# Foundations of DL

Deep Learning

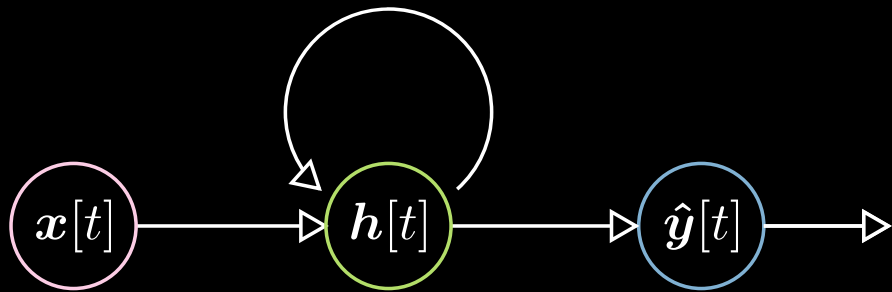Alfredo Canziani, Ritchie Ng

@alfcnz, @RitchieNg
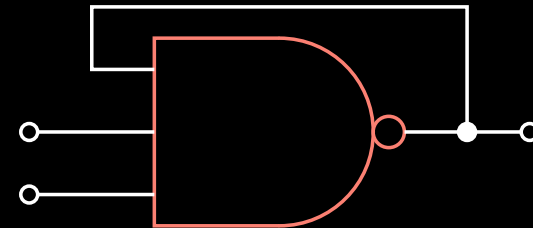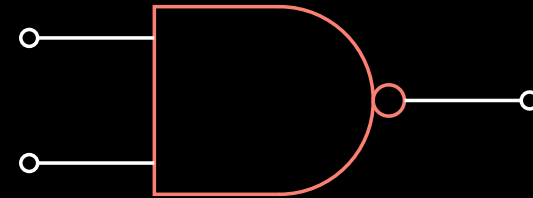
ALF

# Recurrent Neural Nets

Handling sequential data

# Vanilla and Recurrent NN

# Rationale
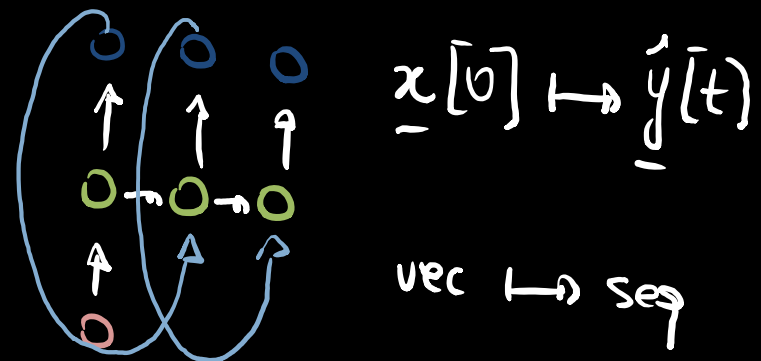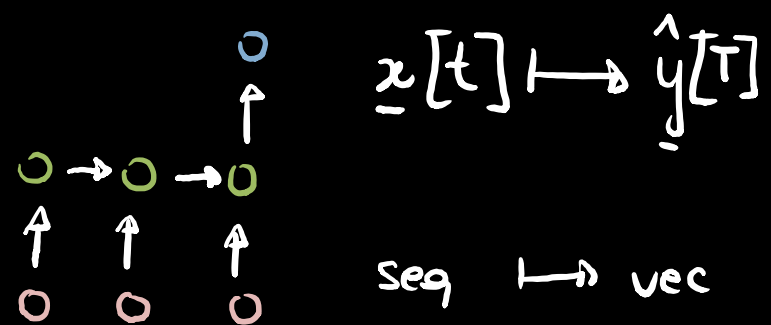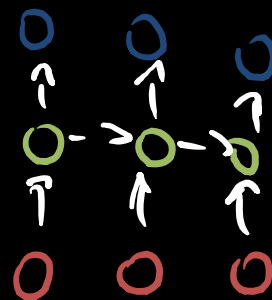


$$x[t] \longmapsto \hat{y}[T]$$

seq $\longmapsto$ vec

$$x[0] \longmapsto \hat{y}(t)$$

vec $\longmapsto$ seq

seq $\longmapsto$ seq

Vinyals *et al.* (2016) Show and Tell: Lessons learned from the 2015 MSCOCO Image Captioning Challenge

# 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.

Cho et al. (2014) Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation

Cho et al. (2014) Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation

Cho et al. (2014) Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation

test.txt      ●      rnn-client.coffee

```
1    The
```

Sloan (2016) rnn-writer – github.com/robinsloan/rnn-writer

# RNN training

Back propagation through time (BPTT)

# 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

# Get batch (I)



$$x[1:T]$$

BPTT period
T

$$y[1:T]$$

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

# Get batch (II)



$$y[1] \quad y[2] \quad y[T]$$

b h n t     c i o u     d j p v

$$h[1] \quad h[2]$$

RNN    RNN    RNN

$$h[0] \quad\quad\quad\quad\quad\quad\quad\quad h[T]$$

a g m s     b h n t     c i o u

$$x[1] \quad x[2] \quad x[T]$$

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

# Vanishing & exploding gradients

Limitations of temporally deep nets

Graves (2012) Supervised sequence labelling

Outputs

Hidden Layer

Inputs

Time   1   2   3   4   5   6   7

Graves (2012) Supervised sequence labelling

# Long Short-Term Memory

Gated RNN

$$\boldsymbol{x}[t]$$
$$\boldsymbol{h}[t-1]$$

th

$$\boldsymbol{h}[t]$$

th

$$\boldsymbol{y}[t]$$

$$\boldsymbol{h}[t] = g\big(\boldsymbol{W}_h \big[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\big] + \boldsymbol{b}_h\big)$$

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

$$\boldsymbol{i}[t] = \sigma\big(\boldsymbol{W}_i \big[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\big] + \boldsymbol{b}_i\big)$$

$$\boldsymbol{f}[t] = \sigma\big(\boldsymbol{W}_f \big[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\big] + \boldsymbol{b}_f\big)$$

$$\boldsymbol{o}[t] = \sigma\big(\boldsymbol{W}_o \big[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\big] + \boldsymbol{b}_o\big)$$
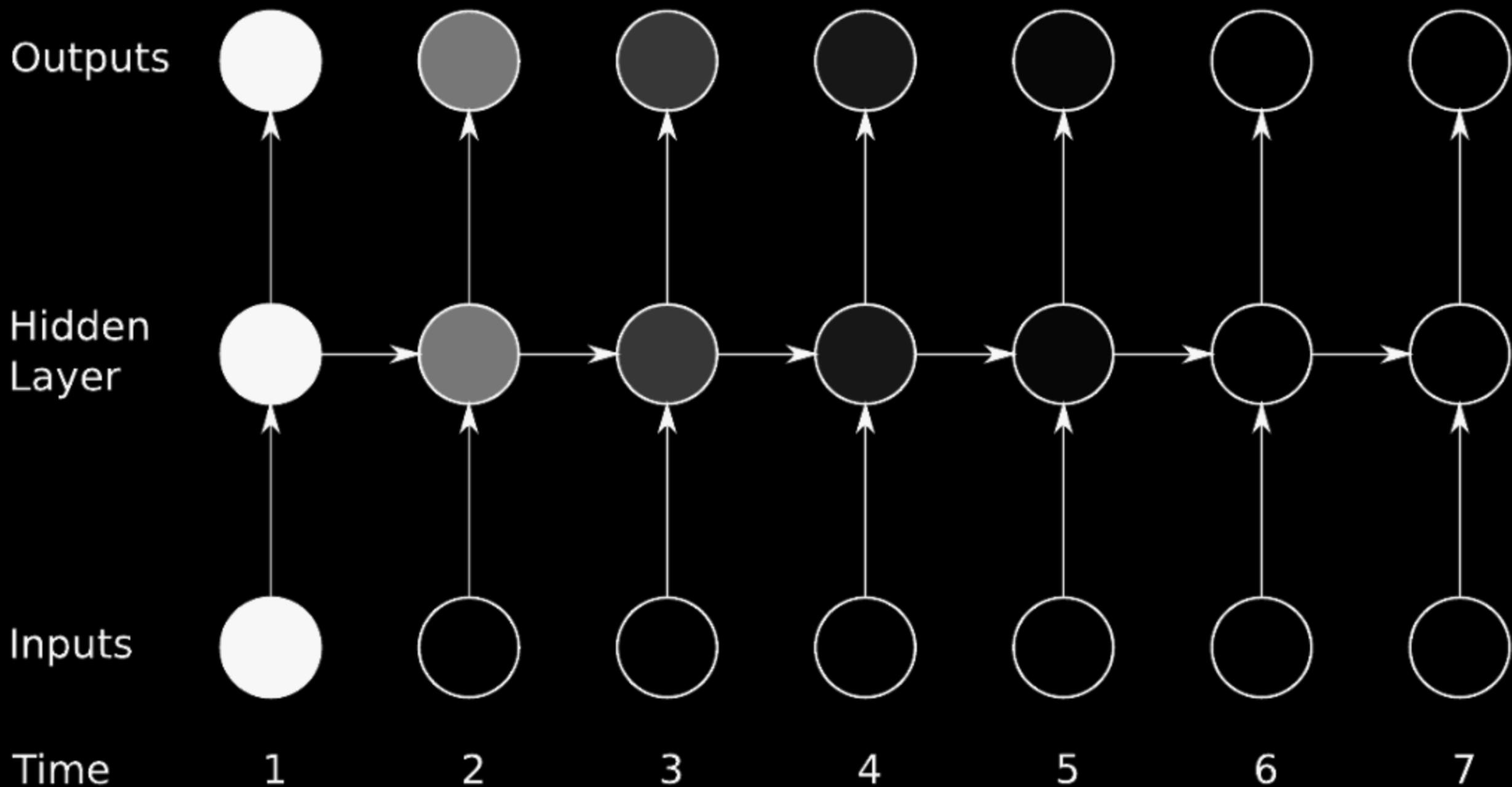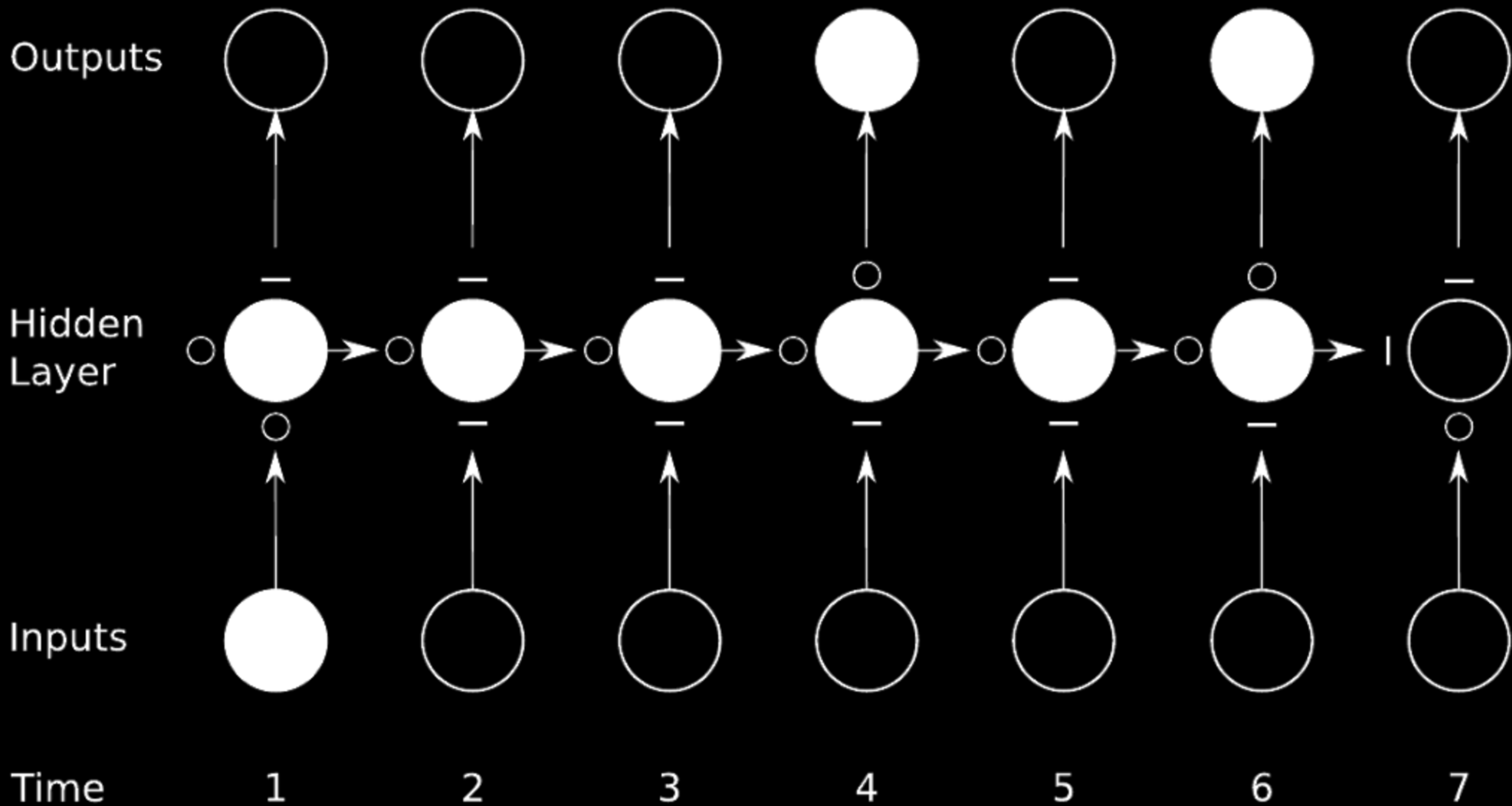
$$\tilde{\boldsymbol{c}}[t] = \tanh\big(\boldsymbol{W}_c \big[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\big] + \boldsymbol{b}_c\big)$$

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

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

input

output

$$\sigma$$

$$\sigma$$

$$\boldsymbol{i}[t]$$

$$\tilde{\boldsymbol{c}}[t]$$

$$\boldsymbol{c}[t]$$

$$\boldsymbol{o}[t]$$

$$\boldsymbol{h}[t]$$

th

th

$$\boldsymbol{f}[t]$$

$$\sigma$$

$$\boldsymbol{x}[t]$$
$$\boldsymbol{h}[t-1]$$

$$\boldsymbol{c}[t-1]$$

don't forget

$$\boldsymbol{c}[t]$$

# Controlling the  output  - ON

Saturated sigmoid

⬤ = 1

⬤ = 0

$$i[t] = \sigma\left(\boldsymbol{W}_i\left[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\right] + \boldsymbol{b}_i\right)$$

$$\boldsymbol{f}[t] = \sigma\left(\boldsymbol{W}_f\left[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\right] + \boldsymbol{b}_f\right)$$

$$\boldsymbol{o}[t] = \sigma\left(\boldsymbol{W}_o\left[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\right] + \boldsymbol{b}_o\right)$$

$$\tilde{\boldsymbol{c}}[t] = \tanh\left(\boldsymbol{W}_c\left[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\right] + \boldsymbol{b}_c\right)$$

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

$$\boldsymbol{h}[t] = \boxed{\boldsymbol{o}[t] \odot \tanh(\boldsymbol{c}[t])}$$

# Controlling the memory - reset

Saturated sigmoid
- 🟢 = 1
- 🔴 = 0

$$\boldsymbol{i}[t] = \sigma\big(\boldsymbol{W}_i\big[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\big] + \boldsymbol{b}_i\big)$$

$$\boldsymbol{f}[t] = \sigma\big(\boldsymbol{W}_f\big[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\big] + \boldsymbol{b}_f\big)$$

$$\boldsymbol{o}[t] = \sigma\big(\boldsymbol{W}_o\big[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\big] + \boldsymbol{b}_o\big)$$

$$\tilde{\boldsymbol{c}}[t] = \tanh\big(\boldsymbol{W}_c\big[\begin{smallmatrix}\boldsymbol{x}[t]\\\boldsymbol{h}[t-1]\end{smallmatrix}\big] + \boldsymbol{b}_c\big)$$

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

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

# Controlling the memory - keep

Saturated sigmoid
- ● = 1
- ● = 0

$$i[t] = \sigma\left(W_i\begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_i\right)$$

$$f[t] = \sigma\left(W_f\begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_f\right)$$

$$o[t] = \sigma\left(W_o\begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_o\right)$$

$$\tilde{c}[t] = \tanh\left(W_c\begin{bmatrix} x[t] \\ h[t-1] \end{bmatrix} + b_c\right)$$

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

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