

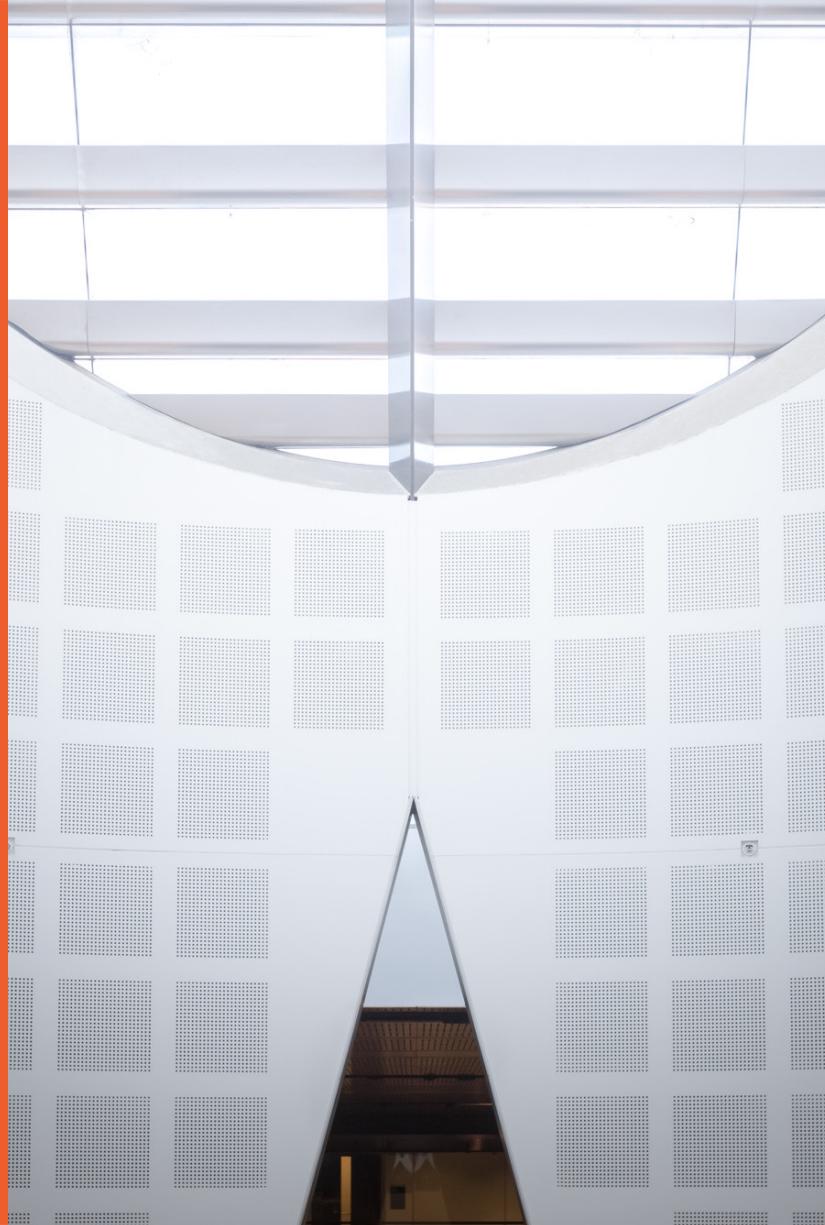
Recurrent Neural Networks

Dr Chang Xu

School of Computer Science



THE UNIVERSITY OF
SYDNEY



Example Application

sequential analysis

Samsung Mobile  @SamsungMobile

关注

There's a new star in town. And it's going to change the way you experience everything. Introducing the #GalaxyS9. Find out more.

翻译自英语



Introducing the
Galaxy S9+
The Camera. Reimagined.

下午1:05 - 2018年2月26日

Samsung Mobile US  @SamsungMobileUS · 2月25日

Nothing wrong with sticking to what you know.

翻译自英语

4 4 19

Emma Driver @emmab1988 · 2月25日

Absolutely. Have nothing  love for my Samsung. Spent 7 years as an iPhone user  but now I'd never look back.

翻译自英语

7



Shubham Kumar @iamshubhmKr · 2月25日

回复 @SamsungMobile

Same old boring Samsung.

翻译自英语

1 1 5

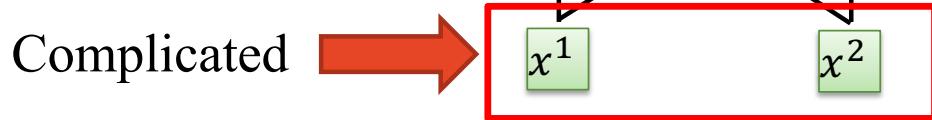


Image credit to <https://twitter.com/SamsungMobile/status/967807667463958531>

Recurrent Neural Network

- Input
 - Word sequence
 - Each word is represented by vector
- Methods
 - One hot encoding
 - Word hash
 - Word embedding

Complicated



Recurrent Neural Network

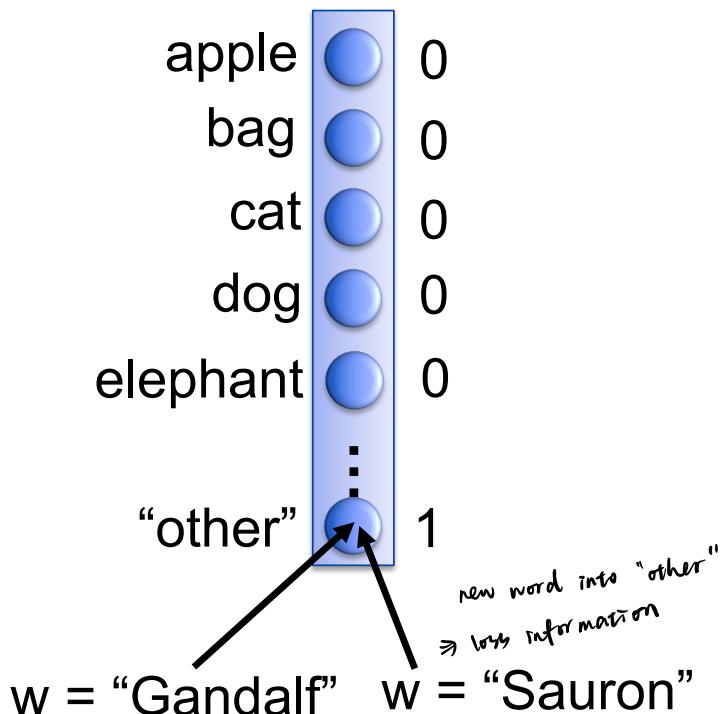
□ One hot encoding

- The length of vector is lexicon size
 - Each dimension corresponds to a word in the lexicon
 - The dimension for the word is 1, and others are 0
- lexicon = {apple, bag, cat, dog, elephant}

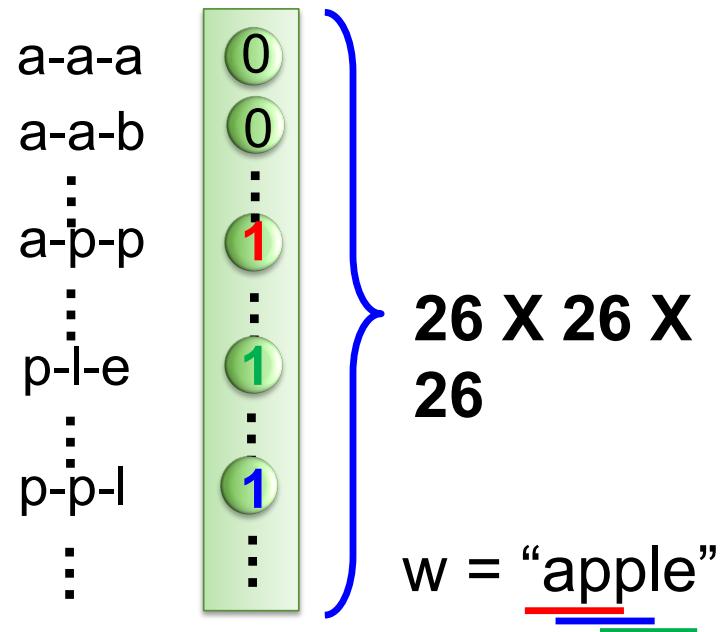
Word	D_1	D_2	D_3	D_4	D_5
Apple	1	0	0	0	0
Bag	0	1	0	0	0
Cat	0	0	1	0	0
Dog	0	0	0	1	0
Elephant	0	0	0	0	1

Recurrent Neural Network

Dimension for “Other”



Word hashing



Recurrent Neural Network

□ Output

□ Probability distribution that the sentence belongs to negative or positive.

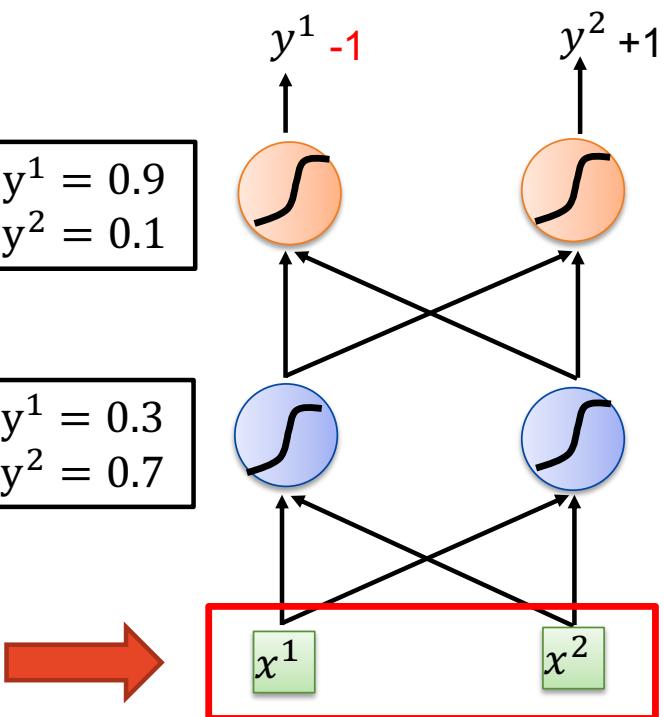
Complicated? Yes,
and a little slow, too.

$$\begin{aligned}y^1 &= 0.9 \\y^2 &= 0.1\end{aligned}$$

But the animation is as
colorful as the story and
emotions do eventually reach
the desired level.

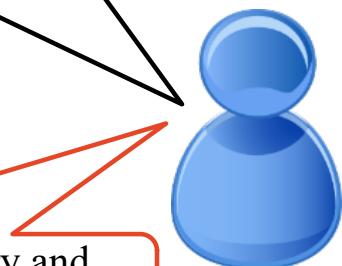
$$\begin{aligned}y^1 &= 0.3 \\y^2 &= 0.7\end{aligned}$$

Complicated? Yes,
and a little slow, too.



Recurrent Neural Network

Complicated? Yes, and a little slow, too.



But the animation is as colorful as the story and emotions do eventually reach the desired level.

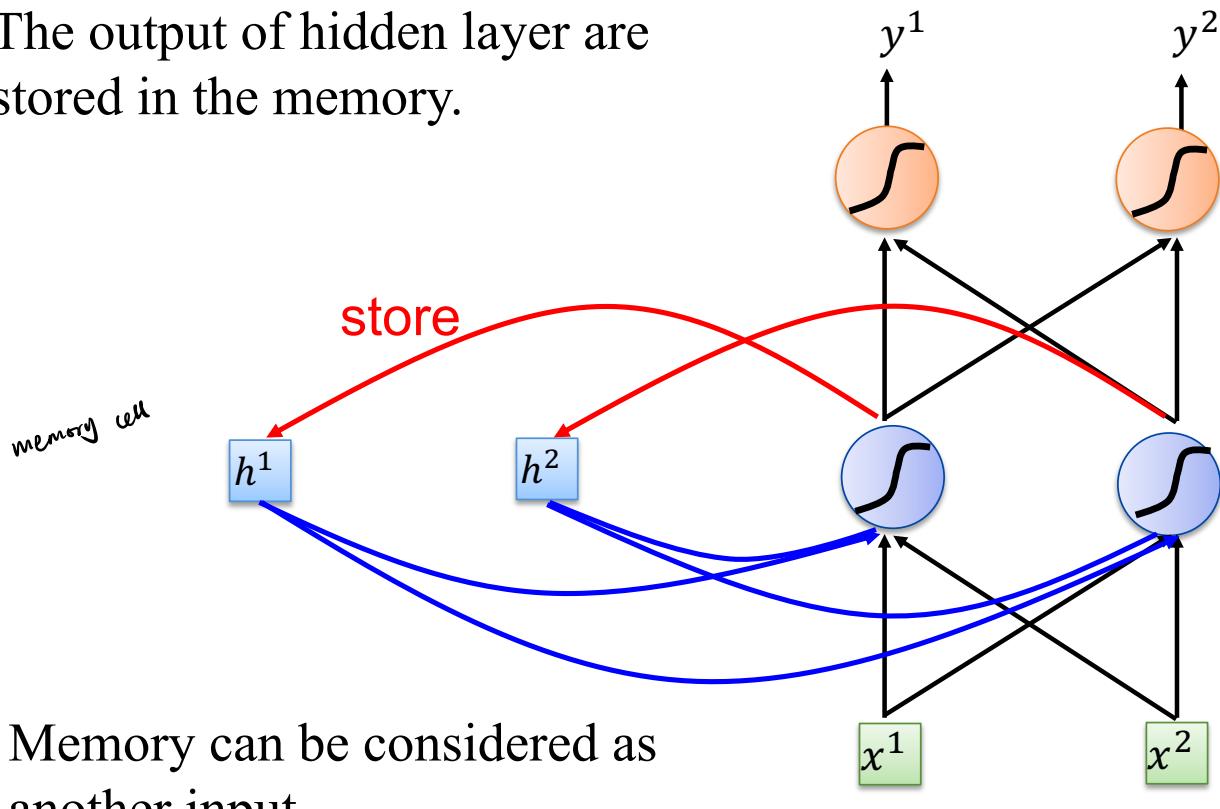
Negative OR Positive



Need
Memory

Recurrent Neural Network

The output of hidden layer are stored in the memory.

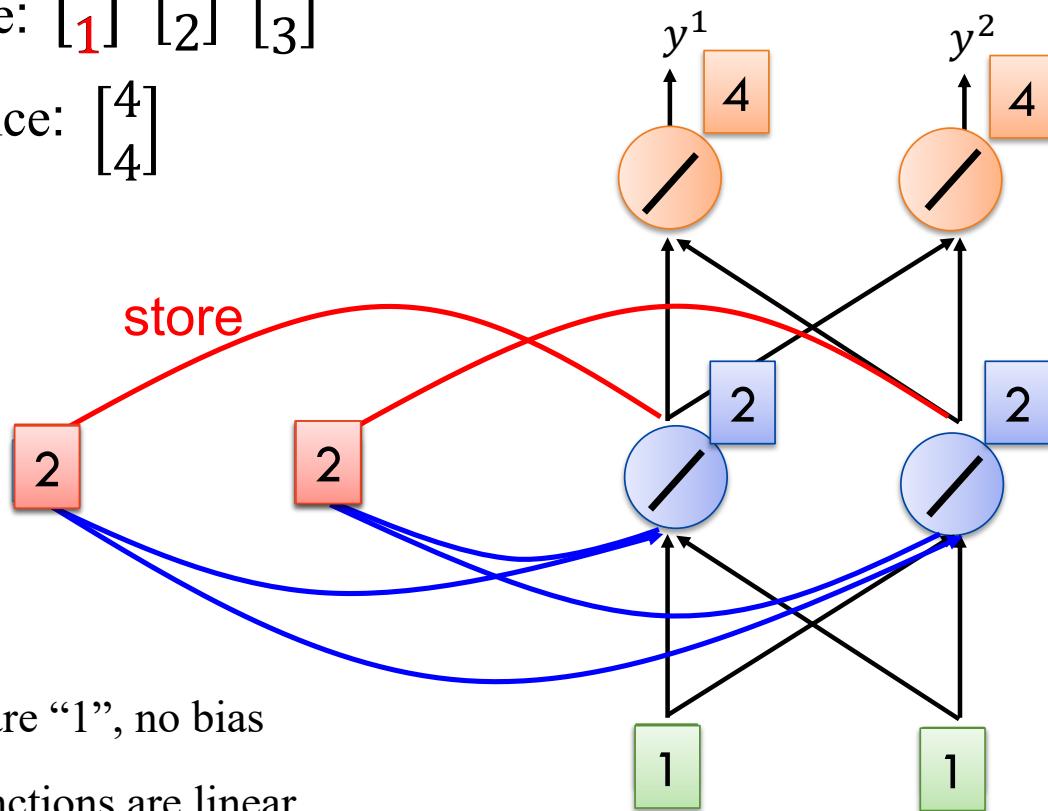


Memory can be considered as another input.

Recurrent Neural Network

Input sequence: $\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$

Output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix}$



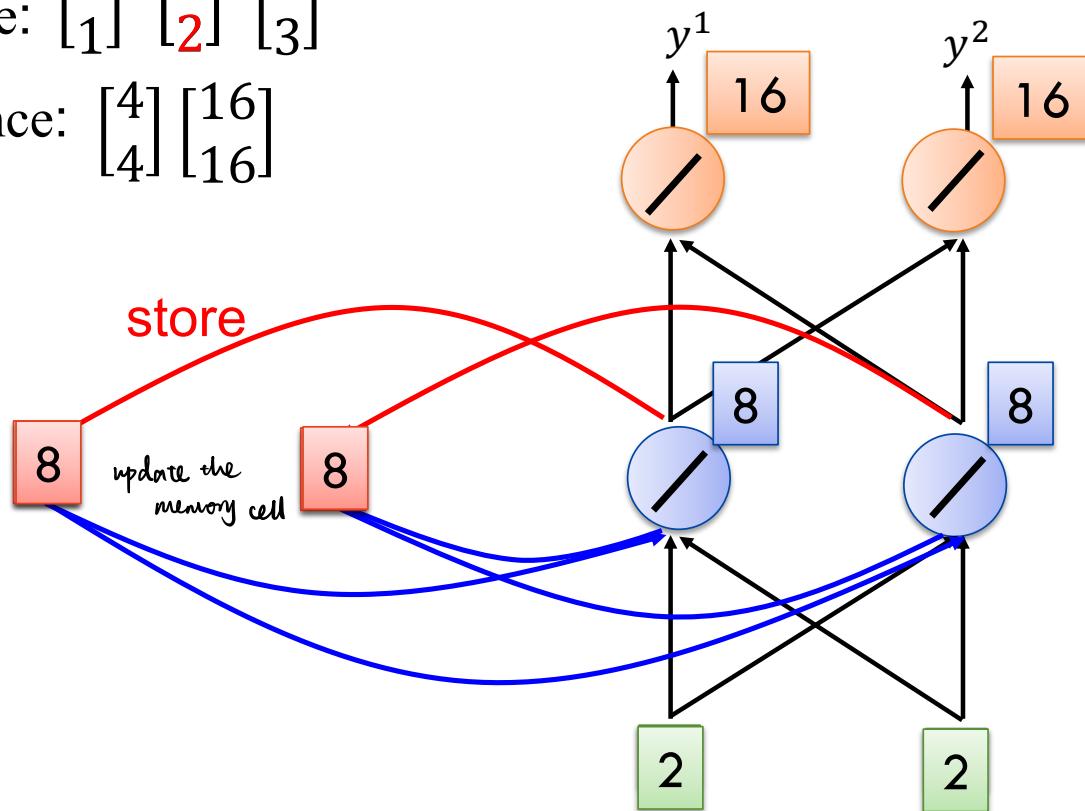
All the weights are “1”, no bias

All activation functions are linear

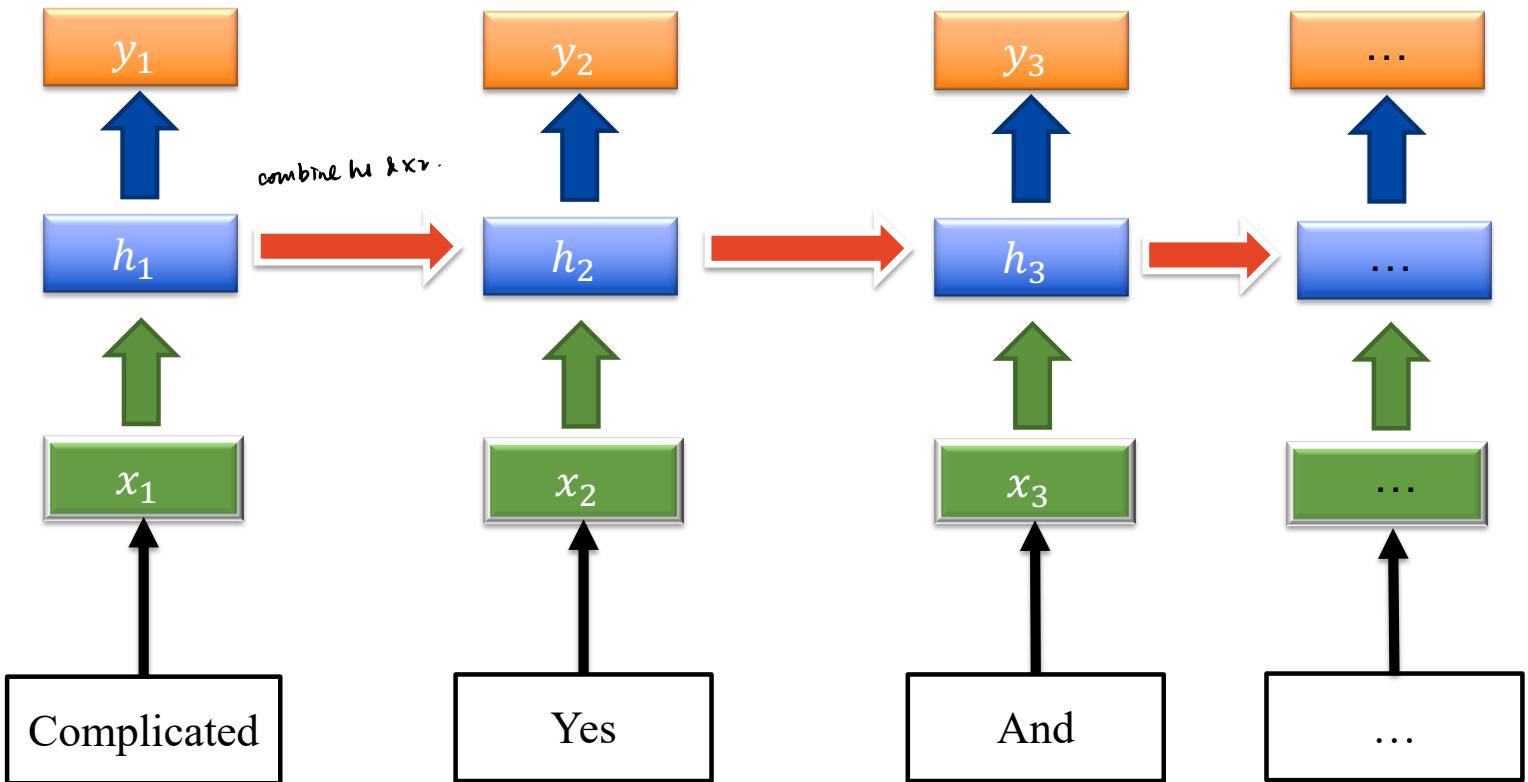
Recurrent Neural Network

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} \begin{bmatrix} 3 \\ 3 \end{bmatrix}$

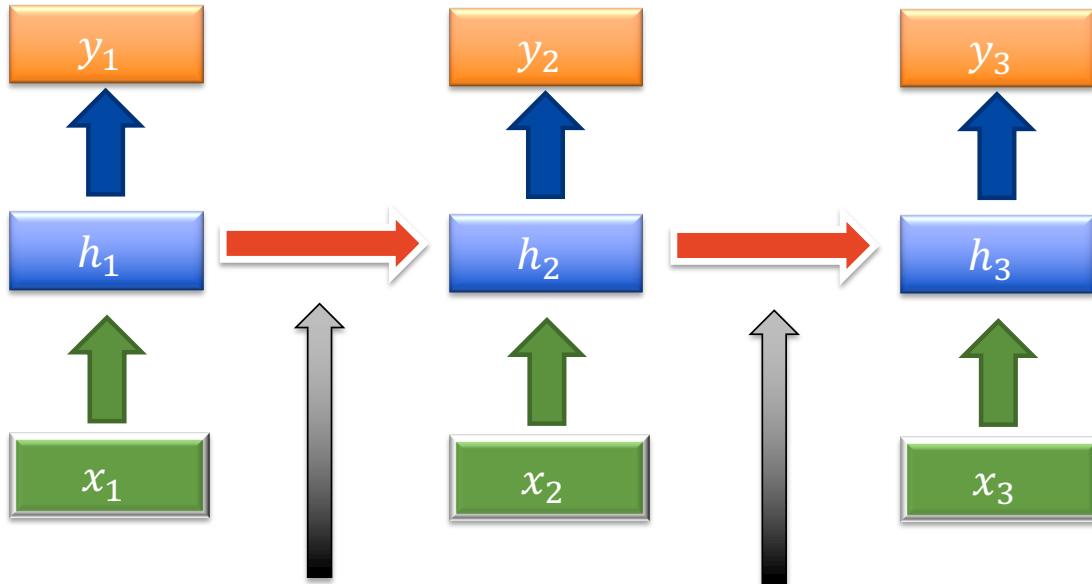
Output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix} \begin{bmatrix} 16 \\ 16 \end{bmatrix}$



Recurrent Neural Network



Recurrent Neural Network



Memory can be considered as another input.

Recurrent Neural Network

Formally

x_t is the input

h_t is the hidden state

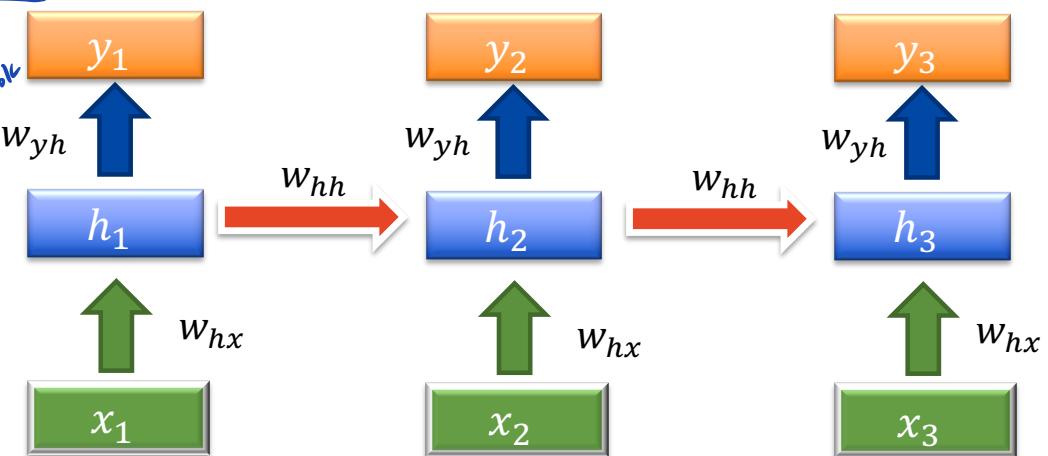
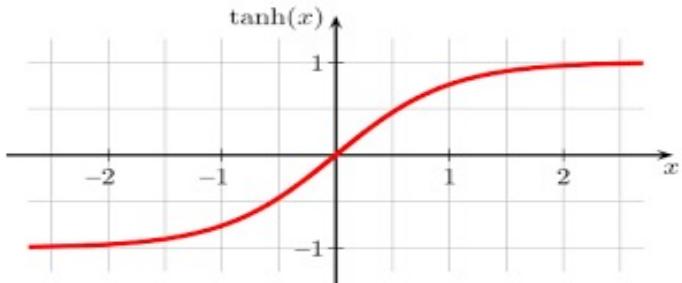
y_t is the output

$h_0 = \vec{0}$ ← initialize memory value.

$h_t = \tanh(w_{hh}h_{t-1} + w_{hx}x_t + b_h)$

$y_t = w_{yh}h_t$

*whether to use
 y_1, \dots, y_t depends on the task*



Recurrent Neural Network

- Vanishing gradient problem

- $h_t = \tanh(w_{hh}h_{t-1} + w_{hx}x_t + b_h)$

- $y_t = w_{yh}h_t$

- For every step, its loss is E_n

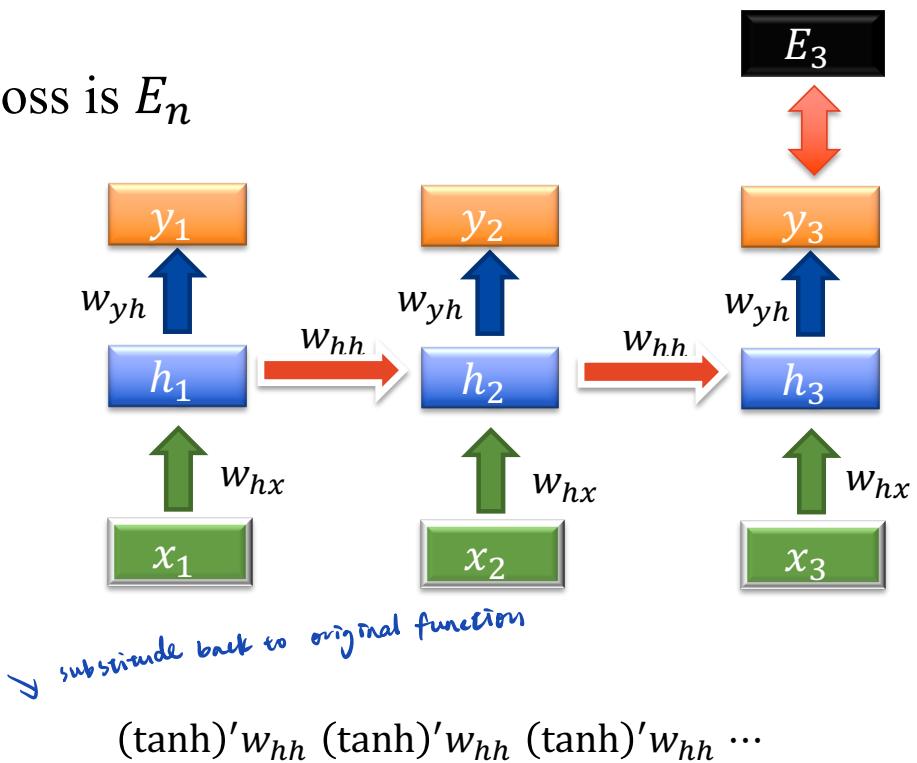
update w_{hh}

$$\frac{\partial E_3}{\partial w_{hh}} = \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial h_3}{\partial w_{hh}}$$

$$h_3 = \tanh(w_{hh}h_2 + w_{hx}x_3 + b_h)$$

$$\frac{\partial h_3}{\partial w_{hh}} = (\tanh)' \left(h_2 + w_{hh} \frac{\partial h_2}{\partial w_{hh}} \right)$$

$$\frac{\partial h_2}{\partial w_{hh}} = (\tanh)' \left(h_1 + w_{hh} \frac{\partial h_1}{\partial w_{hh}} \right)$$



Recurrent Neural Network

- Vanishing gradient problem

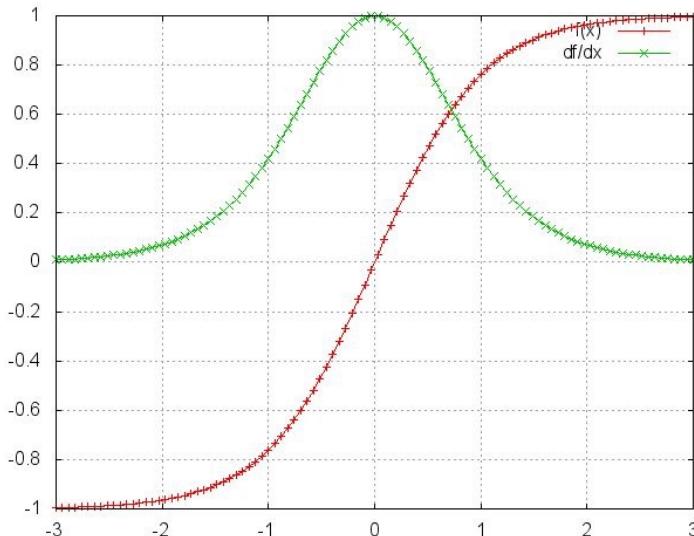
- $h_t = \tanh(w_{hh}h_{t-1} + w_{hx}x_t + b_h)$

- $y_t = w_{yh}h_t$

- For every step, its loss is E_n

$$\frac{\partial E_3}{\partial w_{hh}} = \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial h_3}{\partial w_{hh}}$$

$$(\tanh)'w_{hh}(\tanh)'w_{hh}(\tanh)'w_{hh} \dots$$



$$(\tanh)' \leq 1$$

If $0 < w_{hh} < 1$, Vanishing Gradients

$$(\tanh)'w_{hh}(\tanh)'w_{hh}(\tanh)'w_{hh} \dots \rightarrow 0$$

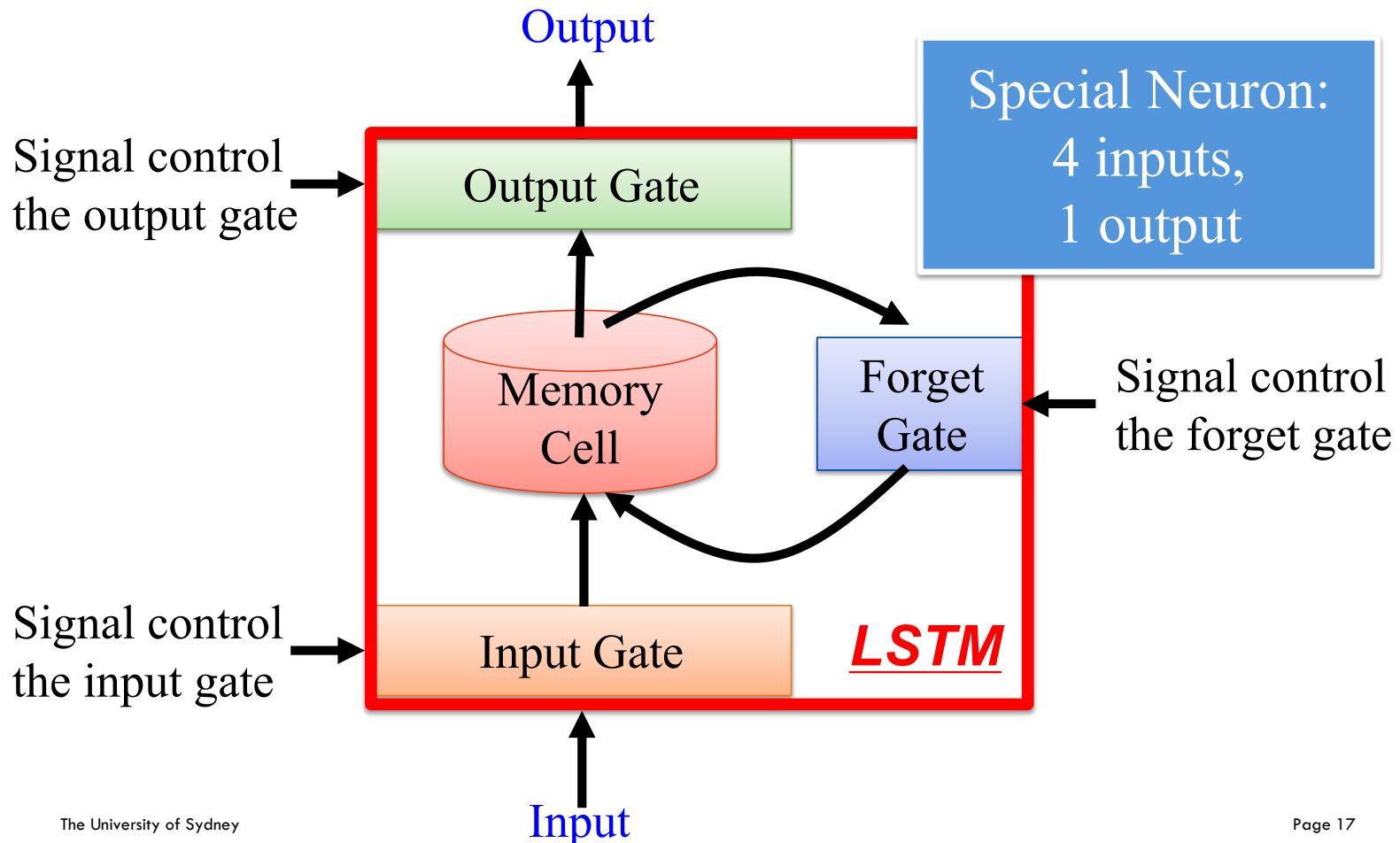
If w_{hh} is larger, Exploding Gradients

$$(\tanh)'w_{hh}(\tanh)'w_{hh}(\tanh)'w_{hh} \dots \rightarrow \infty$$

$$0.9^{1000} \approx 0$$
$$1.01^{1000} \approx 21000$$

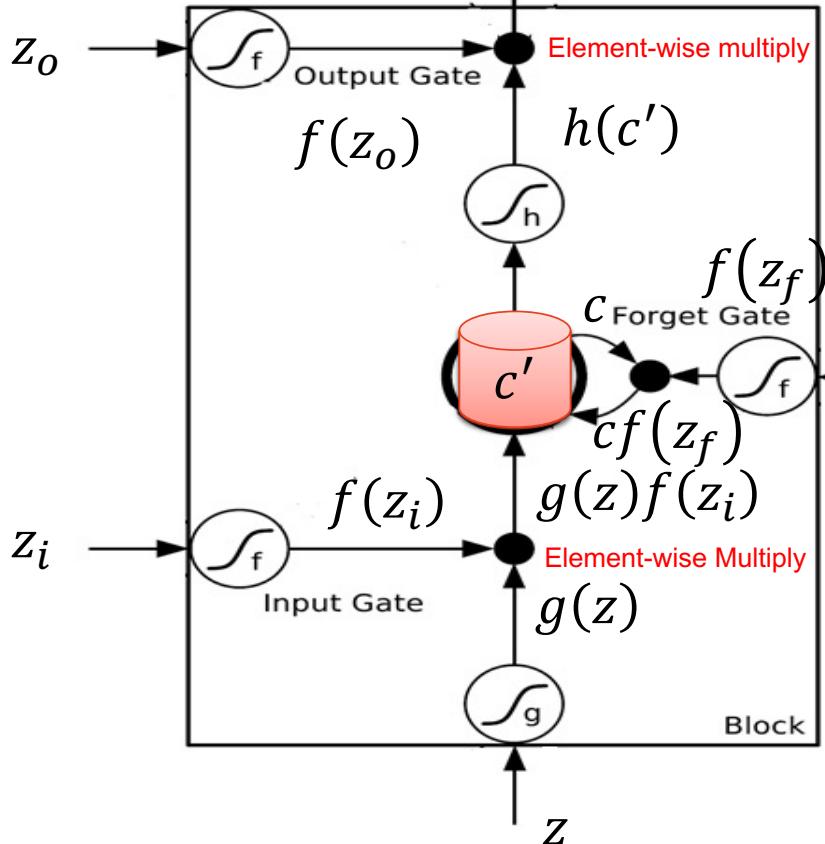
Long Short-Term Memory

Long Short-Term Memory



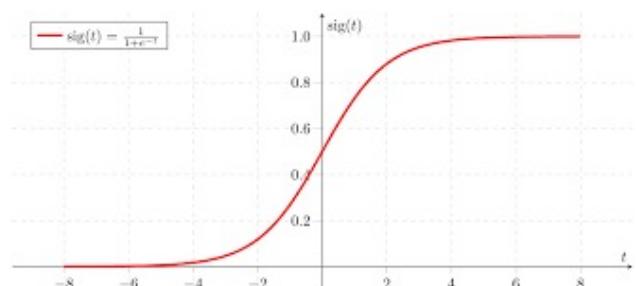
Long Short-Term Memory

$$h = h(c')f(z_o)$$



f is usually a sigmoid function
Value From 0 to 1
Mimic open and close gate
input gate.

$$c' = g(z)f(z_i) + \underset{\substack{\text{new} \\ \text{info}}}{\textcolor{blue}{cf(z_f)}}$$



Long Short-Term Memory

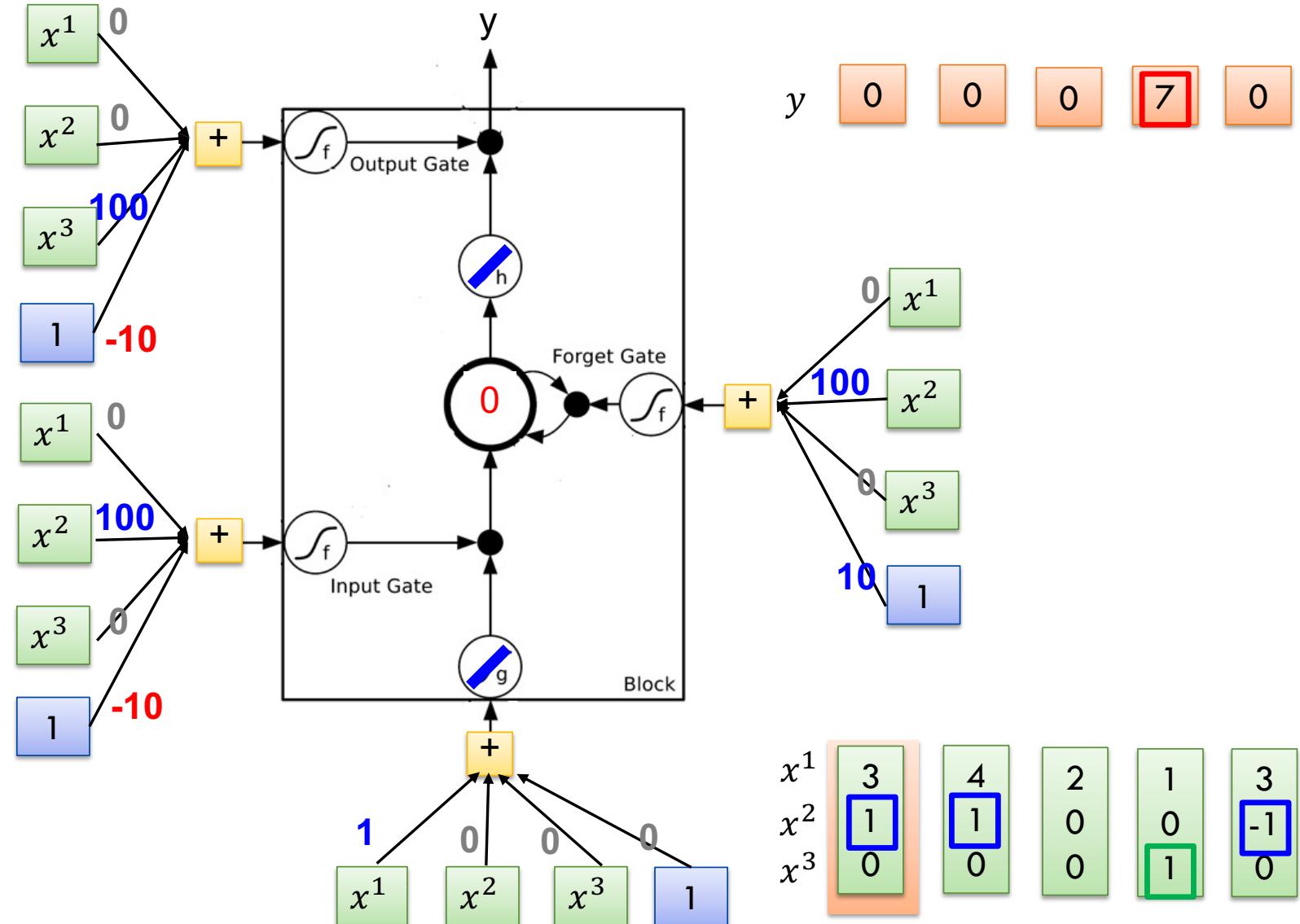
C	0	0	3	3	7	7	7	0	6
x^1	1	3	2	4	2	1	3	6	1
x^2	0	1	0	1	0	0	-1	1	0
x^3	0	0	0	0	0	1	0	0	1

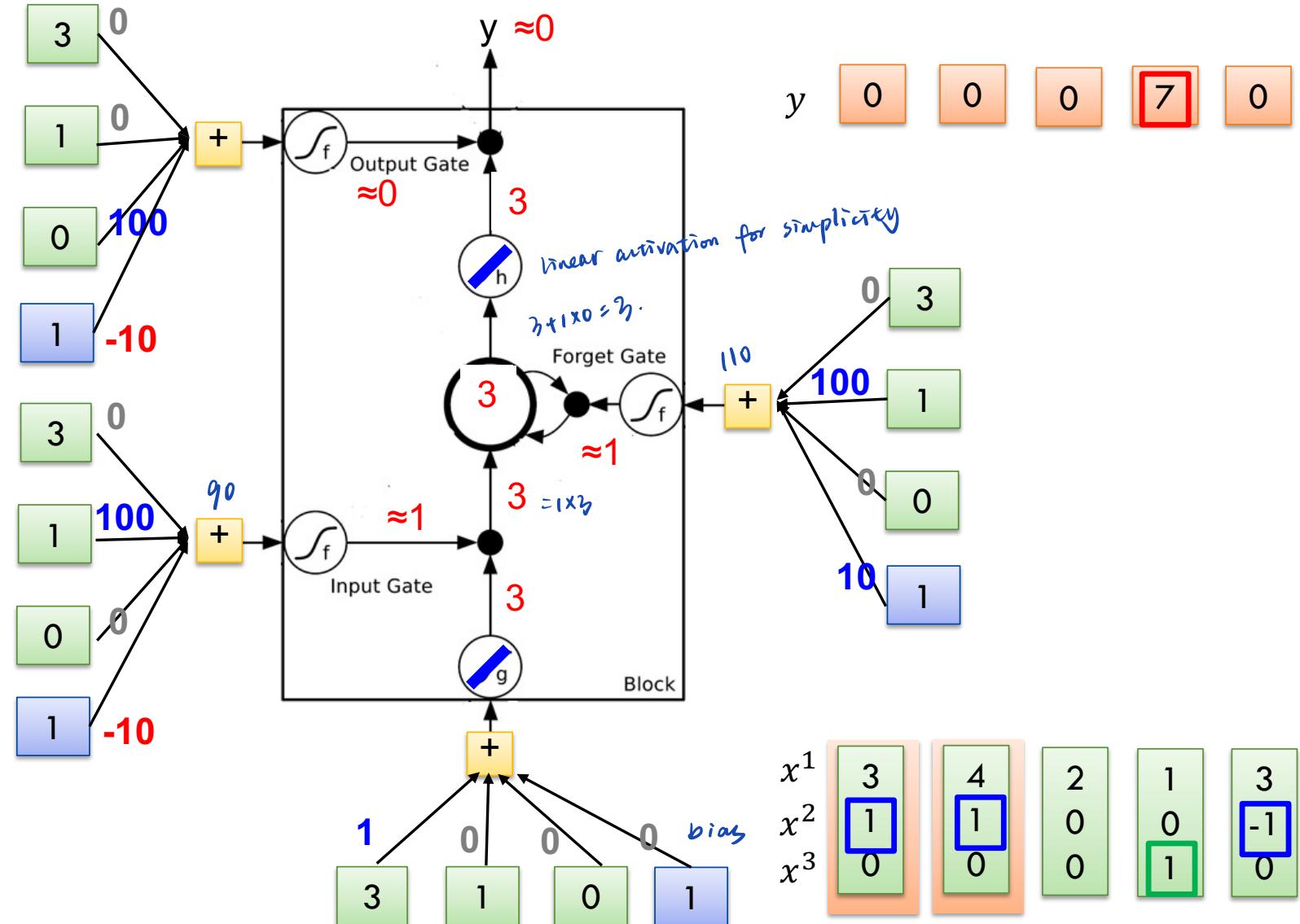
y	0	0	0	0	0	7	0	0	6
-----	---	---	---	---	---	---	---	---	---

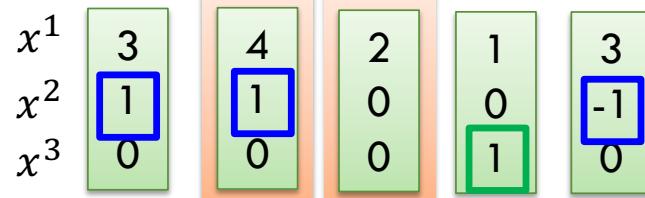
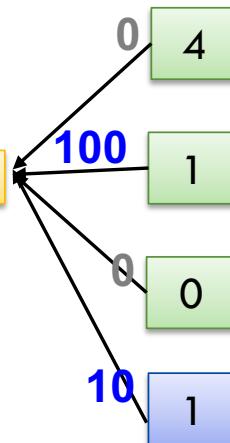
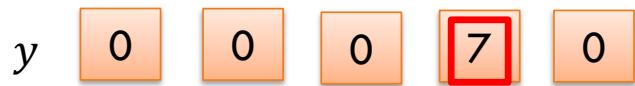
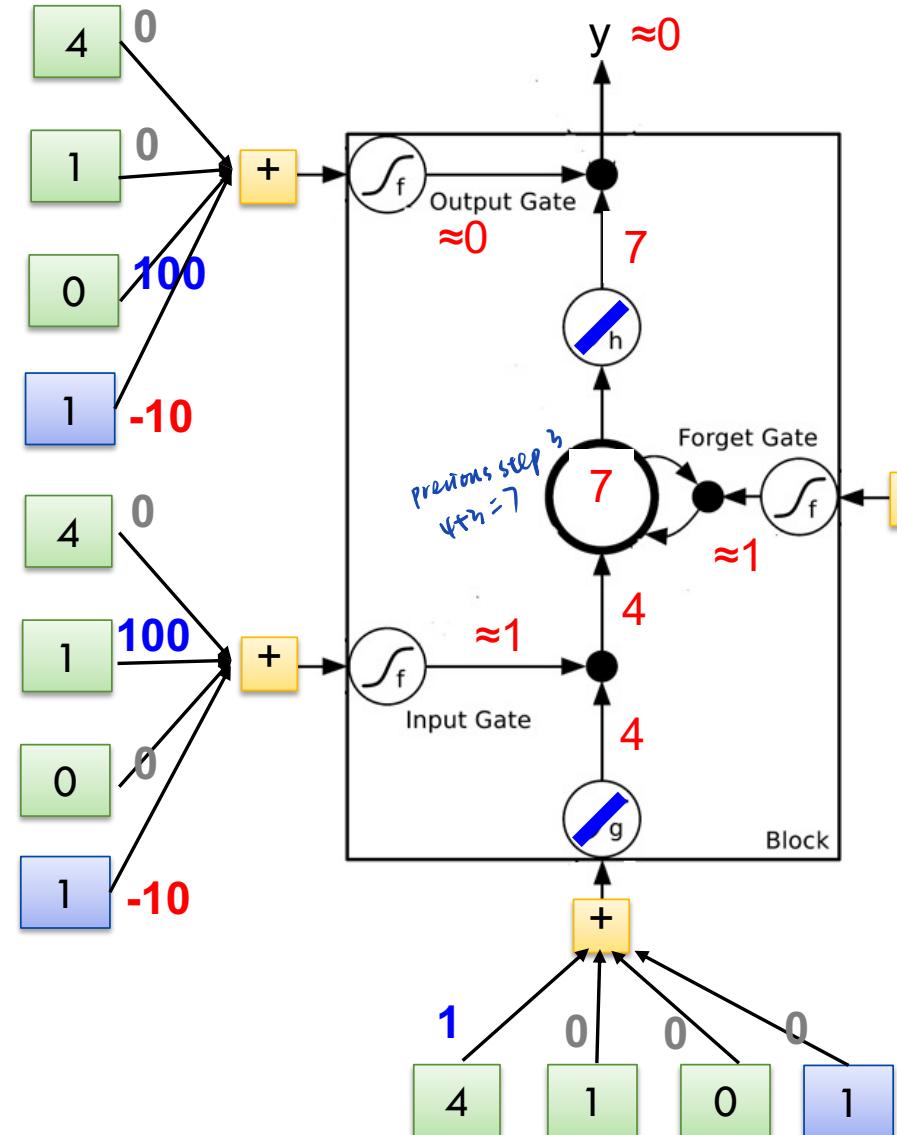
When $x^2 = 1$, add the numbers of x_1 into the memory

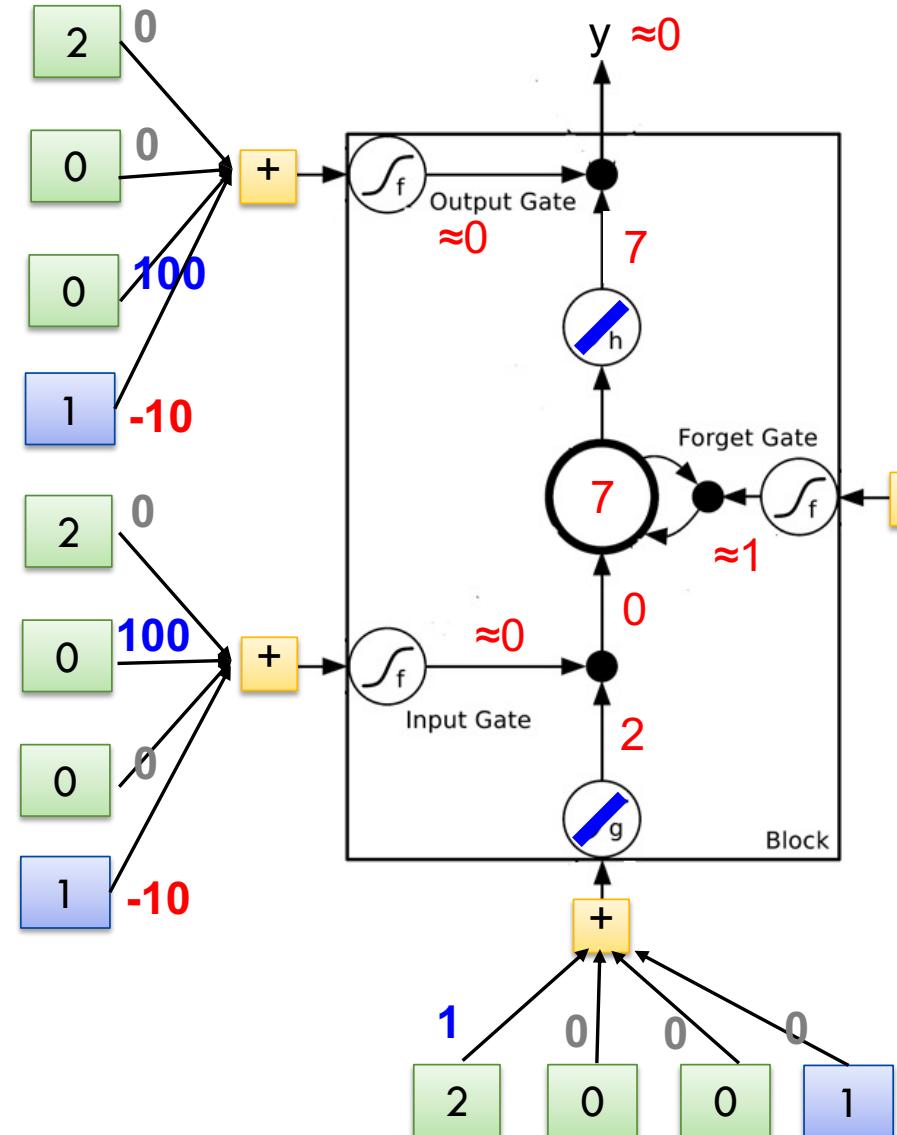
When $x^2 = -1$, reset the memory

When $x^3 = 1$, output the number in the memory.

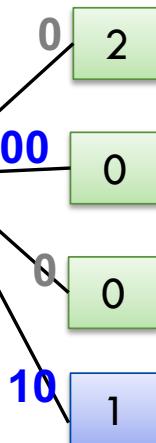




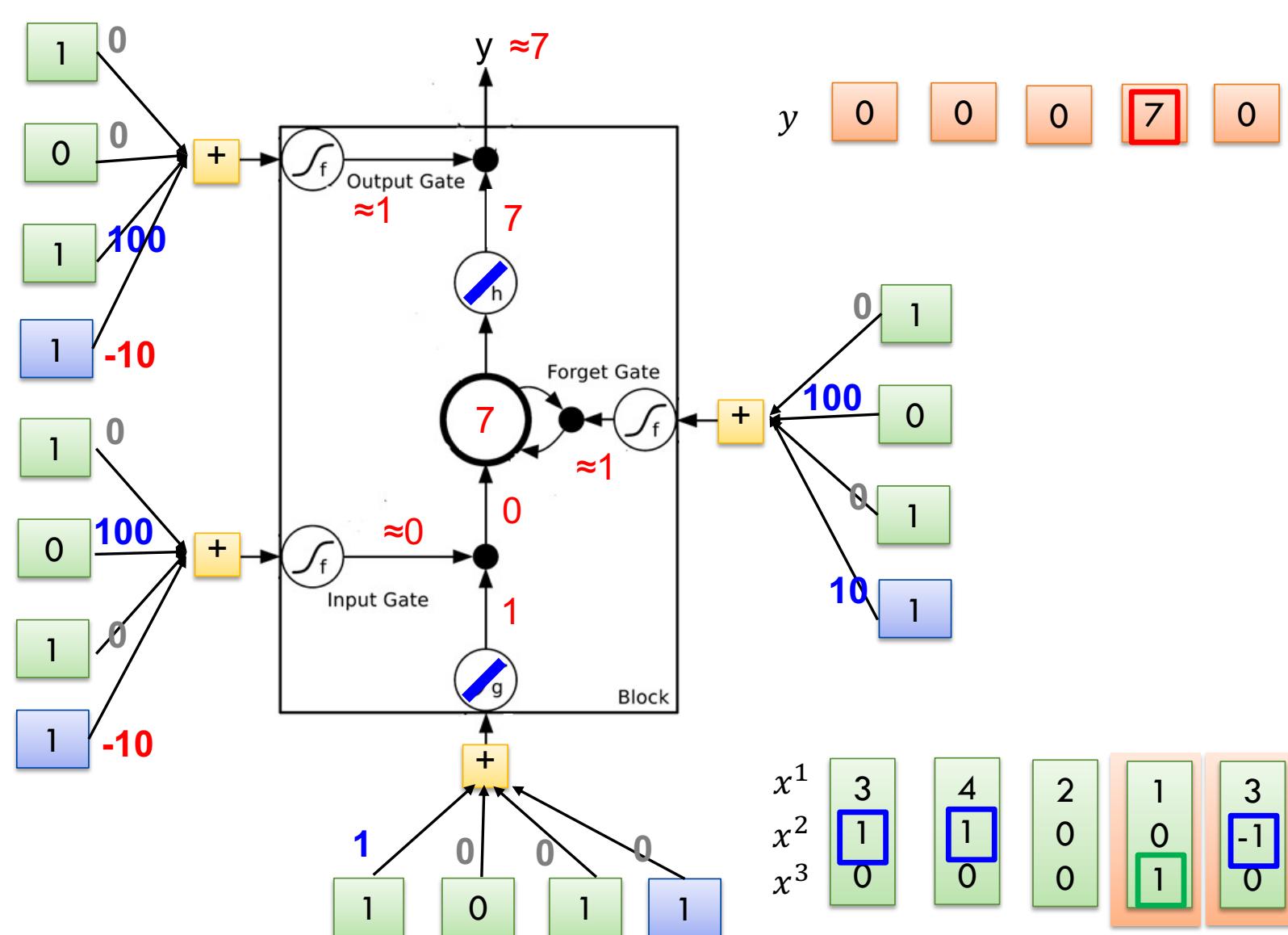


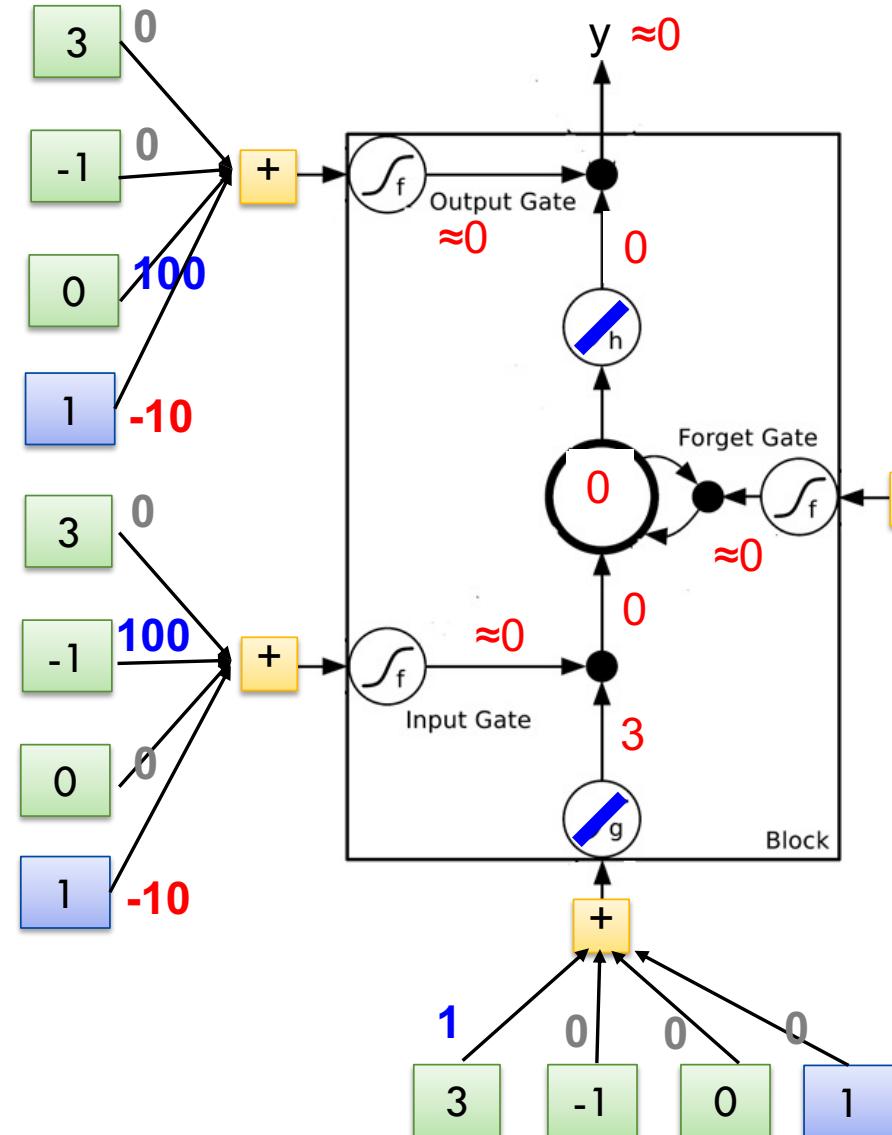


y 0 0 0 7 0

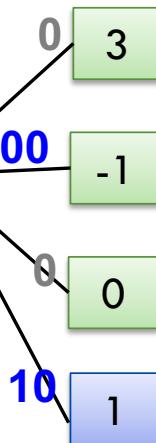


x^1 3 4 2 1 3
 x^2 1 1 0 0 -1
 x^3 0 0 0 1 0



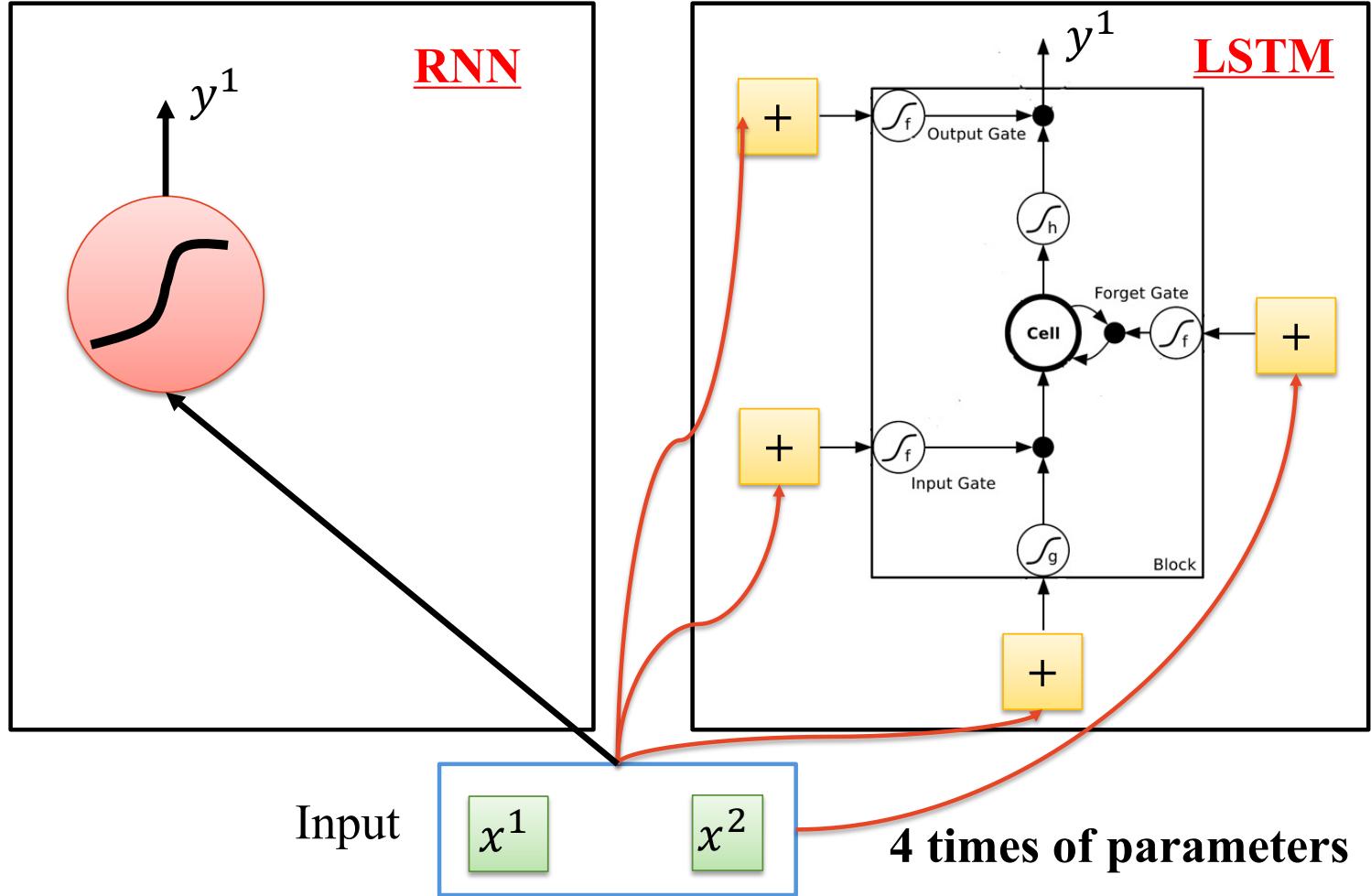


y 0 0 0 7 0



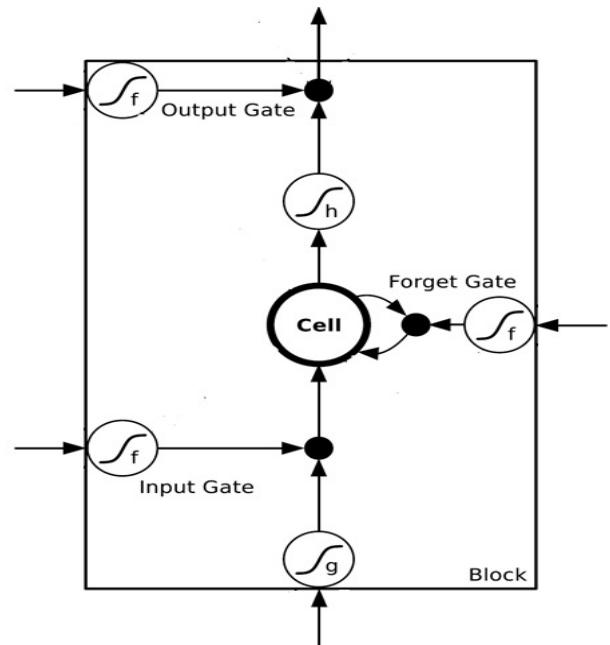
x^1 3 4 2 1 3
 x^2 1 1 0 0 -1
 x^3 0 0 1 1 0

Long Short-Term Memory



Long Short-Term Memory

- Forget gate *hidden state from last time step*
 - $f_t = \sigma(w_{fh}h_{t-1} + w_{fx}x_t + b_f)$
- Input gate
 - $i_t = \sigma(w_{ih}h_{t-1} + w_{ix}x_t + b_i)$
- Output gate
 - $o_t = \sigma(w_{oh}h_{t-1} + w_{ox}x_t + b_o)$



Long Short-Term Memory

□ Candidate cell state

$$\bullet \tilde{c}_t = \tanh(w_{hh}h_{t-1} + w_{hx}x_t + b_h)$$

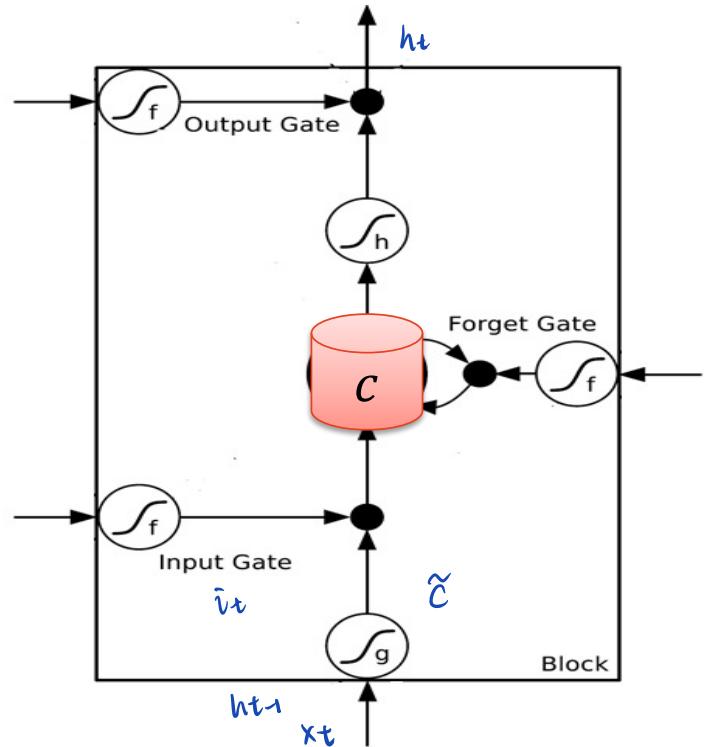
□ Update cell state

$$\bullet c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

□ Output state

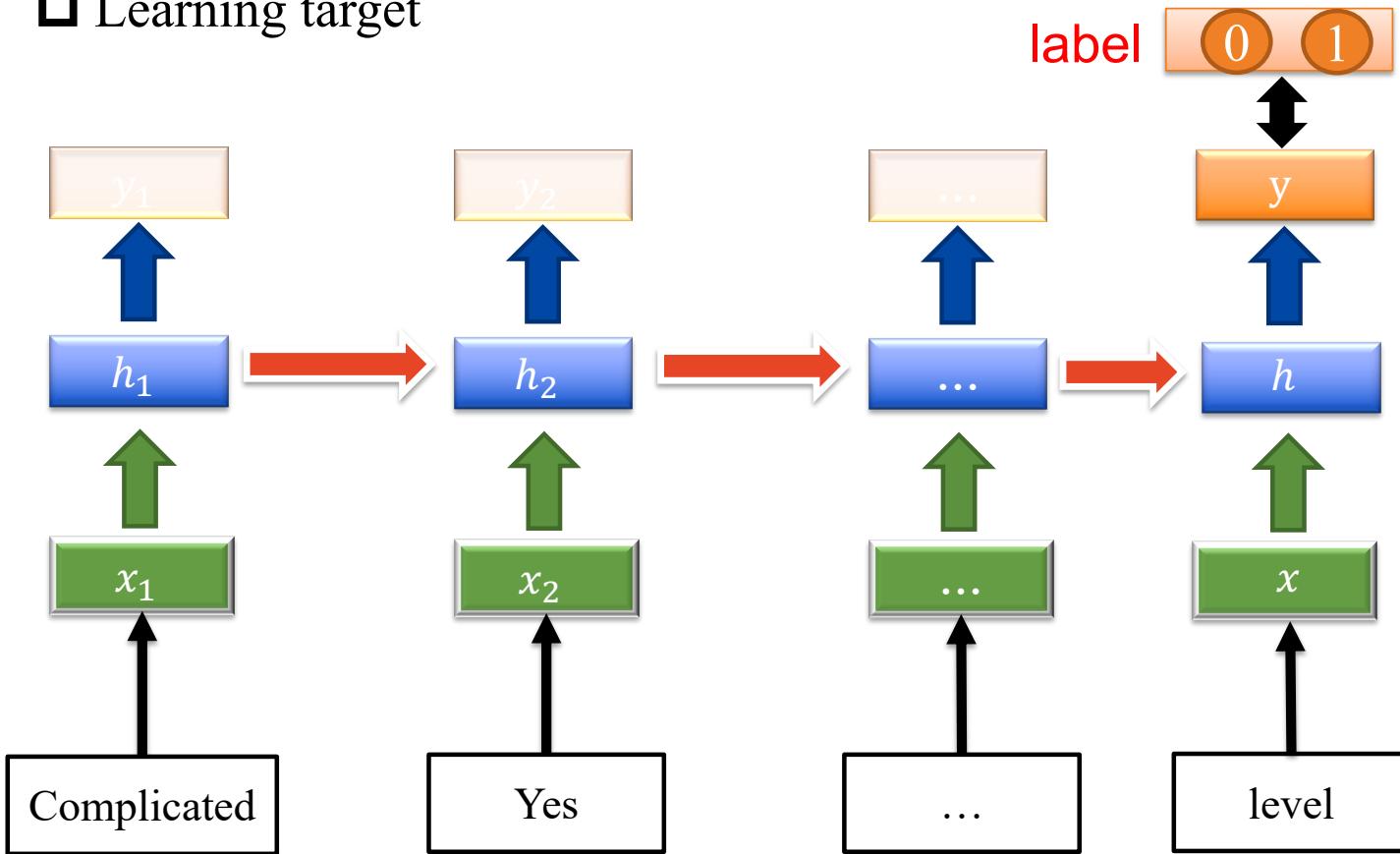
$$\bullet h_t = o_t * \tanh(c_t)$$

↑
memory value



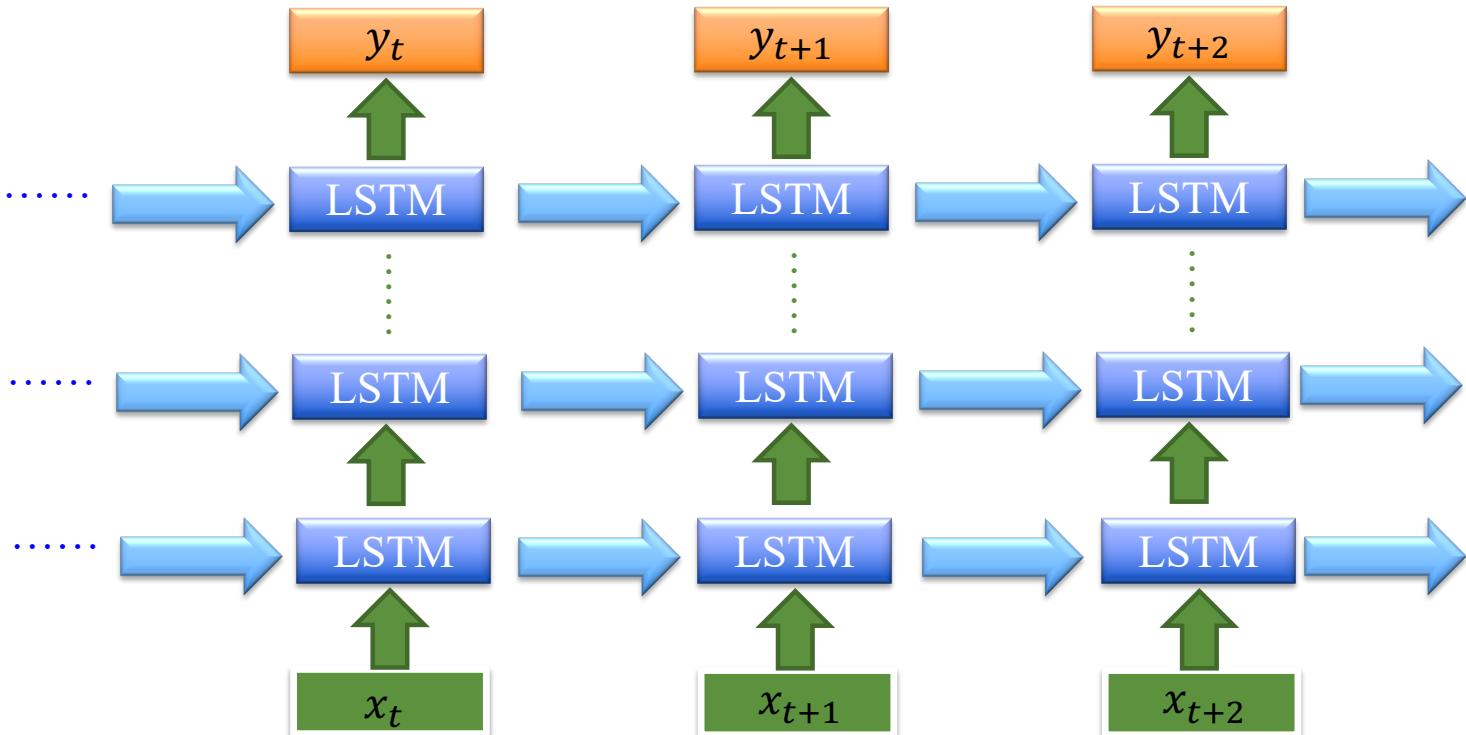
Long Short-Term Memory

- Learning target



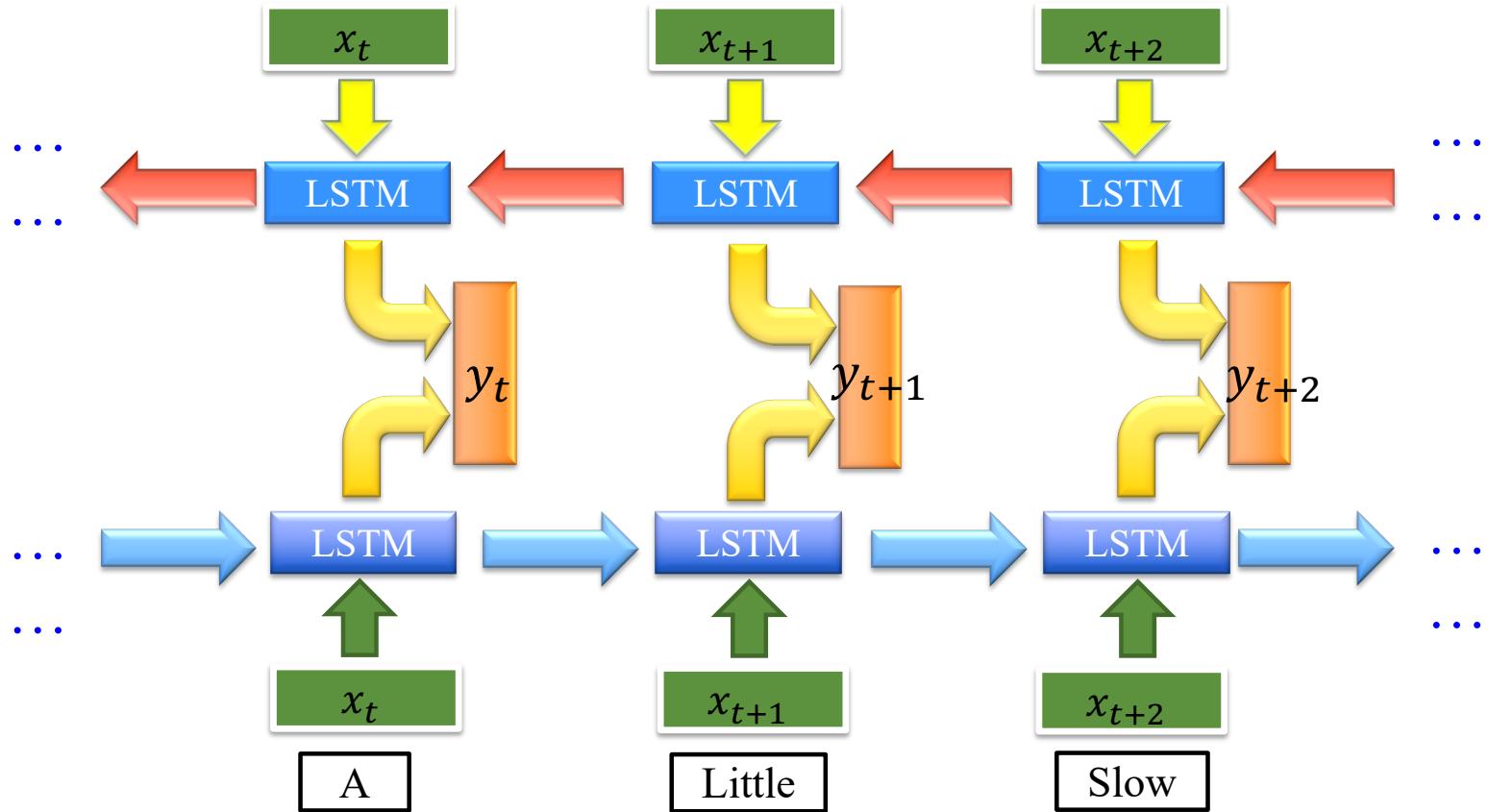
Long Short-Term Memory

□ Stacked LSTM



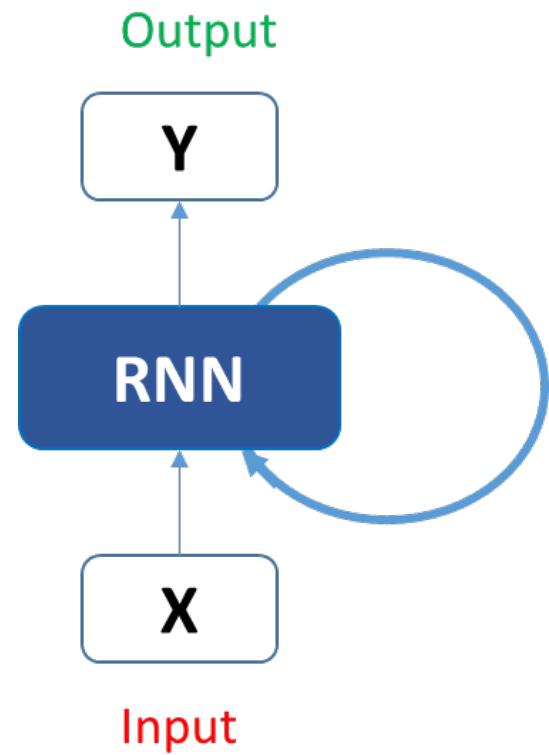
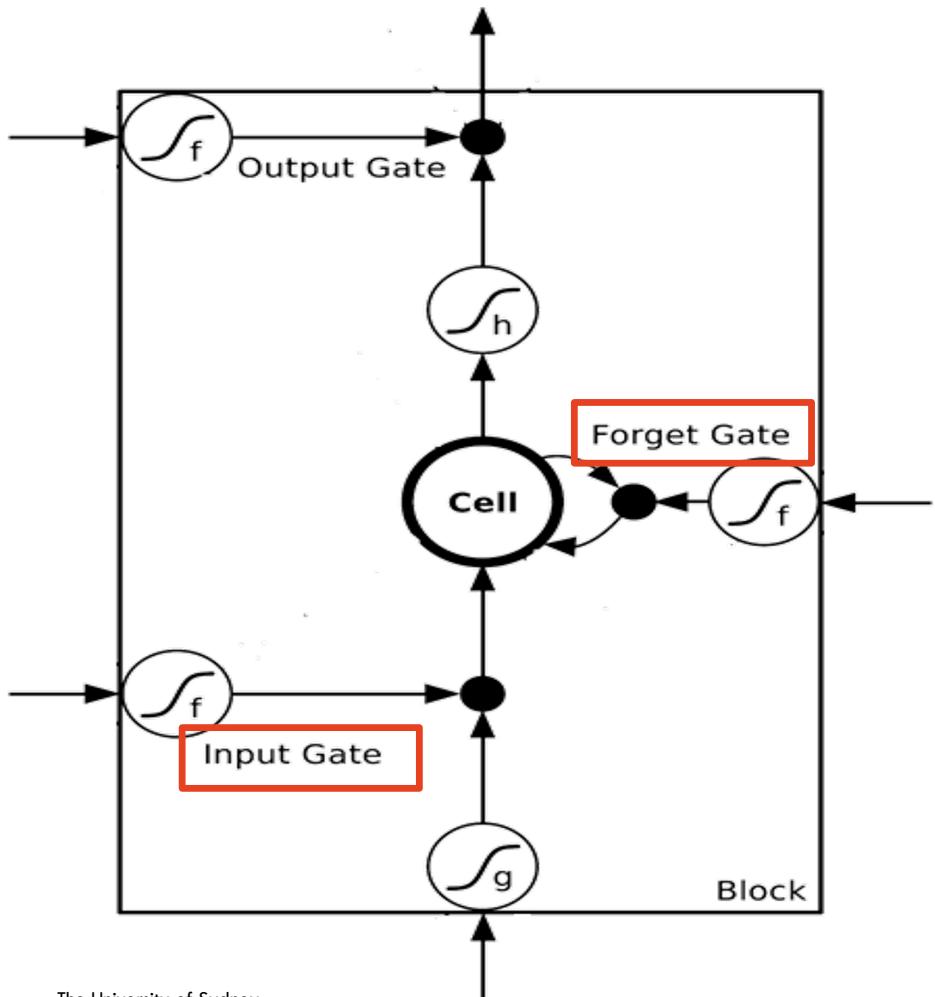
Long Short-Term Memory

□ Bidirectional RNN



LSTM ⚡ maintain more parameters

LSTM V.S. RNN



Prevent vanishing gradient

□ RNN

$$\bullet h_t = \tanh(w_{hh}h_{t-1} + w_{hx}x_t + b_h)$$

□ LSTM

$$\bullet c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

$$\frac{\partial e}{\partial c_t} = \frac{\partial f_t}{\partial c_{t-1}} c_{t-1} + f_t + \frac{\partial i_t}{\partial c_{t-1}} * \tilde{c}_t + \frac{\partial \tilde{c}_t}{\partial c_{t-1}}$$

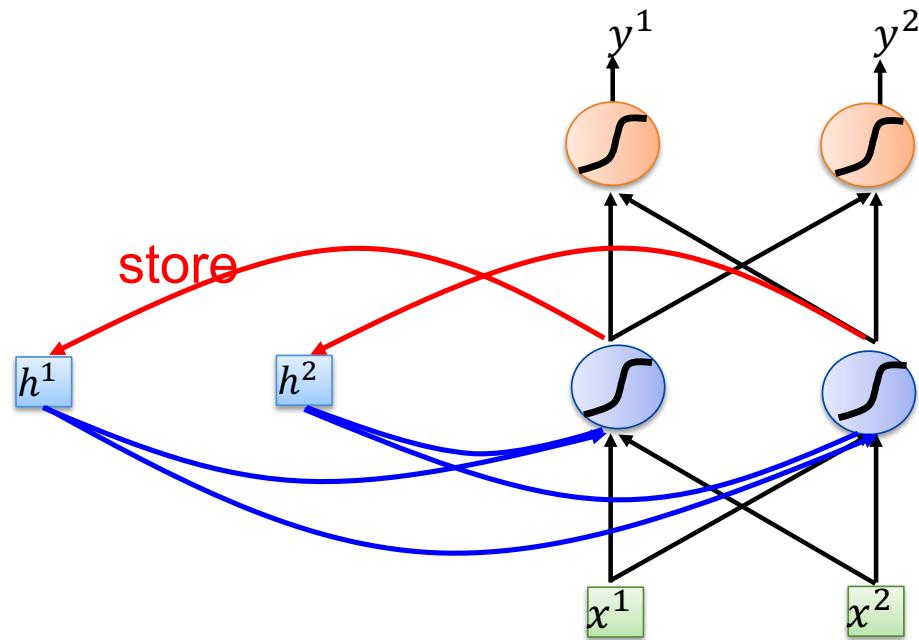
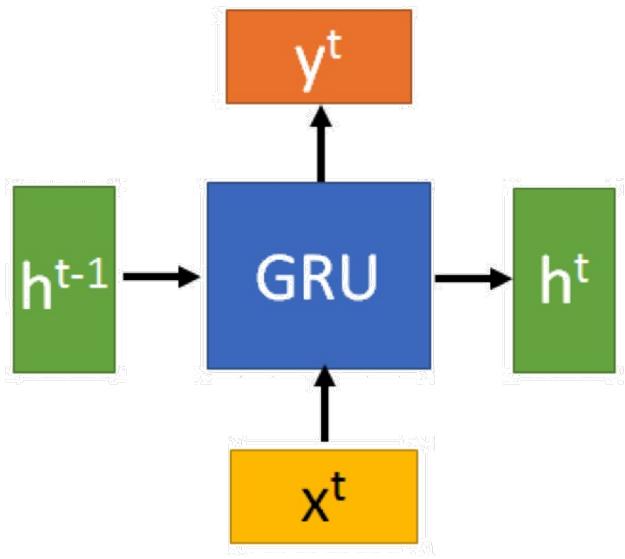
① small chance to have all terms being small values

① prevent vanishing gradient problem, but may have gradient explosion problem

LSTMs solve the problem using a unique additive gradient structure that includes ① direct access to the forget gate's activations, enabling the network to encourage desired behaviour from the error gradient using frequent gates update on every time step of the learning process.

Gate Recurrent Unit

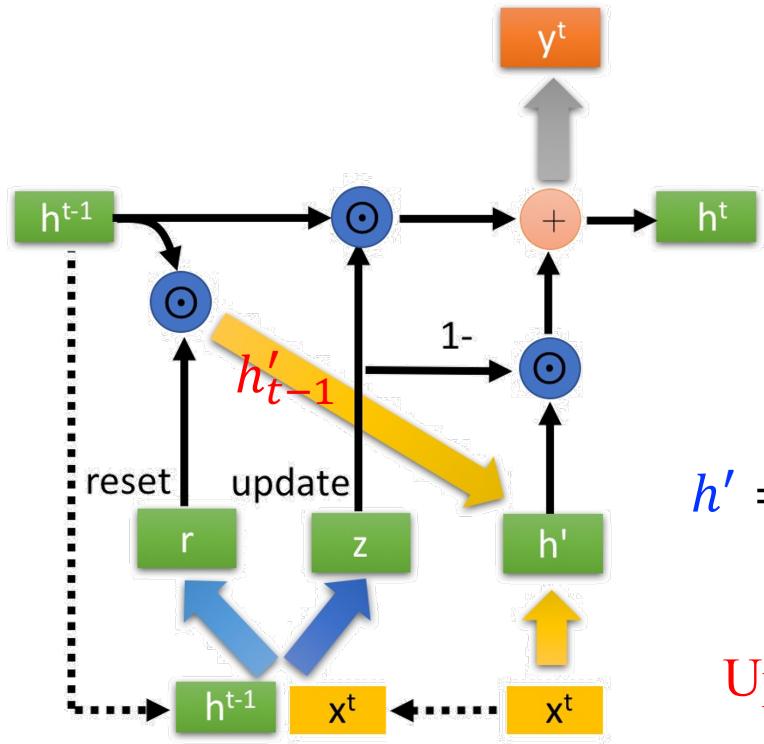
Gate Recurrent Unit



Vanilla RNN

Gate Recurrent Unit

parameter efficient



Reset Gate

$$r_t = \sigma(w_{rh}h_{t-1} + w_{rx}x_t + b_r)$$

$$h'_{t-1} = r_t * h_{t-1}$$

$$h' = \tanh(w_{hh}h'_{t-1} + w_{hx}x_t + b_h)$$

Update Gate (forget gate + input gate)

$$z_t = \sigma(w_{zh}h_{t-1} + w_{zx}x_t + b_z)$$

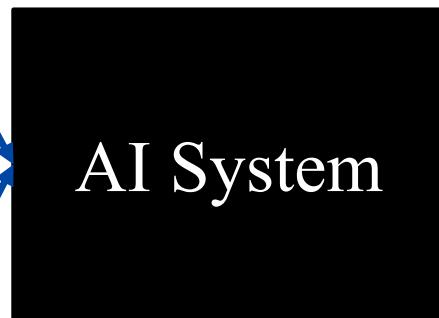


$$h_t = z_t * h_{t-1} + (1 - z_t) * h'_{\substack{\text{new} \\ \text{memory}}}$$

Application

- ❑ Visual Question Answer (VQA)
- ❑ Reading Comprehension
- ❑ Sequence To Sequence (Seq2seq)
 - ❑ Language Translation
 - ❑ Chat Robot

VQA



What's the
mustache made of ?

Banana

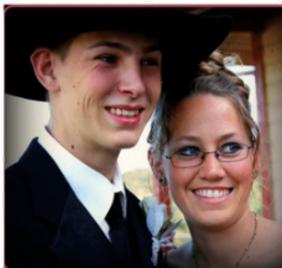
Image credit to <http://www.visualqa.org/challenge.html>

VQA

Who is wearing glasses?

man

woman



Is the umbrella upside down?

yes

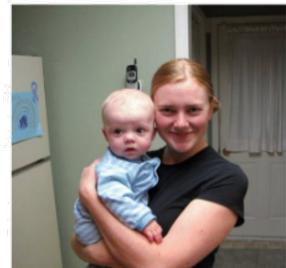
no



Where is the child sitting?

fridge

arms



How many children are in the bed?

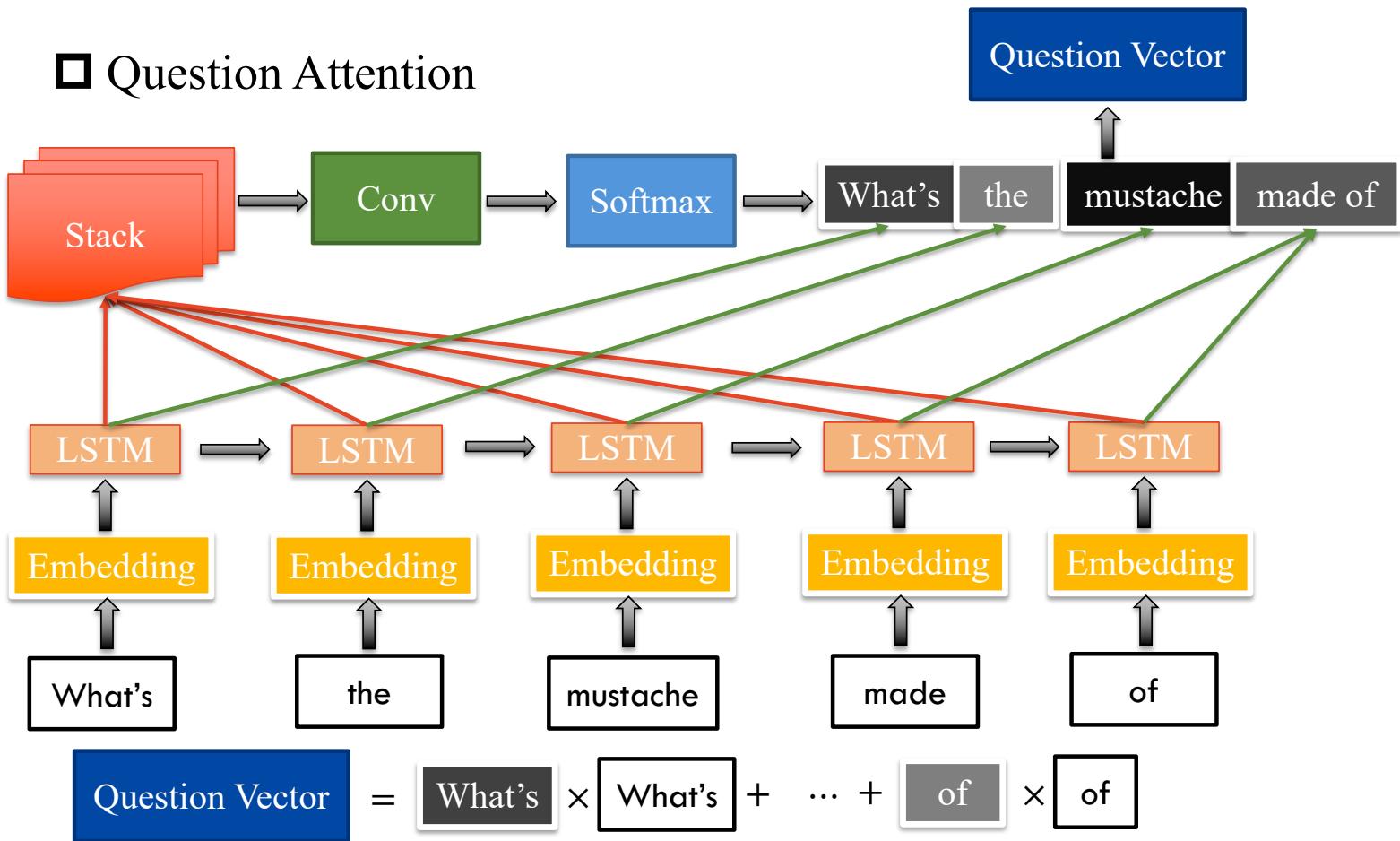
2

1



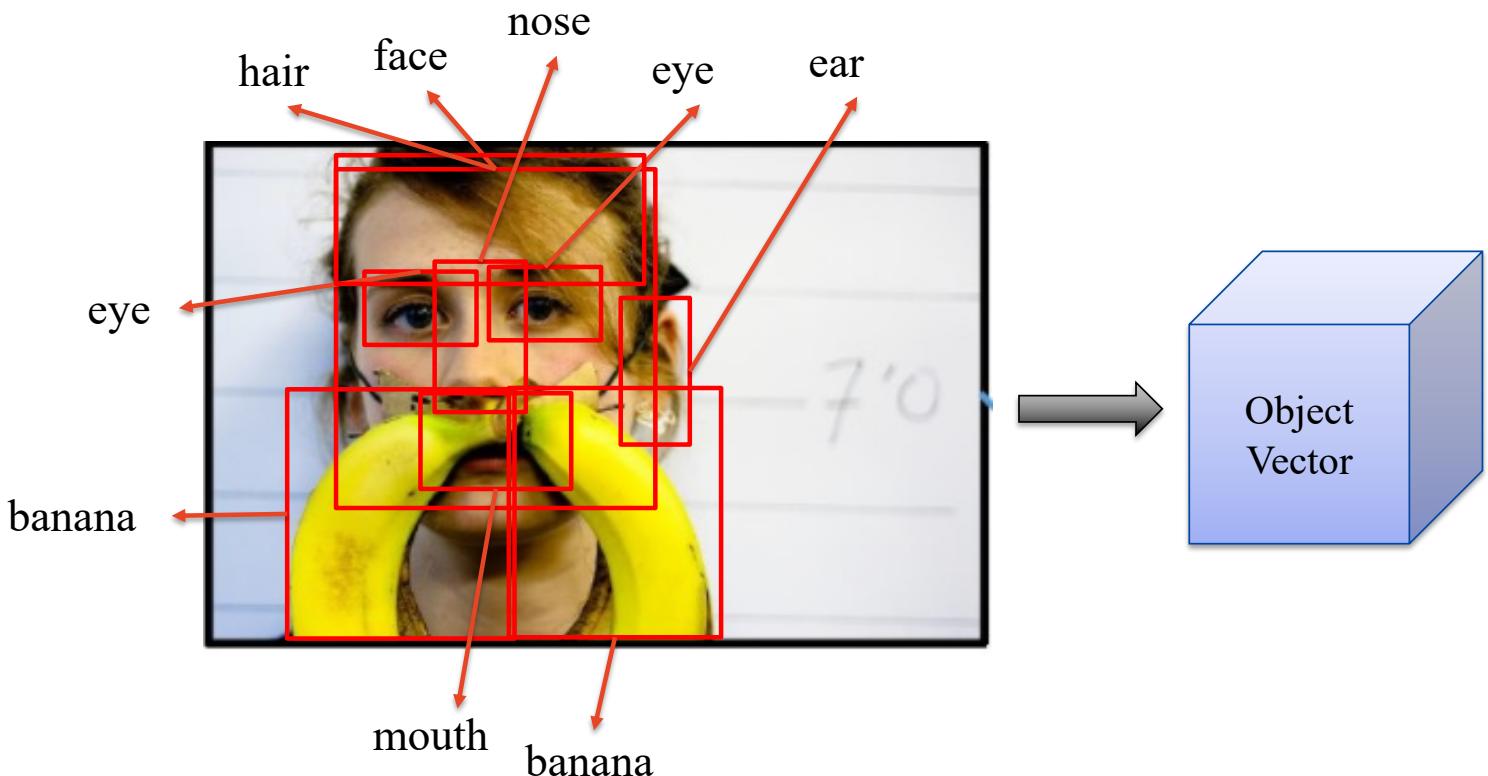
VQA

□ Question Attention



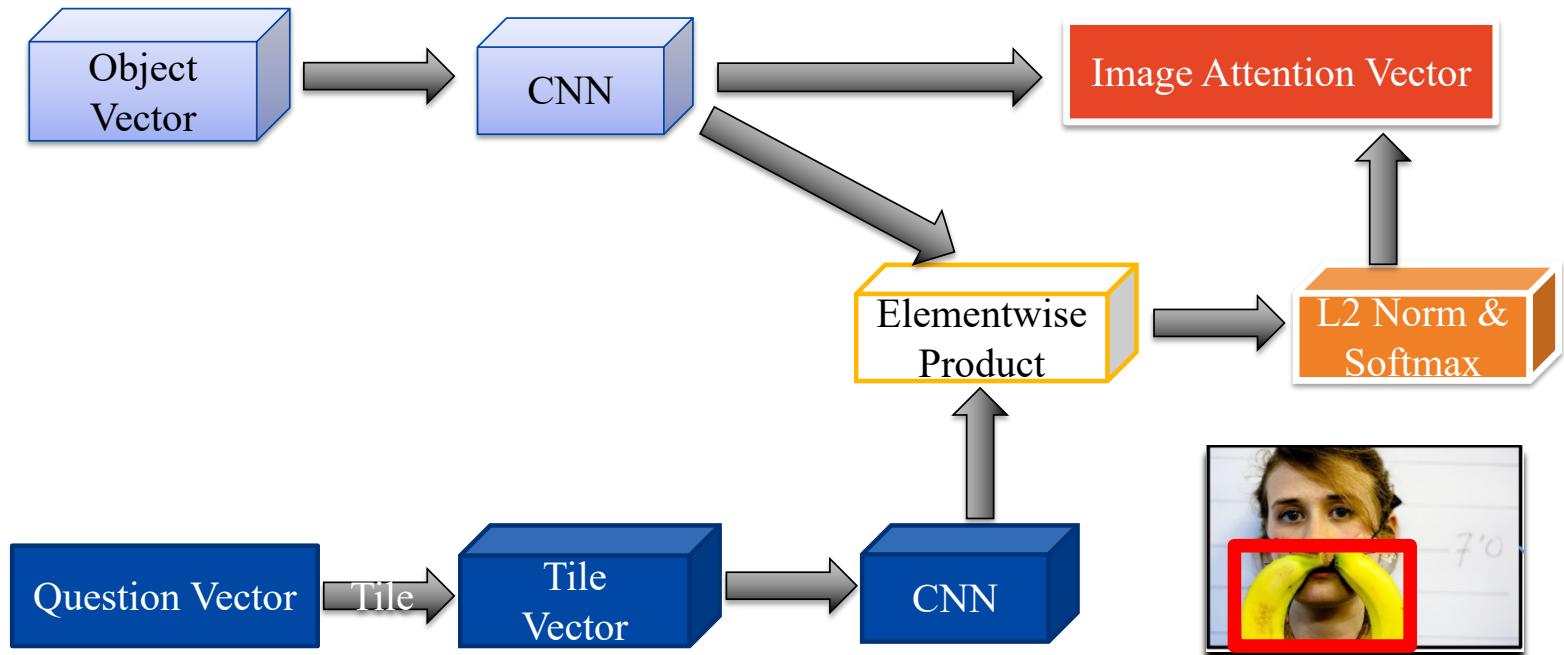
VQA

□ Object Detection

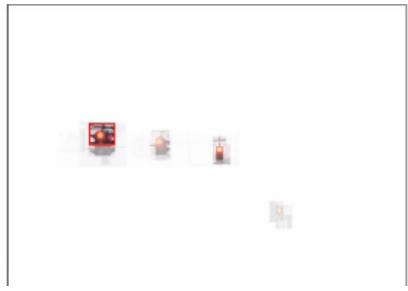


VQA

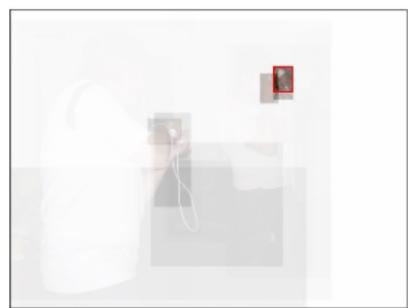
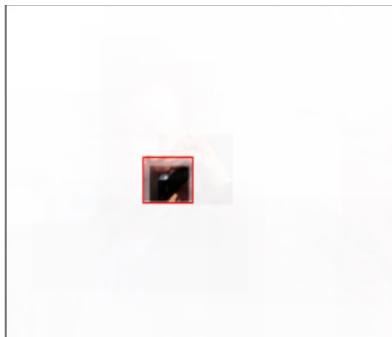
□ Image Attention



VQA



What color is illuminated on the traffic light?

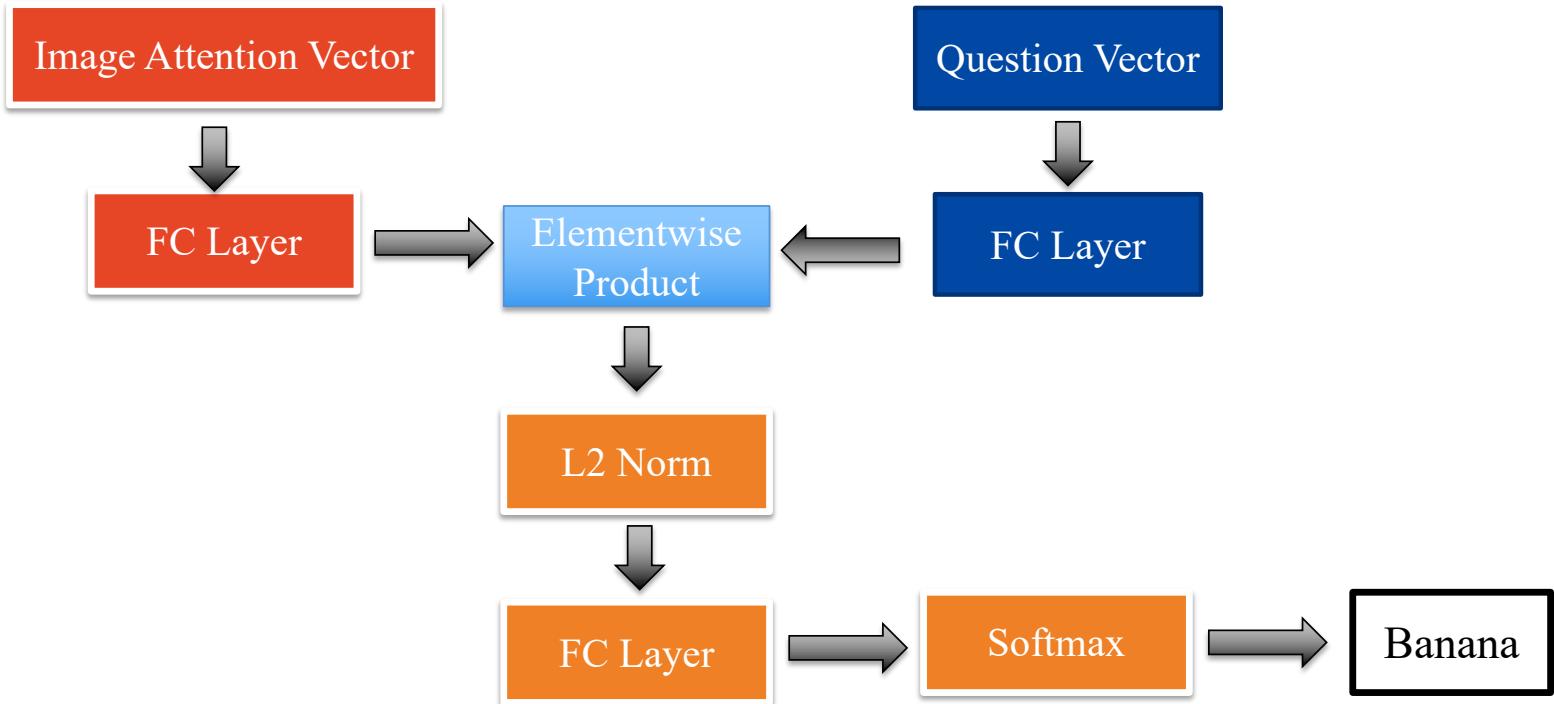


What is the man holding?

Image credit to Bottom-Up and Top-Down
Attention for Image Captioning and Visual
Question Answering

VQA

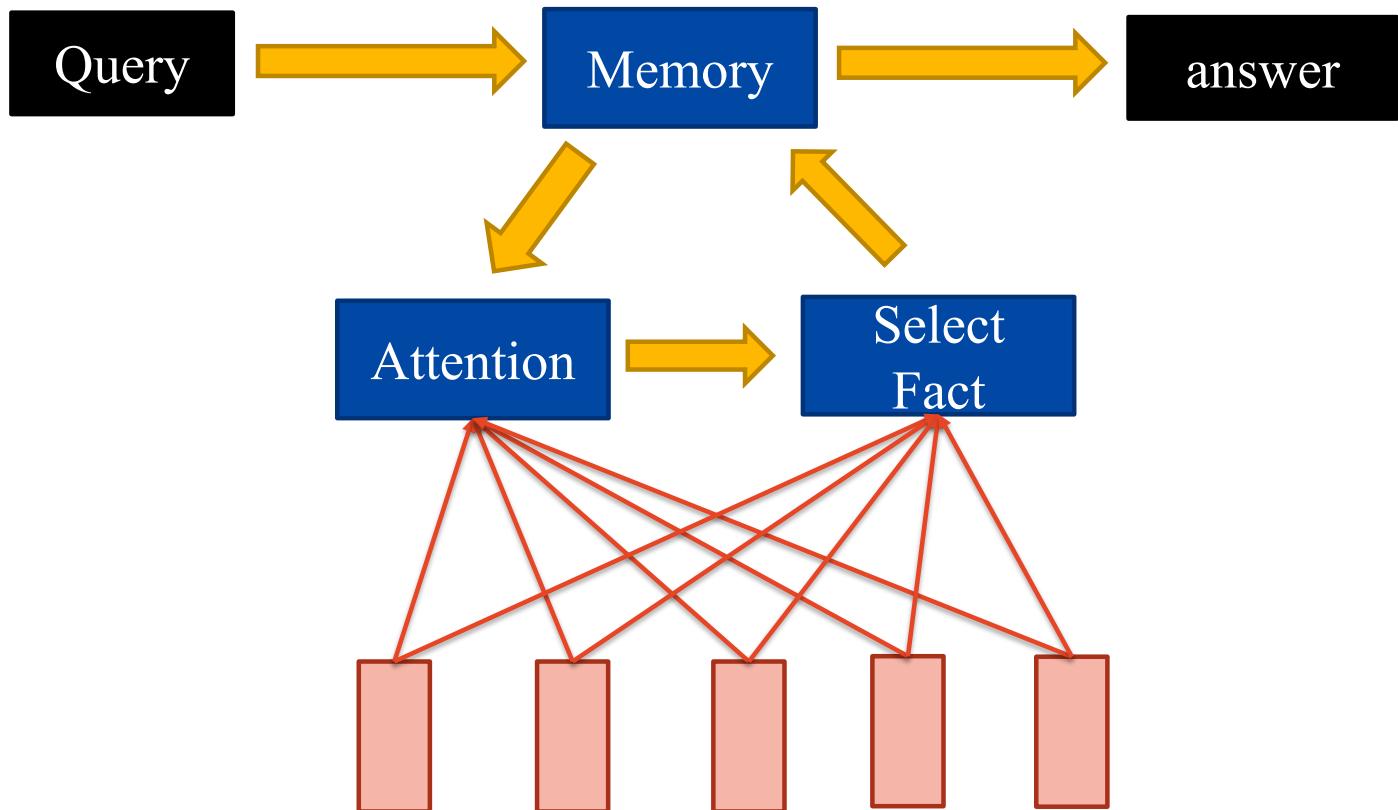
- Question & Image Fusion



Reading Comprehension

- Answering question based on given facts
- Example
 - A. Brian is a frog.
 - B. Lily is gray.
 - C. Brian is yellow
 - D. Julius is green.
 - E. Greg is a frog
- ✓ What color is Greg?

Reading Comprehension



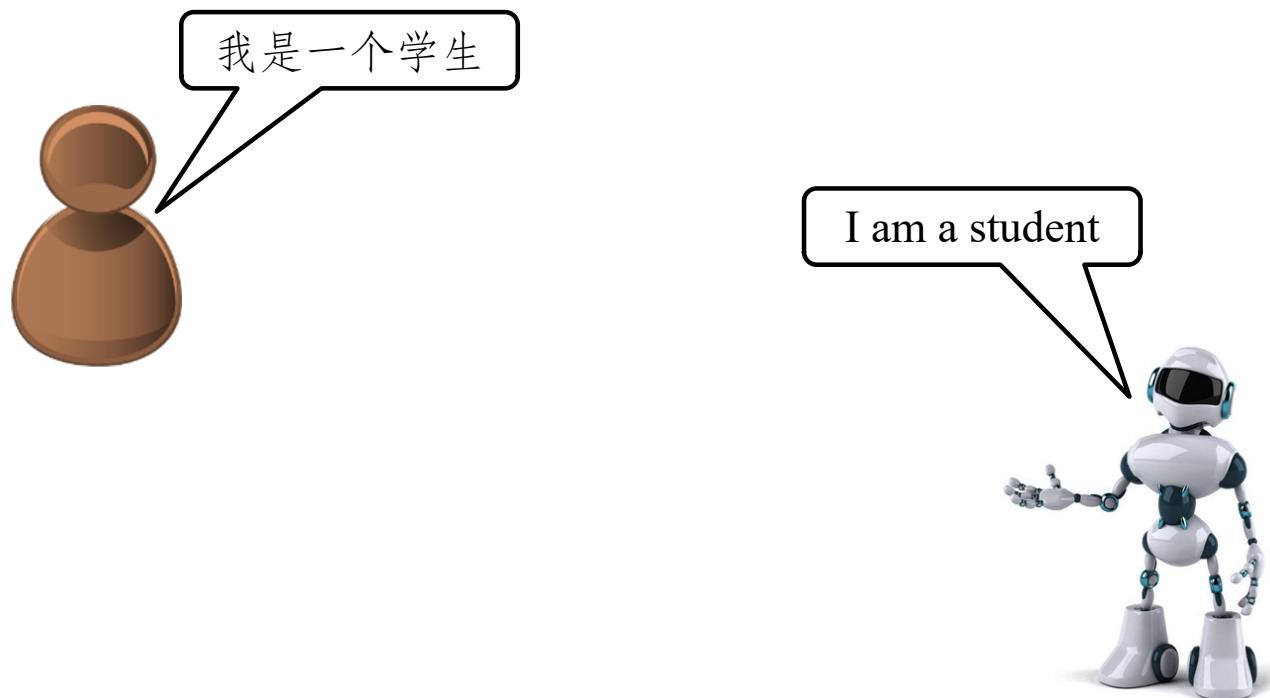
Reading Comprehension

Story (16: basic induction)	Support	Hop 1	Hop 2	Hop 3
Brian is a frog.	yes	0.00	0.98	0.00
Lily is gray.		0.07	0.00	0.00
Brian is yellow.	yes	0.07	0.00	1.00
Julius is green.		0.06	0.00	0.00
Greg is a frog.	yes	0.76	0.02	0.00
What color is Greg? Answer: yellow		Prediction: yellow		

End-To-End Memory Networks. S. Sukhbaatar, A. Szlam, J. Weston, R. Fergus. NIPS, 2015.

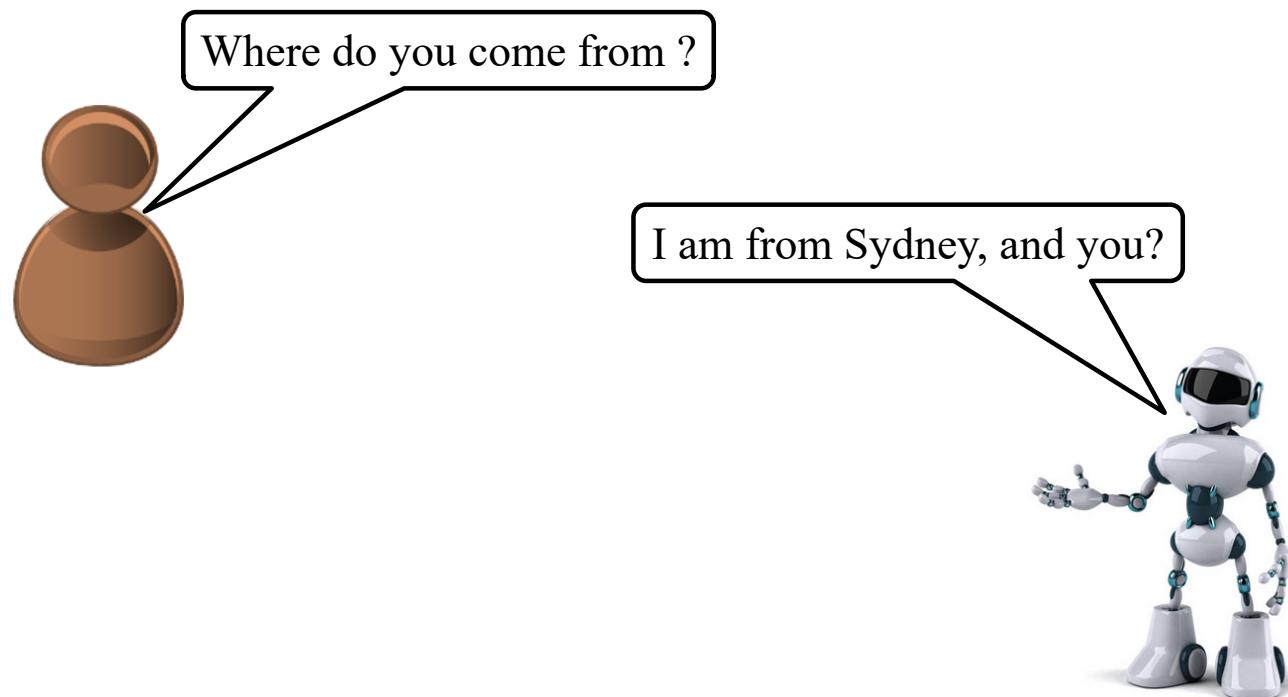
Seq2seq

□ Language Translation



Seq2seq

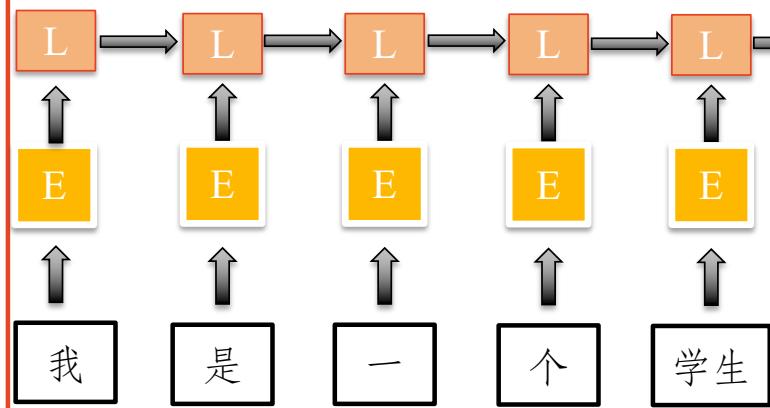
□ Chat Robot



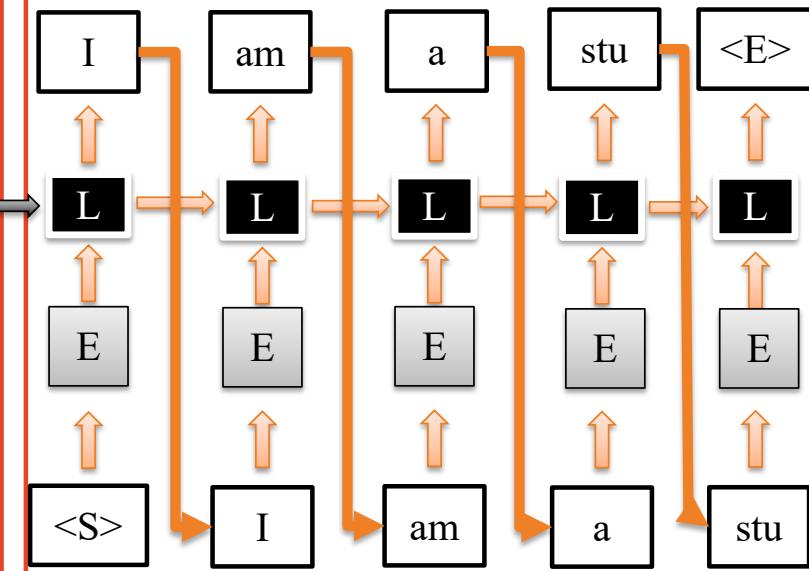
Seq2seq

Encoder

v : LSTM
 E : Embedding

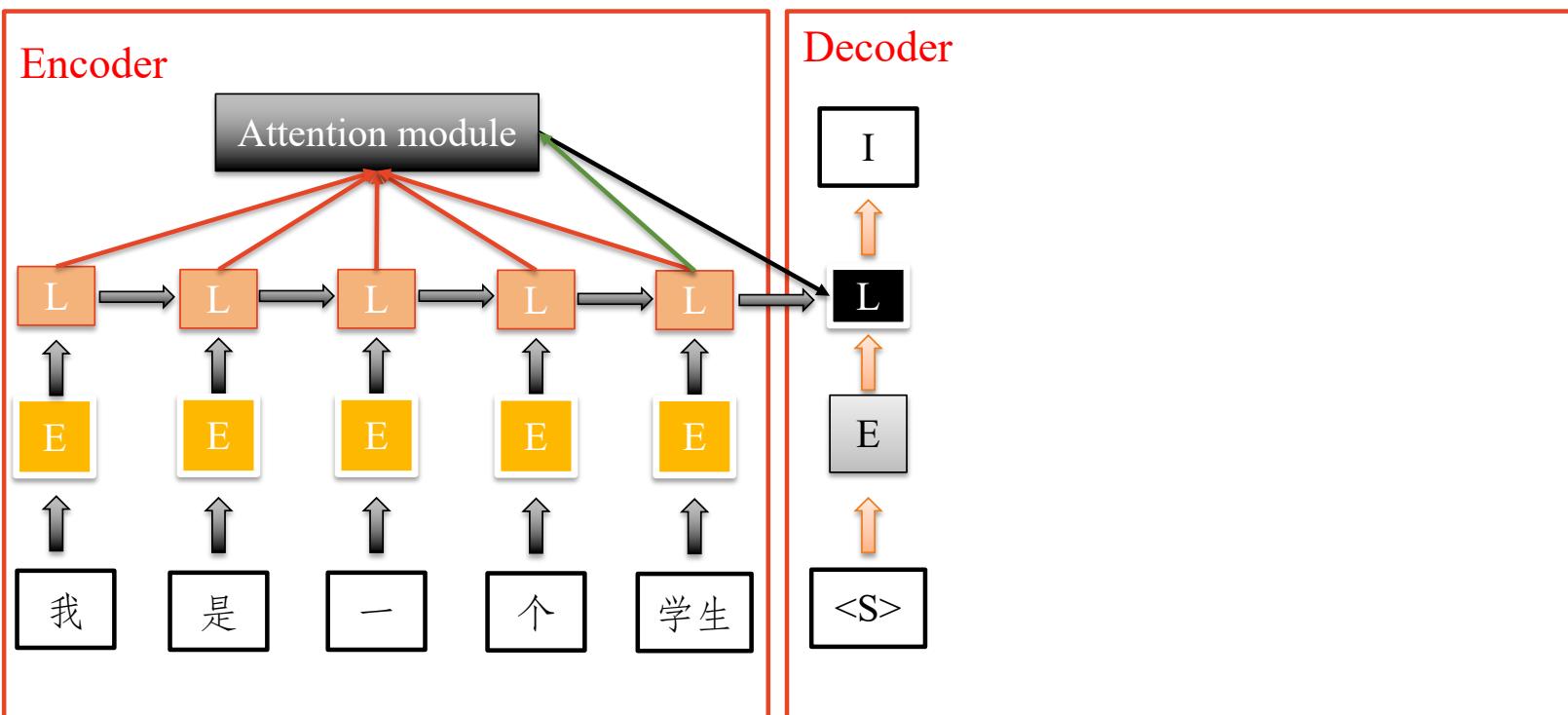


Decoder



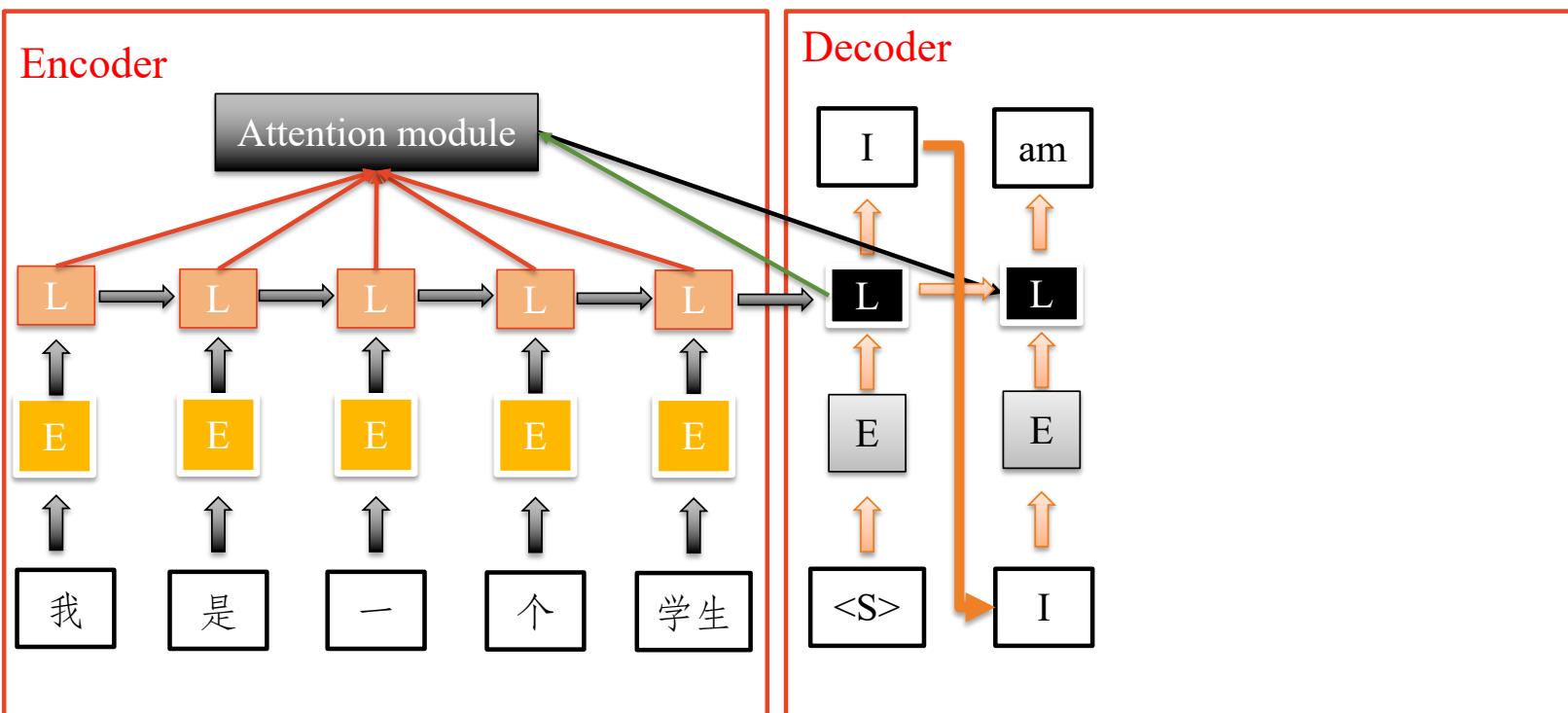
Seq2seq

□ Attention



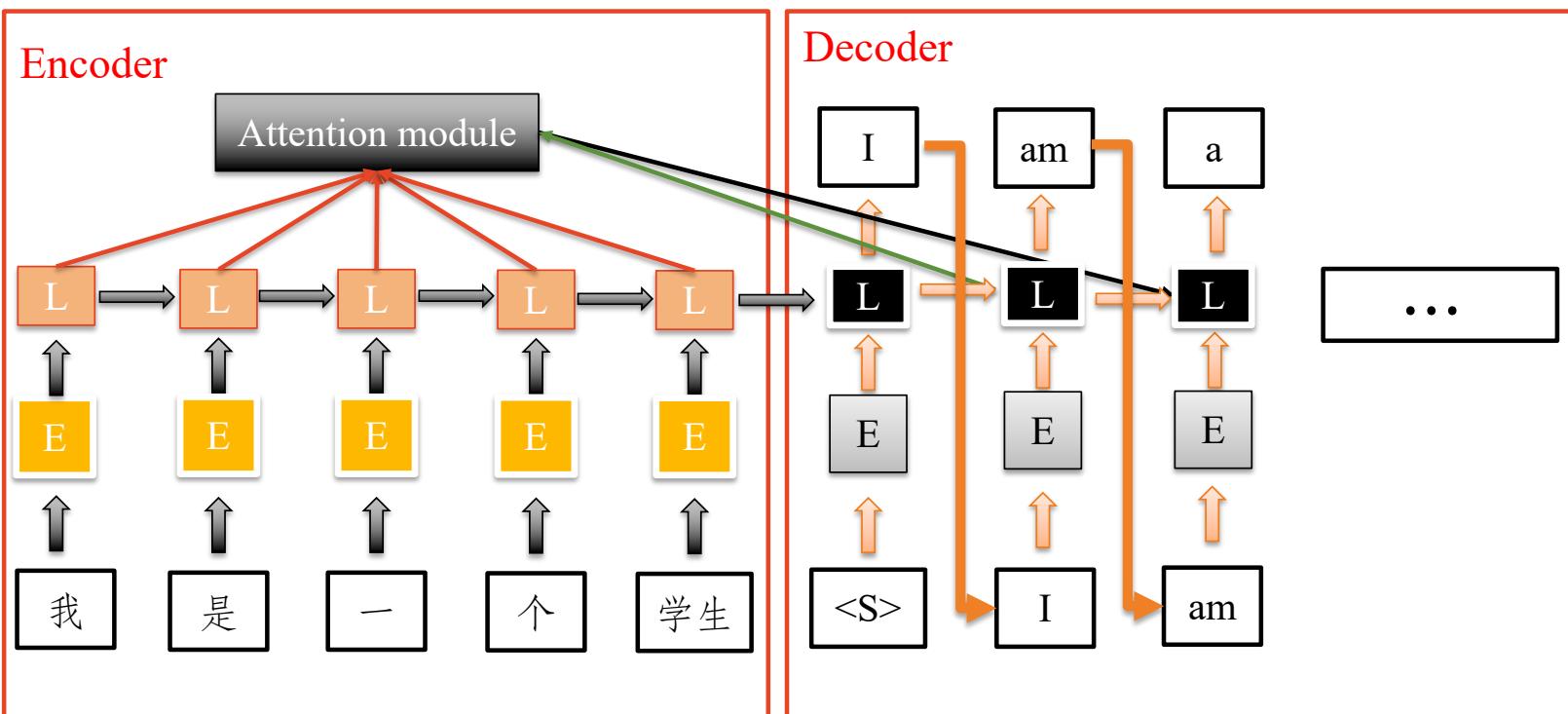
Seq2seq

□ Attention



Seq2seq

□ Attention



Seq2seq

- Two stage
 - Encode stage
 - Decode stage
- Attention
 - Weight for each encode symbol in decode step
 - Another input of LSTM represented as a_t
 - $f_t = \sigma(w_{fh}h_{t-1} + w_{fx}x_t + w_{fa}a_t + b_f)$
 - $i_t = \sigma(w_{ih}h_{t-1} + w_{ix}x_t + w_{ia}a_t + b_i)$
 - $o_t = \sigma(w_{oh}h_{t-1} + w_{ox}x_t + w_{oa}a_t + b_o)$
 - $\tilde{c}_t = \tanh(w_{hh}h_{t-1} + w_{hx}x_t + w_{ha}a_t + b_h)$

Seq2seq

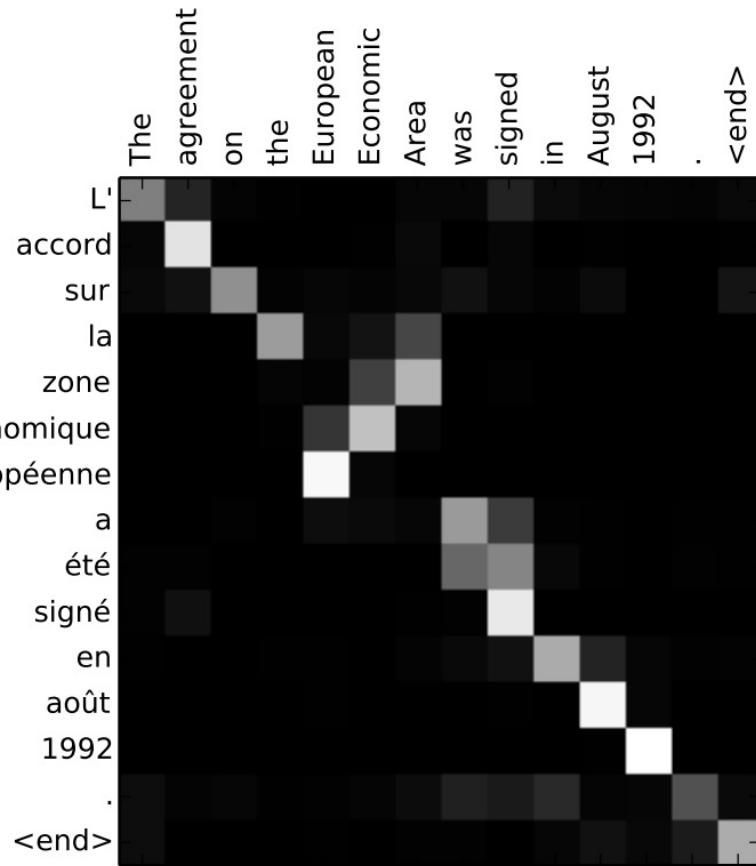


Image credit to Neural Machine Translation by
Jointly Learning to Align and Translate

Thank you !