

# Attention Is All You Need

## Transformer

### I. ARCHITECTURE $\Rightarrow$

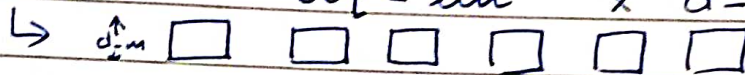
#### 1. Input Embeddings $\Rightarrow$

cat  $\rightarrow$  1234  $\rightarrow$   
Unique Input Id

|      |
|------|
| 1000 |
| 24   |
| 6084 |
| :    |
| 874  |
| 8888 |

$\uparrow$   
d-model  
 $\downarrow$

Whole text  $\rightarrow$  seq-len  $\times$  d-model



(seq-len = 6)

Code  $\rightarrow$

- IE Layer =  $\text{nn.Embedding}(\text{vocab-size}, \text{d-model})$
- Fwd fn  $\rightarrow$   $\text{IE}(u) * \sqrt{\text{d-model}}$  [1]

#### 2. Positional Encoding:

$F^u$  to apply it,  $PE^{[2]} = \begin{cases} \sin\left(\frac{\text{pos}}{10000^{\frac{2i}{\text{d-model}}}}\right), i \rightarrow \text{even} \\ \cos\left(\frac{\text{pos}}{10000^{\frac{2i}{\text{d-model}}}}\right), i \rightarrow \text{odd} \end{cases}$

Input token  $\rightarrow$  PE  $\rightarrow$  Input token + Positional info

Code  $\rightarrow$

- $pe = \text{torch.zeros}(\text{seq-len}, \text{d-model})$   $\rightarrow$  Dummy tensor initially
- $pos = \text{torch.arange}(0, \text{seq-len}).\text{unsqueeze}(1) \Rightarrow [0, 1, 2, \dots, 5]$
- $pe[\text{even columns}] = PE(pos)$  &  $pe[\text{o.c.}] = PE(pos)$
- \* register\_buffer(pe)
- Fwd fn  $\rightarrow$   $u = u + pe$   $\rightarrow$  return dropout(u) [3]

[1], [2], [3]  $\Rightarrow$  According to Paper



### 3. Layer Normalisation:

- $\alpha = \text{torch.ones}(\text{features})$
- $\epsilon = 10^{-9} \Rightarrow \alpha, \beta^{[4]}$  Learnable nn-Parameters()
- Fwd fn  $\Rightarrow$ 

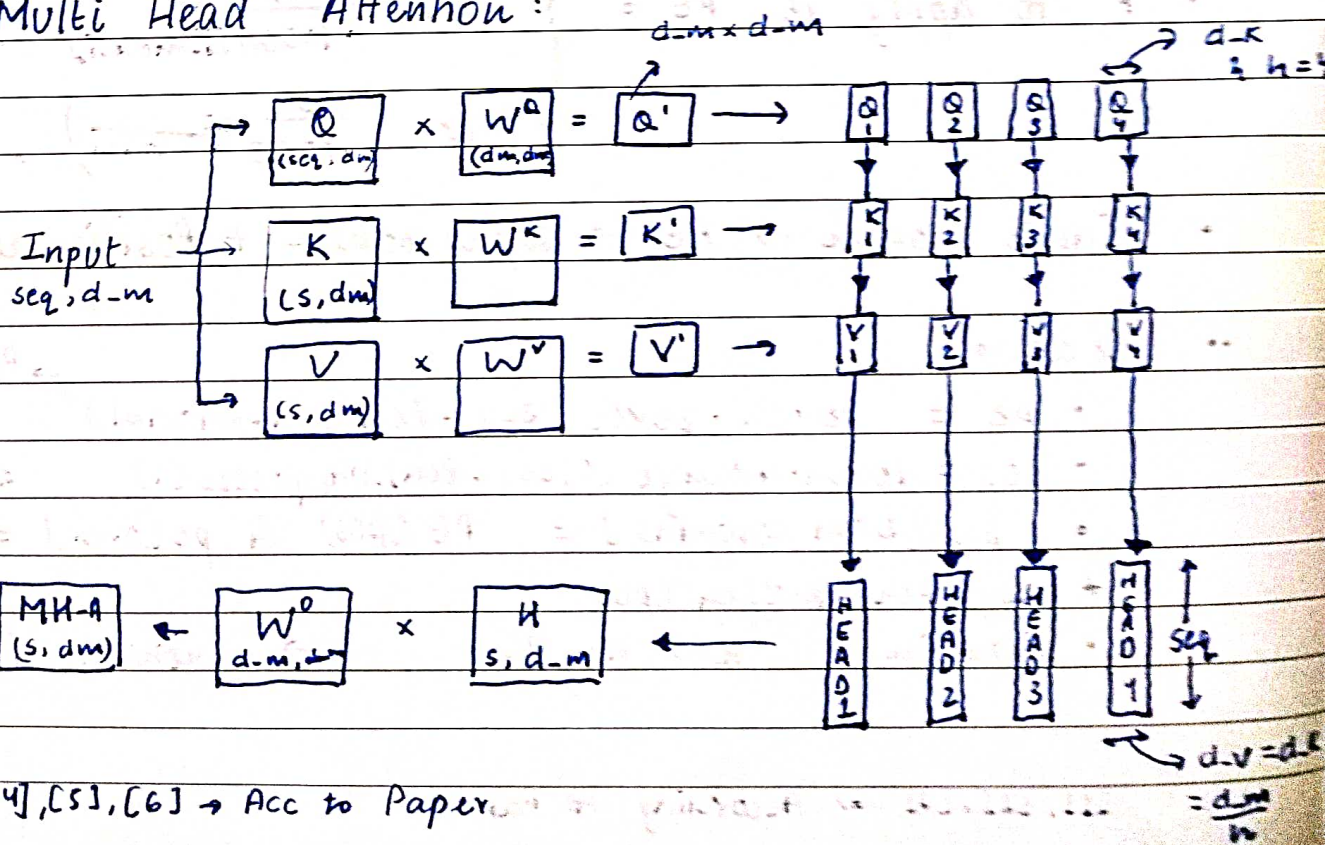
$$\cdot x = \alpha \left( \frac{x_i - \text{mean}_{\text{axis}}}{\sqrt{\text{std dev}(x_i) + \epsilon}} \right) + \beta^{[5]}$$

$\Rightarrow$  This mean is over dim=-1.

### 4. Feed Forward Network:

- Simple NN with dropout to train weights & biases.
- $\text{FFN}(x) = \max(0, xW_1 + b_1) \cdot W_2 + b_2$  [6]
- Code  $\rightarrow$  `FFN_layers = nn.Sequential(Linear(d-m, d-ff), ReLU, dropout, Linear(d-ff, d-model))`

### 5. Multi Head Attention:



[4], [5], [6]  $\rightarrow$  Acc to Paper



$$d_{\text{model}} = d_m \Rightarrow d_m // h = d_k = d_v$$

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) \cdot V \quad [7]$$

Code  $\Rightarrow$

$$\cdot \text{init} \Rightarrow d_k = d_v = d_m // h$$

$$w_q, w_k, w_v, w_o = \text{nn.Linear}(d_m, d_m)$$

$$\cdot \text{attention} \Rightarrow \text{att}^n\text{-scores} = (\text{query} * \text{key}^T) / \sqrt{d_k}$$

$$\text{att}^n\text{-scores} \cdot \text{masked\_fill}(0, 10^{-9})$$

$$\text{return softmax}(\text{att}^n\text{-sc}) * \text{Value}$$

$$\cdot \text{forward} \Rightarrow (q, k, v, \text{mask}) \Rightarrow$$

$$\text{query} = w_q(q), \text{key}, \text{value similarly. (Batch, SL, d)}$$

$$\text{query} = \text{query.reshape}(q.\text{shape}[0], [1], h, d-k) \cdot \text{transpose}(1, 2)$$

$$\# q: d_m \rightarrow h * d-k, \text{Final shape} \rightarrow (B, h, s_k, d-k)$$

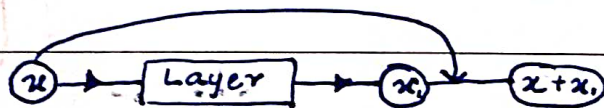
$$\text{key \& value}$$

$$x = \text{attention}(\text{query}, \text{key}, \text{mask})$$

$$(B, SL, d_m) \leftarrow x = x \cdot \text{transpose}(1, 2) \cdot \text{contiguous}(). \text{view}(x.\text{shape}[0], -1, d)$$

$$\text{return } w_o(x)$$

## 6. Residual Connection:



$$\cdot \text{Code} \Rightarrow \text{return } x + \text{dropout}(\text{sublayer}(\text{LayerNorm}(x)))$$

↓  
Layer being applied



## 7. Encoder:

### \* Encoder Block $\Rightarrow$

init  $\Rightarrow$  (features, attn<sup>n</sup>-block, ffn, dropout)  $\Rightarrow$

self-attn<sup>n</sup>-block = attn<sup>n</sup>-block, feed\_fwd\_net = ffn

residual\_conns = ModuleList(RN(2)) [8]

fwd  $\Rightarrow$  (x, src\_mask)  $\xrightarrow{q, k, v}$

attn<sup>n</sup>-output = attn-block (x, x, x, src\_mask)

x = residual\_conns[0] (x, attn<sup>n</sup>-output)

x = residual\_conns[1] (x, ffn)

return x

### \* Encoder $\Rightarrow$

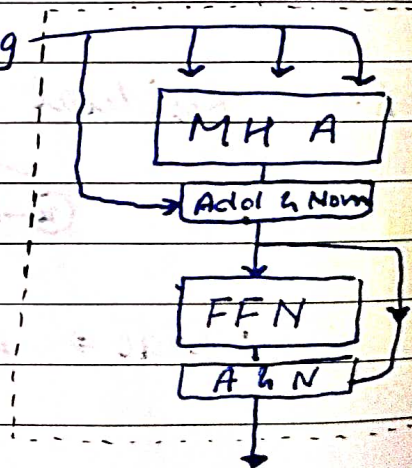
init  $\Rightarrow$  layers\_in\_encoder, LayerNorm

fwd  $\Rightarrow$  for layer in l\_in\_e  $\Rightarrow$  x = layer(x, mask)  
return LN(x)

## Encoder Block Diagram $\Rightarrow$

Input  $\rightarrow$  Input Embed  $\rightarrow$  Pos<sup>n</sup> Encoding

x N  $\leftarrow$  Encoder







9. Projection Layer  $\Rightarrow$ 

- init  $\Rightarrow$  proj = nn.Linear (d-m, vocab-size)
- fwd  $\Rightarrow$  return log\_softmax (proj(u), dim=-1)

10. TRANSFORMER  $\Rightarrow$  Creates structure for Transformer.

- init  $\Rightarrow$  encoder, decoder, src-embed, tgt-embed, src-pos, tgt-pos, proj<sup>n</sup> layer
- encode  $\Rightarrow$  (src, src-mask)  $\Rightarrow$   
src  $\rightarrow$  Input Emb  $\rightarrow$  Pos<sup>n</sup> Enc  $\rightarrow$  Masking & Encoder F<sup>n</sup>
- decode (encoder-opp, src-mask, tgt-mask,  
tgt  $\rightarrow$  IE  $\rightarrow$  PE  $\rightarrow$  Masking src & tgt & Decoder F<sup>n</sup>
- project (u)  $\Rightarrow$   
project (u)

11. BUILDING TRANSFORMER  $\Rightarrow$  Initialises Transformer

- build Transformer (src-vocab-len, t-v-l, src-seq-len,  
t-s-l, d-model = 512<sup>[10]</sup>, N=6, h=8  
dropout = 0.1, d-ff = 2048<sup>[10]</sup>)

$\Rightarrow$  src-embed, tgt-embed from src, tgt-vocab-size  
src-pos, tgt-pos from src-seq-len & t-s-l  
encoder-blocks = []  $\rightarrow$  append encoder block after init  
decoder-block = []  $\rightarrow$  append decoder block after init

encoder from e-b & decoder from d-b

proj-layer using tgt-vocab-size.

transformer initial<sup>n</sup> & setting parameters (xavier)

## II. DATA

### 12. Bilingual Dataset :

- ↑ class :
  - init ( ds, tokenizer = src, t = t, src\_lang, t = t, seq\_len ) ⇒  
initialise all these ↑ +  
sos\_token, eos\_token, pad\_token using token-to-id
  - len-- → return len(ds)
  - getitem-- (idx) ⇒  
src\_tgt\_pair = ds[idx] → src\_text, t = t  
enc\_input\_token, d-i-t &  
enc & dec pad tokens.  
→ encoder-ip = [sos, enc\_input, pad-enc]  
lllly dec & label.

### 13. Build Tokenizer, Dataloaders & Transformers ⇒

- get\_all\_sentences → yield all items
- get\_or\_build\_tokenizer(config, ds, lang) ⇒  
check tokenizer else  
tokenizer = sub-word / word --- level
- get\_ds(config) ⇒  
ds load, get src & tgt tokenizer, train-test split  
train-ds, test-ds, max\_len\_src & max\_len\_tgt  
building train, test dataloaders
- get\_model() ⇒ model = build\_transformer(.....) <sup>ENJ</sup>



## II. DATA

### 12. Bilingual Dataset :

- $\uparrow$  Class :
  - init ( ds, tokenizer=src, t=t, src-lang, t-l, seq-len )  $\Rightarrow$  initialise all these  $\uparrow$  +
  - sos-token, eos-token, pad-token using token-to-id
  - len--  $\rightarrow$  return len(ds)
  - getitem--(idx)
    - src-tgt-pair = ds[idx]  $\rightarrow$  src-text, t-t
    - enc-input-token, d-i-t &
    - enc & dec pad tokens.
    - $\rightarrow$  encoder-ip = [sos, enc\_input, pad-enc]
    - lll<sup>yy</sup> dec & label.

### 13. Build Tokenizer, Dataloaders & Transformers $\Rightarrow$

- get\_all\_sentences  $\rightarrow$  yield all items
- get\_or\_build\_tokenizer(config, ds, lang)  $\Rightarrow$ 
  - check tokenizer else
  - tokenizer = sub-word / word --- level
- get\_ds(config)  $\Rightarrow$ 
  - ds load, get src & tgt tokenizer, train-test split
  - train-ds, test-ds, max-len-src & max-len-tgt
  - building train, test dataloaders
- get\_model()  $\Rightarrow$  model = build\_transformer(....) <sup>CH3</sup>



### III. TRAINING =>

#### 14. Config File =>

- get-config => returns config
- get wt file path => Makes file path for each epoch
- latest-weights-file-path => Finds wts of newest save

#### 15. Training =>

Initialises everything, loss, optim, path, device.

Loads model if saved model present.

Training loop.

Saving Model

#### 16. Evaluation =>

Evaluates the last saved model.