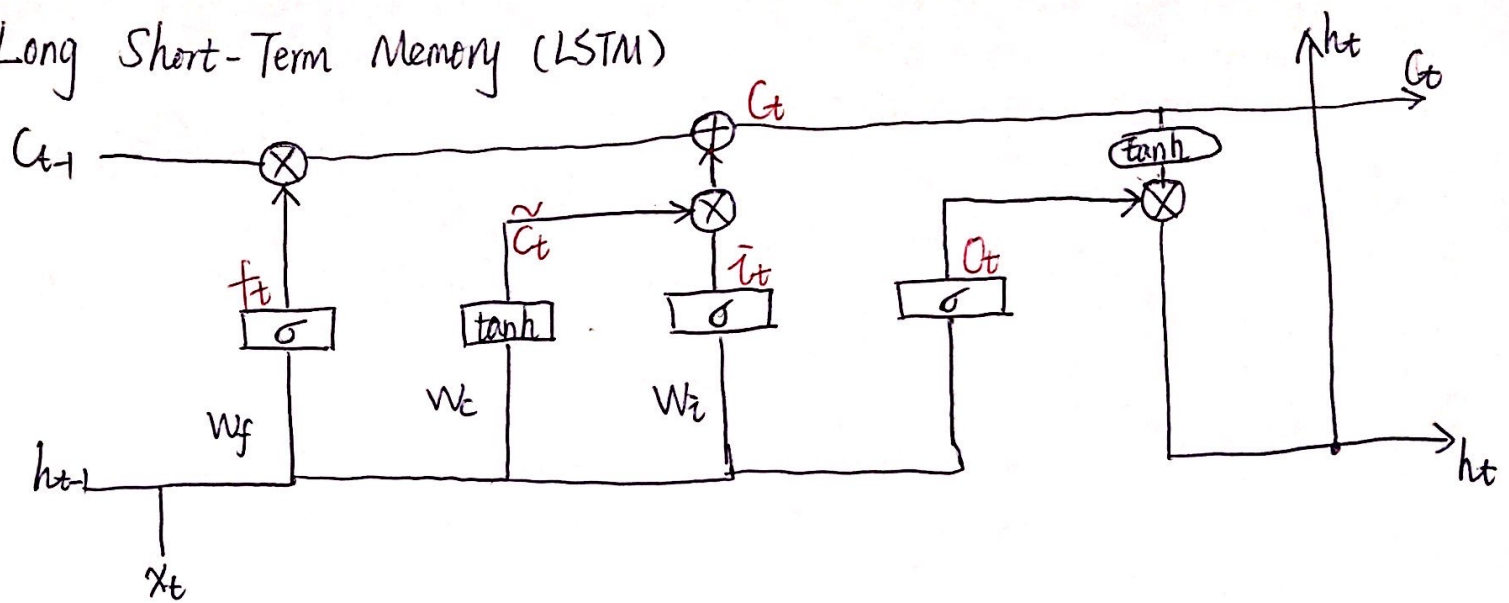# Long Short-Term Memory (LSTM)



$C_t$ ~~Att~~ : knowledge encoded in $C_t$ captures long-term dependencies and relations in the sequential ~~order~~ data

$h_t$ : predictive vectors (hidden state)

3 gates : ~~forget~~$\overset{f_t}{forget}$ / $\overset{i_t}{input}$ / $\overset{O_t}{output}$ gates

$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t]) \qquad\qquad \widetilde{C_t} = \tanh (W_c \cdot [h_{t-1}, x_t])$$

$$i_t = \sigma (W_i \cdot [h_{t-1}, x_t]) \qquad\qquad C_t = f_t \otimes C_{t-1} \oplus i_t \otimes \widetilde{C_t}$$

$$O_t = \sigma (W_o \cdot [h_{t-1}, x_t]) \qquad\qquad h_t = O_t \otimes \tanh (C_t)$$

BPTT :

$$\frac{\partial E_k}{\partial W} = \frac{\partial E_k}{\partial h_k} \cdot \frac{\partial h_k}{\partial C_k} \cdot \left( \prod_{t=2}^{k} \frac{\partial C_t}{\partial C_{t-1}} \right) \cdot \frac{\partial C_1}{\partial W} \quad ②$$

$$\frac{\partial C_t}{\partial C_{t-1}} = \frac{\partial}{\partial C_{t-1}} (C_{t-1} \otimes f_t) + \frac{\partial}{\partial C_{t-1}} (\widetilde{C_t} \otimes i_t)$$

$$= \boxed{\frac{\partial f_t}{\partial C_{t-1}} C_{t-1}} + \boxed{\frac{\partial C_{t-1}}{\partial C_{t-1}} \cdot f_t} + \boxed{\frac{\partial i_t}{\partial C_{t-1}} \cdot \widetilde{C_t}} + \boxed{\frac{\partial \widetilde{C_t}}{\partial C_{t-1}} \cdot i_t}$$
$$\qquad\qquad\qquad ① \qquad\qquad ② \qquad\qquad ③ \qquad\qquad ④$$

$$= \sigma'(W_f \cdot [h_{t-1}, x_t]) \cdot W_f \cdot O_{t-1} \otimes \tanh'(C_{t-1}) \cdot C_{t-1} \quad ① = A$$

$$+ f_t \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad ② = B$$

$$+ \sigma'(W_i \cdot [h_{t-1}, x_t]) \cdot W_i \cdot O_{t-1} \otimes \tanh'(C_{t-1}) \cdot \widetilde{C_t} \quad ③ = C$$

$$+ \tanh'(W_c \cdot [h_{t-1}, x_t]) \cdot W_c \cdot O_{t-1} \otimes \tanh'(C_{t-1}) \cdot i_t \quad ④ = D$$