

Week 8 – attention and BERT

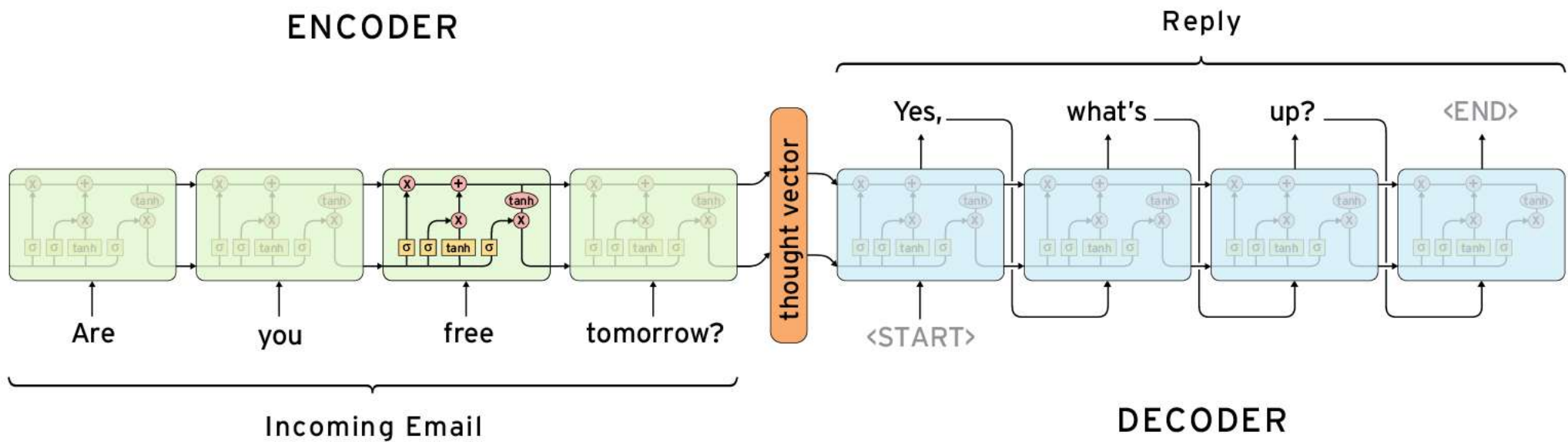
EGCO467 Natural Language and Speech Processing

Sources

- <https://arxiv.org/abs/1706.03762> Vaswani, Ashish, et al. "**Attention is all you need.**" Advances in neural information processing systems. 2017.
- <http://jalammar.github.io/illustrated-transformer/>
- <https://www.youtube.com/watch?v=iDulhoQ2pro>
- <http://nlp.seas.harvard.edu/2018/04/03/attention.html>

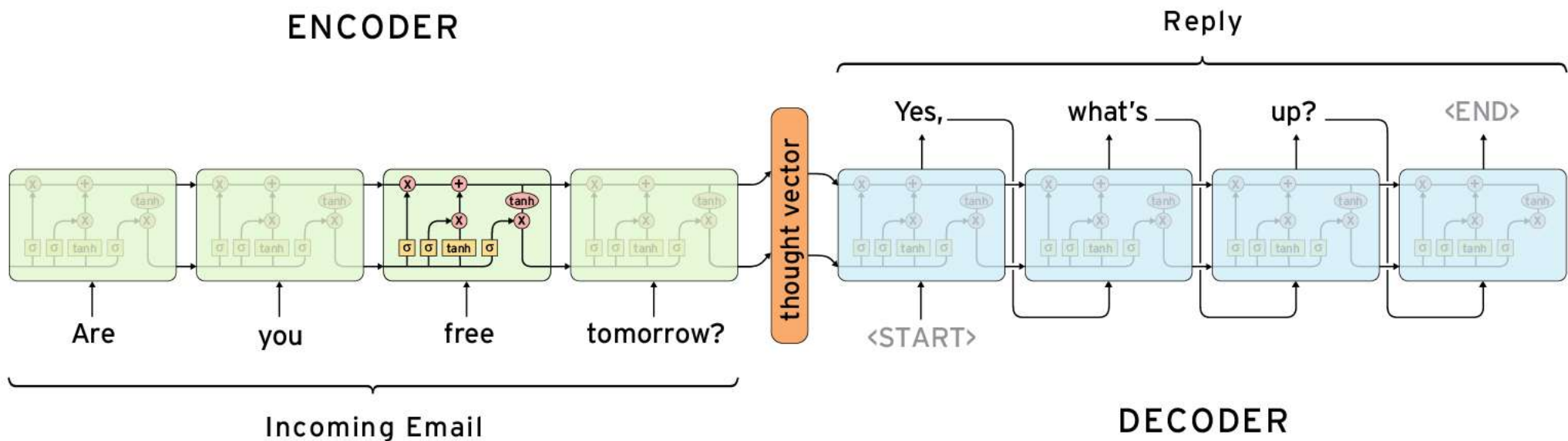
seq2seq

- sequence to sequence model



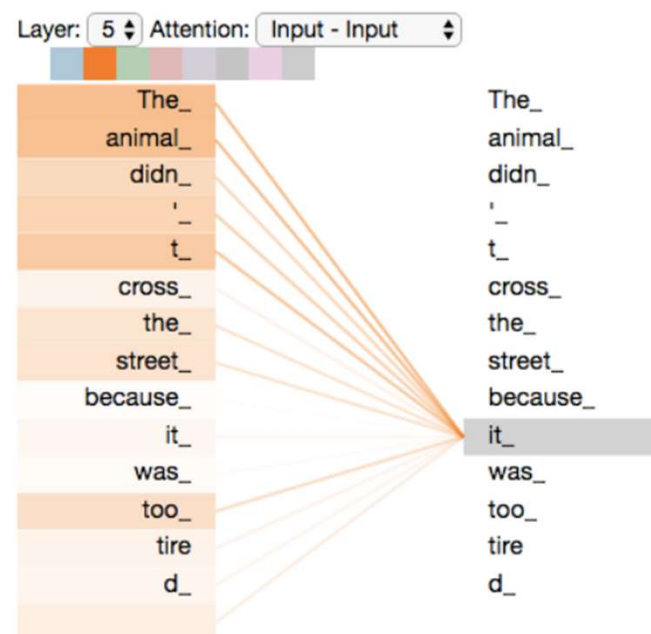
disadvantage of LSTM

- sequential mode – computation can't be parallelized
- all texts that comes before position t already absorbed into state vector h



Attention Mechanism

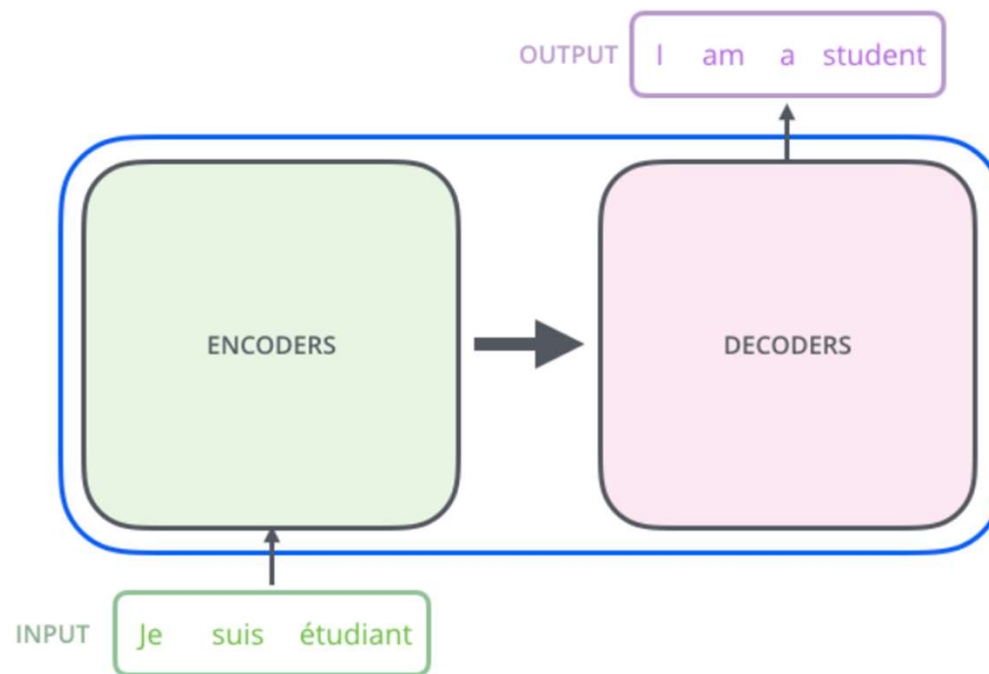
- Consider every pair of token
- If n tokens, there are n^2 pairs



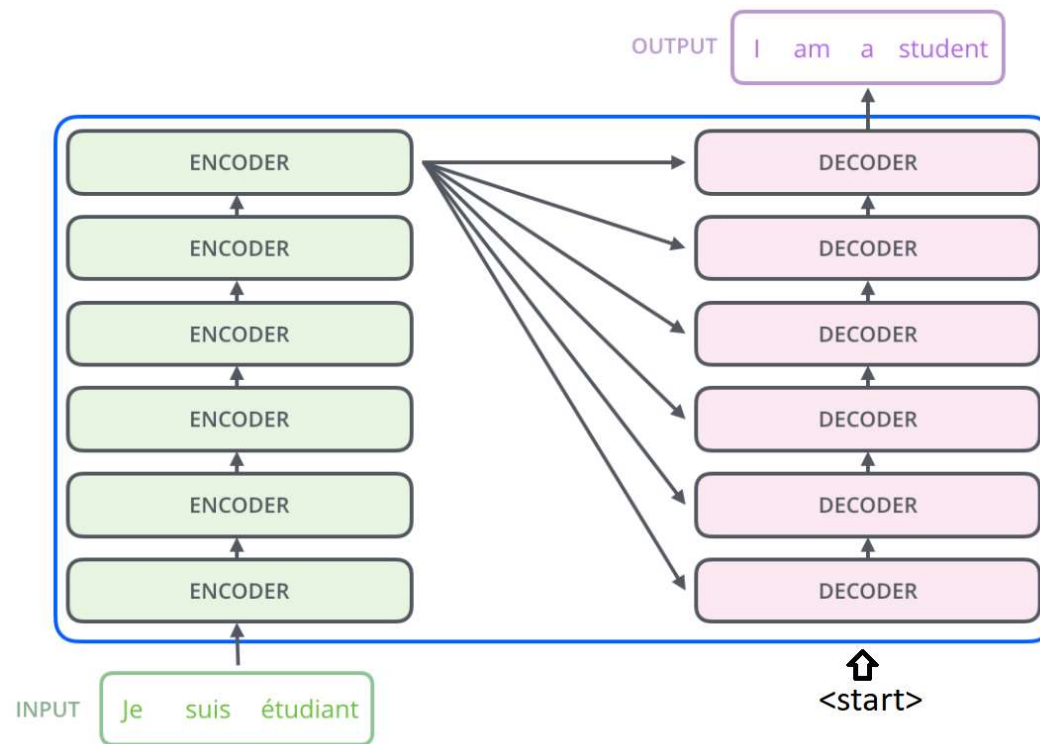
encoder-decoder (seq2seq)



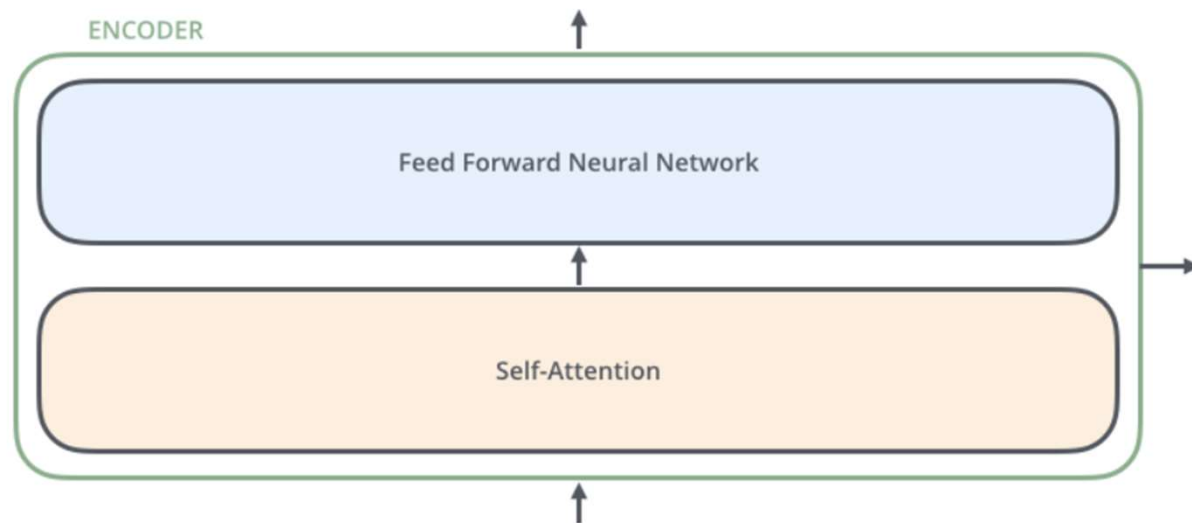
encoder-decoder



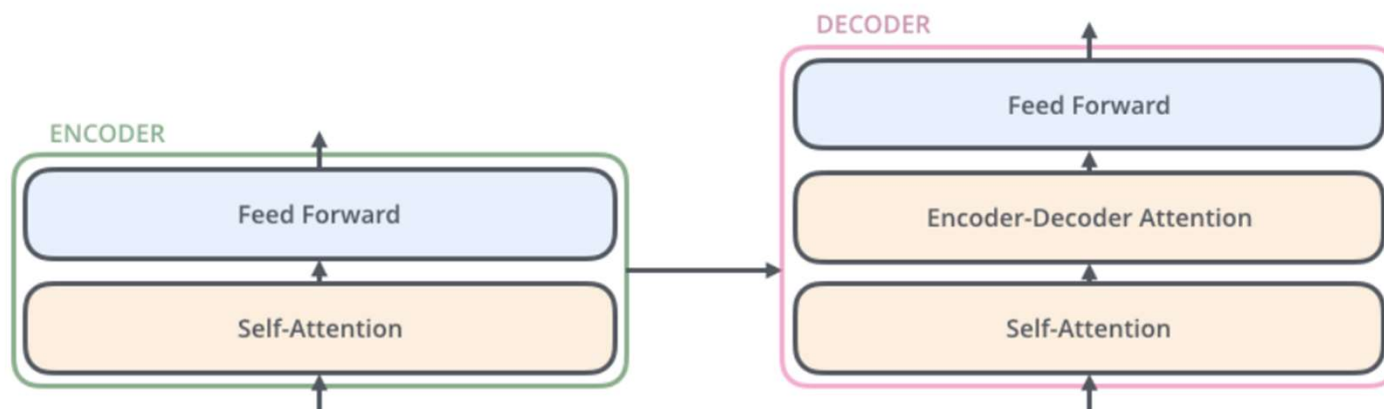
encoder-decoder



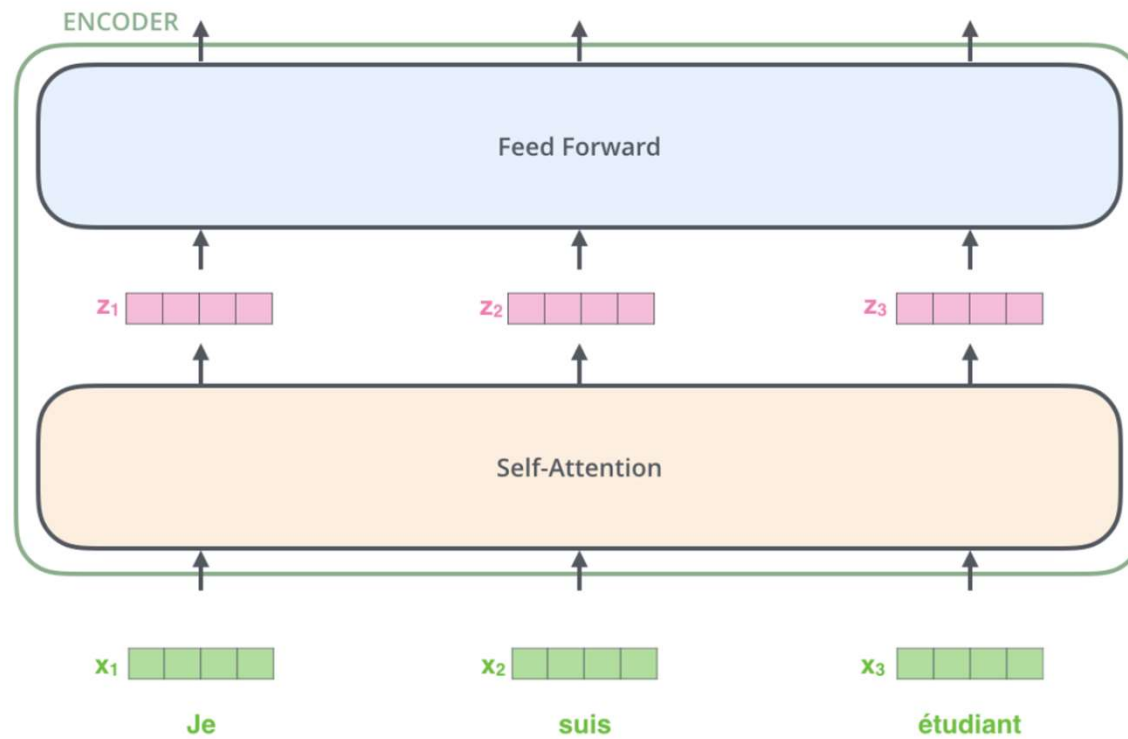
encoder



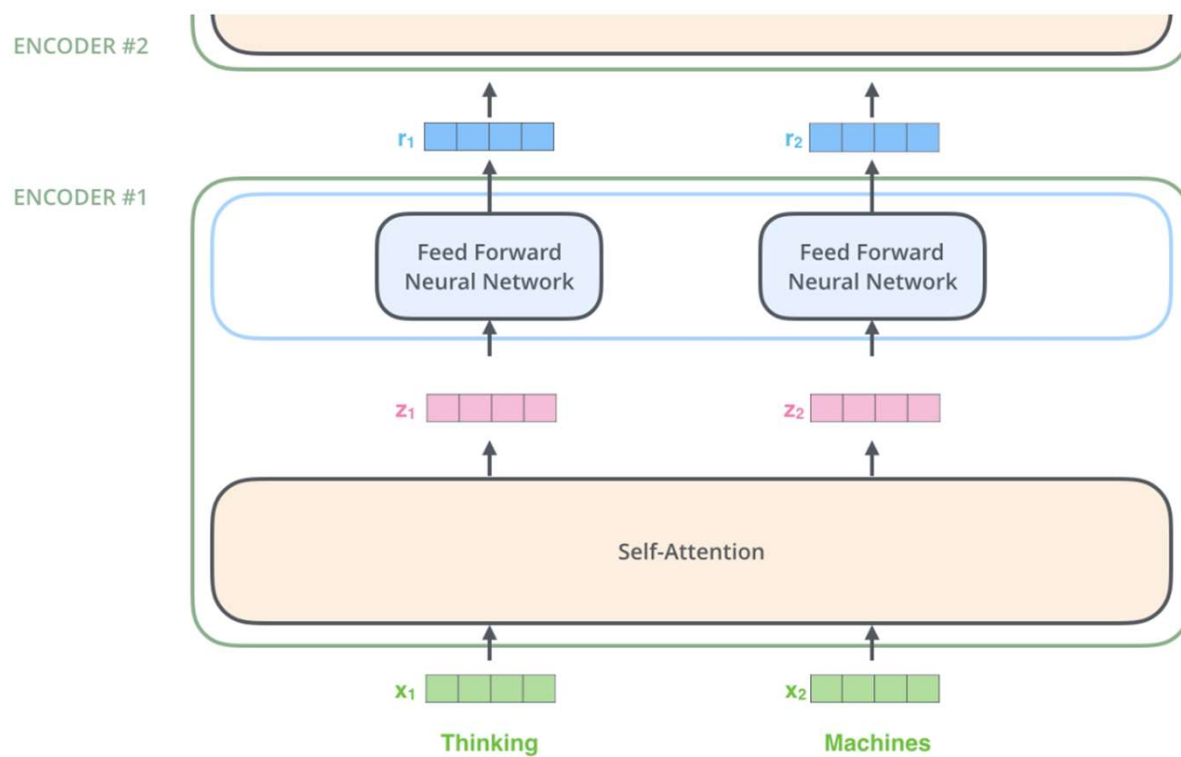
self-attention vs. "normal" attention



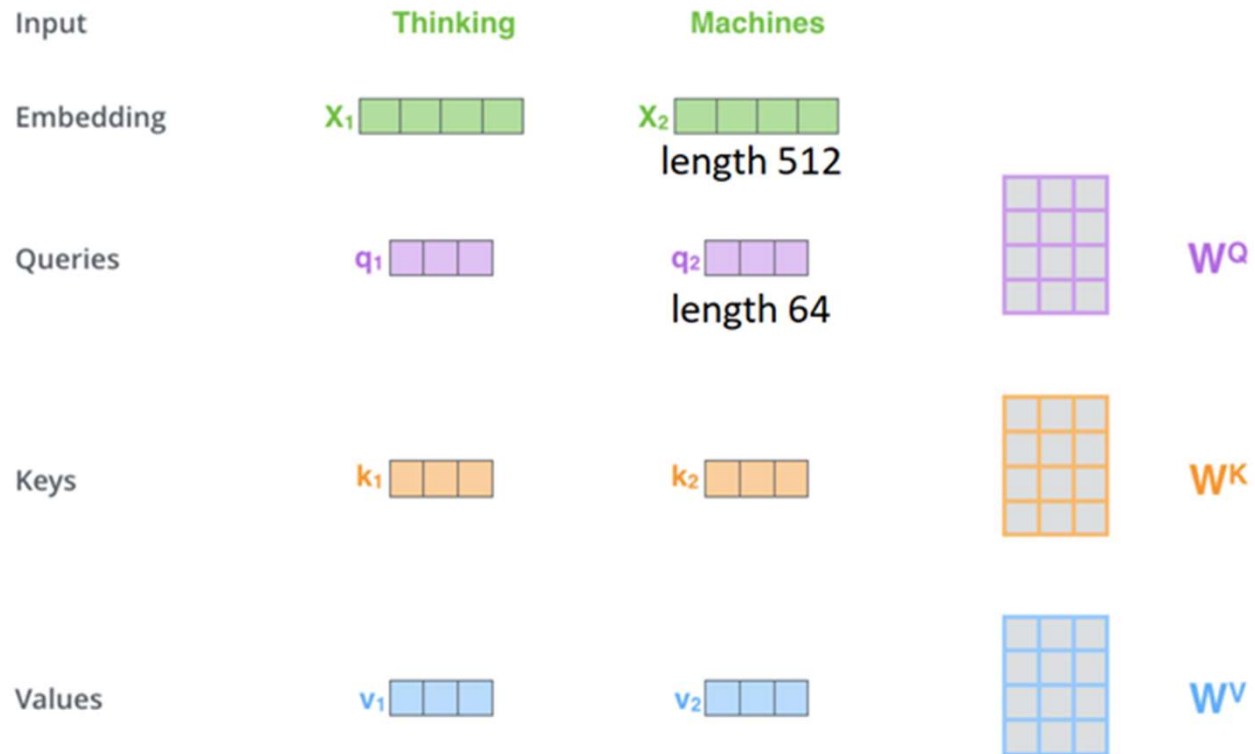
First Layer



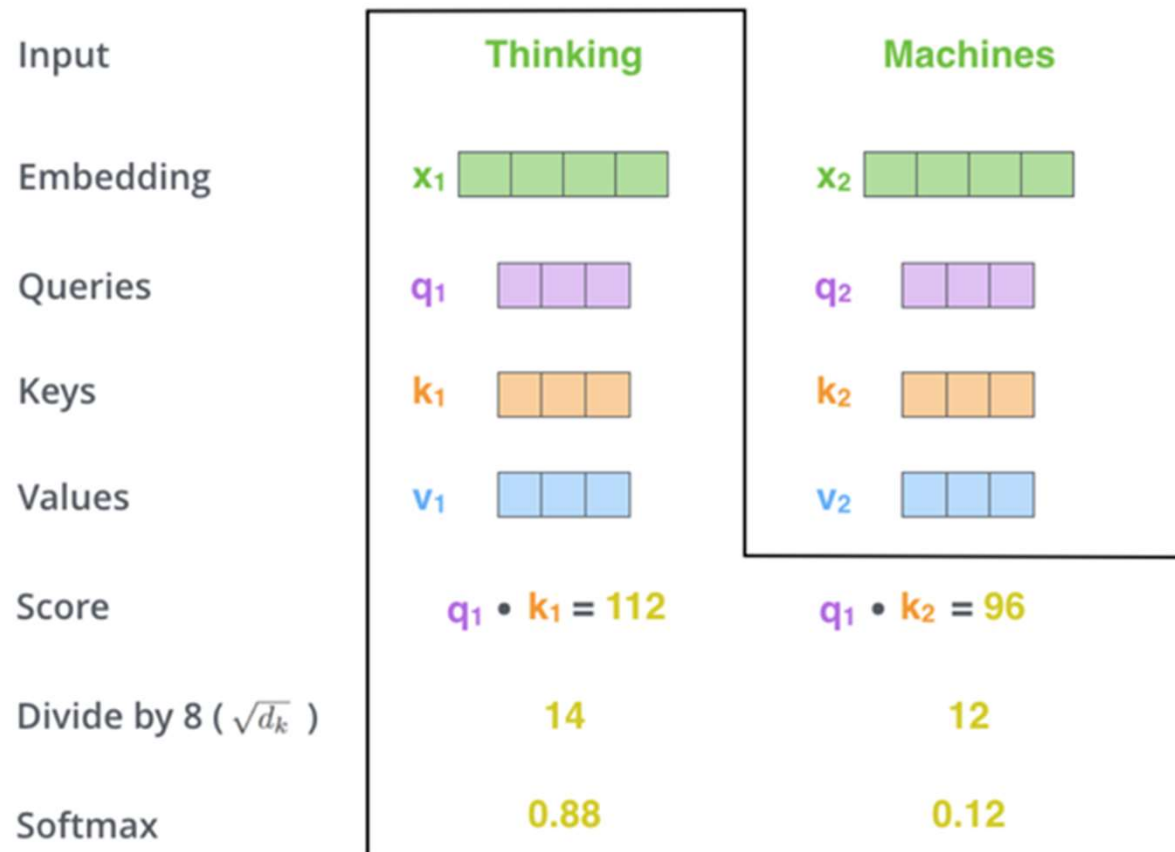
Inside First Layer



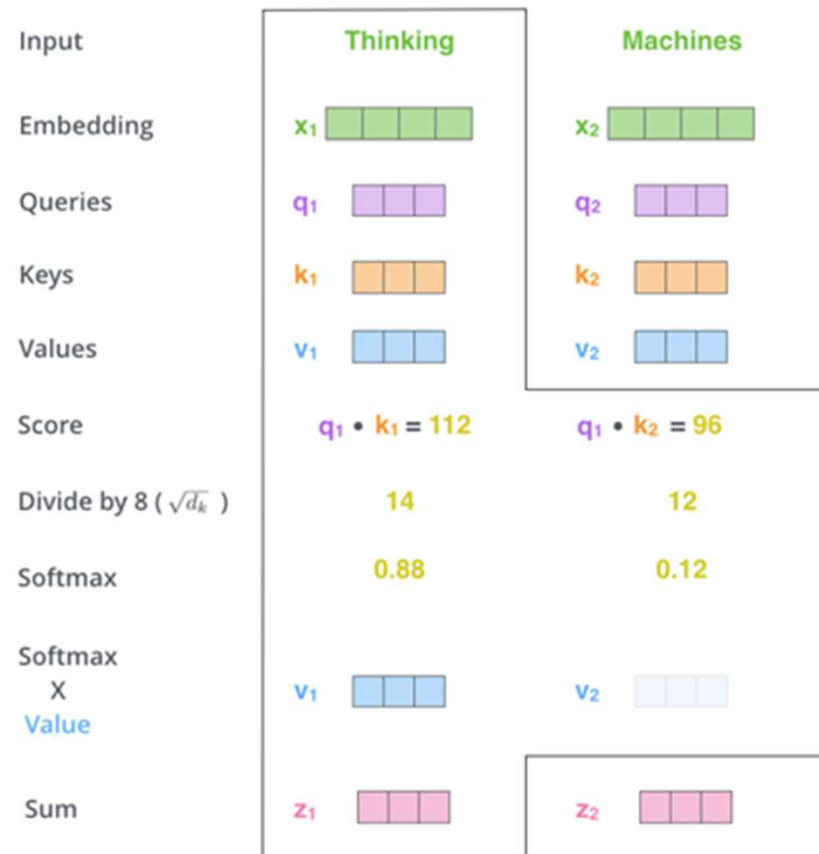
Self attention



Self attention



Self attention



Matrix form

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


A diagram illustrating matrix multiplication. On the left, a green 2x4 matrix labeled \mathbf{X} is multiplied by a purple 4x4 matrix labeled \mathbf{W}^Q . The result is a purple 2x4 matrix labeled \mathbf{Q} . The matrices are represented as grids of colored squares: green for \mathbf{X} , purple for \mathbf{W}^Q and \mathbf{Q} . The multiplication is indicated by a 'x' symbol, and the result is indicated by an '=' symbol.

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


A diagram illustrating matrix multiplication. On the left, a green 2x4 matrix labeled \mathbf{X} is multiplied by an orange 4x4 matrix labeled \mathbf{W}^K . The result is an orange 2x4 matrix labeled \mathbf{K} . The matrices are represented as grids of colored squares: green for \mathbf{X} , orange for \mathbf{W}^K and \mathbf{K} . The multiplication is indicated by a 'x' symbol, and the result is indicated by an '=' symbol.

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


A diagram illustrating matrix multiplication. On the left, a green 2x4 matrix labeled \mathbf{X} is multiplied by a blue 4x4 matrix labeled \mathbf{W}^V . The result is a blue 2x4 matrix labeled \mathbf{V} . The matrices are represented as grids of colored squares: green for \mathbf{X} , blue for \mathbf{W}^V and \mathbf{V} . The multiplication is indicated by a 'x' symbol, and the result is indicated by an '=' symbol.

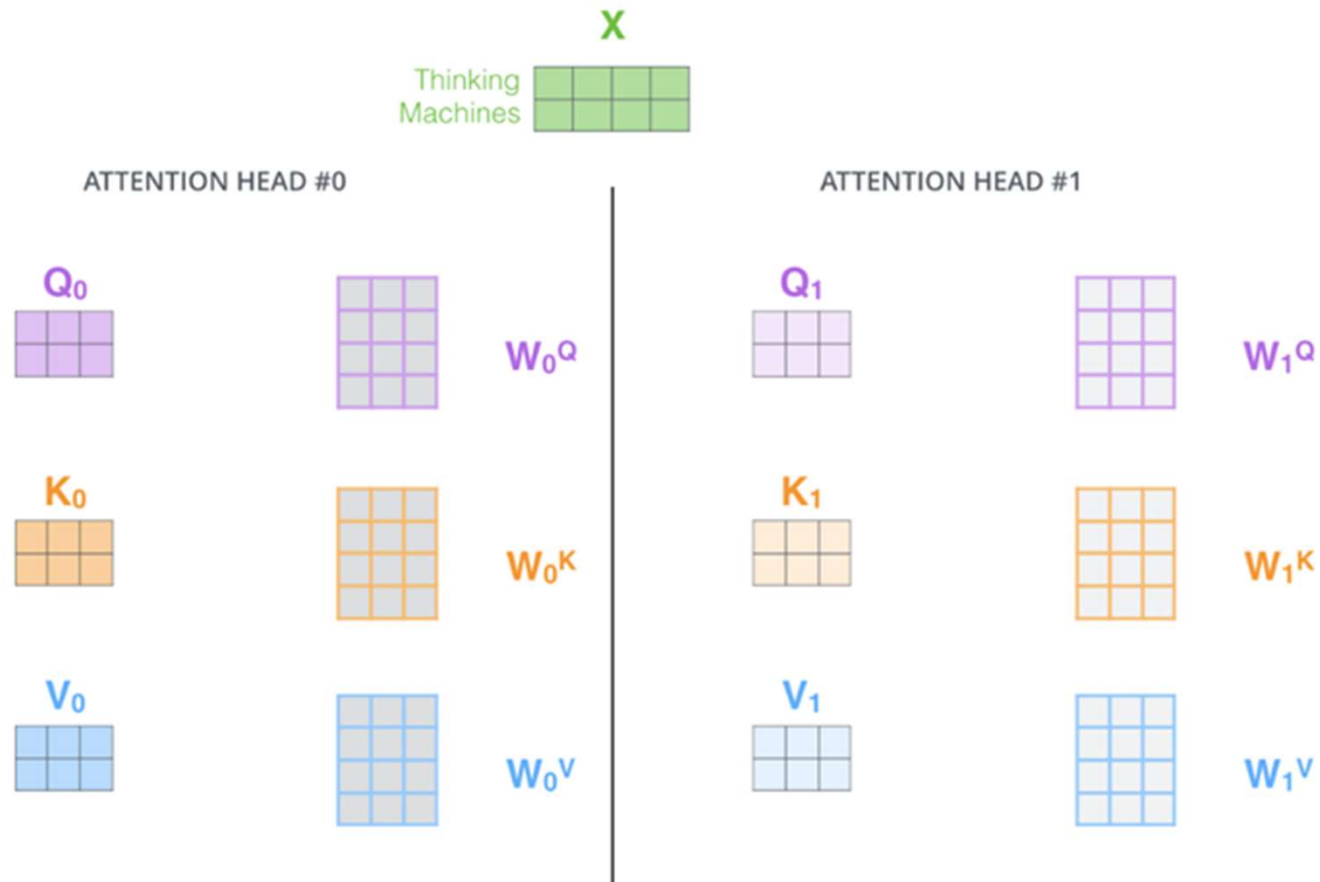
Matrix form

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

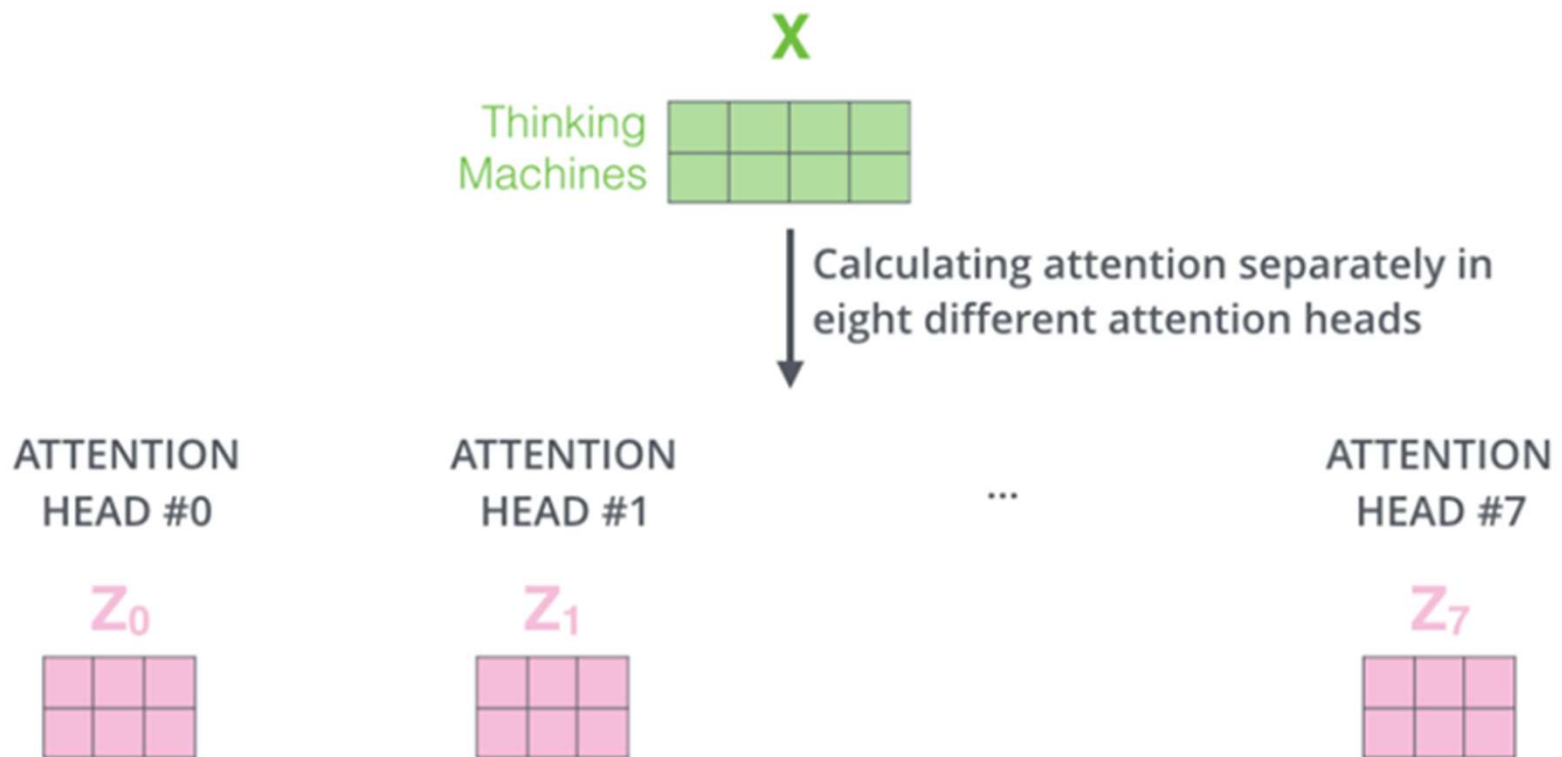
=

$$\begin{matrix} \text{Z} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix}$$

Multi-head



Multi-head



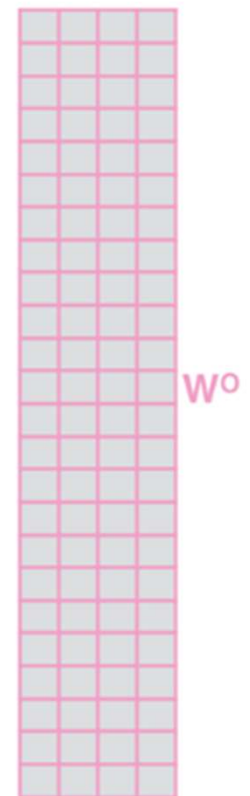
Multi-head

1) Concatenate all the attention heads



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

\times



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



Summary of self-attention

1) This is our input sentence*

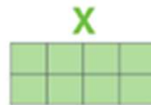
2) We embed each word*

3) Split into 8 heads. We multiply X or R with weight matrices

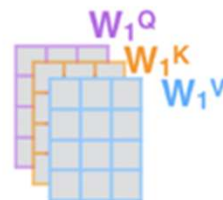
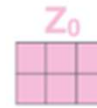
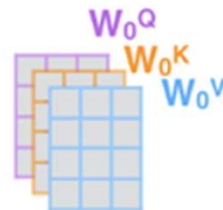
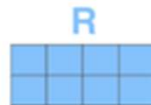
4) Calculate attention using the resulting $Q/K/V$ matrices

5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer

Thinking
Machines



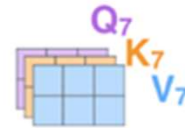
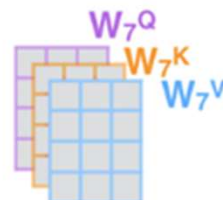
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



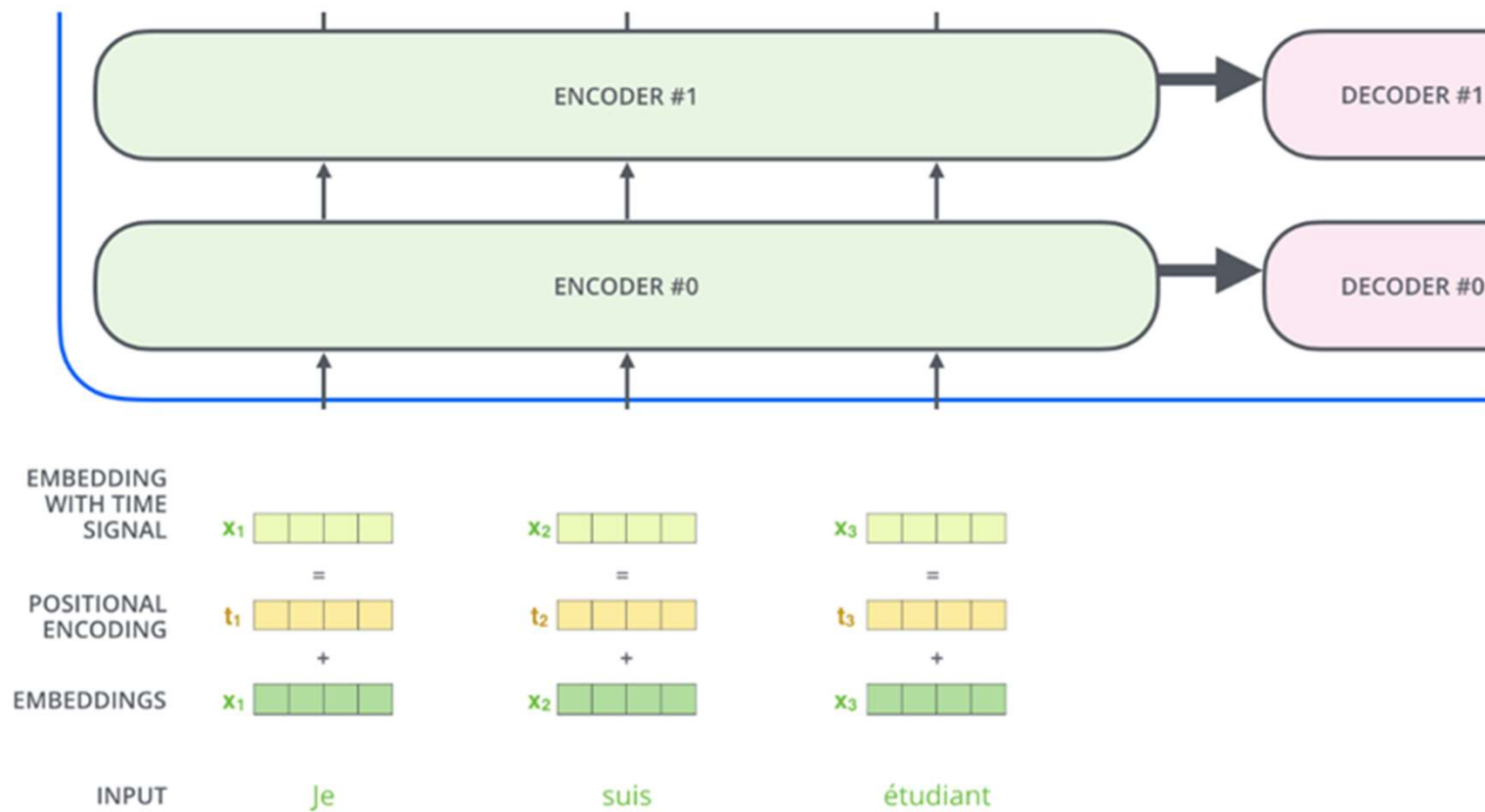
...

...

...



Order of sequence

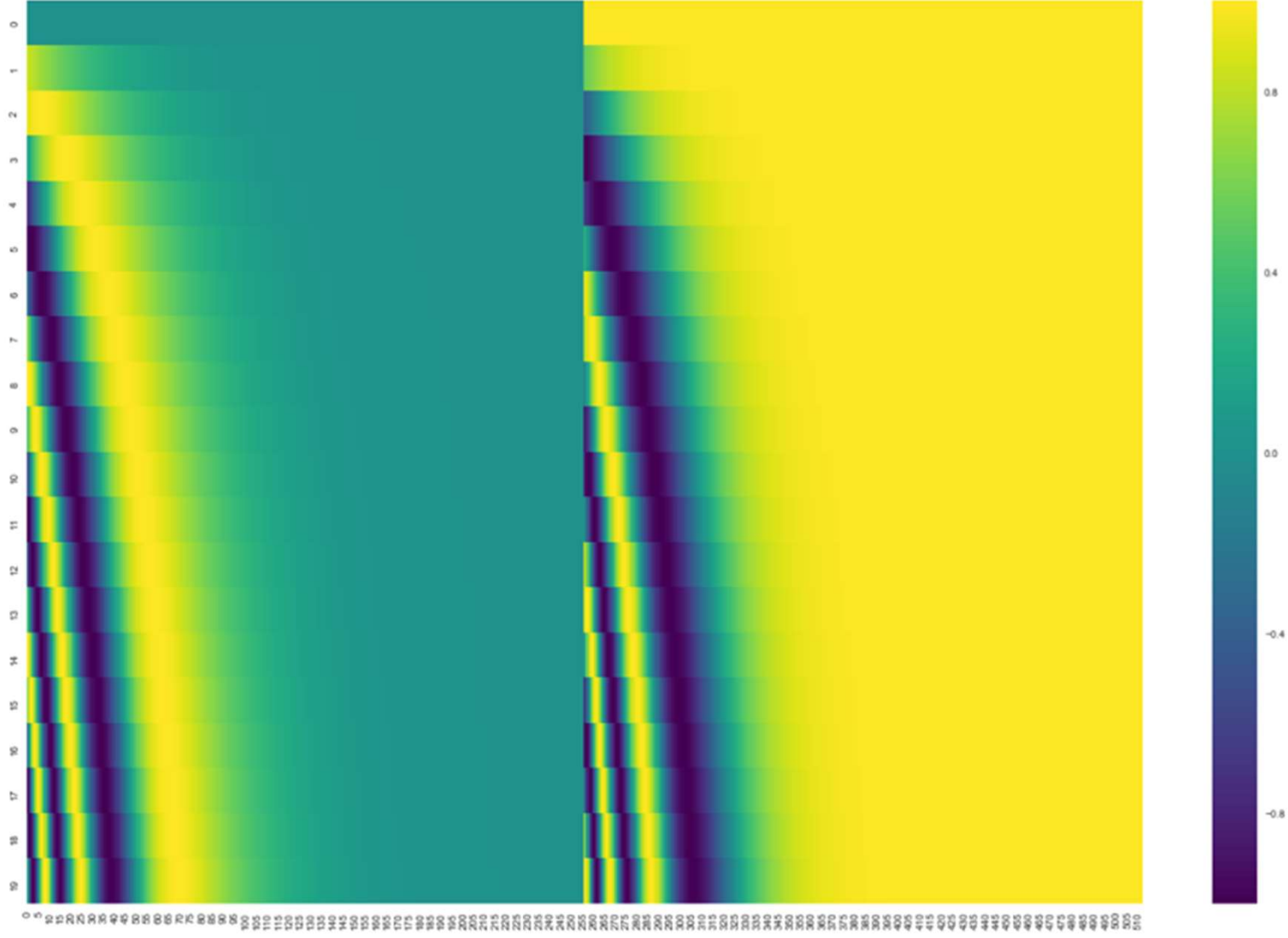


Positional encoding

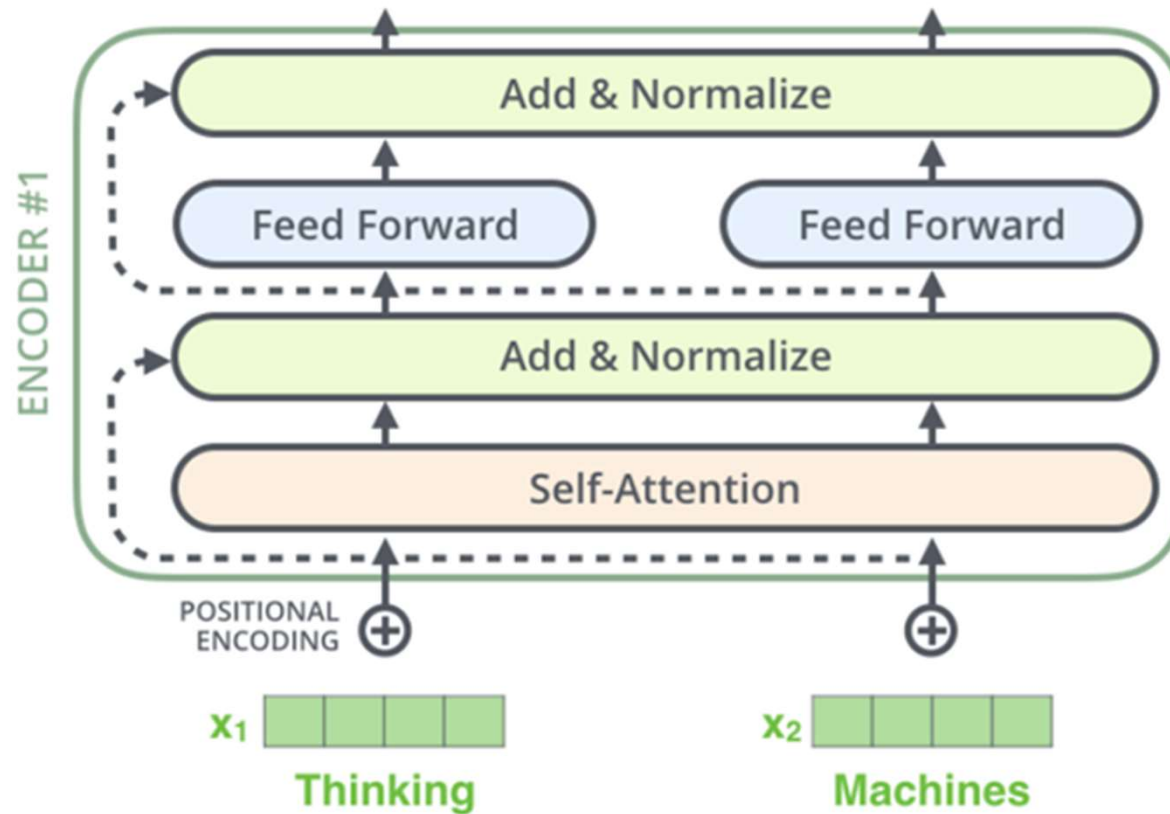
$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

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

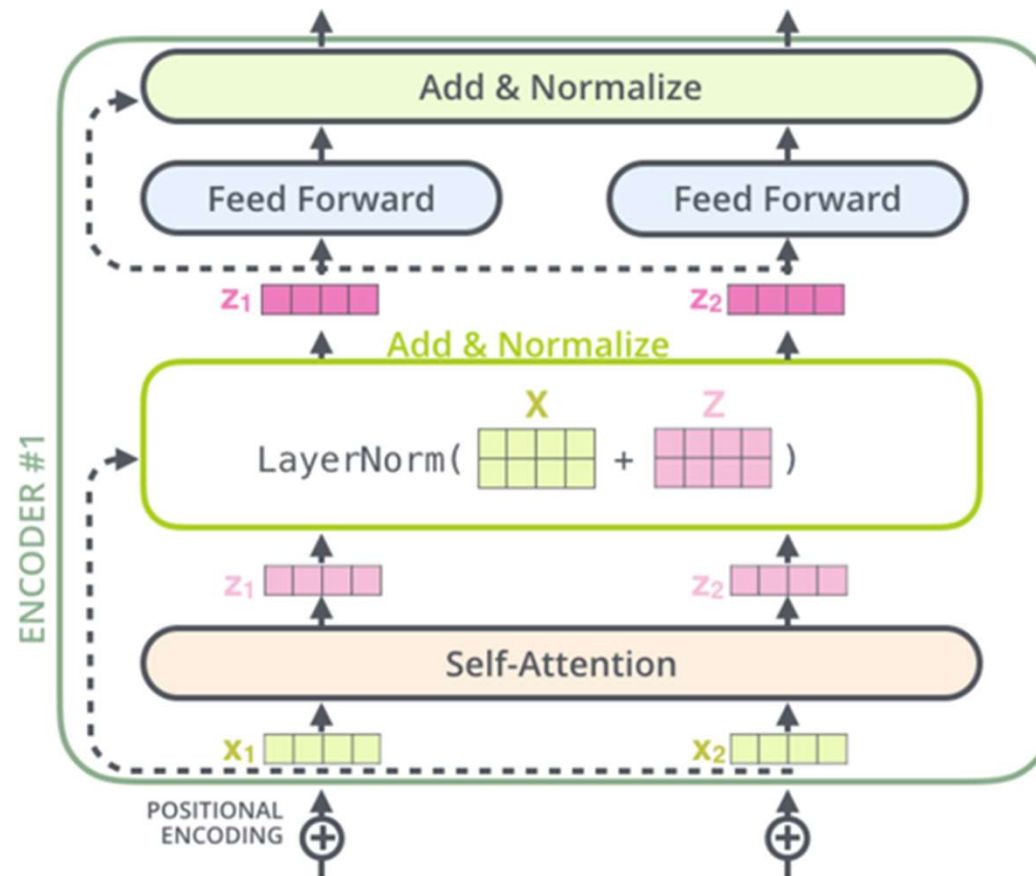




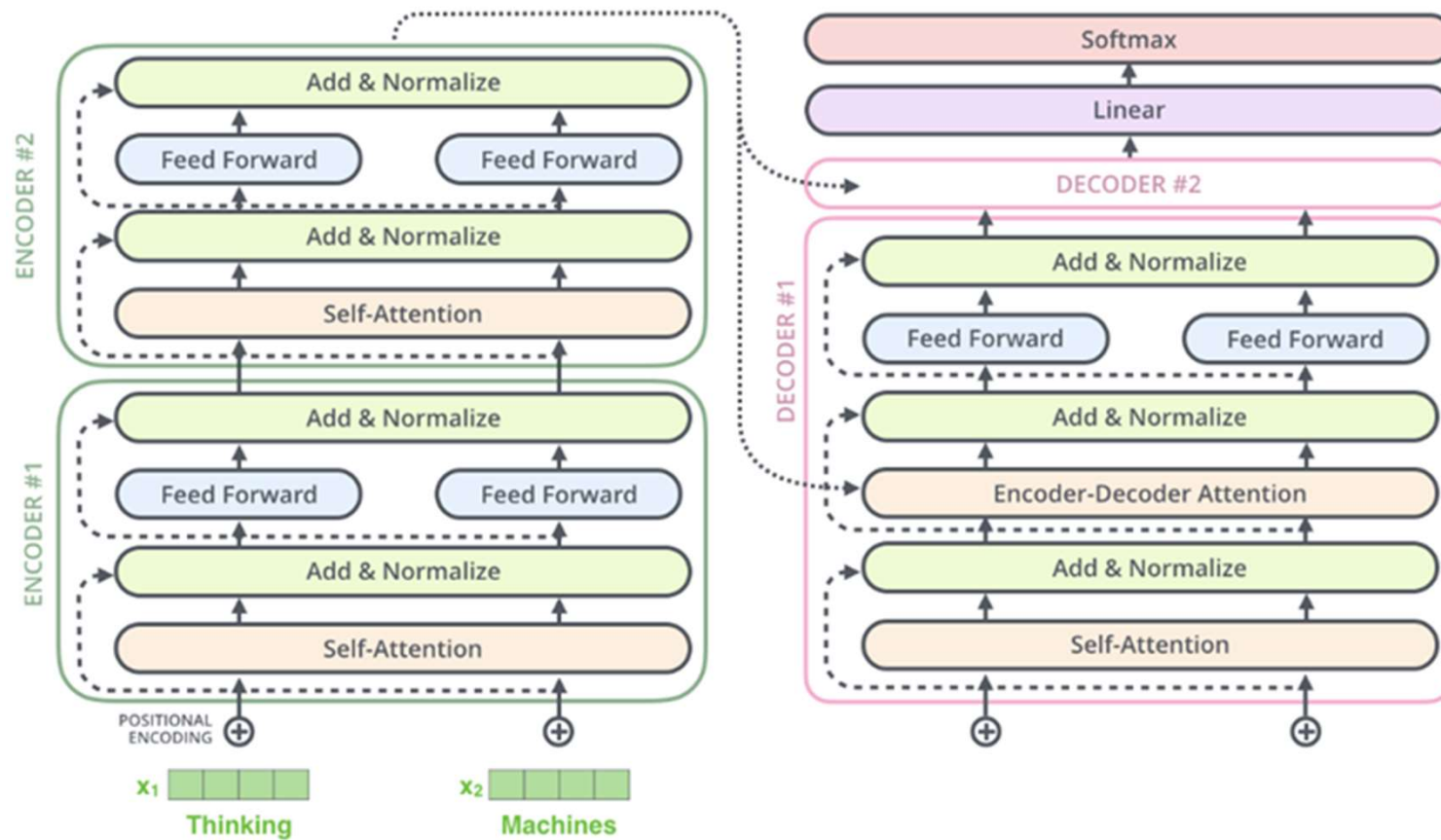
Residue connection

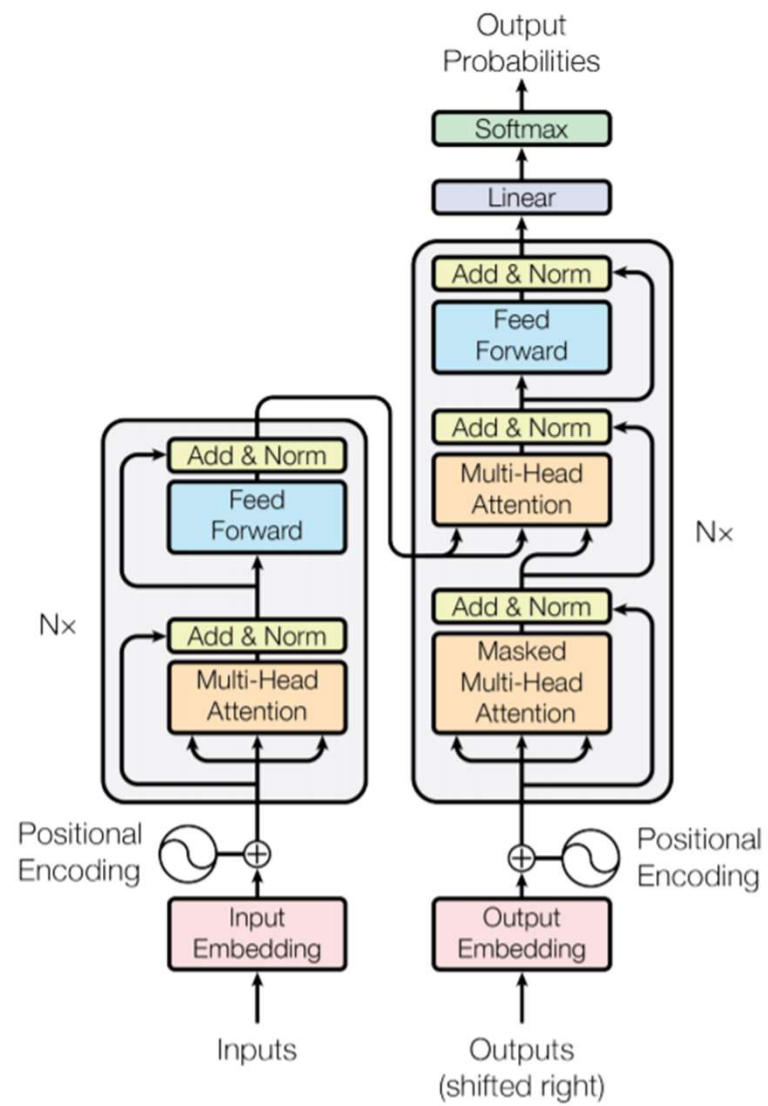


Complete encoder

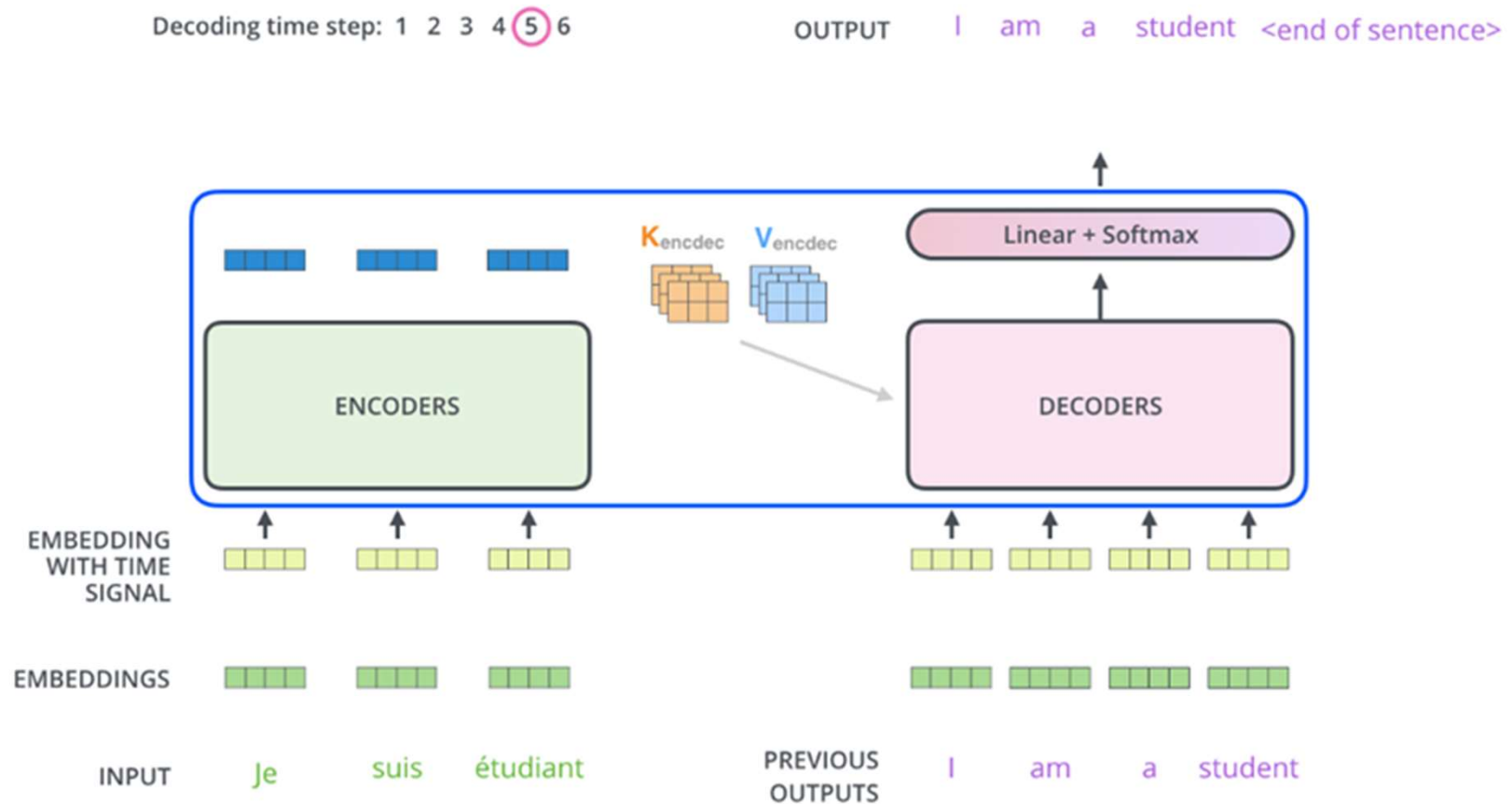


Decoder side





Decoding



Implementation

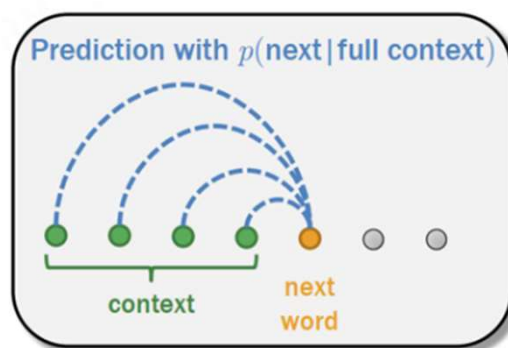
- <https://nlp.seas.harvard.edu/2018/04/03/attention.html>

BERT

<https://arxiv.org/pdf/1810.04805.pdf>

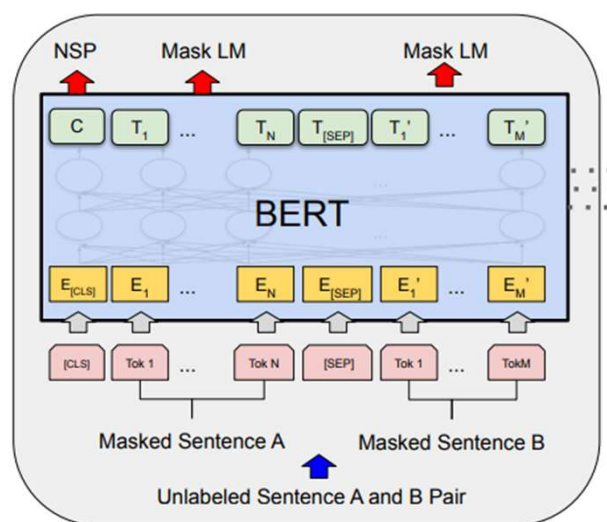
<https://arxiv.org/pdf/1907.11692.pdf>

LM vs. Masked LM



$$\max_{\theta} \log p_{\theta}(\mathbf{x}) = \sum_{t=1}^T \log p_{\theta}(x_t | \mathbf{x}_{<t}) = \sum_{t=1}^T \log \frac{\exp(h_{\theta}(\mathbf{x}_{1:t-1})^{\top} e(x_t))}{\sum_{x'} \exp(h_{\theta}(\mathbf{x}_{1:t-1})^{\top} e(x'))},$$

Masked LM



Special Tokens:

[CLS]

[SEP]

[MASK]

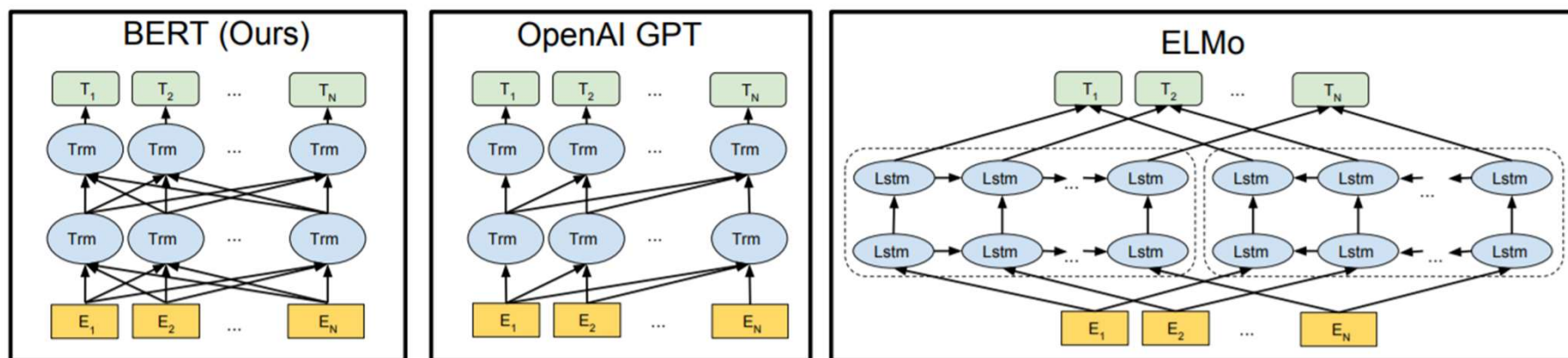
[PAD]

NSP = next sentence prediction

$$\max_{\theta} \log p_{\theta}(\bar{\mathbf{x}} \mid \hat{\mathbf{x}}) \approx \sum_{t=1}^T m_t \log p_{\theta}(x_t \mid \hat{\mathbf{x}}) = \sum_{t=1}^T m_t \log \frac{\exp(H_{\theta}(\hat{\mathbf{x}})_t^{\top} e(x_t))}{\sum_{x'} \exp(H_{\theta}(\hat{\mathbf{x}})_t^{\top} e(x'))},$$

masking: 15% of tokens is replaced with [MASK],
 m_t is 0 if [MASK] and 1 if not

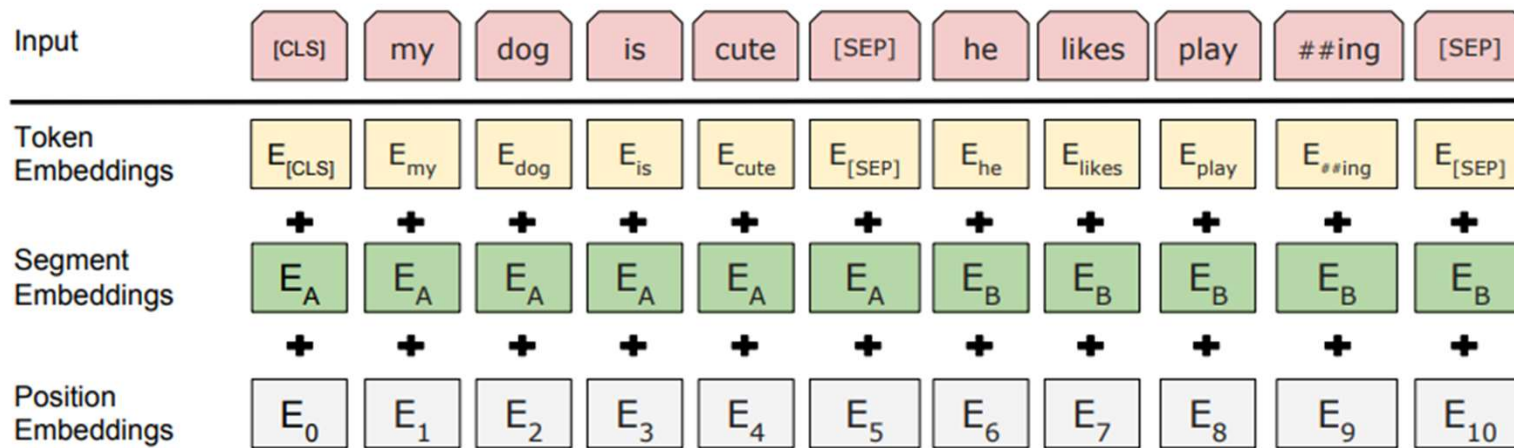
Pre-training model architectures



Pretraining Tasks

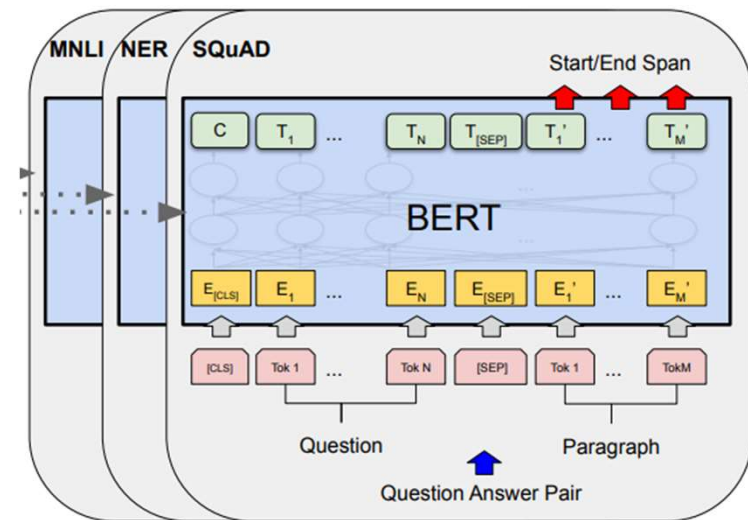
- LM task: predict the masked token.
- NSP: next sentence prediction – predict if sentence B is actually a sentence that follows A in corpus.
- For LM task: token are masked with 15% probability.
- NSP: 50-50 split between actual A followed by B, and replace B with a random sentence from corpus.

Encoding the input



Fine-tuning

- Change the head according to task
- E.g.
 - text recognition – softmax head
 - squad (question answering) two softmax heads, one for start token position and for end token position



Roberta

- A variation on BERT
- remove NSP loss
- dynamic masking (masks are randomized for each batch)
- large batch size
- byte-level BPE tokenizer (BERT uses Wordpiece)
- more data and train for longer

Roberta vs. Bert

Model	data	bsz	steps	SQuAD (v1.1/2.0)	MNLI-m	SST-2
RoBERTa						
with BOOKS + WIKI	16GB	8K	100K	93.6/87.3	89.0	95.3
+ additional data (§3.2)	160GB	8K	100K	94.0/87.7	89.3	95.6
+ pretrain longer	160GB	8K	300K	94.4/88.7	90.0	96.1
+ pretrain even longer	160GB	8K	500K	94.6/89.4	90.2	96.4
BERT _{LARGE}						
with BOOKS + WIKI	13GB	256	1M	90.9/81.8	86.6	93.7
XLNet _{LARGE}						
with BOOKS + WIKI	13GB	256	1M	94.0/87.8	88.4	94.4
+ additional data	126GB	2K	500K	94.5/88.8	89.8	95.6

Example - Huggingface

- <https://github.com/huggingface/transformers/tree/v4.23-release/examples/pytorch>