Sequence to
sequence models

Attention model

deeplearning.ai

# Attention model



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

$c^{<2>} = \sum_{t'} \alpha^{<2, t'>} a^{<t'>}$

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

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

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

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

$c^{<2>}$

$\alpha^{<1,1>}$  $\alpha^{<1,2>}$  $\alpha^{<1,3>}$

$a^{<0>} \rightarrow$

$\overleftarrow{a}^{<6>}$

$x^{<1>}$ jane   $x^{<2>}$ visite   $x^{<3>}$ l'Afrique   $x^{<4>}$ en   $x^{<5>}$ septembre

$t'$

[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]
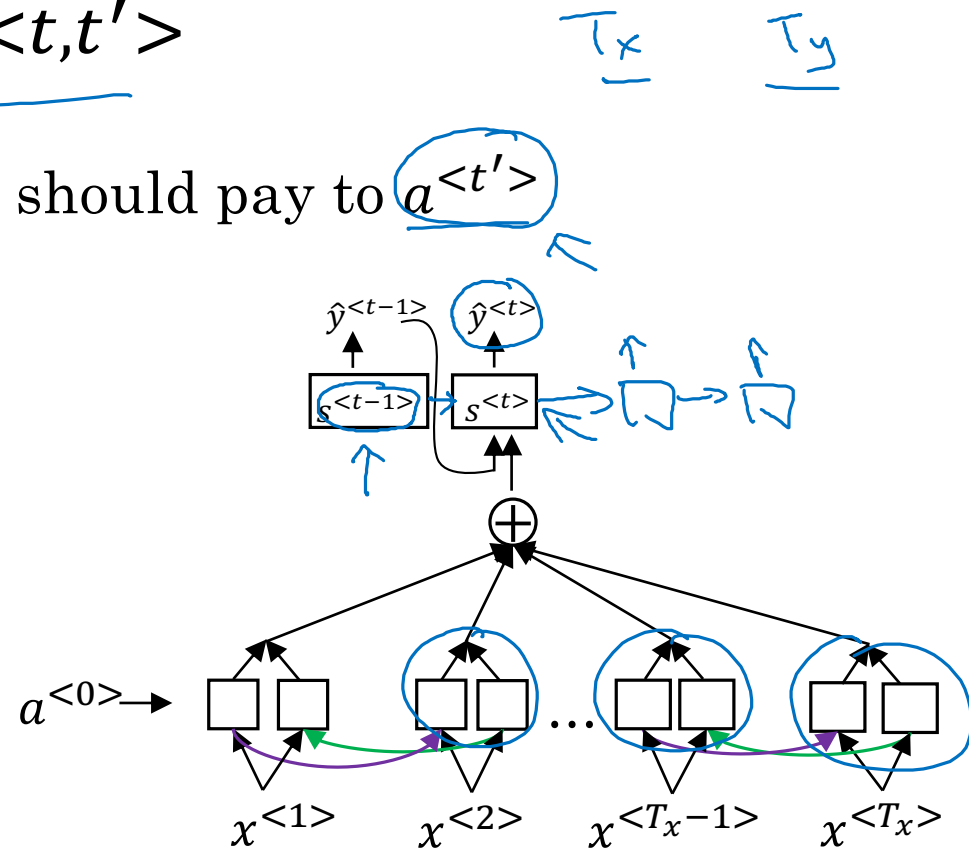
Andrew Ng

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

$T_x$    $T_y$

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

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

$s^{<t-1>}$

$a^{<t'>}$

$e^{<t,t'>}$

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

$\hat{y}^{<t-1>}$    $\hat{y}^{<t>}$

$s^{<t-1>}$    $s^{<t>}$

$\oplus$

$a^{<0>} \rightarrow$

$x^{<1>}$    $x^{<2>}$    $x^{<T_x-1>}$    $x^{<T_x>}$

[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

[Xu et. al., 2015. Show, attend and tell: Neural image caption generation with visual attention]

Andrew Ng

# Attention examples

July 20th 1969 $\longrightarrow$ $1969 - 07 - 20$

23 April, 1564 $\longrightarrow$ $1564 - 04 - 23$

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