

* Transformer Network

Although LSTMs overcome the problem of RNNs, they're still sequential models increased complexity

RNN

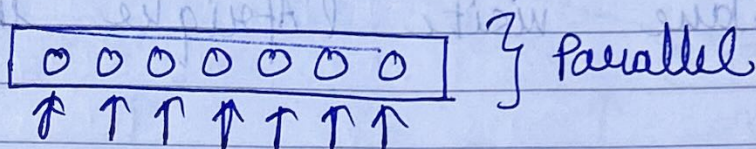
GRU

LSTM

• Attention + CNN processing

- Self Attention

- Multi Head Attention



* Self Attention:

$A(q, k, v)$ = attention-based vector representation of a word
 ↳ calculate for each word $A^{<1>}, \dots, A^{<6>}$

$$A(q, k, v) = \frac{\sum_i \exp(q \cdot k^{<i>})}{\sum_j \exp(q \cdot k^{<j>})} v^{<i>}$$

Main difference from RNN attention:

$q^{< >}$, $k^{< >}$, $v^{< >}$

query

key

value

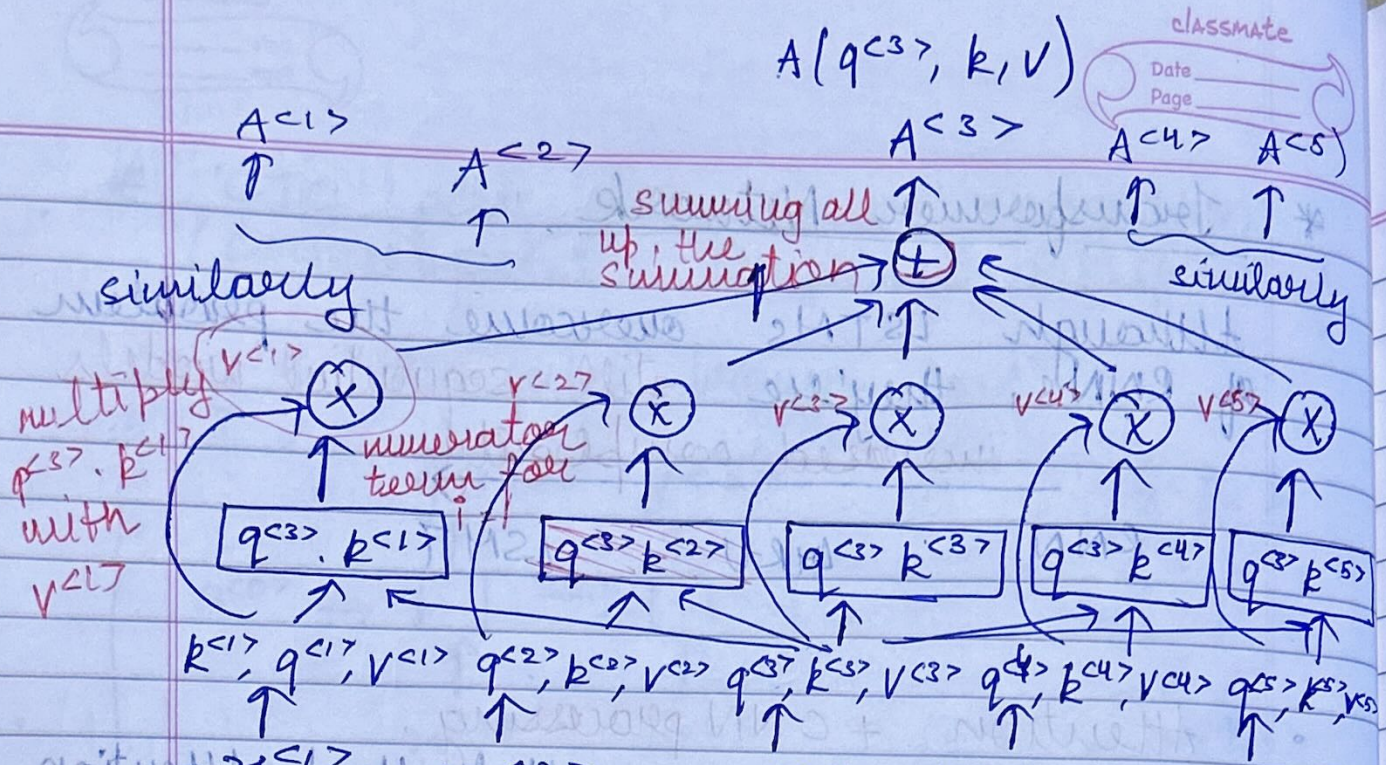
question you get to ask.

most appropriate answer to query, words within

allows representation to plug in, how words should be represented

Eg: what's happening there?
 (Africa)

query, words within $A^{<3>}$ (here).



Jane visite l'Afrique in septembre.

Query (Q)	Key (K)	Value (V)
$q^{(1)}$	$k^{(1)}$	$v^{(1)}$
$q^{(2)}$	$k^{(2)}$	$v^{(2)}$
$q^{(3)}$	$k^{(3)}$	$v^{(3)}$
$q^{(4)}$	$k^{(4)}$	$v^{(4)}$
$q^{(5)}$	$k^{(5)}$	$v^{(5)}$

$$\begin{aligned}
 q^{(3)} &= W^Q \times x^{(3)} \\
 k^{(3)} &= W^K \times x^{(3)} \\
 v^{(3)} &= W^V \times x^{(3)}
 \end{aligned}$$

Scaled Dot Product Attention

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

(compressed/vectorised way)
another way of writing the formula on the previous page
term to scale the dot product.

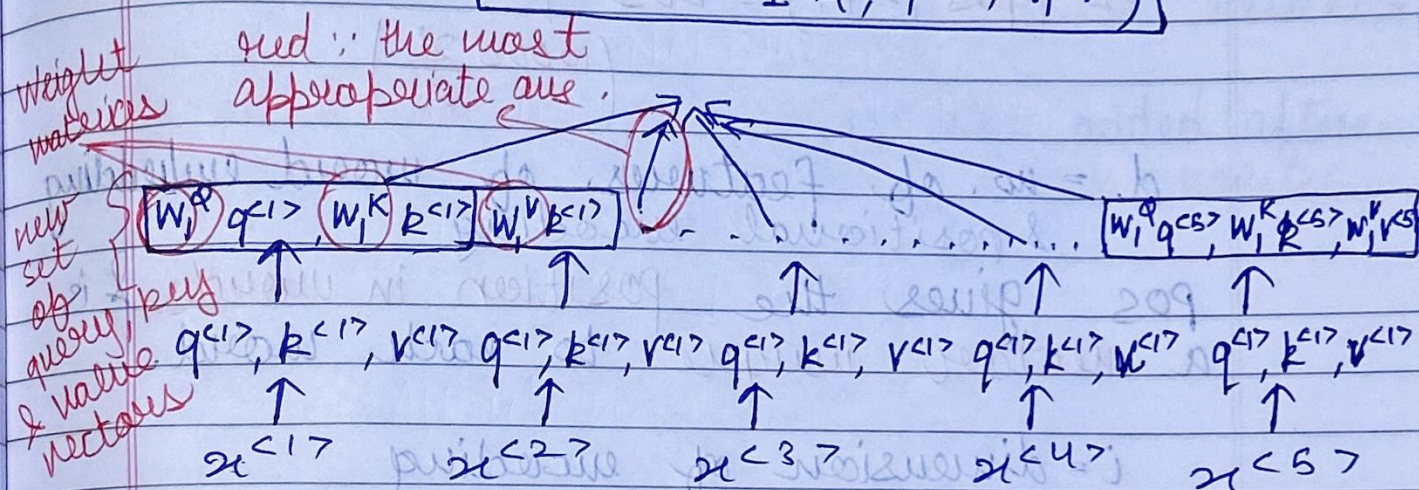
(for loop (for understanding) over self attention)

* **Multi-Head Attention:**
(multiple questions for every word)

Multi Head $(\Phi, K, V) = \text{concat.}(\text{head}_1, \text{head}_2, \dots, \text{head}_h) W_o$
 $\text{head}_i = \text{Attention}(W_i^\Phi Q, W_i^K K, W_i^V V)$
 $h = \text{no. of heads}$

→ All the **heads** can be computed **parallelly**.
 Multihead (Φ, K, V)

Attention $(W_1^\Phi Q, W_1^K K, W_1^V V)$



Intuition: $W_1^\Phi, W_1^K, W_1^V \Rightarrow$ help to ask & answer the question what's happening.

W_1^Φ, W_1^K, W_1^V were the weights for one question (or one self attention example), similarly we have W_2^Φ, W_2^K, W_2^V to answer when, then W_3^Φ, W_3^K, W_3^V to answer who, and so on.

Also, note that $q^{(1)}, k^{(1)}, v^{(1)}$ here are not same as self attention $(W^\Phi x^{(i)})$, here they're actually equal to word embedding $x^{(i)}$. They're treated different due to architecture of transformer.

i = dimension of encoding

→ repeats twice for sin & cos

* Transformer: it uses a k -halperin to calculate i

Positional Embeddings:

$$PE(pos, 2i) = \sin \frac{pos}{10000 \frac{2i}{d}}$$

$$PE(pos, 2i+1) = \cos \frac{pos}{10000 \frac{2i}{d}}$$

d = no. of features of word embedding & positional encoding

pos gives the position in words & is a vector unique to each word.

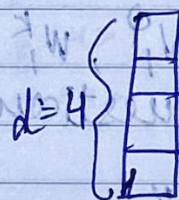
i = dimension of encoding

1. $k = 2i = 0$ $i = 0$

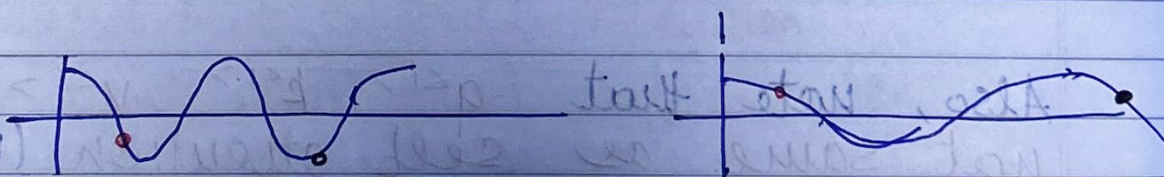
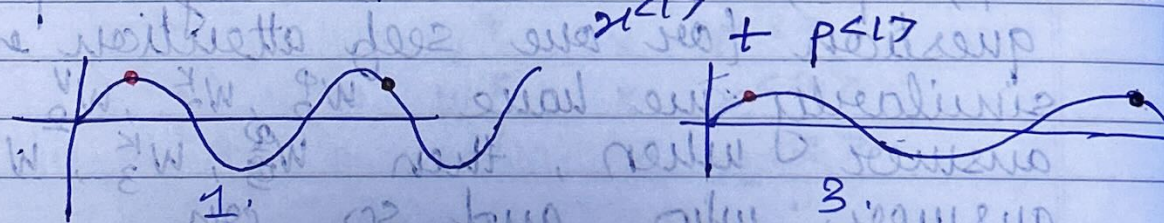
2. $k = 2i + 1 = 1$ $i = 0$

3. $k = 2i = 2$ $i = 1$

4. $k = 2i + 1 = 3$ $i = 1$



= input
fixed
encoder



$p < 3$

Input for decoder :

1. $\langle \text{SOS} \rangle$ (': $\langle \text{SOS} \rangle$ & $\langle \text{EOS} \rangle$ are predefined)

2. Jane

3. Jane visits

⋮

- Both the encoder & decoder layers repeat N times to give the optimum result.
- Add & Norm ~~to~~ layers are added after every layer to act as Batch-norm.