



deeplearning.ai

# Sequence to sequence models

---

## Basic models

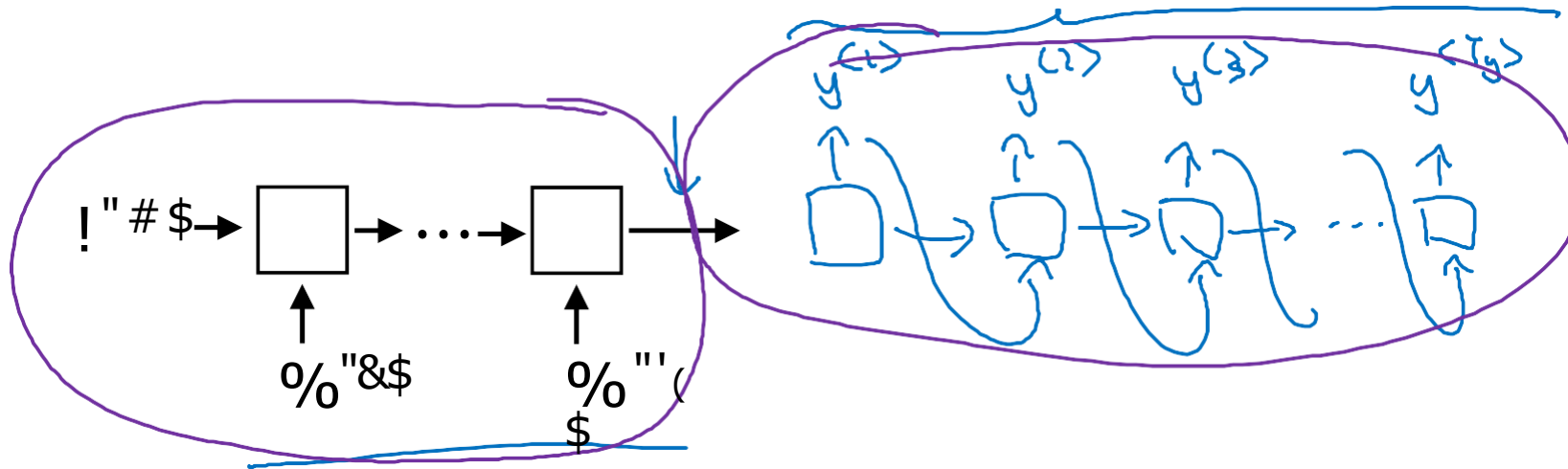
# Sequence to sequence model

%"&\$ %"\*\$ %"+\$ %", \$ %"- \$

Jane visite l'Afrique en septembre

→ Jane is visiting Africa in September.

. "&\$ . "\*\$ . "+\$ . ", \$ . "- \$ . "/ \$

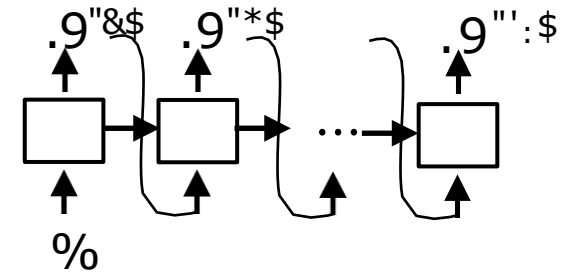
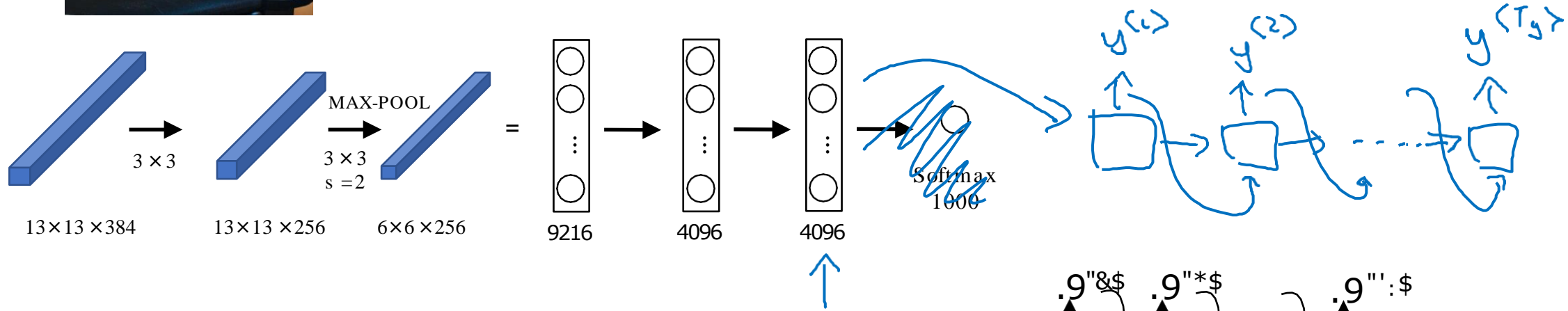
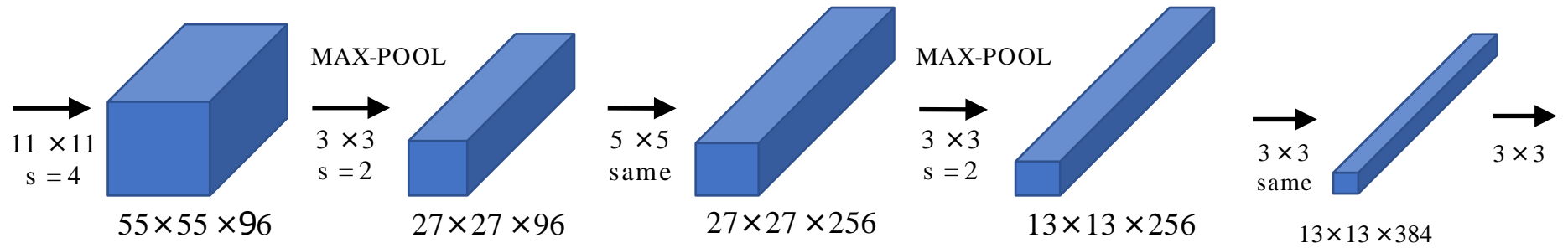
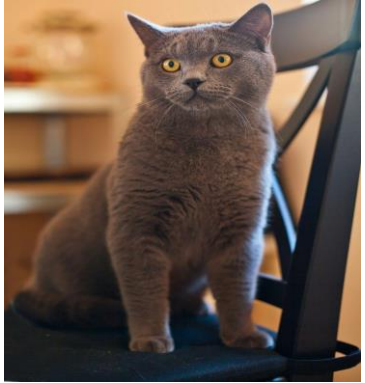


[Sutskever et al., 2014. Sequence to sequence learning with neural networks] ←

[Cho et al., 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation] ←

# Image captioning

. "&\$ . "\*" \$ . "+" \$ . ", \$ . "- \$ . "/" \$ }  
A cat sitting on a chair



- [Mao et. al., 2014. Deep captioning with multimodal recurrent neural networks]
- [Vinyals et. al., 2014. Show and tell: Neural image caption generator]
- [Karpathy and Li, 2015. Deep visual-semantic alignments for generating image descriptions]



deeplearning.ai

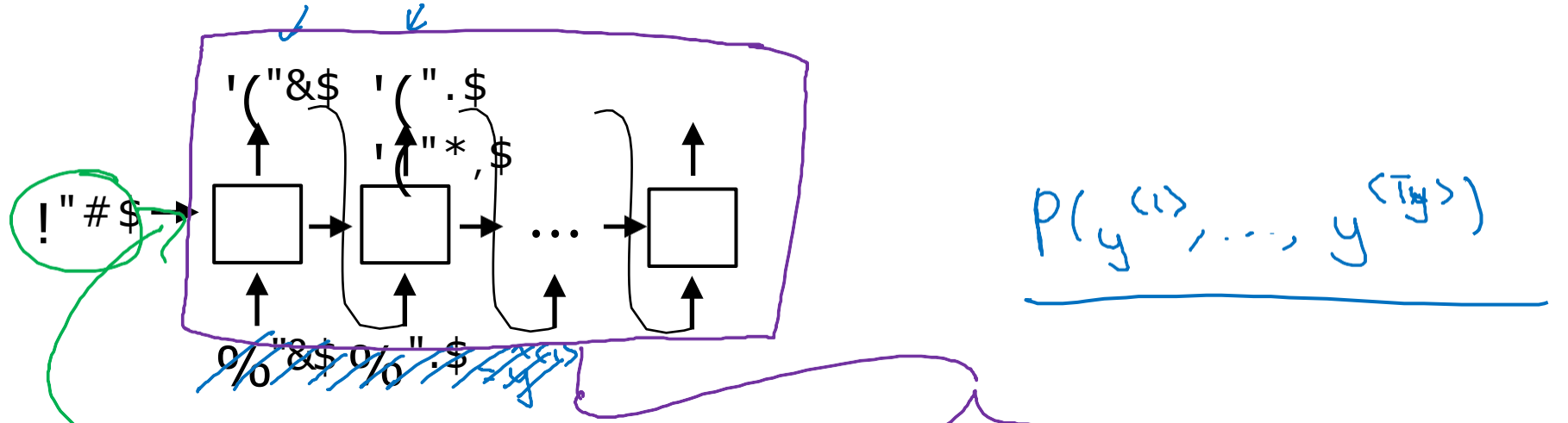
# Sequence to sequence models

---

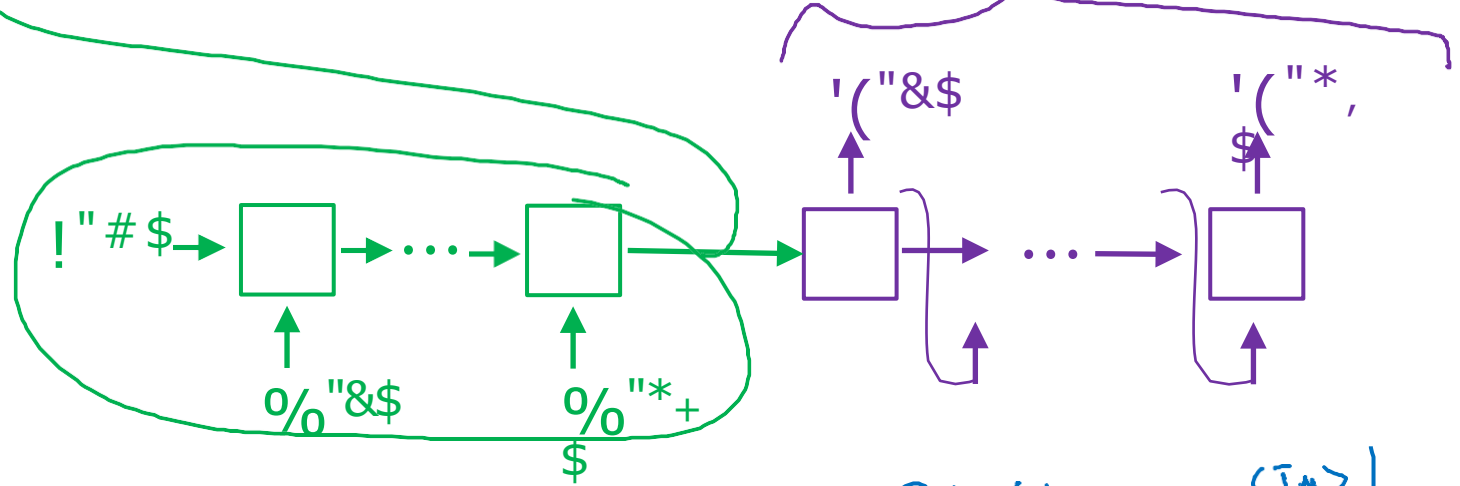
## Picking the most likely sentence

# Machine translation as building a conditional language model

Language model:



## Machine translation:



"Conditional language model"

$$P(y^{(1)}, \dots, y^{(T_y)} | \underline{x^{(1)}, \dots, x^{(T_x)}})$$

# Finding the most likely translation

Jane visite l'Afrique en septembre.

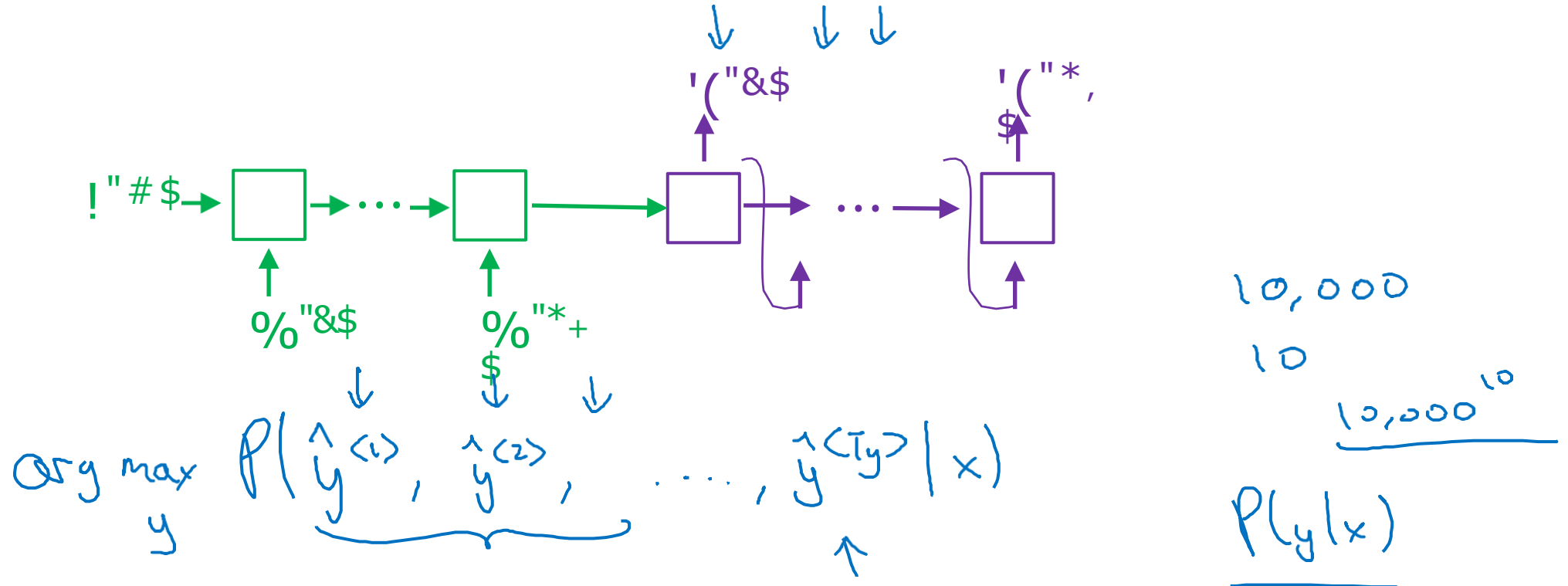
English  
French  
/ ( ' " & \$ , ... , ' " \* , \$ |  
% )

- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.
- In September, Jane will visit Africa.
- Her African friend welcomed Jane in September.

$$\arg \max_{i: i \leq n, j: j \leq m} \frac{P(i, j)}{P(i)P(j)}$$

# Why not a greedy search?

$$p(\hat{y}^{(1)} | x)$$



→ Jane is visiting Africa in September.

→ Jane is going to be visiting Africa in September.

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



deeplearning.ai

# Sequence to sequence models

---

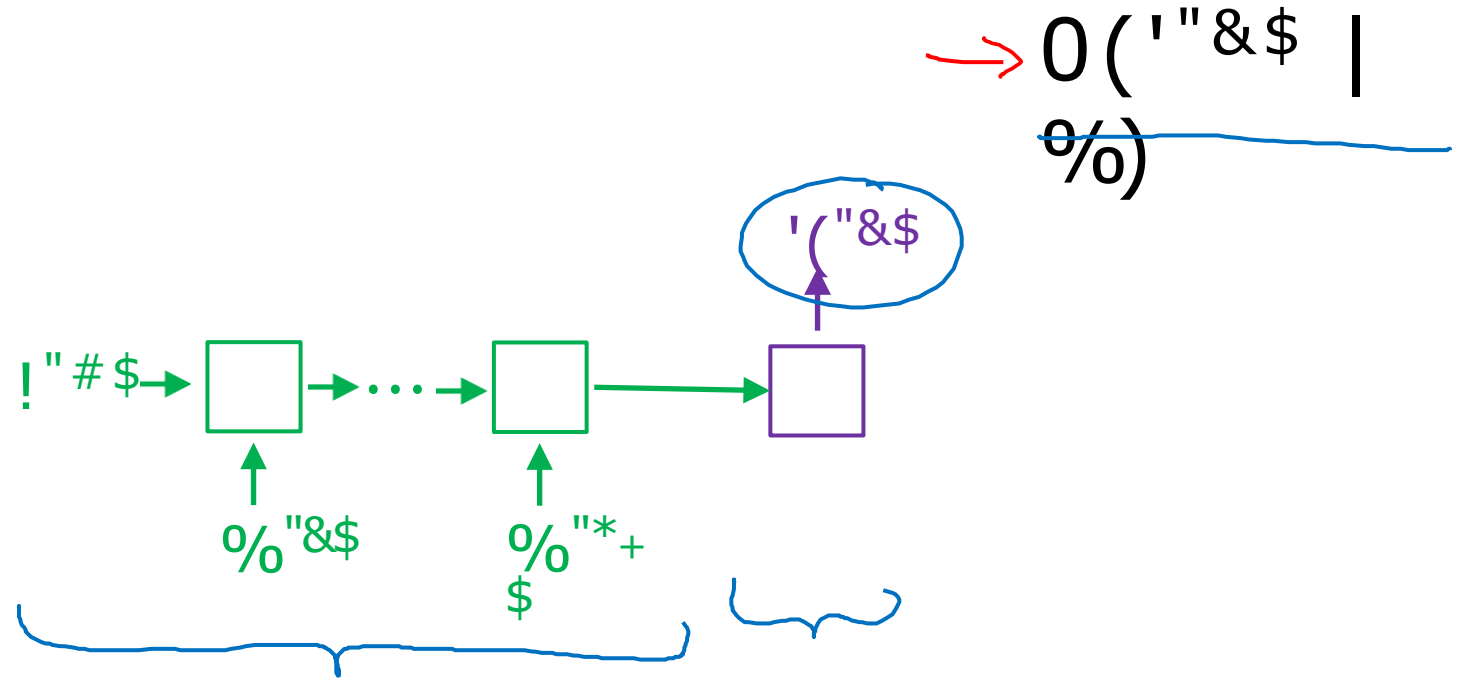
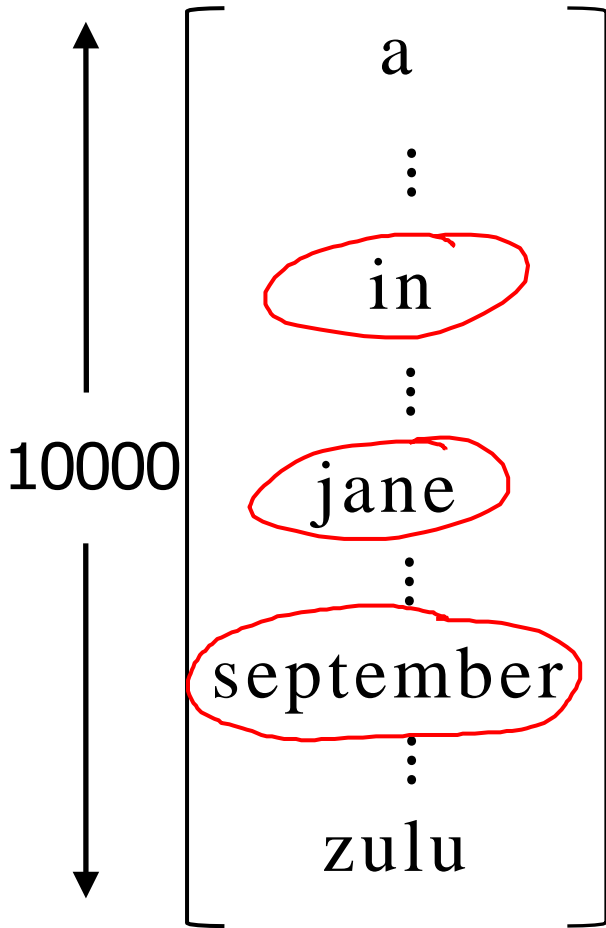
## Beam search



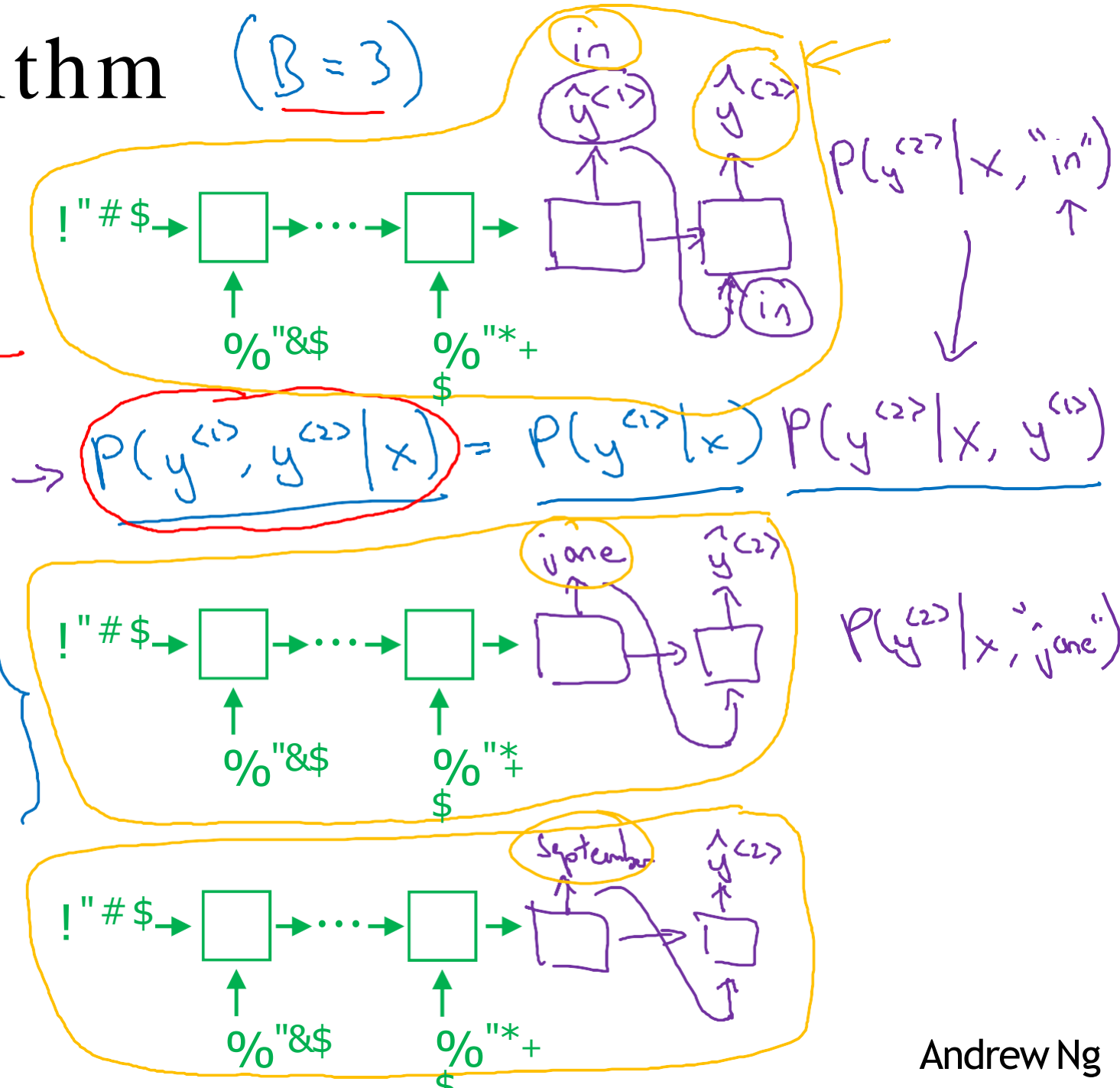
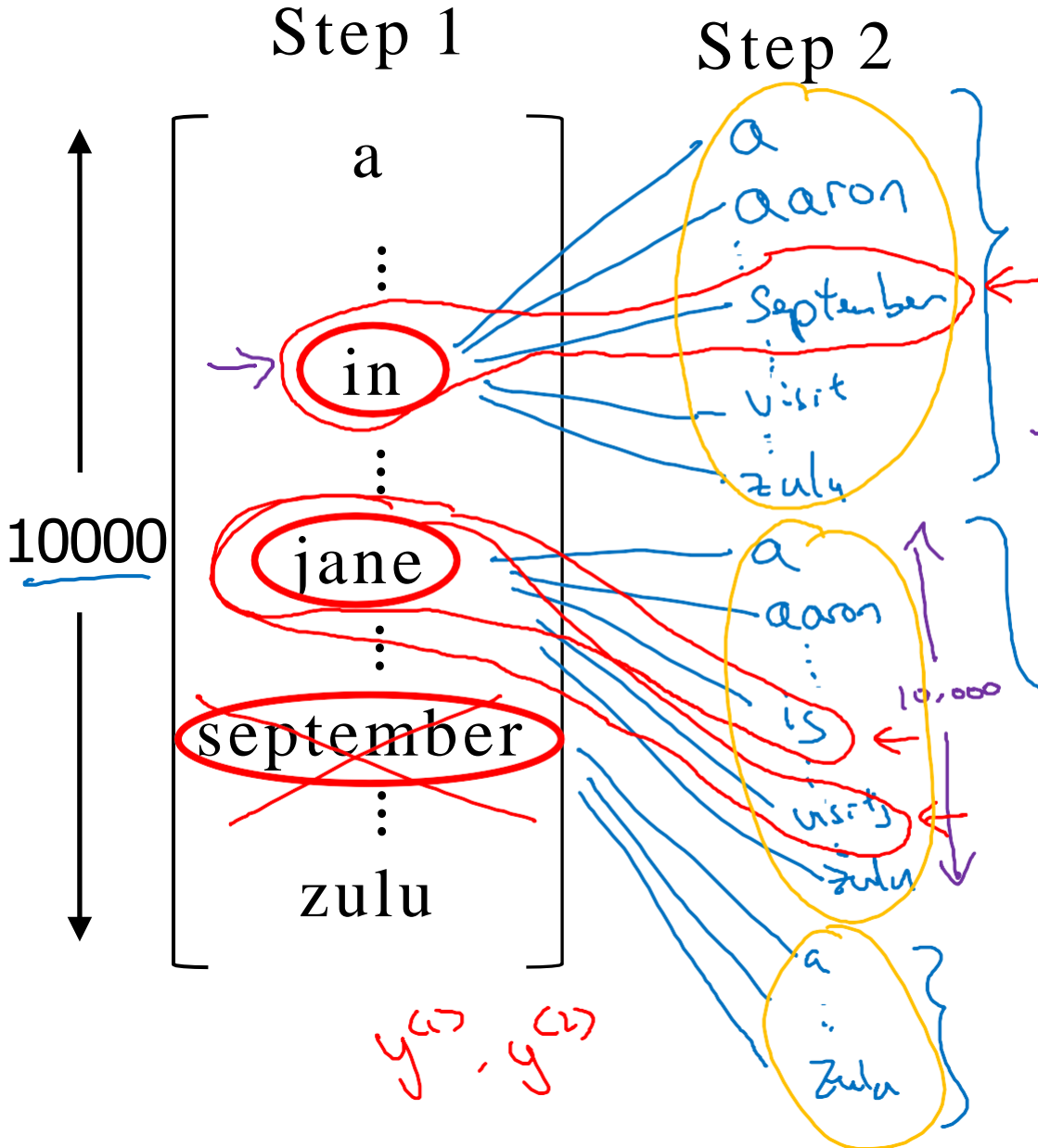
# Beam search algorithm

B = 3 (beam width)

Step 1



# Beam search algorithm

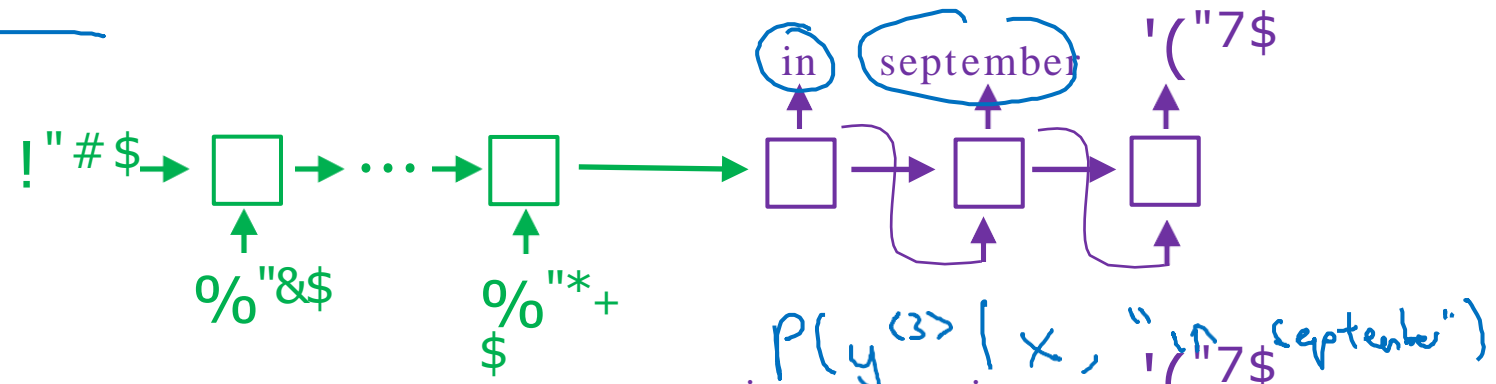
$$(B = 3)$$


# Beam search (4 =

$B=1 \rightarrow$  greedy search

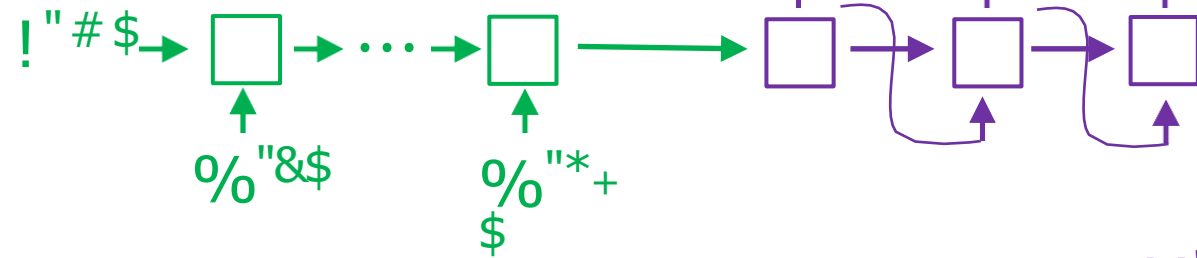
3) in september

*a*  
*aaron*  
*jane*  
*zulu*



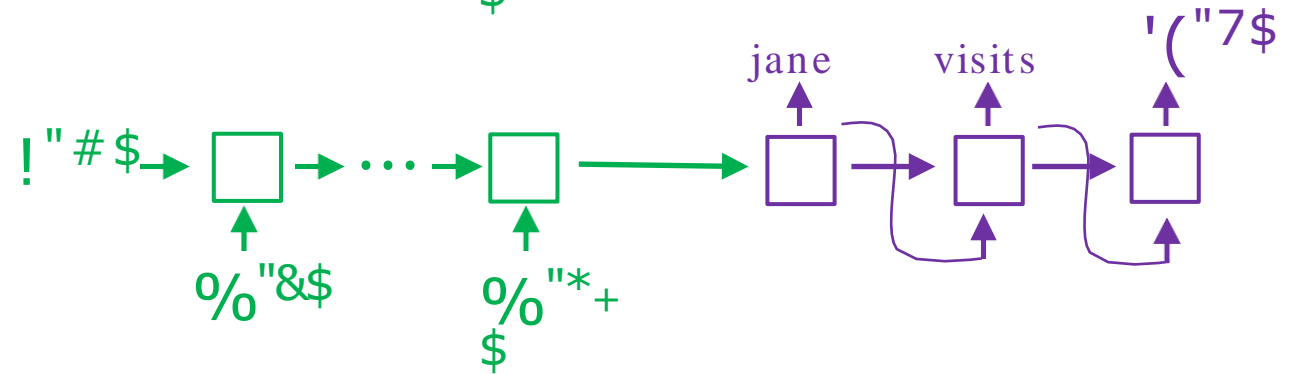
jane is

*a*  
*visits*  
*zulu*



jane visits

*a*  
*africa*  
*zulu*



$O('"&\$, ' "9\$ |$   
 $\%0)$

jane visits africa in september. <EOS>



deeplearning.ai

# Sequence to sequence models

---

## Refinements to beam search

# Length normalization

$$p(y^{(1)} \dots y^{(T_y)} | x) = \frac{p(y^{(1)} | x)}{p(y^{(T_y)} | x, y^{(1)} \dots, y^{(T_y-1)})} \dots$$

max( )

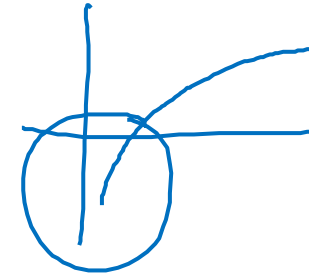
log

45

60

45

(\*, |, \*, \*\*, \*)



$$\log p(y|x) \leftarrow$$

$$p(y|x) \leftarrow$$

max7 9

60

(\*, |, \*, \*\*, \*)

y 1, 2, 3, ..., 30.

→

45

7 9

60

(\*, |, \*, \*\*, \*)

$\frac{1}{T_y \alpha}$

$$\underline{\alpha = 0.7}$$

$$\underline{\alpha = 1}$$

$$\underline{\alpha = 0}$$

# Beam search discussion

Beam width  $B$ ?

$1 \rightarrow 3 \rightarrow 10, \quad 100, \quad 1000 \rightarrow 3000$

large  $B$ : better result, slower  
small  $B$ : worse result, faster

Unlike exact search algorithms like BFS (Breadth First Search) or DFS (Depth First Search), Beam Search runs faster but is not guaranteed to find exact maximum for  $g_{max}(h)$ .



deeplearning.ai

# Sequence to sequence models

---

## Error analysis on beam search

# Example

Jane visite l'Afrique en septembre.

→ RNN

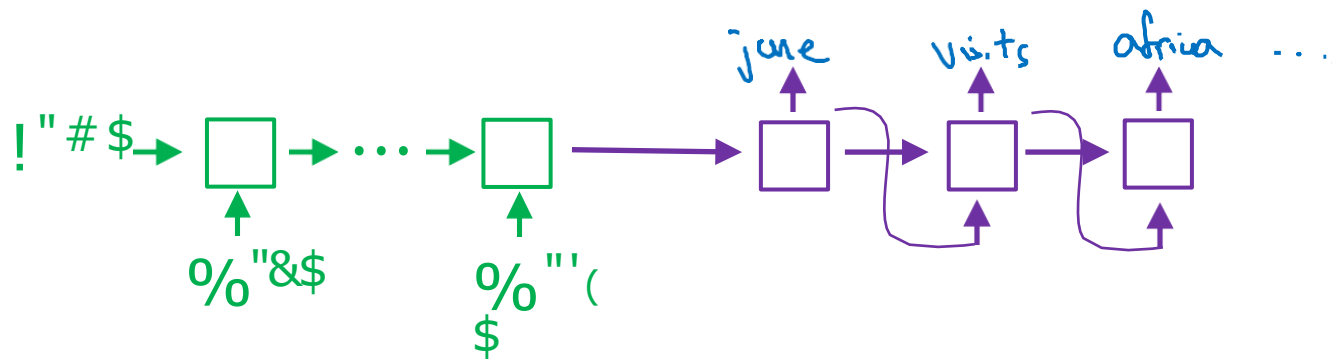
→ Beam Search

BT

Human: Jane visits Africa in September. ( $y^*$ )

Algorithm: Jane visited Africa last September. ( $\hat{y}$ ) ←

RNN computes  $P(y^*|x) \geq P(\hat{y}|x)$





# Error analysis on beam search

Human: Jane visits Africa in September. (+\*)

$$p(y^*|x)$$

$$p(\hat{y}|x)$$

Algorithm: Jane visited Africa last September. (+.)

Case 1:  $p(y^*|x) > p(\hat{y}|x)$  ←

$$\arg \max_y p(y|x)$$

Beam search chose +. But +\* attains higher  $p(+*|x)$

Conclusion: Beam search is at fault.

Case 2:  $p(y^*|x) \leq p(\hat{y}|x)$  ←

+\* is a better translation than +. But RNN predicted  $p(+*|x)$  <  $p(+.|x)$

%

Conclusion: RNN model is at fault

# Error analysis process

| Human   | Algorithm   | / ( +*  )   | / ( +. )  | At fault?   |
|---|---|---|---|---|
| Jane visits Africain<br>September.<br><br>..... | Jane visited Africa<br>last September.<br><br>..... | $\frac{0\%}{2 \times 10^{-10}}$<br>_____<br><br>_____ | $\frac{0\%}{1 \times 10^{-10}}$<br>_____<br><br>_____ | <div>B</div> <div>R</div> <div>R</div> <div>R</div> <div>R</div> <div>...</div> |

Figures out what faction of errors are “due to” beam search vs. RNN model



deeplearning.ai

# Sequence to sequence models

---

Bleu score  
(optional)

# Evaluating machine translation

French: Le chat est sur le tapis.

Reference 1: The cat is on the mat. ←

Reference 2: There is a cat on the mat. ←

MT output: the the the the the the the.

Precision:

Modified precision:

Bleu  
bilingual evaluation understudy

# Bleu score on bigrams

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

Reference 2: There is a cat on the mat. ←

MT output: The cat the cat on the mat. ←

|         | Count | Countclip |       |
|---------|-------|-----------|-------|
| the cat | 2 ←   | 1 ←       |       |
| cat the | 1 ←   | 0         | 4     |
| cat on  | 1 ←   | 1 ←       | <hr/> |
| on the  | 1 ←   | 1 ←       | 6     |
| the mat | 1 ←   | 1 ←       |       |
|         | ↑     |           |       |

# Bleu score on unigrams

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

Reference 2: There is a cat on the mat.

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

$$P_1 \cdot P_2 = \underline{1.0}$$

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

↑  
unigram

$$P_n = \frac{\sum_{n\text{-gram} \in \hat{y}} \text{Count}_{\text{clip}}(n\text{-gram})}{\sum_{n\text{-gram} \in \hat{y}} \text{Count}(n\text{-gram})}$$

# Bleu details

!=Bleu score on n-grams only

$P_1, P_2, P_3, P_4$

Combined Bleu score:

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

BP: brevity penalty

$$BP = \begin{cases} 1 & \text{if } \underline{\text{MT\_output\_length}} > \underline{\text{reference\_output\_length}} \\ \sqrt[n]{1 - \frac{\text{MT\_output\_length}}{\text{reference\_output\_length}}} & \text{otherwise} \end{cases}$$



deeplearning.ai

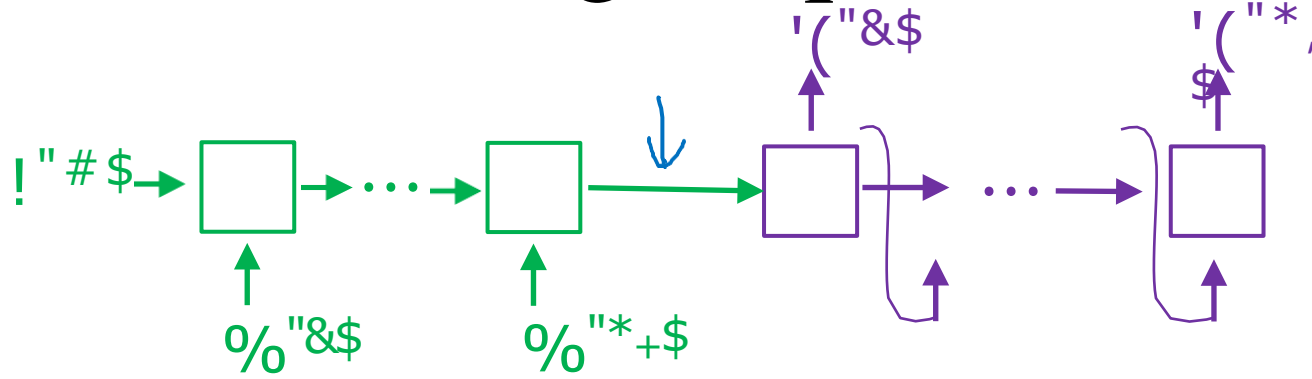
# Sequence to sequence models

---

## Attention model intuition

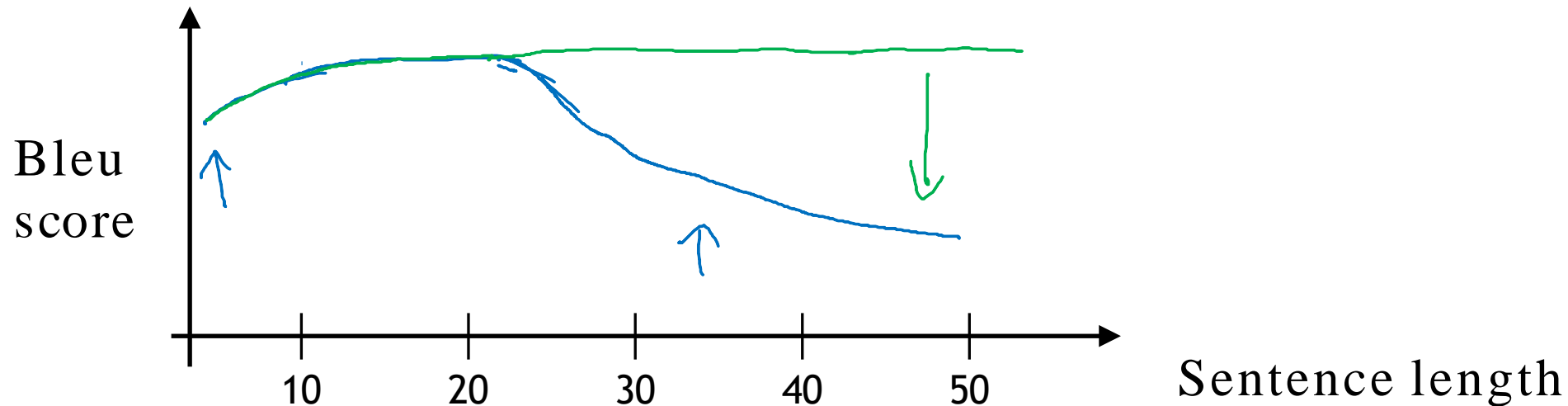


# The problem of long sequences

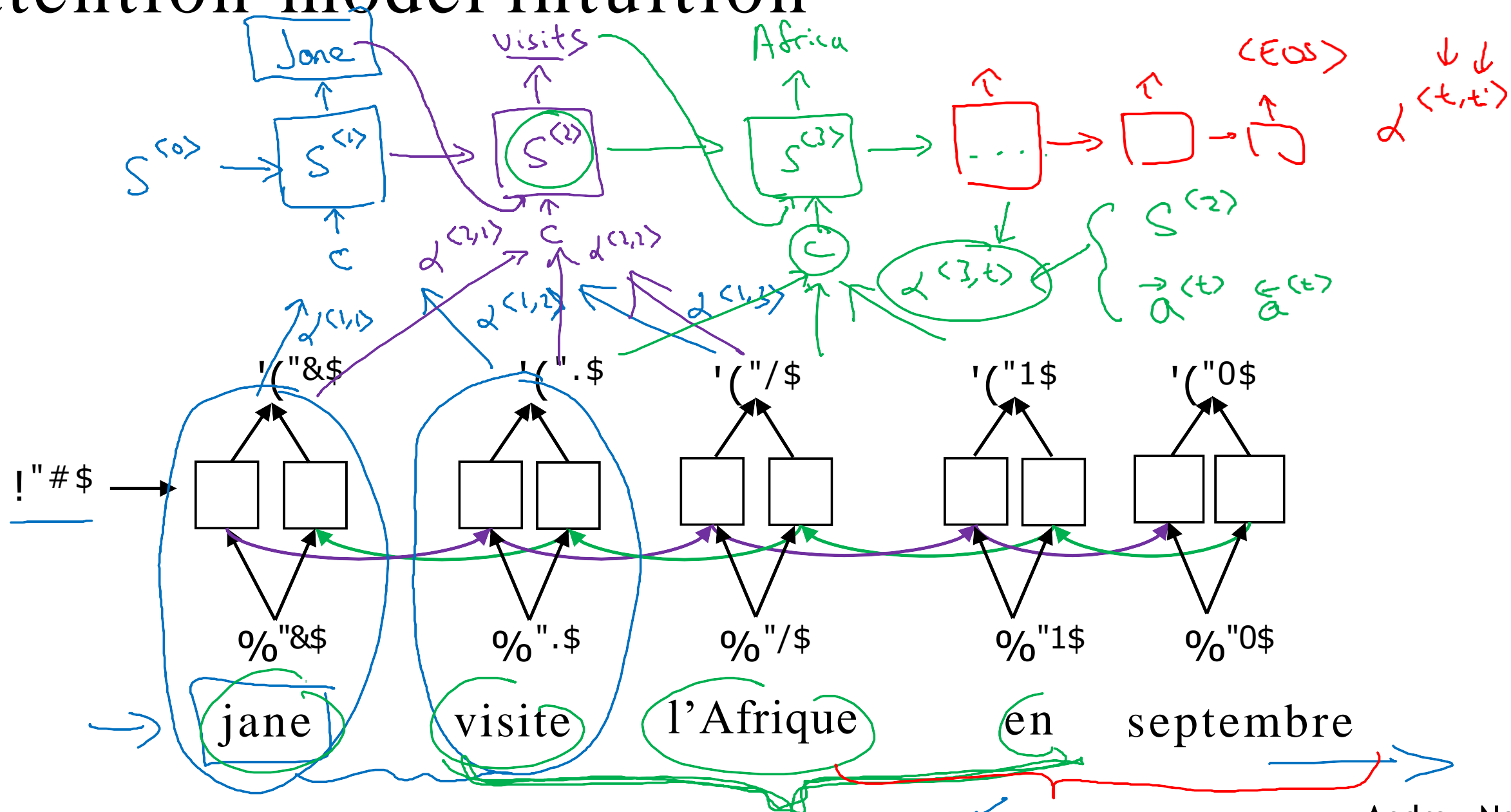


Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.

Jane went to Africa last September, and enjoyed the culture and met many wonderful people; she came back raving about how wonderful her trip was, and is tempting me to go too.



# Attention model intuition





deeplearning.ai

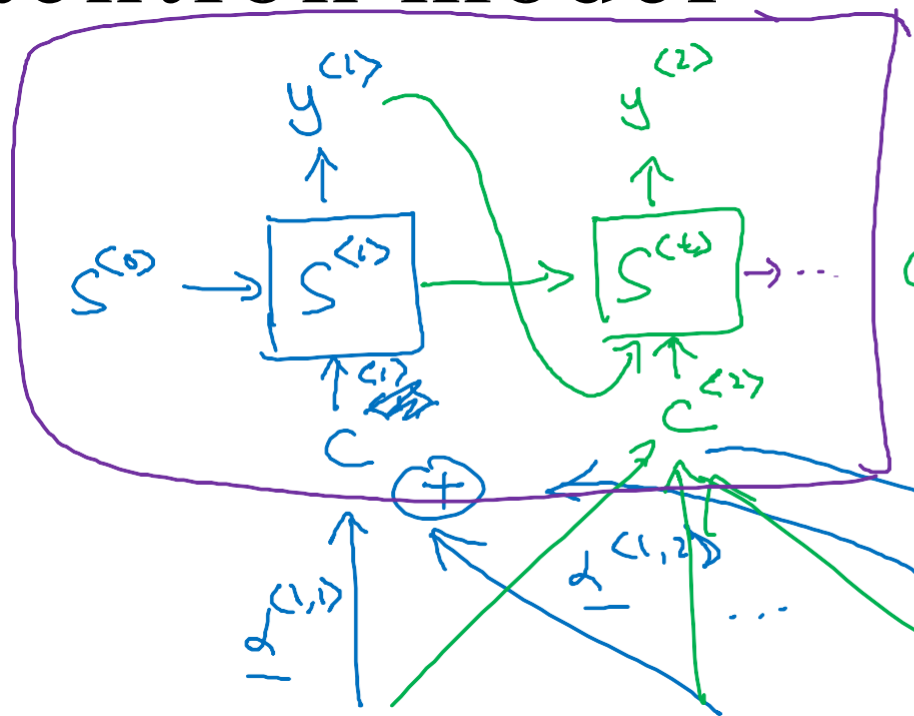
# Sequence to sequence models

---

## Attention model

# Attention model

$\alpha^{(t,t')}$  - amount of "attention"  $y^{(t)}$  should pay to  $a^{(t')}$ .

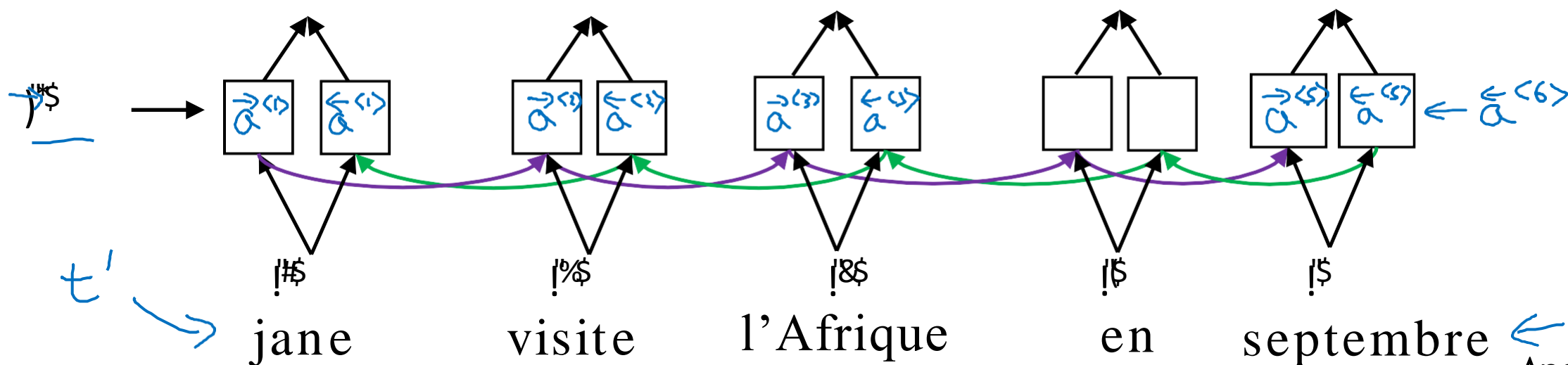


$$c^{(2)} = \sum_{t'} \alpha^{(2,t')} y^{(t')}$$

$$a^{(t')} = (\vec{a}^{(t')}, \leftarrow a^{(t')})$$

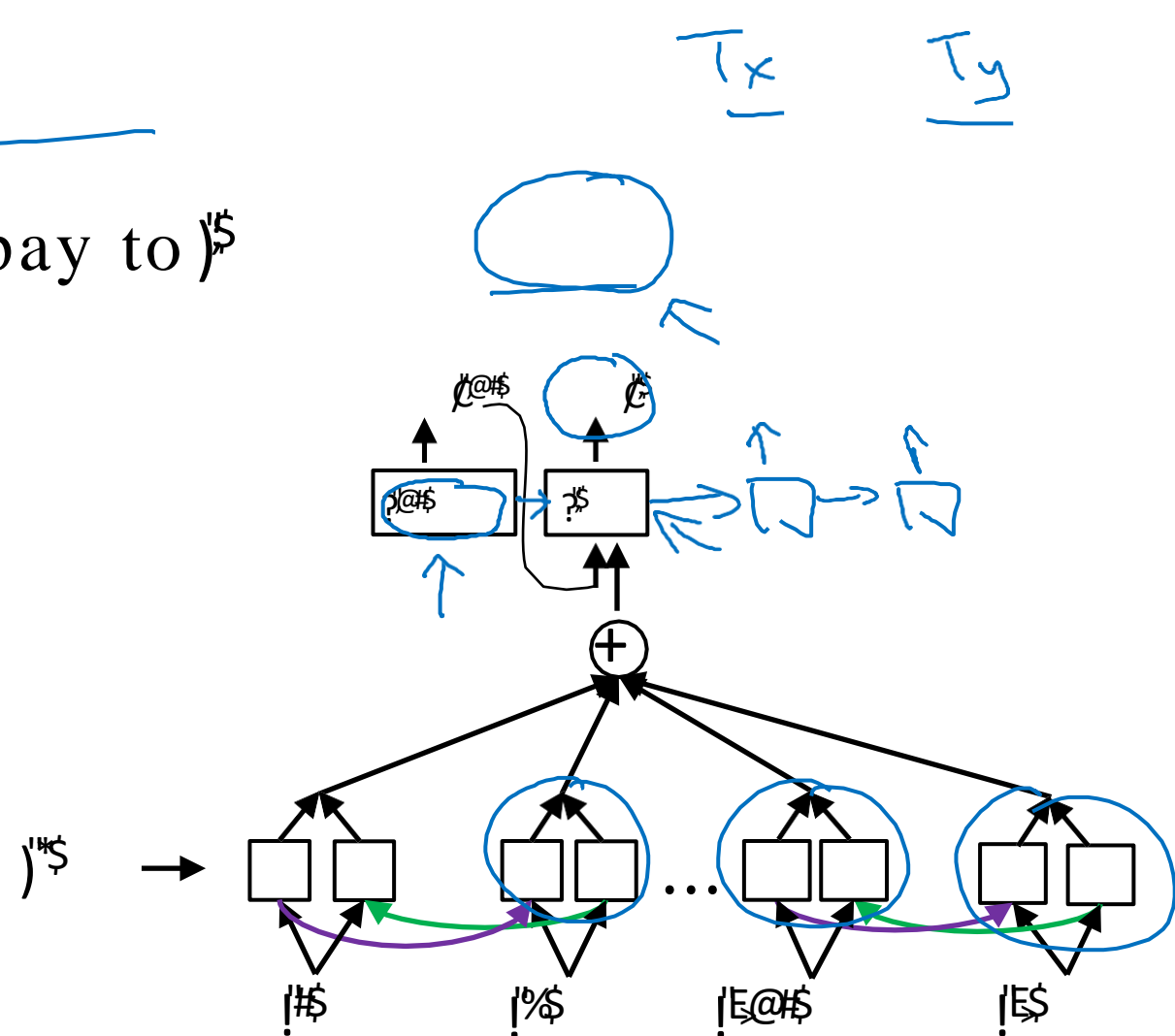
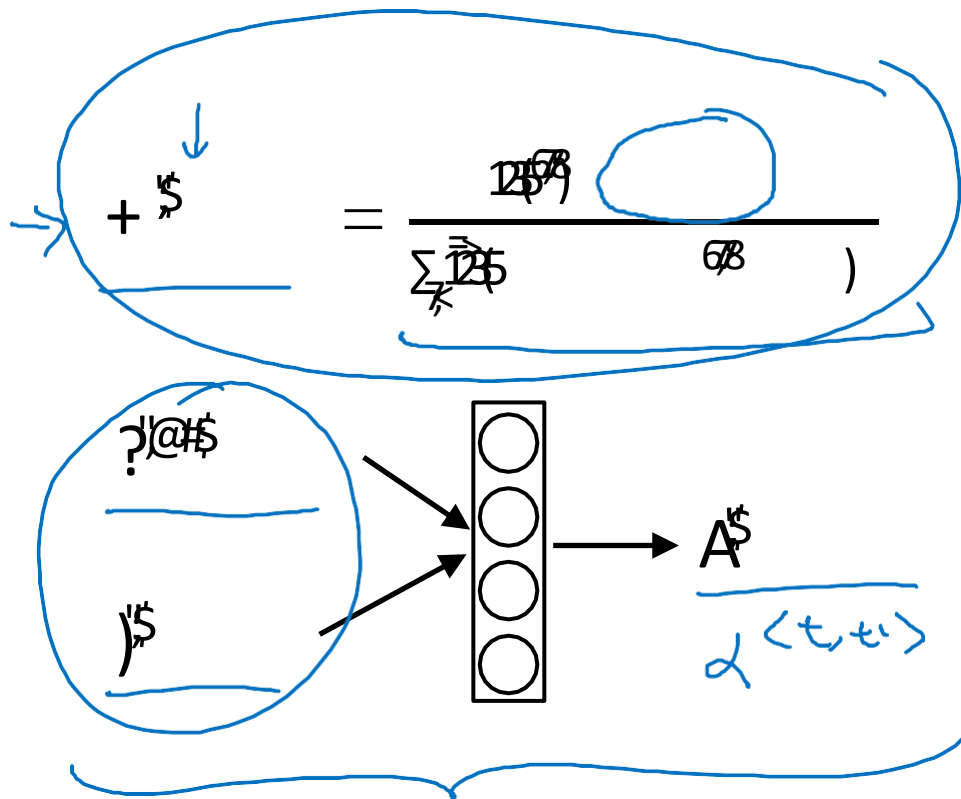
$$\sum_{t'} \alpha^{(1,t')} = 1$$

$$c^{(1)} = \sum_{t'} \alpha^{(1,t')} a^{(t')}$$



# Computing attention $\alpha^s$

$\alpha^s$  = amount of attention  $s$  should pay to  $j^s$



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

[Xu et. al., 2015. Show, attend and tell: Neural image caption generation with visual attention]

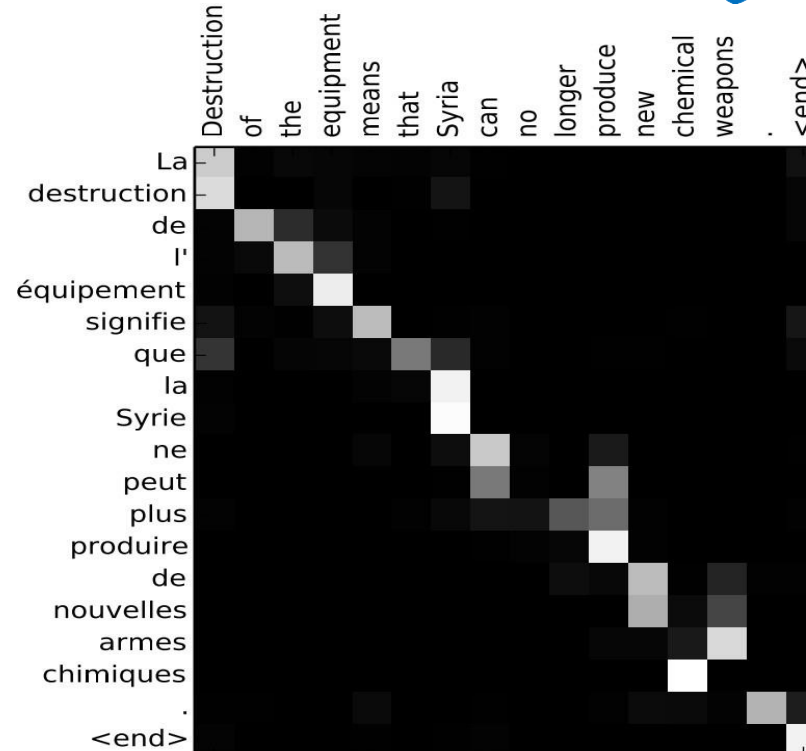
Andrew Ng

# Attention examples

July 20th 1969 → 19-07-20

23 April, 1964 → 1964-04-23

Visualization of  $\alpha$ .





deeplearning.ai

# Audio data

---

# Speech recognition

# Speech recognition problem

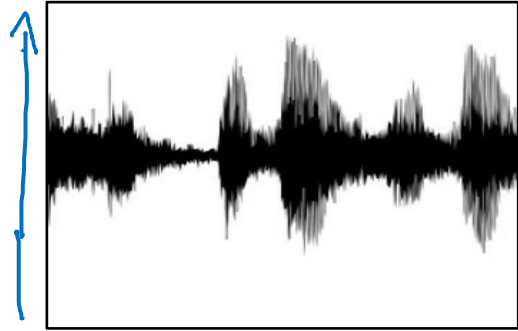
!

audio clip



#

transcript



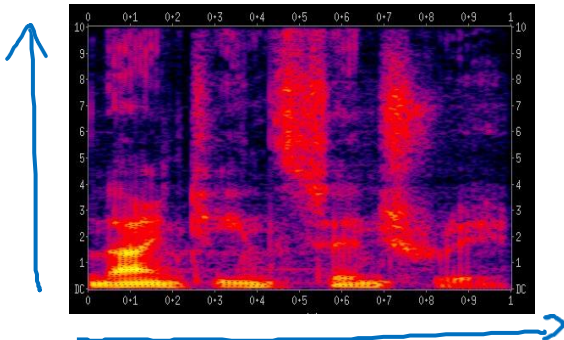
“the quick brown fox”

→ phonemes: de kwi braun

300h

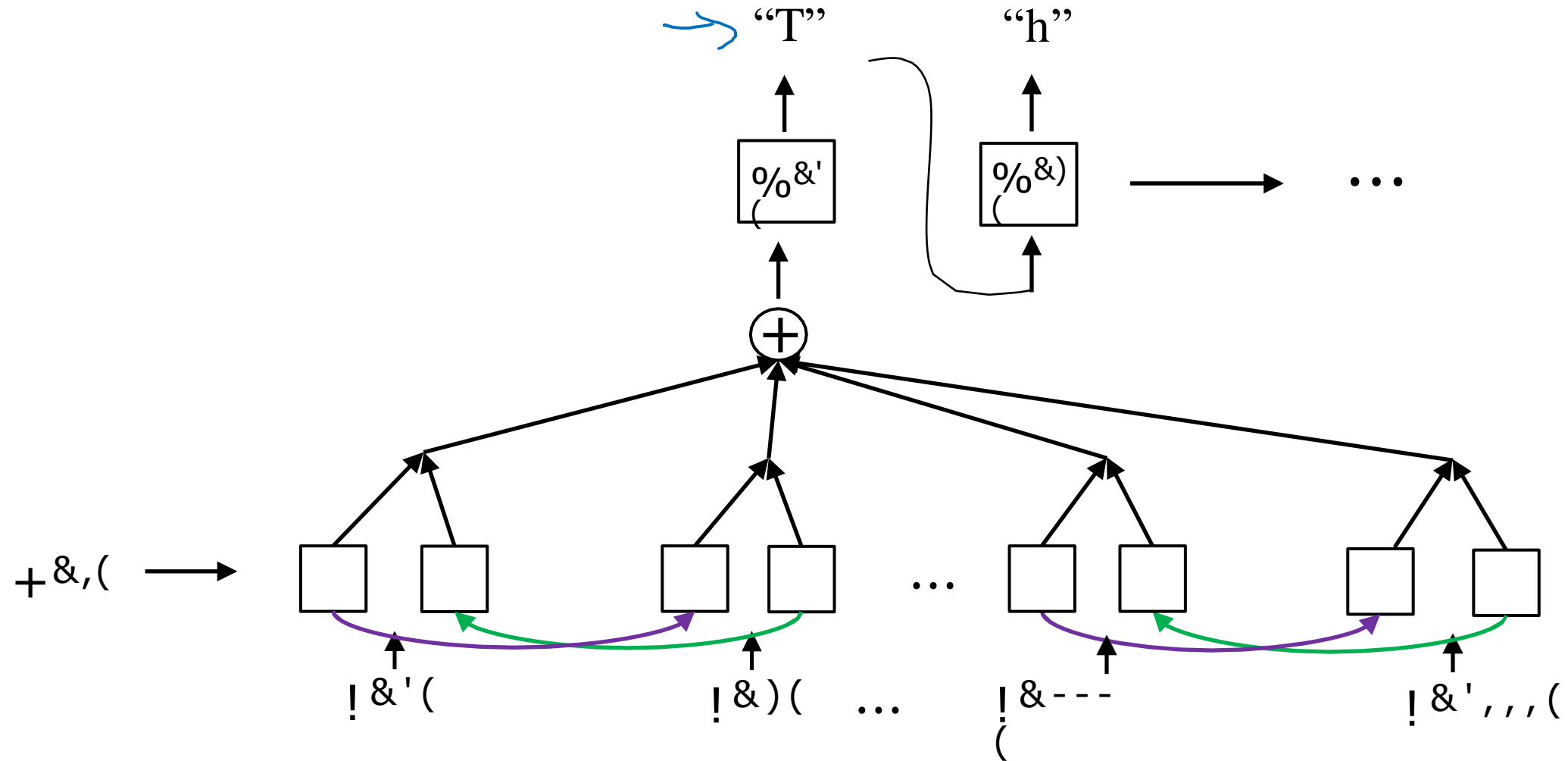
3000h

100,000h



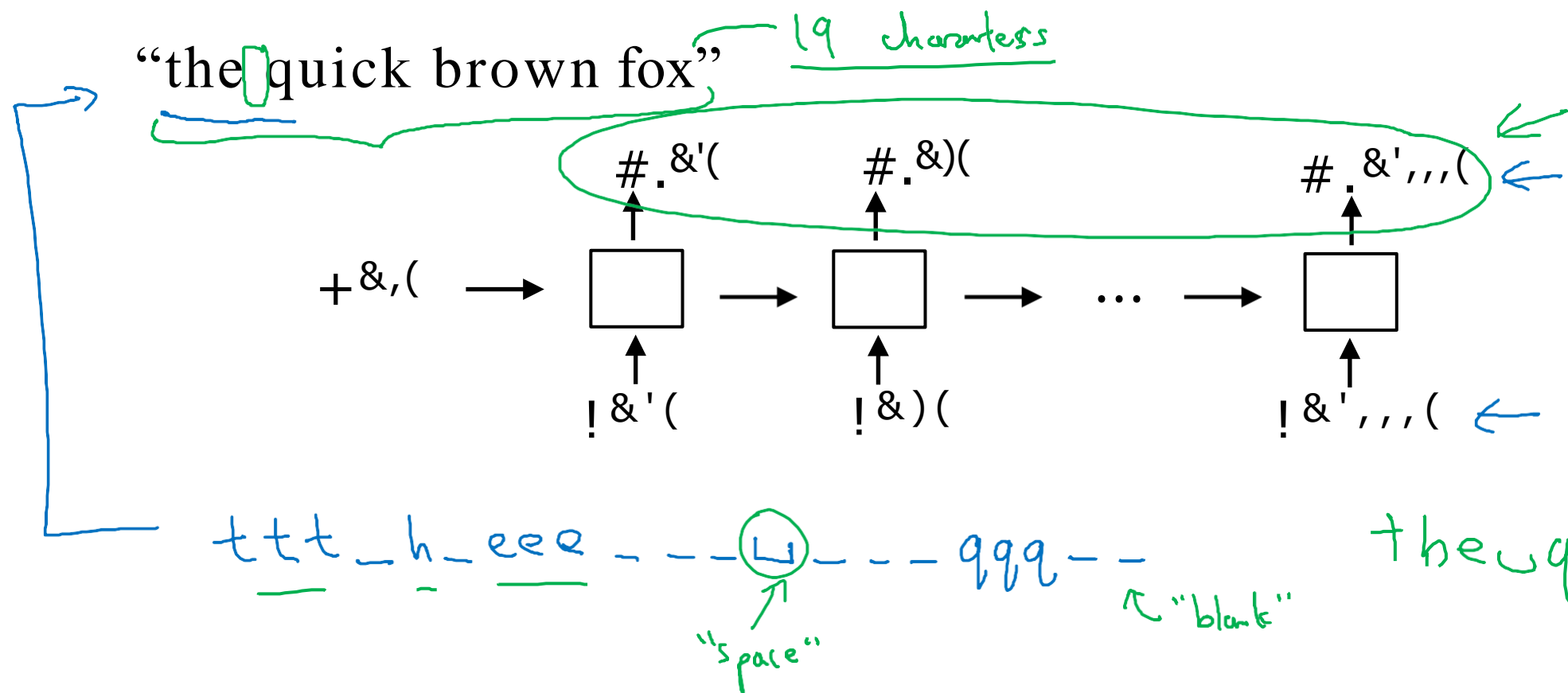


# Attention model for speech recognition



# CTC cost for speech recognition

(Connectionist temporal classification)



Basic rule: collapse repeated characters not separated by “blank”



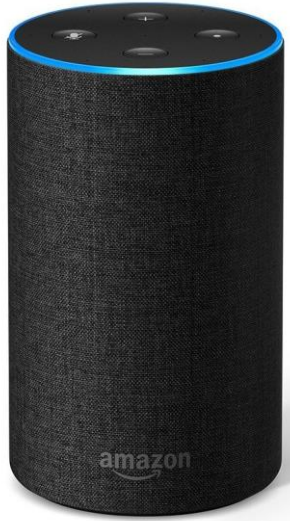
deeplearning.ai

# Audio data

---

# Trigger word detection

# What is trigger word detection?



Amazon Echo  
(Alexa)



Baidu DuerOS  
(xiaodunihao)



Apple Siri  
(Hey Siri)



Google Home  
(Okay Google)

# Trigger word detection algorithm

