



Recurrent Neural Networks

And Long Short Term Memory

“어제의 생각이 오늘의 당신을 만들고 오늘의 생각이 내일의 당신을 만든다.”
블레즈 파스칼



Outline



- ▶ Introduction
- ▶ Motivation
- ▶ RNN architecture
- ▶ RNN problems
- ▶ LSTM
- ▶ How LSTM solves the problem
- ▶ Paper experiments
- ▶ Conclusions



Introduction



- ▶ RNN were introduced in the late 80's.
- ▶ Hochreiter discovers the 'vanishing gradients' problem in 1991.
- ▶ Long Short Term Memory published in 1997.
- ▶ LSTM a recurrent network to overcome the problem.



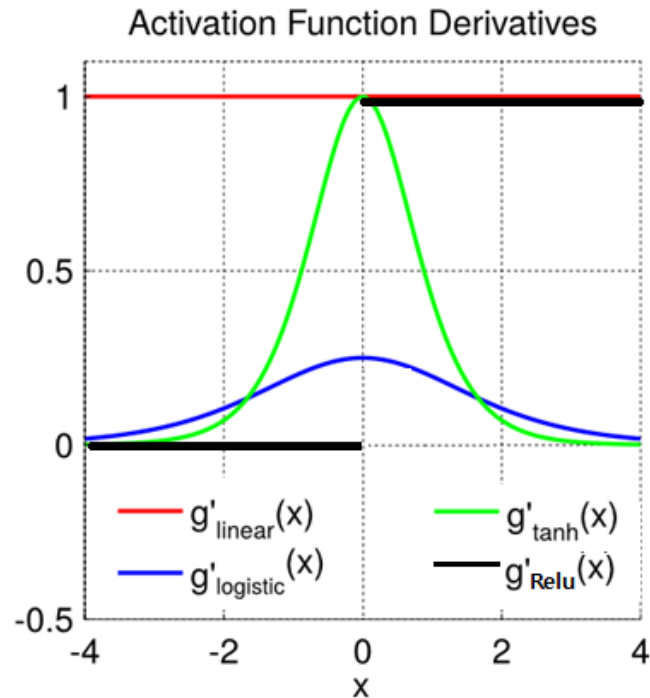
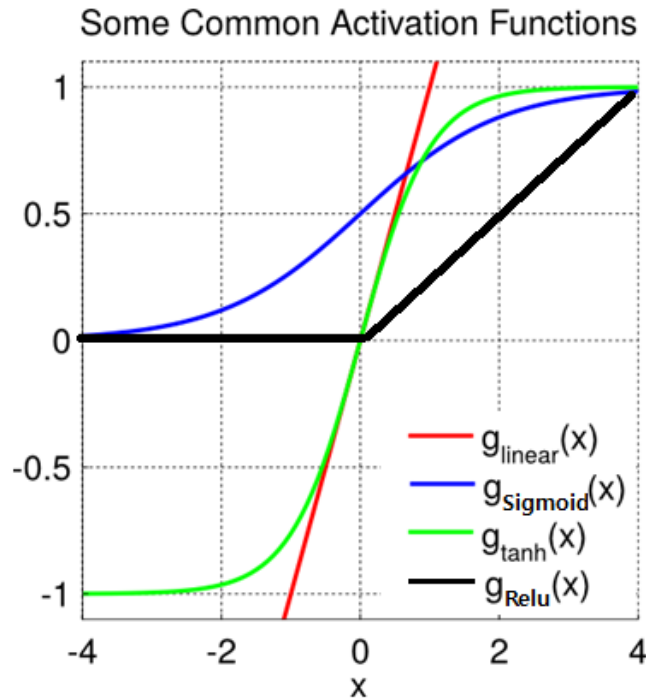
Motivation



- ▶ Feed forward networks accept a fixed-sized vector as input and produce a fixed-sized vector as output
- ▶ fixed amount of computational steps
- ▶ recurrent nets allow us to operate over *sequences of vectors*

Activation Functions

인공신경망에서 활성, 비활성을 결정하는데 사용되는 함수



▶ Activation Functions과 미분

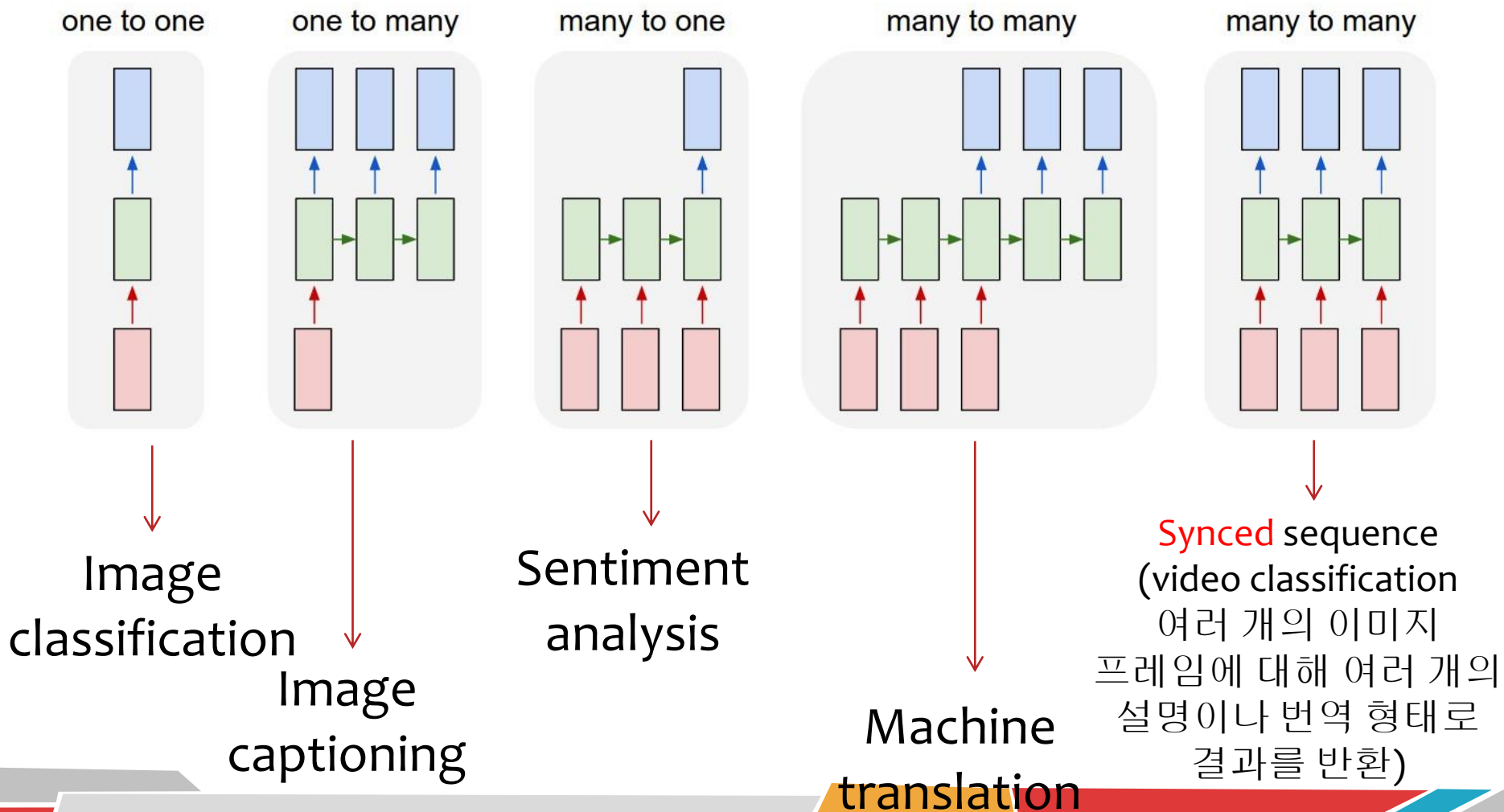
sigmoid 함수 : $\sigma(x) = 1 / (1 + e^x) \rightarrow \sigma'(x) = \sigma(x) (1 - \sigma(x))$

하이퍼볼릭 탄젠트 : $\tanh(x) \rightarrow \tanh'(x) = 2\sigma(2x) - 1 = (e^x - e^{-x}) / (e^x + e^{-x}) = 1 - \tanh^2(x)$

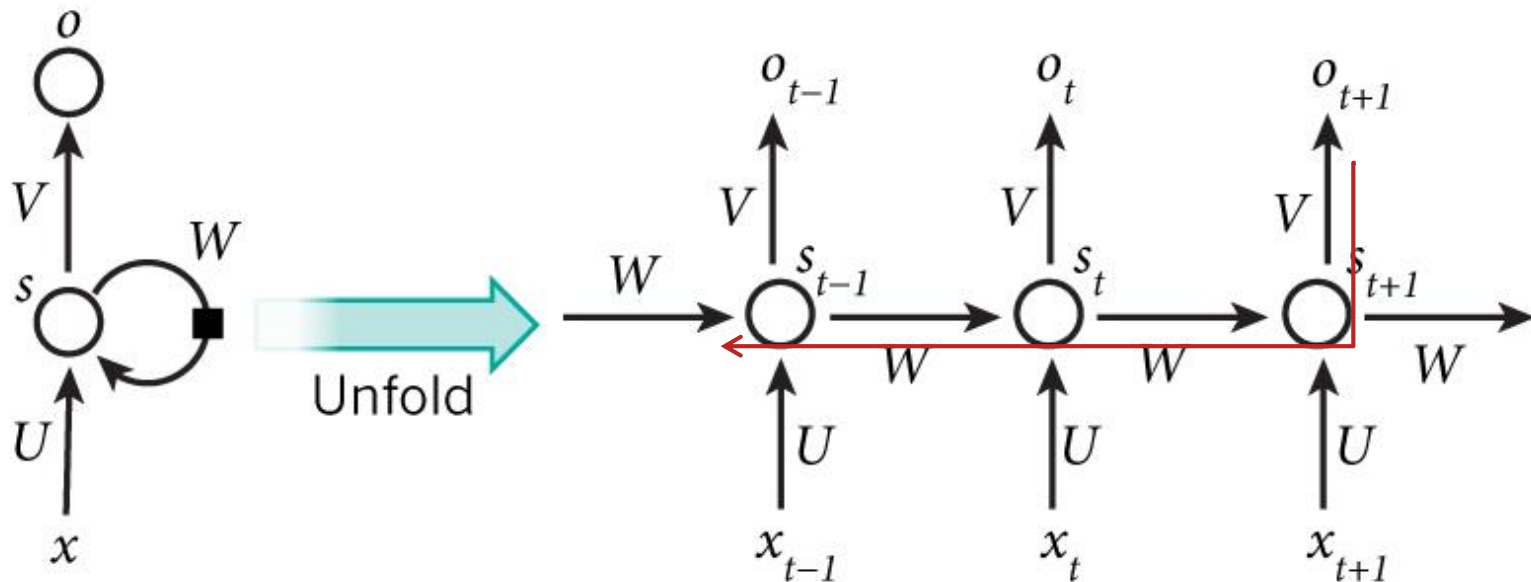
Rectified Linear Unit (ReLU) : $f(x) = \max(0, x) \rightarrow f'(x) = 0$ if $x \leq 0$, $f'(x) = 1$ if $x > 0$



Motivation



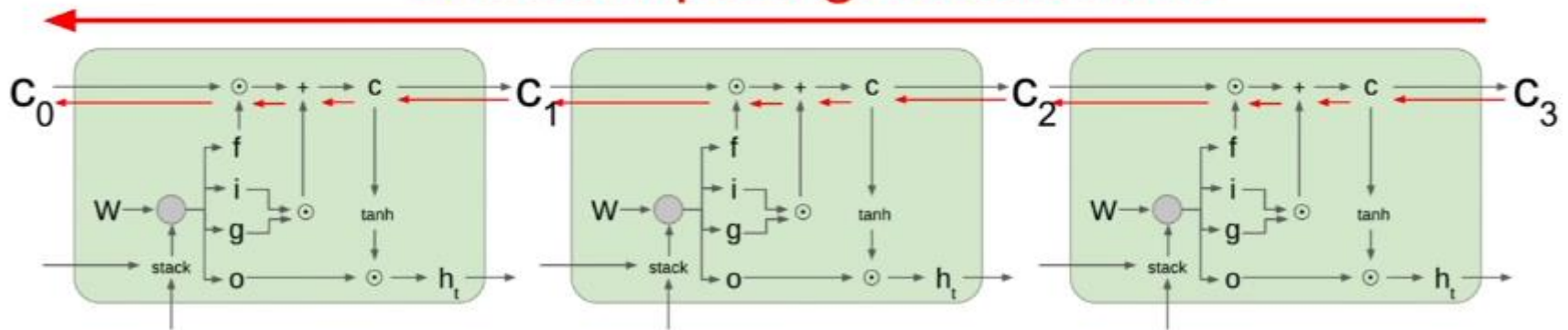
RNN Architecture



- The recurrent network can be converted into a feed forward network by **unfolding over time**
- long-term 정보 손실 문제 - LSTM(Long Short Term Memory)으로 해결

RNN Architecture

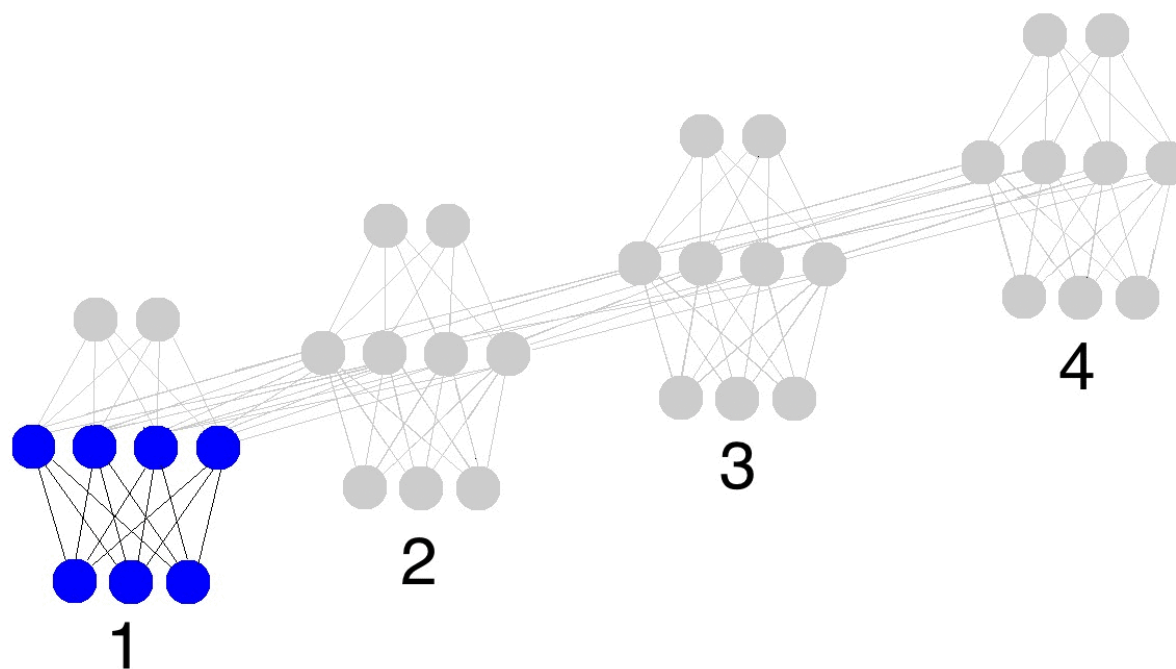
Uninterrupted gradient flow!



Gradient flows smoothly during Backprop



RNN Architecture



MakeAGIF.com

RNN Architecture

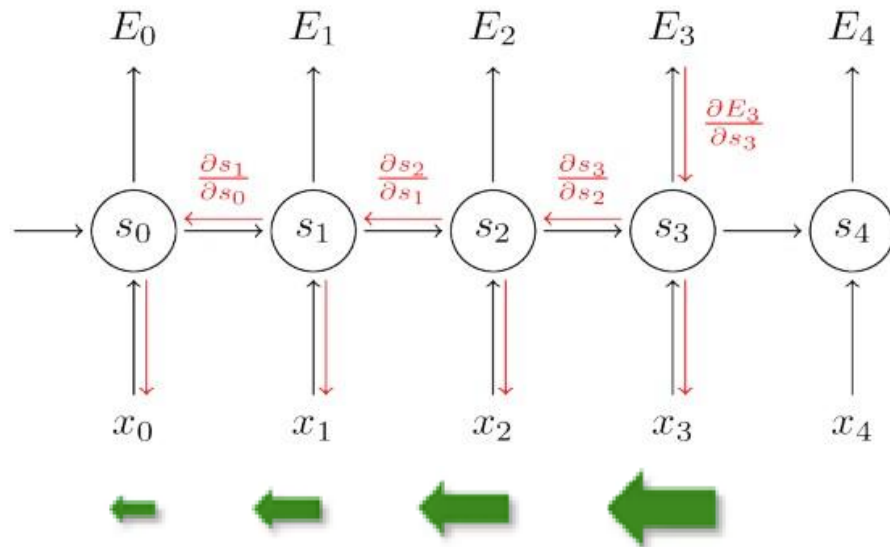
$$\frac{\partial E}{\partial \mathbf{W}} = \sum_t \frac{\partial E_t}{\partial \mathbf{W}}$$

$$\frac{\partial E_3}{\partial \mathbf{W}} = \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial \mathbf{W}}$$

But $s_3 = \tanh(Ux_t + Ws_2)$

s_3 depends on s_2 , which depends on \mathbf{W} and s_1 , and so on.

$$\frac{\partial E_3}{\partial \mathbf{W}} = \sum_{k=0}^3 \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial s_k} \frac{\partial s_k}{\partial \mathbf{W}}$$





Outline - LSTM



- ▶ Introduction
- ▶ Motivation
- ▶ RNN architecture
- ▶ RNN problems
- ▶ **Long Short Term Memory**
- ▶ How LSTM solves the problem
- ▶ Paper experiments
- ▶ Conclusions

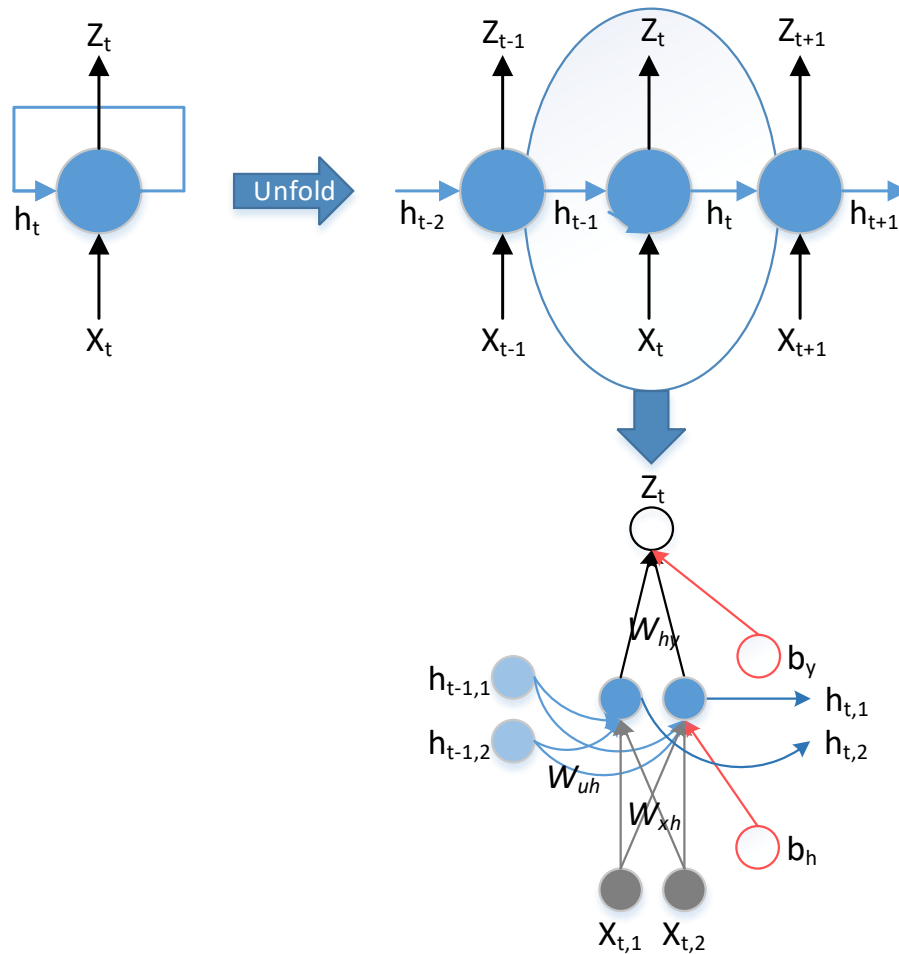


LSTM - introduction

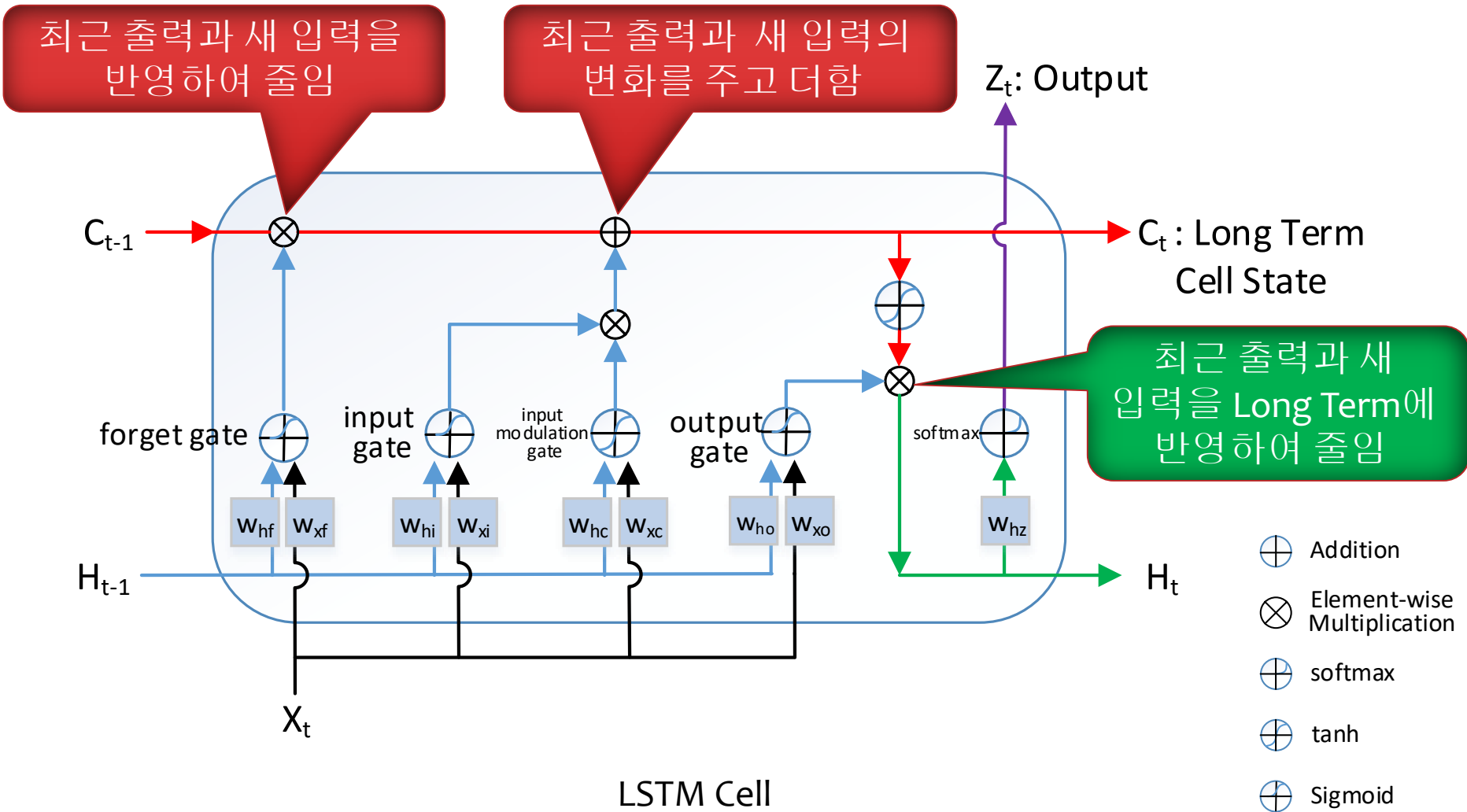


- ▶ LSTM was invented to solve the **vanishing gradients** problem.
- ▶ LSTM maintain a more constant error flow in the backpropagation process.
- ▶ LSTM can learn over more than 1000 time steps , and thus can handle large sequences that are linked remotely.

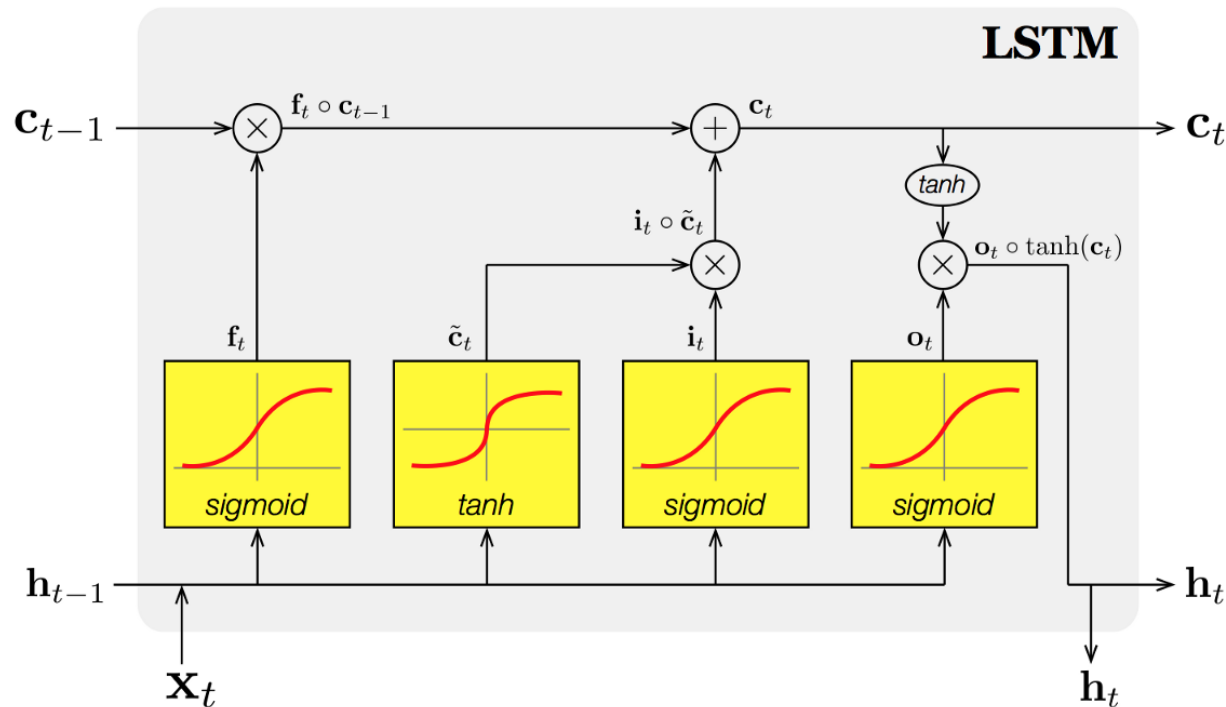
RNN Architecture



○ ○ ○ LSTM Architecture - overview ○ ○ ○



○ ○ ○ LSTM Architecture - overview ○ ○ ○



Gating variables

$$f_t = \sigma(\mathbf{W}_f[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_f)$$

$$i_t = \sigma(\mathbf{W}_i[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_i)$$

$$o_t = \sigma(\mathbf{W}_o[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_o)$$

Candidate (memory) cell state

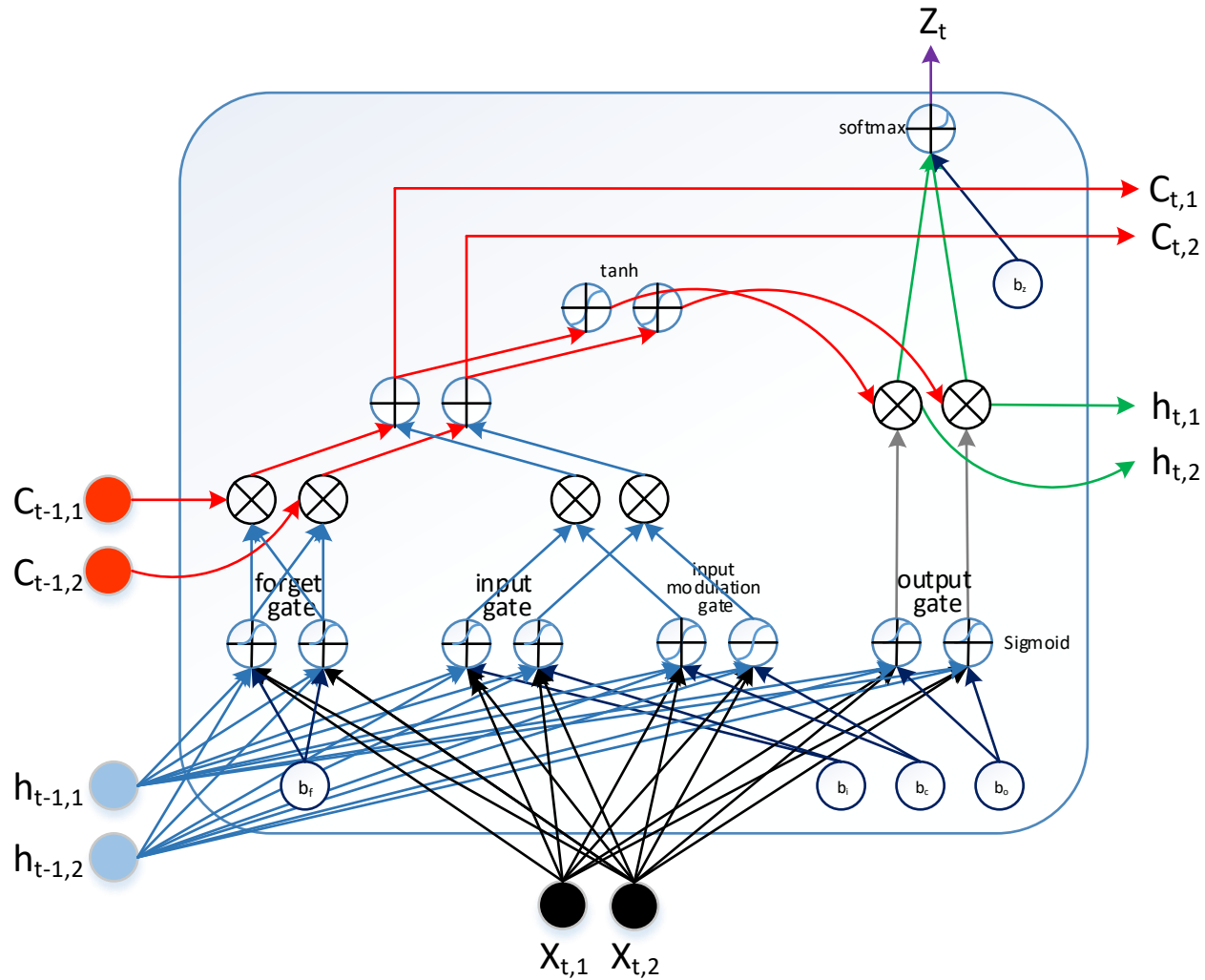
$$\tilde{c}_t = \tanh(\mathbf{W}_c[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_c)$$

Cell & Hidden state

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

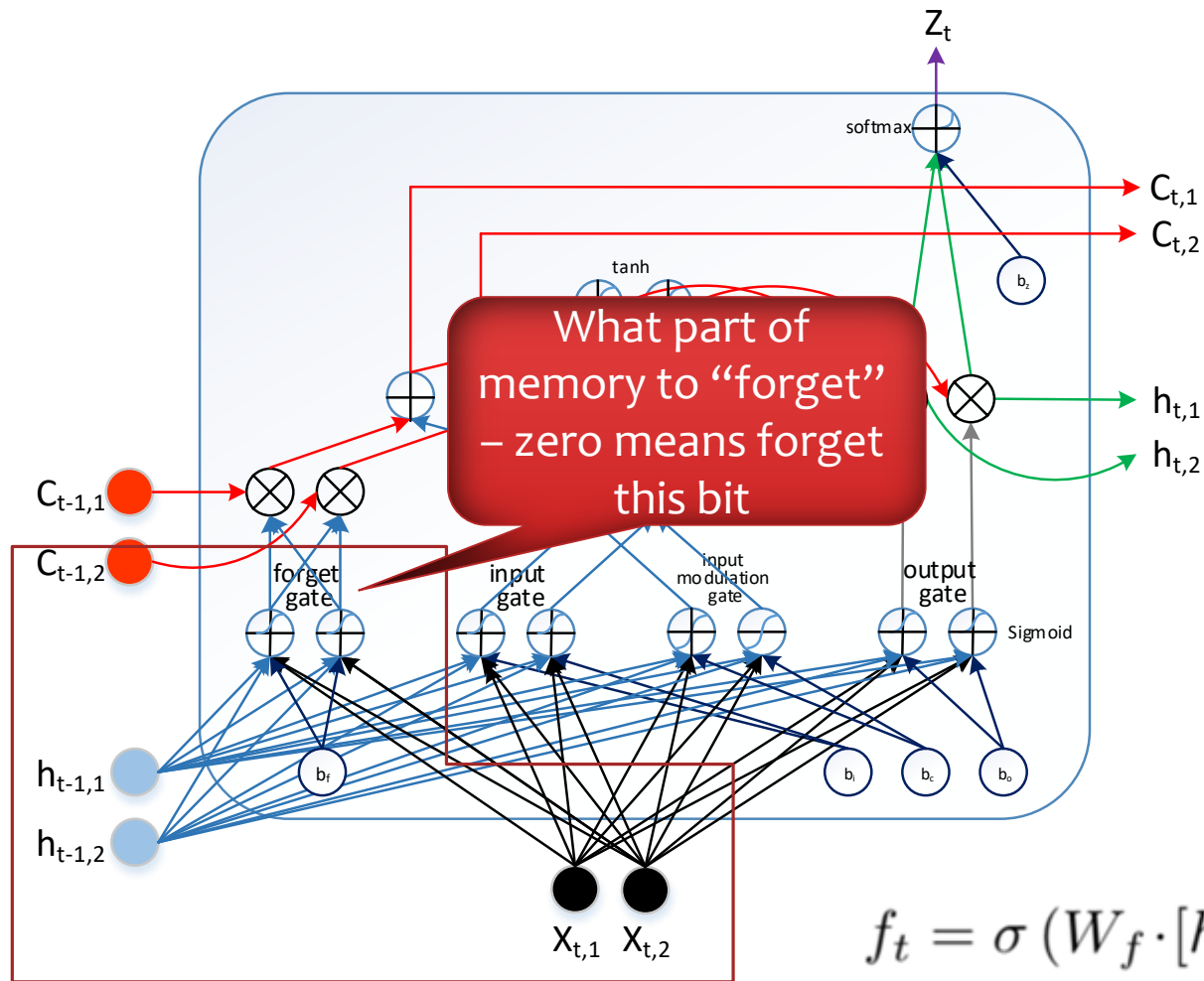
$$h_t = o_t \circ \tanh(c_t)$$

○ ○ ○ LSTM Architecture - overview ○ ○ ○

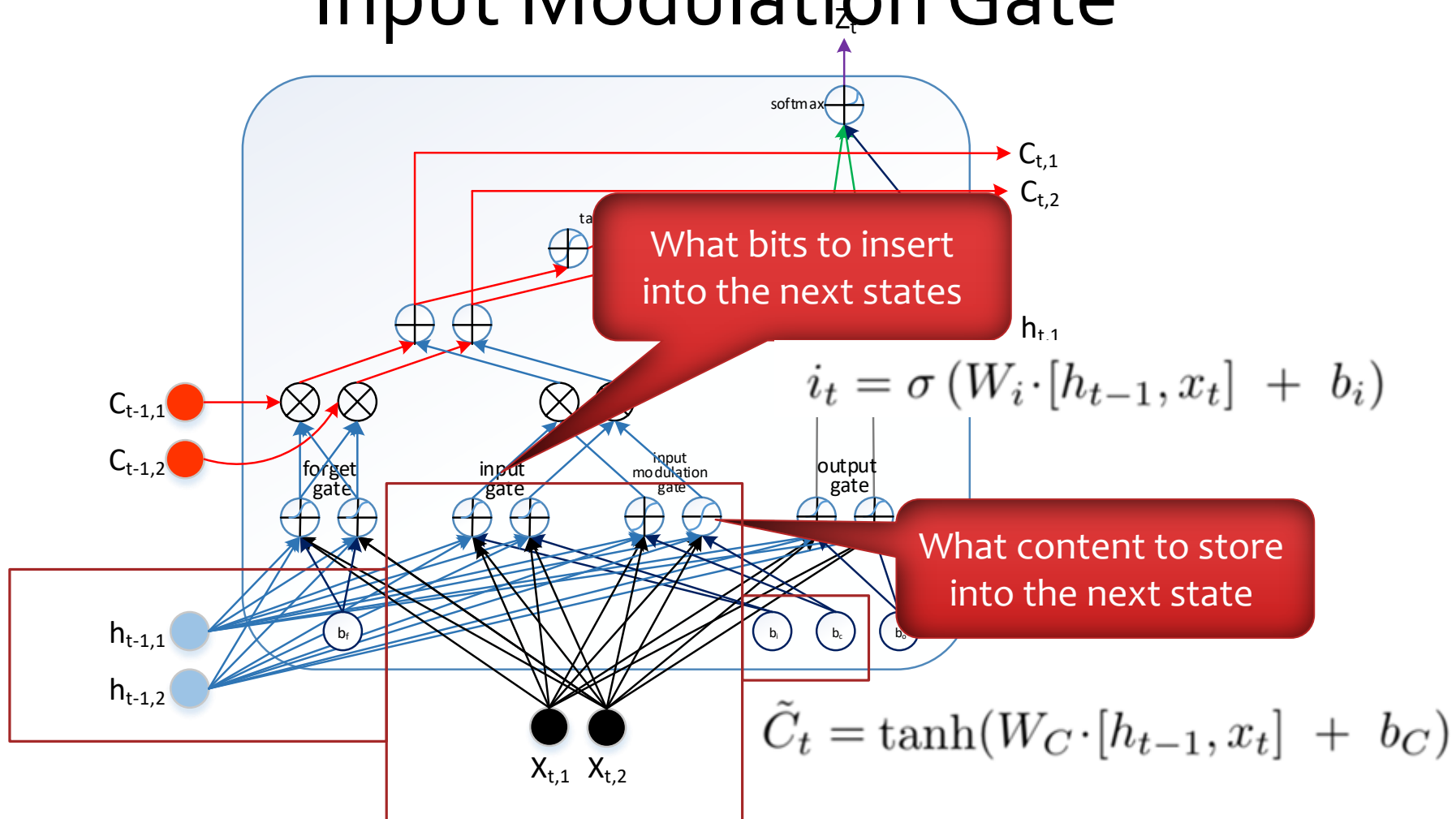


LSTM Cell

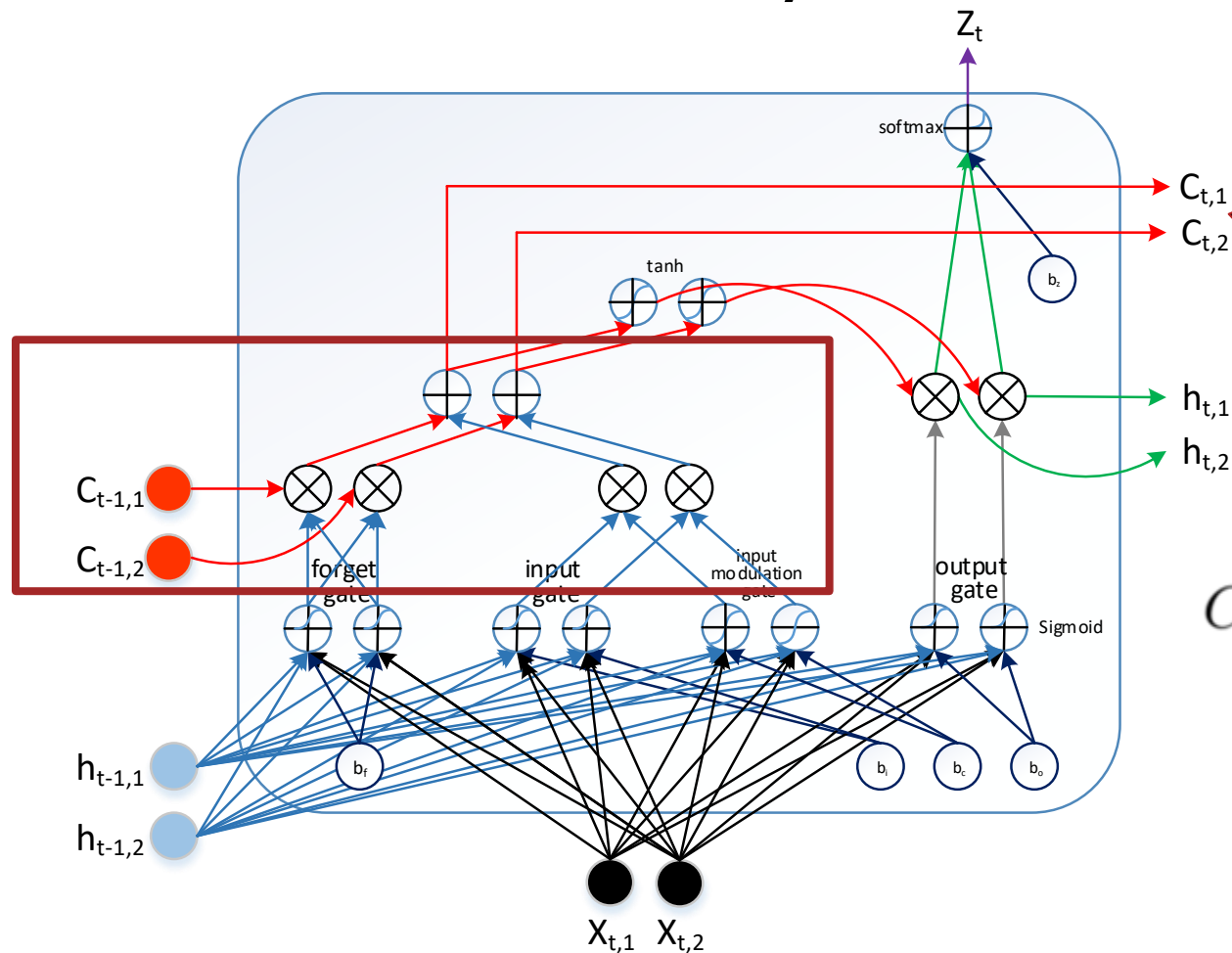
LSTM – Forget Gate



LSTM – Input Gate and Input Modulation Gate



LSTM – Long Term memory cell content



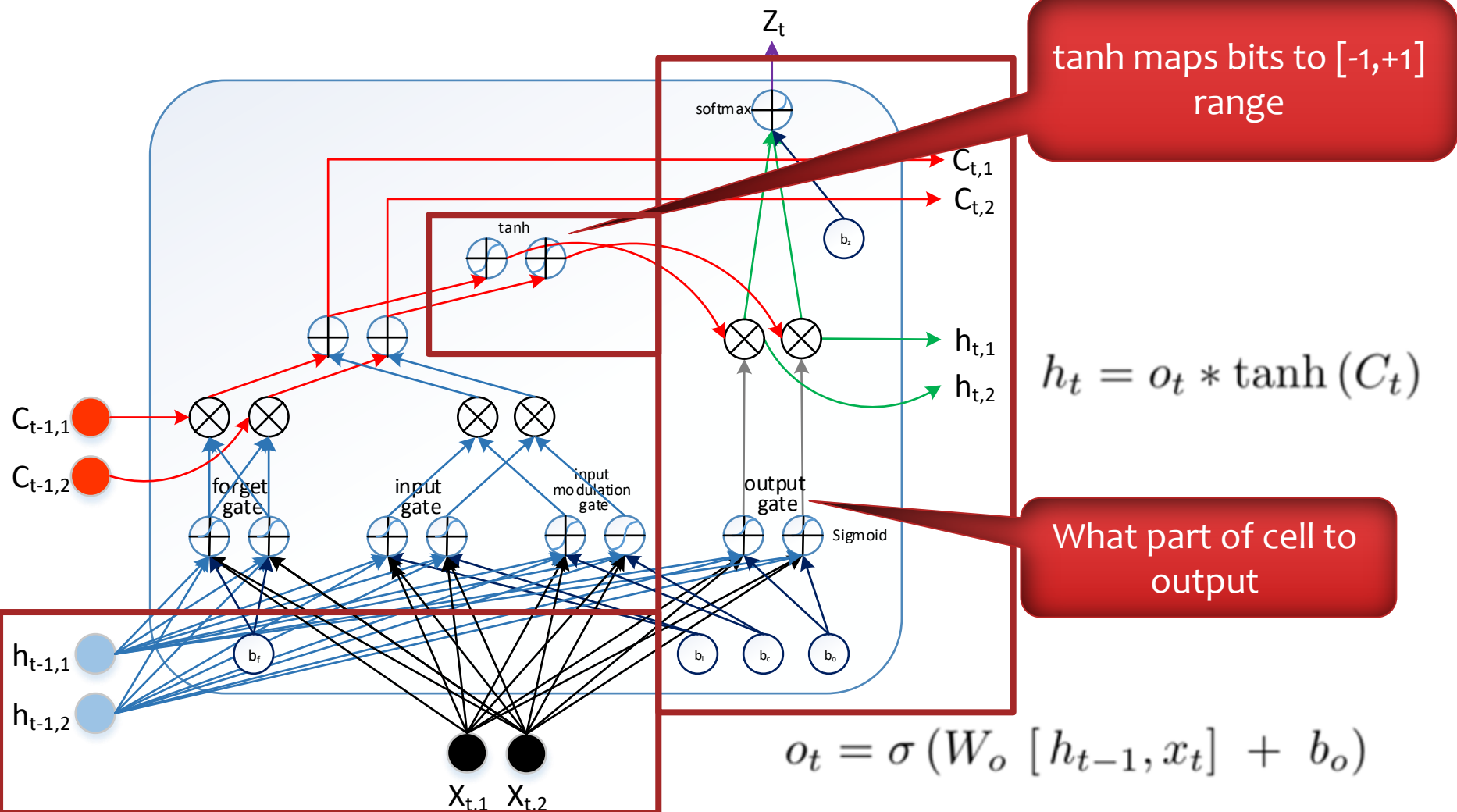
Next memory cell
content – mixture of
not-forgotten part of
previous cell and
insertion

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

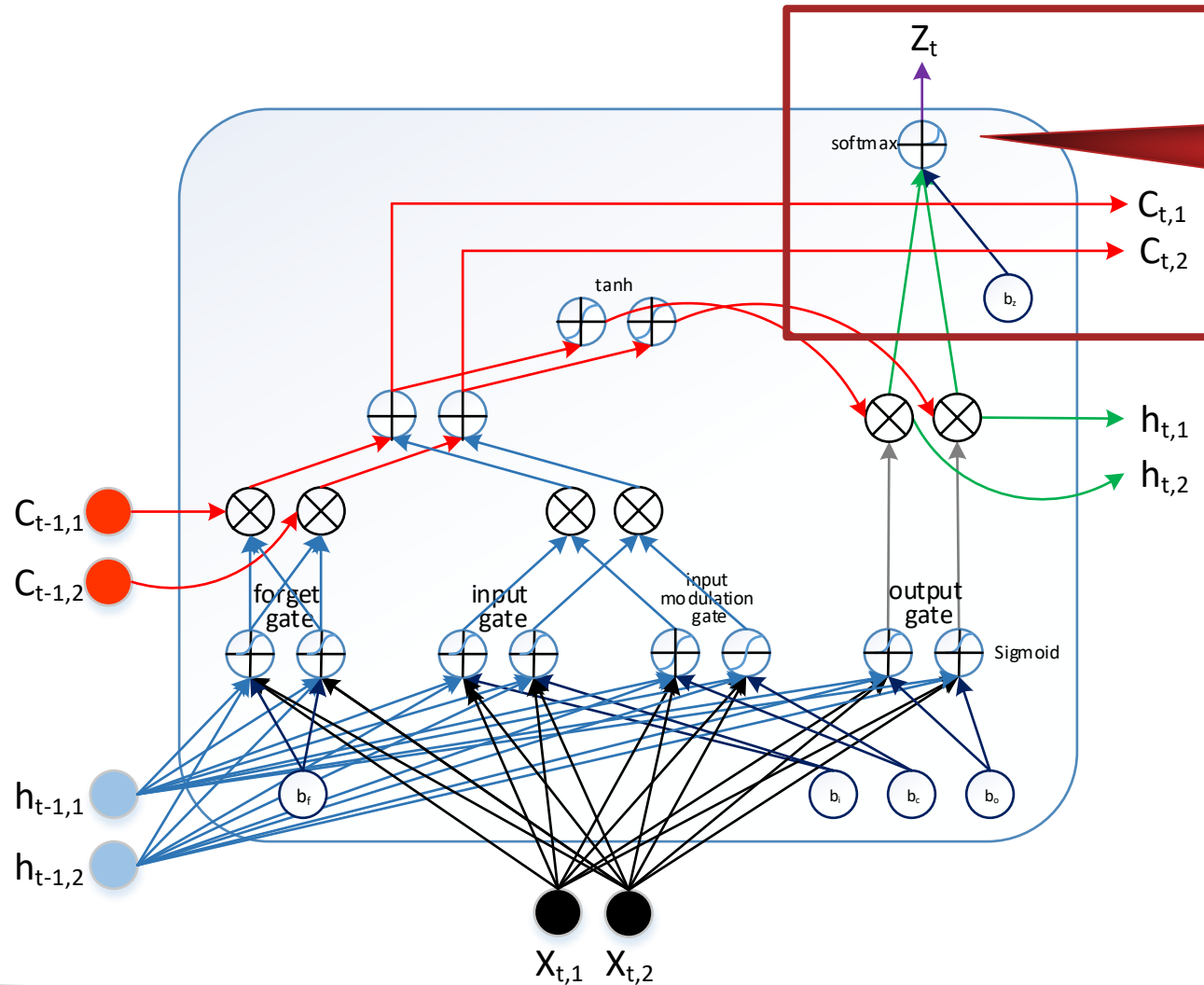
\oplus 상승 효과

\otimes 억제 효과

○ ○ ○ LSTM – Output to Next Cell ○ ○ ○



○○○ LSTM – Softmax Output ○○○

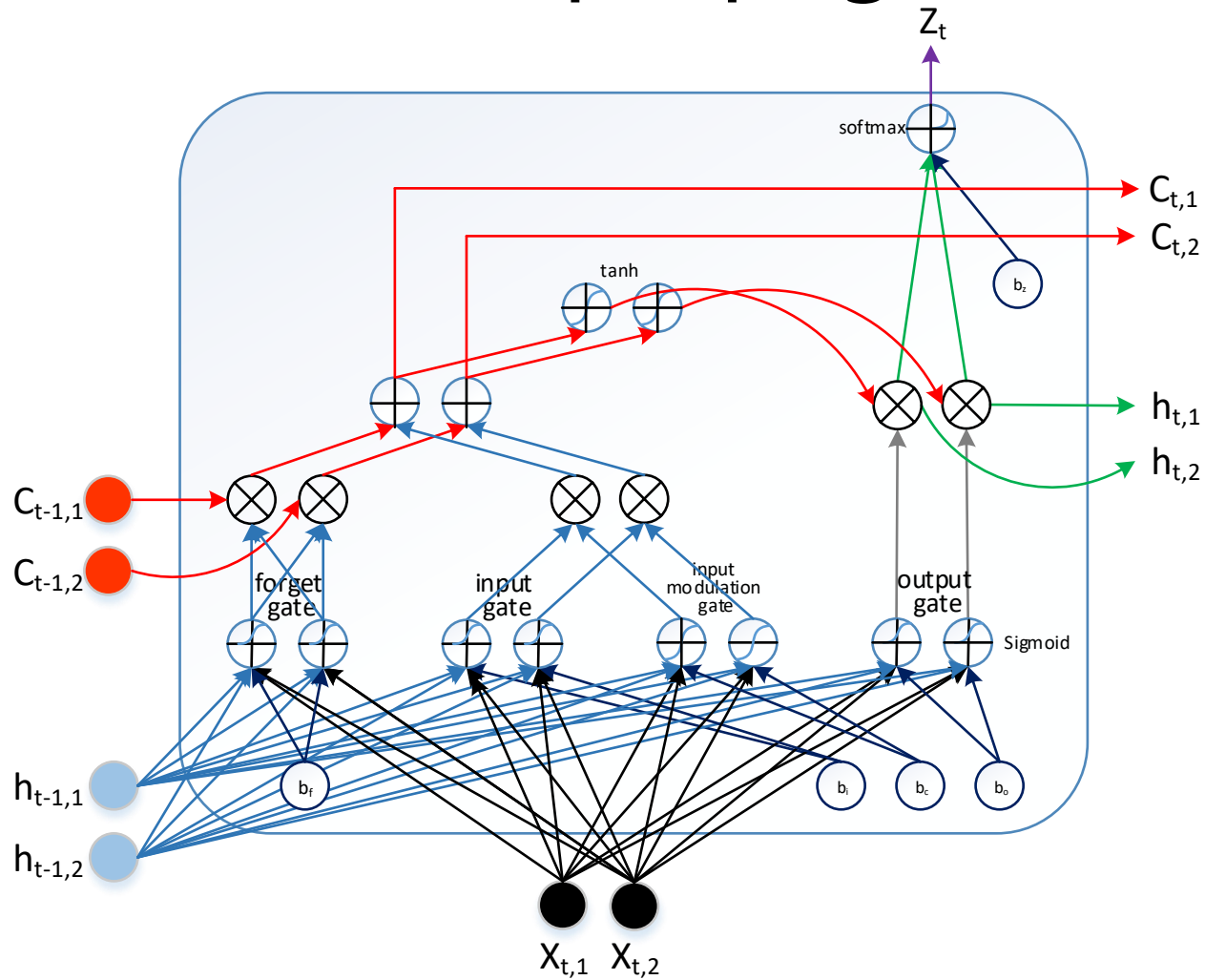


Softmax를 구하여
목표값과의 오차 계산

$$Z_t = s(w_{hz} + b_z)$$

$$s(x_j) = \frac{e^{x_j}}{\sum_i e^{x_i}}$$

○○○ LSTM Backpropogation ○○○



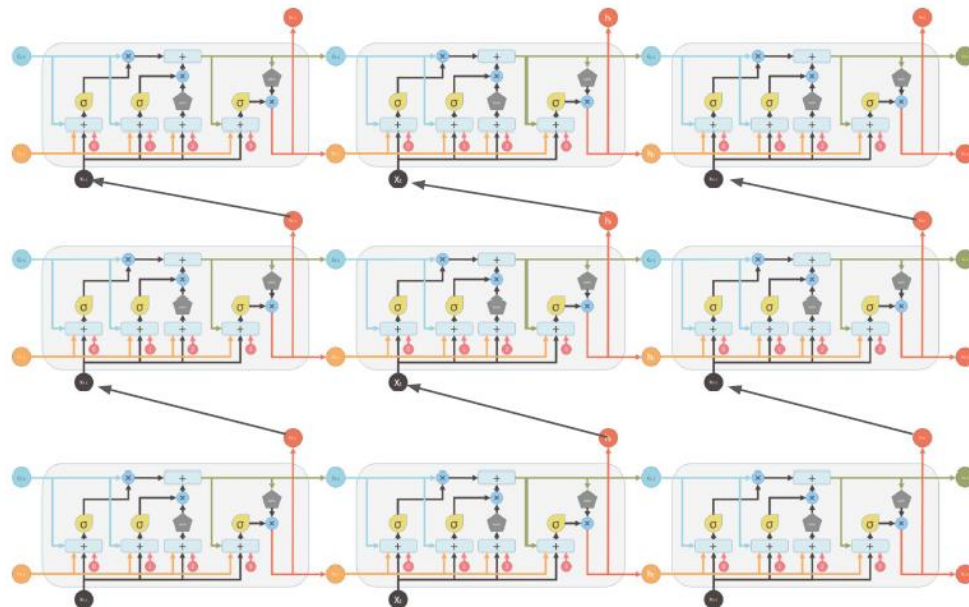
LSTM Cell



Deep LSTM



- Deep LSTMs can be created by stacking multiple LSTM layers vertically, with the output sequence of one layer forming the input sequence of the next (in addition to recurrent connections within the same layer)
- Increases the number of parameters - but given sufficient data, performs significantly better than single-layer LSTMs (Graves et al. 2013)
- Dropout usually applied only to non-recurrent edges, including between layers



LSTM Cell

Advantages of LSTM

- ▶ For long time lag problems, LSTM can handle noise and continuous values.
- ▶ No parameter fine tuning.
- ▶ Memory for long time periods