

Recurrent Neural Nets

Introduction to NLP

Rahul Mishra

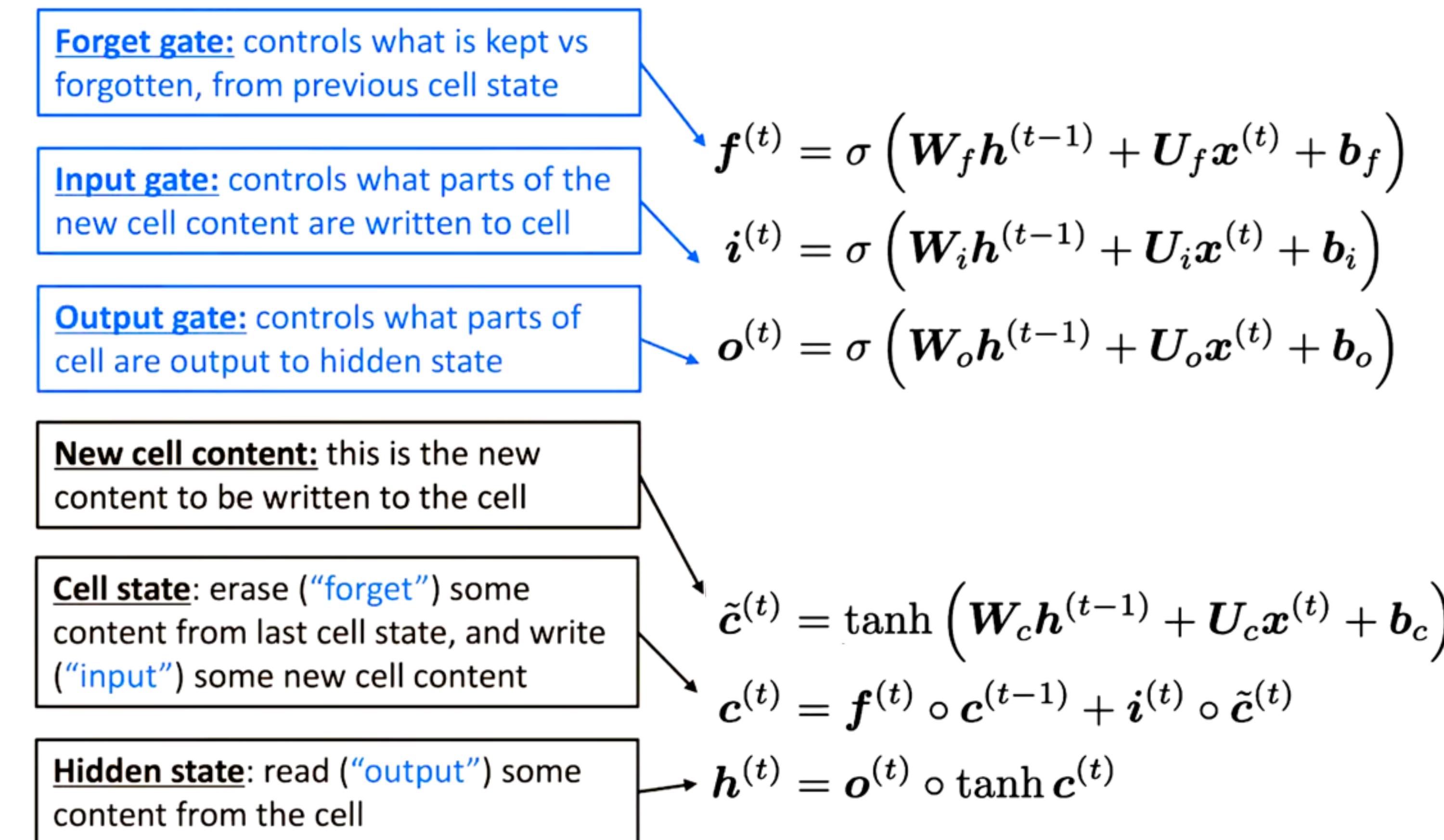
IIIT-Hyderabad
March 12, 2024

INTERNATIONAL INSTITUTE OF
INFORMATION TECHNOLOGY

H Y D E R A B A D

LSTM in pictures

We have a sequence of inputs $x^{(t)}$, and we will compute a sequence of hidden states $h^{(t)}$ and cell states $c^{(t)}$. On timestep t :



C'mon, it's
been around
for 20 years!



LSTM

GRU

$$\tilde{c}^{<t>} = \tanh(W_c [\Gamma_r * \underline{c}^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u [c^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r [c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$a^{<t>} = c^{<t>}$$

LSTM

$$\tilde{c}^{<t>} = \tanh(W_c [a^{<t-1>}, x^{<t>}] + b_c)$$

$$(update) \quad \Gamma_u = \sigma(W_u [a^{<t-1>}, x^{<t>}] + b_u)$$

$$(forget) \quad \Gamma_f = \sigma(W_f [a^{<t-1>}, x^{<t>}] + b_f)$$

$$(output) \quad \Gamma_o = \sigma(W_o [a^{<t-1>}, x^{<t>}] + b_o)$$

$$c^{<t>} = \Gamma_u * \underline{\tilde{c}^{<t>}} + \Gamma_f * \underline{c^{<t-1>}}$$

$$a^{<t>} = \Gamma_o * \underline{c^{<t>}}$$

LSTM in pictures

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

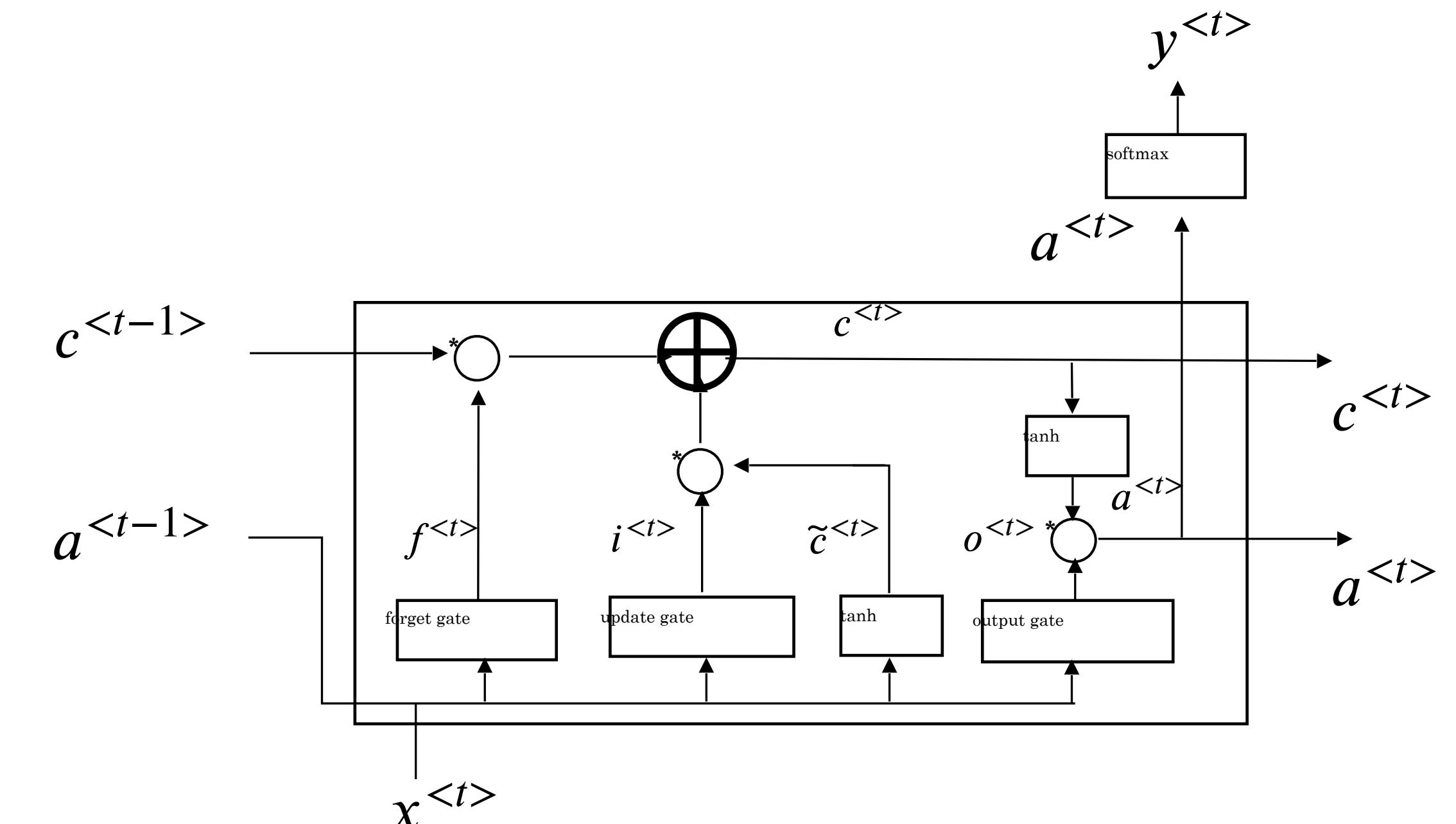
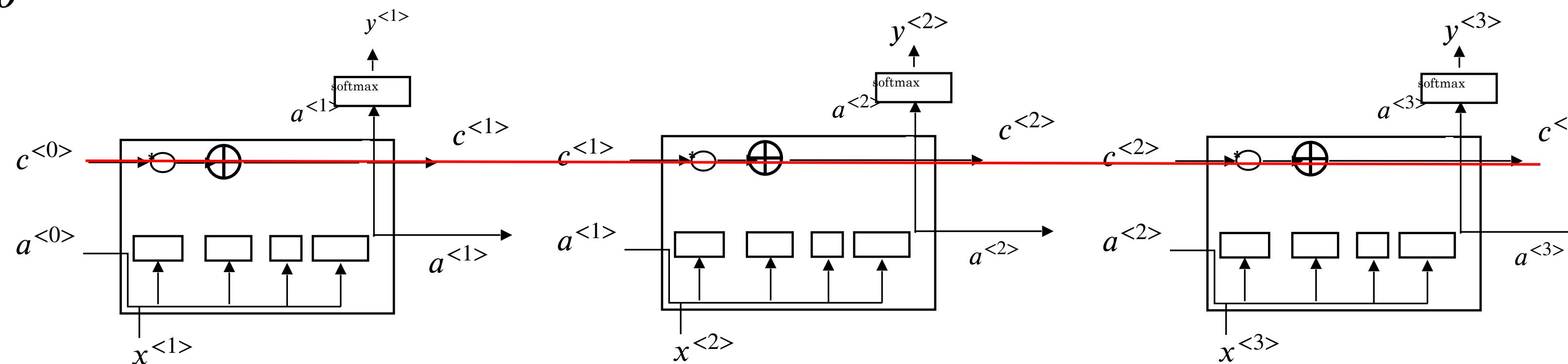
$$\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

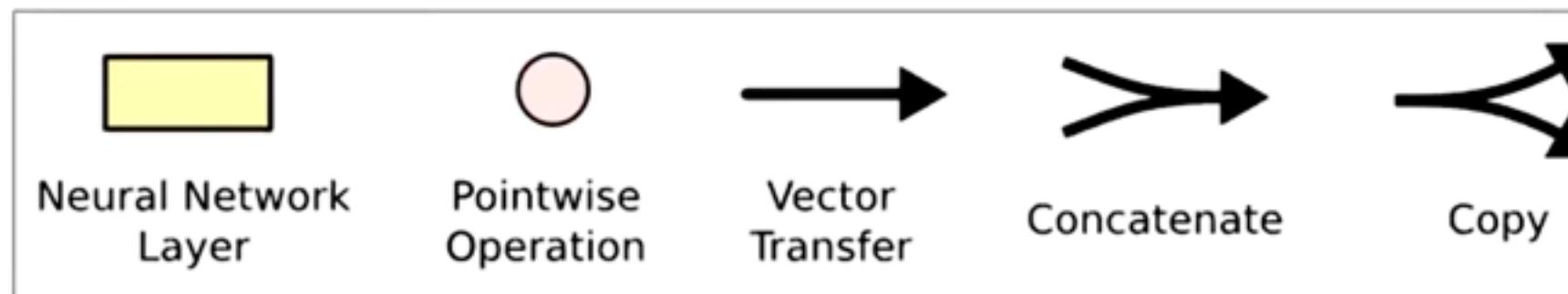
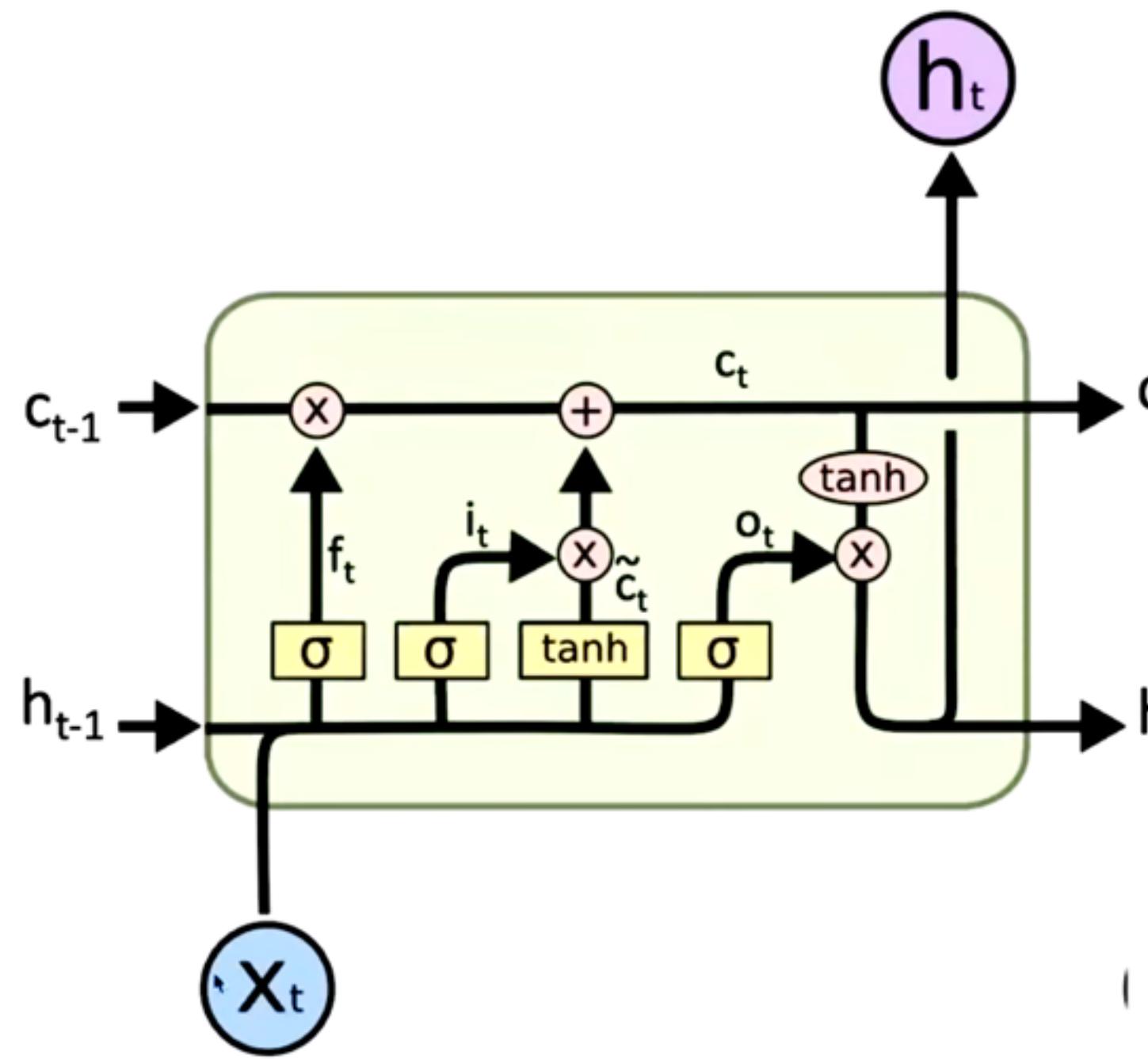
$$\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * c^{<t>}$$

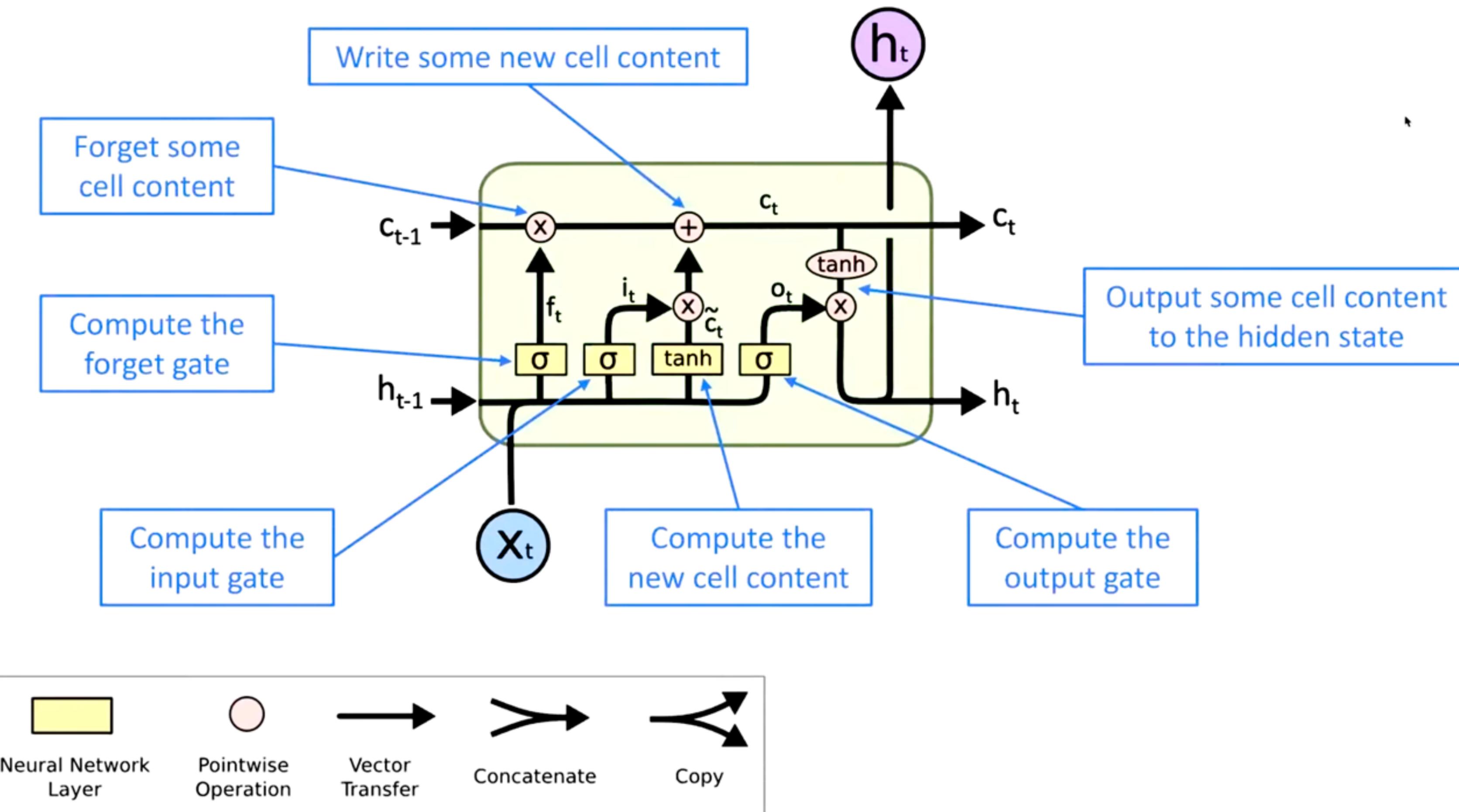


LSTM in pictures



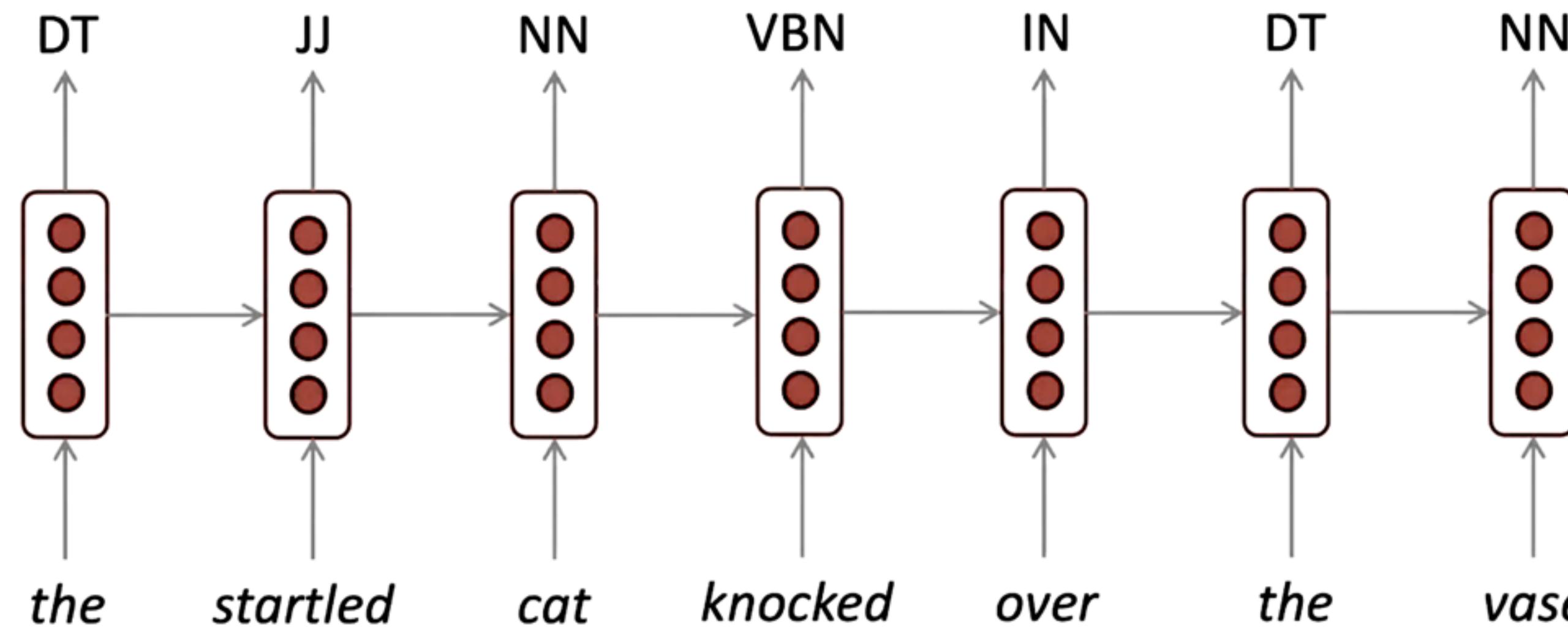
Source: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

LSTM in pictures



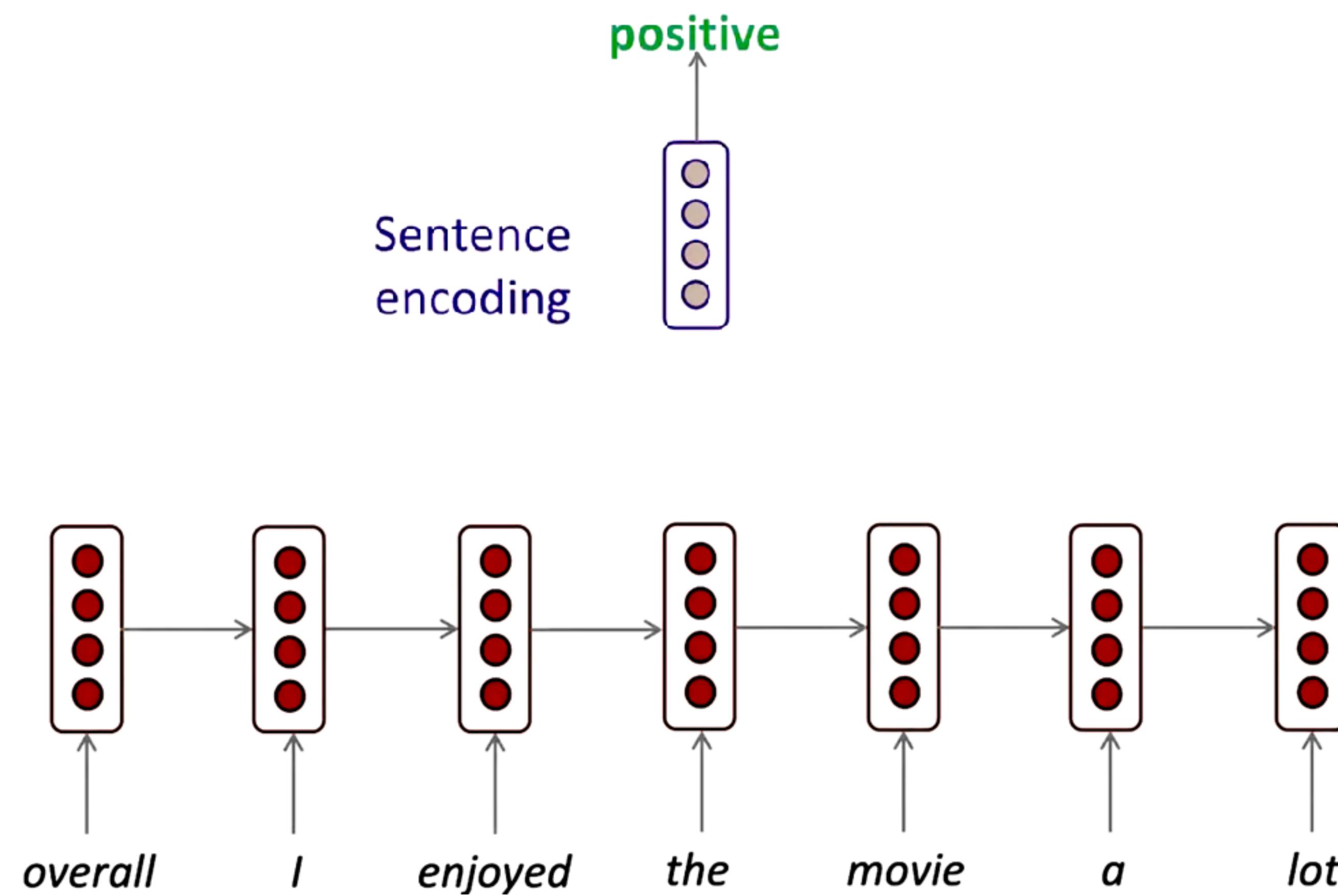
Source: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Applications: Sequence tagging

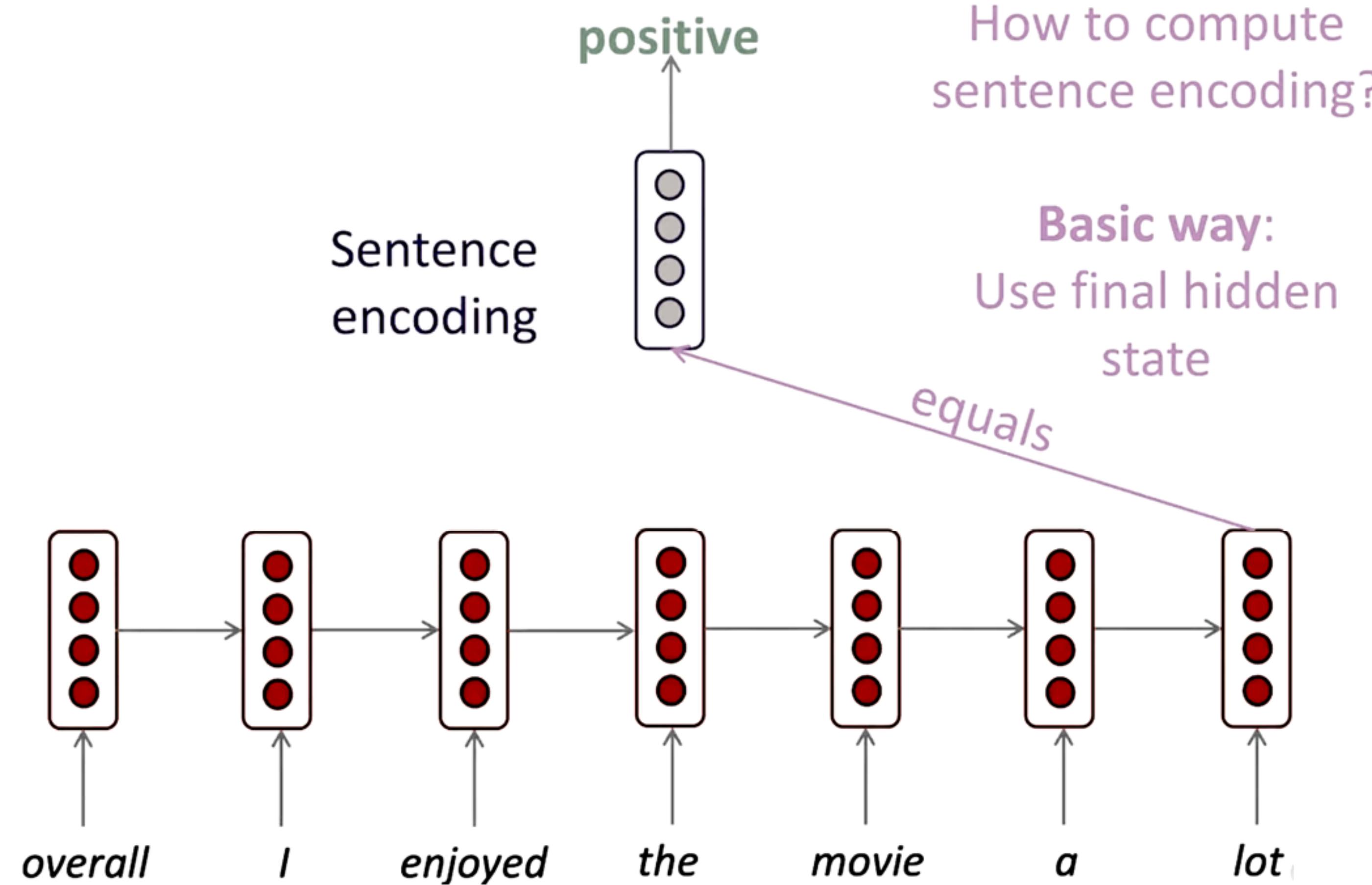


e.g., part-of-speech tagging, named entity recognition

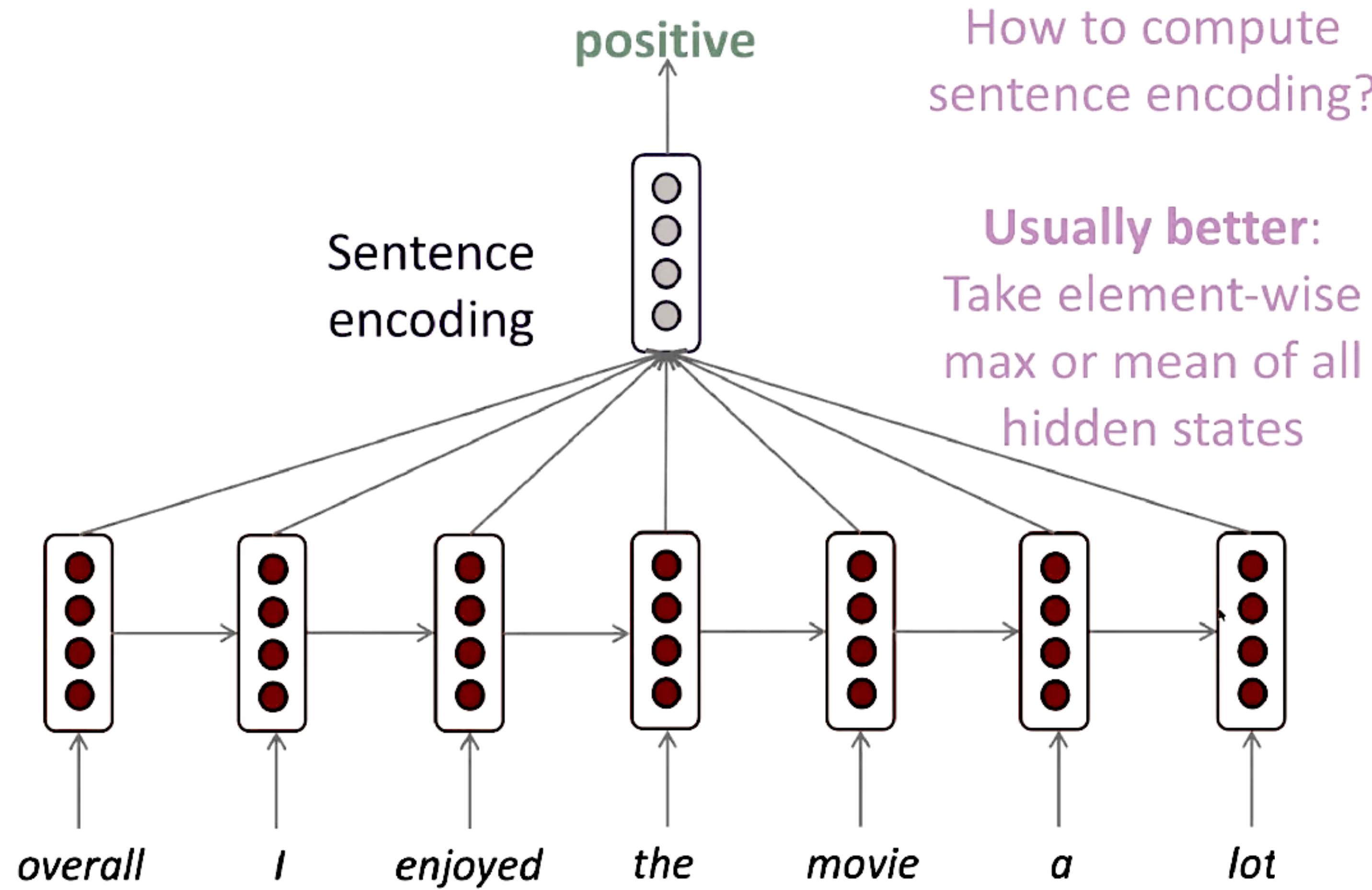
Applications: Sequence Classification



Applications: Sequence Classification

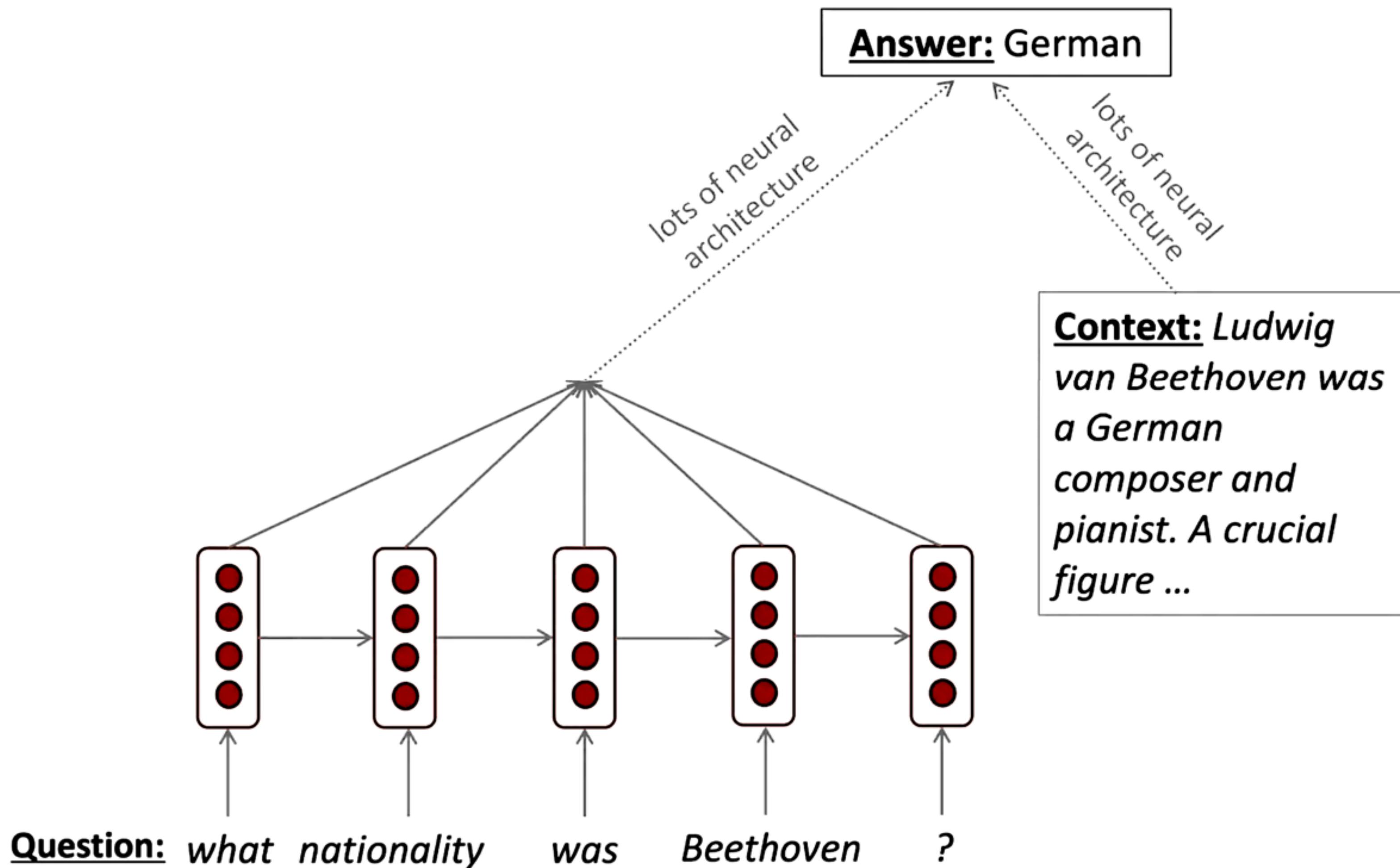


Applications: Sequence Classification



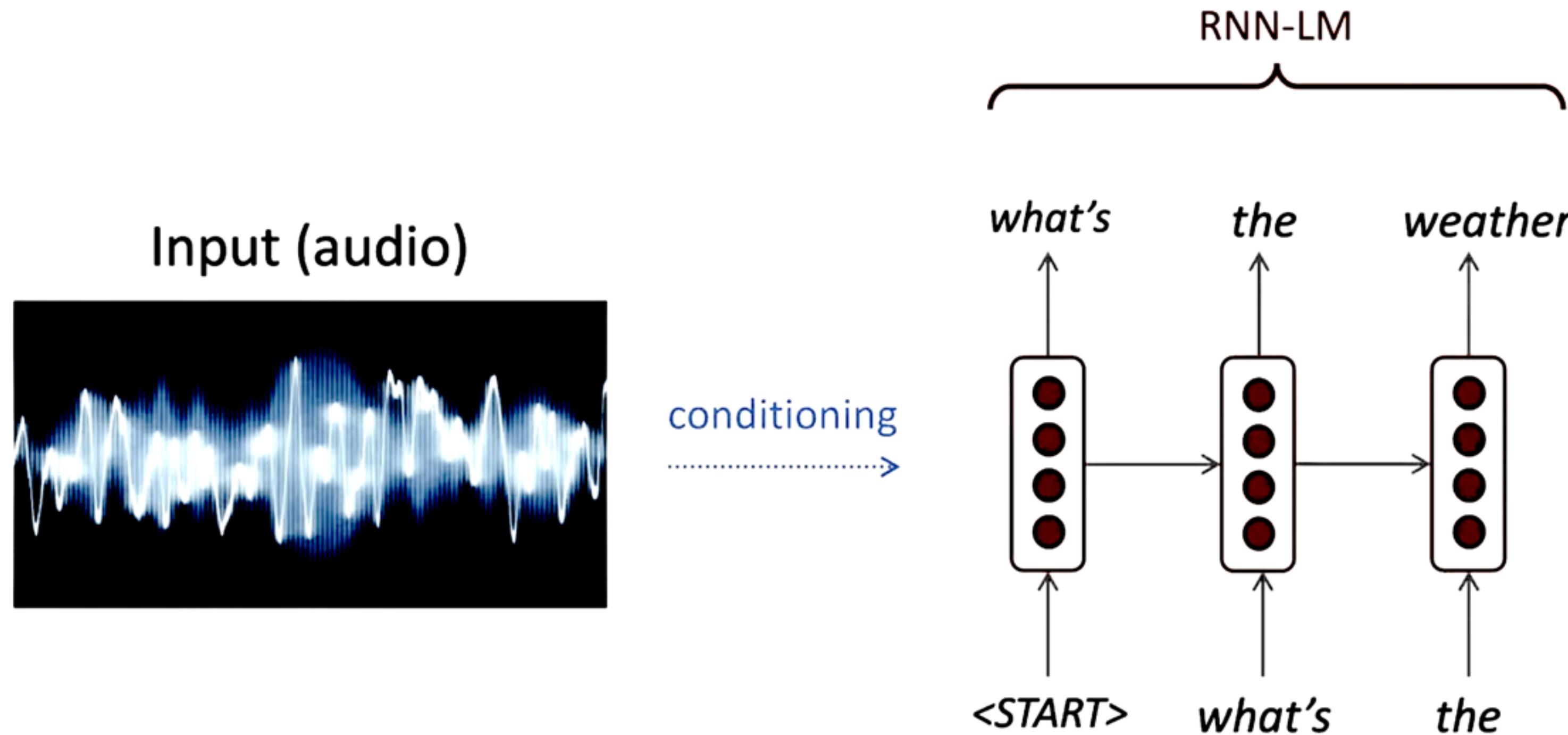
Applications: as Encoder

e.g., **question answering**, machine translation, *many other tasks!*



Applications: as Decoder

e.g., speech recognition, machine translation, summarization

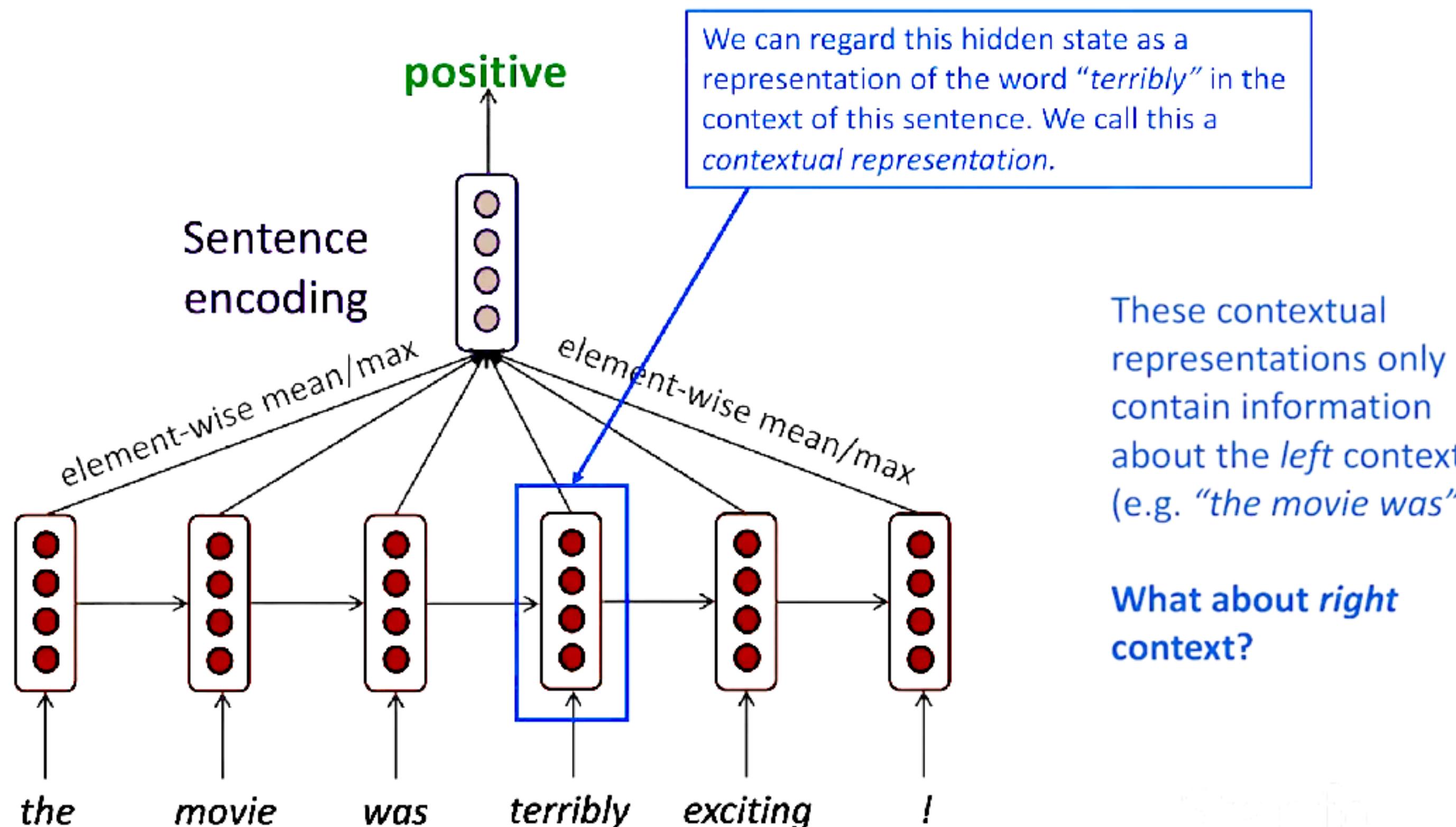


Bidirectional-RNN/LSTM

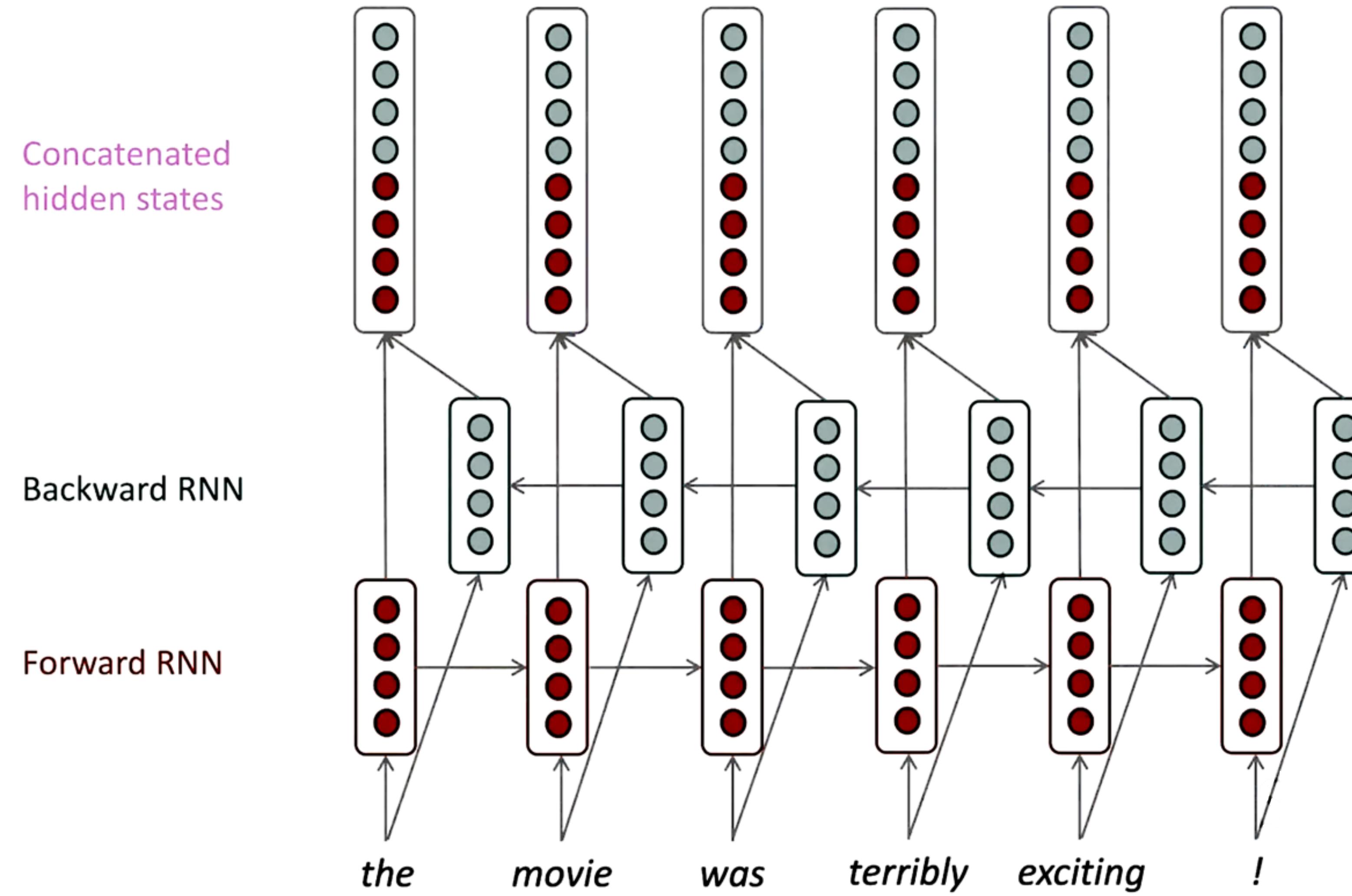
He said, “Teddy bears are on sale!”

He said, “Teddy Roosevelt was a great President!”

Task: Sentiment Classification



Bidirectional-RNN/LSTM



Bidirectional-RNN/LSTM

On timestep t :

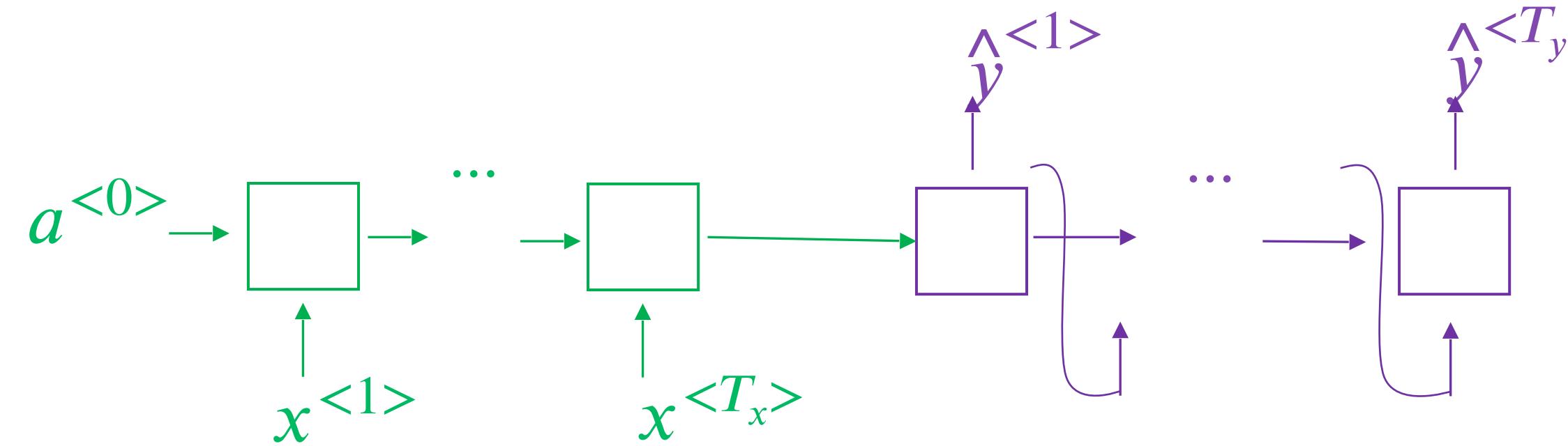
This is a general notation to mean “compute one forward step of the RNN” – it could be a simple, LSTM, or other (e.g., GRU) RNN computation.

Forward RNN $\vec{h}^{(t)} = \text{RNN}_{\text{FW}}(\vec{h}^{(t-1)}, \mathbf{x}^{(t)})$

Backward RNN $\overleftarrow{h}^{(t)} = \text{RNN}_{\text{BW}}(\overleftarrow{h}^{(t+1)}, \mathbf{x}^{(t)})$

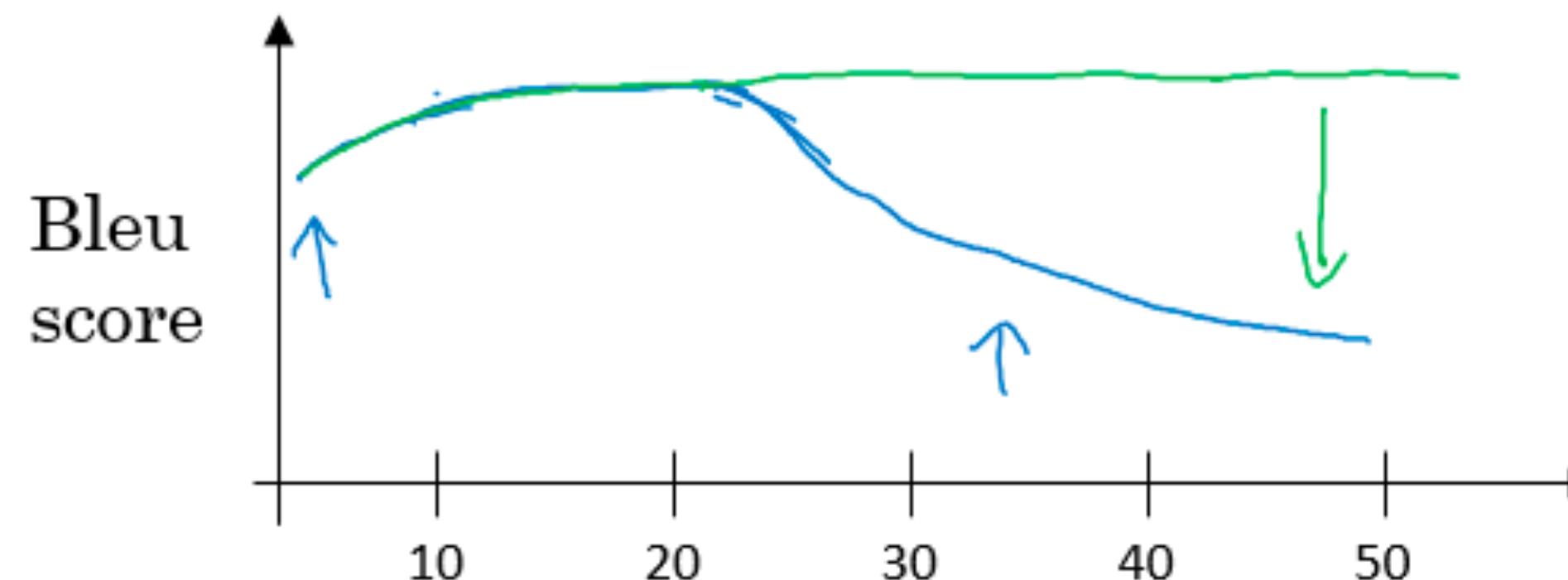
Concatenated hidden states $\mathbf{h}^{(t)} = [\vec{h}^{(t)}; \overleftarrow{h}^{(t)}]$

The problem of long sequences

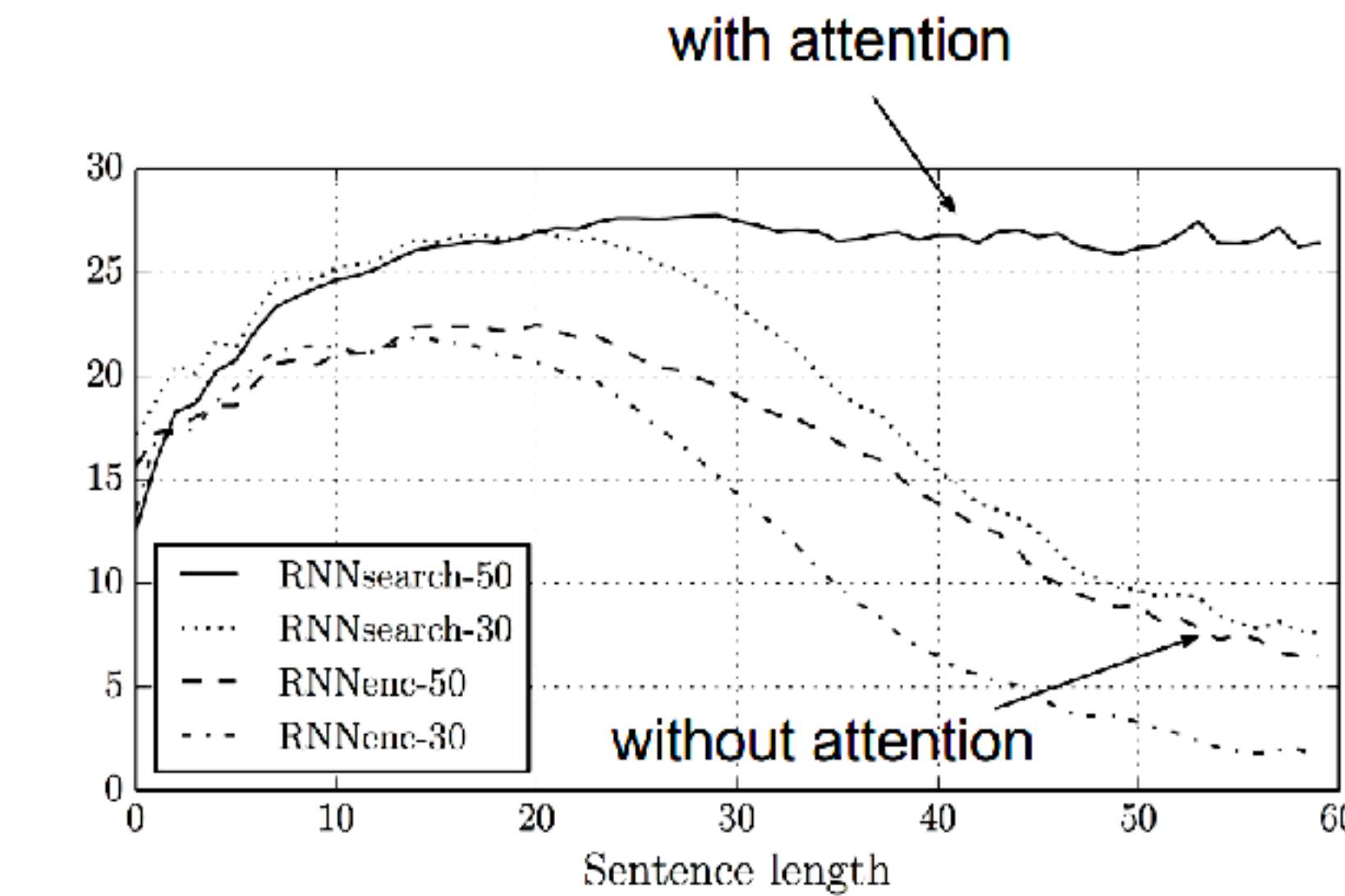
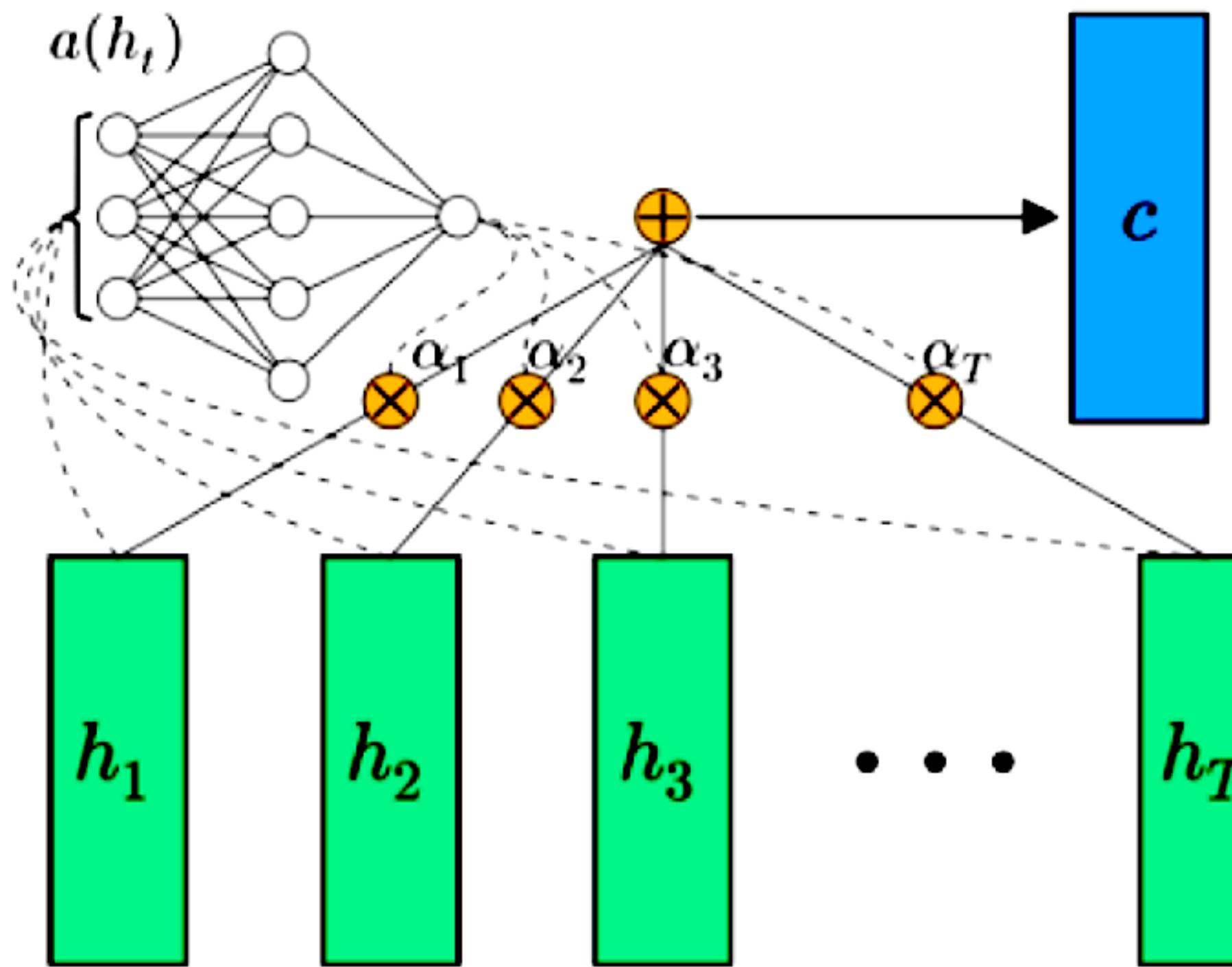


Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.

Jane went to Africa last September, and enjoyed the culture and met many wonderful people; she came back raving about how wonderful her trip was, and is tempting me to go



Attention with RNN/LSTM



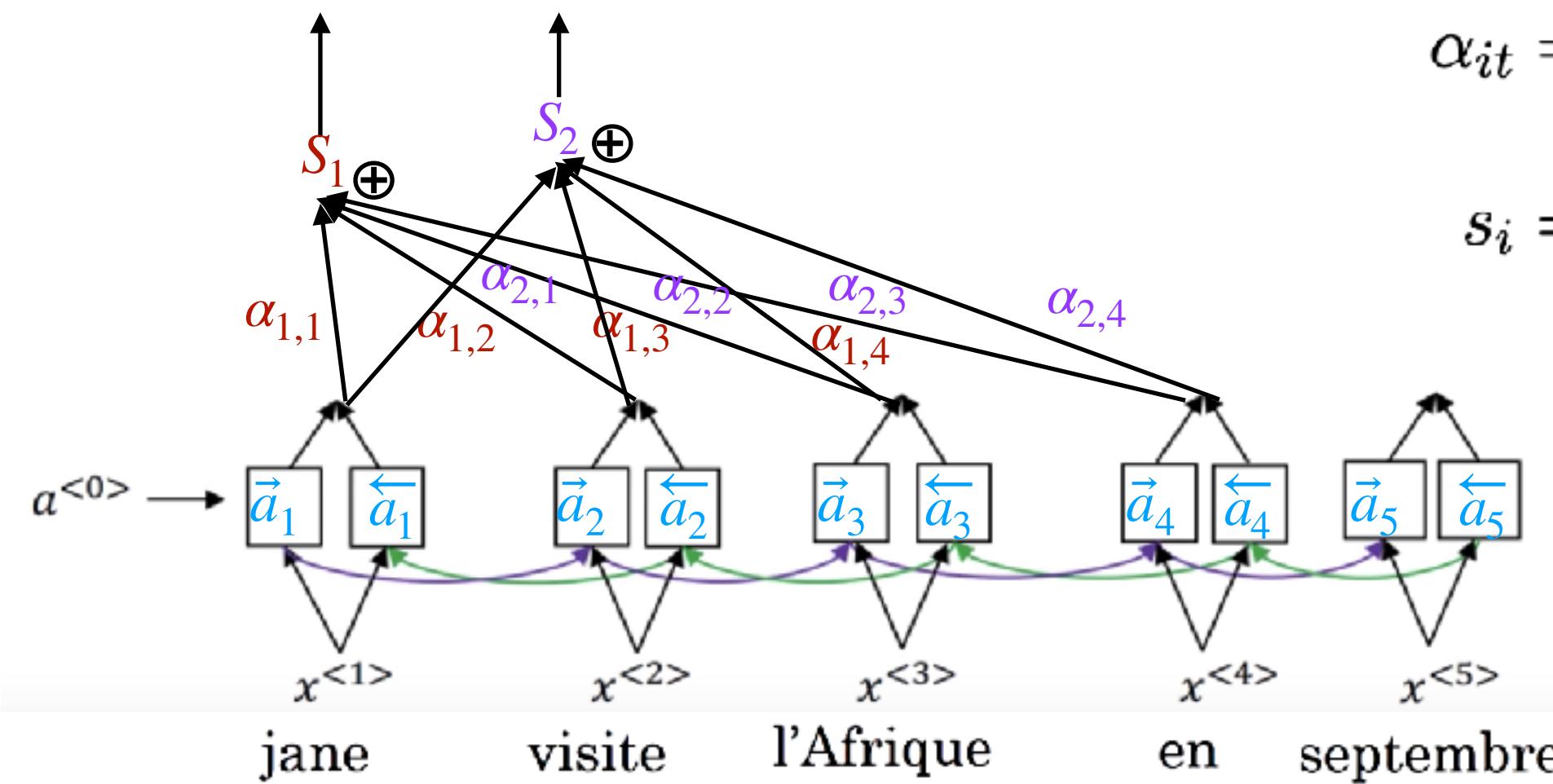
Attention with RNN/LSTM

$$\sum_t \alpha_{1,t} = 1$$

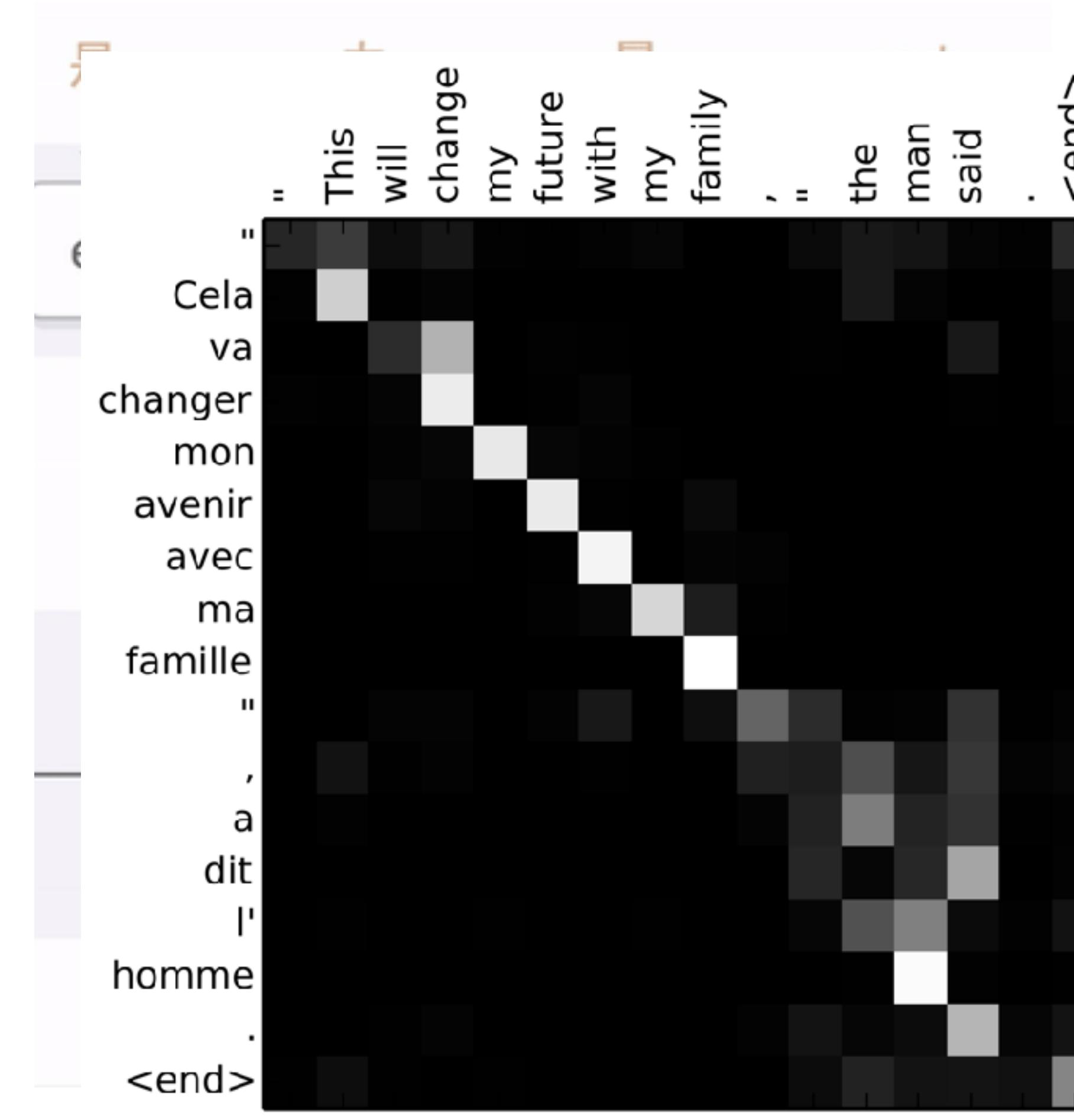
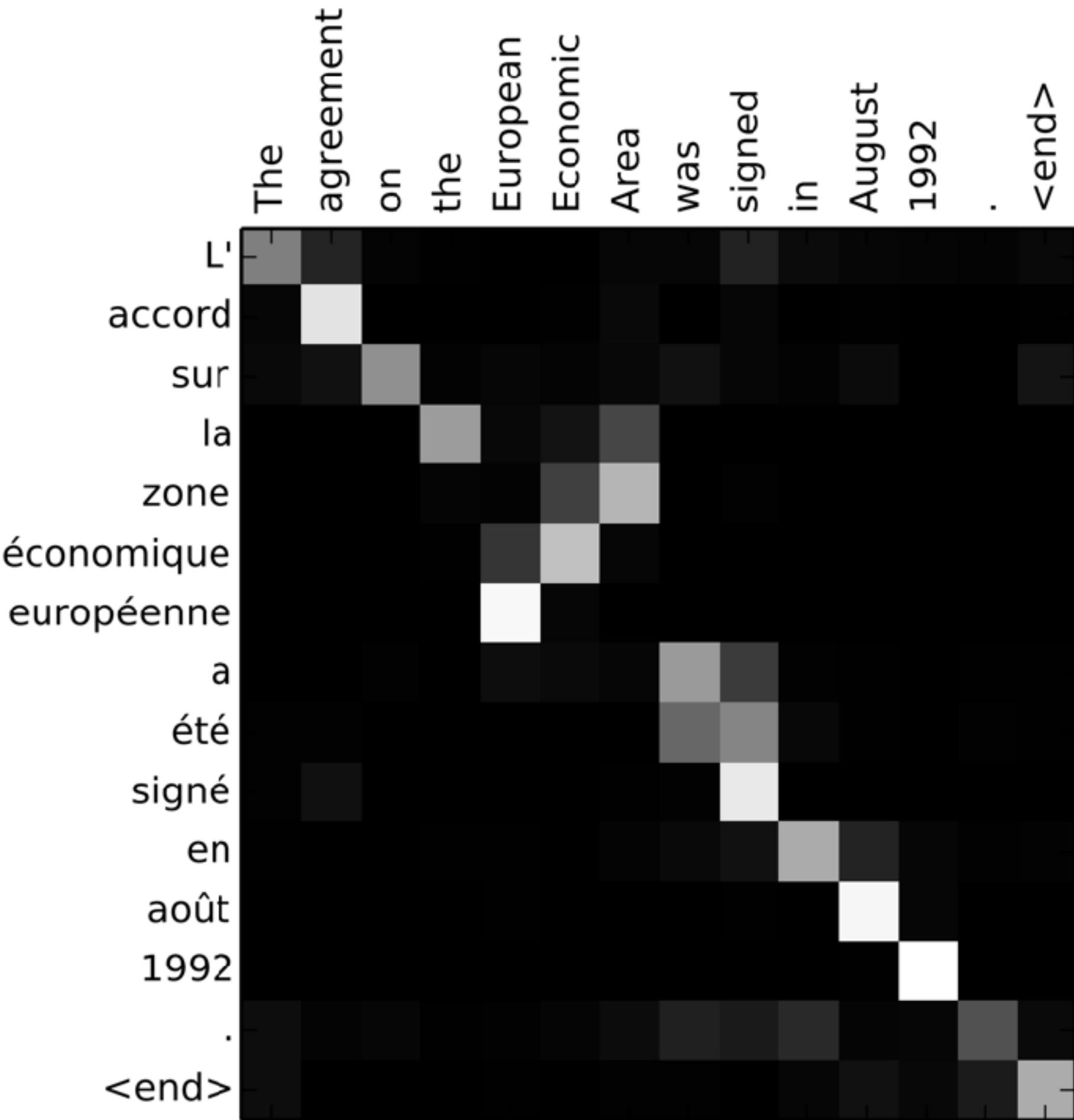
$$u_{it} = \tanh(W_w h_{it} + b_w)$$

$$\alpha_{it} = \frac{\exp(u_{it}^\top u_w)}{\sum_t \exp(u_{it}^\top u_w)}$$

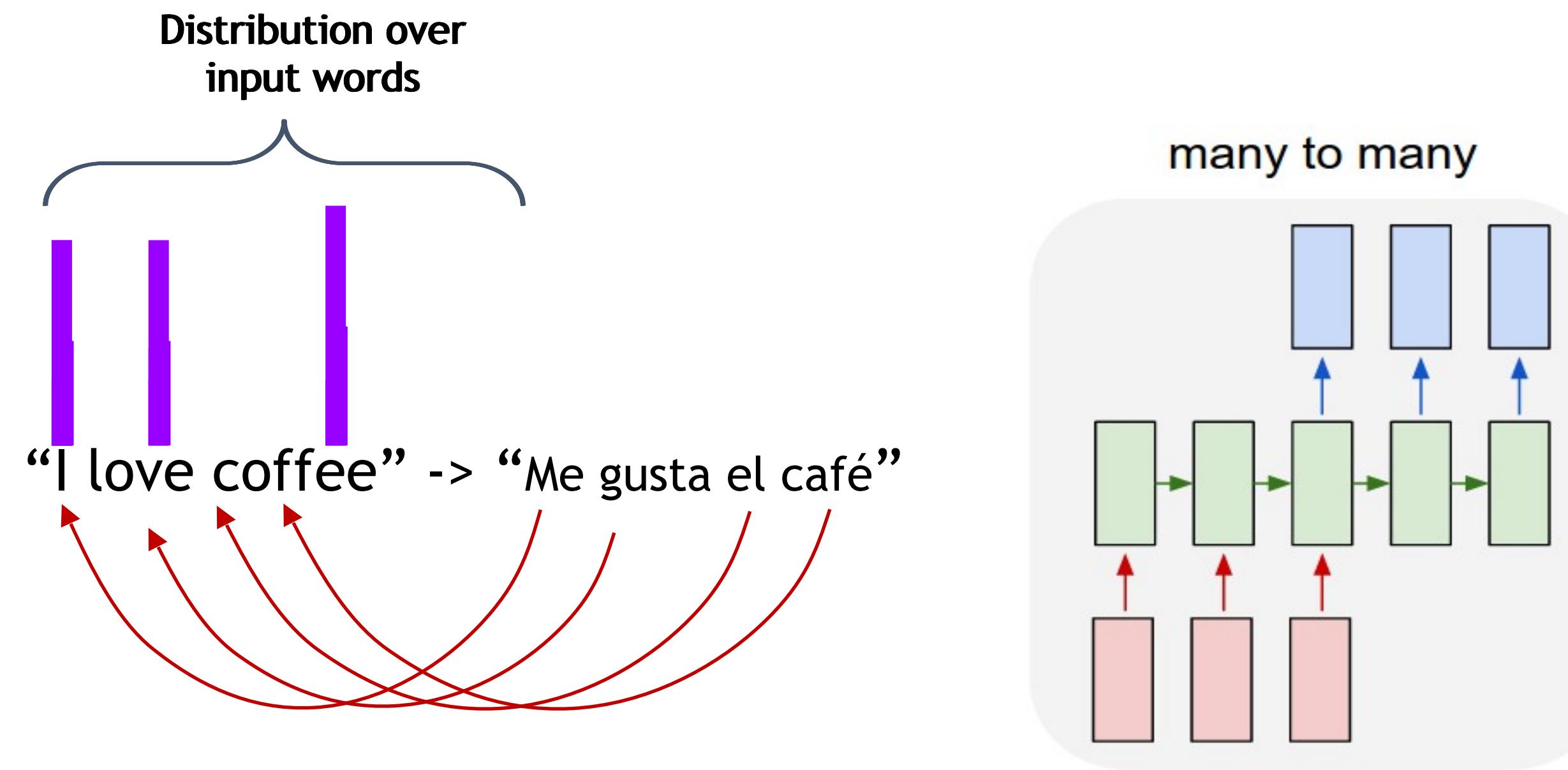
$$s_i = \sum_t \alpha_{it} h_{it}.$$



Attention with RNN/LSTM



LSTM with Self Attention



Bahdanau et al, “Neural Machine Translation by Jointly Learning to Align and Translate”, ICLR 2015

Who is who?

