

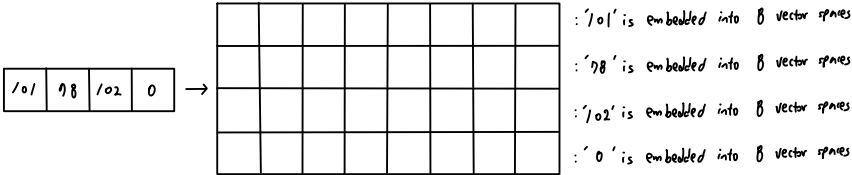
- ① Input : (Batch, Sequence length)
- ② Embedding: (Batch, Sequence length, d_{model})
- ③ Positional Encoding: (Batch, Sequence length, d_{model})

Let Batch = 1, Sequence length = 4, $d_{\text{model}} = 8$

Input size : (1, 4)

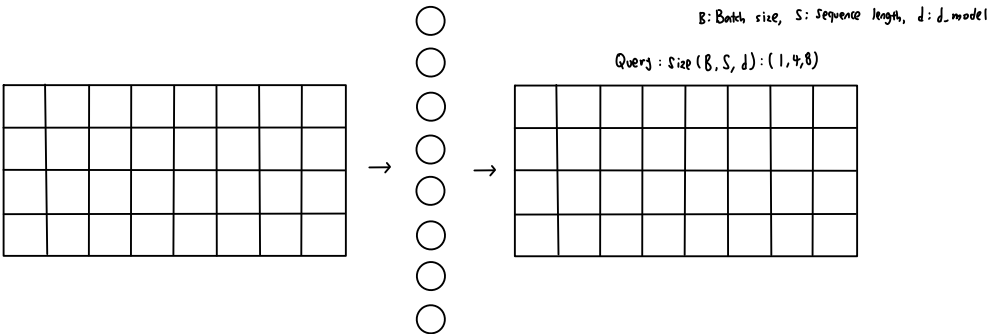
/0/	98	/02	0
-----	----	-----	---

After Input Embedding & Positional Encoding : (1, 4, 8)

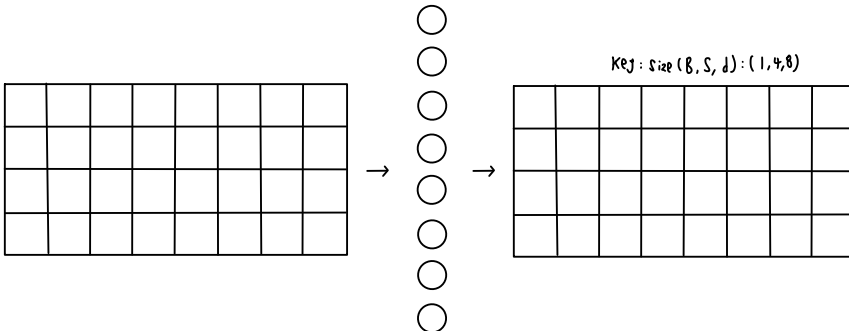


Multi-Head Attention

Query MLP



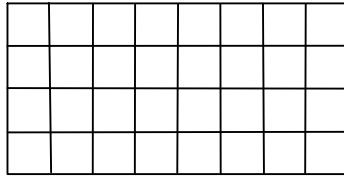
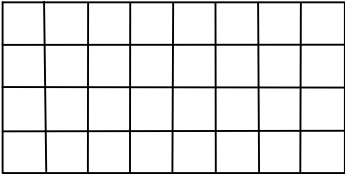
Key MLP



Value MUX

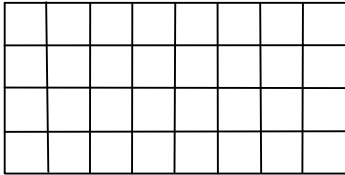


Value : $\text{Size}(B, S, d) : (1, 4, 8)$

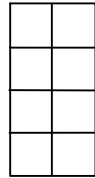


$\text{Size}(B \times h, S, d/h) : (4, 4, 2)$, h : Number of heads

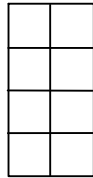
Query : $\text{Size}(B, S, d) : (1, 4, 8)$



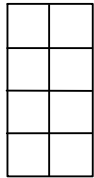
Query's Head 1



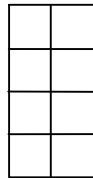
Query's Head 2



Query's Head 3

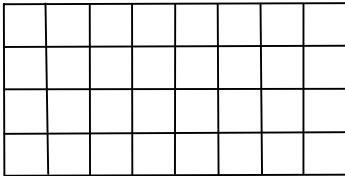


Query's Head 4

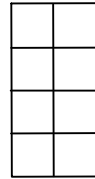


$(4, 4, 2)$

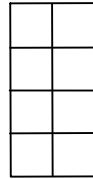
Key : $\text{Size}(B, S, d) : (1, 4, 8)$



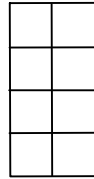
Head 1



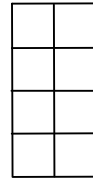
Head 2



Head 3

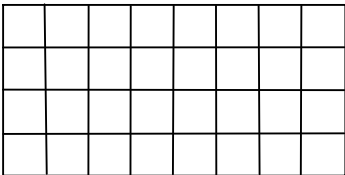


Head 4

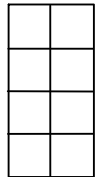


$(4, 4, 2)$

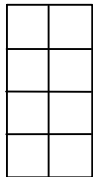
Value : $\text{Size}(B, S, d) : (1, 4, 8)$



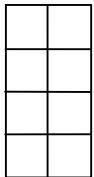
Head 1



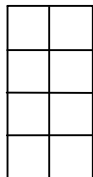
Head 2



Head 3



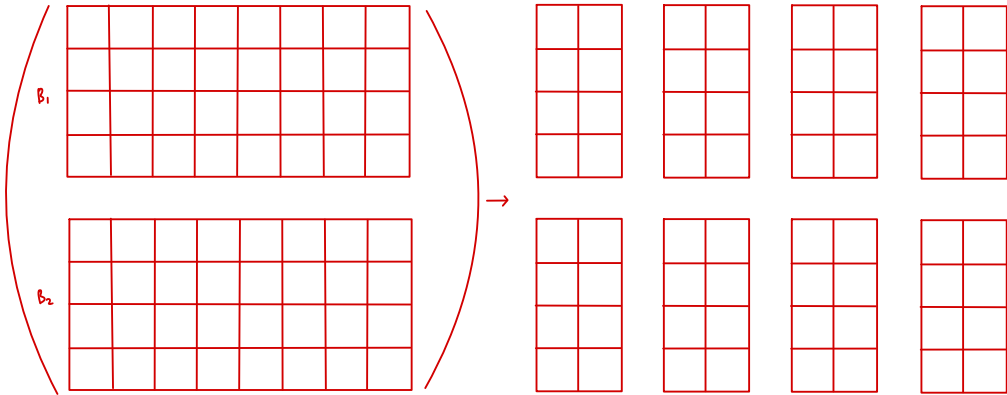
Head 4



If $B=2$?

$(B, S, d) : (2, 4, 8)$

$(B, 4, 2)$

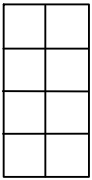


Calculating Attention (W) in Multi-Head Attention

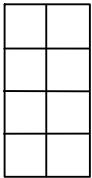
① $W = \text{Query} \cdot \text{Key.T}$

Query : $(4, 4, 2)$

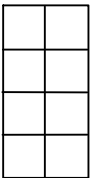
Query's Head 1



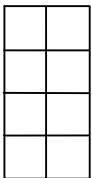
Query's Head 2



Query's Head 3

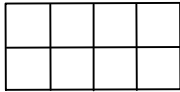


Query's Head 4

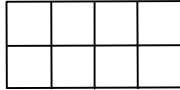


Key.T : $(4, 2, 4)$

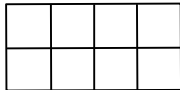
Key's Head 1



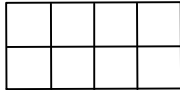
Key's Head 2



Key's Head 3



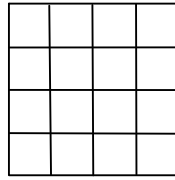
Key's Head 4



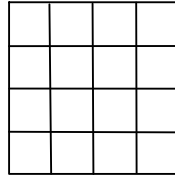
=

$W : (B^{\text{th}}, S, 5) : \text{Attention for each word in sentence.}$

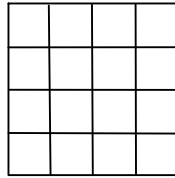
$W^1 : \text{Head 1}$



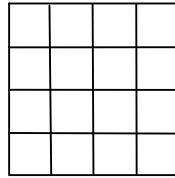
$W^2 : \text{Head 2}$



$W^3 : \text{Head 3}$



$W^4 : \text{Head 4}$



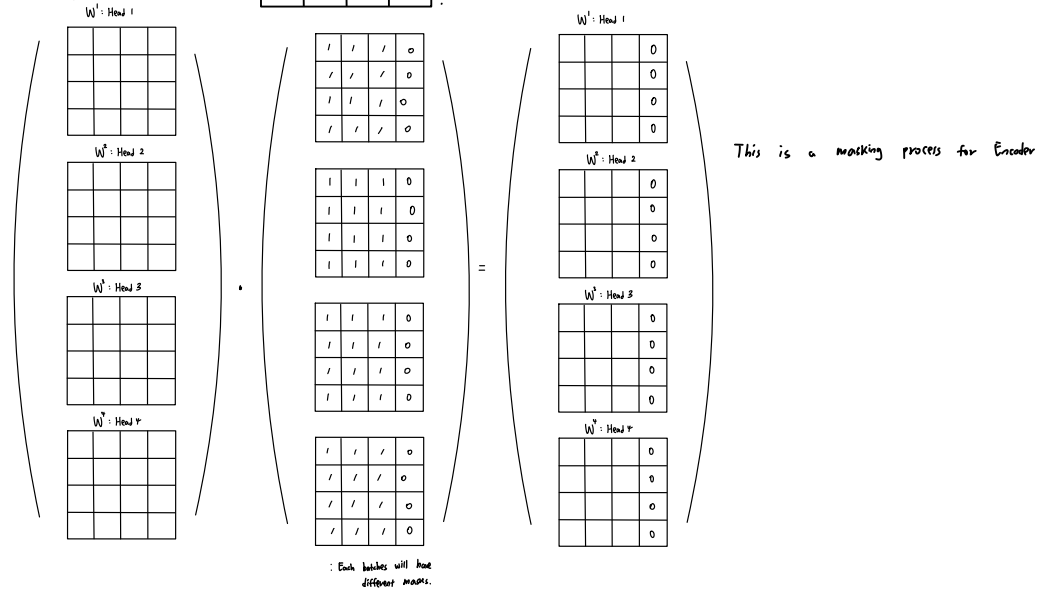
Each head will have different attention values.

From "Attention is All You Need!",

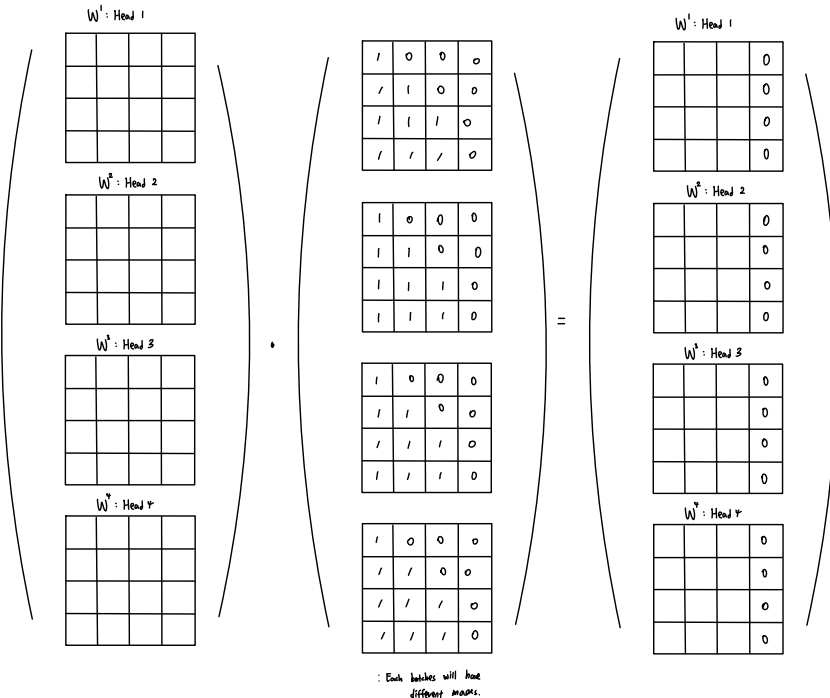
"We found it beneficial to linearly project the queries, keys, and values h times with different, learned linear projections to d_k, d_q , and d_v ."

② $W = W/\sqrt{d}$

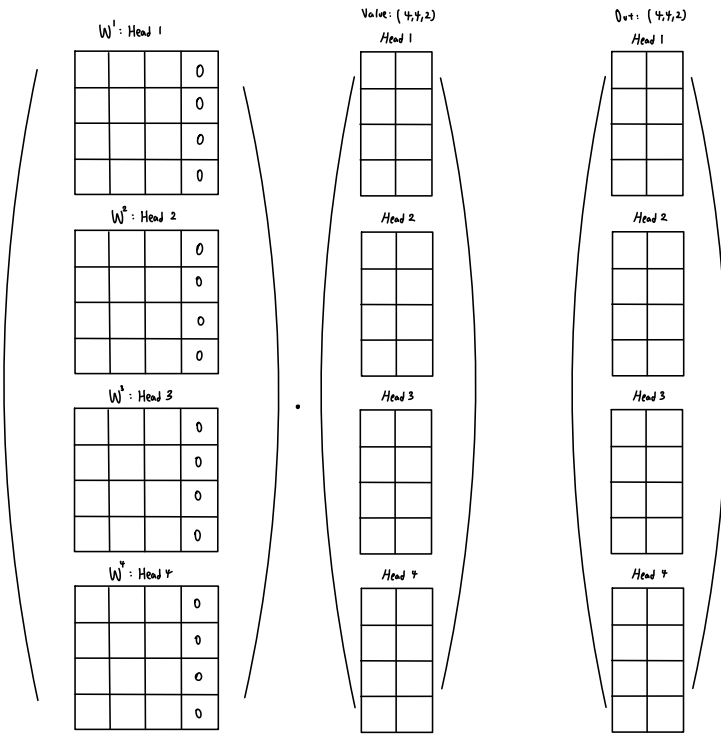
② Masking W : Let's say mask $\begin{bmatrix} / & / & / & 0 \end{bmatrix}$. For masking W , we need same size.



③ - 1: Causally Masking for Decoder. Let's say mask $\begin{bmatrix} / & / & / & 0 \end{bmatrix}$. For masking W , we need same size.



④ Out = W Value



⑤ Reshape output. : (B, S, d)

