

Attention

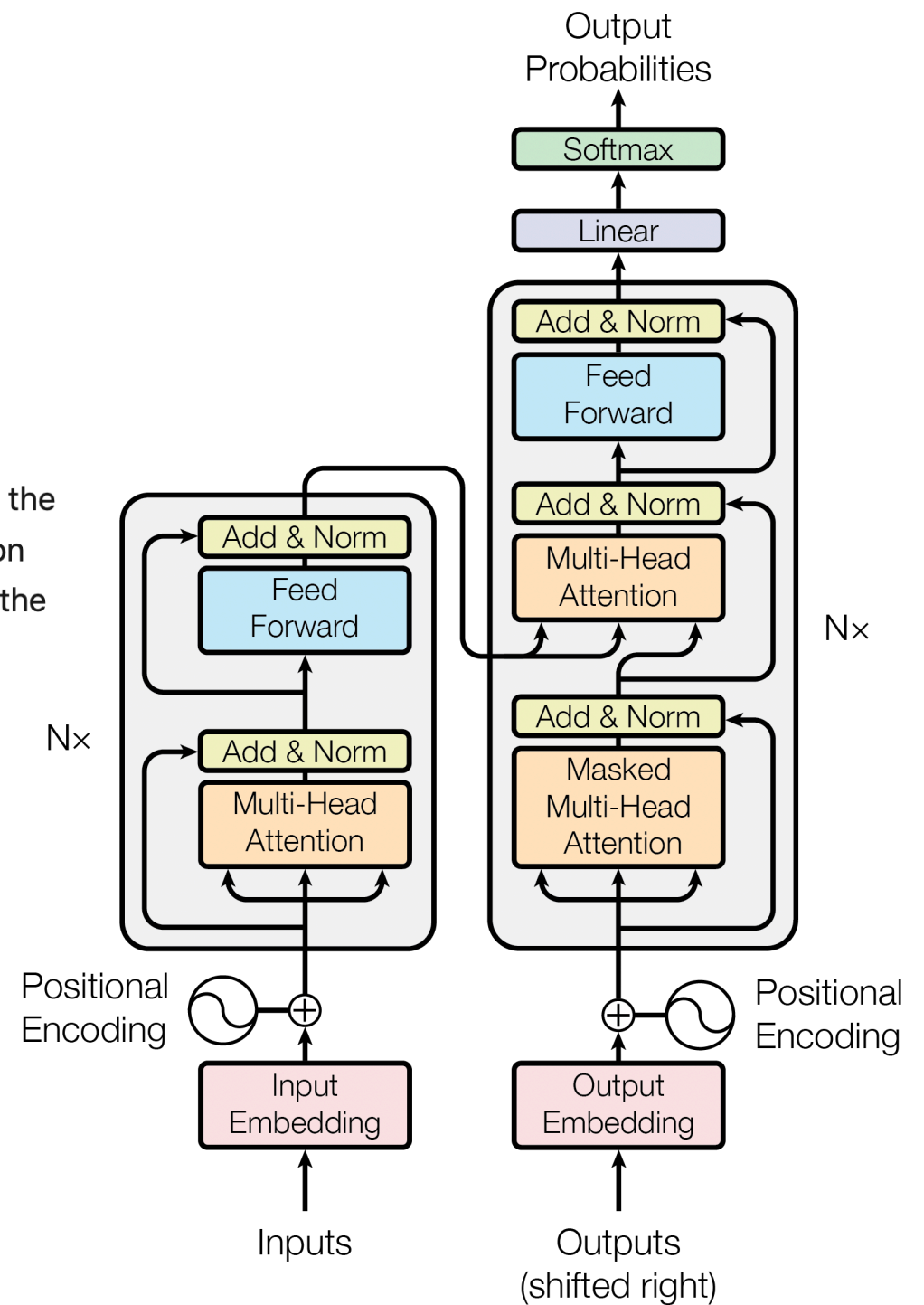
Introduction

$$h_0 = UW_{\text{embed}} + W_{\text{position}}$$

$$h_l = \text{transformer-block}(h_{l-1}) \forall l \in [1, n]$$

$$P(u) = \text{softmax}(h_n W_{\text{embed}}^T)$$

where $U = (u_{-k}, \dots, u_{-1})$ is the context vector of tokens (one hot encodings), n is the number of layers, W_{embed} is the token embedding matrix, and W_{position} is the position embedding matrix. Let d_n be the length of the input and d_{model} be the dimension of the embedding. $W_{\text{embed}} \in \mathbb{R}^{d_{\text{ntokens}} \times d_{\text{model}}}$ and $U \in \mathbb{R}^{d_n \times d_{\text{ntokens}}}$. So $h_0 \in \mathbb{R}^{d_n \times d_{\text{model}}}$.



TransFormer Block

transformer-block:

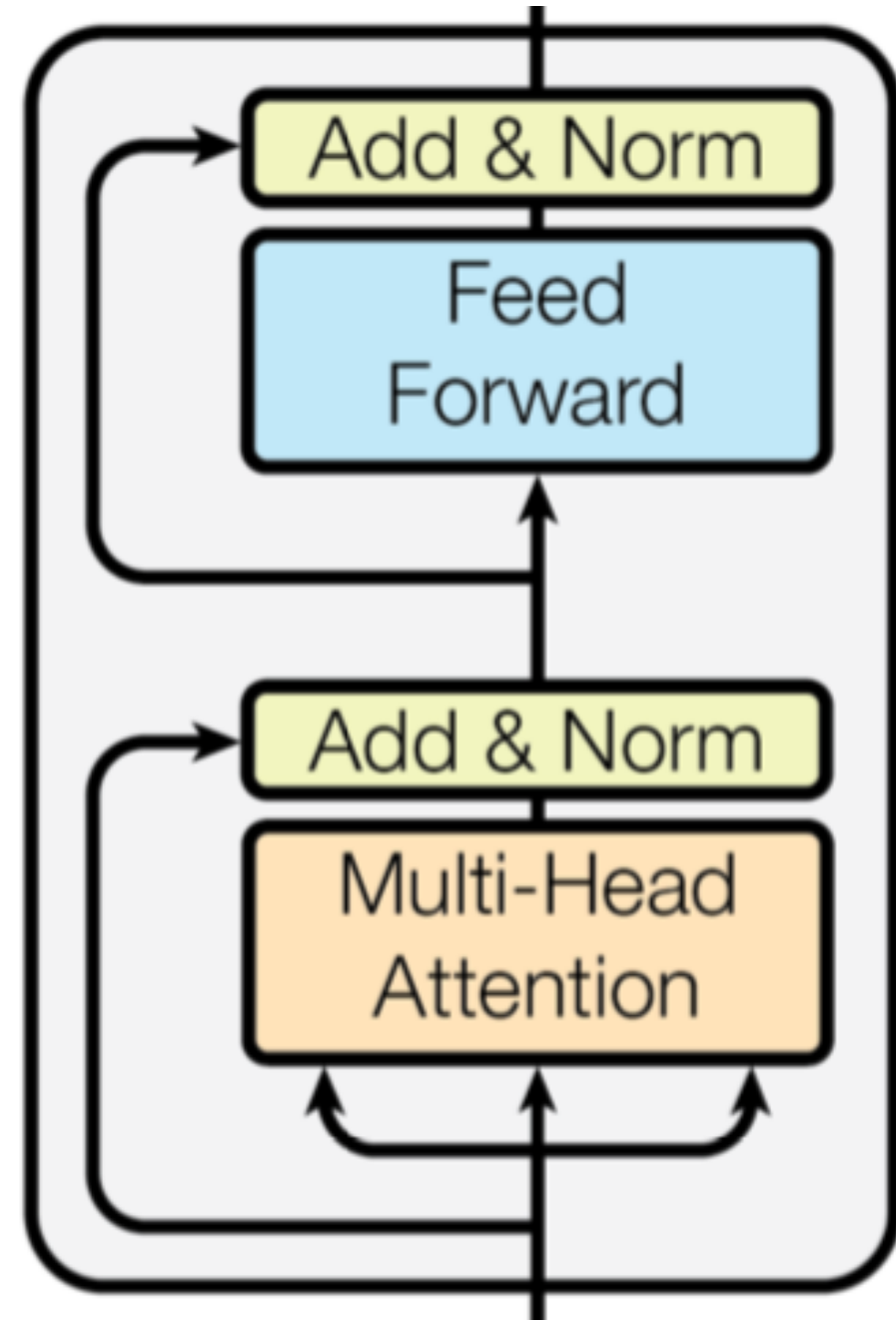
input: h_{in}

output: h_{out}

$h_{mid} = \text{LayerNorm}(h_{in} + \text{MultiHead}(h_{in}))$

$h_{out} = \text{LayerNorm}(h_{mid} + \text{FFN}(h_{mid}))$

where $h_{in}, h_{out} \in \mathbb{R}^{d_n \times d_{\text{model}}}$. d_n is the length of the input and d_{model} is the dimension of the model (e.g., embedding dimension). LayerNorm is layer normalization.



Attention

Scaled Dot-Product Attention

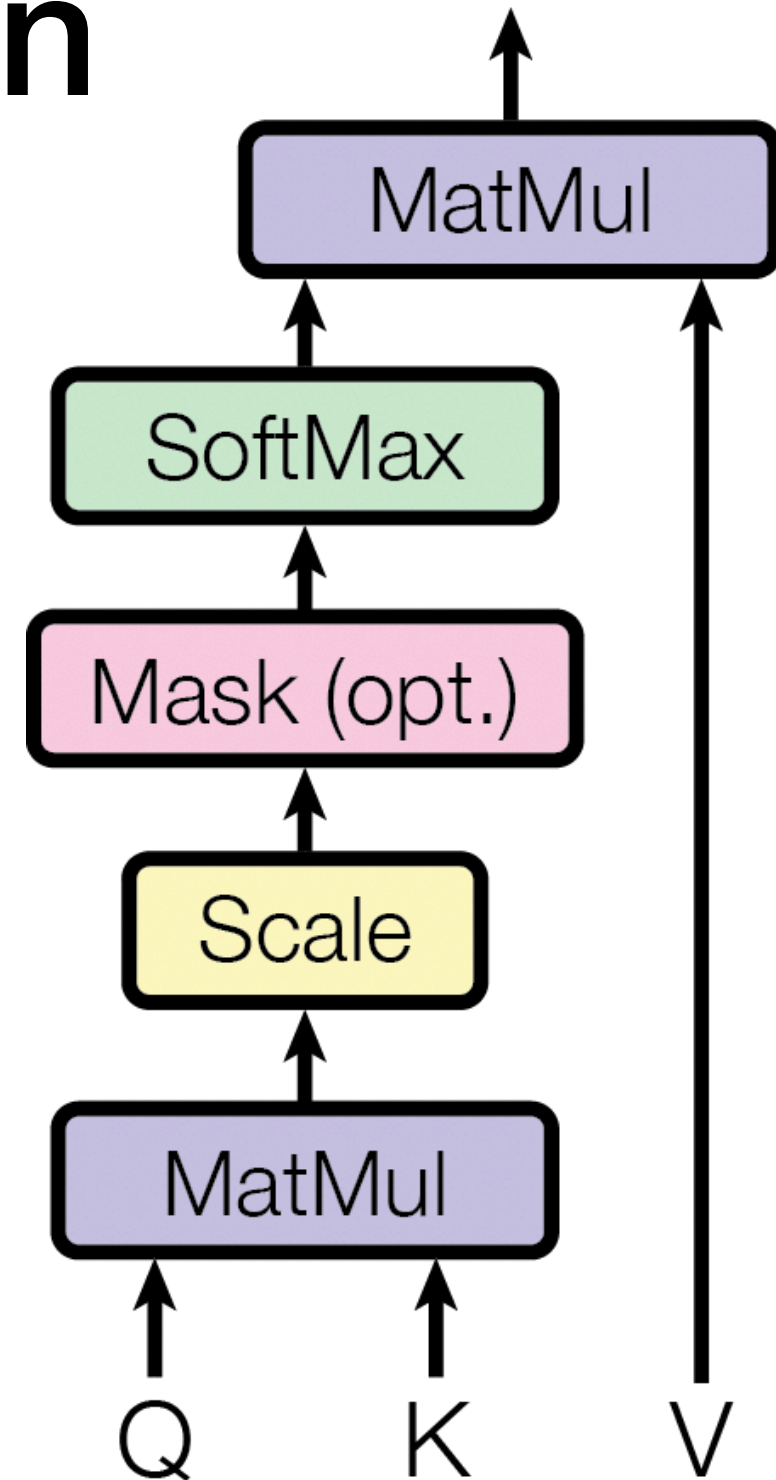
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

where the $Q, K \in \mathbb{R}^{d_n \times d_k}$ and $V \in \mathbb{R}^{d_n \times d_v}$

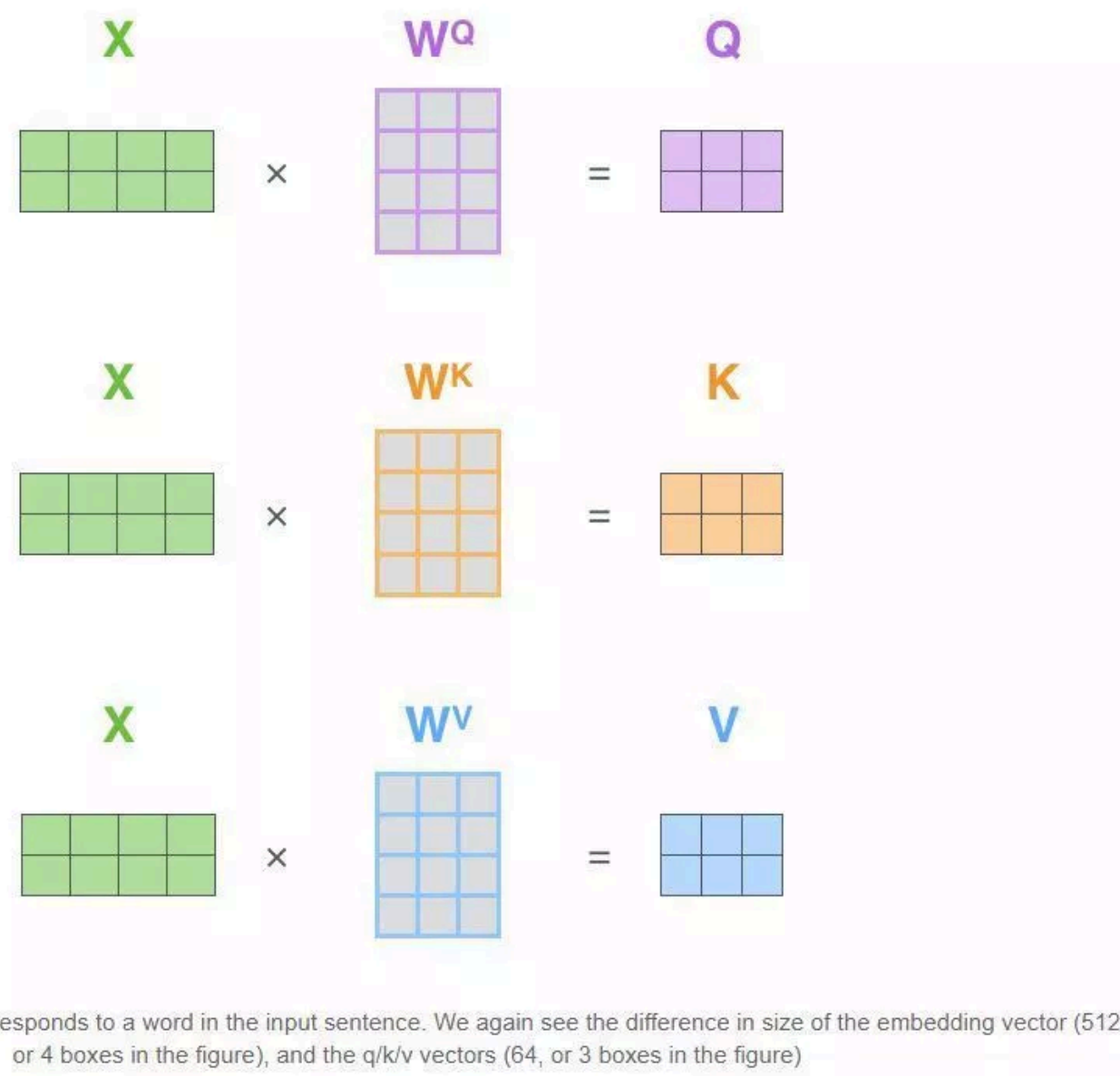
$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) \times V$$

$$= Z$$

The self-attention calculation in matrix form

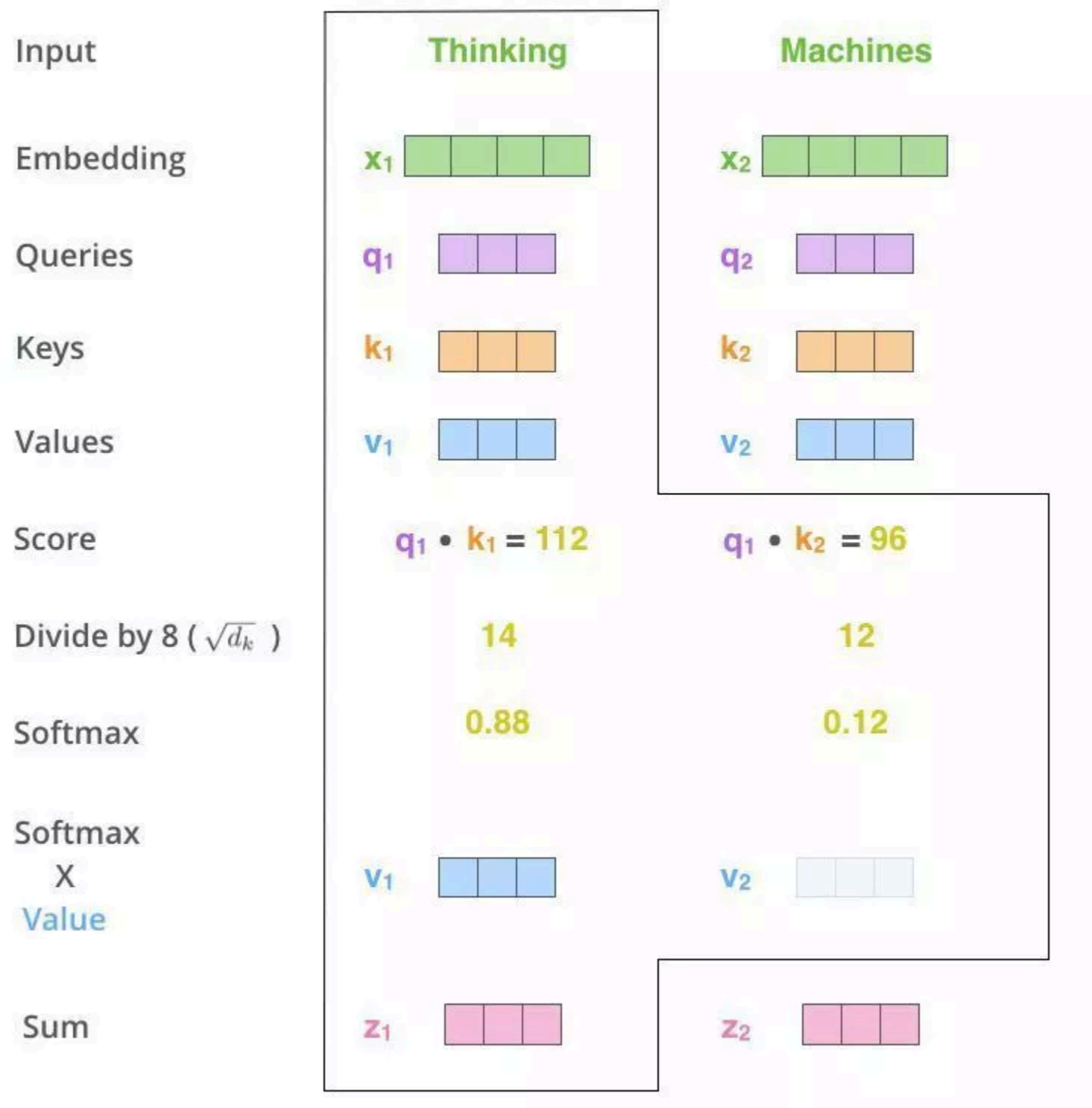


输入X经历3次线性变换得到Q,K,V



具体运算细节

由
X
到
Z



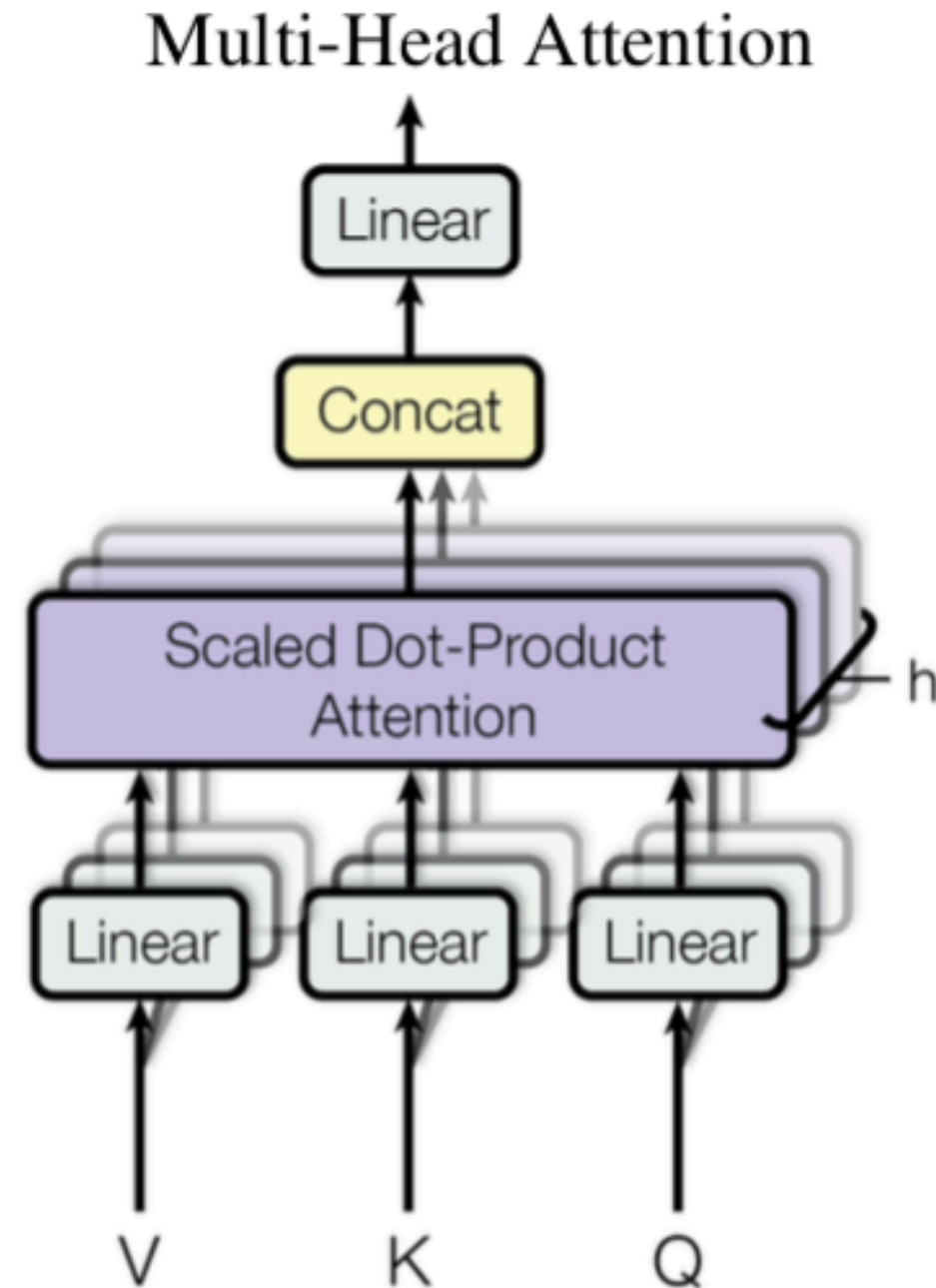
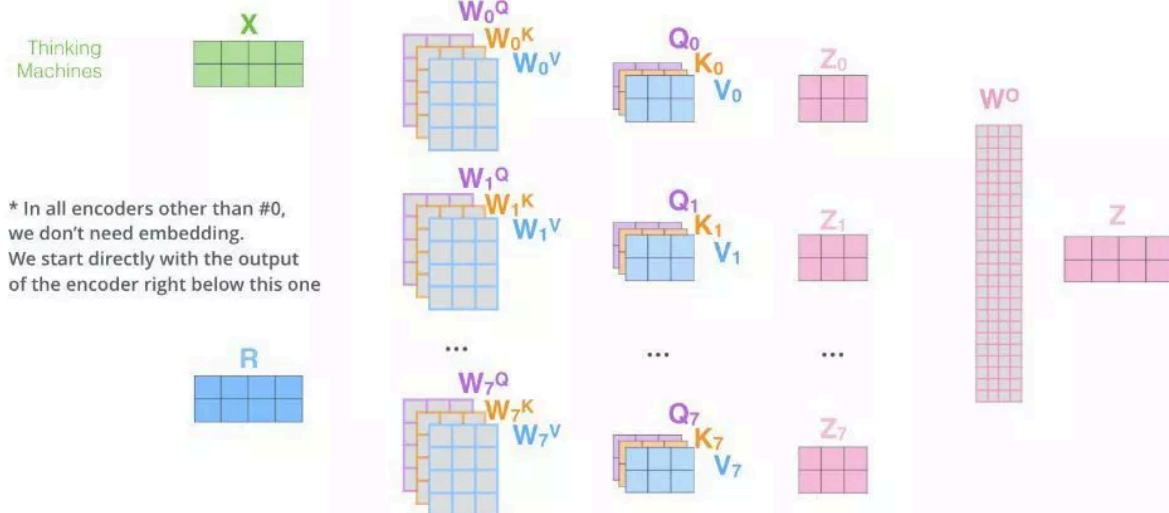
Multi-Head Attention

$$\text{MultiHead}(h) = \text{Concat}[\text{head}_1, \dots, \text{head}_m] W^O$$

where $\text{head}_i = \text{Attention}(Q, K, V)$
 where $Q, K, V = hW_i^Q, hW_i^K, hW_i^V$

where m is the number of heads, $h \in \mathbb{R}^{d_n \times d_{\text{model}}}$ is the input, the $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$, and $W^O \in \mathbb{R}^{m \cdot d_v \times d_{\text{model}}}$. The output of the Attention is $\text{head}_i \in \mathbb{R}^{d_n \times d_v}$ and the output of the MultiHead is $\in \mathbb{R}^{d_n \times d_{\text{model}}}$.

- 1) This is our input sentence*
- 2) We embed each word*
- 3) Split into 8 heads. We multiply X or R with weight matrices
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer



THX

And any questions?