

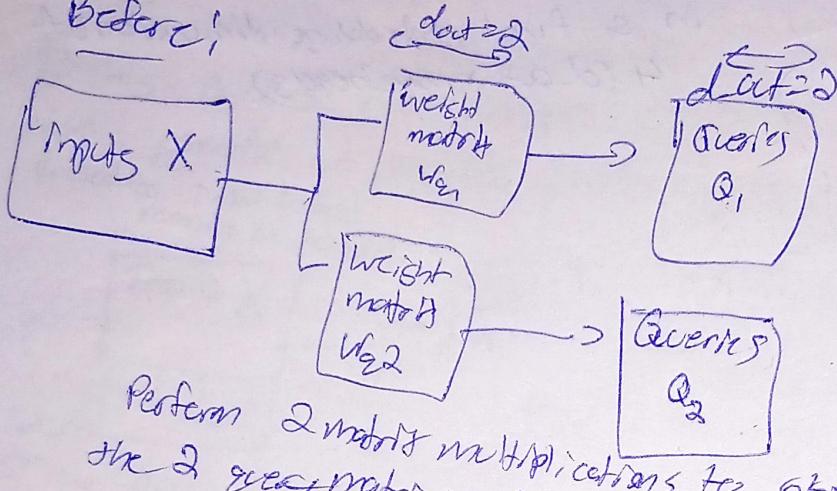
Mult. Head Attention part 2

b) Implementing multi-head attention w/ weight splits

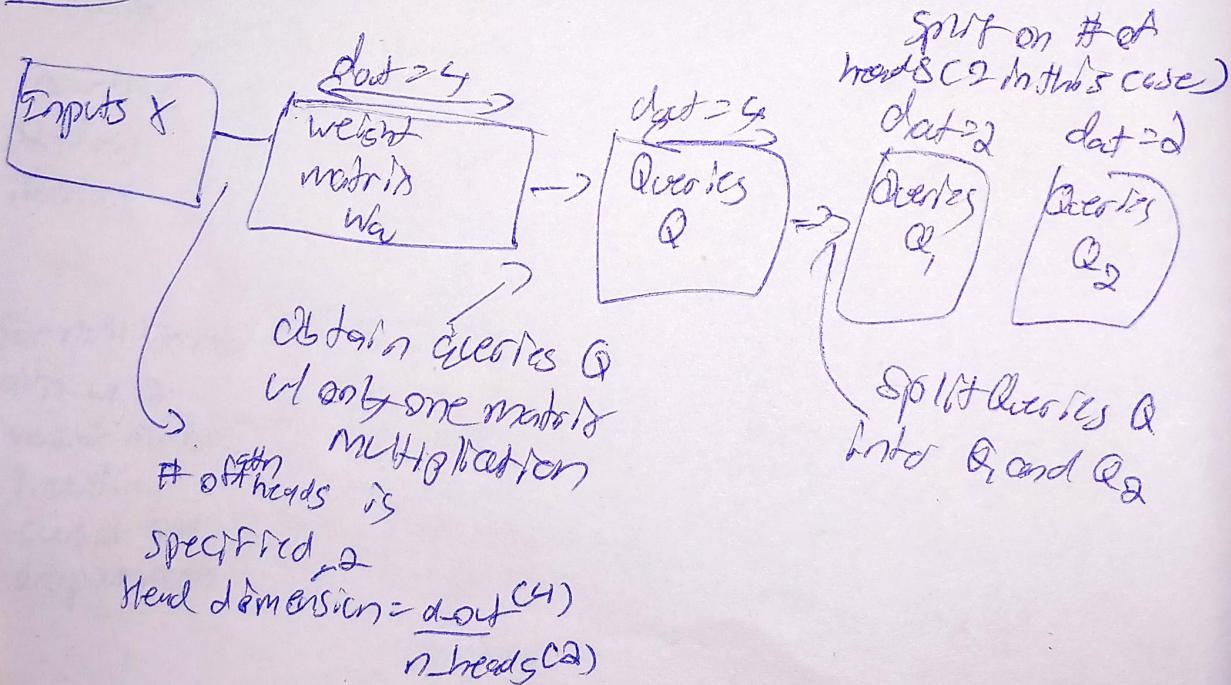
c) Instead of maintaining 2 separate classes;

MultHeadAttentionWrapper and CrossAttention,
we combine both of these into a single
MultiAttention class

Before:



After



Let us start w/ simple example:

Step ①: Start w/ input

b, num_tokens, d-in = 1, 3, 6

The ~ 6 (first row)

cat ~ 6 (second row)

sleeps ~ 6 (third row)

Step ② Decide d_out, num_heads

In GPT
models
 $d_{in} = d_{out}$
typically

$x = \text{torch.tensor}([[[\dots]]])$
(batch=1, num_tokens=3, d-in=6)
CV
[]
3x6

Step ③: Initialize trainable weight matrices

repeat for locs, query, value (W_q, W_k, W_v)

(6x6) d_in x d_out

(6x6)

$$W_q = \begin{bmatrix} \dots \end{bmatrix}$$

(6x6)

$$W_k = \begin{bmatrix} \dots \end{bmatrix}$$

(6x6)

$$W_v = \begin{bmatrix} \dots \end{bmatrix}$$

(6x6)

Step ④: Calculate keys, queries, value matrix

(1x3x6)

(6x6)

$$\begin{bmatrix} \dots \end{bmatrix}$$

(6x6) (6x3x6)

b num tokens d_out

Queries
(1x3x6)

Values
(1x3x6)

Multiplying

(1x3x6) @ (6x6)

compare

(3x6) @ (6x6) -> (3x6)

Add batch,

(1x3x6)

Step ⑤: Unroll last dimension of logits, queries, and values to include num_heads and head_dim

$$\text{head_dim} = \text{d_act} / \text{num_heads} = 6/2 = 3$$

$(b, \text{num_tokens}, \text{d_act}) \rightarrow (b, \text{num_tokens}, \text{num_heads}, \text{head_dim})$

$$(1, 3, 6) \rightarrow (1, 3, 2, 3)$$

"Essentially splitting

that 6-dim-d_act into

2 attn heads, each

with 3 dimensions (check dim)

Reshaped Queries matrix

Queries matrix

tensor($\Sigma \Sigma \Sigma$, shape=(1, 3, 2, 3))

word1 ← [--- token] word2 ← [--- token]

Queries, logits, values of

were same dim.

word1 ← [--- token 3]
word2 ← [--- token 2]

word1 ← [--- token 3]
word2 ← [--- token 2],
3dim head

Step ⑥: Group matrices by "number of heads"

$(b, \text{num_tokens}, \text{num_heads}, \text{head_dim}) \rightarrow (b, \text{num_heads} \times \text{tokens}, \text{head_dim})$

$$(1, 3, 2, 3) \rightarrow (1, 2, 3, 3)$$

Grouping w.r.t
of tokens

Grouping w.r.t
of heads

Transposed Queries matrix

[---] → head 1
[---] → head 2
[---] → head 3

Queries, logits, values
all have same
dimensions

[---] → head 1
[---] → head 2
[---] → head 3

Now we can now group w.r.t heads, we can compute attn score
for each head separately

Step 7: Find attention scores

Queries @ keys, transpose(2,3)

Transpose the last 2

dimensions (num tokens and head-dim)

keys (old)

[[[1 2 3],

[4 5 6],

[7 8 9]],

[[1 2 3],

[4 5 6],

[7 8 9]]]

keys, transpose new

[[[[+ 4 7],

[2 5 8],

[3 6 9]]],

[[1 4 2],

[2 5 8],

[3 6 9]]],

head 1

Draw from old scores
(column)

Queries x keys. T

(1, 2, 3, 3) & (1, 2, 3, 3)

check:

(b, #head, #tokens, head-dim) & (b, numhead,
head-dim,

numtokens)

(1, 2, 3, 3)

Queries @ keys, transpose(2,3)

fields attn scores

matrix

(b, numheads, numtokens,
numtokens)

the cat sleeps

the EEEC - - - - -
cat - - - - -
sleeps - - - - -

first
head)

row is attn-score
of first word
between all other
ones

Step 8: Find attention weights

↳ Mask attention scores to
implement causal attention

[[[-inf, -inf]]]

Divide by Jhead-dim

= $\sqrt{\text{dim}}$ / number

= $\sqrt{8/2} = \sqrt{3}$

[- - - - -]

[- - - - -]

[- - - - -]

stable gradients

Make sure that
variance scales up w/
of dimensions and

scale for prevent variance
from blowing up

!!

apply softmax
on next page

~~SOKEWAT~~ $\Sigma \Sigma [1, 0, 0]$, Attention weights

$[0.964, 0.035, 0]$, $(1, 2, 3, 3)$

$[0.041, 0.2603, 0.6988]$ ($b, \text{num_heads}, \text{num_tokens}, \text{num_tokens}$)

$[1, 0, 0]$, ↗ can also improvement depend after softmax.

$[0.122, 0.8273, 0]$,

$[0.3897, 0.4627, 0.15662222]$

Step (9), Context vector = Attention \odot values

Attention weights \odot values
 $b, \text{num_heads}, \text{num_tokens}, \text{head_dim}$

Attention weights

~~values~~

$(1, 2, 3, 3)$ \odot

↓

$(1, 2, 3, 3)$

$(b, \text{num_heads}, \text{num_tokens}, \text{head_dim})$

$[0, 0, 0, \dots]$ Each row corresponds

$[0, 0, 0, \dots]$ to one token, and

$[0, 0, 0, \dots]$ represents context vec for that particular

token, dimensions

= to head_dim

$[0, 0, 0, \dots]$

$[0, 0, 0, \dots]$

$[0, 0, 0, \dots]$

we want about as

resultant dim. after

context_vec so

need to flip

back to my tokens

and num_heads w/ transpose

Step (10)! Reformatted context vectors

$C_b, \text{num_heads}, \text{num_tokens}, \text{head_dim} \rightarrow C_b, \text{num_tokens}, \text{num_heads}$
head-dim

$\{\{\{1 + 2\} 3\} 3\}$ token 1

$\{4 5 6\} 3$

$\{\{2 8 9\} 3\}$ token 2

$\{10 11 12\} 3$ $C_b, \text{num_tokens}, \text{num_heads}, \text{head_dim}$

$\{13 14 15\} 3$ tokens 3

$\{16 17 18\} 3$ tokens 4

Step (11)! Combine heads

Flatten each token output into each row

First row consists of merging the 2 heads, for first token

Second row consists of merging the 2 heads, for second token

same w/ 3rd token/row nr.

$\{\{1 2 3 4 5 6\},$
 $2 8 9 10 11 12\}$
 $\{13 14 15 16 17 18\}\}$

$C_b, \text{num_tokens}, \text{d_out}$
 $(1, 3, 6)$

↓
Final result!