

seq = sequence length

d_{model} = size of the embedding vector

h = number of heads

$d_k = d_v$ = d_{model} / h

$$MultiHead(Q, K, V) = \text{Concat}(\text{head}_1 \dots \text{head}_h)W^O$$

Multi-Head Attention

Given a sequence of embedding vectors $X \in \mathbb{R}^{N \times D}$, we obtain 3 different projections of it as:

$$Q = X W_Q \in \mathbb{R}^{N \times D}$$

$$K = X W_K \in \mathbb{R}^{N \times D}$$

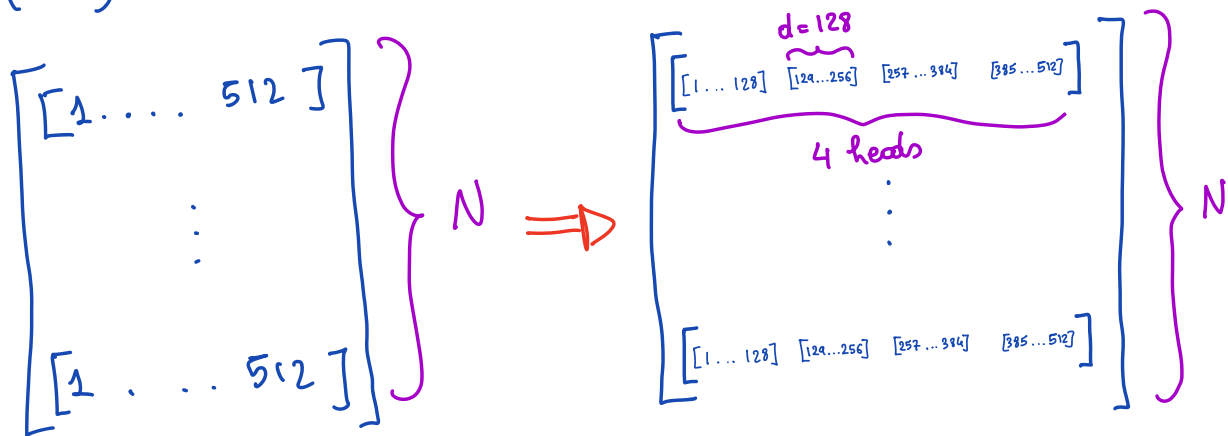
$$V = X W_V \in \mathbb{R}^{N \times D}$$

where W_Q , W_K and W_V are learnable parameter matrices.

Note: in case of cross-attention we use another sequence $Y \in \mathbb{R}^{N' \times D}$ to compute K and V

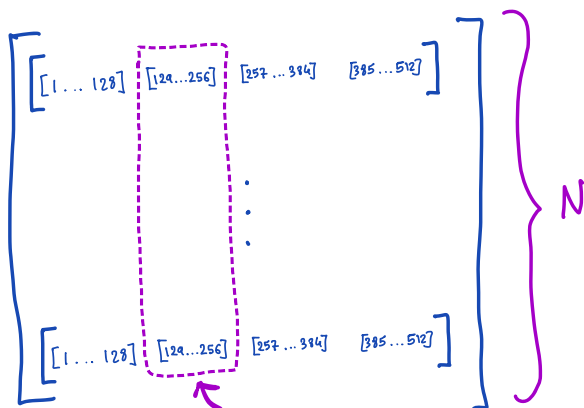
Why is it called Multi-Head?

we are given a sequence of embedding vectors (N, D) and we split it into multiple heads



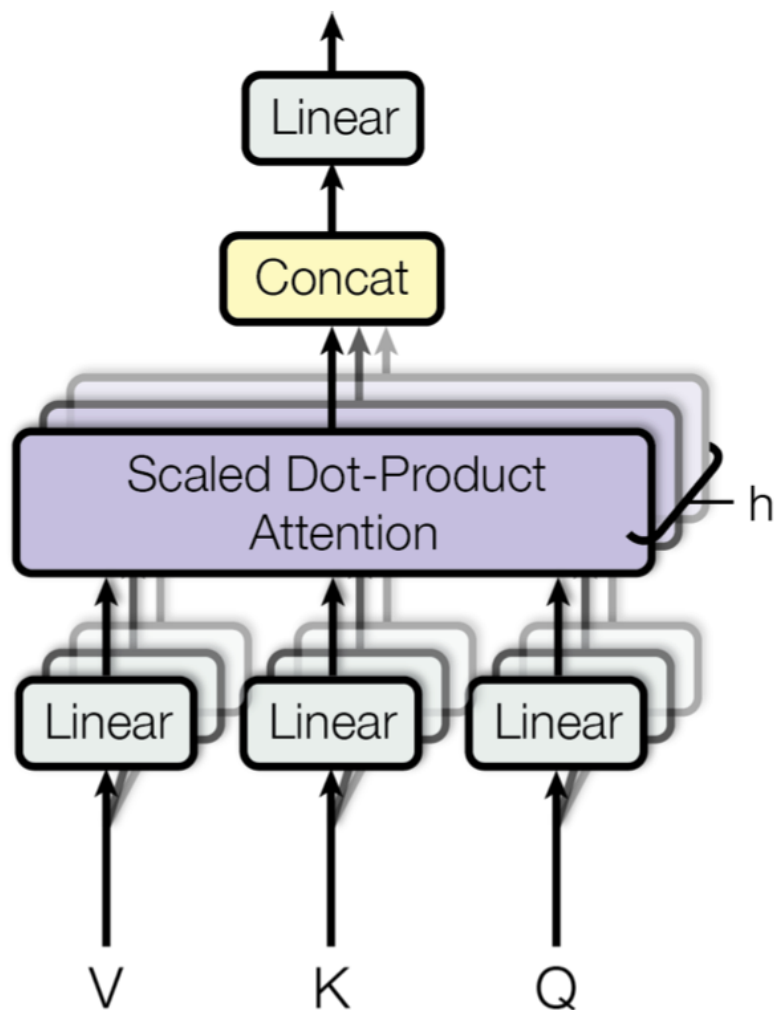
$$(N, D) \longrightarrow (N, 4, d) \text{ where } d = \frac{D}{4}$$

Note: D is commonly called d_{model} while d is commonly called d_{head}



The sequence
used by the head m. 2

Multi-Head Attention



Q : input sequence with shape (N, d)

K : input sequence with shape (N, d)

V : input sequence with shape (N, d)

N = sequence length

d = head dimension

We compute the following:

$$S = QK^T \in \mathbb{R}^{N \times N}$$

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_h}}\right)V$$

$$P = \text{softmax}(S) \in \mathbb{R}^{N \times N}$$

$$O = PV \in \mathbb{R}^{N \times d}$$

Note: we usually scale S by $\frac{1}{\sqrt{d}}$ but it can be absorbed by Q because $aQK^T = (aQ)K^T$ where a is a scalar.