

# A Mathematical View of Attention Models in Deep Learning

Shuiwang Ji  
Department of Computer Science & Engineering  
Texas A&M University

- 1 Given a set of  $n$  query vectors  $q_1, q_2, \dots, q_n \in \mathbb{R}^d$ ,  $m$  key vectors  $k_1, k_2, \dots, k_m \in \mathbb{R}^d$ , and  $m$  value vectors  $v_1, v_2, \dots, v_m \in \mathbb{R}^p$ , the attention mechanism computes a set of output vectors  $o_1, o_2, \dots, o_n \in \mathbb{R}^q$  by linearly combining the  $g$ -transformed value vectors  $g(v_i) \in \mathbb{R}^q$  using the relations between the corresponding query vector and each key vector as coefficients.
- 2 Formally,

$$o_j = \frac{1}{C} \sum_{i=1}^m f(q_j, k_i) g(v_i), \quad (1)$$

where  $f(q_j, k_i)$  characterizes the relation (e.g., similarity) between  $q_j$  and  $k_i$ ,  $g(\cdot)$  is commonly a linear transformation as  $g(v_i) = \mathbf{W}_v v_i \in \mathbb{R}^q$ , where  $\mathbf{W}_v \in \mathbb{R}^{q \times p}$ , and  $C = \sum_{i=1}^m f(q_j, k_i)$  is a normalization factor.

# Attention Model

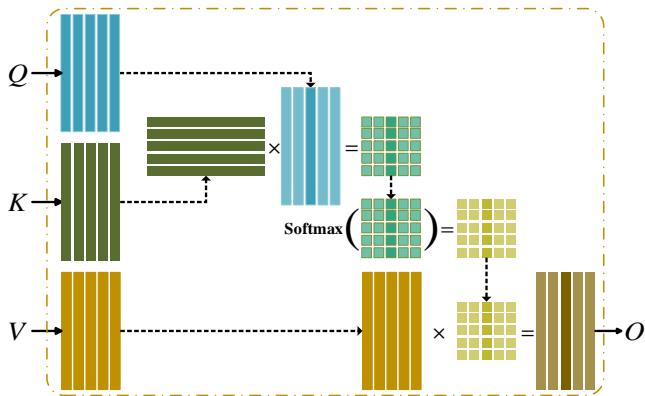
- 1 A commonly used similarity function is the embedded Gaussian, defined as  $f(q_j, k_i) = \exp(\theta(q_j)^T \phi(k_i))$ , where  $\theta(\cdot)$  and  $\phi(\cdot)$  are commonly linear transformations as  $\theta(q_j) = \mathbf{W}_q q_j$  and  $\phi(k_i) = \mathbf{W}_k k_i$ .
- 2 Note that if we treat the value vectors as inputs, each output vector  $o_j$  is dependent on all input vectors. When the embedded Gaussian similarity and linear transformation are used, these computations can be expressed succinctly in matrix form as

$$\mathbf{O} = \mathbf{W}_v \mathbf{V} \times \text{softmax} \left( (\mathbf{W}_k \mathbf{K})^T \mathbf{W}_q \mathbf{Q} \right), \quad (2)$$

where  $\mathbf{Q} = [q_1, q_2, \dots, q_n] \in \mathbb{R}^{d \times n}$ ,  $\mathbf{K} = [k_1, k_2, \dots, k_m] \in \mathbb{R}^{d \times m}$ ,  $\mathbf{V} = [v_1, v_2, \dots, v_m] \in \mathbb{R}^{p \times m}$ ,  $\mathbf{O} = [o_1, o_2, \dots, o_n] \in \mathbb{R}^{q \times n}$ , and  $\text{softmax}(\cdot)$  computes a normalized version of the input matrix, where each column is normalized using the softmax function to sum to one.

- 3 Note that the number of output vectors is equal to the number of query vectors. In self-attention, we have  $\mathbf{Q} = \mathbf{K} = \mathbf{V}$ .

# Attention Model



**Figure:** An illustration of the attention operator. Here,  $\times$  denotes matrix multiplication, and  $\text{Softmax}(\cdot)$  is the column-wise softmax operator.  $Q$ ,  $K$ , and  $V$  are input matrices. A similarity score is computed between each query vector as a column of  $Q$  and each key vector as a column in  $K$ .  $\text{Softmax}(\cdot)$  normalizes these scores and makes them sum to 1. Multiplication between normalized scores and the matrix  $V$  yields the corresponding output vector.

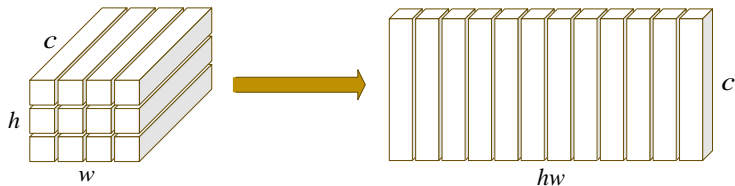
# Attention for higher order data

- 1 The attention mechanism was originally developed in natural language processing to process 1-D data.
- 2 It has been extended to deal with 2-D images and 3-D video data recently.
- 3 When deal with 2-D data, the inputs to the attention operator can be represented as 3-D tensors  $\mathbf{Q} \in \mathbb{R}^{h \times w \times c}$ ,  $\mathbf{K} \in \mathbb{R}^{h \times w \times c}$ , and  $\mathbf{V} \in \mathbb{R}^{h \times w \times c}$ , where  $h$ ,  $w$ , and  $c$  represent the height, width, and number of channels, respectively. Note that for notational simplicity, we have assumed the three tensors having the same size.

# Attention for higher order data

- 1 These tensors are first unfolded into matrices along mode-3, resulting in  $\mathbf{Q}_{(3)}, \mathbf{K}_{(3)}, \mathbf{V}_{(3)} \in \mathbb{R}^{c \times hw}$ .
- 2 Columns of these matrices are the mode-3 fibers of the corresponding tensors.
- 3 These matrices are used to compute output vectors as in regular attention described above. The output vectors are then folded back to a 3-D tensor  $\mathbf{O} \in \mathbb{R}^{h \times w \times q}$  by treating them as mode-3 fibers of  $\mathbf{O}$ .
- 4 Note that the height and width of  $\mathbf{O}$  are equal to those of  $\mathbf{Q}$ . That is, we can obtain an output with larger/smaller spatial size by providing an input  $\mathbf{Q}$  of correspondingly larger/smaller spatial size.
- 5 Again, we have  $\mathbf{Q} = \mathbf{K} = \mathbf{V}$  in self-attention.

# Attention for higher order data



**Figure:** Conversion of a third-order tensor into a matrix by unfolding along mode-3. In this example, a  $h \times w \times c$  tensor is unfolded into a  $c \times hw$  matrix.

THANKS!