# Sequence to Sequence Architectures

**Translation**  $x^{<1>}$  $x^{<2>}$  $x^{<3>}$  $x^{<4>}$  $x^{<5>}$

e.j.  Jane  visite  l'Afrique  en  septembre

$\rightarrow$  Jane  is  visiting  Africa  in  September

$y^{<1>}$  $y^{<2>}$  $y^{<3>}$  $y^{<4>}$  $y^{<5>}$  $y^{<6>}$



Encoder          Decoder

$y^{<1>}$ ...

**Image to Caption**      e.j.   A  cat  is  sitting  on  a  chair



**language model**



**Machine translation**



"Conditioned language model"

$$P(y^{<1>}, \dots, y^{<T_y>} | x^{<1>}, \dots, x^{<T_x>})$$

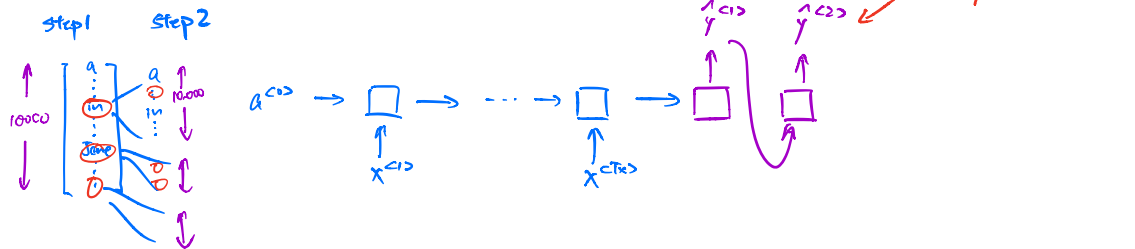$\Rightarrow$ Pick the most likely sentence

Why not greedly search?

e.f.  Jane is visiting Africa  ✓

Jane is going to visit Africa  ✗

$P(\text{Jane is going} | x)$ higher $\Rightarrow$ less good sentence

## Beam Search

$B\nearrow$ slower
$B\searrow$ faster. worse result

$B=3$ (beam width)
$\Rightarrow P(y^{<1>} | x)$

3 copies



step1    step2

10000

a<0> → □ → ... → □ → □ □

$x^{<1>}$    $x^{<Tx>}$

$\hat{y}^{<1>}$    $\hat{y}^{<2>}$

$$P(y^{<1>}, y^{<2>} | x) = P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>})$$

step 3

in    september    

jane    is    

jane    visits    

## Improvements to Beam Search

$P(y^{<2>} | x) P(y^{<3>} | x, y^{<1>}) \cdots$
— can be small
— prefer short sentences

— length normalization

$$\arg\max_y \prod_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>}, \cdots, y^{<t-1>})$$
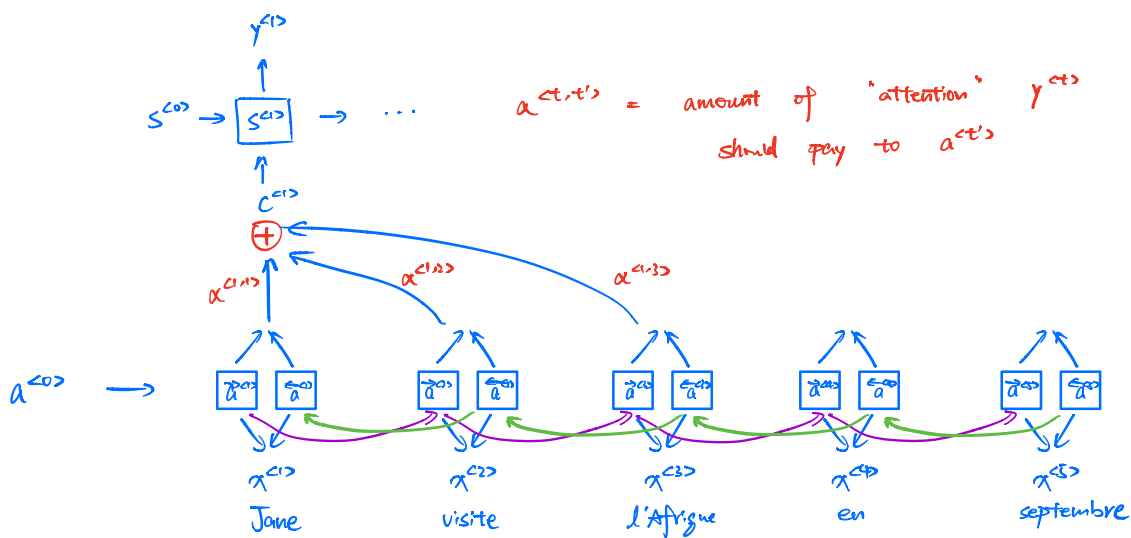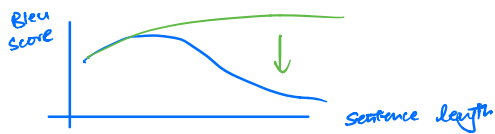
$$\hookrightarrow \arg\max_y \frac{1}{T_y^\alpha} \prod_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \cdots, y^{<t-1>})$$

normalize    (e.f. $\alpha = 0.7$)

## Attention Model

The problem w/ long sequence
$\Rightarrow$ hard to memorize



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



Jane     visite     l'Afrique     en     septembre

$$\sum_{t'} \alpha^{<1,t'>} = 1 \qquad a^{<t>} = (\overrightarrow{a}^{<t'>}, \overleftarrow{a}^{<t'>})$$

$$c^{<1>} = \sum_{t'} \alpha^{<1,t'>} a^{<t'>}$$

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

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



$$s^{<t-1>} \searrow \quad \boxed{\vdots} \rightarrow e^{<t,t'>}$$
$$a^{<t'>} \nearrow$$

train a small NN