

parameters W_q, W_k, W_v

$$3(6 \times 64) \times 4 + (6 \times 120) \times 2$$

3) Calculate the number of parameters in a single Transformer encoder layer given the following:

Input dimension: 64

Number of heads in multi-head attention: 4

Dimension of each head: 16

Dimension of feed-forward network: 120

Assume the weight matrices for the linear transformations in multi-head attention and the feed-forward network are the primary contributors to the parameter count.

27648

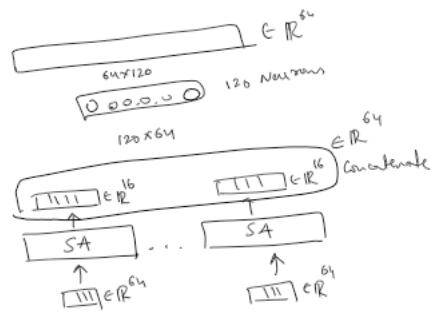
Yes, the answer is correct.

Score: 1

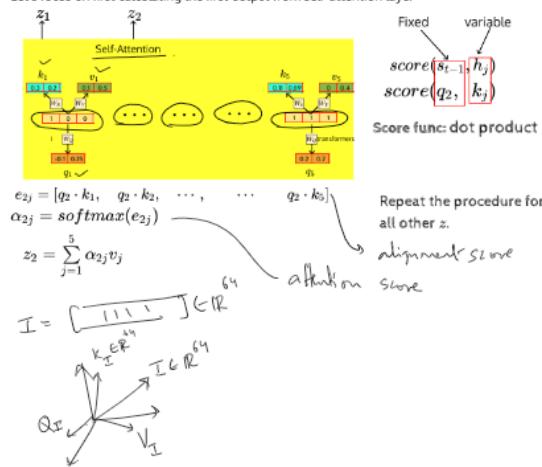
Accepted Answers:

(Type Numeric) 27648

$$\begin{aligned} Q & \in \mathbb{R}^{16 \times 1} & W_q & = \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \in \mathbb{R}^{16 \times 64} & \text{encoder} \\ & \quad \left\{ \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \right\} \in \mathbb{R}^{6 \times 1} & & & \\ \text{by } k & \in \mathbb{R}^{16 \times 1} & & \quad \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \in \mathbb{R}^{16 \times 64} & \quad \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \in \mathbb{R}^{6 \times 1} \\ V & \in \mathbb{R}^{16 \times 1} & & \quad \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \in \mathbb{R}^{16 \times 64} & \quad \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \in \mathbb{R}^{6 \times 1} \end{aligned}$$



Let's focus on first calculating the first output from self-attention layer



$$\begin{aligned} 64 \left\{ \begin{matrix} q_1 \\ q_2 \\ \vdots \\ q_n \end{matrix} \right\} &= \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \in \mathbb{R}^{64 \times 64} & W_Q & \left[\begin{matrix} h_1 \\ h_2 \\ \vdots \\ h_5 \end{matrix} \right] \in \mathbb{R}^{64 \times 5} & \text{lets say each word } \in \mathbb{R}^{64} \\ Q & \in \mathbb{R}^n & Q & \in \mathbb{R}^{64 \times 7} & \\ 64 \left\{ \begin{matrix} k_1 \\ k_2 \\ \vdots \\ k_n \end{matrix} \right\} &= \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \in \mathbb{R}^{64 \times 64} & W_K & \left[\begin{matrix} h_1 \\ h_2 \\ \vdots \\ h_5 \end{matrix} \right] \in \mathbb{R}^{64 \times 5} & \\ K & \in \mathbb{R}^n & K & \in \mathbb{R}^{64 \times 7} & \\ 64 \left\{ \begin{matrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{matrix} \right\} &= \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \in \mathbb{R}^{64 \times 64} & W_V & \left[\begin{matrix} h_1 \\ h_2 \\ \vdots \\ h_5 \end{matrix} \right] \in \mathbb{R}^{64 \times 5} & \end{aligned}$$

- Sequence Length : $t=32$
- Number of Heads : $h=2$
- Embedding dimension : d_{model}
- Input $X \in \mathbb{R}^{d_{\text{model}} \times t}$
- $d_k = d_q = \frac{d_{\text{model}}}{h} = 32$
- $W_Q \in \mathbb{R}^{d_k \times d_{\text{model}}}$
- $W_K \in \mathbb{R}^{d_k \times d_{\text{model}}}$
- $W_V \in \mathbb{R}^{d_v \times d_{\text{model}}}$
- $W_o \in \mathbb{R}^{d_{\text{model}} \times (h \times d_v)}$
 $\in \mathbb{R}^{64 \times (2 \times 16)}$ 64×32

Suppose $t = 32$, $d_{\text{model}} = 64$, $h = 2$ and $d_v = 16$. What will be the shape of the output of the scaled dot-product attention operation for a single head, given by:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{Q^T K}{\sqrt{d_k}}\right) V^T$$

Compute the resulting output dimension and report the total number of elements in the resulting attention output.

$$Q = \begin{bmatrix} \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots \end{bmatrix}_{32 \times 32} = \begin{bmatrix} W_Q \\ \cdots \\ \cdots \end{bmatrix}_{32 \times 64} \begin{bmatrix} \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots \end{bmatrix}_{64 \times 32}$$

 $K = \begin{bmatrix} \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots \end{bmatrix}_{32 \times 32} = \begin{bmatrix} W_K \\ \cdots \\ \cdots \end{bmatrix}_{32 \times 64} \begin{bmatrix} \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots \end{bmatrix}_{64 \times 32}$

 $V = \begin{bmatrix} \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots \end{bmatrix}_{16 \times 32} = \begin{bmatrix} W_V \\ \cdots \\ \cdots \end{bmatrix}_{16 \times 64} \begin{bmatrix} \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots \end{bmatrix}_{64 \times 32}$

