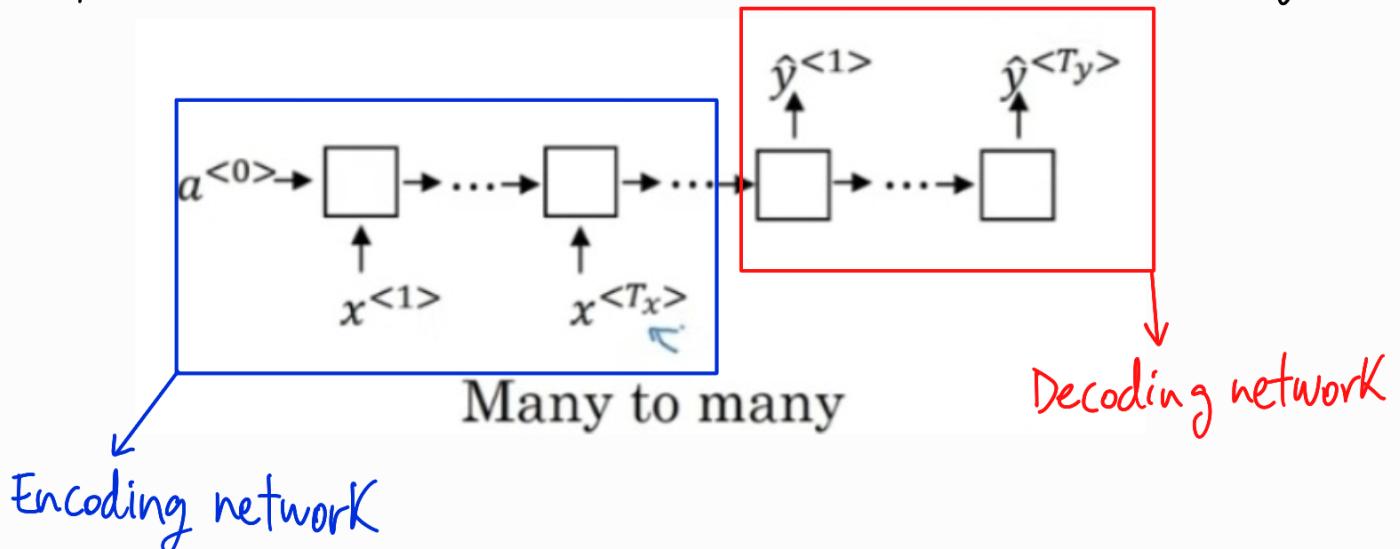


The sequence to sequence model, is a many to many RNN.



So, we are computing  $P(\underbrace{y^{<1>} \dots y^{<T_y>}}_{\text{Output sentence}} | \underbrace{x^{<1>} \dots x^{<T_x>}}_{\text{Input sentence}})$ ,

and try to get the  $\underset{y^{<1>} \dots y^{<T_y>}}{\operatorname{argmax}} P(y^{<1>} \dots y^{<T_y>} | x)$ .

We don't do a random search, but a beam search.

• Beam Search: We specify a parameter  $B$ , which represents the number of most probable words to select from  $\hat{y}^{<i>}$  and pass to  $\hat{y}^{<i+1>}$ .

For instance, after the encoding layer computations have finished, we get the first  $\hat{y}$  from softmax. Select the  $B$  max, and pass them to the next.

Consequently, at each step we instantiate  $B$  copies of the network.

It can be optimized:

- Length normalization: avoid really short translations

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

$\Downarrow$   
prone to numerical underflow!!

Take logs:

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

As the logarithmic function is a strictly monotonically increasing function, we know that maximizing  $\log(P(y|x))$  = maximizing  $P(y|x)$

We can then normalize (divide by  $T_y^\alpha$ ):

$$\frac{1}{T_y^\alpha} \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>}) \quad \underline{\alpha = 0.7}$$

The  $\alpha$  is for softening.

- Choosing a great  $B$ :

↳ Large  $B \Rightarrow$  better results, slower

↳ Small  $B \Rightarrow$  worse results, faster

Note  $\Rightarrow$  Beam search does NOT always guarantee obtaining the best result.

## Error analysis on beam search

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

$$P(y^*|x)$$

$$P(\hat{y}|x)$$

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

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

$$\arg \max_y P(y|x)$$

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

Conclusion: Beam search is at fault.

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

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

Conclusion: RNN model is at fault.

(Being  $y^*$  the output of a human).

A sentence may be translated in different equally correct ways. How to choose?

BLEU score: how similar it is to human translation.

For n-grams;  $\rightarrow$  max times that the n-gram appears

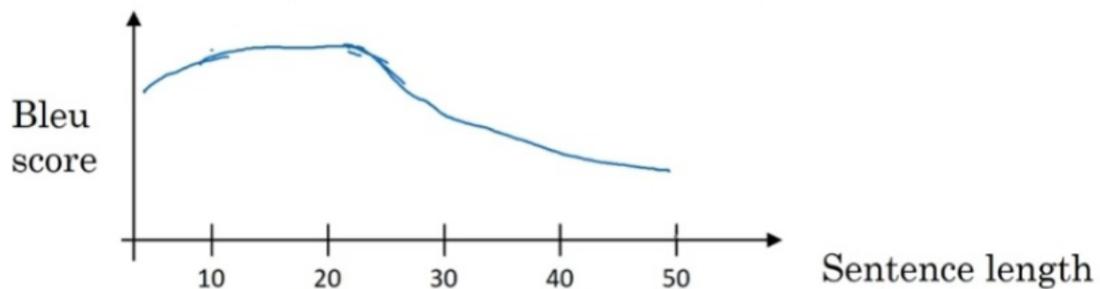
$$P_n = \frac{\sum_{n\text{-grams} \in \hat{y}} \text{Count}_{\text{clip}}(\text{n-gram})}{\sum_{n\text{-grams} \in \hat{y}} \text{Count}(\text{n-gram})} \quad \text{for all the human-translated sample}$$

$\sum_{n\text{-grams} \in \hat{y}} \text{Count}(\text{n-gram})$   $\hookrightarrow$  how many times a n-gram appears

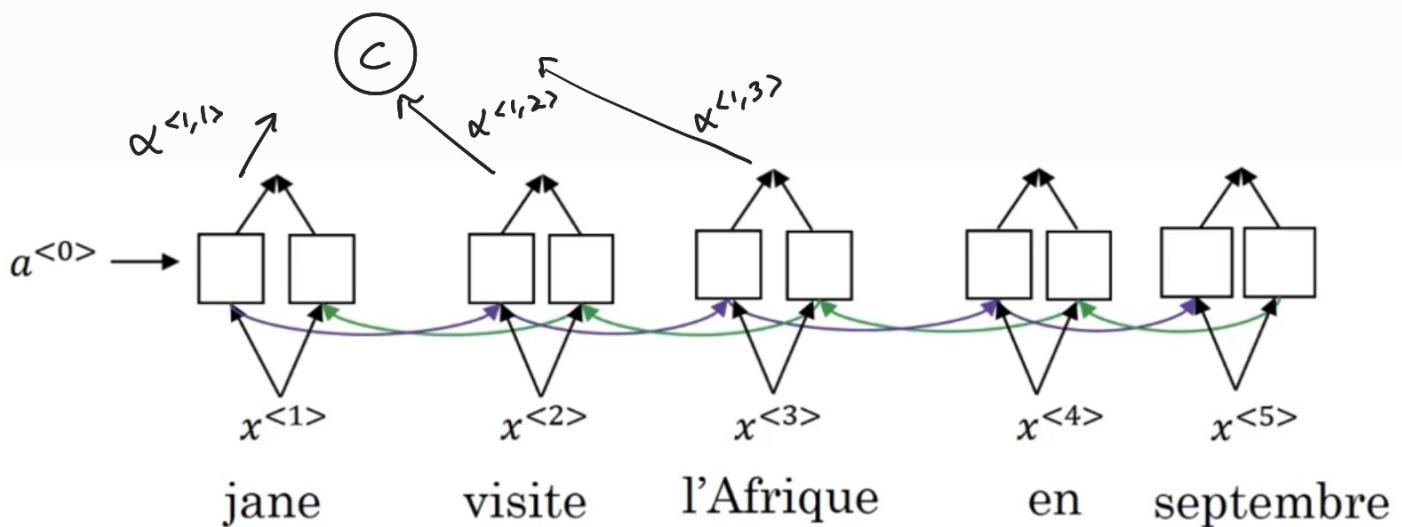
Combining Bleus :  $B.P.e^{(\text{mean of bleus})}$

$$BP = \begin{cases} 1 & \text{if } MT\_output\_length > reference\_output\_length \\ \exp(1 - ref.output.length / Mach.trans.output.length) & \text{otherwise} \end{cases}$$

The encoder-decoder model decreases performance as the length of the sentence to translate is bigger.



> Attention models: with a BRNN, we compute some  $\alpha$  weights for each node, and then compute a context var ( $c$ ).



This way, we are just paying attention at a part of the sentence.

Each alpha expresses how much of the context will depend on the activations of each time step ( $\bar{\alpha}^{<t>}, \bar{\alpha}^{<t>}$ ). therefore, it is the attention  $y^{<t>}$  should pay to  $a^{<t>}$ .

$$C^{<i>} = \sum_{t'} \alpha^{<i,t'} a^{<t>} ; \text{ being } a^{<t>} = (\bar{\alpha}^{<t>}, \bar{\alpha}^{<t>})$$

But, how do we compute each  $\alpha$ ?

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

