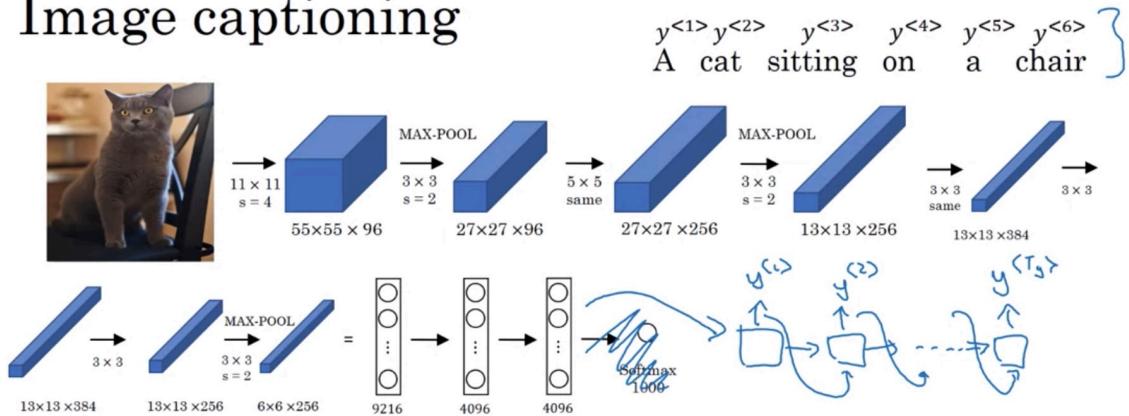
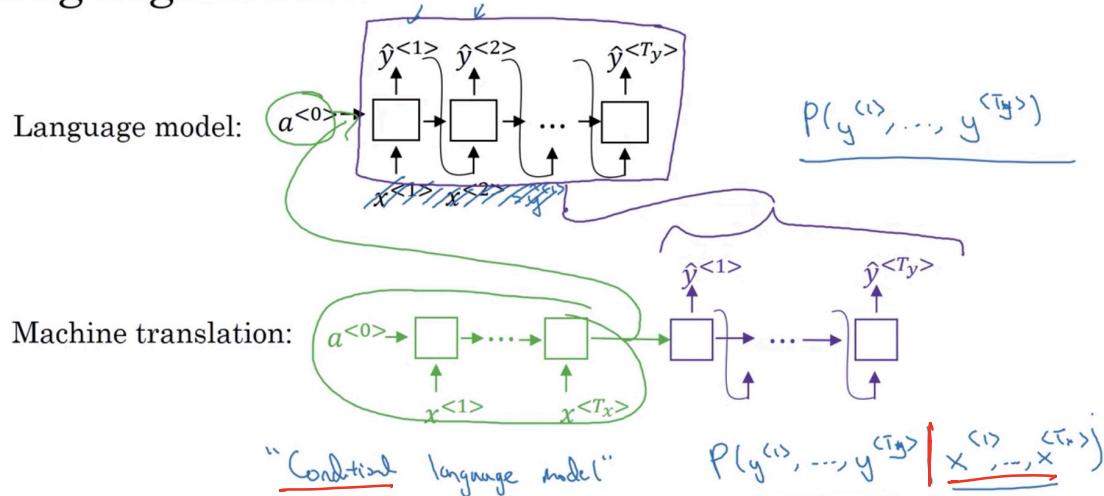


# Image captioning



Machine translation as building a conditional language model



Why not greedy Search? i.e. Pick one optimal word at a time.

Consider :  $\rightarrow$  Jane is visiting Africa in Sept.  
 $\rightarrow$  Jane is going to be visiting Africa in Sept.  
 more common word, more likely chosen.

first sentence is optimal as it's less verbose.

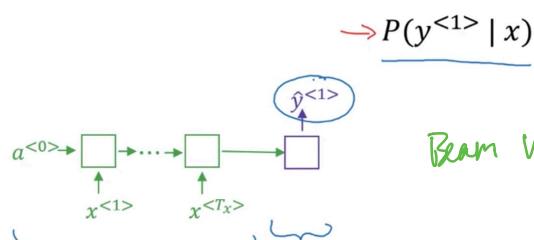
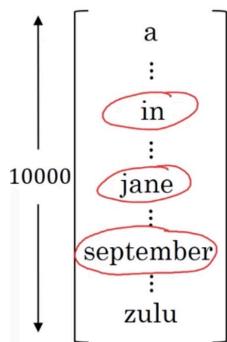
but  $P(\text{Jane is going } | x) > P(\text{Jane is visiting } | x)$   
 therefore picked 2nd sentence.

Brutal Search :  $10000^{10}$  possibility

Solution : Approximate Search Algorithm.

## Beam Search

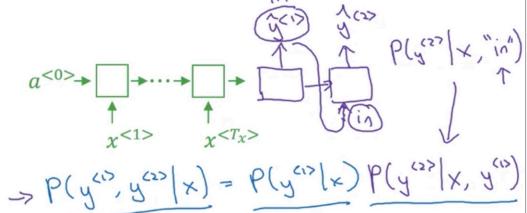
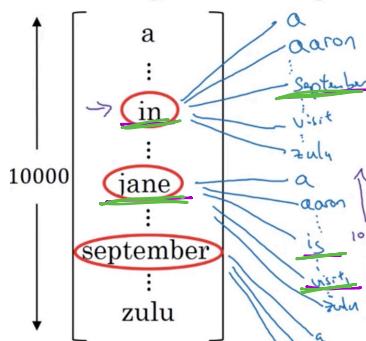
Step 1



Beam Width = 3

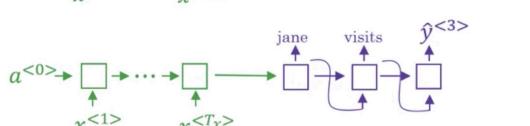
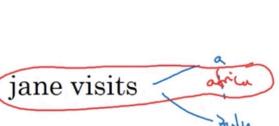
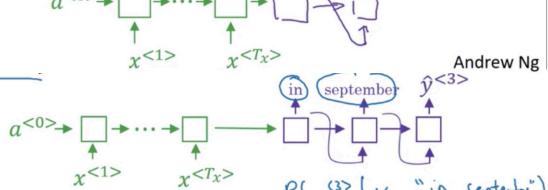
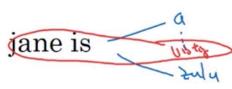
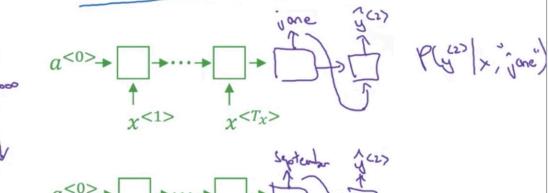
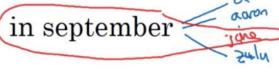
run through 1st estimate  $\rightarrow$  generate probability for 10,000  $\rightarrow$  keep top 3

Step 1

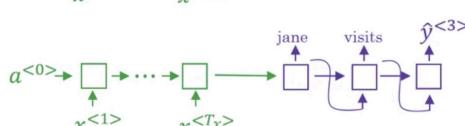


3 copies of  
the same network

Step 2



jane visits



$P(y^{<1>}, y^{<2>} | x)$

jane visits africa in september. <EOS>

## Length normalization

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

$p(y^{<1>} \dots y^{<T_y>} | x) = \frac{P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>}) \dots}{P(y^{<T_y>} | x, y^{<1>}, \dots, y^{<T_y-1>})}$   
 $\log p(y | x) \leftarrow$   
 $p(y | x) \leftarrow$

$$\arg \max_y \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

$\frac{1}{T_y^\alpha} \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$   
 $\alpha = 0.7 \quad \frac{\alpha=1}{\alpha=0}$

compute the score of searcher

The larger the beam width  $B$ , the more computationally expensive.  
 Beam Search does not guarantee optimum like a BFS/DFS.

eg: ✓  $y^*$  Human: Jane visited Africa in Sept.  
 ✗  $\hat{y}$  machine: Jane visited Africa last Sept.

→ RNN ?  $\Rightarrow P(y^* | x)$  [test if  $P(y^* | x) \geq P(\hat{y} | x)$ ]  
 → Beam ?  $\Rightarrow \arg \max_y P(y | x)$

Case 1:  $P(y^* | x) > P(\hat{y} | x)$

Beam choose  $\hat{y}$ , but  $y^*$  attains higher  $P(y | x)$

→ Beam Search at fault.

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

$y^*$  is a better translation. But RNN predicted  $P(y^* | x) < P(\hat{y} | x)$

→ RNN at fault.

## Error analysis process

Human	Algorithm	$P(y^* x)$	$P(\hat{y} x)$	At fault?
Jane visits Africa in September.	Jane visited Africa last September.	$2 \times 10^{-10}$	$1 \times 10^{-10}$	(B)
...	...	—	—	R
...	...	—	—	B
				R
				R
				:

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

Andrew Ng

Given a French Sentence, there could be a few version that is equally as good. How do we evaluate?

Bleu: Bilingual evaluate understanding

Evaluating machine translation (unigram)

French: Le chat est sur le tapis.

- Reference 1: The cat is on the mat. 2 appears
- Reference 2: There is a cat on the mat.
- MT output: the the the the the the the

$$\text{Precision: } \frac{7}{7}$$

$$\text{Modified precision: } \frac{2}{7} \leftarrow \begin{matrix} \text{Count}_{\text{clip}}(\text{"the"}) \\ 2 \end{matrix} \quad \leftarrow \begin{matrix} \text{Count}(\text{"the"}) \\ 7 \end{matrix}$$

Bleu  
bilingual evaluation understanding

$$P_i = \frac{\sum_{\text{unigram}_i} \text{Count}_{\text{clip}}(\text{unigram})}{\sum_{\text{unigram}_i} \text{Count}(\text{unigram})}$$

[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation] Andrew Ng

## 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	Count clip	
the cat	2 ←	1 ←	
cat the	1 ←	0	
cat on	1 ←	1 ←	
on the	1 ←	1 ←	
the mat	1 ←	1 ←	
			4 6.

$$P_n = \frac{\sum_{n\text{-gram } g} \text{Count clip}(n\text{-gram})}{\sum_{n\text{-gram } g} \text{Count}(n\text{-gram})}$$

## Bleu details

$p_n$  = Bleu score on n-grams only

$P_1, P_2, P_3, P_4$

Combined Bleu score:  $\text{BP} \exp\left(\frac{1}{4} \sum_{n=1}^4 p_n\right)$

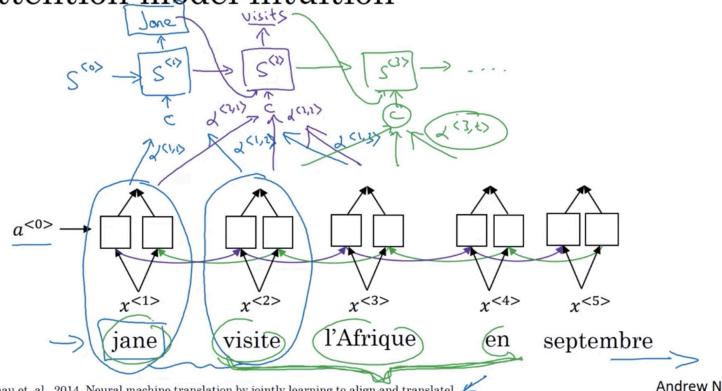
$\text{BP} = \text{brevity penalty}$

$$\text{BP} = \begin{cases} 1 & \text{if MT\_output\_length} > \text{reference\_output\_length} \\ \exp(1 - \text{MT\_output\_length}/\text{reference\_output\_length}) & \text{otherwise} \end{cases}$$

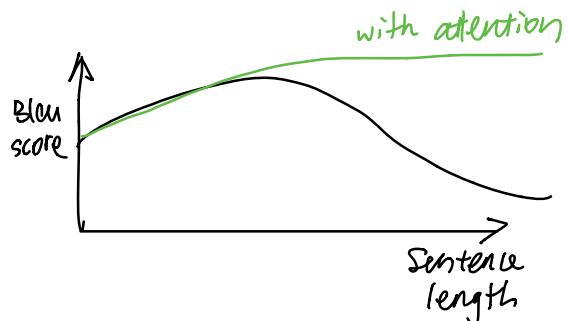
## Attention Model

The problem of long sequence:

Attention model intuition



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



Intuition

$S$  - hidden states of another RNN (for English translation)

$\alpha$  - attention weights

## Attention Model

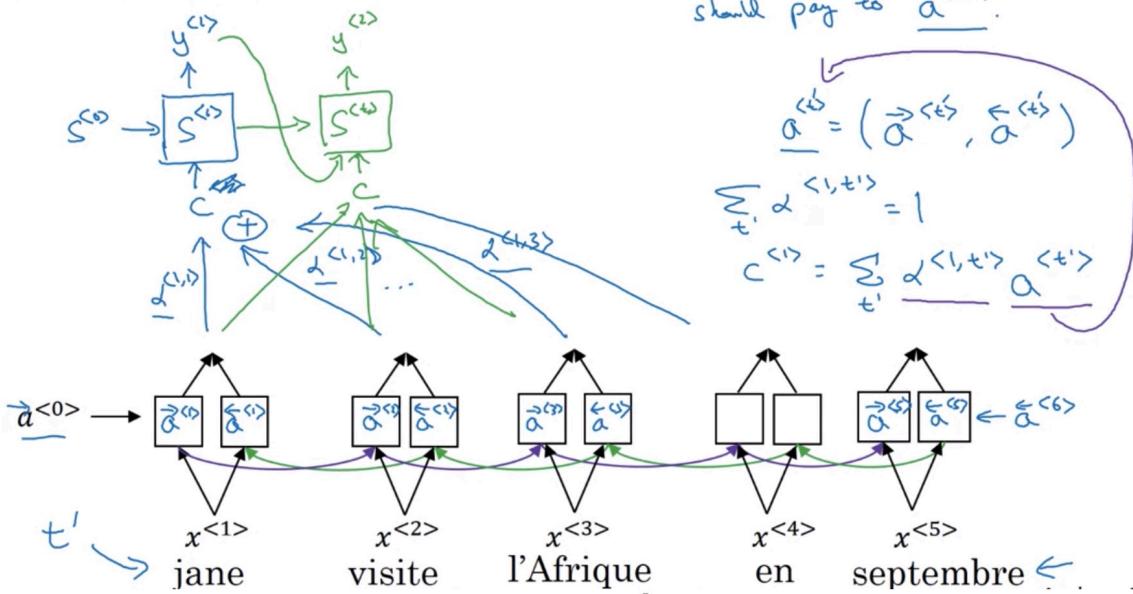
$$a^{<t>} = (\vec{a}^{<t>}, \vec{a}^{<t>})$$

- concat forward/backward RNN

$\alpha^{<1,2>}$  how much the 1st word translated depends on amount of attention the context 2nd word in original sentence

context  $\alpha^{<t>}$  weighted sum of activations

## Attention model



$\alpha^{<t,t>} = \text{amount of "attention" } y^{<t>} \text{ should pay to } \underline{a^{<t>}}$

$$\underline{a^{<t>}} = (\vec{a}^{<t>}, \vec{a}^{<t>})$$

$$\sum_{t'} \alpha^{<1,t'>} = 1$$

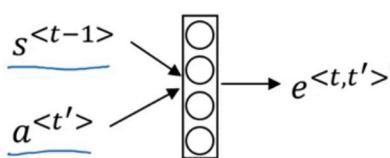
$$c^{<1>} = \sum_{t'} \alpha^{<1,t'>} \underline{a^{<t'>}}$$

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

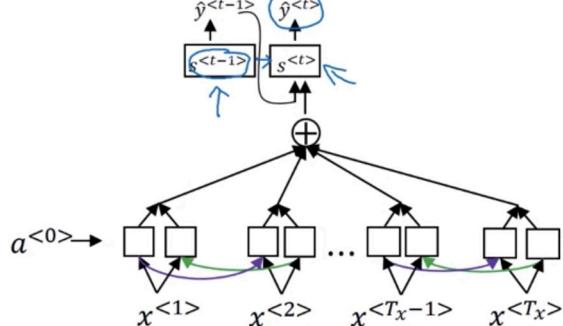
$\alpha^{<t,t'>} = \text{amount of attention } y^{<t>} \text{ should pay to } \underline{a^{<t'>}}$

sum to 1

$$\rightarrow \underline{\alpha^{<t,t'>}} = \frac{\exp(e^{<t,t'>})}{\sum_{t'=1}^{T_x} \exp(e^{<t,t'>})} \text{ softmax}$$



usually simple neural network to combine the two

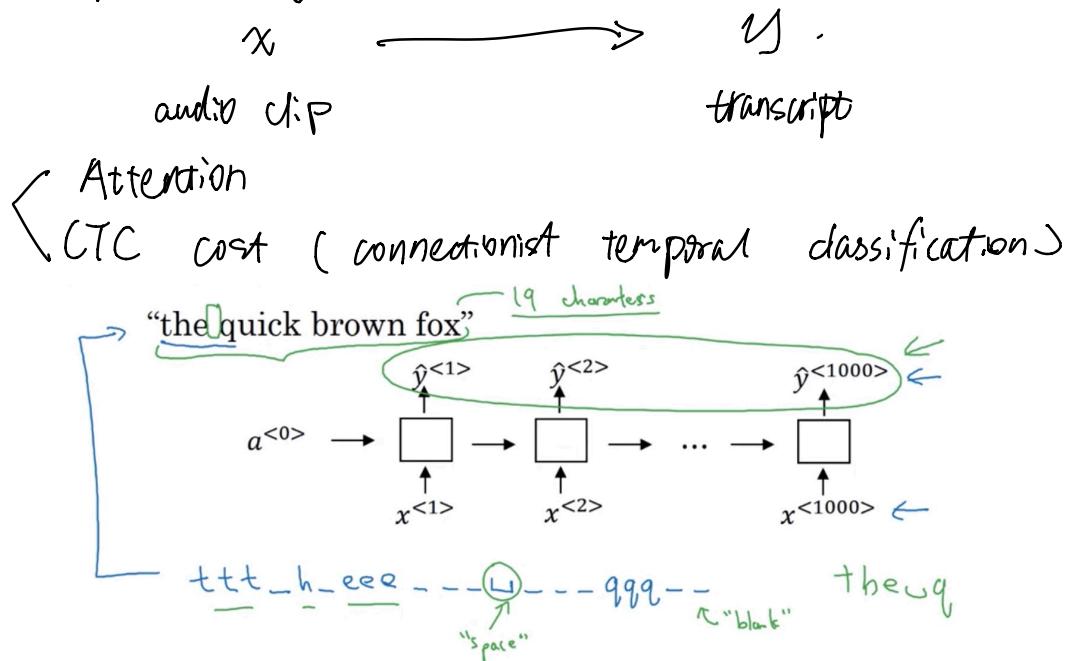


$\Rightarrow$  Because we don't know how much should we depend on prev translation/context

Downside : Runs quadratic times  $\alpha^{<t,t>} = O(T_x \times T_y)$   
 input output

Attention can apply to image captions as well.

## Speech Recognition



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

## Trigger Words

