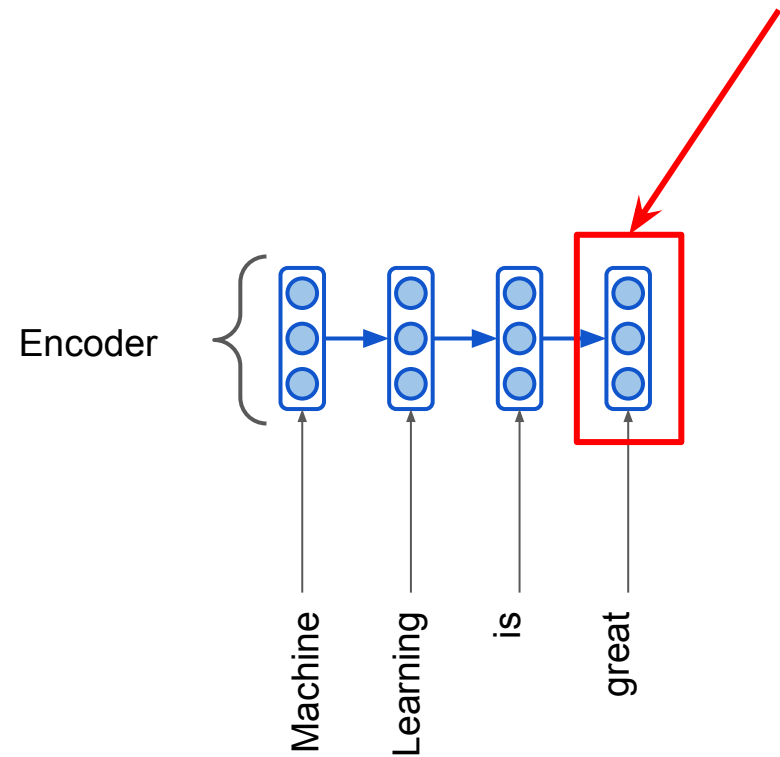
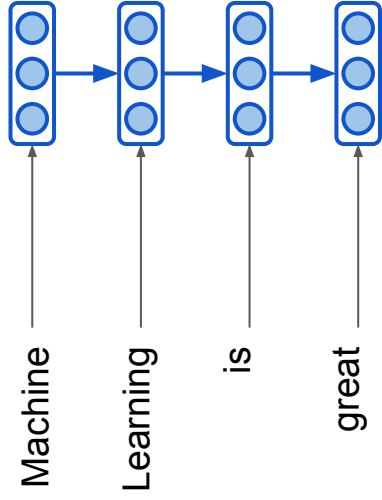


# Seq2seq NMT

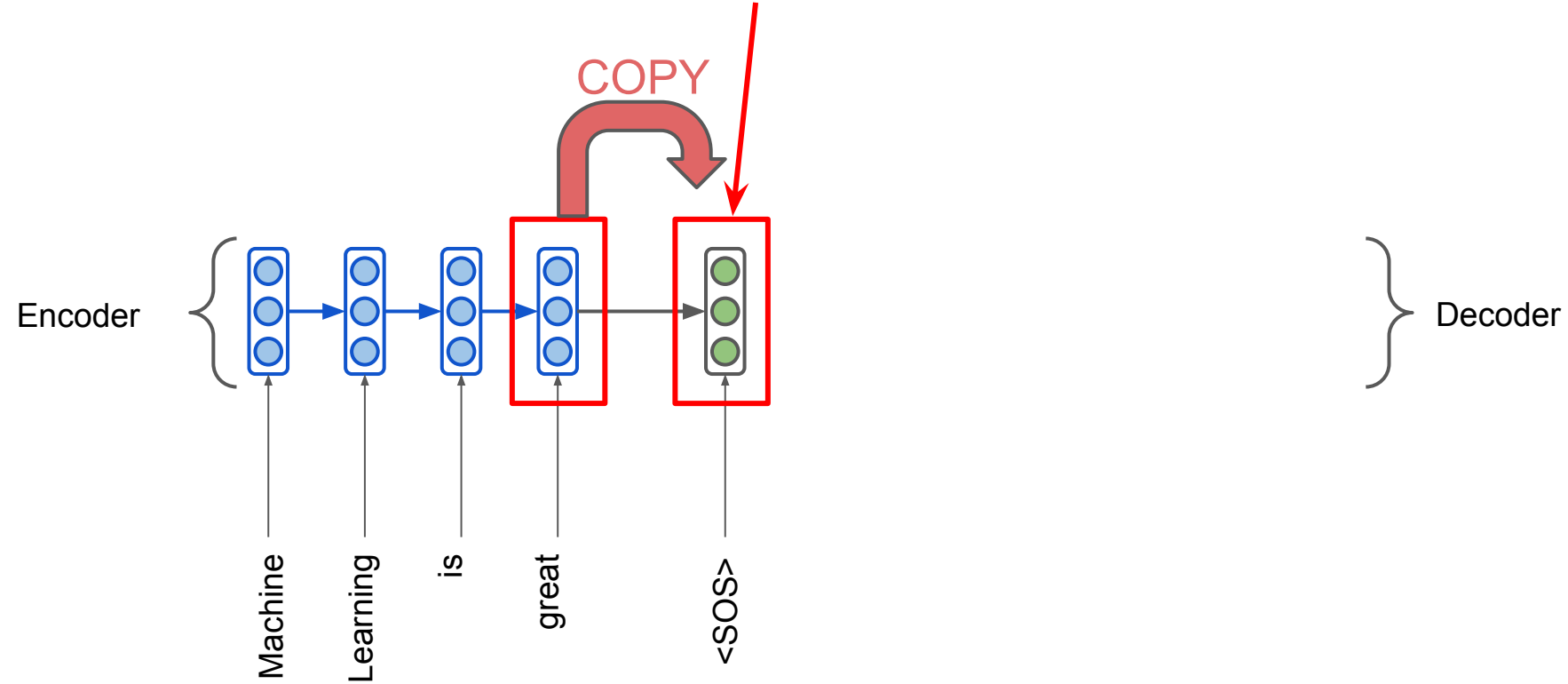
This state encodes  
the whole sentence



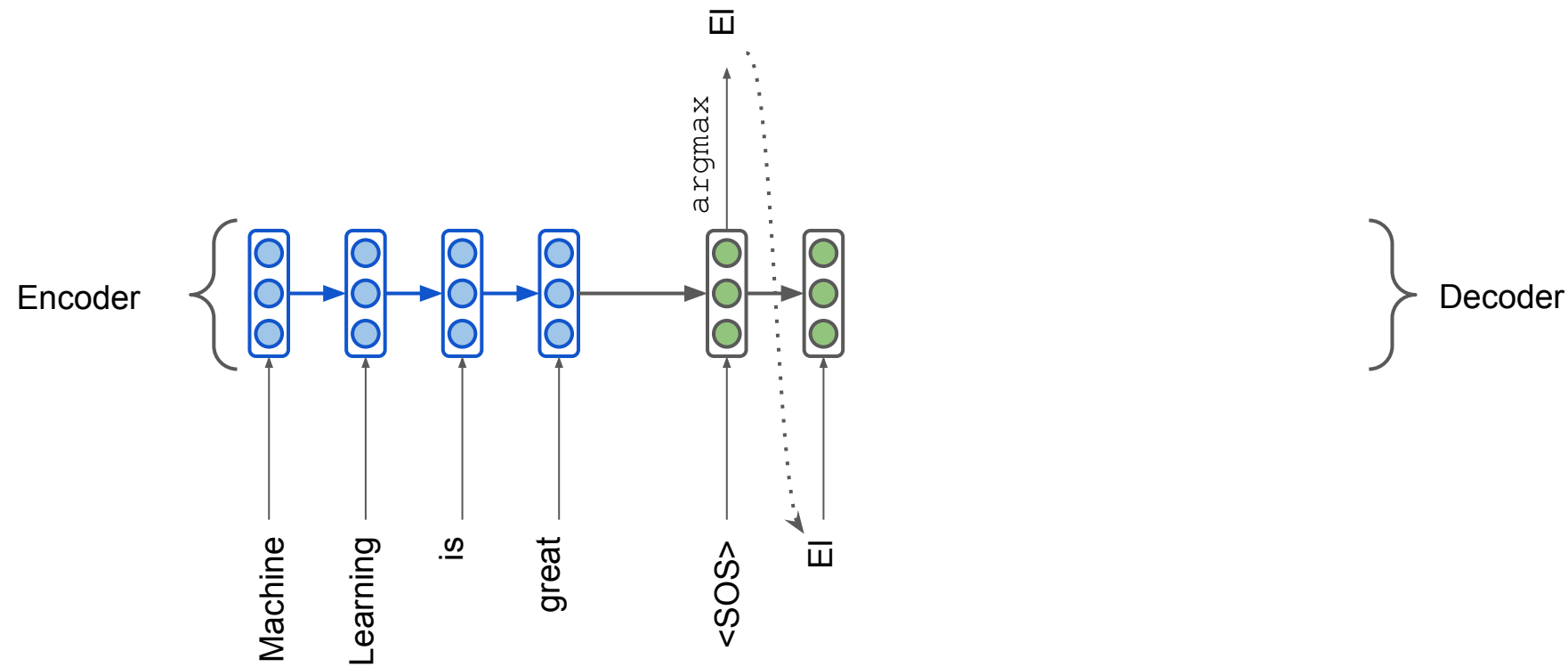


# Seq2seq NMT

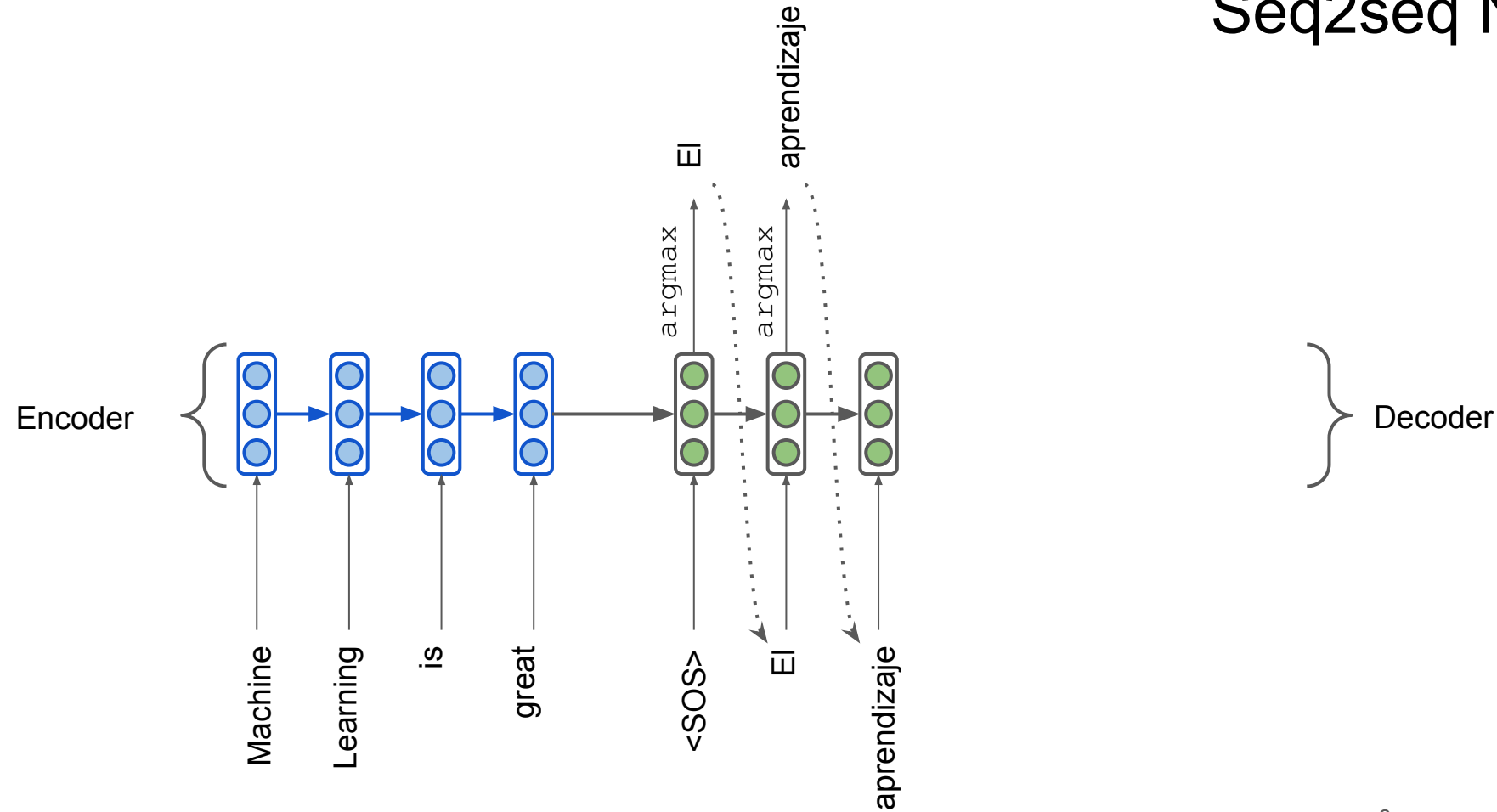
Forwarded as initial  
hidden state to decoder



# Seq2seq NMT

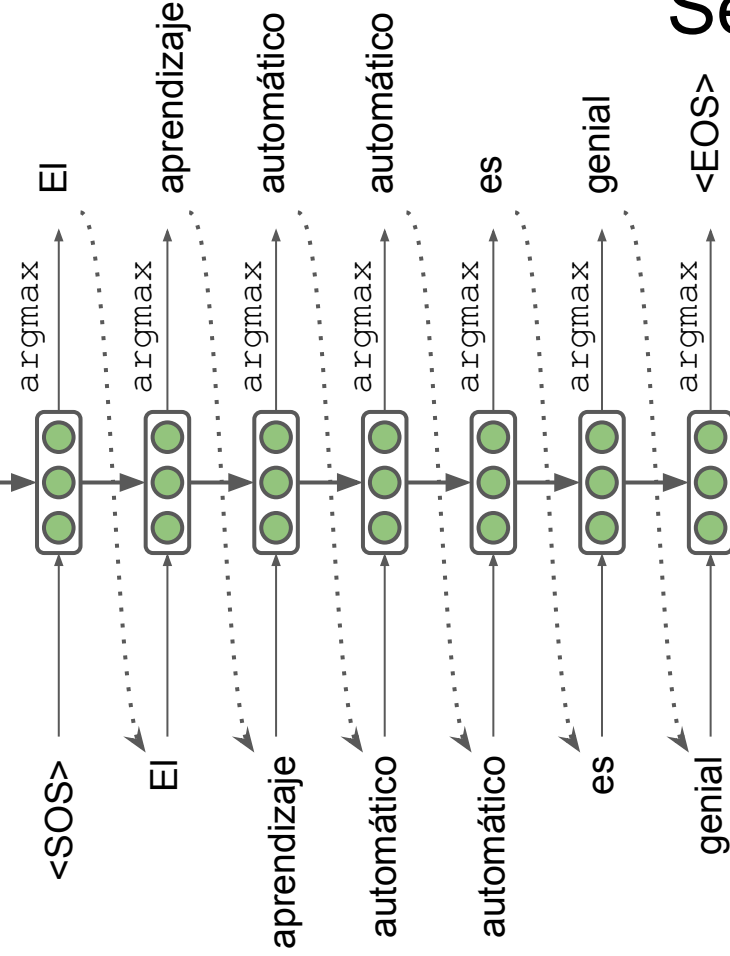
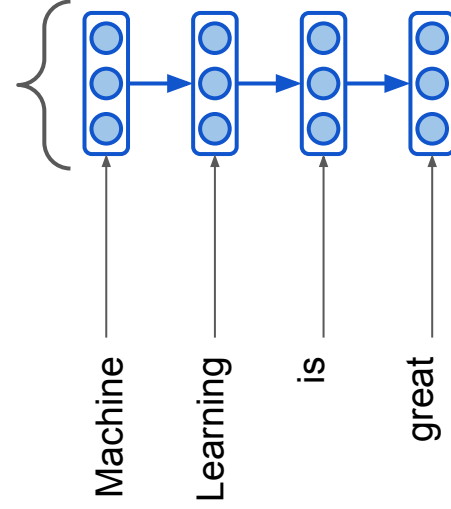


# Seq2seq NMT



# Seq2seq NMT

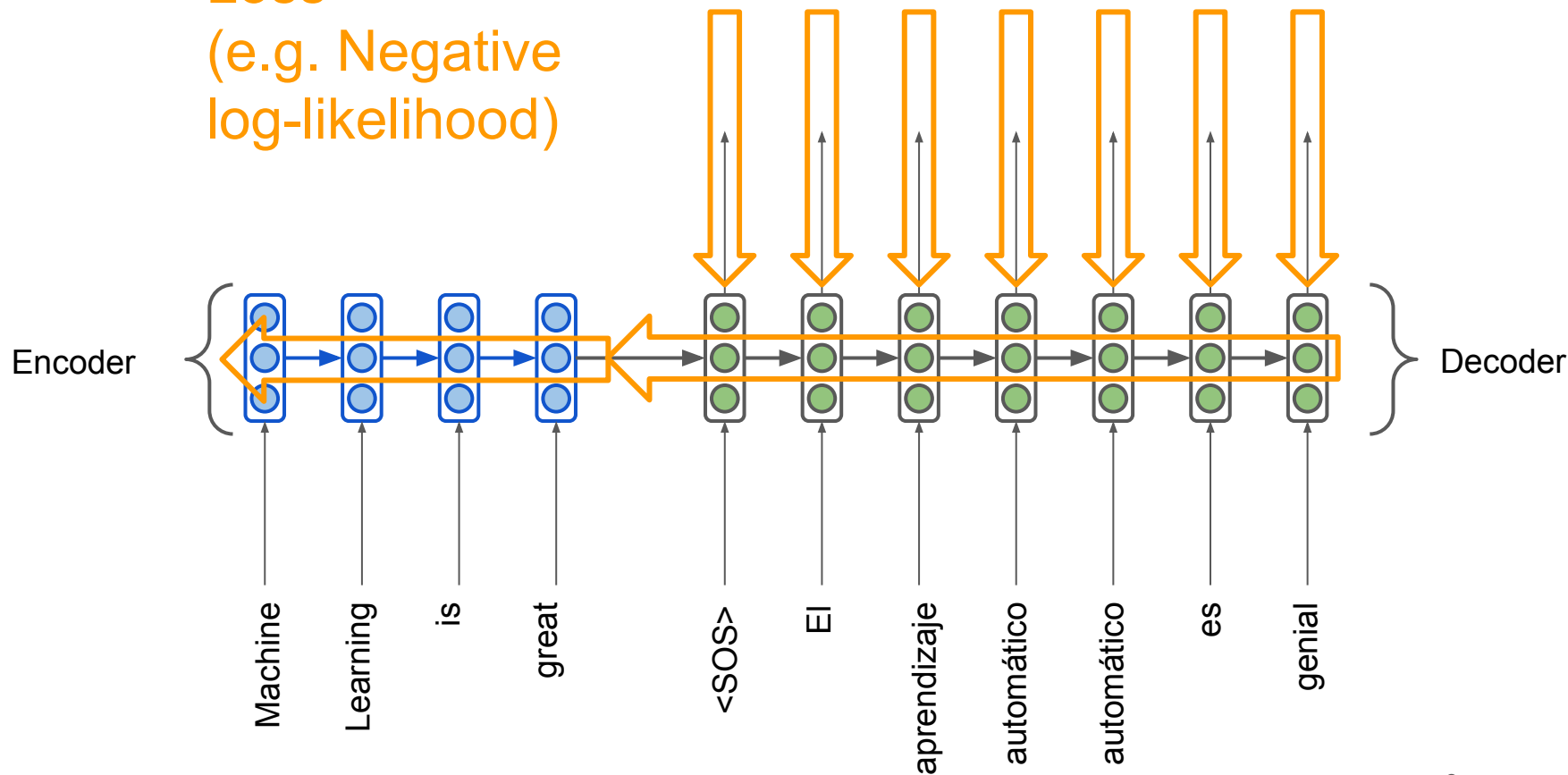
Encoder



Decoder

# Seq2seq is trained end-to-end

Loss  
(e.g. Negative  
log-likelihood)

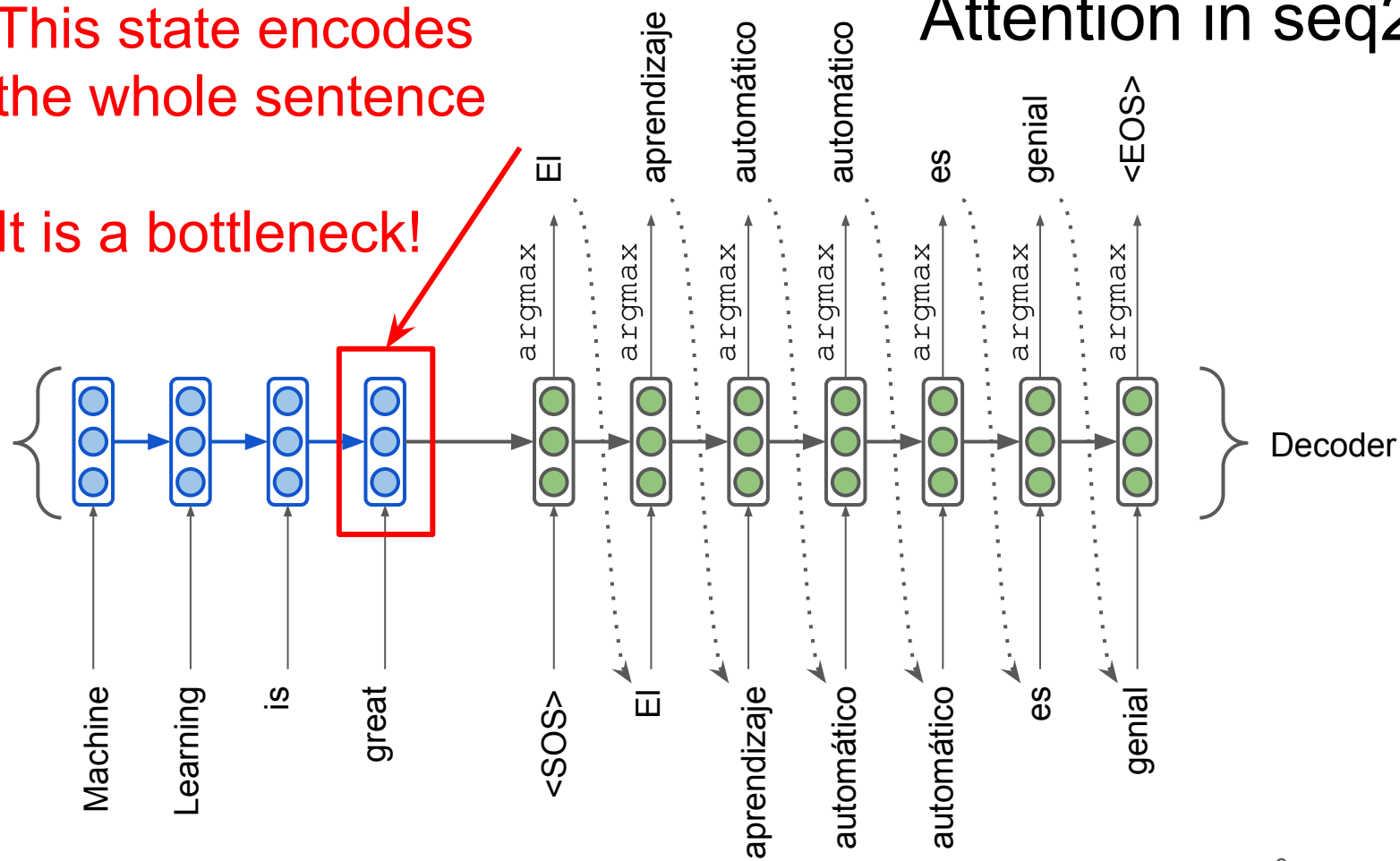


# Attention in seq2seq

This state encodes the whole sentence

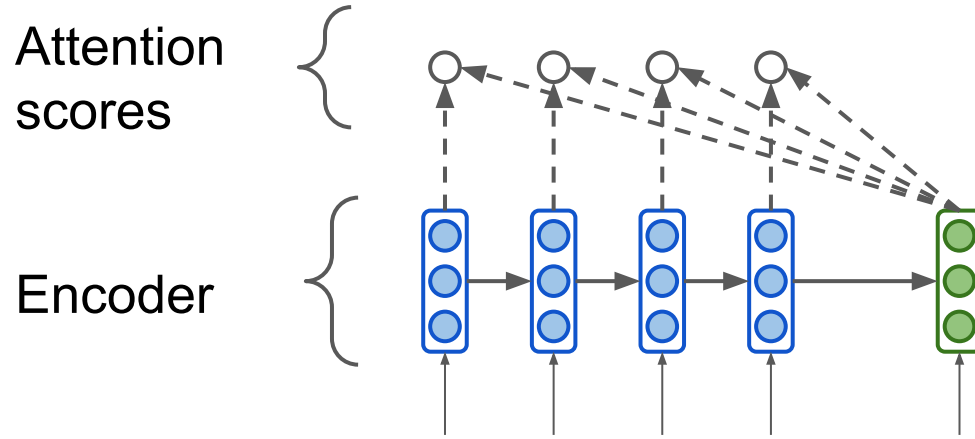
It is a bottleneck!

Encoder

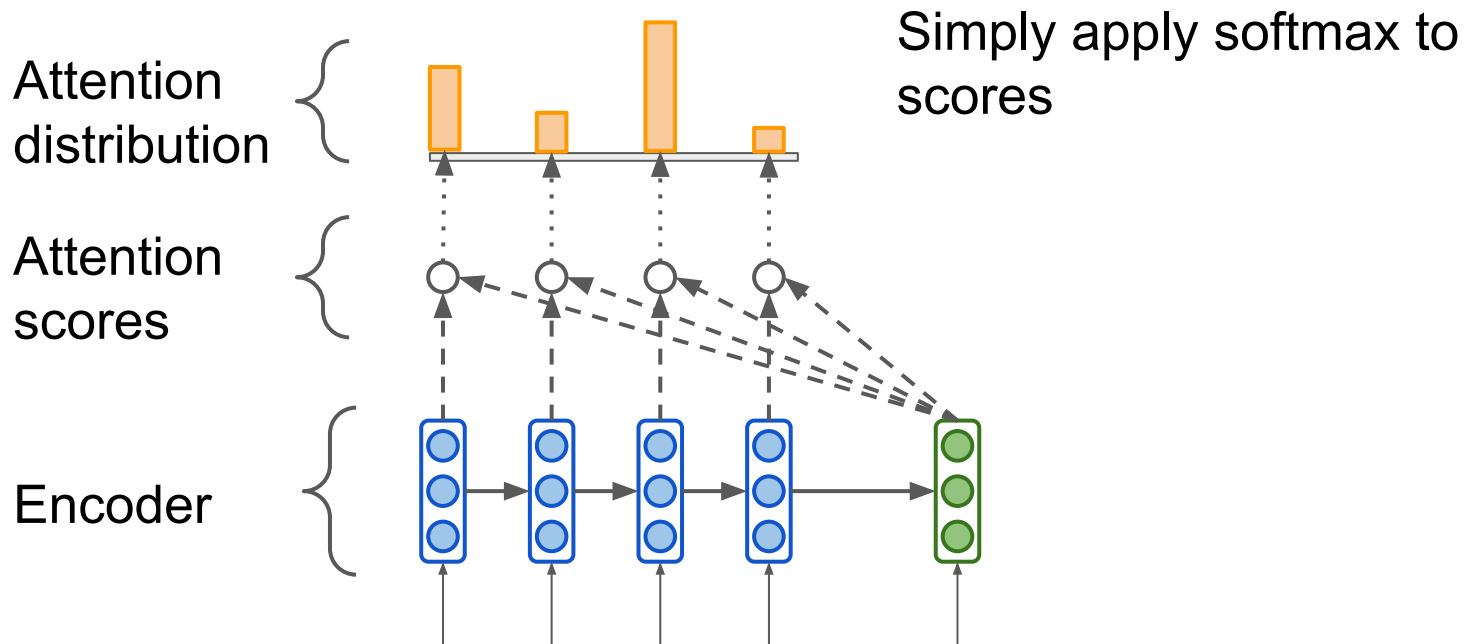




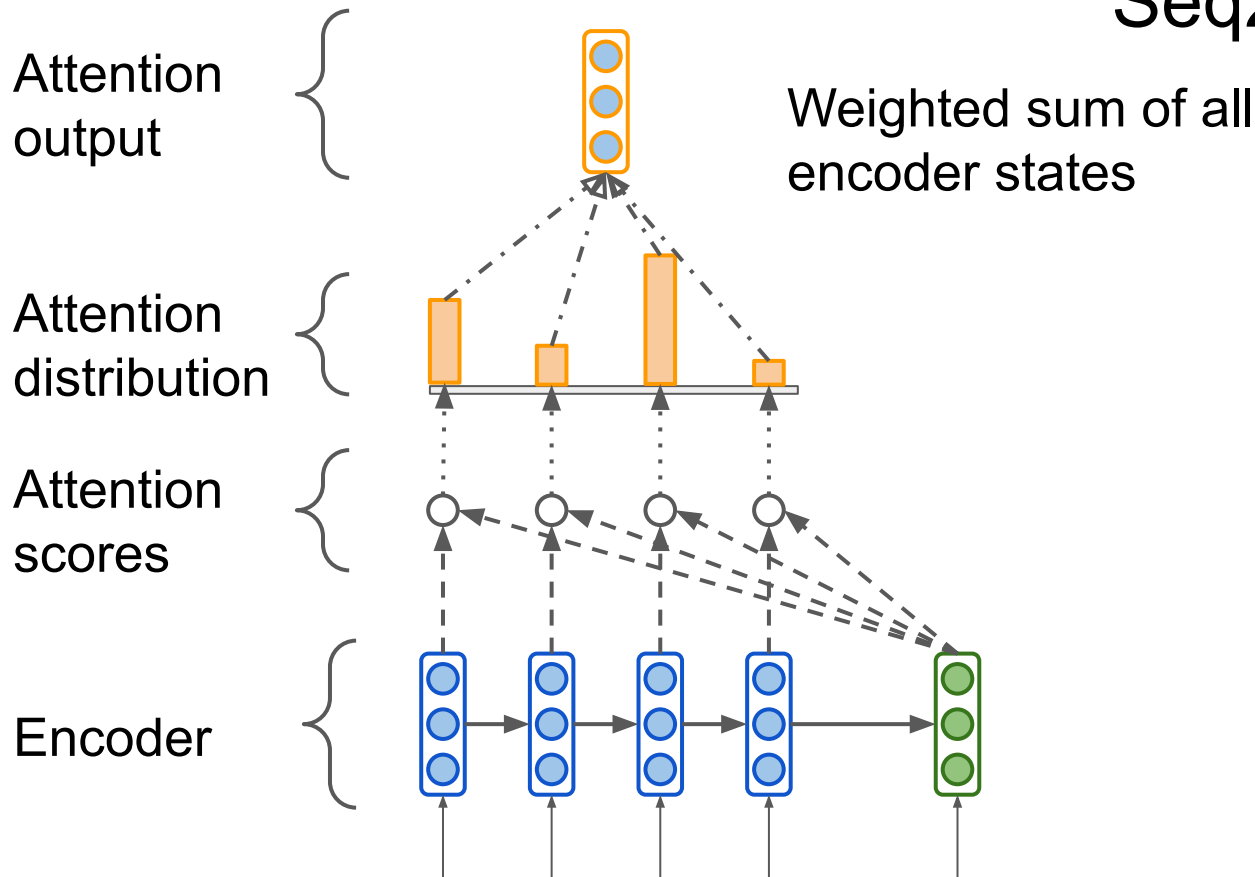
# Seq2seq with attention



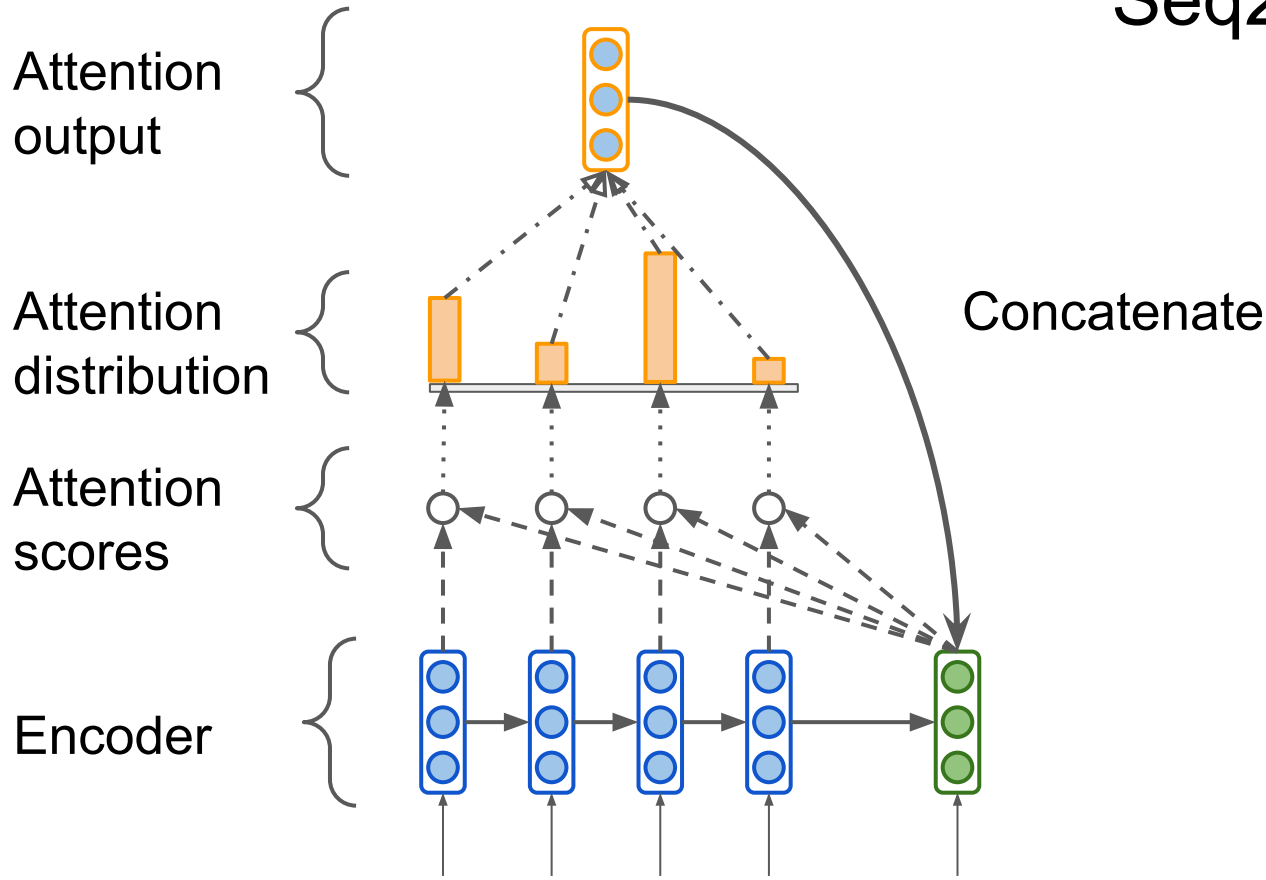
# Seq2seq with attention



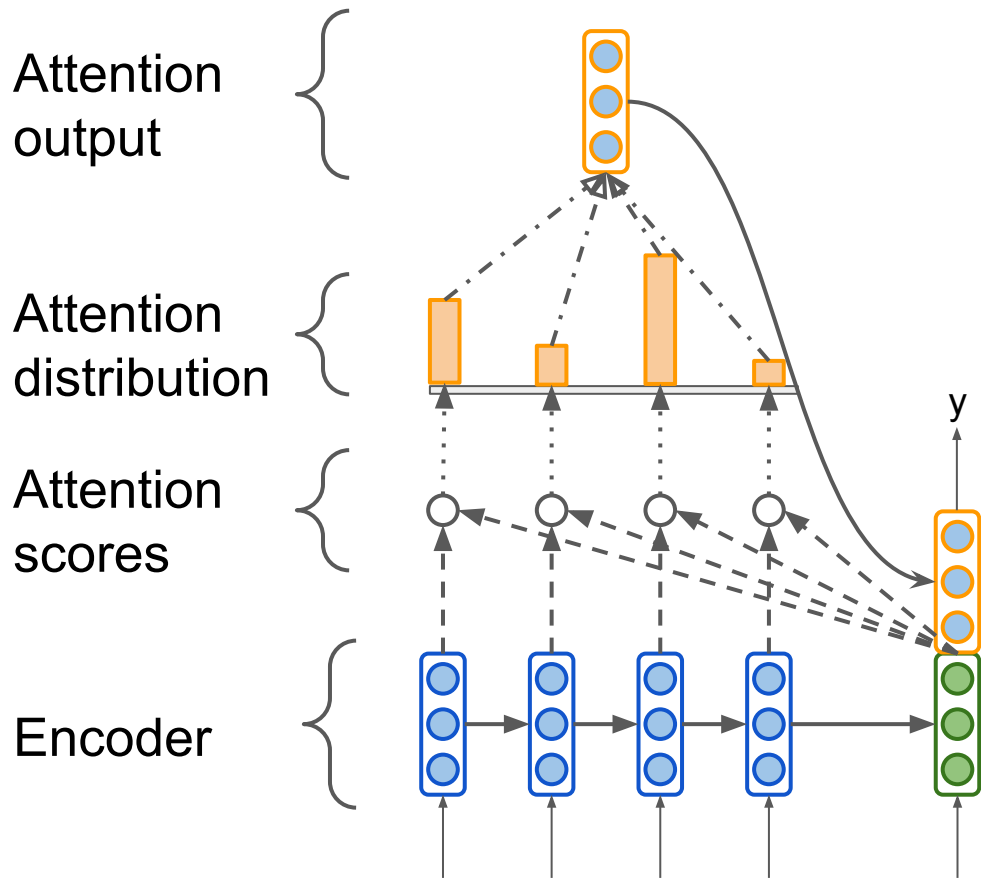
# Seq2seq with attention



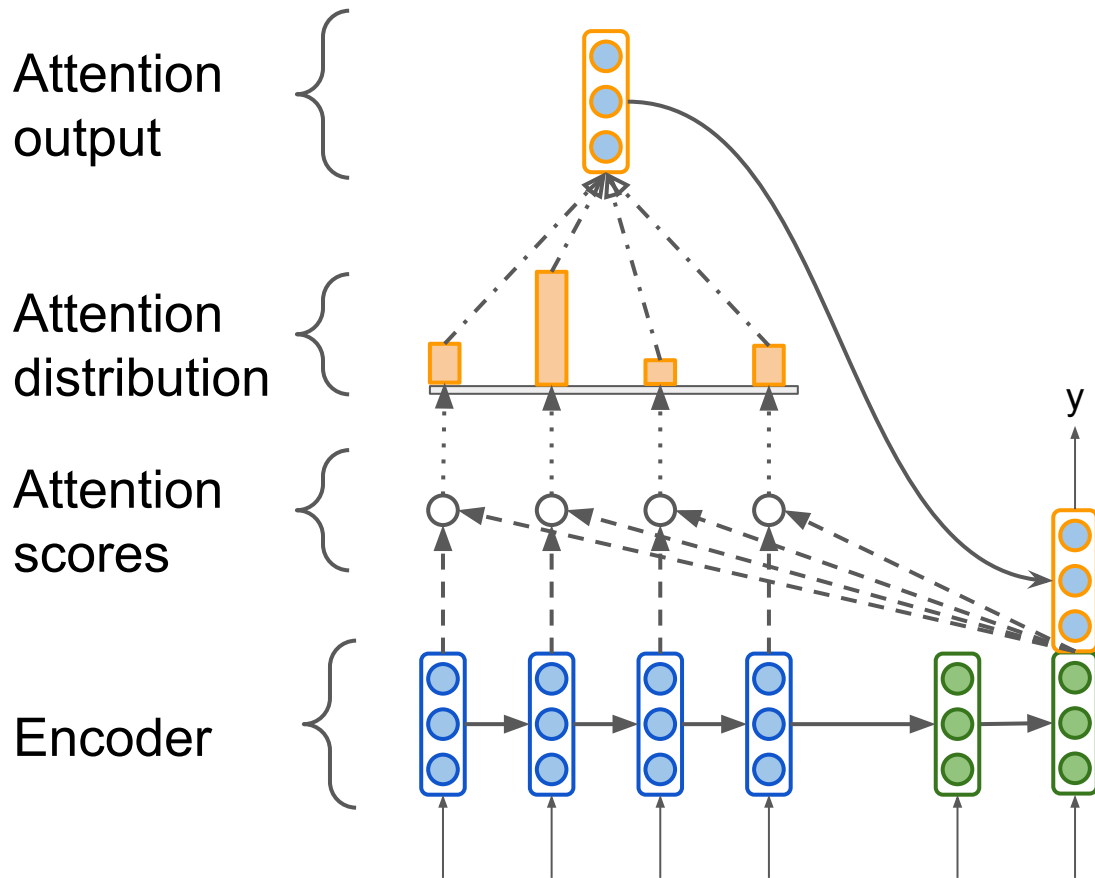
# Seq2seq with attention



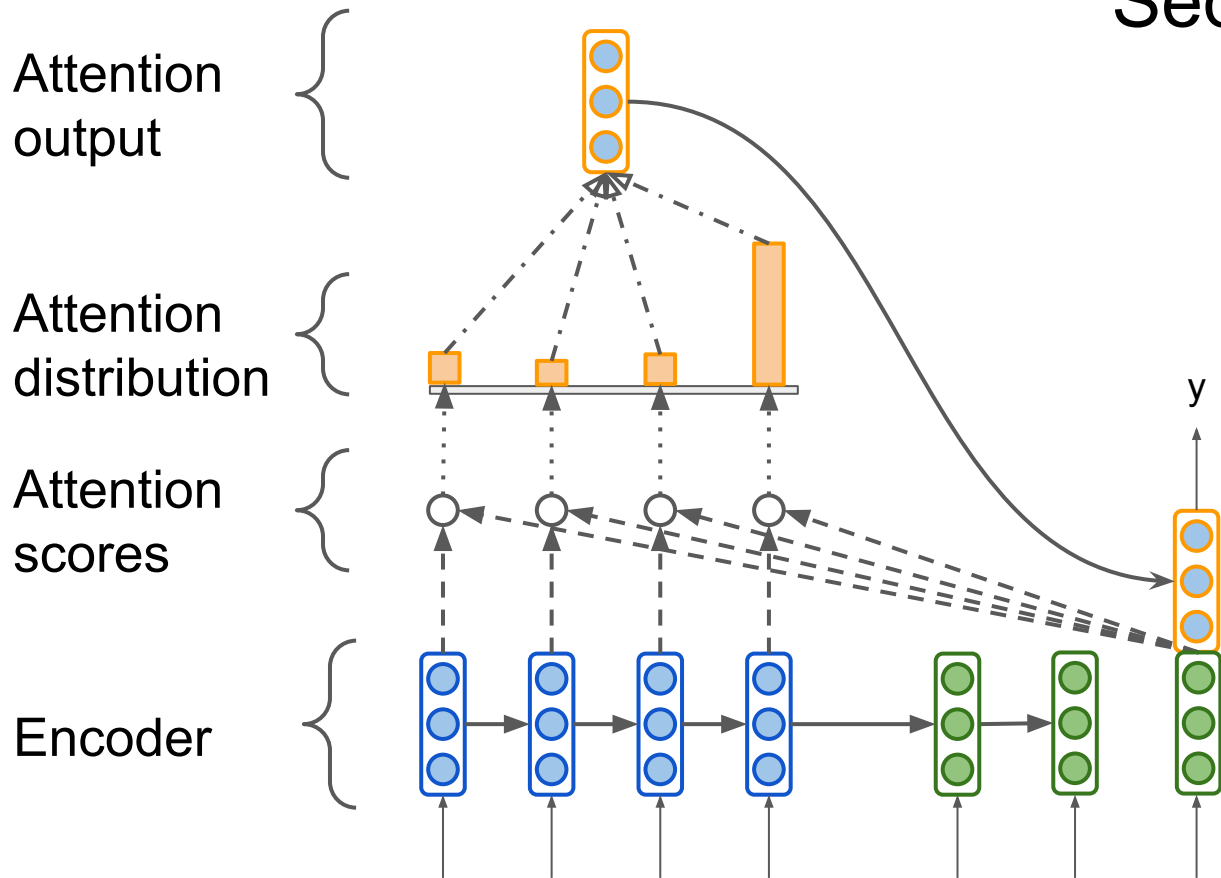
# Seq2seq with attention



# Seq2seq with attention



# Seq2seq with attention



## Attention in equations

Denote encoder hidden states  $\mathbf{h}_1, \dots, \mathbf{h}_N \in \mathbb{R}^k$   
and decoder hidden state at time step  $t$   $\mathbf{s}_t \in \mathbb{R}^k$

The attention scores  $\mathbf{e}^t$  can be computed as dot product

$$\mathbf{e}^t = [\mathbf{s}^T \mathbf{h}_1, \dots, \mathbf{s}^T \mathbf{h}_N]$$

Then the attention vector is a linear combination of encoder states

$$\mathbf{a}_t = \sum_{i=1}^N \alpha_i^t \mathbf{h}_i \in \mathbb{R}^k, \text{ where } \boldsymbol{\alpha}_t = \text{softmax}(\mathbf{e}_t)$$

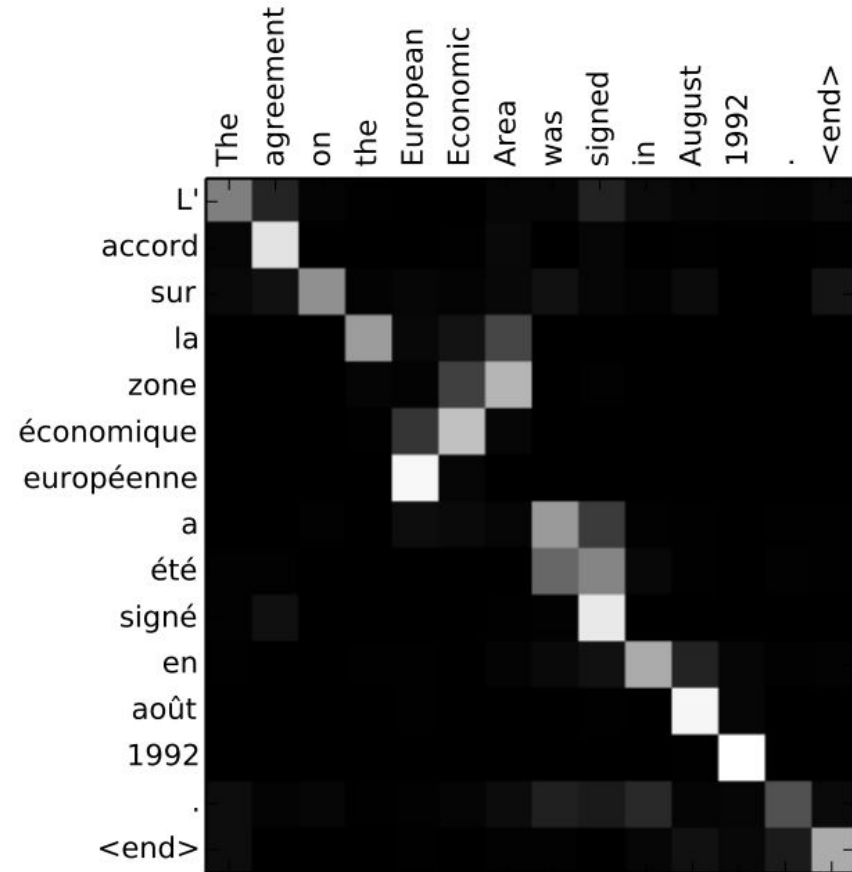
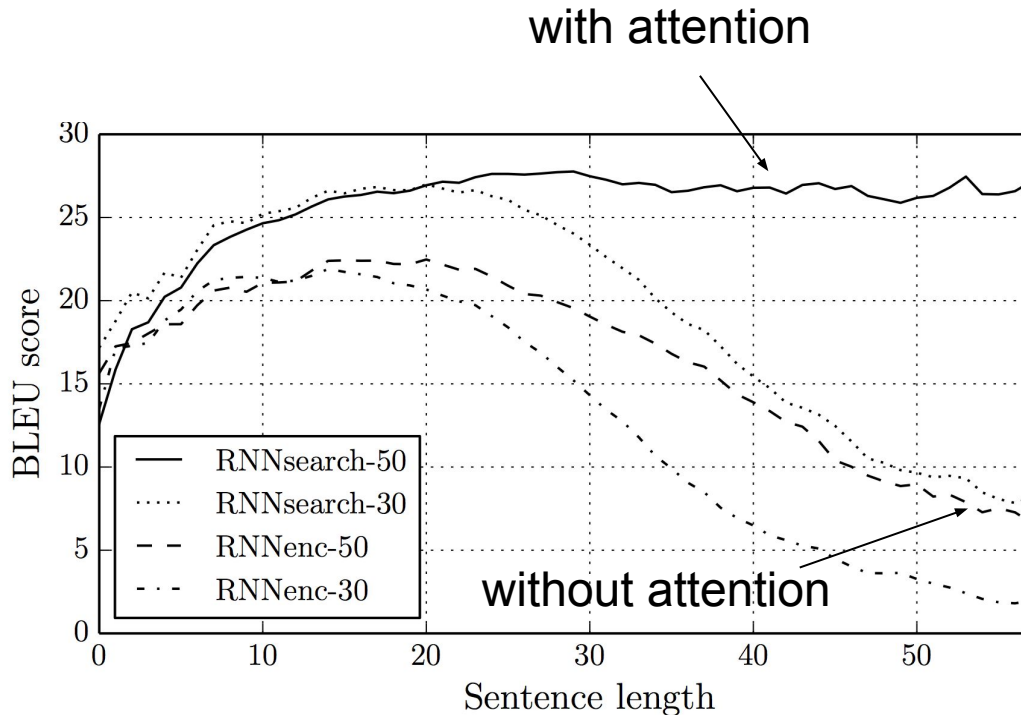


# Attention variants

- Basic dot-product (the one discussed before):  $e_i = \mathbf{s}^T \mathbf{h}_i \in \mathbb{R}$
- Multiplicative attention:  $e_i = \mathbf{s}^T \mathbf{W} \mathbf{h}_i \in \mathbb{R}$ 
  - $\mathbf{W} \in \mathbb{R}^{d_2 \times d_1}$  - weight matrix
- Additive attention:  $e_i = \mathbf{v}^T \tanh(\mathbf{W}_1 \mathbf{h}_i + \mathbf{W}_2 \mathbf{s}) \in \mathbb{R}$ 
  - $\mathbf{W}_1 \in \mathbb{R}^{d_3 \times d_1}, \mathbf{W}_2 \in \mathbb{R}^{d_3 \times d_2}$  - weight matrices
  - $\mathbf{v} \in \mathbb{R}^{d_3}$  - weight vector

# Attention advantages

- “Free” word alignment
- Better results on long sequences

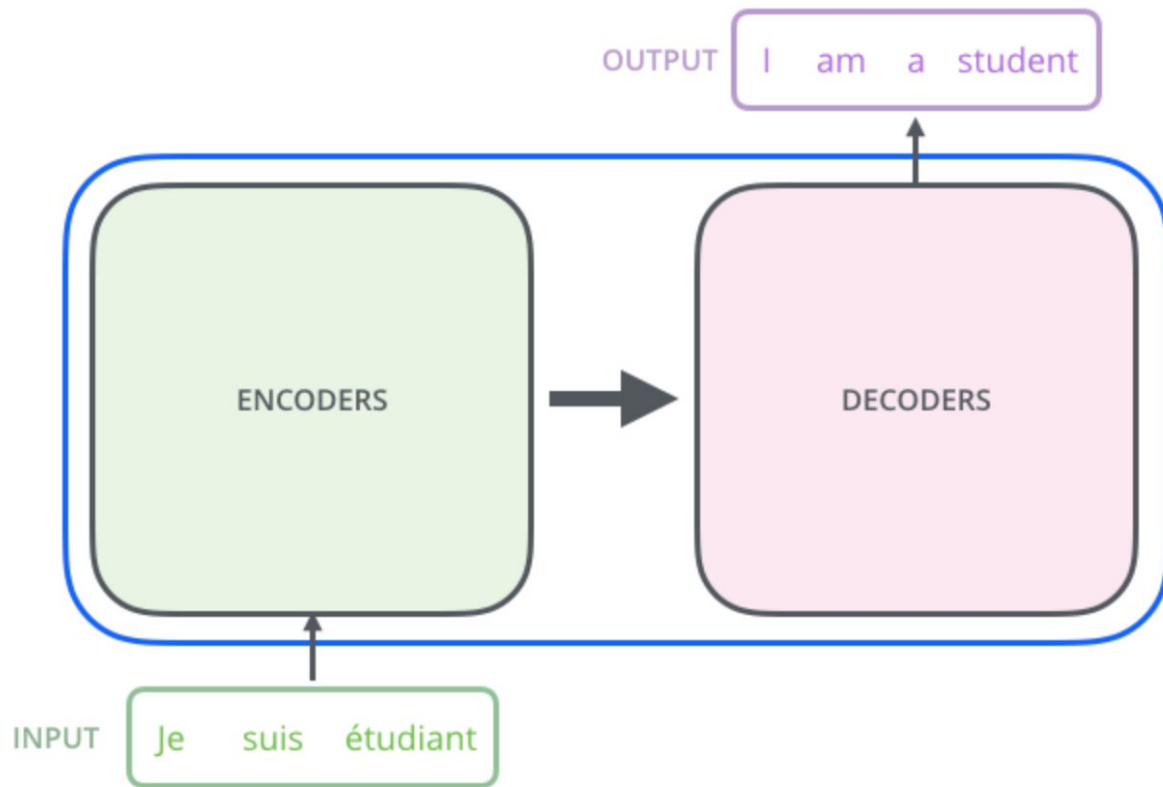


# The Transformer

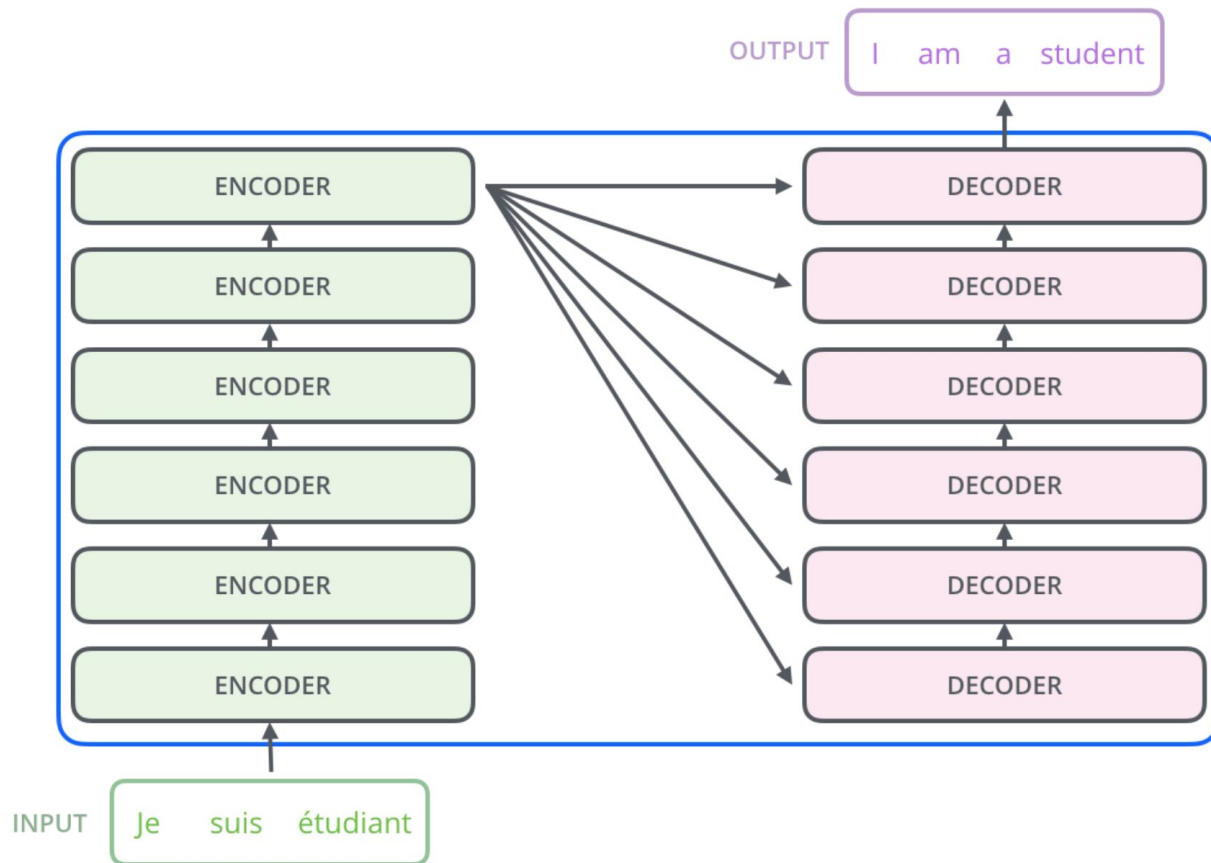
# The Transformer



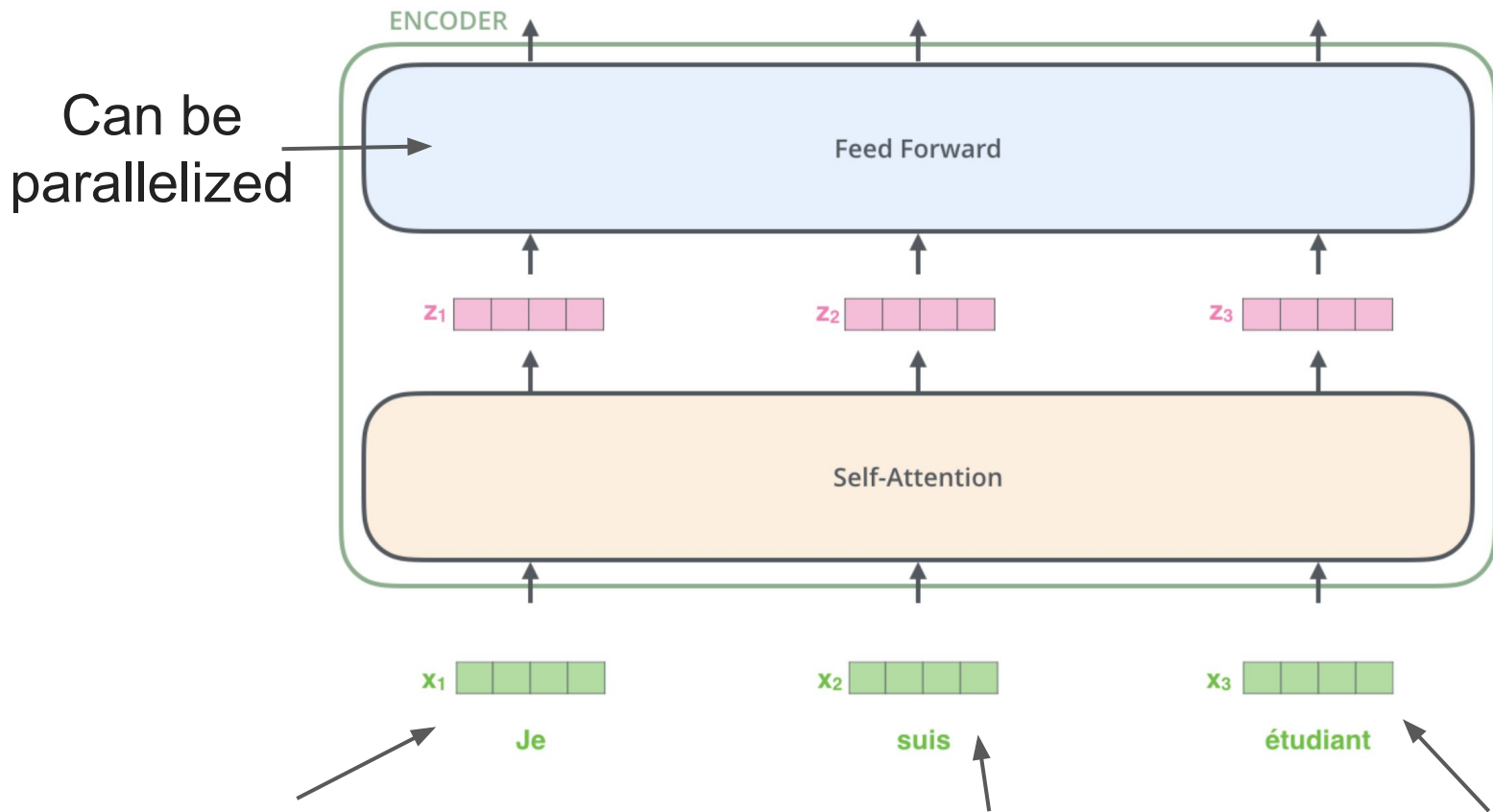
# The Transformer



# The Transformer



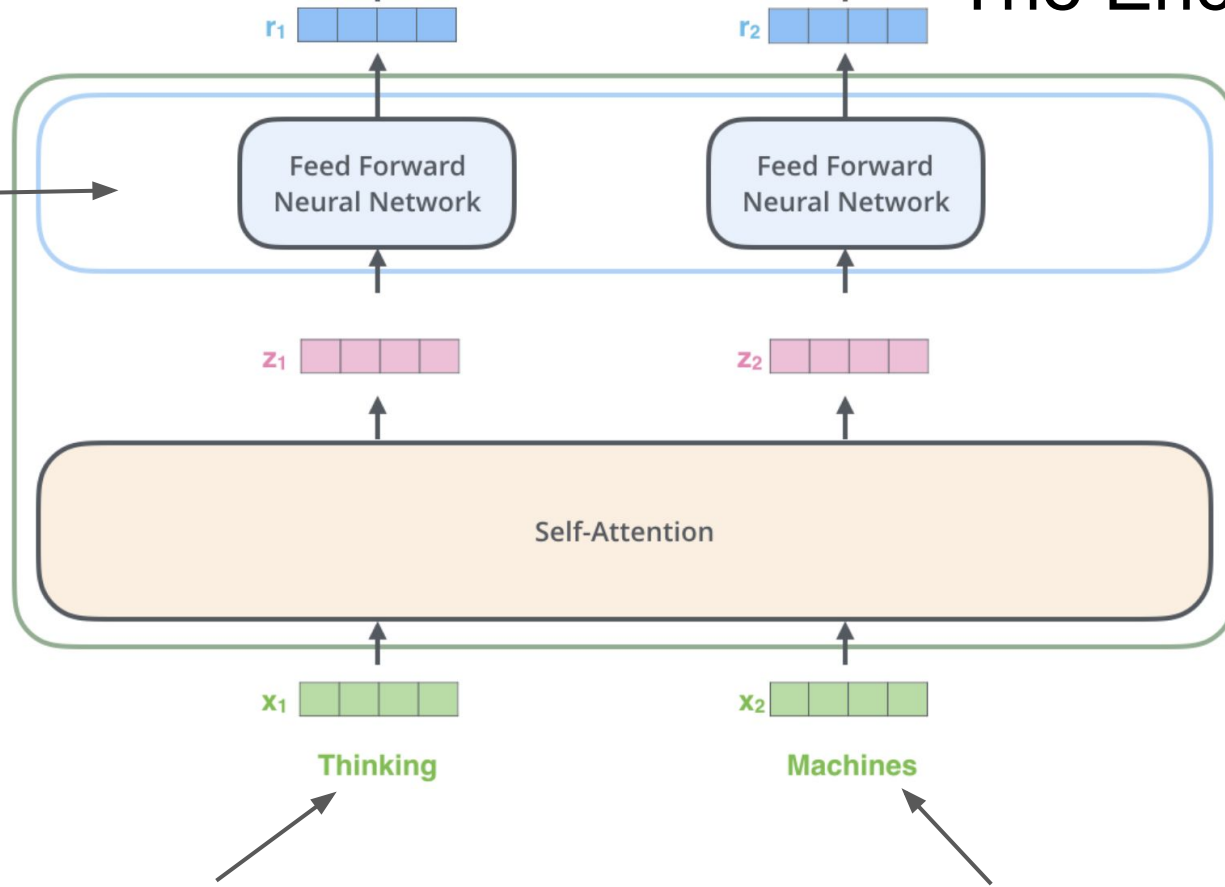
# The Encoder Side



the word in each position flows through its own path in the encoder

# The Encoder Side

Can be  
parallelized



the word in each position flows through its own path in the encoder



# The Transformer: quick overview

- Proposed in 2017 in paper [Attention is All You Need](#) by Ashish Vaswani et al.
- No recurrent or convolutional layers, only attention
- Beats seq2seq in machine translation task
  - *28.4 BLEU on the WMT 2014 English-to-German translation task*
- Much faster
- Uses **self-attention** concept

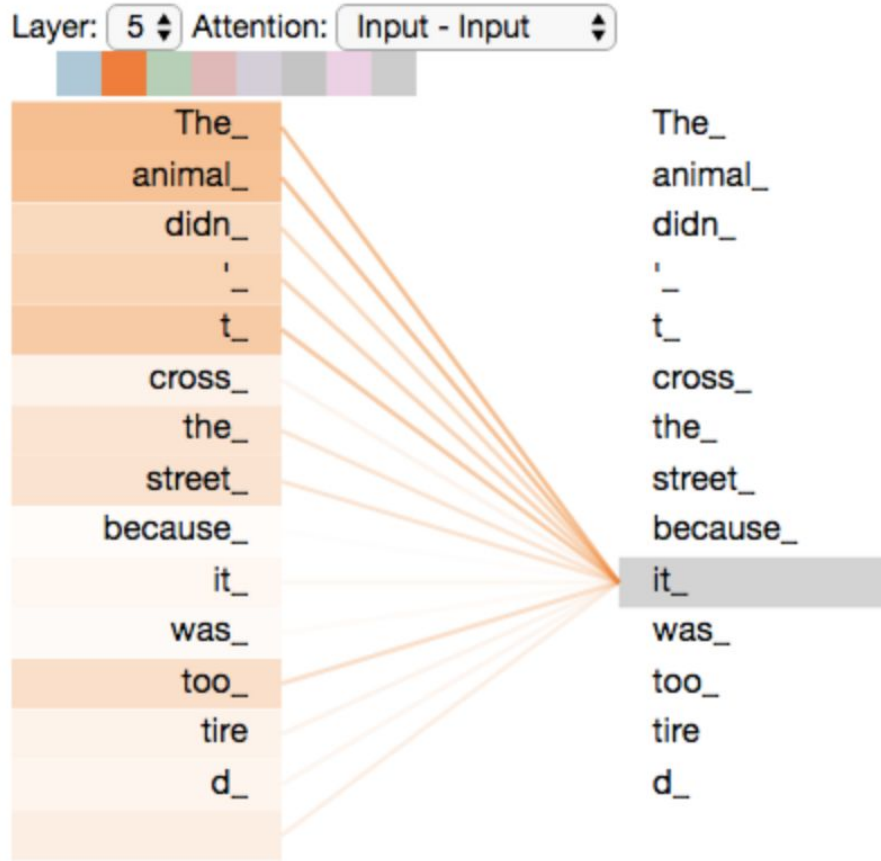
# Self-Attention

# Self-Attention at a High Level

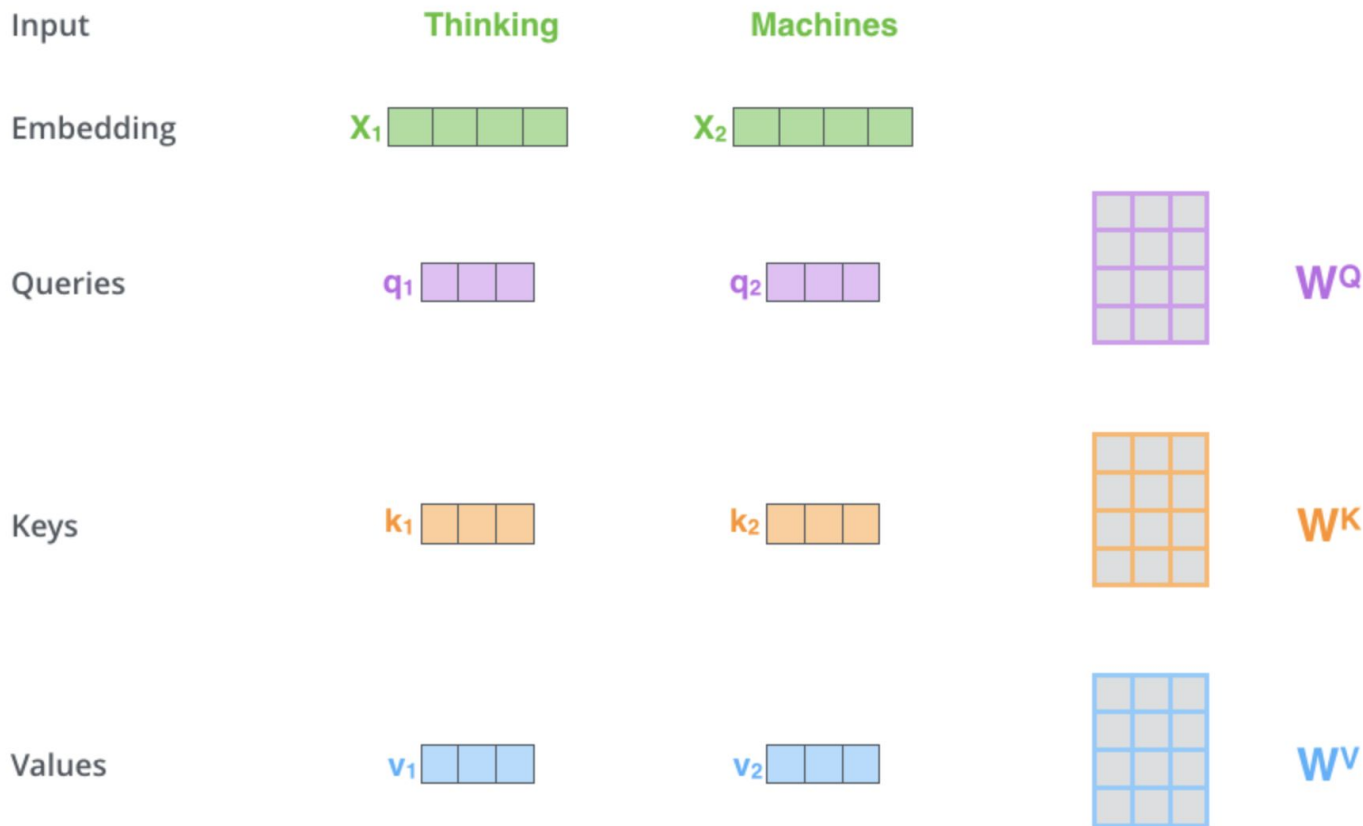
”The animal didn't cross the street because it was too tired”

- What does “it” in this sentence refer to?
- We want self-attention to associate “it” with “animal”
- Self-attention is the method the Transformer uses to bake the “understanding” of other relevant words into the one we’re currently processing

# Self-Attention at a High Level



# Self-Attention: detailed explanation

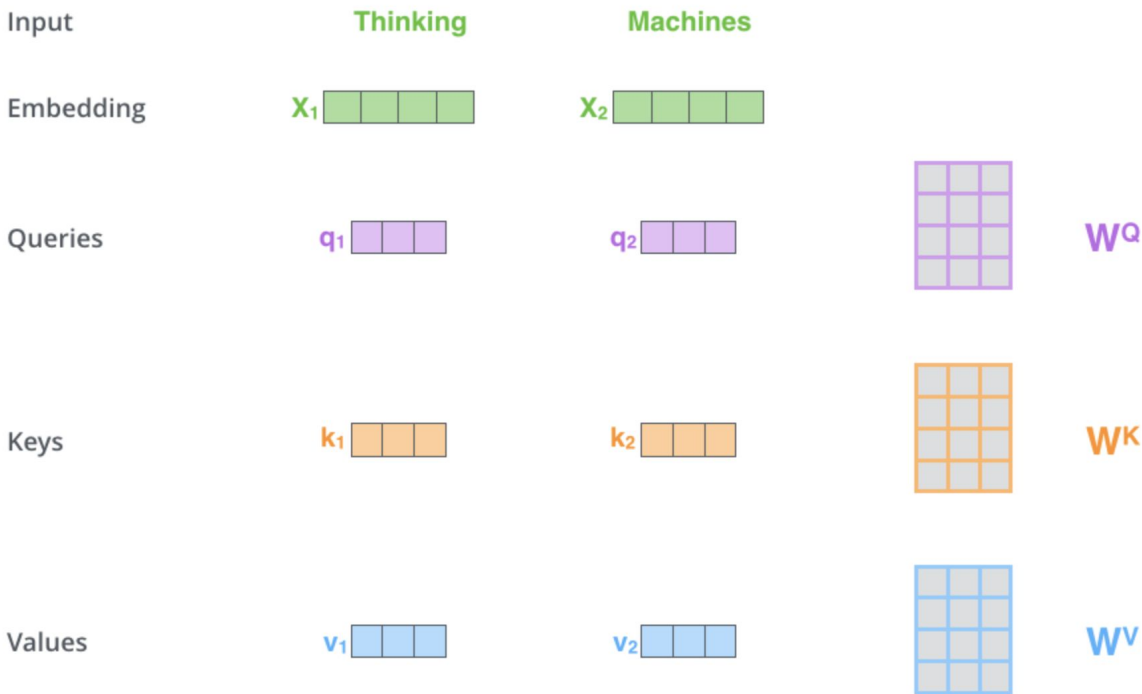


# Self-Attention: detailed explanation

## STEP 1:

create 3 vectors  
(**query**, **key**, **value**)

from each of the encoder's  
input vectors



# Self-Attention: detailed explanation

What are the **query**, **key**, **value** vectors?

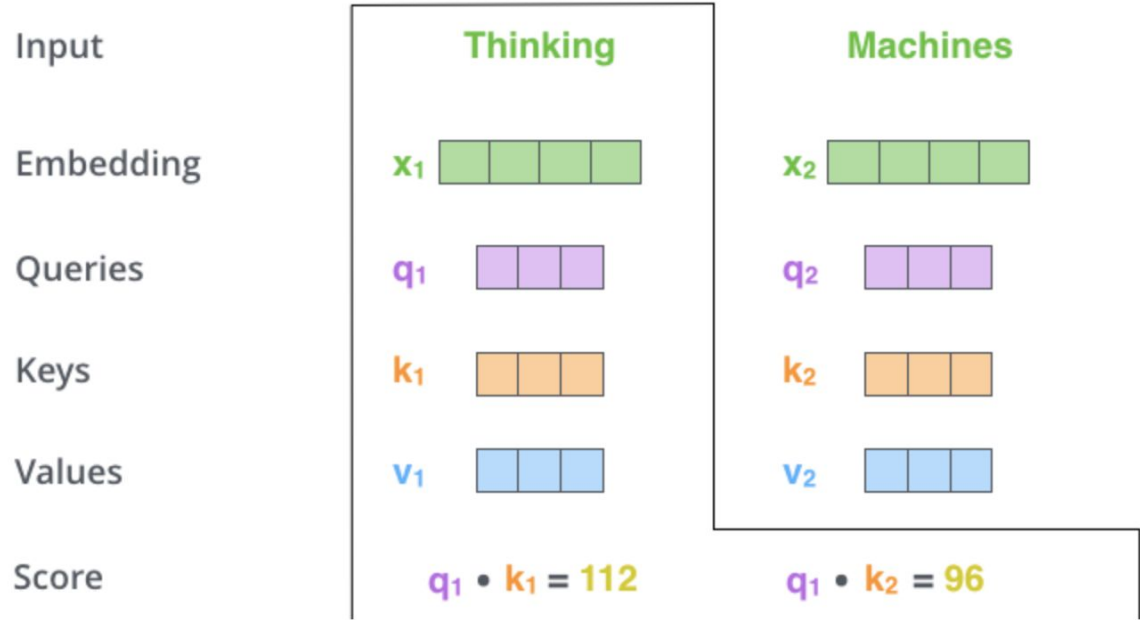
They're abstractions that are useful for calculating and thinking about attention.

# Self-Attention: detailed explanation

## STEP 2:

calculate a score

(score each word of the input sentence against the current word)





# Self-Attention: detailed explanation

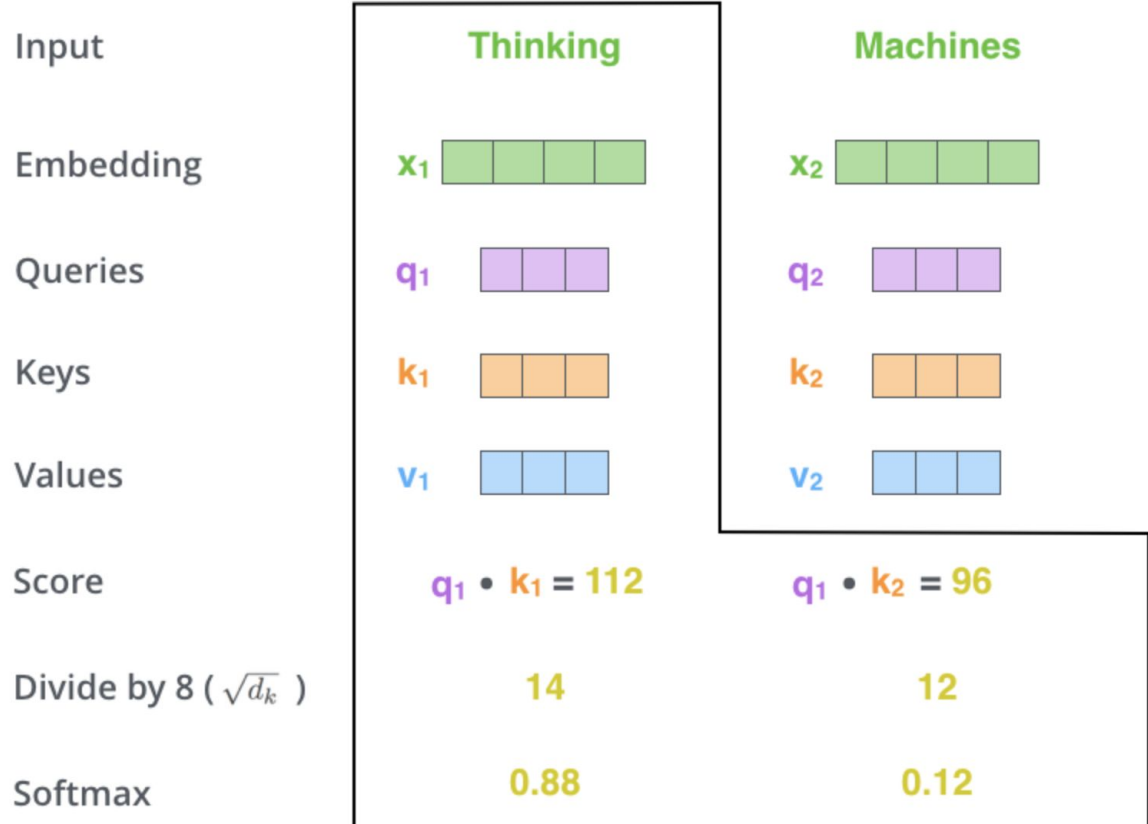
## STEP 3:

divide the scores by 8

(the square root of the  
dimension of the key vectors)

## STEP 4:

softmax



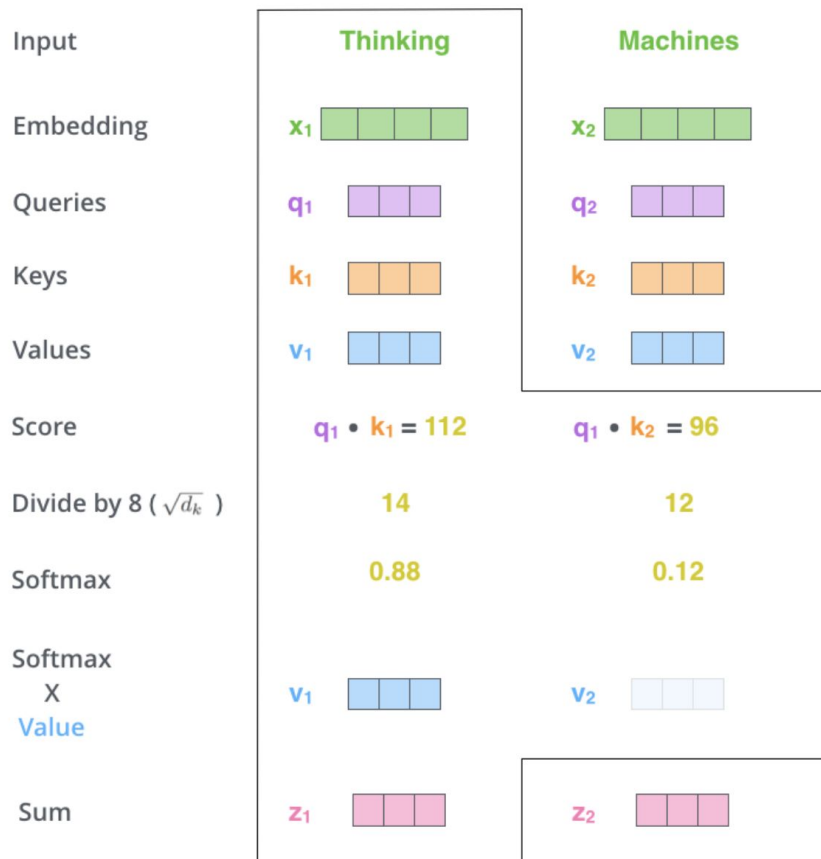
# Self-Attention: detailed explanation

## STEP 5:

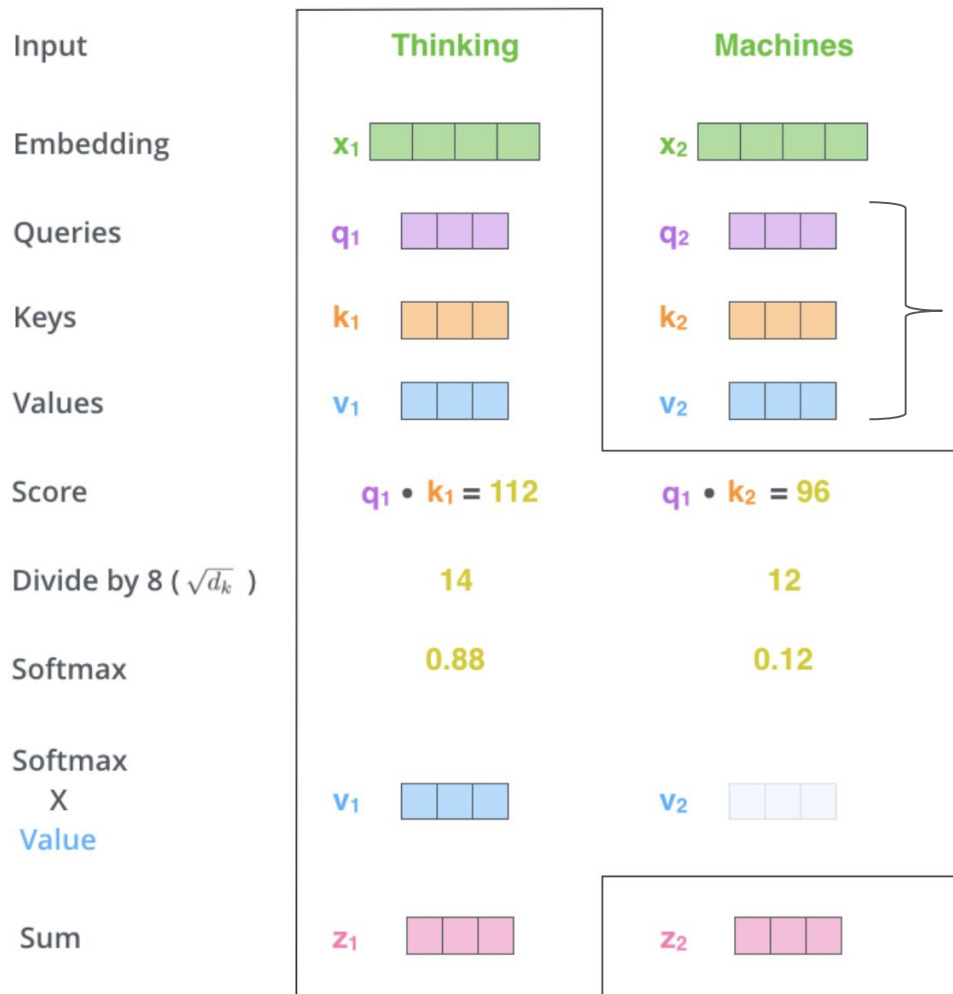
multiply each value vector by the softmax score

## STEP 6:

sum up the weighted value vectors



# Self-Attention



**STEP 1:** create Query, Key, Value

**STEP 2:** calculate scores

**STEP 3:** divide by  $\sqrt{d_k}$

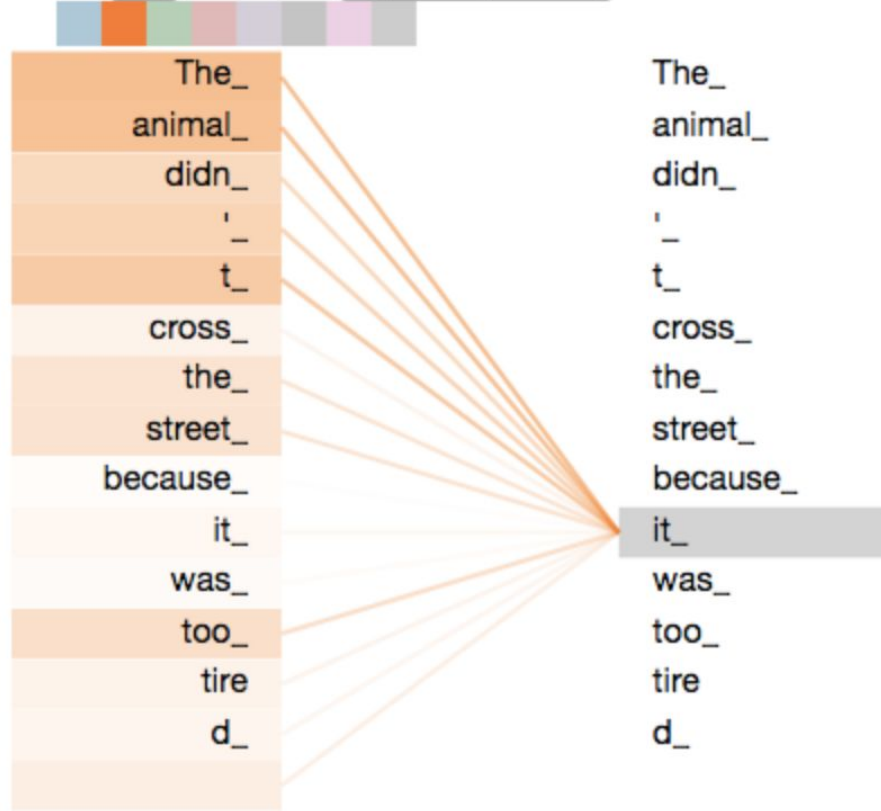
**STEP 4:** softmax

**STEP 5:** multiply each value vector by the softmax score

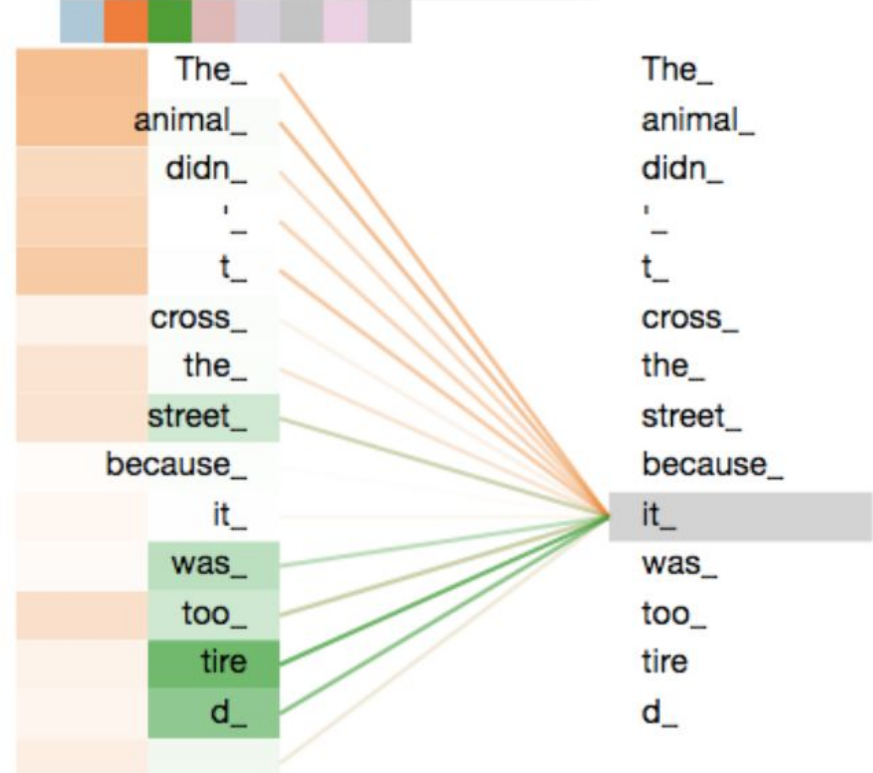
**STEP 6:** sum up the weighted value vectors

# Multi-Head Attention

Layer: 5 Attention: Input - Input

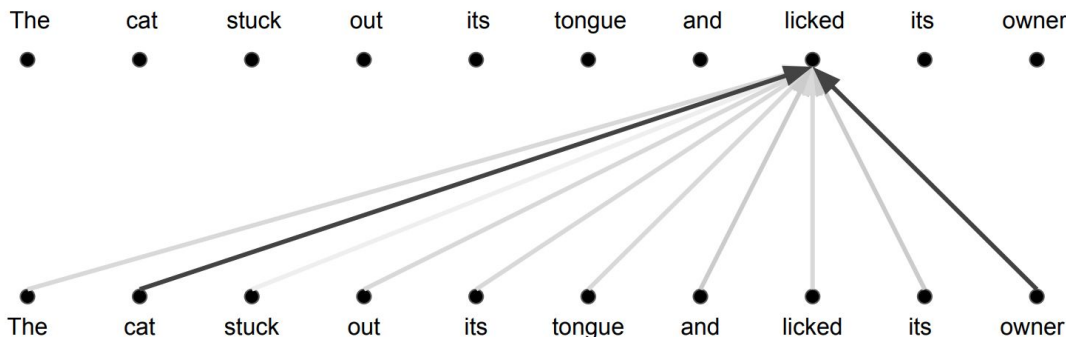


Layer: 5 Attention: Input - Input



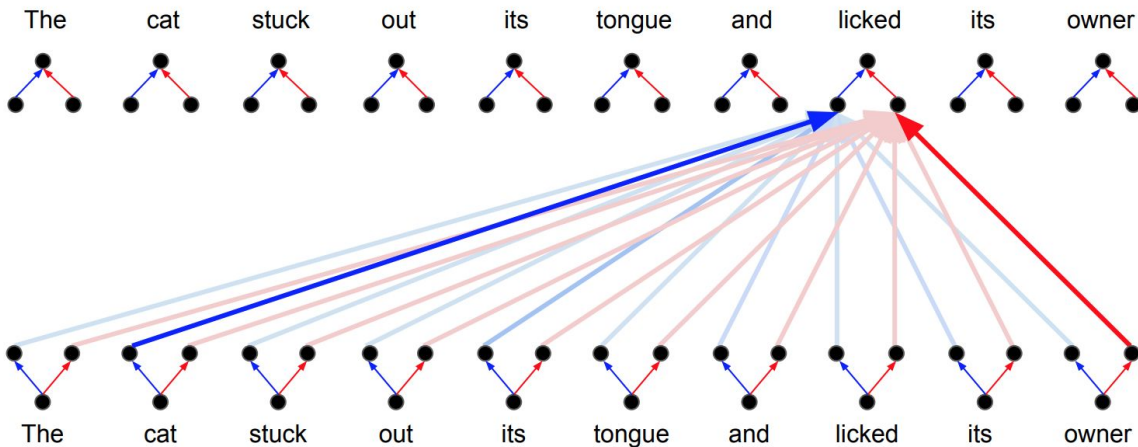
# Attention vs. Multi-Head Attention

**Attention:** a weighted average



**Multi-Head Attention:**

parallel attention layers  
with different linear  
transformations on input  
and output.



# Performance: WMT 2014 BLEU

	EN-DE	EN-FR
GNMT (orig)	24.6	39.9
ConvSeq2Seq	25.2	40.5
Transformer*	<b>28.4</b>	<b>41.8</b>

\*Transformer models trained >3x faster than the others.

# Research Challenges

- Constant 'path length' between any two positions.
- Unbounded memory.
- Trivial to parallelize (per layer).
- Models Self-Similarity.
- Relative attention provides expressive timing, equivariance, and extends naturally to graphs.

# Positional Encoding



# Positional Encoding: why sin and cos?

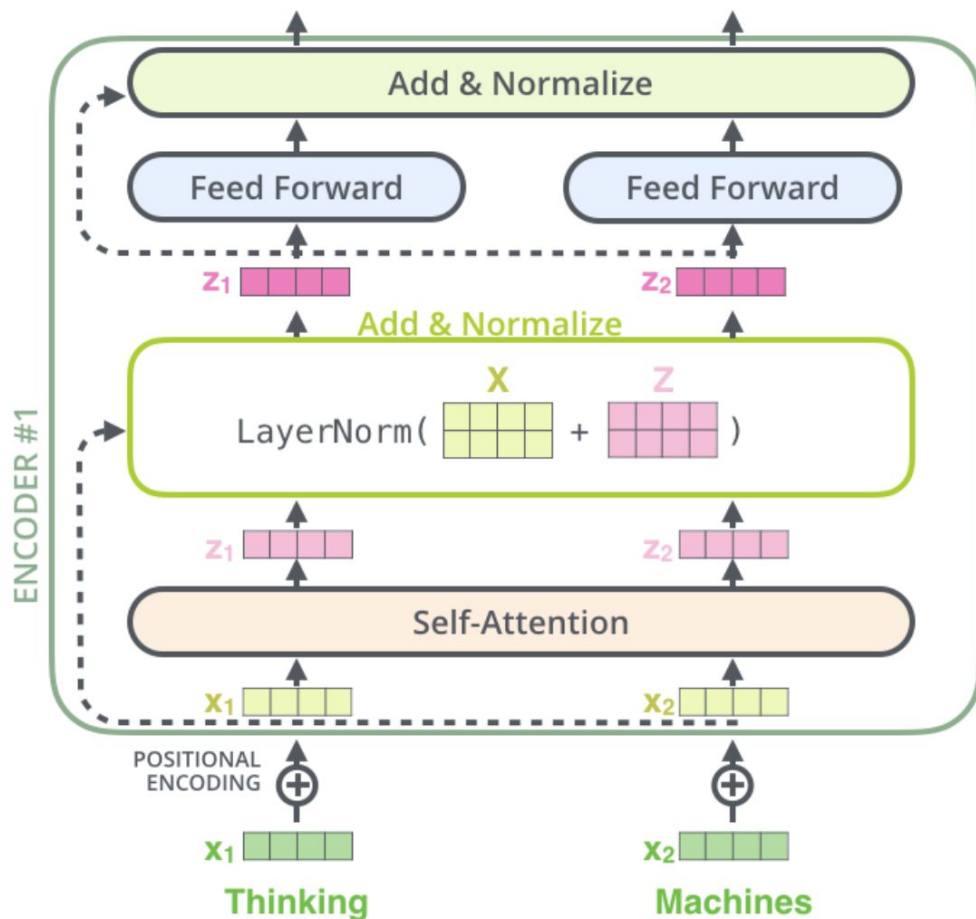
$$\vec{p}_t^{(i)} = f(t)^{(i)} = \begin{cases} \sin(\omega_k t), & \text{if } i = 2k \\ \cos(\omega_k t), & \text{if } i = 2k + 1 \end{cases}$$
$$\omega_k = \frac{1}{10000^{2k/d}}$$
$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1.t) \\ \cos(\omega_1.t) \\ \\ \sin(\omega_2.t) \\ \cos(\omega_2.t) \\ \\ \vdots \\ \\ \sin(\omega_{d/2}.t) \\ \cos(\omega_{d/2}.t) \end{bmatrix}_{d \times 1}$$

t stays for position in the original sequence

k is the index of the element in the positional vector

# Layer Normalization

# Layer Normalization



Like BatchNorm

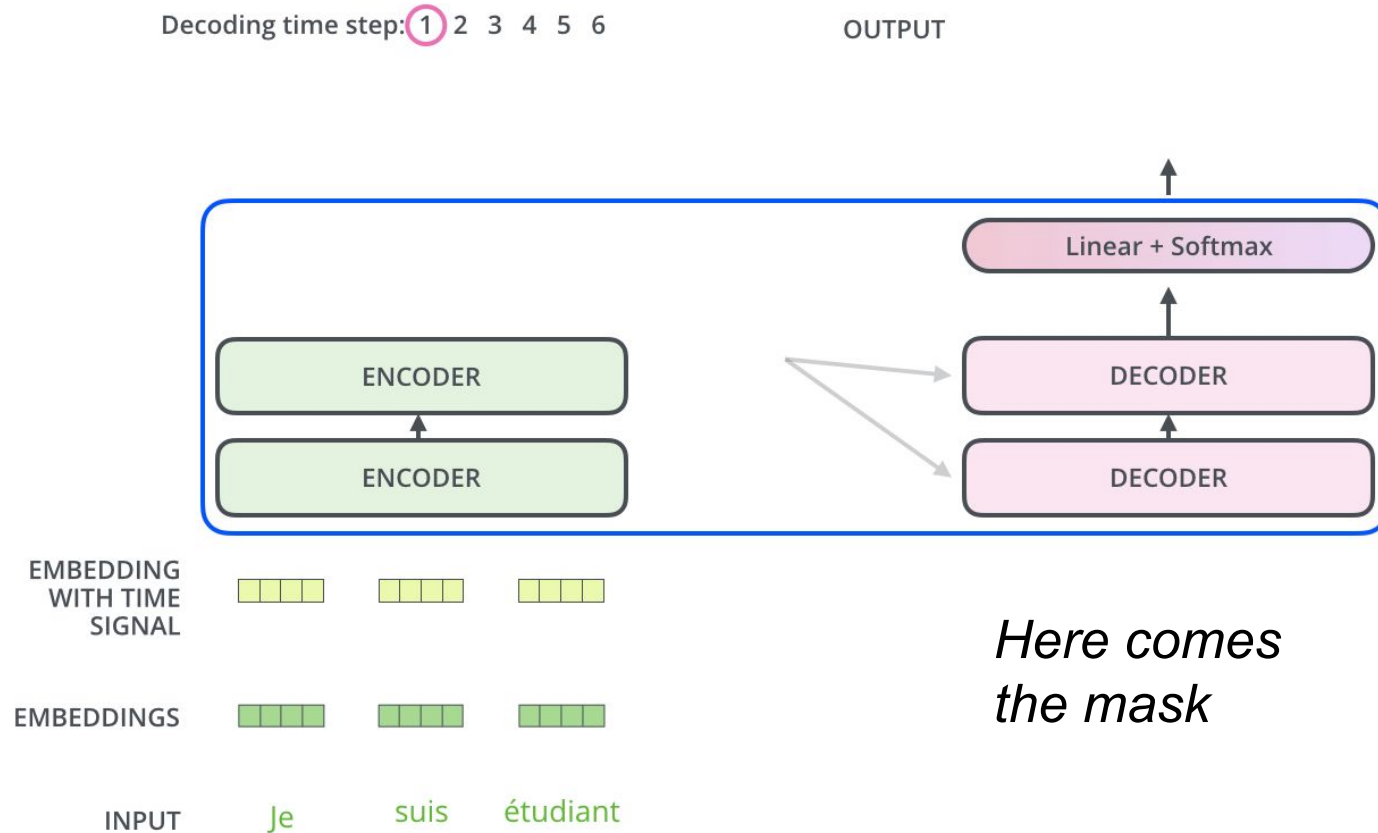
but normalize along  
all features  
representing latent  
vector

More info:

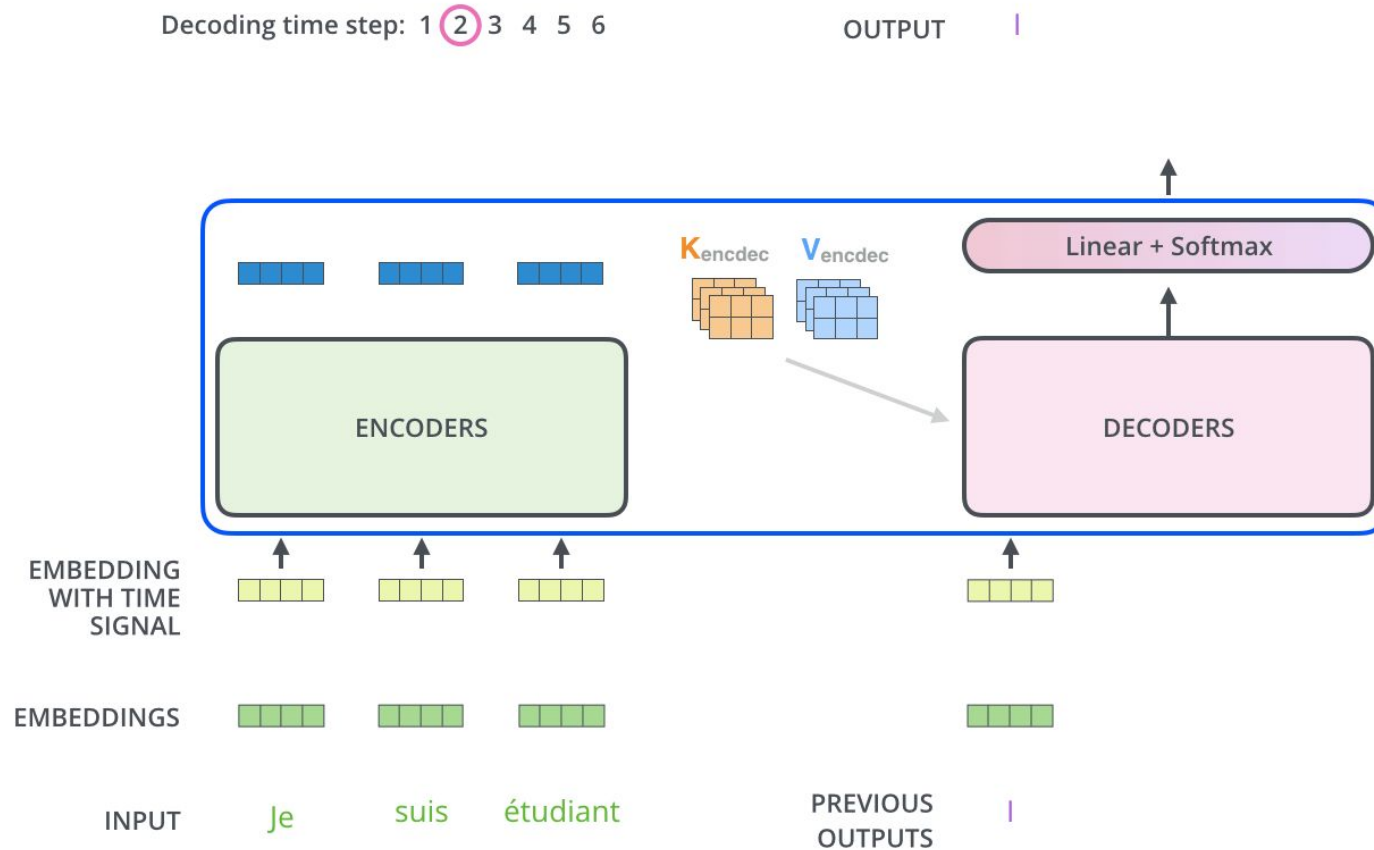
[Layer Normalization](https://jalammar.github.io/illustrated-transformer/)

# The Decoder

# The Decoder Side



# The Decoder Side



# 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.

### 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)



**Classifier**

75% Spam  
25% Not Spam

**Dataset:**

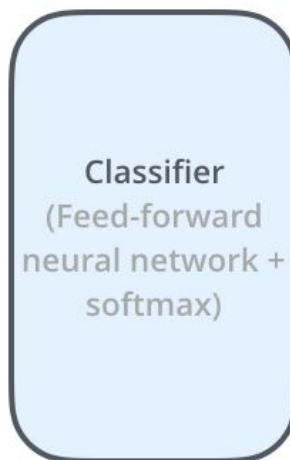
Email message	Class
Buy these pills	Spam
Win cash prizes	Spam
Dear Mr. Atreides, please find attached...	Not Spam



# BERT

Input  
Features

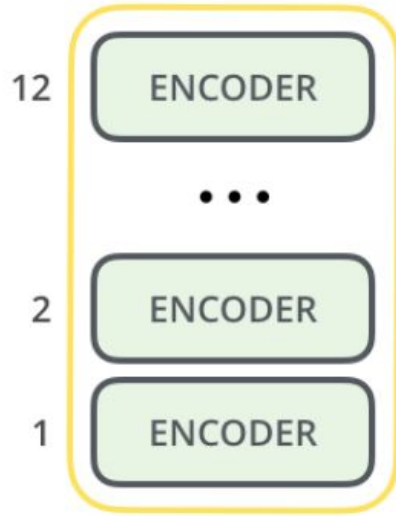
Help Prince Mayuko Transfer  
Huge Inheritance



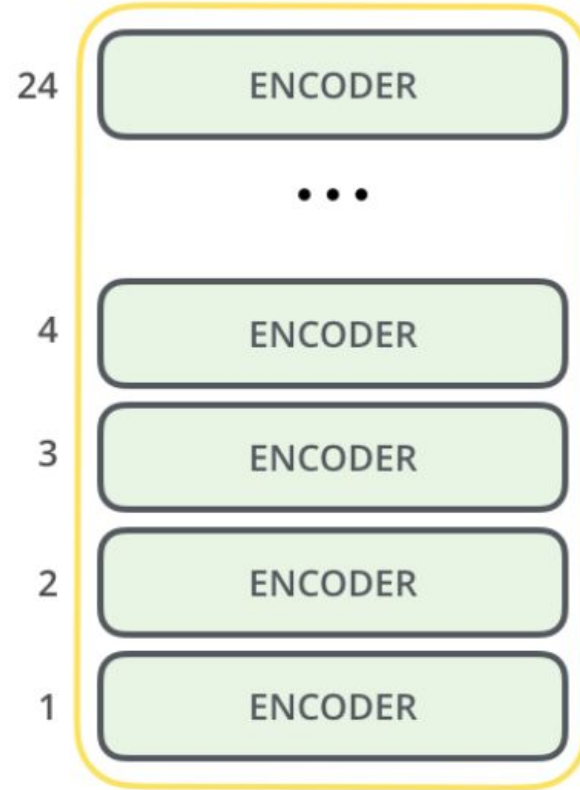
Output  
Prediction



# BERT: base and large





BERT<sub>BASE</sub>

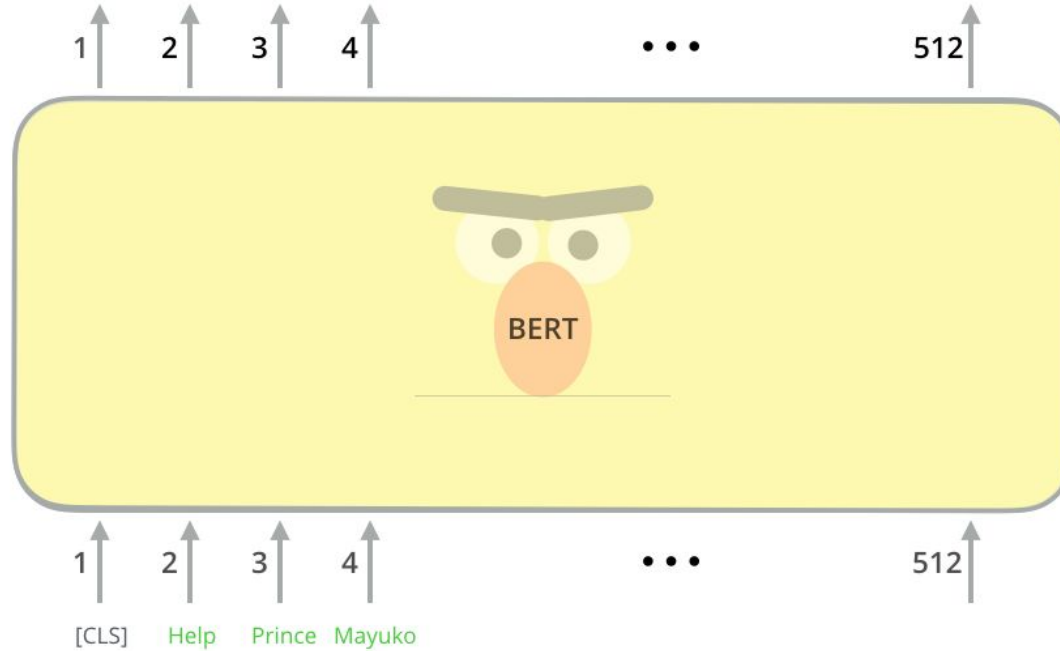


BERT<sub>LARGE</sub>

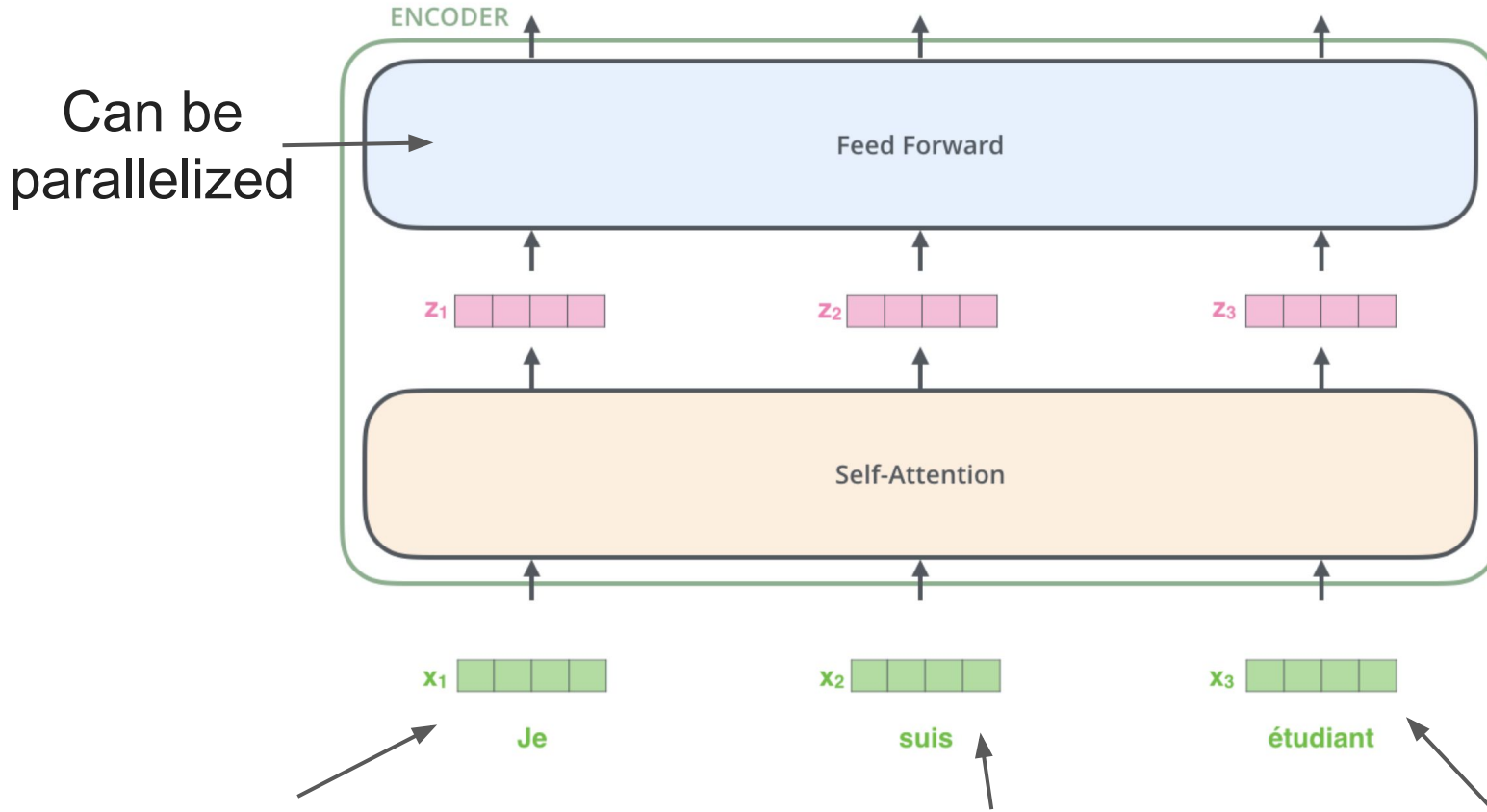
# BERT vs. Transformer

			
		Base BERT	Large BERT
Encoders	6	12	24
Units in FFN	512	768	1024
Attention Heads	8	12	16

# Model inputs

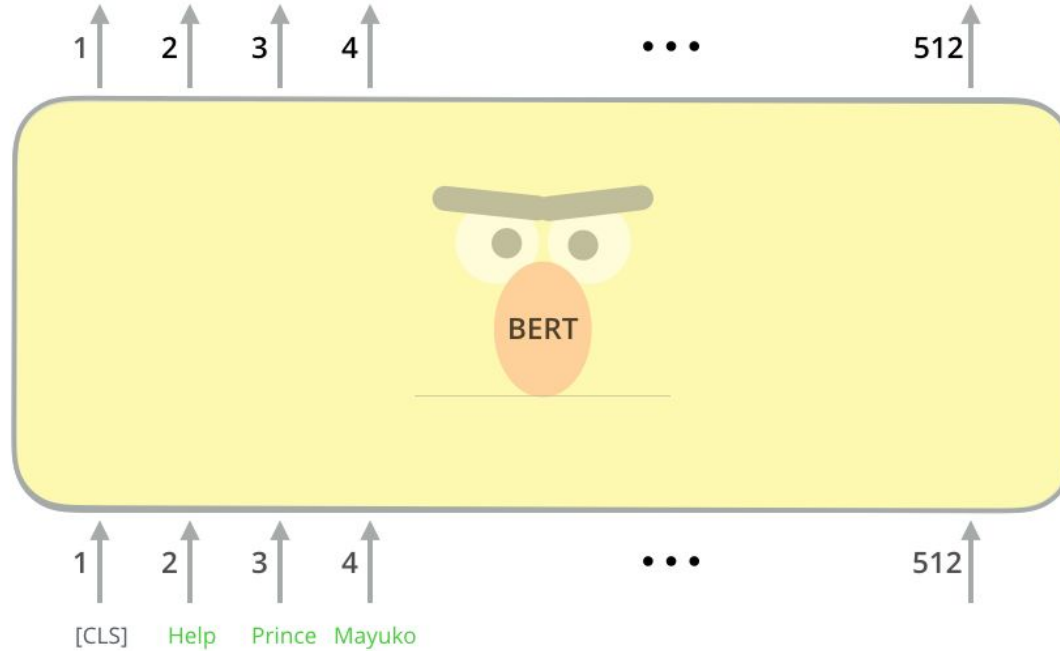


# Transformer Block in BERT



the word in each position flows through its own path in the encoder

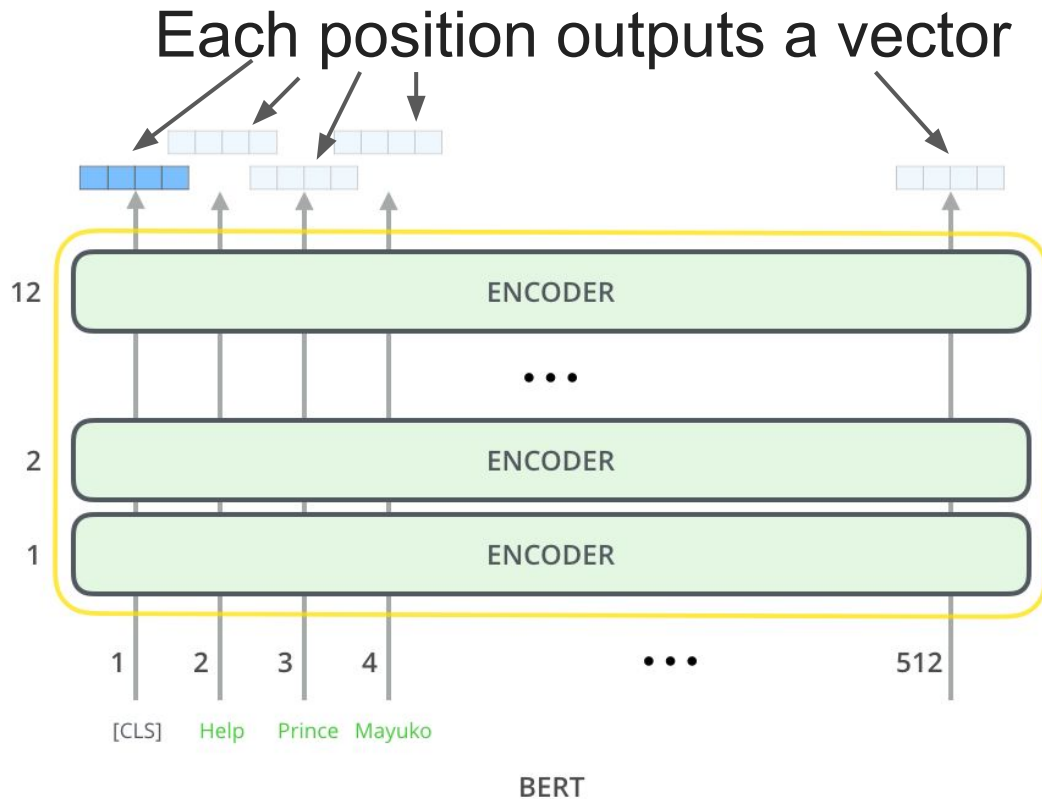
Model inputs



Identical to the Transformer up until this point

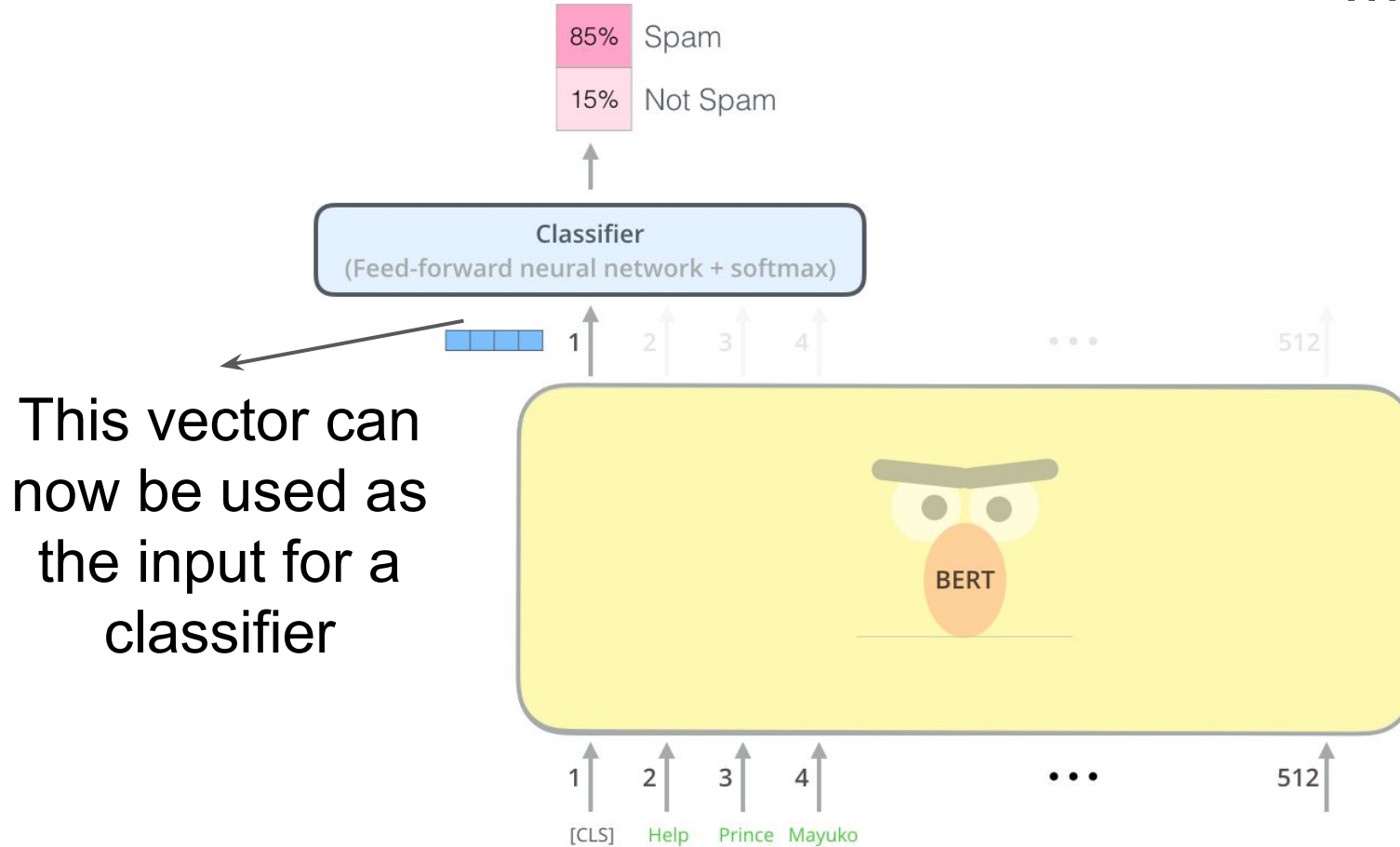
Why is BERT so special?

# Model outputs



For sentence classification we focus on the first position  
(that we passed [CLS] token to)

# Model inputs



This vector can now be used as the input for a classifier



# BERT: pre-training

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

1 2 3 4 5 6 7 8 ... 512



Randomly mask  
15% of tokens

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

Input

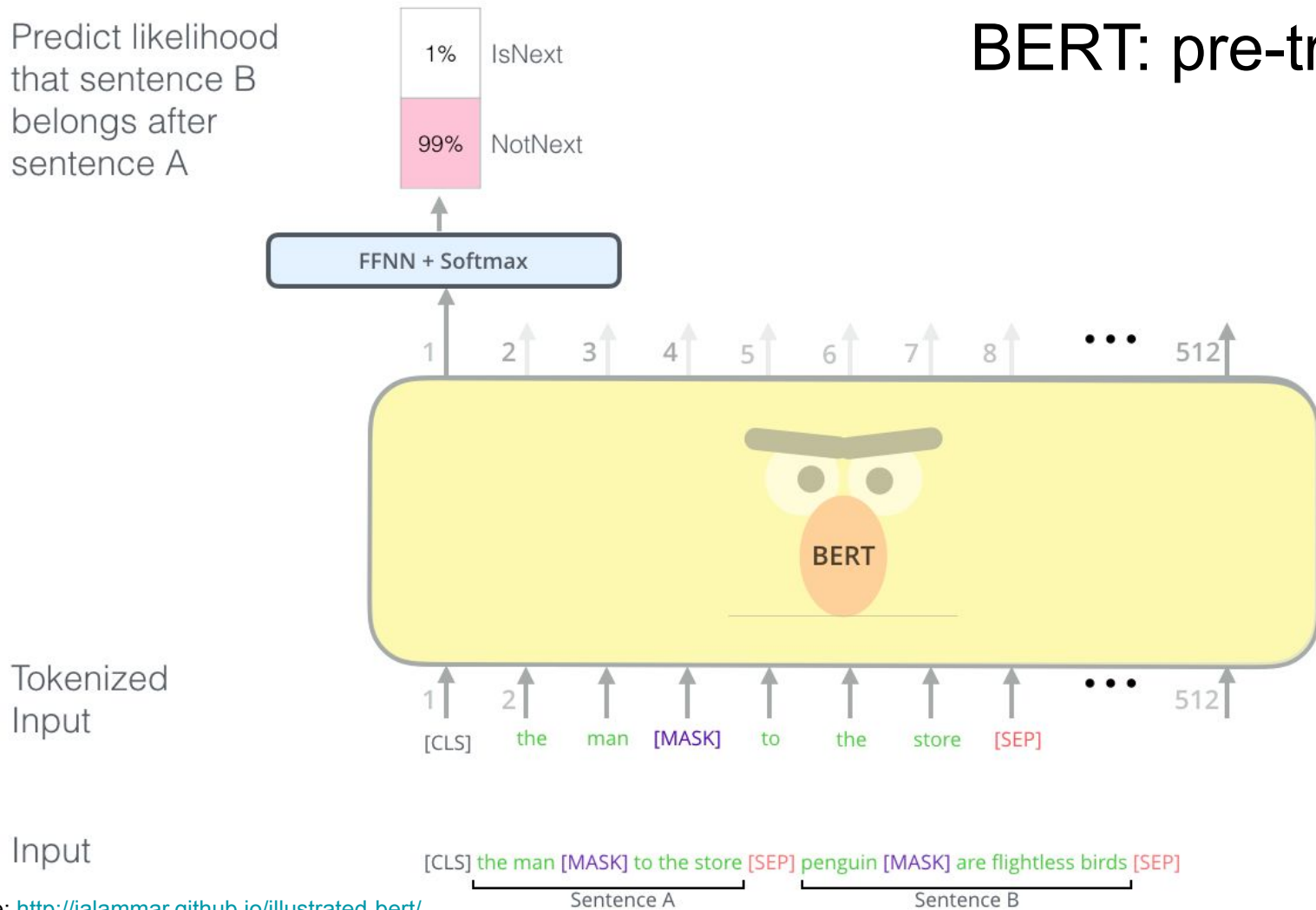
↑ ↑ ↑ ↑ ↑ ↑ ↑ ↑  
[CLS] Let's stick to improvisation in this skit

# BERT: pre-training

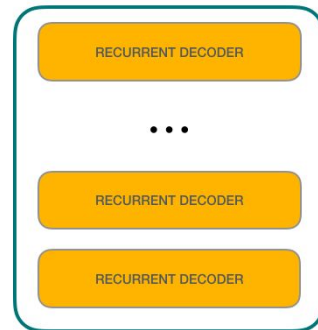
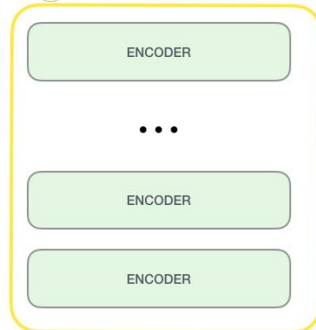
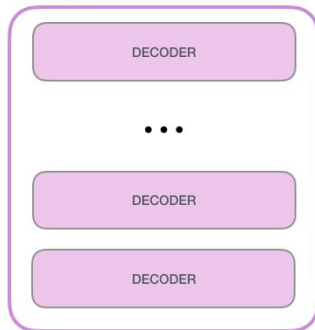
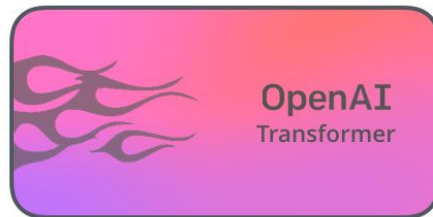
- “Masked Language Model” approach
- To make BERT better at handling relationships between multiple sentences, the pre-training process includes an additional task:  
*“Given two sentences (A and B), is B likely to be the sentence that follows A, or not?”*

# BERT: pre-training

Predict likelihood  
that sentence B  
belongs after  
sentence A



- Transformer
- OpenAI Transformer
- ELMO
- BERT
- BERTology
- GPT
- GPT-2
- GPT-3



- Transformer is novel and very powerful architecture
- It is worth it to understand how Self-Attention works
- BERT is variant of Decoders from Transformer for variety of tasks
- GPT are even bigger and better in metrics but they are made by corporations