

# Attention

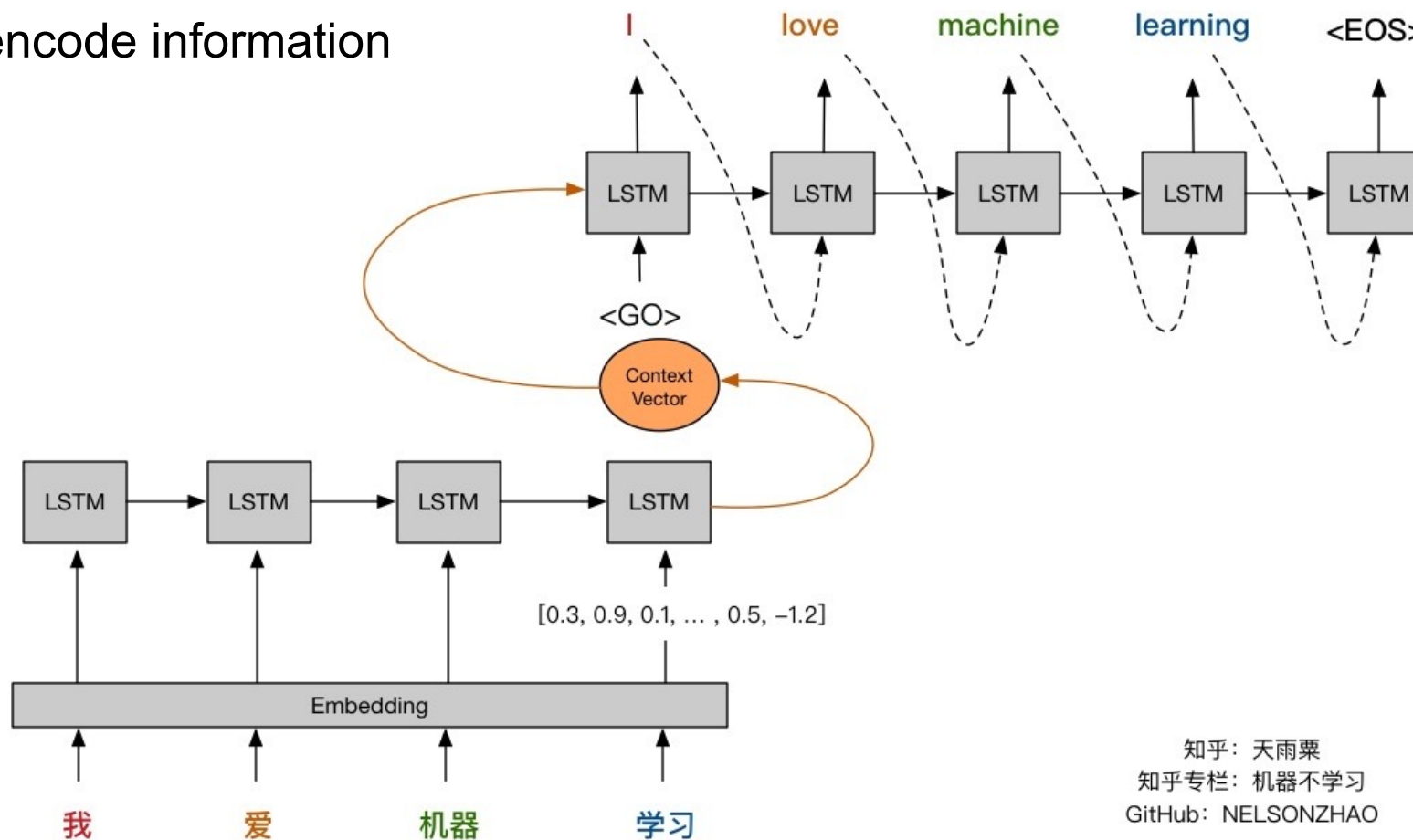
**Kun Yuan**

**Center for Machine Learning Research @ Peking University**

Oct. 24, 2023

# Traditional seq2seq model

RNN is hard to encode information that is far away

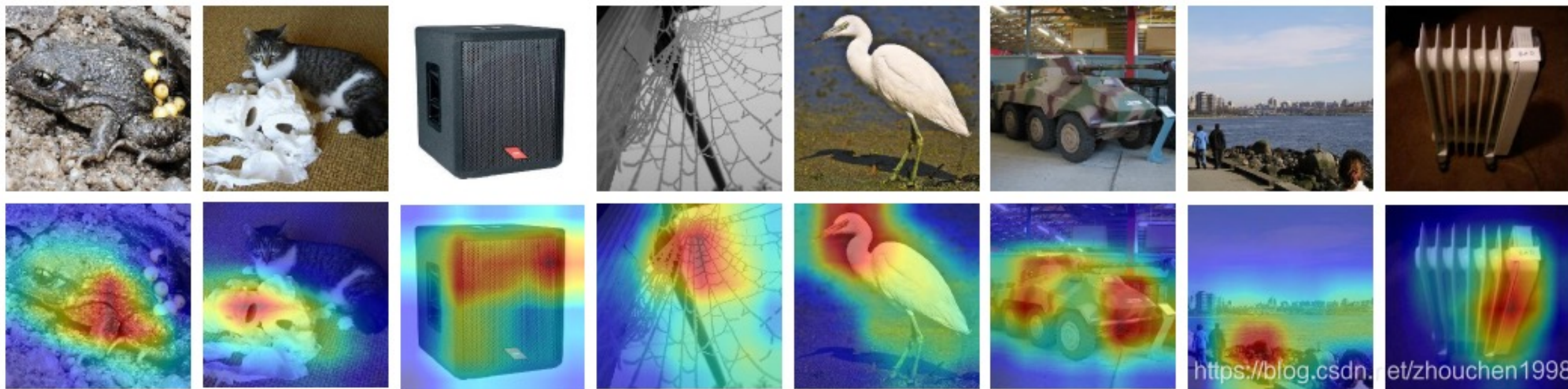


知乎：天雨粟  
知乎专栏：机器不学习  
GitHub：NELSONZHAO

知乎 @天雨粟

# Attention

All you need is attention!



How to capture the most valuable information from a pool of candidate?

Consider a pool of candidate information  $D = (k_1, v_1), (k_2, v_2), \dots, (k_m, v_m)$

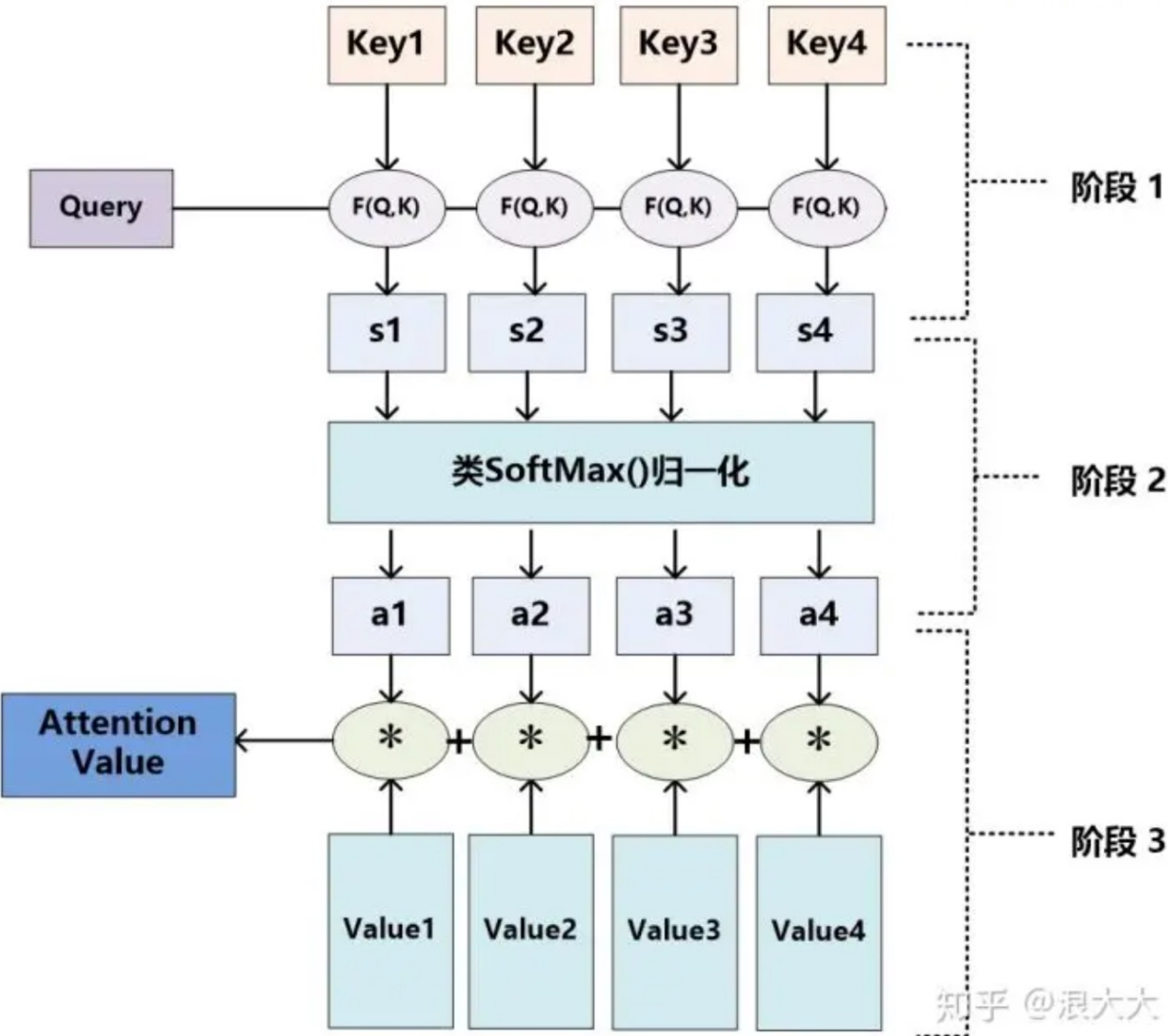
Given a query  $q$ , we can capture the most valuable information as follows

$$\text{Attention}(\mathbf{q}, \mathcal{D}) \stackrel{\text{def}}{=} \sum_{i=1}^m \alpha(\mathbf{q}, \mathbf{k}_i) \mathbf{v}_i,$$

where weight  $\alpha$  is to evaluate how close the query  $q$  is to key  $k_i$

$$\alpha(\mathbf{q}, \mathbf{k}_i) = \text{softmax}(a(\mathbf{q}, \mathbf{k}_i)) = \frac{\exp(\mathbf{q}^\top \mathbf{k}_i / \sqrt{d})}{\sum_{j=1} \exp(\mathbf{q}^\top \mathbf{k}_j / \sqrt{d})}.$$

# Attention



# Seq2Seq with attention

