



Attention and Transformers

November 3, 2021

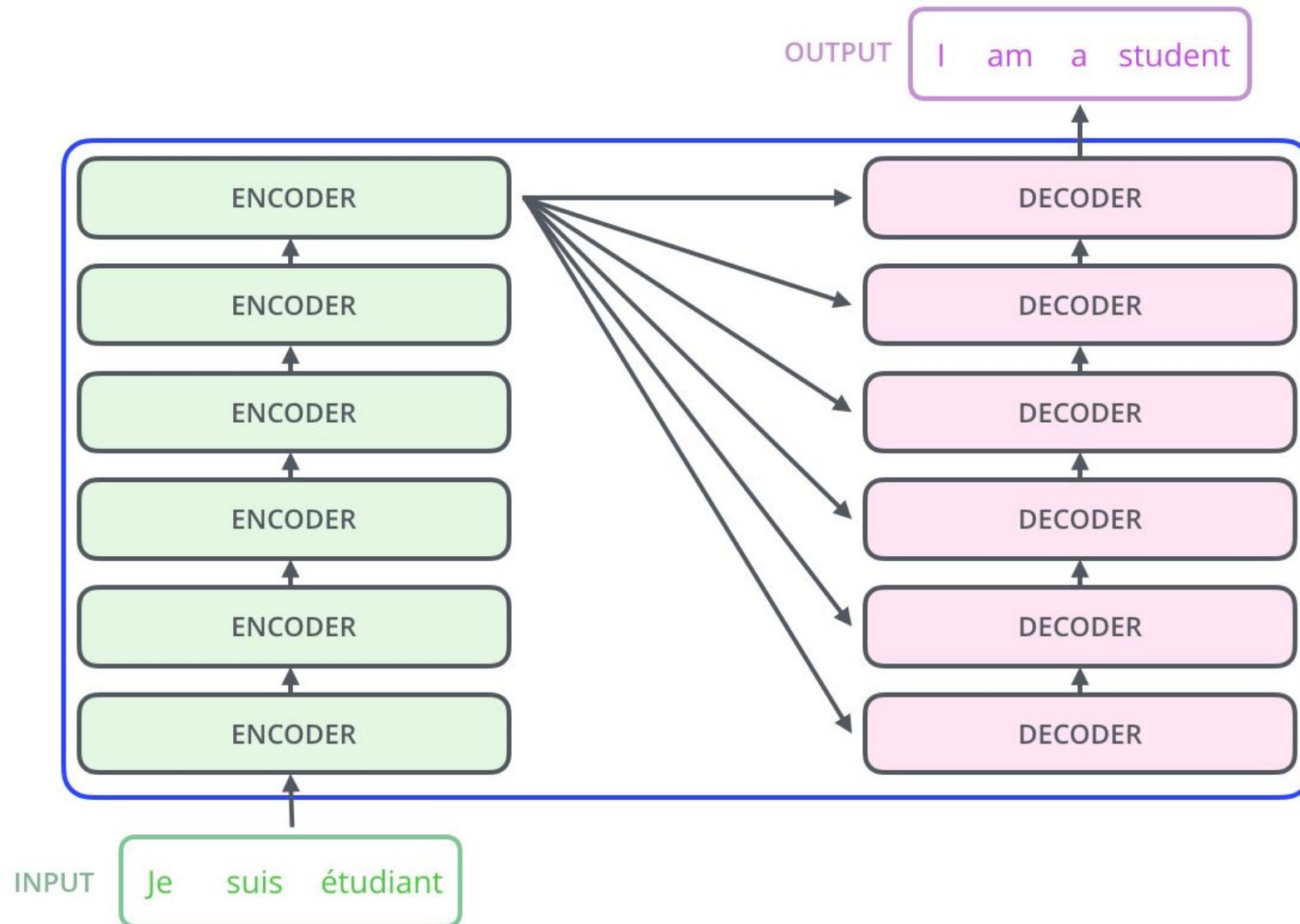


Many ideas for Deep Learning in NLP

1. ULM-FiT, pre-training, transfer learning in NLP
2. Recurrent models require linear sequential computation, hard to parallelize. ELMo, bidirectional LSTM.
3. In order to reduce such sequential computation, several models based on CNN are introduced, such as ConvS2S and ByteNet. Dependency for ConvS2S needs linear depth, and ByteNet logarithmic.
4. The transformer is the first transduction model relying entirely on self-attention to compute the representations of its input and output without using RNN or CNN.

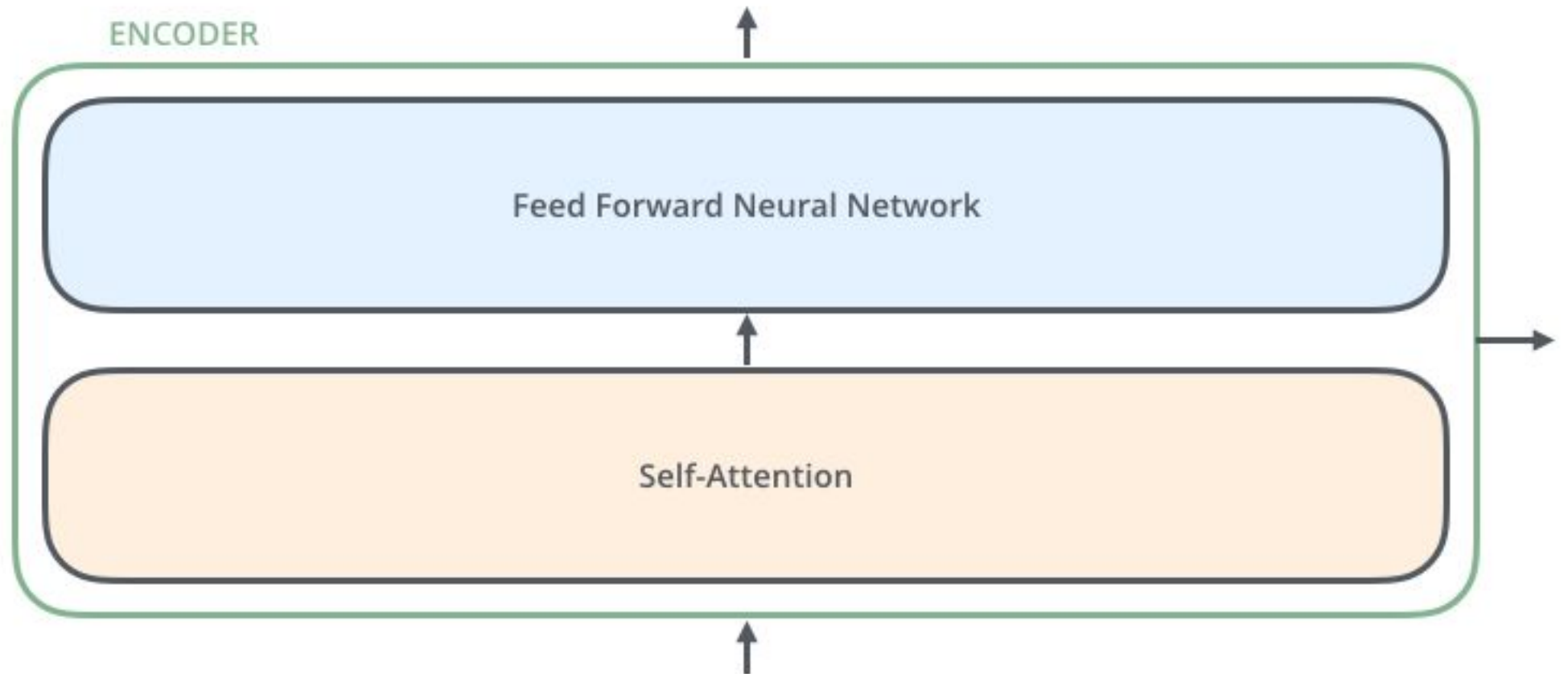


Transformer





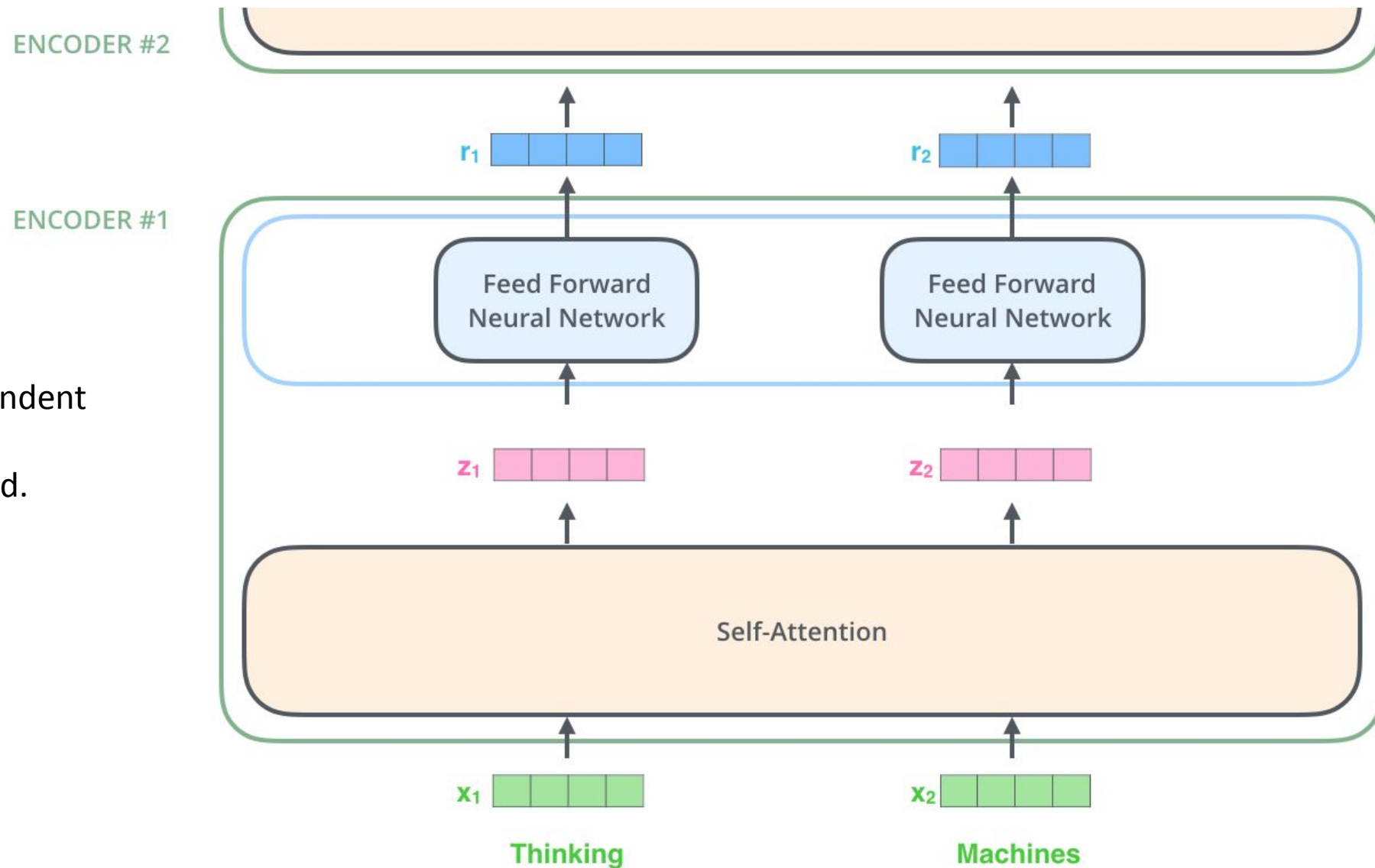
An Encoder Block: same structure, different parameters





Encoder

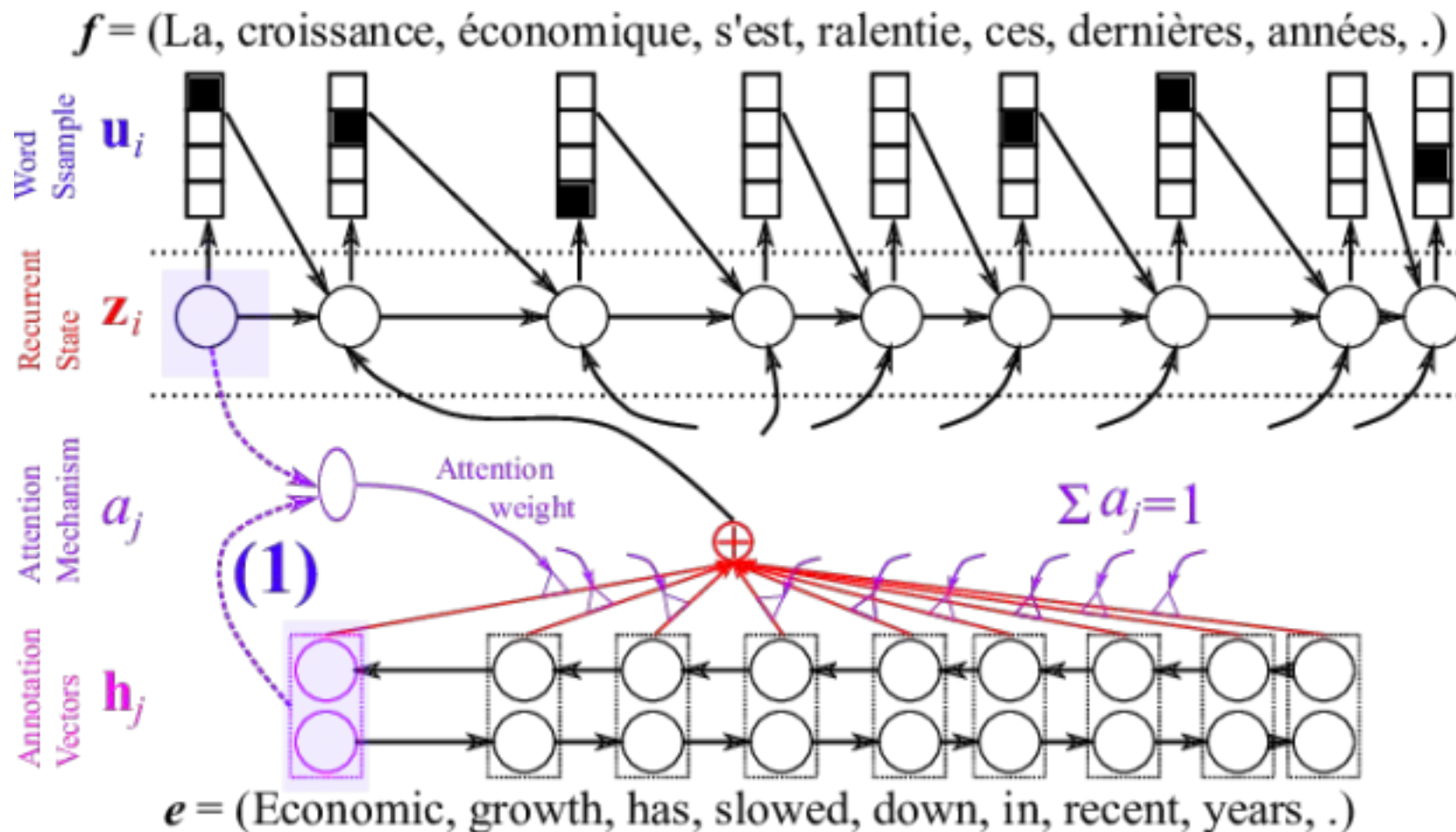
Note: The ffnn is independent for each word.
Hence can be parallelized.





Attention

Given a set of vector values and a vector query, the **attention** is a technique to compute a weighted sum of the values, dependent on the queries.





Self Attention

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

First we create three vectors
by multiplying input embedding
(1x512)

x_i with three matrices (64x512):

$$q_i = x_i W^Q$$

$$K_i = x_i W^K$$

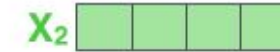
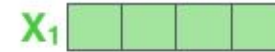
$$V_i = x_i W^V$$

Input

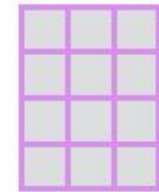
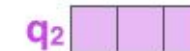
Thinking

Machines

Embedding

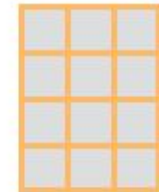
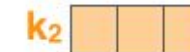
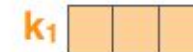


Queries



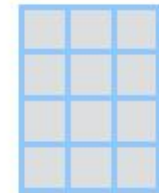
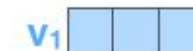
W^Q

Keys



W^K

Values



W^V



Self Attention

Now we need to calculate a score to determine how much focus to place on other parts of the input.

Input

Embedding

Queries

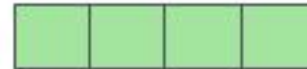
Keys

Values

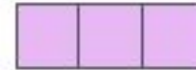
Score

Thinking

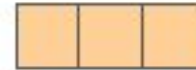
x_1



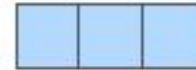
q_1



k_1



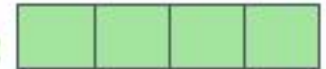
v_1



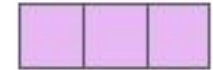
$$q_1 \cdot k_1 = 112$$

Machines

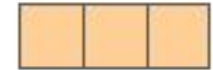
x_2



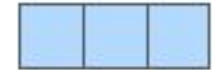
q_2



k_2



v_2



$$q_1 \cdot k_2 = 96$$



Self Attention

Formula

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

$d_k=64$ is dimension of key vector

Input

Embedding

Queries

Keys

Values

Score

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

Softmax

Softmax

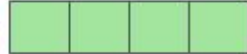
X

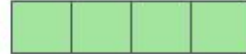
Value

Sum

Thinking

Machines

x_1 

x_2 

q_1 

q_2 

k_1 

k_2 

v_1 

v_2 

$q_1 \cdot k_1 = 112$

$q_1 \cdot k_2 = 96$

14

12

0.88

0.12

\tilde{v}_1 

\tilde{v}_2 

$z_1 = 0.88v_1 + 0.12v_2$

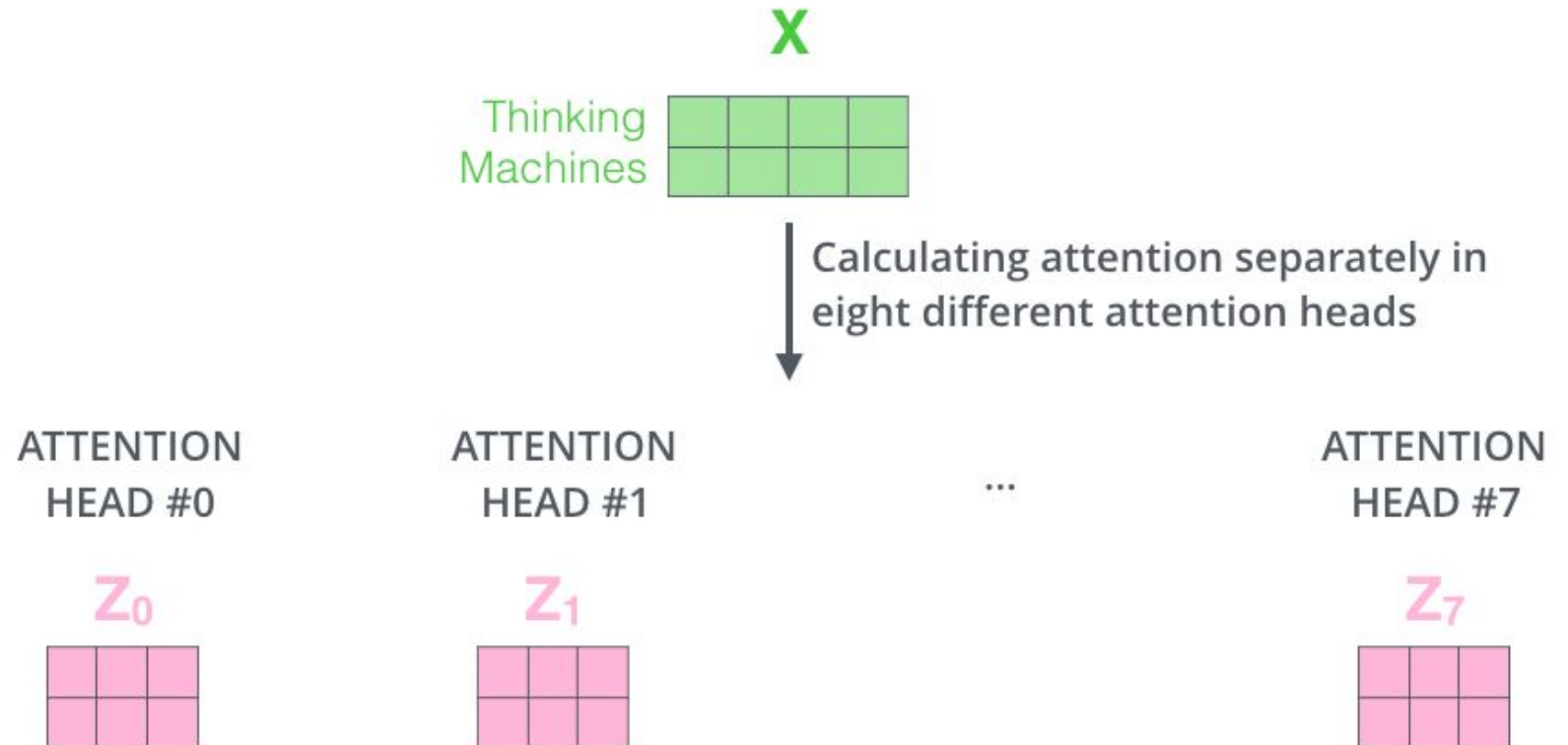
z_1 

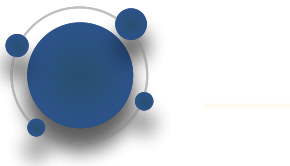
z_2 



Multiple heads

1. It expands the model's ability to focus on different positions.
2. It gives the attention layer multiple "representation subspaces"





The output
is
expecting
only a 2x4
(|input|x6
4) matrix,
hence,

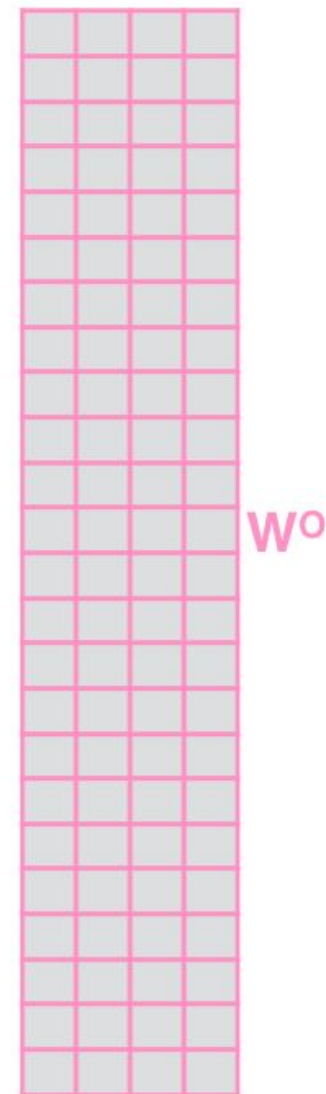
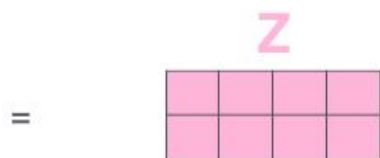
1) Concatenate all the attention heads



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

x

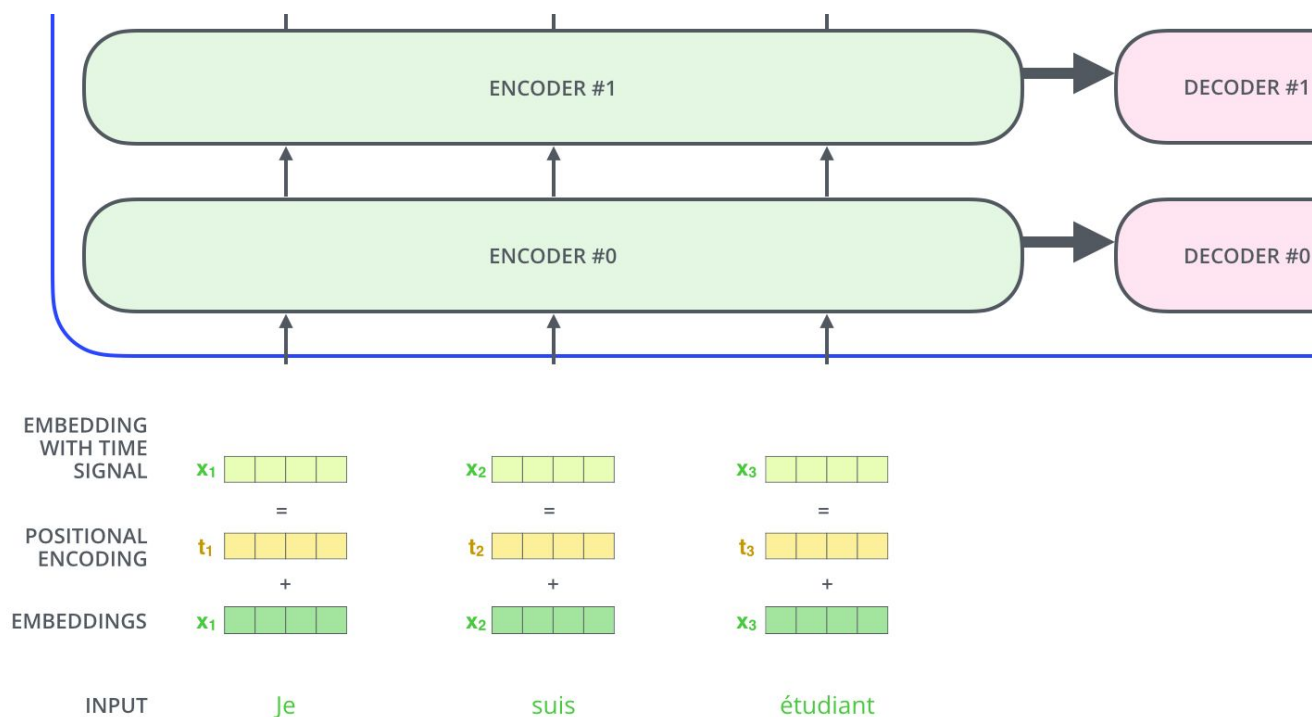
3) The result would be the Z matrix that captures information
from all the attention heads. We can send this forward to the FFNN





Representing the input order (positional encoding)

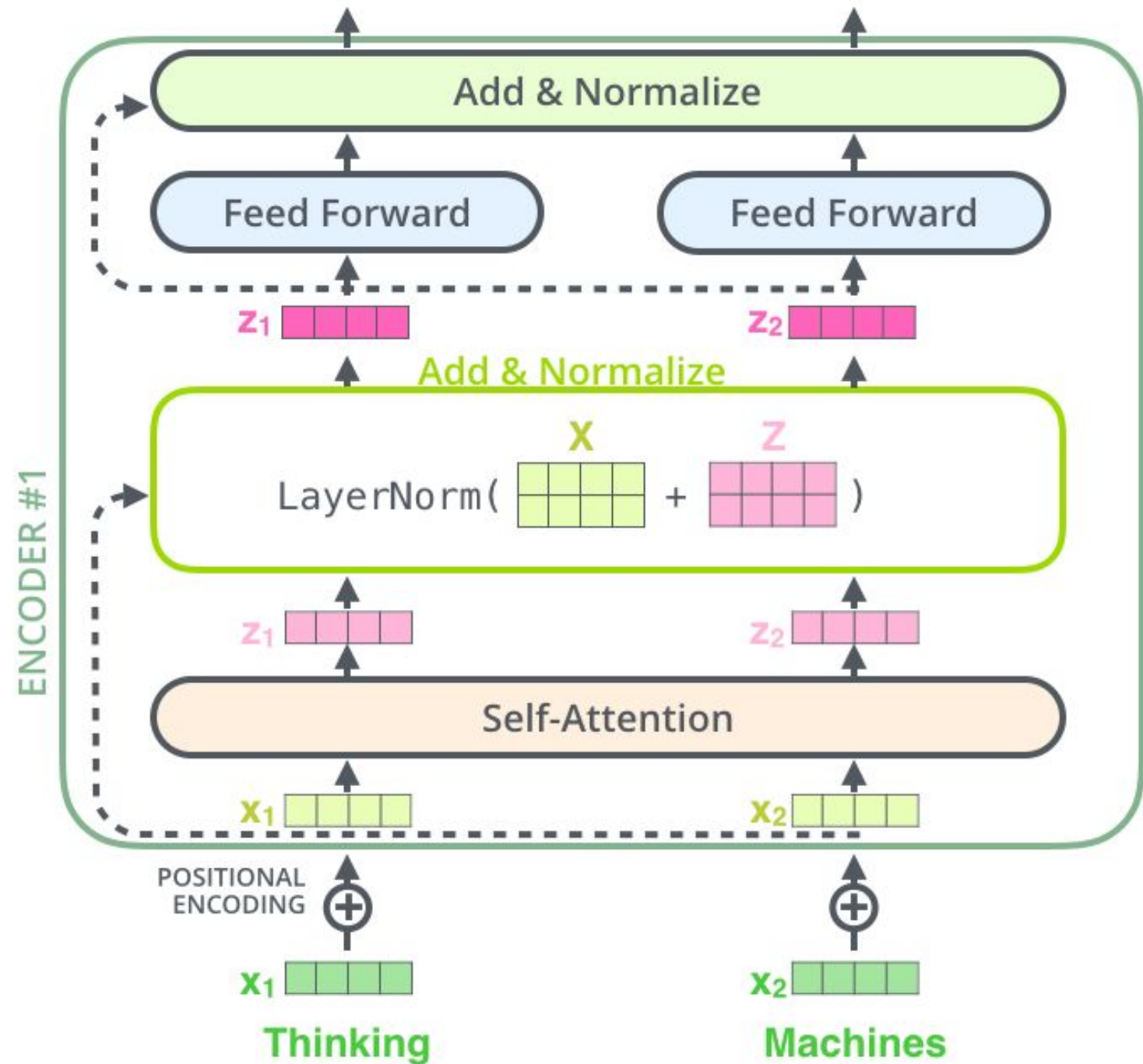
The transformer adds a vector to each input embedding. These vectors follow a specific pattern that the model learns, which helps it determine the position of each word, or the distance between different words in the sequence. The intuition here is that adding these values to the embeddings provides meaningful distances between the embedding vectors once they're projected into Q/K/V vectors and during dot-product attention.





Add and Normalize

In order to regulate the computation, this is a normalization layer so that each feature (column) have the same average and deviation.

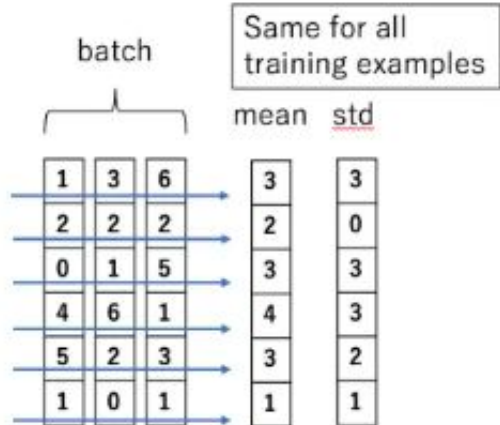




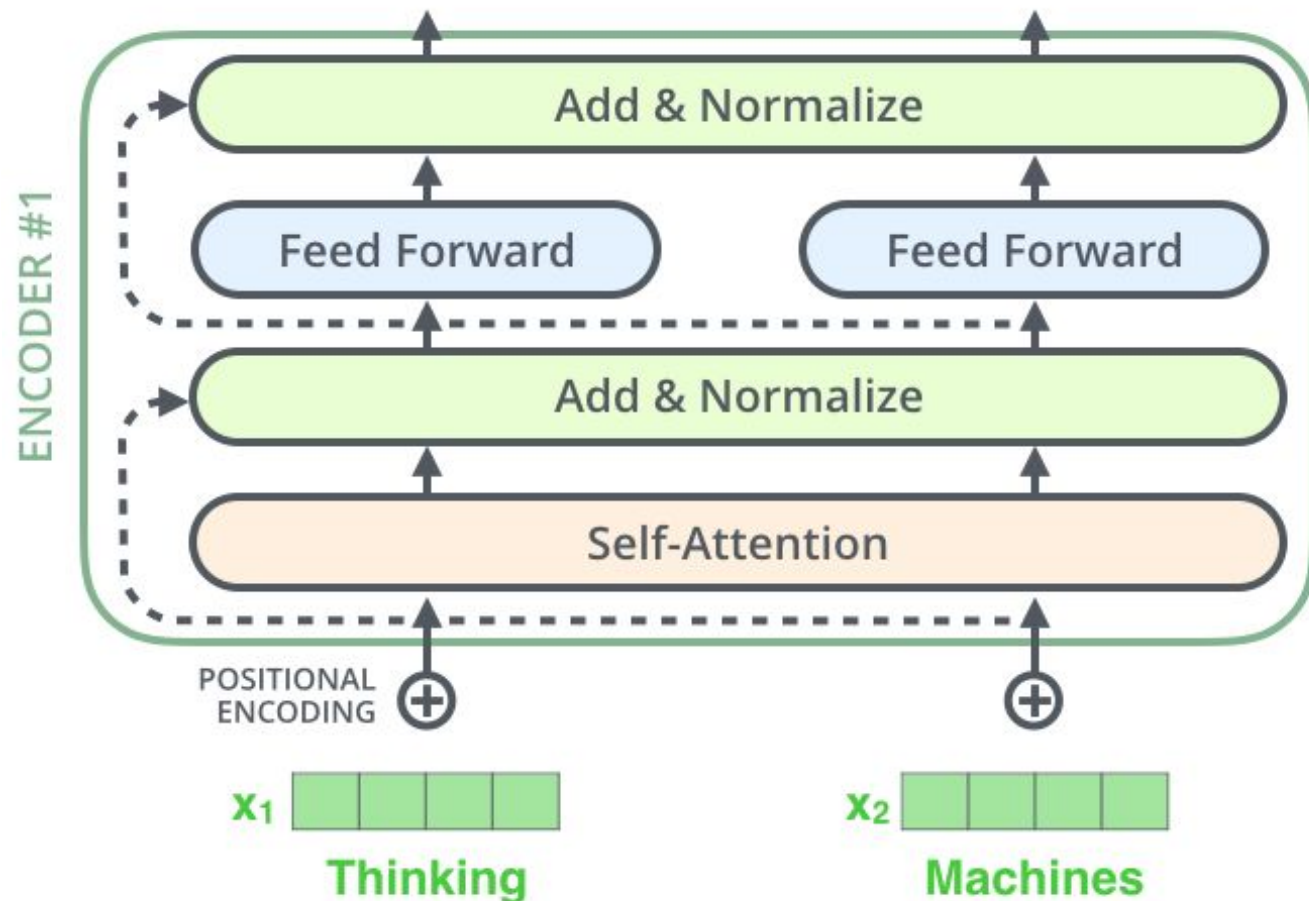
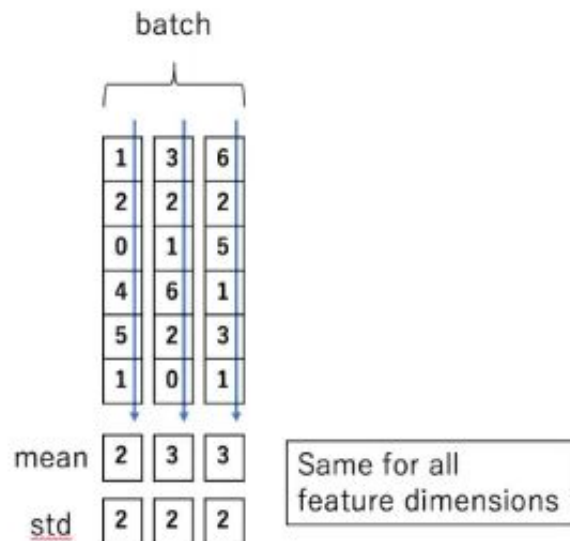
Layer Normalization (Hinton)

Layer normalization normalizes the inputs across the features.

Batch Normalization



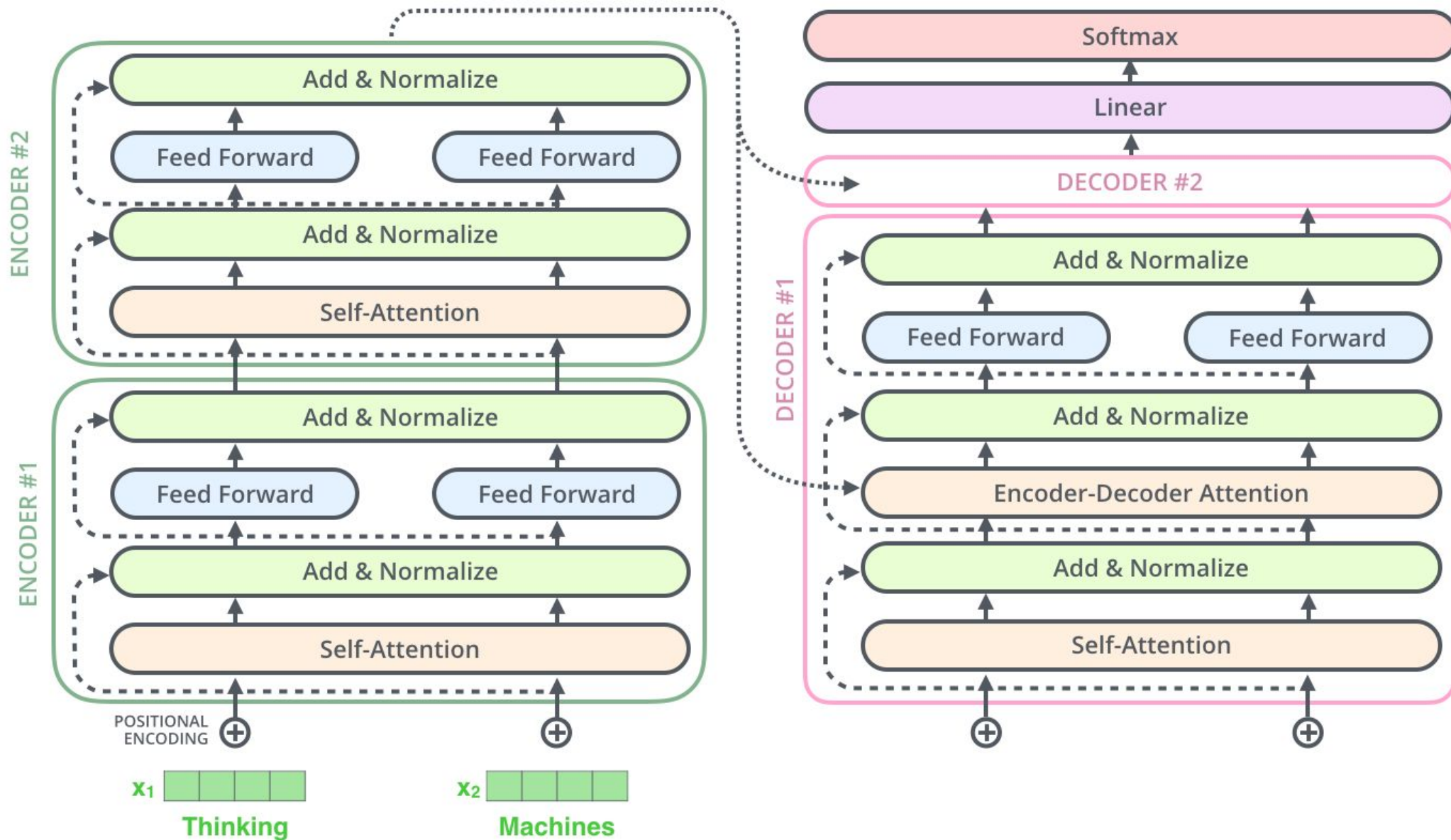
Layer Normalization

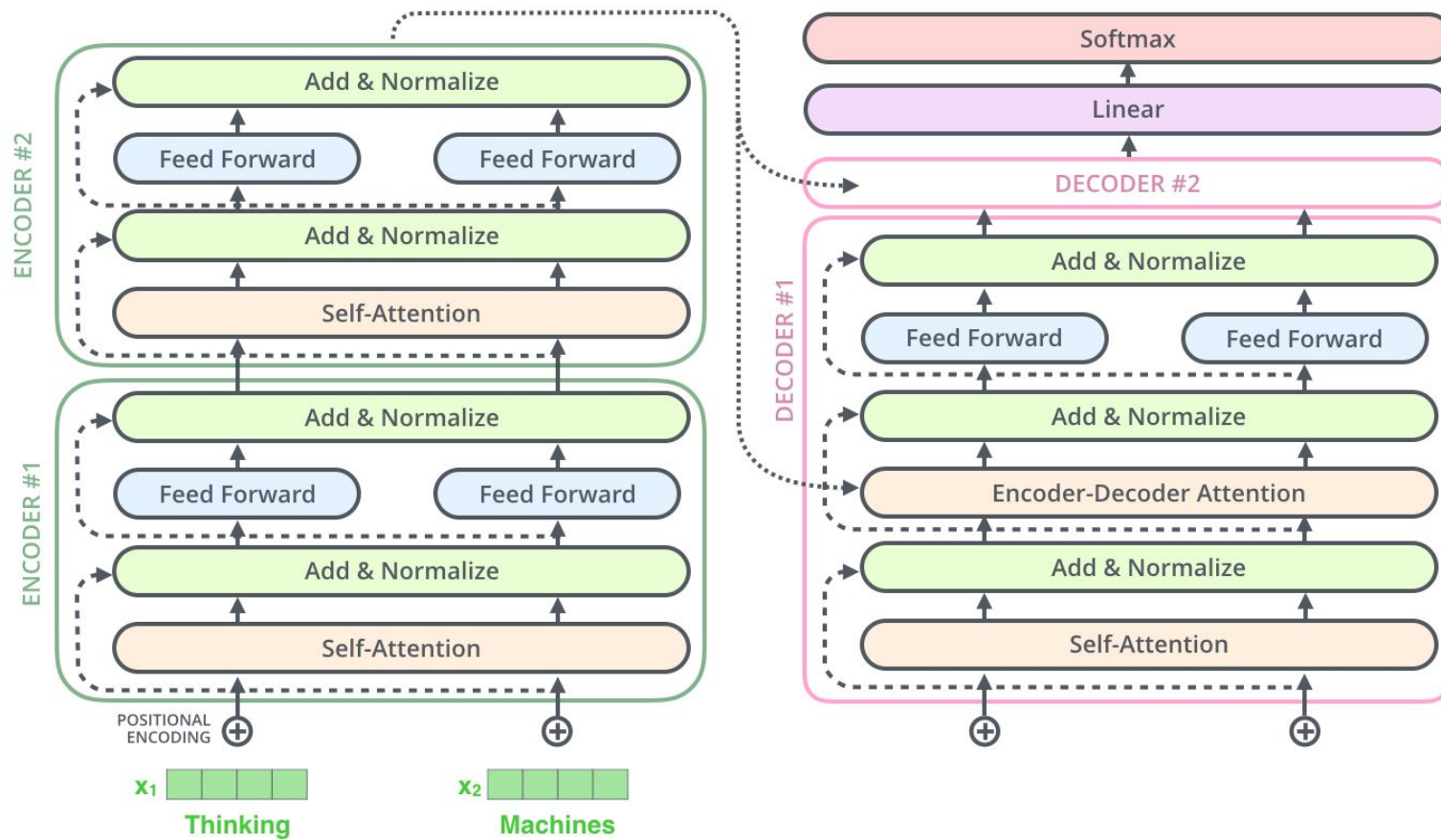




The complete transformer

The encoder-decoder attention is just like self attention, except it uses K, V from the top of encoder output, and its own Q





BERT

<https://www.kaggle.com/residentmario/notes-on-gpt-2-and-bert-models>

GPT





Literature & Resources for Transformers

Vaswani et al. “Attention is all you need.” 2017.

Resources:

<https://nlp.seas.harvard.edu/2018/04/03/attention.html> (Excellent explanation of transformer model with codes.)

Jay Alamar, The illustrated transformer (where many pictures in these slides came from):

<http://jalamar.github.io/illustrated-transformer/>

Kate Lognina: Attention in NLP, summarizes all sorts of attentions.

<https://medium.com/@joealato/attention-in-nlp-734c6fa9d983>