

Feed Forward Block -

→ Moves the things ahead

↳ Sequence ↴

linear layer (weights)



Activation (GELU)



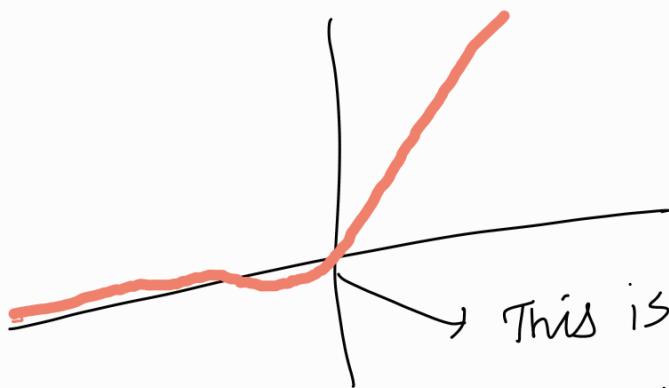
linear layer (weights)



Dropout

GELU Activation →

Gaussian Error Linear Unit



This is -ve for values
very near to zero &
helps in convergence
some how?

$$= (0.5 \times x) \left(1 + \tanh \left[\sqrt{\frac{2}{\pi}} (x + 0.044 \cdot x^3) \right] \right)$$

→ Layer Normalizations

Shift & Scale $\begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$
 ↓
 Matrix with I $\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$
 Matrix with zero

→ makes for mean zero &
 variance 1.

→ so let's say you have weight
 matrix → it makes the mean
 zero & var = 1 for the 1st
 axis to step live it & we add
 some shift & scale to that,
 the shift & scale can be learned
 by the model during back prop

→ Causal Attention

Basically attention mech which only lets the model to see the prev. tokens during training so that it can predict the next token. This is done by assigning K, Q, V matrices.

We have each token is

query] related to other
keys token

values ↓

→ query @ keys^T



This is like each word is a query & has many keys (all other tokens) to relate, it relates with them $\alpha @ k \cdot T \rightarrow \text{atten-scores}$

we get attention scores

+

lets say we have a input
of size $[2, 6, 3]$
batch tokens $\downarrow + \curvearrowright$ in-dim

* lets say you want

DP dim of 2

final DP will be $(2, 6, 2)$

context = 6

$\rightarrow q$ shape = $(2, 6, 2)$

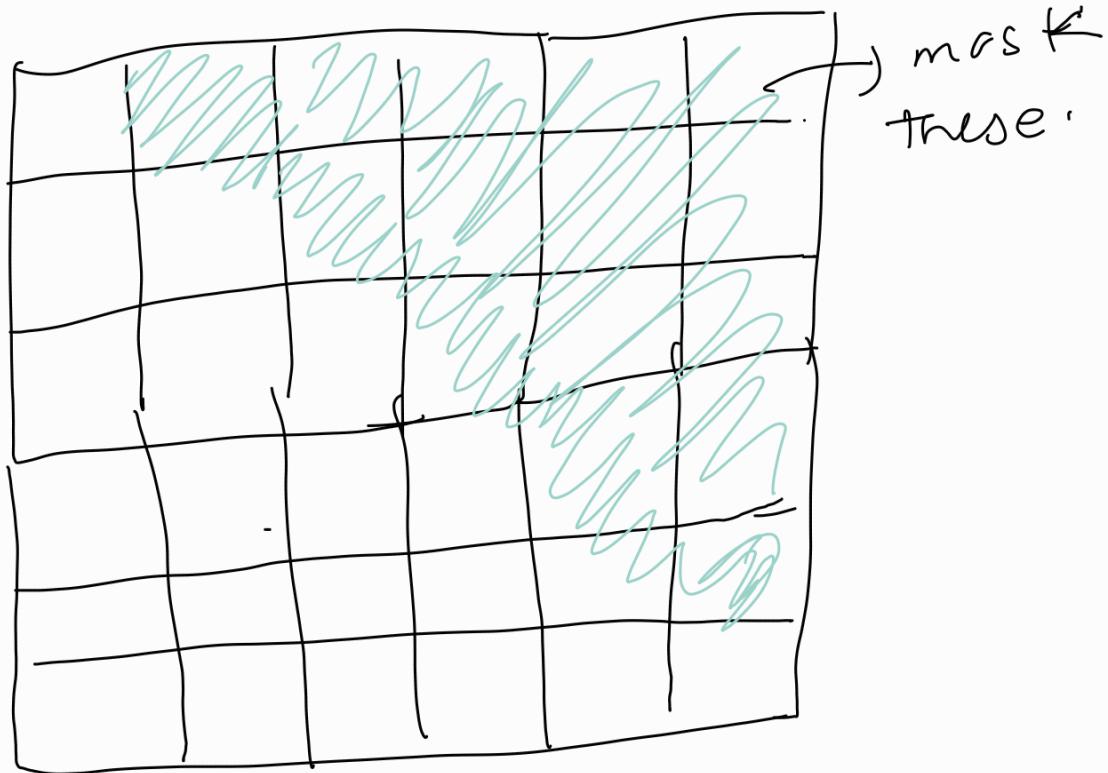
k shape = $(2, 6, 2)$

v shape = $(2, 6, 2)$

$(2, 6, 2) \cdot ^T$
 $(2, 2, 6)$

$\rightarrow (2, 6, 2) @ (2, 2, 6)$
 $6 \times 6 \rightarrow \text{attn-scores}$

now we mask the future ops.
→ add dropout.



mask the upper portion of the
matrix triu.

→ make it as register buffer so
that it can be on GPU wise
weights. → Not on same device
error.

register.mask buffer[{"name"},
Torch.triu(torch.zeros(CL, CL),
diagonal=1)]

↳ apply softmax → to basically round up the attn scores b/w 0 & 1 value. } where their sum = 1.

↳ its like fixing the value b/w 2 numbers

→ These are now attn-weights

↓
dropout random values.

↓
now we matmul attn-weights with values. to calculate context-vectors

$$CV = \begin{matrix} \text{attn-weights} \\ (6, 6) \end{matrix} @ \begin{matrix} Q \text{ values} \\ (M, 6, 2) \end{matrix}$$

$\rightarrow [2, 6, 2]$

↓
this aligns
(it not transpose the result
to match (Batch, Tokens,
 d_{out})).

↳ out-proj (optional)
↳ embed-dim \rightarrow embed-dim
↳ helps with OP's.

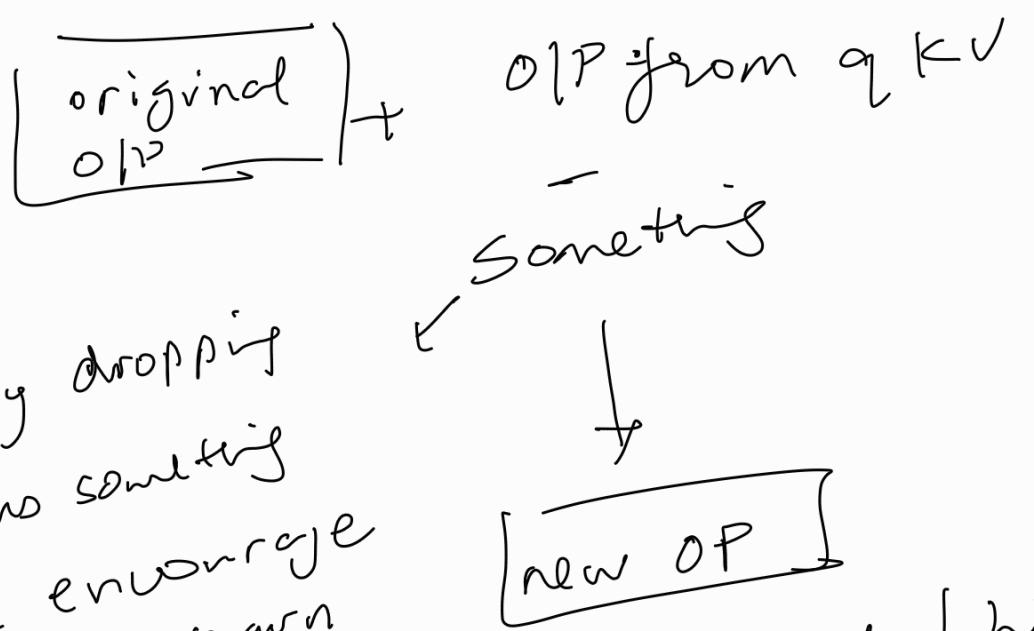
↳ resid-drop } (optional)

Dropout
→ interpret it like this



We subtract or
drop out from the
DOP from qKV

so it becomes



→ Transformer Block

IPs
↓
layer Norm 1

↓
attention

↓
dropout
↓
Layer Norm 2 → Feed Forward

↓
output

↓

↓
dropout

↑

↓
Layer Norm 2 → Feed Forward

↓
This is a point where we can add shortcuts to remove problem of vanishing gradient, basically we bypass the feedback & directly give it to the FF-module.

attention block contains the causal after block \rightarrow can add multiheads there -

\rightarrow feed forward block is linear

↓
linear

↓
relu

↓

dropout

↓

linear

→ Putting it together ~

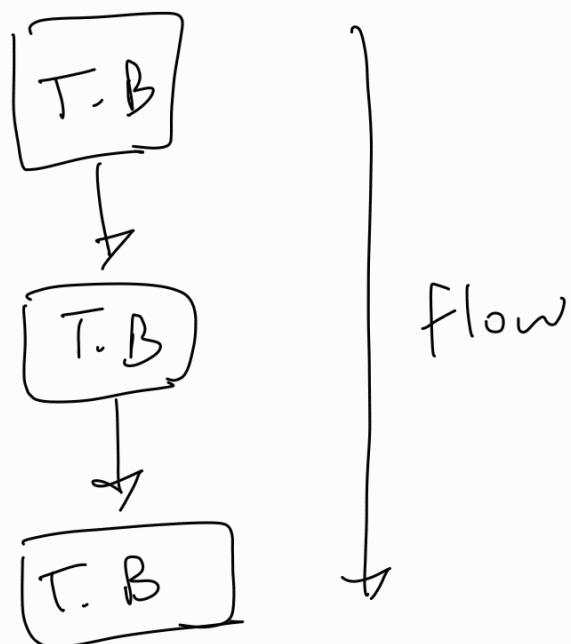
→ word embeds

→ pos embeds

→ configs.

we make multiple transformer
blocks which is layers

in the model



we make a list \rightarrow Module List +

with transformer blocks &

configurations like

embeddings, layers, heads,

head-dims, dropout,

context length, etc --

\rightarrow add final Norm layer

for finl op's.

\rightarrow make a layer that gives out
probab distribution of words.

(this is logits) \rightarrow helps

get the next word \rightarrow one with
highest probab.

→ flow
→ inputs → tokens → to device
+
pos embeds (create)
arrange → to device

↓
combine all embeddings

$x = \text{token} + \text{pos}$
+
dropout

↓
pass through layer of ^{Trans} former
Blocks

final ↓ Norm layer

↓
final Prob dist (logits)

↓
returns logits & when to
get the next word $\rightarrow \max(\text{logits})$

\downarrow
Max - logits \rightarrow decode with
token embeds

\downarrow
Get the Next
word!

— Picture $\cup \underbrace{\dots}_{\text{...}}$