

Sequence models → Week 03 (Attention mechanism)

 aakashgoel12.medium.com/sequence-models-week-03-attention-mechanism-111345e5a8c0

February 21, 2021



Machine translation as building a conditional language model

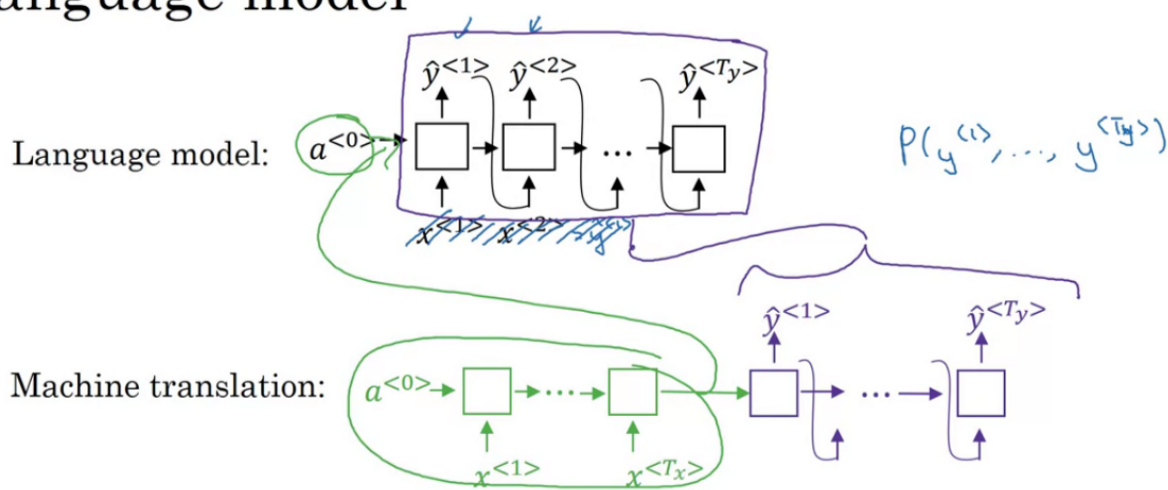


Image 01

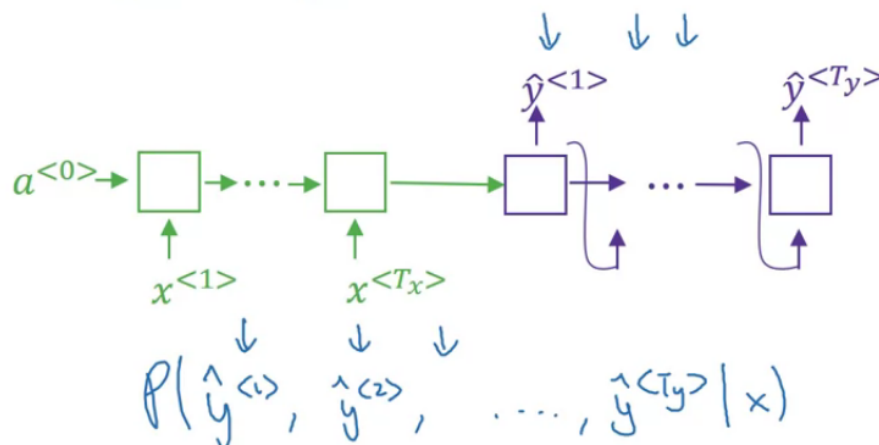
In Language Model, we find probability of sentence.

Decoder in Machine Translation system is same as Language Model and $a_{<0>}$ in language model is similar to Encoder in Machine Translation.

In M/C Translation, we use beam search instead of greedy search.

Why not a greedy search?

$P(\hat{y})$

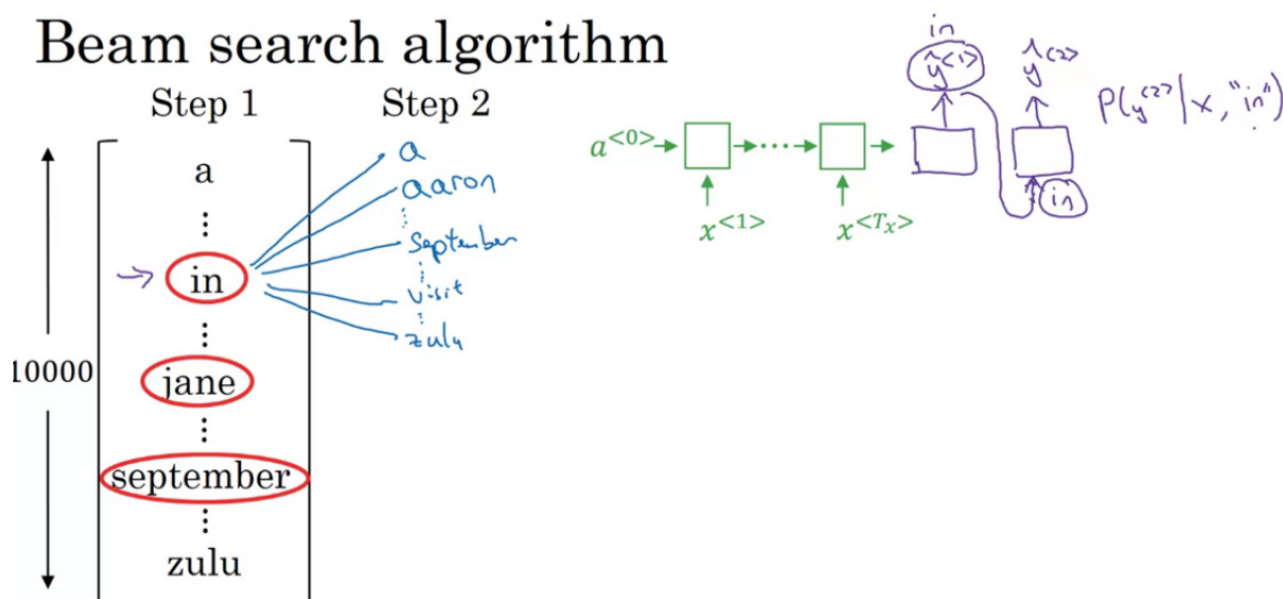


- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.

$P(\text{Jane is going}/X) > P(\text{Jane is visiting}/X)$ but sentence 1 is more optimal

Beam Search

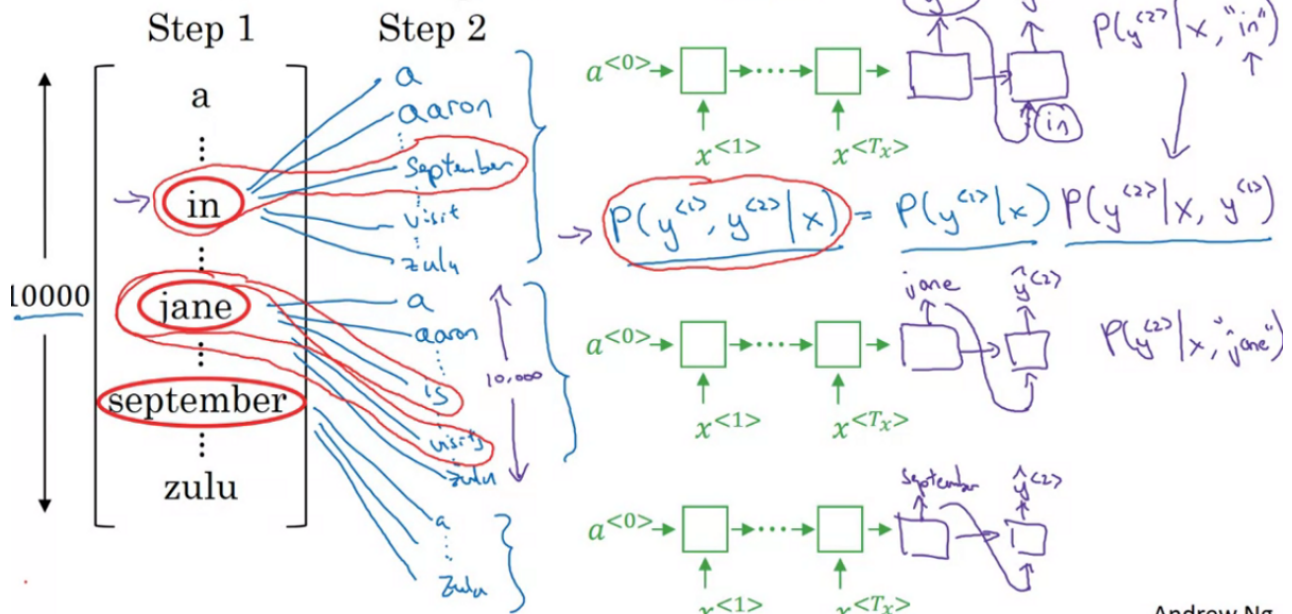
Beam search algorithm



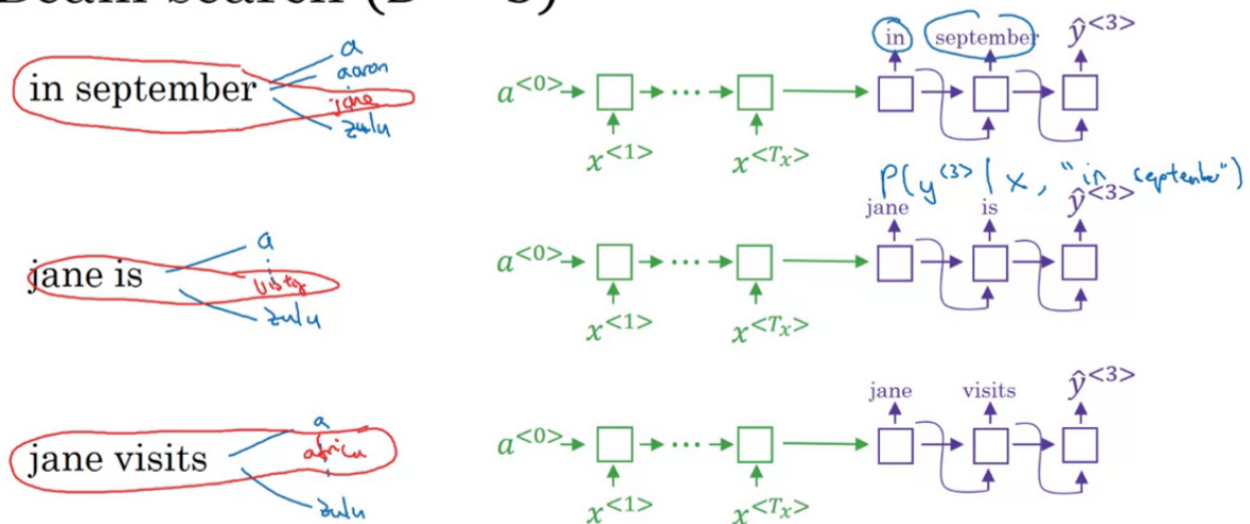
Beam Width considered is 3 i.e. Top 3 words will be considered as candidate..

Say, Word1= "in", need to find $P(Y_2/X, 'in')$ i.e. Prob. of Y_2 given X and "in".

Beam search algorithm ($B=3$)



Beam search ($B = 3$)



$P(y^{<1>}, y^{<2>} | x)$ jane visits africa in september. <EOS>

Log is strictly **monotonically increasing** function i.e. maximizing $P(Y/X)$ is same as maximizing $\text{Log}(p(Y/X))$..

Length normalization

$$p(y^{<1>} \dots y^{<T_y>} | x) = \frac{p(y^{<1>} | x) p(y^{<2>} | x, y^{<1>}) \dots}{p(y^{<T_y>} | x, y^{<1>}, \dots, y^{<T_y-1>})}$$

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

log

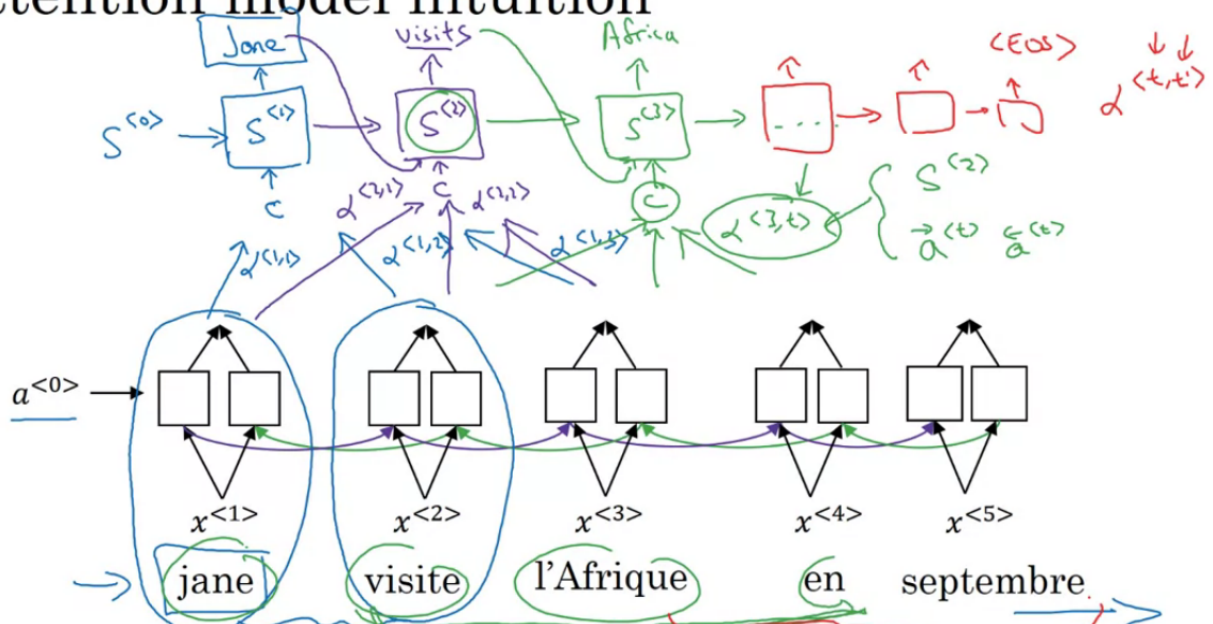
$$\arg \max_y \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

log $P(y|x) \leftarrow P(y|x)$

Above $P(Y_t/X, Y_1, \dots, Y_{t-1})$, Unnaturally tends/prefer short translations as multiplying no less than 1 will give short tiny number ..

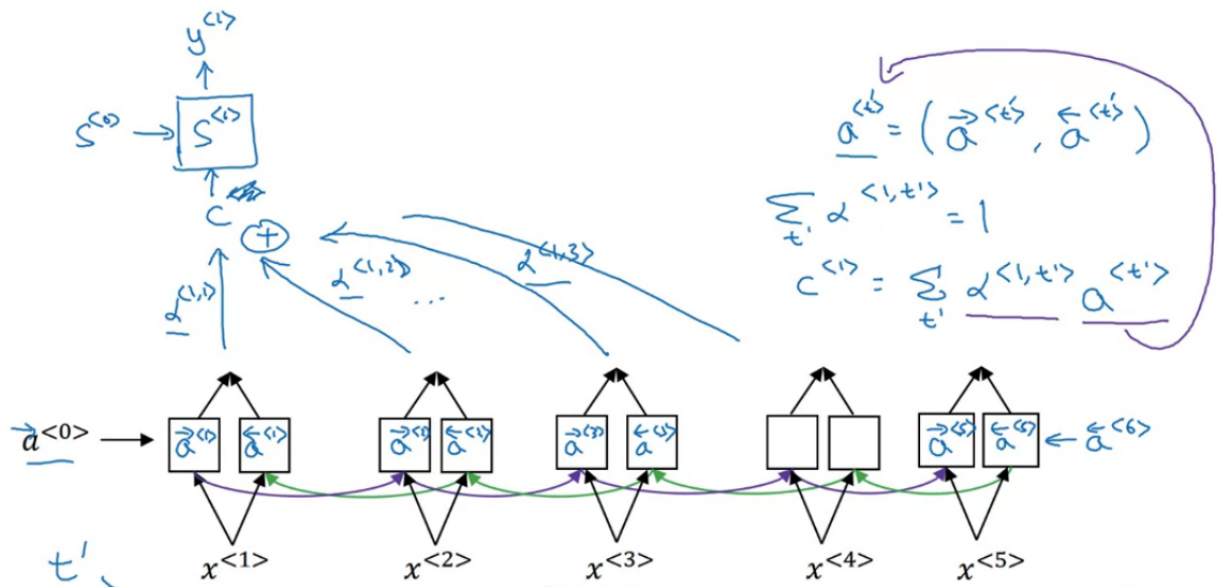
Attention (Alpha (t,x)) → How much weight to be used for generating t word using time-stamp x

Attention model intuition



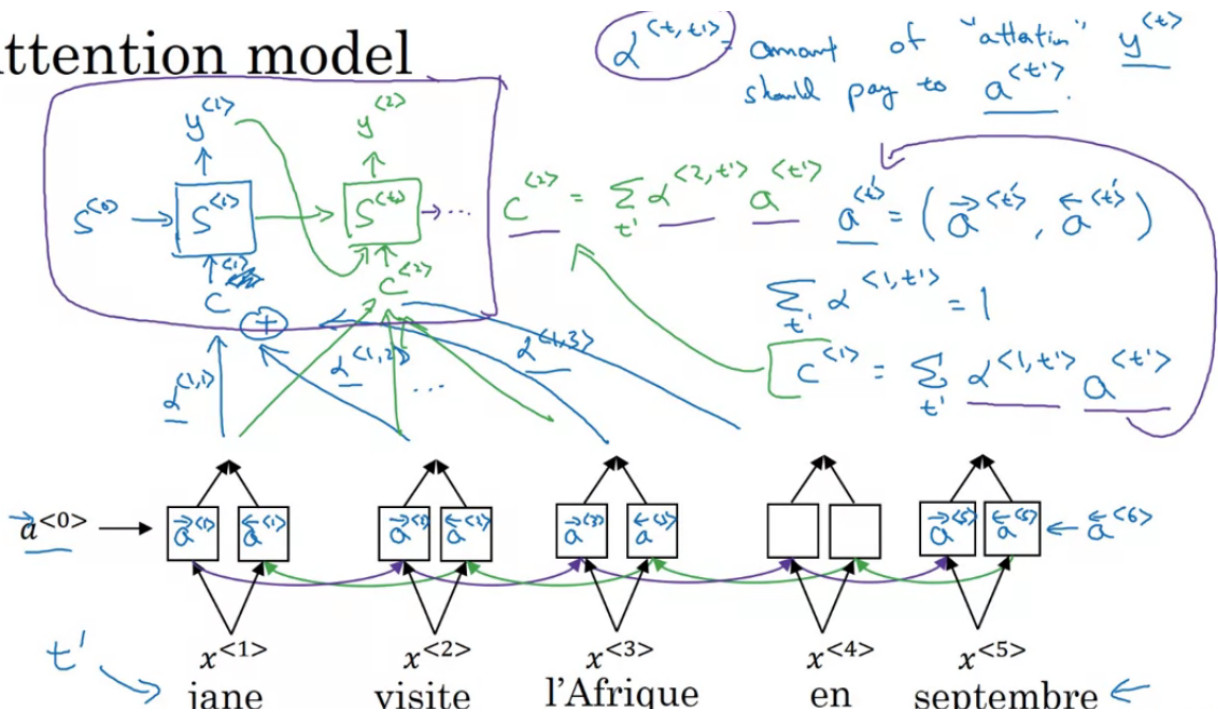
Part — 01 (Attention) → a is combination of backward and forward propagation .. For 1st word, will have 5 timestamp alphas i.e. attention weights and its summation will be 1. C (Context Vectors) is summation of different timestamps.

Attention model



PART — 02 (Attention) → A

Attention model

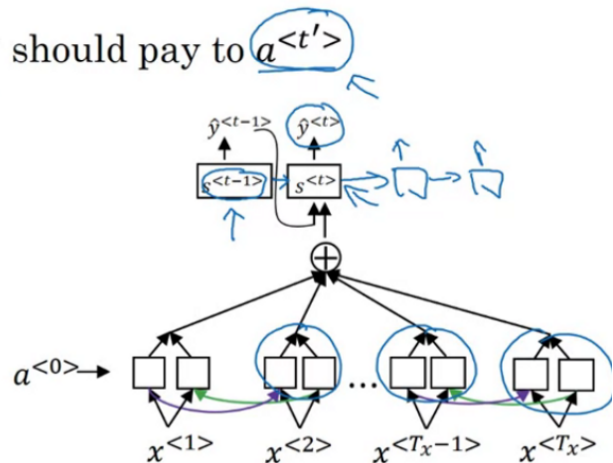


Now, how to calculate Alpha (t,t') i.e. Amount of attention Y(t) should pay to a(t').

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

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

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



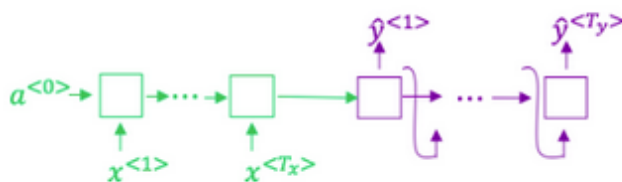
What is NEXT ? → [https://workera.ai/?](https://workera.ai/?utm_source=coursera_sequence_models&utm_medium=Coursera&utm_campaign=coursera_sequence_models)

[utm_source=coursera_sequence_models&utm_medium=Coursera&utm_campaign=cou](https://workera.ai/?utm_source=coursera_sequence_models&utm_medium=Coursera&utm_campaign=coursera_sequence_models)
[rsera_sequence_models](https://workera.ai/?utm_source=coursera_sequence_models&utm_medium=Coursera&utm_campaign=coursera_sequence_models)

https://drive.google.com/file/d/1099XMofOen_QfoNL3qqLUOXy-CdMyJQ4/view

QUIZ

1. Consider using this encoder-decoder model for machine translation.



This model is a "conditional language model" in the sense that the encoder portion (shown in green) is modeling the probability of the input sentence x .

☐ True

☒ False

✓ Correct

2. In beam search, if you increase the beam width B , which of the following would you expect to be true? Check all that apply.

☒ Beam search will run more slowly.

✓ Correct

☒ Beam search will use up more memory.

✓ Correct

☒ Beam search will generally find better solutions (i.e. do a better job maximizing $P(y | x)$)

✓ Correct

☐ Beam search will converge after fewer steps.

3. In machine translation, if we carry out beam search without using sentence normalization, the algorithm will tend to output overly short translations.

☒ True

☐ False

✓ Correct

4. Suppose you are building a speech recognition system, which uses an RNN model to map from audio clip x to a text transcript y . Your algorithm uses beam search to try to find the value of y that maximizes $P(y | x)$.

On a dev set example, given an input audio clip, your algorithm outputs the transcript \hat{y} = "I'm building an A Eye system in Silly con Valley.", whereas a human gives a much superior transcript y^* = "I'm building an AI system in Silicon Valley."

According to your model,

$$P(\hat{y} | x) = 1.09 * 10^{-7}$$

$$P(y^* | x) = 7.21 * 10^{-8}$$

Would you expect increasing the beam width B to help correct this example?

- ☒ No, because $P(y^* | x) \leq P(\hat{y} | x)$ indicates the error should be attributed to the RNN rather than to the search algorithm.
- ☐ No, because $P(y^* | x) \leq P(\hat{y} | x)$ indicates the error should be attributed to the search algorithm rather than to the RNN.
- ☐ Yes, because $P(y^* | x) \leq P(\hat{y} | x)$ indicates the error should be attributed to the RNN rather than to the search algorithm.
- ☐ Yes, because $P(y^* | x) \leq P(\hat{y} | x)$ indicates the error should be attributed to the search algorithm rather than to the RNN.

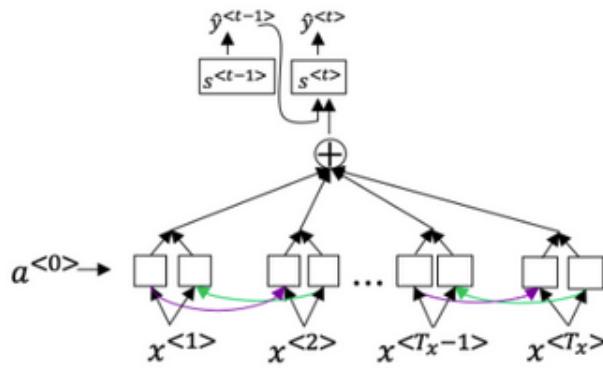
✓ Correct

5. Continuing the example from Q4, suppose you work on your algorithm for a few more weeks, and now find that for the vast majority of examples on which your algorithm makes a mistake, $P(y^* | x) > P(\hat{y} | x)$. This suggests you should focus your attention on improving the search algorithm.

- ☒ True.
- ☐ False.

✓ Correct

6. Consider the attention model for machine translation.



Further, here is the formula for $\alpha^{<t,t'>}$.

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

Which of the following statements about $\alpha^{<t,t'>}$ are true? Check all that apply.

- ☒ We expect $\alpha^{<t,t'>}$ to be generally larger for values of $a^{<t'>}$ that are highly relevant to the value the network should output for $y^{<t>}$. (Note the indices in the superscripts.)

✓ Correct

- ☐ We expect $\alpha^{<t,t'>}$ to be generally larger for values of $a^{<t'>}$ that are highly relevant to the value the network should output for $y^{<t'>}$. (Note the indices in the superscripts.)
- ☐ $\sum_t \alpha^{<t,t'>} = 1$ (Note the summation is over t .)
- ☒ $\sum_{t'} \alpha^{<t,t'>} = 1$ (Note the summation is over t' .)

✓ Correct

7. The network learns where to “pay attention” by learning the values $e^{<t,t'>}$, which are computed using a small neural network:

We can't replace $s^{<t-1>}$ with $s^{<t>}$ as an input to this neural network. This is because $s^{<t>}$ depends on $\alpha^{<t,t'>}$ which in turn depends on $e^{<t,t'>}$; so at the time we need to evaluate this network, we haven't computed $s^{<t>}$ yet.

- ☒ True
☐ False

✓ Correct

8. Compared to the encoder-decoder model shown in Question 1 of this quiz (which does not use an attention mechanism), we expect the attention model to have the greatest advantage when:

- ☒ The input sequence length T_x is large.
☐ The input sequence length T_x is small.

✓ Correct

9. Under the CTC model, identical repeated characters not separated by the “blank” character () are collapsed. Under the CTC model, what does the following string collapse to?

_c_oo_o_kk__b_ooooo_oo_kkk

- ☐ cokbok
☒ cookbook
☐ cook book
☐ coookkboooooookkk

✓ Correct

10. In trigger word detection, $x^{<t>}$ is:

- ☒ Features of the audio (such as spectrogram features) at time t .
☐ The t -th input word, represented as either a one-hot vector or a word embedding.
☐ Whether the trigger word is being said at time t .
☐ Whether someone has just finished saying the trigger word at time t .

✓ Correct

Assignment → Jupyter Notebook

1 - Translating human readable dates into machine readable dates

- The model you will build here could be used to translate from one language to another, such as translating from English to Hindi.
- However, language translation requires massive datasets and usually takes days of training on GPUs.
- To give you a place to experiment with these models without using massive datasets, we will perform a simpler "date translation" task.
- The network will input a date written in a variety of possible formats (e.g. "the 29th of August 1958", "03/30/1968", "24 JUNE 1987")
- The network will translate them into standardized, machine readable dates (e.g. "1958-08-29", "1968-03-30", "1987-06-24").
- We will have the network learn to output dates in the common machine-readable format YYYY-MM-DD.

1.1 - Dataset

We will train the model on a dataset of 10,000 human readable dates and their equivalent, standardized, machine readable dates. Let's run the following cells to load the dataset and print some examples.

```
dataset[:10]
```

```
[('9 may 1998', '1998-05-09'),  
( '10.11.19', '2019-11-10'),  
( '9/10/70', '1970-09-10'),  
( 'saturday april 28 1990', '1990-04-28'),  
( 'thursday january 26 1995', '1995-01-26'),  
( 'monday march 7 1983', '1983-03-07'),  
( 'sunday may 22 1988', '1988-05-22'),  
( '08 jul 2008', '2008-07-08'),  
( '8 sep 1999', '1999-09-08'),  
( 'thursday january 1 1981', '1981-01-01')]
```

You've loaded:

- `dataset`: a list of tuples of (human readable date, machine readable date).
- `human_vocab`: a python dictionary mapping all characters used in the human readable dates to an integer-valued index.
- `machine_vocab`: a python dictionary mapping all characters used in machine readable dates to an integer-valued index.
 - **Note:** These indices are not necessarily consistent with `human_vocab`.
- `inv_machine_vocab`: the inverse dictionary of `machine_vocab`, mapping from indices back to characters.

Let's preprocess the data and map the raw text data into the index values.

- We will set $T_x=30$
 - We assume T_x is the maximum length of the human readable date.
 - If we get a longer input, we would have to truncate it.
- We will set $T_y=10$
 - "YYYY-MM-DD" is 10 characters long.

You now have:

- `X`: a processed version of the human readable dates in the training set.
 - Each character in `X` is replaced by an index (integer) mapped to the character using `human_vocab`.
 - Each date is padded to ensure a length of T_x using a special character (< pad >).
 - `X.shape = (m, T_x)` where m is the number of training examples in a batch.
- `Y`: a processed version of the machine readable dates in the training set.
 - Each character is replaced by the index (integer) it is mapped to in `machine_vocab`.
 - `Y.shape = (m, T_y)`.
- `Xoh`: one-hot version of `X`
 - Each index in `X` is converted to the one-hot representation (if the index is 2, the one-hot version has the index position 2 set to 1, and the remaining positions are 0).
 - `Xoh.shape = (m, T_x , len(human_vocab))`
- `Yoh`: one-hot version of `Y`
 - Each index in `Y` is converted to the one-hot representation.
 - `Yoh.shape = (m, T_y , len(machine_vocab))`.
 - `len(machine_vocab) = 11` since there are 10 numeric digits (0 to 9) and the - symbol.
- Let's also look at some examples of preprocessed training examples.
- Feel free to play with `index` in the cell below to navigate the dataset and see how source/target dates are preprocessed.

Source date: 9 may 1998
Target date: 1998-05-09

[illegible]

```

36 36 36 36 36]
Target after preprocessing (indices): [ 2 10 10  9  0  1  6  0  1 10]

```

```
Source after preprocessing (one-hot): [[ 0.  0.  0. ...,  0.  0.  0.]
 [ 1.  0.  0. ...,  0.  0.  0.]
 [ 0.  0.  0. ...,  0.  0.  0.]
```

```

...,
[ 0.  0.  0. ..., 0.  0.  1.]
[ 0.  0.  0. ..., 0.  0.  1.]
[ 0.  0.  0. ..., 0.  0.  1.]]

```

```
Target after preprocessing (one-hot): [[ 0.  0.  1.  0.  0.  0.  0.  0.  0.  0.]
 [ 0.  0.  0.  0.  0.  0.  0.  0.  0.  1.]
 [ 0.  0.  0.  0.  0.  0.  0.  0.  0.  1.]
 [ 0.  0.  0.  0.  0.  0.  0.  0.  1.  0.]
 [ 1.  0.  0.  0.  0.  0.  0.  0.  0.  0.]
 [ 0.  1.  0.  0.  0.  0.  0.  0.  0.  0.]
 [ 0.  0.  0.  0.  0.  0.  1.  0.  0.  0.]
 [ 1.  0.  0.  0.  0.  0.  0.  0.  0.  0.]
 [ 0.  1.  0.  0.  0.  0.  0.  0.  0.  0.]
 [ 0.  0.  0.  0.  0.  0.  0.  0.  0.  1.]]
```

2.1 - Attention mechanism

In this part, you will implement the attention mechanism presented in the lecture videos.

- Here is a figure to remind you how the model works.
 - The diagram on the left shows the attention model.
 - The diagram on the right shows what one "attention" step does to calculate the attention variables $\alpha^{(t')}$.
 - The attention variables $\alpha^{(t')}$ are used to compute the context variable $context^{(t)}$ for each timestep in the output ($t = 1, \dots, T_Y$).

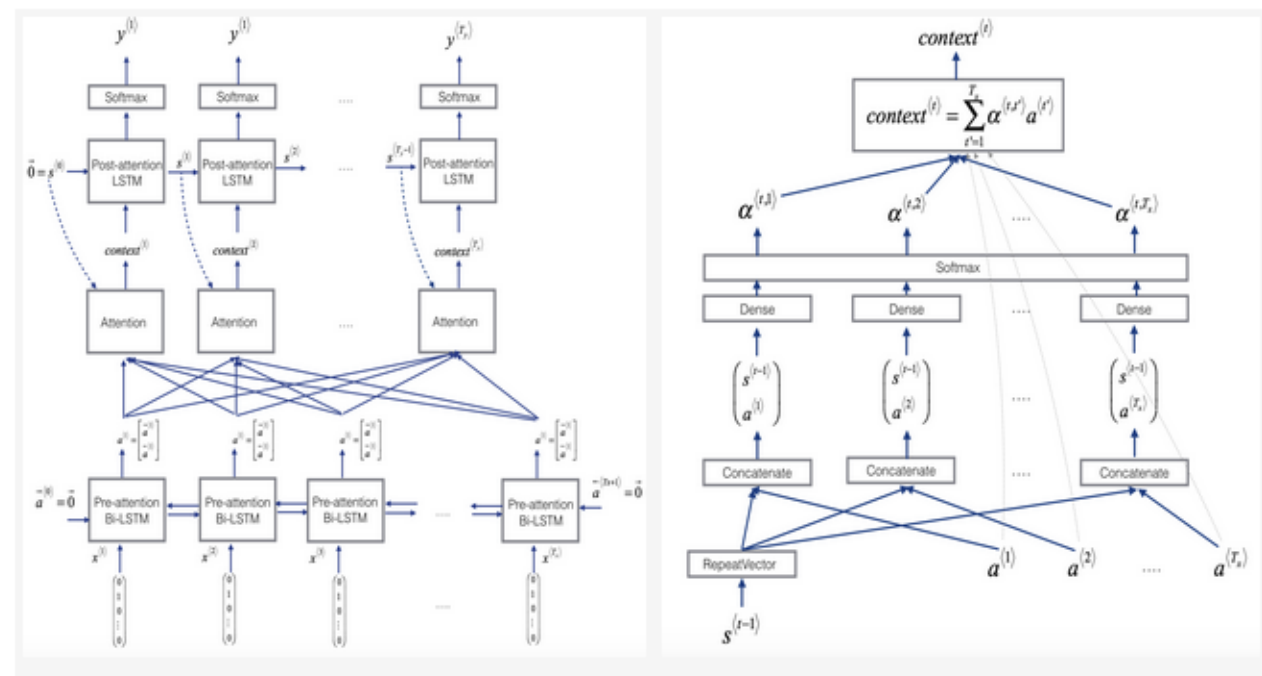


Figure 1: Neural machine translation with attention