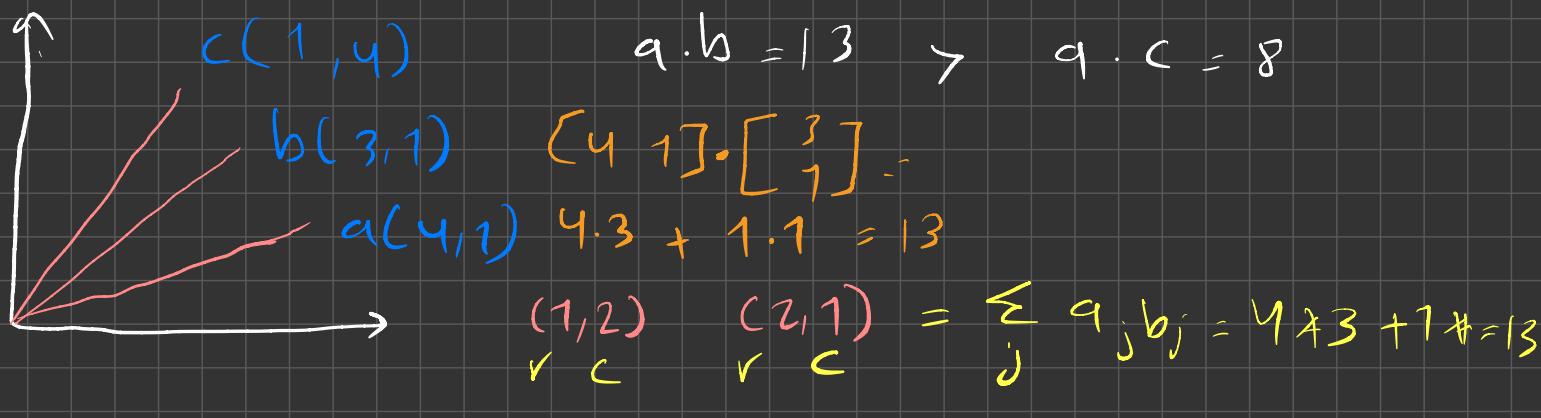


* Dot product :-



* Word Embedding :-

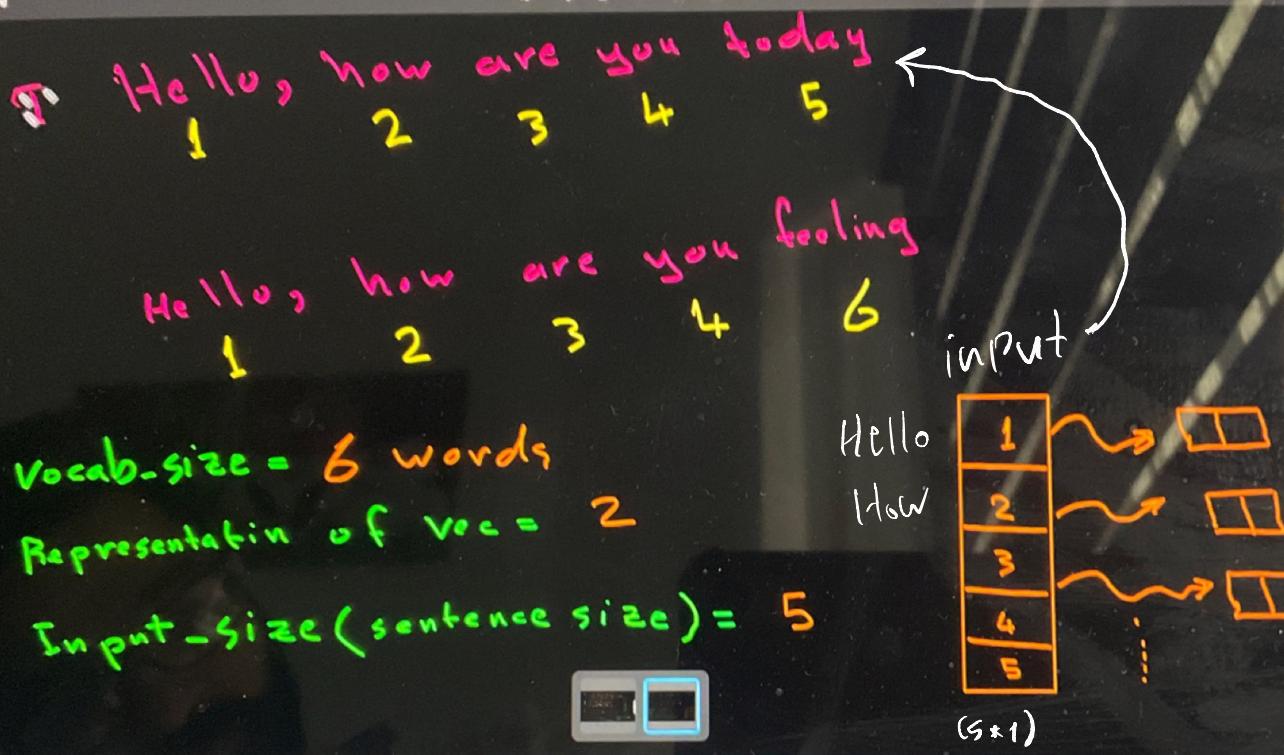
	gender	royal	age	food	
man	-1	0.07	0.03	0.04	...
woman	1	0.02	0.02	0.07	...
king	-0.95	0.93	0.7	0.02	...
queen	0.95	0.97	0.69	0.01	...
;	;	;	;	;	;

Annotations:

- A vertical bracket on the left indicates the matrix has $10K$ rows.
- A horizontal bracket at the bottom indicates the matrix has 300 columns.
- A yellow bracket labeled D spans the entire matrix.
- An arrow points from the word "man" to its corresponding row in the matrix.
- An arrow points from the matrix to a small box containing the vector representation of "man".

$$E_{(10K, 300)}$$

Each word shape = $(1, 300)$



Tokenization

"Where can I find a pizzeria?"

Word-Level Tokenization

[where, can, i, find, a, pizzeria]

- + Intuitive.
- Doesn't handle OOV words including new words, slang, play on words, misspellings, etc.
- Huge vocabulary for large corpora; especially for languages with rich morphology (e.g. Hungarian).
- Handling punctuation is challenging (e.g. "don't" vs "N.Y.C").

Char-Level Tokenization

[a, c, d, e, f, h, i, n, p, r, w, z]

- + Small memory footprint.
- + Handles OOV words.
- Needs to go over all characters and learn a particular sequence for a given word.
- Loss of performance.

Subword Tokenization

"where can I find a pizzeria?"

where, can, I, find, a, pi, zz, eria

listeria

[list, eria]

Subword tokenization has a better chance of handling OOV words while reducing vocabulary size and maintaining performance.

Byte Pair Encoding (BPE)

"she sells seashells by the seashore"

Corpus

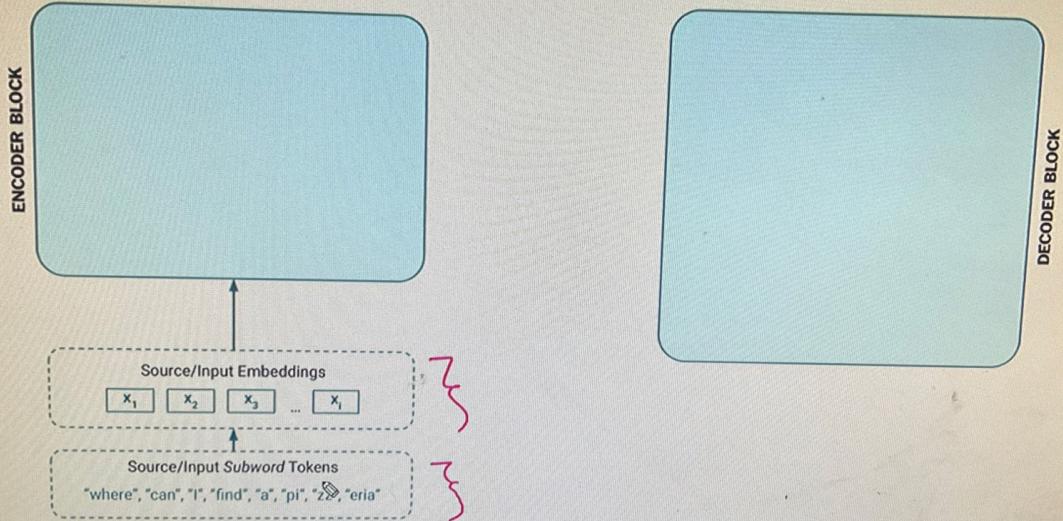
she_
sells_
seashells_
by_
the_
seashore_

Vocabulary

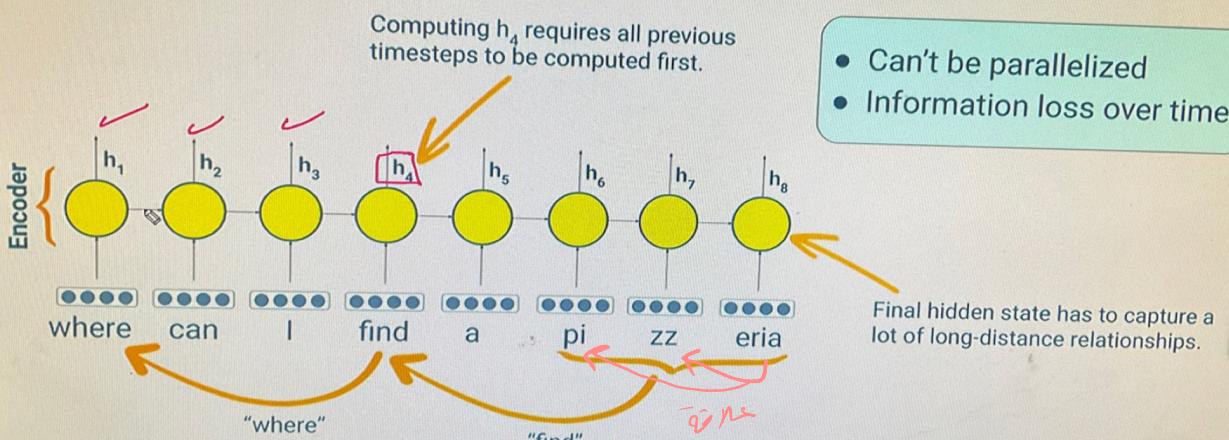
-, a, b, e, h, l, o, r, s, t, y
-, a, b, e, h, l, o, r, s, t, y, sh
-, a, b, e, h, l, o, r, s, t, y, sh, he
-, a, b, e, h, l, o, r, s, t, y, sh, he, e_
-, a, b, e, h, l, o, r, s, t, y, sh, he, e_, se
-, a, b, e, h, l, o, r, s, t, y, sh, he, e_, se, she
.
.
.

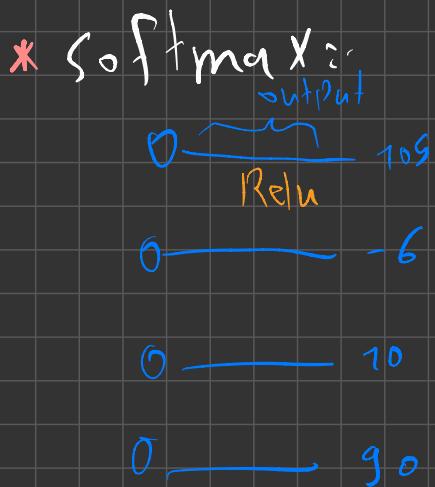
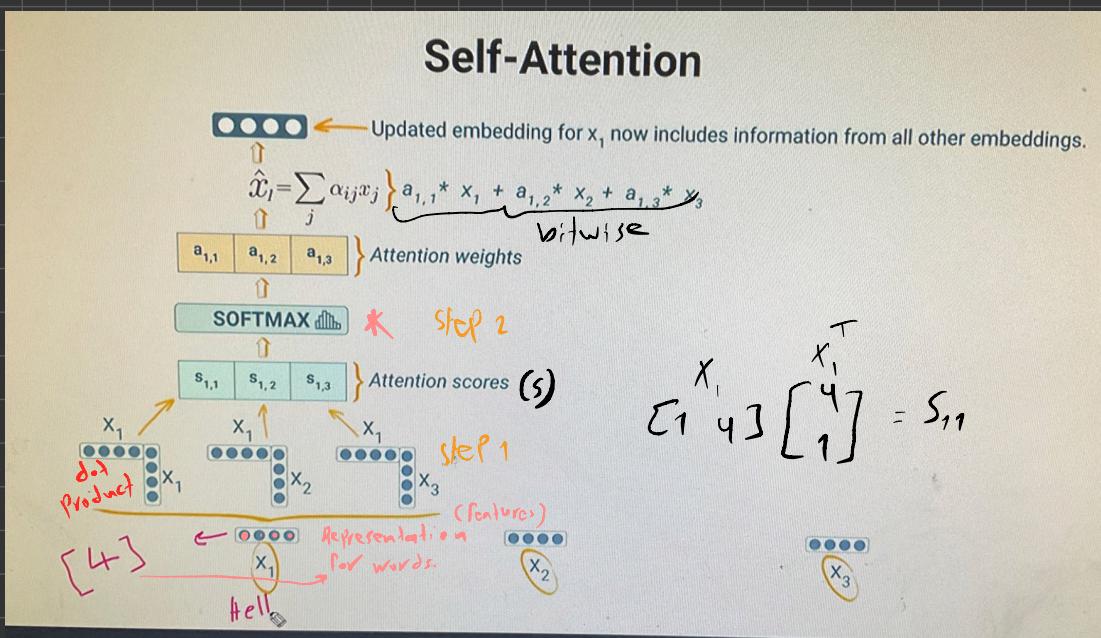
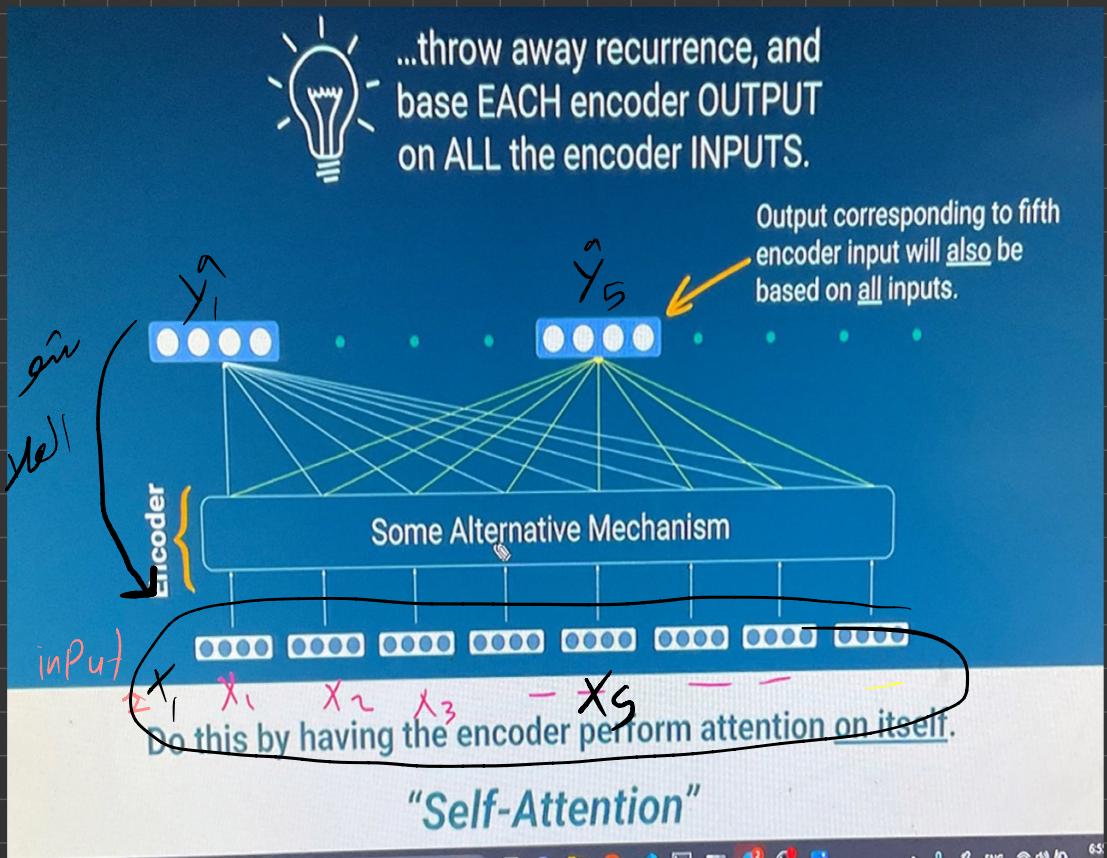
continue until N merges are performed.

Transformers



Why Transformers





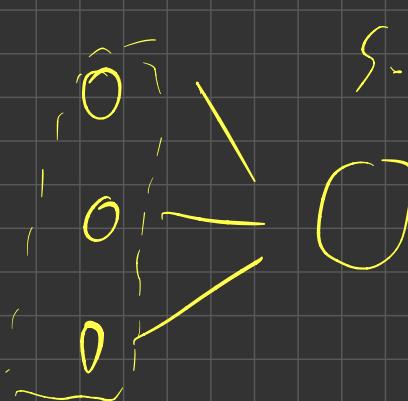
$$S = \begin{bmatrix} S_1 \\ S_2 \\ S_3 \\ \vdots \\ S_n \end{bmatrix}$$

Score

$S \in [-\infty, +\infty]$

Goal: $s_{src} \rightarrow p_{probabilities}$

Need $\left\{ \begin{array}{l} p_i \in \{0, 1\} \\ \sum p_i = 1 \end{array} \right.$



$$p_i = \frac{s_i}{\sum_j s_j}$$

↑ when we use this formula

$$\begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1/3 \\ 2/3 \end{bmatrix}, \begin{bmatrix} 100 \\ 101 \\ 102 \end{bmatrix} = \begin{bmatrix} 0,33 \\ 0,33 \\ 0,33 \end{bmatrix}$$

$$p_i = \frac{e^{s_i}}{\sum_j e^{s_j}}$$

$$\hat{p}_i = \frac{e^{s_i + 100}}{\sum_j e^{s_j + 100}} = \frac{e^{s_i} \times e^c}{\sum_j e^{s_j} \times e^c}$$

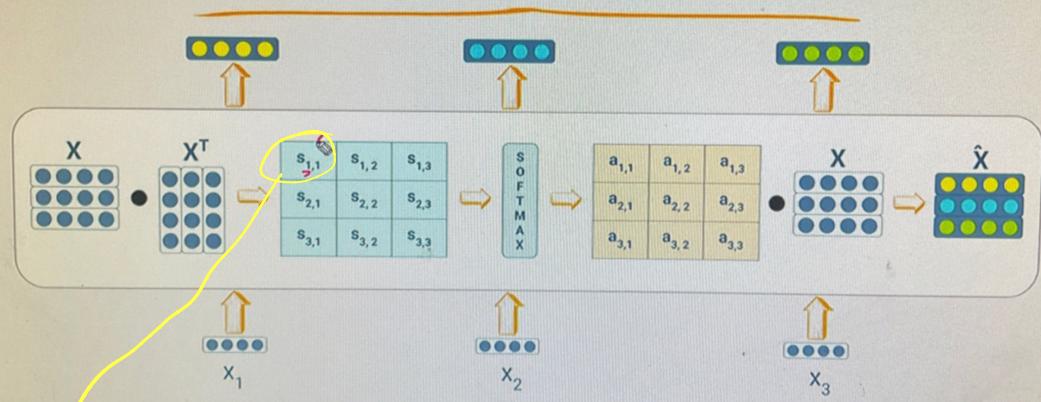
constant of prob.

a negative numbers that output from node

if go inside softmax will be this number positive.

Self-Attention

Each updated embedding now includes information from all other embeddings.

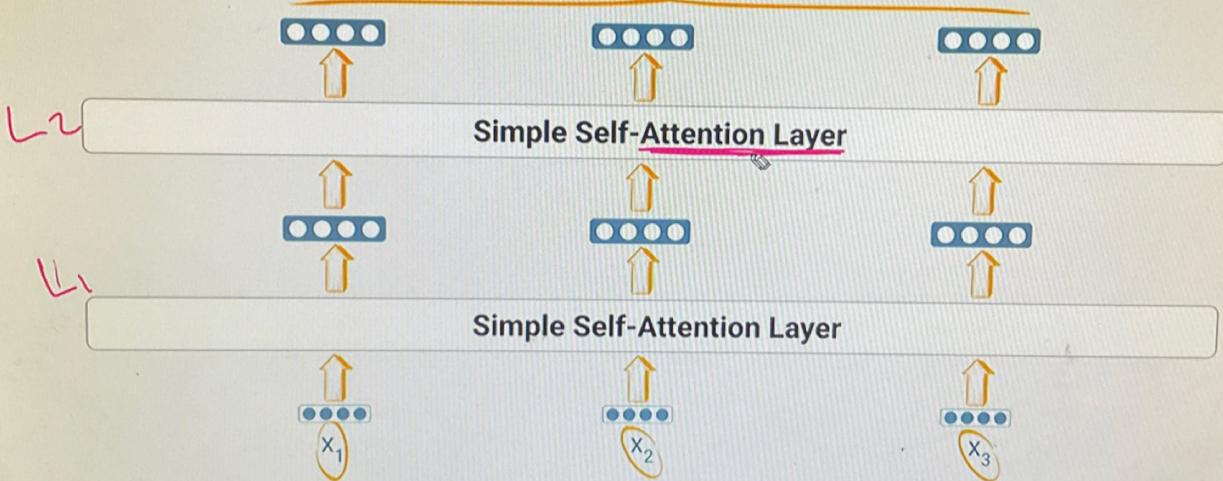


$s_{1,1}$
position
 x_1
input

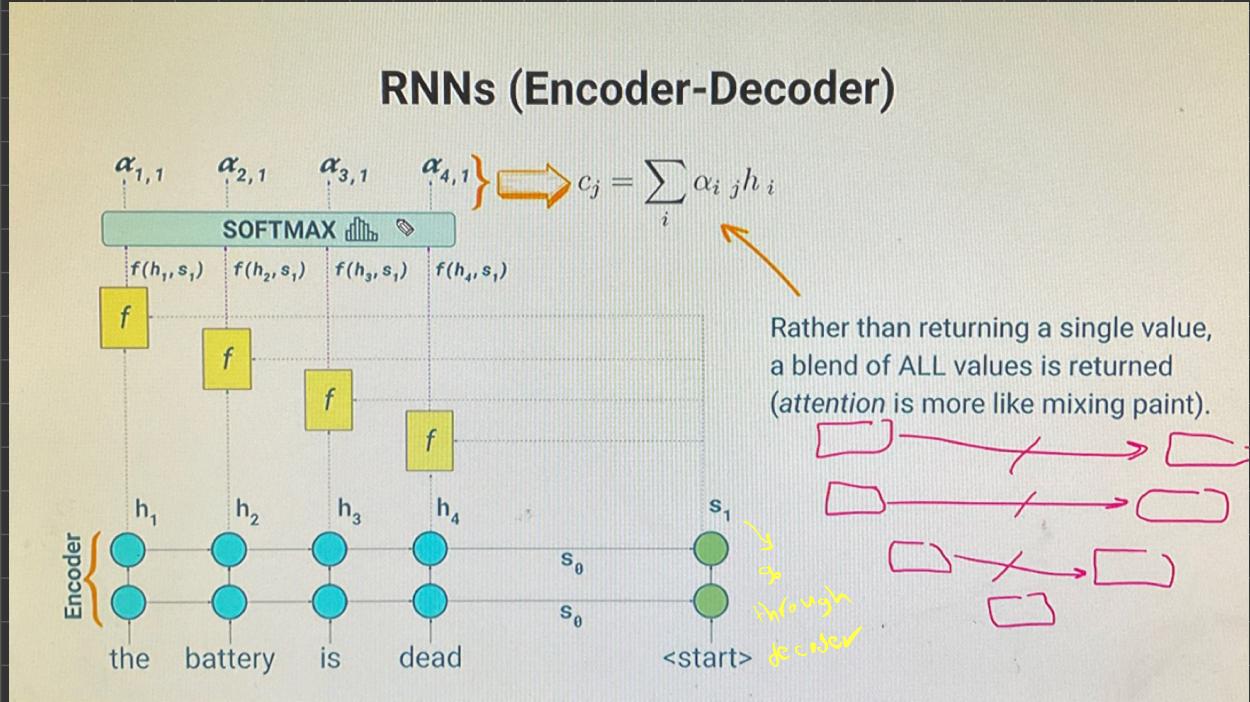
6-Block

Self-Attention

Each updated embedding now includes information from all other embeddings.

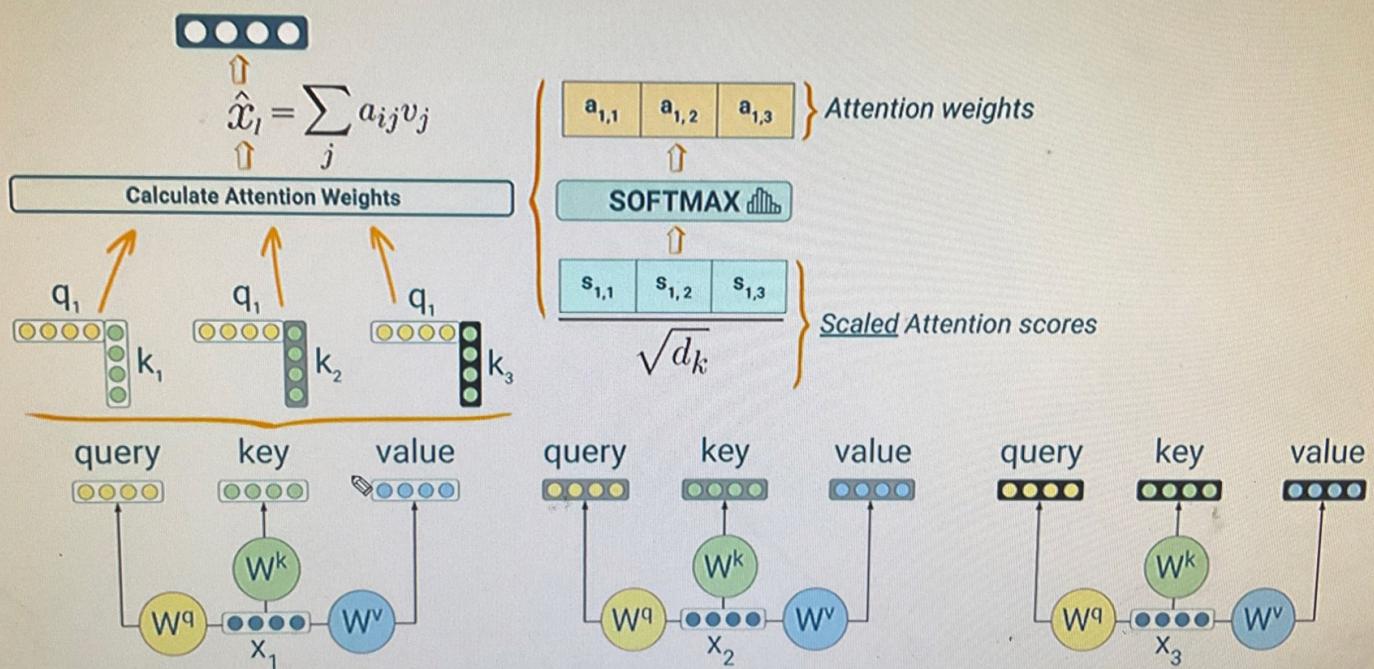


RNNs (Encoder-Decoder)



* the parameters that are need to learn during training
to get good performance

Scale Dot Product



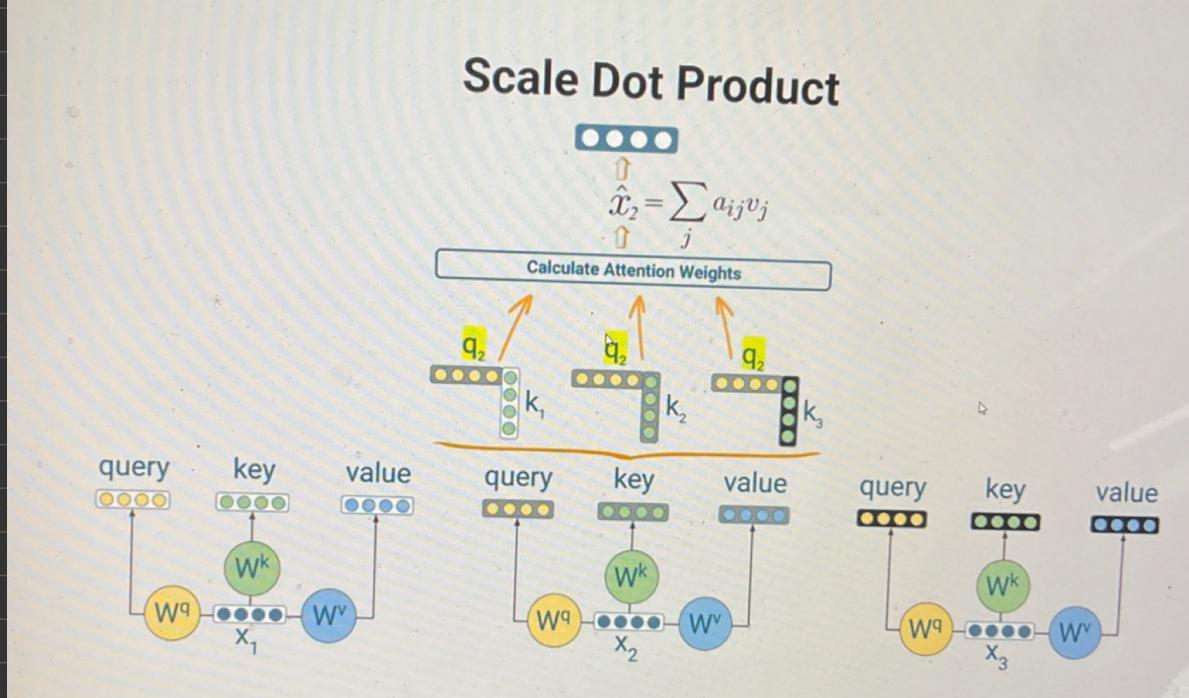
example.. find transformers videos.

* **query**: what you want. (searching)

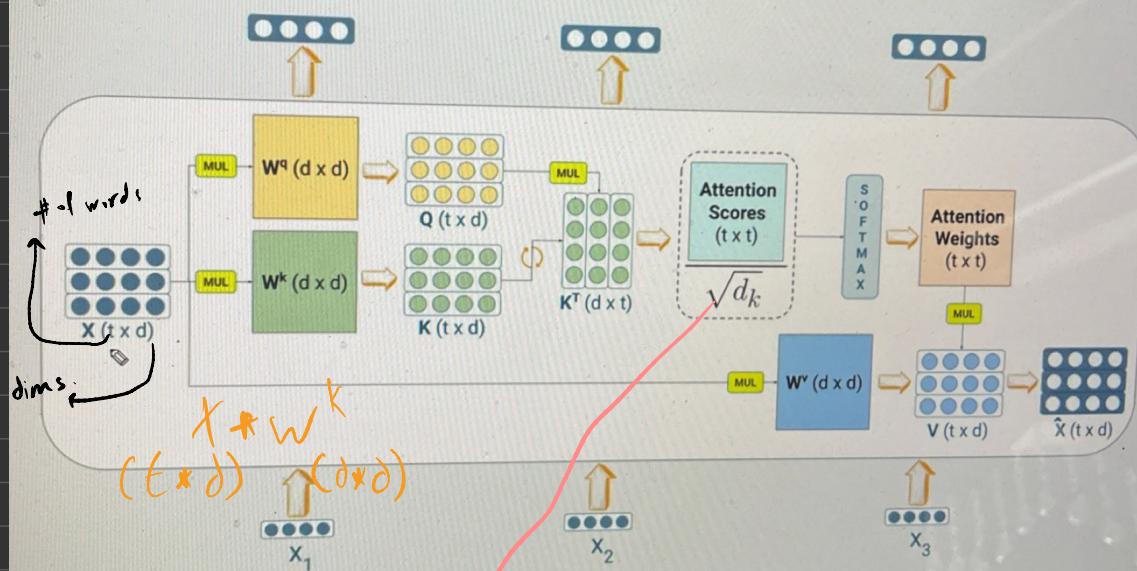
* **key**: person who doing it on youtube

* **value**: where is stored.

Scale Dot Product



Scale Dot Product Self-Attention



Attention

$$\text{Scores} = Q \cdot K$$

(Ex: $(t \times d) \times (d \times d)$)

example

$$(5, 5, 5) = \sqrt{5^2 + 5^2 + 5^2}$$

$$\Rightarrow \sqrt{3 \times 5^2}$$

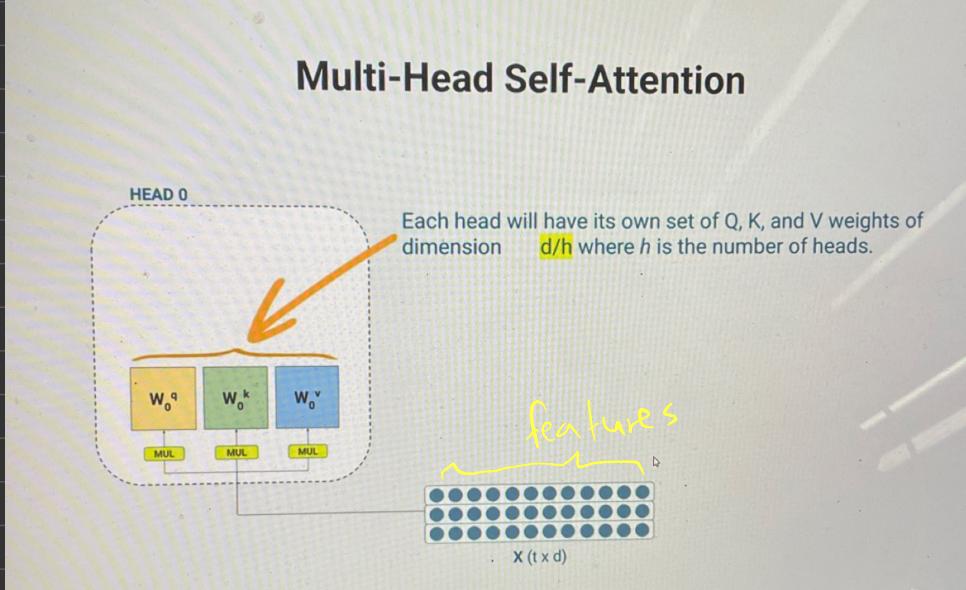
$$= \sqrt{3} \times 5^2$$

d_k : dimension of key.

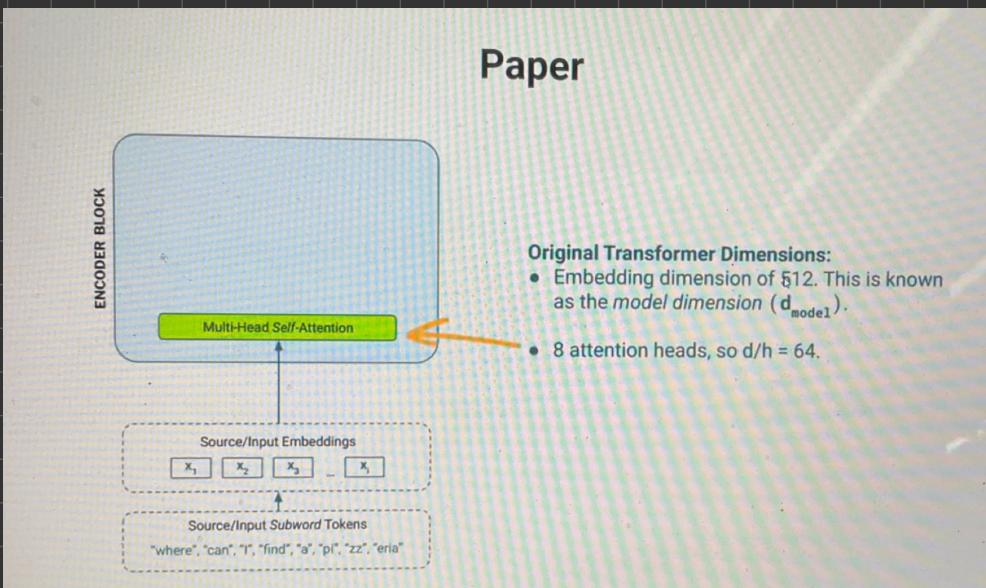
$$\Rightarrow \text{Attention } (Q, K, V) = \text{softmax} \left(\frac{Q \cdot K^T}{\sqrt{d_k}} \right) \cdot V$$

buli jis mao
عوچ چو چو لدی

Multi-Head Self-Attention



Paper



* Positional information :-

- index of the words .

- to know the sequences of words .

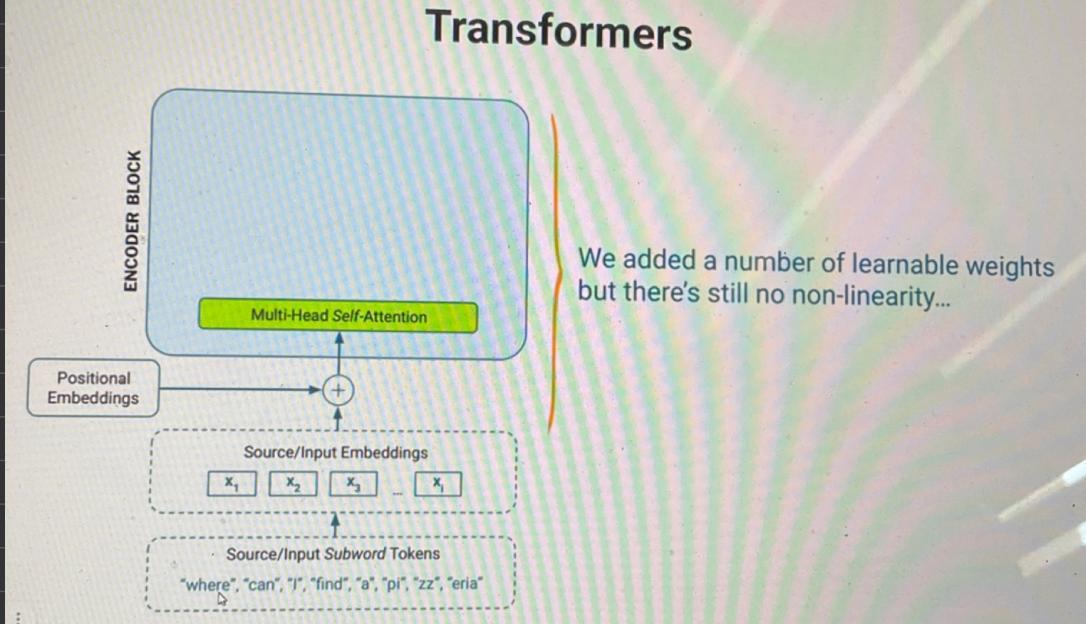
$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$$

embedding size
even

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

odd

Transformers



* in Transformers we can not use Batch normalization



Batch normalization is not typically used in transformers because it can disrupt the self-attention mechanism and hinder the model's performance. Transformers rely on the attention mechanism to capture dependencies between different positions in the sequence, and batch normalization can introduce dependencies between positions by normalizing each position independently. This breaks the ability of the self-attention mechanism to attend to different positions adaptively.

Additionally, transformers already have layer normalization, which is a similar technique to batch normalization but operates independently for each position in the sequence. Layer normalization helps stabilize the training process and improves the generalization of the model. It is used in transformers to normalize the input embeddings and the outputs of each transformer layer.

In summary, batch normalization is not used in transformers due to its interference with the self-attention mechanism, while layer normalization is employed to achieve the necessary normalization within transformers.