



Natural Language Processing: Transformers

HSE Faculty of Computer Science
Machine Learning and Data-Intensive Systems

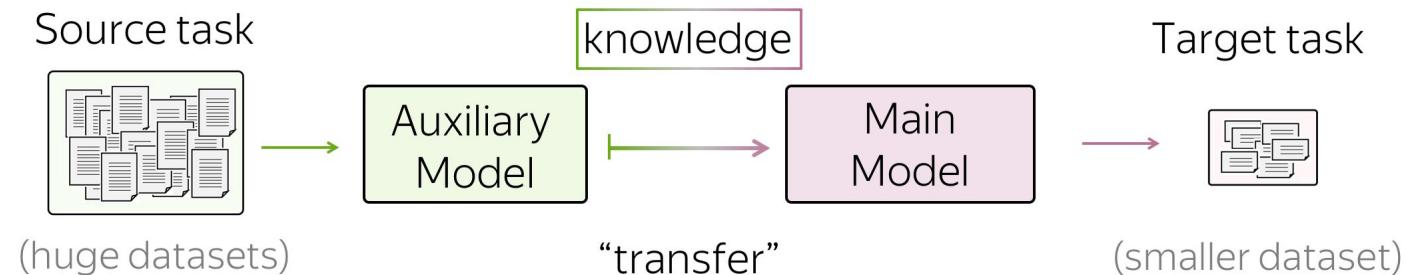
Murat Khazhgeriev



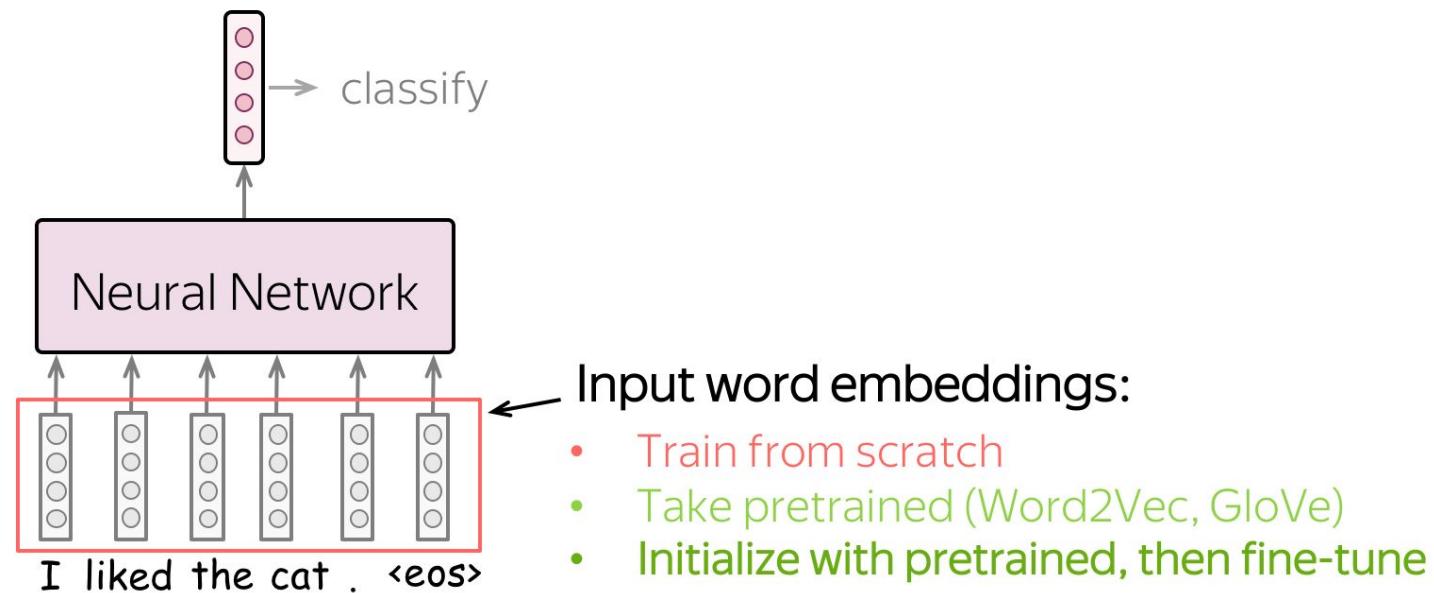
Table of Content

- **The power of transfer learning**
- From word-specific to contextual embeddings
- Transformer architecture overview
- BERT
- GPT

Training a model on a general task can benefit a downstream one



Training a model on a general task can benefit a downstream one



Training a model on a general task can benefit a downstream one

- Train from scratch

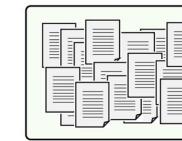
What they will know:



May be not enough to learn relationships between words

- Take pretrained (Word2Vec, GloVe)

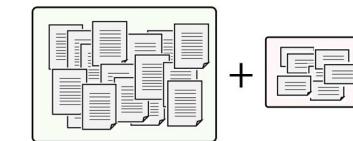
What they will know:



Know relationships between words, but are **not** specific to the task

- Initialize with pretrained, then fine-tune

What they will know:



Know relationships between words and adapted for the task

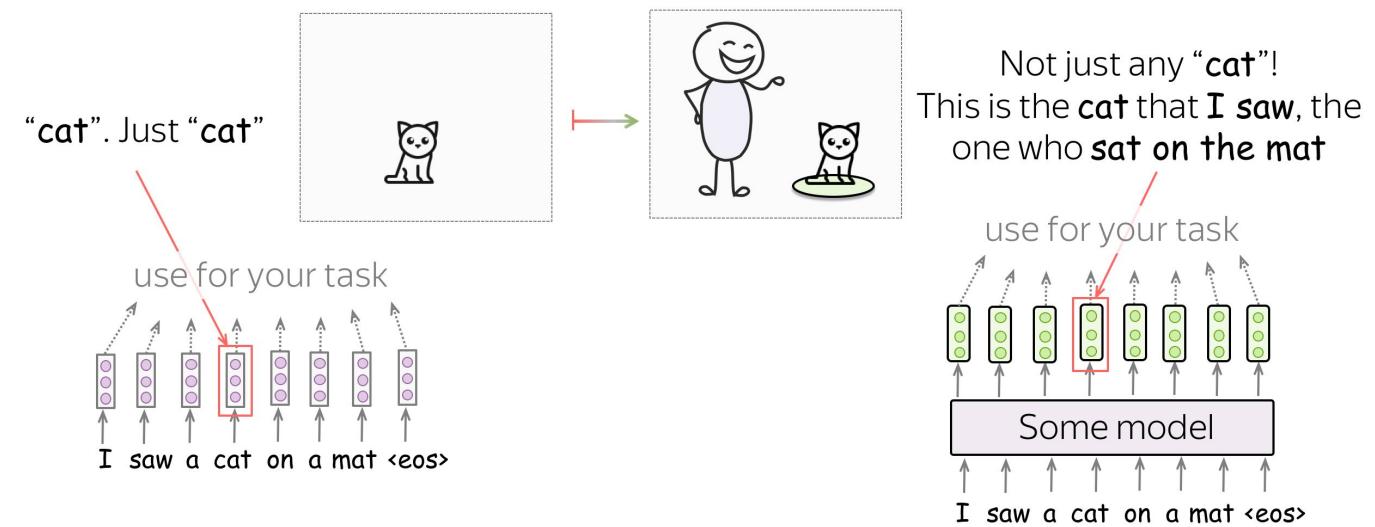
"Transfer" knowledge from a huge unlabeled corpus to your task-specific model



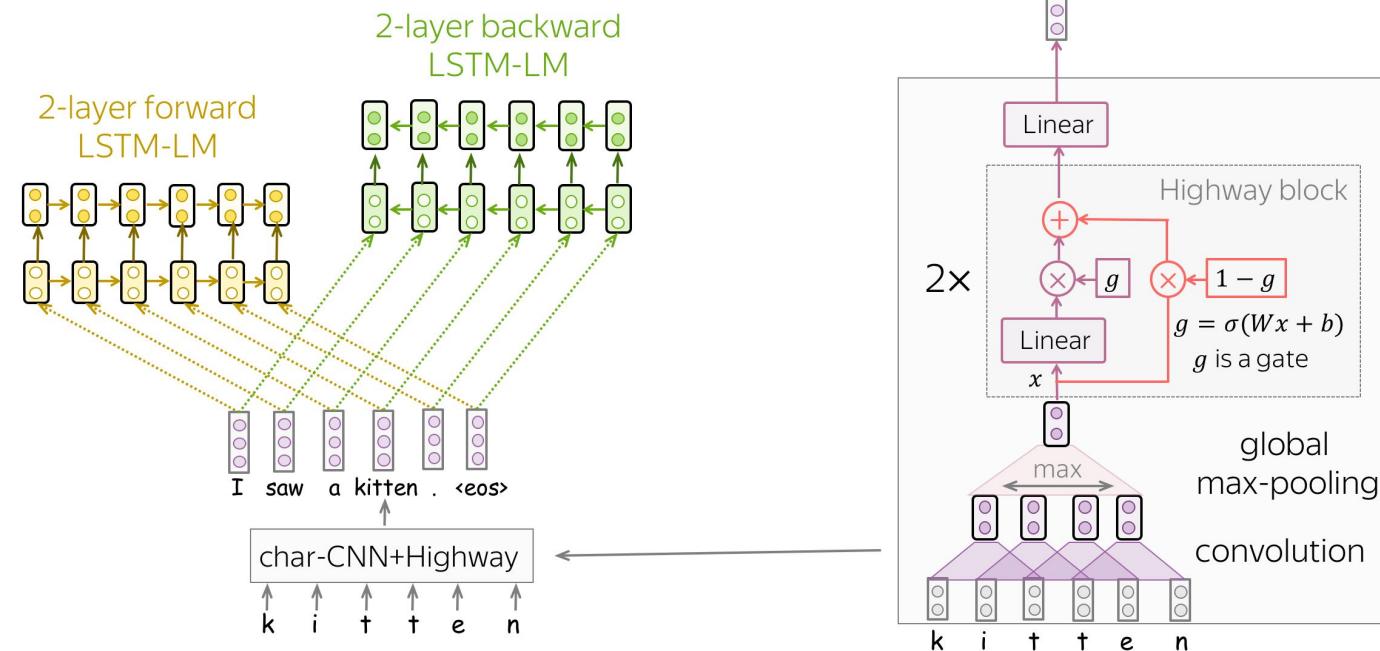
Table of Content

- The power of transfer learning
- **From word-specific to contextual embeddings**
- Transformer architecture overview
- BERT
- GPT

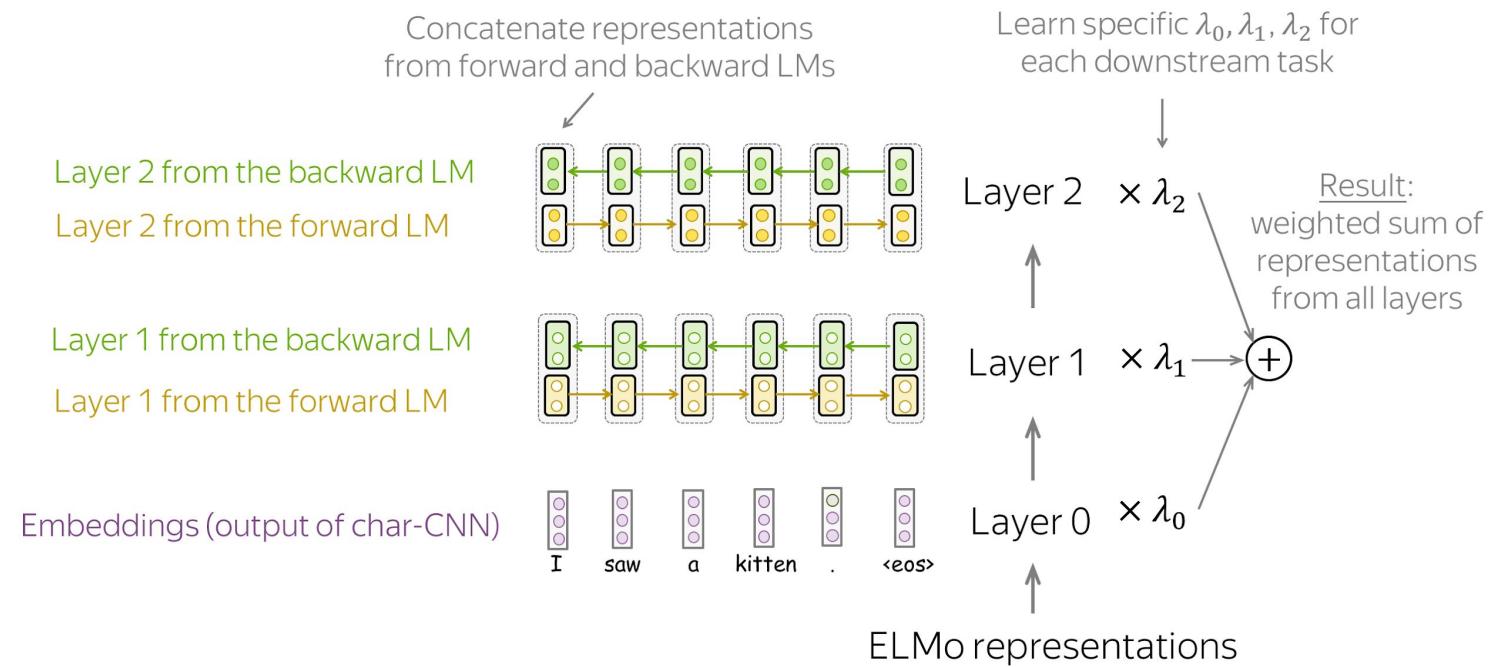
Not just a cat, but the cat!



Train a “translator” from word-specific to “contextual” space



Multiple layers to capture low-level and high-level context



From an embedding generator to a universal model

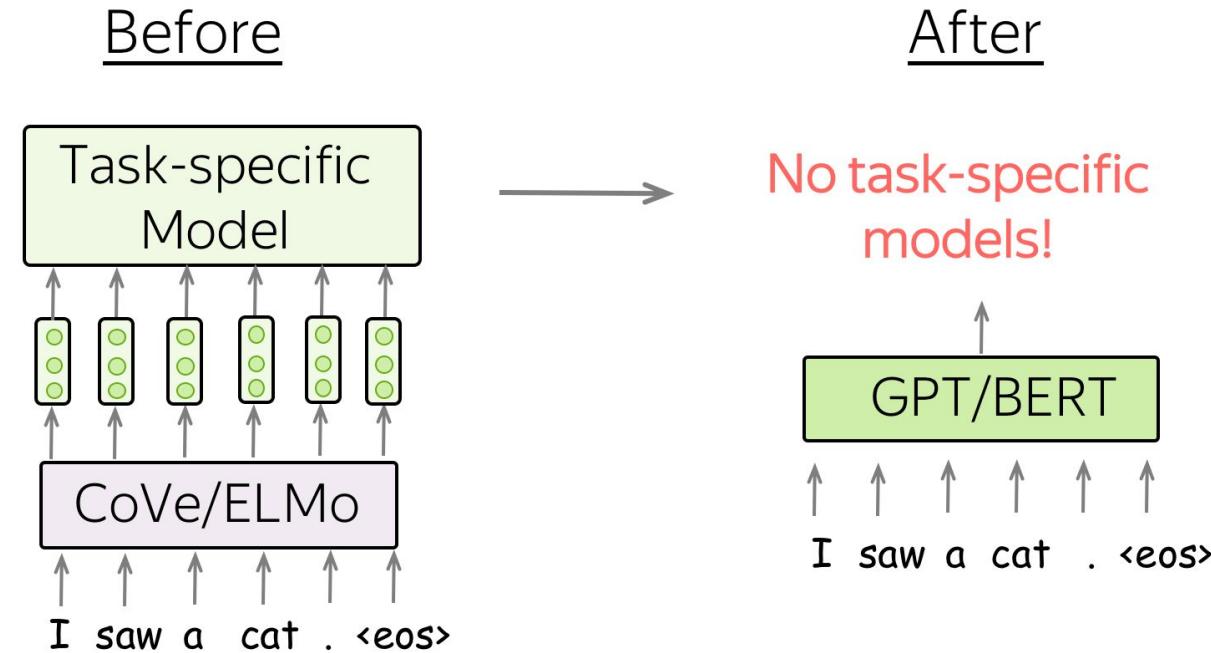


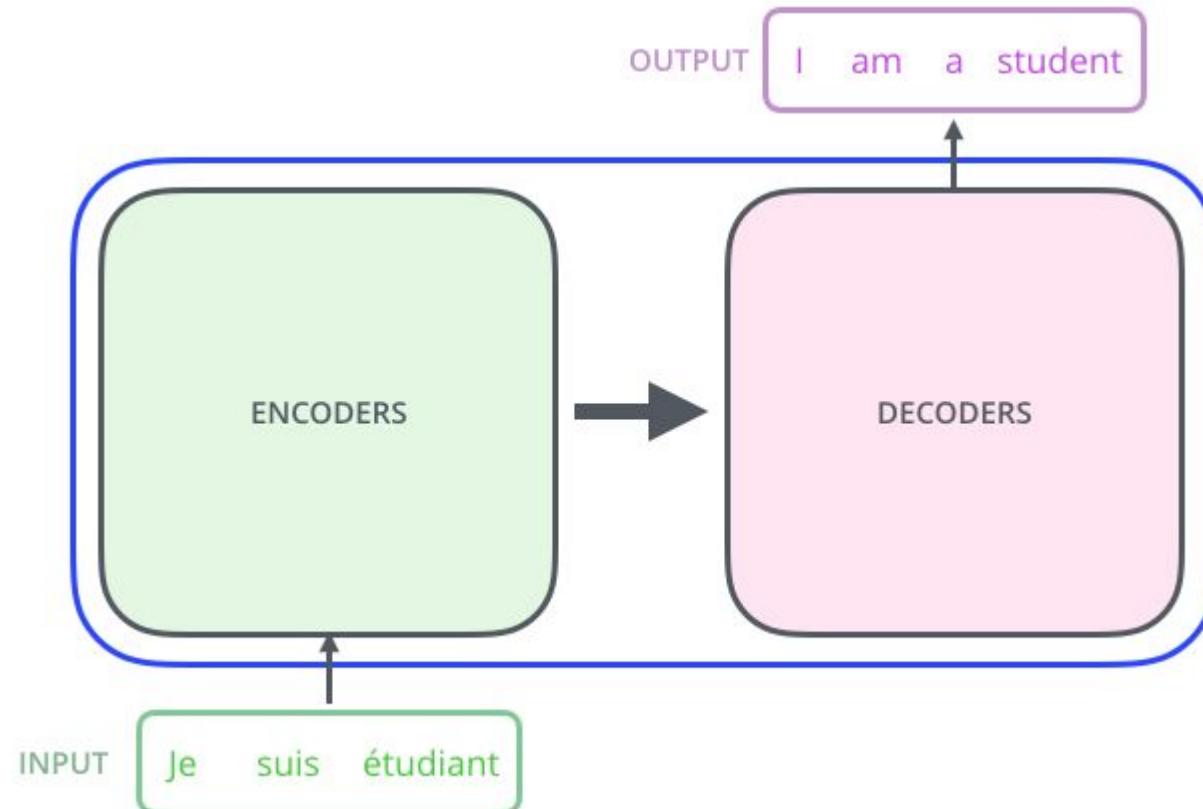
Table of Content

- The power of transfer learning
- From word-specific to contextual embeddings
- **Transformer architecture overview**
- BERT
- GPT

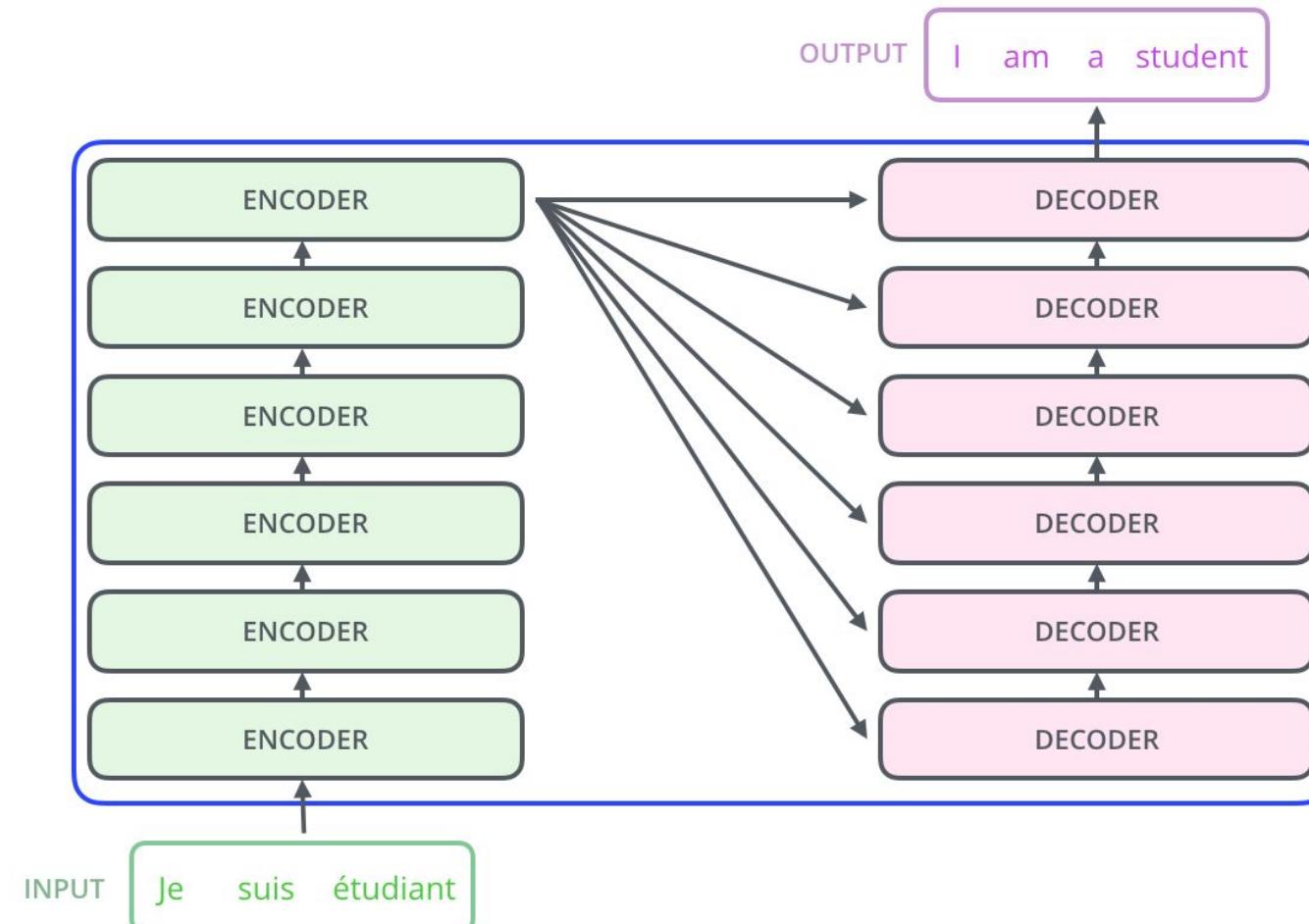
Transformer is an example of Encoder-Decoder architecture



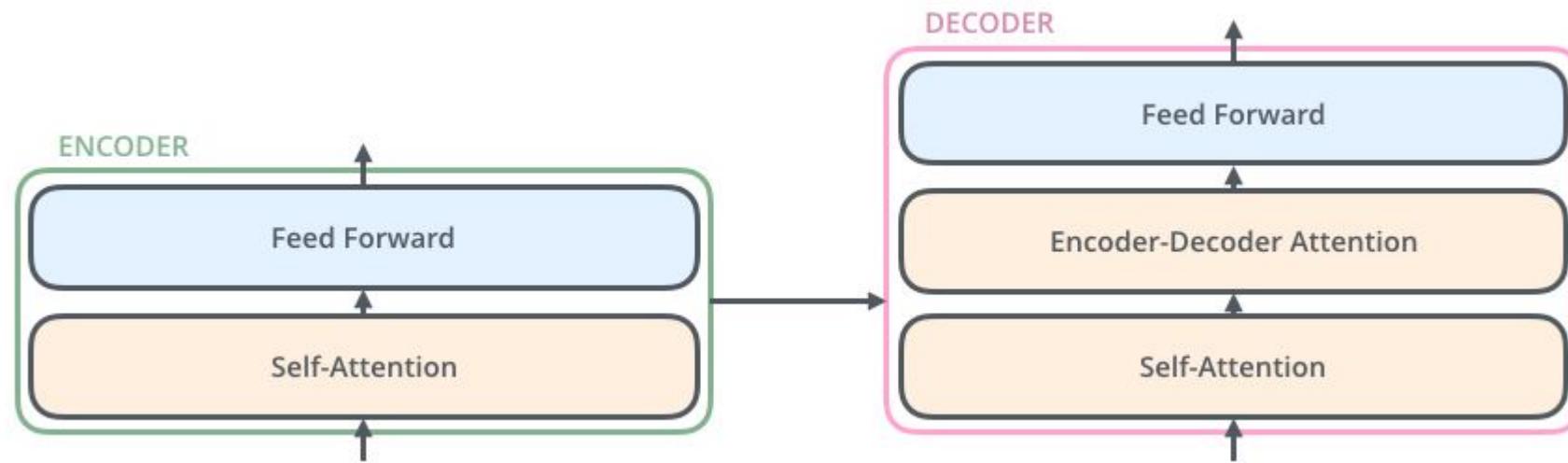
Transformer is an example of Encoder-Decoder architecture



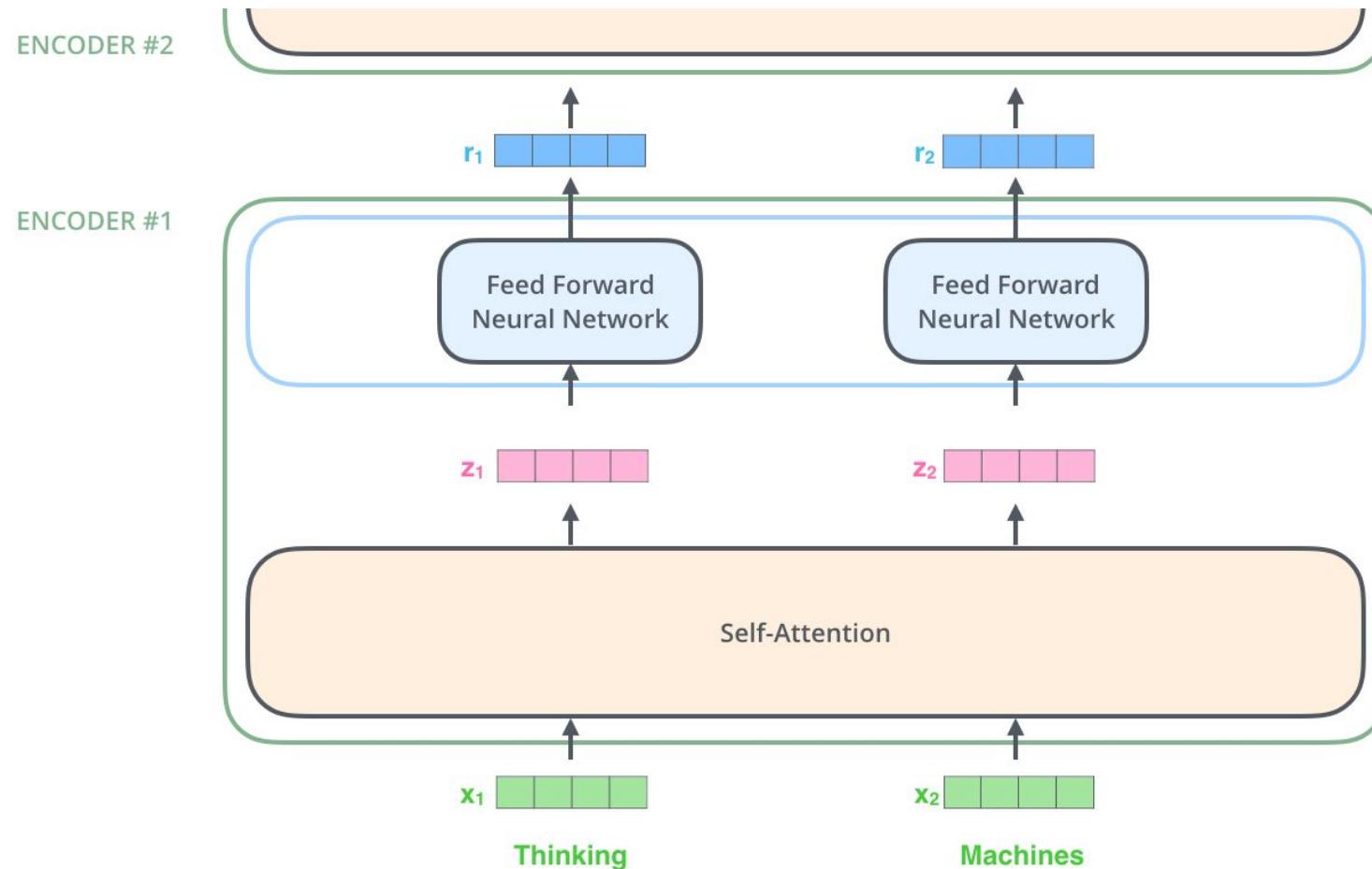
Transformer is an example of Encoder-Decoder architecture



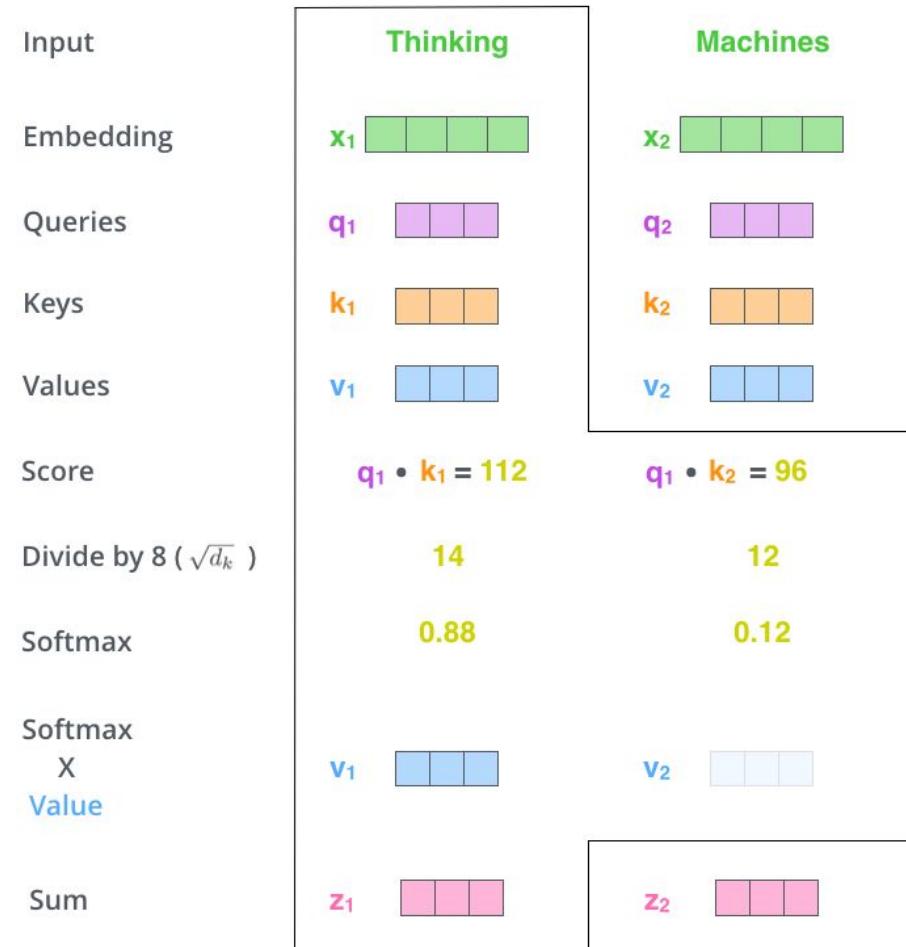
Always two there are. No more, no less. An Attention and an FFN. © Yoda



Always two there are. No more, no less. An Attention and an FFN. © Yoda



Inside a Self-Attention



Inside a Self-Attention: Matrix View

$$\mathbf{X} \quad \times \quad \mathbf{W^Q} \quad = \quad \mathbf{Q}$$

A diagram illustrating the computation of Query (Q) vectors. An input matrix \mathbf{X} (green, 4x4) is multiplied by a weight matrix $\mathbf{W^Q}$ (purple, 4x4). The result is a query matrix \mathbf{Q} (purple, 4x4).

$$\mathbf{X} \quad \times \quad \mathbf{W^K} \quad = \quad \mathbf{K}$$

A diagram illustrating the computation of Key (K) vectors. An input matrix \mathbf{X} (green, 4x4) is multiplied by a weight matrix $\mathbf{W^K}$ (orange, 4x4). The result is a key matrix \mathbf{K} (orange, 4x4).

$$\mathbf{X} \quad \times \quad \mathbf{W^V} \quad = \quad \mathbf{V}$$

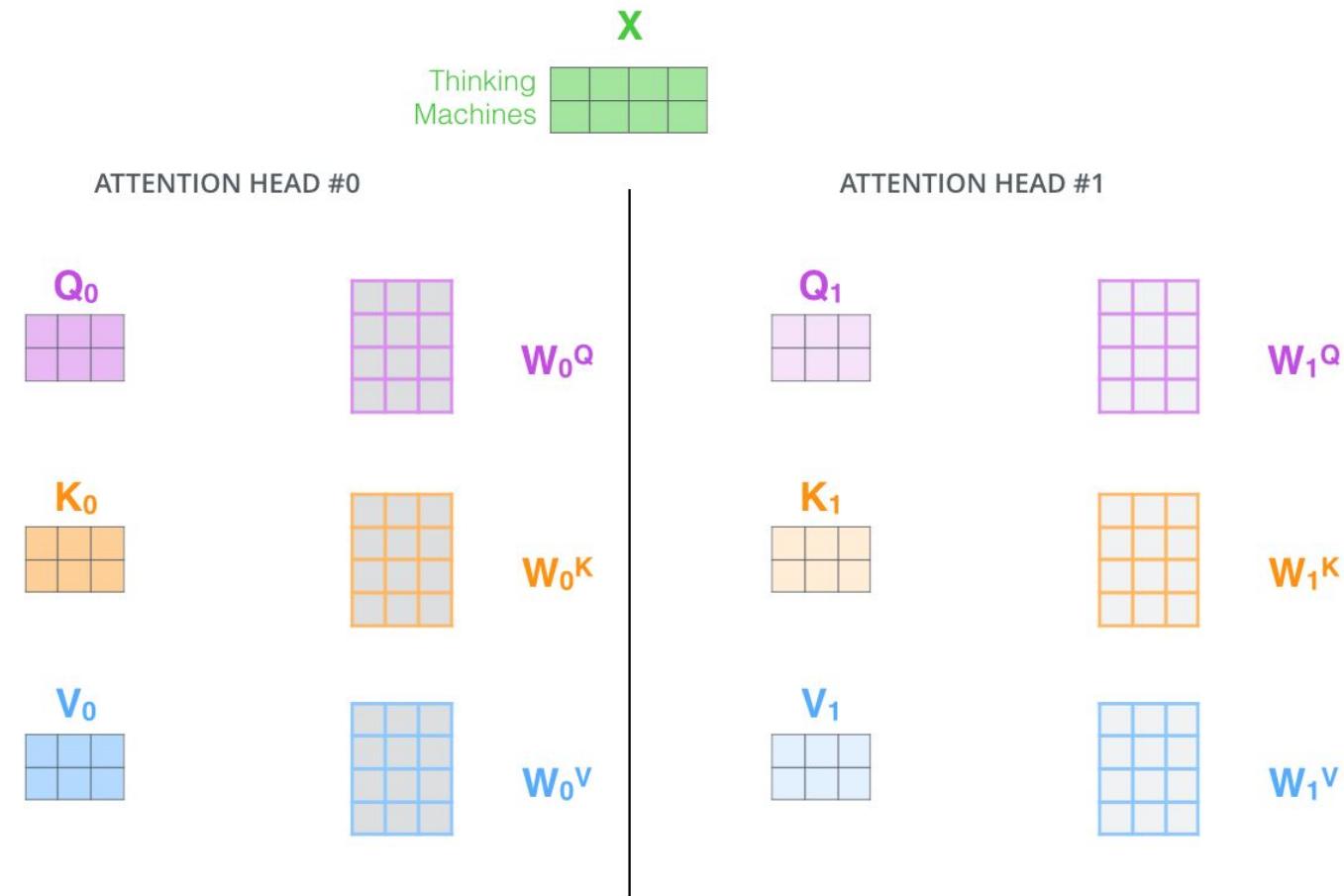
A diagram illustrating the computation of Value (V) vectors. An input matrix \mathbf{X} (green, 4x4) is multiplied by a weight matrix $\mathbf{W^V}$ (blue, 4x4). The result is a value matrix \mathbf{V} (blue, 4x4).

Inside a Self-Attention: Matrix View

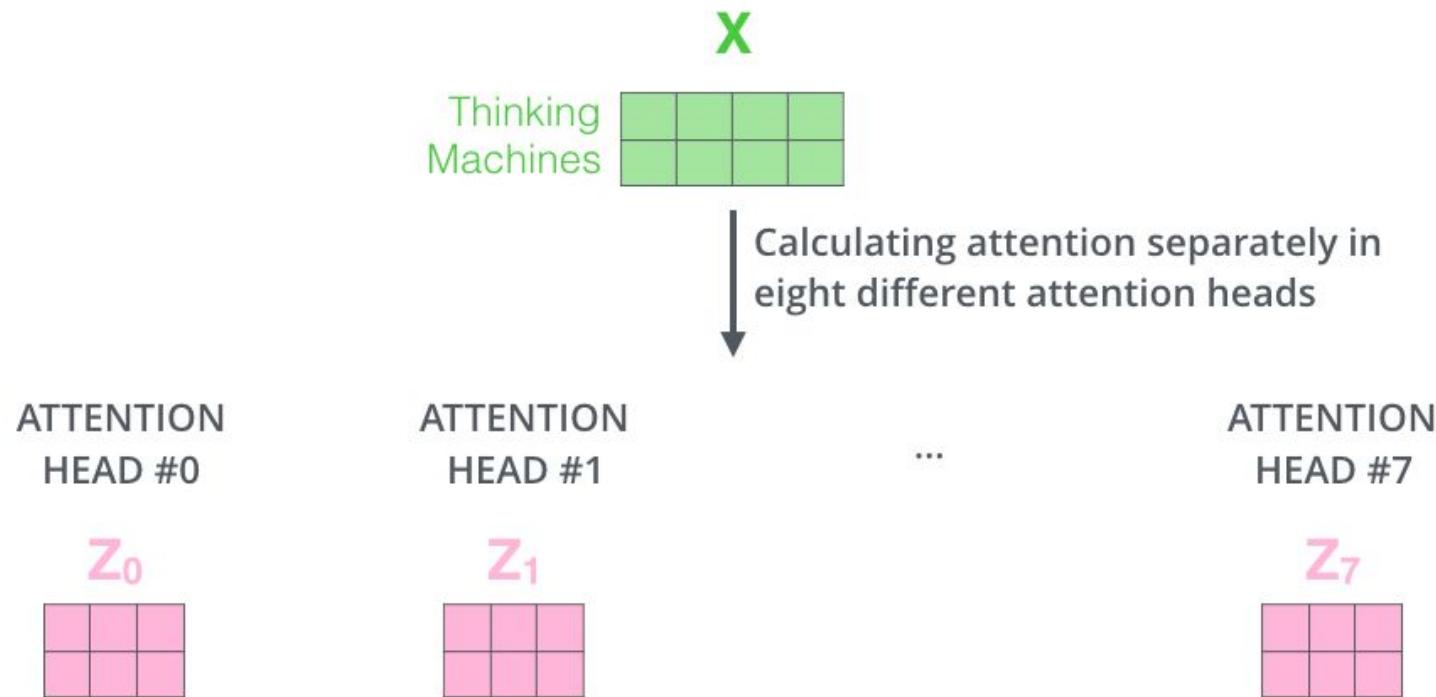
$$\text{softmax}\left(\frac{\mathbf{Q} \times \mathbf{K}^T}{\sqrt{d_k}}\right) \mathbf{V} = \mathbf{Z}$$

The diagram illustrates the computation of self-attention weights. It shows three input matrices: \mathbf{Q} (purple, 2x3), \mathbf{K}^T (orange, 3x2), and \mathbf{V} (blue, 2x3). The product of \mathbf{Q} and \mathbf{K}^T is scaled by $\sqrt{d_k}$ before being passed through a softmax function to produce the attention weights matrix \mathbf{Z} (pink, 2x3).

A beast with many heads



A beast with many heads



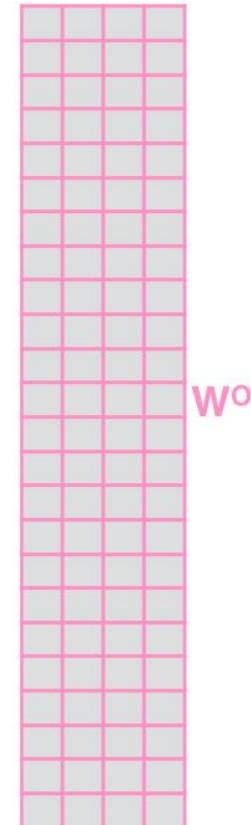
A beast with many heads

1) Concatenate all the attention heads



2) Multiply with a weight matrix W^o that was trained jointly with the model

X

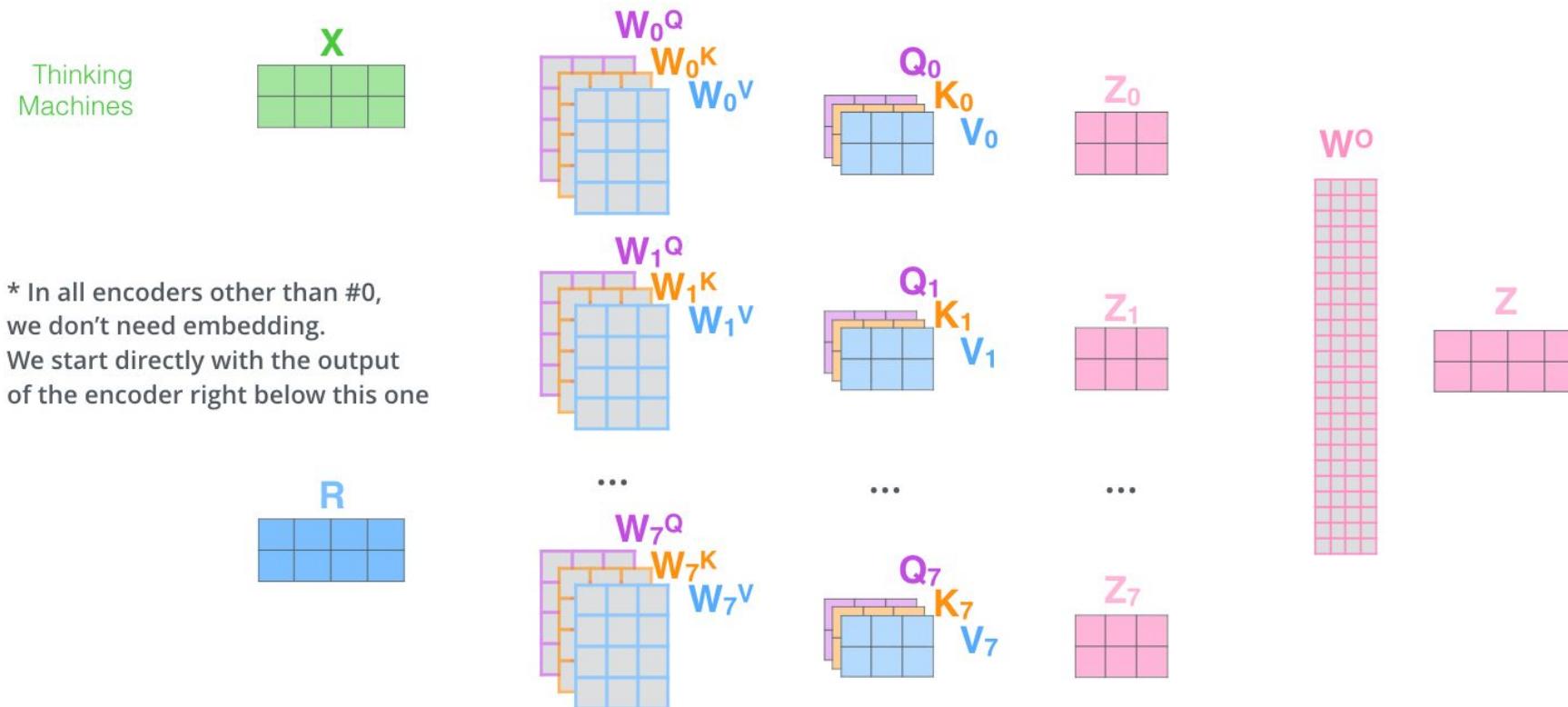


3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

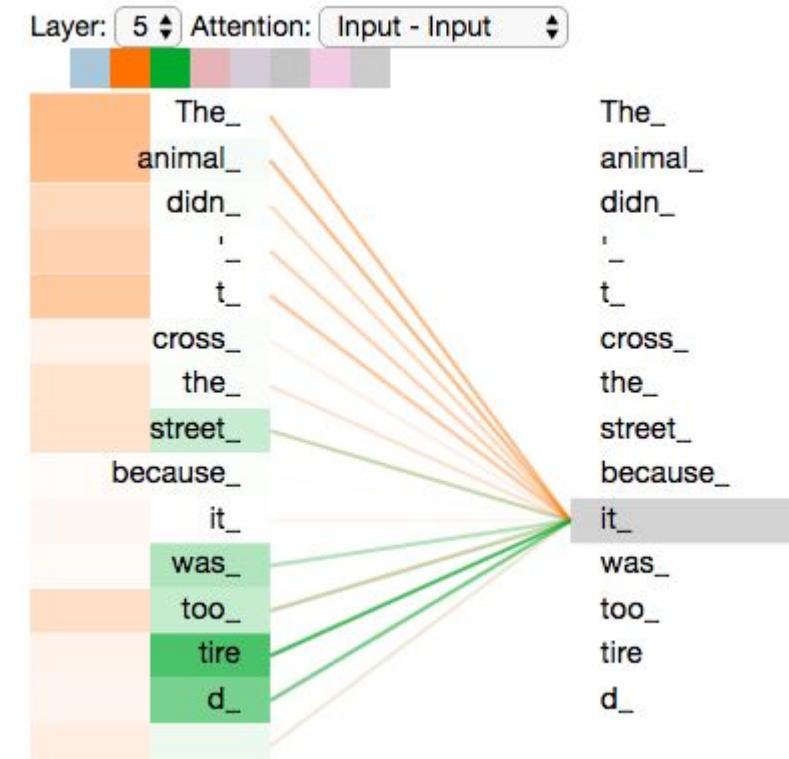
$$= \begin{matrix} Z \\ \hline \text{---} \\ \text{---} \end{matrix}$$

A multi-head attention overview

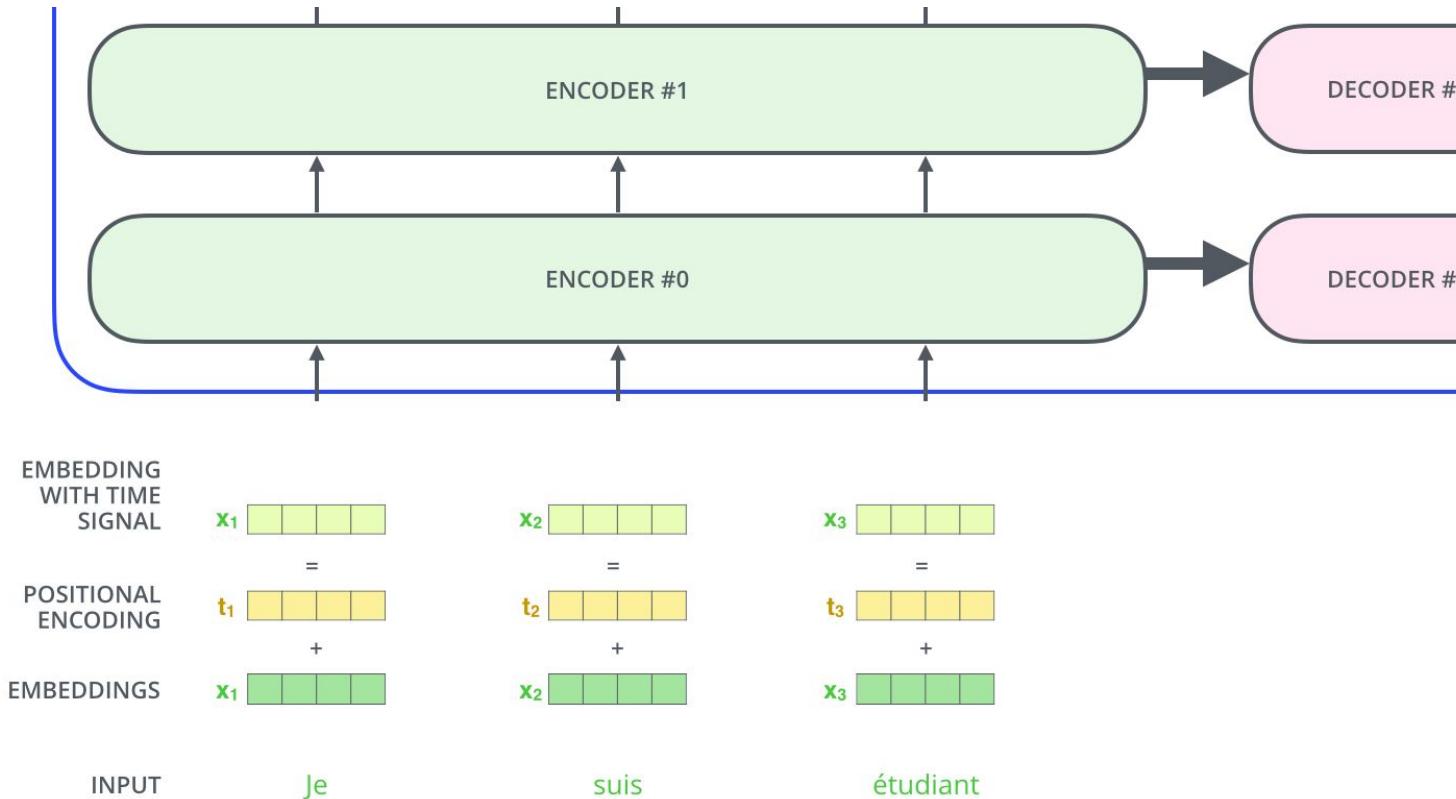
- 1) This is our input sentence* X
- 2) We embed each word*
- 3) Split into 8 heads. We multiply X or R with weight matrices
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^o to produce the output of the layer



Each head focuses on a specific representation



As opposed to RNNs, Transformers do not track the position implicitly

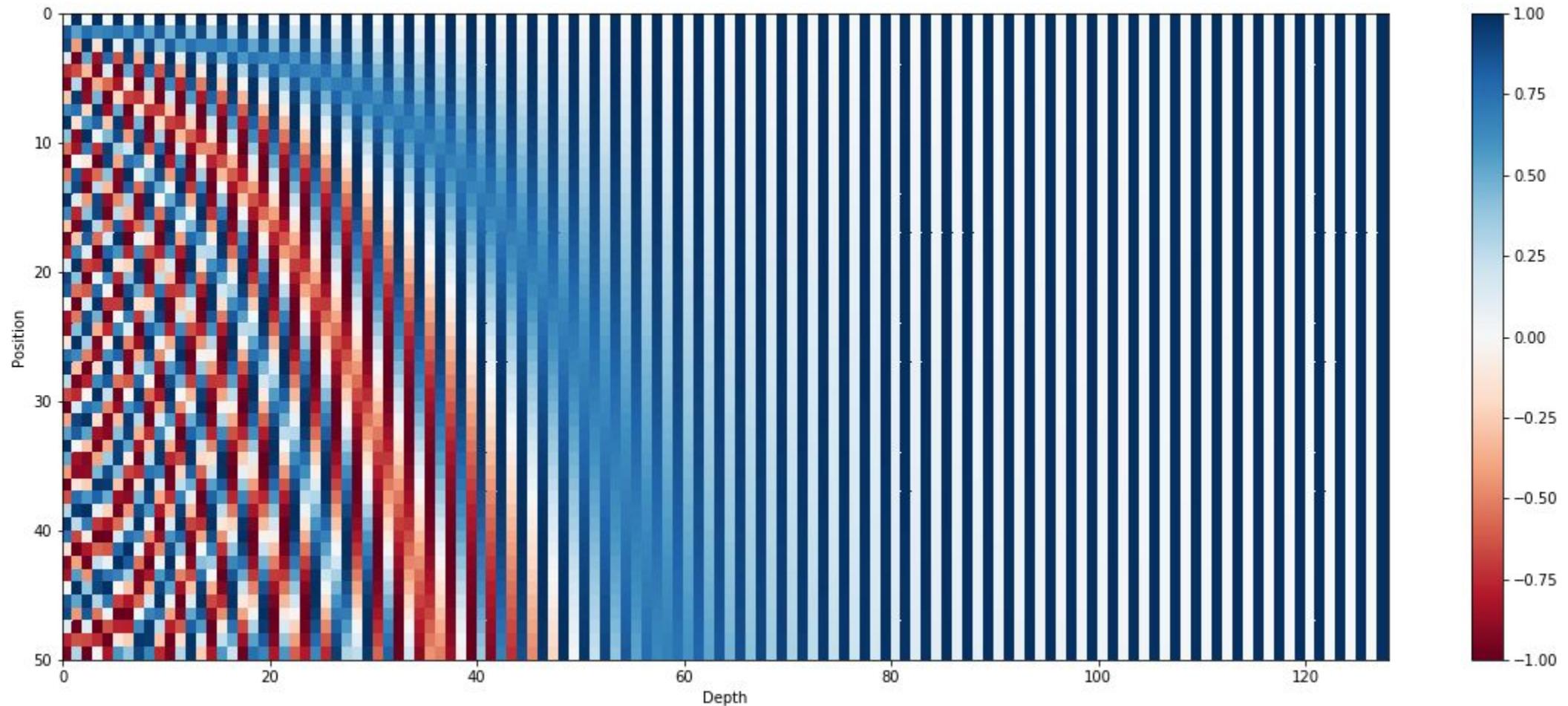


Absolute Positional Encoding

$$\text{PE}(pos, 2i) = \sin\left(pos/10000^{2i/d_{model}}\right)$$

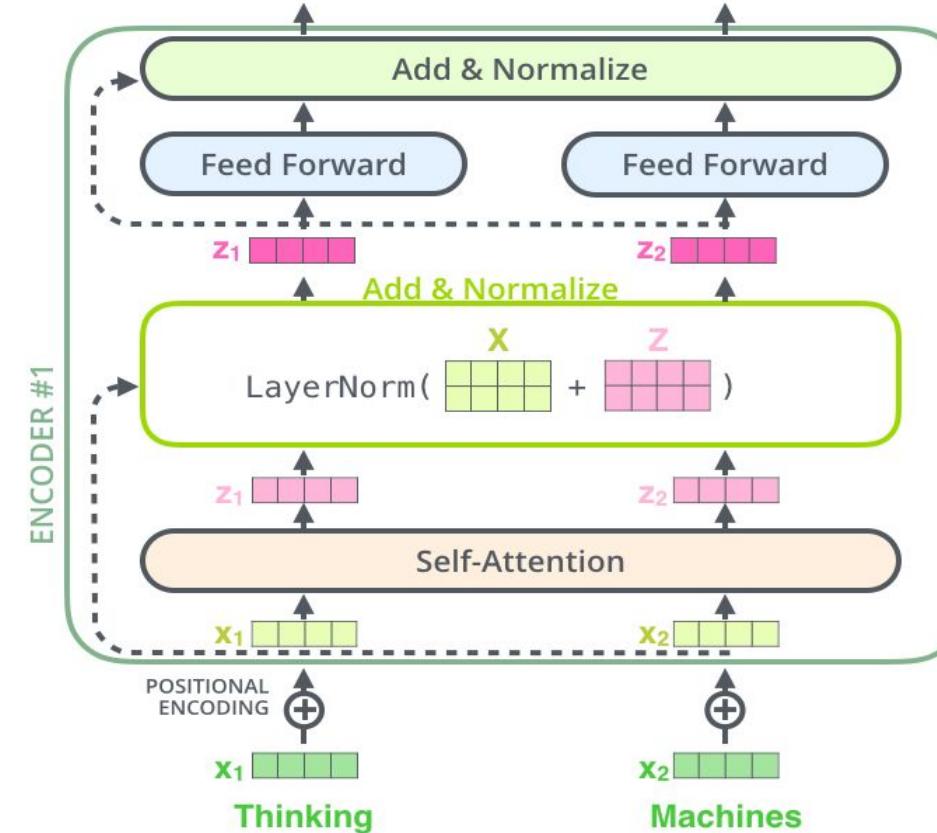
$$\text{PE}(pos, 2i + 1) = \cos\left(pos/10000^{2i/d_{model}}\right)$$

What it looks like

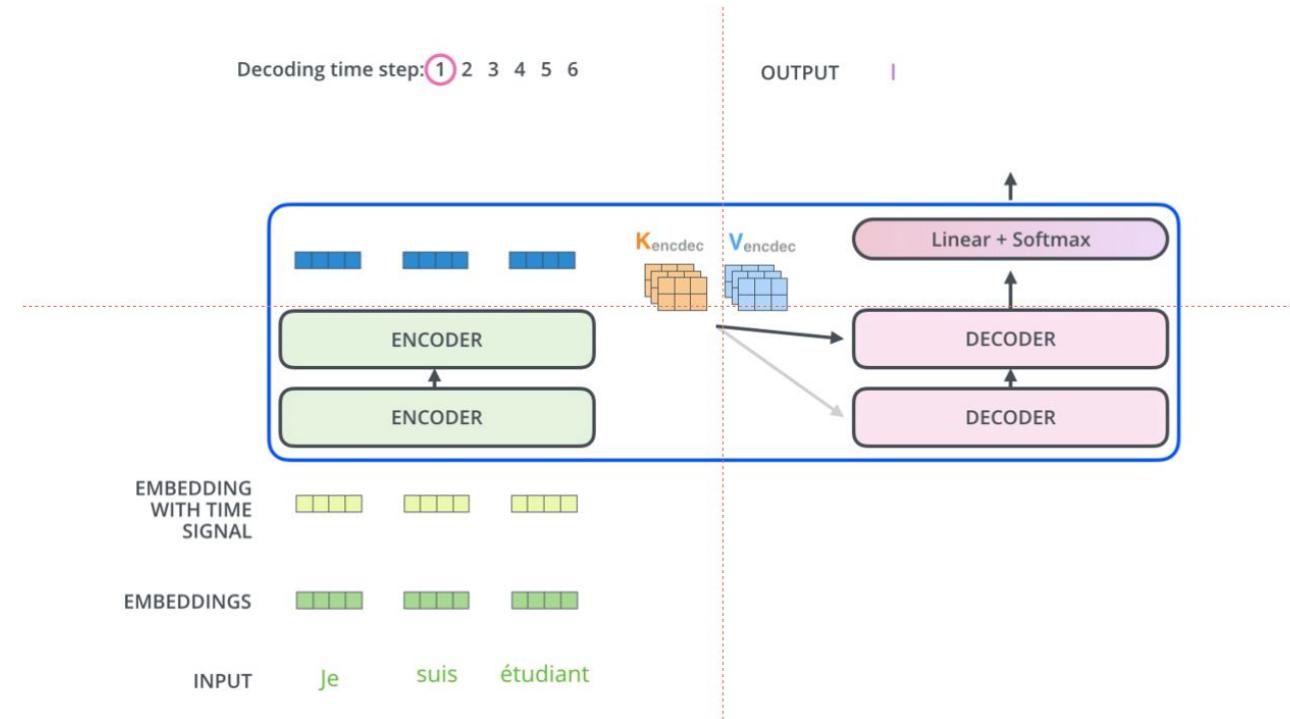


Source: https://kazemnejad.com/blog/transformer_architecture_positional_encoding/

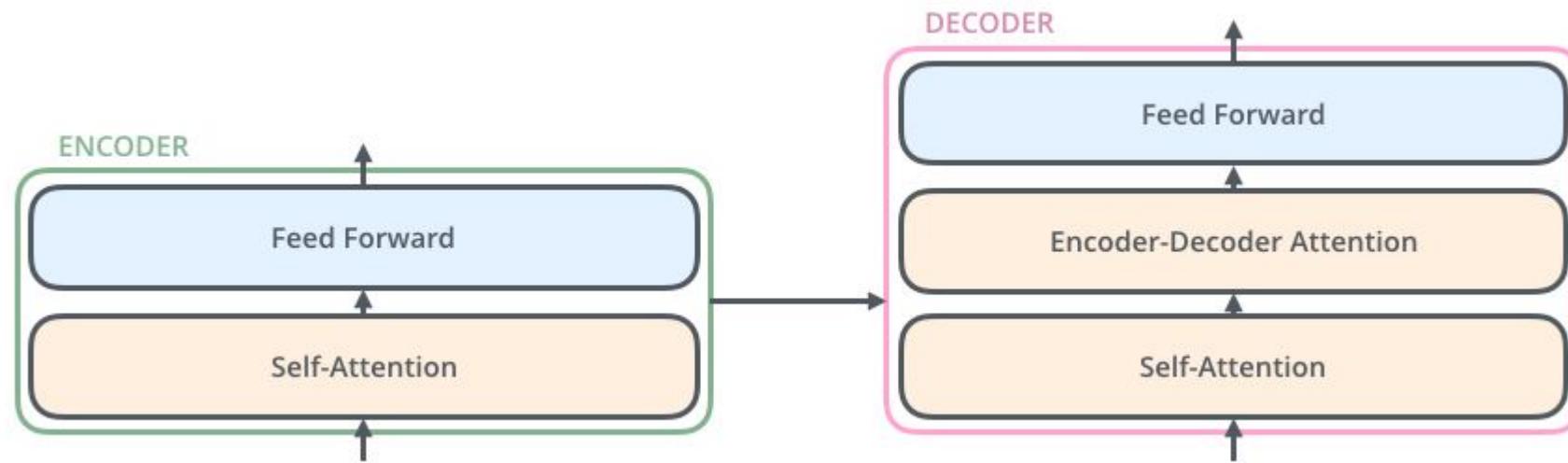
Skip Connection and Layer Normalization for a robust training



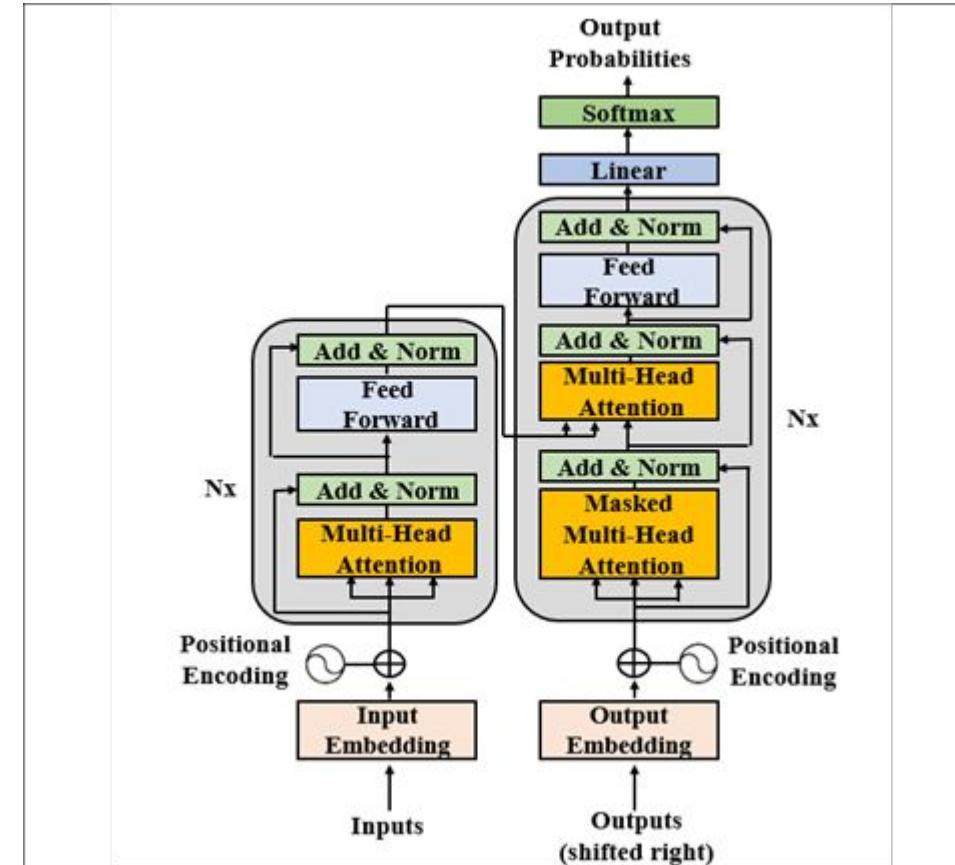
Bringing it all together



Always two there are. No more, no less. An Attention and an FFN. © Yoda



Attention is all you need



Attention is all you need (but not really)

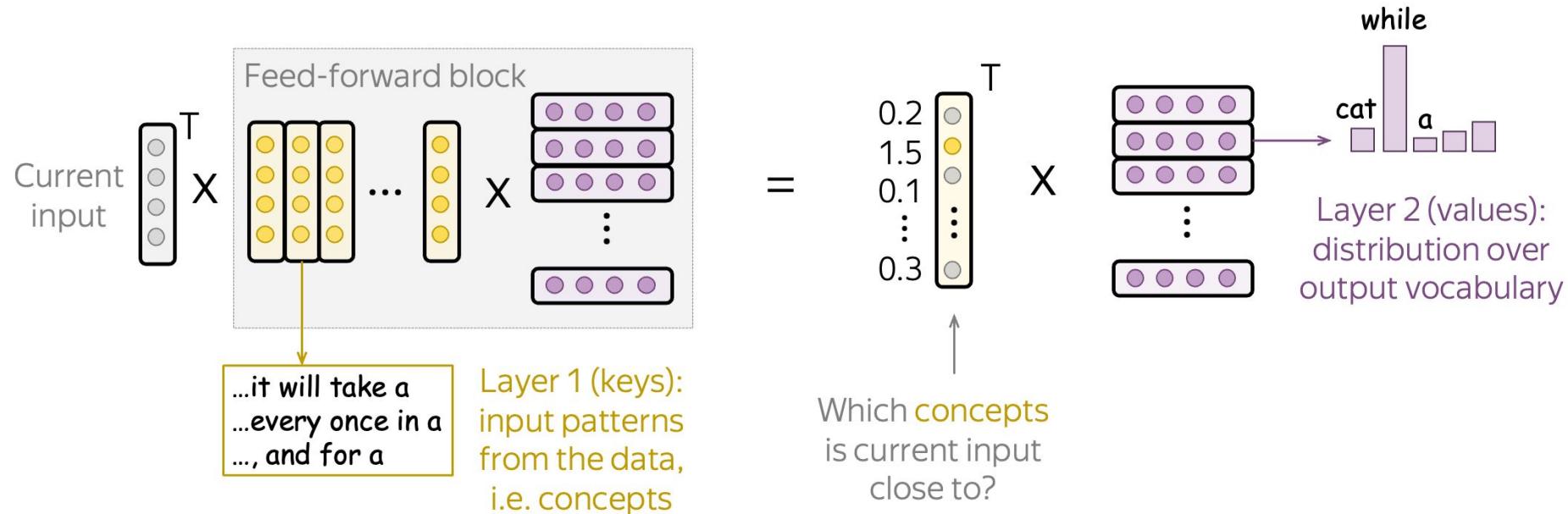
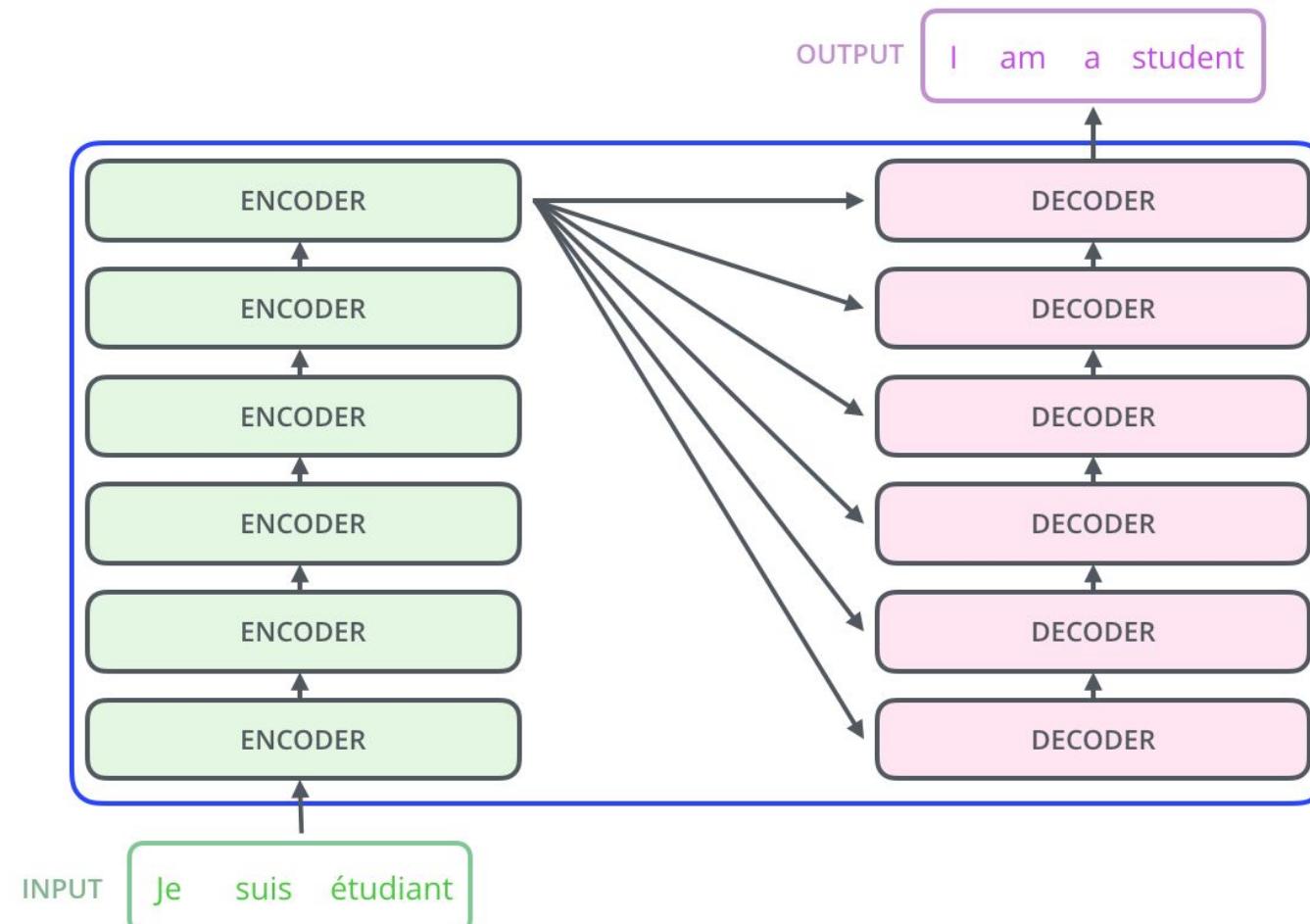


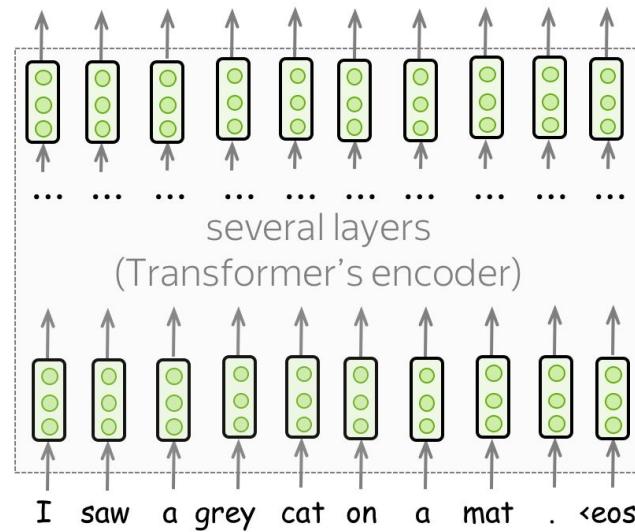
Table of Content

- The power of transfer learning
- From word-specific to contextual embeddings
- Transformer architecture overview
- **BERT**
- GPT

BERT is just an Encoder part



Using encoder as an embedding generator



Model architecture:

- Transformer's encoder

What is special about it:

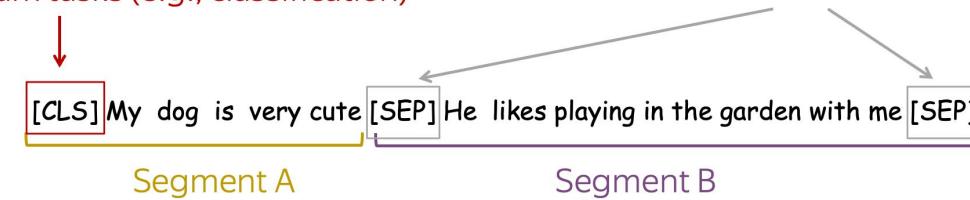
- Training objectives
 - MLM: Masked language modeling
 - NSP: Next sentence prediction
- The way it is used
 - No task-specific models

Objective one: tell whether the two sequences are consecutive

[CLS]: Special token

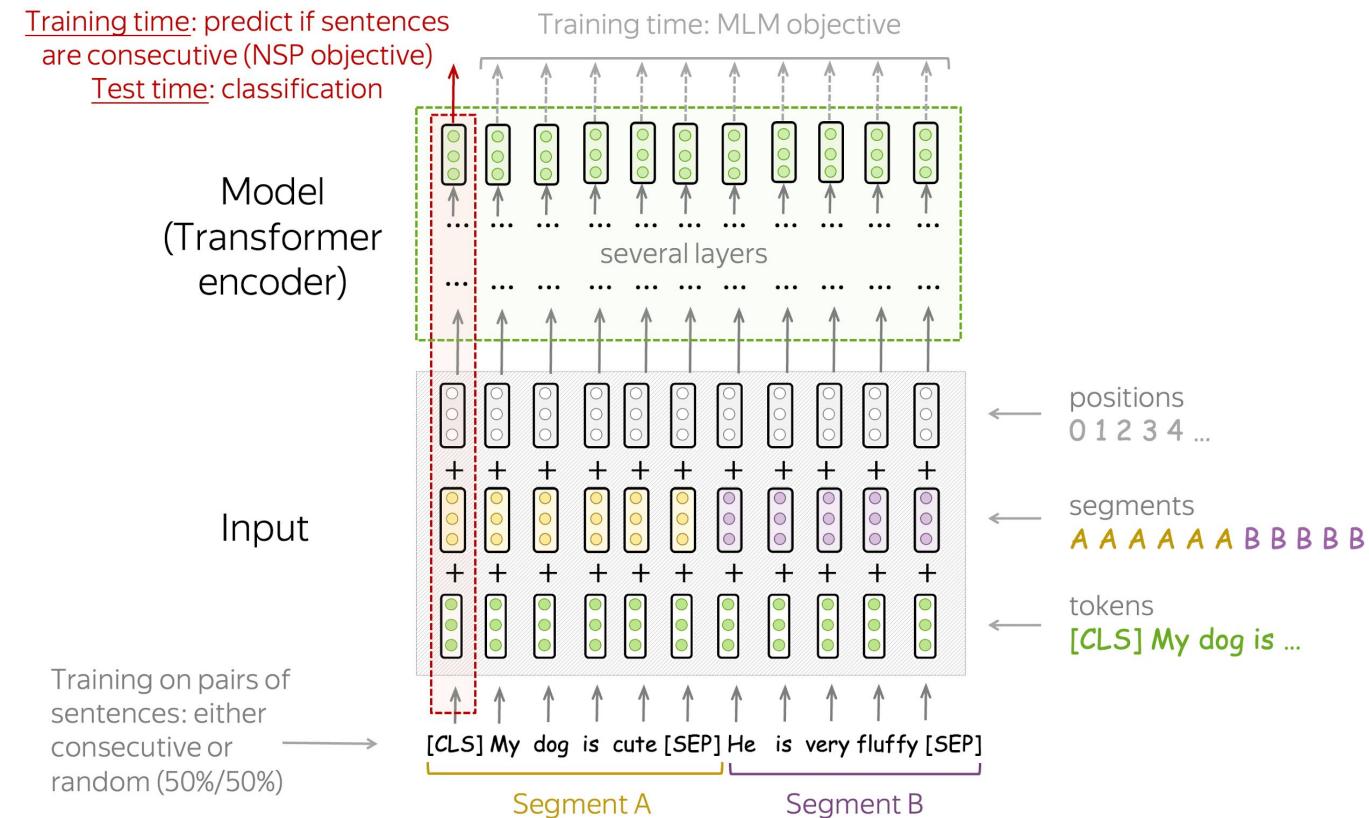
- Training time: predict if sentences are consecutive or not (Next Sentence Prediction /NSP objective)
 - Test time: downstream tasks (e.g., classification)

[SEP]: Special token-separator

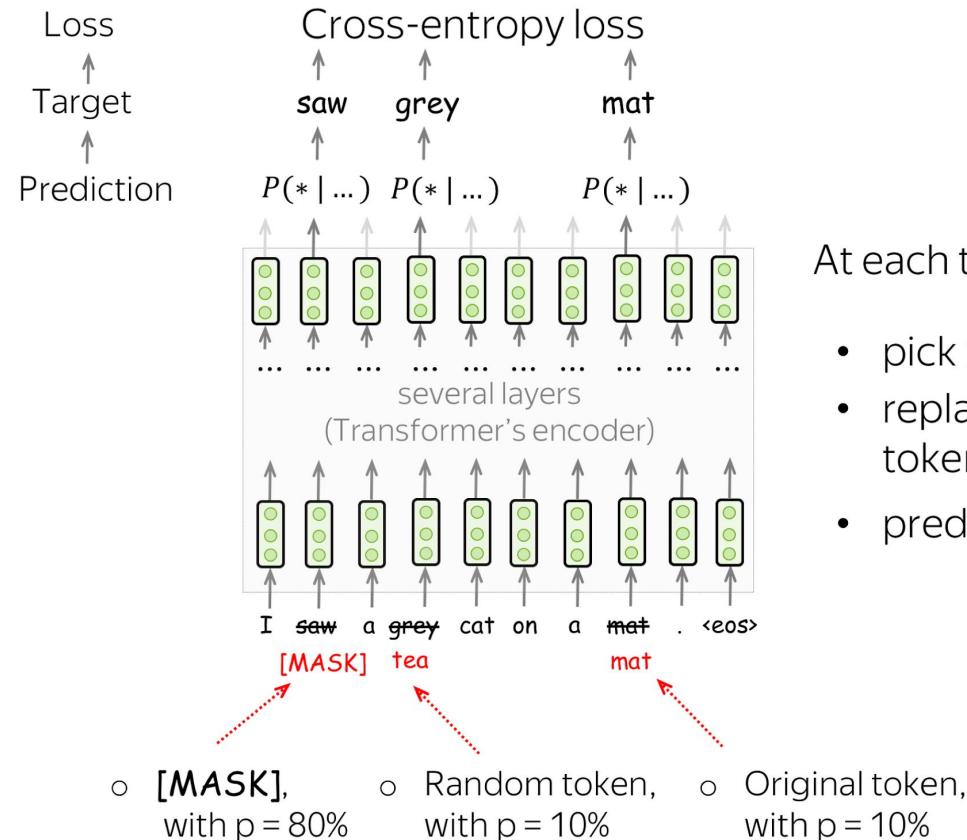


Training on pairs of sentences: either consecutive or random (50%/50%)

Using encoder as an embedding generator



Objective two: Masked Language Modeling



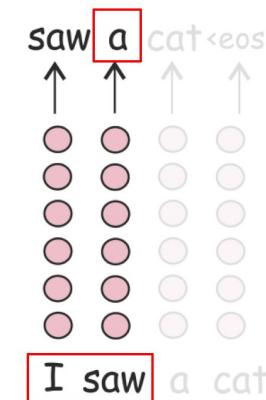
At each training step:

- pick randomly 15% of tokens
- replace each of the chosen tokens with something
- predict original chosen tokens

LM vs. MLM

Language Modeling

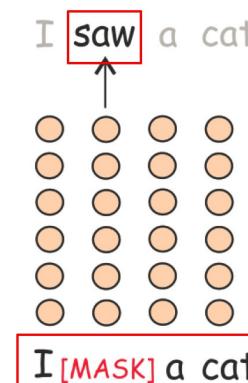
- Target: next token
- Prediction: $P(*) | I \text{ saw}$



left-to-right, does
not see future

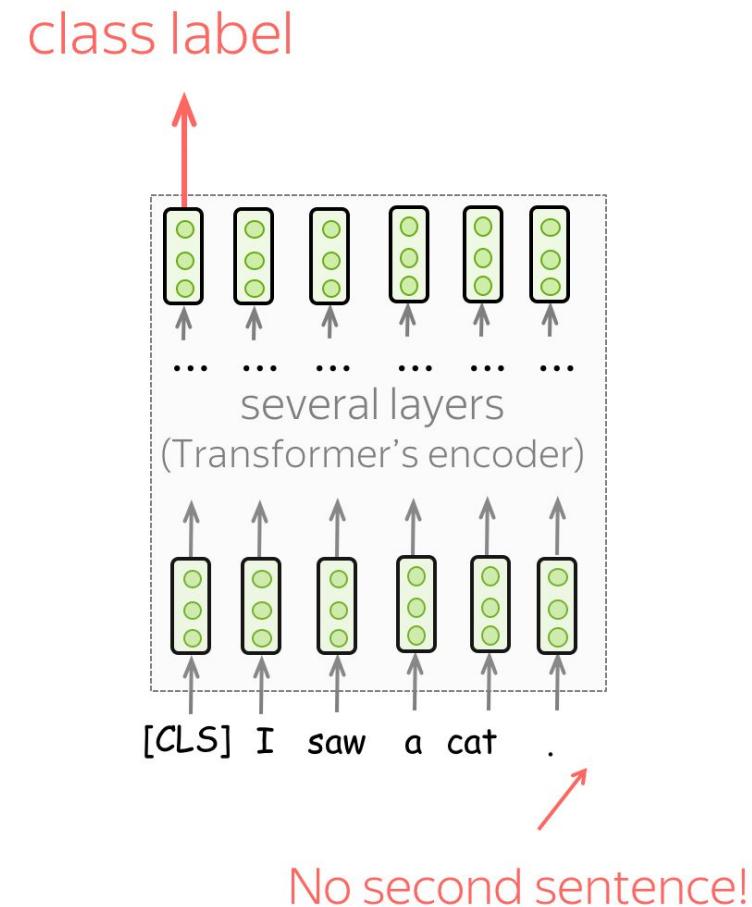
Masked Language Modeling

- Target: current token (the true one)
- Prediction: $P(*) | I [\text{MASK}] a \text{ cat}$

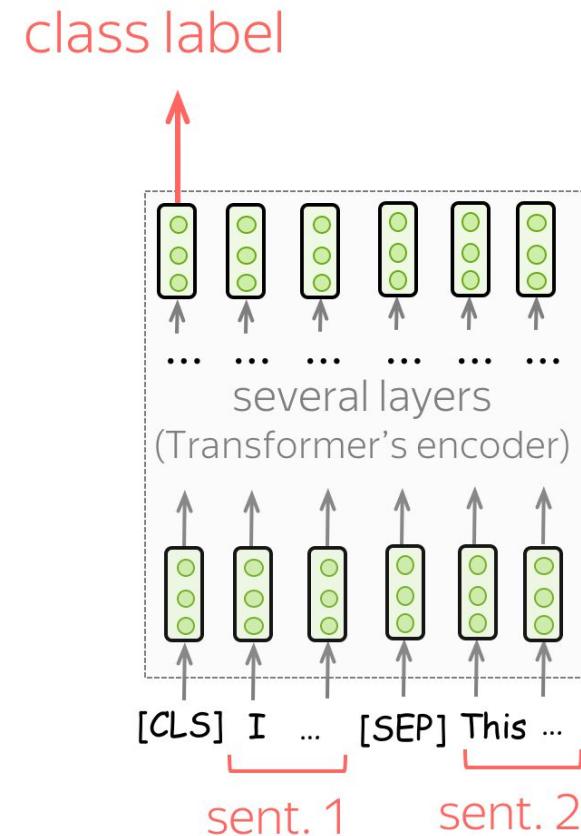


sees the whole text, but
something is corrupted

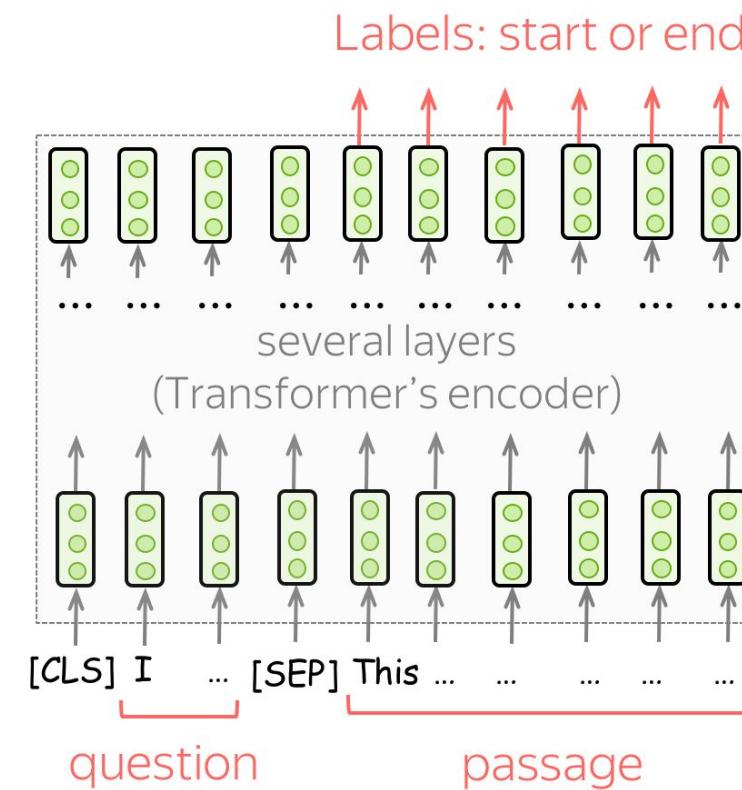
Single Sentence Classification



Sentence Pair Classification



Question Answering



Input tagging

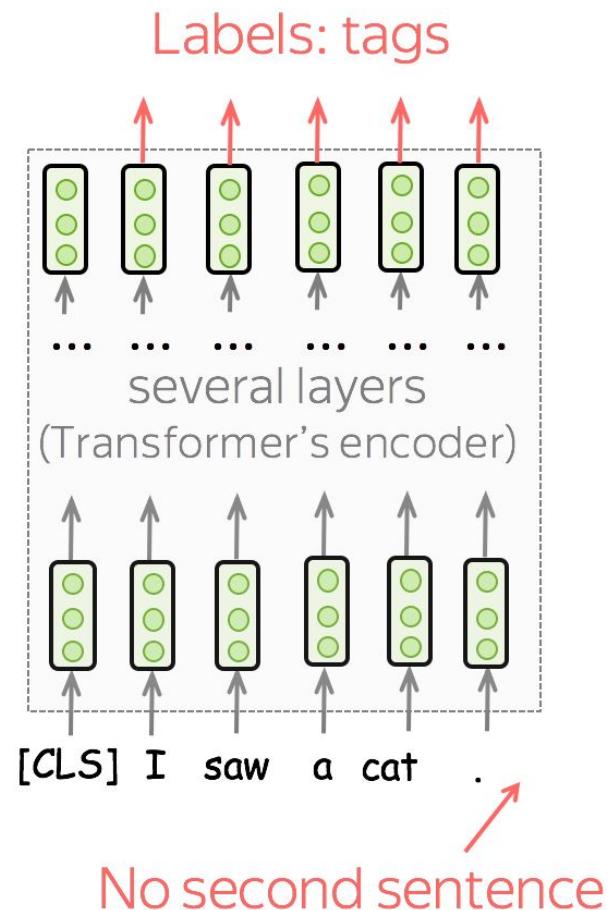
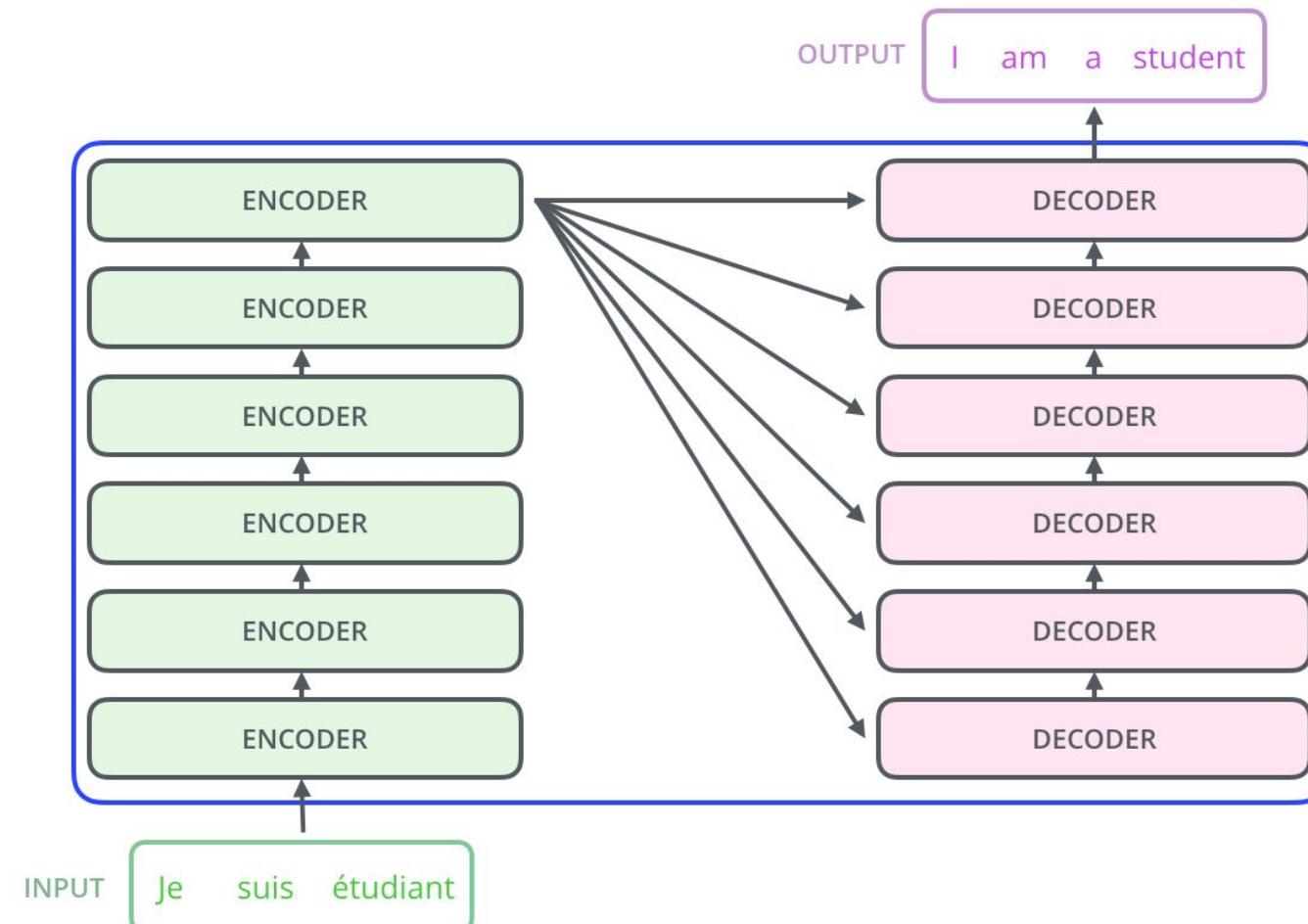


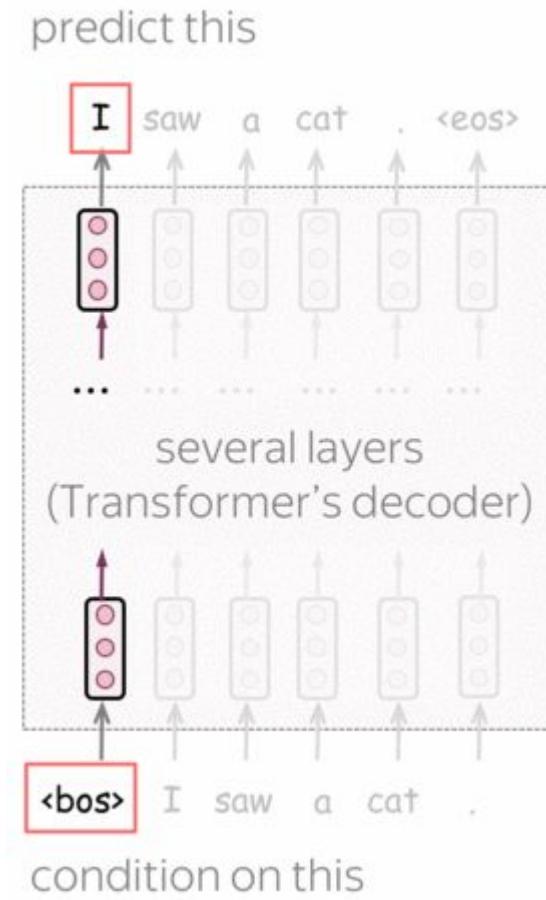
Table of Content

- The power of transfer learning
- From word-specific to contextual embeddings
- Transformer architecture overview
- BERT
- GPT

GPT is just a Decoder part



Decoder as a universal model



Decoder as a universal model

