

HCT NLP Week 3

问答摘要与推理
Seq2Seq（一）

Outline

- Encoder-Decoder结构
- Attention机制
- 模型Layer、Model构建
- Seq2Seq训练

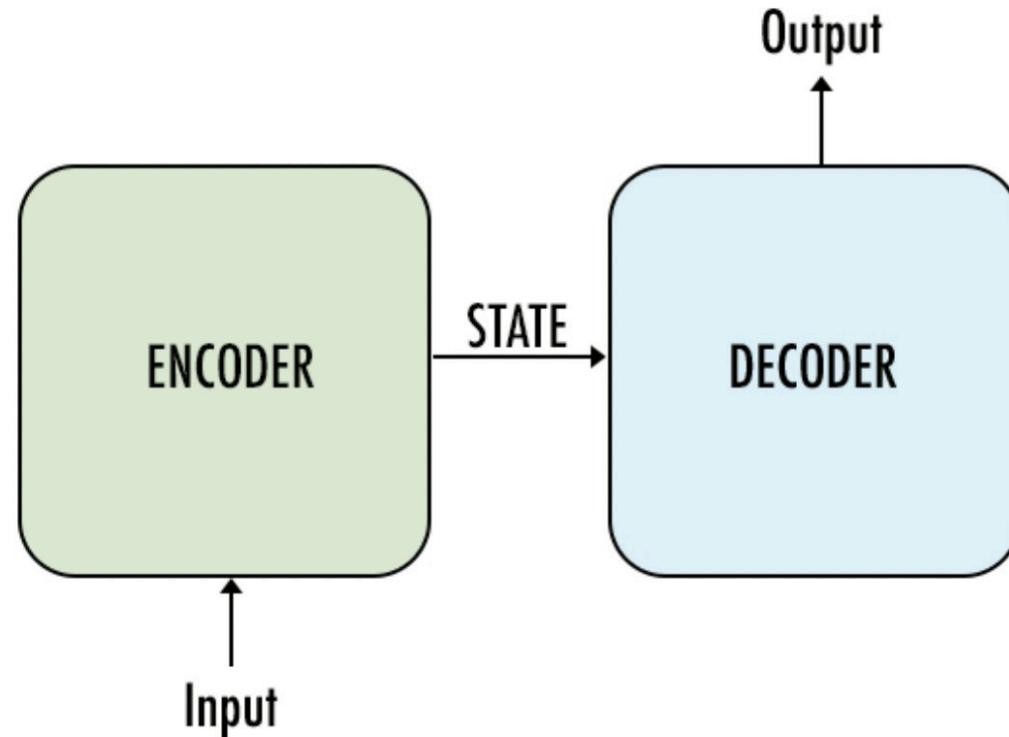
Outline

- Encoder-Decoder结构
- Attention机制
- 模型Layer、Model构建
- Seq2Seq训练

Seq2Seq

Encoder-Decoder结构

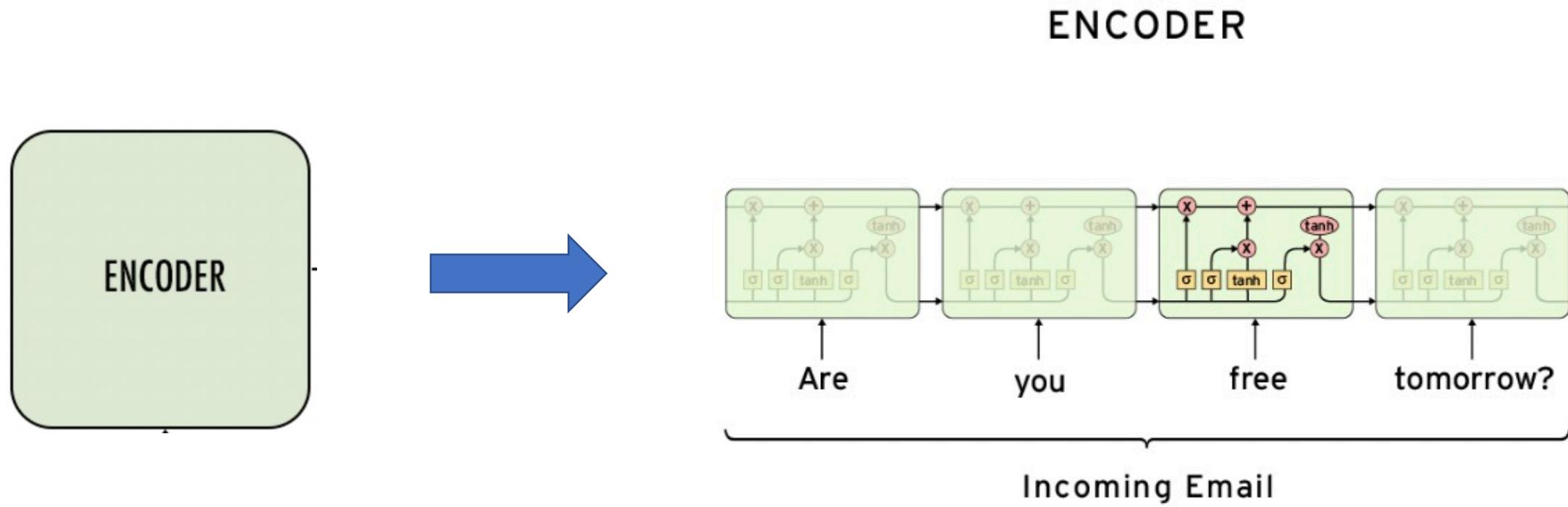
Sequence to sequence was first introduced by Google in 2014



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

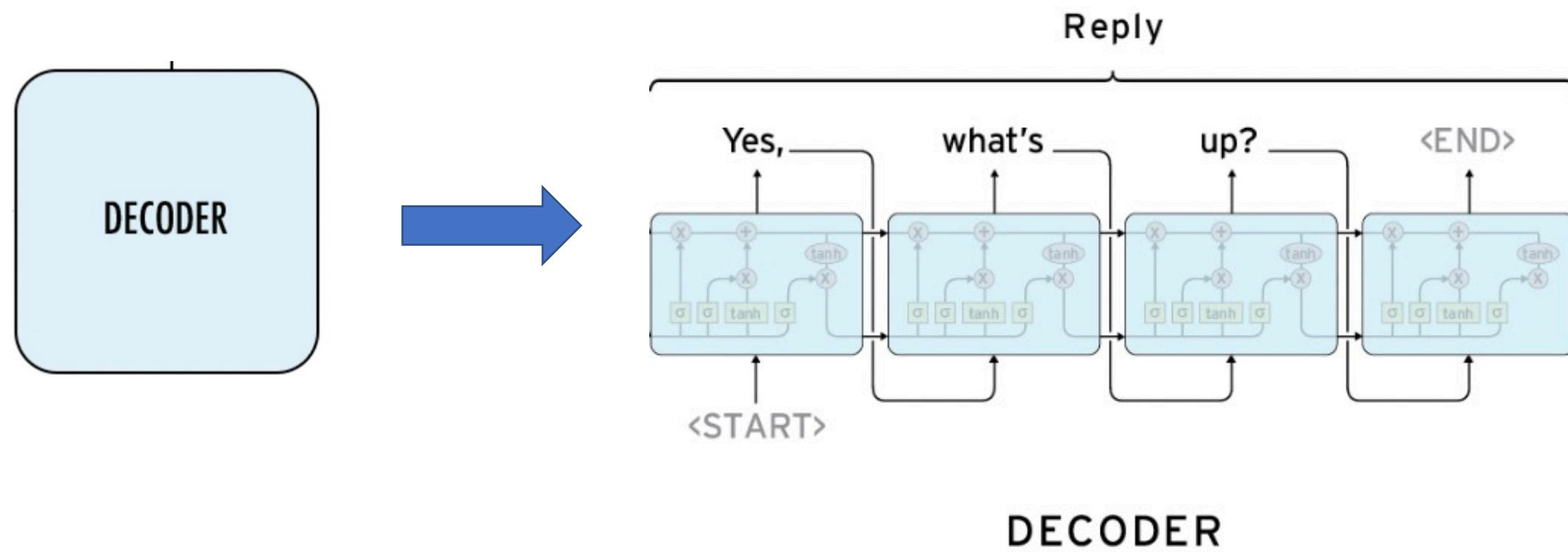
Seq2Seq

Encoder-Decoder结构



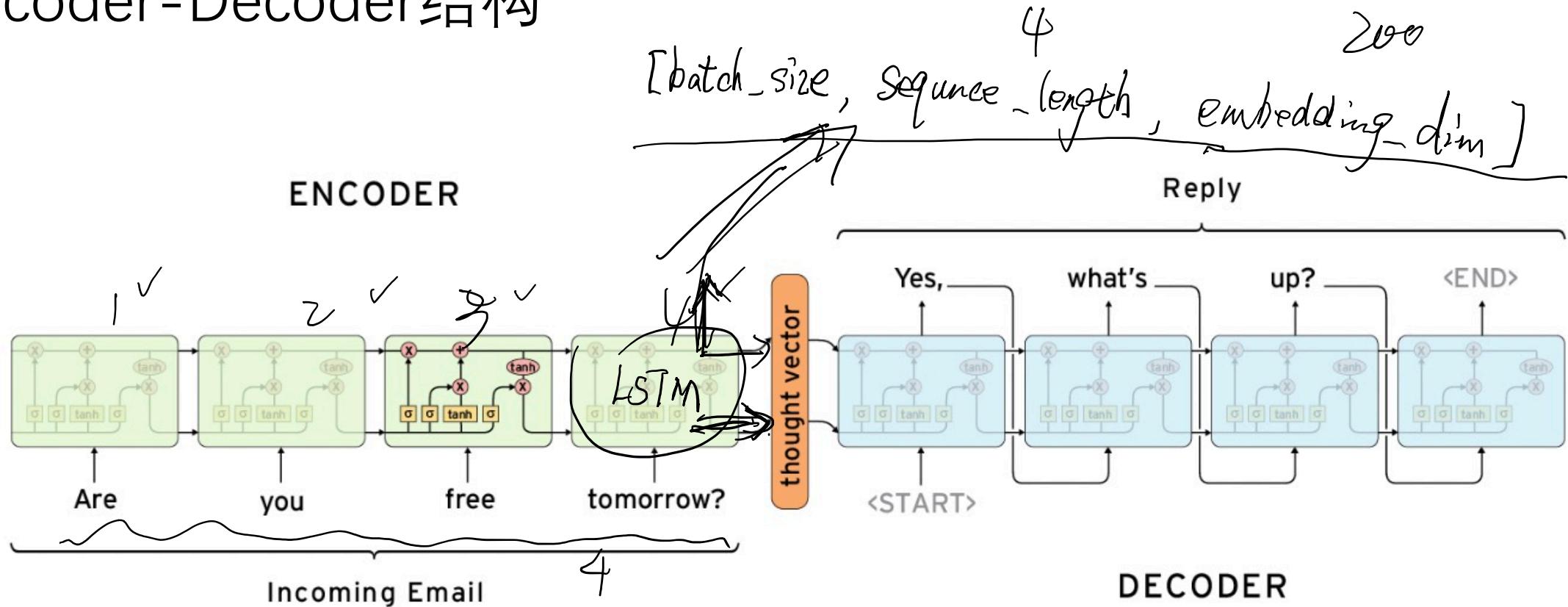
Seq2Seq

Encoder-Decoder结构



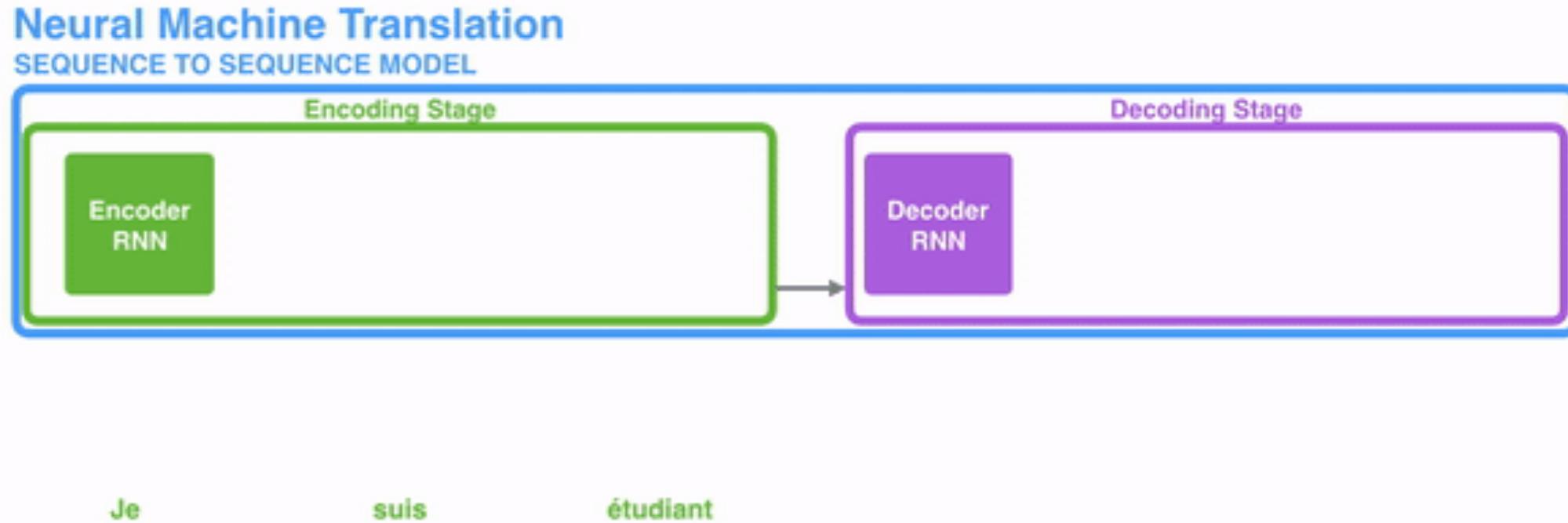
Seq2Seq

Encoder-Decoder结构



Seq2Seq

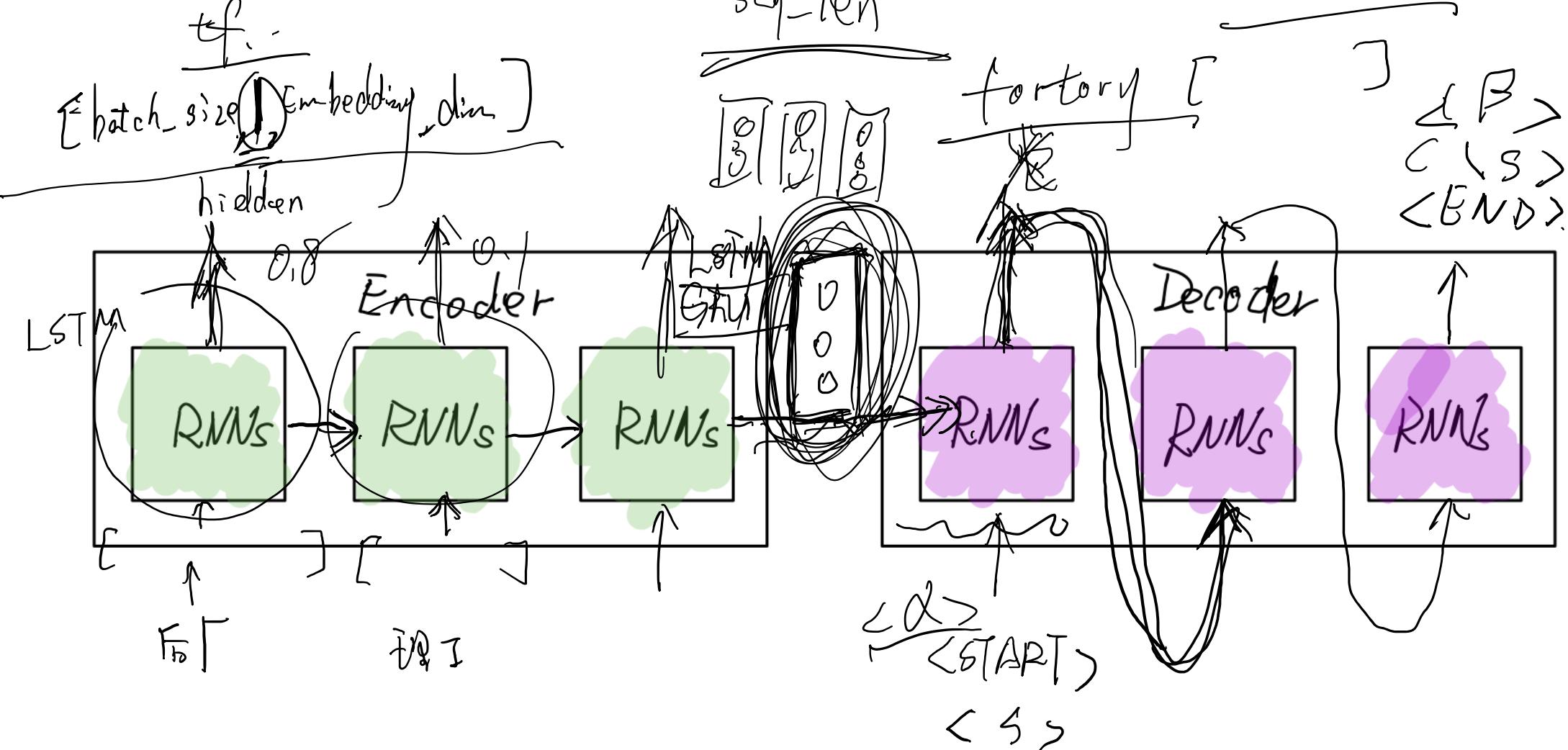
Encoder-Decoder结构



(source: [Jay. Alammar, 2018](#))

Seq2Seq

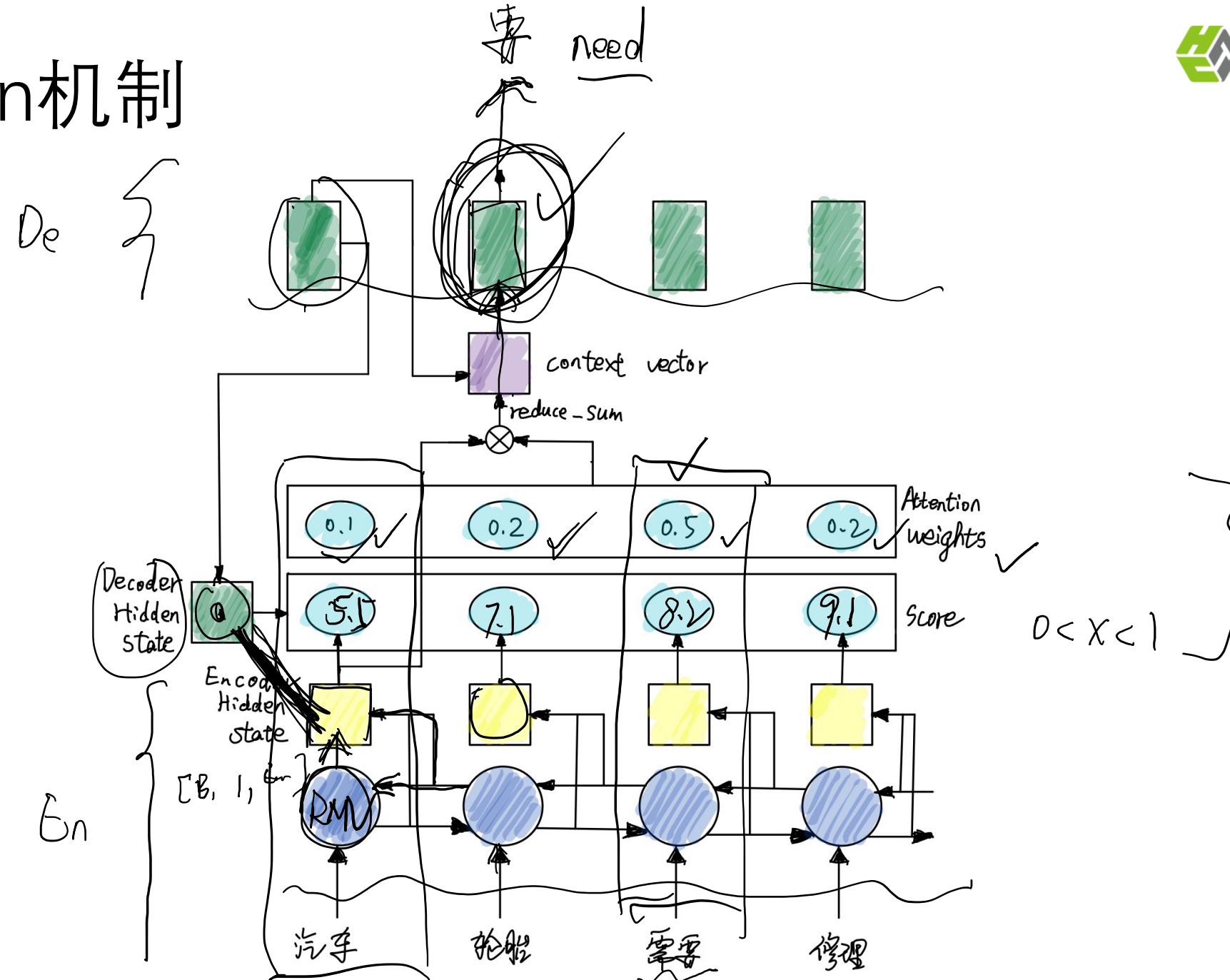
Encoder-Decoder结构



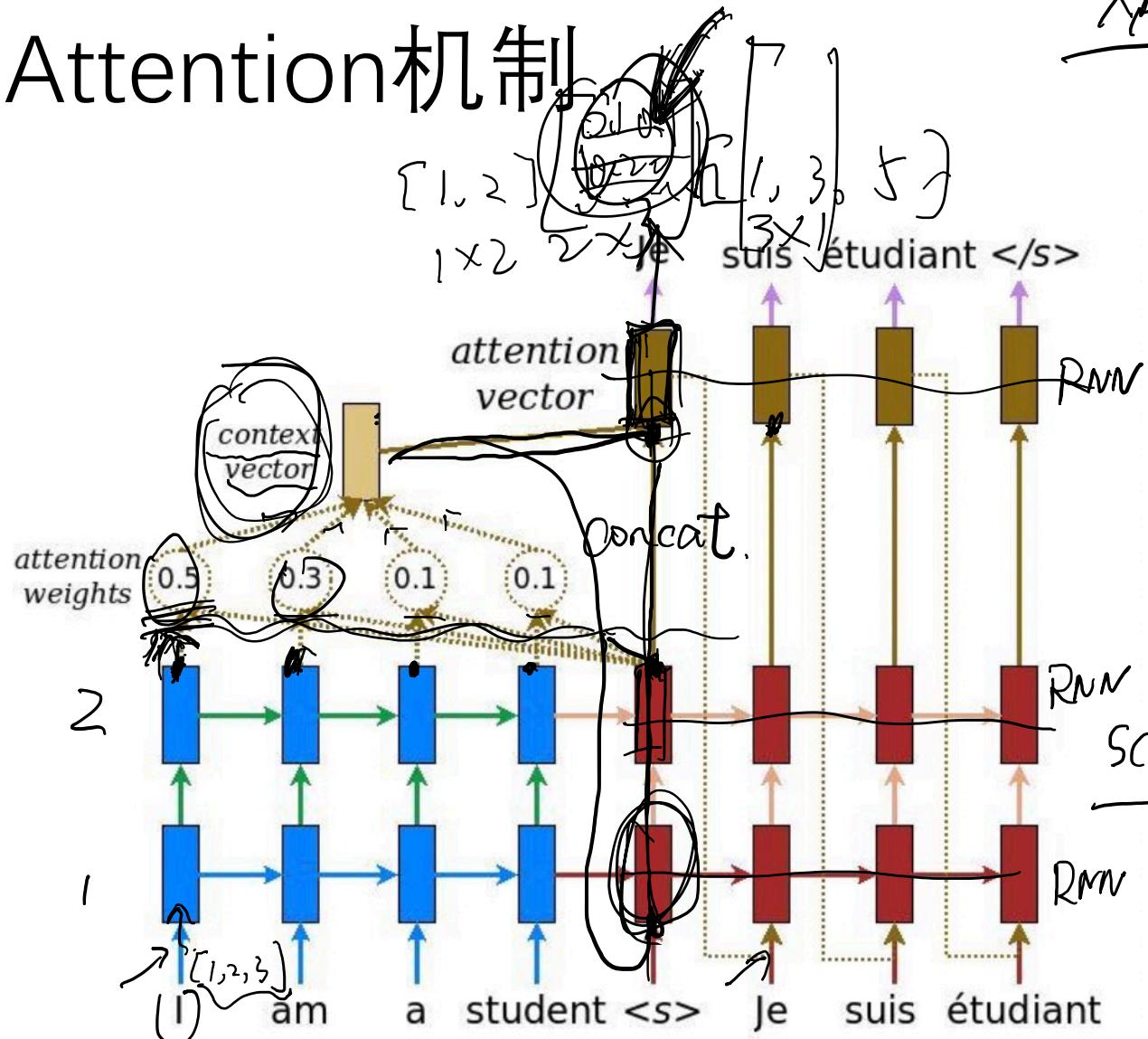
Outline

- Encoder-Decoder结构
- Attention机制
- 模型Layer、Model构建
- Seq2Seq训练

Attention机制



Attention机制



XAI

~~decoder~~

decoder hidden

encoder hidden

$$\frac{h_t}{h_s}$$

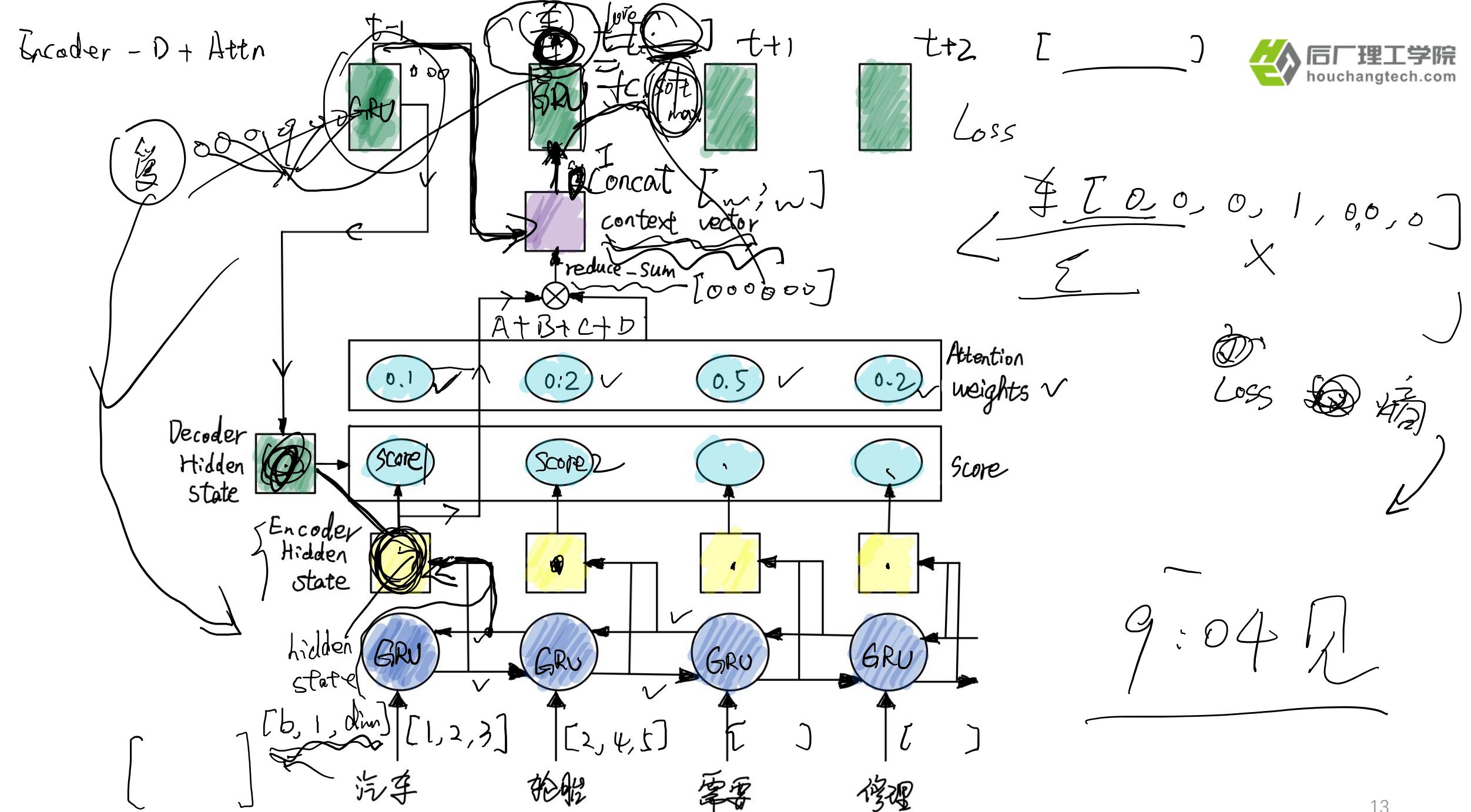


$$\alpha_{ts} = \frac{\exp(score)}{\sum_{s=1}^S \exp(score)}$$

$$score(h_t, \bar{h}_s) = \begin{cases} h_t^T W h_s & [\text{Luong's multiplicative style}] \\ \frac{1}{\pi} \tanh(w_1^T h_t + w_2^T \bar{h}_s) & [\text{Bahdanau's additive style}] \end{cases}$$

$$\text{Context} = \sum \alpha_{ts} \bar{h}_s = 0.5 \bar{h}_{s1} + 0.3 \times \bar{h}_{s2} + \dots +$$

[Bahdanau's additive style]



Outline

- Encoder-Decoder结构
- Attention机制
- 模型Layer、Model构建
- Seq2Seq训练

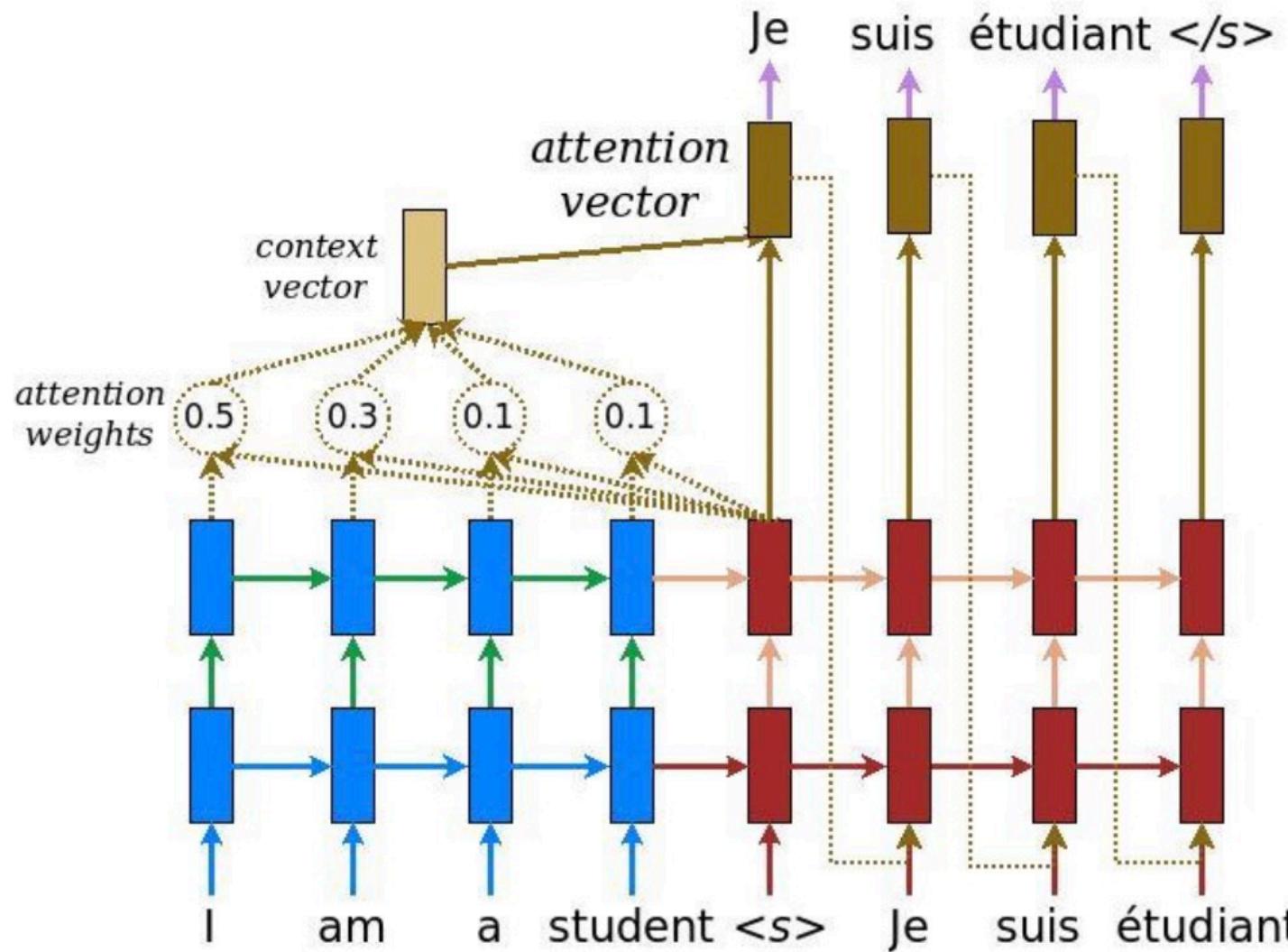
模型Layer、Model构建

代码见jupyter

Outline

- Encoder-Decoder结构
- Attention机制
- 模型Layer、Model构建
- Seq2Seq训练

Seq2Seq训练



How is it trained?

Seq2Seq训练

$[x_1, x_2, x_3]$

cross-entropy loss 交叉熵

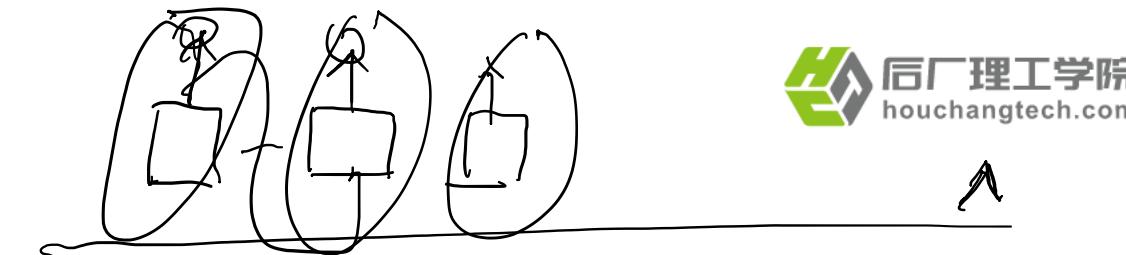
$$z_1 \xrightarrow{3} @ \rightarrow e^{z_1} = \frac{2.0}{2.0 + 2.7 + 0.05} = 0.88$$

$$z_2 \xrightarrow{1} @ \rightarrow e^{z_2} = \frac{2.7}{2.0 + 2.7 + 0.05} = 0.12$$

$$z_3 \xrightarrow{-3} @ \rightarrow e^{z_3} = \frac{0.05}{2.0 + 2.7 + 0.05} \approx 0$$

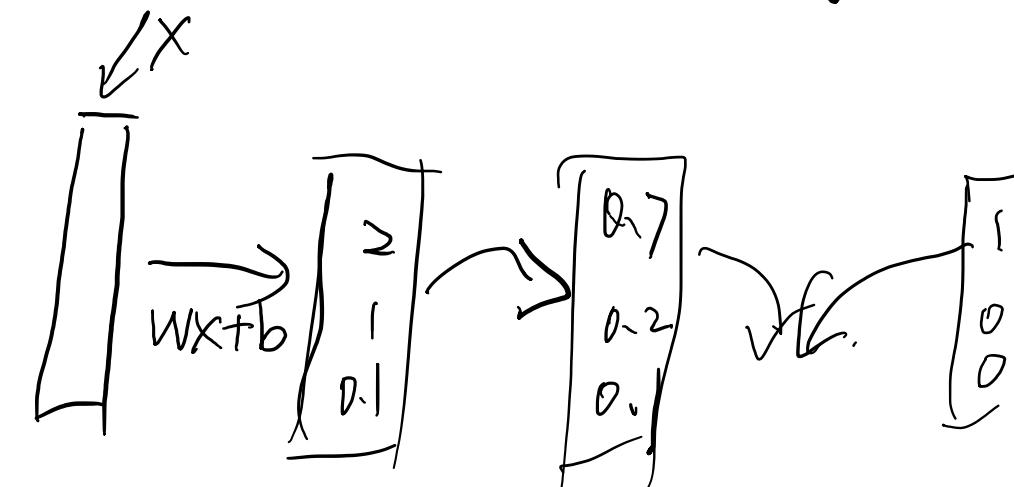
$S(y)$

$$\begin{bmatrix} 0.88 \\ 0.12 \\ 0.0 \end{bmatrix} D(S, L) = - \sum_i L_i \log(S_i) \quad \begin{bmatrix} 1.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$



$$J = \frac{1}{N} \left(\sum_{i=1}^N y_i \log(\hat{y}_i) \right)$$

$$J(w) = -\frac{1}{N} \sum_{n=1}^N [y_n \log \hat{y}_n + (1-y_n) \log (1-\hat{y}_n)]$$



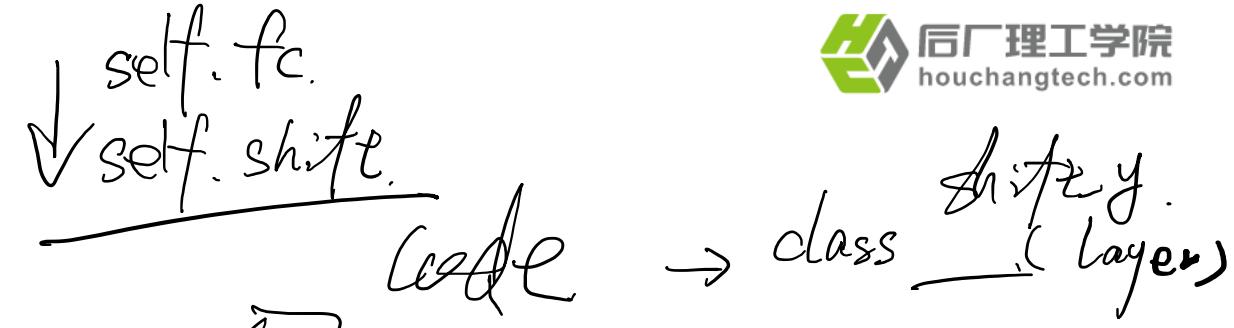
Seq2Seq训练

cross-entropy loss 交叉熵

10:1012

Seq2Seq训练

其他训练技巧-先验知识



$y \Rightarrow \text{softmax}(y) \rightarrow P_i = \frac{e^{y_i}}{\sum e^{y_i}}$
 先验知识, (IV) o/i 向量 $X = f(x_1, x_2, \dots, x_{IV}) \rightarrow$

~~def call~~
 $y = s_i x_i + t$
 $x_i = 0 \text{ no}$

~~$y = s_i x_i + t$~~
 $y = (s_1 x_1 + t_1, s_2 x_2 + t_2, \dots, s_{IV} x_{IV} + t_{IV})$
 $x_i = 1 \text{ 代表文中出现过}$

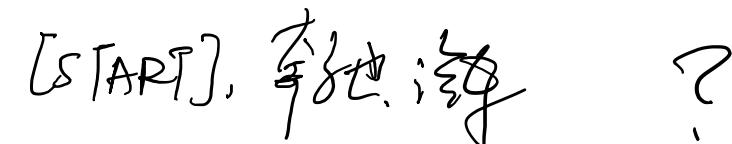
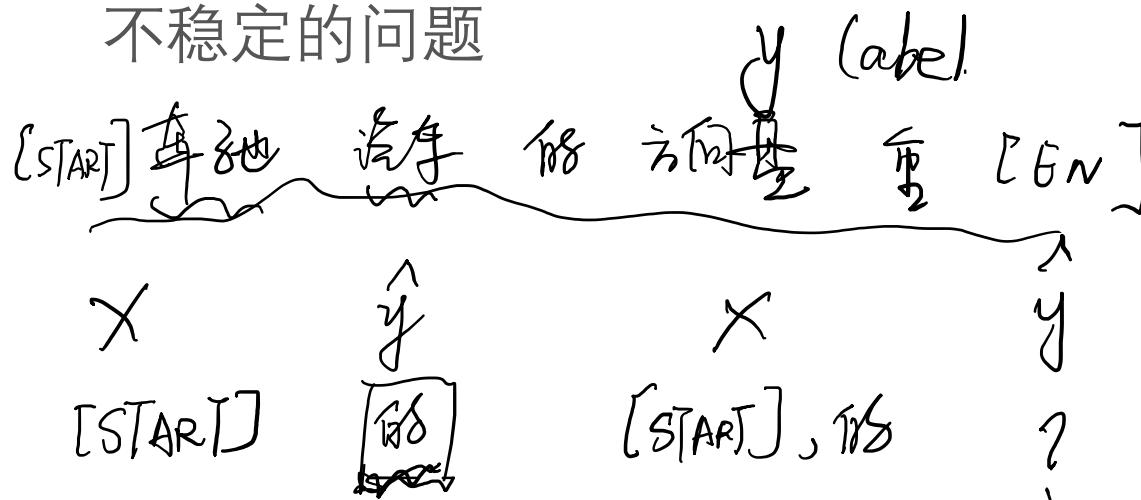
$\frac{y+y}{\sum} \rightarrow \hat{y} \leftarrow P_i = \frac{e^{y_i}}{\sum e^{y_i}}$

Seq2Seq训练

Teacher Forcing

RNN模型是用前一步的输出作为输入

Teacher Forcing方法能够解决收敛速度慢和不稳定的问题



Seq2Seq训练

Teacher Forcing

Seq2Seq训练

Exposure Bias

$$P(y|x) = P(y_1|x) P(y_2|x, y_1) \dots P(y_n|x, y_1, \dots, y_{n-1})$$

$\sim \log n$



Baseline

20

90
80

20

20

30

30

40

Bert

30

40

40

i. Scheduled Sampling

ii. RL

iii. 构造负样本 (mask) ✓

iv. 对抗训练 (梯度惩罚)

Bert

50%

20%

20%

20%

20%

20%

20%

20%

20%

20%

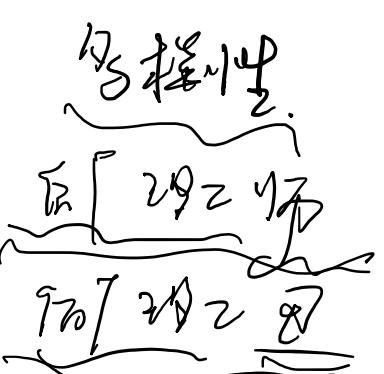
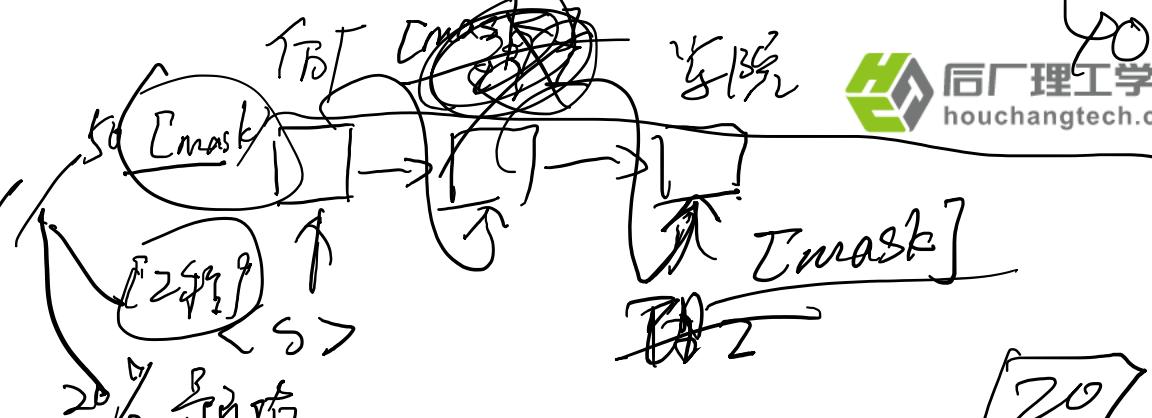
20%

20%

20%

20%

20%



作业

Bye !