LSTM $(h_t, C_t)$

Short ↙  ↘ long.

① Outcome depends on short-term

$$O_t = \sigma \left( W_o \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_o \right)$$  Output gate

$$\mathbb{R}^{d_c \times (d_h + d_x)}$$

the same,

② Short-term State Depends on
long term State

$$h_t = O_t * \tanh(C_t)$$

↑ $\mathbb{R}^{d_c}$      ↑ $\mathbb{R}^{d_c}$

$h_t \longrightarrow$ produce a distribution on $V$.

③ long-term States depend on short-term state & previous periods

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

↳ forget gate    ↳ input gate

④
$$f_t = \mathcal{b}\left(W_f \cdot \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_f\right)$$

⑤
$$i_t = \mathcal{b}\left(W_i \cdot \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_i\right)$$

$$\tilde{C}_t = \tanh\left(W_c \cdot \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_c\right)$$

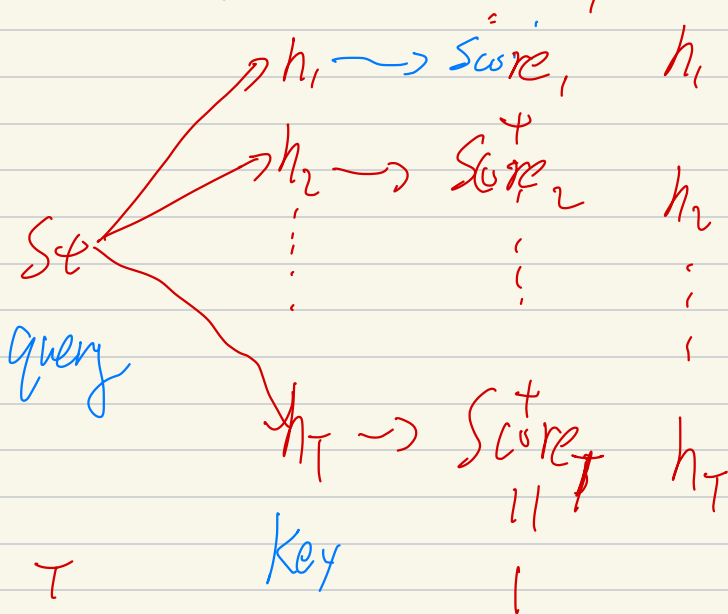# Attention Mechanism

Encoder
$$h_1 \to h_2 \to h_3 \to \cdots \to h_T$$

Decoder
$$S_t$$

$h_1 \rightsquigarrow Score_1 \quad h_1$

$h_2 \rightsquigarrow Score_2 \quad h_2$

$S_t$

query

$h_T \rightsquigarrow Score_T \quad h_T$

Key

$$a_t = \sum_{j=1}^{T} Score_j \cdot h_j$$

$\uparrow$
$\mathbb{R}^d$

why is Attention great?

$$Score = Softmax(e)$$

$\uparrow$
$\mathbb{R}^T$

$$e_i = h_i^T \cdot S_t \in \mathbb{R}$$

$$\binom{a_t}{S_t}$$

$$x_1, x_2 \cdots \cdots, x_n \in \mathbb{R}^d$$

$$W_q, \quad W_k, \quad W_v \in \mathbb{R}^{d \times d}$$

$$\in \mathbb{R}^d$$

$$q_i = W_q \cdot x_i, \quad k_i = W_k \cdot x_i, \quad v_i = W_v \cdot x_i$$

query                   key              value

$$W'$$

$$W'_{ij} = \frac{q'_i \cdot k_j}{\sqrt{d}} \qquad i, j \in [n]$$

Softmax

$$W_{ij} = \frac{\exp(W'_{ij})}{\sum_{j'=1}^{n} \exp(W'_{i,j'})}$$

$$y_i = \sum_{j=1}^{n} W_{ij} \, v_j \in \mathbb{R}^d \qquad \text{linear in } v.$$

How can we do better?

① Better parallelization?

Multi-head attention.

② non-linear $v$?

Add an MLP layer.

③ How to keep the sequence information?

Position encoding / embedding.

$input = \vec{X} + Position\ Encoding$

④ No future information? / Auto-regressive

masked attention

⑤ Can we pray to Optimization God better?

(i) Skip-Connection

$$f(x) + x$$

② Layer Normalization