

## (Week 3)

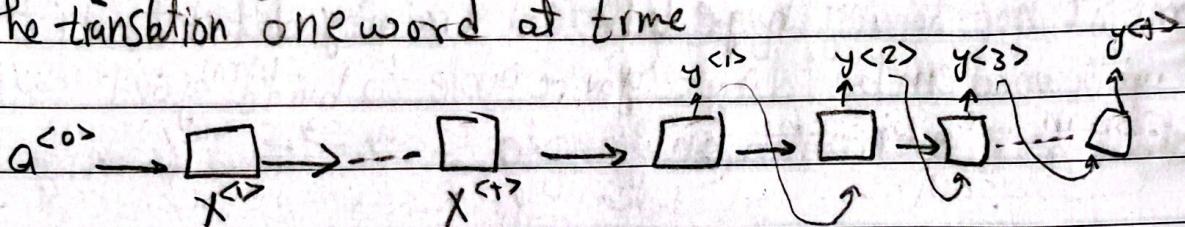
### → Basic models

- Sequence to sequence models which are useful for everything from machine translation to speech recognition.

Seq to seq takes Jane visite L'Afrique → Jane is visiting

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

- first the encoder network is built as a RNN & this could be GRU feeding the input french words one word at time & after them output a vector that represents the input sentence
- second we can build a decoder network which takes input (encoding output) & then can be trained to output the translation one word at time



### → Picking the most likely sentence.

Seq to seq. models vs language models.

- When using Model in we try to use Model for Machine translation, we don't try to sample at random from distribution instead we would find the English sentence  $y$ , that maximizes that Conditional probability  $P(y^{<1>}; \dots; y^{<n>} | X)$

English      French

→ Beam Search algorithm (used during decoding based on seq2seq)  
 • French → English @ speech recognition → Text transcript

ex of sentence:

10000 [ a  
in  
Jane  
September ]

of first thing beam search try

to pick the first words of  
the English translation that's  
going to output

First we use network fragment encoder & decoder to try

$y^{<1>} \rightarrow$  evaluate

$Q <0> \rightarrow \boxed{\quad} \rightarrow \dots \rightarrow \boxed{\quad} \rightarrow \boxed{\quad}$  The prob. of  
 $x^{<1>} \quad x^{<2>} \quad x^{<3>}$  that first word

$P(y^{<1>} | X)$

What prob. of  $y$  given  $X$

so in the first step we run the French sentence through the encoder network & then in decoder the softmax output overall 6000 possibilities, then we take those 6000 possible outputs & keep in memory top 3

$B=3$  beam width

Second, as we picked Jane & September as most likely choice of first word what beam search will do now is consider what should be in the second word

(Notes) Initialization: The model starts from initial input seq ex <start> & computes all prob. for tokens that could appear after this

② Selection: retain top K most prob. tokens

③ Expansion: for each retain token compute prob. for the next token

④ Pruning: After expanding we have g possible sentences so we

Select top 3 sequences

So beam search processes multiple sentences & choose one final seq based on cumulative prob. of all possible paths

## → Refinements to beam Search

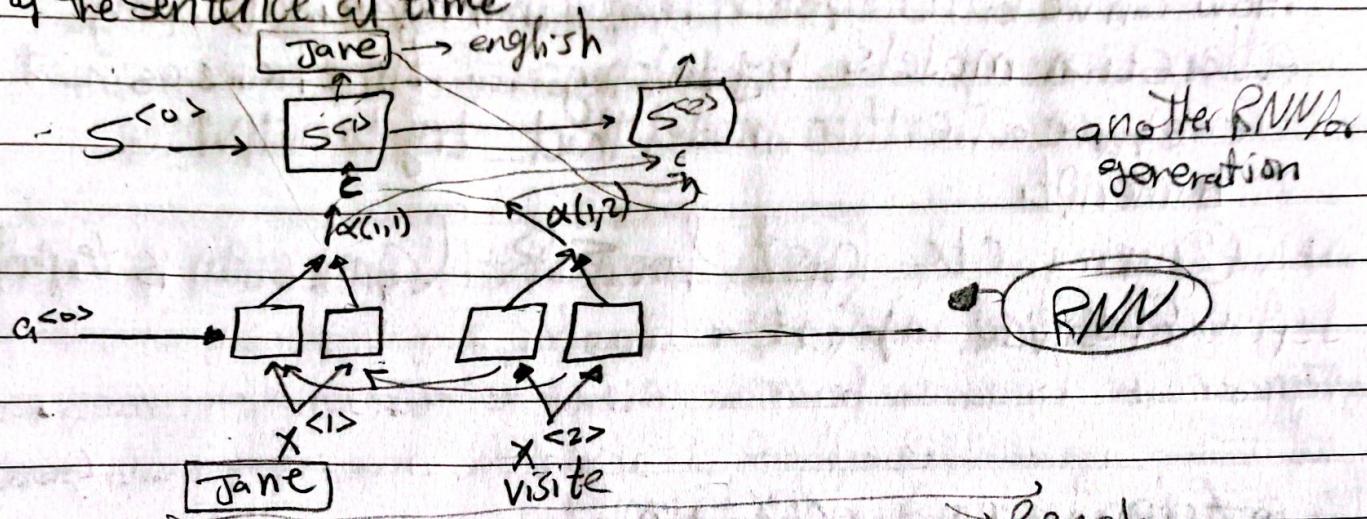
- if the beamwidth is too large then we consider a lot of possibilities & tend to get better result Because we consider a lot of different options but it will be slower & memory req. will grow also computationally slower & vice versa with low beamwidth
- Beam search is unlike search alg. like BFS or DFS  
Beam search runs faster but is not guaranteed to find exact max arg  $\text{Max } P(y|x)$

## → error analysis in Beam Search

- How error analysis interacts with beam search & how you can figure out whether it is the beam search alg. causing prob.
- might be the RNN model
  - if we have french sentence human translate  $y^*$  & alg. translates  $\hat{y}$ . So RNN(encoder+decoder) computes  $P(y^*|x)$   $P(\hat{y}|x)$ 
    - if  $P(y^*|x) > P(\hat{y}|x)$  Beamsearch chose  $\hat{y}$  But  $y^*$  has higher prob. So this means that Beam alg. didn't do its job max prob.
    - if  $P(y^*|x) < P(\hat{y}|x)$   $y^*$  is better than  $\hat{y}$  This means that RNN model is at fault

## → Attention model intuition.

- In the encoder-decoder architecture it works well for short sentences so we might achieve high Bleu score but for long sentences the performance comes down b/c it is difficult for the network to memorize super long sequences. So the attention model translates maybe like humans looking at part of the sentence at time.



Since we are trying to generate Jane in English we need to focus on the part that has Jane in French. So what the attention model would be computing is a set of attention weights ( $\alpha_{i,j}$ ) to denote when you're generating the first words how much should you pay attention to the first piece of info.

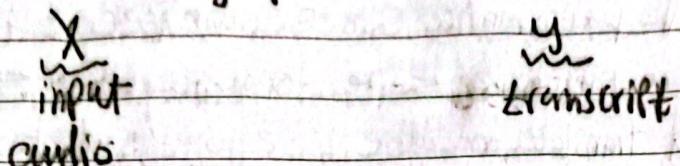
## → Attention model.

- The alpha ( $\alpha$ ) tells us how much attention @ the context ( $C$ ) would depend on the features we're getting @ activations we're getting
- $\alpha^{(t,t+1)}$  is the amount of attention  $y^{t+1}$  should pay to  $a^{t+1}$

→ Bleu Score is a metric that evaluate the quality of machine translated text by comparing it to a reference translation.

## → Speech recognition

- How seq2seq are applied to audio data
- in speech recg. problem



• How can we build speech recognition system (1) using attention models by taking diff. time frame as input  
& we have attention model that try to output the transcript.

(2) using ctc cost for SR (Connection is temporal classification) & it

## → trigger word detection.

• we take audio clip, Compute spectrogram-features, Then generates  $x_1, x_2, x_3$  audio features, Then we pass through RNN and so all that remains to be done is to define the target labels.

• In training set you can set all the target labels to 0 for every thing before the target point (siri) then all after to be 1