# ATTENTION

- **Self-attention** (intra-attention): attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence.
  + Usages: reading comprehension, abstractive summarization, textual entailment, and learning task-independent sequence representations.

- Attention function: mapping a query and a set of key-value to comp pairs to an output
  + Query; keys, values, and output: vectors
  + Output: weighted sum of the values
    ↳ compatibility function of query and keys

① Scaled dot Dot Product Attention:

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

  + $Q$: matrix of queries
  + $k$: matrix of keys
  + $V$: matrix of values
  + $d_k$: keys of dimension
  + softmax: $\sigma(\vec{z})_i = \dfrac{e^{z_i}}{\sum\limits_{j=1}^{k} e^{z_j}}$  (*)

  + $q$ and $k$ are independent random variable with mean 0 and variance 1 ⟹ Dot-product: $q \cdot k = \sum\limits_{i=1}^{d_k} q_i k_i$ has mean 0 and variance $d_k$
  ⟹ Extremely small gradients ⟹ Scale factor: $\dfrac{1}{\sqrt{d_k}}$

② Multi-Head Attention:

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

  + $\text{head}_i = \text{Attention}(QW_i^Q, kW_i^k, VW_i^V)$
  + $W_i^Q \in \mathbb{R}^{d_{model} \times d_k}$
  + $W_i^k \in \mathbb{R}^{d_{model} \times d_k}$
  + $W_i^V \in \mathbb{R}^{d_{model} \times d_k}$
  + $W^O \in \mathbb{R}^{hd_v \times d_{model}}$

  + $d_{model}$: dimension keys, values, and queries
  + $d_k, d_v$: dimension
  + $d_v$: dimension output
  + $h$: parallel attention layers

  - Perform the attention layer in parallel

③ Application:
  - Decoder → queries
    Encoder → keys, values  } → Position in the decoder to attend over all positions in the input sequence