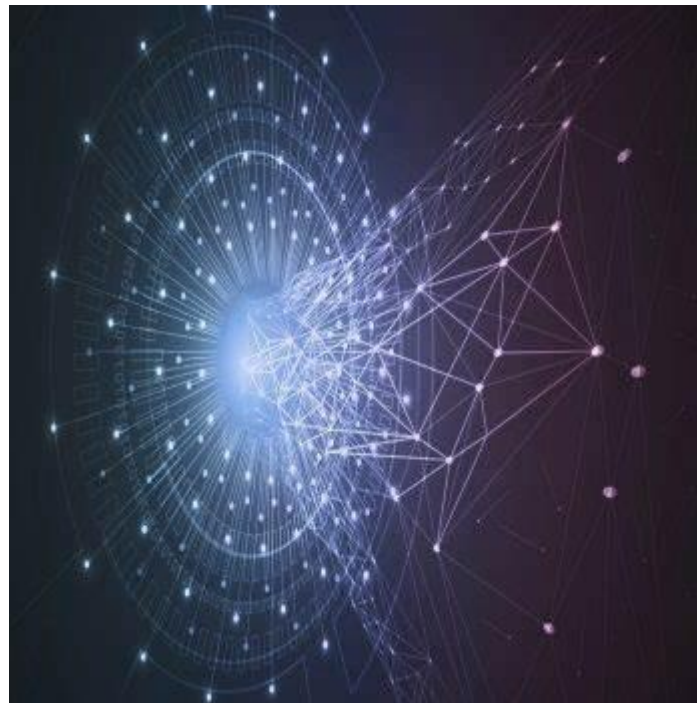


# Transformers



Prof Dr Marko Robnik-Šikonja

Natural Language Processing, Edition 2022

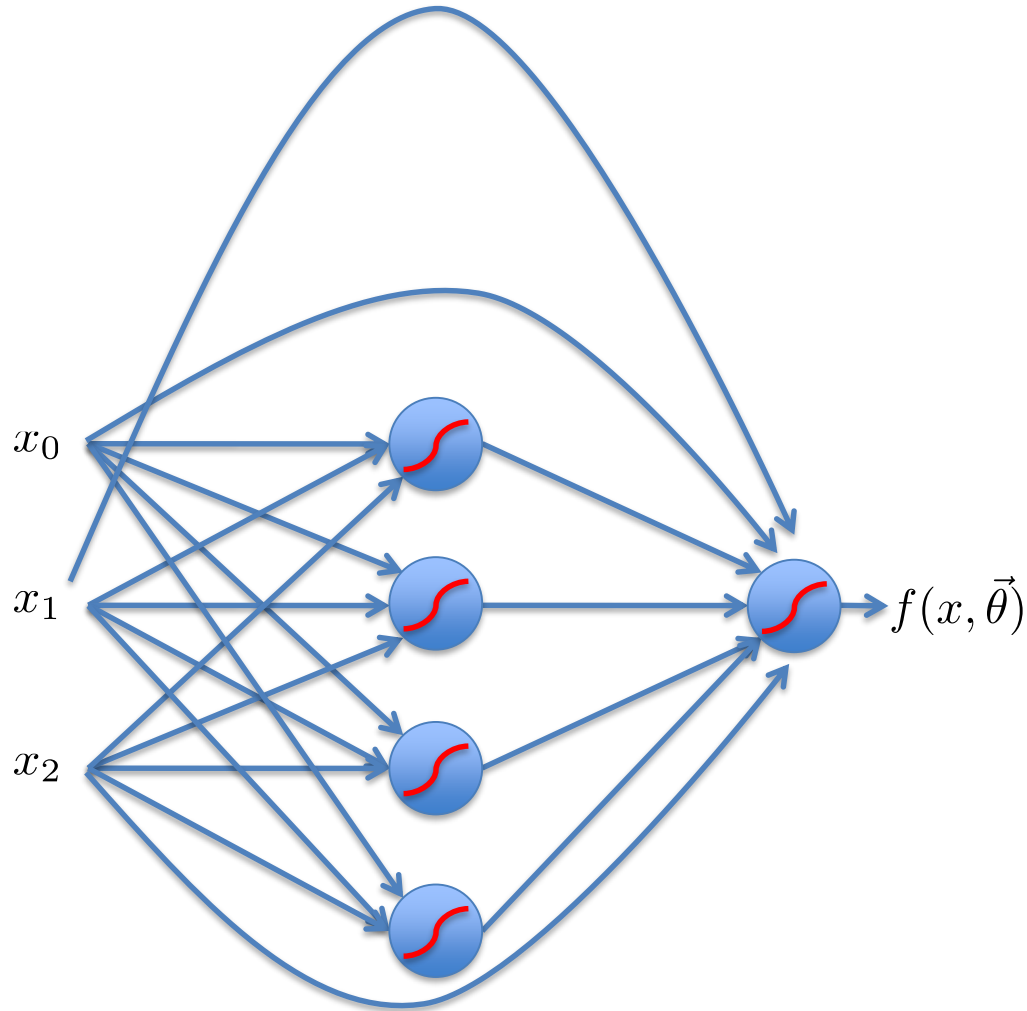
# Contents

- transformer networks
- BERT
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł. and Polosukhin, I., 2017. [Attention is all you need](#). In *Advances in neural information processing systems* (pp. 5998-6008).
- Jay Alammam: [The Illustrated Transformer](#). Blog, 2019.
- Sasha Rush: [Annotated Transformer](#). Jupyter Notebook using PyTorch.
- Elvis Saravia: <https://github.com/dair-ai/Transformers-Recipe>
- Devlin, J., Chang, M.W., Lee, K. and Toutanova, K., 2019. [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Volume 1, pp. 4171-4186.
- some slides by Jay Alammam, Jacob Devlin and Andrej Miščič

# Problems with RNNs

- We want parallelization but RNNs are inherently sequential
- For parallelization we need fixed input size (relatively short)
- Despite GRUs and LSTMs, RNNs still need attention mechanism to deal with long range dependencies – path length between states grows with sequence
- If attention gives us access to any state... maybe we can just use attention and don't need the RNN?

# Problems with long sequences: one solution level is jumping

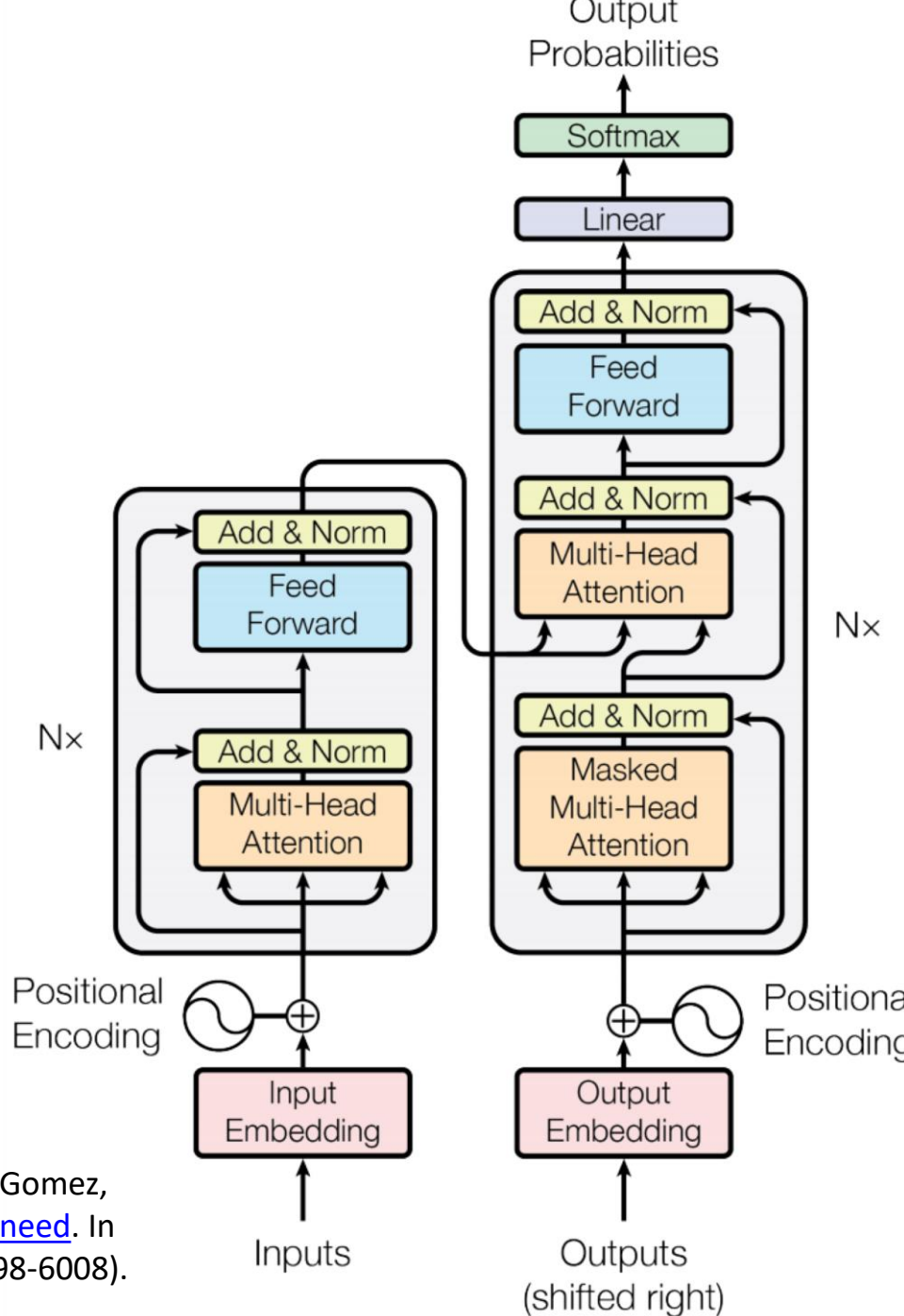


# Transformer model

- currently the most successful DNN
- non-recurrent
- architecturally it is an encoder-decoder model
- fixed input length (relatively long but projected to short)
- adapted for parallelization
- adapted for GPU (TPU) processing
- based on extreme use of attention
- well scalable
- includes level jumping

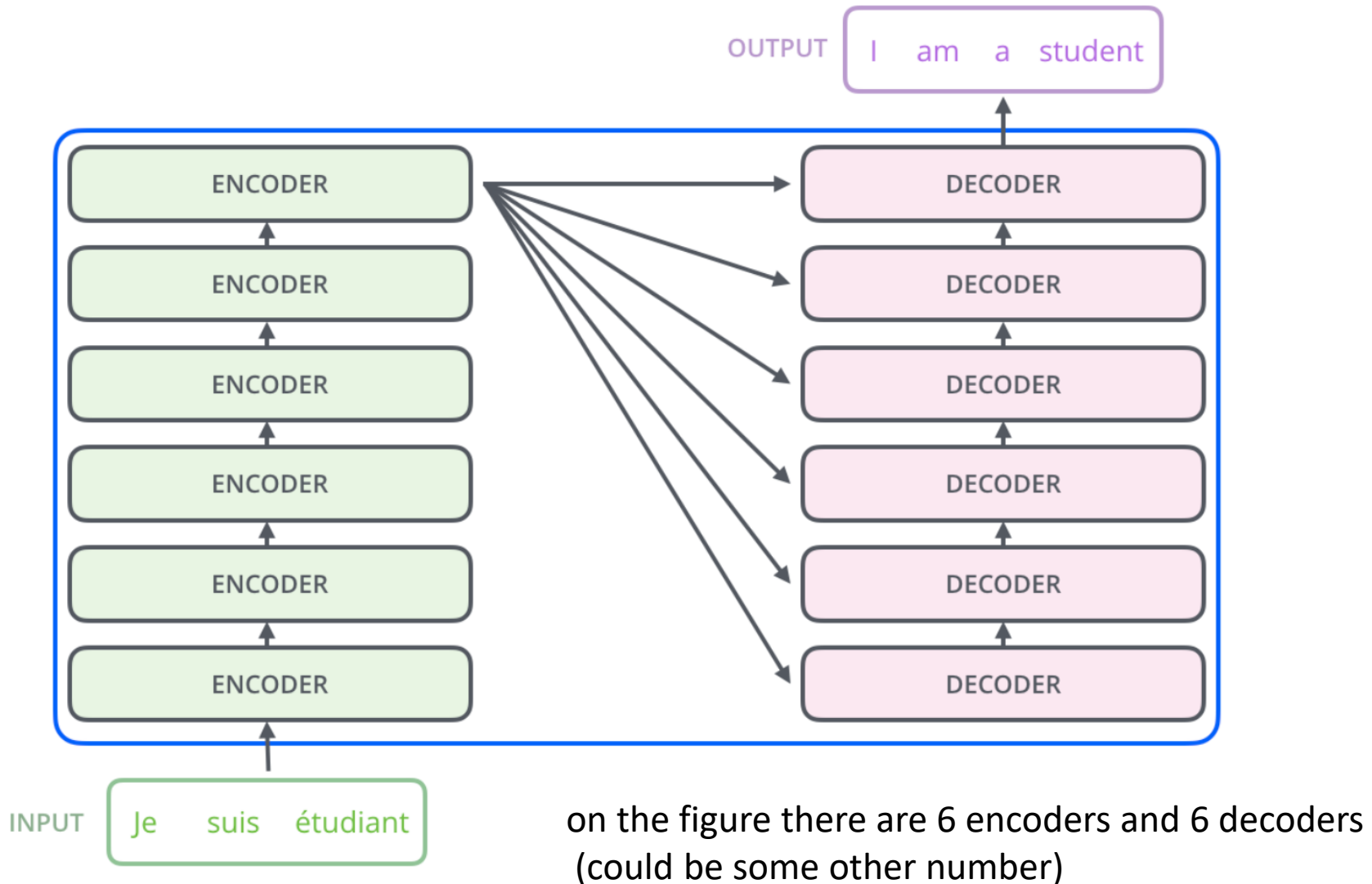
# Transformer overview

- Initial task: machine translation with parallel corpus
- Predict each translated word
- Final cost/loss/error function was standard cross-entropy error on top of a softmax classifier



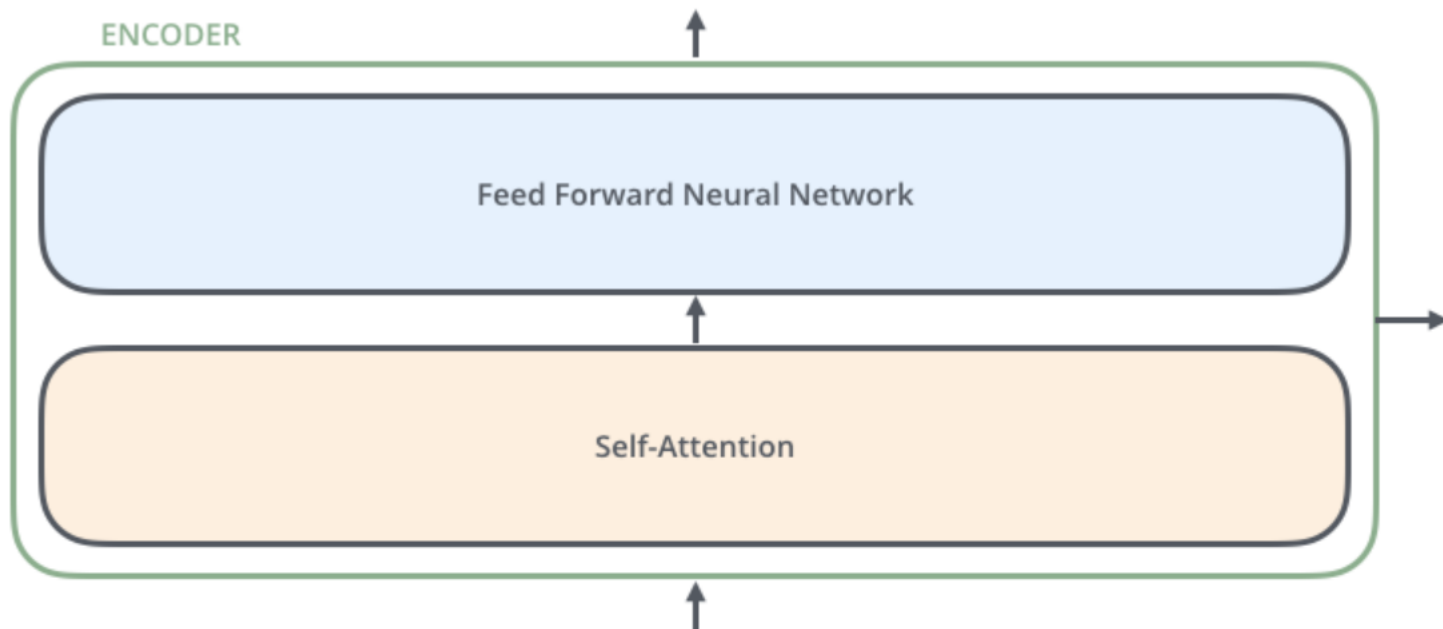
Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł. and Polosukhin, I., 2017. [Attention is all you need](#). In *Advances in neural information processing systems* (pp. 5998-6008).

# Transformer is encoder-decoder model



# Transformer: encoder

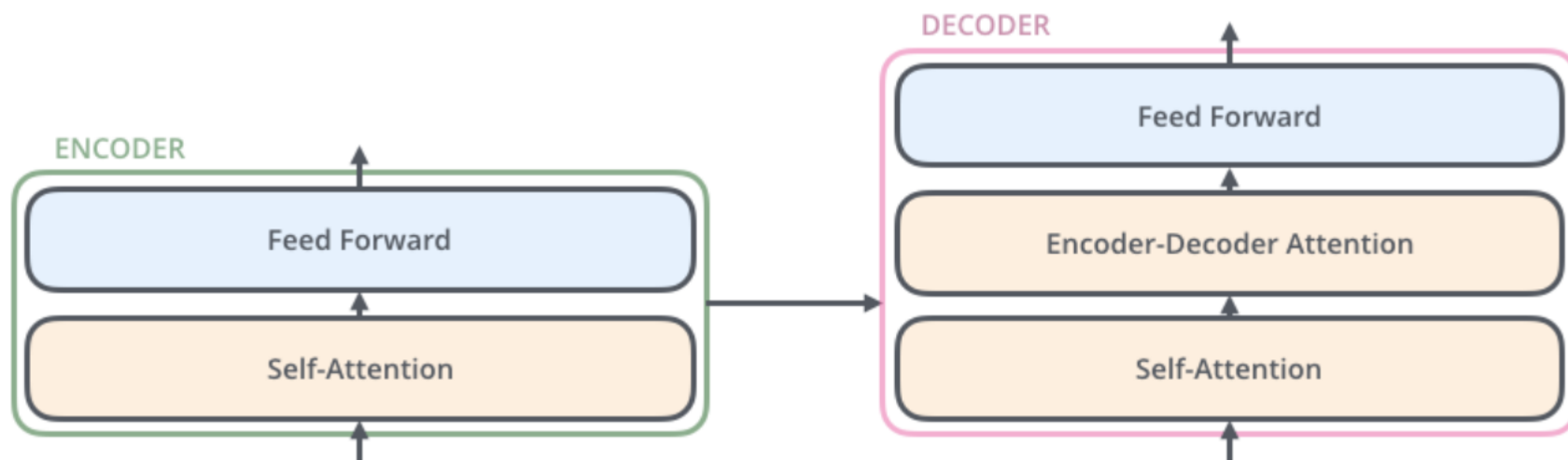
- two layers
- no weight sharing between different encoders
- self-attention helps to focus on relevant part of input



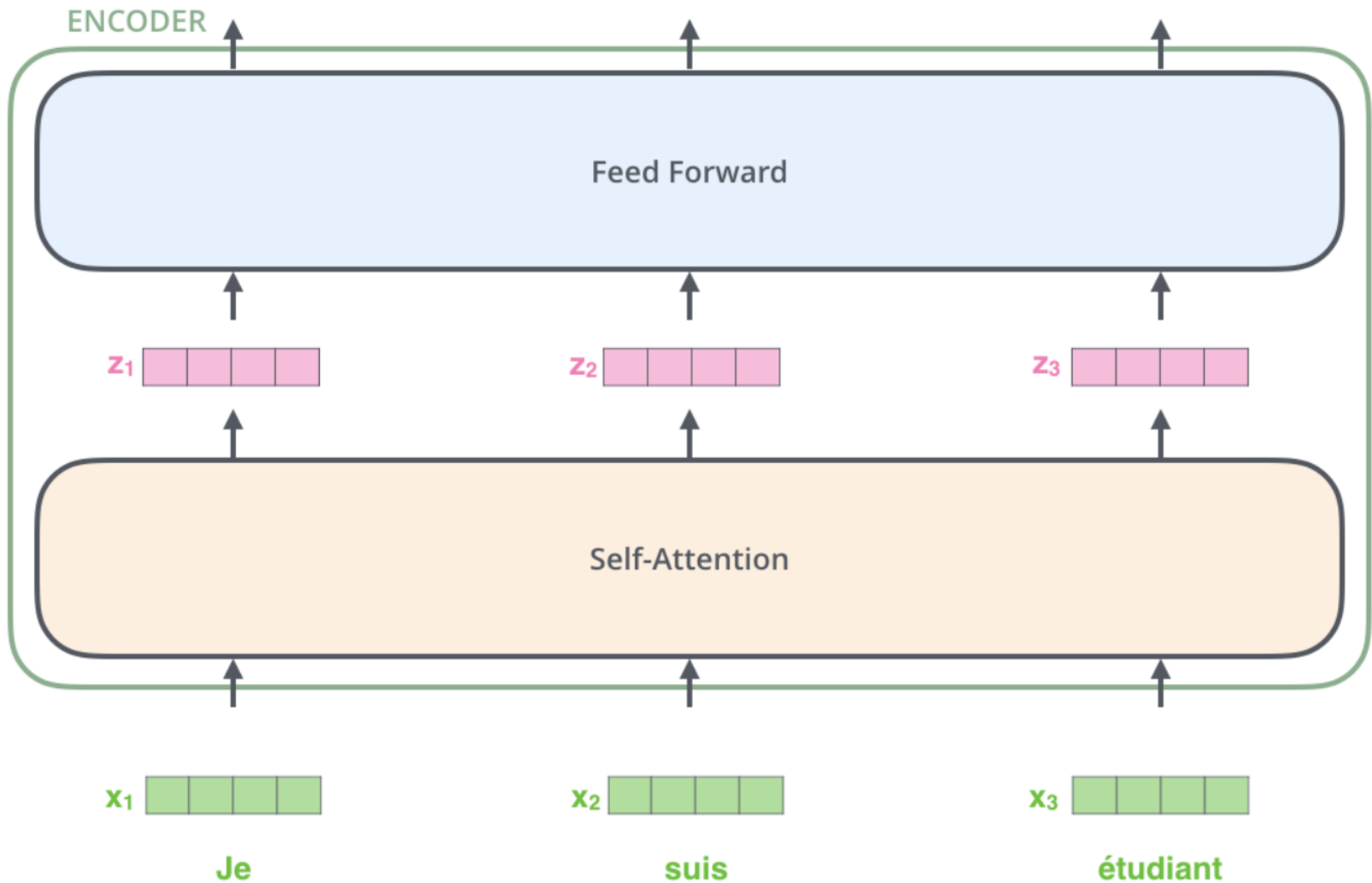


# Transformer: decoder

- the same as encoder but with an additional attention layer in between, receiving input from encoder



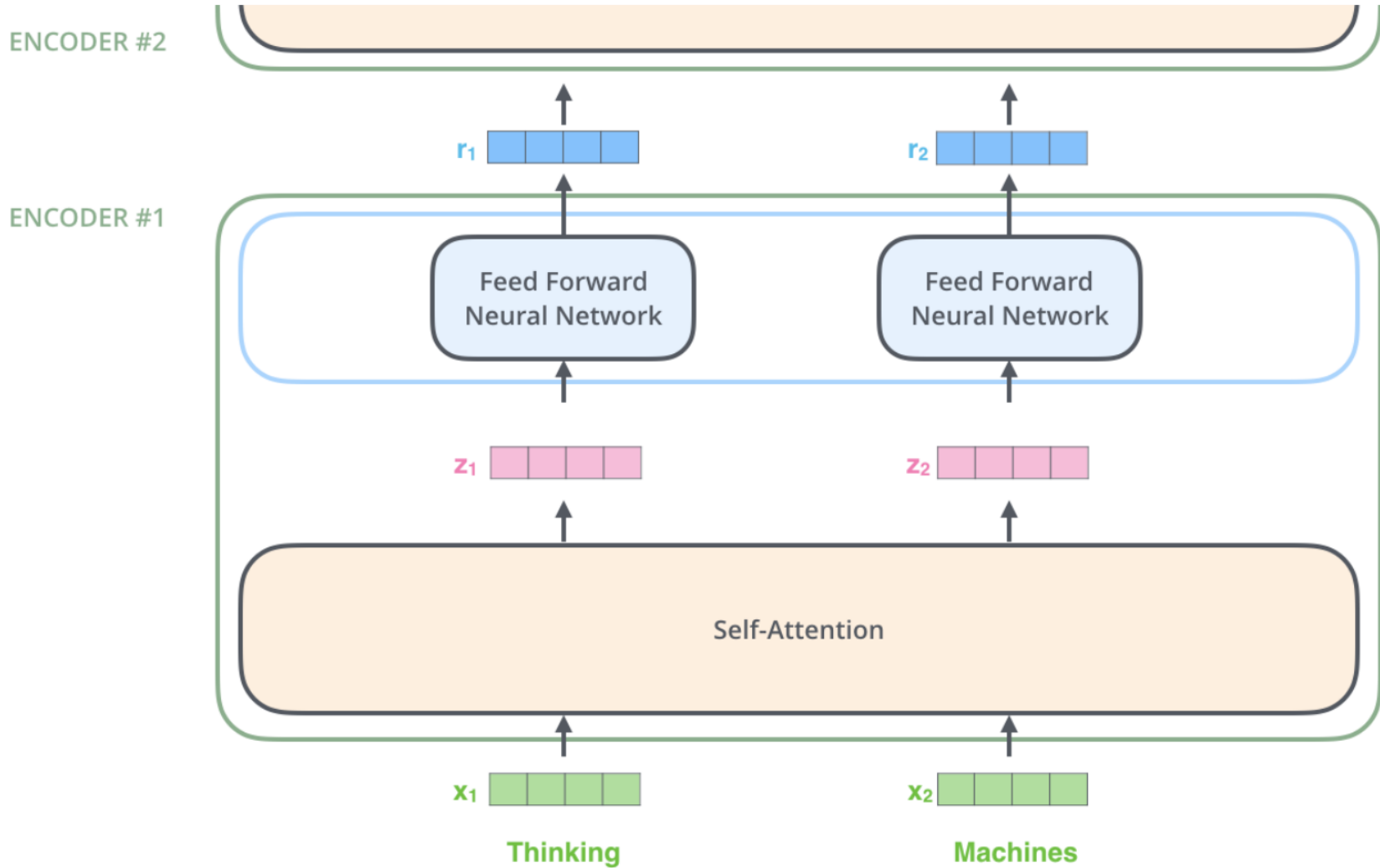
# Start with embeddings



# Input to transformer

- embeddings, e.g., 512 dimensional vectors (special, we will discuss that later)
- fixed length, e.g., max 128 tokens
- dependencies between inputs are only in the self-attention layer, no dependencies in feed forward layer – good for parallelization
- **Let us first present the working of the transformer with illustration of the prediction, later we will cover also training.**

# Encoding



# Self-attention

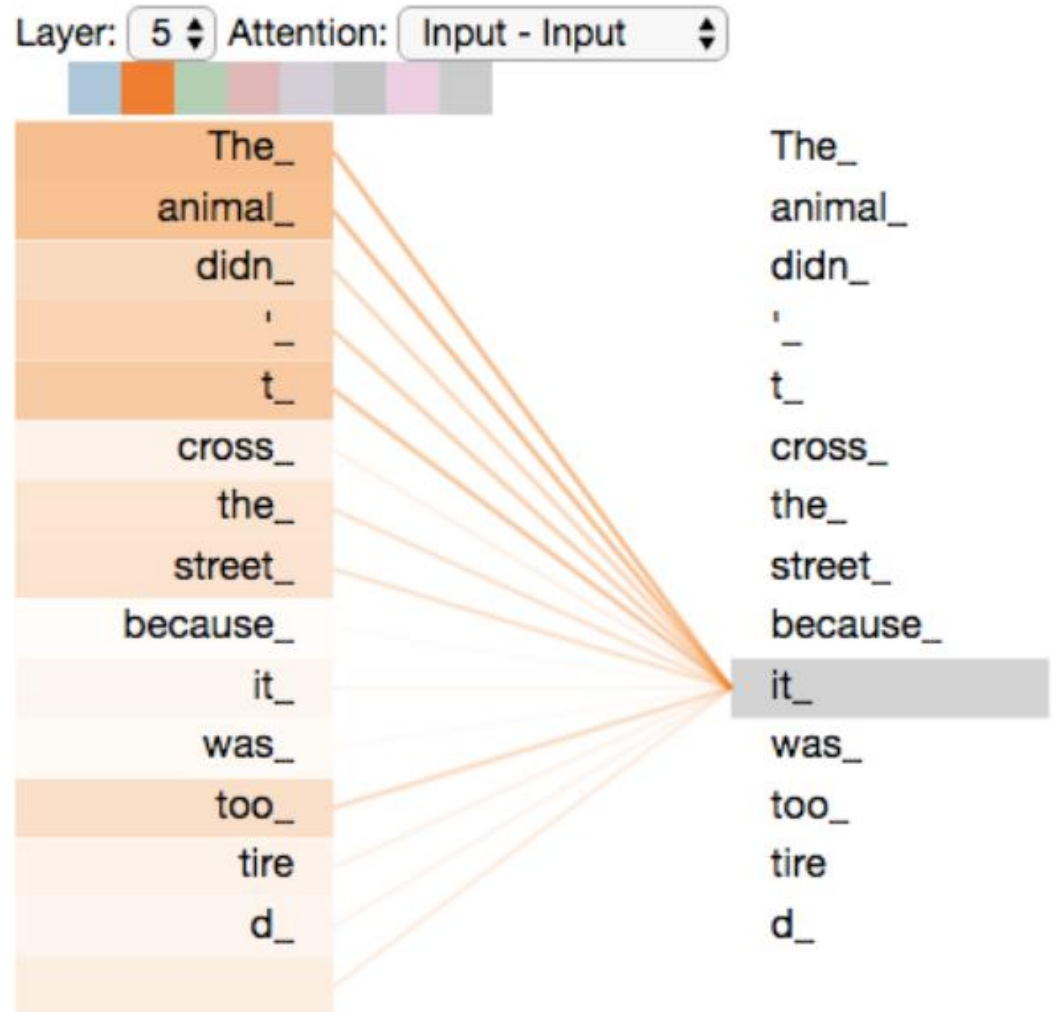
- As the model processes each word (each position in the input sequence), self-attention allows it to look at other positions in the input sequence for clues that can help lead to a better encoding for this word

*“The animal didn't cross the street because it was too tired”*

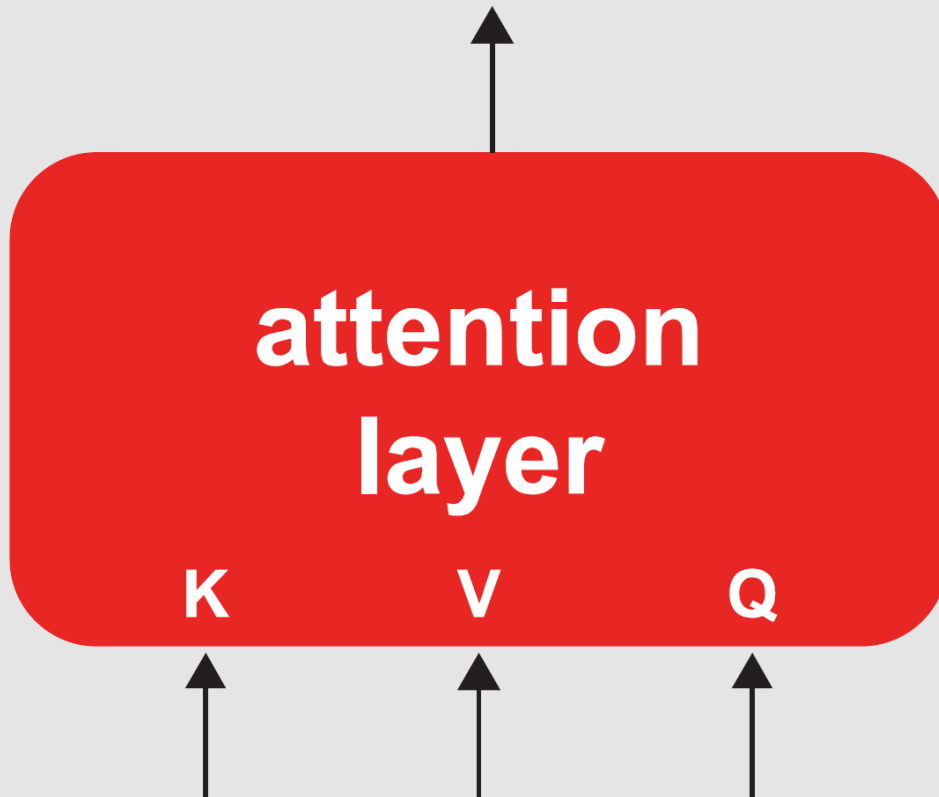
- What does “it” in this sentence refer to? Is it referring to the street or to the animal? It’s a simple question to a human, but not as simple to an algorithm.
- *“The animal didn't cross the street because it was too wide”*
- When the model is processing the word “it”, self-attention allows it to associate “it” with “animal”.

# Illustrating self-attention

- As we are encoding the word "it" in encoder #5 (the top encoder in the stack), part of the attention mechanism was focusing on "The Animal", and baked a part of its representation into the encoding of "it".



$$Y_i = \sum_{j=1}^n \alpha_{ij} V_j$$



$$sim_{ij} = \phi(Q_i, K_j)$$

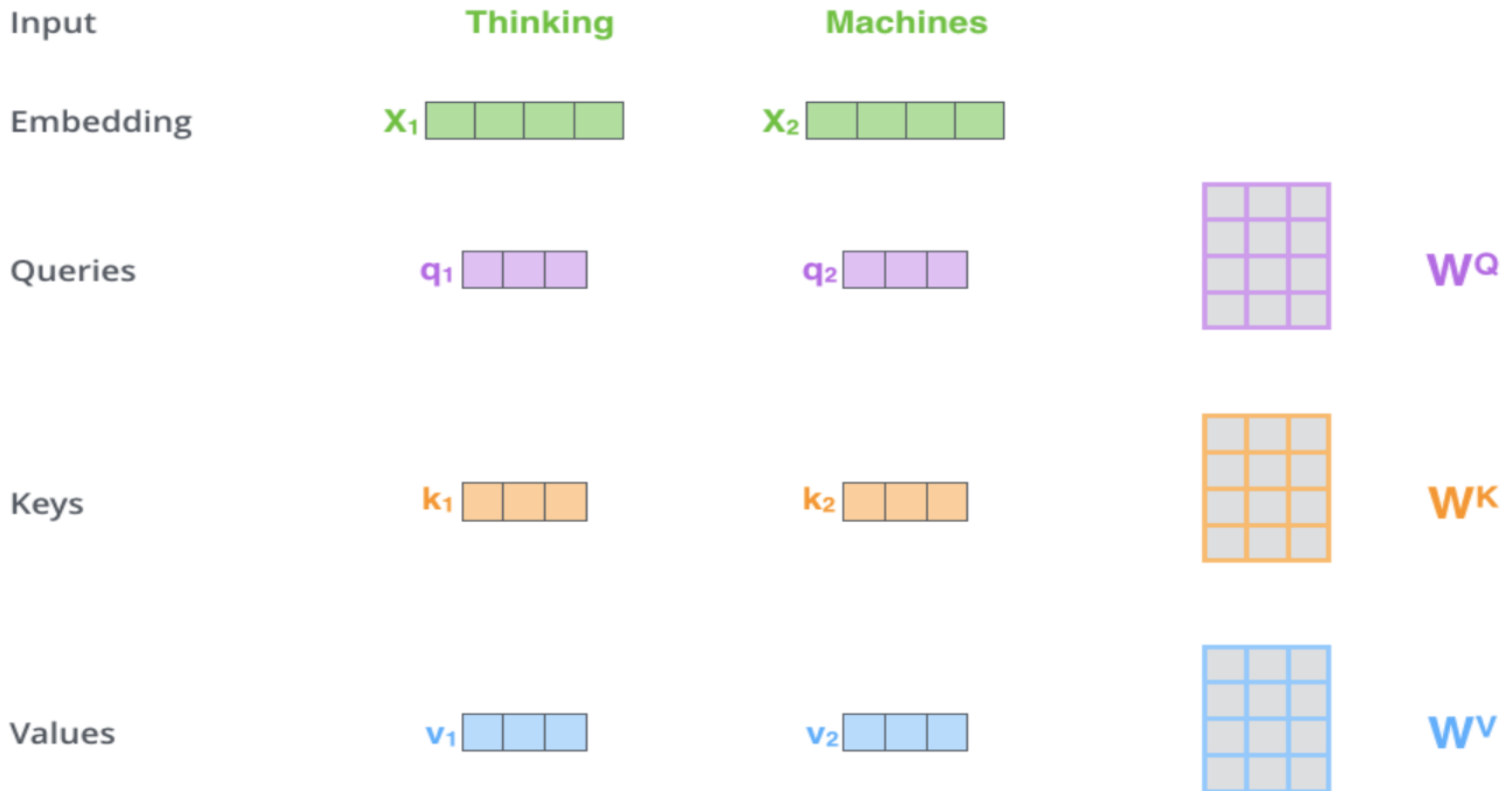
$$\alpha_{ij} = \frac{\exp(sim_{ij})}{\sum_{k=1}^n \exp(sim_{ik})}$$

# Self-attention details 1/4

- create three vectors from each of the encoder's input vectors (in the first layer these are the embedding of each word).
- Query vector  $Q$ , Key vector  $K$ , Value vector  $V$  are created by multiplying the embedding by three matrices that are trained during the training process.
- $Q$ ,  $K$ , and  $V$  are smaller than the embedding vector, typically 64, while the embedding and encoder input/output vectors have dimensionality of 512.
- They are smaller to make the computation of multiheaded attention (mostly) constant.



# Self-attention details 1/4

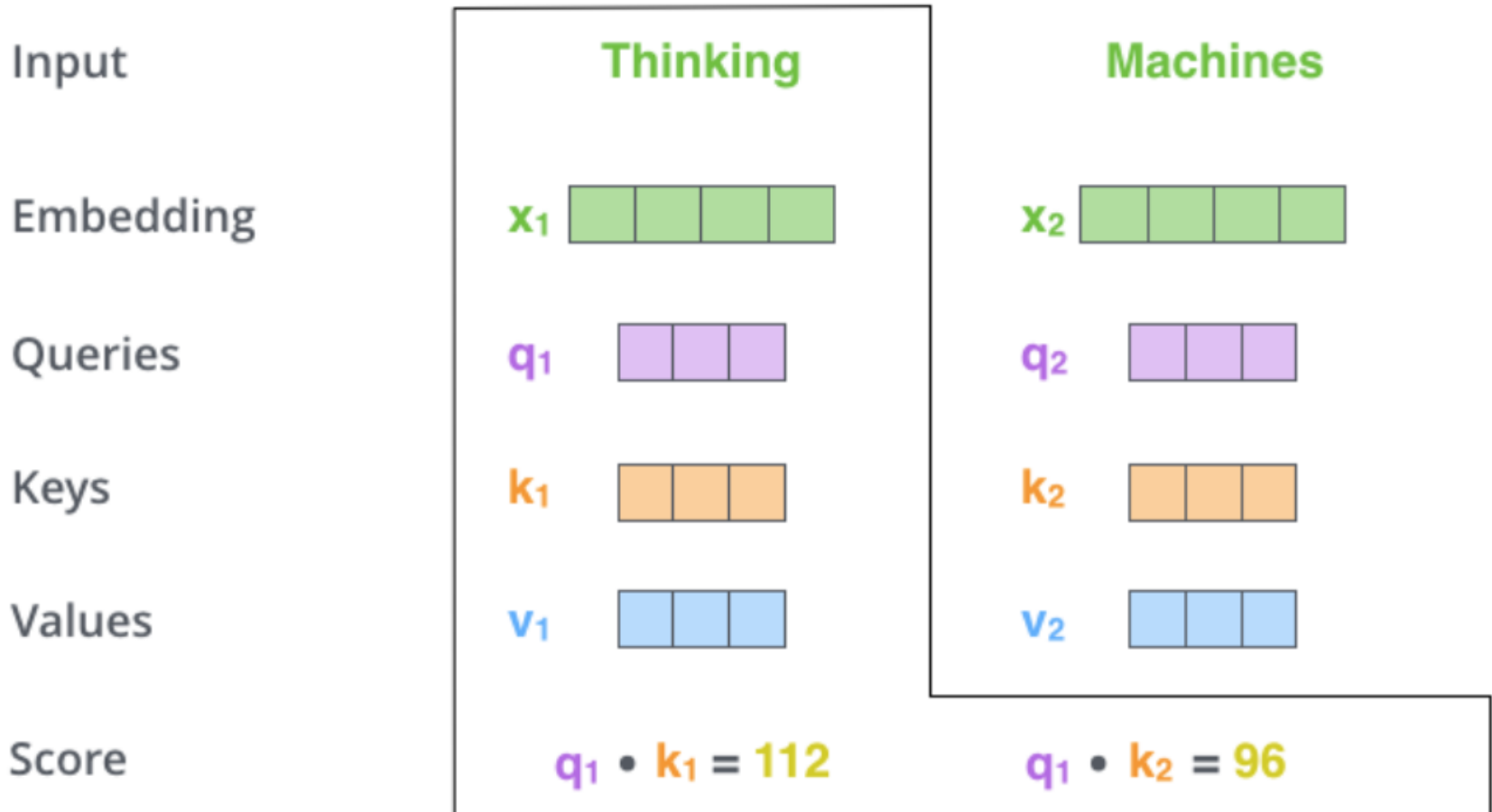


Multiplying  $x_1$  by the  $W^Q$  weight matrix produces  $q_1$ , the "query" vector associated with that word. We end up creating a "query", a "key", and a "value" projection of each word in the input sentence.

# Details 2/4: attention vectors Q, K, V

- The “query” Q, “key” K, and “value” V vectors are abstractions that are useful for calculating and thinking about attention.
- To calculate self-attention for a given word (e.g., “Thinking”), we score each word of the input sentence against this word. The score determines how much focus to place on other parts of the input sentence as we encode a word at a certain position.
- The score is calculated by taking the dot product of the query vector Q with the key vector K of the respective word.
- E.g., on computing the self-attention for the word in position #1, we would compute dot product of  $q_1$  and  $k_1$ , and  $q_1$  and  $k_2$ .

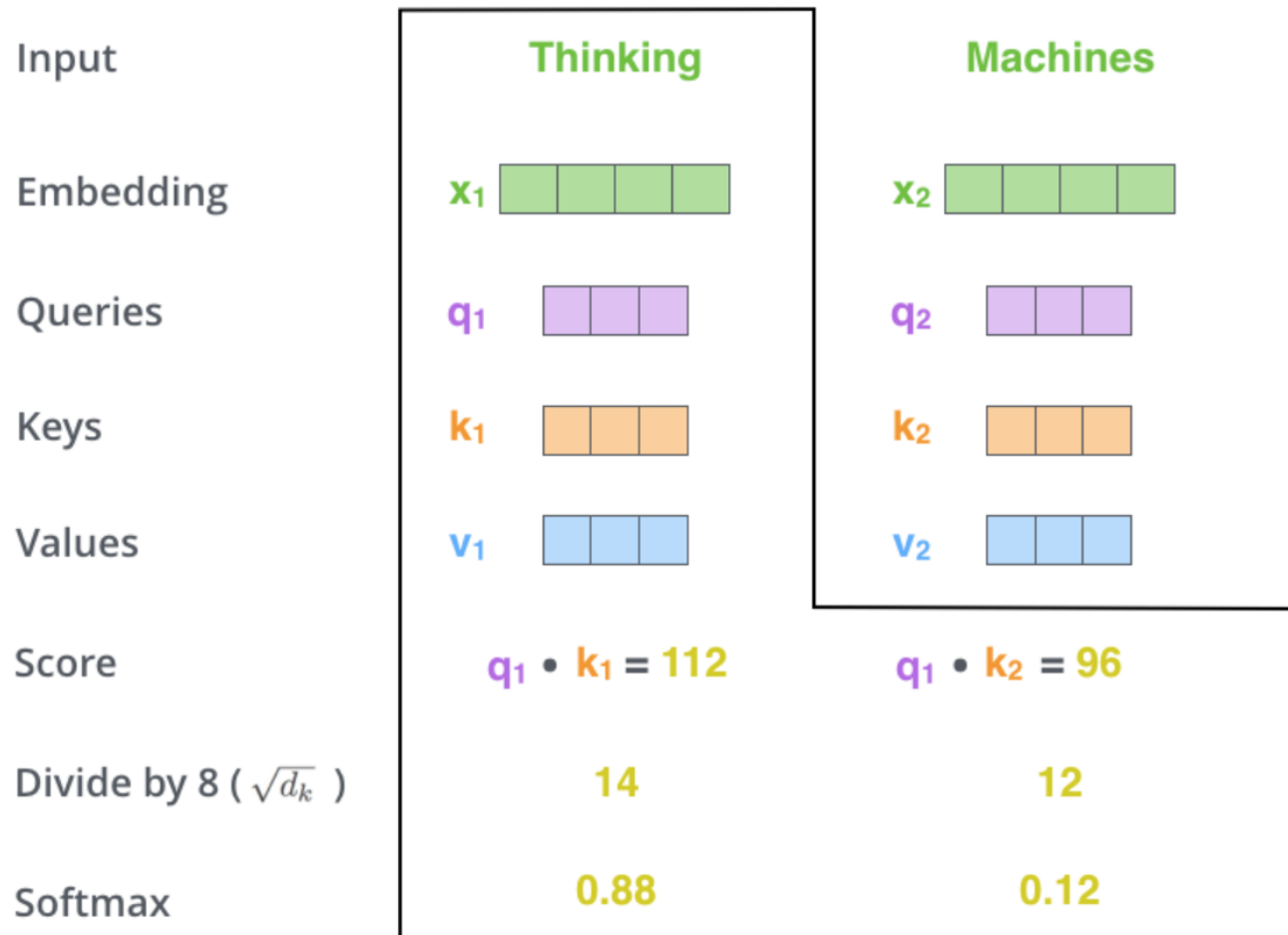
# Details 2/4: scoring



# Details 3/4: normalize scores

- divide the scores by the square root of the dimension of the key vectors used (in example, the vectors are of dimension 64, therefore divide by 8)
- This leads to more stable gradients.
- Then pass the result through a softmax operation. Softmax normalizes the scores so they're all positive and add up to 1.

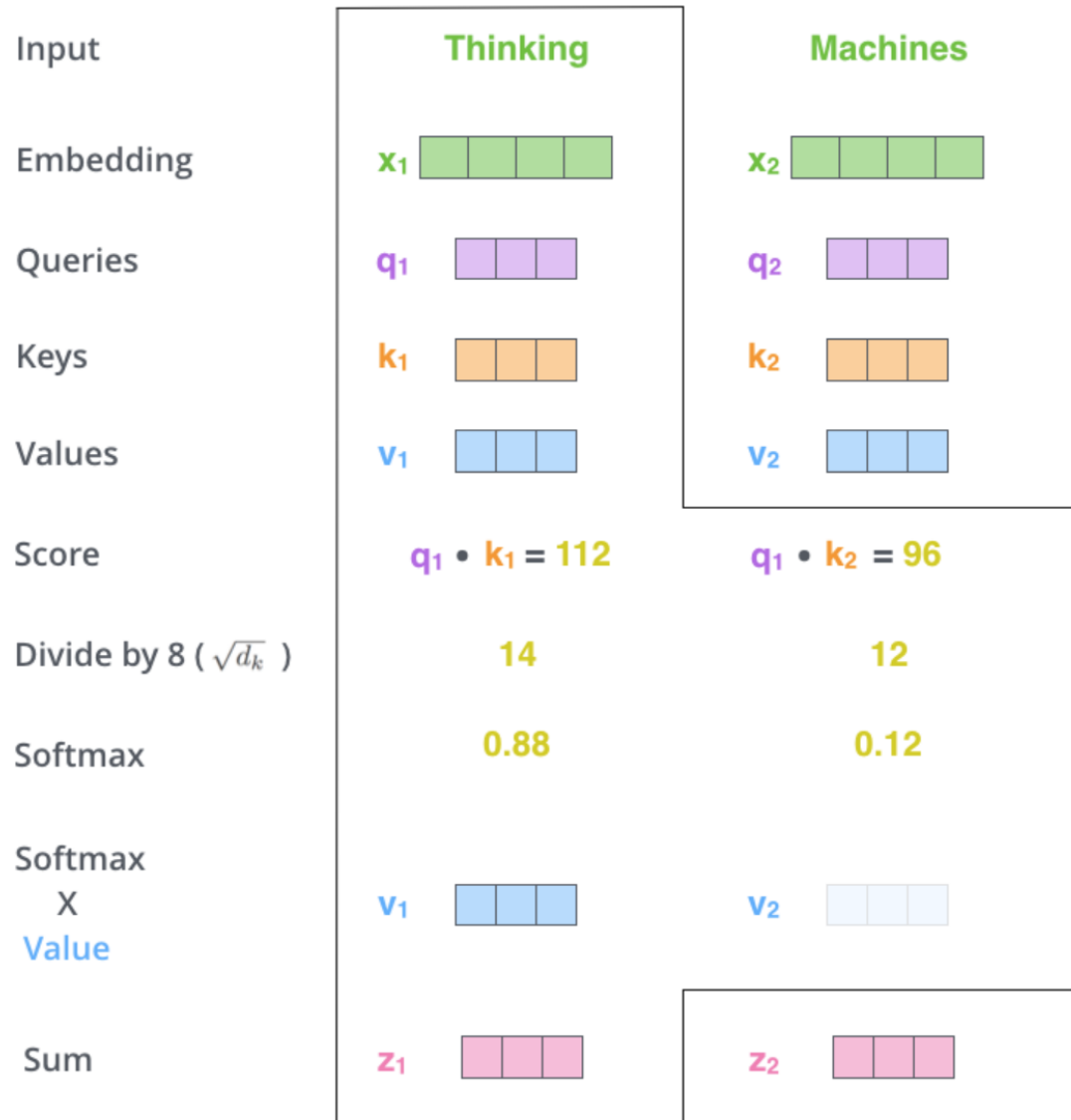
# Details 3/4: normalization of scores

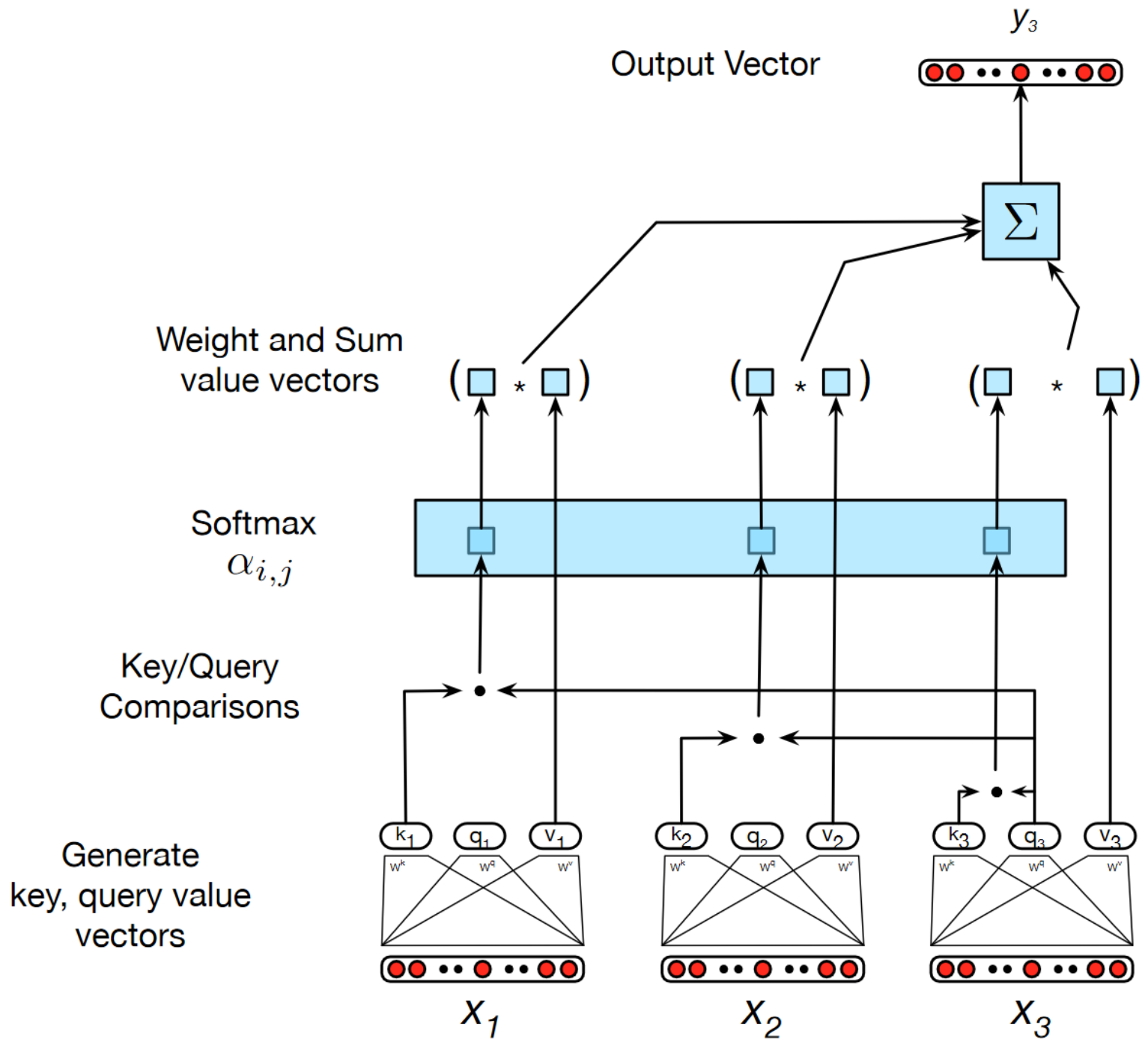


# Details 4/4: apply attention scores

- The softmax score determines how much each word will be expressed at this position. Usually the word at this position will have the highest softmax score, but sometimes it's useful to attend to another word that is relevant to the current word.
- Next multiply each value vector by the softmax score (in preparation to sum them up). The intuition is to keep intact the values of the word(s) we want to focus on, and drown-out irrelevant words (by multiplying them with small scores, e.g., 0.001).
- End computation by summing up the weighted value vectors. This produces the output of the self-attention layer at this position (for the given word – the first one in the example).
- The resulting vector is send to the feed-forward neural network.
- In the actual implementation, however, the calculation is done in matrix form for faster processing.

# Details 4/4: self-attention output







# Matrix calculation of self-attention 1/2

Every row in the X matrix corresponds to a word in the input sentence.  
The embedding vector (512) is larger than the q/k/v vectors (64)



# Matrix calculation of self-attention 2/2

- final calculation

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

=

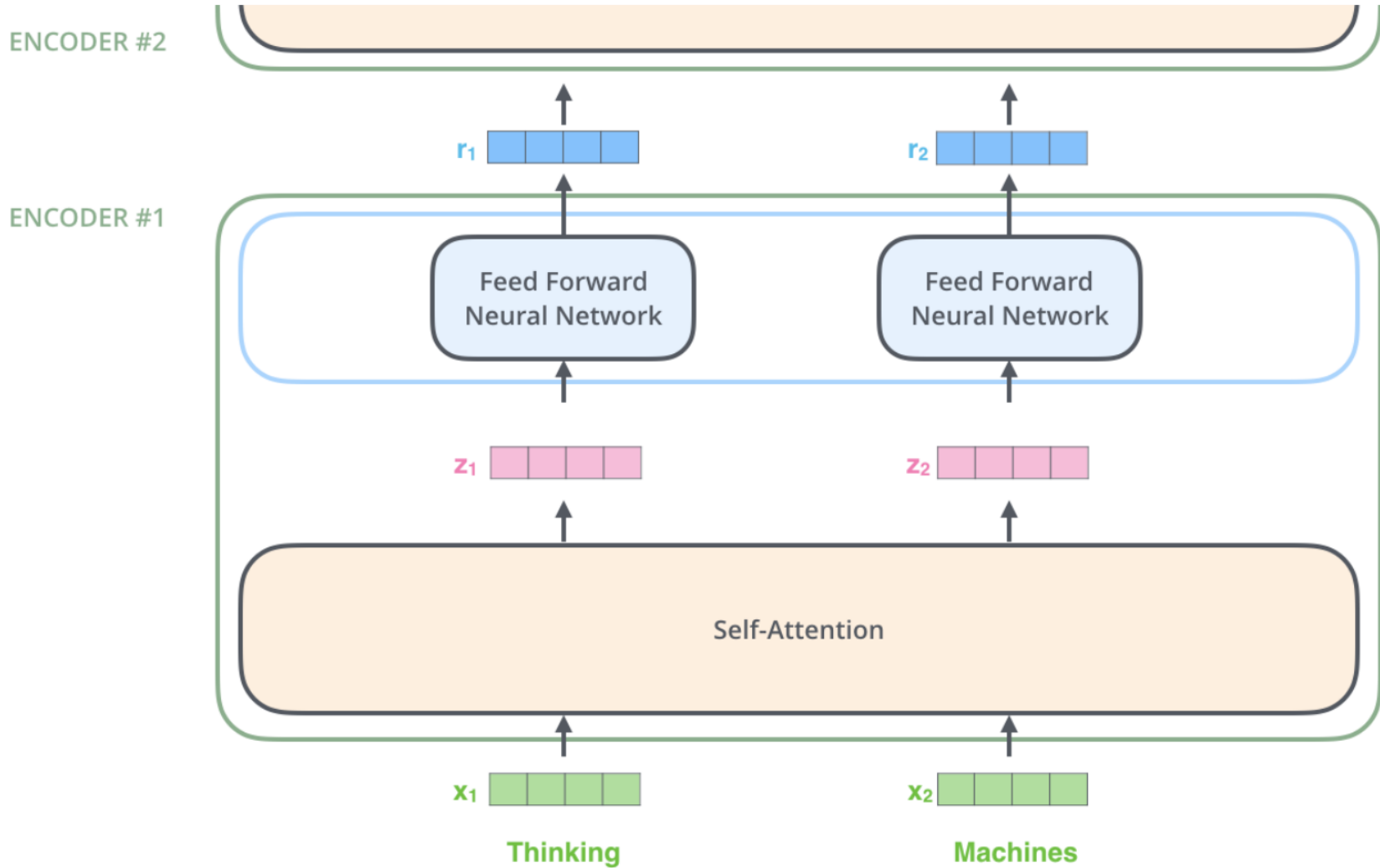
$$\begin{matrix} \text{Z} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$

# Computing attention head

$$A(q, K, V) = \sum_i \frac{e^{q \cdot k_i}}{\sum_j e^{q \cdot k_j}} v_i$$

$$A(Q, K, V) = \text{softmax}(QK^T)V$$

# Encoding

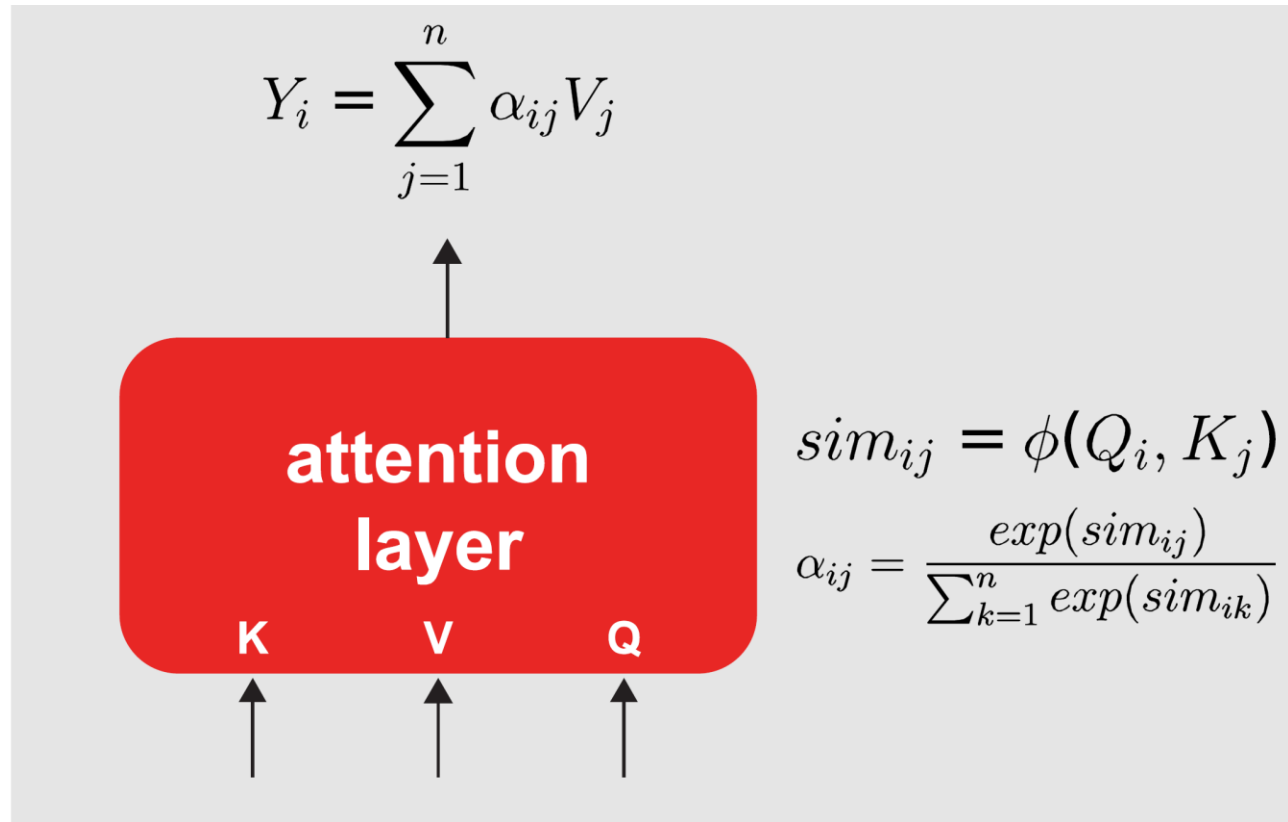


# In summary 1/3

## Self-attention

- recall attention:  
summary of encoder states  
based on similarity between a  
particular decoder state and  
encoder states

- Transformer generalizes:  
summary of values  $V$  based on  
similarity between a particular  
query  $Q_i$  and keys  $K$



# In summary 2/3

## Self-attention

Notation:

$n$  - input length

$d_m$  - representation dim.

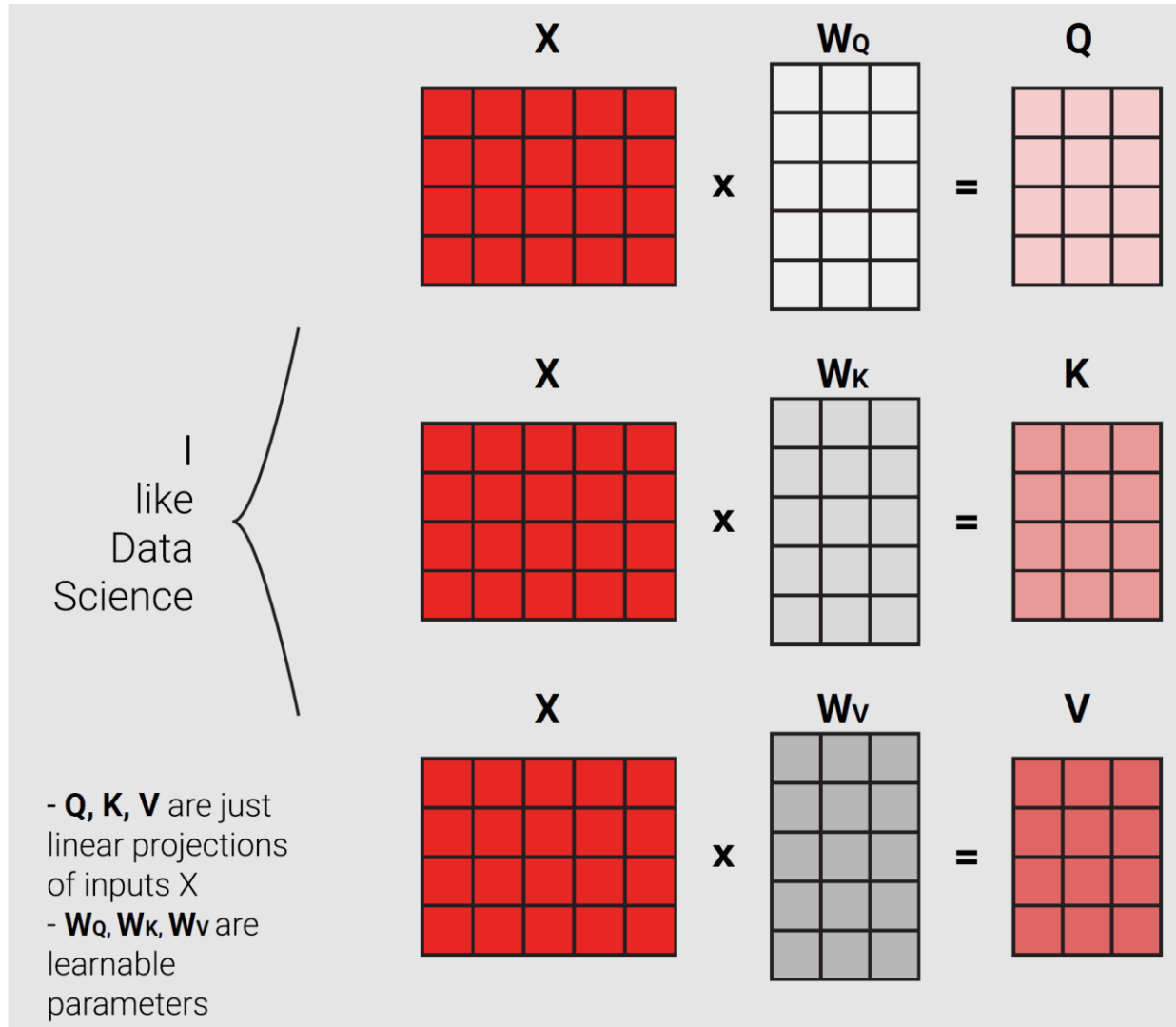
$d$  - dimension of queries,  
keys, values

Example:

$n = 4$

$d_m = 5$

$d = 3$



# In summary 3/3

## Self-attention

- scaled dot-product attention;
- with larger  $d$  the variance of dot products increases, softmax gets more peaked, its gradient gets smaller;
- dividing by  $\sqrt{d}$  mitigates this effect.

$$sim_{ij} = \frac{Q_i^T K_j}{\sqrt{d}} \quad \alpha_{ij} = \frac{\exp(sim_{ij})}{\sum_{k=1}^n \exp(sim_{ik})} \quad Y_i = \sum_{j=1}^n \alpha_{ij} V_j$$

Transformer attention

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

matrix formulation

$$\text{softmax}\left(\frac{\begin{matrix} \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \end{matrix} \times \begin{matrix} \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \end{matrix}}{\sqrt{d}}\right) = \begin{matrix} \color{red}\square & \color{red}\square & \color{red}\square & \color{red}\square \\ \color{red}\square & \color{red}\square & \color{red}\square & \color{red}\square \\ \color{red}\square & \color{red}\square & \color{red}\square & \color{red}\square \\ \color{red}\square & \color{red}\square & \color{red}\square & \color{red}\square \end{matrix}$$

attention weights

$$\begin{matrix} \color{red}\square & \color{red}\square & \color{red}\square & \color{red}\square \\ \color{red}\square & \color{red}\square & \color{red}\square & \color{red}\square \\ \color{red}\square & \color{red}\square & \color{red}\square & \color{red}\square \\ \color{red}\square & \color{red}\square & \color{red}\square & \color{red}\square \end{matrix} \times \begin{matrix} \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \end{matrix} = \begin{matrix} \color{red}\square & \color{red}\square & \color{red}\square & \color{red}\square \\ \color{red}\square & \color{red}\square & \color{red}\square & \color{red}\square \\ \color{red}\square & \color{red}\square & \color{red}\square & \color{red}\square \\ \color{red}\square & \color{red}\square & \color{red}\square & \color{red}\square \end{matrix}$$

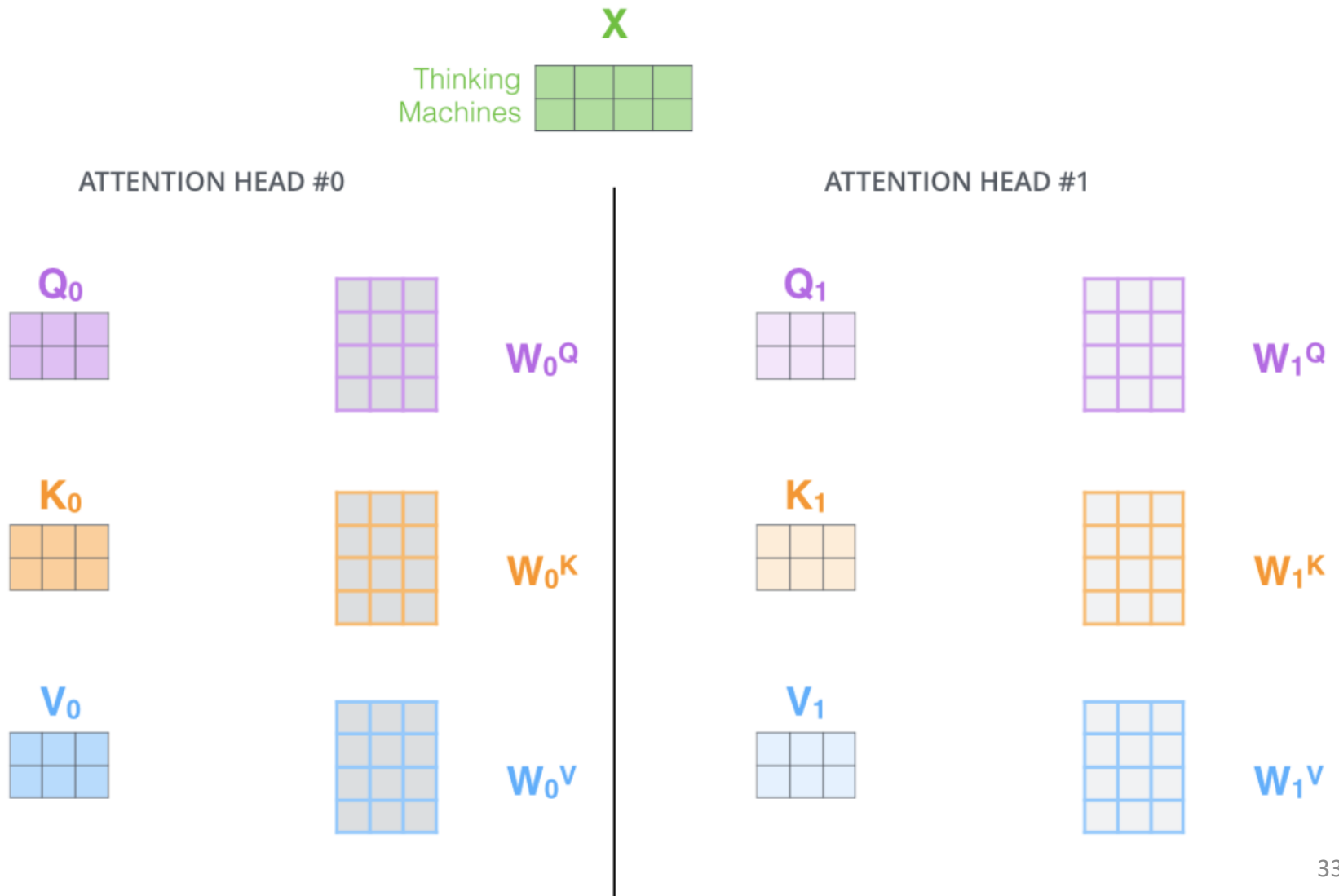
output

# Multi-headed attention

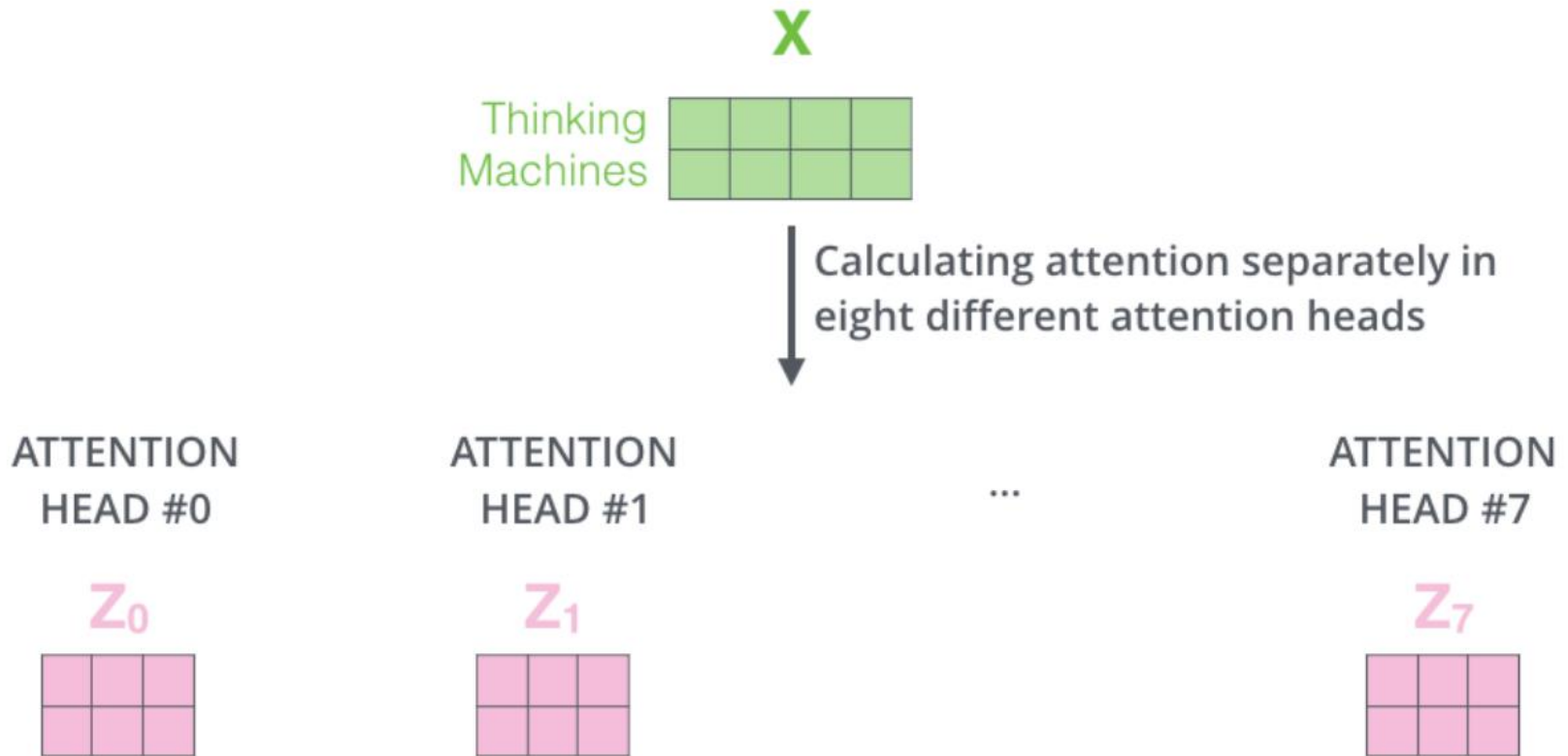
- Self-attention layer is replicated several times, called “multi-headed” attention.
- This expands the model’s ability to focus on different positions. E.g., one attention head might be dominated by the the actual word, but other heads might reveal other important information
- E.g., in translating a sentence like “The animal didn’t cross the street because it was too tired”, we would want to know which word “it” refers to.
- Multi-head attention layer can cover multiple “representation subspaces”
- I.e., we have multiple sets of Query/Key/Value weight matrices (original Transformer uses 8 attention heads)
- Each attention head is randomly initialized. After training, each set is used to project the input embeddings (or vectors from lower encoders/decoders) into a different representation subspace.



# Example: two attention heads



# Example: 8 att. heads



- What to do with 8 Z matrices, the feed-forward layer is expecting a single matrix (one vector for each word). We need to condense all attention heads into one matrix.

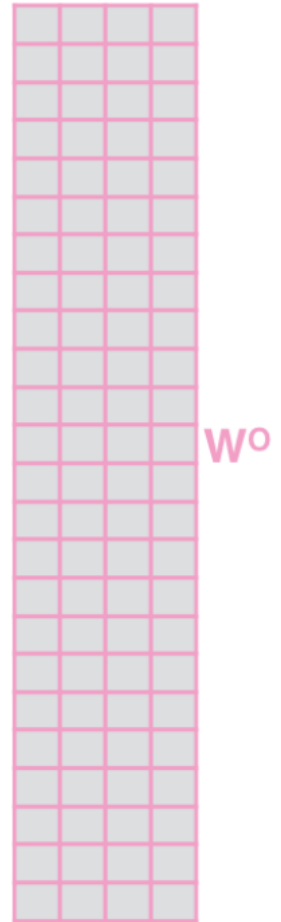
# Condensation of attention heads

1) Concatenate all the attention heads

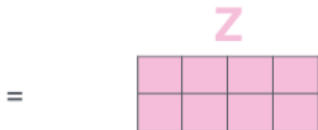


2) Multiply with a weight matrix  $W^O$  that was trained jointly with the model

x



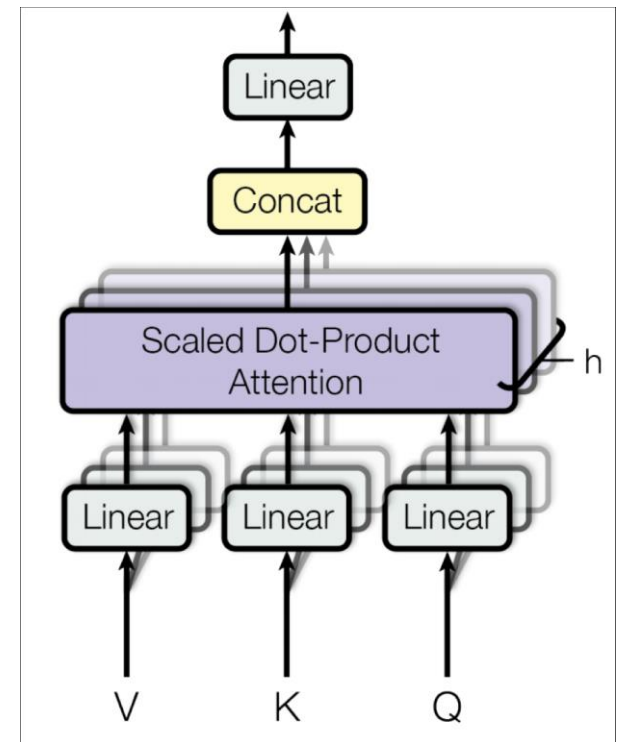
3) The result would be the  $Z$  matrix that captures information from all the attention heads. We can send this forward to the FFNN



# Computing multi-head attention

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

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



# Summary of self-attention

1) This is our input sentence\*

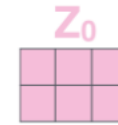
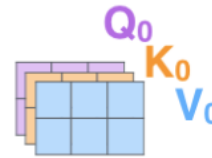
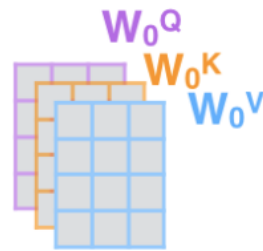
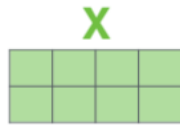
2) We embed each word\*

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

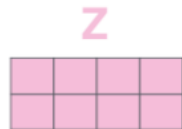
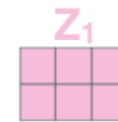
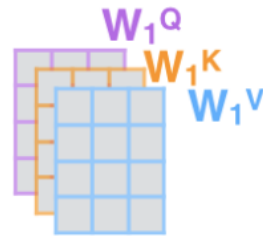
4) Calculate attention using the resulting  $Q/K/V$  matrices

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

Thinking Machines



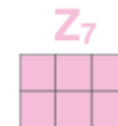
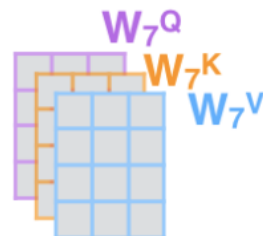
\* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



...

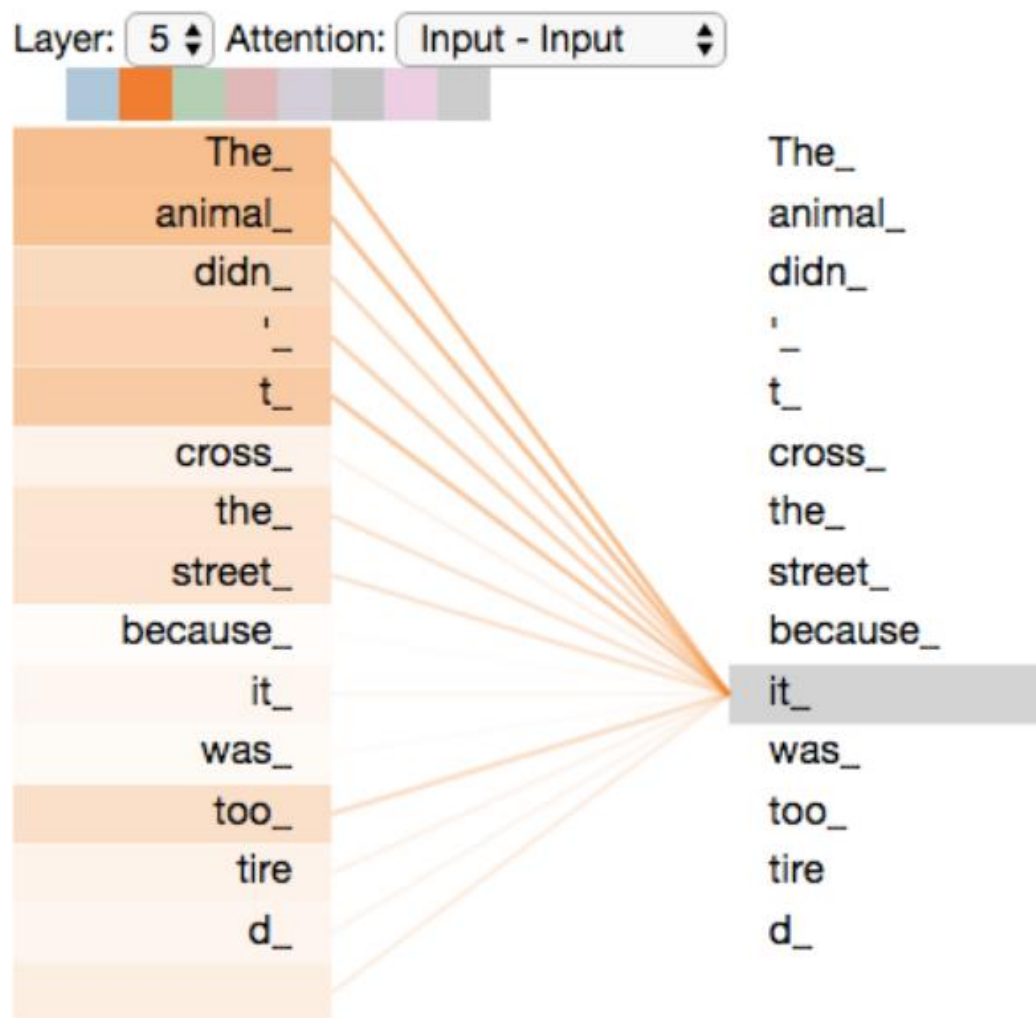
...

...



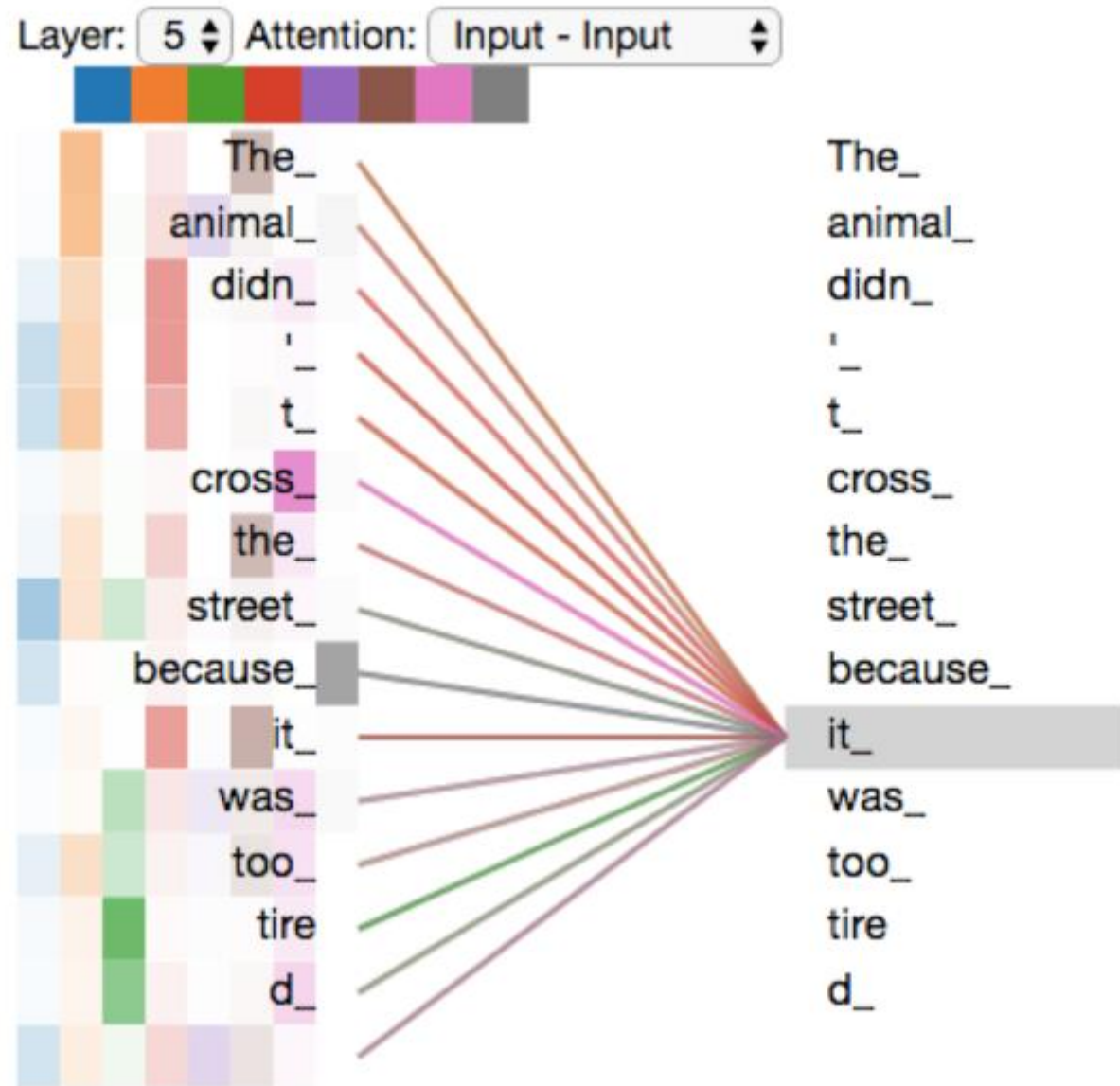
# Illustration of self-attention: 1 head

- encoder #5 (the top encoder in the stack)
- As we encode the word "it", one attention head is focusing most on "the animal", while another is focusing on "tired" -- in a sense, the model's representation of the word "it" bakes in some of the representation of both "animal" and "tired".



# Illustration of self-attention: all heads

- all the attention heads in one picture are harder to interpret

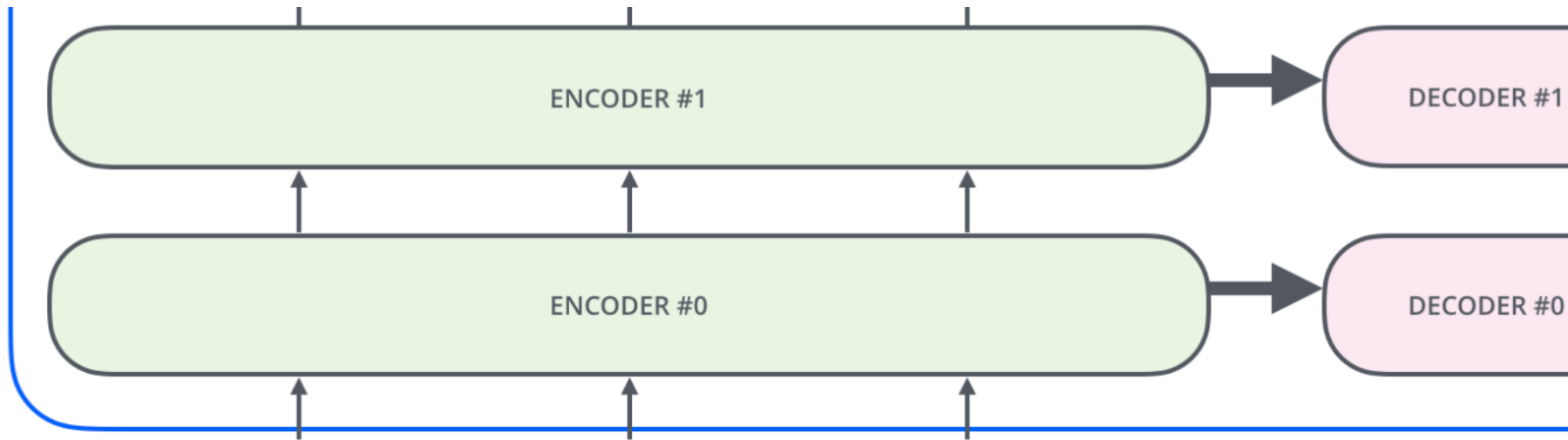


# Representing the order of the sequence using positional encoding

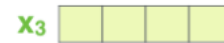
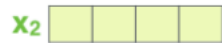
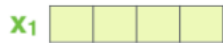
- Order of the sequence is important but it is lost with the described transformation, therefore
- The transformer adds a position vector to each input embedding.
- These vectors follow a specific pattern that the model learns, which helps it determine the position of each word, or the distance between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



# Adding position encoding



EMBEDDING  
WITH TIME  
SIGNAL

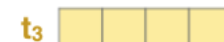
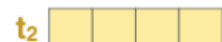
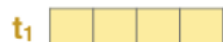


=

=

=

POSITIONAL  
ENCODING

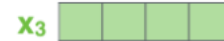
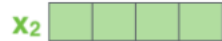
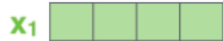


+

+

+

EMBEDDINGS



INPUT

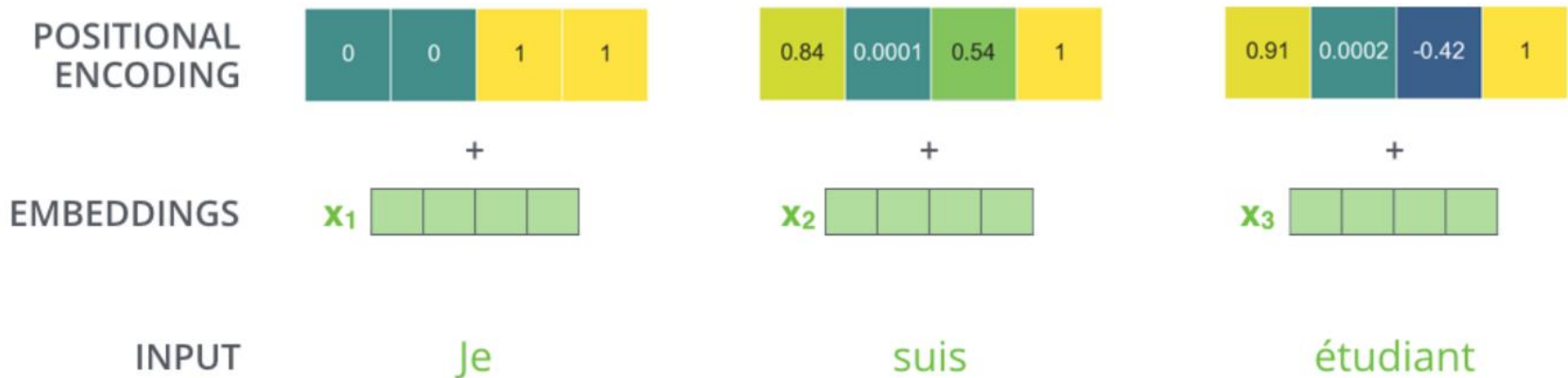
Je

suis

étudiant

# Example: encoding position

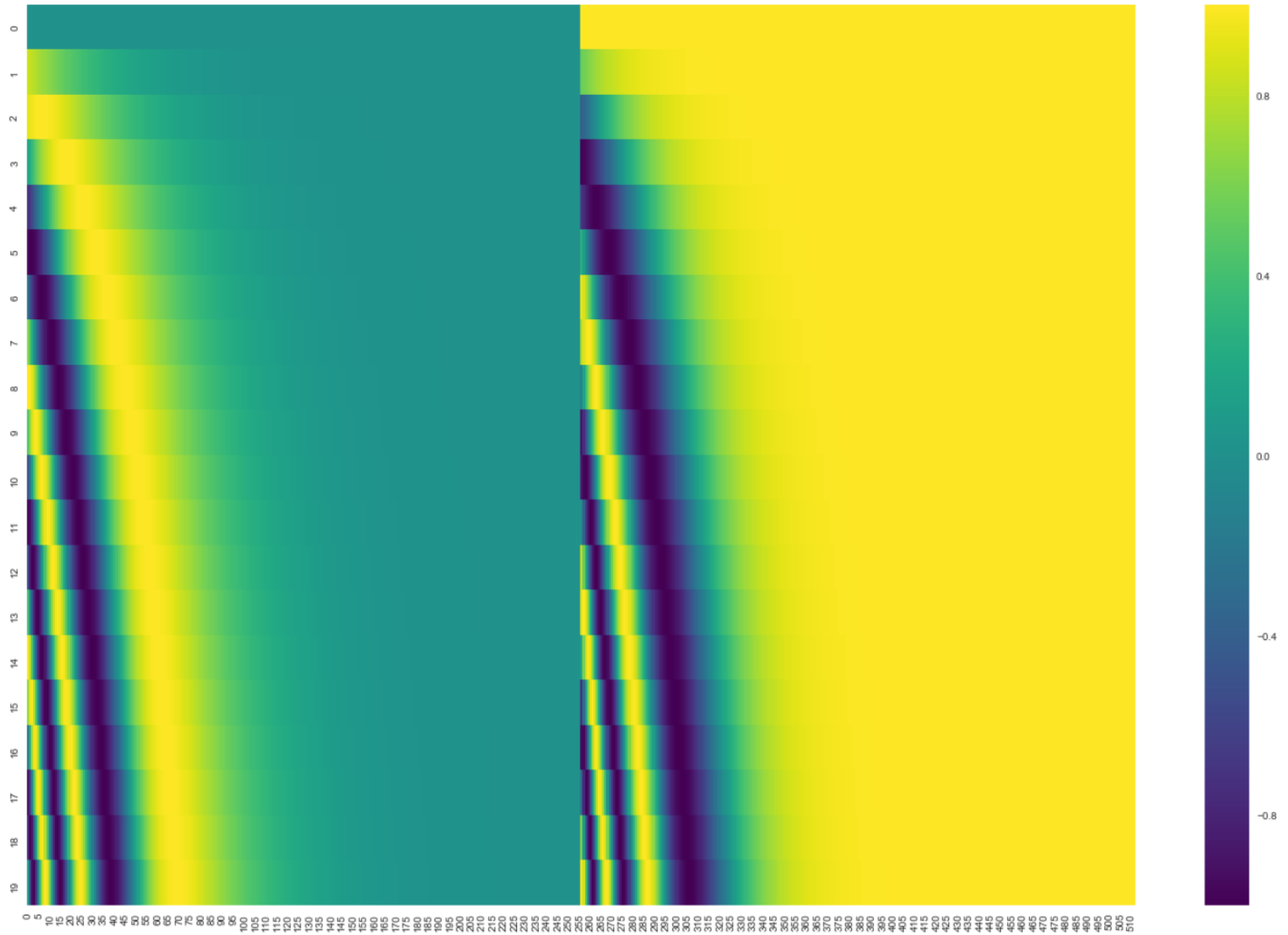
- the values of positional encoding vectors follow a specific pattern.



# Patterns of positional encodings

- many different possibilities how to generate a pattern
- next slide contains an example of positional encoding for 20 words (rows) with an embedding size of 512 (columns).
- the values of the left half are generated by one function (which uses sine), and the right half is generated by another function (which uses cosine). They're then concatenated to form each of the positional encoding vectors

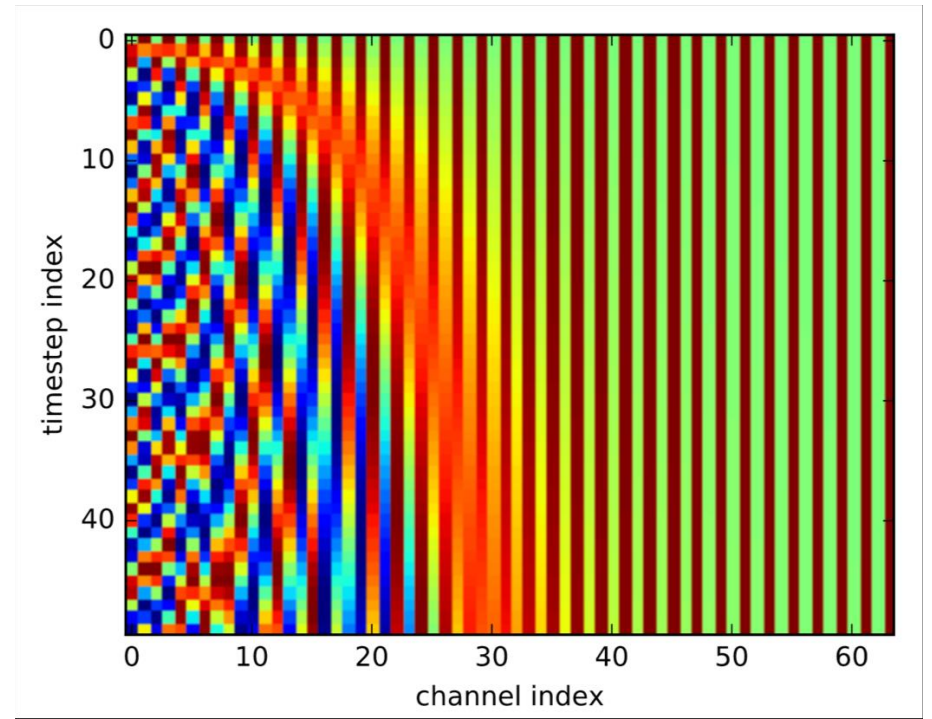
# Example of positional encoding



# Another positional encoding

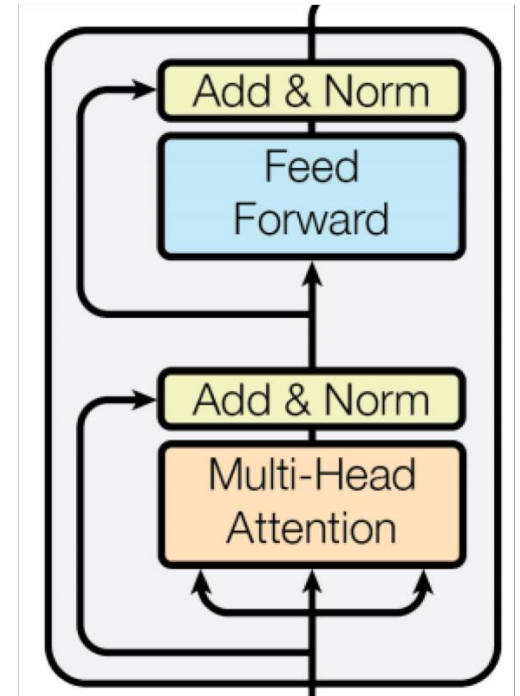
$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

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



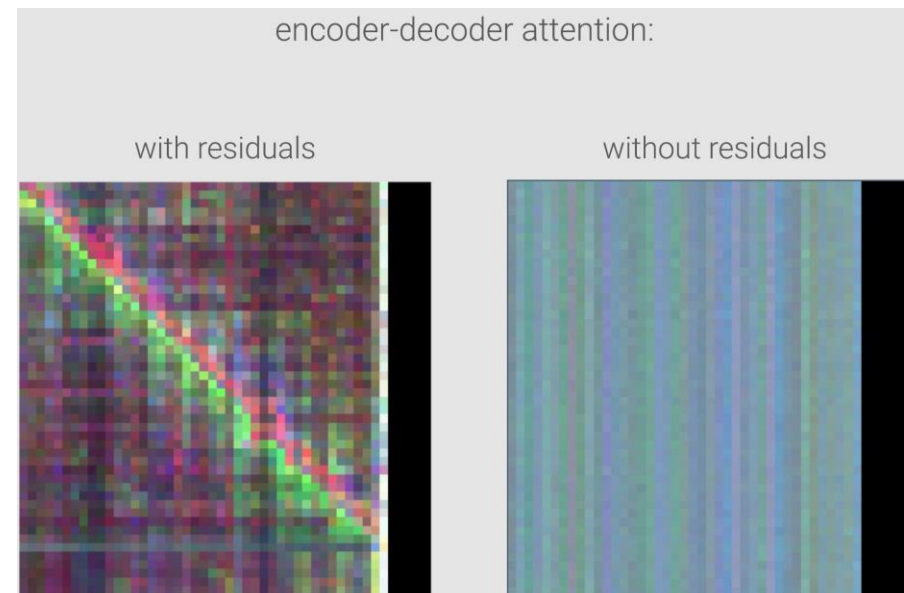
# Encoder blocks

- Each block has two “sublayers”
  - Multihead attention
  - 2-layer feed-forward NNet (with ReLU)
- Each of these two steps also has residual (short-circuit) connection and LayerNorm
  - $\text{LayerNorm}(x + \text{Sublayer}(x))$

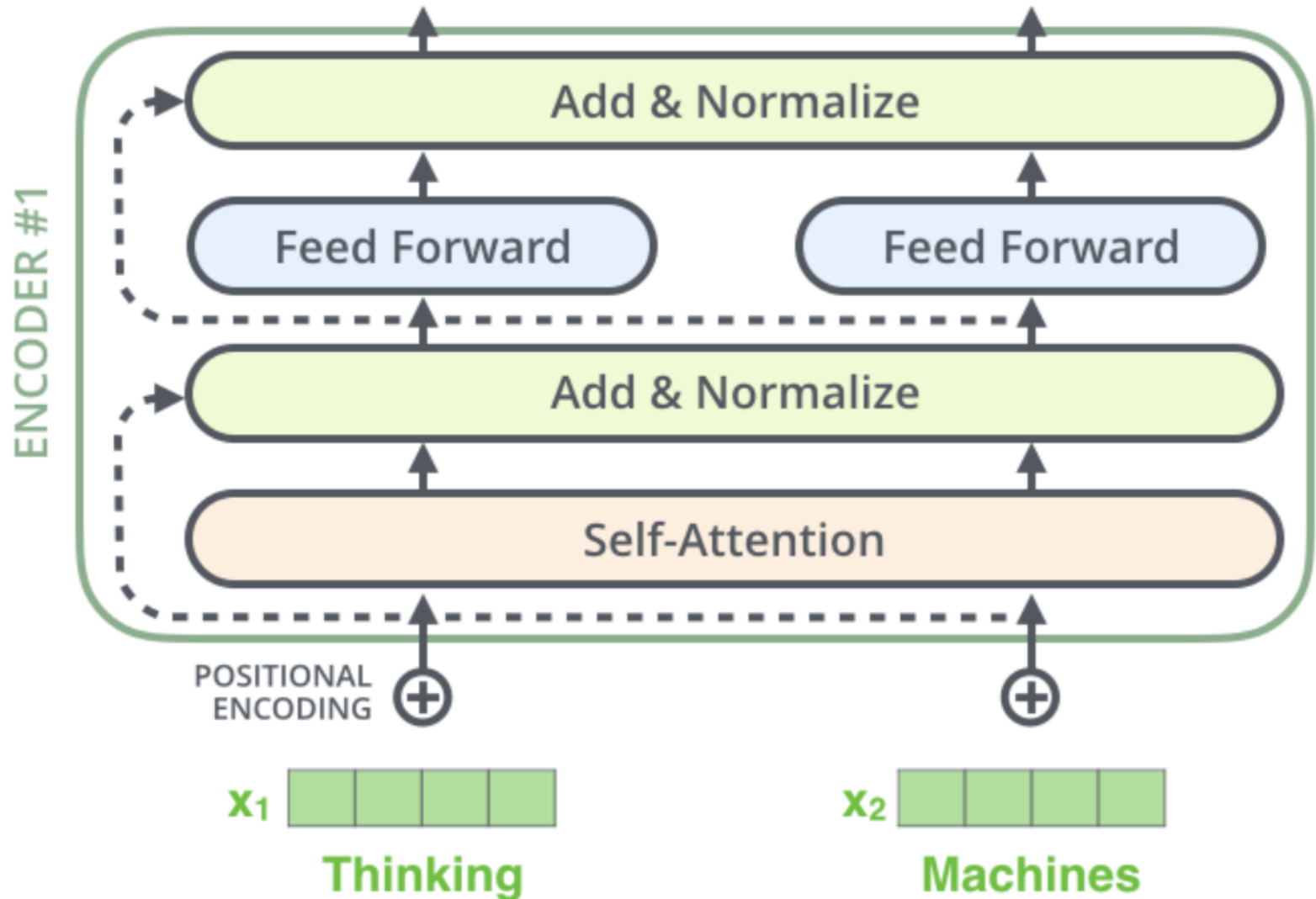


# The Residual connections

- each sub-layer of transformer (self-attention and feed-forward NN) in each encoder has a residual connection around it, and is followed by a layer-normalization step.
- the same for decoder sub-layers
- enable learning of deeper networks by improving a gradient flow
- in transformers they also maintain positional information in higher layers

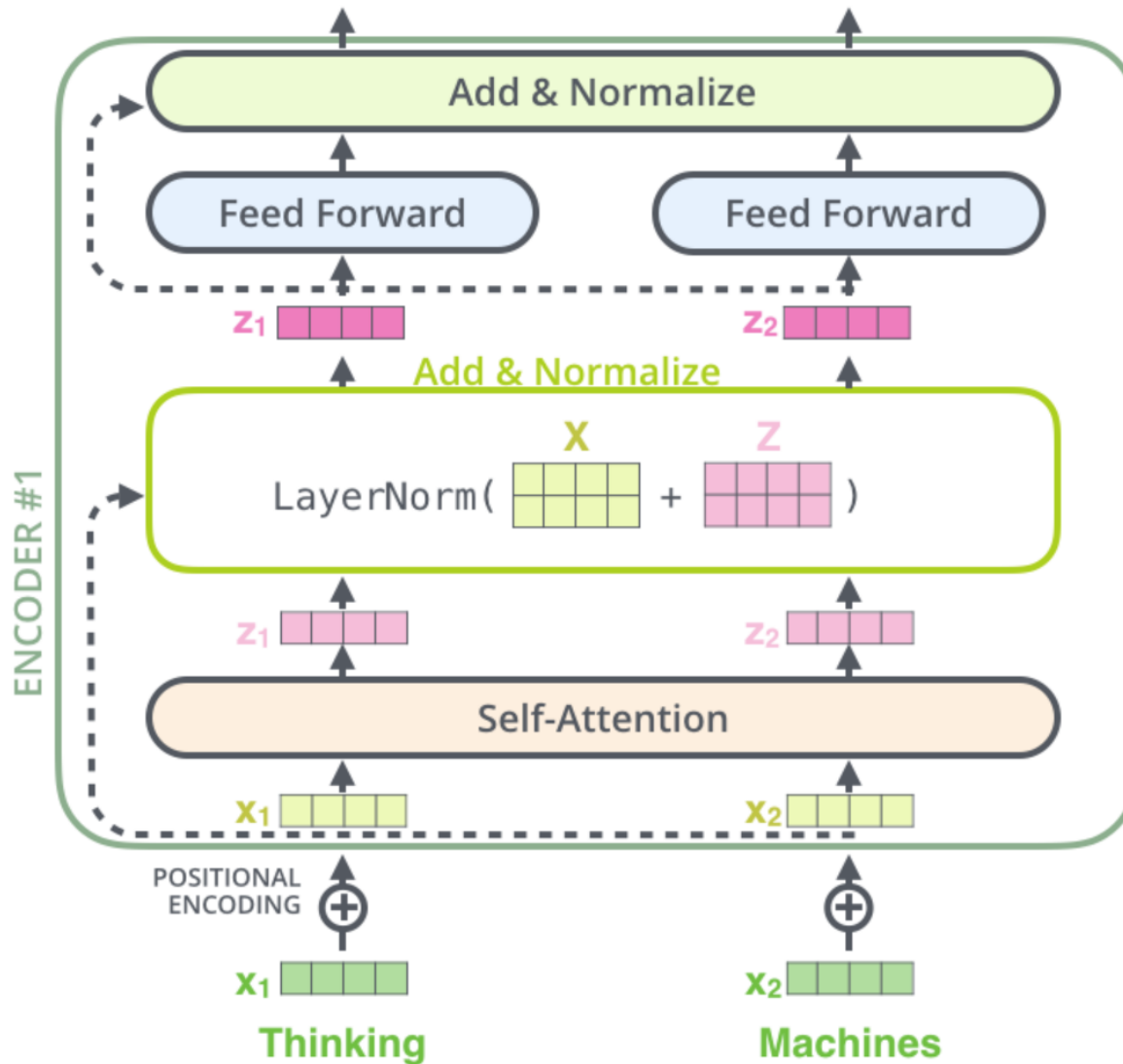


# Architecture with residual connection – top level view

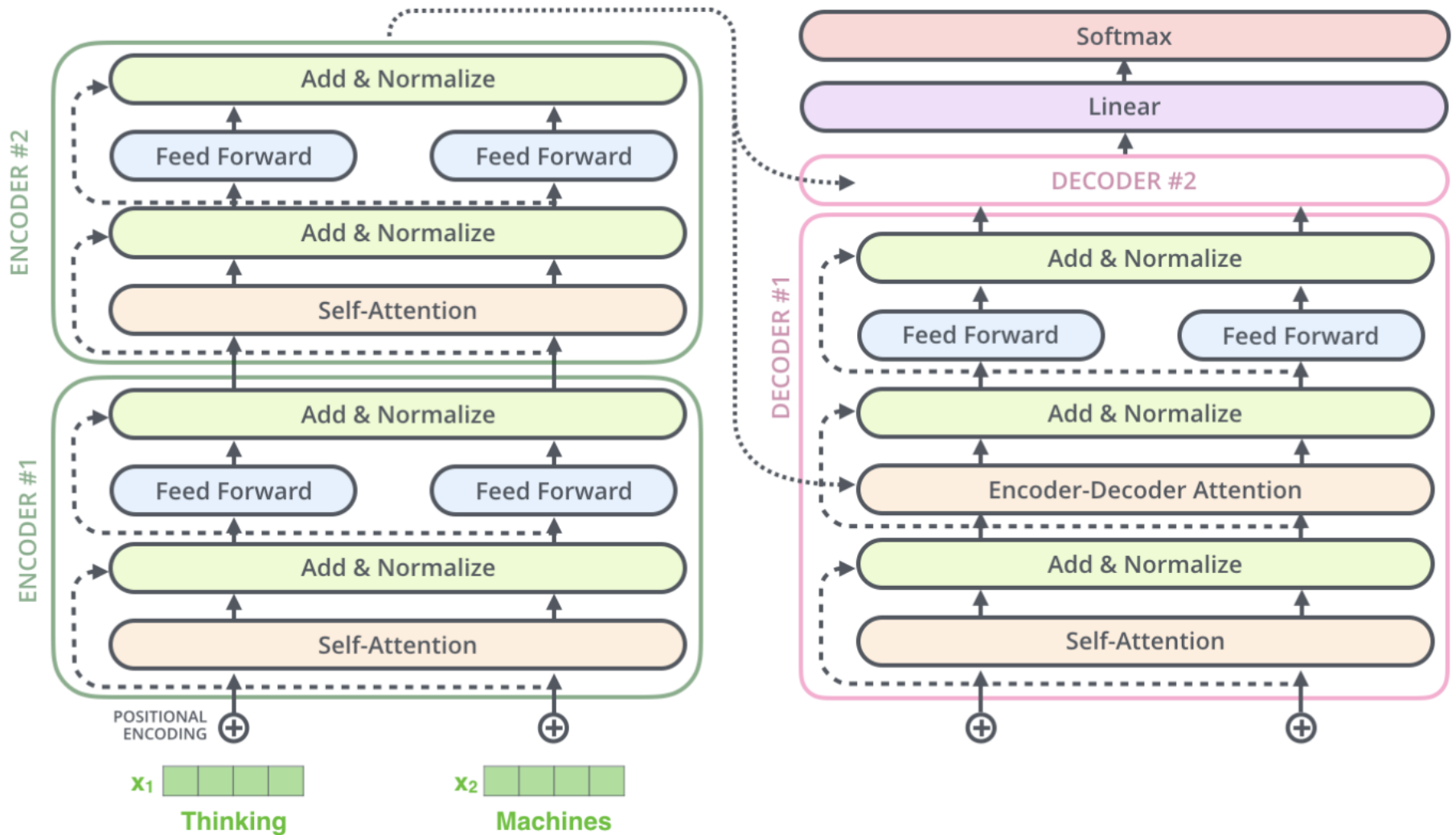




# Architecture with residual connection – example

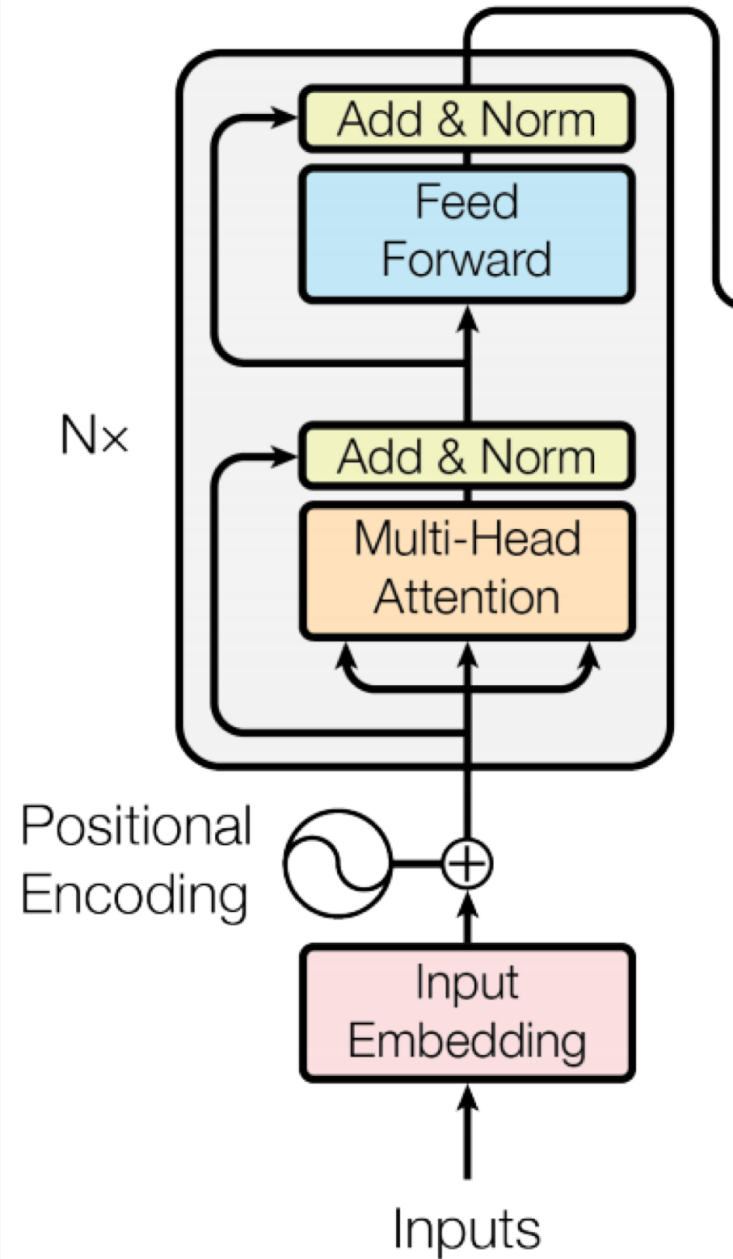


# Example: 2 stacked transformer



# Complete encoder

- each block is repeated several times, e.g., 6 times



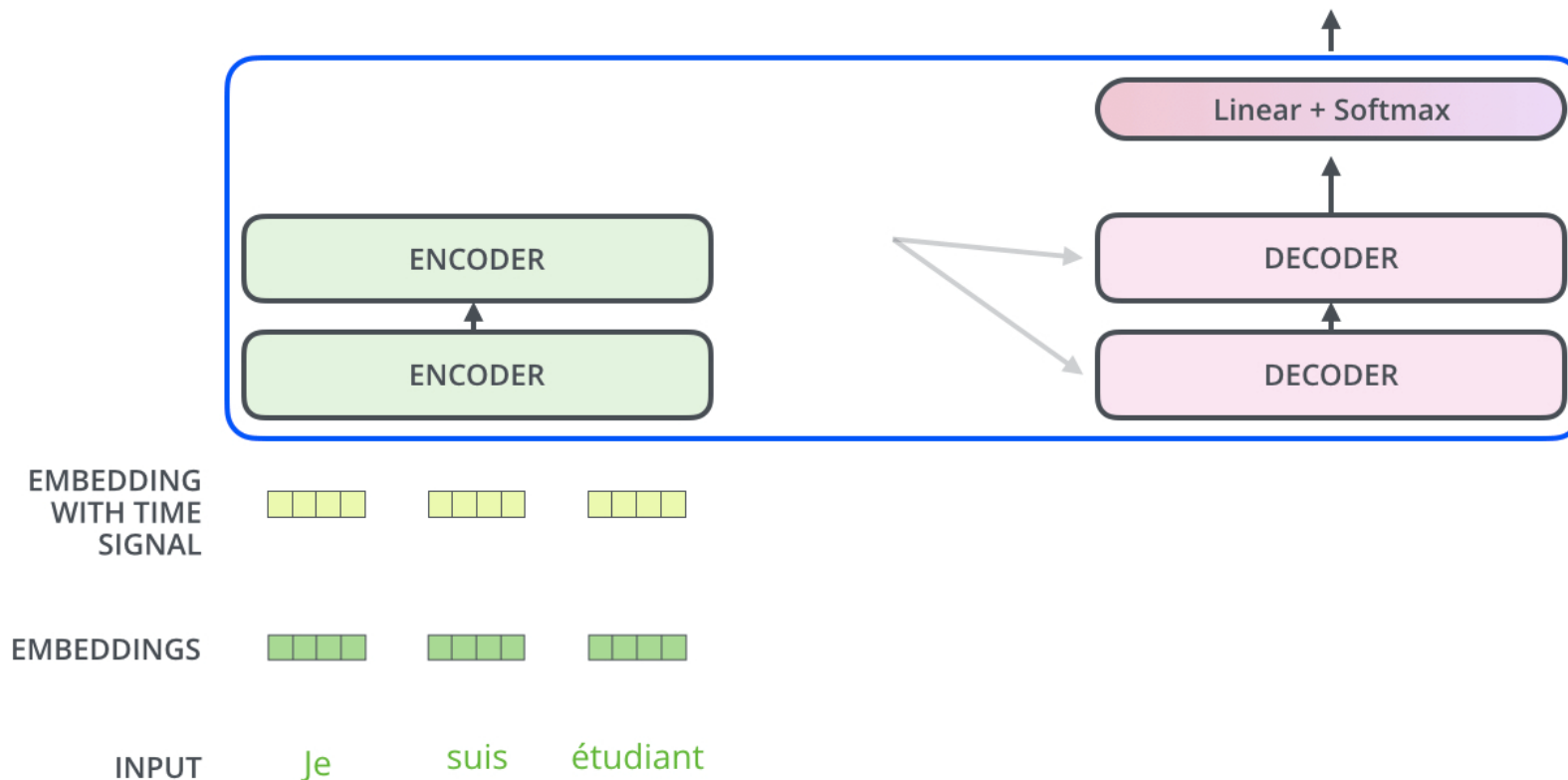
# Decoder

- Decoders have the same components as encoders
- An encoder start by processing the input sequence.
- The output of the top encoder is transformed into a set of attention vectors  $K$  and  $V$ .
- These are used by each decoder in its “encoder-decoder attention” layer which helps the decoder to focus on appropriate places in the input sequence.

# Encoder-decoder in action 1/2

Decoding time step: ① 2 3 4 5 6

OUTPUT

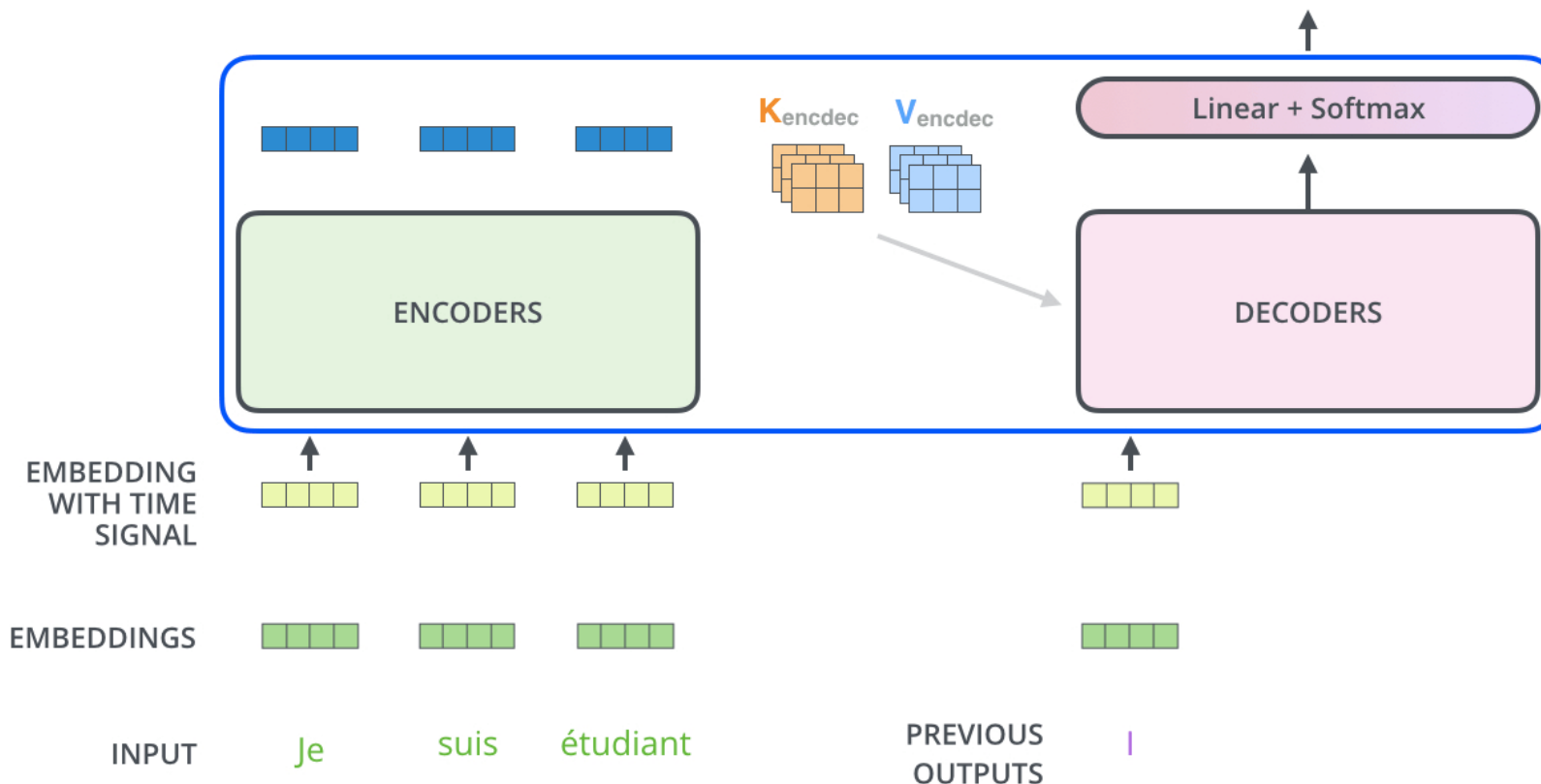


After finishing the encoding phase, we begin the decoding phase. Each step in the decoding phase outputs an element from the output sequence (the English translation sentence in this case).

# Encoder-decoder in action 2/2

Decoding time step: 1 2 3 4 5 6

OUTPUT |



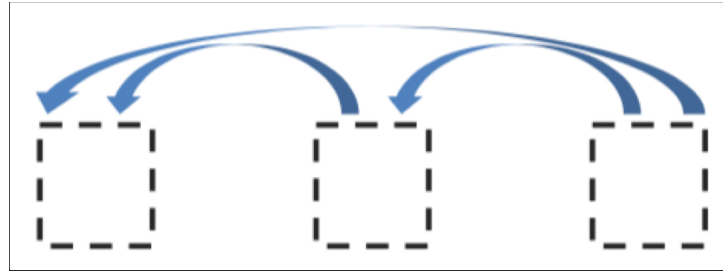
The steps repeat until a special symbol indicating the end of output. The output of each step is fed to the bottom decoder in the next time step. We add positional encoding to decoder inputs to indicate the position of each word.

# Self-attention and encoder-decoder attention in the decoder

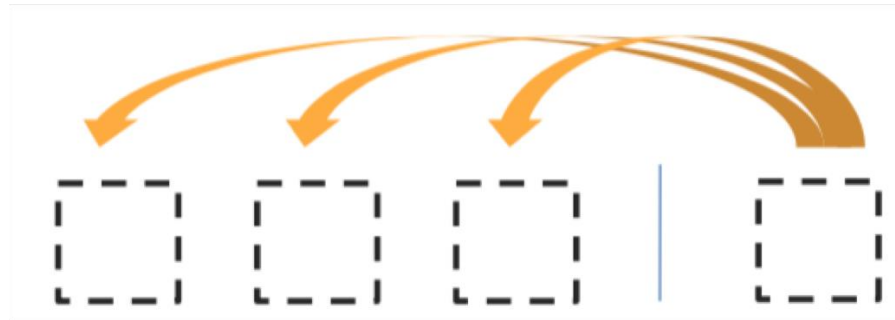
- In the decoder, the self-attention layer is only allowed to attend to itself and earlier positions in the output sequence (to maintain the autoregressive property).
- This is done by masking future positions (setting them to  $-\infty$ ) before the softmax step in the self-attention calculation.
- The “Encoder-Decoder Attention” layer works just like multiheaded self-attention, except it creates its Q (queries) matrix from the layer below it, and takes the K (keys) and V (values) matrix from the output of the encoder stack.

# Attentions in the decoder

1. Masked decoder self-attention on previously generated outputs

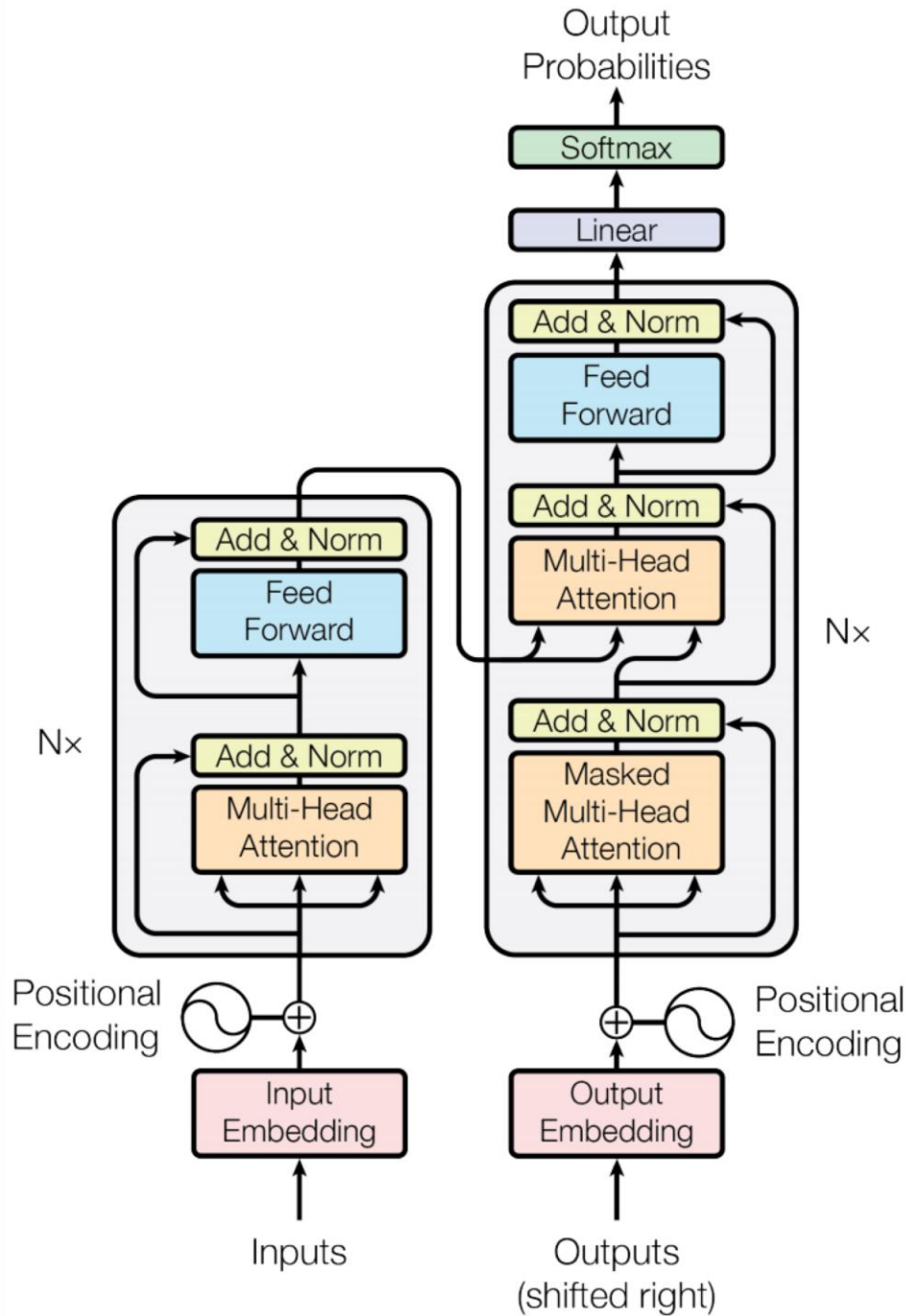


2. Encoder-Decoder Attention, where queries come from previous decoder layer and keys and values come from output of the encoder





# One encoder-decoder block



# Final Linear and Softmax Layer

- The decoder stack outputs a vector of floats.
- The final linear layer which is followed by a softmax layer turns them into words.
- The Linear layer is a simple fully connected neural network that projects the vector produced by the stack of decoders, into a much larger vector called a logits vector (probability scores for each word).
- Example: the model knows 10,000 unique English words (“output vocabulary”) that it’s learned from its training dataset. Therefore, the logits vector is 10,000 cells wide – each cell corresponding to the score of a unique word.
- The softmax layer turns those scores into probabilities (all positive, between 0 and 1, sum to 1.0).  
The cell with the highest probability is chosen, and the word associated with it is produced as the output for this time step.

# Producing the output words

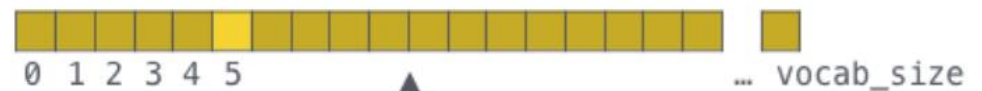
Which word in our vocabulary is associated with this index?

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

am

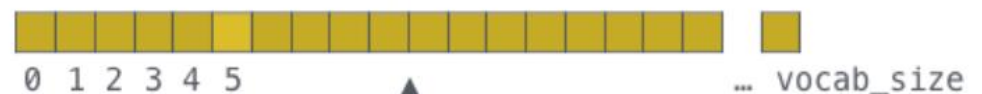
5

Log\_probs



Softmax

logits



Linear

Decoder stack output



# Training the transformer

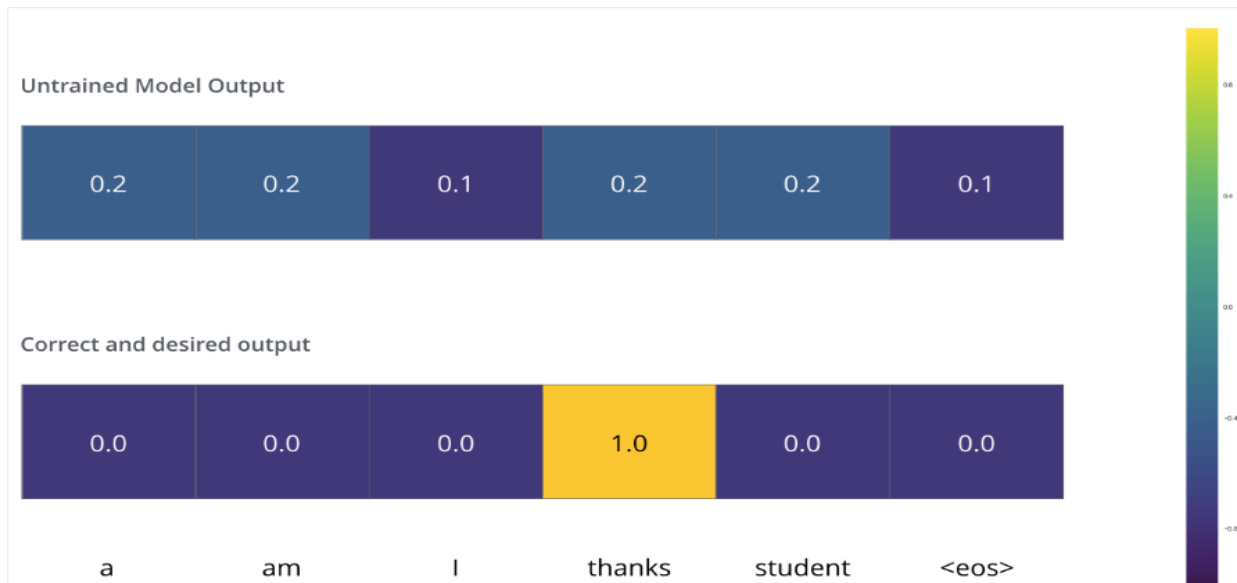
- During training, an untrained model would go through the exact same forward pass. But since we are training it on a labeled training dataset, we can compare its output with the actual correct output.
- For illustration, let's assume that our output vocabulary only contains six words(a, am, i, thanks, student, <eos>)
- The input is typically in the order of  $10^4$  (e.g., 30 000)

Output Vocabulary

WORD	a	am	i	thanks	student	<eos>
INDEX	0	1	2	3	4	5

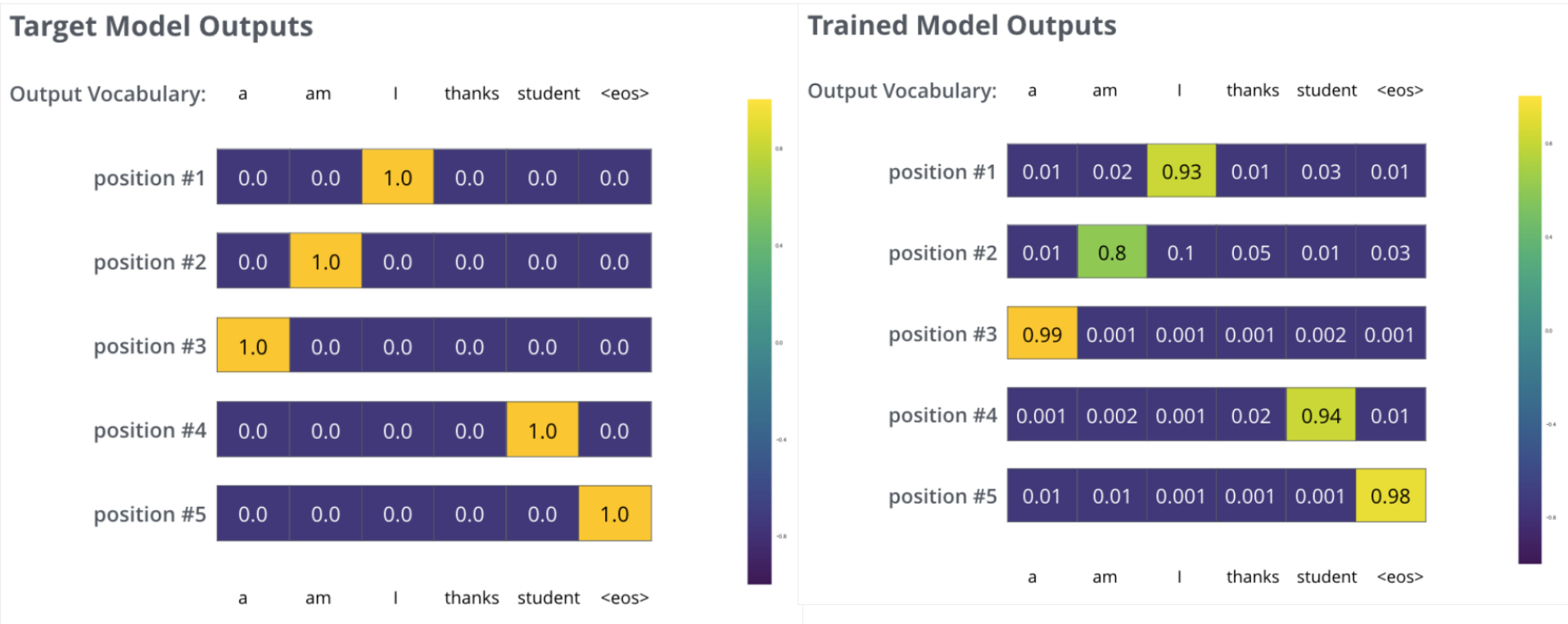
# The Loss Function

- evaluates the difference between the true output and the returned output
- transformer typically uses cross-entropy or Kullback–Leibler divergence.
- both true and returned output are 1-hot encoded, e.g.,



# Loss evaluation for sequences

- loss function has to be evaluated for the whole sentence, not just a single word
- transformers use greedy decoding or beam search



# Animated workings of transformer

# Transformer tricks

- Byte-pair encodings for input tokens
- Checkpoint averaging
- ADAM optimizer with learning rate changes
- Dropout during training at every layer just before adding residual
- Label smoothing
- Auto-regressive decoding with beam search and length penalties
- Use of transformers is spreading but they are hard to optimize and unlike LSTMs don't usually just work out of the box and they don't play well yet with other building blocks on tasks.