

# Transformer Model

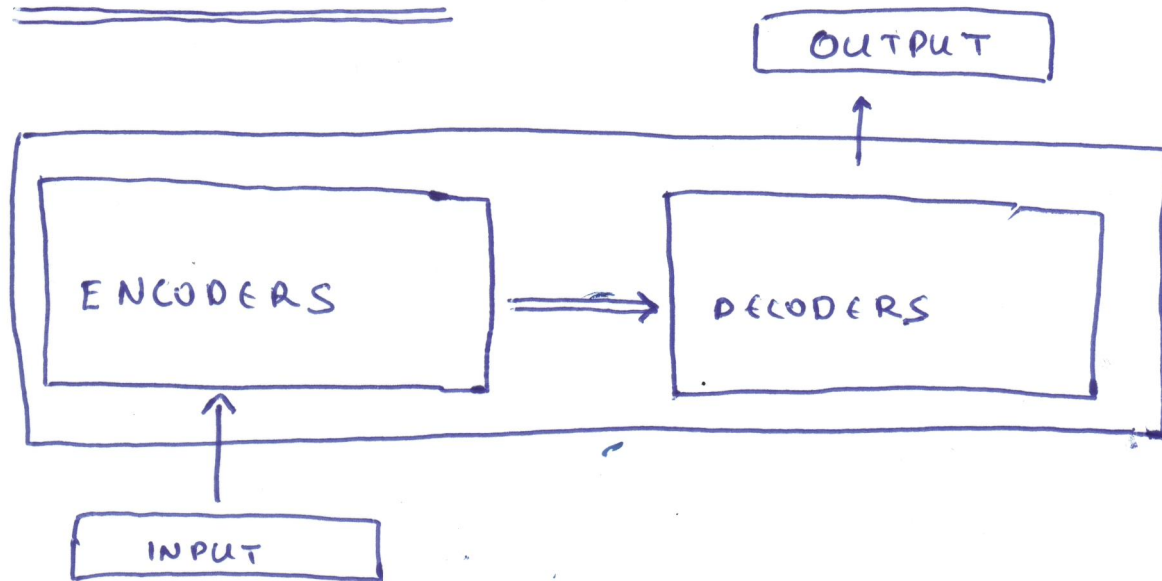
#1

## BIRDS EYE VIEW:



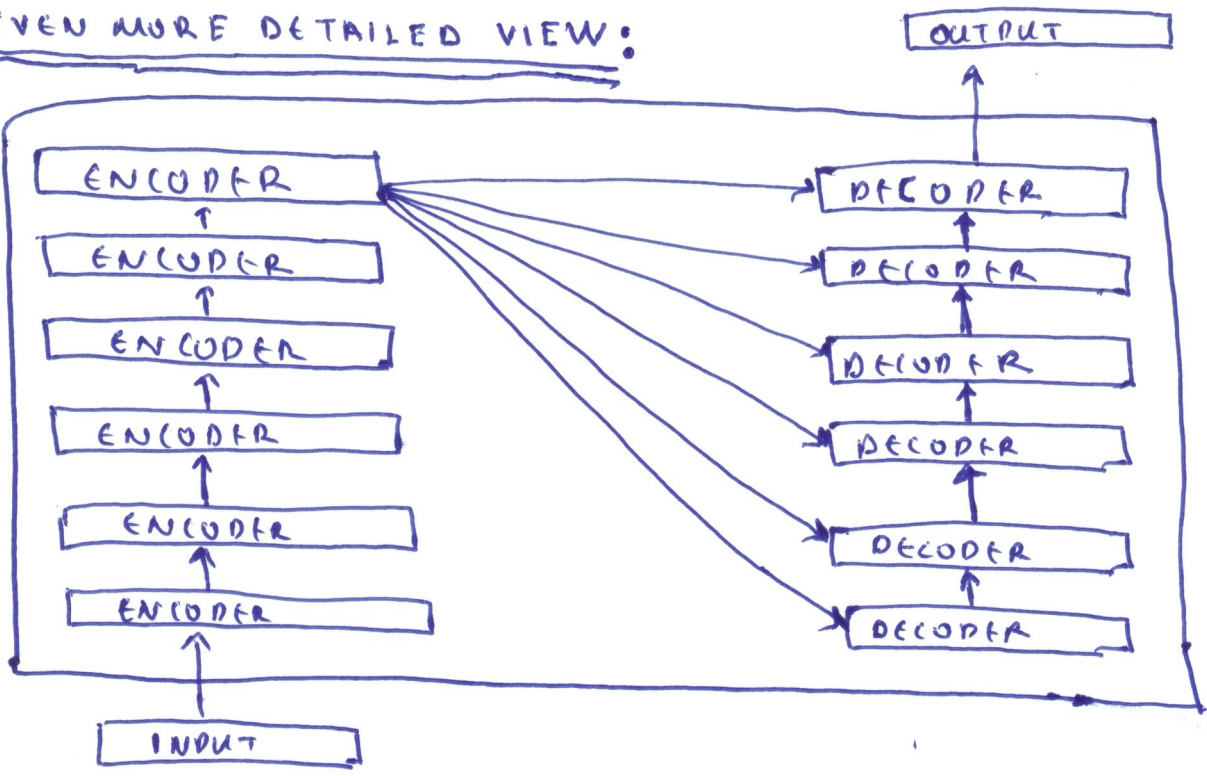
MAIN COMPONENTS : (a) Encoding Component  
(b) Decoding Component

## MORE DETAILED VIEW:

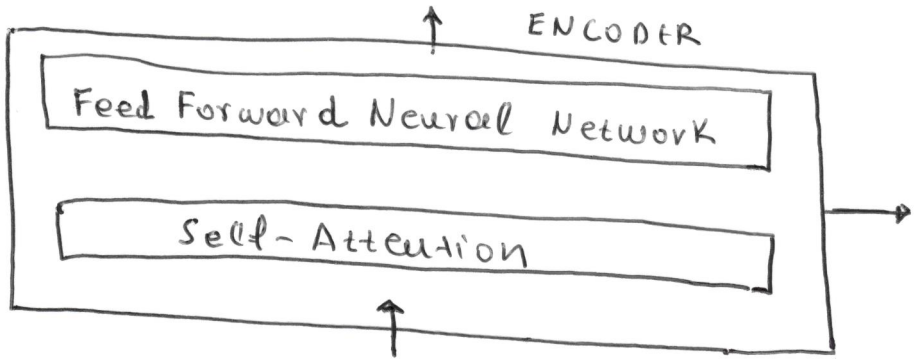


- ★ The encoding component is a stack of encoders. The original paper stacks six of them.
- ★ The decoding component is a stack of decoders of the same number.

EVEN MORE DETAILED VIEW:



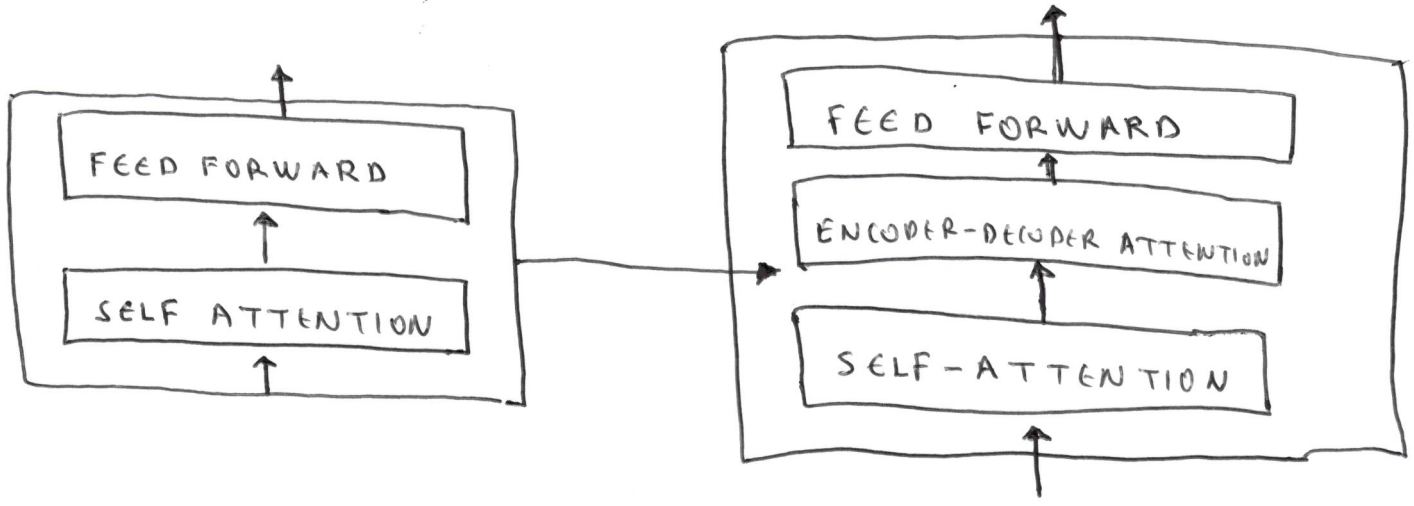
- The encoders are all identical in structure but they do not share weights.
- Each encoder-layer is broken down into two sub-layers:
  - (a): A self-Attention Layer
  - (b): A Feed-Forward Neural Network



- ★ The encoder's inputs first flow through a self-attention layer that help the neural mechanism to look at other instances of the input sequence as it encodes a particular timestep of the input sequence.
- ★ The outputs of the self-attention layer are fed to a feed-forward neural network. The exact same feed-forward network is independently applied to each position.

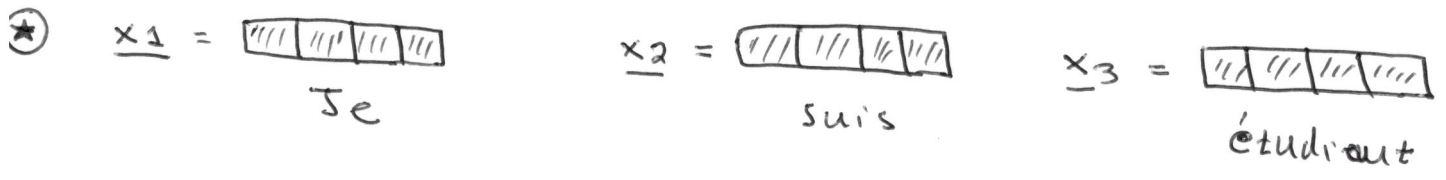
The decoder has both those layers, but between them there exists an attention layer that helps the decoder to focus on relevant parts of the input sequence.

[This attention layer is similar to the attention mechanisms utilized by Sequence to Sequence Models]



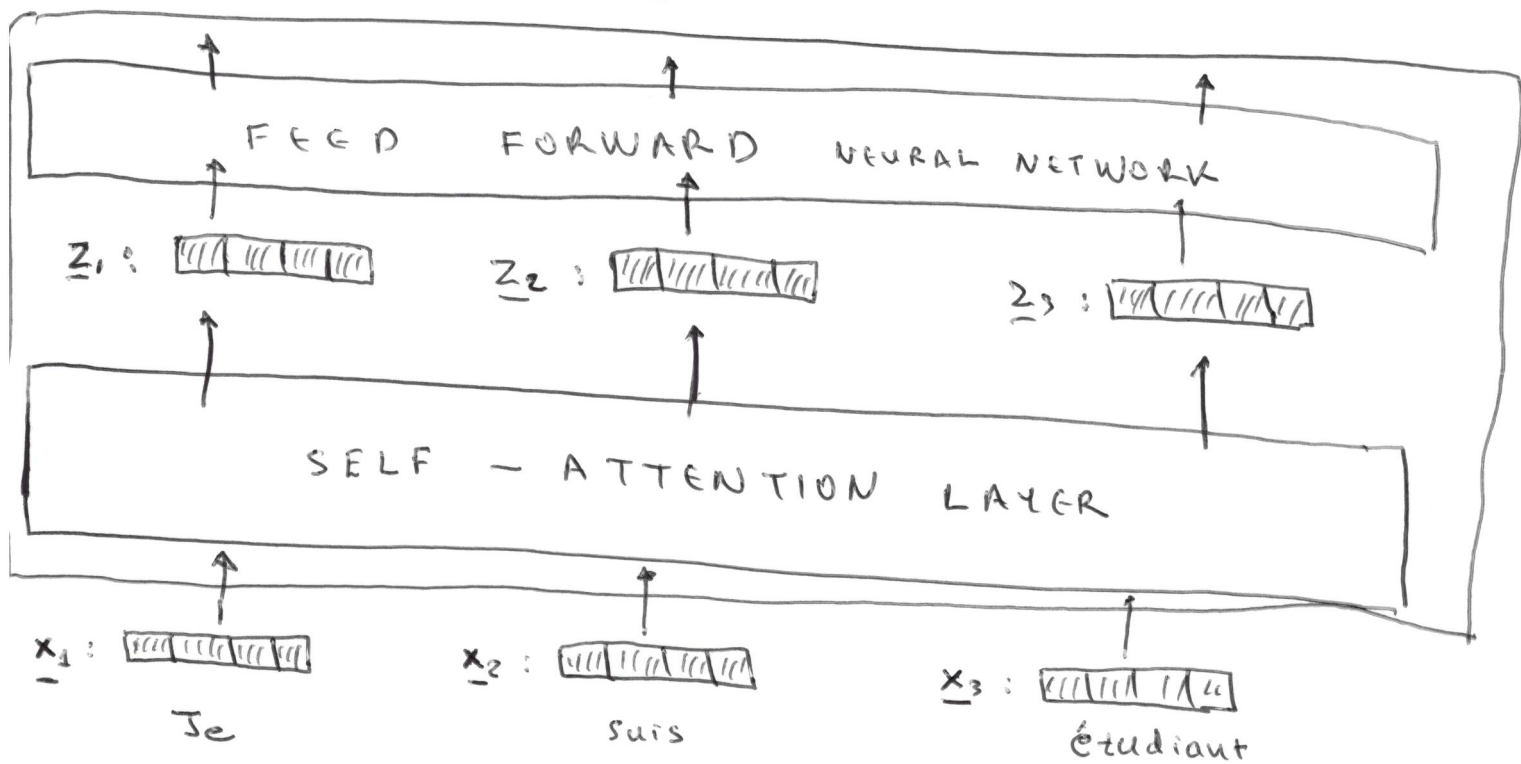
Bringing The Tensors Into The Picture

Let's look at the various vectors/tensors and how they flow between the various components of the transformer model in order to turn the input of a trained model into an output.



Each word is embedded into a vector of size 512.

- ★ The embedding happens only at the bottom-most encoder. The abstraction that is common to all the encoders is that they receive a list of vectors each of size  $512$ .
- ★ The bottom encoder receives the word embeddings while the other encoders receive the output from the directly previous encoder layer.
- ★ The size of the aforementioned list is a hyperparameter that corresponds to the length of the longest input sequence in the training dataset.

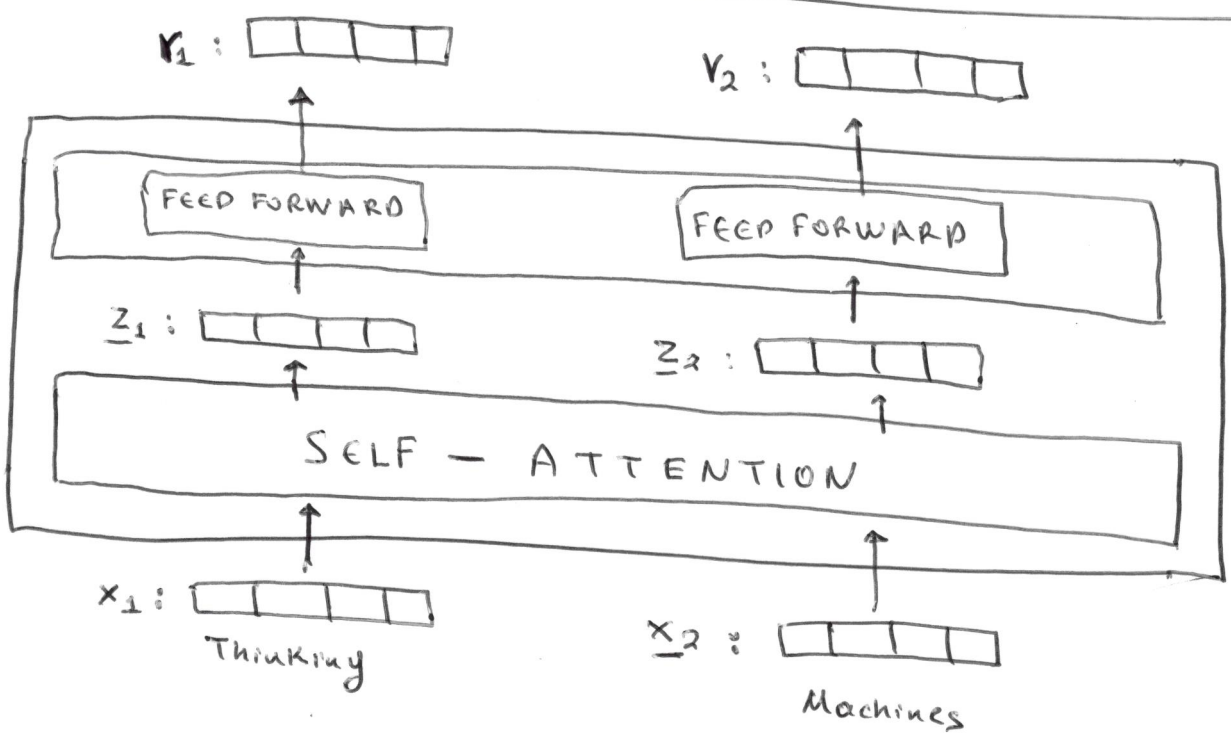


- ★ After embedding the words in the input sequence, each of them flows through each of the two layers of the encoder.

★ One key property of the Transformer, is that each input word (or each input instance at the corresponding position) flows through its own path in the encoder.

• There are dependencies between these paths in the self-attention layer.

★ The feed-forward layer does not have these dependencies and thus, the various paths can be executed in parallel.



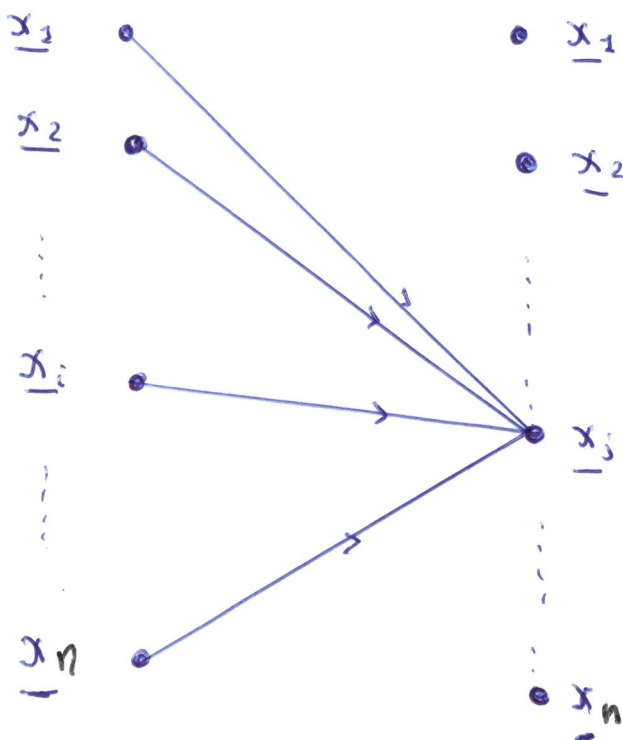
• The word at each position passes through a self-attention process. Then, they each pass through a feed-forward neural network.

[ the exact same, neural network, with each vector flowing through it separately ]

⊛ When the model processes a particular word (or any given feature vector at a specific timestamp), self-attention allows it to look at other positions in the input sequence for clues that can help it to acquire a better encoding for the current token.

• In an analogy with RNNs, one may think that maintaining a particular hidden state allows them to incorporate a representation of the previous/post instances of the input sequence.

• Self-attention is the method employed by the Transformer in order to convey "understanding" of the previous words into the word currently being processed.



⊛ Encoding of the  $j$ -th word may take into consideration the encoding of all other words in the sentence.

Step 1: ▶ The first step in calculating self-attention is to create three vectors from each of the encoder's input vectors (in this case, the embedding of each word).

- ▶ So for each word, a Query vector, a Key vector and a Value vector are being created.
- ▶ These vectors are created by multiplying the embedding by the three matrices that will be acquired during the training process.
- ▶ Notice that these new vectors are smaller in dimension than the original embedding vector. Their dimensionality can be 64 instead of 512. Of course, this is an architectural design choice. They don't have to be smaller. This choice ensures that the computation of the multi-head attention remains constant.

★ Let  $n$  be the number of elements in the sequence such that:

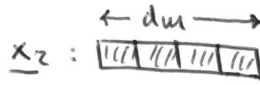
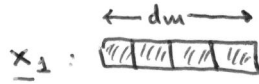
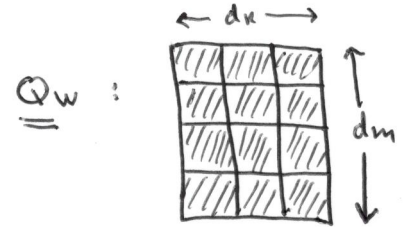
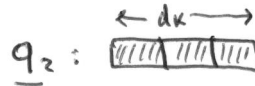
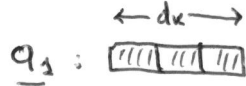
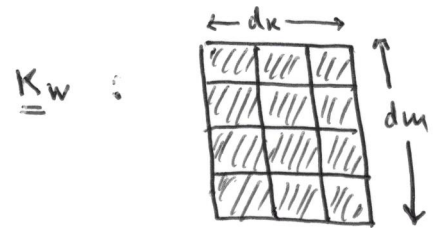
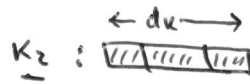
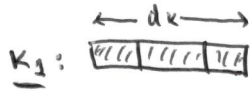
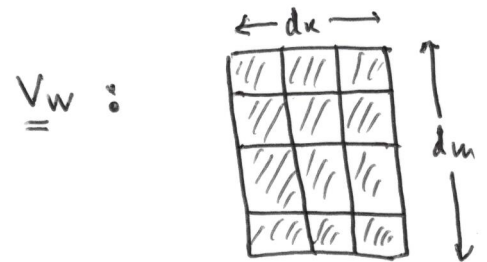
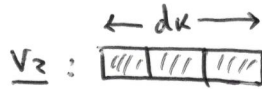
$$X = \{ \underline{x}_1, \underline{x}_2, \dots, \underline{x}_n \} \text{ where } \underline{x}_j \in \mathbb{R}^{d_m} \quad (1)$$

with  $d_m$  denoting the dimensionality of the embedding space.

★ For the illustrated example below, we have that  $n=2$ .

INPUT

Thinking Machines

EMBEDDINGQUERIESKEYSVALUES

★ We may write that:

$$\underline{q}_j = \underline{x}_j \cdot \underline{Q}_w$$

$$[1 \times dk] = [1 \times dm] \cdot [dm \times dk]$$

$$, \forall j \in [n] \quad (2)$$

$$\underline{k}_j = \underline{x}_j \cdot \underline{K}_w$$

$$[1 \times dk] = [1 \times dm] \cdot [dm \times dk]$$

$$, \forall j \in [n] \quad (3)$$

$$\underline{v}_j = \underline{x}_j \cdot \underline{V}_w$$

$$[1 \times dk] = [1 \times dm] \cdot [dm \times dk]$$

$$, \forall j \in [n] \quad (4)$$



Step 2: The second step in calculating self-attention is to calculate a score. Say for example that we are calculating the self-attention score for the first word "Thinking".

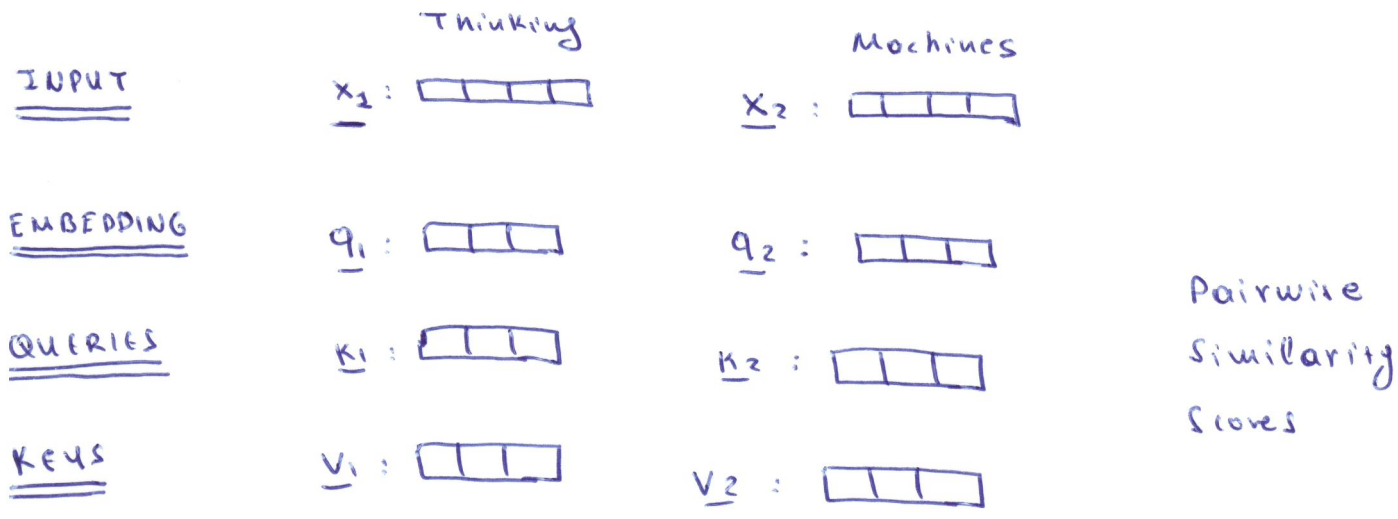
- ▶ Thus, we need to score each word of the input sequence against this word.
- ▶ This score determines how much focus to place on other parts of the input sequence as we encode the feature vector at the current position.

★ We may write that:

$$S_{ij} = q_i \cdot k_j^T \in \mathbb{R}$$

$$\forall i \in [L]$$

$$\forall j \in [L] \quad (5)$$



VALUES

$$S_{11} = q_1 \cdot k_1^T \quad S_{21} = q_2 \cdot k_1^T$$

$$S_{12} = q_1 \cdot k_2^T \quad S_{22} = q_2 \cdot k_2^T$$

$$\underline{S} = \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix}$$

SCORES

★ Let  $\underline{Q} = \begin{bmatrix} \underline{q}_1 \\ \underline{q}_2 \\ \vdots \\ \underline{q}_n \end{bmatrix}, \underline{K} = \begin{bmatrix} \underline{k}_1 \\ \underline{k}_2 \\ \vdots \\ \underline{k}_n \end{bmatrix}, \underline{V} = \begin{bmatrix} \underline{v}_1 \\ \underline{v}_2 \\ \vdots \\ \underline{v}_n \end{bmatrix} \in \mathbb{R}^{n \times dk}$

★  $\underline{S} = \underline{Q} \cdot \underline{K}^T \in \mathbb{R}^{n \times n}$  ★ For this particular example:  $\underline{S} = \begin{bmatrix} \underline{q}_1 \\ \underline{q}_2 \end{bmatrix} \begin{bmatrix} \underline{k}_1^T & \underline{k}_2^T \end{bmatrix} \Rightarrow$

$$\underline{S} = \begin{bmatrix} q_1 \cdot k_1^T & q_1 \cdot k_2^T \\ q_2 \cdot k_1^T & q_2 \cdot k_2^T \end{bmatrix}$$

Steps 3 and 4:  $\odot$  divide scores by  $\sqrt{d_k}$  which leads to more stable gradients.

$\odot$  Apply softmax normalization at each row of the pairwise score matrix  $\underline{S}$ .

$$\hat{\underline{S}} = \text{softmax}(\sqrt{d_k}^{-1} \cdot \underline{S}) \quad (6)$$

$\star$  For this particular example we have that:

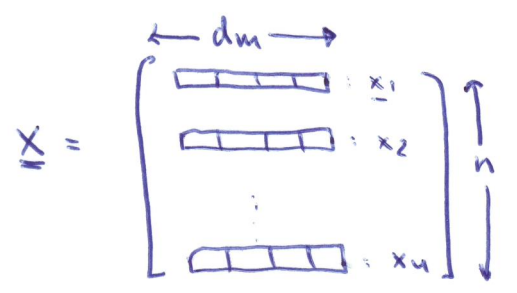
$$\hat{S}_{11} + \hat{S}_{12} = 1 \quad \text{and} \quad \hat{S}_{21} + \hat{S}_{22} = 1$$

Steps 5 and 6: Produce the outputs of the self-attention layer:

$$\underline{z}_j = \sum_{r=1}^n \hat{S}_{jr} \cdot \underline{v}_r, \quad \forall j \in \{1, 2\} \quad (7)$$

where  $\underline{z}_j \in \mathbb{R}^{d_k}$

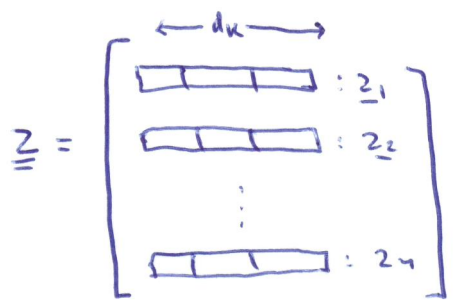
$\star$  Assuming that the sequence  $\mathcal{X}$  of input patterns is stored in a matrix  $\underline{X}$  ( $\underline{x} \in \mathbb{R}^{n \times d_m}$ ) such that:



the output of the attention layer will be given as

$$\underline{z} = \hat{\underline{S}} \cdot \underline{V} \quad (8)$$

such that:



$$[n \times d_k] = [n \times n] \times [n \times d_k]$$

# Matrix calculation of Self-Attention

#11

Letting  $\underline{X} \in \mathbb{R}^{n \times d_m}$  the matrix of the input patterns, we may write that:

$$(I): \quad \underline{Q} = \underline{X} \cdot \underline{Q}_w \in \mathbb{R}^{n \times d_k}$$

$$(II): \quad \underline{K} = \underline{X} \cdot \underline{K}_w \in \mathbb{R}^{n \times d_k}$$

$$(III): \quad \underline{V} = \underline{X} \cdot \underline{V}_w \in \mathbb{R}^{n \times d_k}$$

$$(IV): \quad \hat{S} = \text{softmax} \left( \frac{\underline{Q} \cdot \underline{K}^T}{\sqrt{d_k}} \right) \in \mathbb{R}^{n \times n}$$

$$(V): \quad \underline{Z} = \hat{S} \cdot \underline{V} \in \mathbb{R}^{n \times d_k}$$

Mind that this is actually the output of one attention head. Thus, by considering  $M$  attention heads, we would end up with output vectors of  $d_m$  dimensionality

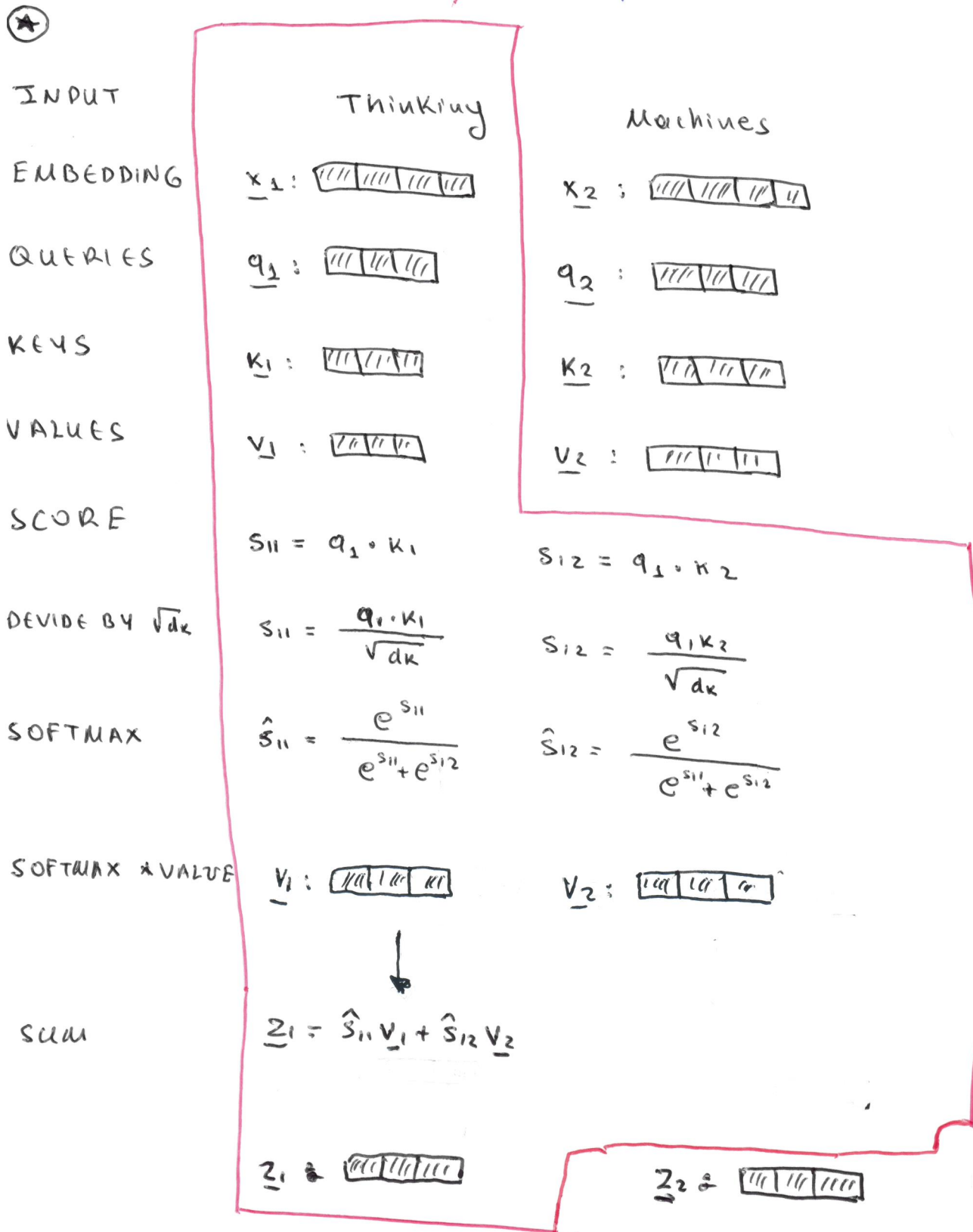
$$M * d_k = d_m$$

by concatenating the outputs of the various heads.

$$\left. \begin{array}{l} M = \frac{d_m}{d_k} \\ \parallel \\ 8 \text{ heads} \end{array} \right\}$$

Multiple heads extend the model's ability to focus on different positions.

Computation block for  $z_1$



★ This figure reveals the attention layer computation for the 1st instance of the input sequence when the number of attention heads is equal to 1

★ In the example above,  $z_1$  contains information about every other encoding, but it could be dominated by the actual word itself.

Multiple heads allow the attention layer to retain multiple "representation subspaces". With multiple attention heads we have multiple sets of Query/Key/Value matrices.

The original Transformer model defined 8 attention heads resulting in 8 sets of each encoder-decoder pair.

Each set of matrices is randomly initialized. After the completion of the training process, each set of matrices is utilized to project the input embeddings (or vectors from lower encoders/decoders) into a different representation subspace.

Let  $m \in \{0, 1, \dots, M-1\}$  indicate a particular attention head associated with weight-matrices  $\underline{Q}_w^{(m)}, \underline{K}_w^{(m)}$  and  $\underline{V}_w^{(m)}$ .

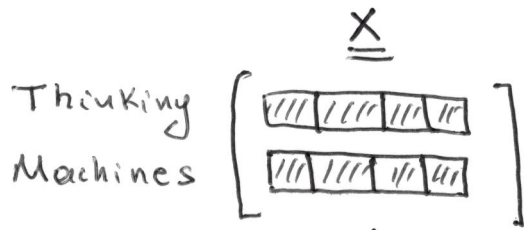
$$\begin{cases} \underline{Q}^{(m)} = \underline{X} \cdot \underline{Q}_w^{(m)} \in \mathbb{R}^{n \times d_k} \\ \underline{K}^{(m)} = \underline{X} \cdot \underline{K}_w^{(m)} \in \mathbb{R}^{n \times d_k} \\ \underline{V}^{(m)} = \underline{X} \cdot \underline{V}_w^{(m)} \in \mathbb{R}^{n \times d_k} \end{cases} \quad \begin{cases} \underline{S}^{(m)} = \frac{\underline{Q}^{(m)} \cdot (\underline{K}^{(m)})^T}{\sqrt{d_k}} \in \mathbb{R}^{n \times n} \\ \hat{\underline{S}}^{(m)} = \text{softmax}(\underline{S}^{(m)}) \in \mathbb{R}^{n \times n} \end{cases}$$

$$\underline{Z}^{(m)} = \hat{\underline{S}}^{(m)} \cdot \underline{V}^{(m)} \in \mathbb{R}^{n \times d_k}, \quad 0 \leq m \leq M-1$$

FORM THE FINAL CONCATENATED OUTPUT:

$$\hat{\underline{Z}} = [ \underline{Z}^{(0)}, \underline{Z}^{(1)}, \dots, \underline{Z}^{(M-1)} ] \in \mathbb{R}^{n \times d_m}$$

Schematic Representation of the FcuzP Computation Step



calculating attention separately in 8 different attention heads.

ATTENTION HEAD #0

ATTENTION HEAD #1

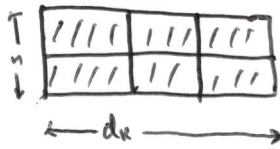
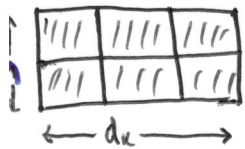
...

ATTENTION HEAD #7

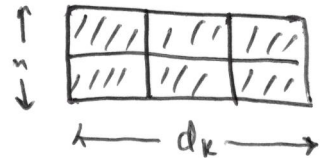
Z<sub>0</sub>

Z<sub>1</sub>

Z<sub>7</sub>



...



Applying the Feed-Forward Neural Network Requires:

(A): since the feed-forward neural network is not expecting 8 matrices (it's expecting a single matrix where a vector is assigned to each word), we need to concatenate the 8 individual heads into a single matrix.

$$\underline{Z} = \left[ \begin{array}{c} \left[ \begin{array}{c} \text{[grid]} \\ \text{[grid]} \end{array} \right] \left[ \begin{array}{c} \text{[grid]} \\ \text{[grid]} \end{array} \right] \dots \left[ \begin{array}{c} \text{[grid]} \\ \text{[grid]} \end{array} \right] \end{array} \right] \begin{array}{l} \uparrow n \\ \downarrow \end{array}$$

$n \times dk = dm$

(B): multiply the concatenated matrix with an additional weight-matrix W<sub>0</sub> such that:

$$\underline{W}_0 \in \mathbb{R}^{dm \times dm}$$

$$\underline{R} = \underline{Z} \cdot \underline{W}_0$$

$[n \times dm] \quad [n \times dm] \quad [dm \times dm]$

# Schematic Representation of the Encoder Layer

