

Week 3

Various Sequence to Sequence Architecture

Basic Models:

Sequence to sequence model:

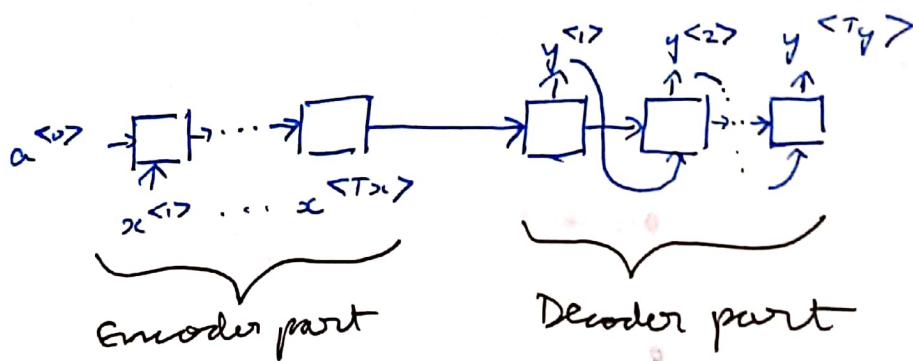
This is used in translators, for example,

French:

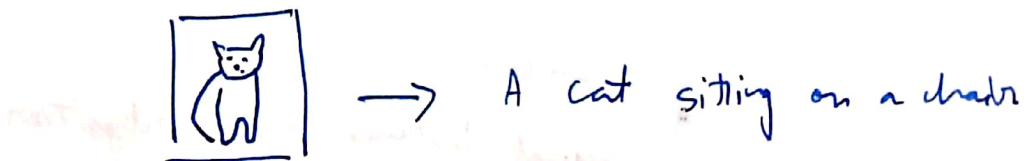
Jane visite l'Afrique en septembre
 $x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$

English:

Jane is visiting Africa in September
 $y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad y^{<4>} \quad y^{<5>} \quad y^{<6>}$



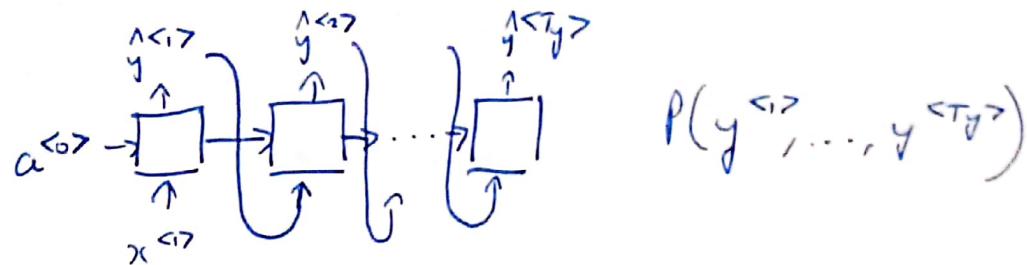
You can even use a similar method for image captioning:



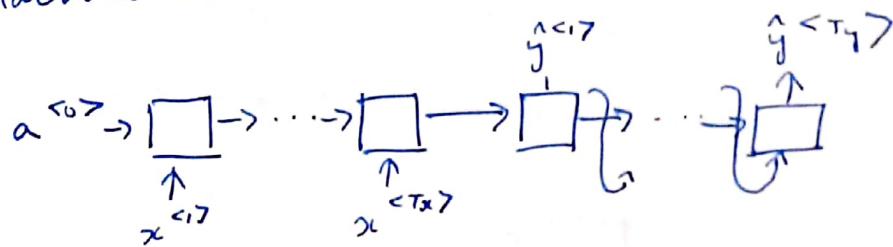
We do this by replacing the encoder part with an AlexNet (the final softmax layer is removed and is instead connected to the decoder part)

Picking the most likely sentence:

Language model (that we learnt in week 2):



Machine translation:



→ Similar to language model, just that rather than inputting $a^{<0>}$ in the decoder, in machine translation we input the output of the encoder.

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

Instead of modelling any sentence (like in language model), it is now modelling the probability of the output English translation, conditions on some input French sentence.

→ In language model, we would select the next word at random (from the list of possible word predictions), but here we pick the best.

$$\underset{y^{<1>} \dots, y^{<T_y>}}{\text{arg max}} P(y^{<1>} \dots, y^{<T_y>} | z)$$

→ Why don't we use greedy search algorithm?

Greedy Search algorithm is a CS algorithm that finds the 1st best word, 2nd best word, 3rd best word and so on. By best word, I mean $P(y^{<1}, y^{<2}, \dots, y^{<T} | x)$ is the highest. But consider the following:

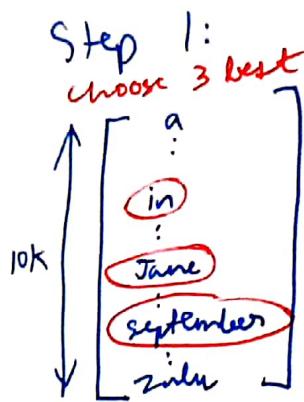
- Jane is visiting Africa in September
- Jane is going to be visiting Africa in September

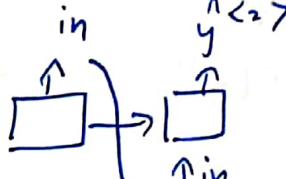
Here $P(\text{Jane is going} | x) > P(\text{Jane is visiting} | x)$ So it'll pick 2 rather than 1, even though 1 is better (~~2~~) (since going is a more probable word than visiting, it'll pick that). So this doesn't work.

Beam Search Algorithm:

Here we use a parameter B .

If $B = 3$, rather than picking the best word like in greedy search, we pick the top 3 words. Then we create 3 parallel networks where we hardcode and feed these 3 words separately for the second cell of the decoder RNN. The output of the 2nd cell will be 10K words for each RNN (So total $3 \times 10K = 30K$). Out of 30K we pick the best 3 two words (1st word + 2nd word), and this goes on.

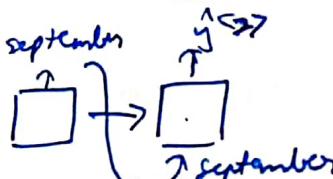


1) Encoder \rightarrow 

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

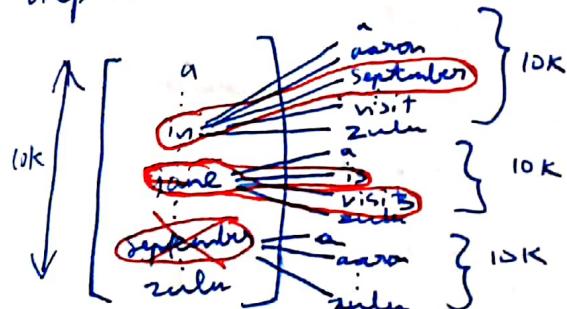
2) Encoder \rightarrow 

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

3) Encoder \rightarrow 

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

Step 2:



Top 3 ~~$P(y^{<1>} | x), P(y^{<2>} | x)$~~

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

so next we find $P(y^{<3>} | x, "in\,september")$

$$P(y^{<3>} | x, "Jane\,is")$$

$$P(y^{<3>} | x, "Jane\,visits")$$

Notice how "September" wasn't in the top 3
so it was relegated.

If we do this for further words we will get:
Jane visits Africa in September. <EOS> ✓

Note: $B=1$ is greedy search algorithm.

Refinements to beam search:

Length normalization:

$$\arg \max_{y} \prod_{t=1}^{Ty} p(y^{} | x, y^{<1>}, \dots, y^{})$$

This is nothing but $p(y^{(1)} \dots y^{(Ty)} | x) = p(y^{(1)} | x) p(y^{(2)} | x, y^{(1)}) \dots p(y^{(Ty)} | x, y^{(1)}, \dots, y^{(Ty-1)})$

~~But when we~~

~~These numbers will~~

→ These terms like $p(y^{(1)} | x)$ will be between 0 and 1, so what happens when we multiply many of these is that we get a very small number that our floating value can't hold (because of low precision). So instead we use this formula:

$$\arg \max_{y} \sum_{t=1}^{Ty} \log p(y^{} | x, y^{<1>}, \dots, y^{})$$

→ But since the more words we have, the lesser the probability becomes, this also tends to choose shorter sentence. So we normalize it.

$$\arg \max_{y} \frac{1}{Ty^\alpha} \sum_{t=1}^{Ty} \log p(y^{} | x, y^{<1>}, \dots, y^{})$$

where $\alpha = 0.7$ usually

Since we are like taking the average, it won't penalize longer sentences.

before we discussed that we will choose the best translation (for $B = 3$), we will score them using this formula and choose the best as final output.

How to choose Beam Width B ?

→ Large B : better results, slower

small B : worse results, faster

considered
big for
this
application

→ In production systems you would use $10 - 100$

In research, where best result is required,

$1000 - 3000$

→ The gain in results from $1 \rightarrow 10$ will be larger than $1000 \rightarrow 3000$ (higher the number, the less performance improvements we will see).

→ 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 any $\max_y P(y|x)$.

Error Analysis in Beam Search:

Beam search is an approximation algo, so it might not give the most accurate result always. If an error occurs, how do you find out if the problem is with beam search or your RNN?

Example:

Jane visite l'Afrique en septembre

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

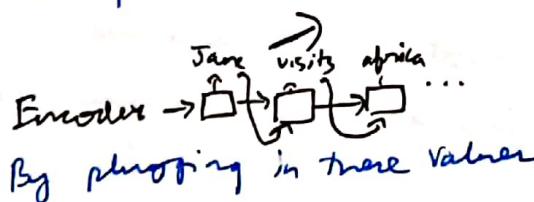
Algorithm: Jane visited Africa last September. (\hat{y})

Was the RNN or the beam search algo to blame for this bad translation?

You can increase the beam width, so it gets better results, but this is like getting more data - it's not the most effective way.

There is a better way.

1. Compute $P(y^*|x)$ and $P(\hat{y}|x)$



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

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

Conclusion: Beam search is at fault

Its job was to find the highest $P(y|x)$ and it failed at that by choosing \hat{y} over y^*

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

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 Boxen:

For all the translations it got wrong,

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^{-5}	B
-	-	-	-	R
-	-	-	-	B
:	:	:	:	.

Find the fraction of errors due to beam search vs RNN model.

- Only if a majority of the errors are due to beam search, ~~then~~ only then its a good idea to increase beam width B.
- If a majority of the errors are due to RNN, you can perform a deeper analysis on whether to add regularization, more complete network, etc.

Bleu score:

One of the challenges of machine translation is that given a French sentence there can be multiple equally good English translations. The BLEU score is a measure of similarity to sort this.

BLEU → Bilingual Evaluation Understudy

Take an example,

French: Le chat est sur le tapis

~~And it has 2 equally good translations~~

Reference 1: The cat is on the mat

Reference 2: There is a cat on the mat

These are two references given by a human

The intuition is that we are going to look at the machine generated output and see if the types of words it generates appear in at least one of the human generated references.

~~native Translation~~

MT output: the the the the the the

Precision: $\frac{7}{7} \leftarrow \text{no. of words in MT output that appear in either ref1 or ref2}$
 $\frac{7}{7} \leftarrow \text{no. of words in MT output}$

So this isn't a good way of measuring precision, so we use a modified precision

Modified Precision: $\frac{2}{7} \leftarrow$ no. of times the word appears in a reference
 The word appears 2 times in ref 1 and 1 time in ref 2 \leftarrow no. of words in MT output
 $\therefore \text{It is } 2$ \leftarrow time the word appears in MT output
 Count ("the")

Den 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
	6	4

$$\therefore \text{Score} = \frac{4}{6}$$

Formulas:

For unigram,

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

↑ means
unigram

↑ g is MT output

For n-gram,

$$P_n = \frac{\sum_{n\text{-grams} \in g} \text{count}_{\text{clip}}(n\text{-gram})}{\sum_{n\text{-grams} \in g} \text{count}(n\text{-gram})}$$

Bleu details:

P_n = Bleu score on n-gram only

Suppose we get P_1, P_2, P_3, P_4

The the combined Bleu Score is,

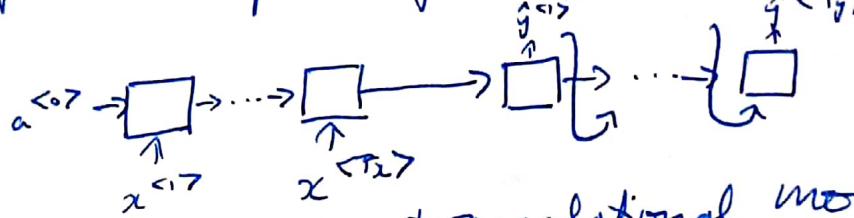
$$BP \exp \left(\frac{1}{n} \sum_{n=1}^4 P_n \right)$$

Here BP is Brevity Penalty. If you output short translations, it's easy to get high precision. But we don't want translations that are very short, so BP is an adjustment factor that penalizes translation systems that output short translations.

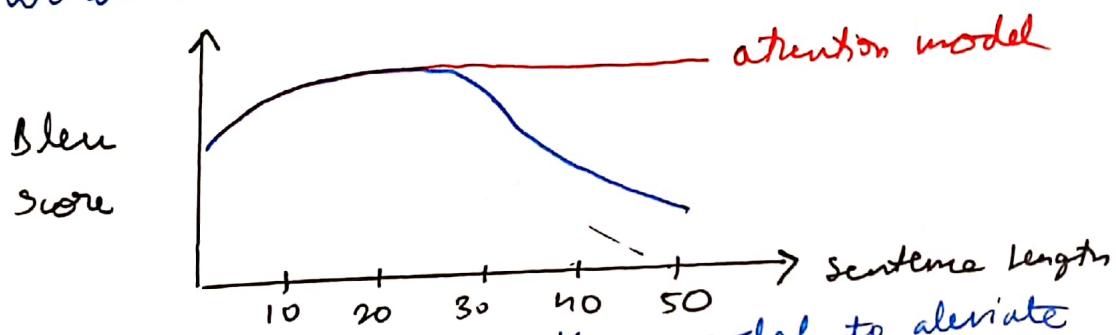
$$BP = \begin{cases} 1 & \text{if MT output length} \geq \text{reference output lengths} \\ \exp \left(-\frac{\text{MT output length}}{\text{reference output lengths}} \right) & \text{otherwise} \\ \exp \left(1 - \frac{\text{reference output length}}{\text{MT output length}} \right) & \text{otherwise} \end{cases}$$

Attention Model:

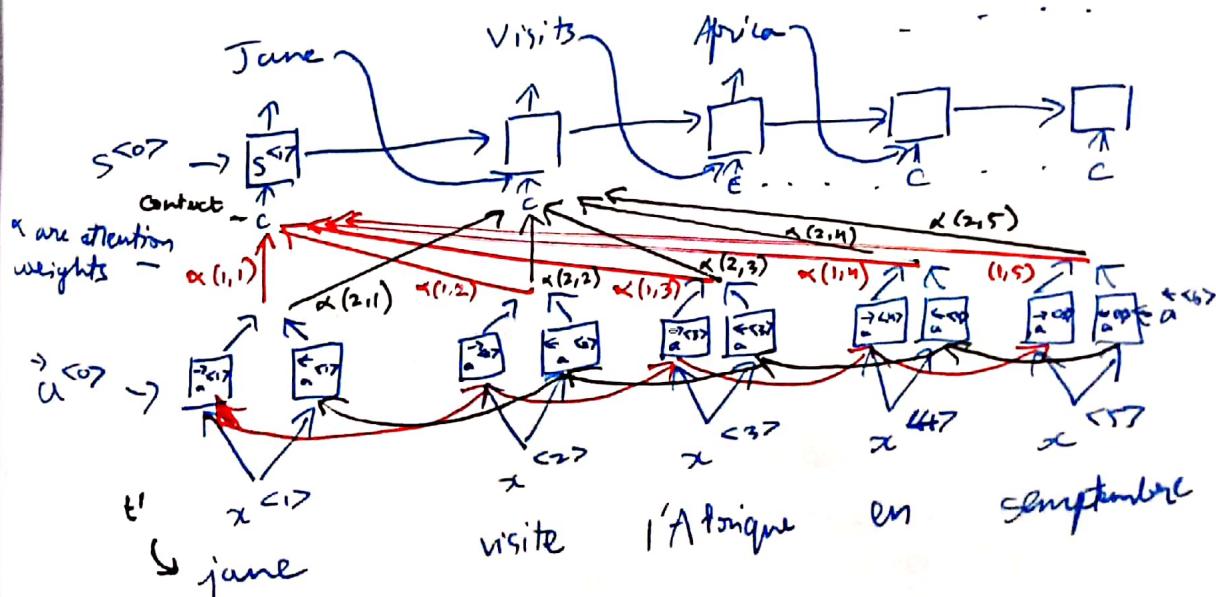
The problem of long sequences:



In the RNN machine translational model, it memorizes the entire large sentence and translates it in one shot. However a real human translator looks at a few words and translates those smaller parts one at a time. Therefore the performance of the encoder-decoder model drops as word count increases.



so we use an attention model to alleviate this.



$$\rightarrow \alpha^{<t'}> = (\vec{\alpha}^{<t'}}, \hat{\alpha}^{<t'}})$$

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

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

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

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

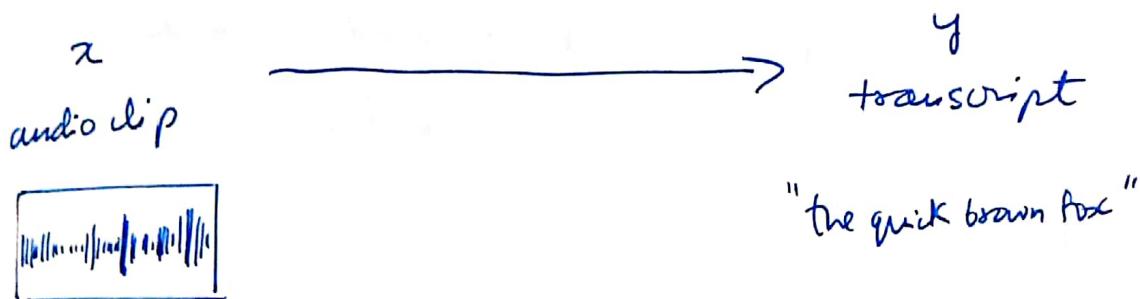
\rightarrow To get $e^{<t, t'}>$:

$$s^{<t-1>} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow e^{<t, t'}>$$
$$a^{<t'}> \rightarrow$$

Speech Recognition - Audio Data

Speech Recognition:

The ~~parrotion~~ speech recognition problem is:



- A common pre-processing step for audio data is to ~~do~~ convert the raw audio clip and generate a spectrogram (plotting the frequency).
- Before speech recognition systems used to be built using phonemes (like the is de, quick is k-wic, etc - written as sound) but due to deep NN, this is no longer required.
- 300h, 3000h (for research), 100,000h (commercial) of training data have been used.

(CTC (connectionist temporal classification))

For speech recognition:

Usually we will have many inputs (since a 1 hr audio clip for example is split into 100 Hz frames) but ~~only~~ ^{letter for that frame} and the outputs we get collapse to form a sentence.

1 frame x^{c1} x^{c2} x^{c3} \dots x^{c1000}

$y^{<1>} \rightarrow$ $y^{<2>} \rightarrow \dots \rightarrow y^{<1000>}$

- can be more complicated like bidirectional LSTM

We get $t t t - h - e e l l - - \frac{1}{space} - - q q q - -$

collapse to $the_\underline{e}_ll_\underline{q}_q$

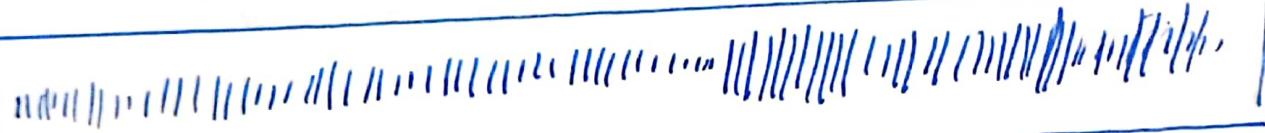
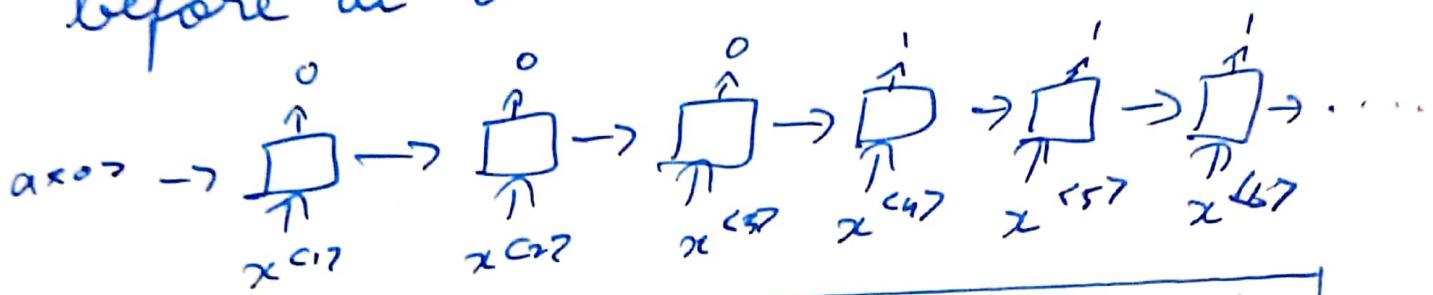
The letters that aren't divided by a - are combined (exa $t t b \rightarrow t$)

Trigger Word Detection!

Example:

Amazon Echo woken up with Alexa

→ When the trigger word is said, mark that frame with a 1, while the ones before an 0



This works reasonably well but it results in a very imbalanced training set where most examples are 0.

Hack: Instead of setting only the example with the trigger word, make subsequent 3-n examples also 1.