

# Lecture 7: Recurrent Neural Networks

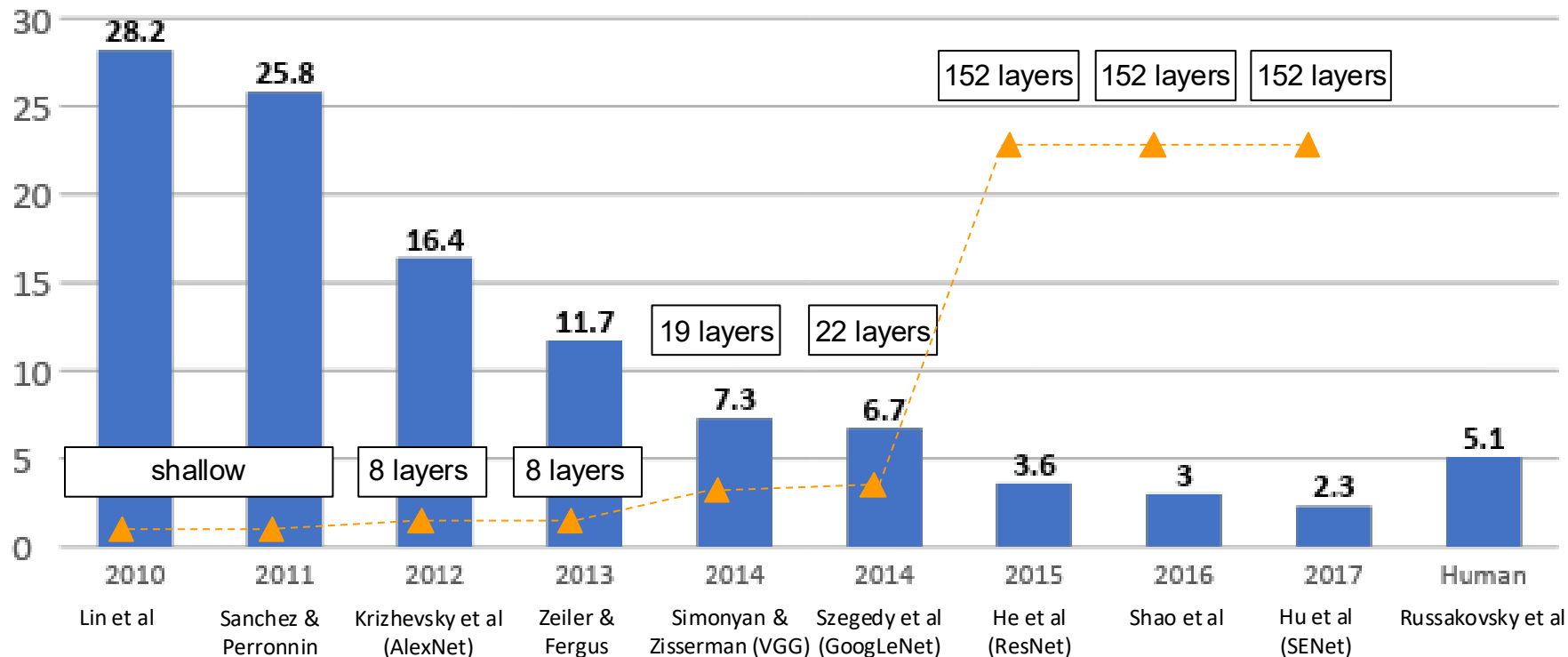
# Course Logistics

- Assignment 2 is going to be **out** this **Thursday** (4/23)!
- Project proposal deadline is **due** this **Thursday** (4/23)!

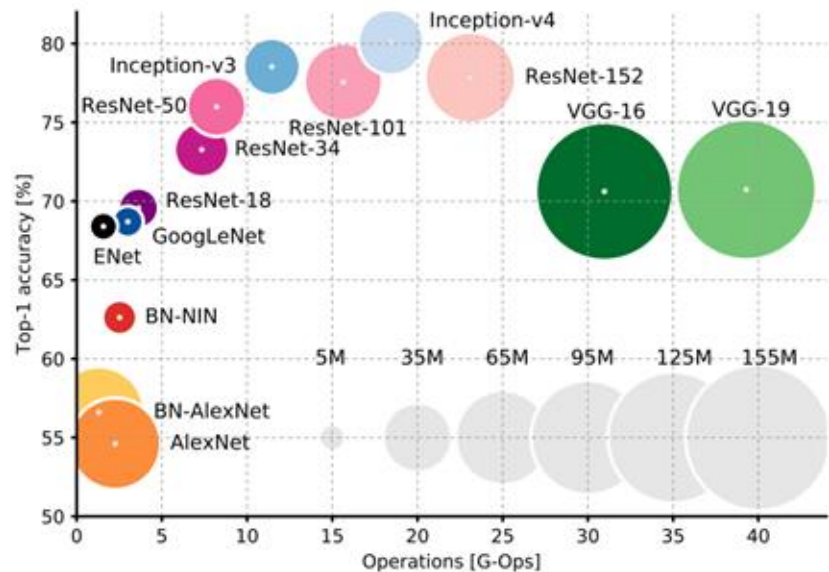
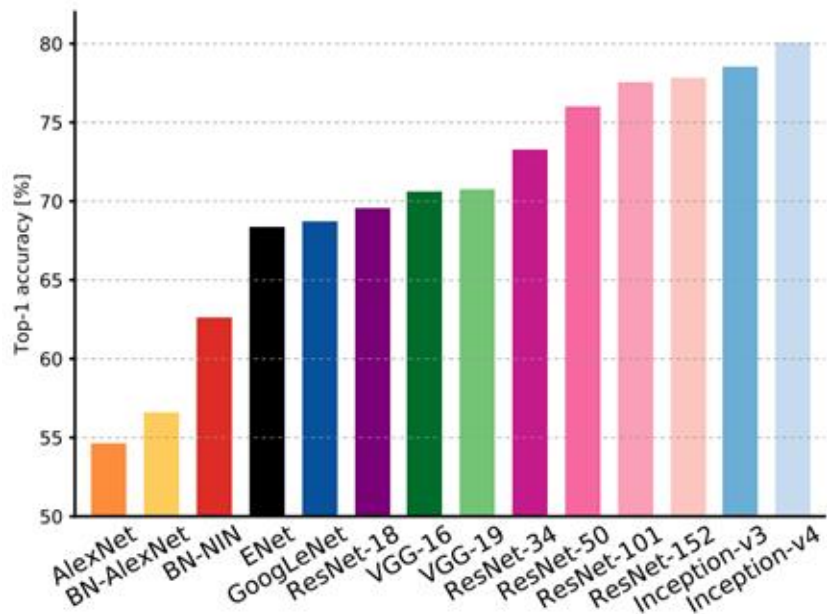
# Training “Feedforward” Neural Networks

- 1. One time setup:** activation functions, preprocessing, weight initialization, regularization, gradient checking
- 2. Training dynamics:** babysitting the learning process, parameter updates, hyperparameter optimization
- 3. Evaluation:** model ensembles, test-time augmentation, transfer learning

# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



# Comparing complexity...

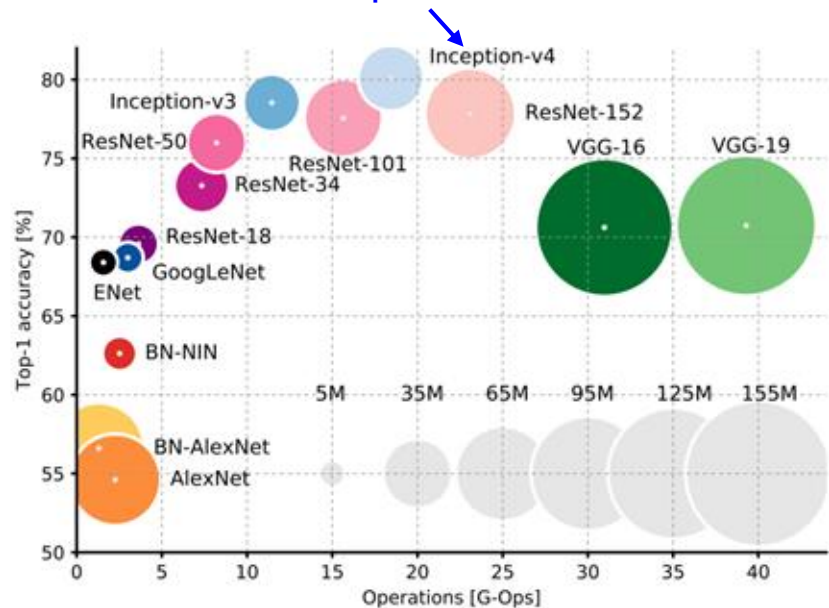
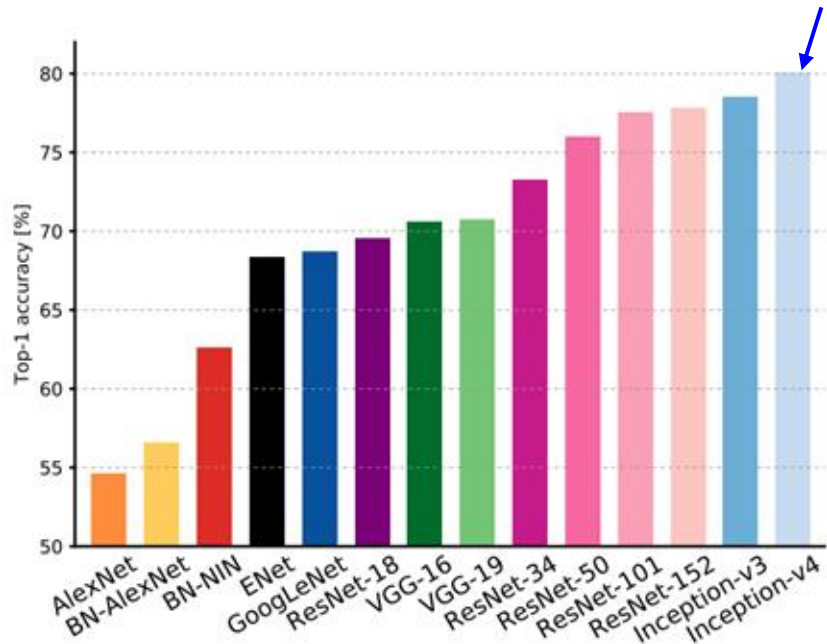


An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

# Comparing complexity...

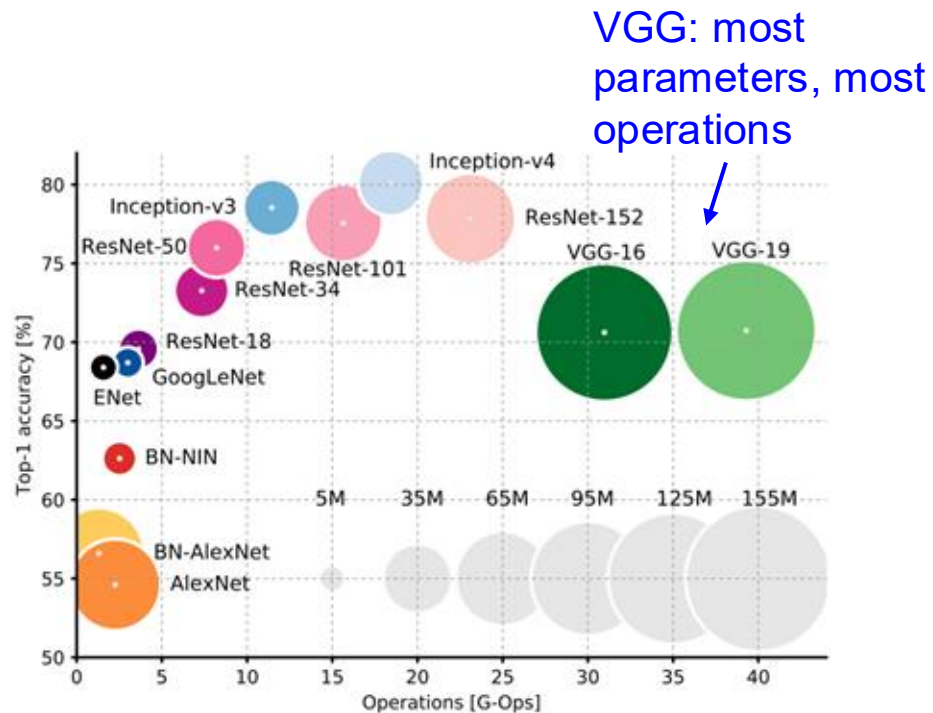
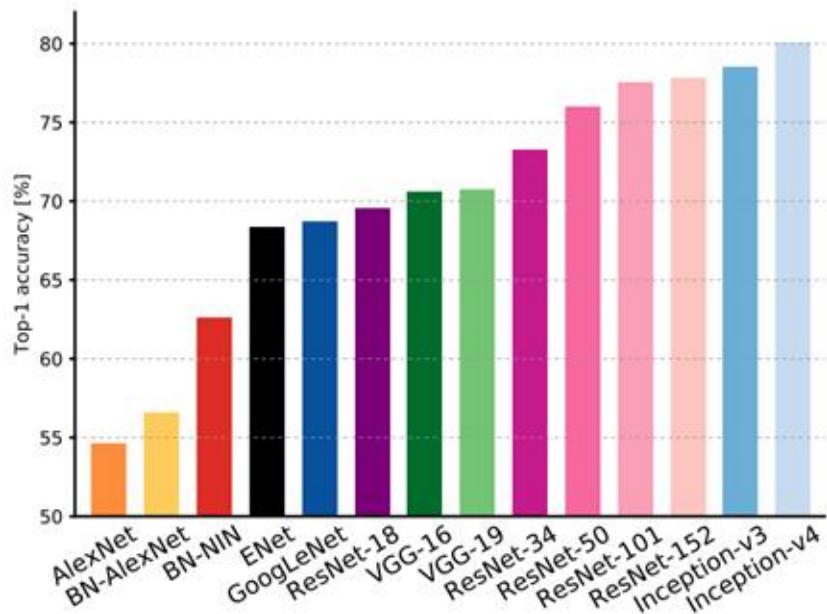
Inception-v4: Resnet + Inception!



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

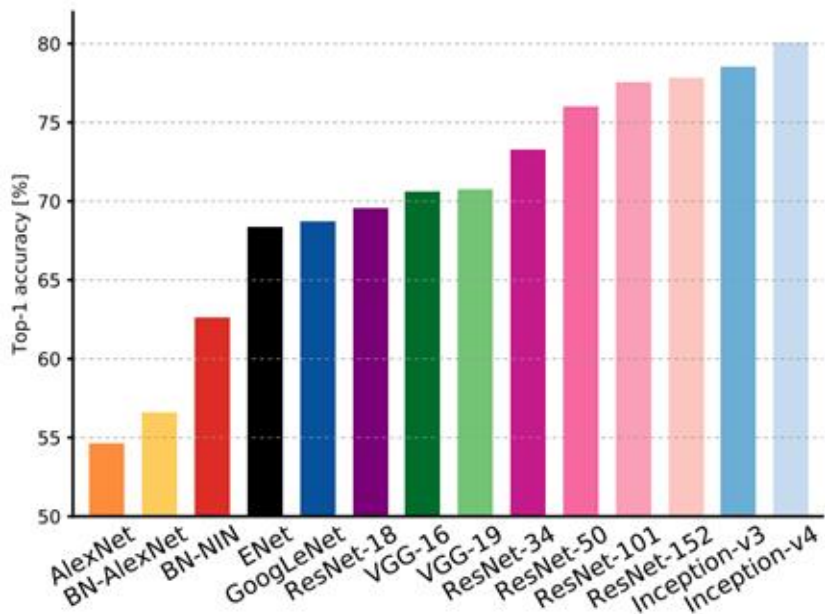
# Comparing complexity...



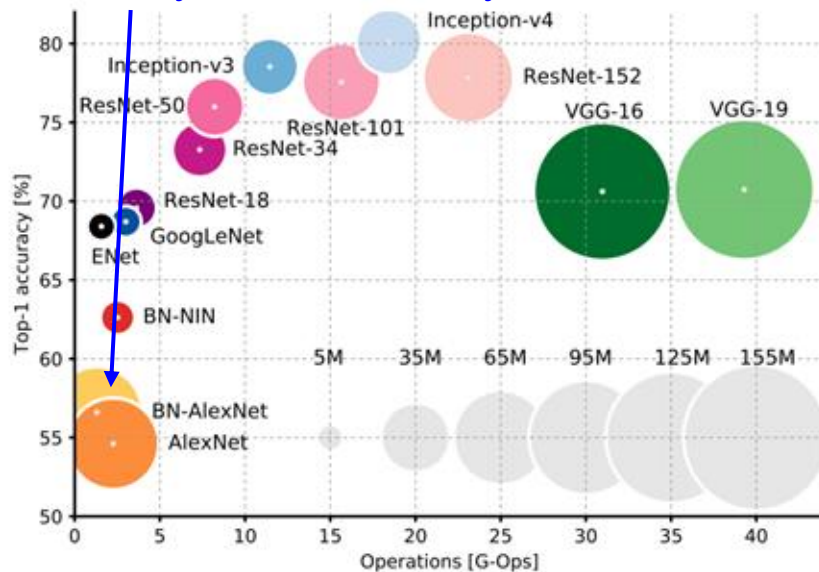
An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

# Comparing complexity...



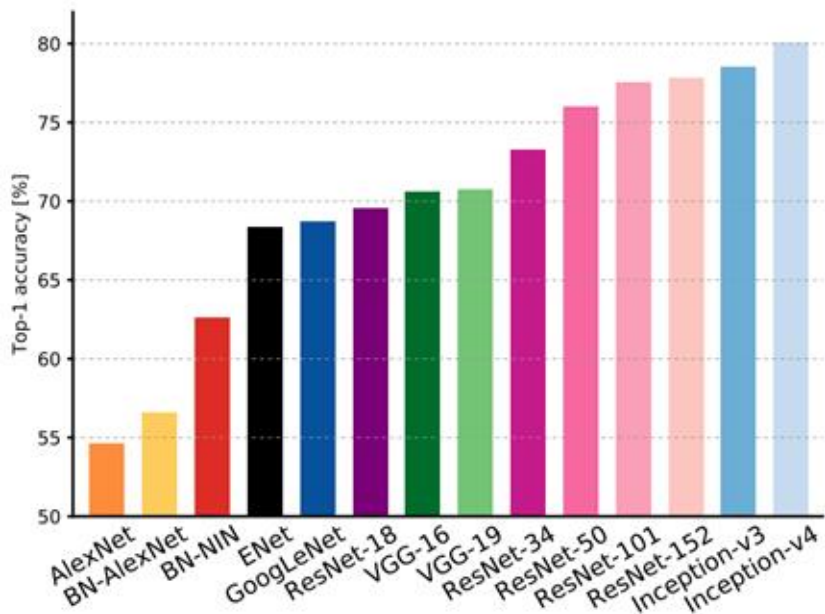
AlexNet:  
Smaller compute, still memory heavy, lower accuracy



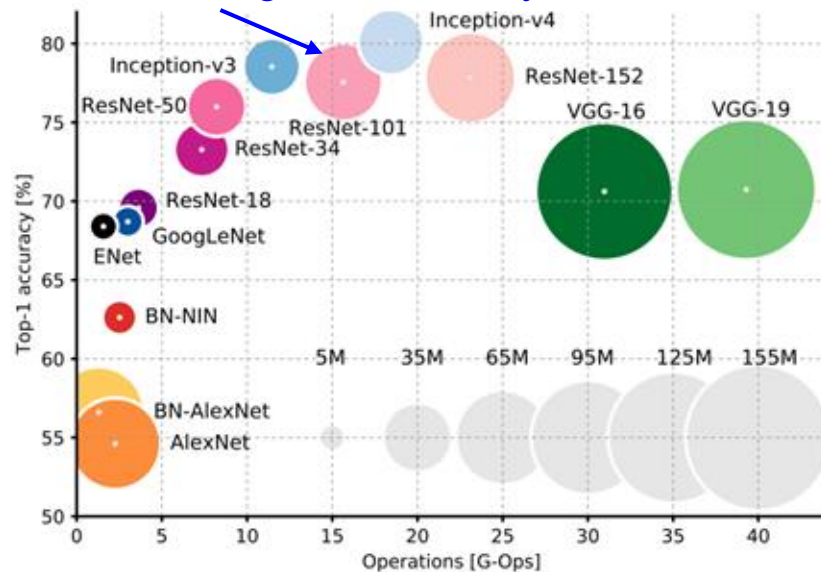
An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

# Comparing complexity...



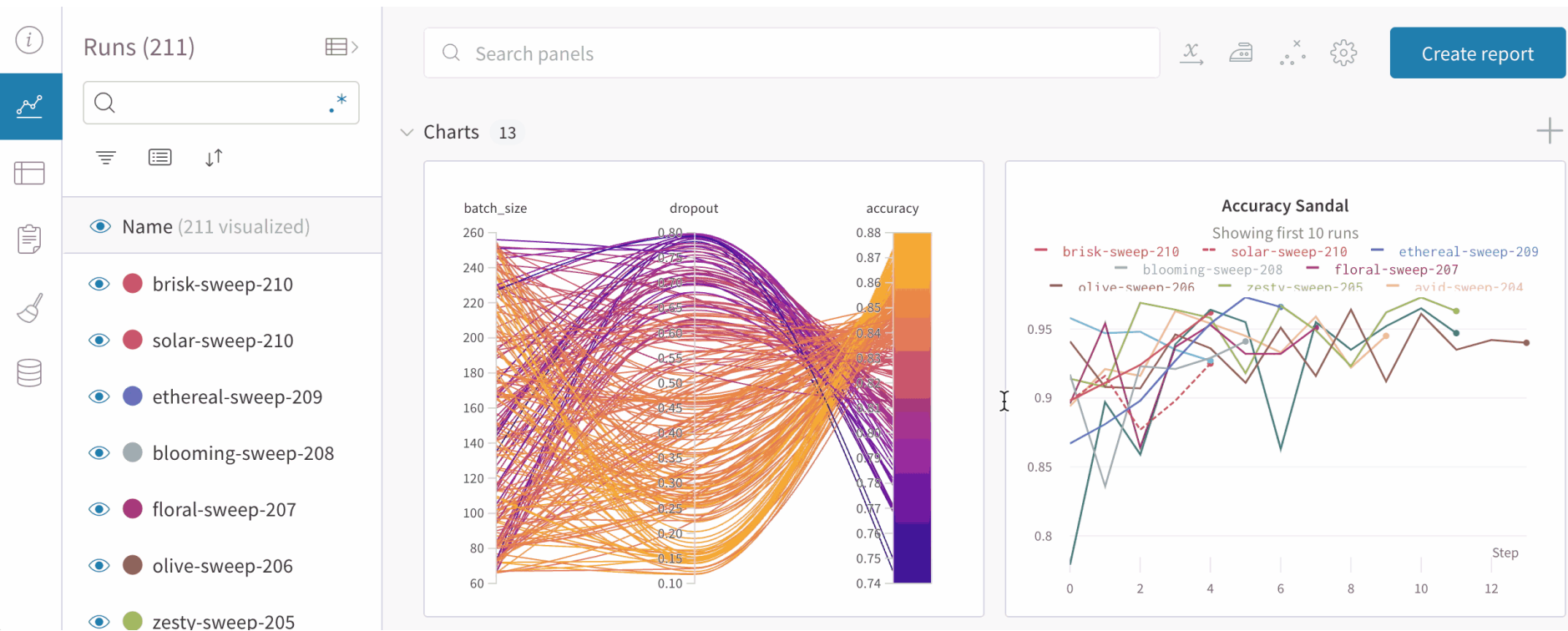
ResNet:  
Moderate efficiency depending on model, highest accuracy



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

# Evaluate models and tune hyperparameters



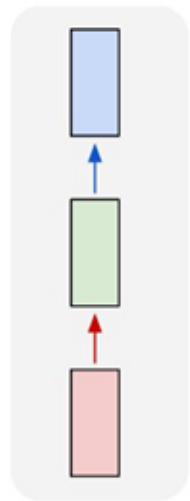
<https://docs.wandb.ai/guides/track/app>

# Today

- Recurrent Neural Networks (RNNs)
- Discuss sequence modeling (assumed fixed-length inputs so far)
- Simple models commonly used before the era of transformers
  - RNNs and some variants
- Relation to modern state-space models

# “Vanilla” Neural Network

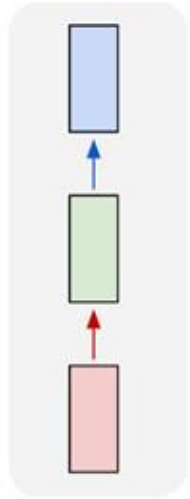
one to one



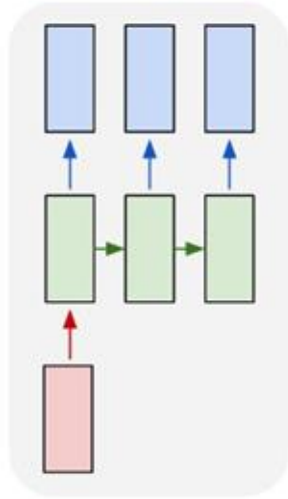
Vanilla Neural Networks

# Recurrent Neural Networks: Process Sequences

one to one



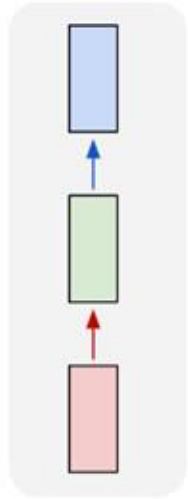
one to many



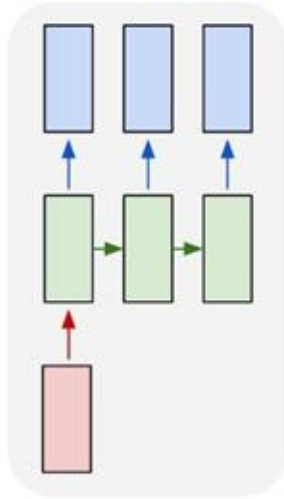
e.g. Image Captioning  
image -> sequence of words

# Recurrent Neural Networks: Process Sequences

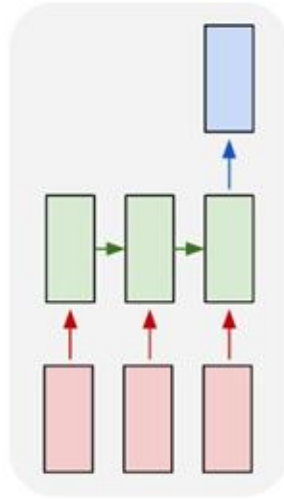
one to one



one to many



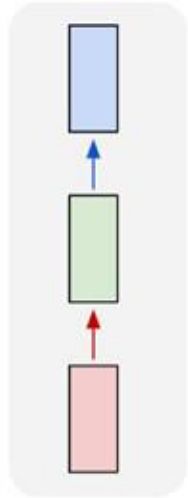
many to one



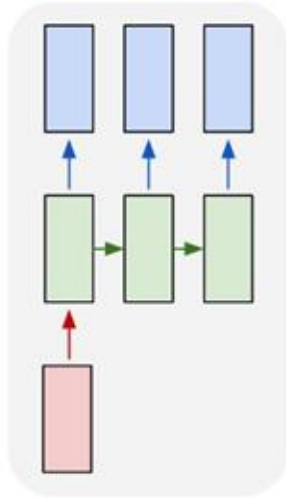
e.g. action prediction  
sequence of video frames -> action class

# Recurrent Neural Networks: Process Sequences

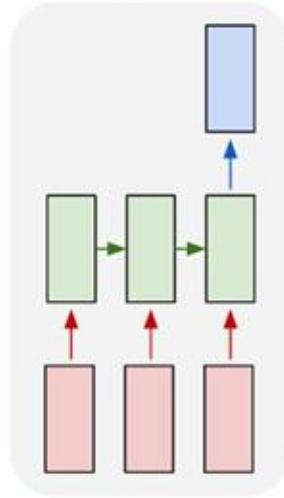
one to one



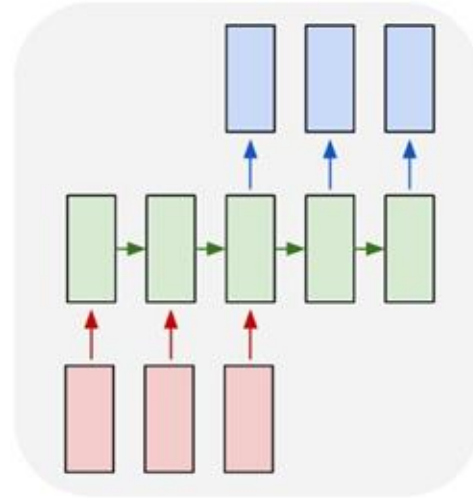
one to many



many to one



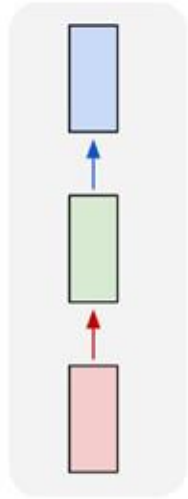
many to many



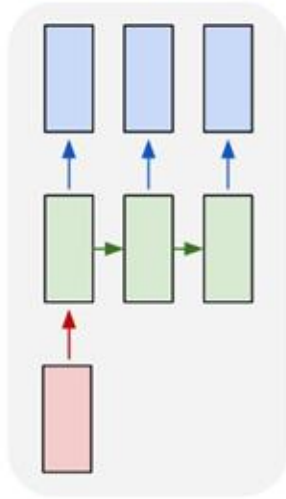
E.g. Video Captioning  
Sequence of video frames -> caption

# Recurrent Neural Networks: Process Sequences

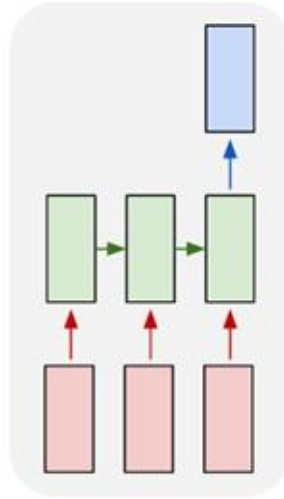
one to one



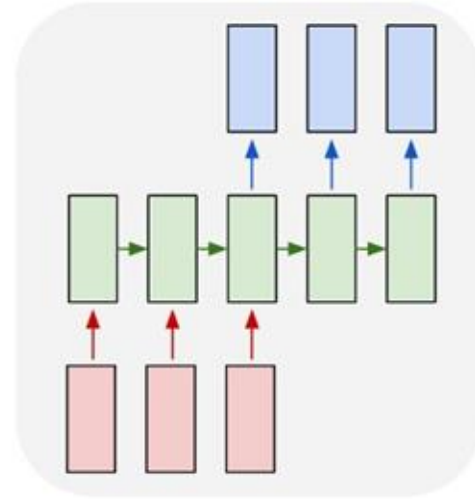
one to many



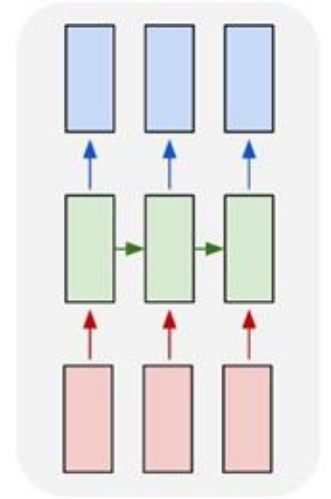
many to one



many to many



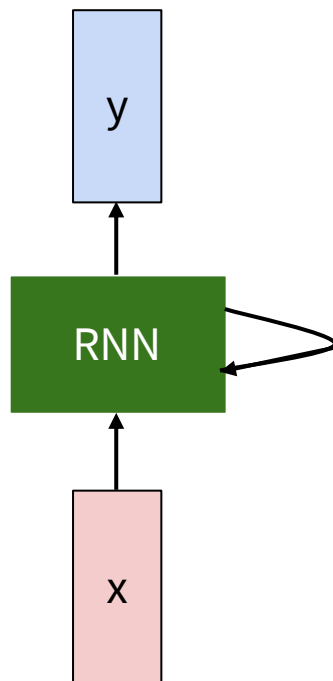
many to many



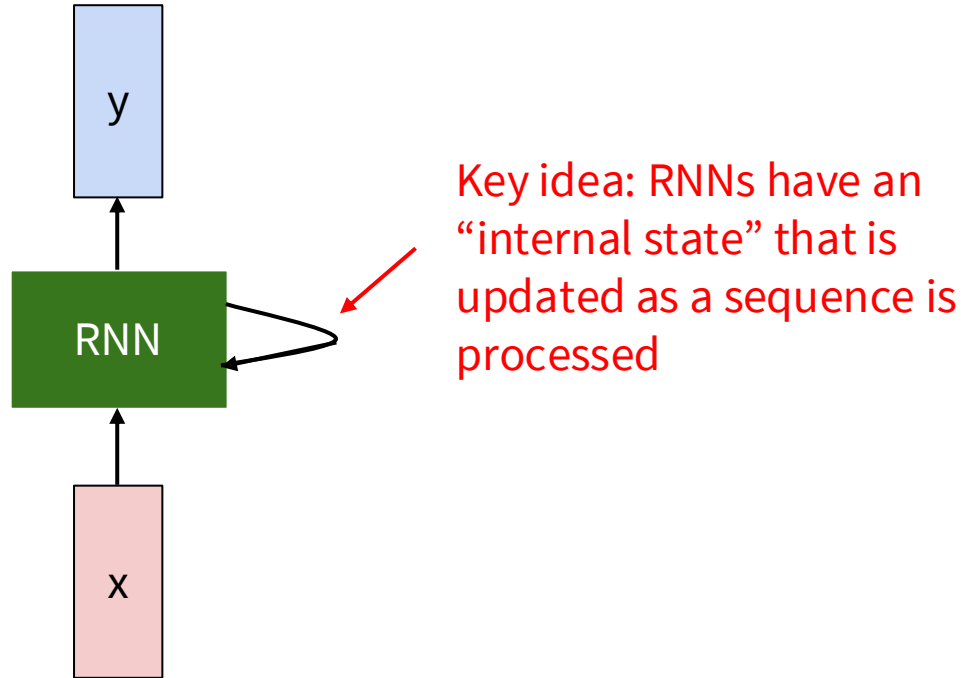
e.g. Video classification on frame level



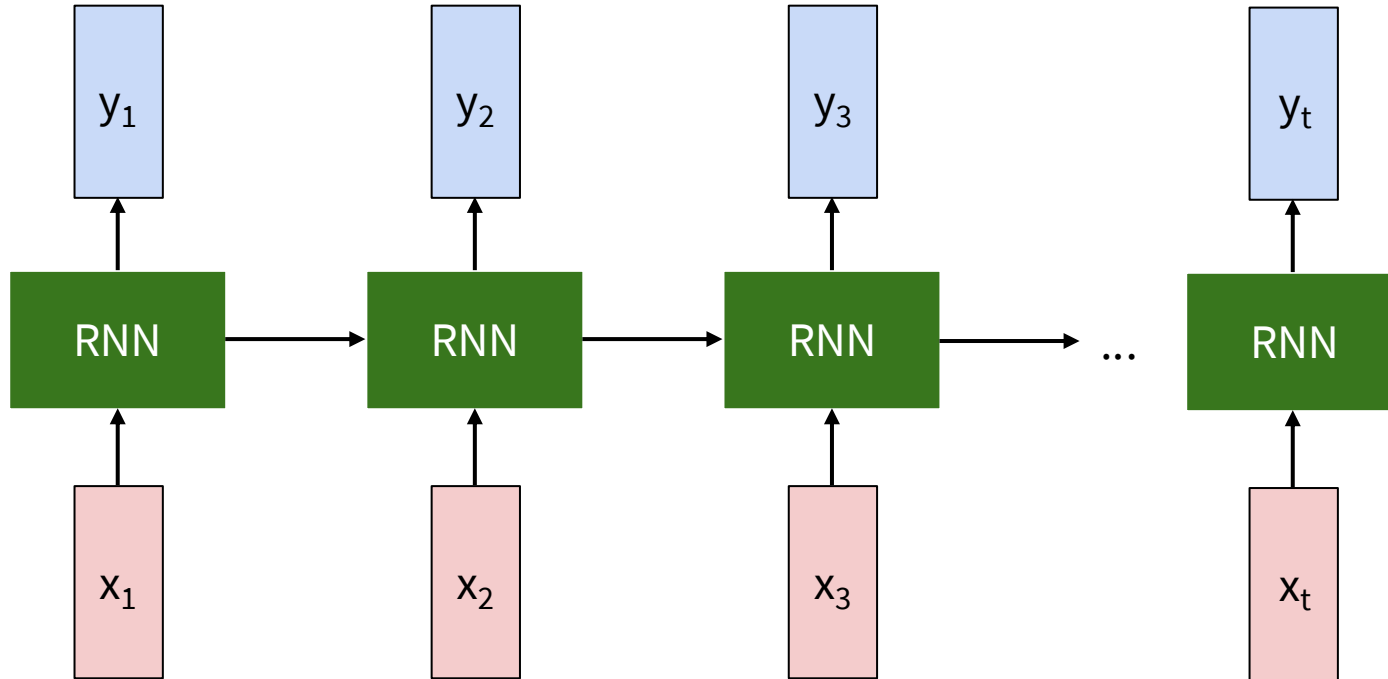
# Recurrent Neural Network



# Recurrent Neural Network



# Unrolled RNN



# RNN hidden state update

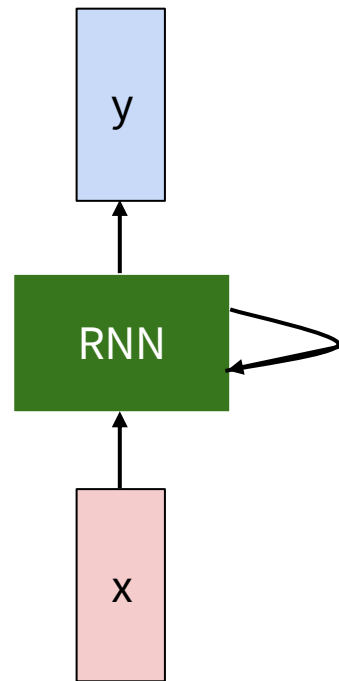
We can process a sequence of vectors  $x$  by applying a recurrence formula at every time step:

$$\boxed{h_t} = \boxed{f_W}(\boxed{h_{t-1}}, \boxed{x_t})$$

new state / some function with parameters  $W$

old state

input vector at some time step



# RNN output generation

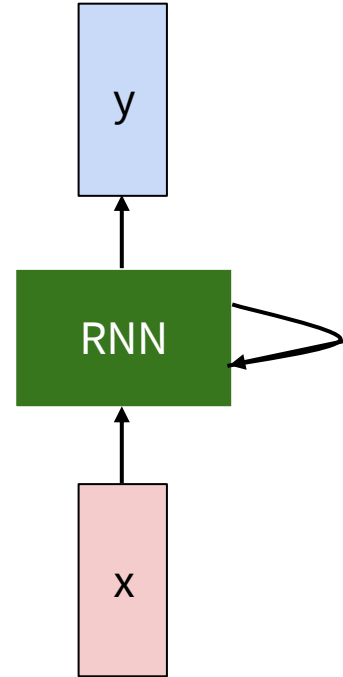
We can process a sequence of vectors  $x$  by applying a recurrence formula at every time step:

$$\boxed{y_t} = \boxed{f_{W_{hy}}}(\boxed{h_t})$$

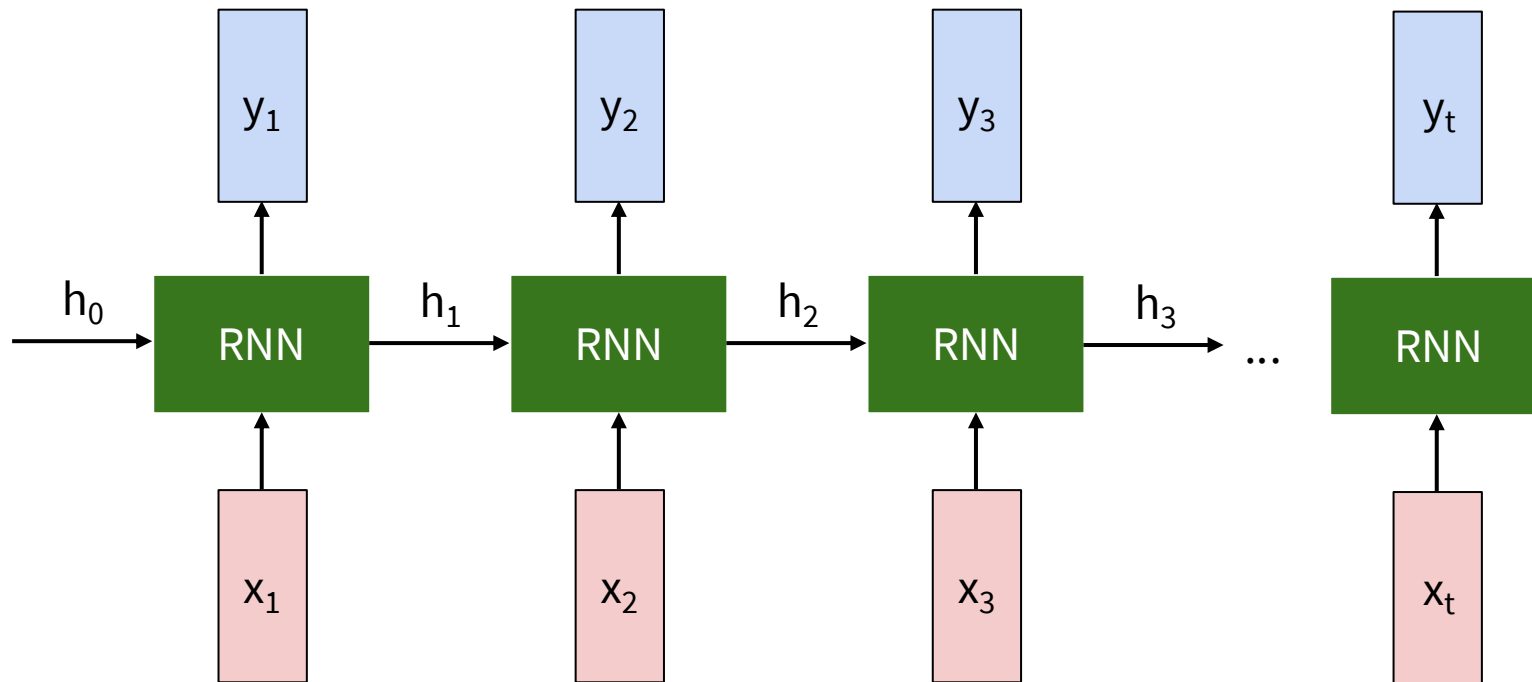
output

another function with parameters  $W_{hy}$

new state



# Recurrent Neural Network

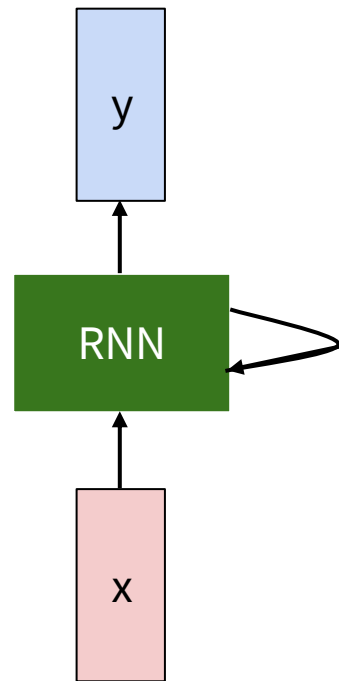


# Recurrent Neural Network

We can process a sequence of vectors  $x$  by applying a recurrence formula at every time step:

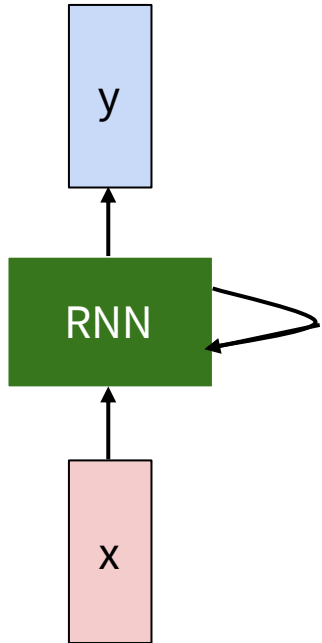
$$h_t = f_W(h_{t-1}, x_t)$$

Notice: the same function and the same set of parameters are used at every time step.



# (Vanilla) Recurrent Neural Network

The state consists of a single “hidden” vector  $h$ :



$$h_t = f_W(h_{t-1}, x_t)$$



$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

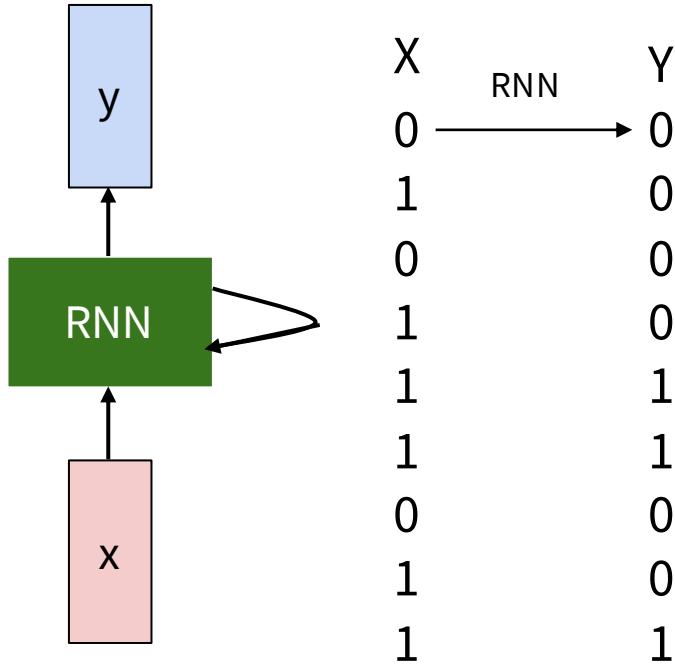
$$y_t = W_{hy}h_t$$

Often, we also have an output function:  $f_y$

Sometimes called a “Vanilla RNN” or an “Elman RNN” after Prof. Jeffrey Elman

# Vanilla(-ish) RNN: Concrete Example

Manually creating a recurrent network for detecting repeated 1s



“Many to many” sequence modeling task

$$h_t = f_W(h_{t-1}, x_t)$$

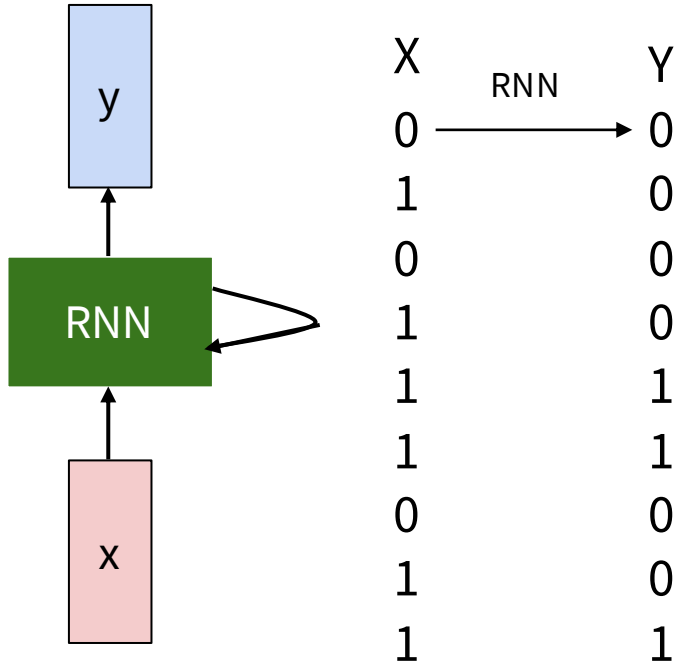
$$y_t = f_y(W_{hy}h_t)$$

How could we create an RNN to do this?

Q: What information should be captured in the “hidden” state?

# Vanilla(-ish) RNN: Concrete Example

Manually creating a recurrent network for detecting repeated 1s



“Many to many” sequence modeling task

$$h_t = f_W(h_{t-1}, x_t)$$

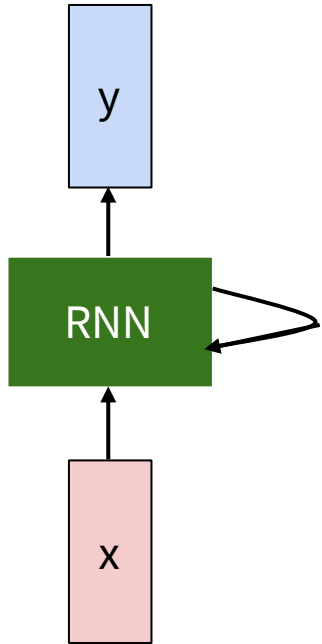
$$y_t = f_y(W_{hy}h_t)$$

How could we create an RNN to do this?

A: Previous input and current value for x

# Vanilla(-ish) RNN: Concrete Example

Manually creating a recurrent network for detecting repeated 1s



| X | RNN | Y |
|---|-----|---|
| 0 | →   | 0 |
| 1 |     | 0 |
| 0 |     | 0 |
| 1 |     | 0 |
| 1 |     | 1 |
| 1 |     | 1 |
| 0 |     | 0 |
| 1 |     | 0 |
| 1 |     | 1 |

“Many to many” sequence modeling task

$$h_t = f_W(h_{t-1}, x_t)$$

$$y_t = f_y(W_{hy}h_t)$$

Set  $f_W$  and  $f_y$  to both be ReLU for simplicity

$$h_t = \begin{pmatrix} \text{Current} \\ \text{Previous} \\ 1 \end{pmatrix}$$

Initialize  $h_0$  to  $(0, 0, 1)$

# Vanilla(-ish) RNN: Concrete Example

Manually creating a recurrent network for detecting repeated 1s

```
w_xh = np.array([[1], [0], [0]])  
  
w_hh = np.array([[0, 0, 0],  
                 [1, 0, 0],  
                 [0, 0, 1]])  
  
w_yh = np.array([1, 1, -1])  
  
x_seq = [0, 1, 0, 1, 1, 1, 0, 1, 1]  
  
h_t_prev = np.array([[0], [0], [1]])  
  
for t, x in enumerate(x_seq):  
    h_t = relu(w_hh @ h_t_prev + (w_xh @ x))  
    y_t = relu(w_yh @ h_t)  
    h_t_prev = h_t
```

$$h_t = \text{ReLU}(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = \text{ReLU}(W_{hy}h_t)$$

$$h_t = \begin{pmatrix} \text{Current} \\ \text{Previous} \\ 1 \end{pmatrix}$$

| X | RNN | Y |
|---|-----|---|
| 0 | →   | 0 |
| 1 |     | 0 |
| 0 |     | 0 |
| 1 |     | 0 |
| 1 |     | 1 |
| 1 |     | 1 |
| 1 |     | 1 |
| 0 |     | 0 |
| 1 |     | 0 |
| 1 |     | 1 |

\* Code is missing parts, for full code see [here](#)

# Vanilla(-ish) RNN: Concrete Example

Manually creating a recurrent network for detecting repeated 1s

```
w_xh = np.array([[1], [0], [0]])
```

```
w_hh = np.array([[0, 0, 0],  
                 [1, 0, 0],  
                 [0, 0, 1]])
```

```
w_yh = np.array([1, 1, -1])
```

```
x_seq = [0, 1, 0, 1, 1, 1, 0, 1, 1]
```

```
h_t_prev = np.array([[0], [0], [1]])
```

```
for t, x in enumerate(x_seq):  
    h_t = relu(w_hh @ h_t_prev + (w_xh @ x))  
    y_t = relu(w_yh @ h_t)  
    h_t_prev = h_t
```

$$h_t = \text{ReLU}(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = \text{ReLU}(W_{hy}h_t)$$

$$h_t = \begin{pmatrix} \text{Current} \\ \text{Previous} \\ 1 \end{pmatrix}$$

Right hand term

X=0 → [0, 0, 0]

X=1 → [1, 0, 0]

| X | RNN | Y |
|---|-----|---|
| 0 | →   | 0 |
| 1 |     | 0 |
| 0 |     | 0 |
| 1 |     | 0 |
| 1 |     | 1 |
| 1 |     | 1 |
| 0 |     | 0 |
| 1 |     | 0 |
| 1 |     | 1 |

\* Code is missing parts, for full code see [here](#)

# Vanilla(-ish) RNN: Concrete Example

Manually creating a recurrent network for detecting repeated 1s

```
w_xh = np.array([[1], [0], [0]])  
w_hh = np.array([[0, 0, 0],  
                 [1, 0, 0],  
                 [0, 0, 1]])  
w_yh = np.array([1, 1, -1])  
x_seq = [0, 1, 0, 1, 1, 1, 0, 1, 1]  
h_t_prev = np.array([[0], [0], [1]])  
  
for t, x in enumerate(x_seq):  
    h_t = relu(w_hh @ h_t_prev + (w_xh @ x))  
    y_t = relu(w_yh @ h_t)  
    h_t_prev = h_t
```

0 for top row of left term (only use value from right term)

$$h_t = \text{ReLU}(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = \text{ReLU}(W_{hy}h_t)$$

$$h_t = \begin{pmatrix} \text{Current} \\ \text{Previous} \\ 1 \end{pmatrix}$$

| X | → RNN → | Y |
|---|---------|---|
| 0 |         | 0 |
| 1 |         | 0 |
| 0 |         | 0 |
| 1 |         | 0 |
| 1 |         | 1 |
| 1 |         | 1 |
| 0 |         | 0 |
| 1 |         | 0 |
| 1 |         | 1 |

\* Code is missing parts, for full code see [here](#)

# Vanilla(-ish) RNN: Concrete Example

Manually creating a recurrent network for detecting repeated 1s

```
w_xh = np.array([[1], [0], [0]])  
w_hh = np.array([[0, 0, 0],  
                 [1, 0, 0],  
                 [0, 0, 1]])  
w_yh = np.array([1, 1, -1])  
x_seq = [0, 1, 0, 1, 1, 1, 0, 1, 1]  
h_t_prev = np.array([[0], [0], [1]])  
  
for t, x in enumerate(x_seq):  
    h_t = relu(w_hh @ h_t_prev + (w_xh @ x))  
    y_t = relu(w_yh @ h_t)  
    h_t_prev = h_t
```

Keep 1 on the bottom  
(helpful for output)

$$h_t = \text{ReLU}(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = \text{ReLU}(W_{hy}h_t)$$

$$h_t = \begin{pmatrix} \text{Current} \\ \text{Previous} \\ 1 \end{pmatrix}$$

| X | → RNN → | Y |
|---|---------|---|
| 0 |         | 0 |
| 1 |         | 0 |
| 0 |         | 0 |
| 1 |         | 0 |
| 1 |         | 1 |
| 1 |         | 1 |
| 0 |         | 0 |
| 1 |         | 0 |
| 1 |         | 1 |

\* Code is missing parts, for full code see [here](#)

# Vanilla(-ish) RNN: Concrete Example

Manually creating a recurrent network for detecting repeated 1s

```
w_xh = np.array([[1], [0], [0]])  
  
w_hh = np.array([[0, 0, 0],  
                [1, 0, 0],  
                [0, 0, 1]])  
  
w_yh = np.array([1, 1, -1])  
  
x_seq = [0, 1, 0, 1, 1, 1, 0, 1, 1]  
  
h_t_prev = np.array([[0], [0], [1]])  
  
for t, x in enumerate(x_seq):  
    h_t = relu(w_hh @ h_t_prev + (w_xh @ x))  
    y_t = relu(w_yh @ h_t)  
    h_t_prev = h_t
```

Max(Current + Previous - 1, 0)

$$h_t = \text{ReLU}(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = \text{ReLU}(W_{hy}h_t)$$

$$h_t = \begin{pmatrix} \text{Current} \\ \text{Previous} \\ 1 \end{pmatrix}$$

| X | RNN | Y |
|---|-----|---|
| 0 | →   | 0 |
| 1 |     | 0 |
| 0 |     | 0 |
| 1 |     | 0 |
| 1 |     | 1 |
| 1 |     | 1 |
| 1 |     | 1 |
| 0 |     | 0 |
| 1 |     | 0 |
| 1 |     | 1 |

\* Code is missing parts, for full code see [here](#)

# Vanilla(-ish) RNN: Concrete Example

Manually creating a recurrent network for detecting repeated 1s

```
w_xh = np.array([[1], [0], [0]])
```

```
w_hh = np.array([[0, 0, 0],  
                [1, 0, 0],  
                [0, 0, 1]])
```

```
w_yh = np.array([1, 1, -1])
```

```
x_seq = [0, 1, 0, 1, 1, 1, 0, 1, 1]
```

```
h_t_prev = np.array([[0], [0], [1]])
```

```
for t, x in enumerate(x_seq):  
    h_t = relu(w_hh @ h_t_prev + (w_xh @ x))  
    y_t = relu(w_yh @ h_t)  
    h_t_prev = h_t
```

$$h_t = \text{ReLU}(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = \text{ReLU}(W_{hy}h_t)$$

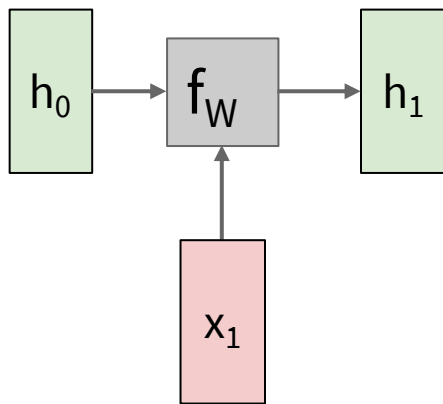
$$h_t = \begin{pmatrix} \text{Current} \\ \text{Previous} \\ 1 \end{pmatrix}$$

And it just works! But how do we find the  $W$ s in practice?

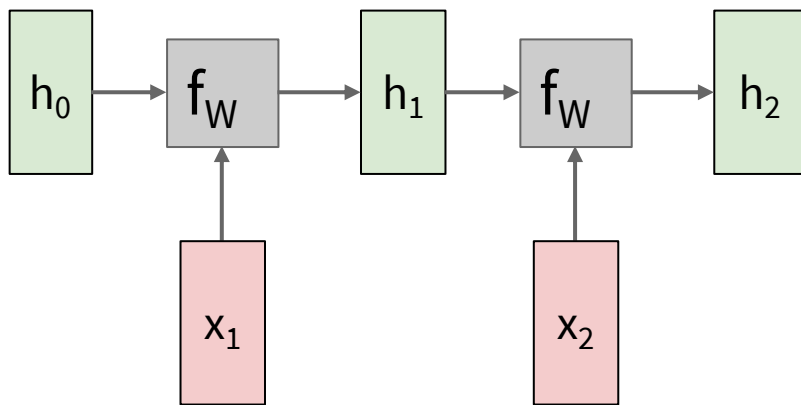
| X | → RNN → | Y |
|---|---------|---|
| 0 |         | 0 |
| 1 |         | 0 |
| 0 |         | 0 |
| 1 |         | 0 |
| 1 |         | 1 |
| 1 |         | 1 |
| 0 |         | 0 |
| 1 |         | 0 |
| 1 |         | 1 |

\* Code is missing parts, for full code see [here](#)

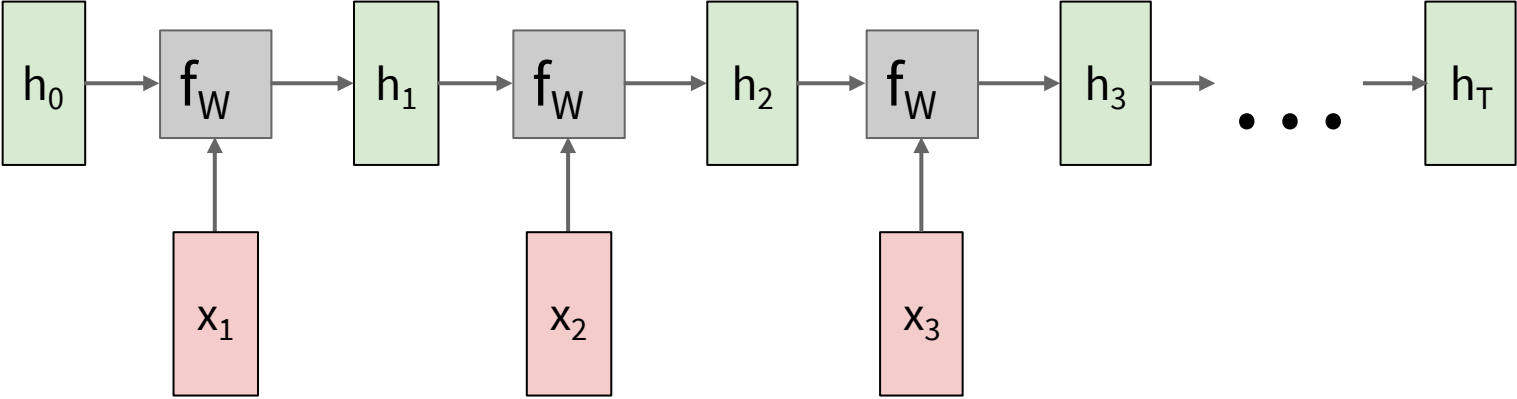
# RNN: Computational Graph



# RNN: Computational Graph

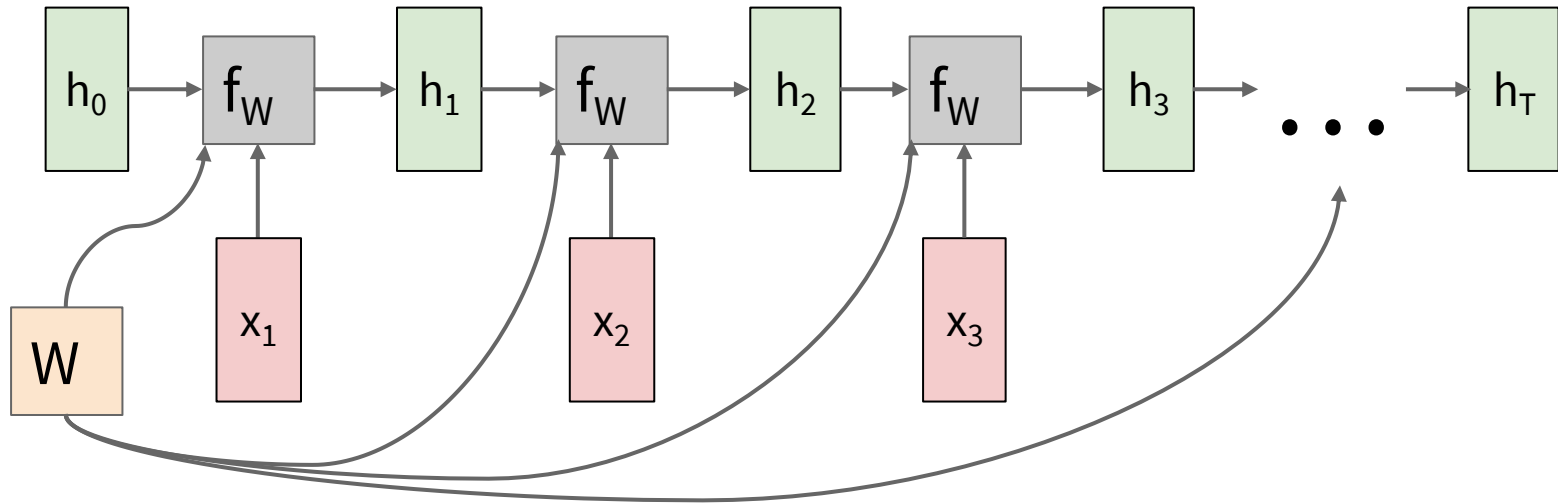


# RNN: Computational Graph

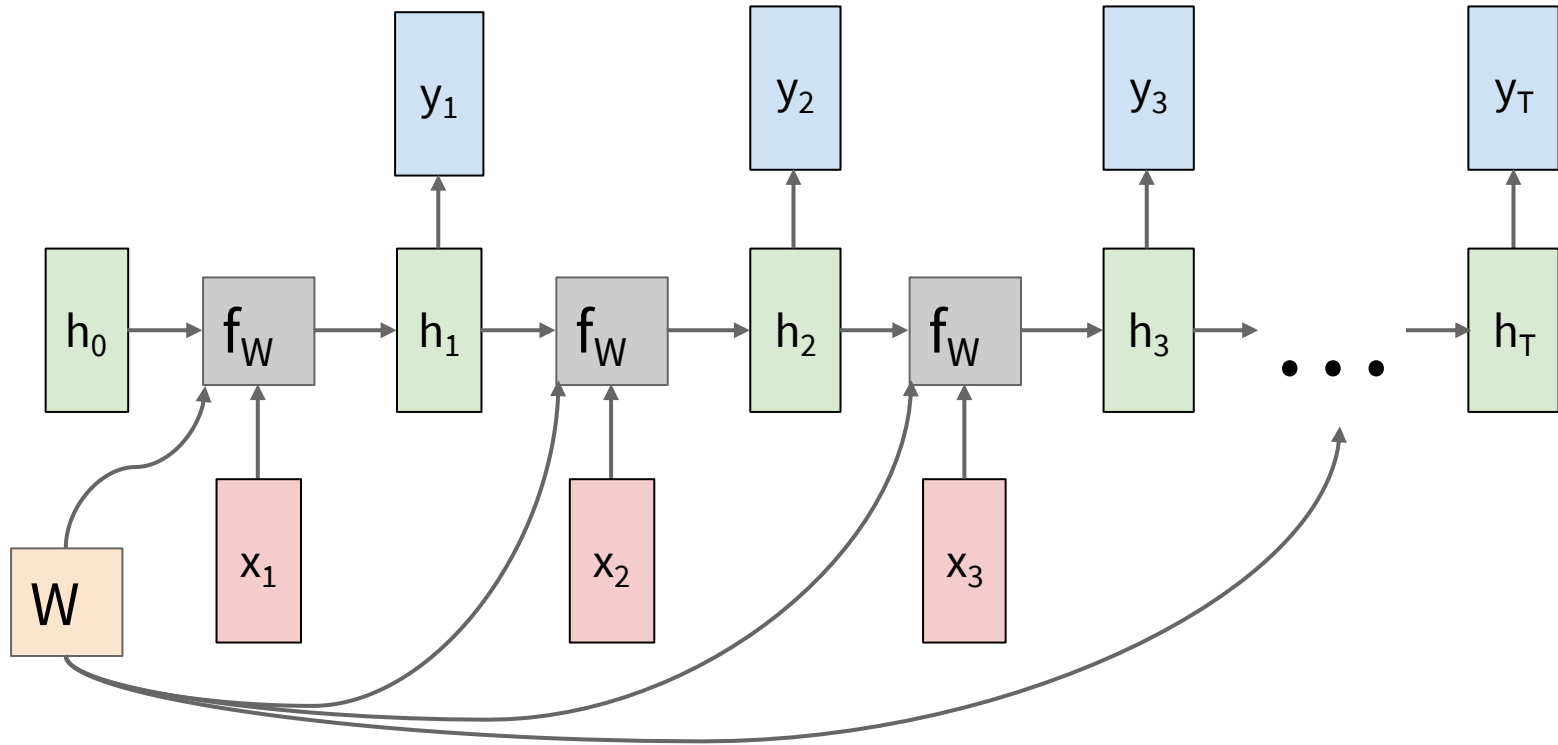


# RNN: Computational Graph

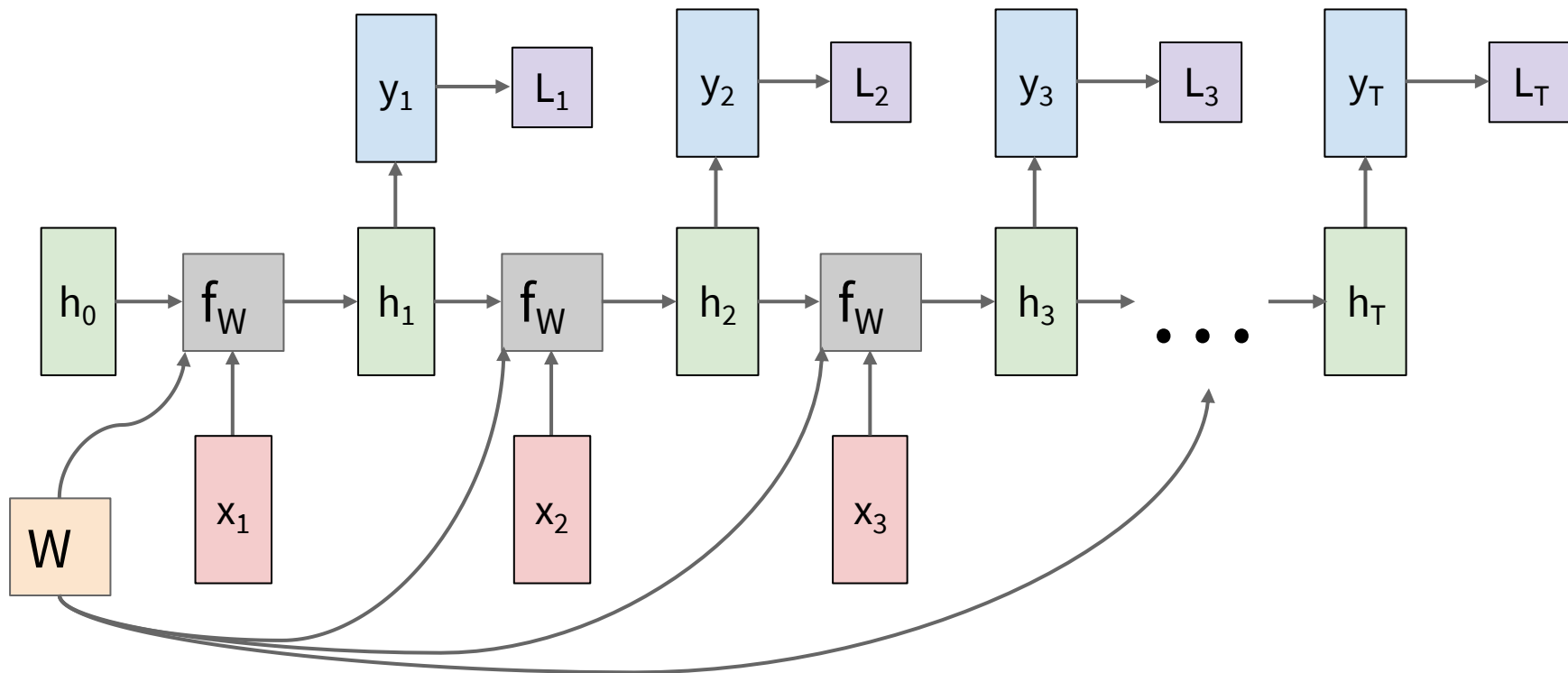
Re-use the same weight matrix at every time-step



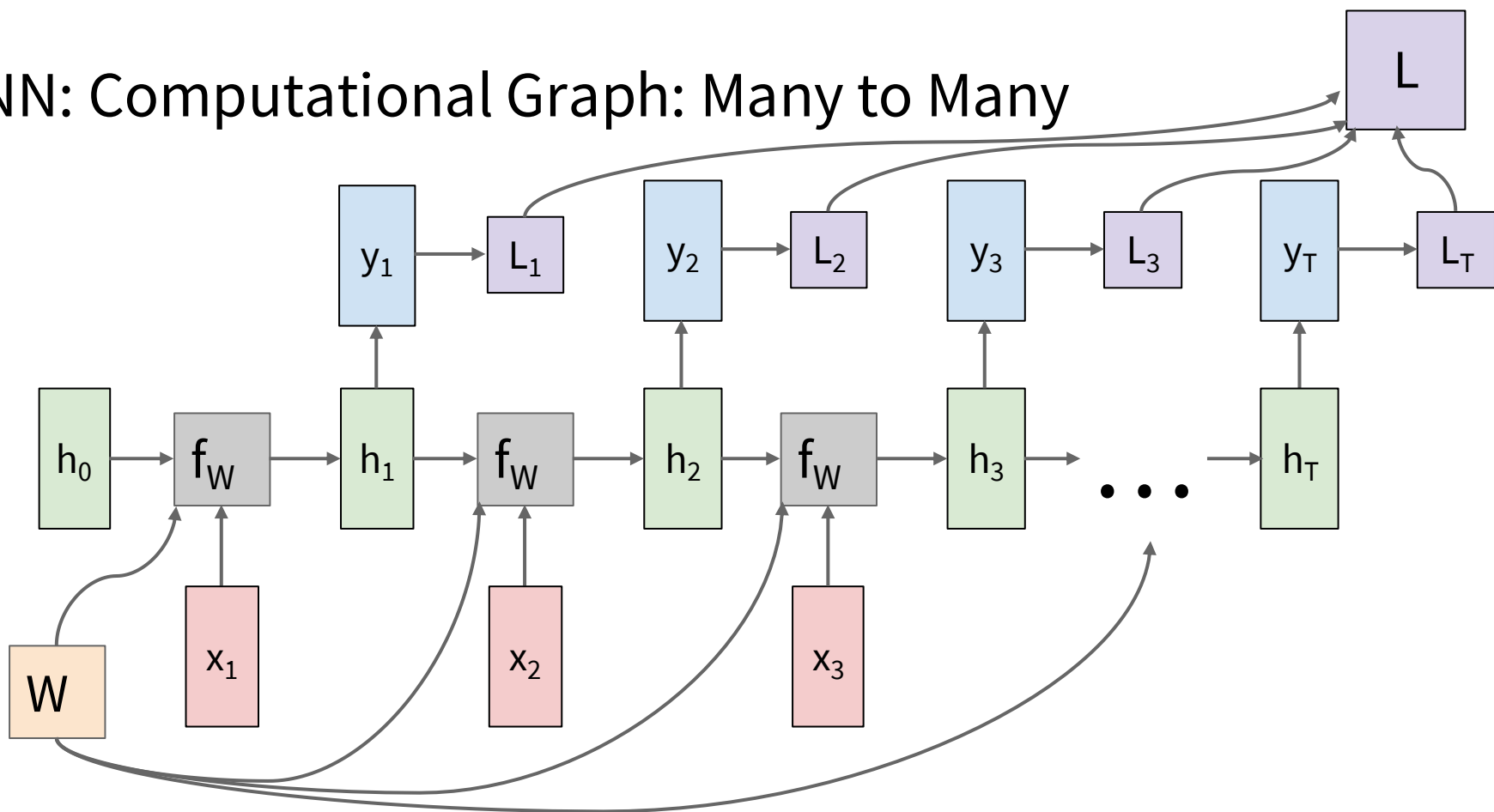
# RNN: Computational Graph: Many to Many



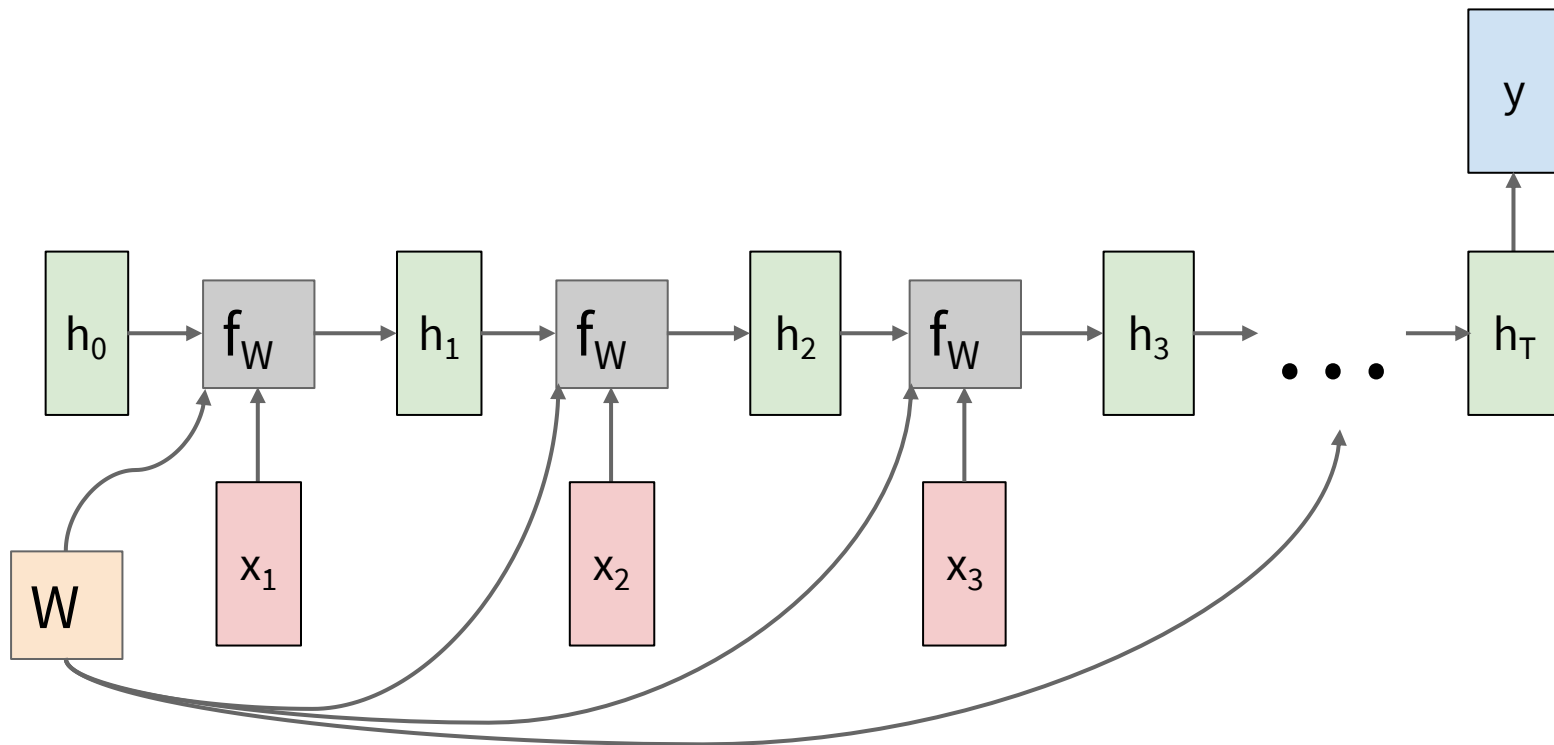
# RNN: Computational Graph: Many to Many



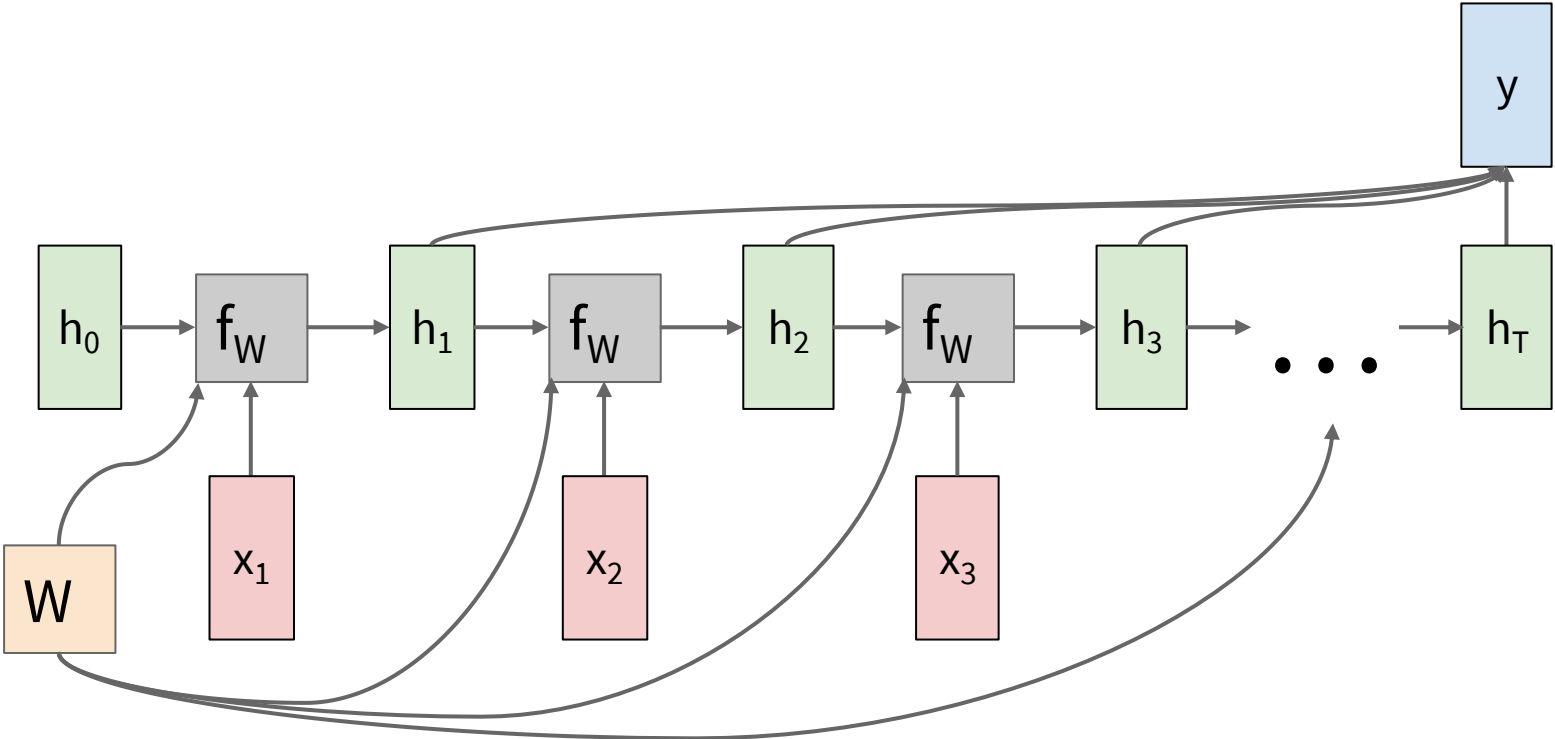
# RNN: Computational Graph: Many to Many



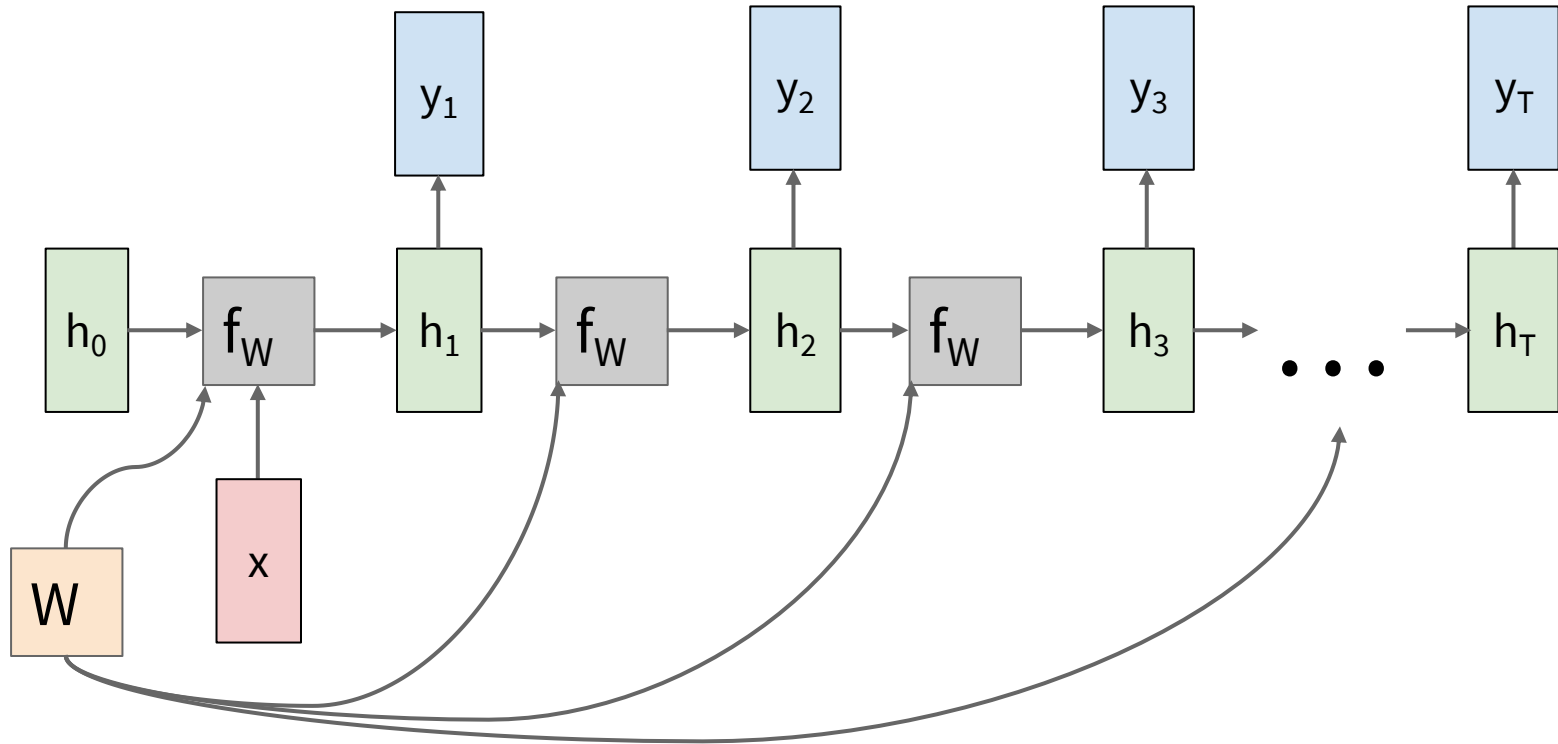
# RNN: Computational Graph: Many to One



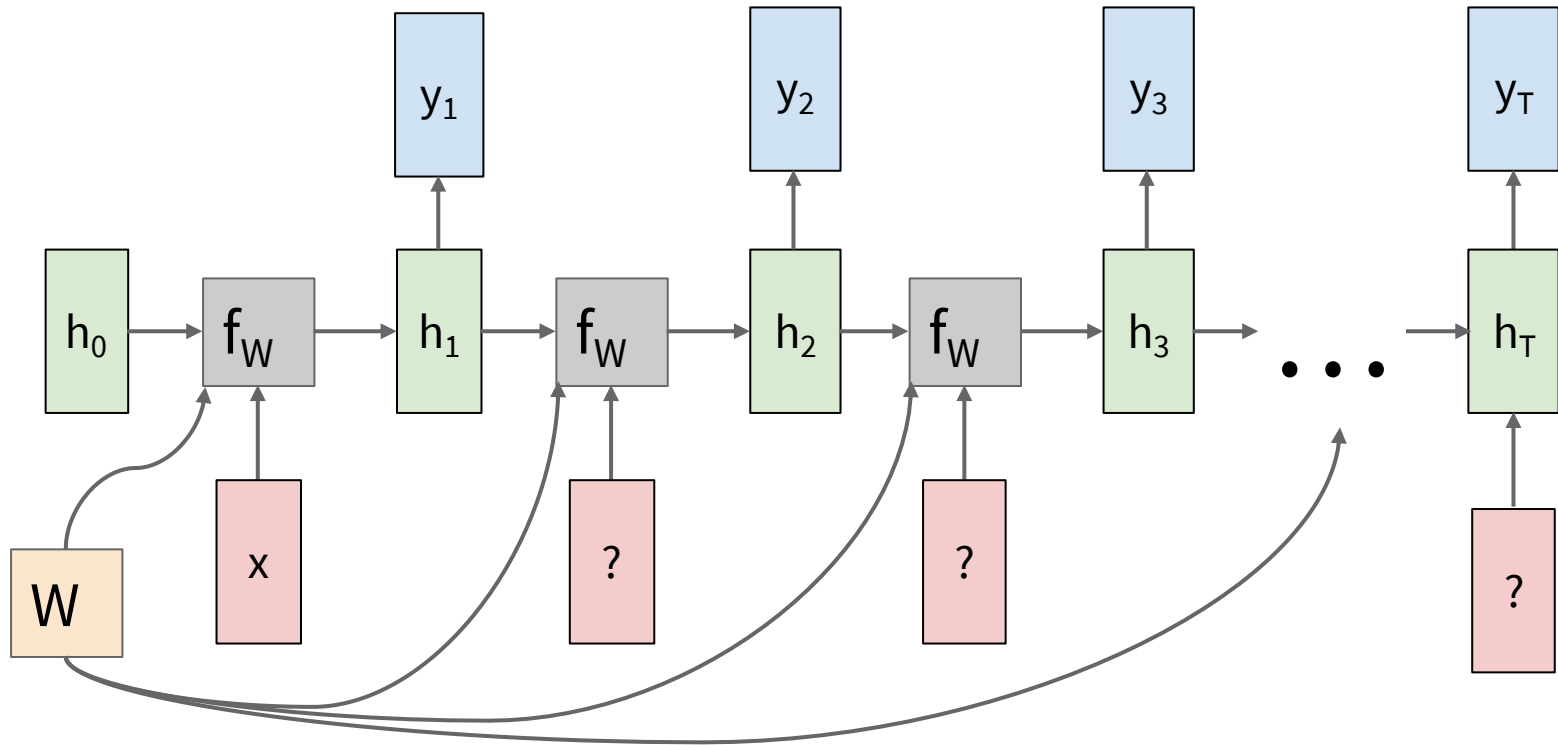
# RNN: Computational Graph: Many to One



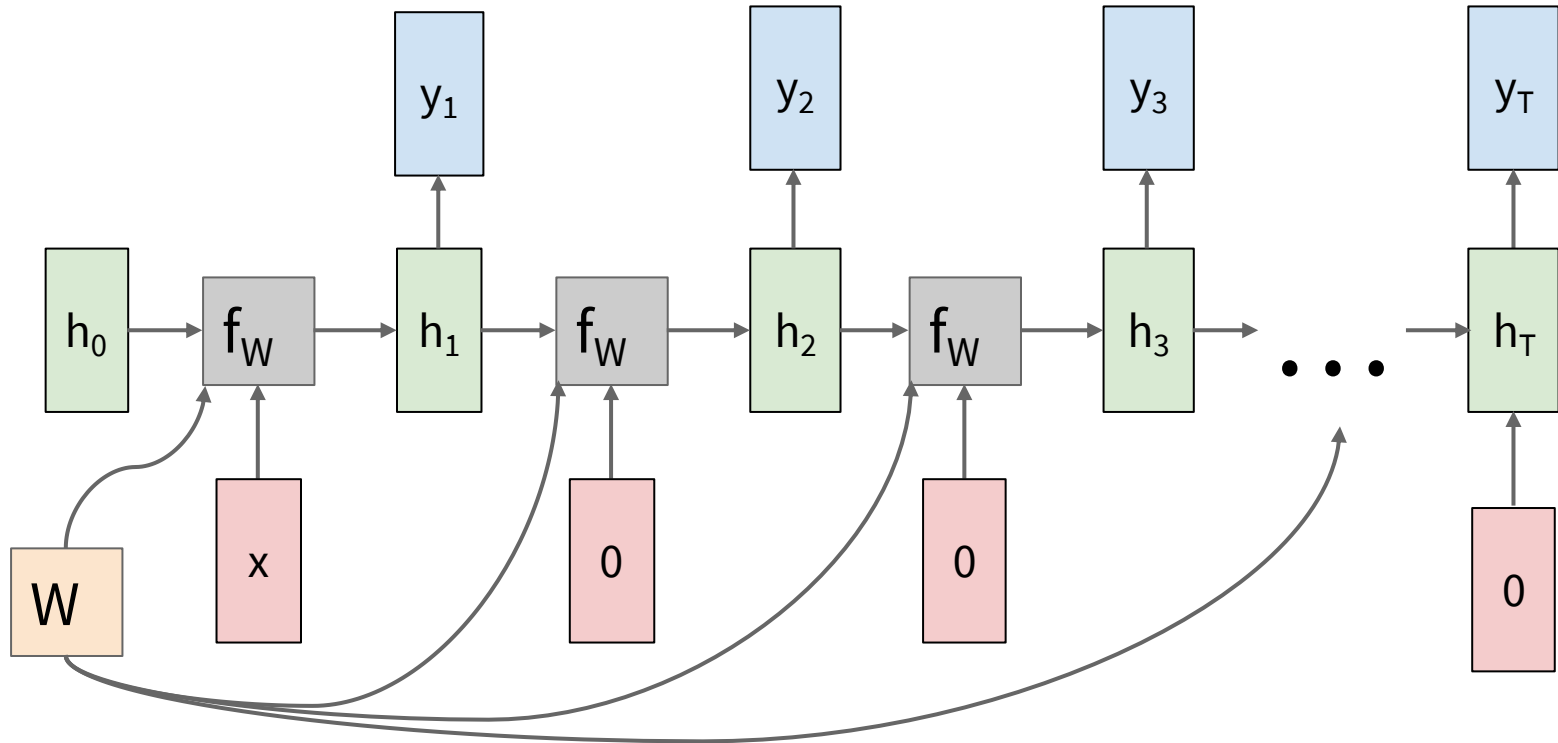
# RNN: Computational Graph: One to Many



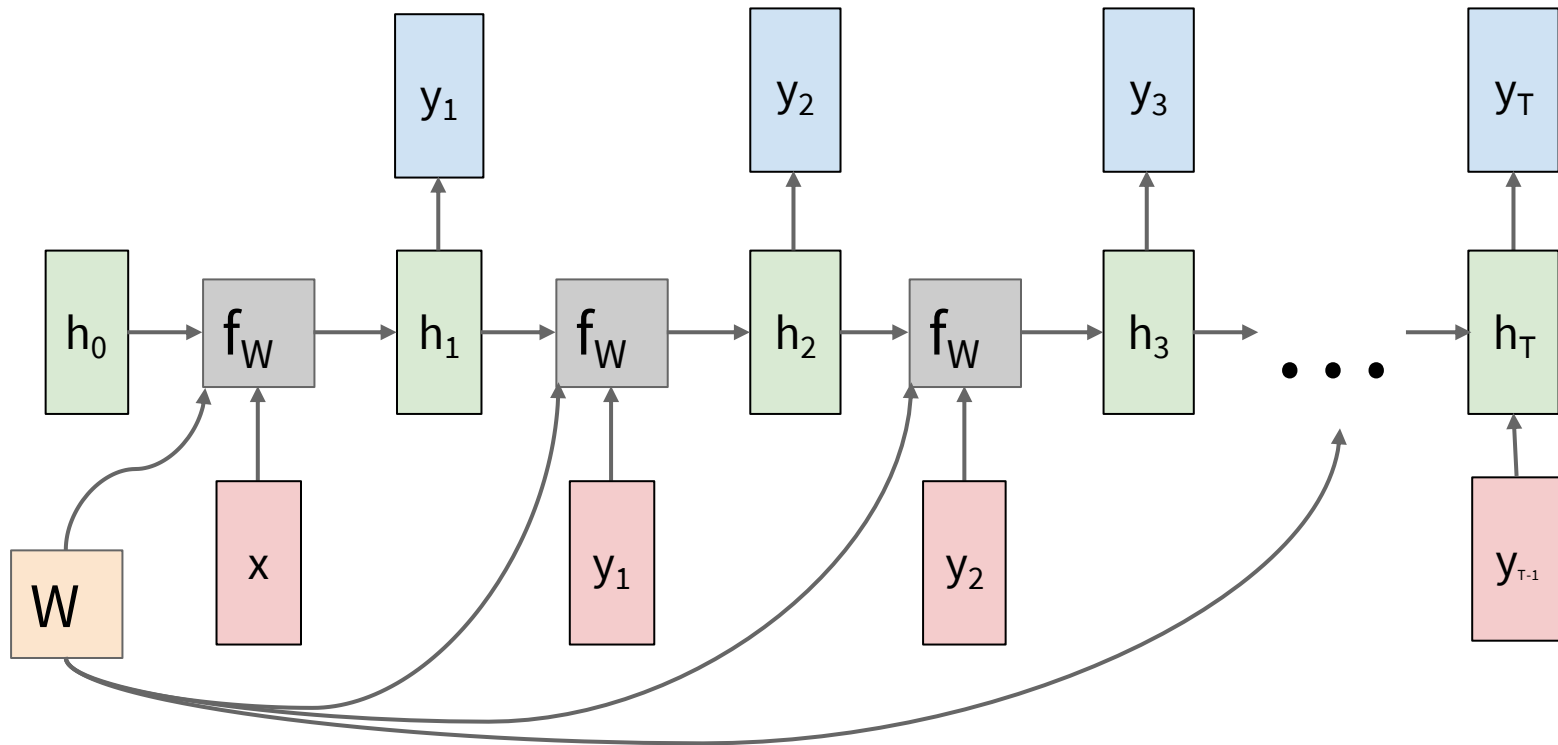
# RNN: Computational Graph: One to Many



# RNN: Computational Graph: One to Many

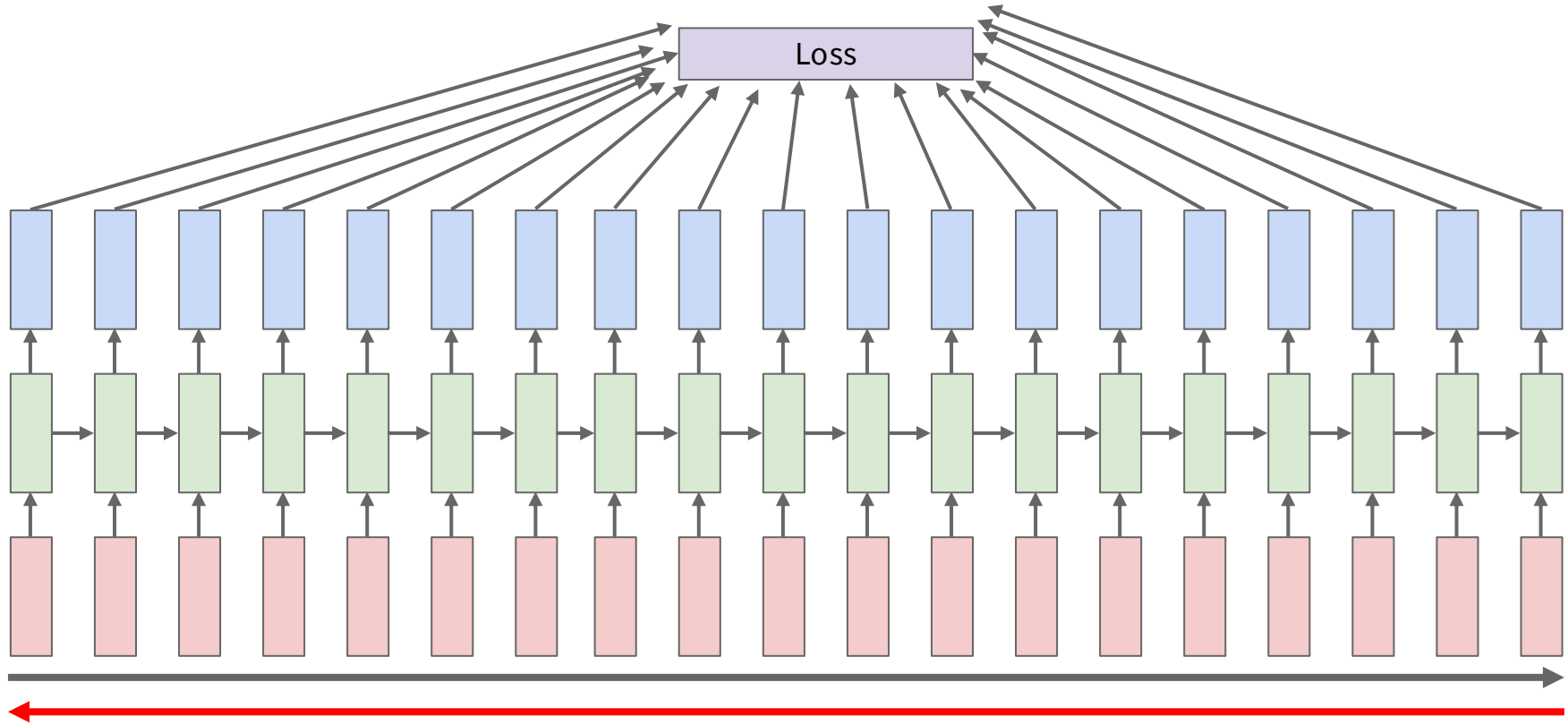


# RNN: Computational Graph: One to Many

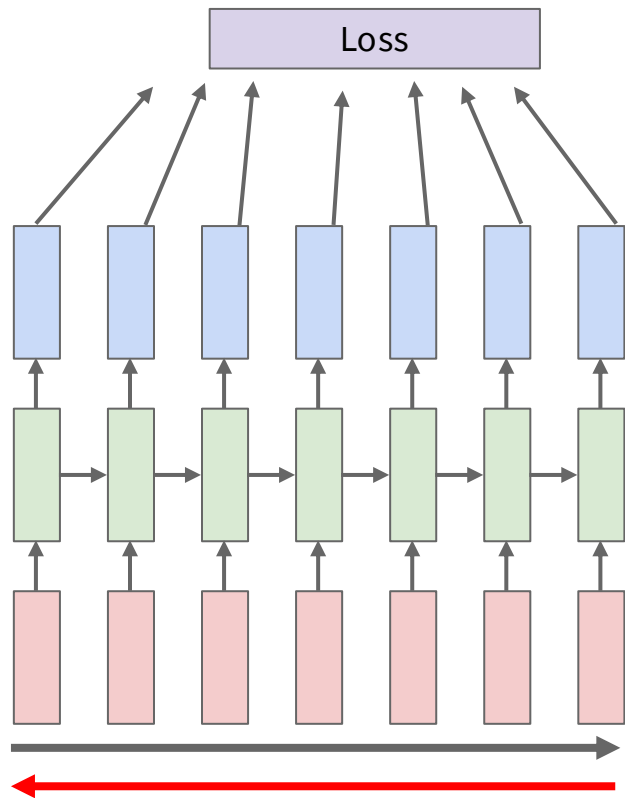


# Backpropagation through time

Forward through entire sequence to compute loss, then backward through entire sequence to compute gradient

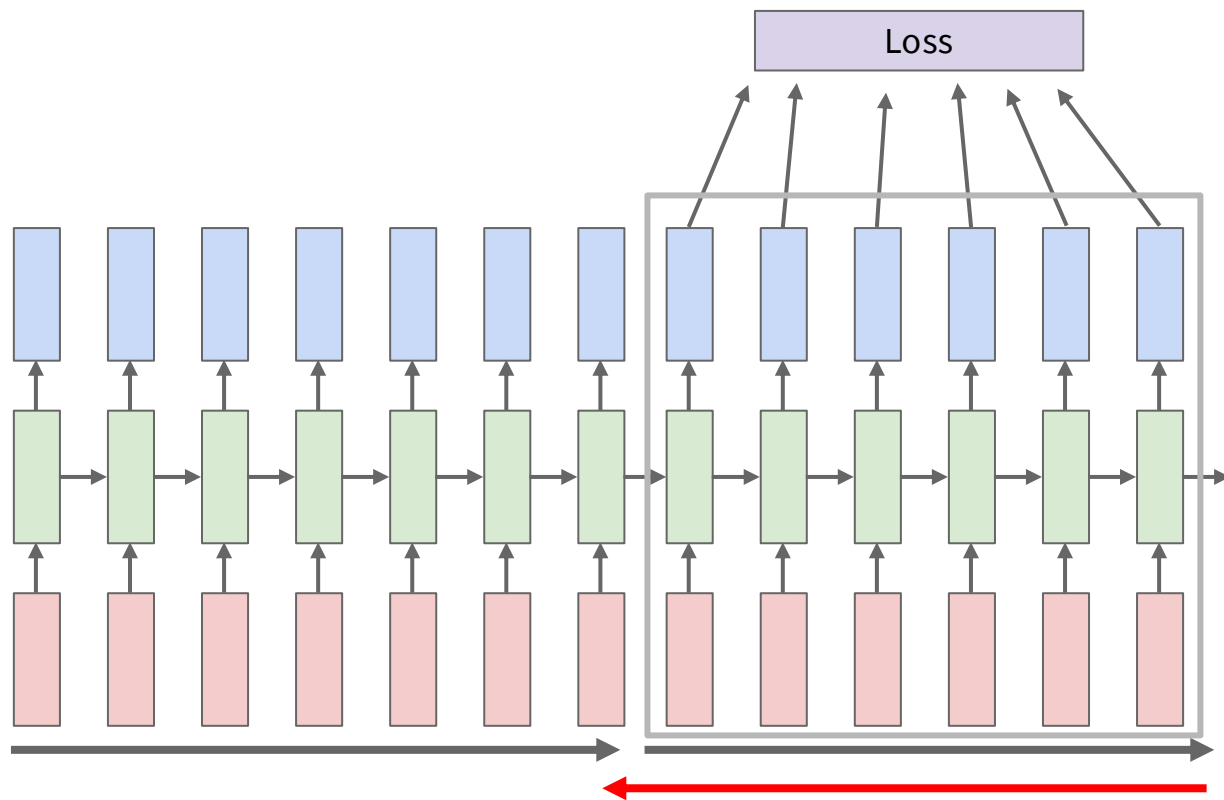


# Truncated Backpropagation through time



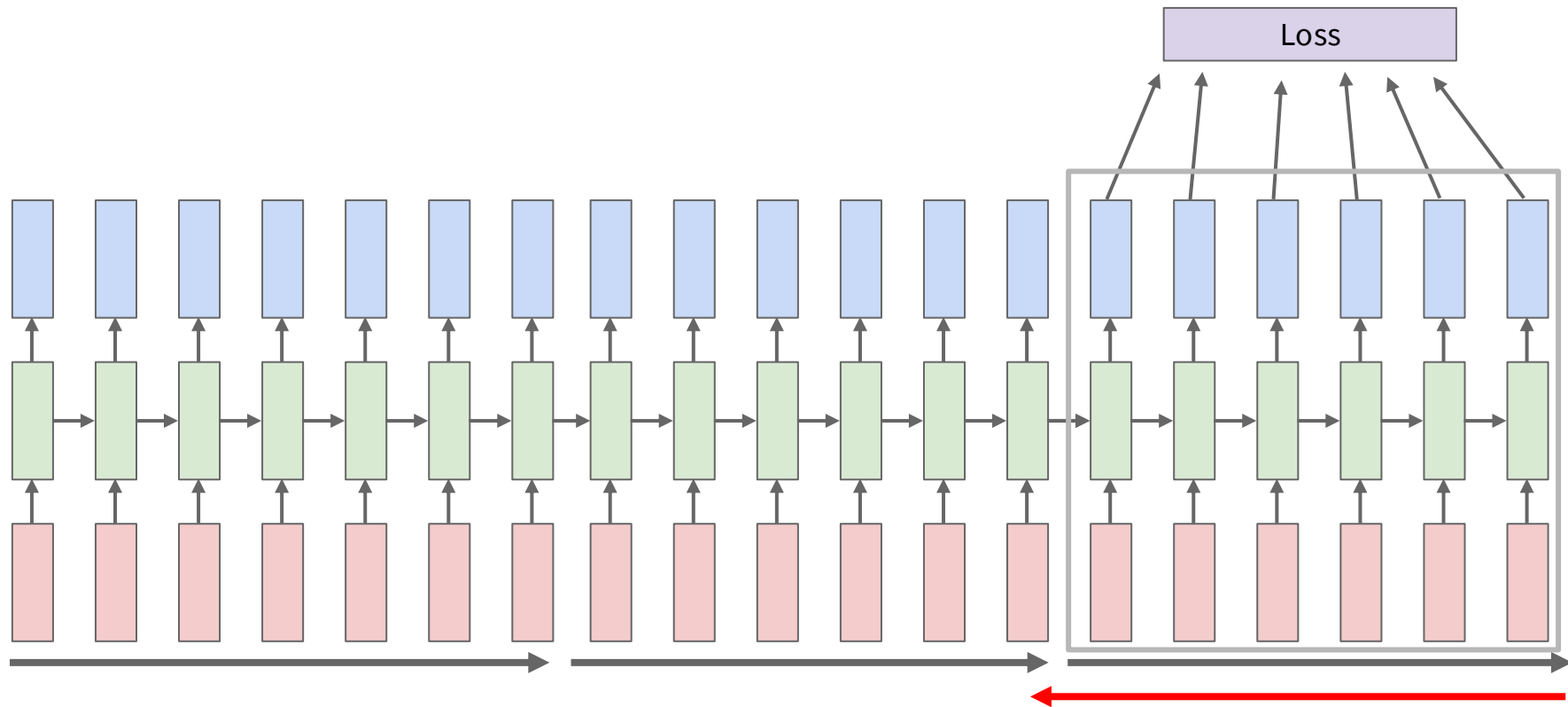
Run forward and backward through chunks of the sequence instead of whole sequence

# Truncated Backpropagation through time



Carry hidden states forward in time forever, but only backpropagate for some smaller number of steps

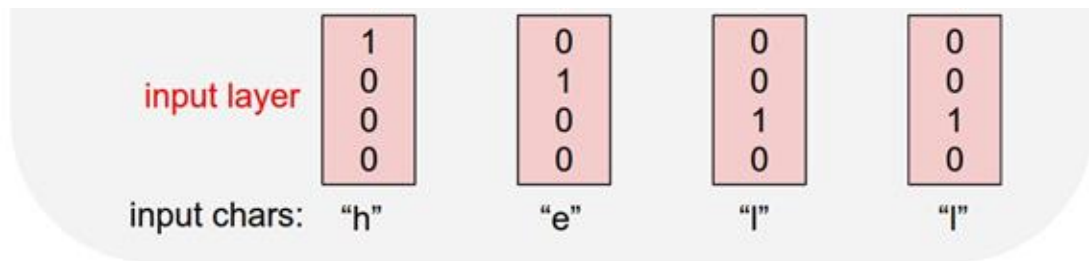
# Truncated Backpropagation through time



A more practical  
example:  
Character-level  
Language Model

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
“hello”

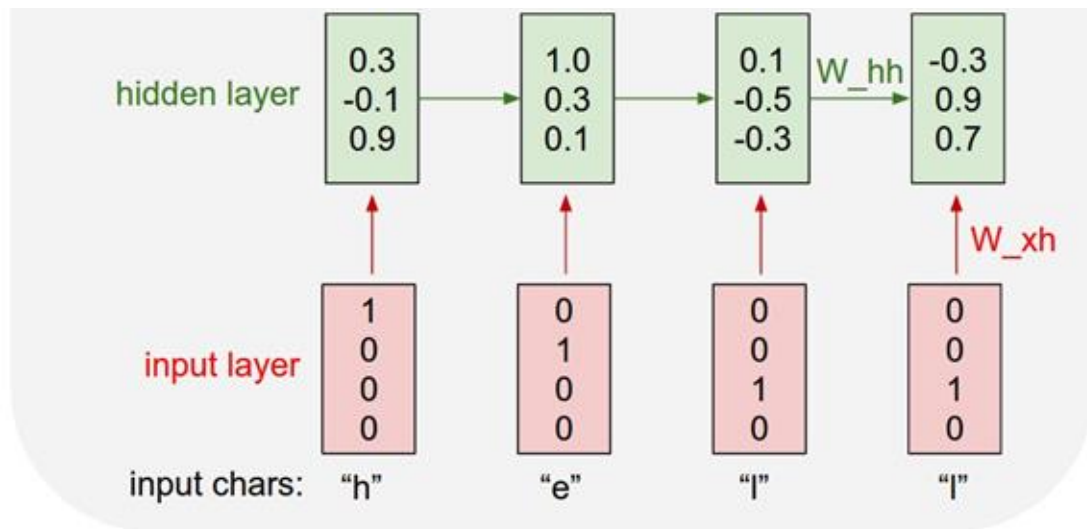


A more practical example:  
Character-level Language Model

Vocabulary:  
[h,e,l,o]

Example training sequence:  
“hello”

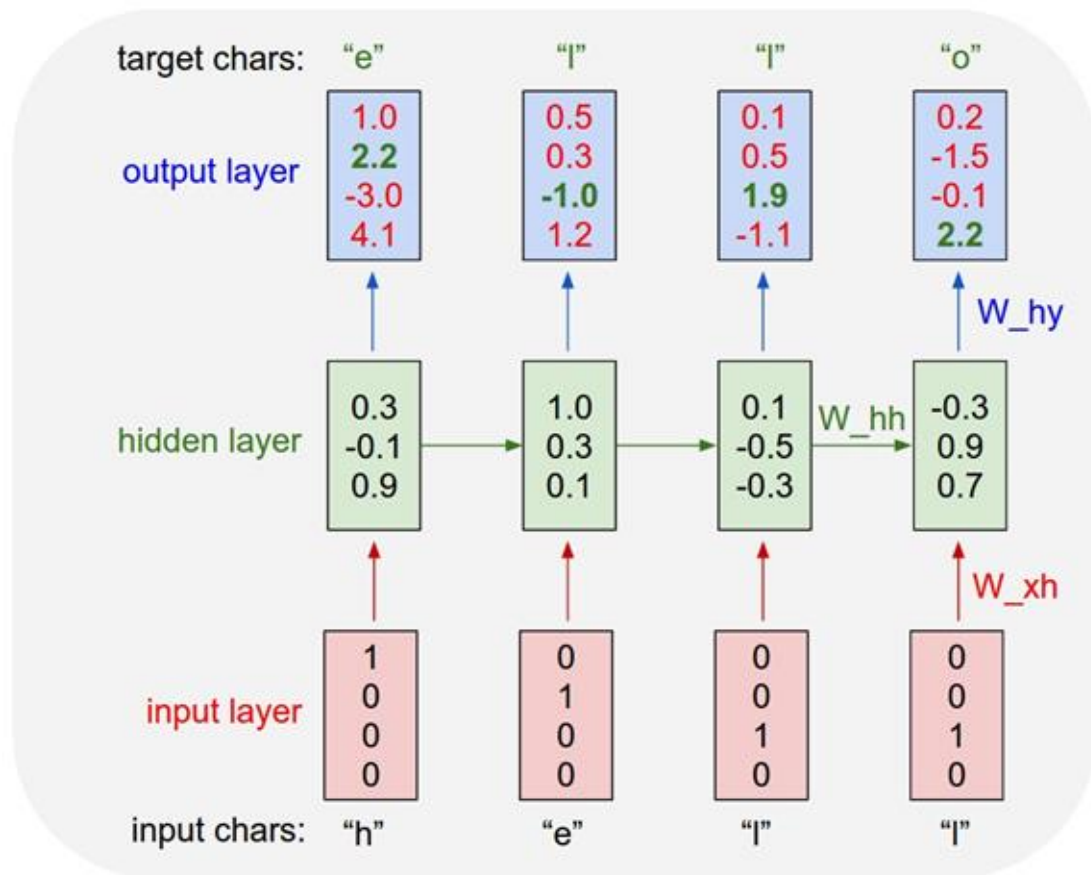
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$



A more practical example:  
Character-level Language Model

Vocabulary:  
[h,e,l,o]

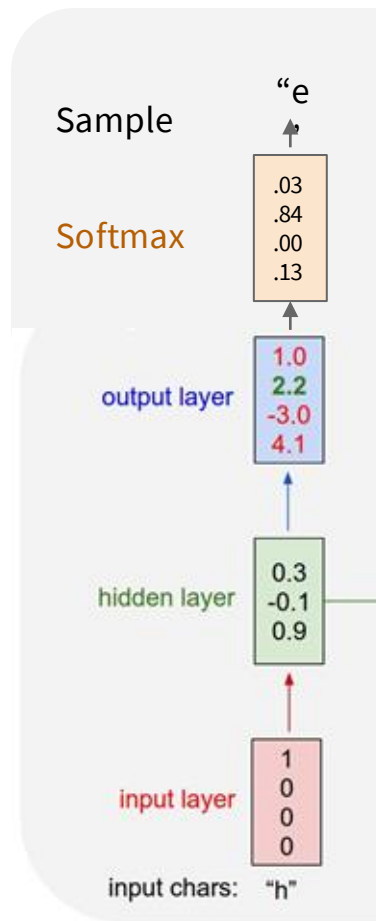
Example training sequence:  
"hello"



# Example: Character-level Language Model Sampling

Vocabulary:  
[h,e,l,o]

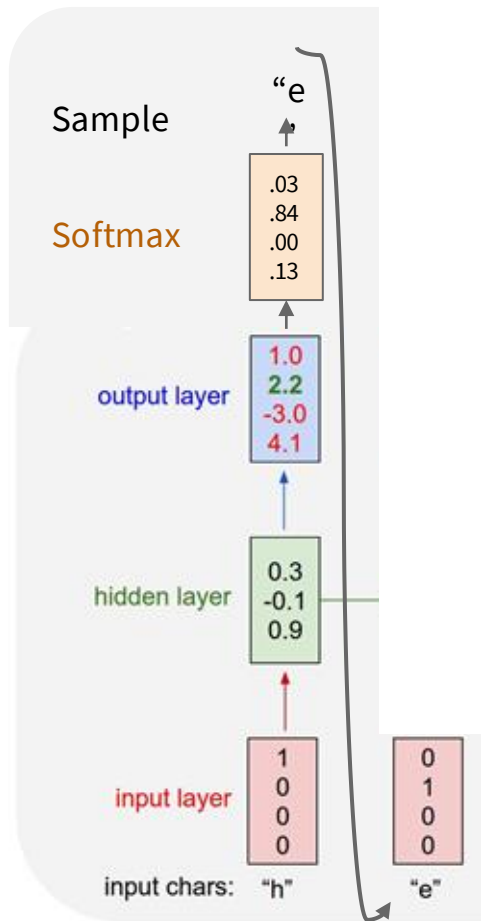
At test-time sample characters one at a time, feed back to model



# Example: Character-level Language Model Sampling

Vocabulary:  
[h,e,l,o]

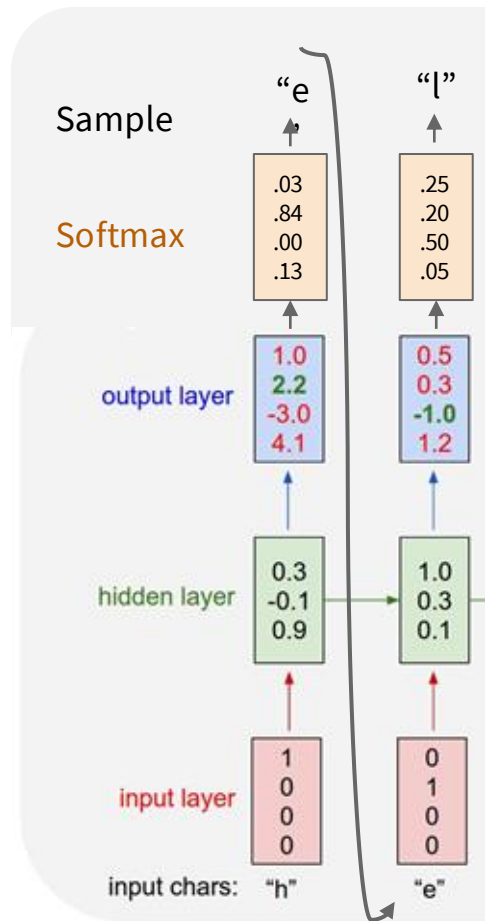
At test-time sample characters one at a time, feed back to model



# Example: Character-level Language Model Sampling

Vocabulary:  
[h,e,l,o]

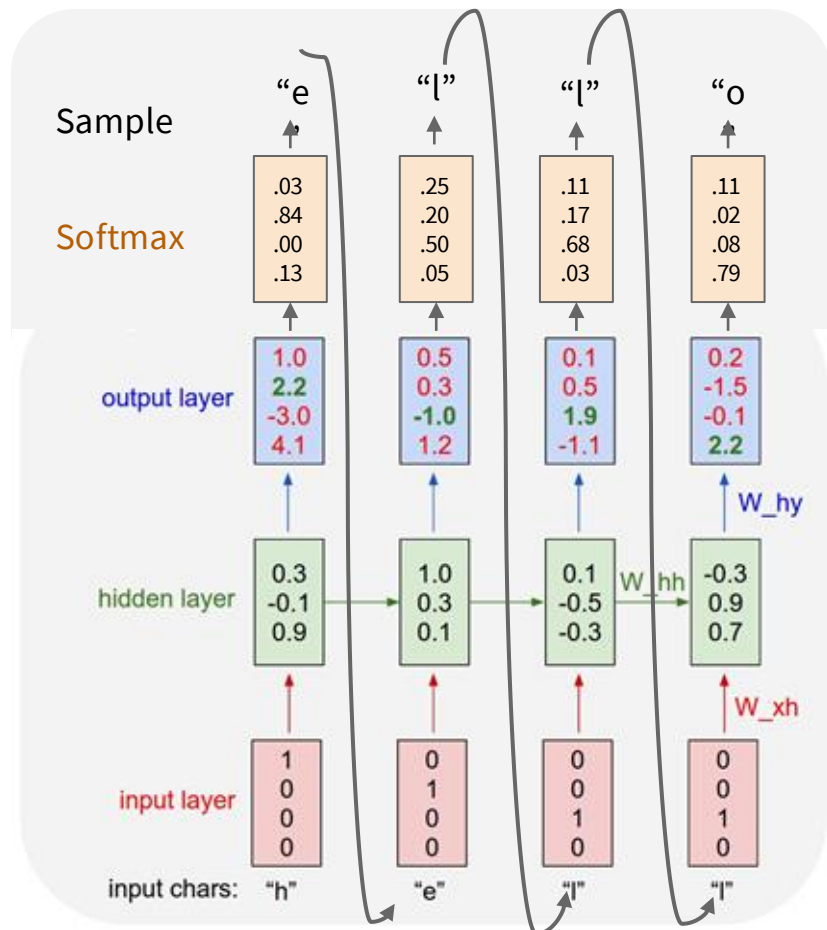
At test-time sample characters one at a time, feed back to model



# Example: Character-level Language Model Sampling

Vocabulary:  
[h,e,l,o]

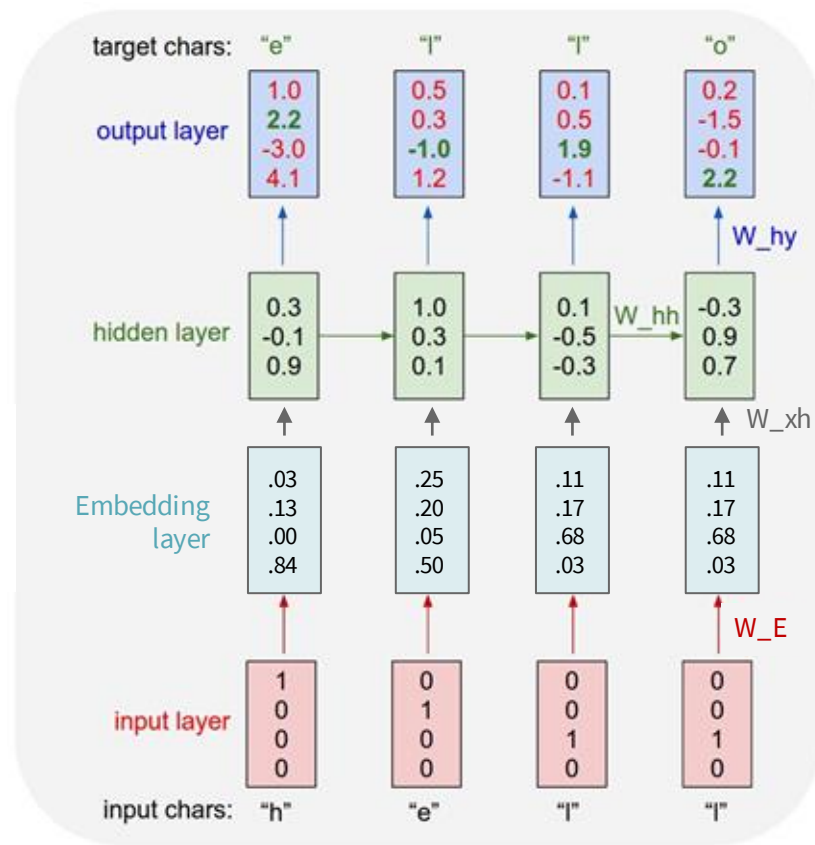
At test-time sample characters one at a time, feed back to model



# Example: Character-level Language Model Sampling

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \\ w_{41} & w_{42} & w_{43} & w_{44} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} w_{11} \\ w_{21} \\ w_{31} \\ w_{41} \end{bmatrix} = W_{\cdot E}$$

Matrix multiplication with a one-hot vector just extracts a column from the weight matrix. We often put a separate embedding layer between the input and hidden layers.



# min-char-rnn.py gist: 112 lines of Python

```
1  """
2  Minimal character-level vanilla RNN model. written by Andrej Karpathy (@karpathy)
3  BSD license
4  """
5  import numpy as np
6
7  # data I/O
8
9  data = open('input.txt', 'r').read() # should be simple plain text file
10 chars = list(set(data))
11 data_size, vocab_size = len(data), len(chars)
12 print "data has %d characters, %d unique." % (data_size, vocab_size)
13 char_to_ix = { ch:i for i,ch in enumerate(chars) }
14 ix_to_char = { i:ch for i,ch in enumerate(chars) }
15
16 # hyperparameters
17 hidden_size = 200 # size of hidden layer of neurons
18 seq_length = 25 # number of steps to unroll the RNN for
19 learning_rate = 0.1
20
21 # model parameters
22 wxx = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
23 whx = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
24 why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
25 bh = np.zeros((hidden_size, 1)) # hidden bias
26 by = np.zeros((vocab_size, 1)) # output bias
27
28 def lossFun(inputs, targets, hprev):
29     """
30     inputs, targets are both list of integers.
31     hprev is Nx1 array of initial hidden state
32     returns the loss, gradients on model parameters, and last hidden state
33     """
34     xs, ys, ys_ga = [], [], []
35     h[1] = np.copy(hprev)
36     loss = 0
37     # forward pass
38     for i in xrange(len(inputs)):
39         xs[i] = np.zeros((vocab_size, 1)) # encode in 1-of-K representation
40         xs[i][inputs[i]] = 1
41         hs[i] = np.tanh(np.dot(wxx, xs[i]) + np.dot(whx, hs[i-1]) + bh) # hidden state
42         ys[i] = np.dot(why, hs[i]) + by # unnormalized log probabilities for next char
43         ps[i] = np.exp(ys[i]) / np.sum(np.exp(ys[i])) # probabilities for next char
44         loss += -np.log(ps[i][targets[i],0]) # softmax (cross-entropy) loss
45     # backward pass: compute gradients going backwards
46     dxx, dwh, dby = np.zeros_like(wxx), np.zeros_like(whx), np.zeros_like(why)
47     dht, dhty = np.zeros_like(hs), np.zeros_like(hs)
48     dparam = np.zeros_like(hs[0])
49     for i in reversed(xrange(len(inputs))):
50         dy = np.copy(dht[i])
51         dhtargets[i] += 1 + dhtargets[i]*y
52         dhty += np.dot(dy, hs[i], T)
53         dht += dy
54         dh = np.dot(why, dy) + dhtnext # backward into h
55         dthrew = (1 - hs[i]**2) * dht[i] # dh = backprop through tanh nonlinearity
56         dht += dthrew
57         dxx += np.dot(dthrew, xs[i], T)
58         dwh += np.dot(dthrew, hs[i-1], T)
59         dhtnext = np.dot(whx, T, dht)
60     for dparam in [dxx, dwh, dby, dht, dhty]:
61         np.clip(dparam, -1, 1, out=dparam) # clip to mitigate exploding gradients
62     return loss, dxx, dwh, dby, dht, dhty, hs[len(inputs)-1]
63
64 def sample(h, seed_ix, n):
65     """
66     sample a sequence of integers from the model
67     h is memory state, seed_ix is seed letter for first time step
68     """
69     x = np.zeros((vocab_size, 1))
70     x[seed_ix] = 1
71     seqs = []
72     for t in xrange(n):
73         h = np.tanh(np.dot(wxx, x) + np.dot(whx, h) + bh)
74         y = np.dot(why, h) + by
75         p = np.exp(y) / np.sum(np.exp(y))
76         ix = np.random.choice(range(vocab_size), p=p.ravel())
77         x = np.zeros((vocab_size, 1))
78         x[ix] = 1
79         seqs.append(ix)
80     return seqs
81
82 n, p, s, g
83 dxx, dwh, dby = np.zeros_like(wxx), np.zeros_like(whx), np.zeros_like(why)
84 dht, dhty = np.zeros_like(hs), np.zeros_like(hs) # memory variables for adjoint
85 smooth_loss = -np.log(1.0/vocab_size)/seq_length # loss at iteration 0
86 while True:
87     # prepare inputs (we're sampling from left to right in steps seq_length long)
88     if p+seq_length > len(data) or n == 0:
89         hprev = np.zeros((hidden_size,1)) # reset RNN memory
90         g = g + 1 # up from start of data
91         inputs = [char_to_ix[ch] for ch in data[p:seq_length]]
92         targets = [char_to_ix[ch] for ch in data[p+1:seq_length+1]]
93     # sample from the model! now and then
94     if n % 100 == 0:
95         sample_ix = sample(hprev, inputs[0], 200)
96         txt = "".join(ix_to_char[ix] for ix in sample_ix)
97         print "----in %s %s----" % (txt, g)
98     # forward seq_length characters through the net and fetch gradients
99     loss, dxx, dwh, dby, dht, dhty, hprev = lossFun(inputs, targets, hprev)
100     smooth_loss = smooth_loss * 0.999 + loss * 0.001
101     if n % 100 == 0: print "iter %d, loss: %f" % (n, smooth_loss) # print progress
102
103     # update parameters with adjgrad
104     for param, dparam, w in zip([wxx, whx, why, bh, by],
105                               [dxx, dwh, dby, dht, dhty],
106                               [wxx, whx, why, bh, by]):
107         w += dparam * dparam
108     param += -learning_rate * dparam / np.sqrt((n + 1) * 1) # adjgrad update
109
110 g += seq_length # new data pointer
111 n += 1 # iteration counter
```

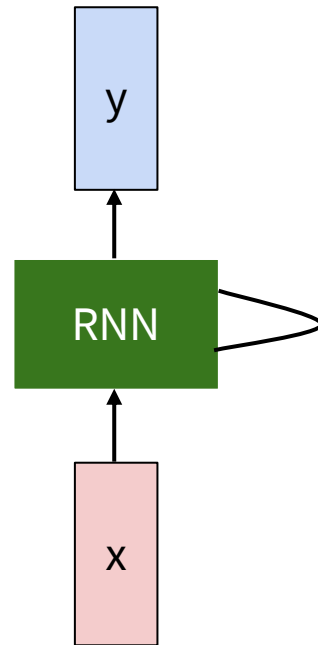
(<https://gist.github.com/karpathy/d4dee566867f8291f086>)

# THE SONNETS

by William Shakespeare

From fairest creatures we desire increase,  
That thereby beauty's rose might never die,  
But as the ripper should by time decrease,  
His tender heir might bear his memory:  
But thou, contracted to thine own bright eyes,  
Feed'st thy light's flame with self-substantial fuel,  
Making a famine where abundance lies,  
Thyself thy foe, to thy sweet self too cruel:  
Thou that art now the world's fresh ornament,  
And only herald to the gaudy spring,  
Within thine own bud buriest thy content,  
And tender churl mak'st waste in niggarding:  
Pity the world, or else this glutton be,  
To eat the world's due, by the grave and thee.

When forty winters shall besiege thy brow,  
And dig deep trenches in thy beauty's field,  
Thy youth's proud livery so gazed on now,  
Will be a tatter'd weed of small worth held:  
Then being asked, where all thy beauty lies,  
Where all the treasure of thy lusty days;  
To say, within thine own deep sunken eyes,  
Were an all-eating shame, and thriftless praise.  
How much more praise deserv'd thy beauty's use,  
If thou couldst answer 'This fair child of mine  
Shall sum my count, and make my old excuse,'  
Proving his beauty by succession thine!  
This were to be new made when thou art old,  
And see thy blood warm when thou feel'st it cold.



[Blogpost from Andrej Karpathy back in 2015!](#)

at first:

tyntd-iafhatawiaoihrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e  
plia tklrqd t o idoe ns,smtt h ne etie h,hregtrs nigtkie,aoaenns lng



train more

"Tmont thithey" fomesscerliund  
Keushey. Thom here  
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwv fil on aseterlome  
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."



train more

Aftair fall unsuch that the hall for Prince Velzonski's that me of  
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort  
how, and Gogition is so overelical and offer.



train more

"Why do what that day," replied Natasha, and wishing to himself the fact the  
princess, Princess Mary was easier, fed in had oftened him.  
Pierre aking his soul came to the packs and drove up his father-in-law women.

PANDARUS:

Alas, I think he shall be come approached and the day  
When little strain would be attain'd into being never fed,  
And who is but a chain and subjects of his death,  
I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,  
Breaking and strongly should be buried, when I perish  
The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

Second Lord:

They would be ruled after this chamber, and  
my fair nudes begun out of the fact, to be conveyed,  
Whose noble souls I'll have the heart of the wars.

Clown:

Come, sir, I will make did behold your worship.

VIOLA:

I'll drink it.

VIOLA:

Why, Salisbury must find his flesh and thought  
That which I am not apt, not a man and in fire,  
To show the reining of the raven and the wars  
To grace my hand reproach within, and not a fair are hand,  
That Caesar and my goodly father's world;  
When I was heaven of presence and our fleets,  
We spare with hours, but cut thy council I am great,  
Murdered and by thy master's ready there  
My power to give thee but so much as hell:  
Some service in the noble bondman here,  
Would show him to her wine.

KING LEAR:

O, if you were a feeble sight, the courtesy of your law,  
Your sight and several breath, will wear the gods  
With his heads, and my hands are wonder'd at the deeds,  
So drop upon your lordship's head, and your opinion  
Shall be against your honour.



This repository Search

Explore Gist Blog Help



karpathy



torvalds / linux

Watch -

3,711

★ Star

23,054

Fork

9,141

### Linux kernel source tree

520,037 commits

1 branch

420 releases

5,039 contributors



branch: master -

linux / +



Merge branch 'drm-fixes' of git://people.freedesktop.org/~airlied/linux



torvalds authored 9 hours ago

latest commit 4b1706927d

|               |   |              |
|---------------|---|--------------|
| Documentation | Merge git://git.kernel.org/pub/scm/linux/kernel/git/hab/target-pending      | 6 days ago   |
| arch          | Merge branch 'x86-urgent-for-linus' of git://git.kernel.org/pub/scm/lin...  | a day ago    |
| block         | block: discard bdi_unregister() in favour of bdi_destroy()                  | 9 days ago   |
| crypto        | Merge git://git.kernel.org/pub/scm/linux/kernel/git/herbert/crypto-2.6      | 10 days ago  |
| drivers       | Merge branch 'drm-fixes' of git://people.freedesktop.org/~airlied/linux     | 9 hours ago  |
| firmware      | firmware/hex2fw.c: restore missing default in switch statement              | 2 months ago |
| fs            | vfs: read file_handle only once in handle_to_path                           | 4 days ago   |
| include       | Merge branch 'perf-urgent-for-linus' of git://git.kernel.org/pub/scm/lin... | a day ago    |
| init          | init: fix regression by supporting devices with major:minor:offset fo...    | a month ago  |
| ipc           | Merge branch 'for-linus' of git://git.kernel.org/pub/scm/linux/kernel...    | a month ago  |



Code



Pull requests

74



Pulse



Graphs

HTTPS clone URL

https://github.c



You can clone with [HTTPS](#), [SSH](#), or [Subversion](#).

Clone in Desktop

Download ZIP

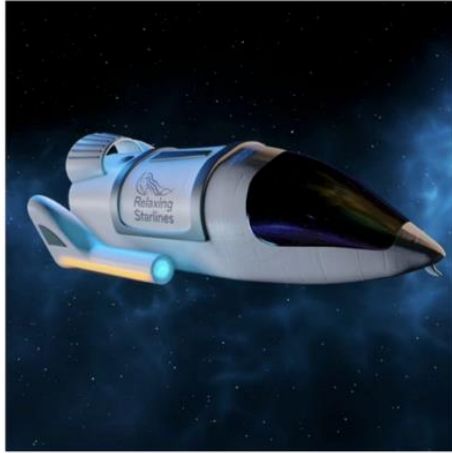
```

static void do_command(struct seq_file *m, void *v)
{
    int column = 32 << (cmd[2] & 0x80);
    if (state)
        cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
    else
        seq = 1;
    for (i = 0; i < 16; i++) {
        if (k & (1 << i))
            pipe = (in_use & UMXTHREAD_UNCCA) +
                ((count & 0x00000000fffffff8) & 0x000000f) << 8;
        if (count == 0)
            sub(pid, ppc_md.kexec_handle, 0x20000000);
        pipe_set_bytes(i, 0);
    }
    /* Free our user pages pointer to place camera if all dash */
    subsystem_info = &of_changes[PAGE_SIZE];
    rek_controls(offset, idx, &soffset);
    /* Now we want to deliberately put it to device */
    control_check_polarity(&context, val, 0);
    for (i = 0; i < COUNTER; i++)
        seq_puts(s, "policy ");
}

```

# Generated C code

# OpenAI Codex, GitHub Copilot, Claude Code, Cursor IDE



Add this image of a rocketship:  
<https://i1.sndcdn.com/artworks-j8xjG7zc1wmTe07b-06l83w-t500x500.jpg>

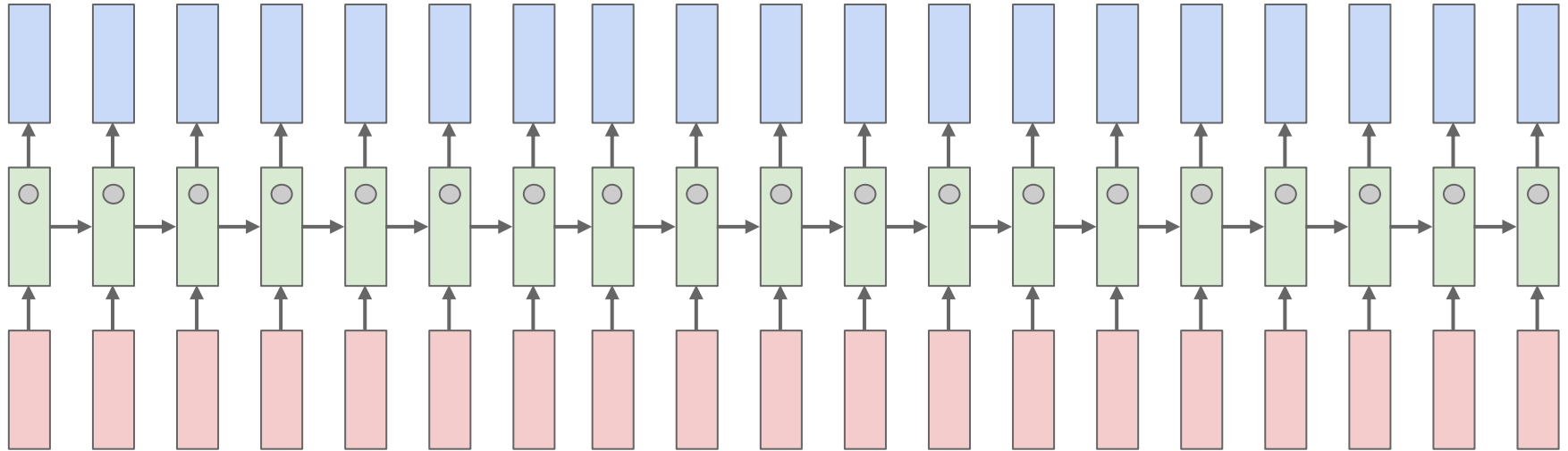


```
/* Add this image of a
rocketship:
https://i1.sndcdn.com/artworks
-j8xjG7zc1wmTe07b-06l83w-
t500x500.jpg */
var rocketship =
document.createElement('img');
rocketship.src =
'https://i1.sndcdn.com/artwork
s-j8xjG7zc1wmTe07b-06l83w-
t500x500.jpg';
document.body.appendChild(rock
etship);
```



<https://openai.com/blog/openai-codex/>

# Searching for interpretable cells



# Searching for interpretable cells

```
/* Unpack a filter field's string representation from user-space
 * buffer. */
char *audit_unpack_string(void **bufp, size_t *remain, size_t len)
{
    char *str;
    if (!*bufp || (len == 0) || (len > *remain))
        return ERR_PTR(-EINVAL);
    /* Of the currently implemented string fields, PATH_MAX
     * defines the longest valid length.
     */
}
```

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016  
Figures copyright Karpathy, Johnson, and Fei-Fei, 2015; reproduced with permission

# Searching for interpretable cells

"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."

quote detection cell

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016  
Figures copyright Karpathy, Johnson, and Fei-Fei, 2015; reproduced with permission

# Searching for interpretable cells

Cell sensitive to position in line:

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.

line length tracking cell

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016  
Figures copyright Karpathy, Johnson, and Fei-Fei, 2015; reproduced with permission

# Searching for interpretable cells

```
static int __dequeue_signal(struct sigpending *pending, sigset_t *mask,
                           siginfo_t *info)
{
    int sig = next_signal(pending, mask);
    if (sig) {
        if (current->notifier) {
            if (sigismember(current->notifier_mask, sig)) {
                if (!(current->notifier)(current->notifier_data)) {
                    clear_thread_flag(TIF_SIGPENDING);
                    return 0;
                }
            }
        }
        collect_signal(sig, pending, info);
    }
    return sig;
}
```

if statement cell

# Searching for interpretable cells

Cell that turns on inside comments and quotes:

```
/* Duplicate LSM field information. The lsm_rule is opaque, so
 * re-initialized. */
static inline int audit_dupe_lsm_field(struct audit_field *df,
                                     struct audit_field *sf)
{
    int ret = 0;
    char *lsm_str;
    /* our own copy of lsm_str */
    lsm_str = kstrdup(sf->lsm_str, GFP_KERNEL);
    if (unlikely(!lsm_str))
        return -ENOMEM;
    df->lsm_str = lsm_str;
    /* our own (refreshed) copy of lsm_rule */
    ret = security_audit_rule_init(df->type, df->op, df->lsm_str,
                                  (void *)&df->lsm_rule);
    /* Keep currently invalid fields around in case they
     * become valid after a policy reload. */
    if (ret == -EINVAL) {
        pr_warn("audit rule for LSM '%s' is invalid\n",
                df->lsm_str);
        ret = 0;
    }
    return ret;
}
```

quote/comment cell

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

Figures copyright Karpathy, Johnson, and Fei-Fei, 2015; reproduced with permission

# Searching for interpretable cells

```
#ifdef CONFIG_AUDITSYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```

code depth cell

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016  
Figures copyright Karpathy, Johnson, and Fei-Fei, 2015; reproduced with permission

# RNN tradeoffs

## RNN Advantages:

- Can process any length of the input (no context length)
- Computation for step  $t$  can (in theory) use information from many steps back
- Model size does not increase for longer input
- The same weights are applied on every timestep, so there is symmetry in how inputs are processed.

## RNN Disadvantages:

- Recurrent computation is slow
- In practice, difficult to access information from many steps back

# Image Captioning

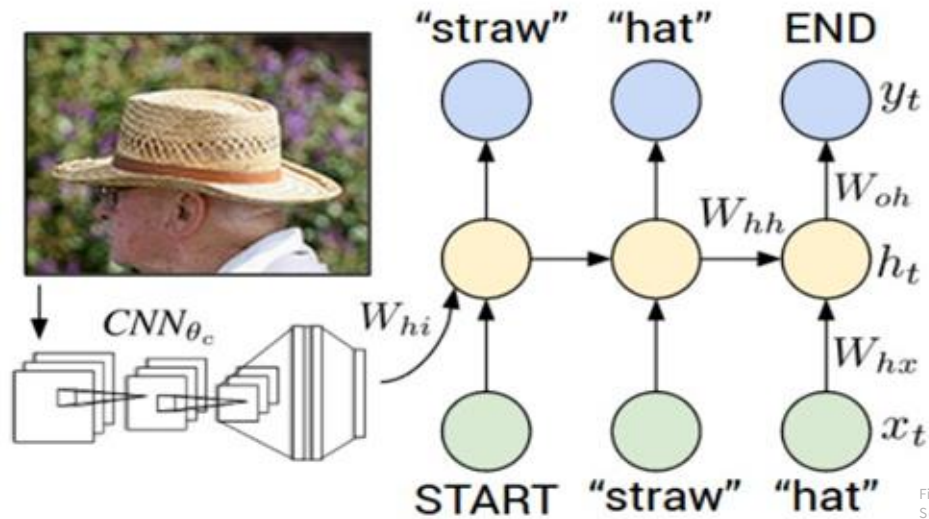


Figure from Karpathy et al, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015; figure copyright IEEE, 2015. Reproduced for educational purposes.

Explain Images with Multimodal Recurrent Neural Networks, Mao et al.

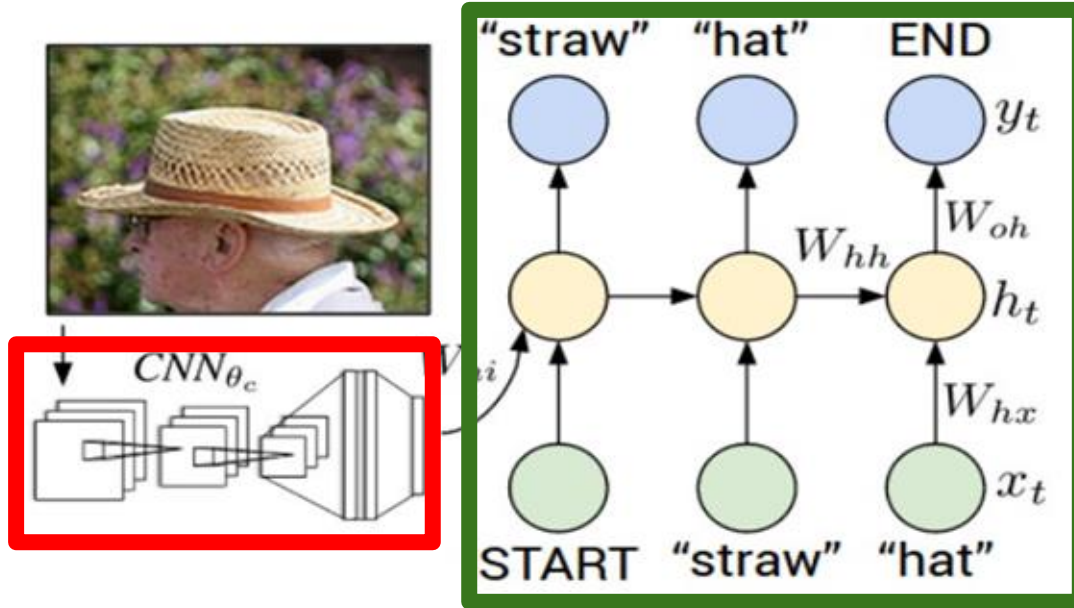
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei

Show and Tell: A Neural Image Caption Generator, Vinyals et al.

Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al.

Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

# Recurrent Neural Network



## Convolutional Neural Network



test image

[This image](#) is [CC0 public domain](#)

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

FC-1000

softmax



test image

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

~~FC-1000~~

~~softmax~~



test image

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

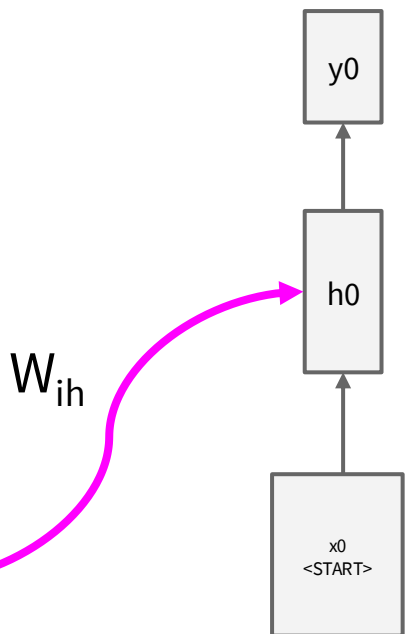


test image

x0  
<START>



test image



before:

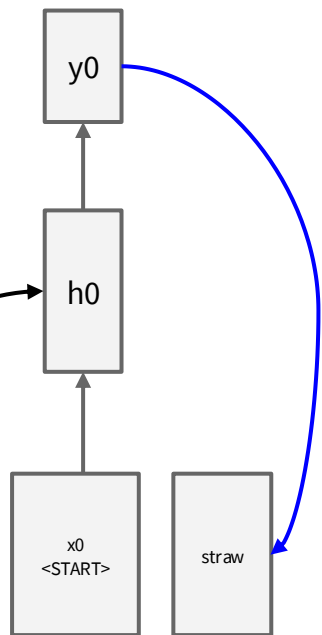
$$h = \tanh(W_{xh} * x + W_{hh} * h)$$

now:

$$h = \tanh(W_{xh} * x + W_{hh} * h + W_{ih} * v)$$



test image



sample!

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

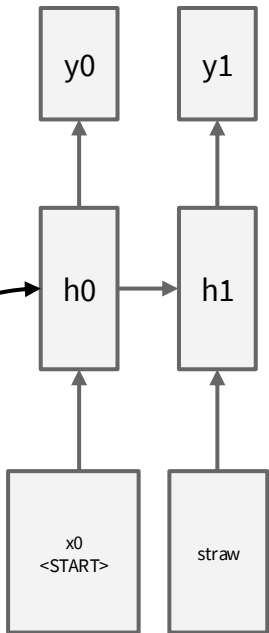
maxpool

FC-4096

FC-4096

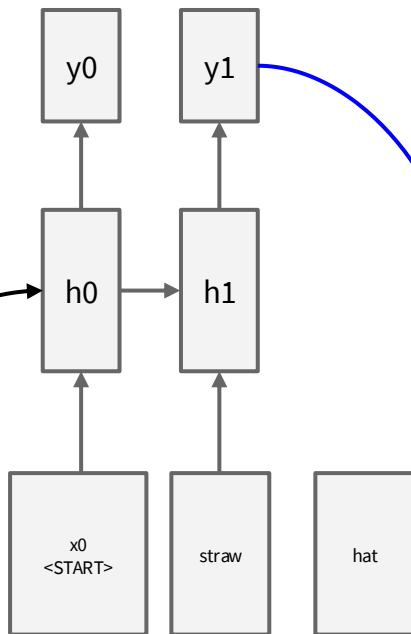


test image





test image



sample!

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

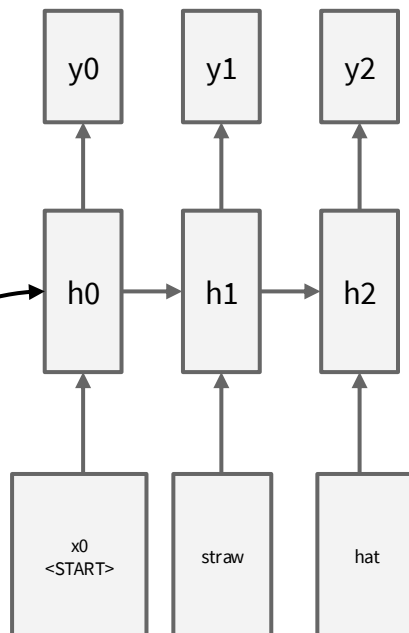
maxpool

FC-4096

FC-4096

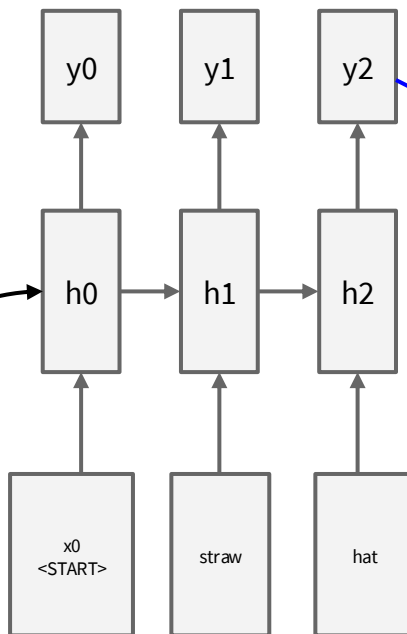


test image





test image



sample  
<END> token  
=> finish.

# Image Captioning: Example Results

Captions generated using [neuraltalk2](#)  
All images are [CC0 Public domain](#): [cat](#)  
[suitcase](#), [cat](#), [tree](#), [dog](#), [bear](#), [surfer](#),  
[tennis](#), [giraffe](#), [motorcycle](#)



A cat sitting on a suitcase on the floor



A cat is sitting on a tree branch



A dog is running in the grass with a frisbee



A white teddy bear sitting in the grass



Two people walking on the beach with surfboards



A tennis player in action on the court



Two giraffes standing in a grassy field



A man riding a dirt bike on a dirt track

# Image Captioning: Failure Cases

Captions generated using [neuraltalk2](#)  
All images are [CC0 Public domain](#): [fur coat](#),  
[handstand](#), [spider web](#), [baseball](#)



A woman is holding a cat  
in her hand



A woman standing on a beach  
holding a surfboard



A bird is perched on a  
tree branch



A person holding a computer  
mouse on a desk



A man in a  
baseball uniform  
throwing a ball

# Visual Question Answering (VQA)



**Q: What endangered animal is featured on the truck?**

- A: A bald eagle.**
- A: A sparrow.
- A: A humming bird.
- A: A raven.



**Q: Where will the driver go if turning right?**

- A: Onto 24 3/4 Rd.**
- A: Onto 25 3/4 Rd.
- A: Onto 23 3/4 Rd.
- A: Onto Main Street.



**Q: When was the picture taken?**

- A: During a wedding.**
- A: During a bar mitzvah.
- A: During a funeral.
- A: During a Sunday church service.

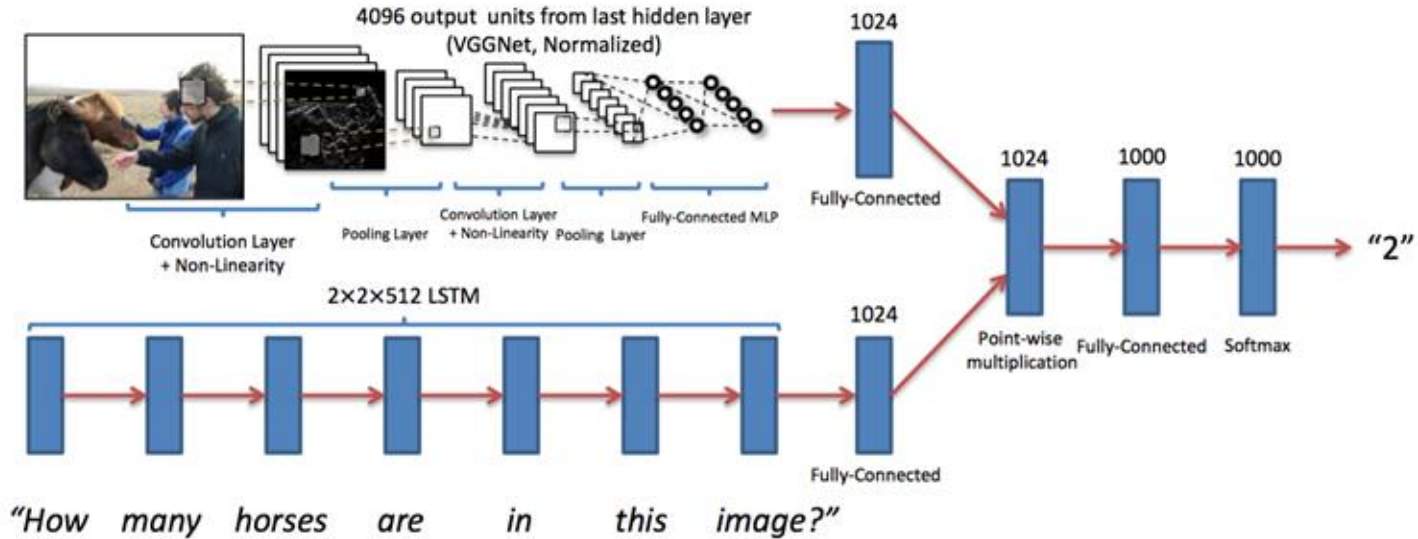


**Q: Who is under the umbrella?**

- A: Two women.**
- A: A child.
- A: An old man.
- A: A husband and a wife.

Agrawal et al, "VQA: Visual Question Answering", ICCV 2015  
Zhu et al, "Visual 7W: Grounded Question Answering in Images", CVPR 2016  
Figure from Zhu et al, copyright IEEE 2016. Reproduced for educational purposes.

# Visual Question Answering (VQA)



Agrawal et al, "Visual 7W: Grounded Question Answering in Images", CVPR 2015  
Figures from Agrawal et al, copyright IEEE 2015. Reproduced for educational purposes.

# Visual Dialog: Conversations about images

The screenshot displays a 'Visual Dialog' interface. At the top, a blue header contains the text 'Visual Dialog'. Below this, on the left, is a photograph of a tabby cat drinking from a blue and white striped mug on a table. To the right of the image, a dark grey box contains the text 'A cat drinking water out of a coffee mug.' Below the image, there are four grey speech bubbles with a small robot icon, representing the system's responses to user questions. The questions are shown in blue speech bubbles on the right. At the bottom, there is a blue bar with a refresh icon, a camera icon, the text 'Start typing question here ...', and a right-pointing arrow.

Visual Dialog

A cat drinking water out of a coffee mug.

What color is the mug?

White and red

Are there any pictures on it?

No, something is there can't tell what it is

Is the mug and cat on a table?

Yes, they are

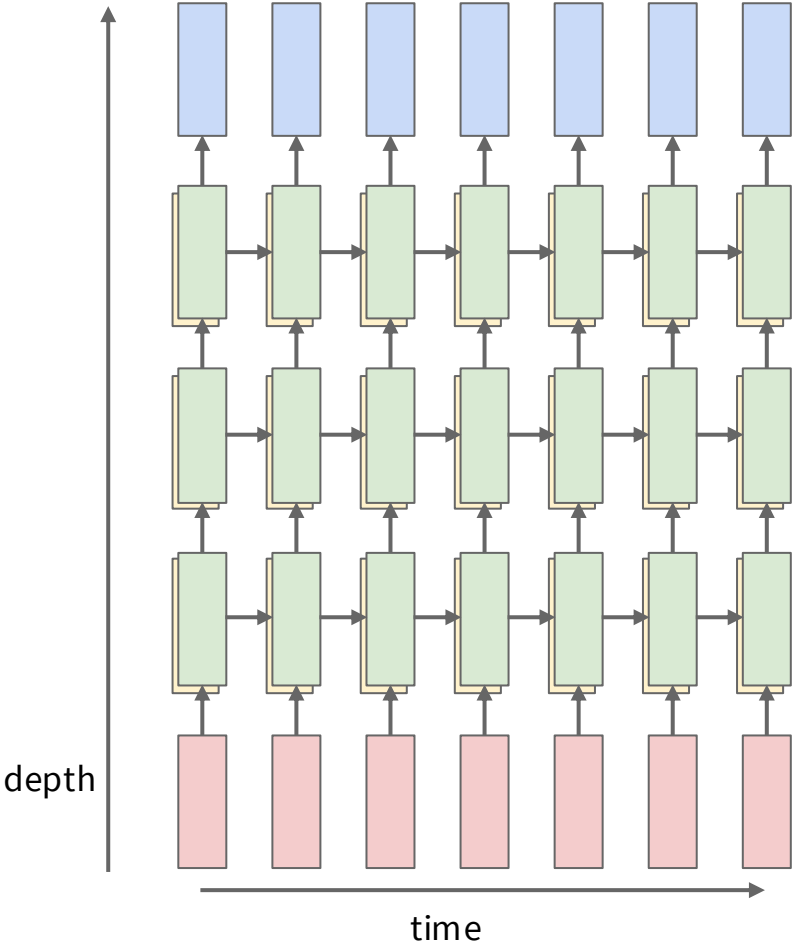
Are there other items on the table?

Yes, magazines, books, toaster and basket, and a plate

Start typing question here ...

Das et al, "Visual Dialog", CVPR 2017  
Figures from Das et al, copyright IEEE 2017. Reproduced with permission.

# Multilayer RNNs



# Long Short Term Memory (LSTM)

Vanilla RNN

$$h_t = \tanh \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

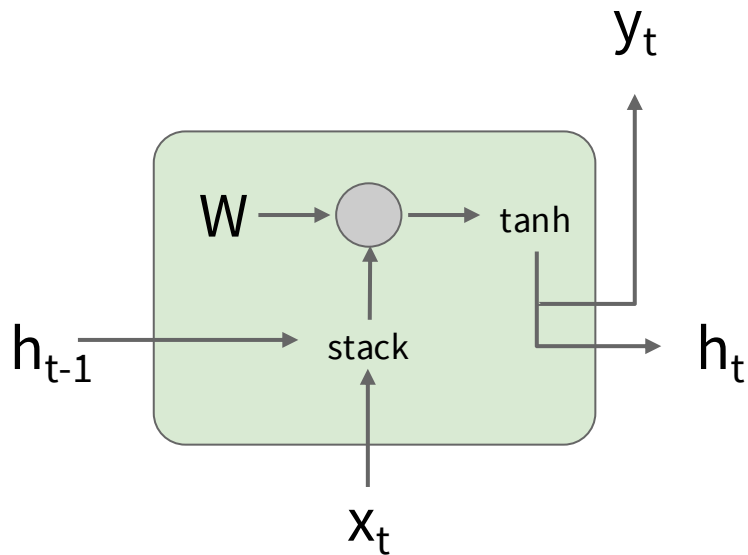
LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

# Vanilla RNN Gradient Flow

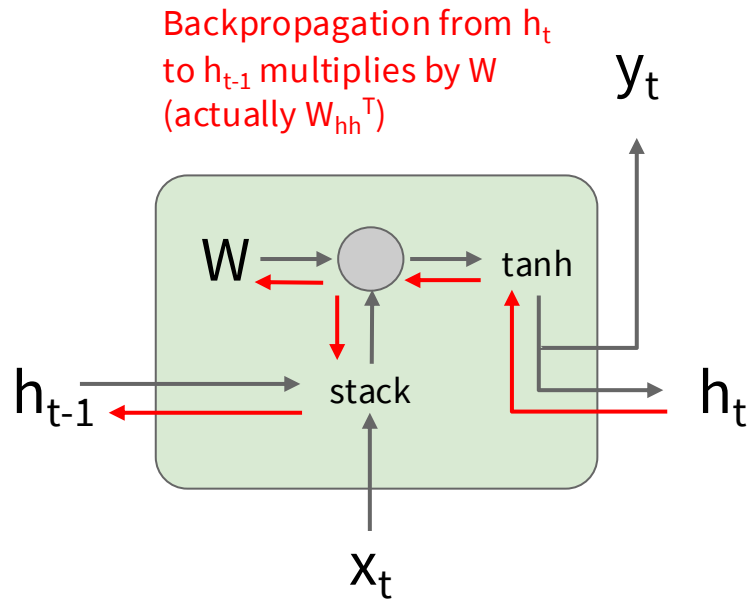
Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994  
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{hx}x_t) \\ &= \tanh\left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \\ &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \end{aligned}$$

# Vanilla RNN Gradient Flow

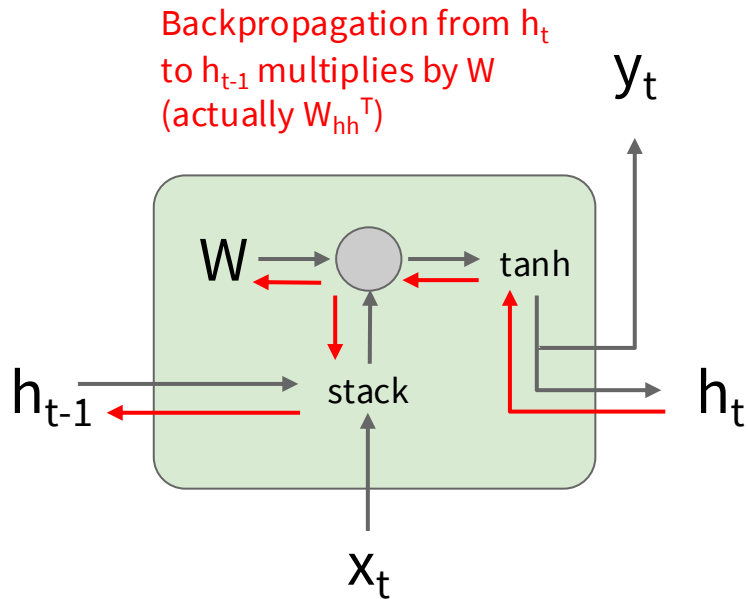
Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994  
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{hx}x_t) \\ &= \tanh\left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \\ &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \end{aligned}$$

# Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994  
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

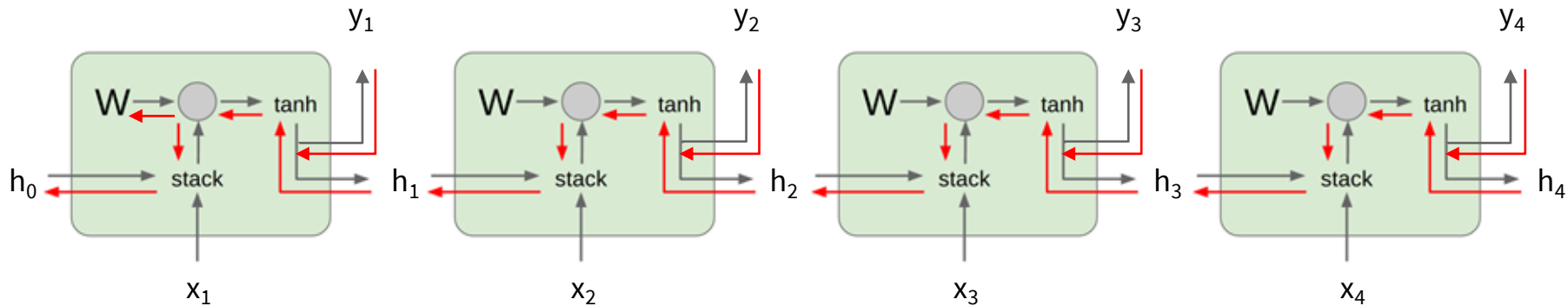


$$\begin{aligned}h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\ &= \tanh\left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \\ &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)\end{aligned}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \tanh'(W_{hh}h_{t-1} + W_{xh}x_t)W_{hh}$$

# Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994  
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

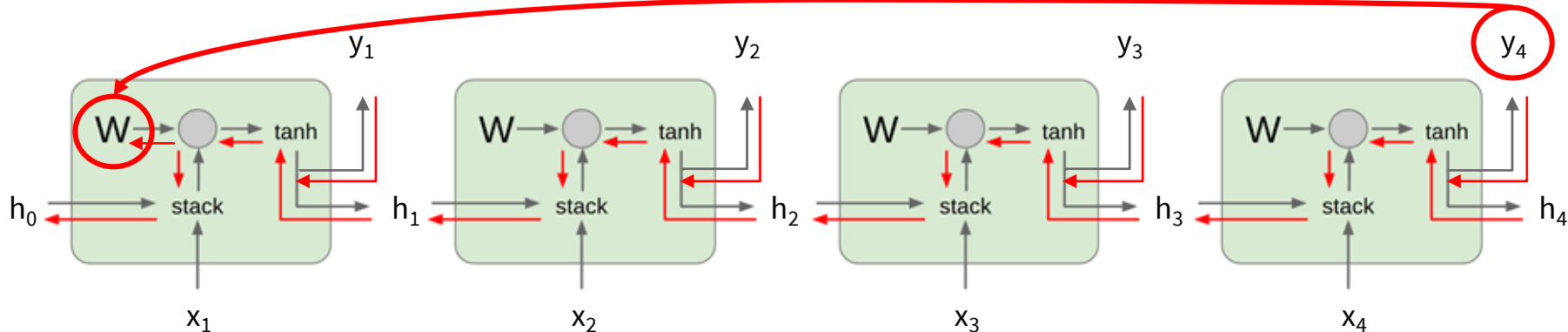


$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

# Vanilla RNN Gradient Flow

Gradients over multiple time steps:

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994  
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



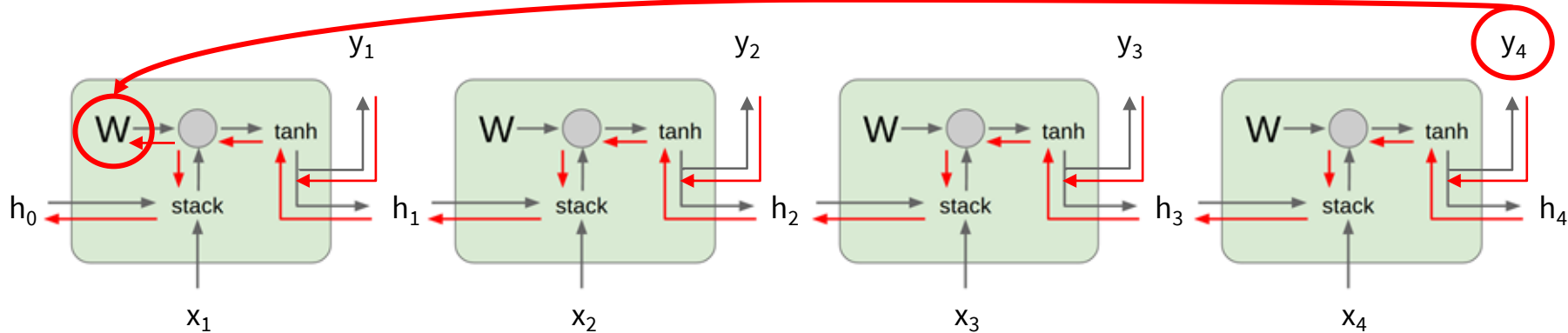
$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \frac{\partial h_T}{\partial h_{T-1}} \cdots \frac{\partial h_1}{\partial W}$$

# Vanilla RNN Gradient Flow

Gradients over multiple time steps:

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994  
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



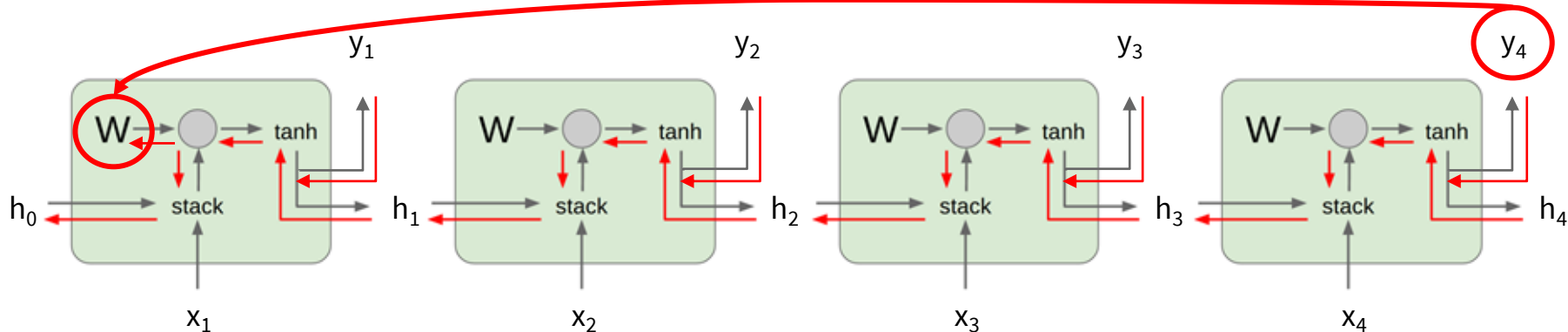
$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \frac{\partial h_T}{\partial h_{T-1}} \cdots \frac{\partial h_1}{\partial W} = \frac{\partial L_T}{\partial h_T} \left( \prod_{t=2}^T \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_1}{\partial W}$$

# Vanilla RNN Gradient Flow

Gradients over multiple time steps:

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994  
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

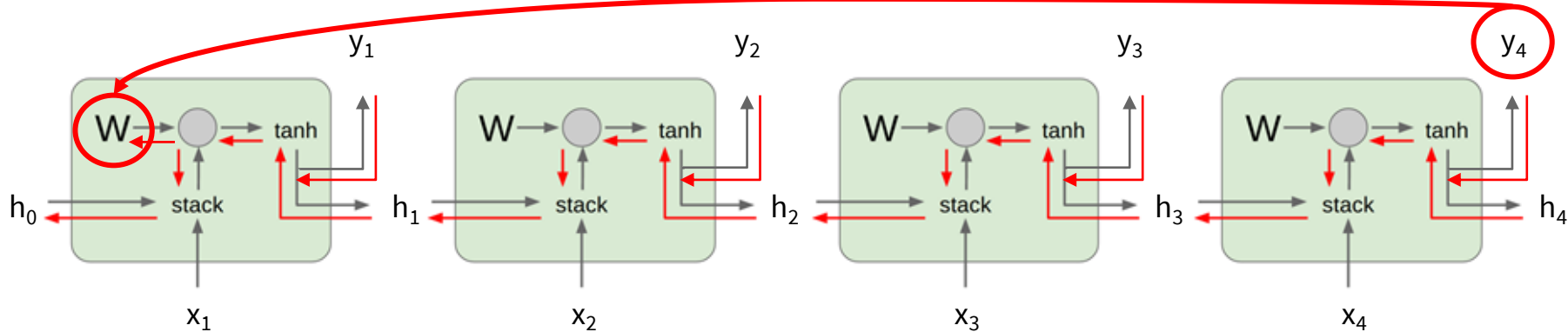
$$\frac{\partial h_t}{\partial h_{t-1}} = \tanh'(W_{hh} h_{t-1} + W_{xh} x_t) W_{hh}$$

$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \frac{\partial h_T}{\partial h_{t-1}} \cdots \frac{\partial h_1}{\partial W} = \frac{\partial L_T}{\partial h_T} \left( \prod_{t=2}^T \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_1}{\partial W}$$

# Vanilla RNN Gradient Flow

Gradients over multiple time steps:

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994  
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

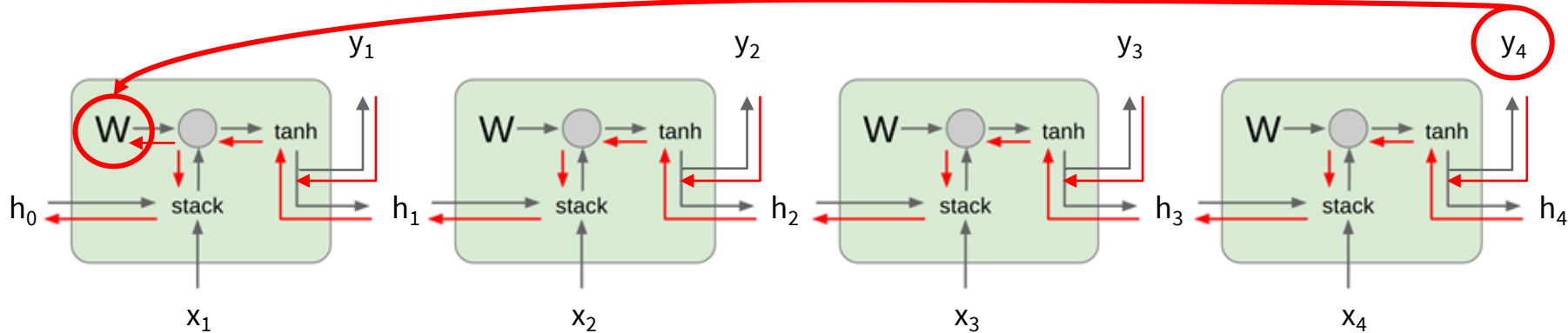
Almost always  $< 1$   
Vanishing gradients

$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \left( \prod_{t=2}^T \tanh'(W_{hh} h_{t-1} + W_{xh} x_t) \right) W_{hh}^{T-1} \frac{\partial h_1}{\partial W}$$

# Vanilla RNN Gradient Flow

Gradients over multiple time steps:

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994  
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



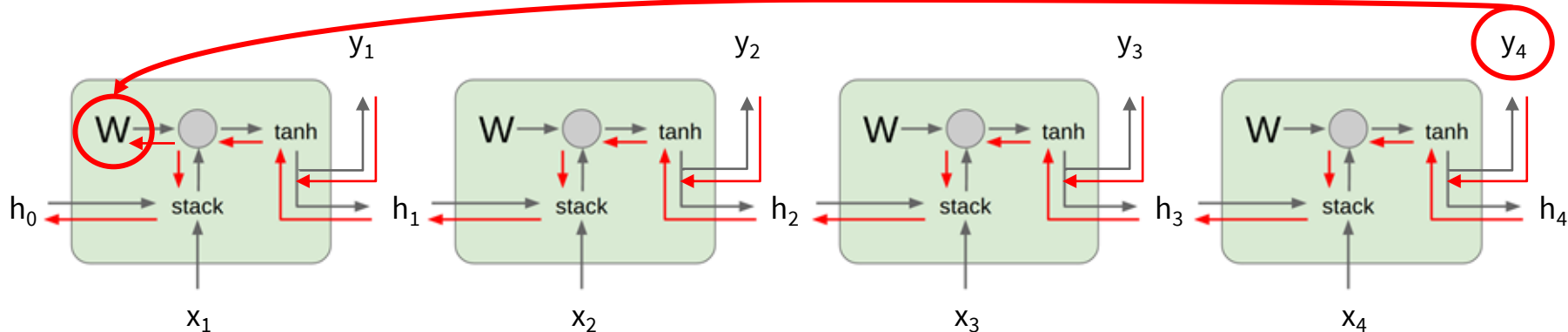
$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

What if we assumed no non-linearity?

# Vanilla RNN Gradient Flow

Gradients over multiple time steps:

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994  
 Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



What if we assumed no non-linearity?

$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

Largest singular value > 1:  
Exploding gradients

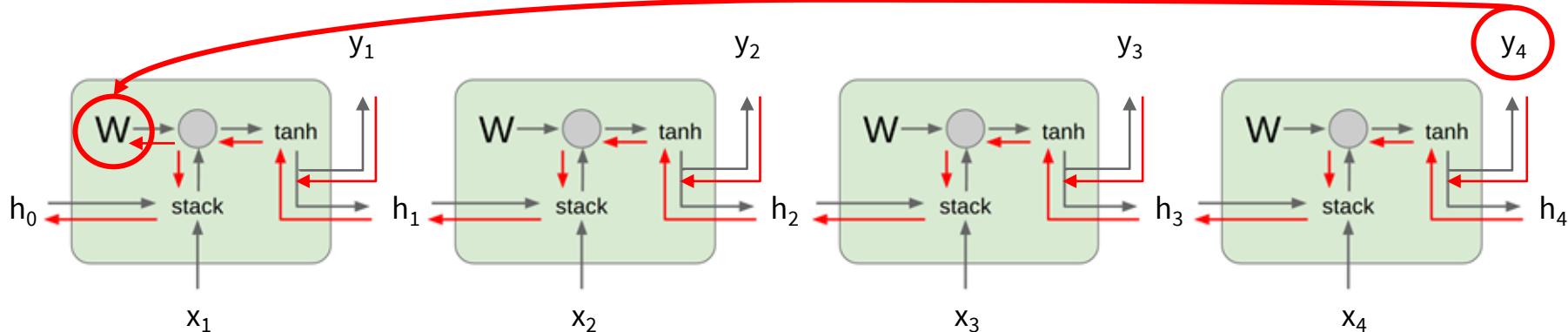
$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \boxed{W_{hh}^{T-1}} \frac{\partial h_1}{\partial W}$$

Largest singular value < 1:  
Vanishing gradients

# Vanilla RNN Gradient Flow

Gradients over multiple time steps:

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994  
 Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



What if we assumed no non-linearity?

$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \boxed{W_{hh}^{T-1}} \frac{\partial h_1}{\partial W}$$

Largest singular value > 1:  
Exploding gradients

Largest singular value < 1:  
Vanishing gradients

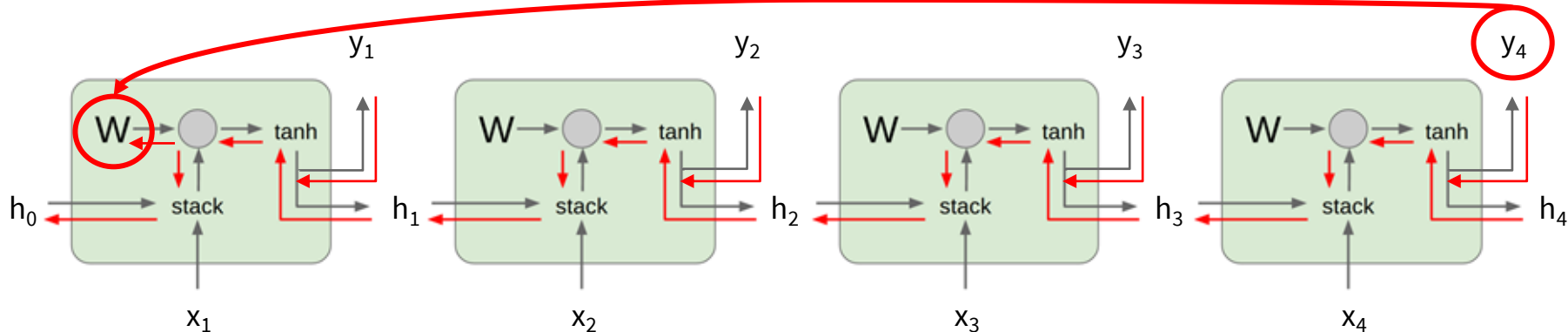
→ Gradient clipping: Scale gradient if its norm is too big

```
grad_norm = np.sum(grad * grad)
if grad_norm > threshold:
    grad *= (threshold / grad_norm)
```

# Vanilla RNN Gradient Flow

Gradients over multiple time steps:

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994  
 Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



What if we assumed no non-linearity?

$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

Largest singular value > 1:  
 Exploding gradients

$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \boxed{W_{hh}^{T-1}} \frac{\partial h_1}{\partial W}$$

Largest singular value < 1:  
 Vanishing gradients

→ Change RNN architecture

# Long Short Term Memory (LSTM) – A Historical Note

Vanilla RNN

$$h_t = \tanh \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

# Long Short Term Memory (LSTM)

Vanilla RNN

$$h_t = \tanh \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

LSTM

Four gates

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

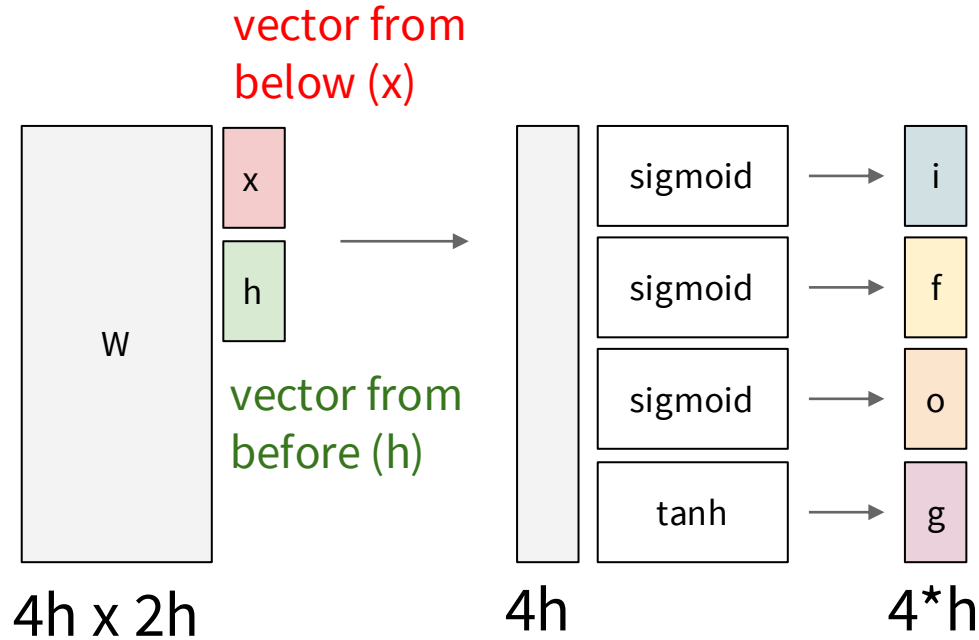
Cell state

Hidden state

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

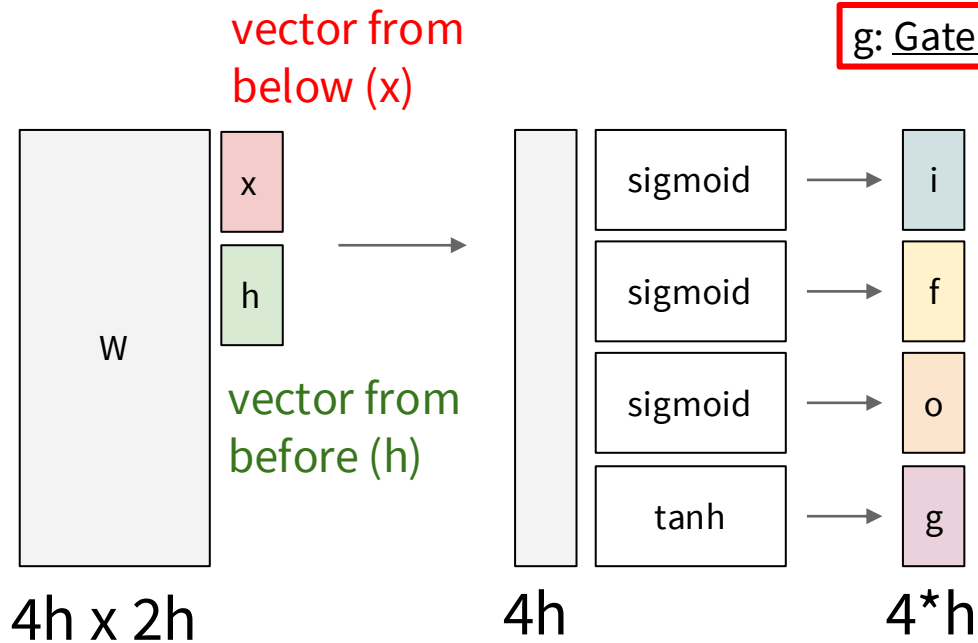
# Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



# Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



$g$ : Gate gate (?), How much to write to cell

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

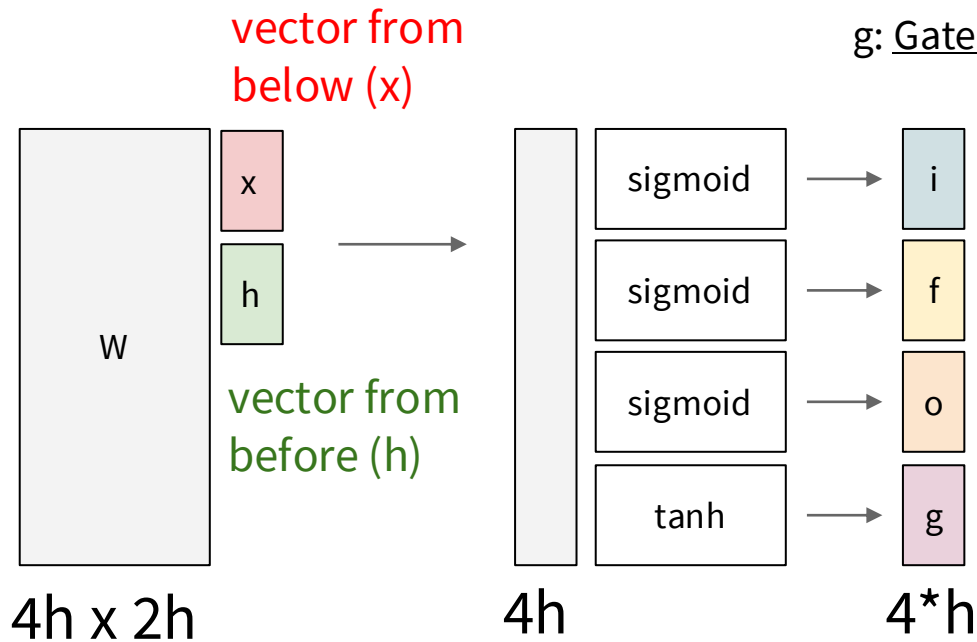
$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

# Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]

i: Input gate, whether to write to cell



g: Gate gate (?), How much to write to cell

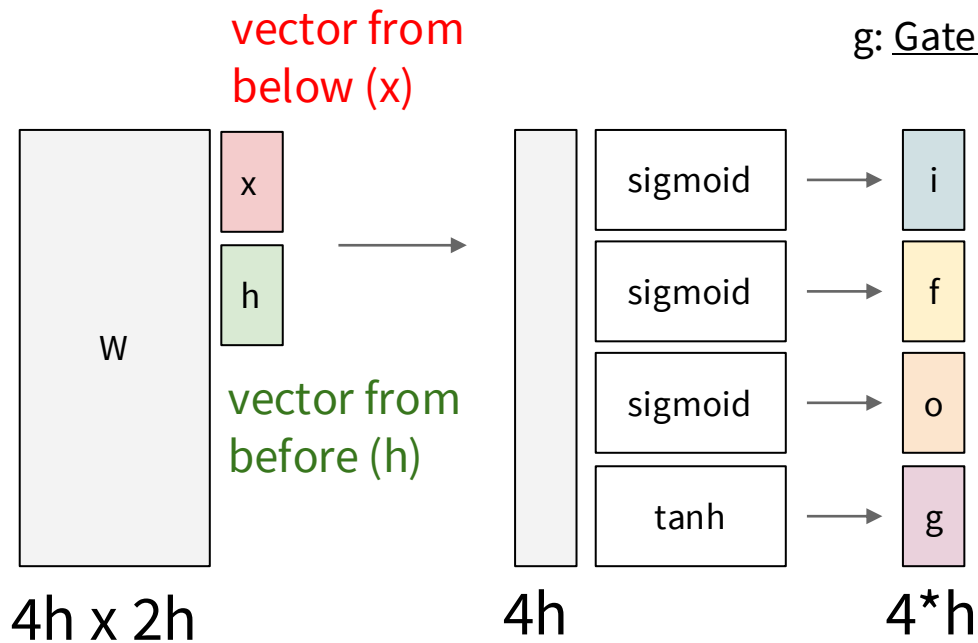
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

# Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



$i$ : Input gate, whether to write to cell

$f$ : Forget gate, Whether to erase cell

$g$ : Gate gate (?), How much to write to cell

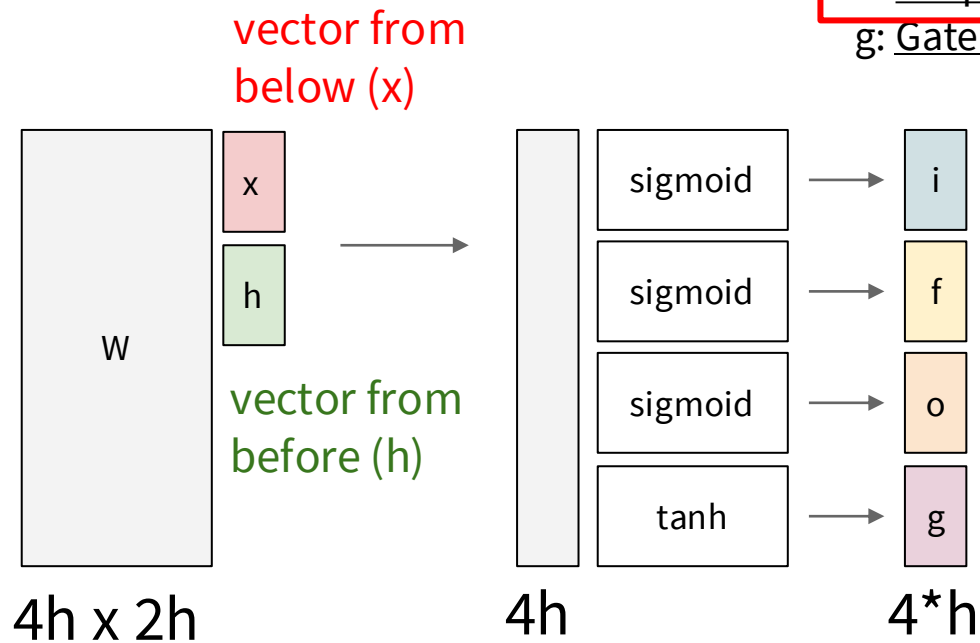
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

# Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



$i$ : Input gate, whether to write to cell

$f$ : Forget gate, Whether to erase cell

$o$ : Output gate, How much to reveal cell

$g$ : Gate gate (?), How much to write to cell

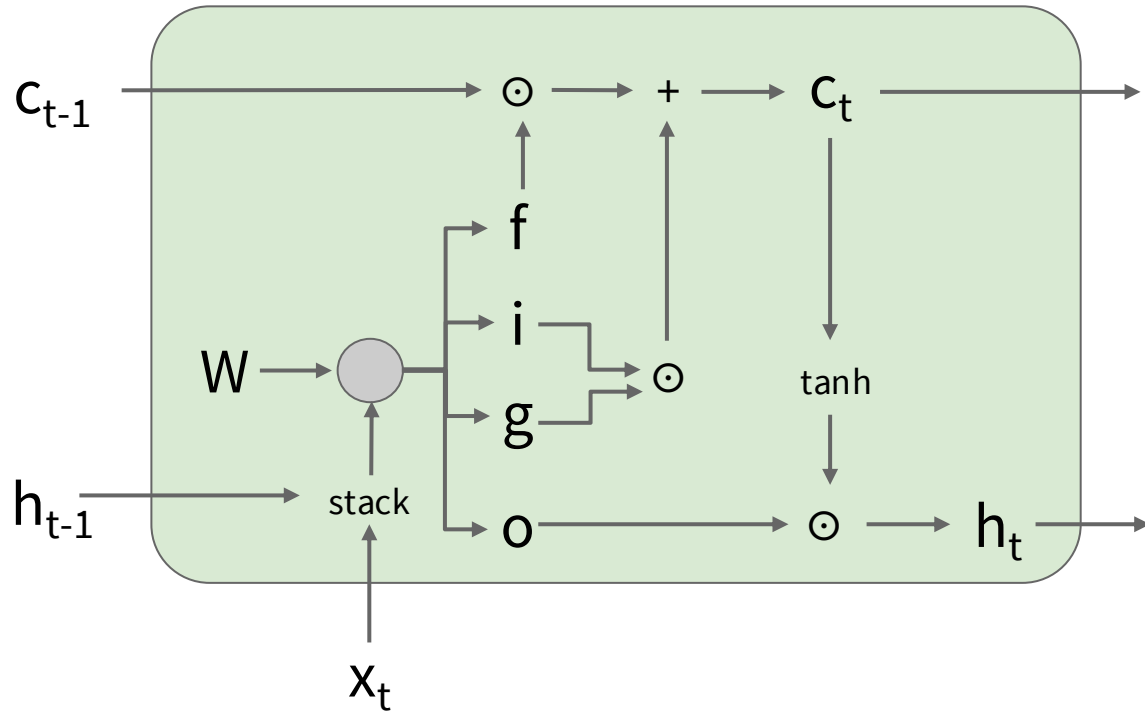
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

# Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



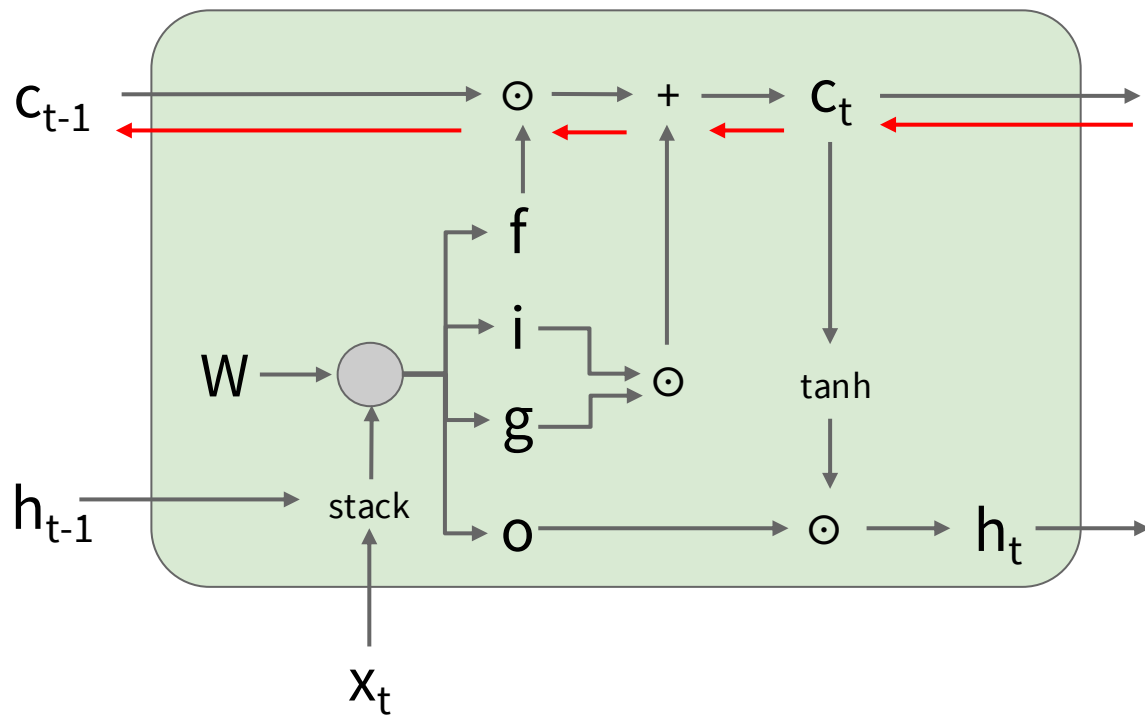
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

# Long Short Term Memory (LSTM): Gradient Flow

[Hochreiter et al., 1997]



Backpropagation from  $c_t$  to  $c_{t-1}$  only elementwise multiplication by  $f$ , no matrix multiply by  $W$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

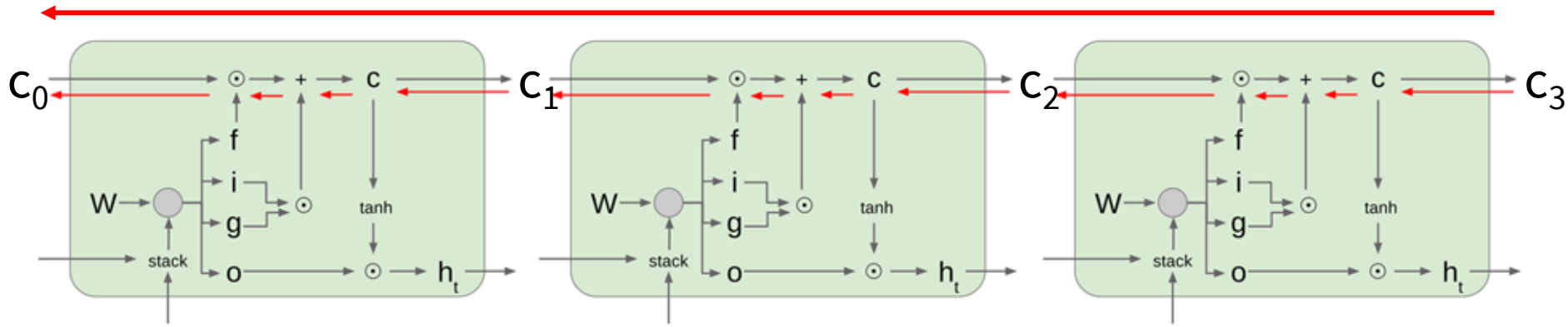
$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

# Long Short Term Memory (LSTM): Gradient Flow

[Hochreiter et al., 1997]

Uninterrupted gradient flow!



# Do LSTMs solve the vanishing gradient problem?

The LSTM architecture makes it easier for the RNN to preserve information over many timesteps

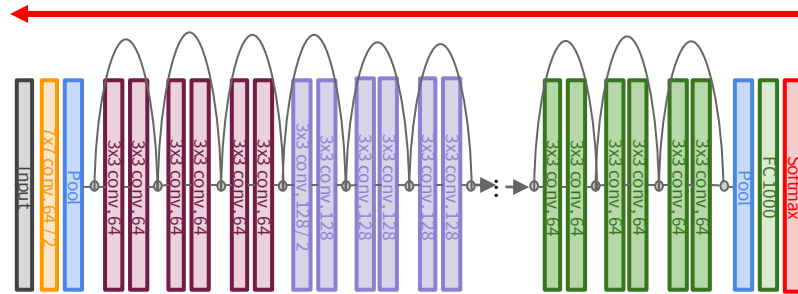
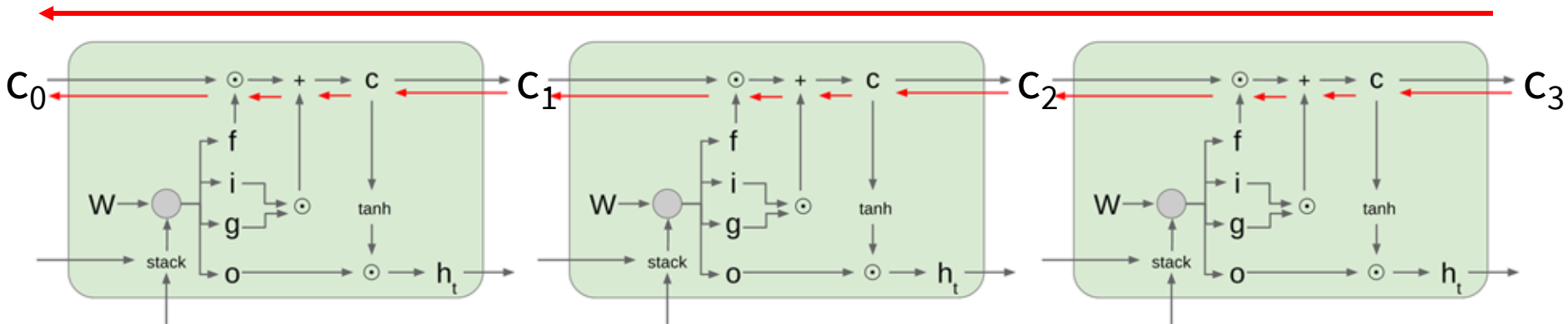
- e.g. if the  $f = 1$  and the  $i = 0$ , then the information of that cell is preserved indefinitely.
- By contrast, it's harder for vanilla RNN to learn a recurrent weight matrix  $W_h$  that preserves info in hidden state

LSTM doesn't guarantee that there is no vanishing/exploding gradient, but it does provide an easier way for the model to learn long-distance dependencies

# Long Short Term Memory (LSTM): Gradient Flow

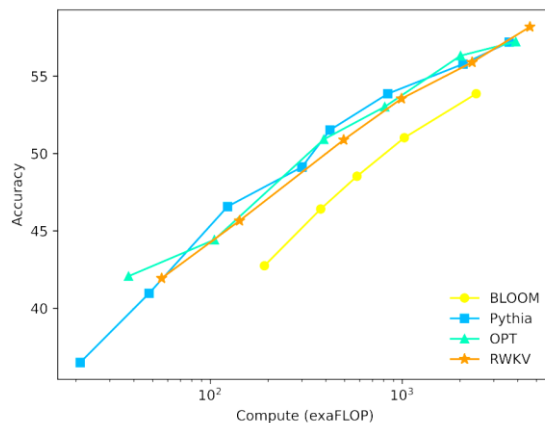
[Hochreiter et al., 1997]

Uninterrupted gradient flow!



# Modern RNNs

- Sometimes called “state space models”
  - Hidden state
- Main advantages:
  - Unlimited context length
  - Compute scales linearly with sequence length



RWKV Scaling ([arXiv](#))

## SIMPLIFIED STATE SPACE LAYERS FOR SEQUENCE MODELING

Mamba: Linear-Time Sequence Modeling with Selective State Spaces

# Summary

- RNNs allow a lot of **flexibility in architecture** design
- Vanilla RNNs are simple but don't work very well
- More complex variants (e.g. LSTMs, Mamba) can introduce ways to **selectively pass** information forward
- Backward flow of gradients in RNN can **explode or vanish**.  
Exploding is controlled with gradient clipping. Backpropagation through time is often needed.

Next time: Attention and Transformers