

Out of the 512 dimension we can add one more embedding to end and make it 513 dimension and add $[1, 2, 3, 4, \dots]$ to it.

PROBLEM

- unbounded \rightarrow no upper limit
- BP hates huge numbers as it creates instability.

- discrete Nos \rightarrow we need continuous numbers.

Relative Position

- subject can be 5th and verb be 6th or it can be 105th and Verb can be 106th.
- their relationship of verb makes the reading part easy for model.

Trigonometric functions \rightarrow

- Instead of adding discrete nos $[1, 2, 3, \dots]$ we add the sine values of it.
- Repetition will occur hence we add the cosine values also

$$\text{PE} = \boxed{\begin{matrix} 1 & 2 & 3 & 4 & 5 \\ \sin(1) & \cos(1) & \sin(2) & \cos(2) & \dots \end{matrix}}$$

To reduce more repetition we add more components of sine & cosine function.

e.g. $(\sin(1), \cos(1), \sin(12), \cos(12), \dots)$

Do decide how many sine and cosine

*batch Normalization

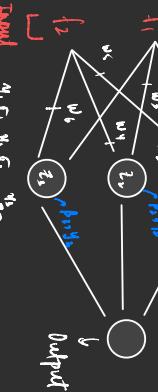
Rescale
Previously Normalized
Batch Normalization



$$\text{PE memory} = \boxed{\begin{matrix} 1 & 2 & 3 & 4 & 5 \\ \vdots & + & \vdots & \vdots & \vdots \end{matrix}}$$

$$\text{PE memory} = \boxed{\begin{matrix} 1 & 2 & 3 & 4 & 5 \\ \vdots & + & \vdots & \vdots & \vdots \end{matrix}}$$

We add this vector and do not normalize



$$\text{BN calculation} \rightarrow \frac{f_i - \bar{w}_i}{\sigma_i} = 0.3649 + \beta_1$$

$$\frac{6 - 4.6}{\sqrt{0.3}} = 0.90 \rightarrow 0.90 \times \beta_2 + \beta_3$$

$$\text{BN calculation} \rightarrow \frac{f_i - \bar{w}_i}{\sigma_i} = 0.3649 + \beta_1$$

$$\frac{6 - 4.6}{\sqrt{0.3}} = 0.90 \rightarrow 0.90 \times \beta_2 + \beta_3$$

$$\text{BN calculation} \rightarrow \frac{f_i - \bar{w}_i}{\sigma_i} = 0.3649 + \beta_1$$

$$\frac{6 - 4.6}{\sqrt{0.3}} = 0.90 \rightarrow 0.90 \times \beta_2 + \beta_3$$

PROBLEM \rightarrow does not work well with self-attention.

$$\text{BN calculation} \rightarrow \frac{f_i - \bar{w}_i}{\sigma_i} = 0.3649 + \beta_1$$

$$\frac{6 - 4.6}{\sqrt{0.3}} = 0.90 \rightarrow 0.90 \times \beta_2 + \beta_3$$

$$\text{BN calculation} \rightarrow \frac{f_i - \bar{w}_i}{\sigma_i} = 0.3649 + \beta_1$$

$$\frac{6 - 4.6}{\sqrt{0.3}} = 0.90 \rightarrow 0.90 \times \beta_2 + \beta_3$$

$$\text{BN calculation} \rightarrow \frac{f_i - \bar{w}_i}{\sigma_i} = 0.3649 + \beta_1$$

$$\frac{6 - 4.6}{\sqrt{0.3}} = 0.90 \rightarrow 0.90 \times \beta_2 + \beta_3$$

$$\text{BN calculation} \rightarrow \frac{f_i - \bar{w}_i}{\sigma_i} = 0.3649 + \beta_1$$

$$\frac{6 - 4.6}{\sqrt{0.3}} = 0.90 \rightarrow 0.90 \times \beta_2 + \beta_3$$

$$\text{BN calculation} \rightarrow \frac{f_i - \bar{w}_i}{\sigma_i} = 0.3649 + \beta_1$$

$$\frac{6 - 4.6}{\sqrt{0.3}} = 0.90 \rightarrow 0.90 \times \beta_2 + \beta_3$$

$$\text{BN calculation} \rightarrow \frac{f_i - \bar{w}_i}{\sigma_i} = 0.3649 + \beta_1$$

$$\frac{6 - 4.6}{\sqrt{0.3}} = 0.90 \rightarrow 0.90 \times \beta_2 + \beta_3$$

$$\text{BN calculation} \rightarrow \frac{f_i - \bar{w}_i}{\sigma_i} = 0.3649 + \beta_1$$

$$\frac{6 - 4.6}{\sqrt{0.3}} = 0.90 \rightarrow 0.90 \times \beta_2 + \beta_3$$

$$\text{BN calculation} \rightarrow \frac{f_i - \bar{w}_i}{\sigma_i} = 0.3649 + \beta_1$$

$$\frac{6 - 4.6}{\sqrt{0.3}} = 0.90 \rightarrow 0.90 \times \beta_2 + \beta_3$$



Transformer Architecture

Input

Output

Probabilities

Softmax

Linear

Layer Norm

Positional Encoding

Inputs

Outputs

(softmax right)

Positional Encoding

Multi-Head Attention

DECODER

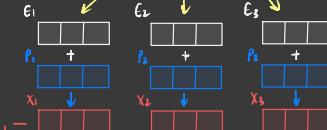
كيف حالك

Right shift

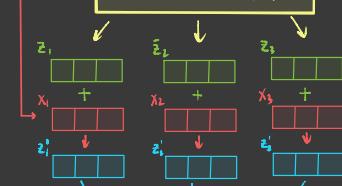
<start> كيف حالك

Tokenizer

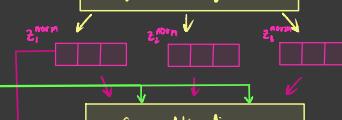
Embedding (512d)



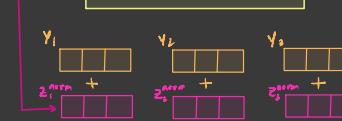
Masked Multi-head Att.



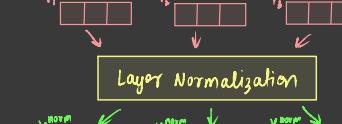
Layer Normalization



Masked Multi-head Att.



Layer Normalization



Layer Normalization

FFNN



Layer Normalization



Layer Normalization

* There are multiple encoders and decoders hence "How Are You" goes through all the encoders and the final output goes to every Decoder.

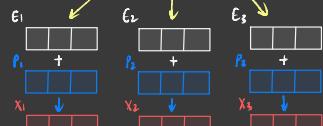
ENCODER

HOW ARE YOU

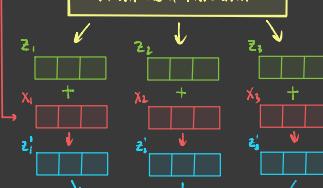
Tokenizer

How are yo

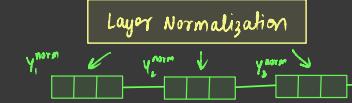
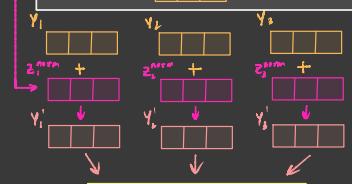
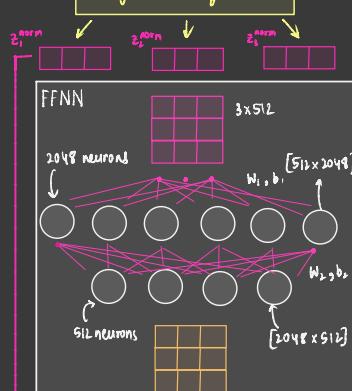
Embedding (512d)



Multi-head Att.



Layer Normalization



Calc. of PE vector Using Formula

$$PE_{(pos, 2)} = \sin(pos / 10000^{2/4m})$$

$$PE_{(pos, 4)} = \cos(pos / 10000^{2/4m})$$



$$pos \rightarrow (0) \quad (1) \quad (2)$$

$i=0$
$PE_{(0,0)} = \sin(0) = 0$
$PE_{(0,1)} = \cos(0) = 1$
$i=1$
$PE_{(1,0)} = \sin(0) = 0$
$PE_{(1,1)} = \cos(0) = 1$
$i=2$
$PE_{(2,0)} = \sin(0) = 0$
$PE_{(2,1)} = \cos(0) = 1$

$A_m \rightarrow i=0$
$PE_{(1,0)} = \sin\left(\frac{1}{10000^0}\right) = 0.841$
$PE_{(1,1)} = \cos\left(\frac{1}{10000^0}\right) = 0.540$
$i=1$
$PE_{(1,2)} = \sin\left(\frac{1}{10000^{2/4}}\right) = 0.009$
$PE_{(1,3)} = \cos\left(\frac{1}{10000^{2/4}}\right) = 0.995$
$i=2$
$PE_{(1,4)} = \sin\left(\frac{1}{10000^{4/4}}\right) = 0.00001$
$PE_{(1,5)} = \cos\left(\frac{1}{10000^{4/4}}\right) = 0.99999$

<u>Nauman</u> $\rightarrow i=0$	$i=1$
$PE_{(2,0)} = \sin\left(\frac{2}{10000^0}\right)$	$PE_{(2,1)} = \sin\left(\frac{2}{10000^{2/4}}\right) = 0.002$
$PE_{(2,1)} = \cos\left(\frac{2}{10000^0}\right)$	$PE_{(2,2)} = \cos\left(\frac{2}{10000^{2/4}}\right) = 0.99999$
$i=2$	
$PE_{(2,4)} = \sin\left(\frac{2}{10000^{4/4}}\right) = 0.00002$	
$PE_{(2,5)} = \cos\left(\frac{2}{10000^{4/4}}\right) = 0.9999999$	