

Transformer 模型簡介

Attention is all you need

電信碩一 謝文崑 r11942078

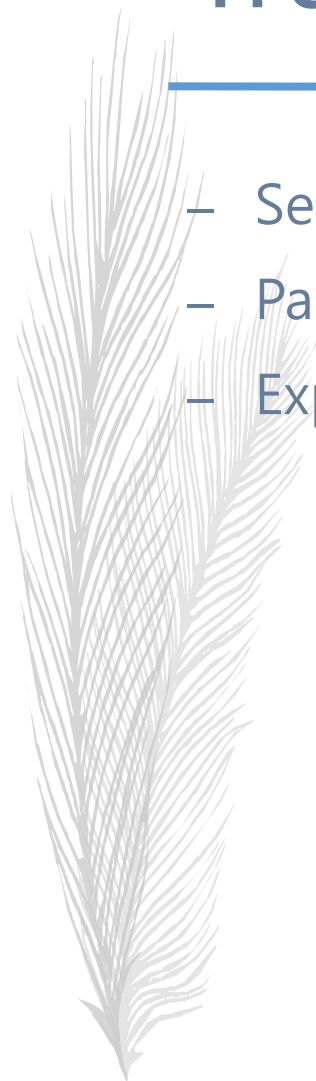
為什麼介紹Transformer？

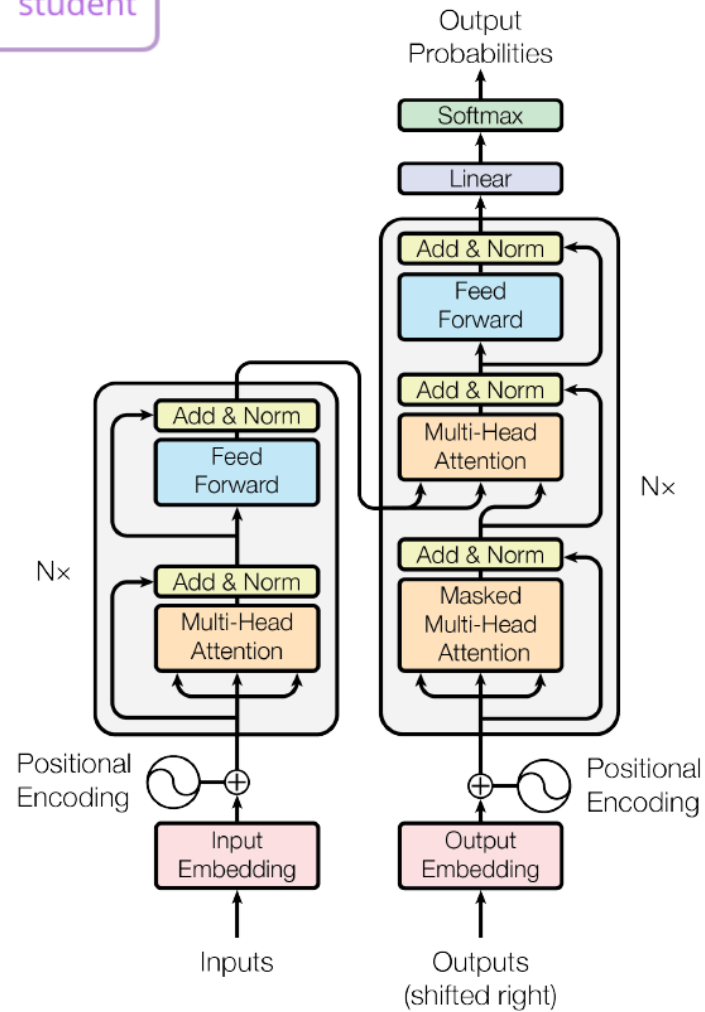
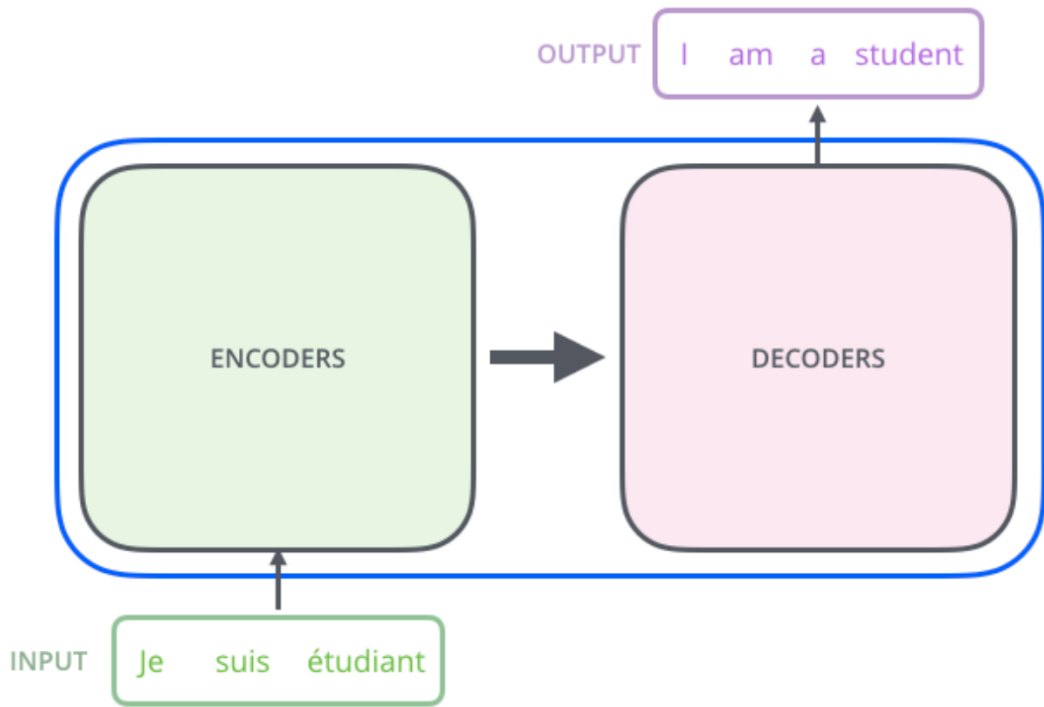
- 現在熱門的大型語言模型如BERT、GPT，都是基於Transformer架構
- 在最近幾年的論文裡，如果有用到機器學習技術解決問題，使用的模型大部分也與Transformer有關
- 我想花點時間認識一下它

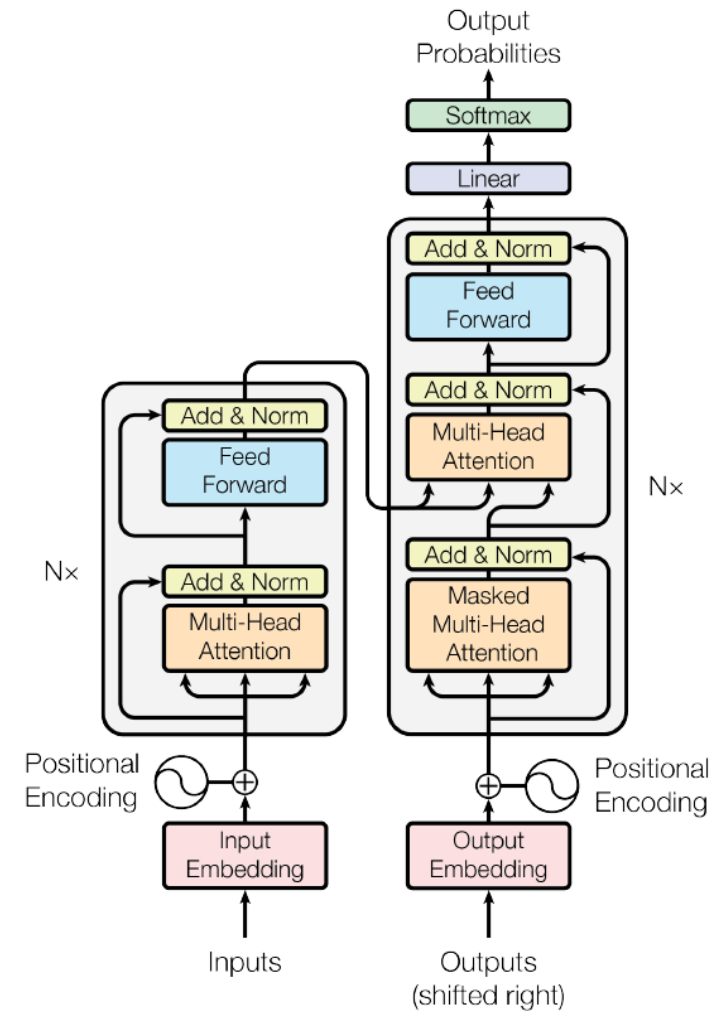
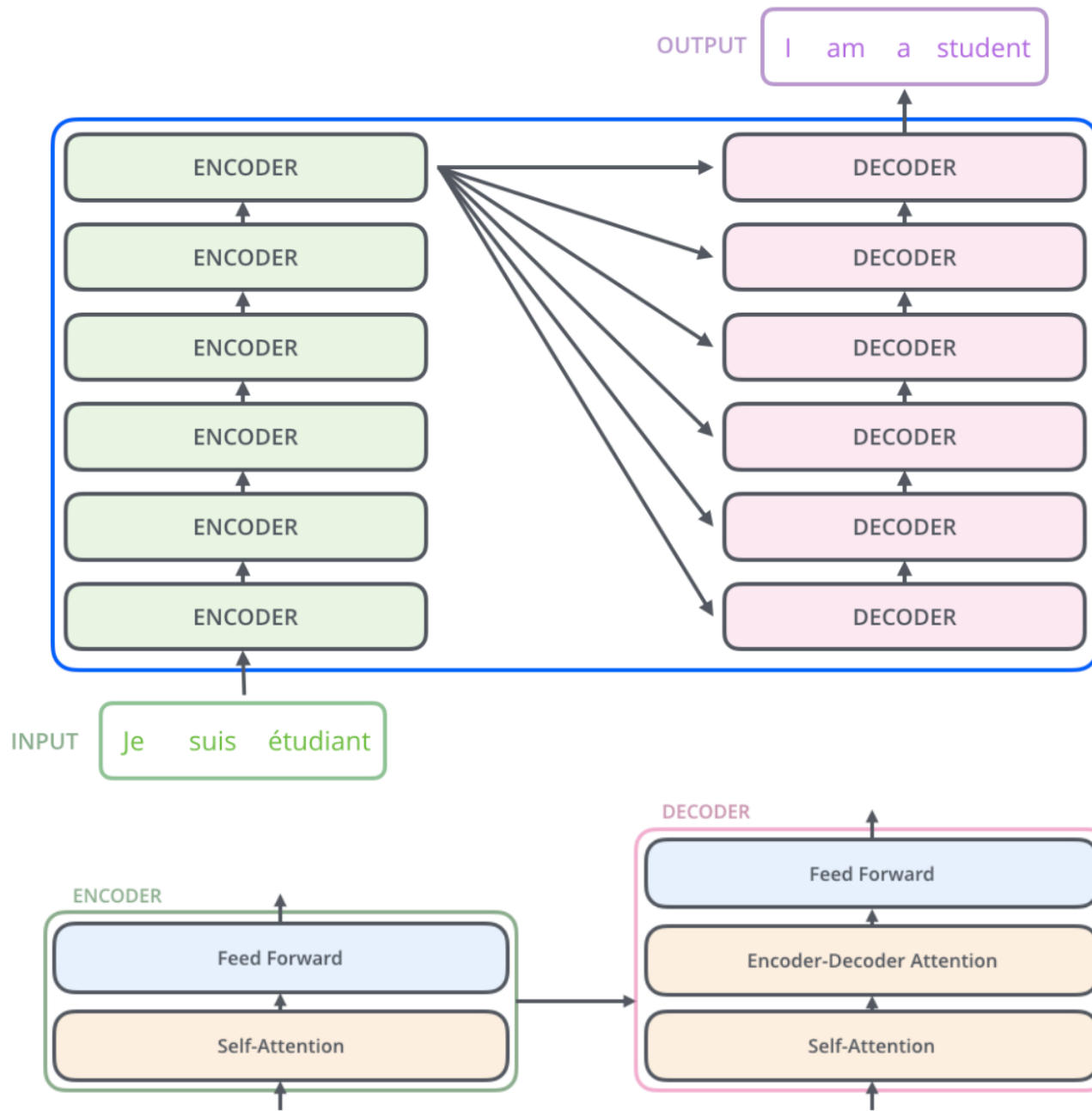


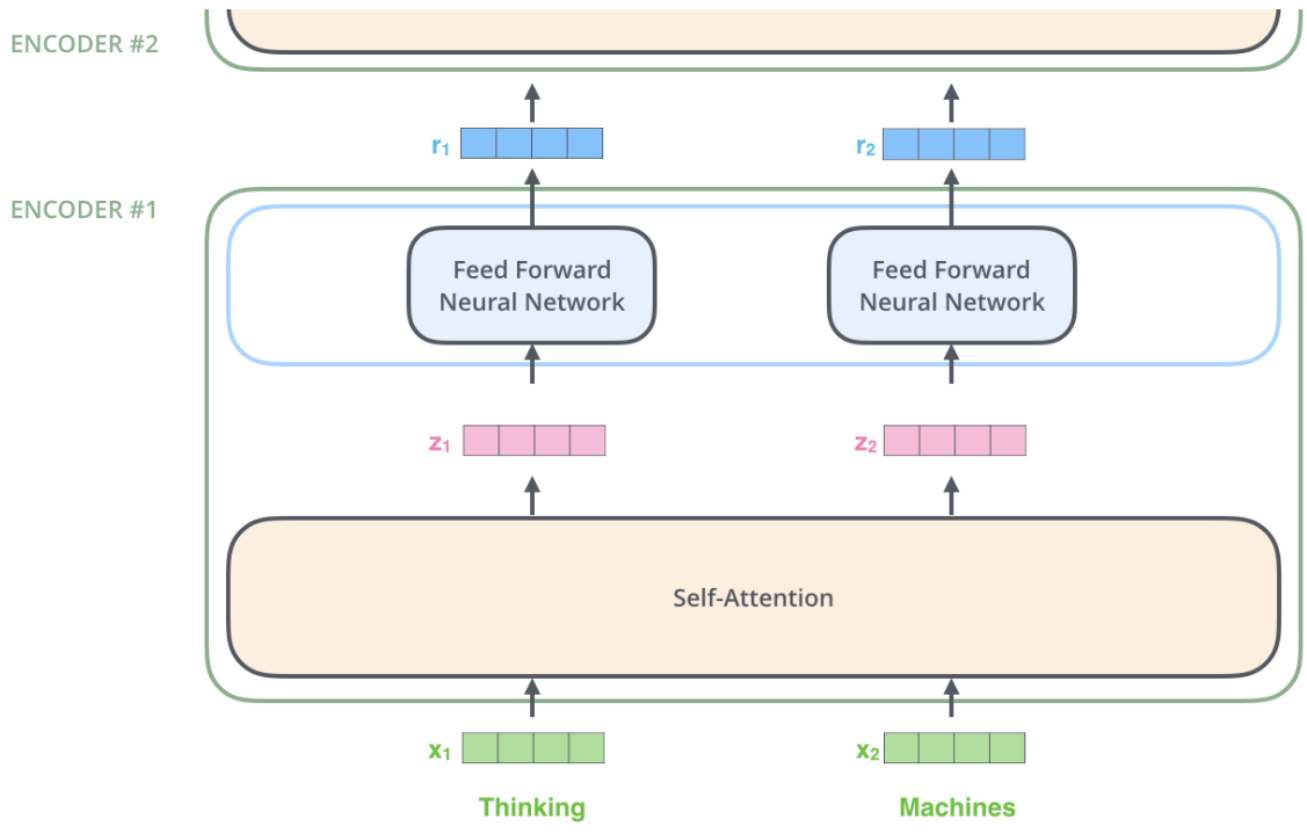
Transformer的特點

- Self-attention mechanism
- Parallel computing
- Explainable

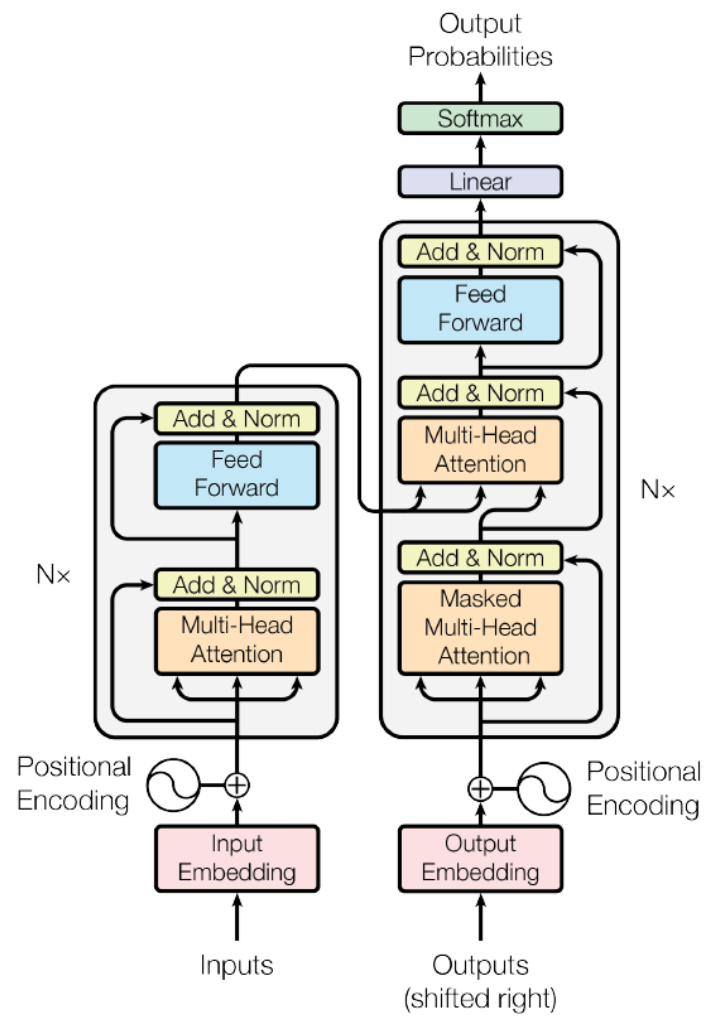




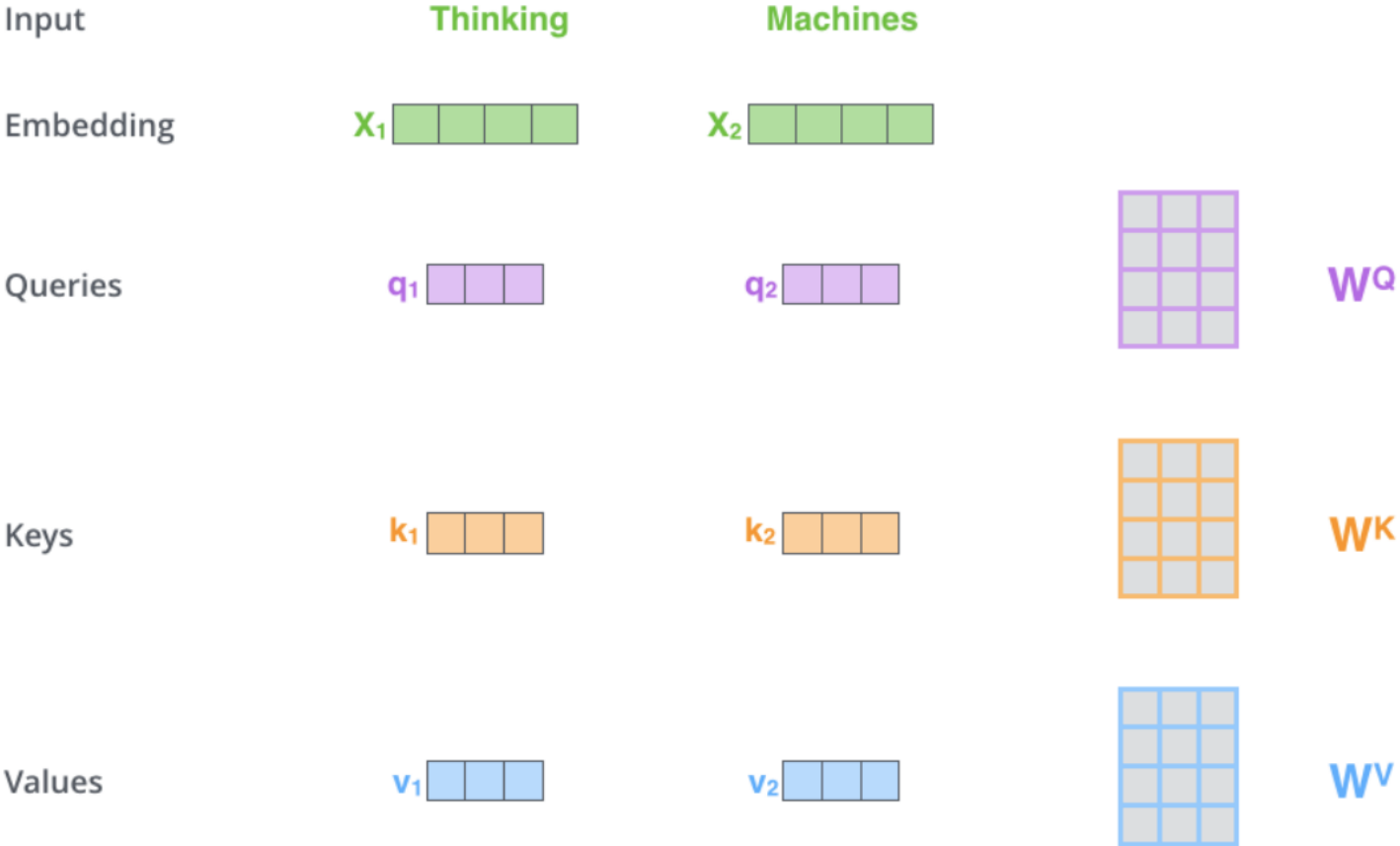




embedding algorithm



Self-attention

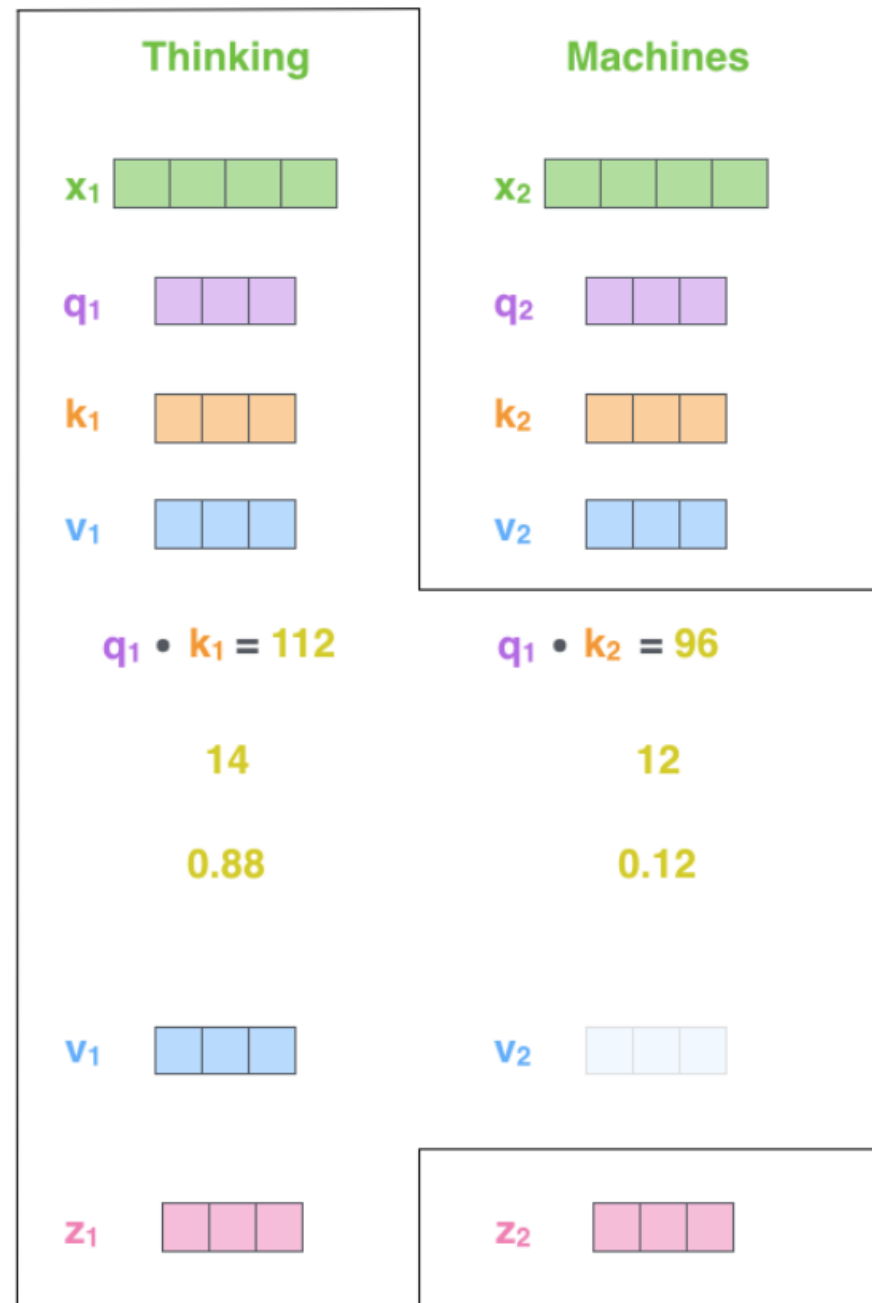


Self-attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



Input
Embedding
Queries
Keys
Values
Score
Divide by 8 ($\sqrt{d_k}$)
Softmax
Softmax
X
Value
Sum



Self-attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



$$X \times W^Q = Q$$

$$X \times W^K = K$$

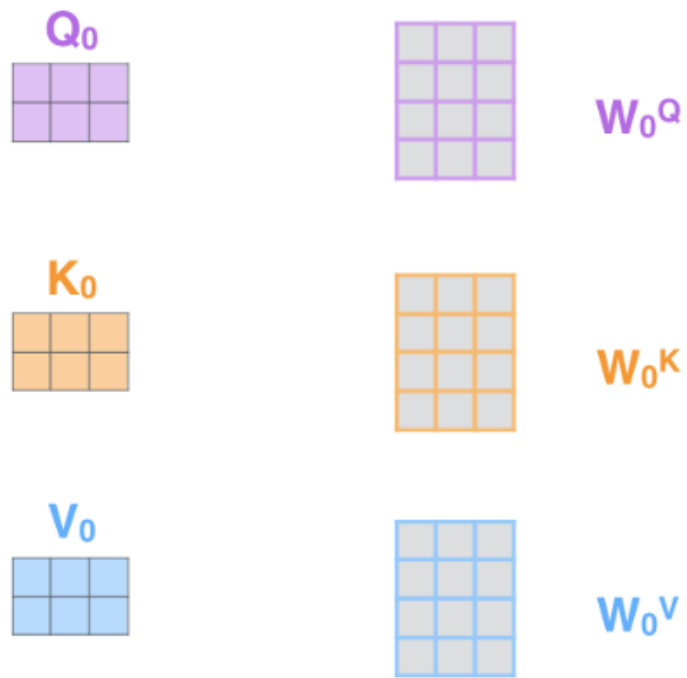
$$X \times W^V = V$$

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

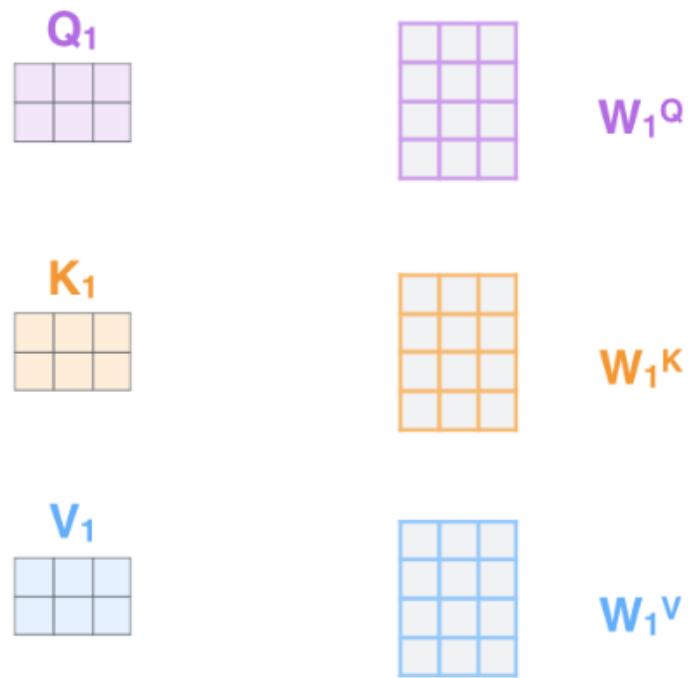
MultiHead



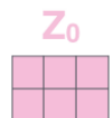
ATTENTION HEAD #0



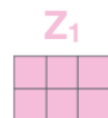
ATTENTION HEAD #1



ATTENTION HEAD #0

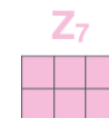


ATTENTION HEAD #1



...

ATTENTION HEAD #7



1) This is our input sentence*

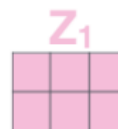
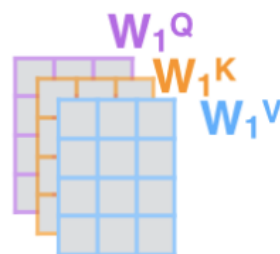
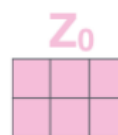
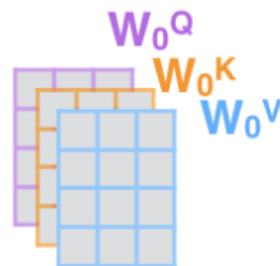
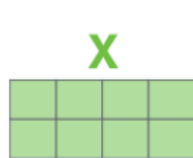
2) We embed each word*

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

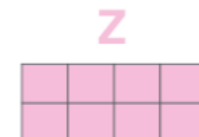
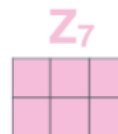
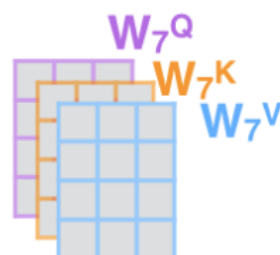
Thinking
Machines



...

...

...



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

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

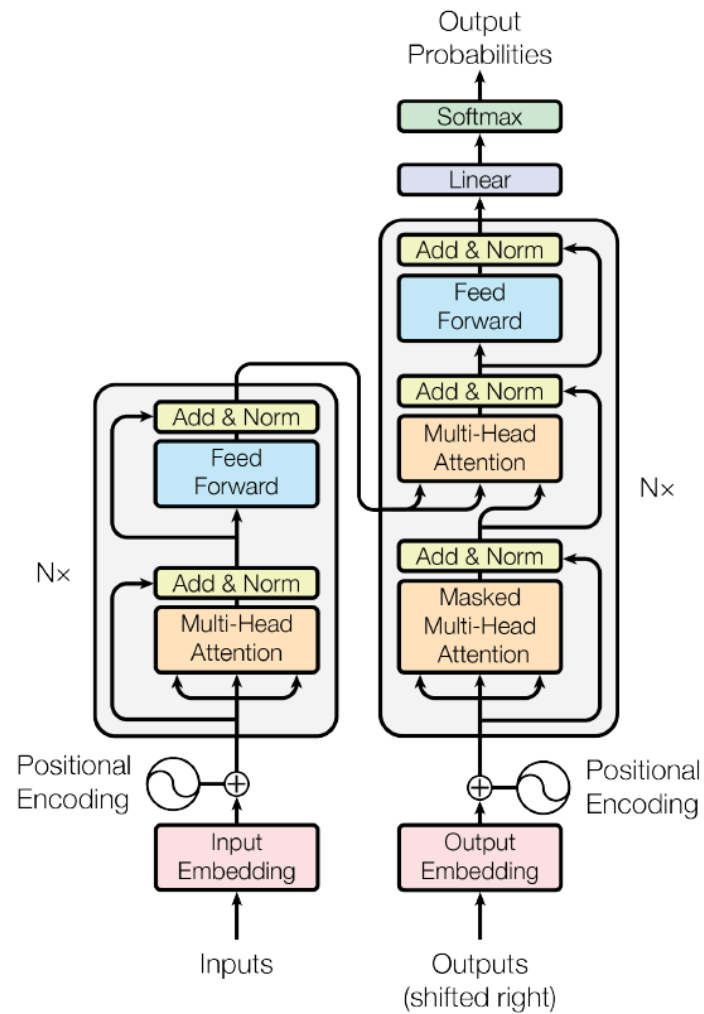
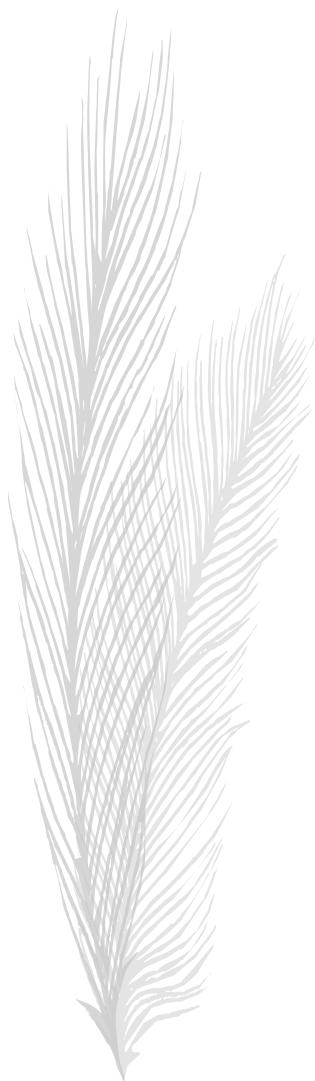
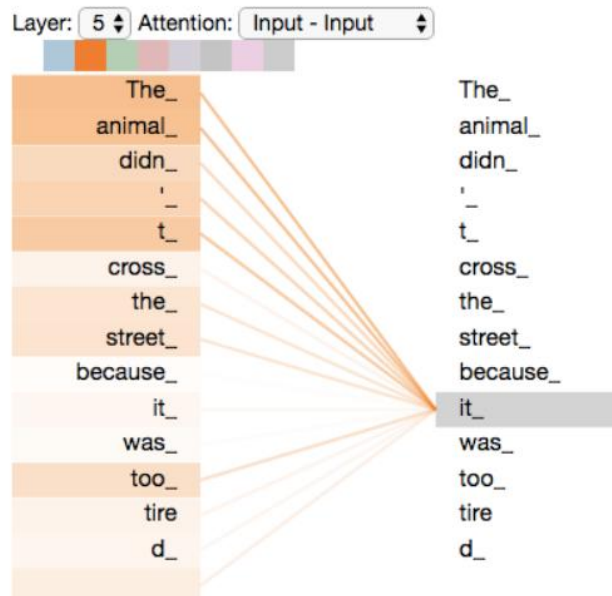
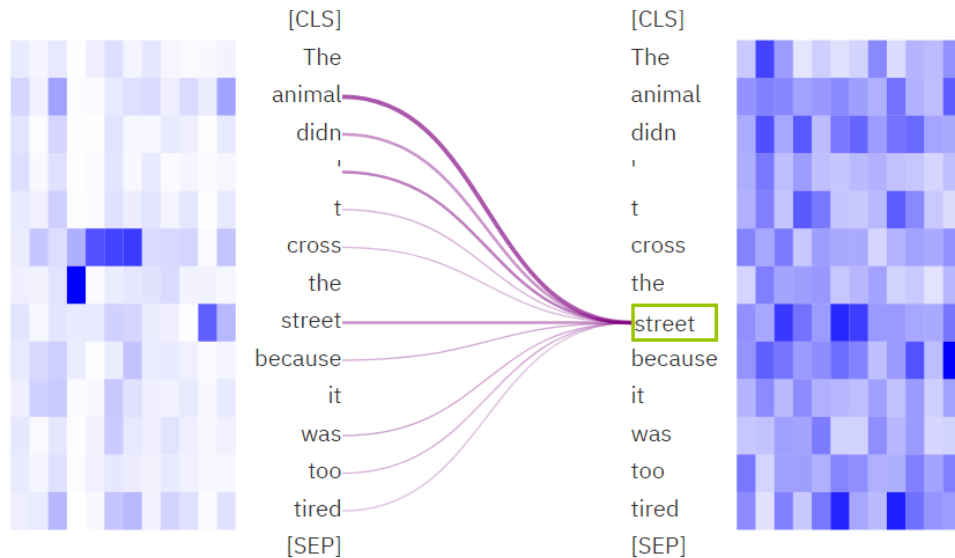


Layer

1 2 3 4 5 6 7 8 9 10 11 12

Selected heads:

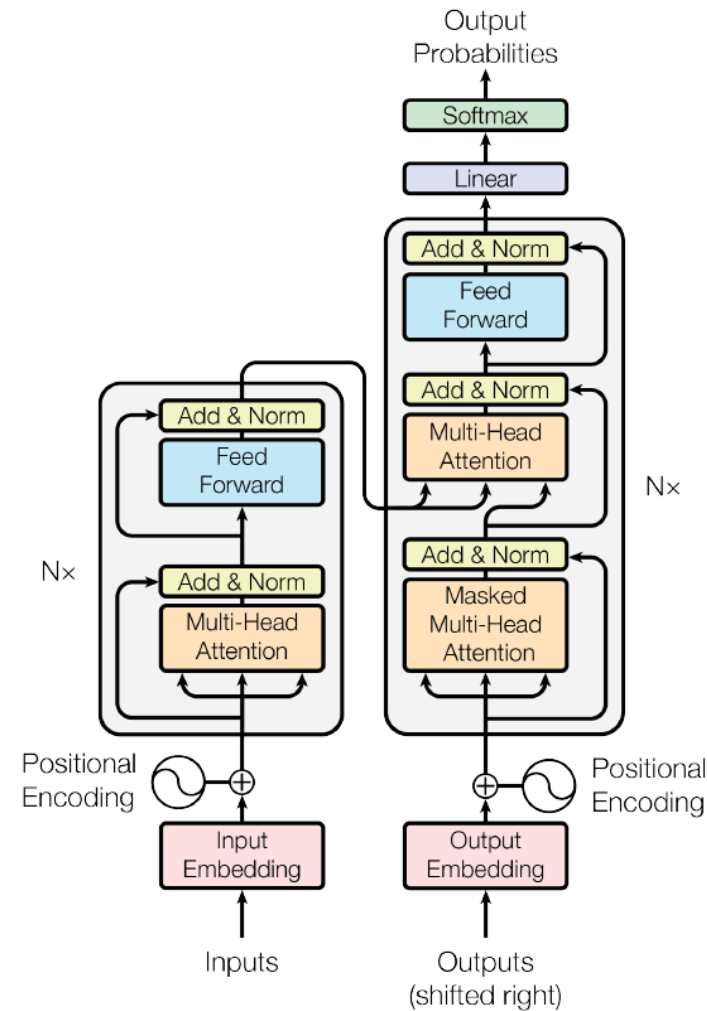
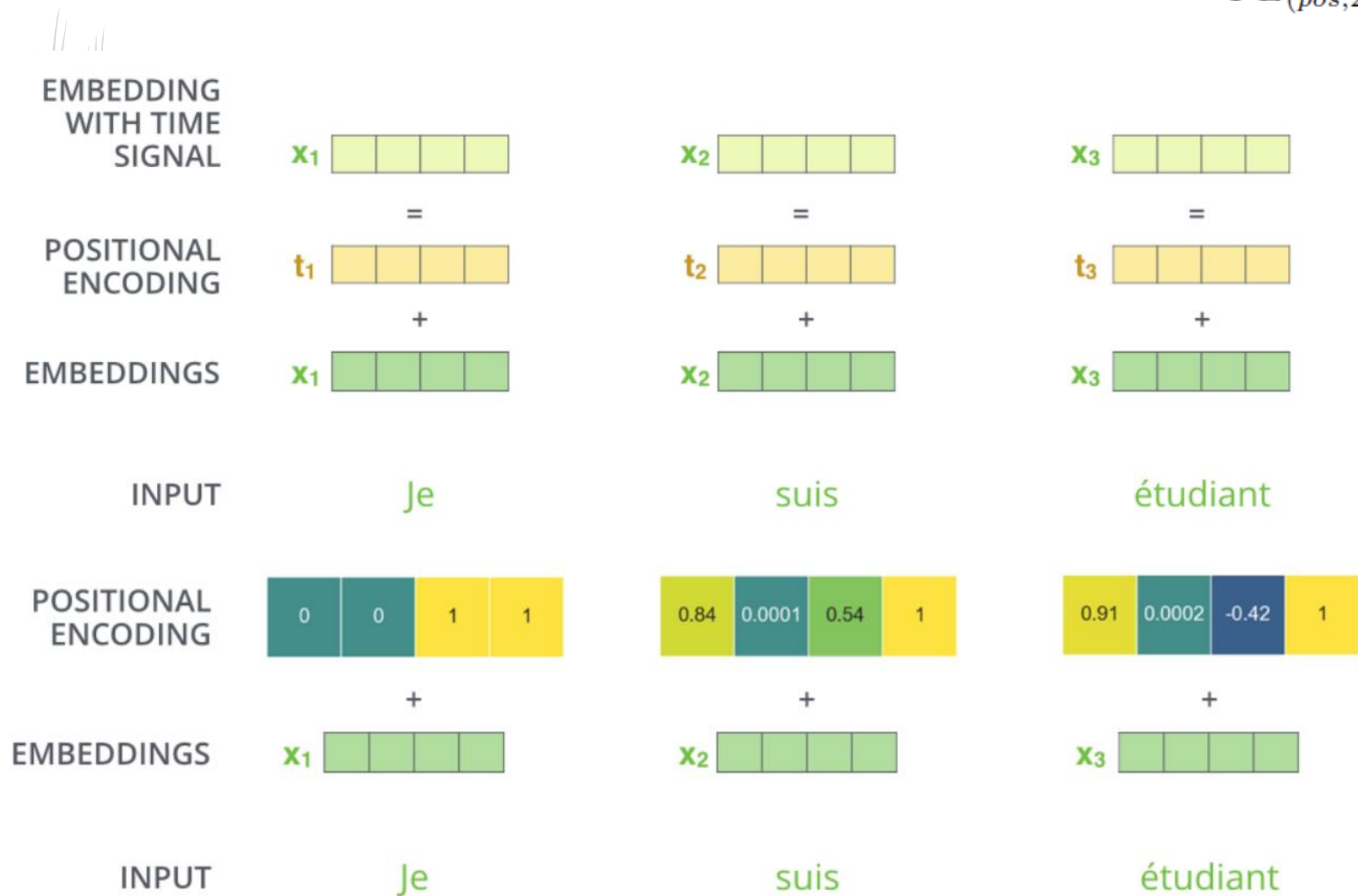
1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12



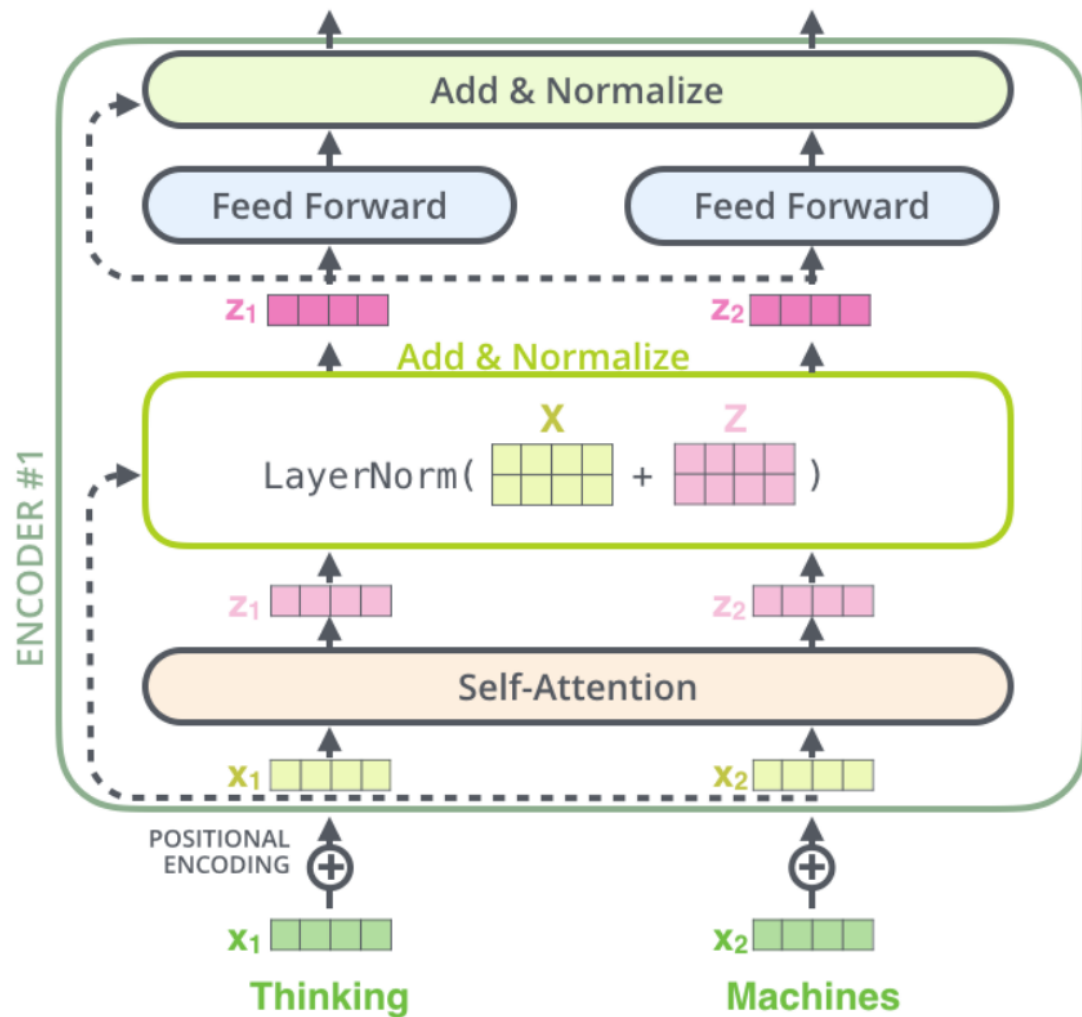
Positional Encoding

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

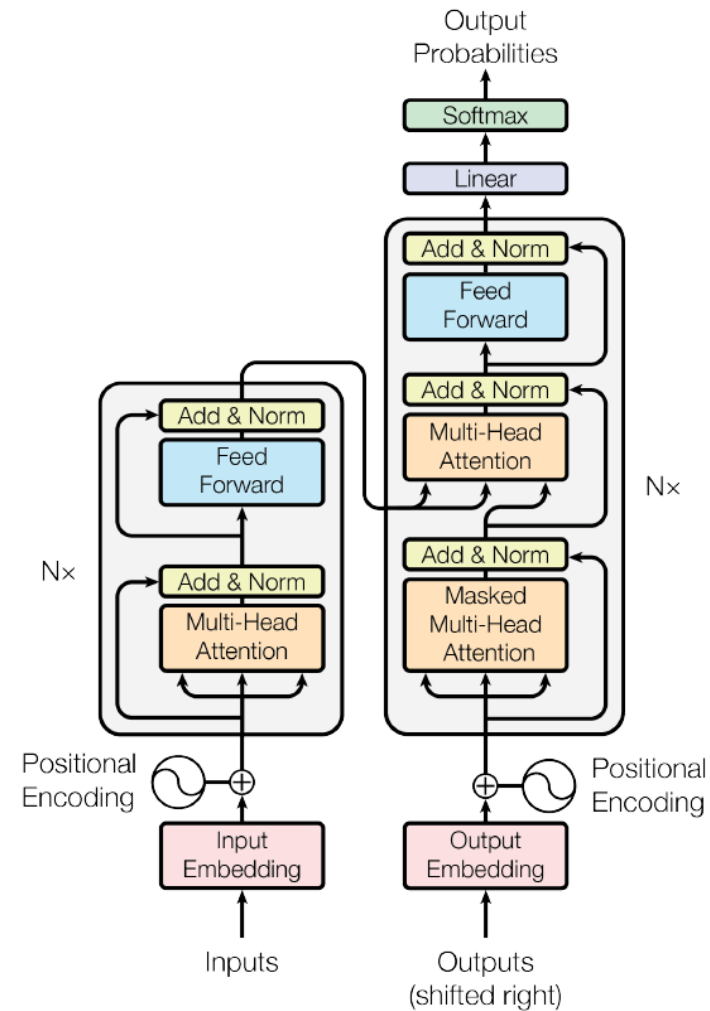
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$$



Add & Norm

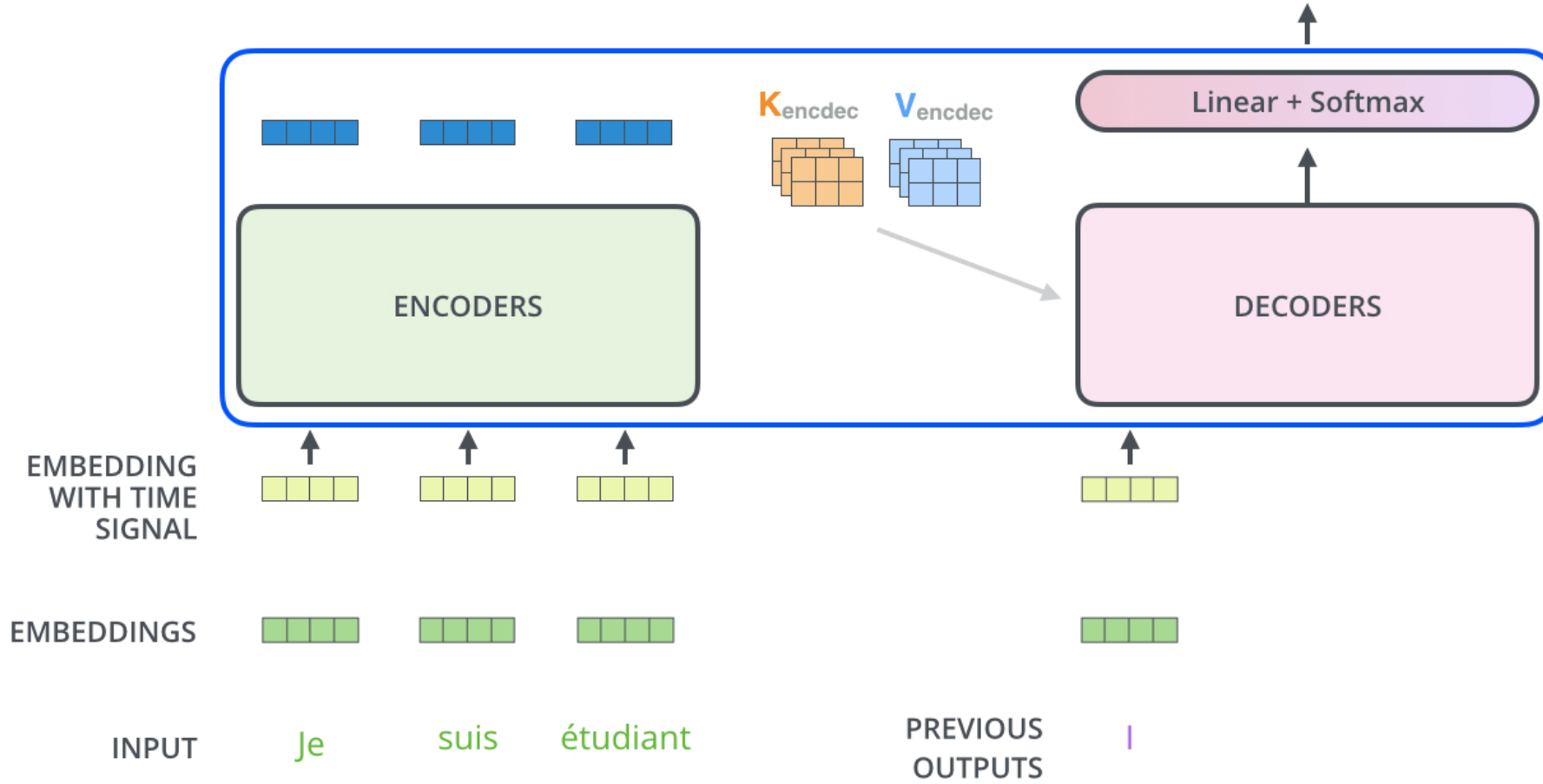
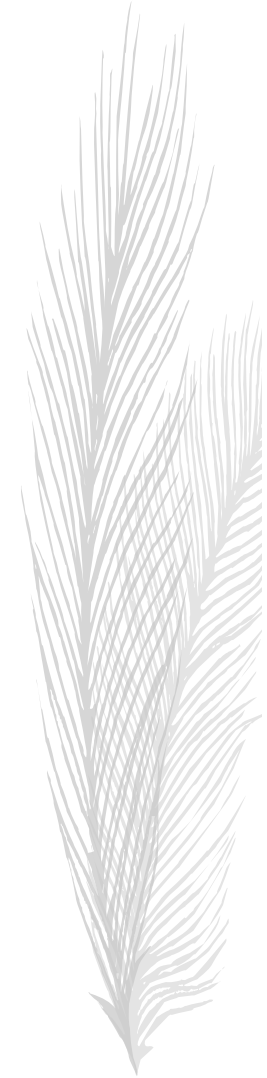


$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$



Decoding time step: 1 2 3 4 5 6

OUTPUT |



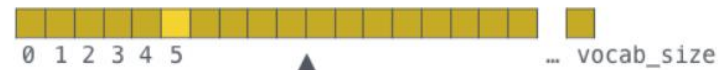
Linear & Softmax



Which word in our vocabulary is associated with this index?

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

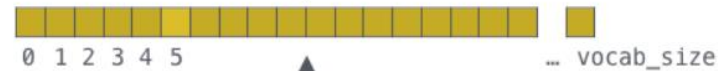
log_probs



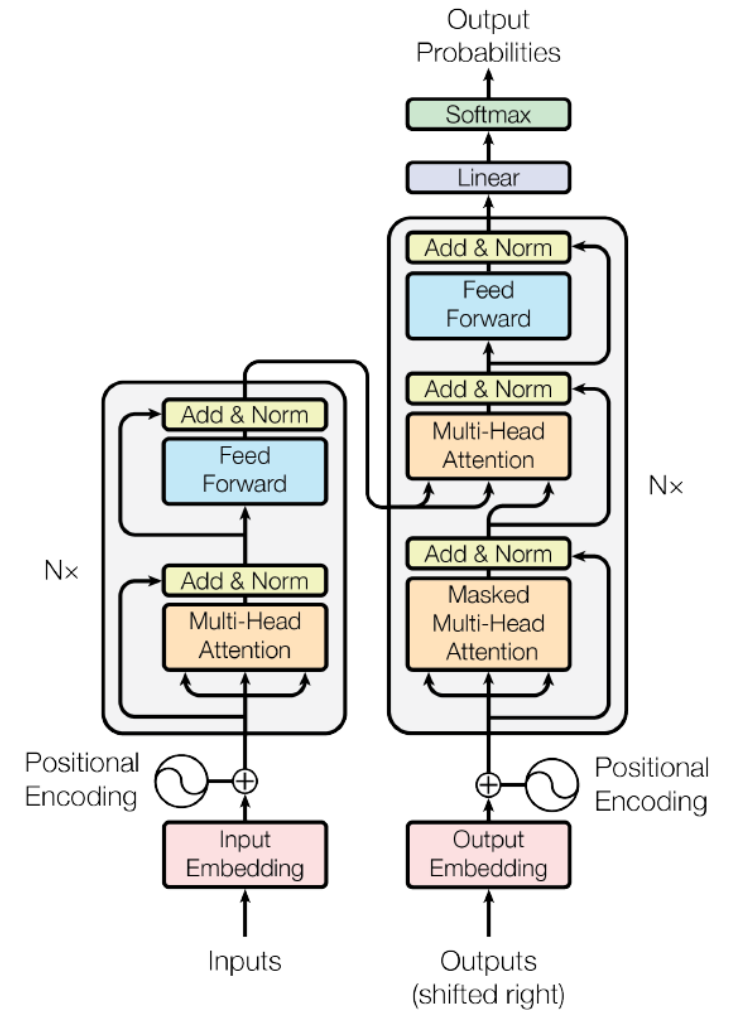
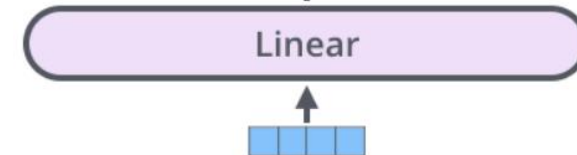
am

5

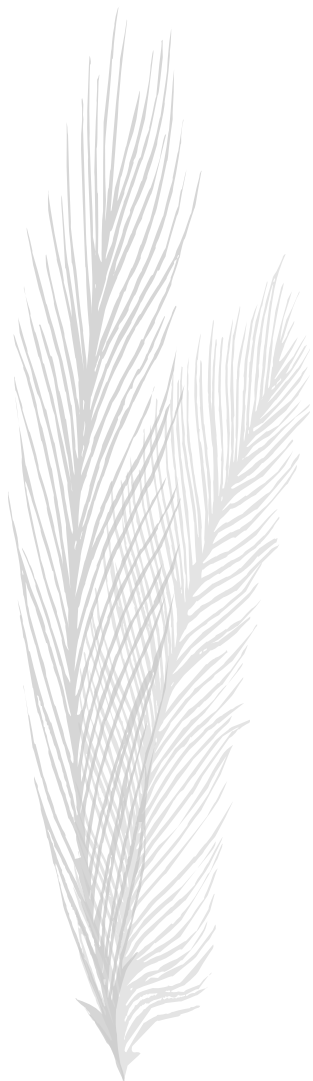
logits



Decoder stack output



Optimization



Target Model Outputs

Output Vocabulary: a am I thanks student <eos>

position #1	0.0	0.0	1.0	0.0	0.0	0.0
position #2	0.0	1.0	0.0	0.0	0.0	0.0
position #3	1.0	0.0	0.0	0.0	0.0	0.0
position #4	0.0	0.0	0.0	0.0	1.0	0.0
position #5	0.0	0.0	0.0	0.0	0.0	1.0
	a	am	I	thanks	student	<eos>

Trained Model Outputs

Output Vocabulary: a am I thanks student <eos>

position #1	0.01	0.02	0.93	0.01	0.03	0.01
position #2	0.01	0.8	0.1	0.05	0.01	0.03
position #3	0.99	0.001	0.001	0.001	0.002	0.001
position #4	0.001	0.002	0.001	0.02	0.94	0.01
position #5	0.01	0.01	0.001	0.001	0.001	0.98
	a	am	I	thanks	student	<eos>

Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types. n is the sequence length, d is the representation dimension, k is the kernel size of convolutions and r the size of the neighborhood in restricted self-attention.

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

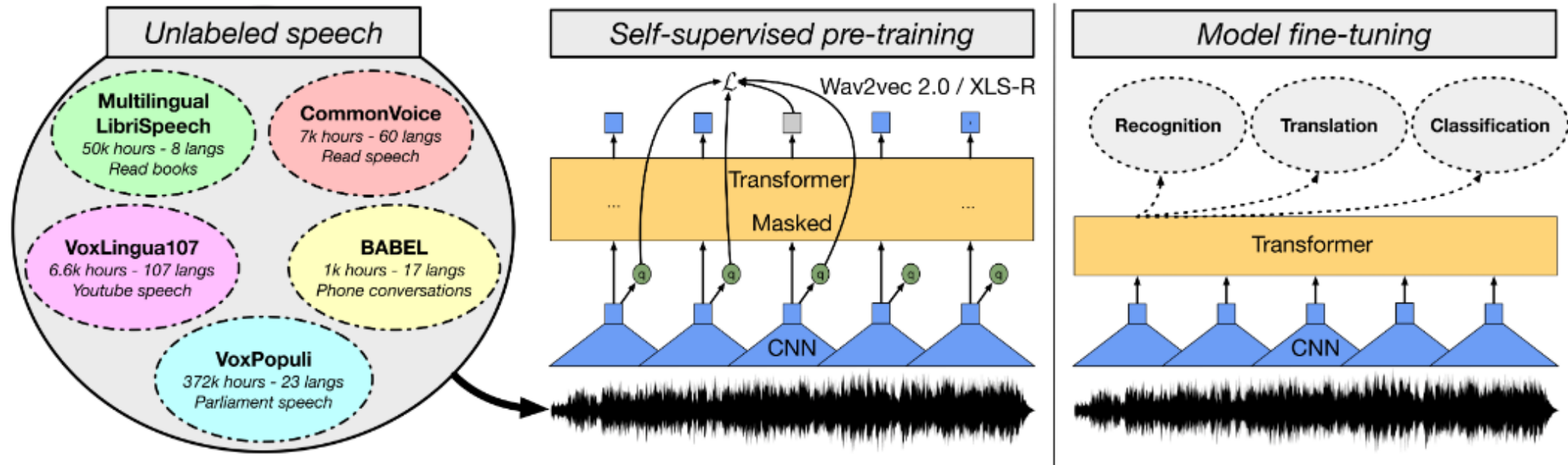
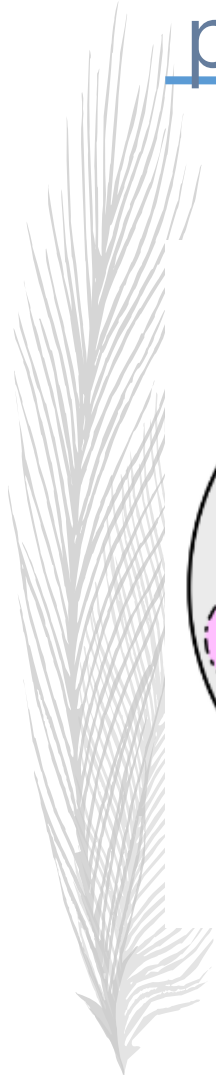
Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	$3.3 \cdot 10^{18}$	
Transformer (big)	28.4	41.8	$2.3 \cdot 10^{19}$	

Transformer的應用

- 預訓練模型：GPT、BERT
- 機器翻譯
- 文本摘要
- 情感分析
- 問答任務
- 語音識別和語音生成
- 電腦視覺

Facebook's Wav2Vec2 XLS-R counting 300 million parameters.



Reference

- Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems 30* (2017).
- <https://medium.com/deeper-learning/glossary-of-deep-learning-word-embedding-f90c3cec34ca>
- <https://huggingface.co/spaces/exbert-project/exbert>
- <https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>
- https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb#scrollTo=r6GPPFy1fL2N
- <http://jalammar.github.io/illustrated-transformer/>