# CS224n LECTURE 7

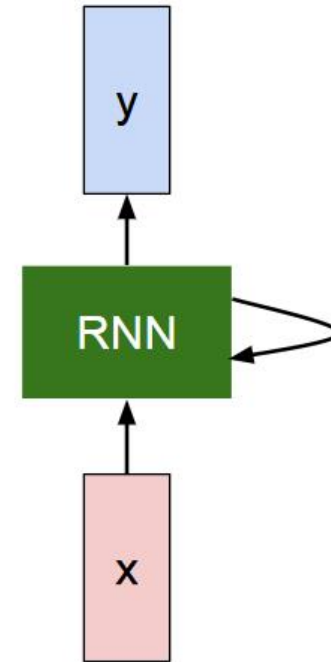## RNN / LSTM / GRU

**UOS STAT NLP Seminar**
**Kang Changgu**

# Vanilla RNN

# RNN (= Recurrent Neural Network)

We can process a sequence of vectors **x** by
applying a **recurrence formula** at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

new state        old state   input vector at
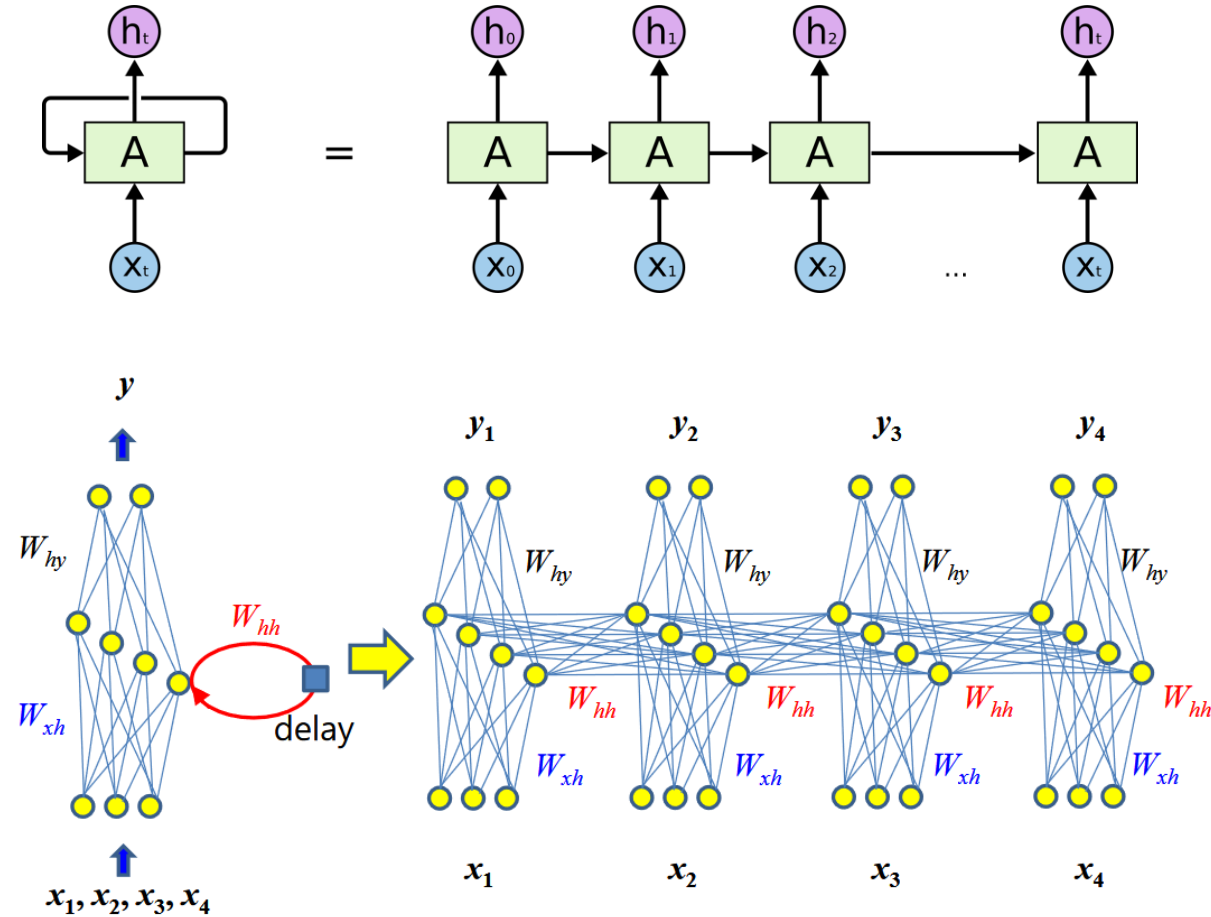some time step

some function
with parameters W

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$
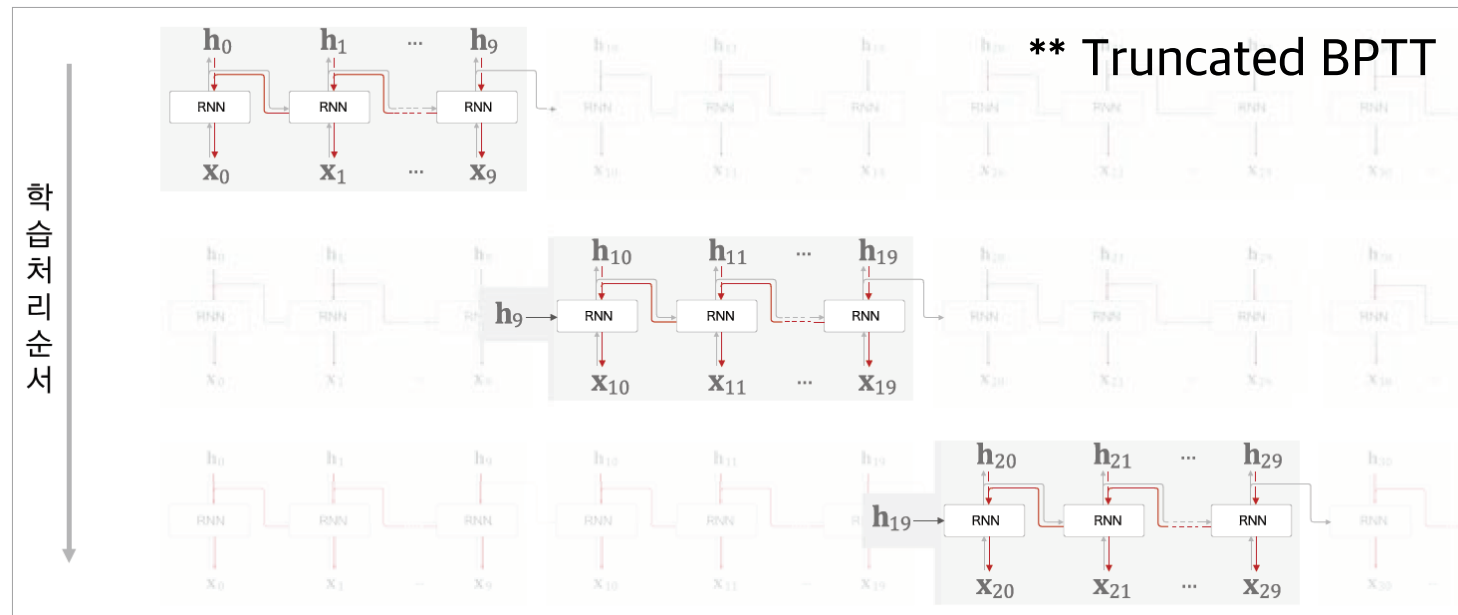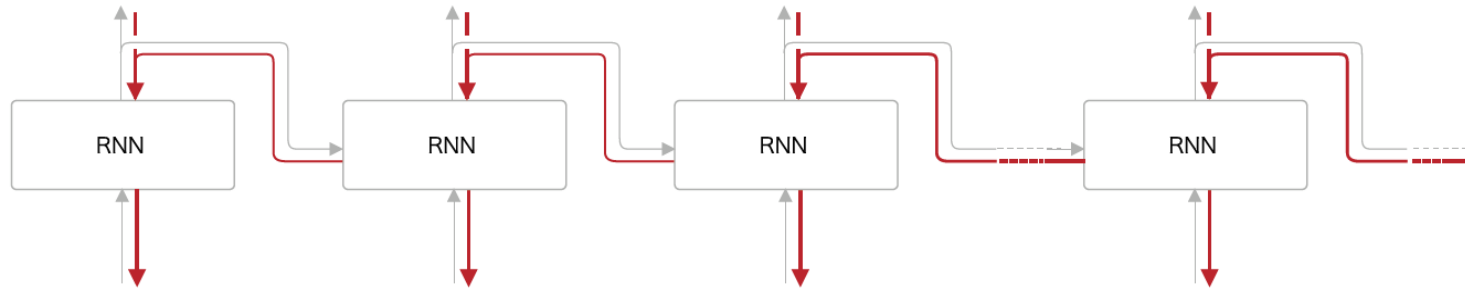
$$y_t = W_{hy}h_t$$

# Unfolding

RNN을 시간 축을 따라 unfolding 했을 때, 입력/은닉 상태를 시간 축을 따라 표현 가능

# BPTT (= Backpropagation Through Time)

RNN을 시간 축을 따라 unfolding 했을 때, backpropagation 시 시간 역순으로 gradient 전파 표현 가능

# End-to-end RNN

$$\hat{y}^{(4)} = P(x^{(5)}|\text{the students opened their})$$

output distribution

$$\hat{y}^{(t)} = \text{softmax}\left(Uh^{(t)} + b_2\right) \in \mathbb{R}^{|V|}$$

hidden states

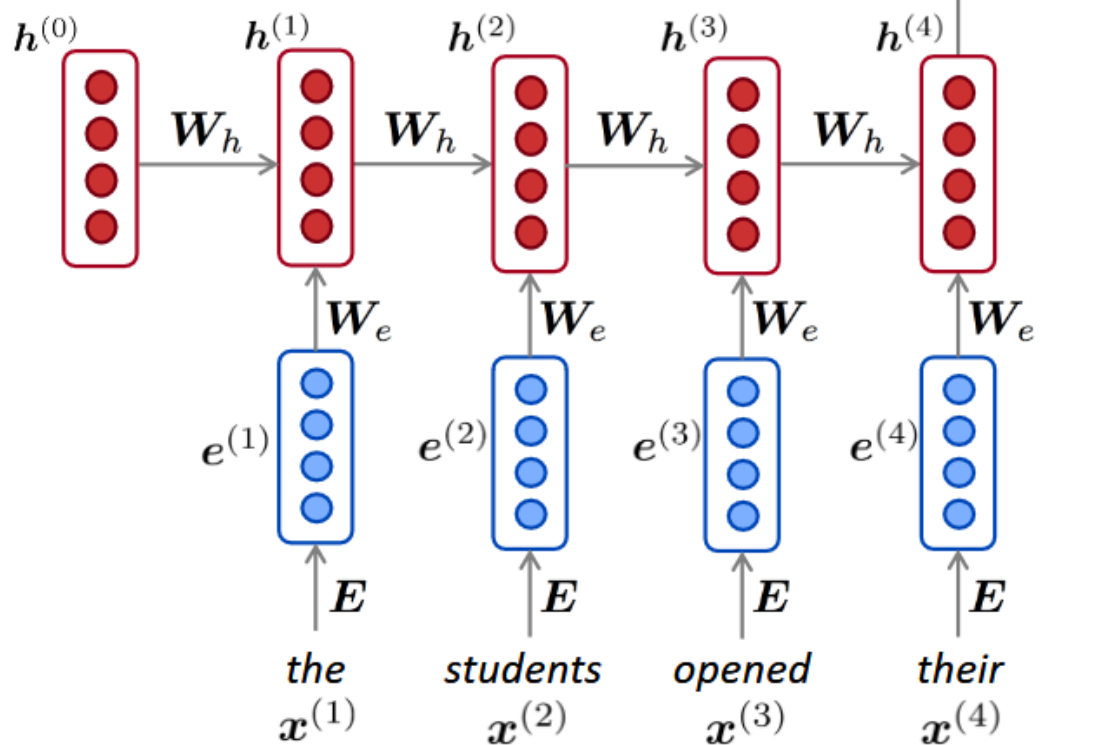$$h^{(t)} = \sigma\left(W_h h^{(t-1)} + W_e e^{(t)} + b_1\right)$$

$h^{(0)}$ is the initial hidden state

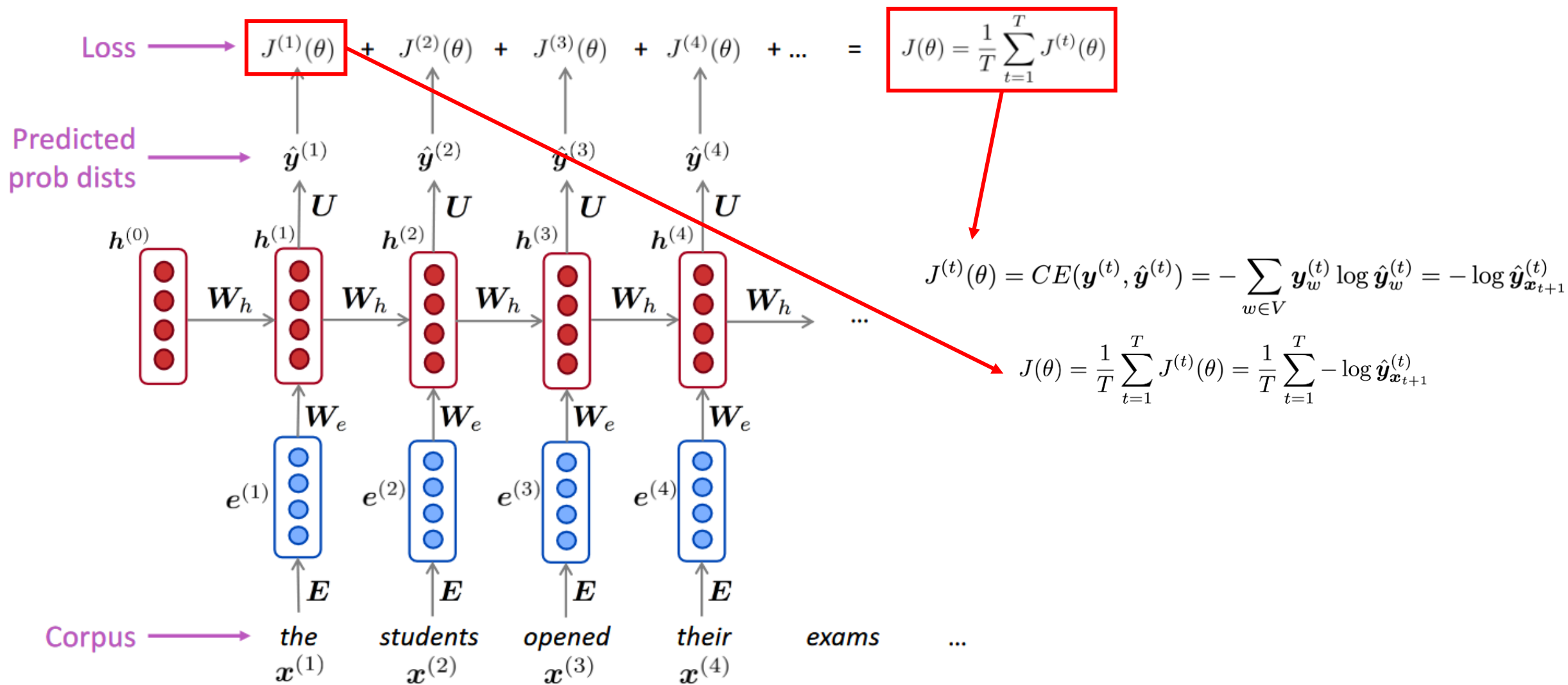word embeddings

$$e^{(t)} = Ex^{(t)}$$

words / one-hot vectors
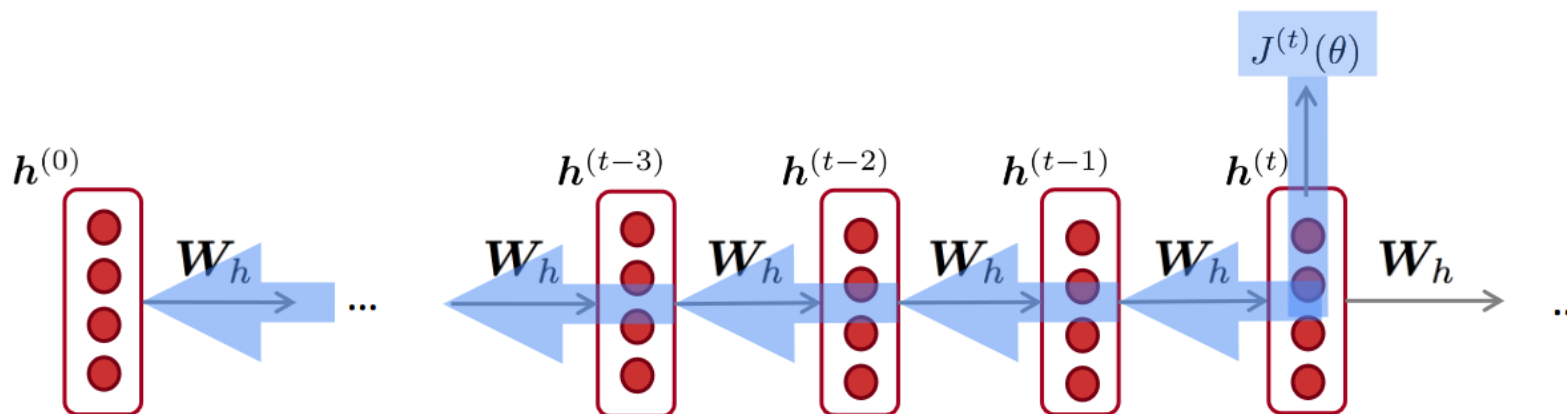
$$x^{(t)} \in \mathbb{R}^{|V|}$$

**Note**: *this input sequence could be much longer, but this slide doesn't have space!*

# Training RNN



Loss $\longrightarrow$ $J^{(1)}(\theta)$ + $J^{(2)}(\theta)$ + $J^{(3)}(\theta)$ + $J^{(4)}(\theta)$ + ... = $J(\theta) = \frac{1}{T}\sum_{t=1}^{T}J^{(t)}(\theta)$

Predicted prob dists $\longrightarrow$ $\hat{\boldsymbol{y}}^{(1)}$ $\hat{\boldsymbol{y}}^{(2)}$ $\hat{\boldsymbol{y}}^{(3)}$ $\hat{\boldsymbol{y}}^{(4)}$

$\boldsymbol{U}$ $\boldsymbol{U}$ $\boldsymbol{U}$ $\boldsymbol{U}$

$\boldsymbol{h}^{(0)}$ $\boldsymbol{h}^{(1)}$ $\boldsymbol{h}^{(2)}$ $\boldsymbol{h}^{(3)}$ $\boldsymbol{h}^{(4)}$

$\boldsymbol{W}_h$ $\boldsymbol{W}_h$ $\boldsymbol{W}_h$ $\boldsymbol{W}_h$ $\boldsymbol{W}_h$ ...

$J^{(t)}(\theta) = CE(\boldsymbol{y}^{(t)}, \hat{\boldsymbol{y}}^{(t)}) = -\sum_{w \in V} \boldsymbol{y}_w^{(t)} \log \hat{\boldsymbol{y}}_w^{(t)} = -\log \hat{\boldsymbol{y}}_{\boldsymbol{x}_{t+1}}^{(t)}$

$J(\theta) = \frac{1}{T}\sum_{t=1}^{T}J^{(t)}(\theta) = \frac{1}{T}\sum_{t=1}^{T} -\log \hat{\boldsymbol{y}}_{\boldsymbol{x}_{t+1}}^{(t)}$

$\boldsymbol{W}_e$ $\boldsymbol{W}_e$ $\boldsymbol{W}_e$ $\boldsymbol{W}_e$

$\boldsymbol{e}^{(1)}$ $\boldsymbol{e}^{(2)}$ $\boldsymbol{e}^{(3)}$ $\boldsymbol{e}^{(4)}$

$\boldsymbol{E}$ $\boldsymbol{E}$ $\boldsymbol{E}$ $\boldsymbol{E}$

Corpus $\longrightarrow$ the students opened their exams ...
$\boldsymbol{x}^{(1)}$ $\boldsymbol{x}^{(2)}$ $\boldsymbol{x}^{(3)}$ $\boldsymbol{x}^{(4)}$

# Backpropagation for RNN

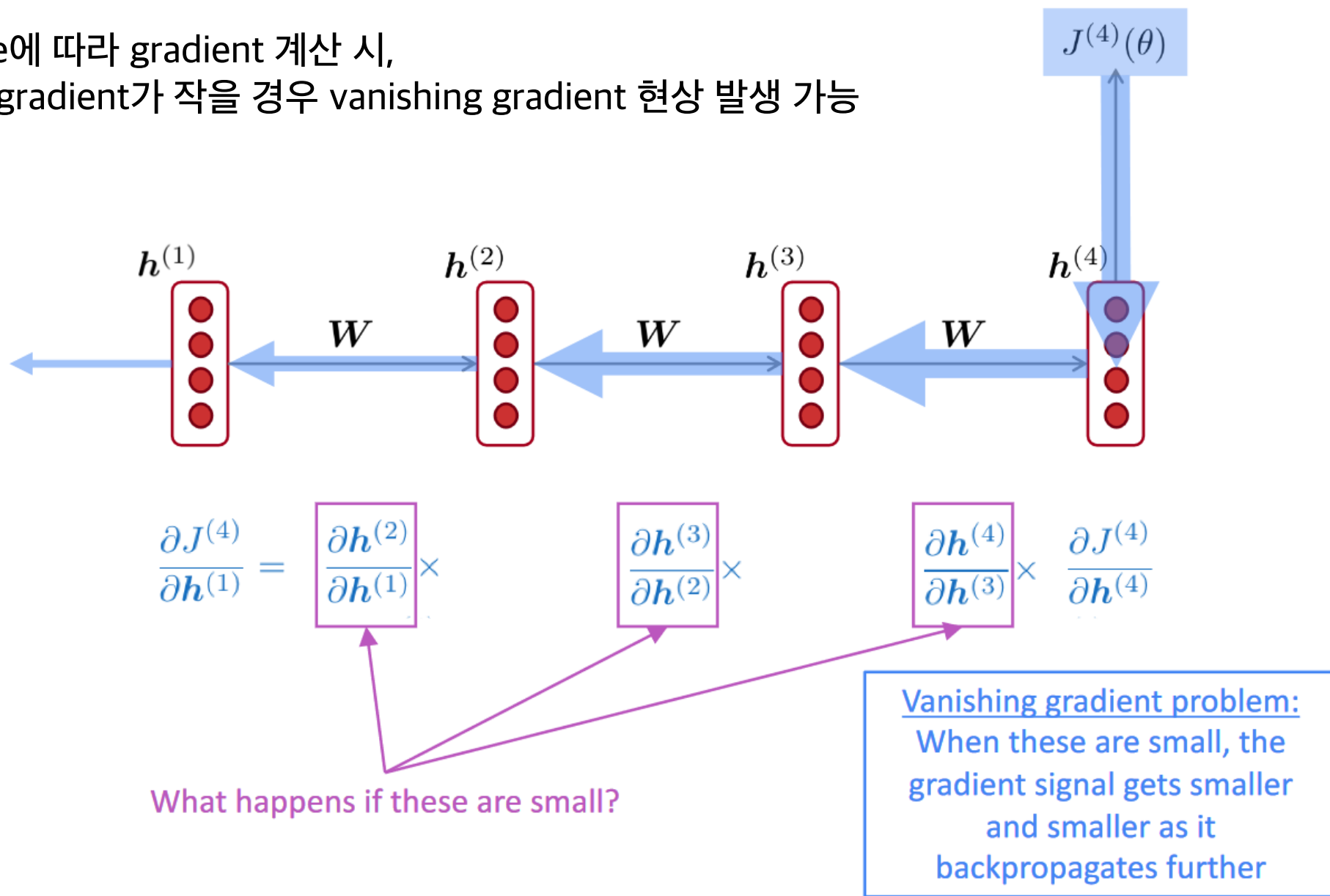BPTT : 시간 역순으로 각 hidden state를 따라 gradient 전달



**Question:** What's the derivative of $J^{(t)}(\theta)$ w.r.t. the repeated weight matrix $W_h$ ?

**Answer:** $\dfrac{\partial J^{(t)}}{\partial W_h} = \displaystyle\sum_{i=1}^{t} \left.\dfrac{\partial J^{(t)}}{\partial W_h}\right|_{(i)}$

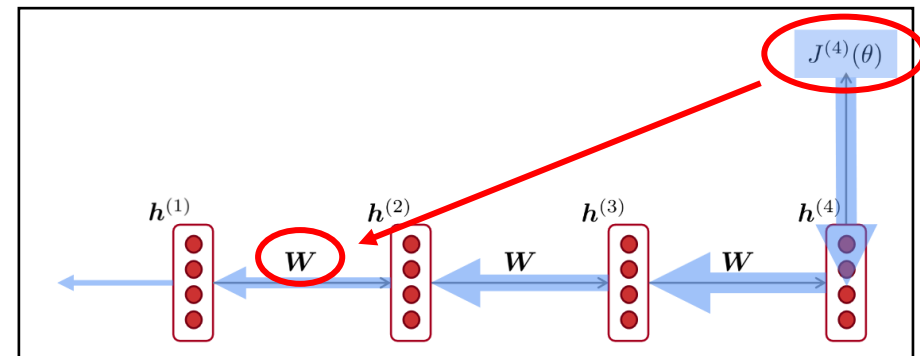"The gradient w.r.t. a repeated weight is the sum of the gradient w.r.t. each time it appears"

# Vanishing gradient

chain rule에 따라 gradient 계산 시,
곱해지는 gradient가 작을 경우 vanishing gradient 현상 발생 가능

$$J^{(4)}(\theta)$$

$$h^{(1)} \qquad h^{(2)} \qquad h^{(3)} \qquad h^{(4)}$$

$$W \qquad W \qquad W$$

$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \frac{\partial h^{(2)}}{\partial h^{(1)}} \times \frac{\partial h^{(3)}}{\partial h^{(2)}} \times \frac{\partial h^{(4)}}{\partial h^{(3)}} \times \frac{\partial J^{(4)}}{\partial h^{(4)}}$$

What happens if these are small?

Vanishing gradient problem:
When these are small, the
gradient signal gets smaller
and smaller as it
backpropagates further

# Vanishing gradient (proof)

https://mmuratarat.github.io/2019-02-07/bptt-of-rnn



$$\frac{\partial E_t}{\partial W} = \sum_{k=1}^{t} \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial W}$$

$$\frac{\partial h_t}{\partial h_k} = \prod_{j=k+1}^{t} \frac{\partial h_j}{\partial h_{j-1}} = \prod_{j=k+1}^{t} W^T \times diag[f'(j_{j-1})]$$
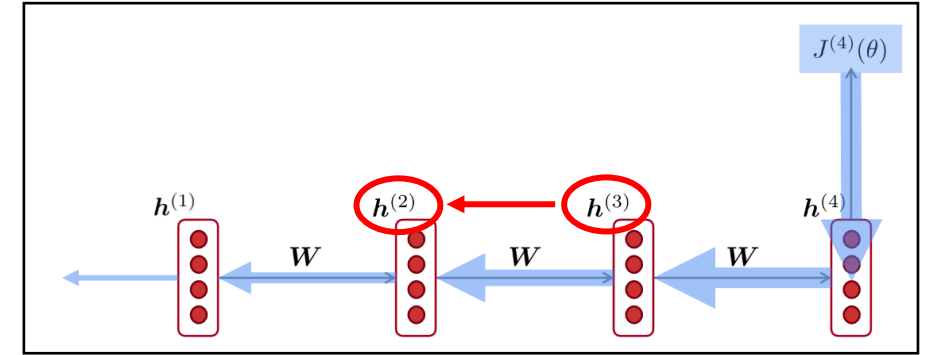
$h_j$을 $h_{j-1}$에만 partial derivative -> Jacobian Matrix

$h_j$는 $f(X_j, h_{j-1})$이므로 chain rule 적용

$$\frac{\partial h_j}{\partial h_{j-1}} = \left[\frac{\partial h_j}{\partial h_{j-1,1}} \cdots \frac{\partial h_j}{\partial h_{j-1,D_n}}\right] = \begin{bmatrix} \frac{\partial h_{j,1}}{\partial h_{j-1,1}} & \cdots & \frac{\partial h_{j,1}}{\partial h_{j-1,D_n}} \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \frac{\partial h_{j,D_n}}{\partial h_{j-1,1}} & \cdots & \frac{\partial h_{j,D_n}}{\partial h_{j-1,D_n}} \end{bmatrix}$$

좌측의 Jacobian Matrix 계산 시,

$W^T \times diag[f'(h_{j-1})]$ 로 분해 가능 (= $W^T$의 각 행에 $f'(h_{j-1,i})$ 곱)

# Vanishing gradient (proof)



$$\frac{\partial h_t}{\partial h_k} = \prod_{j=k+1}^{t} \frac{\partial h_j}{\partial h_{j-1}} = \prod_{j=k+1}^{t} W^T \times diag[f'(j_{j-1})]$$

$$\left\| \frac{\partial h_j}{\partial h_{j-1}} \right\| \leq \| W^T \| \| diag[f'(h_{j-1})] \| \leq \beta_W \beta_h$$

L2 norm에 대하여 upper bound $\beta_W, \beta_h$로 둘 경우,

좌측과 같이 partial derivative의 upper bound 계산 가능

$$\left\| \frac{\partial h_t}{\partial h_k} \right\| = \left\| \prod_{j=k+1}^{t} \frac{\partial h_j}{\partial h_{j-1}} \right\| \leq (\beta_W \beta_h)^{t-k}$$

$\beta_W \beta_h$가 1보다 작고 t-k가 충분히 큰 수일 경우, $(\beta_W \beta_h)^{t-k} \to 0$

$\beta_W \beta_h$가 1보다 크고 t-k가 충분히 큰 수일 경우, $(\beta_W \beta_h)^{t-k} \to large$

=> gradient vanishing/exploding 발생 가능

# Effect of vanishing gradient

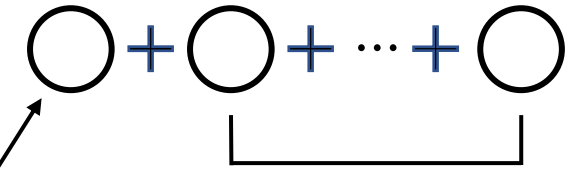$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^{T} -\log \hat{y}_{x_{t+1}}^{(t)}$$

$$\frac{\partial E_t}{\partial W} = \sum_{k=1}^{t} \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial W}$$

$\bigcirc + \bigcirc + \cdots + \bigcirc$

gradient signal
from close-by
-> influence ↑

gradient signal
from faraway
-> influence ↓

$h^{(1)}$   $h^{(2)}$   $h^{(3)}$   $h^{(4)}$

$J^{(2)}(\theta)$   $J^{(4)}(\theta)$

$W$   $W$   $W$

Gradient signal from faraway is lost because it's much
smaller than gradient signal from close-by.

So model weights are updated only with respect to near
effects, not long-term effects.

# Effect of vanishing gradient

If the gradient becomes vanishingly small over longer distances (step $t$ to step $t+n$), then we can't tell whether:

1. There's no dependency between step $t$ and $t+n$ in the data
2. We have wrong parameters to capture the true dependency between $t$ and $t+n$

But if gradient is small, the model can't learn this dependency
- So the model is unable to predict similar long-distance dependencies at test time

model the dependency
= learn the connection between layers

# Way to fix vanishing gradient

In a vanilla RNN, the hidden state is constantly being rewritten

$$h^{(t)} = \sigma \left( W_h h^{(t-1)} + W_x x^{(t)} + b \right)$$

-> not easy to preserve information one hidden state to another,
   particularly putting it through the non-linear function
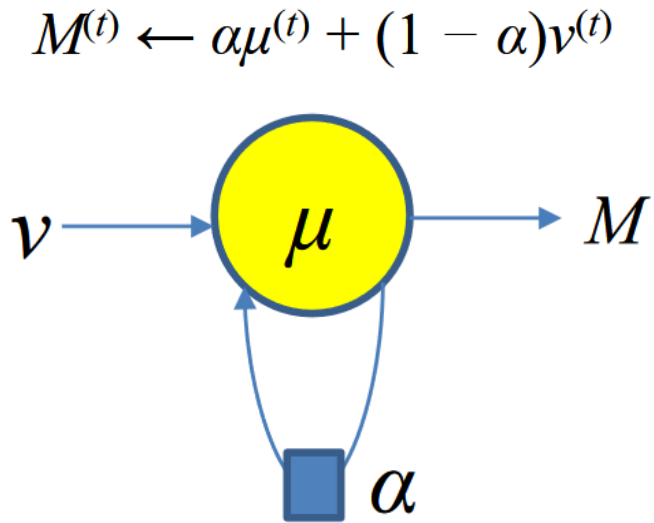
How about a RNN with separate memory?

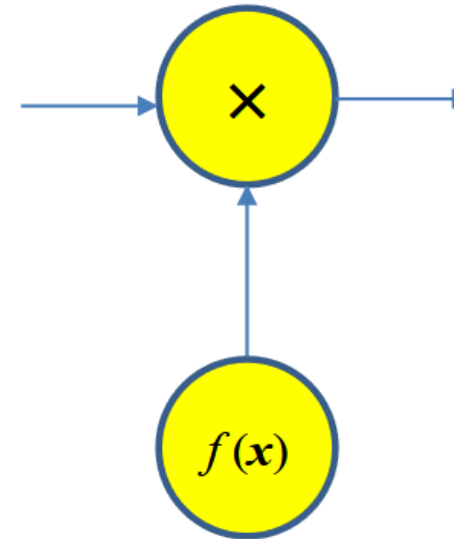-> separate place to store information to preserve

# LSTM

# Gating/Leaky Unit

gating unit : controls the flow of information
leaky unit : determine how much information to preserve (including self-loop)

$$M^{(t)} \leftarrow \alpha\mu^{(t)} + (1 - \alpha)v^{(t)}$$



leaky unit
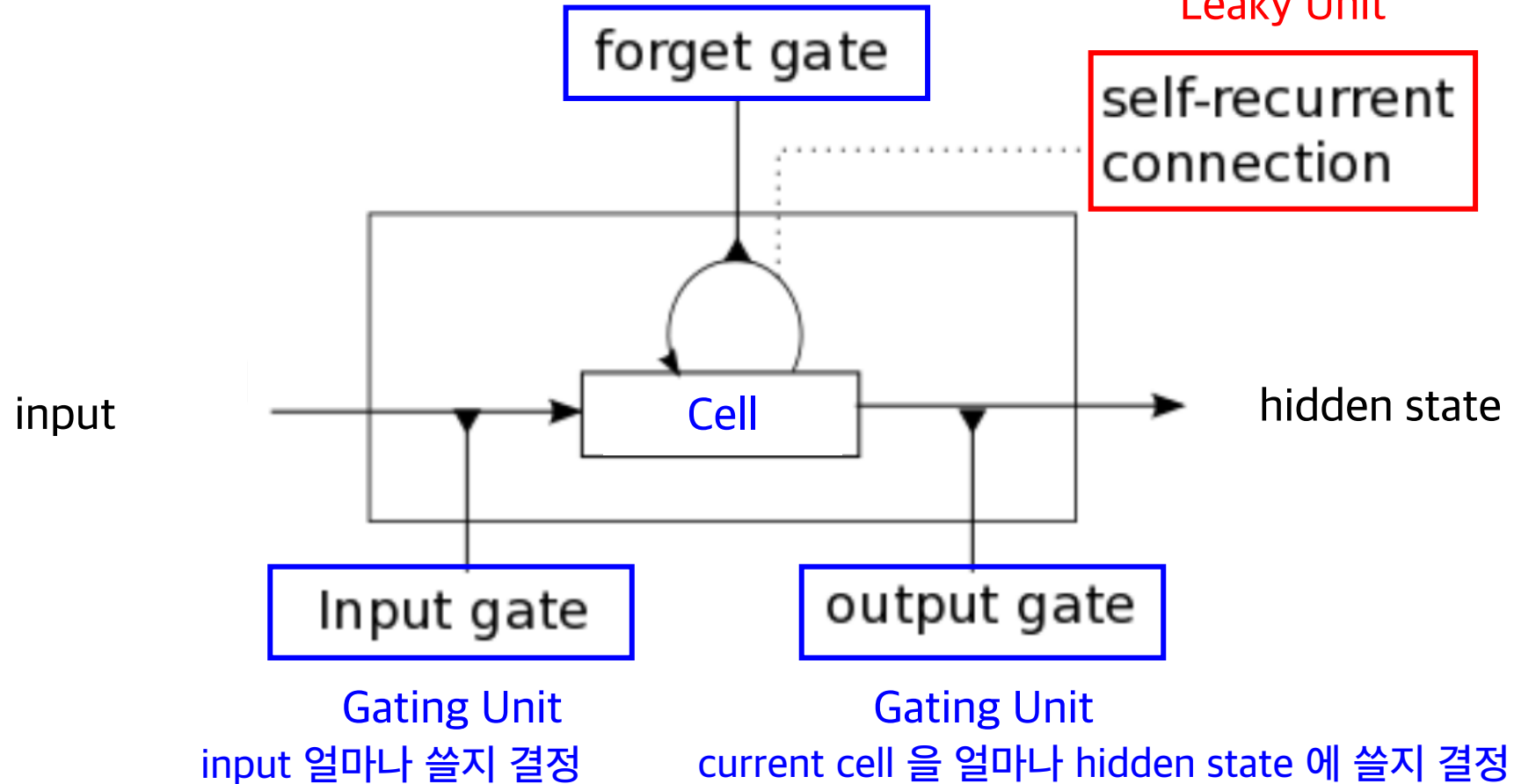
α = 1 일 경우, 정보를 보존
α = 0 일 경우, 정보를 잊음

gating unit

f(x) = 1 일 경우, 정보를 보존
f(x) = 0 일 경우, 정보를 잊음

# LSTM (= Long Short Term Memory)

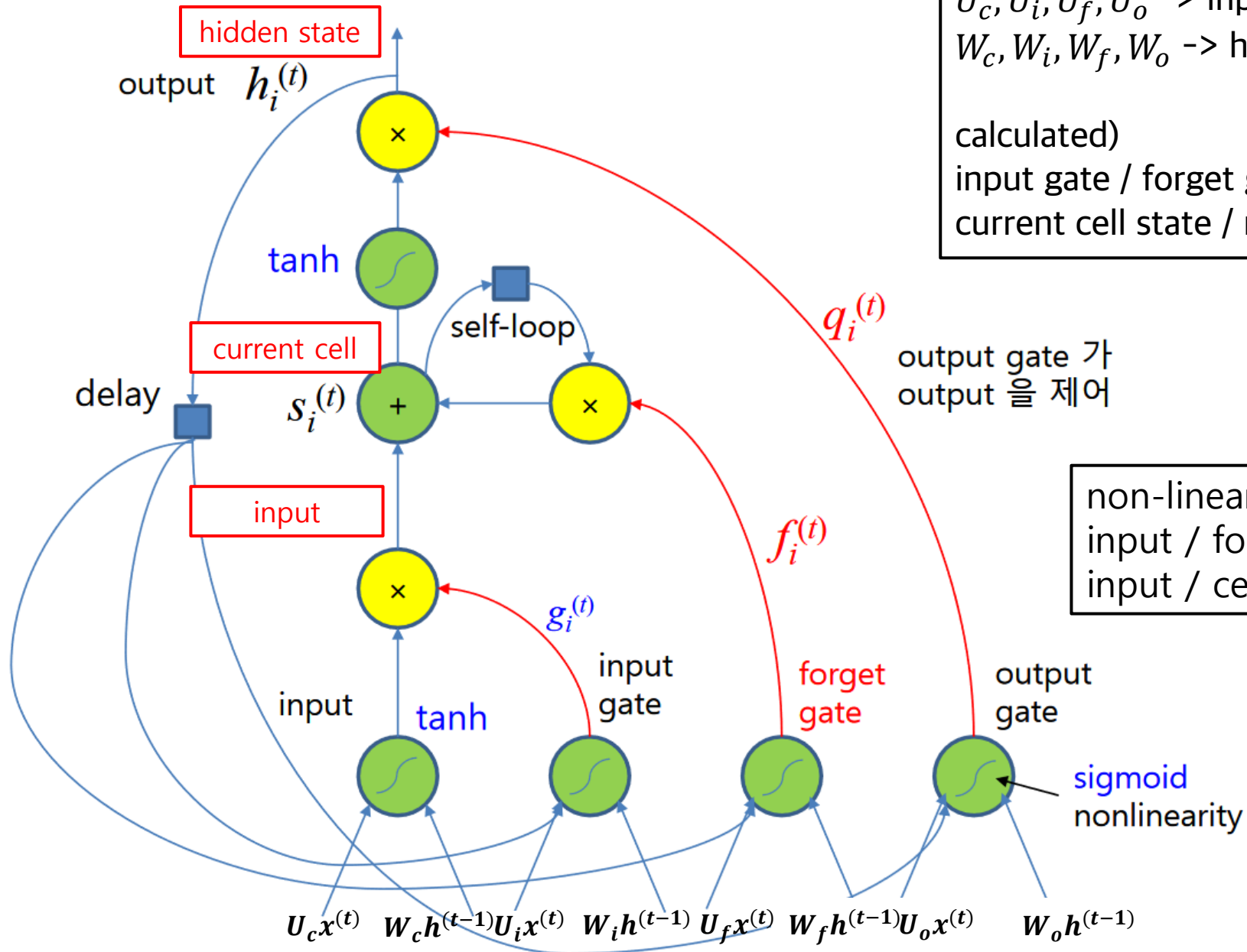input / prev cell 을 얼마나 current cell 에 쓸지 결정

Gating Unit

Leaky Unit

forget gate

self-recurrent connection

input

Cell

hidden state

Input gate

output gate

Gating Unit
input 얼마나 쓸지 결정

Gating Unit
current cell 을 얼마나 hidden state 에 쓸지 결정

# LSTM

weight matrix)
$U_c, U_i, U_f, U_o$ -> input $x^{(t)}$ 와 Affine 연산
$W_c, W_i, W_f, W_o$ -> hidden state $h^{(t)}$ 와 Affine 연산

calculated)
input gate / forget gate / output gate (0~1 사이의 값)
current cell state / next cell state / hidden state

non-linearity function)
input / forget / output gate -> sigmoid
input / cell -> tanh

hidden state

output  $h_i^{(t)}$

tanh

self-loop

current cell

delay    $s_i^{(t)}$

input

$q_i^{(t)}$
output gate 가
output 을 제어

$f_i^{(t)}$

$g_i^{(t)}$

input    tanh    input gate    forget gate    output gate

input

sigmoid nonlinearity

$U_c x^{(t)}$   $W_c h^{(t-1)} U_i x^{(t)}$   $W_i h^{(t-1)}$   $U_f x^{(t)}$   $W_f h^{(t-1)} U_o x^{(t)}$   $W_o h^{(t-1)}$

# LSTM

We have a sequence of inputs $x^{(t)}$, and we will compute a sequence of hidden states $h^{(t)}$ and cell states $c^{(t)}$. On timestep $t$:

**Forget gate:** controls what is kept vs forgotten, from previous cell state

**Sigmoid function:** all gate values are between 0 and 1

**Input gate:** controls what parts of the new cell content are written to cell

**Output gate:** controls what parts of cell are output to hidden state

$$f^{(t)} = \sigma\left(W_f h^{(t-1)} + U_f x^{(t)} + b_f\right)$$

$$i^{(t)} = \sigma\left(W_i h^{(t-1)} + U_i x^{(t)} + b_i\right)$$

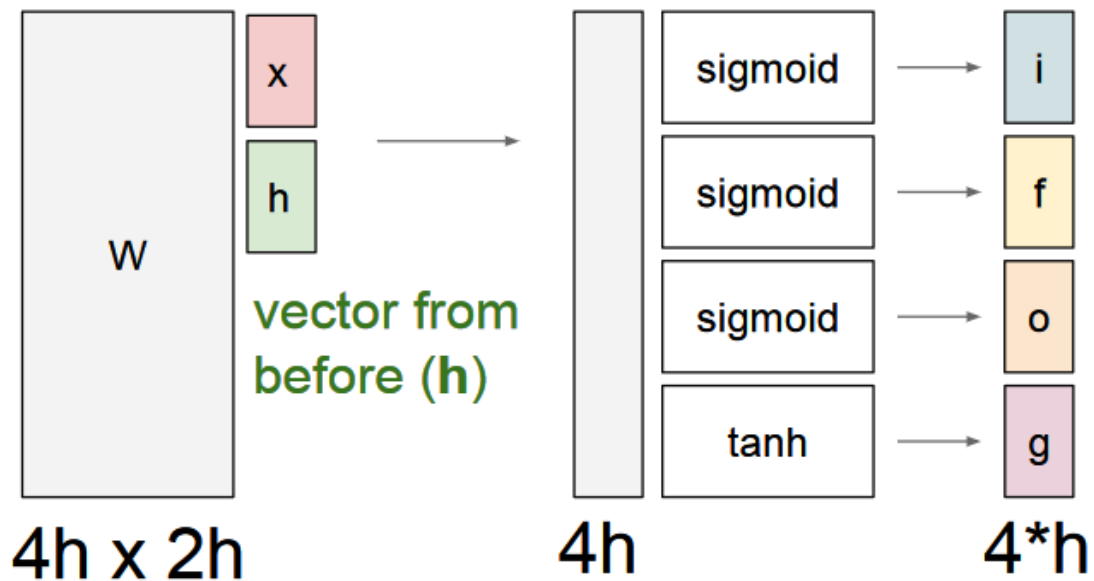$$o^{(t)} = \sigma\left(W_o h^{(t-1)} + U_o x^{(t)} + b_o\right)$$

**New cell content:** this is the new content to be written to the cell

**Cell state:** erase ("forget") some content from last cell state, and write ("input") some new cell content

$$\tilde{c}^{(t)} = \tanh\left(W_c h^{(t-1)} + U_c x^{(t)} + b_c\right)$$

$$c^{(t)} = f^{(t)} \circ c^{(t-1)} + i^{(t)} \circ \tilde{c}^{(t)}$$

**Hidden state:** read ("output") some content from the cell

$$h^{(t)} = o^{(t)} \circ \tanh c^{(t)}$$

Gates are applied using element-wise product

All these are vectors of same length $n$

# Matrix Affine



$4h \times 2h$       $4h$      $4*h$

$U_c, U_i, U_f, U_o$ / $W_c, W_i, W_f, W_o$를
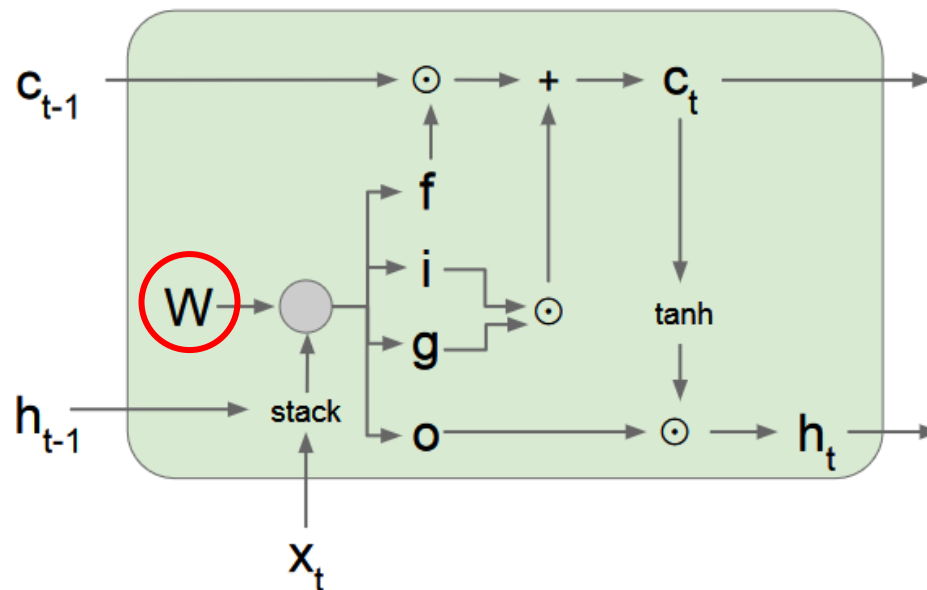하나의 weight matrix W로 통합

-> matrix size $4h \times 2h$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
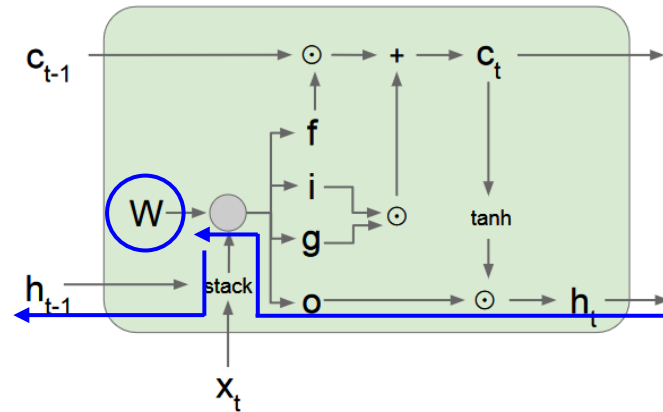
weight matrix W와 Affine 계산 후, hadamard product

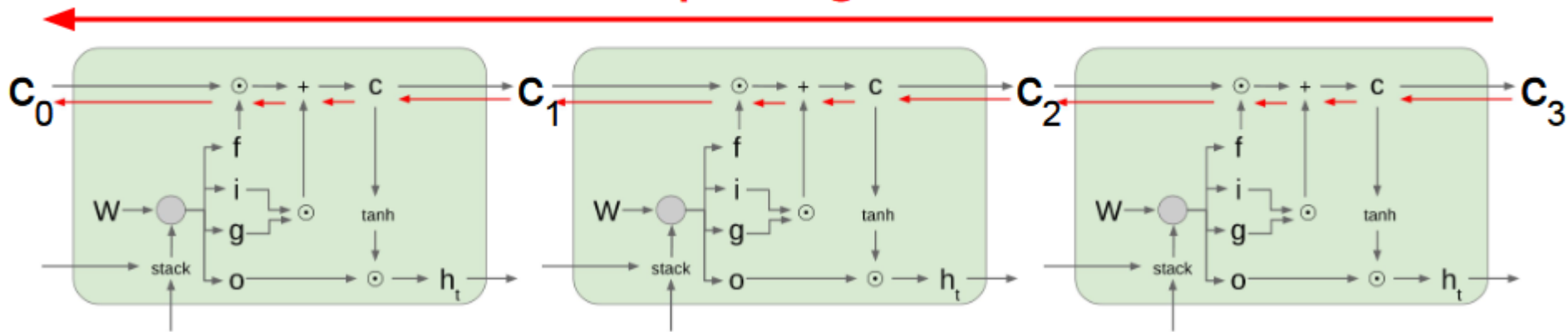$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

# Gradient Flow



hidden state 를 통해 backpropagation 이루어질 경우,

weight matrix W에 대한 partial derivative 가

반복적으로 곱해지면서 vanishing gradient 발생 가능
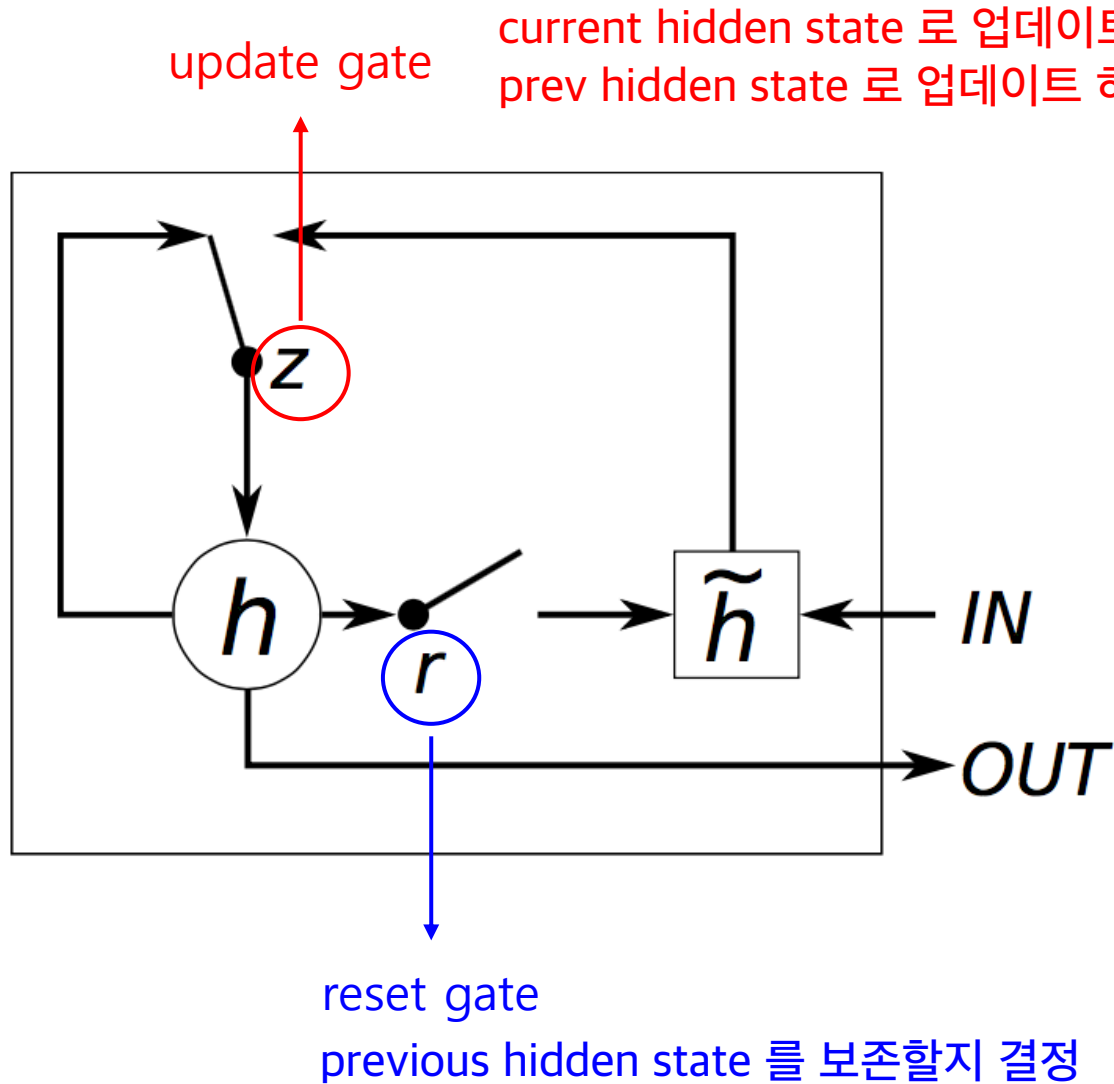
## Uninterrupted gradient flow!



cell state 를 통해 back propagation 이루어질 경우, weight matrix 의 interrupt 없이

Hadamard product 에 대한 partial derivative 만 고려하며 gradient flow 가능 -> gradient vanishing ↓↓

# GRU

# GRU (= Gated Recurrent Unit)



update gate

current hidden state 로 업데이트 할지 =
prev hidden state 로 업데이트 하지 않을지 결정

reset gate
previous hidden state 를 보존할지 결정

LSTM과 달리, cell memory 가 없으므로 complexity ↓

long term memory 를 보존하고 싶다면 update gate z = 0

# GRU

Update gate: controls what parts of hidden state are updated vs preserved

Reset gate: controls what parts of previous hidden state are used to compute new content

New hidden state content: reset gate selects useful parts of prev hidden state. Use this and current input to compute new hidden content.

Hidden state: update gate simultaneously controls what is kept from previous hidden state, and what is updated to new hidden state content

$$u^{(t)} = \sigma\left(W_u h^{(t-1)} + U_u x^{(t)} + b_u\right)$$

$$r^{(t)} = \sigma\left(W_r h^{(t-1)} + U_r x^{(t)} + b_r\right)$$

$$\tilde{h}^{(t)} = \tanh\left(W_h(r^{(t)} \circ h^{(t-1)}) + U_h x^{(t)} + b_h\right)$$

$$h^{(t)} = (1 - u^{(t)}) \circ h^{(t-1)} + u^{(t)} \circ \tilde{h}^{(t)}$$

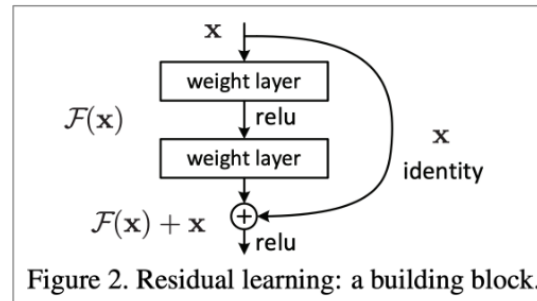How does this solve vanishing gradient?
Like LSTM, GRU makes it easier to retain info long-term (e.g. by setting update gate to 0)

# Is vanishing/exploding gradient just a RNN problem?

No! It can be a problem for all neural architectures (including feed-forward and convolutional), especially deep ones.
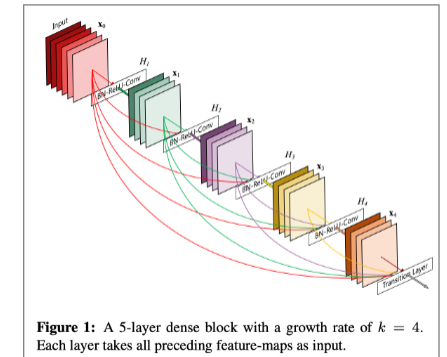
For example:

- Residual connections aka "ResNet"
- Also known as skip-connections
- The identity connection preserves information by default
- This makes deep networks much easier to train



$\mathcal{F}(\mathbf{x})$

weight layer

relu

weight layer

$\mathbf{x}$
identity

$\mathcal{F}(\mathbf{x}) + \mathbf{x}$

relu

Figure 2. Residual learning: a building block.

For example:

- Dense connections aka "DenseNet"
- Directly connect each layer to all future layers!



**Figure 1:** A 5-layer dense block with a growth rate of $k = 4$. Each layer takes all preceding feature-maps as input.

For example:

- Highway connections aka "HighwayNet"
- Similar to residual connections, but the identity connection vs the transformation layer is controlled by a dynamic gate
- Inspired by LSTMs, but applied to deep feedforward/convolutional networks

# Bidirectional RNN



positive

Sentence encoding

element-wise mean/max

element-wise mean/max

the    movie    was    terribly    exciting    !

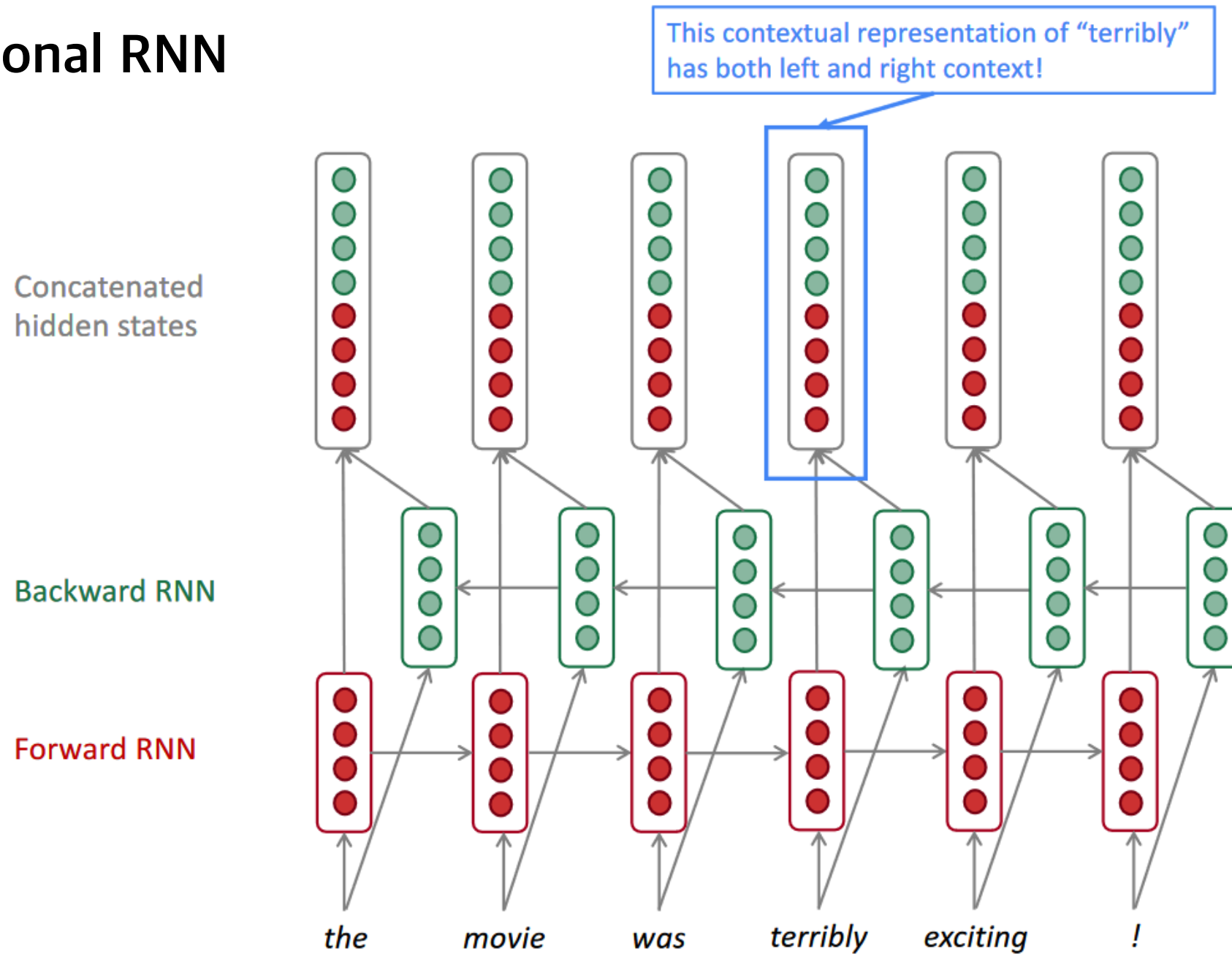We can regard this hidden state as a representation of the word *"terribly"* in the context of this sentence. We call this a *contextual representation.*

These contextual representations only contain information about the *left* context (e.g. *"the movie was"*).

**What about *right* context?**

In this example, *"exciting"* is in the right context and this modifies the meaning of *"terribly"* (from negative to positive)

# Bidirectional RNN



This contextual representation of "terribly" has both left and right context!

Concatenated hidden states

Backward RNN

Forward RNN

the     movie     was     terribly     exciting     !

# Multi-layer RNN



The hidden states from RNN layer *i* are the inputs to RNN layer *i+1*

RNN layer 3

RNN layer 2

RNN layer 1

the    movie    was    terribly    exciting    !