

- [FusedAttention: fusing attention matrix into attention vector](#)
  - [Self-Attention](#)
  - [Fused-Attention](#)

# FusedAttention: fusing attention matrix into attention vector

---

*If some formula are not displayed clearly, chose the pdf format **README.pdf** in this repository.*

## Self-Attention

---

Let  $X \in \mathbb{R}^{h \times L}$  denotes a sequence of  $L$  feature vectors of dimensions  $h$ . Formally,  $X$  is projected by three matrices  $W_Q \in \mathbb{R}^{u \times h}$ ,  $W_K \in \mathbb{R}^{u \times h}$  and  $W_V \in \mathbb{R}^{d \times h}$  to corresponding representations  $Q$ ,  $K$  and  $V$ . The output for all positions is computed as follows

$$\begin{aligned}Q &= W_Q * X, \\K &= W_K * X, \\V &= W_V * X, \\ \alpha &= \text{softmax}(K^T Q / \sqrt{d}), \\Y &= V * \alpha.\end{aligned}$$

Note that the softmax function is applied column-wise. The  $Q$ ,  $K$ ,  $V$  and  $\alpha$  are referred to as the queries, keys, values, and attention matrix respectively, following the common terminology.

## Fused-Attention

---

In self-attention, the output  $Y$  at position  $t$  is computed as a weighted average of the feature representations of all positions with a weight proportional to a similarity score between the representations.

$$\begin{aligned}
Q_t &= W_Q X_t, \\
K &= W_K X, \\
V &= W_V X, \\
\alpha_t &= \text{softmax}(K^T Q_t / \sqrt{d}), \\
Y_t &= V * \alpha_t.
\end{aligned}$$

In short, self-attention maps  $L$  inputs to  $L$  outputs. In fused-attention, we rewrite the formula as

$$\begin{aligned}
Q &= f_Q(X), \\
K &= f_K(X), \\
V &= f_V(X), \\
\alpha &= \text{fuse}(\text{norm}(K^T Q / \sqrt{d})), \\
Y &= g_V(V * \alpha),
\end{aligned}$$

where  $f_Q, f_K$  and  $f_V$  are any functions that are legal to input  $X$ ,  $\text{norm}$  is a normalization function default to softmax applied column-wise,  $\text{fuse}$  is an aggregation function default to mean applied row-wise,  $g_V$  is a transform function related to  $V$ , thus the attention variable  $\alpha \in \mathbb{R}^{L \times 1}$  transforms values  $V$  from a matrix  $\in \mathbb{R}^{d \times L}$  to a vector  $\in \mathbb{R}^{d \times 1}$ . In brief, fused-attention maps  $L$  inputs to 1 output. There are two forms of  $X$ :

- The first case is that  $X$  is a tensor of any size. Ignoring the batch dimension, suppose the size of  $X$  is  $(C, W_1, W_2, \dots, W_n)$ , where  $C$  is the number of channels,  $W_n$  is the spatial width at the  $n$ -th spatial dimension. In this case, take  $f_V$  for example, it first transforms  $X$  into  $V'$  with size  $(C', W'_1, W'_2, \dots, W'_n)$ , then transforms  $V'$  into  $V$  with size  $(d, L)$ , where  $d = C' \prod_{i=1}^n \Delta W'_i$ , of which  $\Delta W'_i$  is the window size at the  $i$ -th spatial dimension,  $L$  is the total number of patches.  $f_Q$  and  $f_K$  are similar to  $f_V$ . If the number of dimensions of  $Y$  is required to be the same with  $X$ , then  $g_V$  transforms  $V * \alpha$  into  $Y$  with size  $(C', \Delta W'_1, \Delta W'_2, \dots, \Delta W'_n)$ , otherwise  $g_V$  does nothing.
- The second case is that  $X$  is a collection of  $L$  tensors of the same size. Also ignoring the batch dimension,  $X_j$  is the  $j$ -th element of  $X$ , suppose the size of  $X_j$  is  $(C, W_1, W_2, \dots, W_n)$ , where  $C$  is the number of channels,  $W_n$  is the spatial width at the  $n$ -th spatial dimension. In this case, take  $f_V$  for example, it first transforms  $X_j$  into  $V'_j$  with size  $(C', W'_1, W'_2, \dots, W'_n)$ , then reshape  $V'_j$  to  $V_j$  of size  $(C' \prod_{i=1}^n W'_i, 1)$ , then concatenate  $V_j, j = 1, 2, \dots, L$  into  $V$  with size

$(C' \prod_{i=1}^n W'_i, L).f_Q$  and  $f_K$  are similar to  $f_V$ . As to  $g_V$ , it transforms  $V * \alpha$  into  $Y$  with size  $(C', W'_1, W'_2, \dots, W'_n)$ .