



deeplearning.ai

# Sequence to sequence models

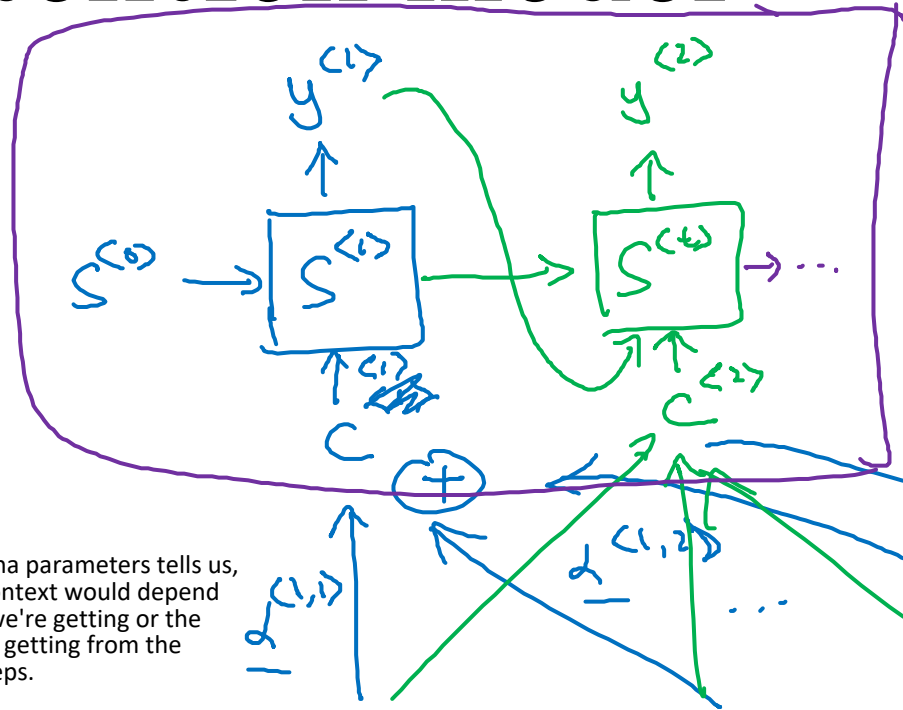
---

## Attention model

Let's now formalize that intuition into the exact detail of how you would implement an attention model.

# Attention model

$\alpha^{(t,t')}$  = amount of 'attention'  $y^{(t)}$  should pay to  $a^{(t')}$ .



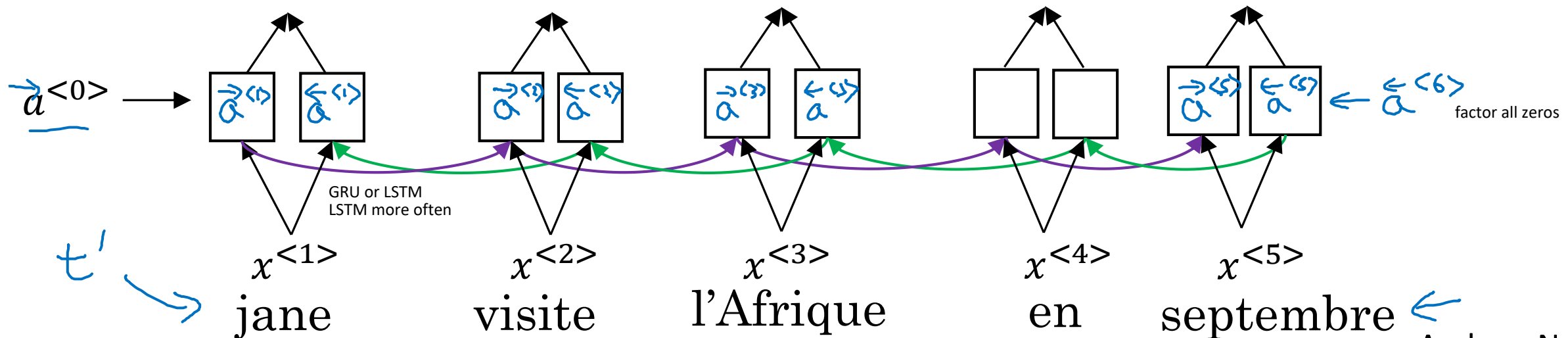
$$C^{(2)} = \sum_{t'} \alpha^{(2,t')} a^{(t')}$$

$$a^{(t')} = (\vec{a}^{(t')}, \leftarrow a^{(t')})$$

$$\sum_{t'} \alpha^{(1,t')} = 1 \quad \text{zero positive sum to 1}$$

$$C^{(1)} = \sum_{t'} \alpha^{(1,t')} a^{(t')}$$

And so these alpha parameters tells us, how much the context would depend on the features we're getting or the activations we're getting from the different time steps.



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

total # of attention parameters are going to be  $T_x$  times  $T_y$

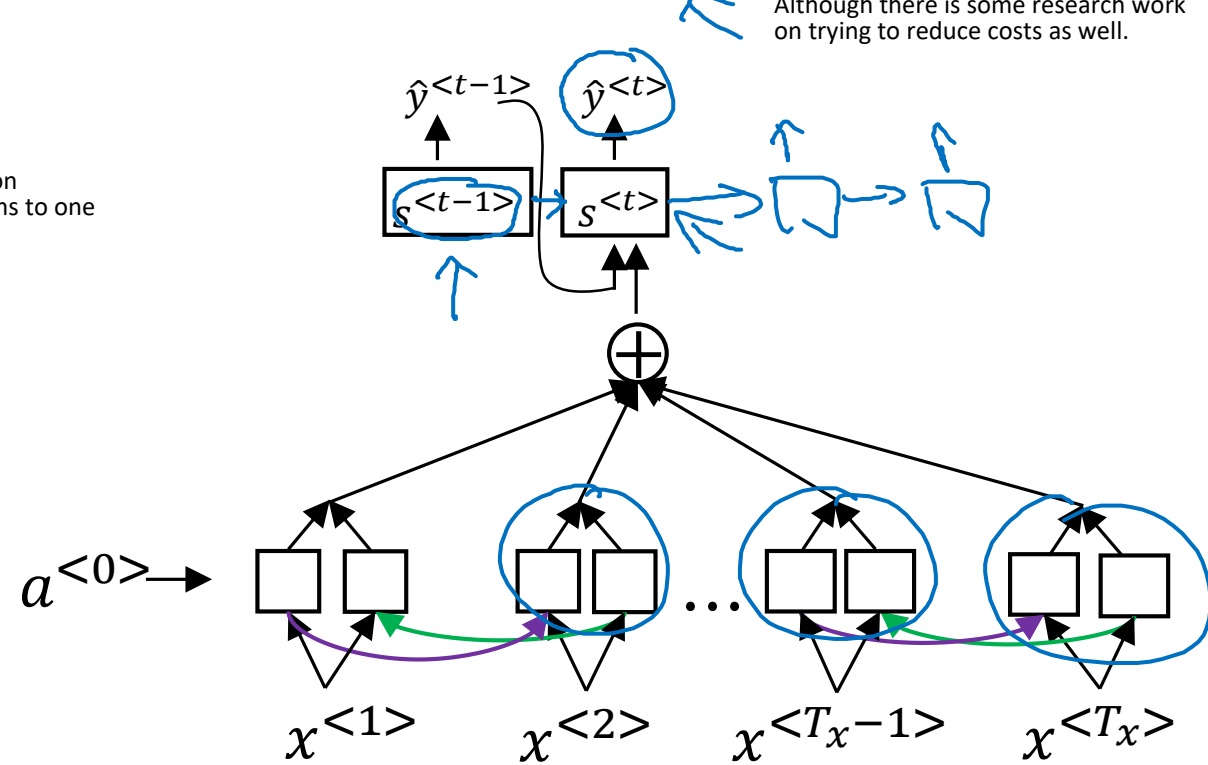
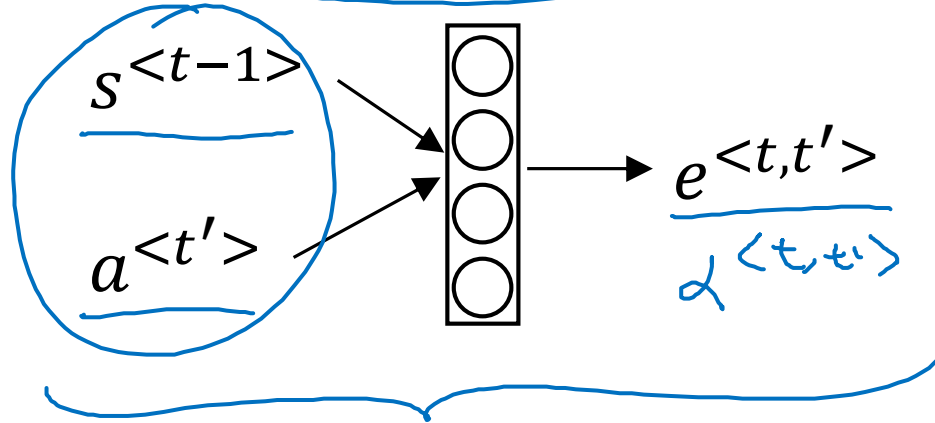
$T_x$   $T_y$

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

Although in machine translations where neither input nor output sentences is usually that long, maybe quadratic is actually acceptable. Although there is some research work on trying to reduce costs as well.

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

using this softmax prioritization just ensure this property sums to one



One downside of this algorithm is that it does take quadratic time or quadratic cost to run this algorithm.

[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 → 1969 – 07 – 20

23 April, 1564 → 1564 – 04 – 23

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

