

AI - FOUNDATION AND APPLICATION

Instructor:

Assoc. Prof. Dr. Truong Ngoc Son

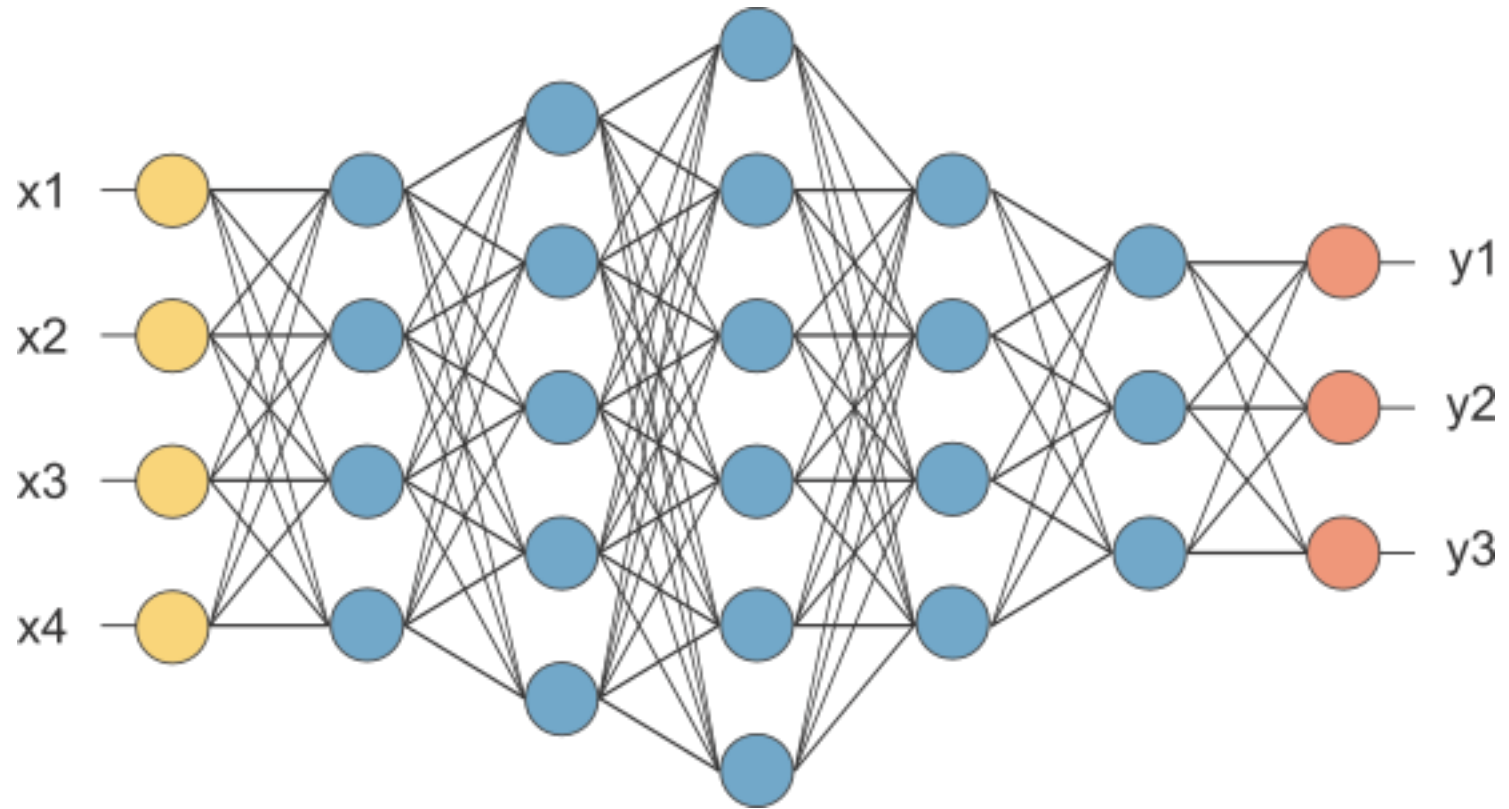
Chapter 5

Recurrent Neural Network

Outline



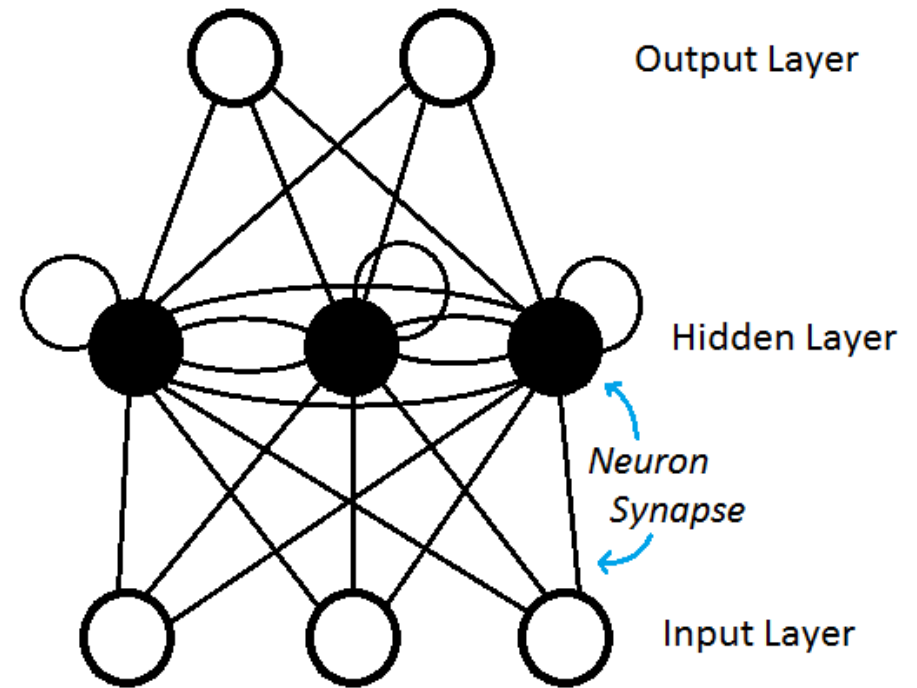
Feed Forward Neural Network



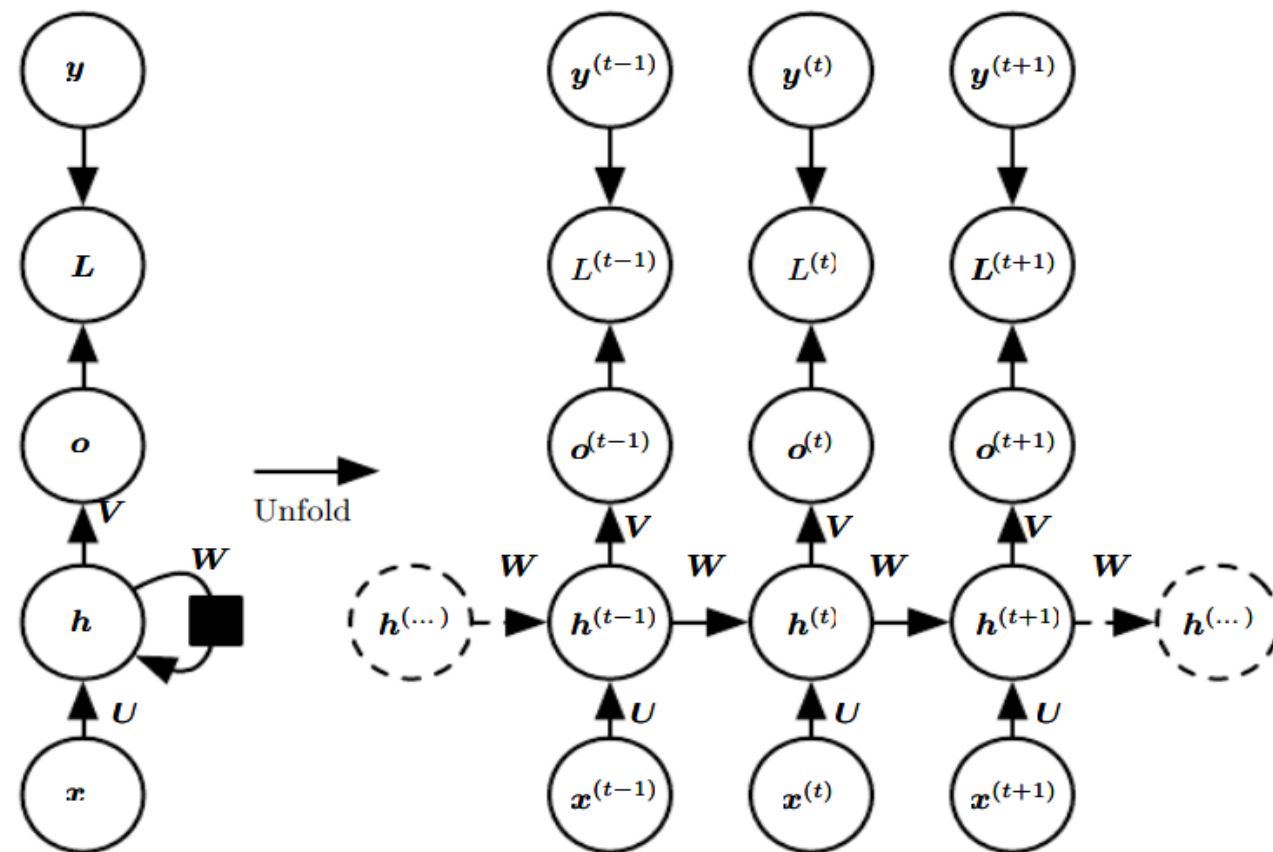
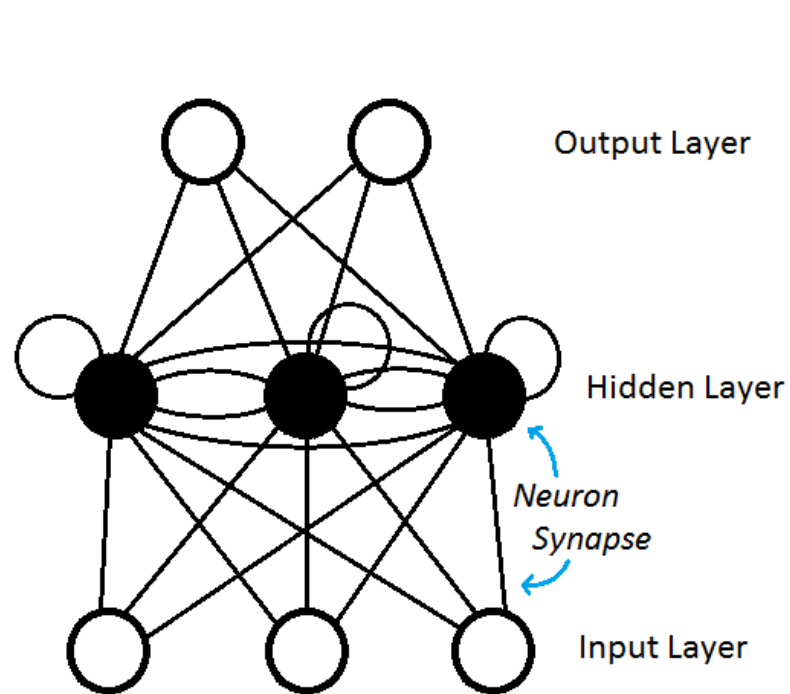
This is our fully connected network. If $x_1 \dots x_n$, n is very large and growing, this network would become too large. We now will input **one x_i at a time**, and **re-use the same edge weights**.

Sequence model – Recurrent Neural Network

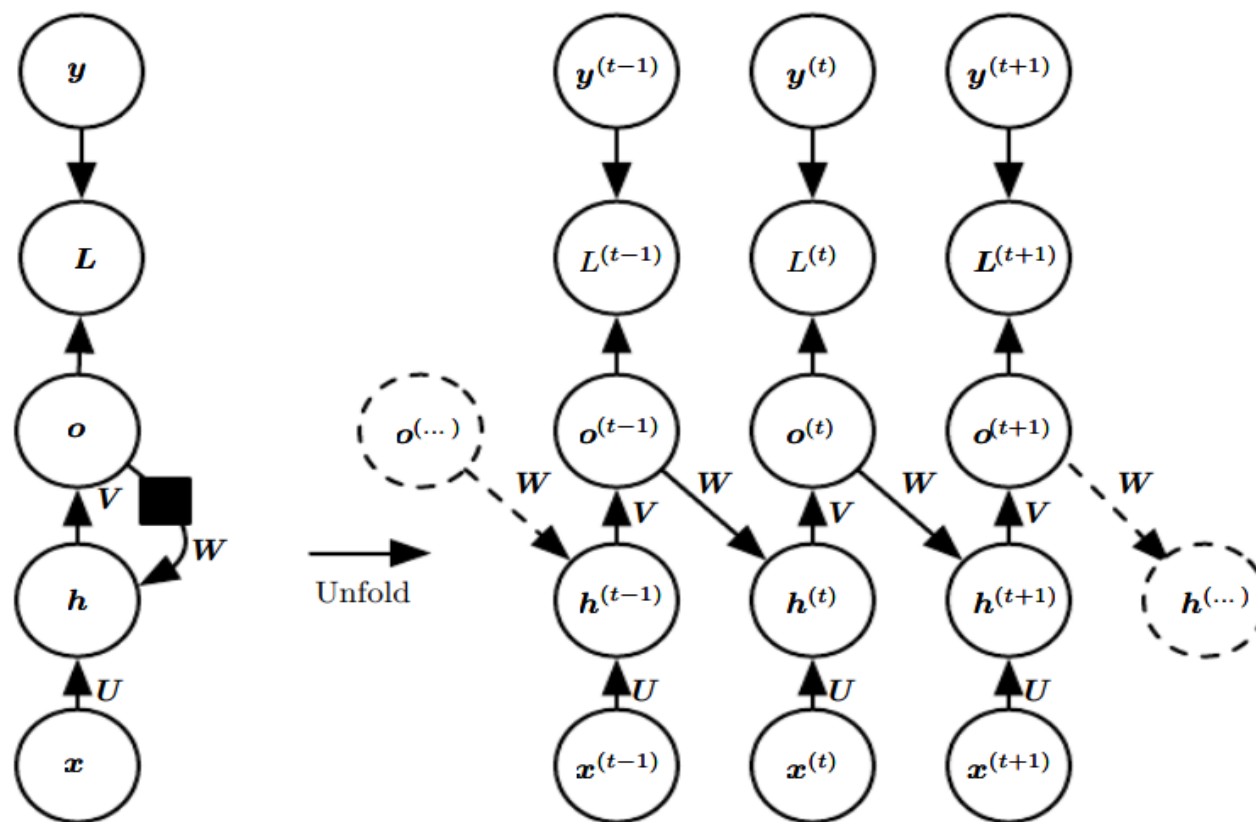
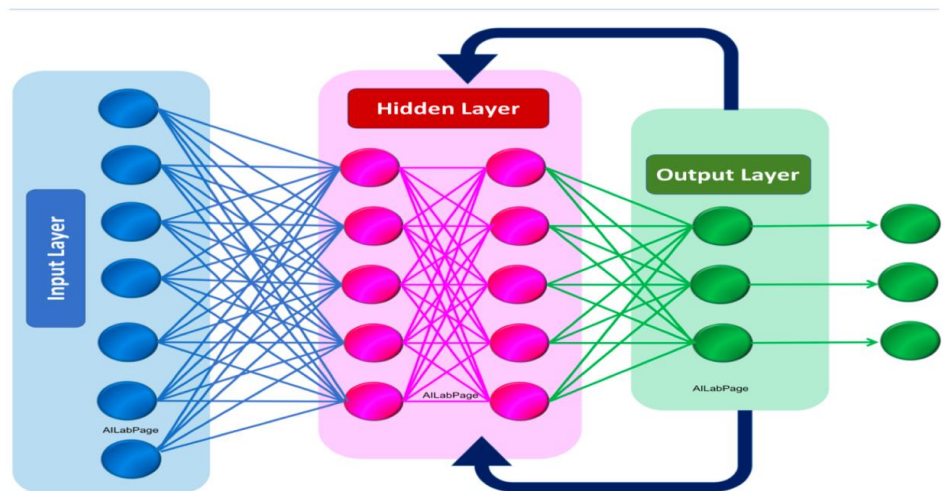
- Speech recognition
- Music generation
- Sentiment classification
- DNA sequence analysis
- Machine translate
- Video activity recognition
- Name entity recognition



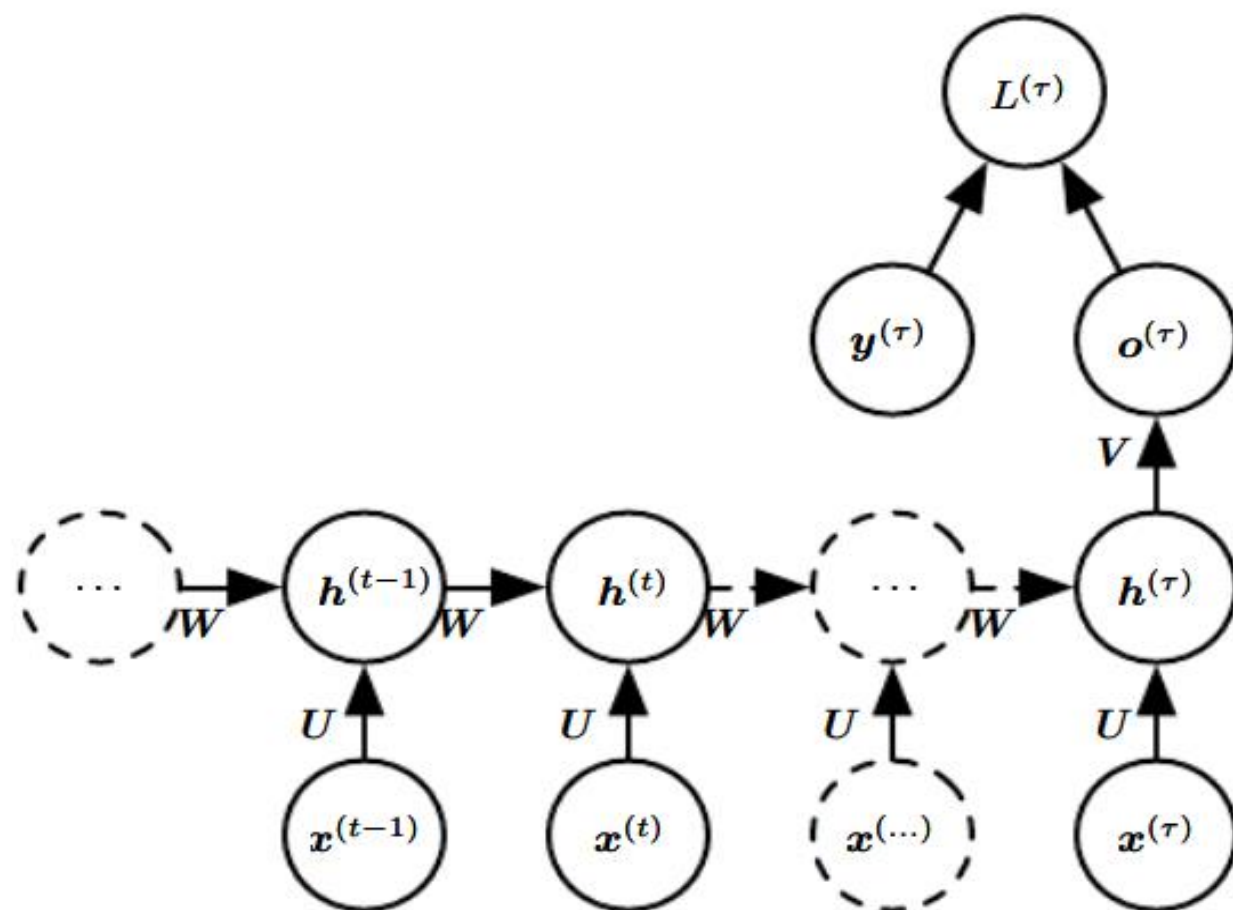
Recurrent Neural Network



Recurrent Neural Network

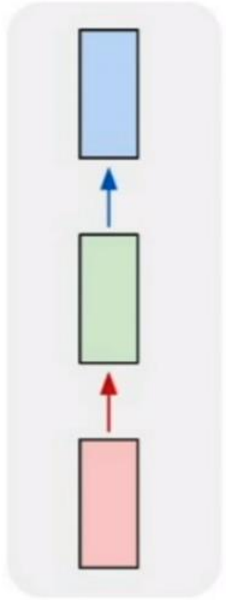


Recurrent Neural Network

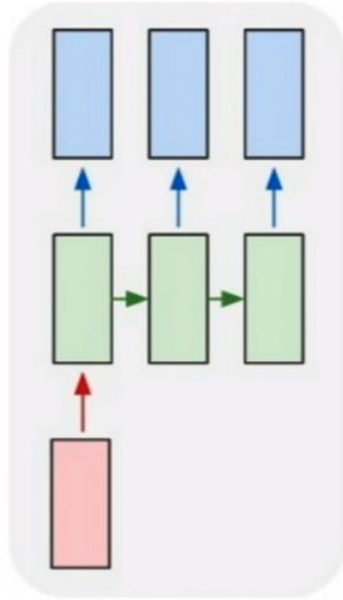


Types of Recurrent Neural Network

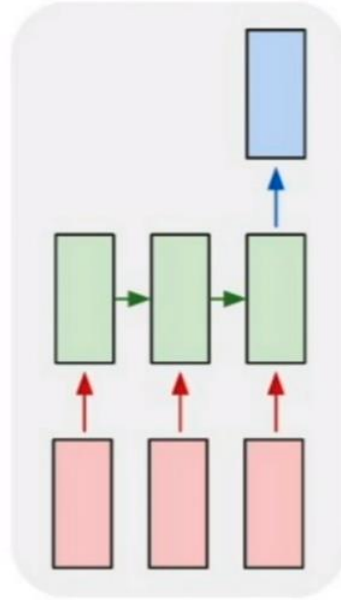
one to one



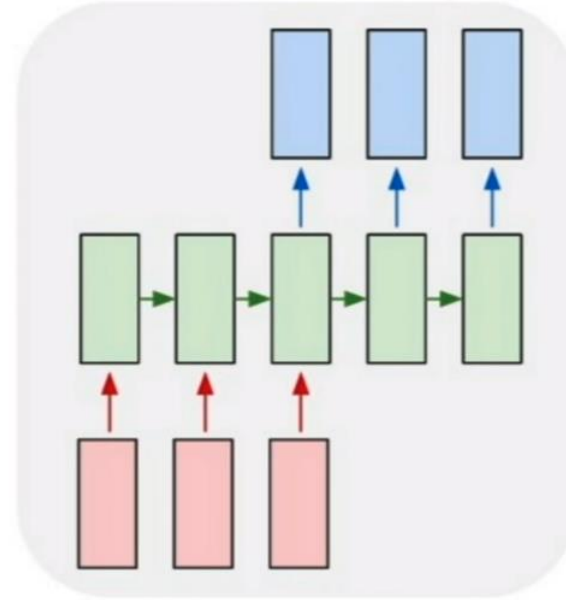
one to many



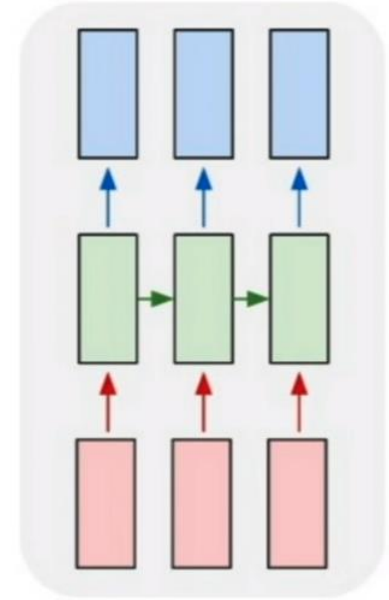
many to one



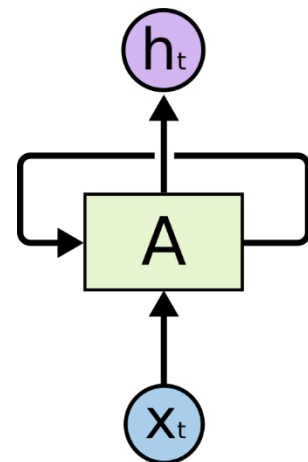
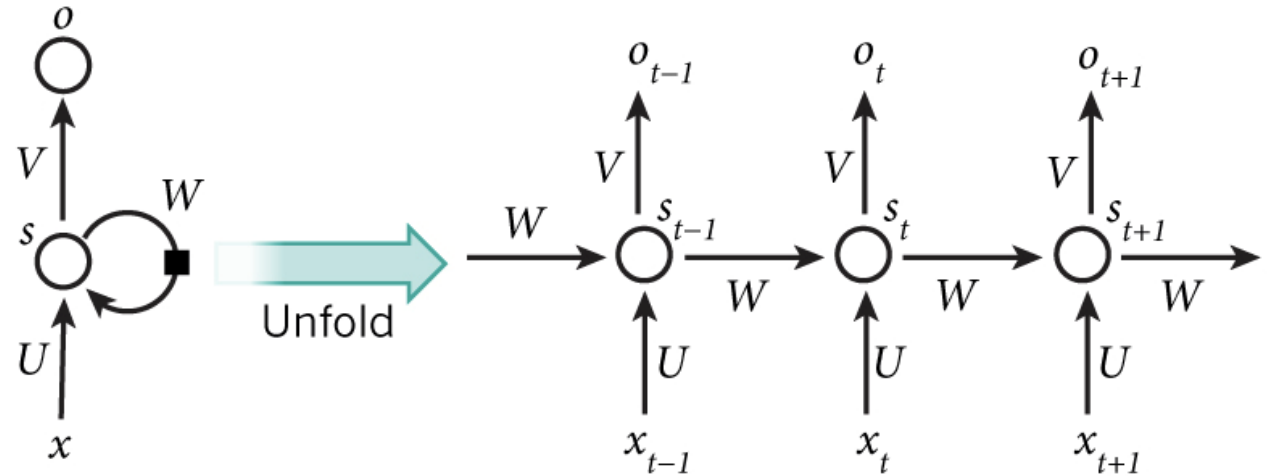
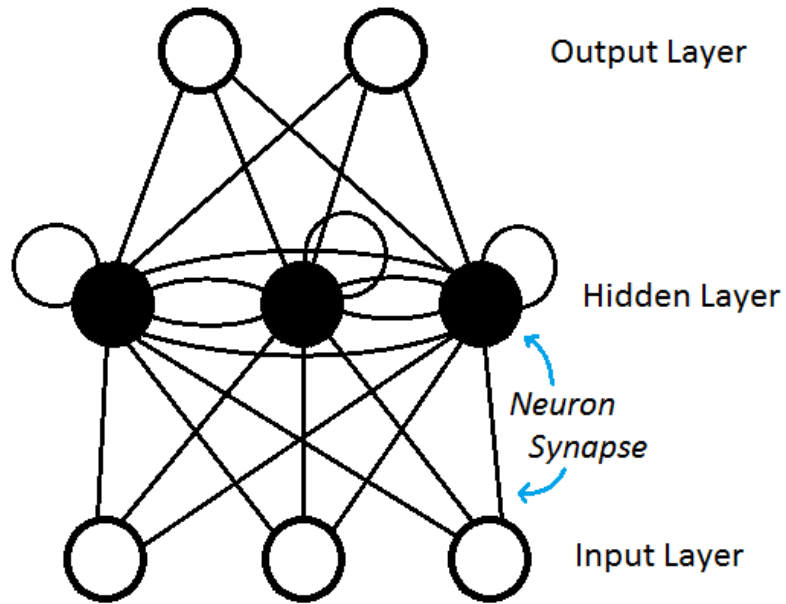
many to many



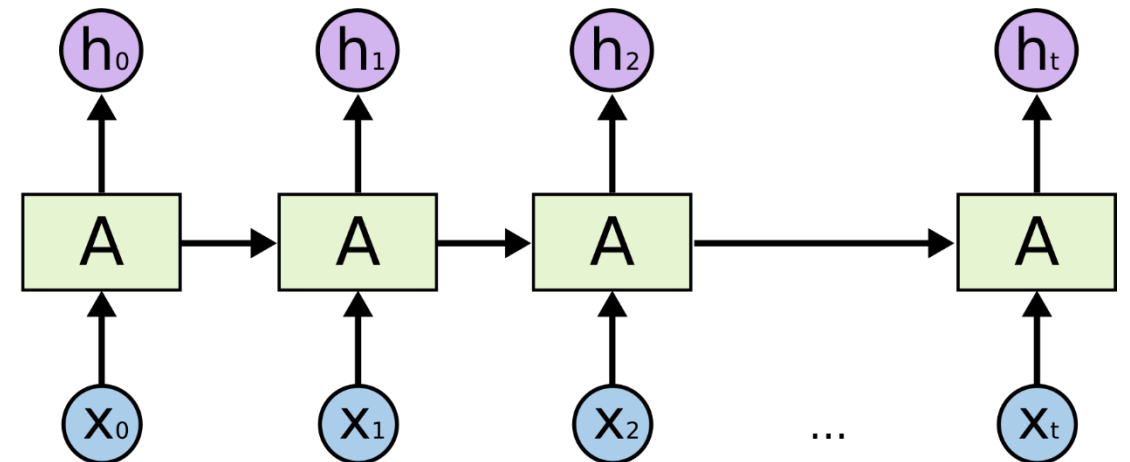
many to many



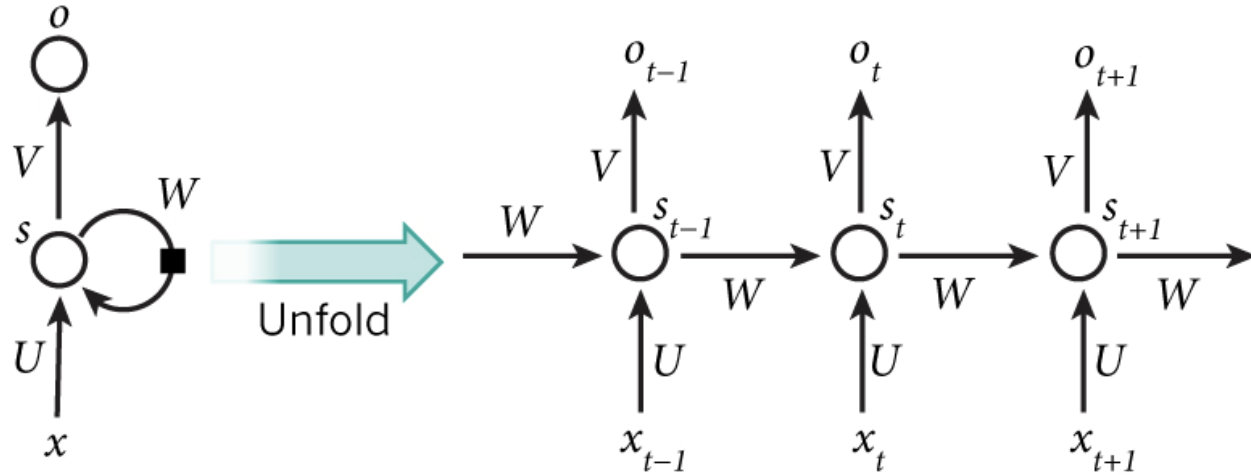
Recurrent Neural Network



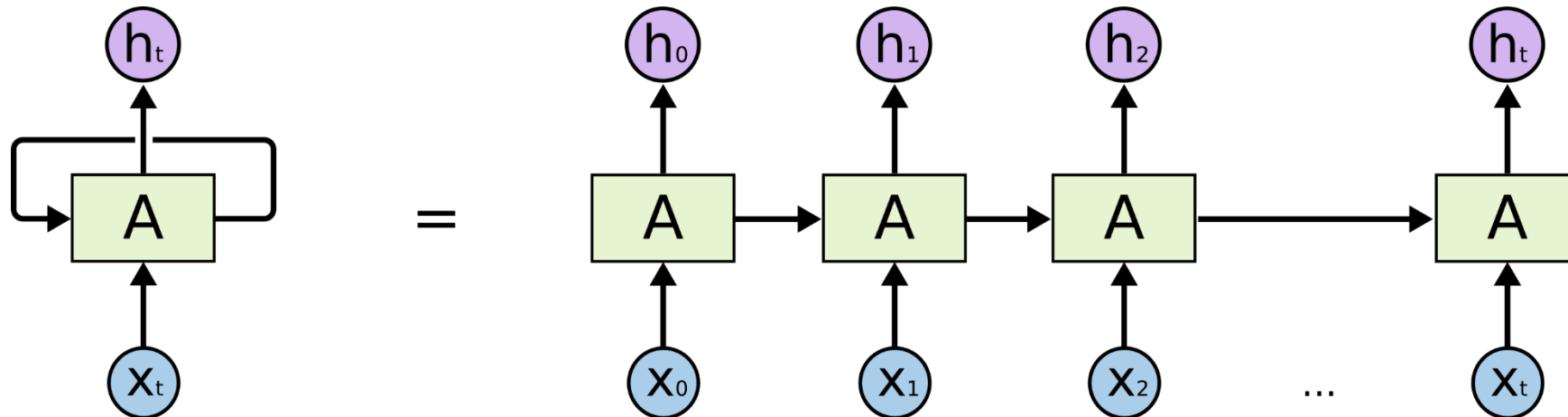
=



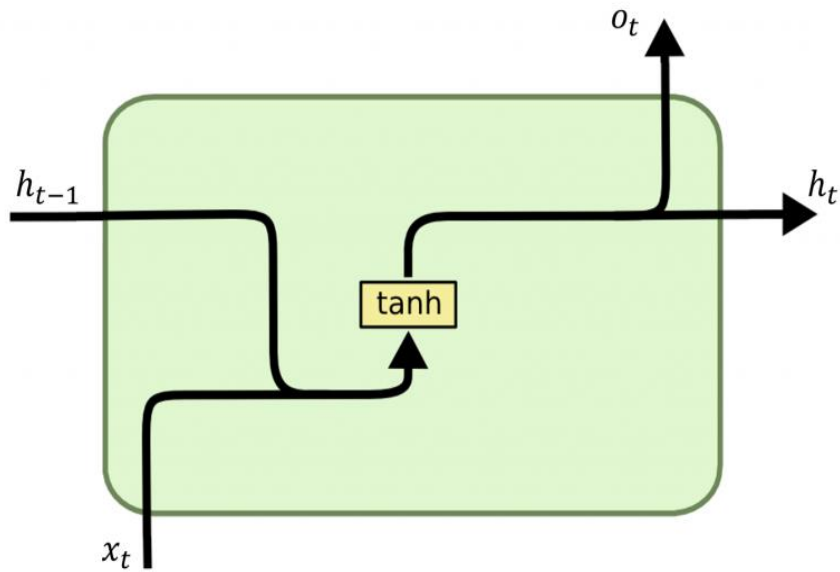
Recurrent Neural Network



$$\begin{aligned}
 a^{(t)} &= b + Wh^{(t-1)} + Ux^{(t)} \\
 h^{(t)} &= \tanh(a^{(t)}) \\
 o^{(t)} &= c + Vh^{(t)} \\
 \hat{y}^{(t)} &= \text{softmax}(o^{(t)})
 \end{aligned}$$



Recurrent layer representation



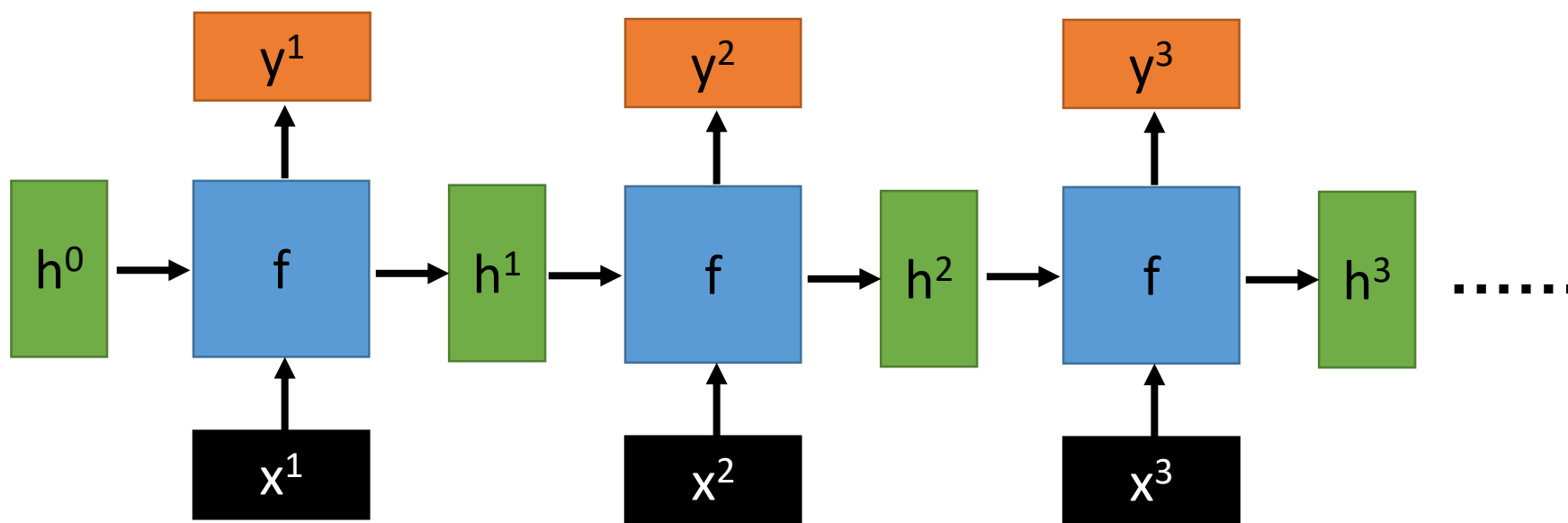
$$h_t = f(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$

$$o_t = \text{SoftMax}(W_{ho}h_t + b_o)$$

How does RNN reduce complexity?

- Given function $f: h', y = f(h, x)$

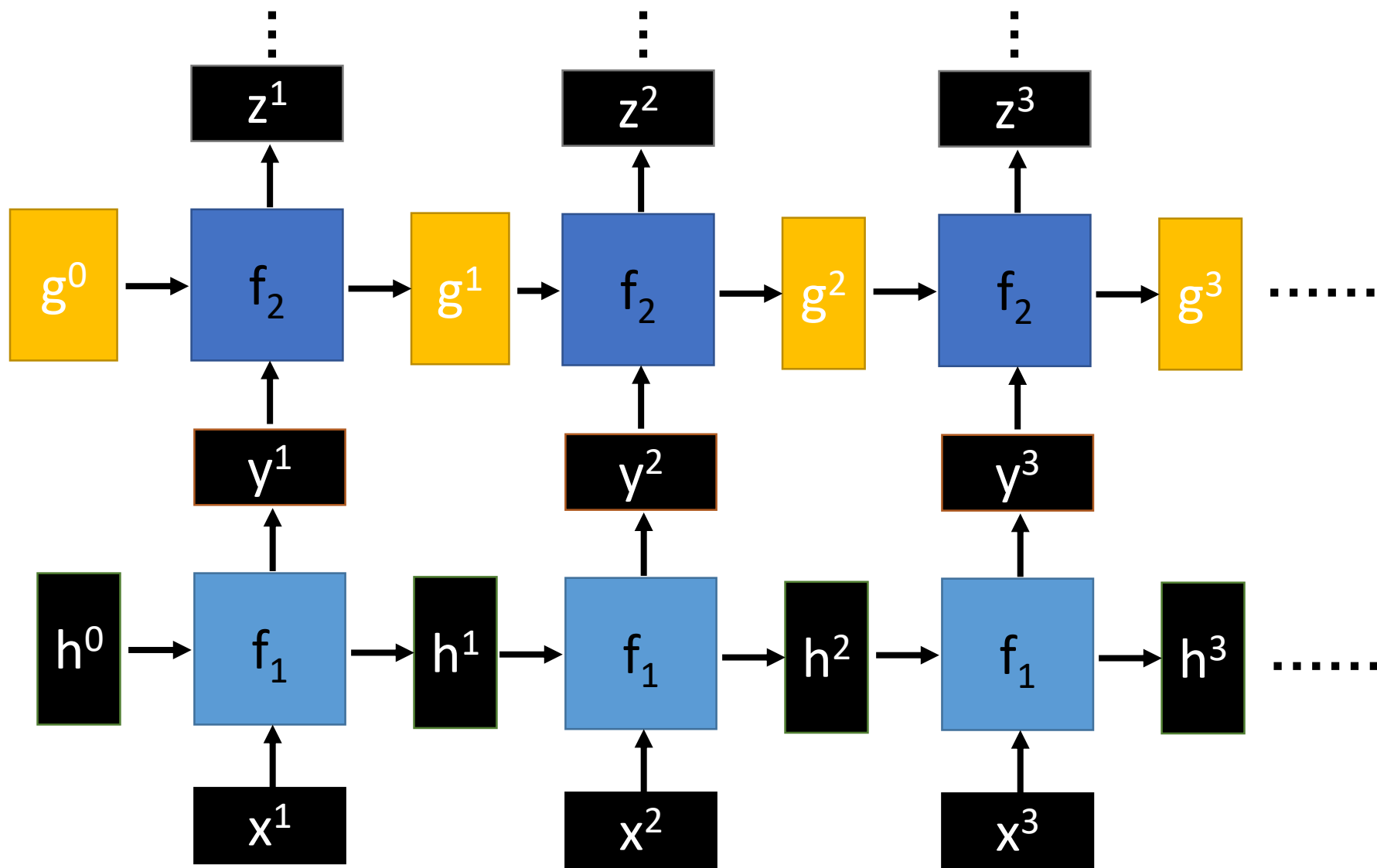
h and h' are vectors with the same dimension



No matter how long the input/output sequence is, we only need one function f . If f 's are different, then it becomes a feedforward NN. This may be treated as another compression from fully connected network.

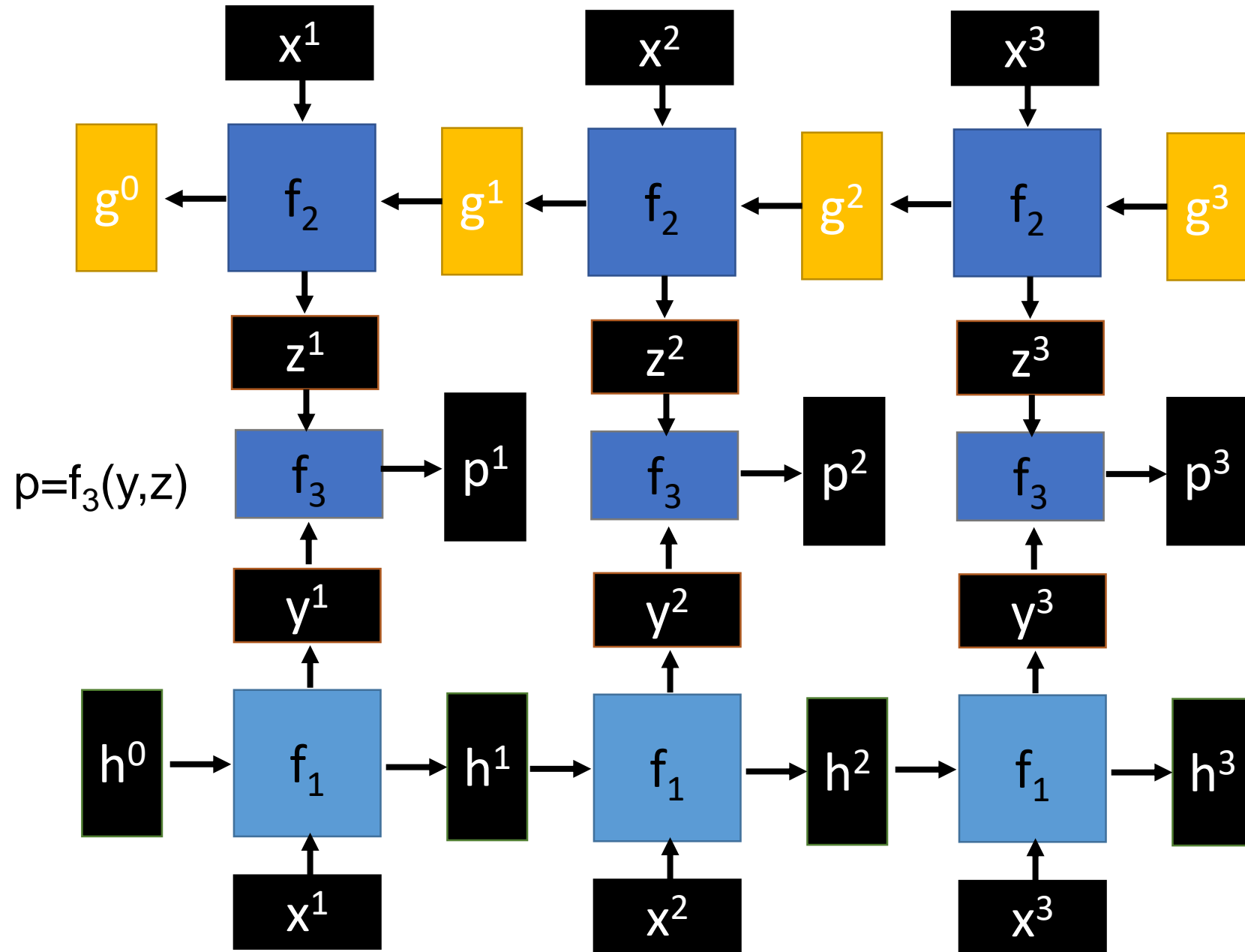
Deep RNN

$$h', y = f_1(h, x), g', z = f_2(g, y) \quad \dots$$



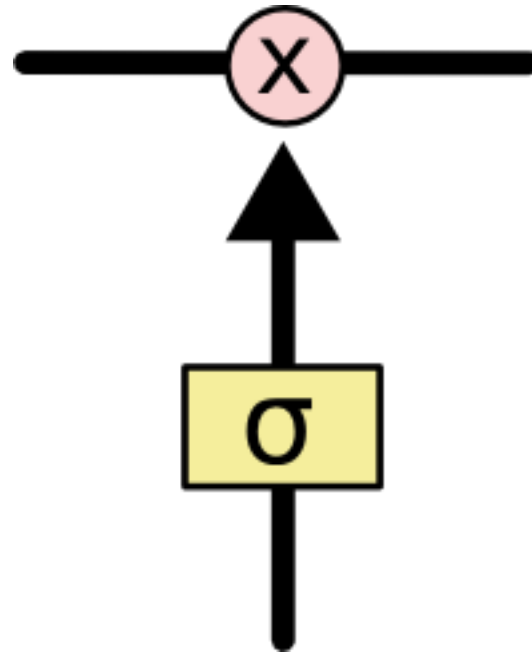
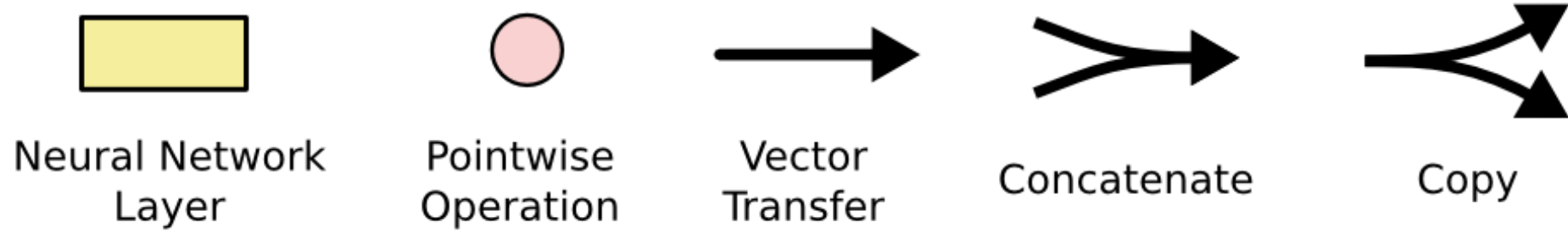
Bidirectional RNN

$$y, h = f_1(x, h) \quad z, g = f_2(g, x)$$



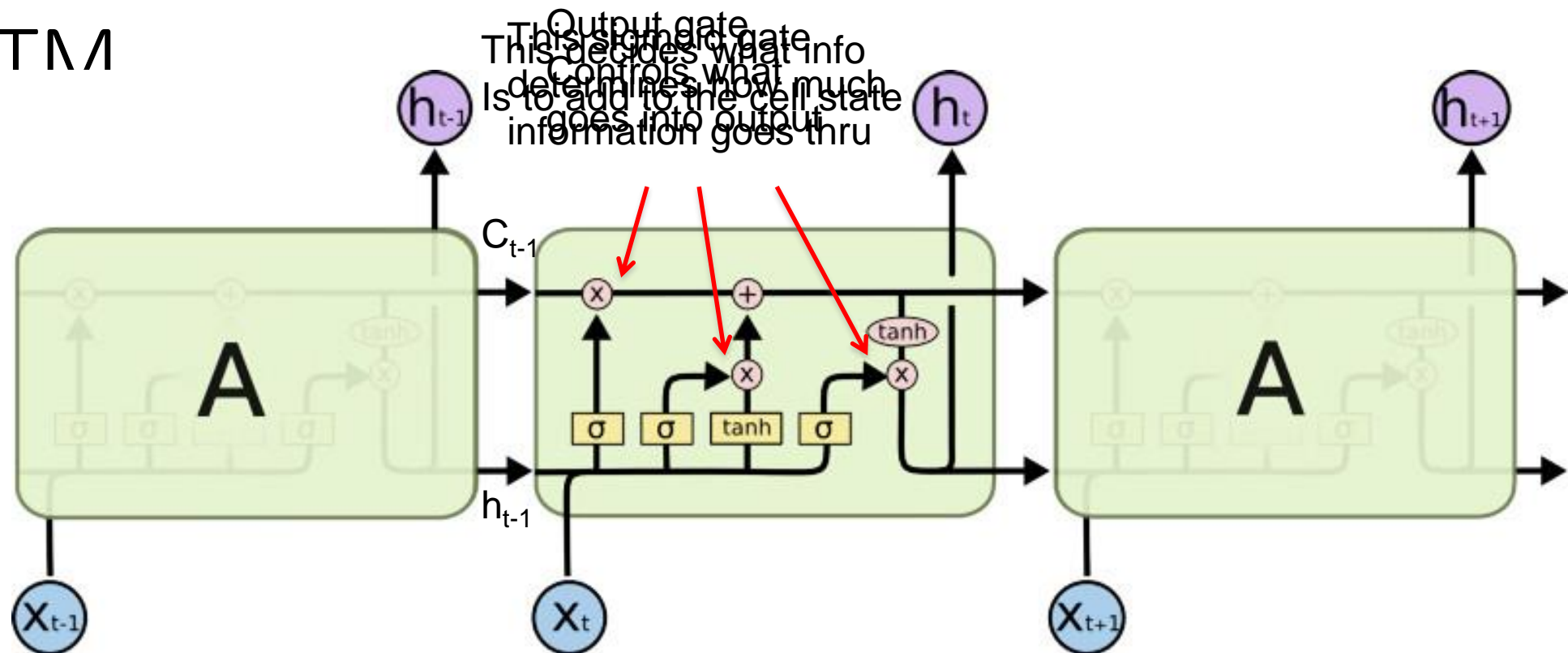
Problems with naive RNN

- When dealing with a time series, it tends to forget old information. When there is a distant relationship of unknown length, we wish to have a “memory” to it.
- Vanishing gradient problem.



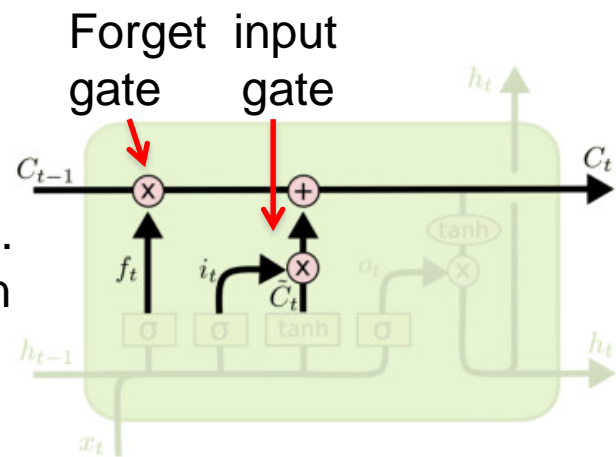
The sigmoid layer outputs numbers between 0-1 determine how much each component should be let through. Pink X gate is point-wise multiplication.

LSTM

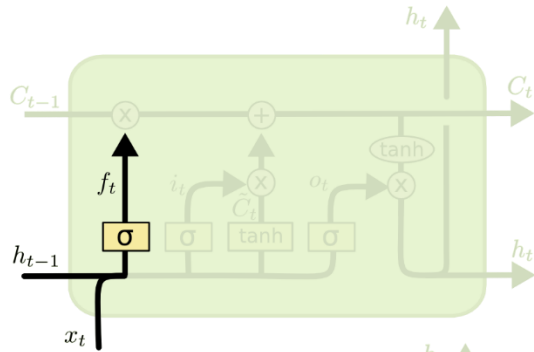


Output gate
This decides what info
goes into the output
Is to add to the cell state
information goes thru

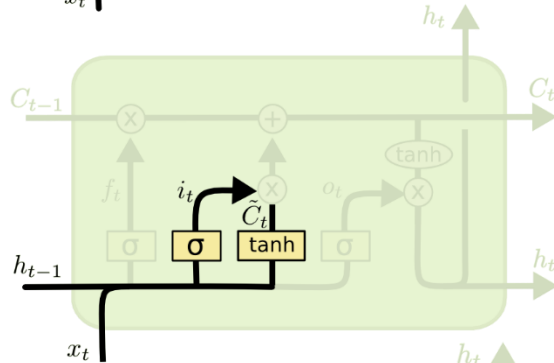
The core idea is this cell
state C_t it is changed
Why sigmoid or tanh:
Sigmoid: 0, 1 gating as switch.
slowly, with only minor
Vanishing gradient problem in
linear interactions. It is very
LSTM is handled already.
easy for information to flow
ReLU replaces tanh ok?
along it unchanged.



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

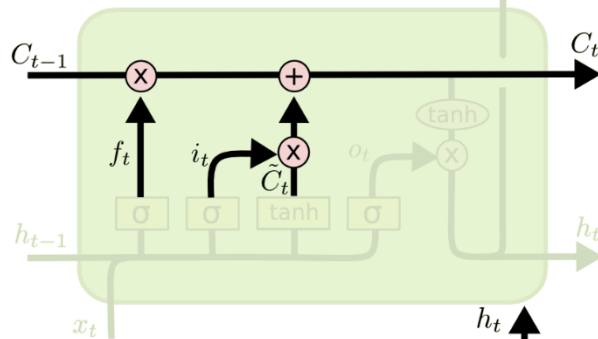


$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

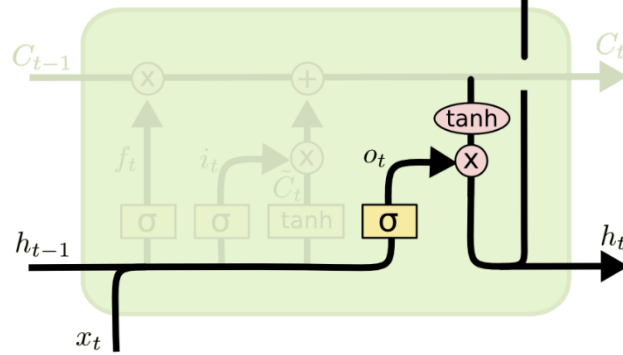


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

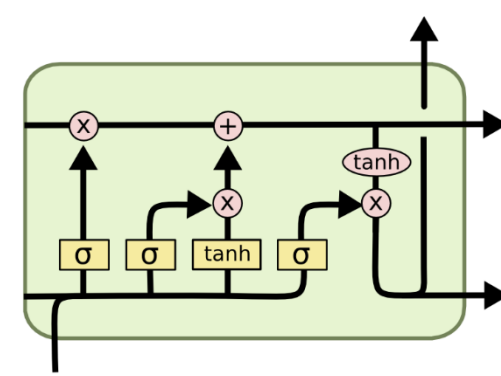


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$



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

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

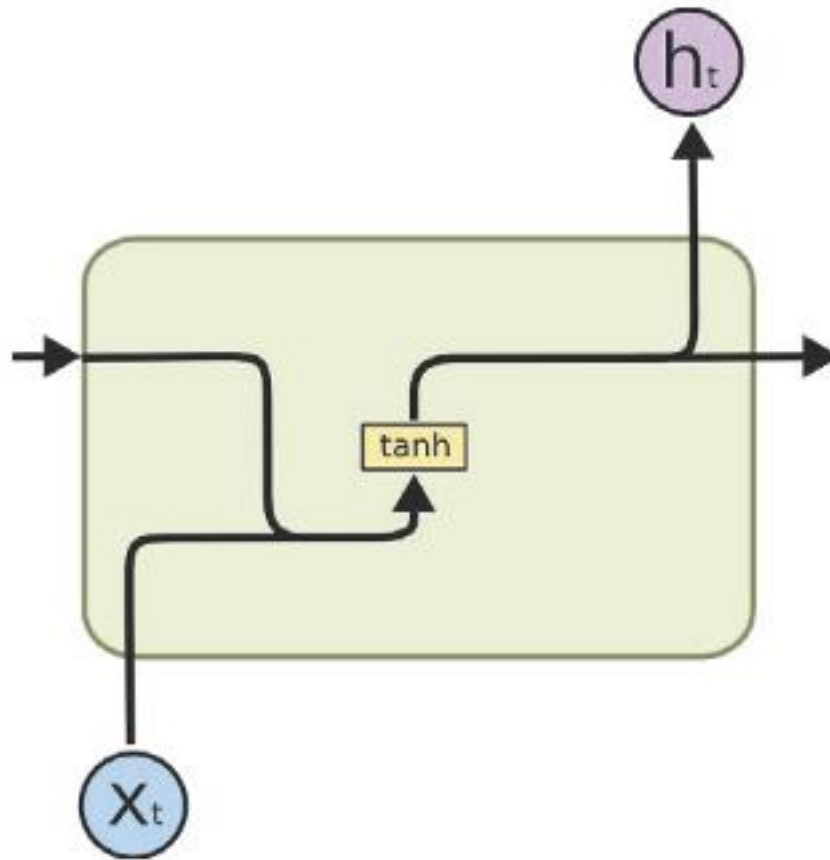


i_t decides what component
is to be updated.
 C'_t provides change contents

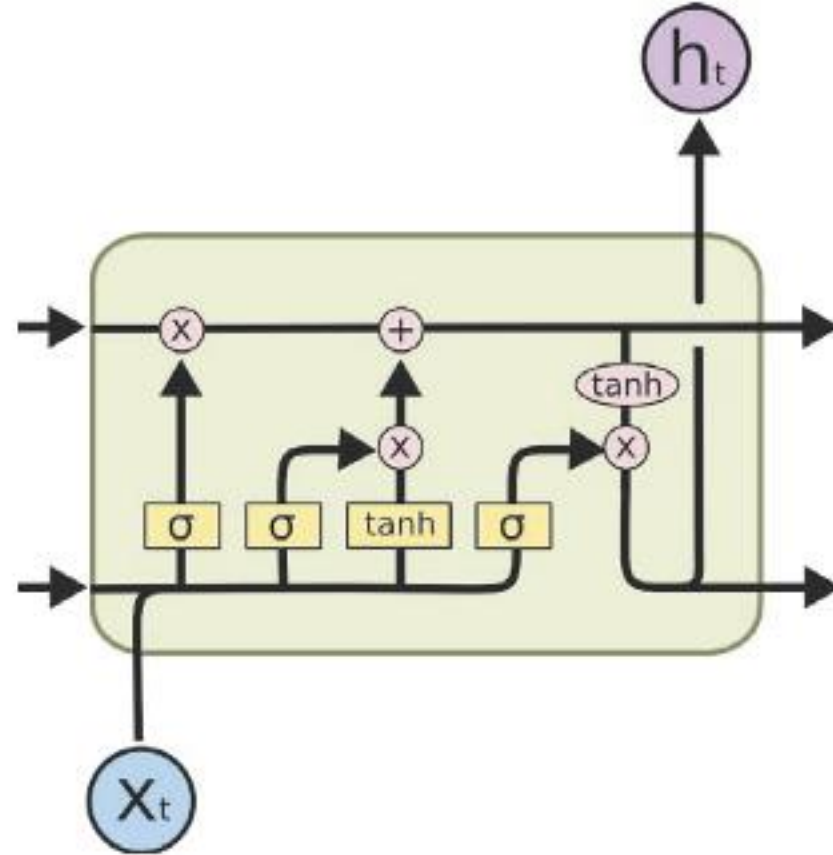
Updating the cell state

Decide what part of the cell
state to output

RNN vs LSTM

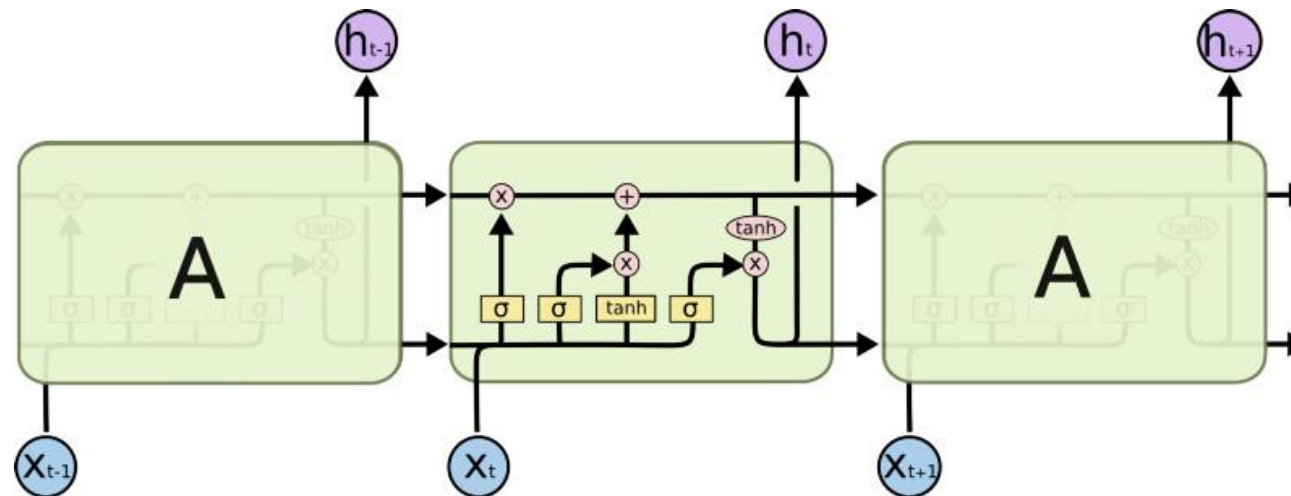


(a) RNN

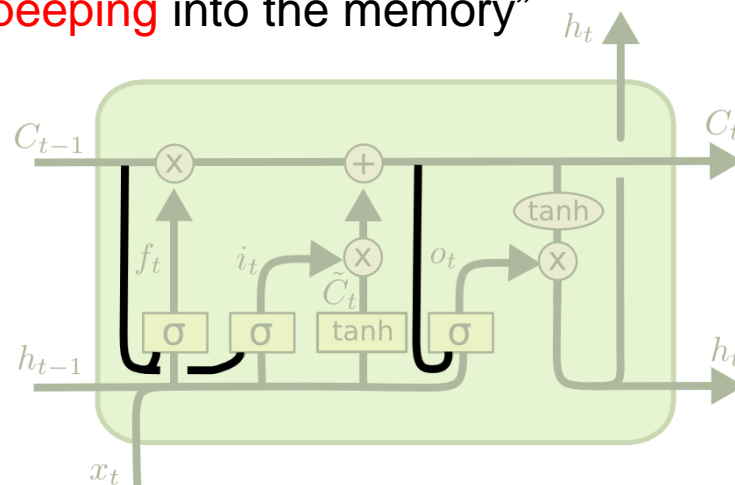


(b) LSTM

Peephole LSTM



Allows “**peeping** into the memory”

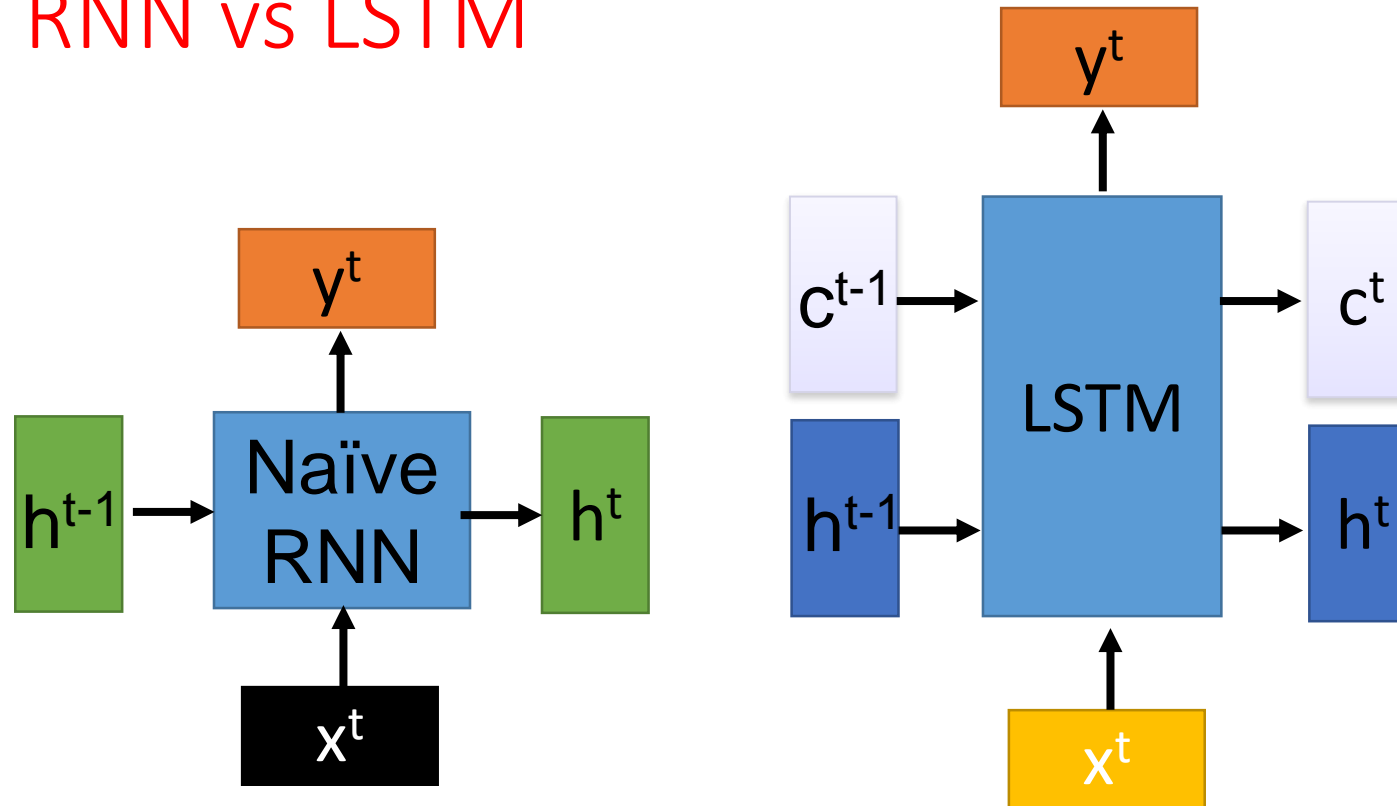


$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

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

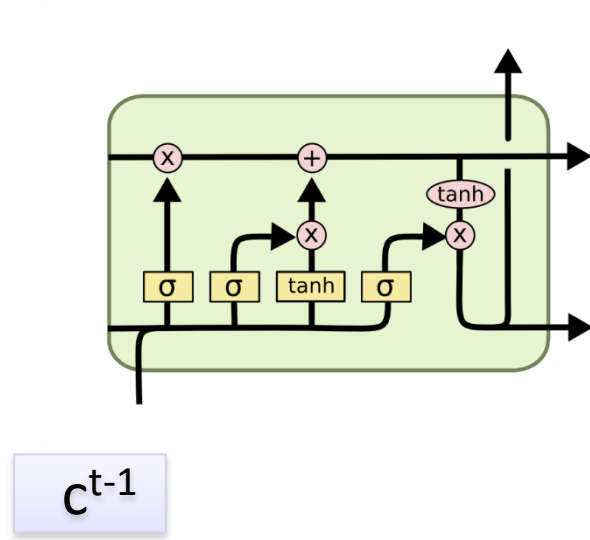
$$o_t = \sigma(W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$

Naïve RNN vs LSTM



c changes slowly $\Rightarrow c^t$ is c^{t-1} added by something

h changes faster $\Rightarrow h^t$ and h^{t-1} can be very different



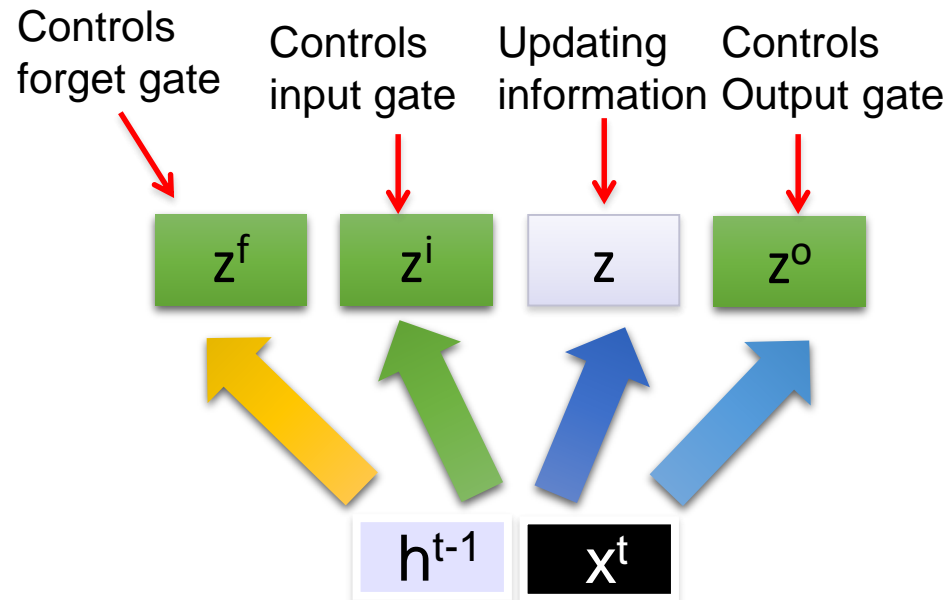
These 4 matrix computation should be done concurrently.

$$z = \tanh(W \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix})$$

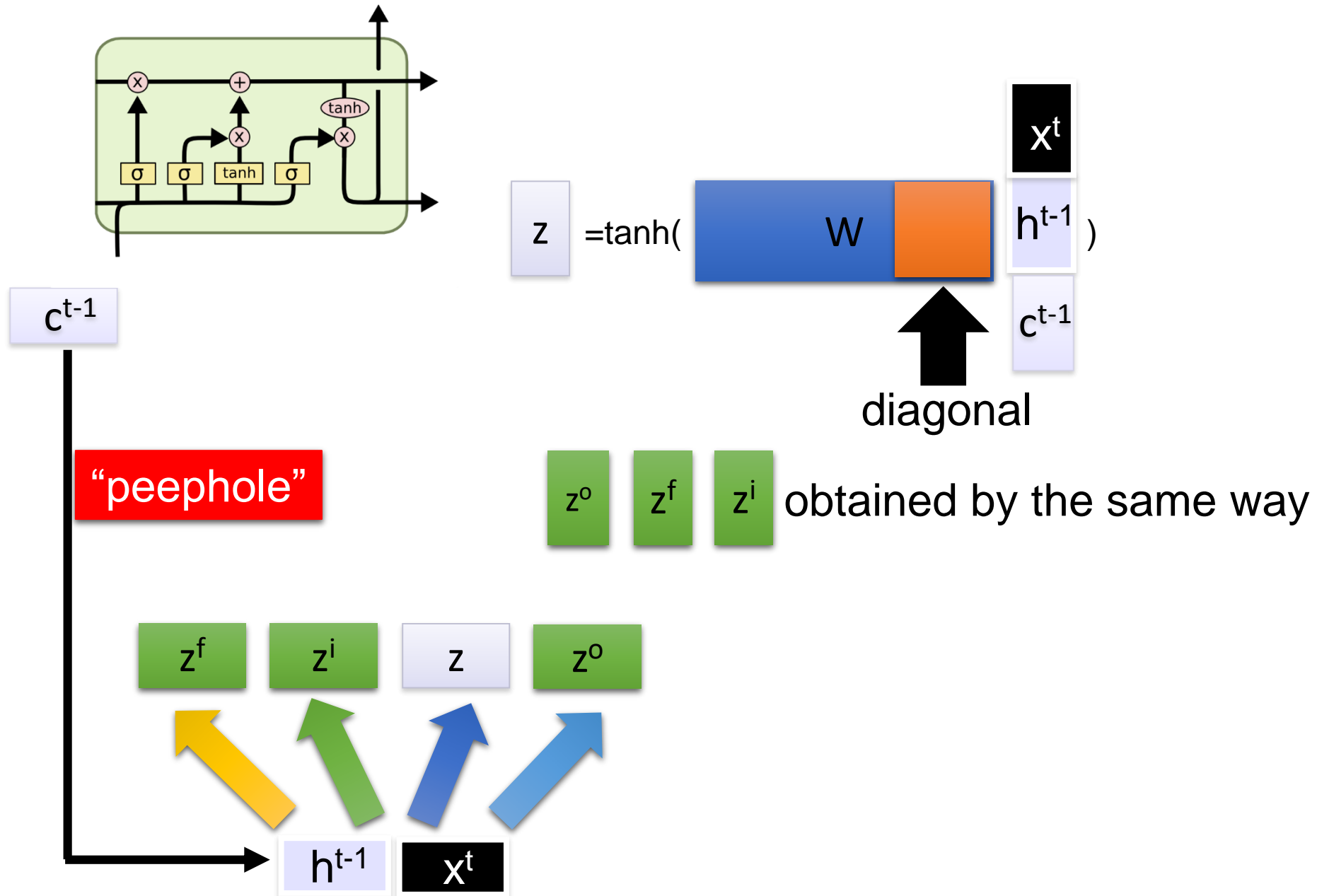
$$z^i = \sigma(W^i \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix})$$

$$z^f = \sigma(W^f \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix})$$

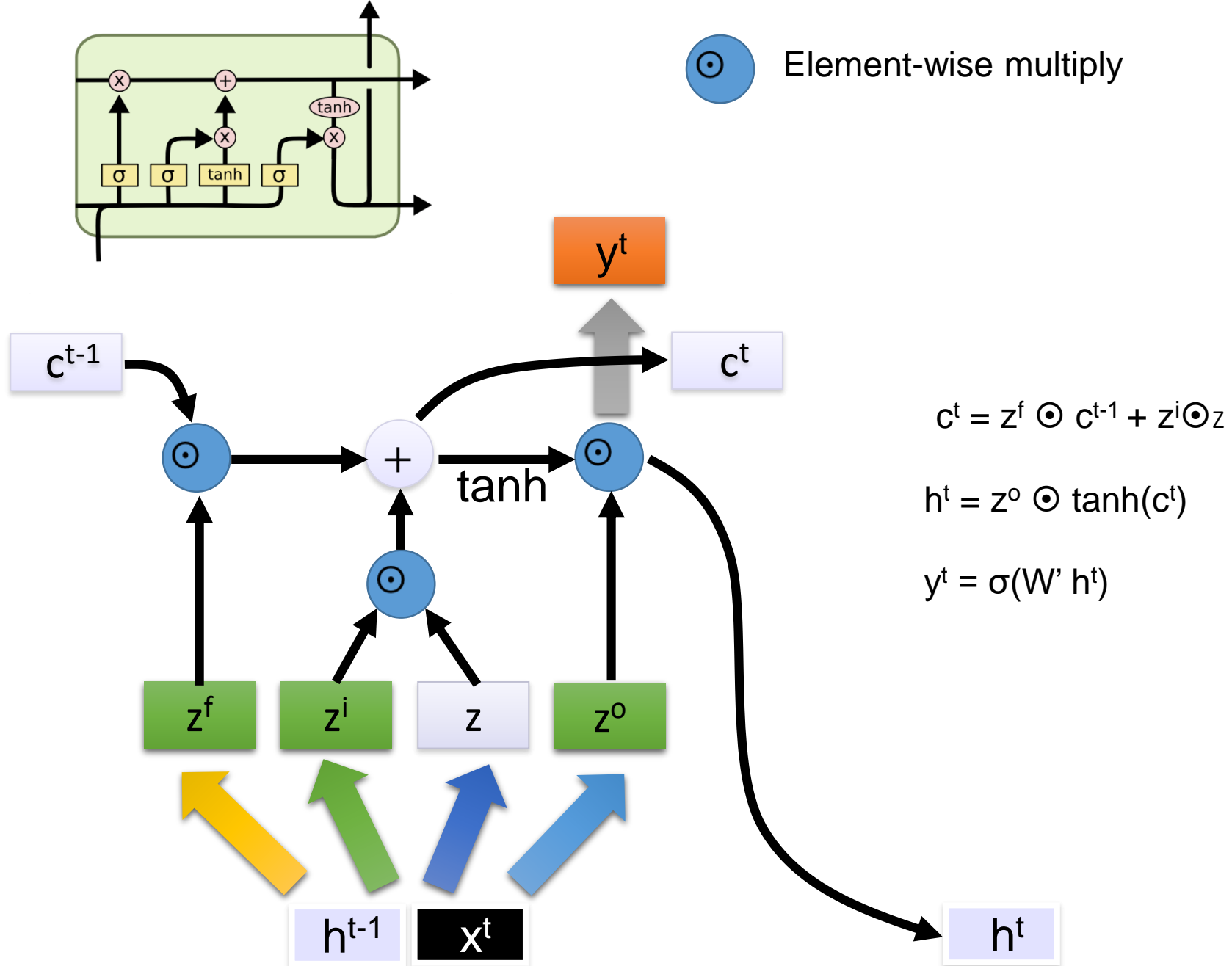
$$z^o = \sigma(W^o \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix})$$



Information flow of LSTM

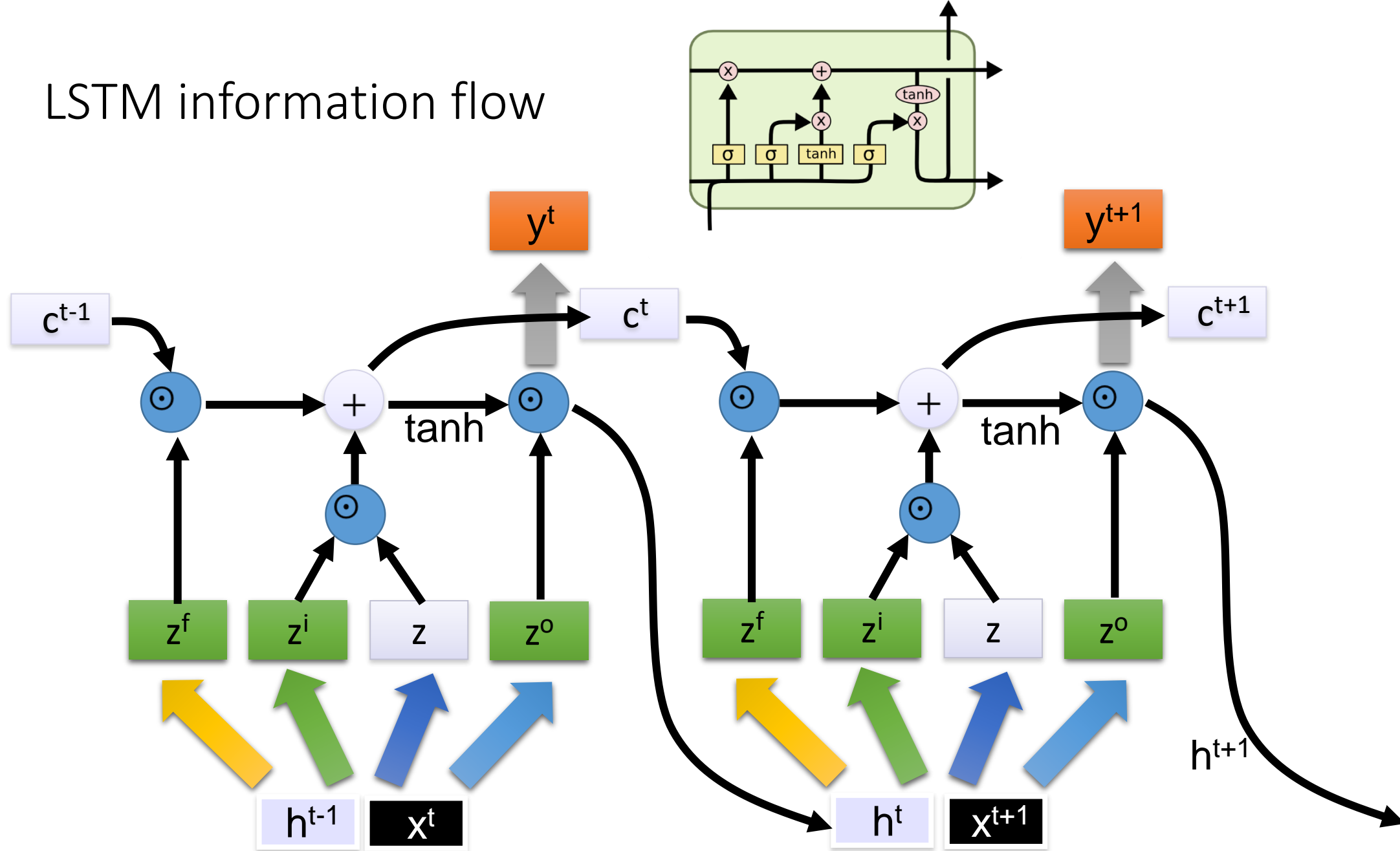


Information flow of LSTM



Information flow of LSTM

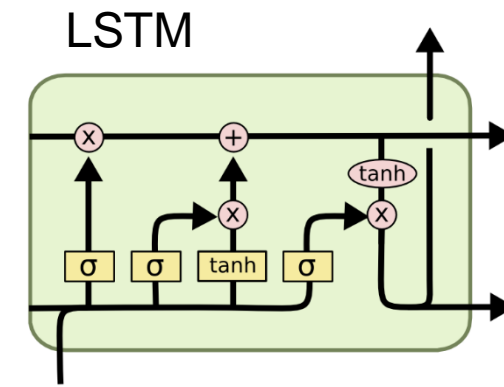
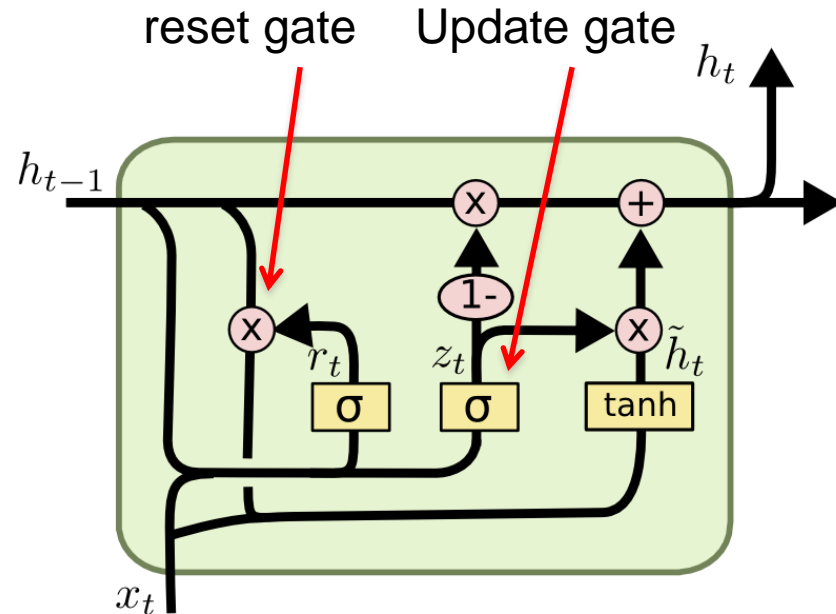
LSTM information flow



Information flow of LSTM

GRU – gated recurrent unit

(more compression)



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

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

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

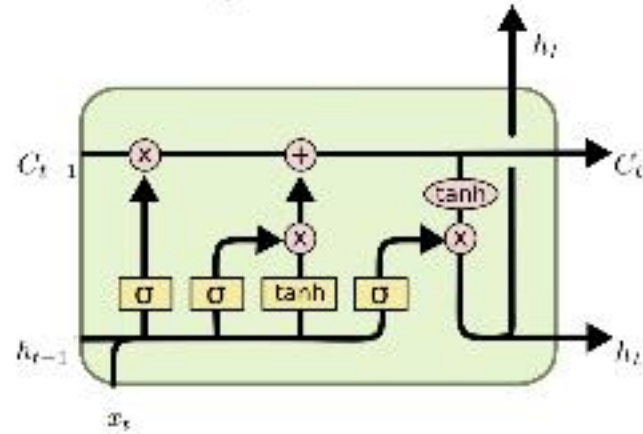
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

It combines the **forget** and **input** into a single **update gate**.
It also merges the cell state and hidden state. This is simpler than LSTM. There are many other variants too.

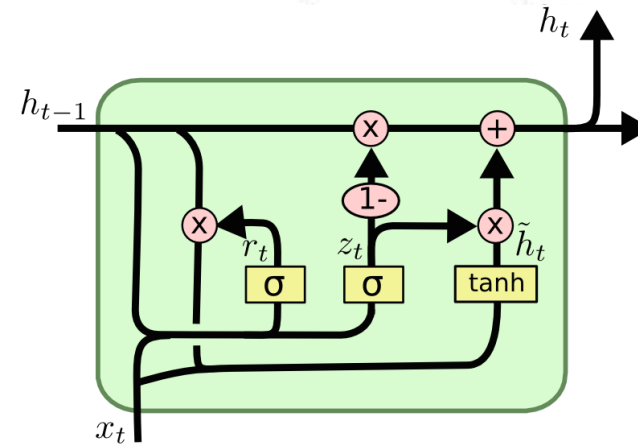
$X, *$: element-wise multiply

LSTM and GRU

- LSTM [Hochreiter&Schmidhuber97]



- GRU [Cho+14]



GRUs also take x_t and h_{t-1} as inputs. They perform some calculations and then pass along h_t . What makes them different from LSTMs is that GRUs don't need the cell layer to pass values along. The calculations within each iteration ensure that the h_t values being passed along either retain a high amount of old information or are jump-started with a high amount of new information.

PYTHON CODE

