

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



یادگیری عمیق

جلسه ۱۹ و ۲۰

شبکه‌های عصبی بازگشتی

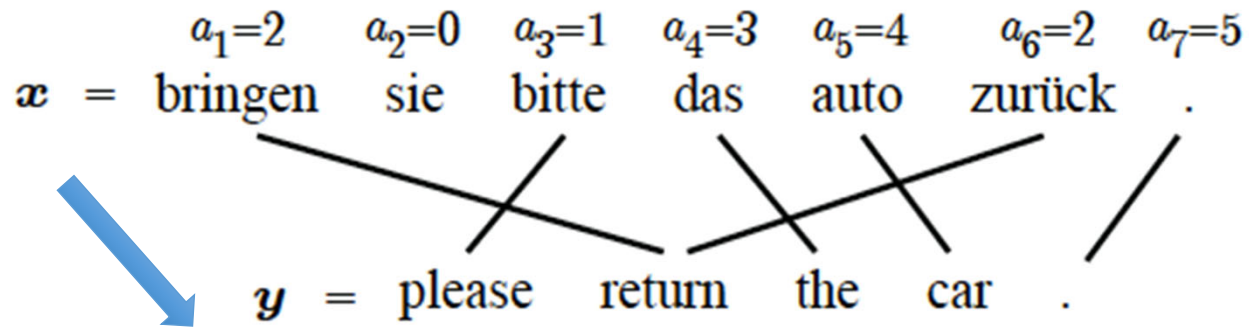
Recurrent Neural Networks (RNN)

کاظم فولادی قلعه
دانشکده مهندسی، پردیس فارابی
دانشگاه تهران

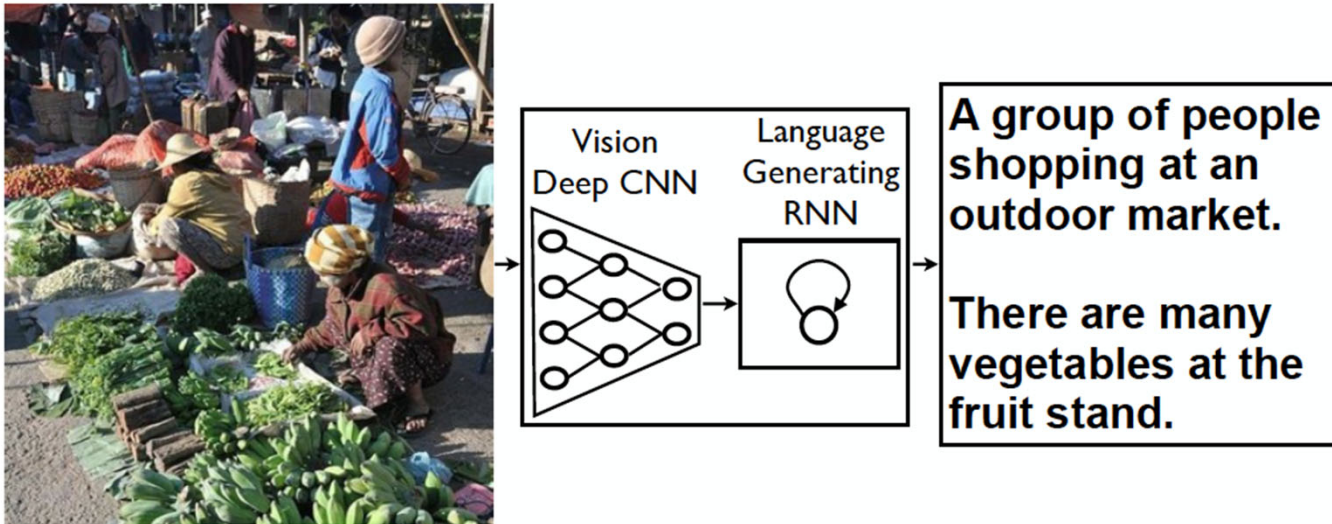
<http://courses.fouladi.ir/deep>

Sequences are everywhere...

Foreign Minister. → FOREIGN MINISTER.

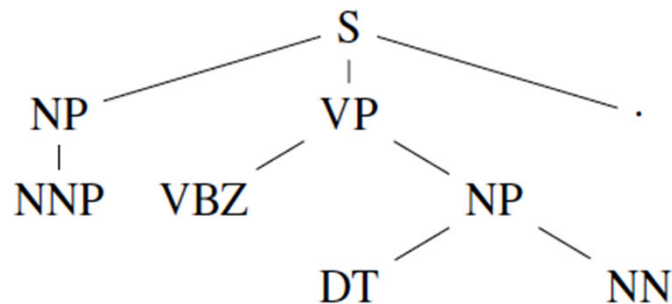


Even where you might not expect a sequence...



John has a dog .

→



John has a dog .

→

(S (NP NNP)_{NP} (VP VBZ (NP DT NN)_{NP})_{VP} .)_S

شبکه‌های عصبی بازگشتی



مقدمات

دنباله‌ها

SEQUENCES

دنباله: داده‌های بعدی به داده‌های قبلی وابستگی دارند.

هدف: پیش‌بینی اینکه چه چیزی بعداً می‌آید؟

$$\Pr(x) = \prod_i \Pr(x_i | x_1, \dots, x_{i-1})$$

در نظر گرفتن تکه‌های x_i ←

تعداد پارامتر کمتر + مدل‌سازی ساده‌تر + امکان تعمیم به طول دلخواه

معمولاً به جای طول دلخواه، یک قاب (frame) به طول T برداشته می‌شود:

$$\Pr(x) = \prod_i \Pr(x_i | x_{i-T}, \dots, x_{i-1})$$

دنباله‌ها

ویژگی‌ها

SEQUENCES

- داده‌های داخل یک دنباله، iid نیستند.
- یعنی مستقل و دارای توزیع یکسان نیستند.
- **کلمه‌ی بعدی**، وابسته به **کلمه‌های قبلی** است.
- به صورت ایده‌آل، به همه‌ی کلمه‌های قبلی وابسته است.
- برای تحلیل دنباله نیازمند **مضمون (context)** و نیز **حافظه (memory)** هستیم.

مدل کردن مضمون و حافظه

مثال (۱ از ۲)

MODELLING CONTEXT AND MEMORY

I am Bond, James

McGuire
Bond
tired
am
!

کلمه‌ی بعدی کدام باید باشد؟

مدل کردن مضمون و حافظه

مثال (۲ از ۲)

MODELLING CONTEXT AND MEMORY

I am Bond, James

Bond

McGuire
Bond
tired
am
!

مدل کردن مضمون و حافظه

بردارهای تک-داغ

ONE-HOT VECTORS $x_i \equiv$ one-hot vector

بردارهای تک-داغ که همه‌ی عناصر آن صفر است، غیر از یک مقدار 1 برای بعد فعال آن

برای مثال: اگر ۱۲ کلمه در یک دنباله داشته باشیم، ۱۲ بردار تک-داغ خواهیم داشت.

VocabularyOne-Hot Vectors

	1	0	0	0	0
am	am 0	am 1	am 0	am 0	am 0
Bond	Bond 0	Bond 0	Bond 1	Bond 0	Bond 0
James	James 0	James 0	James 0	James 1	James 0
tired	tired 0	tired 0	tired 0	tired 0	tired 1
,	,	,	,	,	,
McGuire	McGuire 0	McGuire 0	McGuire 0	McGuire 0	McGuire 0
!	! 0	! 0	! 0	! 0	! 0

پس از ایجاد بردارهای تک-داغ، یک روش جاسازی (مثل Word2Vec یا GloVE) اعمال می‌شود.

مدل زبان

LANGUAGE MODEL

مدل زبان احتمال هر دنباله از کلمات آن زبان را بیان می‌کند.

مدل زبان
Language Model

$$P(w_1, w_2, \dots, w_n) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2) \times \dots \times P(w_n|w_1, w_2, \dots, w_{n-1})$$

کاربردها:

$p(\text{the cat is small}) > p(\text{small the is cat})$
مورد استفاده در ترجمه‌ی ماشینی

ترتیب‌دهی کلمات
Word Ordering

$p(\text{walking home after school}) > p(\text{walking house after school})$
 $p(\text{he likes apple}) > p(\text{he licks apple})$
مورد استفاده در بازشناسی گفتار / تولید گفتار

گزینش کلمات
Word Choice

مدل کردن مضمون و حافظه

حافظه

MEMORY

حافظه، بازنمایی گذشته است.

اطلاعات در گام زمانی t با استفاده از پارامتر θ بر روی یک فضای نهفته c_t افکنده می شود.از اطلاعات افکنده شده، از لحظه t در لحظه $t + 1$ استفاده می شود:

$$c_{t+1} = h(x_{t+1}, c_t; \theta)$$

پارامتر بازگشتی θ برای تمام گام های زمانی به اشتراک گذاشته می شود:

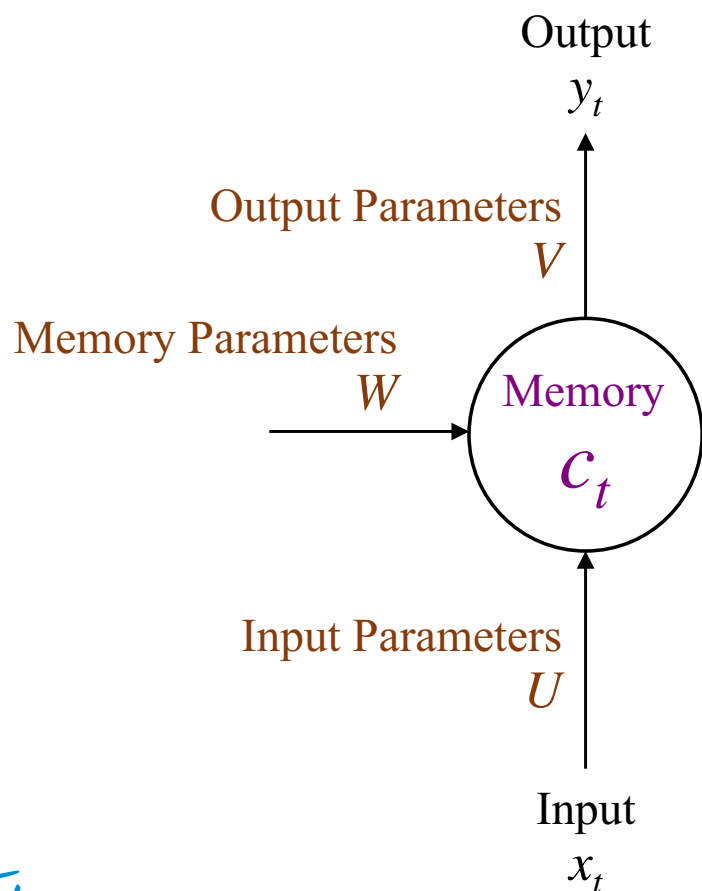
$$c_{t+1} = h(x_{t+1}, h(x_t, h(x_{t-1}, \dots h(x_1, c_0; \theta); \theta); \theta); \theta))$$

مدل کردن مضمون و حافظه

مدل کردن حافظه در قالب یک گراف

MEMORY AS A GRAPH

ساده‌ترین مدل

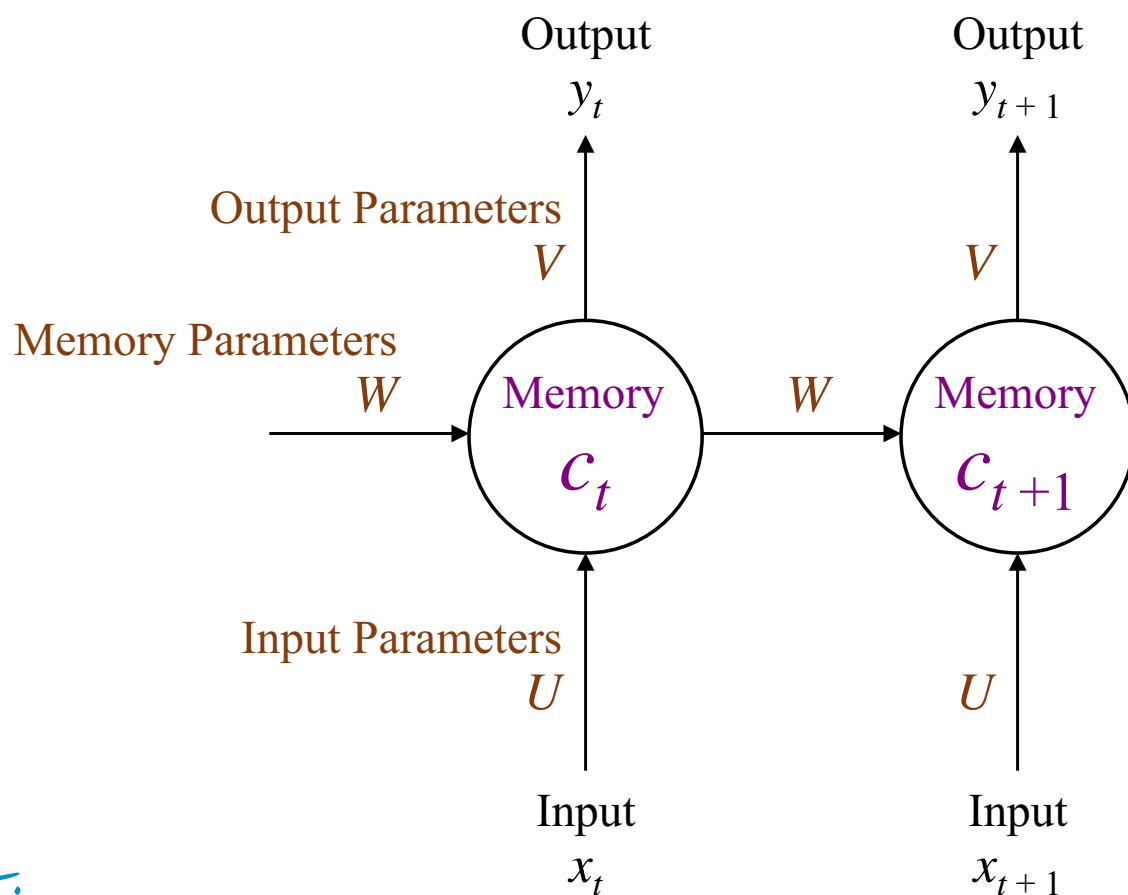


مدل کردن مضمون و حافظه

مدل کردن حافظه در قالب یک گراف

MEMORY AS A GRAPH

ساده‌ترین مدل

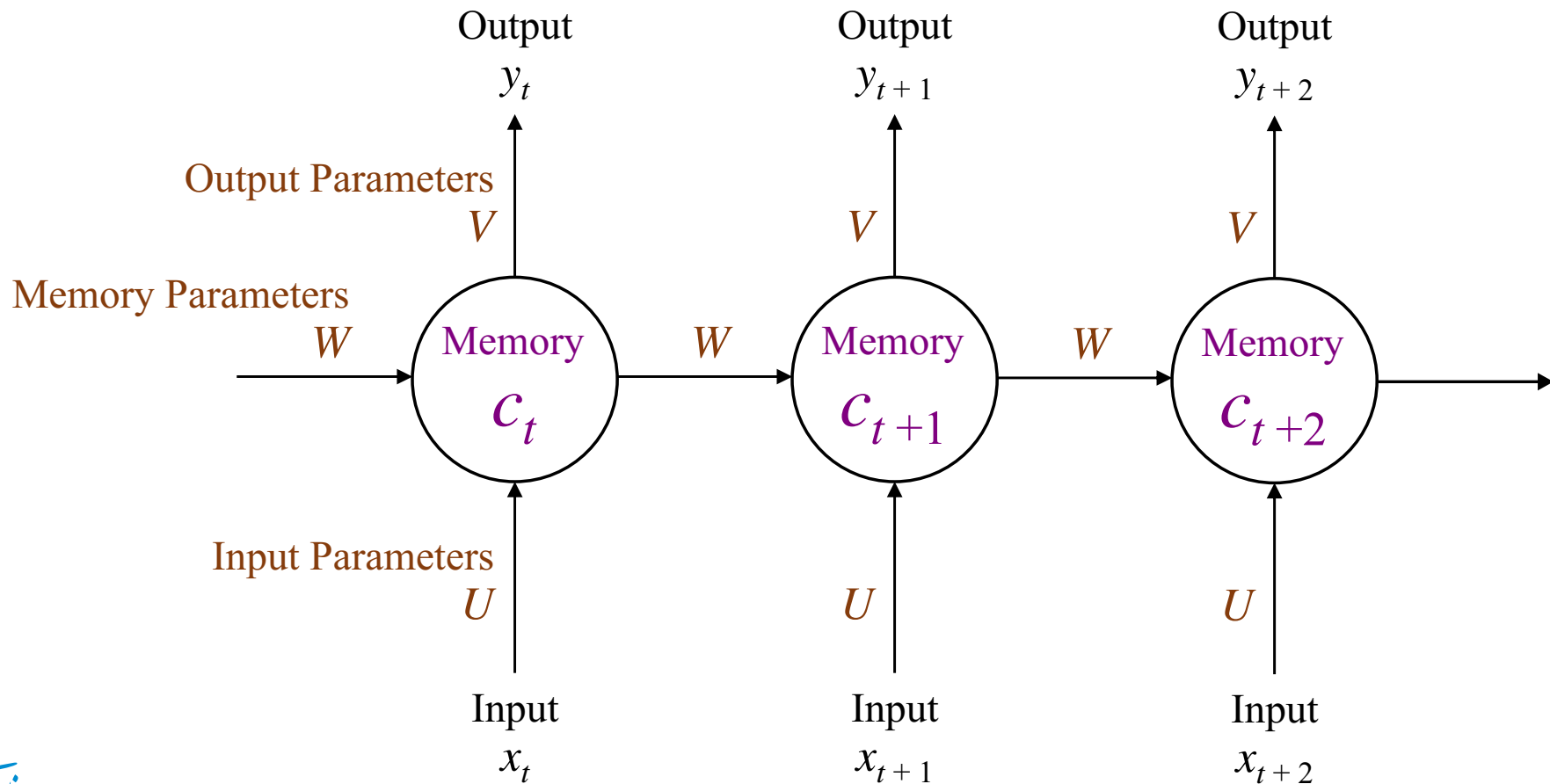


مدل کردن مضمون و حافظه

مدل کردن حافظه در قالب یک گراف

MEMORY AS A GRAPH

ساده‌ترین مدل

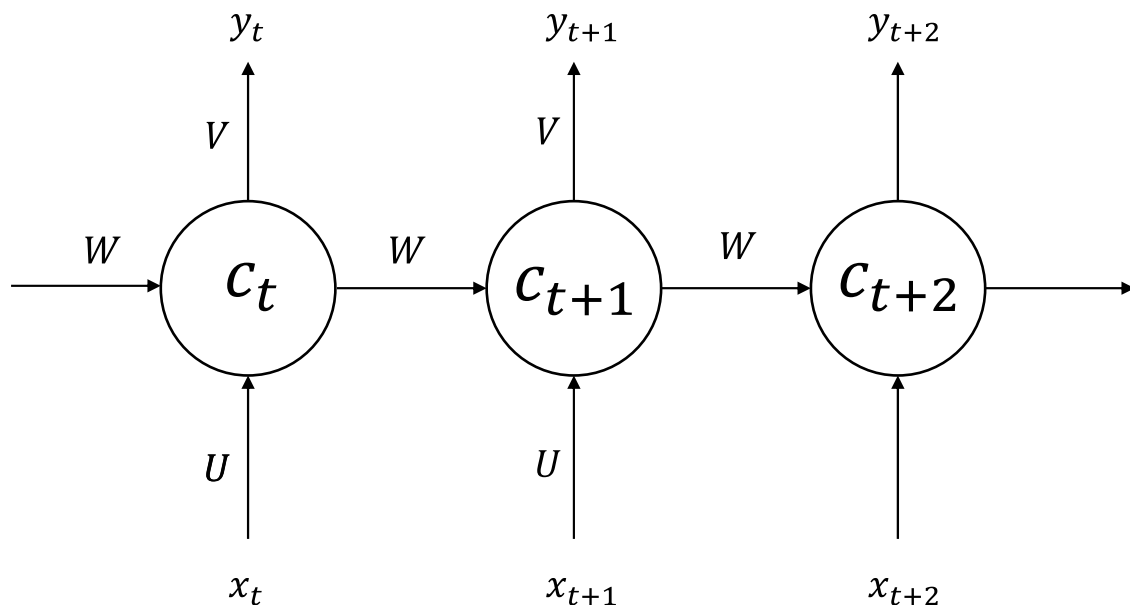


مدل کردن مضمون و حافظه

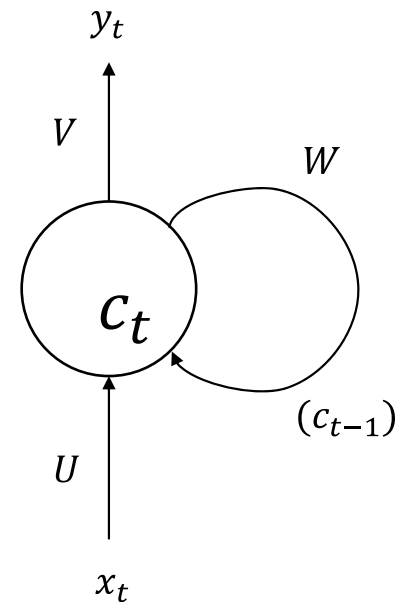
تا کردن حافظه

FOLDING THE MEMORY

شبکه‌ی باز شده / تا نشده
Unrolled / Unfolded Network



شبکه‌ی تا شده
Folded Network



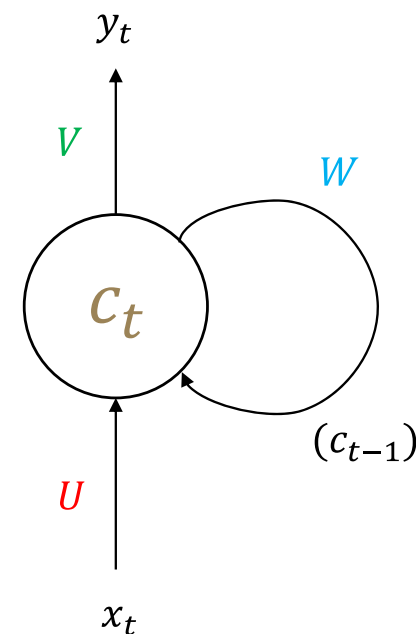
شبکه‌ی عصبی بازگشتی

RECURRENT NEURAL NETWORK (RNN)

شبکه‌ی عصبی بازگشتی، با استفاده از دو معادله تعریف می‌شود:

$$c_t = \tanh(U x_t + W c_{t-1})$$

$$y_t = \text{softmax}(V c_t)$$



شبکه‌ی عصبی بازگشتی

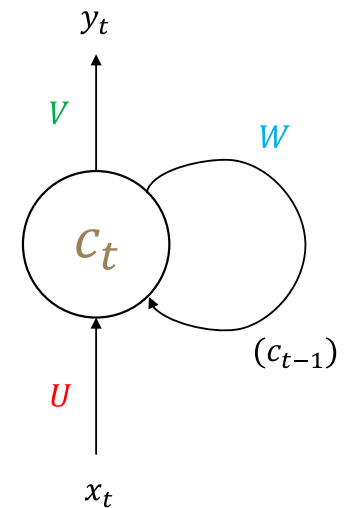
مثال

RECURRENT NEURAL NETWORK (RNN)

- Vocabulary of 5 words
- A memory of 3 units [Hyperparameter that we choose like layer size]
 - c_t : $[3 \times 1]$, W : $[3 \times 3]$
- An input projection of 3 dimensions
 - U : $[3 \times 5]$
- An output projections of 10 dimensions
 - V : $[10 \times 3]$

$$c_t = \tanh(U x_t + W c_{t-1})$$

$$y_t = \text{softmax}(V c_t)$$



$$U \cdot x_{t=4} = \begin{bmatrix} 0.1 & -0.3 & 1.2 & 0.6 & -0.8 \\ -0.2 & 0.4 & 0.5 & 0.9 & -0.1 \\ -0.1 & 0.2 & -0.7 & -0.8 & 0.3 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.6 \\ 0.9 \\ -0.8 \end{bmatrix} = U^{(4)}$$

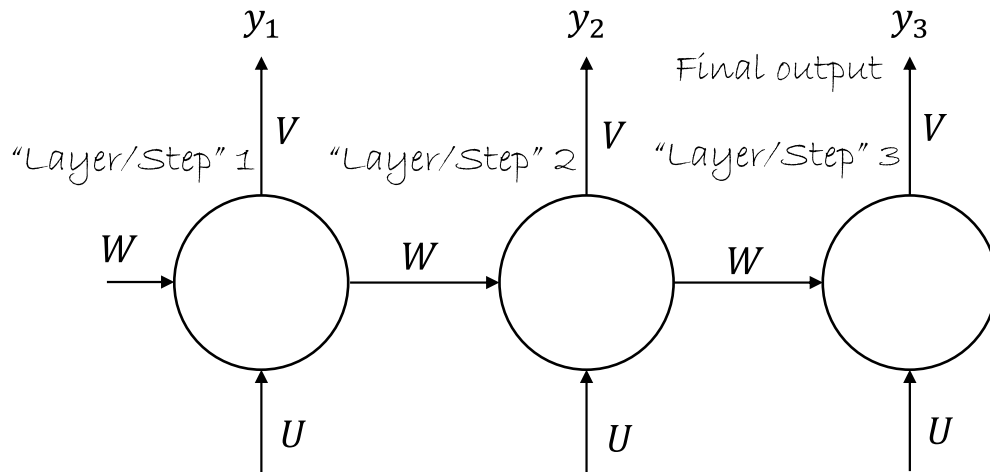
شبکه‌ی عصبی بازگشتی

مقایسه‌ی شبکه‌ی بسته شده با شبکه‌ی چندلایه

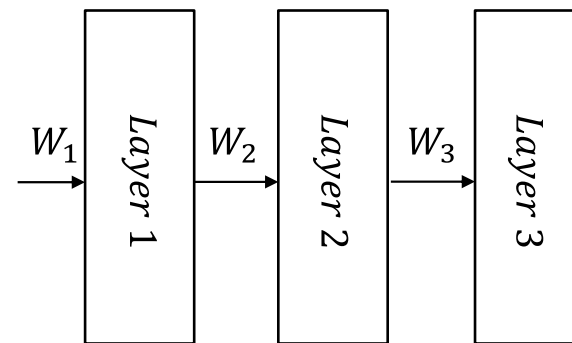
ROLLED NETWORK VS. MULTI-LAYER NETWORK

تفاوت‌ها:

- در شبکه‌ی بسته شده، گام‌ها به جای لایه‌ها نقش بازی می‌کنند.
- در شبکه‌ی بسته شده، پارامترهای گام مشترک است؛ در حالی که در شبکه‌ی چندلایه، پارامترها متفاوت است.



3-gram unrolled recurrent network



3-layer neural network

شبکه‌ی عصبی بازگشتی

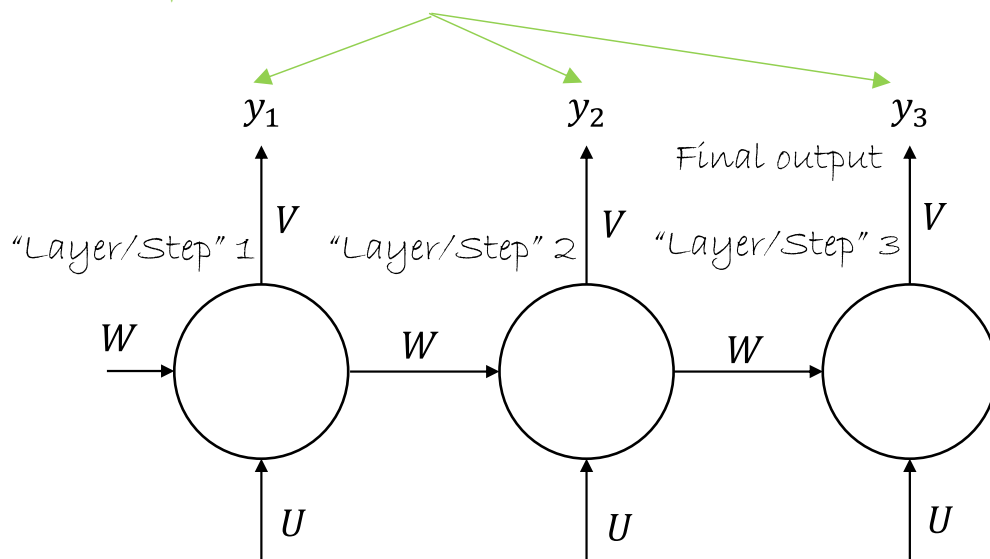
مقایسه‌ی شبکه‌ی بسته شده با شبکه‌ی چندلایه

ROLLED NETWORK VS. MULTI-LAYER NETWORK

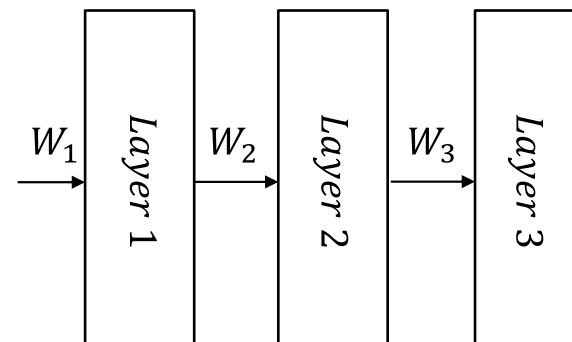
تفاوت‌ها:

- در شبکه‌ی بسته شده، گام‌ها به جای لایه‌ها نقش بازی می‌کنند.
- در شبکه‌ی بسته شده، پارامترهای گام مشترک است؛ در حالی که در شبکه‌ی چندلایه، پارامترها متفاوت است.

گاهی خروجی‌های میانی مورد نیاز نیستند،
با حذف آنها تقریباً به یک مدل استاندارد شبکه‌ی عصبی می‌رسیم:



3-gram unrolled recurrent network



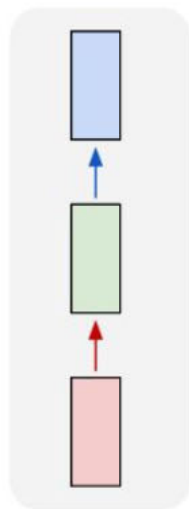
3-layer neural network

پردازش دنباله‌ها با استفاده از شبکه‌های عصبی بازگشتی

یک بردار به یک بردار

RECURRENT NEURAL NETWORKS: PROCESS SEQUENCES

one to one



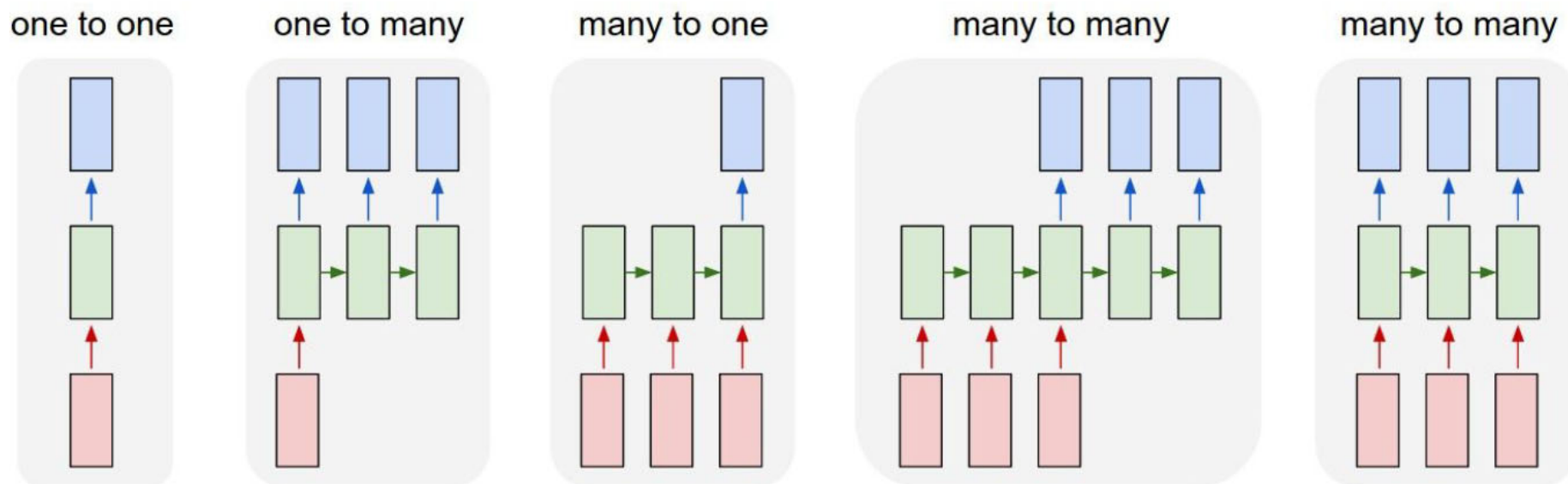
Vanilla Neural Networks

شبکه‌های عصبی ساده:

یک بردار به یک بردار

پردازش دنباله‌ها با استفاده از شبکه‌های عصبی بازگشتی

یک بردار به چند بردار

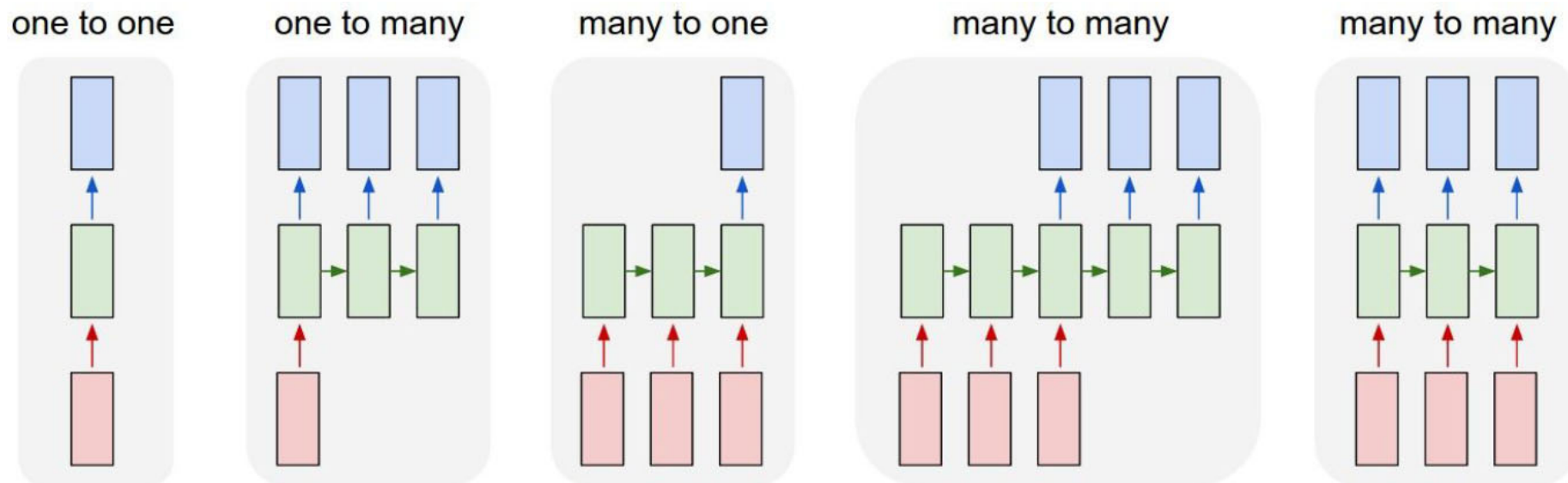


e.g. **Image Captioning**
image -> sequence of words

مانند: عنوان‌گذاری تصاویر
تصویر ← دنباله‌ی کلمات

پردازش دنباله‌ها با استفاده از شبکه‌های عصبی بازگشتی

چند بردار به یک بردار



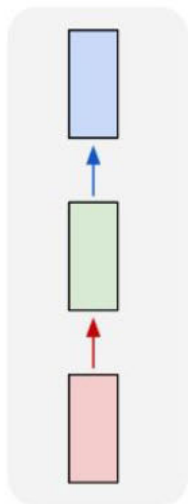
e.g. **Sentiment Classification**
sequence of words -> sentiment

مانند: طبقه‌بندی احساسات (نظرات)
دنباله‌ی کلمات ← احساس (نظر)

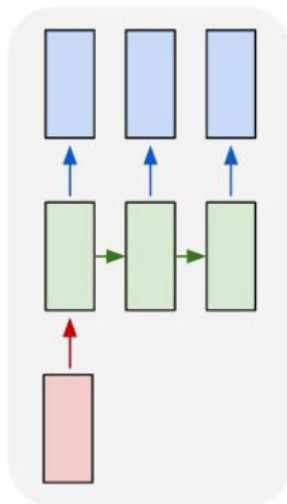
پردازش دنباله‌ها با استفاده از شبکه‌های عصبی بازگشتی

چند بردار به چند بردار

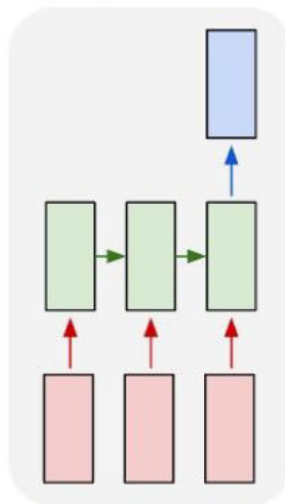
one to one



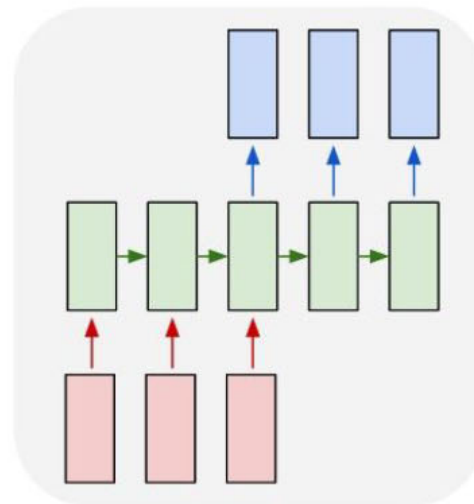
one to many



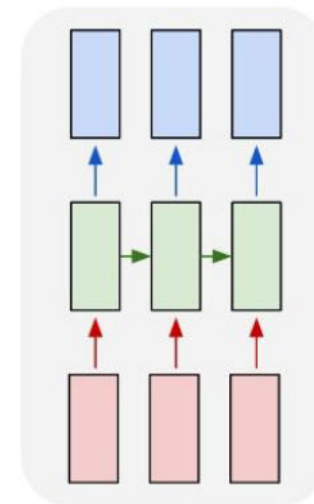
many to one



many to many



many to many

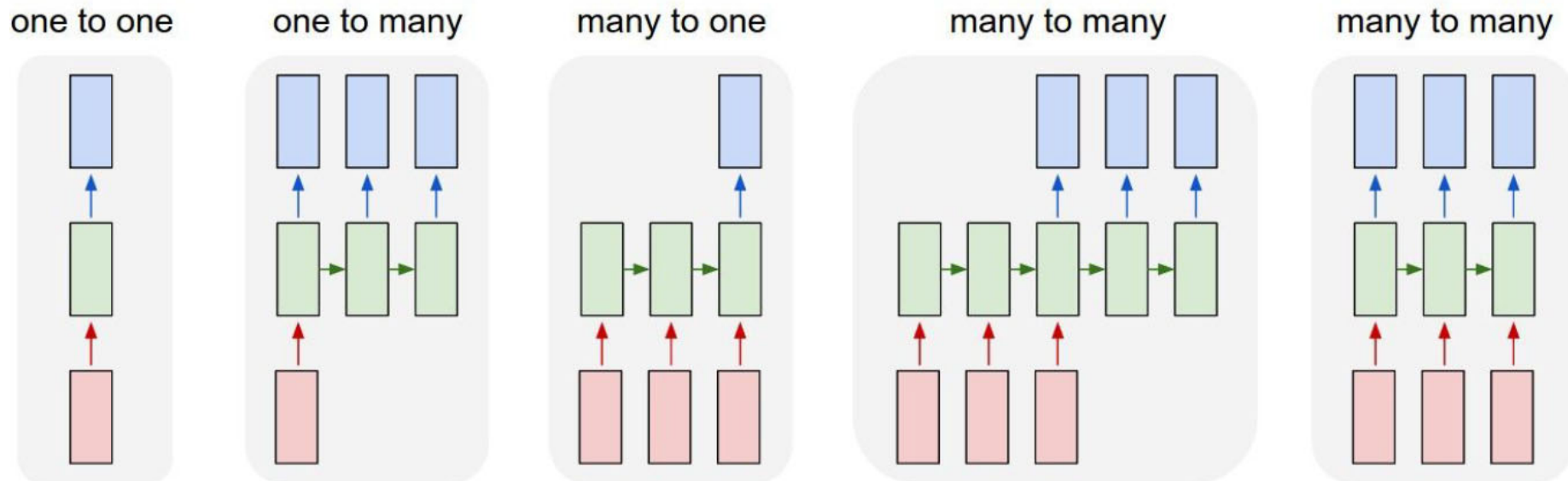


e.g. **Machine Translation**
seq of words -> seq of words

مانند: ترجمه‌ی ماشینی
دنباله‌ی کلمات ← دنباله‌ی کلمات

پردازش دنباله‌ها با استفاده از شبکه‌های عصبی بازگشتی

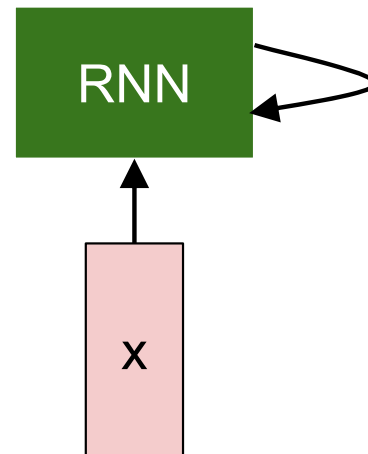
چند بردار به چند بردار



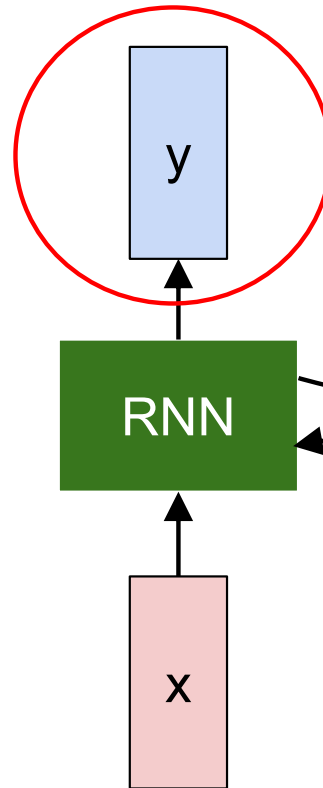
e.g. Video classification on frame level

مانند: طبقه‌بندی ویدئوها در سطح فریم

شبکه‌های عصبی بازگشتی

RECURRENT NEURAL NETWORK

شبکه‌های عصبی بازگشتی

RECURRENT NEURAL NETWORK

معمولاً می‌خواهیم
یک بردار را در چند گام زمانی پیش‌بینی کنیم
usually want to predict
a vector at some time steps

شبکه‌های عصبی بازگشتی

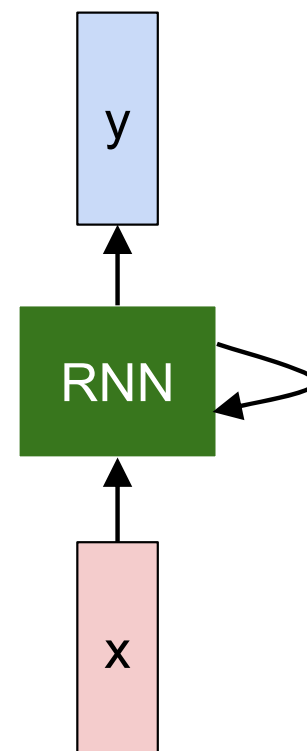
RECURRENT NEURAL NETWORK

می‌توانیم یک دنباله از بردارهای \mathbf{x} را با اعمال یک فرمول بازگشتی در هر گام زمانی پردازش کنیم.

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

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

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



شبکه‌های عصبی بازگشتی

RECURRENT NEURAL NETWORK

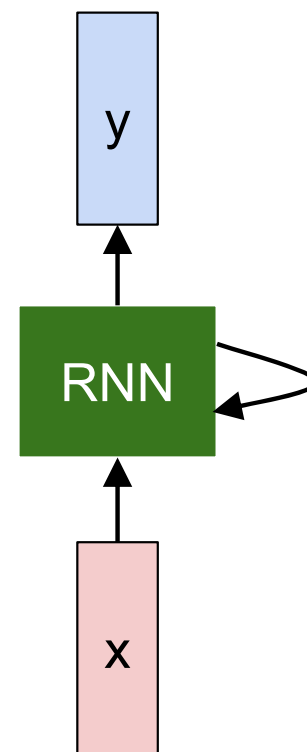
می‌توانیم یک دنباله از بردارهای \mathbf{x} را با اعمال یک فرمول بازگشتی در هر گام زمانی پردازش کنیم.

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

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

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

توجه: تابع یکسان و مجموعه پارامترهای یکسانی در هر گام زمانی استفاده می‌شوند.



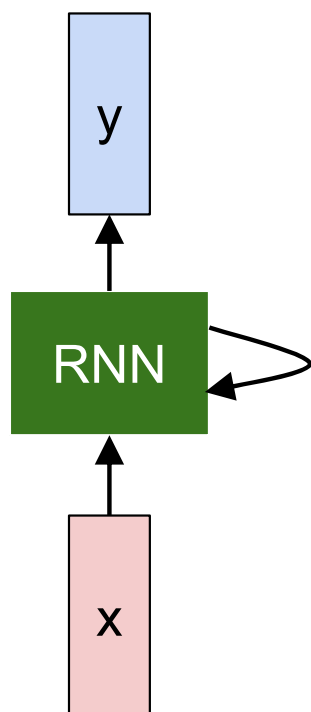
شبکه‌های عصبی بازگشتی

مدل ساده

(VANILLA) RECURRENT NEURAL NETWORK

حالت از یک بردار تنهای پنهان h تشکیل شده است.

The state consists of a single **hidden** vector h :



$$y_t = W_{hy}h_t + b_y$$

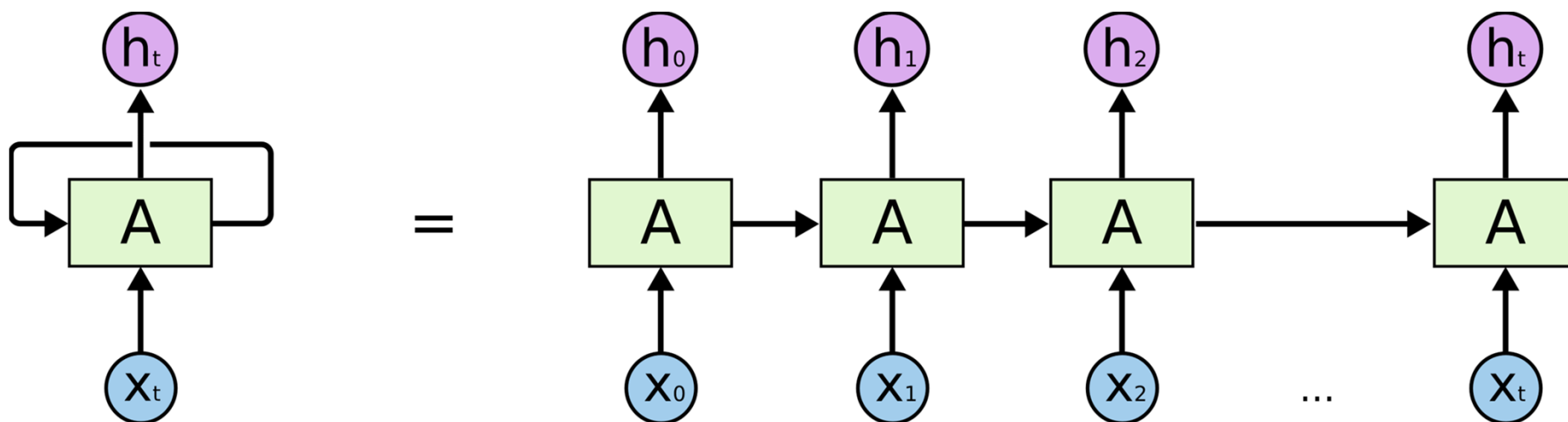
$$h_t = f_W(h_{t-1}, x_t)$$

مانند ↓

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

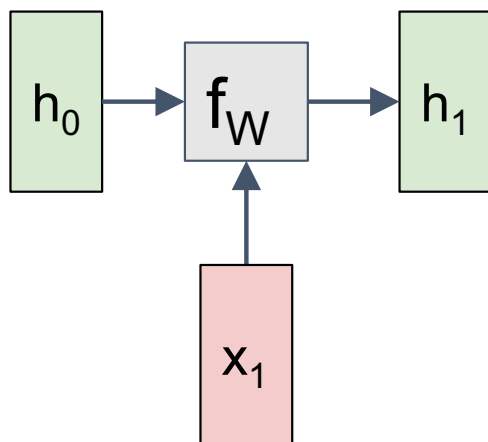
شبکه‌های عصبی بازگشتی

باز کردن مدل ساده



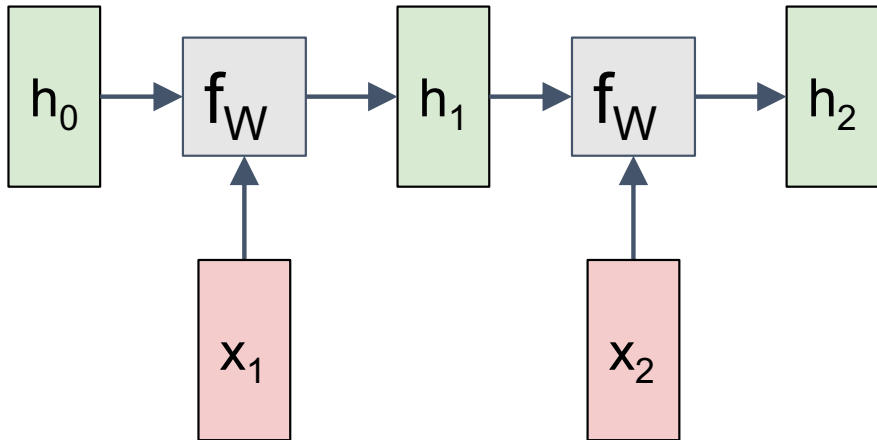
شبکه‌های عصبی بازگشتی

گراف محاسباتی

RNN: COMPUTATIONAL GRAPH

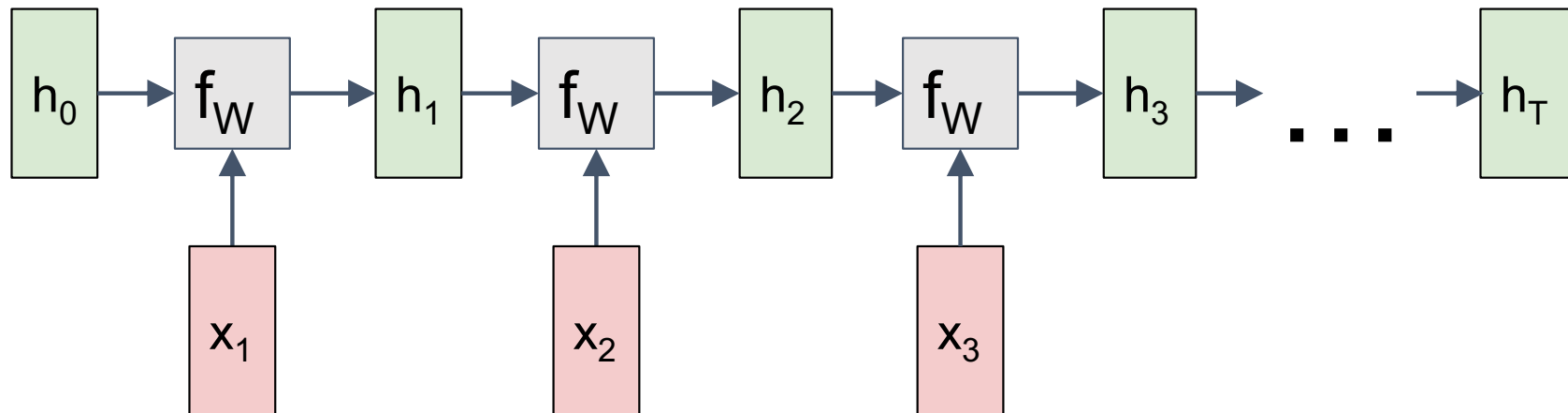
شبکه‌های عصبی بازگشتی

گراف محاسباتی

RNN: COMPUTATIONAL GRAPH

شبکه‌های عصبی بازگشتی

گراف محاسباتی

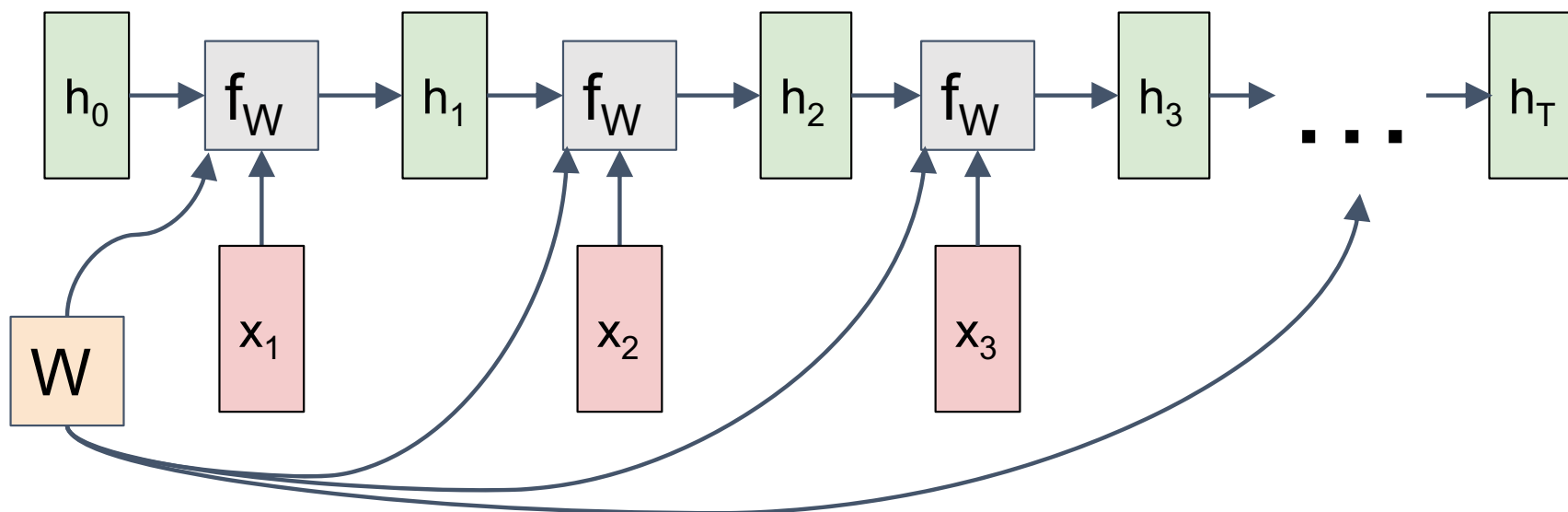
RNN: COMPUTATIONAL GRAPH

شبکه‌های عصبی بازگشتی

گراف محاسباتی: استفاده‌ی مجدد از یک ماتریس وزن یکسان در همه‌ی گام‌های زمانی

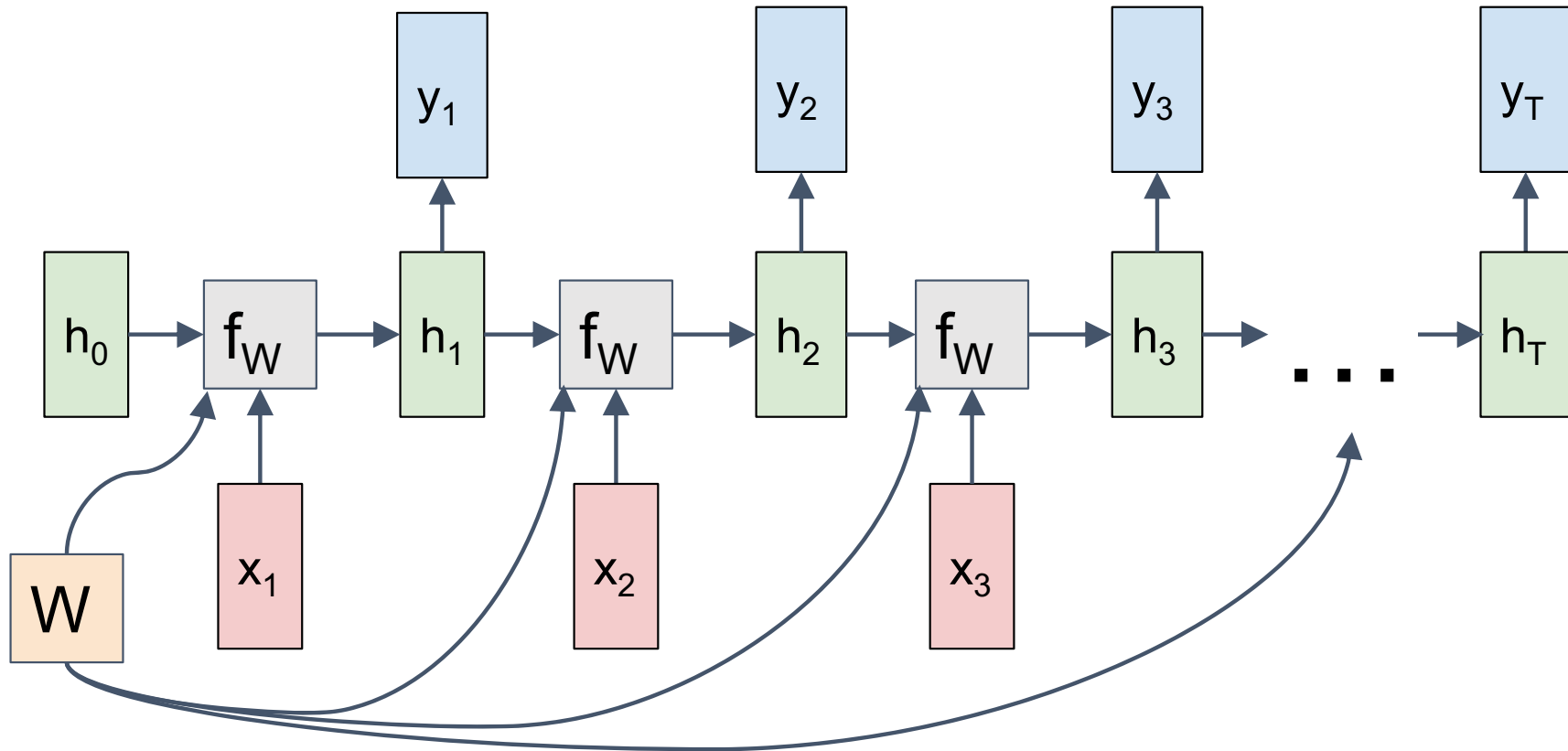
RNN: COMPUTATIONAL GRAPH

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



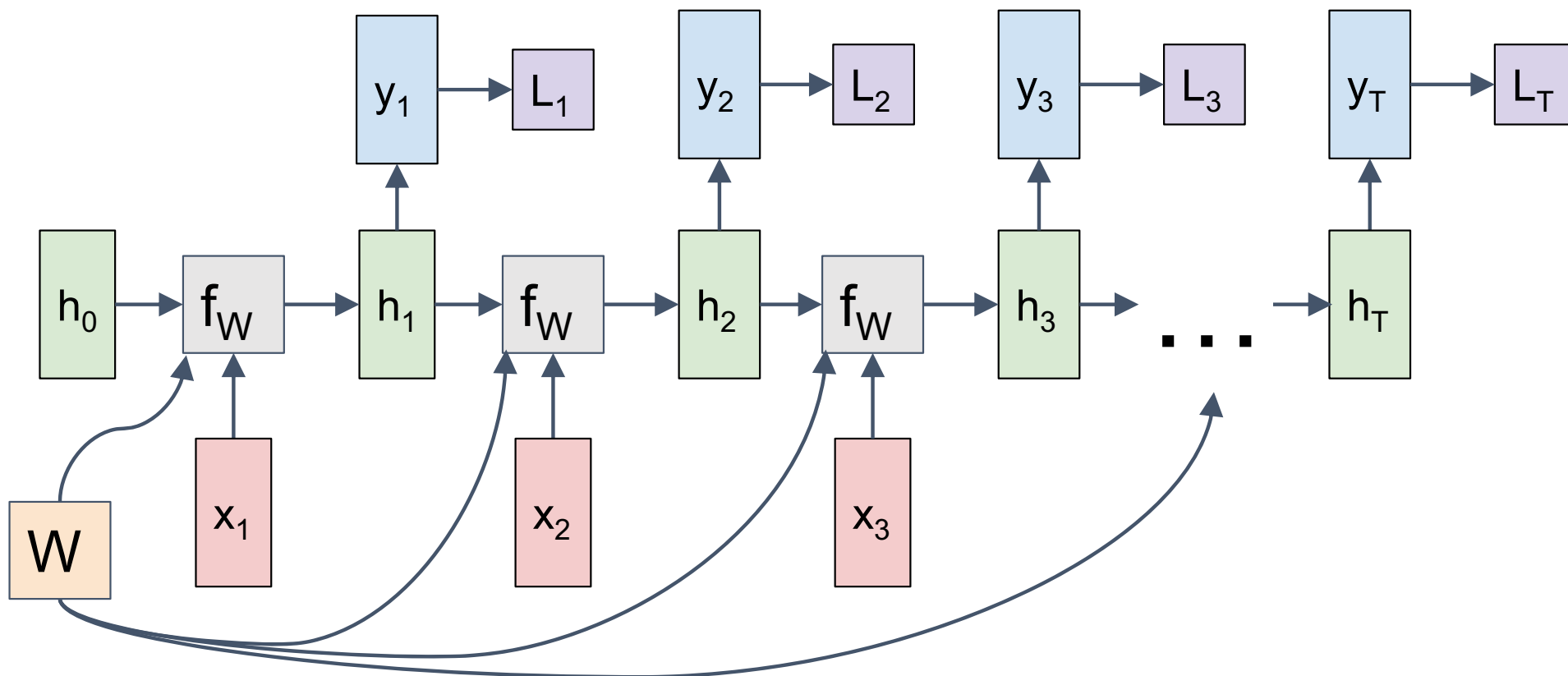
شبکه‌های عصبی بازگشتی

گراف محاسباتی: چند به چند

RNN: COMPUTATIONAL GRAPH: MANY TO MANY

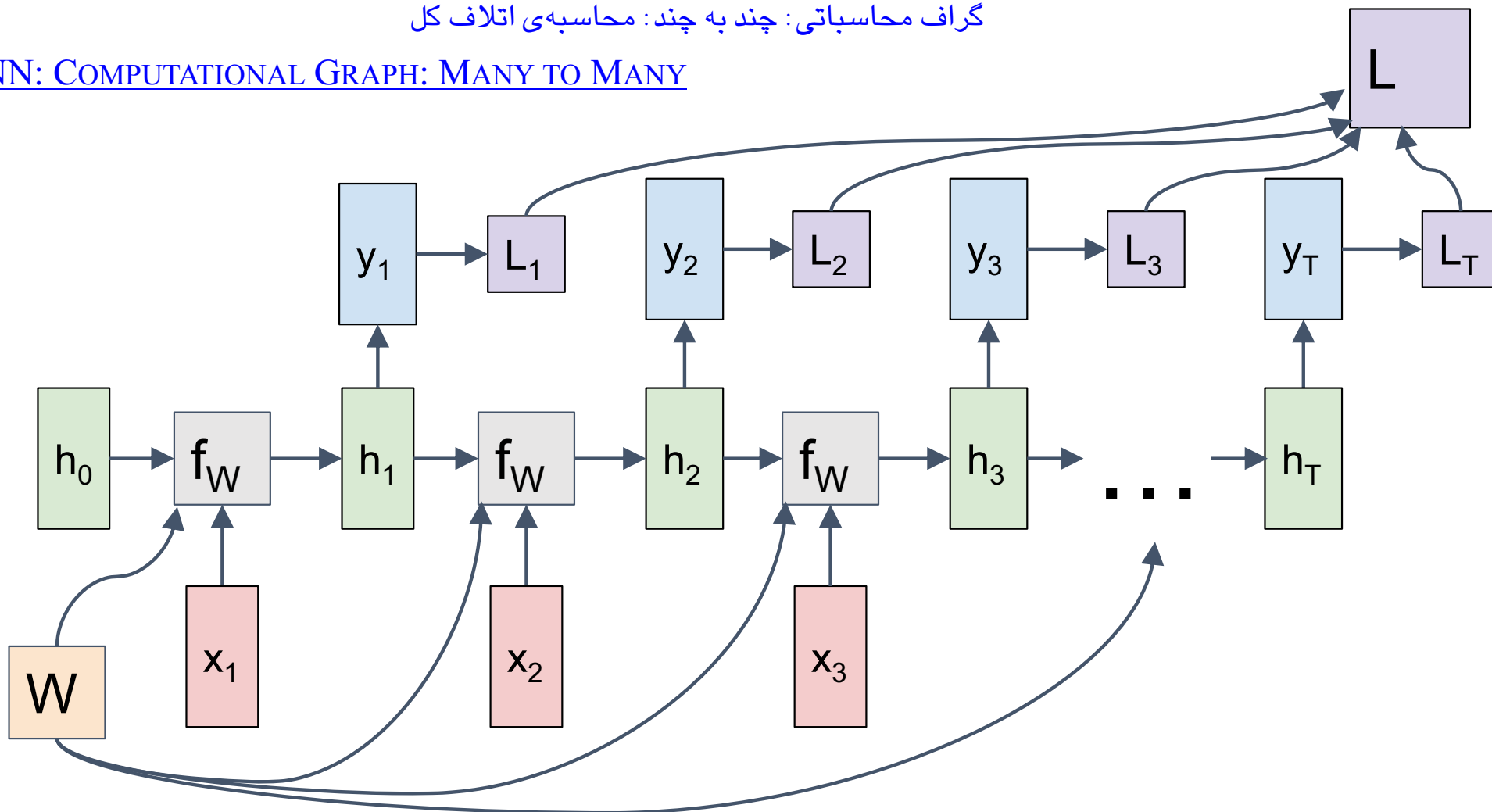
شبکه‌های عصبی بازگشتی

گراف محاسباتی: چند به چند: محاسبه‌ی اتلاف‌ها

RNN: COMPUTATIONAL GRAPH: MANY TO MANY

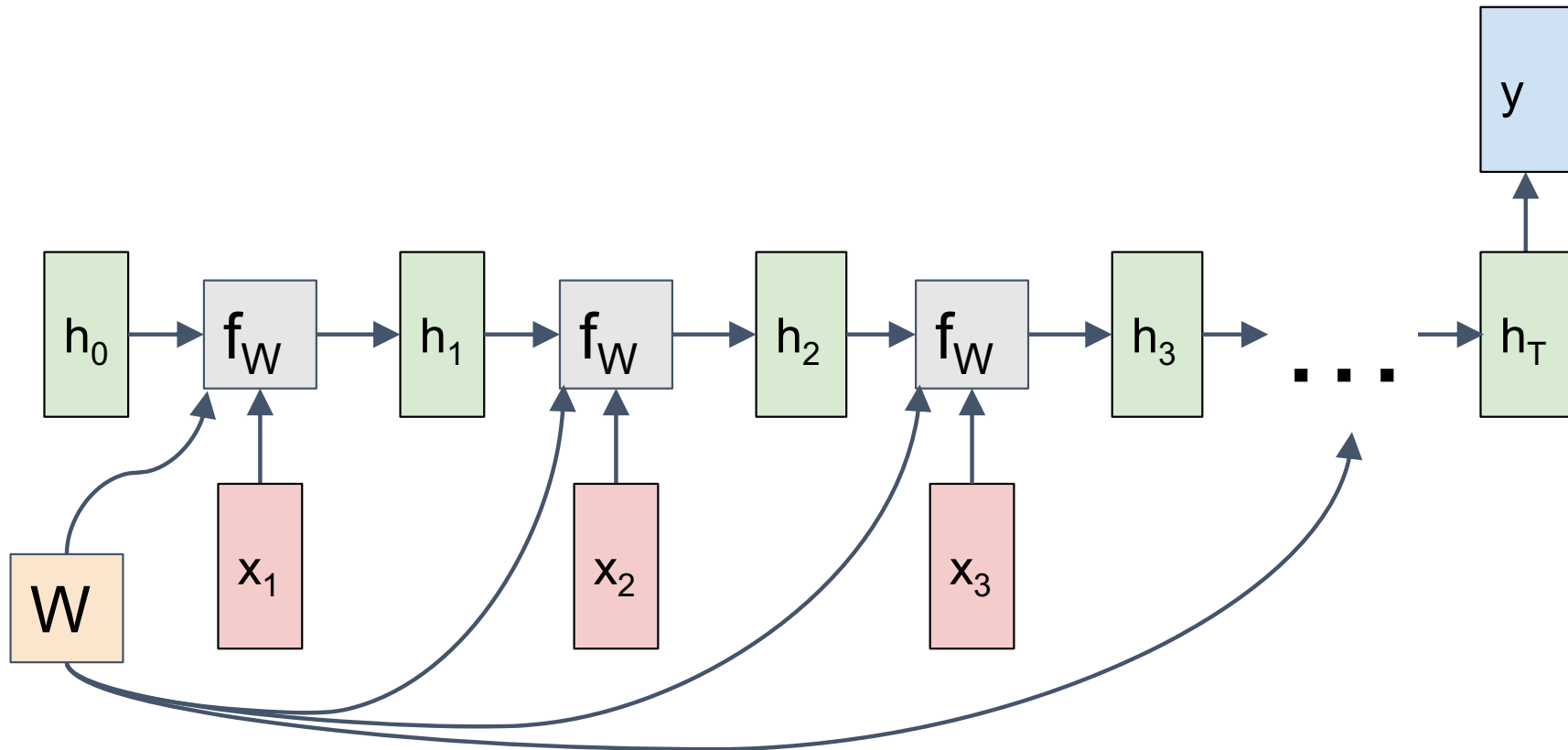
شبکه‌های عصبی بازگشتی

گراف محاسباتی: چند به چند: محاسبه‌ی اتلاف کل

RNN: COMPUTATIONAL GRAPH: MANY TO MANY

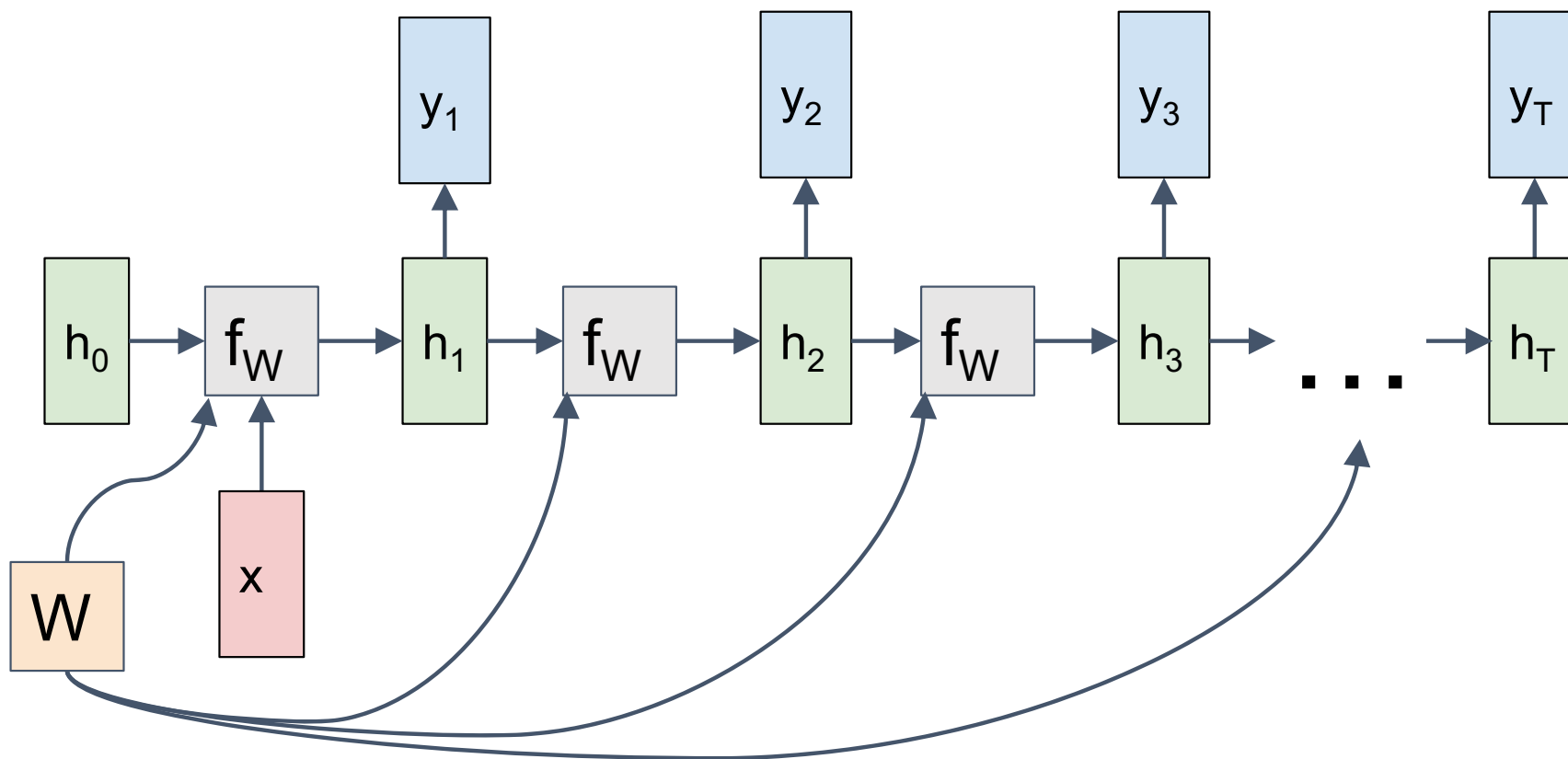
شبکه‌های عصبی بازگشتی

گراف محاسباتی: چند به یک

RNN: COMPUTATIONAL GRAPH: MANY TO ONE

شبکه‌های عصبی بازگشتی

گراف محاسباتی: یک به چند

RNN: COMPUTATIONAL GRAPH: ONE TO MANY

شبکه‌های عصبی بازگشتی

دنباله به دنباله: چند به یک + یک به چند

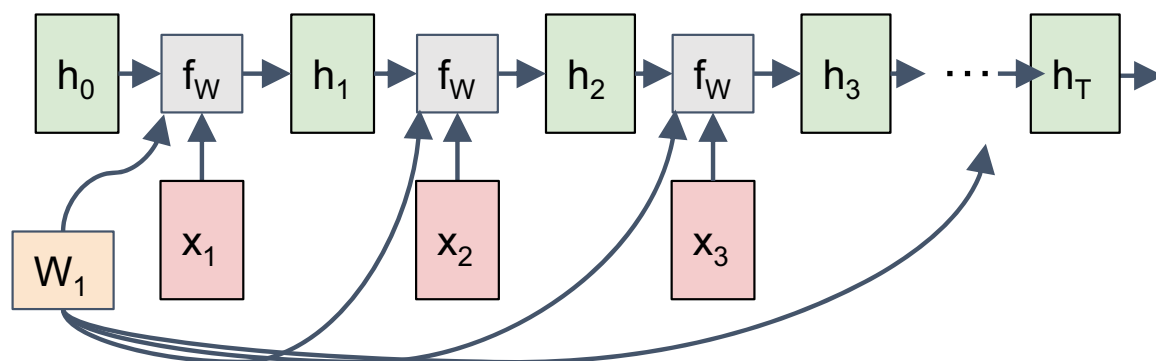
SEQUENCE TO SEQUENCE: MANY-TO-ONE + ONE-TO-MANY

چند به یک:

کدگذاری دنباله‌ی ورودی به یک بردار تنها

Many to one:

Encode input sequence in a single vector



شبکه‌های عصبی بازگشتی

دنباله به دنباله: چند به یک + یک به چند

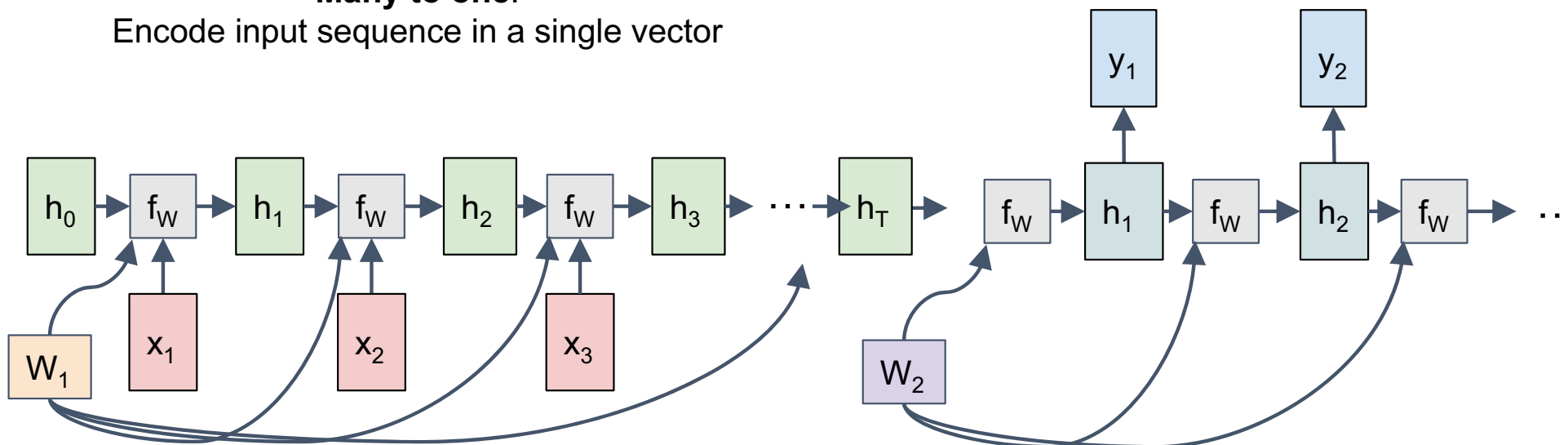
SEQUENCE TO SEQUENCE: MANY-TO-ONE + ONE-TO-MANY

چند به یک:

کدگذاری دنباله‌ی ورودی به یک بردار تنها

Many to one:

Encode input sequence in a single vector



یک به چند:

تولید دنباله‌ی خروجی از روی یک بردار ورودی تنها

One to many:

Produce output sequence from single input vector

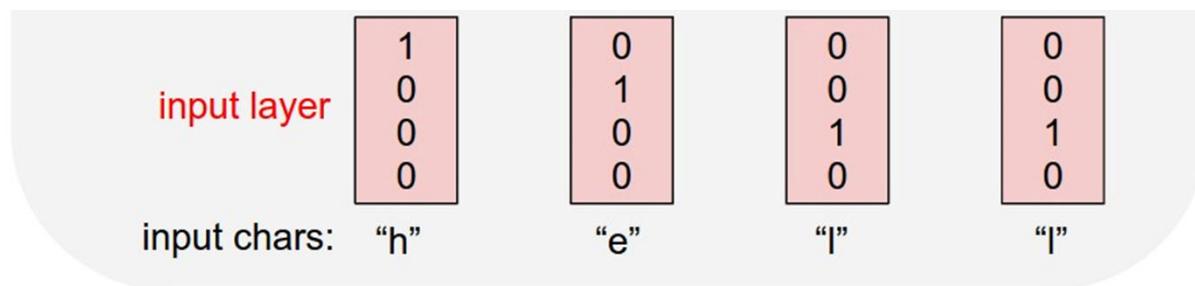
شبکه‌های عصبی بازگشتی

مثال: مدل زبان در سطح کاراکتر

Example: Character-level Language Model

Vocabulary:
[h,e,l,o]

Example training
sequence: **“hello”**



هدف: پیش‌بینی کاراکتر بعدی با داشتن دنباله‌ای از کاراکترها

شبکه‌های عصبی بازگشتی

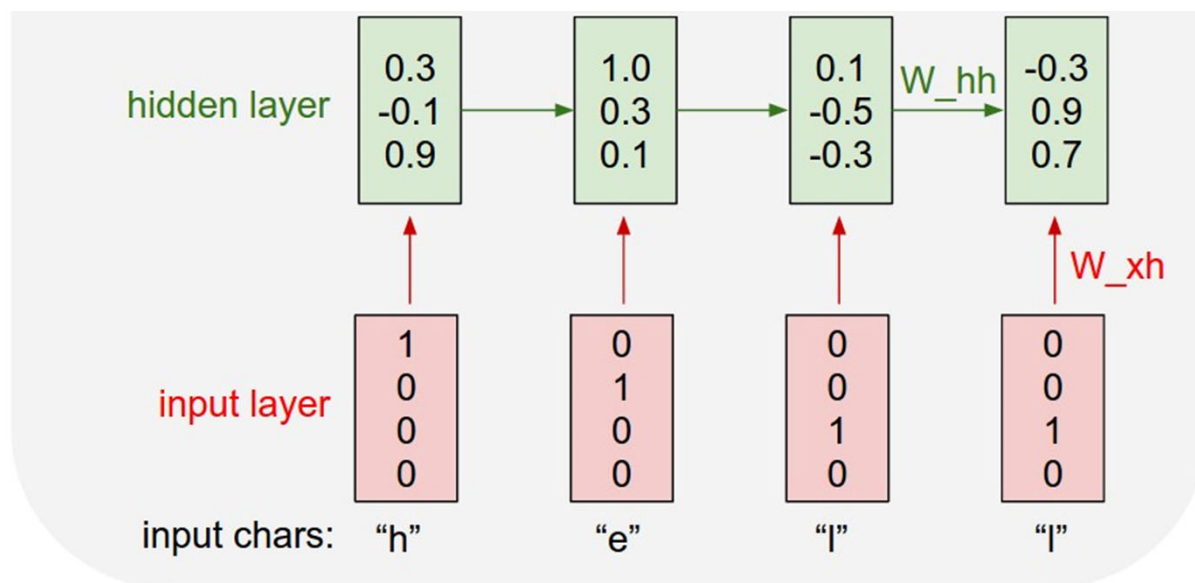
مثال: مدل زبان در سطح کاراکتر

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

Example: Character-level Language Model

Vocabulary:
[h,e,l,o]

Example training
sequence: **“hello”**



هدف: پیش‌بینی کاراکتر بعدی با داشتن دنباله‌ای از کاراکترها

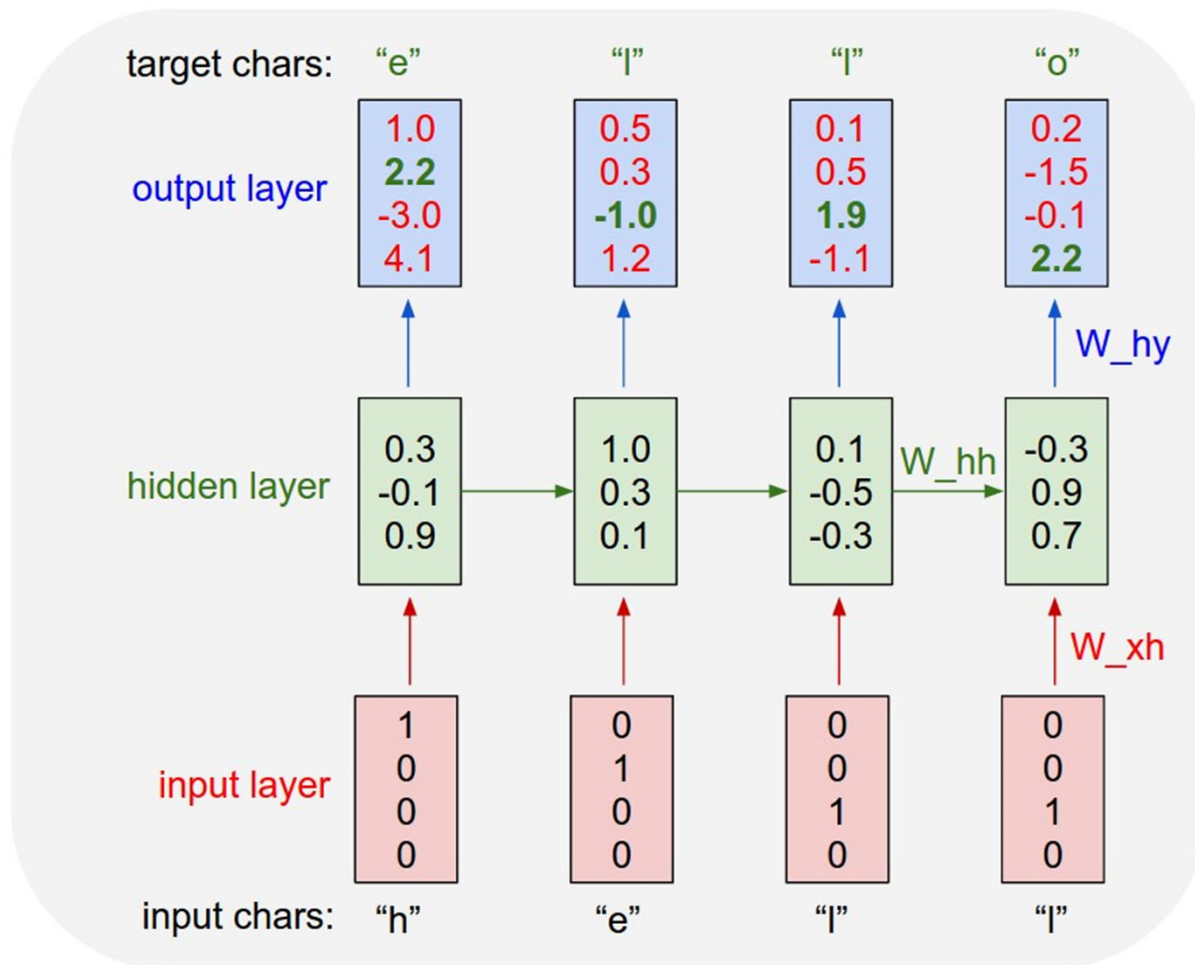
شبکه‌های عصبی بازگشتی

مثال: مدل زبان در سطح کاراکتر

Example: Character-level Language Model

Vocabulary:
[h,e,l,o]

Example training
sequence: “**hello**”



هدف: پیش‌بینی کاراکتر بعدی با داشتن دنباله‌ای از کاراکترها

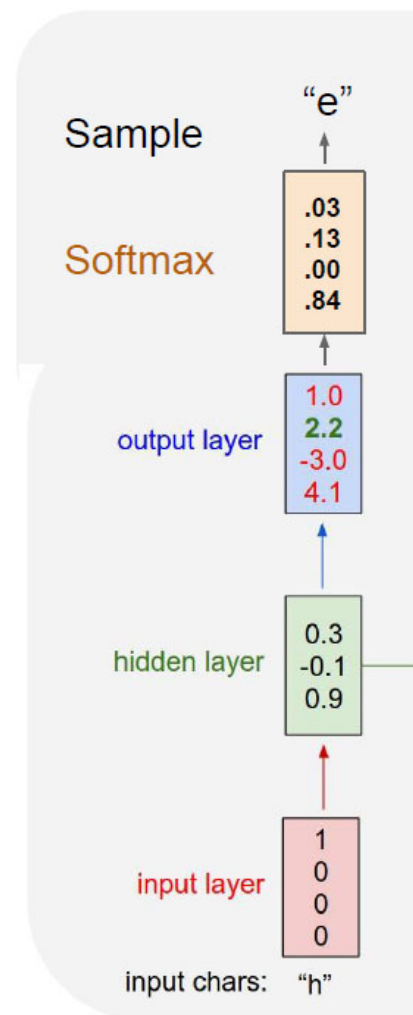
شبکه‌های عصبی بازگشتی

مثال: مدل زبان در سطح کاراکتر

Example: Character-level Language Model

Vocabulary:
[h,e,l,o]

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



در زمان آزمایش، در هر لحظه یکی از کاراکترهای نمونه به مدل فیدبک می‌شود.

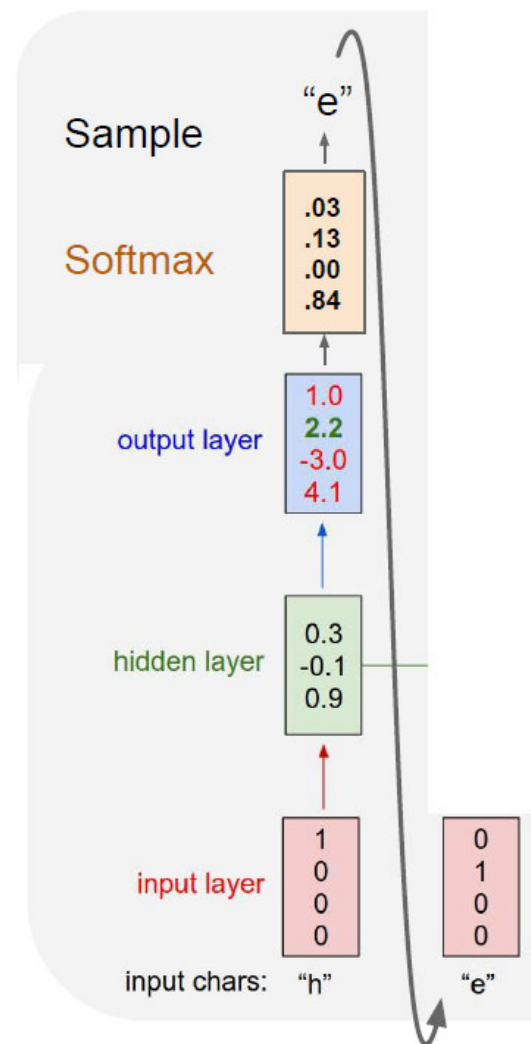
شبکه‌های عصبی بازگشتی

مثال: مدل زبان در سطح کاراکتر

Example: Character-level Language Model

Vocabulary:
[h,e,l,o]

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



در زمان آزمایش، در هر لحظه یکی از کاراکترهای نمونه به مدل فیدبک می‌شود.

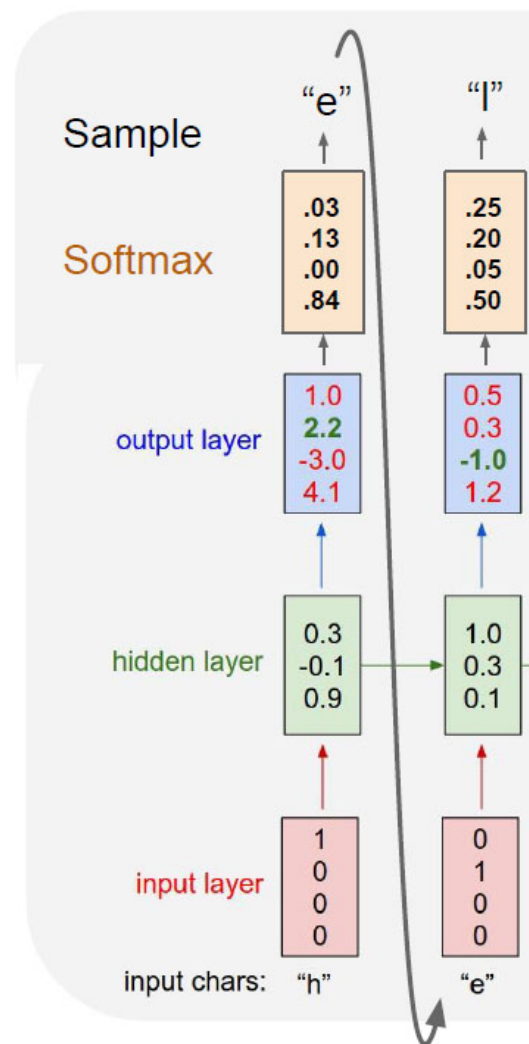
شبکه‌های عصبی بازگشتی

مثال: مدل زبان در سطح کاراکتر

Example: Character-level Language Model

Vocabulary:
[h,e,l,o]

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



در زمان آزمایش، در هر لحظه یکی از کاراکترهای نمونه به مدل فیدبک می‌شود.

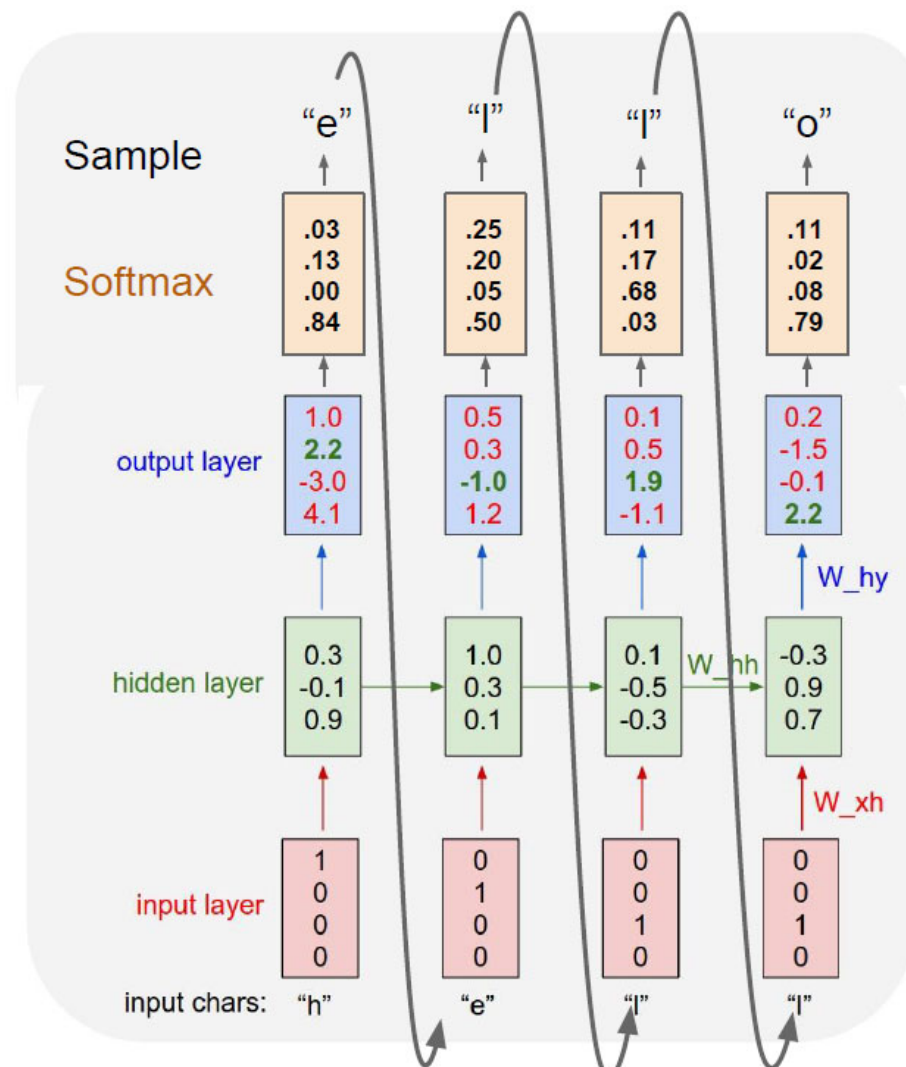
شبکه‌های عصبی بازگشتی

مثال: مدل زبان در سطح کاراکتر

Example: Character-level Language Model

Vocabulary:
[h,e,l,o]

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



در زمان آزمایش، در هر لحظه یکی از کاراکترهای نمونه به مدل فیدبک می‌شود.

شبکه‌های عصبی بازگشتی

۲

نمونه
کاربردهای
شبکه‌های
عصبی
بازگشتی

یادگیری مدل زبان

مثال

```

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(hidden_size, vocab_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 Nx1 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     dwhx, 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     dhnxt = np.zeros_like(hs[0])
48     for t in reversed(xrange(len(inputs))):
49         dy = np.copy(ps[t])
50         dy[targets[t]] -= 1 # backprop into y
51         dwhy += np.dot(dy, hs[t].T)
52         dby += dy
53         dh = np.dot(why.T, dy) + dhnxt # backprop into h
54         dhrw = (1 - hs[t]**2) * dh # backprop through tanh nonlinearity
55         dbh += dhrw
56         dwhx += np.dot(dhrw, xs[t].T)
57         dwhh += np.dot(dhrw, hs[t-1].T)
58         dhnxt = np.dot(whh.T, dhrw)
59     for dparam in [dwhx, dwhh, dwhy, dbh, dby]:
60         np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61     return loss, dwhx, 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 mxh, 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 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, dwhx, 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([wxh, whh, why, bh, by],
106                                  [dwhx, dwhh, dwhy, dbh, dby],
107                                  [mxh, mwhh, mwhy, mbh, mby]):
108        mem += dparam * dparam
109        param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter

```

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

[min-char-rnn.py gist: 112 lines of Python](https://gist.github.com/karpathy/d4dee566867f8291f086)

یادگیری مدل زبان

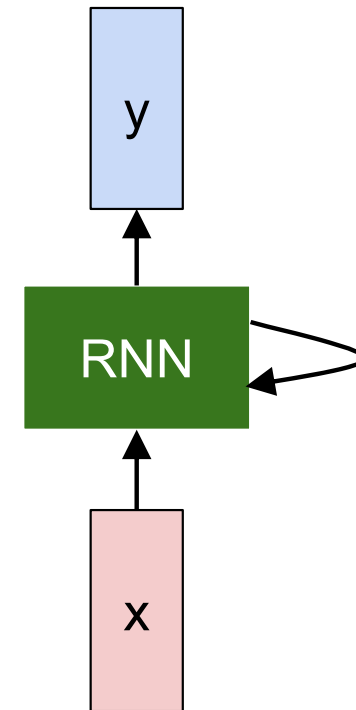
مثال

THE SONNETS

by William Shakespeare

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

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



یادگیری مدل زبان

مثال

at first:

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



train more

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



train more

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



train more

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

یادگیری مدل زبان

مثال

PANDARUS:

Alas, I think he shall be come approached and the day
When little srain 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.

یادگیری مدل زبان





















مثال: پروژه‌ی استکس

THE STACKS PROJECT: OPEN SOURCE ALGEBRAIC GEOMETRY TEXTBOOK

 **The Stacks Project**

[home](#)
[about](#)
[tags explained](#)
[tag lookup](#)
[browse](#)
[search](#)
[bibliography](#)
[recent comments](#)
[blog](#)
[add slogans](#)

Browse chapters

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

Parts

1. [Preliminaries](#)
2. [Schemes](#)
3. [Topics in Scheme Theory](#)
4. [Algebraic Spaces](#)
5. [Topics in Geometry](#)
6. [Deformation Theory](#)
7. [Algebraic Stacks](#)
8. [Miscellany](#)

Statistics

The Stacks project now consists of

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

LaTeX source

<http://stacks.math.columbia.edu/>The stacks project is licensed under the [GNU Free Documentation License](#)

یادگیری مدل زبان

مثال

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

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

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

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

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

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

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

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

is a unique morphism of algebraic stacks. Note that

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

and

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

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

Proof. See discussion of sheaves of sets. \square

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

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

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

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

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

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

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

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

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

Lemma 0.3. In Situation ?? . Hence we may assume $\mathfrak{q}' = 0$.

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

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

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

یادگیری مدل زبان

مثال

Proof. Omitted. □

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

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

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

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

$$\mathcal{O}_X(\mathcal{F}) = \{morph_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

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

Lemma 0.2. *This is an integer Z is injective.*

Proof. See Spaces, Lemma ?? □

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

The following to the construction of the lemma follows.

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

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

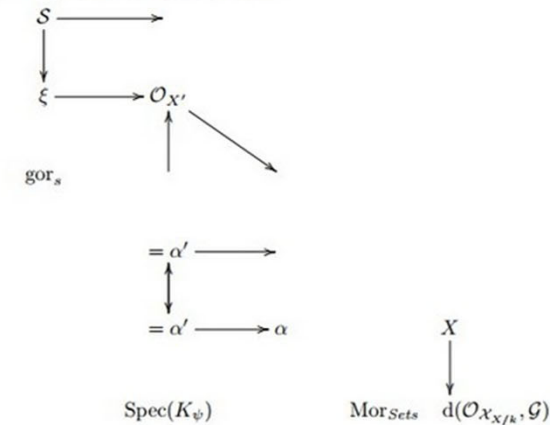
be a morphism of algebraic spaces over S and Y .

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

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

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

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



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

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

□

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

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

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

$$\mathcal{O}_{X,x} \rightarrow \mathcal{F}_{\bar{x}}^{-1}(\mathcal{O}_{X_{\acute{e}tale}}) \rightarrow \mathcal{O}_{X'_t}^{-1} \mathcal{O}_{X_\lambda}(\mathcal{O}_{X_\eta}^{\bar{v}})$$

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

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

If \mathcal{F} is a scheme theoretic image points. □

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

یادگیری مدل زبان

مثال

The screenshot shows the GitHub repository page for 'torvalds / linux'. At the top, there is a search bar and navigation links for 'Explore', 'Gist', 'Blog', and 'Help'. The repository name 'torvalds / linux' is prominently displayed, along with statistics: 3,711 Watchers, 23,054 Stars, and 9,141 Forks. Below this, the 'Linux kernel source tree' is shown with a summary of 520,037 commits, 1 branch, 420 releases, and 5,039 contributors. A recent commit by 'torvalds' is highlighted, showing a merge of the 'drm-fixes' branch. A table of recent commits follows, listing various subdirectories like 'Documentation', 'arch', 'block', 'crypto', 'drivers', 'firmware', 'fs', 'include', 'init', and 'io'. On the right side, there are links for 'Code', 'Pull requests' (74), 'Pulse', and 'Graphs'. At the bottom right, the 'HTTPS clone URL' is provided as 'https://github.c...', along with buttons for 'Clone in Desktop' and 'Download ZIP'.

یادگیری مدل زبان

مثال

Generated
C code

```

static void do_command(struct seq_file *m, void *v)
{
    int column = 32 << (cmd[2] & 0x80);
    if (state)
        cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
    else
        seq = 1;
    for (i = 0; i < 16; i++) {
        if (k & (1 << 1))
            pipe = (in_use & UMXTHREAD_UNCCA) +
                ((count & 0x00000000ffffffff) & 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 ");
}

```

یادگیری مدل زبان

مثال

```
/*
 * 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/setew.h>
#include <asm/pgproto.h>
```

یادگیری مدل زبان

مثال

```

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

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

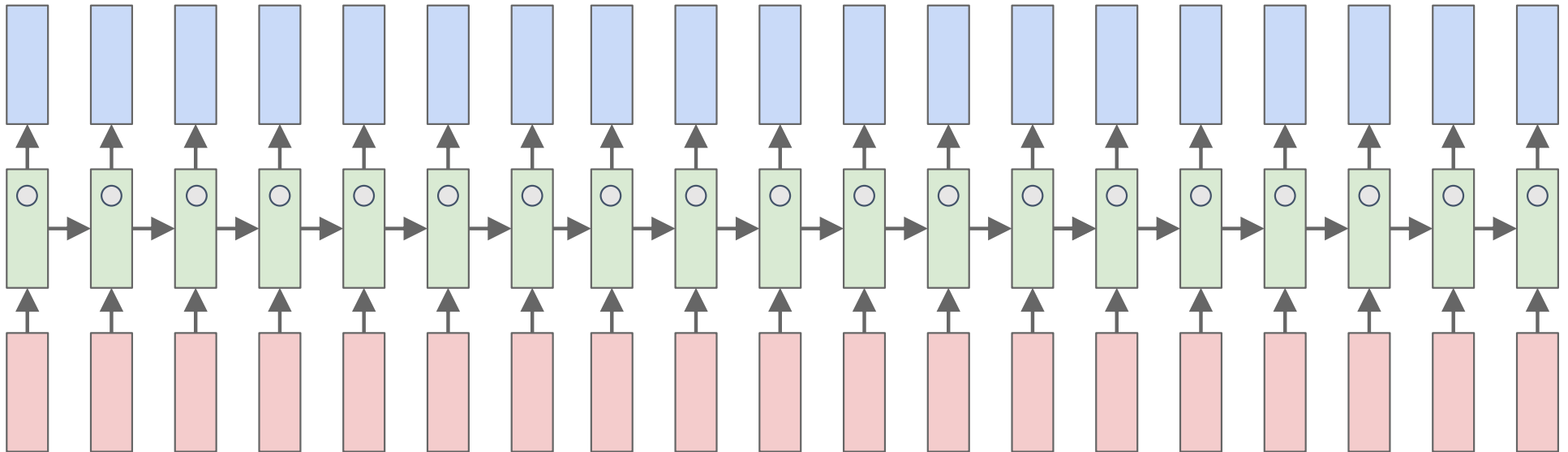
#define SWAP_ALLOCATE(nr)    (e)
#define emulate_sigs()  arch_get_unaligned_child()
#define access_rw(TST)  asm volatile("movd %esp, %0, %3" : : "r" (0)); \
    if (__type & DO_READ)

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

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

```

جستجو به دنبال سلول‌های تفسیرپذیر

SEARCHING FOR INTERPRETABLE CELLS

ttp://www.ynetnews.com/] English-language website of Israel's lar
tp://www.bacahets.com/ - xglis h language sairs site of t sla el i s sing
d : xne . wa ea . . awa to a . s & n t i a c a - s a r d e e l h o a n t b i s a n f a n r e i f ' a a t d
mw - 2 p i i i s o e s s i s . / e r n . c] (d c e e n e p e s a a i k i i e e l e d h , i r t h r a o n s e , c o s e
dr . < : a h b - n p t w t . x i g h / m a) T v d r y z i c o u e d l s u : t h a - o o t u , s t u i f l v e p e r y
s t p , t c o a 2 d r u l w o c l e n s r] p . l l v a o d , , e y t c - n d m - o i b u v s] b b i m s u l t a t t l y b n

gest newspaper ' ' [[Y e d i o t h A h r o n o t h]] ' ' ' H e b r e w - l a n g u a g e p e r i o d
e t a a w s p a p e r s o ' [[T e l t i (f e a n e m t i) ' ' * ' ' [e r r e w s l e n g u a g e : a r o s o d i
i r s c o e e n a i T T h A o a i n n h S r m u w] e y s [' i n e i a ' s i w d d e ' h s o l r i f r :
u s . . s e t l g o r s . a s a t C a r e e g ' a C l r i s z] i e ' : : , # : T A a a a a t B a s e e i l o ' i a n f v l
- t u a e v r t i d , t B A m S u s y u t]] A s a o i g s]] , . : s M B o l o u s : T o u a - n : d w o a p n u
a , d , i i u i t i c p .] (l S v H v t u s u i e D n o e g a n o . ,] : { C C u i b o h e C y b k s l s : r - e p c n t s

icals : ' ' ' * ' ' [[G l o b e s]] ' ' [h t t p : / / w w w . g l o b e s . c o . i l /] b u s i n e s s d a
c a l : ' ' ' * ' ' [T a a b a] ' ' ([t t p : / / w w w . b u o b a l . c o m u n / s A - y t i n e s s a e t
s t l ' [h A e o v e l t s a h a d : x g e . w a o i r . r t o a . e l . i T & a i e g e o o y
t t ' ' ' & [& & m C o e r o n e ' : : , i ' o d w . , : n i i i s a a u e . e n i / o m l c C . (e f t g i r i i u
a ' n : , C : & : # * : a f D r u s u] l , . o m e l p < , d h a ; d e u o o t / i h n c s i f S , u r h o s t , t u n
n k i <] : & 1 1 s T G u i t r s i , : b a c m r - x t p o b - g r e s i s l e r l n a f a D] l o s p t a d , i f r m

ily * ' ' [[H a a r e t z | H a ' A r e t z]] ' ' [h t t p : / / w w w . h a a r e t z . c o . i l /] R e l a t i v
l y * ' ' [[T e r r d n F e r a n t a h]] ' ' ([t t p : / / w w w . b o n m d s t . c o m u n / s - e s a t e o i
r e ' ' ' h A i l n n t t e H a l s r c n o l ' s a h a d : x n e . w a a m r t d h e o h . o l . c & o p i n i v e
k i . : * s C O S a n l t h i T i m ' l i] e : , i m c d w - 2 p h i i s e r d i t . i n a / c m f i . (a f l c a n a
d s - ! [t B T C o m m g d]] W o n a a e , : . b a e r r . < t a i b - d u l c n n c / a r n e s i] l i c e y s t o
n d s # & : G l D u v c c s a o S u c l t e l] z | , : o ' o m t] , : e o a 2 n i v f s r o o e i u n a l a) u v v r o

جستجو به دنبال سلول‌های تفسیرپذیر

SEARCHING FOR INTERPRETABLE CELLS

```

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

```

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

جستجو به دنبال سلول‌های تفسیرپذیر

SEARCHING FOR INTERPRETABLE CELLS

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

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

quote detection cell

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

جستجو به دنبال سلول‌های تفسیرپذیر

SEARCHING FOR INTERPRETABLE CELLS

Cell sensitive to position in line:

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

line length tracking cell

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

جستجو به دنبال سلولهای تفسیرپذیر

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

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

جستجو به دنبال سلول‌های تفسیرپذیر

SEARCHING FOR INTERPRETABLE CELLS

Cell that turns on inside comments and quotes:

```

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

```

quote/comment cell

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

جستجو به دنبال سلولهای تفسیرپذیر

SEARCHING FOR INTERPRETABLE CELLS

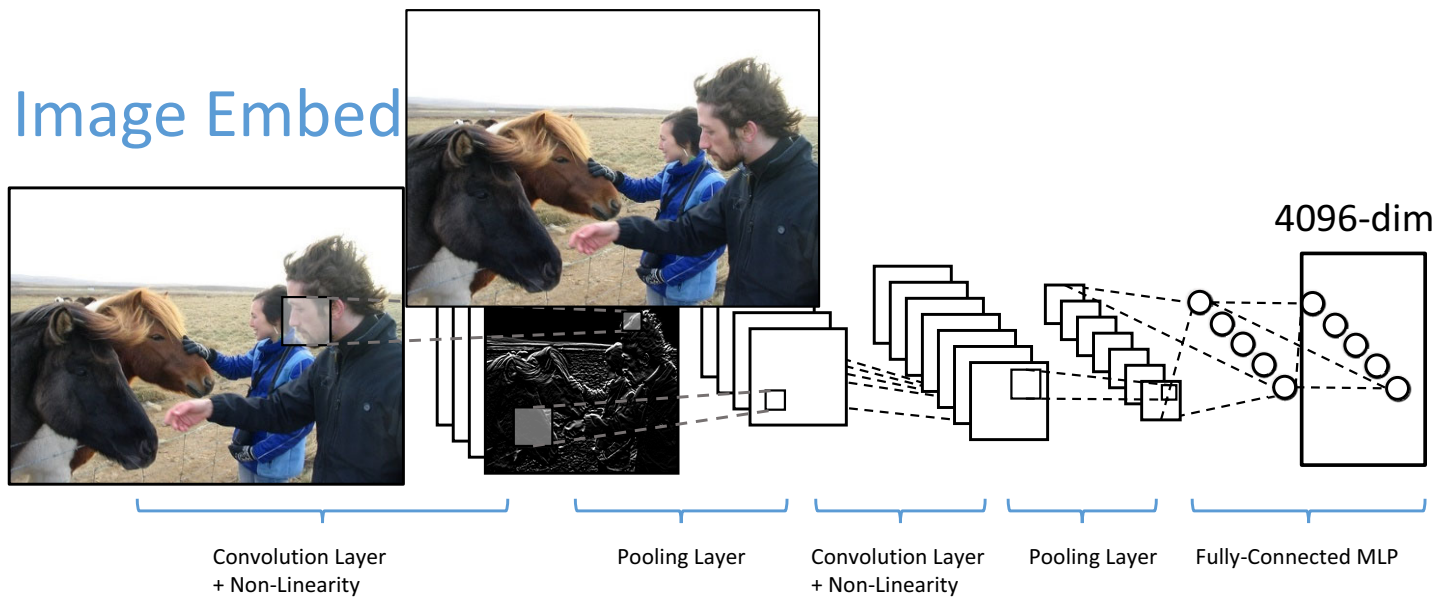
```
#ifdef CONFIG_AUDIT_SYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```

code depth cell

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

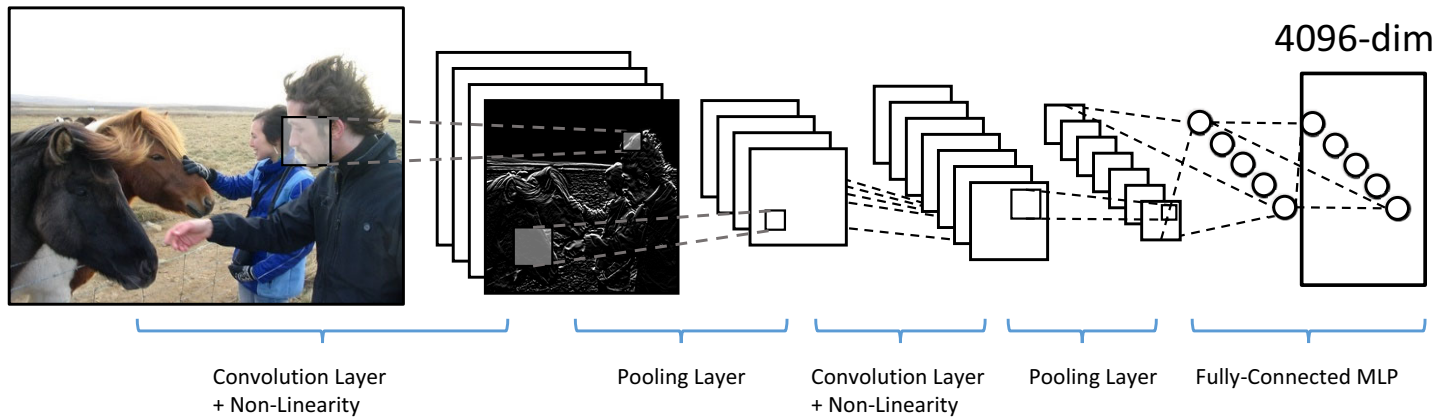
عنوان گذاری عصبی تصاویر

NEURAL IMAGE CAPTIONING

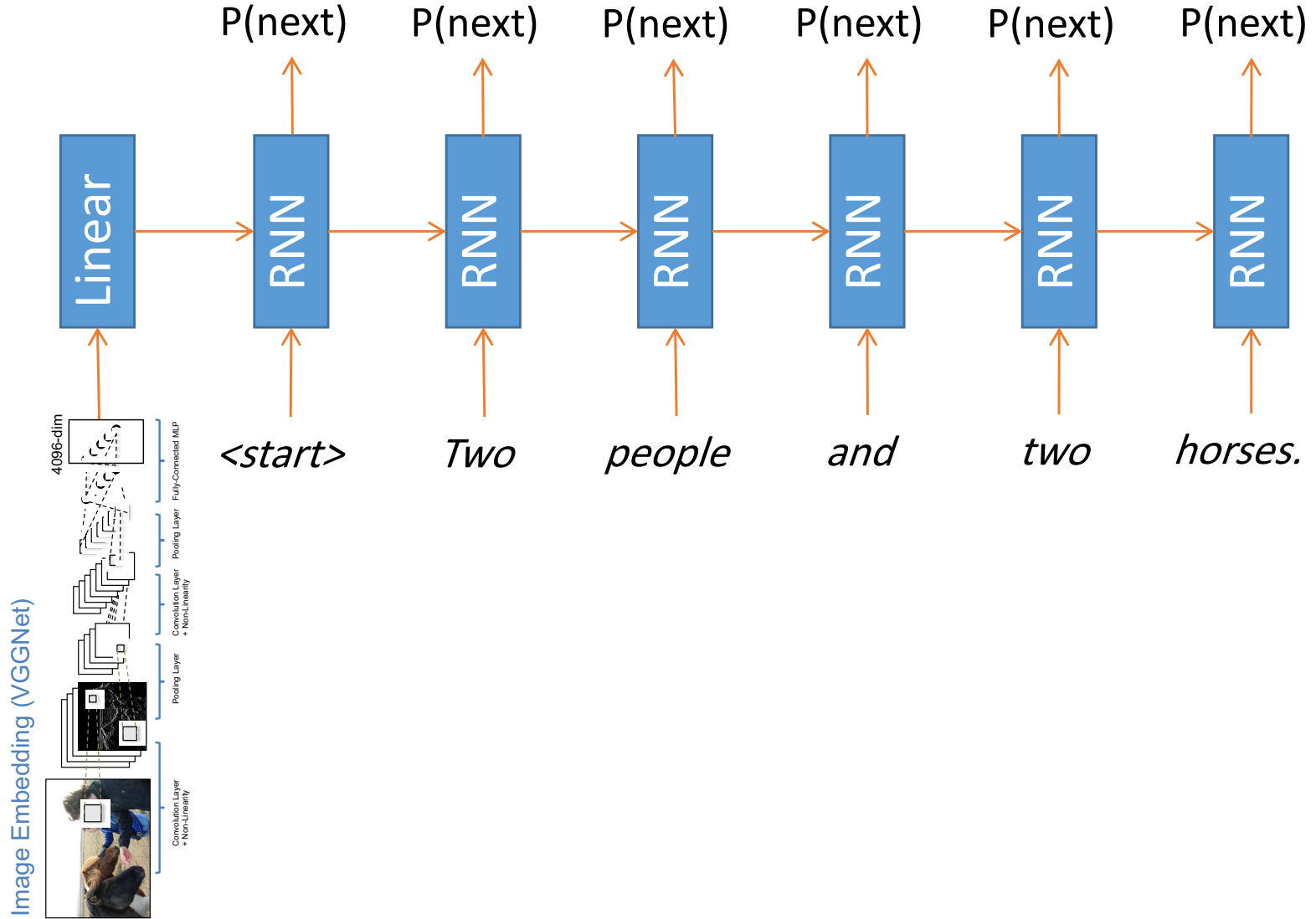


عنوان گذاری عصبی تصاویر

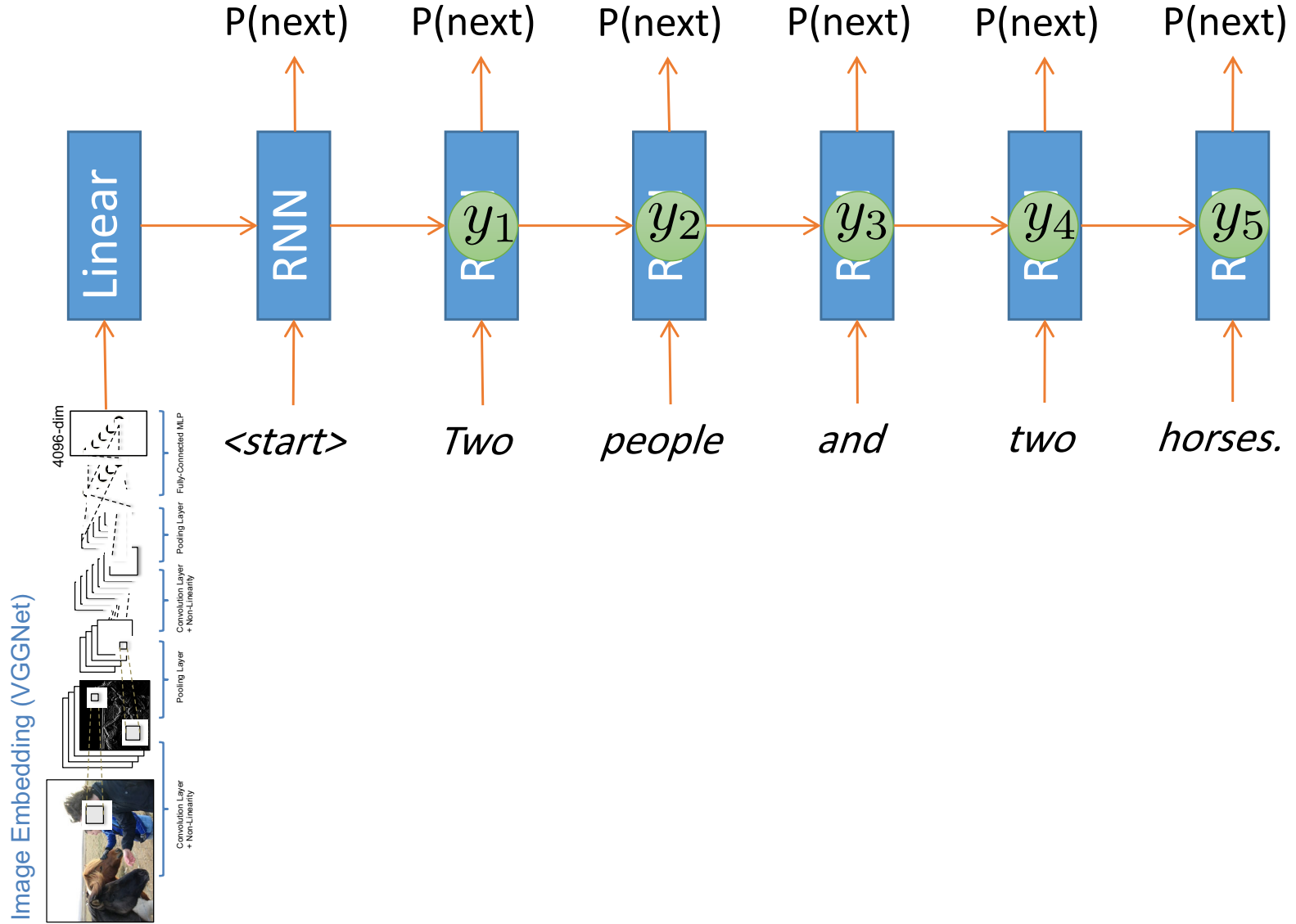
Image Embedding (VGGNet)



عنوان گذاری عصبی تصاویر

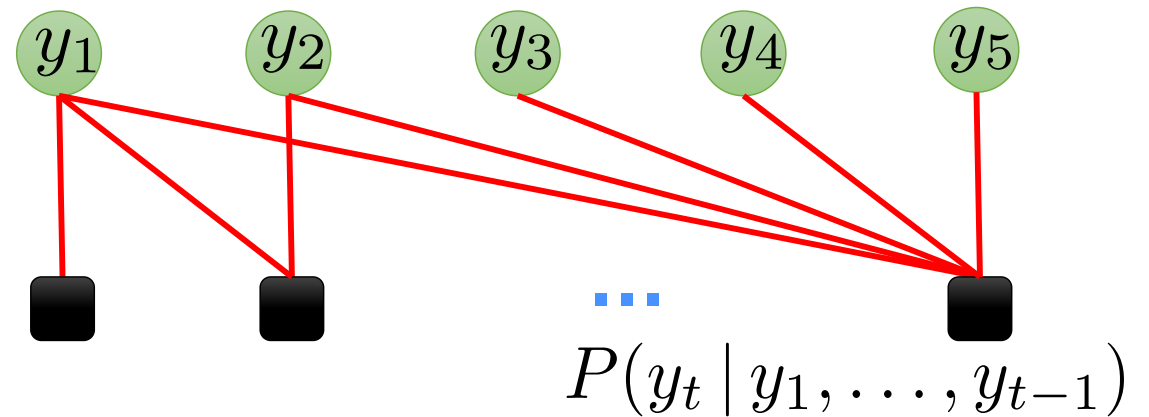


عنوان گذاری عصبی تصاویر



عنوان گذاری عصبی تصاویر

گراف مدل فاکتور دنباله

SEQUENCE MODEL FACTOR GRAPH

<http://dbs.cloudcv.org/captioning&mode=interactive>

عنوان‌گذاری روی تصاویر

IMAGE CAPTIONING

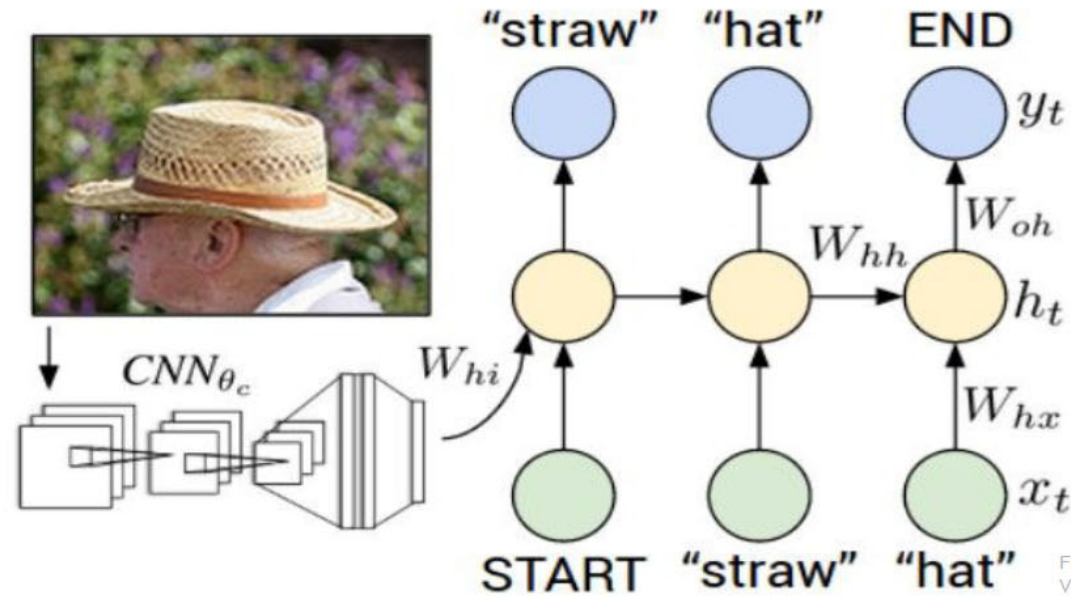


Figure from Karpathy et al, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015; figure copyright IEEE, 2015. Reproduced for educational purposes.

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

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

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

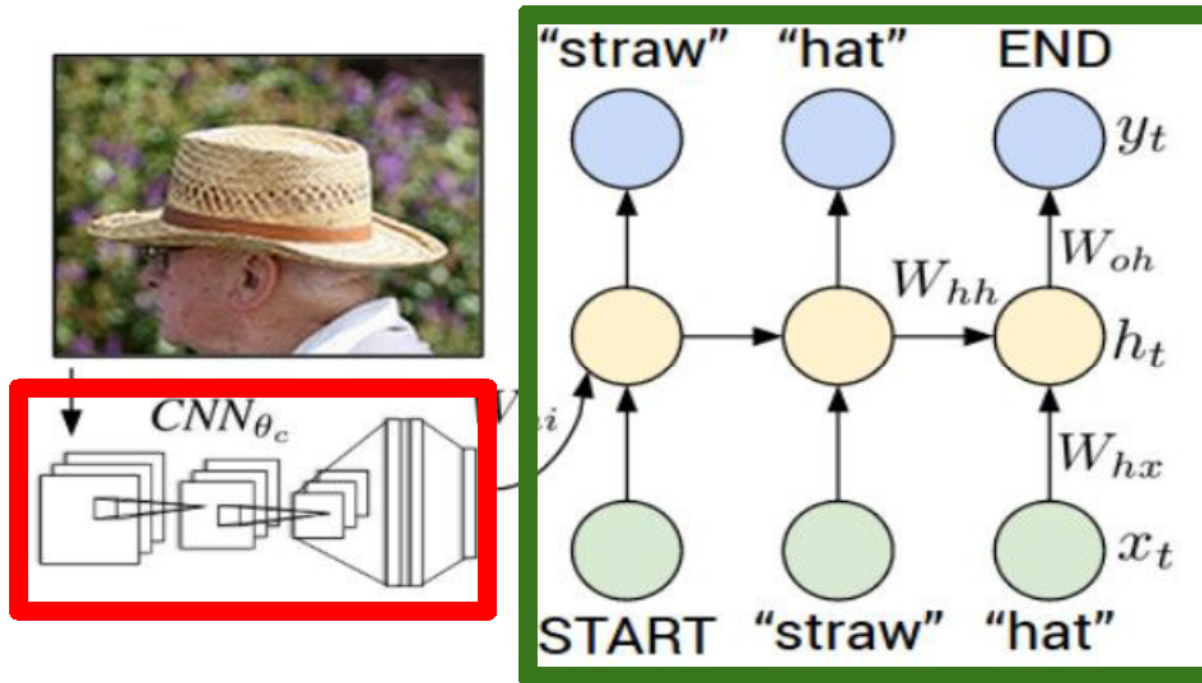
Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al.

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

عنوان‌گذاری روی تصاویر

IMAGE CAPTIONING

Recurrent Neural Network



Convolutional Neural Network

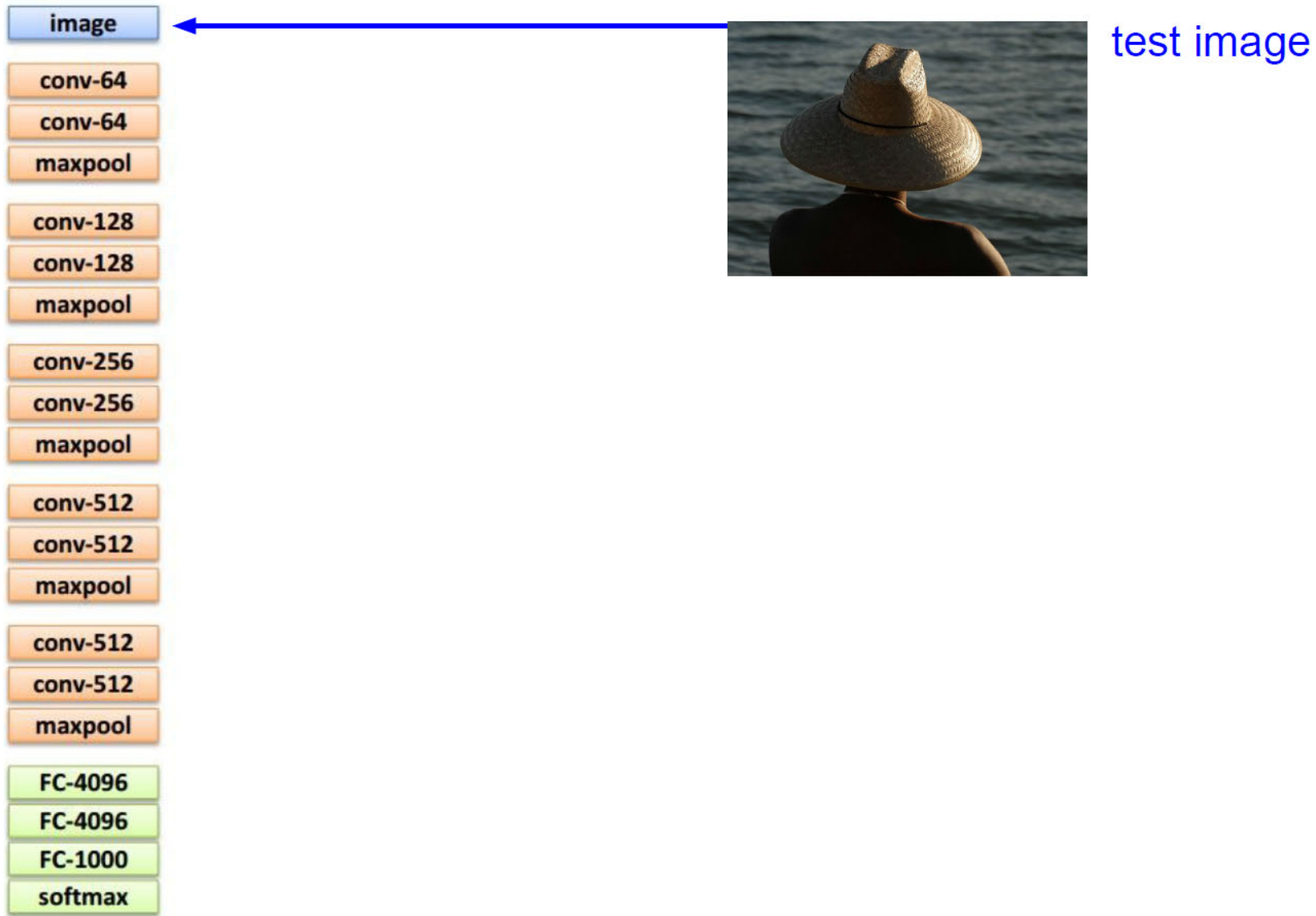
عنوان گذاری روی تصاویر



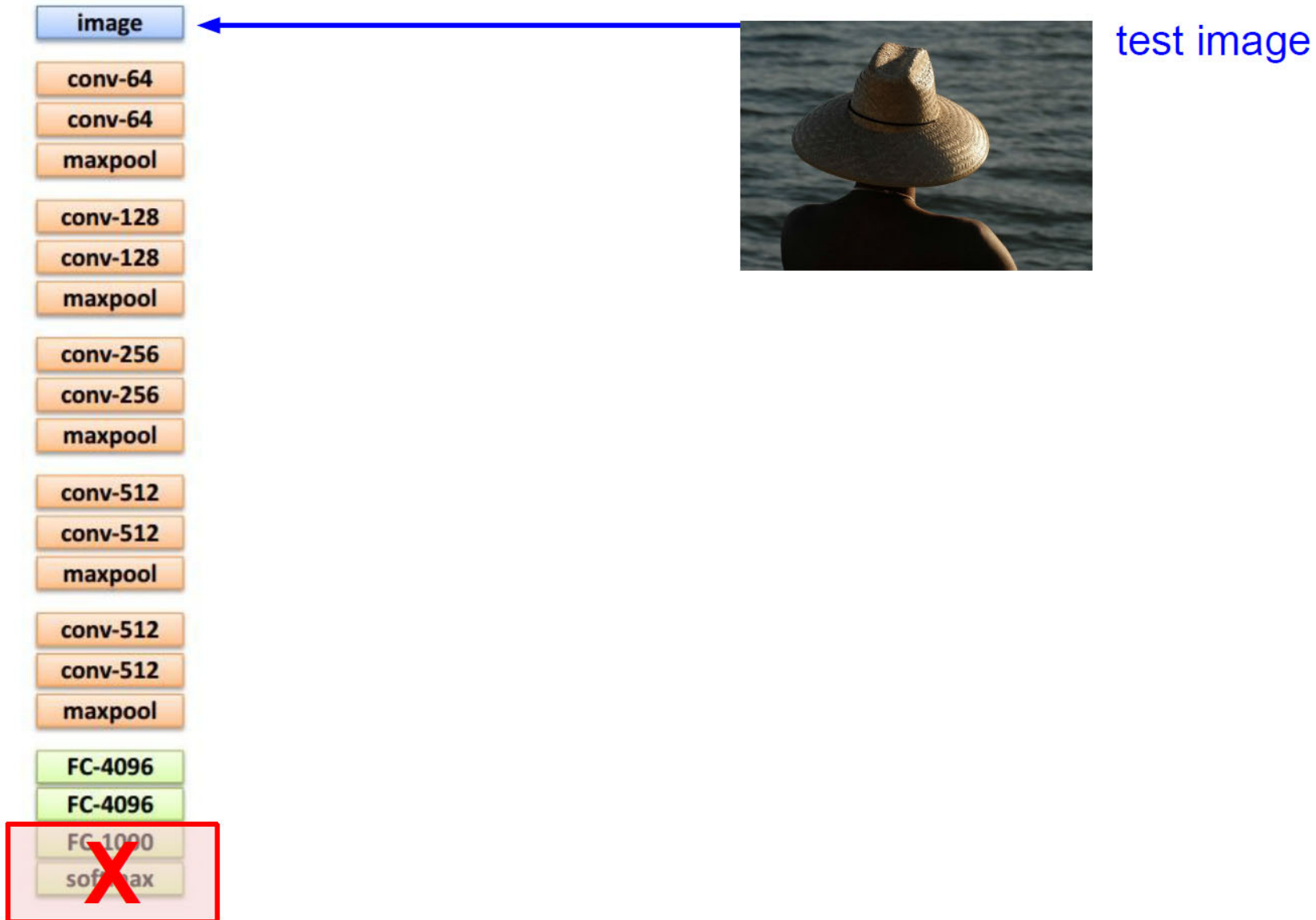
test image

[This image is CC0 public domain](#)

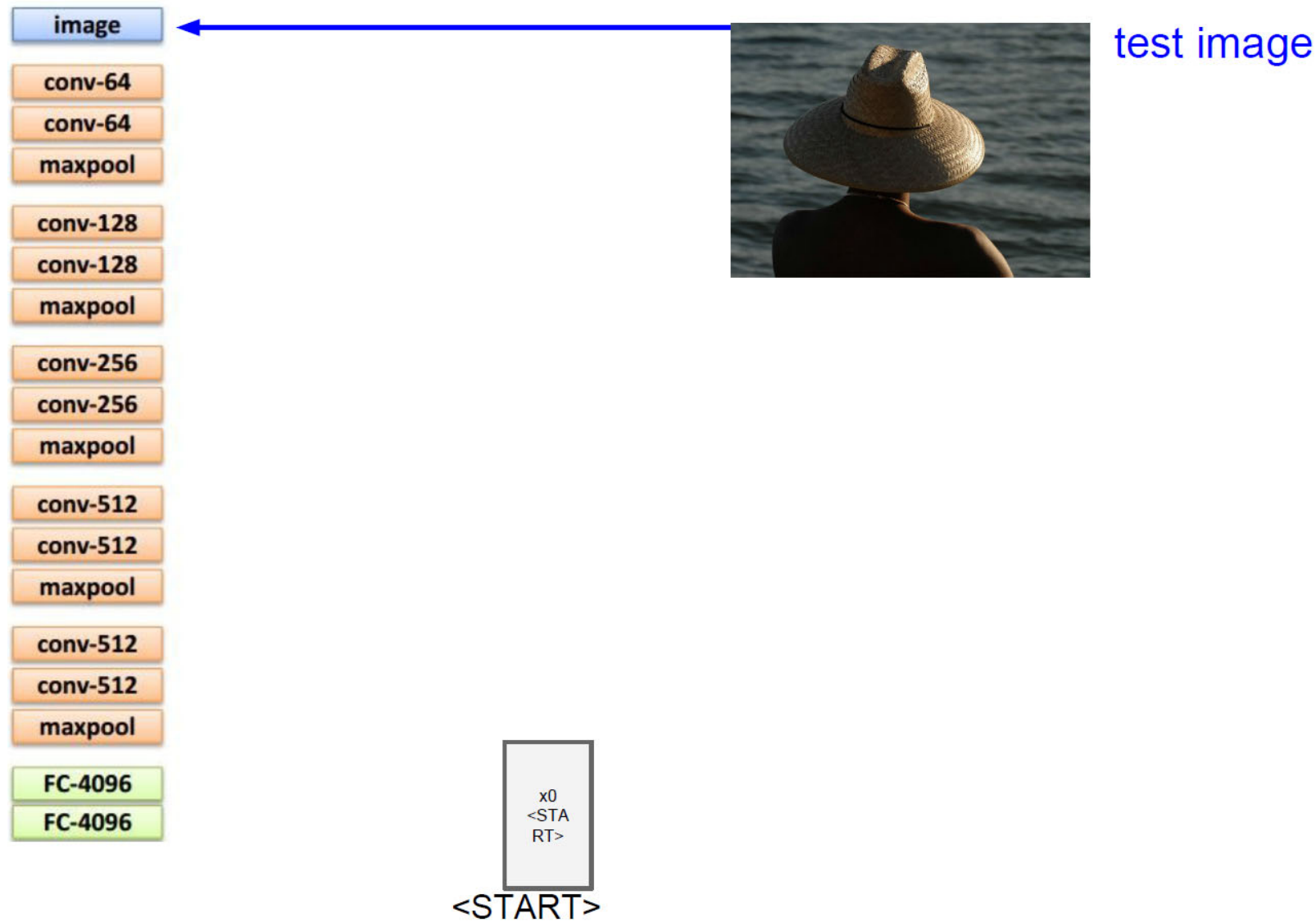
عنوان گذاری روی تصاویر



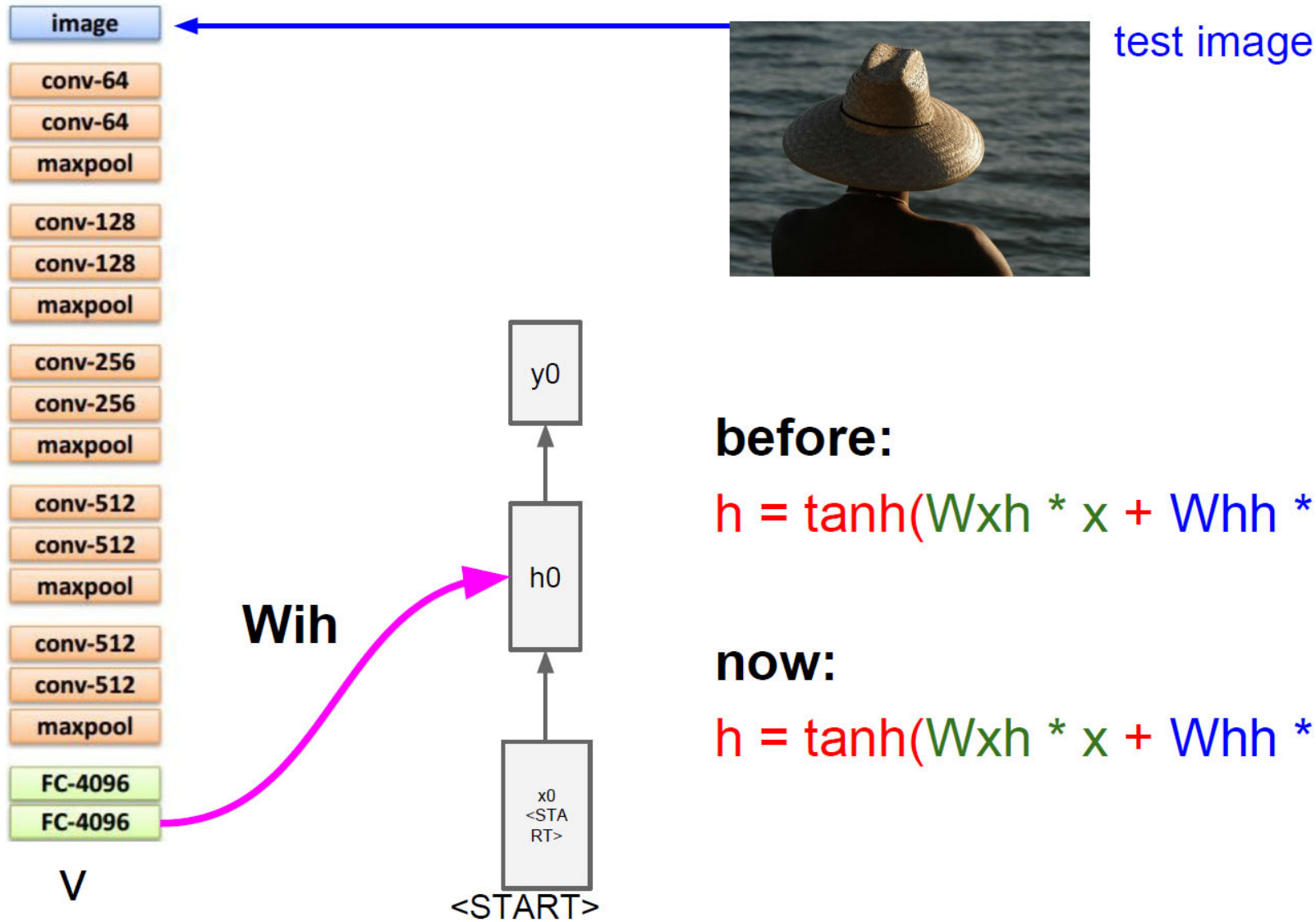
عنوان گذاری روی تصاویر



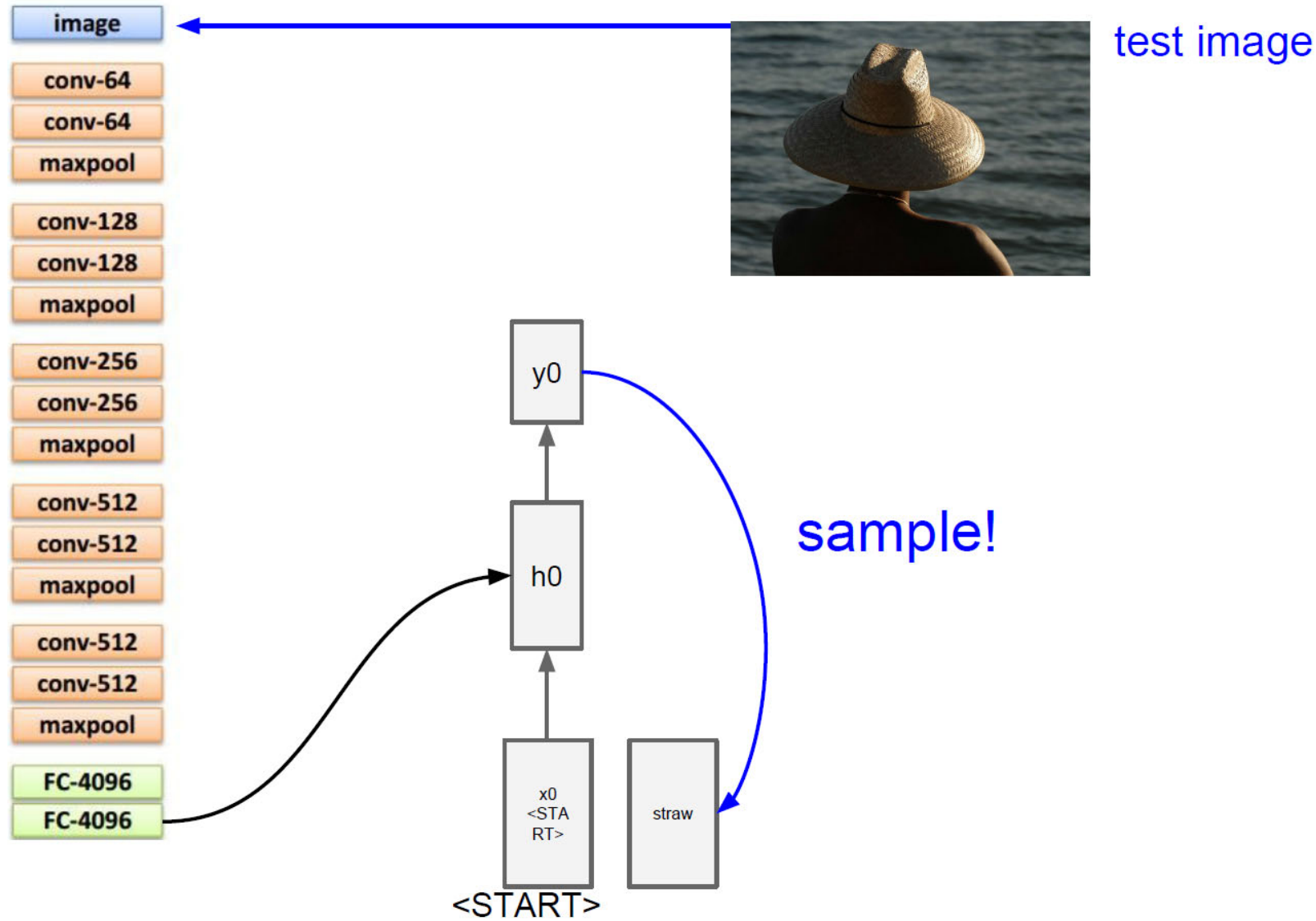
عنوان گذاری روی تصاویر



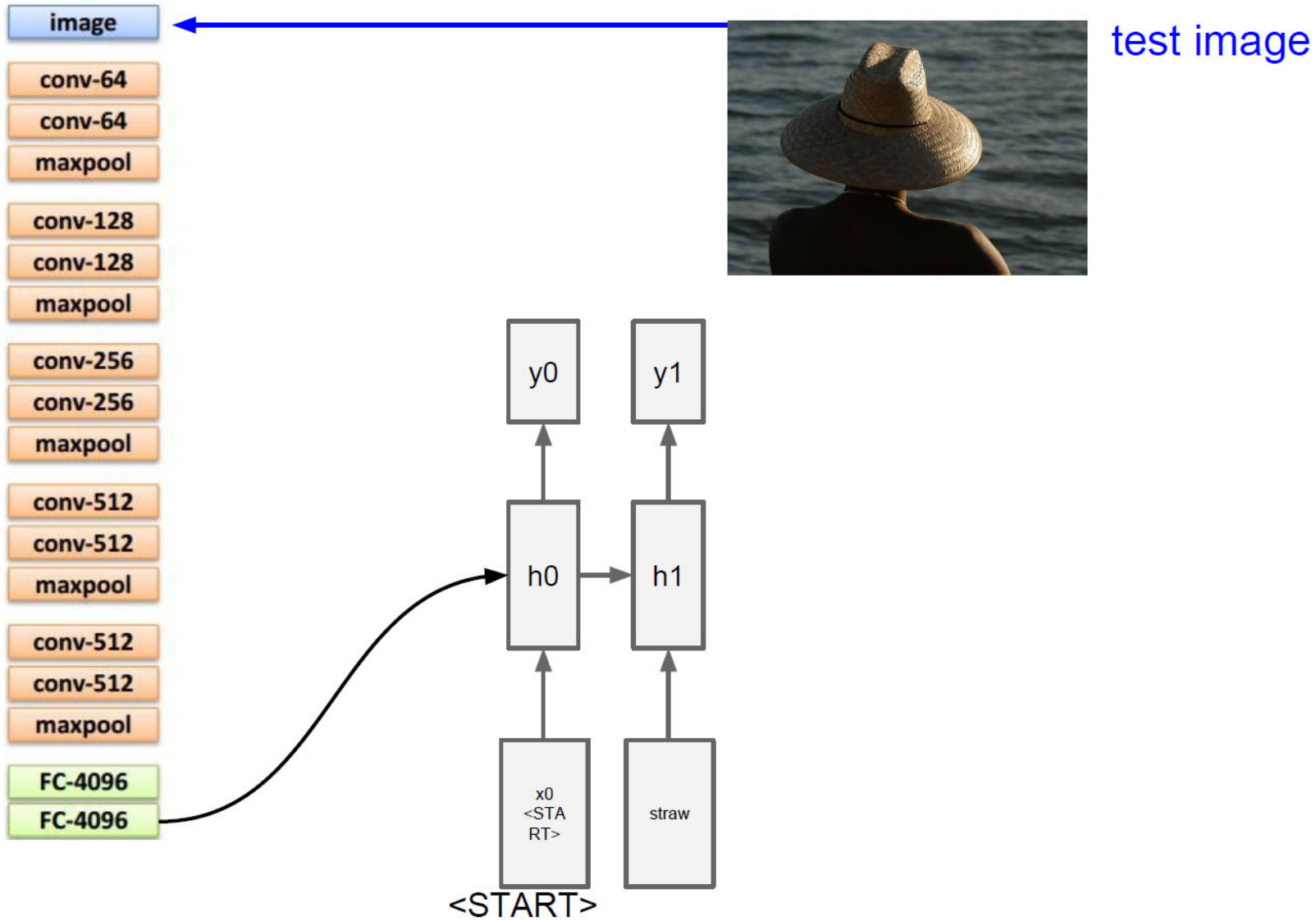
عنوان‌گذاری روی تصاویر



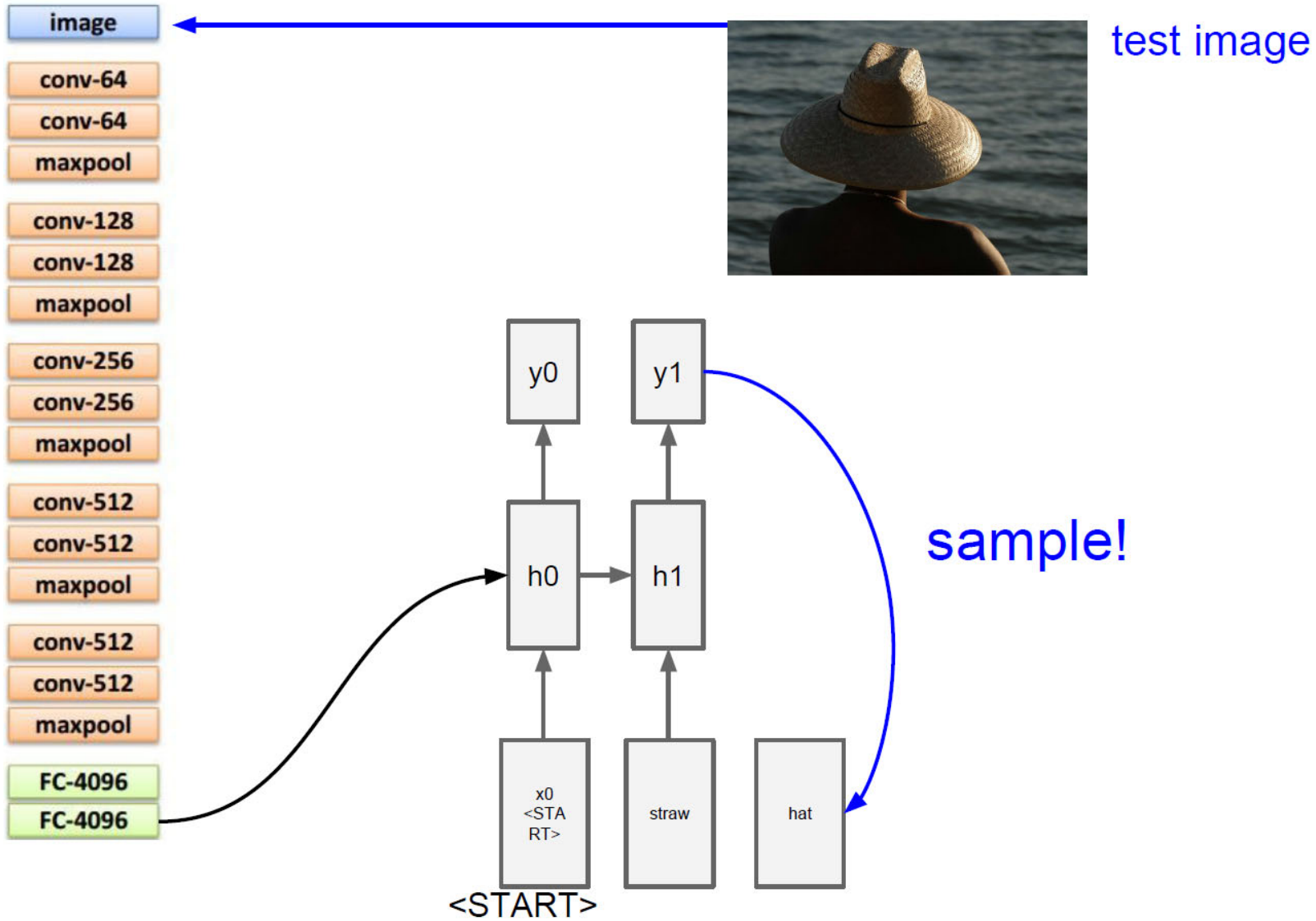
عنوان گذاری روی تصاویر



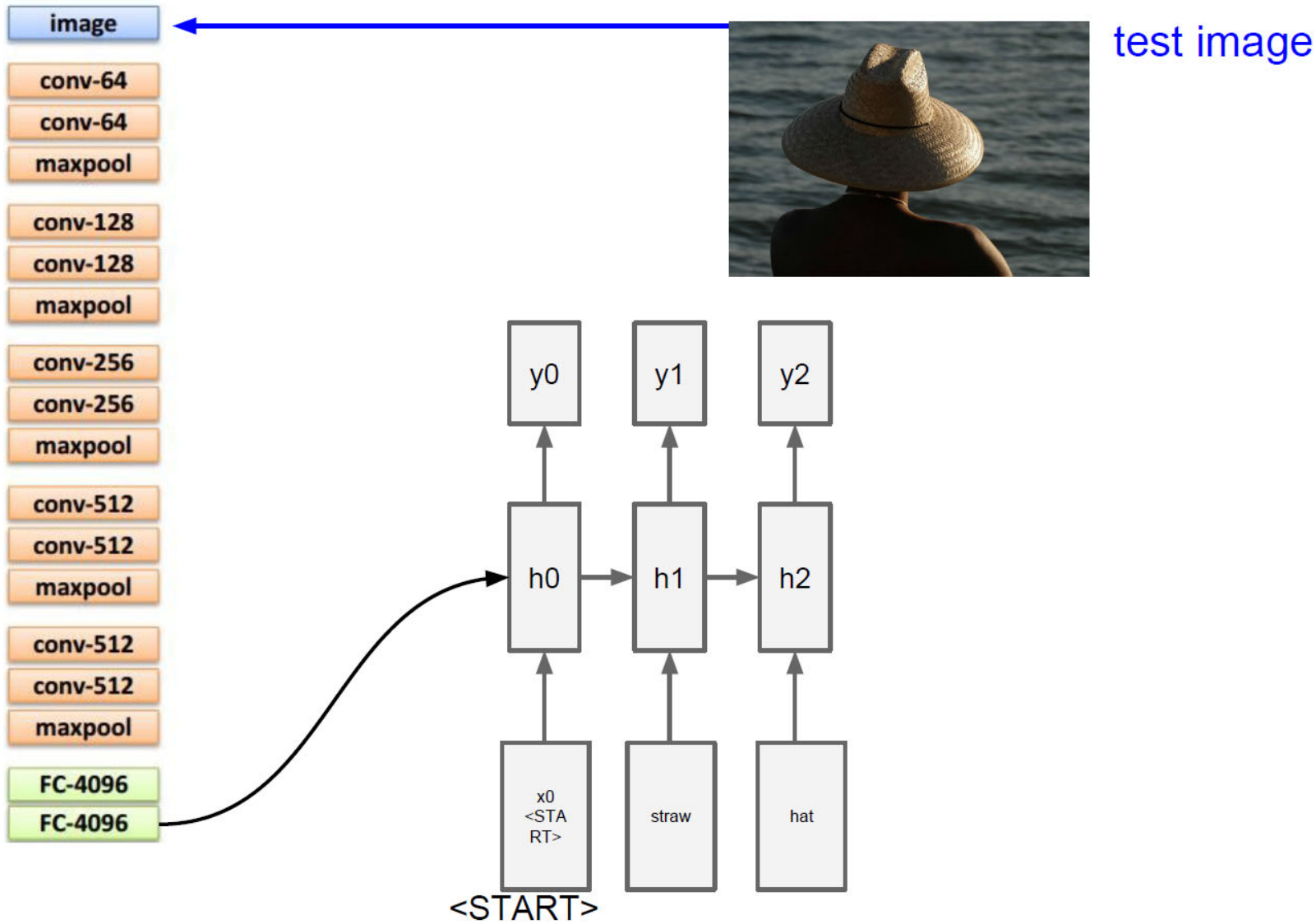
عنوان گذاری روی تصاویر



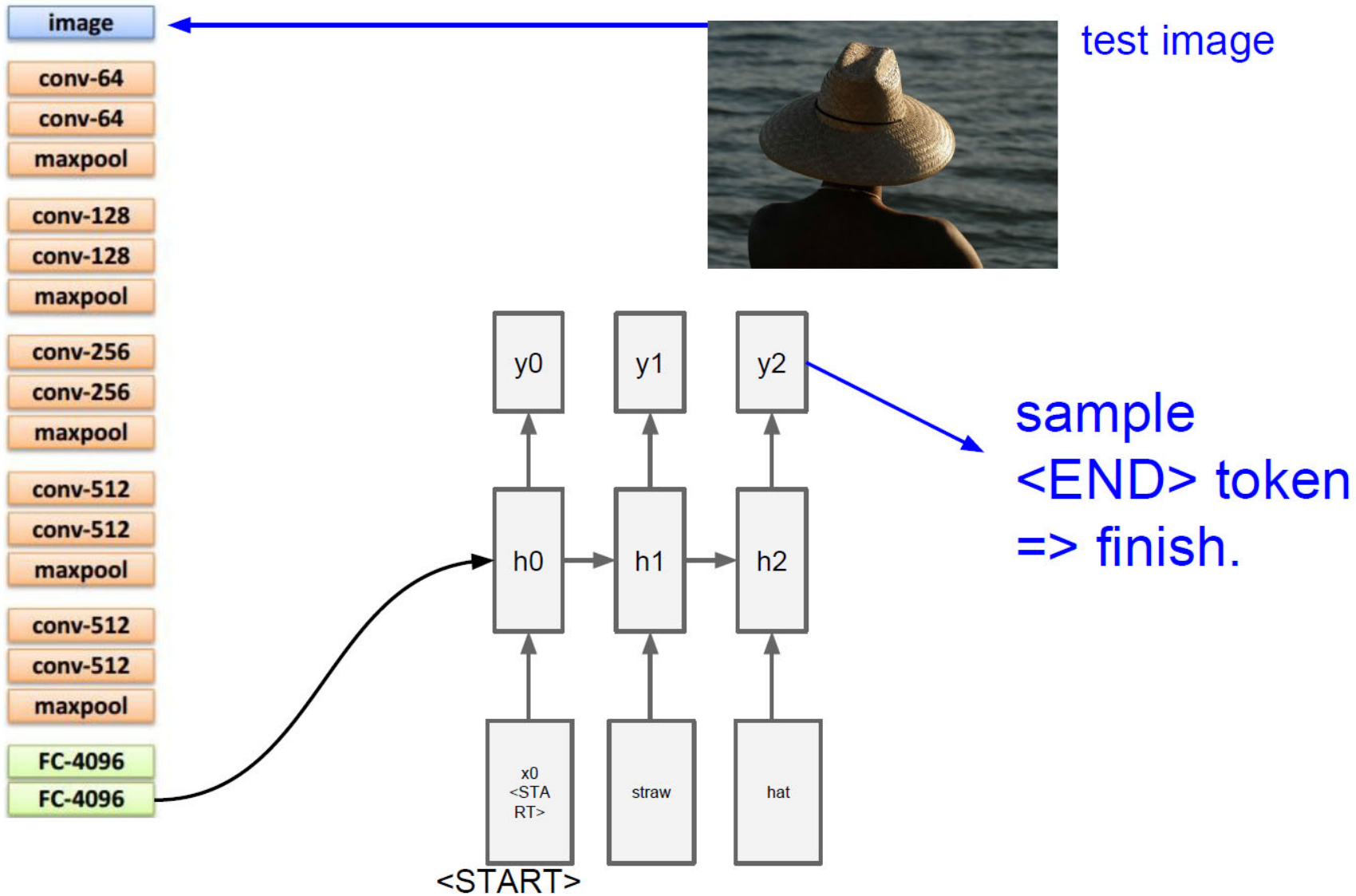
عنوان گذاری روی تصاویر



عنوان گذاری روی تصاویر



عنوان گذاری روی تصاویر



عنوان‌گذاری روی تصاویر

نتایج نمونه

Image Captioning: Example Results

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



A cat sitting on a suitcase on the floor



A cat is sitting on a tree branch



A dog is running in the grass with a frisbee



A white teddy bear sitting in the grass



Two people walking on the beach with surfboards



A tennis player in action on the court



Two giraffes standing in a grassy field



A man riding a dirt bike on a dirt track

عنوان‌گذاری روی تصاویر

نمونه‌های از موارد شکست

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

Image Captioning: Failure Cases



A woman is holding a cat in her hand



A bird is perched on a tree branch



A person holding a computer mouse on a desk

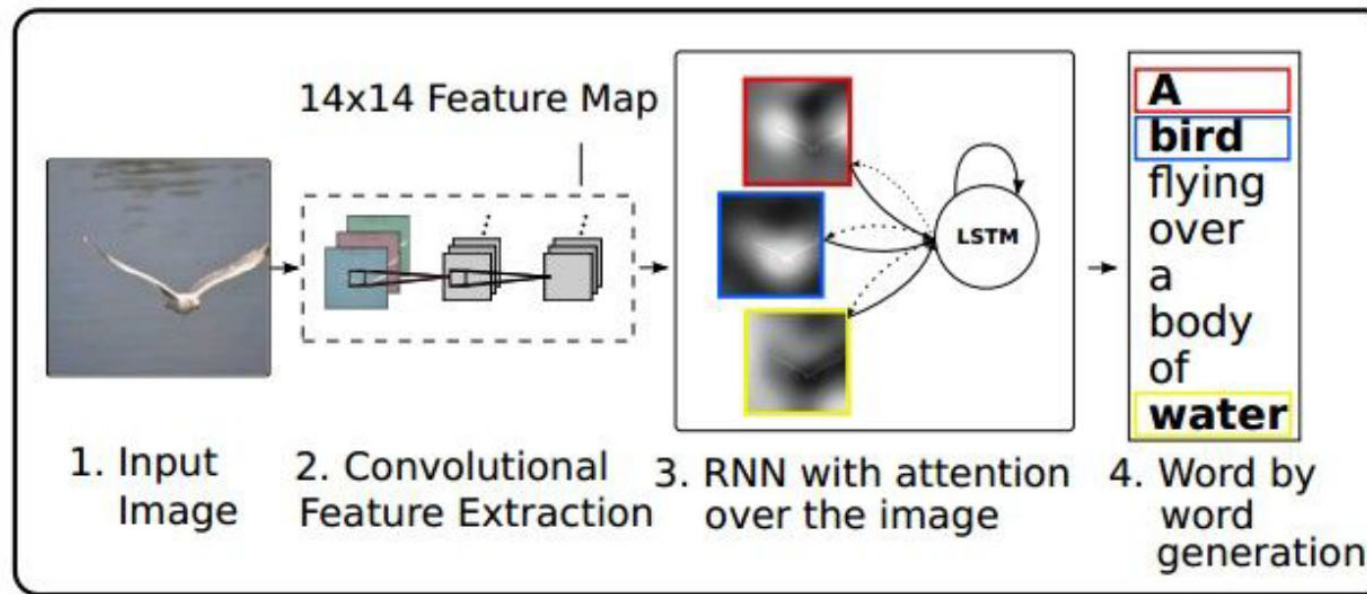


A man in a baseball uniform throwing a ball

عنوان گذاری روی تصاویر با استفاده از توجه

Image Captioning with Attention

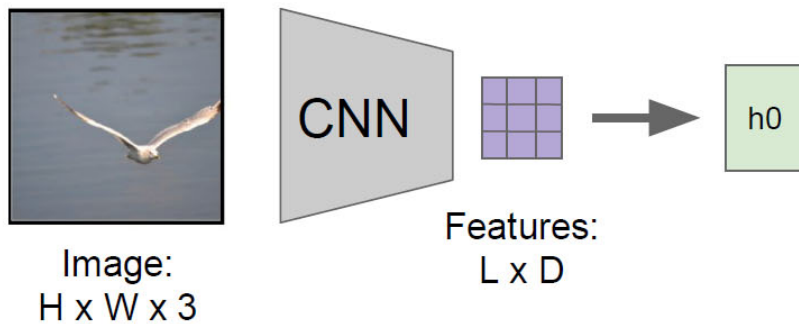
RNN focuses its attention at a different spatial location when generating each word



Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Figure copyright Kelvin Xu, Jimmy Lei Ba, Jamie Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Benchio, 2015. Reproduced with permission.

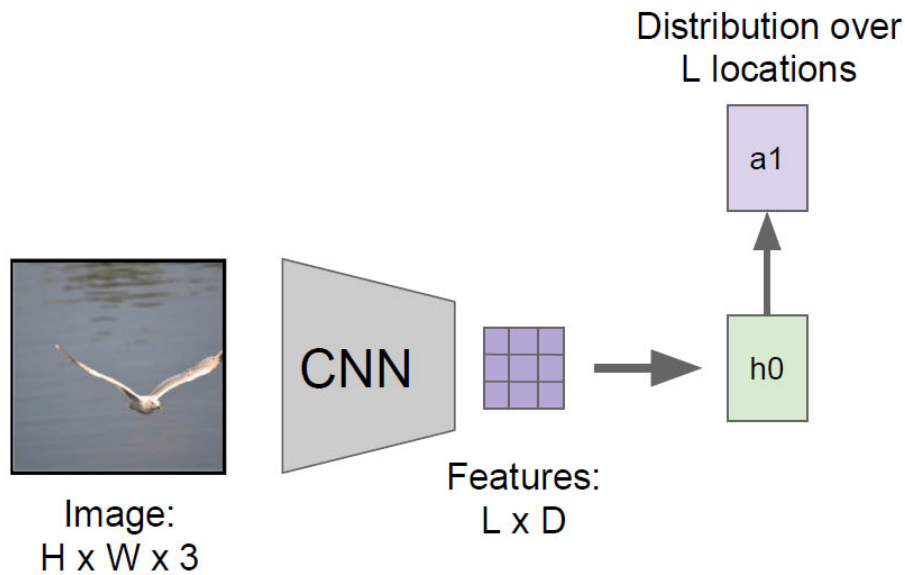
Image Captioning with Attention



Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

عنوان گذاری روی تصاویر با استفاده از توجه

Image Captioning with Attention



Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

عنوان گذاری روی تصاویر با استفاده از توجه

Image Captioning with Attention

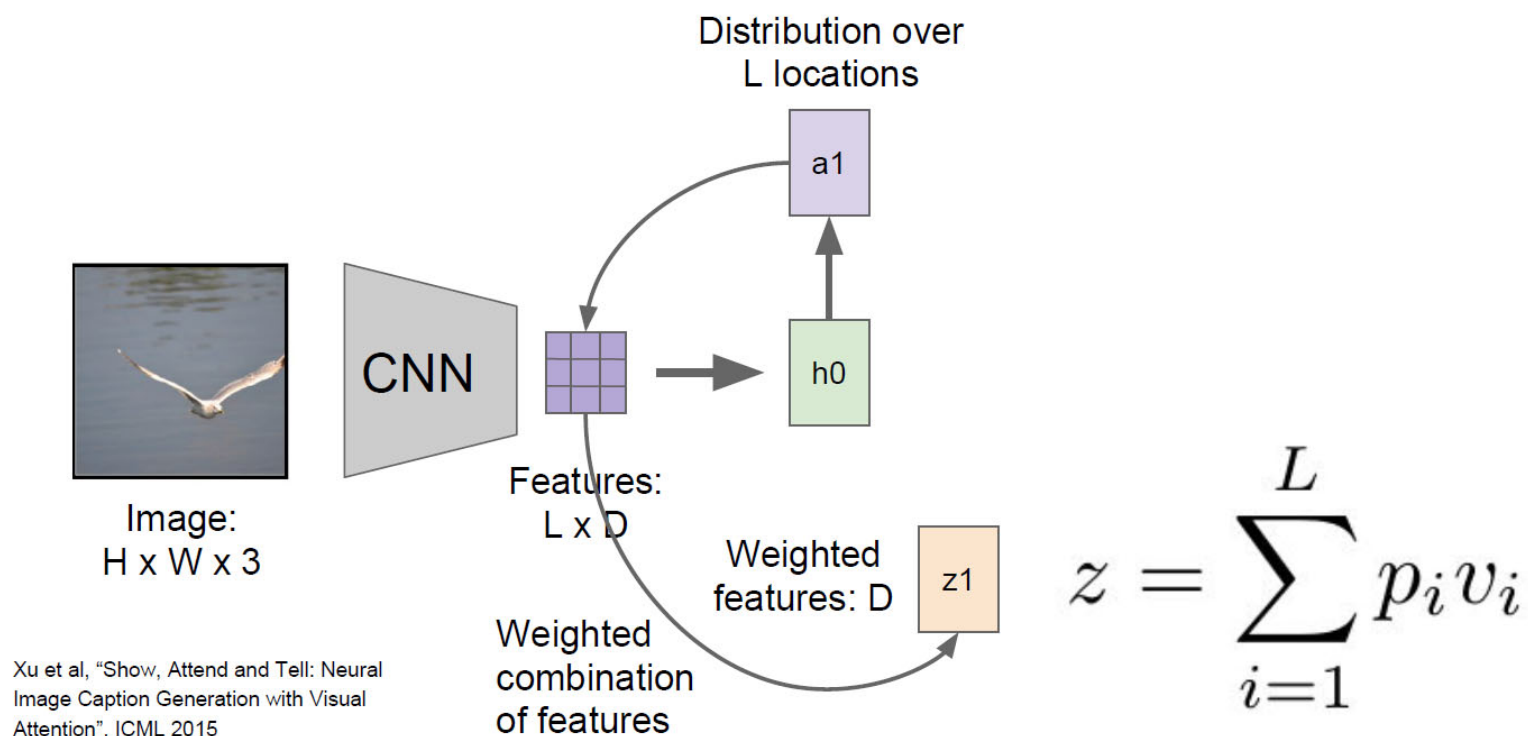
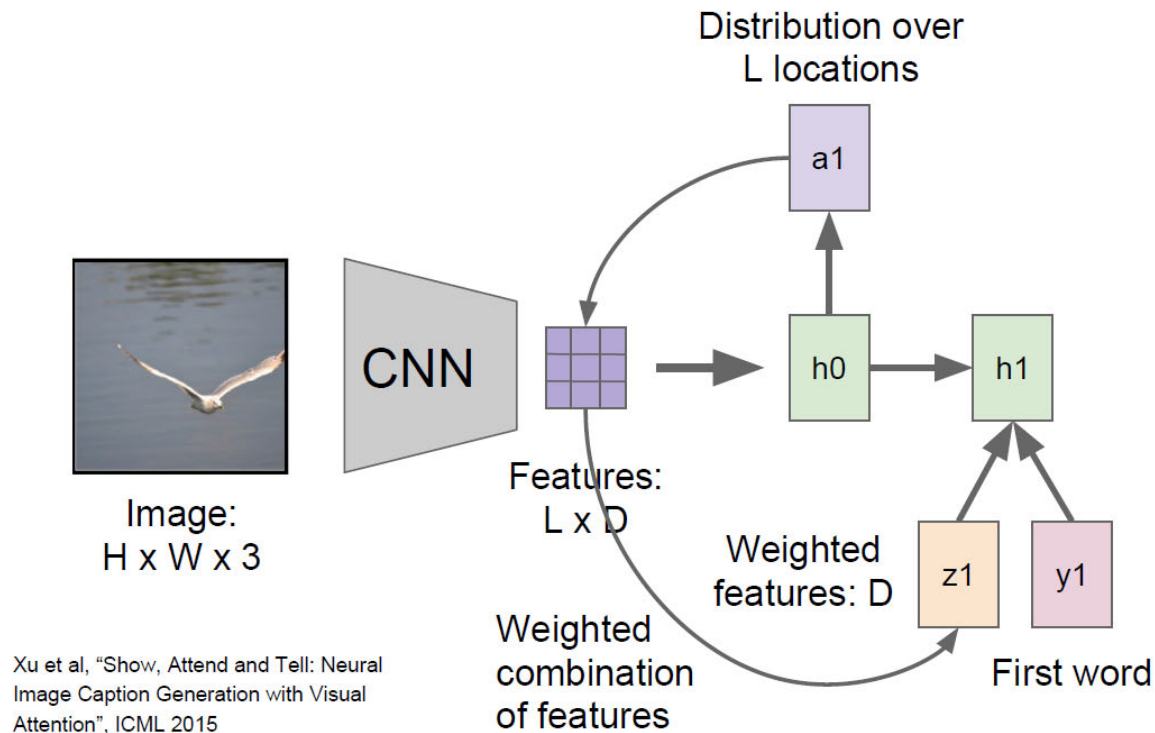


Image Captioning with Attention



عنوان‌گذاری روی تصاویر با استفاده از توجه

Image Captioning with Attention

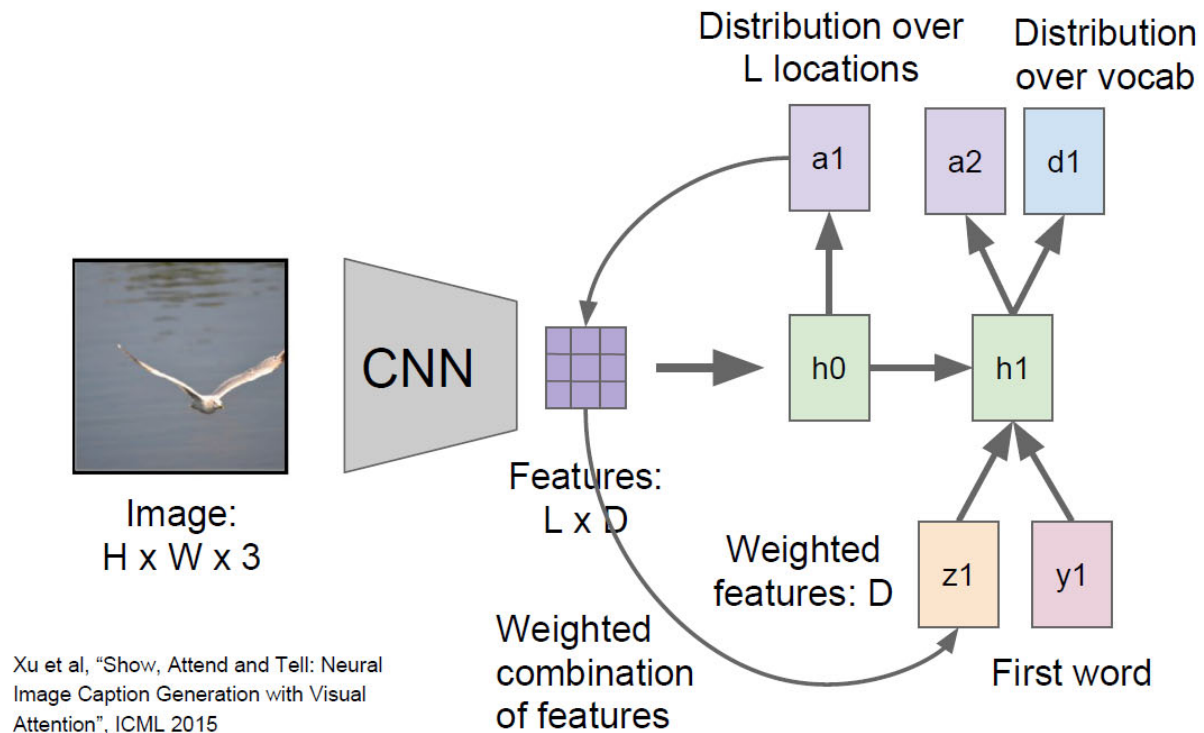
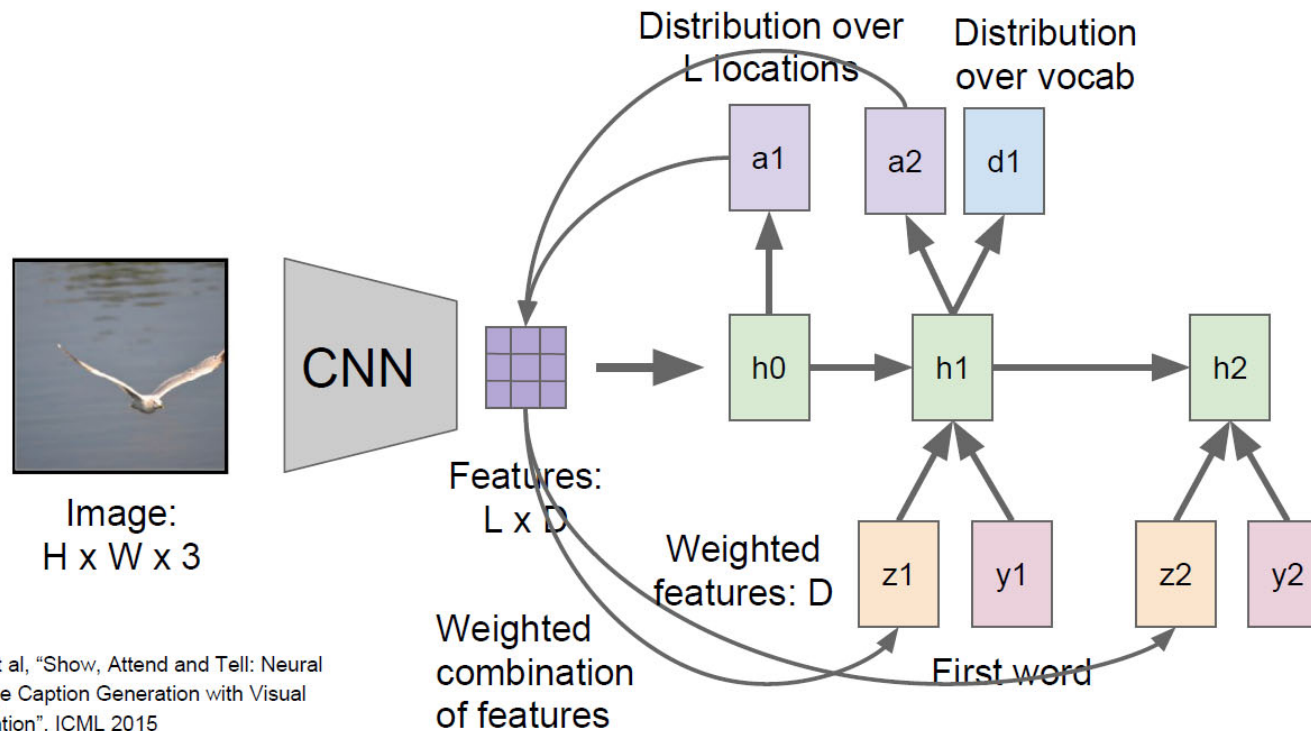
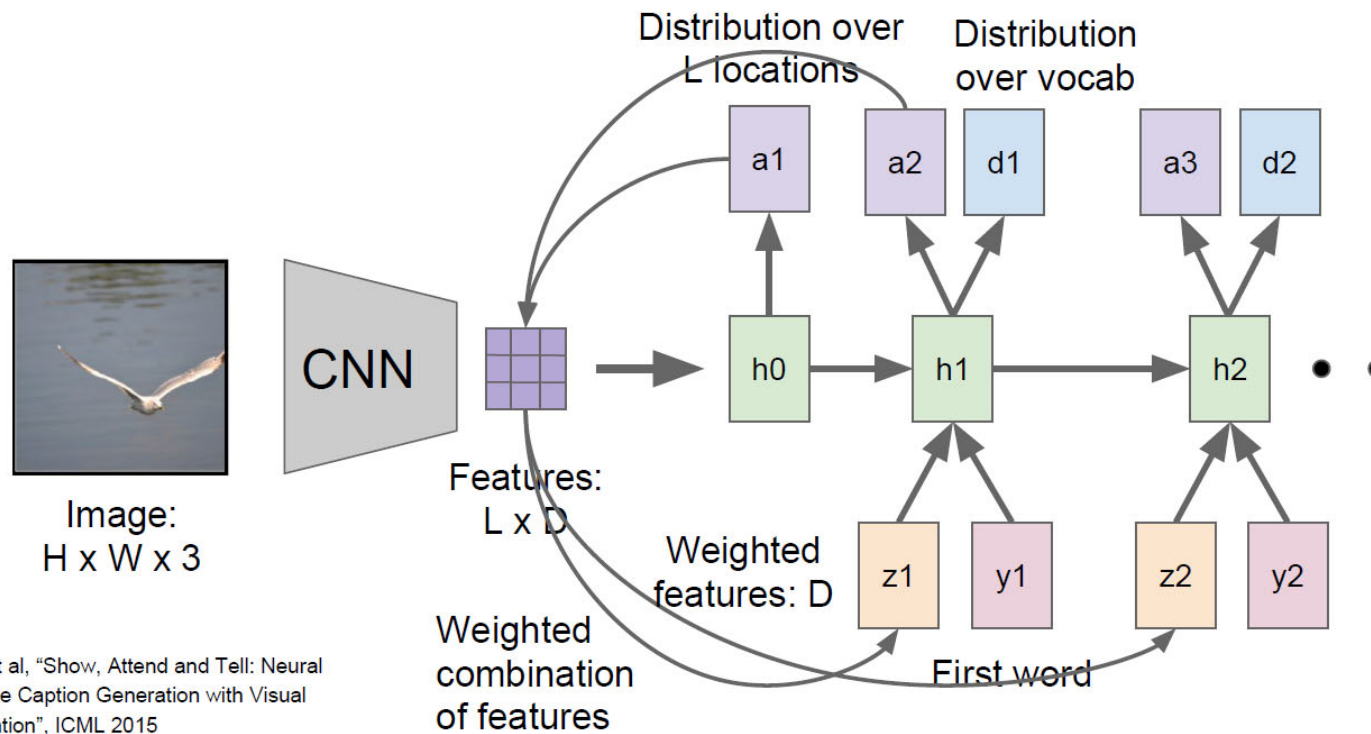


Image Captioning with Attention



Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

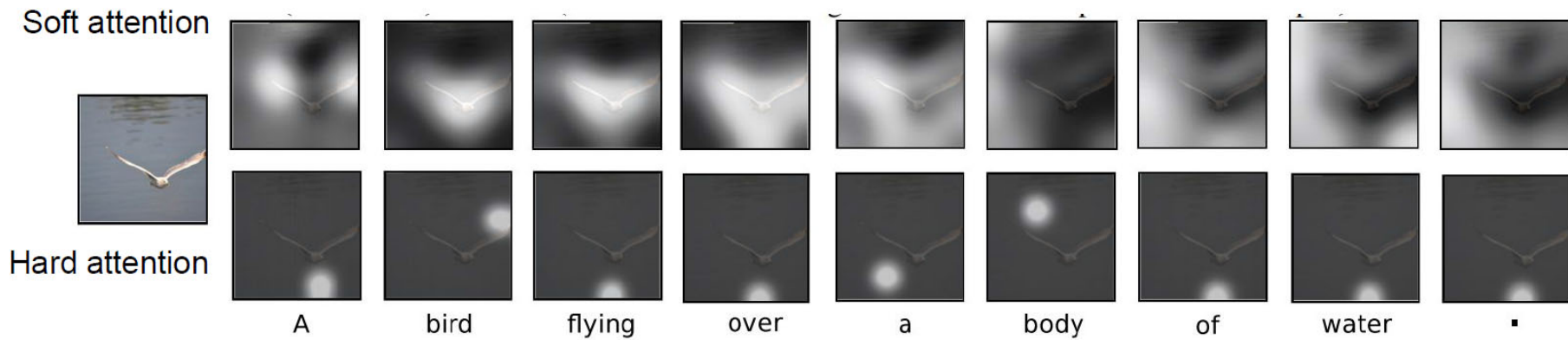
Image Captioning with Attention



Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

عنوان گذاری روی تصاویر با استفاده از توجه

Image Captioning with Attention



Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Figure copyright Kelvin Xu, Jimmy Lei Ba, Jamie Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Bengio, 2015. Reproduced with permission.

عنوان گذاری روی تصاویر با استفاده از توجه

نمونه‌هایی از نتایج

Image Captioning with Attention



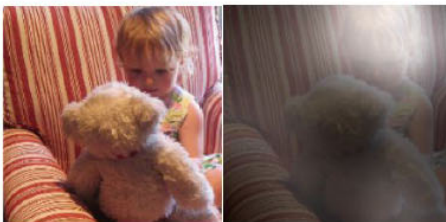
A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

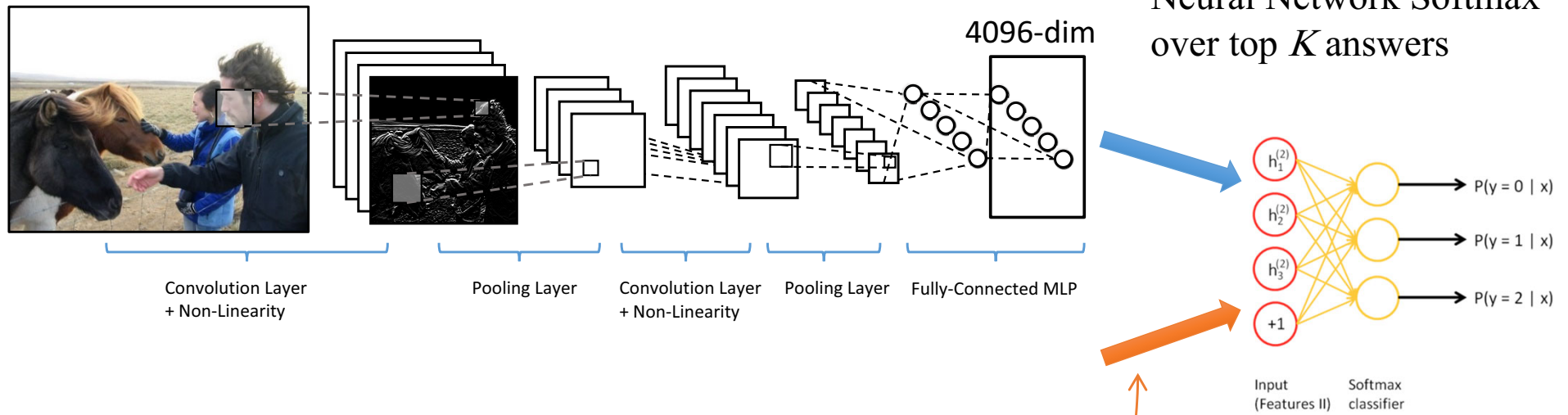
Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Figure copyright Kelvin Xu, Jimmy Lei Ba, Jamie Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Benchio, 2015. Reproduced with permission.

مدل‌های نوعی پاسخ به پرسش دیداری

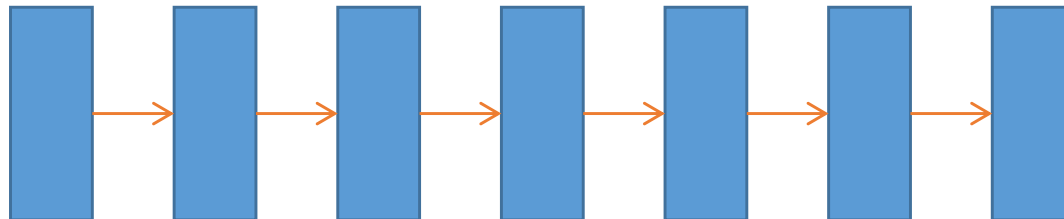
TYPICAL VQA MODELS

Image Embedding (VGGNet)



Question Embedding (LSTM)

“How many horses are in this image?”

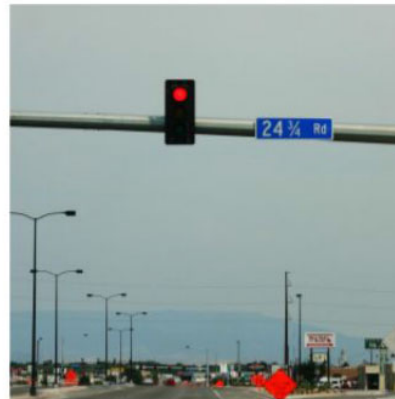


Visual Question Answering



Q: What endangered animal is featured on the truck?

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



Q: Where will the driver go if turning right?

- A: **Onto 24 1/4 Rd.**
- A: Onto 25 1/4 Rd.
- A: Onto 23 1/4 Rd.
- A: Onto Main Street.



Q: Who is under the umbrella?

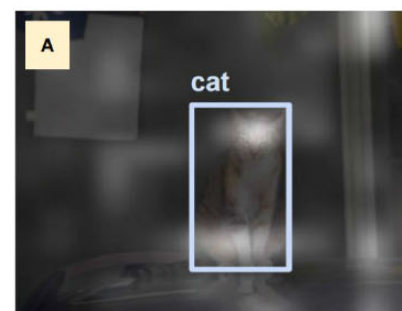
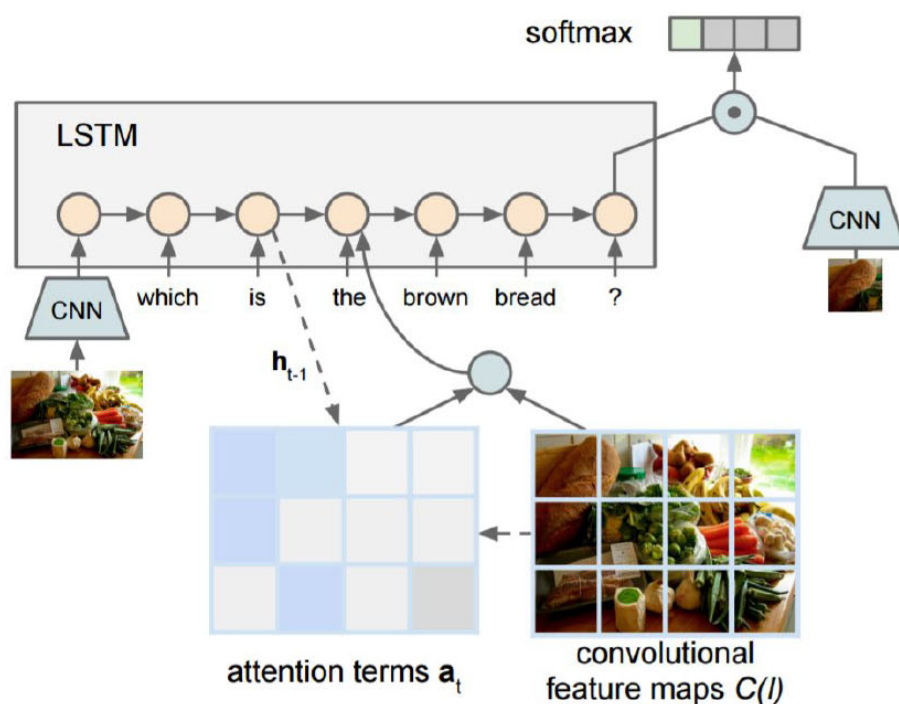
- A: **Two women.**
- A: A child.
- A: An old man.
- A: A husband and a wife.

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

پاسخ به پرسش دیداری

استفاده از شبکه‌های عصبی بازگشتی به همراه توجه

Visual Question Answering: RNNs with Attention



What kind of animal is in the photo?
A **cat**.



Why is the person holding a knife?
To cut the **cake** with.

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

پردازش دنباله‌ای داده‌های غیر دنباله‌ای

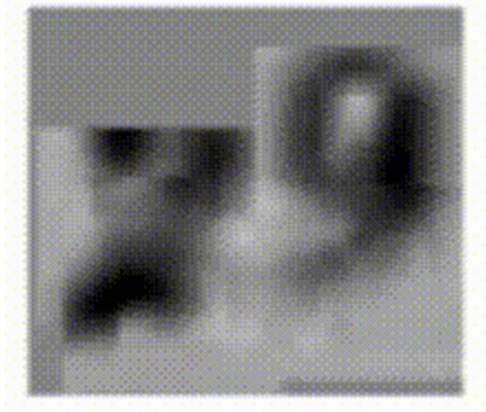
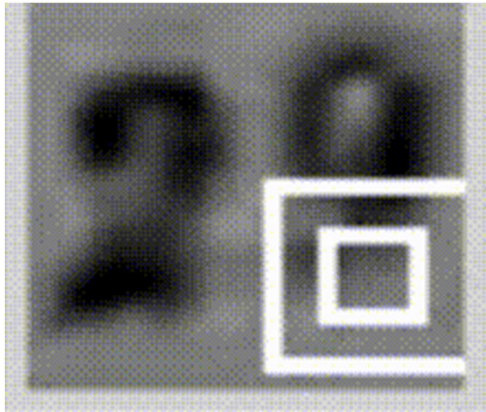
مثال

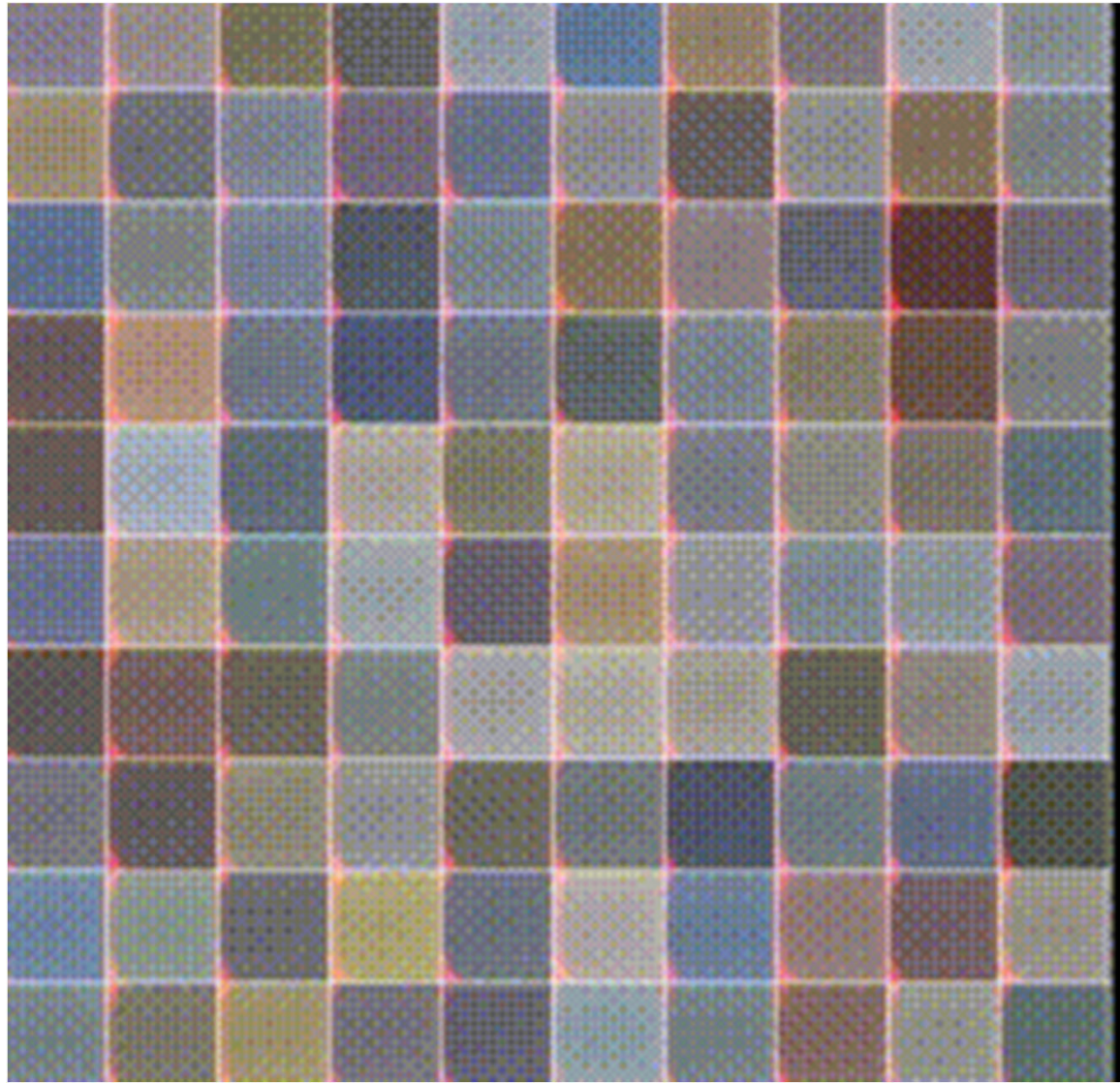
SEQUENTIAL PROCESSING OF NON-SEQUENCE DATA

Classify images by taking a series of “glimpses”

طبقه‌بندی تصاویر با دریافت
یک سری از گلیمپس‌ها







شبکه‌های عصبی بازگشتی

۳

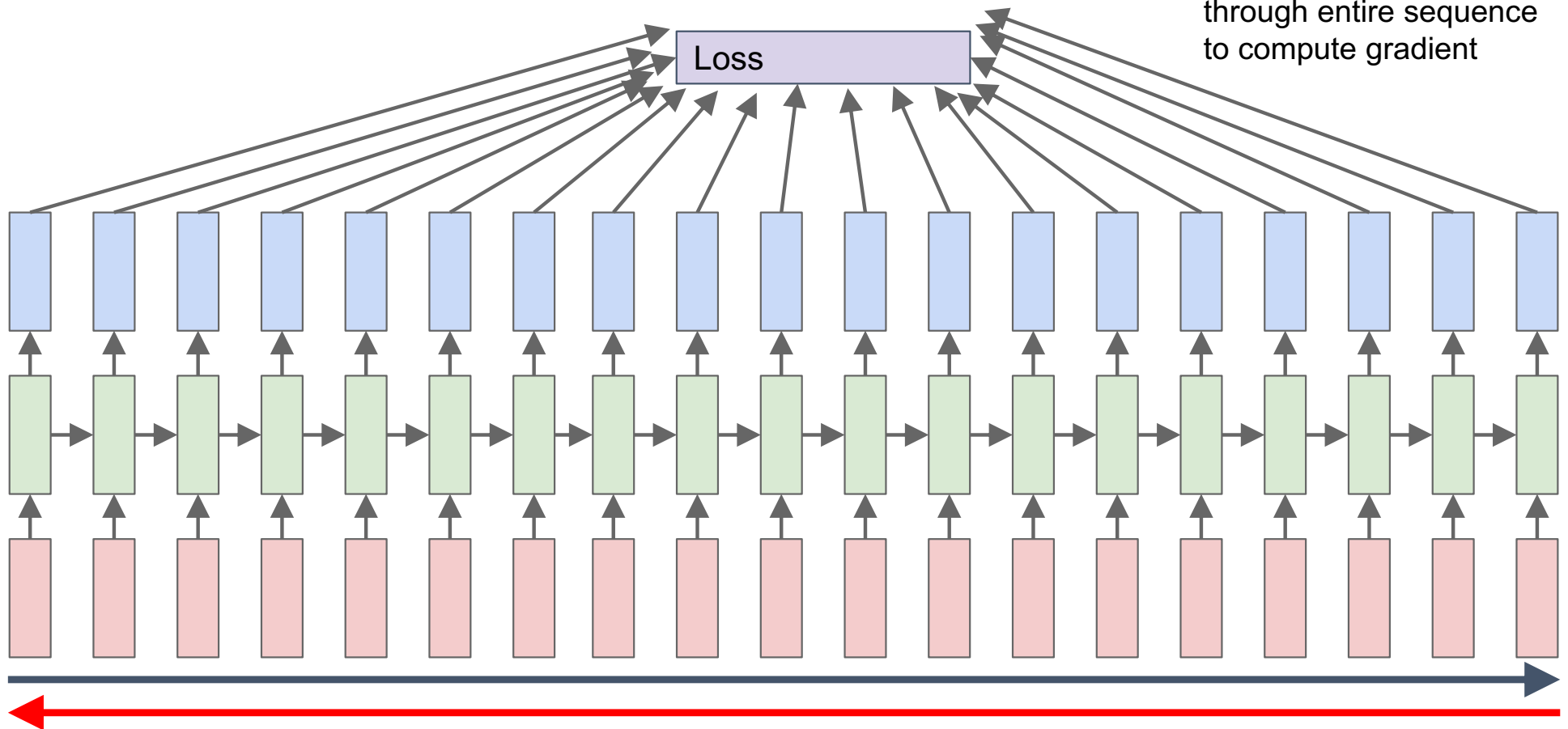
آموزش
شبکه‌های
عصبی
بازگشتی

آموزش شبکه‌های عصبی بازگشتی

پس‌انتشار در امتداد زمان

BACKPROPAGATION THROUGH TIME (BPTT)

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

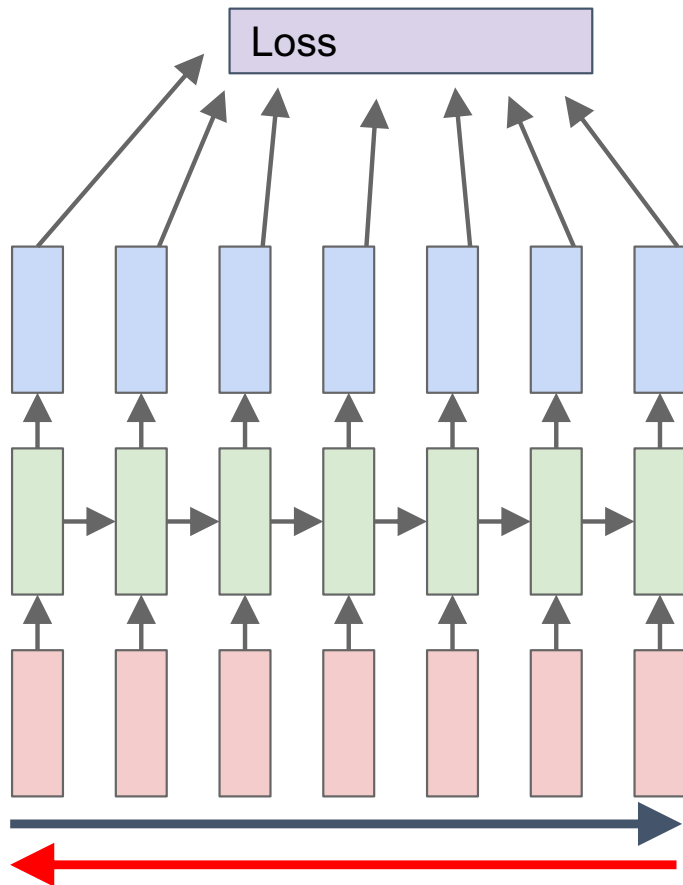


چون دنباله‌های قرار گرفته در یک دسته (batch) ممکن است طول‌های متفاوتی داشته باشند، پردازش دسته‌ای می‌تواند ناکارآمد شود.



آموزش شبکه‌های عصبی بازگشتی

پس‌انتشار برش‌یافته در امتداد زمان

TRUNCATED BACKPROPAGATION THROUGH TIME (TBPTT)

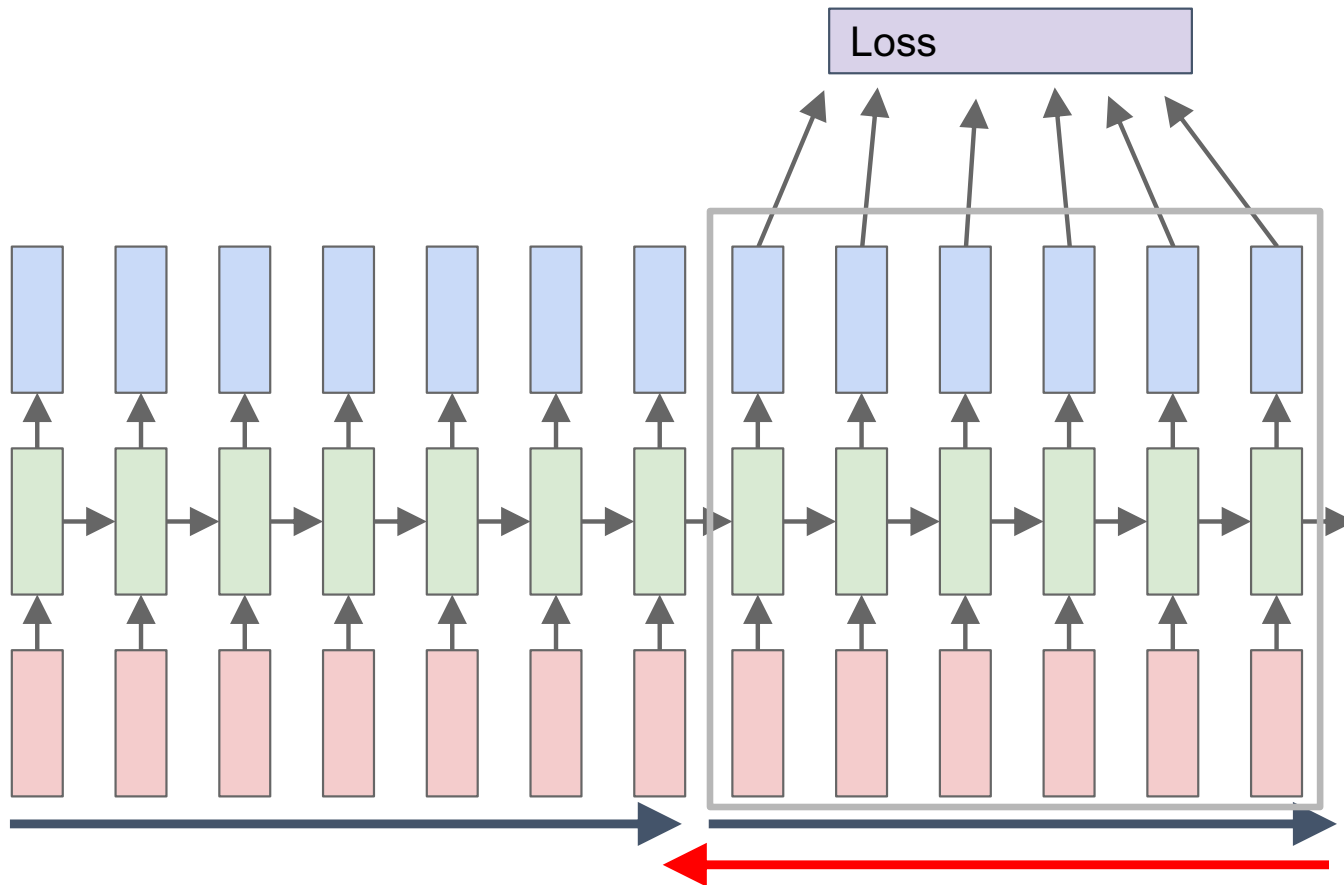
Run forward and backward through chunks of the sequence instead of whole sequence

چون دنباله‌های قرار گرفته در یک دسته (batch) طول‌های یکسانی دارند، پردازش دسته‌ای کارآمد می‌شود.



آموزش شبکه‌های عصبی بازگشتی

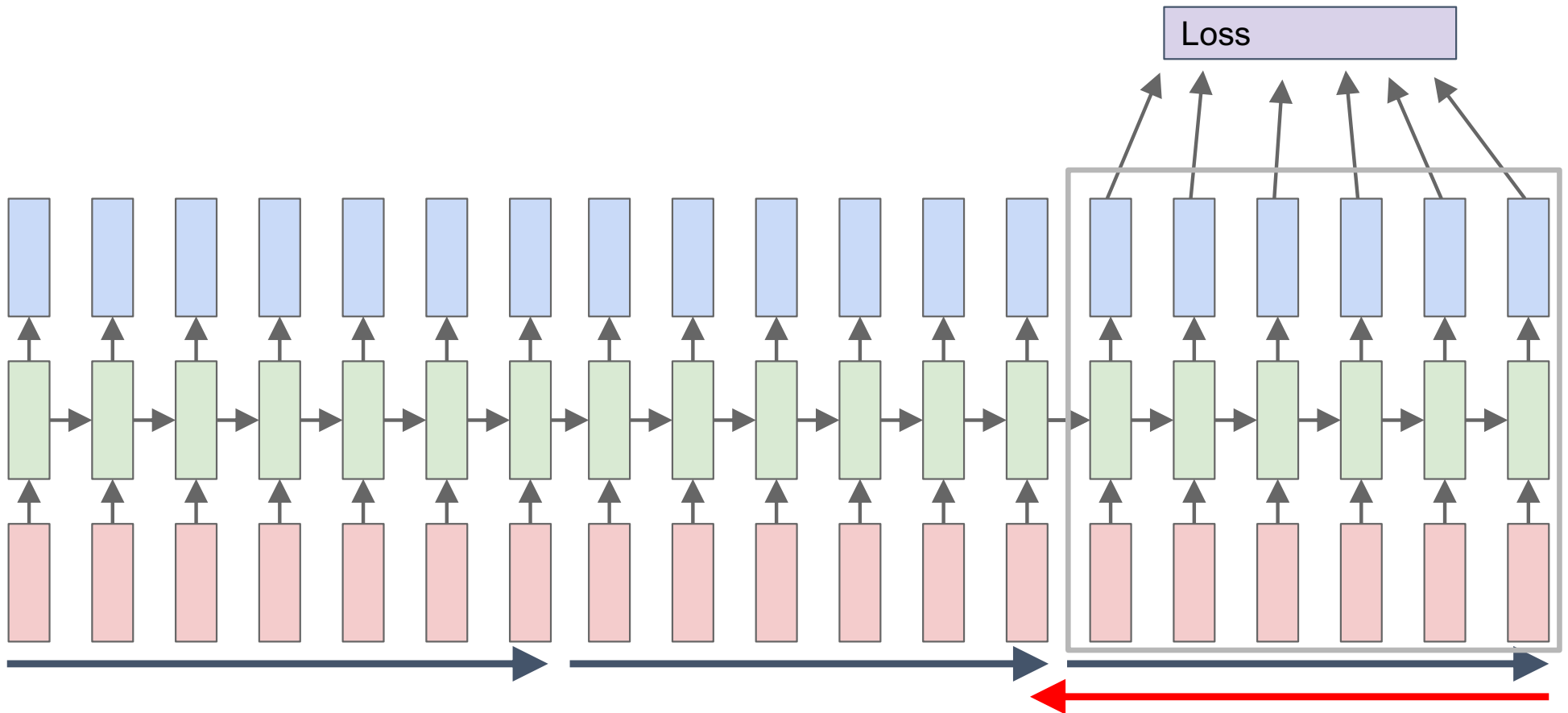
پس‌انتشار برش‌یافته در امتداد زمان

TRUNCATED BACKPROPAGATION THROUGH TIME (TBPTT)

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

آموزش شبکه‌های عصبی بازگشتی

پس‌انتشار برش‌یافته در امتداد زمان

TRUNCATED BACKPROPAGATION THROUGH TIME (TBPTT)

آموزش شبکه‌های عصبی بازگشتی

TRAINING RECURRENT NEURAL NETWORKS

از تابع اتلاف آنروپی متقابل استفاده می‌کنیم:
Cross-Entropy Loss

$$P = \prod_{t,k} y_{tk}^{l_{tk}} \Rightarrow \mathcal{L} = -\log P = \sum_t \mathcal{L}_t = -\frac{1}{T} \sum_t l_t \log y_t$$

الگوریتم مورد استفاده: **پس‌انتشار در طول زمان (Backpropagation Through Time: BPTT)**

- مجدداً از قاعده‌ی زنجیره‌ای استفاده می‌شود.
- تنها تفاوت: گرادیان‌ها بر روی گام‌های زمانی باقی می‌مانند.

آموزش شبکه‌های عصبی بازگشتی

پس‌انتشار در طول زمان: مثال (۱ از ۴)

BACKPROPAGATION THROUGH TIME: BPTT

$$\frac{\partial \mathcal{L}}{\partial V}, \frac{\partial \mathcal{L}}{\partial W}, \frac{\partial \mathcal{L}}{\partial U}$$

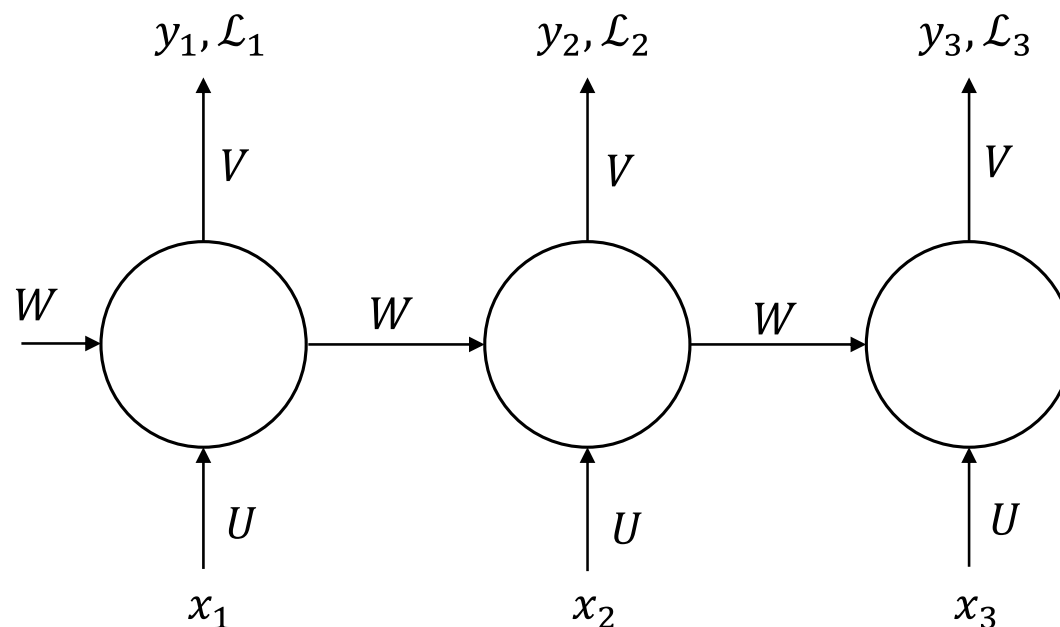
برای ساده‌تر کردن، بر روی گام ۳ تمرکز می‌کنیم:

$$\frac{\partial \mathcal{L}_3}{\partial V}, \frac{\partial \mathcal{L}_3}{\partial W}, \frac{\partial \mathcal{L}_3}{\partial U}$$

$$c_t = \tanh(U x_t + W c_{t-1})$$

$$y_t = \text{softmax}(V c_t)$$

$$\mathcal{L} = - \sum_t l_t \log y_t = \sum_t \mathcal{L}_t$$



آموزش شبکه‌های عصبی بازگشتی

پس‌انتشار در طول زمان: مثال (۲ از ۴)

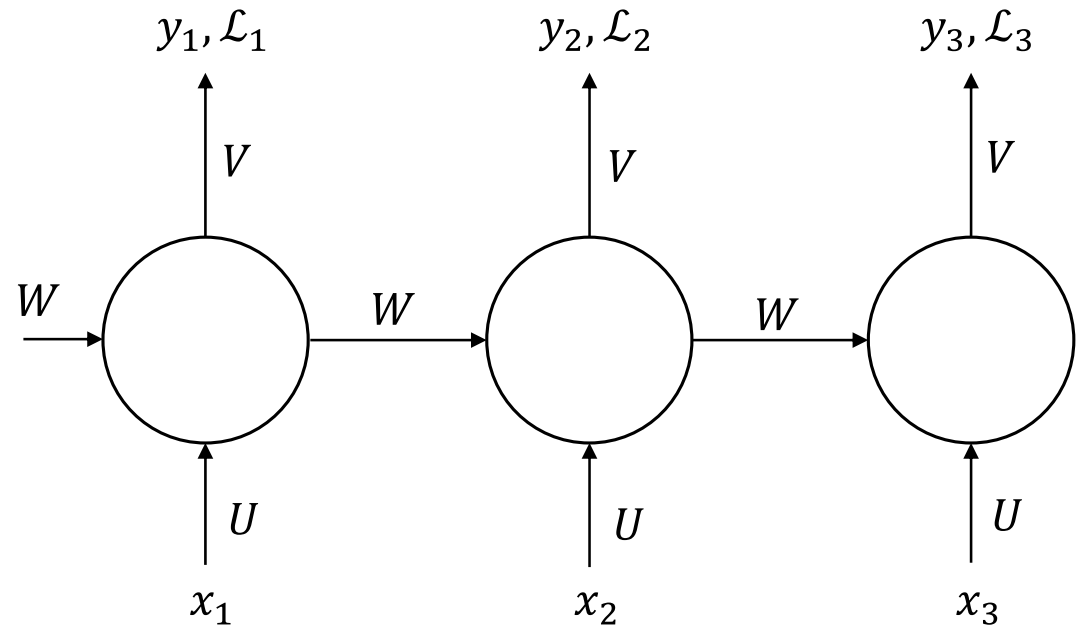
BACKPROPAGATION THROUGH TIME: BPTT

$$\frac{\partial \mathcal{L}_3}{\partial V} = \frac{\partial \mathcal{L}_3}{\partial y_3} \frac{\partial y_3}{\partial V} = (y_3 - l_3) \cdot c_3$$

$$c_t = \tanh(U x_t + W c_{t-1})$$

$$y_t = \text{softmax}(V c_t)$$

$$\mathcal{L} = - \sum_t l_t \log y_t = \sum_t \mathcal{L}_t$$



آموزش شبکه‌های عصبی بازگشتی

پس‌انتشار در طول زمان: مثال (۳ از ۴)

BACKPROPAGATION THROUGH TIME: BPTT

$$\frac{\partial \mathcal{L}_3}{\partial W} = \frac{\partial \mathcal{L}_3}{\partial y_3} \frac{\partial y_3}{\partial c_3} \frac{\partial c_3}{\partial W}$$

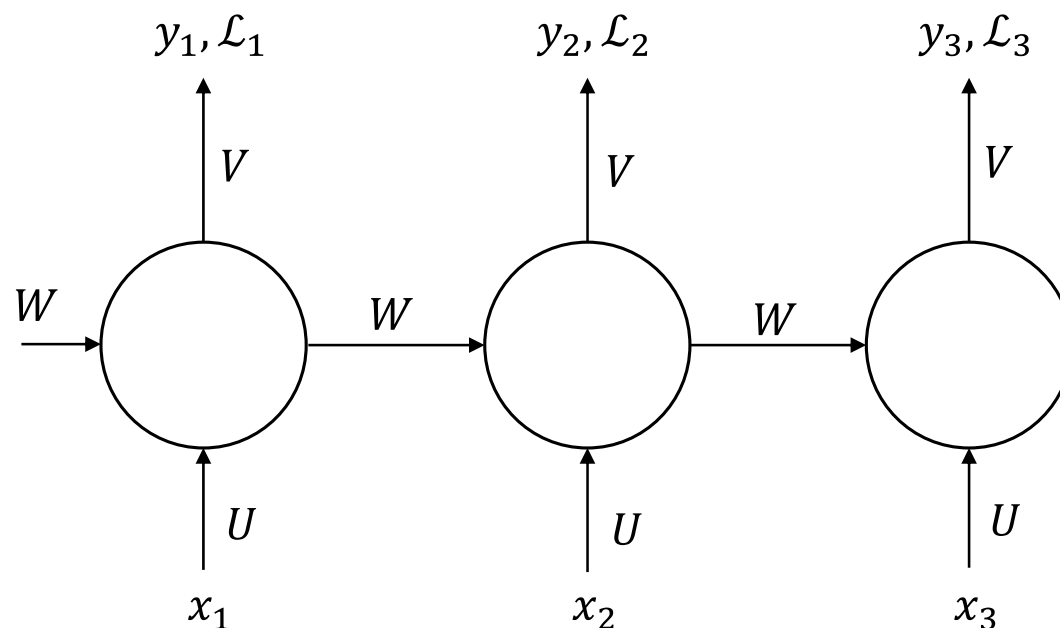
باید رابطه‌ی بین c_3 و W را به دست آوریم:

$$\text{Two-fold: } c_t = \tanh(U x_t + W c_{t-1})$$

$$\frac{\partial f(\varphi(x), \psi(x))}{\partial x} = \frac{\partial f}{\partial \varphi} \frac{\partial \varphi}{\partial x} + \frac{\partial f}{\partial \psi} \frac{\partial \psi}{\partial x}$$

$$\frac{\partial c_3}{\partial W} \propto c_2 + \frac{\partial c_2}{\partial W} \quad \left(\frac{\partial W}{\partial W} = 1 \right)$$

$$\begin{aligned} c_t &= \tanh(U x_t + W c_{t-1}) \\ y_t &= \text{softmax}(V c_t) \\ \mathcal{L} &= - \sum_t l_t \log y_t = \sum_t \mathcal{L}_t \end{aligned}$$



آموزش شبکه‌های عصبی بازگشتی

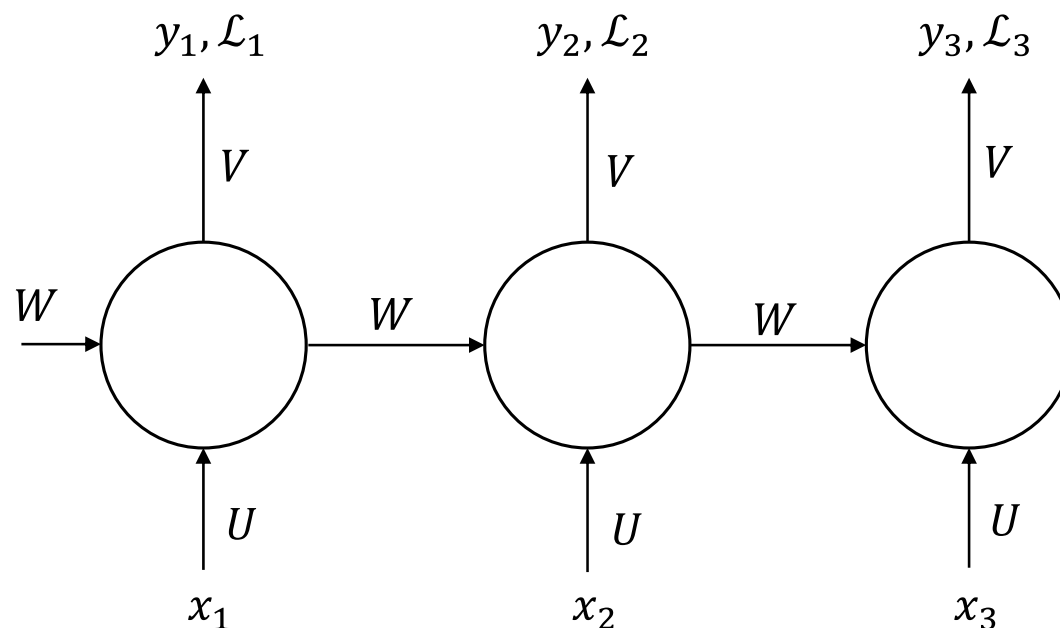
پس‌انتشار در طول زمان: مثال (۴ از ۴)

BACKPROPAGATION THROUGH TIME: BPTT

به صورت بازگشتی:

$$\left. \begin{aligned} \frac{\partial c_3}{\partial W} &= c_2 + \frac{\partial c_2}{\partial W} \\ \frac{\partial c_2}{\partial W} &= c_1 + \frac{\partial c_1}{\partial W} \\ \frac{\partial c_1}{\partial W} &= c_0 + \frac{\partial c_0}{\partial W} \end{aligned} \right\} \frac{\partial c_3}{\partial W} = \sum_{t=1}^3 \frac{\partial c_3}{\partial c_t} \frac{\partial c_t}{\partial W} \Rightarrow \frac{\partial \mathcal{L}_3}{\partial W} = \sum_{t=1}^3 \frac{\partial \mathcal{L}_3}{\partial y_3} \frac{\partial y_3}{\partial c_3} \frac{\partial c_3}{\partial c_t} \frac{\partial c_t}{\partial W}$$

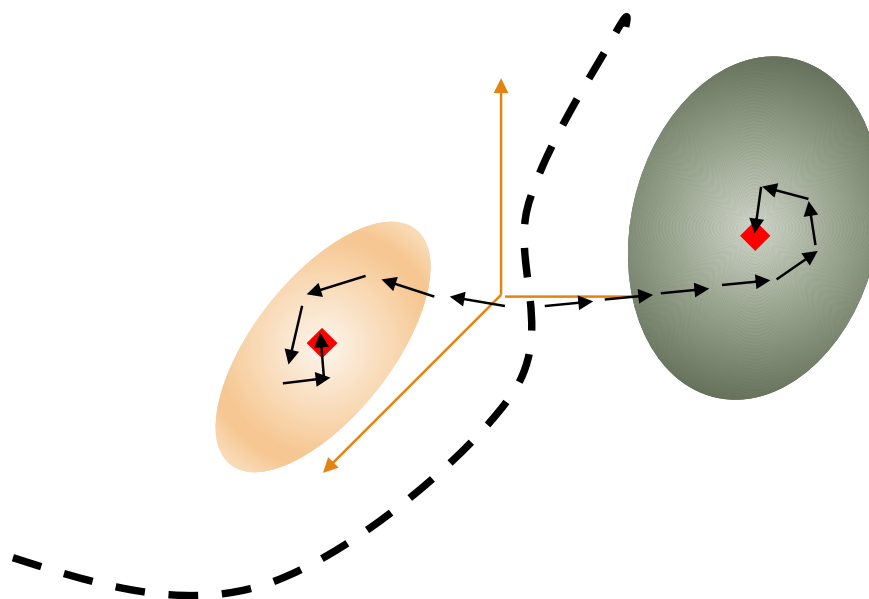
$$\begin{aligned} c_t &= \tanh(U x_t + W c_{t-1}) \\ y_t &= \text{softmax}(V c_t) \\ \mathcal{L} &= - \sum_t l_t \log y_t = \sum_t \mathcal{L}_t \end{aligned}$$



آموزش شبکه‌های عصبی بازگشتی

دشواری‌های آموزش

- فضای حافظه‌ی نهفته، از چندین بعد تشکیل شده است.
- یک زیرفضا از فضای حالت حافظه، می‌تواند اطلاعات را ذخیره کند، اگر چند بستر جذب در برخی ابعاد موجود باشد.
- گرادینان‌ها باید در نزدیک بستر جذب قوی باشند.



آموزش شبکه‌های عصبی بازگشتی

دشواری‌های آموزش

آموزش شبکه‌ی عصبی بازگشتی دشوار است، به دلیل:

- **اضمحلال گرادیان‌ها (Vanishing Gradients)**
پس از چند گام زمانی، گرادیان‌ها تقریباً صفر می‌شوند.
- **انفجار گرادیان‌ها (Exploding Gradients)**
پس از چند گام زمانی، گرادیان‌ها بسیار بزرگ می‌شود.
- **عدم تسخیر وابستگی‌های طولانی-مدت**

فرمول بندی جایگزین برای شبکه‌های عصبی بازگشتی

ALTERNATIVES FORMULATION FOR RNNs

یک فرمول بندی جایگزین:

$$c_t = W \cdot \tanh(c_{t-1}) + U \cdot x_t + b$$

$$\mathcal{L} = \sum_t \mathcal{L}_t(c_t)$$

فرمول بندی جایگزین برای شبکه‌های عصبی بازگشتی

نگاه دیگری به گرادیان‌ها

ANOTHER LOOK AT THE GRADIENTS

$$\mathcal{L} = L(c_T(c_{T-1}(\dots(c_1(x_1, c_0; W); W); W); W))$$

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{\tau=1}^t \frac{\partial \mathcal{L}_t}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} \frac{\partial c_\tau}{\partial W}$$

$$\frac{\partial \mathcal{L}}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} = \underbrace{\frac{\partial \mathcal{L}}{\partial c_t} \cdot \frac{\partial c_t}{\partial c_{t-1}} \cdot \frac{\partial c_{t-1}}{\partial c_{t-2}} \cdot \dots \cdot \frac{\partial c_{\tau+1}}{\partial c_\tau}}_{\text{Rest} \rightarrow \text{short-term factors}} \leq \eta^{t-\tau} \frac{\partial \mathcal{L}_t}{\partial c_t}$$

$t \gg \tau \rightarrow$ long-term factors

η determines the norm of the gradients

گرادیان‌های RNN، یک حاصل ضرب بازگشتی از $\partial c_t / \partial c_{t-1}$ است.

فرمول بندی جایگزین برای شبکه‌های عصبی بازگشتی

گرادیان‌های RNN در یک بعد

RNN GRADIENTS IN 1D

$$\begin{aligned}
 \circ \quad \frac{\partial \mathcal{L}}{\partial c_t} &= \frac{\partial \mathcal{L}}{\partial c_T} \cdot \frac{\partial c_T}{\partial c_{T-1}} \cdot \frac{\partial c_{T-1}}{\partial c_{T-2}} \cdot \dots \cdot \frac{\partial c_{t+1}}{\partial c_c} \left. \begin{array}{l} < 1 & < 1 & < 1 \end{array} \right\} \frac{\partial \mathcal{L}}{\partial W} \ll 1 \Rightarrow \text{Vanishing gradient} \\
 \circ \quad \frac{\partial \mathcal{L}}{\partial c_t} &= \frac{\partial \mathcal{L}}{\partial c_T} \cdot \frac{\partial c_T}{\partial c_{T-1}} \cdot \frac{\partial c_{T-1}}{\partial c_{T-2}} \cdot \dots \cdot \frac{\partial c_1}{\partial c_c} \left. \begin{array}{l} > 1 & > 1 & > 1 \end{array} \right\} \frac{\partial \mathcal{L}}{\partial W} \gg 1 \Rightarrow \text{Exploding gradient}
 \end{aligned}$$

فرمول بندی جایگزین برای شبکه های عصبی بازگشتی

گرادیان های RNN در چند بعد

RNN GRADIENTS IN N-D

When $c_T \in \mathbb{R}^N$ then $\frac{\partial c_t}{\partial c_{t-1}}$ is a Jacobian

$$\begin{aligned} \circ \frac{\partial \mathcal{L}}{\partial c_t} &= \frac{\partial \mathcal{L}}{\partial c_T} \cdot \frac{\partial c_T}{\partial c_{T-1}} \cdot \frac{\partial c_{T-1}}{\partial c_{T-2}} \cdot \dots \cdot \frac{\partial c_{t+1}}{\partial c_t} \left. \vphantom{\frac{\partial \mathcal{L}}{\partial c_t}} \right\} \frac{\partial \mathcal{L}}{\partial \theta} \ll 1 \Rightarrow \text{Vanishing gradient} \\ &\qquad \qquad \qquad < 1 \qquad < 1 \qquad < 1 \\ \circ \frac{\partial \mathcal{L}}{\partial c_t} &= \frac{\partial \mathcal{L}}{\partial c_T} \cdot \frac{\partial c_T}{\partial c_{T-1}} \cdot \frac{\partial c_{T-1}}{\partial c_{T-2}} \cdot \dots \cdot \frac{\partial c_{t+1}}{\partial c_t} \left. \vphantom{\frac{\partial \mathcal{L}}{\partial c_t}} \right\} \frac{\partial \mathcal{L}}{\partial \theta} \gg 1 \Rightarrow \text{Exploding gradient} \\ &\qquad \qquad \qquad > 1 \qquad > 1 \qquad > 1 \end{aligned}$$

$$y \in \mathbb{R}^2, x \in \mathbb{R}^3: \frac{dy}{dx} = \begin{bmatrix} \frac{\partial y^{(1)}}{\partial x^{(1)}} & \frac{\partial y^{(1)}}{\partial x^{(2)}} & \frac{\partial y^{(1)}}{\partial x^{(3)}} \\ \frac{\partial y^{(2)}}{\partial x^{(1)}} & \frac{\partial y^{(2)}}{\partial x^{(2)}} & \frac{\partial y^{(2)}}{\partial x^{(3)}} \end{bmatrix}$$

فرمول بندی جایگزین برای شبکه‌های عصبی بازگشتی

گرادیان‌های RNN در چند بعد

RNN GRADIENTS IN N-D

When $c_T \in \mathbb{R}^N$ then $\frac{\partial c_t}{\partial c_{t-1}}$ is a Jacobian

شعاع طیفی ژاکوبی (= بزرگ‌ترین مقدار ویژه ρ) پارامتر مهمی است:

$$\begin{aligned} \circ \quad \frac{\partial \mathcal{L}}{\partial c_t} &= \frac{\partial \mathcal{L}}{\partial c_T} \cdot \frac{\partial c_T}{\partial c_{T-1}} \cdot \frac{\partial c_{T-1}}{\partial c_{T-2}} \cdot \dots \cdot \frac{\partial c_{t+1}}{\partial c_t} \left. \vphantom{\frac{\partial \mathcal{L}}{\partial c_t}} \right\} \frac{\partial \mathcal{L}}{\partial c_t} \ll 1 \Rightarrow \text{Vanishing gradient} \\ &\quad \rho < 1 \quad \rho < 1 \quad \rho < 1 \\ \circ \quad \frac{\partial \mathcal{L}}{\partial c_t} &= \frac{\partial \mathcal{L}}{\partial c_T} \cdot \frac{\partial c_T}{\partial c_{T-1}} \cdot \frac{\partial c_{T-1}}{\partial c_{T-2}} \cdot \dots \cdot \frac{\partial c_{t+1}}{\partial c_t} \left. \vphantom{\frac{\partial \mathcal{L}}{\partial c_t}} \right\} \frac{\partial \mathcal{L}}{\partial c_t} \gg 1 \Rightarrow \text{Exploding gradient} \\ &\quad \rho > 1 \quad \rho > 1 \quad \rho > 1 \end{aligned}$$

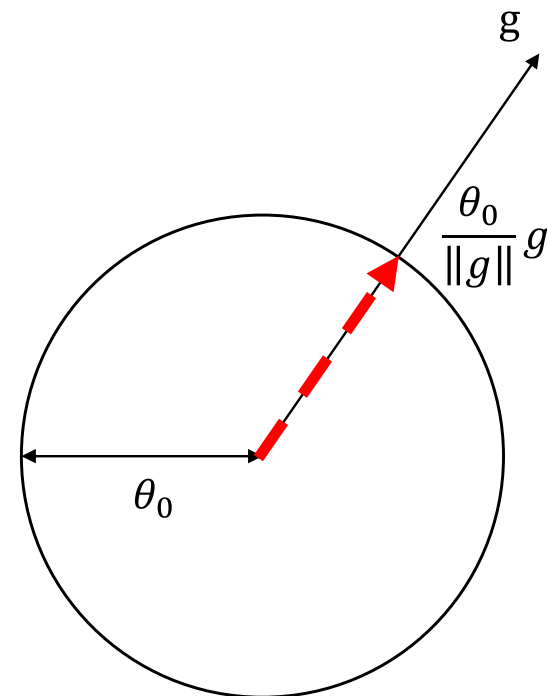
برش گرادیان برای جلوگیری از انفجار گرادیان

GRADIENT CLIPPING

مقیاس‌بندی گرادیان با یک مقدار آستانه:

Pseudocode

1. $g \leftarrow \frac{\partial \mathcal{L}}{\partial W}$
2. if $\|g\| > \theta_0$:
 $g \leftarrow \frac{\theta_0}{\|g\|} g$
 else:
 print('Do nothing')



این الگوریتم ساده است، اما به خوبی کار می‌کند!

اضمحلال گرادیان‌ها

VANISHING GRADIENTS

گرادیان خطا نسبت به سلول میانی:

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{\tau=1}^t \frac{\partial \mathcal{L}_r}{\partial y_t} \frac{\partial y_t}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} \frac{\partial c_\tau}{\partial W}$$

$$\frac{\partial c_t}{\partial c_\tau} = \prod_{t \geq k \geq \tau} \frac{\partial c_k}{\partial c_{k-1}} = \prod_{t \geq k \geq \tau} W \cdot \partial \tanh(c_{k-1})$$

اضمحلال گرادیان‌ها

VANISHING GRADIENTS

گرادیان خطا نسبت به سلول میانی:

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{\tau=1}^t \frac{\partial \mathcal{L}_r}{\partial y_t} \frac{\partial y_t}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} \frac{\partial c_\tau}{\partial W}$$

$$\frac{\partial c_t}{\partial c_\tau} = \prod_{t \geq k \geq \tau} \frac{\partial c_k}{\partial c_{k-1}} = \prod_{t \geq k \geq \tau} W \cdot \partial \tanh(c_{k-1})$$

- For $t = 1, r = 2 \Rightarrow \frac{\partial \mathcal{L}_2}{\partial W} \propto \frac{\partial c_2}{\partial c_1}$
- For $t = 1, r = 3 \Rightarrow \frac{\partial \mathcal{L}_3}{\partial W} \propto \frac{\partial c_3}{\partial c_1} = \frac{\partial c_3}{\partial c_2} \cdot \frac{\partial c_2}{\partial c_1}$
- For $t = 1, r = 4 \Rightarrow \frac{\partial \mathcal{L}_4}{\partial W} \propto \frac{\partial c_4}{\partial c_1} = \frac{\partial c_4}{\partial c_3} \cdot \frac{\partial c_3}{\partial c_2} \cdot \frac{\partial c_2}{\partial c_1}$

اضمحلال گرادیان‌ها

VANISHING GRADIENTS

گرادیان خطا نسبت به سلول میانی:

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{\tau=1}^t \frac{\partial \mathcal{L}_r}{\partial y_t} \frac{\partial y_t}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} \frac{\partial c_\tau}{\partial W}$$

$$\frac{\partial c_t}{\partial c_\tau} = \prod_{t \geq k \geq \tau} \frac{\partial c_k}{\partial c_{k-1}} = \prod_{t \geq k \geq \tau} W \cdot \partial \tanh(c_{k-1})$$

وابستگی‌های طولانی-مدت موجب می‌شوند وزن‌ها به صورت نمایی کوچک و کوچک‌تر شوند.

اضمحلال گرادیانها

VANISHING GRADIENTS

بازمقیاس دهی گرادیانهای مضمحل شده راه حل خوبی نیست!

وزنها بین گامهای زمانی مشترک هستند \Leftarrow اتلافها بر روی گامهای زمانی جمع می شوند:

$$\mathcal{L} = \sum_t \mathcal{L}_t \Rightarrow \frac{\partial \mathcal{L}}{\partial W} = \sum_t \frac{\partial \mathcal{L}_t}{\partial W}$$

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{\tau=1}^t \frac{\partial \mathcal{L}_t}{\partial c_\tau} \frac{\partial c_\tau}{\partial W} = \sum_{\tau=1}^t \frac{\partial \mathcal{L}_t}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} \frac{\partial c_\tau}{\partial W}$$

بازمقیاس دهی برای یک گام زمانی، بر همه‌ی گامهای زمانی تأثیر می گذارد

↓

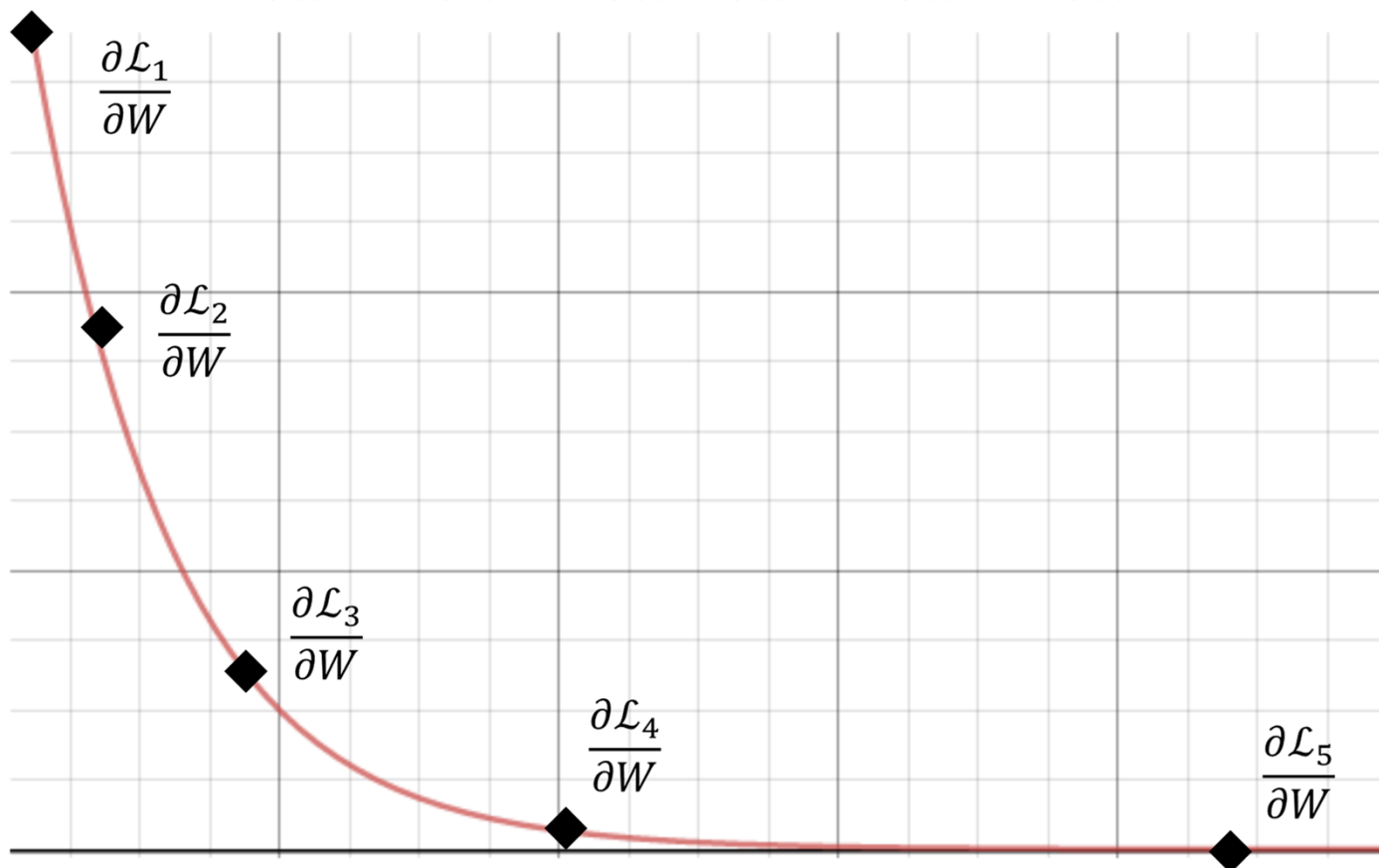
ضریب بازمقیاس دهی برای یک گام زمانی، برای گام زمانی دیگر کار نمی کند!

اضمحلال گرادیانها

مثال (۱ از ۲)

VANISHING GRADIENTS

$$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}_1}{\partial W} + \frac{\partial \mathcal{L}_2}{\partial W} + \frac{\partial \mathcal{L}_3}{\partial W} + \frac{\partial \mathcal{L}_4}{\partial W} + \frac{\partial \mathcal{L}_5}{\partial W}$$



اضمحلال گرادیانها

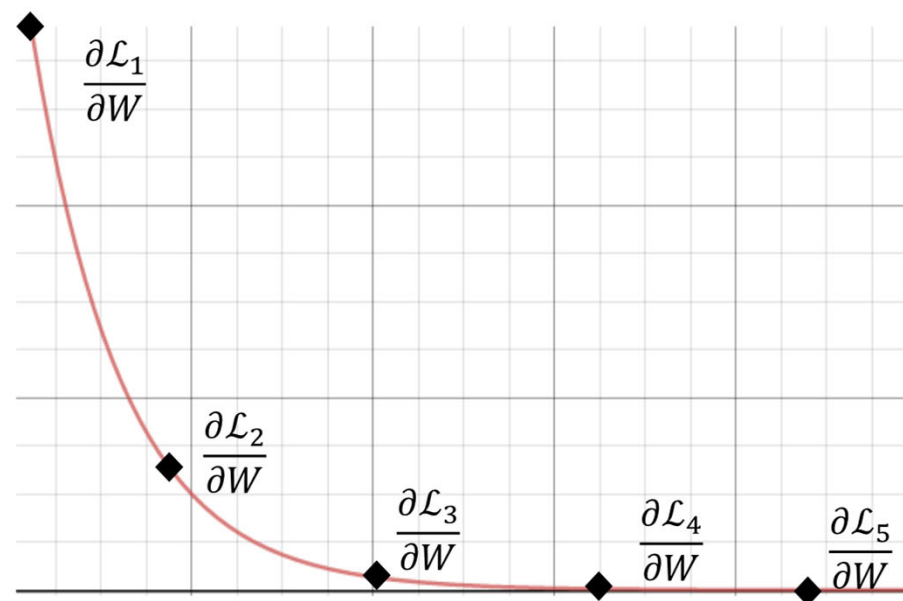
مثال (۲ از ۲)

VANISHING GRADIENTS

- Let's say $\frac{\partial \mathcal{L}_1}{\partial W} \propto 1, \frac{\partial \mathcal{L}_2}{\partial W} \propto 1/10, \frac{\partial \mathcal{L}_3}{\partial W} \propto 1/100, \frac{\partial \mathcal{L}_4}{\partial W} \propto 1/1000, \frac{\partial \mathcal{L}_5}{\partial W} \propto 1/10000$
- $\frac{\partial \mathcal{L}}{\partial W} = \sum_r \frac{\partial \mathcal{L}_r}{\partial W} = 1.1111$
- If $\frac{\partial \mathcal{L}}{\partial W}$ rescaled to 1 $\rightarrow \frac{\partial \mathcal{L}_5}{\partial W} \propto 10^{-5}$

وابستگی های طولانی-مدت قابل چشم پوشی است:
(یادگیری فقط بر روی کوتاه-مدت تمرکز می کند.)

$$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}_1}{\partial W} + \frac{\partial \mathcal{L}_2}{\partial W} + \frac{\partial \mathcal{L}_3}{\partial W} + \frac{\partial \mathcal{L}_4}{\partial W} + \frac{\partial \mathcal{L}_5}{\partial W}$$



اضمحلال گرادیانها

رفع مشکل

FIXING VANISHING GRADIENTS

* رگولاریزاسیون بر روی وزنهای بازگشتی (سیگنال خطا را مجبور می کند که مضمحل نشود).

$$\Omega = \sum_t \Omega_t = \sum_t \left(\frac{\left| \frac{\partial \mathcal{L}}{\partial c_{t+1}} \frac{\partial c_{t+1}}{\partial c_t} \right|}{\left| \frac{\partial \mathcal{L}}{\partial c_{t+1}} \right|} - 1 \right)^2$$

* استفاده از ماژولهای بازگشتی پیشرفته، مانند:

- ماژول حافظه‌ی کوتاه مدت طولانی (Long Short-Term Memory Module)
- ماژول واحد بازگشتی دروازه‌گذاری شده (Gated Recurrent Unit Module)

اضمحلال گرادیان‌ها

رفع مشکل

FIXING VANISHING GRADIENTS

سیگنال خطا در طول زمان، باید دارای نُرم خیلی بزرگ یا خیلی کوچک نباشد.

راه‌حل: استفاده از یک تابع فعال‌سازی که مشتق آن مساوی با 1 باشد.



گرادیان‌ها نه خیلی بزرگ می‌شوند و نه خیلی کوچک

شبکه‌های عصبی بازگشتی

۴

شبکه‌های
عصبی
بازگشتی
عمیق

شبکه‌های عصبی بازگشتی عمیق

شبکه‌های عصبی بازگشتی چندلایه / حافظه‌ی کوتاه-مدت طولانی

DEEP RNNs

Multilayer RNNs

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

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

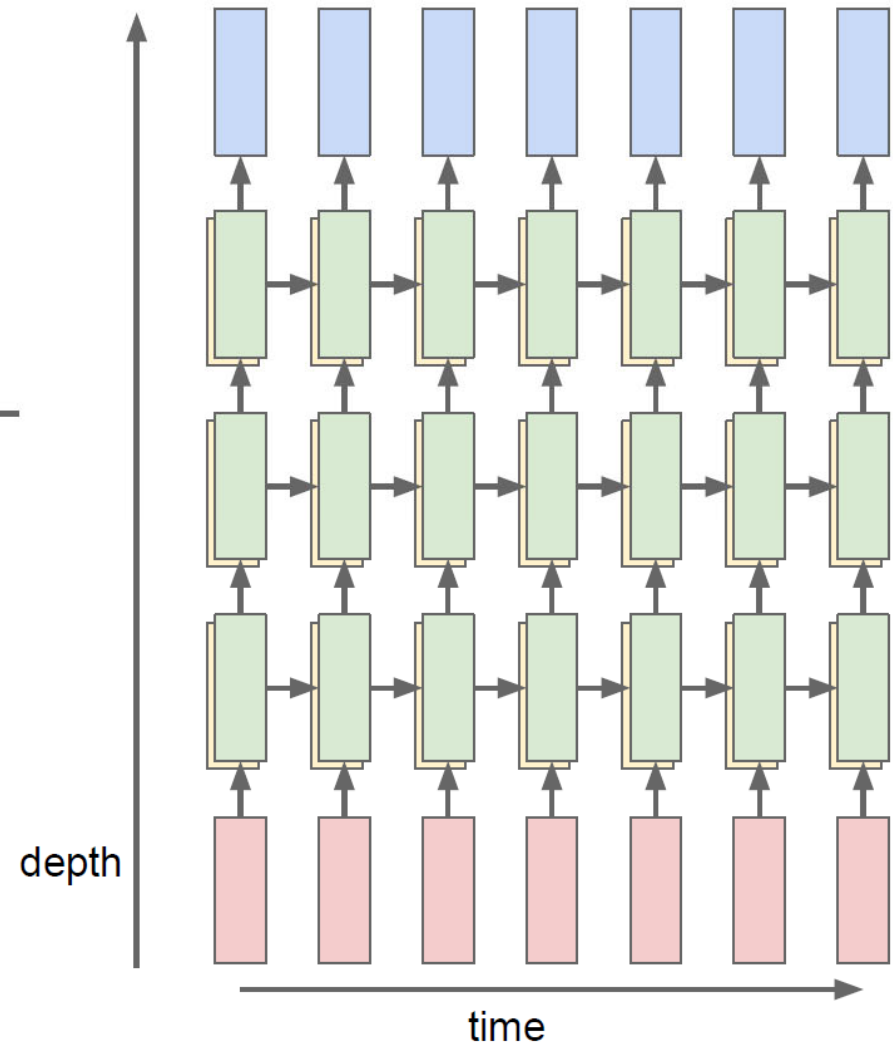
LSTM:

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

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

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

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



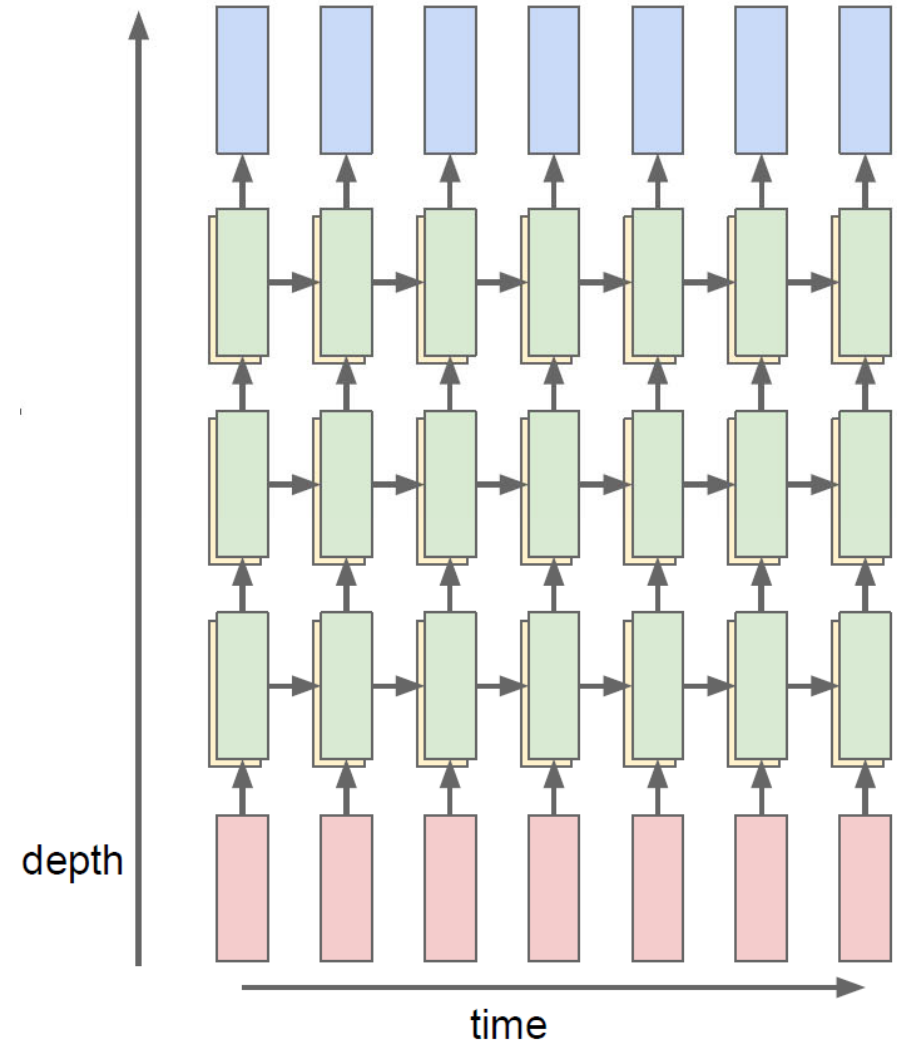
شبکه‌های عصبی بازگشتی

شبکه‌های عصبی بازگشتی چندلایه

MULTILAYER RNNs

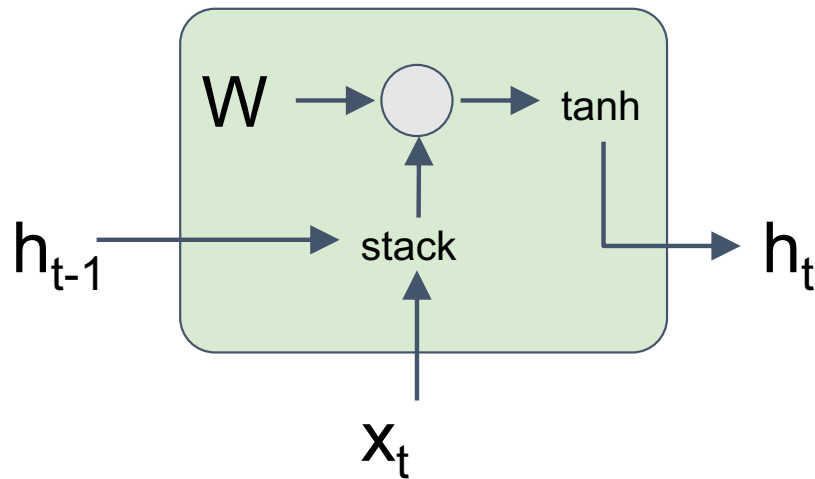
Multilayer RNNs

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

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


شبکه‌های عصبی بازگشتی ساده

جریان گرادیان

VANILLA RNN GRADIENT FLOW

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

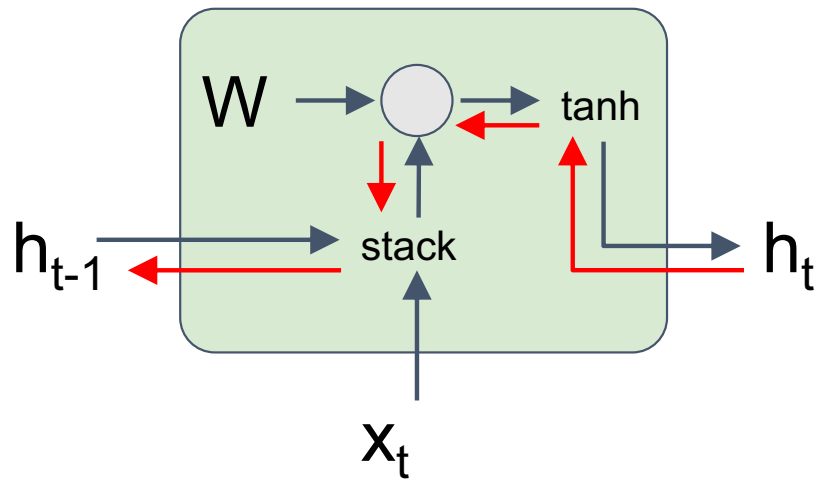
Bengio et al, "Learning long-term dependencies with gradient descent is difficult",
 IEEE Transactions on Neural Networks, 1994
 Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

شبکه‌های عصبی بازگشتی ساده

جریان گرادیان

VANILLA RNN GRADIENT FLOW

Backpropagation from h_t
to h_{t-1} multiplies by W
(actually W_{hh}^T)

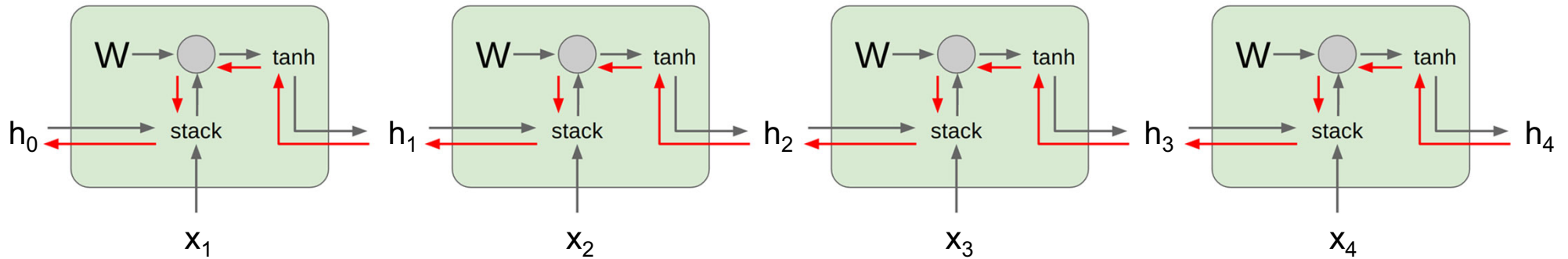


$$\begin{aligned}
 h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\
 &= \tanh\left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \\
 &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)
 \end{aligned}$$

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

شبکه‌های عصبی بازگشتی ساده

جریان گرادیان

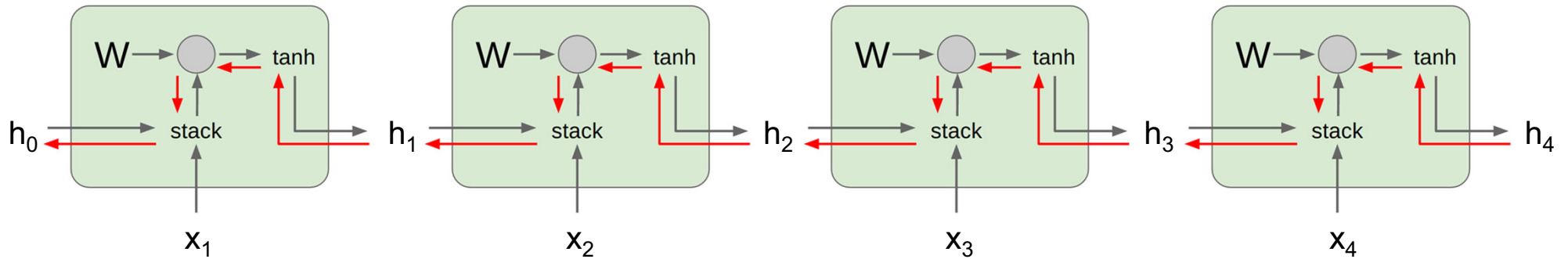
VANILLA RNN GRADIENT FLOW

Computing gradient of h_0 involves many factors of W (and repeated tanh)

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

شبکه‌های عصبی بازگشتی ساده

جریان گرادیان

VANILLA RNN GRADIENT FLOW

Computing gradient of h_0 involves many factors of W (and repeated tanh)

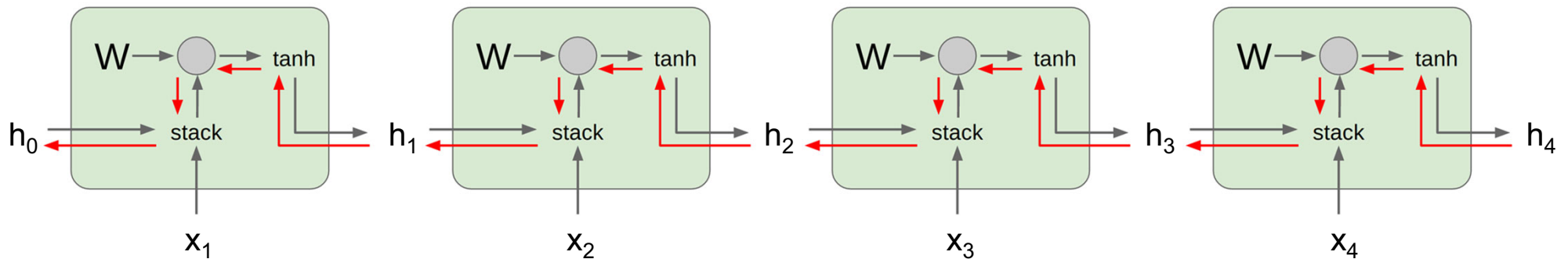
Largest singular value > 1 :
Exploding gradients

Largest singular value < 1 :
Vanishing gradients

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

شبکه‌های عصبی بازگشتی ساده

جریان گرادیان

VANILLA RNN GRADIENT FLOW

Computing gradient of h_0 involves many factors of W (and repeated tanh)

Largest singular value > 1 :
Exploding gradients

Largest singular value < 1 :
Vanishing gradients

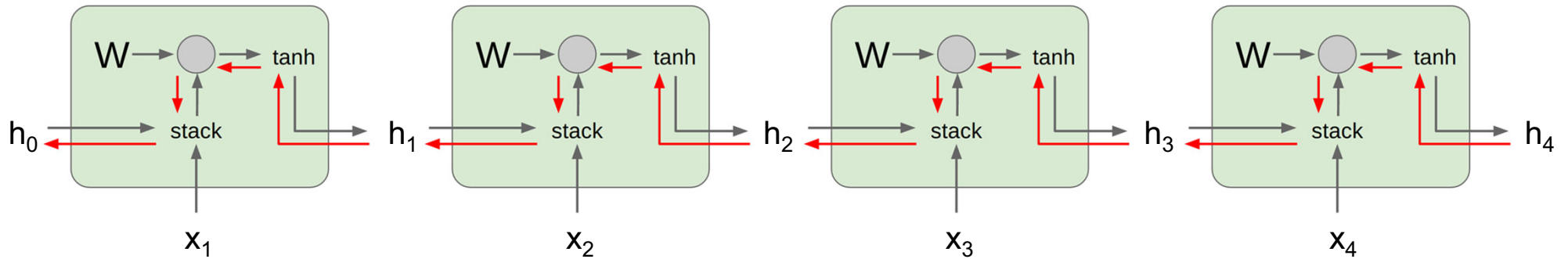
Gradient clipping: Scale gradient if its norm is too big

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

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

شبکه‌های عصبی بازگشتی ساده

جریان گرادیان

VANILLA RNN GRADIENT FLOW

Computing gradient of h_0 involves many factors of W (and repeated \tanh)

Largest singular value > 1 :
Exploding gradients

Largest singular value < 1 :
Vanishing gradients

→ Change RNN architecture

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

شبکه‌های عصبی بازگشتی

۵

حافظه‌ی
کوتاه-مدت
طولانی

حافظه‌ی کوتاه-مدت طولانی

مقایسه

LONG SHORT-TERM MEMORY (LSTM)**Vanilla RNN**

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

LSTM

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

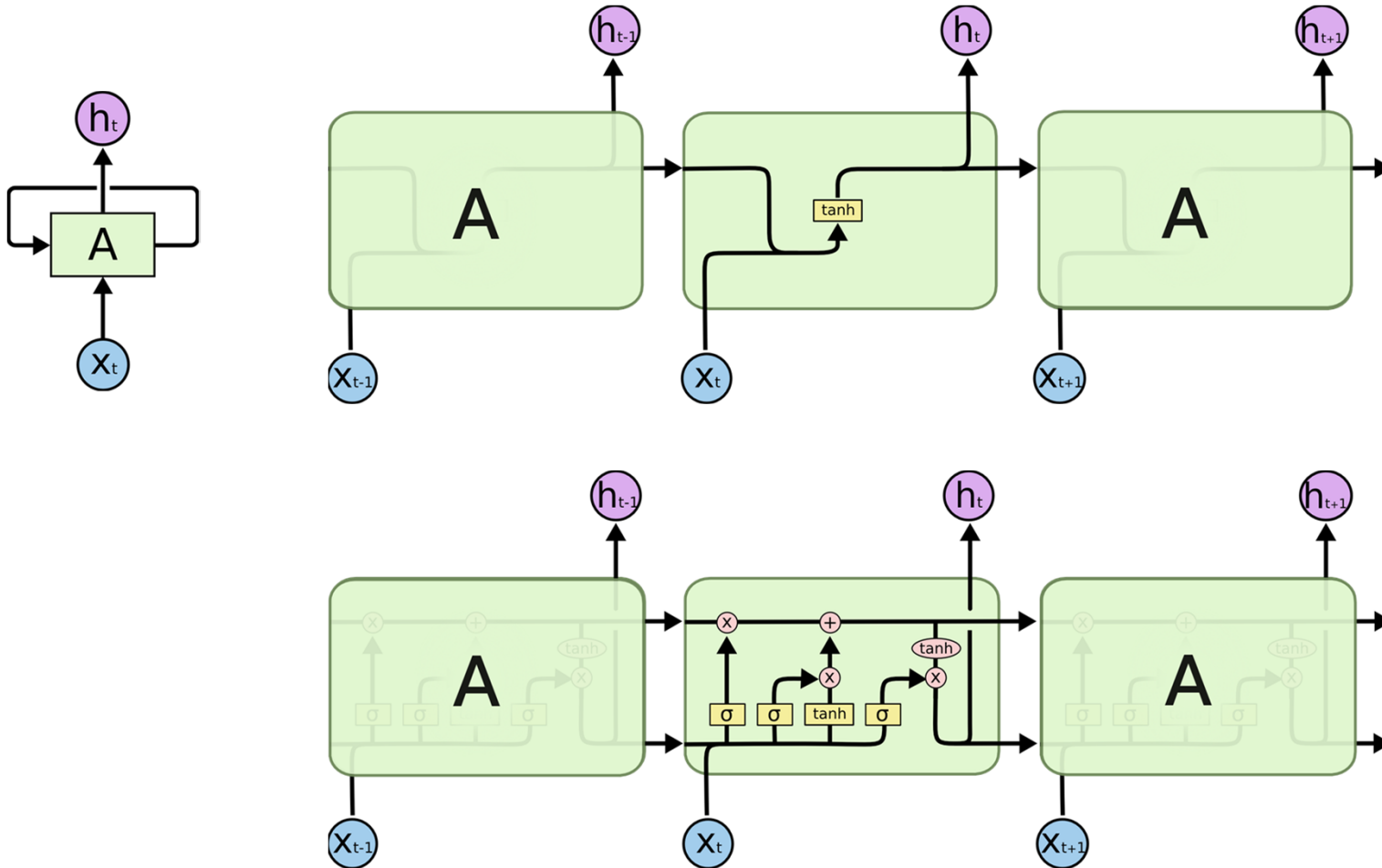
$$c_t = f \odot c_{t-1} + i \odot g$$

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

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

حافظه‌ی کوتاه-مدت طولانی

مقایسه

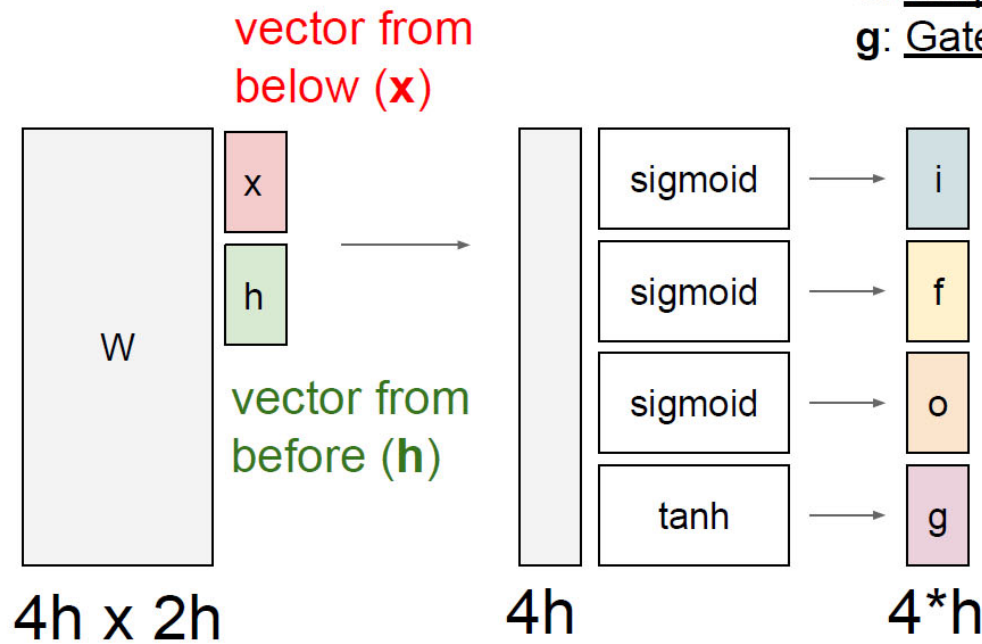
LONG SHORT-TERM MEMORY (LSTM)

حافظه‌ی کوتاه-مدت طولانی

Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]

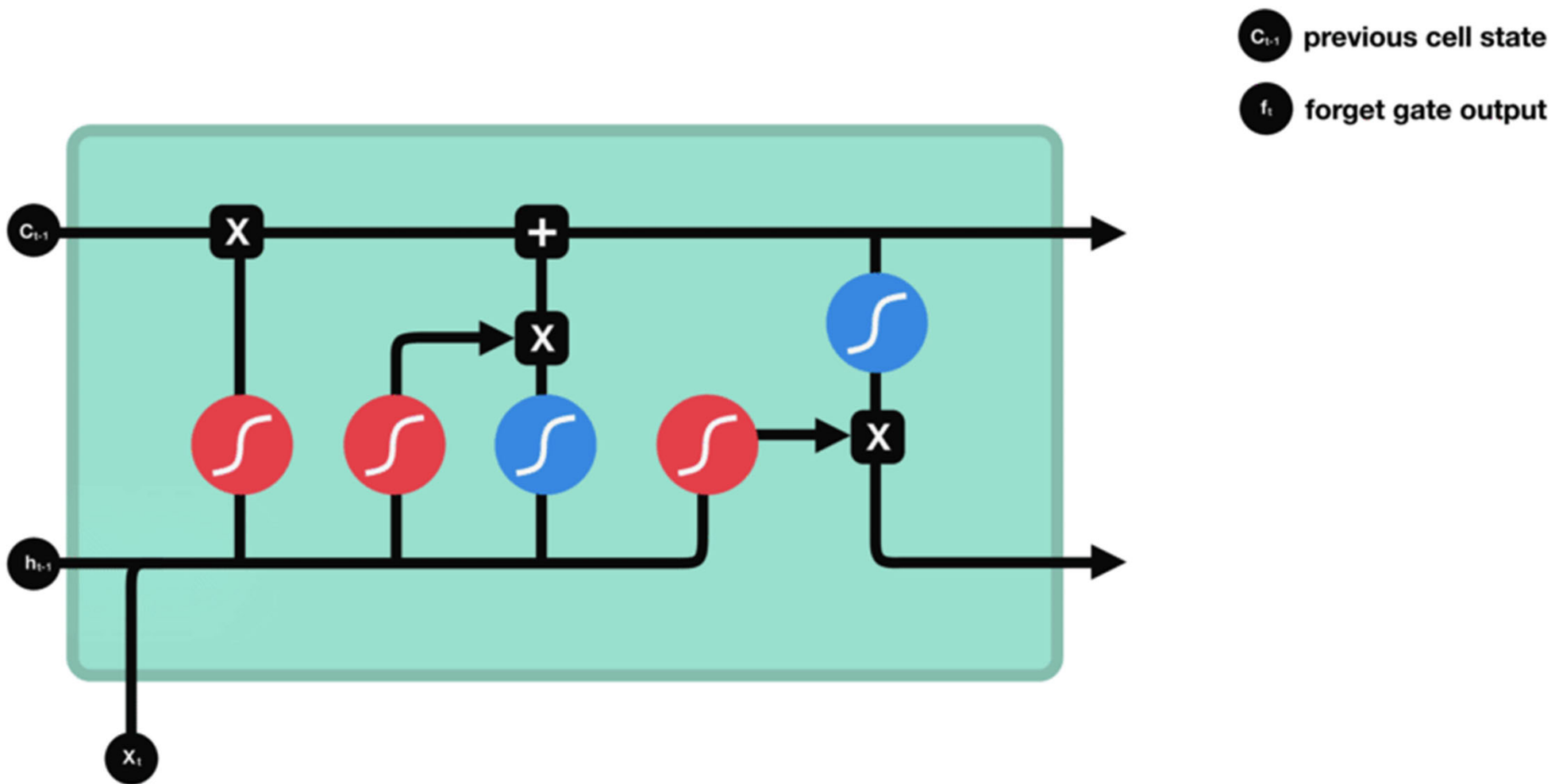
- i: Input gate, whether to write to cell
- f: Forget gate, Whether to erase cell
- o: Output gate, How much to reveal cell
- g: Gate gate (?), How much to write to cell

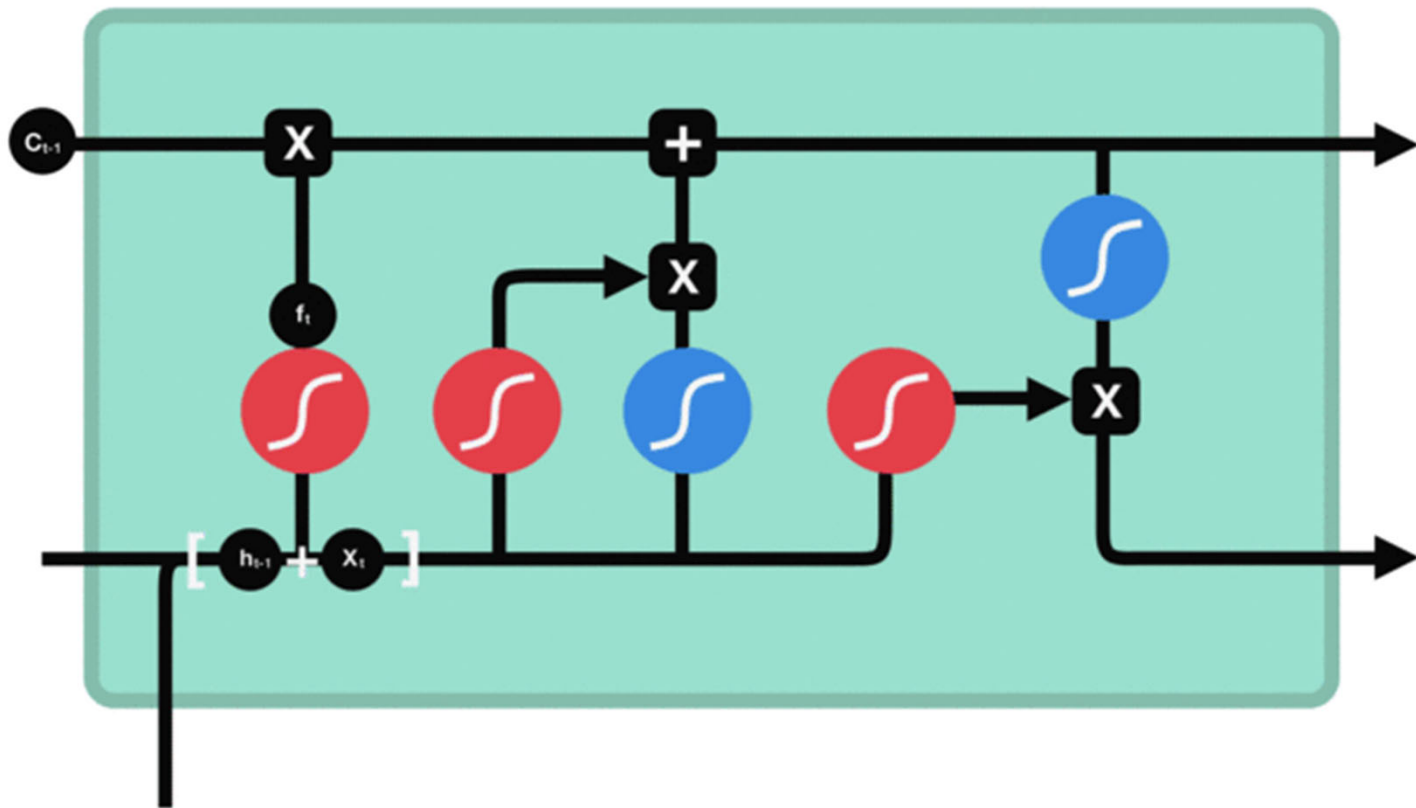


$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

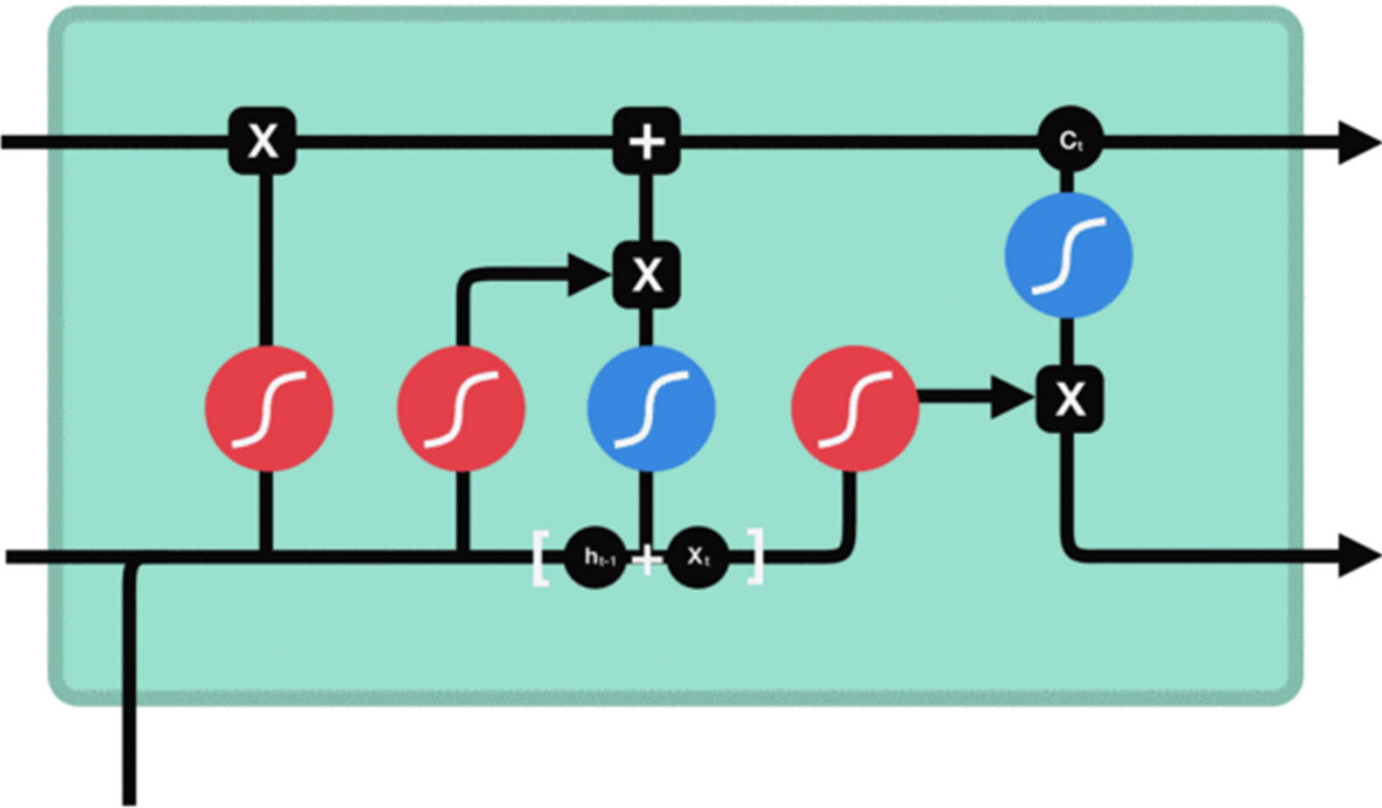
$$c_t = f \odot c_{t-1} + i \odot g$$

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





- c_{t-1} previous cell state
- f_t forget gate output
- i_t input gate output
- \hat{c}_t candidate
- (input gate == update gate)

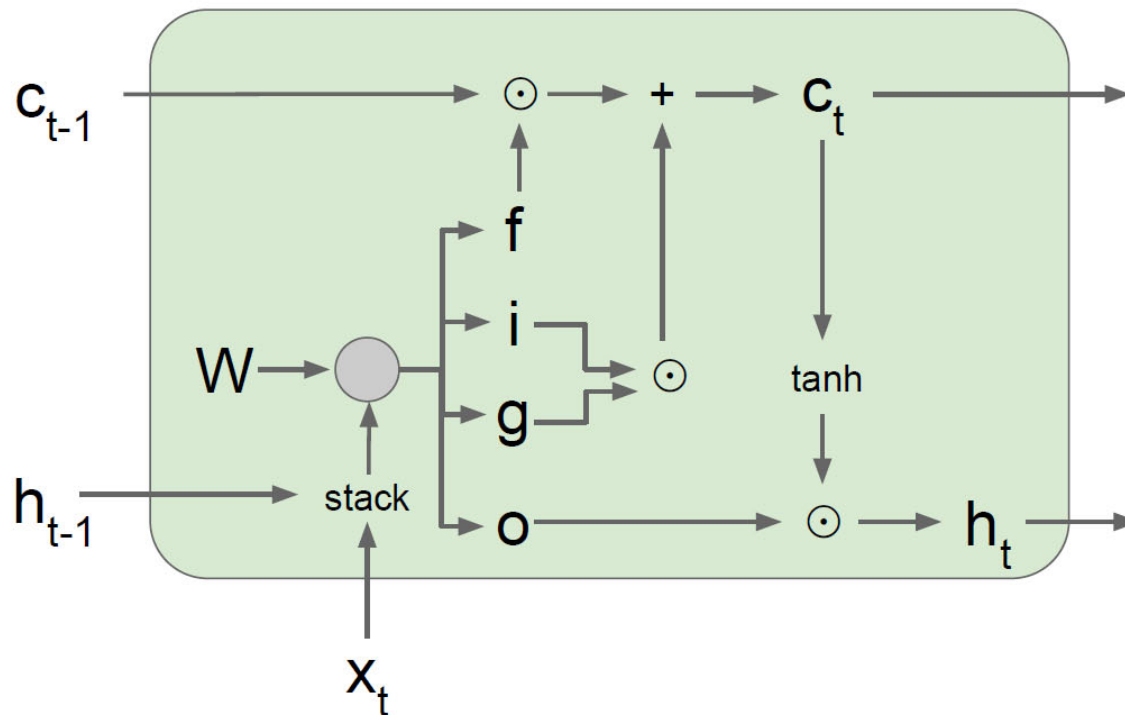


- c_{t-1} previous cell state
 - f_t forget gate output
 - i_t input gate output
 - \tilde{c}_t candidate
 - c_t new cell state
 - o_t output gate output
 - h_t hidden state
- (input gate == update gate)

حافظه‌ی کوتاه-مدت طولانی

Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



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

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

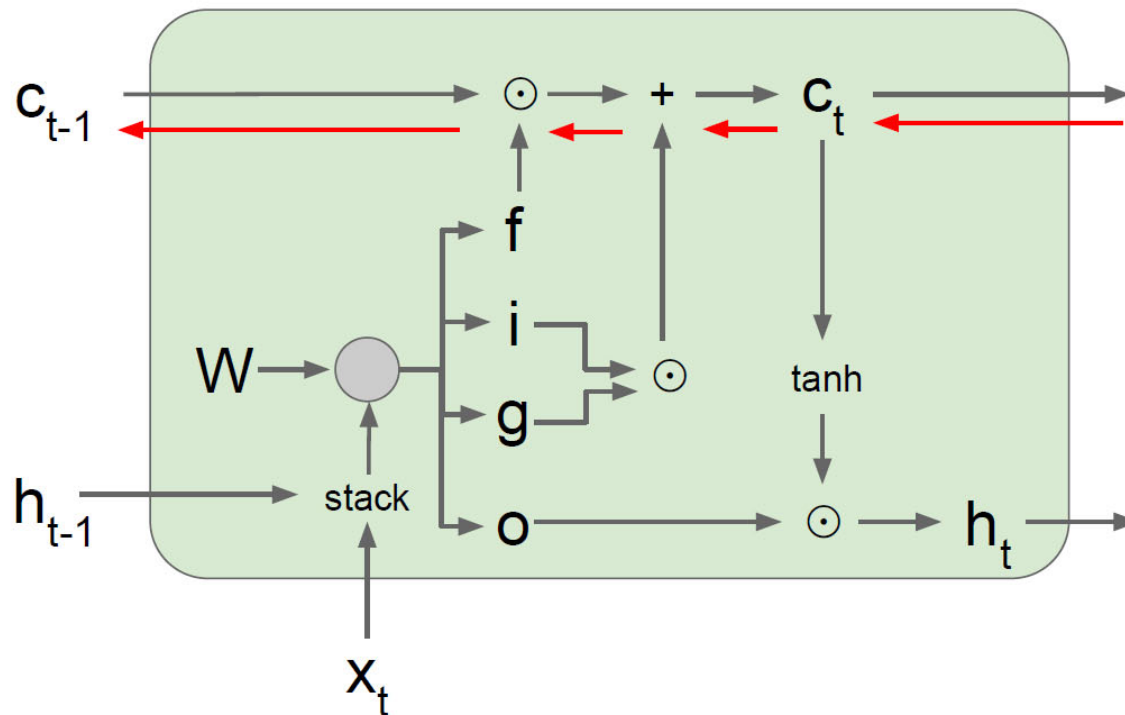
$$h_t = o \odot \tanh(c_t)$$

حافظه‌ی کوتاه-مدت طولانی

جریان گرادیان

Long Short Term Memory (LSTM): Gradient Flow

[Hochreiter et al., 1997]



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

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

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

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

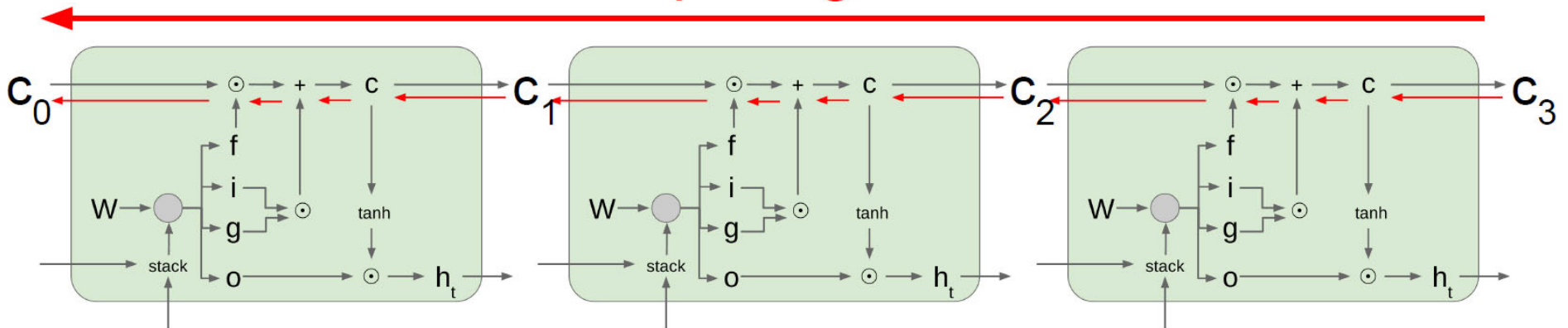
حافظه‌ی کوتاه-مدت طولانی

جریان گرادیان

Long Short Term Memory (LSTM): Gradient Flow

[Hochreiter et al., 1997]

Uninterrupted gradient flow!



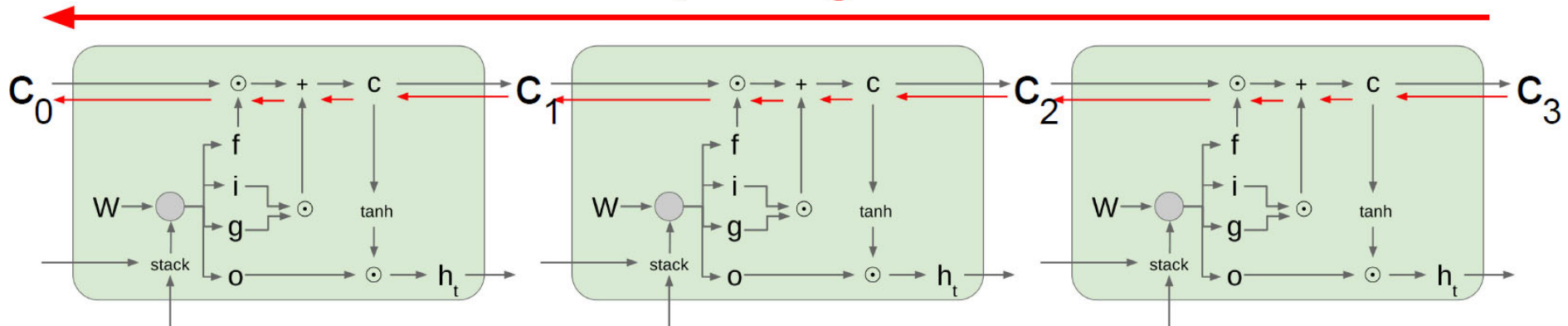
حافظه‌ی کوتاه-مدت طولانی

جریان گرادیان

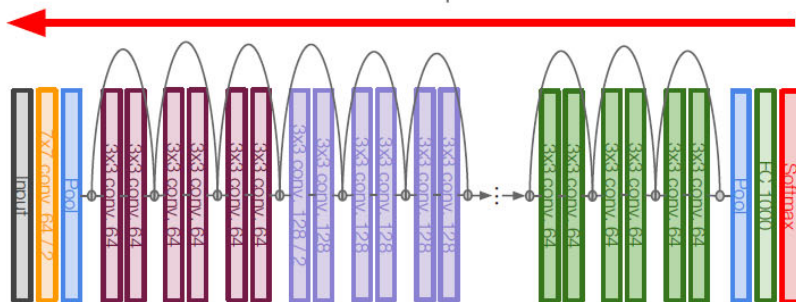
Long Short Term Memory (LSTM): Gradient Flow

[Hochreiter et al., 1997]

Uninterrupted gradient flow!



Similar to ResNet!



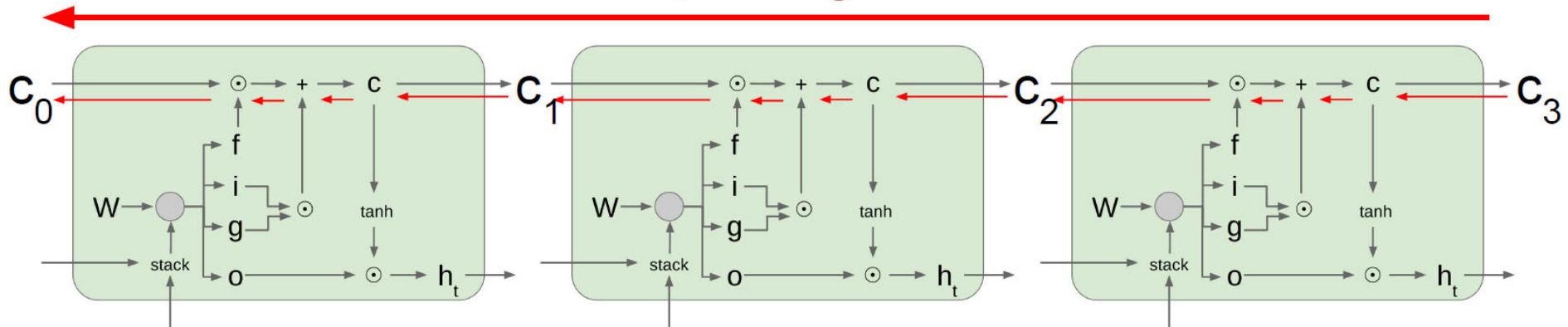
حافظه‌ی کوتاه-مدت طولانی

جریان گرادینان

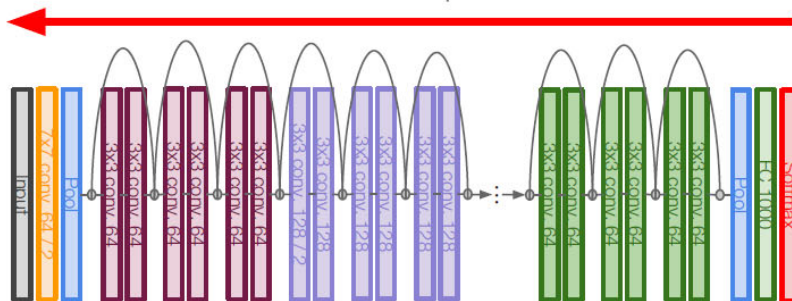
Long Short Term Memory (LSTM): Gradient Flow

[Hochreiter et al., 1997]

Uninterrupted gradient flow!



Similar to ResNet!



In between:
Highway Networks

$$g = T(x, W_T)$$

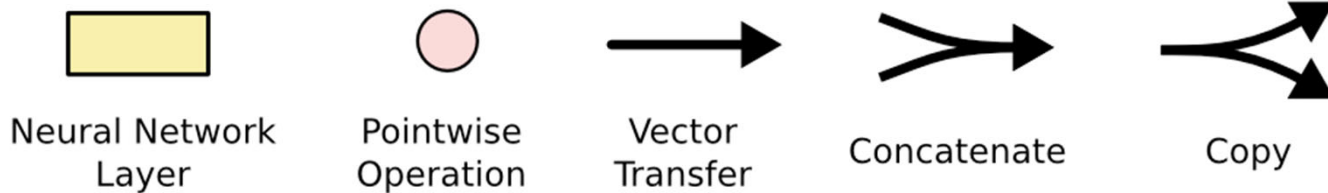
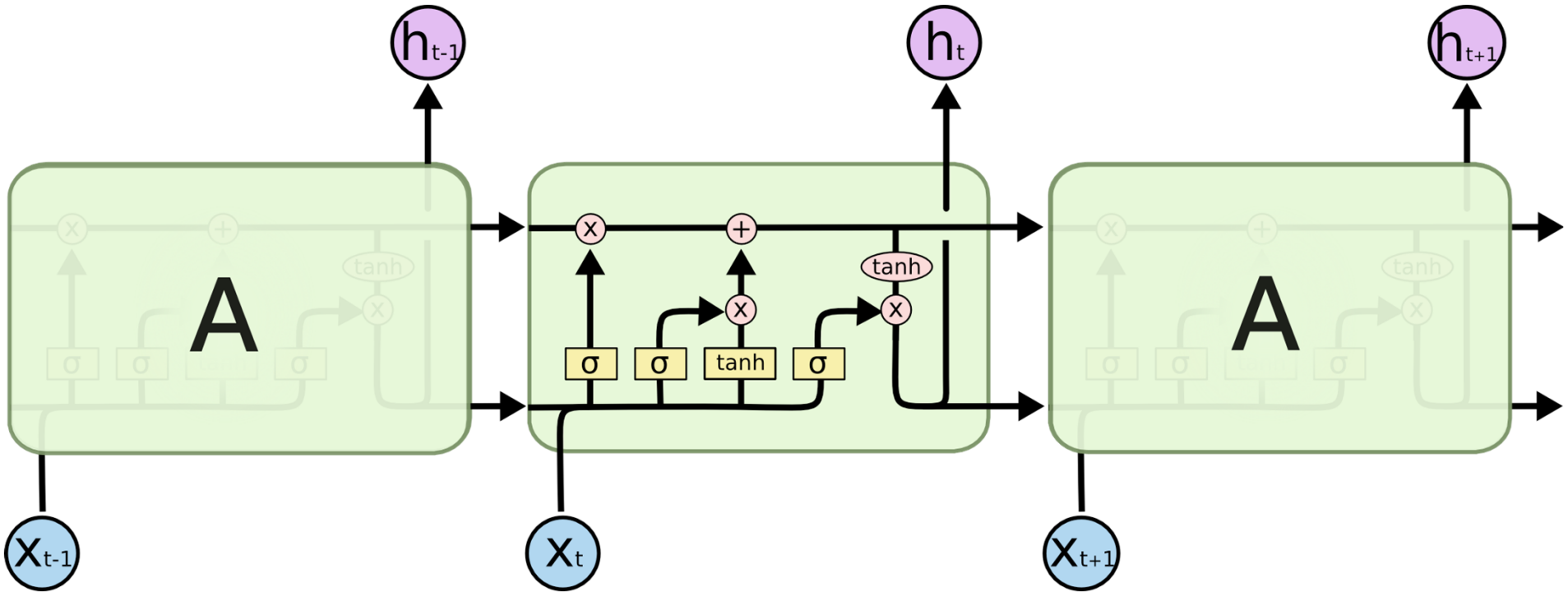
$$y = g \odot H(x, W_H) + (1 - g) \odot x$$

Srivastava et al, "Highway Networks",
ICML DL Workshop 2015

حافظه‌ی کوتاه-مدت طولانی

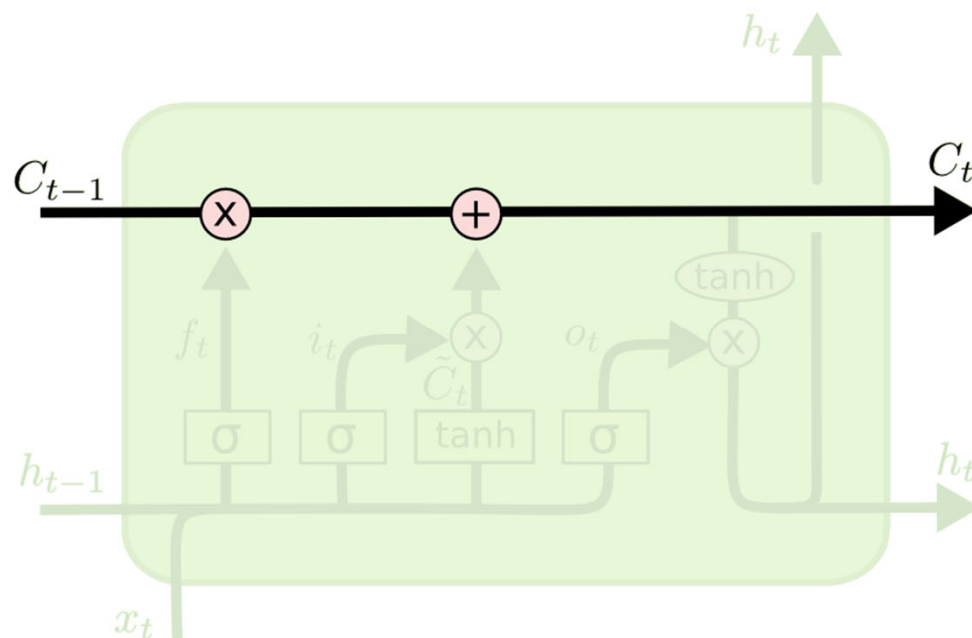
نشانه‌ها

MEET LSTMs



حافظه‌ی کوتاه-مدت طولانی

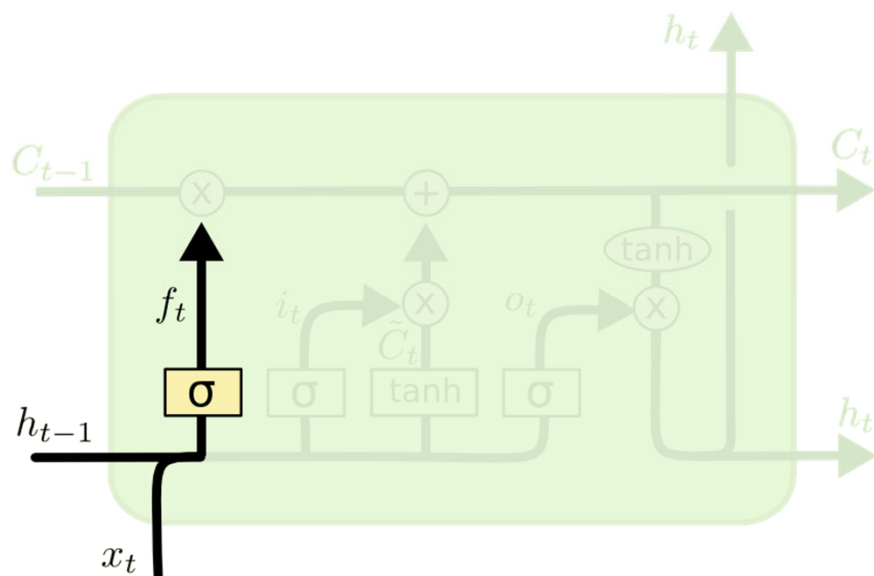
درک شهودی LSTM: حافظه

LSTMs INTUITION: MEMORY

حالت سلول / حافظه
Cell State / Memory

حافظه‌ی کوتاه-مدت طولانی

درک شهودی LSTM: گیت فراموشی

LSTMs INTUITION: FORGET GATE

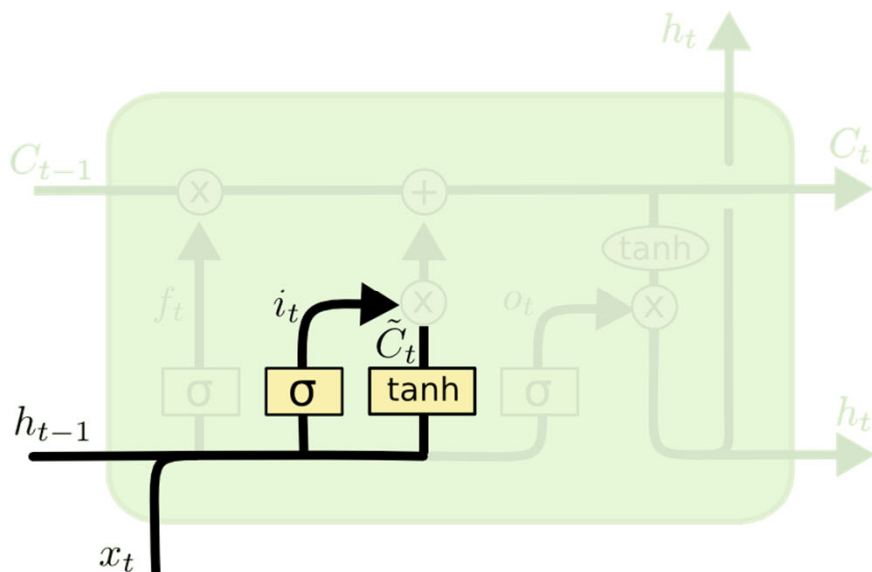
$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

آیا باید این بیت از اطلاعات را به یاد بیاوریم یا خیر؟
Should we continue to remember this “bit” of information or not?

مثال: با دیدن یک فاعل جدید در جمله، می‌خواهیم اطلاعات مربوط به فاعل قبلی (مانند جنسیت و تعداد) را فراموش کنیم.

حافظه‌ی کوتاه-مدت طولانی

درک شهودی LSTM: گیت ورودی

LSTMs INTUITION: INPUT GATE

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

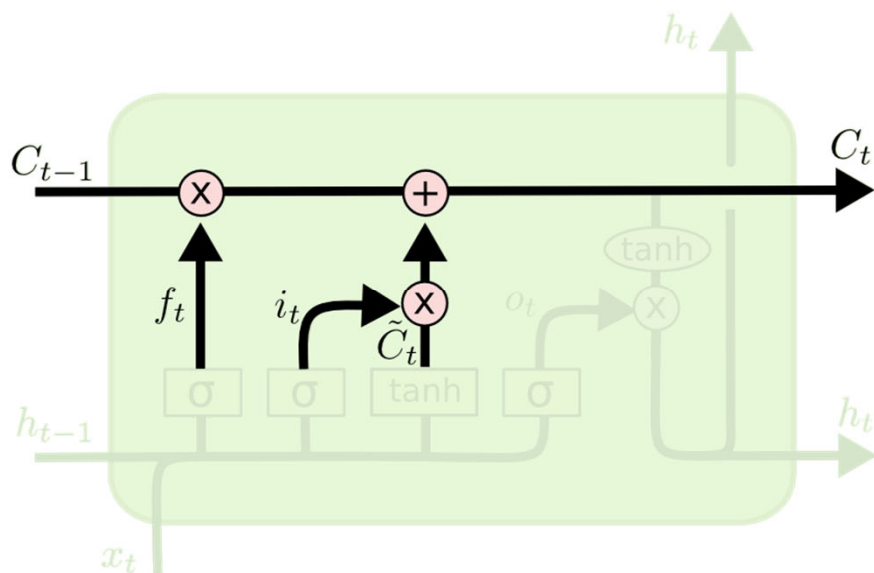
آیا باید این بیت از اطلاعات را به یاد بیاوریم یا خیر؟
اگر این طور است، با کدام ورودی؟

Should we continue to remember this “bit” of information or not?
If so, with what?

مثال: با دیدن یک فاعل جدید در جمله، می‌خواهیم اطلاعات مربوط به فاعل جدید (مانند جنسیت و تعداد) را به خاطر بسپاریم.

حافظه‌ی کوتاه-مدت طولانی

درک شهودی LSTM: به روزرسانی حافظه

LSTMs INTUITION: MEMORY UPDATE

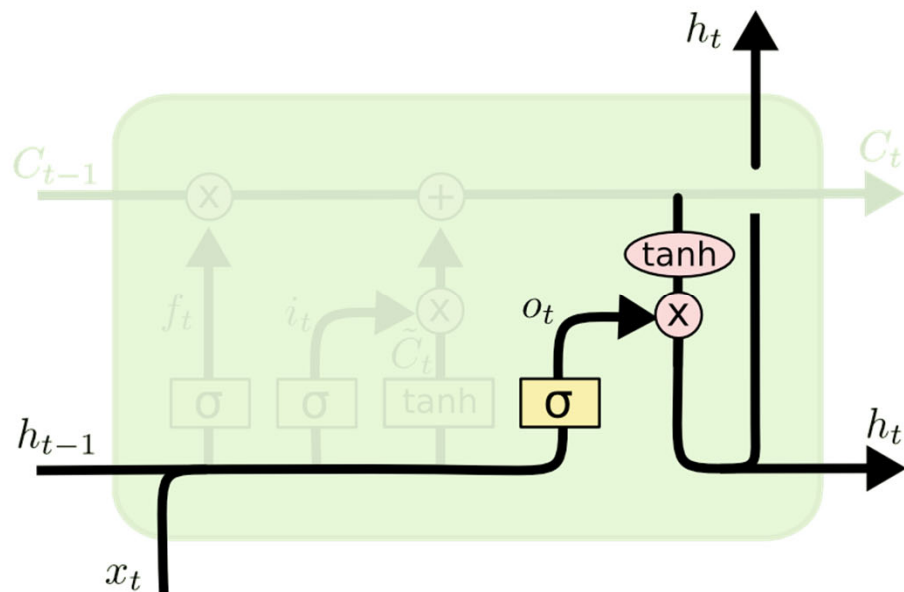
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

آن را فراموش کن + این را به خاطر بسپار
Forget that + Memorize this

مثال: حذف اطلاعات غیرضروری از حافظه و افزودن اطلاعات مفید جدید به آن

حافظه‌ی کوتاه-مدت طولانی

درک شهودی LSTM: گیت خروجی

LSTMs INTUITION: OUTPUT GATE

$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

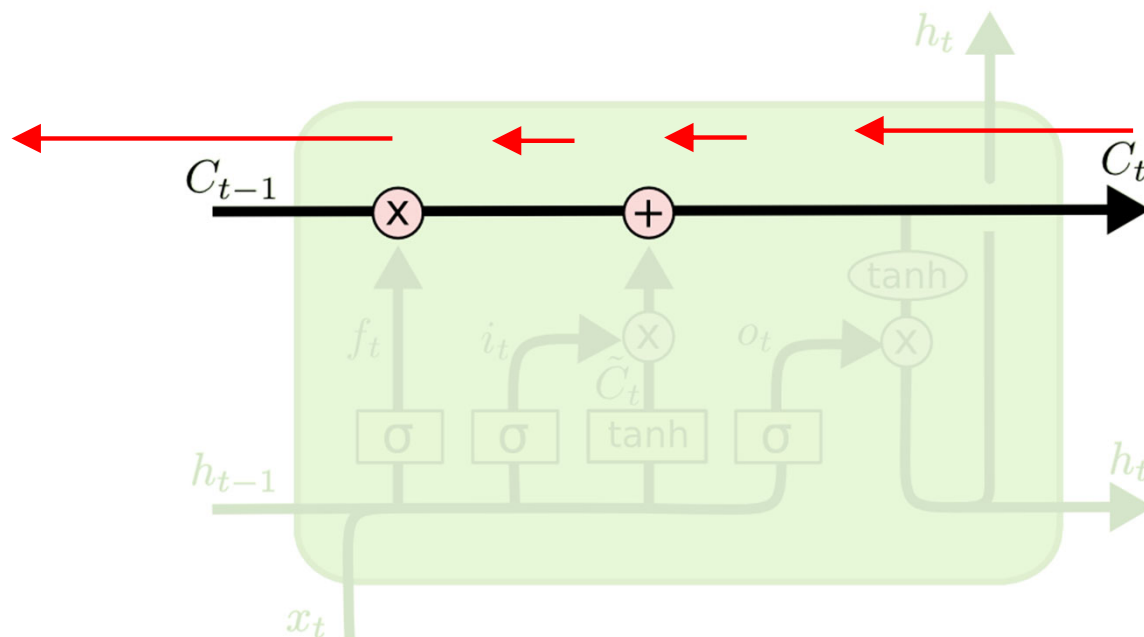
$$h_t = o_t * \tanh (C_t)$$

آیا باید این بیت از اطلاعات را به سمت لایه‌های عمیق‌تر شبکه خارج کنیم؟
Should we output this “bit” of information to “deeper” layers?

مثال: فرستادن اطلاعات فاعل جمله (مثل جنسیت و تعداد) به خروجی به‌منظور صرف درست فعل پس از آن

حافظه‌ی کوتاه-مدت طولانی

درک شهودی LSTM: به روزرسانی‌های جمعی

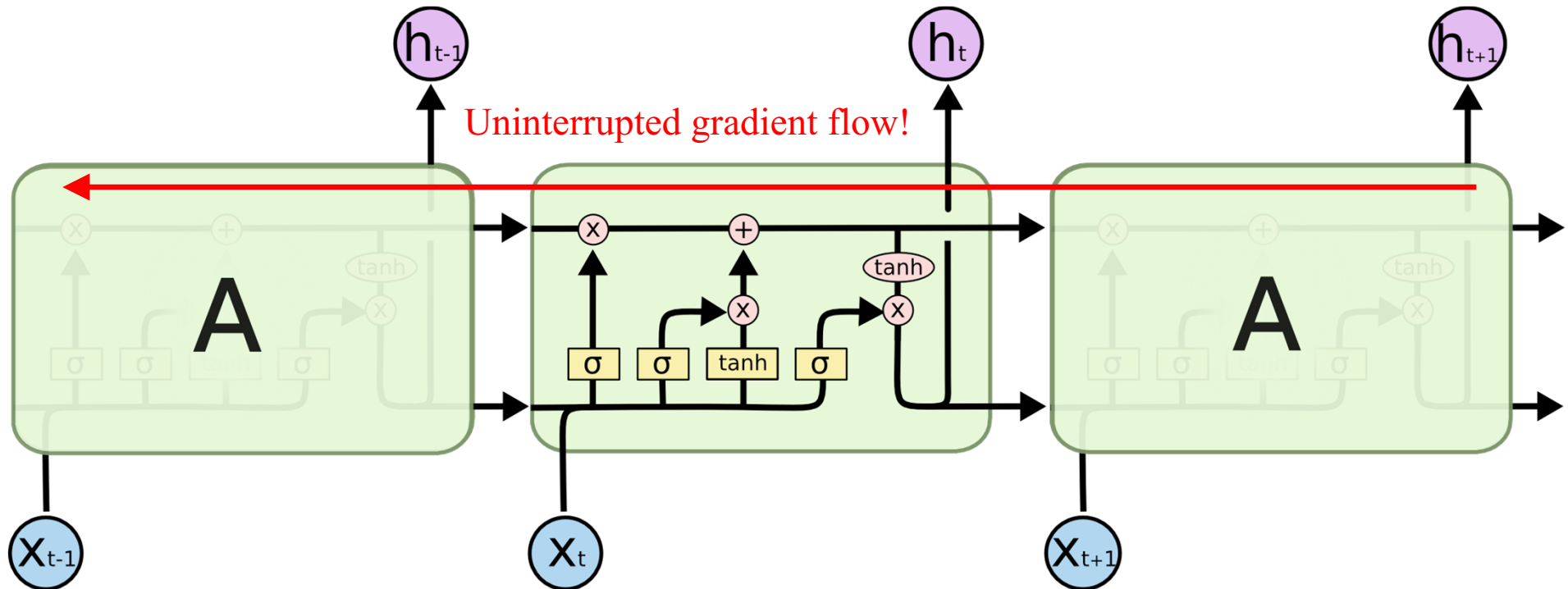
LSTMs INTUITION: ADDITIVE UPDATES

Backpropagation from C_t to C_{t-1} only elementwise multiplication by f , no matrix multiply by W

پس انتشار از C_t به C_{t-1} تنها با ضرب عنصر به عنصر در f انجام می‌شود، هیچ ضرب ماتریسی در W نداریم.

حافظه‌ی کوتاه-مدت طولانی

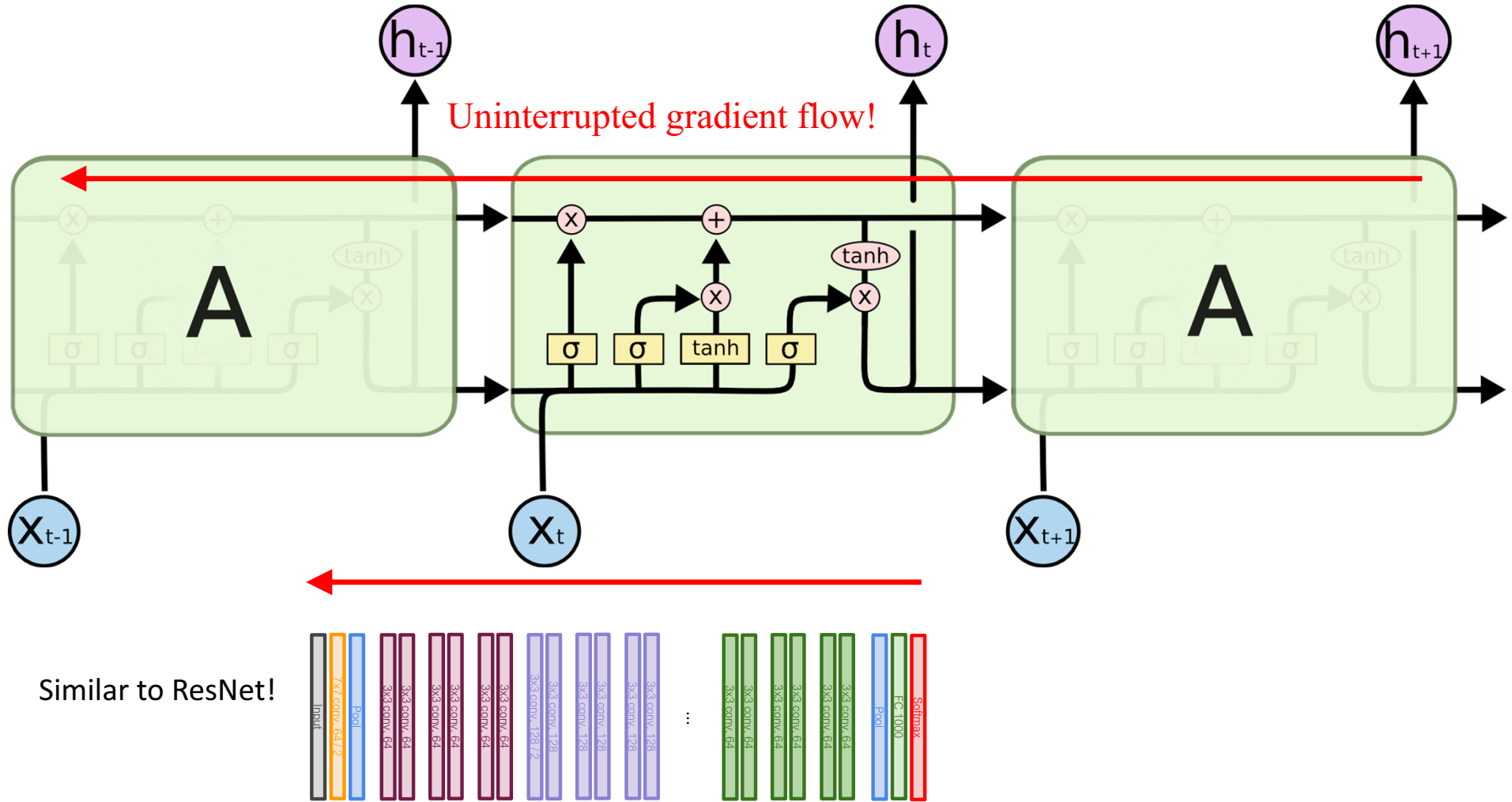
درک شهودی LSTM: به‌روزرسانی‌های جمعی

LSTMs INTUITION: ADDITIVE UPDATES

حافظه‌ی کوتاه-مدت طولانی

درک شهودی LSTM: به‌روزرسانی‌های جمعی

LSTMs INTUITION: ADDITIVE UPDATES



Similar to ResNet!

حافظه‌ی کوتاه-مدت طولانی

LONG SHORT-TERM MEMORY

$$i = \sigma(x_t U^{(i)} + m_{t-1} W^{(i)})$$

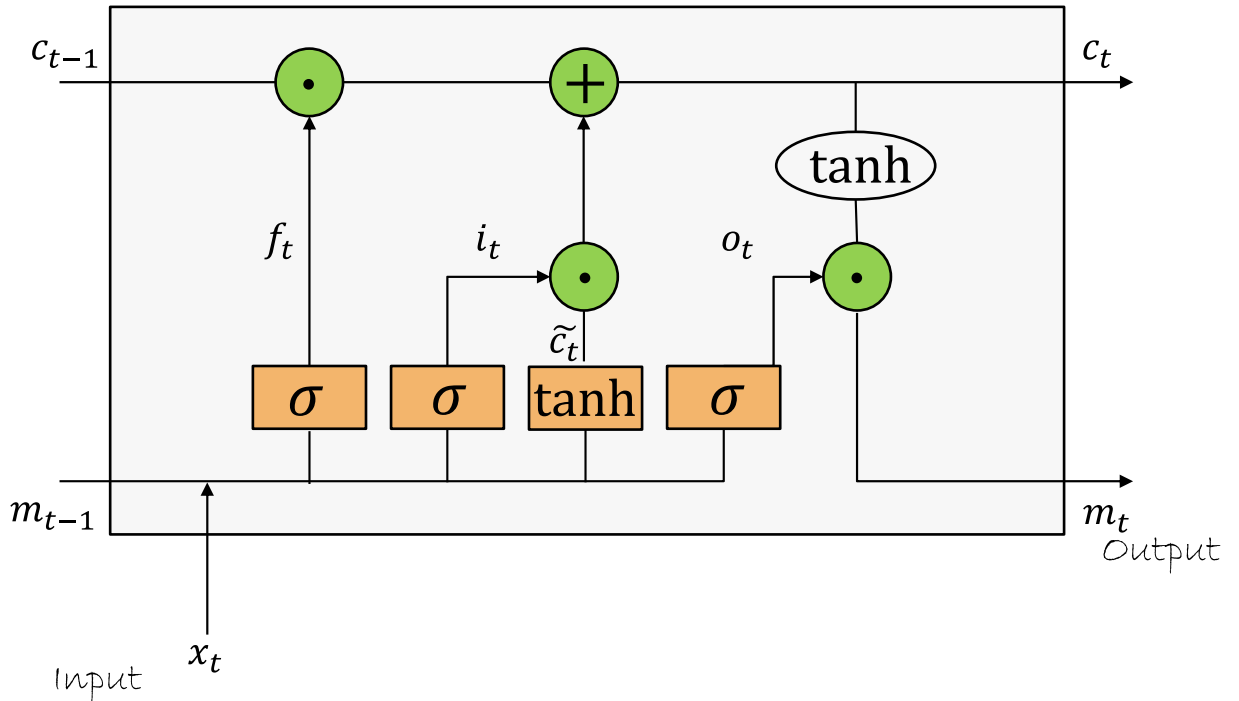
$$f = \sigma(x_t U^{(f)} + m_{t-1} W^{(f)})$$

$$o = \sigma(x_t U^{(o)} + m_{t-1} W^{(o)})$$

$$\tilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$

$$c_t = c_{t-1} \odot f + \tilde{c}_t \odot i$$

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



حافظه‌ی کوتاه-مدت طولانی

حالت سلول

CELL STATE

حالت سلول، اطلاعات اساسی را در طول زمان حمل می‌کند.

Cell state line

$$i = \sigma(x_t U^{(i)} + m_{t-1} W^{(i)})$$

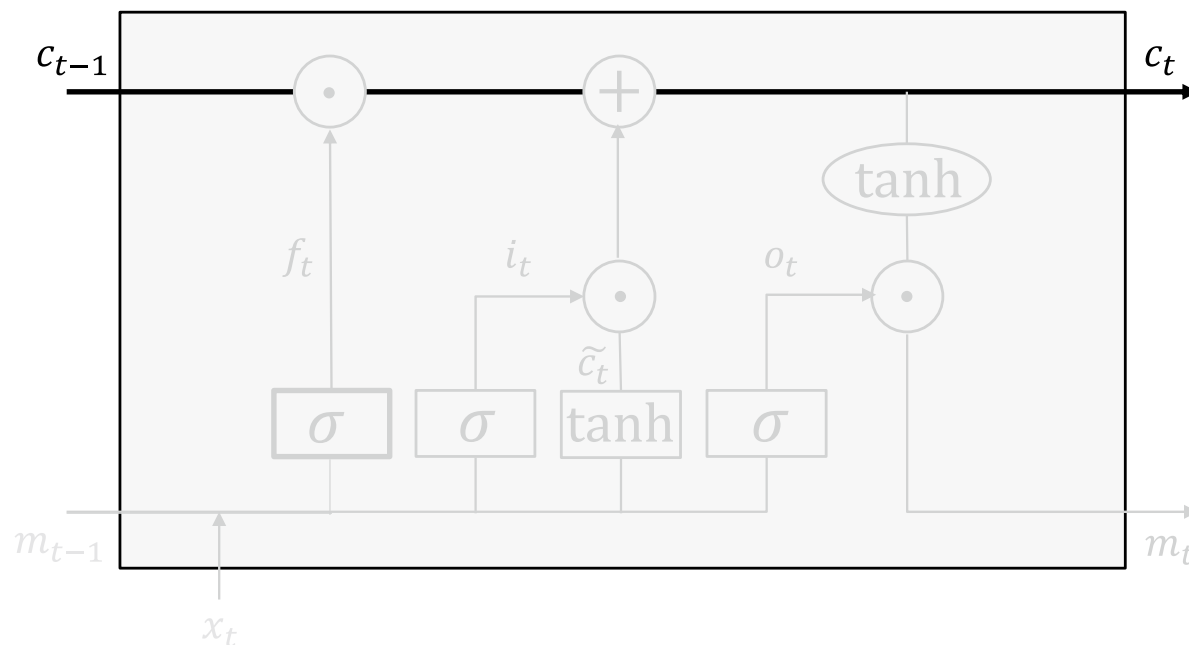
$$f = \sigma(x_t U^{(f)} + m_{t-1} W^{(f)})$$

$$o = \sigma(x_t U^{(o)} + m_{t-1} W^{(o)})$$

$$\tilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$

$$c_t = c_{t-1} \odot f + \tilde{c}_t \odot i$$

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



حافظه‌ی کوتاه-مدت طولانی

غیرخطی‌های LSTM

LSTM NON-LINEARITIES

$\sigma \in (0,1)$ گیت (دروازه) کنترل (control gate): شبیه یک سوئیچ عمل می‌کند.

$\tanh \in (-1,1)$ غیرخطی بازگشتی

$$i = \sigma(x_t U^{(i)} + m_{t-1} W^{(i)})$$

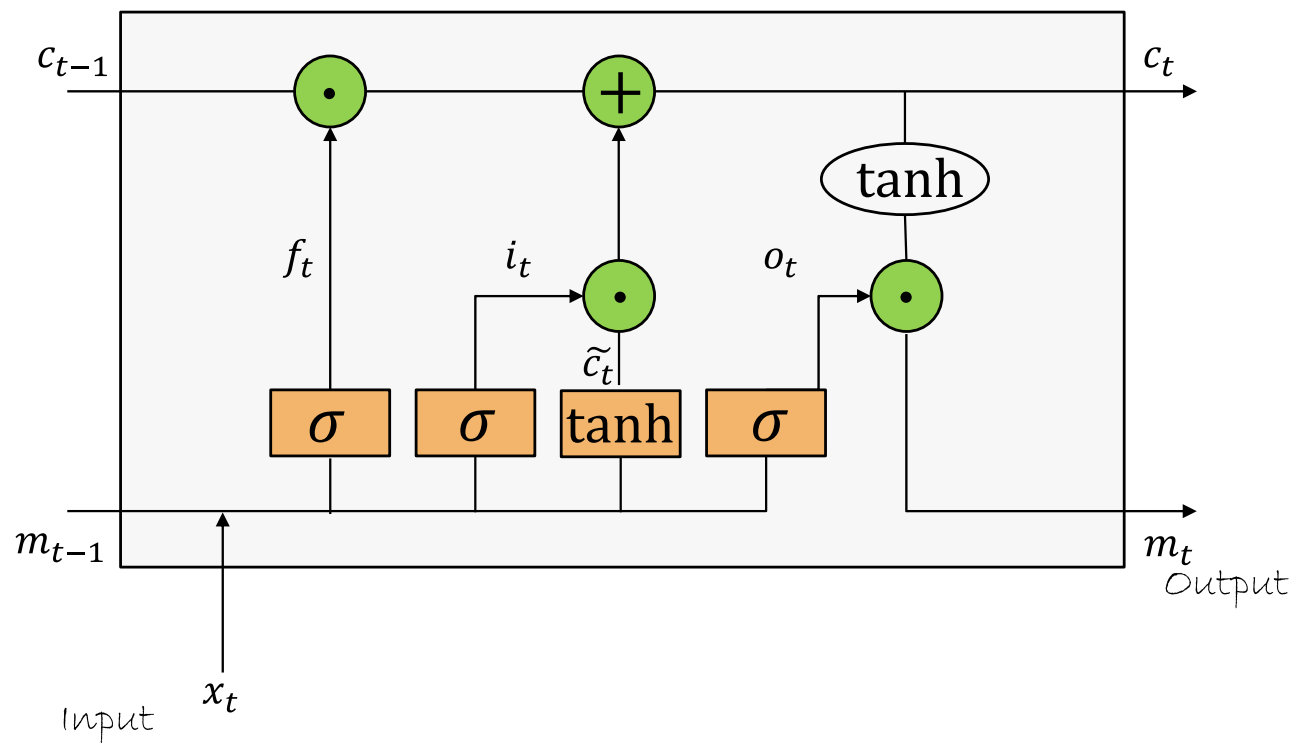
$$f = \sigma(x_t U^{(f)} + m_{t-1} W^{(f)})$$

$$o = \sigma(x_t U^{(o)} + m_{t-1} W^{(o)})$$

$$\tilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$

$$c_t = c_{t-1} \odot f + \tilde{c}_t \odot i$$

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



حافظه‌ی کوتاه-مدت طولانی

گام به گام با LSTM: گام ۱

LSTM STEP-BY-STEP: STEP 1

برای مثال: می‌خواهیم یک جمله را مدل کنیم.

باید تصمیم بگیریم برای حافظه‌ی جدید چه چیزی را فراموش کنیم و چه چیزی را به خاطر بیاوریم.

سیگموئید = 1 \Leftarrow به یادآوری همه چیز
 سیگموئید = 0 \Leftarrow فراموش کردن همه چیز

$$i_t = \sigma(x_t U^{(i)} + m_{t-1} W^{(i)})$$

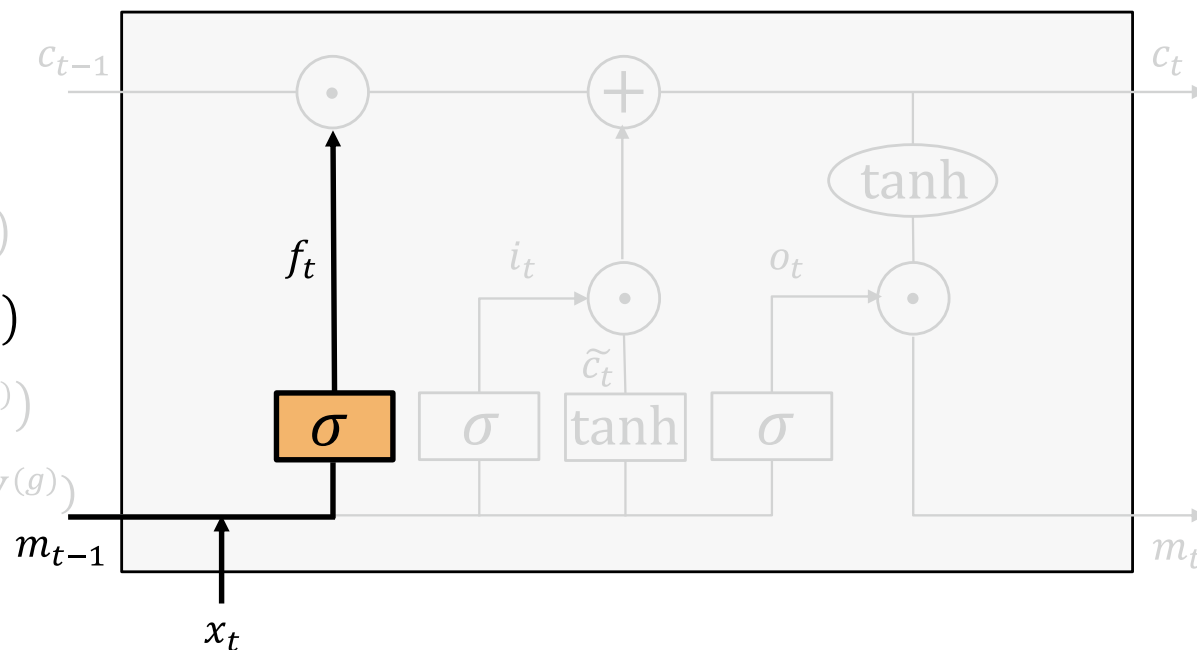
$$f_t = \sigma(x_t U^{(f)} + m_{t-1} W^{(f)})$$

$$o_t = \sigma(x_t U^{(o)} + m_{t-1} W^{(o)})$$

$$\tilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$

$$c_t = c_{t-1} \odot f + \tilde{c}_t \odot i$$

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



حافظه‌ی کوتاه-مدت طولانی

گام به گام با LSTM: گام ۲

LSTM STEP-BY-STEP: STEP 2

باید تصمیم بگیریم چه اطلاعات جدیدی باید به حافظه‌ی جدید اضافه شود.

- ورودی i_t را مدوله می‌کنیم.
- حافظه‌های کاندیدیای \tilde{c}_t را تولید می‌کنیم.

$$i_t = \sigma(x_t U^{(i)} + m_{t-1} W^{(i)})$$

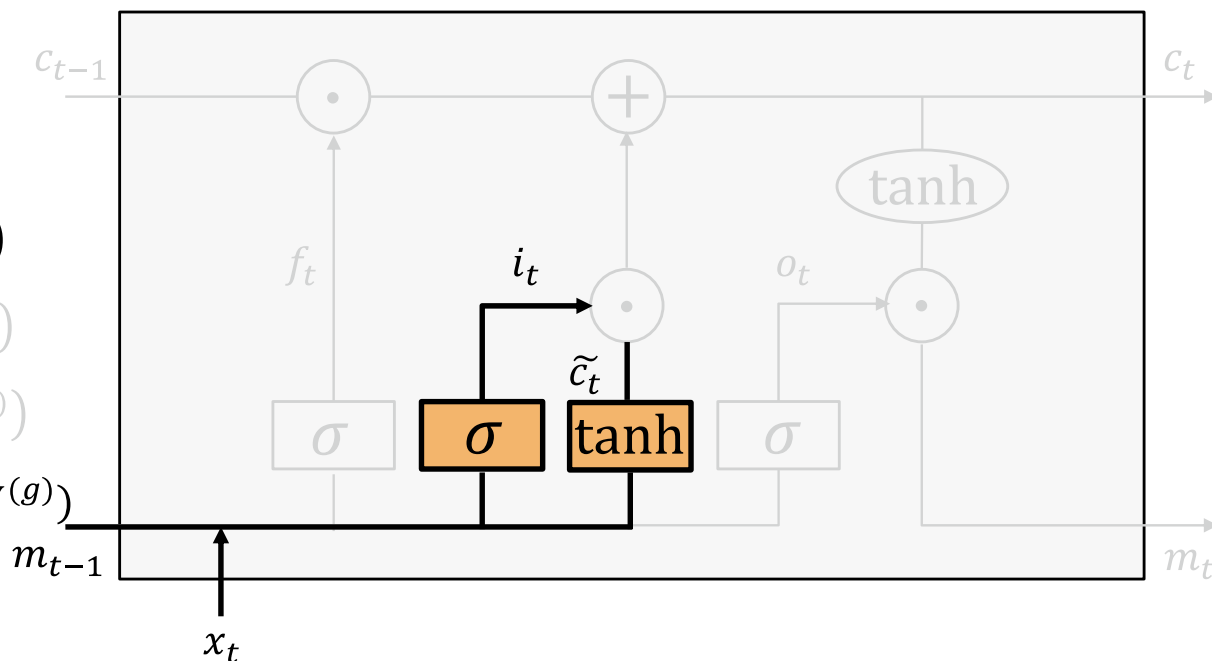
$$f_t = \sigma(x_t U^{(f)} + m_{t-1} W^{(f)})$$

$$o_t = \sigma(x_t U^{(o)} + m_{t-1} W^{(o)})$$

$$\tilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$

$$c_t = c_{t-1} \odot f + \tilde{c}_t \odot i$$

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



حافظه‌ی کوتاه-مدت طولانی

گام به گام با LSTM: گام ۳

LSTM STEP-BY-STEP: STEP 3حالت فعلی سلول c_t را محاسبه و به‌هنگام می‌کنیم، بر اساس:

- حالت قبلی سلول
- آنچه تصمیم داریم فراموش کنیم
- آنچه به‌عنوان ورودی مجاز می‌دانیم
- حافظه‌های کاندیدا

$$i_t = \sigma(x_t U^{(i)} + m_{t-1} W^{(i)})$$

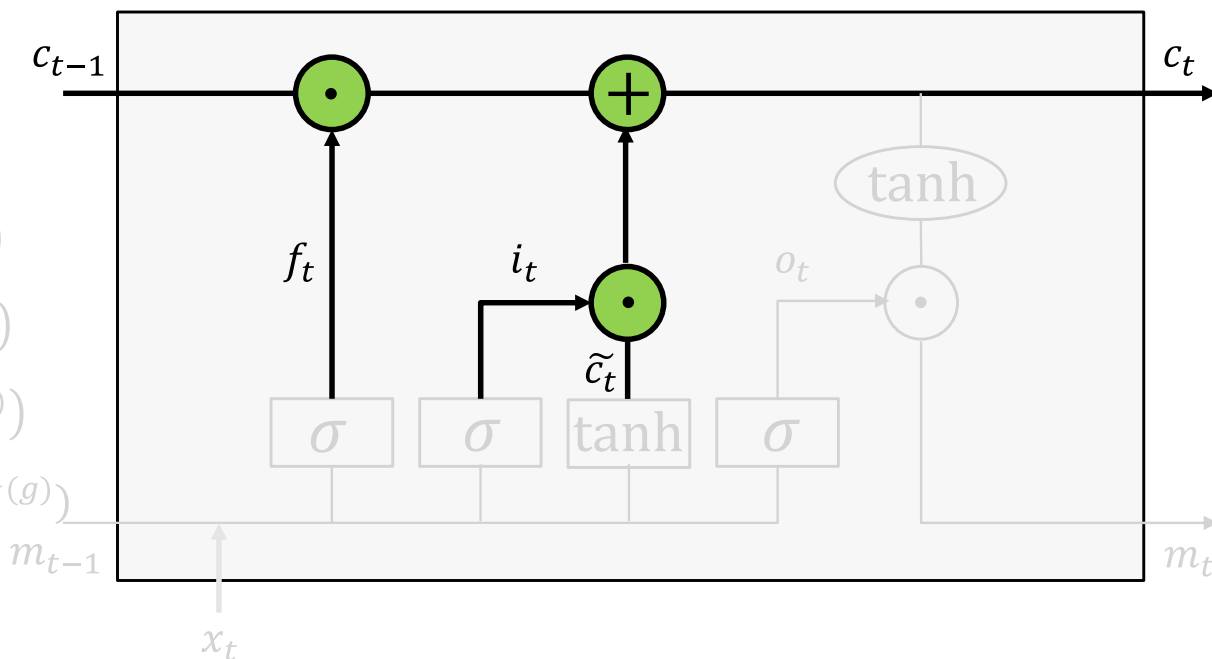
$$f_t = \sigma(x_t U^{(f)} + m_{t-1} W^{(f)})$$

$$o_t = \sigma(x_t U^{(o)} + m_{t-1} W^{(o)})$$

$$\tilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$

$$c_t = c_{t-1} \odot f + \tilde{c}_t \odot i$$

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



حافظه‌ی کوتاه-مدت طولانی

گام به گام با LSTM: گام ۴

LSTM STEP-BY-STEP: STEP 4

خروجی را مدوله می‌کنیم؛ حافظه‌ی جدید را تولید می‌کنیم.

اگر حالت سلول حاوی چیز مربوطی است
 \Leftarrow سیگموئید = 1

$$i_t = \sigma(x_t U^{(i)} + m_{t-1} W^{(i)})$$

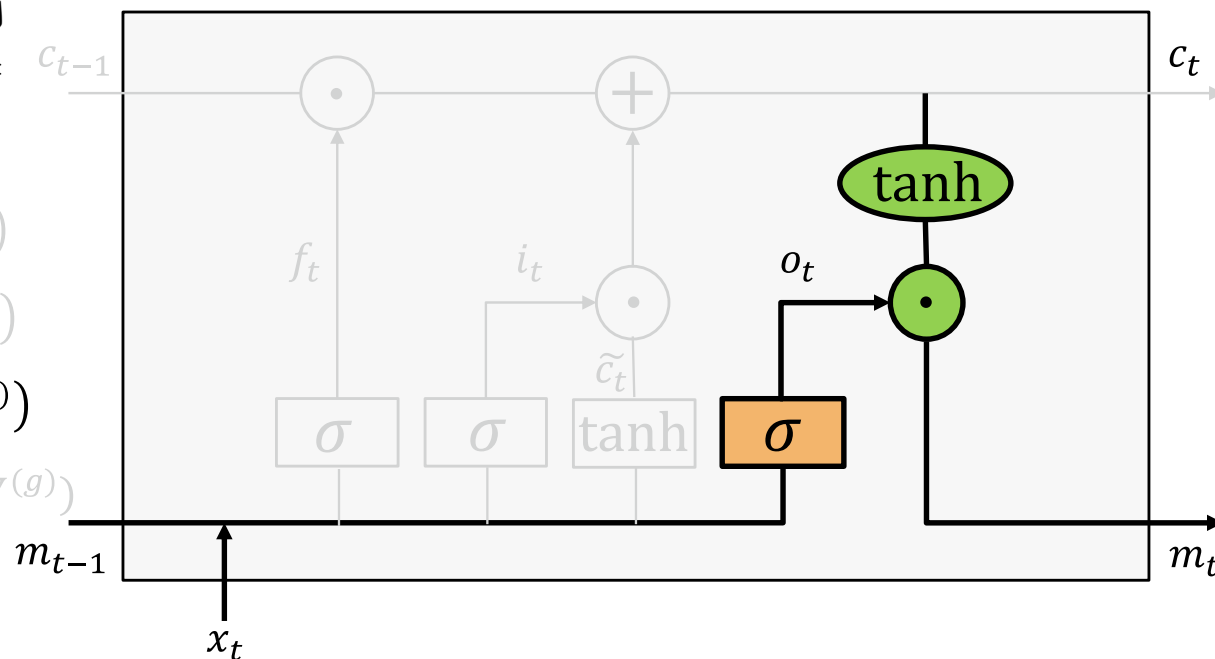
$$f_t = \sigma(x_t U^{(f)} + m_{t-1} W^{(f)})$$

$$o_t = \sigma(x_t U^{(o)} + m_{t-1} W^{(o)})$$

$$\tilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$

$$c_t = c_{t-1} \odot f + \tilde{c}_t \odot i$$

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

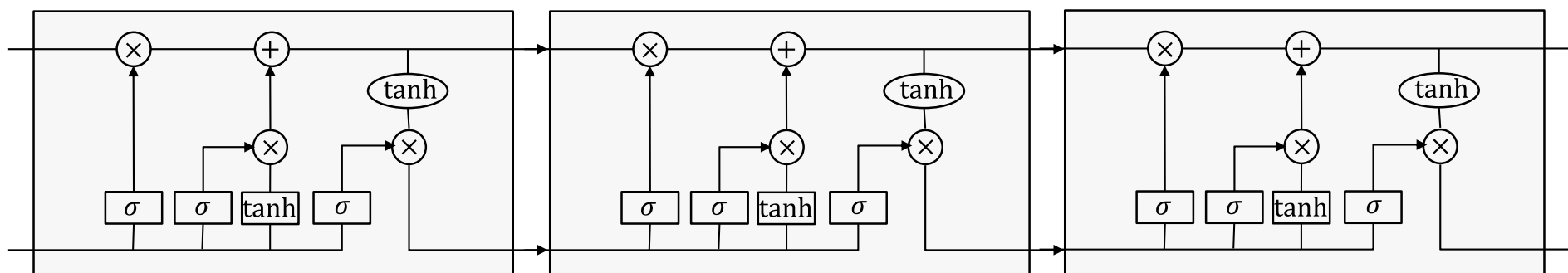


حافظه‌ی کوتاه-مدت طولانی

شبکه‌ی باز شده

LSTM UNROLLED NETWORK

به لحاظ ماکروسکوپی، بسیار شبیه به شبکه‌های عصبی بازگشتی استاندارد است؛
 اما موتور آن کمی متفاوت است (پیچیده‌تر)
 [زیرا گیت‌های LSTM‌های آن وابستگی‌های کوتاه مدت و بلند را تسخیر می‌کند.]



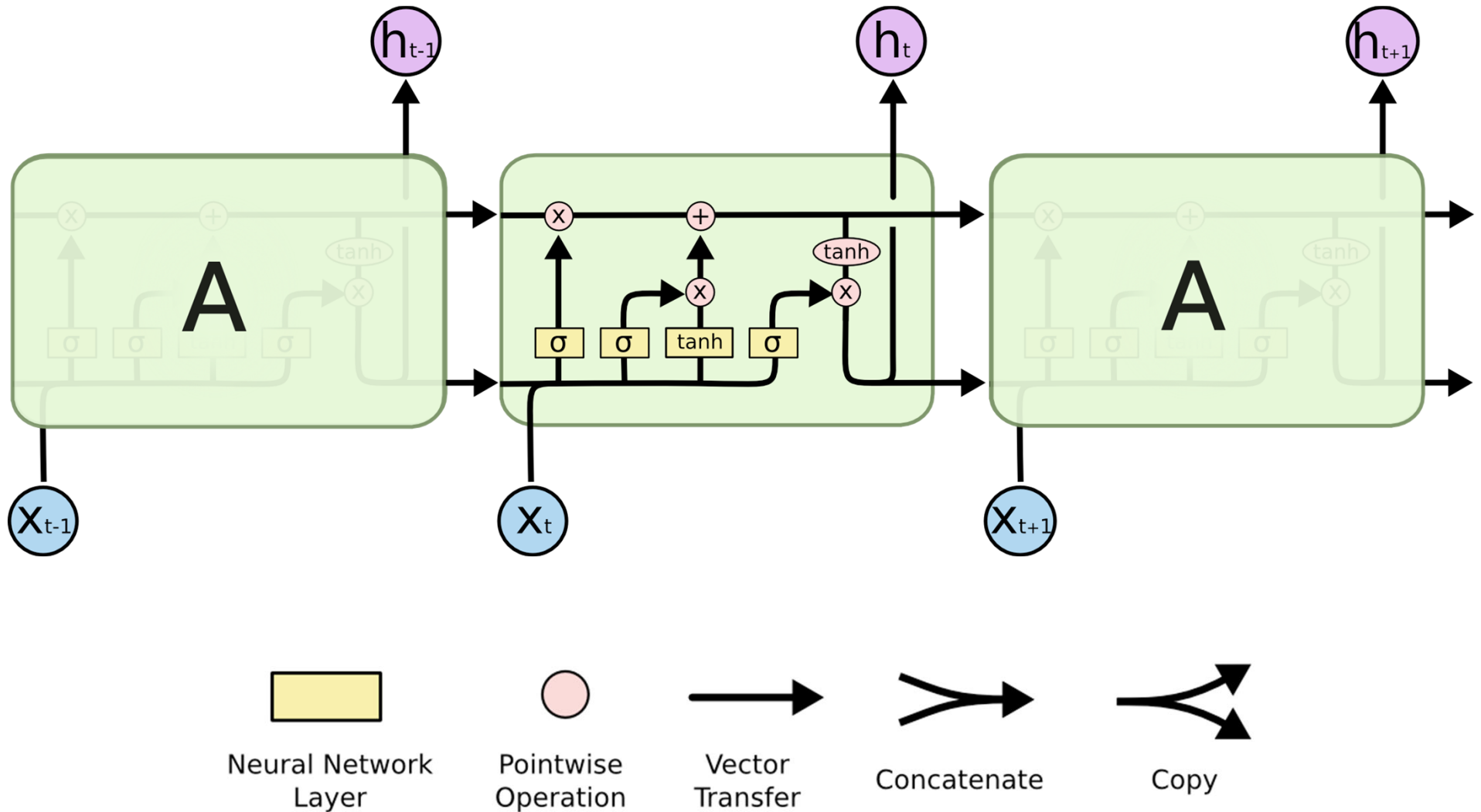
شبکه‌های عصبی بازگشتی

۶

دیگر انواع
شبکه‌های
عصبی
بازگشتی

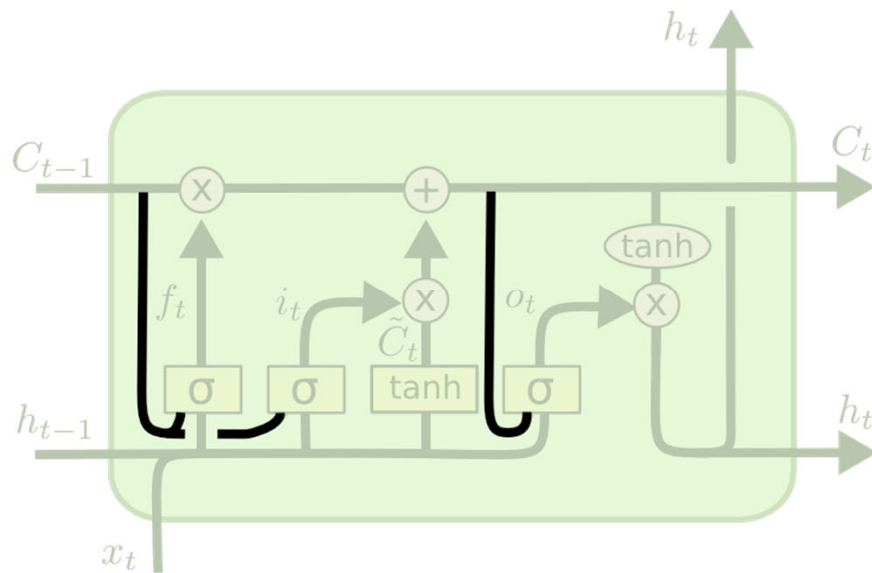
حافظه‌ی کوتاه-مدت طولانی

LSTMs



حافظه‌ی کوتاه-مدت طولانی

تغییر شماره‌ی ۱: اتصالات روزنه‌ای

LSTM VARIANTS #1: PEEPHOLE CONNECTIONS

$$f_t = \sigma (W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

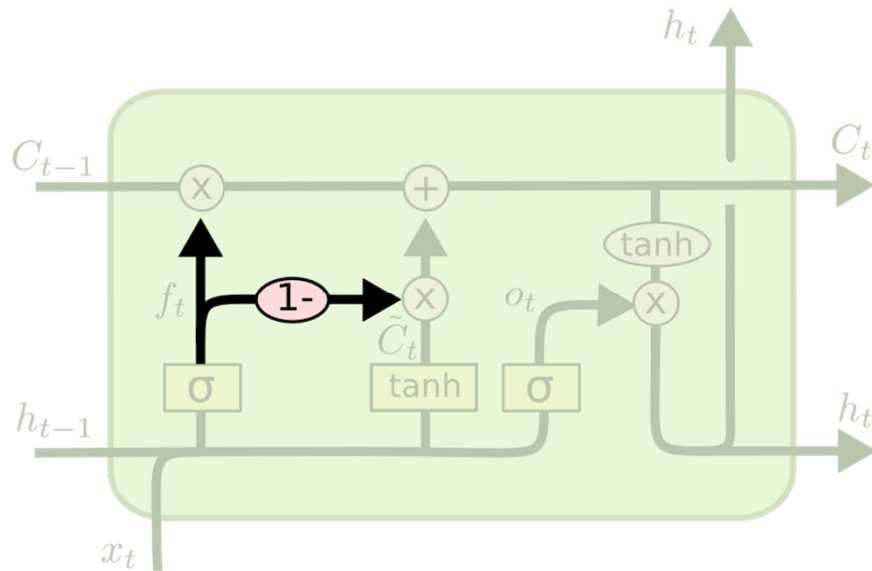
$$i_t = \sigma (W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$

$$o_t = \sigma (W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$

اجازه دادن به گیت‌ها تا بتوانند حالت / حافظه‌ی سلول را ببینند.
Let gates see the cell state / memory.

حافظه‌ی کوتاه-مدت طولانی

تغییر شماره‌ی ۲: گیت‌های جفت‌شده

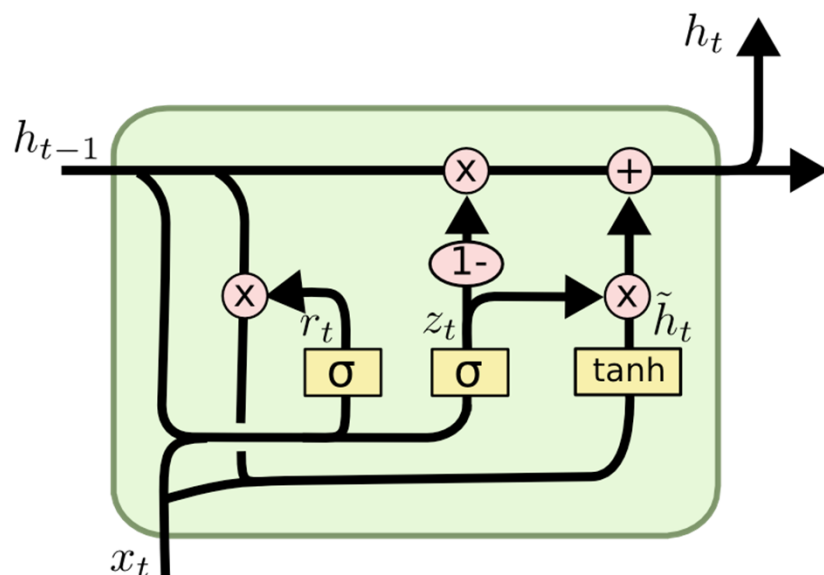
LSTM VARIANTS #2: COUPLED GATES

$$C_t = f_t * C_{t-1} + (1 - f_t) * \tilde{C}_t$$

اطلاعات جدید در صورتی به خاطر سپرده می‌شوند که اطلاعات قدیمی فراموش شوند
Only memorize new if forgetting old

حافظه‌ی کوتاه-مدت طولانی

تغییر شماره‌ی ۳: واحدهای بازگشتی دروازه‌گذاری شده

LSTM VARIANTS #3: GATED RECURRENT UNITS

$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

تغییرات:

حافظه‌ی صریحی وجود ندارد؛ حافظه = خروجی پنهان
 z = اطلاعات جدید را به خاطر بسپار و اطلاعات قدیمی را فراموش کن

Changes:

No explicit memory; memory = hidden output
 Z = memorize new and forget old

Other RNN Variants

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

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

$$r_t = \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r)$$

$$z_t = \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z)$$

$$\tilde{h}_t = \tanh(W_{xh}x_t + W_{hh}(r_t \odot h_{t-1}) + b_h)$$

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t$$

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

MUT1:

$$z = \text{sigm}(W_{xz}x_t + b_z)$$

$$r = \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r)$$

$$h_{t+1} = \tanh(W_{hh}(r \odot h_t) + \tanh(x_t) + b_h) \odot z + h_t \odot (1 - z)$$

MUT2:

$$z = \text{sigm}(W_{xz}x_t + W_{hz}h_t + b_z)$$

$$r = \text{sigm}(x_t + W_{hr}h_t + b_r)$$

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

MUT3:

$$z = \text{sigm}(W_{xz}x_t + W_{hz} \tanh(h_t) + b_z)$$

$$r = \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r)$$

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

شبکه‌های عصبی بازگشتی

خلاصه


- ❖ 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.

شبکه‌های عصبی بازگشتی



منابع


منبع اصلی

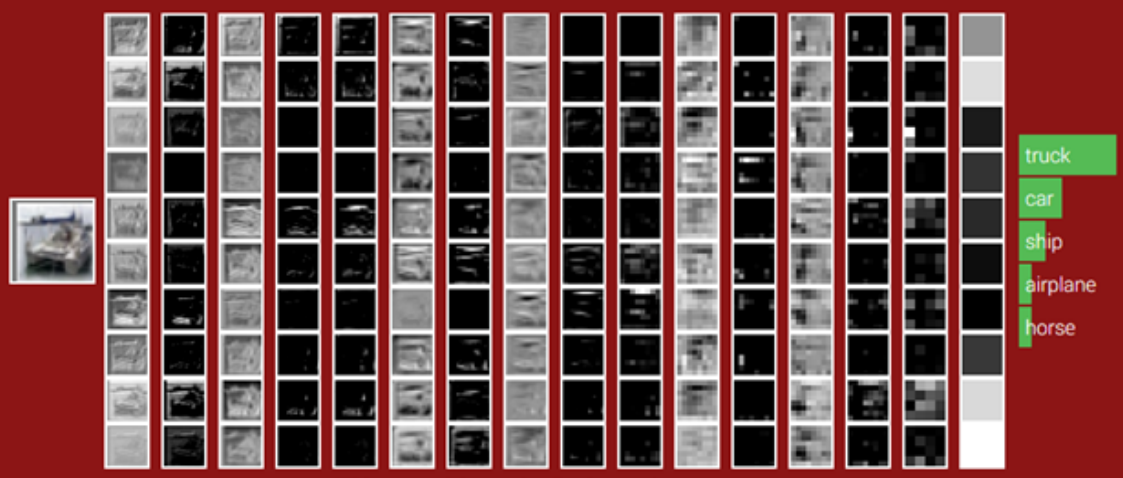


CS231n: Convolutional Neural Networks for Visual Recognition

Spring 2019

Previous Years: [\[Winter 2015\]](#) [\[Winter 2016\]](#) [\[Spring 2017\]](#) [\[Spring 2018\]](#)





*This network is running live in your browser

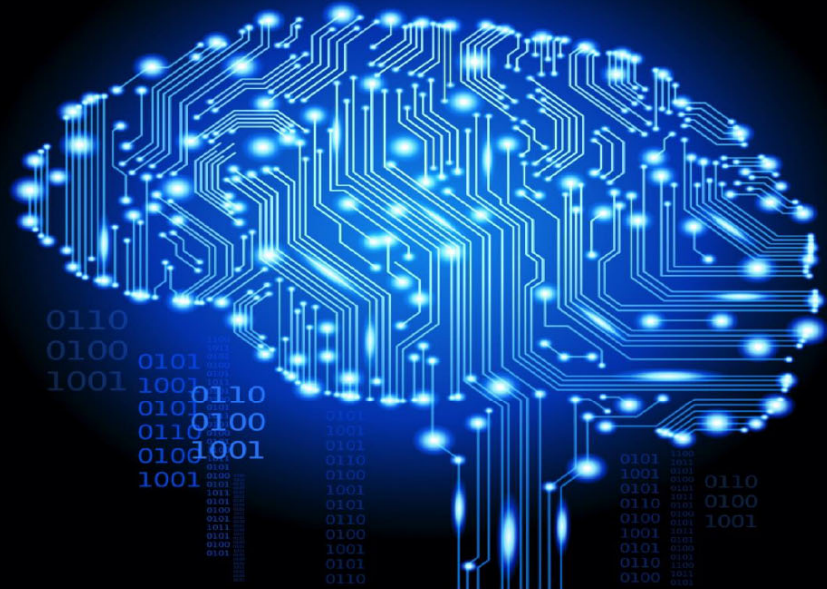
Course Description

Computer Vision has become ubiquitous in our society, with applications in search, image understanding, apps, mapping, medicine, drones, and self-driving cars. Core to many of these applications are visual recognition tasks such as image classification, localization and detection. Recent developments in neural network (aka "deep learning") approaches have greatly advanced the performance of these state-of-the-art visual recognition systems. This course is a deep dive into details of the deep learning architectures with a focus on learning end-to-end models for these tasks, particularly image classification. During the 10-week course, students will learn to implement, train and debug their own neural networks and gain a detailed understanding of cutting-edge research in computer vision. The final assignment will involve training a multi-million parameter convolutional neural network and applying it on the largest image classification dataset (ImageNet). We will focus on teaching how to set up the problem of image recognition, the learning algorithms (e.g. backpropagation), practical engineering tricks for training and fine-tuning the networks and guide the students through hands-on assignments and a final course project. Much of the background and materials of this course will be drawn from the [ImageNet Challenge](#).

<http://cs231n.stanford.edu>

<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

منبع کمکی



Lecture 8: Recurrent Neural Networks

Deep Learning @ UvA

UVA DEEP LEARNING COURSE – EFSTRATIOS GAVVES

RECURRENT NEURAL NETWORKS - 1

<https://uvadlc.github.io/>

colah's blog

[Blog](#)
[About](#)
[Contact](#)

Understanding LSTM Networks

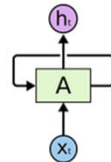
Posted on August 27, 2015

Recurrent Neural Networks

Humans don't start their thinking from scratch every second. As you read this essay, you understand each word based on your understanding of previous words. You don't throw everything away and start thinking from scratch again. Your thoughts have persistence.

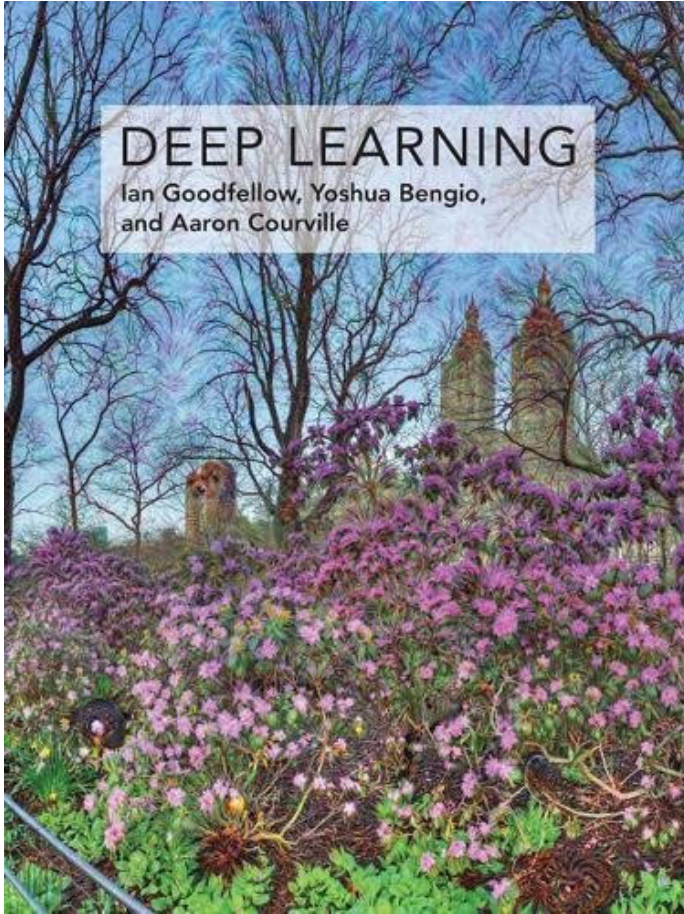
Traditional neural networks can't do this, and it seems like a major shortcoming. For example, imagine you want to classify what kind of event is happening at every point in a movie. It's unclear how a traditional neural network could use its reasoning about previous events in the film to inform later ones.

Recurrent neural networks address this issue. They are networks with loops in them, allowing information to persist.



Recurrent Neural Networks have loops.

<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>



I. Goodfellow, Y. Bengio, A. Courville,
Deep Learning,
 MIT Press, 2016.

Chapter 10

Chapter 10

Sequence Modeling: Recurrent and Recursive Nets

Recurrent neural networks or RNNs (Rumelhart *et al.*, 1986a) are a family of neural networks for processing sequential data. Much as a convolutional network is a neural network that is specialized for processing a grid of values \mathbf{X} such as an image, a recurrent neural network is a neural network that is specialized for processing a sequence of values $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$. Just as convolutional networks can readily scale to images with large width and height, and some convolutional networks can process images of variable size, recurrent networks can scale to much longer sequences than would be practical for networks without sequence-based specialization. Most recurrent networks can also process sequences of variable length.

To go from multi-layer networks to recurrent networks, we need to take advantage of one of the early ideas found in machine learning and statistical models of the 1980s: sharing parameters across different parts of a model. Parameter sharing makes it possible to extend and apply the model to examples of different forms (different lengths, here) and generalize across them. If we had separate parameters for each value of the time index, we could not generalize to sequence lengths not seen during training, nor share statistical strength across different sequence lengths and across different positions in time. Such sharing is particularly important when a specific piece of information can occur at multiple positions within the sequence. For example, consider the two sentences “I went to Nepal in 2009” and “In 2009, I went to Nepal.” If we ask a machine learning model to read each sentence and extract the year in which the narrator went to Nepal, we would like it to recognize the year 2009 as the relevant piece of information, whether it appears in the sixth