$S_1$: Hello world haha

$S_2$: How old are you

$S_3$: Hey you

d_model = 4 , max_len = 4

$S_1$ embedding :

$$
\begin{array}{c}
\text{hello} \\
\text{world} \\
\text{haha} \\
<\cdot>
\end{array}
\overset{\xleftarrow{\quad d\_model \quad}}{
\begin{bmatrix}
1 & 3 & 4 & 1 \\
3 & 2 & 1 & 0 \\
4 & 5 & 7 & 6 \\
-1 & -1 & -1 & -1
\end{bmatrix}}
$$

$S_2$ embedding :

$$
\begin{array}{c}
\text{how} \\
\text{old} \\
\text{are} \\
\text{you}
\end{array}
\begin{bmatrix}
1 & 4 & 6 & 1 \\
3 & 1 & 5 & 4 \\
1 & 10 & 20 & 10 \\
1 & 2 & 0 & 1
\end{bmatrix}
$$

$S_3$ embedding :

$$
\begin{array}{c}
\text{Hey} \\
\text{you} \\
<\cdot> \\
<\cdot>
\end{array}
\begin{bmatrix}
3 & 1 & 4 & 5 \\
1 & 2 & 0 & 1 \\
-1 & -1 & -1 & -1 \\
-1 & -1 & -1 & -1
\end{bmatrix}
$$

# Positional encoding:

$$PE_{(pos, 2i)} = \sin\left(\frac{pos}{1000^{2i/d\text{-model}}}\right)$$

Pos : $[0, 1, 2, 3]$

$i : 0$ to $3$   according to max-len

$PE_{(pos, i=0)} : [\sin(0), \sin(0), \sin(0), \sin(0)]$

$PE_{(pos, i=1)} = [\cos(\cdot), \cos(\cdot), \cos(\cdot), \cos(\cdot)]$

$PE_{(pos, i=2)} = [\sin(0), \sin(0), \sin(\cdot), \sin(\cdot)]$

$EPE_1:$

$$
\begin{array}{c}
\text{hello} \\
\text{world} \\
\text{haha} \\
<\cdot>
\end{array}
\overset{\xleftarrow{\hspace{1cm} d\text{-model} \hspace{1cm}}}{
\begin{bmatrix}
1.1 & 3.2 & 4.1 & 10 \\
3.2 & 2.1 & 1.2 & 0.1 \\
4.1 & 5.5 & 7.1 & 10 \\
-1 & -1 & -1 & -1
\end{bmatrix}}
$$

$\sum PE_2:$ $\overset{\text{Max}}{\text{Len}}$

$$
\begin{array}{c}
\text{how} \\
\text{old} \\
\text{are} \\
\text{you}
\end{array}
\begin{bmatrix}
1 & 4 & 6 & 1 \\
3 & 1 & 5 & 4 \\
1 & 10 & 20 & 10 \\
1 & 2 & 0 & 1
\end{bmatrix}
$$

$EPE_3:$

$$
\begin{array}{c}
\text{Hey} \\
\text{you} \\
<\cdot> \\
<\cdot>
\end{array}
\begin{bmatrix}
3 & 1 & 4 & 5 \\
1 & 2 & 0 & 1 \\
-1 & -1 & -1 & -1 \\
-1 & -1 & -1 & -1
\end{bmatrix}
$$

3

data shape: $(3, 4, 4)$ // samples

samples

# Feed-Forward :



token
embbeding
siz $n$

$x W + b$

$x G$

$d\text{-model} = 16$

$d\text{-ff} = 64$

input matrix
of words
**6 , 512**

A

B

C

D

E

F

X

512 , 6

A   B   C   D   E F

dot product $= \sum_{i=1}^{n} x_i y_i$

6 , 6

$=$

| A·A | A·B | A·C |  |  | A·f |
|-----|-----|-----|--|--|-----|
| B·A | B·B | B·C |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |

# Encoder:

Start with word embeddings

tokenize:

You    Cat    is    A    lonely    Cat

input IDS
(Positions
in vocab)

105    6587    5478    3578    65    6587

embedding
vector of
size 512

$$\begin{bmatrix} 982.1 \\ 331 \\ \vdots \\ 260.1 \end{bmatrix}$$
512,1

$$\begin{bmatrix} 112 \\ 310 \\ \vdots \\ 2011 \end{bmatrix}$$
512,1

$$\begin{bmatrix} 3 \\ 2 \\ 9 \\ 3 \\ 100 \end{bmatrix}$$
512,1

$$\begin{bmatrix} 105 \\ 120 \\ 3 \\ \vdots \\ 100 \end{bmatrix}$$
512,1

$$\begin{bmatrix} 2 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$
512,1

$$\begin{bmatrix} 3 \\ 3 \\ \vdots \\ 1 \end{bmatrix}$$
512,1

embedding
**layer**

# Positional encoding

Carry information about the position of the words

You    Cat    is    A    lonely    Cat

**Embedding layer**

$$\begin{bmatrix} 952.1 \\ 33.1 \\ \vdots \\ 260.1 \end{bmatrix}_{512,1}$$

$$\begin{bmatrix} 112 \\ 310 \\ \vdots \\ 2.11 \end{bmatrix}_{512,1}$$

$$\begin{bmatrix} 3 \\ 2 \\ 1 \\ \vdots \\ 100 \end{bmatrix}_{512,1}$$

$$\begin{bmatrix} 100 \\ 120 \\ 3 \\ \vdots \\ 100 \end{bmatrix}_{512,1}$$

$$\begin{bmatrix} 2 \\ 1 \\ 1 \\ 0 \end{bmatrix}_{512,1}$$

$$\begin{bmatrix} 3 \\ 3 \\ \vdots \\ 1 \end{bmatrix}_{512,1}$$

+    +    +    +    +    +

**Position embedding**

**not learned**

**Computed once**

$$\begin{bmatrix} 1010 \\ 2.0 \\ \vdots \\ 3 \end{bmatrix}_{512,1}$$

$$\begin{bmatrix} 4 \\ 5 \\ \vdots \\ 6 \end{bmatrix}_{512,1}$$

$$\begin{bmatrix} 700 \\ 800 \\ \vdots \\ 991 \end{bmatrix}_{512,1}$$

$$\begin{bmatrix} 990 \\ 50 \\ \vdots \\ 101 \end{bmatrix}_{512,1}$$

$$\begin{bmatrix} 2.1 \\ 320 \\ \vdots \\ 400 \end{bmatrix}_{512,1}$$

$$\begin{bmatrix} 1.1 \\ 3.0 \\ \vdots \\ 600 \end{bmatrix}_{512,1}$$

= = = = = =

encoder
input

$S12,1$    $S12,1$    $S12,1$    $S12,1$    $S12,1$    $S12,1$

How to calculate positional encodings?

assume input text:

**Your**      **Cat**      **is**

| $PE(0,0)$ | | $PE(1,0)$ | | $PE(2,0)$ |
| $PE(0,1)$ | | $PE(1,1)$ | | $PE(2,1)$ |
| $PE(0,2)$ | | $PE(1,2)$ | | $PE(2,2)$ |
| $\vdots$ | | $\vdots$ | | $\vdots$ |
| $PE(0,S10)$ | | $PE(1,S10)$ | | $PE(2,S10)$ |
| $PE(0,S11)$ | | $PE(1,S11)$ | | $PE(2,S11)$ |

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) : \text{even}$$

$$PE(pos, 2i+1) = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) : \text{odd}$$

every other sentence will have
the same Positional encodings

## Self_attention

allows to relate words to each other.

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right) \cdot V$$

$$\text{softmax}\left[ \begin{array}{c} Q \quad (6, 512) \\ \end{array} \right. \quad X \quad \left. K^T \quad (512,6) \right]$$

$d_K = 512$

# (6×6) matrix

| | You | Cat | is | a | lovely | Cat |
|---|---|---|---|---|---|---|
| **your** | 0.268 | 0.119 | 0.124 | 0.148 | 0.174 | 0.152 |
| **Cat** | 0.124 | 0.278 | 0.21 | 0.128 | 0.154 | 0.115 |
| **is** | 0.147 | 0.132 | 0.262 | 0.097 | 0.218 | 0.145 |
| **a** | 0.210 | 0.128 | 0.206 | 0.212 | 0.119 | 0.125 |
| **lovely** | 0.146 | 0.158 | 0.152 | 0.143 | 0.227 | 0.174 |
| **Cat** | 0.195 | 0.114 | 0.203 | 0.203 | 0.157 | 0.229 |

$\sum = 1$

the values represent the
the relation between words
low intense
between words

## (6×6) matrix

| | You | Cat | is | a | lovely | Cat |
|---|---|---|---|---|---|---|
| You | 0.268 | 0.119 | 0.134 | 0.148 | 0.174 | 0.152 |
| Cat | 0.124 | 0.278 | 0.21 | 0.128 | 0.154 | 0.115 |
| is | 0.147 | 0.132 | 0.262 | 0.097 | 0.218 | 0.145 |
| a | 0.210 | 0.128 | 0.206 | 0.212 | 0.114 | 0.125 |
| lovely | 0.146 | 0.158 | 0.152 | 0.143 | 0.227 | 0.174 |
| Cat | 0.195 | 0.114 | 0.203 | 0.203 | 0.157 | 0.229 |

$\sum = 1$

$$X$$

6.512

$$V \qquad =$$

$z$

$6, 512$

attention

$$embedding \begin{bmatrix} [[1,2,1], & [3,1,2], [3,1,3]] \\ [[2,1,0], [3,4,1], [2,1,0]] \\ [[3,0,0], [4,1,2], [3,3,3]] \\ [[1,0,1], [2,3,7], [7,1,0]] \end{bmatrix}$$

4, 3, 3

$$PE = \begin{bmatrix} [[1,0,0], [2,1,3], [4,1,2]] \end{bmatrix}$$

1, 3, 3

# Multi-head Attention

d-model = 4

num heads = 2

$$\text{depth} = \frac{d\text{-model}}{num\text{-heads}} = 2$$

assume seq-len = 3

$$q = \begin{bmatrix} \begin{bmatrix} [1,0,0,1] \\ [0,1,1,0] \\ [1,1,1,1] \end{bmatrix} \end{bmatrix}$$

1 x 3 x 4

head for q: will be shape (1, 2, 3, 2)

batch-size, heads, seq-len, depth

$$q\text{-heads} = \begin{bmatrix} \begin{bmatrix} [1,0] \\ [0,1] \\ [1,1] \end{bmatrix} , \begin{bmatrix} [0,1] \\ [1,0] \\ [1,1] \end{bmatrix} \end{bmatrix}$$

# Steps:

**(1)** embeddings + PE:

shape: (batch_size, seq_len, d_model)

**(2)** Multi head attention:

$$q = dense(d\_model) \quad // \quad q \cdot w_q$$

$$v = dense(d\_model) \quad // \quad v \cdot w_v$$

$$k = dense(d\_model) \quad // \quad k \cdot w_k$$

shapes:

(batch_size, seq_len, d_model)

- split into heads:

q-heads
k-heads
v-heads

$\Bigg\{$ (batch_size, num_heads, seq_len, depth)

— dt. Product attention:

(*) $A = (q\text{-heads} \cdot k\text{-heads}^T) / \sqrt{d_k}$

<span style="color:red">shape $(k)[-1]$</span>

(*) $S = \text{softmax}(A)$

(*) $S \cdot v\text{-heads}$