

Various Sequence to Sequence Architectures

Basic Models : Sequence to Sequence

French:

$x^{(1)}$ $x^{(2)}$ $x^{(3)}$ $x^{(4)}$ $x^{(5)}$

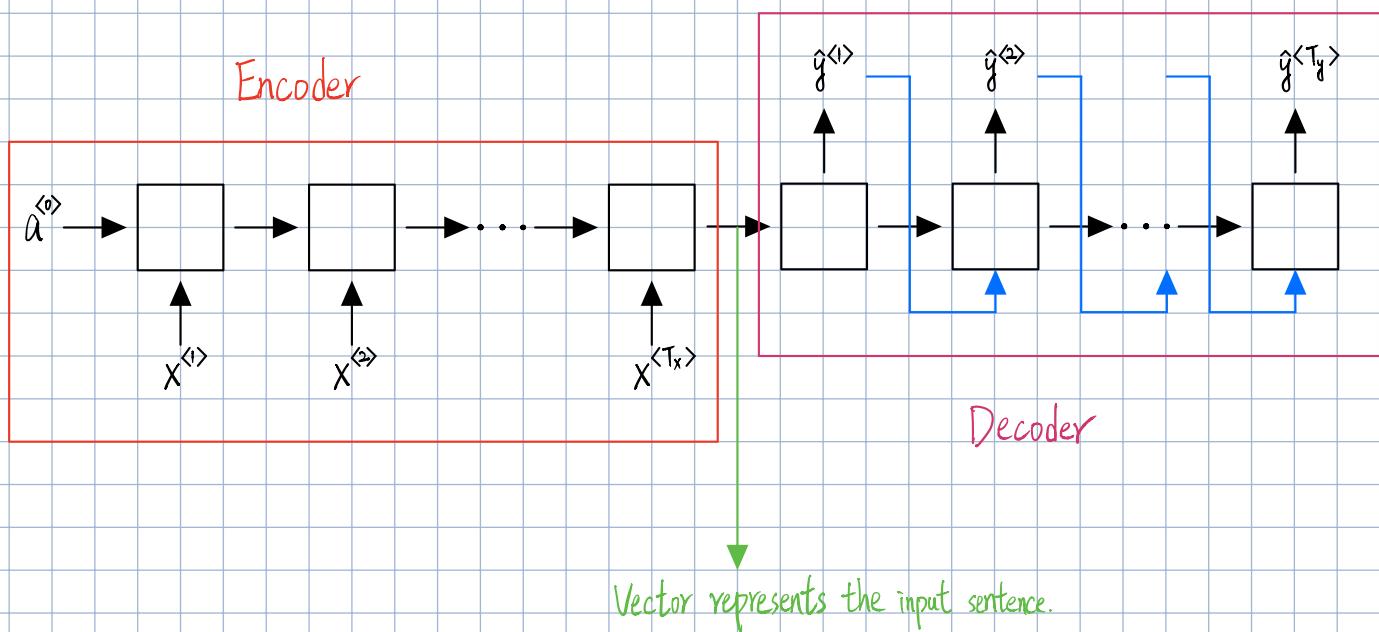
Input: Jane visite l'Afrique en septembre.

English:

$y^{(1)}$ $y^{(2)}$ $y^{(3)}$ $y^{(4)}$ $y^{(5)}$ $y^{(6)}$

Output: Jane is visiting Africa in September.

Task: French to English translation



Read: 2014 "Sequence to Sequence Learning with Neural Networks."

2014 "Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation"

Image captioning

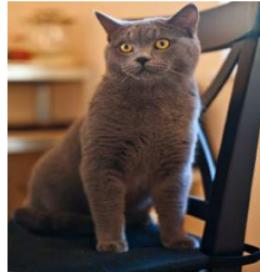
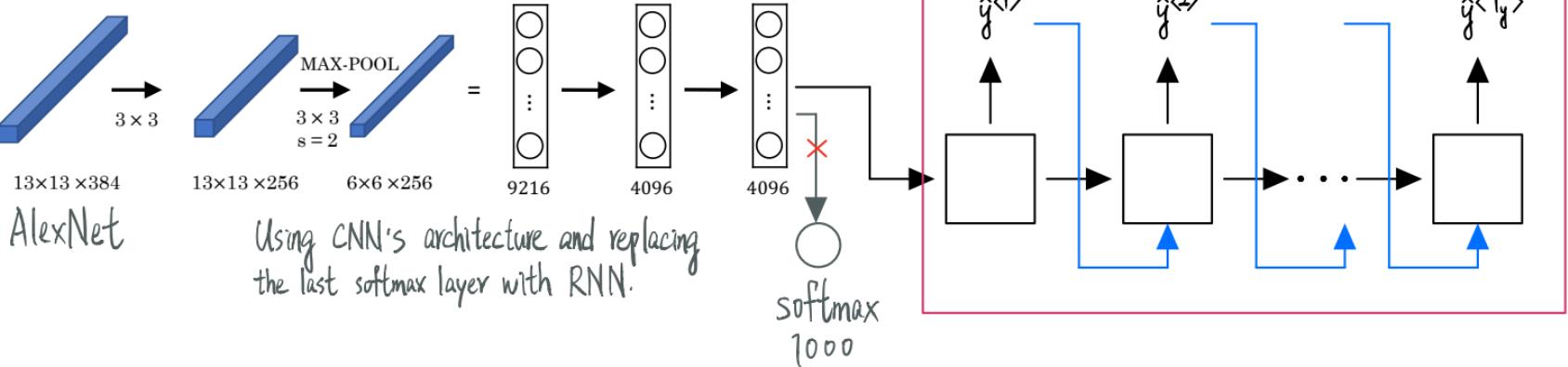
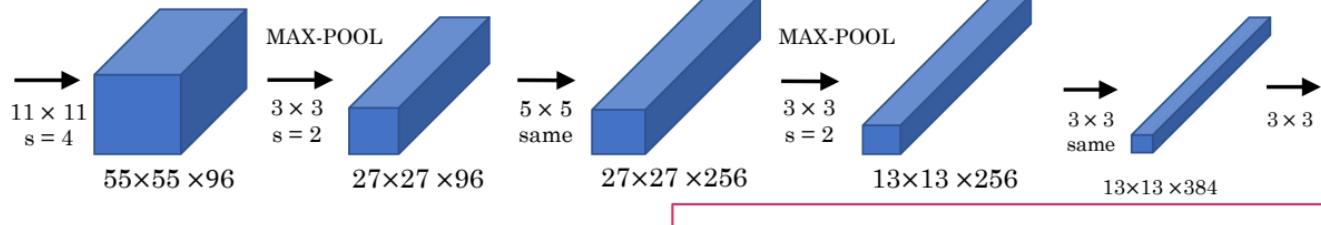


Image to Sequence model



[Mao et. al., 2014. Deep captioning with multimodal recurrent neural networks]

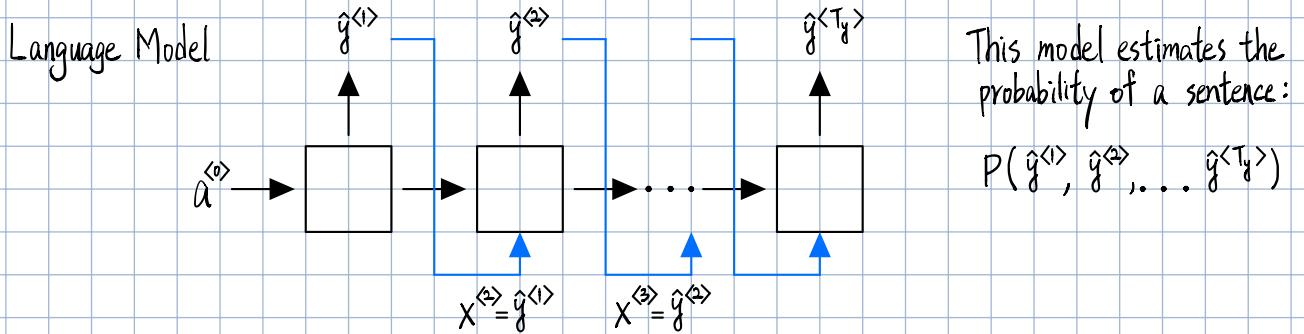
[Vinyals et. al., 2014. Show and tell: Neural image caption generator]

[Karpathy and Li, 2015. Deep visual-semantic alignments for generating image descriptions]

Andrew Ng

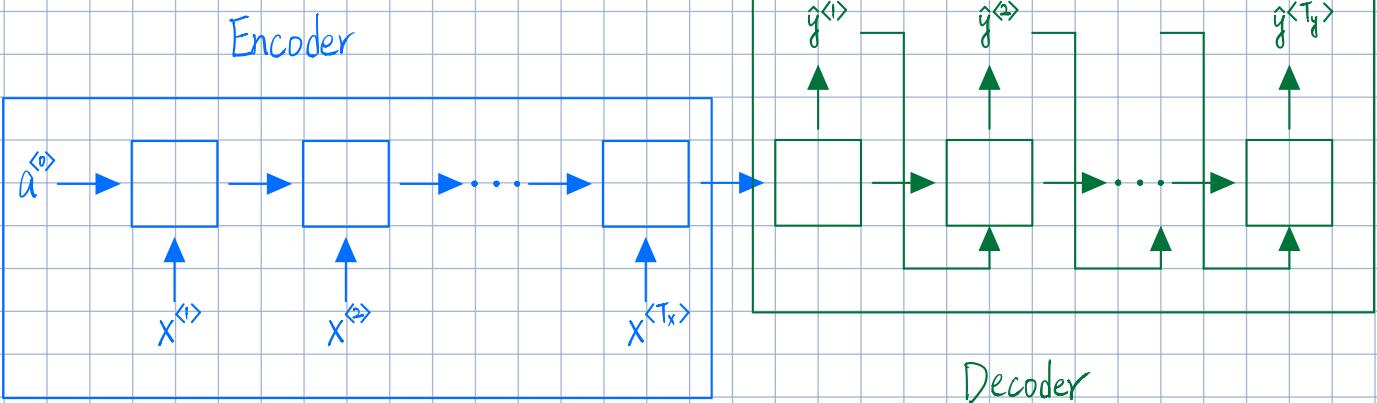
Picking the most likely sentence

Previously, we built a language model as follows:



Machine Translation as building a **conditional** language model

Machine Translation Model: **Conditional** Language Model



This model estimates the probability of the English output sentence, conditioning on an French input sentence.

$$P(\hat{g}^{(1)}, \hat{g}^{(2)}, \dots, \hat{g}^{(T_y)} | X^{(1)}, X^{(2)}, \dots, X^{(T_x)}) = P(\hat{g}^{(1)}, \hat{g}^{(2)}, \dots, \hat{g}^{(T_y)} | X)$$

English French

Now, how do we find the most likely translation?

Input: Jane visite l'Afrique en septembre.

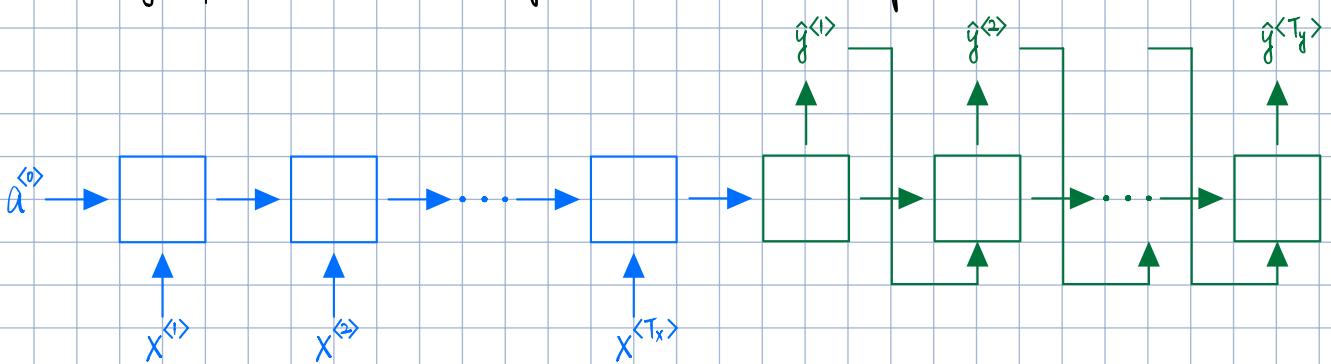
Possible Outputs:

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

$$\underset{\hat{g}^{(1)}, \dots, \hat{g}^{(T_y)}}{\operatorname{argmax}} P(\hat{g}^{(1)}, \hat{g}^{(2)}, \dots, \hat{g}^{(T_y)} | X)$$

Why not a greedy search to find the most likely translation?

A greedy search is first calculating $P(\hat{y}^{(1)}|X)$, then calculating $P(\hat{y}^{(2)}|X)$, and so on.



But what we want is the joint possibility:

$$\underset{\hat{y}^{(1)}, \hat{y}^{(2)}, \dots, \hat{y}^{(T_y)}}{\operatorname{argmax}} P(\hat{y}^{(1)}, \hat{y}^{(2)}, \dots, \hat{y}^{(T_y)} | X)$$

Here is an example to illustrate:

Possible Outputs:

- Jane is visiting Africa in September. ✓ The better translation.
- Jane is going to be visiting Africa in September. ✗ but this is chosen.
 $P(\text{Jane is going} | X) > P(\text{Jane is visiting} | X)$

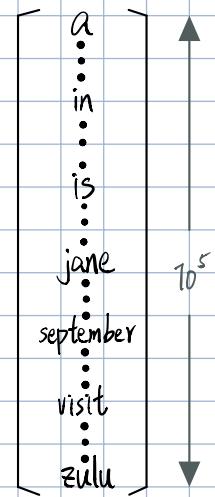
Also, the possible sentences are $(\# \text{ vocabulary})^{\# \text{ words in sentence}}$

This problem can be solved by Beam Search algorithm.

Beam Search

Beam search is a heuristic search algorithm that explores a graph by expanding the most promising node in a limited set.

Dictionary =



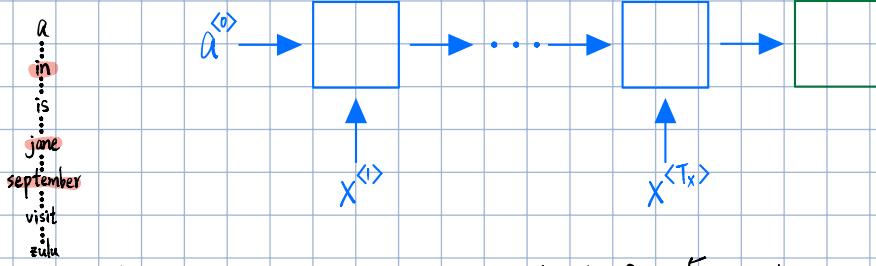
Given an input sentence and a dictionary with size = 10^5

Input: Jane visite l'Afrique en septembre.

beam width = 3. Number of the most possible words to consider in each step.

If beam width = 1, then it's exactly greedy search?

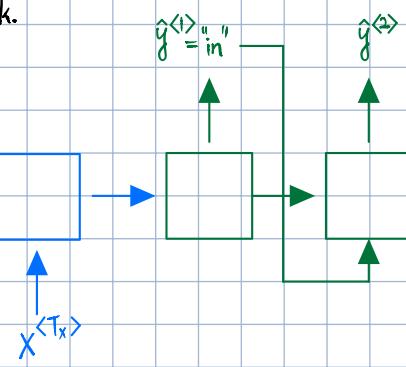
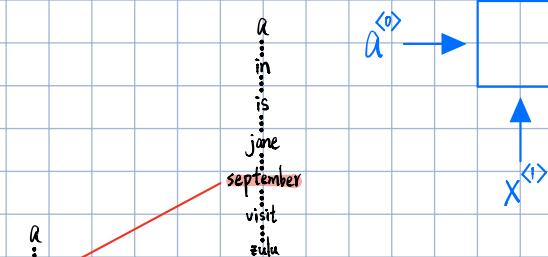
Step 1: Calculate $P(\hat{y}^{(1)} | X)$.



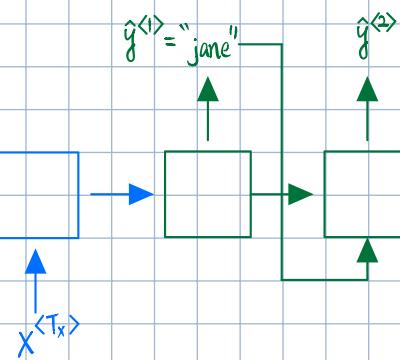
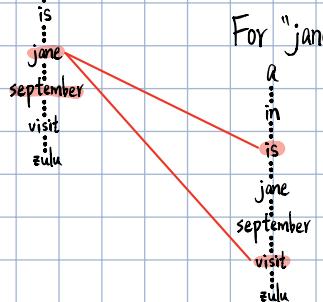
It selects "in", "jane", and "september" out of 10^5 outputs.

Step 2: Calculate $P(\hat{y}^{(1)}, \hat{y}^{(2)} | X, \hat{y}^{(1)})$ from the neural network.

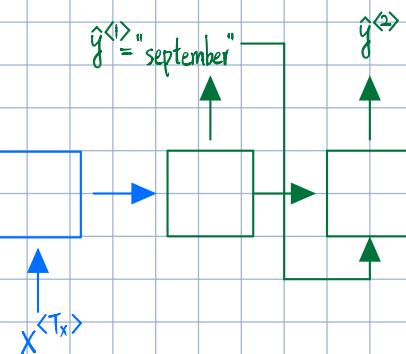
For "in"



For "jane"



For "september"

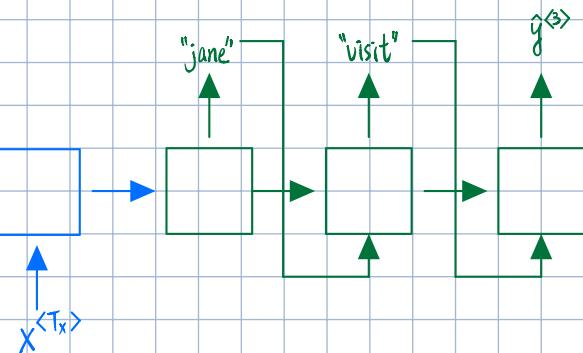
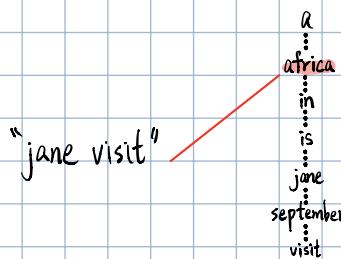
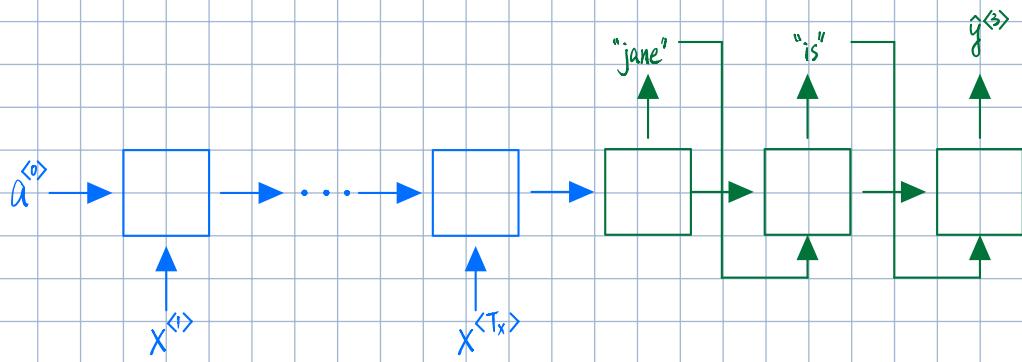
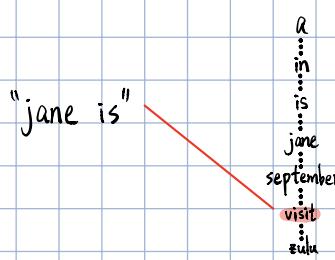
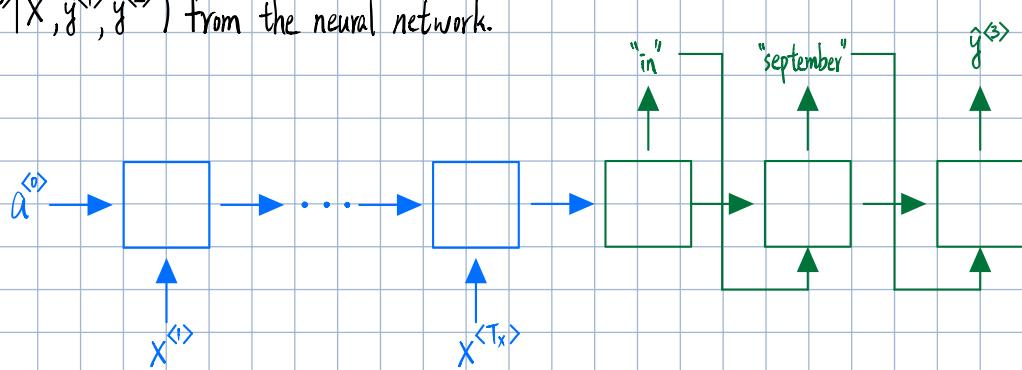
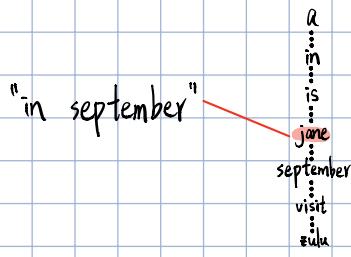


What we want is $P(\hat{y}^{(1)}, \hat{y}^{(2)} | X) = P(\hat{y}^{(1)} | X) P(\hat{y}^{(2)} | X, \hat{y}^{(1)})$

There are $3 \cdot 10^4$ choices in step 2. We choose top 3 $P(\hat{y}^{(1)}, \hat{y}^{(2)} | X)$ out of them.

from step 1.

Step 3: Calculate $P(\hat{g}^{(3)}|X, \hat{g}^{(1)}, \hat{g}^{(2)})$ from the neural network.



Again, we pick 3 out of 3×10^4 possibilities.

$$\begin{aligned}
 P(\hat{g}^{(1)}, \hat{g}^{(2)}, \hat{g}^{(3)} | X) &= P(\hat{g}^{(1)}, \hat{g}^{(2)} | X) P(\hat{g}^{(3)} | X, \hat{g}^{(1)}, \hat{g}^{(2)}) \\
 &= P(\hat{g}^{(1)} | X) P(\hat{g}^{(2)} | X, \hat{g}^{(1)}) P(\hat{g}^{(3)} | X, \hat{g}^{(1)}, \hat{g}^{(2)})
 \end{aligned}$$

from step 2

Refinements to Beam Search

$$\begin{aligned}
 P(\hat{y}^{(1)}, \hat{y}^{(2)}, \dots, \hat{y}^{(T_y)} | X) \\
 &= P(\hat{y}^{(1)} | X) P(\hat{y}^{(2)} | X, \hat{y}^{(1)}) P(\hat{y}^{(3)} | X, \hat{y}^{(1)}, \hat{y}^{(2)}) \cdots P(\hat{y}^{(T_y)} | X, \hat{y}^{(1)}, \dots, \hat{y}^{(T_y-1)}) \\
 &= \prod_{t=1}^{T_y} P(\hat{y}^{(t)} | X, \hat{y}^{(1)}, \dots, \hat{y}^{(t-1)})
 \end{aligned}$$

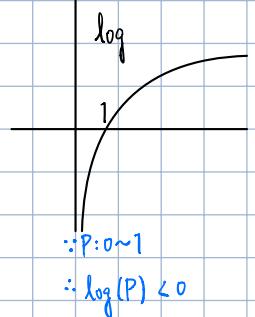
Length Normalization

$$\operatorname{argmax}_{\hat{y}} \prod_{t=1}^{T_y} P(\hat{y}^{(t)} | X, \hat{y}^{(1)}, \dots, \hat{y}^{(t-1)}) \quad \text{This equation is numerical unstable due to probability being float.}$$

$$\rightarrow \operatorname{argmax}_{\hat{y}} \sum_{t=1}^{T_y} \log P(\hat{y}^{(t)} | X, \hat{y}^{(1)}, \dots, \hat{y}^{(t-1)}) \neq \operatorname{argmax}_{\hat{y}} \sum_{t=1}^{T_y} P(\hat{y}^{(t)} | X, \hat{y}^{(1)}, \dots, \hat{y}^{(t-1)})$$

This model tends to make the output sentence shorter. More output, more negative numbers are added together, and smaller the model.

Why? $\because \log(P) < 0$



Thus, we can normalize this function by:

$$\operatorname{argmax}_{\hat{y}} \frac{1}{T_y} \sum_{t=1}^{T_y} \log P(\hat{y}^{(t)} | X, \hat{y}^{(1)}, \dots, \hat{y}^{(t-1)})$$

$\alpha = 0$, no normalization

$\alpha = 1$, normalize by the length of the output sentence

$\alpha = 0.7$, heuristic value

Discussion:

Large beam width: better result, more memory, running slower.

Small beam width: worse result, less memory, running faster.

Beam width =	10	100	1000+
Product			Publish papers

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 $\operatorname{argmax}_{\hat{y}} P(\hat{y} | X)$.

Error Analysis in Beam Search

At this point, we can use RNN models and beam search to run machine translation.

But if the translation result is bad, how do I know whom to blame. RNN models or beam search or me.

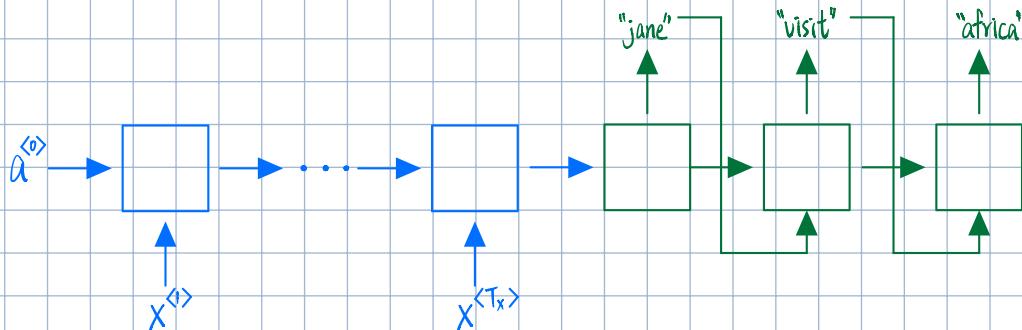
Eg:

French: Jane visite l'Afrique en septembre.

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

Algorithm: Jane visited Africa in September. (\hat{y}) A quite bad translation.

I can use the RNN model to compute $P(y^*|x)$ and $P(\hat{y}|x)$, and see which one is larger.



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

Conclusion: Beam search is to be blamed.

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

Conclusion: RNN model is to be blamed.

Now, acquire more samples and go through this process and the sinner will be found.

Bleu Score

Read "BLEU: a Method for Automatic Evaluation of Machine Translation"

One of the challenges of machine translation is that, given a French sentence, there could be multiple English translations that are equally good translations of that French sentence. How do I evaluate a machine translation system if there are multiple equally good answers.

BiLingual Evaluation Understudy (BLEU)

Given French: Le chat est sur le tapis.

Reference 1: The cat is on the mat. 2 "the"

$2 > 1 \therefore \text{select 2}$

Reference 2: There is a cat on the mat. 1 "the"

MT (Machine Translation) output: the the the the the the Terrible Translation

(Wrong) Precision: $\frac{7}{7} \text{ Every "the" is in either one of Reference 1 and 2.} = 1$

The score is high but it's a terrible translation. Don't use it.

Each word has only the credit up to the maximum

Modified Precision: $\frac{\text{number of times it appears in the reference sentences}}{7 \text{ "the" in output}} = \frac{2}{7}$

BLEU score on bigram: 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.

Bigram "the cat" "cat the" "cat on" "on the" "the mat"

Count in MT: 2 1 1 1 1

Count_{clip} (max times of occurrences in references) : 1 0 1 1 1

$$\text{Thus, the Precision} = \frac{1+1+1+1}{2+1+1+1+1}$$

$$= \frac{4}{6}$$

BLEU score on unigram: 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. (\hat{y})

Define P_1 as the Precision for unigram.

P_2 as the Precision for bigram.

\vdots

P_n as the Precision for n-gram.

$$P_1 = \frac{\sum_{\text{unigram} \in \hat{y}} \text{Count}_{\text{clip}}(\text{unigram})}{\sum_{\text{unigram} \in \hat{y}} \text{Count}(\text{unigram})}$$

$$P_1 = \frac{\sum_{n\text{-gram} \in \hat{y}} \text{Count}_{\text{clip}}(n\text{-gram})}{\sum_{n\text{-gram} \in \hat{y}} \text{Count}(n\text{-gram})}$$

Given P_n as BLEU score on n-gram only. $n=1, 2, 3, 4$

Combine BLEU score: $BP \cdot \exp\left(\sum_{n=1}^4 P_n\right)$ BP: Brevity Penalty

$$BP = \begin{cases} 1 & , \text{ if } \text{MT_output_length} > \text{reference_output_length} \\ \exp\left(1 - \frac{\text{reference_output_length}}{\text{MT_output_length}}\right) & , \text{ otherwise} \end{cases}$$

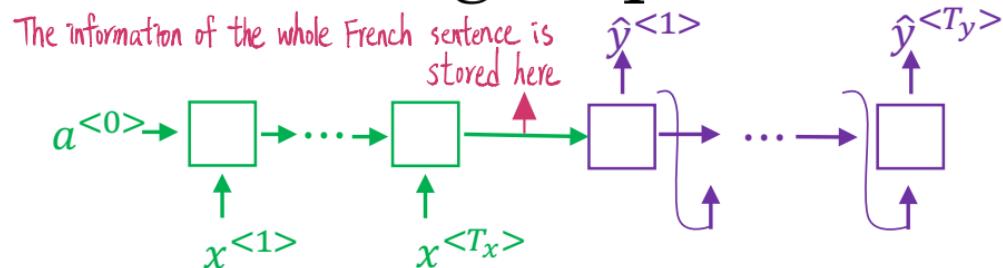


deeplearning.ai

Sequence to sequence models

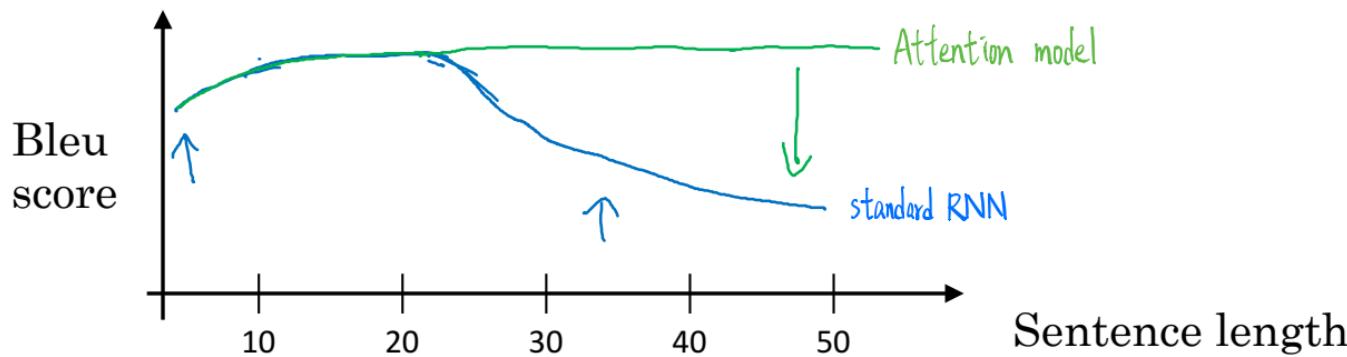
Attention model intuition

The problem of long sequences



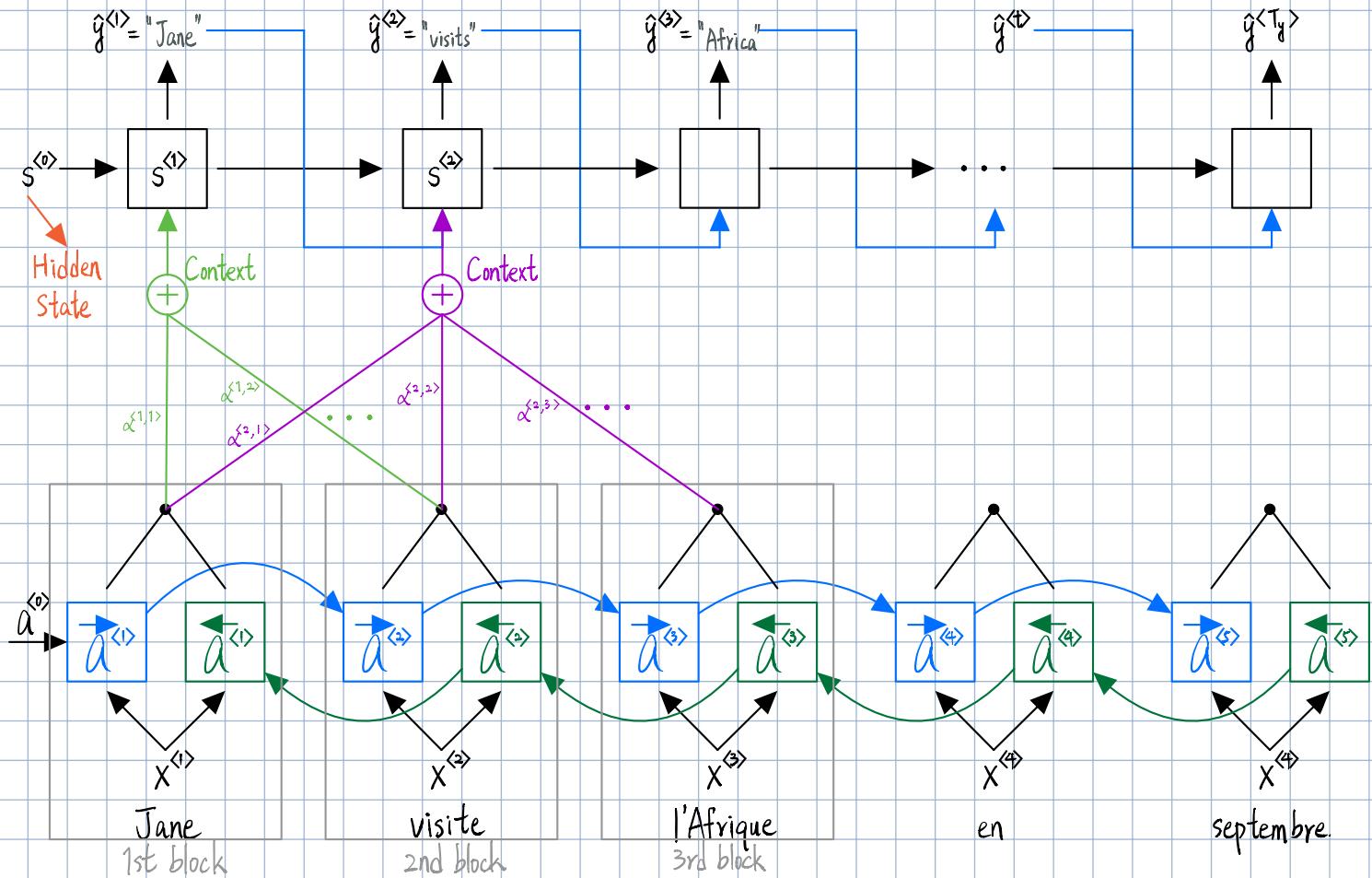
Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.

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.



Attention Model Intuition

Read 2014 "Neural Machine Translation by Jointly Learning to Align and Translate"

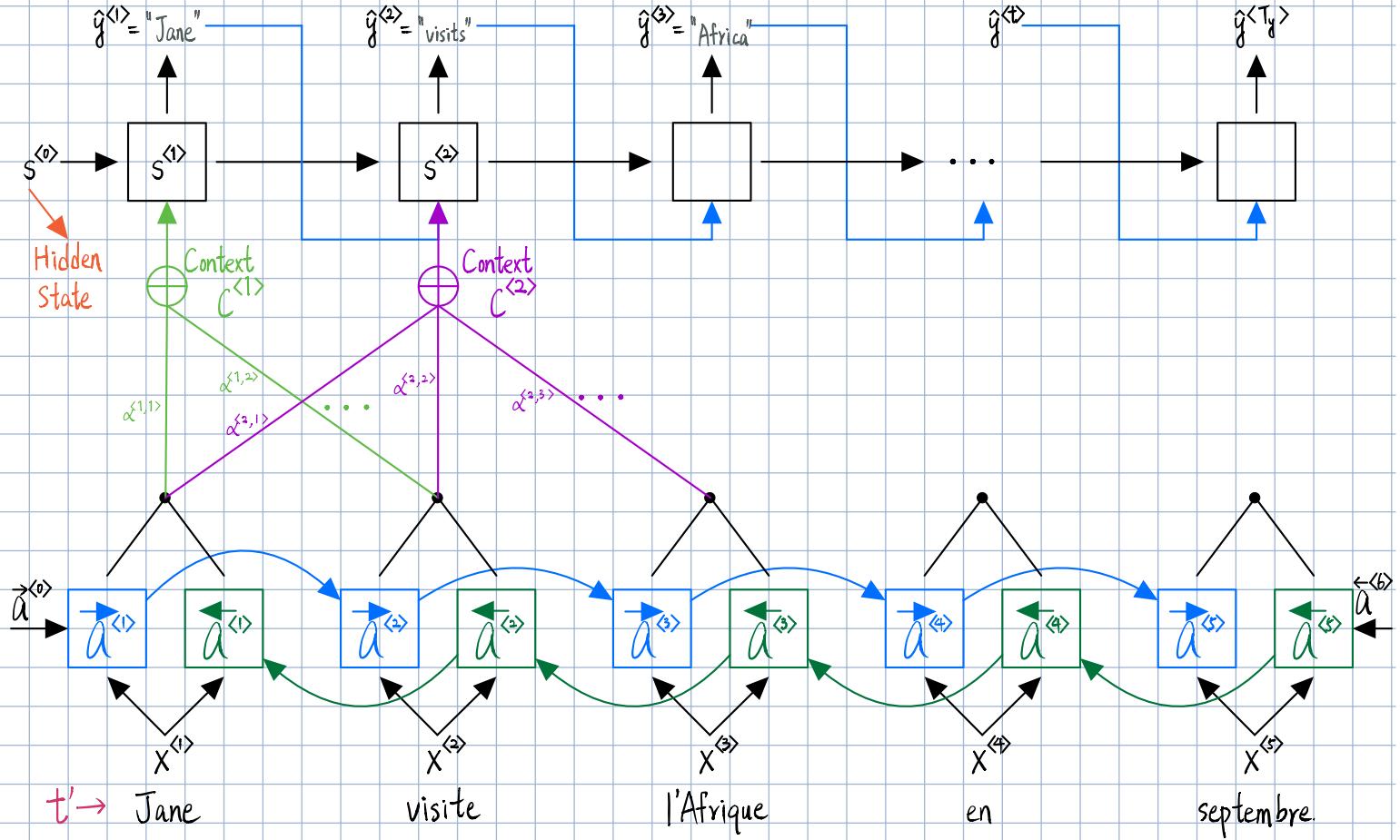


Attention Weight $\alpha^{(1,1)}$: When the model is generating the first output, how much it should be paying attention to the first block.

$\alpha^{(1,2)}$: When the model is generating the first output, how much it should be paying attention to the second block.

$\alpha^{(2,1)}$: When the model is generating the second output, how much it should be paying attention to the first block.

Attention Model Implementation



Redefine the time step of the input sentence as t' to distinguish from the time step t of the output sentence.

$$\text{Thus, } \hat{a}^{(t)} = (\hat{a}^{(t)}, \hat{a}^{(t)})$$

The attention weights for the first output word $\sum_{t'=1}^{T_x} \alpha^{(1, t')} = 1$

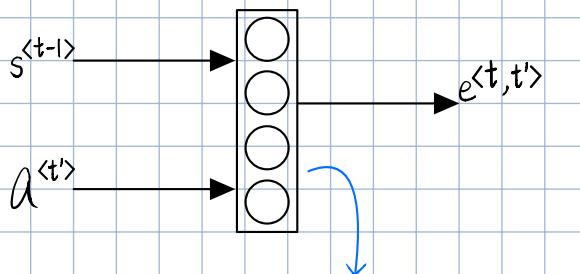
The context C for the first output word $C^{(1)} = \sum_{t'=1}^{T_x} \alpha^{(1, t')} \hat{a}^{(t')}$. For the second output word $C^{(2)} = \sum_{t'=1}^{T_x} \alpha^{(2, t')} \hat{a}^{(t')}$

$\alpha^{(t, t')}$: the amount of "Attention" that $\hat{y}^{(t)}$ should pay to $\hat{a}^{(t')}$. In other words, when the model is generating the t -th output word, how much attention it should be paying to the t' -th input word.

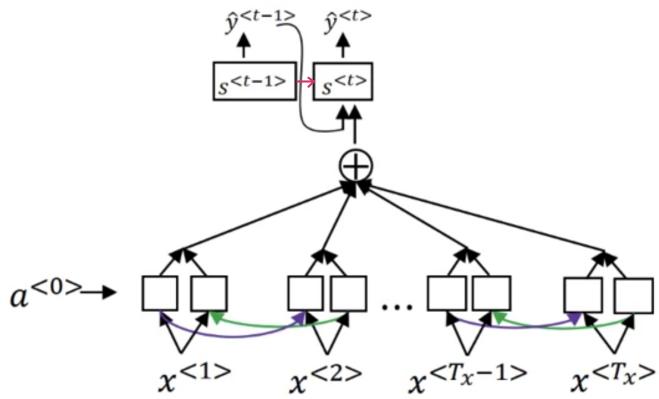
Computing Attention $\alpha^{t,t'}$

$$\alpha^{t,t'} = \frac{\exp(e^{t,t'})}{\sum_{t'=1}^{T_x} \exp(e^{t,t'})} \quad \text{This is why } \sum_{t'=1}^{T_x} \alpha^{t,t'} = 1$$

Train a small neural network to generate $e^{t,t'}$.



This neural network tells how much attention $\hat{y}^{t'}$ should pay to $a^{t'}$.



The intuitive explanation is that $e^{t,t'}$ is depending on the state of the previous time step s^{t-1} and the associated $a^{t'}$, which makes sense.

$s^{t'}$ depends on $\alpha^{t,t'}$ which in turn depends on $e^{t,t'}$.

One downside of attention model is that it takes quadratic cost to run this algorithm. If there are T_x words in the input and T_y words in the output, then the total number of attention parameters $\alpha^{t,t'}$ is $T_x \cdot T_y$.

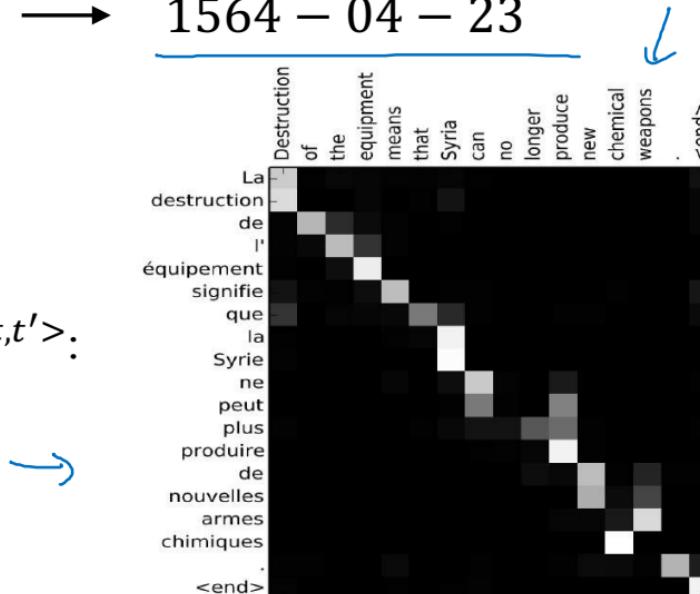
Read 2015 "Show, Attention and Tell: Neural Image Caption Generation with Visual Attention".

Attention examples

July 20th 1969 → 1969 – 07 – 20

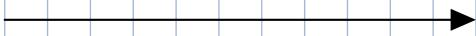
23 April, 1564 → 1564 – 04 – 23

Visualization of $\alpha^{<t,t'>}$:

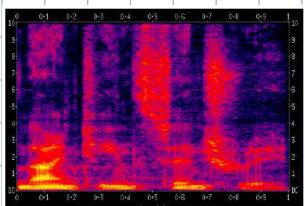
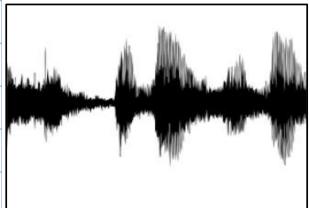


Speech Recognition

Audio clip X



Transcript Y



"the quick brown fox"

Once upon a time, speech recognition system
is used to be built using **Phonemes**.

Phonemes are hand-crafted basic units of sound.

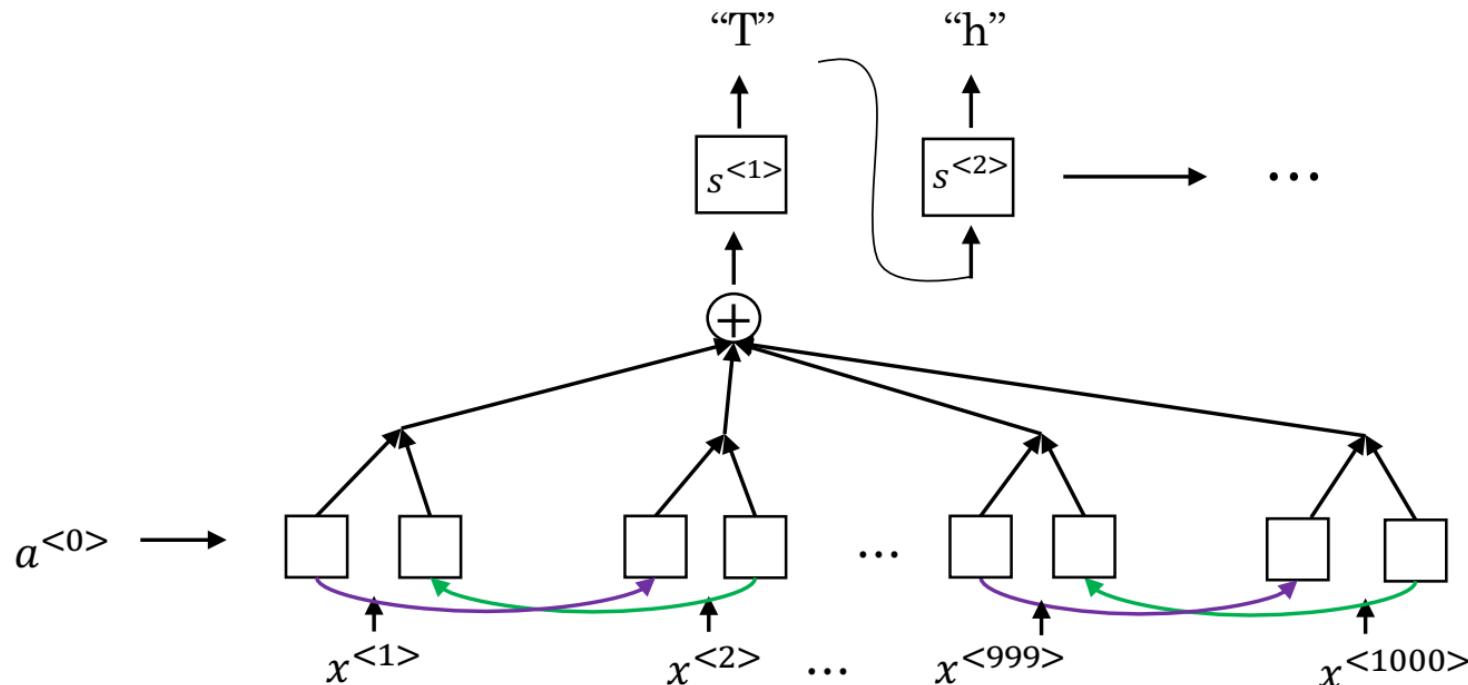
People used to decompose audio into phonemes
to conduct speech recognition.
But with end-to-end deep learning, people
find that phonemes representations are no longer
necessary.

We just need More data.

In academia, $300h \sim 3000h$ data.

In commercial, $10^4h \sim 10^5h$ data.

Attention model for speech recognition

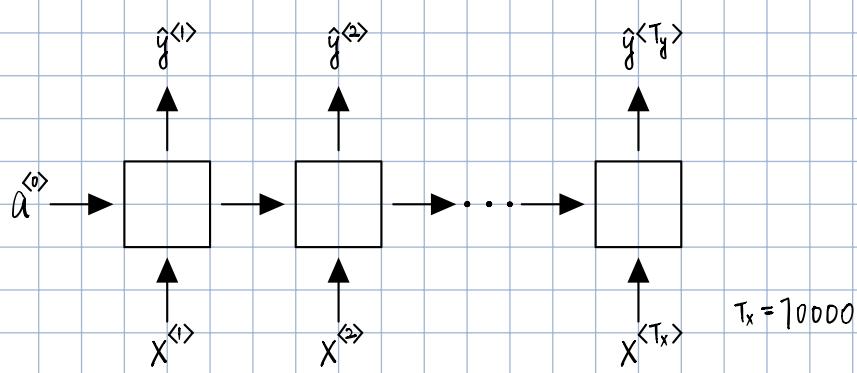


CTC cost for speech recognition

CTC: Connectionist Temporal Classification

Read 2006 "Connectionist Temporal Classification : Labeling Unsegmented Sequence Data with Recurrent Neural Networks"

audio = "the quick brown fox"



In practice, bidirectional RNN with LSTM or GRU is used.

In speech recognition, the number of input time steps is much bigger than the number of output time steps.
Eg, 70 seconds audio, sampled with 1k Hz, has 70k input time step, while the output words are only few words.

Basic rule: collapse repeated characters not separated by "blank"

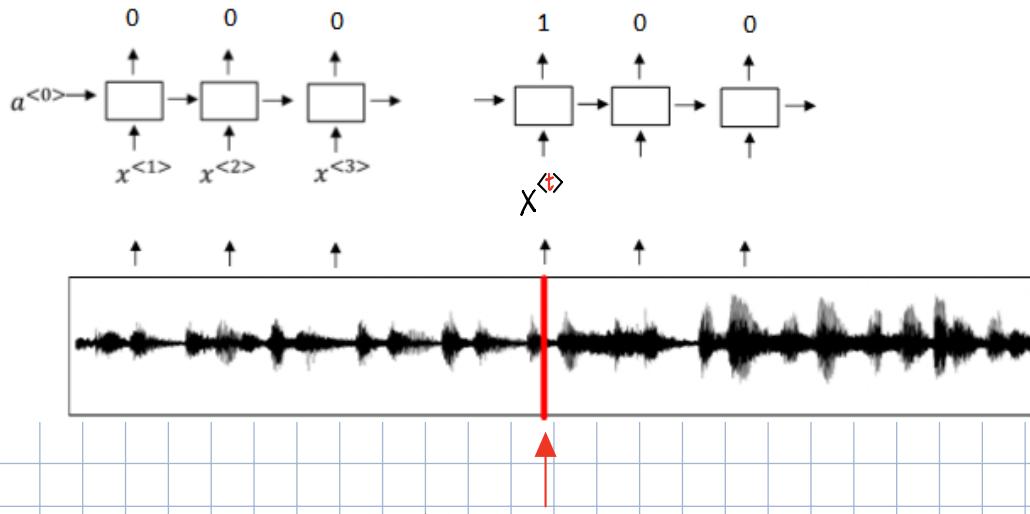
\hat{y} : ttt_heeee space ffff --> the l q
blank character

Trigger Word Detection

Eg: Alexa

Hey Siri

Okay Google



When the target label of $x^{(t)}$ is set to 1, it means a trigger word has been said.

$x^{(t)}$ represents the feature of the audio at time t.

However, the above scheme creates a **imbalanced training set**.

So, we can set the output target label a few 1s for a fixed period of time before resetting to 0.

