

AI in Built Environment

DCP4300

Week 12: Natural Language Processing

Transformers

Dr. Chaofeng Wang
Jianhao Gao (TA)

University of Florida
College of Design Construction and Planning

Word embeddings

Word2vec, 2013

Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. arXiv preprint arXiv:1301.3781.

<https://arxiv.org/abs/1301.3781>

GloVe, 2014

Pennington, J., Socher, R., & Manning, C. D. (2014, October). Glove: Global vectors for word representation. In Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP) (pp. 1532-1543).

<https://aclanthology.org/D14-1162.pdf>

RNNs

LSTM, 2000

Gers, F. A., Schmidhuber, J., & Cummins, F. (2000). Learning to forget: Continual prediction with LSTM. Neural computation, 12(10), 2451-2471.
<https://ieeexplore.ieee.org/abstract/document/6789445>

GRU, 2014

Chung, J., Gulcehre, C., Cho, K., & Bengio, Y. (2014). Empirical evaluation of gated recurrent neural networks on sequence modeling. arXiv preprint arXiv:1412.3555. <https://arxiv.org/abs/1412.3555>

Transformers (BERT, GPT)

2017/2018 -

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). **Attention is all you need**. Advances in neural information processing systems, 30. <https://arxiv.org/abs/1706.03762>

Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
<https://arxiv.org/abs/1810.04805>

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. Advances in neural information processing systems, 33, 1877-1901.
<https://arxiv.org/abs/2005.14165>

Self-attention layer:

Help encoder look at other words in the input sentences as it encodes a specific word.

Language models

ChatGPT Decoder:

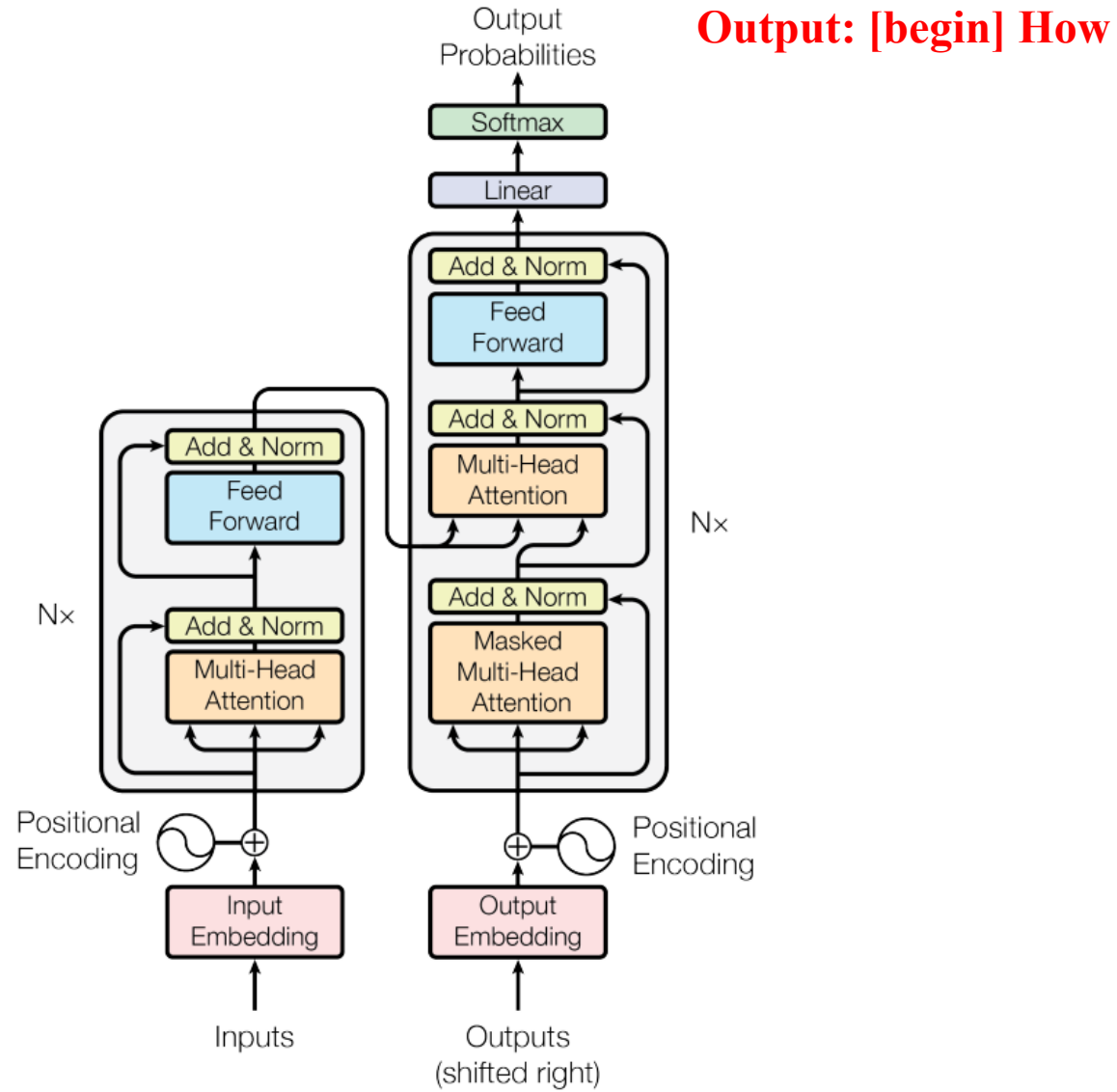
predict the next word in a sentence or complete a partial sentence

$$P(w_1 w_2 \dots w_n) = P(w_n \mid w_1 \dots w_{n-1}) * P(w_{n-1} \mid w_1 \dots w_{n-2}) \dots P(w_2 \mid w_1) * P(w_1)$$

$$P(\text{The quick brown fox jumps over the lazy dog}) =$$

$$P(\text{The}) * P(\text{quick} \mid \text{The}) * P(\text{brown} \mid \text{The quick}) * P(\text{fox} \mid \text{The quick brown}) * P(\text{jumps} \mid \text{The quick brown fox}) * P(\text{over} \mid \text{The quick brown fox jumps}) * P(\text{the} \mid \text{The quick brown fox jumps over}) * P(\text{lazy} \mid \text{The quick brown fox jumps over the}) * P(\text{dog} \mid \text{The quick brown fox jumps over the lazy})$$

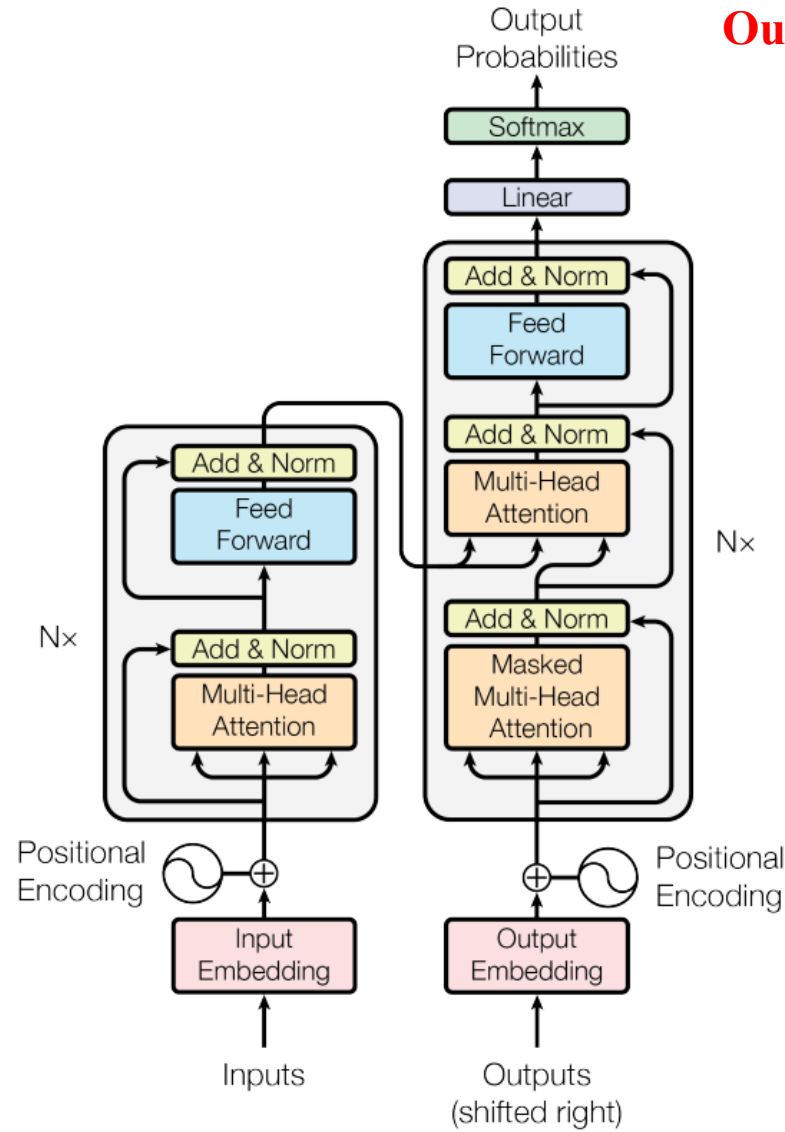
predict the next word



Input: Give me a random sentence

Output: [begin]

predict the next word

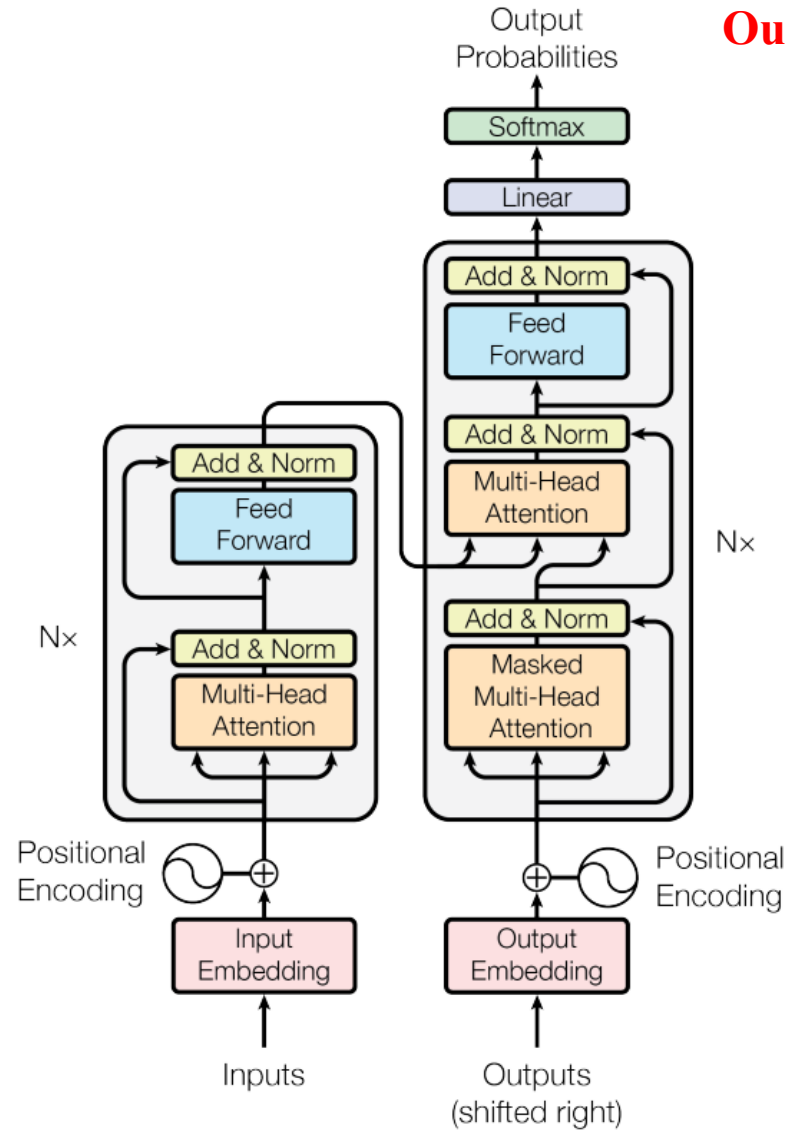


Output: [begin] How may

Input: Give me a random sentence

Output: [begin] How

predict the next word

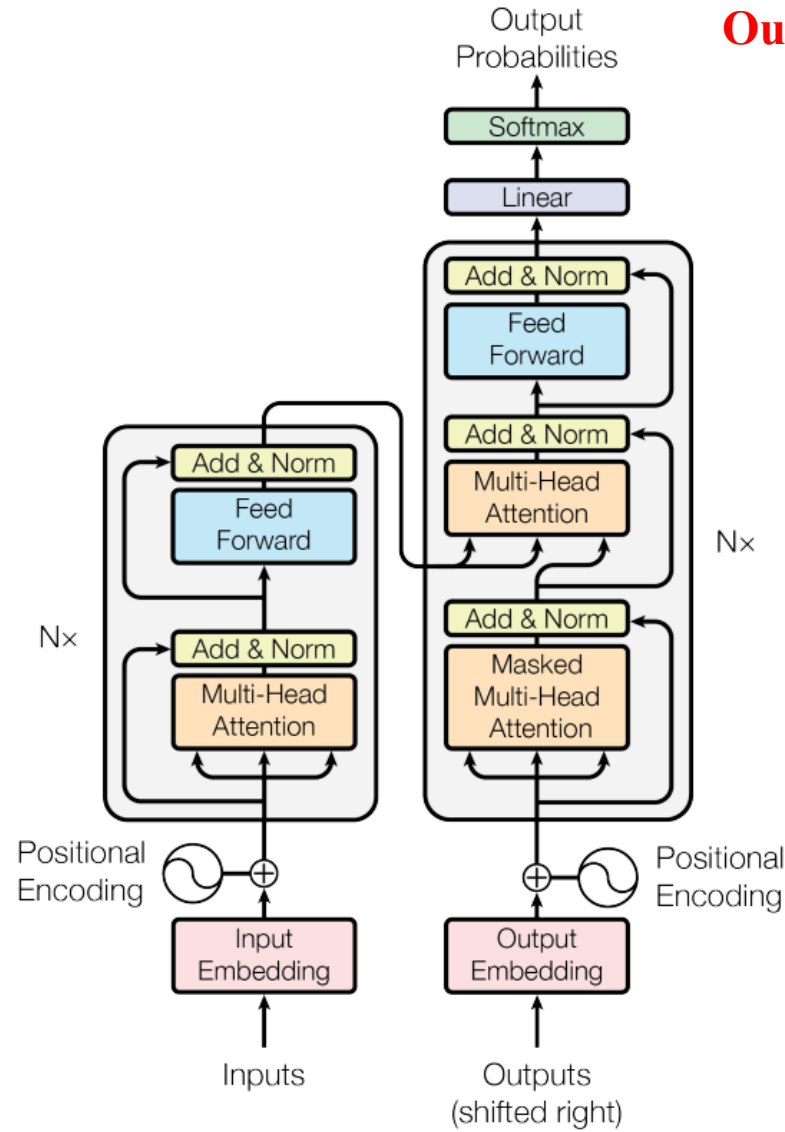


Output: [begin] How may I

Input: Give me a random sentence

Output: [begin] How may

predict the next word

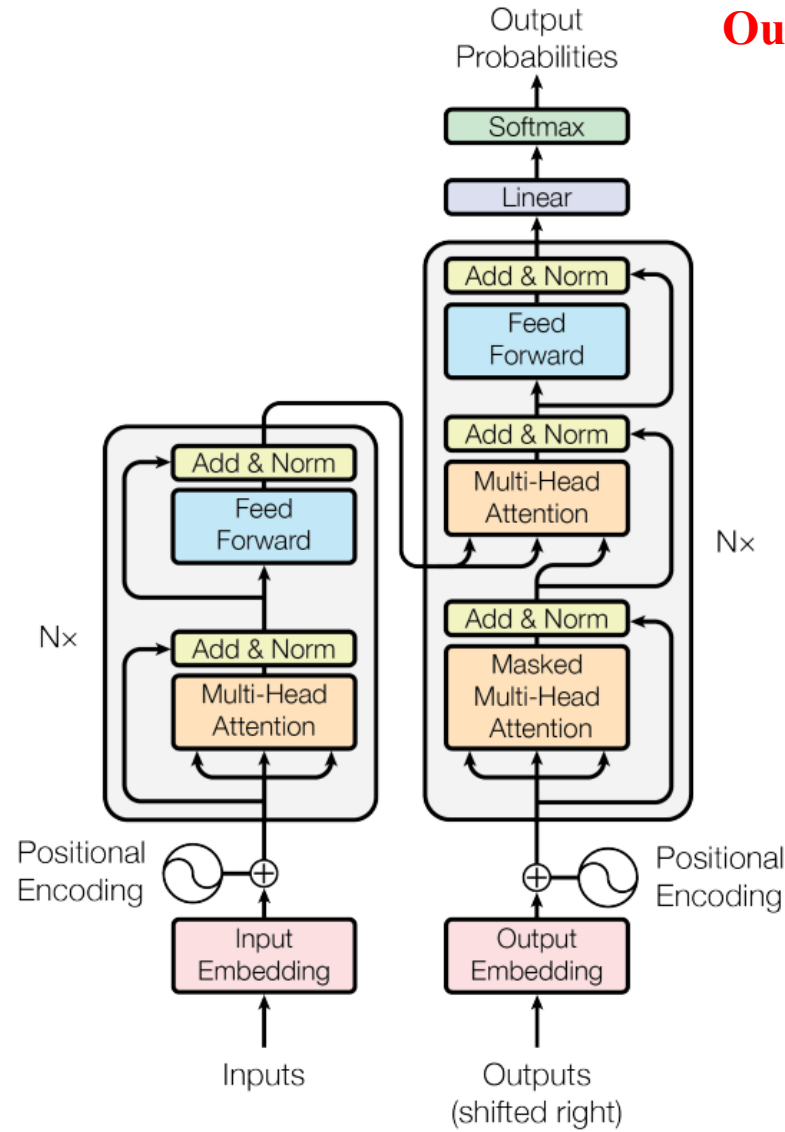


Output: [begin] How may I help

Input: Give me a random sentence

Output: [begin] How may I

predict the next word

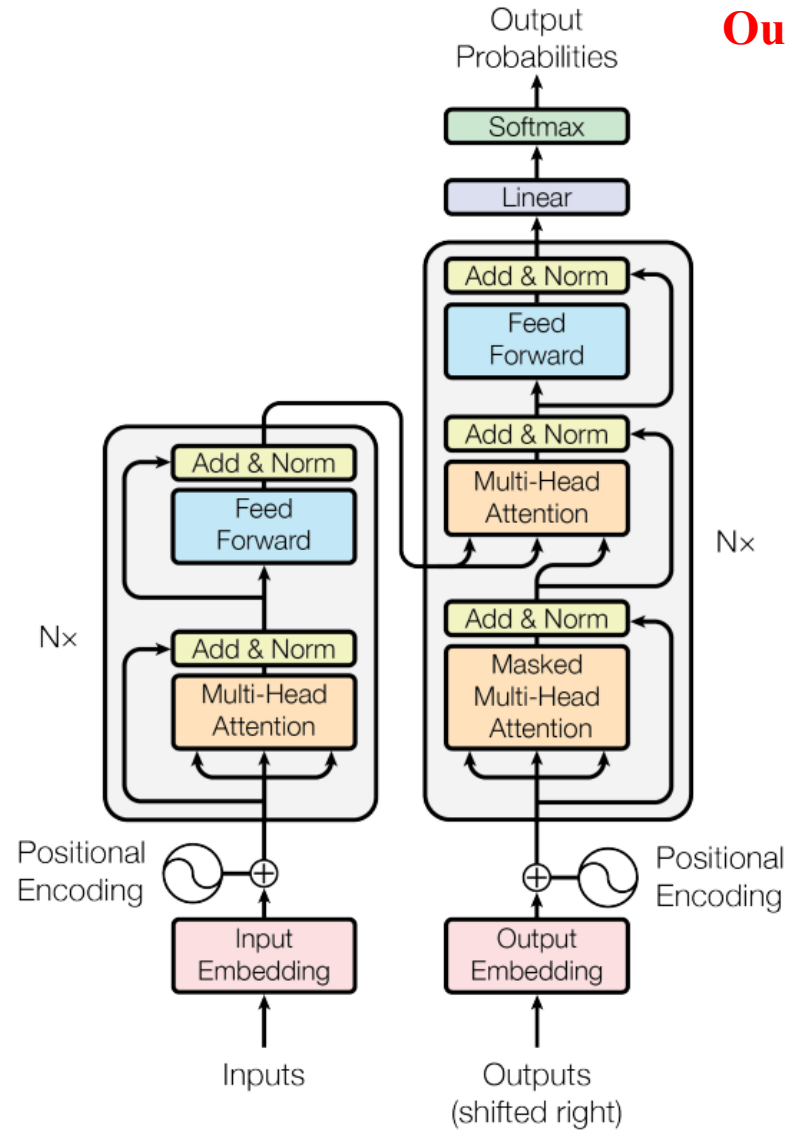


Output: [begin] How may I help you

Input: Give me a random sentence

Output: [begin] How may I help

predict the next word

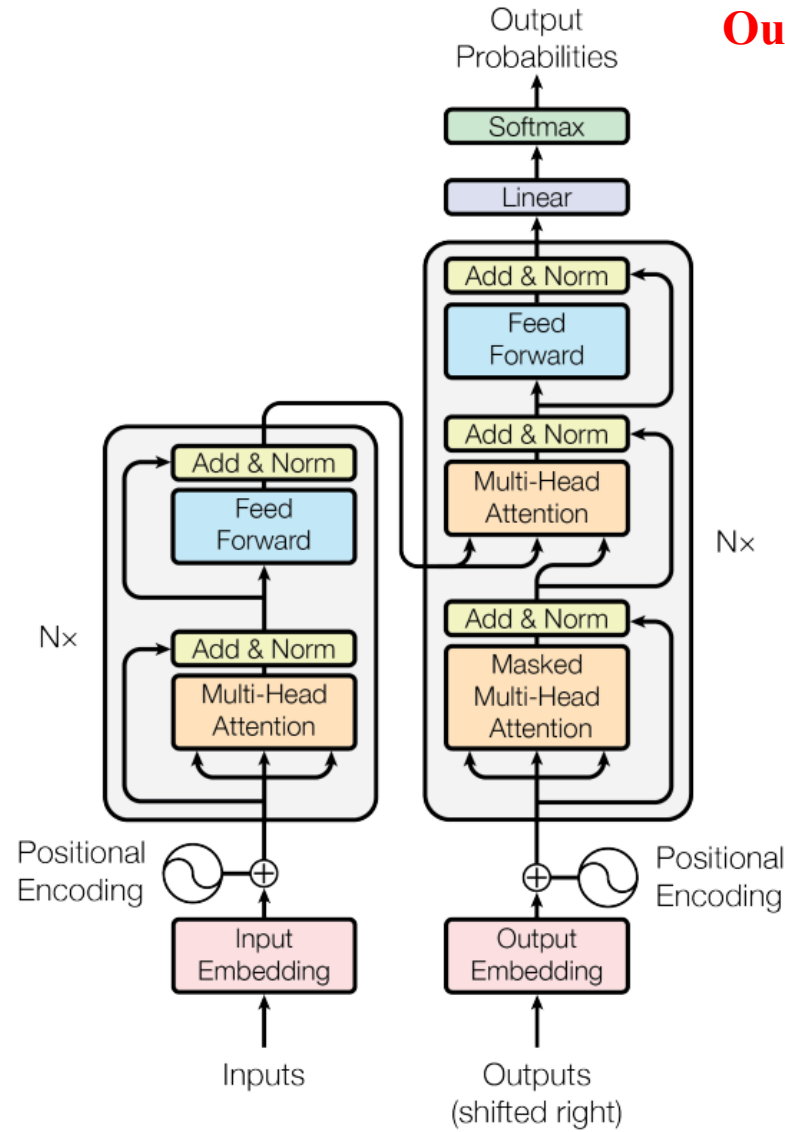


Output: [begin] How may I help you?

Input: Give me a random sentence

Output: [begin] How may I help you

predict the next word



Output: [begin] How may I help you? [end]

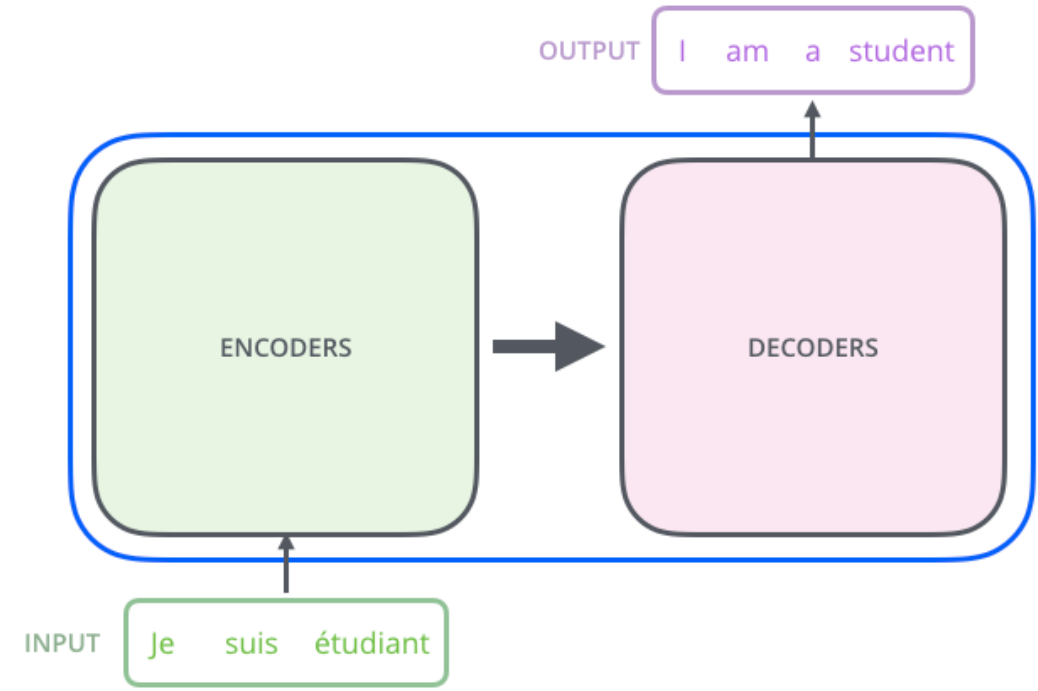
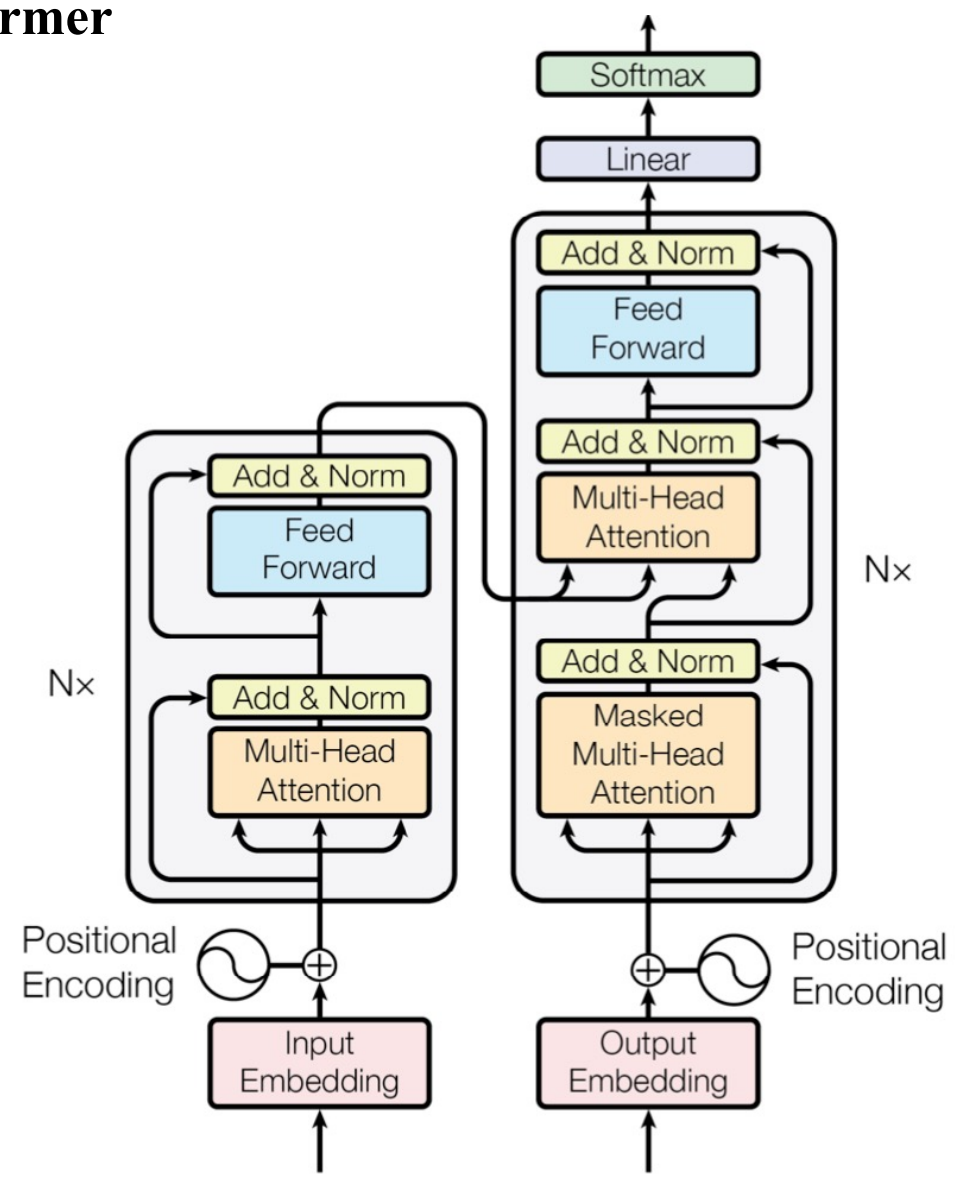
Input: Give me a random sentence

Output: [begin] How may I help you?

- A good video to illustrate the structure:

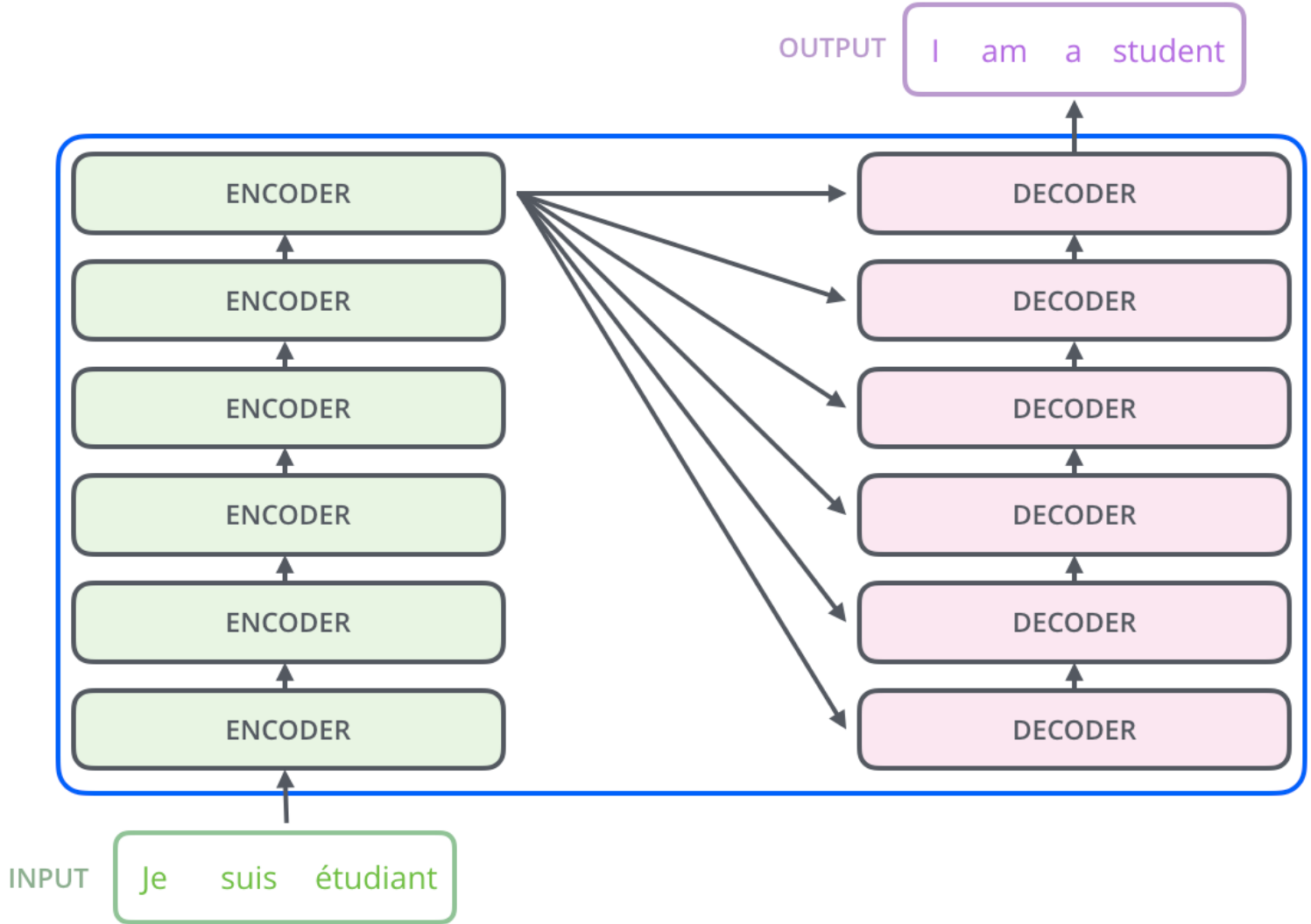
Illustrated Guide to Transformers Neural Network: A step by step explanation <https://www.youtube.com/watch?v=4Bdc55j80l8>

Transformer



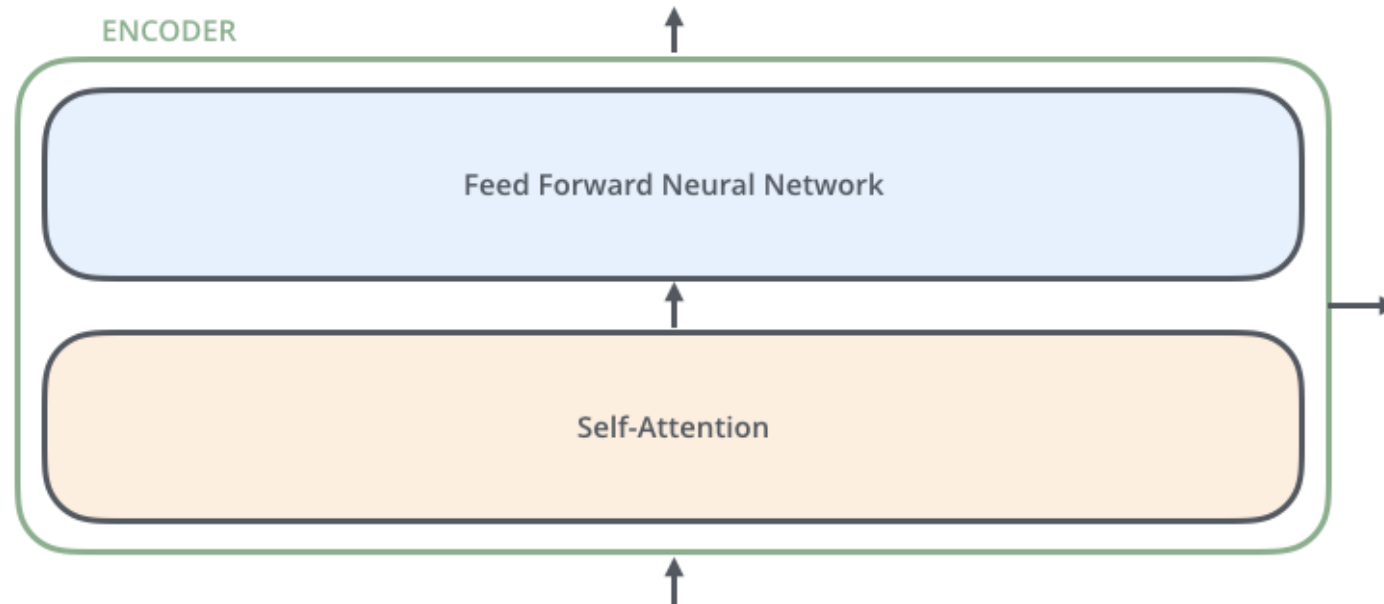
Transformer

Do not share weights.



Transformer

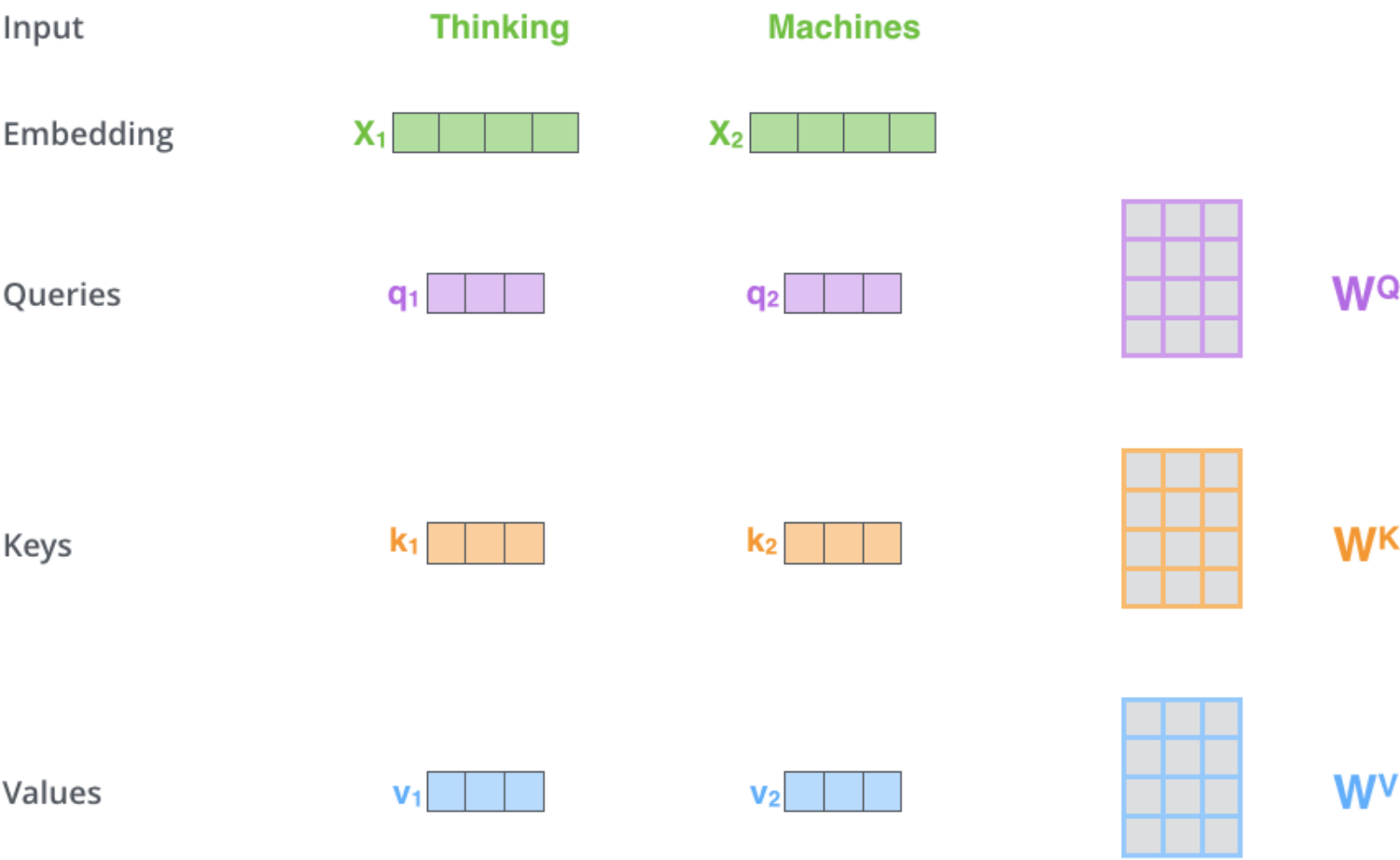
Take a closer look into an encoder.



Self-attention layer: A layer that helps the encoder look at other words in the input sentence as it encodes a specific word

Transformer

Self-attention



Transformer

Self-attention

Input

Embedding

Queries

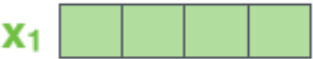
Keys

Values

Score

Thinking

Machines



$q_1 \cdot k_1 = 112$

$q_1 \cdot k_2 = 96$

Transformer

Self-attention

Input

Embedding

Queries

Keys

Values


Score

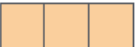
Divide by 8 ($\sqrt{d_k}$)


Softmax

Thinking

x_1 

q_1 

k_1 

v_1 

$q_1 \cdot k_1 = 112$

14

0.88

Machines

x_2 

q_2 

k_2 

v_2 

$q_1 \cdot k_2 = 96$

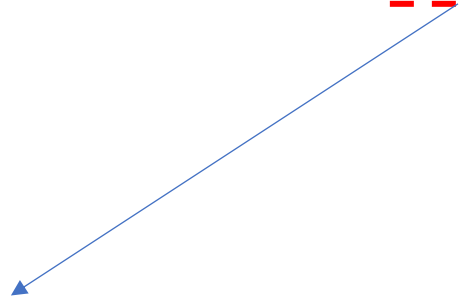
12

0.12

Transformer

Self-attention

I saw Anchorage flying back to China.



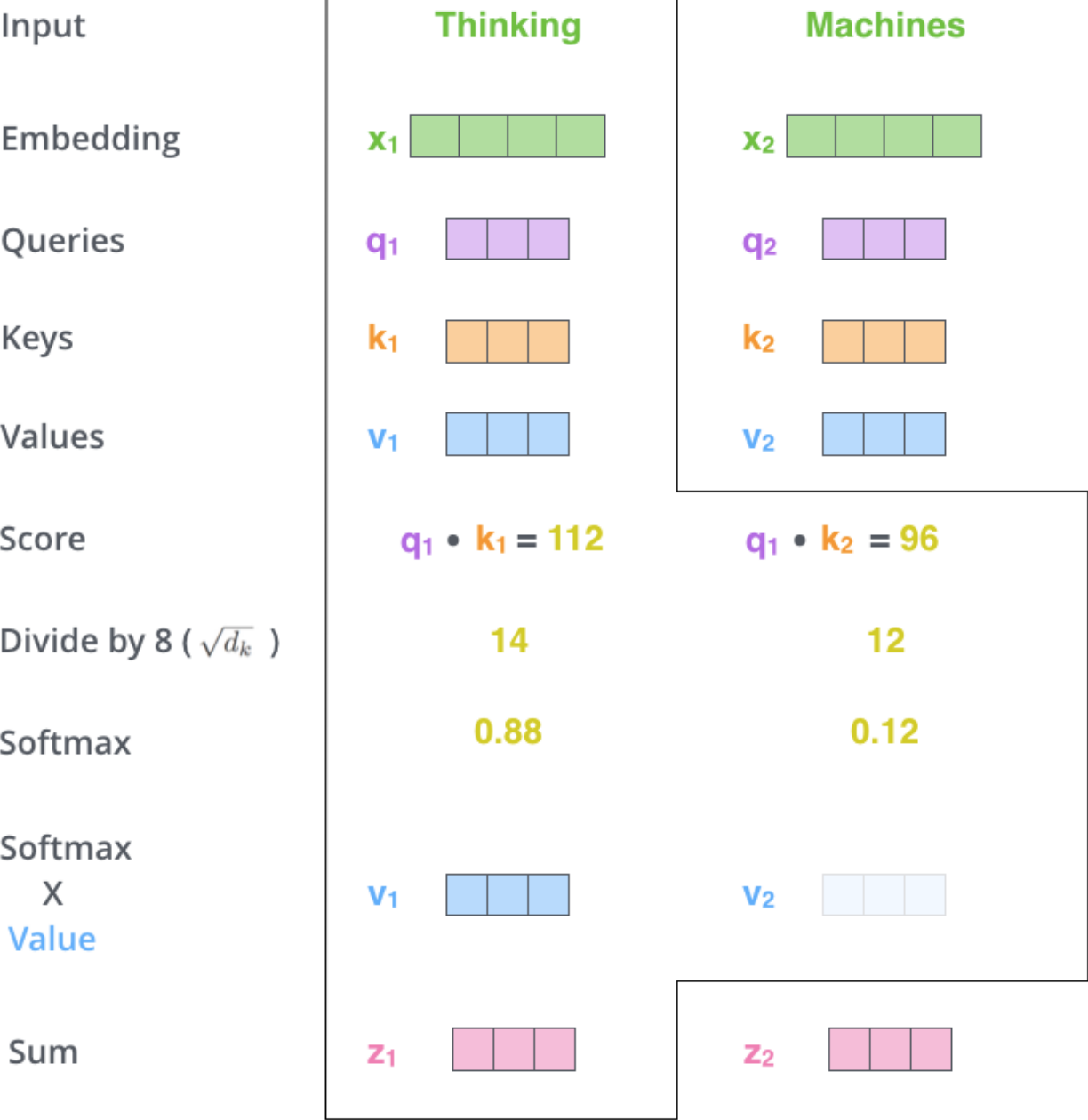
I saw Anchorage flying back to China.



Attention vector used to encode/embed/understand the word 'flying' , in the context of this sentence.

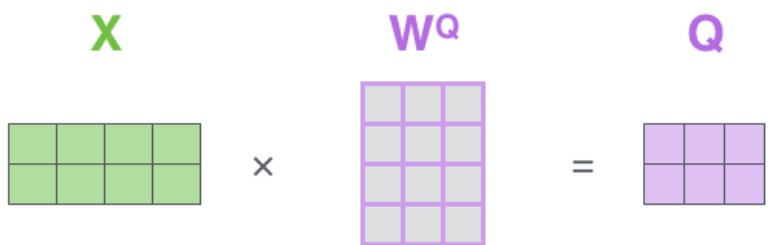
Transformer

Self-attention



Transformer

Self-attention
Matrix Calculation

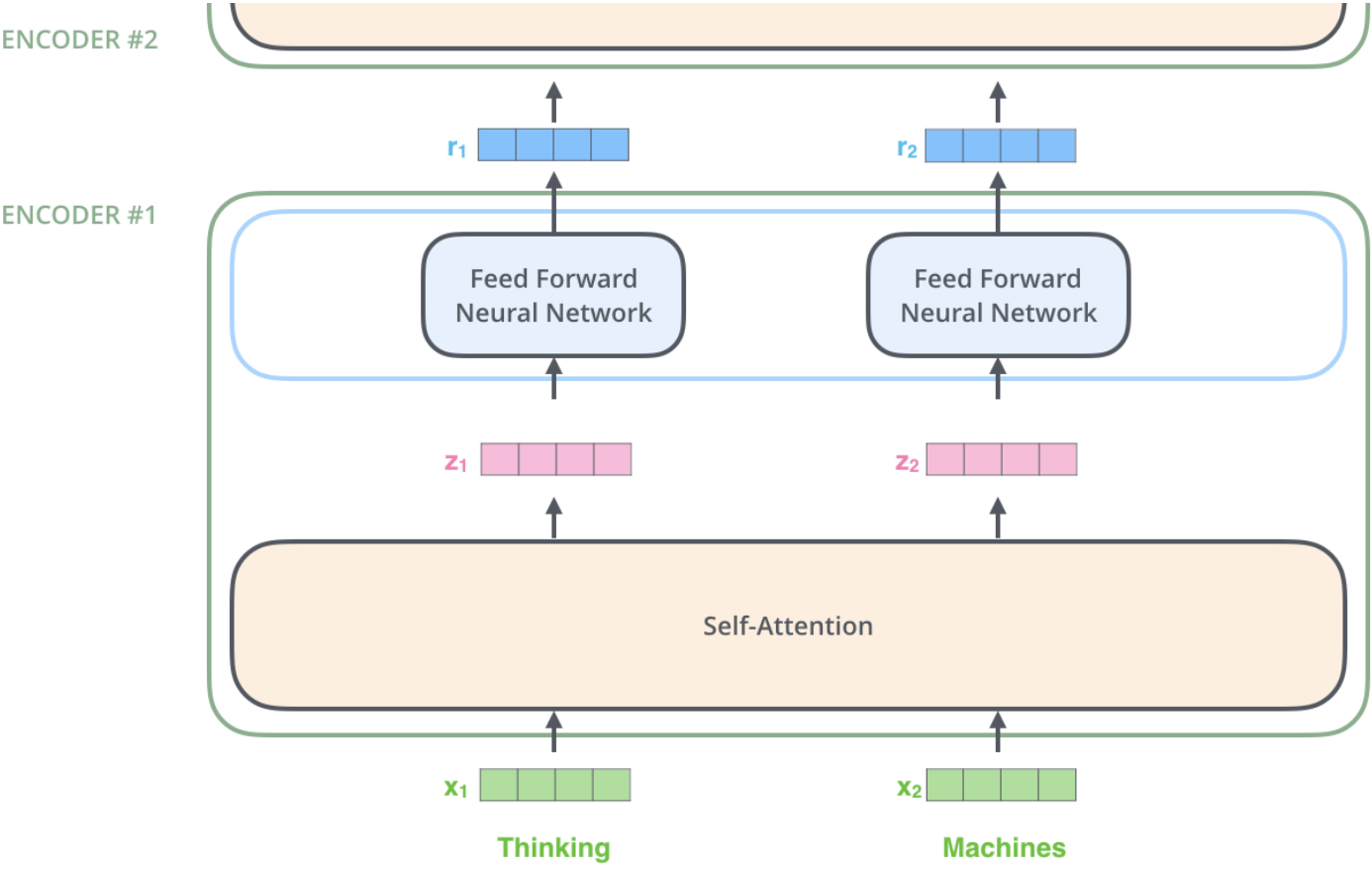


Transformer

Self-attention
Matrix Calculation

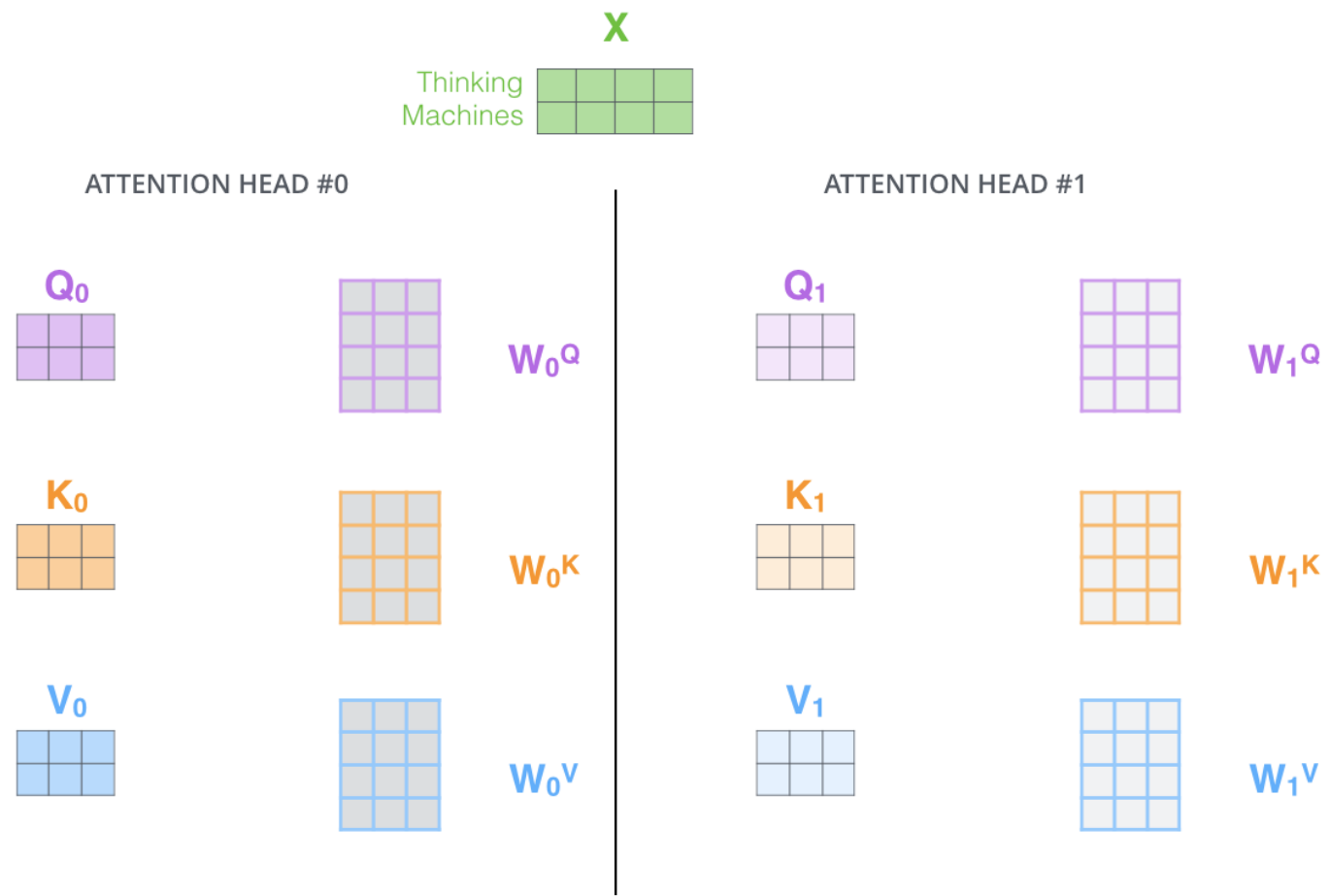
$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \times \begin{matrix} \text{K}^T \\ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline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

Transformer



Transformer

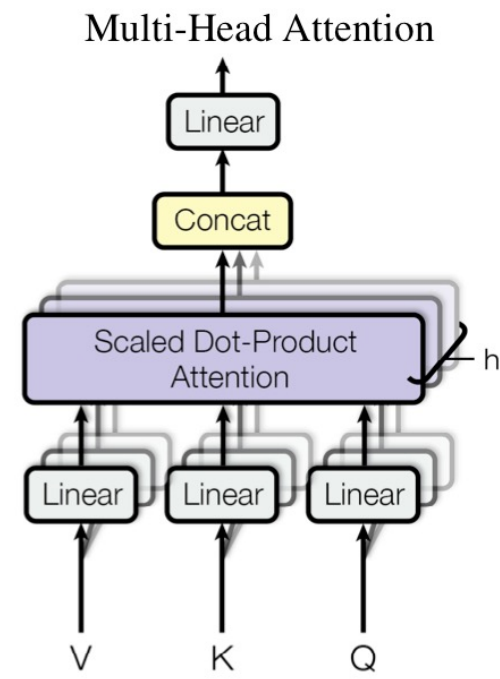
Multi-head attention



Each head can project the input into a **different representation subspace**

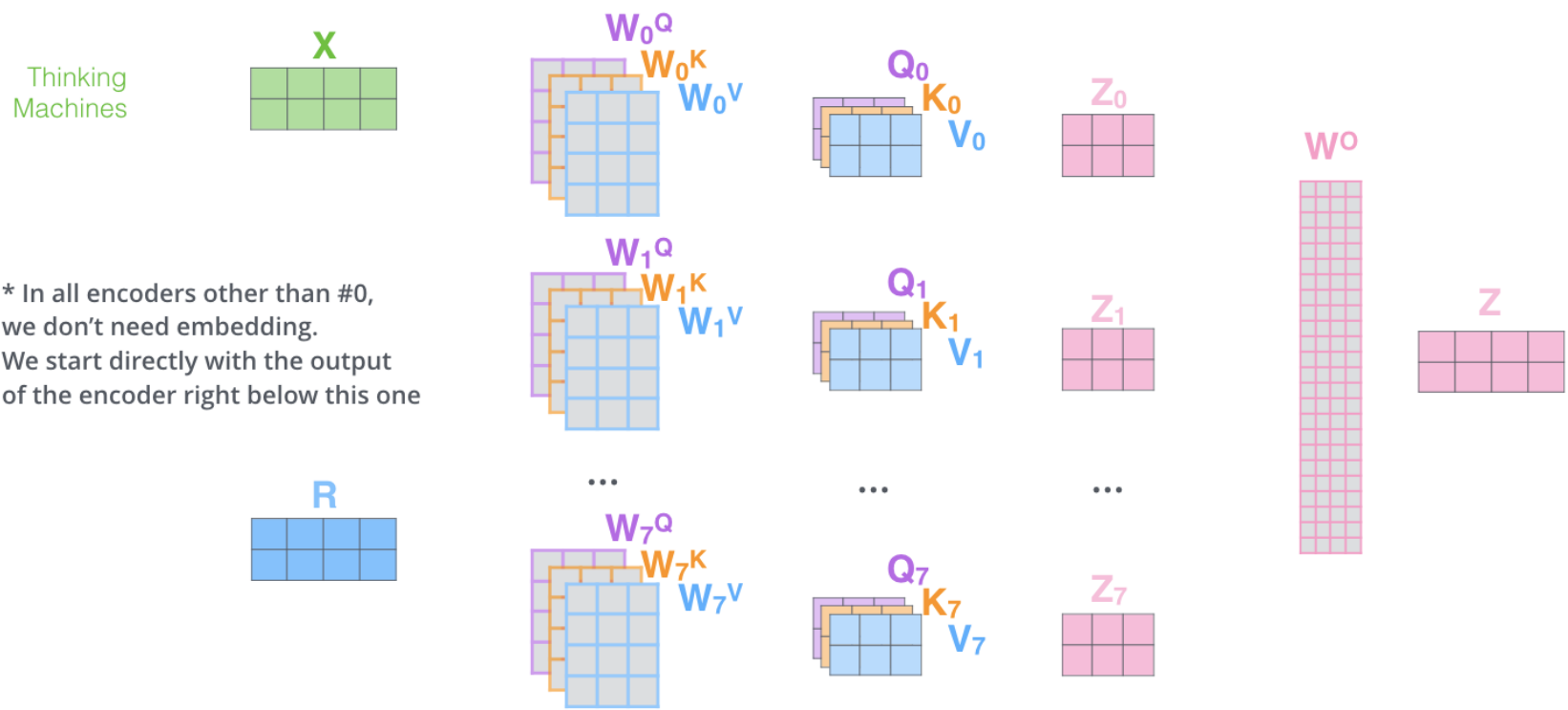
Transformer

Multi-head attention



Detail->

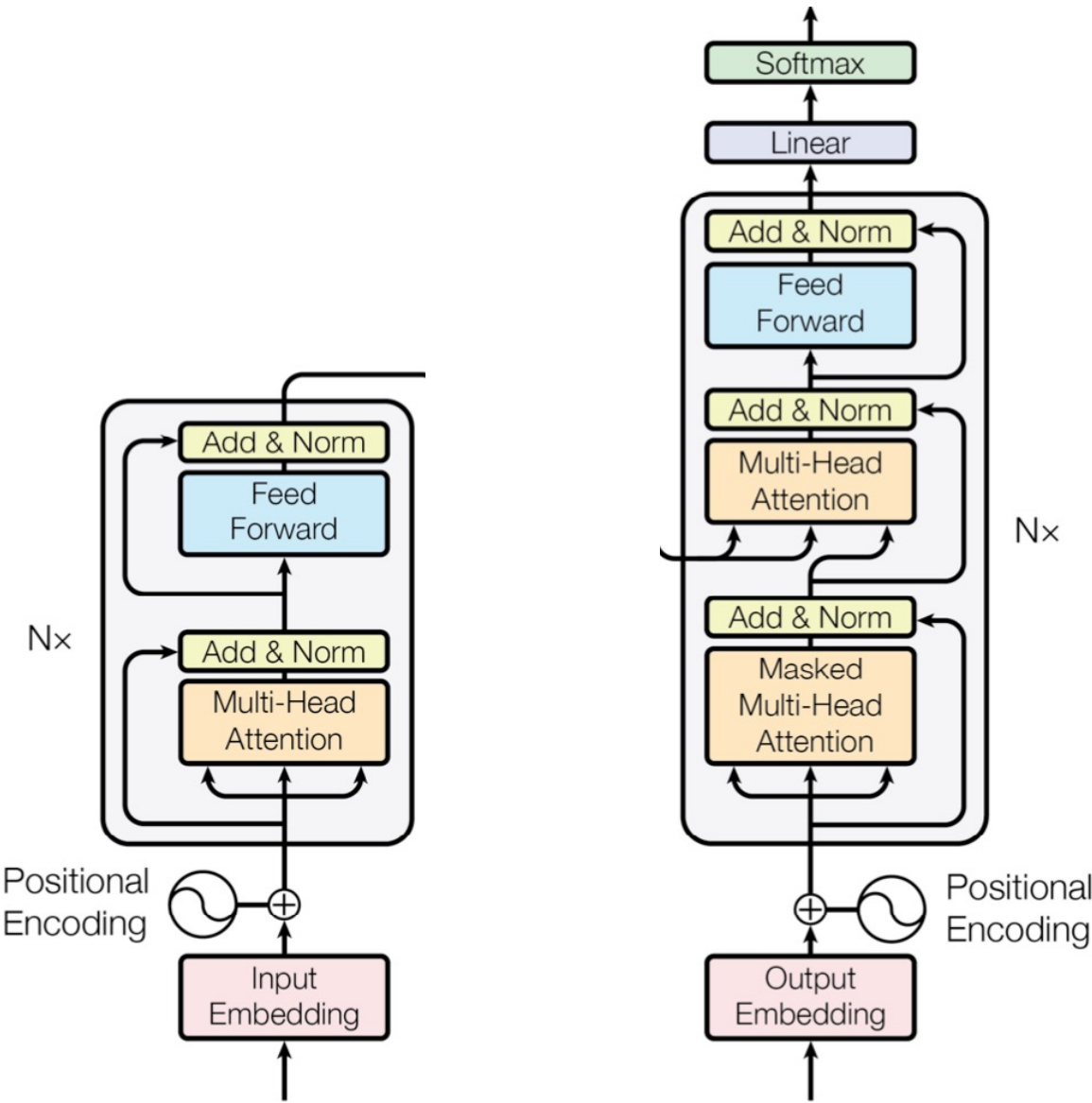
- 1) This is our input sentence*
- 2) We embed each word*
- 3) Split into 8 heads. We multiply X or R with weight matrices
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer



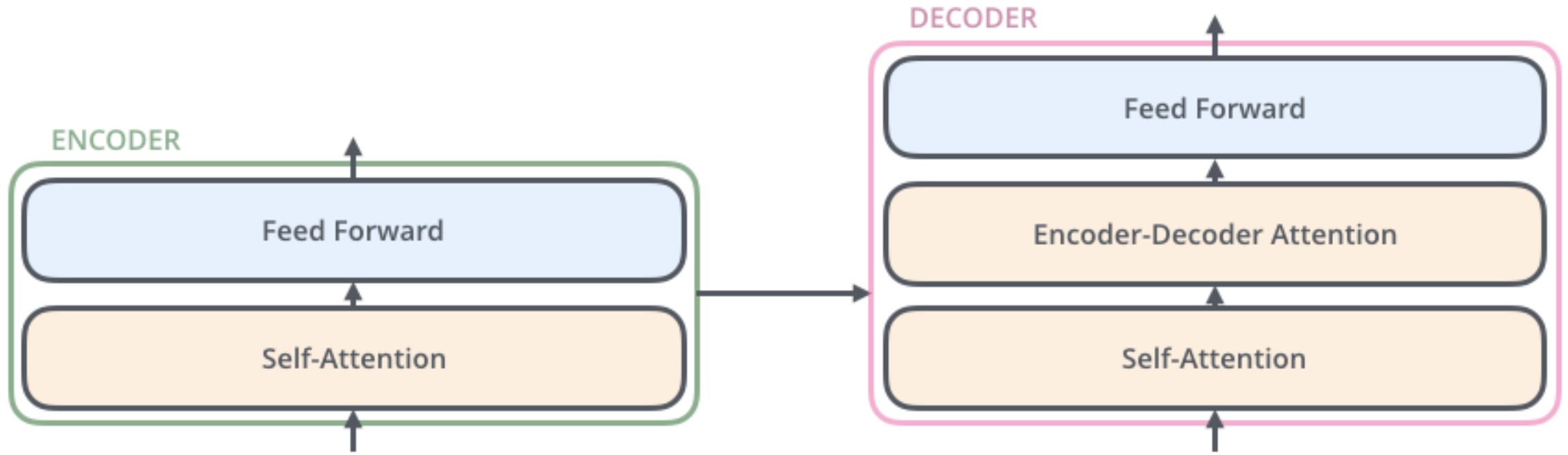
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

Transformer

Decoder



Transformer



Encoder-Decoder Attention layer: An attention layer that helps the decoder focus on relevant parts of the input sentence

Transformer

Decoder

