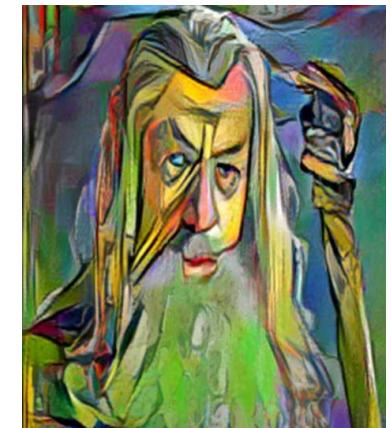
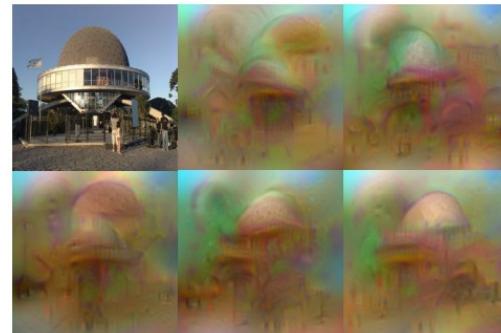
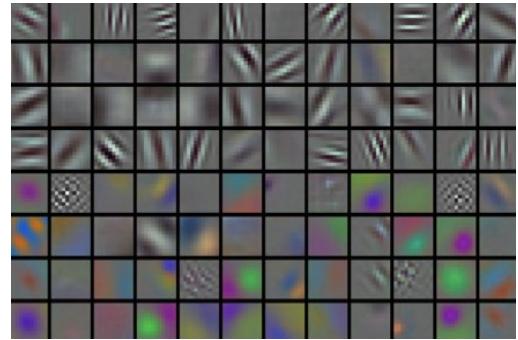


Lecture 10:

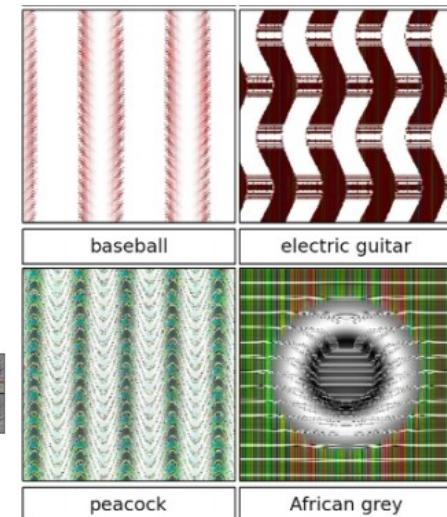
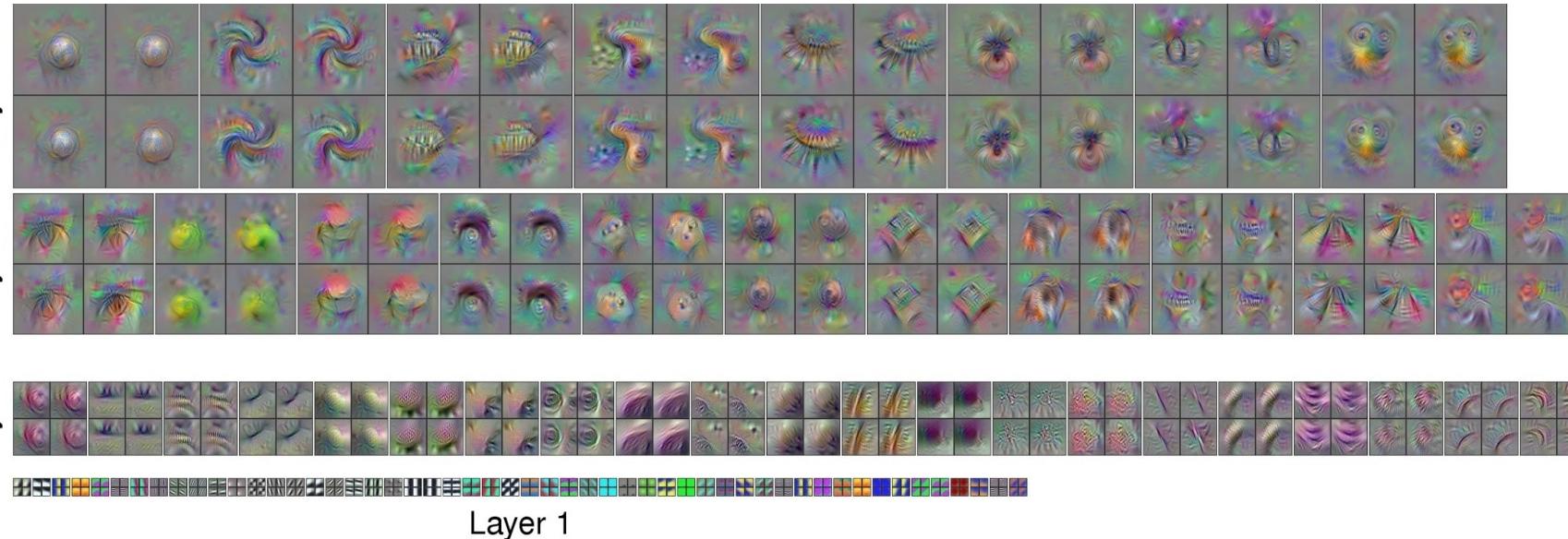
Recurrent Neural Networks

Administrative

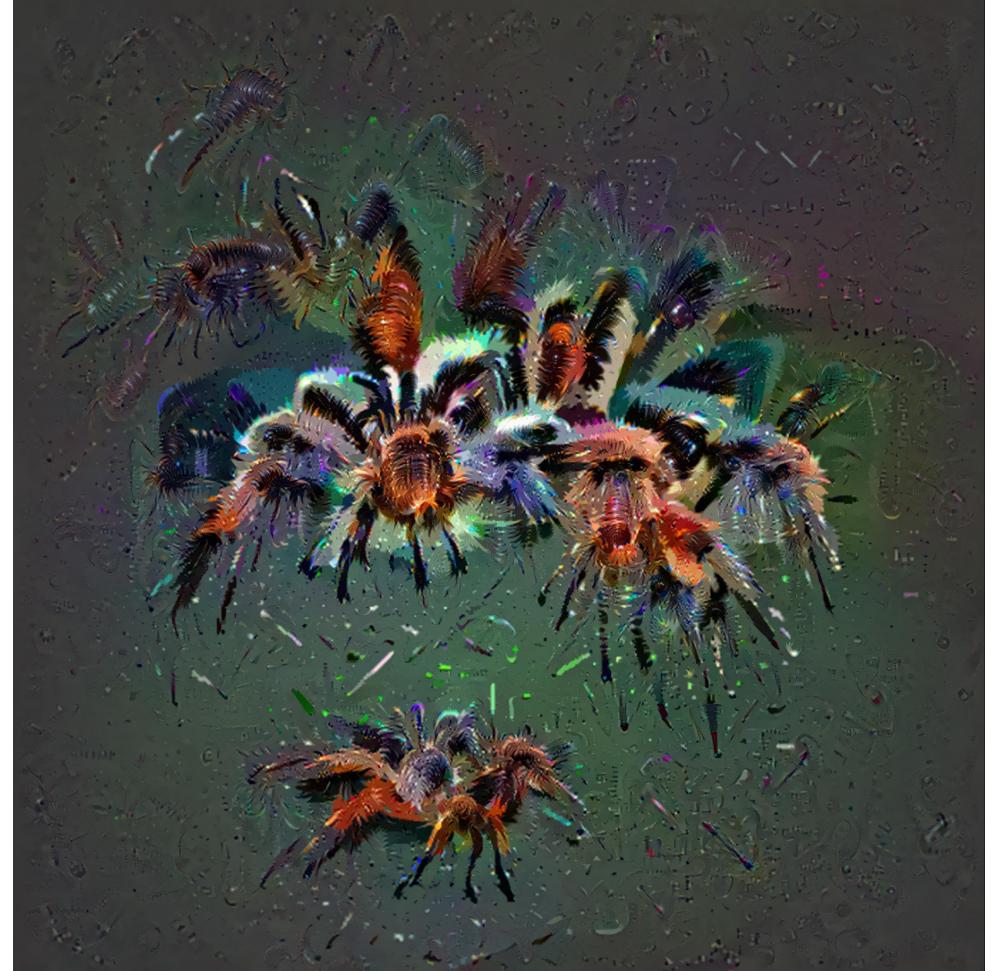
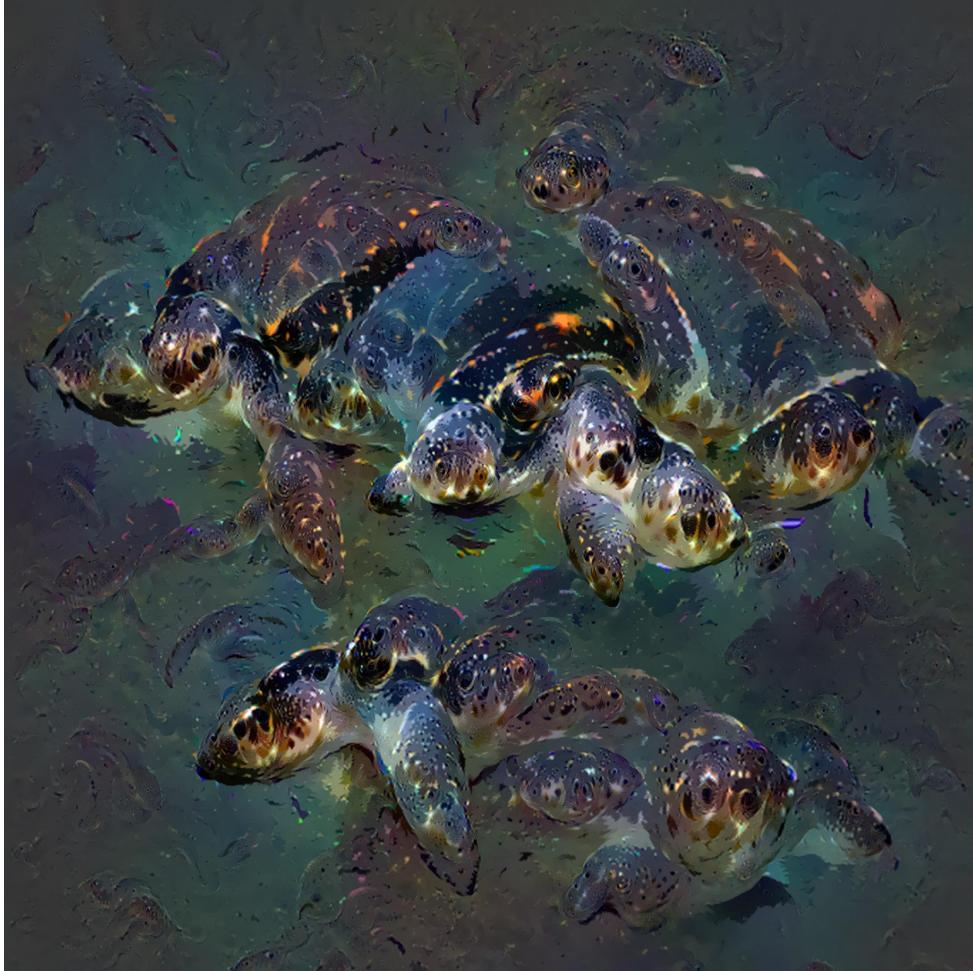
- Midterm this Wednesday! woohoo!
- A3 will be out ~Wednesday



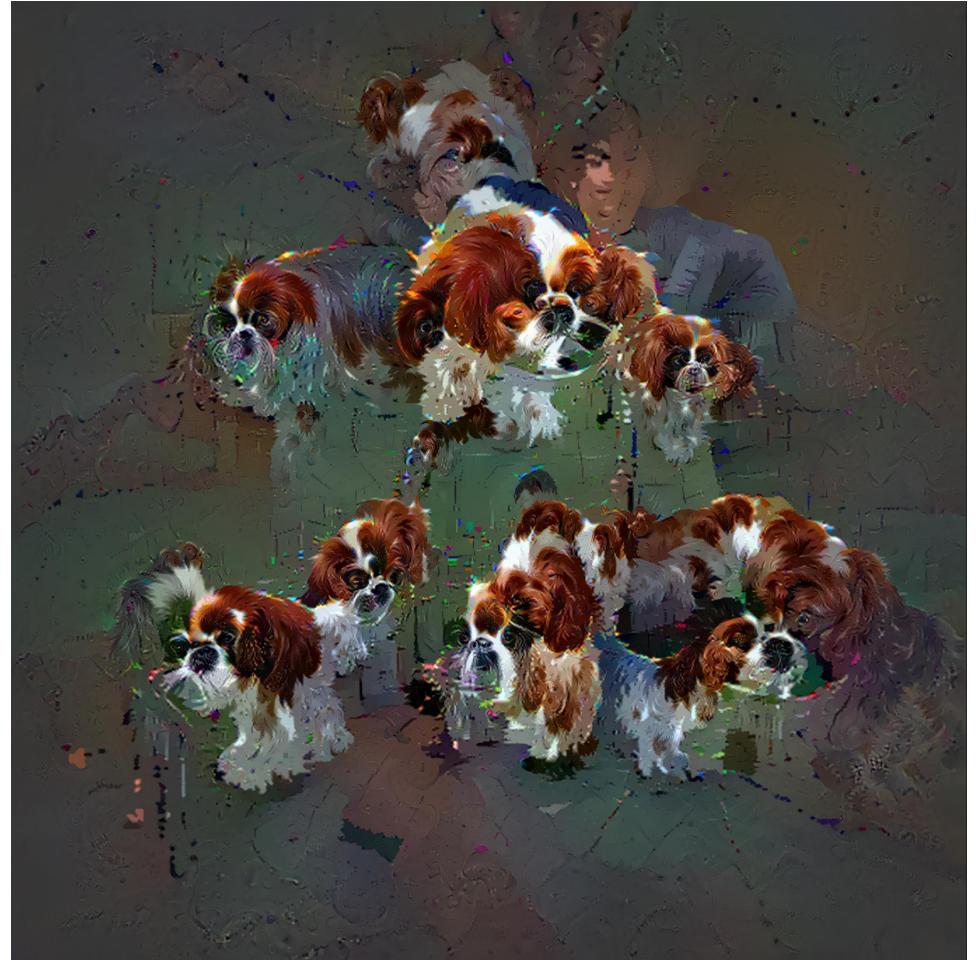
Layer 4
Layer 3
Layer 2
Layer 1



<http://mtyka.github.io/deepdream/2016/02/05/bilateral-class-vis.html>

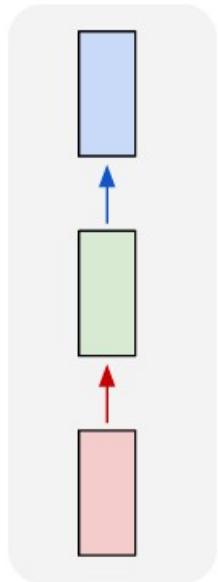


<http://mtyka.github.io/deepdream/2016/02/05/bilateral-class-vis.html>

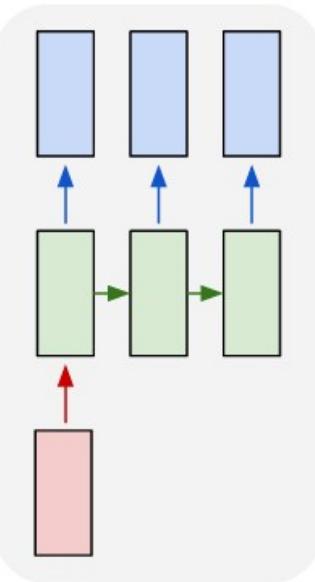


Recurrent Networks offer a lot of flexibility:

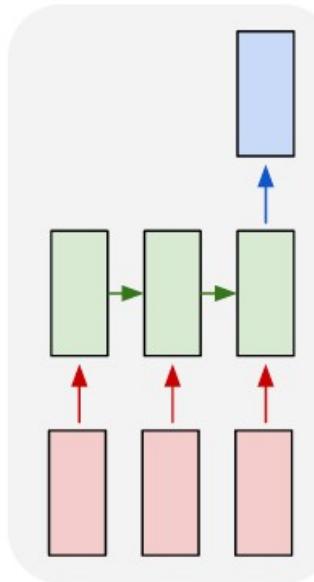
one to one



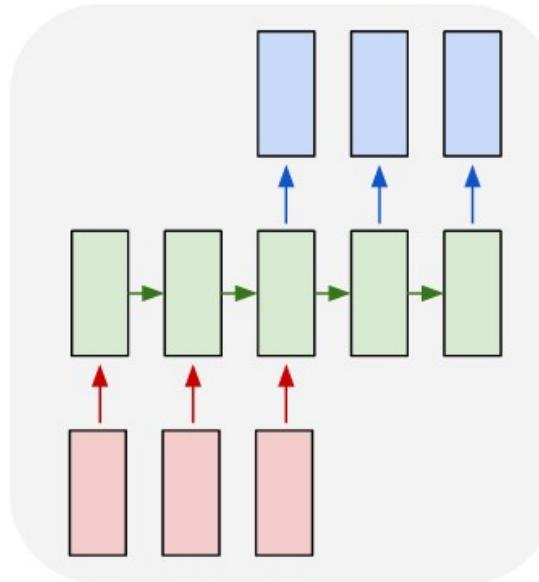
one to many



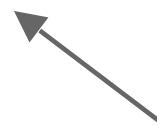
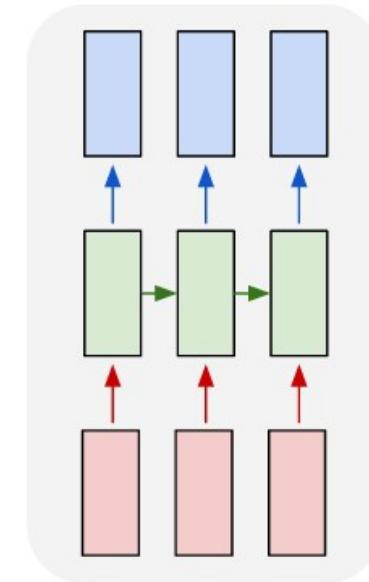
many to one



many to many



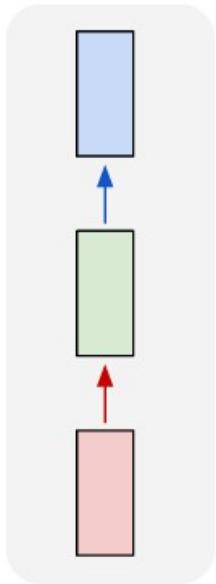
many to many



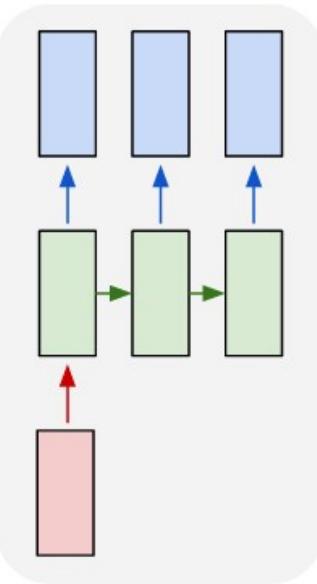
Vanilla Neural Networks

Recurrent Networks offer a lot of flexibility:

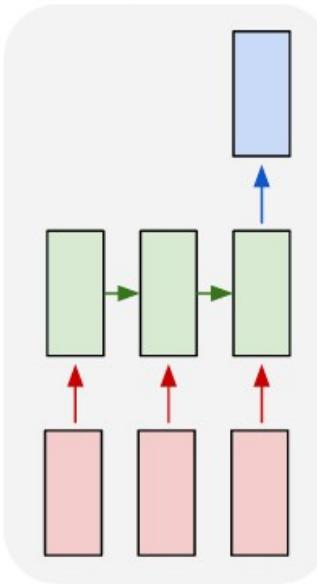
one to one



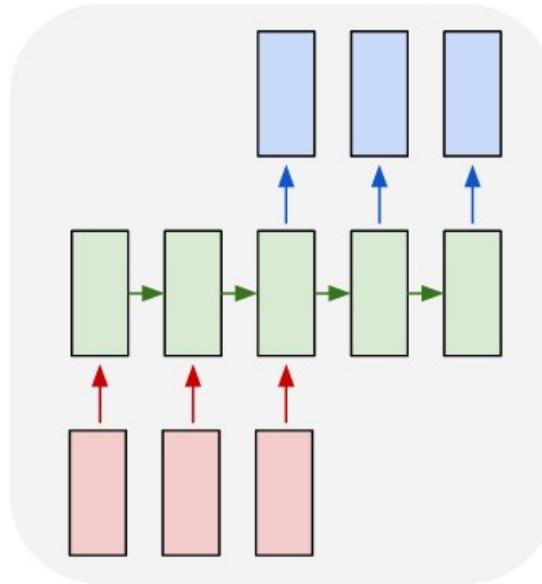
one to many



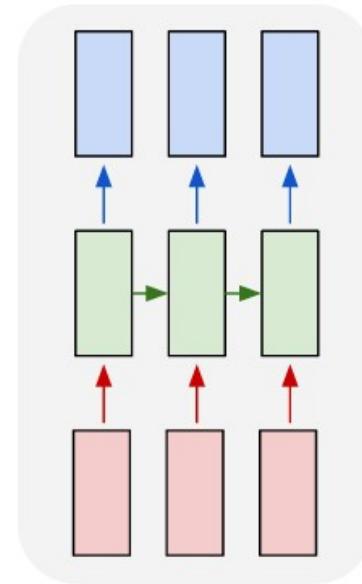
many to one



many to many

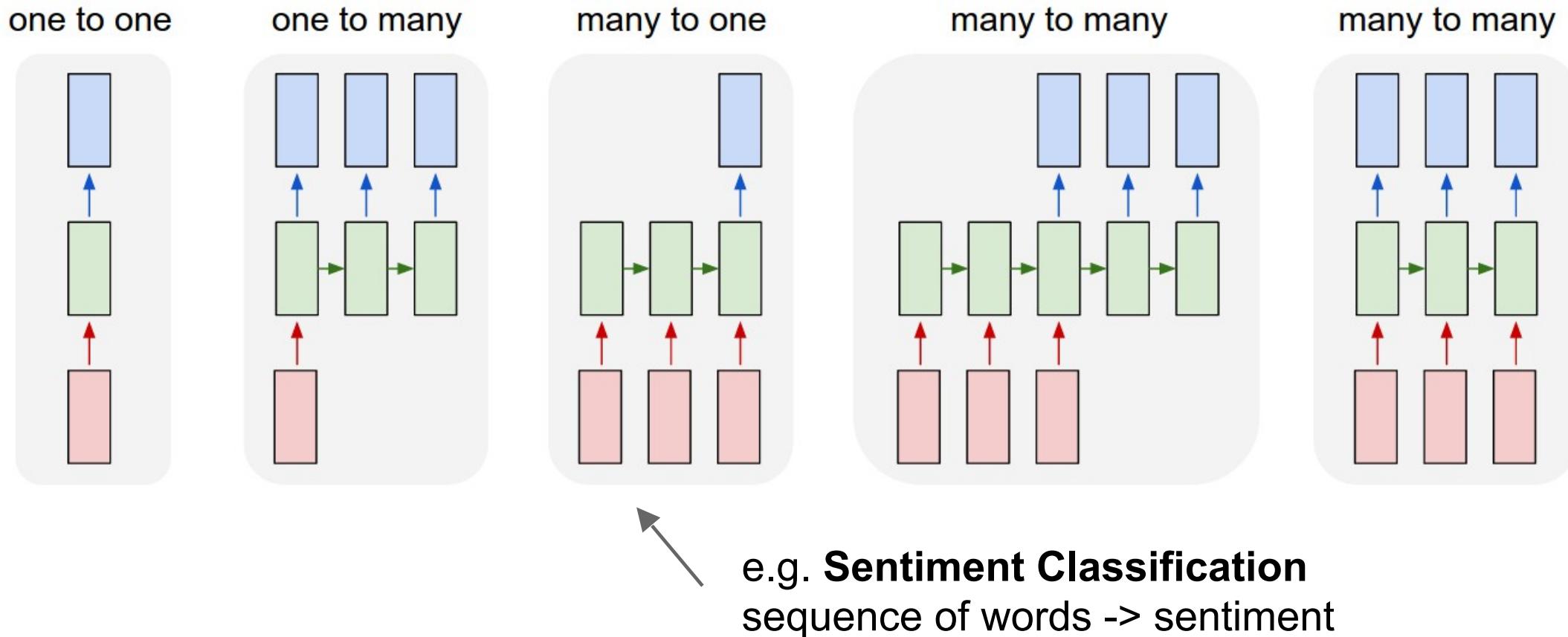


many to many



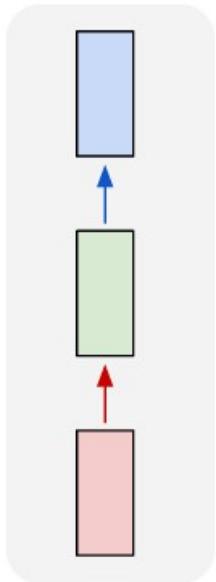
e.g. **Image Captioning**
image -> sequence of words

Recurrent Networks offer a lot of flexibility:

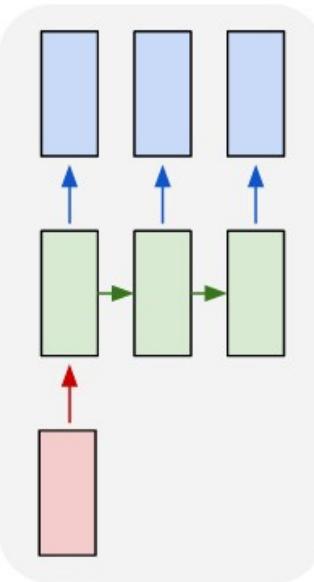


Recurrent Networks offer a lot of flexibility:

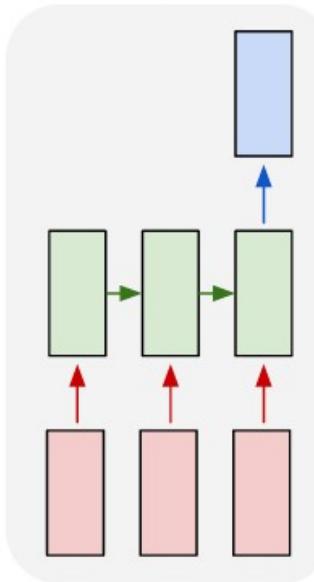
one to one



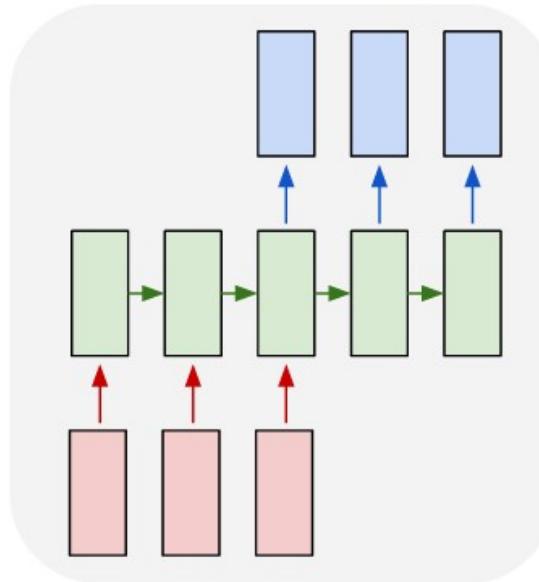
one to many



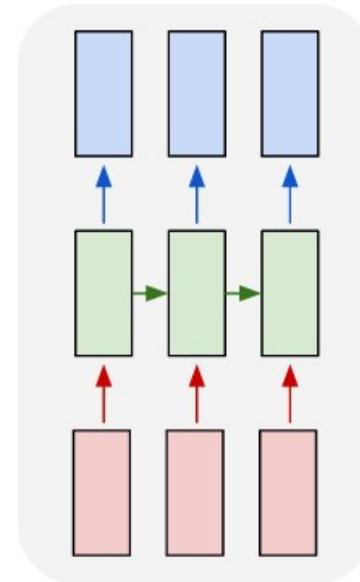
many to one



many to many



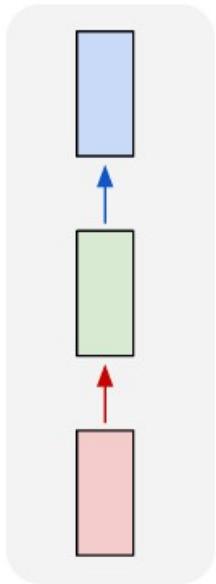
many to many



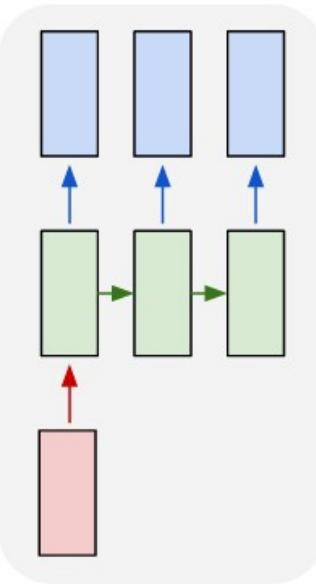
↑
e.g. **Machine Translation**
seq of words -> seq of words

Recurrent Networks offer a lot of flexibility:

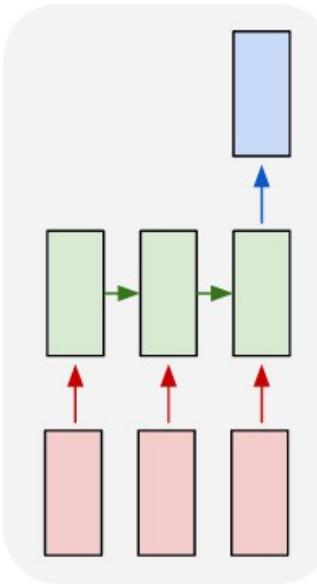
one to one



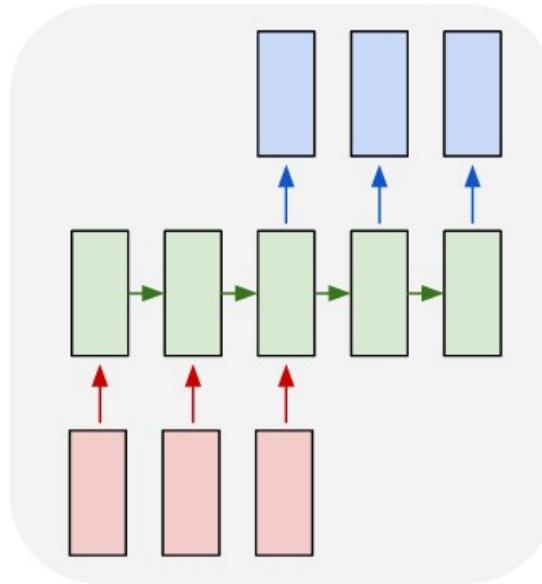
one to many



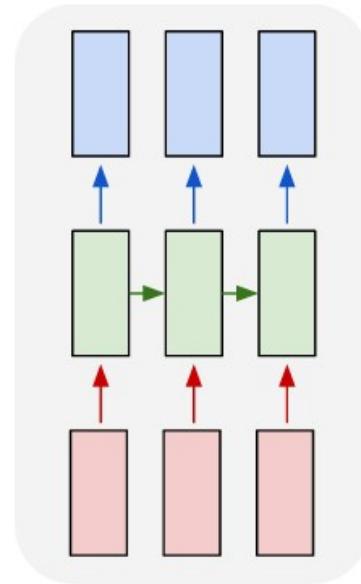
many to one



many to many



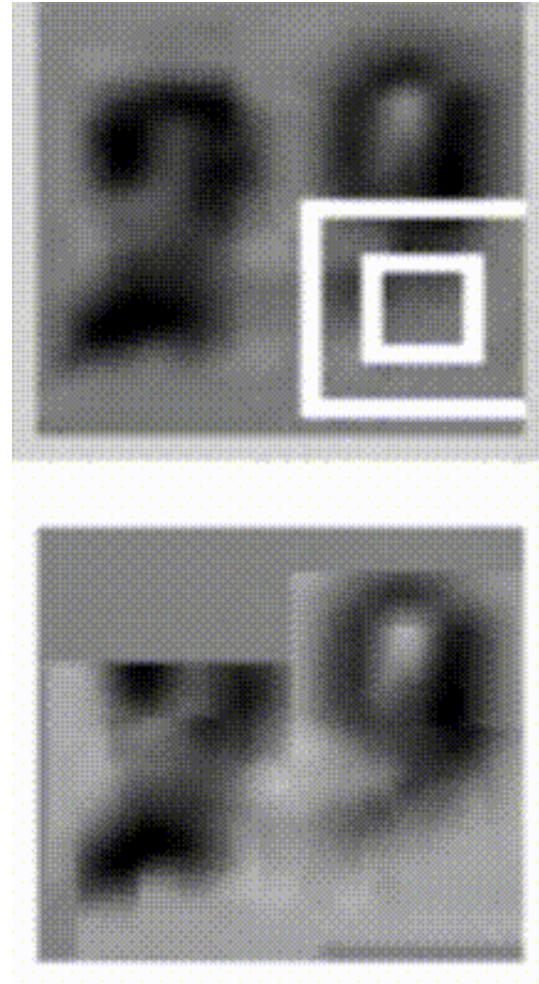
many to many



e.g. **Video classification on frame level**

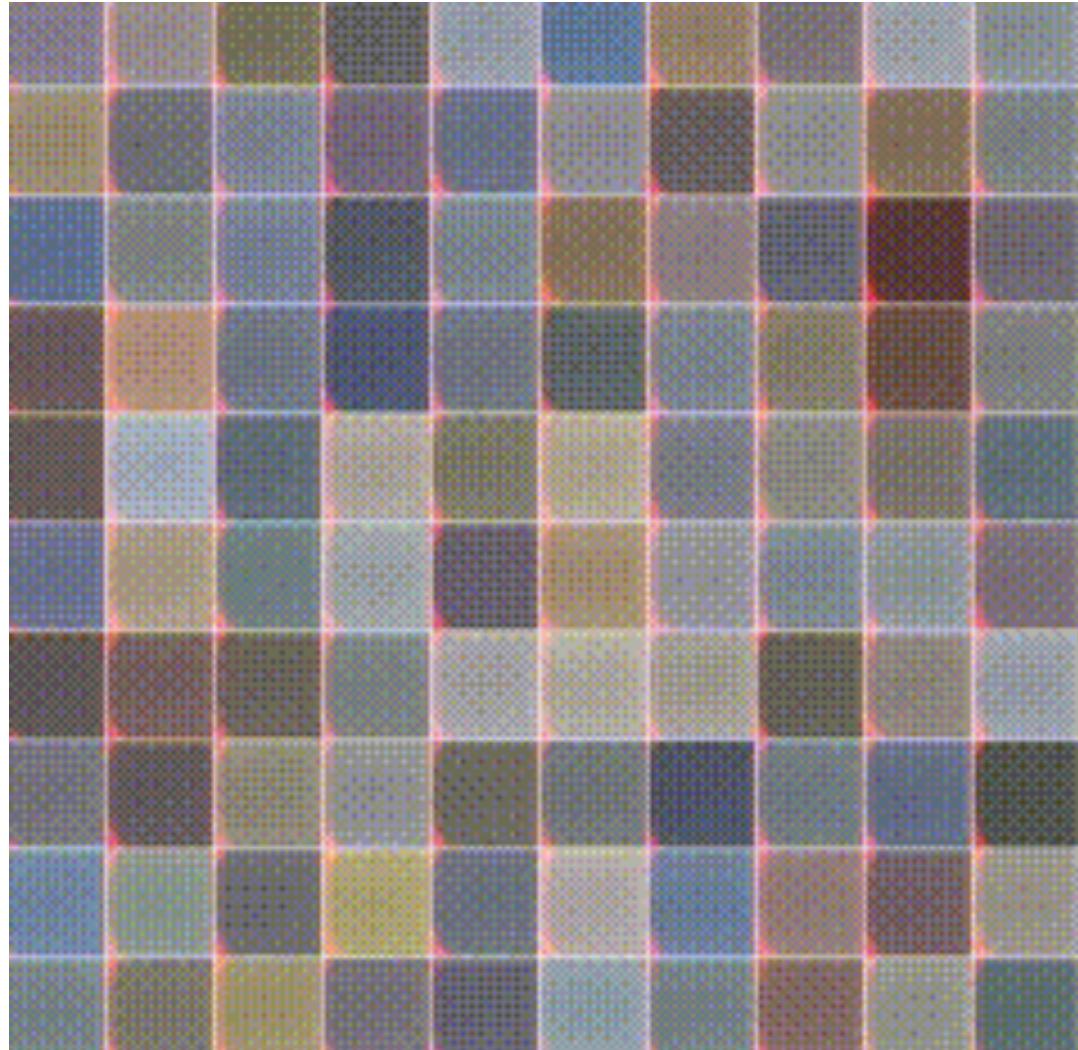
Sequential Processing of fixed inputs

Multiple Object Recognition with
Visual Attention, Ba et al.

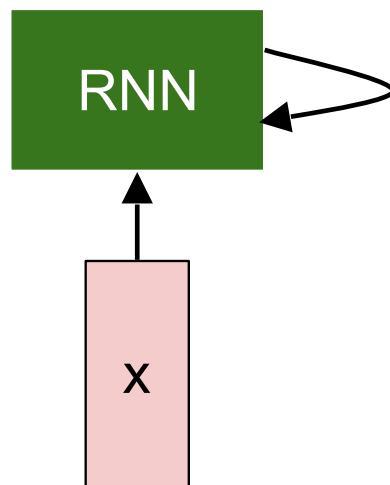


Sequential Processing of fixed outputs

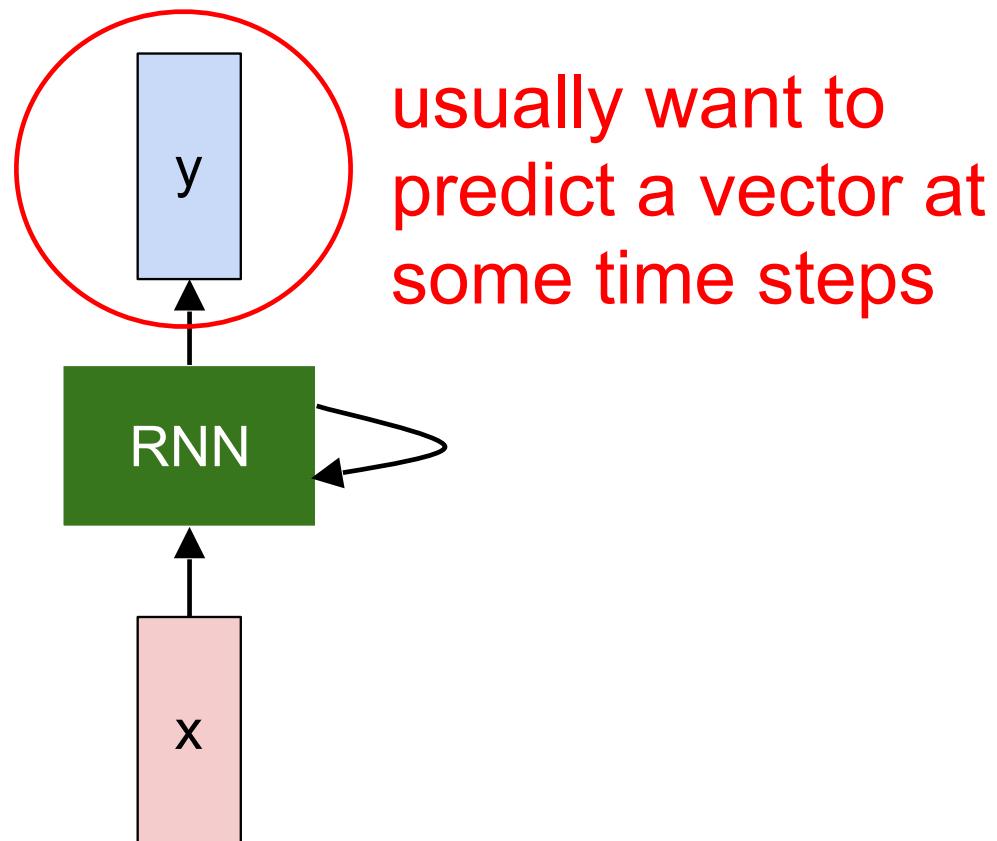
DRAW: A Recurrent
Neural Network For
Image Generation,
Gregor et al.



Recurrent Neural Network



Recurrent Neural Network

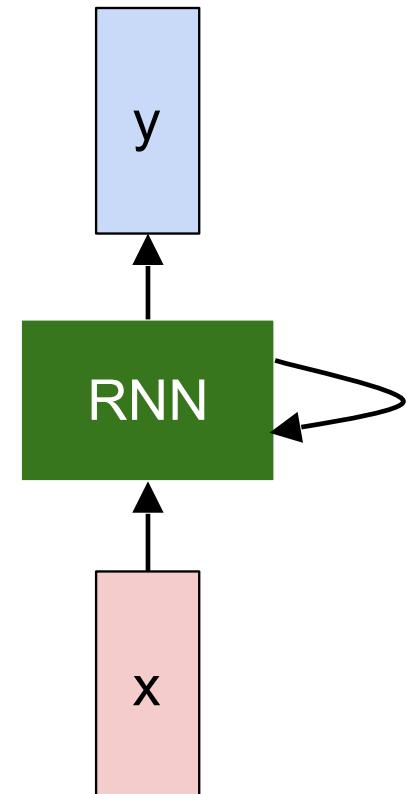


Recurrent Neural Network

We can process a sequence of vectors \mathbf{x} by applying a recurrence formula at every time step:

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

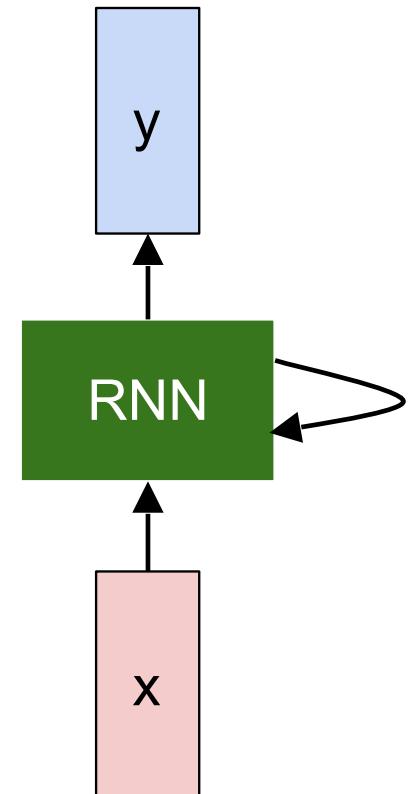
new state / old state input vector at
 \ some time step
 some function
 with parameters W



Recurrent Neural Network

We can process a sequence of vectors \mathbf{x} by applying a recurrence formula at every time step:

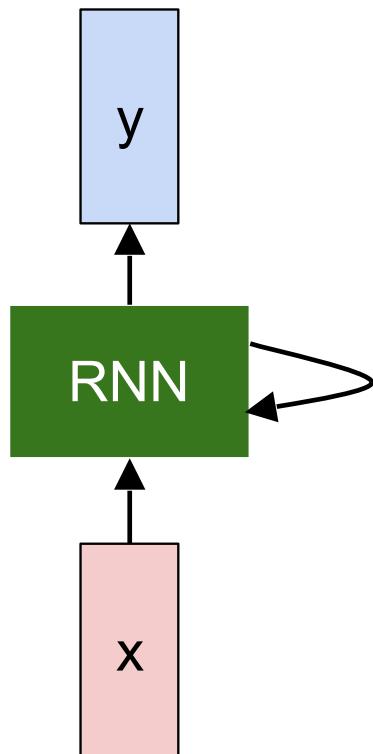
$$h_t = f_W(h_{t-1}, x_t)$$



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 \mathbf{h} :



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



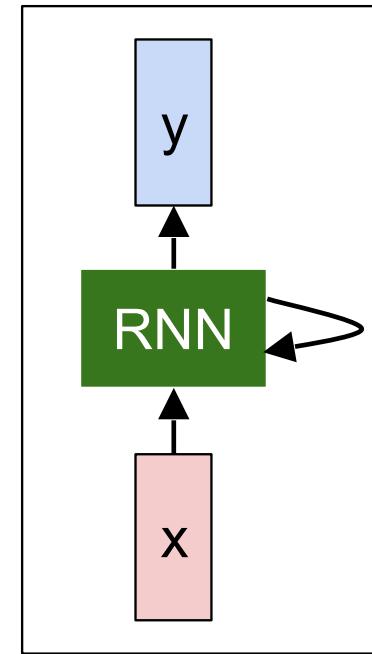
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

Character-level language model example

Vocabulary:
[h,e,l,o]

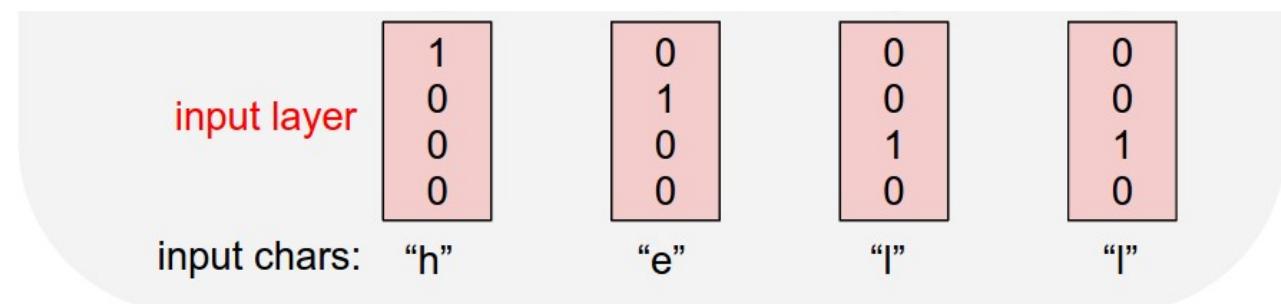
Example training
sequence:
“hello”



Character-level language model example

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”

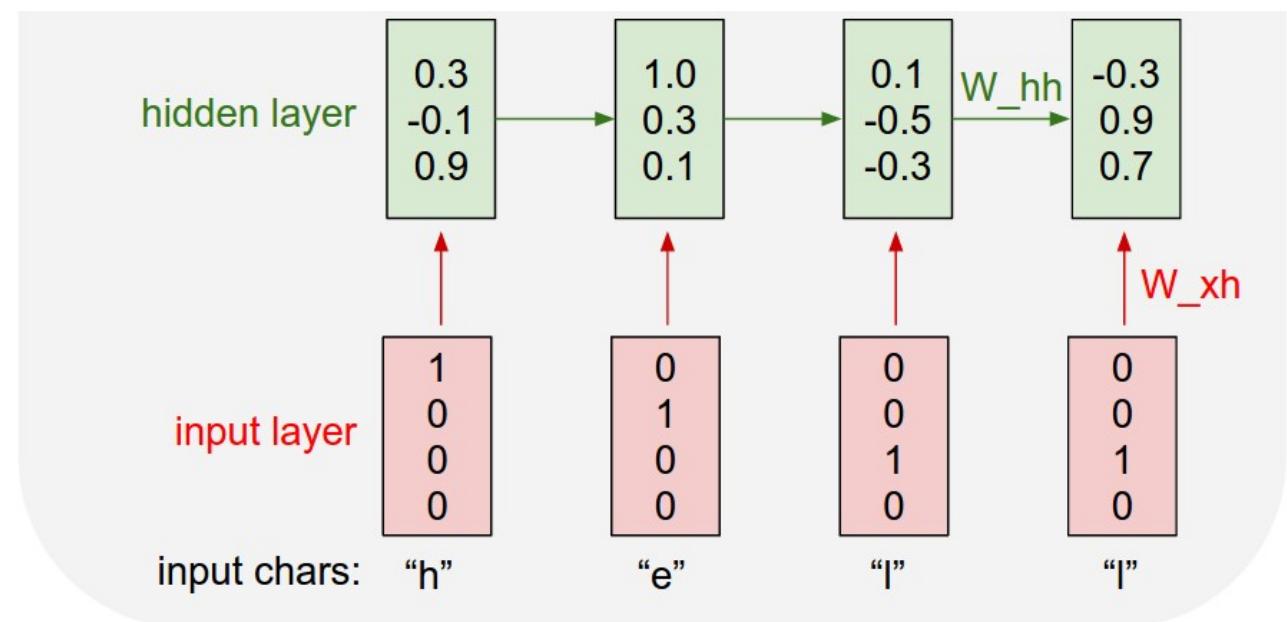


Character-level language model example

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”

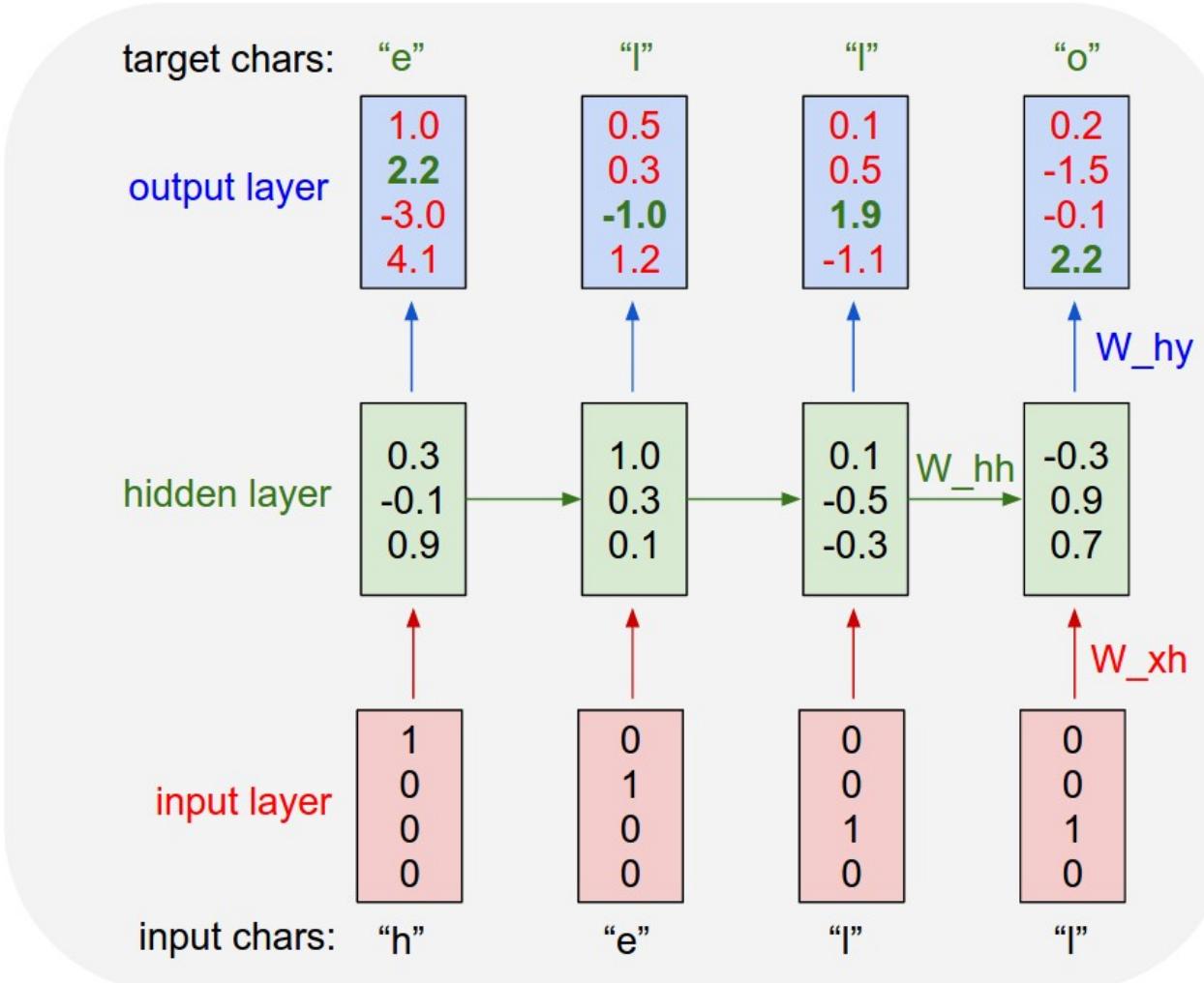
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$



Character-level language model example

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”



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

```
"""
Minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)
BSD License
"""

import numpy as np

# data I/O
data = open('input.txt', 'r').read() # should be simple plain text file
chars = list(set(data))
data_size, vocab_size = len(data), len(chars)
print('data has %d characters, %d unique.' % (data_size, vocab_size))
char_to_ix = { ch:i for i,ch in enumerate(chars) }
ix_to_char = { i:ch for i,ch in enumerate(chars) }

# hyperparameters
hidden_size = 100 # size of hidden layer of neurons
seq_length = 25 # number of steps to unroll the RNN for
learning_rate = 1e-1

# model parameters
wxh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
bh = np.zeros((hidden_size, 1)) # hidden bias
by = np.zeros((vocab_size, 1)) # output bias

def lossFun(inputs, targets, hprev):
    """
    inputs,targets are both list of integers.
    hprev is Hx1 array of initial hidden state
    returns the loss, gradients on model parameters, and last hidden state
    """
    xs, hs, ys, ps = {}, {}, {}, {}
    hs[-1] = np.copy(hprev)
    loss = 0
    # forward pass
    for t in xrange(len(inputs)):
        xs[t] = np.zeros((vocab_size,1)) # encode in 1-of-k representation
        xs[t][inputs[t]] = 1
        hs[t] = np.tanh(np.dot(wxh, xs[t]) + np.dot(whh, hs[t-1]) + bh) # hidden state
        ys[t] = np.dot(why, hs[t]) + by # unnormalized log probabilities for next chars
        ps[t] = np.exp(ys[t]) / np.sum(np.exp(ys[t])) # probabilities for next chars
        loss += -np.log(ps[t][targets[t],0]) # softmax (cross-entropy loss)
        # backward pass: compute gradients going backwards
        dwxh, dwhh, dwhy = np.zeros_like(wxh), np.zeros_like(whh), np.zeros_like(why)
        dbh, dby = np.zeros_like(bh), np.zeros_like(by)
        dhnext = np.zeros_like(hs[0])
        for t in reversed(xrange(len(inputs))):
            dy = np.copy(ps[t])
            dy[targets[t]] -= 1 # backprop into y
            dhy += np.dot(dy, hs[t].T)
            dby += dy
            dh = np.dot(why.T, dy) + dhnext # backprop into h
            dhran = (1 - hs[t] * hs[t].T) * dh # backprop through tanh nonlinearity
            dbh += dhran
            dwxh += np.dot(dhran, xs[t].T)
            dwhh += np.dot(dhran, hs[t-1].T)
            dhnext = np.dot(whh.T, dhran)
            for dparam in [dwxh, dwhh, dwhy, dbh, dby]:
                np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
        return loss, dwxh, dwhh, dwhy, dbh, dby, hs[len(inputs)-1]

def sample(h, seed_ix, n):
    """
    sample a sequence of integers from the model
    h is memory state, seed_ix is seed letter for first time step
    """
    x = np.zeros((vocab_size, 1))
    x[seed_ix] = 1
    ixes = []
    for t in xrange(n):
        h = np.tanh(np.dot(wxh, x) + np.dot(whh, h) + bh)
        y = np.dot(why, h) + by
        p = np.exp(y) / np.sum(np.exp(y))
        ix = np.random.choice(range(vocab_size), p=p.ravel())
        x = np.zeros((vocab_size, 1))
        x[ix] = 1
        ixes.append(ix)
    return ixes

n, p = 0, 0
mxvh, mwvh, mhvh = np.zeros_like(wxh), np.zeros_like(whh), np.zeros_like(why)
mbh, mbv = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
while True:
    if p+seq_length+1 >= len(data) or n == 0:
        hprev = np.zeros((hidden_size,1)) # reset RNN memory
        p = 0 # go from start of data
    inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
    targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
    # sample from the model now and then
    if n % 100 == 0:
        sample_ix = sample(hprev, inputs[0], 200)
        txt = ''.join(ix_to_char[ix] for ix in sample_ix)
        print '----\n%s\n----' % (txt, )
    # forward seq_length characters through the net and fetch gradient
    loss, dwxh, dwhh, dwhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
    smooth_loss = smooth_loss * 0.999 + loss * 0.001
    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
    # perform parameter update with Adagrad
    for param, dparam, mem in zip([wxh, whh, why, bh, by],
                                   [dwxh, dwhh, dwhy, dbh, dby],
                                   [mxvh, mwvh, mhvh, mbh, mbv]):
        mem += dparam * dparam
        param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
    p += seq_length # move data pointer
    n += 1 # iteration counter
```

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

min-char-rnn.py gist

```
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 data = open('input.txt', 'r').read() # should be simple plain text file
9 chars = list(set(data))
10 data_size, vocab_size = len(data), len(chars)
11 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
12 char_to_ix = { ch:i for i,ch in enumerate(chars) }
13 ix_to_char = { i:ch for i,ch in enumerate(chars) }
14
15 # hyperparameters
16 hidden_size = 100 # size of hidden layer of neurons
17 seq_length = 25 # number of steps to unroll the RNN for
18 learning_rate = 1e-1
19
20 # model parameters
21 Wh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
22 Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
23 Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
24 bh = np.zeros(hidden_size, 1) # hidden bias
25 by = np.zeros(vocab_size, 1) # output bias
26
27 def lossFun(inputs, targets, hprev):
28     """
29     inputs, targets are both list of integers.
30     hprev is hxt array of initial hidden state
31     returns the loss, gradients on model parameters, and last hidden state
32     """
33     xs, hs, ys, ps = [], (), (), []
34     hprev[1] = np.copy(hprev)
35     loss = 0
36     for t in range(len(inputs)):
37         x = inputs[t]
38         x1 = np.zeros((vocab_size, 1)) # encode in 1-of-k representation
39         x1[x] = 1
40         hprev[0] = np.dot(Wh, x1) + np.dot(bh, hprev[1]) + bh # hidden state
41         y1 = np.tanh(np.dot(Whh, hprev[1]) + by) # unnormalized log probabilities for next chars
42         ps1 = np.exp(y1) / np.sum(np.exp(y1)) # probabilities for next chars
43         loss += -np.log(ps1[targets[t], 1]) # softmax (cross-entropy loss)
44         # backward pass: compute gradients going backwards
45         dprev_dh = np.zeros_like(Whh) # gradients for hidden-to-hidden
46         dprev_dy = np.zeros_like(Wh) # gradients for input-to-hidden
47         dbh = np.zeros_like(bh) # gradients for hidden bias
48         dy = np.zeros_like(by) # gradients for output bias
49         dhnext = np.zeros_like(hprev[1])
50         for t in reversed(xrange(len(inputs))):
51             dy += np.copy(ps1)
52             dy[targets[t]] -= 1 # backprop into y
53             dby += dy
54             dh = np.dot(Why, T, dy) + dhnext # backprop into h
55             ddraw = (1 - hprev[1] * hprev[1]) * dh # backprop through tanh nonlinearity
56             dWhh += np.dot(ddraw, hprev[1].T)
57             dWh += np.dot(ddraw, x1.T)
58             dbh += np.dot(ddraw, hprev[1])
59             dhnext = np.dot(Whh, hprev[1])
60         for dparam in [dWh, dWhh, dby, dbh, ddy]:
61             np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
62         return loss, dxh, dh, dmy, dbh, dby, hs[-1]
63     def sample(seed_ix, n):
64         """
65         sample a sequence of integers from the model
66         n is memory state, seed_ix is seed letter for first time step
67         """
68         x = np.zeros((vocab_size, 1))
69         x[seed_ix] = 1
70         ixes = []
71         for t in range(n):
72             h = np.tanh(np.dot(Wh, x) + np.dot(bh, h) + bh)
73             y = np.dot(Why, h) + by
74             p = np.exp(y) / np.sum(np.exp(y))
75             ix = np.random.choice(range(vocab_size), p=p.ravel())
76             x[1] = 0
77             ixes.append(ix)
78         return ixes
79
80 n, p = 0, 0
81 dxh, dWh, dmy, dbh, dby = np.zeros_like(Wh), np.zeros_like(Why), np.zeros_like(bh)
82 dWhh, dbyh, dbhy = np.zeros_like(Whh), np.zeros_like(Wh), np.zeros_like(by) # memory variables for Adagrad
83 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
84 while True:
85     # prepare inputs (we're sweeping from left to right in steps seq_length long)
86     if p+seq_length >= len(data) or n == 0:
87         hprev[0] = np.zeros((hidden_size, 1)) # reset RNN memory
88         p = 0 # go from start of data
89     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
90     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
91
92     # sample from the model now and then
93     if n % 100 == 0:
94         s = np.random.sample(hprev, inputs[0], 200)
95         txt = ''
96         for ix in sample(hprev, inputs[0], 200):
97             txt += ix_to_char[ix] # fetch next char
98         print '-----%s-----%s' % (txt, )
99
100    # forward seq_length characters through the net and fetch gradient
101    loss, dxh, dWhh, dWh, dby, dbh, dbyh, hprev = lossFun(inputs, targets, hprev)
102    smooth_loss = smooth_loss * 0.999 + loss * 0.001
103    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
104
105    # perform parameter update with Adagrad
106    for param, dparam, mem in zip([dxh, dWh, dmy, dbh, dby],
107                                 [dWhh, dWh, dbhy, dbyh, dby],
108                                 [mem, mem, mem, mem, mem]):
109        mem += dparam * dparam
110        param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adgrad update
111
112    p += seq_length # move data pointer
113    n += 1 # iteration counter
```

Data I/O

```
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 data = open('input.txt', 'r').read() # should be simple plain text file
9 chars = list(set(data))
10 data_size, vocab_size = len(data), len(chars)
11 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
12 char_to_ix = { ch:i for i,ch in enumerate(chars) }
13 ix_to_char = { i:ch for i,ch in enumerate(chars) }
14
15 # sample from the model now and then
16 if n % 100 == 0:
17     s = np.random.sample(hprev, inputs[0], 200)
18     txt = ''
19     for ix in sample(hprev, inputs[0], 200):
20         txt += ix_to_char[ix] # fetch next char
21     print '-----%s-----%s' % (txt, )
22
23    # forward seq_length characters through the net and fetch gradient
24    loss, dxh, dWhh, dWh, dby, dbh, dbyh, hprev = lossFun(inputs, targets, hprev)
25    smooth_loss = smooth_loss * 0.999 + loss * 0.001
26    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
27
28    # perform parameter update with Adagrad
29    for param, dparam, mem in zip([dxh, dWh, dmy, dbh, dby],
30                                 [dWhh, dWh, dbhy, dbyh, dby],
31                                 [mem, mem, mem, mem, mem]):
32        mem += dparam * dparam
33        param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adgrad update
34
35    p += seq_length # move data pointer
36    n += 1 # iteration counter
```

min-char-rnn.py gist

```

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 data = open('input.txt', 'r').read() # should be simple plain text file
9 chars = list(set(data))
10 data_size, vocab_size = len(data), len(chars)
11 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
12 char_to_ix = { ch:i for i,ch in enumerate(chars) }
13 ix_to_char = { i:ch for i,ch in enumerate(chars) }

# hyperparameters
hidden_size = 100 # size of hidden layer of neurons
seq_length = 25 # number of steps to unroll the RNN for learning
learning_rate = 1e-1

# model parameters
Wxh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
bh = np.zeros((hidden_size, 1)) # hidden bias
by = np.zeros((vocab_size, 1)) # output bias

def lossFun(inputs, targets, hprev):
    """
    inputs, targets are both list of integers.
    hprev is Hx1 array of initial hidden state
    returns the loss, gradients on model parameters, and last hidden state
    """
    xs, hs, ys, ps = [], (), (), []
    loss = 0
    for t in range(len(inputs)):
        # forward pass
        for t in xrange(len(inputs)):
            xs[t] = inputs[t] # encode in 1-of-k representation
            hs[t] = np.zeros(hidden_size) if t == 0 else np.dot(Wh, hs[t-1]) + bh # hidden state
            ys[t] = np.dot(Why, hs[t]) + by # unnormalized log probabilities for next chars
            ps[t] = np.exp(ys[t]) / np.sum(np.exp(ys[t])) # probabilities for next chars
            loss += -np.log(ps[t][targets[t], 0]) # softmax (cross-entropy loss)
        # backward pass: compute gradients going backwards
        dprev_h = np.zeros_like(hprev)
        dh = np.zeros_like(bh)
        dbh = np.zeros_like(bh)
        dyh = np.zeros_like(Why)
        dhy = np.zeros_like(Why)
        dby = np.zeros_like(by)
        dnext_h = np.zeros_like(hs[0])
        for t in reversed(xrange(len(inputs))):
            dy = np.copy(ps[t])
            dy[targets[t]] -= 1 # backprop into y
            dhy += np.dot(dy, Why.T)
            dby += dy
            dh = np.dot(dyh, Wh.T) + dnext_h # backprop into h
            dhs = (1 - hs[t] * hs[t]) * dh # backprop through tanh nonlinearity
            dWh = np.dot(dhs, hs[t-1].T)
            dnext_h = np.dot(Wh, dhs)
        for dparam in [dparam, -5, out=dparam]: # clip to mitigate exploding gradients
            np.clip(dparam, -5, 5, out=dparam)
        return loss, dxh, ddmh, ddyh, dbh, dby, hs[0], inputs[-1]

def sample(hseed, seed_ix, n):
    """
    sample a sequence of integers from the model
    h is memory state, seed_ix is seed letter for first time step
    """
    x = np.zeros((vocab_size, 1))
    x[seed_ix] = 1
    ixes = []
    for t in range(n):
        h = np.tanh(np.dot(Wxh, x) + np.dot(Wh, h) + bh)
        y = np.dot(Why, h) + by
        p = np.exp(y) / np.sum(np.exp(y))
        ix = np.random.choice(range(vocab_size), p=p.ravel())
        x = np.zeros((vocab_size, 1))
        x[ix] = 1
        ixes.append(ix)
    return ixes

n, p = 0, 0
smooth, smooth_m, mhy = np.zeros_like(Wh), np.zeros_like(Why), np.zeros_like(Why)
dbh, dby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
smooth_mss = np.zeros((vocab_size, seq_length)) # loss at iteration 0
while True:
    # prepare inputs (we're sweeping from left to right in steps seq_length long)
    if p+seq_length+1 >= len(data) or n == 0:
        hprev = np.zeros((hidden_size, 1)) # reset RNN memory
        p = 0 # go from start of data
        inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
        targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
    else:
        # sample from the model now and then
        if n % 100 == 0:
            smooth += smooth_mss
            smooth_mss = np.zeros((vocab_size, seq_length))
            print '-'*50, 'n', n, '-'*50
            txt = '...'.join(ix_to_char[i] for i in sample_ix)
            print '-'*50, 'n', n, '-'*50, '(txt, )'
        # forward seq_length characters through the net and fetch gradient
        loss, dxh, ddmh, ddyh, dbh, dby, hprev = lossFun(inputs, targets, hprev)
        smooth_mss += smooth * 0.999 + loss * 0.001
        if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_mss)
        print progress
        n += 1
        # perform parameter update with Adagrad
        for param, dparam, mem in zip([Wh, Why, bh, by],
                                      [dWh, dWhy, dbh, dby],
                                      [smooth, smooth_m, dbh, dby]):
            mem += dparam * dparam
            param -= learning_rate * dparam / np.sqrt(mem + 1e-8) # adgrad update
        p += seq_length # move data pointer
        n += 1 # iteration counter
    
```

Initializations

```

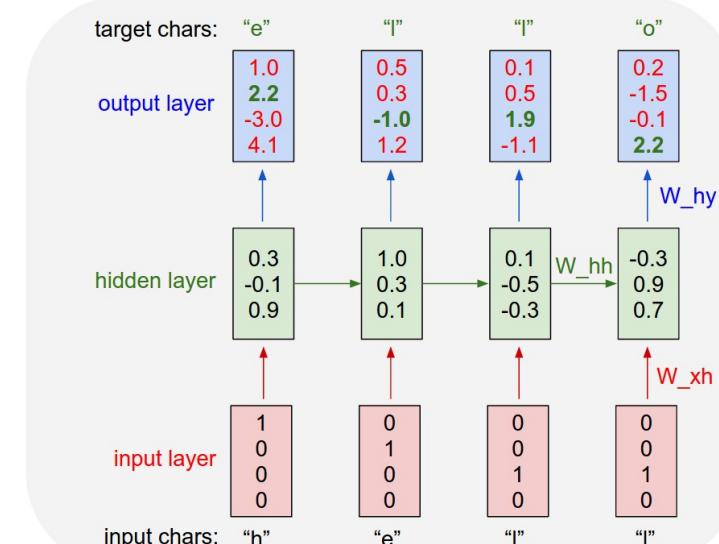
15 # hyperparameters
16 hidden_size = 100 # size of hidden layer of neurons
17 seq_length = 25 # number of steps to unroll the RNN for learning
18 learning_rate = 1e-1

19 # model parameters
20 Wxh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
21 Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
22 Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
23 bh = np.zeros((hidden_size, 1)) # hidden bias
24 by = np.zeros((vocab_size, 1)) # output bias

25

```

recall:



min-char-rnn.py gist

```
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 data = open('input.txt', 'r').read() # should be simple plain text file
9 chars = list(set(data))
10 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
11 char_to_ix = { ch:i for i,ch in enumerate(chars) }
12 ix_to_char = { i:ch for i,ch in enumerate(chars) }
13
14 # hyperparameters
15 hidden_size = 100 # size of hidden layer of neurons
16 seq_length = 25 # number of steps to unroll the RNN for
17 learning_rate = 1e-1
18
19 # model parameters
20 Wkh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
21 Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
22 Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
23 bh = np.zeros((hidden_size, 1)) # hidden bias
24 by = np.zeros((vocab_size, 1)) # output bias
25
26 def lossFun(inputs, targets, hprev):
27     """
28     inputs, targets are both list of integers.
29     hprev is hxi array of initial hidden state
30     returns the loss, gradients on model parameters, and last hidden state
31     """
32     xs, hs, ys, ps = [], [], [], []
33     hs[-1] = np.copy(hprev)
34     loss = 0
35     for t in xrange(len(inputs)):
36         x = inputs[t]
37         # forward pass
38         xh = np.zeros((vocab_size, 1)) # encode in 1-of-k representation
39         xh[inputs[t]] = 1
40         whh_t = np.dot(Whh, hs[-1]) + bh # hidden state
41         why_t = np.dot(Why, whh_t) + by # unnormalized log probabilities for next chars
42         ps_t = np.exp(why_t) / np.sum(np.exp(why_t)) # probabilities for next chars
43         loss += -np.log(ps_t[targets[t]]) # softmax (cross-entropy loss)
44         # backward pass: compute gradients going backwards
45         dprev_dh = np.zeros_like(whh_t)
46         dprev_dbh = np.zeros_like(bh), np.zeros_like(by)
47         dnext = np.zeros_like(hs[-1])
48         for t in reversed(xrange(len(inputs))):
49             dy = np.copy(ps_t)
50             dy[targets[t]] -= 1 # backprop into y
51             dby += dy
52             dhy += dy
53             dh = np.dot(Why.T, dy) + dnext # backprop into h
54             dhraw = (1 - hs[-1]**2) * dh # backprop through tanh nonlinearity
55             dbh += np.dot(dhraw, hs[-1])
56             dwhh += np.dot(dhraw, xs[t].T)
57             dwhh += np.dot(dhraw, hs[-1].T)
58             dnext = np.dot(Whh, dhraw)
59         for dparam in [dwhh, dhy, dby, dbh]:
60             np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61     return loss, dxh, ddmh, dhy, dbh, dby, hs[-1]
62
63 def sample(hseed, seed_ix, n):
64     """
65     sample a sequence of integers from the model
66     h is memory state, seed_ix is seed letter for first time step
67     """
68     x = np.zeros((vocab_size, 1))
69     x[seed_ix] = 1
70     sample_ix = []
71     for t in xrange(n):
72         h = np.tanh(np.dot(Wkh, x) + np.dot(bh, h) + bh)
73         y = np.dot(Why, h) + by
74         p = np.exp(y) / np.sum(np.exp(y))
75         ix = np.random.choice(range(vocab_size), p=p.ravel())
76         x = np.zeros((vocab_size, 1))
77         x[ix] = 1
78         sample_ix.append(ix)
79     return sample_ix
80
81 n, p, q, a
82 mwhh, mwhh, mwhy = np.zeros_like(Whh), np.zeros_like(Whh), np.zeros_like(Why)
83 mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
84 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
85 while True:
86     """
87     # prepare inputs (we're sweeping from left to right in steps seq_length long)
88     if p+seq_length+1 >= len(data) or n == 0:
89         hprev = np.zeros((hidden_size,1)) # reset RNN memory
90         p = 0 # go from start of data
91     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
92     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
93
94     # sample from the model now and then
95     if n % 100 == 0:
96         sample_ix = sample(hprev, inputs[0], 200)
97         txt = ''.join(ix_to_char[ix] for ix in sample_ix)
98         print '----\n%s ----' % (txt, )
99
100    # forward seq_length characters through the net and fetch gradient
101    loss, dWkh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
102    smooth_loss = smooth_loss * 0.999 + loss * 0.001
103
104    # perform parameter update with Adagrad
105    for param, dparam, mem in zip([Wkh, Whh, Why, bh, by],
106                                 [dWkh, dWhh, dWhy, dbh, dby],
107                                 [mWkh, mWhh, mWhy, mbh, mby]):
108        mem += dparam * dparam
109        param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter
```

Main loop



min-char-rnn.py gist

```

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 data = open('input.txt', 'r').read() # should be simple plain text file
9 chars = list(set(data))
10 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
11 char_to_ix = { ch:i for i,ch in enumerate(chars) }
12 ix_to_char = { i:ch for i,ch in enumerate(chars) }
13
14 # hyperparameters
15 hidden_size = 100 # size of hidden layer of neurons
16 seq_length = 25 # number of steps to unroll the RNN for
17 learning_rate = 1e-1
18
19 # model parameters
20 Wkh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
21 Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
22 Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
23 bh = np.zeros((hidden_size, 1)) # hidden bias
24 by = np.zeros((vocab_size, 1)) # output bias
25
26 def lossFun(inputs, targets, hprev):
27     """
28     inputs, targets are both list of integers.
29     hprev is hxi array of initial hidden state
30     returns the loss, gradients on model parameters, and last hidden state
31     """
32     xs, hs, ys, ps = [], [], [], []
33     hs[-1] = np.copy(hprev)
34     loss = 0
35     for t in xrange(len(inputs)):
36         x = inputs[t]
37         # forward pass
38         for t in xrange(len(inputs)):
39             xs[t] = np.zeros((vocab_size,1)) # encode in 1-of-k representation
40             xs[t][x] = 1
41             wh = np.dot(Wkh, xs[t]) + np.dot(bh, hs[-1]) + bh # hidden state
42             wh += np.dot(Why, hs[t]) + b # unnormalized log probabilities for next chars
43             ps[t] = np.exp(ps[t]) / np.sum(np.exp(ps[t])) # softmax (cross-entropy loss)
44             loss += -np.log(ps[t][targets[t], 0])
45             # backward pass: compute gradients going backwards
46             dprev_dh = np.zeros_like(whh) # gradients for previous hidden state
47             dbh = np.zeros_like(bh) # gradients for hidden bias
48             dby = np.zeros_like(by) # gradients for output bias
49             dhnext = np.zeros_like(hs[0])
50             for t in reversed(xrange(len(inputs))):
51                 dy = np.copy(ps[t])
52                 dy[targets[t]] -= 1 # backprop into y
53                 dby += dy
54                 dyh = np.dot(Why.T, dy) + dhnext # backprop into h
55                 dhraw = (1 - hs[t] * hs[t]) * dyh # tanh nonlinearity
56                 dh = dhraw * dprev_dh # backprop through tanh nonlinearity
57                 dwh = np.dot(xs[t].T, dh) # gradients for weight h
58                 dbh += np.dot(dhraw, hs[t-1].T) # gradients for hidden bias
59                 dhnext = np.dot(Whh, hs[t-1].T) # gradients for next hidden state
60             for dparam, param in zip(dparam, [dwh, dbh, dhy, dbh, dby]):
61                 np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
62             return loss, dxh, ddmh, ddmh, dbh, dby, hs[-1]
63
64 def sample(hseed, ix):
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     for t in xrange(0):
72         h = np.tanh(np.dot(Wkh, x) + np.dot(bh, h) + bh)
73         y = np.dot(Why, h) + by
74         p = np.exp(y) / np.sum(np.exp(y))
75         ix = np.random.choice(range(vocab_size), p=p.ravel())
76         x = np.zeros((vocab_size, 1))
77         x[ix] = 1
78         if ix in appendix:
79             return ix
80
81     n, p, q, a
82     mWxh, mWhh, mWhy = np.zeros_like(Wkh), np.zeros_like(Whh), np.zeros_like(Why)
83     mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
84     smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
85     while True:
86         """
87         # prepare inputs (we're sweeping from left to right in steps seq_length long)
88         if p+seq_length+1 >= len(data) or n == 0:
89             hprev = np.zeros((hidden_size,1)) # reset RNN memory
90             p = 0 # go from start of data
91             inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
92             targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
93
94             # sample from the model now and then
95             if n % 100 == 0:
96                 sample_ix = sample(hprev, inputs[0], 200)
97                 txt = ''.join(ix_to_char[ix] for ix in sample_ix)
98                 print '----\n%s ----' % (txt, )
99
100            # forward seq_length characters through the net and fetch gradient
101            loss, dWxh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
102            smooth_loss = smooth_loss * 0.999 + loss * 0.001
103
104            # perform parameter update with Adagrad
105            for param, dparam, mem in zip([Wkh, Whh, Why, bh, by],
106                                         [dWxh, dWhh, dWhy, dbh, dby],
107                                         [mWxh, mWhh, mWhy, mbh, mby]):
108                mem += dparam * dparam
109                param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
110
111            p += seq_length # move data pointer
112            n += 1 # iteration counter
113
114        # perform parameter update with Adagrad
115        for param, dparam, mem in zip([Wkh, Whh, Why, bh, by],
116                                     [dWxh, dWhh, dWhy, dbh, dby],
117                                     [mWxh, mWhh, mWhy, mbh, mby]):
118            mem += dparam * dparam
119            param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
120
121        p += seq_length # move data pointer
122        n += 1 # iteration counter

```

Main loop



min-char-rnn.py gist

```

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 data = open('input.txt', 'r').read() # should be simple plain text file
9 chars = list(set(data))
10 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
11 char_to_ix = { ch:i for i,ch in enumerate(chars) }
12 ix_to_char = { i:ch for i,ch in enumerate(chars) }
13
14 # hyperparameters
15 hidden_size = 100 # size of hidden layer of neurons
16 seq_length = 25 # number of steps to unroll the RNN for
17 learning_rate = 1e-1
18
19 # model parameters
20 Wkh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
21 Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
22 Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
23 bh = np.zeros((hidden_size, 1)) # hidden bias
24 by = np.zeros((vocab_size, 1)) # output bias
25
26 def lossFun(inputs, targets, hprev):
27     """
28     inputs, targets are both list of integers.
29     hprev is hxi array of initial hidden state
30     returns the loss, gradients on model parameters, and last hidden state
31     """
32     xs, hs, ys, ps = [], [], [], []
33     hs[-1] = np.copy(hprev)
34     loss = 0
35     for t in reversed(xrange(len(inputs))):
36         xs.append(inputs[t]) # encode in 1-of-k representation
37         x = np.zeros((vocab_size, 1))
38         x[inputs[t]] = 1
39         ws = np.dot(Wkh, x) + np.dot(bh, hs[-1]) + bh # hidden state
40         wh = np.dot(Why, hs[-1]) + b # unnormalized log probabilities for next chars
41         ps[t] = np.exp(ws) / np.sum(np.exp(ws)) # softmax (cross-entropy loss)
42         loss += -np.log(ps[t][targets[t], 0])
43         # backward pass: compute gradients going backwards
44         dws = np.dot(xs[t].T, ps[t]) # gradients for weights
45         dbh = np.zeros_like(bh) # gradients for hidden bias
46         dby = np.zeros_like(by) # gradients for output bias
47         dhnext = np.zeros_like(hs[-1])
48         for t in reversed(xrange(len(inputs))):
49             dy = np.copy(ps[t])
50             dy[targets[t]] -= 1 # backprop into y
51             dby += dy
52             dhy += np.dot(Why.T, dy) # backprop into h
53             dhraw = (1 - hs[t] * hs[t]) * dhy # backprop through tanh nonlinearity
54             dh = dhraw * dhnext
55             dws += np.dot(xs[t].T, dh)
56             dwh = np.dot(dhraw, xs[t].T)
57             dhy += np.dot(dhraw, hs[t-1].T)
58             dhnext = np.dot(Why.T, dhraw)
59         for dparam in [dws, dwh, dhy, dbh, dby]:
60             np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61     return loss, dxh, dmh, dhy, dbh, dby, hs[-1]
62
63 def sample(hseed, ix):
64     """
65     sample a sequence of integers from the model
66     h is memory state, seed_ix is seed letter for first time step
67     """
68     x = np.zeros((vocab_size, 1))
69     x[seed_ix] = 1
70     h = hseed
71     for t in xrange(seq_length):
72         h = np.tanh(np.dot(Wkh, x) + np.dot(bh, h) + bh)
73         y = np.dot(Why, h) + by
74         p = np.exp(y) / np.sum(np.exp(y))
75         ix = np.random.choice(range(vocab_size), p=p.ravel())
76         x = np.zeros((vocab_size, 1))
77         x[ix] = 1
78         hseed = h
79     return ix
80
81 n, p = 0, 0
82 mWxh, mWhh, mWhy = np.zeros_like(Wkh), np.zeros_like(Whh), np.zeros_like(Why)
83 mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
84 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
85 while True:
86     # prepare inputs (we're sweeping from left to right in steps seq_length long)
87     if p+seq_length+1 >= len(data) or n == 0:
88         hprev = np.zeros((hidden_size, 1)) # reset RNN memory
89         p = 0 # go from start of data
90     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
91     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
92
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 '----\n%s ----' % (txt, )
98
99     # forward seq_length characters through the net and fetch gradient
100    loss, dWxh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
101    smooth_loss = smooth_loss * 0.999 + loss * 0.001
102
103    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
104
105    # perform parameter update with Adagrad
106    for param, dparam, mem in zip([Wkh, Whh, Why, bh, by],
107                                 [dWxh, dWhh, dWhy, dbh, dby],
108                                 [mWxh, mWhh, mWhy, mbh, mby]):
109        mem += dparam * dparam
110        param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
111
112    p += seq_length # move data pointer
113    n += 1 # iteration counter

```

Main loop



min-char-rnn.py gist

```

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 data = open('input.txt', 'r').read() # should be simple plain text file
9 chars = list(set(data))
10 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
11 char_to_ix = { ch:i for i,ch in enumerate(chars) }
12 ix_to_char = { i:ch for i,ch in enumerate(chars) }
13
14 # hyperparameters
15 hidden_size = 100 # size of hidden layer of neurons
16 seq_length = 25 # number of steps to unroll the RNN for
17 learning_rate = 1e-1
18
19 # model parameters
20 Wkh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
21 Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
22 Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
23 bh = np.zeros((hidden_size, 1)) # hidden bias
24 by = np.zeros((vocab_size, 1)) # output bias
25
26 def lossFun(inputs, targets, hprev):
27     """
28     inputs, targets are both list of integers.
29     hprev is hxi array of initial hidden state
30     returns the loss, gradients on model parameters, and last hidden state
31     """
32     xs, hs, ys, ps = [], [], [], []
33     hs[-1] = np.copy(hprev)
34     loss = 0
35     for t in xrange(len(inputs)):
36         x = inputs[t]
37         # forward pass
38         for t in xrange(len(inputs)):
39             xs[t] = np.zeros((vocab_size,1)) # encode in 1-of-k representation
40             xs[t][inputs[t]] = 1
41             wh = np.dot(Wkh, xs[t]) + np.dot(bh, hs[-1]) + bh # hidden state
42             why = np.dot(Why, hs[-1]) + b # unnormalized log probabilities for next chars
43             ps[t] = np.exp(why) / np.sum(np.exp(why)) # softmax (cross-entropy loss)
44             loss += -np.log(ps[t][targets[t]]) # softmax (cross-entropy loss)
45             # backward pass: compute gradients going backwards
46             dwh, dbh = -np.gradient(why, [Wkh, bh], np.zeros_like(why), np.zeros_like(bh))
47             dbh = np.zeros_like(bh)
48             dnext = np.zeros_like(hs[-1])
49             for t in reversed(xrange(len(inputs))):
50                 dy = np.copy(ps[t])
51                 dy[targets[t]] -= 1 # backprop into y
52                 dby = np.dot(dy, hs[t])
53                 dby += dy
54                 dh = np.dot(Why.T, dy) + dnext # backprop into h
55                 dhraw = (1 - hs[t] * hs[t]) * dh # backprop through tanh nonlinearity
56                 dwh = -np.dot(xs[t].T, dhraw)
57                 dwh += np.dot(dhraw, xs[t].T)
58                 dwh += np.dot(Whh, hs[-1].T)
59                 dnext = np.dot(Whh, dhraw)
60             for dparam, param in zip(dparam, [dwh, dbh, dhy, dbh, dby]):
61                 np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
62             return loss, dwh, dbh, dhy, dbh, dby, hs[-1]
63
64 def sample(h, n, seed_ix, h):
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     sample_ix = seed_ix
71     for t in xrange(n):
72         h = np.tanh(np.dot(Wkh, x) + np.dot(bh, h) + bh)
73         y = np.dot(Why, h) + by
74         p = np.exp(y) / np.sum(np.exp(y))
75         ix = np.random.choice(range(vocab_size), p=p.ravel())
76         x = np.zeros((vocab_size, 1))
77         x[ix] = 1
78         sample_ix = ix
79     return sample_ix
80
81 n, p = 0, 0
82 mWxh, mWhh, mWhy = np.zeros_like(Wkh), np.zeros_like(Whh), np.zeros_like(Why)
83 mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
84 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
85 while True:
86     """
87     # prepare inputs (we're sweeping from left to right in steps seq_length long)
88     if p+seq_length+1 >= len(data) or n == 0:
89         hprev = np.zeros((hidden_size,1)) # reset RNN memory
90         p = 0 # go from start of data
91     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
92     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
93
94     # sample from the model now and then
95     if n % 100 == 0:
96         sample_ix = sample(hprev, inputs[0], 200)
97         txt = ''.join(ix_to_char[ix] for ix in sample_ix)
98         print '----\n%s ----' % (txt, )
99
100    # forward seq_length characters through the net and fetch gradient
101    loss, dWxh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
102    smooth_loss = smooth_loss * 0.999 + loss * 0.001
103
104    # perform parameter update with Adagrad
105    for param, dparam, mem in zip([Wkh, Whh, Why, bh, by],
106                                 [dWxh, dWhh, dWhy, dbh, dby],
107                                 [mWxh, mWhh, mWhy, mbh, mby]):
108        mem += dparam * dparam
109        param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter

```

Main loop



min-char-rnn.py gist

```

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 data = open('input.txt', 'r').read() # should be simple plain text file
9 chars = list(set(data))
10 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
11 char_to_ix = { ch:i for i,ch in enumerate(chars) }
12 ix_to_char = { i:ch for i,ch in enumerate(chars) }
13
14 # hyperparameters
15 hidden_size = 100 # size of hidden layer of neurons
16 seq_length = 25 # number of steps to unroll the RNN for
17 learning_rate = 1e-1
18
19 # model parameters
20 Wkh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
21 Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
22 Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
23 bh = np.zeros((hidden_size, 1)) # hidden bias
24 by = np.zeros((vocab_size, 1)) # output bias
25
26 def lossFun(inputs, targets, hprev):
27     """
28     inputs, targets are both list of integers.
29     hprev is hxi array of initial hidden state
30     returns the loss, gradients on model parameters, and last hidden state
31     """
32     xs, hs, ys, ps = [], [], [], []
33     hs[-1] = np.copy(hprev)
34     loss = 0
35     for t in xrange(len(inputs)):
36         x = inputs[t]
37         # forward pass
38         for t in xrange(len(inputs)):
39             xs[t] = np.zeros((vocab_size,1)) # encode in 1-of-k representation
40             xs[t][x] = 1.0
41             wh = np.dot(Wkh, xs[t]) + bh # hidden state
42             why = np.dot(Why, wh) + by # unnormalized log probabilities for next chars
43             ps[t] = np.exp(why) / np.sum(np.exp(why)) # softmax (cross-entropy loss)
44             loss += -np.log(ps[t][targets[t]]) # softmax (cross-entropy loss)
45             # backward pass: compute gradients going backwards
46             dwh, dbh = np.zeros_like(wh), np.zeros_like(why)
47             dby = np.zeros_like(by)
48             dhnext = np.zeros_like(hs[0])
49             for t in reversed(xrange(len(inputs))):
50                 dy = np.copy(ps[t])
51                 dy[targets[t]] -= 1 # backprop into y
52                 dby += dy
53                 dhy = np.dot(Why.T, dy) + dhnext # backprop into h
54                 dhraw = (1 - hs[t] * hs[t]) * dhy # backprop through tanh nonlinearity
55                 dwh += np.dot(xs[t].T, dhraw)
56                 dwh += np.dot(dhraw, xs[t].T)
57                 dwh += np.dot(Whh, hs[t-1].T)
58                 dhnext = np.dot(Whh, dhraw)
59             for dparam in [dwh, dbh, dhy, dby]:
60                 np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61             return loss, dwh, dbh, dhy, dby, hs[-1]
62
63 def sample(hseed, seed_ix, n):
64     """
65     sample a sequence of integers from the model
66     h is memory state, seed_ix is seed letter for first time step
67     """
68     x = np.zeros((vocab_size, 1))
69     x[seed_ix] = 1
70     sample_ix = []
71     for t in xrange(n):
72         h = np.tanh(np.dot(Wkh, x) + np.dot(Whh, h) + bh)
73         y = np.dot(Why, h) + by
74         p = np.exp(y) / np.sum(np.exp(y))
75         ix = np.random.choice(range(vocab_size), p=p.ravel())
76         x = np.zeros((vocab_size, 1))
77         x[ix] = 1.0
78         sample_ix.append(ix)
79    return sample_ix
80
81 n, p, q, a
82 mwhxh, mwhh, mwhy = np.zeros_like(Wkh), np.zeros_like(Whh), np.zeros_like(Why)
83 mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
84 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
85 while True:
86     """
87     # prepare inputs (we're sweeping from left to right in steps seq_length long)
88     if p+seq_length+1 >= len(data) or n == 0:
89         hprev = np.zeros((hidden_size,1)) # reset RNN memory
90         p = 0 # go from start of data
91     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
92     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
93
94     # sample from the model now and then
95     if n % 100 == 0:
96         sample_ix = sample(hprev, inputs[0], 200)
97         txt = ''.join(ix_to_char[ix] for ix in sample_ix)
98         print '----\n%s ----' % (txt, )
99
100    # forward seq_length characters through the net and fetch gradient
101    loss, dWkh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
102    smooth_loss = smooth_loss * 0.999 + loss * 0.001
103
104    # perform parameter update with Adagrad
105    for param, dparam, mem in zip([Wkh, Whh, Why, bh, by],
106                                 [dWkh, dWhh, dWhy, dbh, dby],
107                                 [mwhxh, mwhh, mwhy, mbh, mby]):
108        mem += dparam * dparam
109        param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter

```

Main loop



```

104    # perform parameter update with Adagrad
105    for param, dparam, mem in zip([Wkh, Whh, Why, bh, by],
106                                 [dWkh, dWhh, dWhy, dbh, dby],
107                                 [mwhxh, mwhh, mwhy, mbh, mby]):
108        mem += dparam * dparam
109        param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter

```

min-char-rnn.py gist

```

2 Minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)
3 BSD license
4
5 import numpy as np
6
7 # data
8 data = open('input.txt', 'r').read() # should be simple plain text file
9 chars = list(set(data))
10 data_size, vocab_size = len(data), len(chars)
11 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
12 char_to_ix = { ch:i for i,ch in enumerate(chars) }
13 ix_to_char = { i:ch for i,ch in enumerate(chars) }
14
15 # hyperparameters
16 hidden_size = 100 # size of hidden layer of neurons
17 seq_length = 25 # number of steps to unroll the RNN for
18 learning_rate = 1e-1
19
20 # model parameters
21 Wxh = np.random.rand(hidden_size, vocab_size)*0.01 # input to hidden
22 Whh = np.random.rand(hidden_size, hidden_size)*0.01 # hidden to hidden
23 Why = np.random.rand(vocab_size, hidden_size)*0.01 # hidden to output
24 bh = np.zeros(hidden_size, 1) # hidden bias
25 by = np.zeros(vocab_size, 1) # output bias
26
27 def lossFun(inputs, targets, hprev):
28     """
29     inputs, targets are both list of integers.
30     hprev is Hx1 array of initial hidden state
31     returns the loss, gradients on model parameters, and last hidden state
32     """
33     xs, hs, ys, ps = {}, {}, {}, {}
34     hs[-1] = np.copy(hprev)
35     loss = 0
36     for t in xrange(len(inputs)):
37         # forward pass
38         for t in xrange(len(inputs)):
39             xs[t] = np.zeros((vocab_size,1)) # encode in 1-of-k representation
40             xs[t][inputs[t]] = 1
41             wh = np.dot(Wxh, xs[t]) + bh # hidden state
42             ws[t] = np.dot(why, hs[t]) + by # unnormalized log probabilities for next chars
43             ps[t] = np.exp(ws[t]) / np.sum(np.exp(ws[t])) # probabilities for next chars
44             loss += -np.log(ps[t][targets[t],0]) # softmax (cross-entropy loss)
45             # backward pass: compute gradients going backwards
46             dws[t] = np.zeros_like(wh) # gradients for w
47             dbh = np.zeros_like(bh) # gradients for b
48             dby = np.zeros_like(by) # gradients for b
49             dhnext = np.zeros_like(hs[t])
50             for t in reversed(xrange(len(inputs))):
51                 dy = np.copy(ps[t])
52                 dy[targets[t]] -= 1 # backprop into y
53                 dyb = np.dot(dy, hs[t].T)
54                 dhy = np.dot(dy, why.T, dy) + dhnext # backprop into h
55                 ddraw = (1 - hs[t] * hs[t].T) * dyb # backprop through tanh nonlinearity
56                 dwh = np.dot(ddraw, xs[t].T)
57                 dwh += np.dot(draw, hs[t-1].T)
58                 dhnext = np.dot(why.T, dhy)
59             for dparam in [dwh, ddraw, dhy, dbh, dby]:
60                 np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61             return loss, dxh, dwh, dhy, dbh, dby, hs[t-1]
62
63 def sample(h0, seed_ix, n):
64     """
65     sample a sequence of integers from the model
66     h0 is memory state, seed_ix is seed letter for first time step
67     """
68     x = np.zeros((vocab_size, 1))
69     x[seed_ix] = 1
70     hs, ys, ps = {}, {}, {}
71     for t in xrange(n):
72         h = np.tanh(np.dot(Wxh, x) + np.dot(bh, h) + bh)
73         y = np.dot(why, h) + by
74         ws[t] = np.dot(why, h) + by # unnormalized log probabilities for next chars
75         ix = np.random.choice(range(vocab_size), p=p.ravel())
76         x = np.zeros((vocab_size, 1))
77         x[ix] = 1 # append ix to x
78     return hs, ys, ps
79
80 n, p = 0, 0
81 smoth, smoth_m, smoth_b, smoth_y = np.zeros_like(wh), np.zeros_like(Why), np.zeros_like(bh), np.zeros_like(by)
82 dbh, dby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
83 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
84 while True:
85     # prepare inputs (we're sweeping from left to right in steps seq_length long)
86     if p+seq_length >= len(data) or p == 0:
87         hprev = np.zeros((hidden_size, 1)) # reset RNN memory
88         p = 0 # go from start of data
89     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
90     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
91
92     # sample from the model now and then
93     if n % 100 == 0:
94         smoth += sample(hprev, inputs[0], 200)
95         txt = '...'.join(ix_to_char[i] for i in sample_ix)
96         print '...'.join('%.3f' % n, smoth_loss, '(%t, %t)' % (txt, ))
97
98     # forward seq_length characters through the net and fetch gradient
99     loss, dxh, dwh, dhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
100    smoth_loss += smoth.loss * 0.999 + loss * 0.001
101
102    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smoth_loss) # print progress
103
104    # perform parameter update with Adagrad
105    for param, dparam, mem in zip([xh, wh, why, bh, by],
106                                 [dxh, dwh, dhy, dbh, dby],
107                                 [smoth, smoth_m, smoth_b, smoth_y]):
108        mem += dparam * dparam
109        param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adgrad update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter

```

Loss function

- forward pass (compute loss)
- backward pass (compute param gradient)

```

27 def lossFun(inputs, targets, hprev):
28     """
29     inputs,targets are both list of integers.
30     hprev is Hx1 array of initial hidden state
31     returns the loss, gradients on model parameters, and last hidden state
32     """
33     xs, hs, ys, ps = {}, {}, {}, {}
34     hs[-1] = np.copy(hprev)
35     loss = 0
36
37     # forward pass
38     for t in xrange(len(inputs)):
39         xs[t] = np.zeros((vocab_size,1)) # encode in 1-of-k representation
40         xs[t][inputs[t]] = 1
41         wh = np.dot(Wxh, xs[t]) + bh # hidden state
42         ws[t] = np.dot(why, hs[t]) + by # unnormalized log probabilities for next chars
43         ps[t] = np.exp(ws[t]) / np.sum(np.exp(ws[t])) # probabilities for next chars
44         loss += -np.log(ps[t][targets[t],0]) # softmax (cross-entropy loss)
45
46     # backward pass: compute gradients going backwards
47     dWxh, dWhh, dWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
48     dbh, dby = np.zeros_like(bh), np.zeros_like(by)
49     dhnext = np.zeros_like(hs[0])
50
51     for t in reversed(xrange(len(inputs))):
52         dy = np.copy(ps[t])
53         dy[targets[t]] -= 1 # backprop into y
54         dWhy += np.dot(dy, hs[t].T)
55         dyb = np.dot(dy, why.T)
56         dhy = np.dot(dy, why.T, dy) + dhnext # backprop into h
57         ddraw = (1 - hs[t] * hs[t].T) * dyb # backprop through tanh nonlinearity
58         dWxh += np.dot(ddraw, xs[t].T)
59         dbh += ddraw
60         dWhy += np.dot(ddraw, hs[t-1].T)
61         dhnext = np.dot(why.T, dhy)
62
63     for dparam in [dWxh, dWhh, dWhy, dbh, dby]:
64         np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
65
66     return loss, dWxh, dWhh, dWhy, dbh, dby, hs[len(inputs)-1]

```

min-char-rnn.py gist

```

2 Minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)
3 BSD license
4
5 import numpy as np
6
7 # data
8 data = open('input.txt', 'r').read() # should be simple plain text file
9 chars = list(set(data))
10 vocab_size = len(chars)
11 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
12 char_to_ix = { ch:i for i,ch in enumerate(chars) }
13 ix_to_char = { i:ch for i,ch in enumerate(chars) }
14
15 # hyperparameters
16 hidden_size = 100 # size of hidden layer of neurons
17 seq_length = 25 # number of steps to unroll the RNN for
18 learning_rate = 1e-1
19
20 # model parameters
21 Wxh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
22 Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
23 Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
24 bh = np.zeros(hidden_size, 1) # hidden bias
25 by = np.zeros(vocab_size, 1) # output bias
26
27 def lossFun(inputs, targets, hprev):
28     """
29     inputs, targets are both list of integers.
30     hprev is Hx1 array of initial hidden state
31     returns the loss, gradients on model parameters, and last hidden state
32     """
33     xs, hs, ys, ps = {}, {}, {}, {}
34     hs[-1] = np.copy(hprev)
35     loss = 0
36     # forward pass
37     for t in xrange(len(inputs)):
38         xs[t] = np.zeros((vocab_size,1)) # encode in 1-of-k representation
39         xs[t][inputs[t]] = 1
40         hs[t] = np.tanh(np.dot(Wxh, xs[t]) + np.dot(Whh, hs[t-1]) + bh) # hidden state
41         ys[t] = np.dot(Why, hs[t]) + by # unnormalized log probabilities for next chars
42         ps[t] = np.exp(ys[t]) / np.sum(np.exp(ys[t])) # probabilities for next chars
43         loss += -np.log(ps[t][targets[t],0]) # softmax (cross-entropy loss)
44
45     # backward pass: compute gradients going backwards
46     dws, dwh, dbh, dy, dby = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(bh),
47     dby = np.zeros_like(Why), np.zeros_like(by)
48     dhnext = np.zeros_like(hs[0])
49     for t in reversed(xrange(len(inputs))):
50         dy[t] = np.copy(ps[t]) - 1 # backprop into y
51         dby += dy[t]
52         dy += dby
53         dh = np.dot(Why.T, dy) + dhnext # backprop through Why
54         dhraw = (1 - hs[t] * hs[t]) * dh # backprop through tanh nonlinearity
55         dwh += np.dot(dhraw, hs[t-1].T)
56         dbh += np.dot(dhraw, hs[t-1])
57         dhnext = np.dot(Whh.T, dhraw)
58     for dparam, dhx, dhh, dhy, dbx, dbh, dby:
59         np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients!
60     return loss, hs[-1], dws, dwh, dbh, dy, dby, hs[1:][inputs[1:]]
61
62 def sample(h0, seed_ix):
63     """
64     sample a sequence of integers from the model
65     h0 is memory state, seed_ix is seed letter for first time step
66     """
67     x = np.zeros((vocab_size, 1))
68     x[seed_ix] = 1
69     ixes = []
70     for t in xrange(n):
71         h = np.tanh(np.dot(Wxh, x) + np.dot(Whh, h0) + bh)
72         y = np.dot(Why, h) + by
73         p = np.exp(y) / np.sum(np.exp(y))
74         ix = np.random.choice(range(vocab_size), p=p.ravel())
75         x = np.zeros((vocab_size, 1))
76         x[ix] = 1 # append ix to x
77     return ixes
78
79 n, p = 0, 0
80
81 Wxh, Whh, Why, bh, dbh, dy, dby = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why),
82 dbh, dby = np.zeros_like(bh), np.zeros_like(dy) # memory variables for Adagrad
83 smooth_loss = np.inf # average loss
84 seq_length = len(data) # seq_length = loss at iteration 0
85 while True:
86     # prepare inputs (we're sweeping from left to right in steps seq_length long)
87     if p+seq_length+1 >= len(data) or p == 0:
88         h0 = np.zeros(hidden_size) # reset RNN memory
89         p = 0 # go from start of data
90     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
91     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
92
93     # sample from the model now and then
94     if n % 100 == 0:
95         s = sample(h0, inputs[0])
96         txt = ''.join(ix_to_char[ix] for ix in s)
97         print '-----%s-----' % (txt, )
98
99     # forward seq_length characters through the net and fetch gradient
100    loss, dxh, dhh, dbh, dy, dby, hprev = lossFun(inputs, targets, h0)
101    smooth_loss = smooth_loss * 0.999 + loss * 0.001
102    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
103
104    # perform parameter update with Adagrad
105    for param, dparam, mem in zip([Wxh, Whh, Why, bh],
106                                 [dxh, dhh, dbh, dy],
107                                 [dbh, dby, dbh, dby]):
108        mem += dparam * dparam
109        param -= learning_rate * dparam / np.sqrt(mem + 1e-8) # adgrad update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter

```

```

27 def lossFun(inputs, targets, hprev):
28     """
29     inputs,targets are both list of integers.
30     hprev is Hx1 array of initial hidden state
31     returns the loss, gradients on model parameters, and last hidden state
32     """
33     xs, hs, ys, ps = {}, {}, {}, {}
34     hs[-1] = np.copy(hprev)
35     loss = 0
36     # forward pass
37     for t in xrange(len(inputs)):
38         xs[t] = np.zeros((vocab_size,1)) # encode in 1-of-k representation
39         xs[t][inputs[t]] = 1
40         hs[t] = np.tanh(np.dot(Wxh, xs[t]) + np.dot(Whh, hs[t-1]) + bh) # hidden state
41         ys[t] = np.dot(Why, hs[t]) + by # unnormalized log probabilities for next chars
42         ps[t] = np.exp(ys[t]) / np.sum(np.exp(ys[t])) # probabilities for next chars
43         loss += -np.log(ps[t][targets[t],0]) # softmax (cross-entropy loss)

```

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

$$y_t = W_{hy}h_t$$

Softmax classifier

min-char-rnn.py gist

```

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

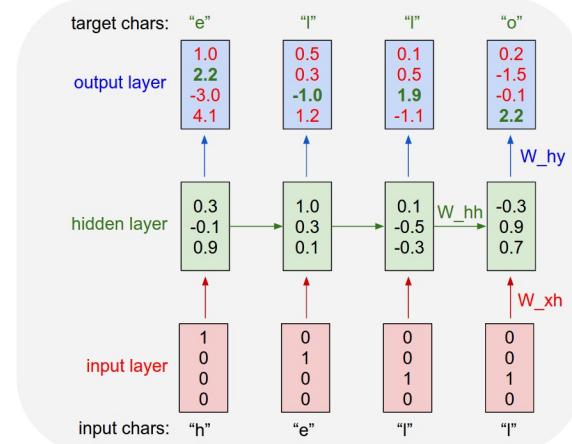
```

```

44     # backward pass: compute gradients going backwards
45     dWxh, dWhh, dWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
46     dbh, dby = np.zeros_like(bh), np.zeros_like(by)
47     dhnnext = np.zeros_like(hs[0])
48     for t in reversed(xrange(len(inputs))):
49         dy = np.copy(ps[t])
50         dy[targets[t]] -= 1 # backprop into y
51         dWhy += np.dot(dy, hs[t].T)
52         dhy += dy
53         dh = np.dot(Why.T, dy) + dhnnext # backprop into h
54         ddraw = (1 - hs[t] * hs[t]) * dh # backprop through tanh nonlinearity
55         dbh += ddraw
56         dWxh += np.dot(ddraw, xs[t].T)
57         dWhh += np.dot(ddraw, hs[t-1].T)
58         dhnnext = np.dot(Whh.T, ddraw)
59         for dparam in [dWxh, dWhh, dWhy, dbh, dby]:
60             np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61     return loss, dWxh, dWhh, dWhy, dbh, dby, hs[-1]
62
63 def sample(hseed, ix, n):
64     """
65     sample a sequence of integers from the model
66     h is in memory state, seed_ix is seed letter for first time step
67     """
68     x = np.zeros((vocab_size, 1))
69     x[seed_ix] = 1
70     ixes = []
71     for t in xrange(n):
72         h = np.tanh(np.dot(Wxh, x) + np.dot(Whh, h) + Whh)
73         y = np.dot(Why, h) + Why
74         p = np.exp(y) / np.sum(np.exp(y))
75         ix = np.random.choice(range(vocab_size), p=p.ravel())
76         x[1] = np.zeros((vocab_size, 1))
77         x[1][ix] = 1
78         ixes.append(ix)
79     return ixes
80
81 n, p = 0, 0
82 smoth, smoth_m, smoth_b = np.zeros_like(Wxh), np.zeros_like(Wxh), np.zeros_like(Wxh)
83 dbh, dby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
84 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
85 while True:
86     # prepare inputs (we're sweeping from left to right in steps seq_length long)
87     if p+seq_length >= len(data) or n == 0:
88         hprev = np.zeros_like(h) # reset RNN memory
89         p = 0 # go from start of data
90     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
91     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
92
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 '-----%s-----%s' % (txt, )
98
99     # forward seq_length characters through the net and fetch gradient
100    loss, dWxh, dWhh, dWhy, dbh, dhy, dhh, dWhh, dbyh, dbyw = lossFun(inputs, targets, hprev)
101    smooth_loss = smooth_loss * 0.999 + loss * 0.001
102    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
103
104    # perform parameter update with Adagrad
105    for param, dparam, mem in zip([Wxh, bh, Why, by],
106                                 [dWxh, dbh, dWhy, dby],
107                                 [smoth_m, smoth_b, smoth_m, smoth_b, smoth_m, smoth_b, smoth_m, smoth_b, smoth_m]):
108        mem += dparam * dparam
109        param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adgrad update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter

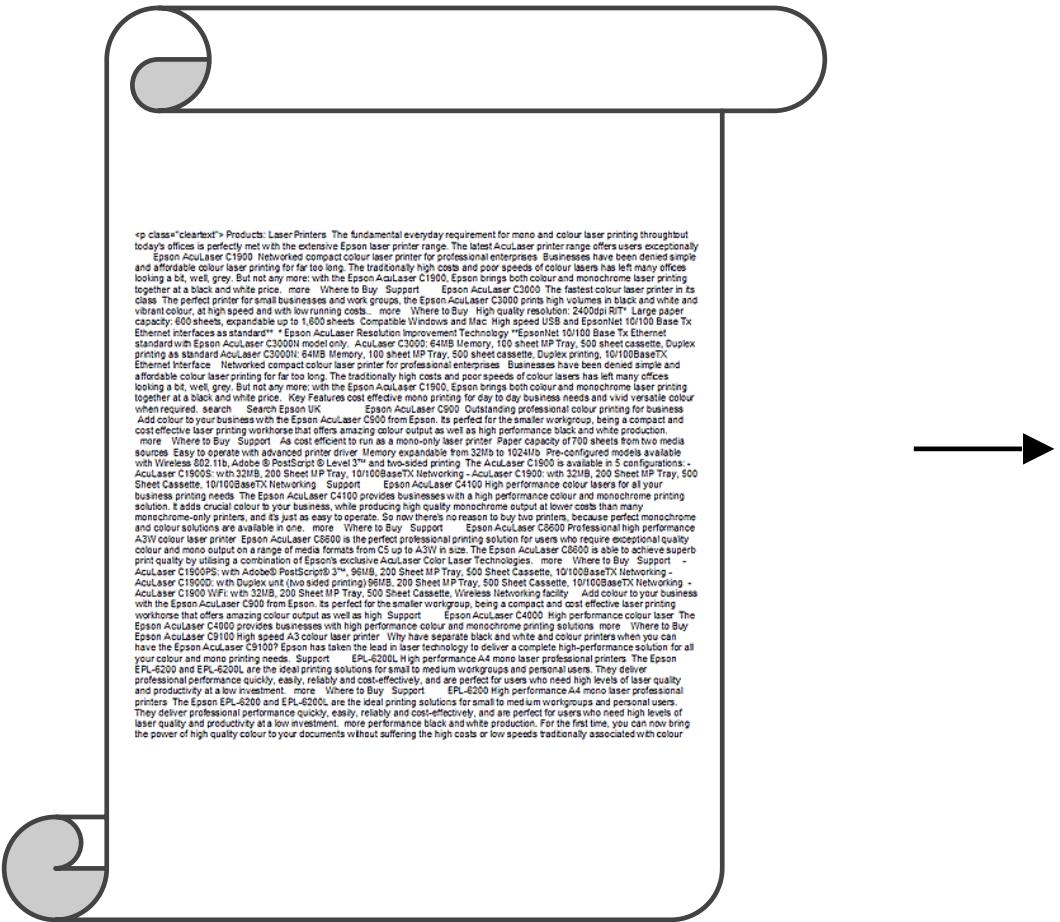
```

recall:



min-char-rnn.py gist

```
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 data = open('input.txt', 'r').read() # should be simple plain text file
9 chars = list(set(data))
10 data_size, vocab_size = len(data), len(chars)
11 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
12 char_to_ix = { ch:i for i,ch in enumerate(chars) }
13 ix_to_char = { i:ch for i,ch in enumerate(chars) }
14
15 # hyperparameters
16 hidden_size = 100 # size of hidden layer of neurons
17 seq_length = 25 # number of steps to unroll the RNN for
18 learning_rate = 1e-1
19
20 # model parameters
21 Whh = np.random.rand(hidden_size, vocab_size)*0.01 # input to hidden
22 Whx = np.random.rand(hidden_size, hidden_size)*0.01 # hidden to hidden
23 Why = np.random.rand(vocab_size, hidden_size)*0.01 # hidden to output
24 bh = np.zeros((hidden_size, 1)) # hidden bias
25 by = np.zeros((vocab_size, 1)) # output bias
26
27 def lossFun(inputs, targets, hprev):
28     """
29     inputs, targets are both list of integers.
30     hprev is hxt array of initial hidden state
31     returns the loss, gradients on model parameters, and last hidden state
32     """
33     xs, hs, ys, ps = [], (), (), {}
34     hprev[0] = np.copy(hprev)
35     loss = 0
36     for t in xrange(len(inputs)):
37         for t in xrange(len(inputs)):
38             xs[t] = np.zeros((vocab_size,1)) # encode in 1-of-k representation
39             xs[t][inputs[t]] = 1
40             hprev[0] = np.dot(Whh, hprev[0]) + np.dot(Whx, xs[t]) + bh # hidden state
41             ys[t] = np.dot(Why, hprev[0]) + by # unnormalized log probabilities for next chars
42             ps[t] = np.exp(ys[t]) / np.sum(np.exp(ys[t])) # probabilities for next chars
43             loss += -np.log(ps[t][targets[t], 0]) # softmax (cross-entropy loss)
44             # backward pass: compute gradients going backwards
45             dprev_dh = np.zeros_like(hprev[0])
46             dprev_dy = np.zeros_like(why)
47             dprev_bh = np.zeros_like(bh)
48             dprev_dx = np.zeros_like(xs[t])
49             for t in reversed(xrange(len(inputs))):
50                 dy = np.copy(ps[t])
51                 dy[targets[t]] -= 1 # backprop into y
52                 ddy = -dy
53                 dby += ddy
54                 dh = np.dot(Why.T, dy) + dprev_dh # backprop into h
55                 dprev_dhraw = (1 - hprev[0]**2) * dh # backprop through tanh nonlinearity
56                 dprev_dx = np.dot(Whx.T, dprev_dhraw)
57                 dprev_dx += np.dot(dprev_dhraw, xs[t].T)
58                 dprev_dx = np.dot(dprev_dx, hprev[0].T)
59                 dprev_dx = np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
60             return loss, dprev_dx, dprev_dy, dprev_bh, dby, h[1:len(inputs)-1]
61
62 def sample(h, seed_ix, n):
63     """
64     sample a sequence of integers from the model
65     h is memory state, seed_ix is seed letter for first time step
66     """
67
68     x = np.zeros((vocab_size, 1))
69     x[seed_ix] = 1
70     ixes = []
71     for t in xrange(n):
72         h = np.tanh(np.dot(Whh, x) + np.dot(Whx, h) + bh)
73         y = np.dot(Why, h) + by
74         p = np.exp(y) / np.sum(np.exp(y))
75         ix = np.random.choice(range(vocab_size), p=p.ravel())
76         x[0] = 0
77         ixes.append(ix)
78
79     return ixes
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
```



Sonnet 116 – Let me not ...

by William Shakespeare

Let me not to the marriage of true minds
Admit impediments. Love is not love
Which alters when it alteration finds,
Or bends with the remover to remove:
O no! it is an ever-fixed mark
That looks on tempests and is never shaken;
It is the star to every wandering bark,
Whose worth's unknown, although his height be taken.
Love's not Time's fool, though rosy lips and cheeks
Within his bending sickle's compass come:
Love alters not with his brief hours and weeks,
But bears it out even to the edge of doom.
If this be error and upon me proved,
I never writ, nor no man ever loved.

at first:

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

↓ train more

"Tmont thithey" fomesscerliund
Keushey. Thom here
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwyl 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 ofter.

↓ 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 nues 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 aps, 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.

open source textbook on algebraic geometry

The Screenshot shows the homepage of The Stacks Project. The main content area is titled "Browse chapters" and lists chapters under the "Preliminaries" part. Each chapter entry includes links for "online", "TeX source", and "view pdf". To the right, there is a sidebar titled "Parts" which lists numbered sections: 1. Preliminaries, 2. Schemes, 3. Topics in Scheme Theory, 4. Algebraic Spaces, 5. Topics in Geometry, 6. Deformation Theory, 7. Algebraic Stacks, and 8. Miscellany. Below the parts is a section titled "Statistics" which provides project metrics.

Part	Chapter	online	TeX source	view pdf
Preliminaries	1. Introduction	online	tex	pdf
	2. Conventions	online	tex	pdf
	3. Set Theory	online	tex	pdf
	4. Categories	online	tex	pdf
	5. Topology	online	tex	pdf
	6. Sheaves on Spaces	online	tex	pdf
	7. Sites and Sheaves	online	tex	pdf
	8. Stacks	online	tex	pdf
	9. Fields	online	tex	pdf
	10. Commutative Algebra	online	tex	pdf

Parts

- [Preliminaries](#)
- [Schemes](#)
- [Topics in Scheme Theory](#)
- [Algebraic Spaces](#)
- [Topics in Geometry](#)
- [Deformation Theory](#)
- [Algebraic Stacks](#)
- [Miscellany](#)

Statistics

The Stacks project now consists of

- 455910 lines of code
- 14221 tags (56 inactive tags)
- 2366 sections

Latex source

For $\bigoplus_{n=1,\dots,m} \mathcal{L}_{m,n} = 0$, hence we can find a closed subset \mathcal{H} in \mathcal{H} and any sets \mathcal{F} on X , U is a closed immersion of S , then $U \rightarrow T$ is a separated algebraic space.

Proof. Proof of (1). It also start we get

$$S = \text{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by $\coprod Z \times_U U \rightarrow V$. Consider the maps M along the set of points Sch_{fppf} and $U \rightarrow U$ is the fibre category of S in U in Section ?? and the fact that any U affine, see Morphisms, Lemma ???. Hence we obtain a scheme S and any open subset $W \subset U$ in $\text{Sh}(G)$ such that $\text{Spec}(R') \rightarrow S$ is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that f_i is of finite presentation over S . We claim that $\mathcal{O}_{X,x}$ is a scheme where $x, x' \in S'$ such that $\mathcal{O}_{X,x'} \rightarrow \mathcal{O}'_{X',x'}$ is separated. By Algebra, Lemma ?? we can define a map of complexes $\text{GL}_{S'}(x'/S'')$ and we win. \square

To prove study we see that $\mathcal{F}|_U$ is a covering of X' , and \mathcal{T}_i is an object of $\mathcal{F}_{X/S}$ for $i > 0$ and \mathcal{F}_p exists and let \mathcal{T}_i be a presheaf of \mathcal{O}_X -modules on \mathcal{C} as a \mathcal{F} -module. In particular $\mathcal{F} = U/\mathcal{F}$ we have to show that

$$\widetilde{M}^\bullet = \mathcal{I}^\bullet \otimes_{\text{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1} \mathcal{F}$$

is a unique morphism of algebraic stacks. Note that

$$\text{Arrows} = (\text{Sch}/S)^{\text{opp}}_{fppf}, (\text{Sch}/S)_{fppf}$$

and

$$V = \Gamma(S, \mathcal{O}) \longmapsto (U, \text{Spec}(A))$$

is an open subset of X . Thus U is affine. This is a continuous map of X is the inverse, the groupoid scheme S .

Proof. See discussion of sheaves of sets. \square

The result for prove any open covering follows from the less of Example ???. It may replace S by $X_{\text{spaces},\text{étale}}$ which gives an open subspace of X and T equal to S_{Zar} , see Descent, Lemma ???. Namely, by Lemma ?? we see that R is geometrically regular over S .

Lemma 0.1. Assume (3) and (3) by the construction in the description.

Suppose $X = \lim |X|$ (by the formal open covering X and a single map $\underline{\text{Proj}}_X(\mathcal{A}) = \text{Spec}(B)$ over U compatible with the complex

$$\text{Set}(\mathcal{A}) = \Gamma(X, \mathcal{O}_{X,\mathcal{O}_X}).$$

When in this case of to show that $\mathcal{Q} \rightarrow \mathcal{C}_{Z/X}$ is stable under the following result in the second conditions of (1), and (3). This finishes the proof. By Definition ?? (without element is when the closed subschemes are catenary. If T is surjective we may assume that T is connected with residue fields of S . Moreover there exists a closed subspace $Z \subset X$ of X where U in X' is proper (some defining as a closed subset of the uniqueness it suffices to check the fact that the following theorem

(1) f is locally of finite type. Since $S = \text{Spec}(R)$ and $Y = \text{Spec}(R)$.

Proof. This is form all sheaves of sheaves on X . But given a scheme U and a surjective étale morphism $U \rightarrow X$. Let $U \cap U = \coprod_{i=1,\dots,n} U_i$ be the scheme X over S at the schemes $X_i \rightarrow X$ and $U = \lim_i X_i$. \square

The following lemma surjective restrocomposes of this implies that $\mathcal{F}_{x_0} = \mathcal{F}_{x_0} = \mathcal{F}_{x,\dots,0}$.

Lemma 0.2. Let X be a locally Noetherian scheme over S , $E = \mathcal{F}_{X/S}$. Set $\mathcal{I} = \mathcal{J}_1 \subset \mathcal{I}'_n$. Since $\mathcal{I}^n \subset \mathcal{I}^n$ are nonzero over $i_0 \leq p$ is a subset of $\mathcal{J}_{n,0} \circ \mathcal{A}_2$ works.

Lemma 0.3. In Situation ???. Hence we may assume $q' = 0$.

Proof. We will use the property we see that p is the next functor (??). On the other hand, by Lemma ?? we see that

$$D(\mathcal{O}_{X'}) = \mathcal{O}_X(D)$$

where K is an F -algebra where δ_{n+1} is a scheme over S . \square

Proof. Omitted. \square

Lemma 0.1. Let \mathcal{C} be a set of the construction.

Let \mathcal{C} be a gerber covering. Let \mathcal{F} be a quasi-coherent sheaves of \mathcal{O} -modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

Proof. This is an algebraic space with the composition of sheaves \mathcal{F} on $X_{\text{étale}}$ we have

$$\mathcal{O}_X(\mathcal{F}) = \{\text{morph}_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where \mathcal{G} defines an isomorphism $\mathcal{F} \rightarrow \mathcal{F}$ of \mathcal{O} -modules. \square

Lemma 0.2. This is an integer \mathcal{Z} is injective.

Proof. See Spaces, Lemma ??.

Lemma 0.3. Let S be a scheme. Let X be a scheme and X is an affine open covering. Let $\mathcal{U} \subset \mathcal{X}$ be a canonical and locally of finite type. Let X be a scheme. Let X be a scheme which is equal to the formal complex.

The following to the construction of the lemma follows.

Let X be a scheme. Let X be a scheme covering. Let

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

be a morphism of algebraic spaces over S and Y .

Proof. Let X be a nonzero scheme of X . Let X be an algebraic space. Let \mathcal{F} be a quasi-coherent sheaf of \mathcal{O}_X -modules. The following are equivalent

- (1) \mathcal{F} is an algebraic space over S .
- (2) If X is an affine open covering.

Consider a common structure on X and X the functor $\mathcal{O}_X(U)$ which is locally of finite type. \square

This since $\mathcal{F} \in \mathcal{F}$ and $x \in \mathcal{G}$ the diagram

$$\begin{array}{ccccc}
 S & \xrightarrow{\quad} & & & \\
 \downarrow & & & & \\
 \xi & \xrightarrow{\quad} & \mathcal{O}_{X'} & \xleftarrow{\quad} & \\
 \text{gor}_s & & \uparrow & & \\
 & & =\alpha' & \longrightarrow & \\
 & & \downarrow & & \\
 & & =\alpha' & \longrightarrow & \alpha \\
 & & & & \\
 \text{Spec}(K_\psi) & & \text{Mor}_{\text{Sets}} & & d(\mathcal{O}_{X_{f/k}}, \mathcal{G}) \\
 & & & & \\
 & & & & X \\
 & & & & \downarrow \\
 & & & &
 \end{array}$$

is a limit. Then \mathcal{G} is a finite type and assume S is a flat and \mathcal{F} and \mathcal{G} is a finite type f_* . This is of finite type diagrams, and

- the composition of \mathcal{G} is a regular sequence,
- $\mathcal{O}_{X'}$ is a sheaf of rings.

\square

Proof. We have see that $X = \text{Spec}(R)$ and \mathcal{F} is a finite type representable by algebraic space. The property \mathcal{F} is a finite morphism of algebraic stacks. Then the cohomology of X is an open neighbourhood of U . \square

Proof. This is clear that \mathcal{G} is a finite presentation, see Lemmas ??.

A reduced above we conclude that U is an open covering of \mathcal{C} . The functor \mathcal{F} is a “field”

$$\mathcal{O}_{X,x} \rightarrow \mathcal{F}_{\bar{x}} \dashv (\mathcal{O}_{X_{\text{étale}}}) \rightarrow \mathcal{O}_{X_\ell}^{-1} \mathcal{O}_{X_\lambda}(\mathcal{O}_{X_\eta}^{\bar{v}})$$

is an isomorphism of covering of \mathcal{O}_{X_i} . If \mathcal{F} is the unique element of \mathcal{F} such that X is an isomorphism.

The property \mathcal{F} is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme \mathcal{O}_X -algebra with \mathcal{F} are opens of finite type over S . If \mathcal{F} is a scheme theoretic image points. \square

If \mathcal{F} is a finite direct sum \mathcal{O}_{X_λ} is a closed immersion, see Lemma ?? . This is a sequence of \mathcal{F} is a similar morphism.

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/nab/target-pending 6 days ago

arch Merge branch 'x86-urgent-for-linus' of git://git.kernel.org/pub/scm/l... 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/lhex2fw.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/... a day ago

init init: fix regression by supporting devices with major:minor:offset fo... a month ago

ipc

HTTPS clone URL <https://github.com/torvalds/linux>

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

Clone in Desktop Download ZIP

Code

Pulse

Graphs

```

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 << 1))
            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

```
/*
 * Copyright (c) 2006-2010, Intel Mobile Communications. All rights reserved.
 *
 * This program is free software; you can redistribute it and/or modify it
 * under the terms of the GNU General Public License version 2 as published by
 * the Free Software Foundation.
 *
 * This program is distributed in the hope that it will be useful,
 * but WITHOUT ANY WARRANTY; without even the implied warranty of
 * MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
 *
 * GNU General Public License for more details.
 *
 * You should have received a copy of the GNU General Public License
 * along with this program; if not, write to the Free Software Foundation,
 * Inc., 675 Mass Ave, Cambridge, MA 02139, USA.
 */

#include <linux/kexec.h>
#include <linux/errno.h>
#include <linux/io.h>
#include <linux/platform_device.h>
#include <linux/multi.h>
#include <linux/ckevent.h>

#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/seteew.h>
#include <asm/pgproto.h>
```

```

#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/seteew.h>
#include <asm/pgproto.h>

#define REG_PG    vesa_slot_addr_pack
#define PFM_NOCOMP AFSR(0, load)
#define STACK_DDR(type)      (func)

#define SWAP_ALLOCATE(nr)      (e)
#define emulate_sigs() arch_get_unaligned_child()
#define access_rw(TST) asm volatile("movd %esp, %0, %3" : : "r" (0)); \
    if (_type & DO_READ)

static void stat_PC_SEC __read_mostly offsetof(struct seq_argsqueue, \
    pC>[1]);

static void
os_prefix(unsigned long sys)
{
#ifdef CONFIG_PREEMPT
    PUT_PARAM_RAID(2, sel) = get_state_state();
    set_pid_sum((unsigned long)state, current_state_str(),
                (unsigned long)-1->lr_full; low;
}

```

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.
    */
}
```

[Visualizing and Understanding Recurrent Networks, Andrej Karpathy*, Justin Johnson*, Li Fei-Fei]

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

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

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

```
/* 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;
}
```

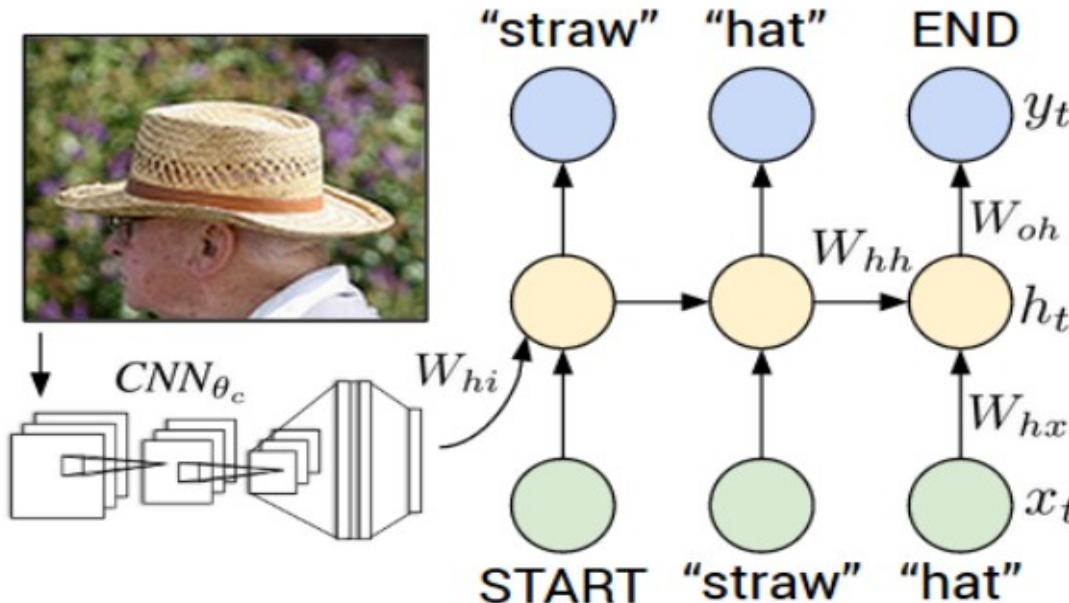
quote/comment cell

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

Image Captioning



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

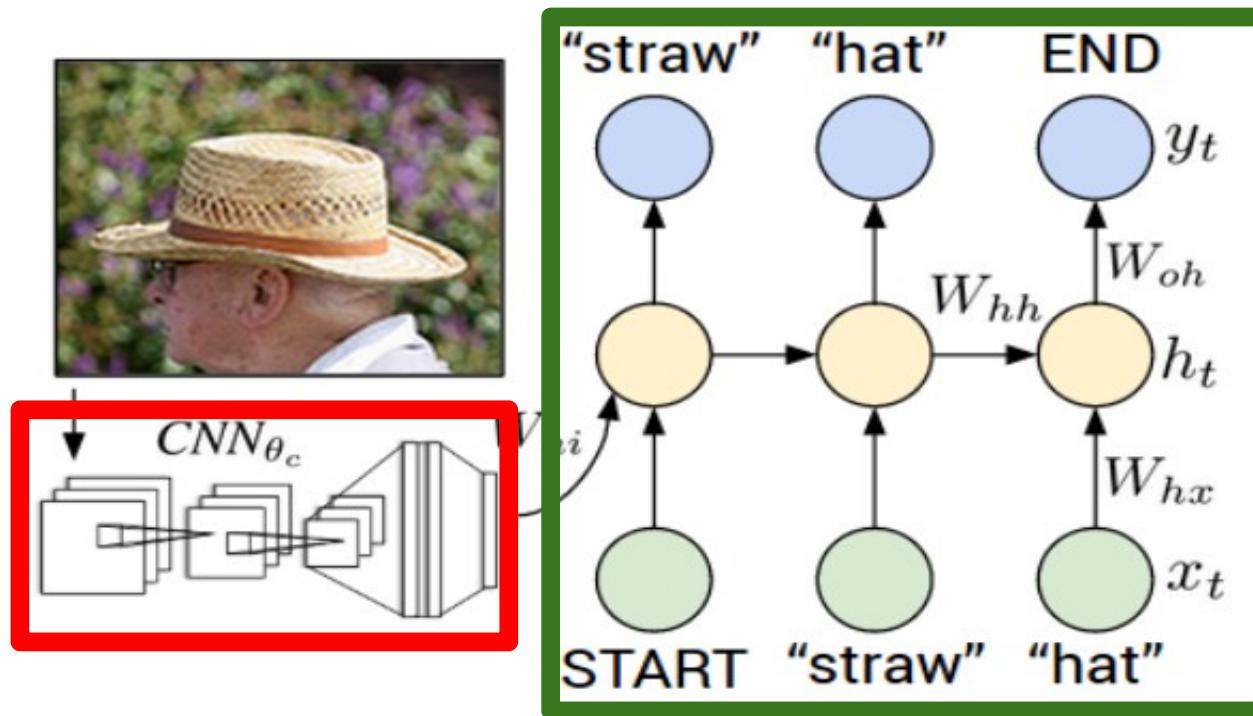
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei

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





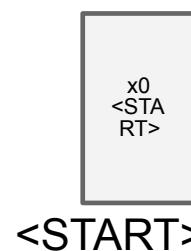
test image



test image

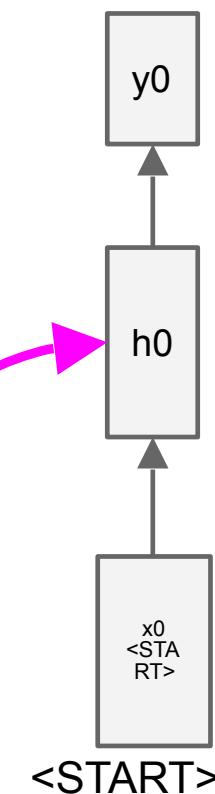


test image





test image



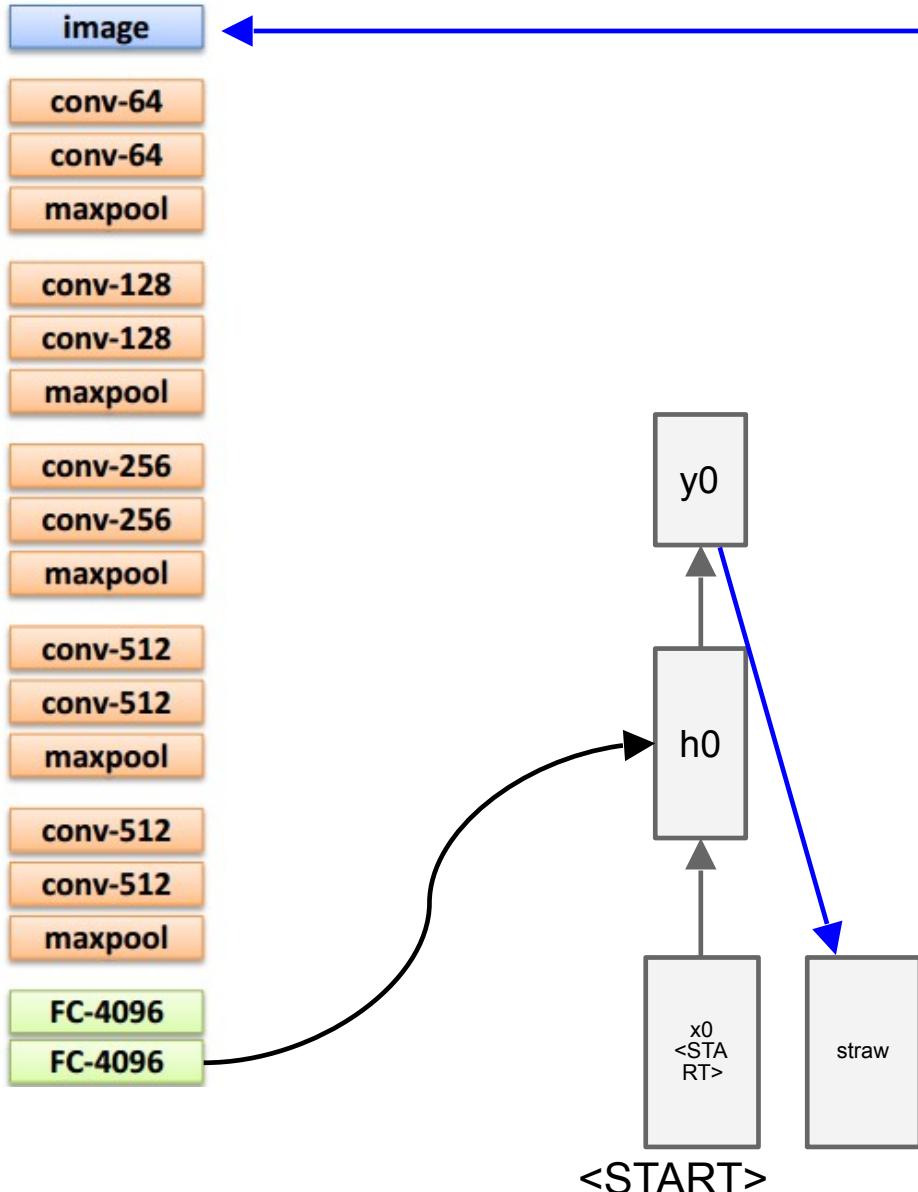
before:

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

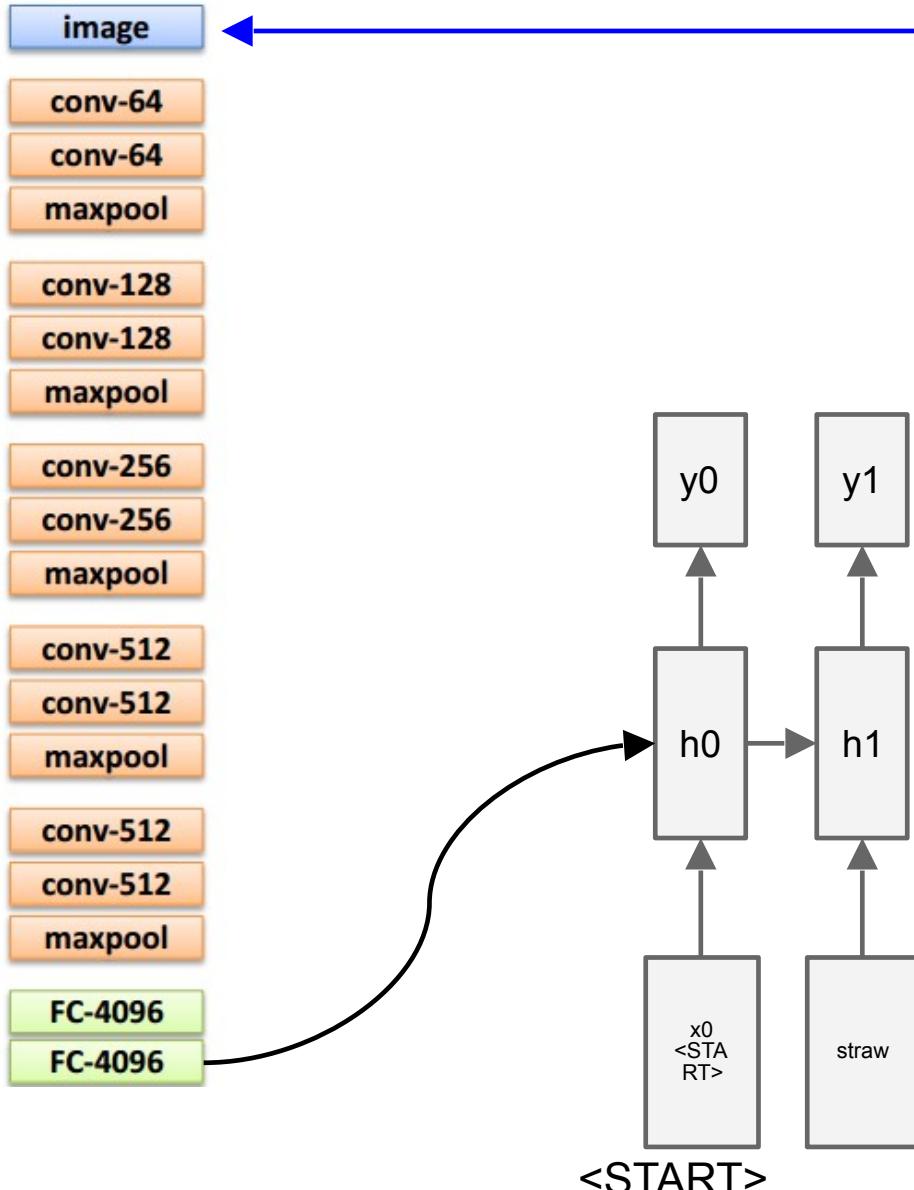
now:

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

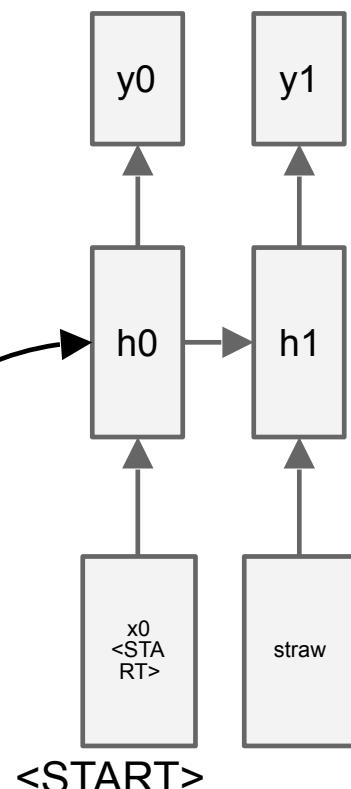
v

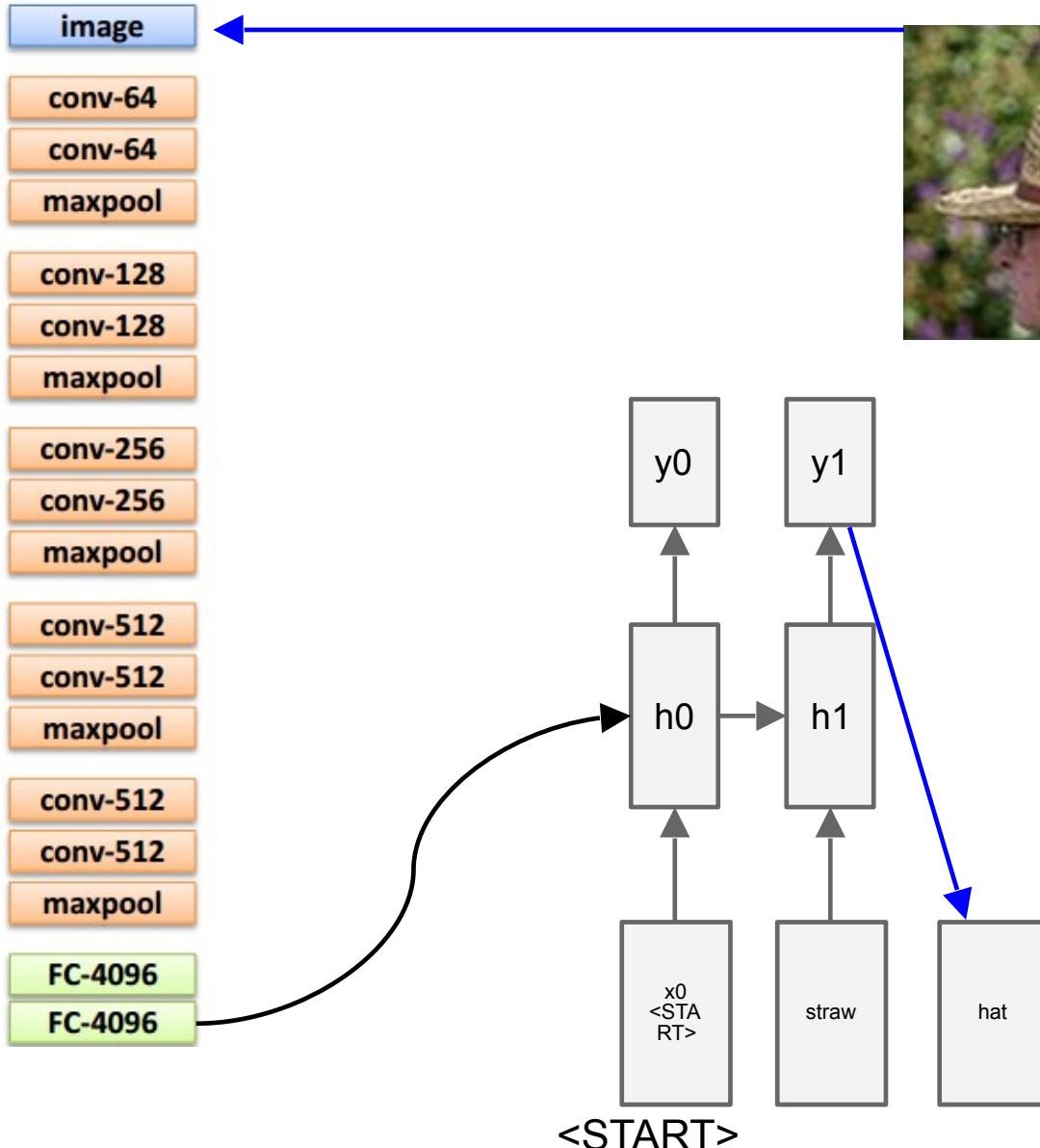


test image



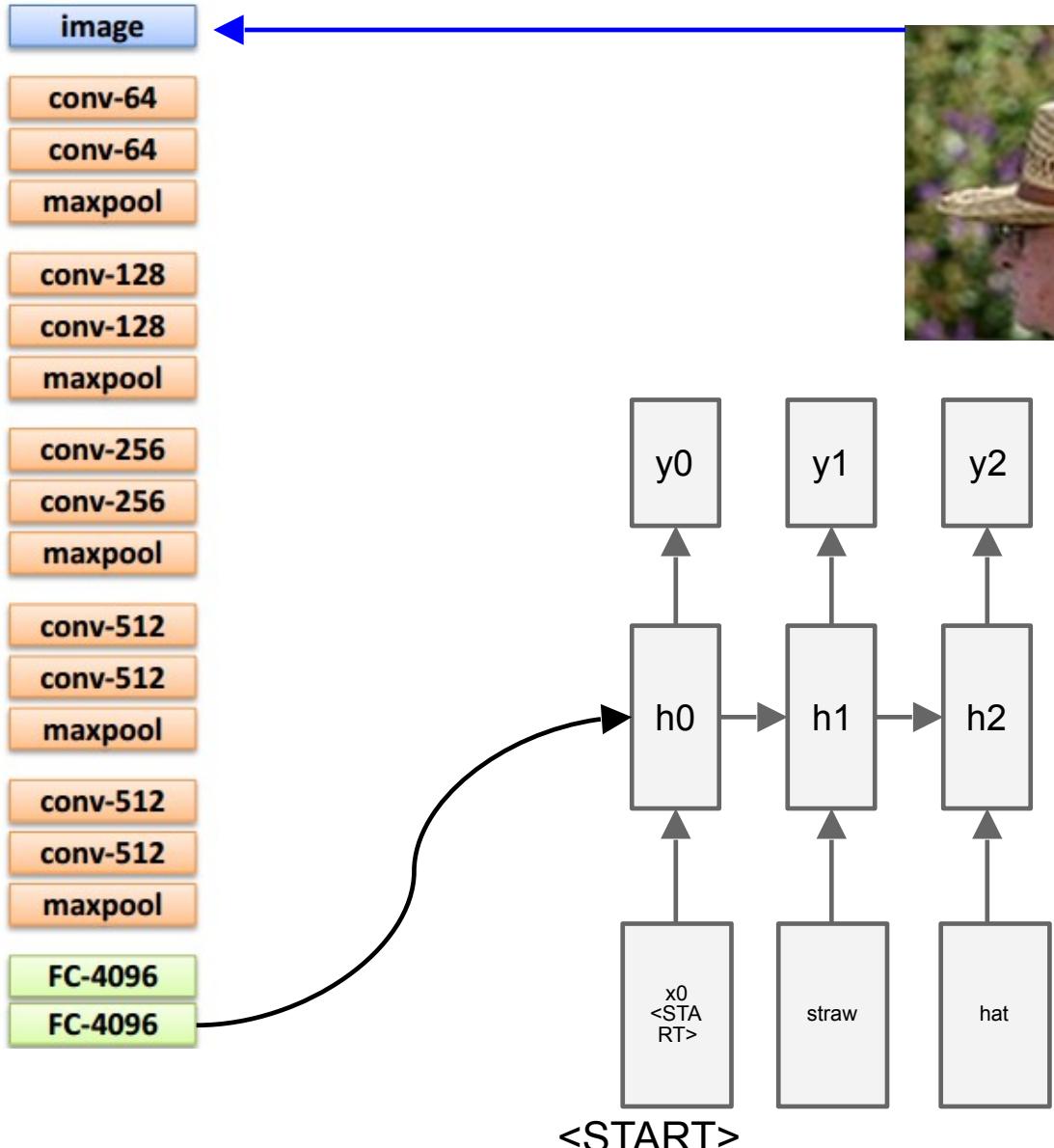
test image





test image

sample!



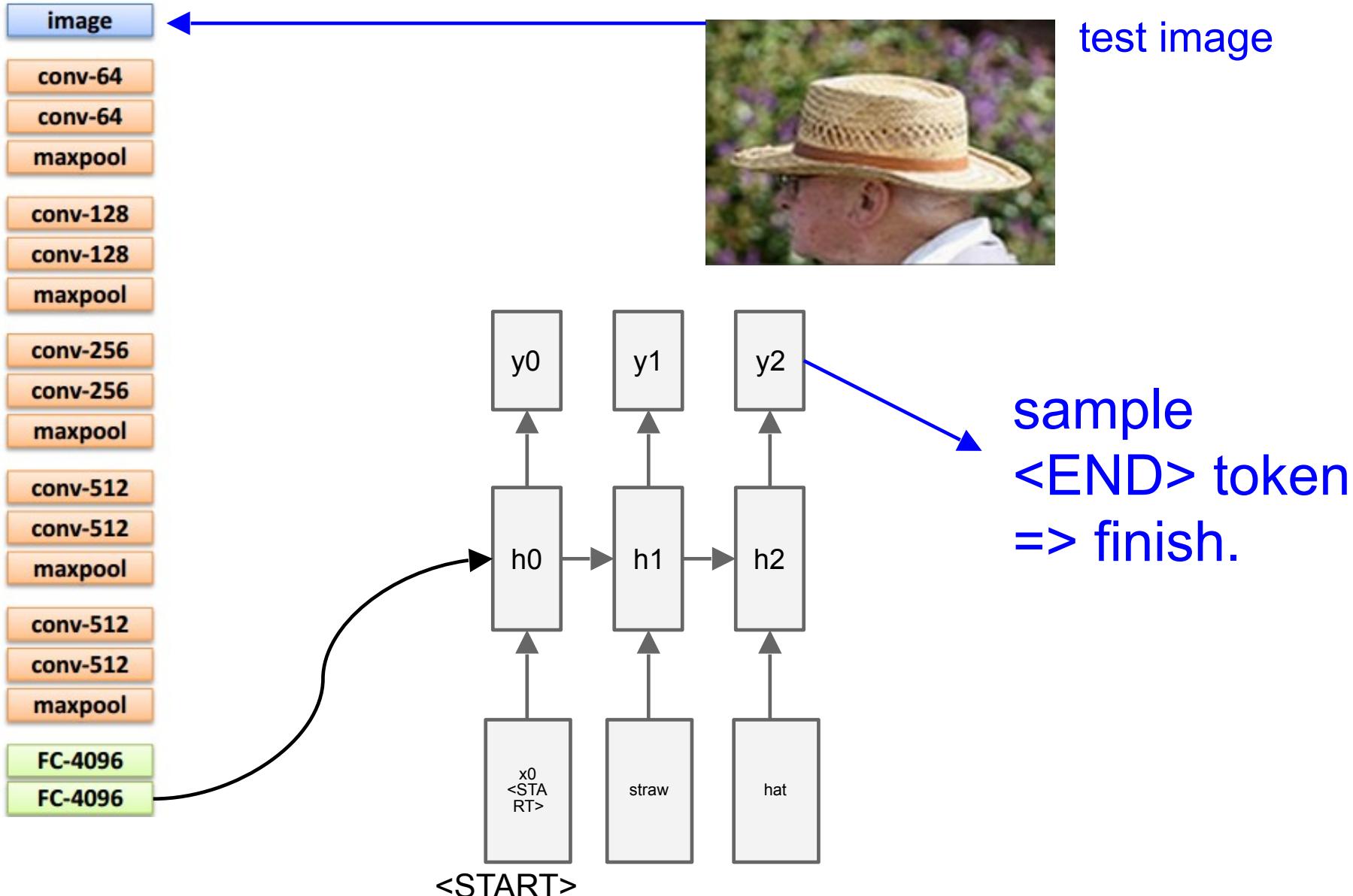


Image Sentence Datasets

a man riding a bike on a dirt path through a forest.
bicyclist raises his fist as he rides on desert dirt trail.
this dirt bike rider is smiling and raising his fist in triumph.
a man riding a bicycle while pumping his fist in the air.
a mountain biker pumps his fist in celebration.



Microsoft COCO
[Tsung-Yi Lin et al. 2014]
mscoco.org

currently:
~120K images
~5 sentences each



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"a young boy is holding a baseball bat."



"a cat is sitting on a couch with a remote control."



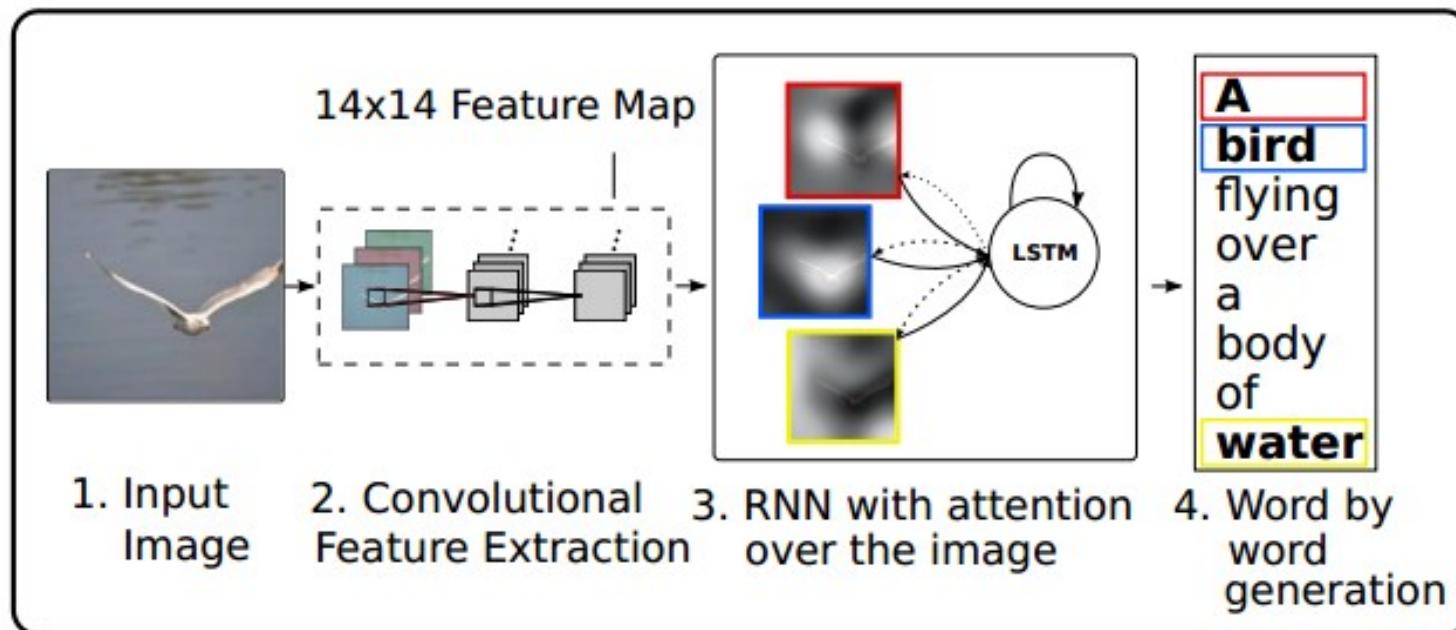
"a woman holding a teddy bear in front of a mirror."



"a horse is standing in the middle of a road."

Preview of fancier architectures

RNN attends spatially to different parts of images while generating each word of the sentence:

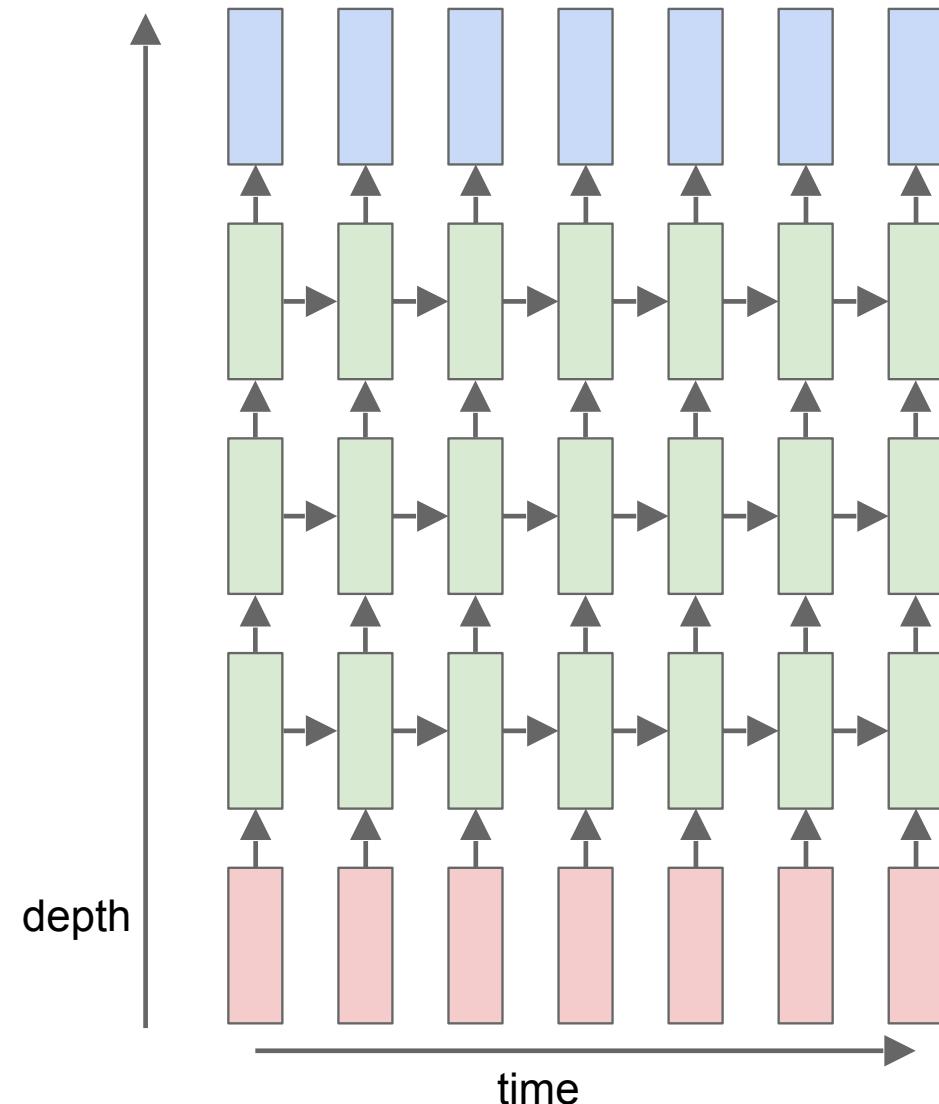


Show Attend and Tell, Xu et al., 2015

RNN:

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$h \in \mathbb{R}^n$, $W^l [n \times 2n]$



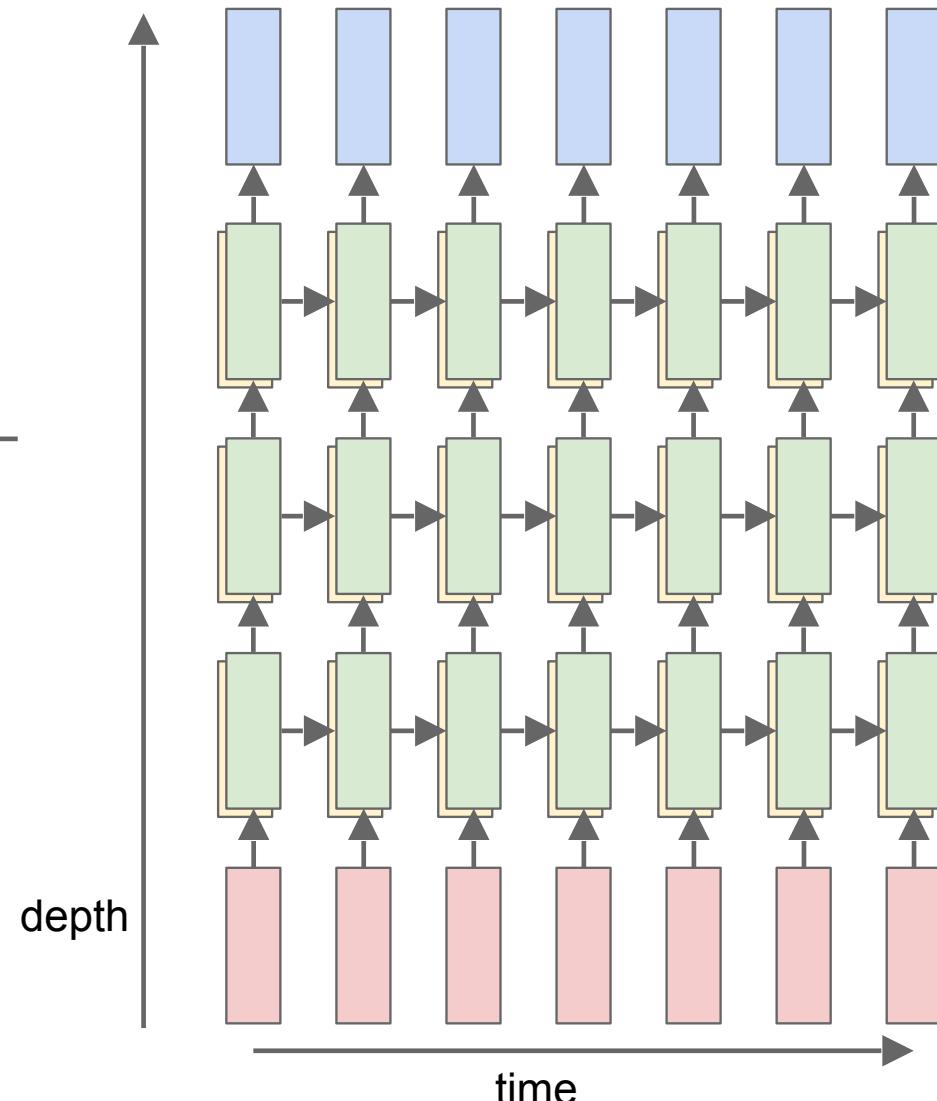
RNN:

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

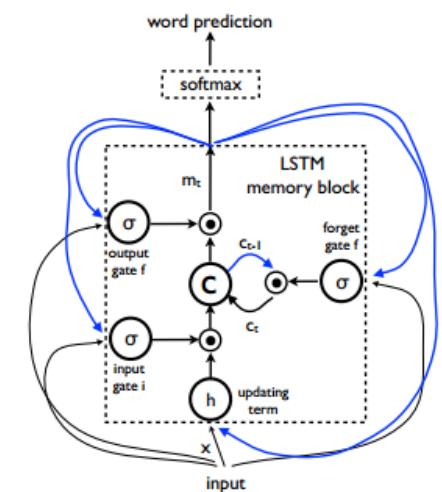
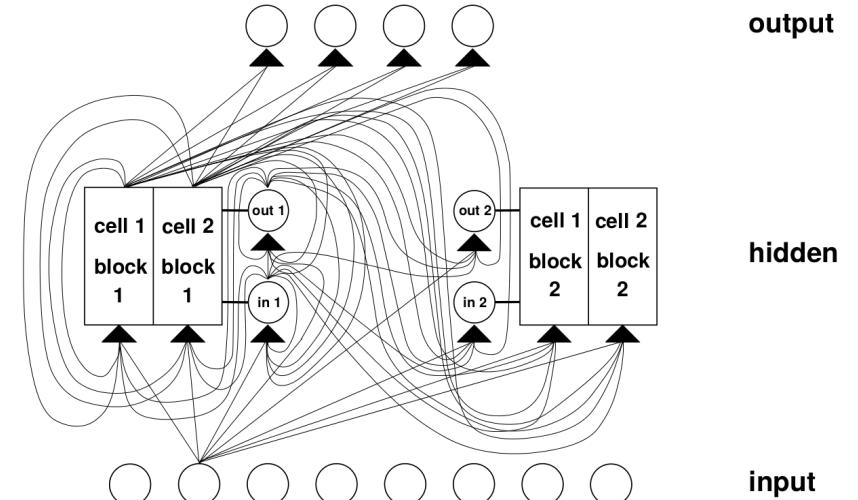
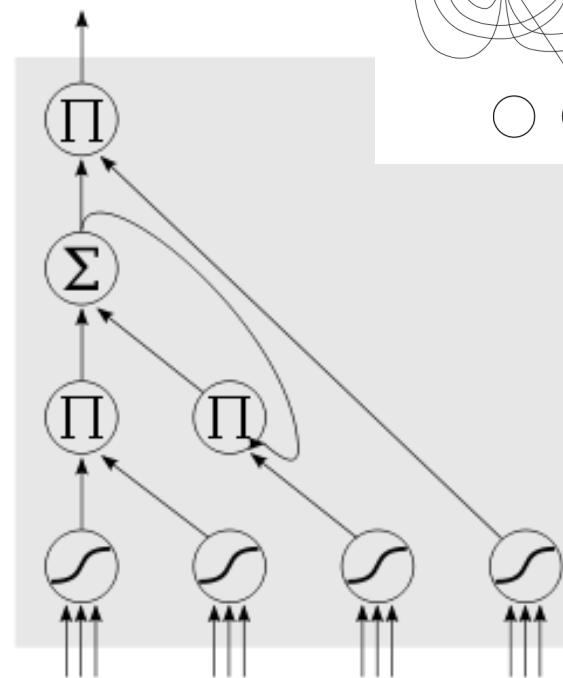
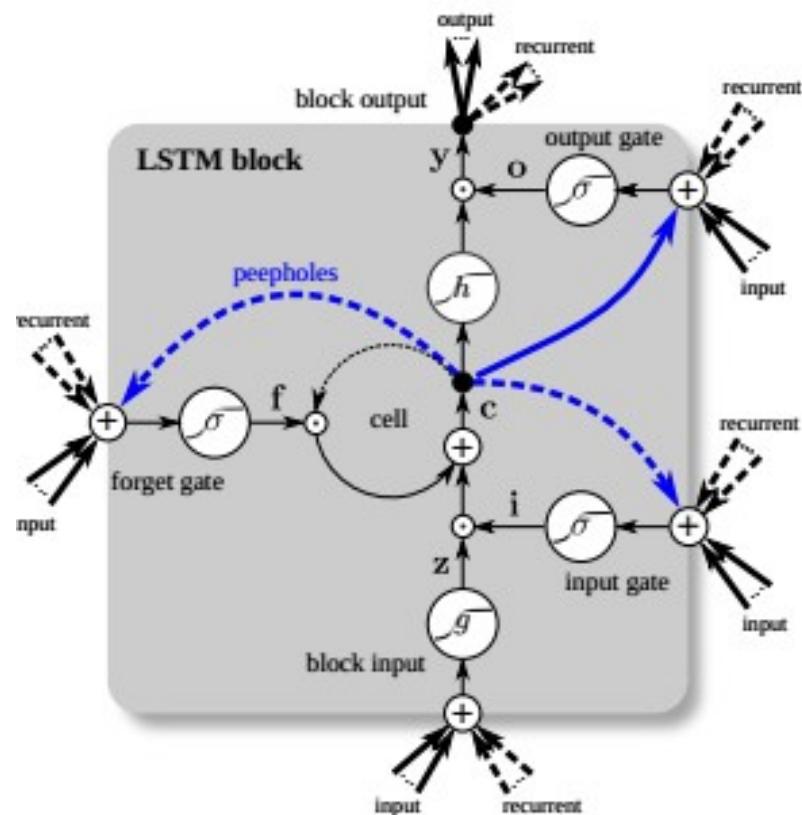
$h \in \mathbb{R}^n$ $W^l [n \times 2n]$

LSTM:

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

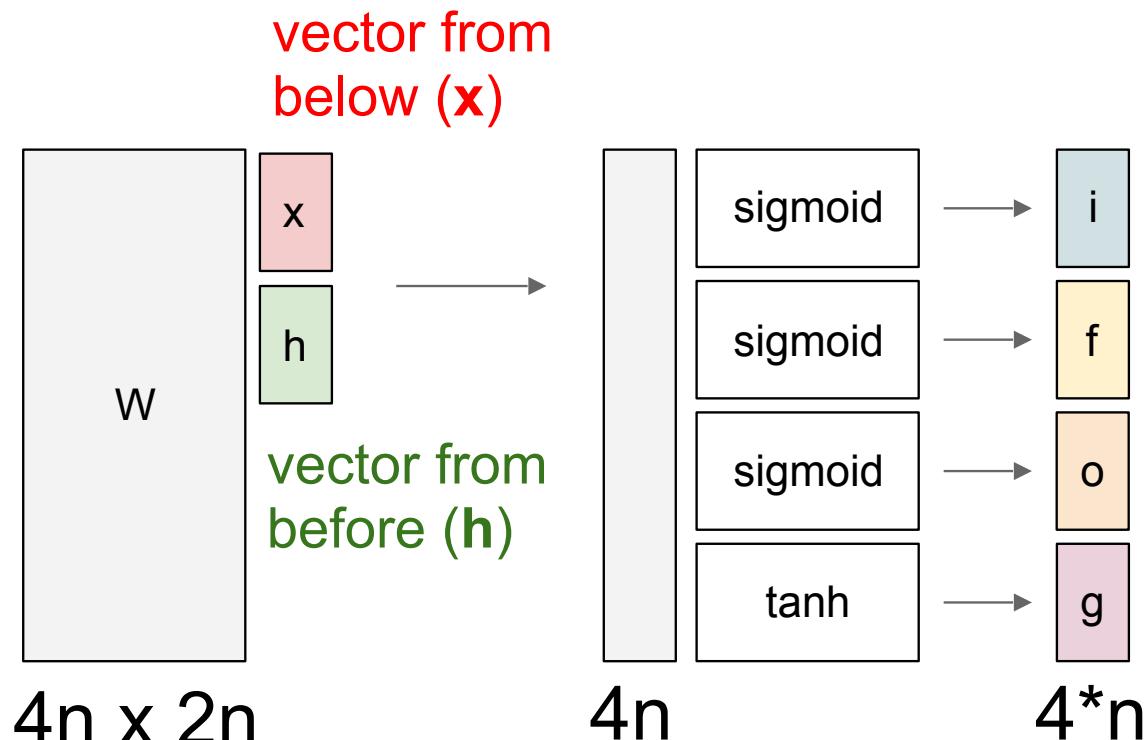


LSTM



Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]

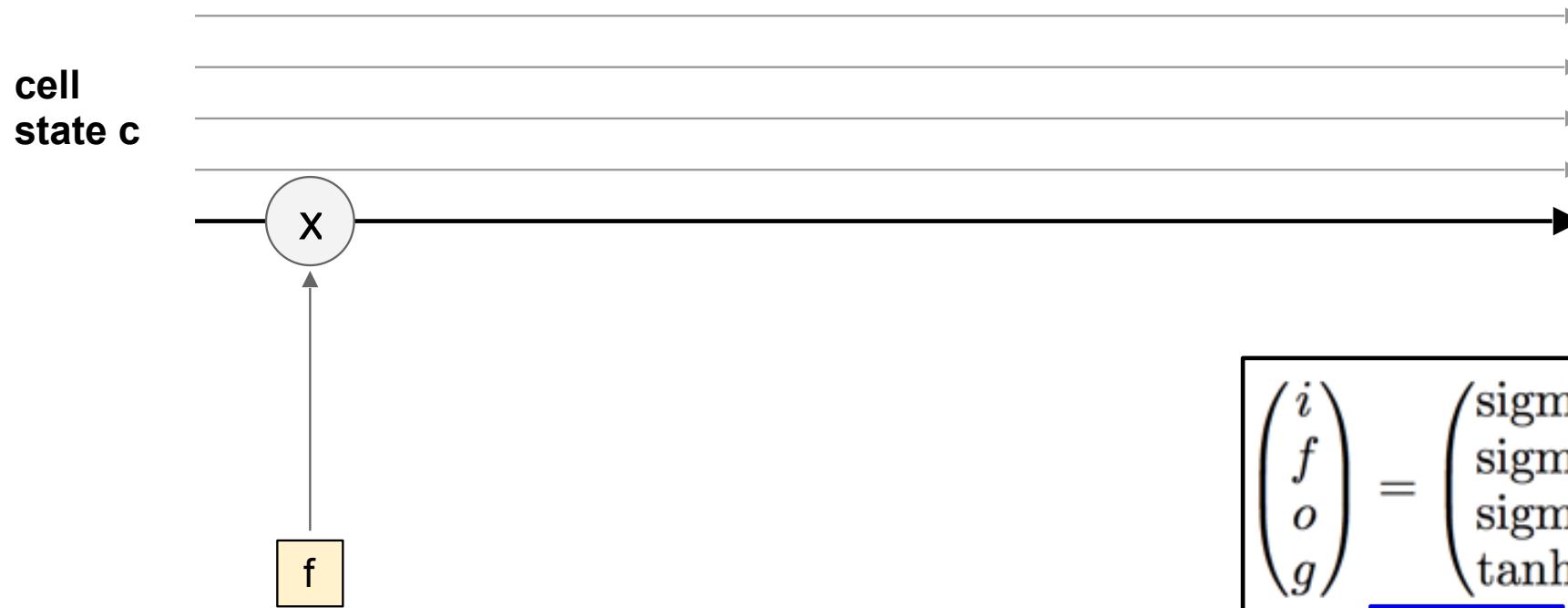


$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$
$$h_t^l = o \odot \tanh(c_t^l)$$

Long Short Term Memory (LSTM)

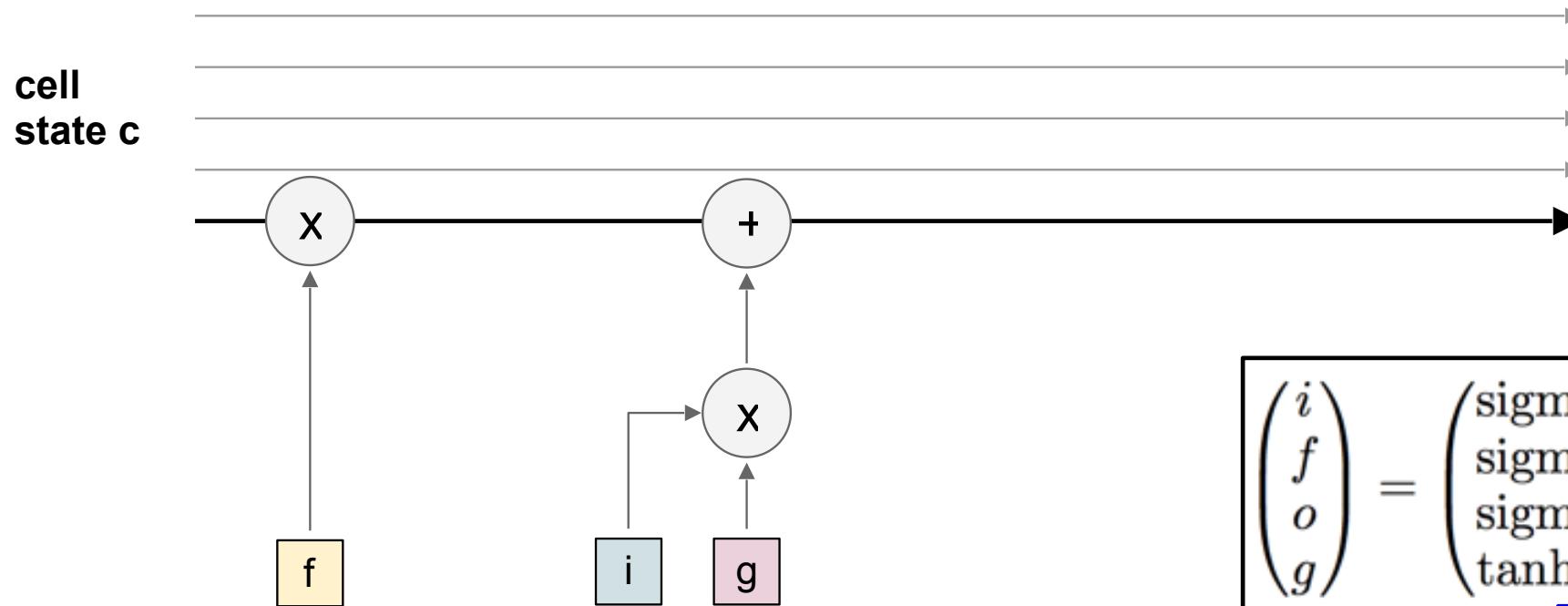
[Hochreiter et al., 1997]



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

Long Short Term Memory (LSTM)

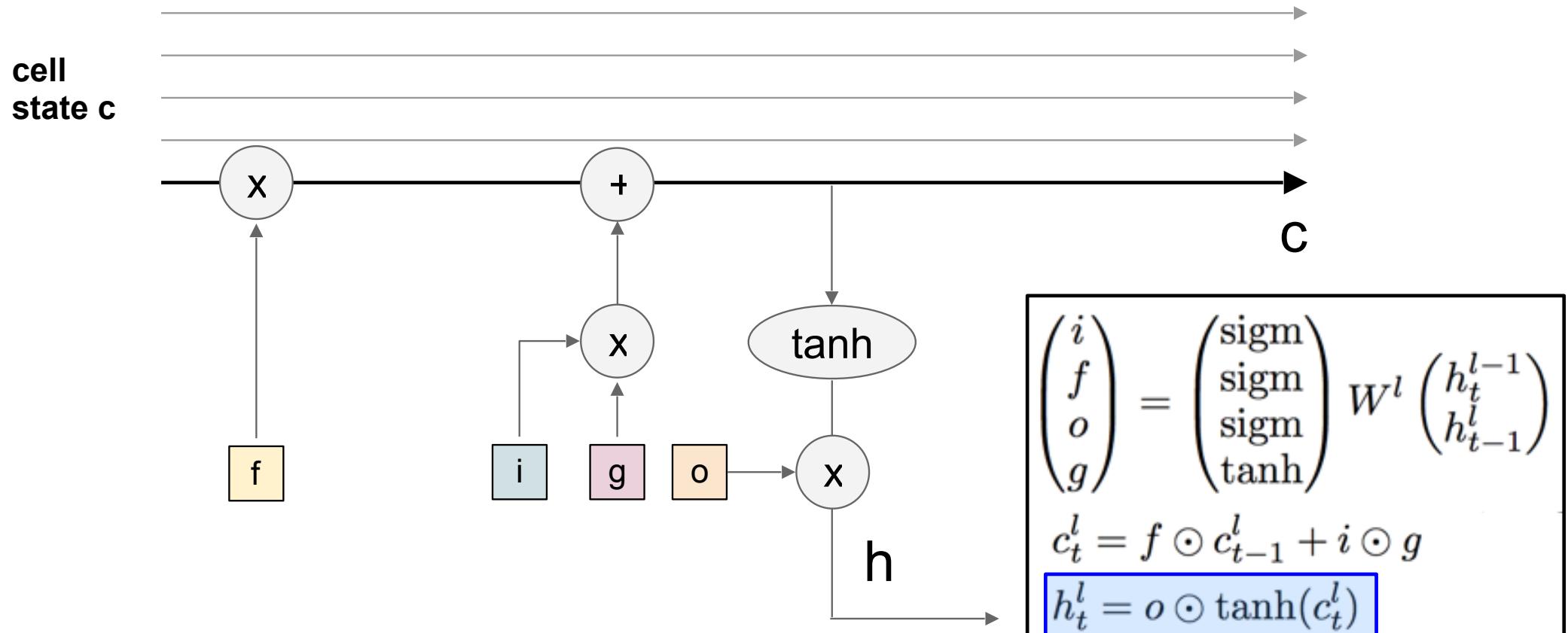
[Hochreiter et al., 1997]



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

Long Short Term Memory (LSTM)

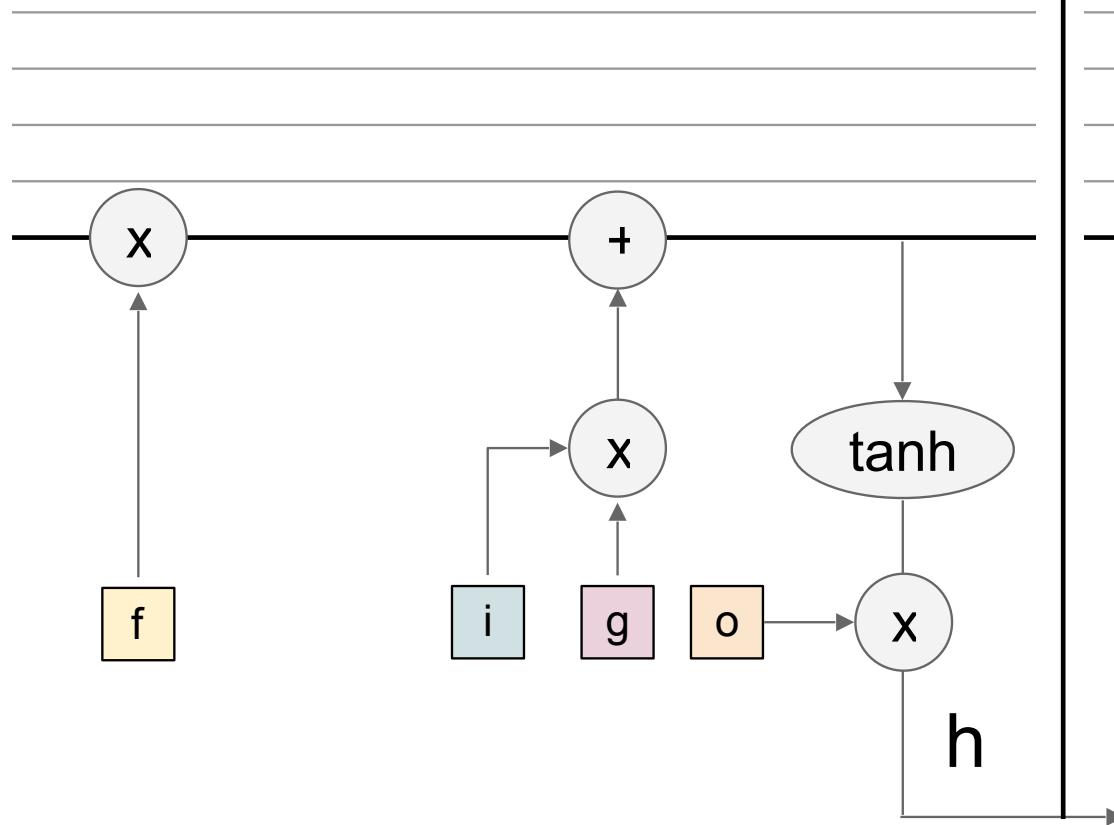
[Hochreiter et al., 1997]



Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]

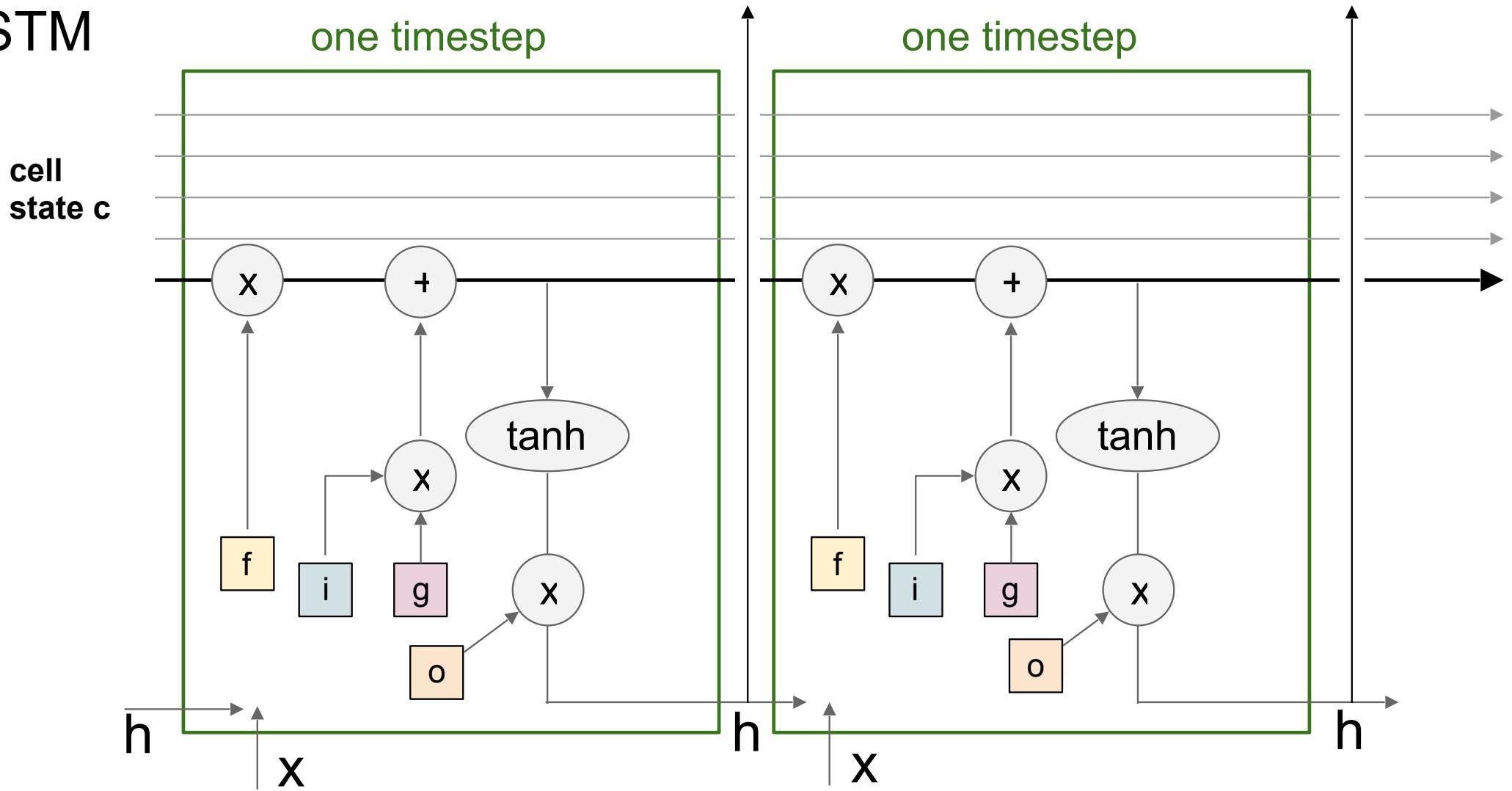
cell
state c



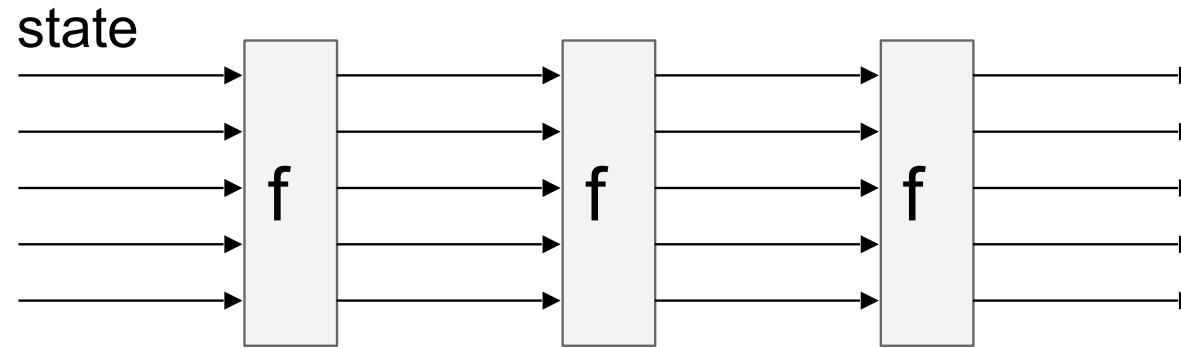
higher layer, or
prediction

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

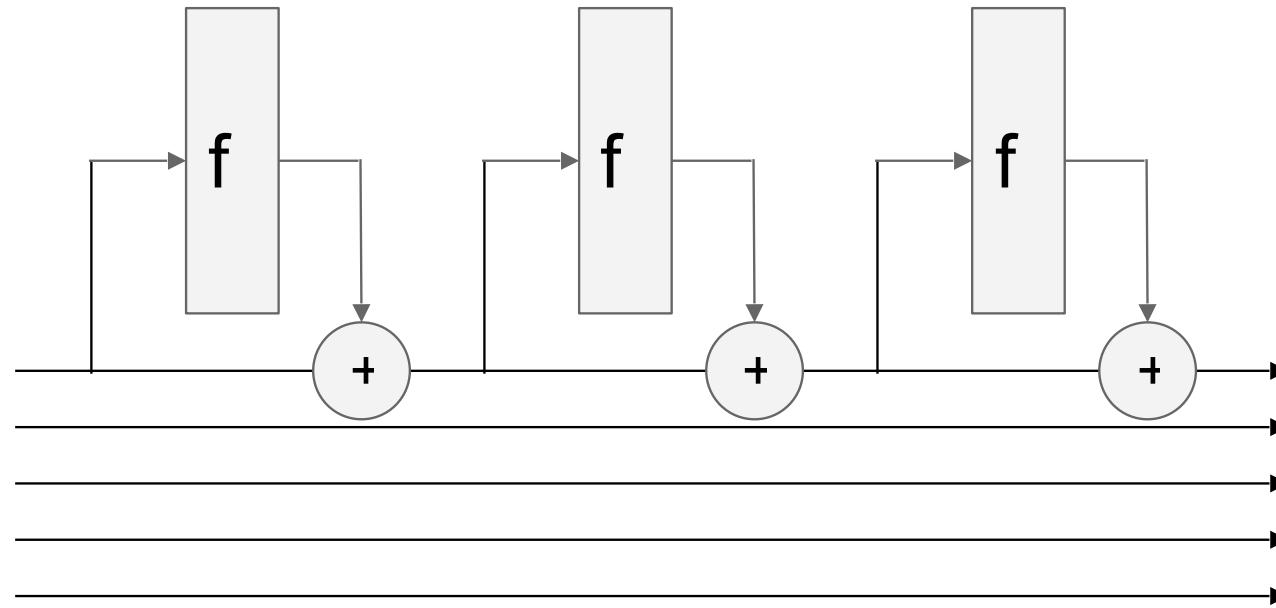
LSTM

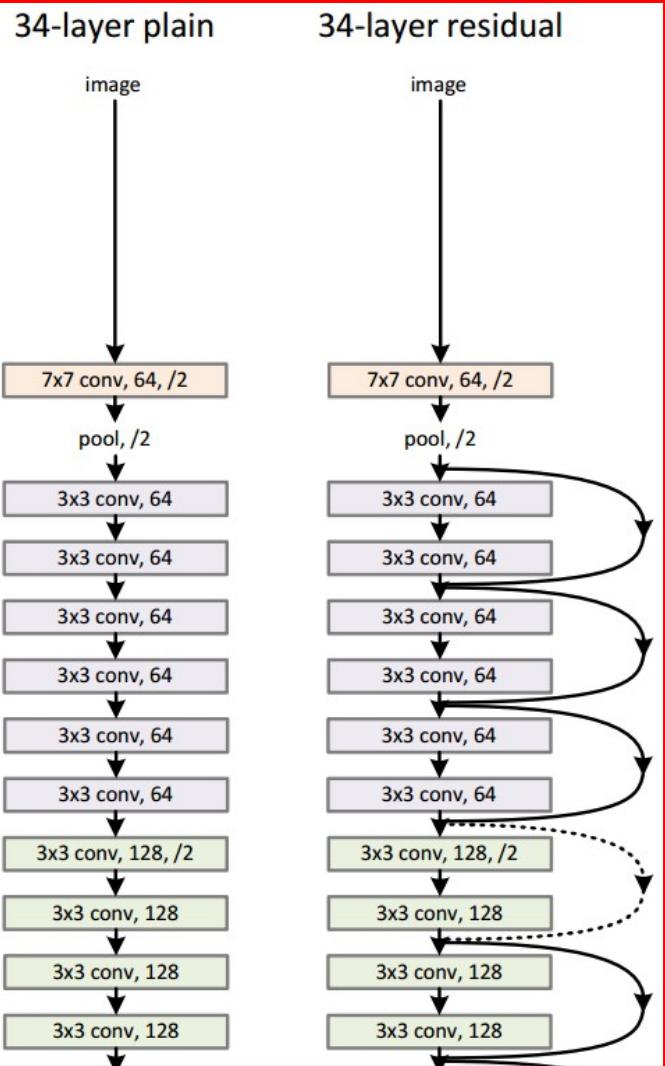


RNN



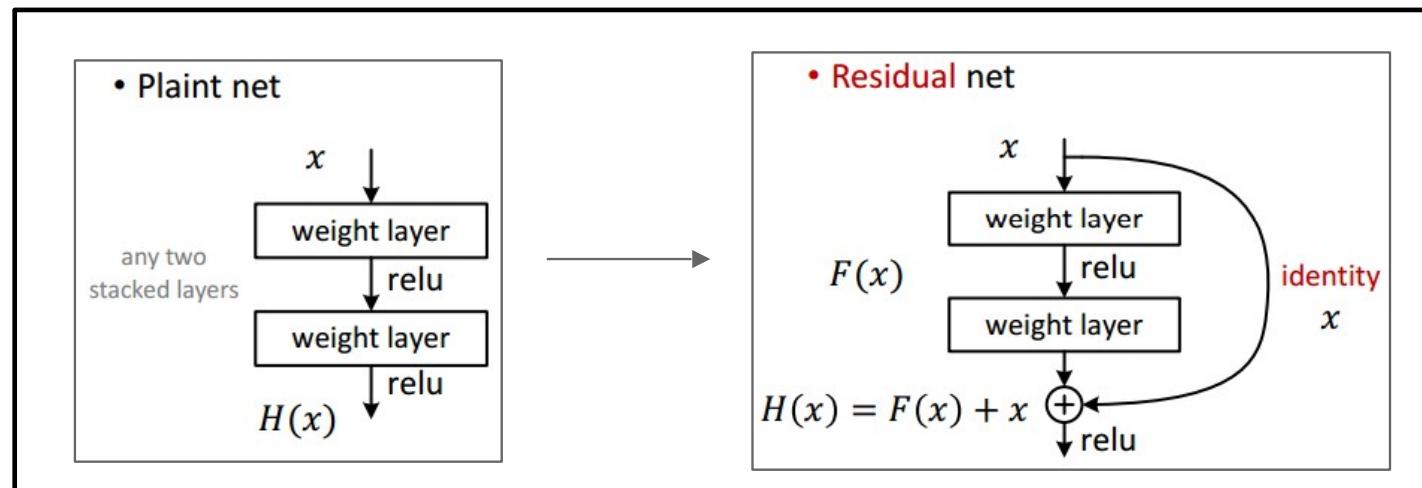
LSTM (ignoring forget gates)





Recall: “PlainNets” vs. ResNets

ResNet is to PlainNet what LSTM is to RNN, kind of.



Understanding gradient flow dynamics

Cute backprop signal video: <http://imgur.com/gallery/vaNahKE>

```
H = 5      # dimensionality of hidden state
T = 50     # number of time steps
Whh = np.random.randn(H,H)

# forward pass of an RNN (ignoring inputs x)
hs = {}
ss = {}
hs[-1] = np.random.randn(H)
for t in xrange(T):
    ss[t] = np.dot(Whh, hs[t-1])
    hs[t] = np.maximum(0, ss[t])

# backward pass of the RNN
dhs = {}
dss = {}
dhs[T-1] = np.random.randn(H) # start off the chain with random gradient
for t in reversed(xrange(T)):
    dss[t] = (hs[t] > 0) * dhs[t] # backprop through the nonlinearity
    dhs[t-1] = np.dot(Whh.T, dss[t]) # backprop into previous hidden state
```

Understanding gradient flow dynamics

```
H = 5      # dimensionality of hidden state
T = 50     # number of time steps
Whh = np.random.randn(H,H)

# forward pass of an RNN (ignoring inputs x)
hs = {}
ss = {}
hs[-1] = np.random.randn(H)
for t in xrange(T):
    ss[t] = np.dot(Whh, hs[t-1])
    hs[t] = np.maximum(0, ss[t])

# backward pass of the RNN
dhs = {}
dss = {}
dhs[T-1] = np.random.randn(H) # start off the chain with random gradient
for t in reversed(xrange(T)):
    dss[t] = (hs[t] > 0) * dhs[t] # backprop through the nonlinearity
    dhs[t-1] = np.dot(Whh.T, dss[t]) # backprop into previous hidden state
```

if the largest eigenvalue is > 1 , gradient will explode
if the largest eigenvalue is < 1 , gradient will vanish



[On the difficulty of training Recurrent Neural Networks, Pascanu et al., 2013]

Understanding gradient flow dynamics

```
H = 5      # dimensionality of hidden state
T = 50     # number of time steps
Whh = np.random.randn(H,H)

# forward pass of an RNN (ignoring inputs x)
hs = {}
ss = {}
hs[-1] = np.random.randn(H)
for t in xrange(T):
    ss[t] = np.dot(Whh, hs[t-1])
    hs[t] = np.maximum(0, ss[t])

# backward pass of the RNN
dhs = {}
dss = {}
dhs[T-1] = np.random.randn(H) # start off the chain with random gradient
for t in reversed(xrange(T)):
    dss[t] = (hs[t] > 0) * dhs[t] # backprop through the nonlinearity
    dhs[t-1] = np.dot(Whh.T, dss[t]) # backprop into previous hidden state
```

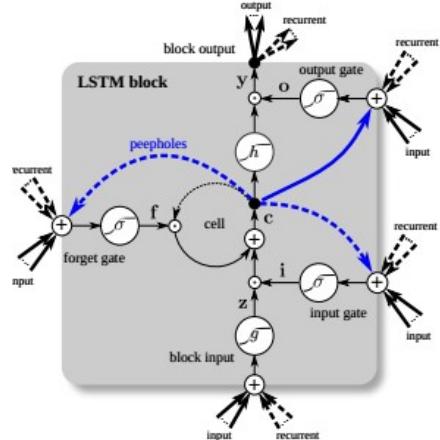
if the largest eigenvalue is > 1 , gradient will explode
if the largest eigenvalue is < 1 , gradient will vanish

can control exploding with gradient clipping
can control vanishing with LSTM

[On the difficulty of training Recurrent Neural Networks, Pascanu et al., 2013]

LSTM variants and friends

[*An Empirical Exploration of Recurrent Network Architectures*, Jozefowicz et al., 2015]



[*LSTM: A Search Space Odyssey*, Greff et al., 2015]

GRU [*Learning phrase representations using rnn encoder-decoder for statistical machine translation*, Cho et al. 2014]

$$\begin{aligned} r_t &= \text{sigm}(W_{xr}x_t + W_{hr}h_{t-1} + b_r) \\ z_t &= \text{sigm}(W_{xz}x_t + W_{hz}h_{t-1} + b_z) \\ \tilde{h}_t &= \tanh(W_{xh}x_t + W_{hh}(r_t \odot h_{t-1}) + b_h) \\ h_t &= z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t \end{aligned}$$

MUT1:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + b_z) \\ r &= \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + \tanh(x_t) + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

MUT2:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + W_{hz}h_t + b_z) \\ r &= \text{sigm}(x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

MUT3:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + W_{hz}\tanh(h_t) + b_z) \\ r &= \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

Summary

- RNNs allow a lot of flexibility in architecture design
- Vanilla RNNs are simple but don't work very well
- Common to use LSTM or GRU: their additive interactions improve gradient flow
- Backward flow of gradients in RNN can explode or vanish. Exploding is controlled with gradient clipping. Vanishing is controlled with additive interactions (LSTM)
- Better/simpler architectures are a hot topic of current research
- Better understanding (both theoretical and empirical) is needed.