

Long-short term memory (LSTM) 1997.

GRU: 2014.

1. long-term dependencies (LTD)

Q: Can this standard RNN achieve the LTD?

Yes: \rightarrow in theory

No: \rightarrow in practice.

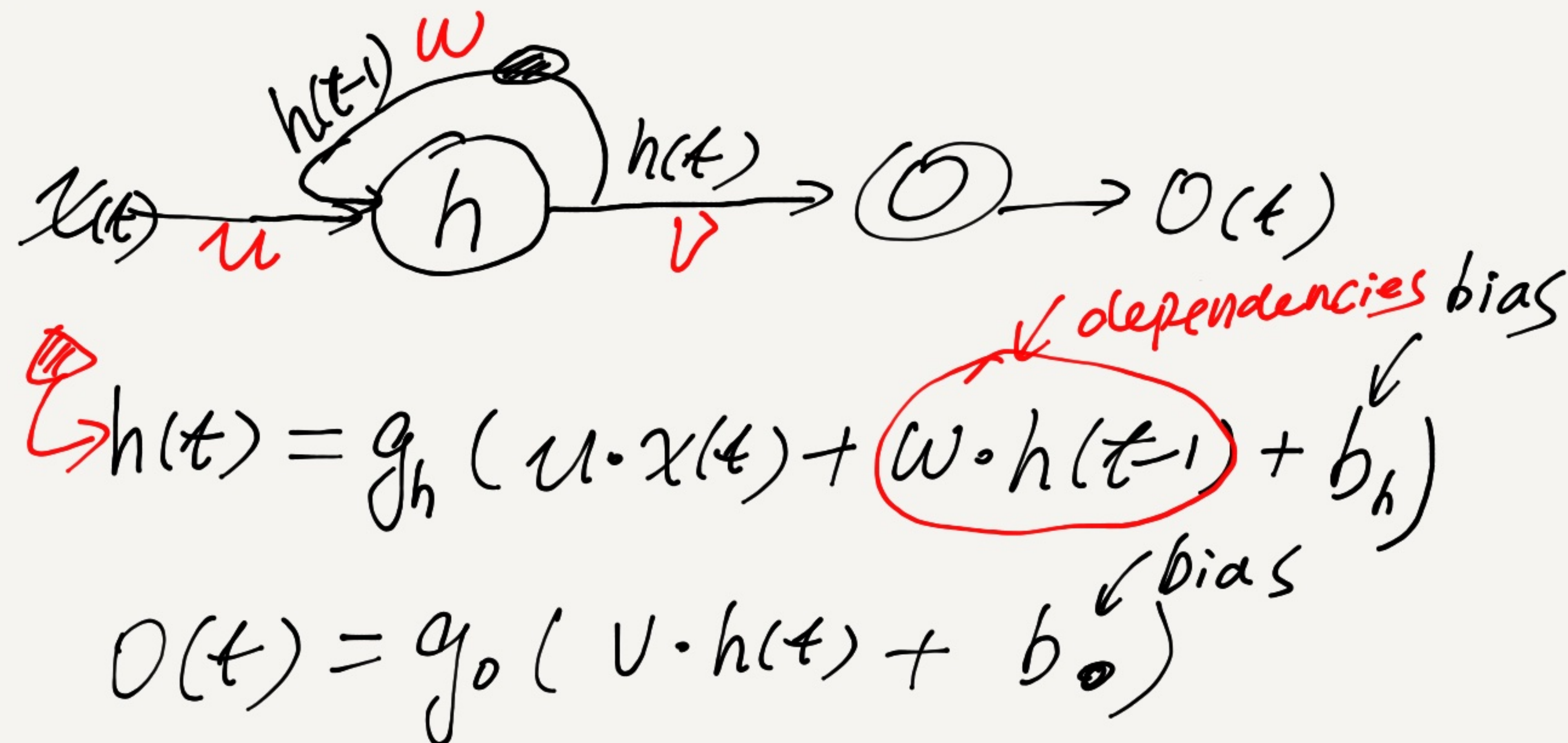
\downarrow If t is large (long chain),

we apply g_h many times,

$$h(t) \approx g_h(g_h(g_h(\dots)))$$

suppose g_h is sigmoid ($g_h \in [0,1]$)

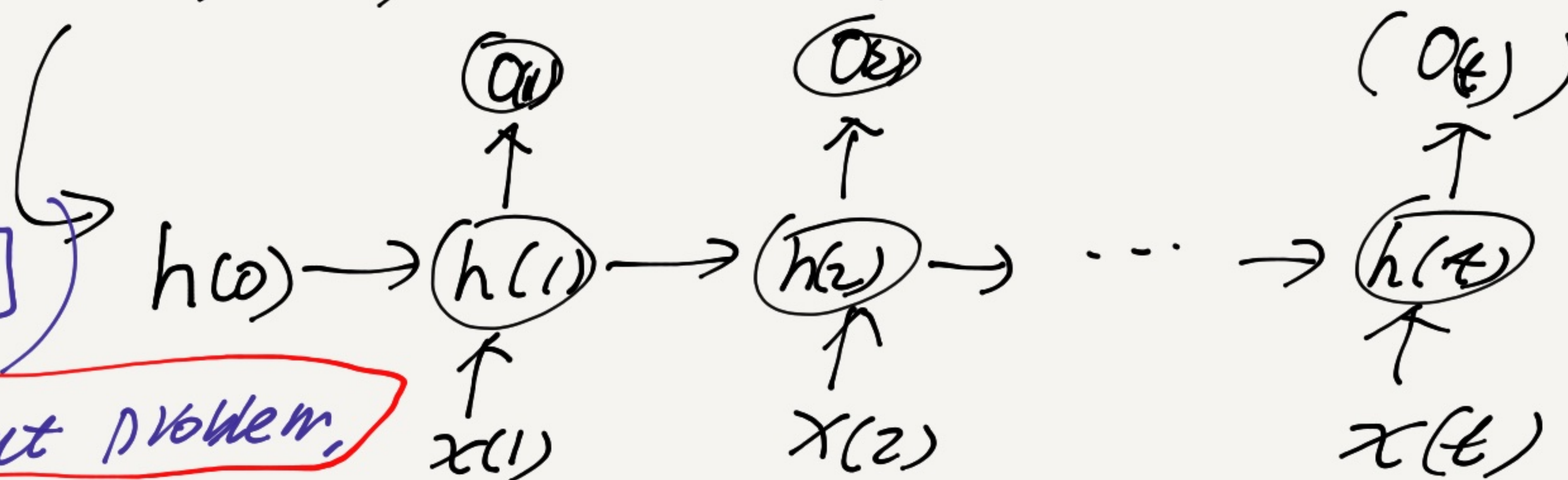
We will have the Vanishing gradient problem.



$h(t-1)$ to encode information from past sequence.

$h(t)$ calculation is a long chain from.

$h(0), x(1), h(1), x(2), \dots$ to $h(t-1), x(t)$

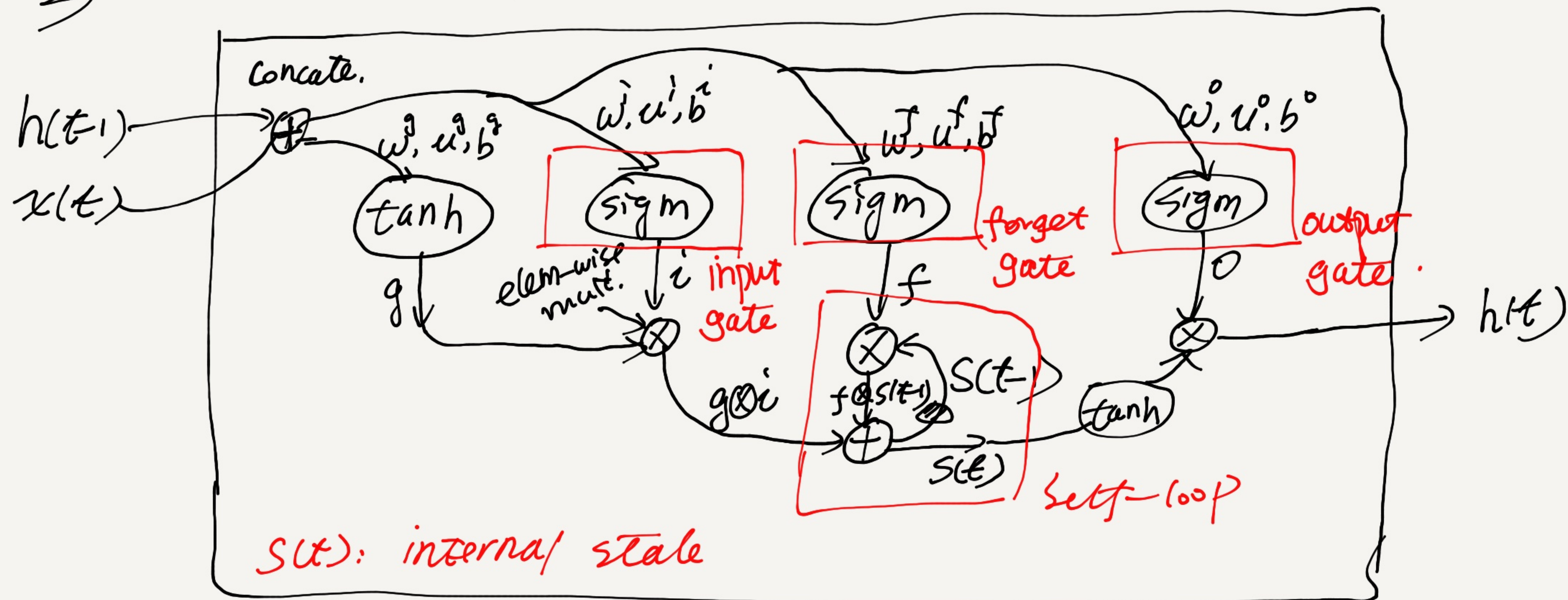


2. LSTM (1997)

D. main idea of LSTM:

- (1) Introduced 3 gates: input gate, forget gate, output gate
- (2) Introduced an internal state and self-loop to produce path to let gradient flow for a long term.

2) LSTM cell.



\oplus : concatenation \odot : element-wise multiplication

$$\underline{g} = \tanh(\omega^g \cdot h(t-1) + u^g \cdot x(t) + b^g) \in [-1, 1] \rightarrow \text{Candidate}$$

new information

$$\text{input gate: } \hat{i} = \text{sigm}(\omega^i \cdot h(t-1) + u^i \cdot x(t) + b^i) \in [0, 1]$$

$$\text{forget gate: } f = \text{sigm}(\omega^f \cdot h(t-1) + u^f \cdot x(t) + b^f) \in [0, 1]$$

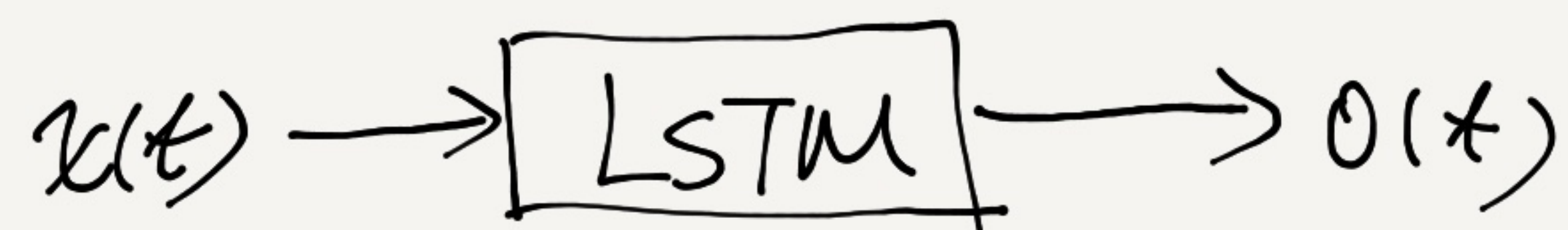
Controls the fraction of data to forget,

$$\text{internal state: } S(t) = \underbrace{(f \otimes S(t-1))}_{\text{forgetting}} + \underbrace{(\hat{i} \otimes g)}_{\text{input new information}}$$

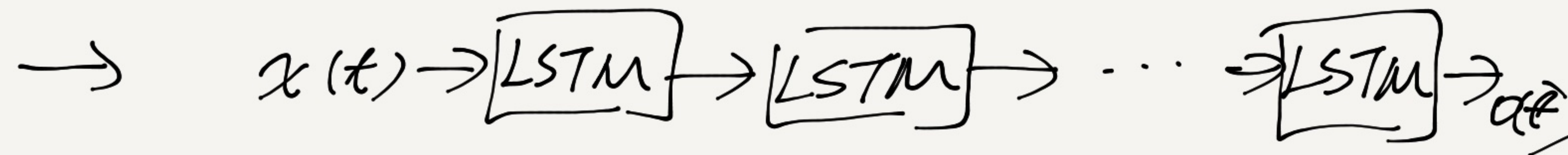
$$\begin{matrix} f & S \\ \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \otimes & \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \end{matrix}$$

$$\text{Output gate: } O = \text{sigm}(\omega^o \cdot h(t-1) + u^o \cdot x(t) + b^o)$$

3. LSTM-based RNN.



single LSTM cell



multiple LSTM cells