

$$L = -\ln(p_c) = -\ln(\text{softmax}(y_c))$$

$$\frac{\partial L}{\partial p_c} = -\frac{1}{p_c}$$

$$\frac{\partial p_c}{\partial y_c} = \frac{e^{y_c}}{\sum e^{y_i}} - \frac{e^{y_c} \cdot e^{y_c}}{(\sum e^{y_i})^2} = \frac{e^{y_c} (\sum e^{y_i} - e^{y_c})}{(\sum e^{y_i})^2} = p_c \cdot \frac{\sum e^{y_i} - e^{y_c}}{\sum e^{y_i}} = p_c (1 - p_c) \quad (p_c = \frac{e^{y_c}}{\sum e^{y_i}})$$

$$\frac{\partial L}{\partial y_c} = \frac{\partial L}{\partial p_c} \cdot \frac{\partial p_c}{\partial y_c} = p_c - 1$$

$$\frac{\partial p_c}{\partial y_i} = -\frac{e^{y_c} \cdot e^{y_i}}{(\sum e^{y_i})^2} = -p_c \cdot p_i$$

$$\frac{\partial L}{\partial y_i} = \frac{\partial L}{\partial p_c} \cdot \frac{\partial p_c}{\partial y_i} = p_i$$

$$\therefore \frac{\partial L}{\partial y_i} = \begin{cases} p_i & , i \neq c \\ p_{i-1} & , i = c \end{cases}, \frac{\partial L}{\partial y} = \begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_{c-1} \\ p_c \\ \vdots \\ p_n \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial y_1} \\ \vdots \\ \frac{\partial L}{\partial y_n} \end{bmatrix} = \begin{bmatrix} p_1 \\ p_{c-1} \end{bmatrix} \otimes \begin{bmatrix} p_{i-1} \\ p_n \end{bmatrix}$$

$$Y = W_{hy} h + b_y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}, \quad y_1 = w_{hy_1} h + b_{y_1}, \quad y_2 = w_{hy_2} h + b_{y_2}$$

$$\frac{\partial Y}{\partial W_{hy}} = \begin{bmatrix} \frac{\partial y_1}{\partial w_{hy_1}} \\ \frac{\partial y_1}{\partial w_{hy_2}} \\ \vdots \\ \frac{\partial y_n}{\partial w_{hy_n}} \end{bmatrix} = \begin{bmatrix} c & h & \dots \\ c & h & \dots \end{bmatrix}$$

$$\frac{\partial L}{\partial W_{hy}} = \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial W_{hy}} = \begin{bmatrix} \frac{\partial L}{\partial w_{hy_1}} \\ \vdots \\ \frac{\partial L}{\partial w_{hy_n}} \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial y_1} \frac{\partial y_1}{\partial w_{hy_1}} \\ \vdots \\ \frac{\partial L}{\partial y_n} \frac{\partial y_n}{\partial w_{hy_n}} \end{bmatrix} = \begin{bmatrix} p_1 \cdot [ & h & ] \\ (p_{c-1}) \cdot [ & h & ] \end{bmatrix} = \frac{\partial L}{\partial Y} \otimes \frac{\partial Y}{\partial W_{hy}} = \begin{bmatrix} p_1 \\ p_{c-1} \end{bmatrix} \otimes \begin{bmatrix} h \\ \vdots \\ h \end{bmatrix}$$

$$\frac{\partial Y}{\partial b_y} = 1$$

$$\frac{\partial L}{\partial b_y} = \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial b_y} = \begin{bmatrix} p_1 \\ p_{c-1} \end{bmatrix}$$

$$\frac{\partial Y}{\partial h} = W_{hy} = \begin{bmatrix} \frac{\partial y_1}{\partial h} \\ \vdots \\ \frac{\partial y_n}{\partial h} \end{bmatrix} = \begin{bmatrix} w_{hy_1} \\ \vdots \\ w_{hy_n} \end{bmatrix} = \begin{bmatrix} c & \dots \\ c & \dots \end{bmatrix} = \frac{\partial Y}{\partial h_n}$$

2x64

$$h_i = \tanh(w_{xh} x_i + w_{hh} h_{i-1} + b_h)$$

$$h_h = \tanh(w_{xh} x_h + w_{hh} h_{h-1} + b_h), \quad h_{hi} = \tanh(w_{xh_i} x_h + w_{hh_i} h_{h-1} + b_{hi}) = \begin{bmatrix} \vdots \end{bmatrix}_{64 \times 1}$$

$$\frac{\partial h_n}{\partial b_h} = 1 - h_n^2 = \begin{bmatrix} \frac{\partial h_n}{\partial b_h} \\ \vdots \\ \frac{\partial h_n}{\partial b_h} \end{bmatrix} = \begin{bmatrix} 1 - h_{n1}^2 \\ 1 - h_{n2}^2 \\ \vdots \\ 1 - h_{n64}^2 \end{bmatrix}_{64 \times 1}$$

$$\frac{\partial L}{\partial b_h} = \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial h_n} \frac{\partial h_n}{\partial b_h} = p_i \otimes (p_{i-1}) \cdot \begin{bmatrix} w_{hy_1} \\ \vdots \\ w_{hy_n} \end{bmatrix} \cdot \begin{bmatrix} 1 - h_{n1}^2 \\ 1 - h_{n2}^2 \\ \vdots \\ 1 - h_{n64}^2 \end{bmatrix} = \frac{1}{m} \sum_{i=1}^m \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial h_n} \frac{\partial h_n}{\partial b_h} = \frac{1}{2} (p_1 \cdot \begin{bmatrix} w_{hy_1} \\ \vdots \\ w_{hy_n} \end{bmatrix} \cdot \begin{bmatrix} 1 - h_{n1}^2 \\ 1 - h_{n2}^2 \\ \vdots \\ 1 - h_{n64}^2 \end{bmatrix} + (p_{c-1}) \cdot \begin{bmatrix} w_{hy_1} \\ \vdots \\ w_{hy_n} \end{bmatrix} \cdot \begin{bmatrix} 1 - h_{n1}^2 \\ 1 - h_{n2}^2 \\ \vdots \\ 1 - h_{n64}^2 \end{bmatrix})$$

$$= \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial h_j} \frac{\partial h_j}{\partial b_h} = \frac{1}{N} \sum_{j=1}^N \left( \frac{1}{m} \sum_{i=1}^m \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial h_j} \frac{\partial h_j}{\partial b_h} \right)$$

$$\frac{\partial h_n}{\partial w_{xh}} = (1 - h_n^2) X_h = \begin{bmatrix} \frac{\partial h_n}{\partial w_{xh}} \\ \vdots \\ \frac{\partial h_n}{\partial w_{xh}} \end{bmatrix} = \begin{bmatrix} (1 - h_{n1}^2) X_h \\ (1 - h_{n2}^2) X_h \\ \vdots \\ (1 - h_{n64}^2) X_h \end{bmatrix} = \begin{bmatrix} c & \dots \\ c & \dots \\ \vdots & \vdots \end{bmatrix}_{64 \times 1}$$

64x1

$$\frac{\partial L}{\partial w_{xh}} = \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial h_n} \frac{\partial h_n}{\partial w_{xh}} = \frac{1}{m} \sum_{i=1}^m \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial h_n} \frac{\partial h_n}{\partial w_{xh}} = \frac{1}{2} ([p_1, p_{c-1}] \otimes \begin{bmatrix} w_{hy_1} \\ \vdots \\ w_{hy_n} \end{bmatrix}) \otimes \begin{bmatrix} c (1 - h_{n1}^2) X_h \\ c (1 - h_{n2}^2) X_h \\ \vdots \\ c (1 - h_{n64}^2) X_h \end{bmatrix}_{64 \times 1}$$

$$= \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial h_j} \frac{\partial h_j}{\partial w_{xh}} = \frac{1}{N} \sum_{j=1}^N \left( \frac{1}{m} \sum_{i=1}^m \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial h_j} \frac{\partial h_j}{\partial w_{xh}} \right)$$

$$\frac{\partial h_n}{\partial w_{hh}} = (1 - h_n^2) h_{h-1} = \begin{bmatrix} \frac{\partial h_n}{\partial w_{hh}} \\ \vdots \\ \frac{\partial h_n}{\partial w_{hh}} \end{bmatrix} = \begin{bmatrix} (1 - h_{n1}^2) h_{h-1} \\ (1 - h_{n2}^2) h_{h-1} \\ \vdots \\ (1 - h_{n64}^2) h_{h-1} \end{bmatrix} = \begin{bmatrix} c & \dots \\ c & \dots \\ \vdots & \vdots \end{bmatrix}_{64 \times 1}$$

64x1

$$\frac{\partial L}{\partial w_{hh}} = \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial h_n} \frac{\partial h_n}{\partial w_{hh}} = \frac{1}{m} \sum_{i=1}^m \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial h_n} \frac{\partial h_n}{\partial w_{hh}} = \frac{1}{2} ([p_1, p_{c-1}] \otimes \begin{bmatrix} w_{hy_1} \\ \vdots \\ w_{hy_n} \end{bmatrix}) \otimes \begin{bmatrix} c (1 - h_{n1}^2) h_{h-1} \\ c (1 - h_{n2}^2) h_{h-1} \\ \vdots \\ c (1 - h_{n64}^2) h_{h-1} \end{bmatrix}_{64 \times 1}$$

$$= \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial h_j} \frac{\partial h_j}{\partial w_{hh}} = \frac{1}{N} \sum_{j=1}^N \left( \frac{1}{m} \sum_{i=1}^m \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial h_j} \frac{\partial h_j}{\partial w_{hh}} \right)$$

$$\frac{\partial L}{\partial h_{h-1}} = \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial h_n} \frac{\partial h_n}{\partial h_{h-1}} = \frac{\partial L}{\partial h_n} (1 - h_n^2) w_{hh}$$

随时间反向传播算法 BPTT  
Backpropagation Through Time

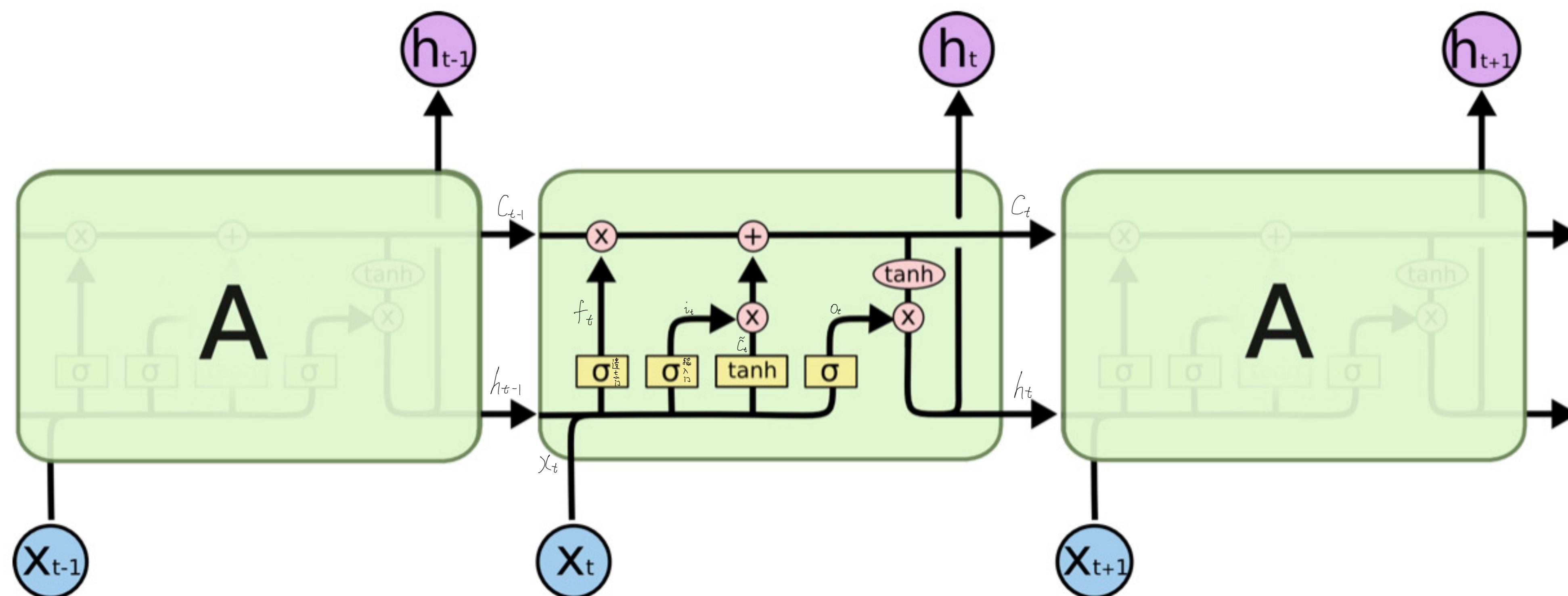
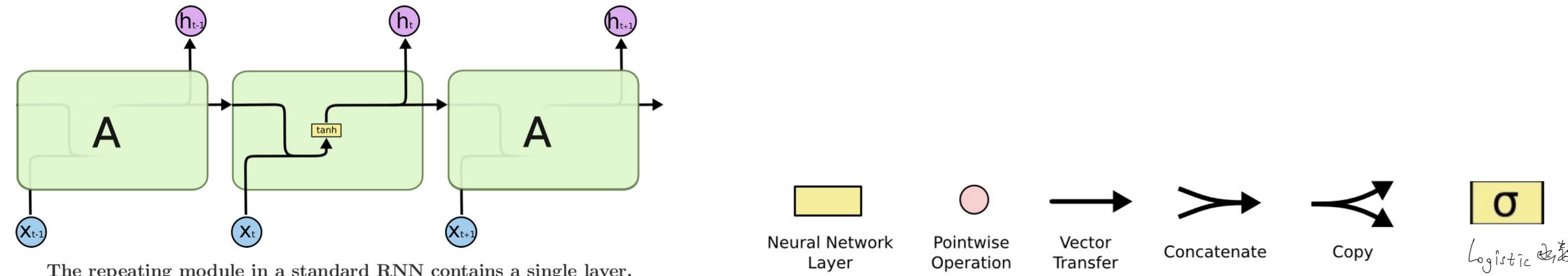
计算梯度的反向传播随时间反向传播 BPTT

实时循环学习RTL：通过前向传播来计算梯

# 长颈鹿板块问题与梯度消失

小林度爆炸：杖重衰減、極度微弱

基于门控的循环神经网络  
(Gated RNN) < 长短期记忆网络 LSTM, Long short-term Memory Network  
门控循环单元网络 GRU, Gated Recurrent Unit



The repeating module in an LSTM contains four interacting layers.

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) = \sigma(c \cdot W_f \cdot \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} + [b_f]) = \sigma(\lambda_f) = \frac{1}{1+e^{-\lambda_f}} \quad (\text{激活函数})$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) = \sigma(c \cdot W_i \cdot \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} + [b_i]) = \sigma(\lambda_i) = \frac{1}{1+e^{-\lambda_i}} \quad (\text{激活函数})$$

$$\tilde{c}_t = \tanh(W_{\tilde{c}} \cdot [h_{t-1}, x_t] + b_{\tilde{c}}) = \tanh(c \cdot W_{\tilde{c}} \cdot \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} + [b_{\tilde{c}}]) = \tanh(\lambda_{\tilde{c}}) = \frac{e^{\lambda_{\tilde{c}}}}{1+e^{\lambda_{\tilde{c}}}} - 1$$

$$\begin{bmatrix} \tilde{\mathbf{c}}_t \\ \mathbf{o}_t \\ \mathbf{i}_t \\ f_t \end{bmatrix} = \begin{bmatrix} \tanh \\ \sigma \\ \sigma \\ \sigma \end{bmatrix} \left( \mathbf{W} \begin{bmatrix} \mathbf{x}_t \\ \mathbf{h}_{t-1} \end{bmatrix} + \mathbf{b} \right)$$

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

$$Q_t = \sigma(W_0[h_{t-1}, x_t] + b_0) = \sigma(c - W_0 \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} + [b_0]) = \sigma(\lambda_0) = \frac{1}{1+e^{-\lambda_0}} \quad (\text{from } 1)$$

$$h_t = O_t \times \tanh(c_t)$$

「おは空れり」

- (1) 遗忘门  $f_t$  控制上一个时刻的内部状态  $c_{t-1}$  需要遗忘多少信息.
  - (2) 输入门  $i_t$  控制当前时刻的候选状态  $\tilde{c}_t$  有多少信息需要保存.  $C_t = f_t \star C_{t-1} + i_t \star \tilde{C}_t$
  - (3) 输出门  $o_t$  控制当前时刻的内部状态  $c_t$  有多少信息需要输出给外部状  $b_t = o_t \star \text{tanh}(c_t)$