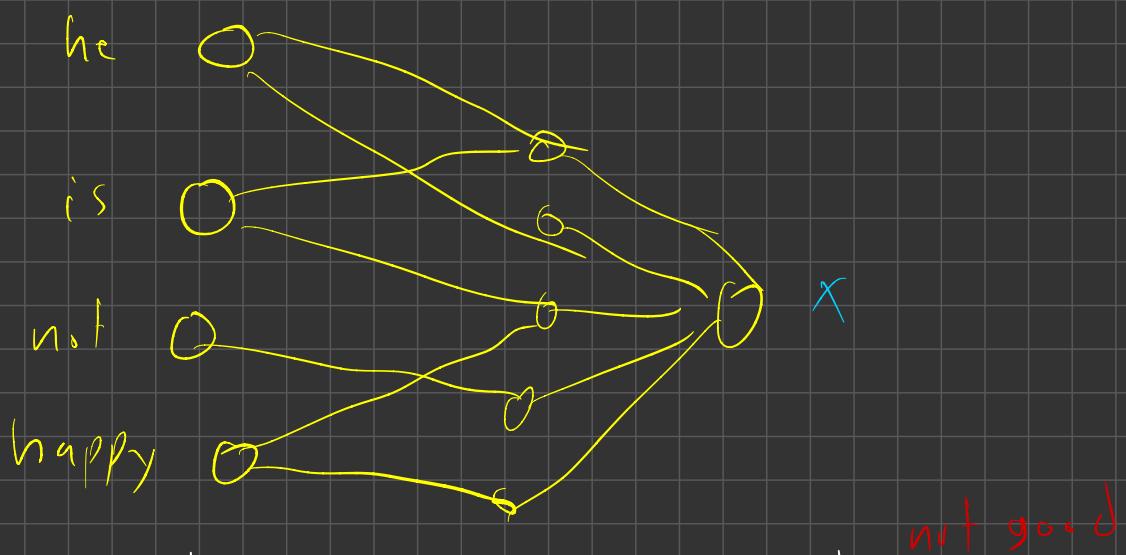


• RNN: example: we can do this method (neural network)



- First we need to do for each word (One hot encoding).  
if we do that we have too much vectors.

example:

he	is	not	happ	—	maybe more words
1	0	0	0	—	—
0	1	0	0	—	—
0	0	1	0	—	—
0	0	0	1	—	—

⇒ Payable to high computational power.

• RNNs have short memory.

- remember just current things and forget previous.

ex - there is 30 words. Only the last 3 word are remembered.

Soul: LSTM

## \* Word Embedding:-

The main advantage of using word Embedding is that it allows words of similar context to be grouped together and dissimilar words are positioned far away from each other. With help of an Embedding matrix.

- The similarity of 2 words can be found  $\Rightarrow$  (Cosine Similarity).

Embedding Matrix		King	Queen	Princess	Boy
Royal		0,99	0,99	0,99	0,01
Male		0,99	0,02	0,01	0,98
Female	→	0,02	0,99	0,99	0,01
Age		0,7	0,6	0,1	0,2

## \* Pre-processing for Embedding matrix:-

There are a lot of algorithms for this:

1. One Hot Encoding
2. Term Frequency-Inverse Document Frequency
3. Tokenization (Text to Sequence)

- Tokenization:- Assigning a number for each unique word.

example:-

[what is your name", "how are you", "where are you"]

After tokenize:-

what:1 , is:2 , your:3 , name:4 , how:5 , are:6 , you:7 , where:8

- Tokenized of first sentence : [1, 2, 3, 4]
- " " second : [5, 6, 7]
- " " third : [8, 6, 7]

**Word2Vec**: is a neural network that creates Word Embeddings (a vector that represents a word in numeric form). To present all the words in a database of document. A word embedding will capture many different parts of word including its semantic, syntactical similarity and relation to other words.

- We have 50,000 words with using word2vec we can represent 30 word instead of 50,000 words.

### \*Cosine Similarity:

- $\cos(0^\circ) = 1$   $\therefore$  Correlation is high.
- $\cos(90^\circ) = 0$   $\therefore$  there is no correlation.
- $\cos(180^\circ) = -1$   $\therefore$  correlation is opposite.

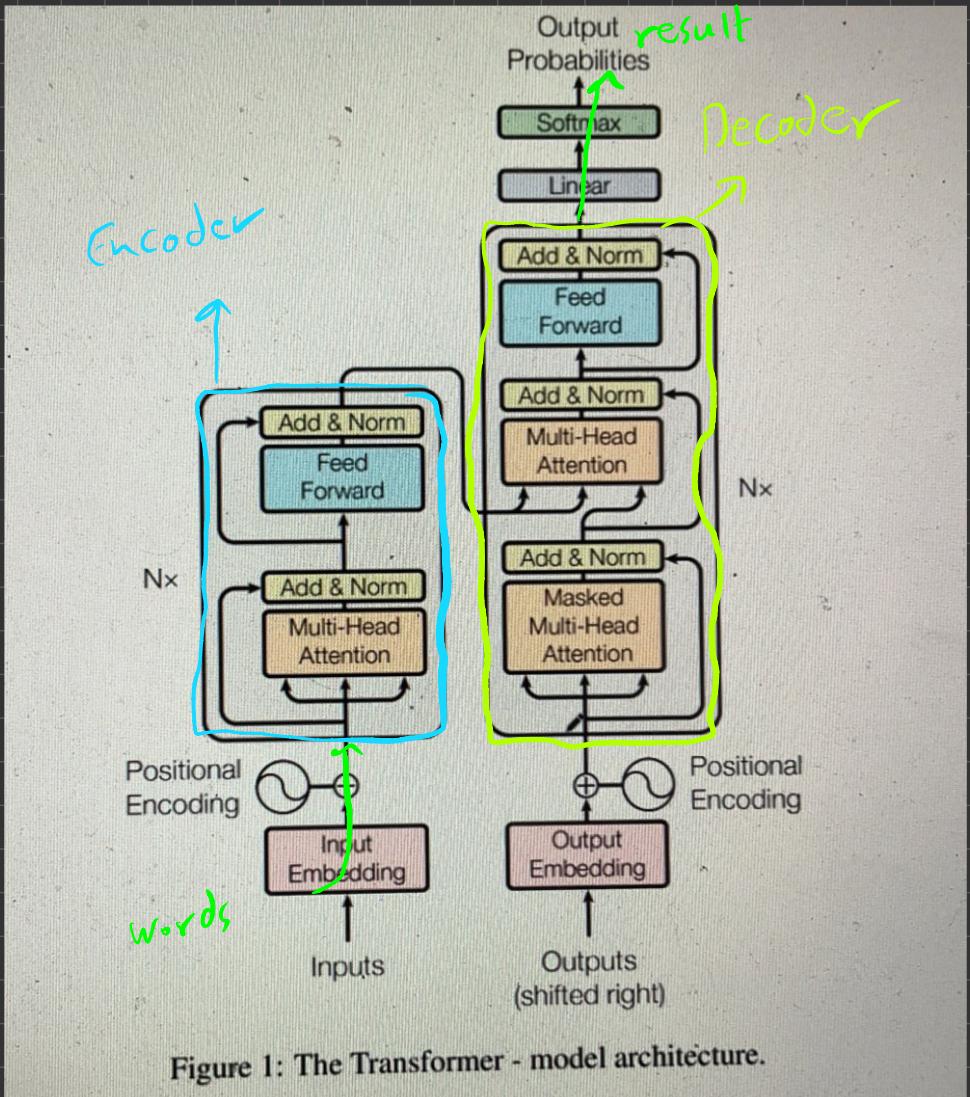


Figure 1: The Transformer - model architecture.

1. input embedding:- (embedding layers).  
returns sentences to words using tokenization.  
Then to vectors.

2. Positional Encoding. (index of vector)

$$PE_i = \sin\left(\frac{pos}{1000^{\frac{i}{d}}}\right) \quad \text{even}$$

embedding size

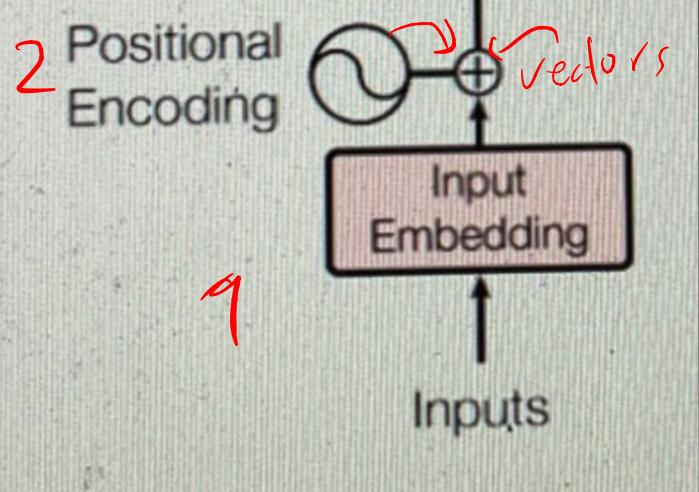
$$PE_{i+1} = \cos\left(\frac{pos}{1000^{\frac{i+1}{d}}}\right) \quad \text{odd}$$

• Uniqueness for position vectors

1. Periodic function. ( $\sin, \cos$ )

2. different frequencies (i). (shift for words)

3.  $D_{\text{model}}(10,000) \Rightarrow$  large periodicity



### 3. Self Attention: (multi-head Attention)

$\text{dict} = \{\text{"a": } \underbrace{\begin{bmatrix} s \\ r \end{bmatrix}, \text{"b": } \begin{bmatrix} 1 \\ 6 \end{bmatrix}\}$

dict["a"] = 15

Query

• dot product :-  $\begin{bmatrix} a \\ b \end{bmatrix} \cdot \begin{bmatrix} b^T \end{bmatrix}$

$$\begin{bmatrix} 1 \times 10 \\ 1 \times 10 \end{bmatrix} \left[ \begin{array}{c} \\ b^T \end{array} \right]_{10 \times 1} = \begin{bmatrix} 1 \times 1 \\ 1 \times 1 \end{bmatrix}$$

$$a \cdot b^T$$

\* cos similarity :-

$$\cos(a, b) = \frac{a \cdot b}{\|a\| \|b\|} \Rightarrow \text{Sim}(a, b)$$

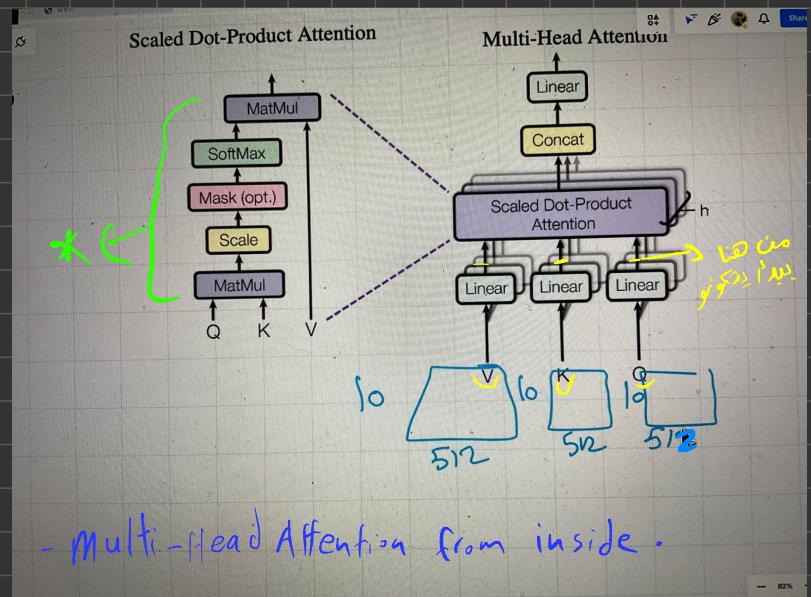
example:-  $a \cdot b = 50$      $a \cdot c = 70$  }  $\rightarrow a$  closer to  $b$  because  $a \cdot b = 50 \uparrow$

$$\Rightarrow \underbrace{\frac{Q \cdot K^T}{\text{Scale}}}_{10 \times 10} = \underbrace{\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}}_{10 \times 10} \rightarrow \text{Value by row} \rightarrow \text{row } q^T$$

$$\Rightarrow \boxed{\frac{Q \cdot K}{\sqrt{d_K}}} \quad (S, S, S)_3 = \sqrt{S^2 + S^2 + S^2} = \sqrt{3} * S = \sqrt{3} * S \quad \boxed{\begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \rightarrow \text{Attention filter}}$$

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

\* answer



matrix  $10 \times 512$

جاء

\* الارقام موجودة  
بالورقة المحدث

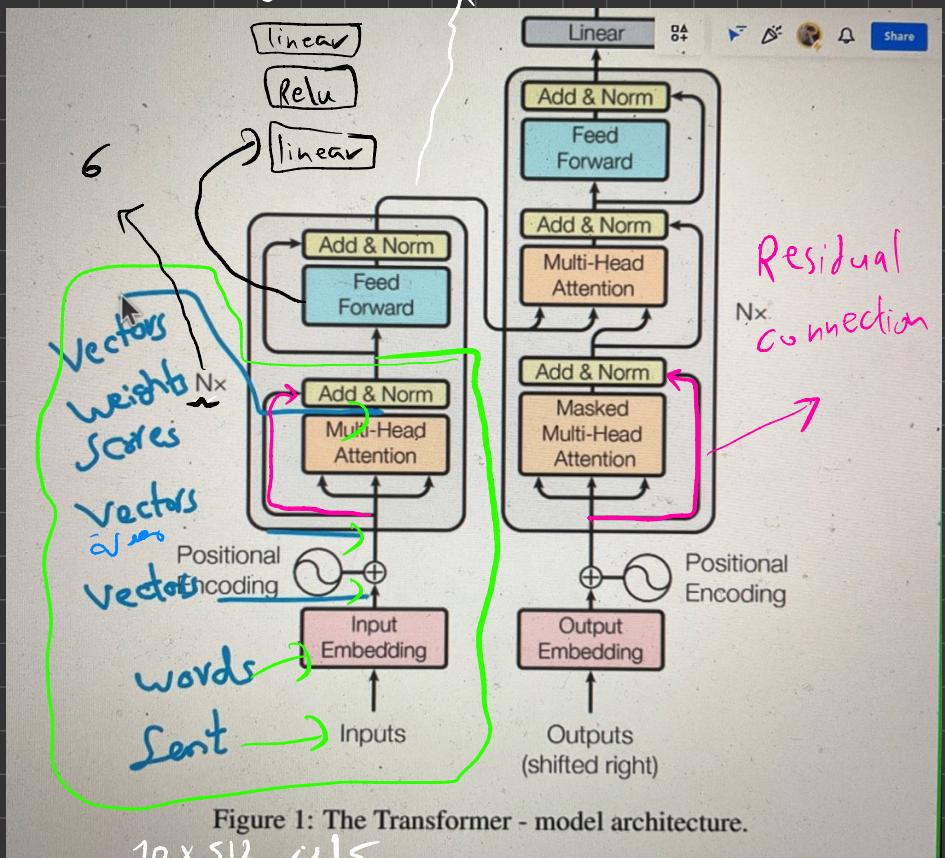
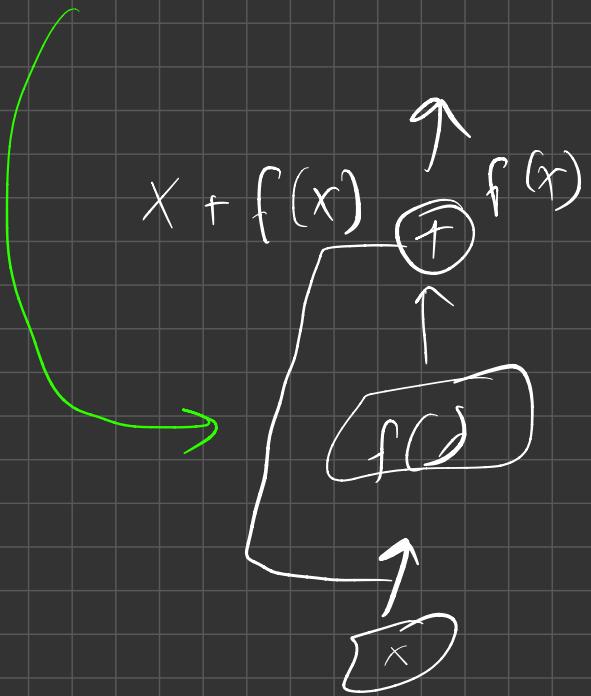


Figure 1: The Transformer - model architecture.

10 x 512

كان

## \* Residual Connection :-



دخل RNN  
دخل Tran.

\* Feed Forward :- performed more encoding for values.

Linear

ReLU

Linear

في عبارة

