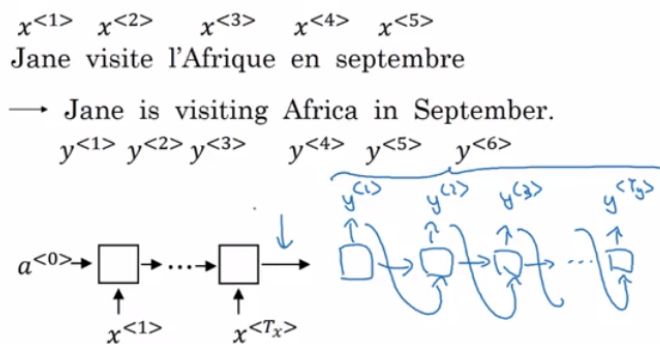


Various sequence to sequence architectures

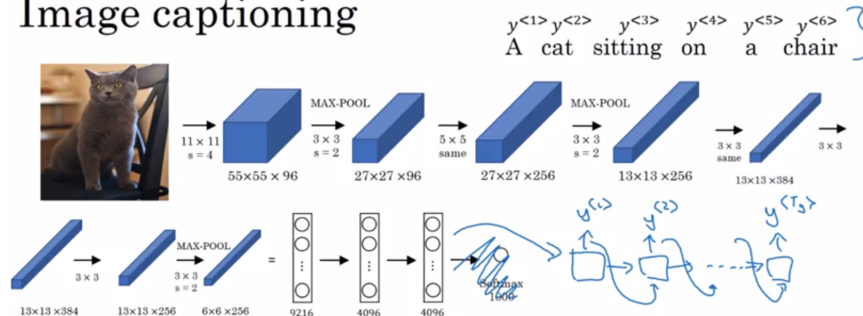
- basic model
 - [encoder] → [decoder]
 - for machine translation task:
 - NN architecture

Sequence to sequence model

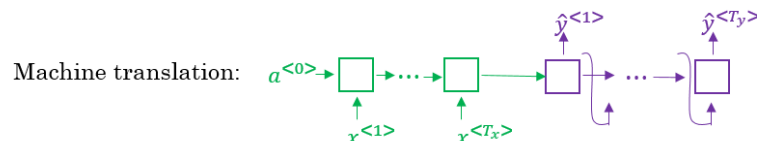
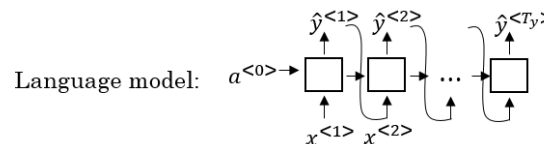


- for image captioning
 - uses a pretrained CNN (like AlexNet) as an encoder for the image, and the decoder is an RNN.

Image captioning



- Picking the most likely model
 - comparison between language model and machine translation
 - NN



- Problems formulations are different:
 - In language model: $P(y^{<1>} \dots y^{<Ty>} | x^{<1>} \dots x^{<Tx>})$
 - In machine translation: $P(y^{<1>} \dots y^{<Ty>} | x^{<1>} \dots x^{<Tx>})$

- problem:

- we want the best translation instead of random answer

Finding the most likely translation

Jane visite l'Afrique en septembre.

$$P(y^{<1>}, \dots, y^{<T_y>} | x)$$

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.

$$\arg \max_{y^{<1>}, \dots, y^{<T_y>}} P(y^{<1>}, \dots, y^{<T_y>} | x)$$

- aim:

- best translation

- why not a greedy search?

- doesn't really work
- the words in sentence have relation in long distance instead of just adjacent one.
- distribution space is too large for greedy search.

- Beam Search

- a heuristic search algorithm

- by iterating and selecting the best results in each "B" outcomes.
- optimize:

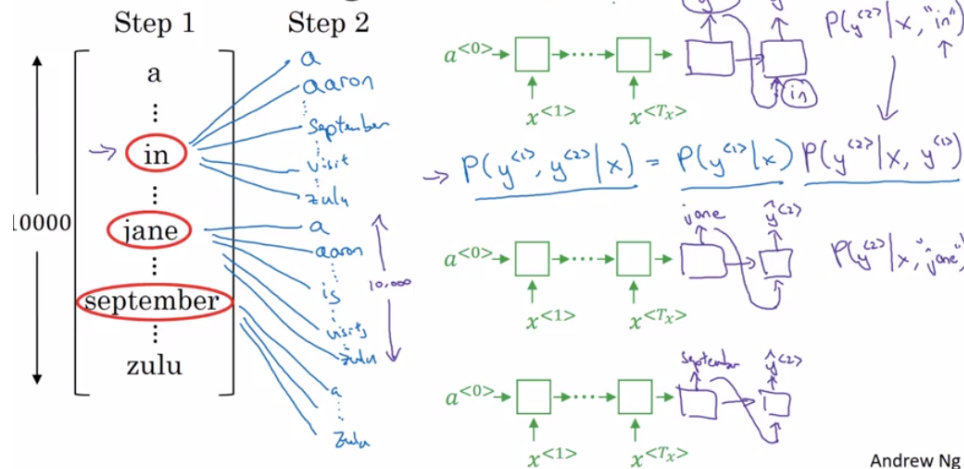
$$\arg \max_{y^{<1>}, \dots, y^{<T_y>}} P(y^{<1>}, \dots, y^{<T_y>} | x)$$

- caution. It's a bit computational expensive as **B** growing bigger.

- parameter B : the beam width.

- means that the top "B" number of possible outcomes.
- if B=1, it turns to be a greedy search

Beam search algorithm (B=3)



Andrew Ng

- refinements to beam search

- optimization:

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

may be too small

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

MLE

- using Maximum Likelihood Estimate
- another problem: The optimization function prefers small sequences rather than long ones.
 - Because multiplying more fractions gives a smaller value, so fewer fractions - bigger result.
- So: dividing by the number of elements in the sequence.

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

- how to choose B?

Beam search discussion

Beam width B?

1 → 10, 100, 1000, → 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 $\arg \max_y P(y|x)$.

- what's more, error analysis could help.
- error analysis
 - help to find out which limit model performance, RNN or Beam Width
 - EX.

Example

→ RNN

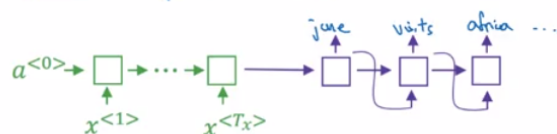
→ Beam Search

Jane visite l'Afrique en septembre.

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

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

RNN computes $P(y|x)$



- calculate $P(y^* | X)$ and $P(\hat{y} | X)$
 - if $(P(y^* | X) > P(\hat{y} | X))$:
 - Beam search is at fault.
 - else: $\# (P(y^* | X) \leq P(\hat{y} | X))$

- RNN model is at fault.
- so, make a table, then get counts and decide what to work on next.

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

• BLEU score.

- BLEU stands for *bilingual evaluation understudy*.
- task:
 - given a sentence in a language there are many possible good translation. How to evaluate our results?
- intuition:
 - choose the result which pretty closes to any of the references provided by humans
- algorithm:
 - quite like a modified kind of "Template matching"
 - compute p_n :

$$p_n = \frac{\sum_{ngram \in \hat{y}} count_{clip}(ngram)}{\sum_{ngram \in \hat{y}} count(ngram)}$$

- combined Bleu Score:
 - BLEU score = BP * exp(1/n * sum(Pn))
 - BP here means:

BP: brevity penalty

$$BP = \begin{cases} 1 & \text{if } \underline{MT_output_length} > \underline{reference_output_length} \\ \exp(1 - MT_output_length/reference_output_length) & \text{otherwise} \end{cases}$$

- It turns out that if a machine outputs a small number of words it will get a better score so we need to handle that.
- Ex.

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	4
cat on	1 ←	1 ←	6
on the	1 ←	1 ←	
the mat	1 ←	1 ←	

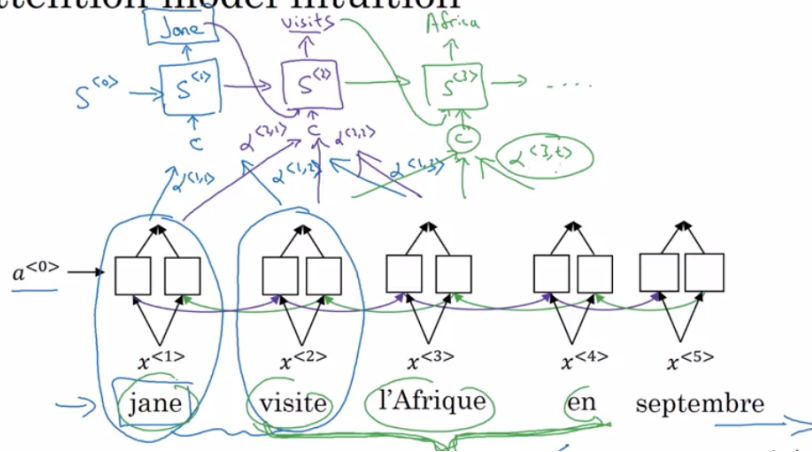
here n-grams, n=2

- application:
 1. machine translation
 2. image captioning

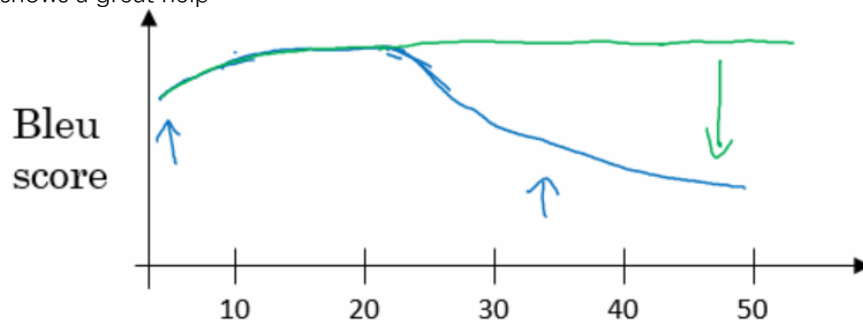
Attention Mechanism

- intuition:
 - for a long sequence, instead of understanding it as a whole, we may pay more attention on certain part of it.
 - in fact, as sequence growing longer, the Bleu scores would go down conspicuously.
 - to adapt long sequence, we add a attention mechanism part to our model.

Attention model intuition

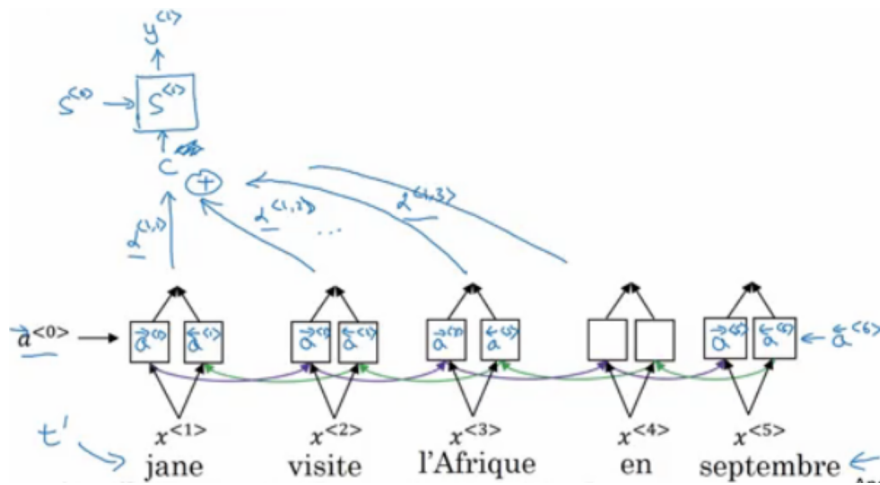


- which shows a great help



- Blue is the normal model, while green is the model with attention mechanism.

- details
 - structure:



- add on a RNN/LSTM/GRU
- formula
 - for each \mathbf{a} :

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

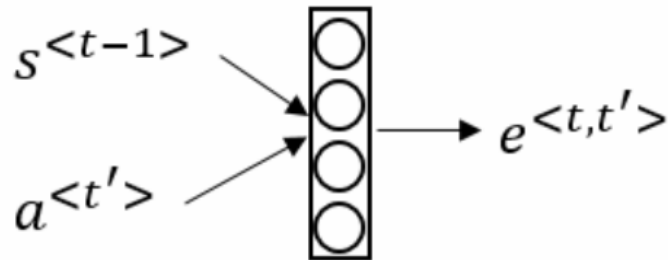
- and by summing up, it would be

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

- for context \mathbf{c} :

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

- for error \mathbf{e}

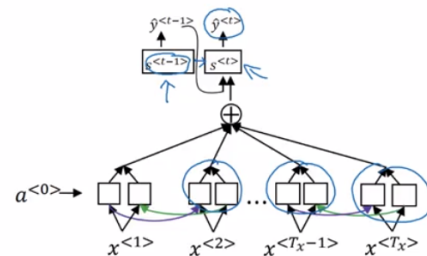
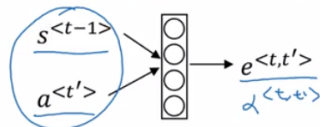


- wrap up:

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

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

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



- visualizing weights

