

# NLP Study Week 1

# Sequence

- sequence는 시계열적인 특성을 가짐
  - > time에 따라 변해가는 양상을 모델에 표현할 수 있음
- 자연어 처리도 마찬가지로 단어의 순서에 따라 영향을 미치기 때문에 적용할 필요 있음
- 이미지 모델에서는 sequence를 어떻게 적용할까?
  - > 동영상과 같은 움직임이 들어가 있는 이미지의 경우 적용 가능
- 시퀀스 모델의 타입은 3가지로 볼수 있음
  - > one to sequence
  - > sequence to one
  - > sequence to sequence

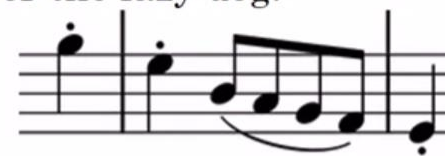
# Examples of sequence data

Speech recognition



“The quick brown fox jumped over the lazy dog.”

Music generation



Sentiment classification

“There is nothing to like in this movie.”



DNA sequence analysis

AGCCCCTGTGAGGAACTAG



AGCCCCTGTGAGGAACTAG

Machine translation

Voulez-vous chanter avec moi?



Do you want to sing with me?

Video activity recognition



Running

Name entity recognition

Yesterday, Harry Potter met Hermione Granger.



Yesterday, **Harry Potter** met **Hermione Granger**.  
Andrew Ng

# Notation

X : Harry Potter and Harmione Granger invented a new spell  
 $x^{<1>}$     $x^{<2>}$    .....    $x^{<9>}$

Y :   1        1        0        1        1        .....        0  
 $y^{<1>}$     $y^{<2>}$    .....    $y^{<9>}$

$X^{(i)<t>}$       $Y^{(i)<t>}$  : i번째 문장 t번째 단어

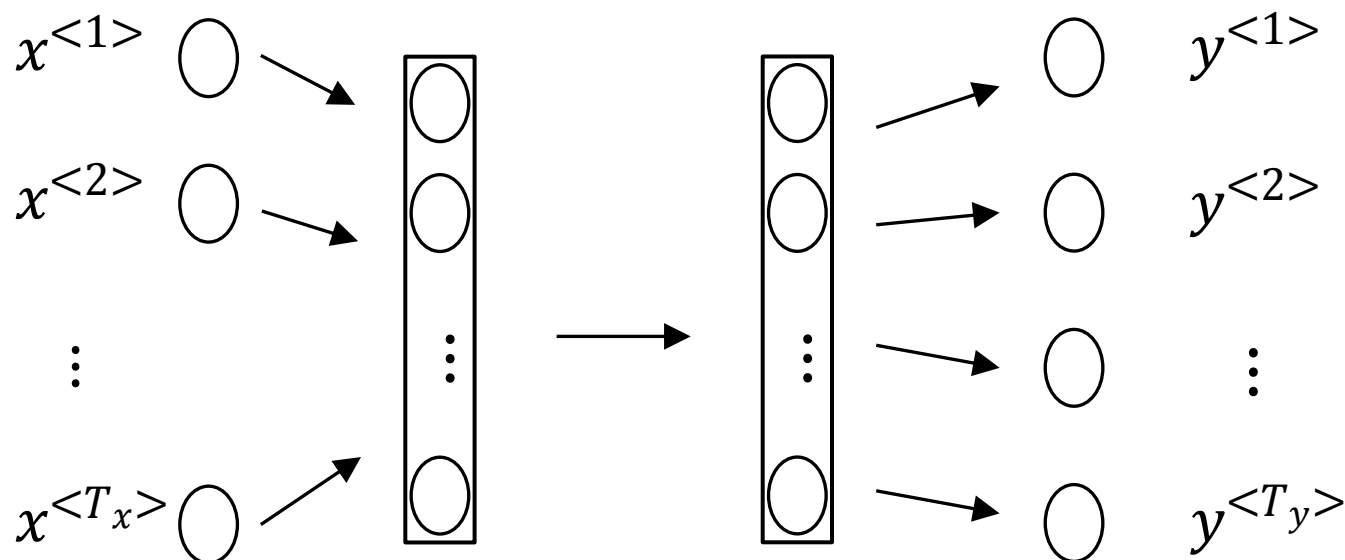
$T_x^{(i)} = 9$       $T_y^{(i)} = 9$  : 갯수

# Notation

x: Harry Potter and Hermione Granger invented a new spell.

$x^{<1>}$     $x^{<2>}$     $x^{<3>}$     $\dots$     $x^{<9>}$

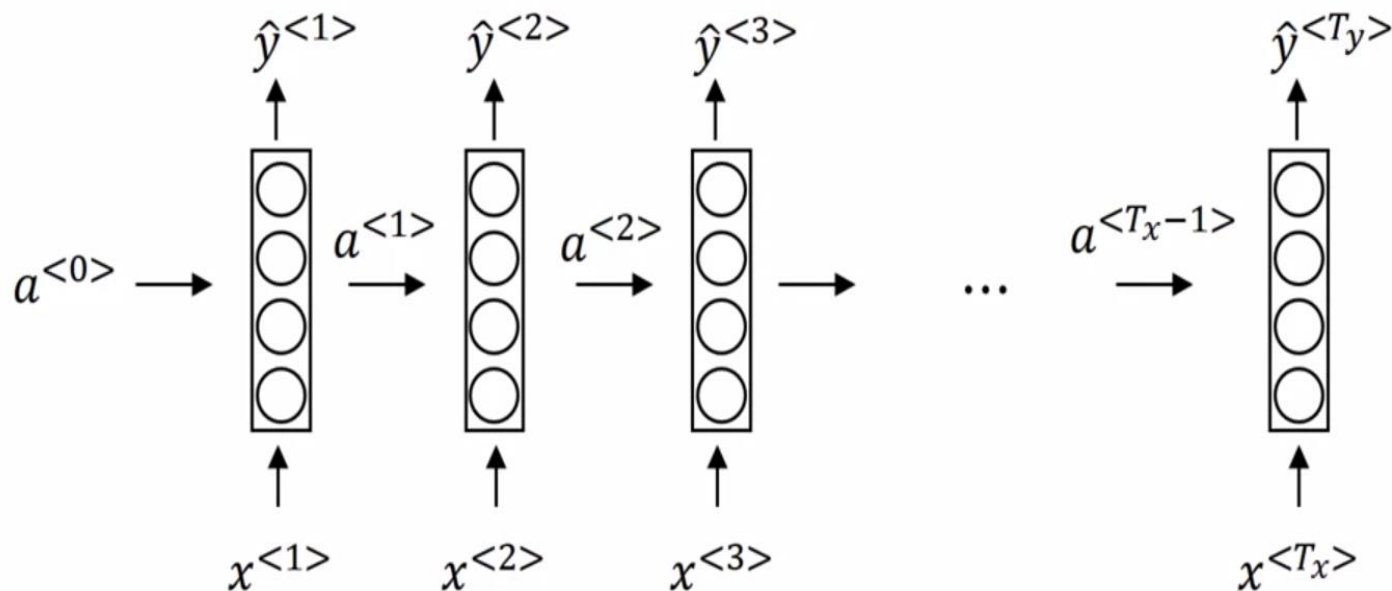
# Why not a standard network?



## Problems:

- Inputs, outputs can be different lengths in different examples.
- Doesn't share features learned across different positions of text.
- Too many parameters

# RNN



$$a^{<0>} = \vec{0}$$

$$\underline{a^{<1>}} = g_1(W_{aa} a^{<0>} + \underline{W_{ax}} x^{<1>} + b_a) \leftarrow \tanh / \text{Relu}$$

$$\underline{\hat{y}^{<1>}} = g_2(\underline{W_{ya}} \underline{a^{<1>}} + b_y) \leftarrow \text{sigmoid}$$

$$\boxed{\begin{aligned} a^{<t>} &= g(W_{aa} a^{<t-1>} + W_{ax} x^{<t>} + b_a) \\ \hat{y}^{<t>} &= a(W_{ya} a^{<t>} + b_y) \end{aligned}}$$

# RNN

$$a^{<t>} = g(W_{aa}a^{<t>} + W_{ax}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(W_{ya}a^{<t>} + b_y)$$



$$\hat{y}^{<t>} = g(W_y a^{<t>} + b_y)$$

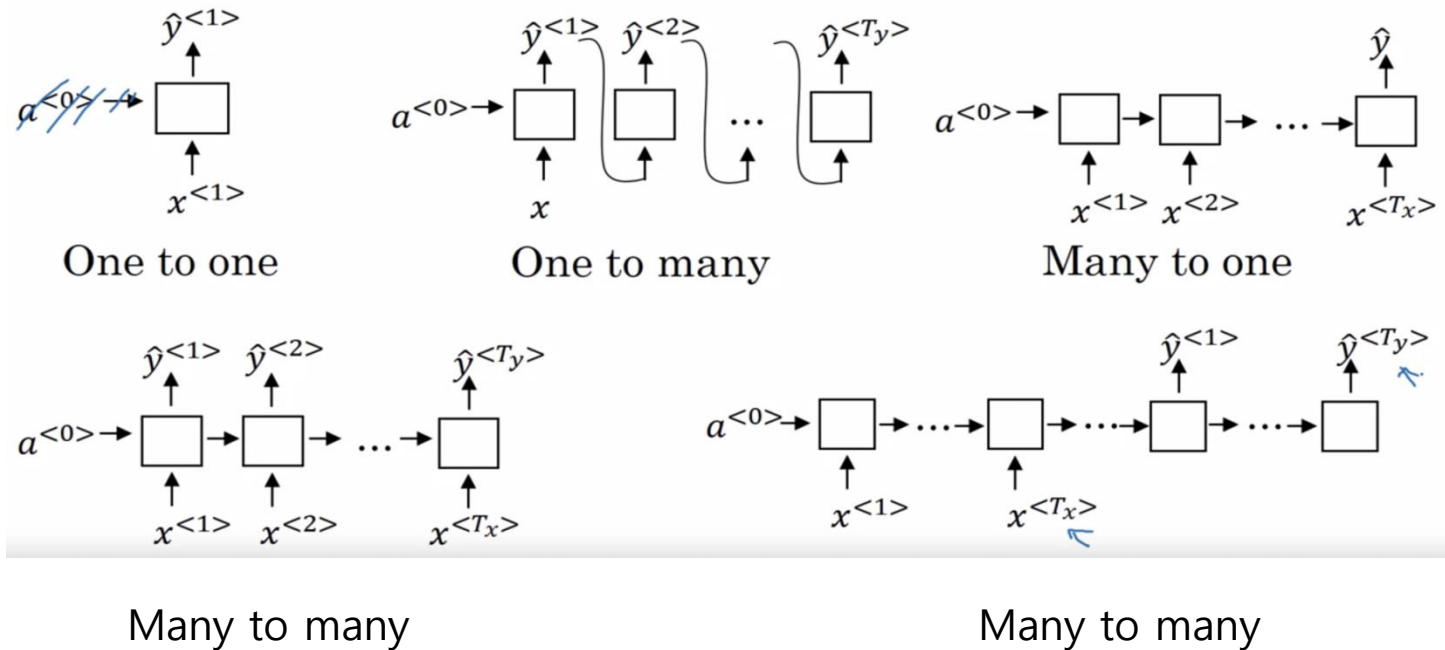
$$a^{<t>} = g(W_a [a^{<t-1>}, x^{<t>}] + b_a)$$

$$\overset{(100)}{\updownarrow} \left[ \overset{(100)}{\underbrace{W_{aa}}}, \overset{(10000)}{\underbrace{W_{ax}}} \right] = W_a \quad (100, 10100)$$

$$[a^{<t-1>}, x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} \quad \begin{matrix} \updownarrow 100 \\ \updownarrow 10000 \end{matrix} \quad \updownarrow 10100$$



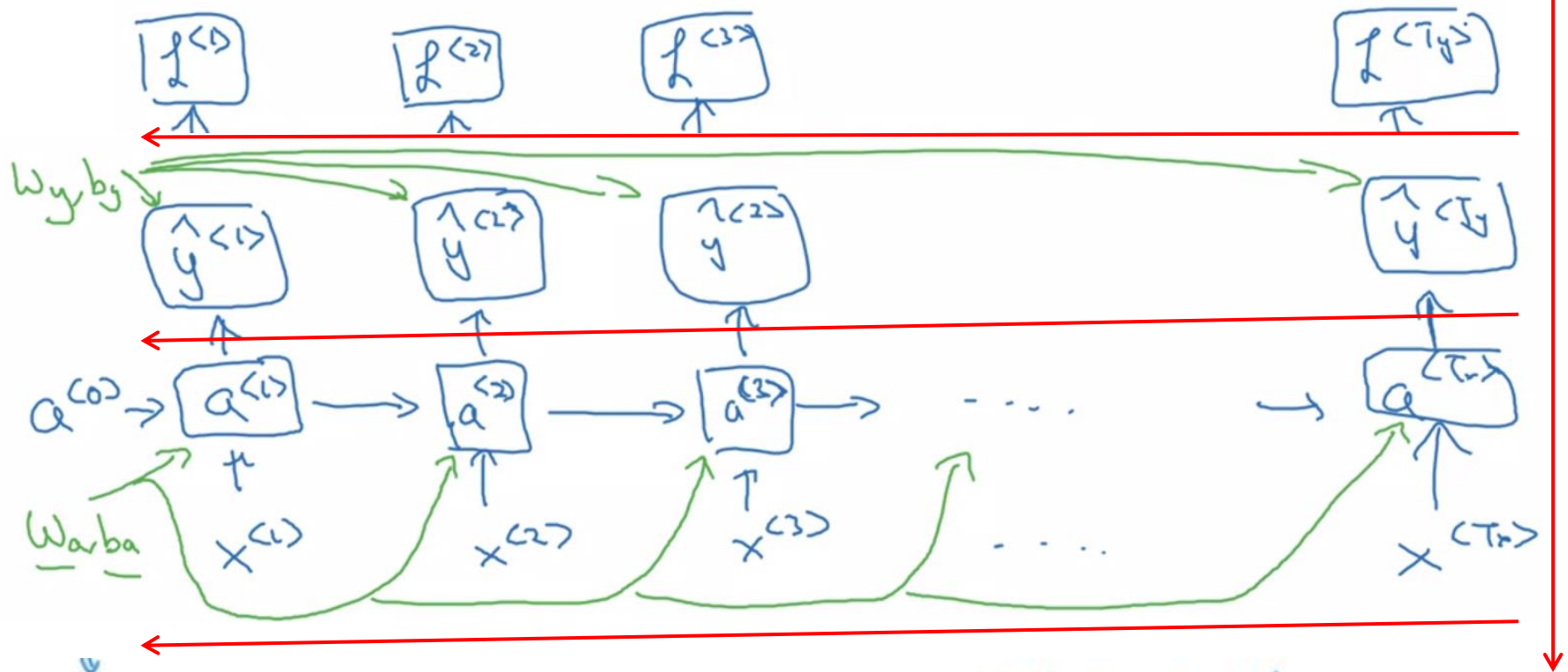
# RNN



- one to one : 일반적인 신경 네트워크
- one to many : (음악 생성), (언어모델링), (이미지캡션),
- many to one : 동영상 및 텍스트 입력하여 감성분석
- many to many : 이름 엔터티 인식
- many to many (encoder-decoder) : 기계번역

# Backpropagation

$$\mathcal{L}(\hat{y}, y)$$



$$\mathcal{L}^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -y^{(t)} \log \hat{y}^{(t)} - (1 - y^{(t)}) \log (1 - \hat{y}^{(t)})$$

$$\mathcal{L}(\hat{y}, y) = \sum_{t=1}^{T_y} \mathcal{L}^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

# Language Model with RNN

Speech recognition

The apple and pair salad.

→ The apple and pear salad.

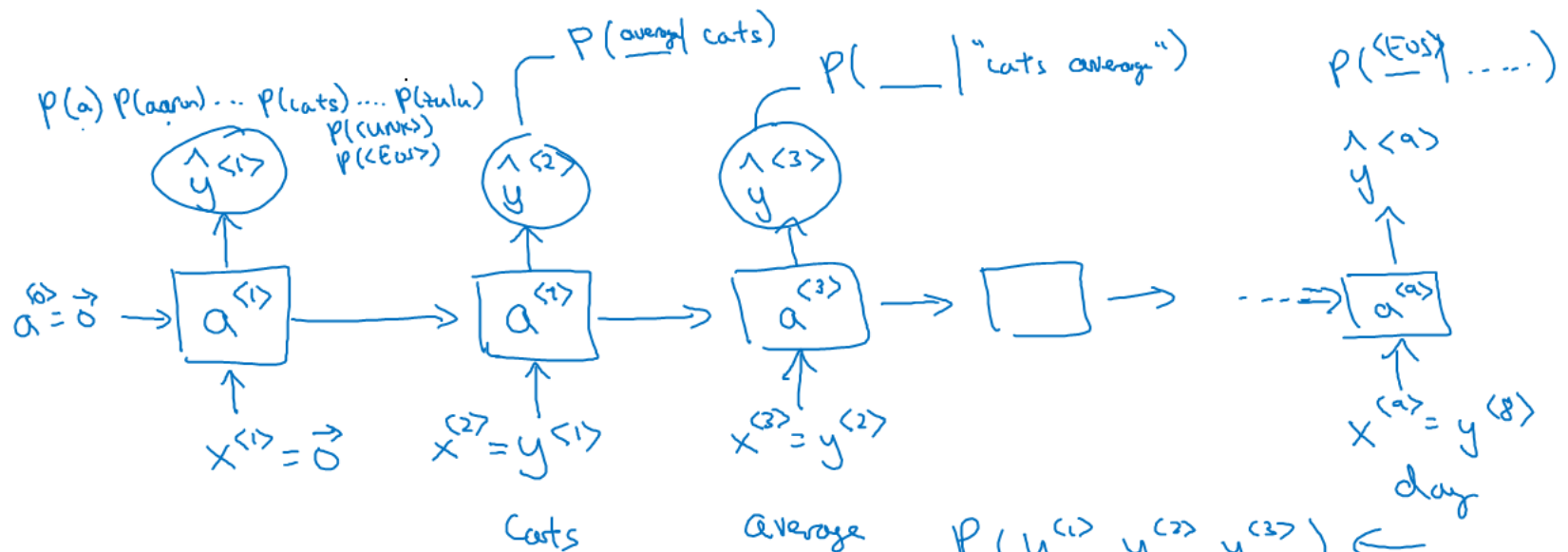
$$P(\text{The apple and pair salad}) = 3.2 \times 10^{-13}$$

$$P(\text{The apple and pear salad}) = 5.7 \times 10^{-10}$$

$$P(\text{Sentence}) = ?$$

$$P(y^{(1)}, y^{(2)}, \dots, y^{(T)})$$

# Language Model with RNN



→ Cats average 15 hours of sleep a day. <EOS>

$$\mathcal{L}(\hat{y}^{<t>}, y^{<t>}) = - \sum_i y_i^{<t>} \log \hat{y}_i^{<t>} \quad \leftarrow$$

$$\mathcal{L} = \sum_t \mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

$$\begin{aligned} P(y^{(1)}, y^{(2)}, y^{(3)}) &\leftarrow \\ &= \frac{P(y^{(1)}) P(y^{(2)} | y^{(1)})}{P(y^{(3)} | y^{(1)}, y^{(2)})} \end{aligned}$$

# Language Model with RNN

## Character level language model

- 장점

"알 수 없는 단어토큰"에 대해 걱정할 필요 없음

- 단점

훨씬 더 긴 배열로 모델이 끝난다.

=> 문장의 초기 부분이 문장의 뒷부분에도 어떻게 영향을 주는지,

long range dependency 캡처할 때 word level만큼 좋지 않음

computationally expensive

Word level Vocabulary = [a, aaron, ..., zulu, <UNK>]

Character level Vocabulary = [a, b, ... , A, B, ..., 0, ... 9, ...]

# Problems of RNN

## Vanishing Gradient

- Long range dependency capturing 약하다
- Layer 많이 쌓기 힘들다

⇒ GRU, LSTM 등으로 해결

## Exploding Gradient

- Parameter 엉망
- NaN 등, 숫자 아닌게 나올수도

⇒ Gradient clipping 통해 한계점보다 큰 벡터 조정

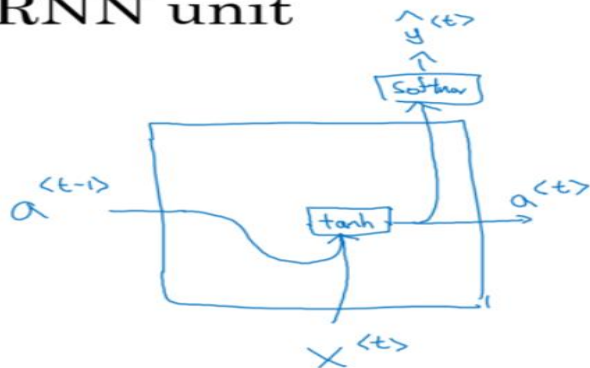
Vanishing이 더 해결하기 어렵고 중요한 문제

# GRU

The **cat**, which already ate ..., **was** full.

The **cats**, which already ate ..., **were** full.

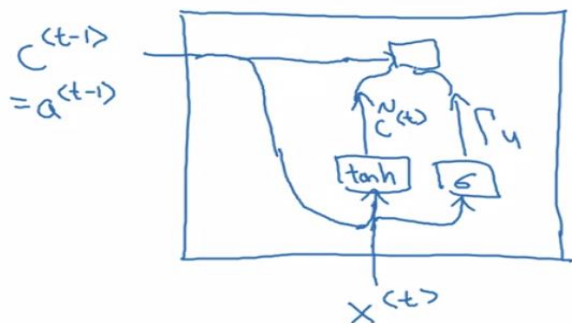
## RNN unit



$$\underline{a}^{<t>} = \underline{g}(\underline{W}_a[\underline{a}^{<t-1>}, x^{<t>}] + \underline{b}_a)$$

(Note: A handwritten 'tanh' with an arrow points to the 'g' function in the equation.)

## GRU (simplified)



$c$  = memory cell

$$\rightarrow \underline{c}^{<t>} = \underline{a}^{<t>}$$

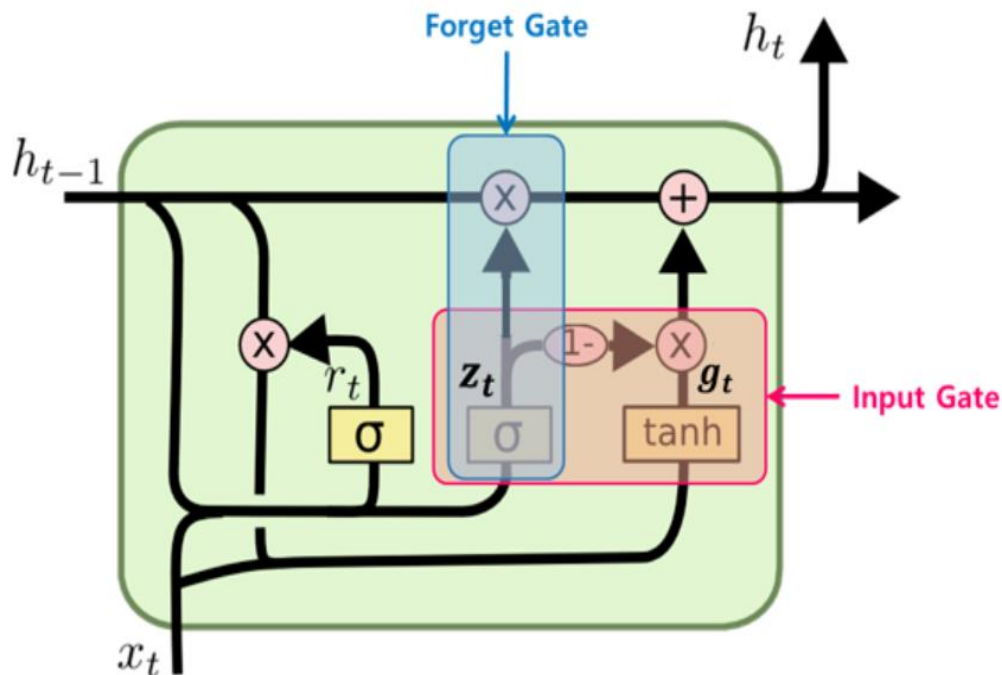
$$\rightarrow \tilde{c}^{<t>} = \tanh(W_c[c^{<t-1>}, x^{<t>}] + b_c)$$

$$\rightarrow \Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

(Note: A handwritten 'update' with an arrow points to the  $\Gamma_u$  term.)

$$\underline{c}^{<t>} = \underline{\Gamma_u} * \underline{\tilde{c}^{<t>}} + (1 - \underline{\Gamma_u}) * \underline{c^{<t-1>}}$$

# GRU



$$r_t = \sigma(W_{xr}^T \cdot x_t + W_{hr}^T \cdot h_{t-1} + b_r)$$

$$z_t = \sigma(W_{xz}^T \cdot x_t + W_{hz}^T \cdot h_{t-1} + b_z)$$

$$g_t = \tanh(W_{xg}^T \cdot x_t + W_{hg}^T \cdot (r_t \otimes h_{t-1}) + b_g)$$

$$h_t = z_t \otimes h_{t-1} + (1 - z_t) \otimes g_t$$

- 하나의 gate controller인  $z$ 가 **forget**과 **input** 게이트(gate)를 모두 제어한다.  
 $z$ 가 1을 출력하면 forget 게이트가 열리고 input 게이트가 닫히며,  
 $z$ 가 0일 경우 반대로 forget 게이트가 닫히고 input 게이트가 열린다.  
 즉, 이전( $t-1$ )의 기억이 저장 될때 마다 타임 스텝 의 입력은 삭제된다.
- GRU 셀은 output 게이트가 없어 전체 상태 벡터  $h$ 가 타임 스텝마다 출력되며,  
 이전 상태  $h$ 의 어느 부분이 출력될지 제어하는 새로운 gate controller인  $r$ 이 있다.



# LSTM

## GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$a^{<t>} = c^{<t>}$$

## LSTM

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

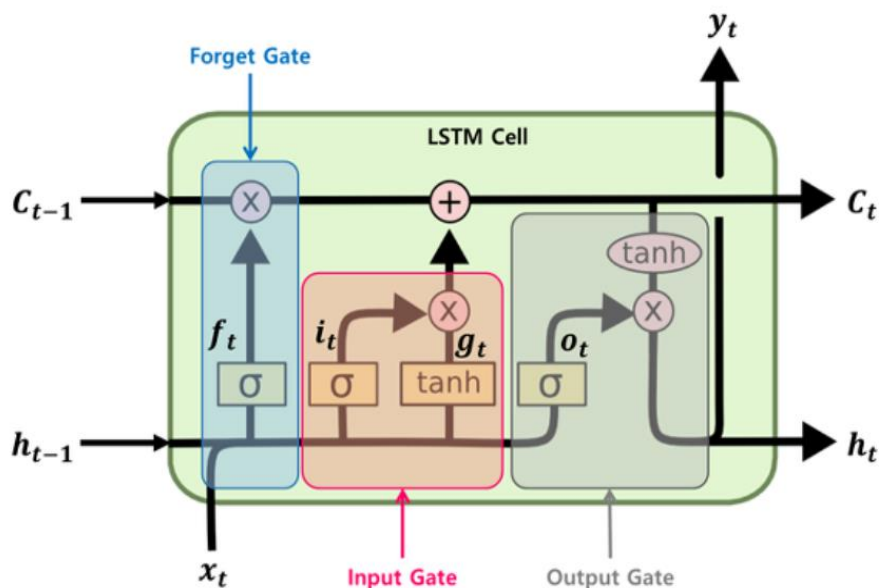
$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

$$\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \tanh(c^{<t>})$$

# LSTM



$$\begin{aligned}
 f_t &= \sigma(W_{xf}^T \cdot x_t + W_{hf}^T \cdot h_{t-1} + b_f) \\
 i_t &= \sigma(W_{xi}^T \cdot x_t + W_{hi}^T \cdot h_{t-1} + b_i) \\
 o_t &= \sigma(W_{xo}^T \cdot x_t + W_{ho}^T \cdot h_{t-1} + b_o) \\
 g_t &= \tanh(W_{xg}^T \cdot x_t + W_{hg}^T \cdot h_{t-1} + b_g) \\
 c_t &= f_t \otimes c_{t-1} + i_t \otimes g_t \\
 y_t, h_t &= o_t \otimes \tanh(c_t)
 \end{aligned}$$

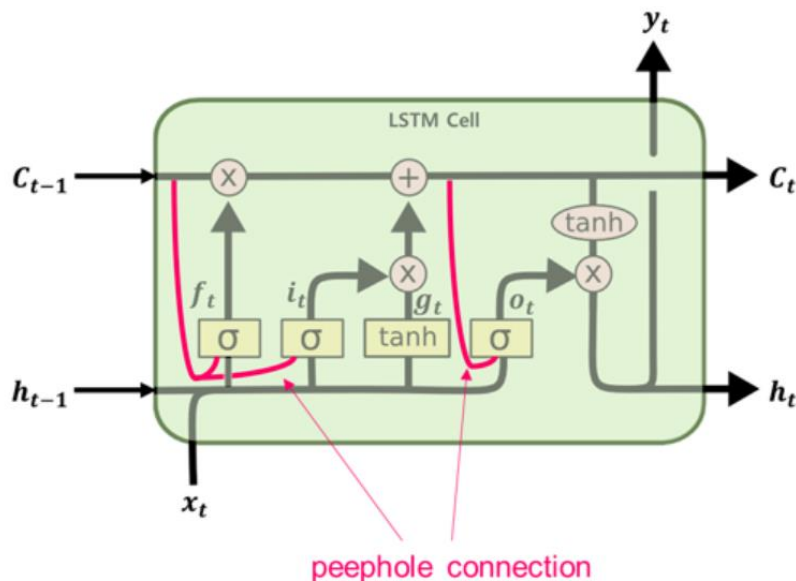
- $h$ (단기 상태),  $c$ (장기 상태) 나누어짐
- 장기 상태에서 기억할 부분, 삭제할 부분( $f$ ), 그리고 읽어 들일 부분( $i$ )을 학습

**Forget gate** :  $f$ 에 의해 제어되며 장기 상태의 어느 부분을 삭제할지 제어.

**Input gate** :  $i$ 에 의해 제어되며  $g$ 의 어느 부분이 장기 상태에 더해져야 하는지 제어

**Output gate** :  $o$ 는 장기 상태의 어느 부분을 읽어서  $h$ 로 출력해야 하는지 제어

# LSTM



$$f_t = \sigma(W_{cf}^T \cdot c_{t-1} + W_{xf}^T \cdot x_t + W_{hf}^T \cdot h_{t-1} + b_f)$$

$$i_t = \sigma(W_{ci}^T \cdot c_{t-1} + W_{xi}^T \cdot x_t + W_{hi}^T \cdot h_{t-1} + b_i)$$

$$o_t = \sigma(W_{co}^T \cdot c_t + W_{xo}^T \cdot x_t + W_{ho}^T \cdot h_{t-1} + b_o)$$

- 기존의 LSTM에서 gate controller( $f$ ,  $i$ ,  $o$ )는 입력  $x$ 와 이전 타임스텝의 단기 상태  $h$ 만 입력으로 받는다. 하지만 **핍홀(peephole)** 연결 해주면 gate controller에 이전 타임스텝의 장기 상태  $c$ 가 입력으로 추가되며, 좀 더 많은 맥락(context)를 인식할 수 있다.

# GRU Vs. LSTM

## 2) LSTM vs GRU

### - GRU

간단한 모델 : 쉽게 큰 네트워크 형성,

2개의 게이트만 사용 : 계산량 적다. 그래서 빠르다.

### - LSTM

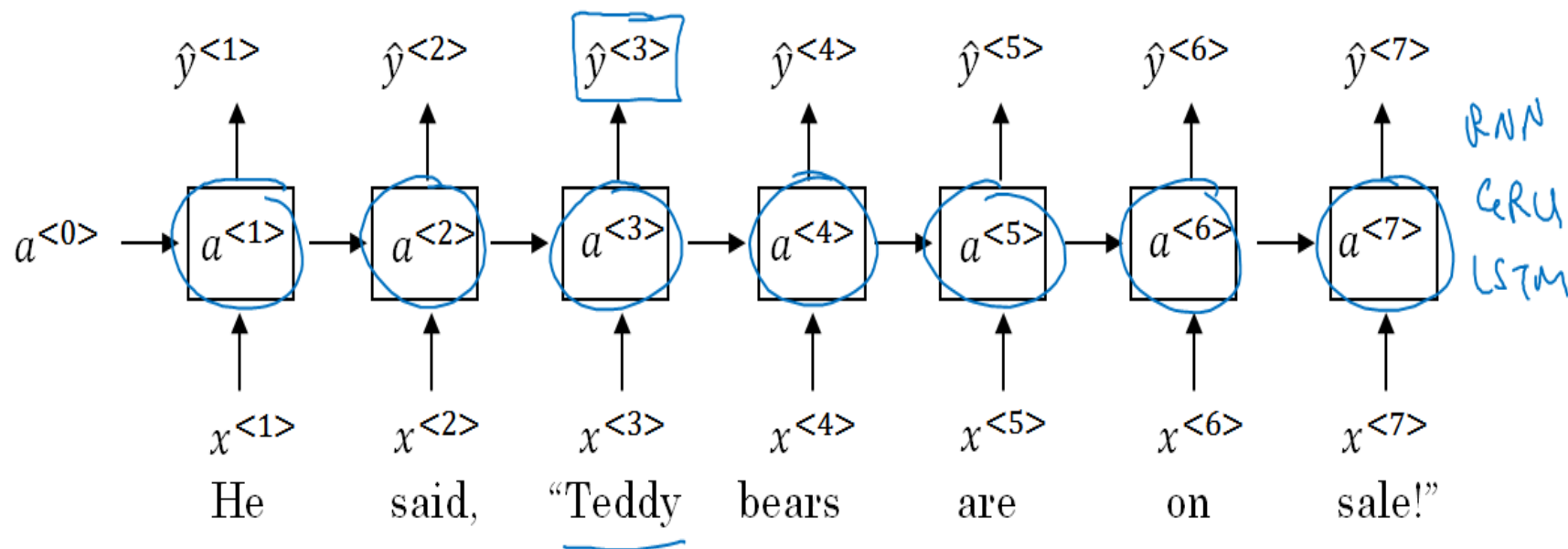
큰 모델 : 강하다.

3개의 게이트 사용 : 더 효과적이다.

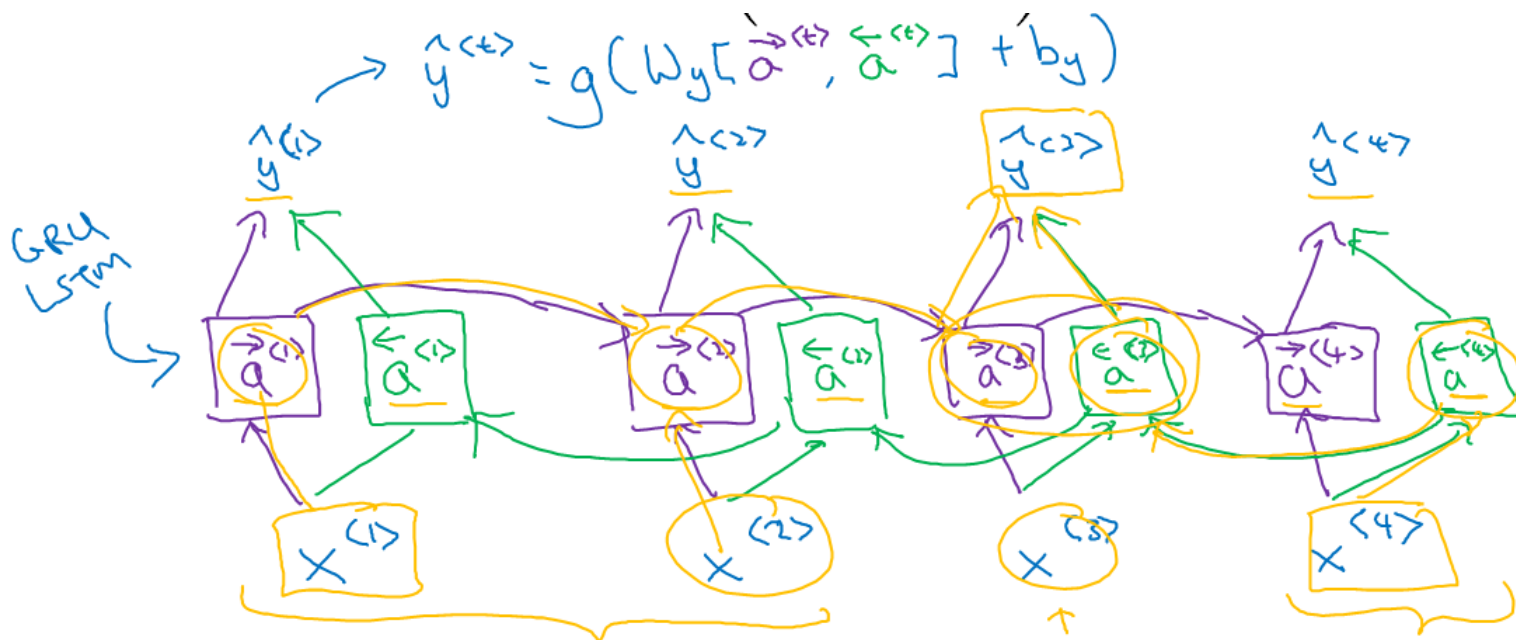
# BRNN

He said, "Teddy bears are on sale!"

He said, "Teddy Roosevelt was a great President!"

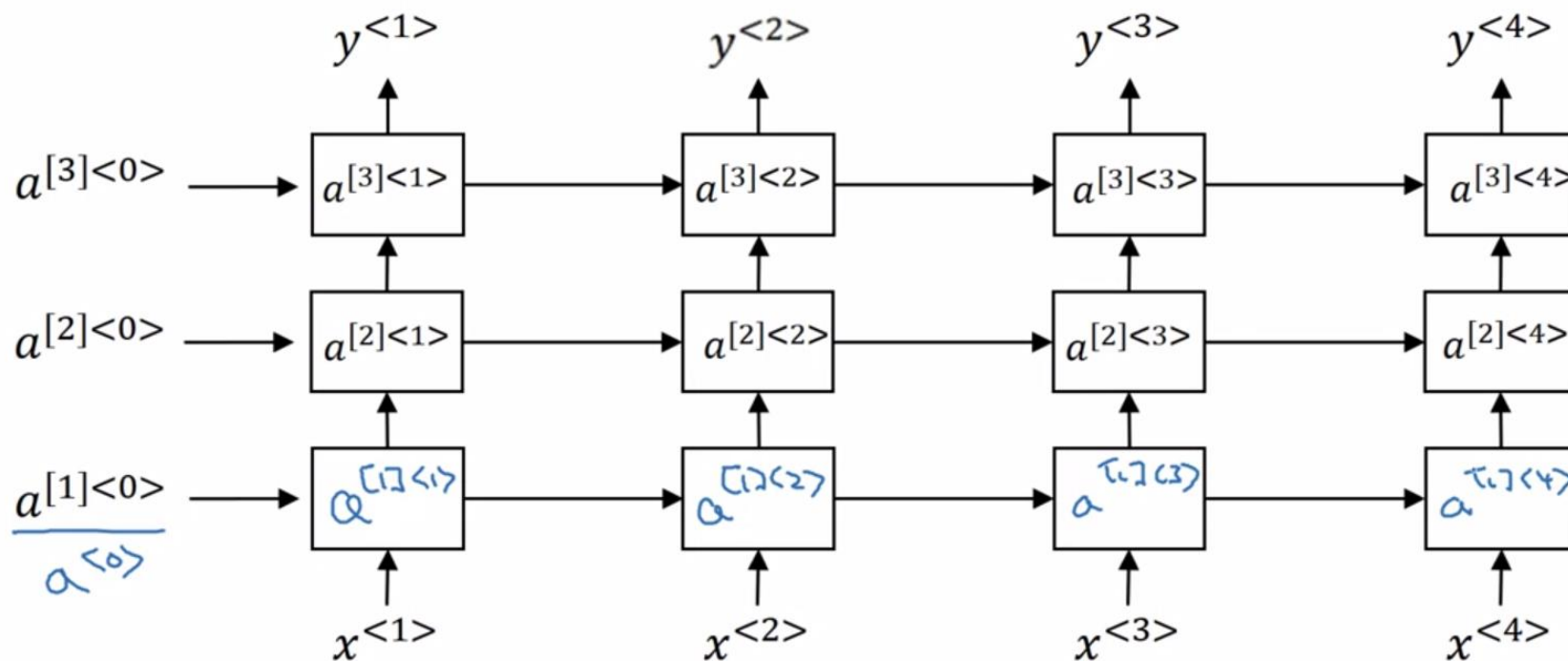


# BRNN



음성 인식의 경우, 멈추는 지점이 없으면 BRNN 성능 안 좋다

# Deep RNN



너무 많이 쌓으면 안 좋다  
3개도 충분히 많은 편

$$a^{[2]3} = g(w_a^{[2]} [a^{[1]2}, a^{[1]3}] + b_a^{[2]})$$

# Reference

Coursera Sequence Models – Andrew Ng

(<https://www.coursera.org/learn/nlp-sequence-models/home/week/1>)

욱이의 따뜻한 감성

(<https://m.blog.naver.com/ehdsnck/221778218877>)

EXCELSIOR

(<https://excelsior-cjh.tistory.com/185>)