

# RNN

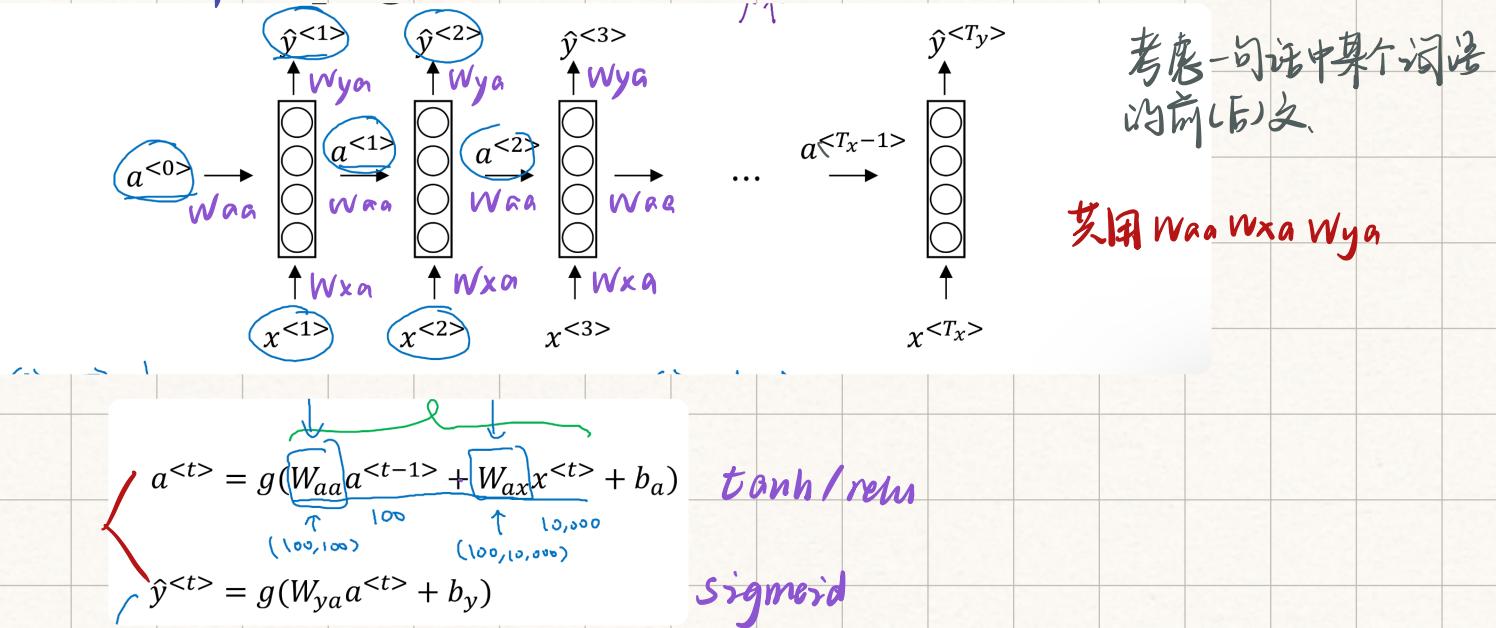
## 一. 循环序列模型 Sequence model.

### 1. Notation

words:  $x^{(1)} x^{(2)} \dots x^{(T_x)}$   $T_x = 9$

dicts:  $y^{(1)} y^{(2)} \dots y^{(T_y)}$   $T_y = 9$

## 2. 循环神经网络 Recurrent Neural Network



### 反向传播

$$-\sum y \sum \log(\hat{y})$$

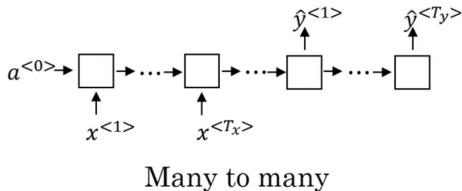
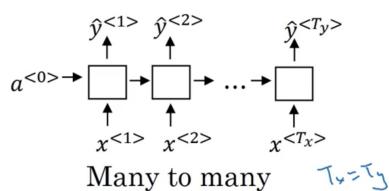
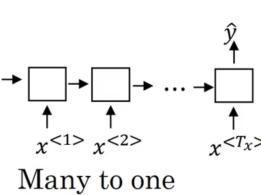
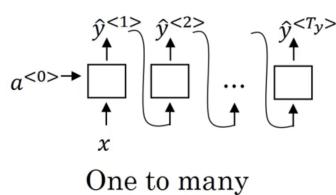
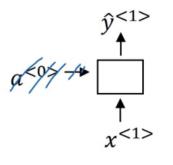
对  $\hat{y}^{<i>}$  使用交叉熵损失 (Logistic)  $\Rightarrow d_i$

总损失  $\lambda = \sum d_i$ . backprop through time

## 3. 多种类型的 RNN

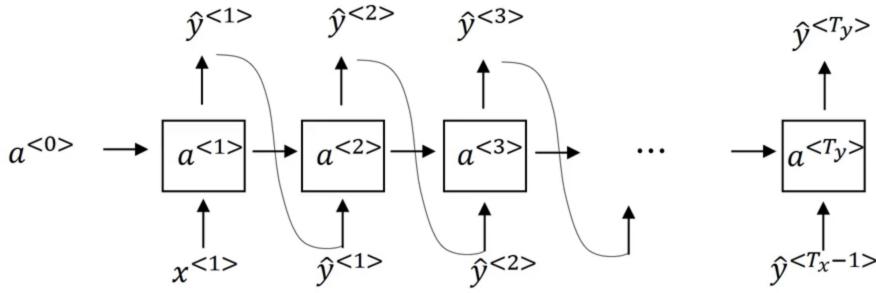
Summary of RNN types

网易云课堂



## 4. 语言模型和序列生成

### 5. 新序列采样



单词库, unk, EOS,  
字符库

### 6. 梯度消失.

靠后的梯度很难影响靠前的层.

即子后段的单词无法长期依赖前段的单词.

梯度爆炸: 剪枝

梯度消失: ? (GRU) 如何在长序列保持有效

### 6. GRU门控循环单元 gated recurrent unit

$c$  = memory cell.

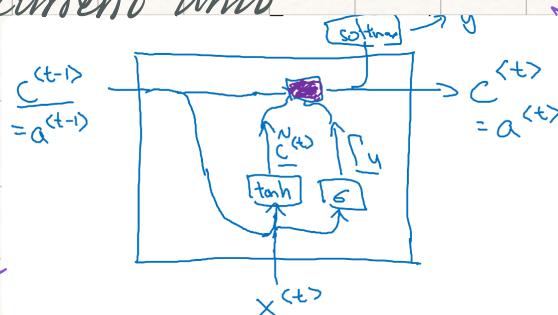
$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

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

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

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

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



$\Gamma_u$  - 取很小，因此  $c^{<t>}$  变化很久

用  $c$  代替  $a$ .

### 7. LSTM长短时记忆网络 Long Short Term Memory

org  
org

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

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

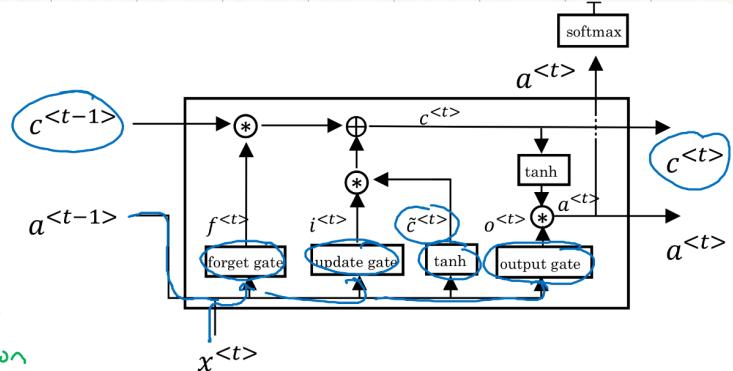
$$\rightarrow \Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f) \text{ Forget}$$

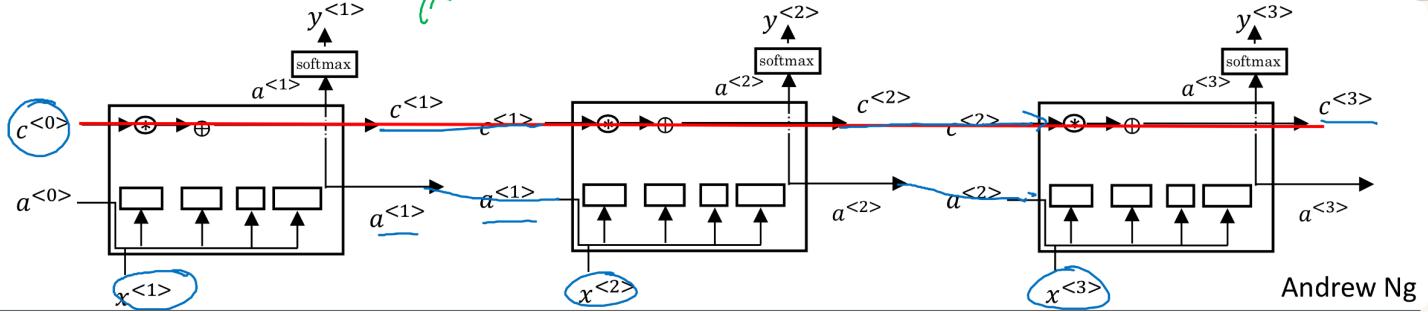
$$\rightarrow \Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o) \text{ Output}$$

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

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

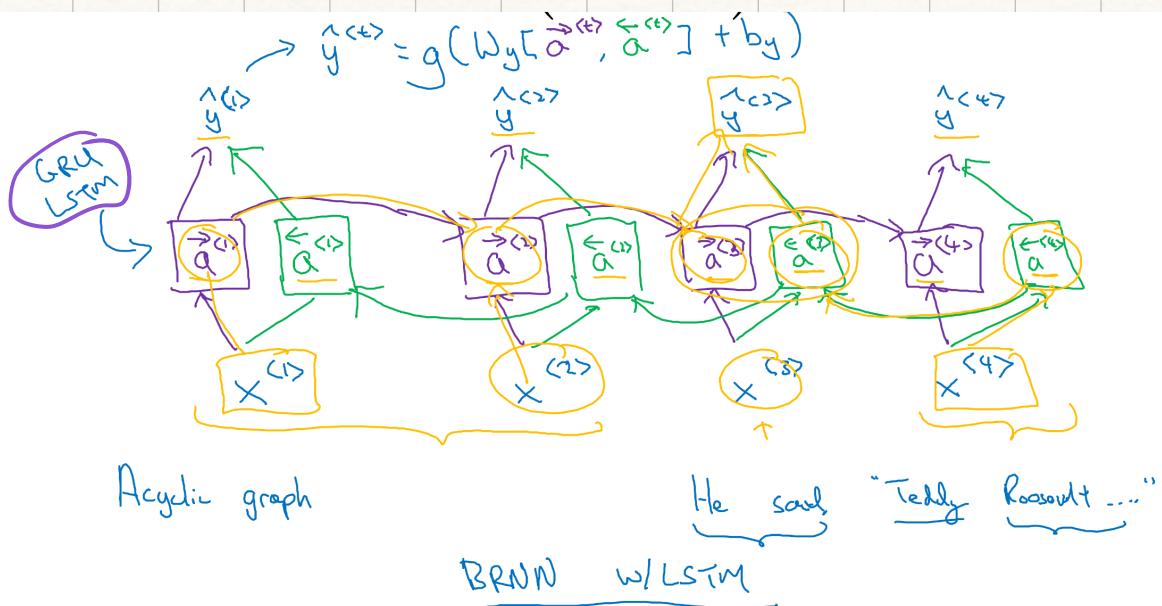
perphole  
connection





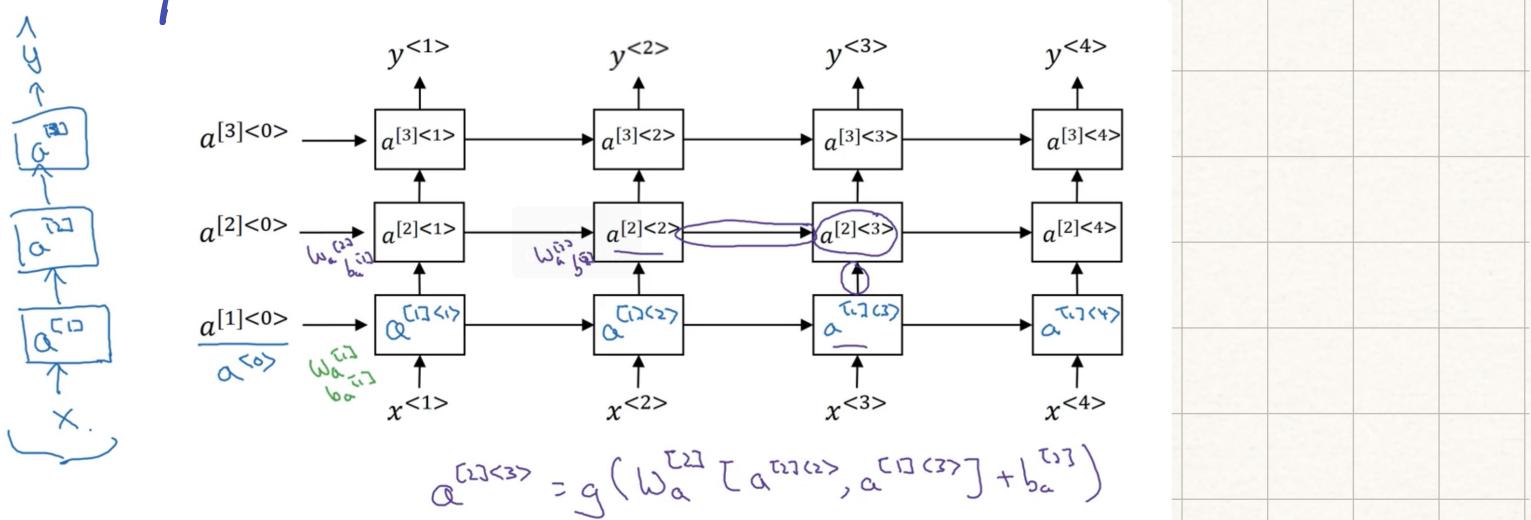
Andrew Ng

## 8. BRNN. 双向循环网络 Bidirectional.



同时考虑的前后的信息来进行预测. 常用 BRNN+LSTM

## 9. Deep RNN



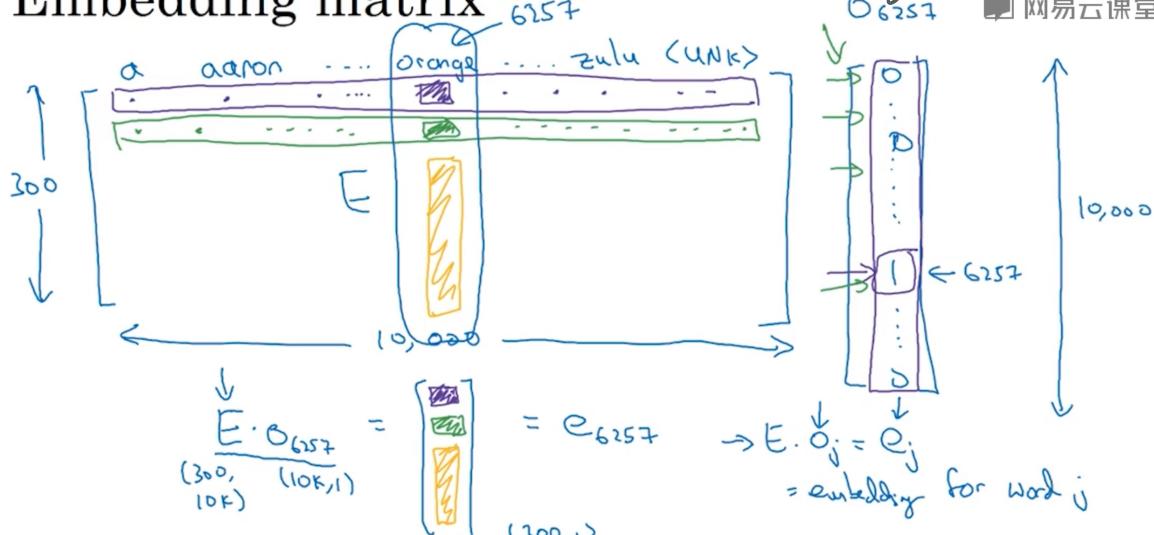
Andrew Ng

层数不会很深

# 三. 自然语言处理

## 1. 词向量和词嵌入 Word Embedding

### Embedding matrix



网易云课堂

Andrew Ng

# 三. 序列模型和注意力机制

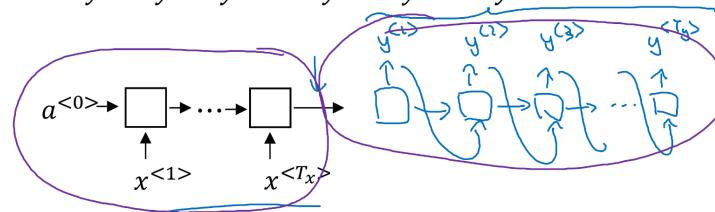
## 1. 基础模型. Seq2Seq.

### Sequence to sequence model

$x^{<1>} x^{<2>} x^{<3>} x^{<4>} x^{<5>}$   
Jane visite l'Afrique en septembre

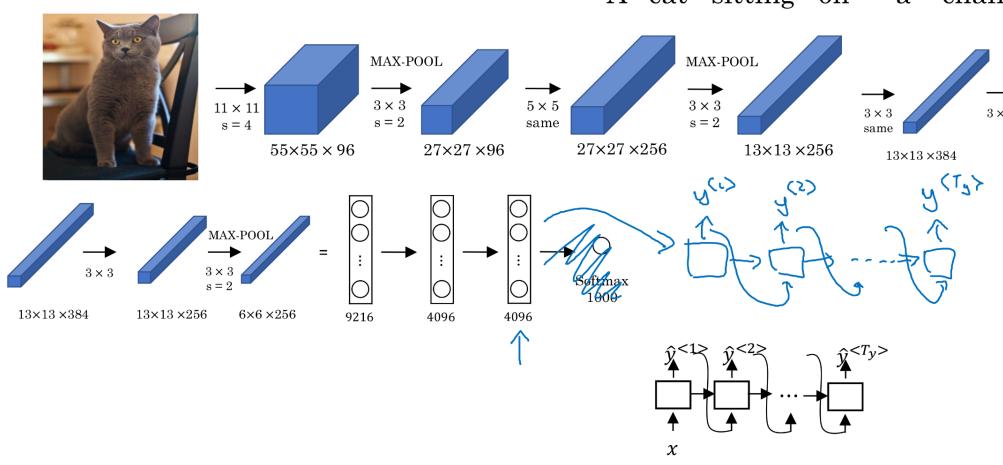
→ Jane is visiting Africa in September.

$y^{<1>} y^{<2>} y^{<3>} y^{<4>} y^{<5>} y^{<6>}$



机器翻译

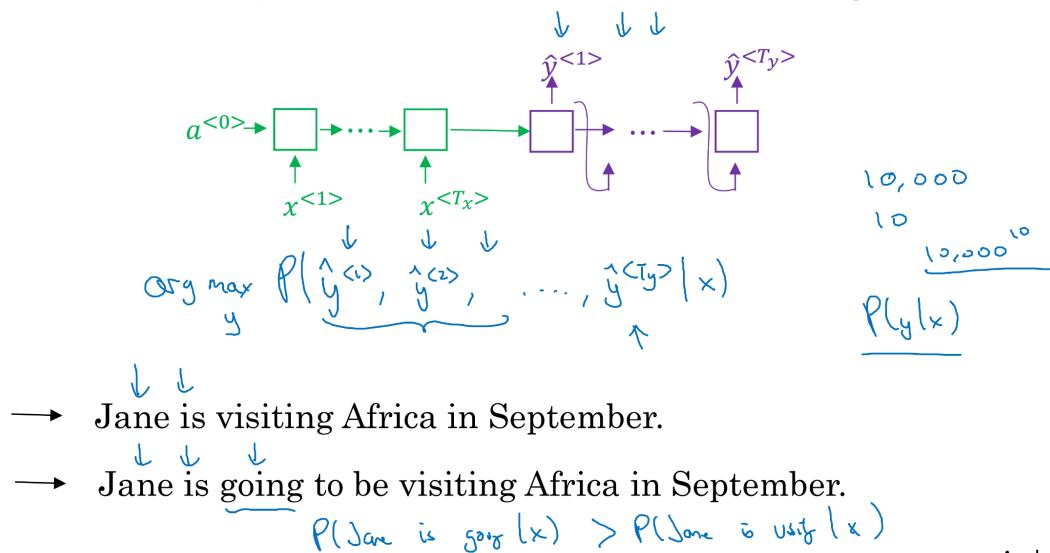
## Image captioning



$y^{<1>} y^{<2>} y^{<3>} y^{<4>} y^{<5>} y^{<6>} \}$

## 2. 选择取尽可能的句子.

Why not a greedy search?



→ Jane is visiting Africa in September.

→ Jane is going to be visiting Africa in September.

$$P(\text{Jane is going}|x) > P(\text{Jane is visit}|x)$$

贪心算法会导致高级词句流以及句子长度过长.

采用束搜索法.  $B=3$ .  $\Rightarrow$  每次保留概率最高的三组词句.

## 3. 改进定向搜索算法

Length normalization

$$\arg \max_y \prod_{t=1}^{T_y} P(y^{\langle t \rangle} | x, y^{\langle 1 \rangle}, \dots, y^{\langle t-1 \rangle})$$

$$\log \sum_{y^{\langle t \rangle}} \log P(y^{\langle t \rangle} | x, y^{\langle 1 \rangle}, \dots, y^{\langle t-1 \rangle})$$

$T_y = 1, 2, 3, \dots, 30.$

$$\rightarrow \frac{1}{T_y^\alpha} \sum_{t=1}^{T_y} \log P(y^{\langle t \rangle} | x, y^{\langle 1 \rangle}, \dots, y^{\langle t-1 \rangle})$$

$\alpha = 0.7$      $d = 1$   
 $\alpha = 0$      $d = 0$

Andrew Ng

减少长句子 (厚参数偏向更短的句子).

## 4 Bleu Bilogical Evaluation Understanding

# Bleu score on unigrams

Example: Reference 1: The cat is on the mat.

$$P_1, P_2 = 1.0$$

Reference 2: There is a cat on the mat.

MT output: The cat the cat on the mat. ( $\hat{y}$ )

$$P_1 = \frac{\sum_{\text{unigrams} \in \hat{y}} \text{count}_{\text{clip}}(\text{unigram})}{\sum_{\text{unigrams} \in \hat{y}} \text{count}(\text{unigram})}$$

↑  
Unigram  
Unigram  $\in \hat{y}$   
Count (unigram)

$$p_n = \frac{\sum_{n\text{-grams} \in \hat{y}} \text{count}_{\text{clip}}(n\text{-gram})}{\sum_{n\text{-grams} \in \hat{y}} \text{count}(n\text{-gram})}$$

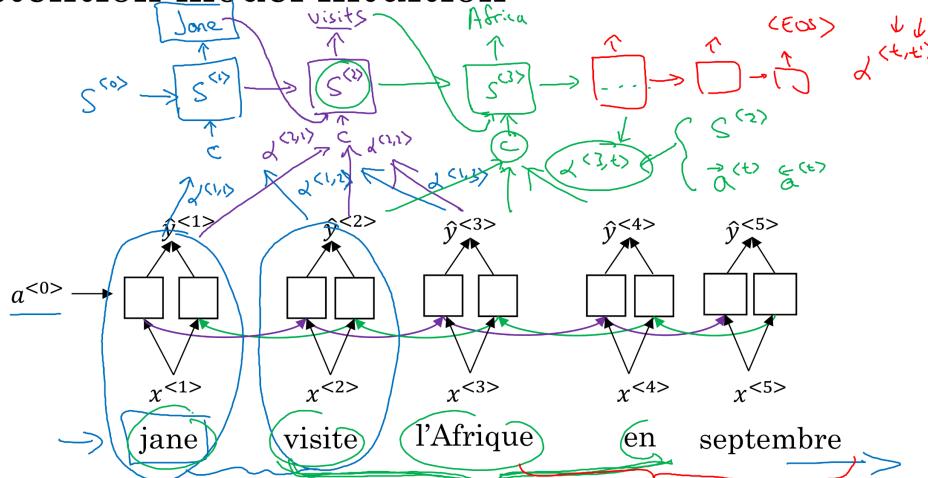
↑  
 $n\text{-gram}$   
 $n\text{-gram} \in \hat{y}$   
Count (n-gram)

$$\text{BP. } \exp\left(\frac{1}{4} \sum_{n=1}^4 p_n\right).$$

## 5. 注意力模型

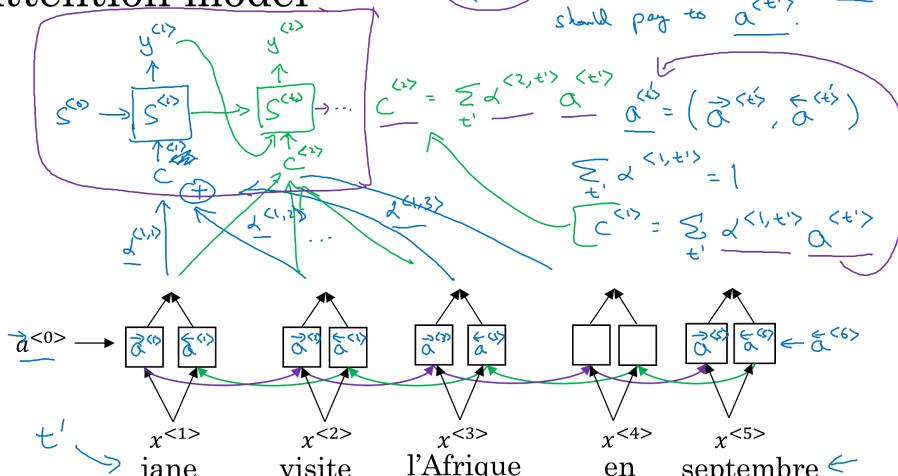
长句中，Bleu会下降，因此需要考虑全局。

### Attention model intuition



[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate] Andrew Ng

### Attention model



[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

Andrew Ng

## Computing attention $\alpha^{<t,t'>}$

