

Let model dimensions (embedding dims) be 512.
 consider an input sentence with 5 words.

$$\tilde{I} \in \mathbb{R}^{5 \times 512}$$

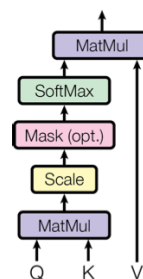
denote randomly initialized query, key, and value matrices by $\tilde{Q}_w, \tilde{K}_w, \tilde{V}_w$. Each of $\tilde{Q}_w, \tilde{K}_w, \tilde{V}_w \in \mathbb{R}^{512 \times 512}$. To obtain queries (\tilde{Q}), keys (\tilde{K}), and values multiply the input (\tilde{I}) with $\tilde{Q}_w, \tilde{K}_w, \tilde{V}_w$. (Assume that we have only one head)

$$\underbrace{\tilde{I}}_{\mathbb{R}^{5 \times 512}} \times \underbrace{\tilde{Q}_w}_{\mathbb{R}^{512 \times 512}} = \underbrace{\tilde{Q}}_{\mathbb{R}^{5 \times 512}}$$

$$\underbrace{\tilde{I}}_{\mathbb{R}^{5 \times 512}} \times \underbrace{\tilde{K}_w}_{\mathbb{R}^{512 \times 512}} = \underbrace{\tilde{K}}_{\mathbb{R}^{5 \times 512}}$$

$$\underbrace{\tilde{I}}_{\mathbb{R}^{5 \times 512}} \times \underbrace{\tilde{V}_w}_{\mathbb{R}^{512 \times 512}} = \underbrace{\tilde{V}}_{\mathbb{R}^{5 \times 512}}$$

Computational graph of the scaled dot product attention is given on the right, we will look into it in depth.



$$\tilde{Q} \in \mathbb{R}^{5 \times 512}, \tilde{K} \in \mathbb{R}^{5 \times 512}, \tilde{V} \in \mathbb{R}^{5 \times 512}$$

$$\text{Attn}(\tilde{Q}, \tilde{K}, \tilde{V}) = \frac{\text{Softmax}(\tilde{Q} \tilde{K}^T, \text{axis}=1) \times \tilde{V}}{\sqrt{d_k}} \quad d_k=512$$

$$\tilde{Q} \tilde{K}^T =$$

$$\begin{bmatrix} q_{11} & \dots & q_{1,512} \\ \vdots & & \vdots \\ q_{51} & \dots & q_{5,512} \end{bmatrix} \begin{bmatrix} k_{11} & \dots & k_{15} \\ \vdots & & \vdots \\ k_{512,1} & \dots & k_{512,5} \end{bmatrix} = \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

$\mathbb{R}^{5 \times 512} \quad \mathbb{R}^{512 \times 5} \quad \mathbb{R}^{5 \times 5}$

let each row of \tilde{Q} be \tilde{q}_i
and each column of \tilde{K} be \tilde{k}_i^T

$$\begin{bmatrix} \tilde{q}_1 \\ \vdots \\ \tilde{q}_5 \end{bmatrix} \begin{bmatrix} \tilde{k}_1^T & \dots & \tilde{k}_5^T \end{bmatrix} = \begin{bmatrix} \tilde{q}_1 \tilde{k}_1^T & \tilde{q}_1 \tilde{k}_2^T & \tilde{q}_1 \tilde{k}_3^T & \tilde{q}_1 \tilde{k}_4^T & \tilde{q}_1 \tilde{k}_5^T \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \tilde{q}_5 \tilde{k}_1^T & \tilde{q}_5 \tilde{k}_2^T & \tilde{q}_5 \tilde{k}_3^T & \tilde{q}_5 \tilde{k}_4^T & \tilde{q}_5 \tilde{k}_5^T \end{bmatrix}$$

$\mathbb{R}^{5 \times 512} \quad \mathbb{R}^{512 \times 5} \quad \mathbb{R}^{5 \times 5} \rightarrow \tilde{A}$

each of the $\tilde{q}_i \tilde{k}_j^T$ yields a scalar as it is a dot

product

$$1 \leq i \leq 5, 1 \leq j \leq 5$$

Now, remember that $\tilde{V} \in \mathbb{R}^{5 \times 12}$ row point of view
(see end for more)

$$\begin{bmatrix} \tilde{A} \\ \tilde{V} \end{bmatrix} = \begin{bmatrix} \tilde{A} \\ \tilde{V}_1 \\ \vdots \\ \tilde{V}_5 \end{bmatrix} = \begin{bmatrix} \tilde{A} \\ \tilde{V}_1 + \tilde{K}_1^T \tilde{V}_1 + \tilde{K}_2^T \tilde{V}_2 + \tilde{K}_3^T \tilde{V}_3 + \tilde{K}_4^T \tilde{V}_4 + \tilde{K}_5^T \tilde{V}_5 \\ \vdots \\ \tilde{V}_5 + \tilde{K}_1^T \tilde{V}_1 + \tilde{K}_2^T \tilde{V}_2 + \tilde{K}_3^T \tilde{V}_3 + \tilde{K}_4^T \tilde{V}_4 + \tilde{K}_5^T \tilde{V}_5 \end{bmatrix} \in \mathbb{R}^{5 \times 12}$$

$\mathbb{R}^{5 \times 5} \times \mathbb{R}^{5 \times 12}$ looks like a $\mathbb{R}^{5 \times 1}$ = $\mathbb{R}^{5 \times 1}$ (looks like)

Note that we have not done softmax yet on the attention matrix (\tilde{A}), $\tilde{A} \in \mathbb{R}^{5 \times 5}$ and now we have to do a softmax operation on \tilde{A} .

Each row of \tilde{A} linearly scales each row of \tilde{V} to obtain a scaled version of the input, this scaled version of the input (\tilde{O}) contains the attention information because \tilde{V} is scaled by the attention matrix (\tilde{A}).

To obtain \tilde{O} , we must take softmax (\tilde{A} , axis=1), that is softmax operation is performed individually for each of the rows.

to obtain attention matrix $A \in \mathbb{R}^{S \times S}$ from

\tilde{A} . We perform

$$\left[\begin{array}{ccc} \frac{\tilde{q}_1 \tilde{k}_1^T}{e^{\sum_{j=1}^S \tilde{q}_1 \tilde{k}_j^T}} & \dots & \frac{\tilde{q}_1 \tilde{k}_5^T}{e^{\sum_{j=1}^S \tilde{q}_1 \tilde{k}_j^T}} \\ \vdots & & \vdots \\ \frac{\tilde{q}_5 \tilde{k}_1^T}{e^{\sum_{j=1}^S \tilde{q}_5 \tilde{k}_j^T}} & \dots & \frac{\tilde{q}_5 \tilde{k}_5^T}{e^{\sum_{j=1}^S \tilde{q}_5 \tilde{k}_j^T}} \end{array} \right]$$

Matrix Multiplication Interpretations

I) Column point of view

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{pmatrix} = \begin{pmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \\ c_{31} & c_{32} \end{pmatrix}$$

Columns of C are linear combination of columns of A with weights given by columns of B .

$$\begin{pmatrix} c_{11} \\ c_{21} \\ c_{31} \end{pmatrix} = \begin{pmatrix} a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31} \\ a_{21}b_{11} + a_{22}b_{21} + a_{23}b_{31} \\ a_{31}b_{11} + a_{32}b_{21} + a_{33}b_{31} \end{pmatrix} = \begin{pmatrix} a_{11} \\ a_{21} \\ a_{31} \end{pmatrix} b_{11} + \begin{pmatrix} a_{12} \\ a_{22} \\ a_{32} \end{pmatrix} b_{21} + \begin{pmatrix} a_{13} \\ a_{23} \\ a_{33} \end{pmatrix} b_{31}$$

II) Row Point of view
 Rows of C are linear combinations of the rows of B
 with weights given by rows of A.

$$\begin{matrix} & A & & B & & C \\ \begin{pmatrix} \boxed{a_{11} \ a_{12} \ a_{13}} \\ a_{21} \ a_{22} \ a_{23} \\ a_{31} \ a_{32} \ a_{33} \end{pmatrix} & \begin{pmatrix} \boxed{b_{11} \ b_{12}} \\ b_{21} \ b_{22} \\ b_{31} \ b_{32} \end{pmatrix} & = & \begin{pmatrix} c_{11} \ c_{12} \\ c_{21} \ c_{22} \\ c_{31} \ c_{32} \end{pmatrix}
 \end{matrix}$$

$$\begin{aligned}
 (c_{11} \ c_{12}) &= (a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31} \quad a_{11}b_{12} + a_{12}b_{22} + a_{13}b_{32}) \\
 &= (b_{11} \ b_{12})a_{11} + (b_{21} \ b_{22})a_{12} + (b_{31} \ b_{32})a_{13}
 \end{aligned}$$