

Foundations of Deep Learning



ALF

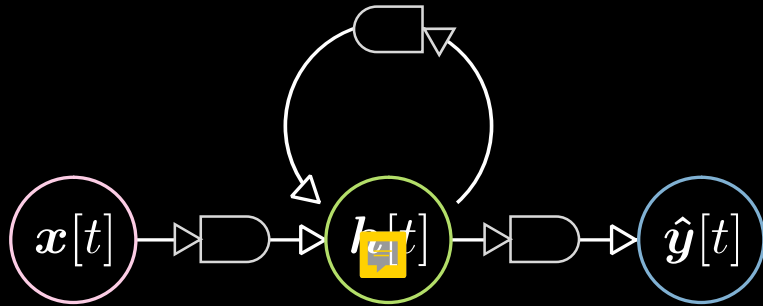
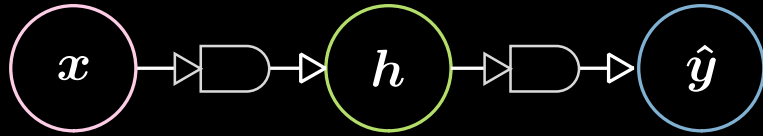
Alfredo Canziani

 @alfcnz

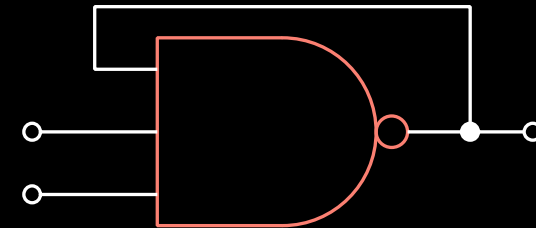
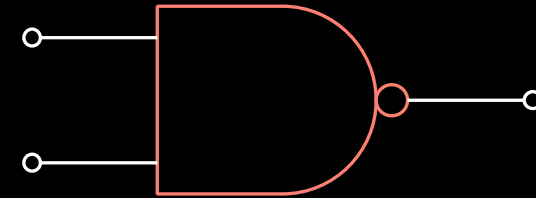
Recurrent Neural Nets

Handling sequential data

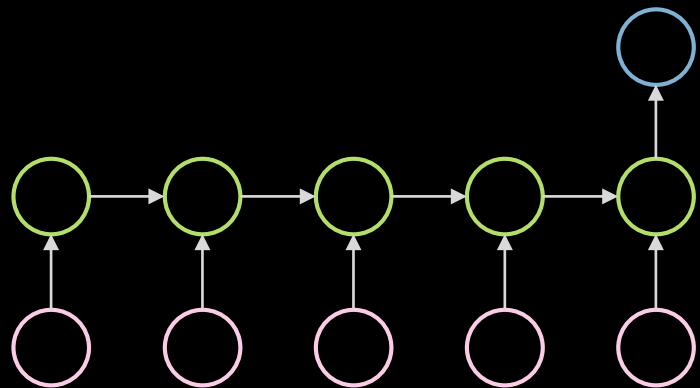
Vanilla and Recurrent NN



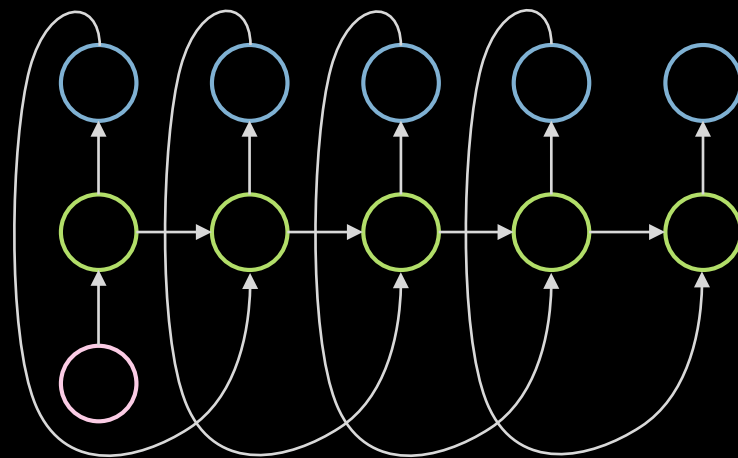
Combinatorial logic



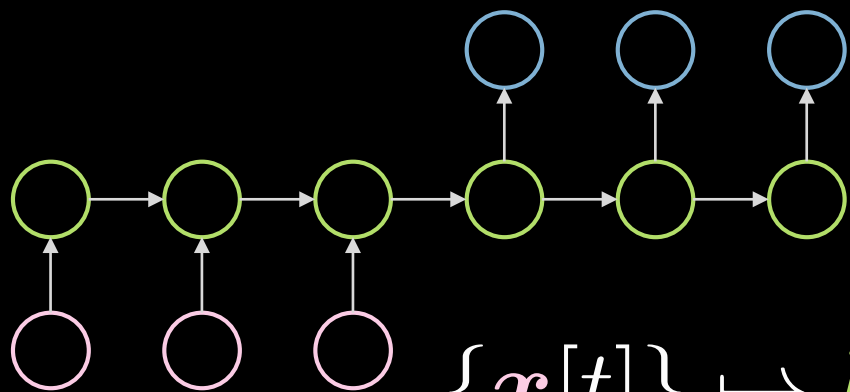
Sequential logic



$$\{\mathbf{x}[t]\} \mapsto \hat{\mathbf{y}}[T] \quad \text{seq} \xrightarrow{\text{vec}} \text{vec}$$

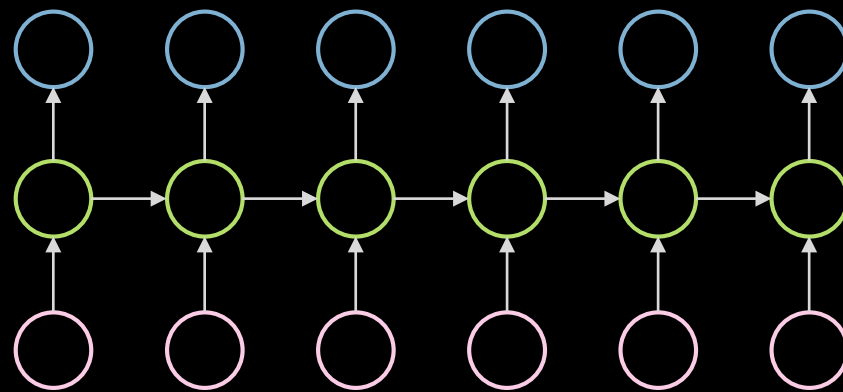


$$\mathbf{x}[1] \mapsto \{\hat{\mathbf{y}}[t]\} \quad \text{vec} \xrightarrow{\text{seq}} \text{seq}$$



$$\{\mathbf{x}[t]\} \mapsto \mathbf{h} \mapsto \{\hat{\mathbf{y}}[t]\}$$

$$\text{seq} \mapsto \text{vec} \xrightarrow{\text{seq}} \text{seq}$$



$$\{\mathbf{x}[t]\} \mapsto \{\hat{\mathbf{y}}[t]\} \quad \text{seq} \mapsto \text{seq}$$

A person riding a motorcycle on a dirt road.



Two dogs play in the grass.



A skateboarder does a trick on a ramp.



A dog is jumping to catch a frisbee.



A group of young people playing a game of frisbee.



Two hockey players are fighting over the puck.



A little girl in a pink hat is blowing bubbles.



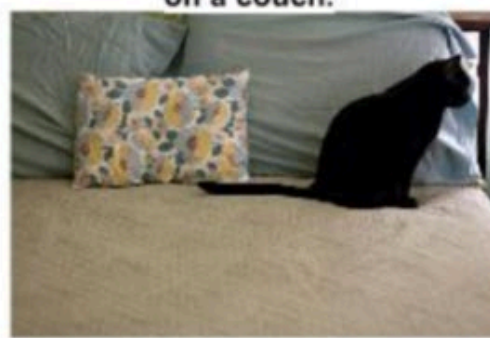
A refrigerator filled with lots of food and drinks.



A herd of elephants walking across a dry grass field.



A close up of a cat laying on a couch.



A red motorcycle parked on the side of the road.



A yellow school bus parked in a parking lot.



Describes without errors

Describes with minor errors

Somewhat related to the image

Unrelated to the image

Learning to execute

- **Input:**

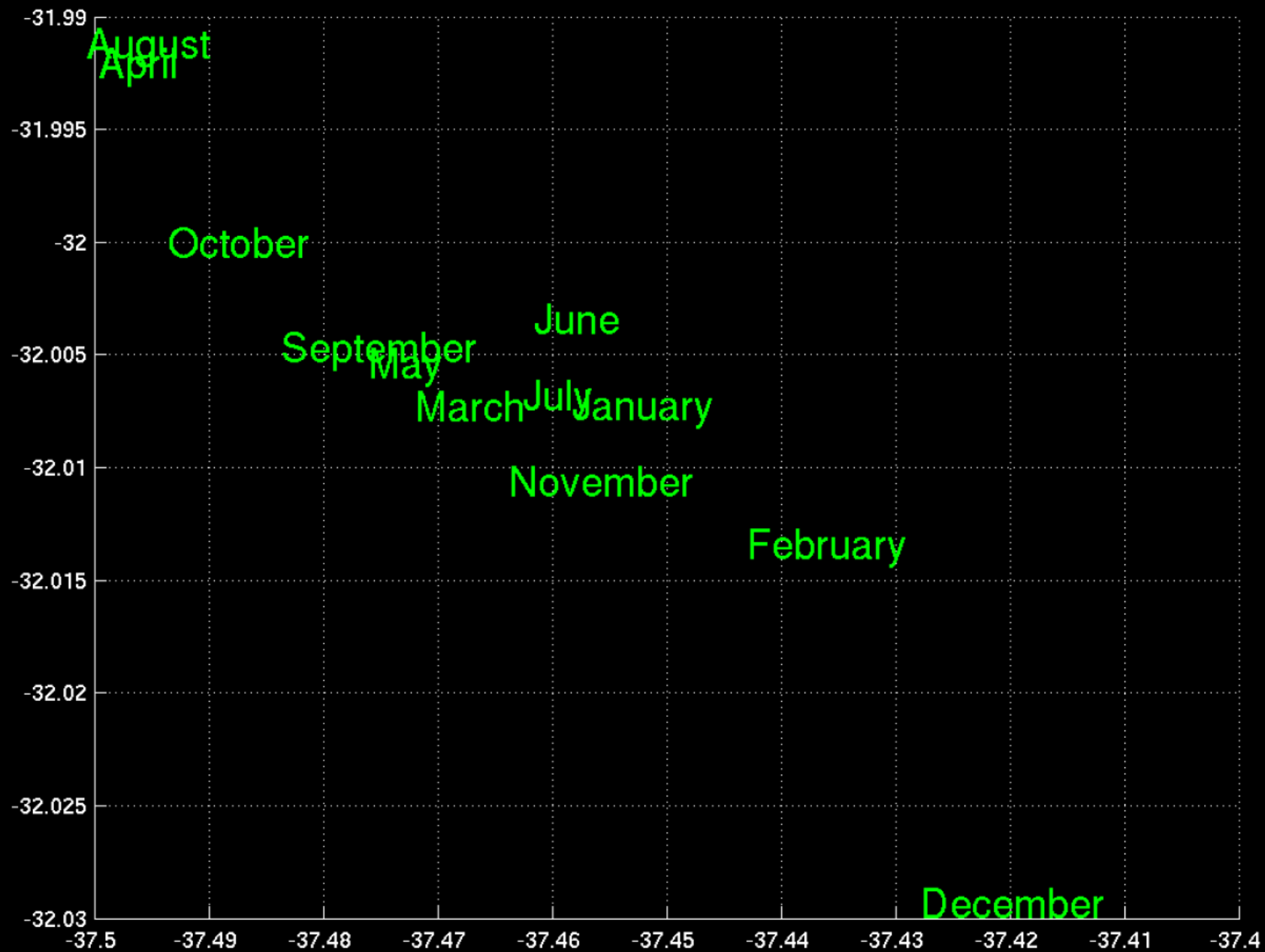
```
j=8584
for x in range(8):
    j+=920
    b=(1500+j)
    print((b+7567))
```

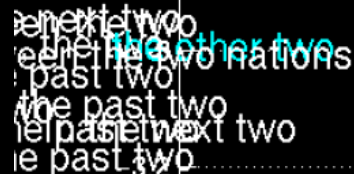
- **Target:** 25011.

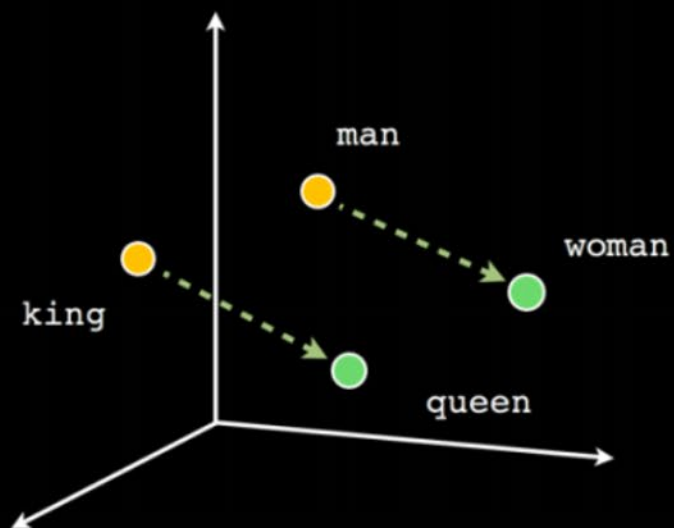
- **Input:**

```
i=8827
c=(i-5347)
print((c+8704) if
2641<8500 else 5308)
```

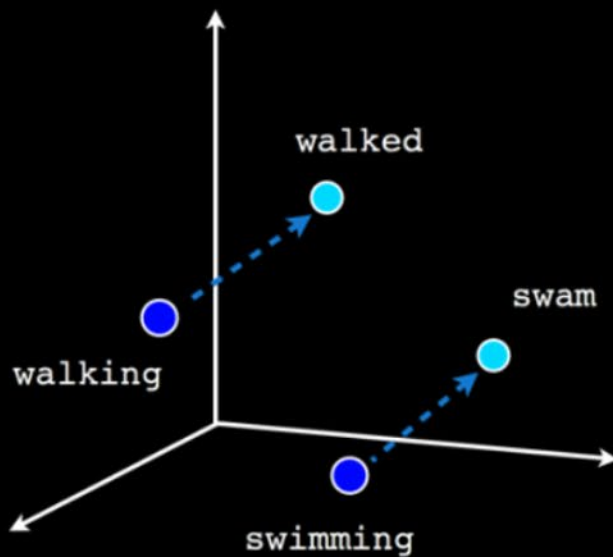
- **Target:** 12184.







Male-Female



Verb tense



Country-Capital

test.txt

rnn-client.coffee

1 The

SYSTEM PROMPT
(HUMAN-WRITTEN)

In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

MODEL COMPLETION
(MACHINE-
WRITTEN, 10 TRIES)

The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

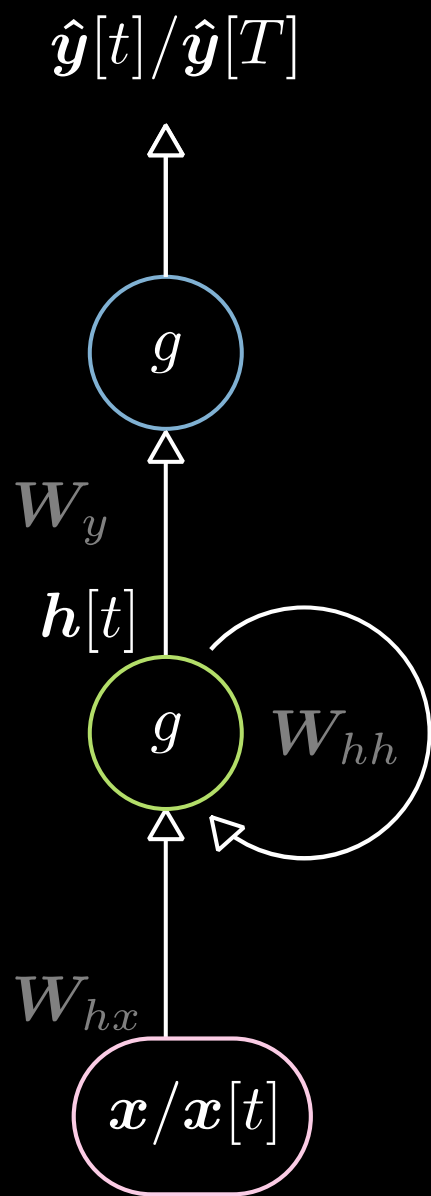
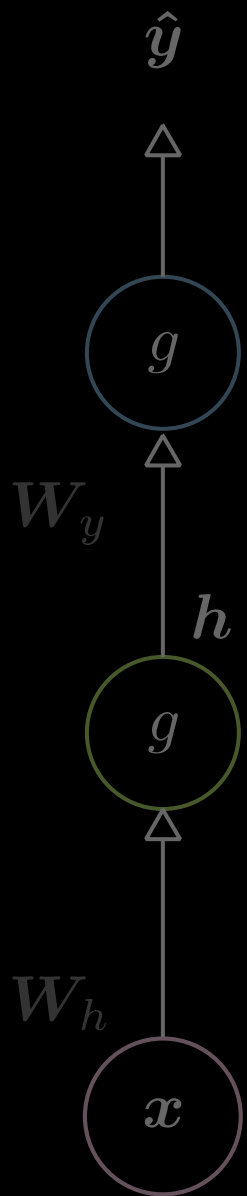
Pérez and the others then ventured further into the valley. "By the time we reached the top of one peak, the water looked blue, with some crystals on top," said Pérez.

Pérez and his friends were astonished to see the unicorn herd. These creatures could be seen from the air without having to move too much to see them – they were so close they could touch their horns.

While examining these bizarre creatures the scientists discovered that the creatures also spoke some fairly regular English. Pérez stated, "We can see, for example, that they have a common 'language,' something like a dialect or dialectic."

RNN training

Back propagation through time (BPTT)



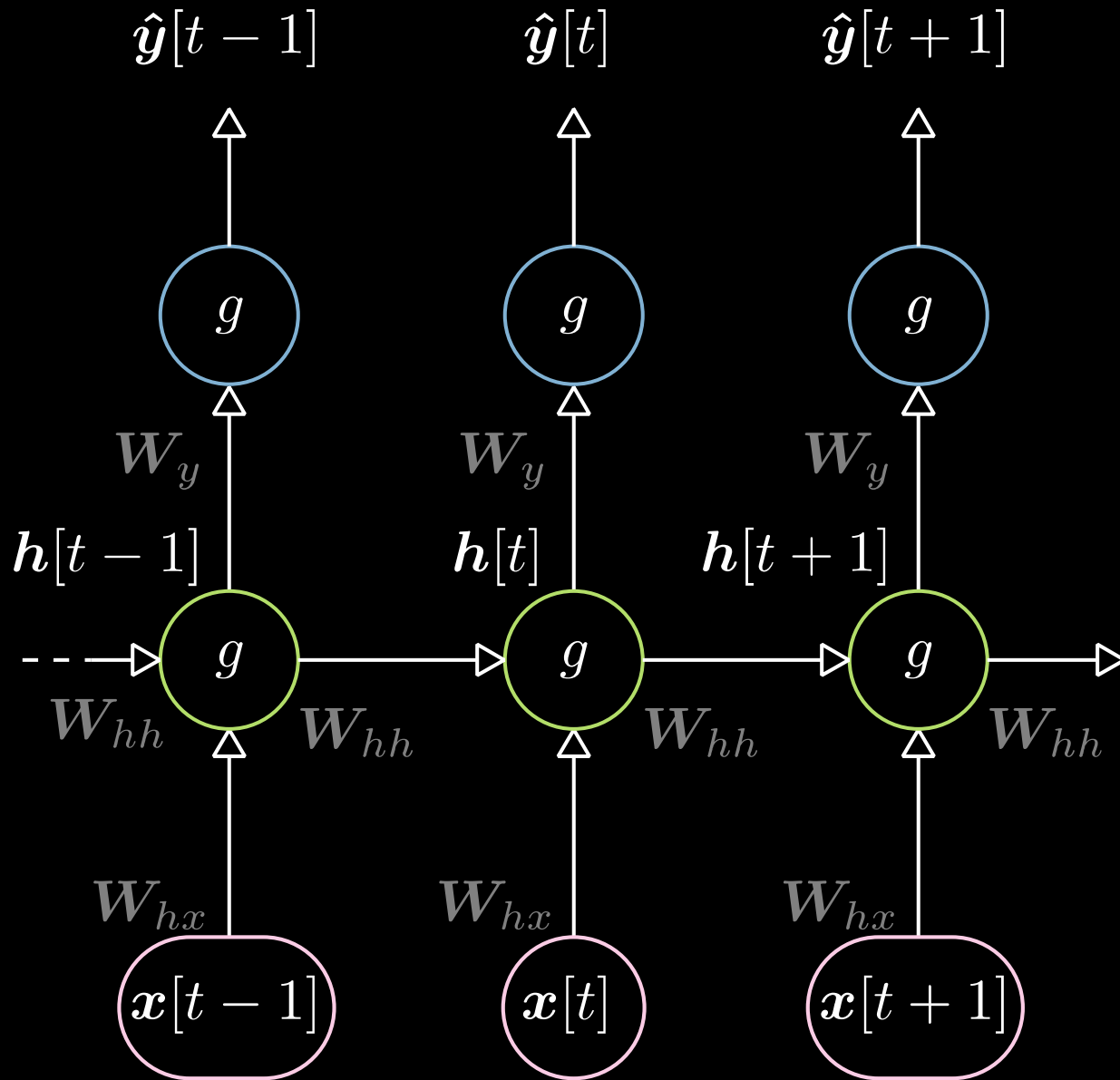
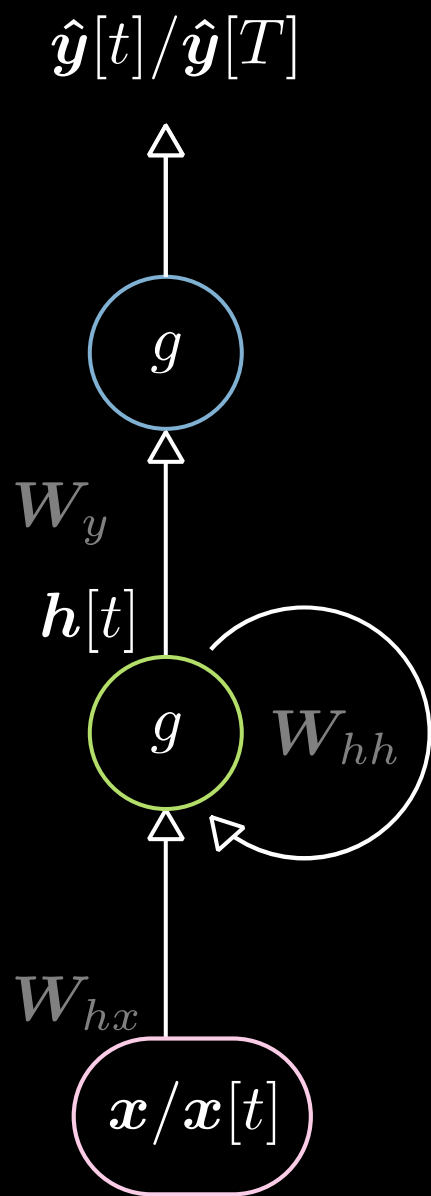
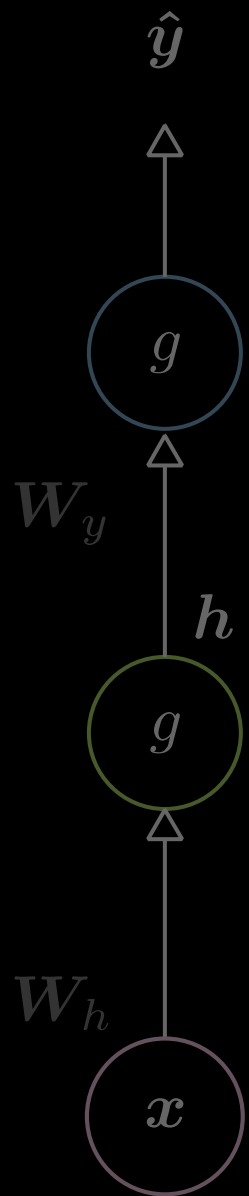
$$h = f(W_h x + b_h)$$

$$\hat{y} = g(W_y h + b_y)$$

$$h[t] = g(W_h \begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_h)$$

$$h[0] \doteq \mathbf{0}, W_h \doteq \begin{bmatrix} W_{hx} & W_{hh} \end{bmatrix}$$

$$\hat{y}[t] = g(W_y h[t] + b_y)$$



Training example

Language modelling

Batch-ification

abcdefghijklmnopqrstuvwxyz



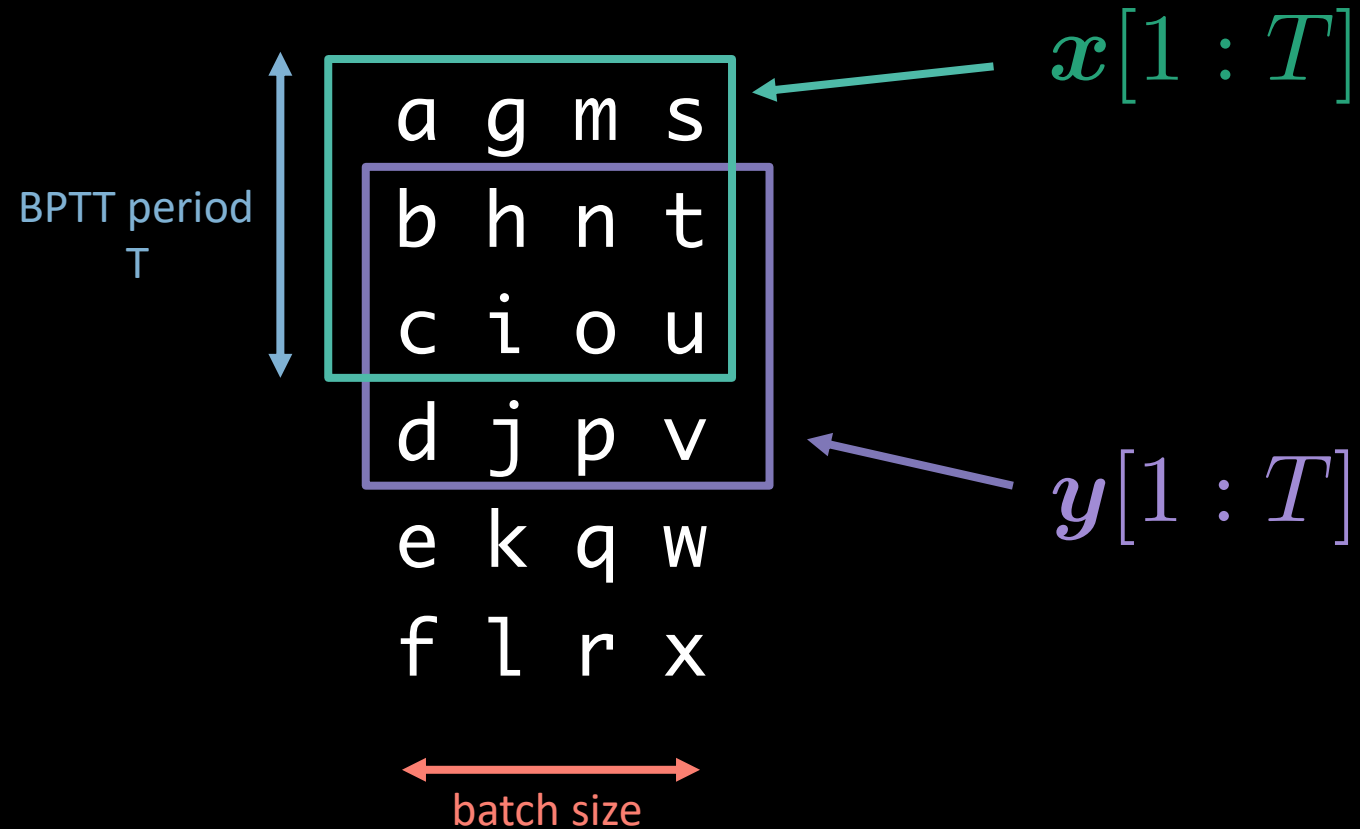
a	g	m	s
b	h	n	t
c	i	o	u
d	j	p	v
e	k	q	w
f	l	r	x



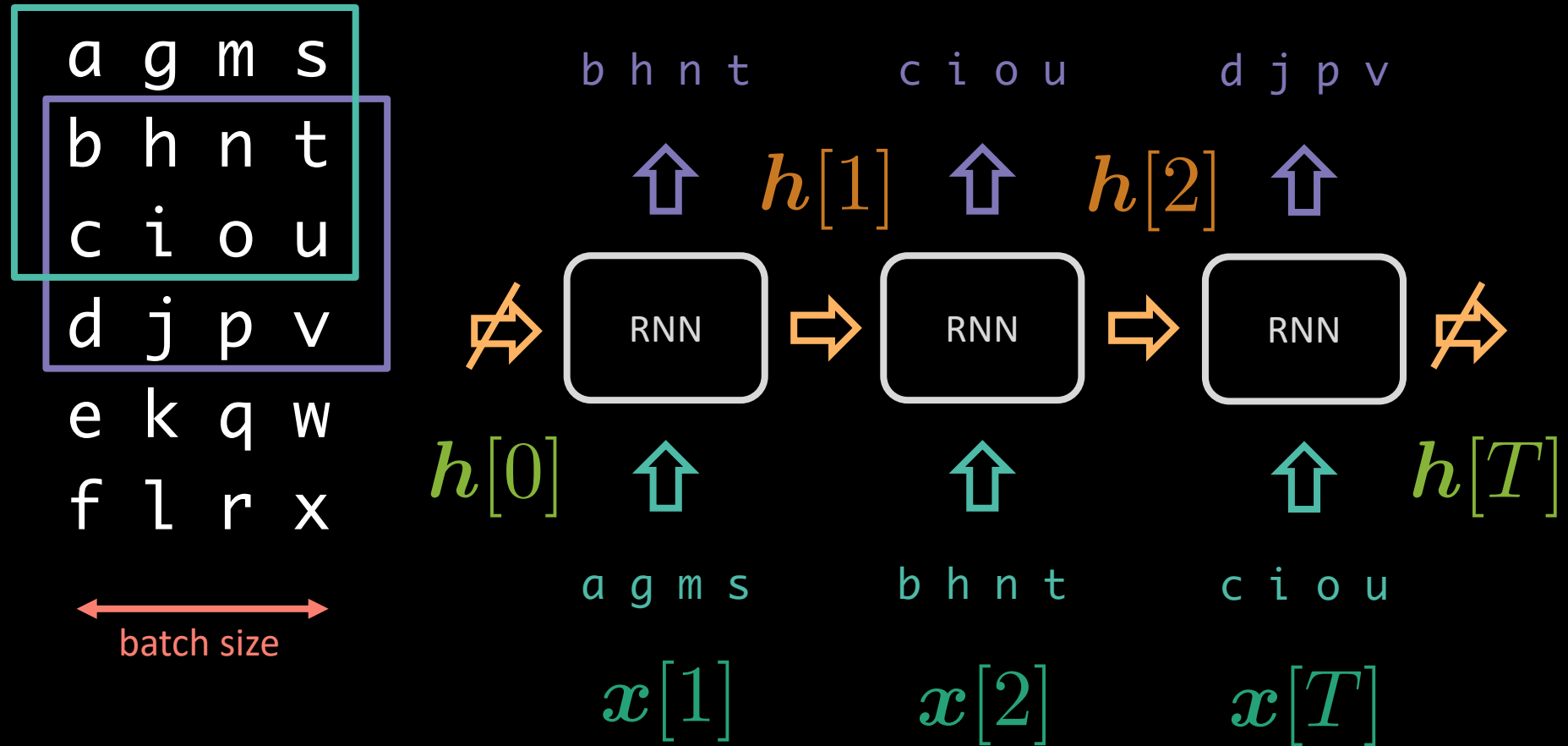
batch size

Check `word_language_model` @ github.com/pytorch/examples/

Get batch (I)

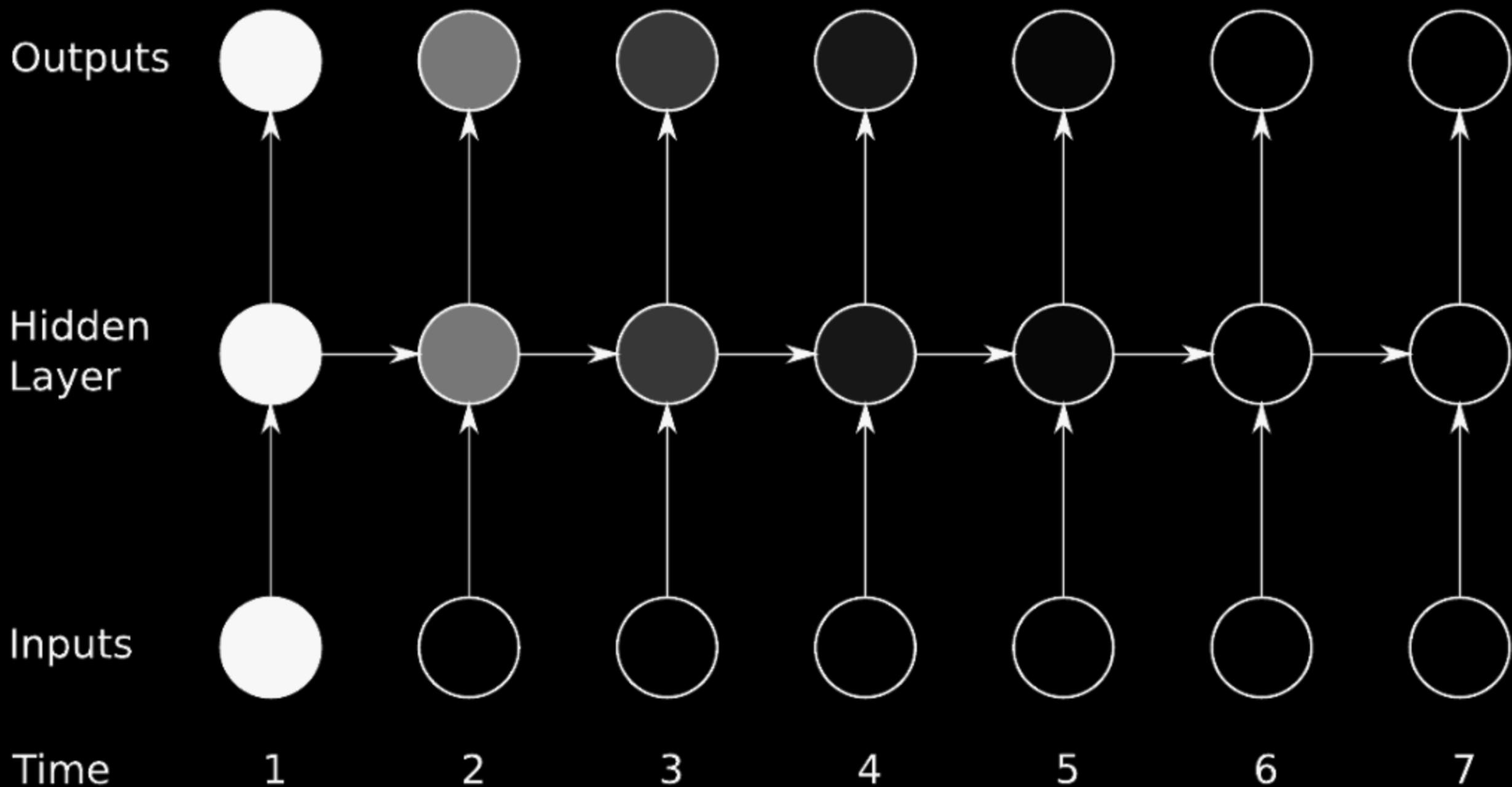


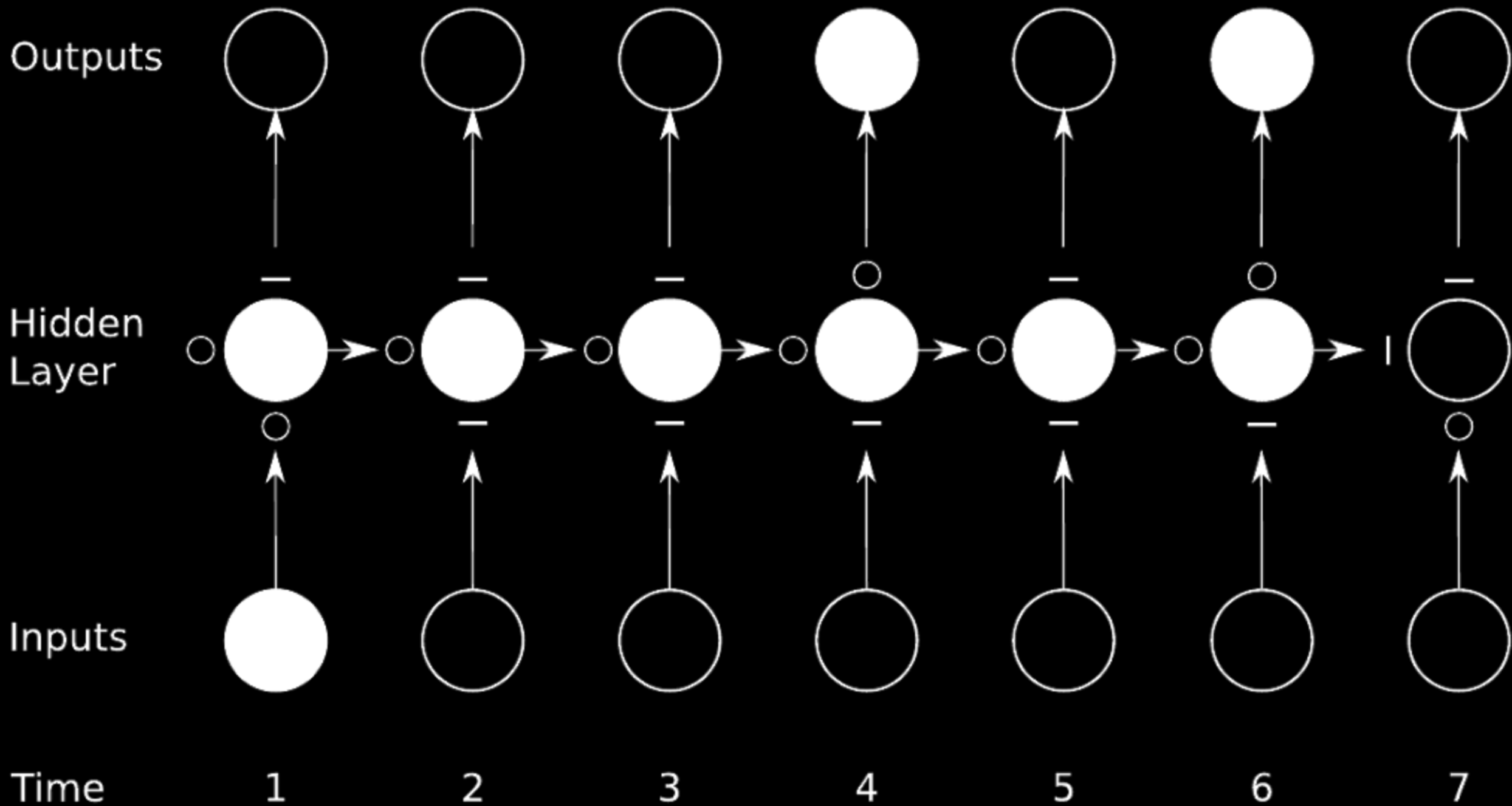
Get batch (II)



Vanishing & exploding gradients

Limitations of temporally deep nets



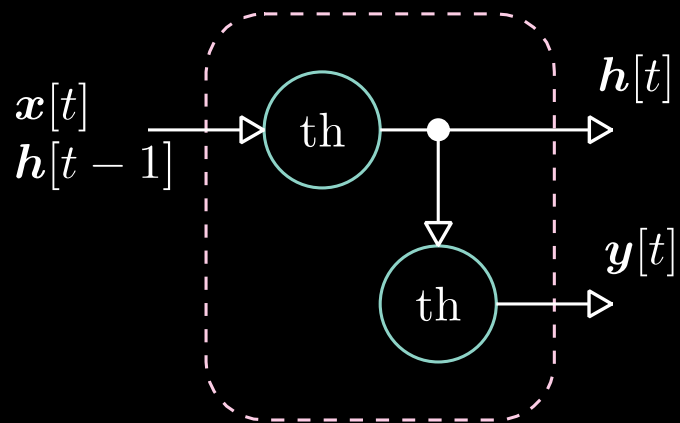


Graves (2012) Supervised sequence labelling

Long Short-Term Memory

Gated RNN





$$h[t] = g(\mathbf{W}_h [h[t-1]] + \mathbf{b}_h)$$

$$\hat{y}[t] = g(\mathbf{W}_y h[t] + \mathbf{b}_y)$$

$$i[t] = \sigma(\mathbf{W}_i [h[t-1]] + \mathbf{b}_i)$$

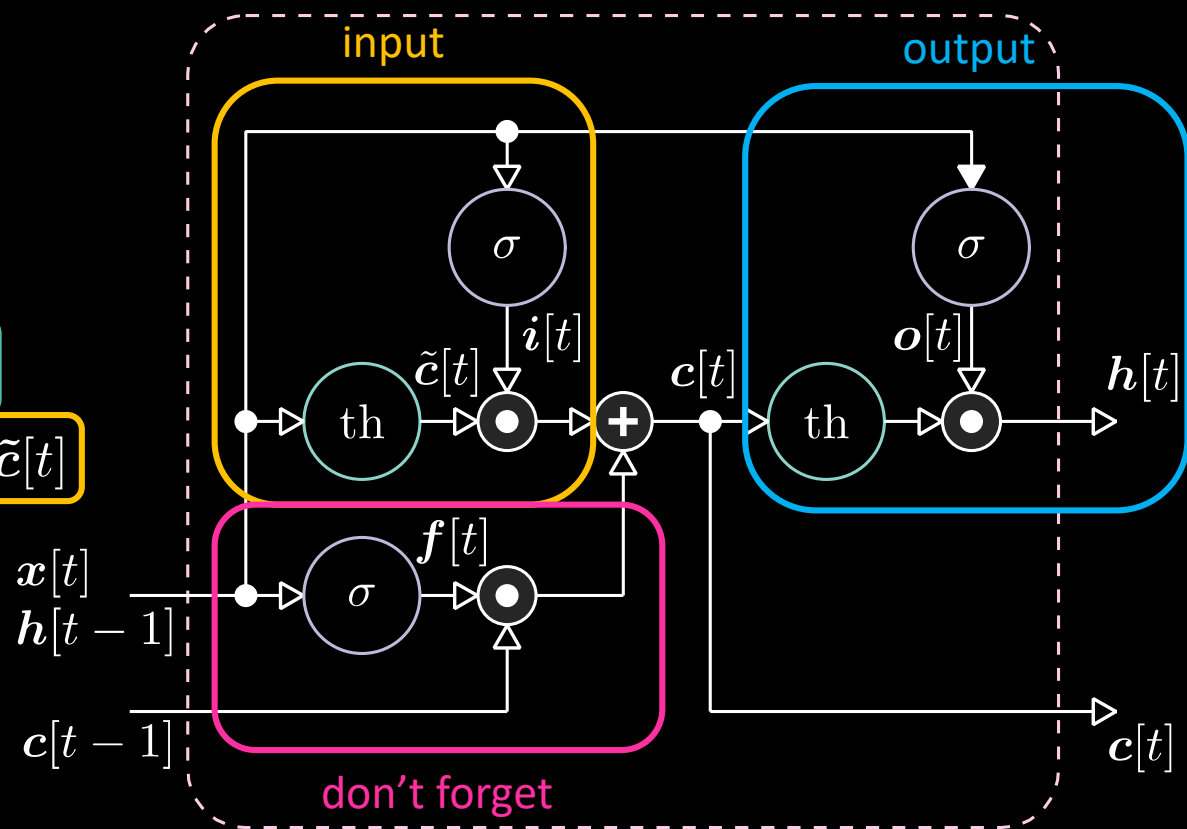
$$f[t] = \sigma(\mathbf{W}_f [h[t-1]] + \mathbf{b}_f)$$

$$o[t] = \sigma(\mathbf{W}_o [h[t-1]] + \mathbf{b}_o)$$

$$\tilde{c}[t] = \tanh(\mathbf{W}_c [h[t-1]] + \mathbf{b}_c)$$

$$c[t] = f[t] \odot c[t-1] + i[t] \odot \tilde{c}[t]$$

$$h[t] = o[t] \odot \tanh(c[t])$$



Controlling the output - OFF

Saturated sigmoid $\begin{cases} \text{green circle} = 1 \\ \text{red circle} = 0 \end{cases}$

$$i[t] = \sigma(\mathbf{W}_i [\mathbf{x}^t] + \mathbf{b}_i)$$

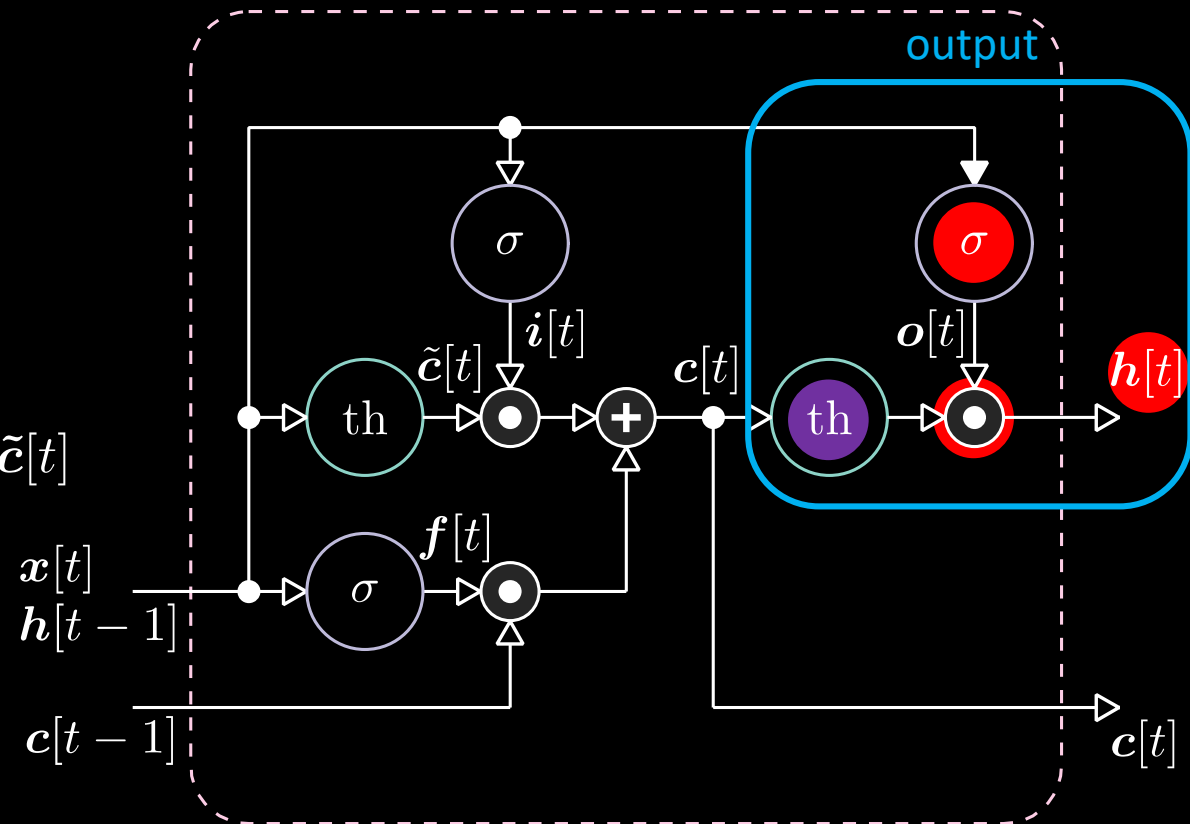
$$f[t] = \sigma(\mathbf{W}_f [\mathbf{x}^t] + \mathbf{b}_f)$$

$$o[t] = \sigma(\mathbf{W}_o [\mathbf{x}^t] + \mathbf{b}_o)$$

$$\tilde{c}[t] = \tanh(\mathbf{W}_c [\mathbf{x}^t] + \mathbf{b}_c)$$

$$c[t] = f[t] \odot c[t-1] + i[t] \odot \tilde{c}[t]$$

$$h[t] = o[t] \odot \tanh(c[t])$$



Controlling the output - ON

Saturated sigmoid $\begin{cases} \text{green circle} = 1 \\ \text{red circle} = 0 \end{cases}$

$$i[t] = \sigma(\mathbf{W}_i [\mathbf{x}^t] + \mathbf{b}_i)$$

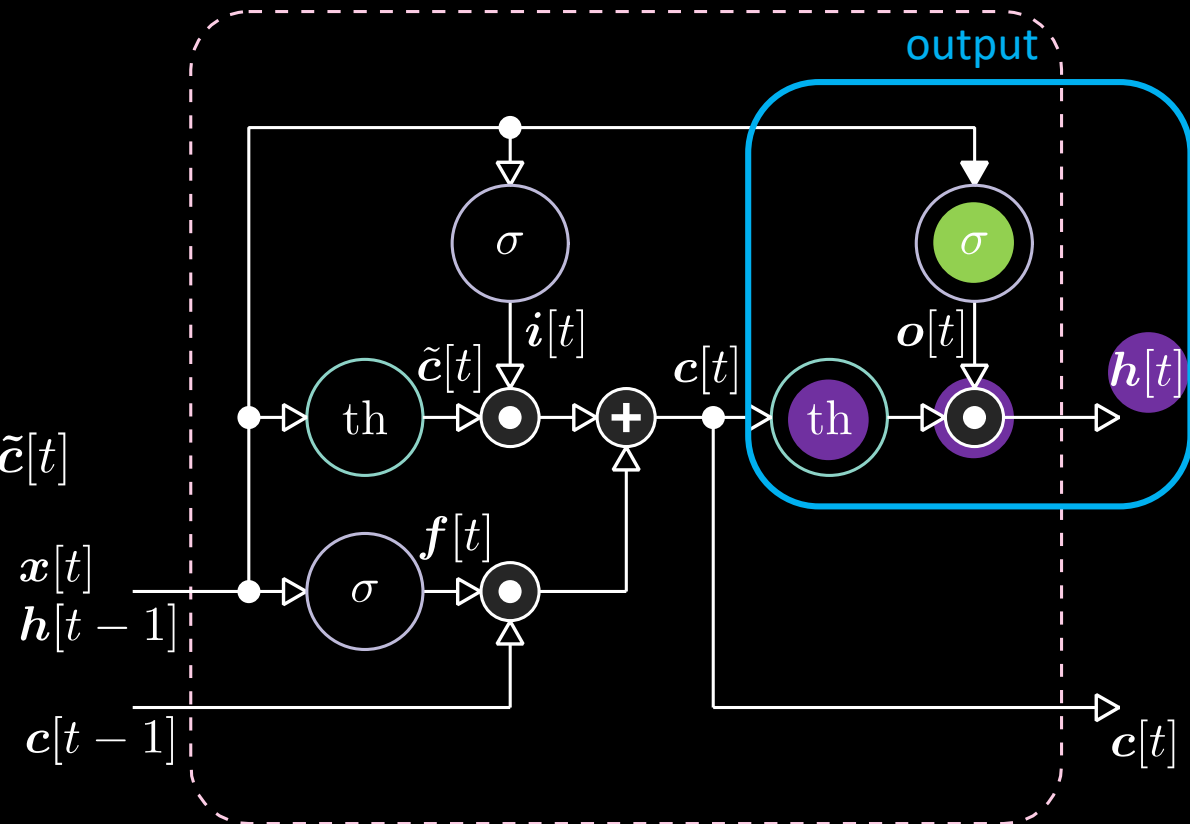
$$f[t] = \sigma(\mathbf{W}_f [\mathbf{x}^t] + \mathbf{b}_f)$$

$$o[t] = \sigma(\mathbf{W}_o [\mathbf{x}^t] + \mathbf{b}_o)$$

$$\tilde{c}[t] = \tanh(\mathbf{W}_c [\mathbf{x}^t] + \mathbf{b}_c)$$

$$c[t] = f[t] \odot c[t-1] + i[t] \odot \tilde{c}[t]$$

$$h[t] = o[t] \odot \tanh(c[t])$$



Controlling the memory - reset

Saturated sigmoid $\left\{ \begin{array}{l} \text{green circle} = 1 \\ \text{red circle} = 0 \end{array} \right.$

$$i[t] = \sigma(\mathbf{W}_i [\mathbf{x}^t] + \mathbf{b}_i)$$

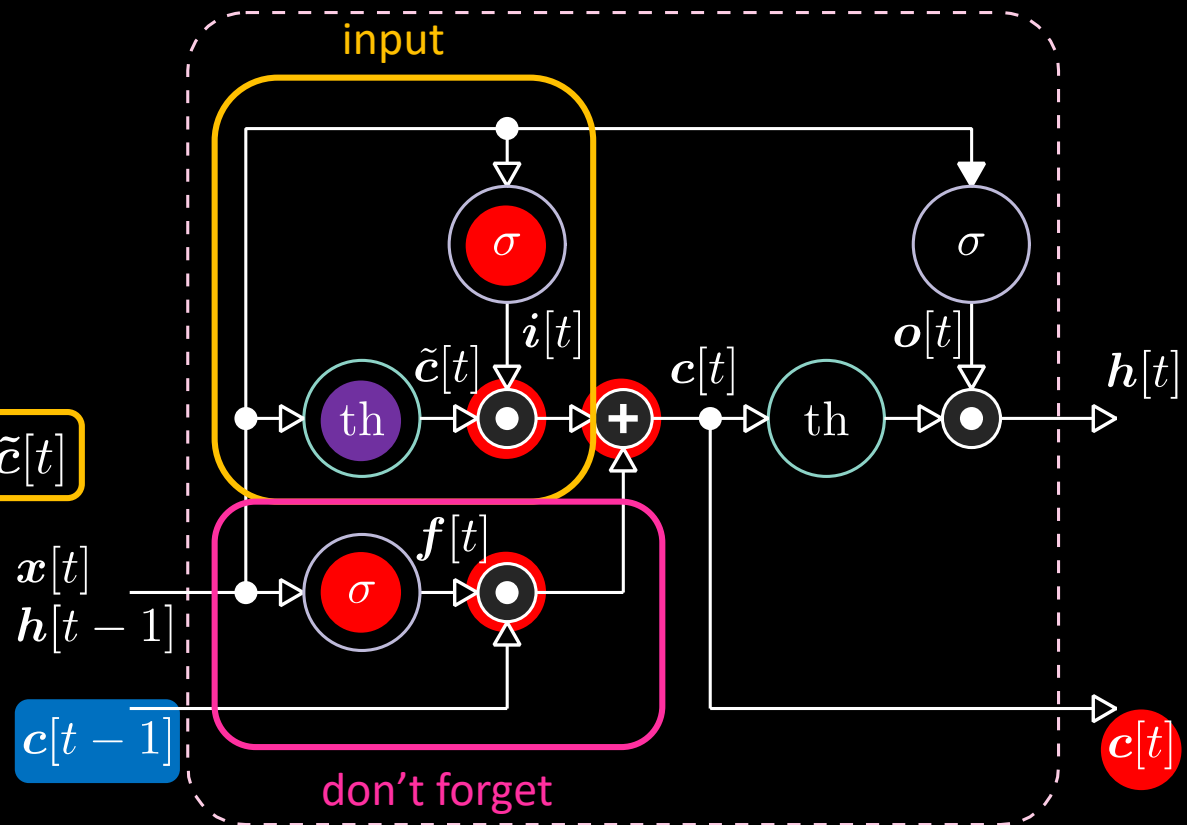
$$f[t] = \sigma(\mathbf{W}_f [\mathbf{x}^t] + \mathbf{b}_f)$$

$$o[t] = \sigma(\mathbf{W}_o [\mathbf{x}^t] + \mathbf{b}_o)$$

$$\tilde{c}[t] = \tanh(\mathbf{W}_c [\mathbf{x}^t] + \mathbf{b}_c)$$

$$c[t] = f[t] \odot c[t-1] + i[t] \odot \tilde{c}[t]$$

$$h[t] = o[t] \odot \tanh(c[t])$$



Controlling the **memory** - keep

Saturated sigmoid $\left\{ \begin{array}{l} \text{green circle} = 1 \\ \text{red circle} = 0 \end{array} \right.$

$$i[t] = \sigma(\mathbf{W}_i [x[t], h[t-1]] + b_i)$$

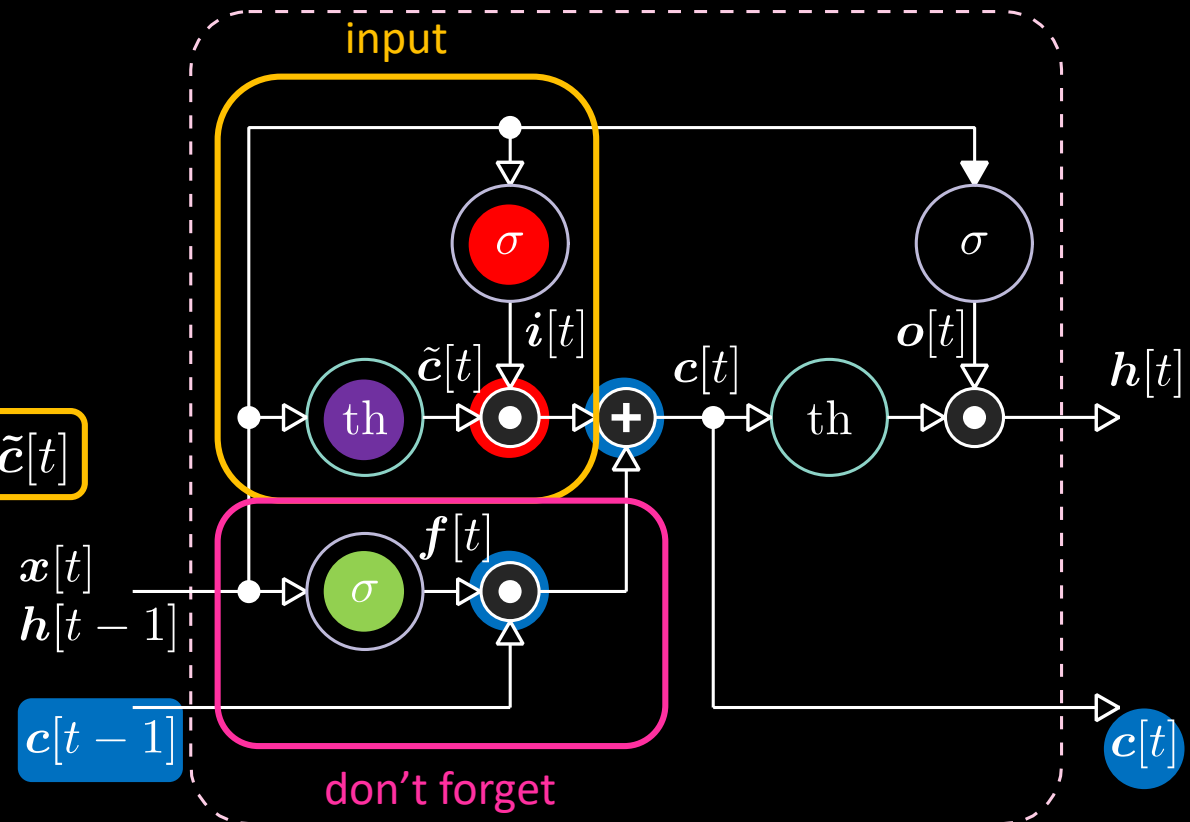
$$f[t] = \sigma(\mathbf{W}_f [x[t], h[t-1]] + b_f)$$

$$o[t] = \sigma(\mathbf{W}_o [x[t], h[t-1]] + b_o)$$

$$\tilde{c}[t] = \tanh(\mathbf{W}_c [x[t], h[t-1]] + b_c)$$

$$c[t] = f[t] \odot c[t-1] + i[t] \odot \tilde{c}[t]$$

$$h[t] = o[t] \odot \tanh(c[t])$$



Controlling the **memory** - write

Saturated sigmoid $\left\{ \begin{array}{l} \text{green circle} = 1 \\ \text{red circle} = 0 \end{array} \right.$

$$i[t] = \sigma(\mathbf{W}_i [\mathbf{x}^{[t]} \parallel \mathbf{h}^{[t-1]}] + \mathbf{b}_i)$$

$$f[t] = \sigma(\mathbf{W}_f [\mathbf{x}^{[t]} \parallel \mathbf{h}^{[t-1]}] + \mathbf{b}_f)$$

$$o[t] = \sigma(\mathbf{W}_o [\mathbf{x}^{[t]} \parallel \mathbf{h}^{[t-1]}] + \mathbf{b}_o)$$

$$\tilde{c}[t] = \tanh(\mathbf{W}_c [\mathbf{x}^{[t]} \parallel \mathbf{h}^{[t-1]}] + \mathbf{b}_c)$$

$$c[t] = f[t] \odot c[t-1] + i[t] \odot \tilde{c}[t]$$

$$h[t] = o[t] \odot \tanh(c[t])$$

