

## (Week 4)

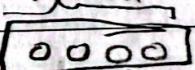
## ➡ Transformers Network intuition

- RNN → GRU → LSTM

The mode increase complexity & These model ingested the input word by word @ one token at time

- In the transformer architecture, it allows us to run a lot more of these computations for an entire seq. in parallel so we can ingest an entire seq. all at the same time

• The major innovation of the transformer architecture is combining the use of Attention based representation & CNN style. The Attention network is a way of computing very rich & useful representations of words but with something more akin to the CNN style of parallel processing



- Self attention The goal of it is if we have a sentence of 5 words,  $A^{<1>} - A^{<5>}$  we will compute 5 representations for these 5 words. & these representations are the attention based way of computing representation for all the words in your seq. in parallel

- Multi head attention is for loop over self attention processes. so it end up with multiple versions of these representations

## → Self Attention:

- to use Attention with a style more like CNN we need to calculate self attention where we create attention-based representations for each of the words in your input sentence
- after making the representations we call the five repres.  
 $A^{(1)}, A^{(2)}, \dots, A^{(5)}$

sentence: Jane visite l'Afrique en Septembre

$x^{(1)}, x^{(2)}, \underline{x^{(3)}}, x^{(4)}, x^{(5)}$

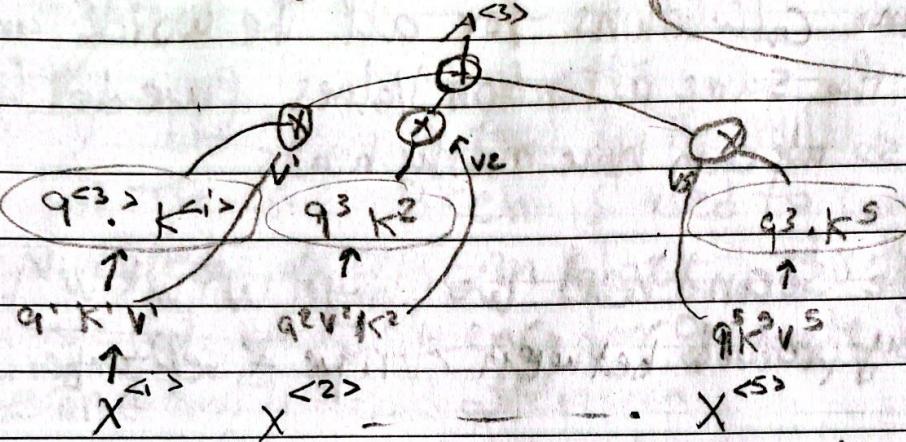
- if we take  $x^{(3)}$  (How self attention allows computing  $A^{(3)}$ )  
depending on how we're thinking of  $x^{(3)}$  we will choose to represent it differently by looking to the surrounding words.  
we will compute these representations in parallel for all 5 words
- what is the computations needed to go from  $A^{(3)}$  → selfatten.  
representation of  $x^{(3)}$ . [First] we associate each of the words with three values query key & value pairs, so if  $x^{(3)}$  is the word embedding for "l'Afrique" the query  $q_3$  is computed is as a learned matrix  $q^{(3)} = W_Q^T x^{(3)}$  & similar with key & value  
 $k^{(3)} = W_K^T x^{(3)}$   
 $v^{(3)} = W_V^T x^{(3)}$

These matrices  $W_Q, W_K$  &  $W_V$  are the parameters of the learning algorithms

- $q_3$  is might be a question ex: what happening there? (l'Afrique)  
& what we going to do is Compute the inner product bet  $q_3$  &  $k^1$   
 $q_3 \cdot k^1$  ---  $q_3 \cdot k^5$  & each one will tell us how good an answer of  $k^i$  to the question  $q_3$  & the goal of this operation is to pull up most info that's needed to compute most useful  $A^{(3)}$

Then we take  $q^{<3>} k^{<1>} \dots q^{<2>} k^{<5>} \dots$  and compute Softmax over them  $\Sigma_j \exp(q_i k_j^{<1>})$

Then finally we take the softmax values & multiply them with  $v^{<1>} \dots$  which is value for word ① & then sum it up to give  $A^{<3>} \dots$  or  $A(q, k, v) = \sum_i \frac{\exp(q_i k_i^{<1>})}{\sum_j \exp(q_i k_j^{<1>})} v_i$



The key advantage of this representation is the word of "l'Afrique" isn't some fixed embedding instead the self-att. mechanism realize that l'Afrique is the distribution of a visit which is more useful representation

The notation used is  $\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)$

## → multi-head Attention

- basically just a big for-loop over the self attention
- each time we calculate self attention for a seq it is called a head
- After the same calculations for all the words and end up with the same attention values if we do this many times so we can have many heads

• for example in the second head we change  $W^Q W^K W^V$  when so the key product between Afrique & September is the highest

- we can consider each of these heads as different features that could be passed to the NN

## → Transformer Network

- the first step in the transformer is these embeddings is sent into encoder block which has a multi-head attention layer, where you feed in the values Q, K, V computed from the embeddings & the weight matrices W
- Second the produced matrix is passed into feedforward NN to determine what interesting features there are in the sentence
- The encoder block is repeated N times Then feed the encoder into a decoder block

In the decoder block which its job to output the English sentence. The first output <SOS> token at every step the decoder block will input the first few words whatever we have already generated of the translation.

First The SOS is feed to the multi head attention block and just the <SOS> token is used to compute  $Q | K | V$  for the multihead attention block. This first multihead generate  $Q$  while the encoder generate  $K, V$  for the second multihead attention.

Then feed forward & repeat the decoder  $N$  times then predict the word that is going to enter with the <SOS> as input for the decoder

Steps makes transformer network work even better

① Positional encoding of input: we encode the position of elements in input by using combination of sine cos equations

ex. if the word embedding is vector of four values.  $\boxed{d=4}$   
 $x^{c_1} \sim x^{c_4}$ , so we create positional embedding vector of the same dimension  $\boxed{p^{c_1}}$

② Add & Norm  $\approx$  batch norm to speed learning.