

Deep learning in NLP

Part II: Introduction to RNN and LSTM

Zeinab Rahimi, Feb 2022

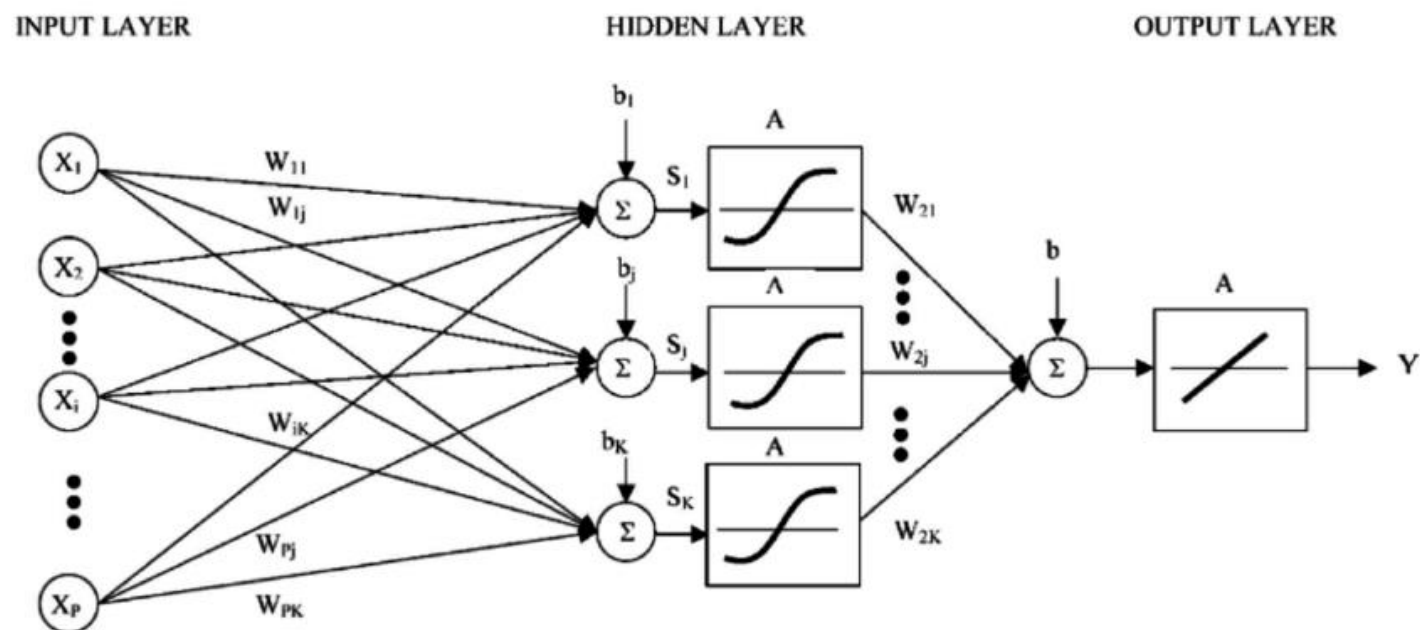
شبکه های عصبی سنتی

□ لایه ورودی

□ لایه هیدن

□ لایه خروجی

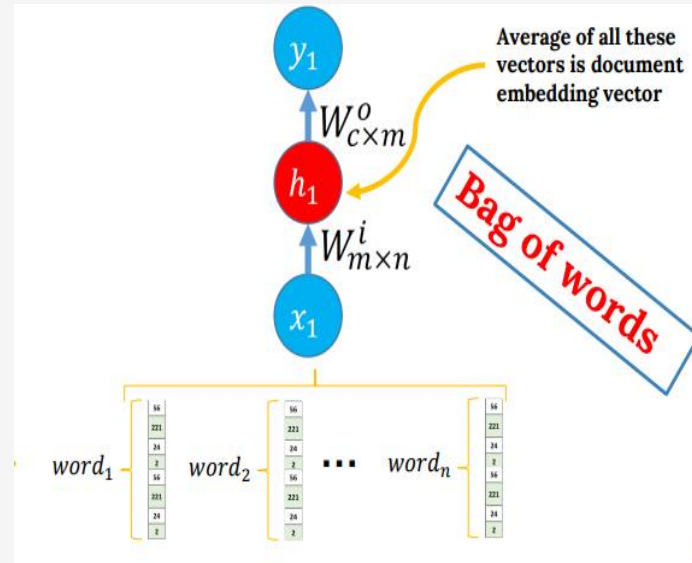
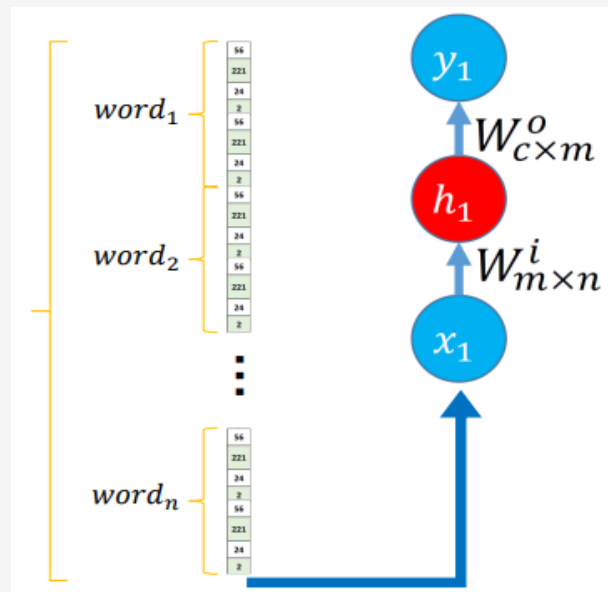
□ فعالسازی



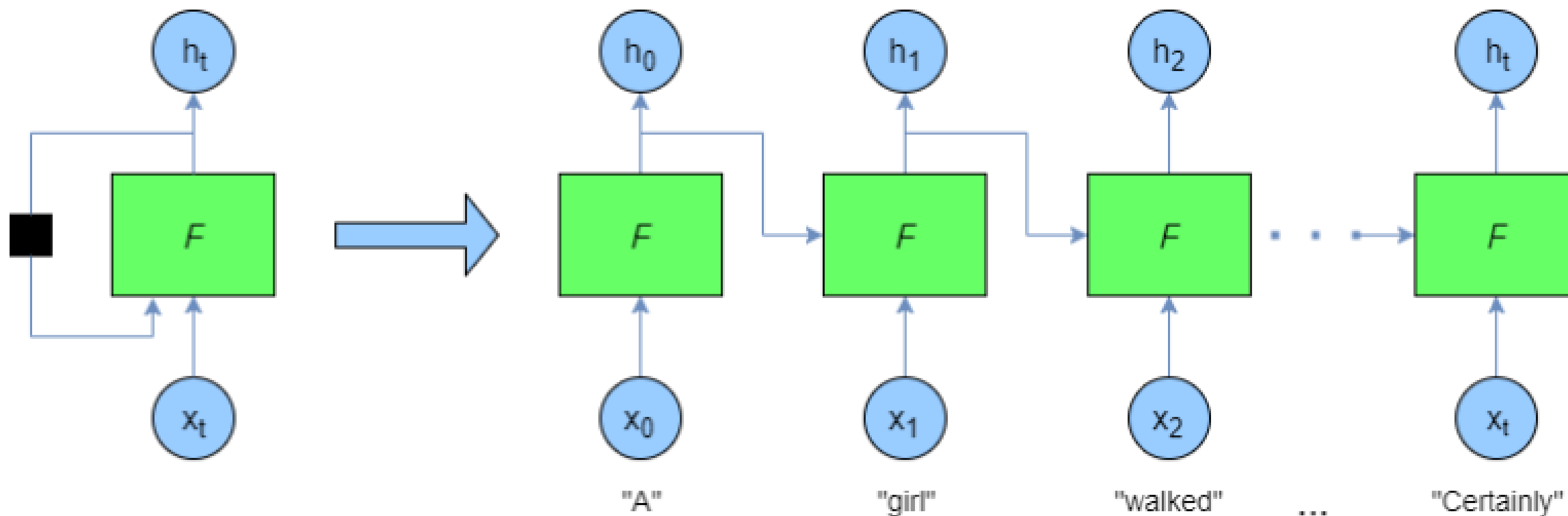
مشکلات شبکه های قبلی برای NLP

- در تسک های قبلی مثل Word2vec از شبکه های FullyConnected استفاده می شد:
- ورود یک کلمه به شبکه و انجام تسک طبقه بندی
- در هر مرحله خروجی بر اساس ورودی همان مرحله مثلا دسته بندی بر اساس آخرین ورودی
- نیاز به تصمیم گیری بعد از دیدن یک رشته توکن با ترتیب (seq)
- نیاز به ورود یک جمله یا سند برای تسک های پیچیده تر
- نیاز به داشتن اطلاعات از کلمات قبل و بعد در رشته ورودی
- نیاز به در نظر گرفتن ترتیب توکن ها
- خروجی باید به ورودی فعلی و خروجی قبلی وابسته باشد.

راه حل برای نواقص شبکه های سنتی



- میانگین گرفتن از بردارهای جاسازی
- در نظر نگرفتن توالی و ساختار جمله
- مناسب نبودن برای تسک های پیچیده تر
- به هم وصل کردن بردارهای جاسازی
- ثابت نبودن طول ورودی شبکه
- زیاد بودن طول ورودی
- معرفی شبکه های بازگشتی (RNN)
- همان FC ها با این تفاوت که لایه های هیدن آنها به هم وصل است.
- در FC اطلاعاتی از کلمات قبل و بعد نداریم اما در RNN در هر لحظه از قبلی ها مطلعیم.
- در هر زمان با ورود هر توکن اطلاعاتی در لایه هیدن ذخیره می شود و به لایه هیدن در زمان بعدی منتقل می شود.

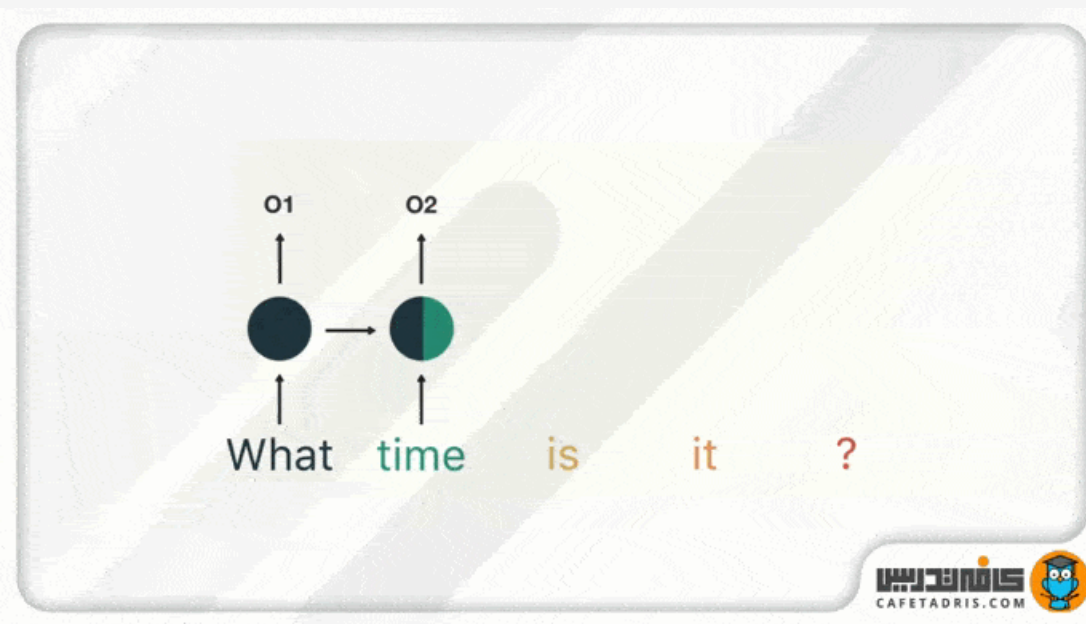
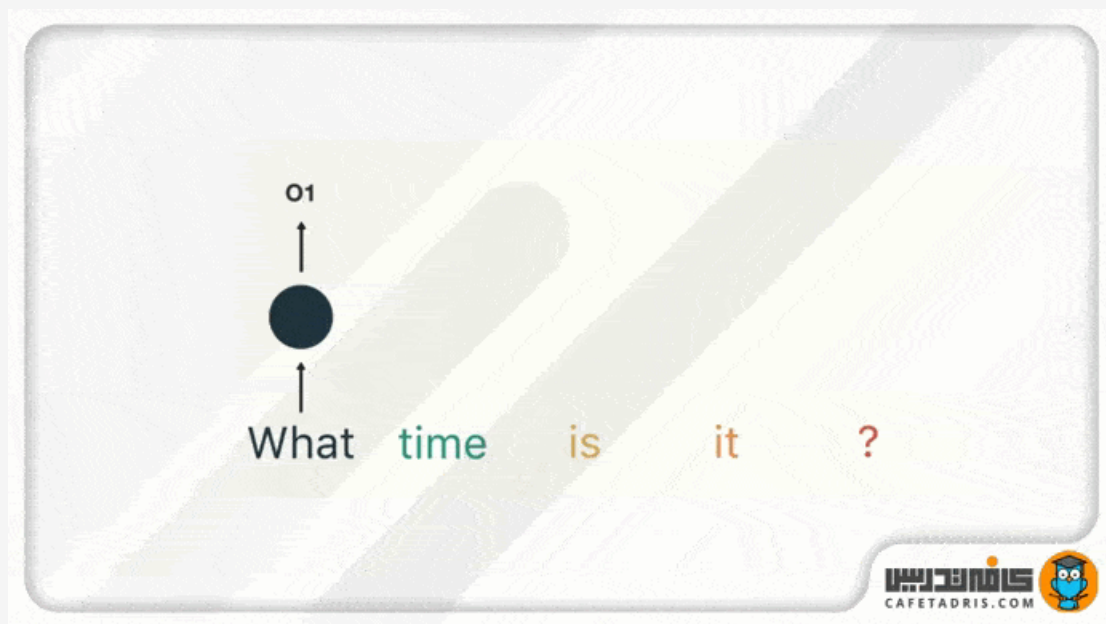


Recurrent Neural Network

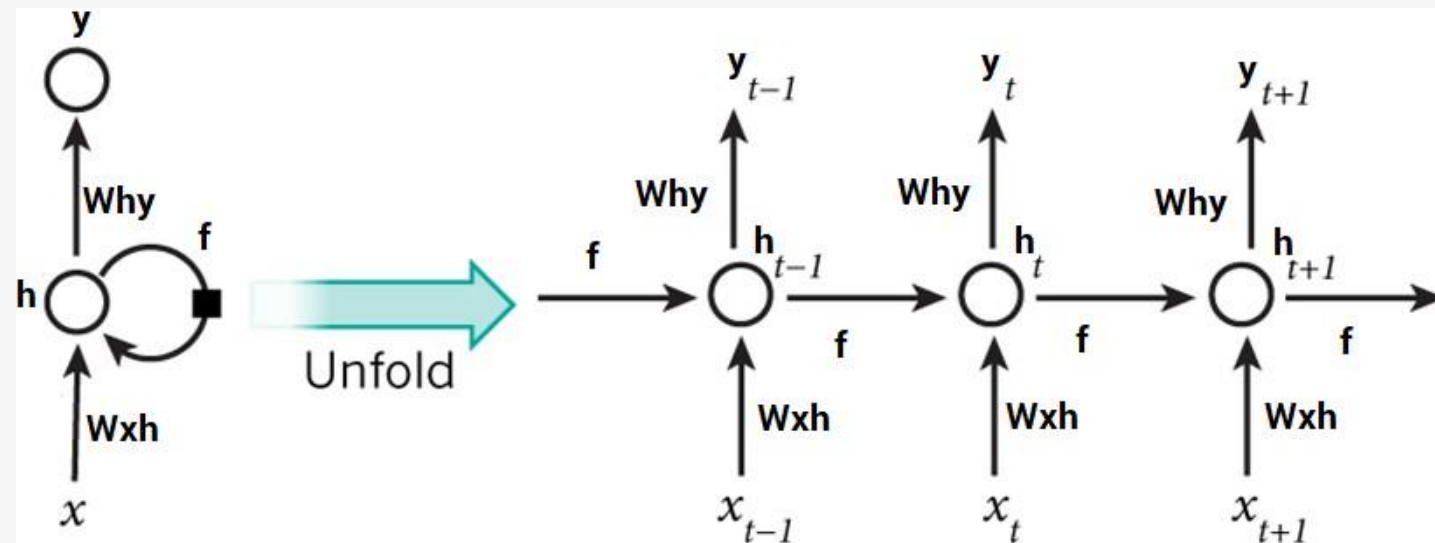
از خروجی قبلی به عنوان ورودی بعدی استفاده می کند.

شبکه های عصبی بازگشتی

- شکستن جمله به کلمات
- شبکه ی RNN به صورت ترتیبی کار می کند، در هر گام یک کلمه را به آن وارد می شود.(بردار کلمه)
- زمانی که به آخر جمله برسیم می بینیم که شبکه ی اطلاعات مربوط به تمامی کلمه های موجود در جمله را پردازش کرده است.



شبکه های عصبی بازگشتی

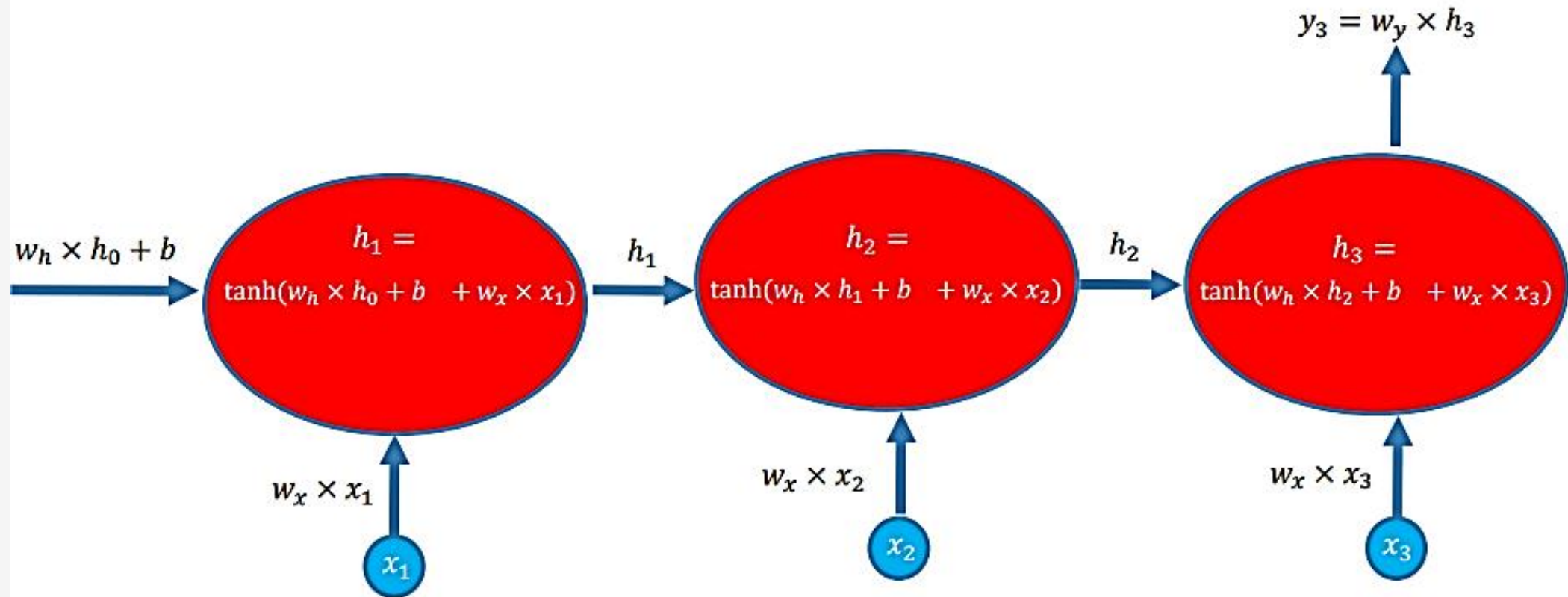


$$h_t = f(h_{t-1}, x_t)$$

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

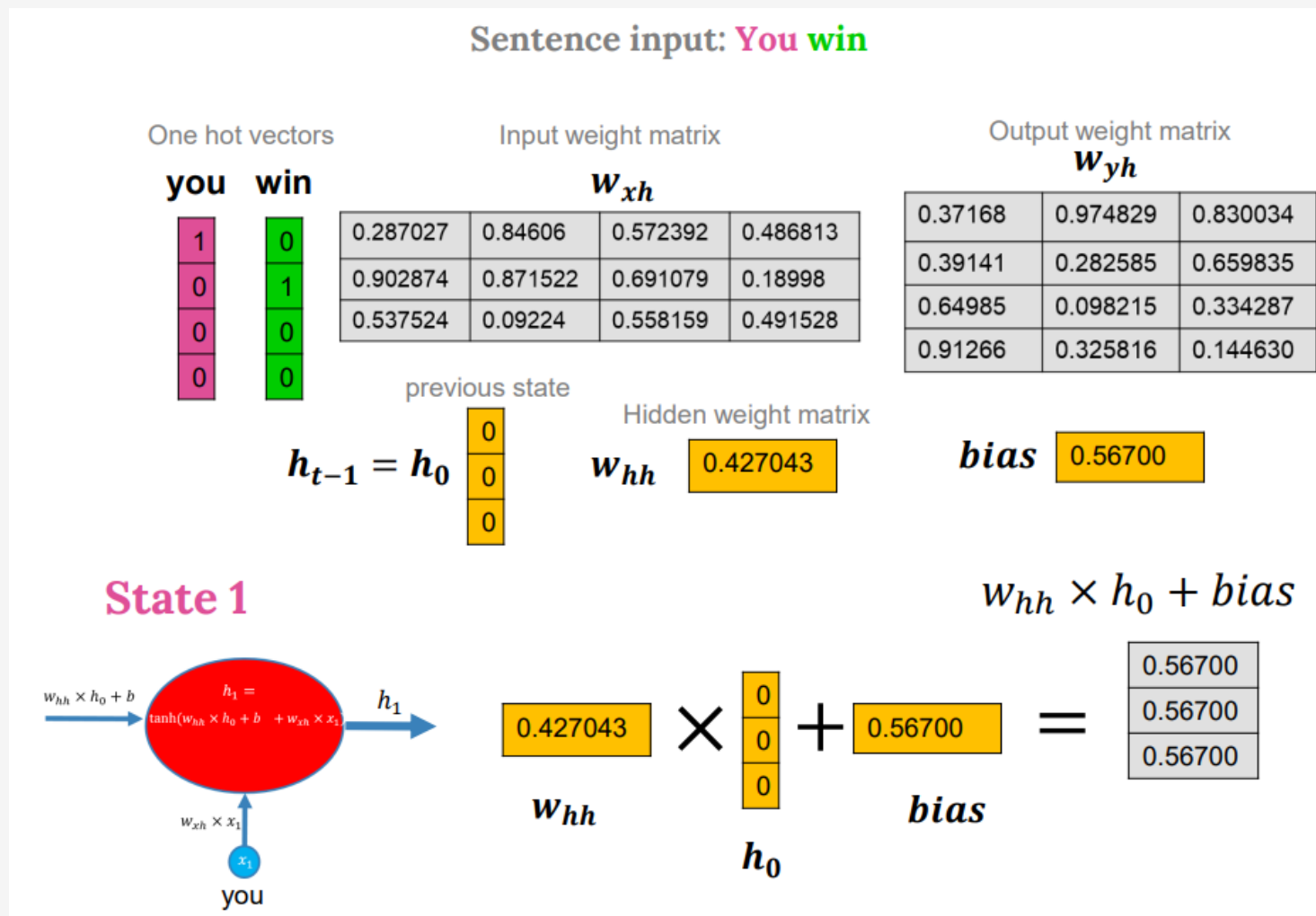
$$y_t = W_{hy}h_t$$

شبکه های عصبی بازگشتی: روابط



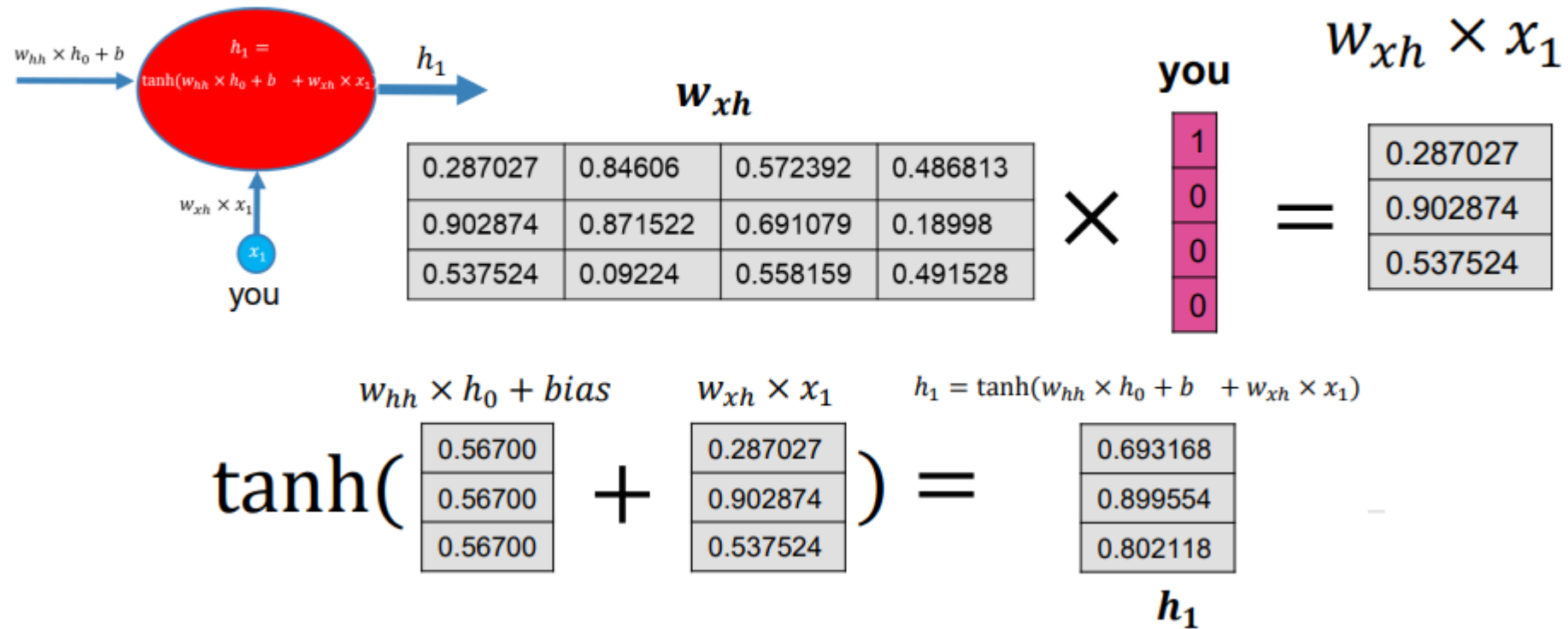
مثال برداشته شده از (Deep NLP workshop (IKT,2021)

شبکه های عصبی بازگشتی: مثال

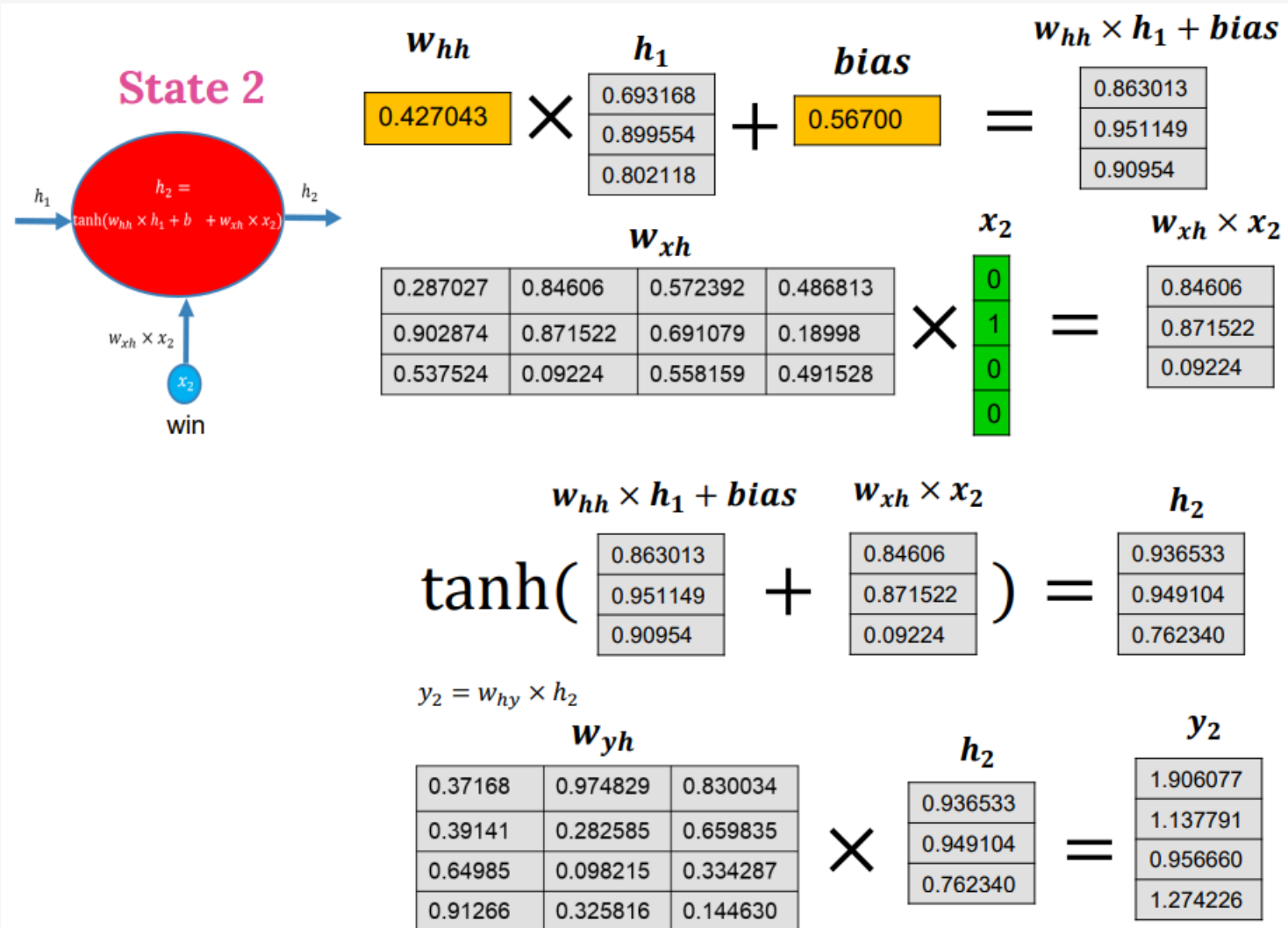


شبکه های عصبی بازگشتی: مثال

State 1



شبکه های عصبی بازگشتی: مثال



$$Softmax(y_2) = \begin{matrix} 1.906077 \\ 1.137791 \\ 0.956660 \\ 1.274226 \end{matrix} = \begin{matrix} 0.41974 \\ 0.1946 \\ 0.1624 \\ 0.2231 \end{matrix} \quad \begin{matrix} 1 \\ 0 \\ 0 \\ 0 \end{matrix}$$

شبکه های عصبی بازگشتی: Back Propagation

Feed Forward

$$h = \tanh(w_h \times x_t + b)$$

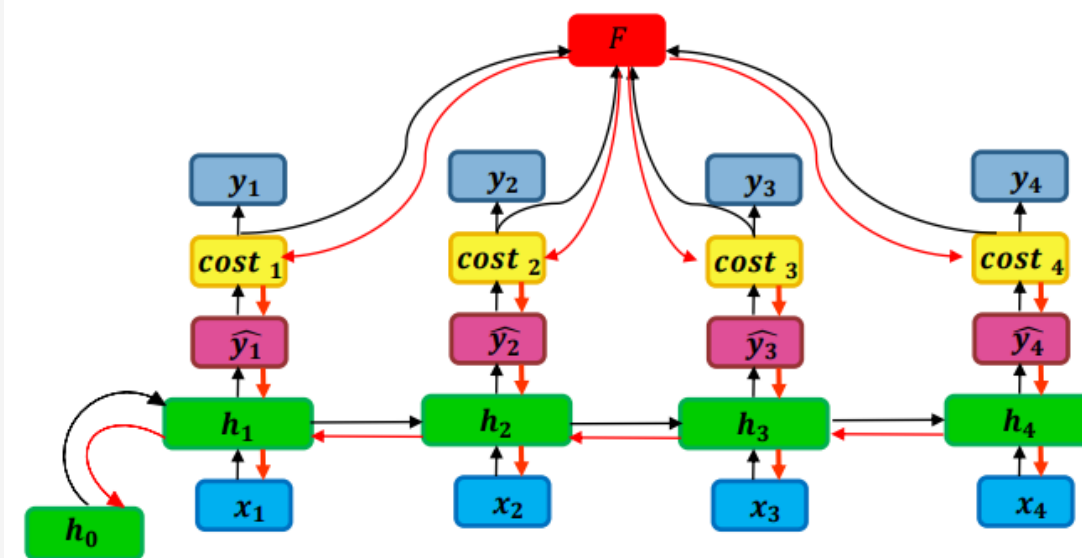
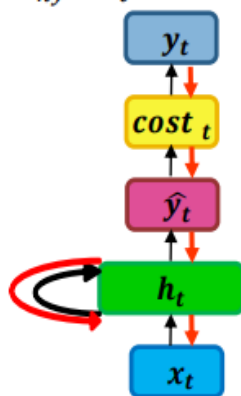
$$\hat{y} = w_{hy} \times h$$



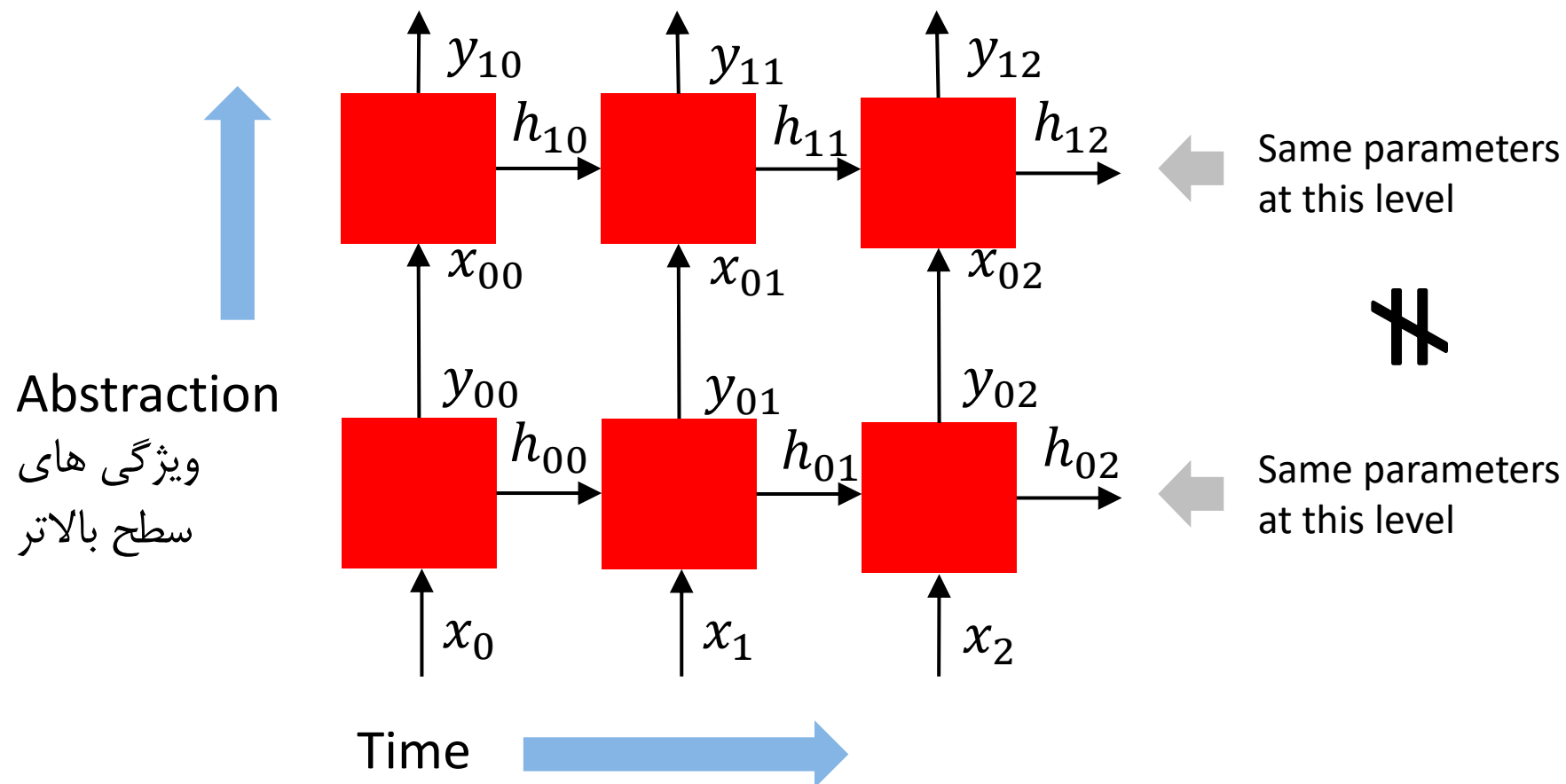
Recurrent Neural Networks

$$h_t = \tanh(w_h \times [h_{t-1}, x_t] + b)$$

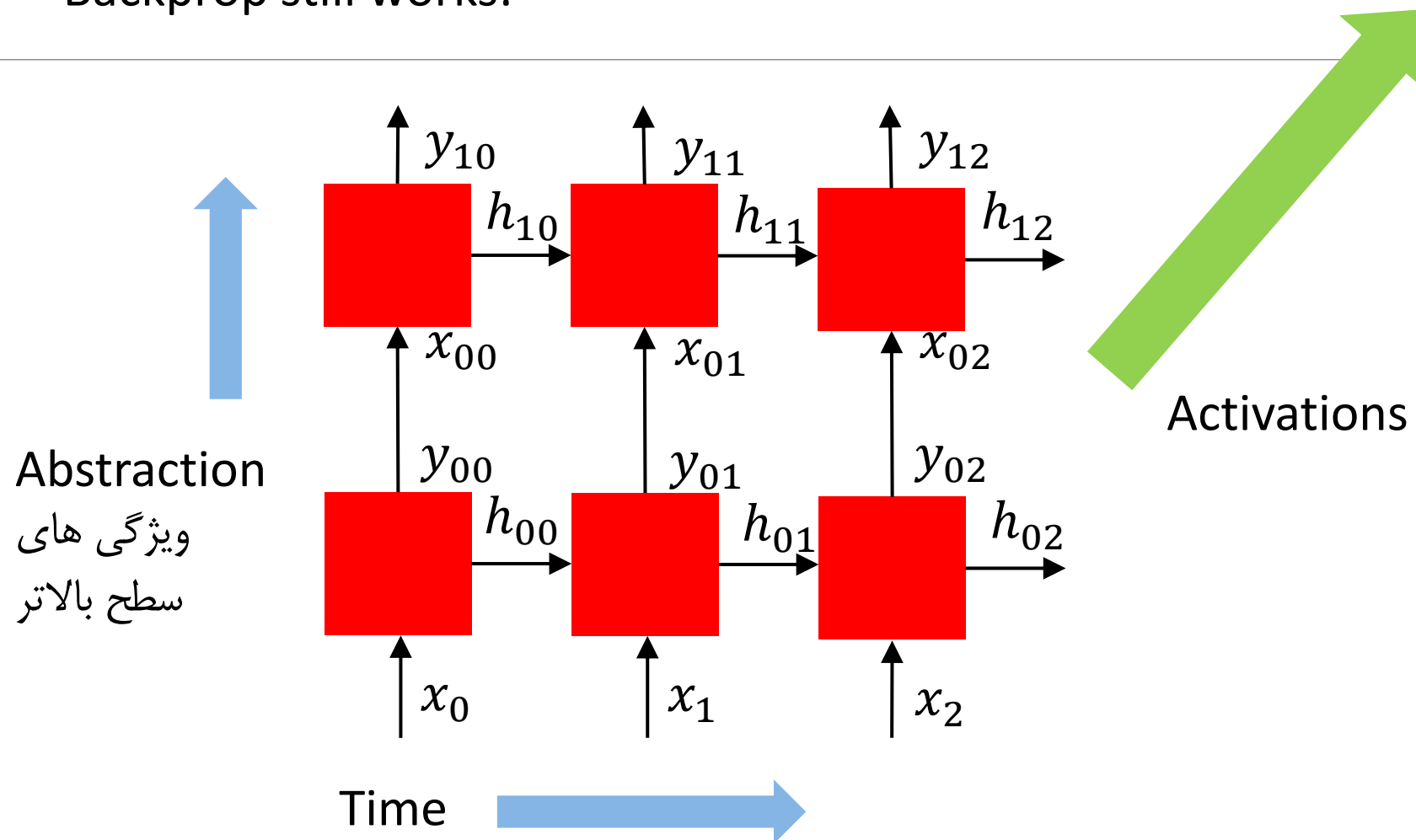
$$\hat{y}_t = w_{hy} \times h_t$$



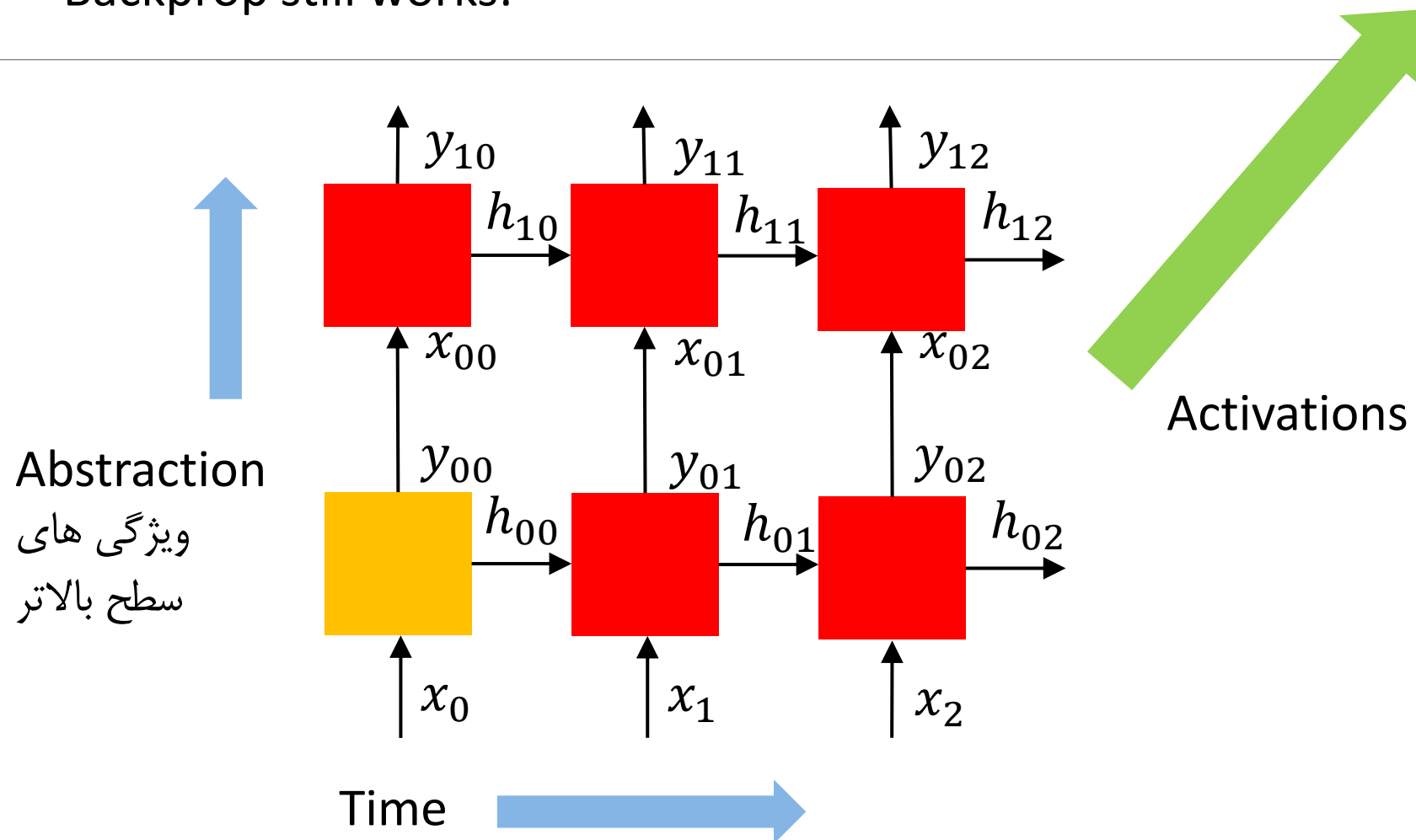
در RNN عمیق غالباً لایه ها به صورت عمودی استک می شوند.



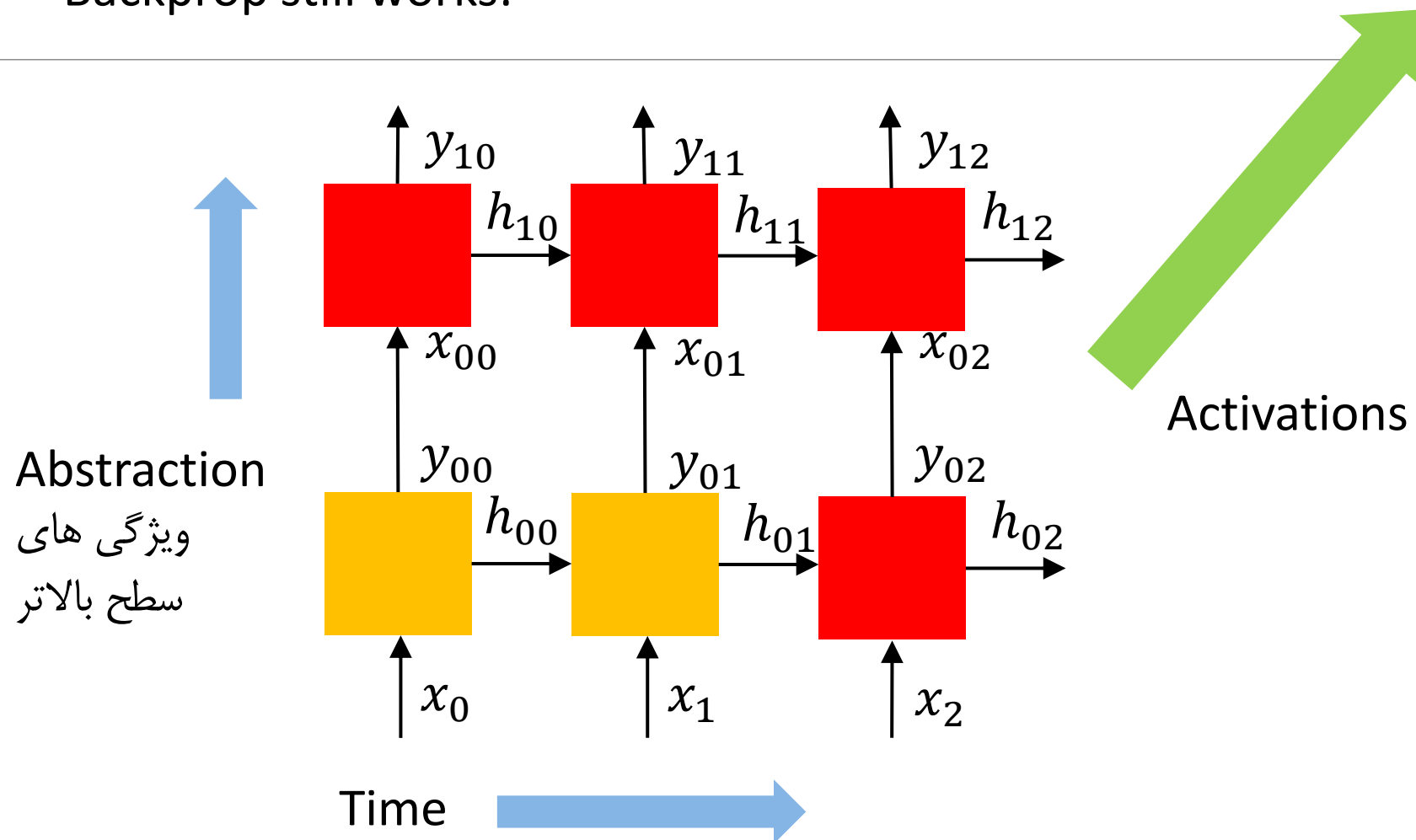
Backprop still works:



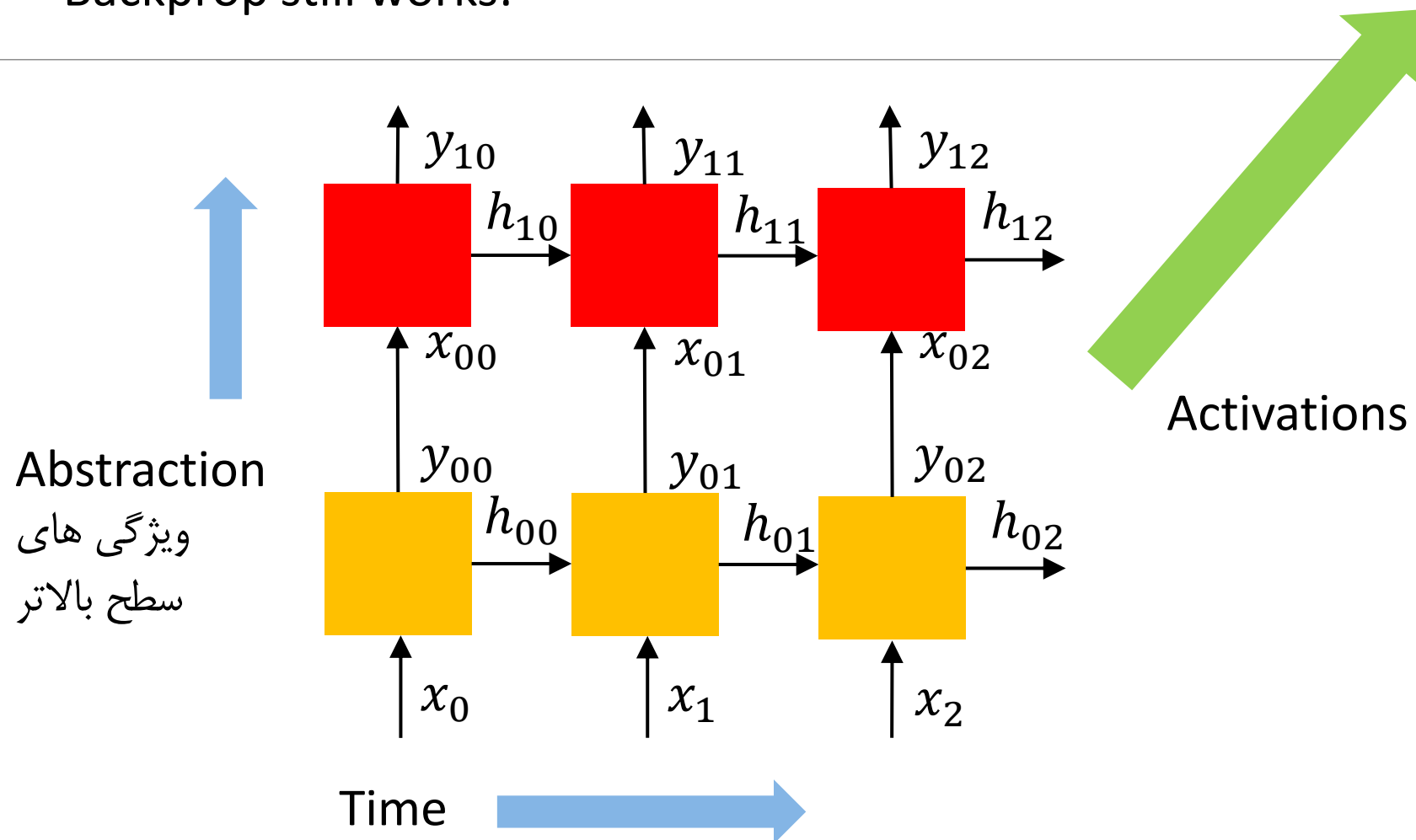
Backprop still works:



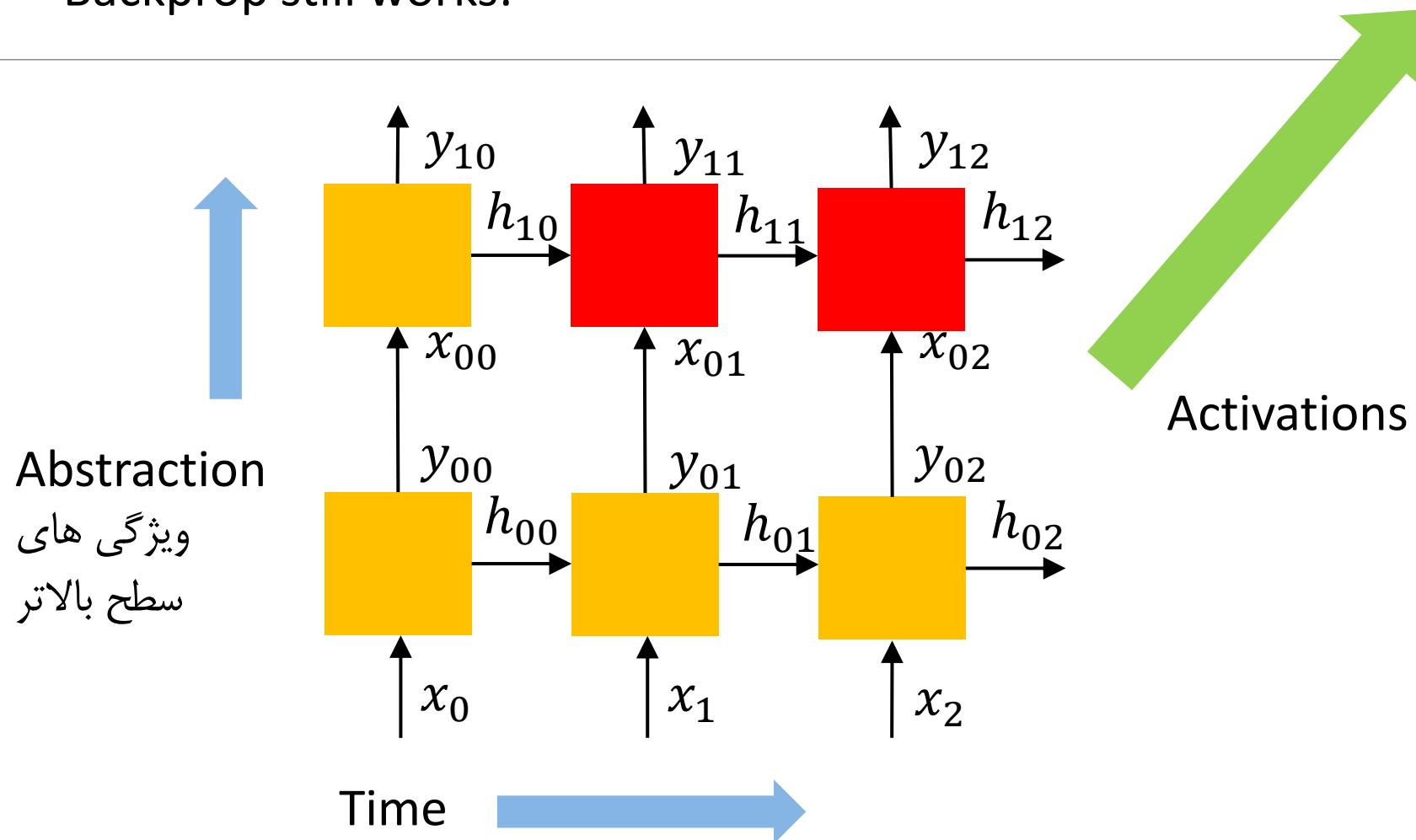
Backprop still works:



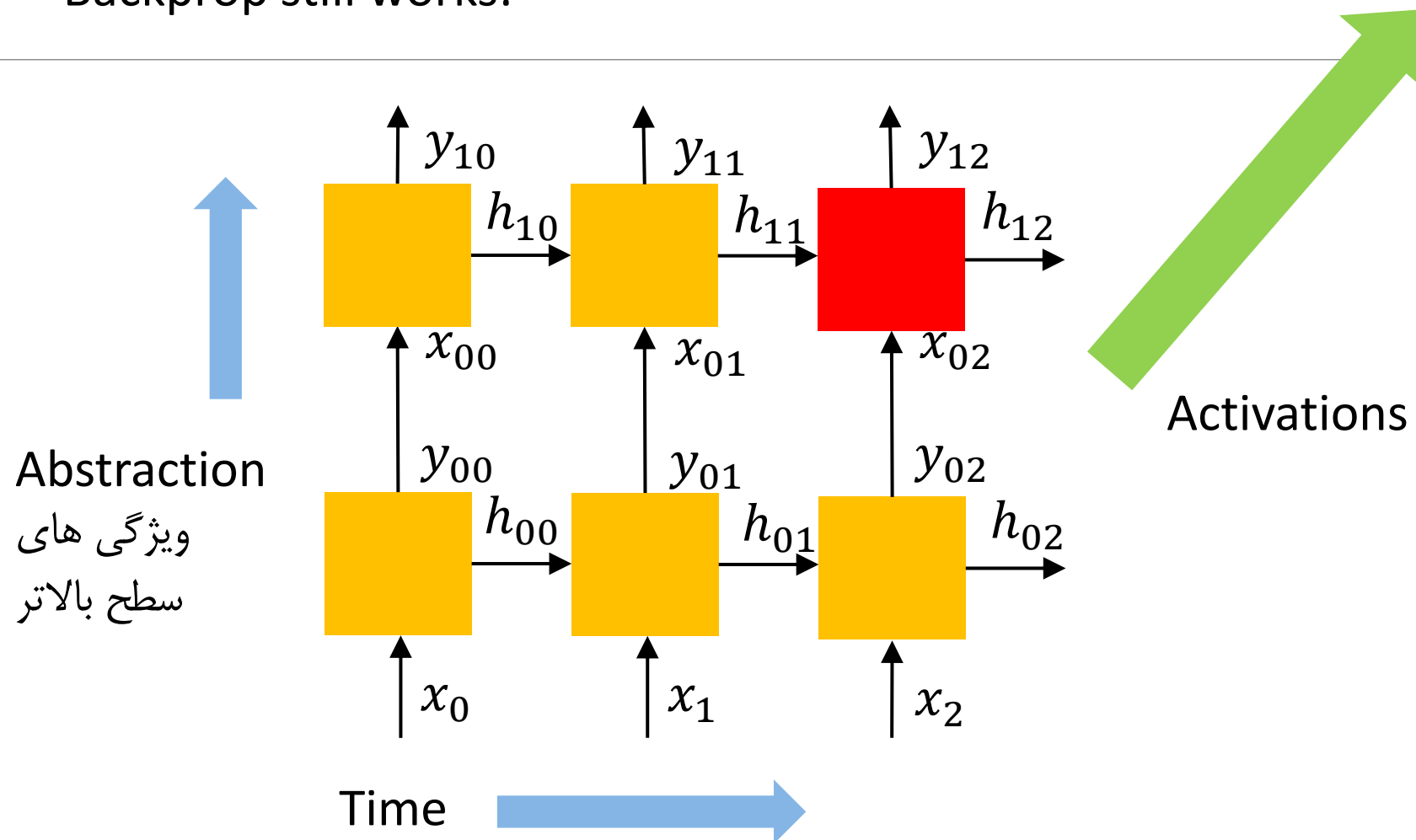
Backprop still works:



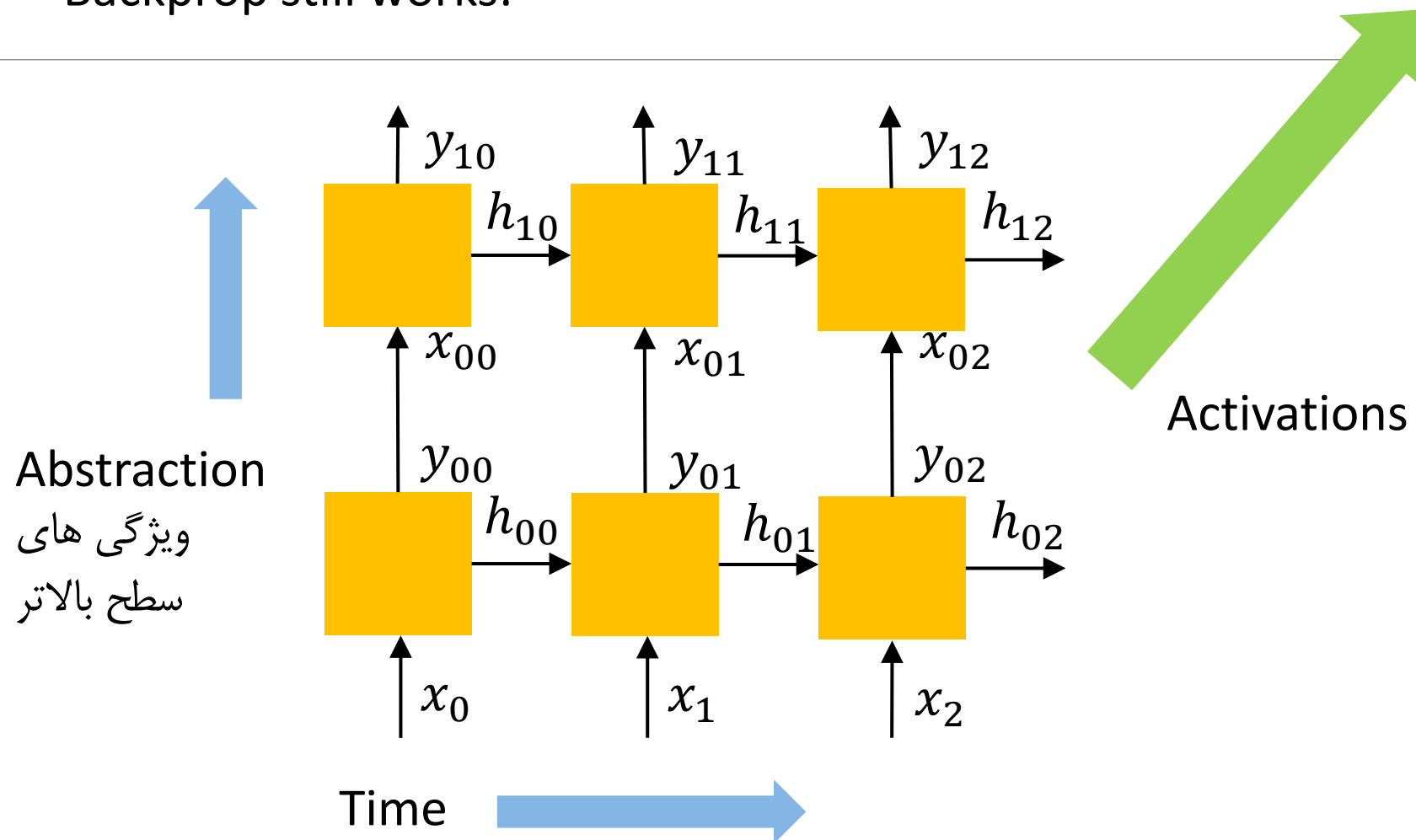
Backprop still works:



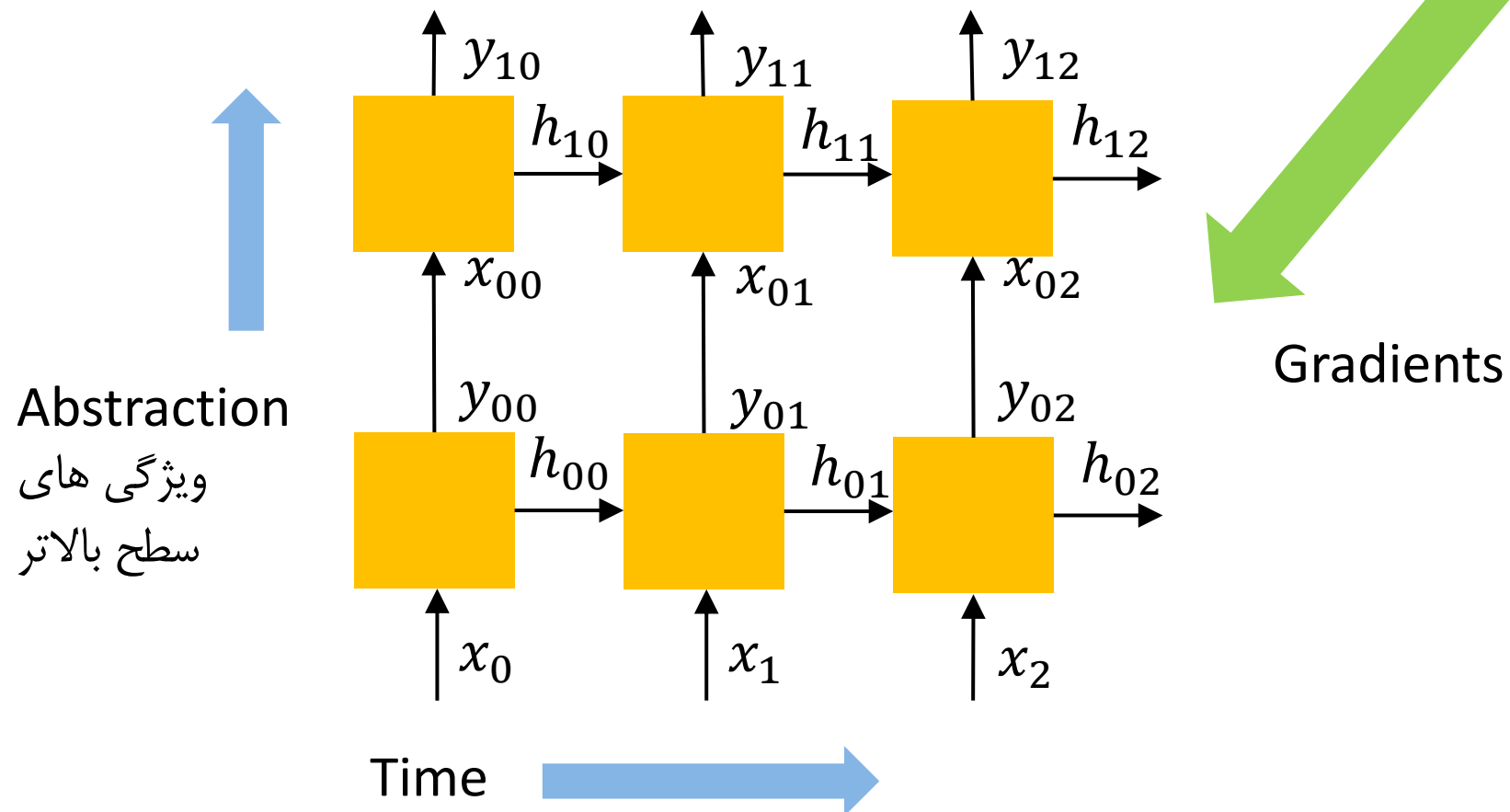
Backprop still works:



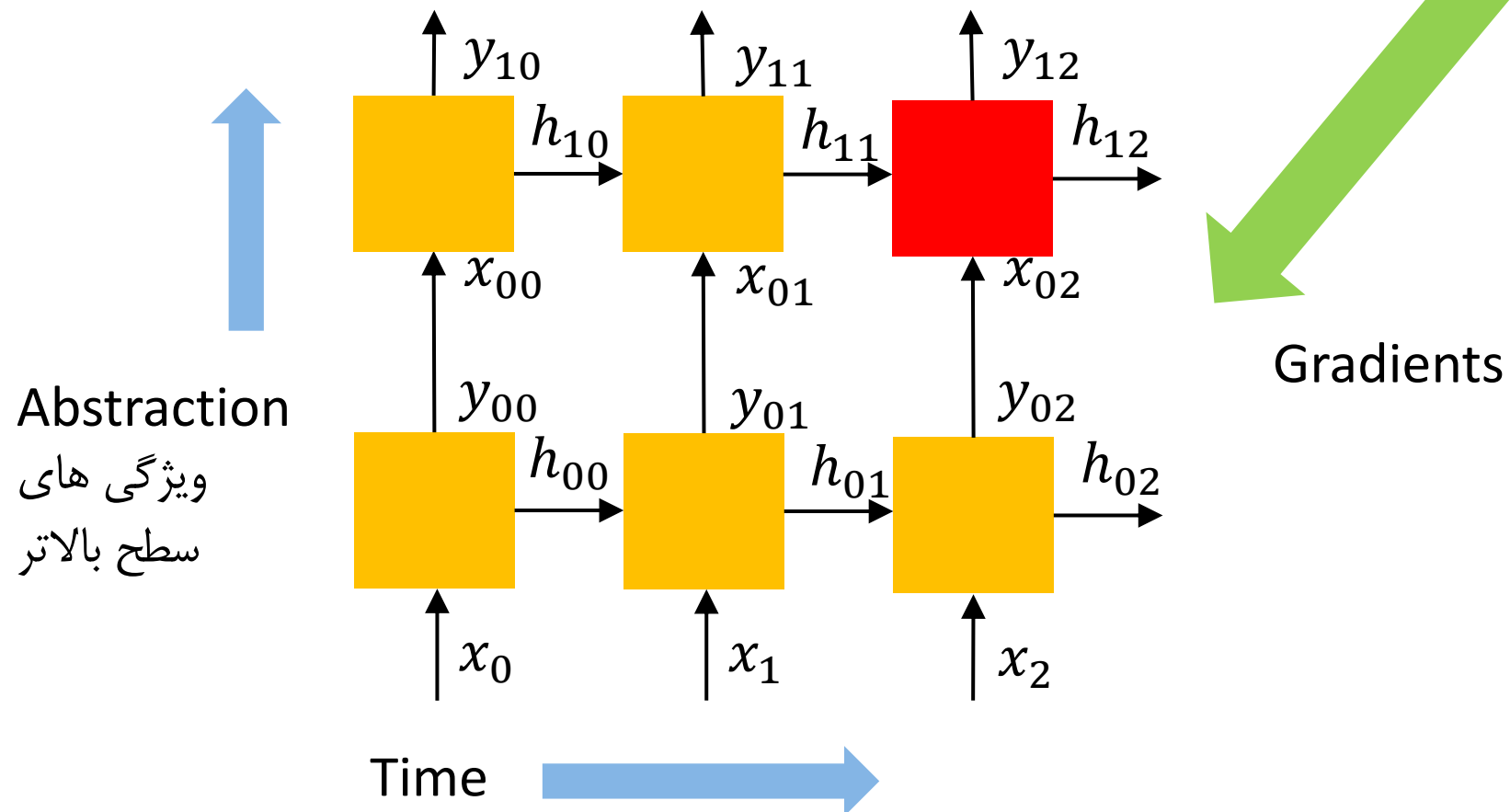
Backprop still works:



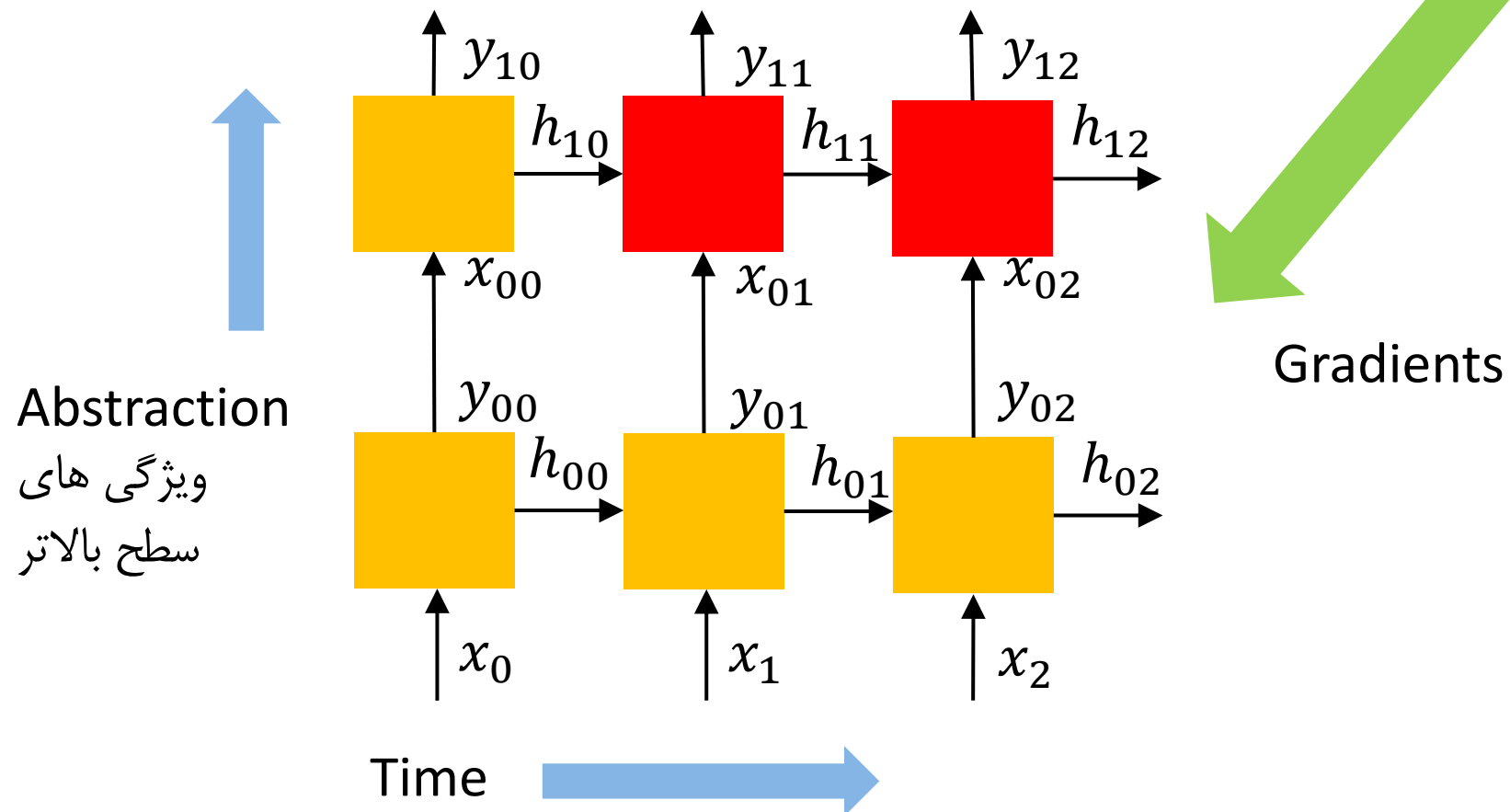
Backprop still works:



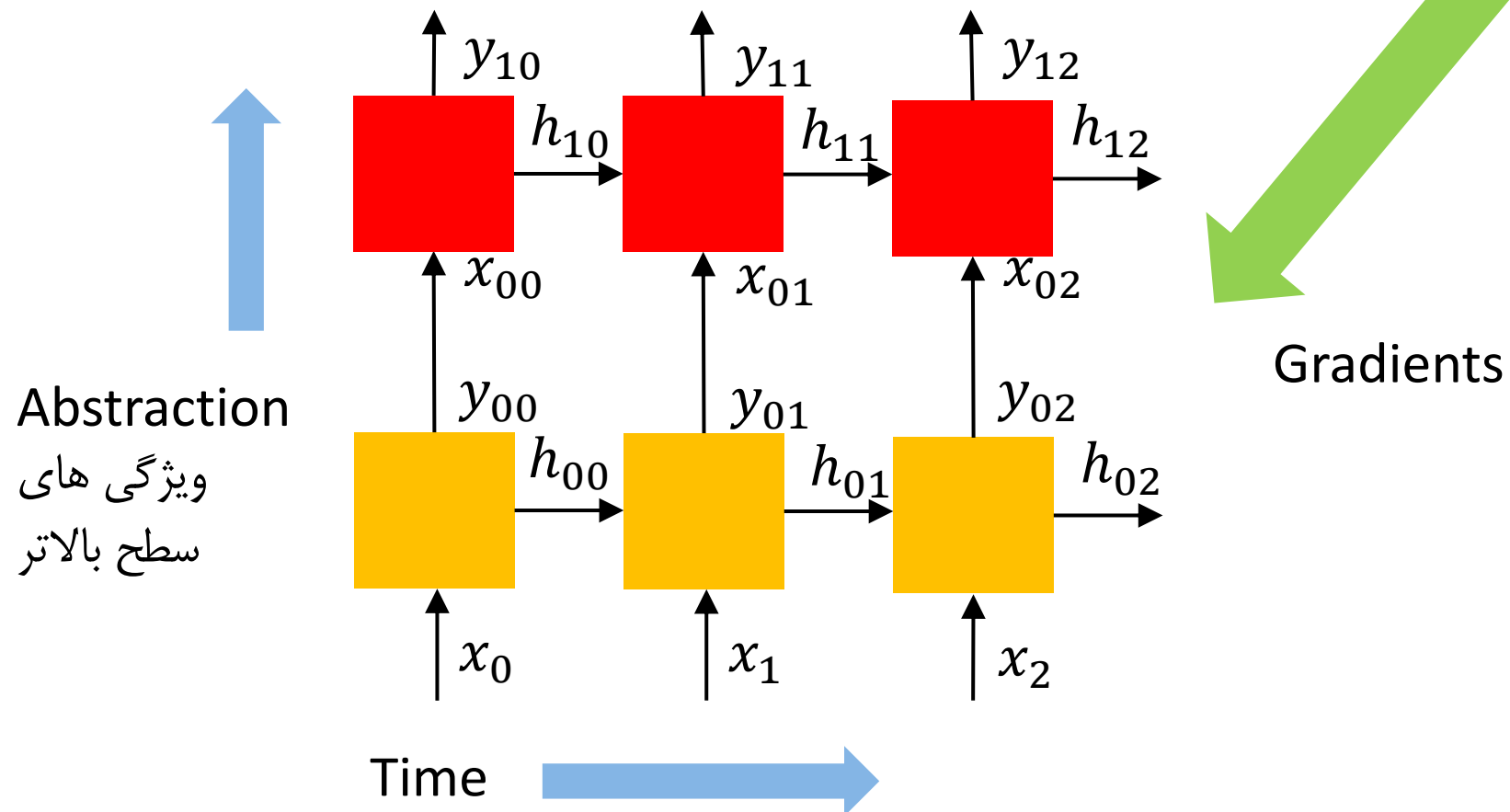
Backprop still works:



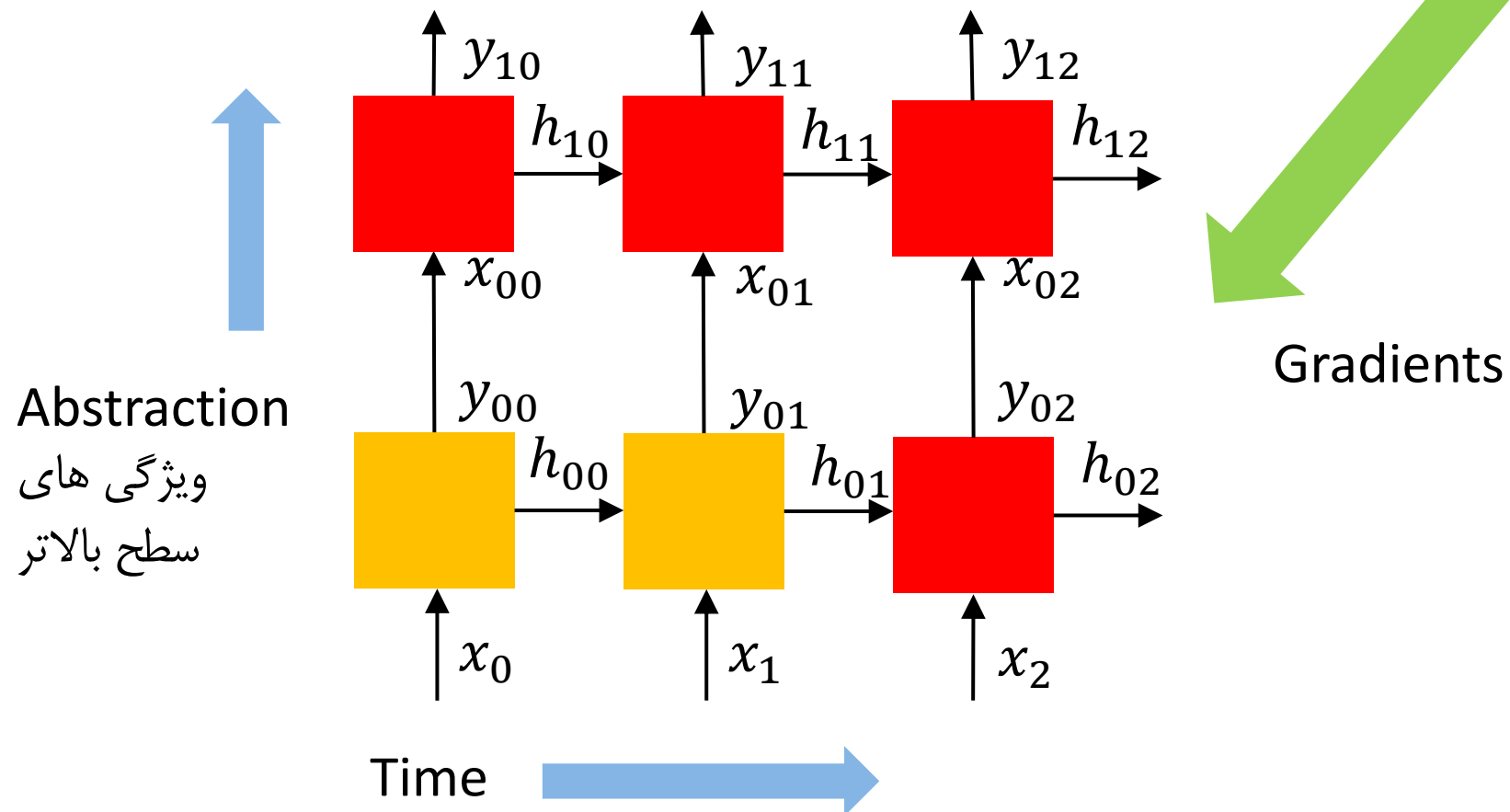
Backprop still works:



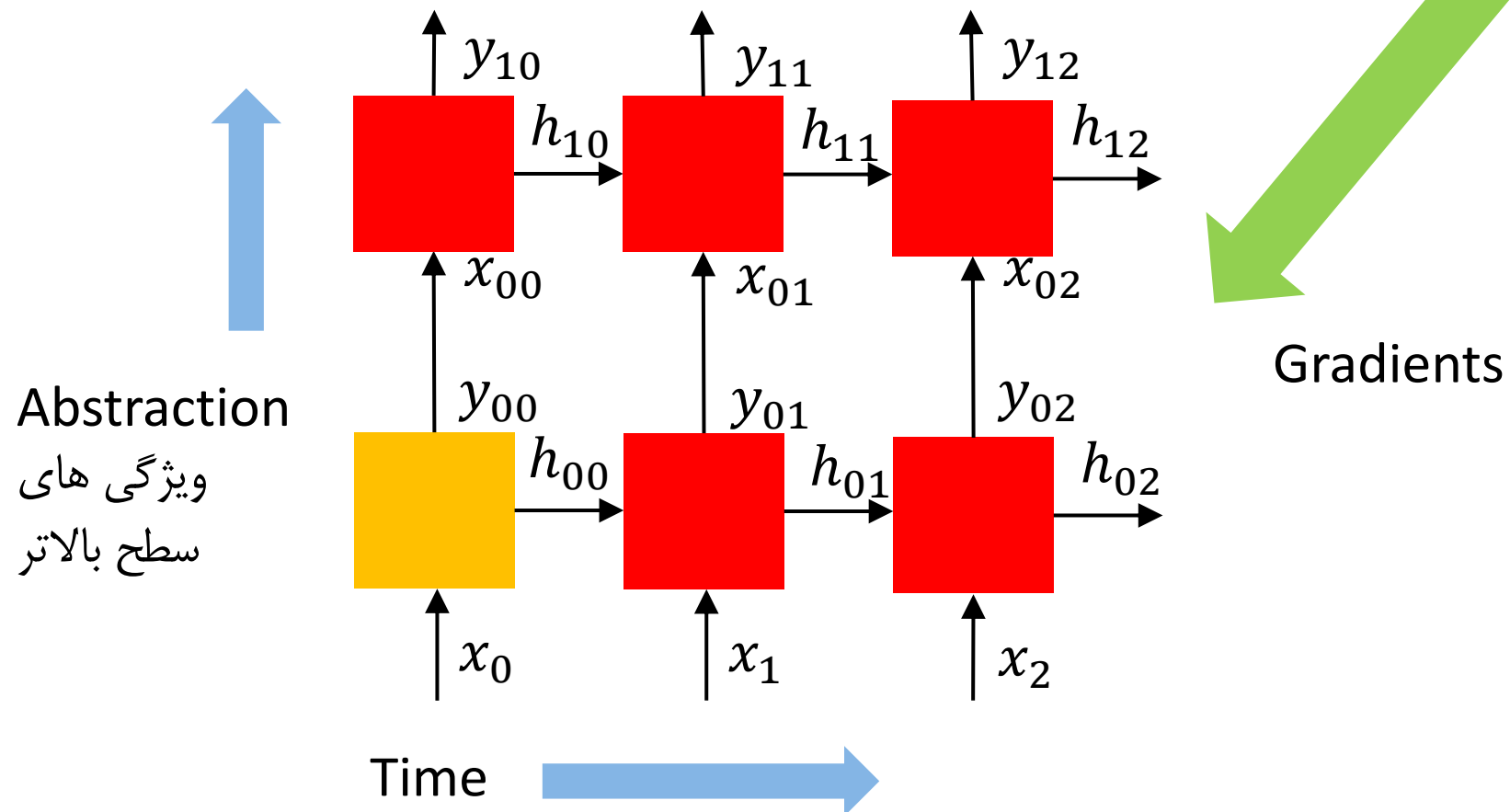
Backprop still works:



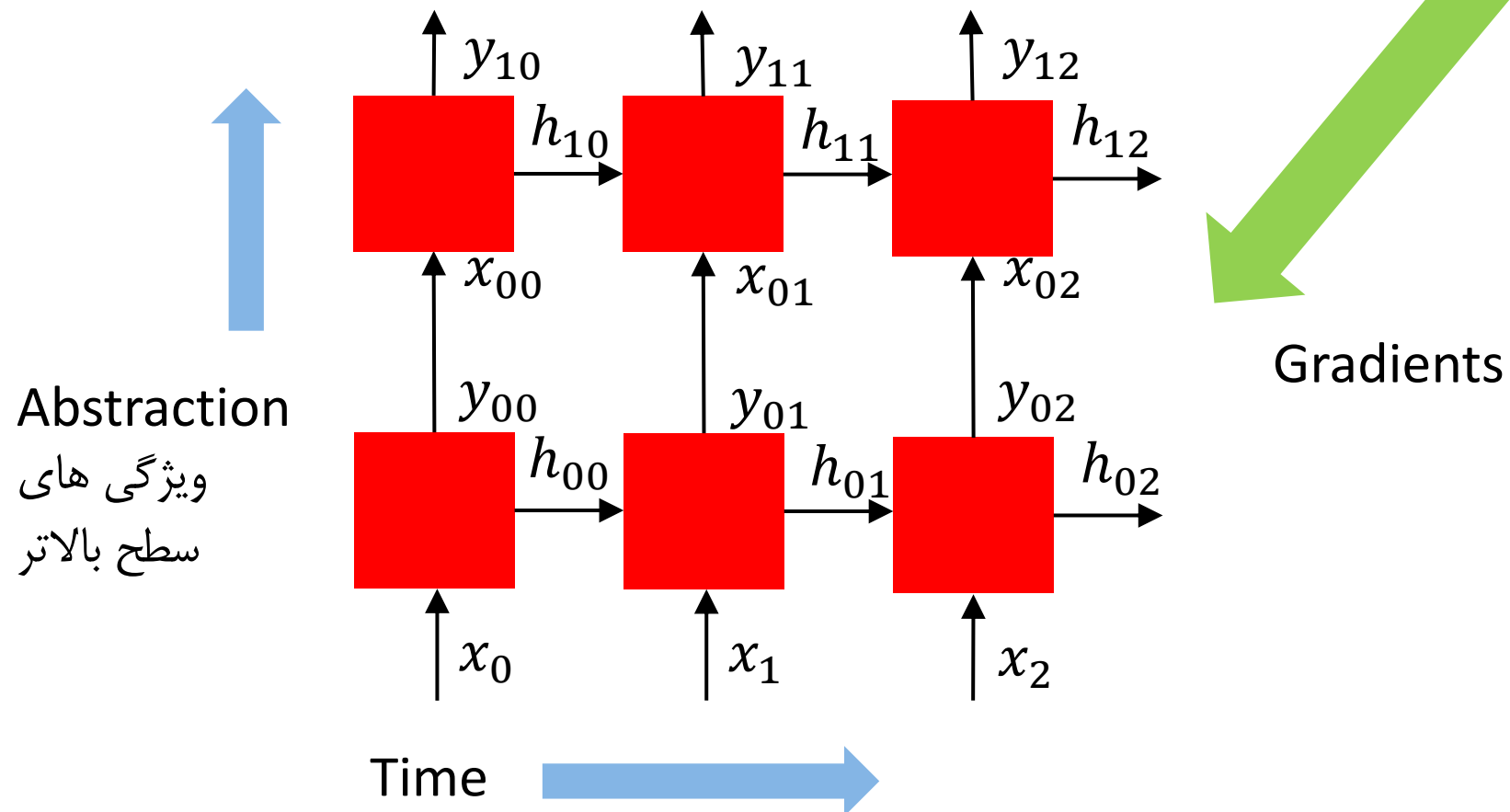
Backprop still works:



Backprop still works:

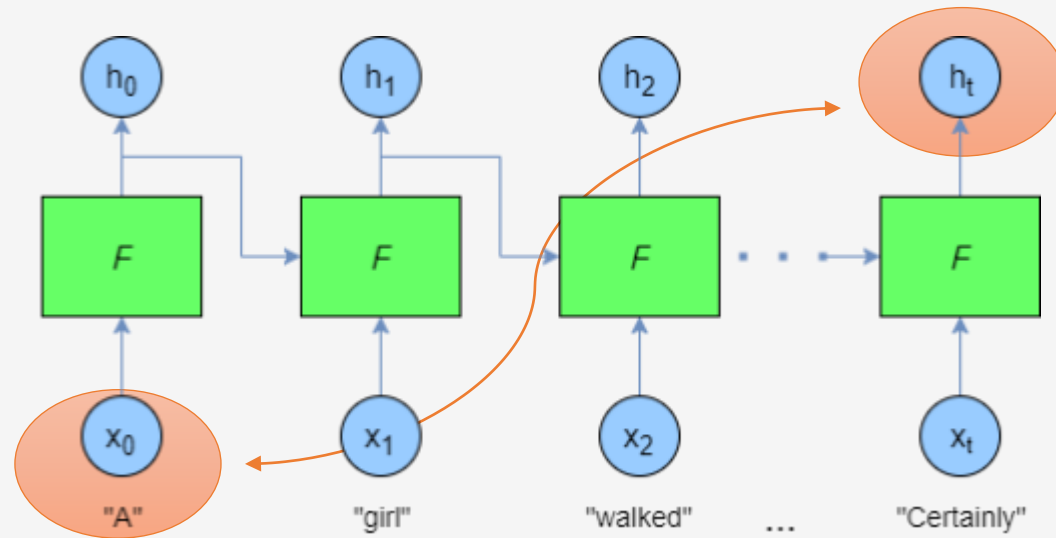


Backprop still works:



مشکلات RNN

- مفهوم و ایده پشت شبکه های بازگشت عالی است.
- مشکل عملیاتی هنگام آموزش وجود دارد.



ناپدید شدن گرادیان

• وابستگی بین h_t و x_0 به صورت زیر قابل تعریف است:

$$\frac{\partial h_t}{\partial x_0}$$

• می تواند به صورت زیر خلاصه شود:

$$\frac{\partial h_t}{\partial x_0} = \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} \frac{\partial h_{t-2}}{\partial h_{t-3}} \dots \frac{\partial h_0}{\partial x_0}$$

$$h_n = \tanh(\underbrace{w_{hh}h_{n-1} + w_{xh}x_n + b}_z)$$

$$h_n = \tanh(z)$$

$$\frac{\partial z}{\partial h_{n-1}} = w_{hh}$$

$$\frac{\partial h_n}{\partial z} = 1 - \tanh^2(z)$$

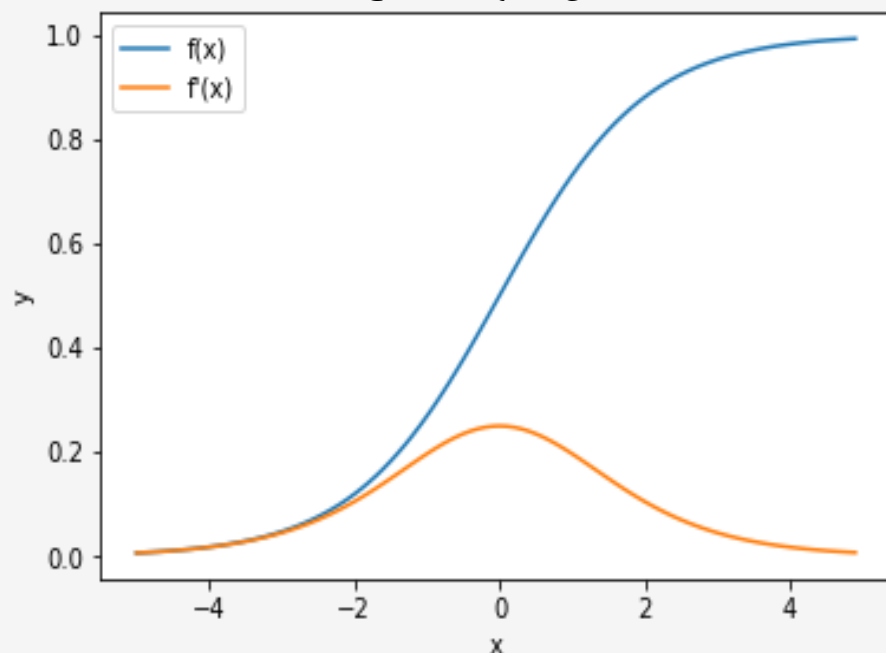
• مشتق \tanh عددی بین صفر و یک است.

• بخش اصلی حاصلضرب مکرر w_{hh} است. اگر بزرگترین مقدار ویژه w_{hh} برابر:

- ۱، سپس گرادیان منتشر می شود،
- < 1 ، نتیجه به صورت تصاعدی رشد می کند (منفجر می شود)،
- > 1 ، نتیجه به صورت تصاعدی کوچک می شود (ناپدید می شود).

مشکل گرادیان

نتیجه مشتق با توجه به تابع فعالسازی



• نتیجه

- مشتق جزئی نزدیک به صفر است.
- خروجی به ورودی های فاصله دور قبلی مربوط نیست.
- وابستگی های طولانی در حال از بین رفتن است.

$$\text{new weight} = \text{weight} - \text{learning rate} * \text{gradient}$$

$$2.0999 = 2.1$$

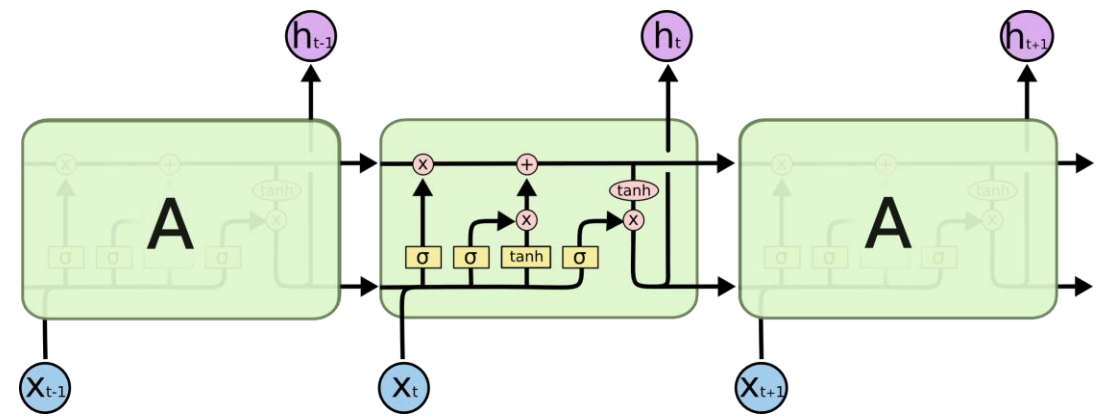
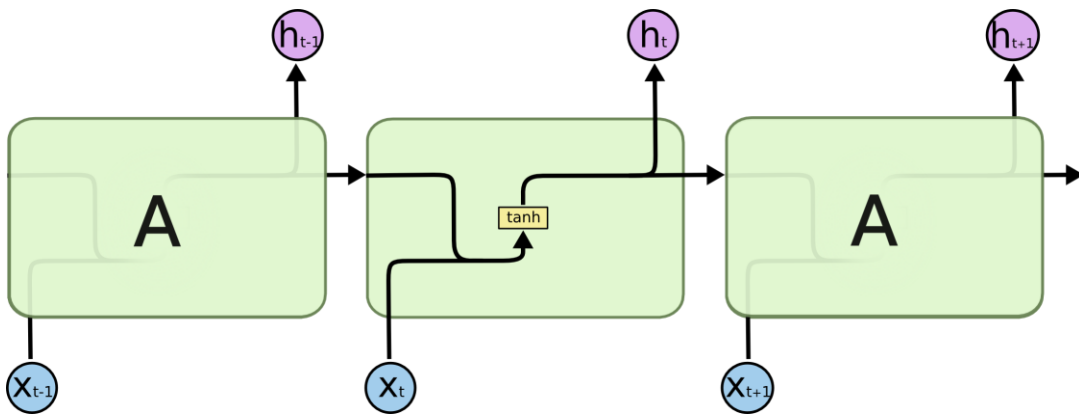
Not much of a difference

$$- 0.001$$

update value

• راه حل

- تغییر توابع فعال سازی
- تعداد کمی از توابع فعالسازی گرادیان دقیقا یک دارند:
- محدودیت زیاد
- مشکل را می توان تا حدودی کاهش داد اما برطرف نشد.
- ممکن است بر عملکرد یادگیری تأثیر بگذارد.
- پیشنهاد تغییر معماری: LSTM



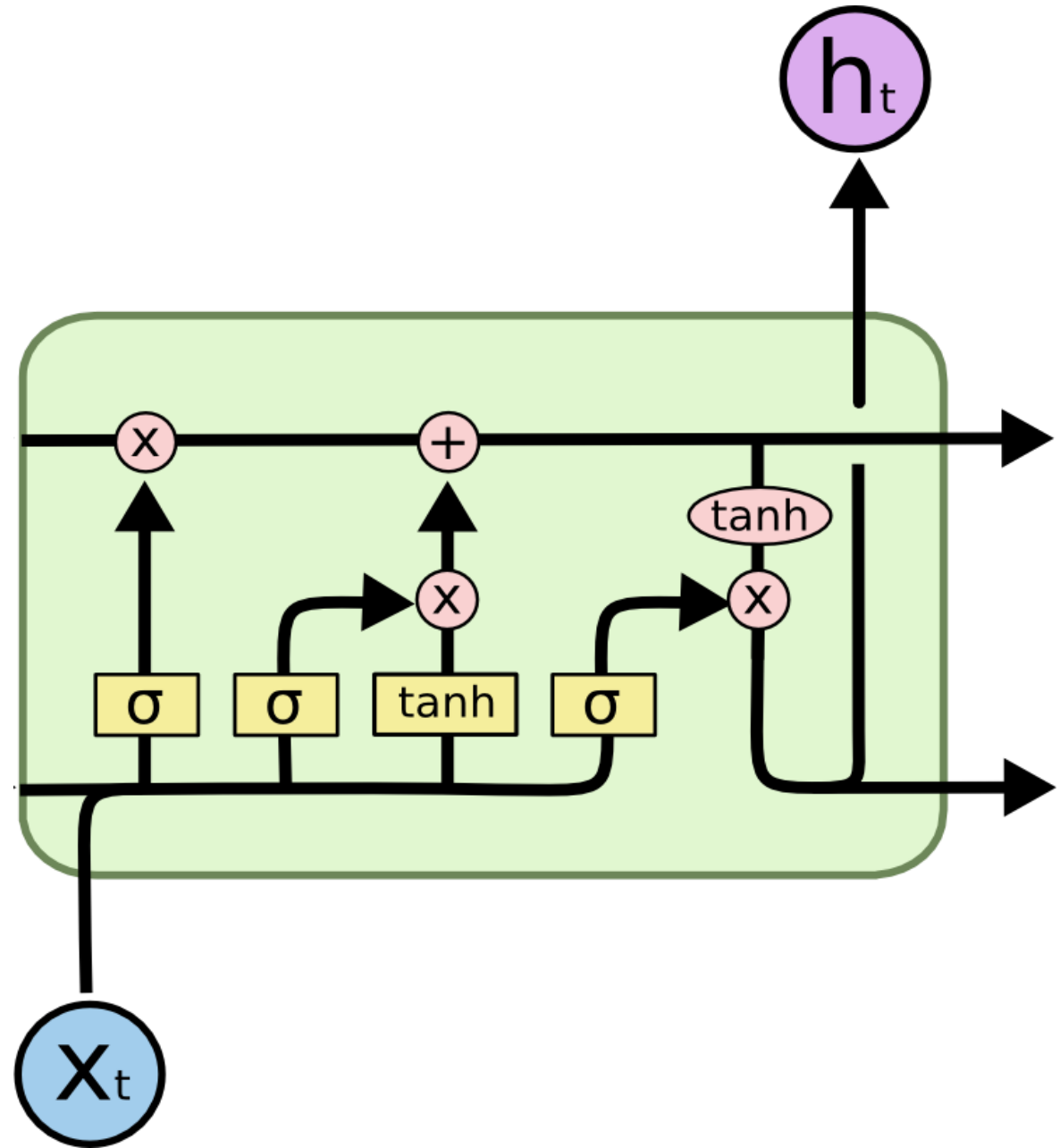
Long Short-Term Memory

ایده اصلی: ایجاد یک مسیر مستقیم از ورودی های قبلی به خروجی فعلی

حفظ اطلاعات مفید و حذف غیرمفیدها

LSTM

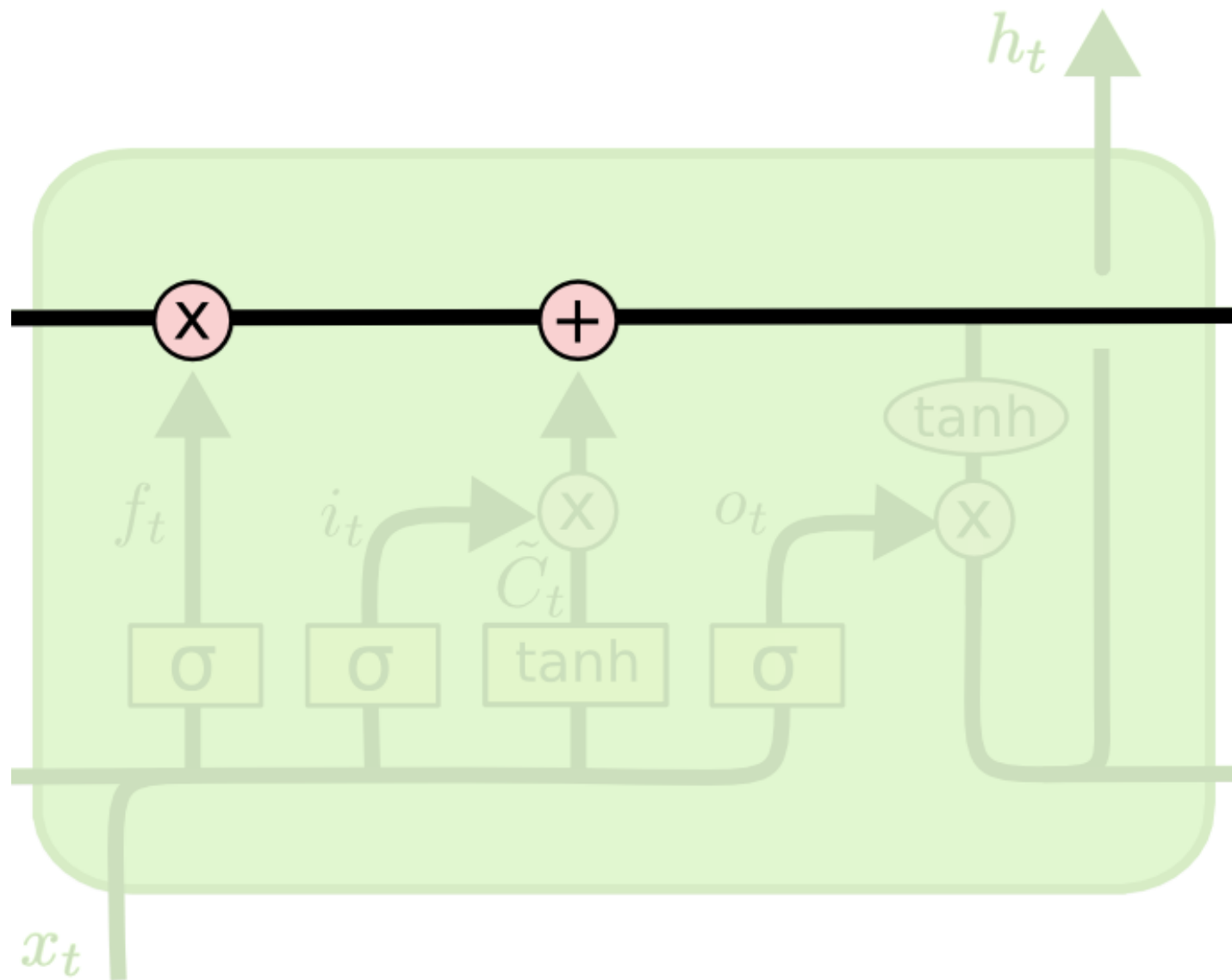
- Hidden states
- Input gates
- Forget gates
- Output gates

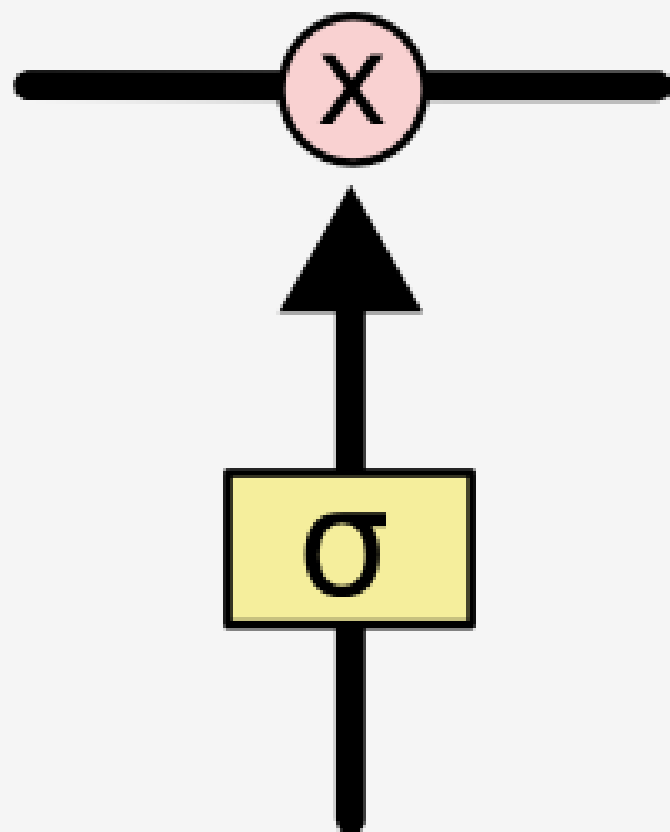


Hidden States

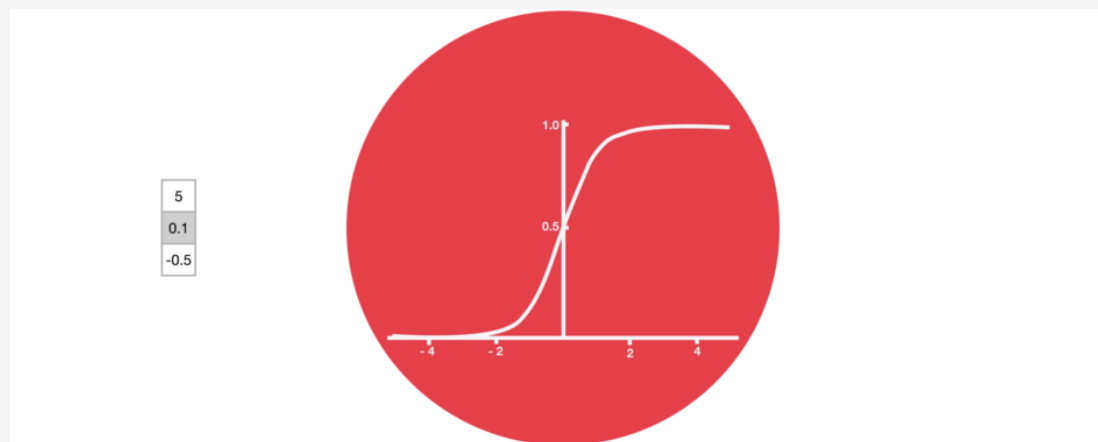
حمل اطلاعات قبلی

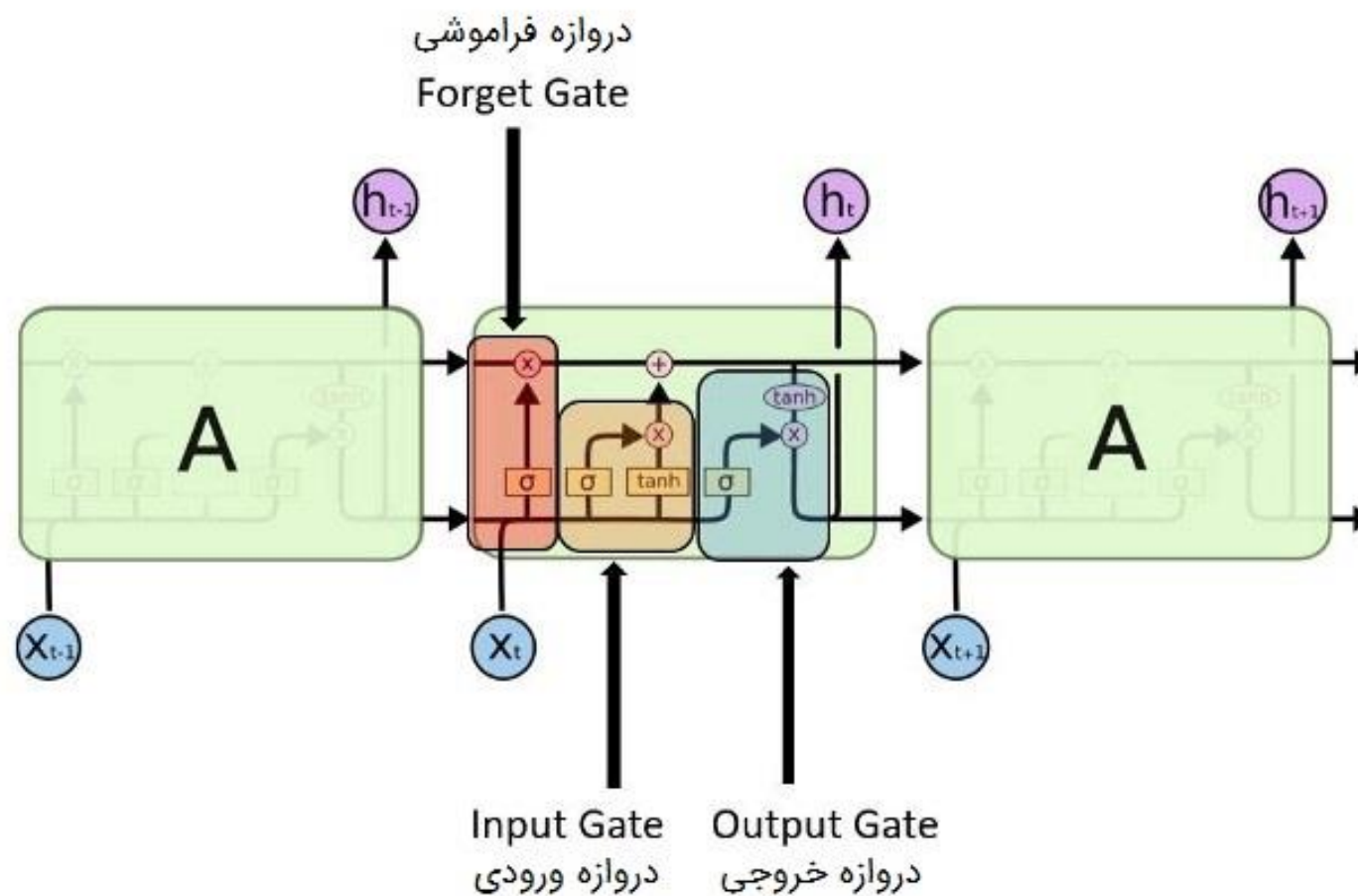
ساخت مسیر مستقیم

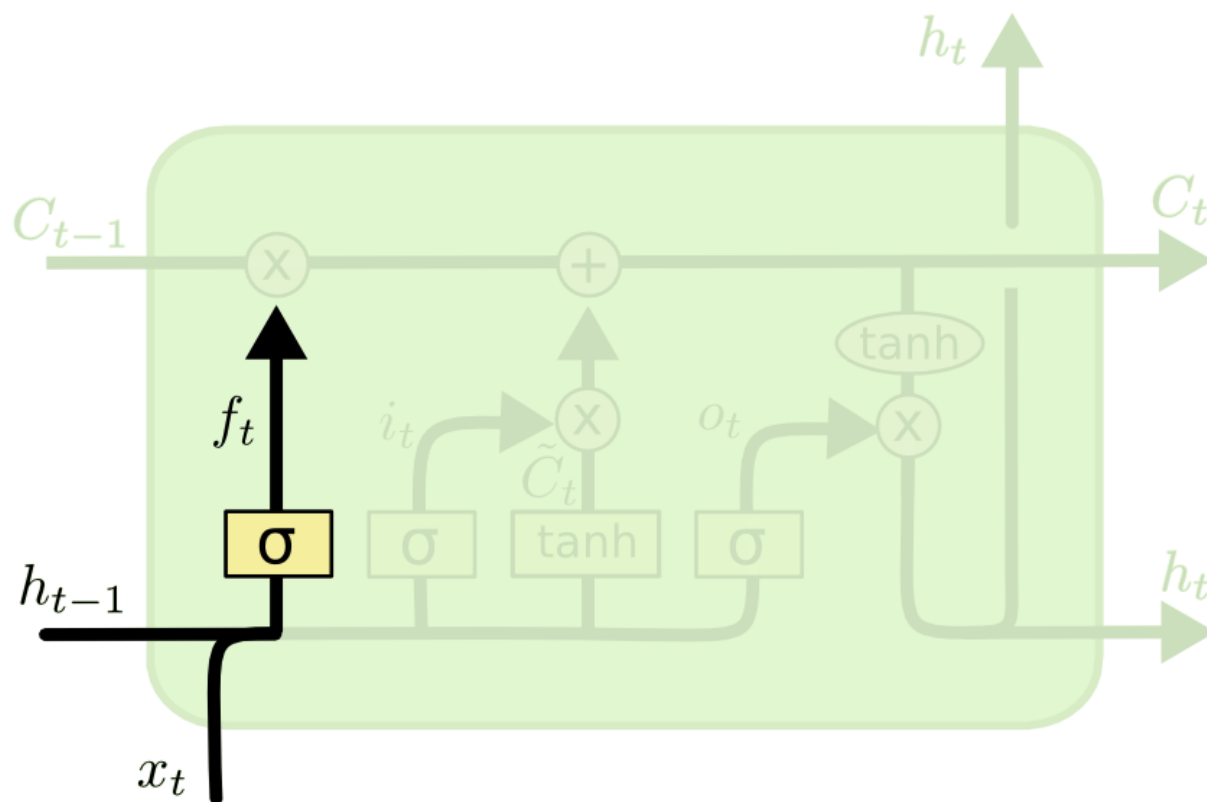




- یک شبکه عصبی سیگموئید
- رنج خروجی $[0,1]$
- ضرب نقطه ای
- تابع فیلترکننده

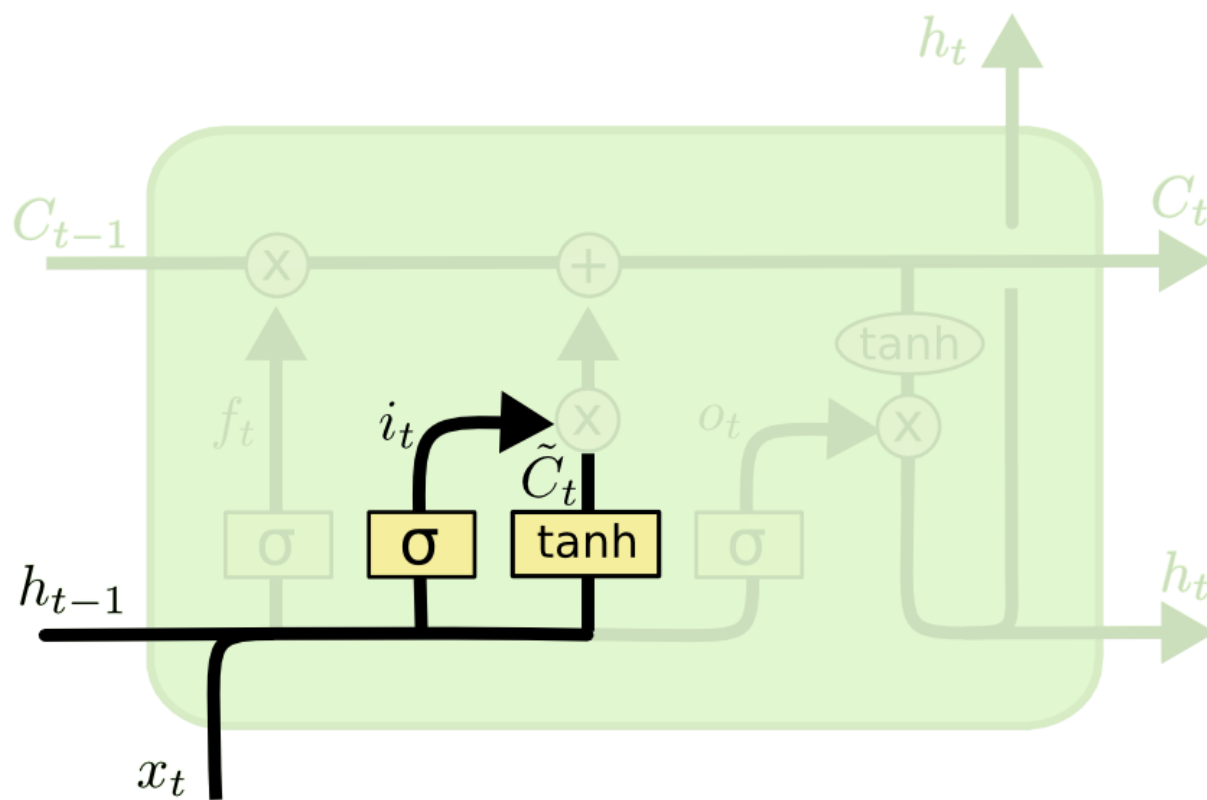






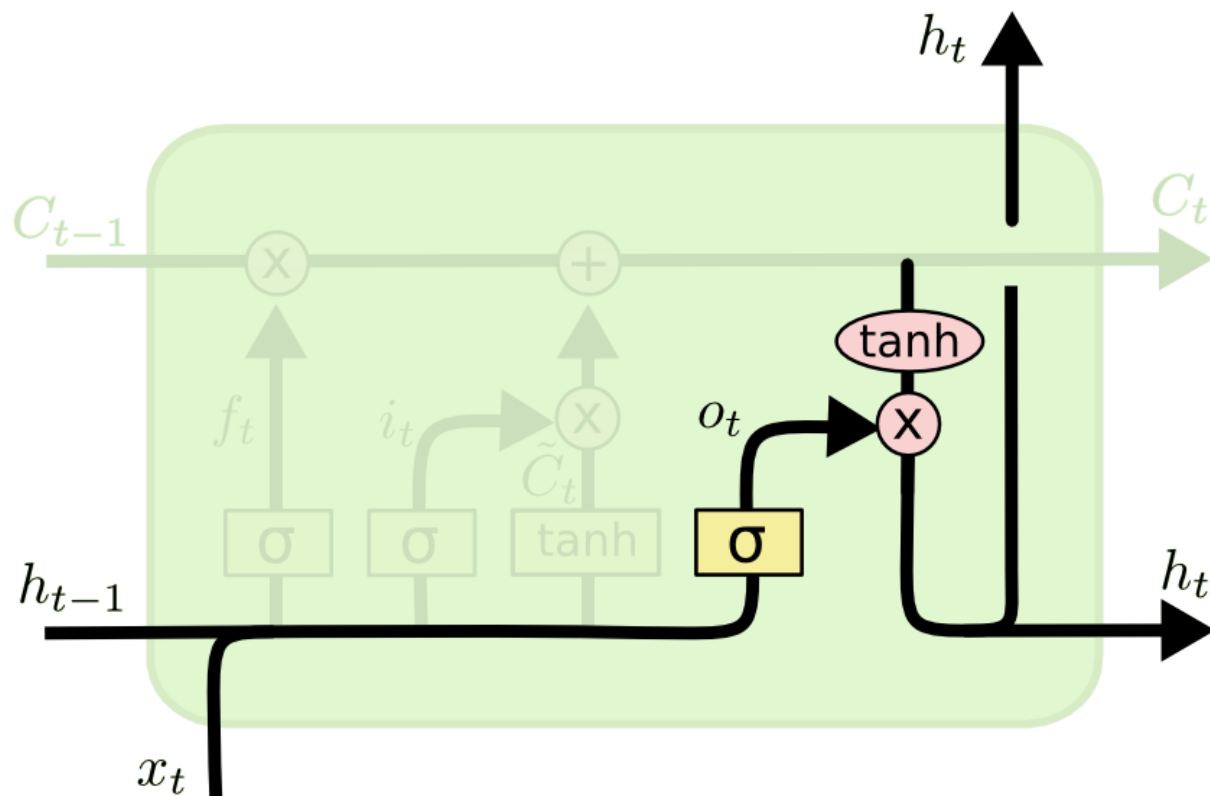
Input Gates

چیزهایی که از حالت های قبلی حفظ شده



Forget Gate

چیزی که به لایه هیدن اضافه می شود.



Output Gates

چیزی که به عنوان خروجی آماده می شود.



Neural Network
Layer



Pointwise
Operation



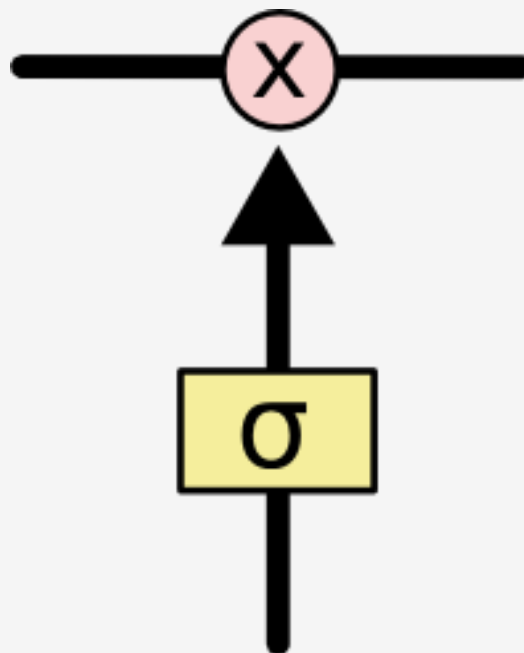
Vector
Transfer



Concatenate

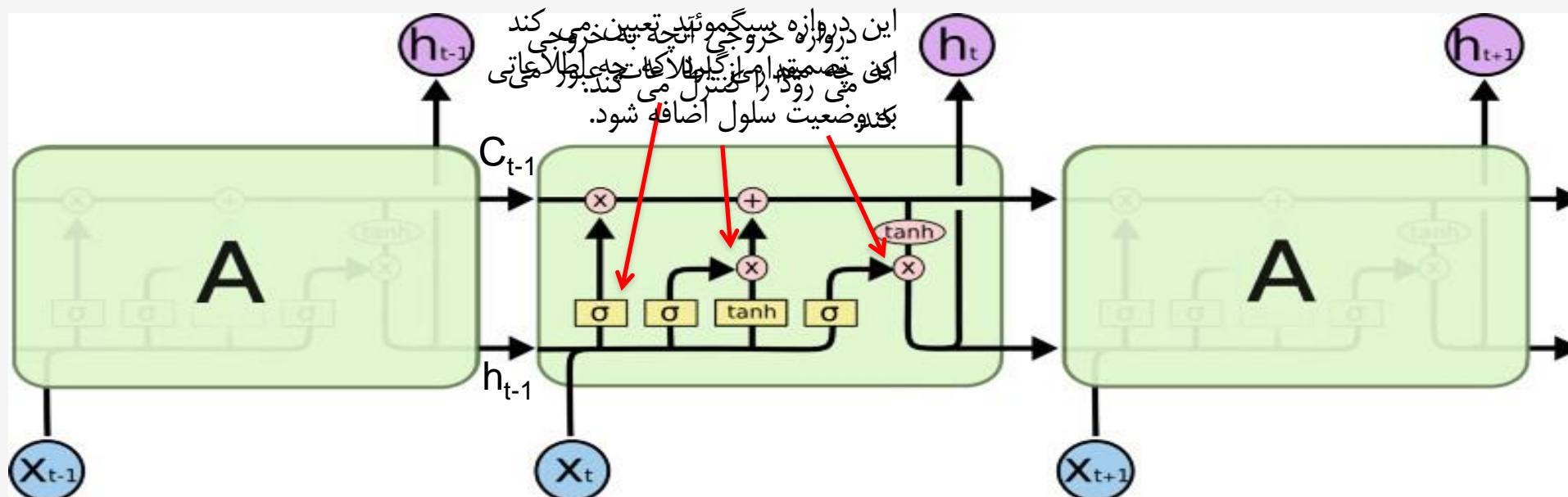


Copy

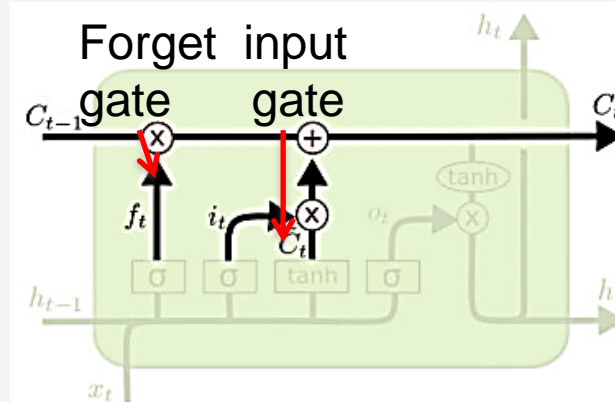


خروجی لایه سیگموئید با اعداد بین ۰-۱ تعیین می کند که هر جزء چقدر باید عبور کند. گیت X صورتی ضرب نقطه ای است.

LSTM

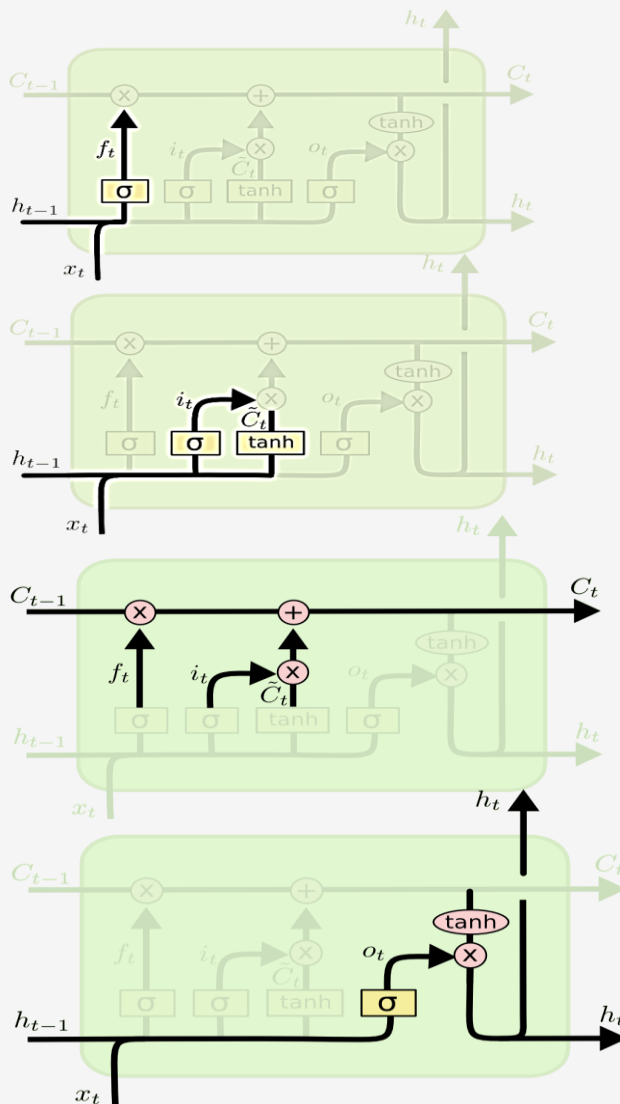


ایده اصلی این حالت سلولی C_t است، که به آرامی تغییر می‌کند، تنها با برهمکنش‌های خطی جزئی. جریان اطلاعات بدون تغییر در امتداد آن بسیار آسان است.



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

LSTM



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

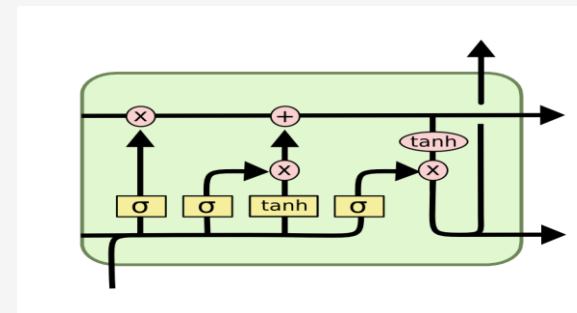
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

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

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

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

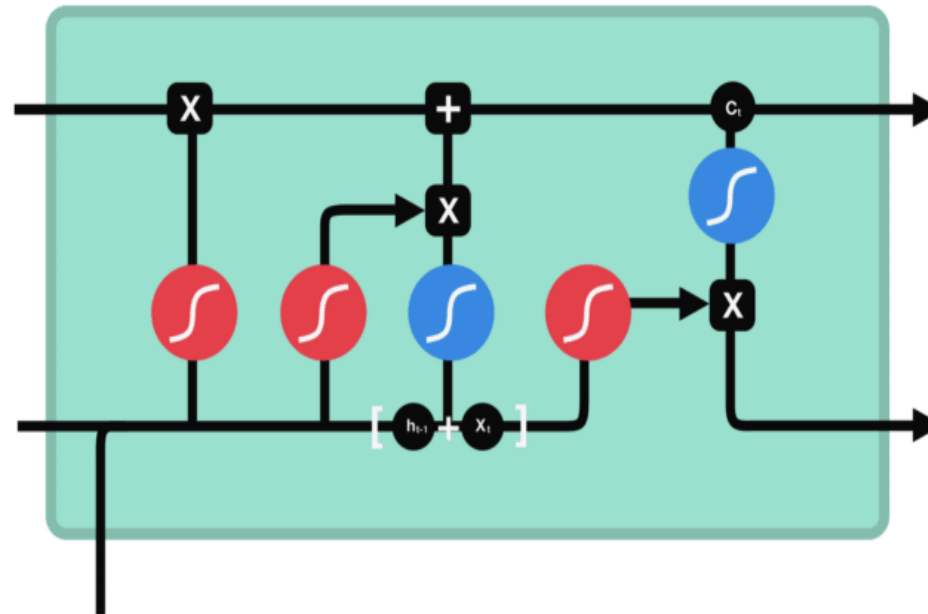


ا_t تصمیم می گیرد چه مولفه ای به روز شود.
C'_t تغییر را اعمال می کند.

به روزرسانی cell state

O_t تصمیم می گیرد چه چیزی به خروجی برود.

LSTM به صورت متحرک



- 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

Candidate input:

$$a_t = \tanh(W_a \cdot x_t + U_a \cdot h_{t-1} + b_a)$$

Input gate:

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

Forget gate:

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

Output gate:

$$o_t = \sigma(W_o \cdot x_t + U_o \cdot h_{t-1} + b_o)$$

Cell state:

$$C_t = a_t \odot i_t + f_t \odot C_{t-1}$$

Hidden state:

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

LSTM: مثال

Internal weights

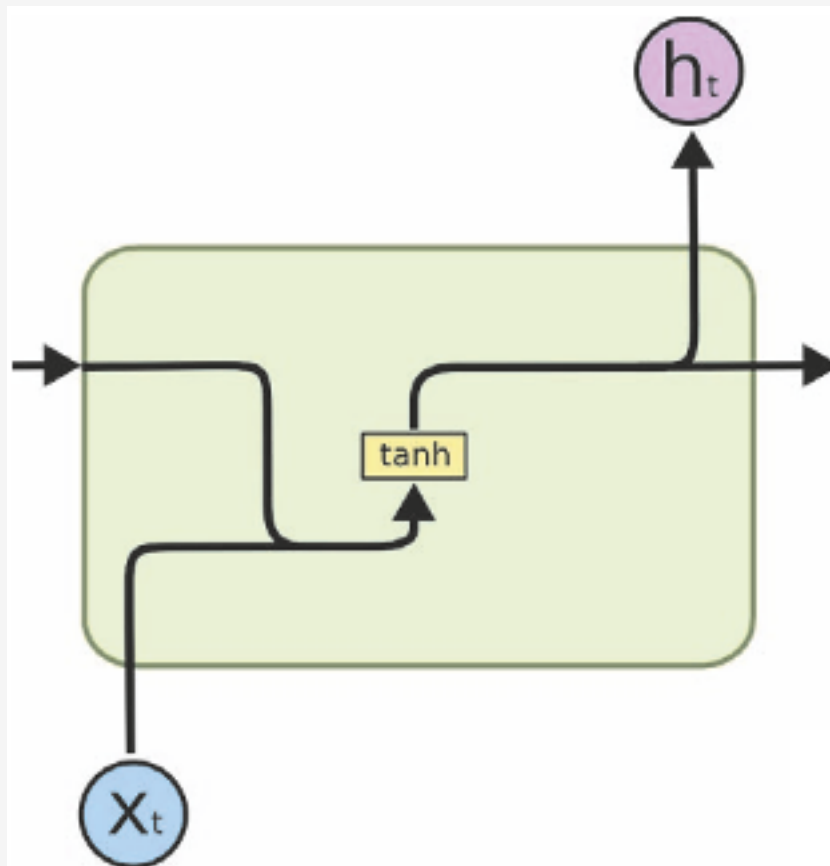
$W_a = \begin{bmatrix} 0.45 \\ 0.25 \end{bmatrix}$	$U_a = \begin{bmatrix} 0.15 \end{bmatrix}$	$b_a = \begin{bmatrix} 0.2 \end{bmatrix}$
$W_i = \begin{bmatrix} 0.95 \\ 0.8 \end{bmatrix}$	$U_i = \begin{bmatrix} 0.8 \end{bmatrix}$	$b_i = \begin{bmatrix} 0.65 \end{bmatrix}$
$W_f = \begin{bmatrix} 0.7 \\ 0.45 \end{bmatrix}$	$U_f = \begin{bmatrix} 0.1 \end{bmatrix}$	$b_f = \begin{bmatrix} 0.15 \end{bmatrix}$
$W_o = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}$	$U_o = \begin{bmatrix} 0.25 \end{bmatrix}$	$b_o = \begin{bmatrix} 0.1 \end{bmatrix}$

Input data

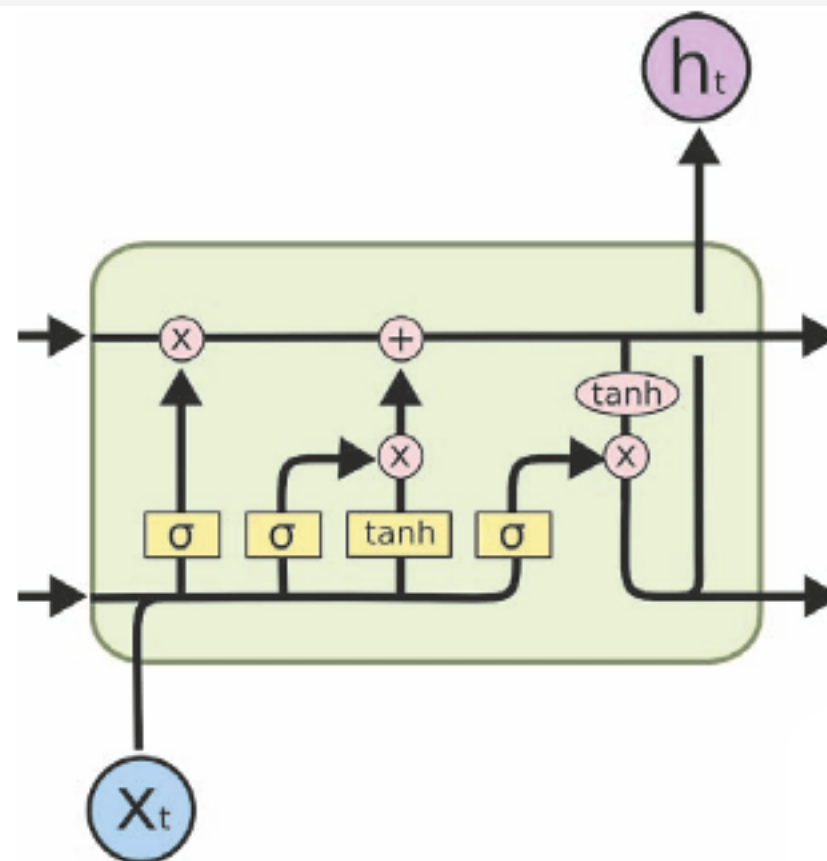
$x_0 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$	With label: 0.5
$x_1 = \begin{bmatrix} 0.5 \\ 3 \end{bmatrix}$	With label: 1.25

Forget gate:

$$f_t = \sigma(W_f \cdot x_t + U_f \cdot h_{t-1} + b_f) = \sigma\left(\begin{bmatrix} 0.7 & 0.45 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} + \begin{bmatrix} 0.1 \end{bmatrix} \begin{bmatrix} 0 \end{bmatrix} + \begin{bmatrix} 0.15 \end{bmatrix}\right) = 0.85195$$



(a) RNN



(b) LSTM

1

RNN برای مسئله
توالی / سری زمانی
استفاده می شود.

2

مشکل ناپدید شدن
گرادیان اثر بخشی آن را
محدود می کند.

3

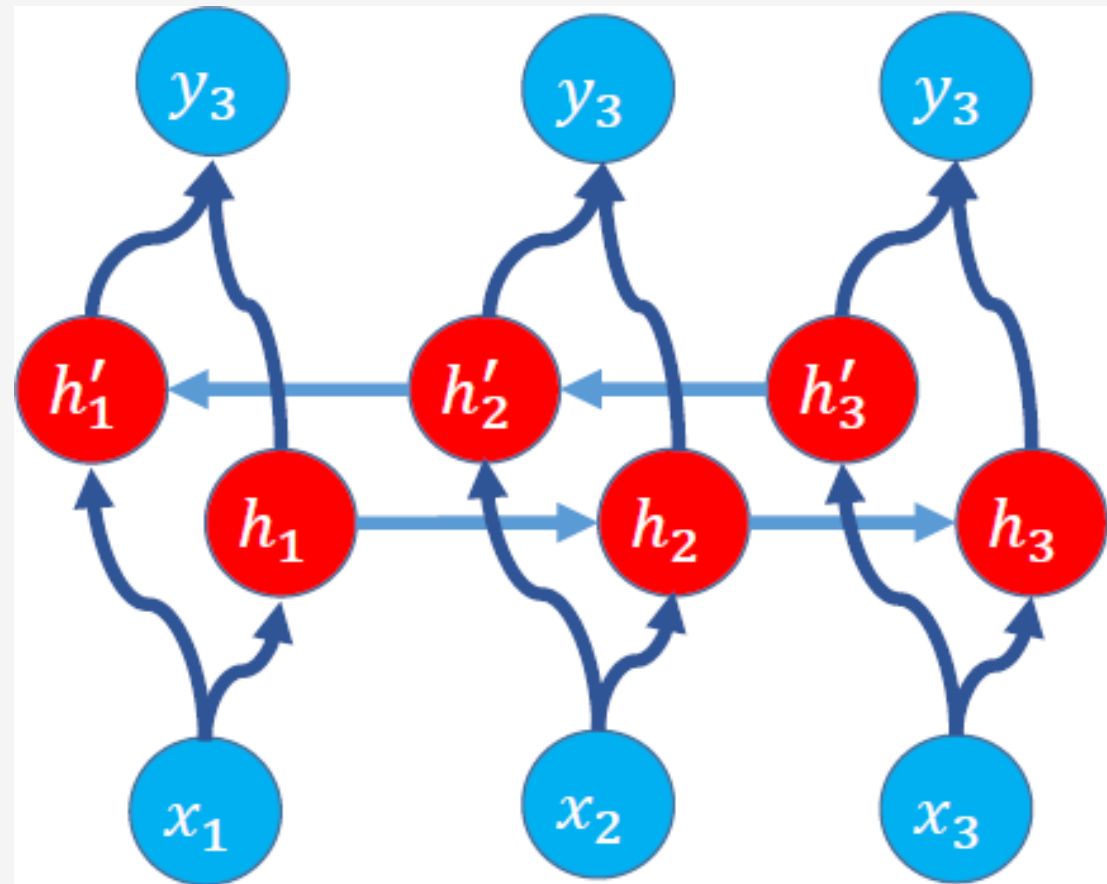
معرفی سلول LSTM
برای از بین بردن
مشکل

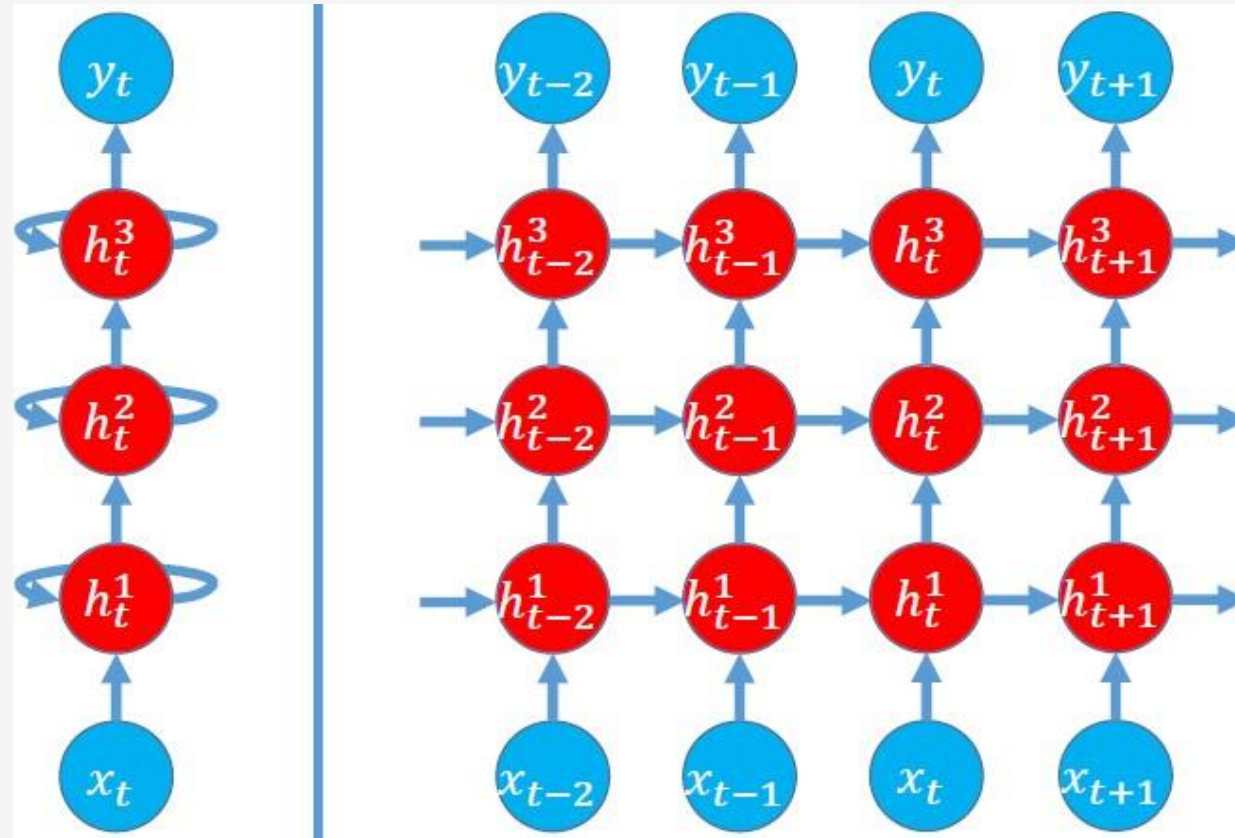
4

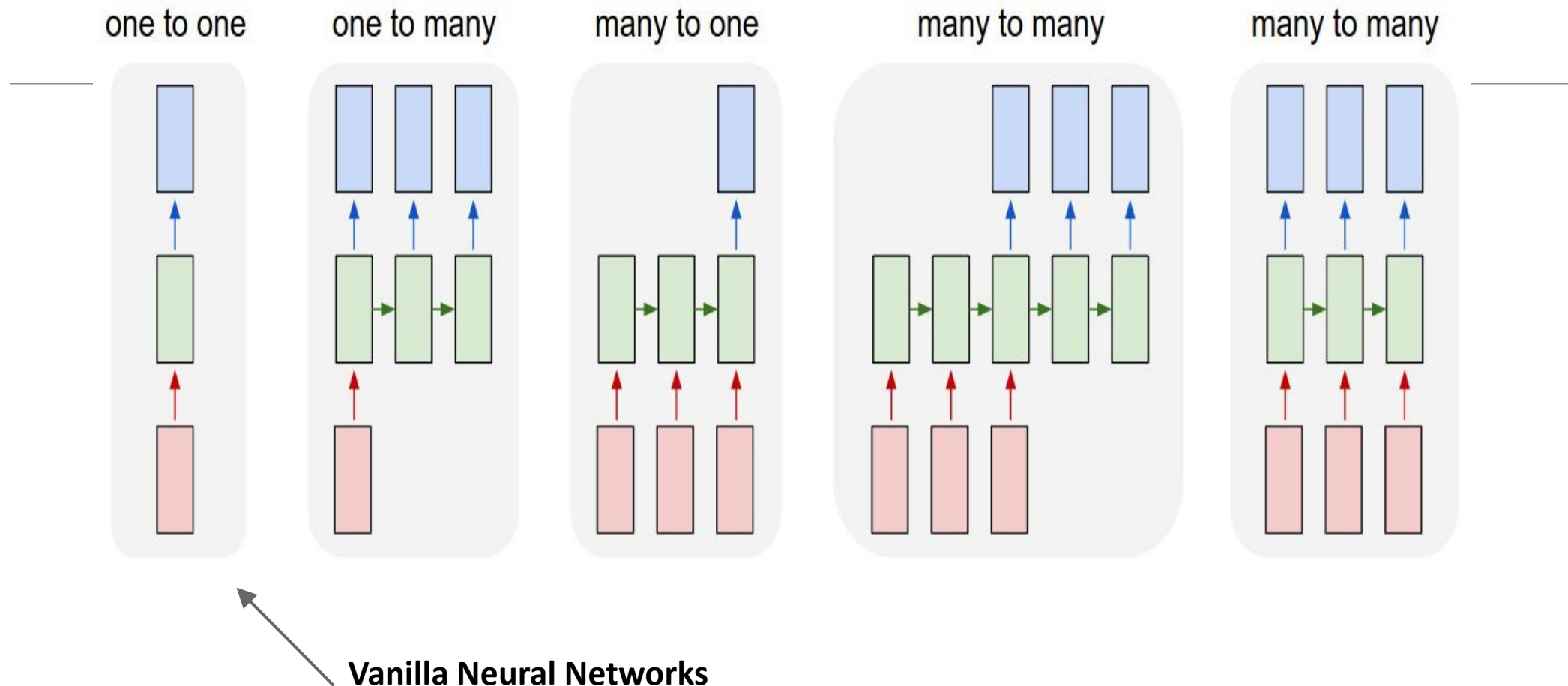
عملای LSTM بسیار
بهتر از RNN کار می
کند.

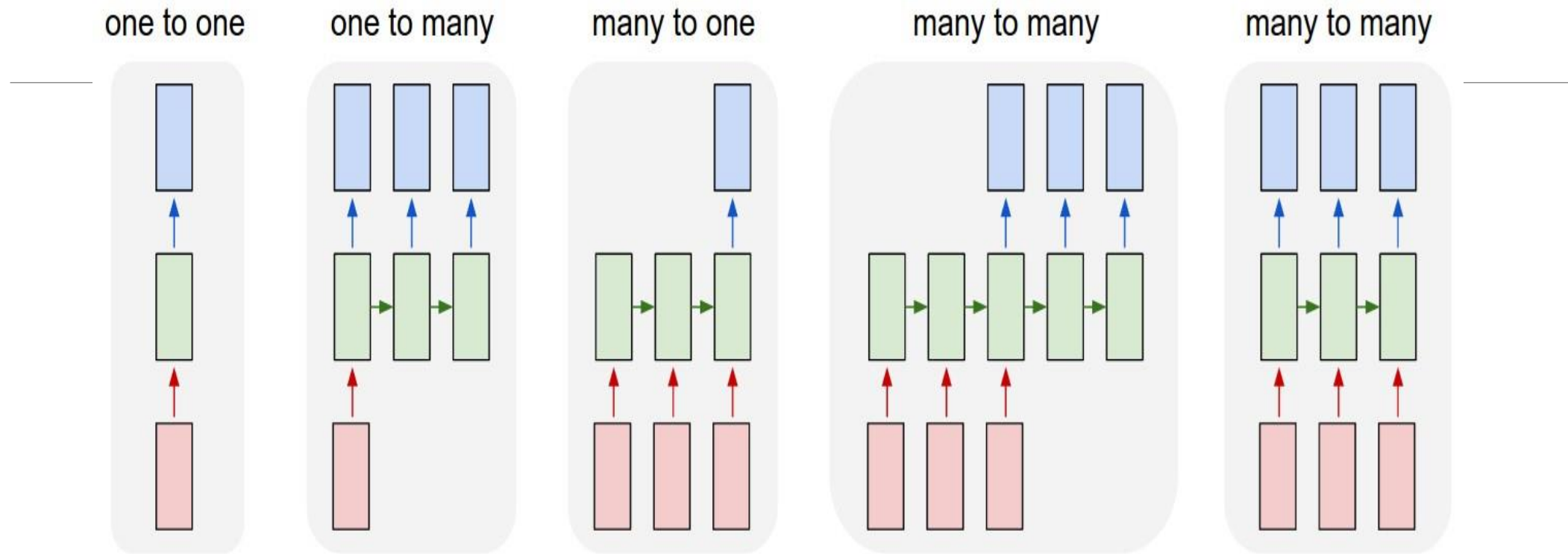
نتیجه

$$y = \text{softmax}(w[h_{t-1}, h'_{t-1}] + b)$$

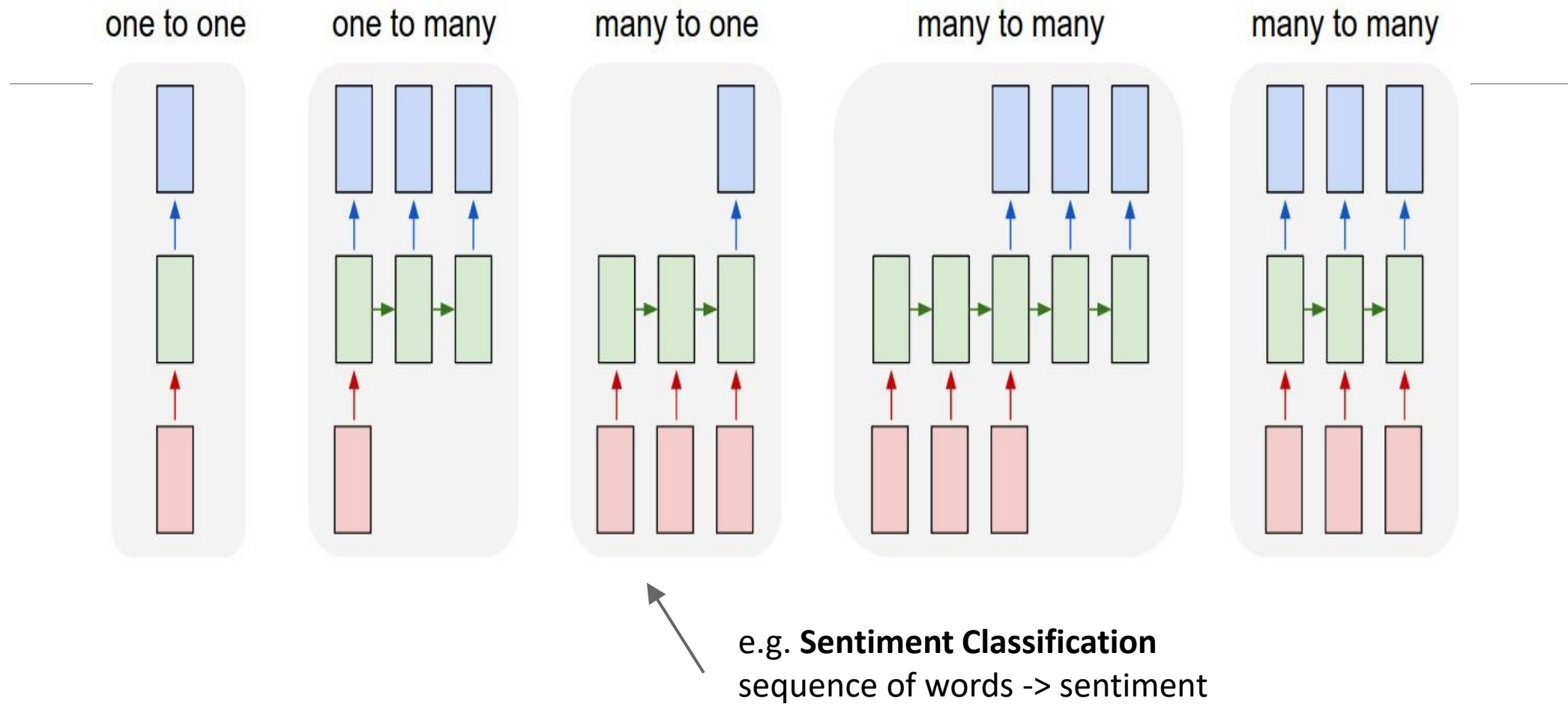




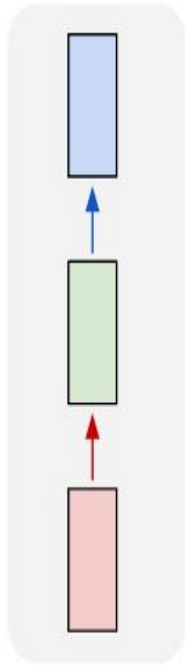




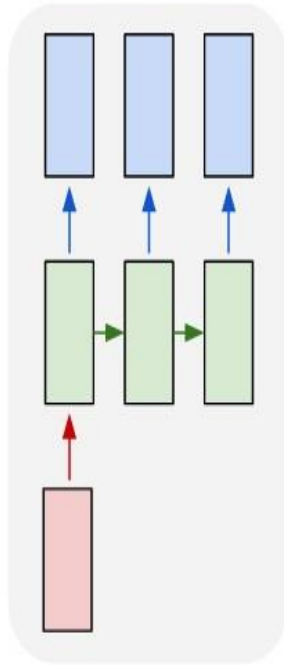
↖ e.g. **Image Captioning**
image -> sequence of words



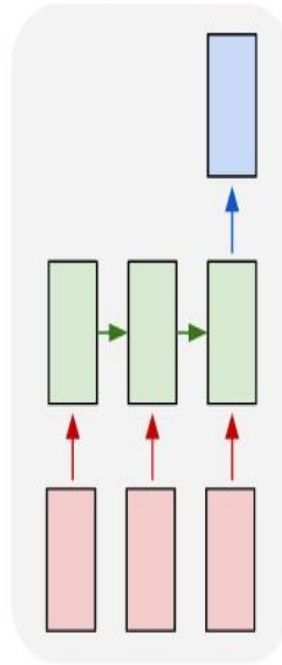
one to one



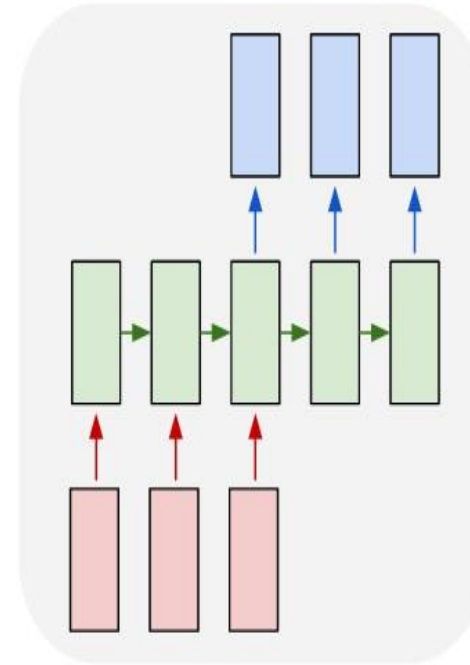
one to many



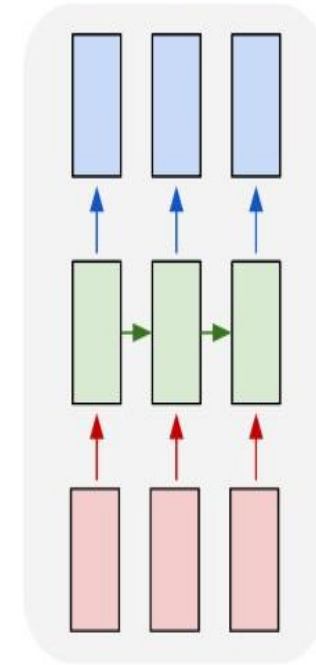
many to one



many to many



many to many



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

