

# Week 4: Transformers

## 1) Transformer Network Intuition

sequential models:

RNN

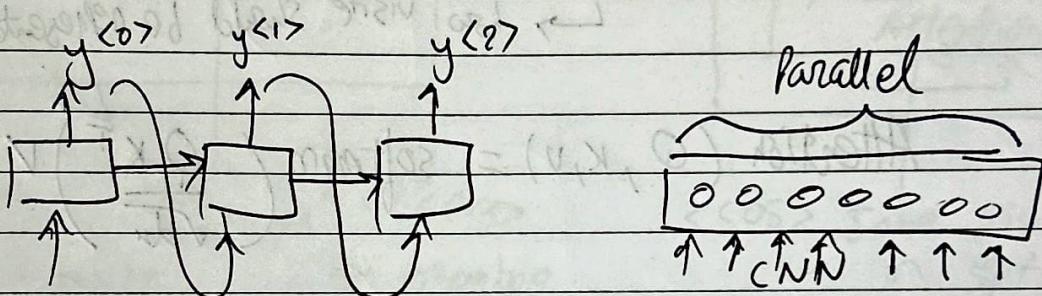
GRU

LSTM

→ increased complexity

- Attention + CNN

→ self-Attention ←  $A^{<1>} A^{<2>} A^{<3>} A^{<4>} A^{<5>}$   
 → multi-head attention (for loop over self attention)



## 2) Self-Attention

Jane visite l'Afrique en Septembre  
 $x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$

$A(q, K, V) = \text{attention-based vector representation of } v$   
 (calculate for each word  $A^{<1>} \dots, A^{<5>}$ )

To get  $A^{<3>}$ , it will look at the surrounding words to get the content of l'Afrique i.e. whether it is a holiday destination or historical representation on the second largest continent in the world.

## RNN attention

$$\alpha^{(t,t')} = \frac{\exp(e^{(t,t')})}{\sum_{t'=1}^T \exp(e^{(t,t')})}$$

## Transformers Attention

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

for l'Afrique  $\rightarrow q^{(3)}, k^{(3)}, v^{(3)}$   
 (query key value pairs)

Eg:

(what's happening)  
 (in l'Afrique?)  $q^{(3)} = W^Q \cdot n^{(3)}$   $\rightarrow$  asks a question (vectorized version)  
 (visite)  $k^{(3)} = W^K \cdot n^{(3)}$   $\rightarrow$  figures out most relevant answer to q  
 $v^{(3)} = W^V \cdot n^{(3)}$   
 $\hookrightarrow$  how visite should be represented in A

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

## 3) Multi-Head Attention

$$h = \# \text{heads}$$

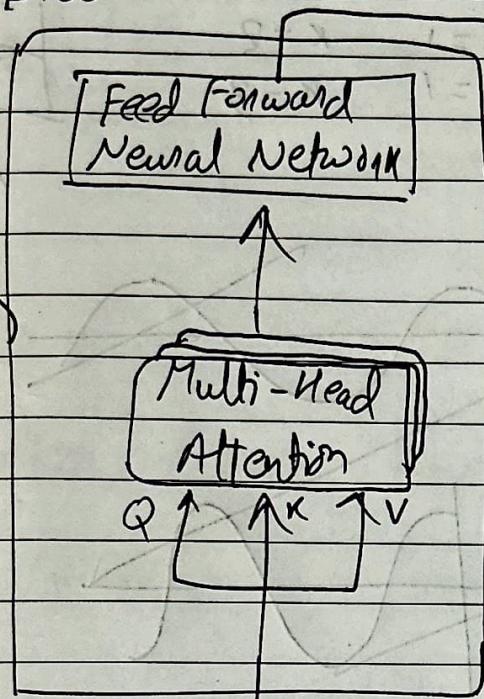
All the heads are computed in parallel instead of sequentially.

$$\text{head}_i = \text{Attention}(W_i^Q Q, W_i^K K, W_i^V V)$$

$$\text{MultiHead}(Q, K, V) = \text{concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h) \cdot W_0$$

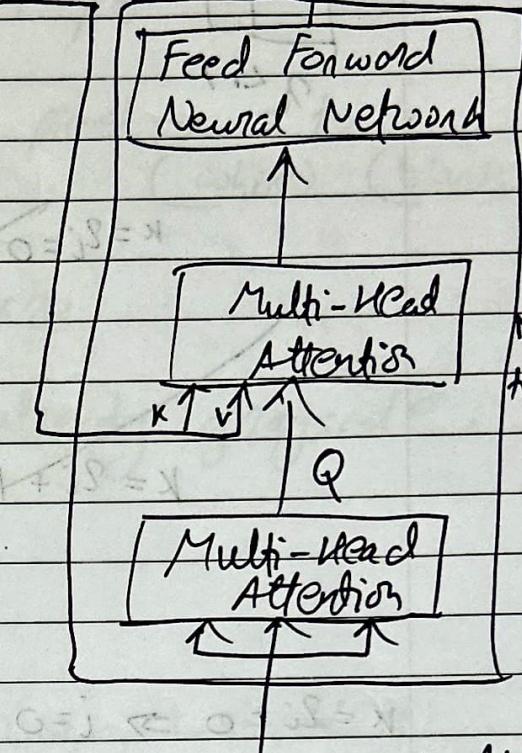
## 4) Transformer Network

Encoder



$\langle \text{SOS} \rangle$  Jane visits Africa in  
September  $\langle \text{EOS} \rangle$

Decoder



$\langle \text{SOS} \rangle, \langle \text{T}_1 \rangle, \dots, \langle \text{T}_m \rangle \langle \text{EOS} \rangle$   
Jane visite ... en septembre

$\langle \text{SOS} \rangle$  Jane visits Africa  
in september

Both encoder and decoder blocks iterate  $N$  times. Each of the output word generated is feed again into the decoder as input.

### ① Positional Encoding

$$\text{Equations: } PE_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{2i/d}}\right)$$

$$PE_{(pos, 2i+1)} = \cos\left(\frac{pos}{10000^{2i/d}}\right)$$

$d \Rightarrow$  dimension of the word embedding expand positional encoding  
 $pos \Rightarrow$  position of the word

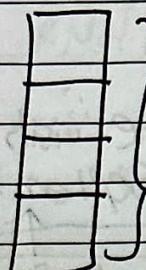
Page No. \_\_\_\_\_  
Date \_\_\_\_\_

$$\Leftrightarrow k/12 = i$$

word Embedding (vector)

$$pos = 1$$

3



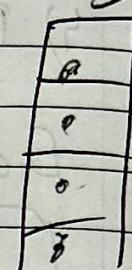
$$d=4$$

$$x^{<1>}$$



$$p^{<1>}$$

$i=0$	$k=0$
$i=0$	$k=1$
$i=1$	$k=2$
$i=1$	$k=3$



$$p^{<3>}$$

$$k=2i=0 \Rightarrow i=0$$

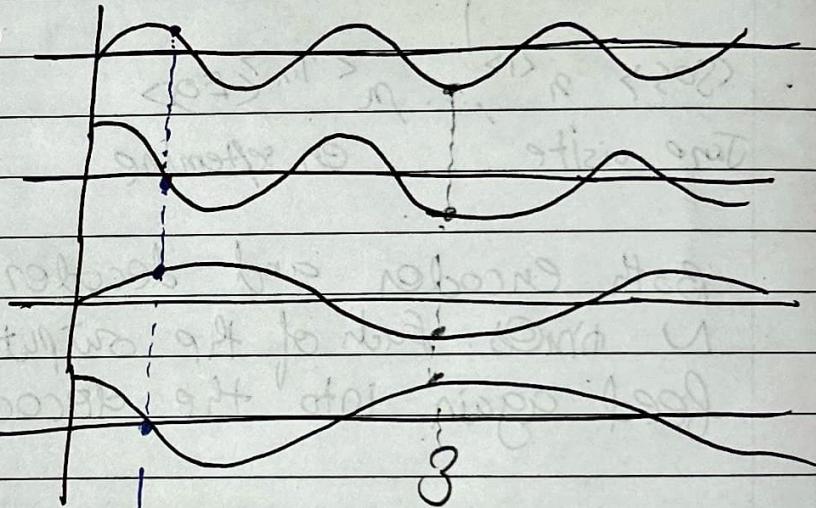
$$k=2i+1=1 \Rightarrow i=0$$

$$k=2i=0 \Rightarrow i=0$$

$$k=2i+1=1 \Rightarrow i=0$$

$$k=2i=2 \Rightarrow i=1$$

$$k=2i+1=3 \Rightarrow i=1$$



$p^{<3>}$  will be different than  $p^{<1>}$

The transformer also uses batch Norm and Add Norm type of layers to speed up learning.