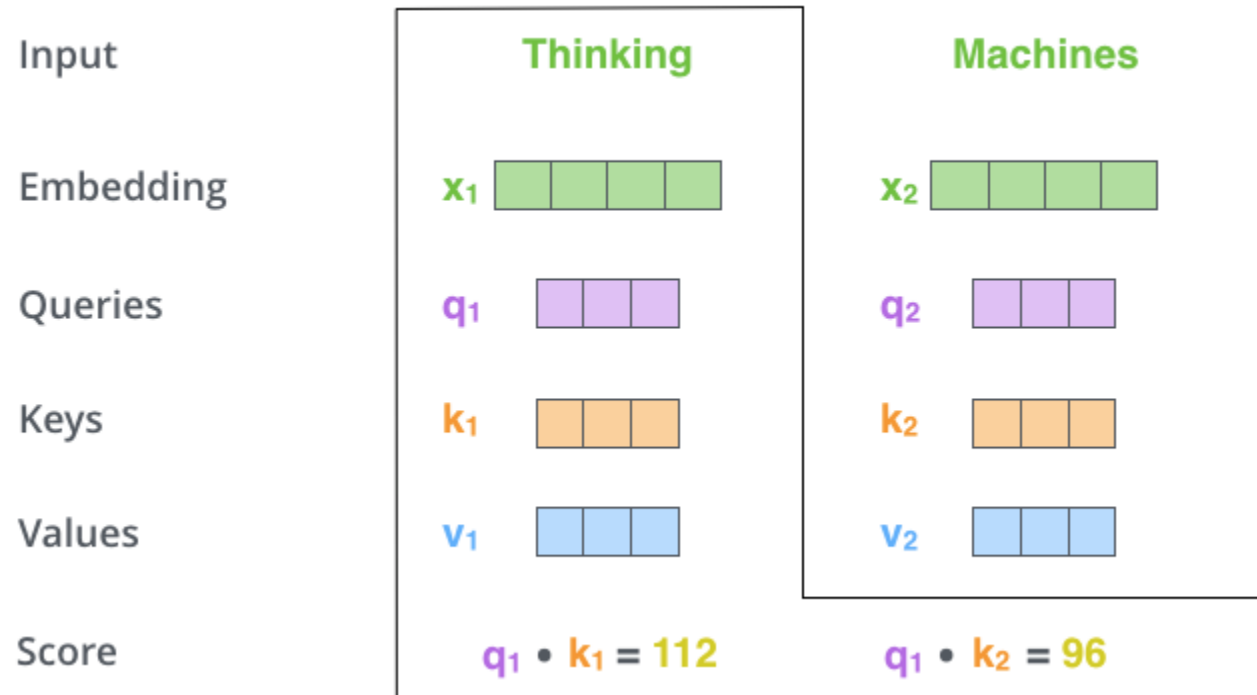
# Calculating $b^2$

$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$

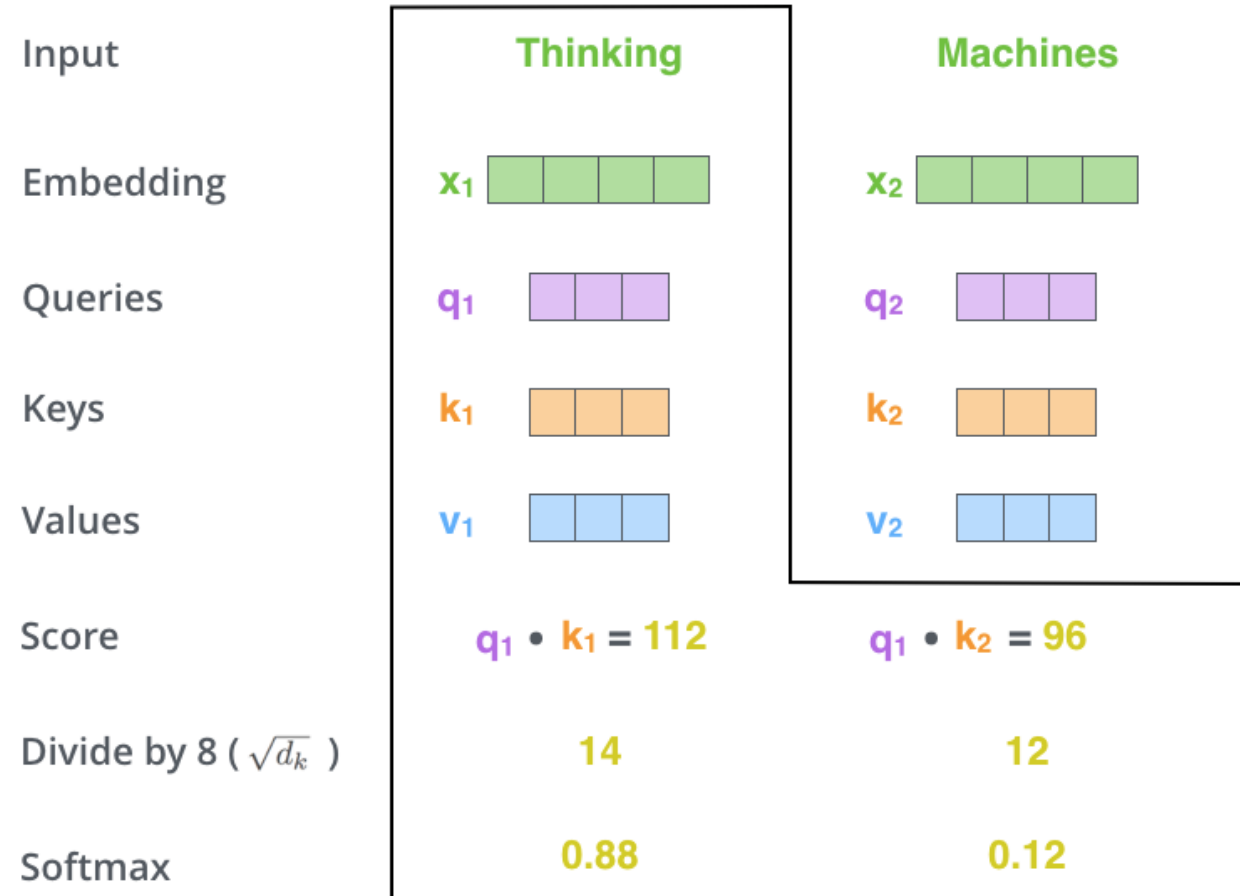


Created with EverCam.  
<http://www.camdemy.com>

# Calculate Dot Product of Query & Value



# Softmax Normalization



# Calculate Values

- Multiply each value vector by the Softmax score
- Sum up the weighted value vectors

Input

Embedding

Queries

Keys

Values

Score

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

Softmax

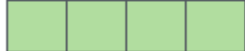
Softmax

X

Value

Sum

Thinking

$x_1$  

$q_1$  

$k_1$  

$v_1$  

$q_1 \cdot k_1 = 112$

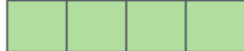
14

0.88

$v_1$  

$z_1$  

Machines

$x_2$  

$q_2$  

$k_2$  

$v_2$  

$q_1 \cdot k_2 = 96$

12

0.12

$v_2$  

$z_2$  

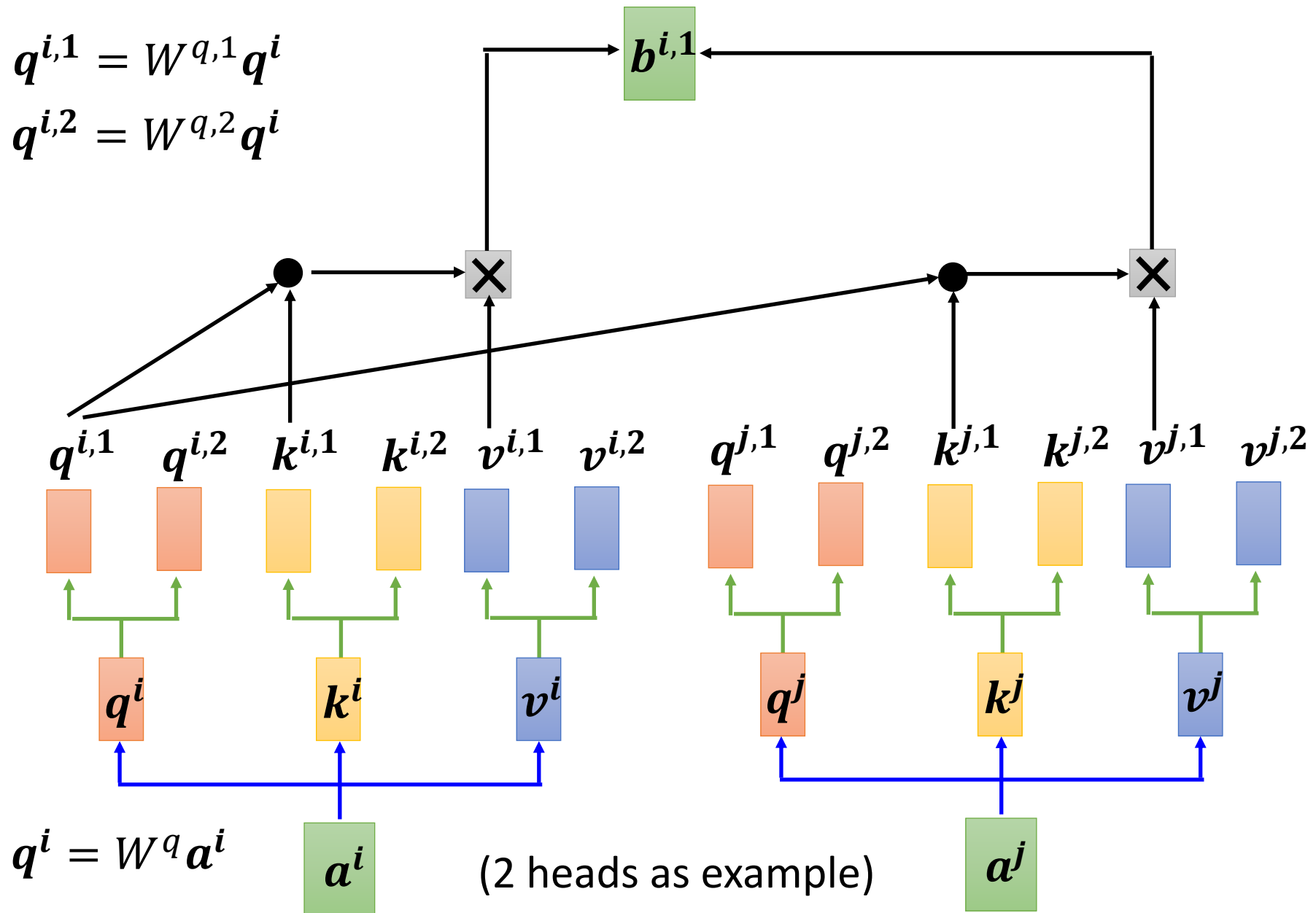
# Final Output of Self-attention

$$\text{softmax} \left( \frac{\begin{matrix} \mathbf{Q} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix} \times \begin{matrix} \mathbf{K}^T \\ \begin{array}{|c|c|} \hline \square & \square \\ \hline \square & \square \\ \hline \square & \square \\ \hline \end{array} \end{matrix} \right) \begin{matrix} \mathbf{V} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$
$$= \begin{matrix} \mathbf{Z} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$



# The Beast With Multiple Heads

# Multi-head Self-attention



1) This is our input sentence\*

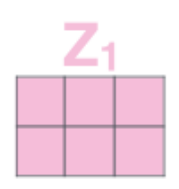
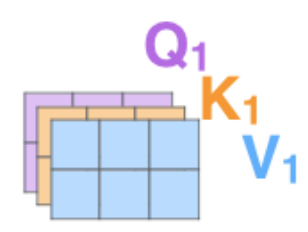
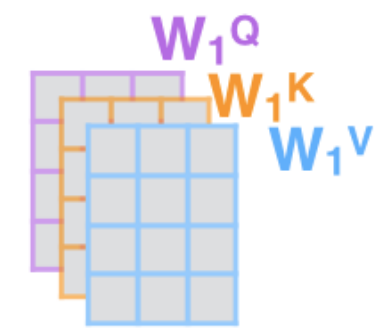
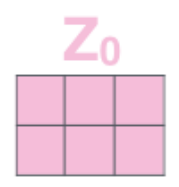
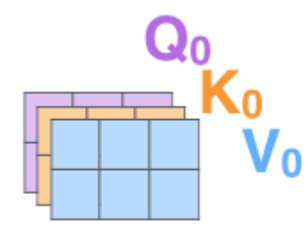
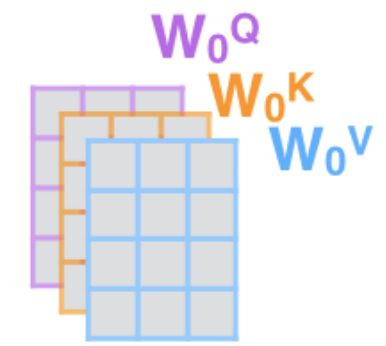
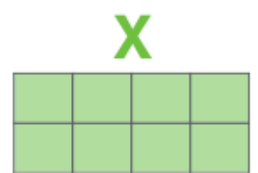
2) We embed each word\*

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

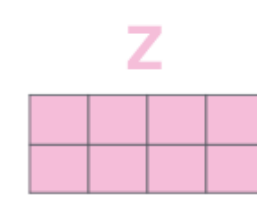
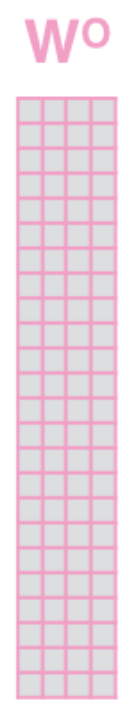
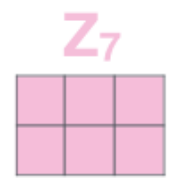
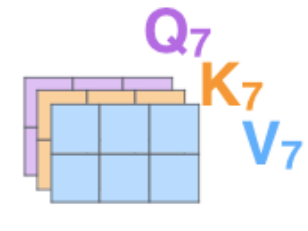
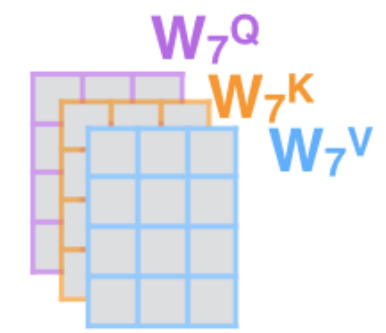
Thinking Machines



...

...

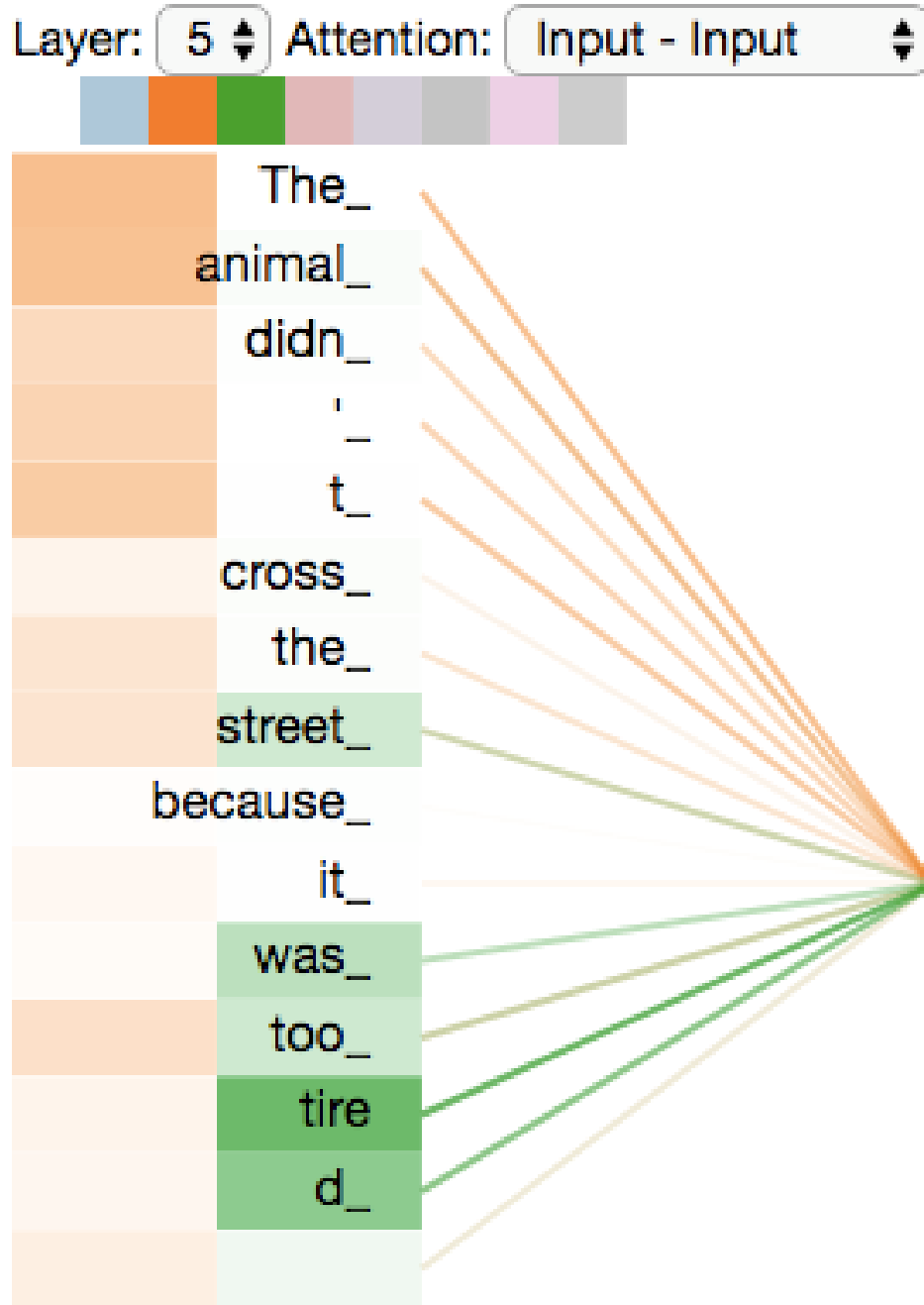
...



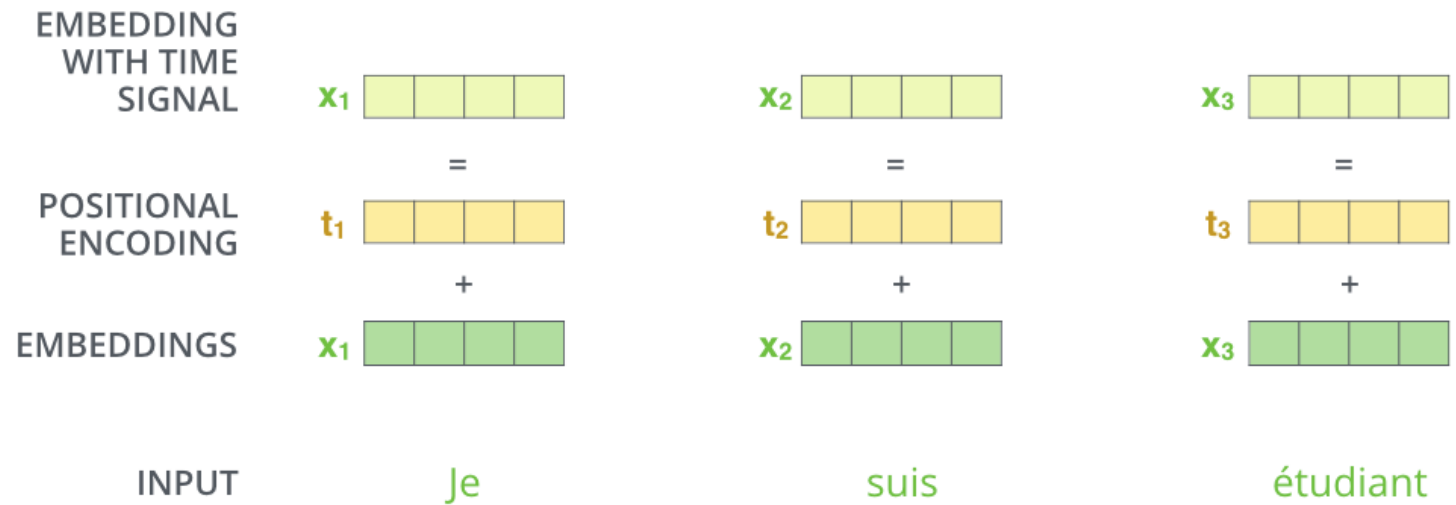
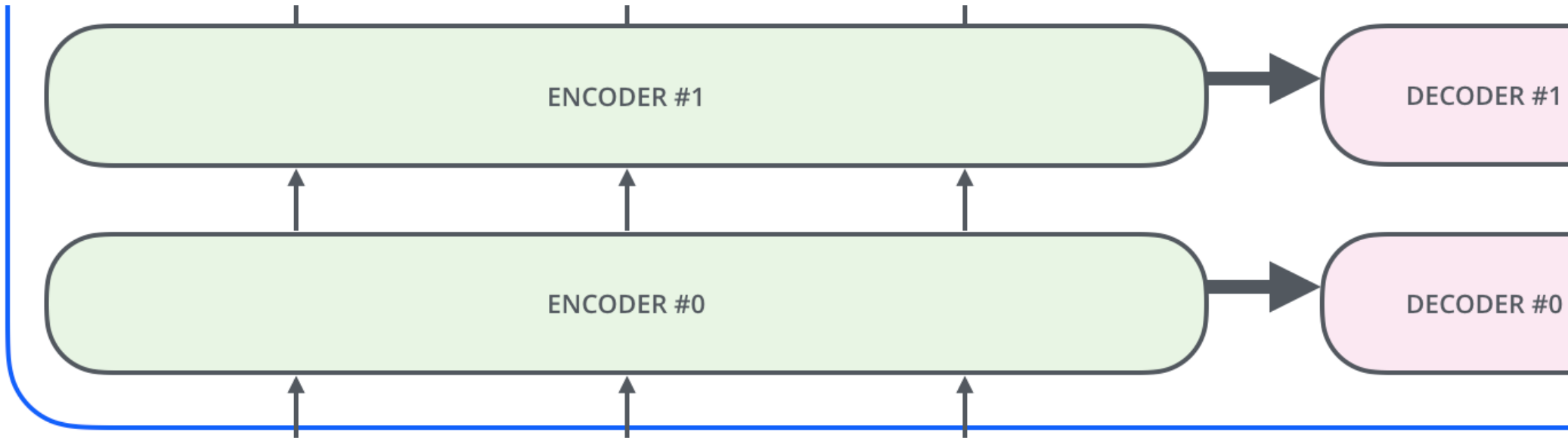
\* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



# Two-head Self-attention

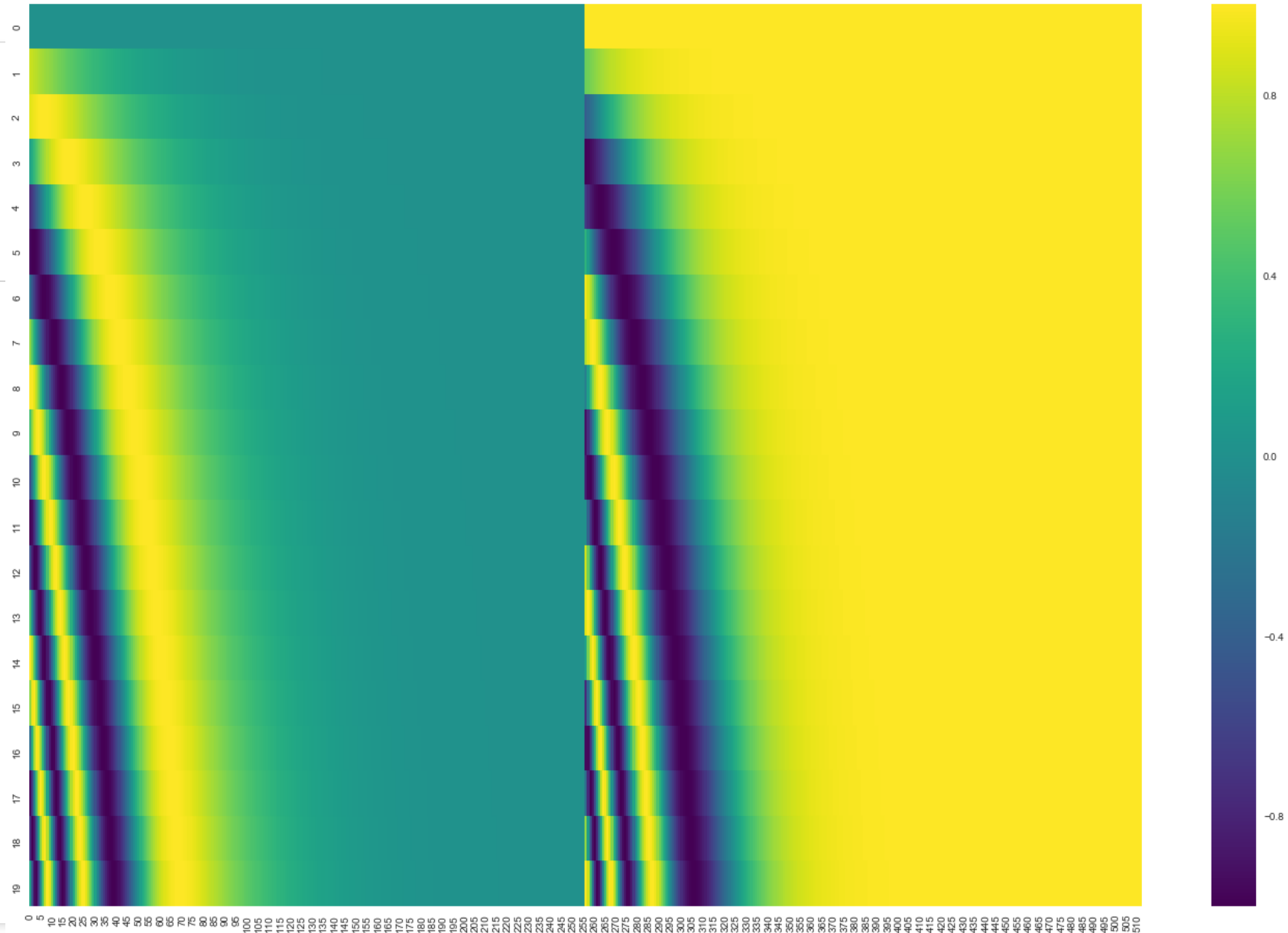


# Positional Encoding

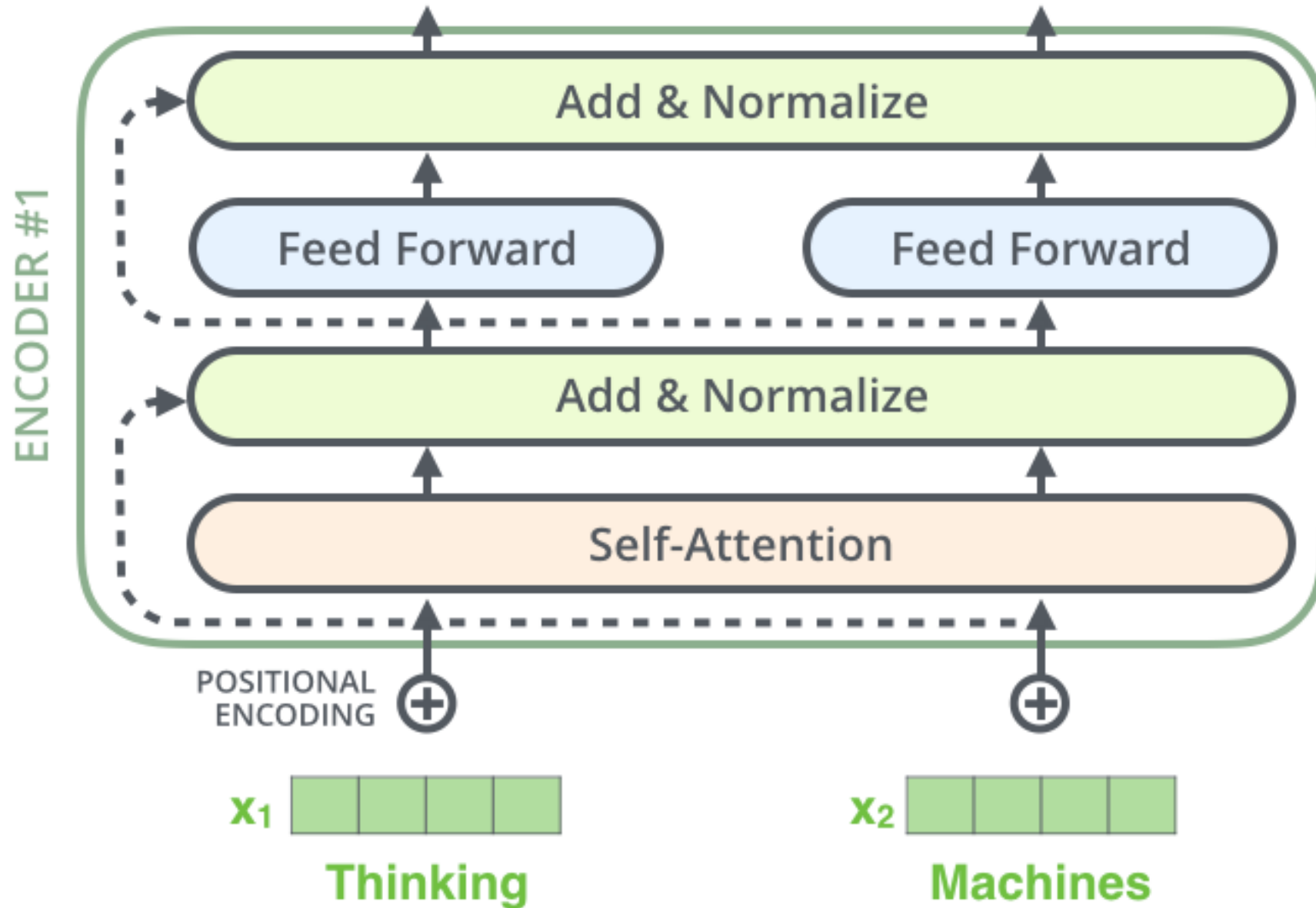


# Positional Encoding Visualization

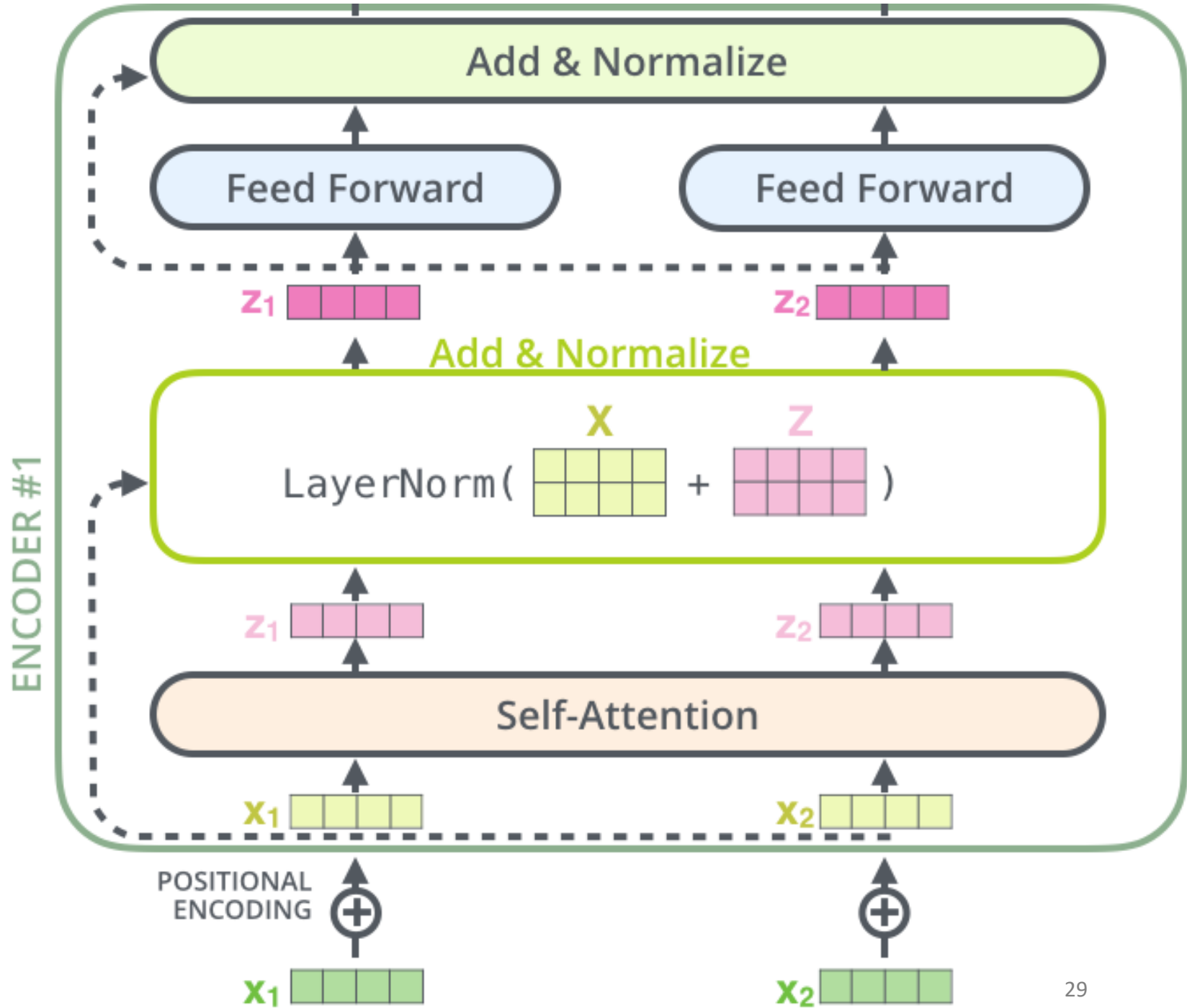
- Visualize positional encoding for 20 words (rows) with an embedding size of 512 (columns)

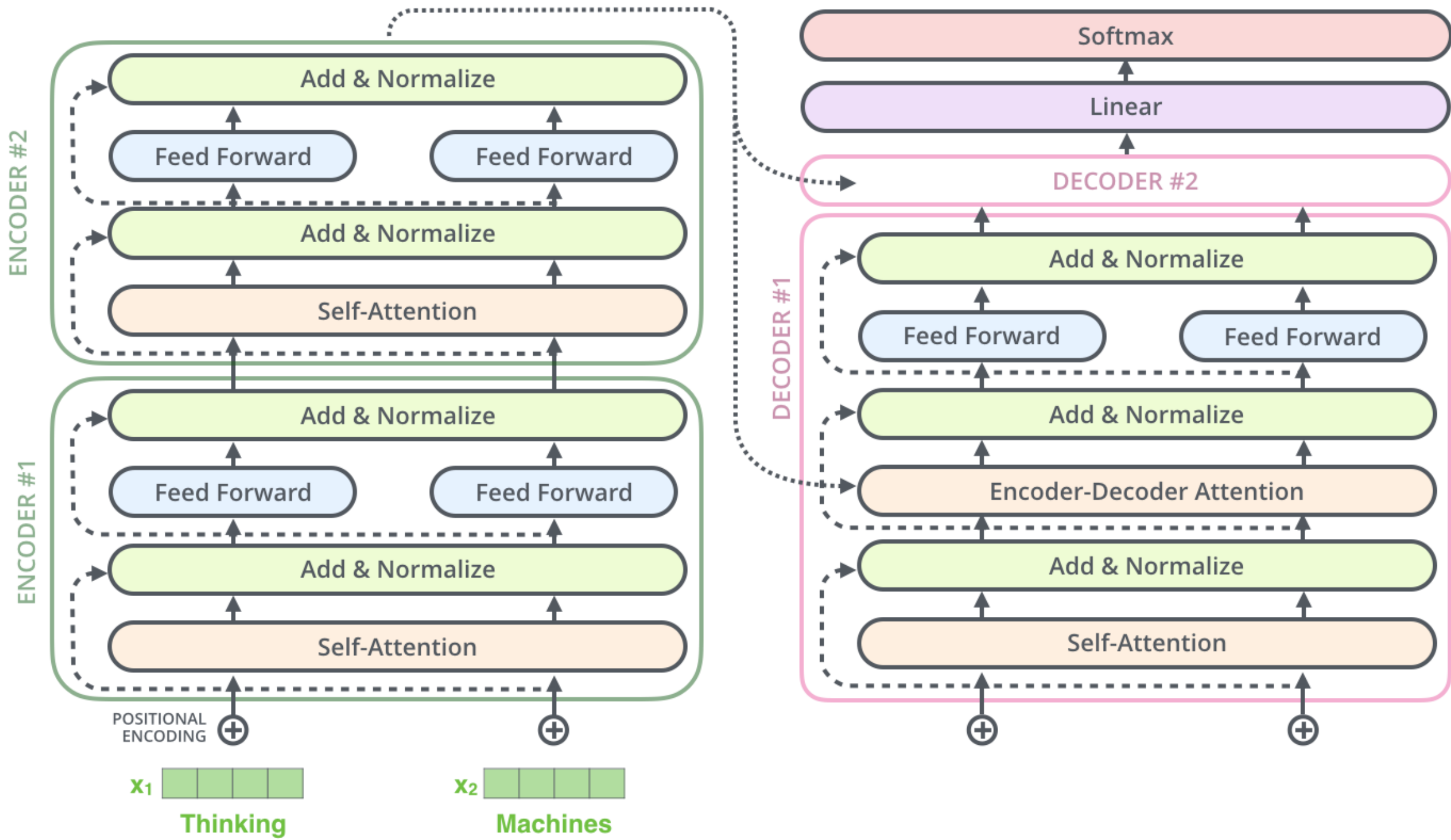


# Adding Residual Connections



# Layer Normalization



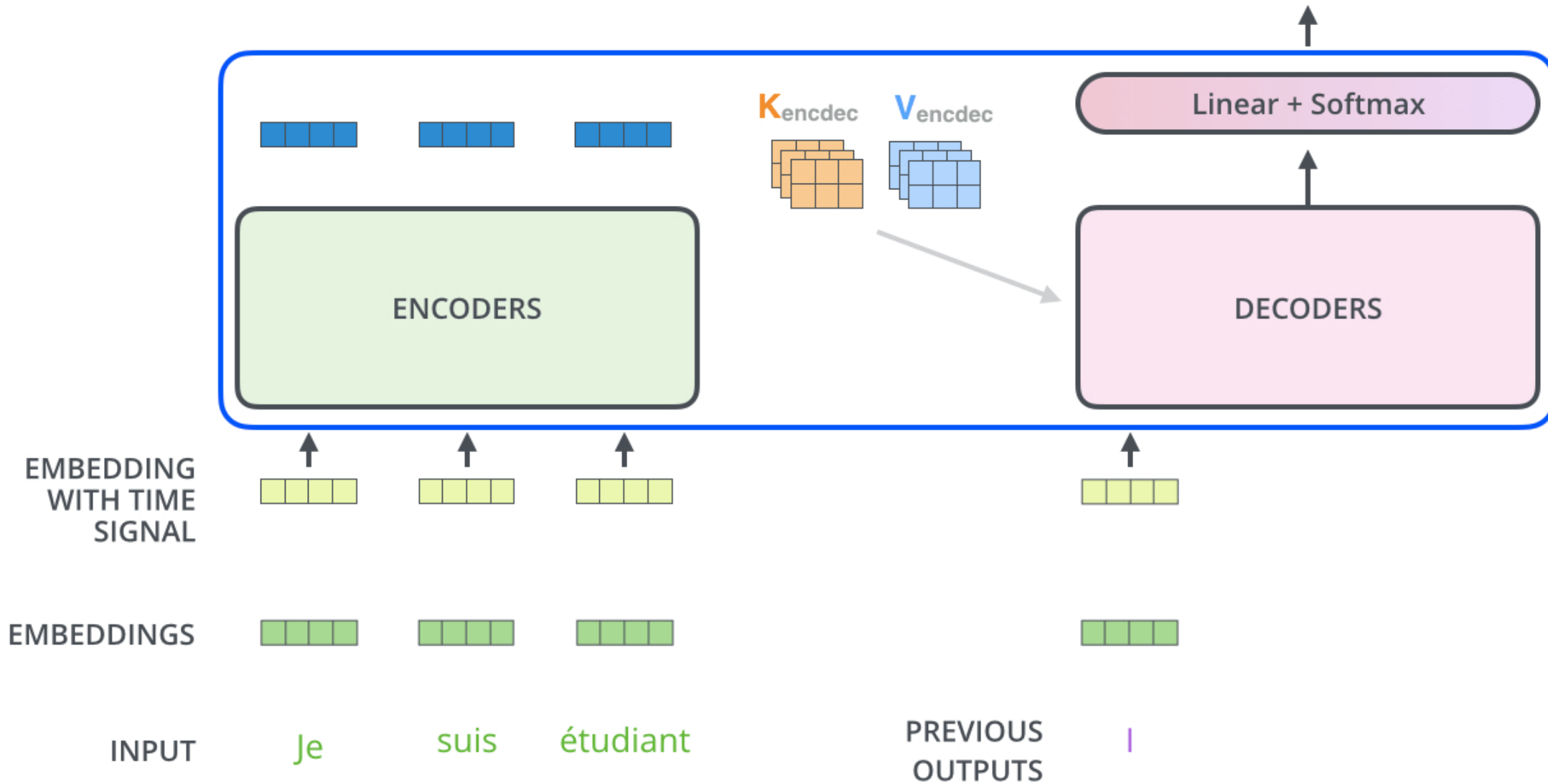


# Decoder

Decoding time step: 1 2 3 4 5 6

OUTPUT |

# Decoder





Which word in our vocabulary  
is associated with this index?

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

**log\_probs**



am

5

Softmax

**logits**



Linear

Decoder stack output



# Output of Decoder

Output Vocabulary

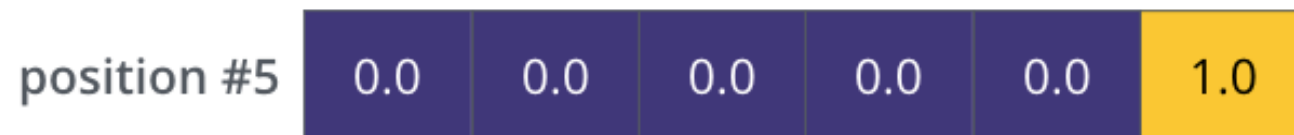
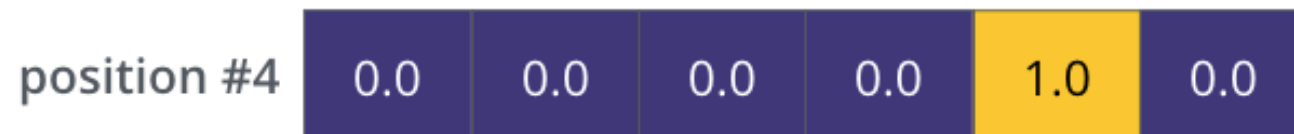
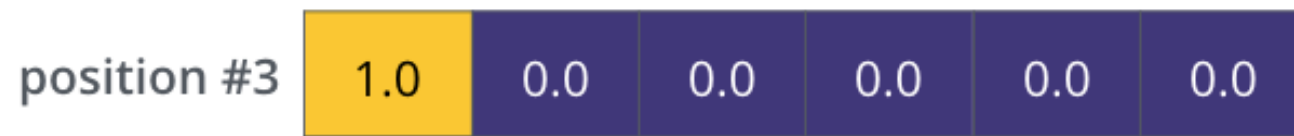
WORD	a	am	I	thanks	student	<eos>
INDEX	0	1	2	3	4	5

One-hot encoding of the word "am"



# Target Model Outputs

Output Vocabulary: a am I thanks student <eos>

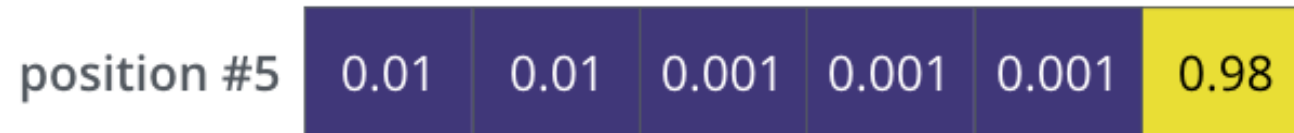
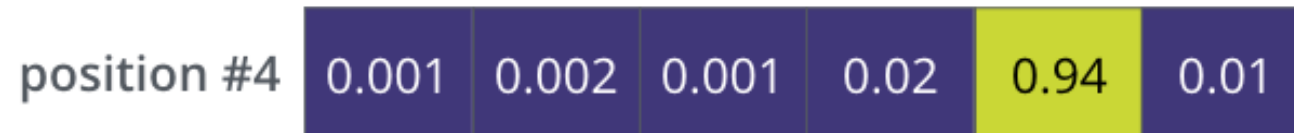
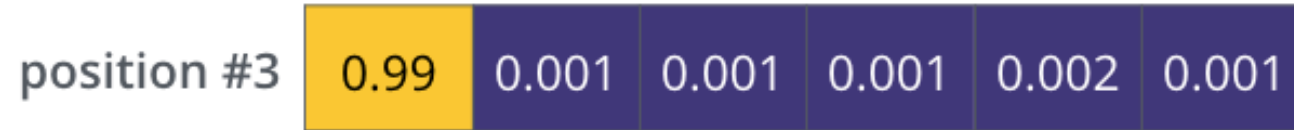


a am I thanks student <eos>



# Trained Model Outputs

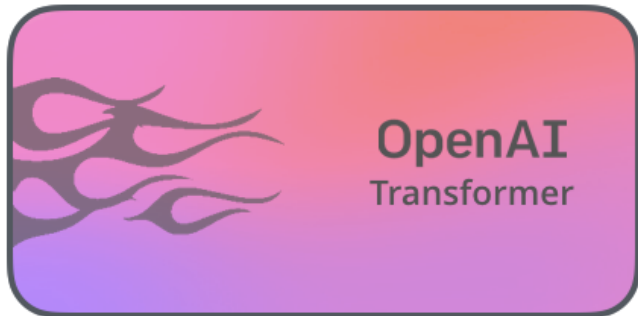
Output Vocabulary: a am I thanks student <eos>



a am I thanks student <eos>



# BERT, ELMo, GPT



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.

### Semi-supervised Learning Step

**Model:**



**Dataset:**



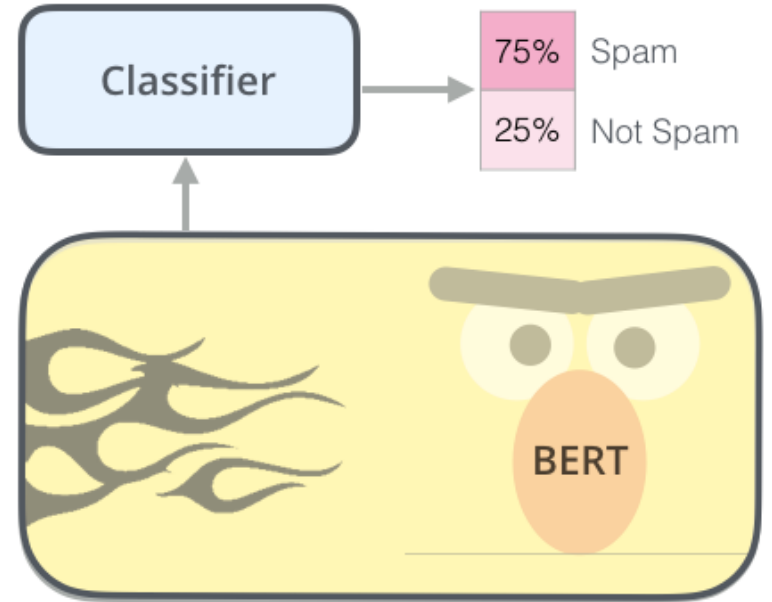
**Objective:**

Predict the masked word (language modeling)

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

### Supervised Learning Step

**Model:**  
(pre-trained in step #1)



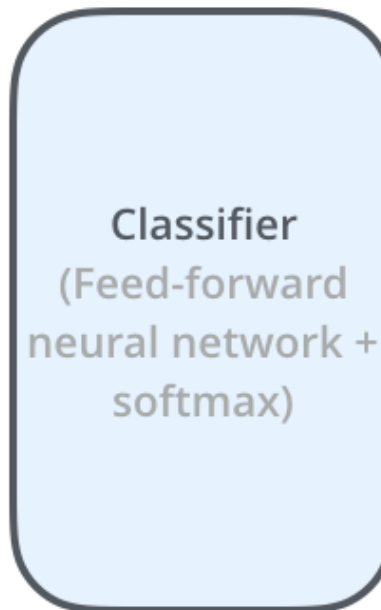
**Dataset:**

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

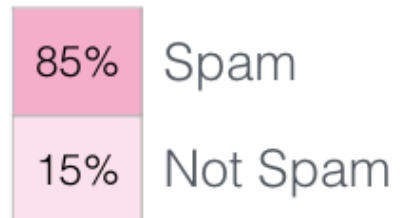
# Sentence Classification

Input  
Features

Help Prince Mayuko Transfer  
Huge Inheritance

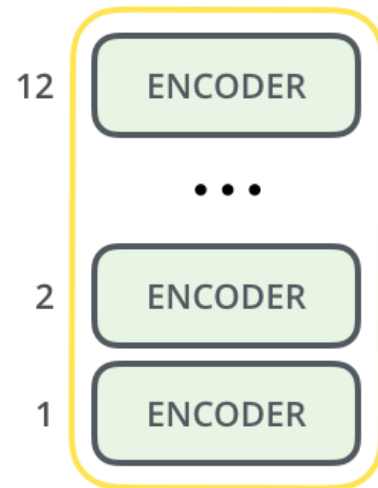


Output  
Prediction

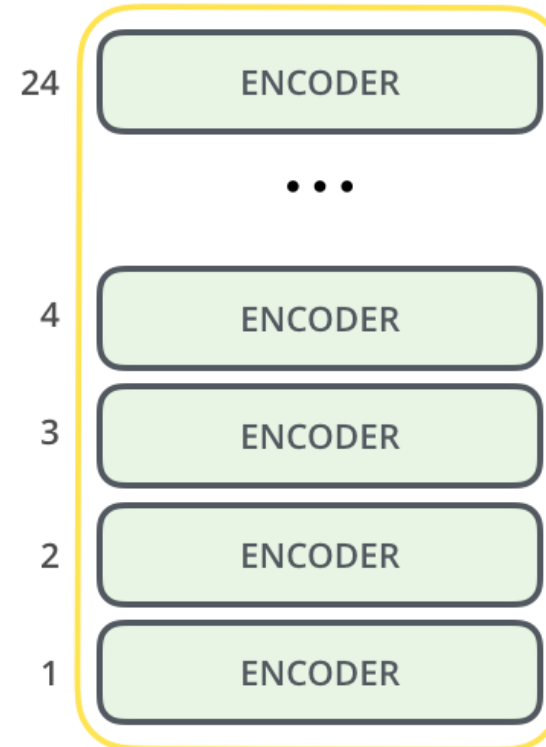


# Model Architecture

- BERT BASE – Comparable in size to the OpenAI Transformer in order to compare performance
- BERT LARGE – A ridiculously huge model which achieved the state of the art results reported in the paper



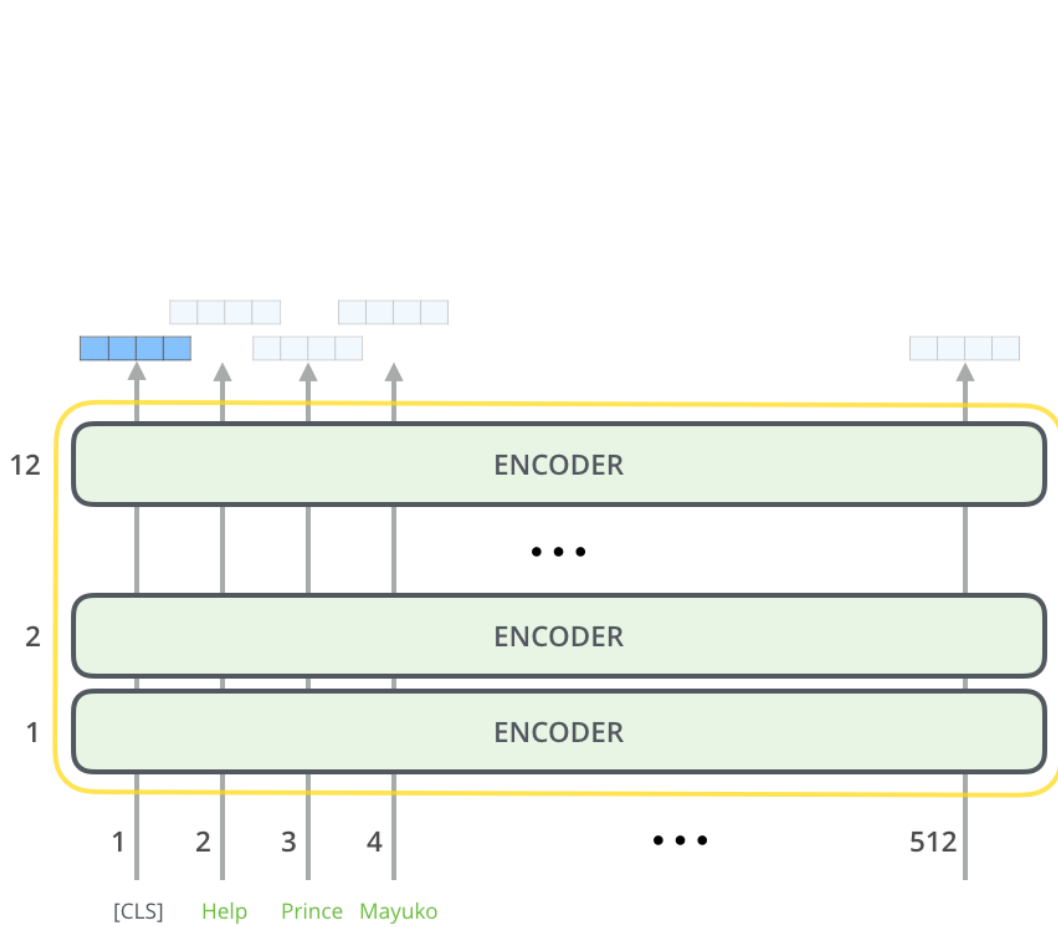
BERT<sub>BASE</sub>



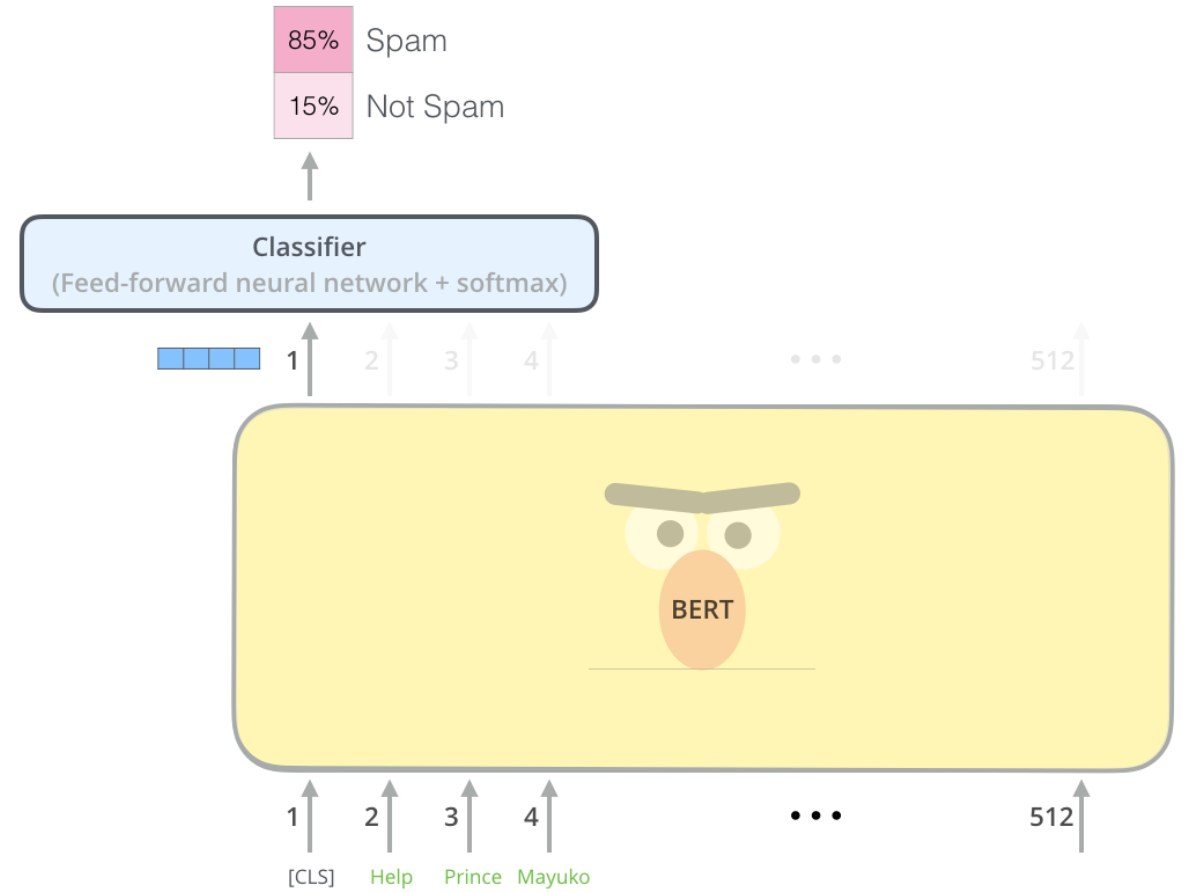
BERT<sub>LARGE</sub>



# Use Pre-trained Encoders



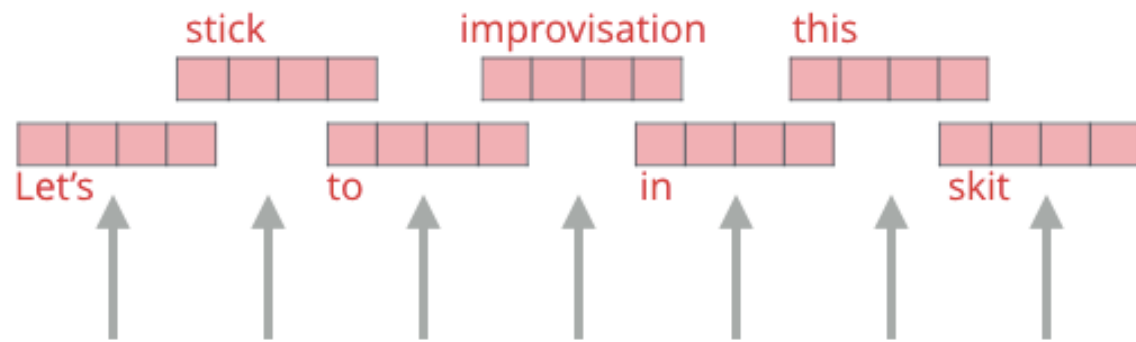
BERT



# ELMo: Context Matters



# ELMo Embeddings



Words to embed



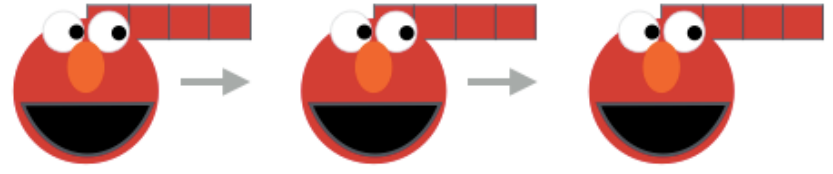
Possible classes:  
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzzzyva

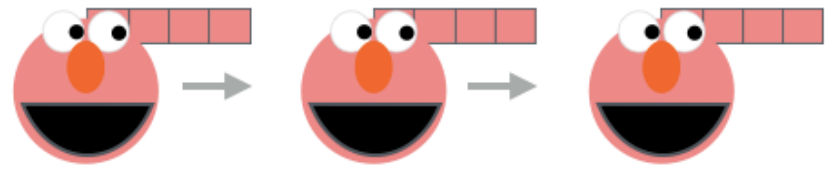
Output  
Layer



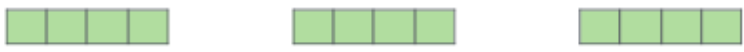
LSTM  
Layer #2



LSTM  
Layer #1



Embedding



Let's

stick

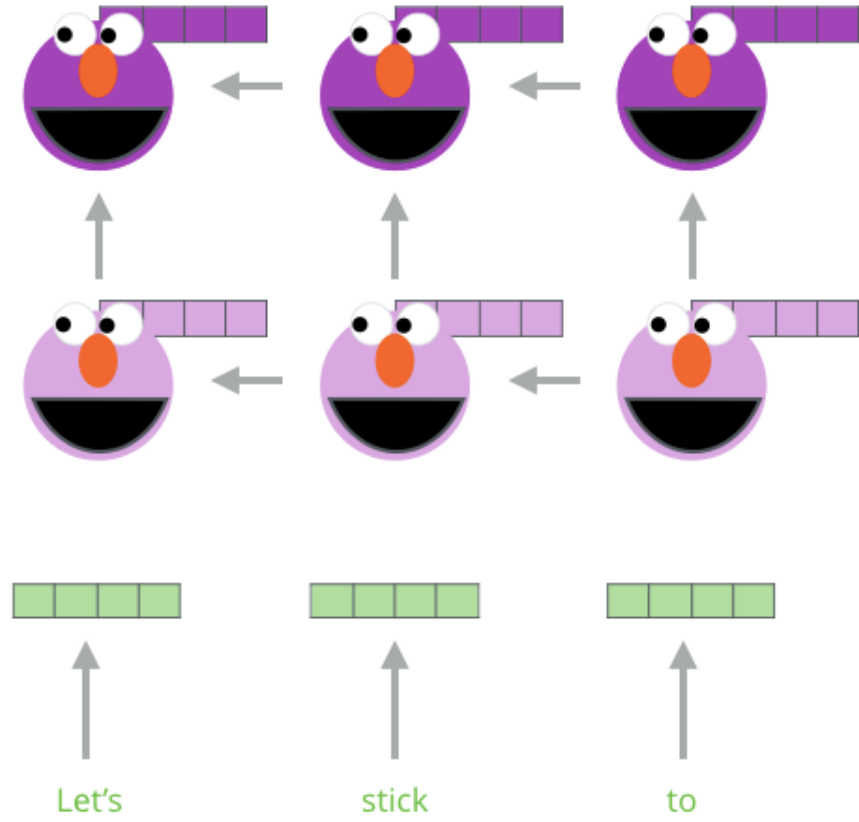
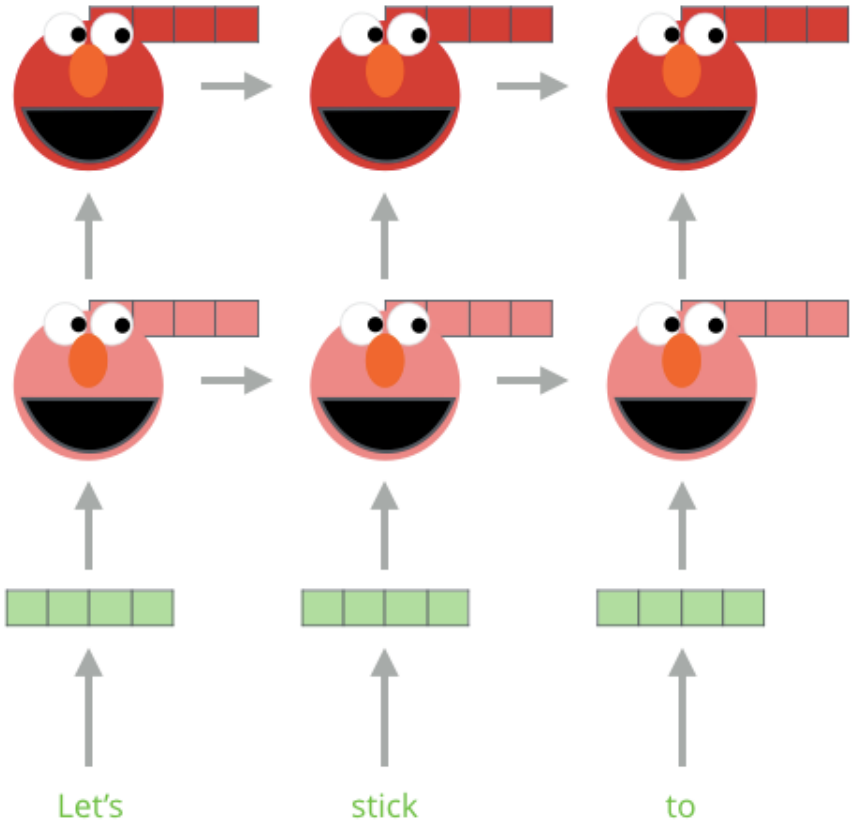
to

# Embedding of "stick" in "Let's stick to" - Step #1

### Forward Language Model

### Backward Language Model

LSTM Layer #2  
LSTM Layer #1  
Embedding



# Embedding of "stick" in "Let's stick to" - Step #2

1- Concatenate hidden layers



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

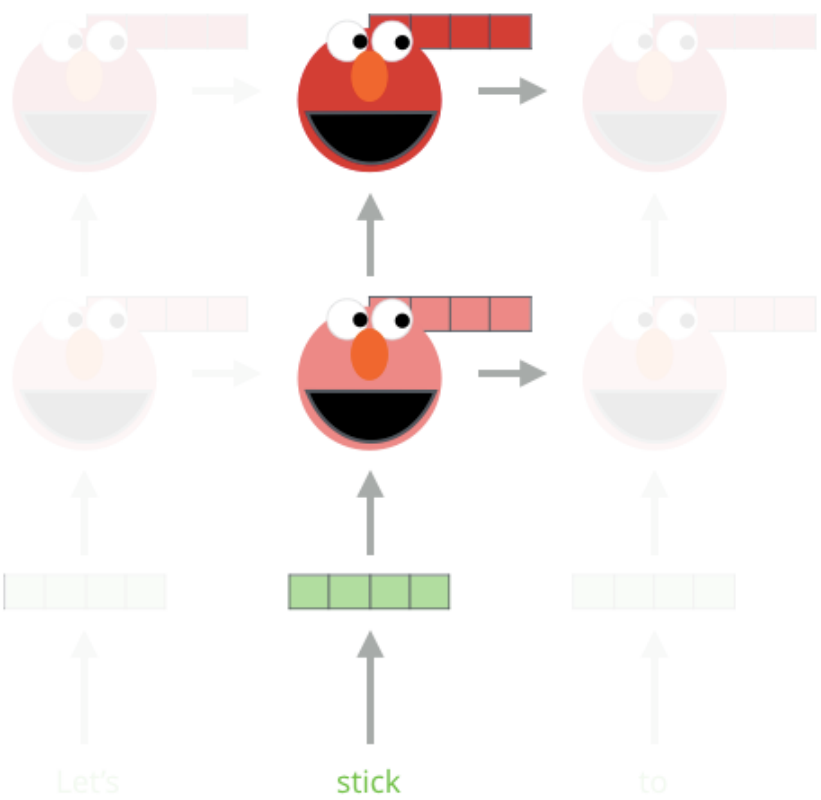


3- Sum the (now weighted) vectors

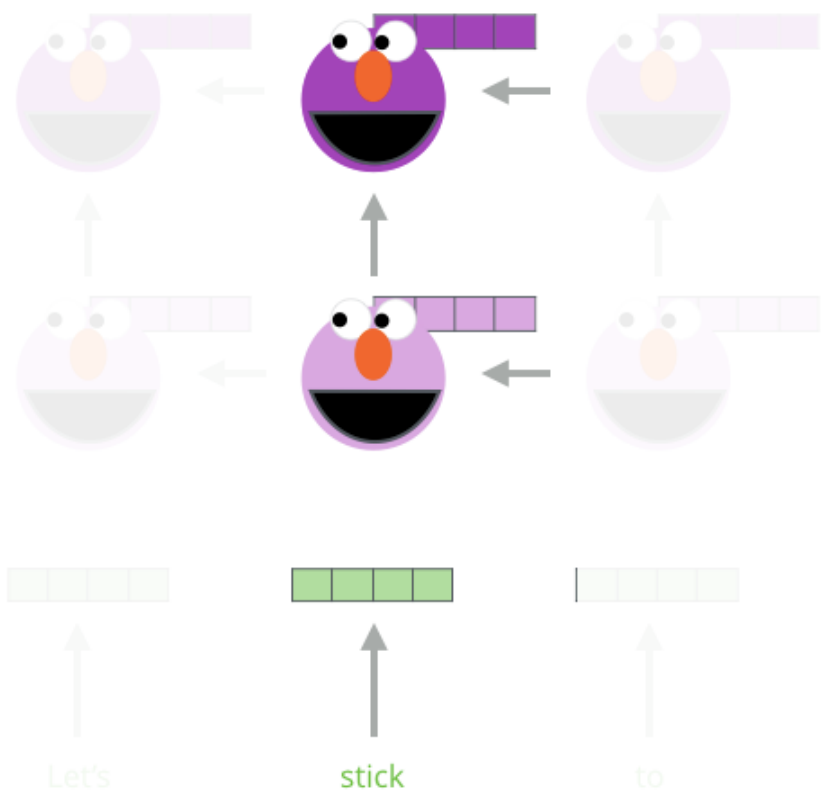


ELMo embedding of "stick" for this task in this context

Forward Language Model



Backward Language Model



# OpenAI Transformer

- Pre-training a Transformer Decoder for Language Modeling

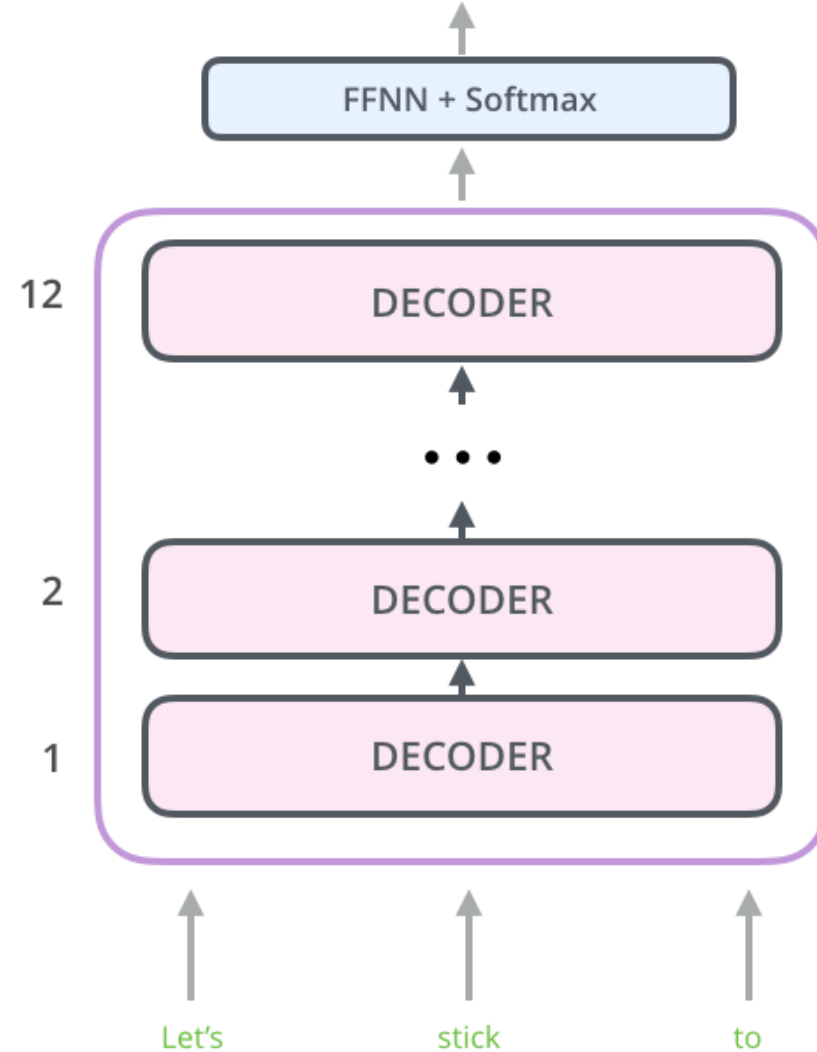


## Re-train Decoder

- No encoder-decoder attention sublayer

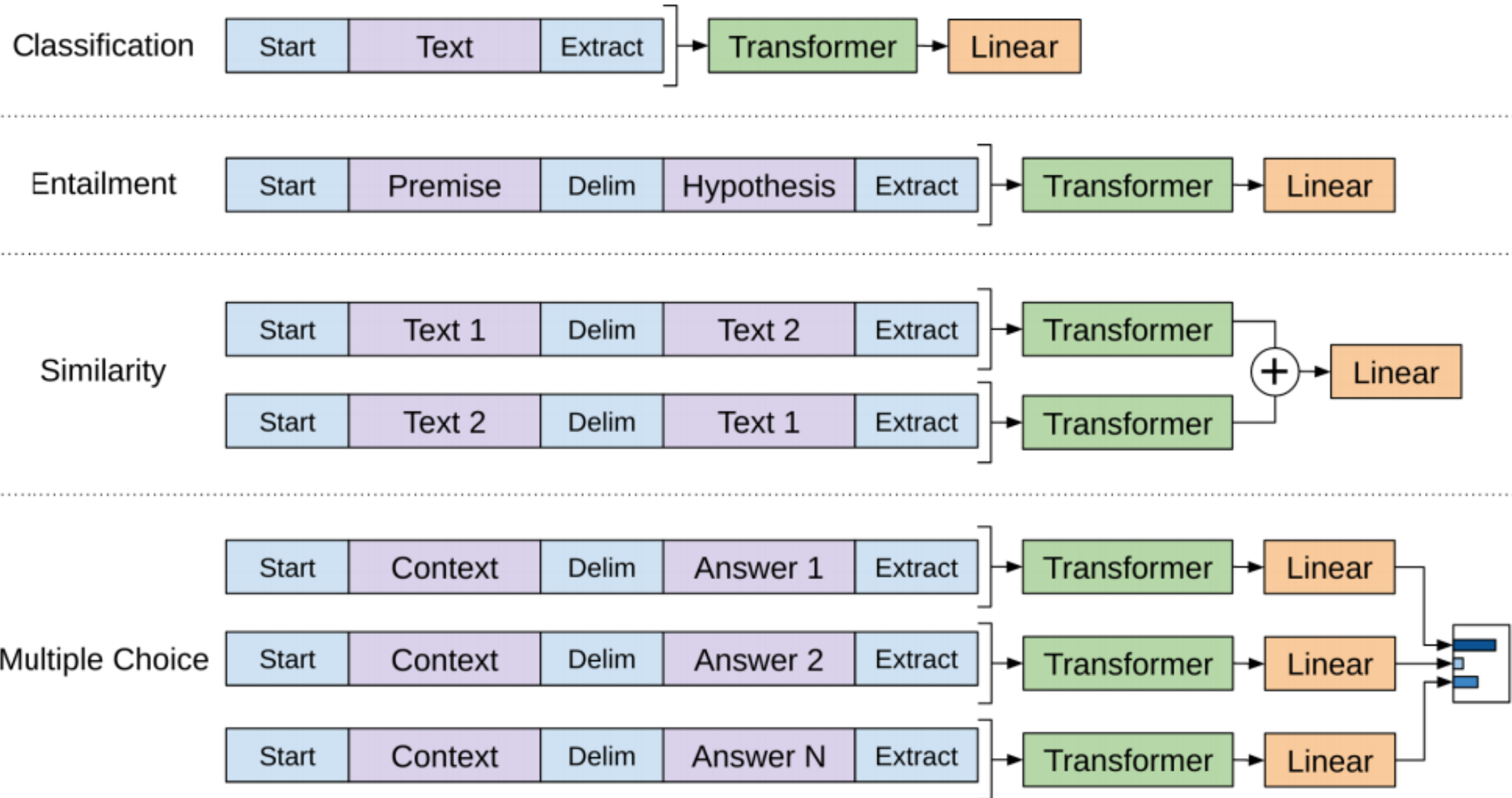
Possible classes:  
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzzzyva





# OpenAI Transformer for Different Tasks



# BERT: From Decoders to Encoders

- ELMo's language model was bi-directional, but the openAI transformer only trains a forward language model.
- Could we build a transformer-based model whose language model looks both forward and backwards, i.e. "is conditioned on both left and right context"?
  - **Masked Language Model**

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

Possible classes:  
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzzzyva

FFNN + Softmax



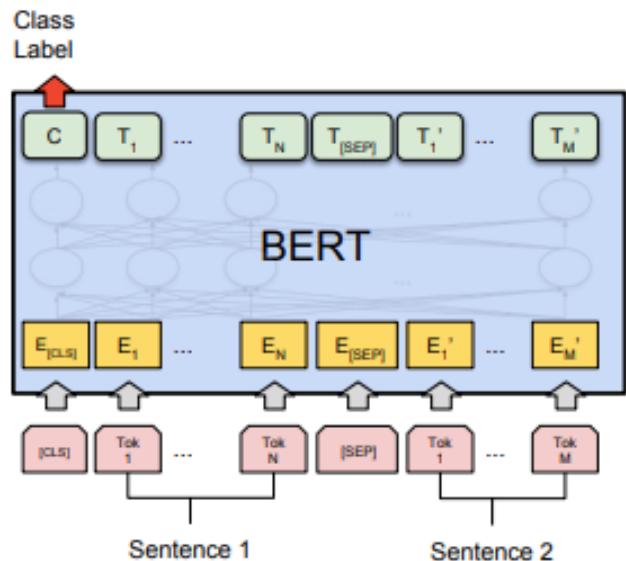
Randomly mask 15% of tokens

1 [CLS] 2 Let's 3 stick 4 to 5 [MASK] 6 in 7 this 8 skit ... 512

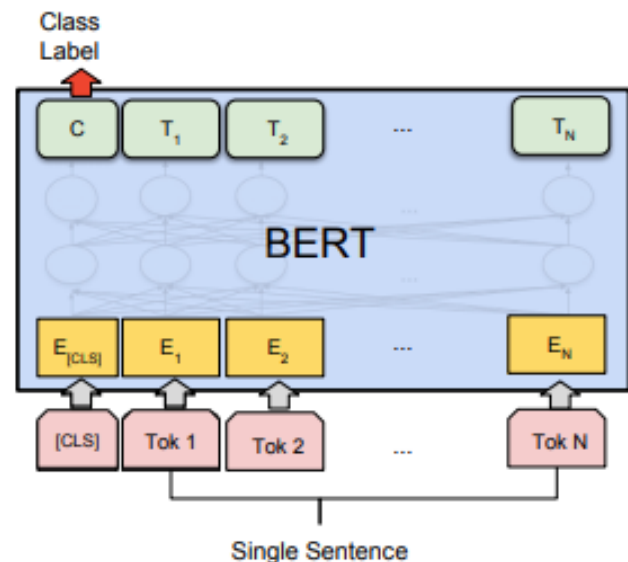
Input

[CLS] Let's stick to improvisation in this skit

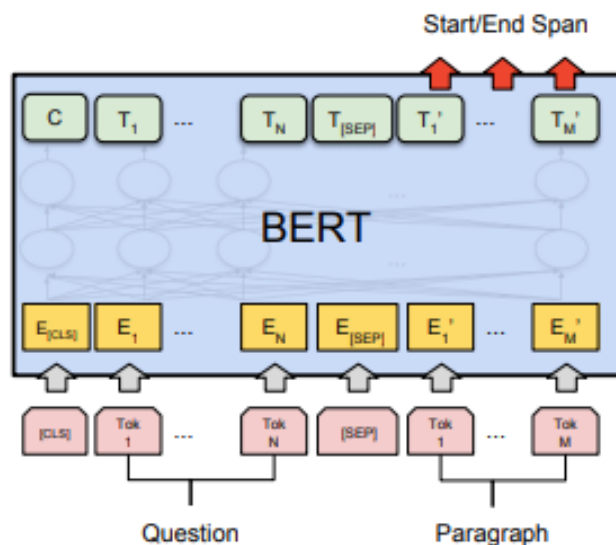
# BERT for Different tasks



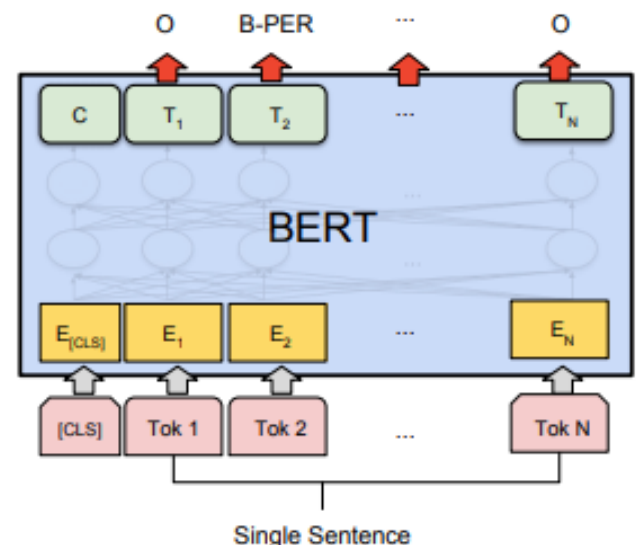
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



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



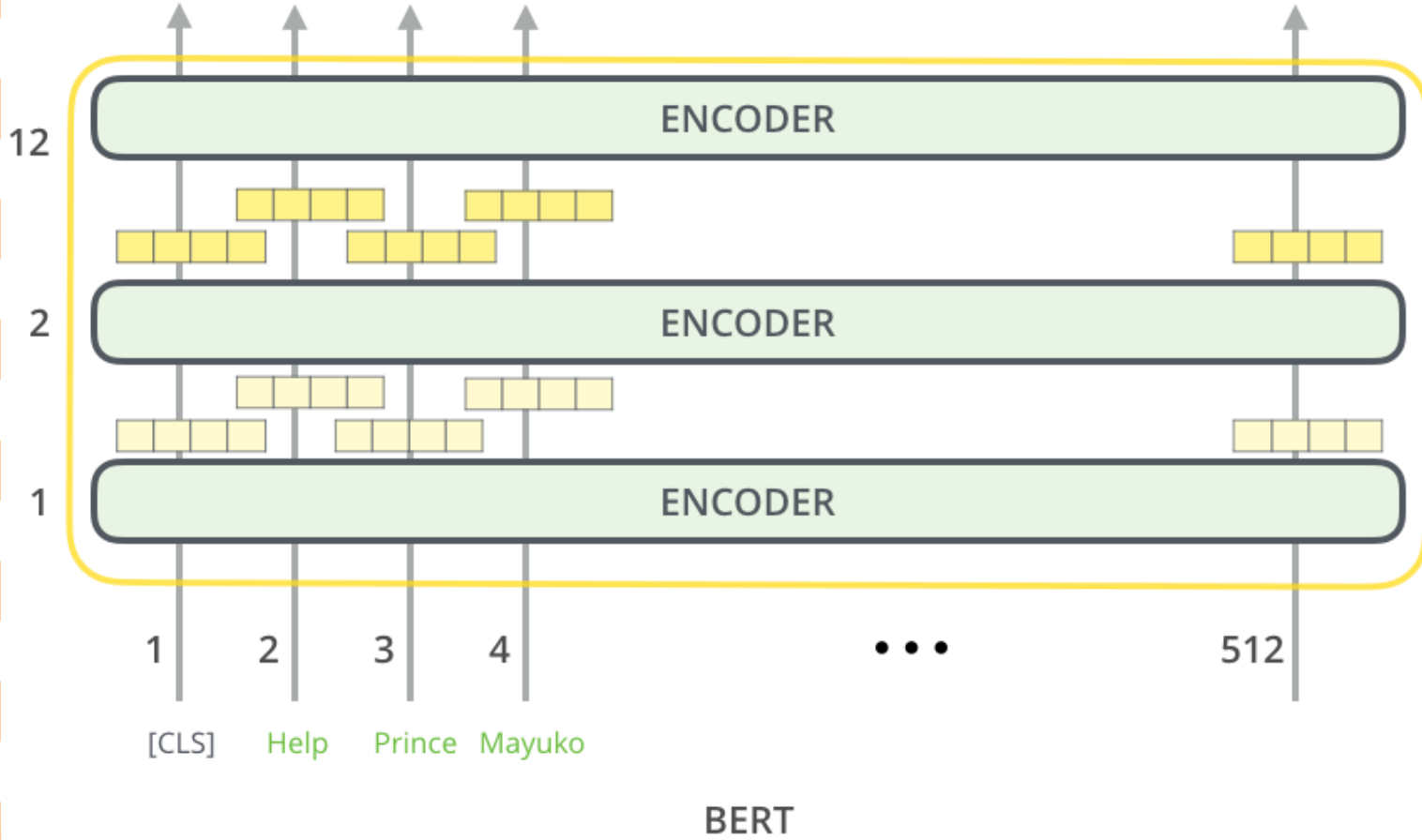
(c) Question Answering Tasks:  
SQuAD v1.1



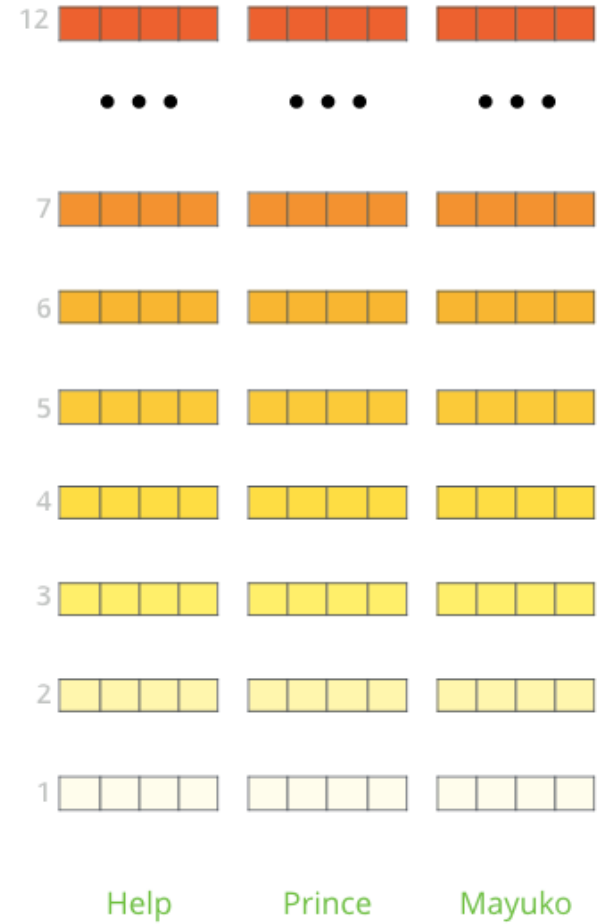
(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

Generate Contextualized Embeddings

# BERT for Feature Extraction



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

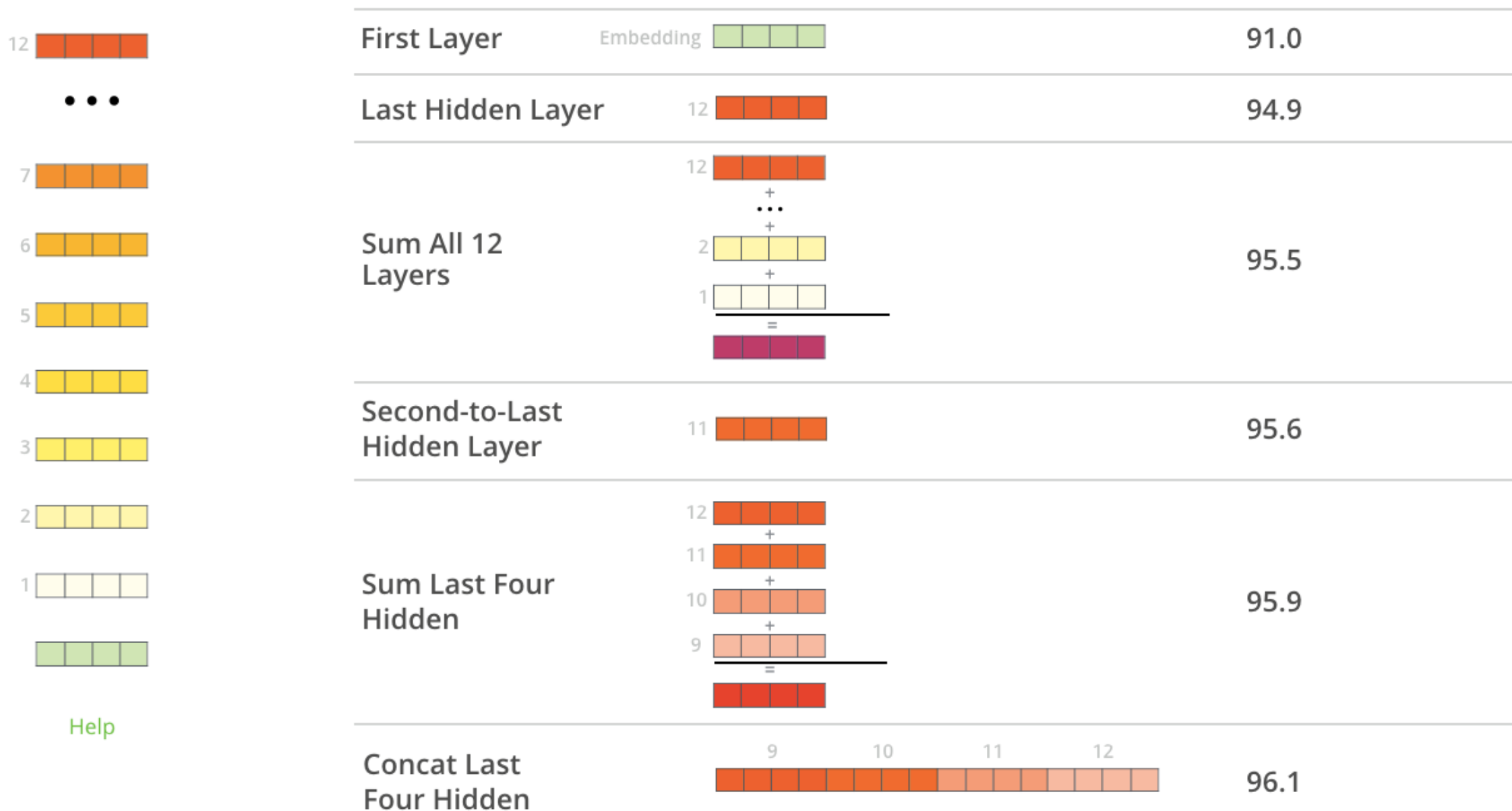


But which one should we use?

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

For named-entity recognition task CoNLL-2003 NER

Dev F1 Score



# References

1. <https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html>
2. <http://jalammarm.github.io/illustrated-transformer/>
3. <http://jalammarm.github.io/illustrated-bert/>
4. Hong-Yi Lee, Transformer, 2019, <https://www.youtube.com/watch?v=ugWDIIOHtPA>