

رسالة محمد



یادگیری عمیق

مدرس: محمدرضا محمدی

بهار ۱۴۰۲

شبکه‌های عصبی بازگشتی

Recurrent Neural Networks

مدل‌های مارکوف

- فرض می‌کنیم نمونه فعلی به تعداد کمی از نمونه‌های قبلی وابسته باشد

- مرتبه ۰ (مستقل):

$$P(x_1, x_2, x_3, x_4) = P(x_1)P(x_2)P(x_3)P(x_4)$$

- مرتبه ۱:

$$P(x_1, x_2, x_3, x_4) = P(x_1)P(x_2|x_1)P(x_3|x_2)P(x_4|x_3)$$

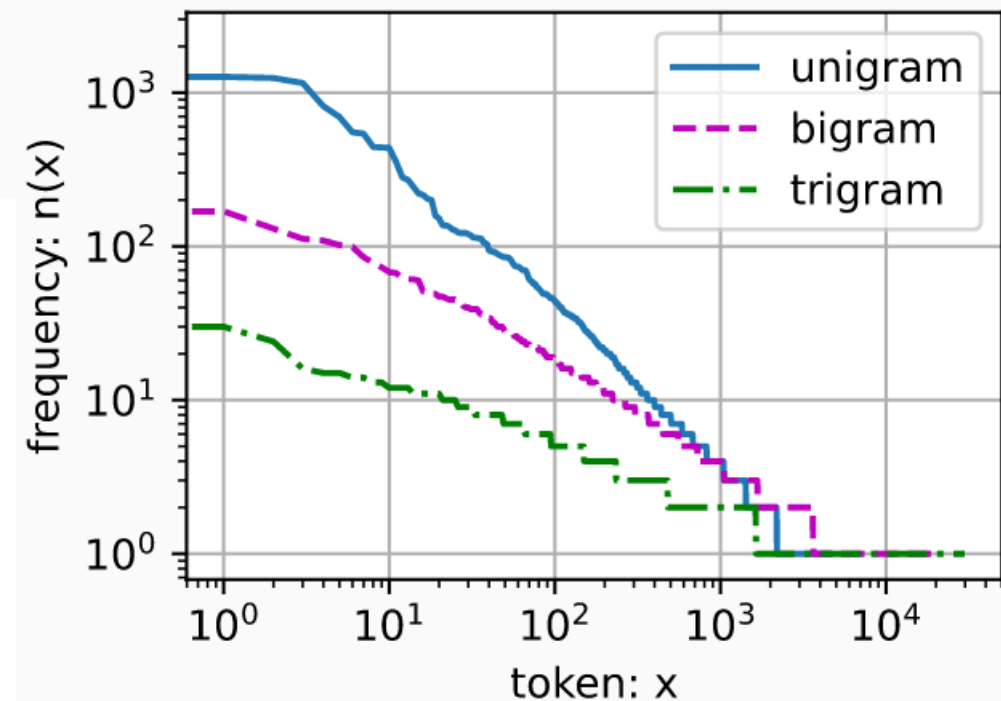
- مرتبه ۲:

$$P(x_1, x_2, x_3, x_4) = P(x_1)P(x_2|x_1)P(x_3|x_1, x_2)P(x_4|x_2, x_3)$$

- به این مدل‌ها به ترتیب unigram، bigram و trigram گفته می‌شود

مثال: مجموعه داده The Time Machine

```
[('the', 2261),  
 ('i', 1267),  
 ('and', 1245),  
 ('of', 1155),  
 ('a', 816),  
 ('to', 695),  
 ('was', 552),  
 ('in', 541),  
 ('that', 443),  
 ('my', 440)]  
[(('of', 'the'), 309),  
 (('in', 'the'), 169),  
 (('i', 'had'), 130),  
 (('i', 'was'), 112),  
 (('and', 'the'), 109),  
 (('the', 'time'), 102),  
 (('it', 'was'), 99),  
 (('to', 'the'), 85),  
 (('as', 'i'), 78),  
 (('of', 'a'), 73)]  
[(('the', 'time', 'traveller'), 59),  
 (('the', 'time', 'machine'), 30),  
 (('the', 'medical', 'man'), 24),  
 (('it', 'seemed', 'to'), 16),  
 (('it', 'was', 'a'), 15),  
 (('here', 'and', 'there'), 15),  
 (('seemed', 'to', 'me'), 14),  
 (('i', 'did', 'not'), 14),  
 (('i', 'saw', 'the'), 13),  
 (('i', 'began', 'to'), 13)]
```



یادگیری یک مدل زبان

- برای تخمین احتمال شرطی می توان از حالت پنهان استفاده کرد

$$P(x_t \mid x_{t-1}, \dots, x_1) \approx P(x_t \mid h_{t-1})$$

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

- با استفاده از یادگیری عمیق می توانیم مدلی آموزش دهیم که بتواند احتمال توکن بعدی را پیش بینی کند

مثال: مدل زبانی سطح کاراکتر

- دنباله آموزشی نمونه:

“hello” -

input chars: “h” “e” “l” “l”

مثال: مدل زبانی سطح کاراکتر

target chars: "e" "l" "l" "o"

• دنباله آموزشی نمونه:

- "hello"

input chars: "h" "e" "l" "l"

مثال: مدل زبانی سطح کاراکتر

target chars: "e" "l" "l" "o"

- دنباله آموزشی نمونه:

- "hello"

- توکن‌ها:

- [h,e,l,o]

input layer

1
0
0
0

0
1
0
0

0
0
1
0

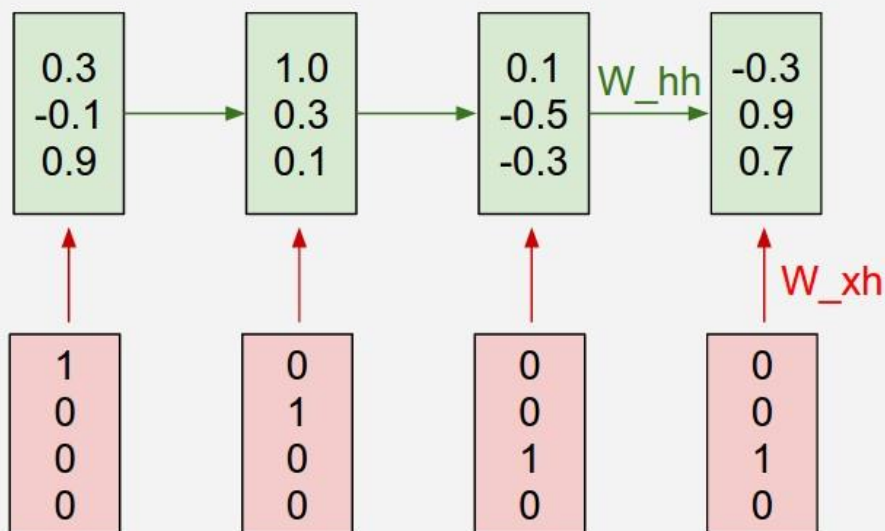
0
0
1
0

input chars: "h" "e" "l" "l"

مثال: مدل زبانی سطح کاراکتر

target chars: "e" "l" "l" "o"

hidden layer



- دنباله آموزشی نمونه:

- "hello"

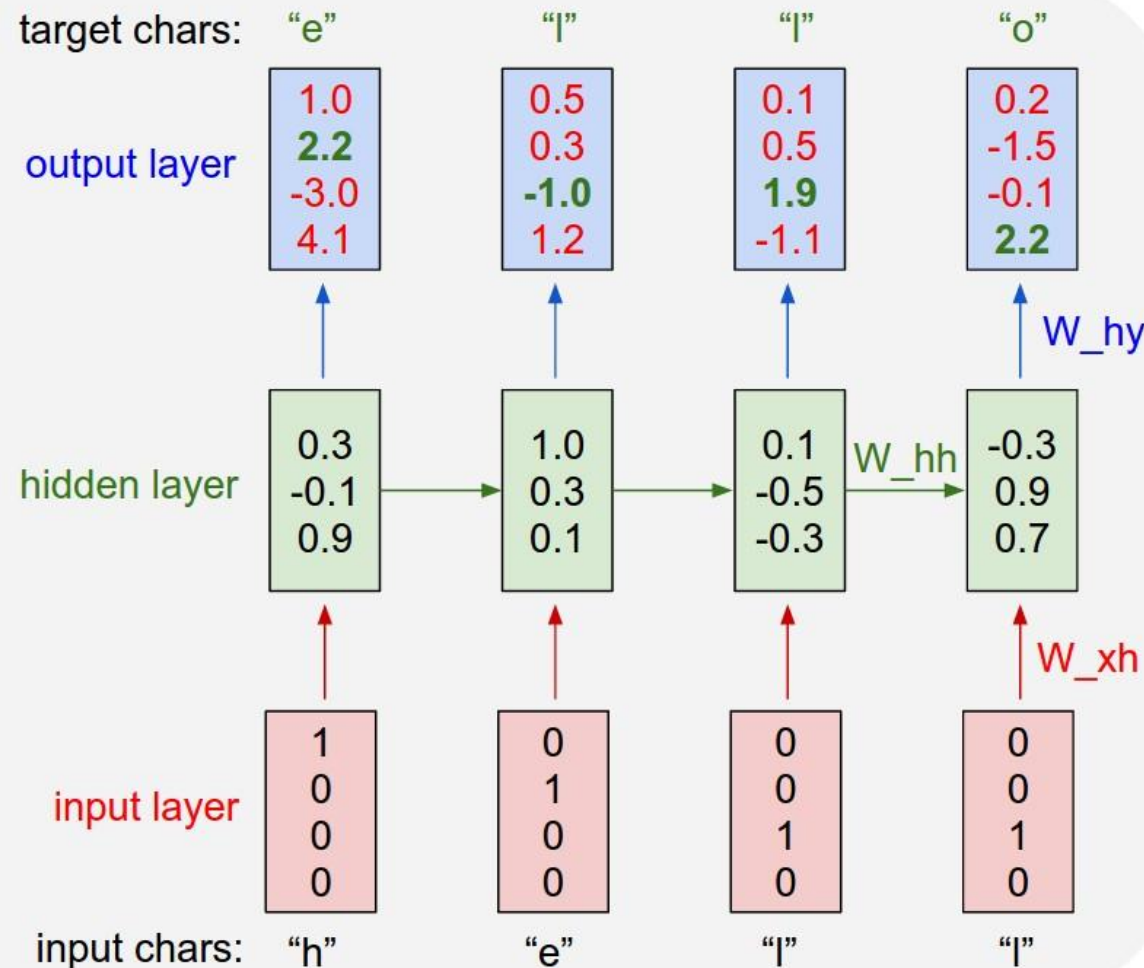
- توکن‌ها:

- [h,e,l,o]

- لایه بازگشتی میانی:

$$\mathbf{h}_t = \phi(\mathbf{W}_{xh}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h) -$$

مثال: مدل زبانی سطح کاراکتر



- دنباله آموزشی نمونه:

- "hello"

- توکن‌ها:

- [h,e,l,o]

- لایه بازگشتی میانی:

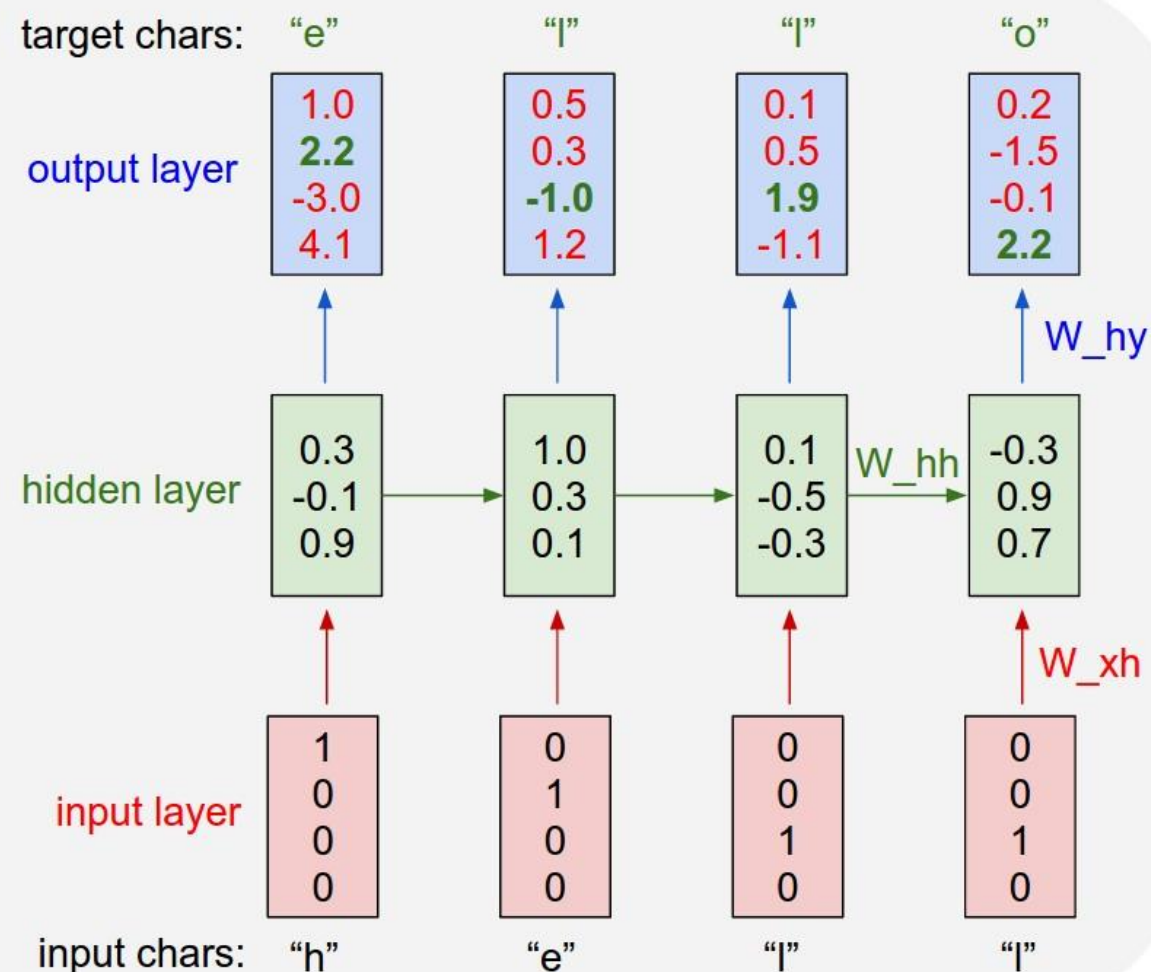
$$\mathbf{h}_t = \phi(\mathbf{W}_{xh}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h) -$$

- لایه کاملاً متصل خروجی:

$$\mathbf{y}_t = \mathbf{W}_{hy}\mathbf{h}_t -$$

- می‌توان از SoftMax هم استفاده کرد

مثال: مدل زبانی سطح کاراکتر



• در زمان تست:

- در هر گام یک کاراکتر نمونه برداری می شود و ورودی گام بعد می شود

min-char-rnn.py

```

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 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     # 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     dhnext = 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. see http://cs231n.github.io/neural-networks-case-study/#grad if confused here
51         dwhy += np.dot(dy, hs[t].T)
52         dby += dy
53         dh = np.dot(why.T, dy) + dhnext # backprop into h
54         dhraw = (1 - hs[t] * hs[t]) * dh # backprop through tanh nonlinearity

```

```

55     dbh += dhraw
56     dwxh += np.dot(dhraw, xs[t].T)
57     dwhh += np.dot(dhraw, hs[t-1].T)
58     dhnext = np.dot(whh.T, dhraw)
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[len(inputs)-1]
62
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(wxh, 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     return ixes
80
81 n, p = 0, 0
82 mwxx, mwxxh, mwxy = 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 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 \n----' % (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    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([dwxh, dwhh, dwhy, dbh, dby],
106                                  [dwxh, dwhh, dwhy, dbh, dby],
107                                  [mwxx, mwxxh, mwxy, 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

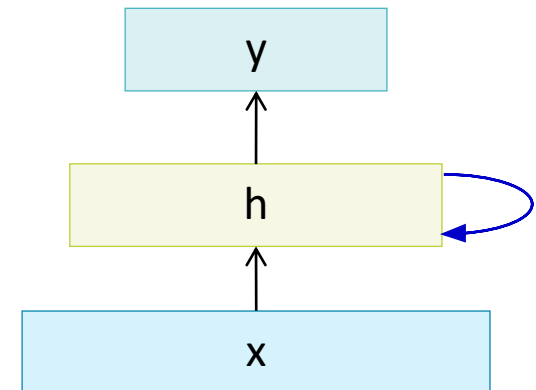
```

Shakespeare

- شامل حدود ۱۰۰,۰۰۰ کلمه

input.txt ✕

```
1 That, poor contempt, or claim'd thou slept so faithful,  
2 I may contrive our father; and, in their defeated queen,  
3 Her flesh broke me and puttance of expedition house,  
4 And in that same that ever I lament this stomach,  
5 And he, nor Butly and my fury, knowing everything  
6 Grew daily ever, his great strength and thought  
7 The bright buds of mine own.  
8  
9 BIONDELLO:  
10 Marry, that it may not pray their patience.'  
11  
12 KING LEAR:  
13 The instant common maid, as we may less be  
14 a brave gentleman and joiner: he that finds us with wax  
15 And owe so full of presence and our fooder at our  
16 staves. It is remorse'd the bridal's man his grace  
17 for every business in my tongue, but I was thinking  
18 that he contends, he hath respected thee.  
19  
20 BIRON:  
21 She left thee on, I'll die to blessed and most reasonable  
22 Nature in this honour, and her bosom is safe, some  
23 others from his speedy-birth, a bill and as  
24 Forestem with Richard in your heart  
25 Be question'd on, nor that I was enough:  
26 Which of a partier forth the obsers d'punish'd the hate
```



تکامل نمونه‌ها در حین آموزش

• تکرار ۱۰۰

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

• تکرار ۳۰۰

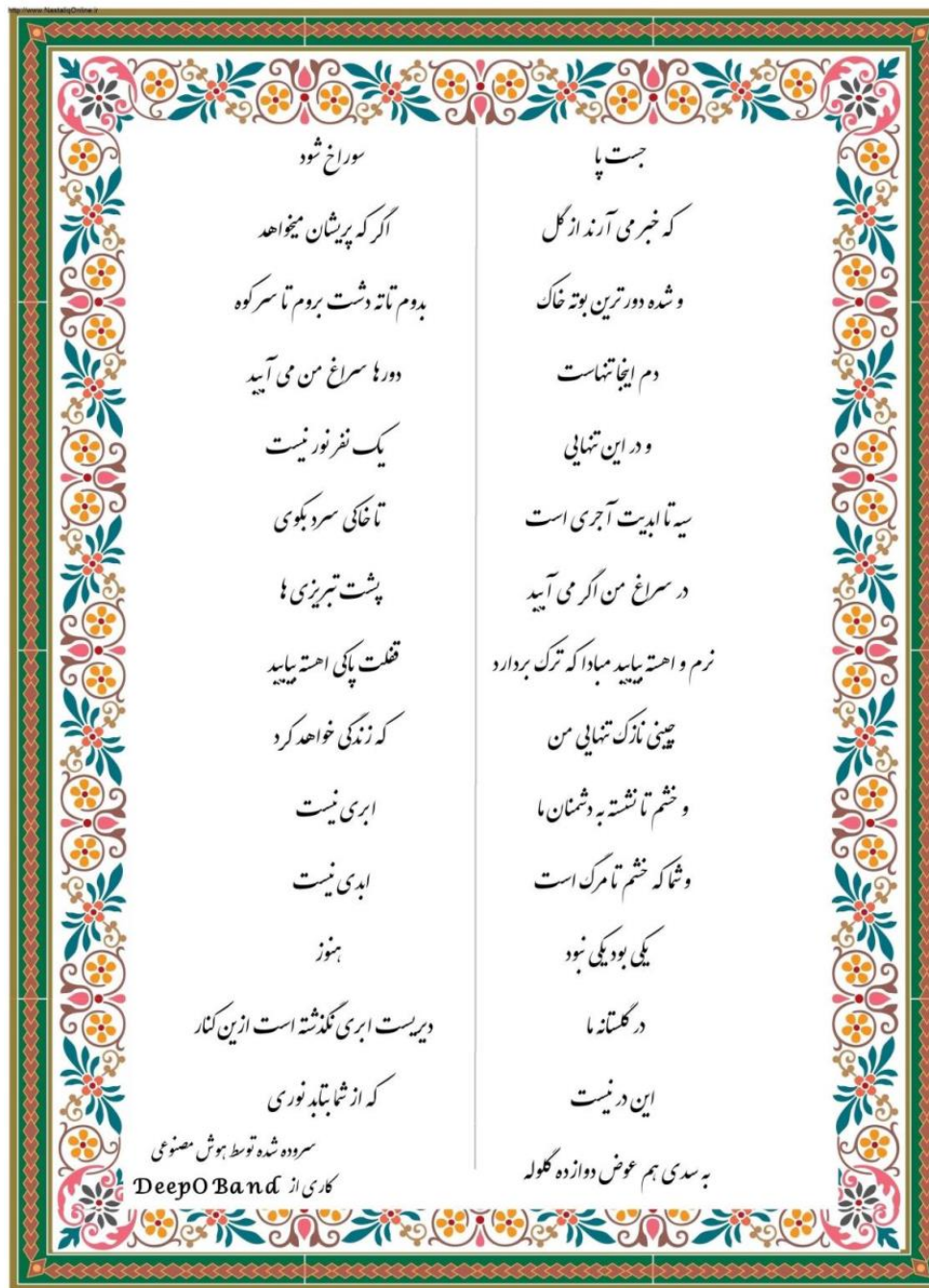
"Tmont thithey" fomesscerliund
Keushey. Thom here
sheulke, anmerenith ol sivh l lalterthend Bleipile shuw y fil on aseterlome
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

• تکرار ۷۰۰

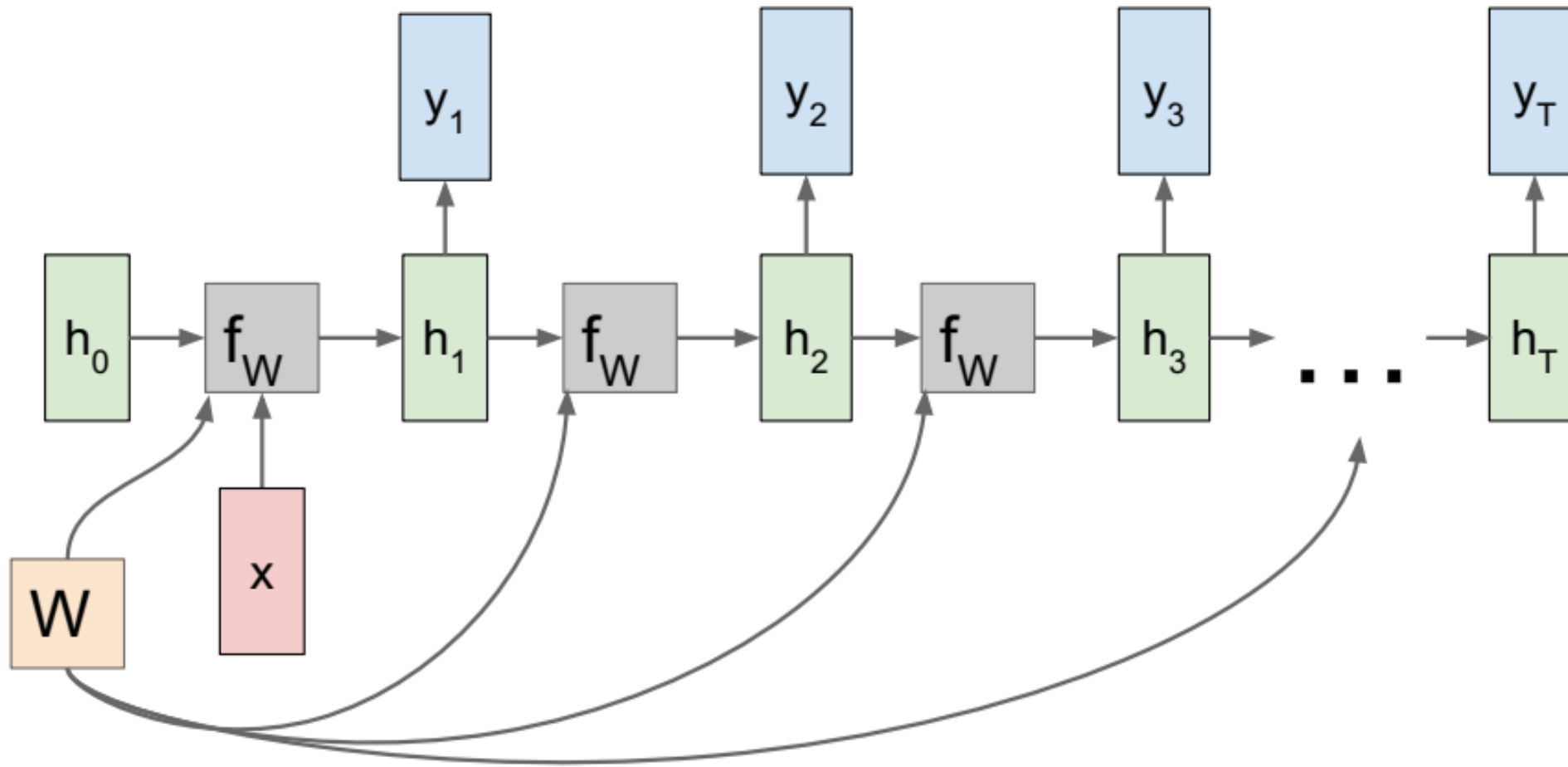
Aftair fall unsuch that the hall for Prince Velzonski's that me of
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort
how, and Gogition is so overelical and ofter.

• تکرار ۲۰۰۰

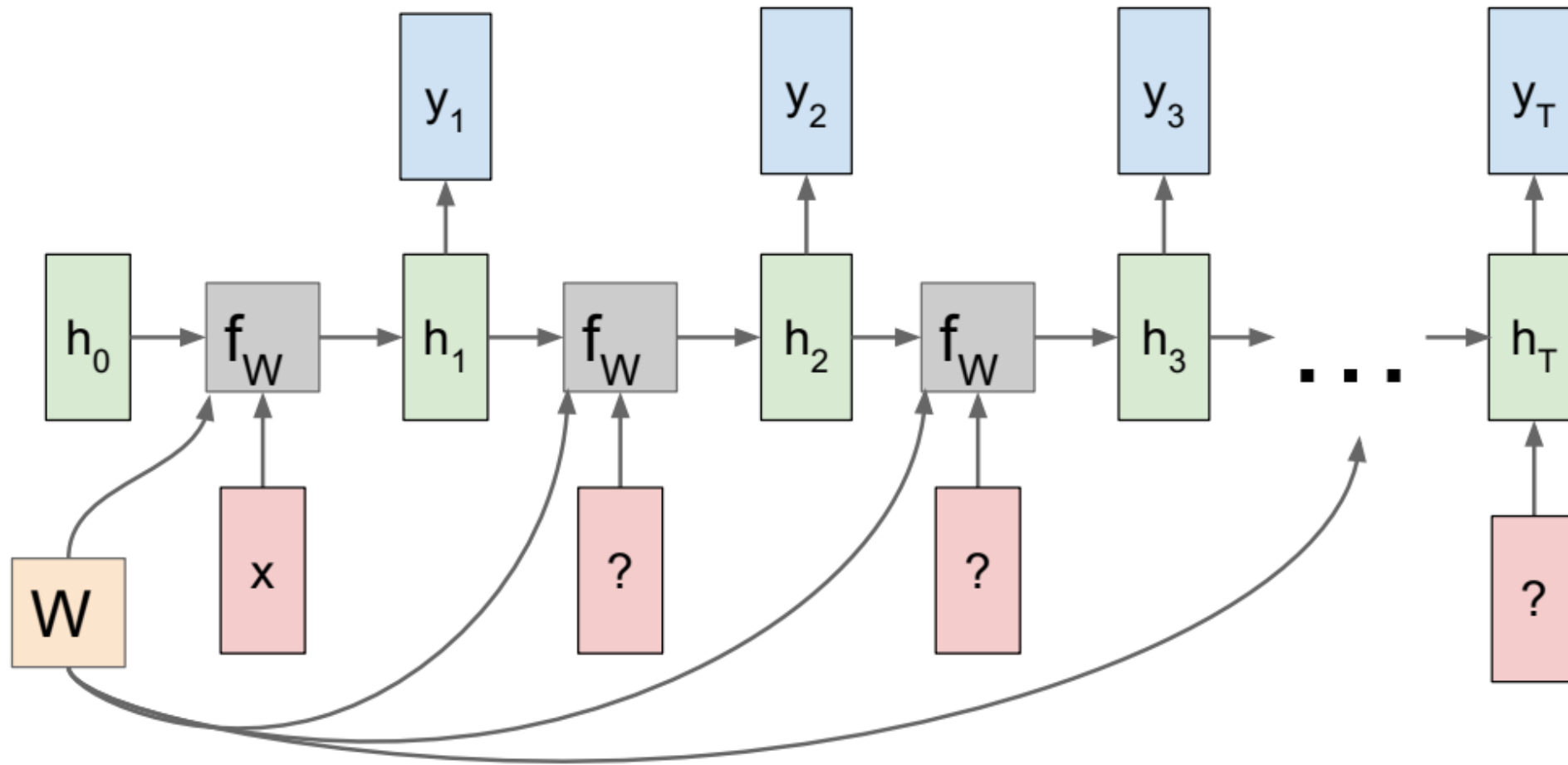
"Why do what that day," replied Natasha, and wishing to himself the fact the
princess, Princess Mary was easier, fed in had oftened him.
Pierre aking his soul came to the packs and drove up his father-in-law women.



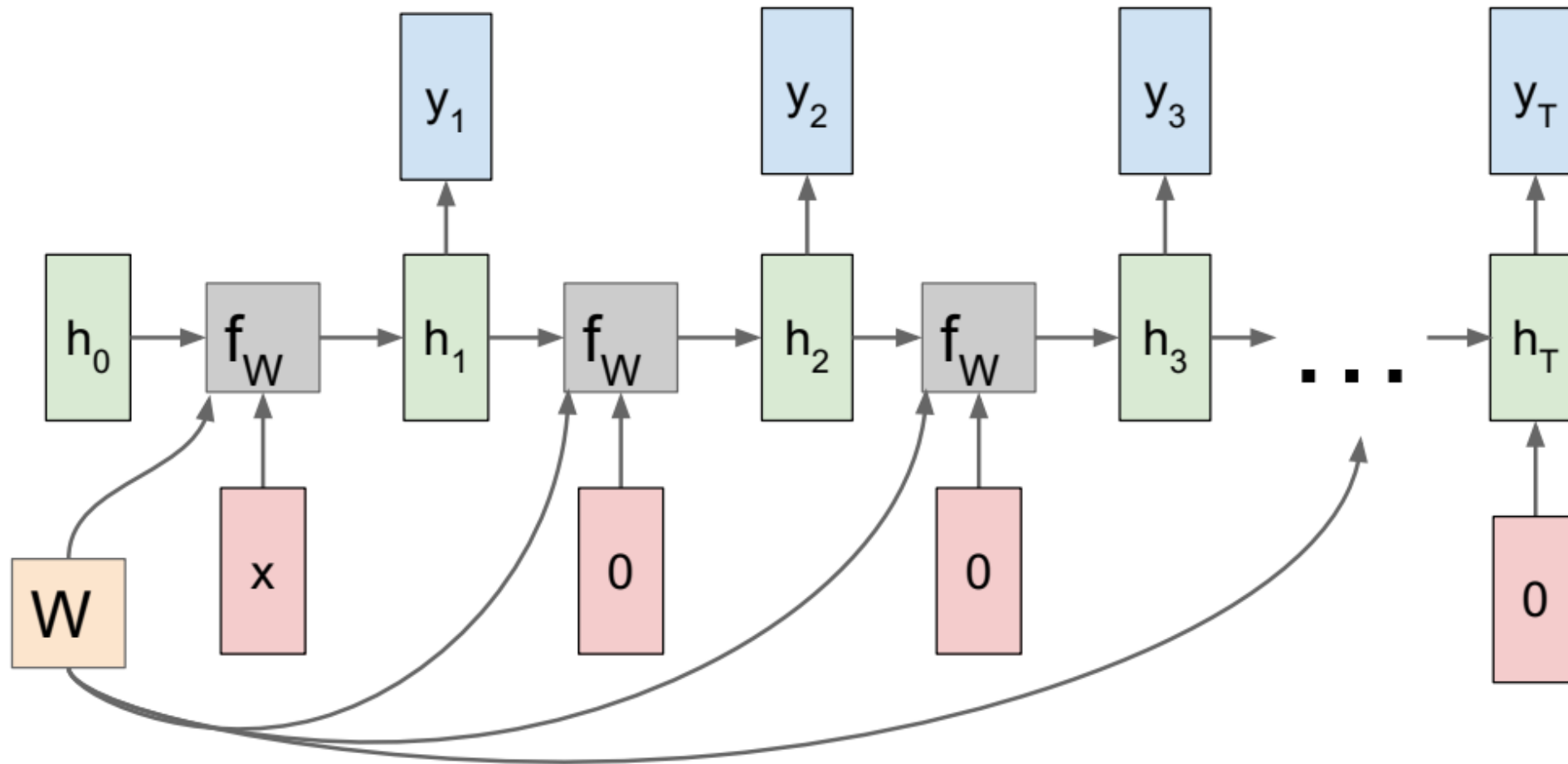
One to Many



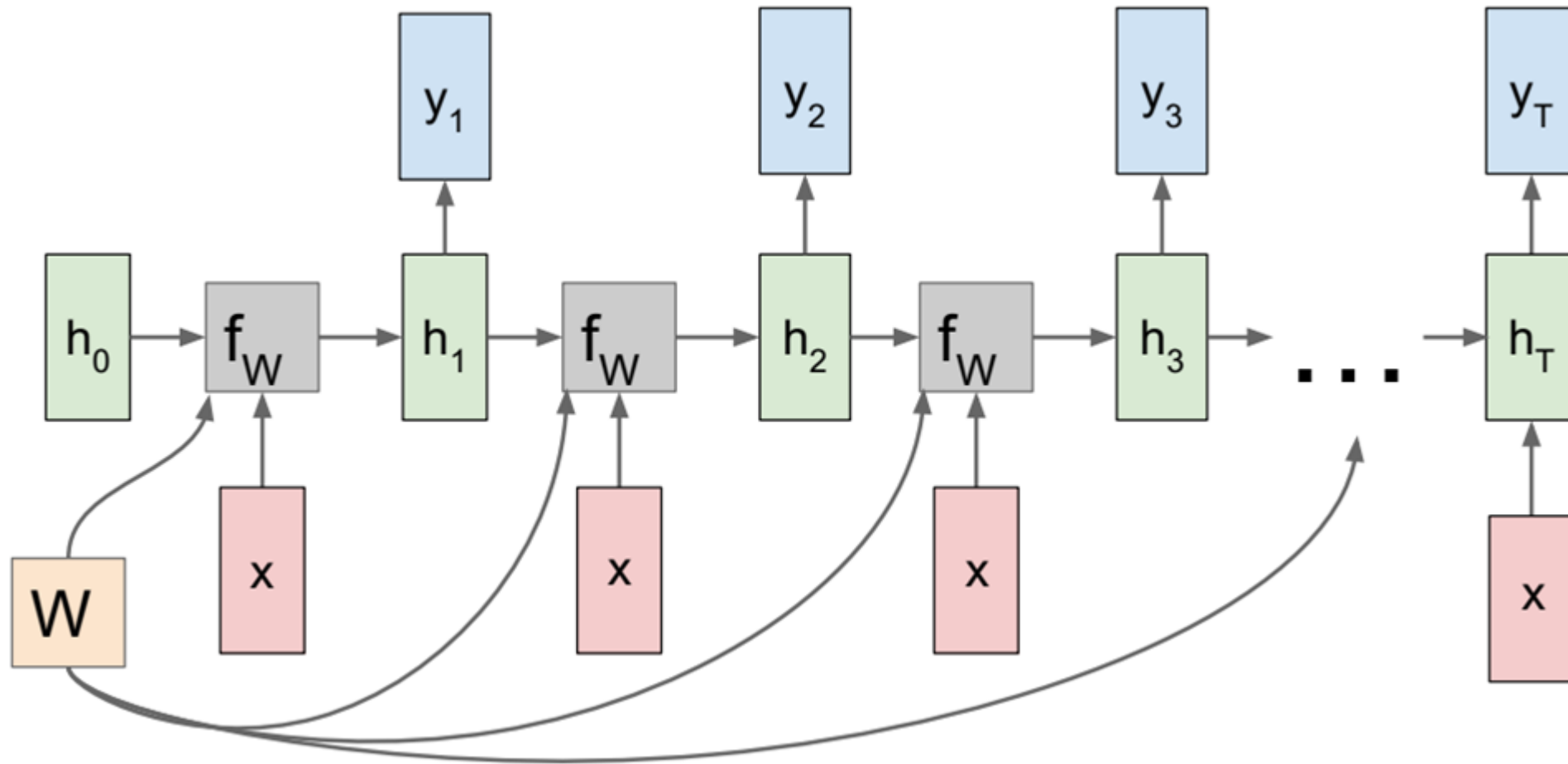
One to Many



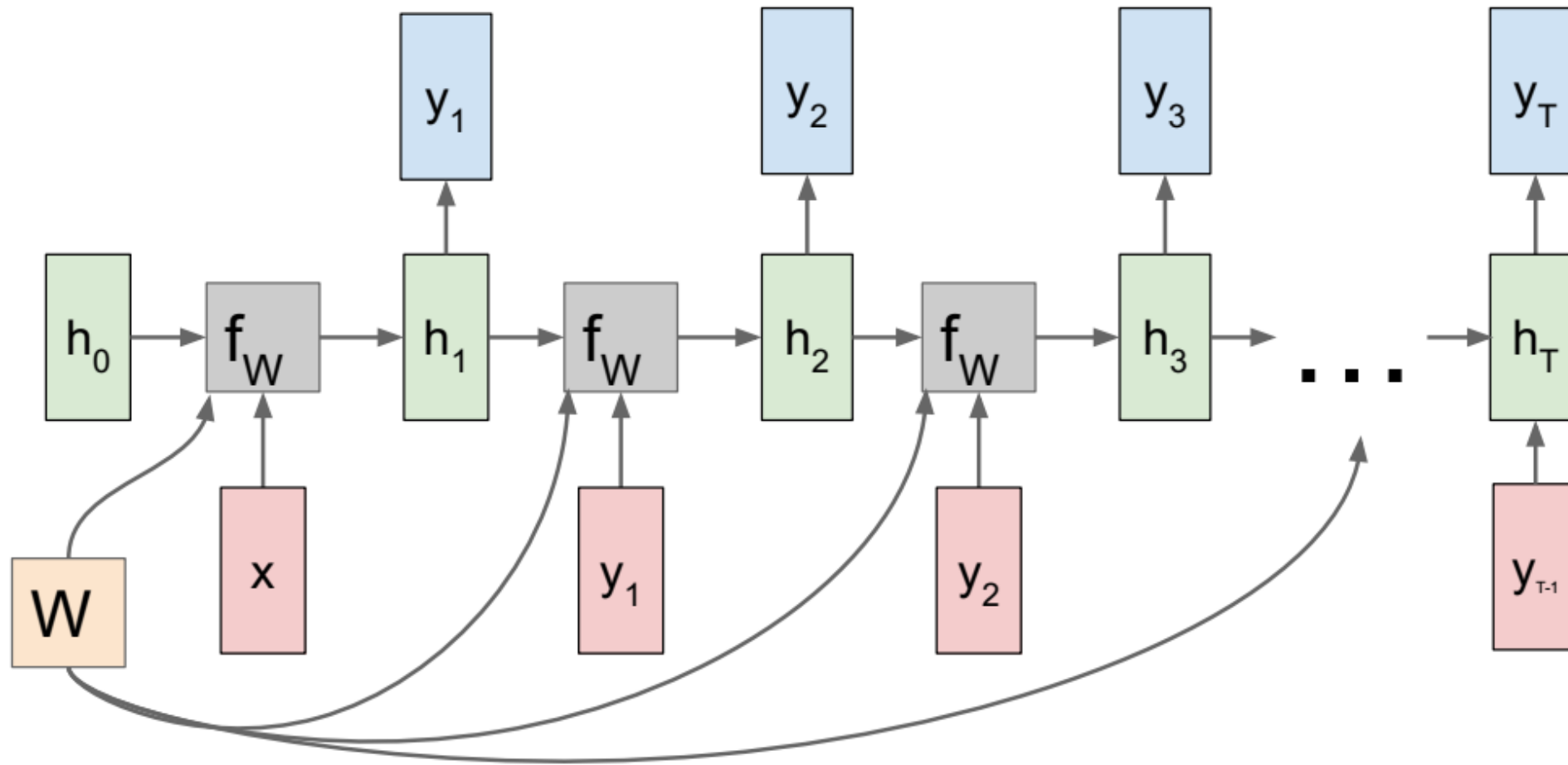
One to Many



One to Many

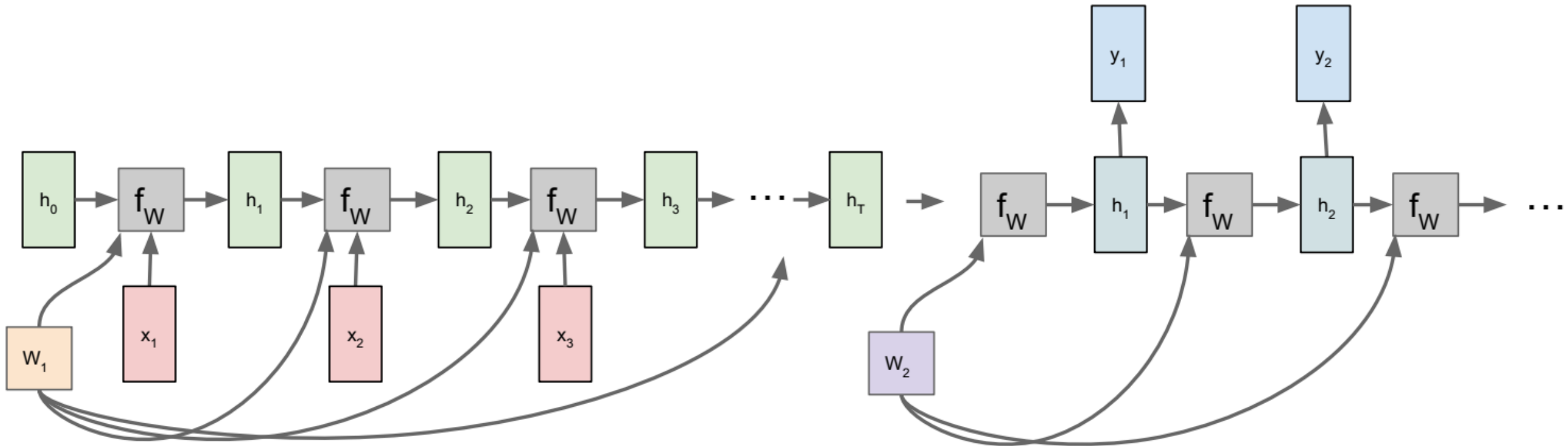


One to Many



Sequence to Sequence

- Many to One + One to Many
 - Many to one: Encode input sequence in a single vector
 - One to many: Produce output sequence from single input vector



پس انتشار در طول زمان

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

$$o_t = g(h_t, w_o)$$

- که f و g تبدیل‌های مربوط به لایه پنهان و لایه خروجی هستند

- w_o و w_h هم تمام پارامترهای این لایه‌ها هستند

- تابع ضرر در حالت many to many به صورت زیر تعریف می‌شود:

$$L(x_1, \dots, x_T, y_1, \dots, y_T, w_h, w_o) = \frac{1}{T} \sum_{t=1}^T l(y_t, o_t)$$

- باید گرادینان این تابع ضرر نسبت به w_o و w_h را محاسبه کنیم

$$\frac{\partial L}{\partial w_o} = \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial w_o} = \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial o_t} \frac{\partial g(h_t, w_o)}{\partial w_o}$$

پس انتشار در طول زمان

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

$$o_t = g(h_t, w_o)$$

$$L(x_1, \dots, x_T, y_1, \dots, y_T, w_h, w_o) = \frac{1}{T} \sum_{t=1}^T l(y_t, o_t)$$

$$\frac{\partial L}{\partial w_h} = \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial w_h} = \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial o_t} \frac{\partial g(h_t, w_o)}{\partial w_h} = \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial o_t} \frac{\partial g(h_t, w_o)}{\partial h_t} \frac{\partial h_t}{\partial w_h}$$

- محاسبه $\frac{\partial h_t}{\partial w_h}$ ساده نیست زیرا در محاسبه h_t به تعداد t بار از w_h به صورت بازگشتی استفاده می شود

پس انتشار در طول زمان

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

$$o_t = g(h_t, w_o)$$

$$L(x_1, \dots, x_T, y_1, \dots, y_T, w_h, w_o) = \frac{1}{T} \sum_{t=1}^T l(y_t, o_t)$$

$$\frac{\partial L}{\partial w_h} = \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial w_h} = \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial o_t} \frac{\partial g(h_t, w_o)}{\partial w_h} = \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial o_t} \frac{\partial g(h_t, w_o)}{\partial h_t} \frac{\partial h_t}{\partial w_h}$$

$$\frac{\partial h_t}{\partial w_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial w_h}$$

پس انتشار در طول زمان

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

$$o_t = g(h_t, w_o)$$

$$L(x_1, \dots, x_T, y_1, \dots, y_T, w_h, w_o) = \frac{1}{T} \sum_{t=1}^T l(y_t, o_t)$$

$$\frac{\partial L}{\partial w_h} = \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial w_h} = \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial o_t} \frac{\partial g(h_t, w_o)}{\partial w_h} = \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial o_t} \frac{\partial g(h_t, w_o)}{\partial h_t} \frac{\partial h_t}{\partial w_h}$$

$$\frac{\partial h_t}{\partial w_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h}$$

پس انتشار در طول زمان

- زنجیره محاسبه گرادیان می تواند خیلی طولانی شود
- محاسبات کامل: می تواند خیلی کند باشد و گرادیان ها می توانند منفجر شوند (مشابه با اثر پروانه ای)
- کوتاه کردن گام های زمانی: در محاسبات فقط از τ زمان گذشته برای تقریب گرادیان استفاده می شود
 - مدل بر تأثیر کوتاه مدت بجای بلند مدت تمرکز می کند و تخمین را به سمت مدل های ساده تر و پایدارتر سوق می دهد
- کوتاه کردن تصادفی: طول دنباله در محاسبه گرادیان به صورت تصادفی تغییر می کند و ضریب می گیرد تا امید ریاضی گرادیان برابر با مقدار درست باشد

$$\frac{\partial h_t}{\partial w_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h}$$

پس انتشار در طول زمان

- کوتاه کردن تصادفی: طول دنباله در محاسبه گرادیان به صورت تصادفی تغییر می کند و ضریب می گیرد تا امید ریاضی گرادیان برابر با مقدار درست باشد
- محاسبات در $\xi_t = 0$ متوقف می شود، طول دنباله ها متغیر خواهد بود و دنباله های طولانی هم به ندرت استفاده خواهند شد

$$0 \leq \pi_t \leq 1 \quad \begin{cases} P(\xi_t = 0) = 1 - \pi_t \\ P(\xi_t = \pi_t^{-1}) = \pi_t \end{cases} \Rightarrow E[\xi_t] = 1$$

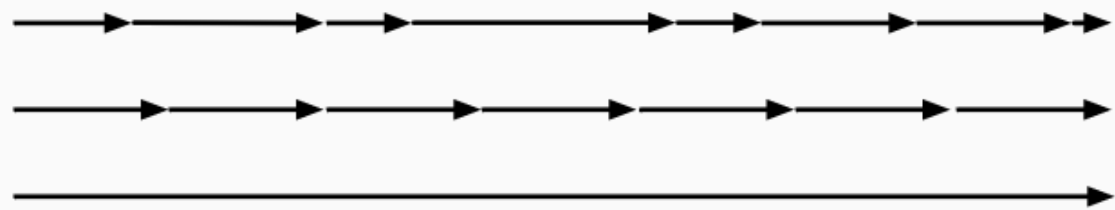
$$z_t = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \xi_t \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial w_h} \Rightarrow E[z_t] = \frac{\partial h_t}{\partial w_h}$$

$$\frac{\partial h_t}{\partial w_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h}$$

پس انتشار در طول زمان

- زنجیره محاسبه گرادیان می تواند خیلی طولانی شود
- محاسبات کامل: می تواند خیلی کند باشد و گرادیان ها می توانند منفجر شوند (مشابه با اثر پروانه ای)
- کوتاه کردن گام های زمانی: در محاسبات فقط از τ زمان گذشته برای تقریب گرادیان استفاده می شود
 - مدل بر تأثیر کوتاه مدت بجای بلند مدت تمرکز می کند و تخمین را به سمت مدل های ساده تر و پایدارتر سوق می دهد
- کوتاه کردن تصادفی: طول دنباله در محاسبه گرادیان به صورت تصادفی تغییر می کند و ضریب می گیرد تا امید ریاضی گرادیان برابر با مقدار درست باشد

the time machine by h g well



BPTT برای SimpleRNN

- برای سادگی از بایاس و تابع غیرخطی استفاده نمی‌کنیم

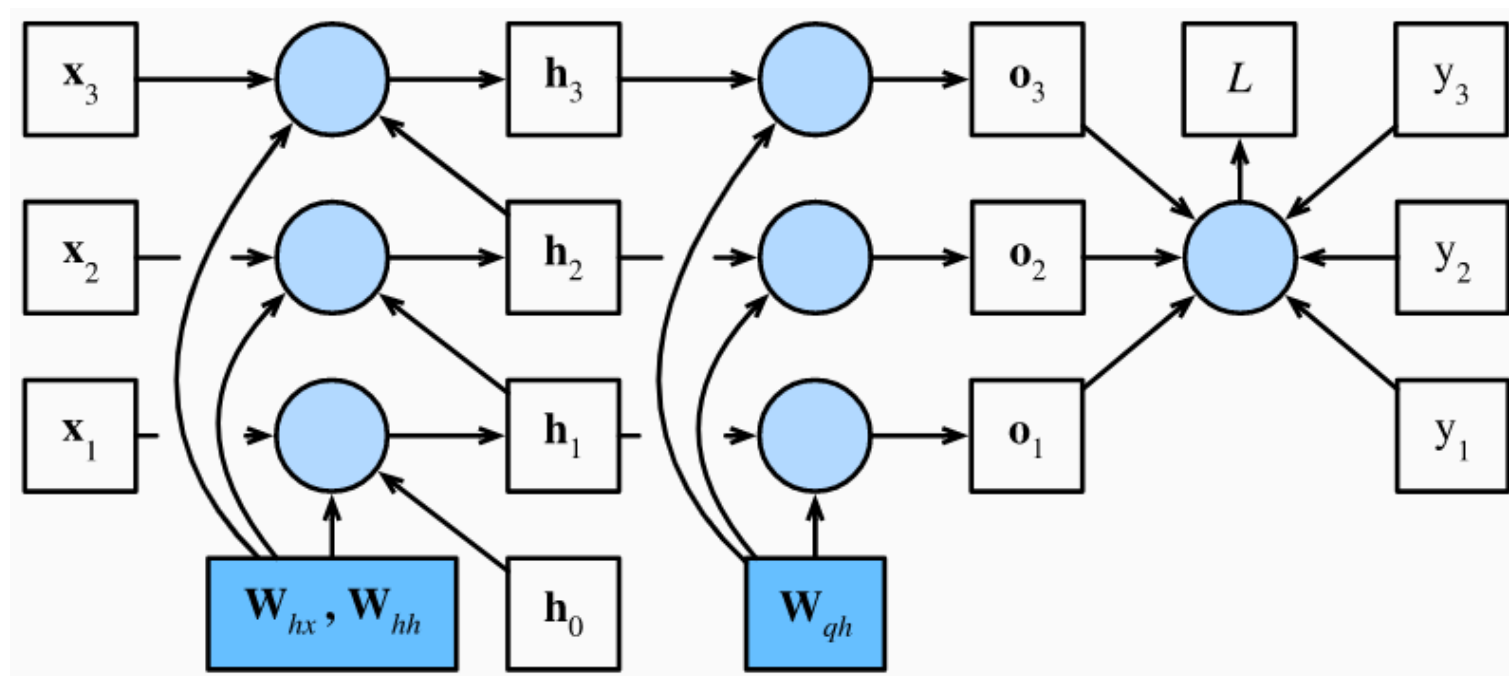
$$\mathbf{h}_t = \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1}$$

$$\mathbf{o}_t = \mathbf{W}_{qh}\mathbf{h}_t$$

$$L = \frac{1}{T} \sum_{t=1}^T l(\mathbf{y}_t, \mathbf{o}_t)$$

$$\frac{\partial L}{\partial \mathbf{o}_t} = \frac{\partial l(\mathbf{y}_t, \mathbf{o}_t)}{T \partial \mathbf{o}_t} \in \mathbb{R}^q$$

$$\frac{\partial L}{\partial \mathbf{W}_{qh}} = \sum_{t=1}^T \frac{\partial L}{\partial \mathbf{o}_t} \mathbf{h}_t^T$$



BPTT برای SimpleRNN

- برای سادگی از بایاس و تابع غیرخطی استفاده نمی‌کنیم

$$\mathbf{h}_t = \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1}$$

$$\mathbf{o}_t = \mathbf{W}_{qh}\mathbf{h}_t$$

$$L = \frac{1}{T} \sum_{t=1}^T l(\mathbf{y}_t, \mathbf{o}_t)$$

$$\frac{\partial L}{\partial \mathbf{o}_t} = \frac{\partial l(\mathbf{y}_t, \mathbf{o}_t)}{T \partial \mathbf{o}_t} \in \mathbb{R}^q$$

$$\frac{\partial L}{\partial \mathbf{W}_{qh}} = \sum_{t=1}^T \frac{\partial L}{\partial \mathbf{o}_t} \mathbf{h}_t^T$$

$$\frac{\partial L}{\partial \mathbf{h}_T} = \text{prod} \left(\frac{\partial L}{\partial \mathbf{o}_t}, \frac{\partial \mathbf{o}_t}{\partial \mathbf{h}_T} \right) = \mathbf{W}_{qh}^T \frac{\partial L}{\partial \mathbf{o}_t}$$

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{h}_t} &= \text{prod} \left(\frac{\partial L}{\partial \mathbf{o}_t}, \frac{\partial \mathbf{o}_t}{\partial \mathbf{h}_t} \right) + \text{prod} \left(\frac{\partial L}{\partial \mathbf{h}_{t+1}}, \frac{\partial \mathbf{h}_{t+1}}{\partial \mathbf{h}_t} \right) \\ &= \mathbf{W}_{qh}^T \frac{\partial L}{\partial \mathbf{o}_t} + \mathbf{W}_{hh}^T \frac{\partial L}{\partial \mathbf{h}_{t+1}} \end{aligned}$$

$$\frac{\partial L}{\partial \mathbf{h}_t} = \sum_{i=t}^T (\mathbf{W}_{hh}^T)^{T-i} \mathbf{W}_{qh}^T \frac{\partial L}{\partial \mathbf{o}_{T+t-i}}$$

BPTT برای SimpleRNN

$$\mathbf{h}_t = \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1}$$

$$\mathbf{o}_t = \mathbf{W}_{qh}\mathbf{h}_t$$

$$L = \frac{1}{T} \sum_{t=1}^T l(\mathbf{y}_t, \mathbf{o}_t)$$

$$\frac{\partial L}{\partial \mathbf{h}_t} = \sum_{i=t}^T (\mathbf{w}_{hh}^T)^{T-i} \mathbf{w}_{qh}^T \frac{\partial L}{\partial \mathbf{o}_{T+t-i}}$$

$$\frac{\partial L}{\partial \mathbf{W}_{hx}} = \sum_{i=1}^T \frac{\partial L}{\partial \mathbf{h}_t} \mathbf{x}_t^T \quad \frac{\partial L}{\partial \mathbf{W}_{hh}} = \sum_{i=1}^T \frac{\partial L}{\partial \mathbf{h}_t} \mathbf{h}_{t-1}^T$$

- برای سادگی از بایاس و تابع غیرخطی استفاده نمی‌کنیم
- زمانی که دنباله طولانی باشد، \mathbf{W}_{hh}^T به توان اعداد بزرگ می‌رسد و اگر مقادیر ویژه آن بزرگتر از ۱ باشند دچار انفجار گرادیان و اگر کوچکتر از ۱ باشند دچار محوشدگی گرادیان خواهیم شد
- با کوتاه کردن گام‌های زمانی می‌توان این مشکل را تا حدی جبران کرد

محوشدگی و انفجار گرادیان

- در شبکه‌های بازگشتی برای دنباله‌های طولانی این مشکل بسیار جدی است
- می‌توان با برش گرادیان از انفجار گرادیان جلوگیری کرد

→ Forward Propagation →

$$I/p \xrightarrow{w_1} (b_1) \xrightarrow{w_2} (b_2) \xrightarrow{w_3} (b_3) \xrightarrow{w_4} (b_4) \xrightarrow{z_4} J$$

$$w_1 a_0 + b_1 = z_1 \quad a_1 = \text{sigmoid}(z_1)$$

$$w_2 z_1 + b_2 = z_2 \quad a_2 = \text{sigmoid}(a_1 w_2 + b_2)$$

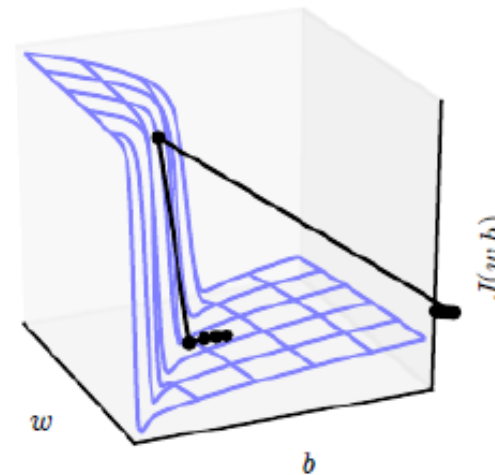
$$w_3 z_2 + b_3 = z_3 \quad a_3 = \text{sigmoid}(a_2 w_3 + b_3)$$

$$w_4 z_3 + b_4 = z_4 \quad a_4 = \text{sigmoid}(a_3 w_4 + b_4)$$

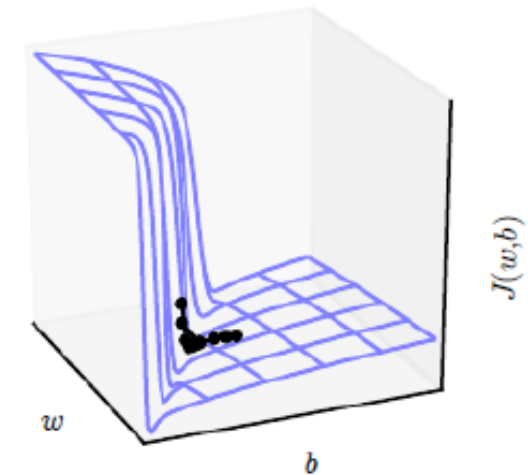
$$\frac{\partial J}{\partial b_1} = \frac{\partial J}{\partial a_4} \times \sigma'(z_4) w_4 \times \sigma'(z_3) w_3 \times \sigma'(z_2) w_2 \times \sigma'(z_1)$$

(Back Propagation Equation)

Without clipping



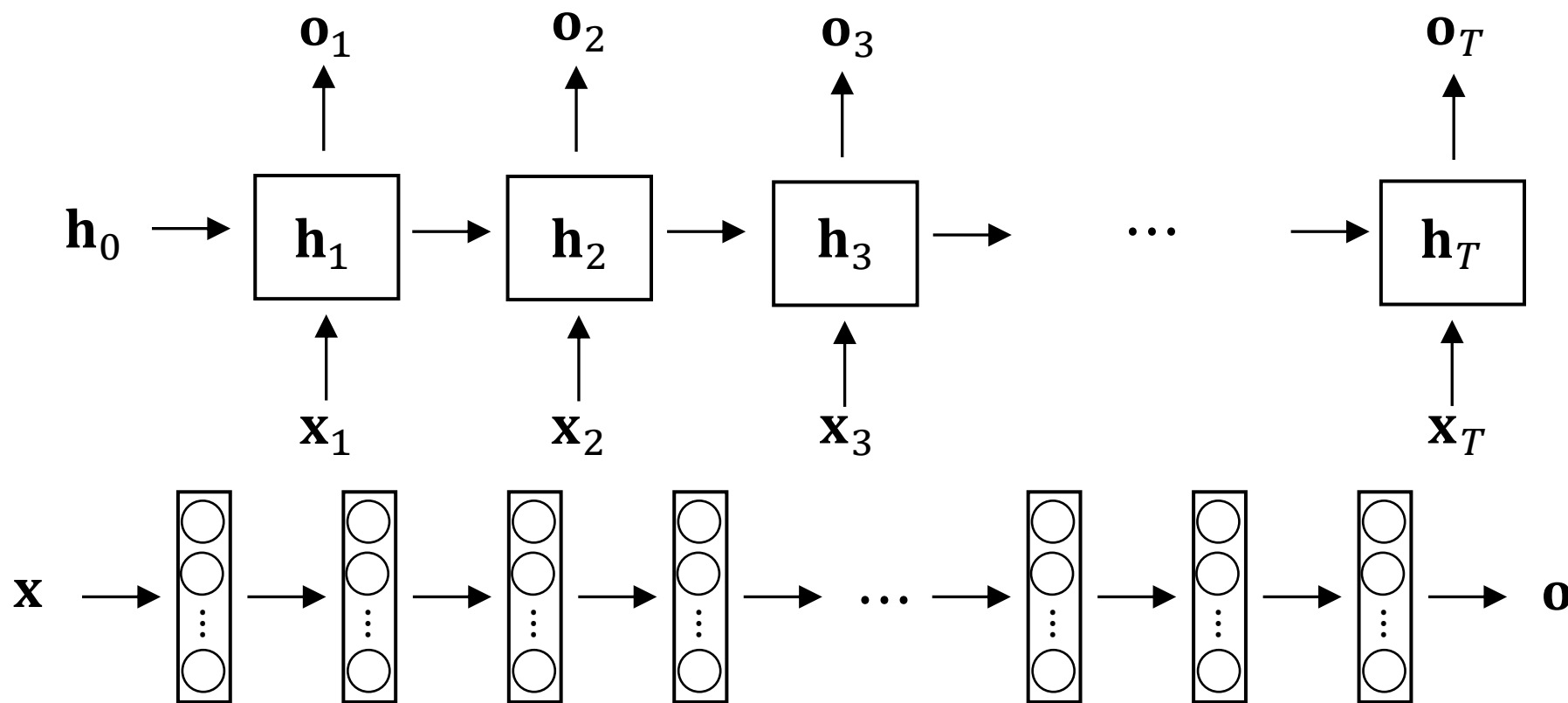
With clipping



محوشدگی گرادیان در RNNها

The **cat**, which already ate ..., **was** full.

The **cats**, which already ate ..., **were** full.

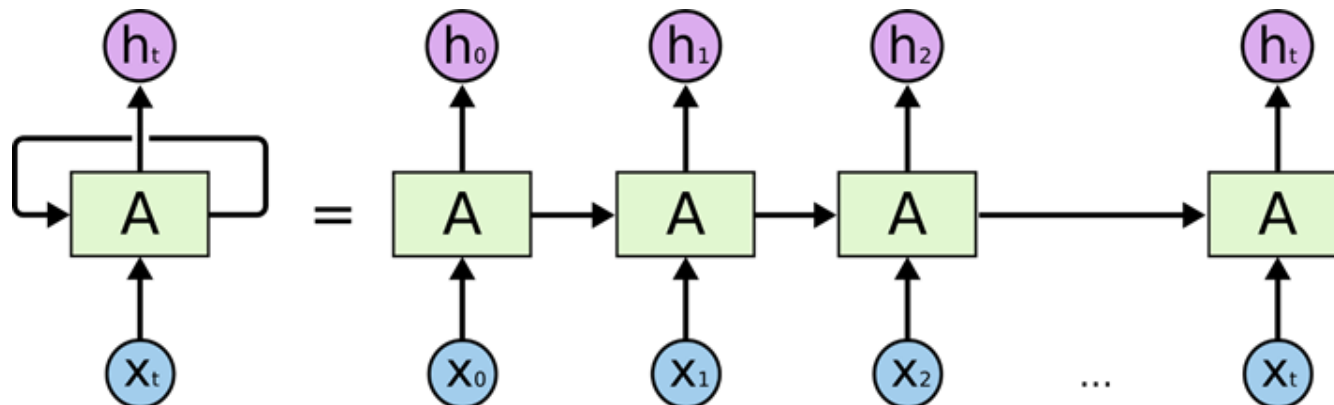


Gated RNNs

- شبکه‌های عصبی بازگشتی از حافظه کوتاه مدت رنج می‌برند
- RNN ممکن است اطلاعات مهم ابتدایی را نادیده بگیرد
- Gated RNNs مبتنی بر ایده ایجاد مسیرهایی در طول زمان هستند که از محوشدگی یا انفجار گرادیان جلوگیری می‌کنند

The **cat**, which already ate ..., **was** full.

The **cats**, which already ate ..., **were** full.



Customers Review 2,491

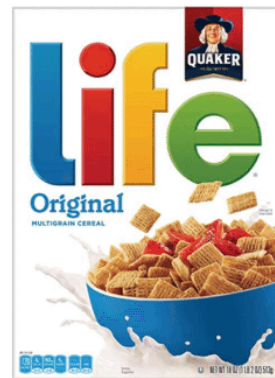


Thanos

September 2018

Verified Purchase

Amazing! This box of cereal gave me a perfectly balanced breakfast, as all things should be. I only ate half of it but will definitely be buying again!



A Box of Cereal
\$3.99

Gated Recurrent Units

GRU (simplified)

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{b}_h)$$

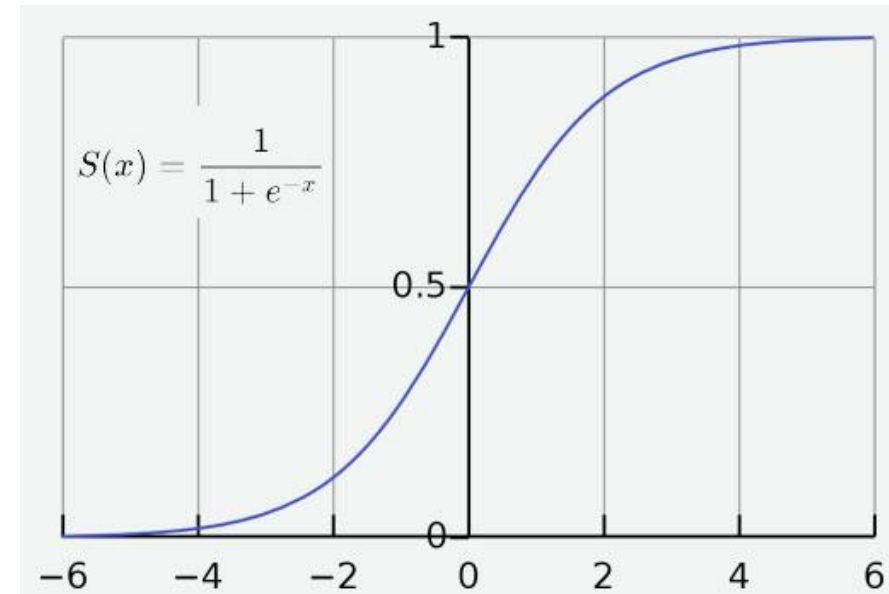
$$\mathbf{h}_t = \mathbf{u}_t \cdot \tilde{\mathbf{h}}_t + (1 - \mathbf{u}_t) \cdot \mathbf{h}_{t-1}$$

$$\mathbf{u}_t = \sigma(\mathbf{W}_{hu}\mathbf{h}_{t-1} + \mathbf{W}_{xu}\mathbf{x}_t + \mathbf{b}_u)$$

The cat, which already ate ..., was full.

Simple RNN

$$\mathbf{h}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{b}_h)$$



Gated Recurrent Units

GRU (simplified)

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{b}_h)$$

$$\mathbf{h}_t = \mathbf{u}_t \cdot \tilde{\mathbf{h}}_t + (1 - \mathbf{u}_t) \cdot \mathbf{h}_{t-1}$$

$$\mathbf{u}_t = \sigma(\mathbf{W}_{uh}\mathbf{h}_{t-1} + \mathbf{W}_{ux}\mathbf{x}_t + \mathbf{b}_u)$$

GRU

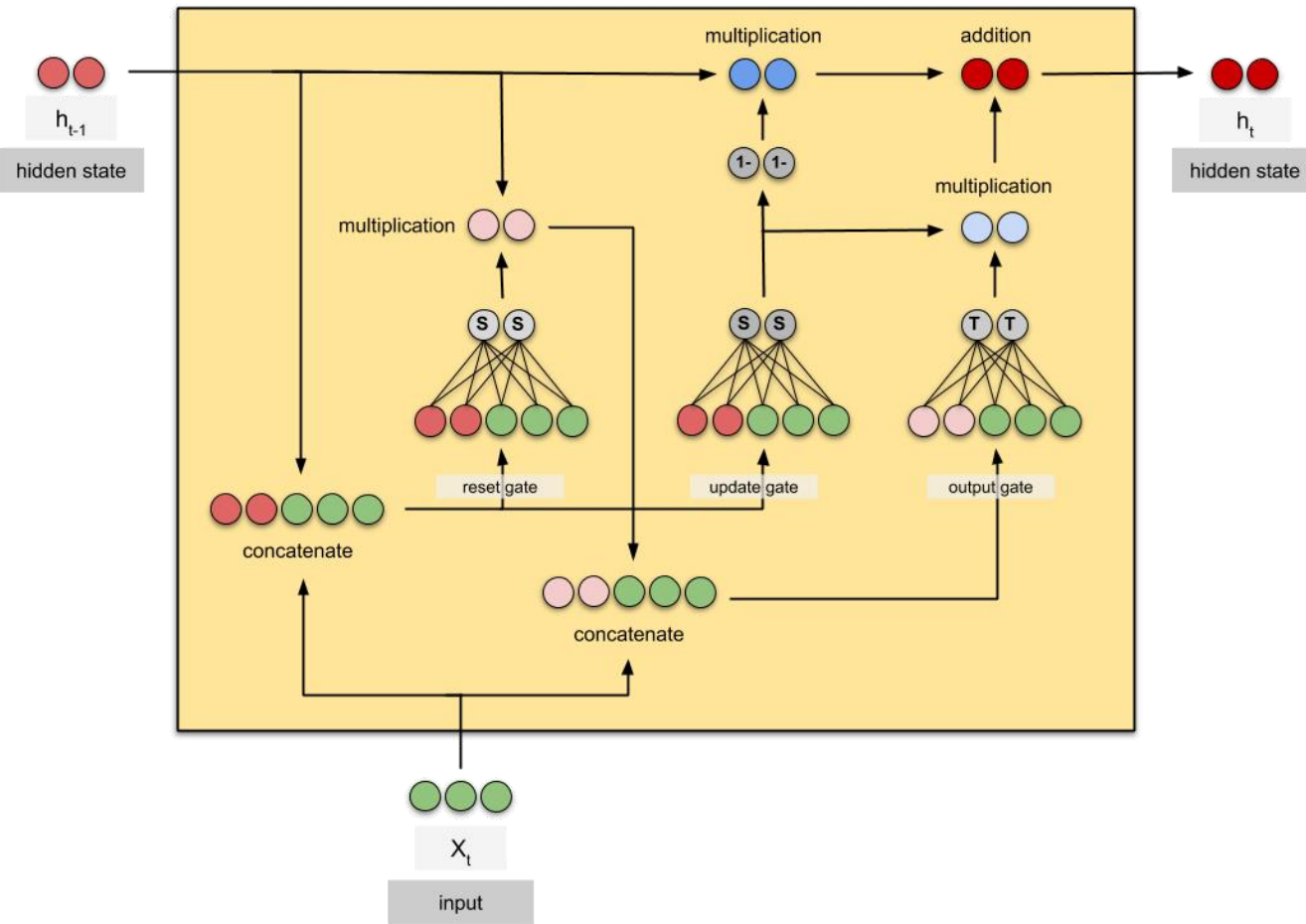
$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_{hh}(\mathbf{r}_t \cdot \mathbf{h}_{t-1}) + \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{b}_h)$$

$$\mathbf{h}_t = \mathbf{u}_t \cdot \tilde{\mathbf{h}}_t + (1 - \mathbf{u}_t) \cdot \mathbf{h}_{t-1}$$

$$\mathbf{u}_t = \sigma(\mathbf{W}_{uh}\mathbf{h}_{t-1} + \mathbf{W}_{ux}\mathbf{x}_t + \mathbf{b}_u)$$

$$\mathbf{r}_t = \sigma(\mathbf{W}_{rh}\mathbf{h}_{t-1} + \mathbf{W}_{rx}\mathbf{x}_t + \mathbf{b}_r)$$

GRU



$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_{hh}(\mathbf{r}_t \cdot \mathbf{h}_{t-1}) + \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{b}_h)$$

$$\mathbf{h}_t = \mathbf{u}_t \cdot \tilde{\mathbf{h}}_t + (1 - \mathbf{u}_t) \cdot \mathbf{h}_{t-1}$$

$$\mathbf{u}_t = \sigma(\mathbf{W}_{uh}\mathbf{h}_{t-1} + \mathbf{W}_{ux}\mathbf{x}_t + \mathbf{b}_u)$$

$$\mathbf{r}_t = \sigma(\mathbf{W}_{rh}\mathbf{h}_{t-1} + \mathbf{W}_{rx}\mathbf{x}_t + \mathbf{b}_r)$$

LSTM

$$\mathbf{h}_t = \mathbf{o}_t \cdot \tanh(\mathbf{c}_t)$$

$$\tilde{\mathbf{c}}_t = \tanh(\mathbf{W}_{ch}\mathbf{h}_{t-1} + \mathbf{W}_{cx}\mathbf{x}_t + \mathbf{b}_c)$$

$$\mathbf{c}_t = \mathbf{u}_t \cdot \tilde{\mathbf{c}}_t + \mathbf{f}_t \cdot \mathbf{c}_{t-1}$$

$$\mathbf{u}_t = \sigma(\mathbf{W}_{uh}\mathbf{h}_{t-1} + \mathbf{W}_{ux}\mathbf{x}_t + \mathbf{b}_u)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{fh}\mathbf{h}_{t-1} + \mathbf{W}_{fx}\mathbf{x}_t + \mathbf{b}_f)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{oh}\mathbf{h}_{t-1} + \mathbf{W}_{ox}\mathbf{x}_t + \mathbf{b}_o)$$

Long Short-Term Memory

GRU

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_{hh}(\mathbf{r}_t \cdot \mathbf{h}_{t-1}) + \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{b}_h)$$

$$\mathbf{h}_t = \mathbf{u}_t \cdot \tilde{\mathbf{h}}_t + (1 - \mathbf{u}_t) \cdot \mathbf{h}_{t-1}$$

$$\mathbf{u}_t = \sigma(\mathbf{W}_{uh}\mathbf{h}_{t-1} + \mathbf{W}_{ux}\mathbf{x}_t + \mathbf{b}_u)$$

$$\mathbf{r}_t = \sigma(\mathbf{W}_{rh}\mathbf{h}_{t-1} + \mathbf{W}_{rx}\mathbf{x}_t + \mathbf{b}_r)$$

LSTM

$$\mathbf{h}_t = \mathbf{o}_t \cdot \tanh(\mathbf{c}_t)$$

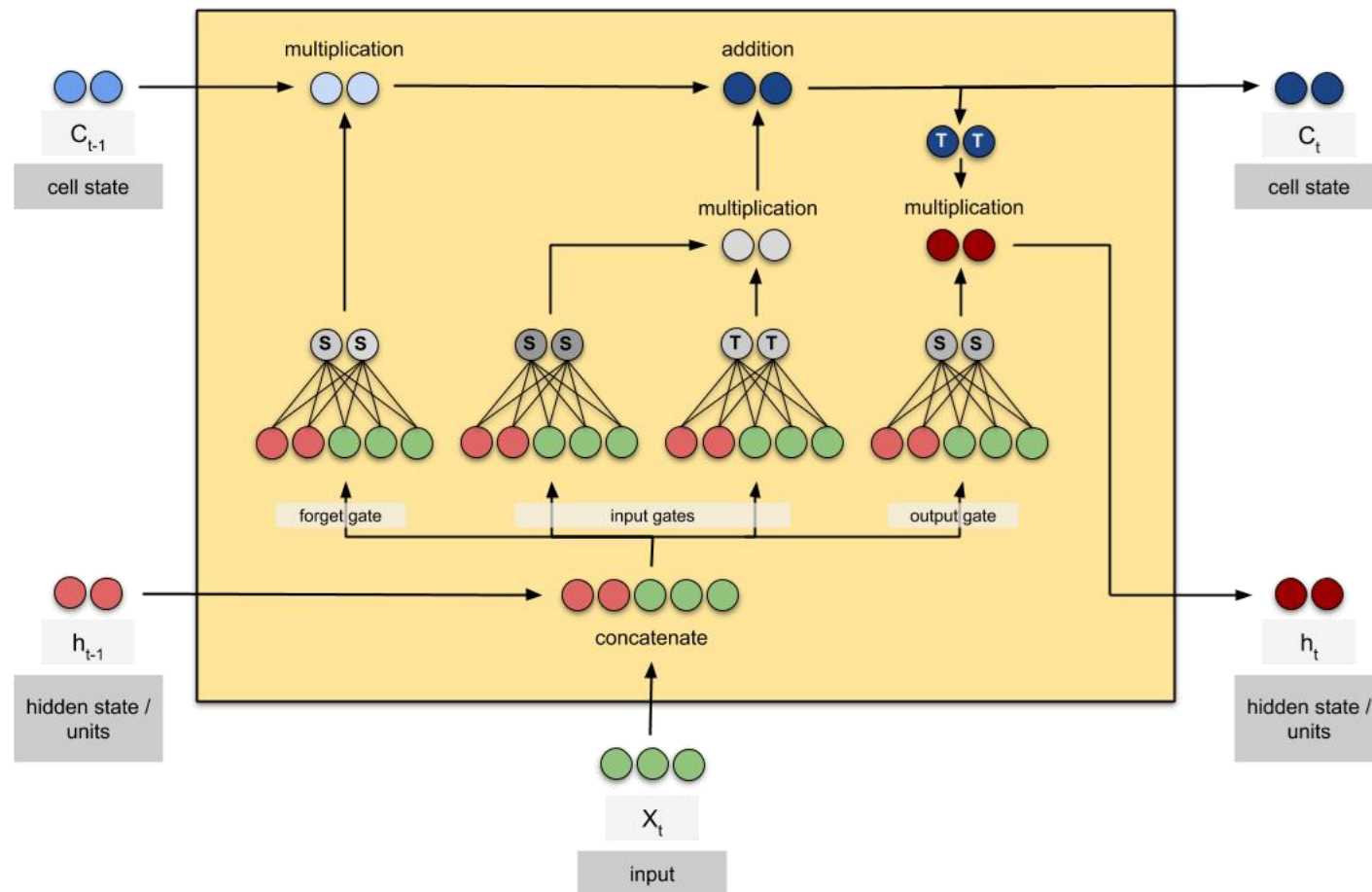
$$\tilde{\mathbf{c}}_t = \tanh(\mathbf{W}_{ch}\mathbf{h}_{t-1} + \mathbf{W}_{cx}\mathbf{x}_t + \mathbf{b}_c)$$

$$\mathbf{c}_t = \mathbf{u}_t \cdot \tilde{\mathbf{c}}_t + \mathbf{f}_t \cdot \mathbf{c}_{t-1}$$

$$\mathbf{u}_t = \sigma(\mathbf{W}_{uh}\mathbf{h}_{t-1} + \mathbf{W}_{ux}\mathbf{x}_t + \mathbf{b}_u)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{fh}\mathbf{h}_{t-1} + \mathbf{W}_{fx}\mathbf{x}_t + \mathbf{b}_f)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{oh}\mathbf{h}_{t-1} + \mathbf{W}_{ox}\mathbf{x}_t + \mathbf{b}_o)$$



LSTM

```
def LSTMCELL(prev_ct, prev_ht, input):  
    combine = prev_ht + input  
    ft = forget_layer(combine)  
    candidate = candidate_layer(combine)  
    it = input_layer(combine)  
    Ct = prev_ct * ft + candidate * it  
    ot = output_layer(combine)  
    ht = ot * tanh(Ct)  
    return ht, Ct
```

```
ct = [0, 0, 0]  
ht = [0, 0, 0]
```

```
for input in inputs:  
    ct, ht = LSTMCELL(ct, ht, input)
```

