



DSBA CS224n 2021 Study

# **[Lecture 05] Language Models and RNNs**

## **[Lecture 06] Simple and LSTM RNNs**

---



고려대학교 산업경영공학과

Data Science & Business Analytics Lab

발표자 : 고유경



- 1 Language Model
- 2 RNN Language Model
- 3 Problem of RNN
- 4 LSTM
- 5 Conclusion





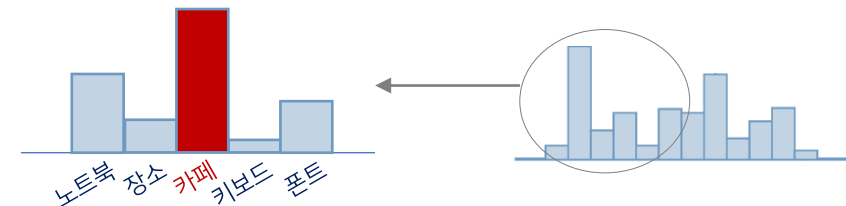
- 확률 분포를 기반으로 주어진 문맥(sequence) 이후에 위치할 단어 예측

- 시퀀스의 결합 확률 by multiplication rule

$$P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}) = P(\mathbf{x}^{(1)}) \times P(\mathbf{x}^{(2)} | \mathbf{x}^{(1)}) \times \dots \times P(\mathbf{x}^{(T)} | \mathbf{x}^{(T-1)}, \dots, \mathbf{x}^{(1)})$$

$$= \prod_{t=1}^T P(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(1)})$$

- $P(x^{(3)} | \text{코딩하기, 좋은})$



## N-gram Language Model

$P(\text{코딩하기, 좋은, 카페, 추천})$

uni-gram

$P(\text{코딩하기}) \quad P(\text{좋은}) \quad P(\text{카페}) \quad P(\text{추천})$

bi-gram

$P(\text{코딩하기}) \quad P(\text{좋은} | \text{코딩하기}) \quad P(\text{카페} | \text{좋은}) \quad P(\text{추천} | \text{카페})$

$$P(\text{카페} | \text{코딩하기, 좋은}) = \frac{\text{count}(\text{코딩하기, 좋은, 카페})}{\text{count}(\text{코딩하기, 좋은})}$$

- ✓ 이전에 등장한 n개의 단어를 바탕으로 예측
- ✓ count(빈도)를 기반으로 확률 계산
- ✓ (-) sparsity problem: n-gram chunk 문서 내 등장 x



## N-gram Language Model

 $P(\text{코딩하기, 좋은, 카페, 추천})$ 

uni-gram

 $P(\text{코딩하기}) \quad P(\text{좋은}) \quad P(\text{카페}) \quad P(\text{추천})$ 

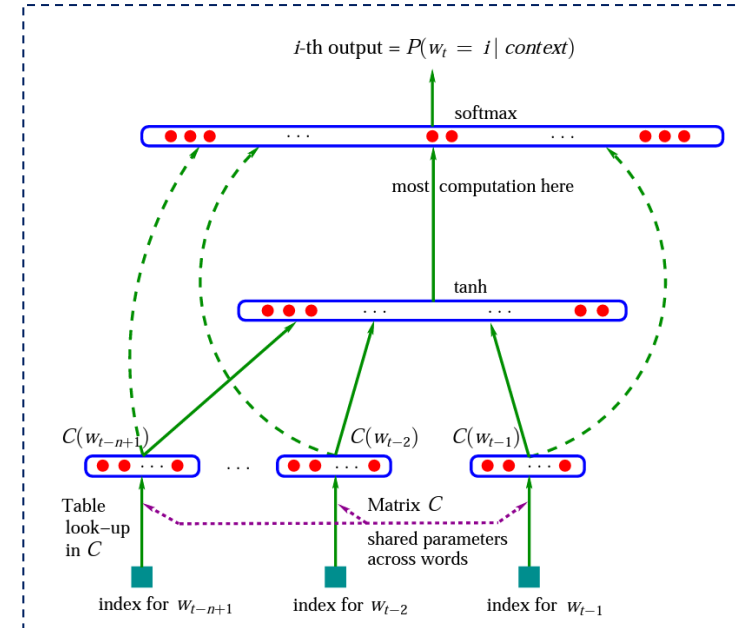
bi-gram

 $P(\text{코딩하기}) \quad P(\text{좋은 | 코딩하기}) \quad P(\text{카페 | 좋은}) \quad P(\text{추천 | 카페})$ 

$$P(\text{카페 | 코딩하기, 좋은}) = \frac{\text{count}(\text{코딩하기, 좋은, 카페})}{\text{count}(\text{코딩하기, 좋은})}$$

- ✓ 이전에 등장한 n개의 단어를 바탕으로 예측
- ✓ count(빈도)를 기반으로 확률 계산
- ✓ (-) sparsity problem: n-gram chunk 문서 내 등장 x

## Neural Language Model



A neural probabilistic language model (Bengio et al, 2003)

- ✓ 미리 지정한 window size 이전 단어를 바탕으로 예측
- ✓ word embedding: 단어의 distributed representation 학습
- ✓ (-) window size가 커지면 파라미터 수 증가 -> 연산 부담, 과적합 가능

## N-gram Language Model

 $P(\text{코딩하기, 좋은, 카페, 추천})$ 

uni-gram

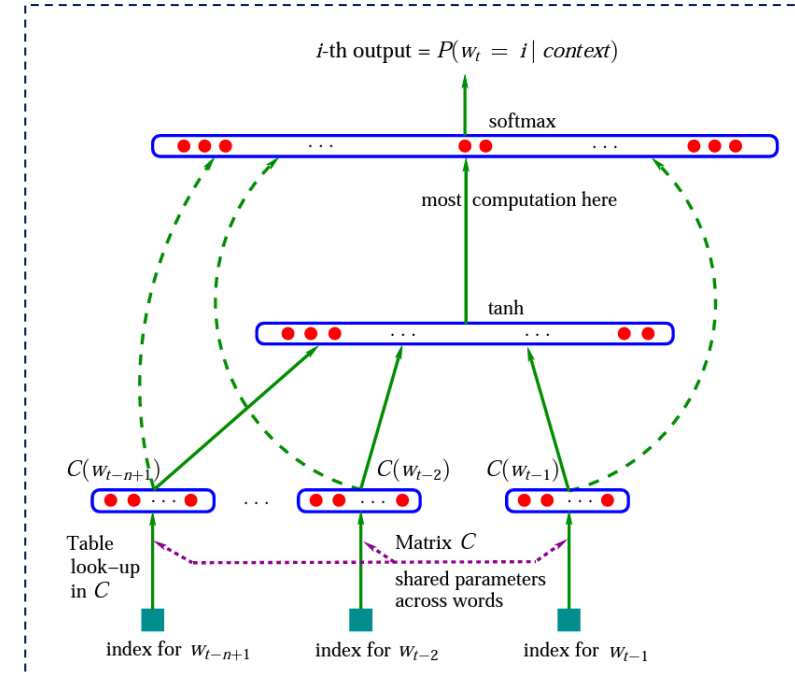
 $P(\text{코딩하기}) \quad P(\text{좋은}) \quad P(\text{카페}) \quad P(\text{추천})$ 

bi-gram

 $P(\text{코딩하기}) \quad P(\text{좋은 | 코딩하기}) \quad P(\text{카페 | 좋은}) \quad P(\text{추천 | 카페})$ 

$$P(\text{카페 | 코딩하기, 좋은}) = \frac{\text{count}(\text{코딩하기, 좋은, 카페})}{\text{count}(\text{코딩하기, 좋은})}$$

## Neural Language Model



A neural probabilistic language model (Bengio et al, 2003)

✓ 입력 길이가 고정되어(N, window size) 이전에 등장하는 모든 단어를 고려할 수 없음



## 02

## RNN Language Model

## Architecture

## ④ Output distribution

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

## ③ Hidden States

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

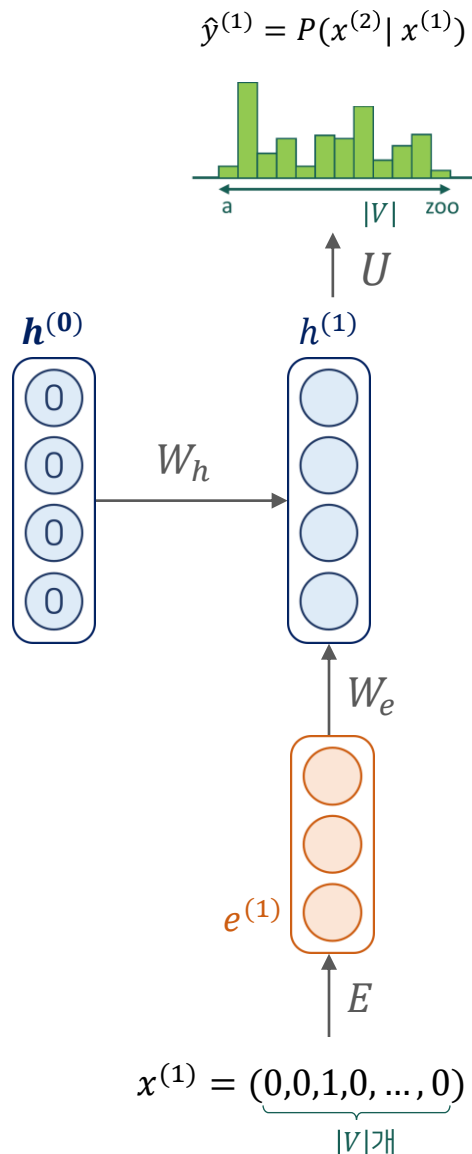
$$W_h: (d_h \times d_h), W_e: (d_e \times d_h)$$

## ② Word embedding

$$e^{(t)} = E x^{(t)} \quad E: (|V| \times d_e)$$

## ① Input word sequence

- 원핫 벡터  $x^{(t)} \in \mathbb{R}^{|V|}$



## 02

## RNN Language Model

## Architecture

## ④ Output distribution

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

## ③ Hidden States

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

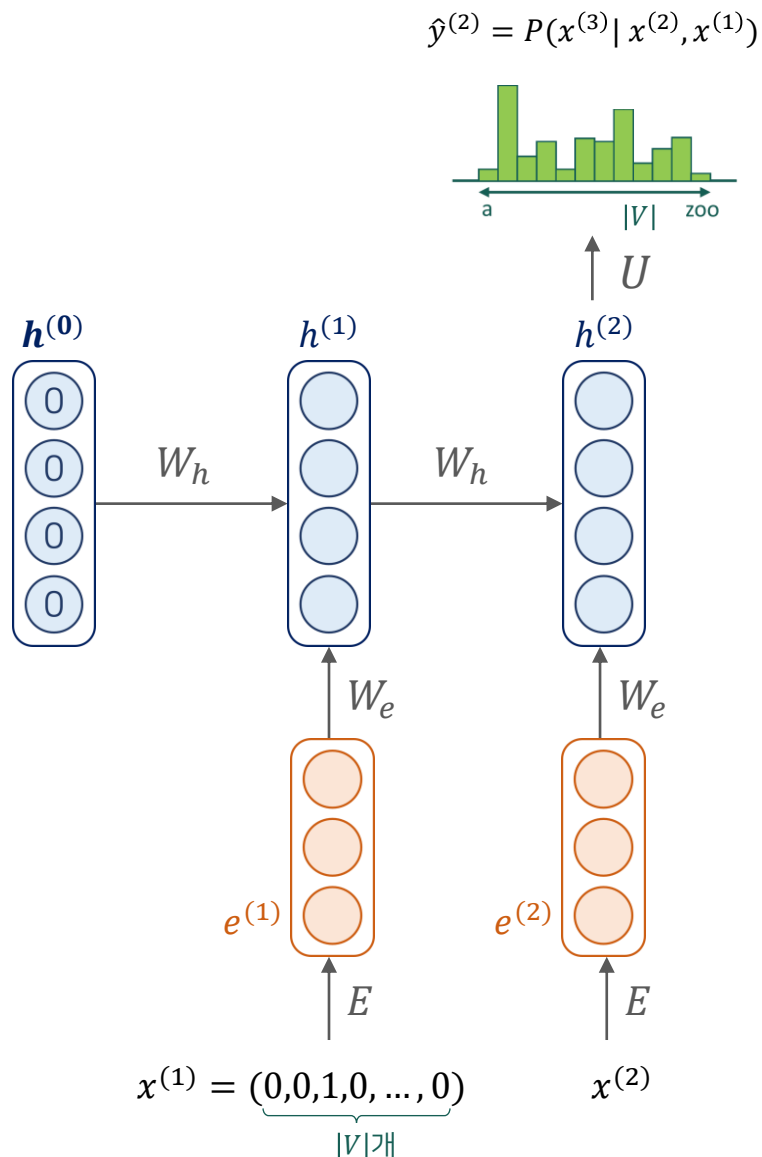
$$W_h: (d_h \times d_h), W_e: (d_e \times d_h)$$

## ② Word embedding

$$e^{(t)} = E x^{(t)} \quad E: (|V| \times d_e)$$

## ① Input word sequence

- 원핫 벡터  $x^{(t)} \in \mathbb{R}^{|V|}$





## 02

## RNN Language Model

## Architecture

## ④ Output distribution

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

## ③ Hidden States

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

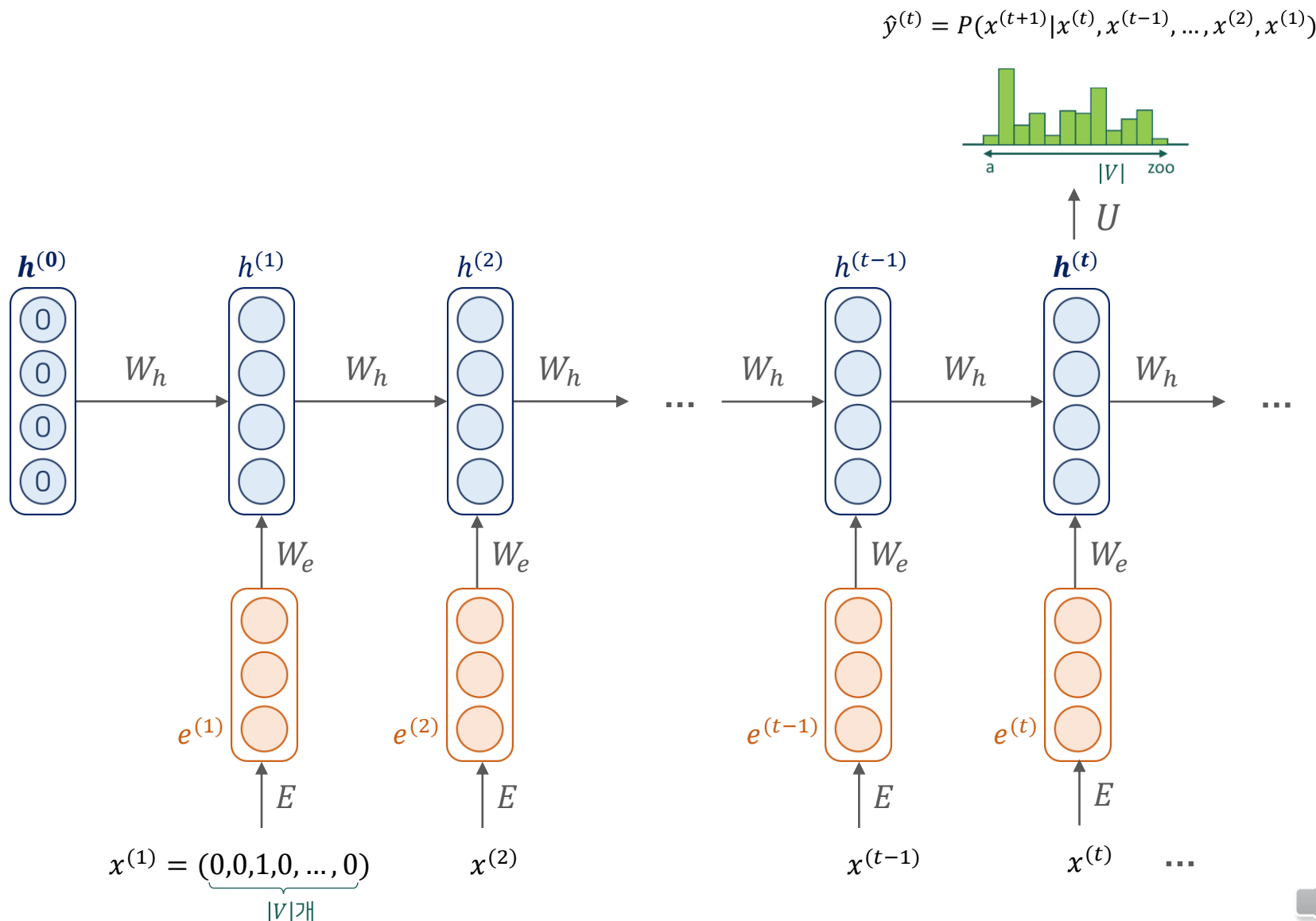
$$W_h: (d_h \times d_h), W_e: (d_e \times d_h)$$

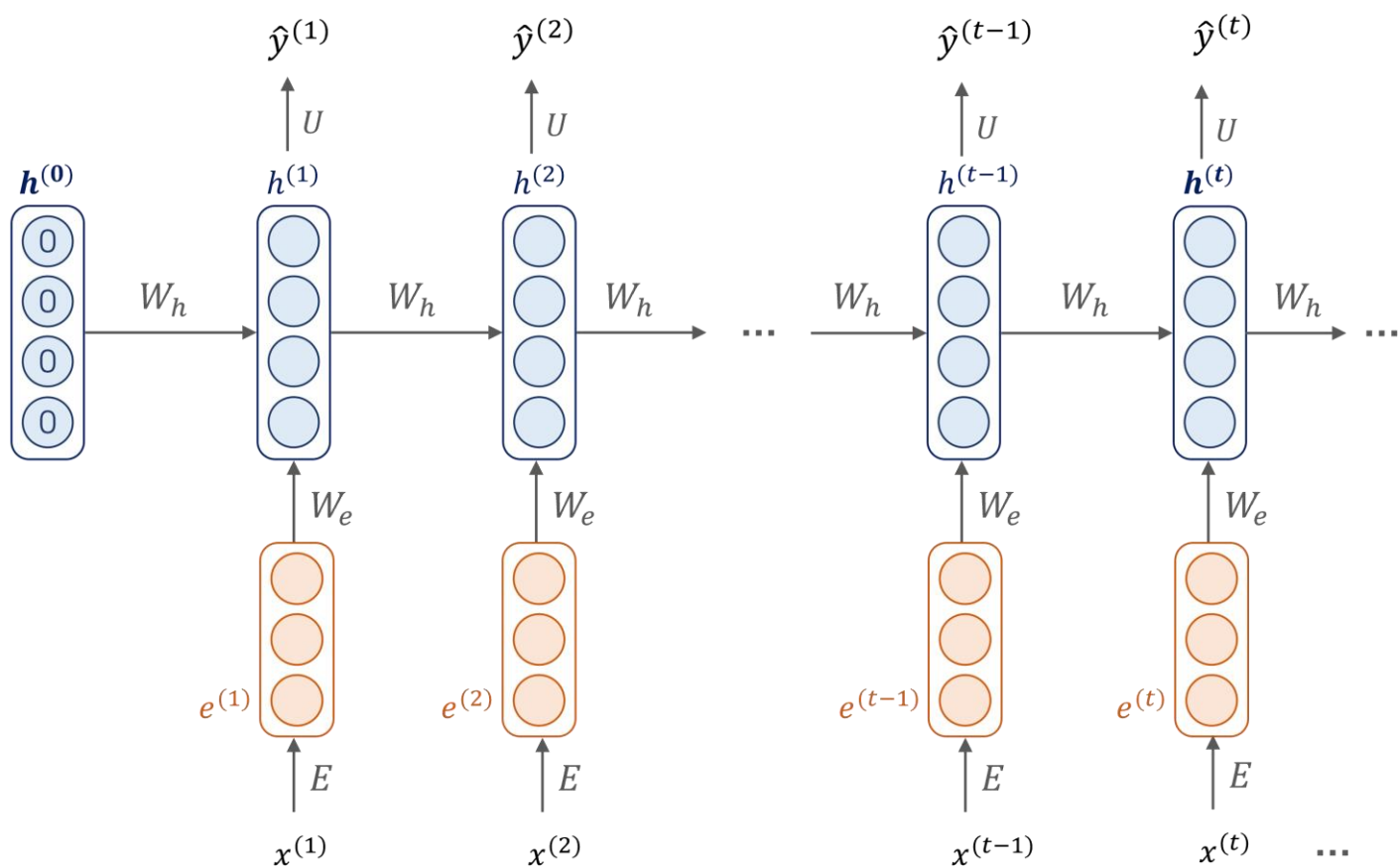
## ② Word embedding

$$e^{(t)} = E x^{(t)} \quad E: (|V| \times d_e)$$

## ① Input word sequence

- 원핫 벡터  $x^{(t)} \in \mathbb{R}^{|V|}$



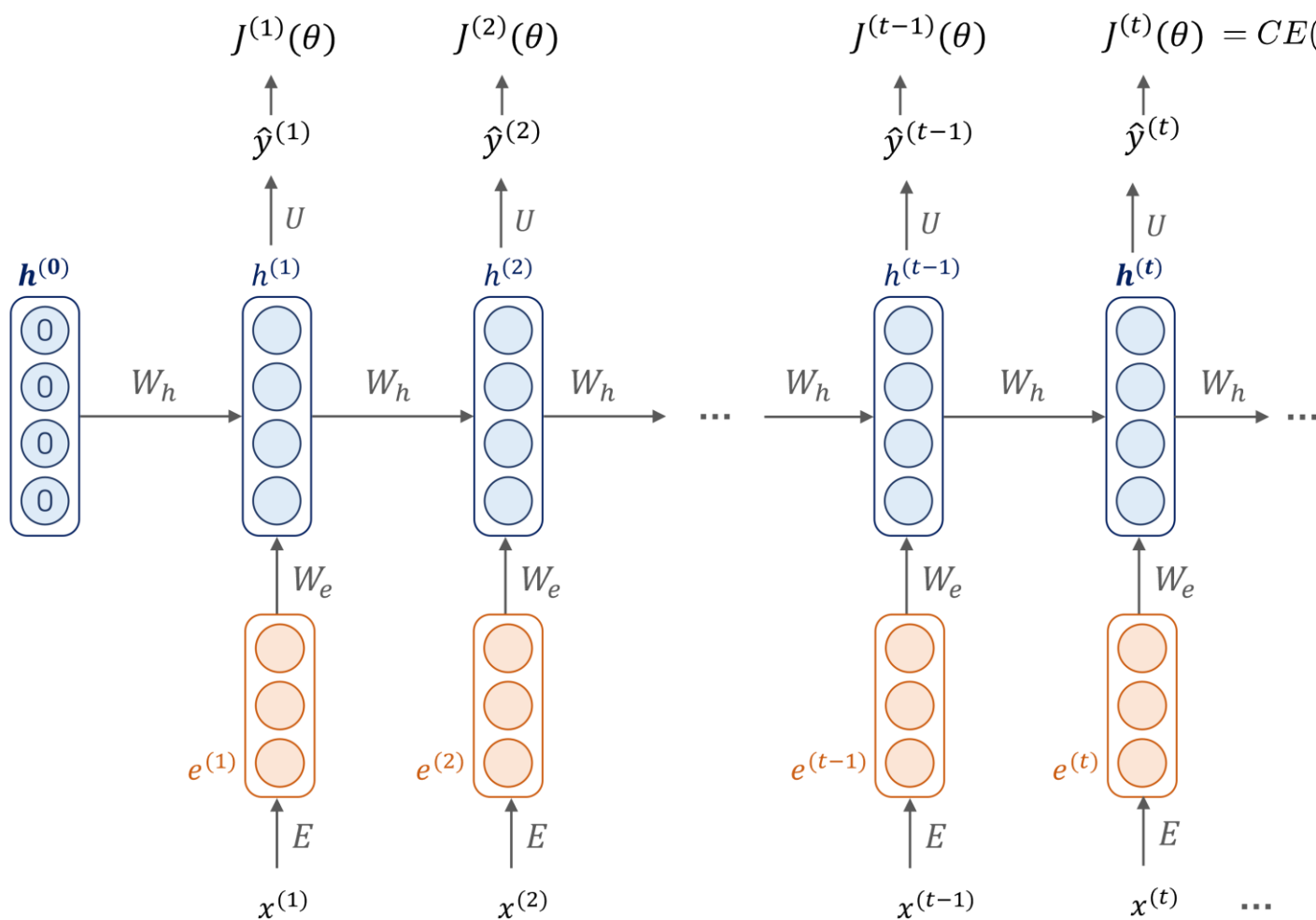


- ✓ many-to-many
- ✓ 매 time step마다 입,출력 존재

## 02

## RNN Language Model

## Loss Function



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

## Cross Entropy Loss

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

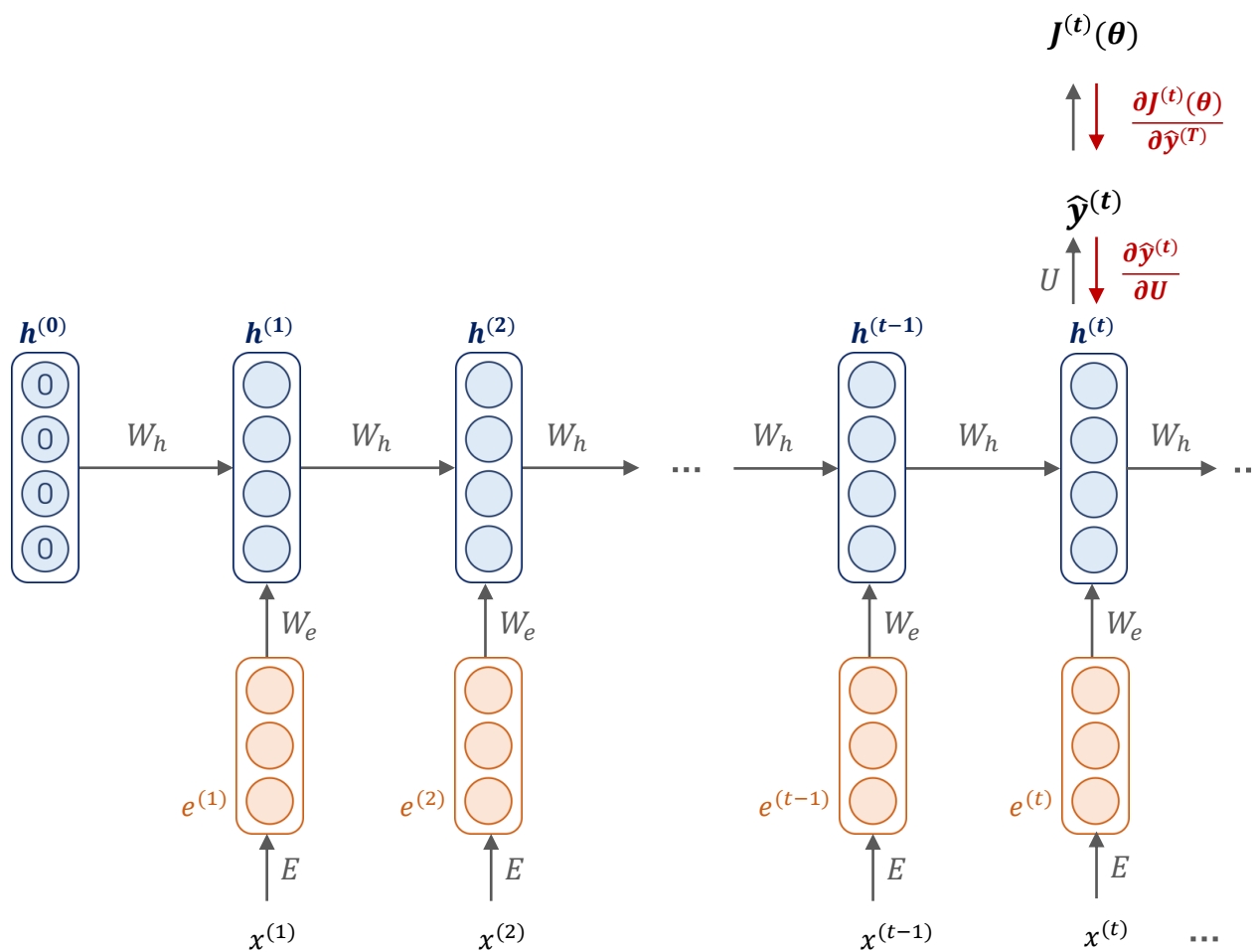
- ✓ many-to-many
- ✓ 매 time step마다 입,출력 존재



## 02

## RNN Language Model

Back Propagation Through Time (BPTT) ①



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

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

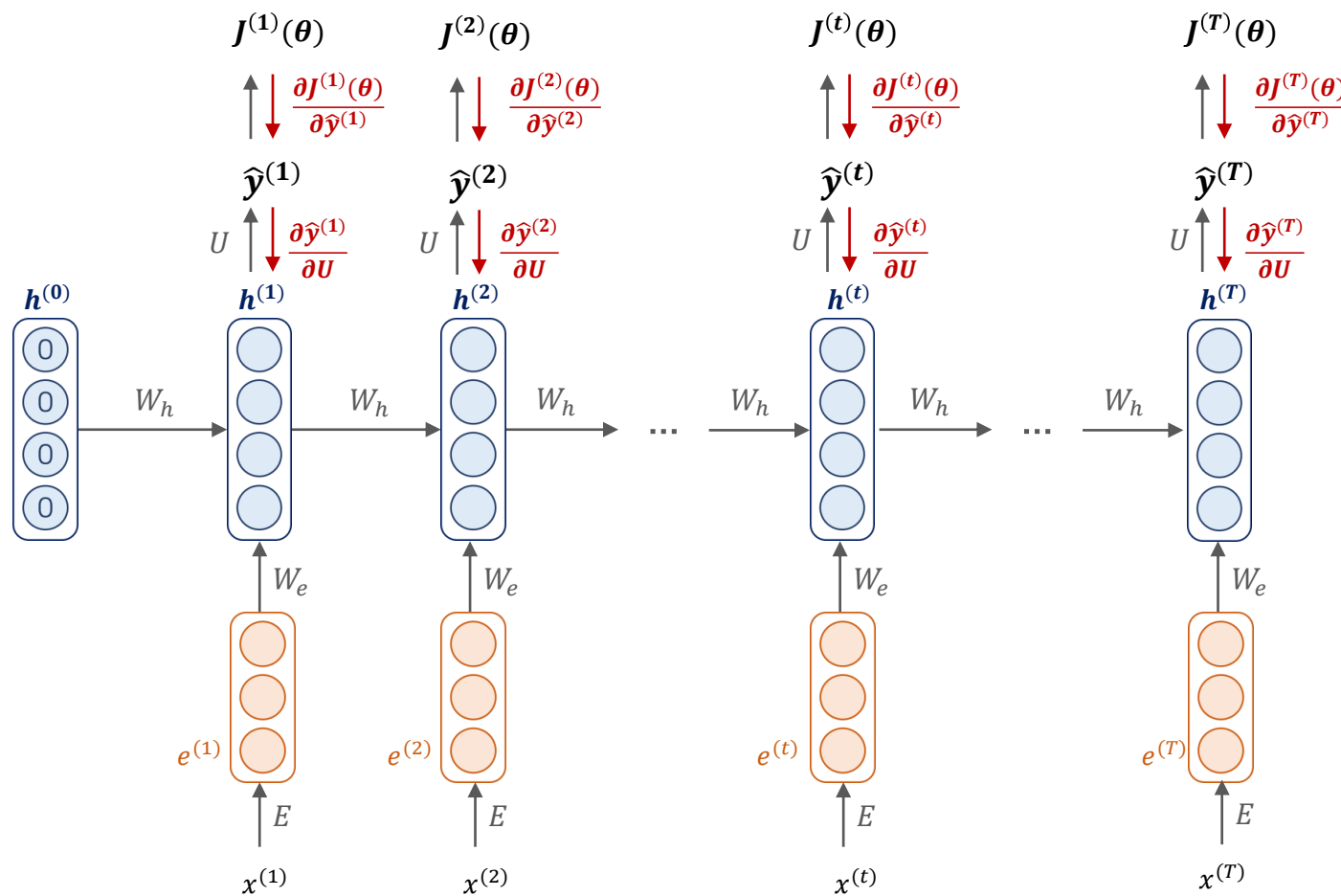
$$\textcircled{1} \frac{\partial J^{(t)}(\theta)}{\partial U} = \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial U}$$



## 02

## RNN Language Model

## Back Propagation Through Time (BPTT) ①



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

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

$$\textcircled{1} \frac{\partial J^{(t)}(\theta)}{\partial U} = \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial U}$$

$$\Rightarrow \frac{\partial J(\theta)}{\partial U} = \sum_{t=1}^T \frac{\partial J^{(t)}(\theta)}{\partial U}$$

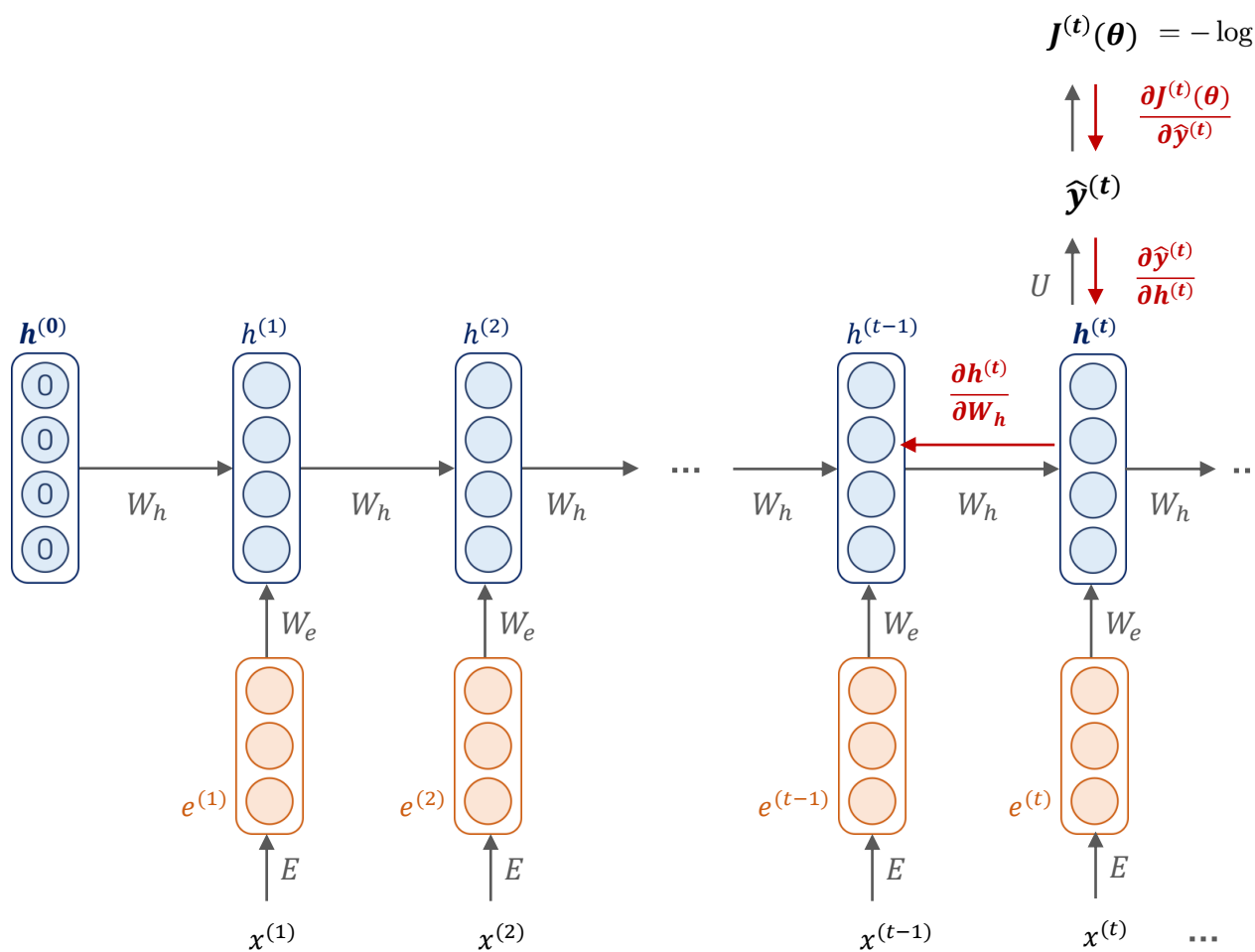
$$U^{\text{new}} = U^{\text{old}} - \alpha \frac{\partial J(\theta)}{\partial U}$$



## 02

## RNN Language Model

## Back Propagation Through Time (BPTT) ②



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

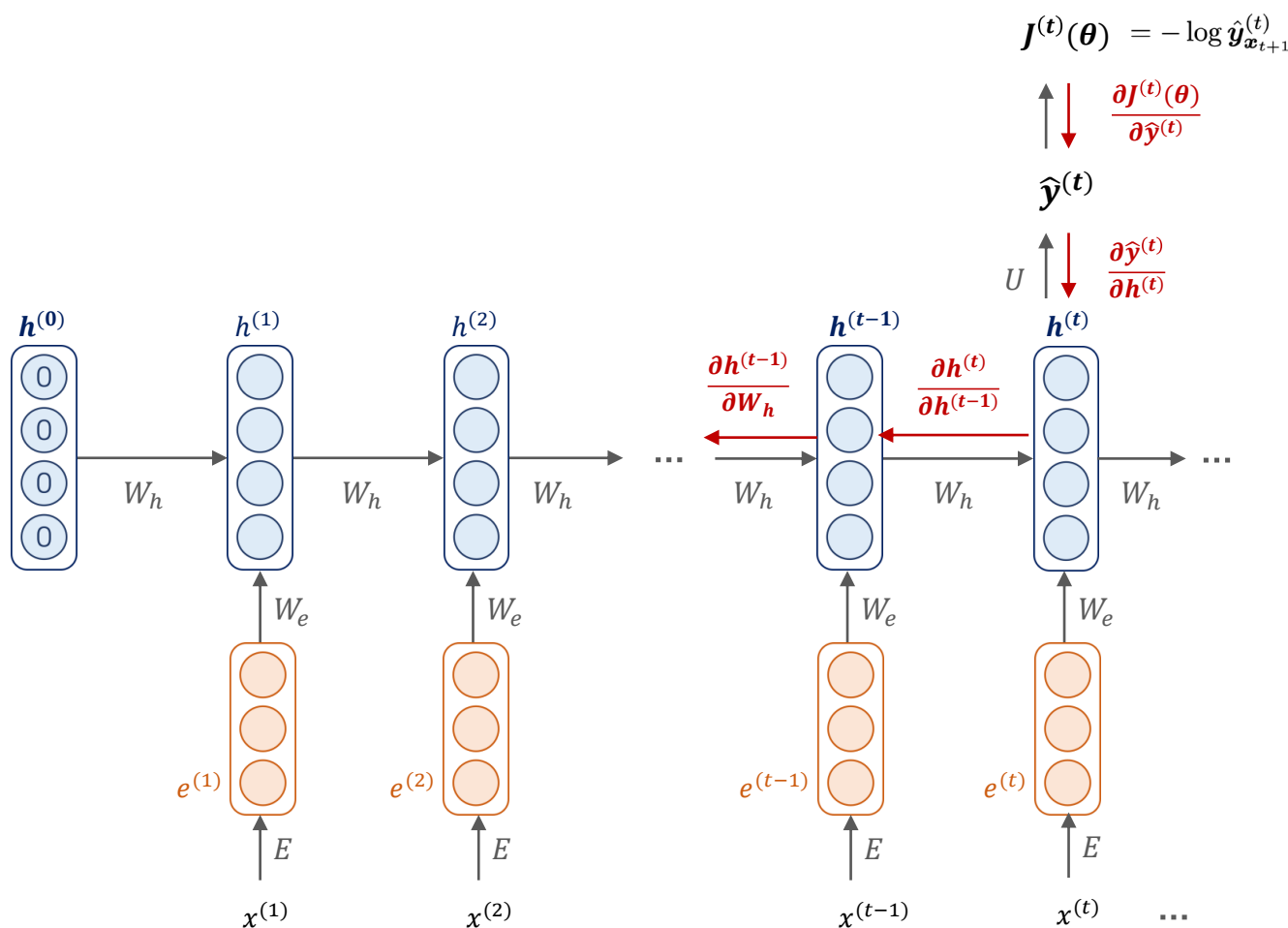
$$h^{(t)} = \sigma \left( \boxed{W_h} h^{(t-1)} + W_e e^{(t)} + b_1 \right)$$

$$\textcircled{2} \frac{\partial J^{(t)}(\theta)}{\partial W_h} = \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial W_h}$$



# RNN Language Model

## Back Propagation Through Time (BPTT) ②



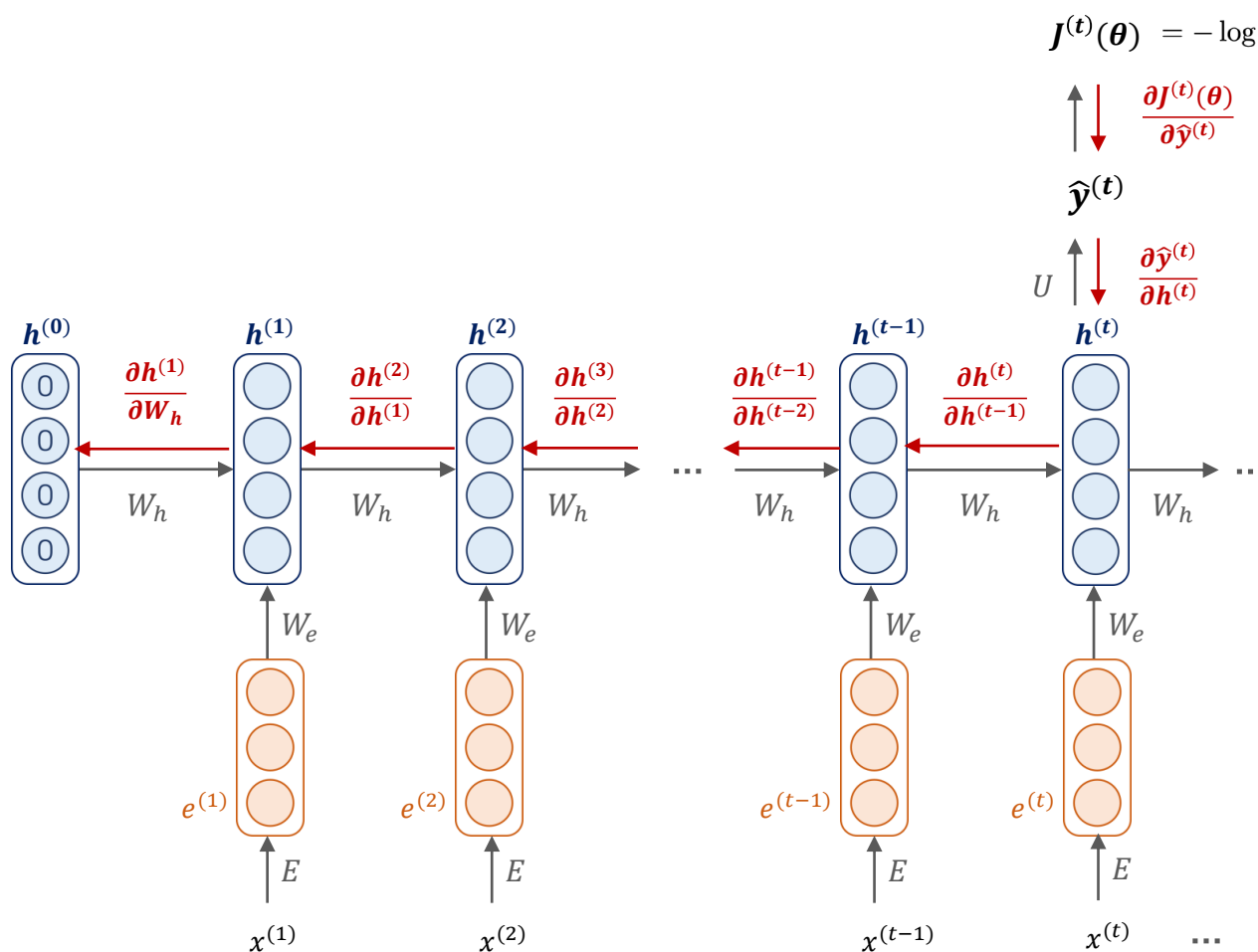
$$\hat{y}^{(t)} = \text{softmax} \left( U h^{(t)} + b_2 \right) \in \mathbb{R}^{|V|}$$

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

$$\begin{aligned} \textcircled{2} \frac{\partial J^{(t)}(\theta)}{\partial W_h} &= \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial W_h} \\ &+ \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial h^{(t-1)}} \cdot \frac{\partial h^{(t-1)}}{\partial W_h} \end{aligned}$$

# RNN Language Model

## Back Propagation Through Time (BPTT) ②



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

$$\frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}}$$

$$\hat{y}^{(t)}$$

$$\frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}}$$

$$\begin{aligned} \textcircled{2} \frac{\partial J^{(t)}(\theta)}{\partial W_h} &= \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial W_h} \\ &+ \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial h^{(t-1)}} \cdot \frac{\partial h^{(t-1)}}{\partial W_h} \end{aligned}$$

$$\vdots$$

$$+ \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial h^{(t-1)}} \cdots \frac{\partial h^{(2)}}{\partial h^{(1)}} \cdot \frac{\partial h^{(1)}}{\partial W_h}$$

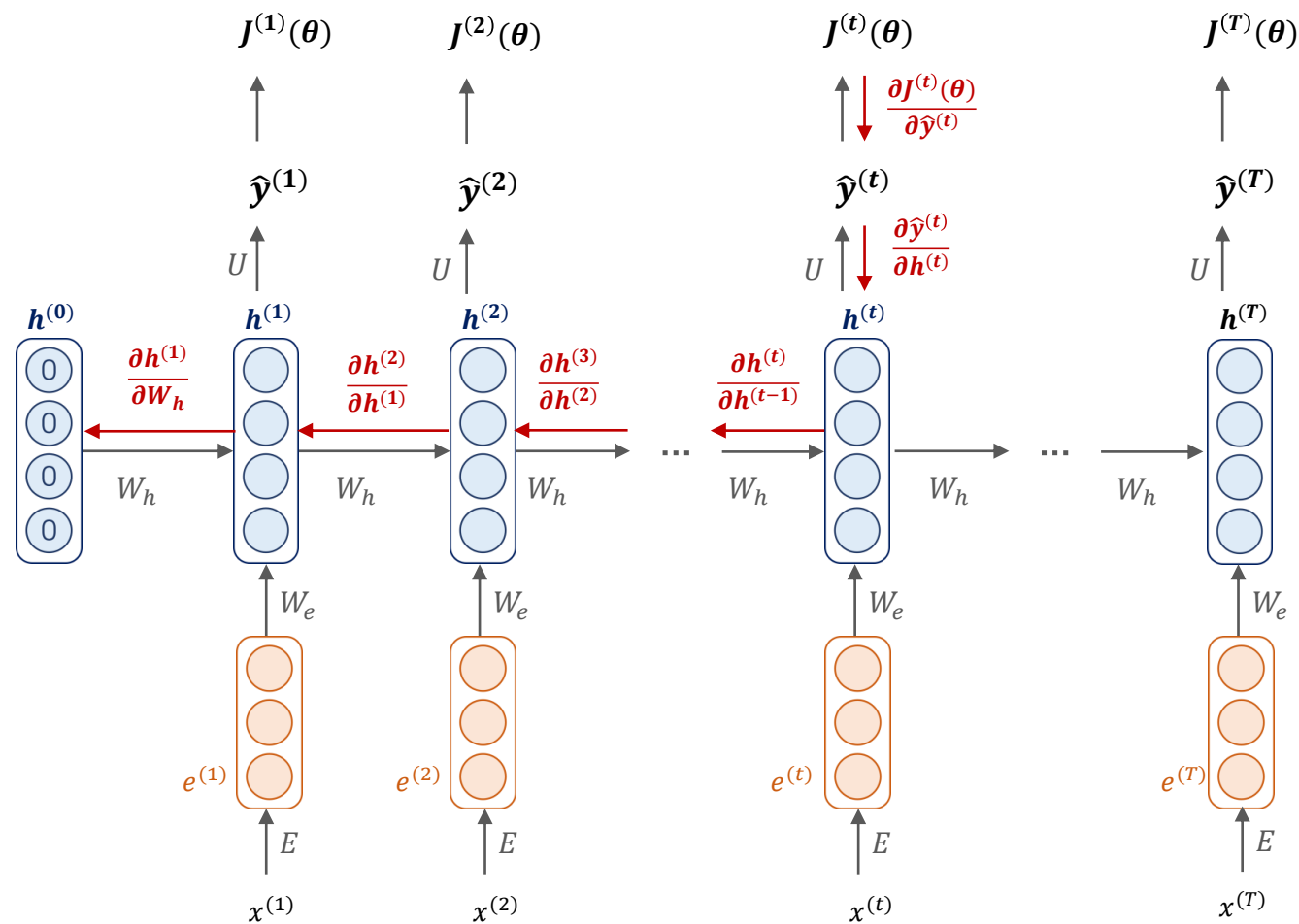
$$= \sum_{i=1}^t \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial h^{(i)}} \cdot \frac{\partial h^{(i)}}{\partial W_h}$$



## 02

## RNN Language Model

## Back Propagation Through Time (BPTT) ②



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

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

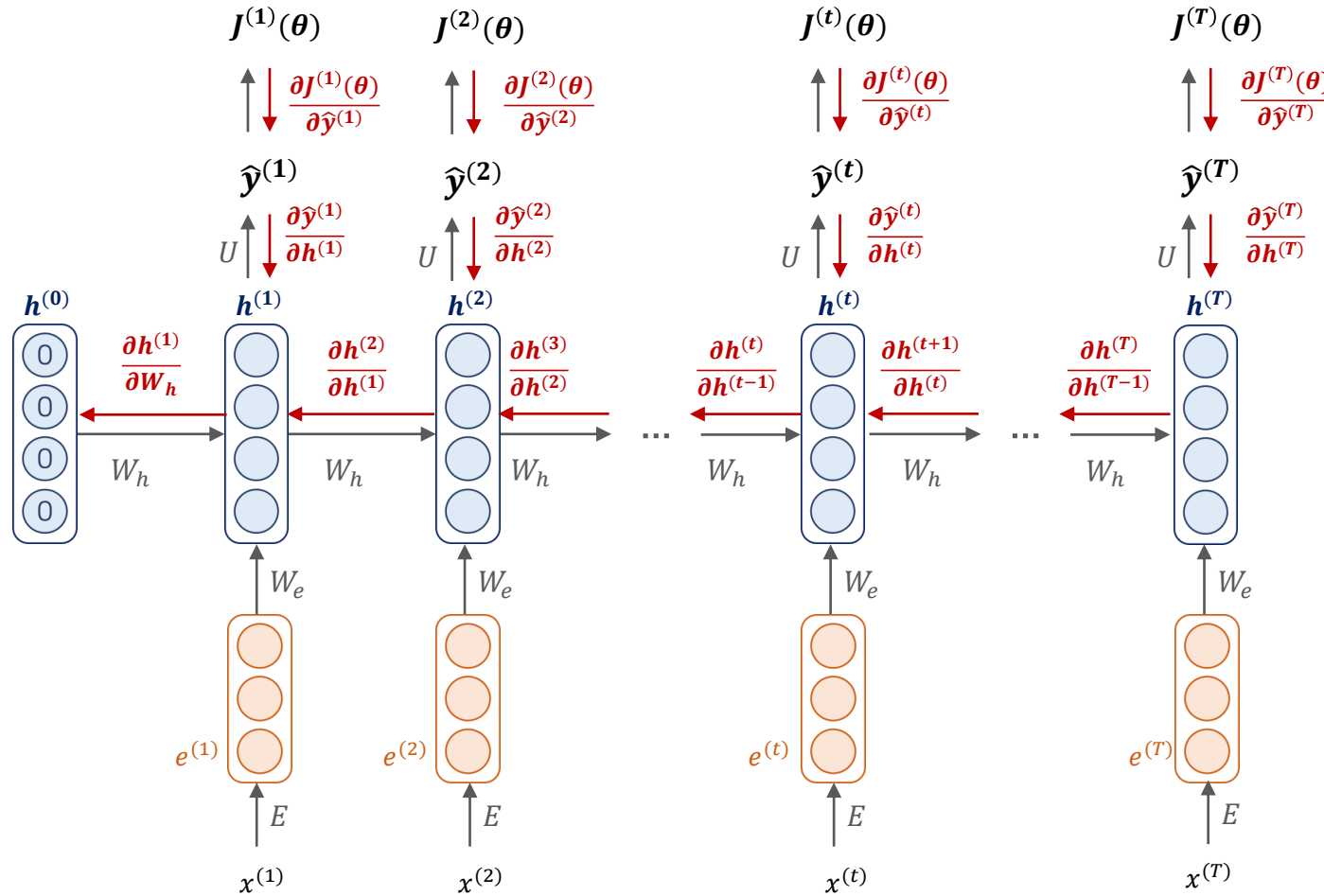
$$\textcircled{2} \frac{\partial J^{(t)}(\theta)}{\partial W_h} = \sum_{i=1}^t \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial h^{(i)}} \cdot \frac{\partial h^{(i)}}{\partial W_h}$$

$$\Rightarrow \frac{\partial J(\theta)}{\partial W_h} =$$



# RNN Language Model

## Back Propagation Through Time (BPTT) ②



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

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

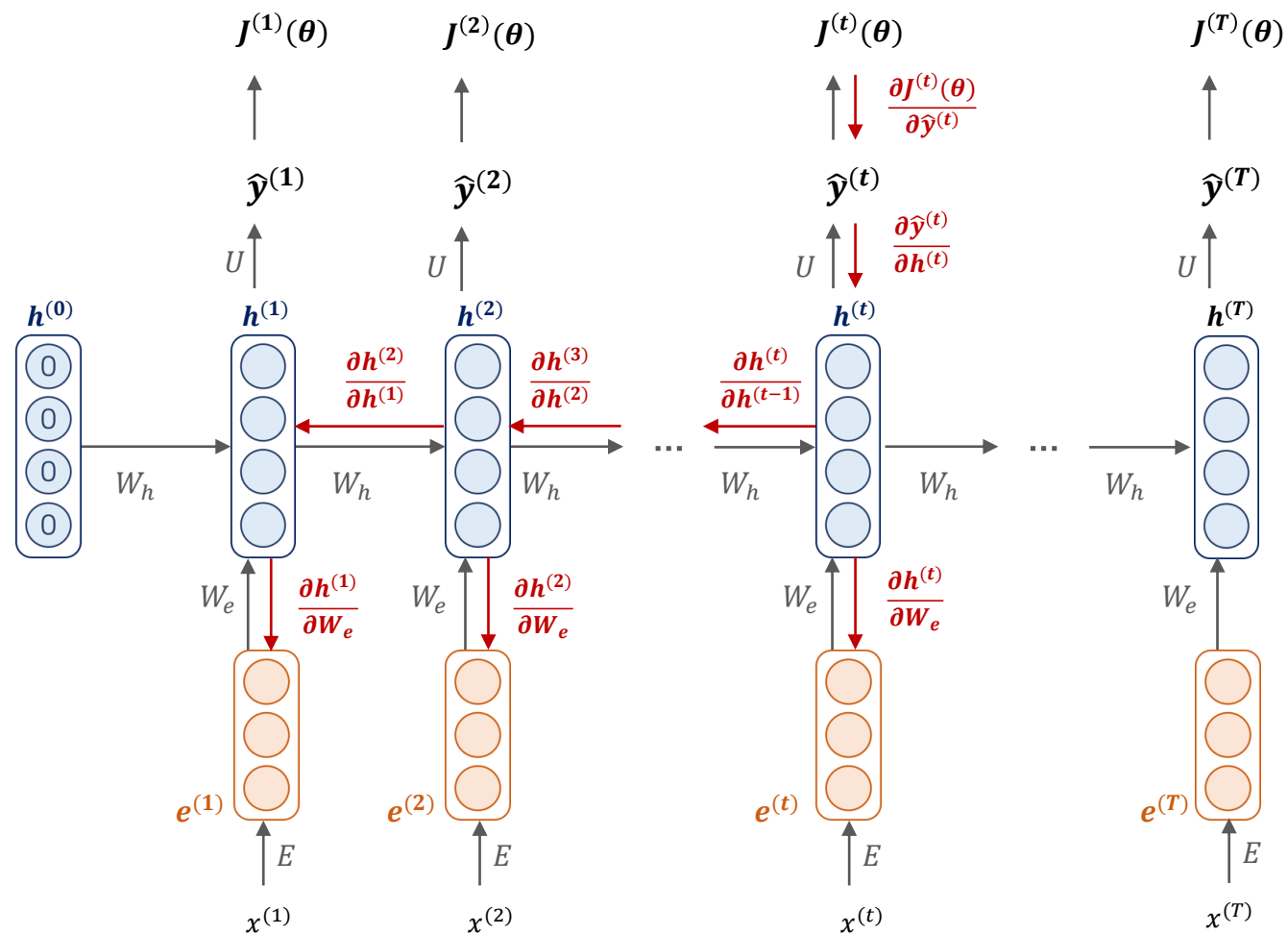
$$\textcircled{2} \frac{\partial J^{(t)}(\theta)}{\partial W_h} = \sum_{i=1}^t \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial h^{(i)}} \cdot \frac{\partial h^{(i)}}{\partial W_h}$$

$$\Rightarrow \frac{\partial J(\theta)}{\partial W_h} = \sum_{t=1}^T \frac{\partial J^{(t)}(\theta)}{\partial W_h}$$

$$W_h^{\text{new}} = W_h^{\text{old}} - \alpha \frac{\partial J(\theta)}{\partial W_h}$$

# RNN Language Model

## Back Propagation Through Time (BPTT) ③



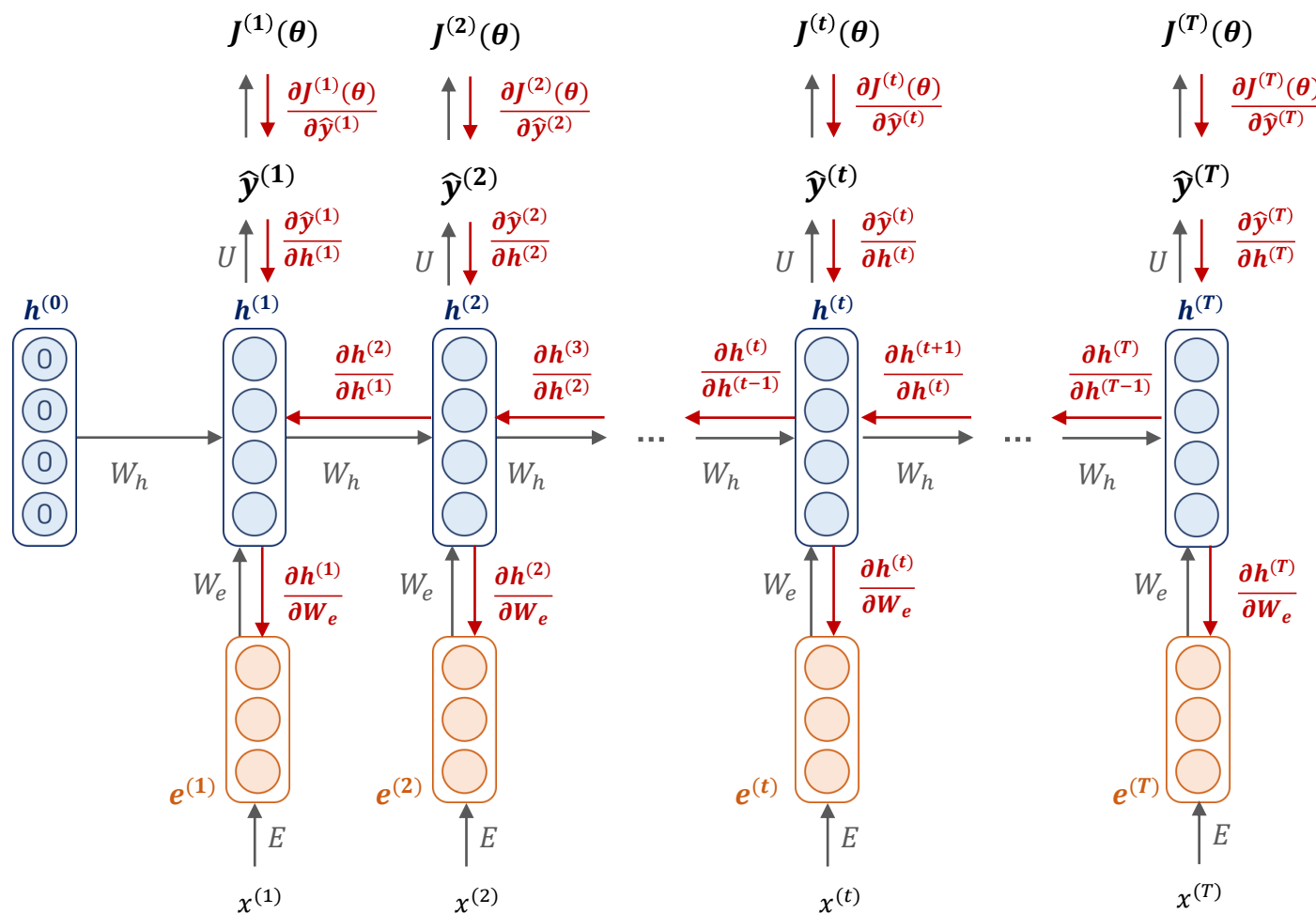
$$\hat{y}^{(t)} = \text{softmax} \left( U h^{(t)} + b_2 \right) \in \mathbb{R}^{|V|}$$

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

$$\textcircled{3} \frac{\partial J^{(t)}(\theta)}{\partial W_e} = \sum_{i=1}^t \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial h^{(i)}} \cdot \frac{\partial h^{(i)}}{\partial W_e}$$

# RNN Language Model

## Back Propagation Through Time (BPTT) ③



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

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

$$\textcircled{3} \frac{\partial J^{(t)}(\theta)}{\partial W_e} = \sum_{i=1}^t \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial h^{(i)}} \cdot \frac{\partial h^{(i)}}{\partial W_e}$$

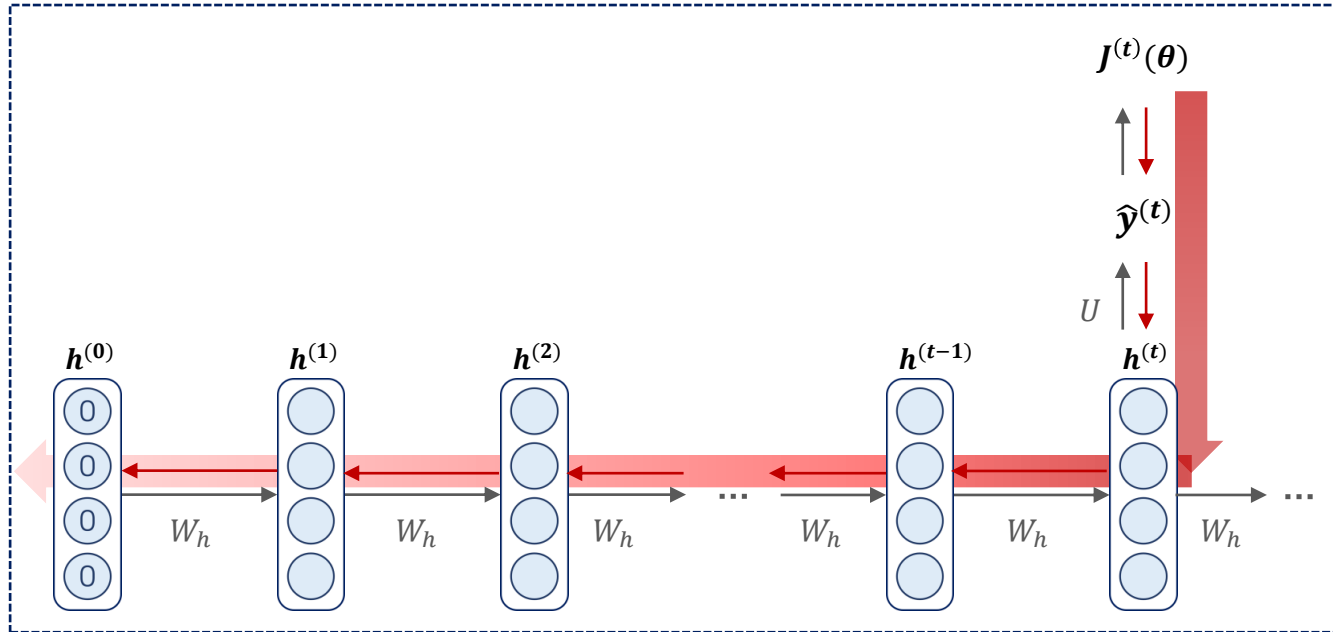
$$\Rightarrow \frac{\partial J(\theta)}{\partial W_e} = \sum_{t=1}^T \frac{\partial J^{(t)}(\theta)}{\partial W_e}$$

$$W_e^{\text{new}} = W_e^{\text{old}} - \alpha \frac{\partial J(\theta)}{\partial W_e}$$



# Problem of RNN

## ① Vanishing/Exploding Gradients



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

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

$$\frac{\partial h^{(t)}}{\partial h^{(t-1)}} = \frac{\partial h^{(t)}}{\partial s^{(t)}} \cdot \frac{\partial s^{(t)}}{\partial h^{(t-1)}} = \frac{\partial h^{(t)}}{\partial s^{(t)}} \cdot W_h$$

$$\frac{\partial J^{(t)}(\theta)}{\partial W_h} = \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial W_h} + \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \boxed{\frac{\partial h^{(t)}}{\partial h^{(t-1)}}} \cdot \frac{\partial h^{(t-1)}}{\partial W_h} \dots + \frac{\partial J^{(t)}(\theta)}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \boxed{\frac{\partial h^{(t)}}{\partial h^{(t-1)}} \dots \frac{\partial h^{(2)}}{\partial h^{(1)}}} \cdot \frac{\partial h^{(1)}}{\partial W_h}$$

동일한 가중치( $W_h$ ) 공유

- $W_h$ 가 작을 수록( $< 1$ ) 반복적으로 곱해지는 값이 0에 가까워져 gradient vanishing
- $W_h$ 가 클 수록( $> 1$ ) 반복적으로 곱해지는 값이 기하급수적으로 커져 gradient exploding

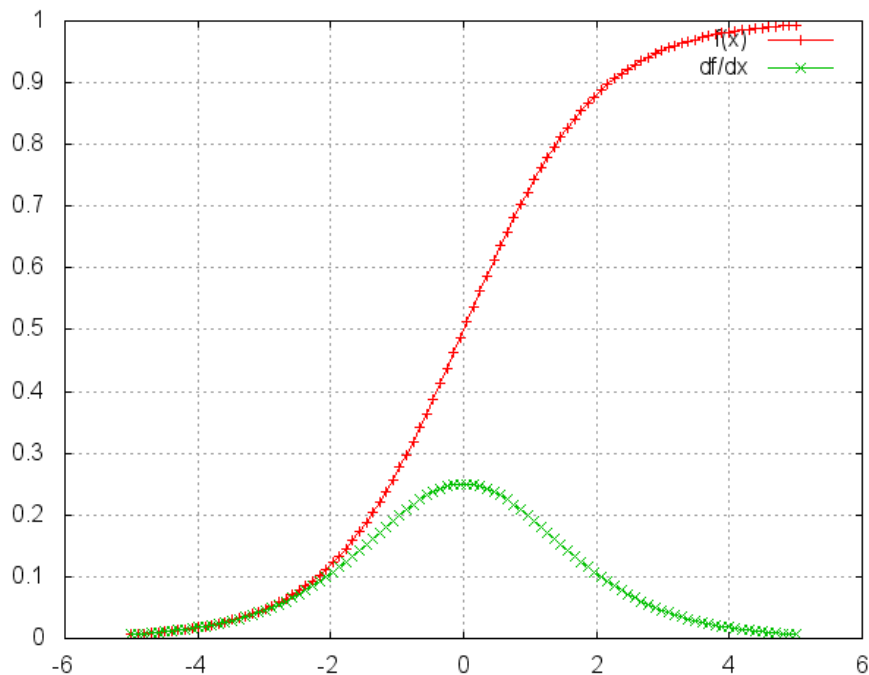
# Problem of RNN

Sigmoid vs. Tanh

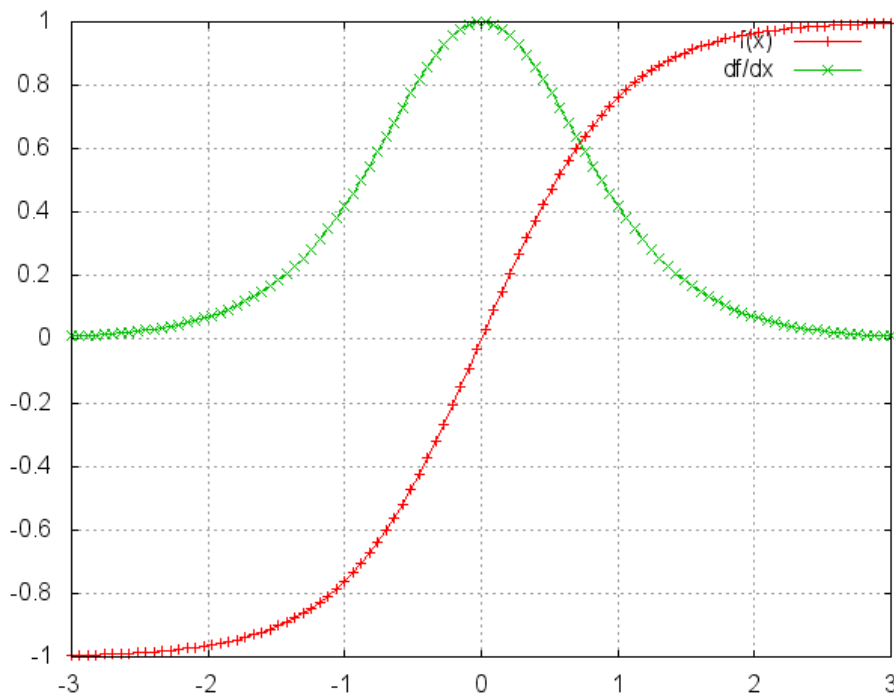
## ③ Hidden States

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

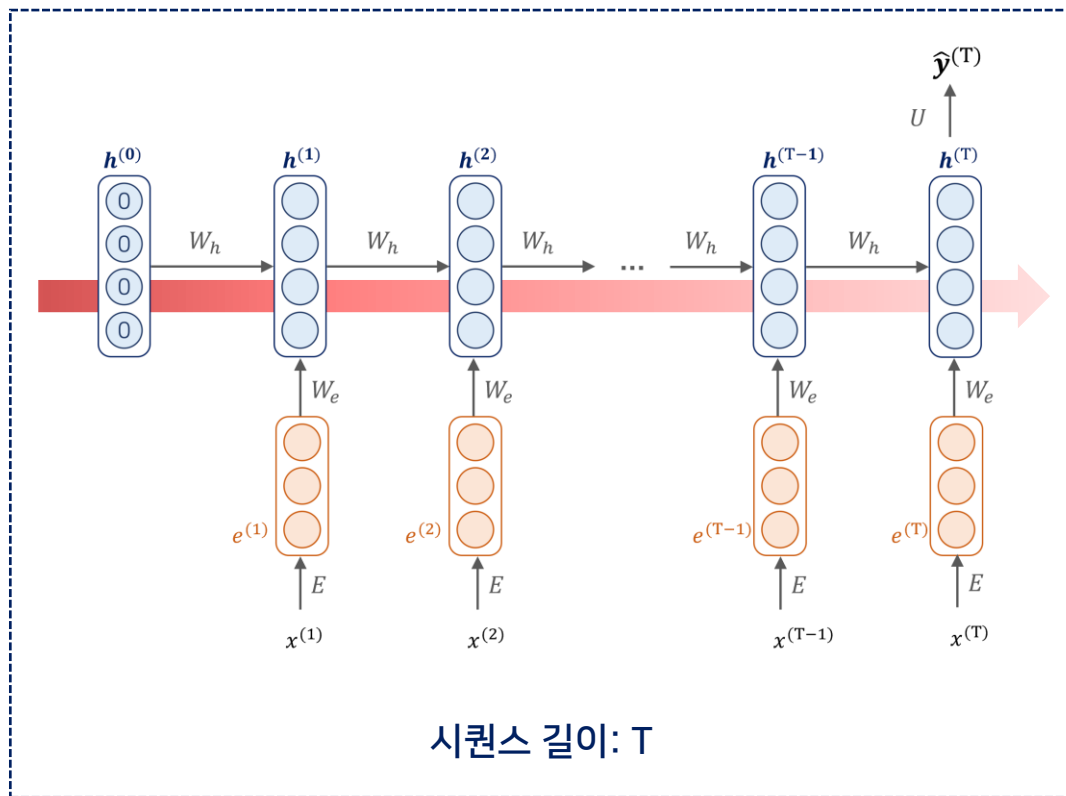
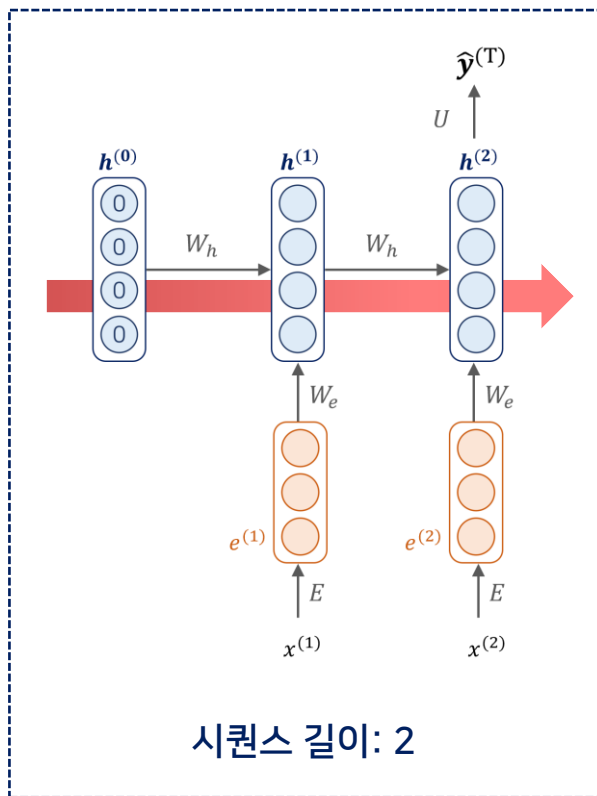
Sigmoid



Tanh



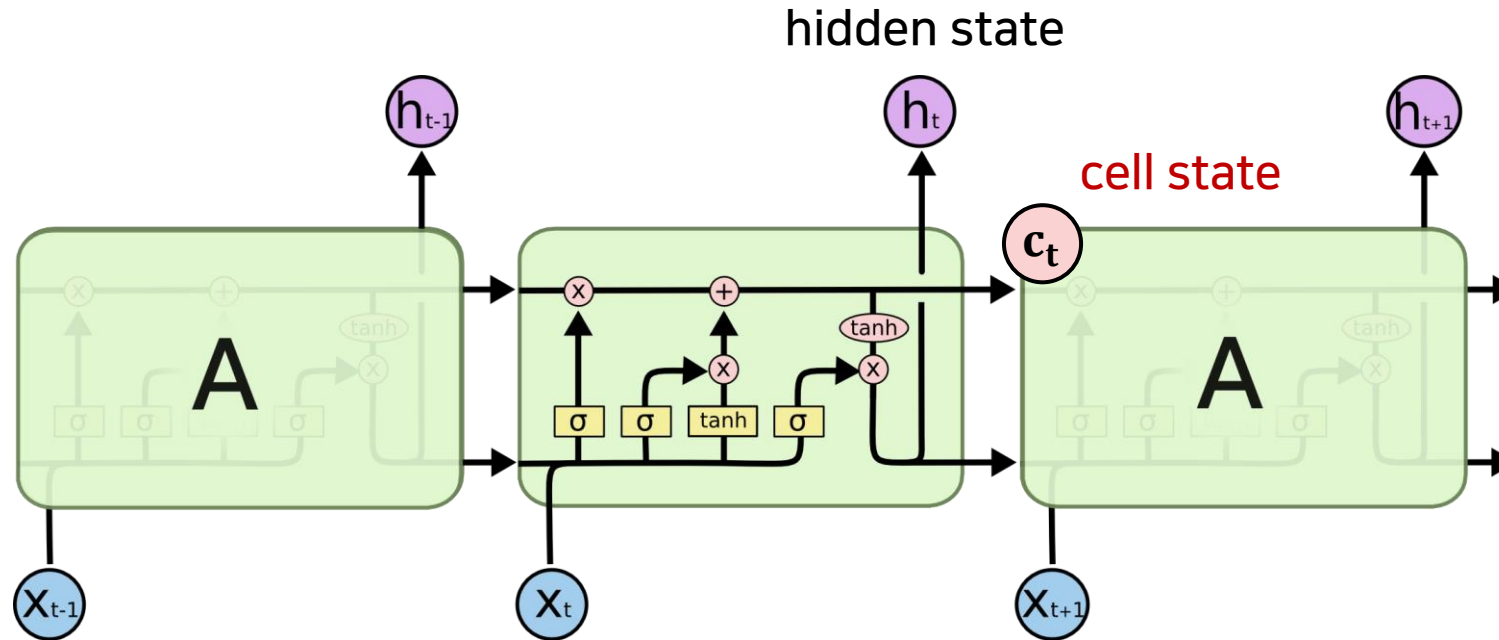
- ✓ Sigmoid: 기울기 0에서 약 0.25 사이
- ✓ Tanh: 기울기 0에서 1 사이 -> gradient vanishing에 더 강함



- input sequence 길이가 긴 경우, 시퀀스 초반의 정보가 후반 time step의 hidden state를 도출하는데까지 전달되기 어려움 -> Long term dependency를 잘 반영하지 못함

# LSTM: Long Short-Term Memory RNN

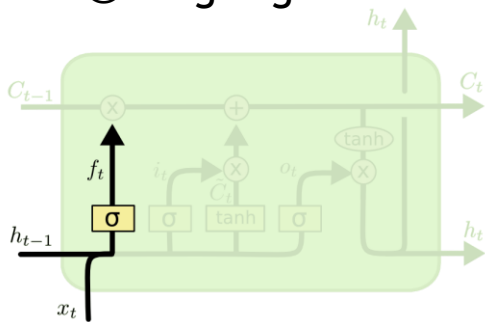
RNN with Separate Memory



- ① hidden state를 통해 short term memory를 조절하고 cell state를 통해 long term memory 보존
- ② forget, input, output, 3개의 gate를 통해 매 time step의 cell state와 hidden state, input에서 취할 정보의 양 결정



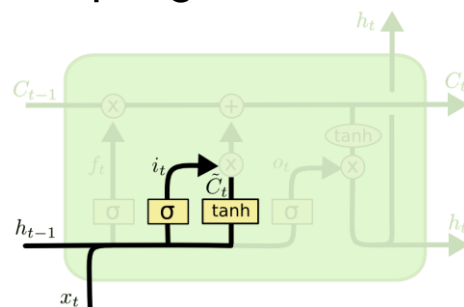
## ① forget gate



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

- 입력 정보(새로운 입력 시퀀스와 이전 시점의 hidden state)에 시그모이드를 취함

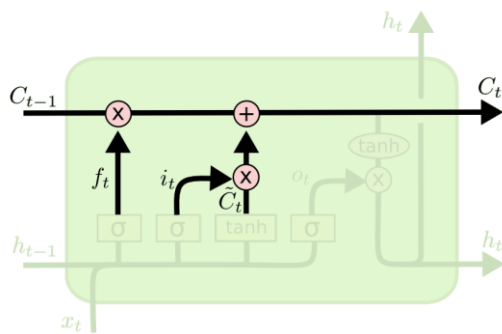
## ② input gate, new cell content



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

- input gate: 입력 정보에 시그모이드를 취함
- $\tilde{C}_t = \tanh(W_c h^{(t-1)} + U_c x^{(t)} + b_c)$
- new cell content: 입력 정보에 tanh를 취함

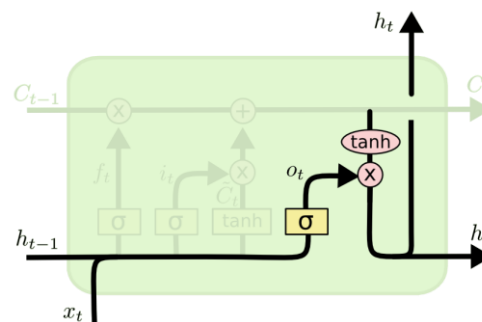
## ③ Cell state



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

- forget gate: 이전 시점의 cell state에서 어느 정도의 정보를 가져갈 것인지 결정
- input gate: 입력 정보에서 장기 기억으로 가져갈 정보의 양 결정
- element wise product

## ④ Output gate, hidden state

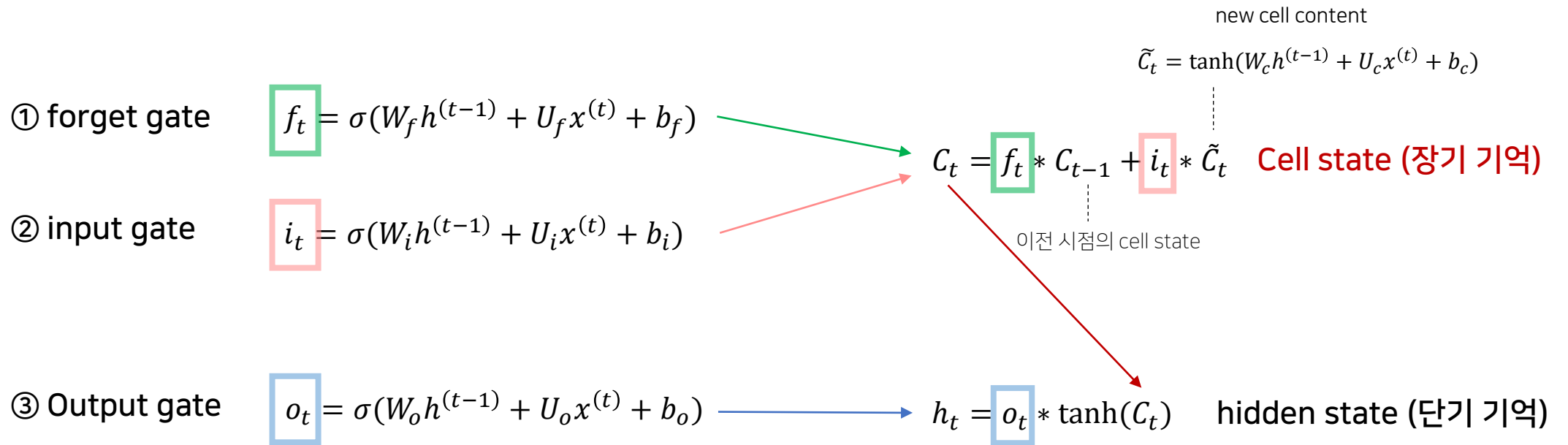


$$o_t = \sigma(W_o h^{(t-1)} + U_o x^{(t)} + b_o)$$

- output gate: 입력 정보에 시그모이드를 취함
- $h_t = o_t * \tanh(C_t)$
- hidden state: 현재 입력과 대비해서 장기 기억에서 어느 정도의 정보를 단기 기억으로 사용할지 결정

# LSTM: Long Short-Term Memory RNN

## Architecture



- ✓ **Language Model** : 확률 분포를 기반으로 주어진 문맥(sequence) 이후에 위치할 단어 예측
  - N-gram LM: 이전에 등장하는 N개의 단어 chunk를 바탕으로 다음에 올 단어 예측 (sparsity problem)
  - Neural LM: 이전에 등장하는 window size만큼의 단어를 바탕으로 다음에 올 단어 예측
  - input length 에 제한 -> 모든 문맥 고려 불가능
  
- ✓ **RNN**
  - input length 에 제한이 없는 sequential neural network
  - gradient vanishing /exploding
  - long term dependency problem
  
- ✓ **LSTM**
  - cell state를 통해 long term memory 보존
  - 3개의 gate를 통해 정보의 flow 조절

**감사합니다**