

# HCT NLP Week 2

问答摘要与推理  
词向量实践与RNNs

# Outline

- 词向量计算两种优化方法
- 词向量在工程中的具体实现
- RNN递归神经网络结构
- RNN、LSTM、GRU

RNNs

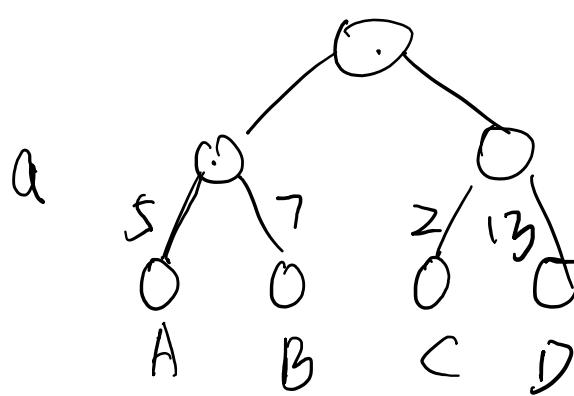
# Outline

- 词向量计算两种优化方法
- 词向量在工程中的具体实现
- RNN递归神经网络结构
- RNN、LSTM、GRU

# Hierarchical Softmax $O(v)$

## Huffman Tree (哈夫曼树)

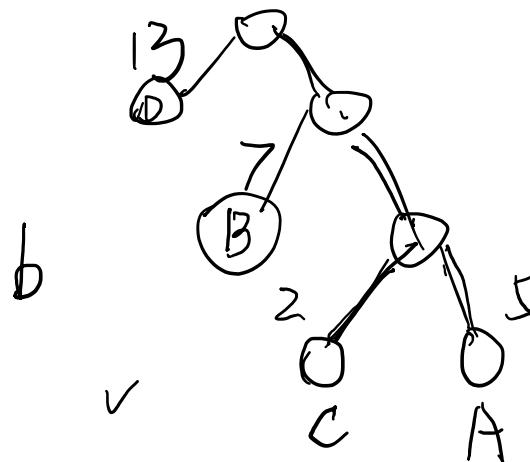
最优二叉树



$$\text{Path\_a} = 5 \times 2 + 7 \times 2 + 2 \times 2 + 13 \times 2 = 54$$

$$\text{Path\_b} = 5 \times 3 + 2 \times 3 + 7 \times 2 + 13 \times 1 = 48$$

$$P_a > P_b$$

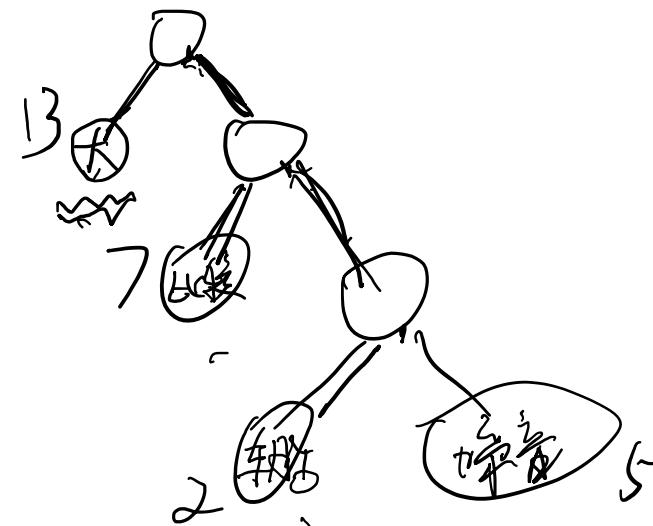


# Hierarchical Softmax

## Huffman Tree (哈夫曼树)

文本：轮胎 噪音 大 ...

词频 2 5 7 13



左 0 右 1

大: 0

~~比~~: 10

轮胎: 110

噪音: 111

1) V-1个中间节点

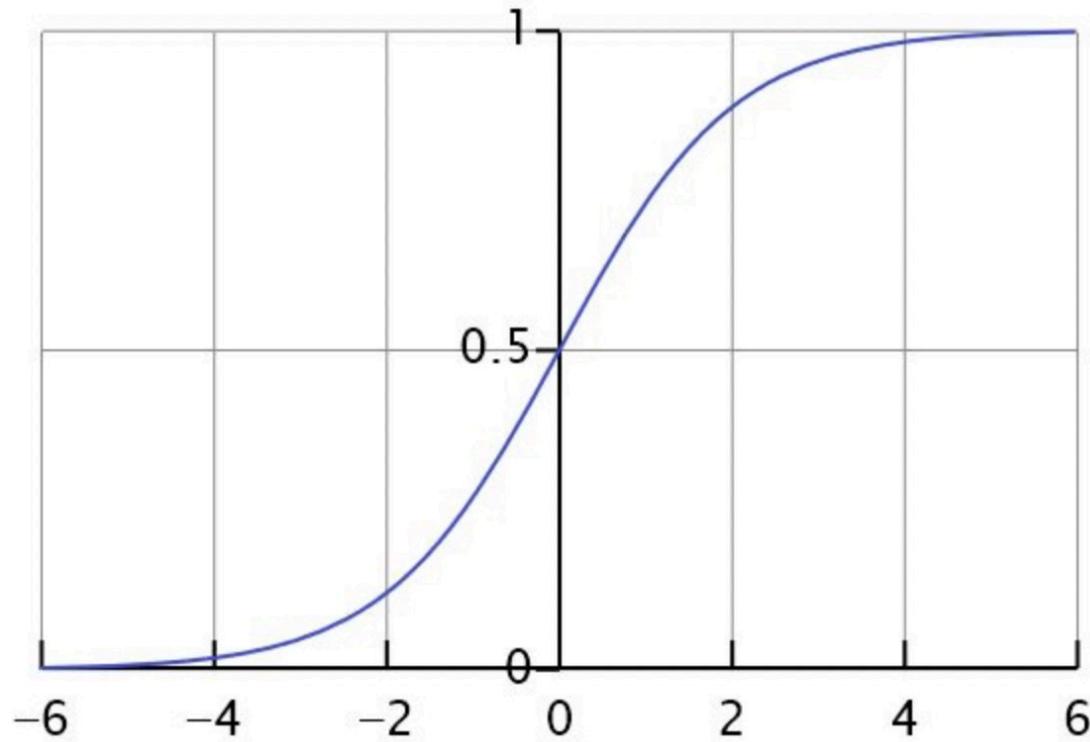
V个叶节点

2) -- 对应

3)

# Hierarchical Softmax

## Logistic Regression



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

包围  
预测/排序

IW

100

1-100

LR.

GBD<sup>T</sup> XGBoost

$$\hat{y} = P(y=1|x) \quad 0 \leq \hat{y} \leq 1$$

$$\hat{y} = \underline{\omega^T x + b} \quad \omega^T x = \sum_{i=1}^n w_i x_i$$

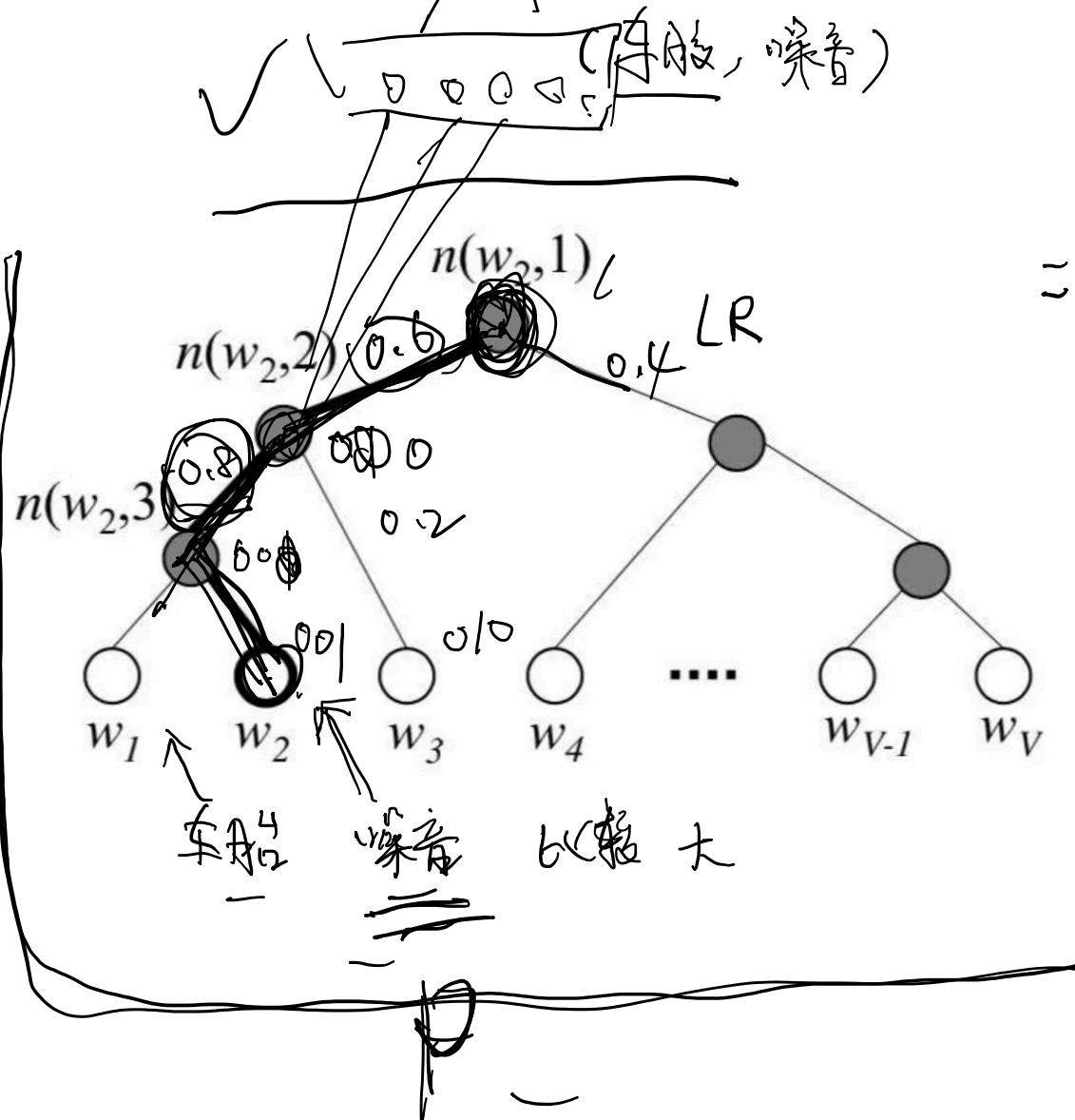
$$= w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$$= \frac{1}{1 + e^{-(\omega^T x + b)}}$$

$$\text{Loss: } P(y|x) = \begin{cases} \hat{y} & y=1 \\ 1-\hat{y} & y=0 \end{cases} \quad y=1 = \hat{y} (1-\hat{y}) (1-y)$$
$$\log(P(y|x))$$

# Hierarchical Softmax



$p_w$

$$p(w | \text{context}(w))$$

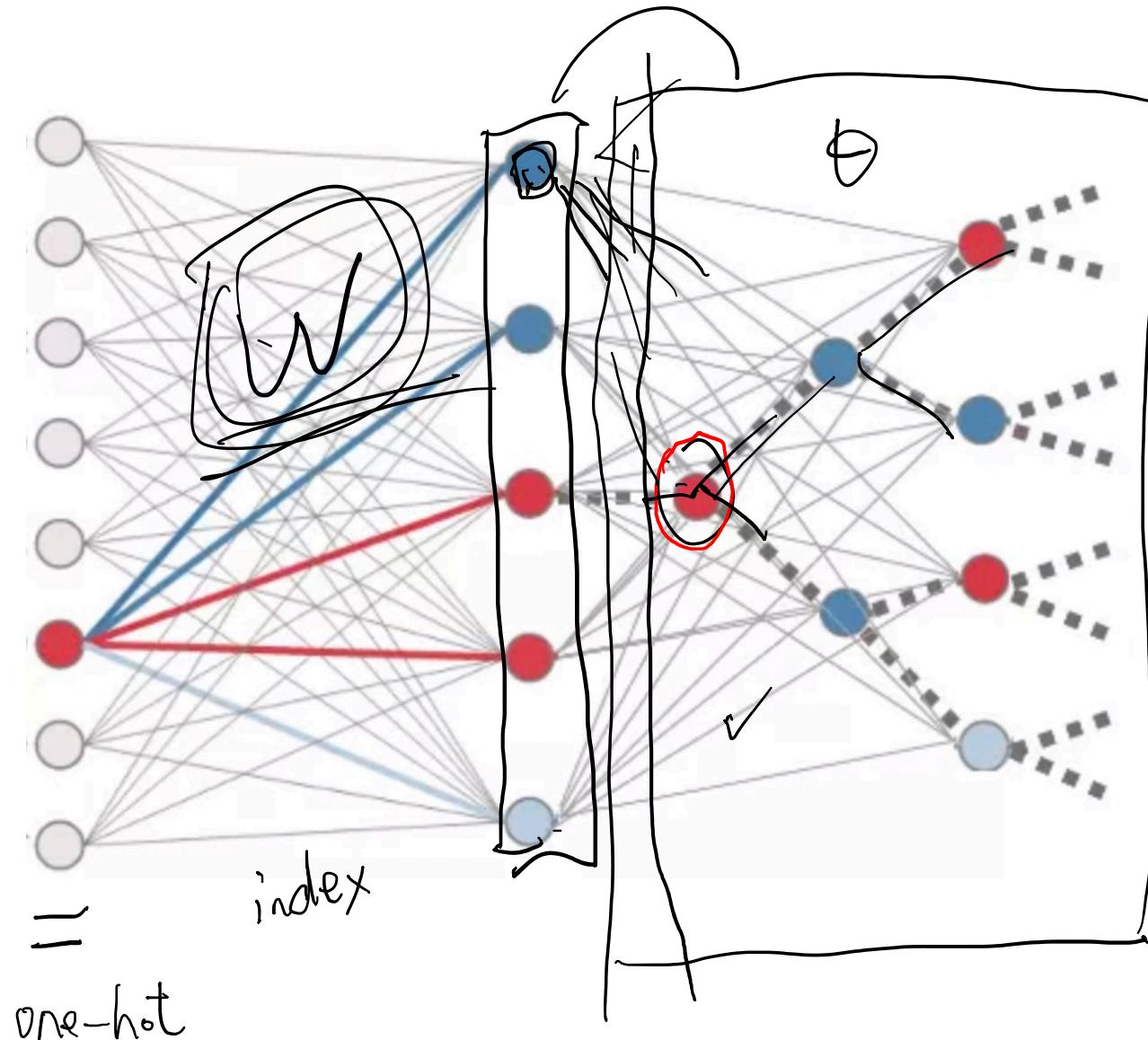
$$= \prod_{j=2}^W p(d_j^w | x_w, \theta_{j-1}^w)$$

$$p(d_j^w | x_w, \theta_{j-1}^w) = \begin{cases} \sigma(\underbrace{x_w^\top \theta_{j-1}^w}_{\parallel}) & d_j^w = 0 \\ 1 - \sigma(x_w^\top \theta_{j-1}^w) & d_j^w = 1 \end{cases}$$

$$J(w, j) = (1 - d_j^w) \cdot \log[\sigma(x_w^\top \theta_{j-1}^w)]$$

$$+ d_j^w \log[1 - \sigma(x_w^\top \theta_{j-1}^w)]$$

# Hierarchical Softmax



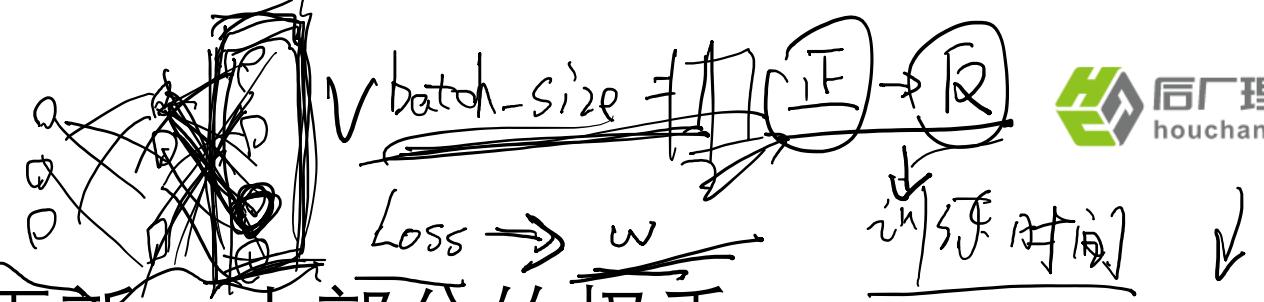
$\log V$

①  $O(|V|) \rightarrow O(\log |V|)$

② 高级 路径更短，计算更快

# Negative Sampling

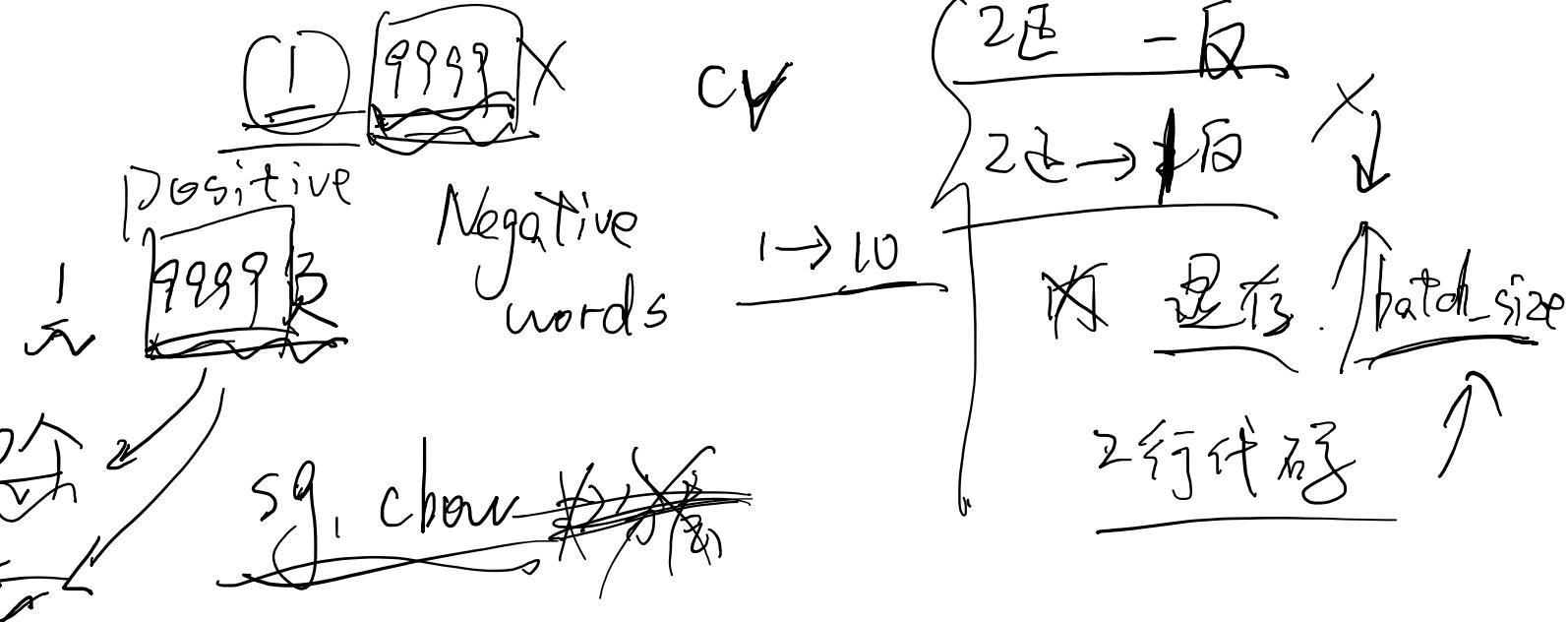
每次让一个训练样本仅仅更新一小部分的权重



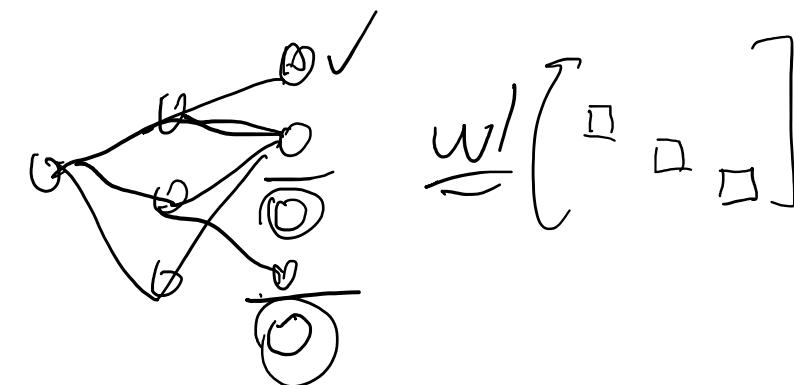
车船/噪音/比值/大

(车船, 噪音)  $V = 10000$

2.13 小规模  
大

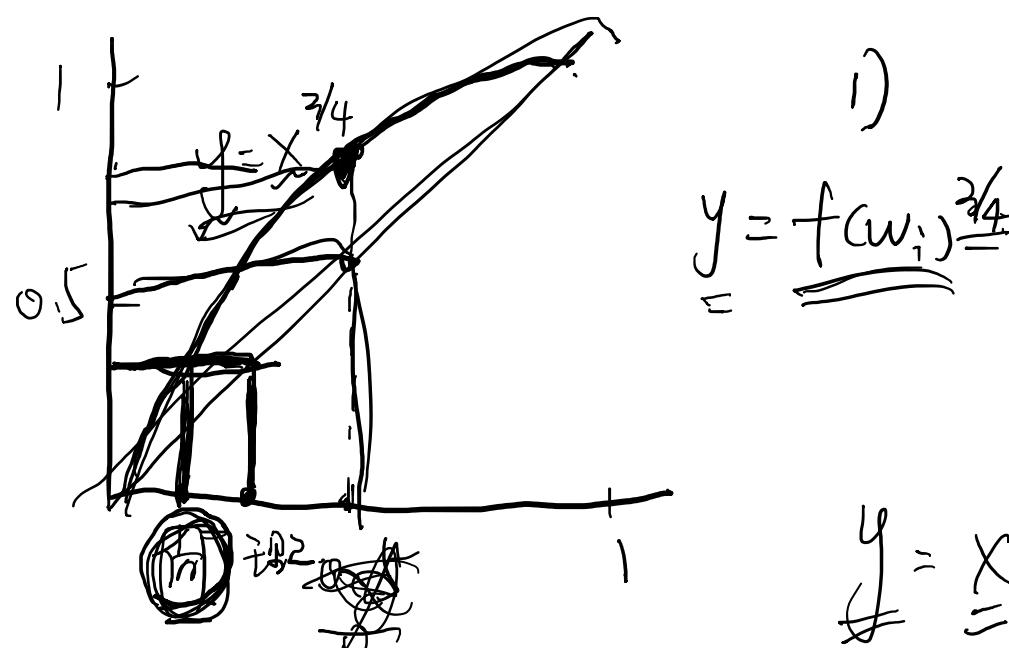


$$\downarrow \begin{matrix} 300 \times 10000 \\ 6 \times 300 \end{matrix} \rightarrow 0.06\%$$



# Negative Sampling

$$p(w_i) = \frac{f(w_i)}{\sum_{j=0}^n f(w_j)^{3/4}}$$



$$y = f(w_i)^{3/4}$$

$$\begin{aligned} y &= X^{3/4} \\ y &= 1 \end{aligned}$$

~~负采样~~: ~~负采样~~

9:10

word2vec → sg. cbow

$$\frac{5}{100}$$

$$\frac{1}{100}$$

$$y = X$$

$$V = [0 \dots 0]$$

↓

2~5

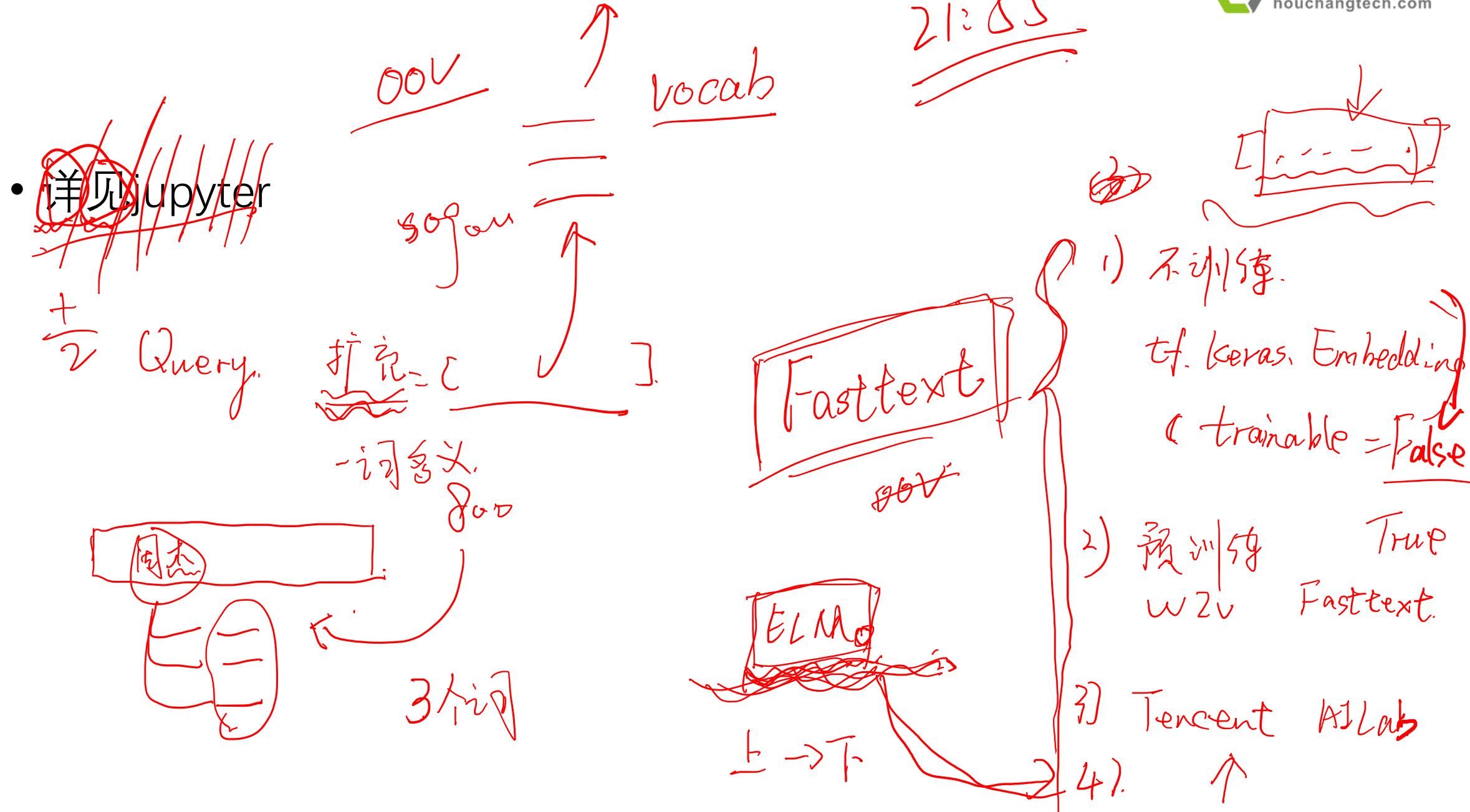
5~20

word2vec trick. 2~10↑  
高频率的正采样

$$p(w_i) = \frac{1}{\sqrt{\text{freq}(w_i)}} \approx 10^{-5} \sim 10^{-2}$$

# Outline

- 词向量计算两种优化方法
- 词向量在工程中的具体实现
- RNN递归神经网络结构
- RNN、LSTM、GRU

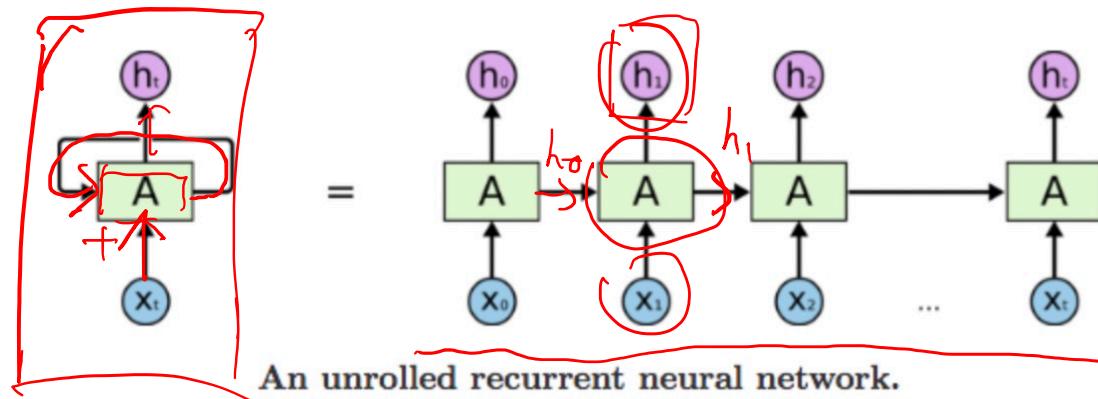


# Outline

- 词向量计算两种优化方法
- 词向量在工程中的具体实现
- RNN递归神经网络结构
- RNN、LSTM、GRU

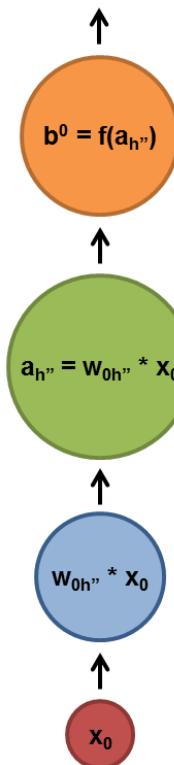
# RNN

## Recurrent Neural Network



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

$b^0$  is fed to next layer

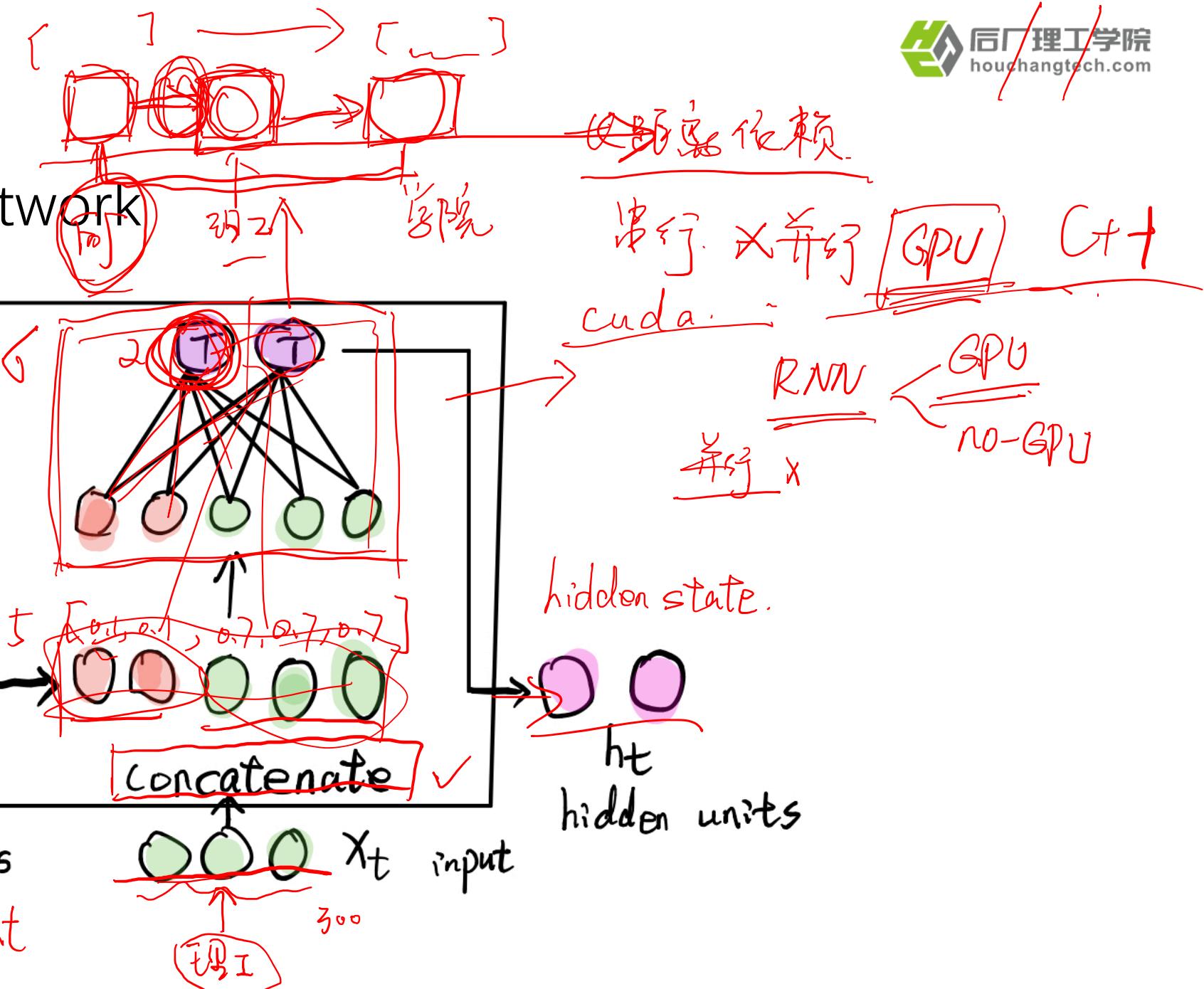
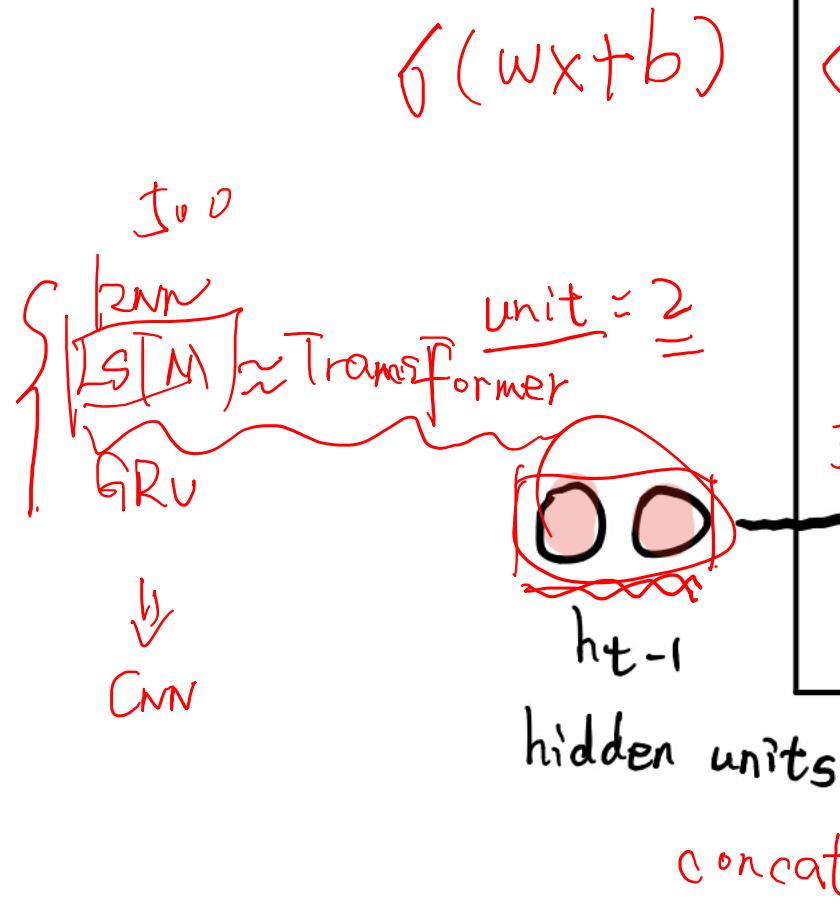


# RNN

## Recurrent Neural Network

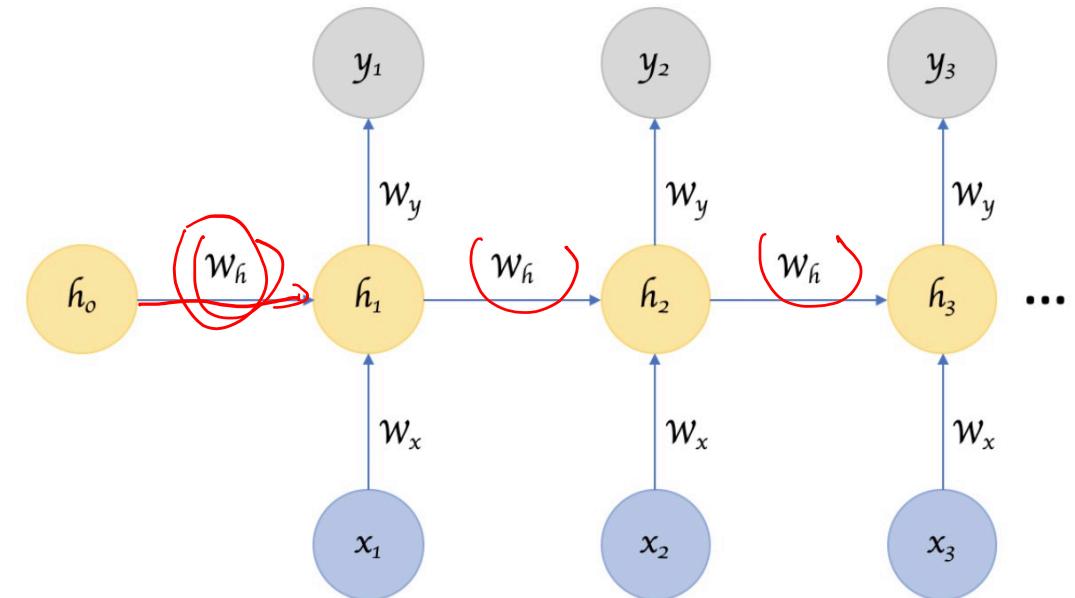
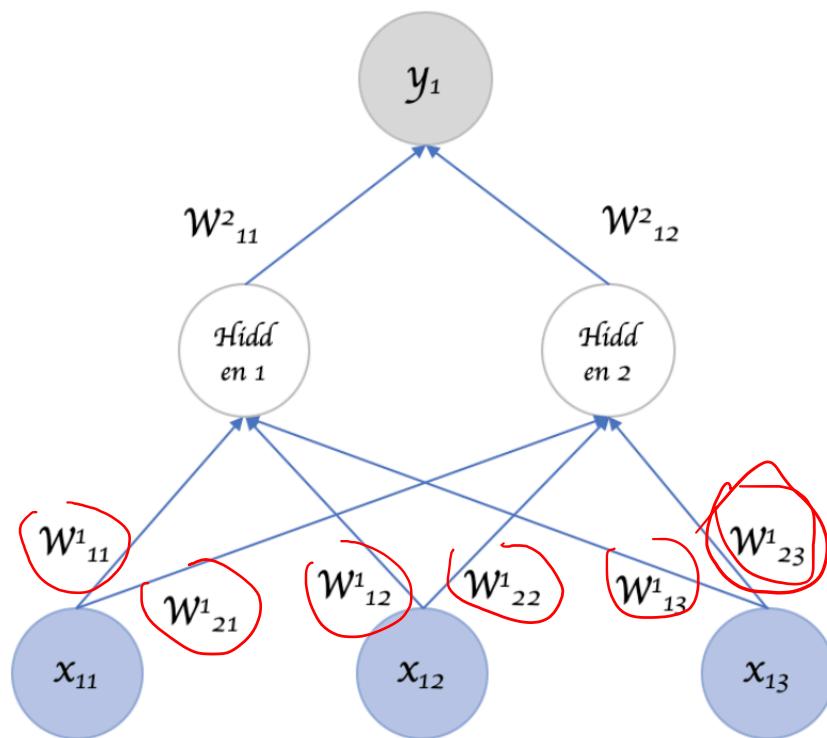
# RNN

## Recurrent Neural Network



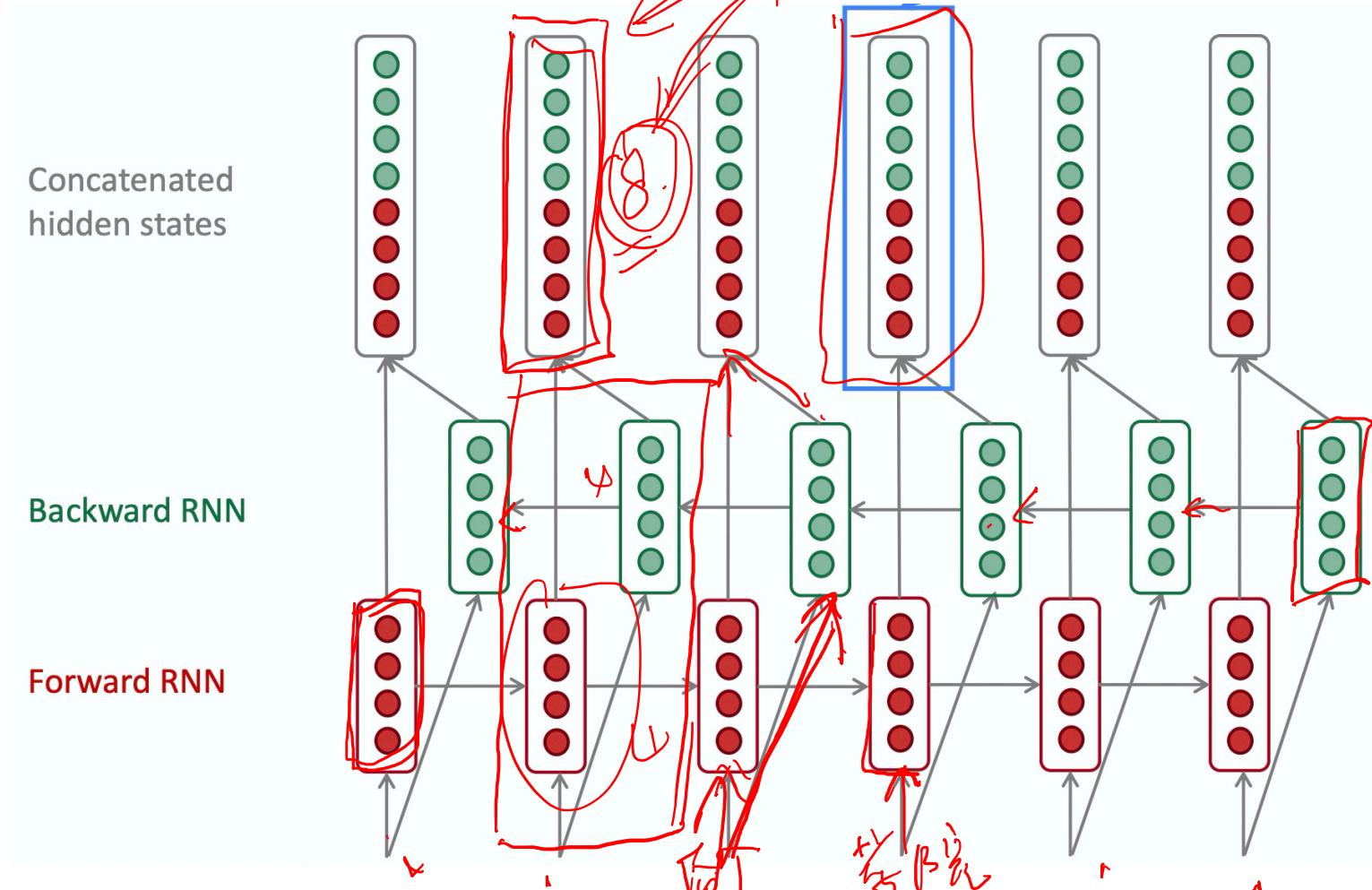
# RNN

parameter sharing ✓



# RNN

## Bidirectional RNNs



$$h^t = [\vec{h}^{(t)}, \overleftarrow{h}^{(t)}] \quad \text{concat.}$$

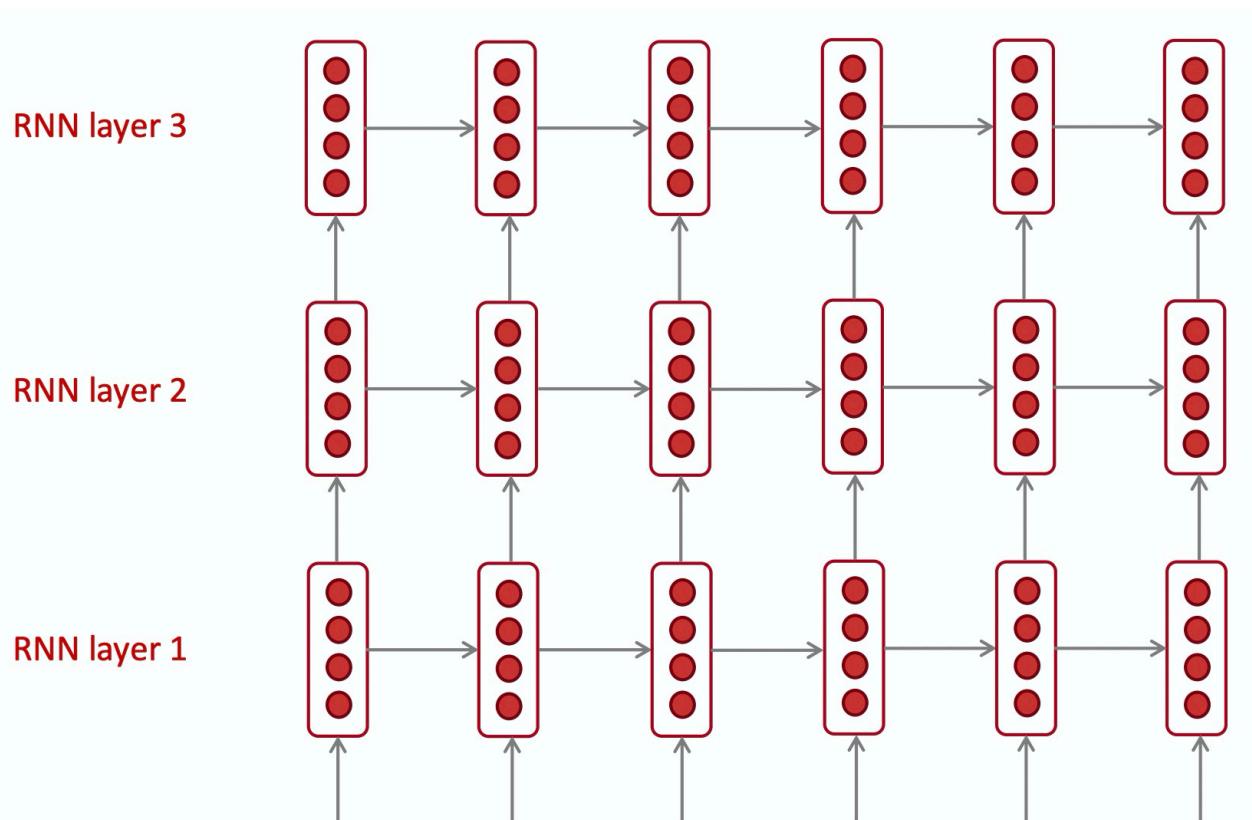
$\overleftarrow{h}^{(t)}$

$$\overleftarrow{h}^{(t)} = \text{RNN}_{BW}(\overleftarrow{h}^{(t+1)}, x^{(t)})$$

$$\vec{h}^{(t)} = \text{RNN}_{Fw}(\vec{h}^{(t-1)}, x^{(t)})$$

# RNN

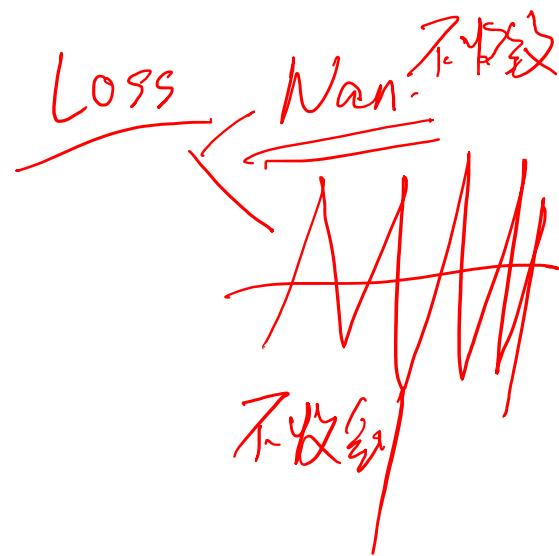
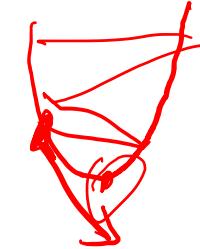
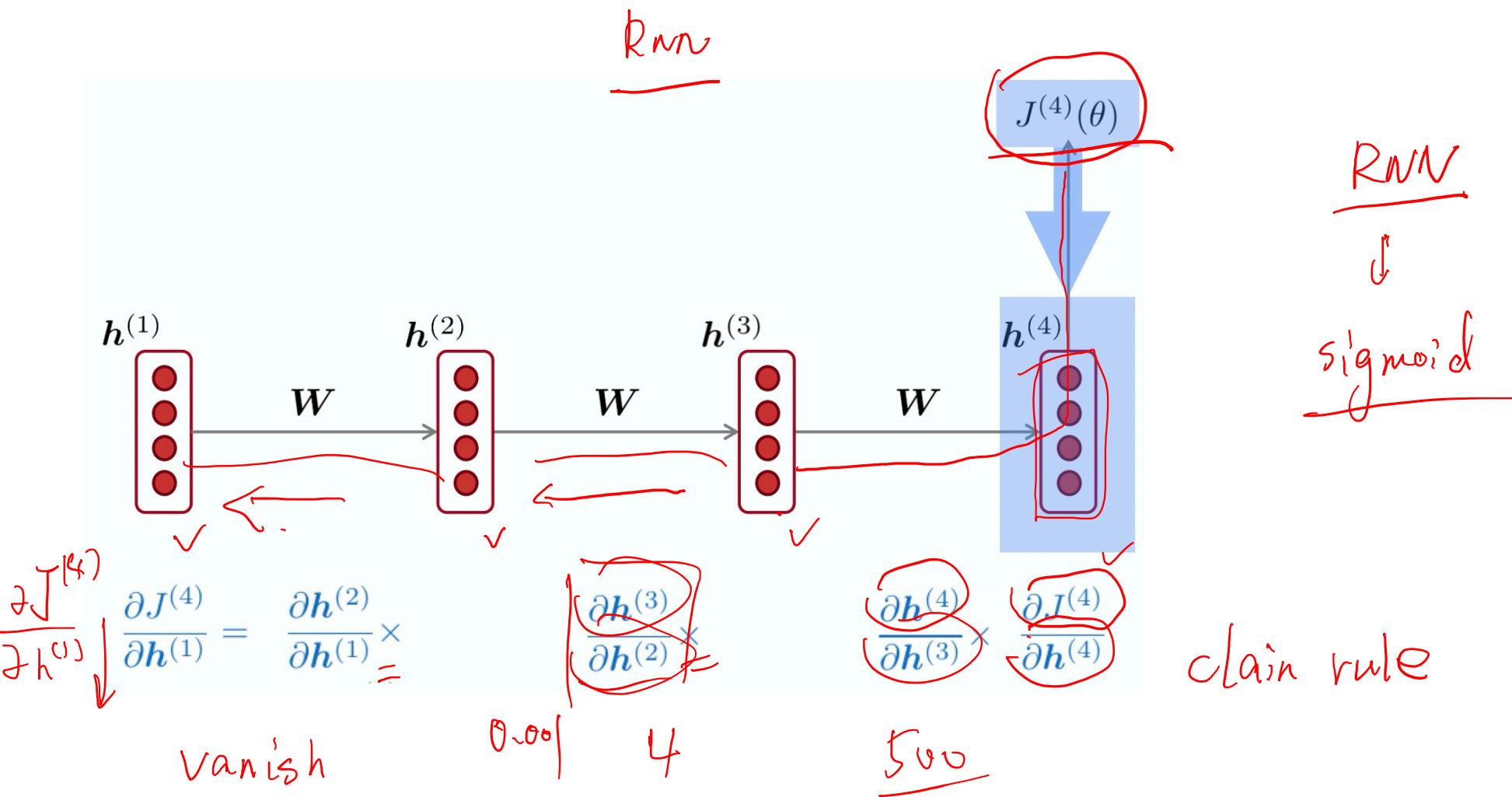
## Deep RNNs



NMT  $\rightarrow$  2-4 layers  
encoder  
decoder  
4 layers

文本顺序      单向    双向  
IRNN    GRU    LSTM

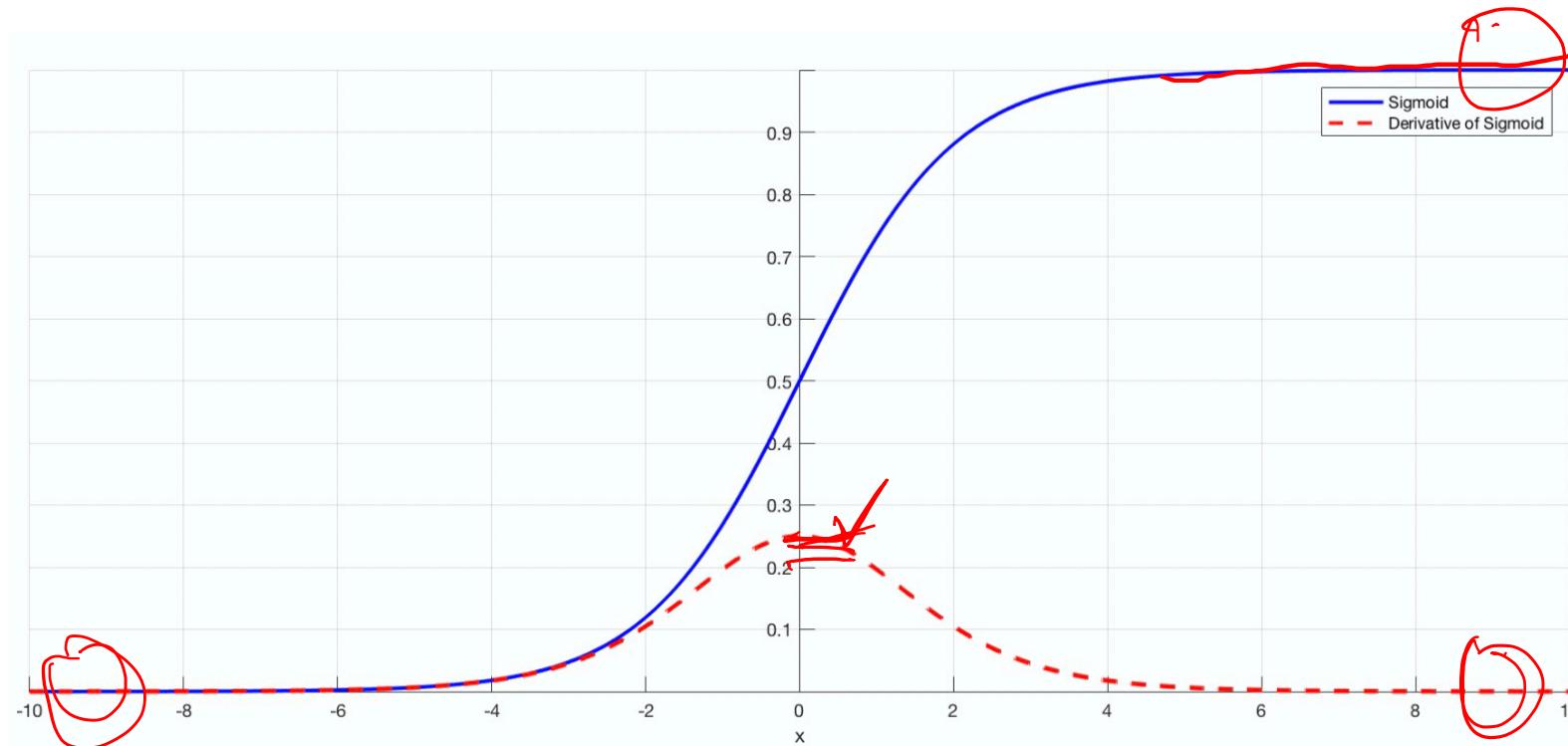
# Exploding and Vanishing Gradient Problem



# Exploding and Vanishing Gradient Problem

$$\sigma(wx+b)$$

$$[0.000] \quad x$$



# Exploding Gradient

gradient clipping

# Vanishing Gradient

Identity Initialization

LSTM

Residual Networks

Batch Normalization

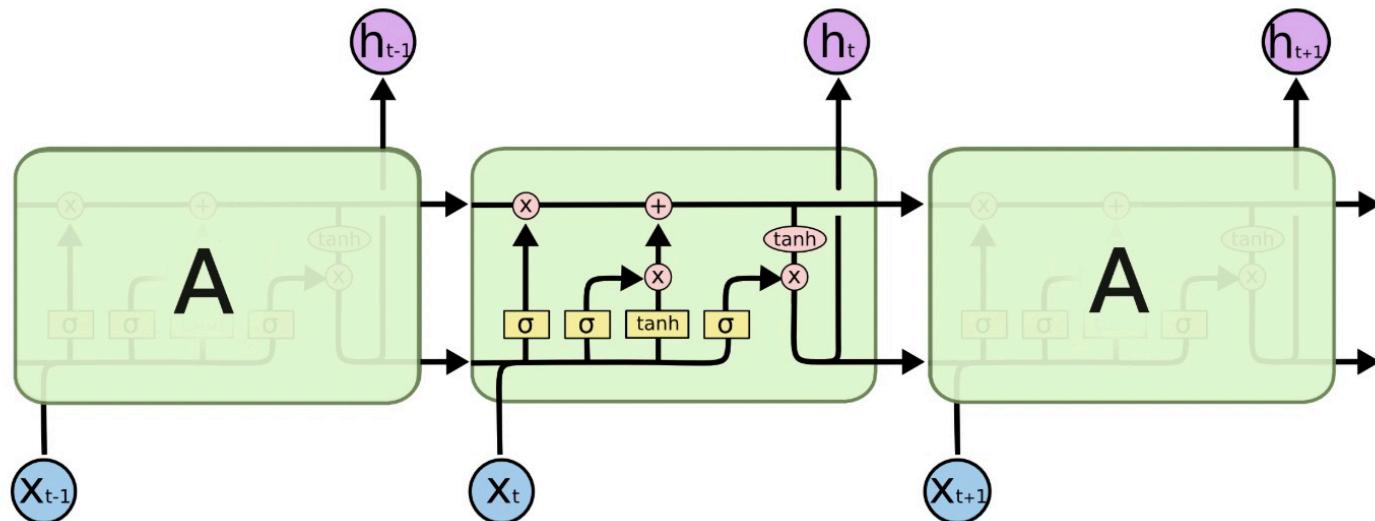
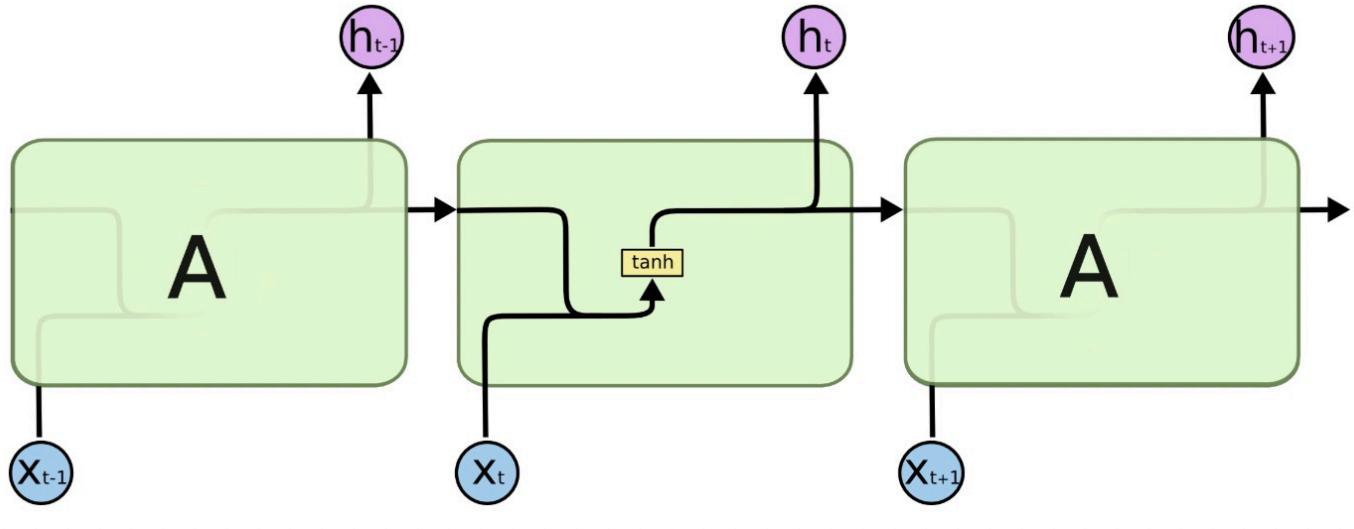
# Outline

- 词向量计算两种优化方法
- 词向量在工程中的具体实现
- RNN递归神经网络结构
- RNN、LSTM、GRU

# LSTM

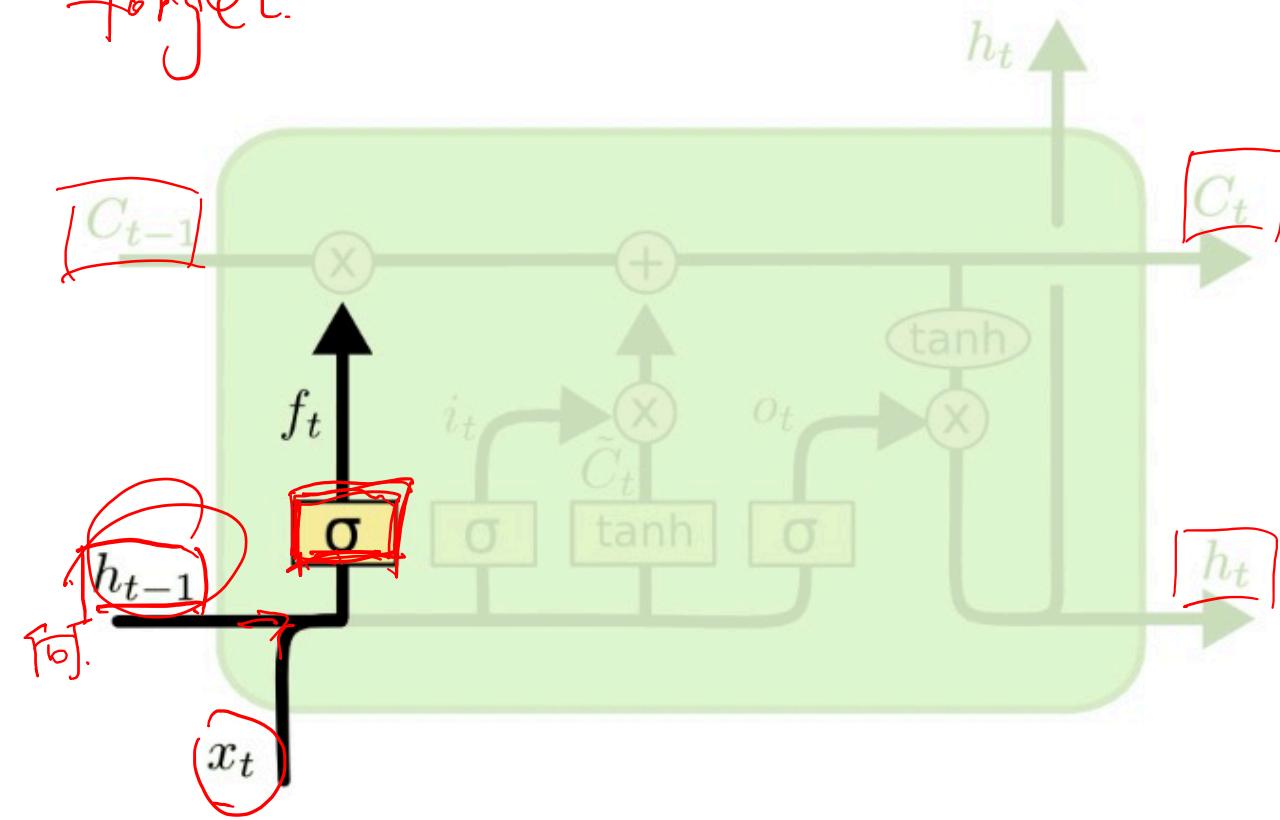
## Long Short Term Memory

□  
○  
→  
⤒ concat  
⤒ copy



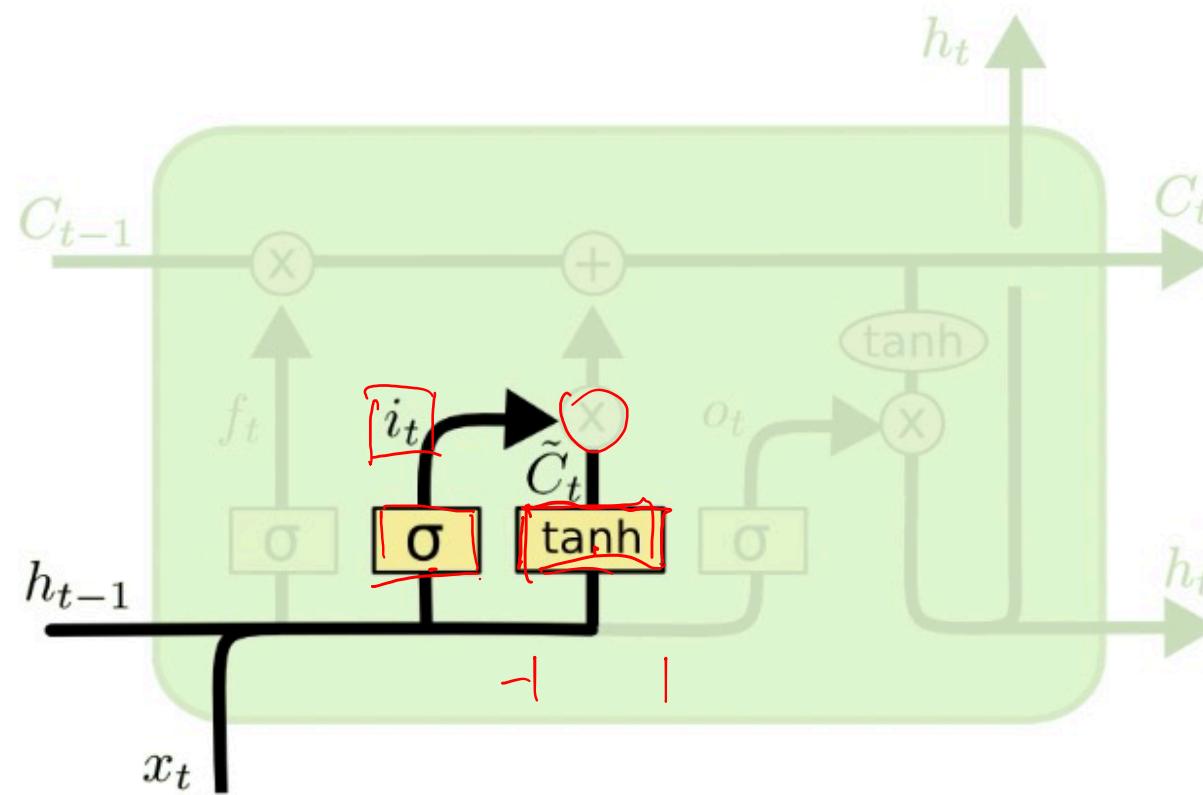
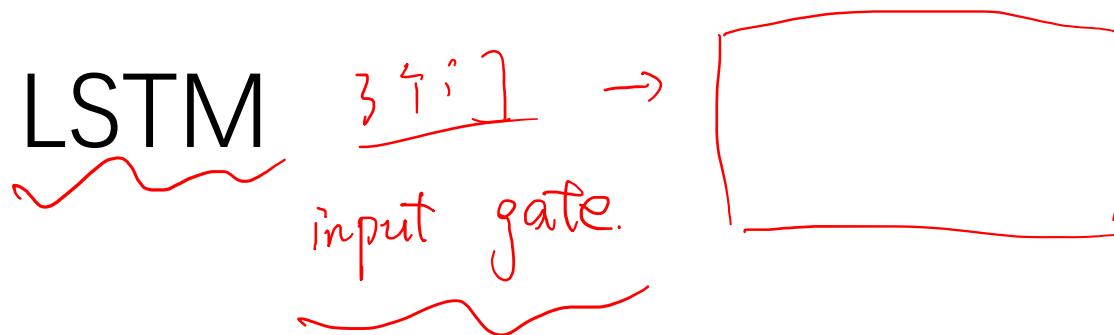
# LSTM

~~forget~~ gate.  
forget.



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

Sigmoid

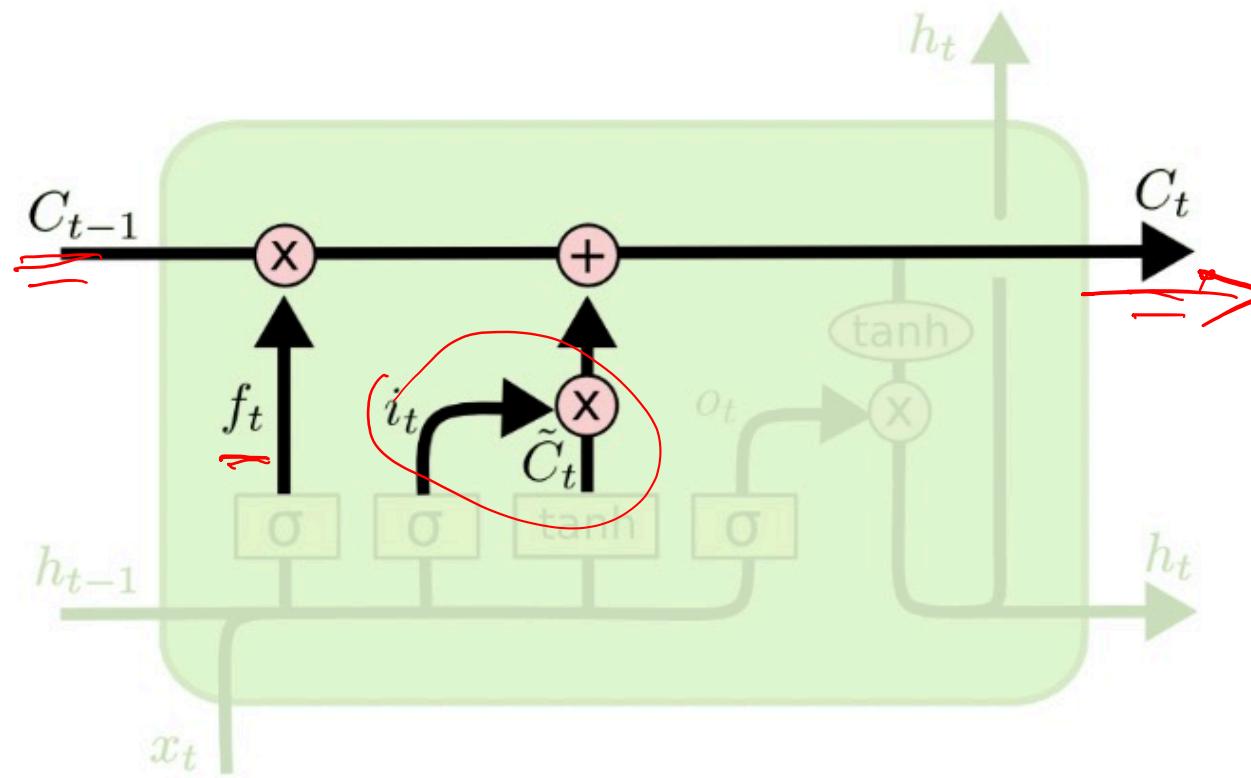


$$\hat{i}_t = \sigma(w_i[h_{t-1}, x_t] + b_i)$$

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

# LSTM

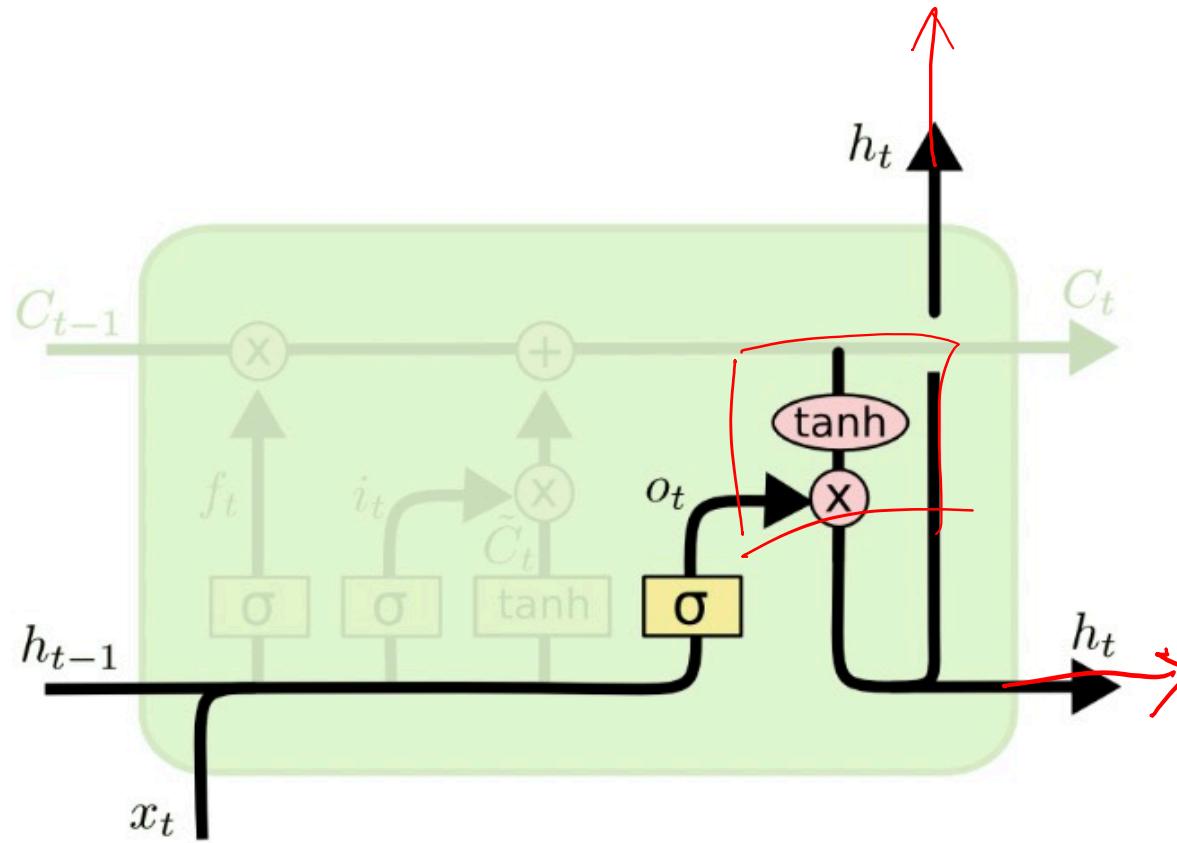
*Memory*



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

# LSTM

Output



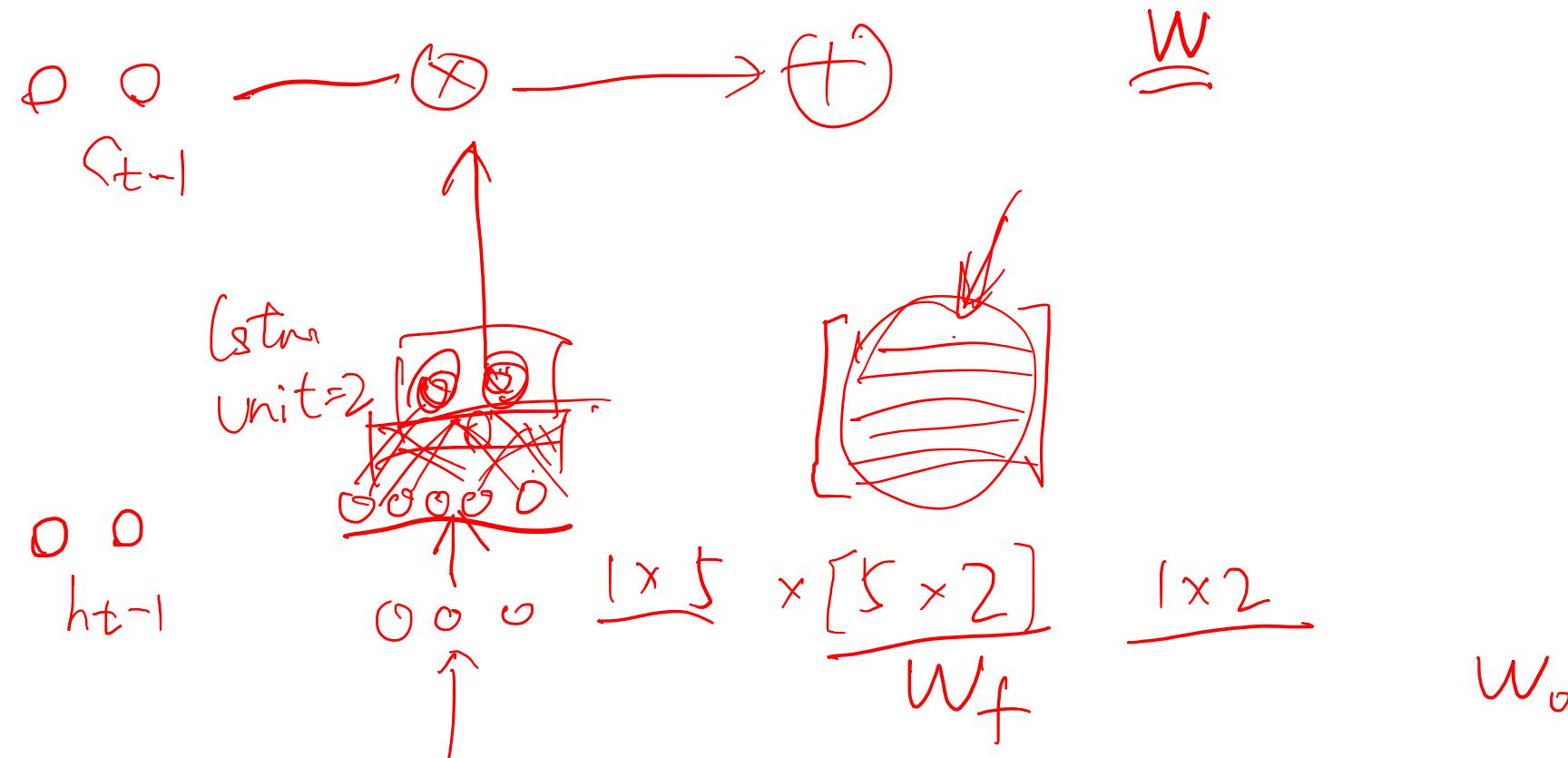
$$O_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$$

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

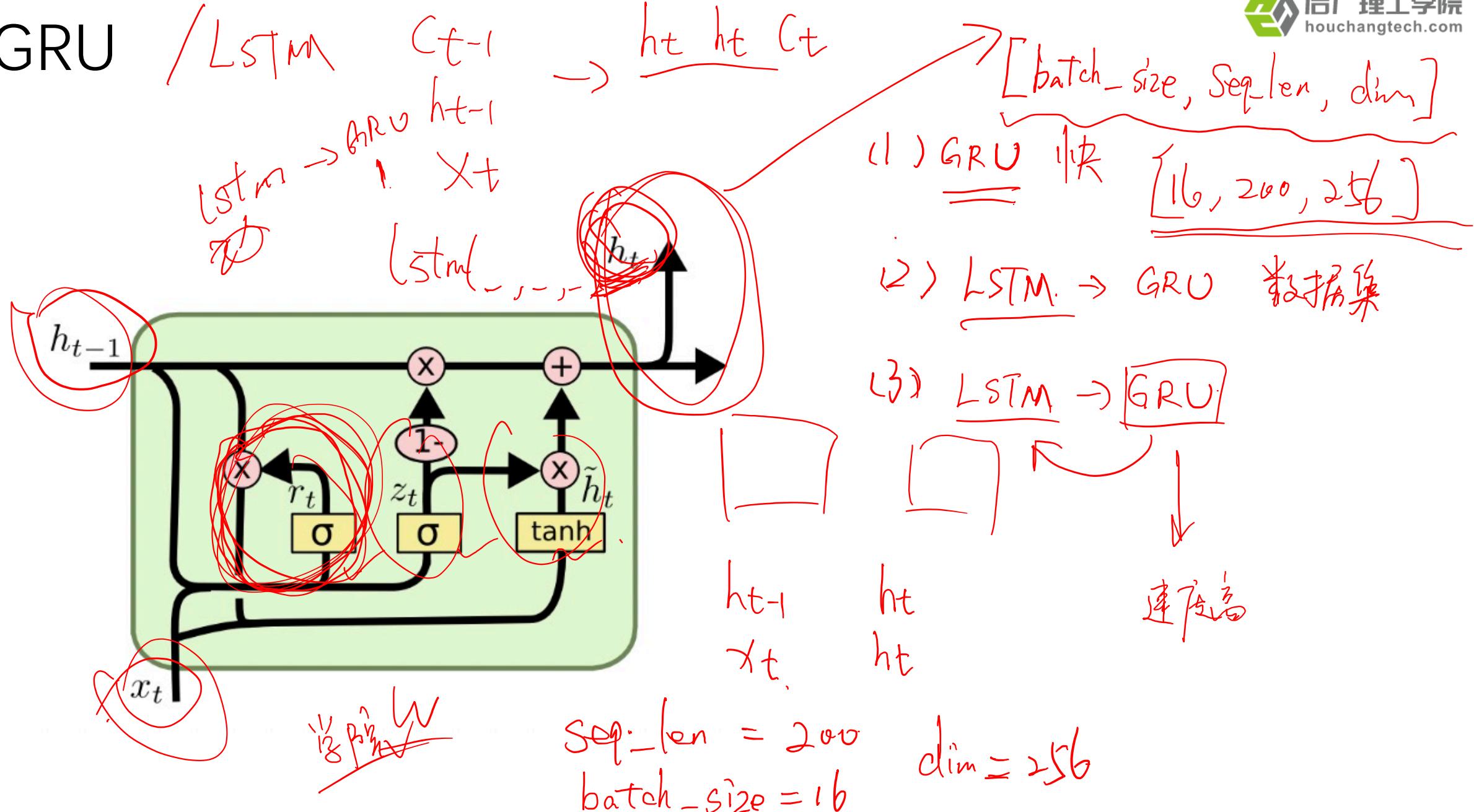
$$2 h_t C_t$$

$$\underline{\underline{h_t, h_t, O_t}} = \underline{\underline{\text{LSTM}(\quad)}}$$

# LSTM



# GRU



# 作业

Bye !