# Multi-head Attention and Transformer

Tirtharaj Dash

Dept. of CS & IS and APPCAIR
BITS Pilani, Goa Campus

November 20, 2021

In the last lecture:

(1) attention — mechanism.

(2) Self-attention.

$$( Q, K, V )$$

"block"

$$\downarrow$$

Next:

Multi-headed attention.

"multi" "head"

$\downarrow$

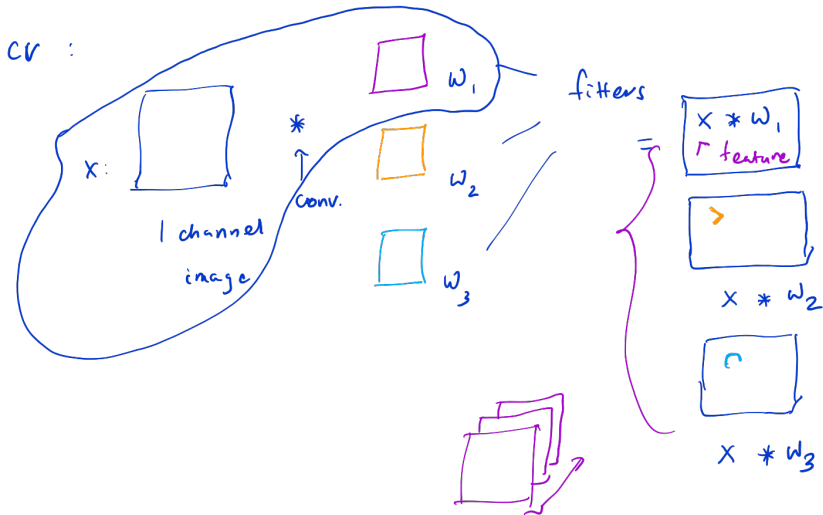"multiple Self-attention"

Understanding / Intuition of multihead attention:-

CV :

$X$: 1 channel image

$*$ conv.

$W_1$ $W_2$ $W_3$ filters

$= \{$

$X * W_1$ $\Gamma$ feature

$X * W_2$

$X * W_3$

Self-attention :-

Multiple self-attention

"Multihead attention"

Sequence:  "$I$  want  to  live  in  Europe. "

$(Q, K, V)$

"head 2"

$\left( W_Q^{(2)}, W_K^{(2)}, W_V^{(2)} \right)$

self-attention

$\left( W_Q^{(1)}, W_K^{(1)}, W_V^{(1)} \right)$ ← first feature

"head 1"

feature:  "Who ?"

"where?"

block

concatenate ( ☐ , ☐ , ☐ )

multi-head vector.

"$L$" - heads

where?

attention 1.

who!

$(W_Q^{(2)} q_1, W_K^{(2)} k_1, W_V^{(2)} v_1)$

$(W_Q^{(4)} q_6, W_K^{(4)} k_6, W_V^{(4)} v_6)$

$(W_Q^{(1)} q_1, W_K^{(1)} k_1, W_V^{(1)} v_1)$

$(W_Q^{(1)} q_6, W_K^{(1)} k_6, W_V^{(1)} v_6)$

$q_1, k_1, v_1$    $q_2, k_2, v_2$    $q_4, k_4, v_4$    $q_6, k_6, v_6$

$x_1$    $x_2$    $x_3$    $x_4$    $x_5$    $x_6$

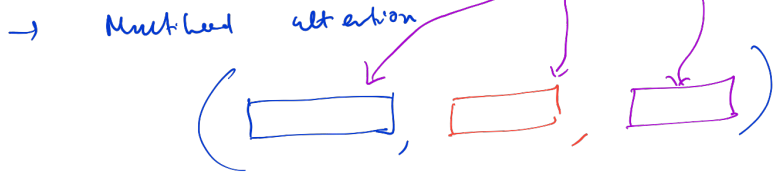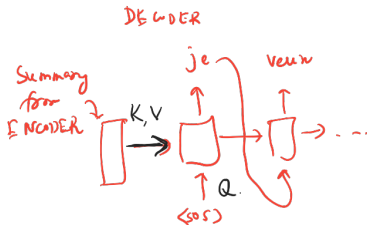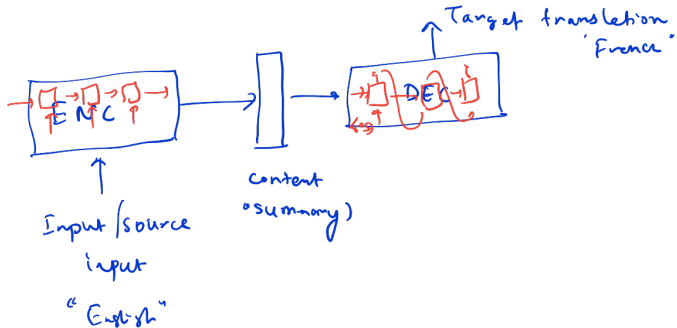I    want    to    live    In    Europe.

$$\text{attention score} \; (Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

→ Self-Attention allows us construct "rich" embeddings representations

for the input sequence

→ Multihead attention

$$\left( \boxed{\phantom{xxxx}}, \; \boxed{\phantom{xxx}} , \; \boxed{\phantom{xxxx}} \right)$$

# Transformers:- Seq2seq (Machine translation)



Target translation 'French'

ENC

DEC

Input/source input "English"

Content (summary)

DECODER

Summary from ENCODER

K,V

je

veux

Q

<SOS>

ENCODER

×N

N-times
(deep)

Feedforward
Network ∷∷ ○○○○ ← which parts of
        ○○○● this representation
             are good?

representation/embedding
for the
input

h-heads

Q ⌊ K ⌊ V

⟨sos⟩  I  want  to  live  in  Europe. ⟨eos⟩
$x_1$  $x_2$  $x_3$  $x_4$         $x_7$  $x_8$  $x_9$

DECODER

⟨eos⟩ je veun
         vinm
              ⟨eos⟩

Softmax

linear

Feedforward

×M

K ↑ ↑ v
       Q

h-heads

Q ⌊ K ⌊ v
      ↑ ↑?

⟨os⟩

⟨sos⟩je
⟨sos⟩ je venn
     (target language?)

Some additional machinery in Transformer (Enc, Dec):-

(1) Model might miss the "importance" of positions of the words in the input sentence.

→ Positional encoding. ( explicitly providing these info.)

(2) Stabilising training.

→ akin "Batch Norm" (CNNs, MLPs)

"add & norm" → Layer Normalisation (sequence models)

# Positional Encoding:
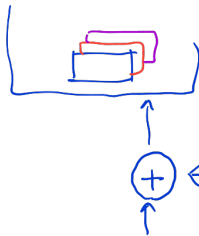
Words    :    word    $\langle$ embedding $\rangle$

word$_1$         word$_2$                                    word$_{T_x}$

$\langle$        $\rangle$   $\langle$        $\rangle$              $\langle$           $\rangle$

$\langle$ pos info $\rangle$   $\langle$ pos info $\rangle$              $\langle$ pos info $\rangle$

Positional encoding using trigonometric functions
(sine, cosine)

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

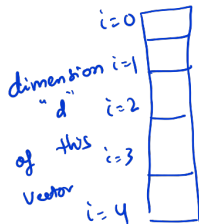$$PE_{pos, 2i+1} = \cos\left(\frac{pos}{10000^{2i/d}}\right)$$

$\oplus$ ← Positional encoding.

$\langle sos \rangle$ I want to live in Europe $\langle eos \rangle$

ENC
DEC

Word embedding vectors
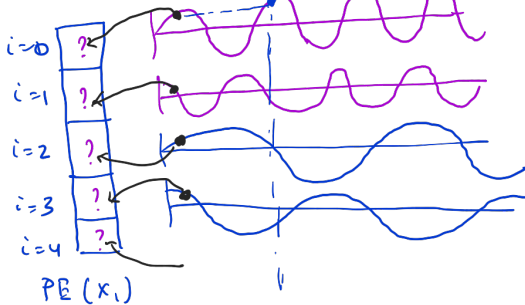
$$\sin\left(\frac{pos \to 1}{10000^{2i/d}}\right) \to n$$

dimension "d" of this vector

i=0
i=1
i=2
i=3
i=4

$X_1$

Word $\boxed{1}$

pos=1

$d=5$

i=0 ?
i=1 ?
i=2 ?
i=3 ?
i=4 ?

$PE(X_1)$

$\boxed{d=5}$

$(+)$

pos=1     pos=2

adding directly.

Word pos = n.

dictionary size = 10000

< 10000 dimen >

Word embedding: encoding the word

Pos- encoding: ordering or positions $\left. \begin{array}{l} \\ \\ \end{array} \right\}$ ⊕

Why sine / cosine ?

→ range [-1, 1]
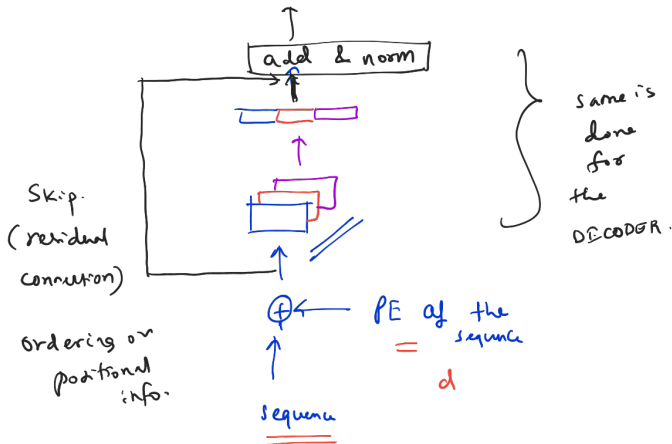       small

→ order: Periodicity in these
                             functions

d ?

what exactly is <u>d ?</u> (max length of the
sentence in the
input/source language)

## Add & Norm:



add & norm

same is
done
for
the
DECODER.

Skip.
(residual
connection)

ordering or
positional
info.

⊕ ← PE of the
= sequence
d

sequence

# Homework:

(1) Layer normalisation

$$( vs. \quad BatchNorm )$$

(2) Decoder: Multi-head attention block.

"Masked Multi-head attention"

↑ (Training)