

Basic model:-

Sequence to sequence model

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$

Jane visite l'Afrique en septembre

→ Jane is visiting Africa in September.

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

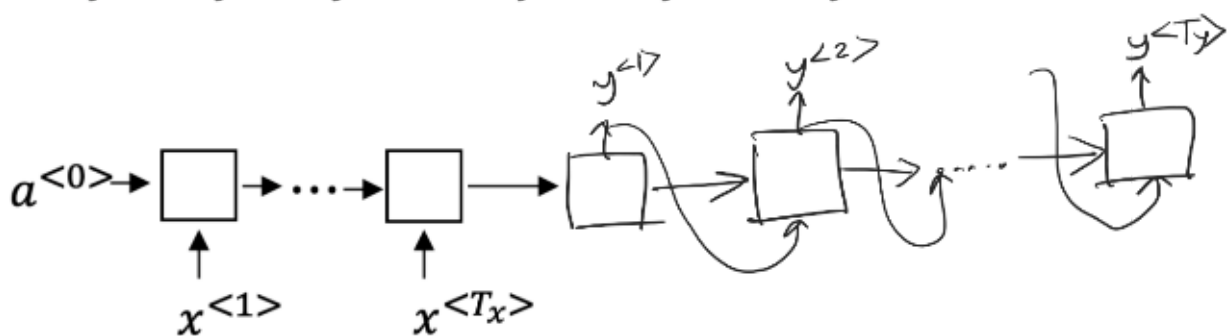
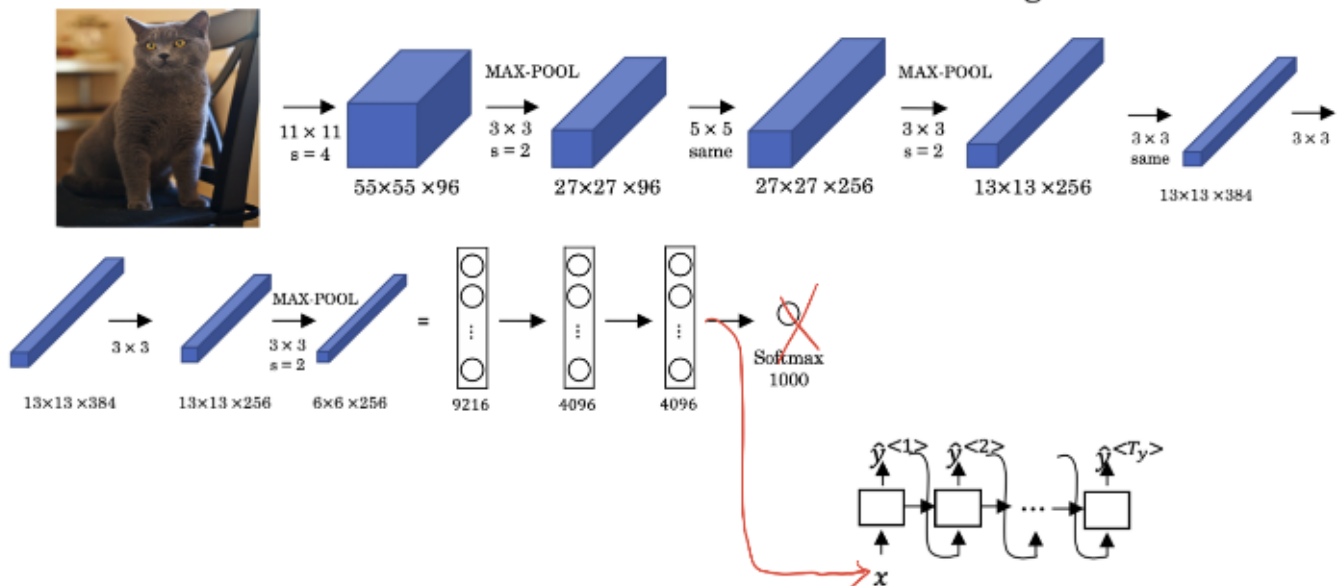


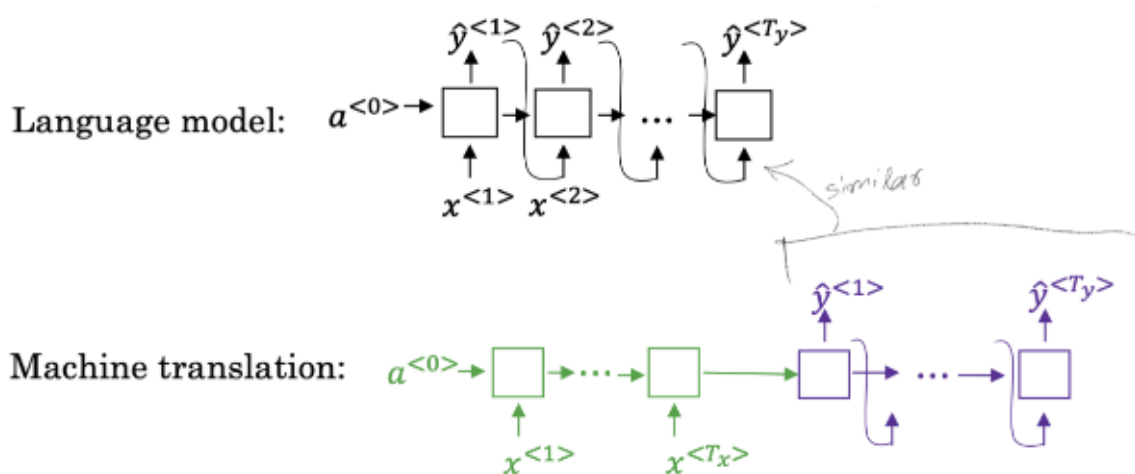
Image captioning

$y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad y^{<4>} \quad y^{<5>} \quad y^{<6>}$
A cat sitting on a chair

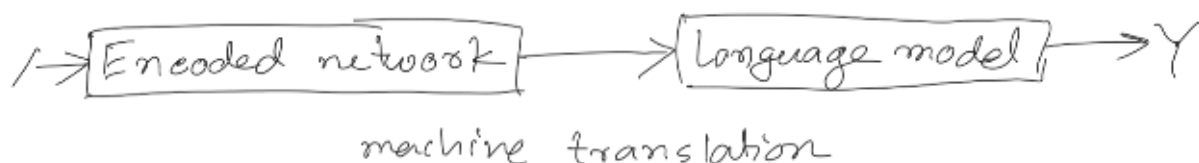


Machine translation as building a conditional language model

estimate the probability of a sentence



Machine translation is pretty much similar to the language model. Only difference is instead of starting the network to all zero, it has an encoded network that figure out some representation of the input sentence. And it use this representation as an input of the decoder network (which is identical to language model.)



Machine translation provides the probability of different sentence being the translation of input sentence. This is why it is called conditional language model.

Finding the most likely translation

Jane visite l'Afrique en septembre.

possible translation:

→ Jane is visiting Africa in September.

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

\downarrow English \downarrow French
 \downarrow

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

Machine translation will provide the probability of each sentence

Our goal $\rightarrow \arg \max_{y^{<1>, \dots, y^{<T_y>}} P(y^{<1>, \dots, y^{<T_y>} | x) \leftarrow$ Done by beam search

Why not greedy search?

Greedy search select the word one at a time. for example:-

→ Jane is visiting africa in september.

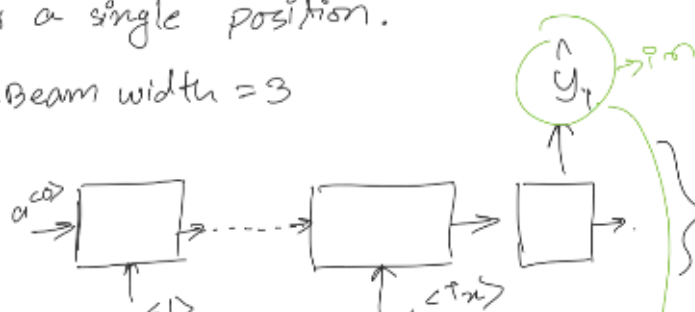
→ Jane is going to be visiting Africa in september.

"going" is more common words than visiting. so greedy search might select the second sentence instead of first one. While the first one is more appropriate.

Beam search:-

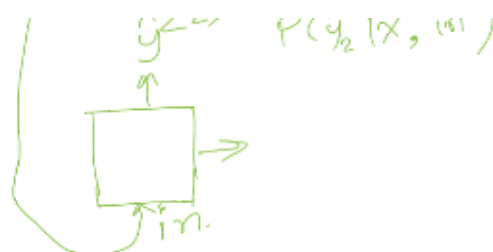
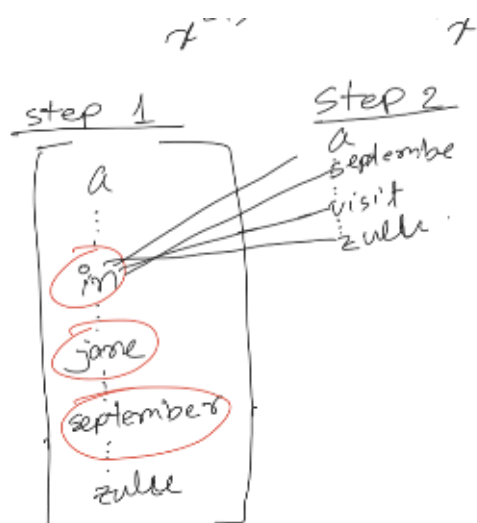
Greedy search just considers one value in each position & move on. where Beam search consider multiple value for a single position.
beam width

Let, Beam width = 3



step 1, beam search remembers 3 most likely words in this case: in, zulu, september

1, 2, 3



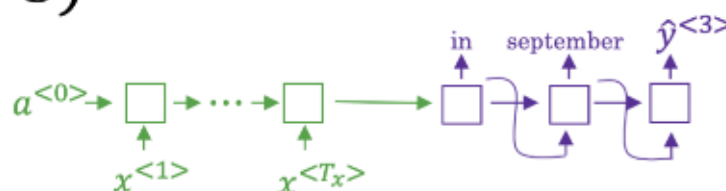
$$P(y^{<1>}, y^{<2>} | x) = P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>})$$

Repeat this for jane & september

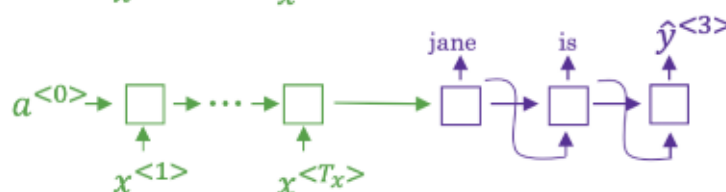
⇒ If our vocabulary size is 10,000 then in the second stage we will calculate for 30,000 possibility

Beam search ($B = 3$)

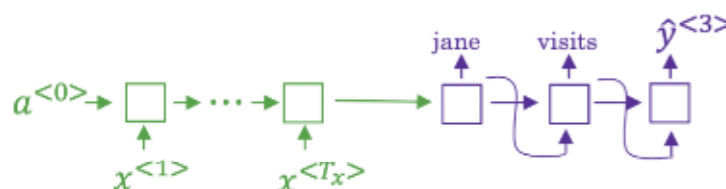
in september



jane is



jane visits



Refinements to beam search:-

Length normalization:-

$$\arg \max_y \prod_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>}) = \frac{P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>})}{P(y^{<3>} | x, y^{<1>}, y^{<2>})} \dots$$

some of the probability are very tiny so it will make the value very small.

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

$t=1$
 This approach prefers shorter sentences than longer sentences. Even if longer sentence is more correct. Cause in longer sentence we will have more probability that will be multiplied together. Hence, make the probability even less. To overcome this, we use following eqⁿ:-

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

→ normalize with length. α = hyperparameters.

$\alpha = 1$ = normalize along the length.

$\alpha = 0$ = Don't normalize at all

$\alpha = 0.7$ = works better in practice

Larger beam-width (B) means lot's of computation.

larger B :- Better result, slower, more memory

small B :- worse result, faster, less memory

$B \rightarrow 10$ is more common in production.

$B \rightarrow 100$ is considered huge in production.

$B \rightarrow 1000, 10,000 \rightarrow$ used in academia

$B \rightarrow 1 \rightarrow$ is greedy search.

Beam search is a heuristic search. So it doesn't always give the correct result.

Error analysis in Beam search:-

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

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

Case 1:

Beam search chose \hat{y} . But y^* attains higher $P(y|x)$.

Conclusion: Beam search is at fault.

Case 2:

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

Conclusion: RNN model is 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.	2×10^{-10}	1×10^{-10}	B
		—	—	R
		—	—	B
				R
				R
				B

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

Andrew Ng

Blue score

If there are multiple sentences that are good translations
How would we choose between them?

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.

Both are good translations

Blue looks for the human reference for this sentence

look each of the words and see if they appeared in the references.

Precision: $\frac{7}{7}$
 7 → support
 7 → total 7 words (blue)

Modified precision: $\frac{2}{7}$
 2 → Count up.
 7 → count

But this is not a good translation still we are getting precision 1. So we will modified precision, where we will give credit upto the maximum numbers of times it appeared in the references. "The" appeared 2 times in the references 1, and 1 times in reference-2. So we will take maximum of this numbers.

Bleu score on bigrams \rightarrow Pairs of words appearing next to each others.

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.

<u>Bigram</u>	<u>Count</u>	<u>count_{clip}</u>
the cat	2	1
cat the	1	0
cat on	1	1
on the	1	1
the mat	$\frac{1}{6}$	$\frac{1}{4}$

modified precision = $\frac{4}{6} = \frac{2}{3}$

$$p_1 = \frac{\sum_{unigram \in \hat{y}} count_{clip}(unigram)}{\sum_{unigram \in \hat{y}} count(unigram)}$$

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

p_n = Bleu score on n-grams only

Combined Bleu score:

BP > Brevity penalty.

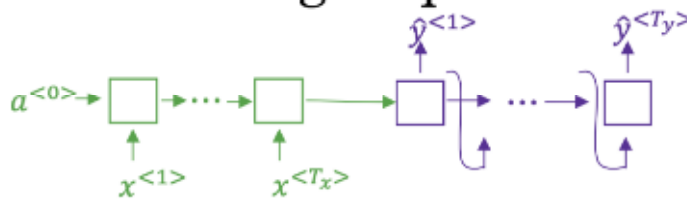
$$BP = \begin{cases} 1 & \text{if MT output length} > \text{reference output length} \end{cases}$$

$$D1 = \begin{cases} \exp(1 - \text{MT_output_length} / \text{reference_output_length}) & \text{otherwise} \\ \text{reference_output_length} / \text{MT_output_length} \end{cases}$$

shorter sentence tends to have better precision. so we will penalize the precision if we get shorter sentence.

Attention model:-

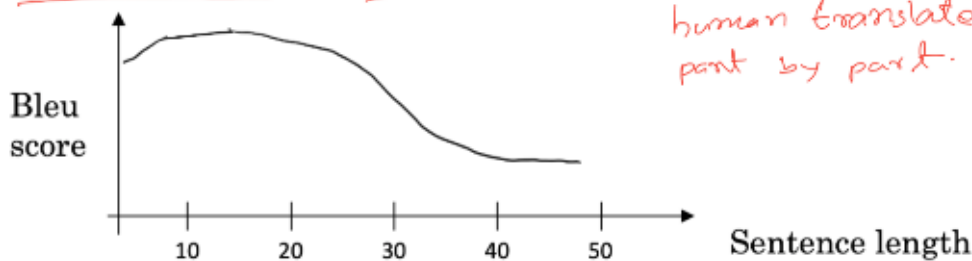
The problem of long sequences



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.

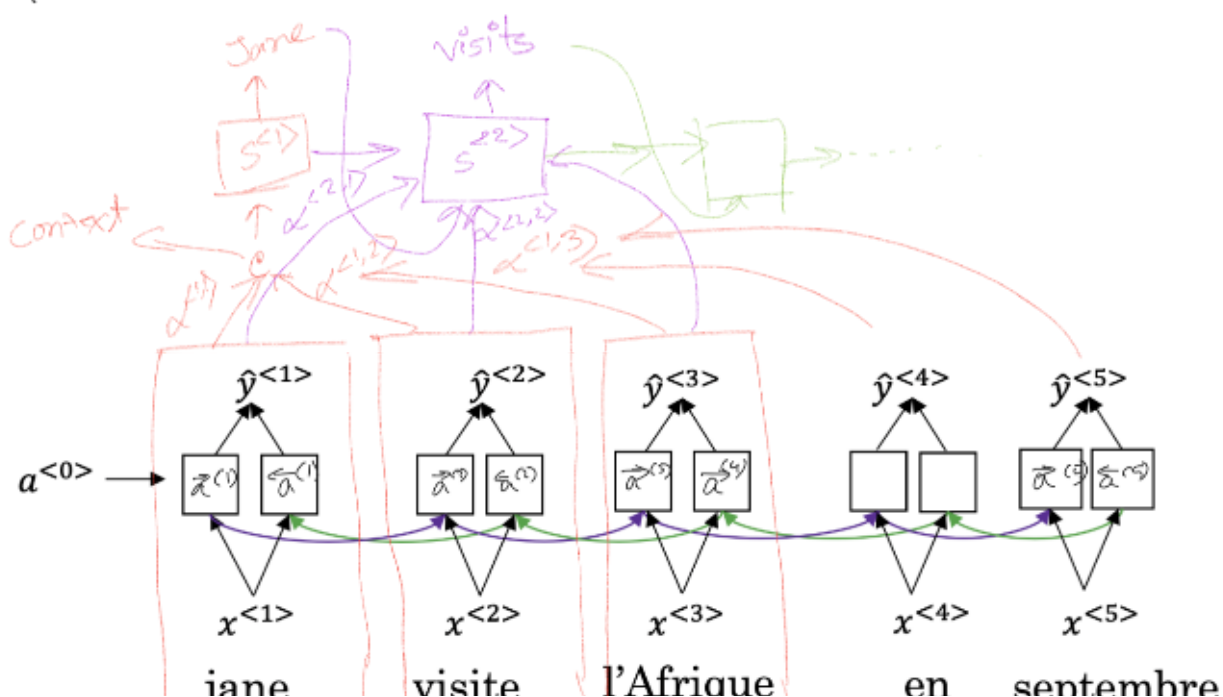
human translate long sentence part by part.



Andrew Ng

very short & very long sentences are harder to train.

Which word to pay attention:-



α = how much we should pay attention to this word =
 α = attention weight

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

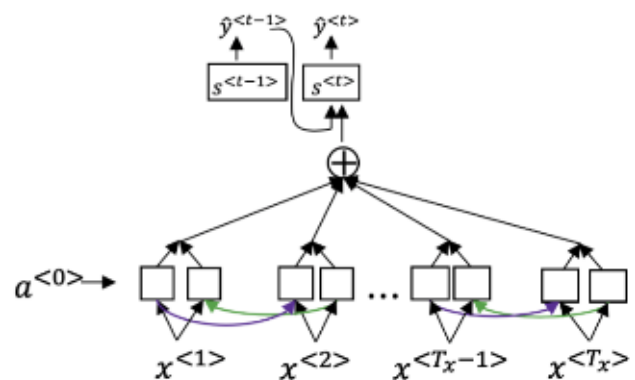
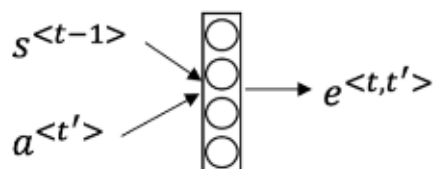
$$\sum_{t'} \alpha^{<t, t'>} = 1.$$

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

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

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

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



Speech recognition:-