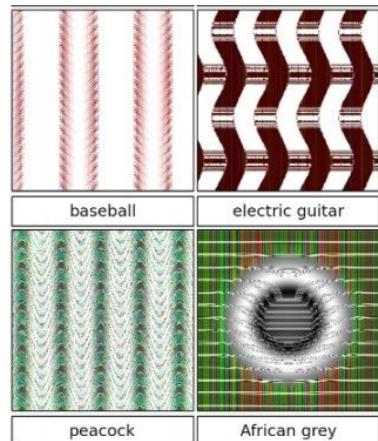
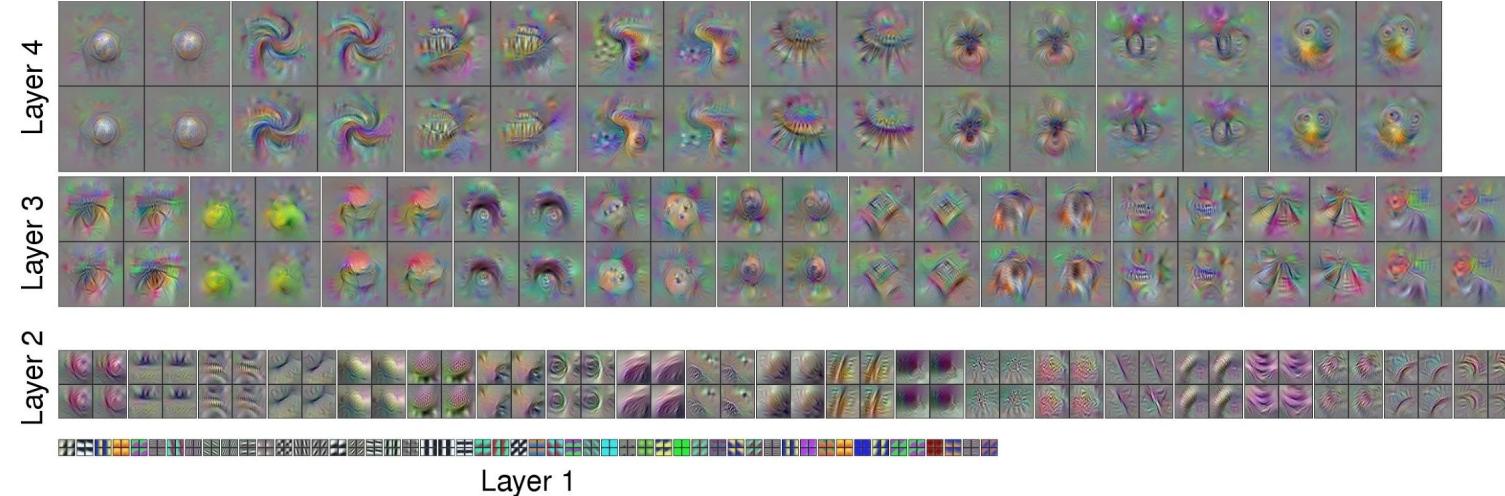
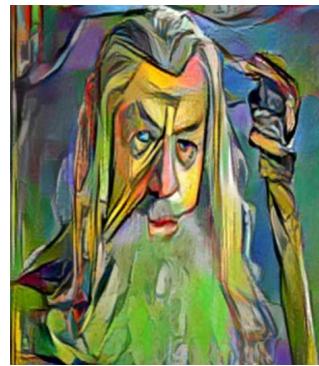
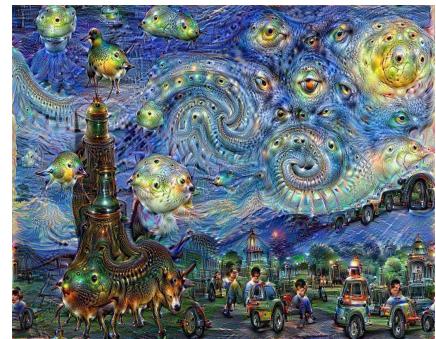
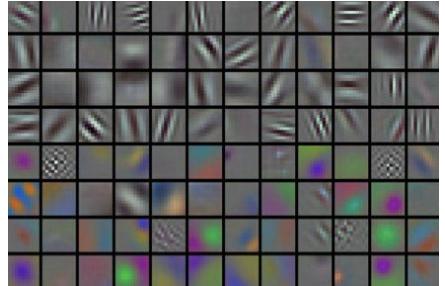
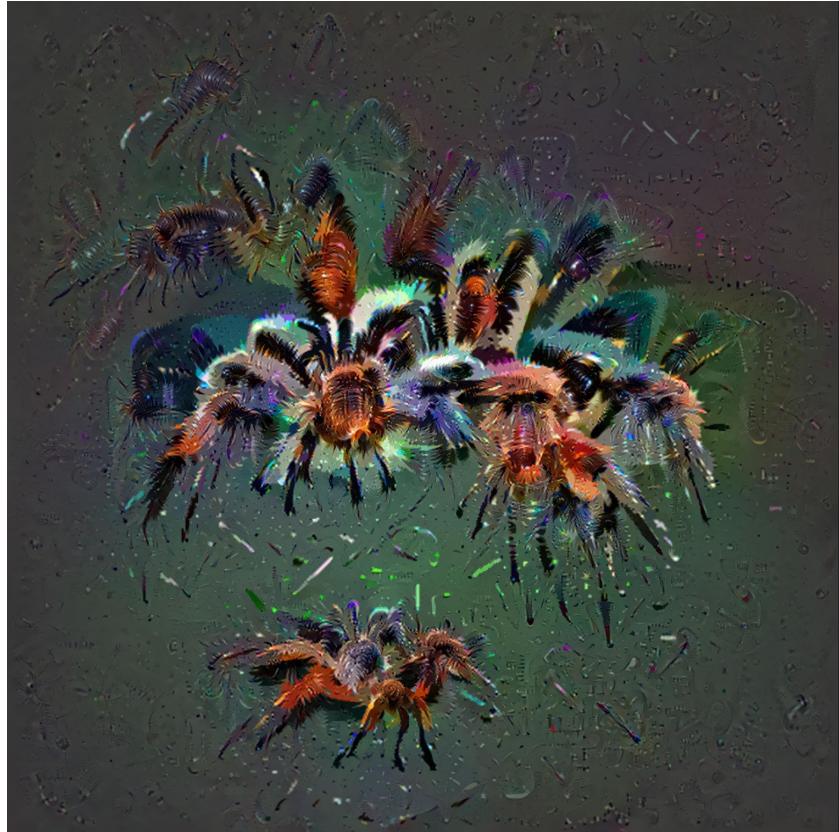
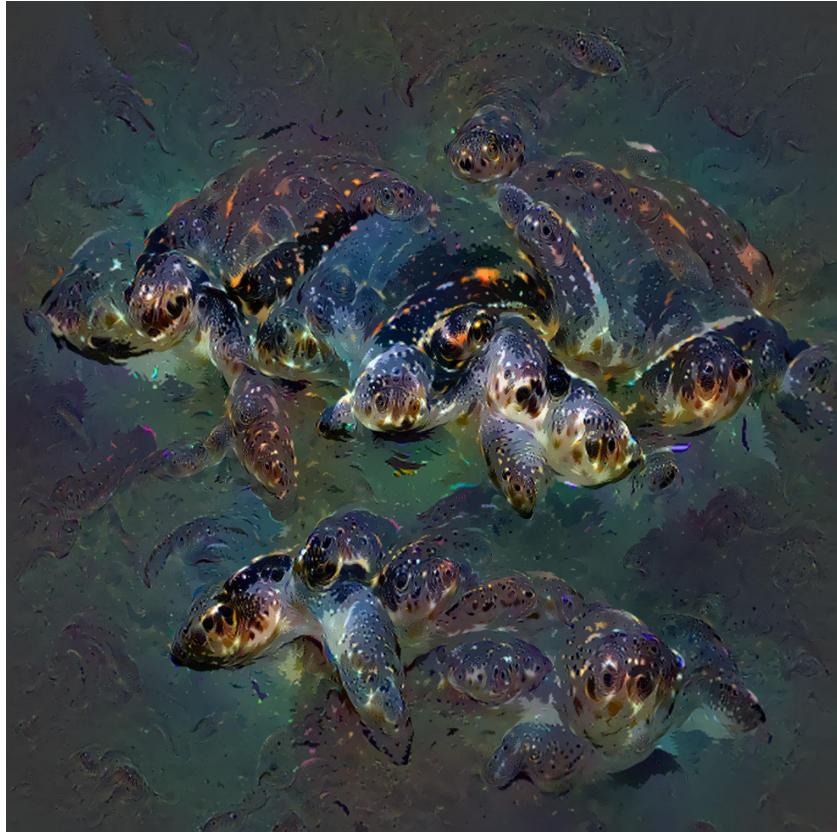


# Lecture 10:

## Recurrent Neural Networks



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

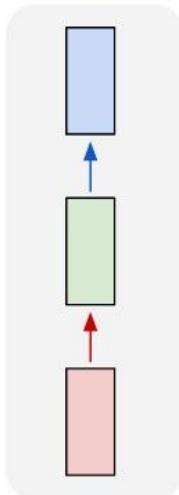


<http://mtyka.github.io/deepdream/2016/02/05/bilateral-class-vis.html>

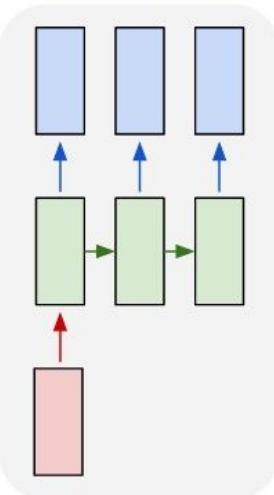


# Recurrent Networks offer a lot of flexibility:

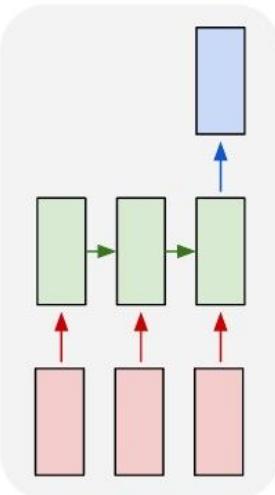
one to one



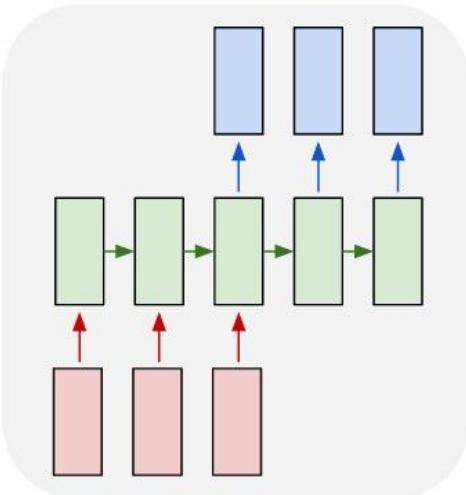
one to many



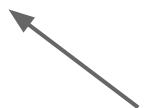
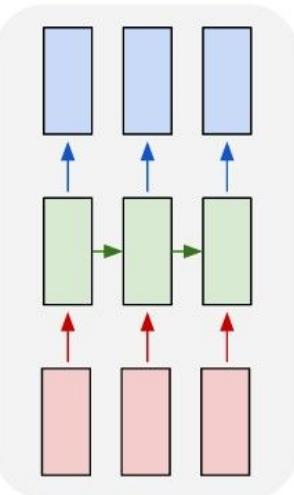
many to one



many to many



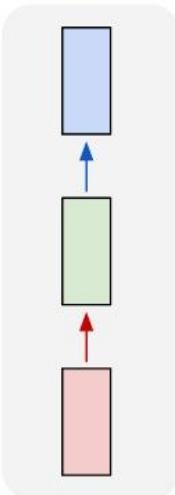
many to many



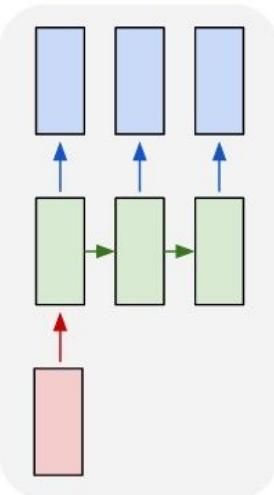
**Vanilla Neural Networks**

# Recurrent Networks offer a lot of flexibility:

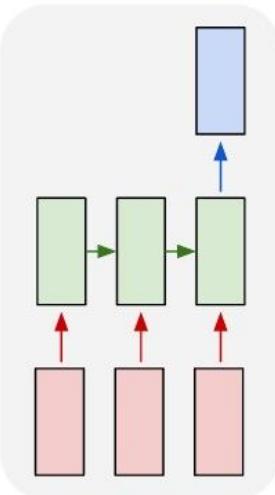
one to one



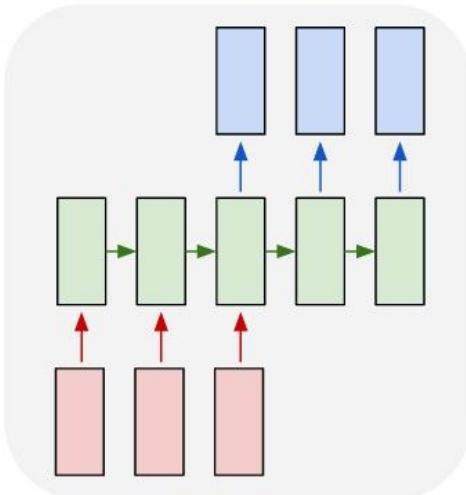
one to many



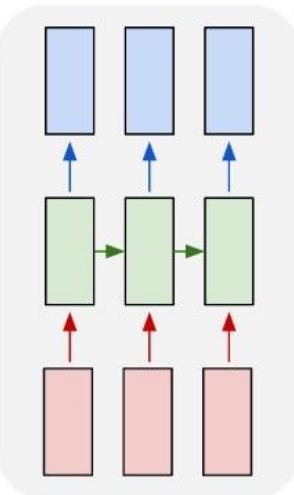
many to one



many to many



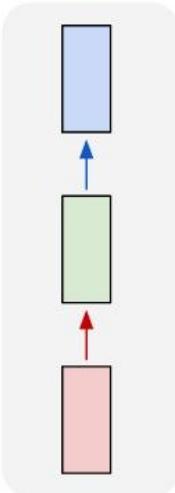
many to many



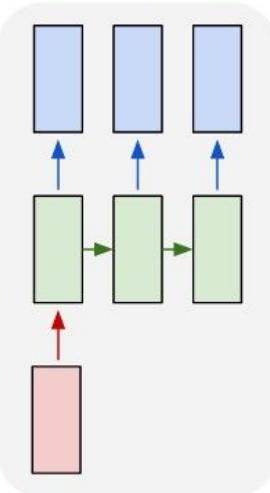
→ e.g. **Image Captioning**  
image -> sequence of words

# Recurrent Networks offer a lot of flexibility:

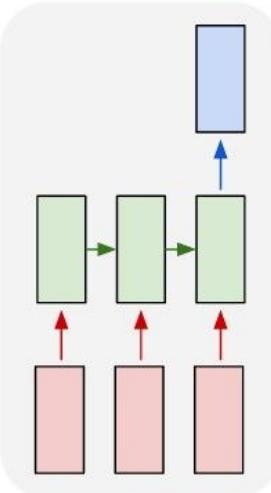
one to one



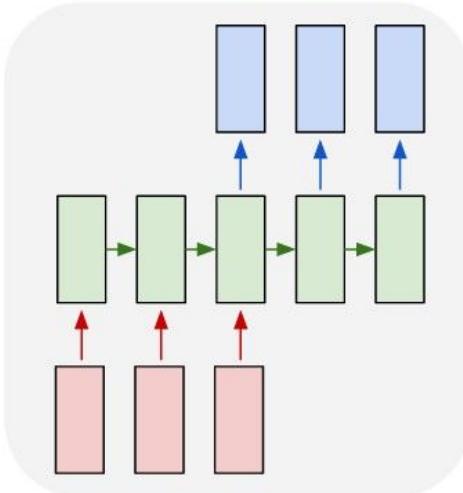
one to many



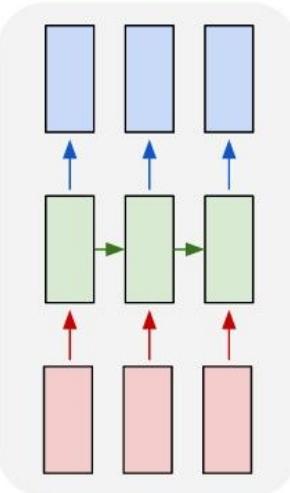
many to one



many to many



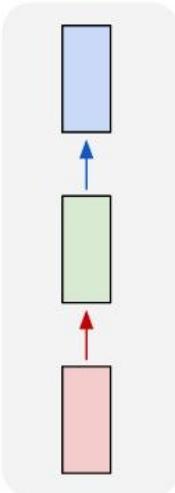
many to many



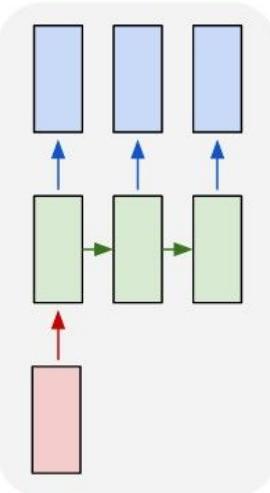
e.g. **Sentiment Classification**  
sequence of words -> sentiment

# Recurrent Networks offer a lot of flexibility:

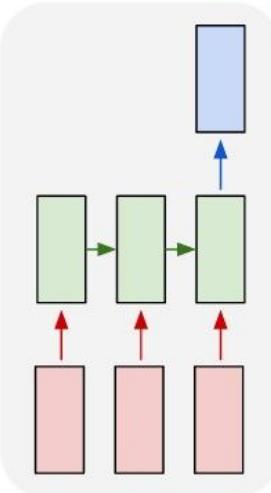
one to one



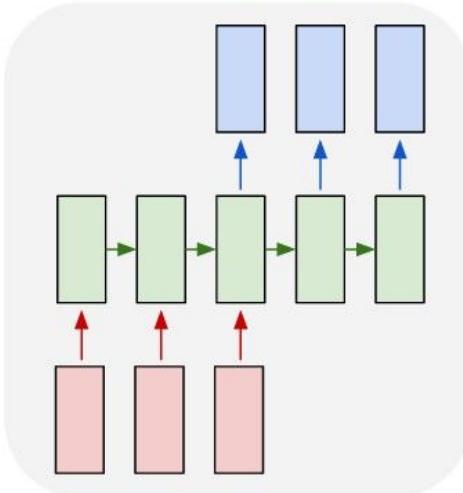
one to many



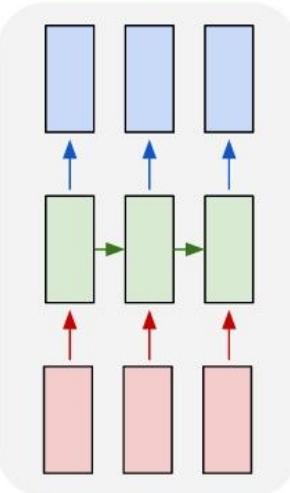
many to one



many to many



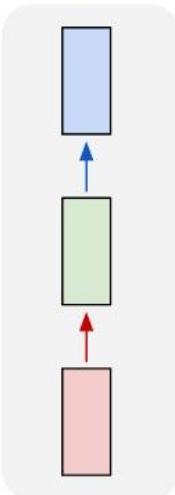
many to many



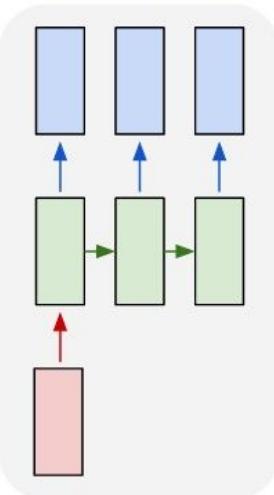
↑  
e.g. **Machine Translation**  
seq of words -> seq of words

# Recurrent Networks offer a lot of flexibility:

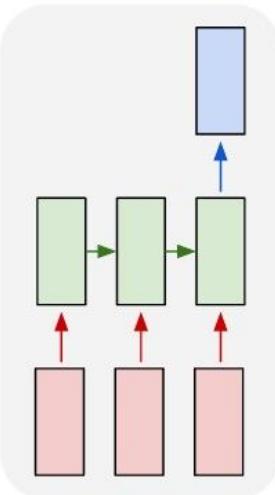
one to one



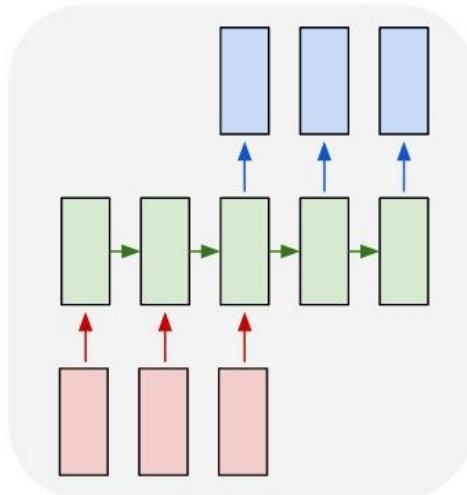
one to many



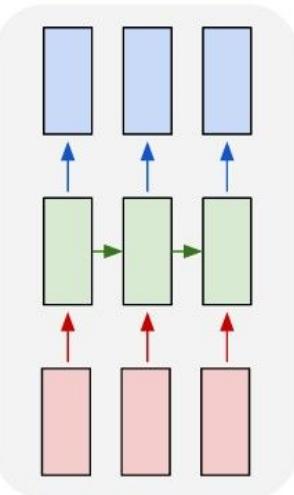
many to one



many to many

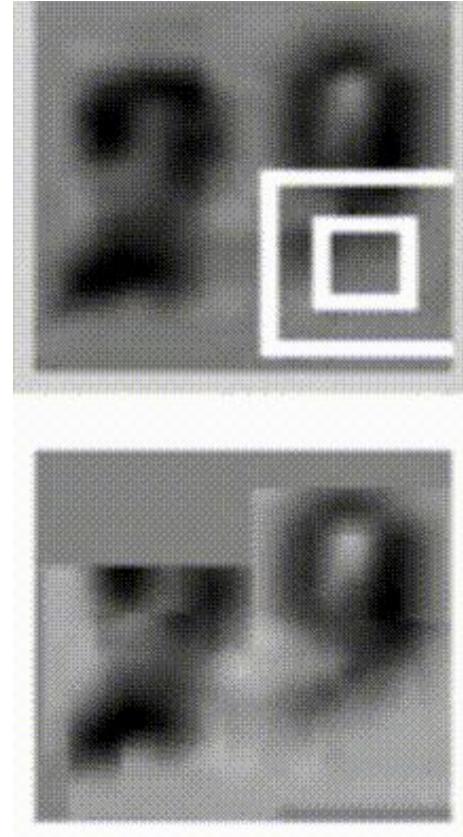


many to many



e.g. Video classification on frame level

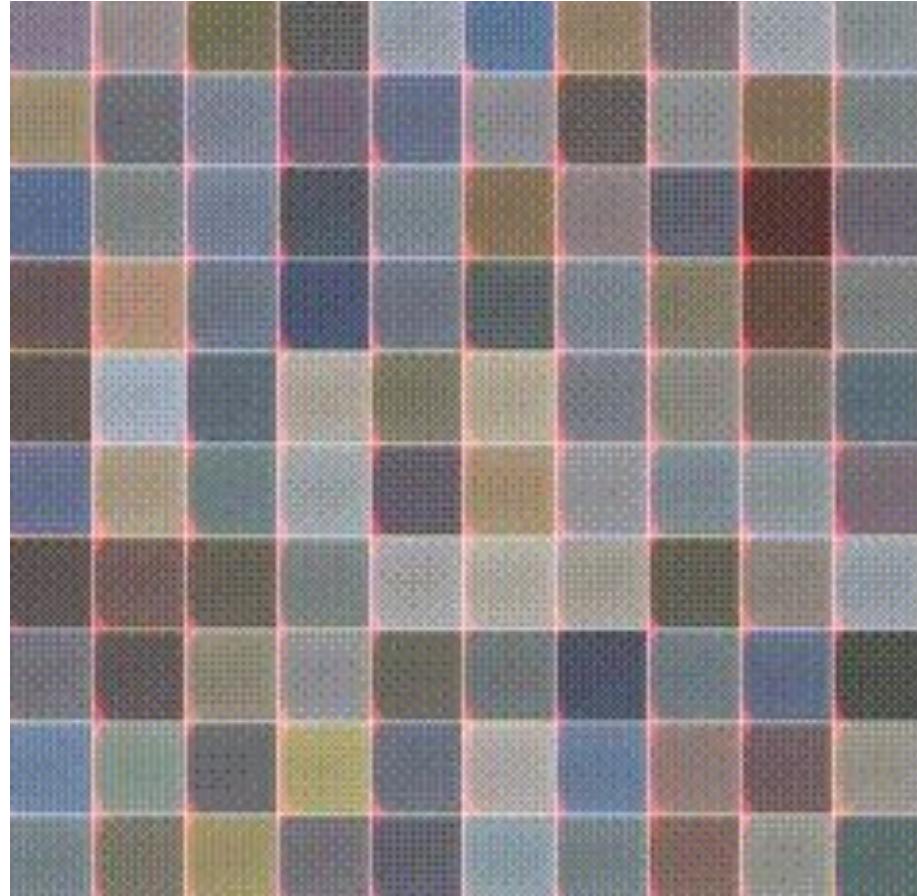
# Sequential Processing of fixed inputs



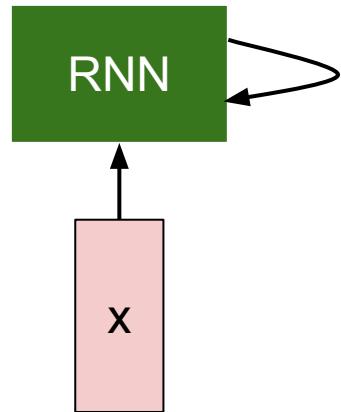
Multiple Object Recognition with  
Visual Attention, Ba et al.

# Sequential Processing of fixed outputs

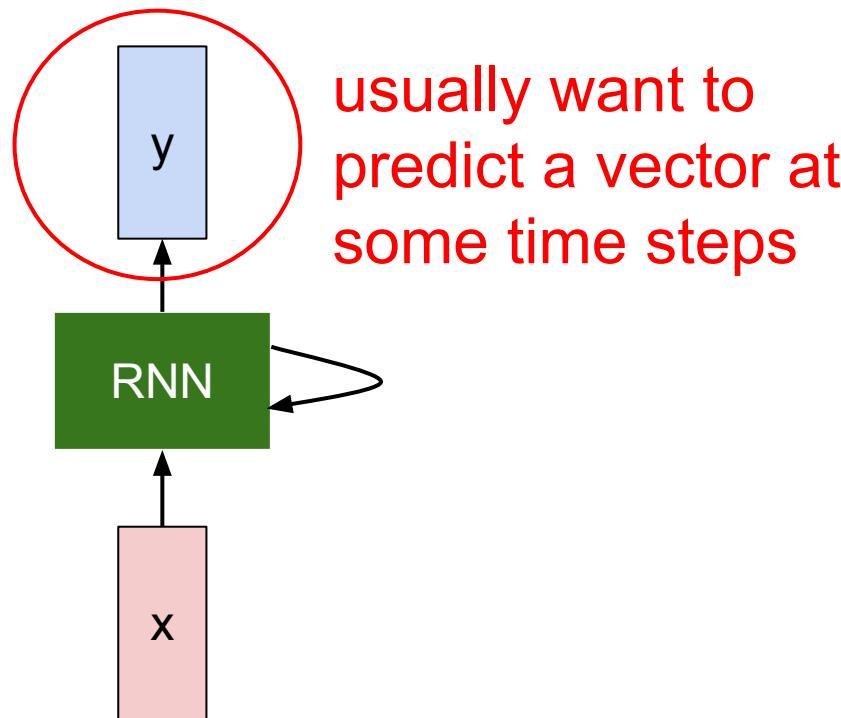
DRAW: A Recurrent  
Neural Network For  
Image Generation,  
Gregor et al.



# Recurrent Neural Network

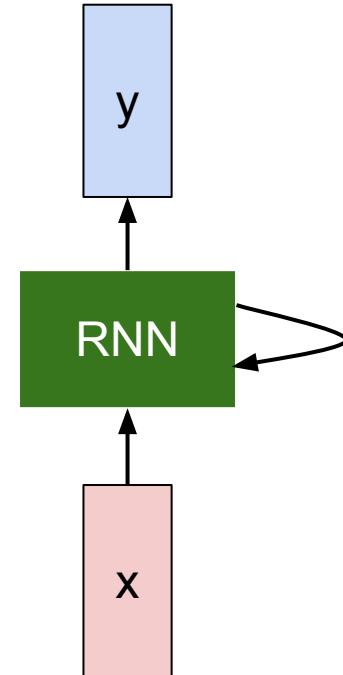


# Recurrent Neural Network



# Recurrent Neural Network

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

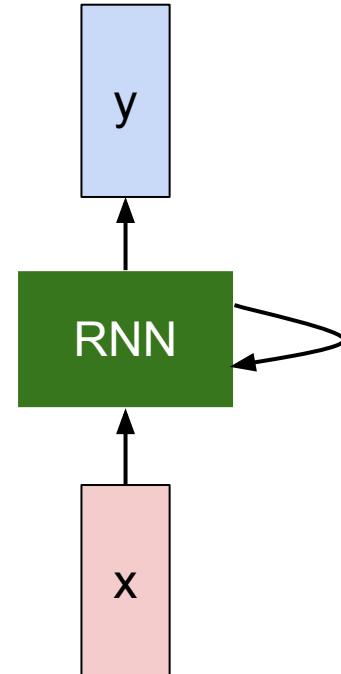


# Recurrent Neural Network

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

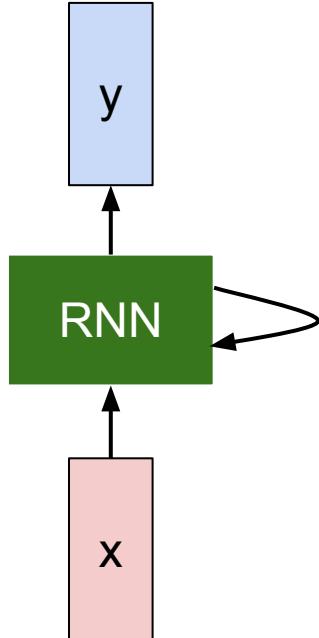
$$h_t = f_W(h_{t-1}, x_t)$$

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



# (Vanilla) Recurrent Neural Network

The state consists of a single “*hidden*” vector  $\mathbf{h}$ :



$$\mathbf{h}_t = f_W(\mathbf{h}_{t-1}, \mathbf{x}_t)$$

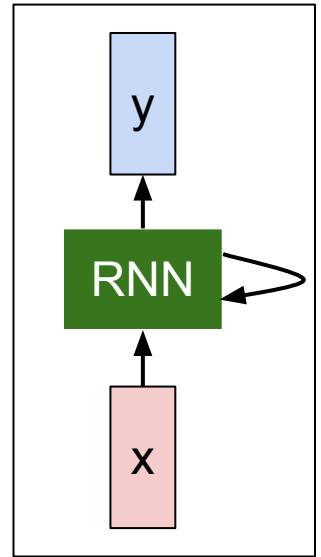
$$\mathbf{h}_t = \tanh(W_{hh}\mathbf{h}_{t-1} + W_{xh}\mathbf{x}_t)$$

$$y_t = W_{hy}\mathbf{h}_t$$

# Character-level language model example

Vocabulary:  
[h,e,l,o]

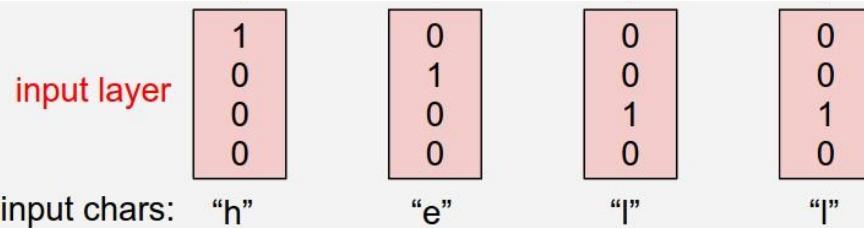
Example training  
sequence:  
“hello”



# Character-level language model example

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
“hello”

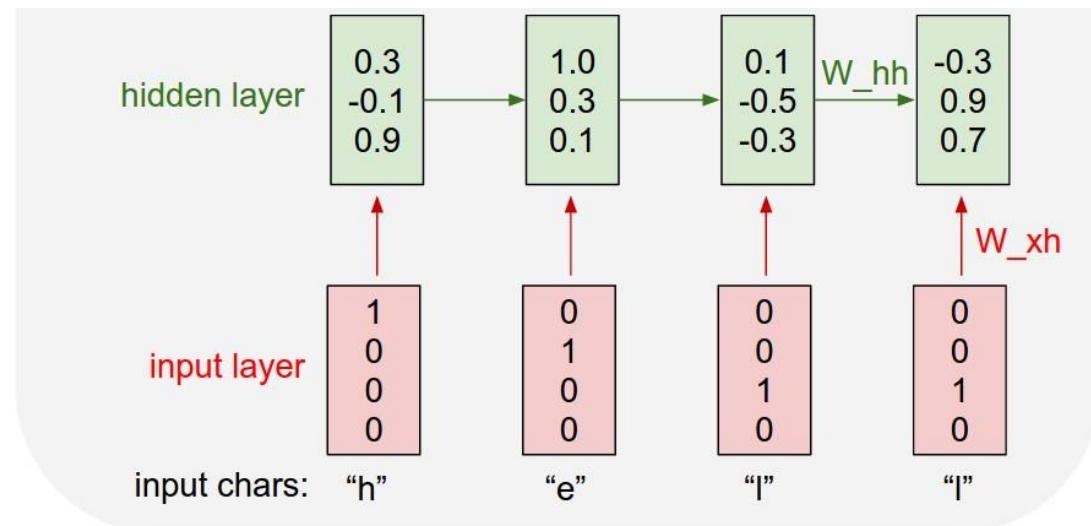


# Character-level language model example

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
“hello”

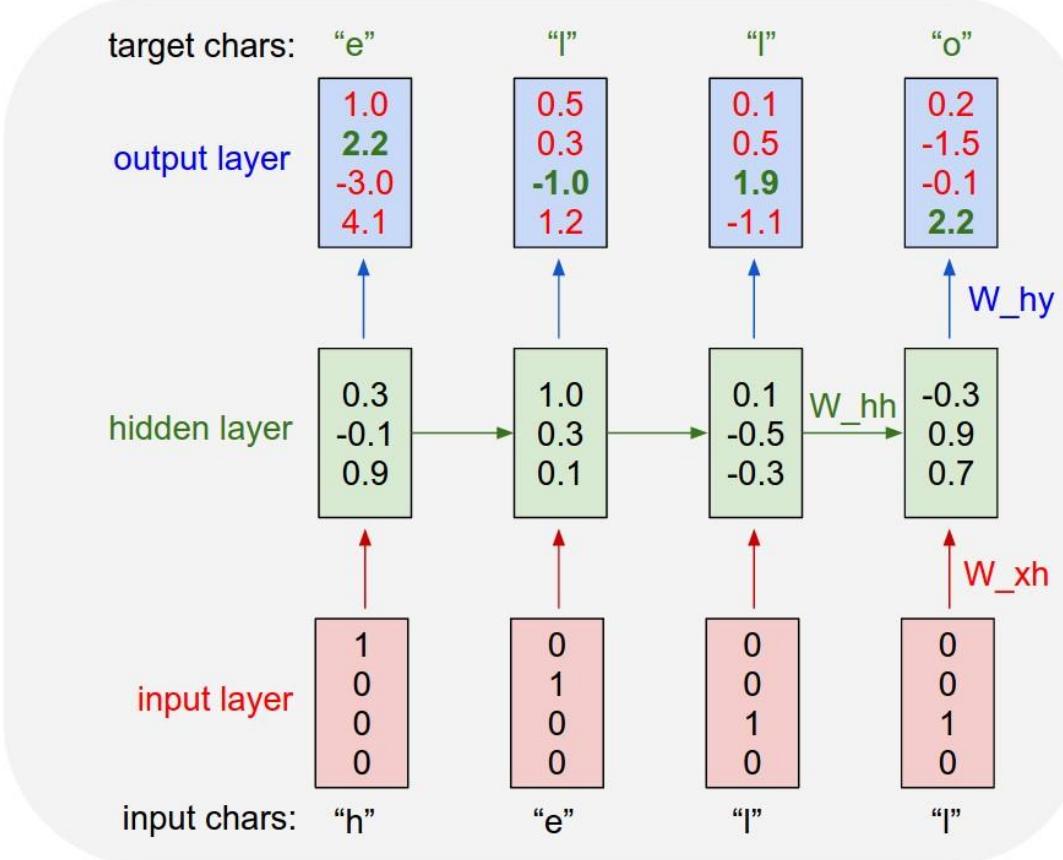
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$



# Character-level language model example

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
“hello”



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

```
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 # 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 wkh = 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     # 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(wkh, 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         dwhx, dwhh, dwhy = np.zeros_like(wkh), np.zeros_like(whh), np.zeros_like(why)
47         dbh, dby = np.zeros_like(bh), np.zeros_like(by)
48         dhnext = np.zeros_like(hs[0])
49         for t2 in reversed(xrange(len(inputs))):
50             dy = np.copy(ps[t2])
51             dy[targets[t2]] -= 1 # backprop into y
52             dyw = np.dot(dy, hs[t2].T)
53             dh = np.dot(why.T, dy) + dhnext # backprop into h
54             ddraw = (i - hs[t2].T) * dh # backprop through tanh nonlinearity
55             ddbh = ddraw
56             dwhh += np.dot(ddraw, xs[t].T)
57             dwhy += np.dot(ddraw, hs[t-1].T)
58             dhnext = np.dot(whh.T, ddraw)
59             for dparam in [dwhx, dwhh, dwhy, dbh, dby]:
60                 np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61
62     return loss, dwhx, dwhh, dwhy, dbh, dby, hs[len(inputs)-1]
```

```
63 def sample(h, 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     ixes = []
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
78         ixes.append(ix)
79
80     return ixes
81
82 n, p = 0, 0
83 mxwh, mwhh, mmhy = np.zeros_like(wkh), np.zeros_like(whh), np.zeros_like(why)
84 mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
85 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
86 while True:
87     # prepare inputs (we're sweeping from left to right in steps seq_length long)
88     if p+seq_length >= len(data) or n == 0:
89         hprev = np.zeros((hidden_size,1)) # reset RNN memory
90         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\n----' % (txt,))
99
100     # forward seq_length characters through the net and fetch gradient
101     loss, dwhx, dwhh, dwhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
102     smooth_loss = smooth_loss * .999 + loss * .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([wkh, whh, why, bh, by],
107                                   [dwhx, dwhh, dwhy, dbh, dby],
108                                   [mxwh, mwhh, mmhy, mbh, mby]):
109         mem += dparam * dparam
110         param += -learning_rate * param / np.sqrt(mem + 1e-8) # adagrad update
111
112     p += seq_length # move data pointer
113     n += 1 # iteration counter
```

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

# Data I/O

## 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: %s' % (data_size, vocab_size, ''.join(chars)))
12
13 # Parameters
14 hidden_size = 100 # HIDDEN_SIZE
15 input_size = vocab_size
16
17 # Hyperparameters
18 K=1000000 # Number of steps to train
19 L=100 # Loss threshold
```

```
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) }
```

# min-char-rnn.py gist

```

***  

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

BSD license  

http://karpathy.github.io/char-rnn/ (MIT license)  

import numpy as np  
  

# Data I/O  

data = open('input.txt', 'r').read() # should be simple plain text file  

data_size, vocab_size = len(data), len(chars)  

print 'data has %d characters, %d unique.' % (data_size, vocab_size)  

chars = sorted(list(set(chars)))  

char_to_ix = {ch:i for i, ch in enumerate(chars)}  
  

# Hyperparameters  

hidden_size = 100 # size of hidden layer of neurons  

n_in = vocab_size # input layer size  

n_out = vocab_size # output layer size  

n_steps = 25 # number of steps to unroll the RNN for  

learning_rate = 1e-1  
  

# Model parameters  

Wxh = np.random.randn(n_in, hidden_size) # input to hidden  

whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden  

why = np.random.randn(hidden_size, n_out)*0.01 # hidden to output  

bh = np.zeros((hidden_size,)) # hidden bias  

bo = np.zeros((n_out,)) # output bias  
  

# Training parameters  

n_hiddens = 100  

n_in = vocab_size  

n_out = vocab_size  

n_steps = 25  

batch_size = 100  

n_iters = 1000000  

learning_rate = 1e-1  
  

# Misc  

w_in = np.zeros((n_in, hidden_size)) # input to hidden  

w_hh = np.zeros((hidden_size, hidden_size)) # hidden to hidden  

w_out = np.zeros((hidden_size, n_out)) # hidden to output  

b_h = np.zeros(hidden_size) # hidden bias  

b_o = np.zeros(n_out) # output bias  
  

# Loss and derivatives  

def lossFun(inputs, targets, hprev):  
    """  
    inputs, targets are both lists of integers.  
    hprev is the array of initial hidden state  
    returns the loss, gradients on model parameters, and last hidden state  
    """  
  
    n_in, x, y, pe = 0, [], [], []  
    hprev = np.copy(hprev)  
    loss = 0  
  
    # Forward pass  
    for t in range(len(inputs)): # loop over time steps  
        x[t] = np.zeros(vocab_size, 1) # encode in 1-of-k representation  
        x[t][inputs[t]] = 1  
        w_in_t = w_in[inputs[t], :]  
        hprev_t = hprev if t == 0 else hprev  
        whh_t = whh if t == 0 else np.dot(w_hh, hprev_t[-1])  
        bh_t = bh if t == 0 else np.zeros_like(bh)  
        y[t] = np.dot(w_out, hprev_t) + np.sum(np.exp(y[t]))  
        why_t = np.exp(y[t]) / np.sum(np.exp(y[t])) * probabilities for next chars  
        pe[t] = -np.log(why_t[targets[t]])  
        loss += -pe[t] # softmax loss per step  
        dhprev = np.dot(w_hh.T, dhyy) + np.dot(w_in.T, dy) # backprop through tanh nonlinearity  
        dwhh = np.dot(dy, dhprev) # backprop through y  
        dbh = np.sum(dy, axis=0)  
        dw_in = np.outer(dy, x[t]) # backprop into x  
        dhraw = (1 - np.tanh(hprev_t)**2) * dhyy # dh raw  
        dhraw += np.dot(dwhh, hprev_t[-1])  
        dhraw += np.dot(dbh, np.ones_like(hprev_t[-1]))  
        dhraw = np.tanh(dhraw) # clip to mitigate exploding gradients  
        for i in range(n_in):  
            if x[t][i] == 1:  
                dw_in[i] = dhraw  
        dw_in = np.clip(dw_in, -5, 5, out=dw_in) # clip to mitigate exploding gradients  
        return loss, dhraw, dwhh, dbh, dy, dhraw[inputs[t]-1]  
    def sample(hprev, seed_ix, n):  
        """  
        sample a sequence of integers from the model  
        hprev is the previous hidden state, seed is letter for first time step  
        """  
  
        x = np.zeros((n, hidden_size, 1))  
        x[seed_ix] = 1  
        iids = []  
        for t in range(n):  
            h = np.tanh(np.dot(w_in, x) + np.dot(w_hh, h) + bh)  
            y = np.exp(np.dot(w_out, h))  
            ix = np.argmax(np.exp(y))  
            x = np.zeros(vocab_size, 1)  
            x[ix] = 1  
            iids.append(ix)  
        return iids  
  
    h = np.zeros((n_in, hidden_size))  
    math, mathy, mathy = np.zeros_like(bh), np.zeros_like(bh), np.zeros_like(bh)  
    mathy = np.zeros_like(bh), np.zeros_like(bh), np.zeros_like(bh)  
    memory_variables for Adagrad  
    for i in range(n_iters):  
        if i % 1000 == 0:  
            print 'step %d' % i  
        if i % 100 == 0:  
            print 'Loss at step %d is %f' % (i, loss)  
  
        inputs = [char_to_ix[ch] for ch in data[:n_in]]  
        targets = [char_to_ix[ch] for ch in data[n_in:n_in+n_steps]]  
  
        # forward pass  
        loss, dhraw, dwhh, dbh, dy, dhyy = lossFun(inputs, targets, h)  
        loss += 0.5*(dwhh**2) + (dbh**2) # regularization term  
        if i % 100 == 0:  
            print 'step %d loss: %f' % (i, loss), print progress  
  
        # perform parameter update with Adagrad  
        for para, paraee in zip([w_in, whh, bh, w_out], [math, mathy, mathy, mathy]):  
            mem = np.sqrt(np.sum(paraee**2)) + 1e-8 # added update  
            paraee += -learning_rate * paraee / mem  
            para += paraee  
  
        x = np.zeros((n_in, hidden_size, 1))  
        i = 1 # iteration counter

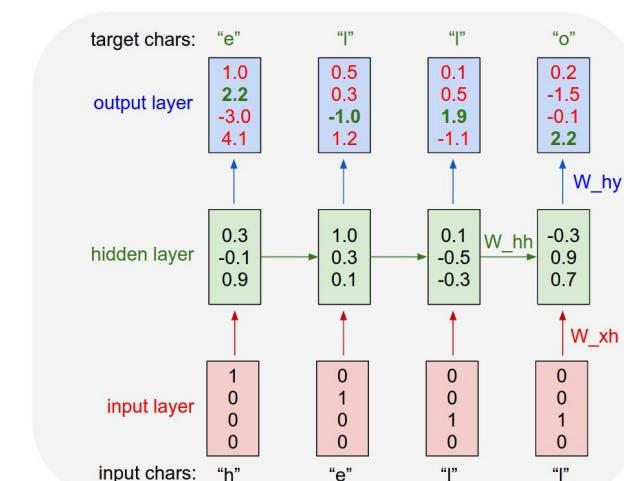
```

## 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 
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
```

## recall



# min-char-rnn.py gist

```
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 vocab_size = len(chars)
11 data_size = len(data)
12 print("data has %d characters, %d unique." % (data_size, vocab_size))
13 char_to_ix = {ch:i for i in range(len(chars))}
14 ix_to_char = {i:ch for ch in range(len(chars))}
15 ix_to_char = {i:ch for ch in range(len(chars)) if ch >= 32}
16 ix_to_char = {i:ch for ch in range(len(chars)) if ch <= 126}
17
18 # hyperparameters
19 hidden_size = 100 # size of hidden layer of neurons
20 seq_length = 20 # number of steps to unroll the RNN for
21 learning_rate = 1e-1
22
23 model_params = {}
24
25 def init_random_params():
26     wh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
27     bh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
28     why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
29     by = np.zeros(vocab_size, 1) # hidden bias
30     bi = np.zeros(vocab_size, 1) # output bias
31
32     return wh, bh, why, by, bi, ix_to_char, char_to_ix
33
34 def lossFun(inputs, targets, hprev):
35
36     inputs,targets = both_list_of_integers(
37         inputs,targets)
38
39     hprev = hM.array of initial hidden state
40     return the loss, gradients on model parameters, and last hidden state
41
42     xs, hs, ys, ps = 0, 0, 0, 0
43     hM = hM.copy(hprev)
44     loss = 0
45
46     for t in range(seq_length):
47         x = np.zeros((vocab_size, 1)) # encode in 1-of-k representation
48         x[inputs[t], 0] = 1 # x[t] is one-hot encoding of the char we want
49
50         hM1 = np.tanh(np.dot(wh, x) + np.dot(bh, hs[-1]) + bh) # hidden state
51         ps1 = np.exp(why * hM1) / np.sum(np.exp(why * hM1)) # probabilities for next chars
52         loss += -np.log(ps1[targets[t], 0]) # softmax (cross-entropy loss)
53
54         # backprop through tanh nonlinearity
55         dwhy = np.zeros_like(why)
56         dbh = np.zeros_like(bh)
57         dwh = np.zeros_like(wh)
58         dx = np.zeros_like(x)
59
60         for dparam in [dwhy, dbh, dwh]:
61             dparam *= 0
62
63         dy = np.copy(ps1)
64         dy[targets[t]] -= 1 # backprop into y
65         dh = dy.dot(why.T)
66         dwhy += dy[:, None].dot(hs[-1])
67         dbh += dh[None, :]
68         dwh += np.dot(dy.T, x)
69
70         for dparam in [dwhy, dbh, dwh, dh]:
71             np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
72
73         dh = np.dot(why, dh) + dwh # backprop into h
74         dh += dh * np.tanh(dh) * (1 - dh**2) # backprop through tanh nonlinearity
75         dbh += dh * dwh
76         dwh += dh * dwh
77         dwhy += dh * dy
78
79         for t in range(seq_length):
80             h = np.tanh(dh, dwh, dby)
81             hM = hM.copy(h)
82             dby = np.zeros_like(by)
83             dy = np.exp(why * h)
84             dx = np.random.choice(range(vocab_size), p=dy.ravel())
85             x[dx, 0] = 1
86             x[inputs[t], 0] = 0
87
88         return loss
89
90
91     sample_ix = sample(hprev, inputs[0], 200)
92     txt = ''.join(ix_to_char[ix] for ix in sample_ix)
93     print '----\n %s \n----' % (txt, )
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
```



# Main loop

```
n, p = 0, 0
mWxh, mWhh, mWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
mbh, mby = 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:
    # 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]]

    # 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],
                                  [mWxh, mWhh, mWhy, mbh, mby]):
        mem += dparam * dparam
        param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update

    p += seq_length # move data pointer
    n += 1 # iteration counter
```

# min-char-rnn.py gist

```
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 vocab_size = len(chars)
11 data_size, vocab_size = len(data), len(chars)
12 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
13 char_to_ix = {ch:i for i in range(len(chars))}
14 ix_to_char = {i:ch for ch in range(len(chars))} 
15
16 # Hyperparameters
17 hidden_size = 100 # size of hidden layer of neurons
18 seq_length = 20 # number of steps to unroll the RNN for
19 learning_rate = 1e-1
20
21 # Model parameters
22 dWxh, dWhh, dWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
23 dbh, mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
24 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
25
26 while True:
27     if p+seq_length+1 >= len(data) or n == 0:
28         hprev = np.zeros((hidden_size,1)) # reset RNN memory
29         p = 0 # go from start of data
30         inputs, targets = [ord(c) for c in data[p:p+seq_length]], [c for c in data[p+1:p+seq_length+1]]
31
32     x, hs, ys, ps = O, O, O, O
33     hprev = np.copy(hprev)
34     loss, dWxh, dWhh, dWhy = 0.0, 0.0, 0.0, 0.0
35
36     # Forward pass
37     for t in range(seq_length):
38         x = np.zeros((vocab_size,1)) # encode in 1-of-k representation
39         x[ord(inputs[t]):] = 1
40         hprev = np.tanh(np.dot(wxh, x) + np.dot(whh, hs[-1]) + bh) # hidden state
41         y = np.exp(hprev) # probabilities for next chars
42         ps = np.log(y[targets[t]] + np.spacing(0.0)) # softmax (cross-entropy loss)
43         loss += -ps # cross-entropy loss
44         dy = np.copy(ps)
45         dy[targets[t]] -= 1 # backprop into y
46         dh = np.dot(dy, dotw) # backprop into h
47         dh += np.dot(dy, (dotw.T) * hprev) # dh = backprop through tanh nonlinearity
48         dh += dhw # dhw = np.dot(dh, whh.T)
49         dWxh += np.dot(x.T, x) * t
50         dWhh += np.dot(hs[-1].T, dh) * t
51         dbh += np.sum(dh, axis=0, keepdims=True) * t
52         dWxh, dWhh, dbh = clip(dWxh, -5, 5, out=dWxh), clip(dWhh, -5, 5, out=dWhh), clip(dbh, -5, 5, out=dbh)
53         if np.isnan(dWxh).any() or np.isnan(dWhh).any() or np.isnan(dbh).any():
54             raise ValueError("Nan error")
55
56     if n % 100 == 0:
57         sample_ix = sample(hprev, inputs[0], 200)
58         txt = ''.join(ix_to_char[i] for i in sample_ix)
59         print '----\n %s ----' % (txt, )
60
61     # Forward seq_length characters through the net and fetch gradient
62     loss, dWxh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
63     smooth_loss = smooth_loss * 0.999 + loss * 0.001
64     if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
65
66     # Perform parameter update with Adagrad
67     for param, dparam, mem in zip([Wxh, Whh, Why, bh, by],
68                                 [dWxh, dWhh, dWhy, dbh, dby],
69                                 [mWxh, mWhh, mWhy, mbh, mby]):
70         mem += dparam * dparam
71         param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
72
73     p += seq_length # move data pointer
74     n += 1 # iteration counter
75
```



# Main loop

```
81 n, p = 0, 0
82 mWxh, mWhh, mWhy = np.zeros_like(Wxh), 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
86 while True:
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[i] for i 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     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([Wxh, 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
```

# min-char-rnn.py gist

# Main loop

```
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(data[1:(1)])
10 vocab_size = len(chars), len(data), len(chars)
11 print("data has %d characters, %d unique." % (data_size, vocab_size))
12 char_to_ix = {ch:i for i in range(len(chars))}
13 ix_to_char = {i:ch for ch in range(len(chars))} # 14
14
15 # hyperparameters
16 hidden_size = 100 # size of hidden layer of neurons
17 seq_length = 20 # number of steps to unroll the RNN for
18 learning_rate = 1e-1
19
20 # model parameters
21 dWxh, dWhh, dWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
22 dbh, mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
23 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
24
25 while True:
26     # prepare inputs (we're sweeping from left to right in steps seq_length long)
27     if p+seq_length+1 >= len(data) or n == 0:
28         hprev = np.zeros((hidden_size,1)) # reset RNN memory
29         p = 0 # go from start of data
30
31     inputs, targets = data[p:p+seq_length], data[p+1:p+seq_length+1]
32
33     xs, hs, ys, ps = [0, 0, 0, 0]
34     hprev = np.copy(hprev)
35     loss = 0
36
37     for t in range(seq_length):
38         # forward pass
39         x = np.zeros((vocab_size, 1)) # encode in 1-of-k representation
40         x[inputs[t]] = 1
41         hprev = np.tanh(np.dot(wh, hprev) + np.dot(bh, hprev[-1]) + bh)
42         y = np.exp(hprev)
43         ps = np.exp(y)/np.sum(np.exp(y)) # probabilities for next chars
44         loss += -np.log(ps[targets[t]]) # softmax (cross-entropy loss)
45
46         # backward pass: compute gradients to params
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         dnext = np.zeros_like(hprev)
50
51         for i in range(len(inputs)):
52             dy = np.copy(ps[i])
53             dy[targets[t]] -= 1 # backprop into y
54             dh = np.dot(Wxh.T, dy) + dnext # backprop into h
55             dh += dh * np.tanh(hprev) * (1-hprev**2) # backprop through tanh nonlinearity
56             dWxh += np.dot(dy, hprev.T)
57             dWhh += np.dot(dy, hprev[-1].T)
58             dbh += np.sum(dy, axis=0, keepdims=True)
59             dnext = np.dot(Whh, dy)
60
61         for param in [dWxh, dWhh, dbh, dby]:
62             np.clip(param, -5, 5, out=param) # clip to mitigate exploding gradients
63
64         if np.isnan(dWxh).any() or np.isnan(dWhh).any() or np.isnan(dbh).any() or np.isnan(dby).any():
65             raise ValueError("NaN in gradients")
66
67         # sample from the model now and then
68         if n % 100 == 0:
69             sample_ix = sample(hprev, inputs[0], 200)
70             txt = ''.join(ix_to_char[ix] for ix in sample_ix)
71             print('----\n%s\n----' % (txt, ))
72
73     # forward seq_length characters through the net and fetch gradient
74     loss, dWxh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
75     smooth_loss = smooth_loss * 0.999 + loss * 0.001
76
77     if n % 100 == 0: print('iter %d, loss: %f' % (n, smooth_loss)) # print progress
78
79     # perform parameter update with Adagrad
80     for param, dparam, mem in zip([Wxh, Whh, Why, bh, by],
81                                 [dWxh, dWhh, dWhy, dbh, dby],
82                                 [mWxh, mWhh, mWhy, mbh, mby]):
83         mem += dparam * dparam
84         param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
85
86     p += seq_length # move data pointer
87     n += 1 # iteration counter
```

# min-char-rnn.py gist

```
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 vocab_size = len(chars)
11 data_size, vocab_size = len(data), len(chars)
12 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
13 char_to_ix = {ch:i for i in range(len(chars))}
14 ix_to_char = {i:ch for ch in range(len(chars))} # 14
15
16 # Hyperparameters
17 hidden_size = 100 # size of hidden layer of neurons
18 seq_length = 20 # number of steps to unroll the RNN for
19 learning_rate = 1e-1
20
21 # Model parameters
22 dWxh, dWhh, dWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
23 dbh, mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
24 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
25
26 while True:
27     # Prepare inputs (we're sweeping from left to right in steps seq_length long)
28     if p+seq_length+1 >= len(data) or n == 0:
29         hprev = np.zeros((hidden_size,1)) # reset RNN memory
30         p = 0 # go from start of data
31         inputs = [c for c in data[p:p+seq_length]] # 31
32         targets = [c for c in data[p+1:p+seq_length+1]]
33
34     x, hs, ys, ps = o, O, O, O
35     hprev = np.copy(hprev)
36     loss = 0
37
38     for t in range(seq_length):
39         x = np.zeros((vocab_size,1)) # encode in 1-of-k representation
40         x[0][char_to_ix[inputs[t]]] = 1
41         hprev = np.tanh(np.dot(wh, hprev) + np.dot(bh, hprev[-1]) + bh) # hidden state
42         y = np.dot(Wxh, hprev) + bh # output neuron
43         yprob = np.exp(y)/np.sum(np.exp(y)) # probabilities for next chars
44         loss += -np.log(yprob[targets[t]]) # softmax (cross-entropy loss)
45
46         dy = np.copy(ps) # 46
47         dy[targets[t]] -= 1 # backprop into y
48         dh = np.dot(dy, Wxh.T) # 48
49         dh += np.dot(dy, bh.T) # dh = backprop through tanh nonlinearity
50         dh += dh * (1-hprev**2) # 50
51         dWxh += np.dot(inputs[t].T, dy) # 51
52         dbh += np.sum(dy, axis=0) # 52
53         dWhh += np.dot(hprev.T, dy) # 53
54         dWhy += np.sum(dy*x, axis=0) # 54
55
56         dx = np.dot(dy, Wxh) # 56
57         dh = np.tanh(dx) # 57
58         dh *= (1-dh**2) # 58
59
60         for dparam in [dWxh, dbh, dWhh, dWhy]:
61             np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
62             dparam += -learning_rate * dparam # 62
63             dbh, dWhh, dWhy = np.zeros_like(dbh, dWhh, dWhy, dbh[0])
64
65         sample_ix = np.random.choice(range(vocab_size), 1) # 65
66         seed_ix = sample_ix[0]
67
68     # Sample from the model now and then
69     h = np.dot(Wxh, x) + np.dot(wh, hprev) # 69
70     h[0] = np.zeros((hidden_size,1)) # 70
71     x = np.zeros((vocab_size,1)) # 71
72     x[seed_ix] = 1 # 72
73
74     for t in range(seq_length):
75         x = np.tanh(np.dot(wh, h) + bh) # 75
76         y = np.dot(Wxh, h) + bh # 76
77         yprob = np.exp(y)/np.sum(np.exp(y)) # 77
78         ix = np.random.choice(range(vocab_size), 1, p=yprob.ravel()) # 78
79         x[0] = np.zeros((vocab_size,1)) # 79
80         x[ix] = 1 # 80
81
82     return loss
83
84
85 n, p = 0, 0
86 mWxh, mWhh, mWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
87 dbh, mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
88 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
89
90 while True:
91     # Prepare inputs (we're sweeping from left to right in steps seq_length long)
92     if p+seq_length+1 >= len(data) or n == 0:
93         hprev = np.zeros((hidden_size,1)) # reset RNN memory
94         p = 0 # go from start of data
95         inputs = [c for c in data[p:p+seq_length]] # 95
96         targets = [c for c in data[p+1:p+seq_length+1]]
97
98
99     # Sample from the model now and then
100    h = np.dot(Wxh, x) + np.dot(wh, hprev) # 100
101    h[0] = np.zeros((hidden_size,1)) # 102
102    x = np.zeros((vocab_size,1)) # 103
103    x[seed_ix] = 1 # 104
104
105    for t in range(seq_length):
106        x = np.tanh(np.dot(wh, h) + bh) # 106
107        y = np.dot(Wxh, h) + bh # 107
108        yprob = np.exp(y)/np.sum(np.exp(y)) # 108
109        ix = np.random.choice(range(vocab_size), 1, p=yprob.ravel()) # 109
110        x[0] = np.zeros((vocab_size,1)) # 110
111        x[ix] = 1 # 112
112
113
114    # Compute gradients, but don't backprop through the first time step
115    loss, dWxh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
116    smooth_loss = smooth_loss * 0.999 + loss * 0.001
117
118    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
119
120
121    # Perform parameter update with Adagrad
122    for param, dparam, mem in zip([Wxh, Whh, Why, bh, by],
123                                 [dWxh, dWhh, dWhy, dbh, dby],
124                                 [mWxh, mWhh, mWhy, mbh, mby]):
125        mem += dparam * dparam
126        param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
127
128    p += seq_length # move data pointer
129    n += 1 # iteration counter
130
```

# Main loop



# min-char-rnn.py gist

```
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(data[1:(1)])
10 vocab_size = len(chars), len(data), len(chars)
11 print("data has %d characters, %d unique." % (data_size, vocab_size))
12 char_to_ix = {c: i for i in range(len(chars))}
13 ix_to_char = {i: c for i in range(len(chars))} 
14
15 # hyperparameters
16 hidden_size = 100 # size of hidden layer of neurons
17 seq_length = 20 # number of steps to unroll the RNN for
18 learning_rate = 1e-1
19
20 # Model parameters
21 dWxh, dWhh, dWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
22 dbh, mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
23 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
24
25 while True:
26     # prepare inputs (we're sweeping from left to right in steps seq_length long)
27     if p+seq_length+1 >= len(data) or n == 0:
28         hprev = np.zeros((hidden_size,1)) # reset RNN memory
29         p = 0 # go from start of data
30         inputs, targets = [char_to_ix[ch] for ch in data[p:p+seq_length]]
31         hprev = np.copy(hprev)
32
33     xs, hs, ys, ps = np.zeros((O, O, O, O))
34     h1_t1 = np.copy(hprev)
35     loss = 0
36
37     # forward pass
38     for t in range(seq_length):
39         # encode in 1-of-k representation
40         x1_t1 = np.zeros((vocab_size,1)) + encode_in_1-of-k_representation(chars[t])
41         h1_t1 = np.tanh(np.dot(wh, x1_t1) + np.dot(bh, h1_t1) + bh) # hidden state
42         ps1_t1 = np.exp(ps1_t1) / np.sum(np.exp(ps1_t1)) = softmax # probabilities for next chars
43         loss += -np.exp(ps1_t1)[targets[t]] / np.sum(np.exp(ps1_t1)) = softmax # (cross-entropy loss)
44
45         # backward pass: compute gradients
46         dWxh, dWhh, dWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
47         dbh, dby = np.zeros_like(bh), np.zeros_like(by)
48         dnext = np.zeros_like(h1_t1)
49         for i in range(len(inputs)):
50             dy = np.copy(ps1_t1)
51             dy[targets[t]] -= 1 # backprop into y
52             dh = np.dot(dy, dotWxh, h1_t1.T)
53             dWxh += np.dot(dy, h1_t1) # dh = backprop into Wxh through tanh nonlinearity
54             dh += dh * np.dot(dy, dotWhh, h1_t1.T) # dh = backprop through tanh nonlinearity
55             dWhh += np.dot(dy, dotWhh, h1_t1.T)
56             dbh += np.sum(dy, axis=0, keepdims=True) * h1_t1.T
57             dnext = np.dot(dy, dotWhy, h1_t1)
58
59         for dparam in [dWxh, dWhh, dbh, dby]:
60             np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61             dparam *= learning_rate # update rule, da/dtheta * learning_rate
62             dparam -= mem * dparam # adam update
63             mem += dparam * dparam
64
65     # sample from the model now and then
66     if n % 100 == 0:
67         sample_ix = sample(hprev, inputs[0], 200)
68         txt = ''.join(ix_to_char[ix] for ix in sample_ix)
69         print '----\n%s\n----' % (txt, )
70
71     # forward seq_length characters through the net and fetch gradient
72     loss, dWxh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
73     smooth_loss = smooth_loss * 0.999 + loss * 0.001
74
75     if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
76
77     # perform parameter update with Adagrad
78     for param, dparam, mem in zip([Wxh, Whh, Why, bh, by],
79                                 [dWxh, dWhh, dWhy, dbh, dby],
80                                 [mWxh, mWhh, mWhy, mbh, mby]):
81         mem += dparam * dparam
82         param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
83
84         p += seq_length # move data pointer
85         n += 1 # iteration counter
```



# Main loop

```
81     n, p = 0, 0
82     mWxh, mWhh, mWhy = np.zeros_like(Wxh), 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
86     while True:
87
88         # prepare inputs (we're sweeping from left to right in steps seq_length long)
89         if p+seq_length+1 >= len(data) or n == 0:
90             hprev = np.zeros((hidden_size,1)) # reset RNN memory
91             p = 0 # go from start of data
92
93             inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
94             targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
95
96
97         # sample from the model now and then
98         if n % 100 == 0:
99             sample_ix = sample(hprev, inputs[0], 200)
100            txt = ''.join(ix_to_char[ix] for ix in sample_ix)
101            print '----\n%s\n----' % (txt, )
102
103
104         # forward seq_length characters through the net and fetch gradient
105         loss, dWxh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
106         smooth_loss = smooth_loss * 0.999 + loss * 0.001
107
108         if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
109
110
111         # perform parameter update with Adagrad
112         for param, dparam, mem in zip([Wxh, Whh, Why, bh, by],
113                                     [dWxh, dWhh, dWhy, dbh, dby],
114                                     [mWxh, mWhh, mWhy, mbh, mby]):
115
116             mem += dparam * dparam
117             param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
118
119             p += seq_length # move data pointer
120             n += 1 # iteration counter
```

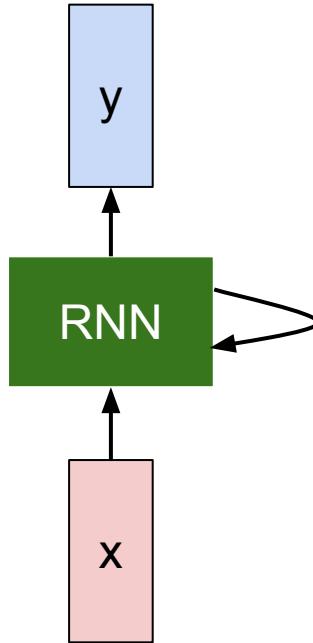
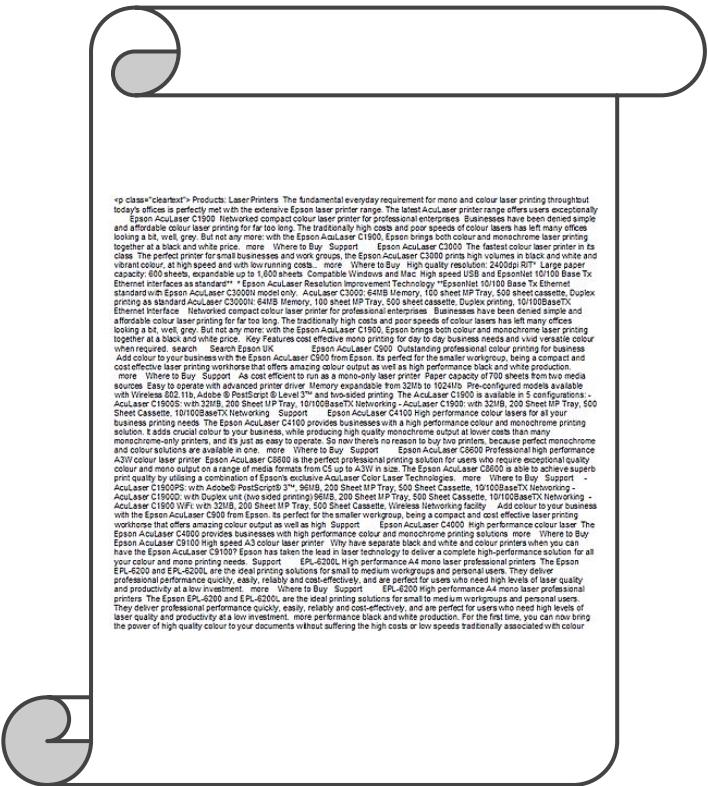
# min-char-rnn.py gist

```
***  
Minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)  
BSD License  
***  
Import numpy as np  
  
# Data I/O  
data = open('train.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 ch in enumerate(chars)}  
  
# Hyperparameters  
hidden_size = 100 # size of hidden layer of neurons  
seq_length = 20 # number of steps to unroll the RNN for  
learning_rate = 1e-1  
  
# Model parameters  
wh = np.random.rand(hidden_size, vocab_size)*0.01 # input to hidden  
bh = np.random.rand(hidden_size, hidden_size)*0.01 # hidden to hidden  
why = np.random.rand(vocab_size, hidden_size)*0.01 # hidden to output  
bh0 = 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  
    for t in xrange(len(inputs)):  
        xs[t] = np.zeros((vocab_size, 1)) # encode in 1-of-k representation  
        xs[t][inputs[t]] = 1  
        wht = np.tanh(np.dot(wh, xs[t]) + np.dot(bh, hs[t-1]) + bh) # hidden state  
        whyt = np.dot(why, wht) + by # unnormalized log probabilities for next chars  
        ps[t] = np.exp(whyt) / np.sum(np.exp(whyt)) # softmax (cross-entropy loss)  
        loss += -np.log(ps[t][targets[t], 0])  
        dy = np.copy(ps[t])  
        dy[targets[t]] -= 1 # backprop into y  
        dhnext = np.zeros_like(hs[t])  
        dy *= dy  
        dhy = np.dot(dy, wh.T)  
        dh = np.dot(dhy, bh) + dhnext # backprop through tanh nonlinearity  
        dbh += np.sum(dhy, axis=0)  
        dwh = np.dot(dy, xs[t].T)  
        dhyt = np.dot(why, dy) + by  
        dhyt += np.sum(dhyt, axis=0)  
        dhyt -= np.sum(dhyt) * np.ones_like(dhyt)  
        dhyt *= dy  
        for dparam in [dwh, dbh, dhy]:  
            np.clip(dpараметre, -5, 5, out=dparam) # clip to mitigate exploding gradients  
        dhnext = dhyt  
    return loss, xs, hs, ys, ps  
  
def sample(hx, ix):  
    """  
    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  
  
    for t in xrange(n):  
        h = np.tanh(np.dot(wh, x) + np.dot(bh, h) + bh)  
        p = np.exp(why * np.dot(h, wh))  
        ix = np.random.choice(range(vocab_size), p=p.ravel())  
        x = np.zeros((vocab_size, 1))  
        x[ix] = 1  
    return ix  
  
n, p, t, b  
meh, moh, mhy = np.zeros_like(by), np.zeros_like(wh), np.zeros_like(why)  
meh0, moh0, mhy0 = np.zeros_like(bh), np.zeros_like(why), np.zeros_like(why)  
smooth_loss = -np.inf  
for i in range(seq_length):  
    if i > 0:  
        print 'resetting from left to right in step', i, 'iteration', t  
    if p+seq_length > len(data) or t == 0:  
        np.zeros((hidden_size, 1)) # reset RNN memory  
    inputs = [char_to_ix[ch] for ch in data[p:seq_length]]  
    targets = [char_to_ix[ch] for ch in data[p+seq_length:]]  
  
    # forward pass: compute scores, through the net and fetch gradient  
    loss, dxh, dmbh, dmy, dby, dprev = lossFun(inputs, targets, hprev)  
    smooth_loss = smooth_loss * 0.999 + loss * 0.001  
    t += 1  
    if t % 100 == 0:  
        print 'iter %d, loss: %f' % (t, smooth_loss) # print progress  
  
    # backward pass: compute gradients going backwards  
    dwhxh, dwhh, dwhy = np.zeros_like(whxh), np.zeros_like(whhh), 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  
        dwhy += np.dot(dy, hs[t].T)  
        dbh += dy  
        dh = np.dot(why.T, dy) + dhnext # backprop into h  
        ddraw = (1 - hs[t] * hs[t]) * dh # backprop through tanh nonlinearity  
        dbh += ddraw  
        dwhxh += np.dot(ddraw, xs[t].T)  
        dwhh += np.dot(ddraw, hs[t-1].T)  
        dhnext = np.dot(why.T, ddraw)  
  
    for dpараметре in [dwhxh, dwhh, dwhy, dbh, dby]:  
        np.clip(dpараметре, -5, 5, out=dpараметре) # clip to mitigate exploding gradients  
  
    return loss, dwhxh, dwhh, dwhy, dbh, dby, hs[len(inputs)-1]
```

# Loss function

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





## 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-iafhatawiaoihrdemot 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 shuwy 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.

 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 agolatest commit 4b1786927d  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/hex2fw.c: restore missing default in switch statement

2 months ago

 fs

vfs: read file\_handle only once in handle\_to\_path

4 days ago

 include

Merge branch 'perl-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

 io

bio: bio\_start: New bio\_start and multi-block bio\_start methods introduced

a month ago

 Code Pull requests  
74 Pulse Graphs

HTTPS clone URL

<https://github.com/torvalds/linux> You can clone with [HTTPS](#),  
[SSH](#), or [Subversion](#).  Clone in Desktop Download ZIP

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

# Generated C code

```
/*
 * 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)
{
#endif 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;
}
```

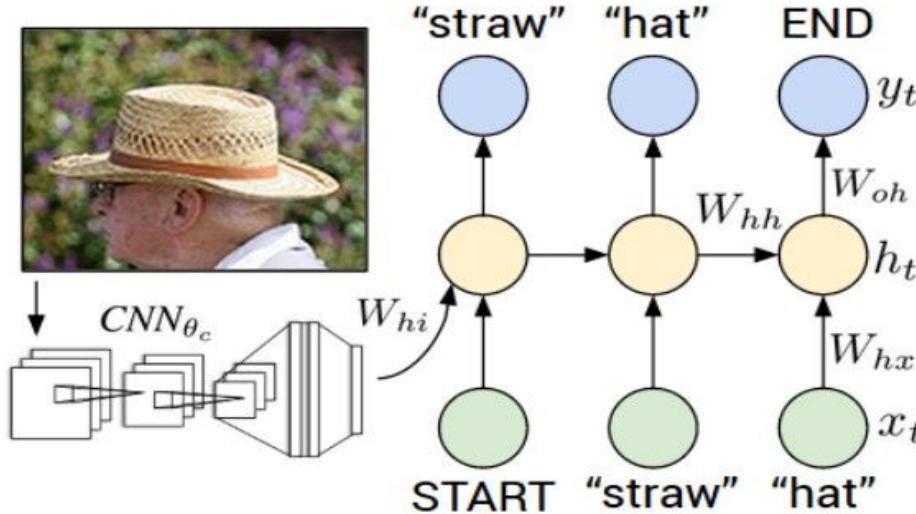
quote/comment cell

# 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.

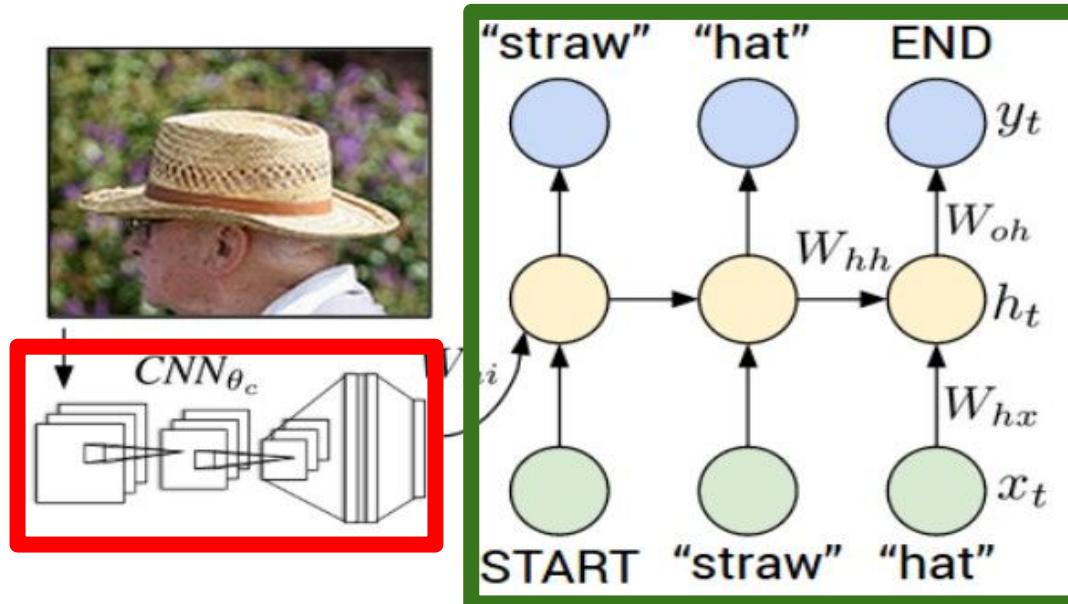
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei

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



image



test image



conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

FC-1000

softmax

image



test image



conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

FC-1000

softmax

X

image



test image



conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

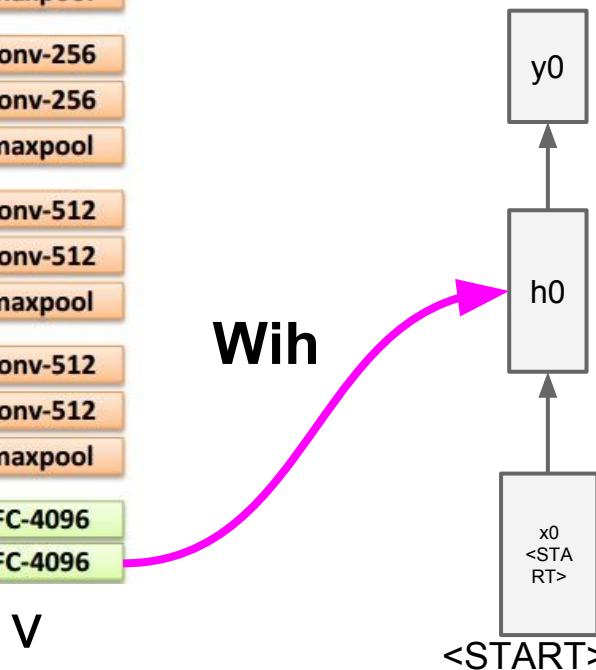
FC-4096



<START>



test image



**before:**

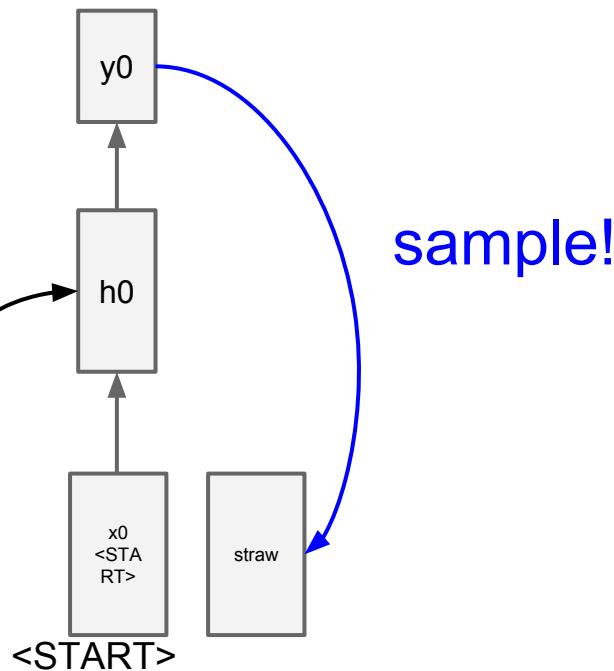
$$h = \tanh(W_{xh} * x + W_{hh} * h)$$

**now:**

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



test image



image



conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

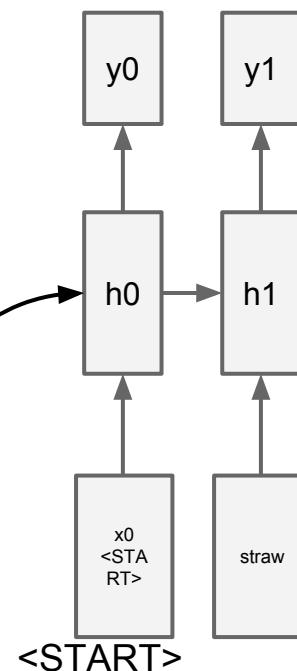
maxpool

FC-4096

FC-4096

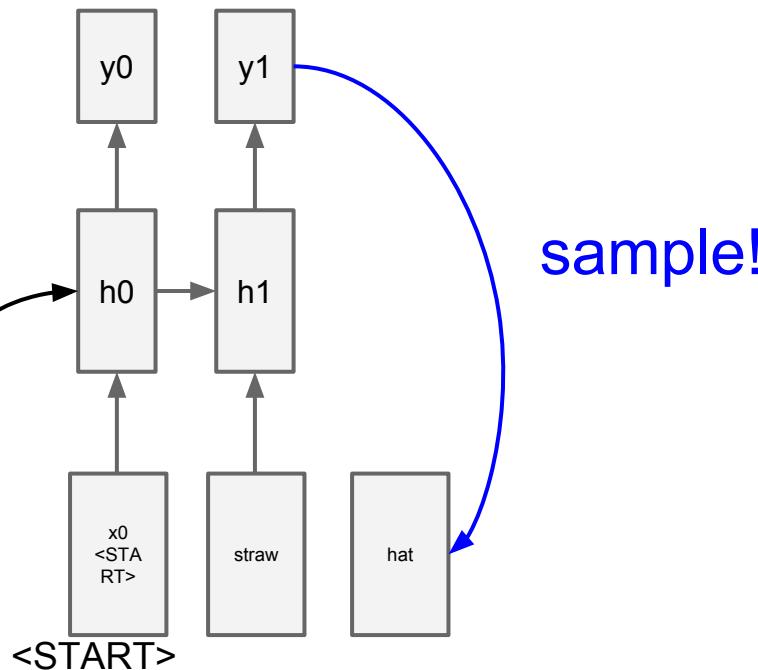


test image





test image



image



test image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

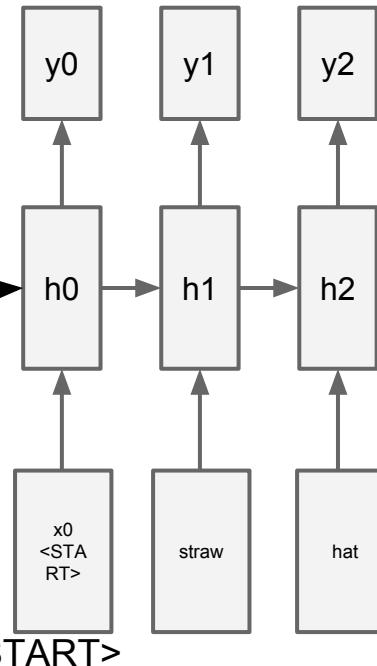
conv-512

conv-512

maxpool

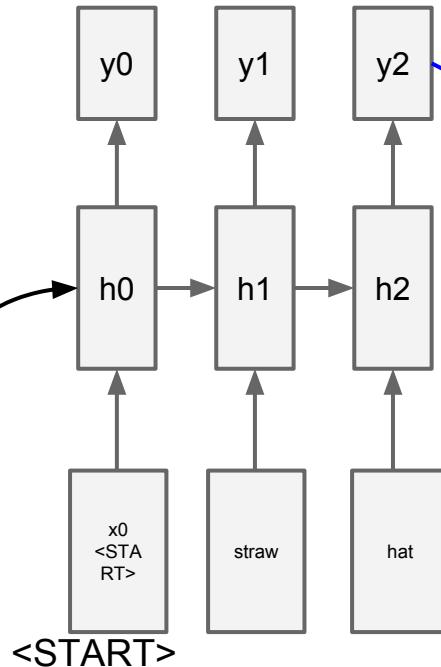
FC-4096

FC-4096





test image



sample  
<END> token  
=> finish.

# 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](http://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."



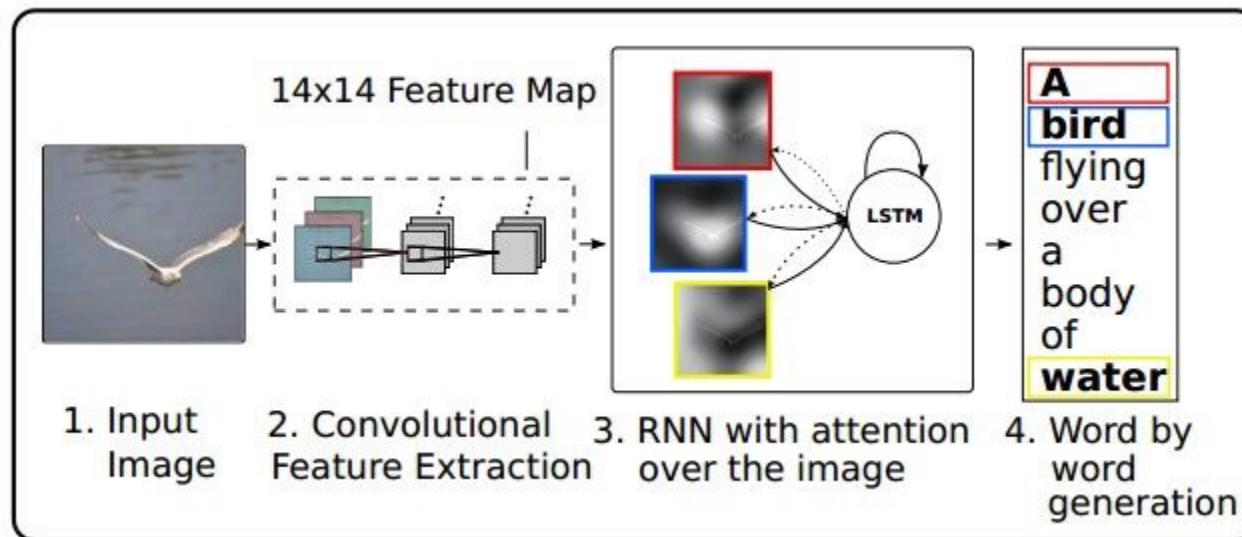
"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."

# Preview of fancier architectures

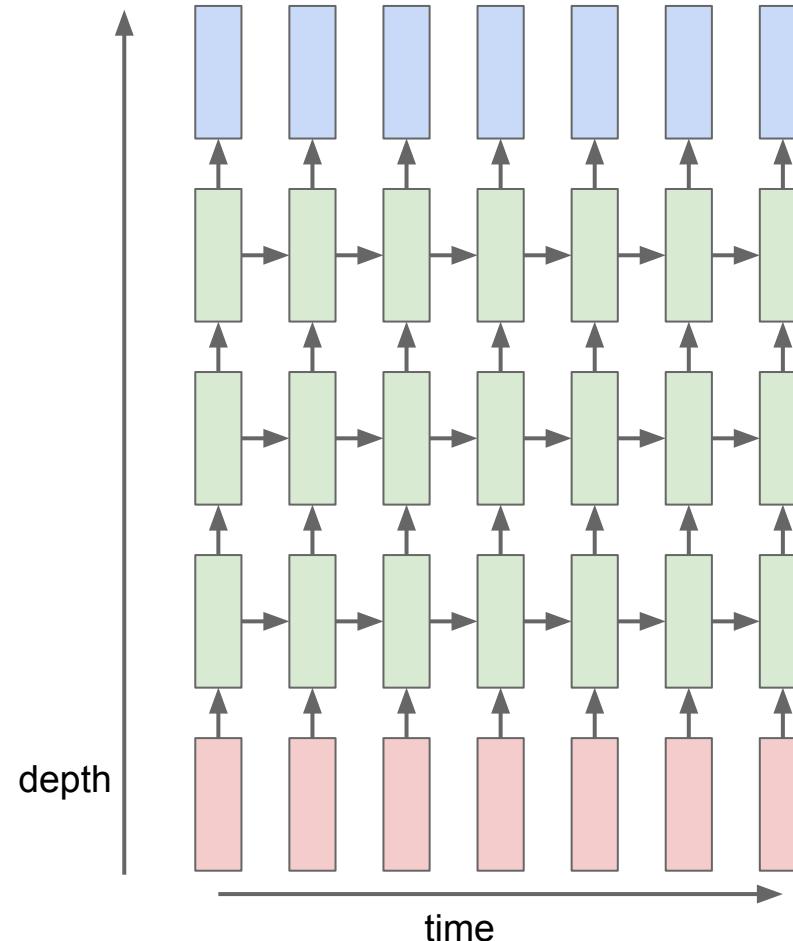
RNN attends spatially to different parts of images while generating each word of the sentence:



# 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:

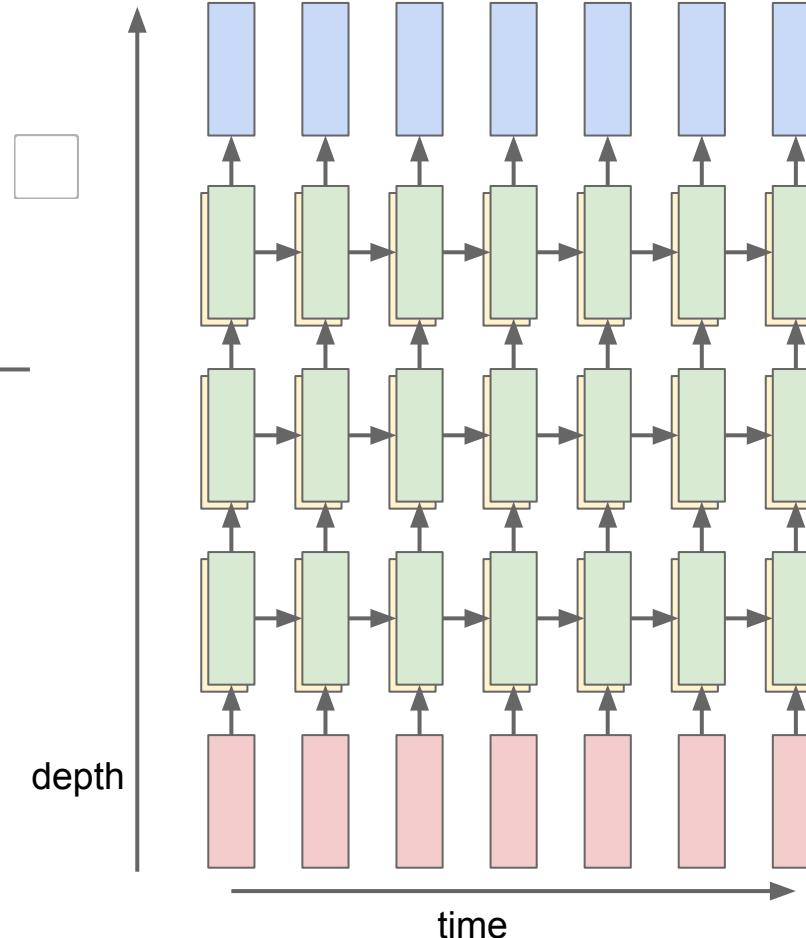
$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

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

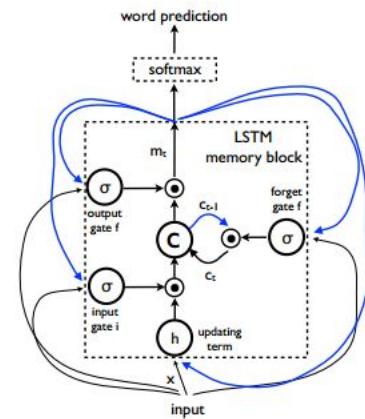
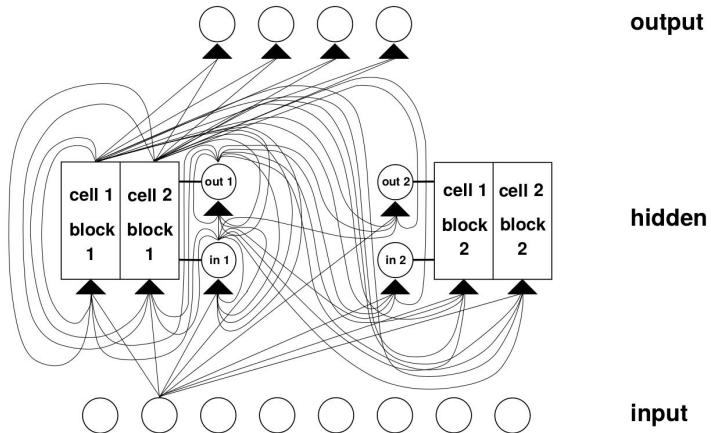
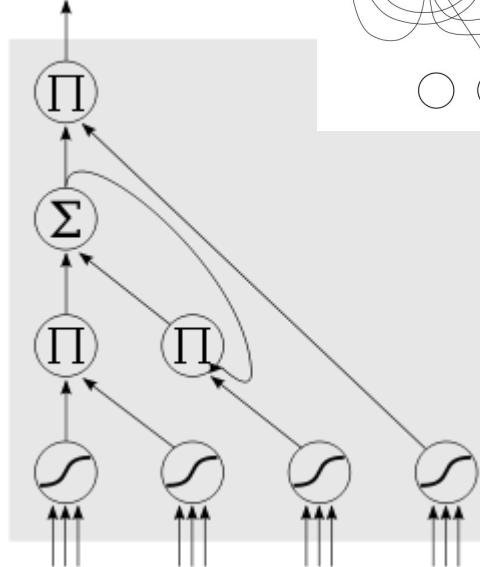
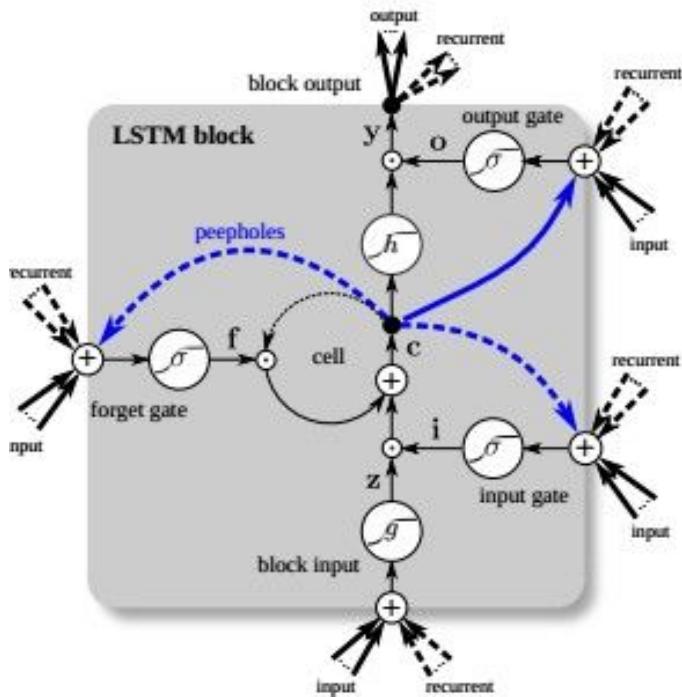
# LSTM:

$$W^l [4n \times 2n]$$

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

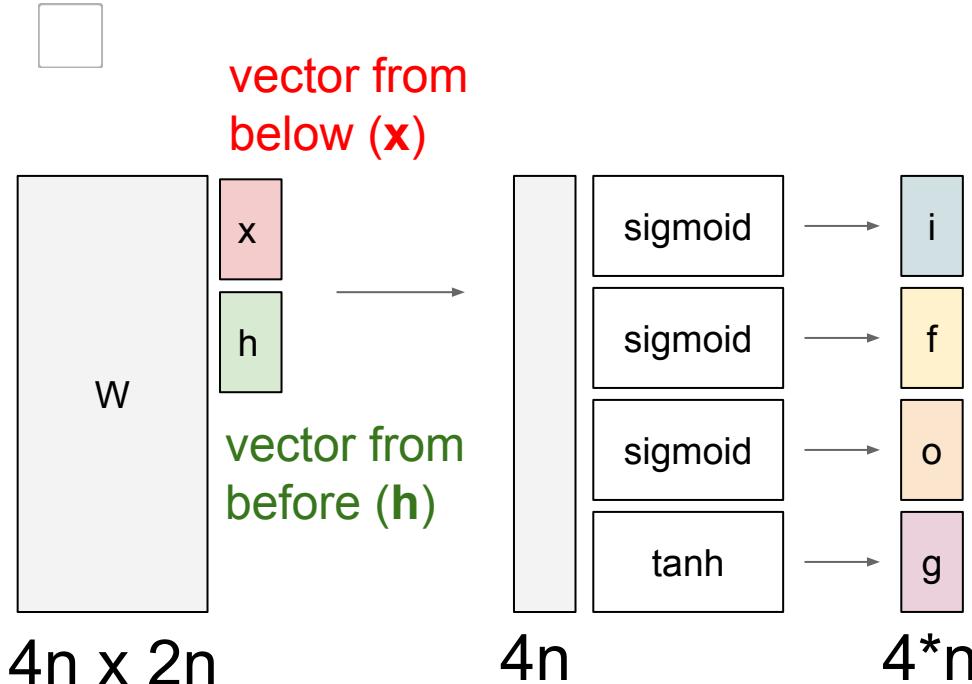


# LSTM



# Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]

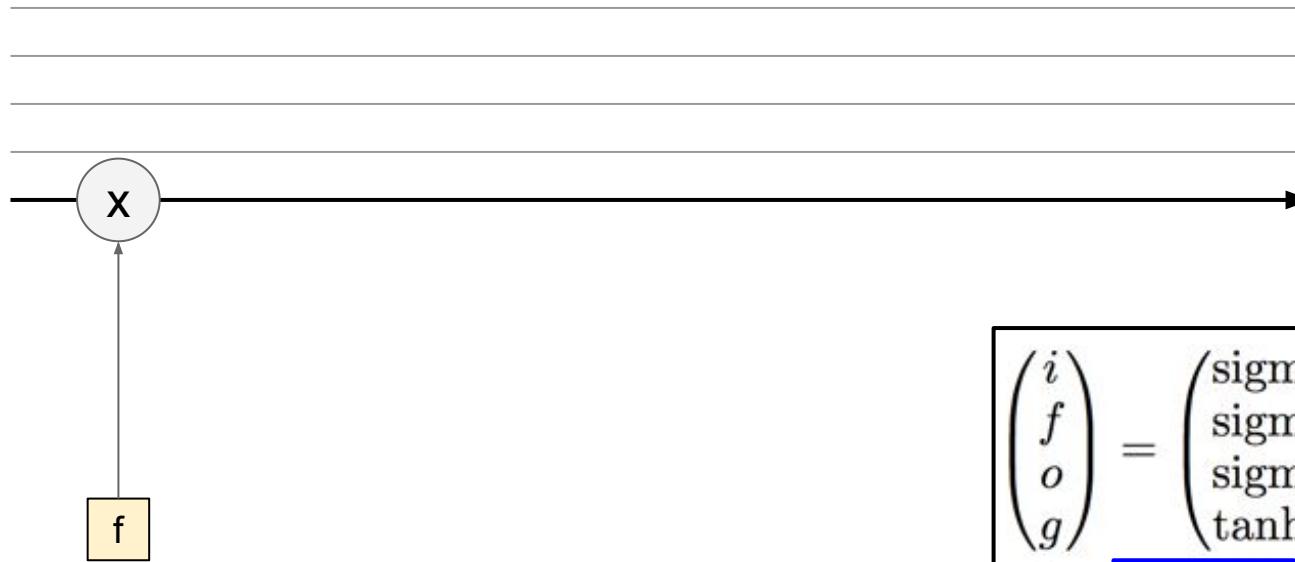


$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_t^l \end{pmatrix}$$
$$c_t^l = f \odot c_{t-1}^l + i \odot g$$
$$h_t^l = o \odot \tanh(c_t^l)$$

# Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]

cell  
state  $c$

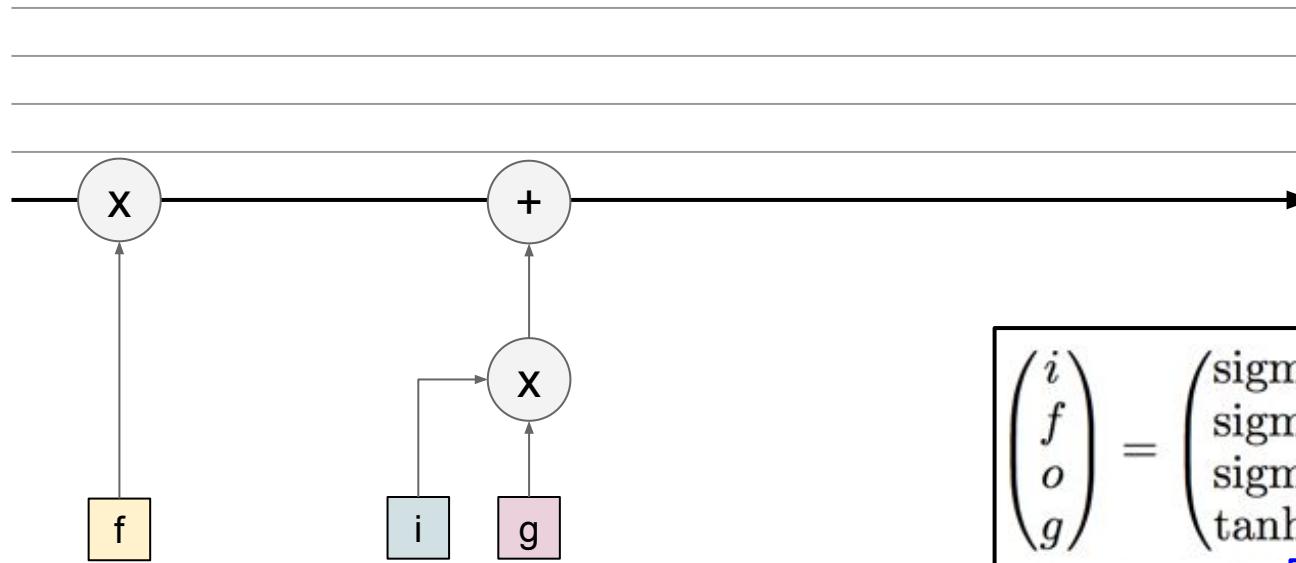


$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$
$$c_t^l = f \odot c_{t-1}^l + i \odot g$$
$$h_t^l = o \odot \tanh(c_t^l)$$

# Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]

cell  
state  $c$

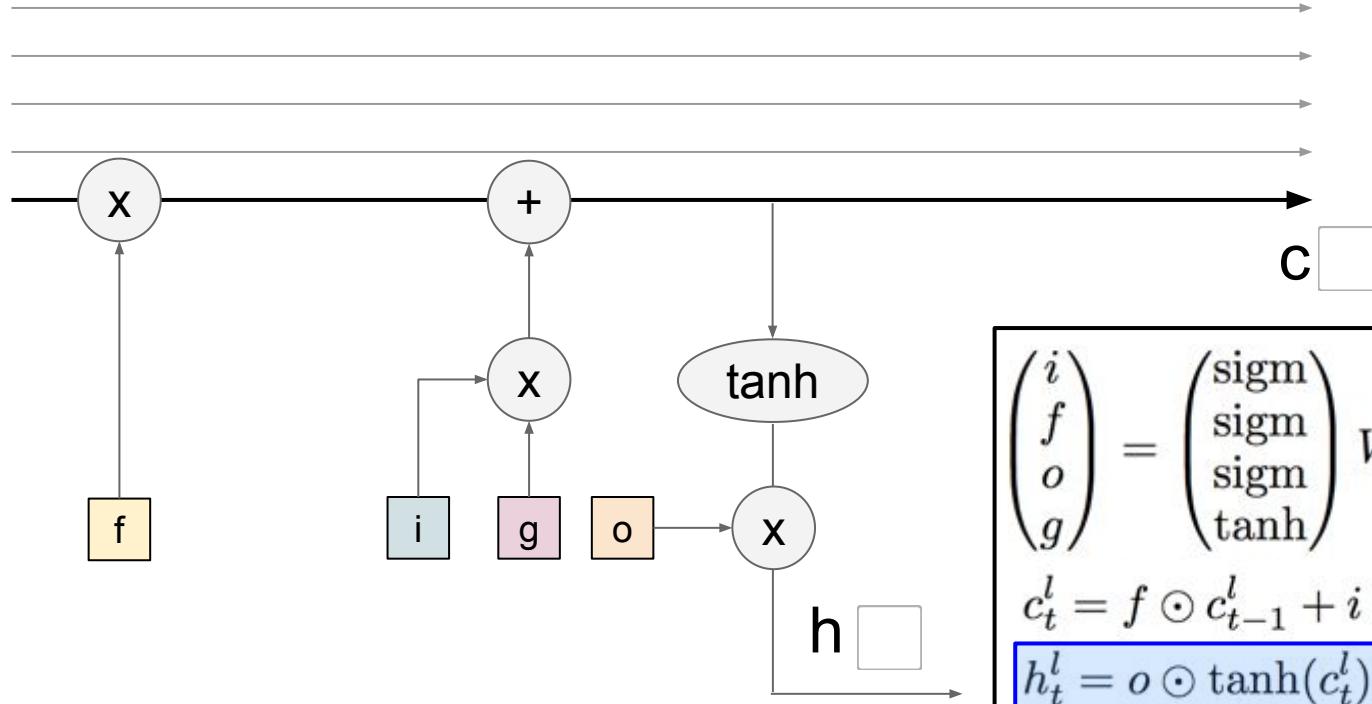


$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$
$$c_t^l = f \odot c_{t-1}^l + i \odot g$$
$$h_t^l = o \odot \tanh(c_t^l)$$

# Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]

cell  
state  $c$



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

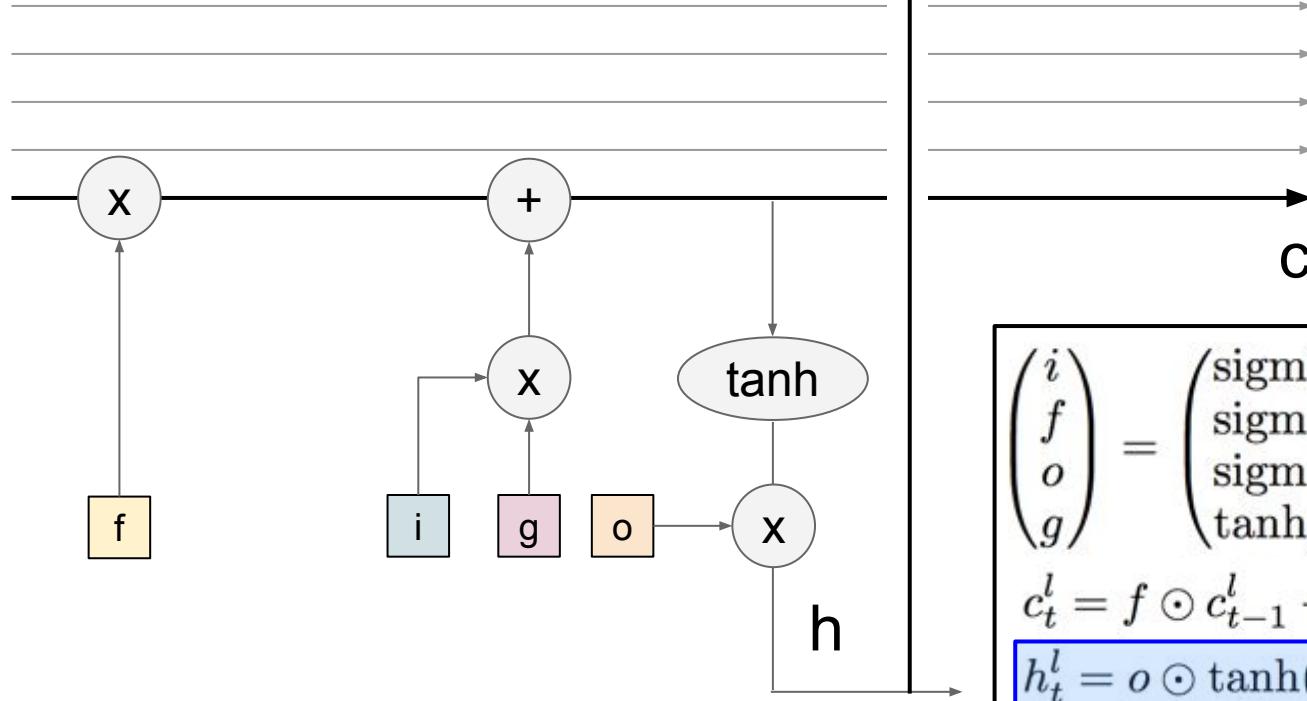
$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

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

# Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]

cell  
state  $c$



higher layer, or  
prediction

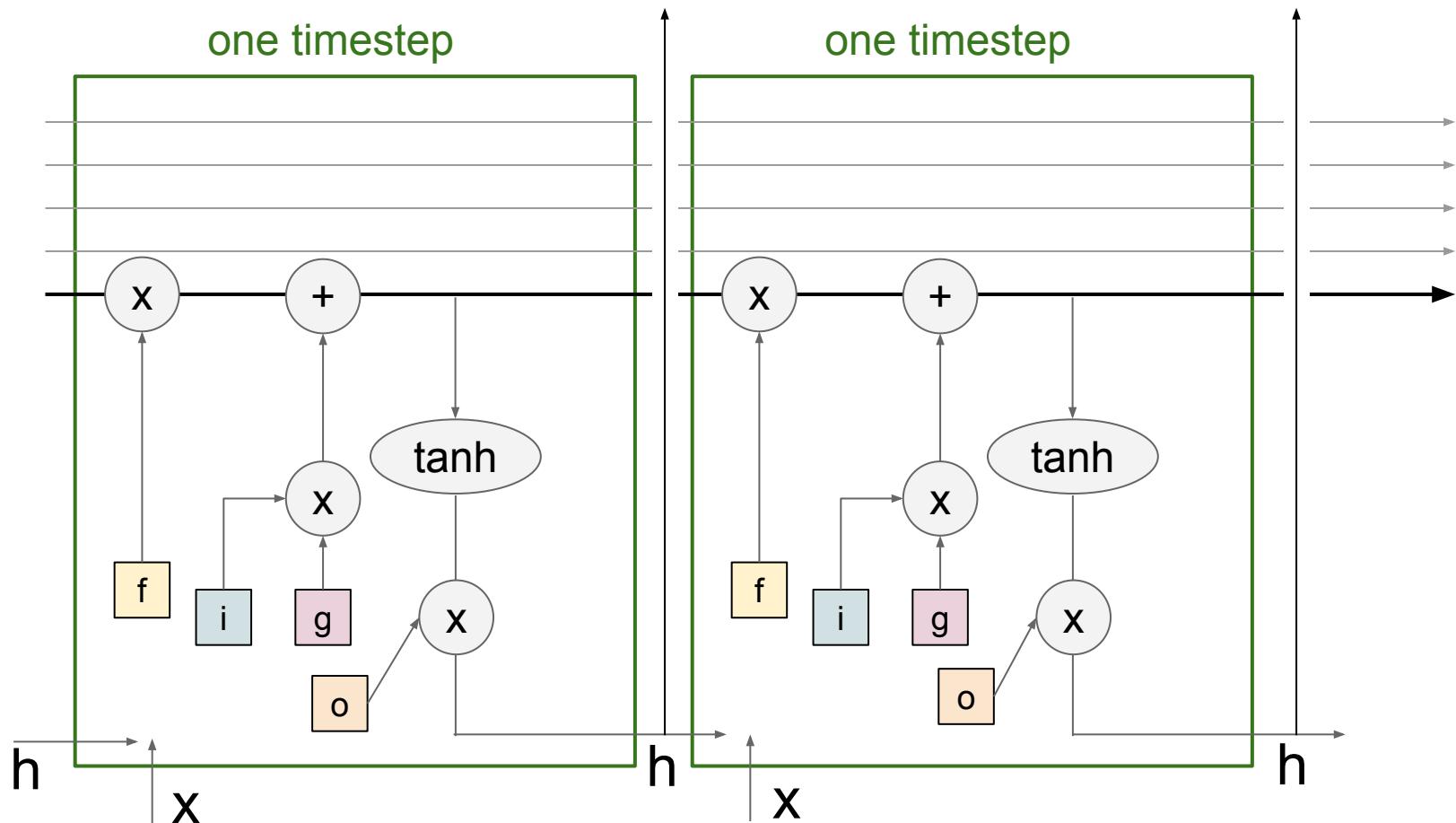
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$
$$c_t^l = f \odot c_{t-1}^l + i \odot g$$
$$h_t^l = o \odot \tanh(c_t^l)$$



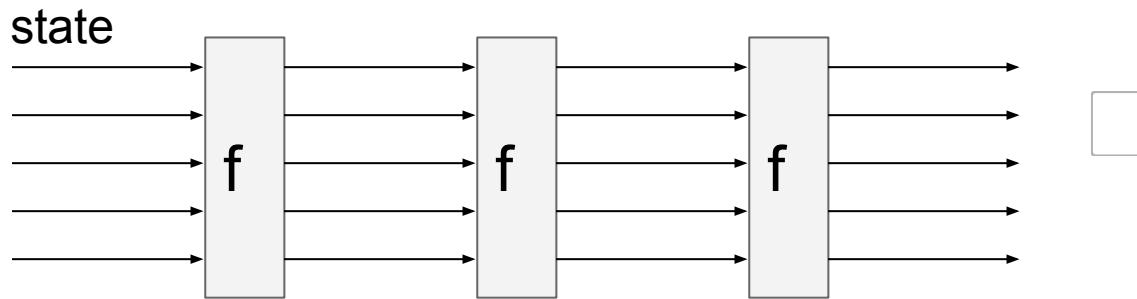
# LSTM

one timestep

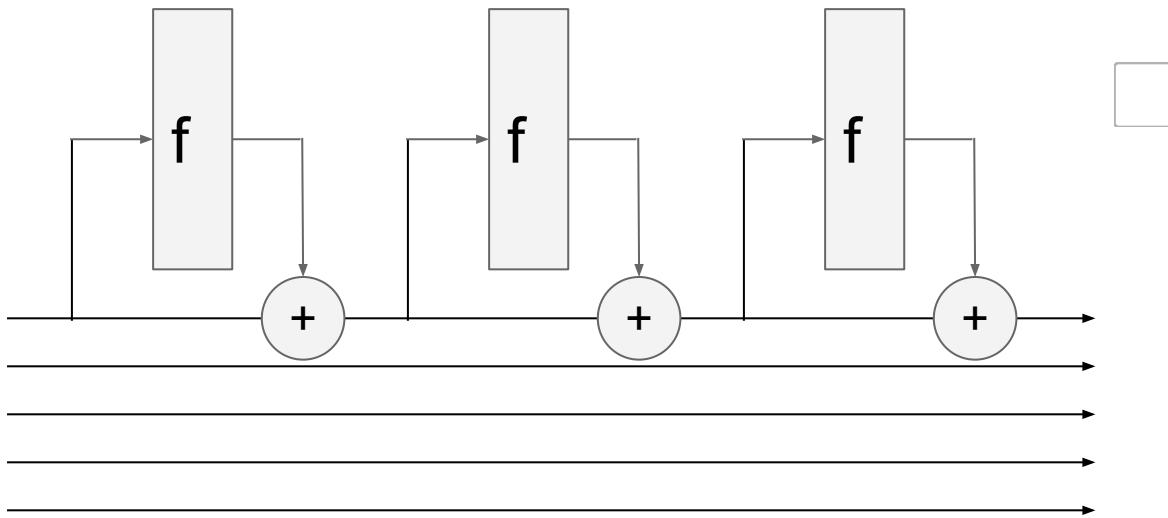
cell  
state  $c$



# RNN

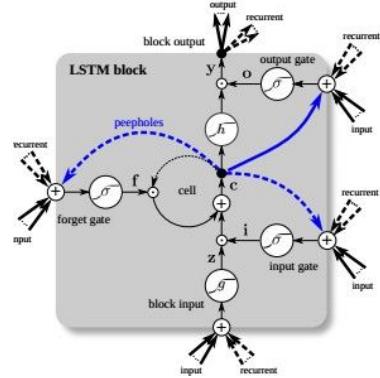


# LSTM (ignoring forget gates)



# 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.