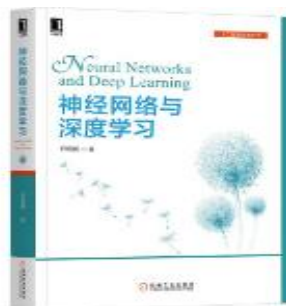


MACHINE LEARNING

机器学习

Neural Networks

循环神经网络与注意力机制



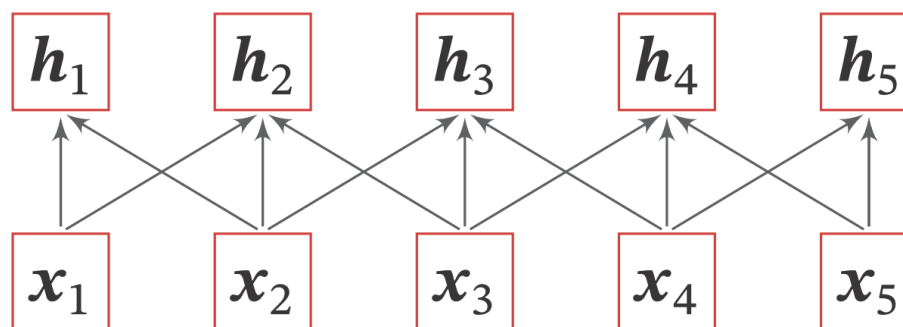
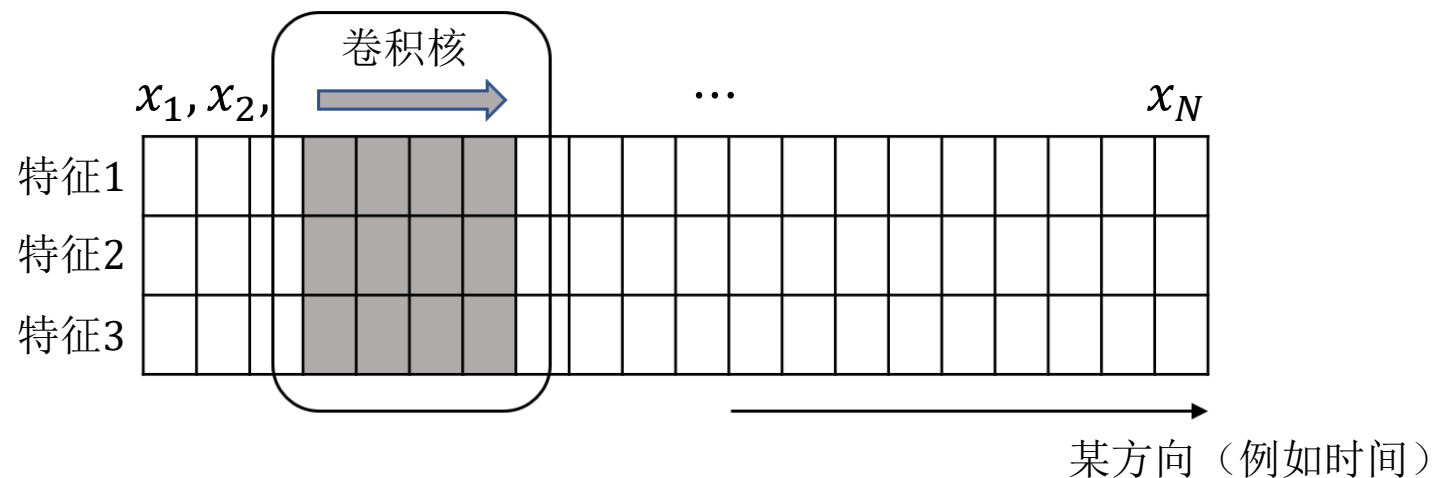
参考：
《神经网络与深度学习》

Machine Learning Course
Copyright belongs to Wenting Tu.

循环神经网络 *Recurrent Neural Network*

- 卷积层实现序列到序列

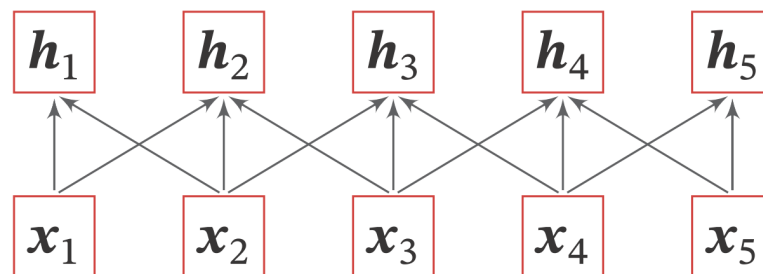
卷积层处理1D序列数据



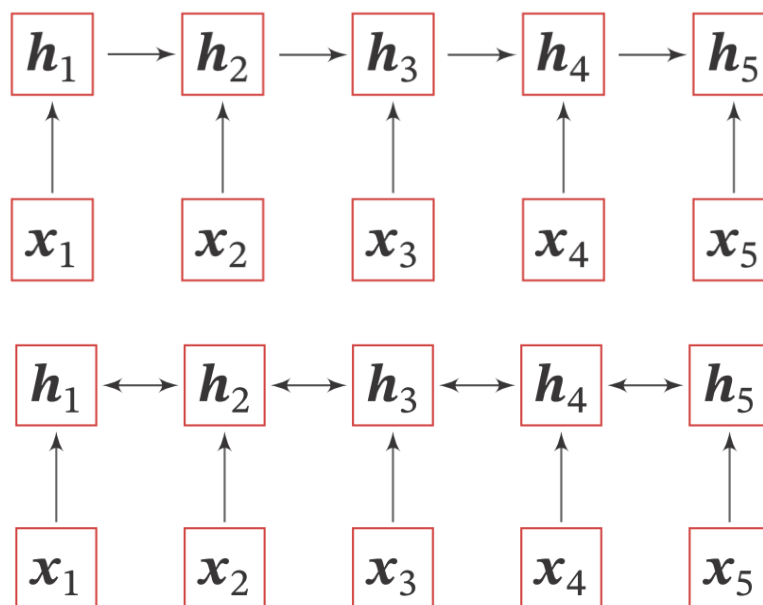
循环神经网络 *Recurrent Neural Network*

- 循环层实现序列到序列

CNN



RNN



$$h_t = f(Uh_{t-1} + Wx_t + b)$$

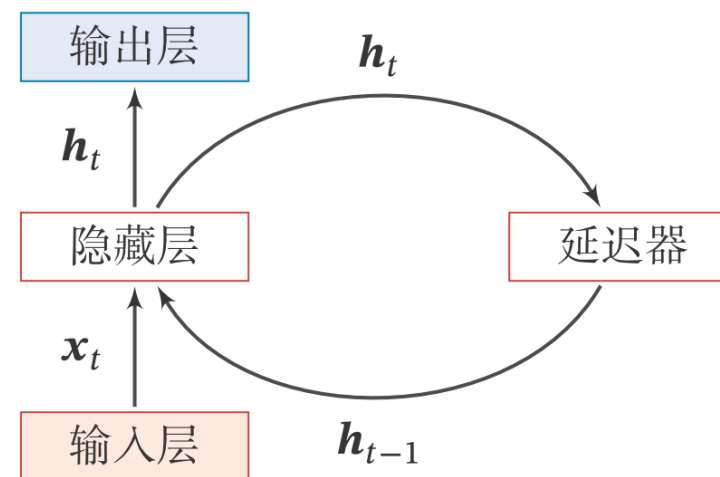
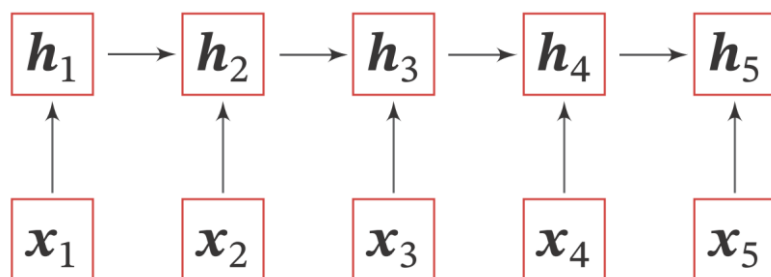
$$h_t^{(1)} = f(U^{(1)}h_{t-1}^{(1)} + W^{(1)}x_t + b^{(1)})$$

$$h_t^{(2)} = f(U^{(2)}h_{t+1}^{(2)} + W^{(2)}x_t + b^{(2)})$$

$$h_t = h_t^{(1)} \oplus h_t^{(2)}$$

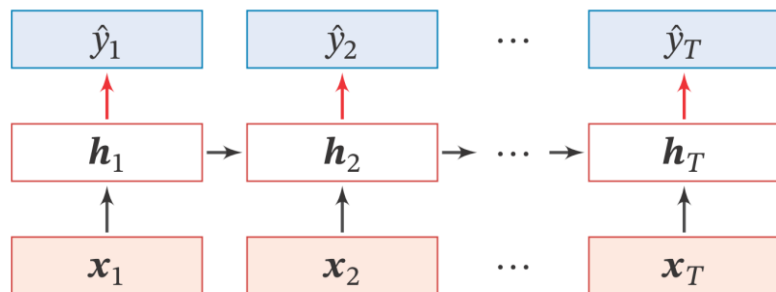
循环神经网络 *Recurrent Neural Network*

- 循环层实现序列到序列

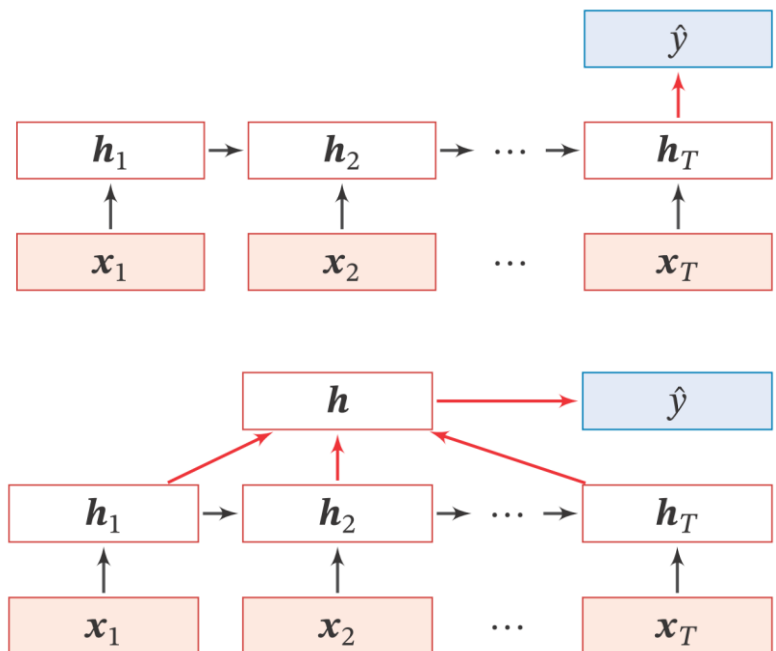


循环神经网络 *Recurrent Neural Network*

- 循环层上的输出



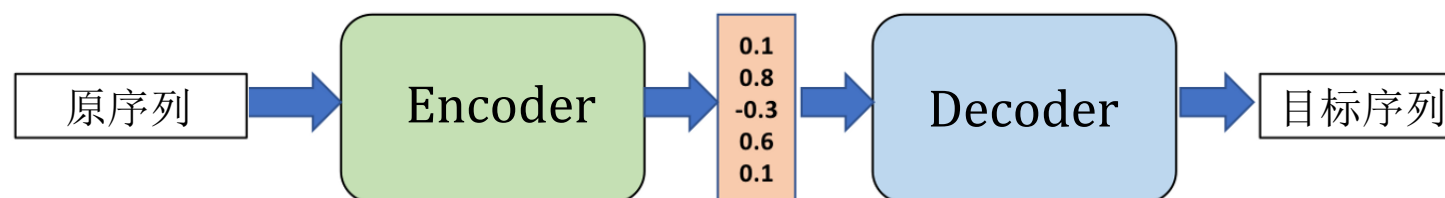
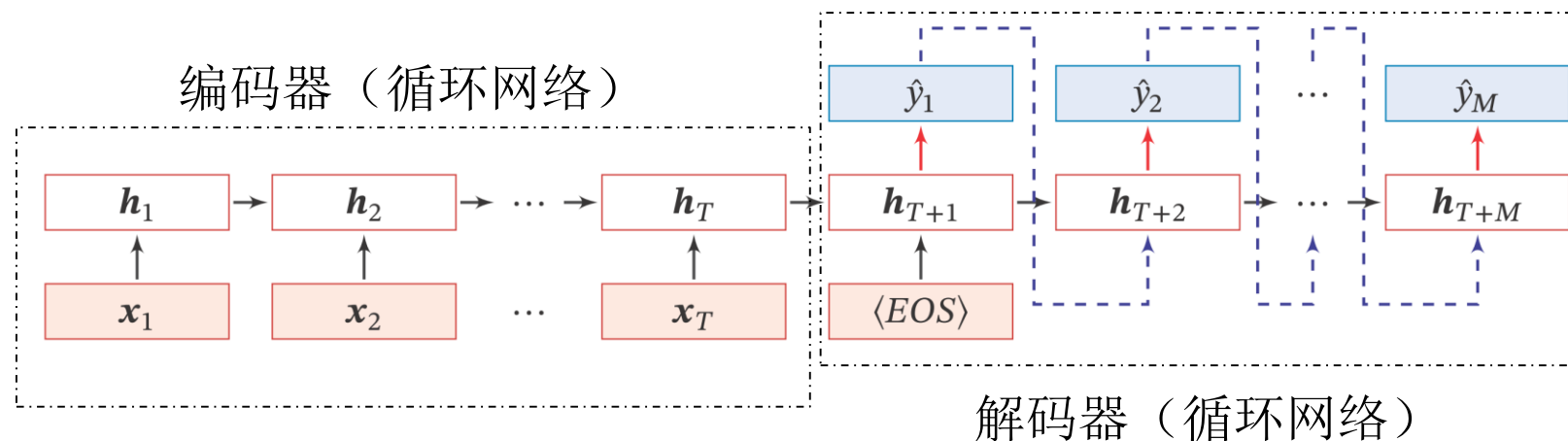
如果输出是一个同步的序列，可以把每步输出建立在每步的隐层上



如果输出是一个标签，可以把输出建立最后一步隐层或所有步的隐层上

循环神经网络 *Recurrent Neural Network*

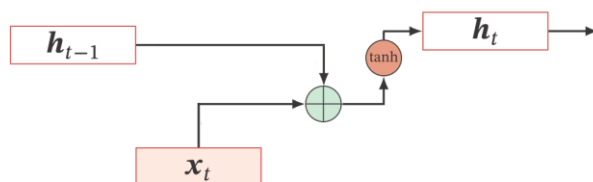
- 异步的序列到序列
 - 编码器-解码器 (Encoder-Decoder)



门控循环单元 *Gated Recurrent Unit*

- 门控制

为了改善循环神经网络的长程依赖问题，流行的方法是引入门控机制来控制信息的累积速度，包括有选择地遗忘之前累积的信息以及有选择地加入新的信息。

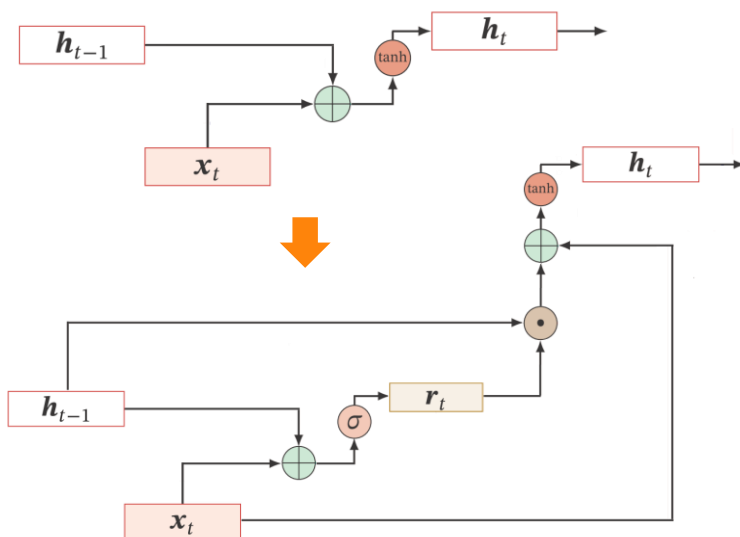


$$h_t = f(Uh_{t-1} + Wx_t + b)$$

门控循环单元 *Gated Recurrent Unit*

• 门控制

为了改善循环神经网络的长程依赖问题，流行的方法是引入门控机制来控制信息的累积速度，包括有选择地遗忘之前累积的信息以及有选择地加入新的信息。



$$h_t = f(Uh_{t-1} + Wx_t + b)$$

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r)$$

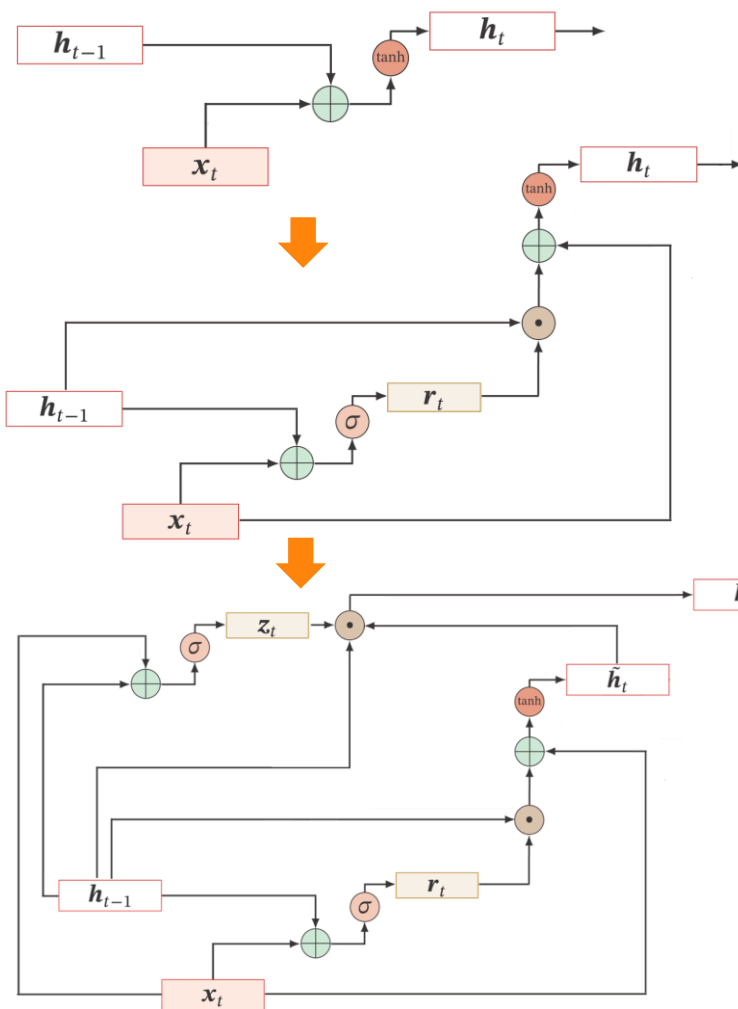
$$h_t = \tanh(W_h x_t + U_h(r_t \odot h_{t-1}) + b_h)$$

r_t 控制 h_t 的计算是否/多大程度依赖上一时刻的状态 h_{t-1}
这允许我们转换当前信息的时候一定程度上遗忘历史信息

门控循环单元 *Gated Recurrent Unit*

• 门控制

为了改善循环神经网络的长程依赖问题，流行的方法是引入门控机制来控制信息的累积速度，包括有选择地遗忘之前累积的信息以及有选择地加入新的信息。



$$h_t = f(Uh_{t-1} + Wx_t + b)$$

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r)$$

$$h_t = \tanh(W_h x_t + U_h(r_t \odot h_{t-1}) + b_h)$$

r_t 控制 h_t 的计算是否/多大程度依赖上一时刻的状态 h_{t-1}
这允许我们转换当前信息的时候一定程度上遗忘历史信息

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r)$$

$$\tilde{h}_t = \tanh(W_h x_t + U_h(r_t \odot h_{t-1}) + b_h)$$

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z)$$

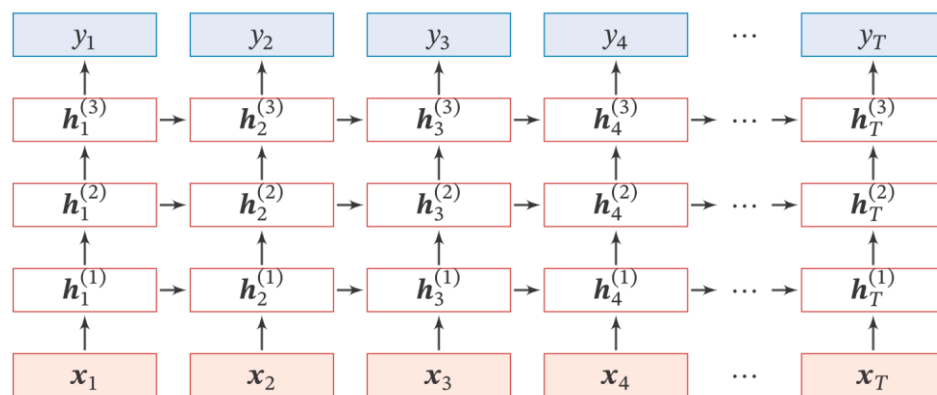
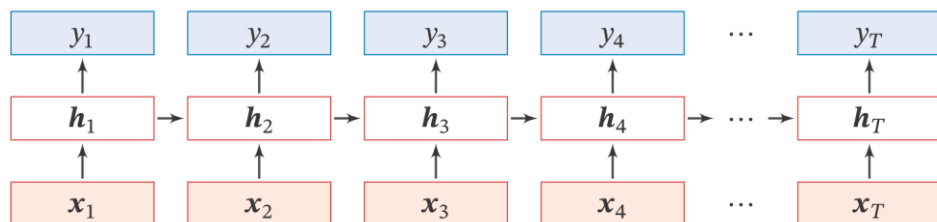
$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t$$

z_t 控制 h_t 的计算是如何融合当前信息与过去信息的
这允许我们在输出信息的时候一定程度上遗忘当前

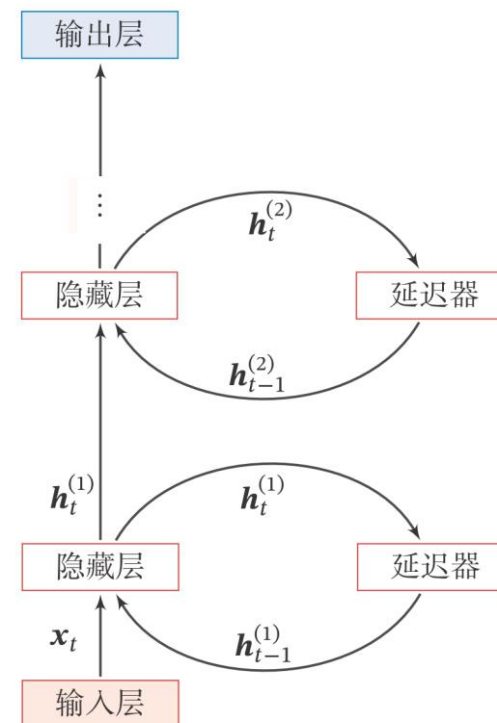
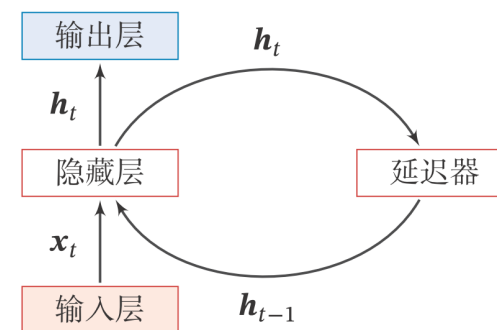
Gated Recurrent Unit, GRU

其他变体

• 堆叠循环神经网络

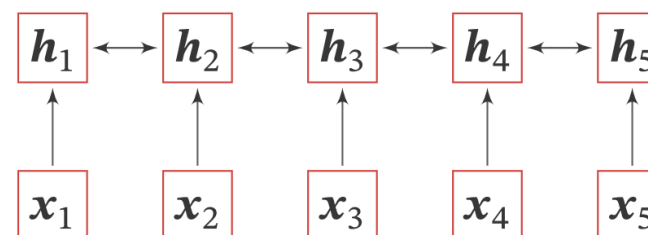
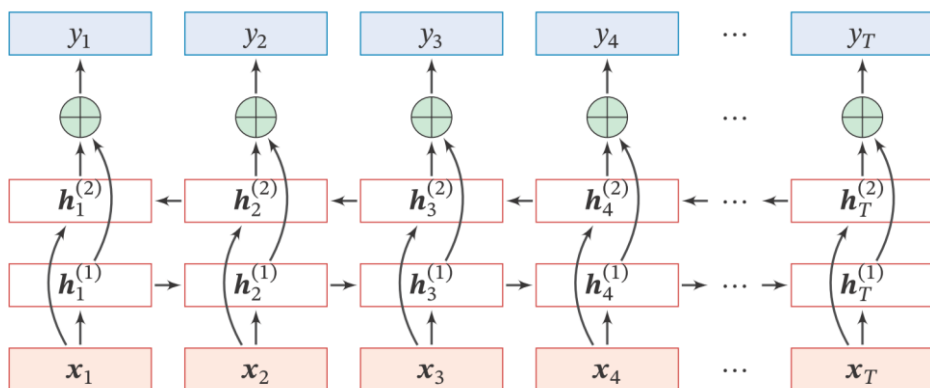
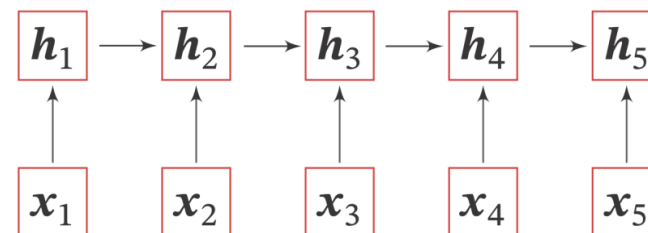
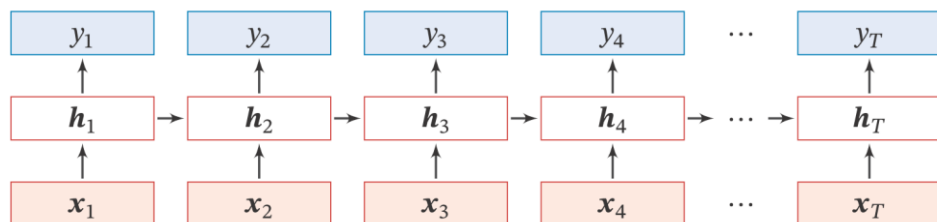


$$h_t^{(l)} = f(U^{(l)}h_{t-1}^{(l)} + W^{(l)}h_t^{(l-1)} + b^{(l)})$$



其他变体

• 双向循环神经网络



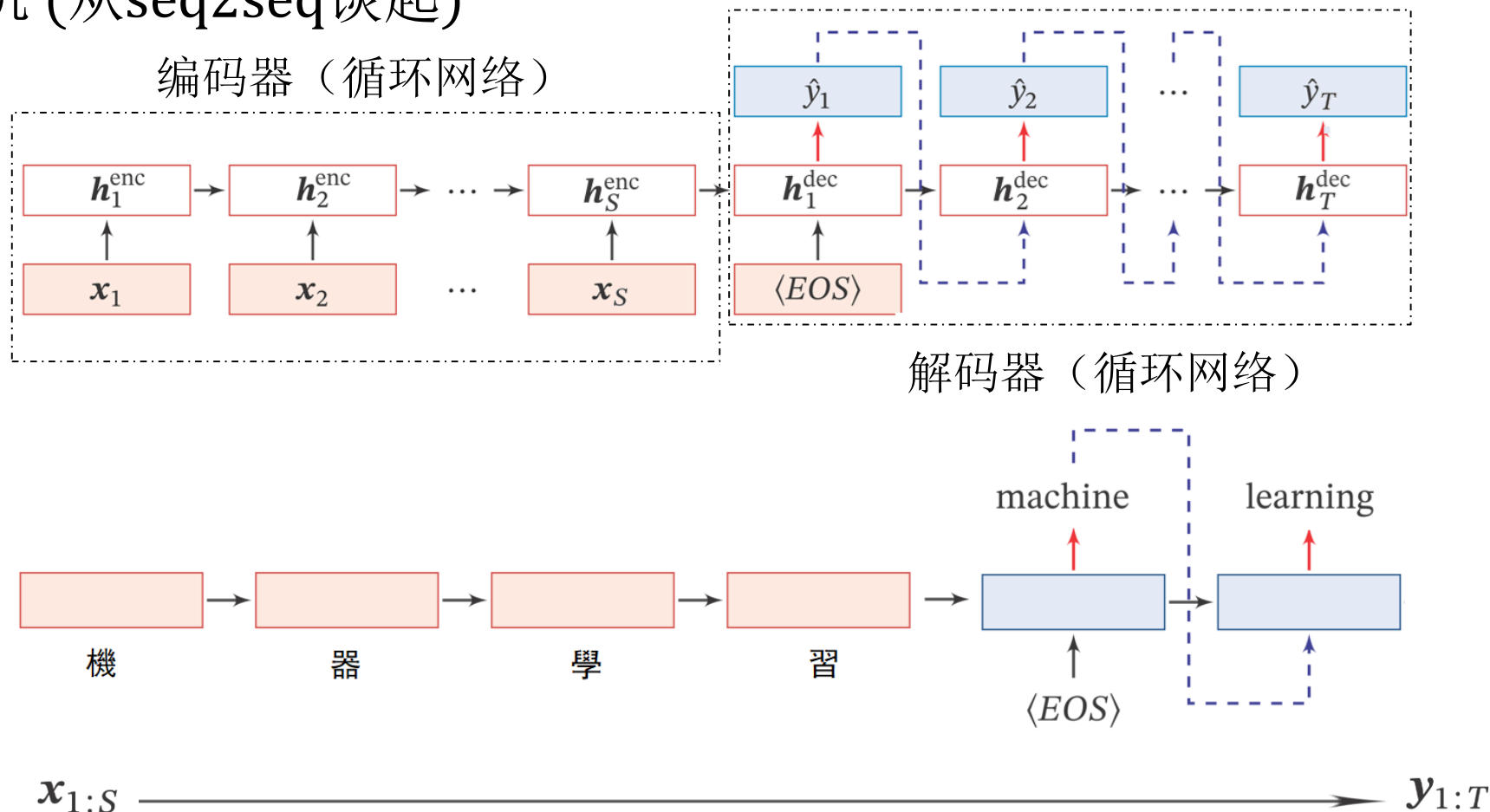
$$h_t^{(1)} = f(\mathbf{U}^{(1)}h_{t-1}^{(1)} + \mathbf{W}^{(1)}x_t + \mathbf{b}^{(1)})$$

$$h_t^{(2)} = f(\mathbf{U}^{(2)}h_{t+1}^{(2)} + \mathbf{W}^{(2)}x_t + \mathbf{b}^{(2)})$$

$$h_t = h_t^{(1)} \oplus h_t^{(2)}$$

注意力机制

• 动机 (从seq2seq谈起)



$$h_t^{\text{enc}} = f_{\text{enc}}(h_{t-1}^{\text{enc}}, \mathbf{e}_{x_t}, \theta_{\text{enc}})$$

$$\forall t \in [1 : S]$$

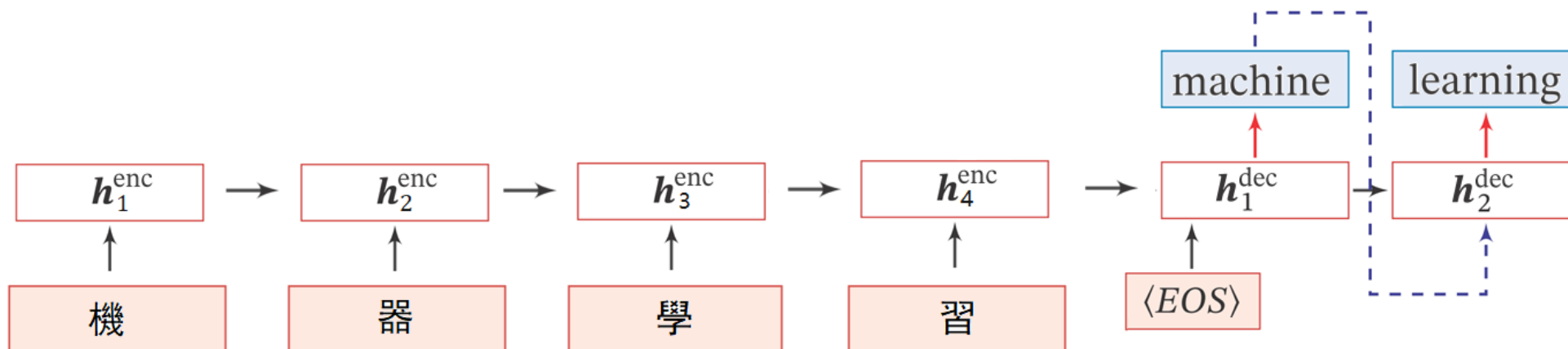
$$\mathbf{u} = h_S^{\text{enc}} \quad h_0^{\text{dec}} = \mathbf{u}$$

$$h_t^{\text{dec}} = f_{\text{dec}}(h_{t-1}^{\text{dec}}, \mathbf{e}_{y_{t-1}}, \theta_{\text{dec}})$$

$$\mathbf{o}_t = g(h_t^{\text{dec}}, \theta_o)$$

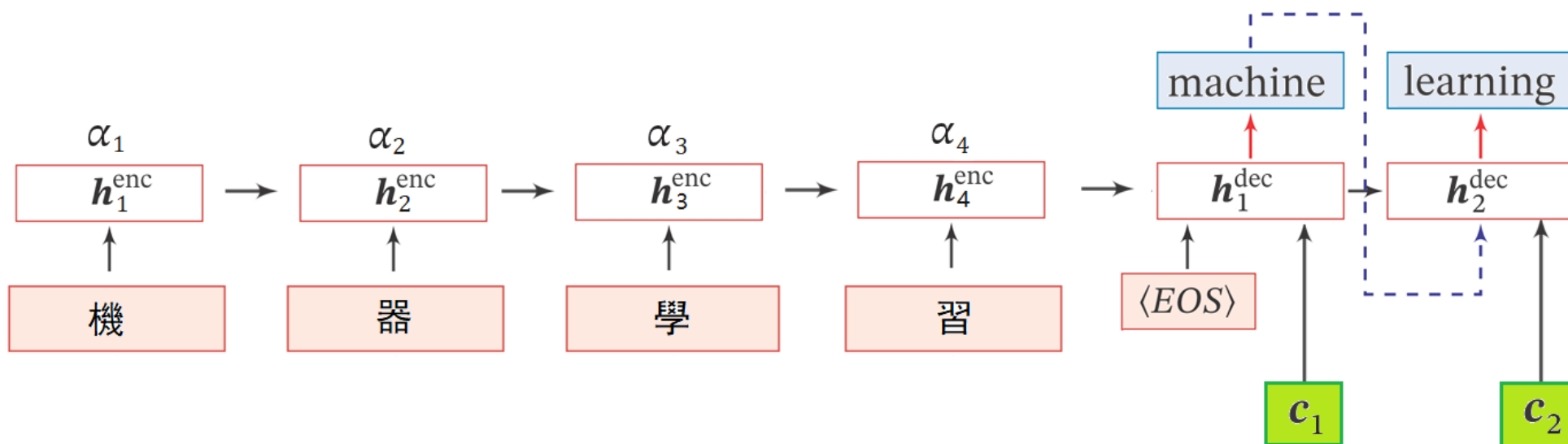
注意力机制

- 动机 (从seq2seq谈起)



注意力机制

- 动机 (从seq2seq谈起)

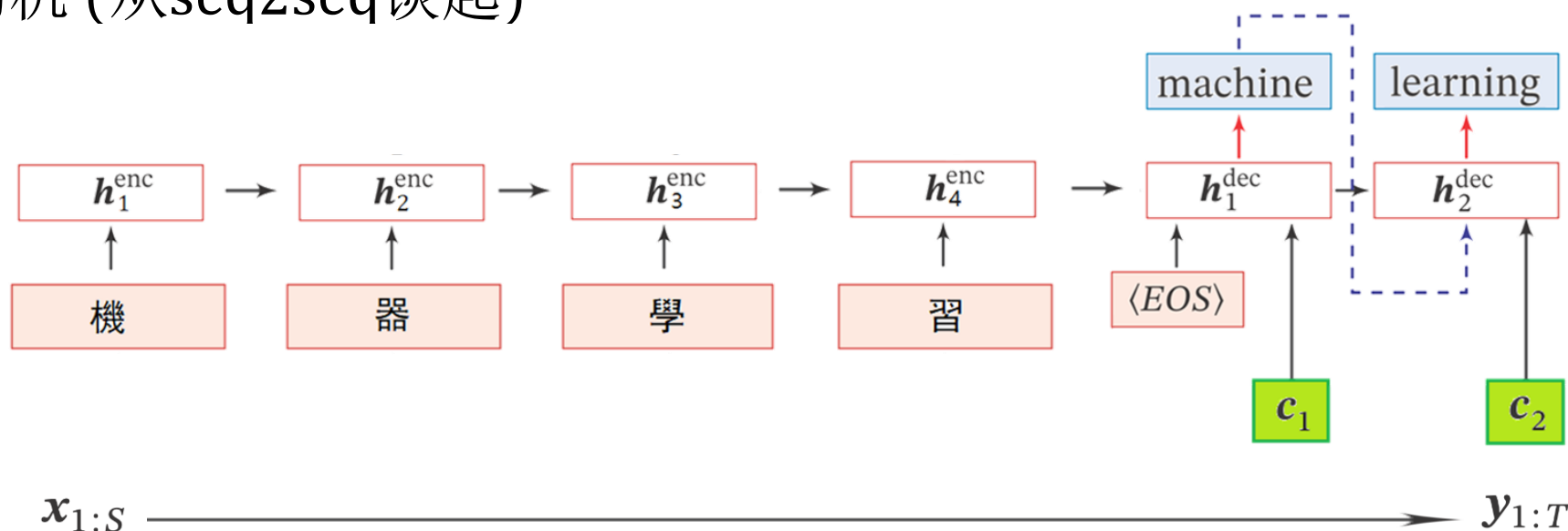


$$c_1 = 0.5 h_1^{\text{enc}} + 0.5 h_2^{\text{enc}} + 0.0 h_3^{\text{enc}} + 0.0 h_4^{\text{enc}}$$

$$c_2 = 0.0 h_1^{\text{enc}} + 0.0 h_2^{\text{enc}} + 0.5 h_3^{\text{enc}} + 0.5 h_4^{\text{enc}}$$

注意力机制

• 动机 (从seq2seq谈起)



$$h_t^{\text{enc}} = f_{\text{enc}}(h_{t-1}^{\text{enc}}, e_{x_t}, \theta_{\text{enc}})$$

$$\forall t \in [1 : S]$$

$$H^{\text{enc}} = [h_1^{\text{enc}}, \dots, h_S^{\text{enc}}]$$

$$u = h_S^{\text{enc}} \quad h_0^{\text{dec}} = u$$

$$c_t = \text{att}(H^{\text{enc}}, h_{t-1}^{\text{dec}}) = \sum_{i=1}^S \alpha_i h_i^{\text{enc}}$$

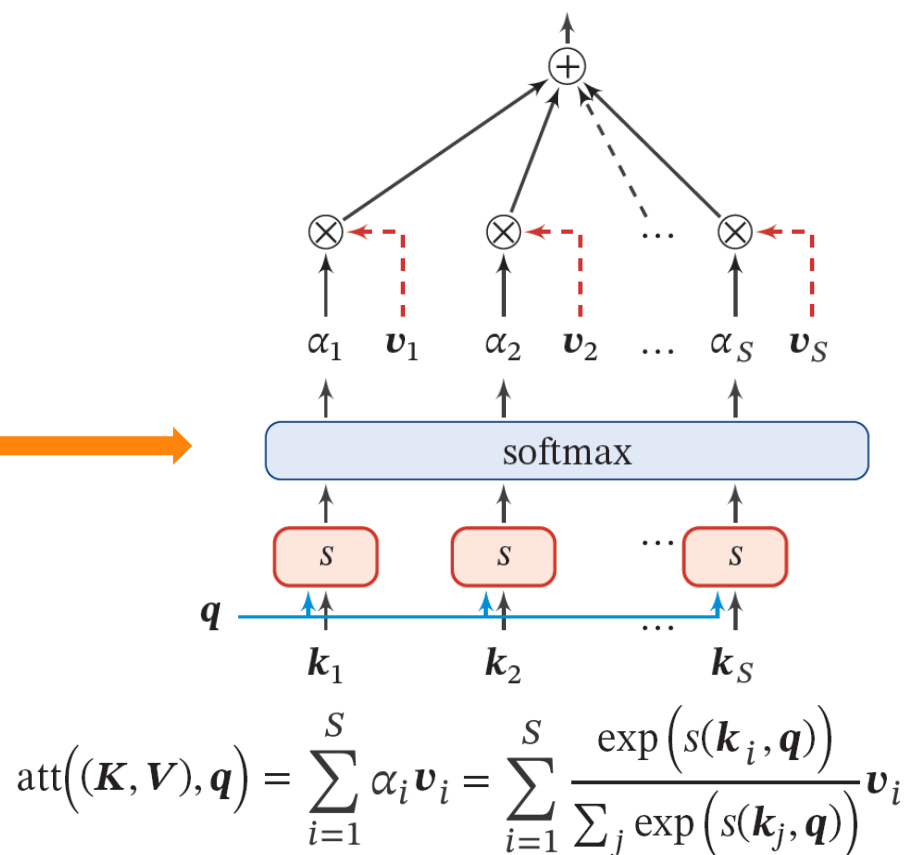
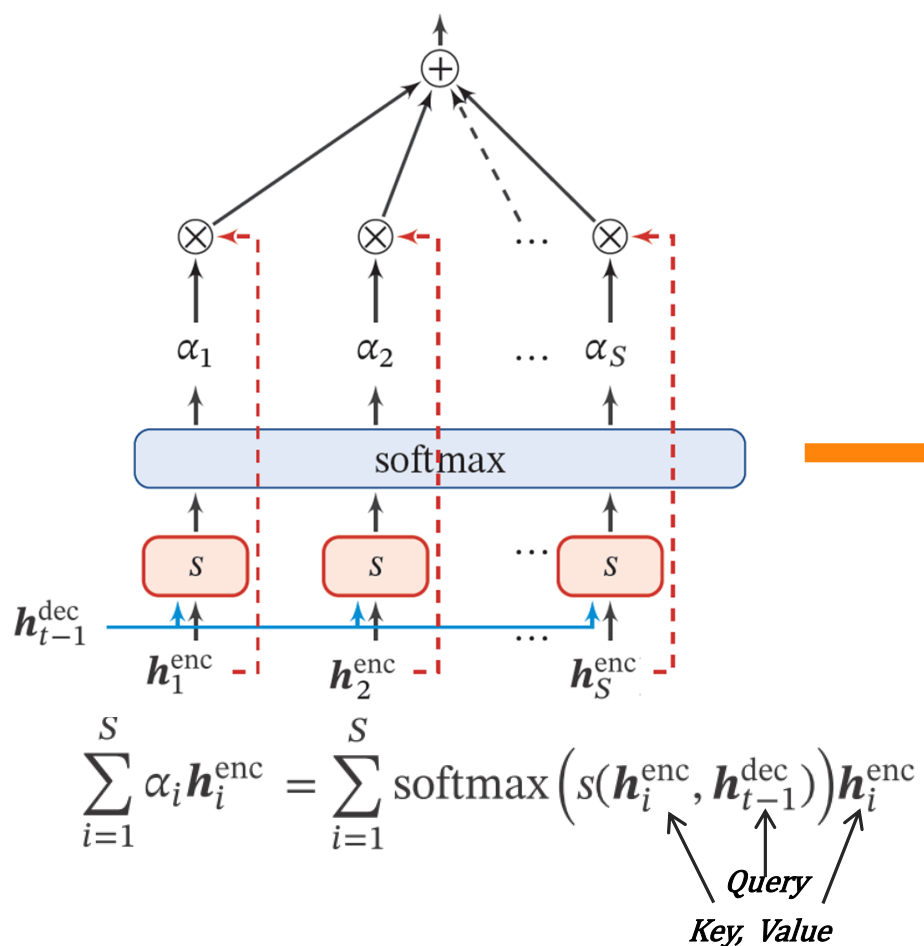
$$= \sum_{i=1}^S \text{softmax}(s(h_i^{\text{enc}}, h_{t-1}^{\text{dec}})) h_i^{\text{enc}}$$

$$h_t^{\text{dec}} = f_{\text{dec}}(h_{t-1}^{\text{dec}}, [e_{y_{t-1}}; c_t], \theta_{\text{dec}})$$

$$o_t = g(h_t^{\text{dec}}, \theta_o)$$

注意力机制

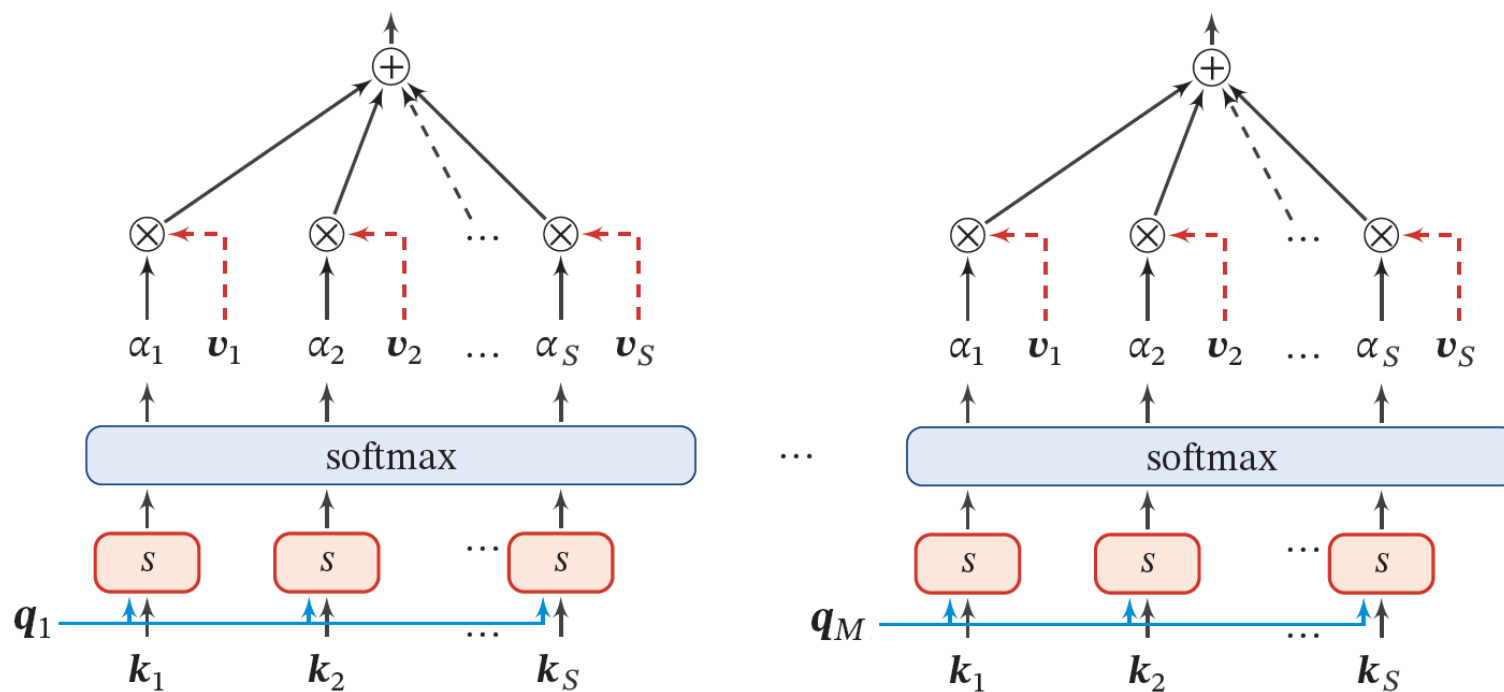
• 键值对注意力



当 $K=V$ 时，键值对模式就等价于普通的注意力机制

注意力机制

• 多头注意力



利用多个查询 $Q = [q_1, \dots, q_M]$

$$\text{att}((K, V), Q) = \text{att}((K, V), q_1) \oplus \dots \oplus \text{att}((K, V), q_M)$$

注意力机制

• 自注意力

