

BERT



(Bi-Directional Encoder Representations from Transformer)

A Short Introduction

By Tariq Jamil
PGD-DS (BII)

Traditional Approach / LSTM

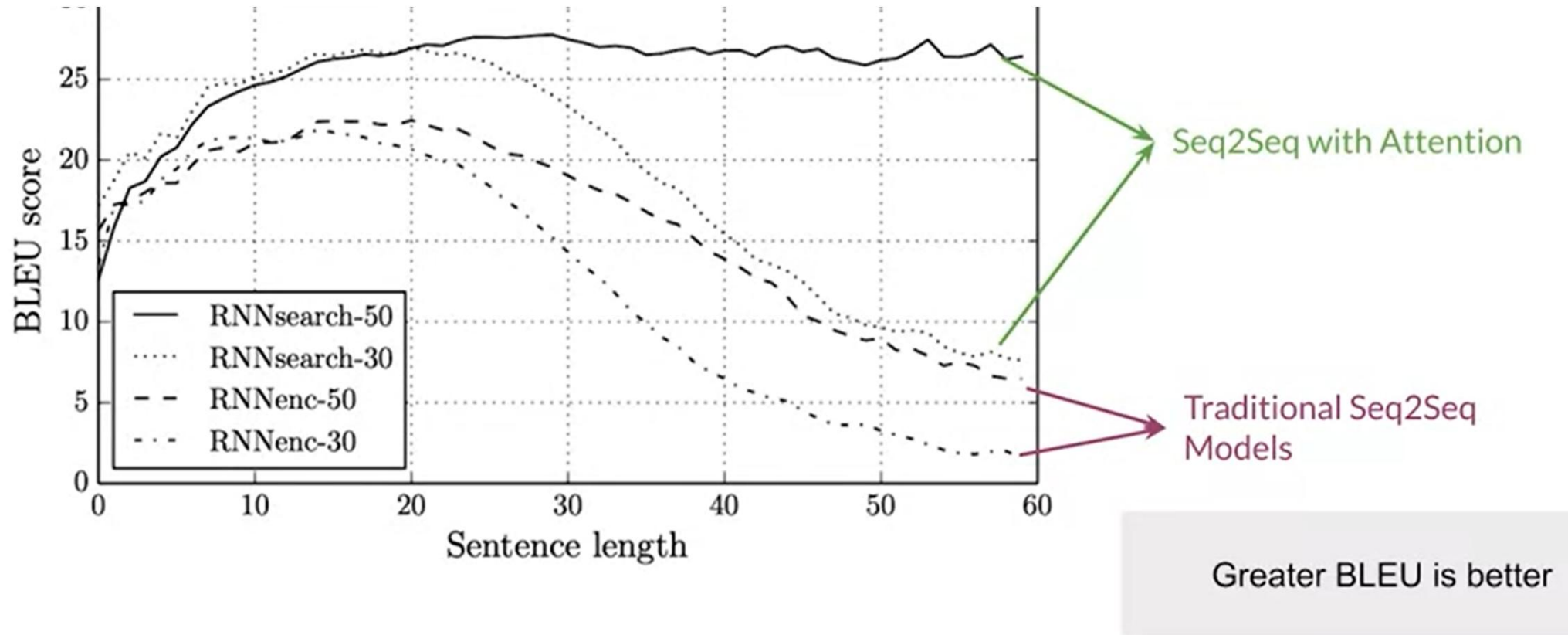
Seq2Seq shortcomings

- Variable-length sentences + fixed-length memory =

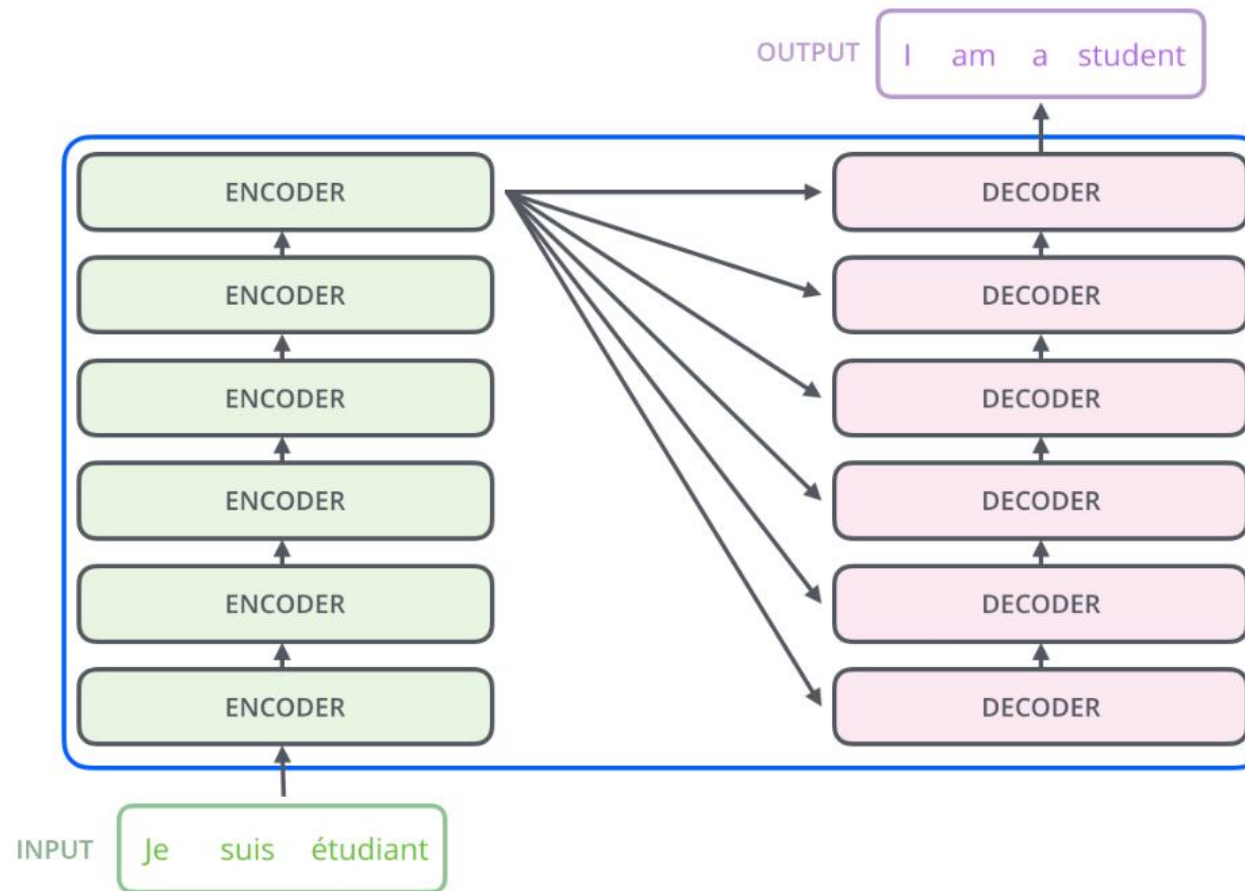


- As sequence size increases, model performance decreases

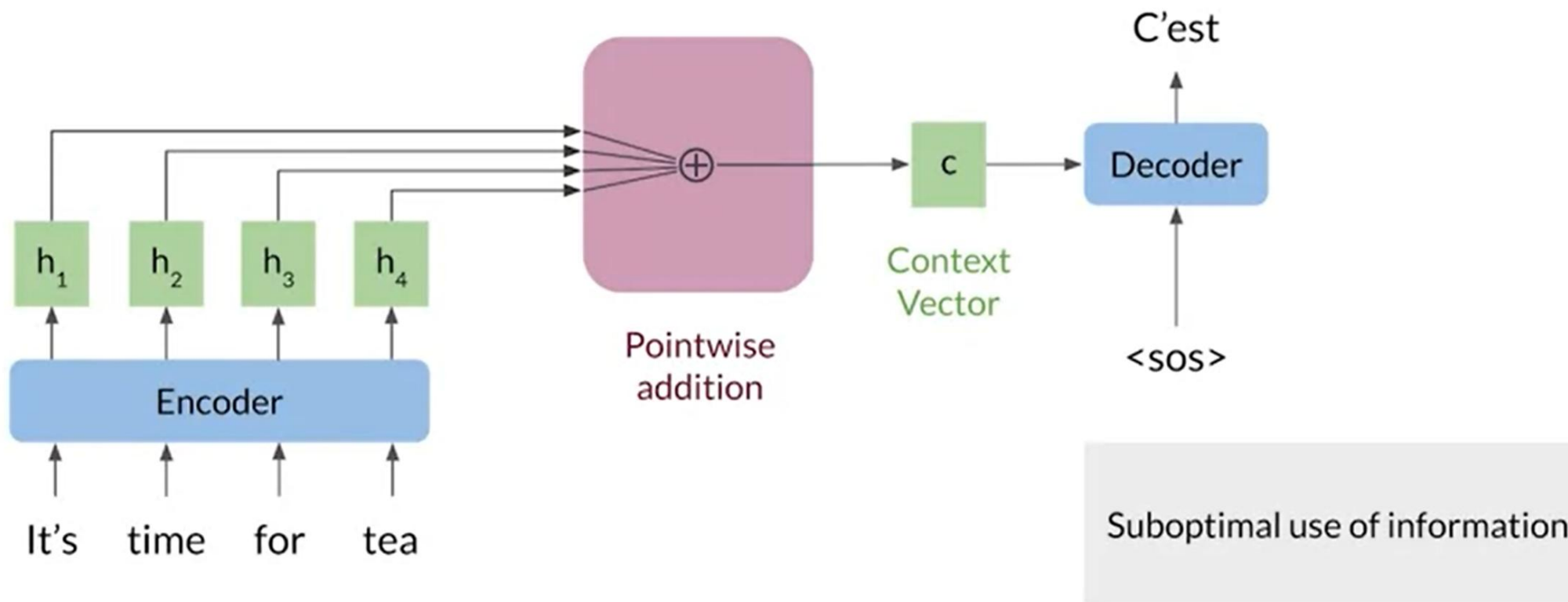
Seq2Seq Models Comparison



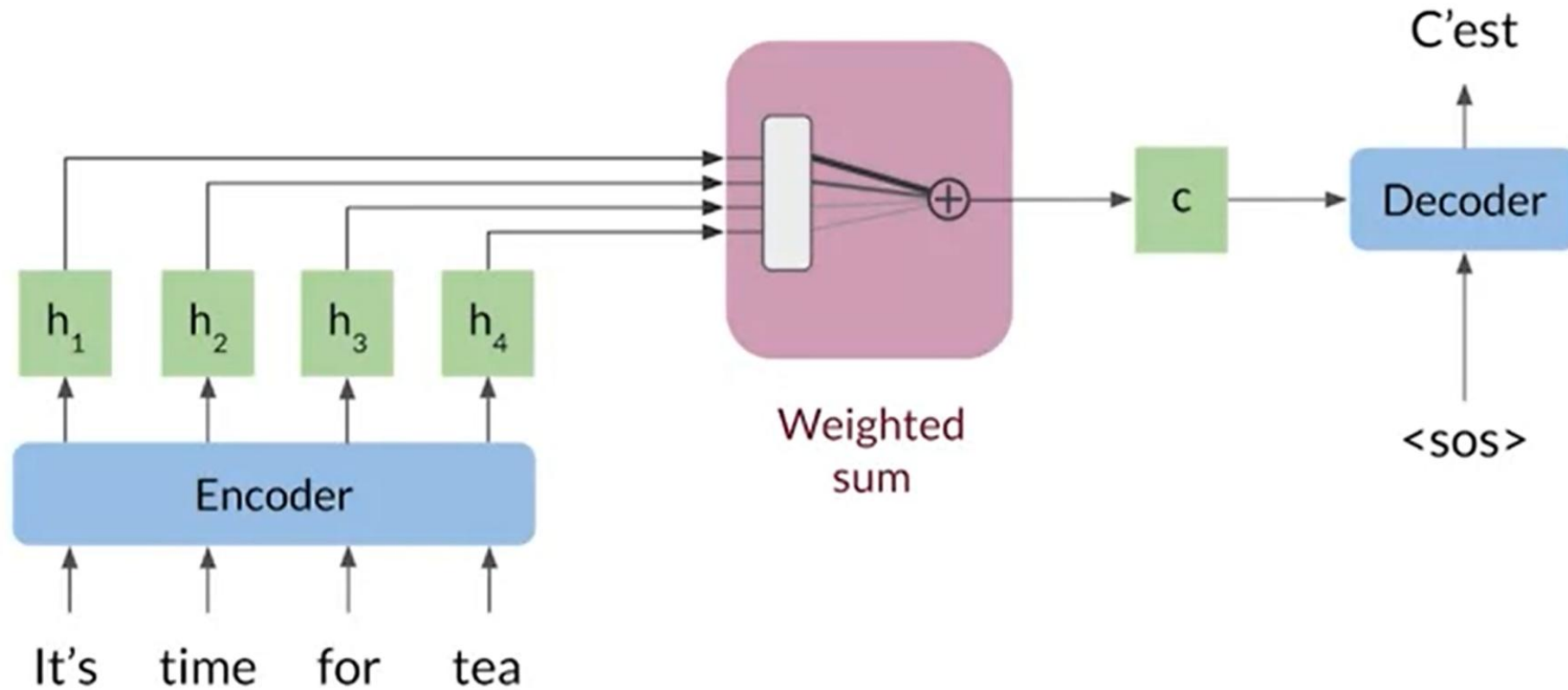
Transformer | Block Diagram



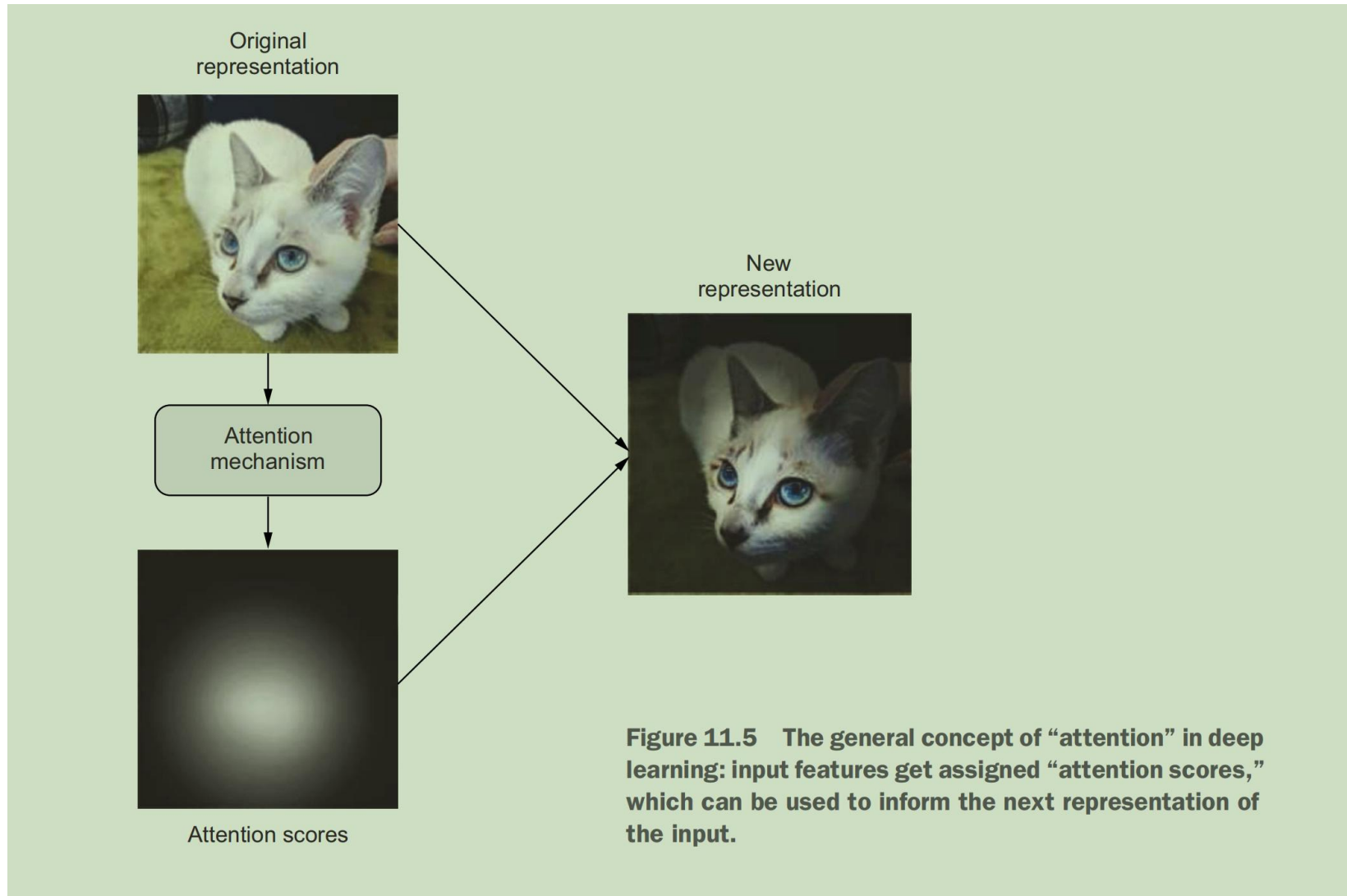
Transformers | Traditional Approach



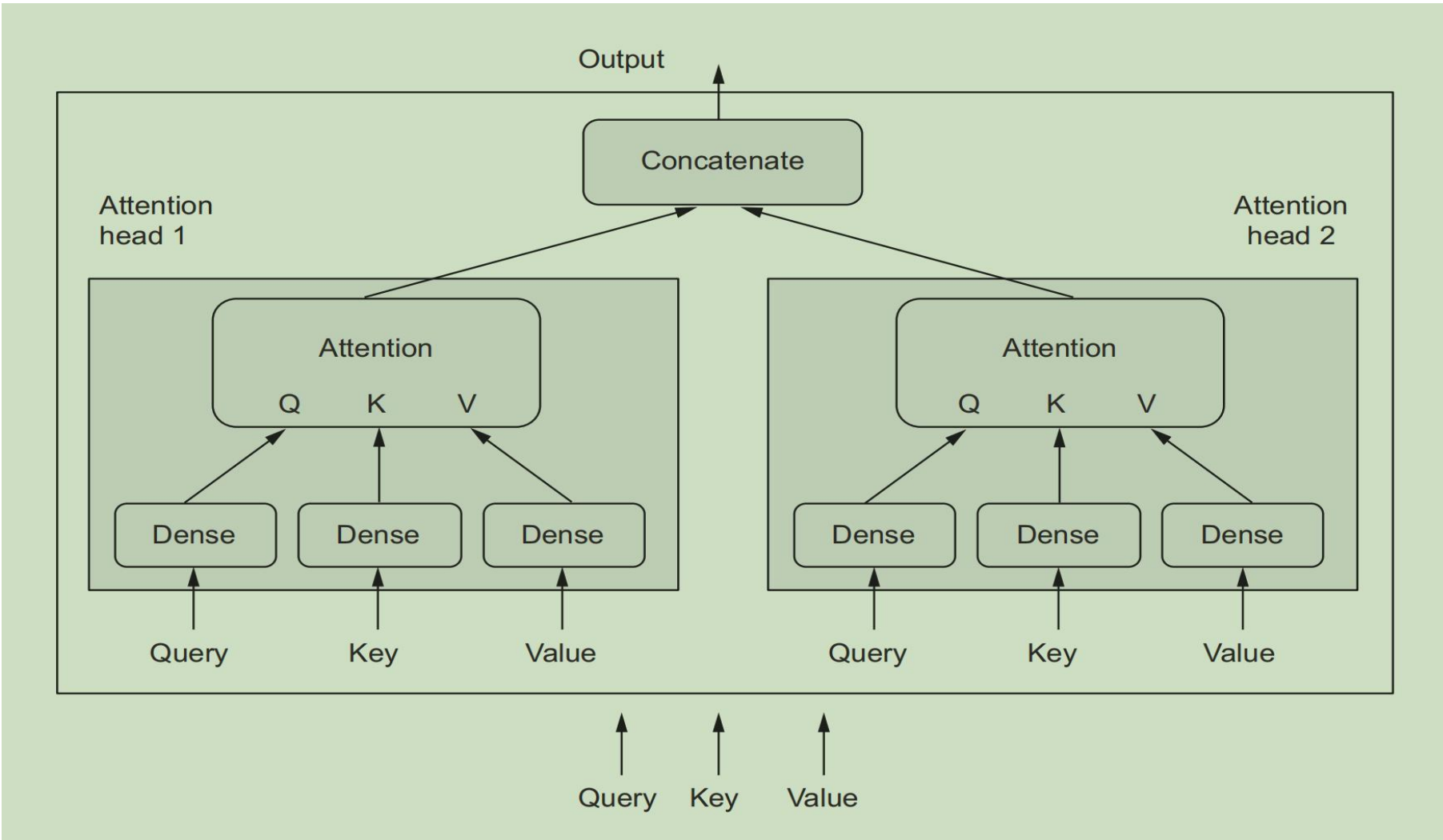
Transformers | with Self Attention



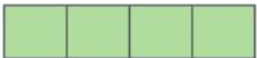
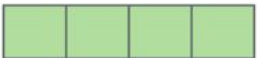
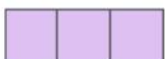
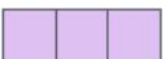
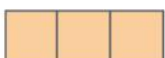
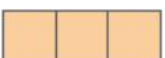
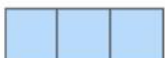
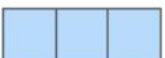
Understanding - Self Attention



Self Attention

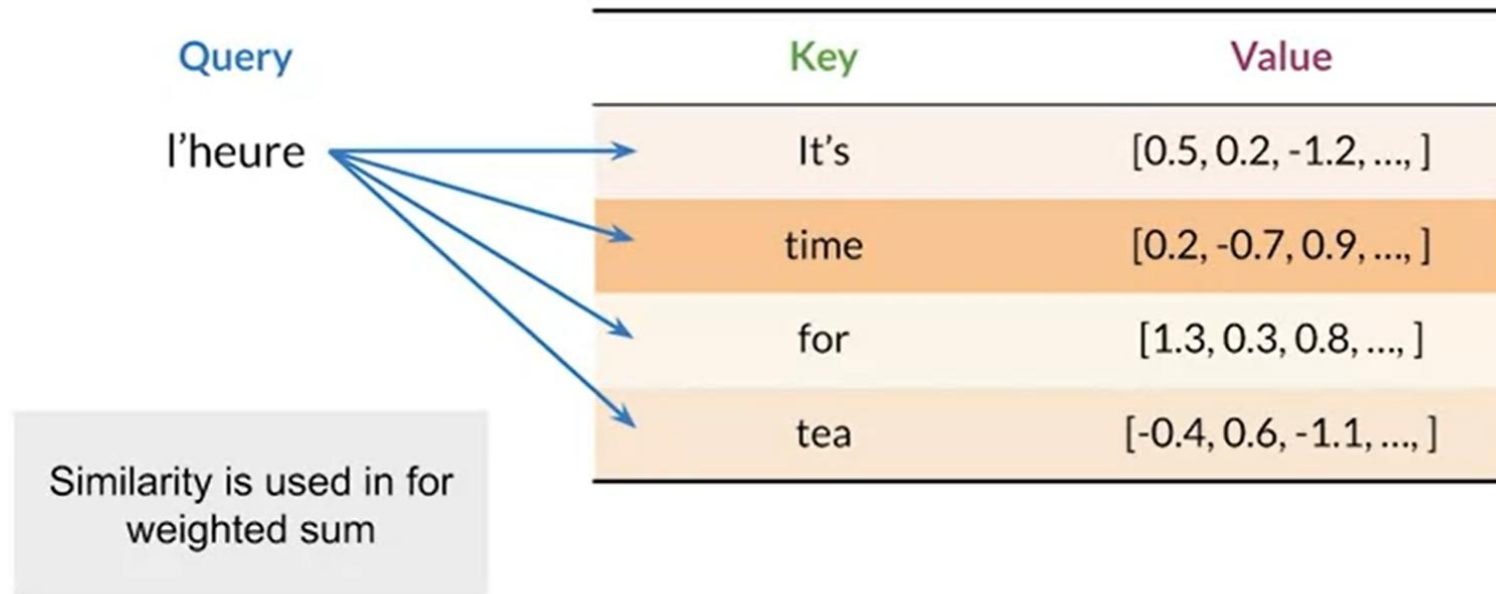


Process Flow

Input	Thinking	Machines
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

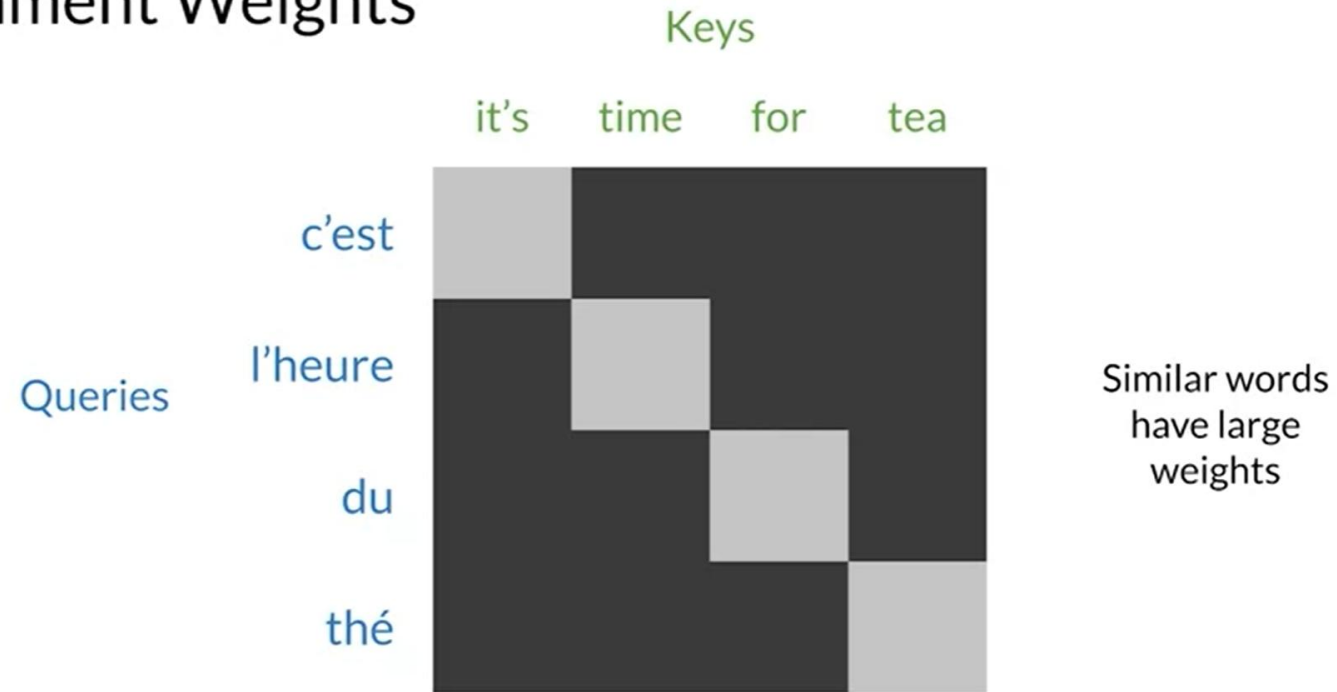
Attention - components

Queries, Keys, Values



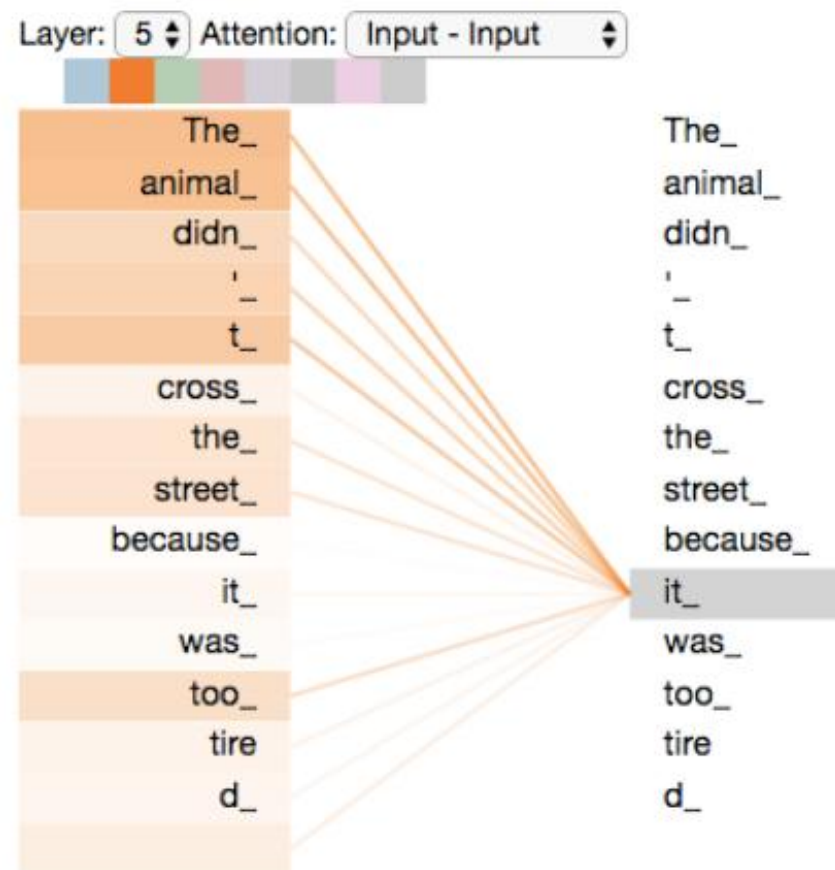
Similarity Score

Alignment Weights



Example - Why need self attention

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

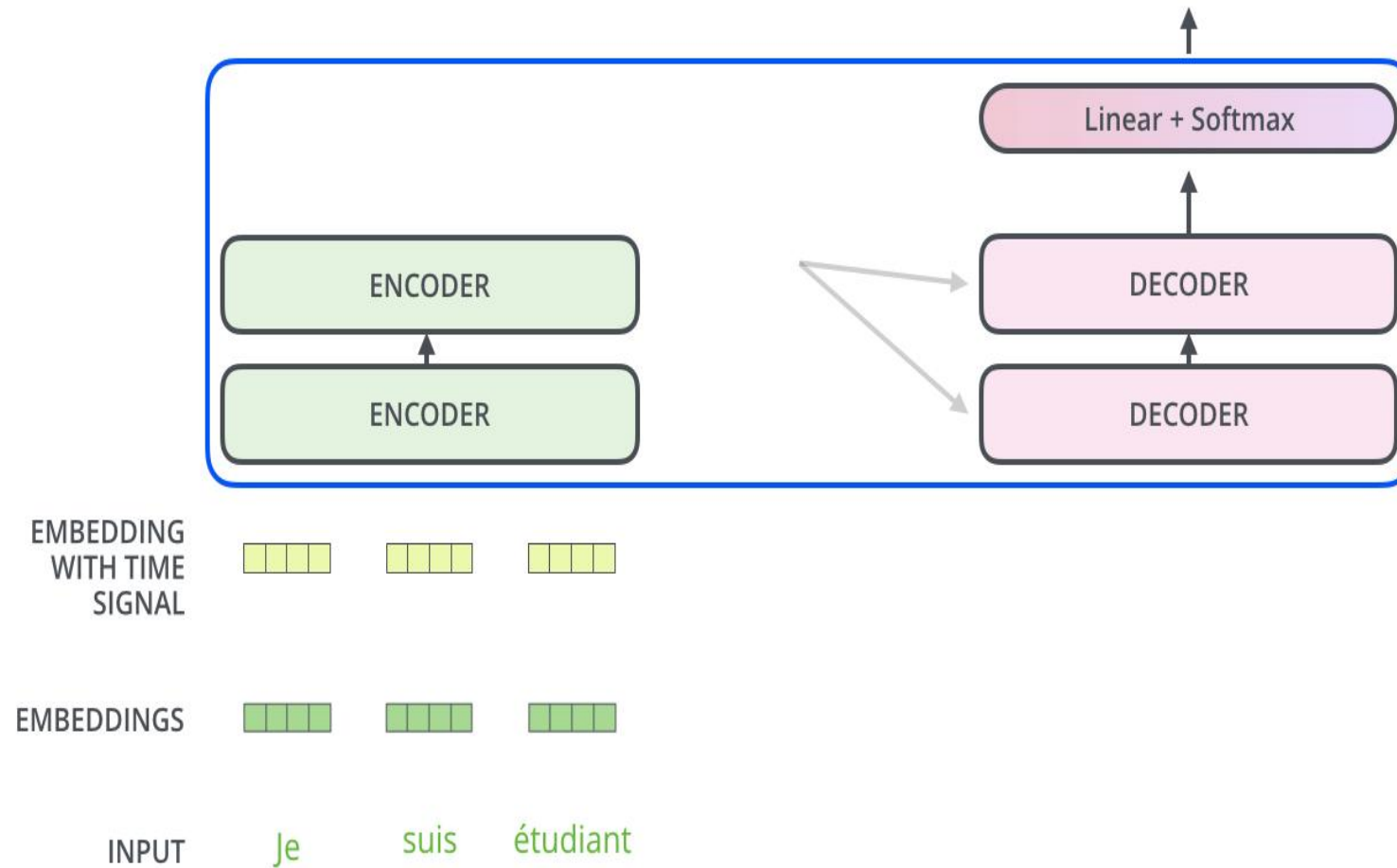


Visualization

<http://jalammar.github.io/im>

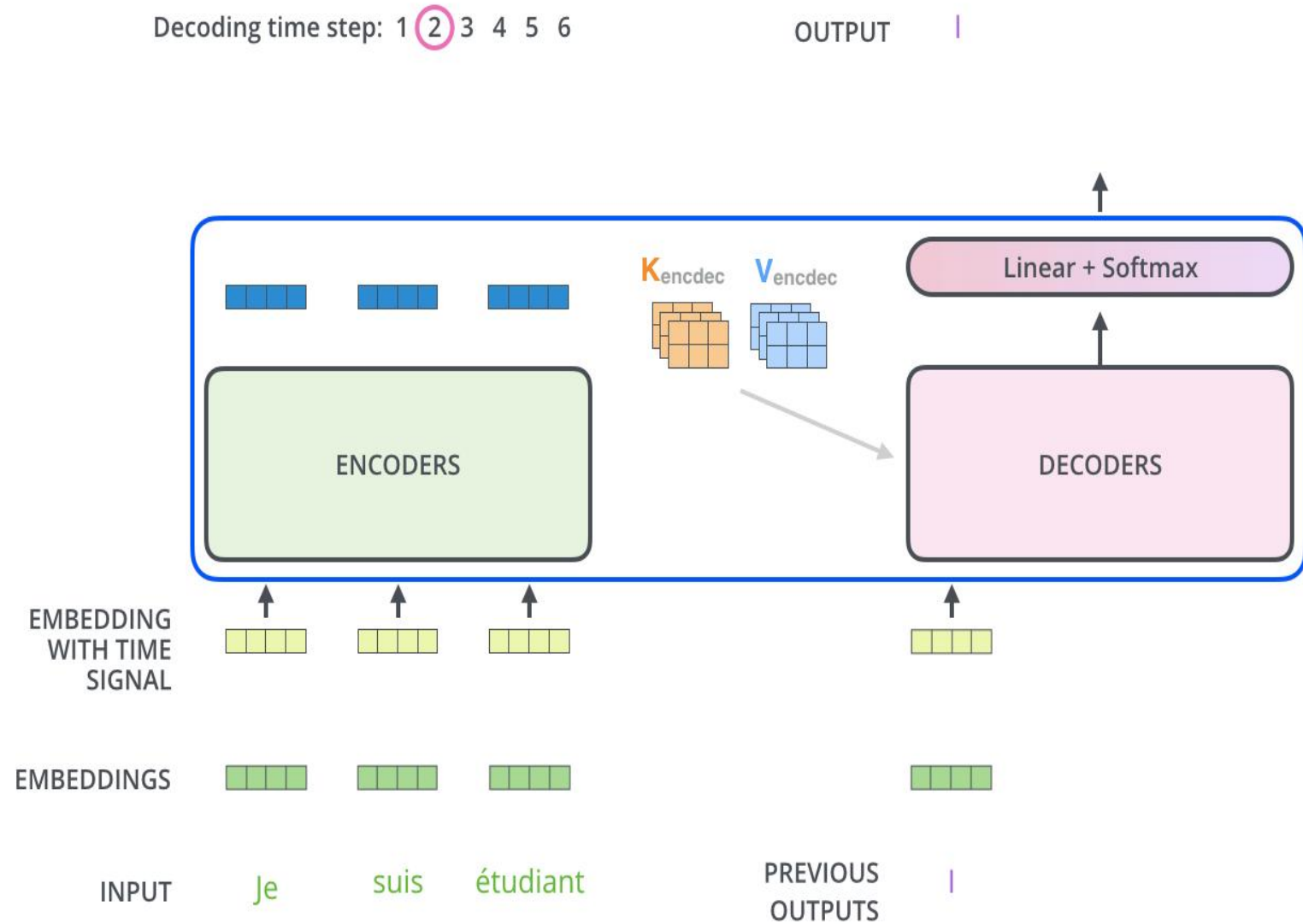
Decoding time step: 1 2 3 4 5 6

OUTPUT



Visualization

<http://jalammar.github.io/image>



BERT | A Transformer with Attention and bidirection Training

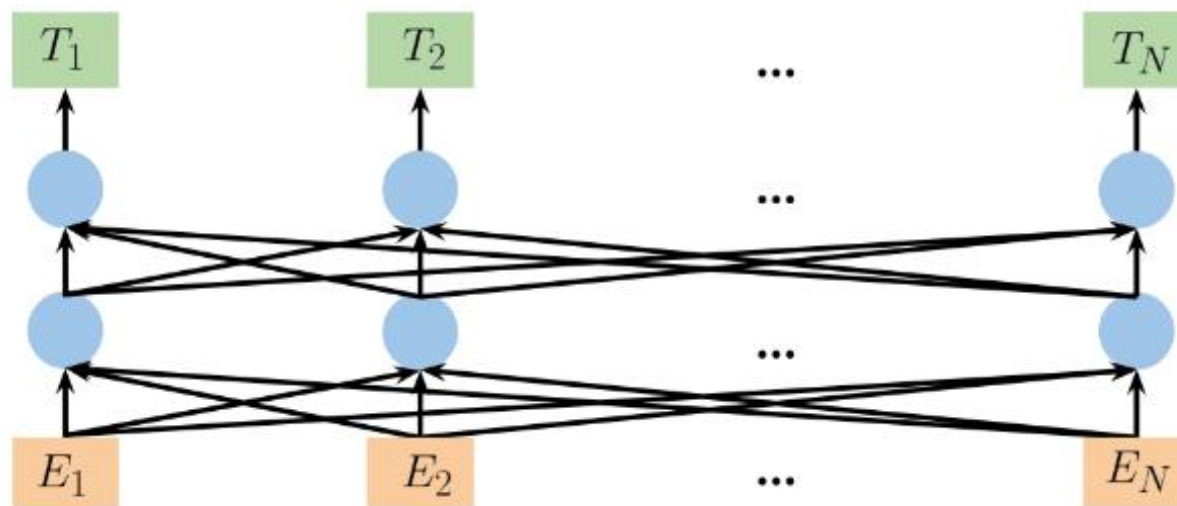
- A multi layer bidirectional transformer
- Positional embeddings
- BERT_base:
 - 12 layers (12 transformer blocks)
 - 12 attentions heads
 - 110 million parameters

BERT | Architecture

Bidirectional Encoder Representations from Transformers (BERT)

You will now learn about the BERT architecture and understand how the pre-training works.

- Makes use of transfer learning/pre-training:

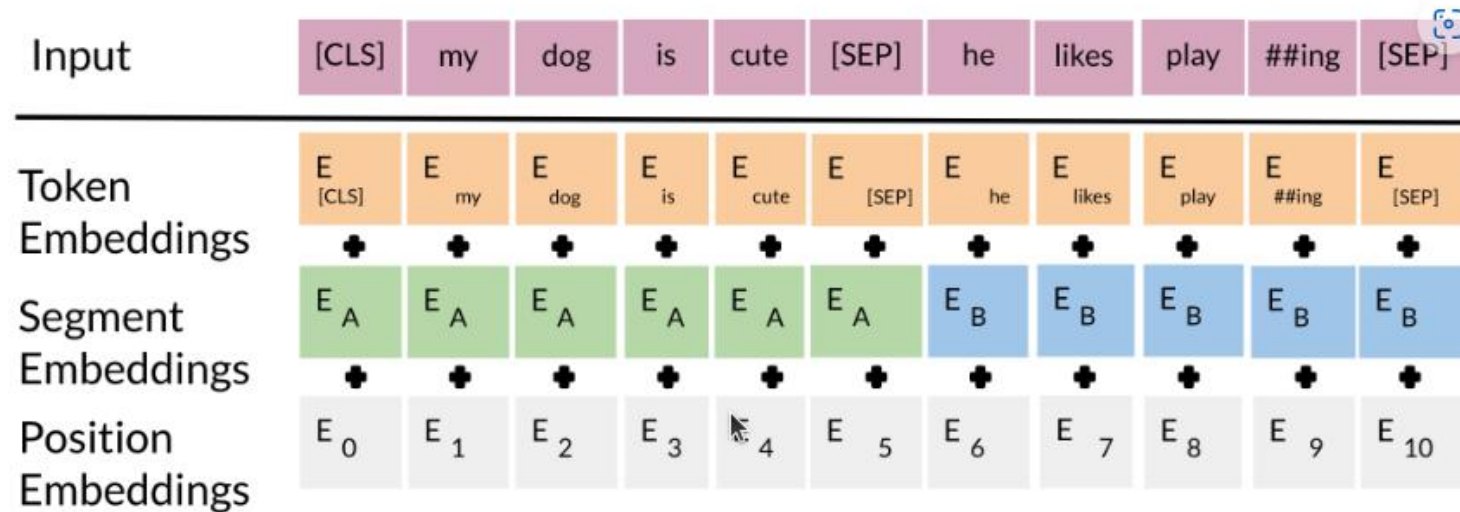


BERT | Training

- Choose 15% of the tokens at random: mask them 80% of the time, replace them with a random token 10% of the time, or keep as is 10% of the time.
- There could be multiple masked spans in a sentence
- Next sentence prediction is also used when pre-training.

BERT | Training

We will first start by visualizing the input.



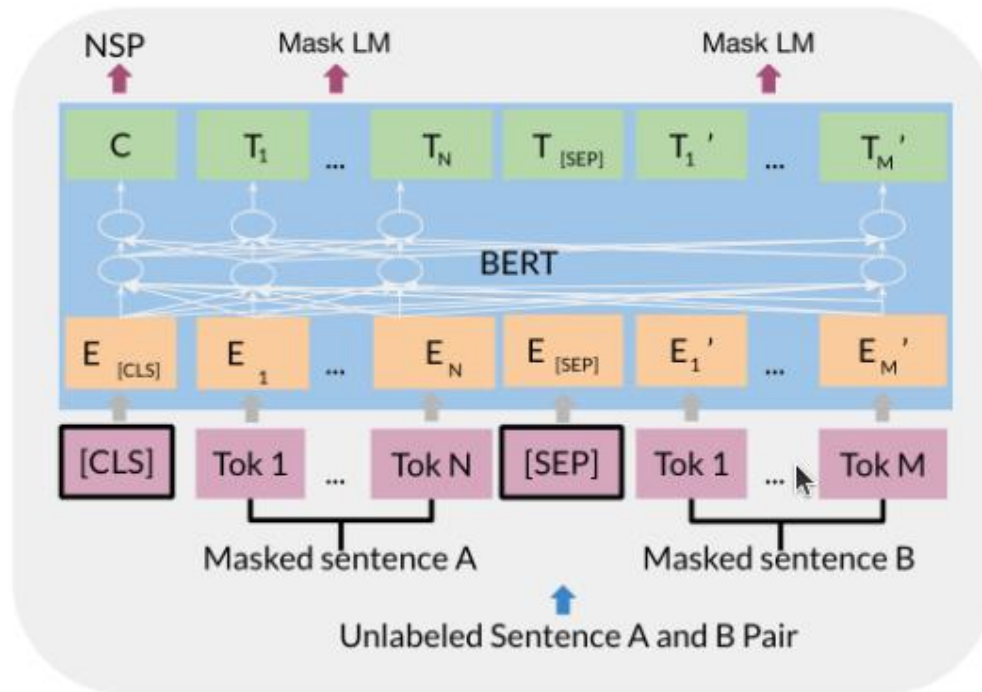
The input embeddings are the sum of the token embeddings, the segmentation embeddings and the position embeddings.

The input embeddings: you have a CLS token to indicate the beginning of the sentence and a sep to indicate the end of the sentence

The segment embeddings: allows you to indicate whether it is sentence a or b.

Positional embeddings: allows you to indicate the word's position in the sentence.

BERT | Training



- **[CLS]**: a special classification symbol added in front of every input
- **[SEP]**: a special separator token

The C token in the image above could be used for classification purposes. The unlabeled sentence A/B pair will depend on what you are trying to predict, it could range from question answering to sentiment. (in which case the second sentence could be just empty). The BERT objective is defined as follows:

BERT | Use Cases



Text encoding

**Text
Summarization**

**Response
Selection**

**Question
Answering**

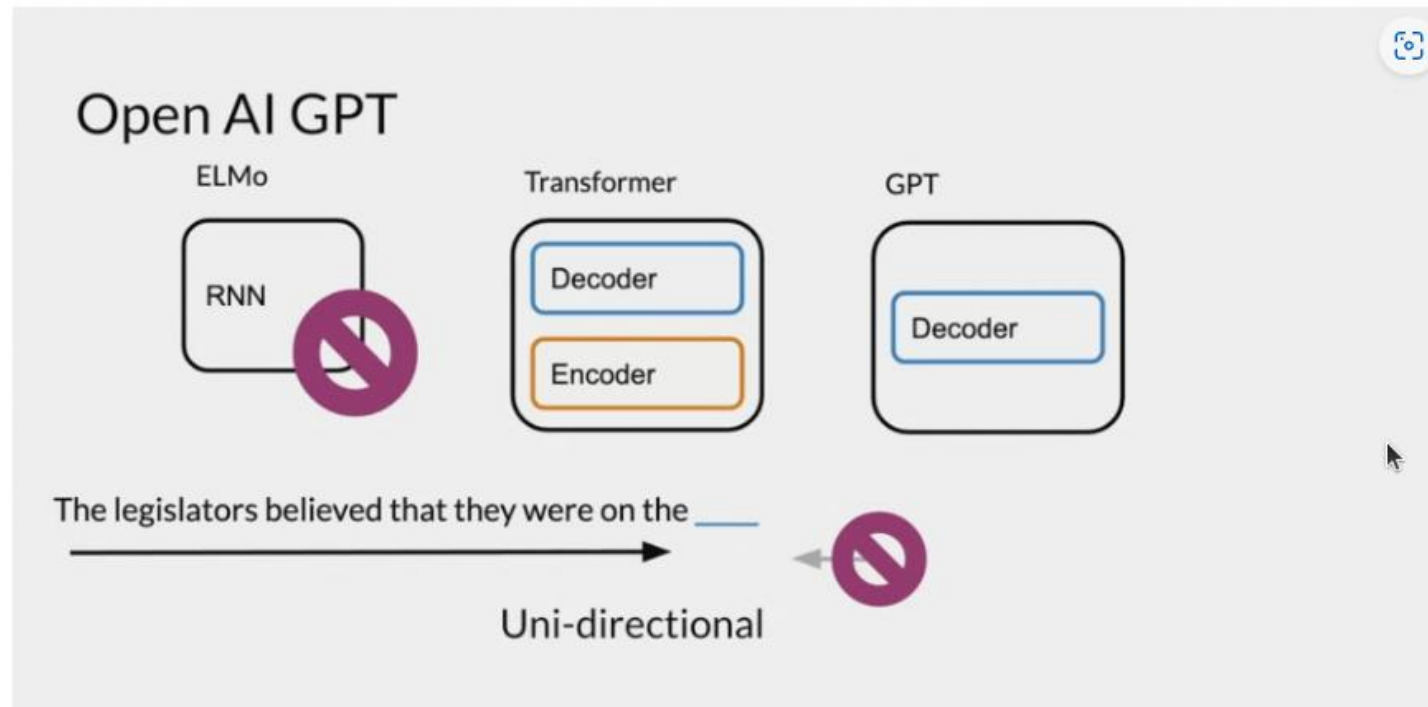
Similarity retrieval

And more...

Language Models | General

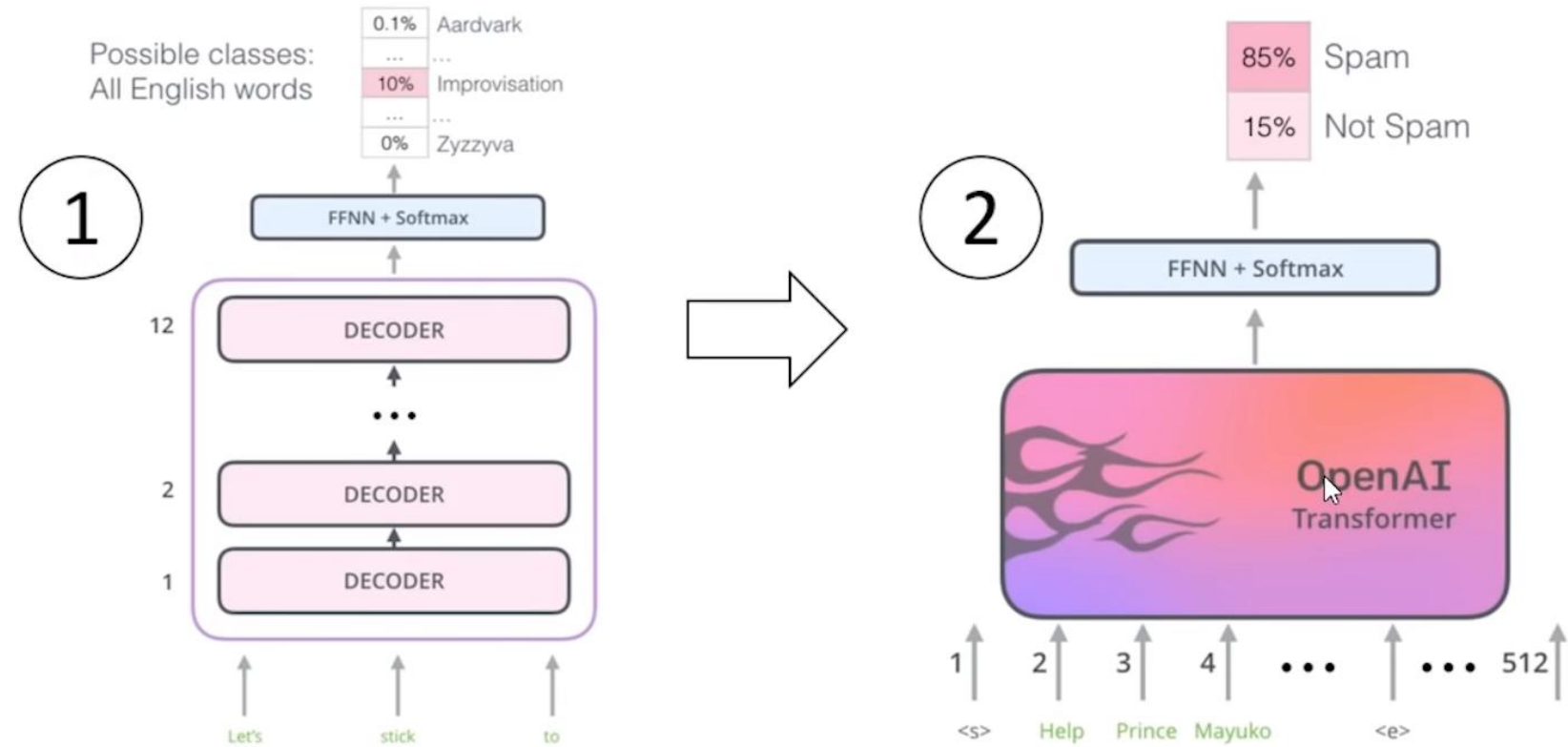
ELMo, GPT, BERT, T5

Around 2:50, Younes *incorrectly* mentions that ELMo is uni-directional. Please note, ELMo is **bi-directional**.



GPT

OpenAI Transformer



1. Pre-train a Transformer's decoder for language modeling
2. Train it on, for example, a sentence classification task

Hugging Face | A BERT & Transformers- Echo System

Hugging Face

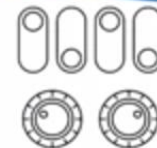
Transformers library

Use it with



Use it for

Applying state of the art
transformer models



Fine-tuning pretrained
transformer models

Hugging Face | Using Transformers

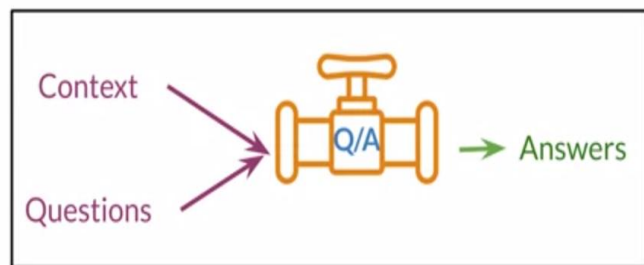
Pipelines



1. Pre-processing your inputs

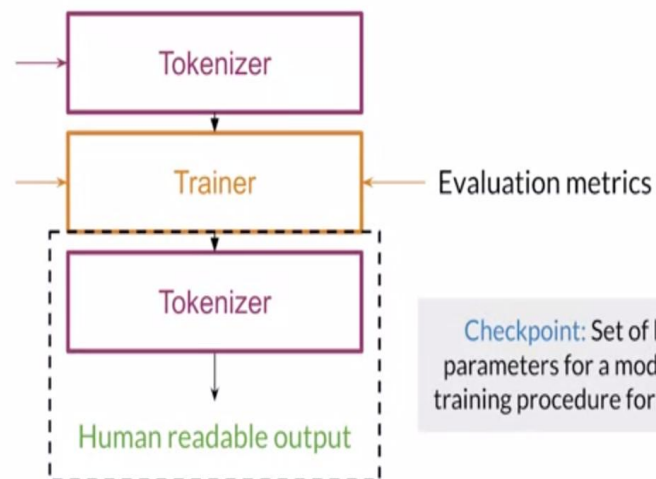
2. Running the model

3. Post-processing the outputs



Datasets:
One Thousand

Model Checkpoints:
More than 14 thousand



Checkpoint: Set of learned parameters for a model using a training procedure for some task

Model Checkpoints

Model Checkpoints:
More than 15 thousand
(and increasing)

Upload the architecture and weights with 1 line of code!

Model	Dataset	Name in 🤗
DistilBERT	Stanford Question Answering Dataset (SQuAD)	distilbert-base-cased-distilled-squad
BERT	Wikipedia and Book Corpus	bert-base-cased

Hugging Face | Code Example

- C4_W3_1_Question_Answering_with_BERT_and_HuggingFace_Pytorch_tydiqa.ipynb - Colaboratory

<https://colab.research.google.com/drive/1O4LvdhHw6Zx7Kd43HK-p5a1rtsHUEia5#scrollTo=oDG5fgap-N7I>

GPT3 , and ChatGPT

- GPT3 - is a Language Model which may be used in variety of Language agnostic applications
- ChatGPT - is a specifically designed for conversational Language Applications.

GPT-3

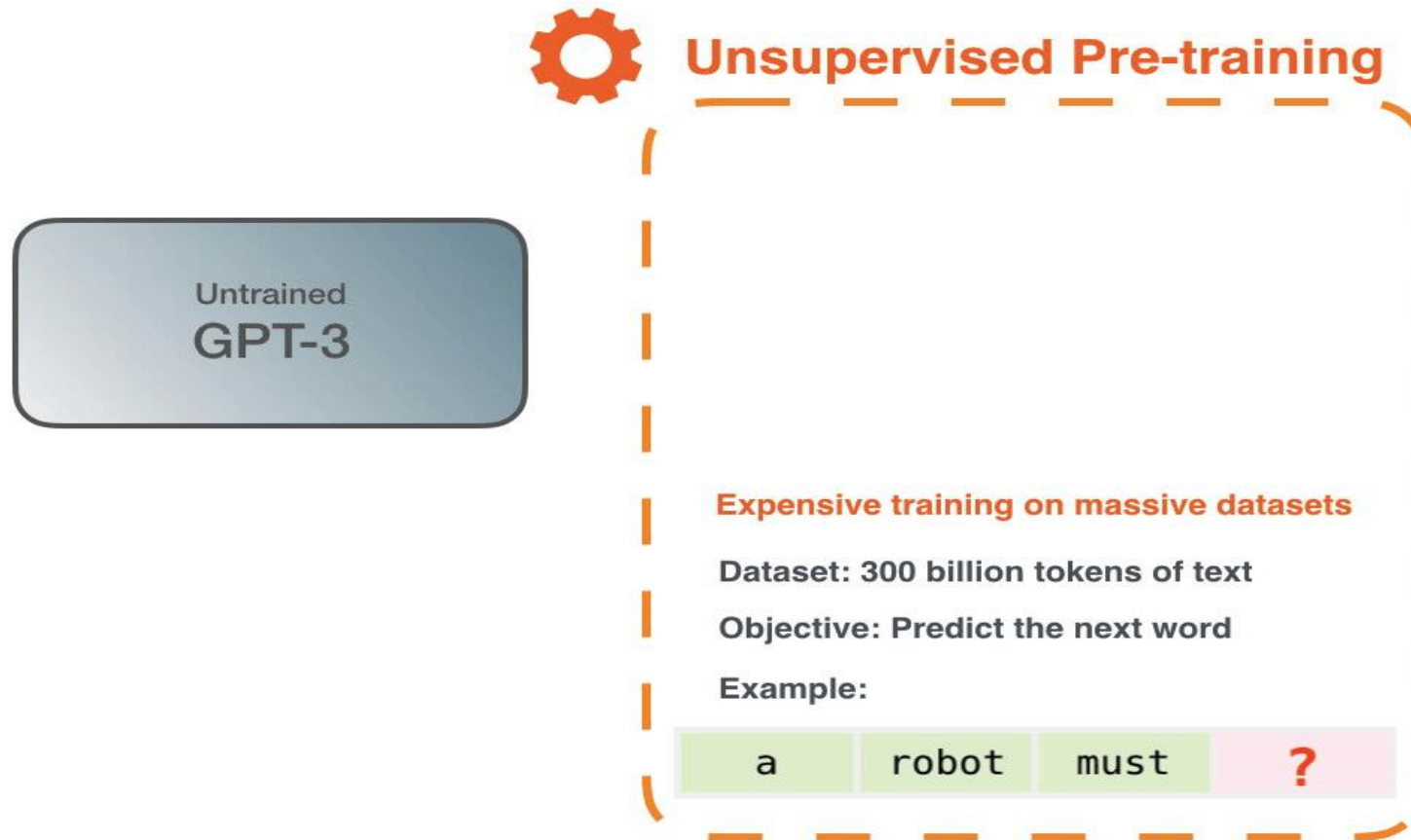
Input Prompt:

Recite the first law of robotics

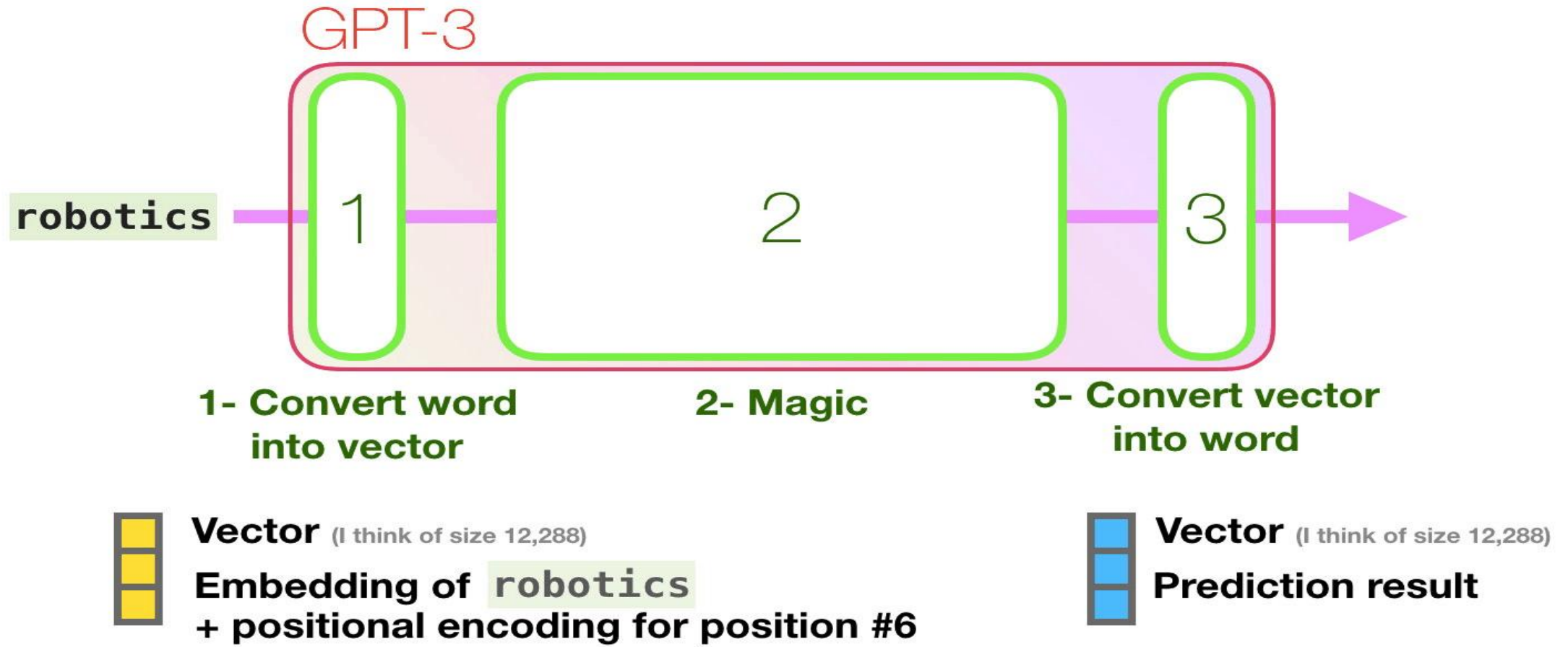


Output:

GPT3

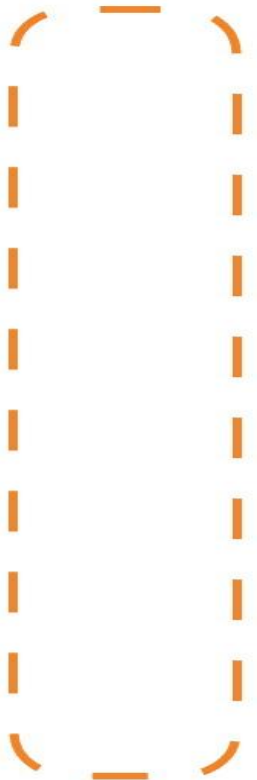


GPT3



GPT - Fine Tuning

Pre-training



Fine-tuning

Additional training to become better at a certain task

Example: English to French Translation

