



## Refinements of Beam Search

→ Length Normalization

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

Take log to avoid numerical underflow

$$= P(y^{(1)}|x) P(y^{(2)}|y^{(1)}) \dots P(y^{(T_y)}|x, y^{(1)}, y^{(2)}, \dots, y^{(T_y-1)})$$

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

normalize sentence length to allow longer output sentences

$$\arg \max_{\substack{y \\ t=1}} \frac{1}{T_y^\alpha} \sum_{t=1}^{T_y} \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

$\alpha = 0$  No Normalization  
 $\alpha = 1$  Full Normalization

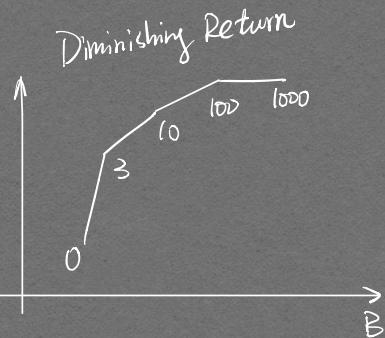
Normalized Log Likelihood

→ Choose Beam Width  $B$

Large  $B \rightarrow$  Better Results, Slower

Small  $B \rightarrow$  Worse Results, Faster

Beam Search is an approximate optimization algorithm, or heuristic algorithm



## Error Analysis of Beam Search

Use RNN to compute  $P(\hat{y}|x), P(y^*|x)$

Case 1:  $P(\hat{y}|x) > P(y^*|x) \rightarrow$  Beam Search doesn't find the optimal  $y$

Case 2:  $P(y^*|x) \leq P(\hat{y}|x) \xrightarrow{\text{RNN computes } P(y|x) \text{ wrong}} \text{But } y^* \text{ is better than } \hat{y}$

## BLEU Score

Bilingual Evaluation Understudy

BLEU : A method for automatic of machine translation [Papineni et al, 2002]

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

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

Combined BLEU Score

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

Brevity Penalty

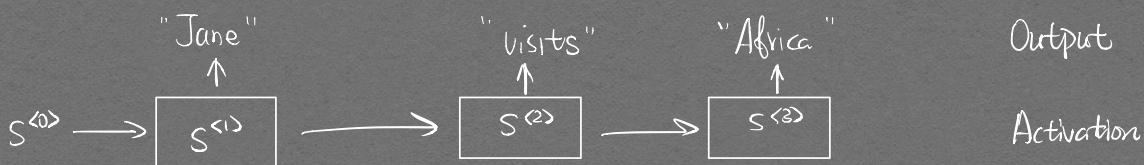
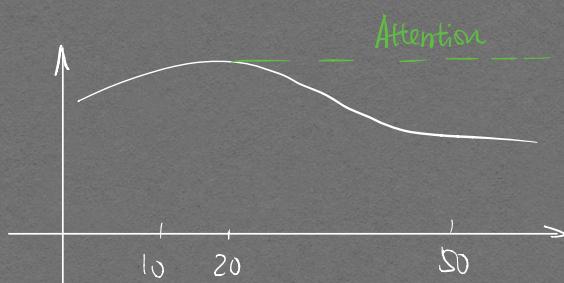
$$BP = \begin{cases} 1 & \text{if Machine Translation Output Length} < \text{Reference Length} \\ \exp\left(1 - \frac{\text{Reference Output Length}}{\text{Machine Translation Output Length}}\right) & \text{otherwise} \end{cases}$$

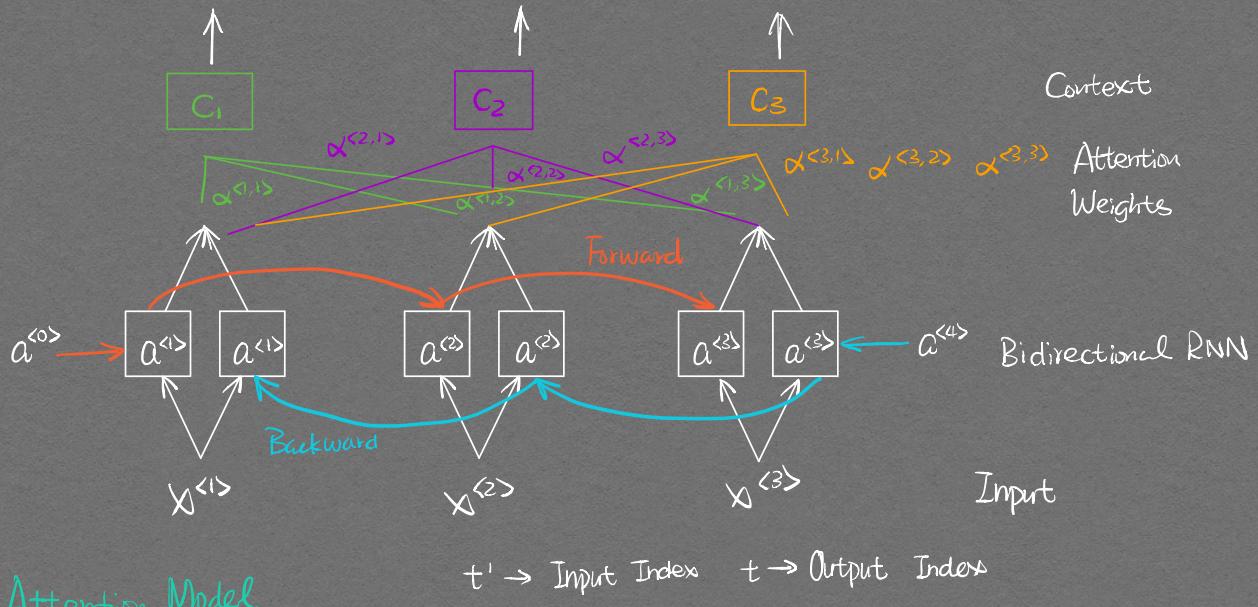
Attention Model Intuition

Problem of Long Sequence.

Neural Machine Translation by jointly learning to align and translate

Bahdanau et al. 2014





### Attention Model

$$a^{<t>} = (\vec{a}^{<t>} , \bar{a}^{<t>}) \quad \text{Feature Vector}$$

$$\sum_{t'} \alpha^{<1,t'} = 1 \quad \alpha^{<t,t'} = \frac{\text{amount of "attention" } y^{(t)} \text{ should pay to } a^{<t>}}{\sum_{t'} \alpha^{<1,t'} a^{<t>}}$$

How to Compute  $\alpha^{<t,t'>}$

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



Attention also applies to :

- 1) Image Caption
- 2) Date Format    23 April, 1564  $\Rightarrow$  1564-04-23 !
- 3) Speech Recognition

## Speech Recognition

$x \rightarrow y$   
audio clip transcript

Before : phonemes (units of sound)

Academia 300h  
3000h

Commercial 10,000h

CTC cost (Connectionist Temporal Classification)

10s 100 Hz Audio  $\rightarrow$  5 words

"the quick brown fox"

Generate output with blank spaces

ttt - -  $\downarrow$  h - e - -  $\uparrow$  - q q - -  
Blank Space

## Triggerword Detection

