# GRU

*Learning Phrase Representations using RNN Encoder–Decoder*
*for Statistical Machine Translation*

황주훈

# INDEX

# 01. 배경지식
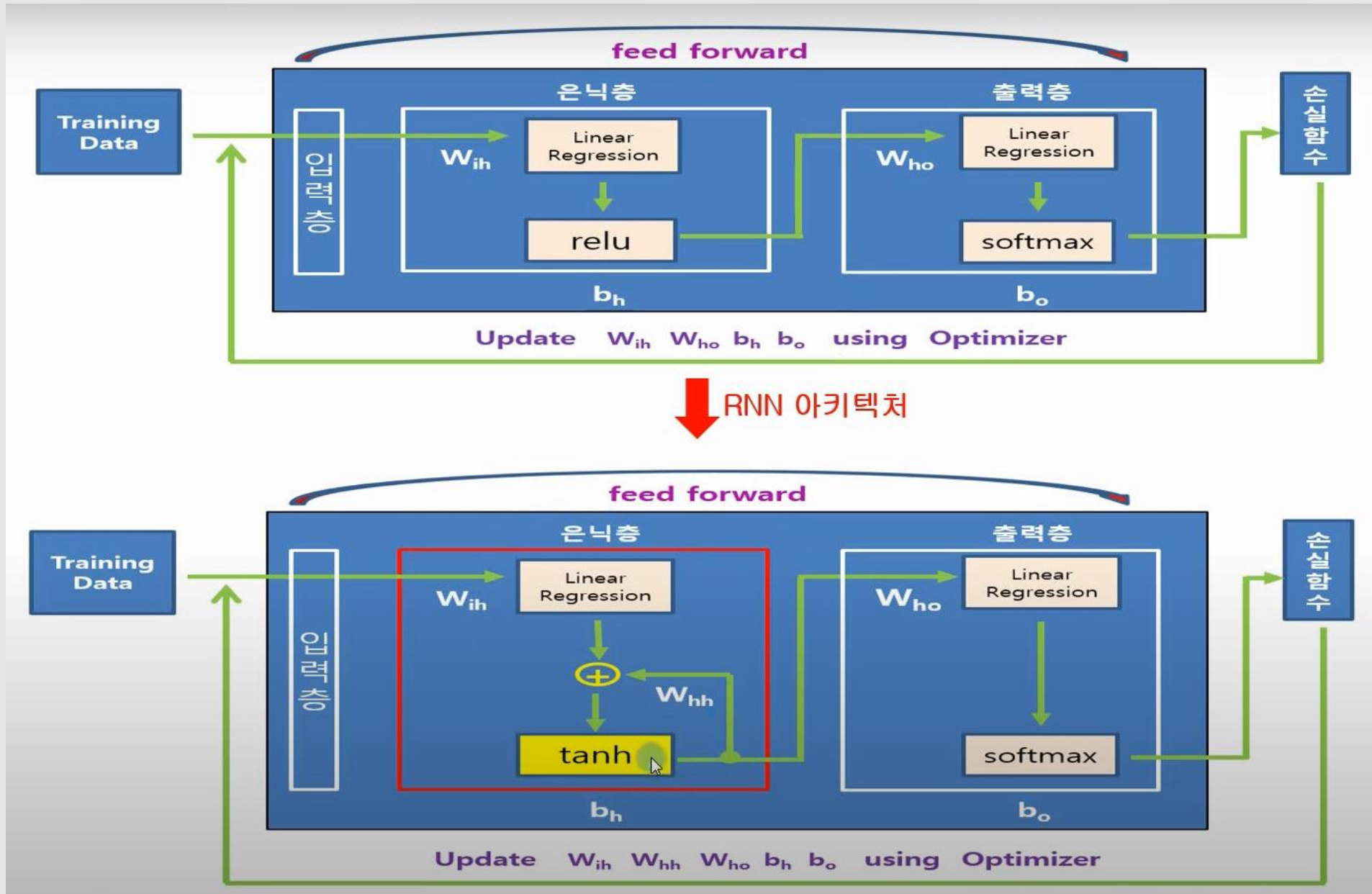
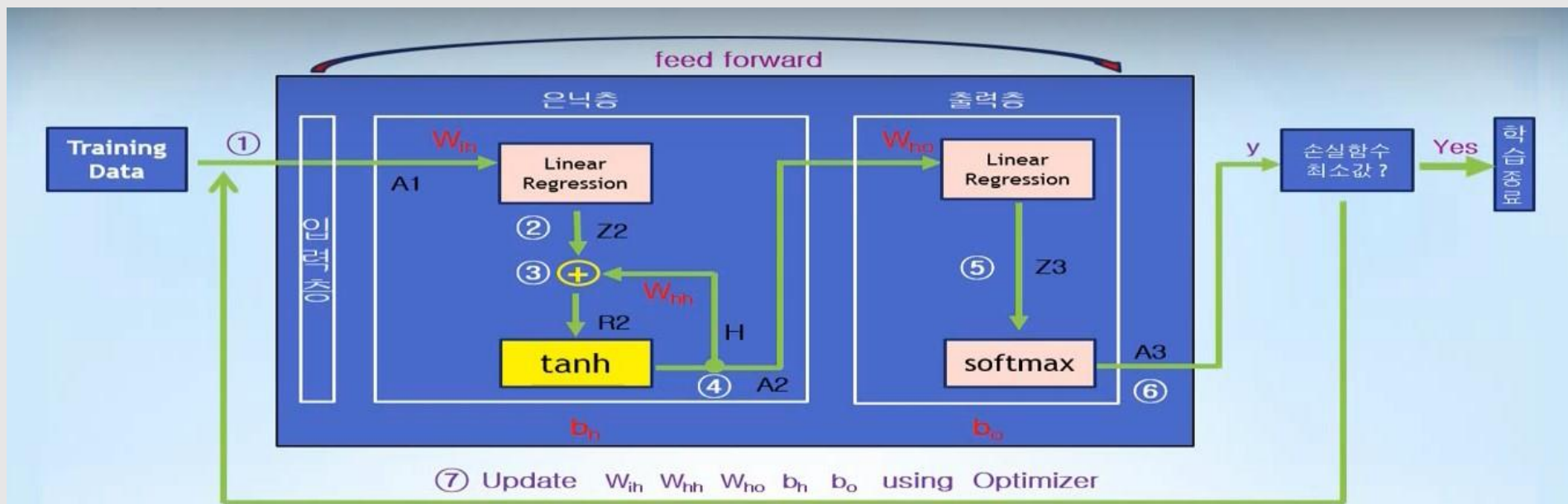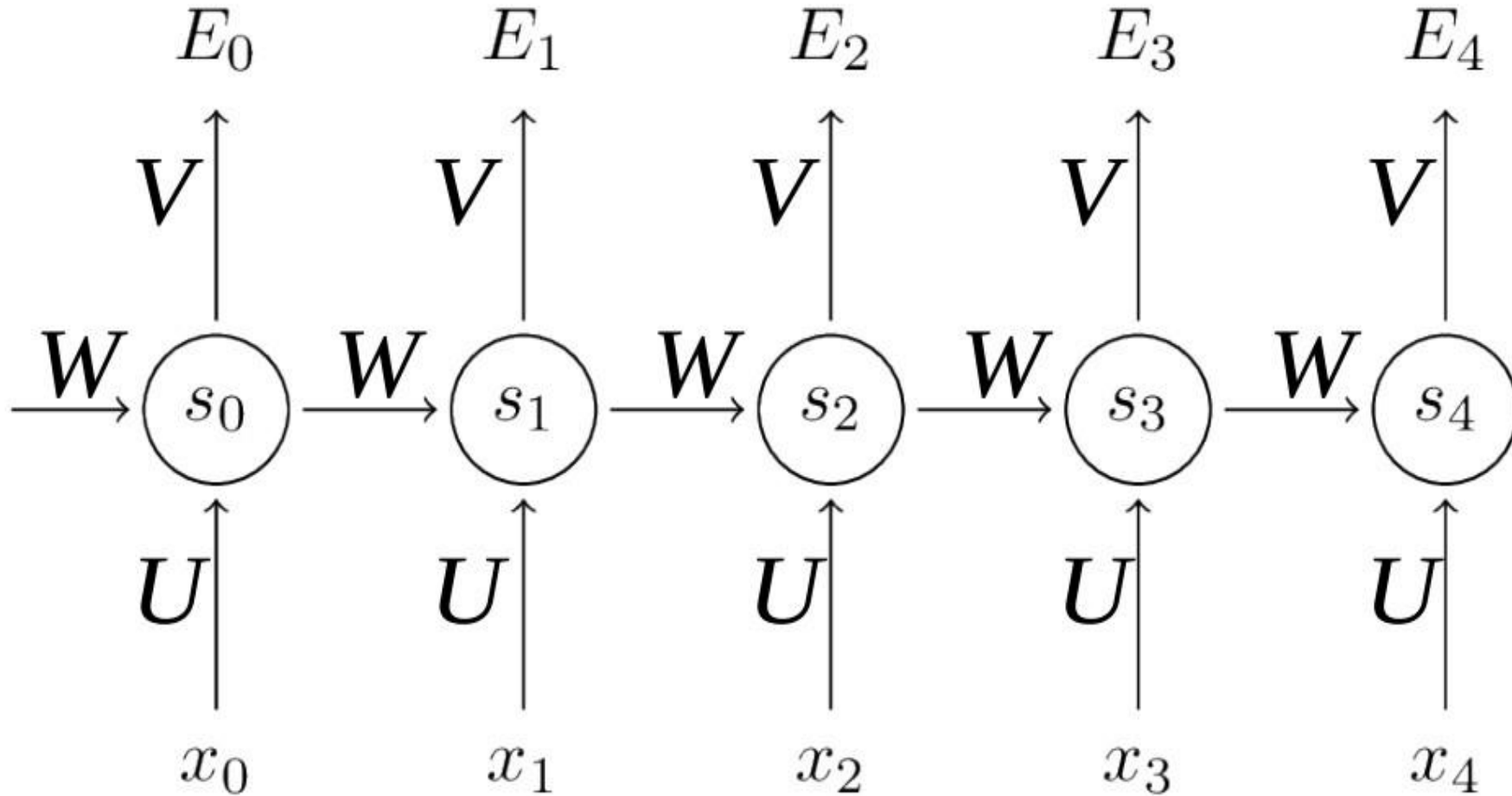시간개념을 포함한 current state $H_t$

현재 입력데이터 A1 에 적용되는 가중치     과거(이전) state에 적용되는 가중치

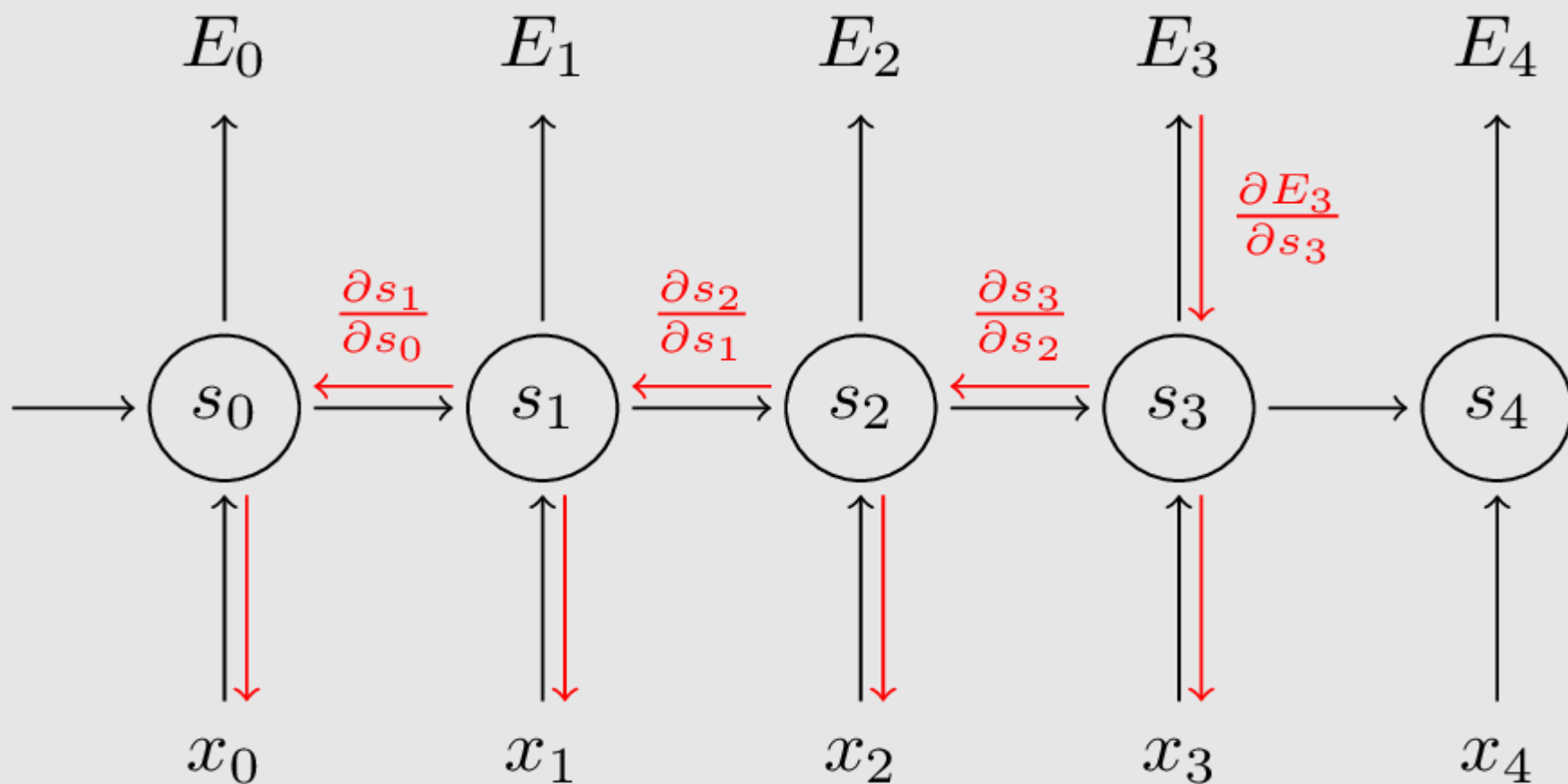$$H_t = A2 = \tanh(A1 \cdot W_{ih} + H_{t-1} \cdot W_{hh} + b_h)$$

현재 입력데이터 A1에 대한 state    현재 입력데이터 A1    과거(이전) 입력데이터 A1에 대한 state    은닉층 바이어스

E=Error

E=Y - ŷ

$$s_t = \tanh(Ux_t + Ws_{t-1})$$
$$\hat{y}_t = softmax(Vs_t)$$
$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3}\frac{\partial \bar{y}_3}{\partial \bar{s}_3}\frac{\partial \bar{s}_3}{\partial W_s}$$
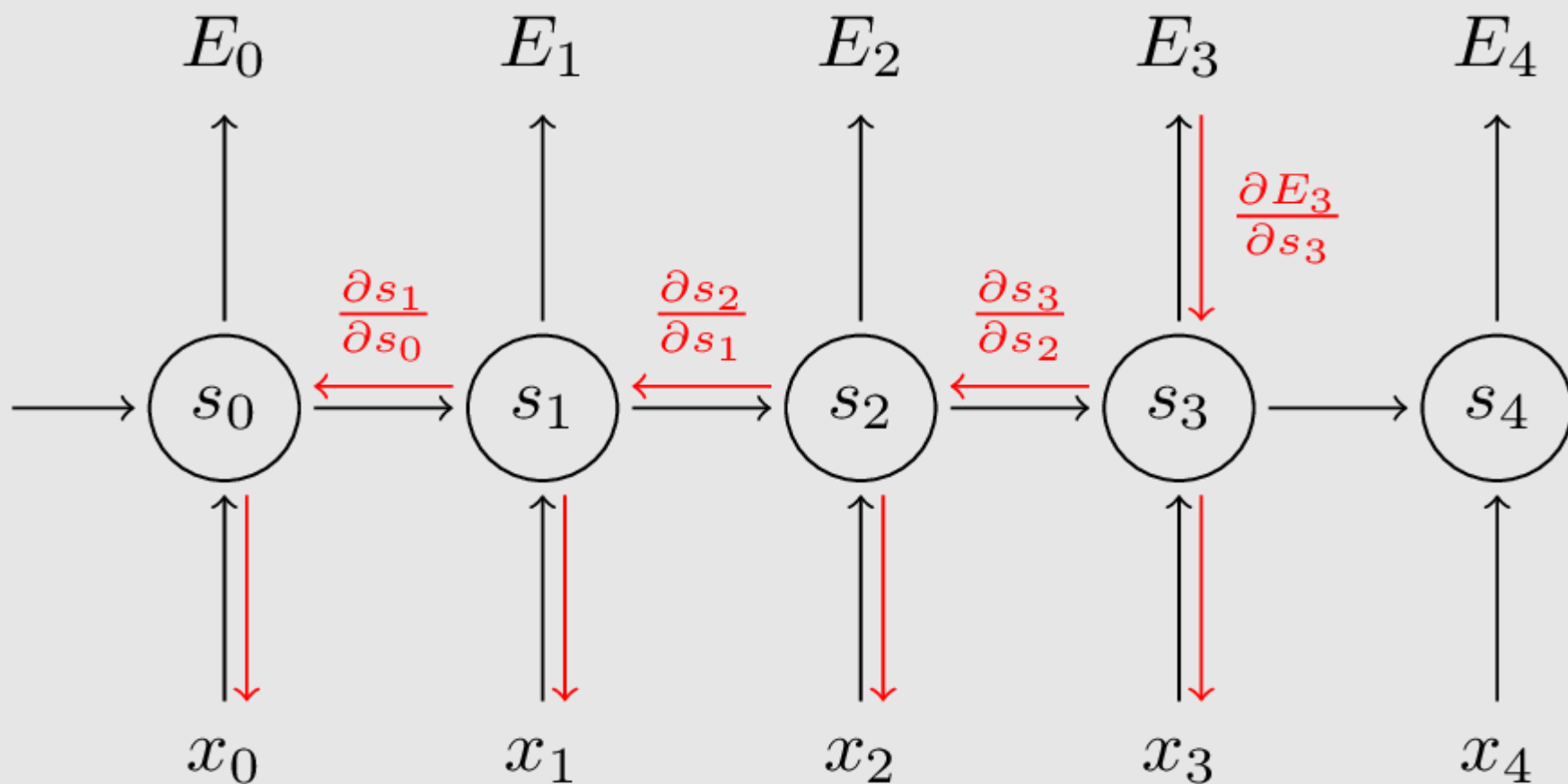
$$s_t = \tanh(Ux_t + Ws_{t-1})$$

$$\hat{y}_t = softmax(Vs_t)$$

$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3} \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \frac{\partial \bar{s}_3}{\partial W_s}$$

$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3} \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \frac{\partial \bar{s}_3}{\partial \bar{s}_2} \frac{\partial \bar{s}_2}{\partial W_s}$$
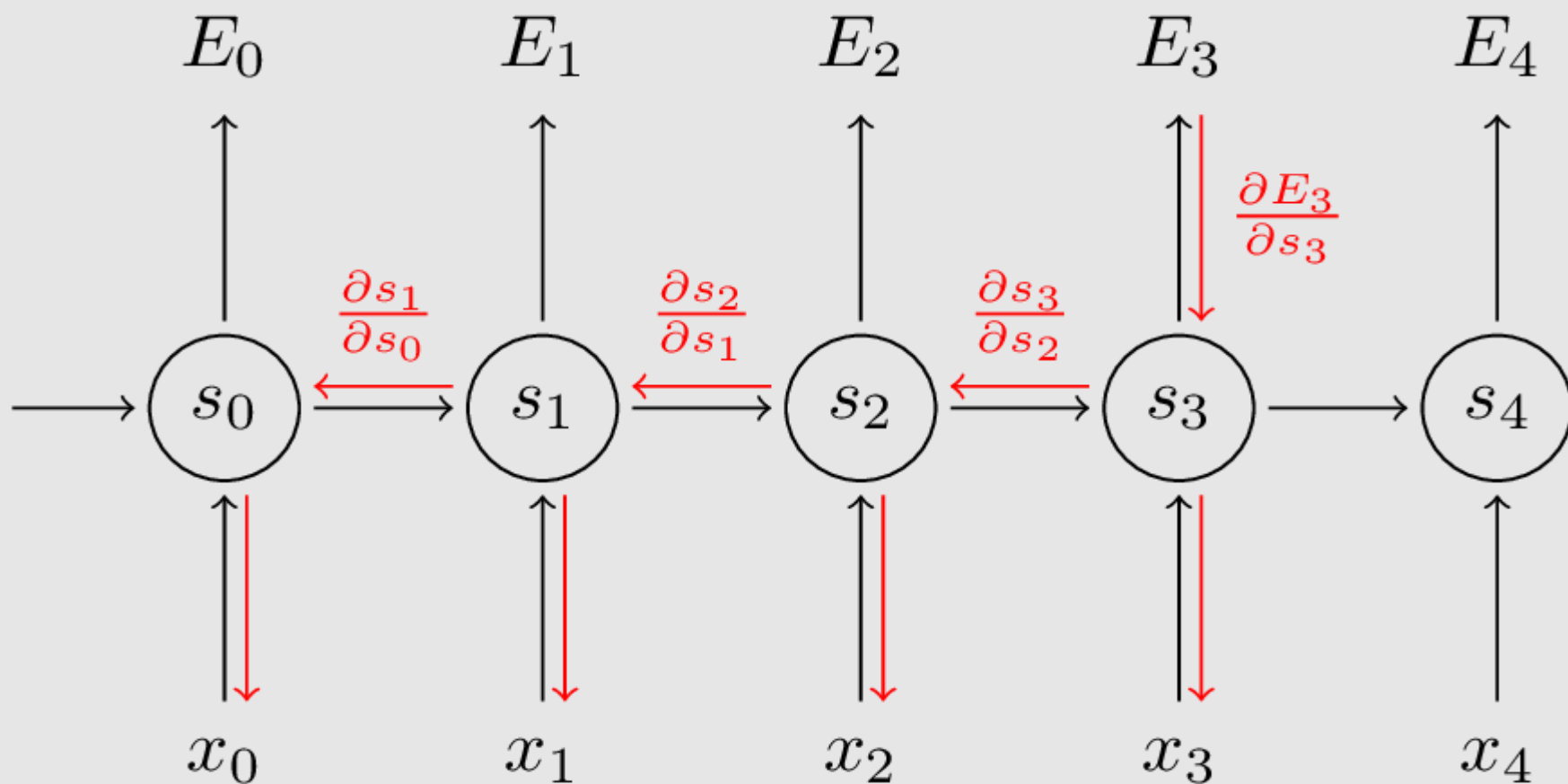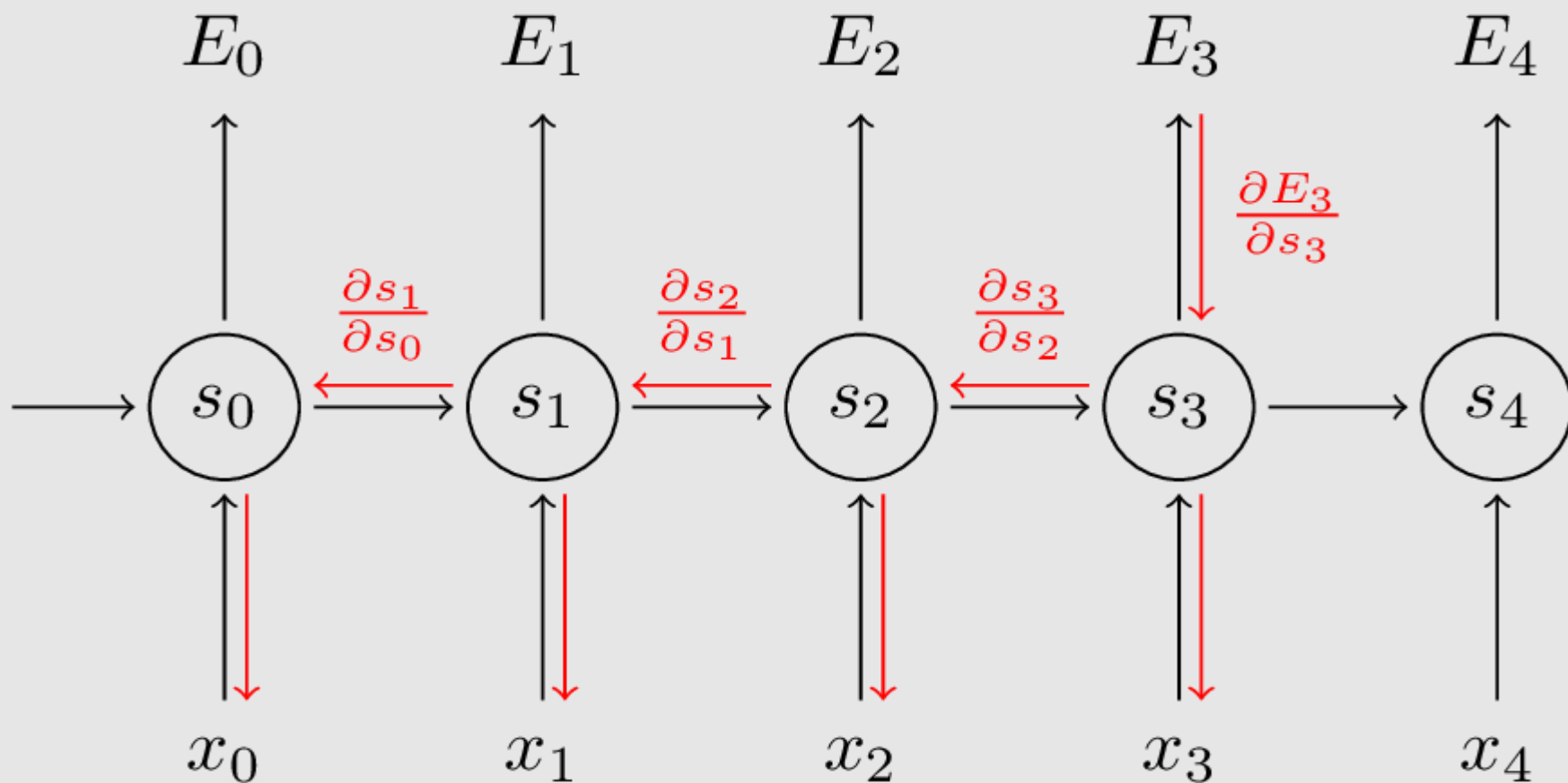
$$s_t = \tanh(Ux_t + Ws_{t-1})$$
$$\hat{y}_t = softmax(Vs_t)$$

$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3}\frac{\partial \bar{y}_3}{\partial \bar{s}_3}\frac{\partial \bar{s}_3}{\partial W_s}$$

$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3}\frac{\partial \bar{y}_3}{\partial \bar{s}_3}\frac{\partial \bar{s}_3}{\partial \bar{s}_2}\frac{\partial \bar{s}_2}{\partial W_s}$$

$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3}\frac{\partial \bar{y}_3}{\partial \bar{s}_3}\frac{\partial \bar{s}_3}{\partial \bar{s}_2}\frac{\partial \bar{s}_2}{\partial \bar{s}_1}\frac{\partial \bar{s}_1}{\partial W_s}$$

$$s_t = \tanh(Ux_t + Ws_{t-1})$$

$$\hat{y}_t = softmax(Vs_t)$$

$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3} \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \frac{\partial \bar{s}_3}{\partial W_s}$$

$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3} \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \frac{\partial \bar{s}_3}{\partial \bar{s}_2} \frac{\partial \bar{s}_2}{\partial W_s}$$
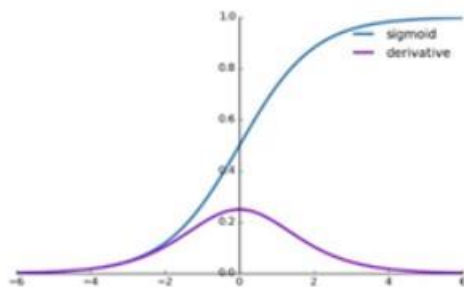
$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3} \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \frac{\partial \bar{s}_3}{\partial \bar{s}_2} \frac{\partial \bar{s}_2}{\partial \bar{s}_1} \frac{\partial \bar{s}_1}{\partial W_s}$$

$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3} \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \frac{\partial \bar{s}_3}{\partial W_s} +$$

$$\frac{\partial E_3}{\partial \bar{y}_3} \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \frac{\partial \bar{s}_3}{\partial \bar{s}_2} \frac{\partial \bar{s}_2}{\partial W_s} +$$

$$\frac{\partial E_3}{\partial \bar{y}_3} \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \frac{\partial \bar{s}_3}{\partial \bar{s}_2} \frac{\partial \bar{s}_2}{\partial \bar{s}_1} \frac{\partial \bar{s}_1}{\partial W_s}$$
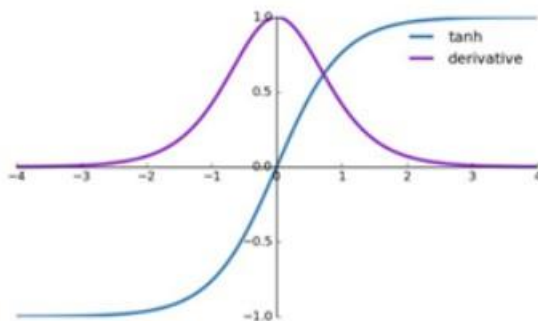
$$W = W - learning\_rate * \frac{\partial E}{\partial W}$$

## 왜 tanh를 쓰는가?



- **_Activation Function_**
  - Sigmoid
  - Tanh

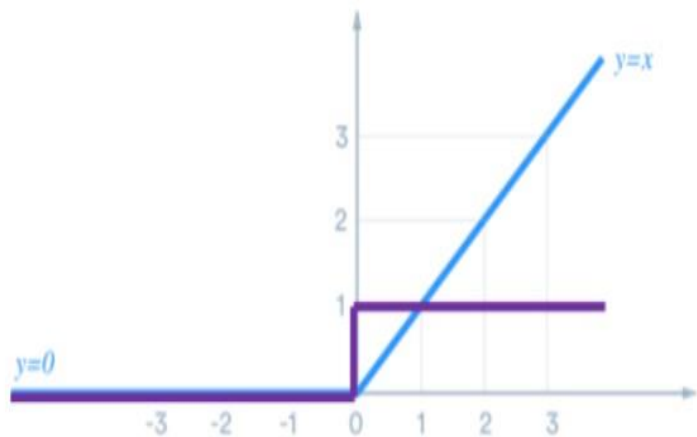| Sigmoid | |
| --- | --- |
| $f(x)$ | $\dfrac{1}{1+e^{-x}}$ ($y$: $0 \sim 1$) |
| $\dfrac{d}{dx}f(x)$ | $\dfrac{1}{1+e^{-x}}\left(1 - \dfrac{1}{1+e^{-x}}\right)$ ($y'$: $0 \sim 0.25$) |

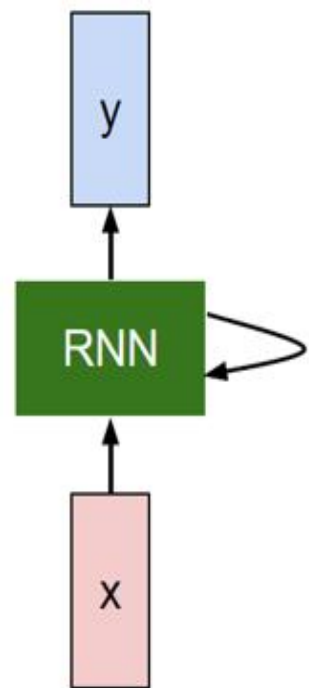| Tanh | |
| --- | --- |
| $f(x)$ | $\dfrac{e^x - e^{-x}}{e^x + e^{-x}}$ ($y$: $-1 \sim 1$) |
| $\dfrac{d}{dx}f(x)$ | $1 - \tanh(x)^2$ ($y'$: $0 \sim 1$) |

왜 tanh를 쓰는가?

Sequence is important for POS tagging

one to many      many to one      many to many

<기본구조>

<사진설명 붙이기>
사진 → 단어들
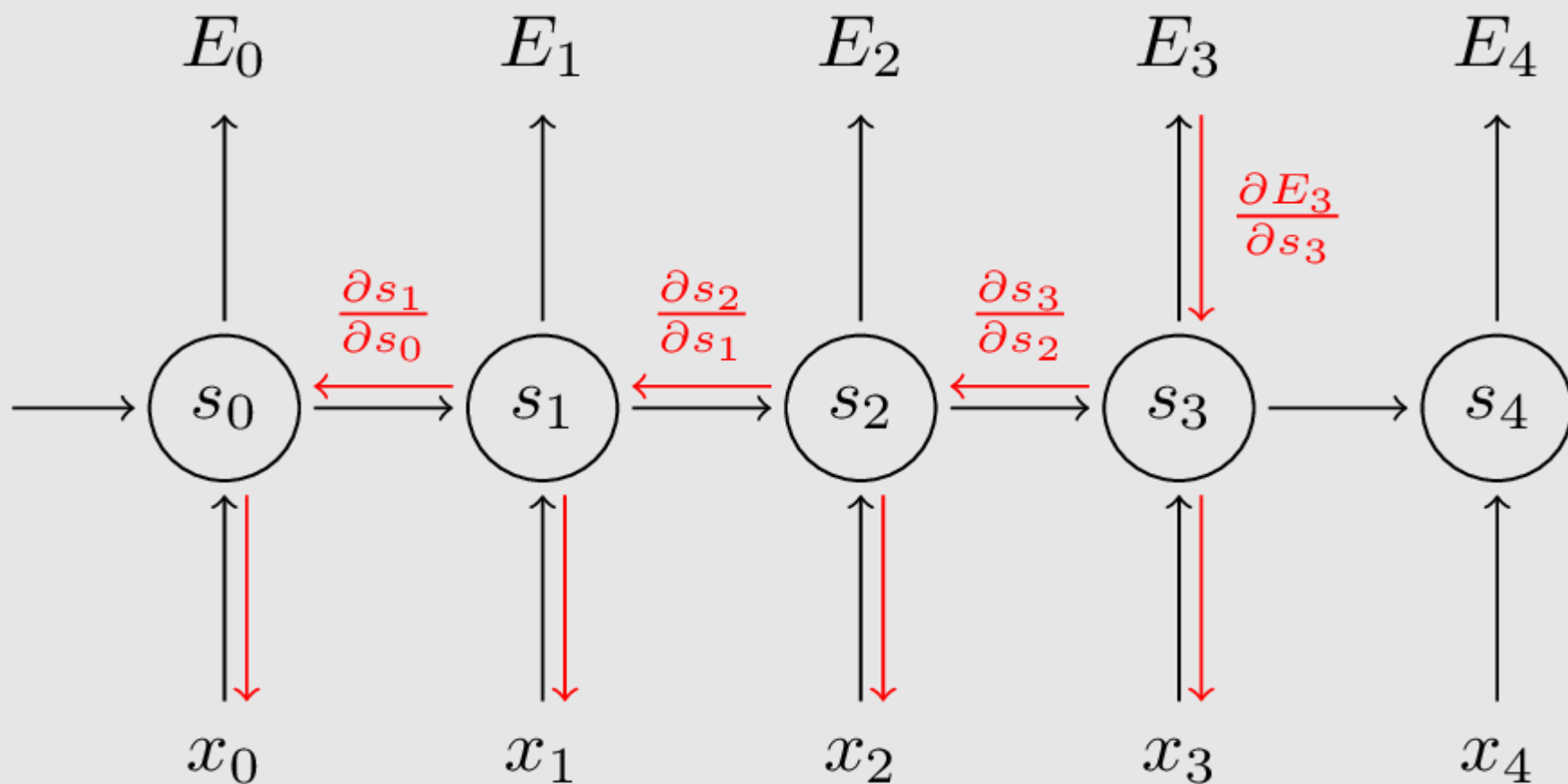
<감성분석>
단어들 → 감성점수

<번역>
단어들 → 단어들

$$s_t = \tanh(Ux_t + Ws_{t-1})$$
$$\hat{y}_t = softmax(Vs_t)$$

$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3} \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \frac{\partial \bar{s}_3}{\partial W_s}$$

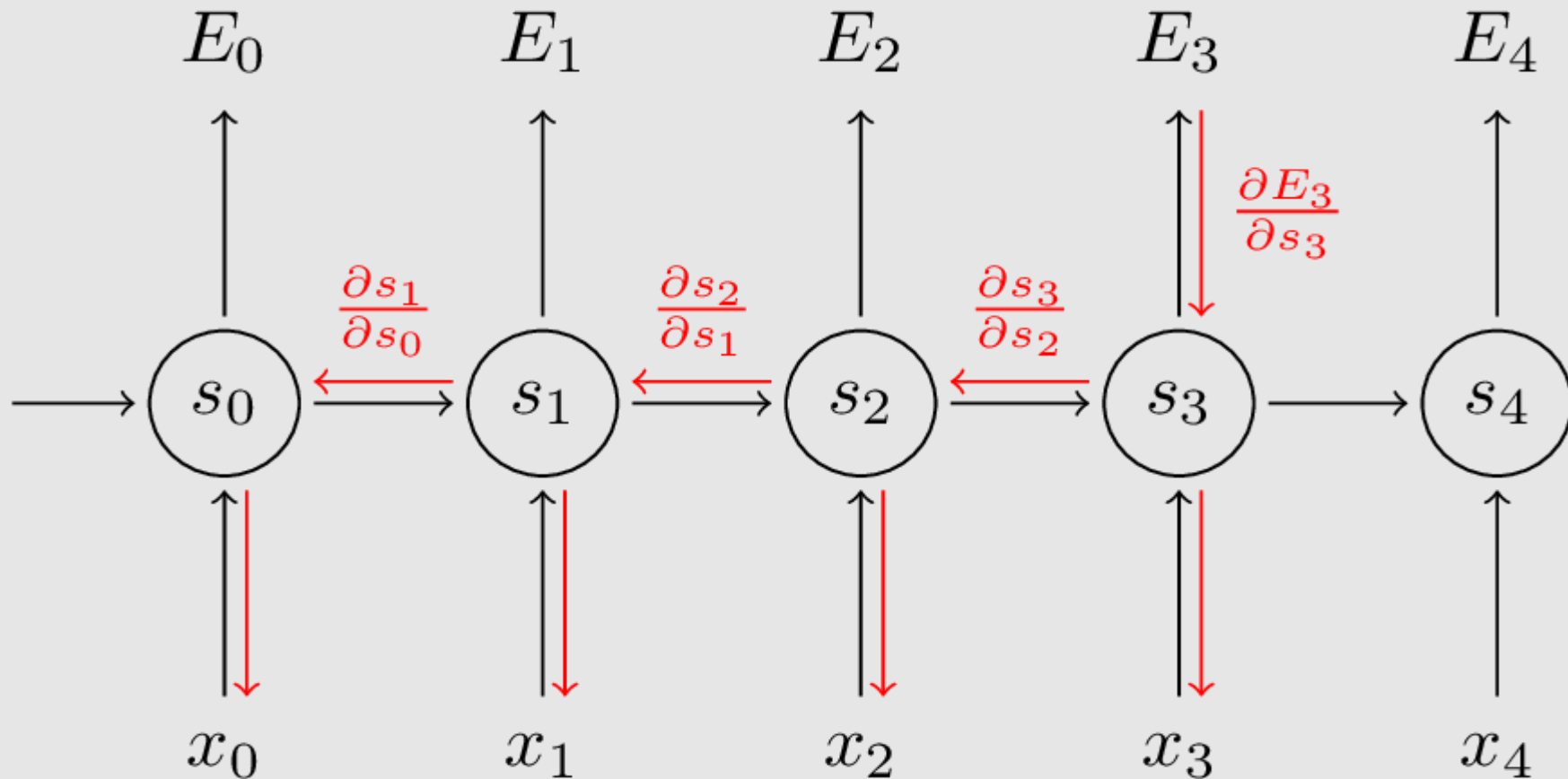$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3} \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \frac{\partial \bar{s}_3}{\partial \bar{s}_2} \frac{\partial \bar{s}_2}{\partial W_s}$$
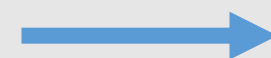
$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3} \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \frac{\partial \bar{s}_3}{\partial \bar{s}_2} \frac{\partial \bar{s}_2}{\partial \bar{s}_1} \frac{\partial \bar{s}_1}{\partial W_s}$$

$$W = W - learning\_rate * \frac{\partial E}{\partial W}$$

1보다 미분값이 클경우
Gradient Exploding

$E_0$     $E_1$     $E_2$     $E_3$     $E_4$

$\frac{\partial E_3}{\partial s_3}$

$\frac{\partial s_1}{\partial s_0}$    $\frac{\partial s_2}{\partial s_1}$    $\frac{\partial s_3}{\partial s_2}$

$s_0$    $s_1$    $s_2$    $s_3$    $s_4$

$x_0$    $x_1$    $x_2$    $x_3$    $x_4$

$$s_t = \tanh(Ux_t + Ws_{t-1})$$
$$\hat{y}_t = softmax(Vs_t)$$

$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3}\frac{\partial \bar{y}_3}{\partial \bar{s}_3}\frac{\partial \bar{s}_3}{\partial W_s}$$

$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3}\frac{\partial \bar{y}_3}{\partial \bar{s}_3}\frac{\partial \bar{s}_3}{\partial \bar{s}_2}\frac{\partial \bar{s}_2}{\partial W_s}$$

$$\frac{\partial E_3}{\partial W_s} = \frac{\partial E_3}{\partial \bar{y}_3}\frac{\partial \bar{y}_3}{\partial \bar{s}_3}\frac{\partial \bar{s}_3}{\partial \bar{s}_2}\frac{\partial \bar{s}_2}{\partial \bar{s}_1}\frac{\partial \bar{s}_1}{\partial W_s}$$

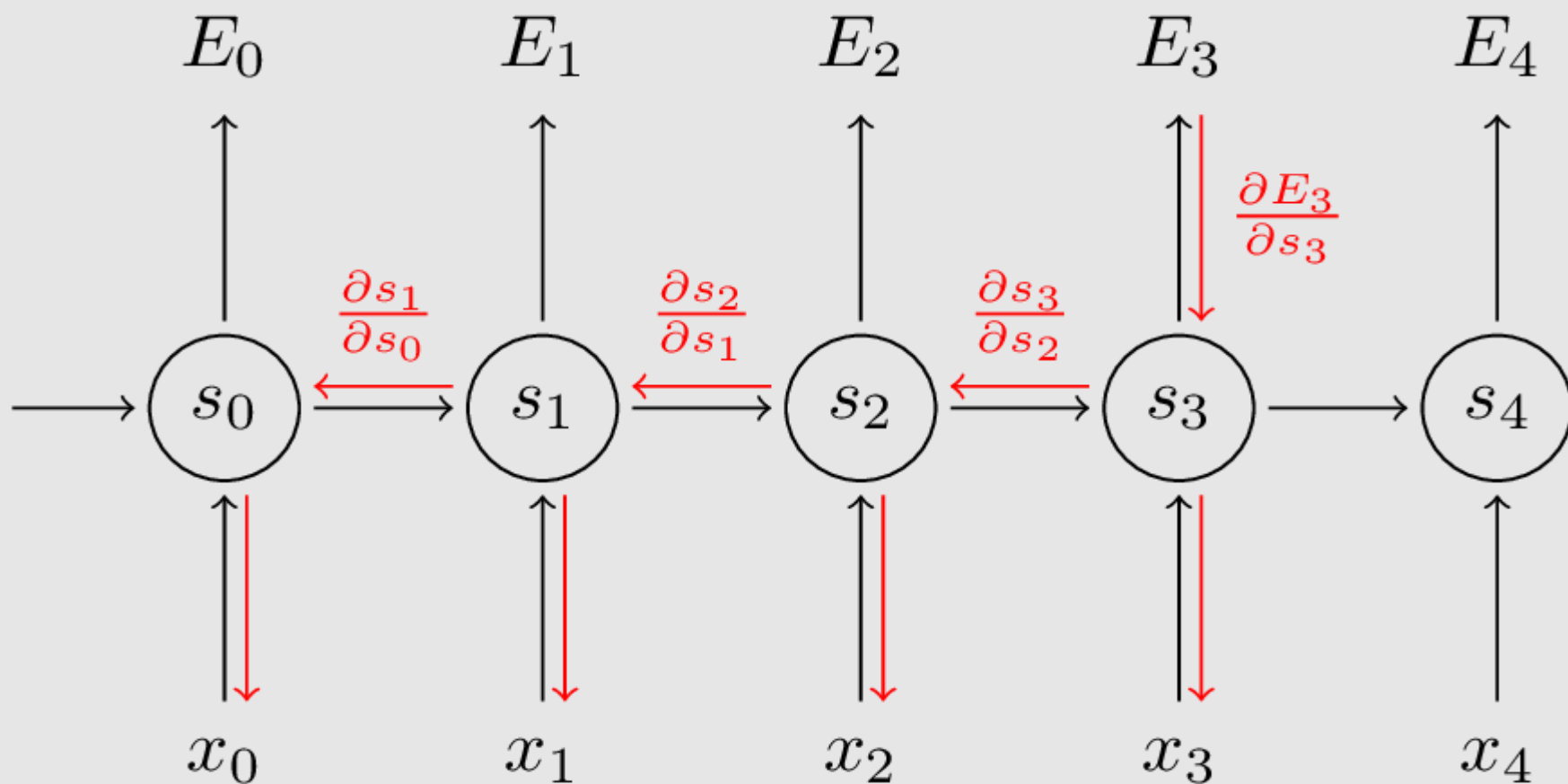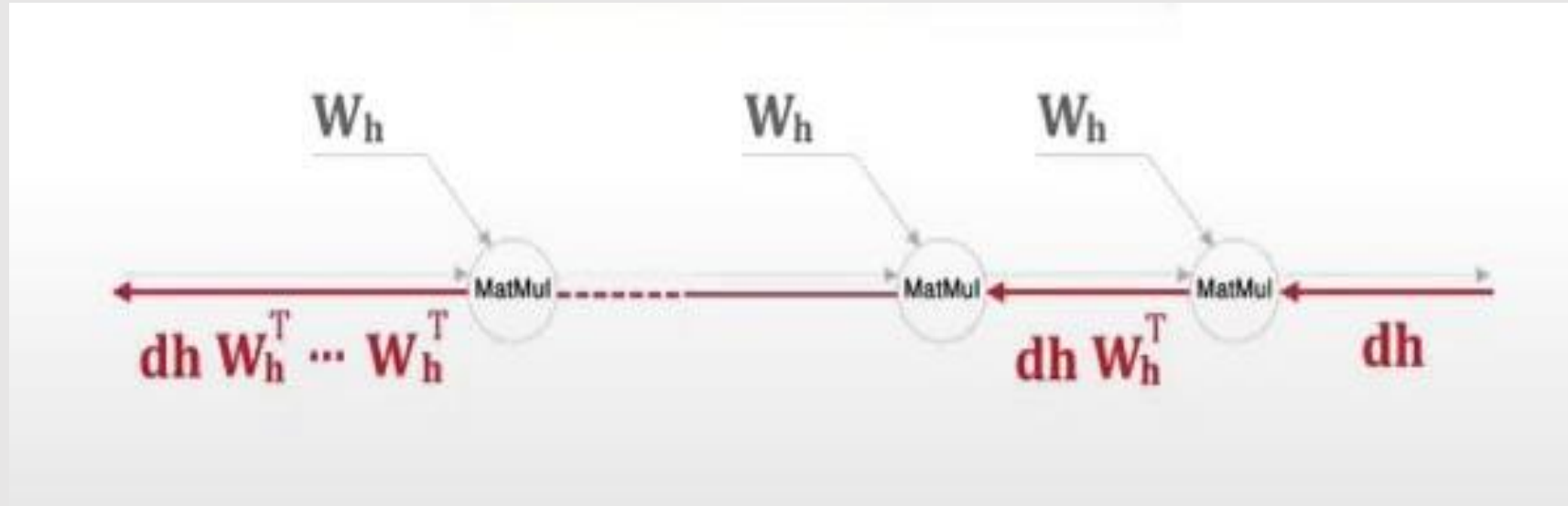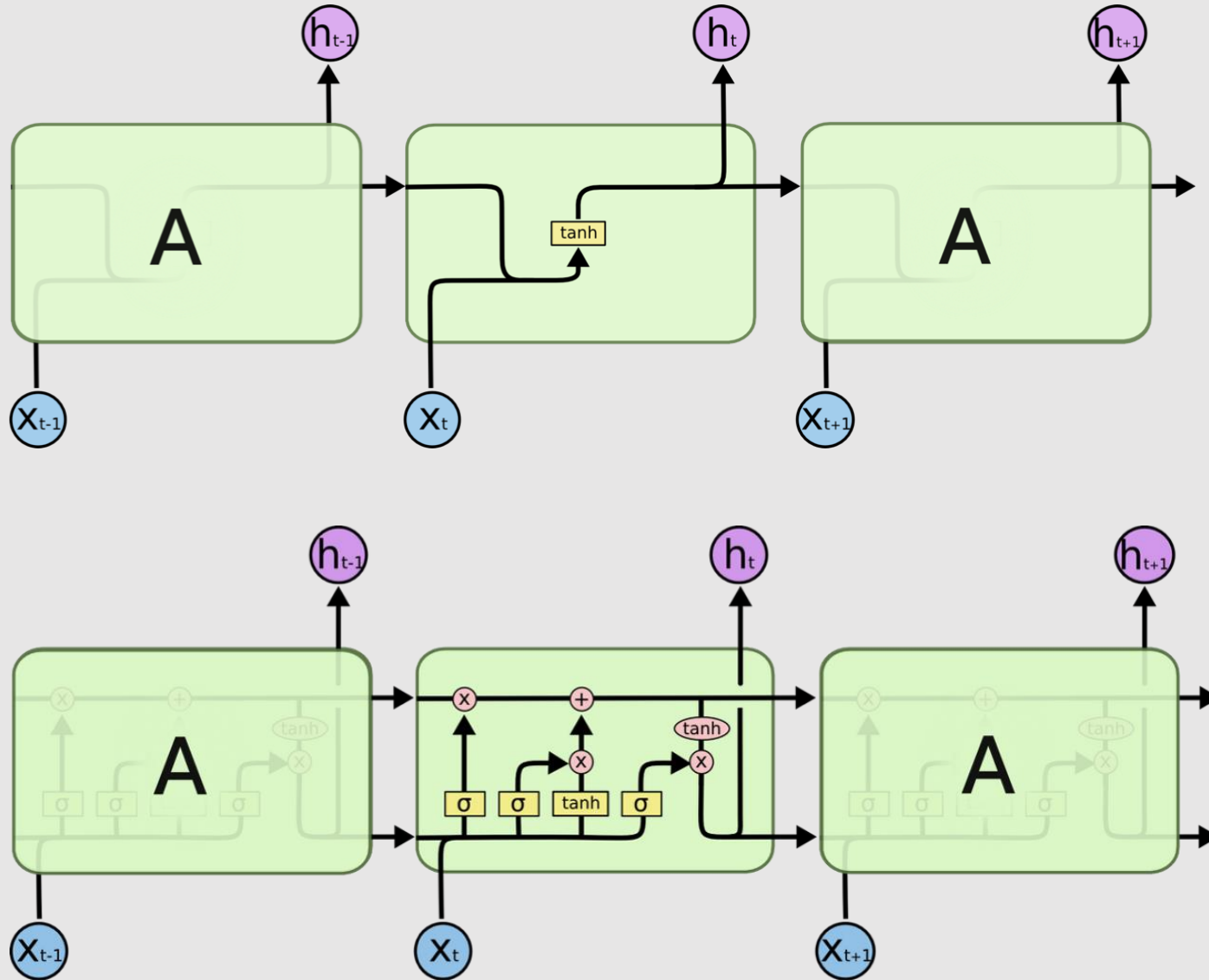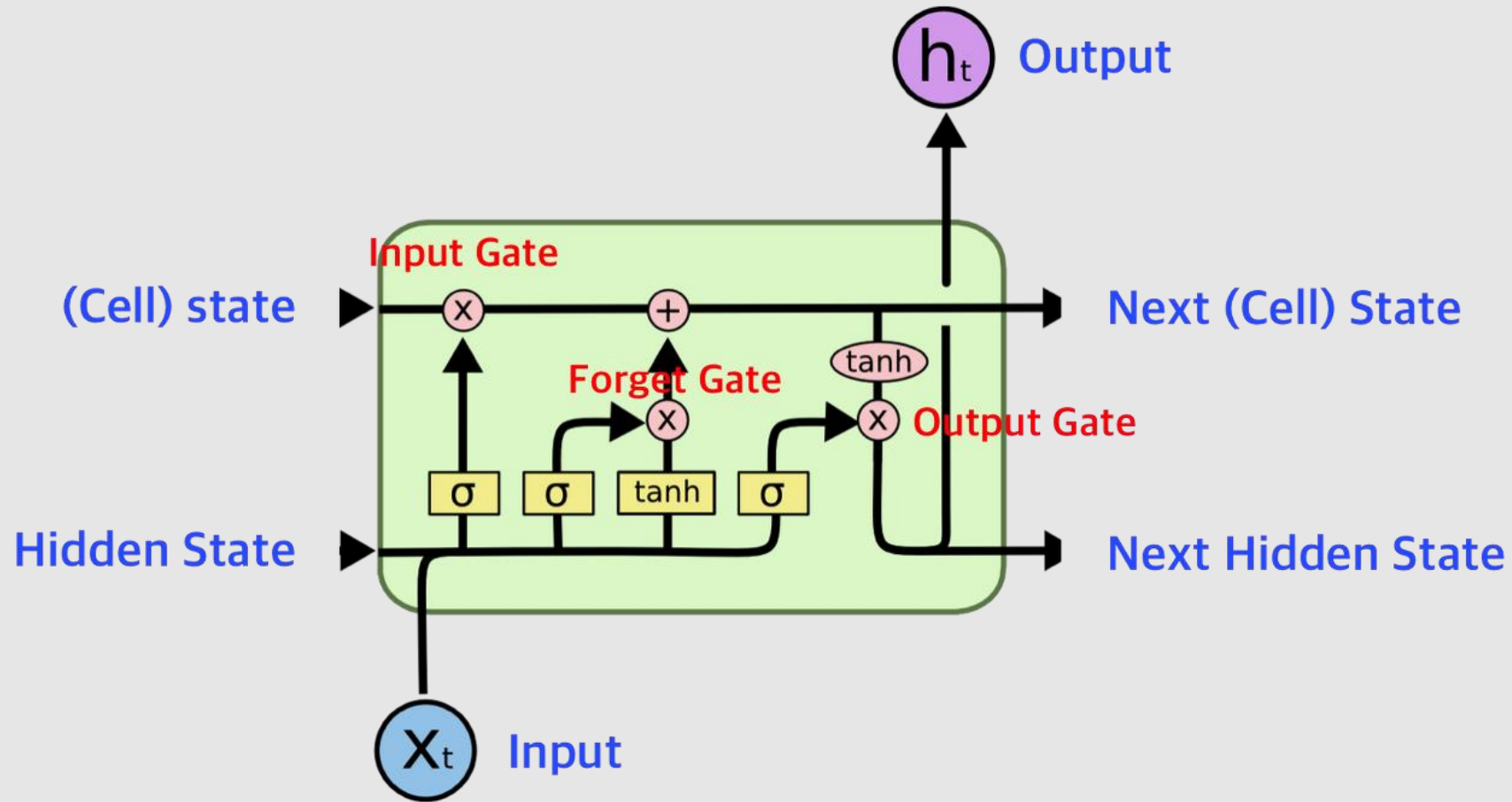$$W = W - learning\_rate * \frac{\partial E}{\partial W}$$

1보다 미분값이 클경우
Gradient Exploding     →     Gradient cliffing

1보다 미분값이 작을경우
Gradient Vanishing

## Forget Gate



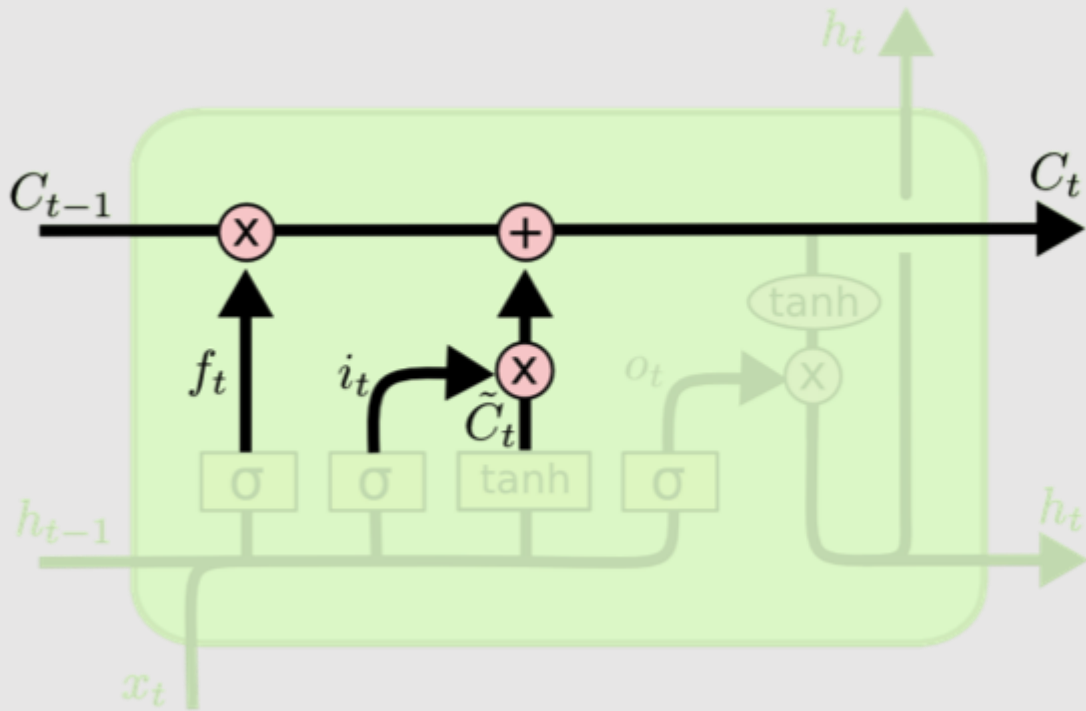$$f_t = \sigma\left(W_f \cdot [h_{t-1}, x_t] + b_f\right)$$

# Input Gate



$$i_t = \sigma\left(W_i \cdot [h_{t-1}, x_t] + b_i\right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

# Cell update



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

## Output Gate



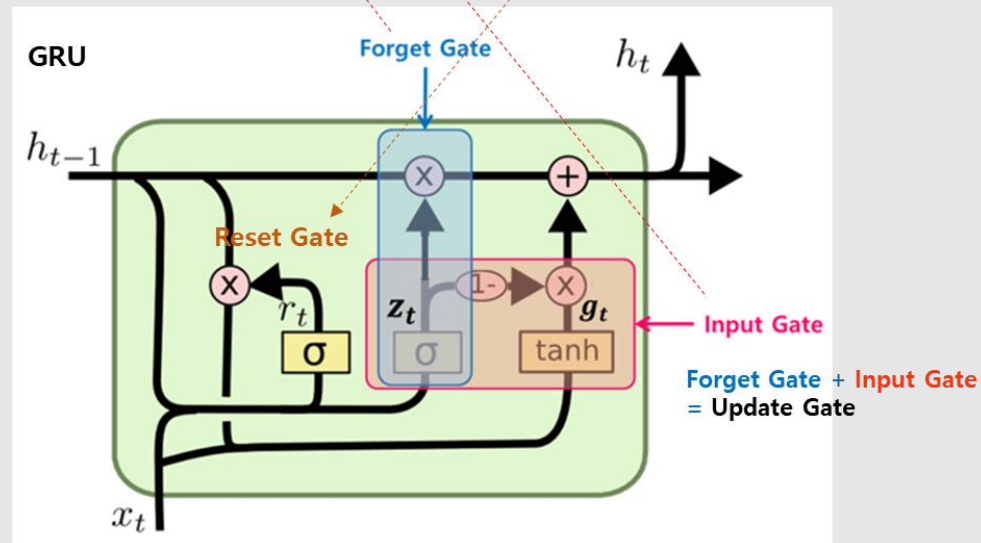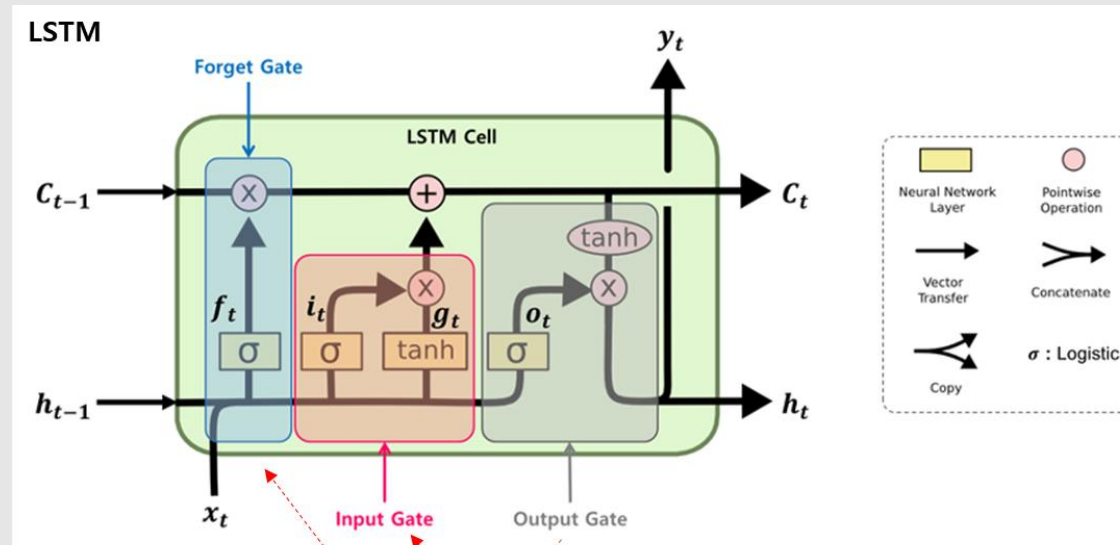$$o_t = \sigma\left(W_o\left[h_{t-1}, x_t\right] + b_o\right)$$

$$h_t = o_t * \tanh\left(C_t\right)$$

# 02.GRU Model

| Gate | Equation |
|------|----------|
| Reset Gate | $r_t = \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r)$ |
| Forget Gate | $z_t = \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z)$ |
| Input Gate | $1 - z_t$ |
| Hidden State | $g_t = tanh(W_{xg}x_t + W_{hg}(r_t \odot h_{t-1}) + b_g)$ |
| | $h_t = z_t \odot h_{t-1} + (1 - z_t) \odot g_t$ |

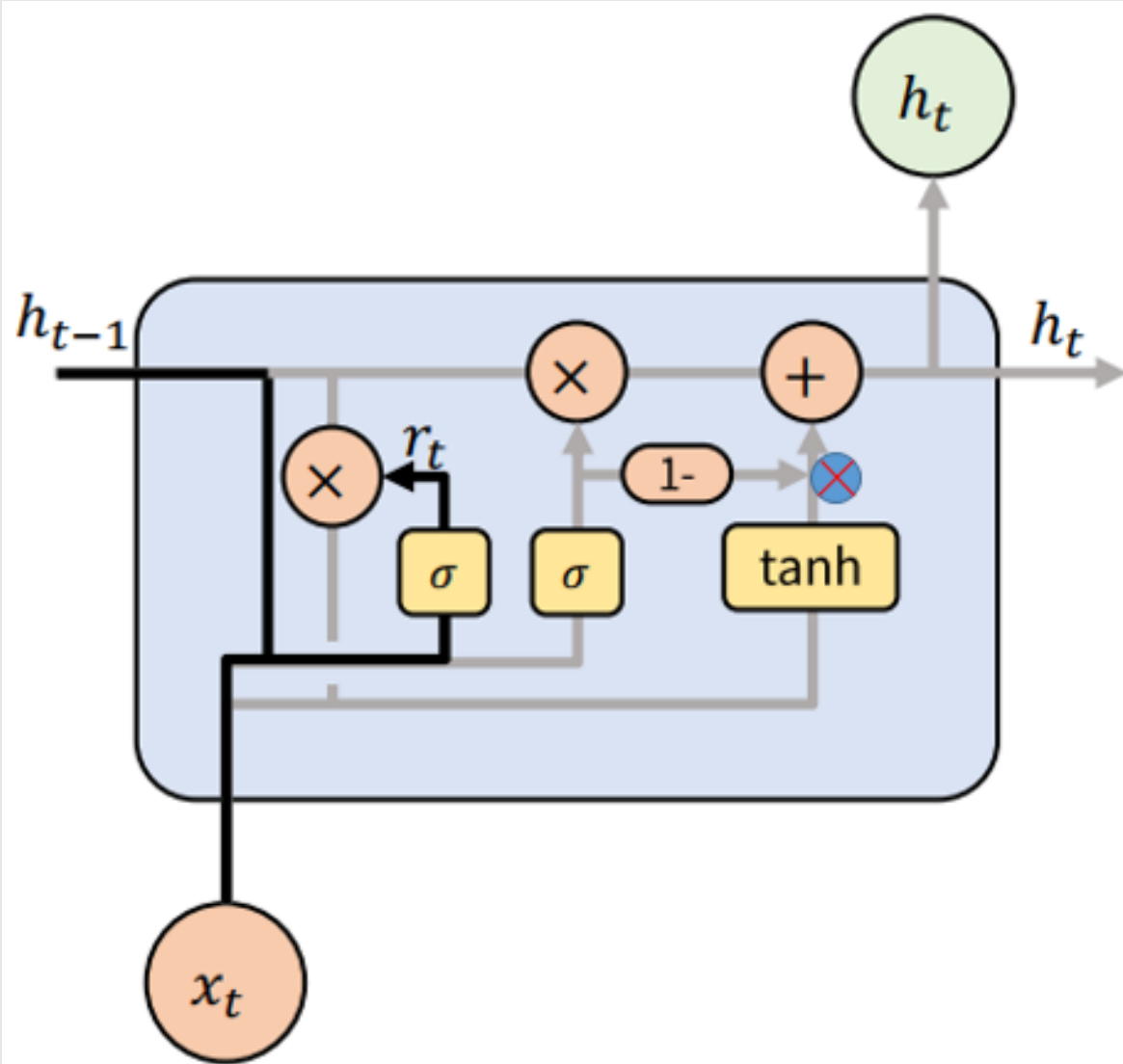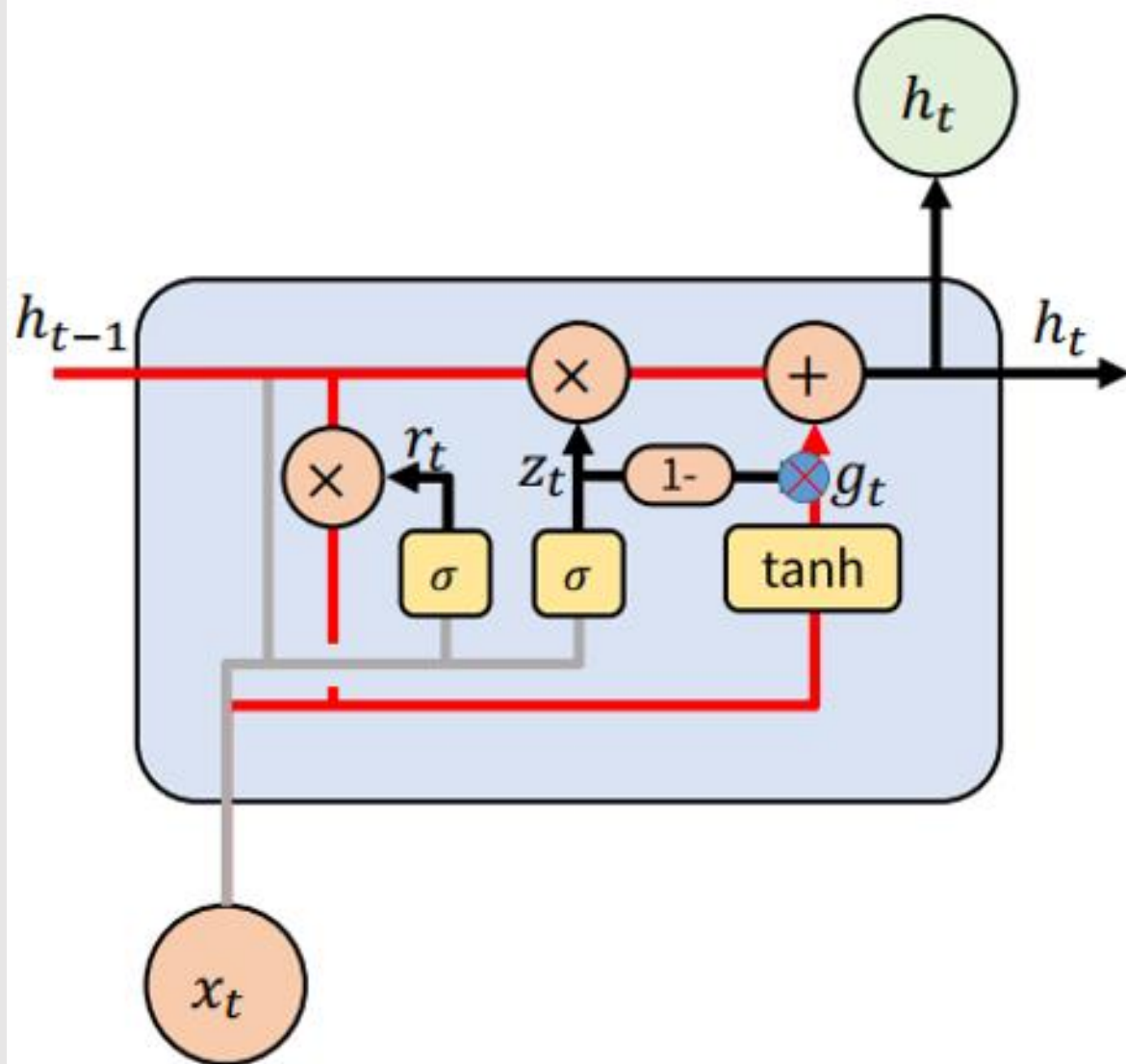| Gate | Equation |
|------|----------|
| Reset Gate | $r_t = \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r)$ |
| Forget Gate | $z_t = \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z)$ |
| Input Gate | $1 - z_t$ |
| Hidden State | $g_t = tanh(W_{xg}x_t + W_{hg}(r_t \odot h_{t-1}) + b_g)$ |
| | $h_t = z_t \odot h_{t-1} + (1 - z_t) \odot g_t$ |

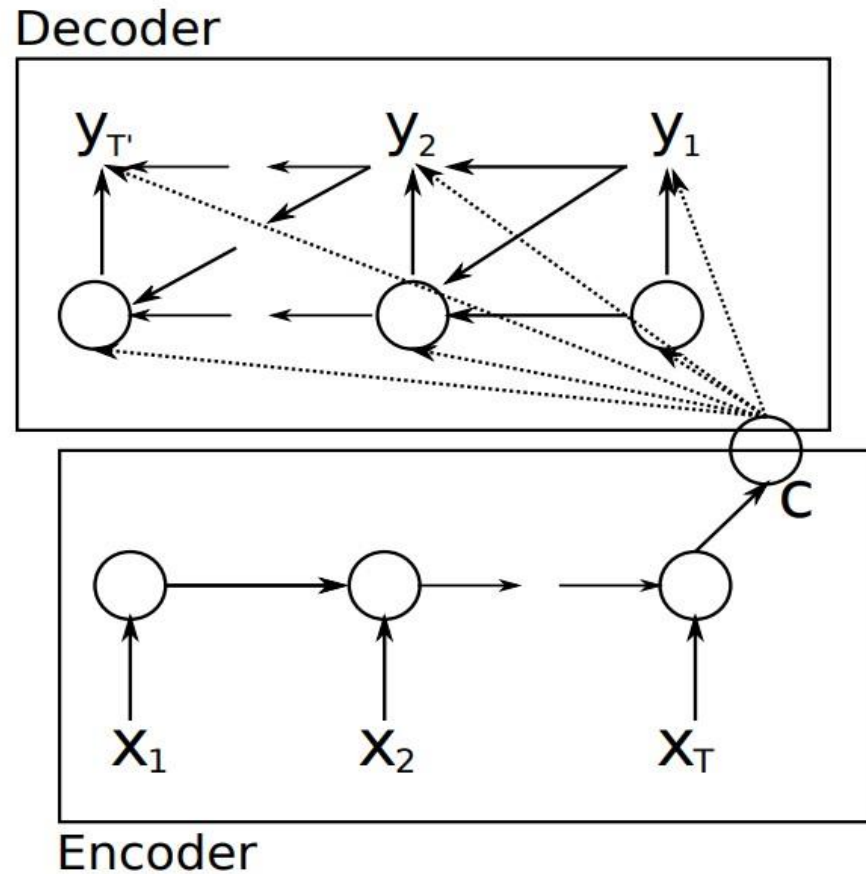| Gate | Equation |
|------|----------|
| Reset Gate | $r_t = \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r)$ |
| Forget Gate | $z_t = \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z)$ |
| Input Gate | $1 - z_t$ |
| Hidden State | $g_t = tanh(W_{xg}x_t + W_{hg}(r_t \odot h_{t-1}) + b_g)$ |
| | $h_t = z_t \odot h_{t-1} + (1 - z_t) \odot g_t$ |

Figure 1: An illustration of the proposed RNN Encoder–Decoder.

$$p(y_1, \ldots, y_{T'} \mid x_1, \ldots, x_T),$$

**Encoder** $\longrightarrow$

$$\mathbf{h}_{\langle t \rangle} = f\left(\mathbf{h}_{\langle t-1 \rangle}, x_t\right),$$

Figure 1: An illustration of the proposed RNN Encoder–Decoder.

**Encoder** $\longrightarrow$

**Decoder** $\longrightarrow$

$$p(y_1, \ldots, y_{T'} \mid x_1, \ldots, x_T),$$

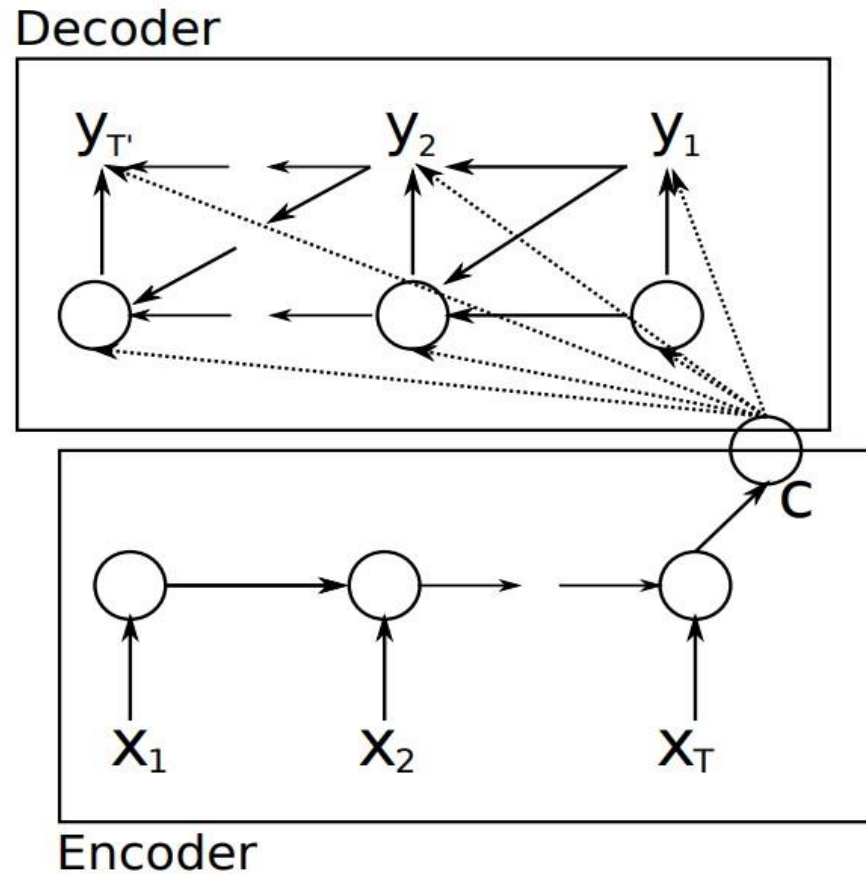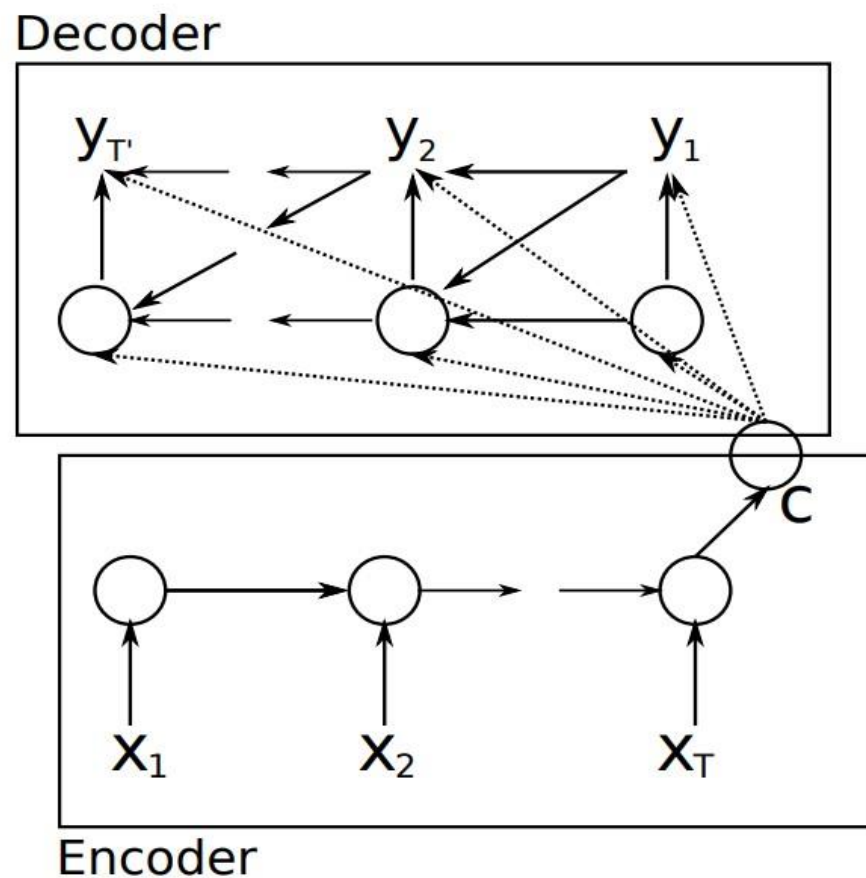$$\mathbf{h}_{\langle t \rangle} = f\left(\mathbf{h}_{\langle t-1 \rangle}, x_t\right),$$

$$\mathbf{h}_{\langle t \rangle} = f\left(\mathbf{h}_{\langle t-1 \rangle}, y_{t-1}, \mathbf{c}\right),$$

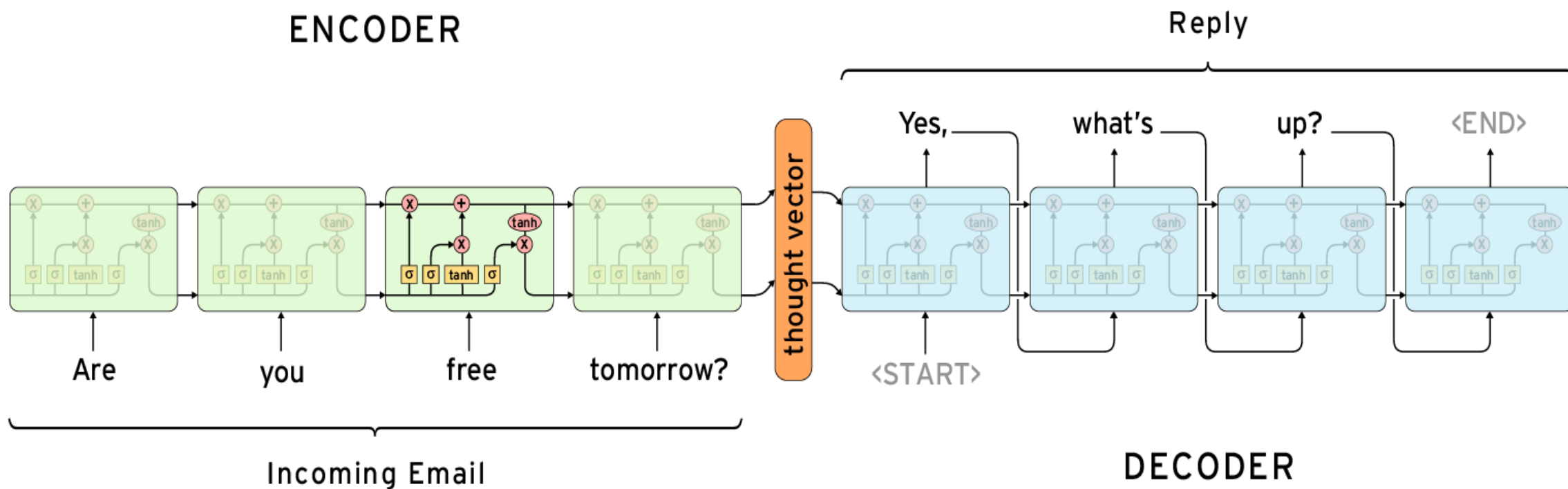Figure 1: An illustration of the proposed RNN Encoder–Decoder.

$$P(y_t|y_{t-1}, y_{t-2}, \ldots, y_1, \mathbf{c}) = g\left(\mathbf{h}_{\langle t \rangle}, y_{t-1}, \mathbf{c}\right).$$

$$\max_{\boldsymbol{\theta}} \frac{1}{N} \sum_{n=1}^{N} \log p_{\boldsymbol{\theta}}(\mathbf{y}_n \mid \mathbf{x}_n),$$

학습목표 ⟶

$$\max_{\boldsymbol{\theta}} \frac{1}{N} \sum_{n=1}^{N} \log p_{\boldsymbol{\theta}}(\mathbf{y}_n \mid \mathbf{x}_n),$$

학습목표 ⟶ $\max_{\theta} \frac{1}{4}\left(\log P_{(Are)} + \log P_{(you)} + \log P_{(free)} + \log P_{(tomorrow)}\right)$

$$p(e|f)$$

F가 주어졌을때 e가 나올 확률

$$p(e|f) \propto p(f|e)p(e)$$

비례관계

$$\tilde{e} = arg\max_{e \in e^*} p(e|f) = arg\max_{e \in e^*} p(f|e)p(e)$$
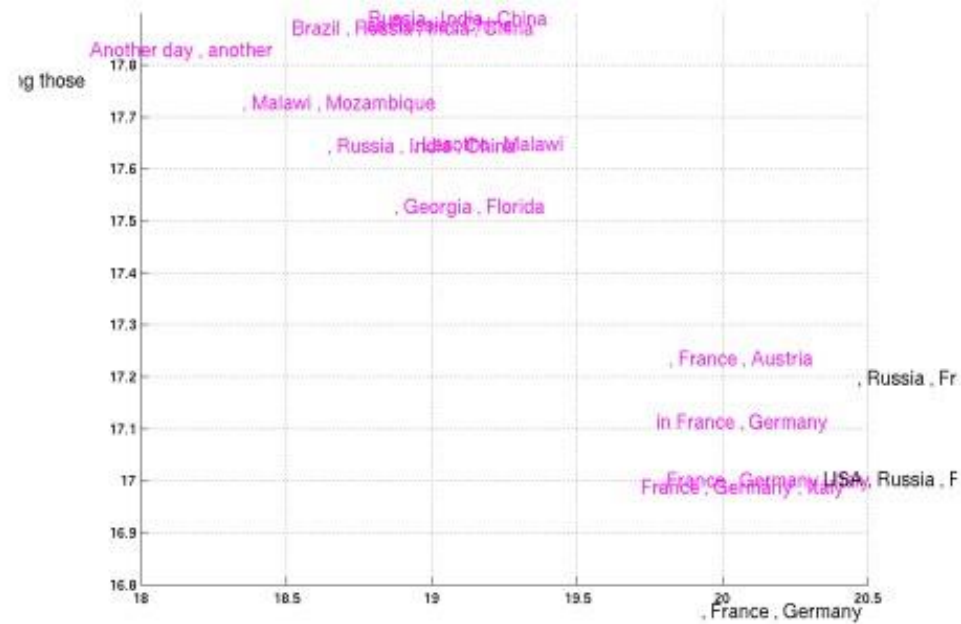
가장 높은 확률로 나오는 표현을 골라 번역 e를 찾는다

$$\log p(\mathbf{f} \mid \mathbf{e}) = \sum_{n=1}^{N} w_n f_n(\mathbf{f}, \mathbf{e}) + \log Z(\mathbf{e}), \quad (9)$$

# 03. Result & Conclusion

| Models | BLEU | |
|---|---|---|
| | dev | test |
| Baseline | 30.64 | 33.30 |
| RNN | 31.20 | 33.87 |
| CSLM + RNN | 31.48 | 34.64 |
| CSLM + RNN + WP | 31.50 | 34.54 |

# THANK YOU!