

# HCT NLP Week 3

问答摘要与推理  
Seq2Seq（一）

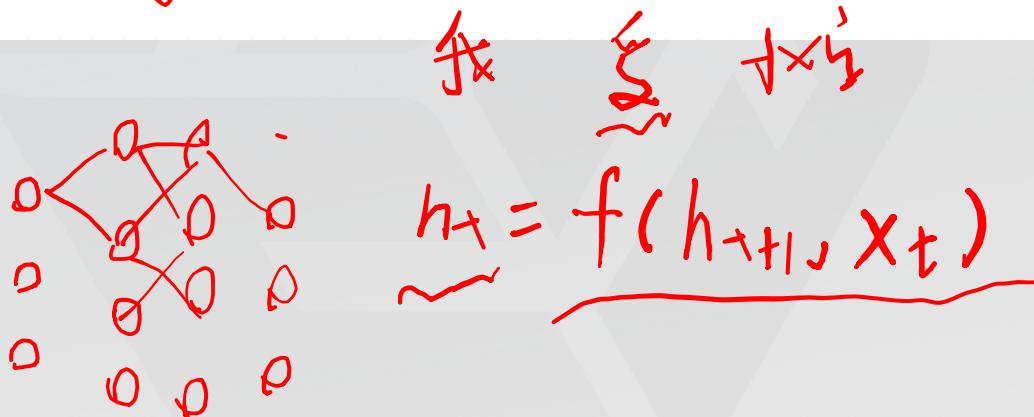
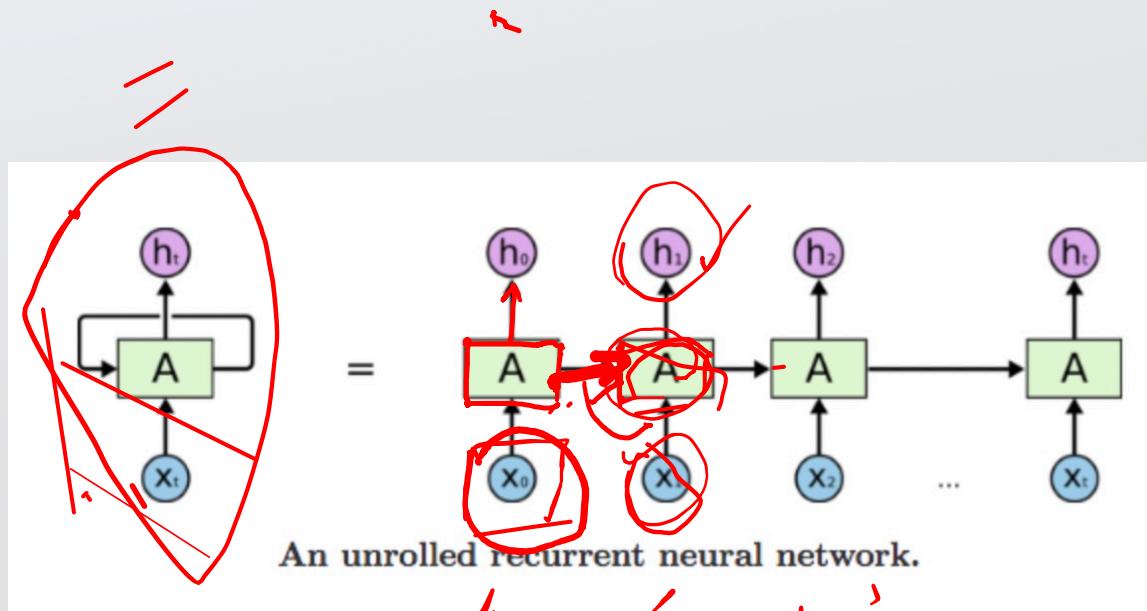
# Outline

1. RNN、LSTM、GRU
2. Encoder-Decoder结构
3. Attention机制
4. 作业

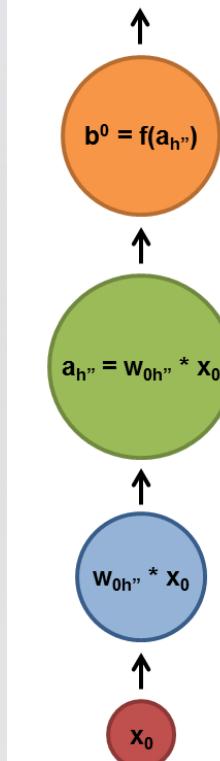
# 1. RNN、LSTM、GRU

# RNN

## Recurrent Neural Network



$b^0$  is fed to next layer



# RNN

## Recurrent Neural Network

multiplication

$$\begin{matrix} 0 & 0 \\ [1, 2] \end{matrix} \quad \begin{matrix} 0 & 0 \\ [4, 6] \end{matrix}$$
$$\begin{matrix} 0 & 0 \\ [4, 3] \end{matrix}$$

addition

$$\begin{matrix} 0 & 0 \\ [5, 5] \end{matrix}$$

concatenate

Concat

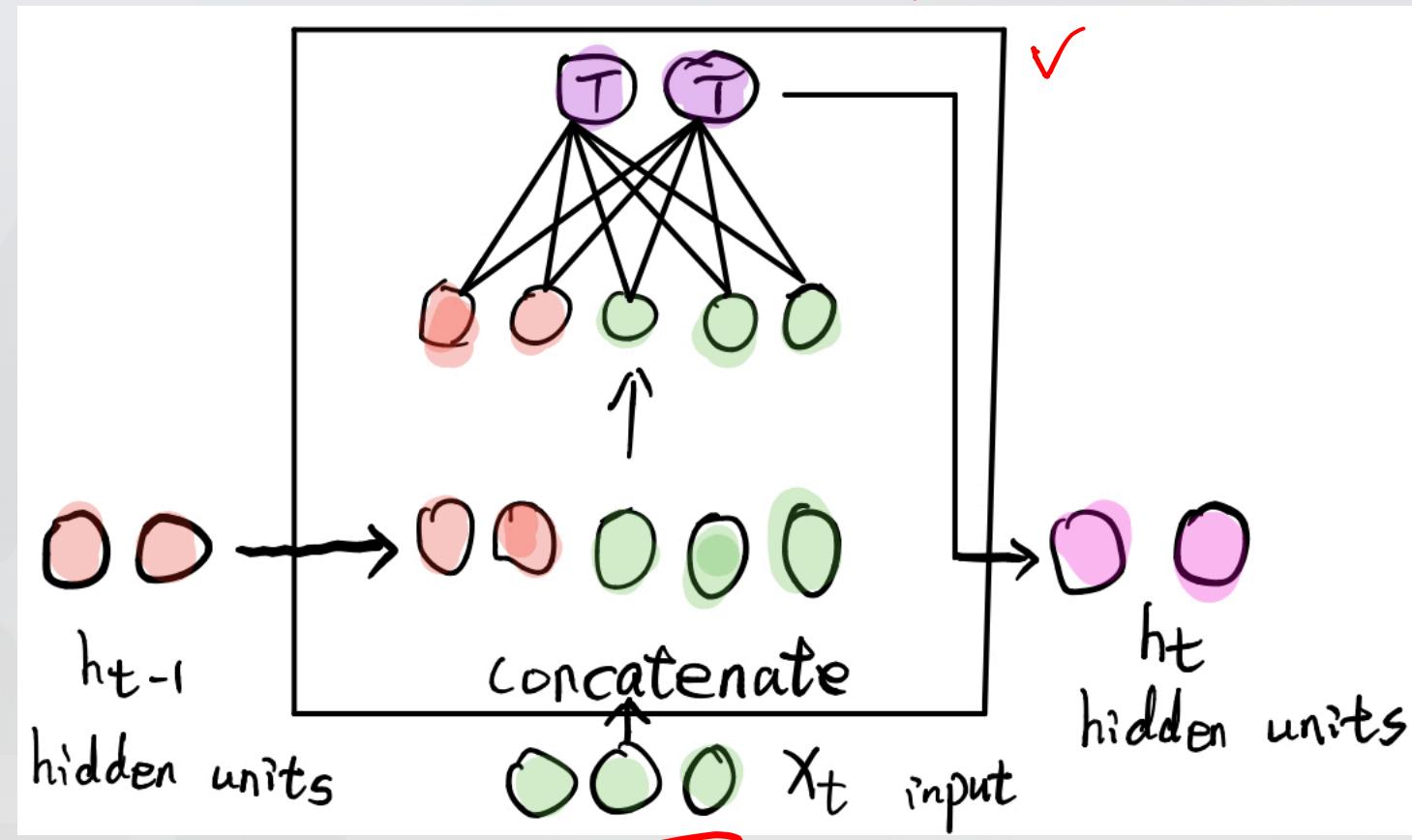
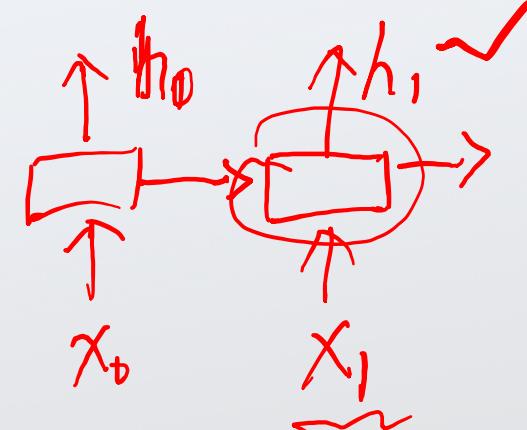
$$\begin{matrix} 0 & 0 \\ [1, 2] \end{matrix} \quad \begin{matrix} 0 & 0 & 0 & 0 \\ [1, 2, 4, 3] \end{matrix}$$
$$\begin{matrix} 0 & 0 \\ [4, 3] \end{matrix}$$

$$0 \ 0 \rightarrow \begin{matrix} 1 \\ 1 \end{matrix}$$
$$[1, 2] \quad [0, 1]$$



# RNN

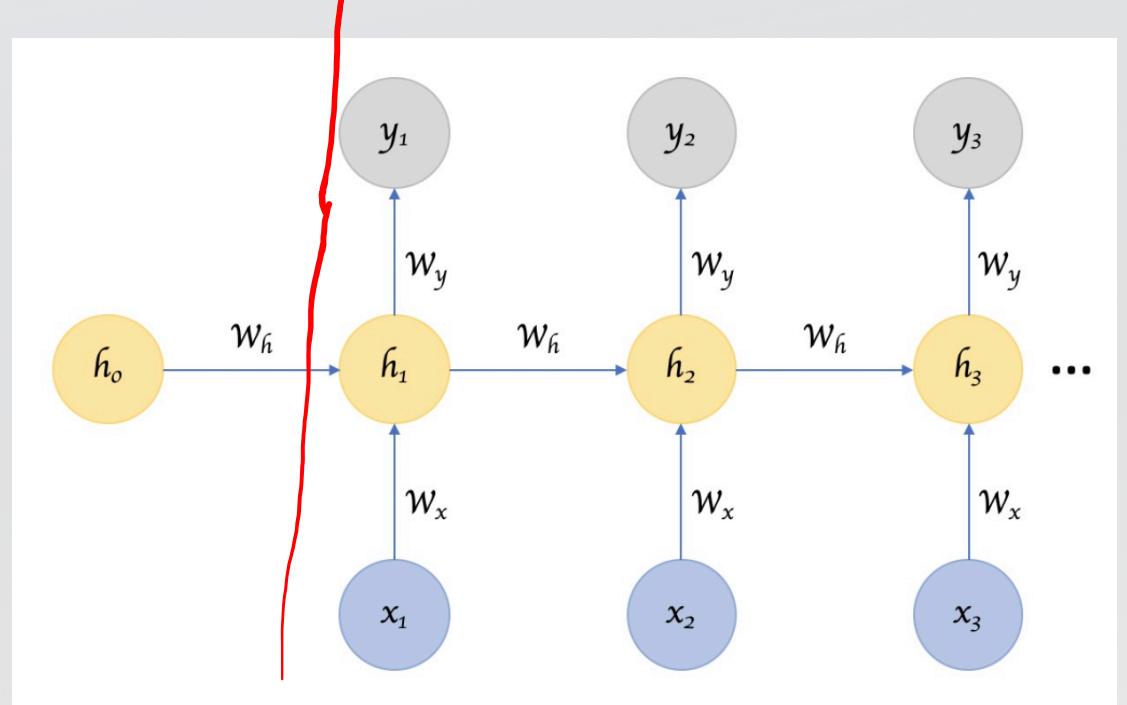
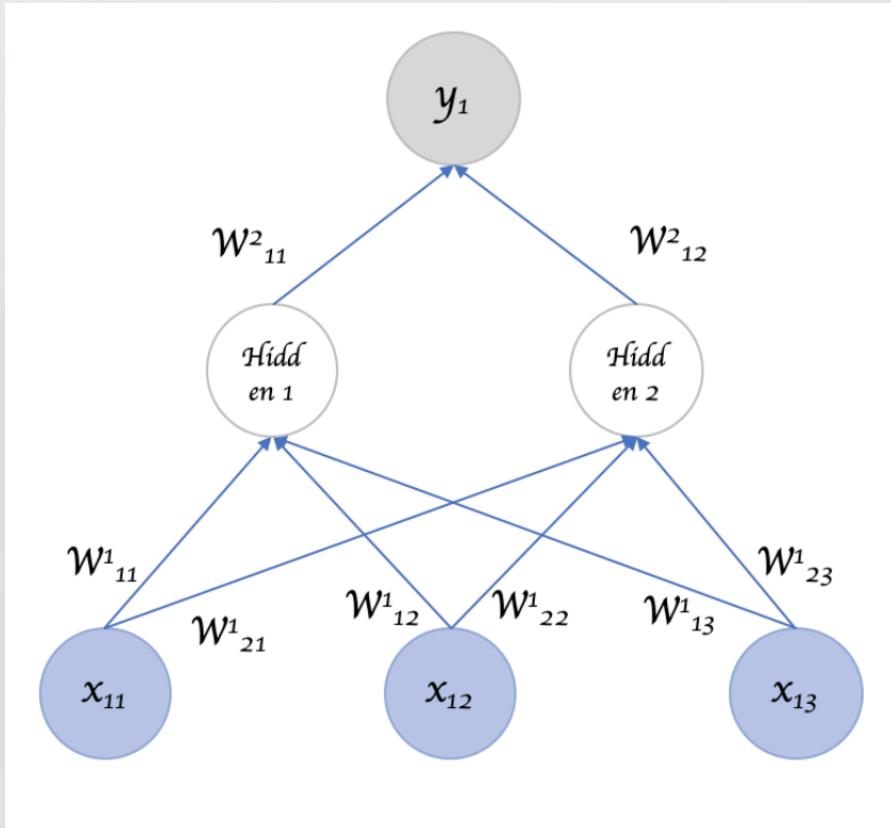
## Recurrent Neural Network



(1, 2, 3) ✓

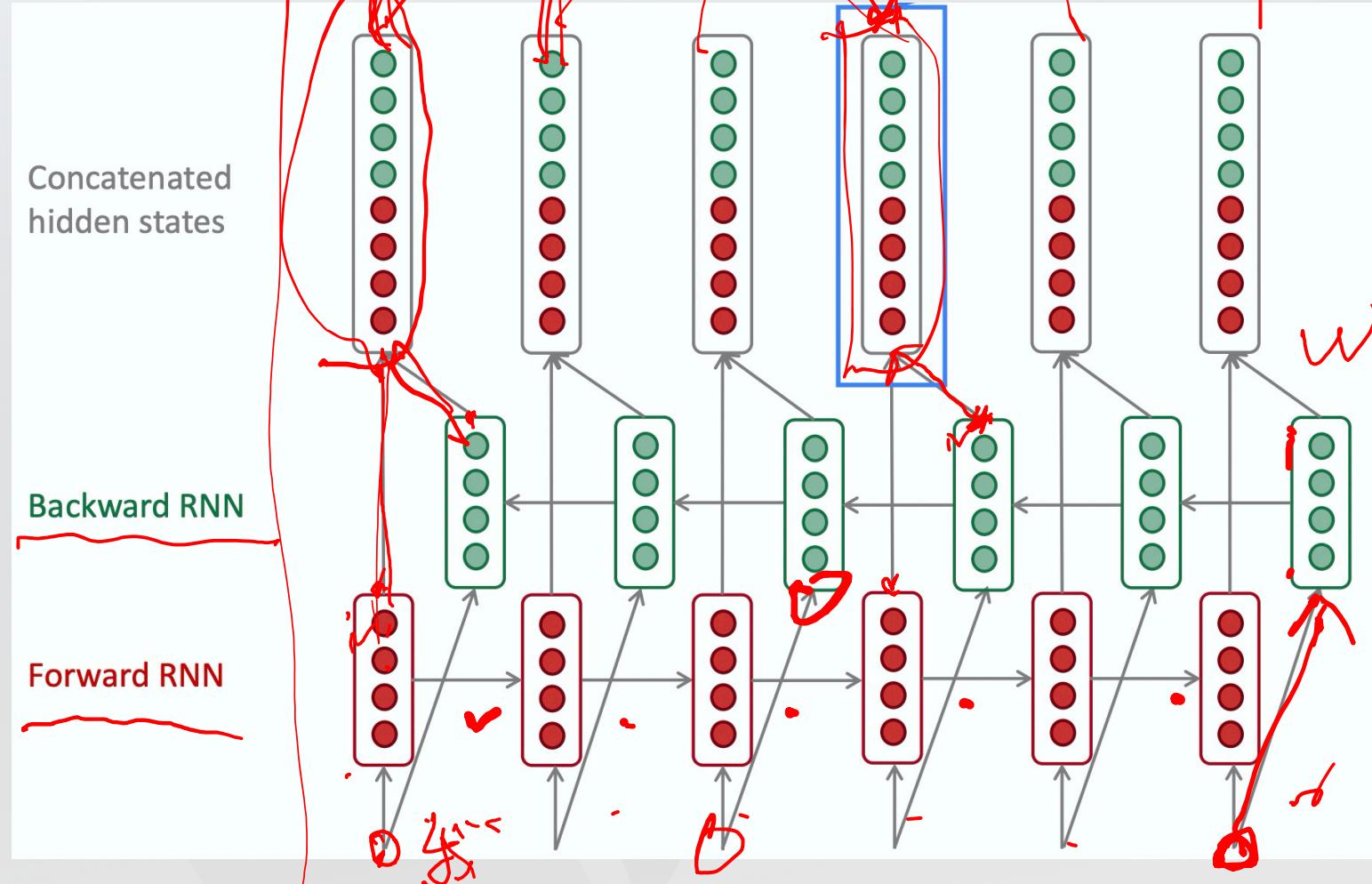
# RNN

parameter sharing



# RNN

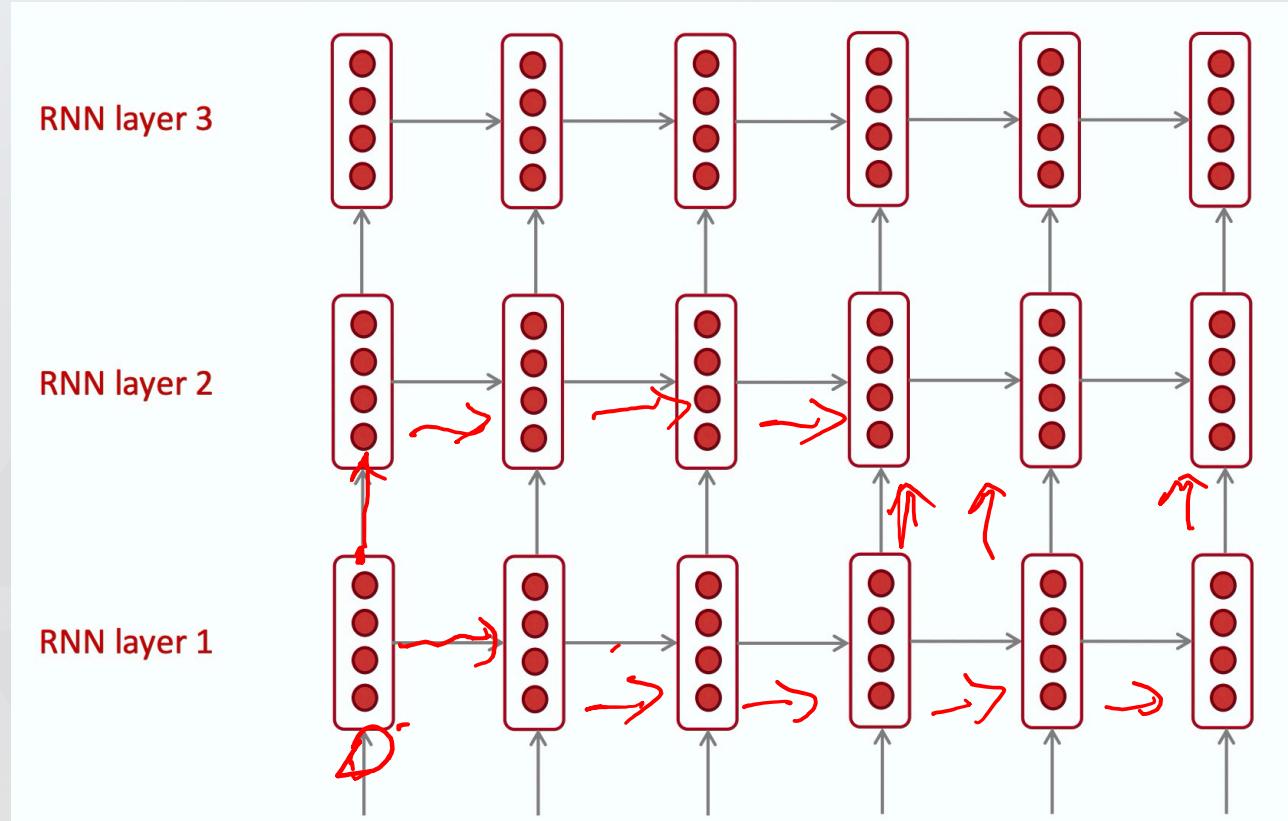
## Bidirectional RNNs



$$F: \vec{h}^{(t)} = RNN_{FW}(h^{(t-1)}, x^{(t)})$$
$$B: \overleftarrow{h}^{(t)} = RNN_{BW}(h^{(t+1)}, x^{(t)})$$
$$h^{(t)} = [\vec{h}^{(t)}, \overleftarrow{h}^{(t)}]$$
$$\text{Loss } W_{ij} = h$$

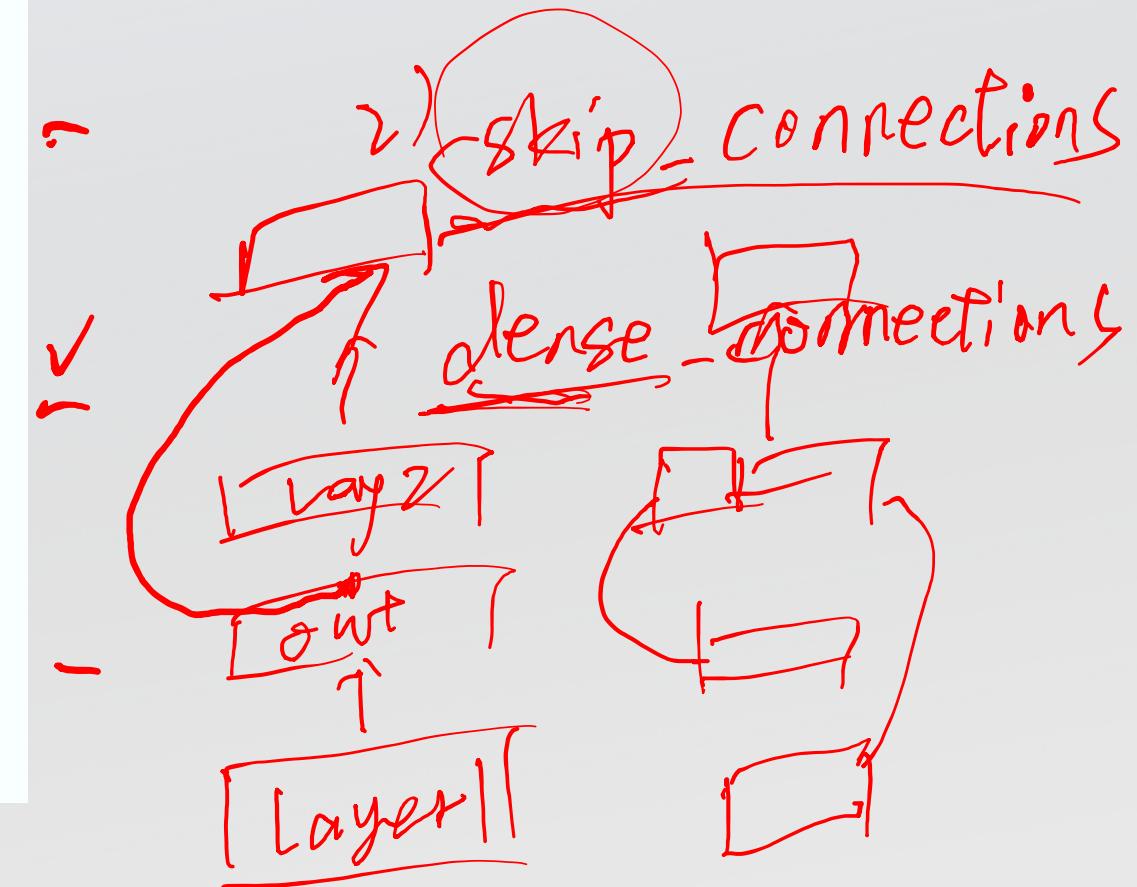
# RNN

## Bi-Deep Deep RNNs

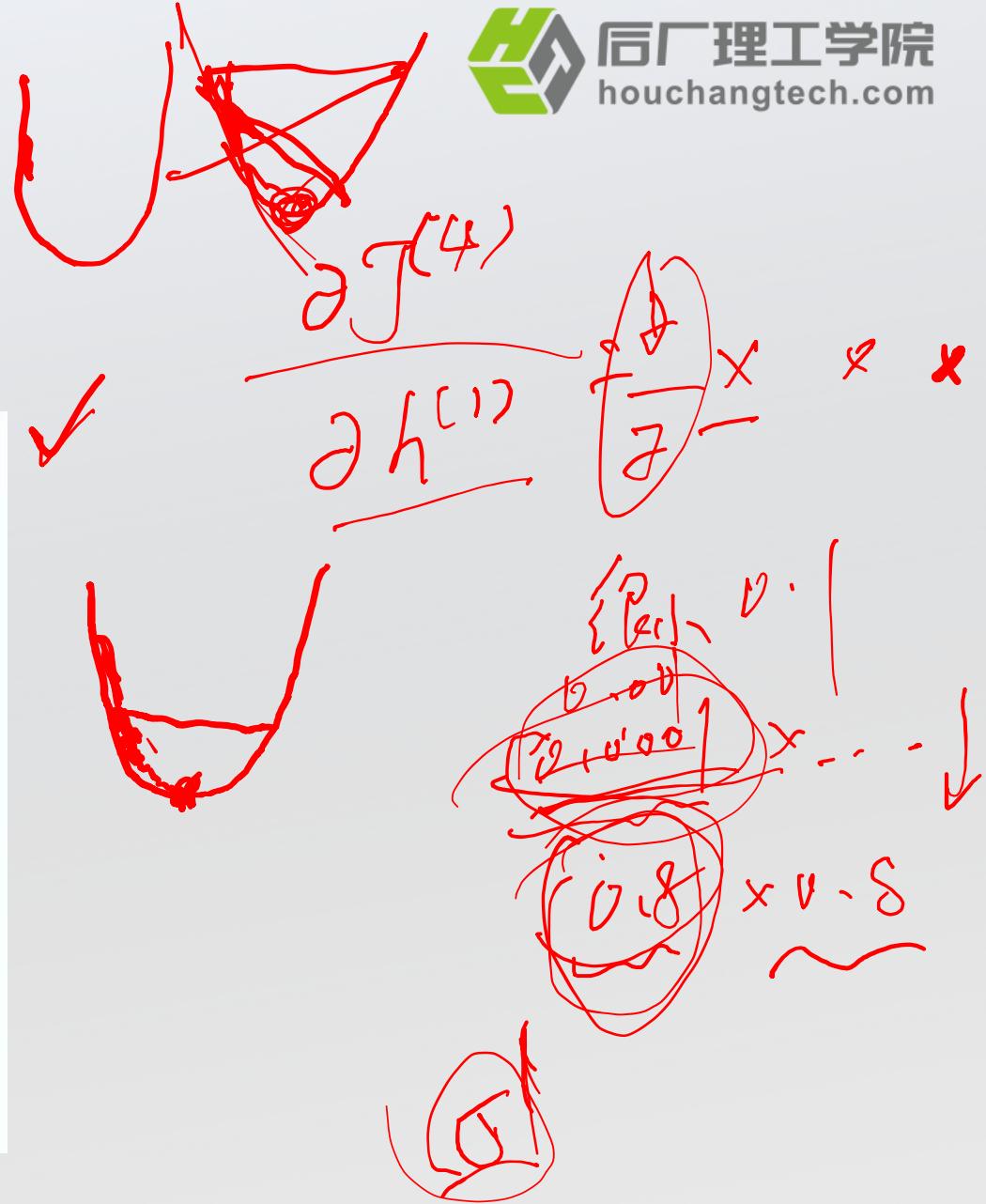
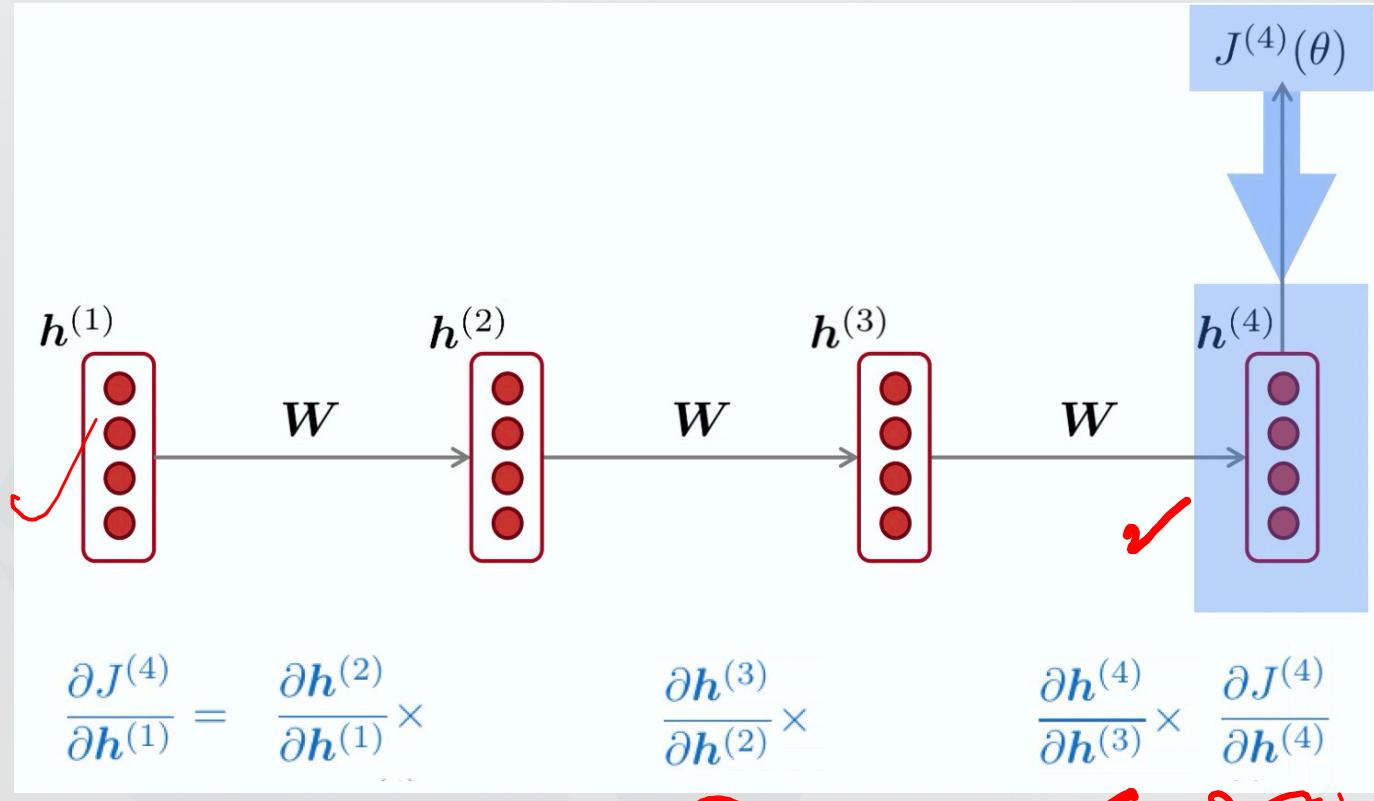


NMT

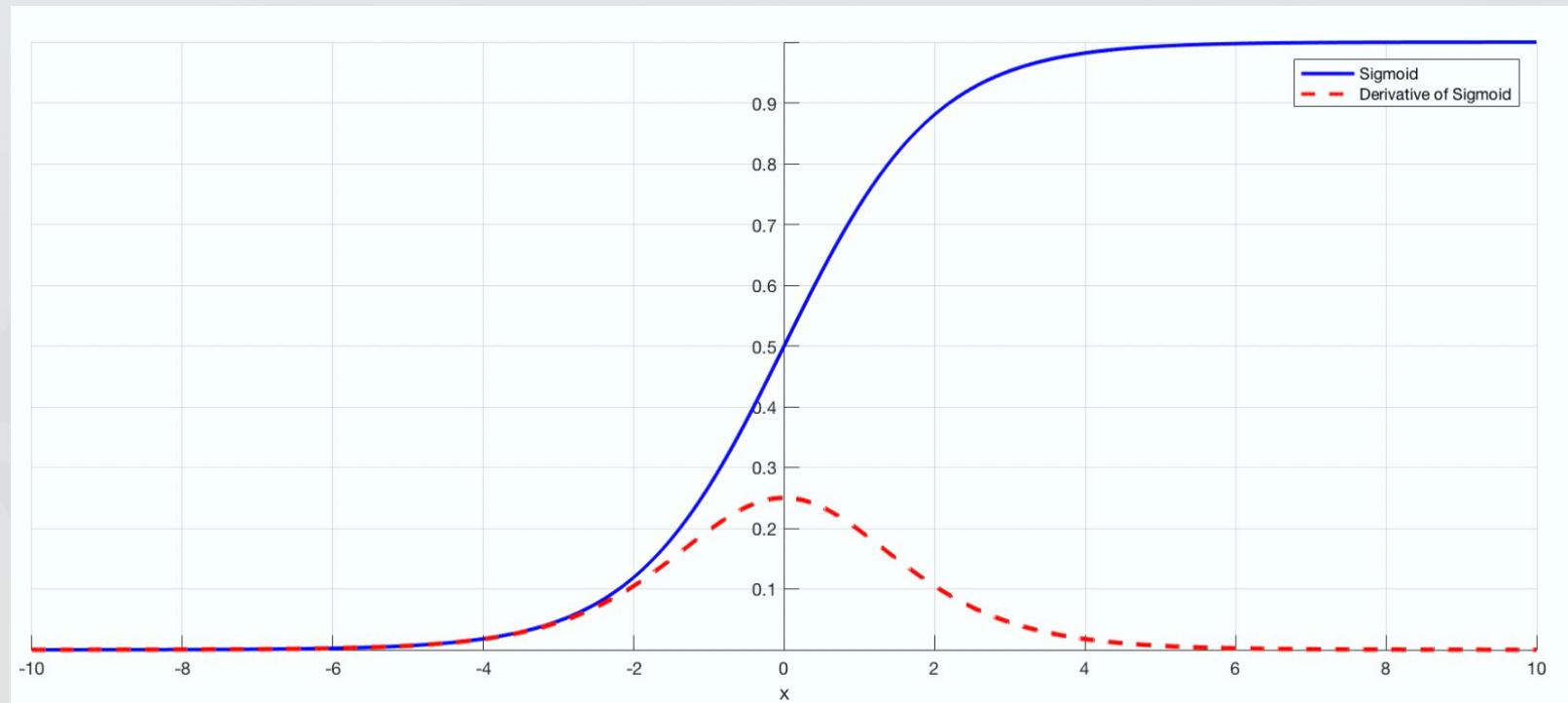
2 to 4 + 4



# Exploding and Vanishing Gradient Problem



# Exploding and Vanishing Gradient Problem



2/15

# Exploding Gradient

gradient clipping

# Vanishing Gradient

Identity Initialization

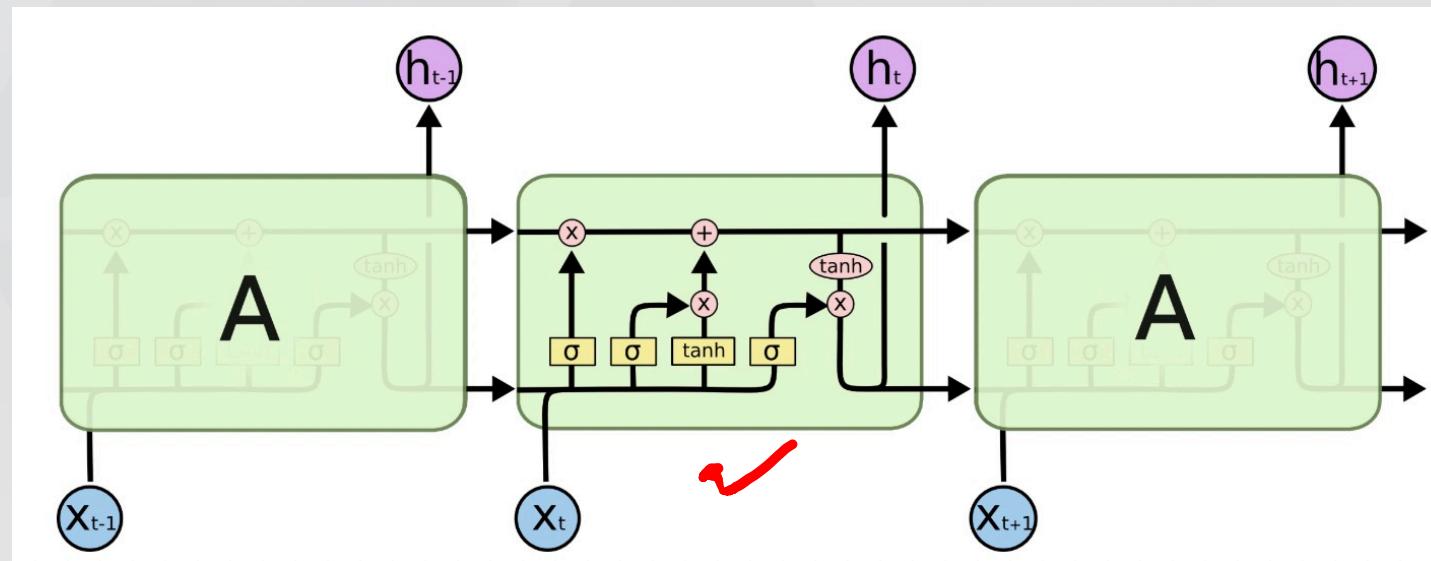
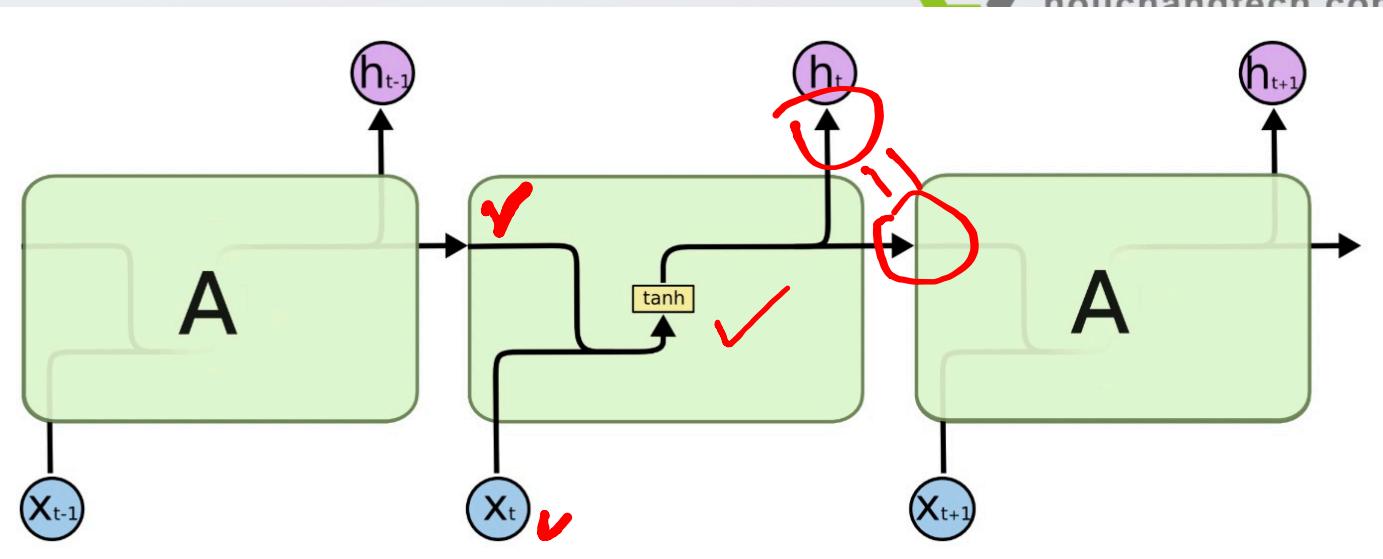
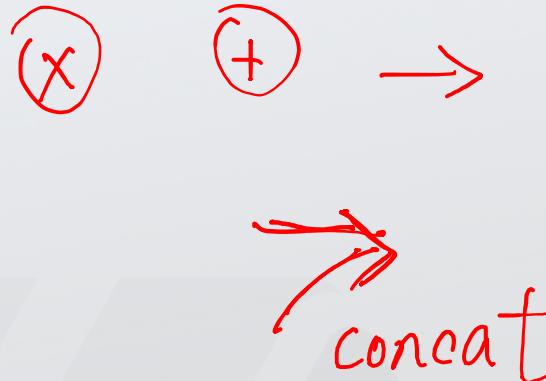
LSTM

Residual Networks

Batch Normalization

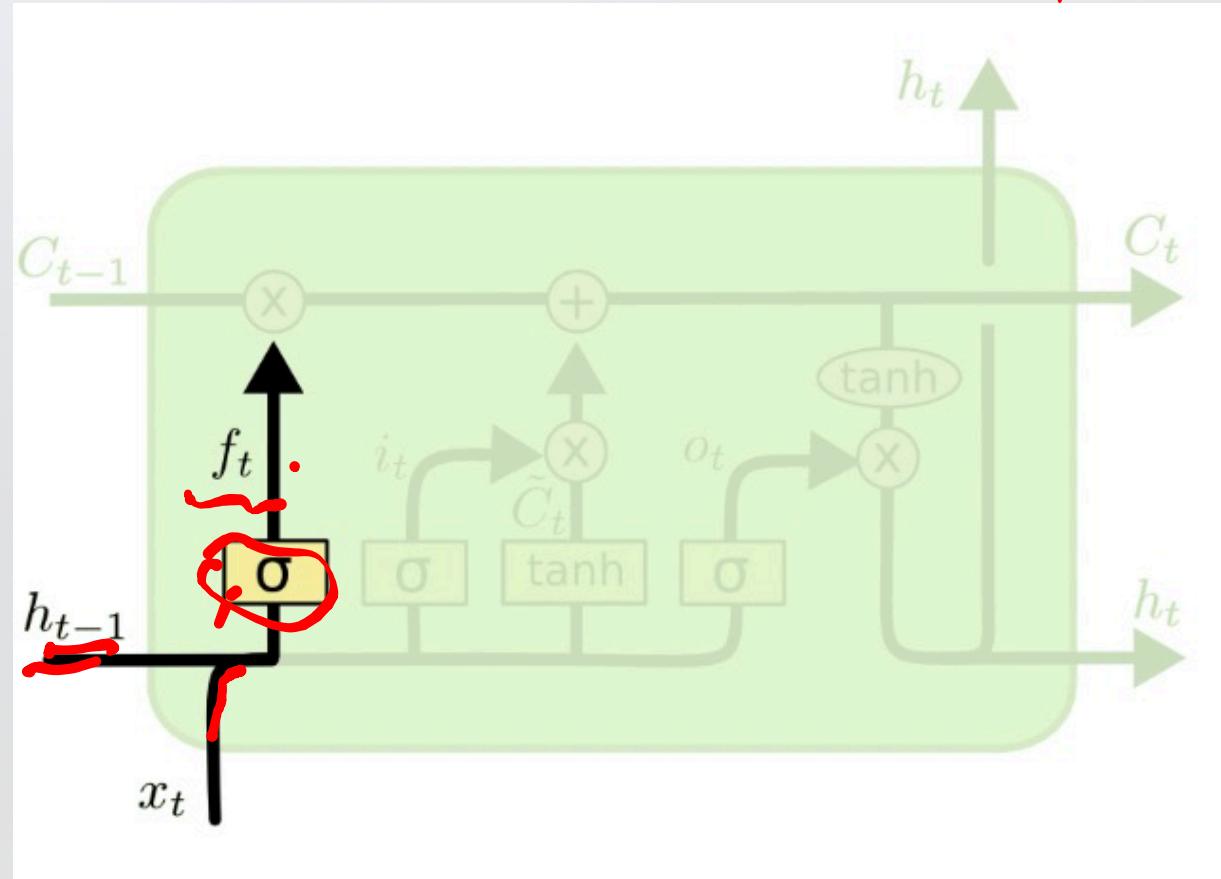
# LSTM

## Long Short Term Memory



# LSTM

3个门 gates  
input  
forget  
output



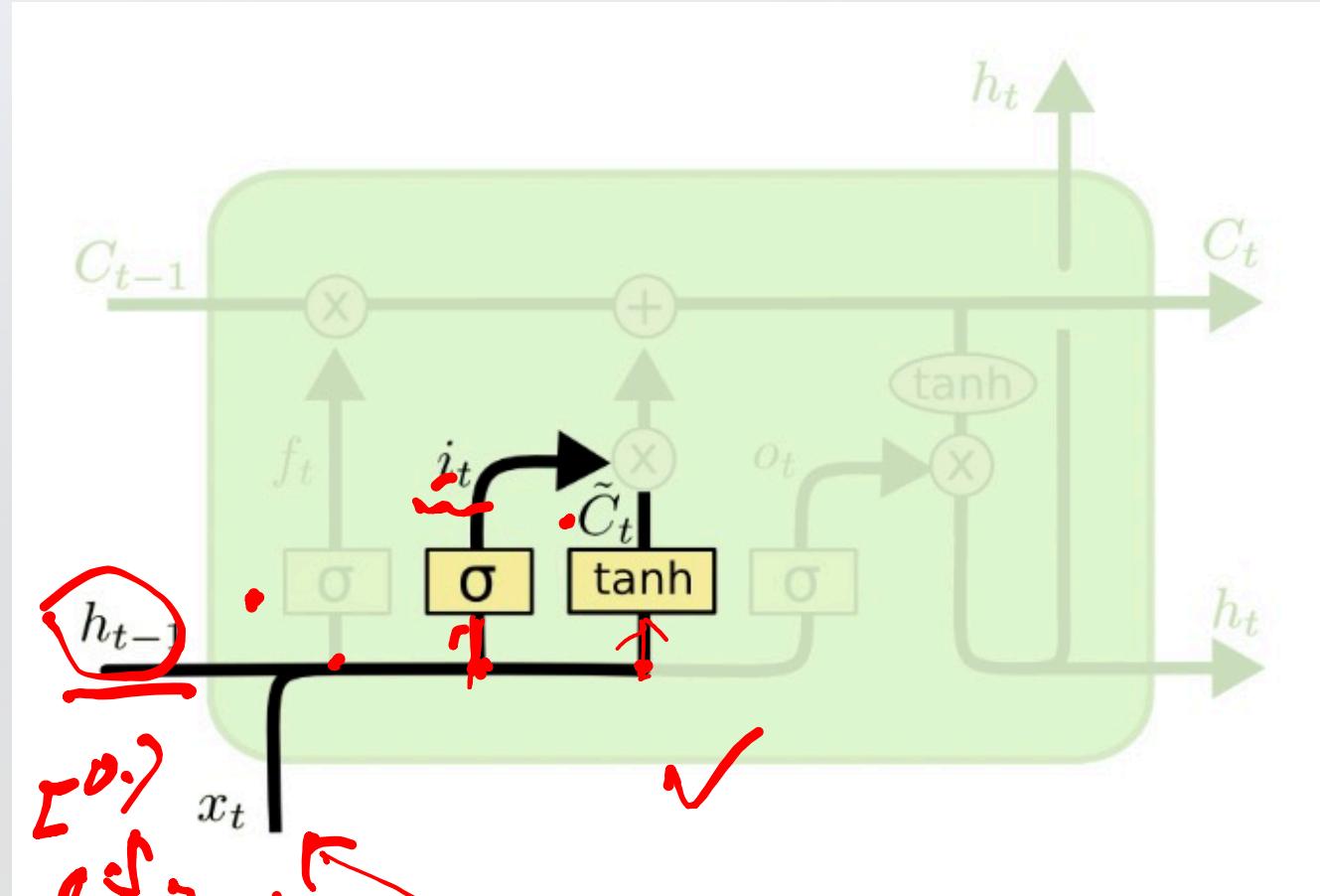
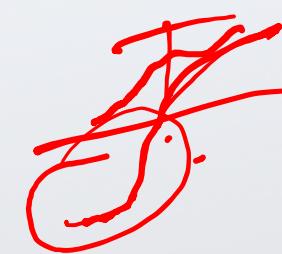
~~input gate~~

forget gate

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

# LSTM

input gate



$[0.1, 0.2, 0.3]$   
 $[0.4, 0.5, 0.6]$

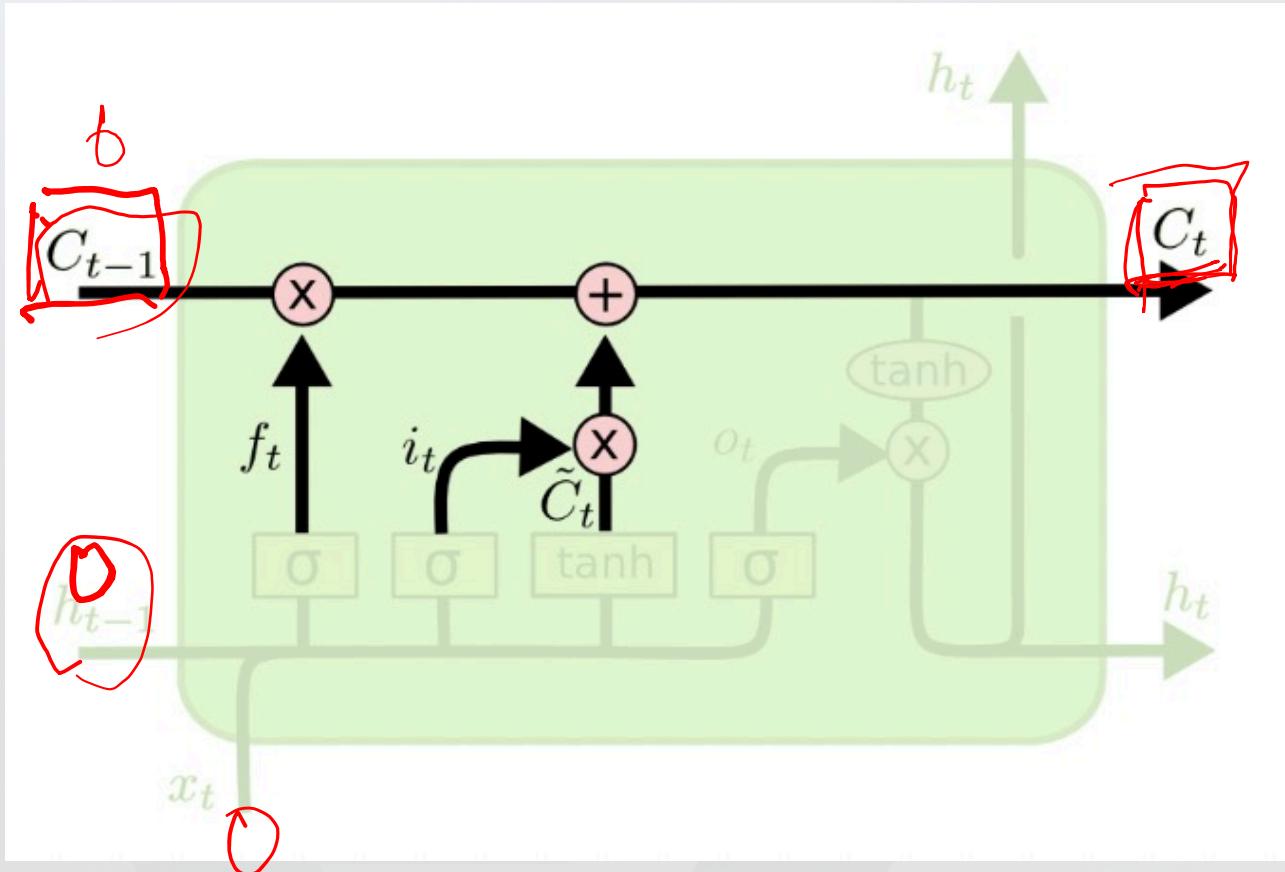
$x_t$

tensor

$$i_t = \sigma(w_i \cdot [h_{t-1}; x_t] + b_i)$$

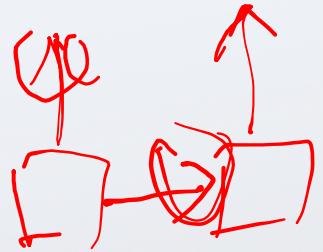
$$\tilde{C}_t = \tanh(w_c \cdot [h_{t-1}; x_t] + b_c)$$

# LSTM

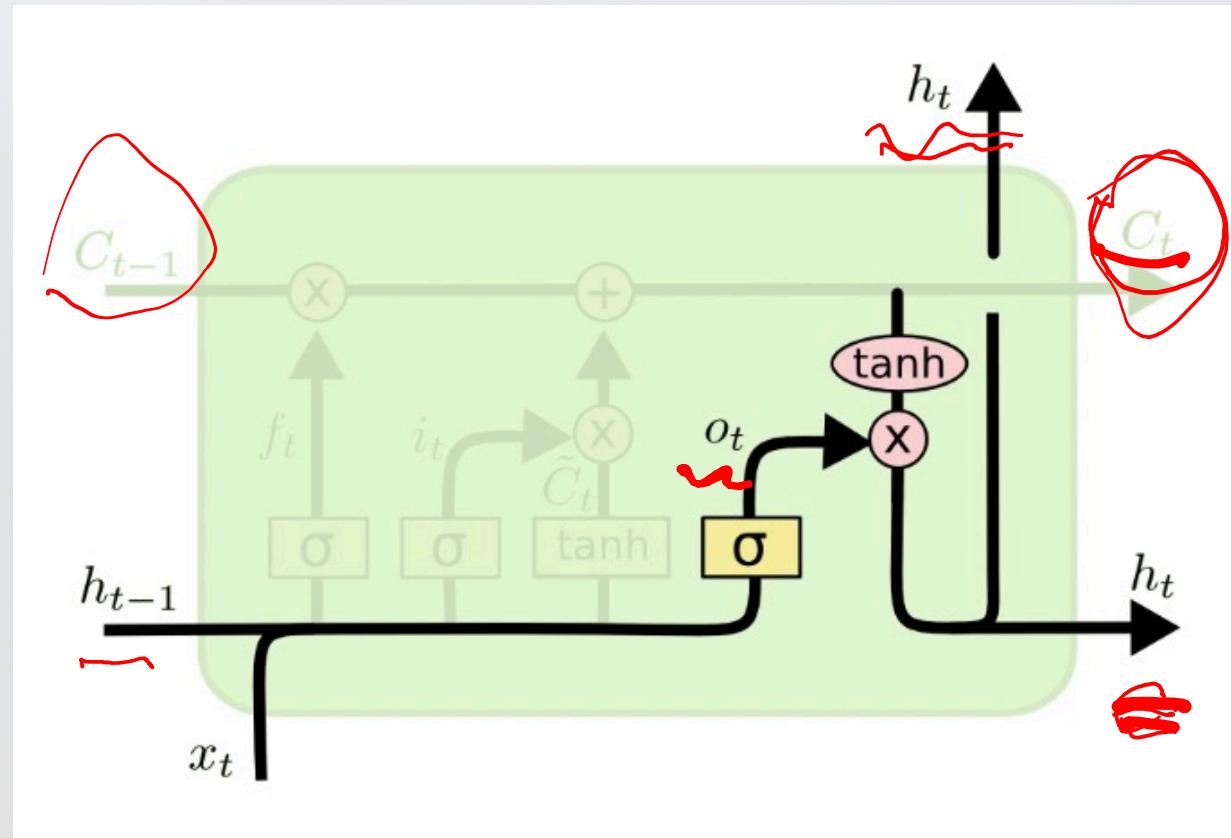
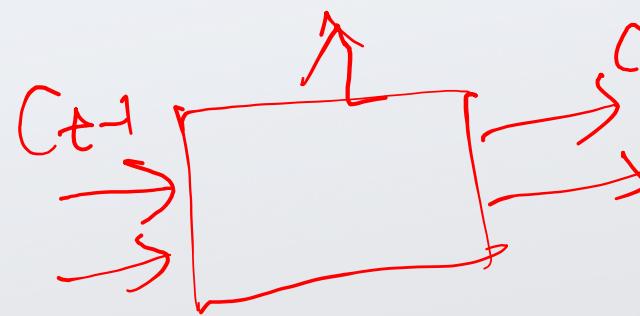


$$C_t = \underbrace{f_t * C_{t-1} +}_{i_t * \tilde{C}_t} \underbrace{i_t * \tilde{C}_t}_{\sim}$$

LSTM



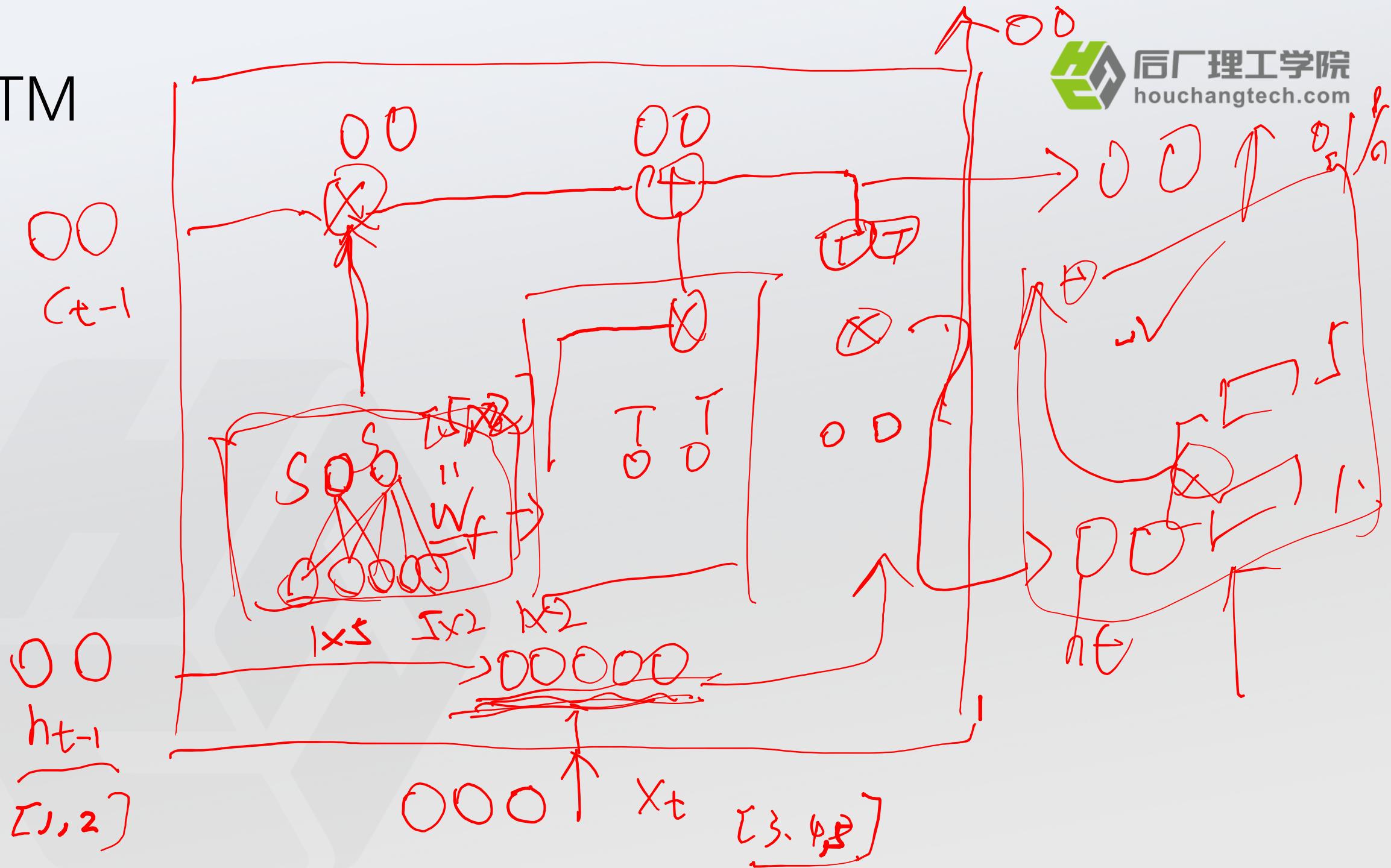
Output gate



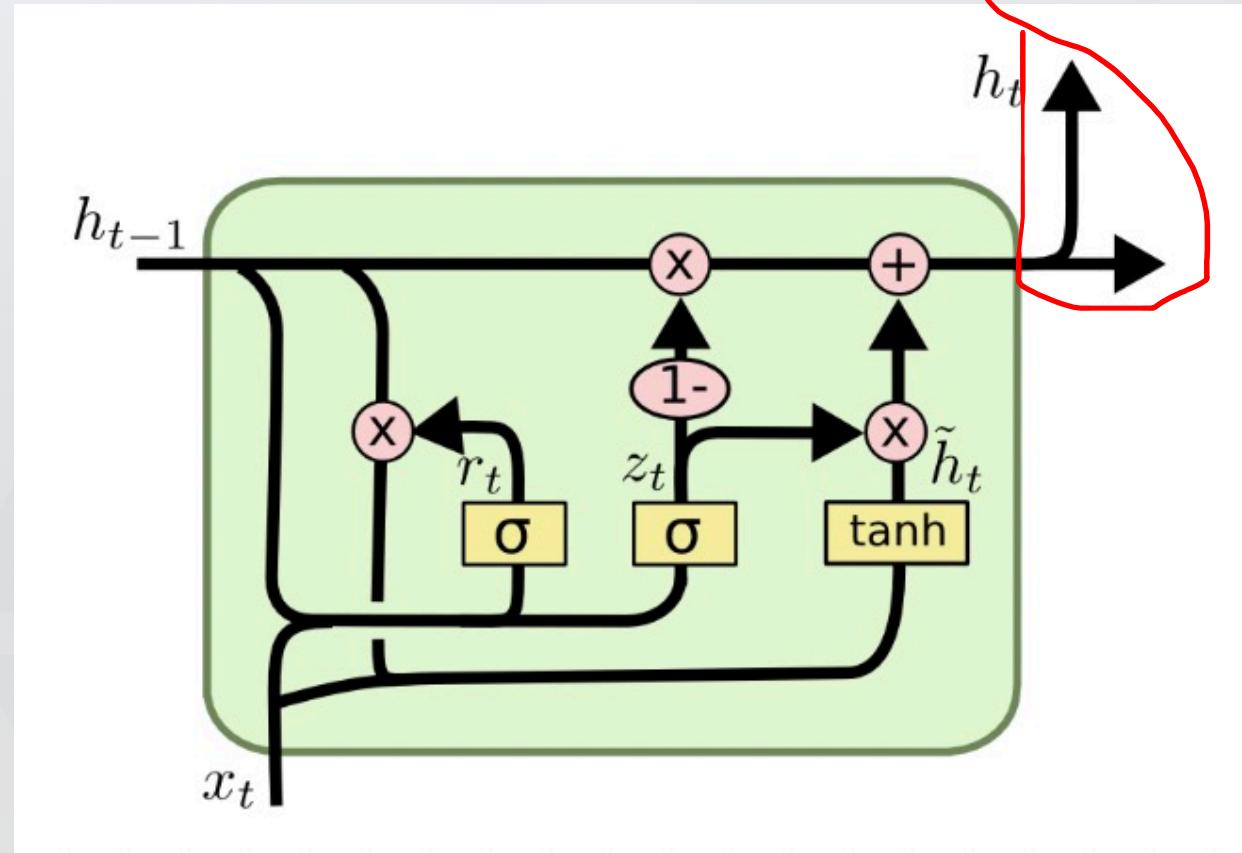
$$O_t = \sigma(W_o[h_{t-1}; x_t] + b_o)$$

$$h_t = O_t \star \tanh(C_t)$$

# LSTM



# GRU



①  $r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$

$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$

$\tilde{h}_t = \tanh(W \cdot [r_t \cdot h_{t-1}, x_t])$

$$h_t = \underbrace{(1 - z_t) \cdot h_{t-1}}_{\text{Forget}} + \underbrace{z_t \cdot \tilde{h}_t}_{\text{Update}}$$

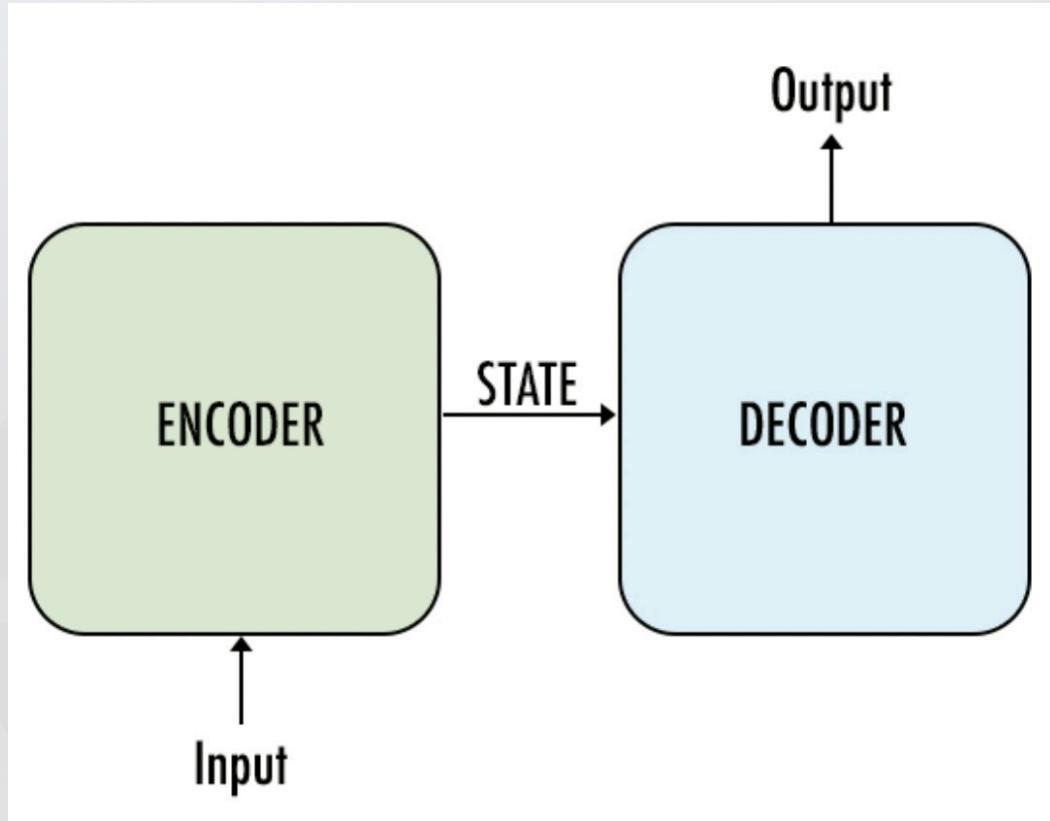
## 2. Encoder-Decoder结构

22:00

# Seq2Seq

## Encoder-Decoder结构

Sequence to sequence was first introduced by Google in 2014

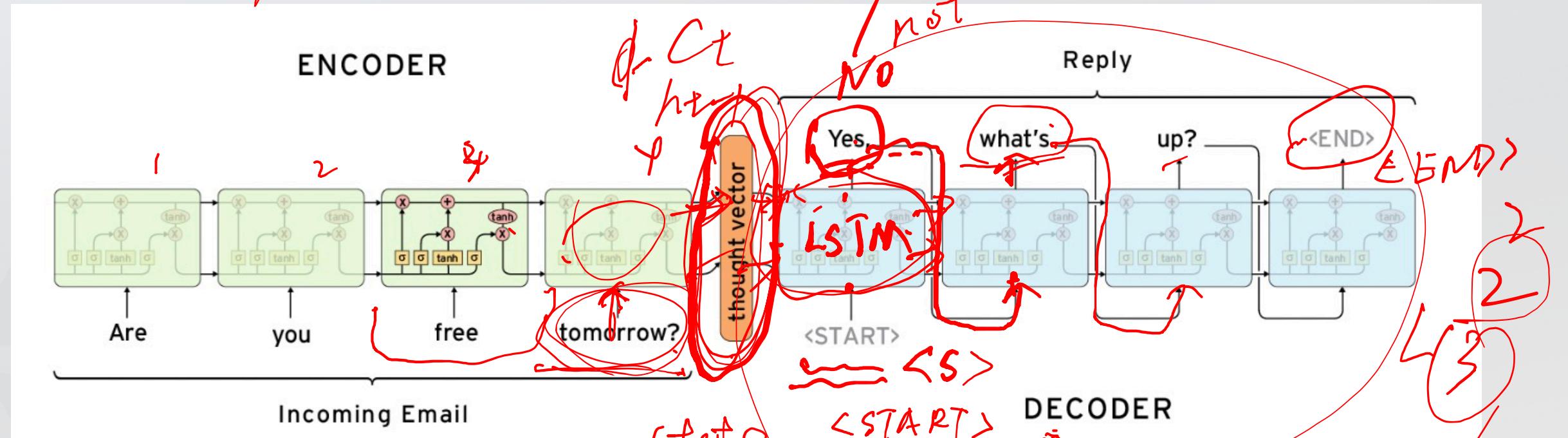


1. Speech Recognition
2. Machine Language Translation
3. Name entity/Subject extraction
4. Relation Classification
5. Path Query Answering
6. Speech Generation
7. Chatbot
8. Text Summarization
9. Product Sales Forecasting

# Seq2Seq

## Encoder-Decoder结构

Sequence ↗



state.

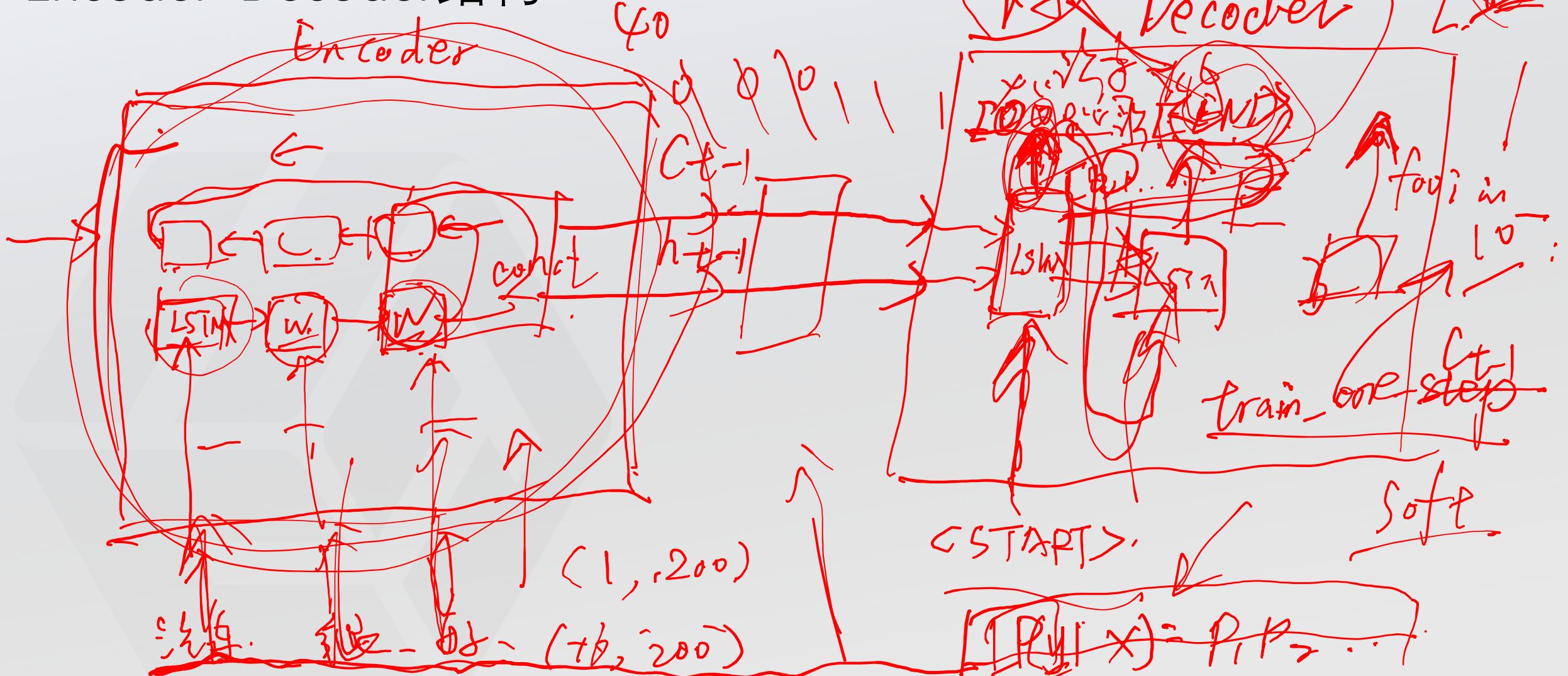
$[1, 200]$   
shape.

$[1, 4, 200]$

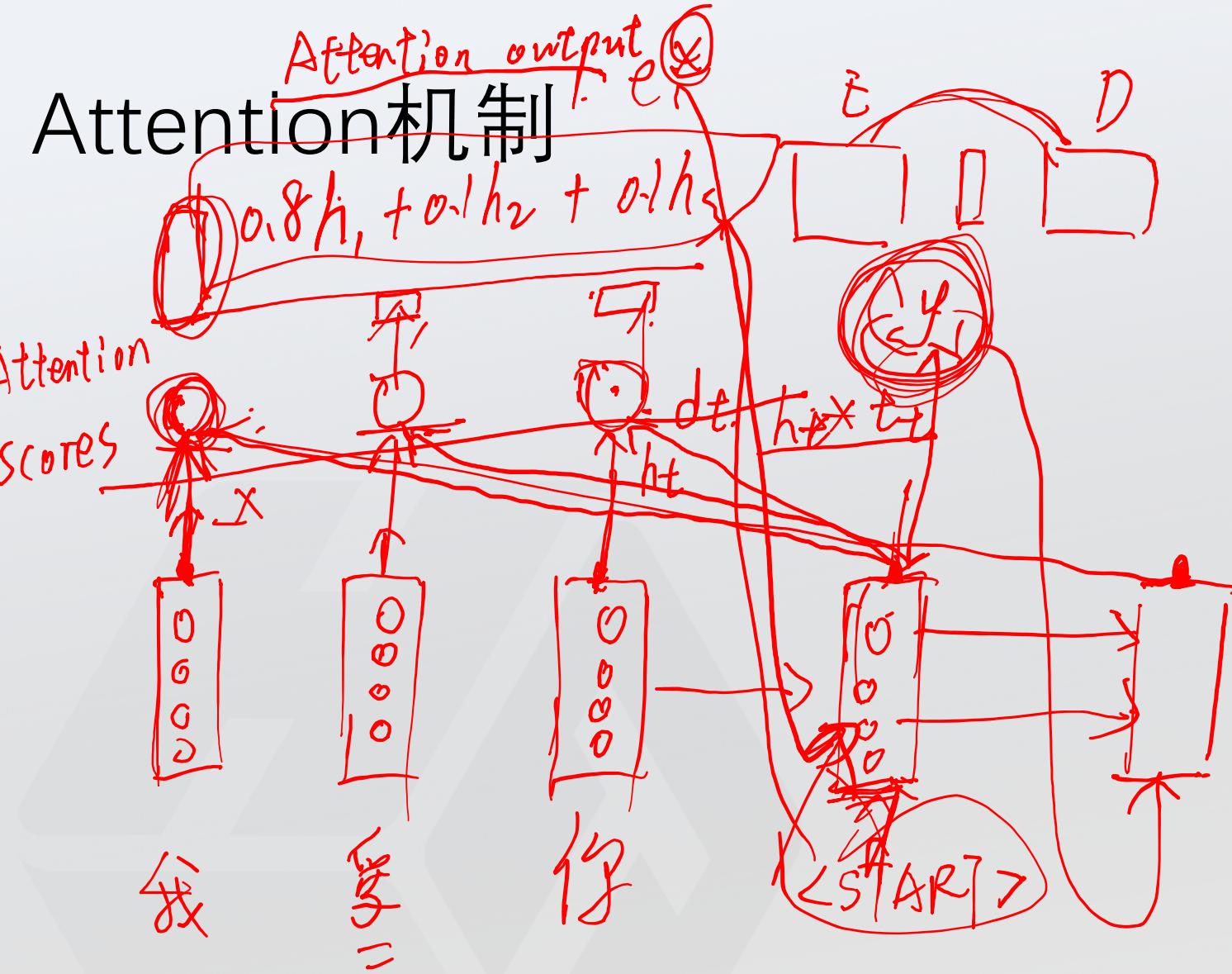
1 2 3 4 5  
XWZ

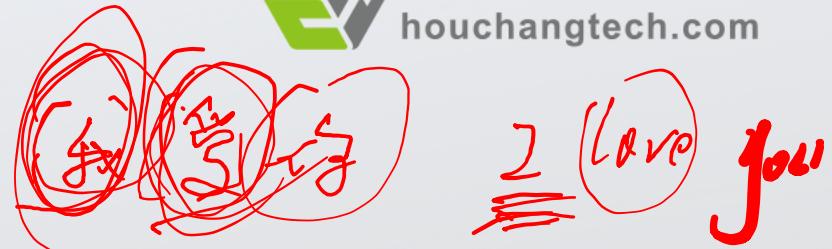
# Seq2Seq

## Encoder-Decoder结构

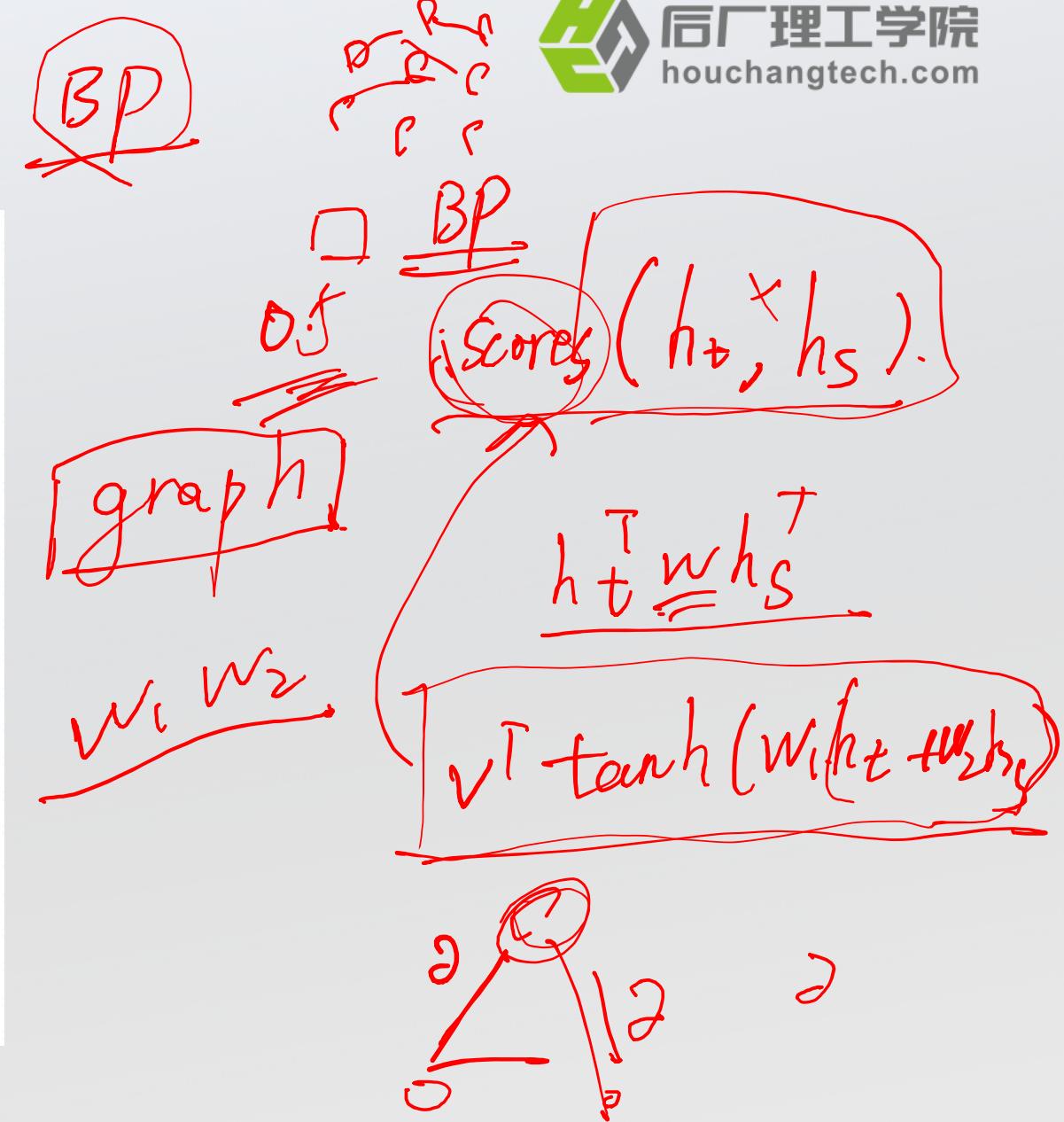
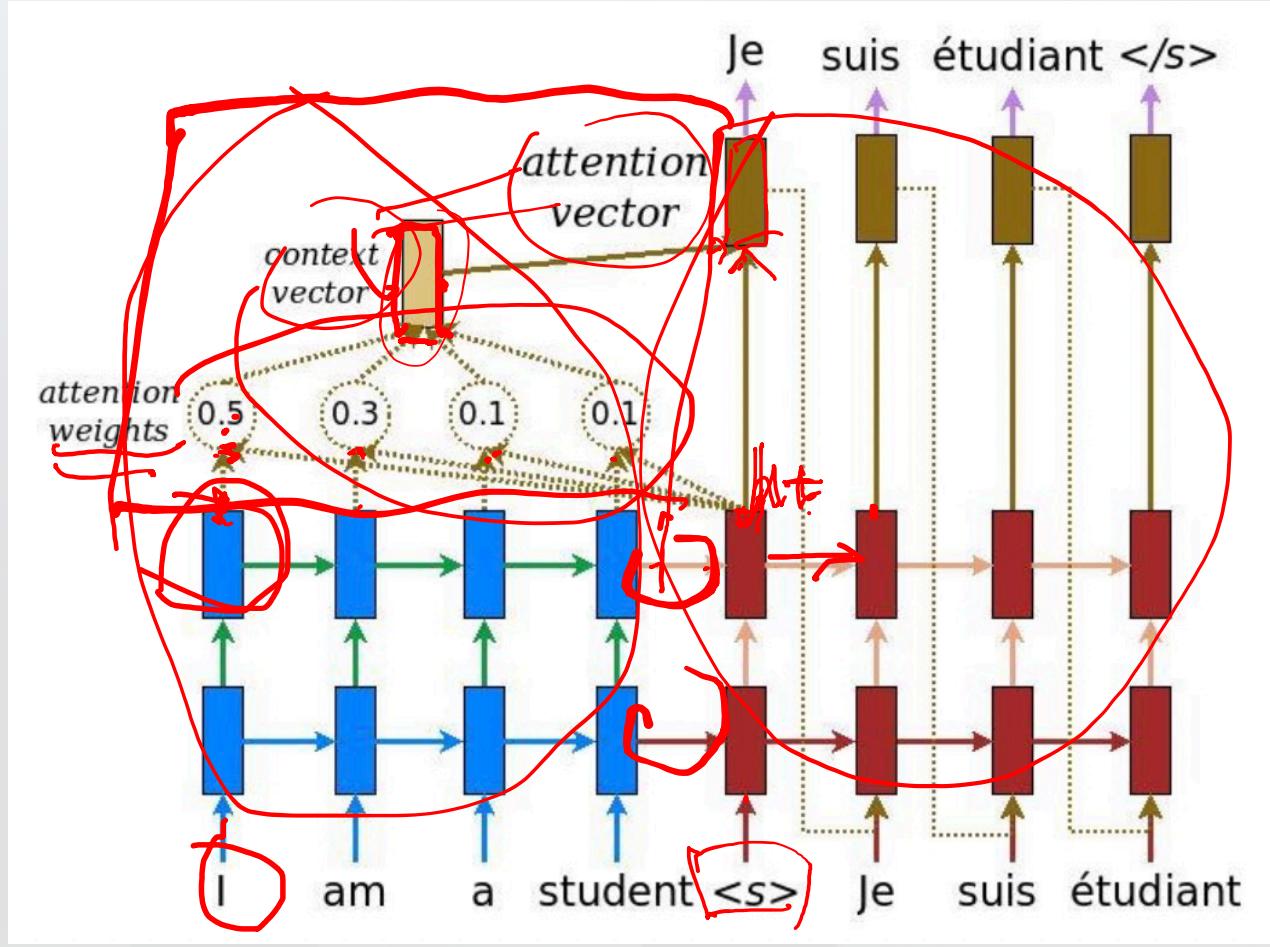


### 3. Attention机制





# Attention机制



## 4. 作业

Bye !