

CS 4973 & 6463 :  
Advanced Topics in Data Science

# Deep Learning



# Recurrent Neural Networks

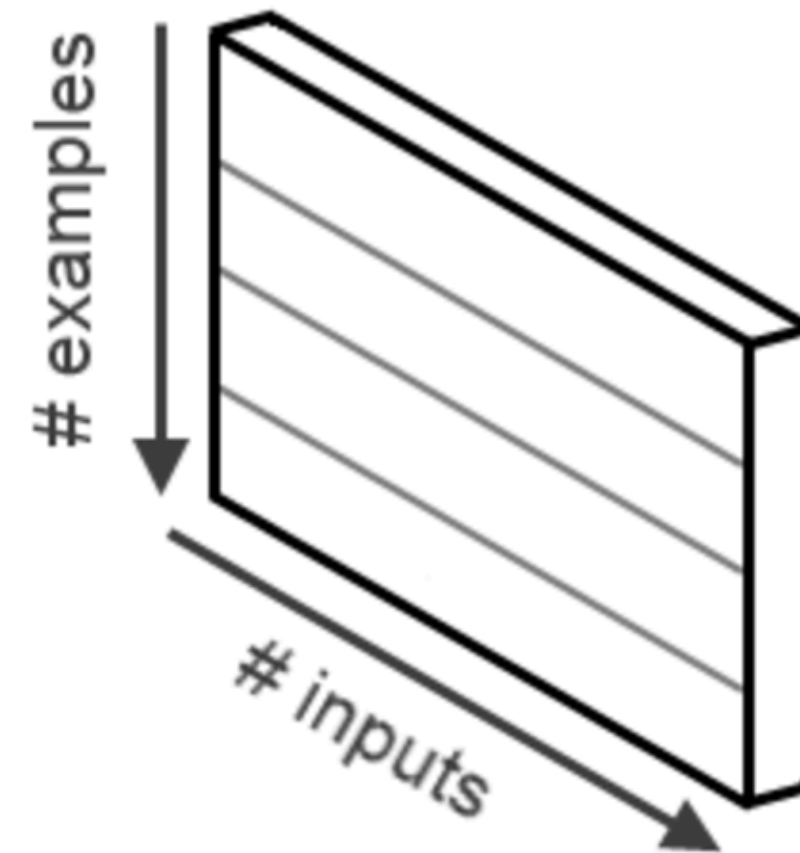
# RNNs

- Problem: **inputs of variable size**
  - *Image data can be normalized. Text is more difficult..*
- Solution: **recurrent neural networks**
  - *Varying length sequences handled by **stateful neurons***

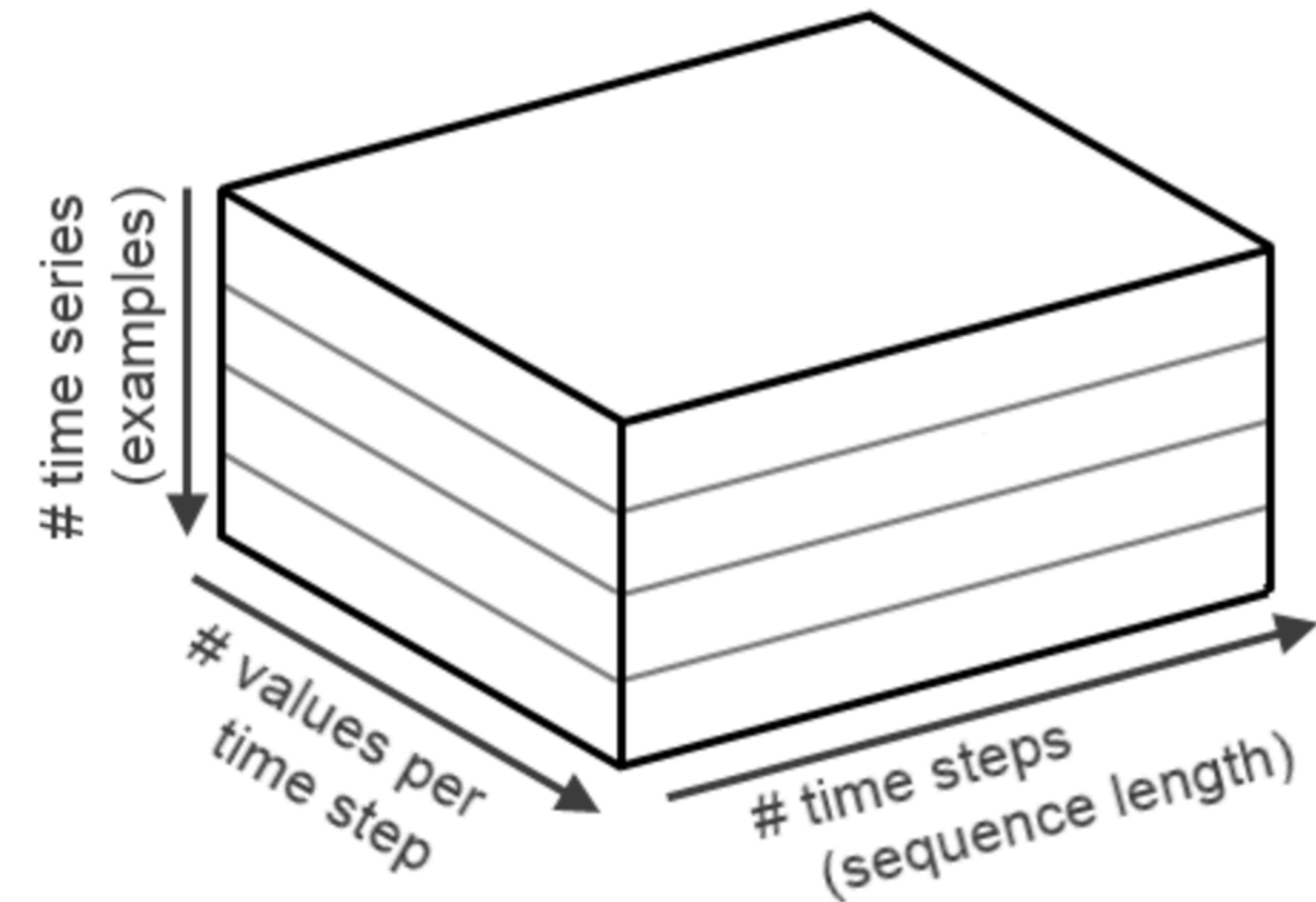
# RNNs: Motivating Applications

- Time-series data
- Language processing
  - Speech recognition
  - Text generation
- Sensor data
- Genome sequences
- Stock prediction

Feed Forward Network Data



Recurrent Network Data



# Beyond CNNs

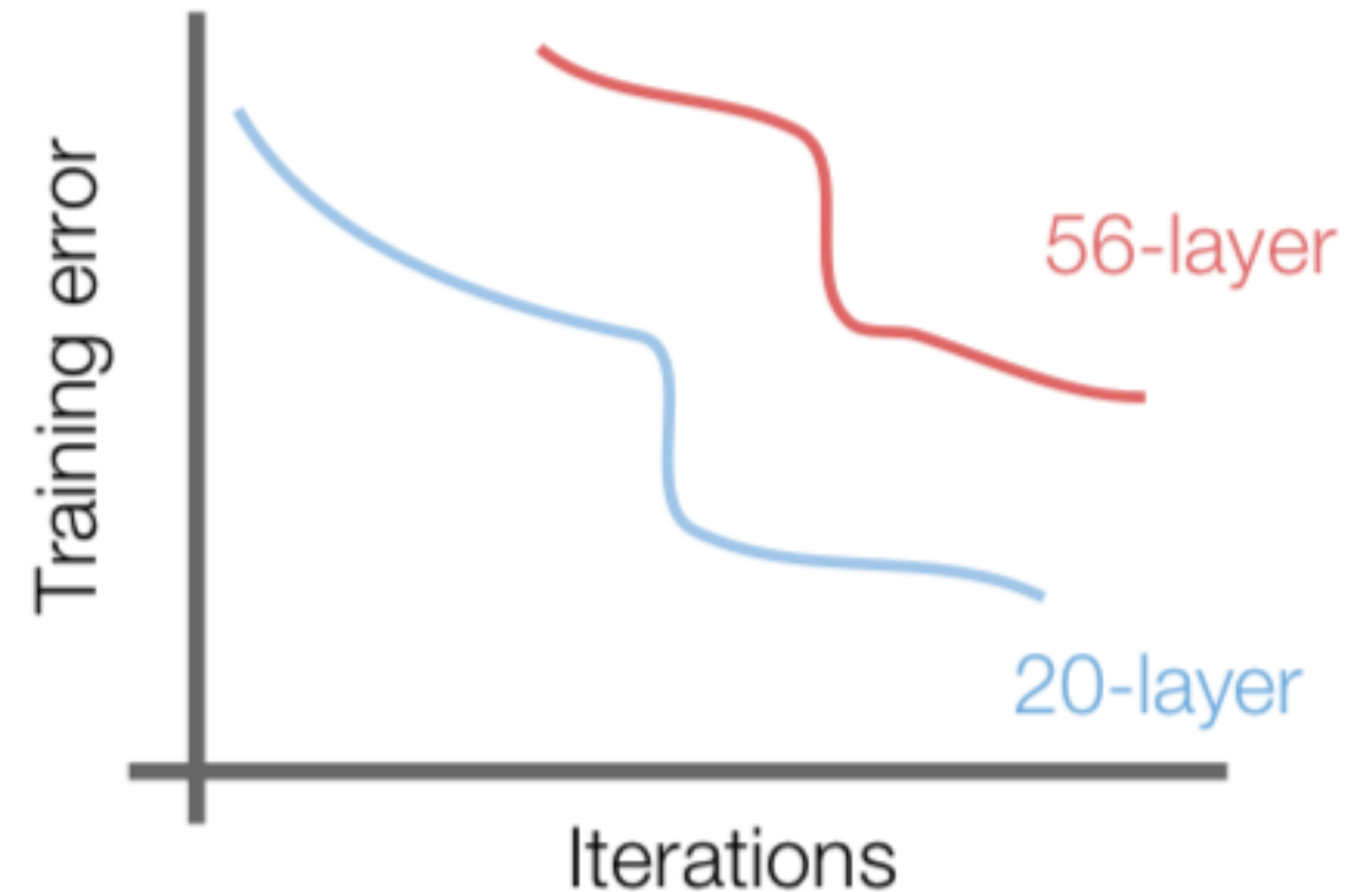
Problem: **Deeper models are hard to optimize.**

Solution: **Construct deep model from learned layers of shallow model.**

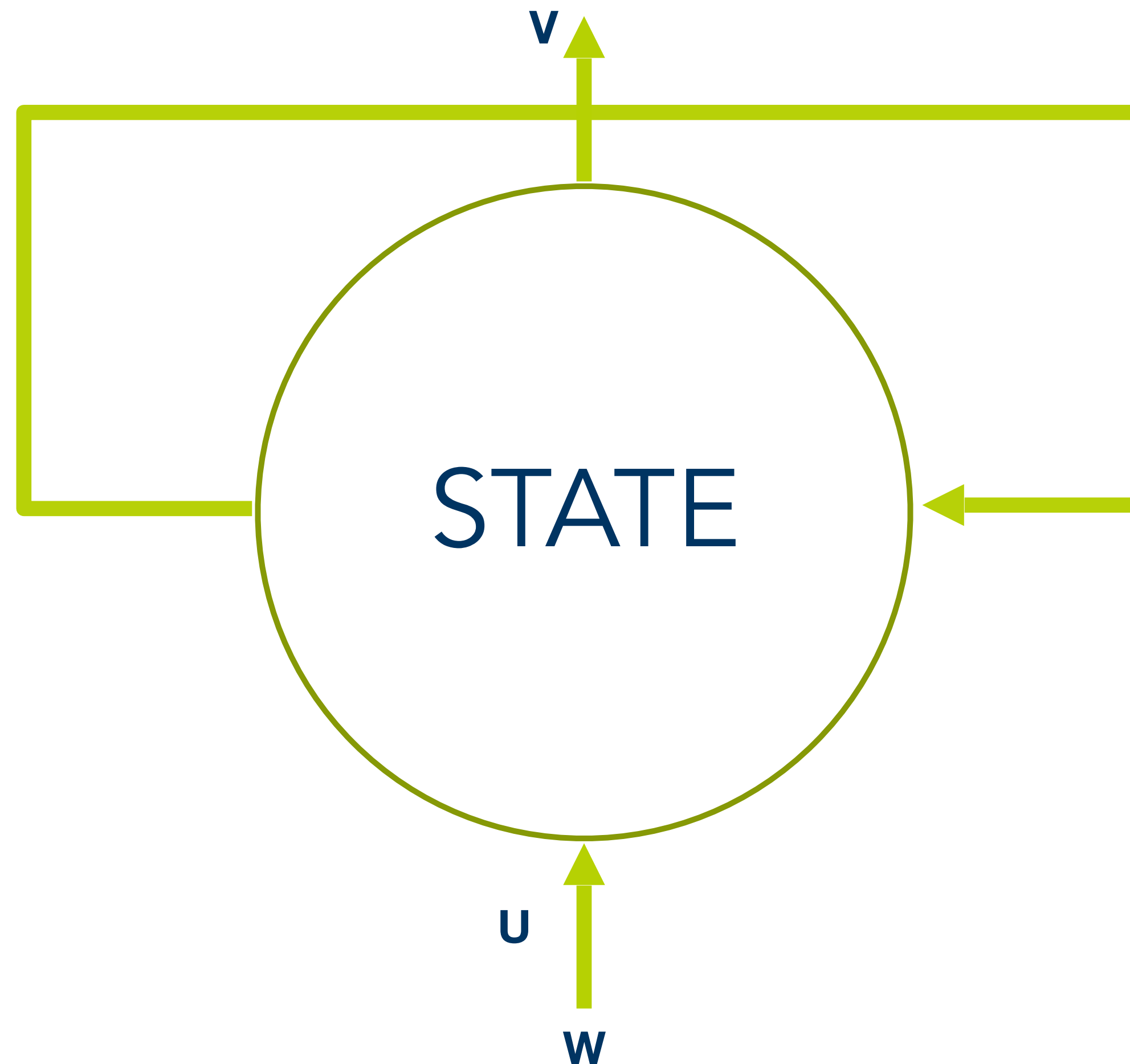
- Stacking deeper layers on CNNs may result in **worsening performance** in both training and testing.

- Is this due to overfitting?

*No! See the pattern?*



# RNNs: Architecture



Let:  
 $r$  = dim of input vector  
 $s$  = dim of hidden state  
 $t$  = dim output vector

$U$  encodes the input

$[s \times r]$

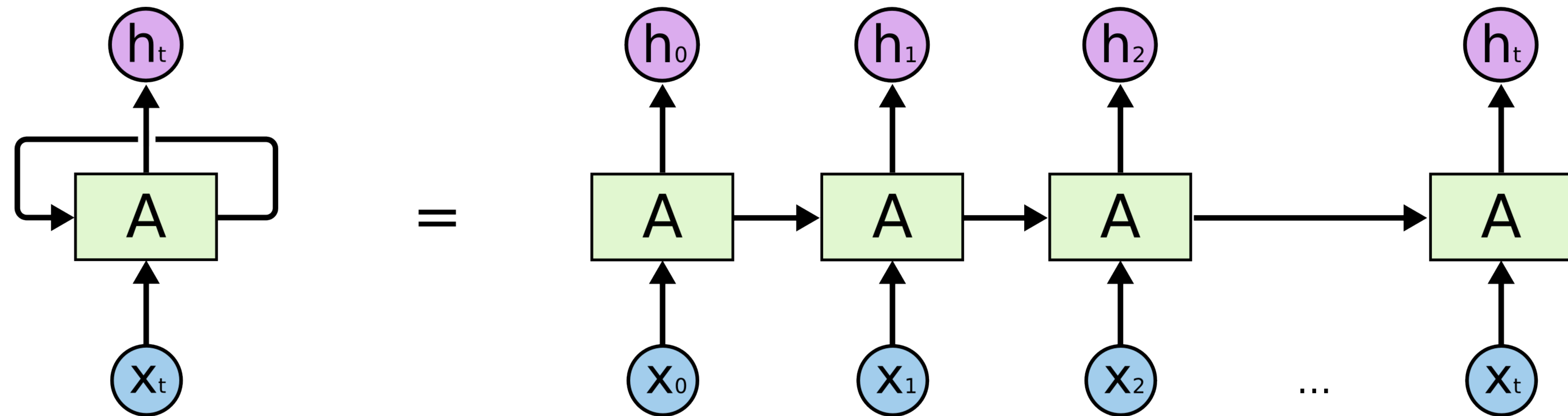
$W$  updates the info from  
the previous state

$[s \times s]$

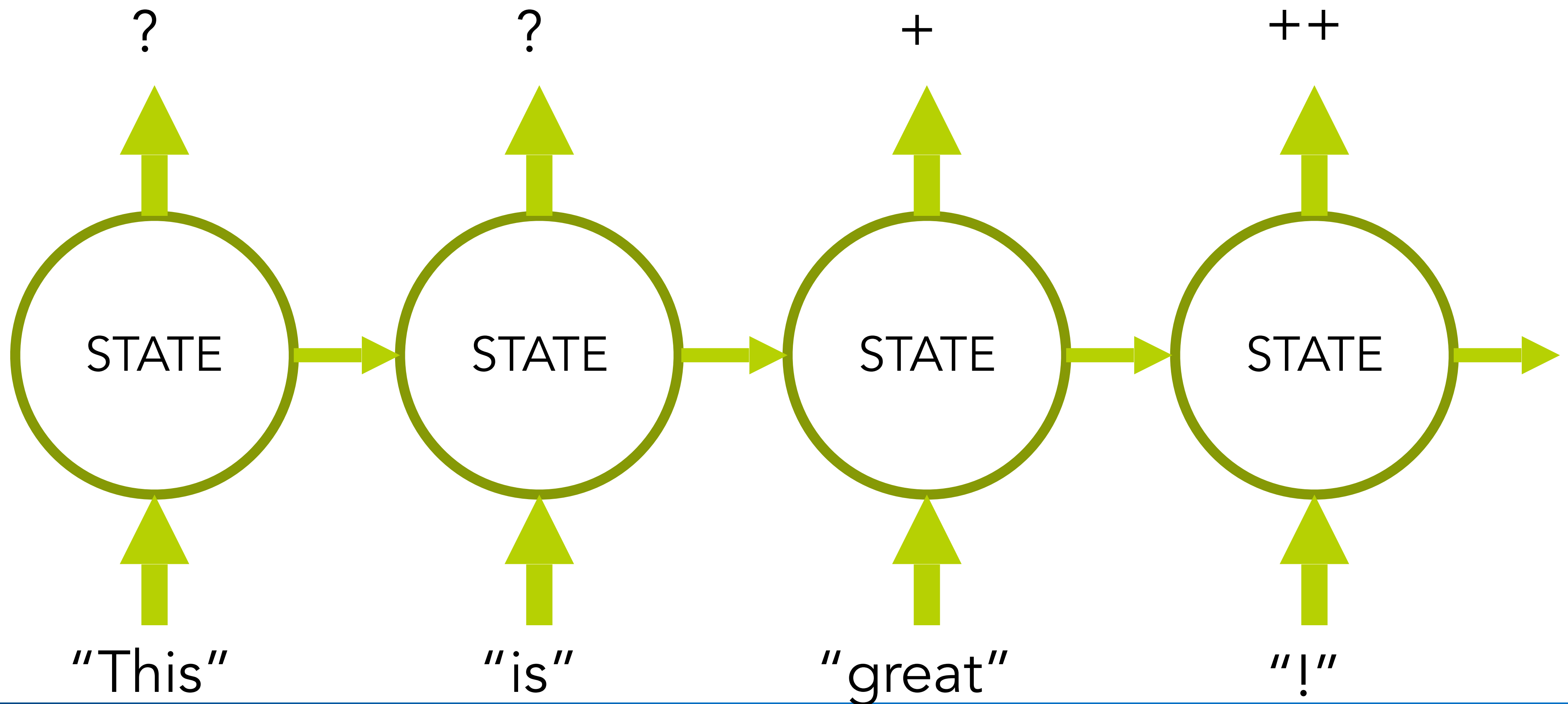
$U+W$  is passed into the  
activation function

$V$  is  $[t \times s]$

# “Unfolding” an RNN



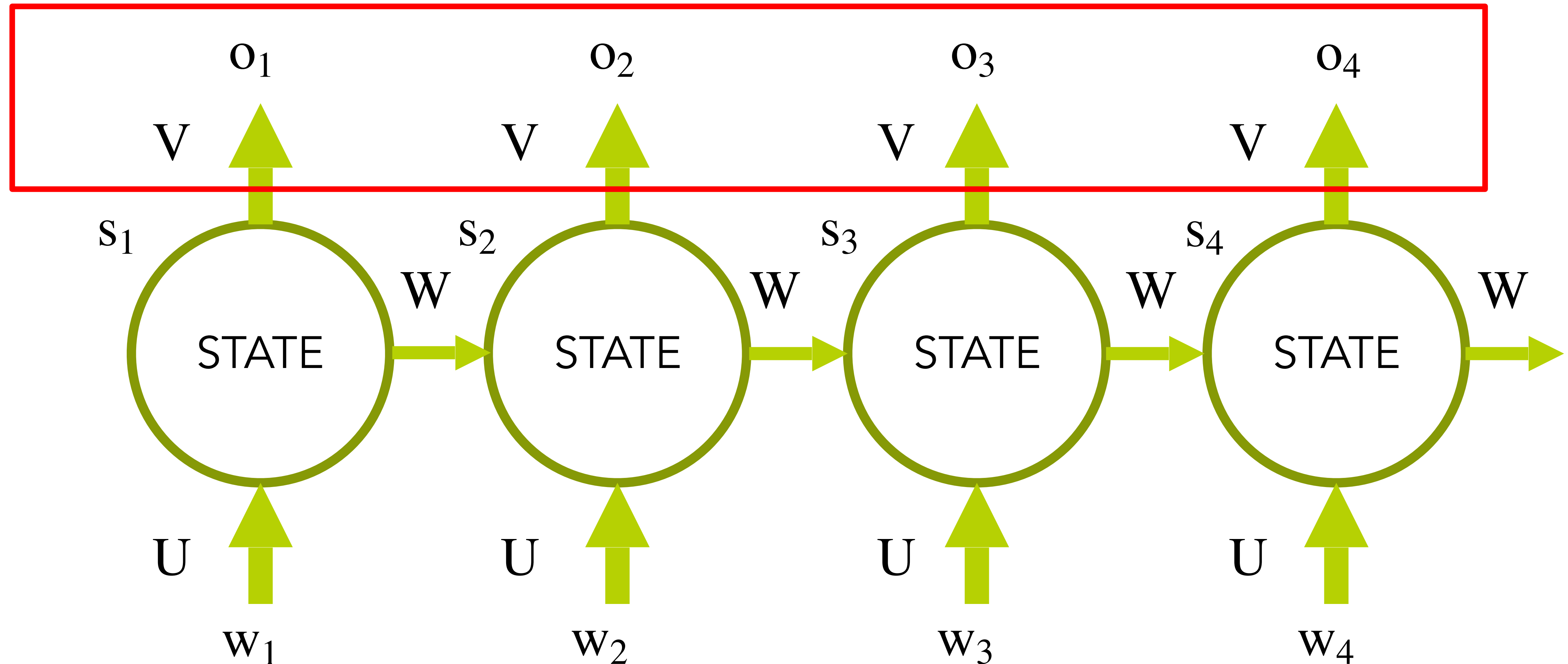
# "Unrolling" the RNN





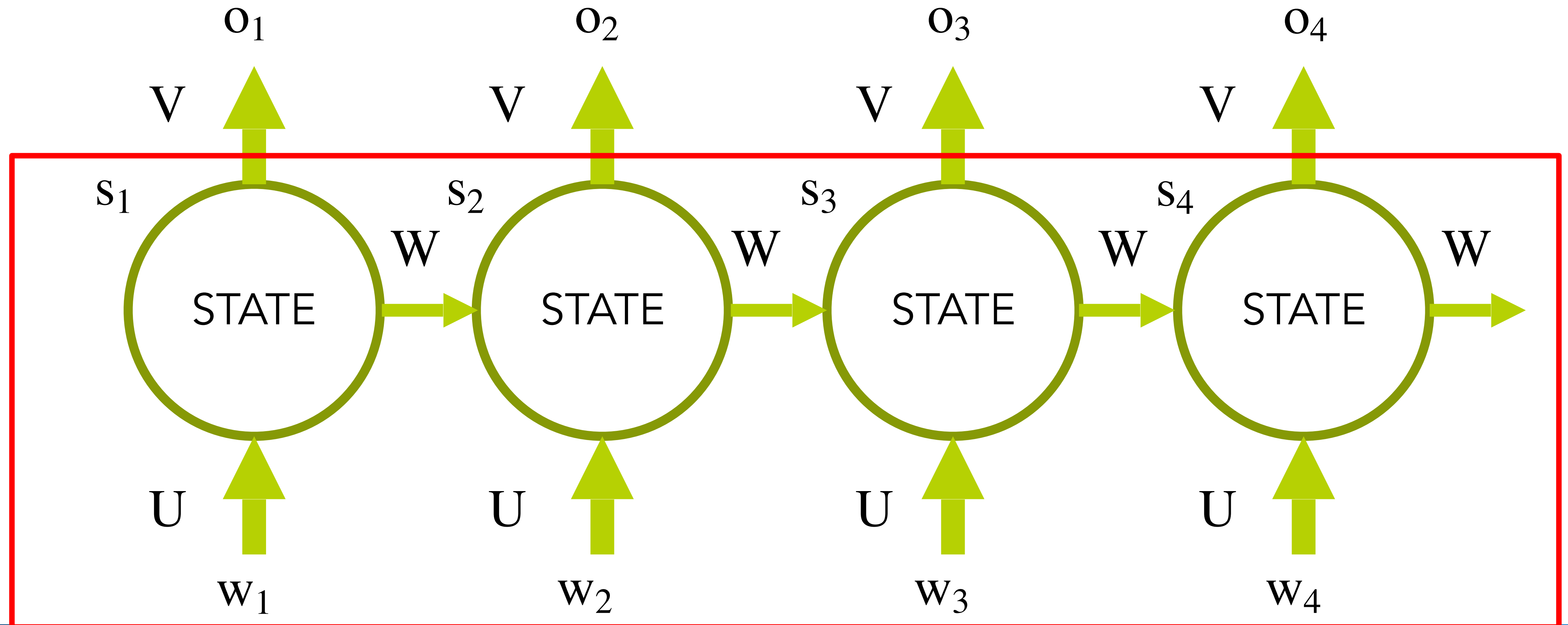
# “Unrolling” the RNN

In Keras, this part is accomplished by a subsequent Dense layer



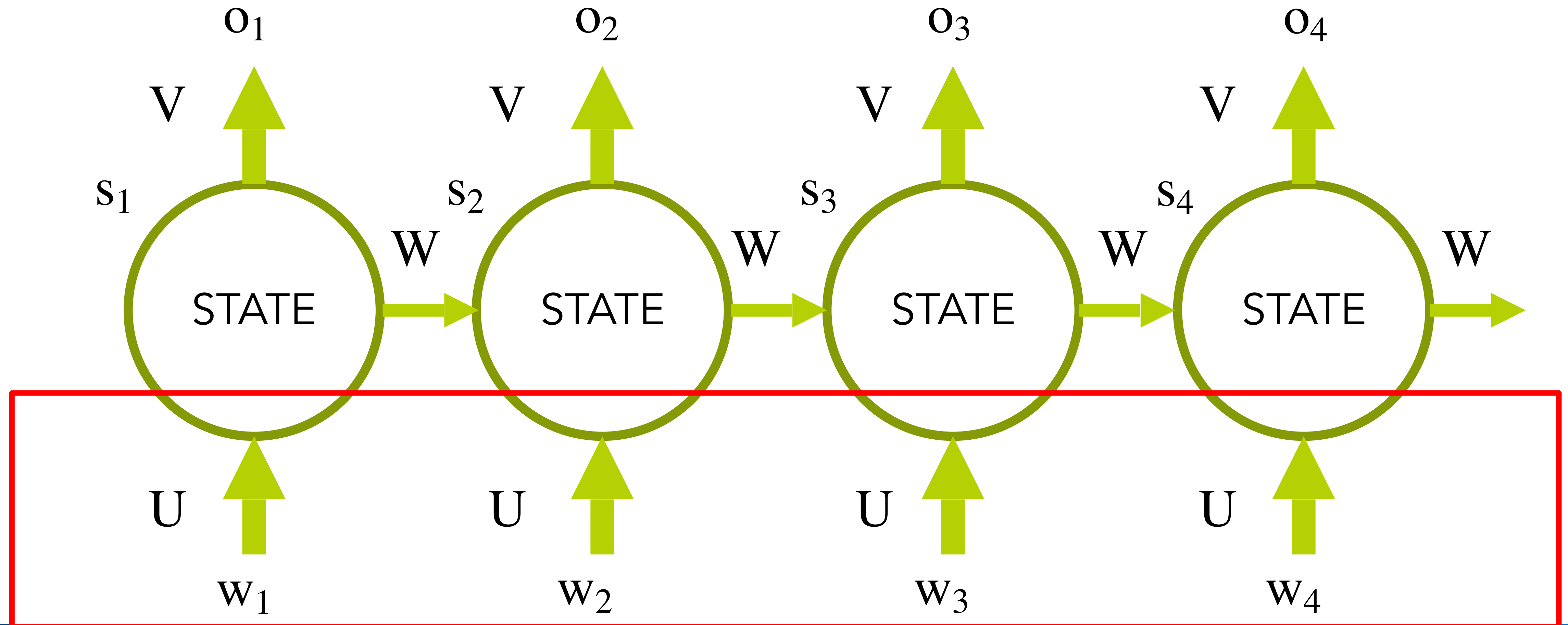
# “Unrolling” the RNN

This part is the core RNN



# “Unrolling” the RNN

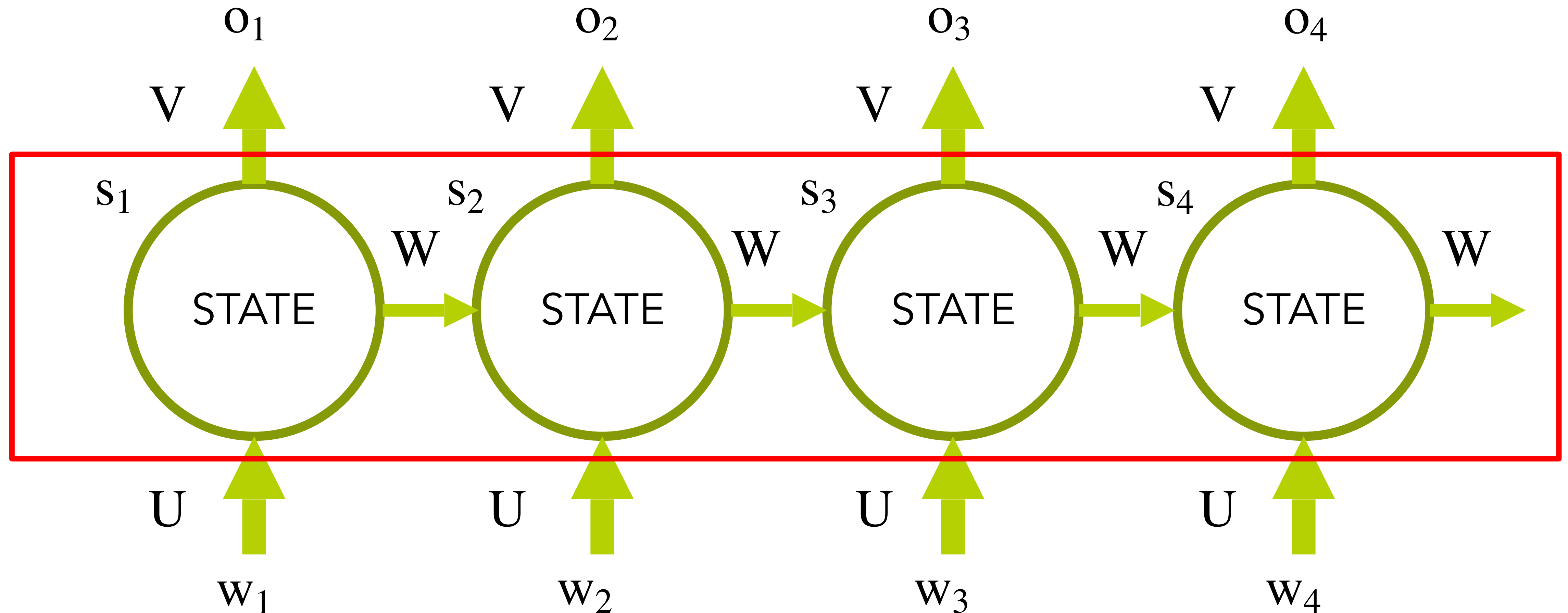
Keras calls this part the “kernel” (e.g. `kernel_initializer,...`)





# "Unrolling" the RNN

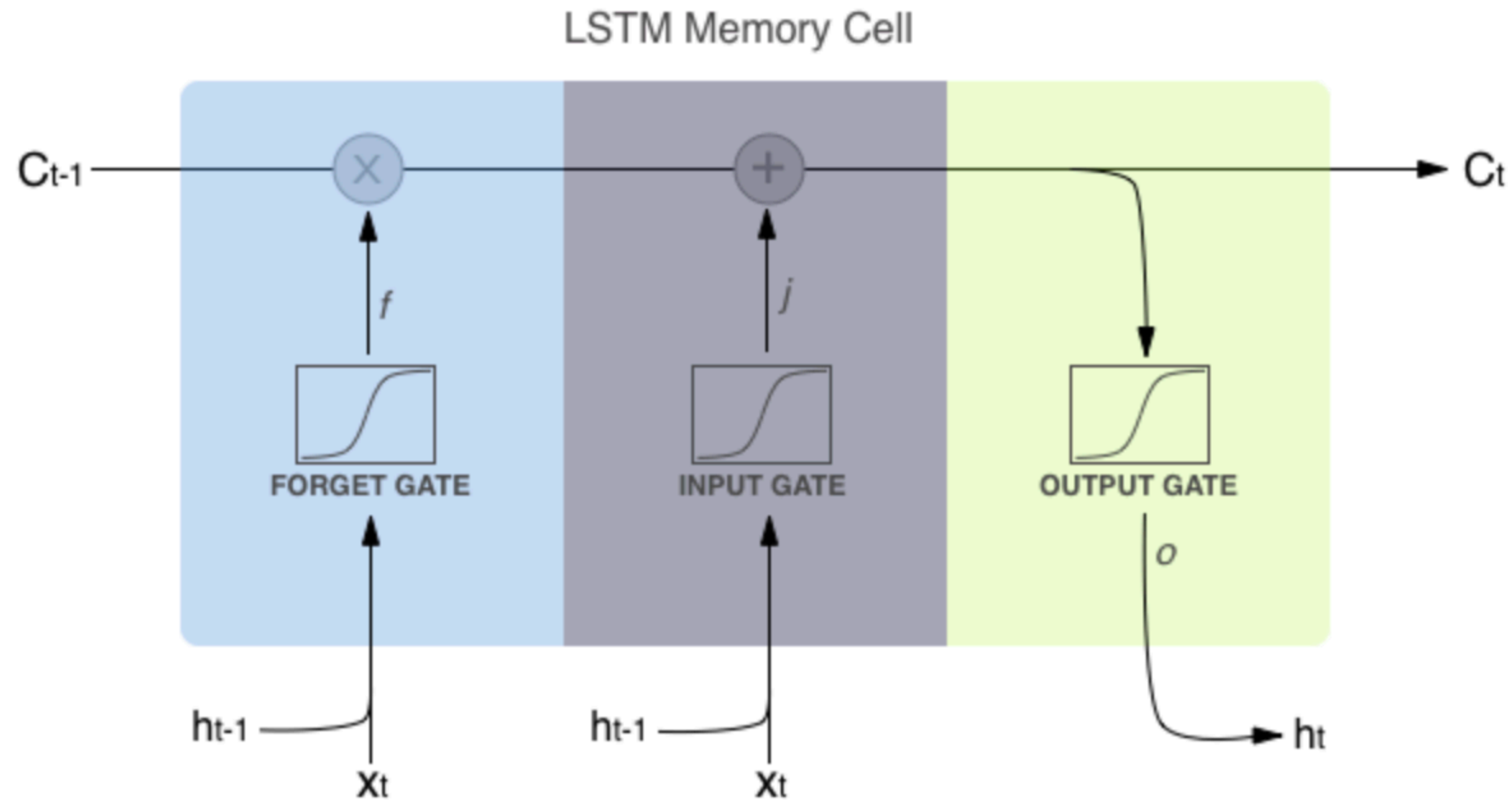
Keras calls this part "recurrent" (`recurrent_initializer,...`)



# RNNs with LSTM

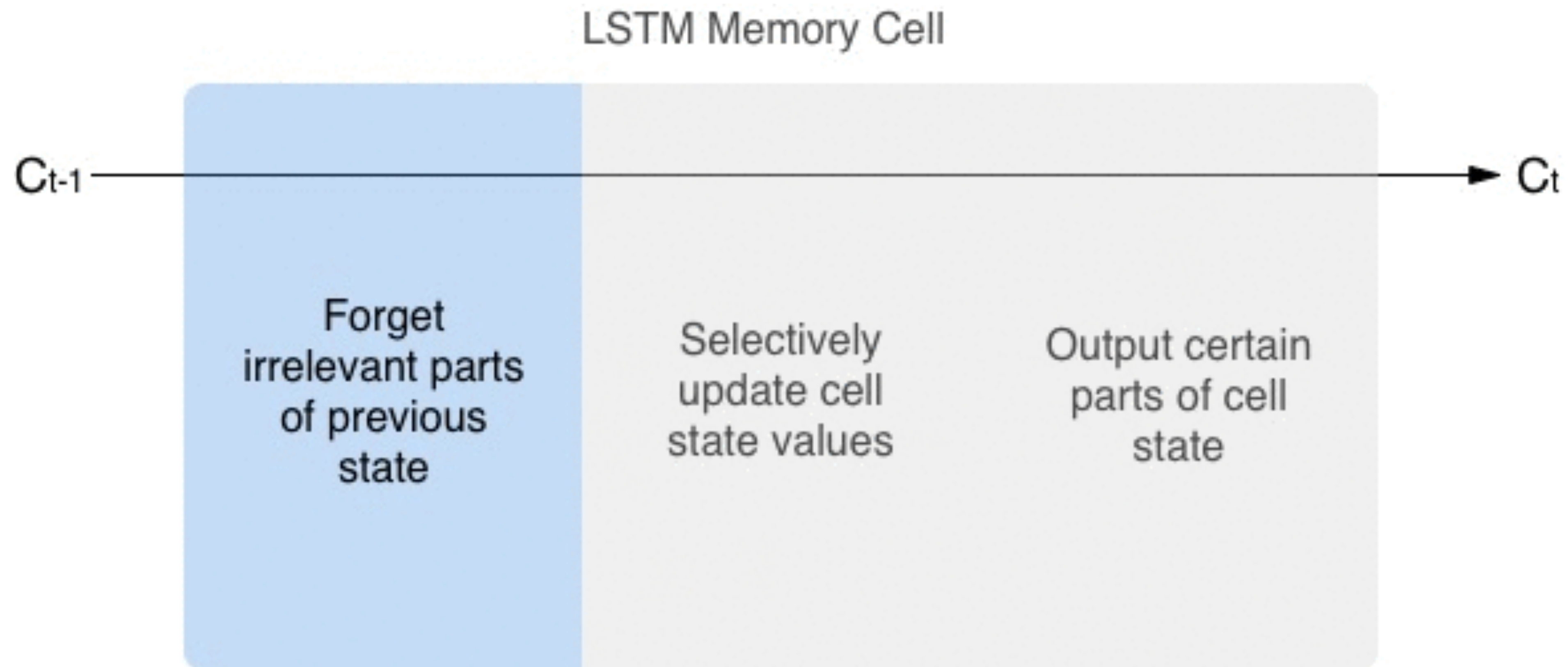
- Problem: **state transitions**
  - *It is hard to keep info from the distant past in current memory (without reinforcement)*
- Solution: **long short-term memory**

# LSTM

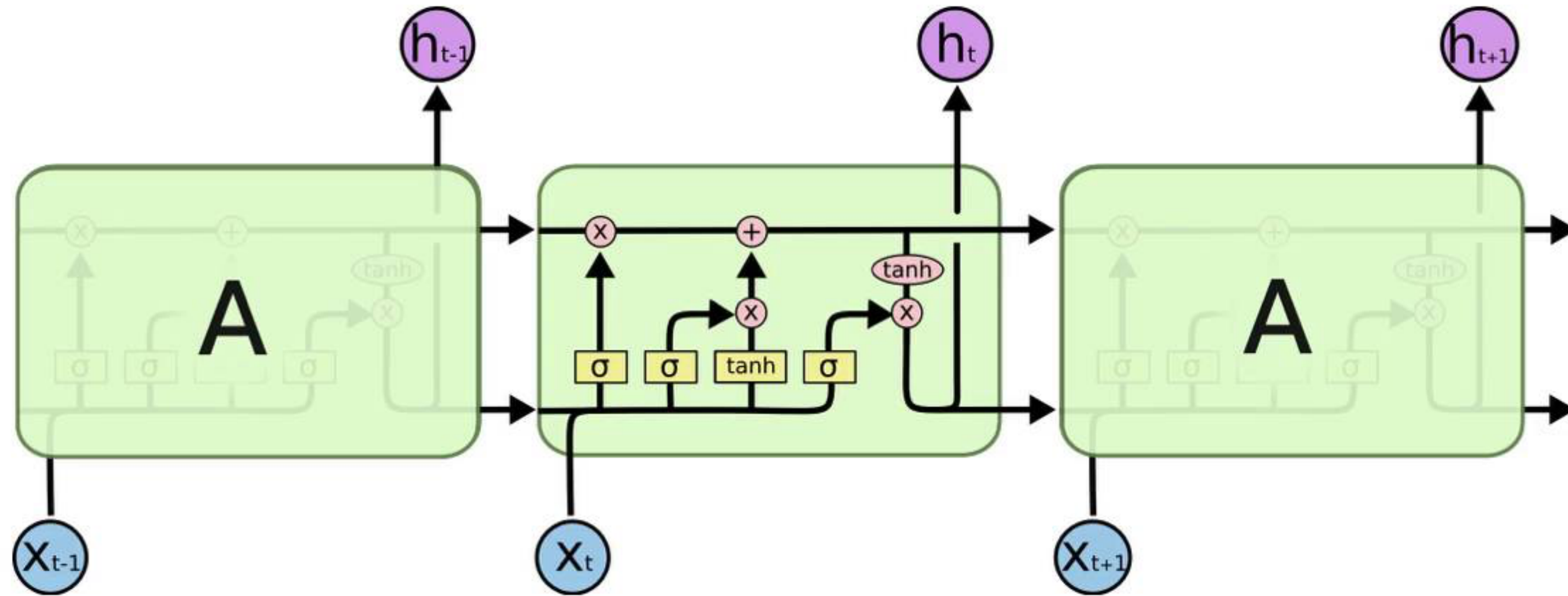




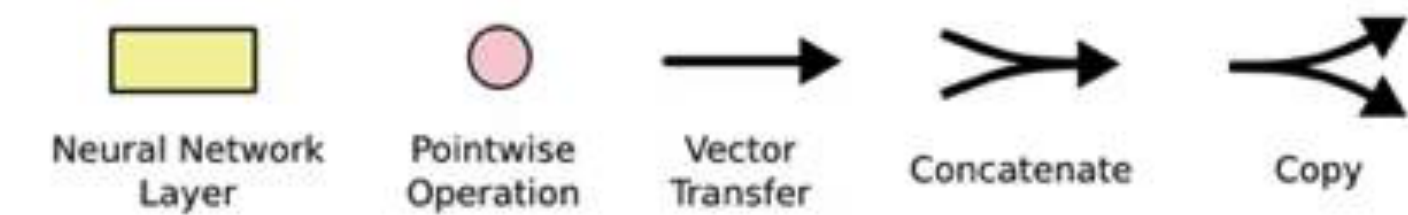
# LSTM



# LSTM

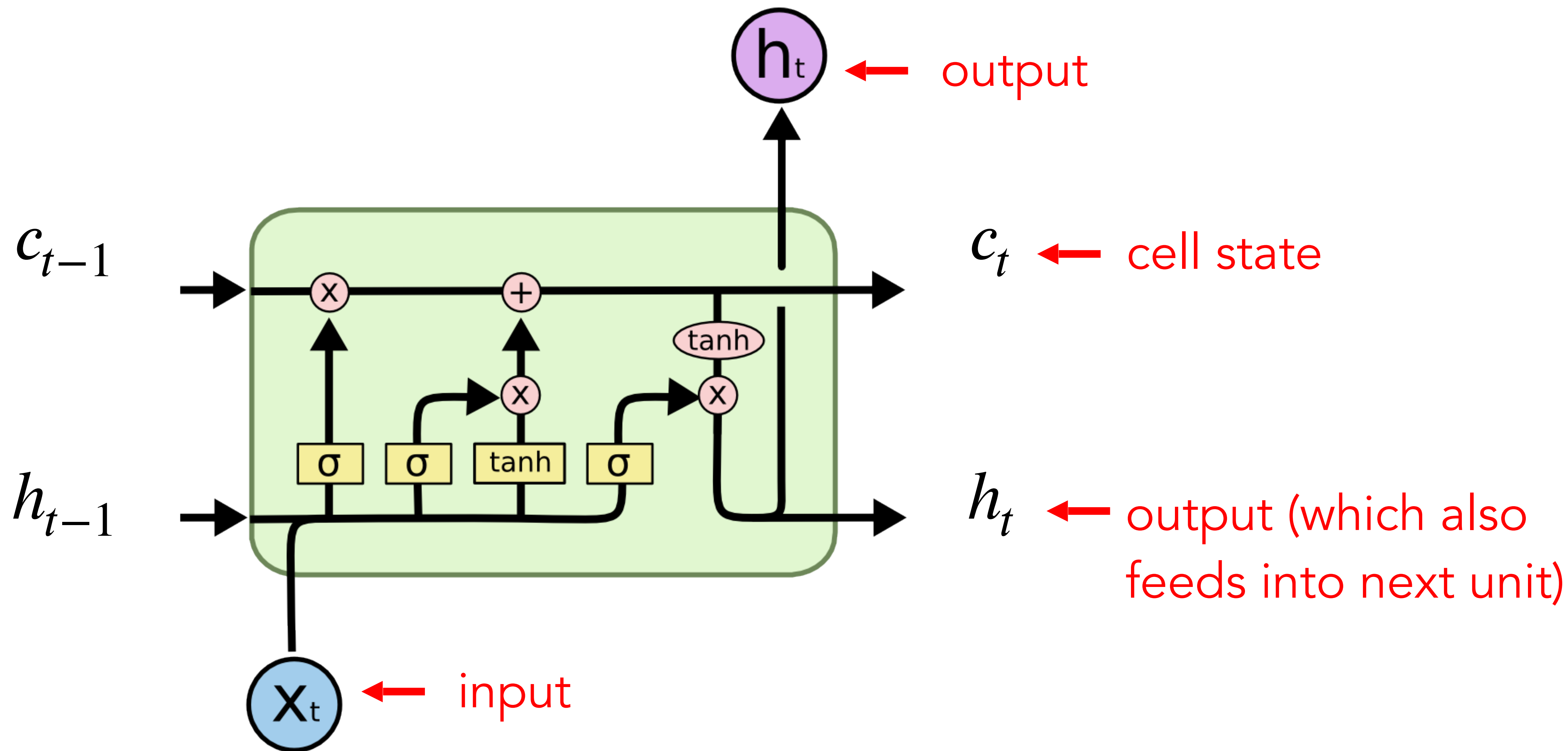


long-short term memory modules used in an RNN



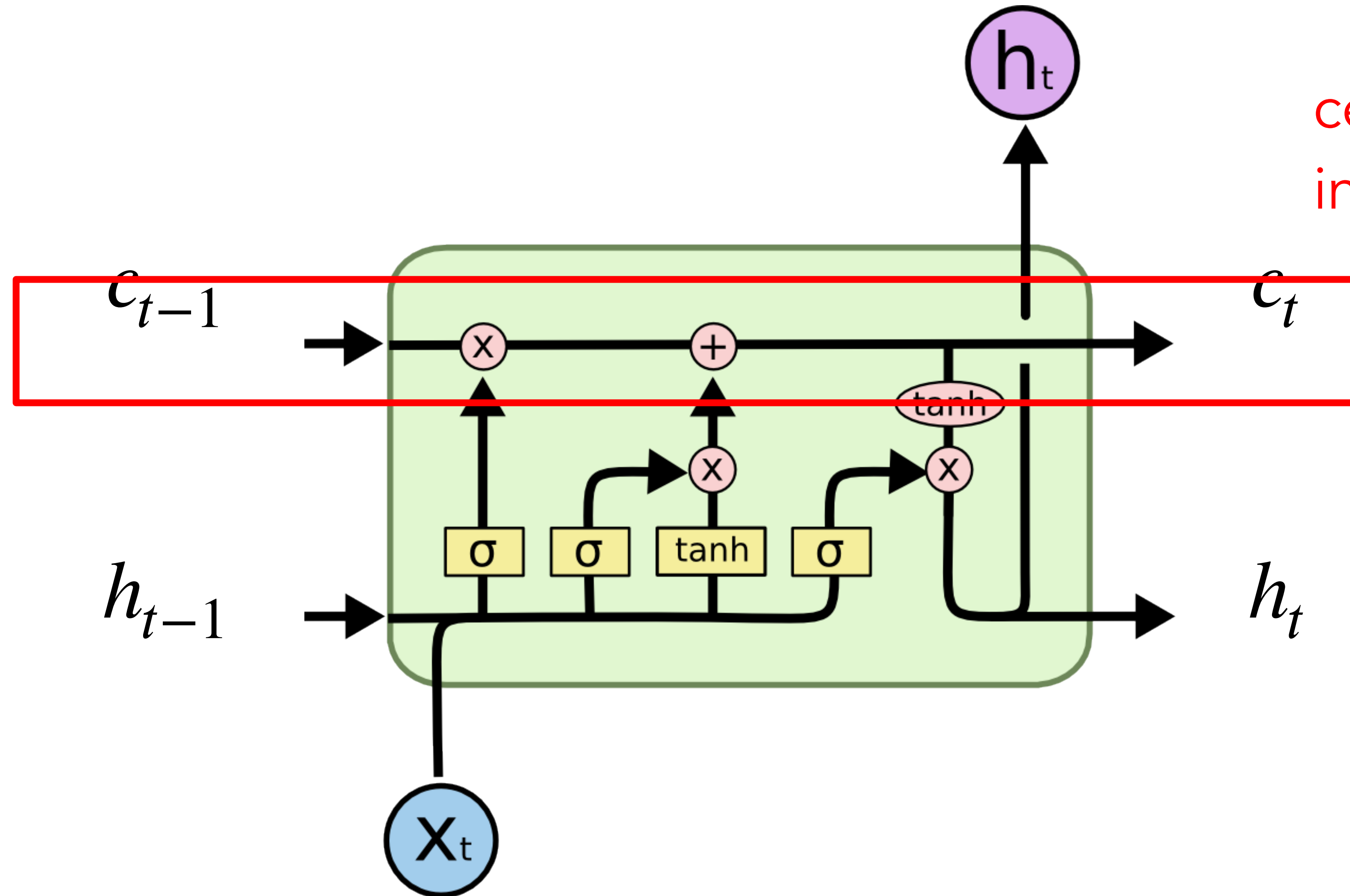
<http://colah.github.io/posts/2015-08-Understanding-LSTMs/> Eugenio Culurciello  
© 2016

# LSTM diagram



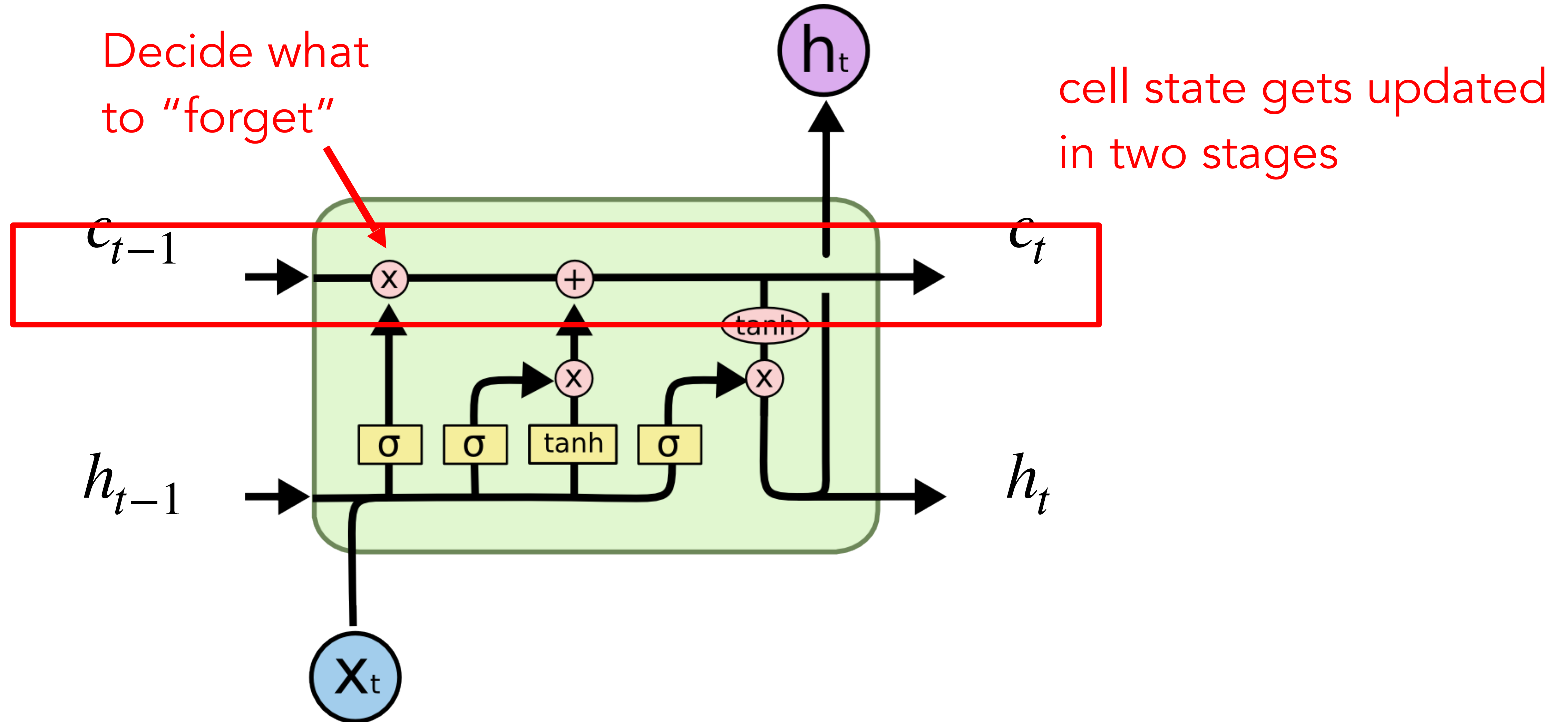


# LSTM diagram

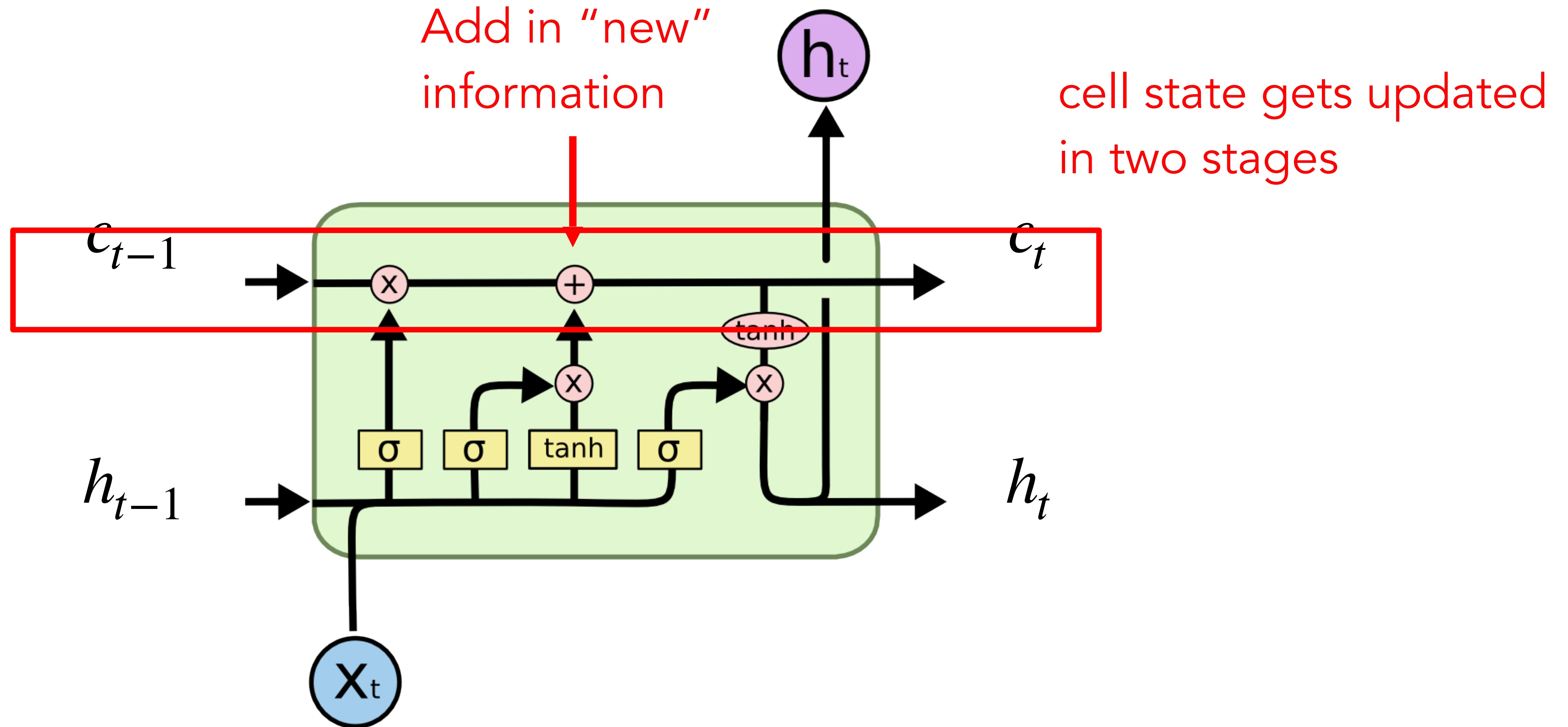


cell state gets updated  
in two stages

# LSTM diagram

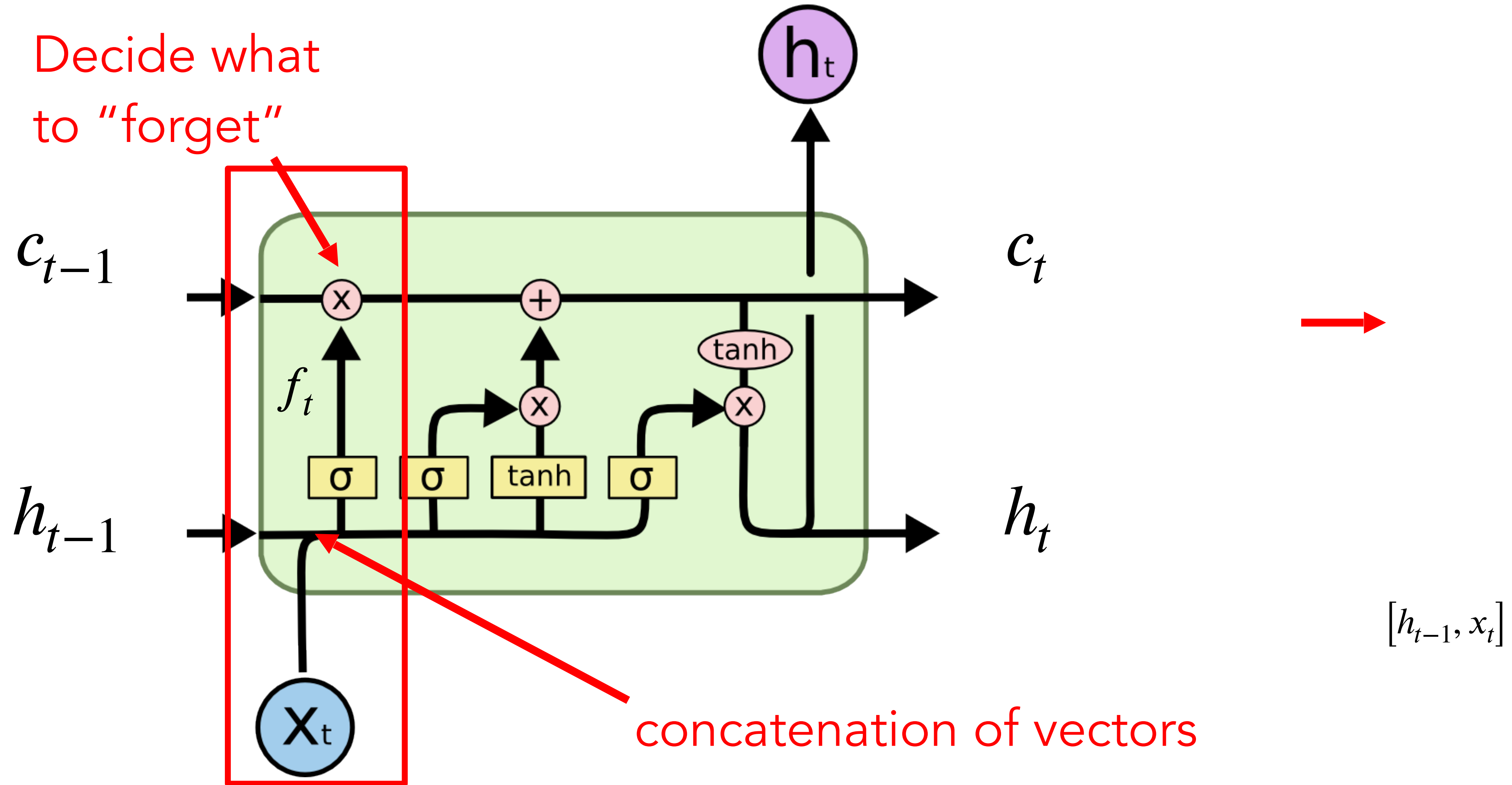


# LSTM diagram

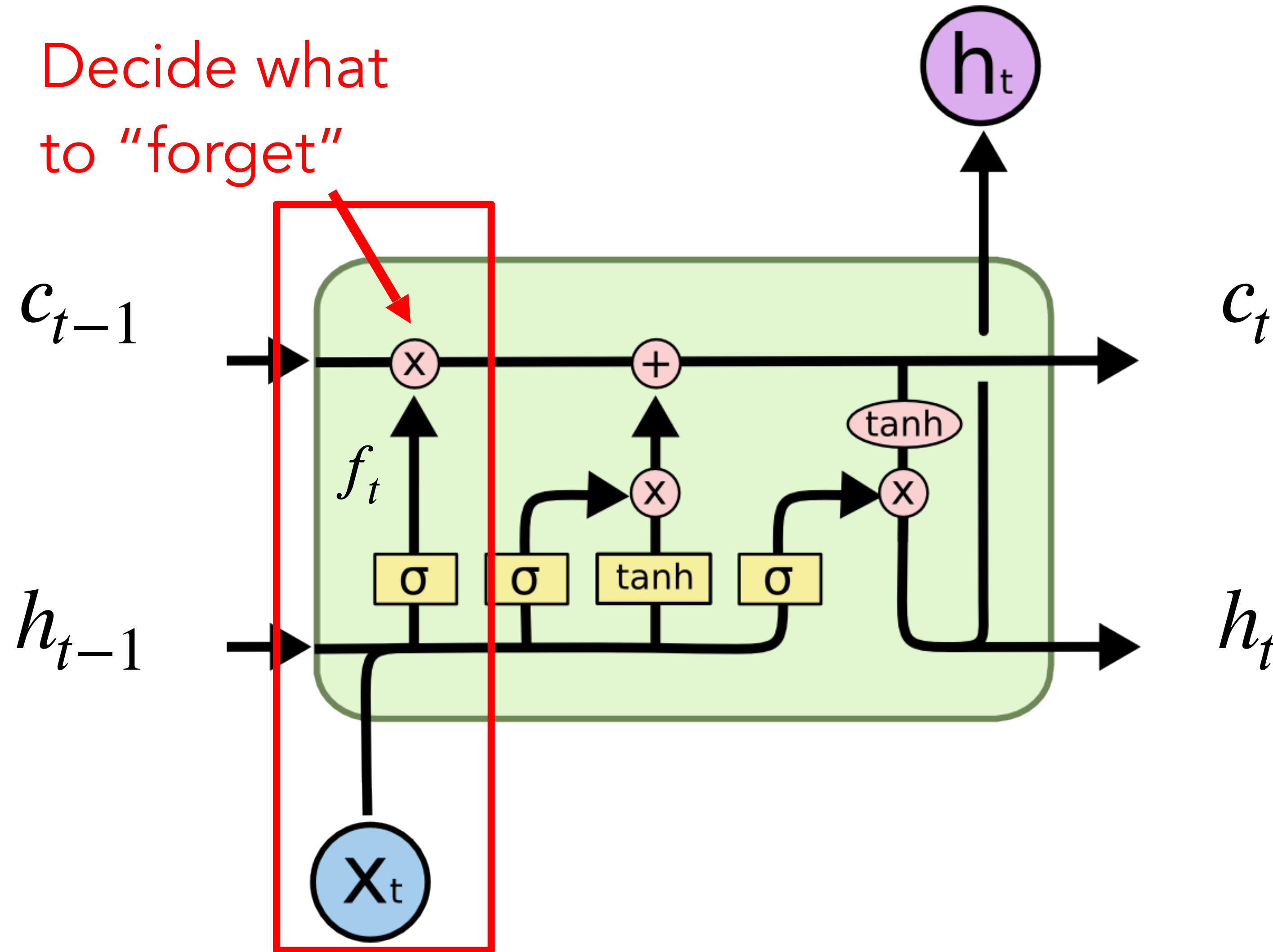




# LSTM diagram



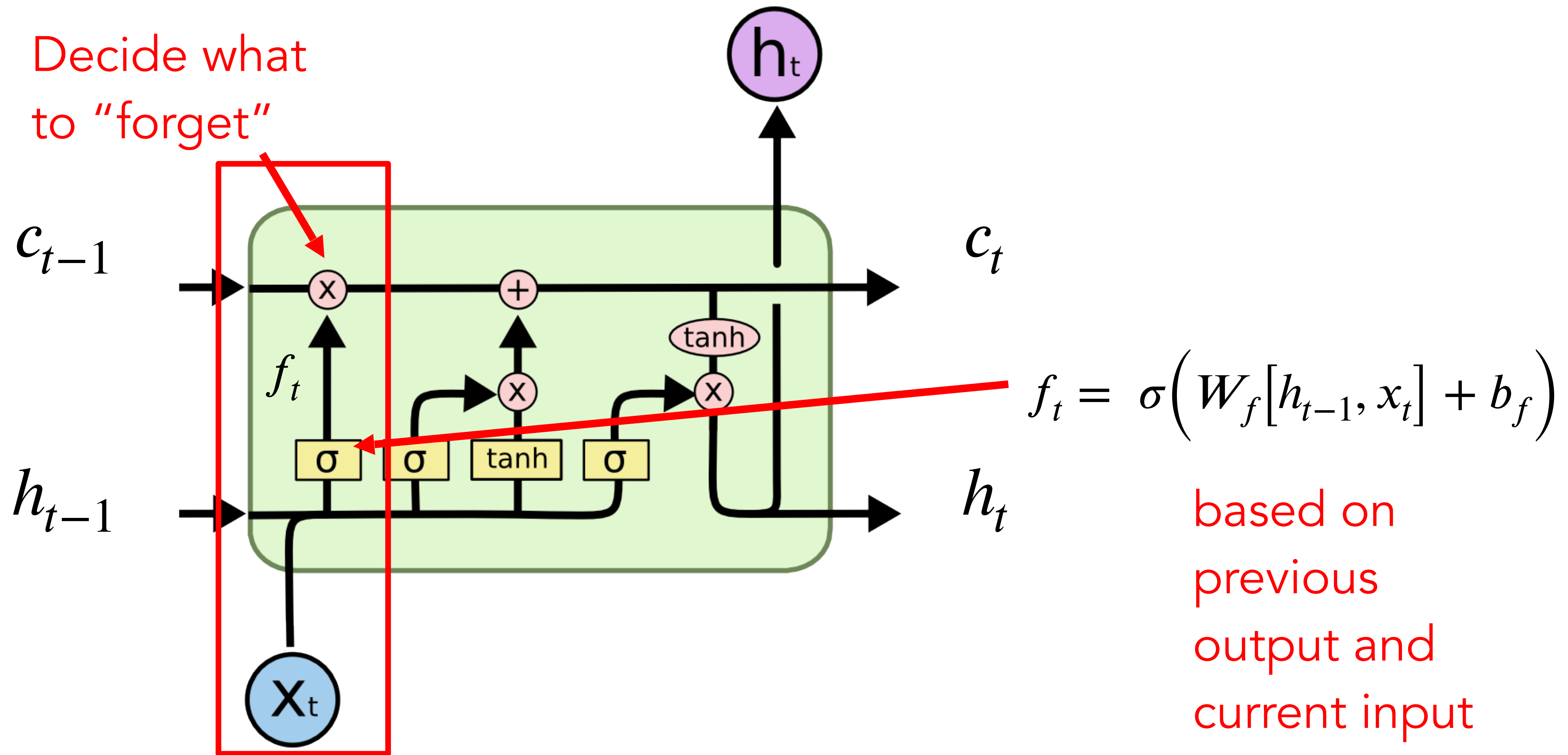
# LSTM diagram



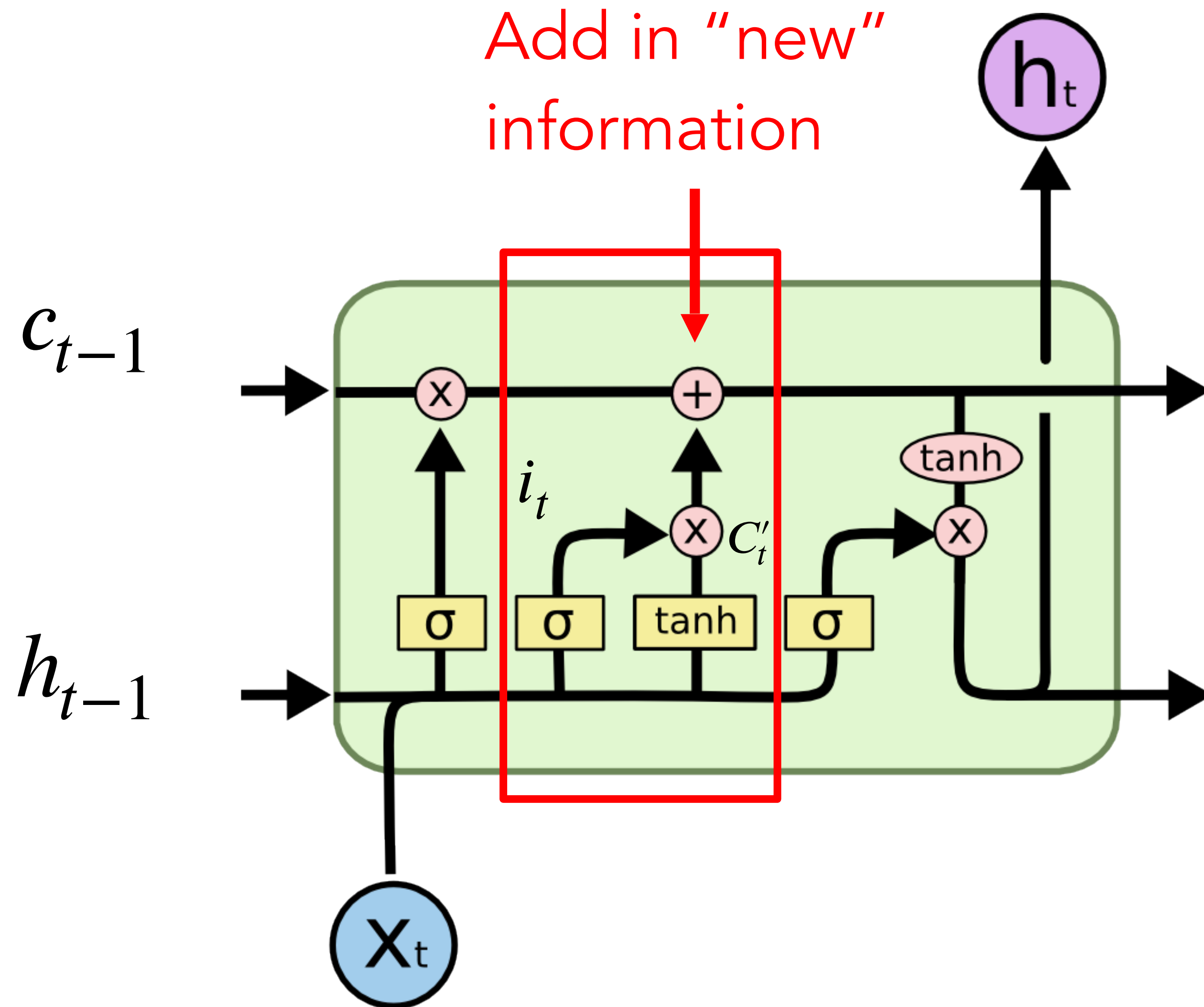
$$f_t = \sigma(W_f[h_{t-1}, x_t] + b_f)$$

based on  
previous  
output and  
current input

# LSTM diagram

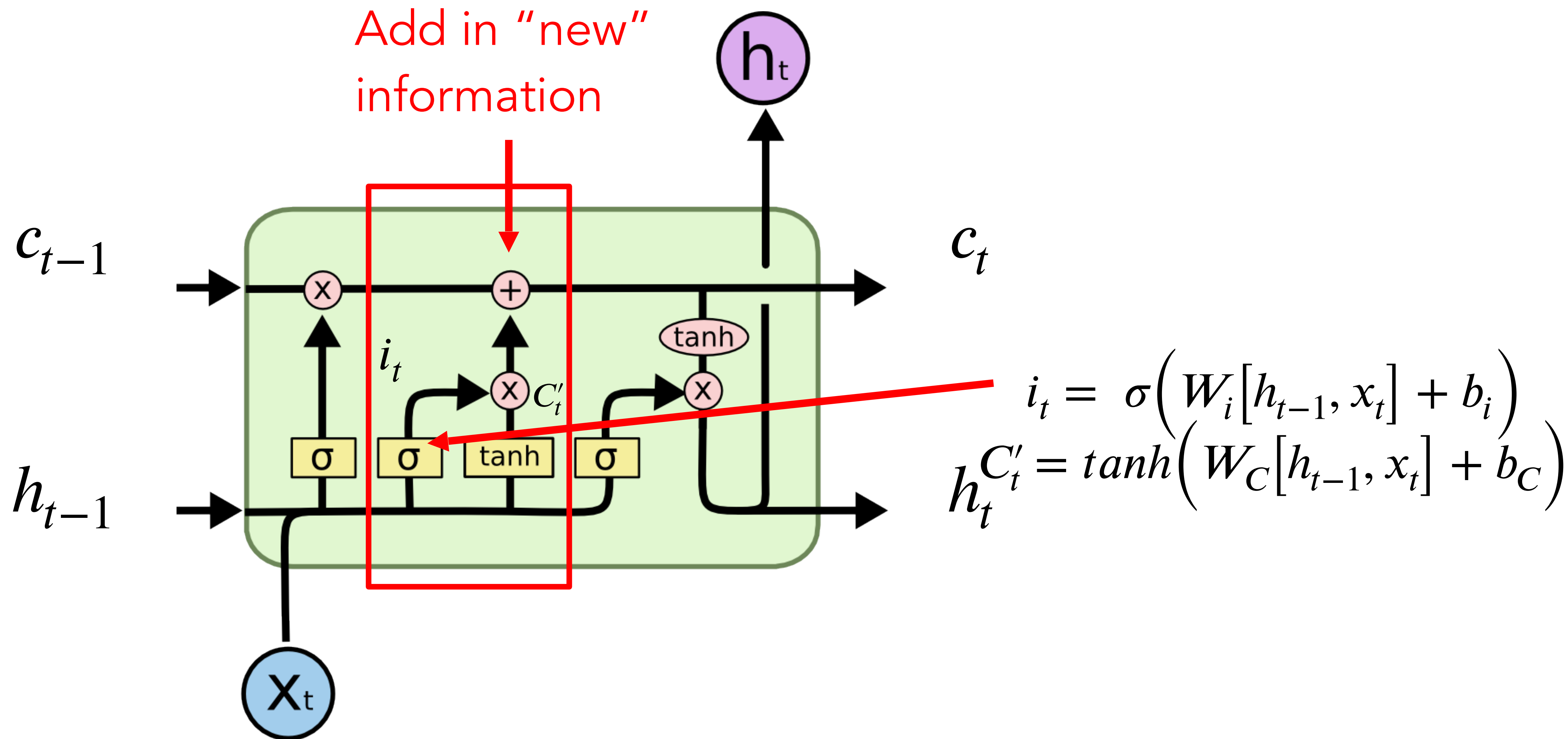


# LSTM diagram



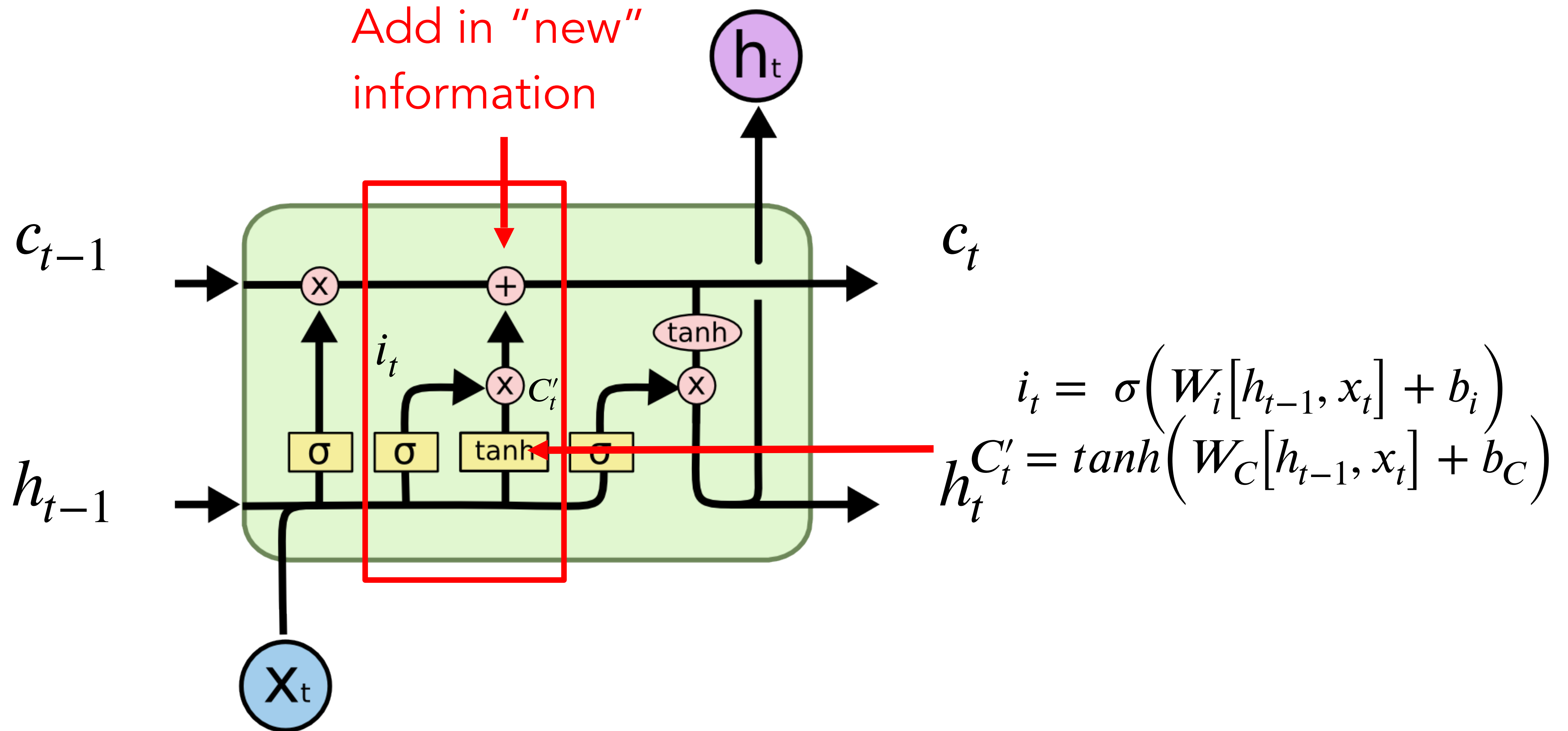
$$i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$$
$$h_t C'_t = \tanh(W_C[h_{t-1}, x_t] + b_C)$$

# LSTM diagram

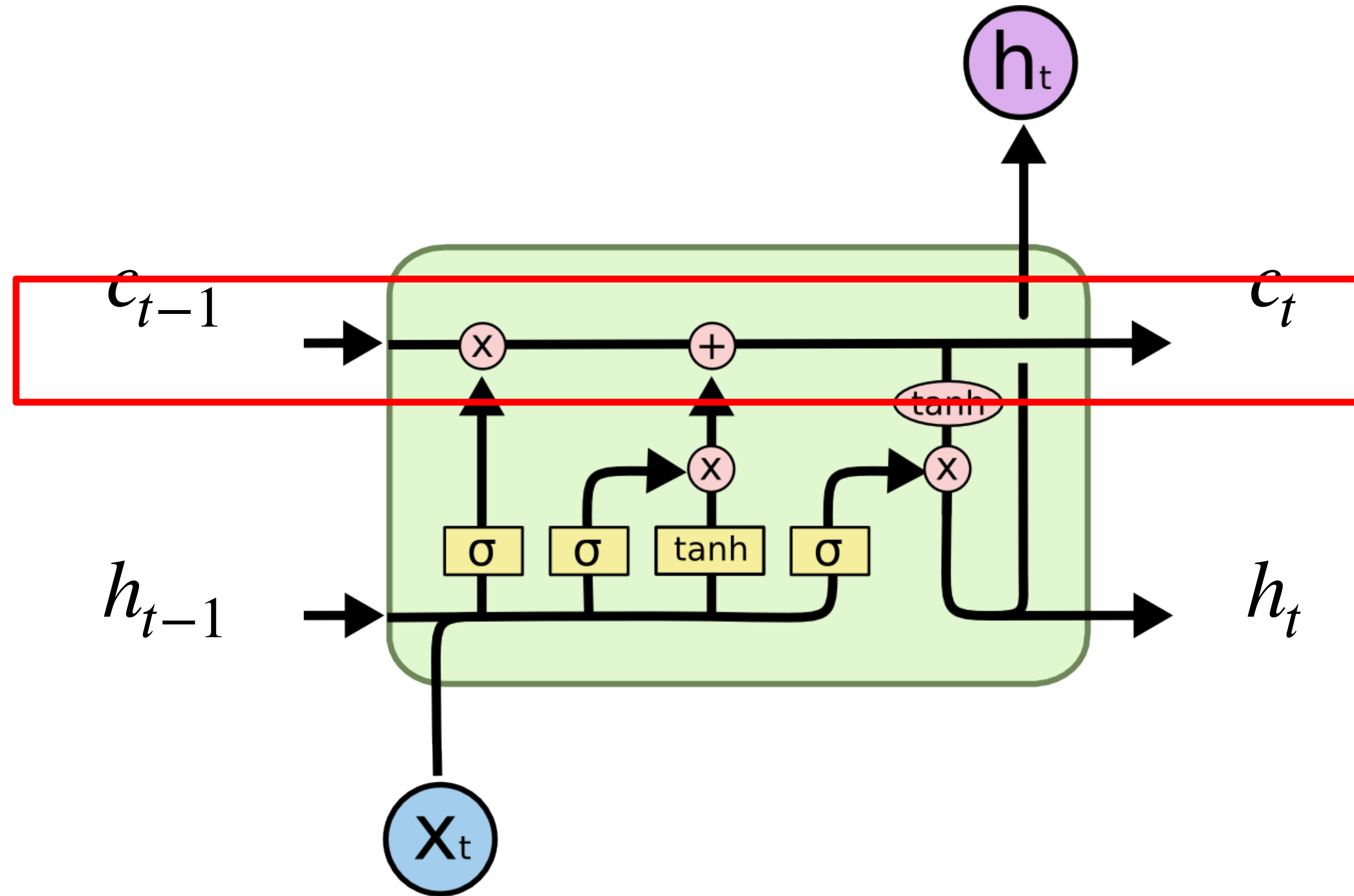




# LSTM diagram



# LSTM diagram



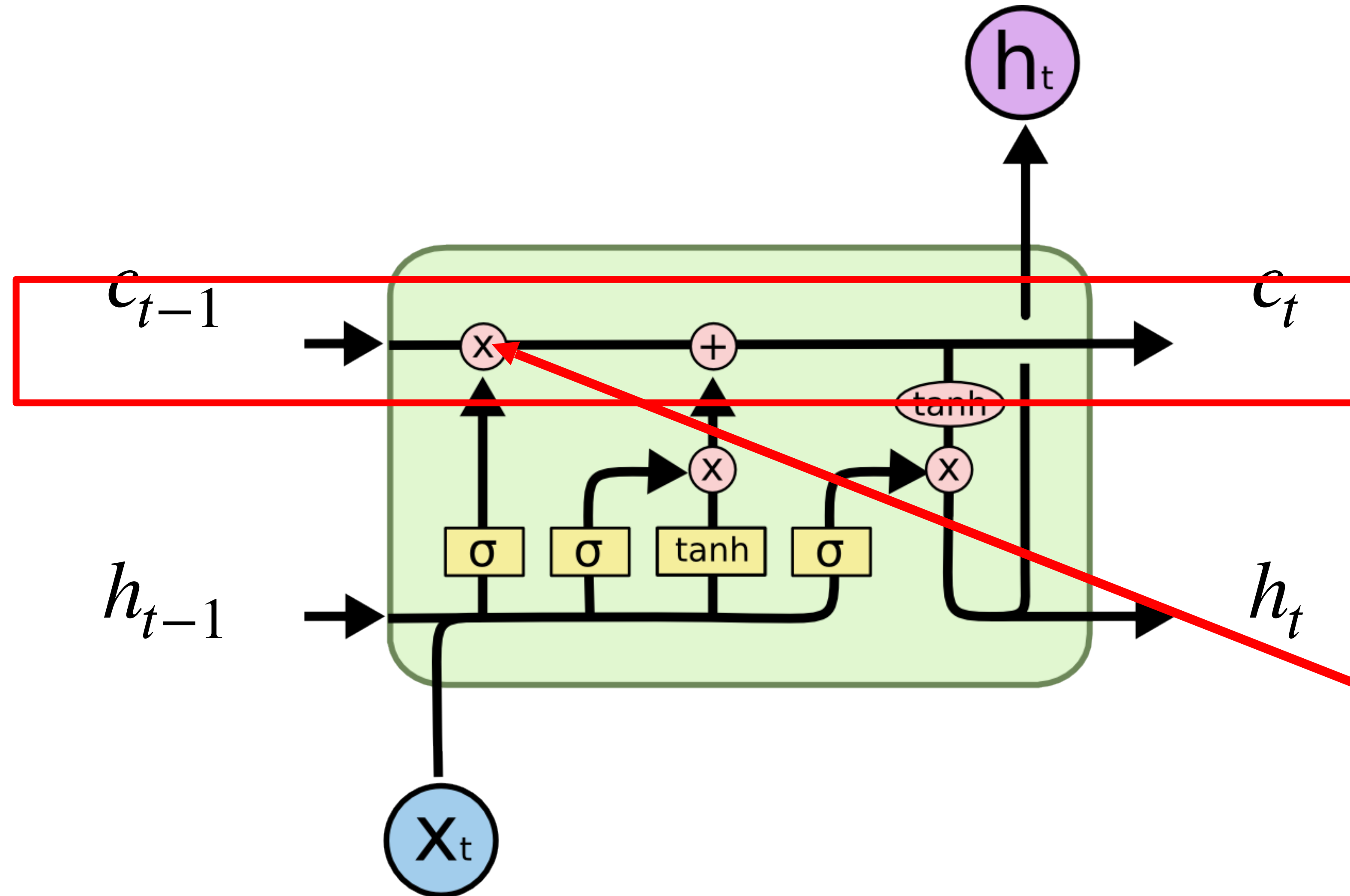
Note: 'x' represents element-wise multiplication

$$C_t = f_i * C_{t-1} + i_t * C'_t$$

forget  
the old  
(or not)

add the  
new (or  
not)

# LSTM diagram



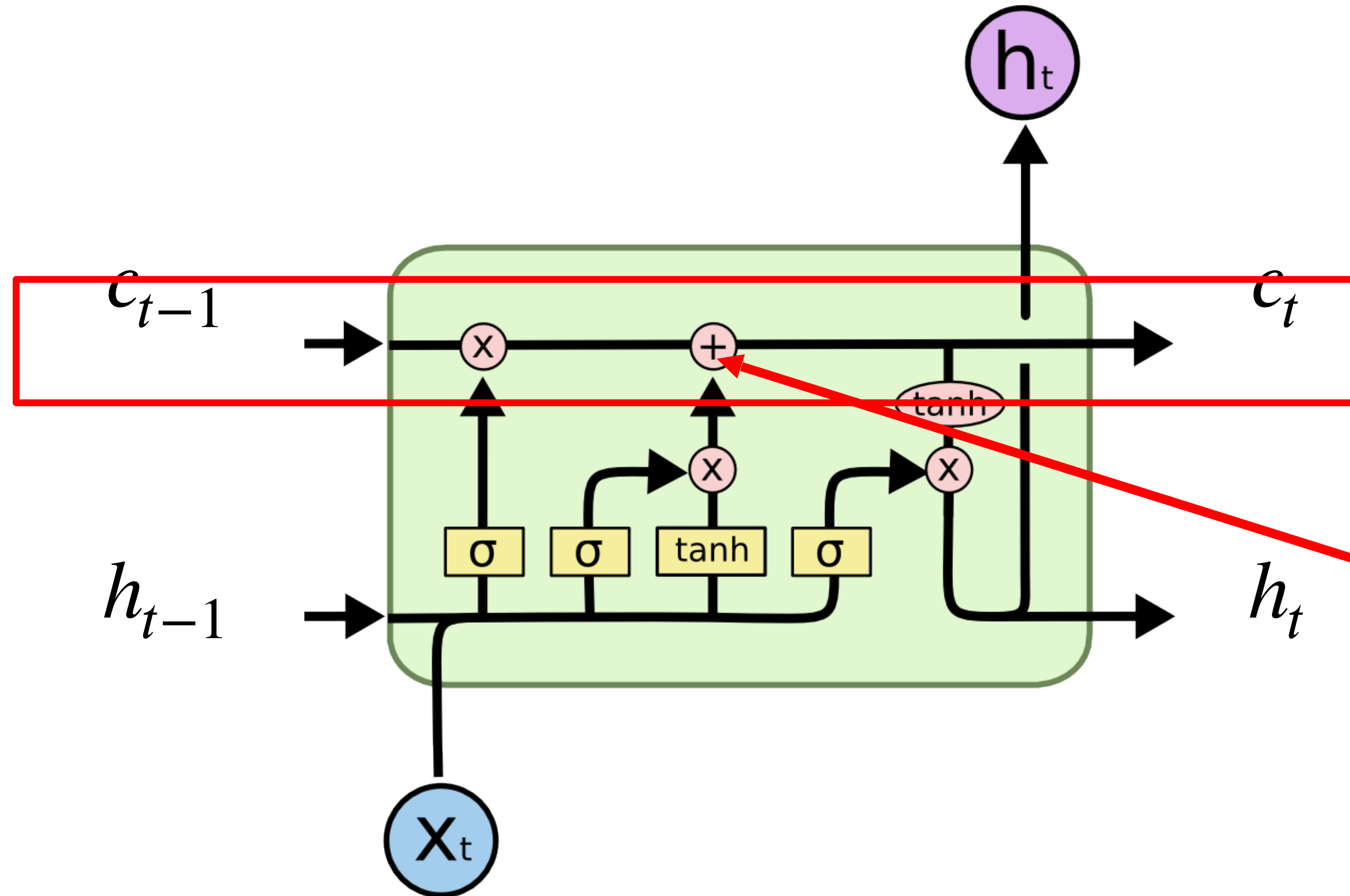
Note: 'x' represents element-wise multiplication

$$C_t = f_i * C_{t-1} + i_t * C'_t$$

forget  
the old  
(or not)

add the  
new (or  
not)

# LSTM diagram



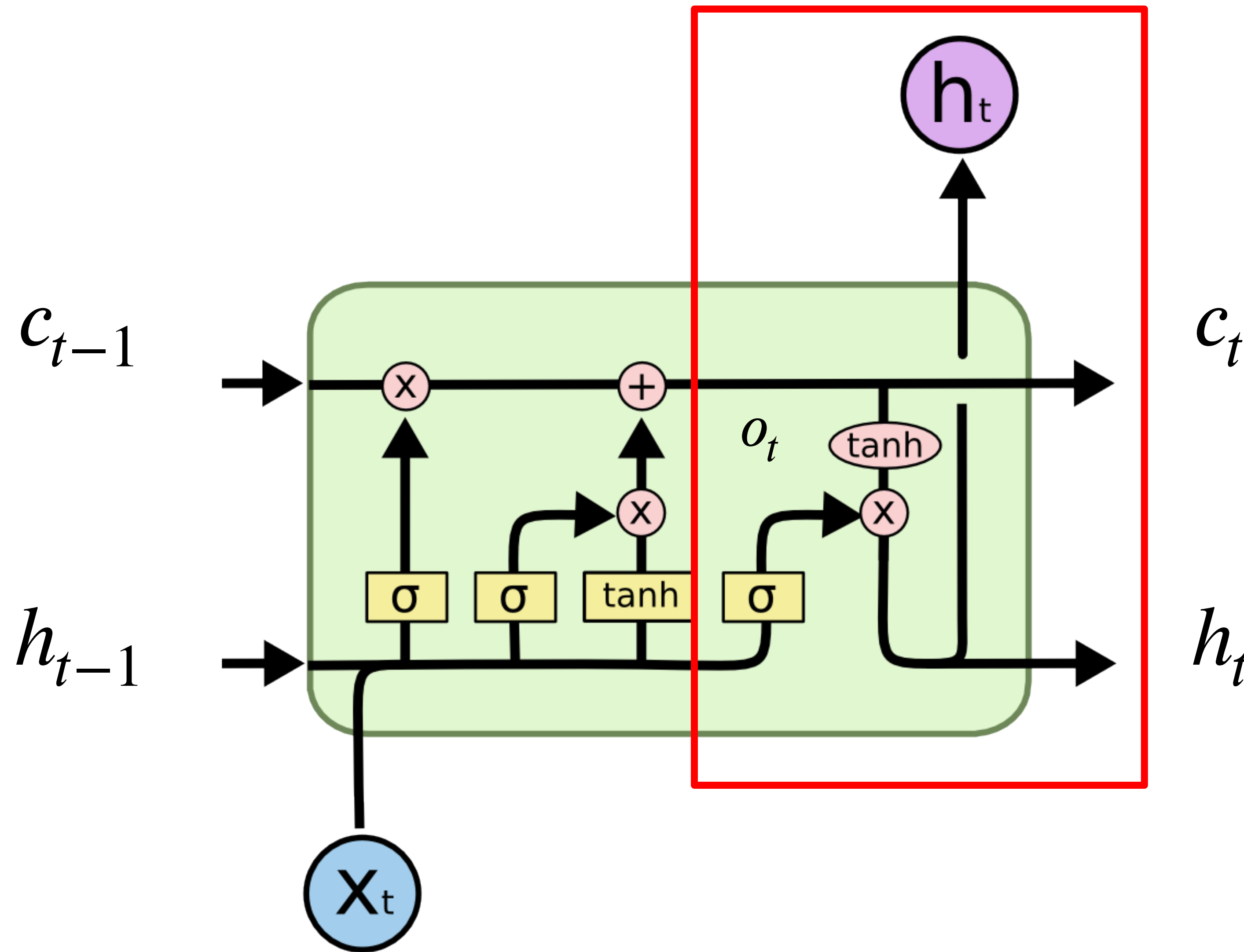
Note: " represents element-wise multiplication

$$C_t = f_i * C_{t-1} + i_t * C'_t$$

forget  
the old  
(or not)

add the  
new (or  
not)

# LSTM diagram

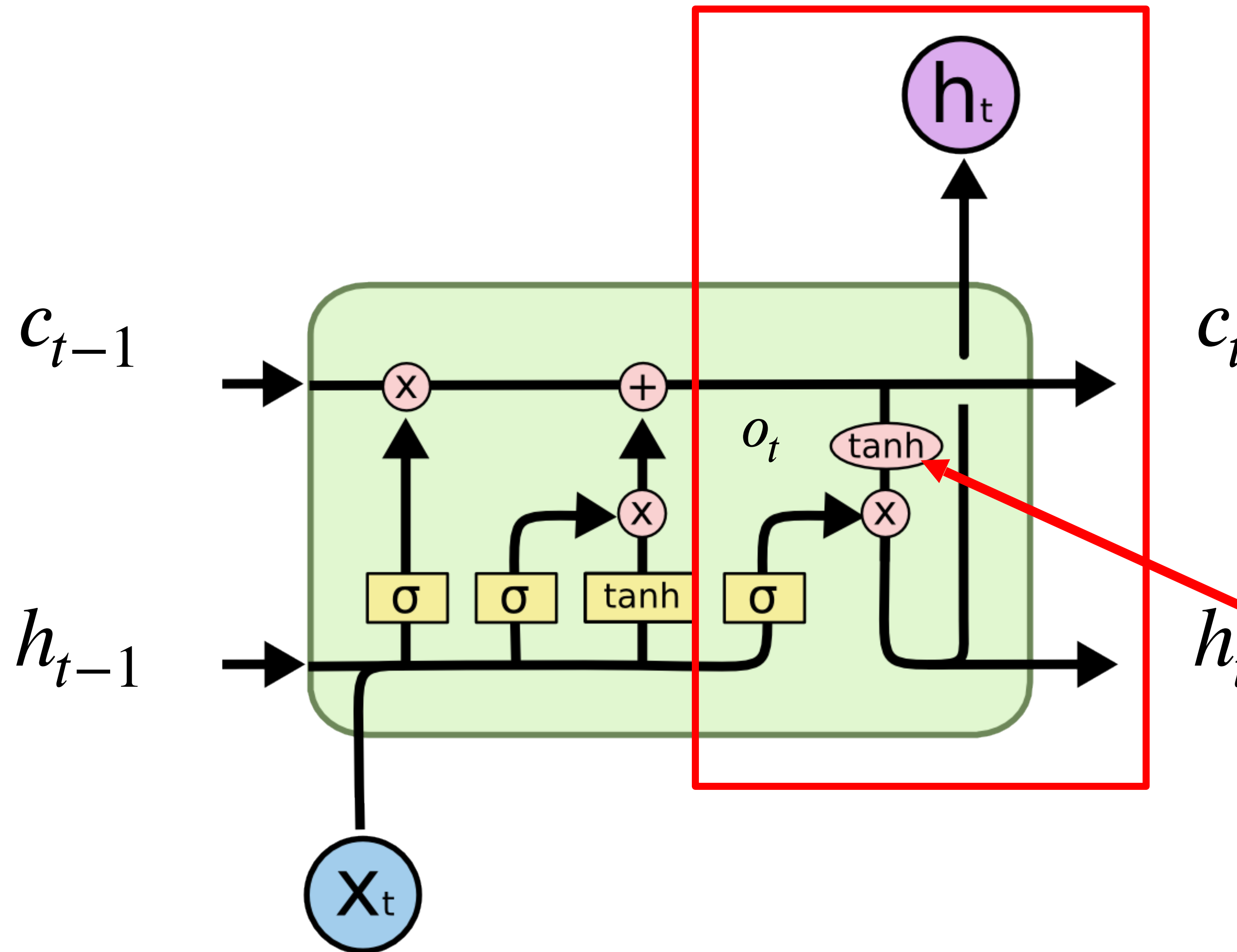


Final stage computes the output

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



# LSTM diagram



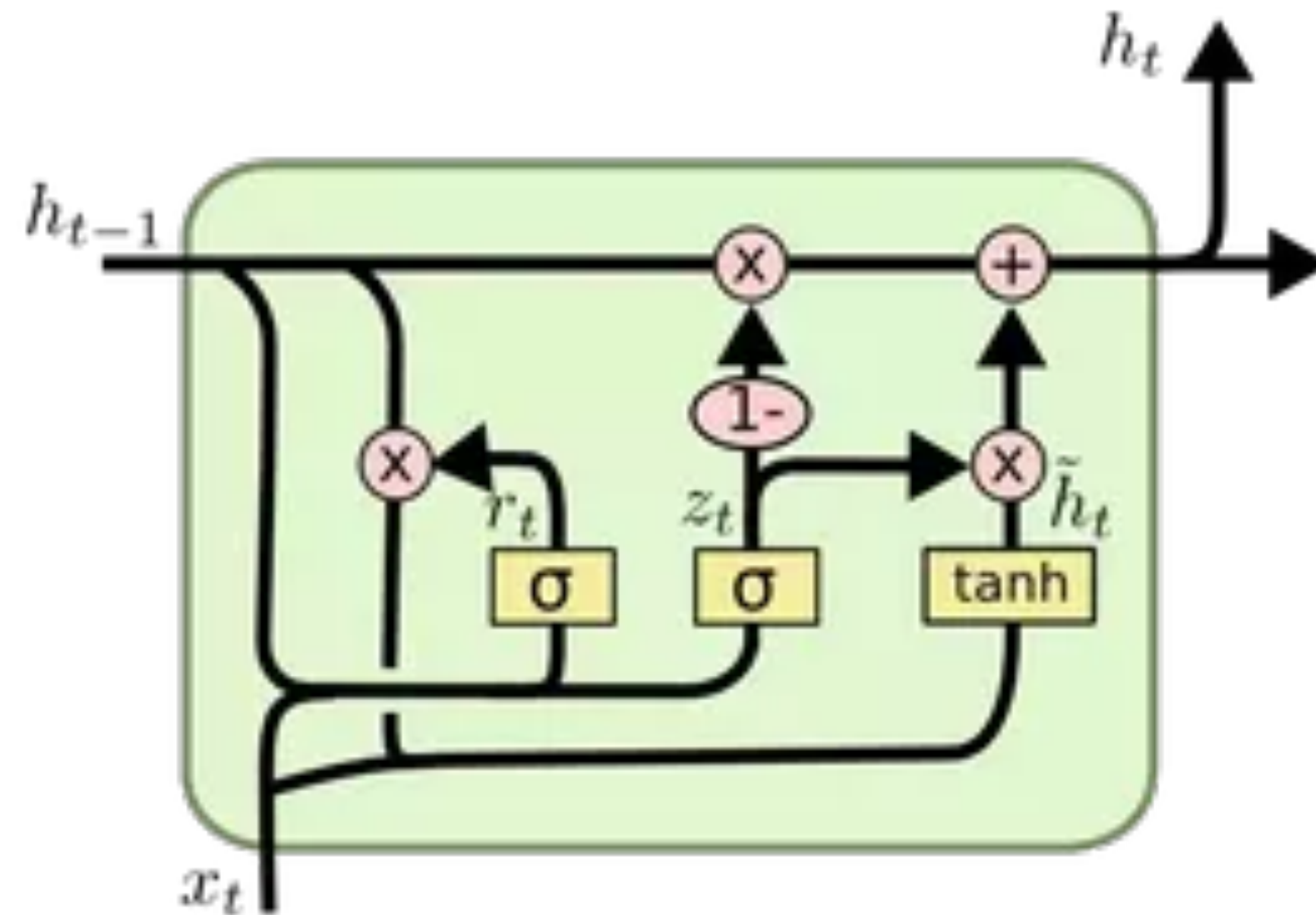
Final stage computes the output

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

Note: No weights here

# Gated Recurrent Unit (GRU)

- Gated Recurrent Unit (GRU)



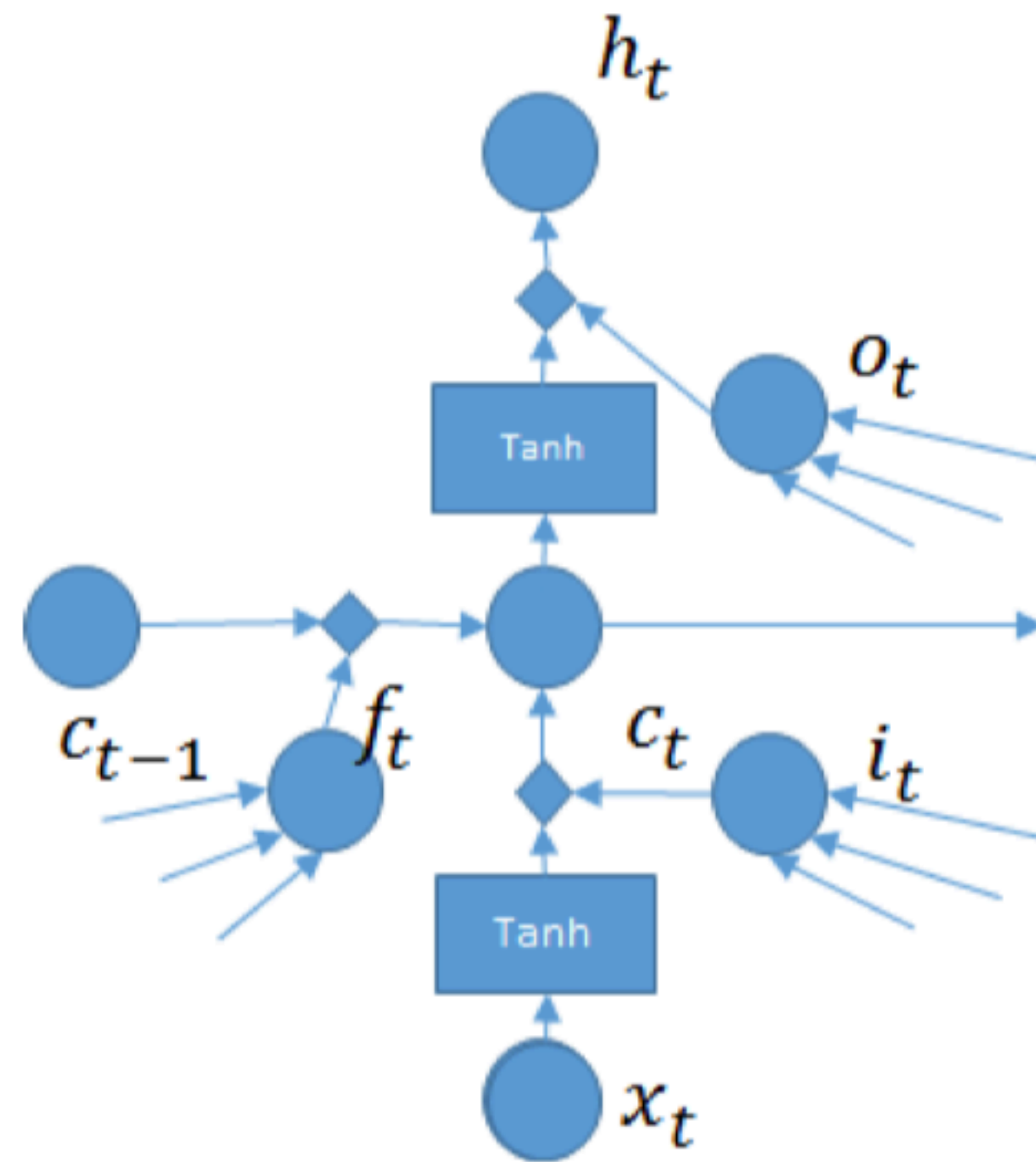
$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

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

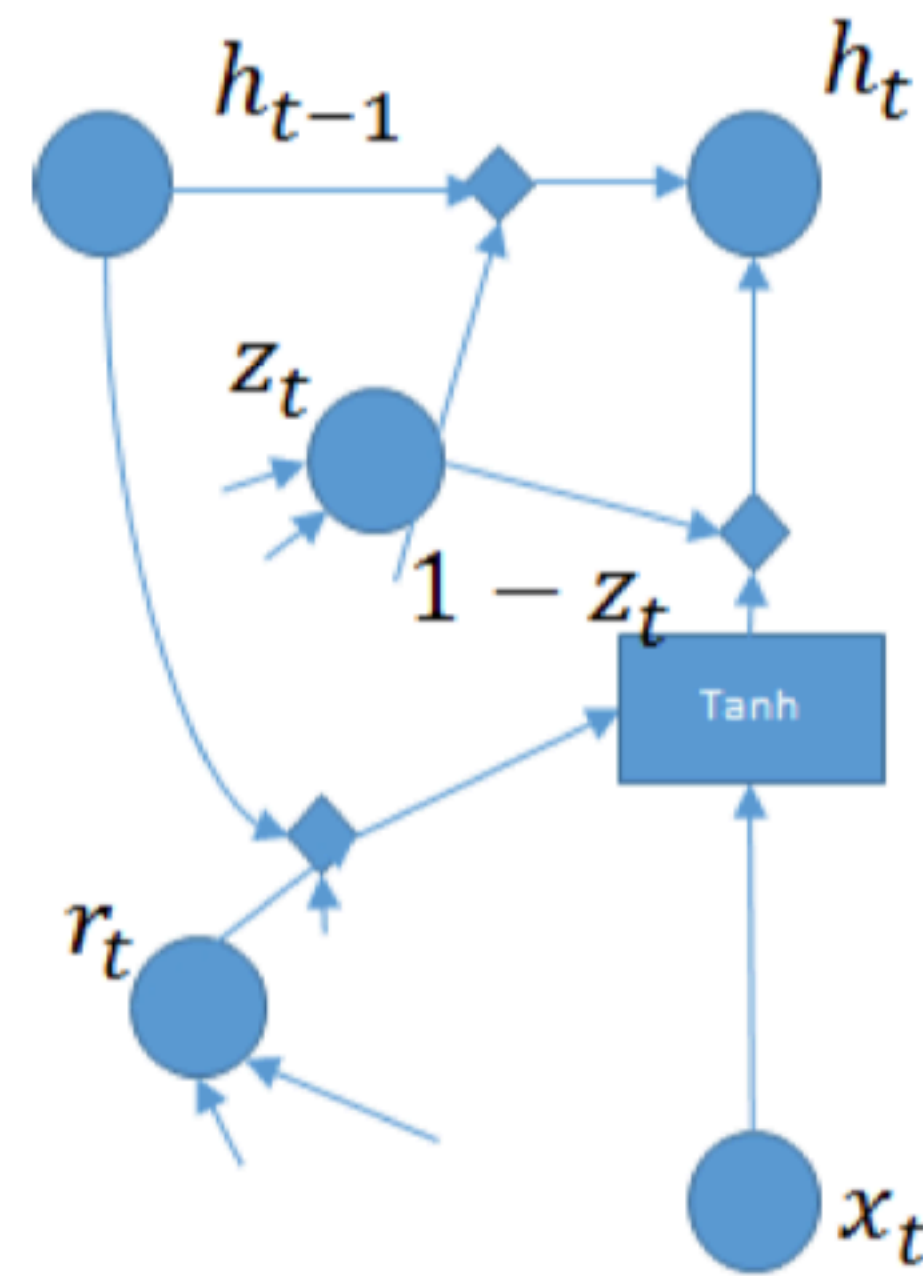
$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

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

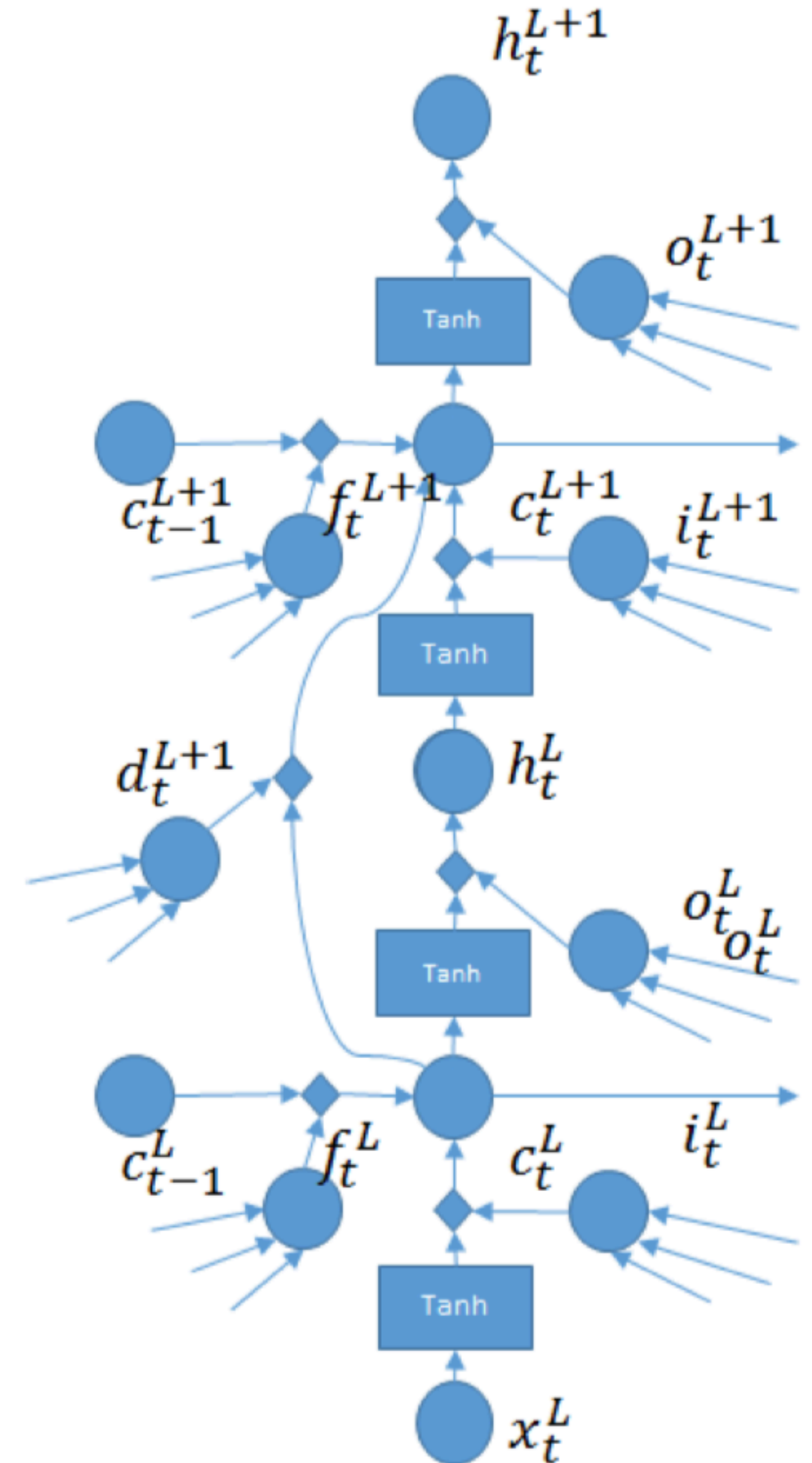
# Variations of LSTM



**LSTM**



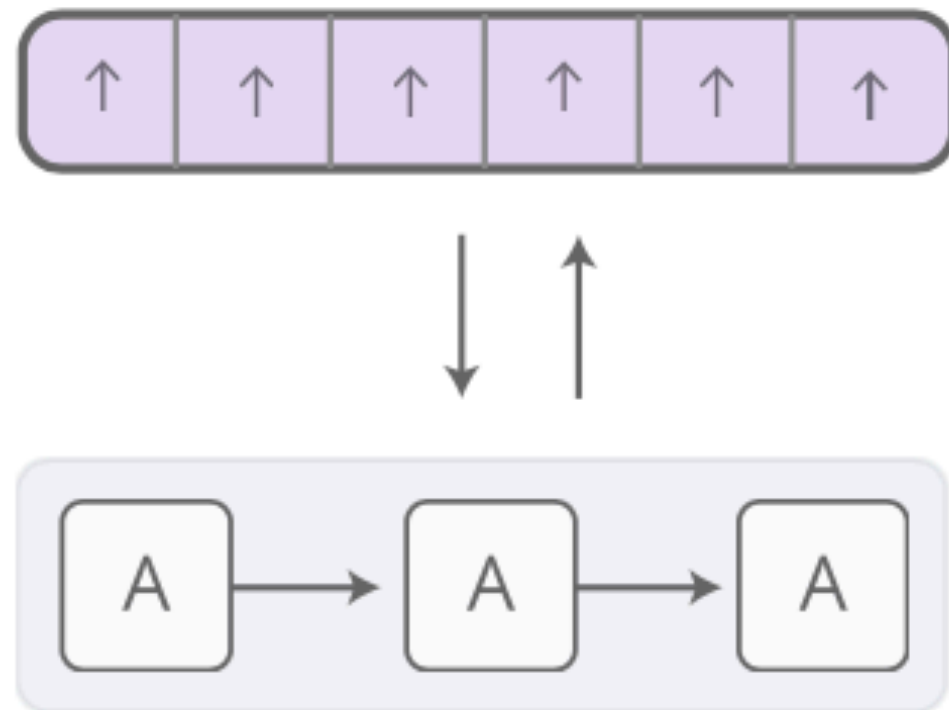
**GRU**



**Depth-Gated LSTM**

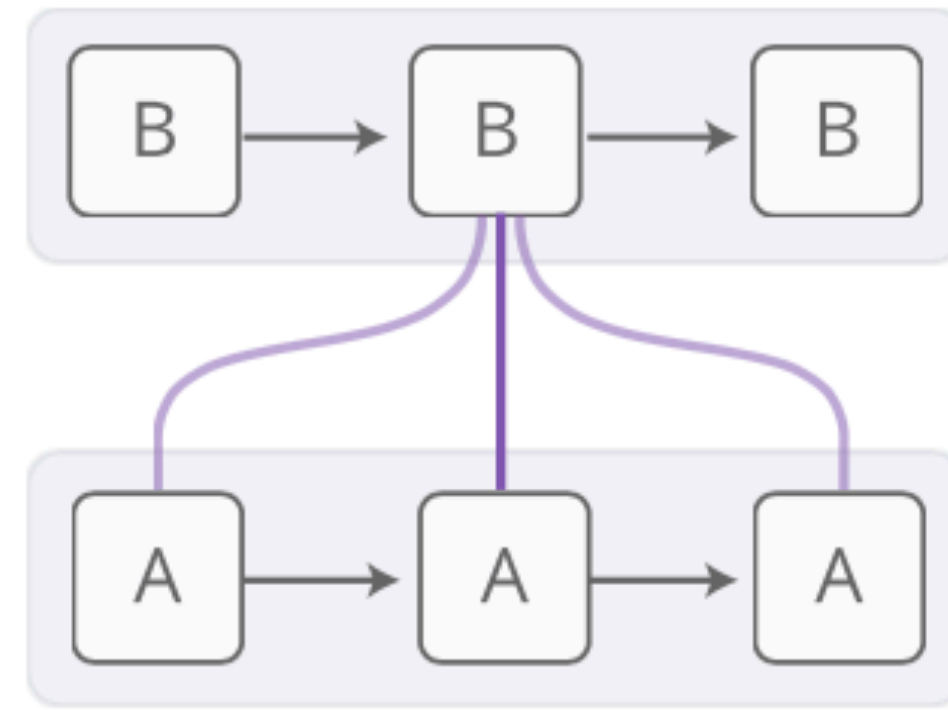
# Improving RNNs

# Augmenting RNNs



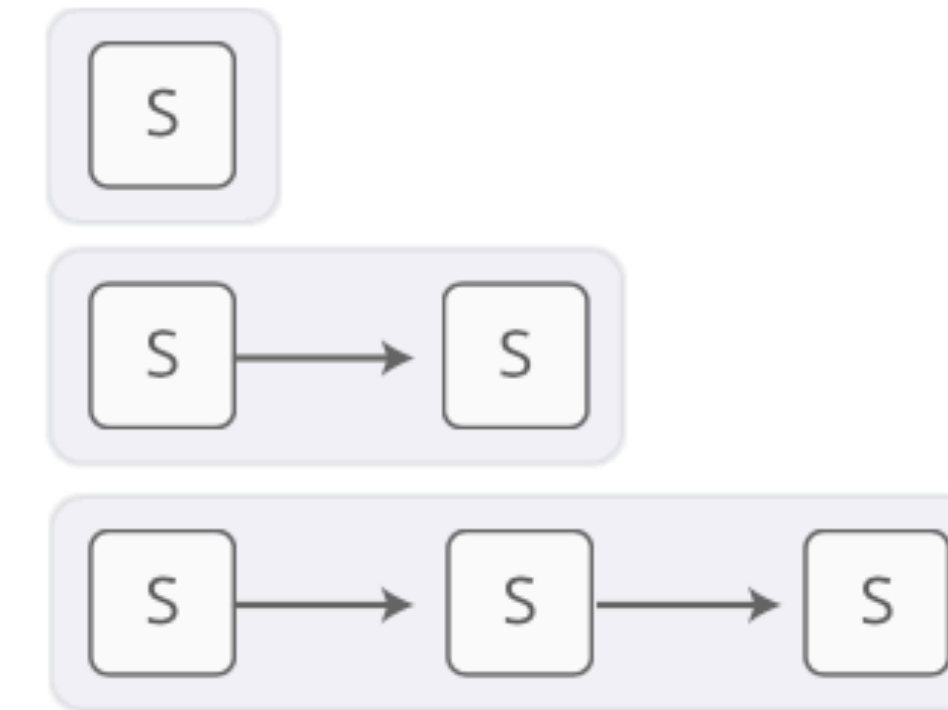
## Neural Turing Machines

have external memory that they can read and write to.



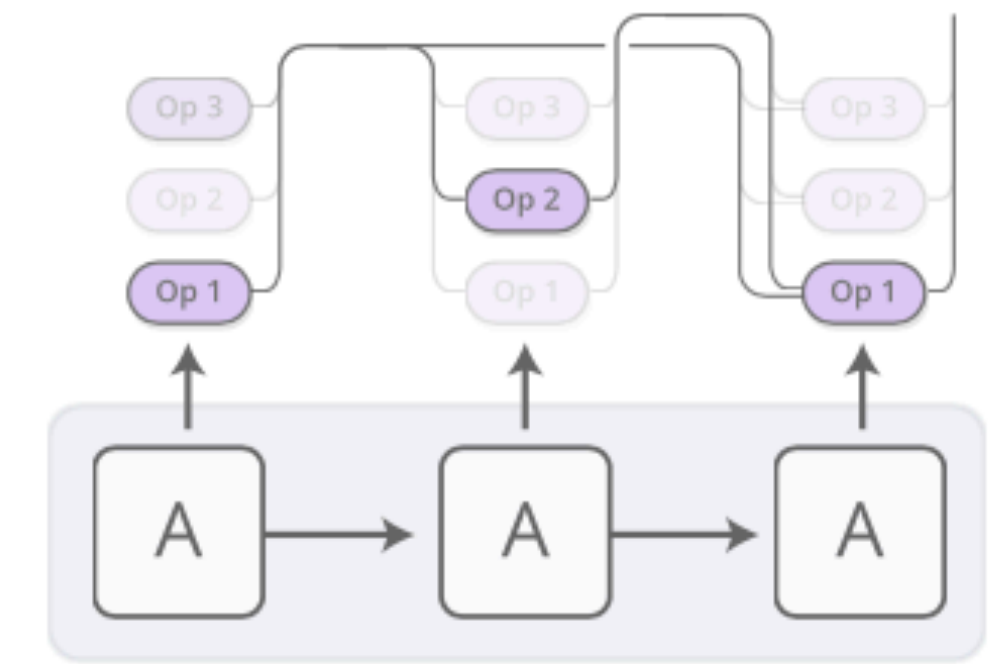
## Attentional Interfaces

allow RNNs to focus on parts of their input.



## Adaptive Computation Time

allows for varying amounts of computation per step.



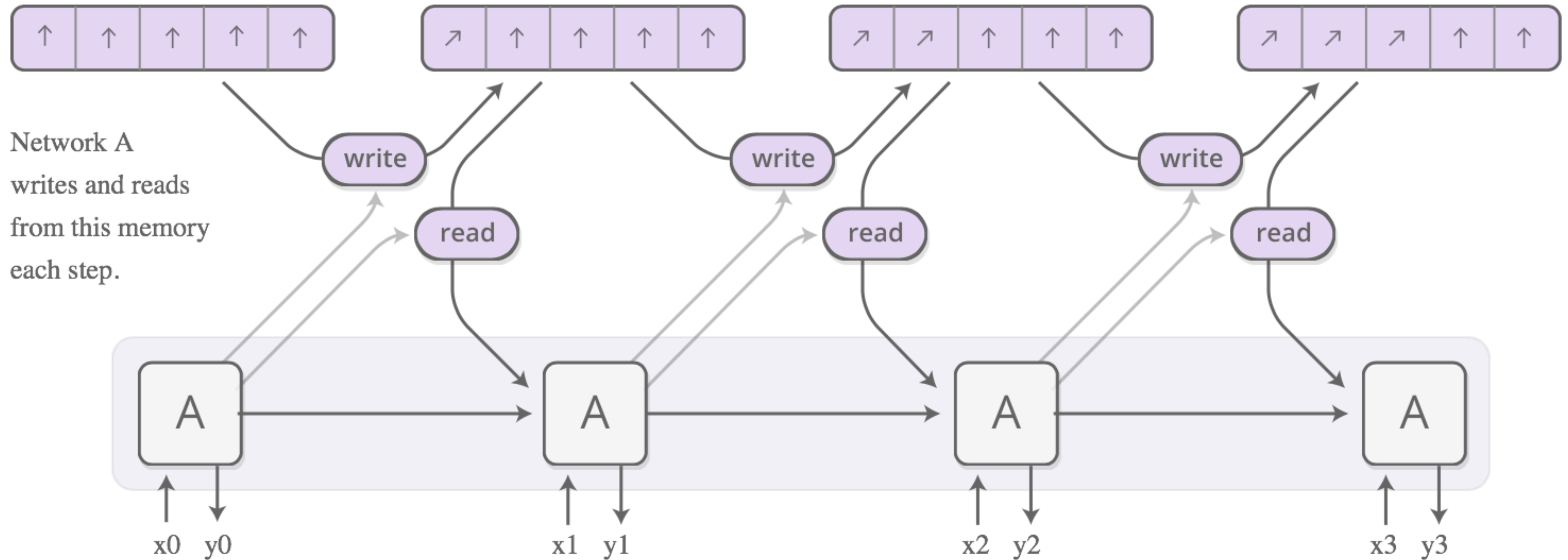
## Neural Programmers

can call functions, building programs as they run.



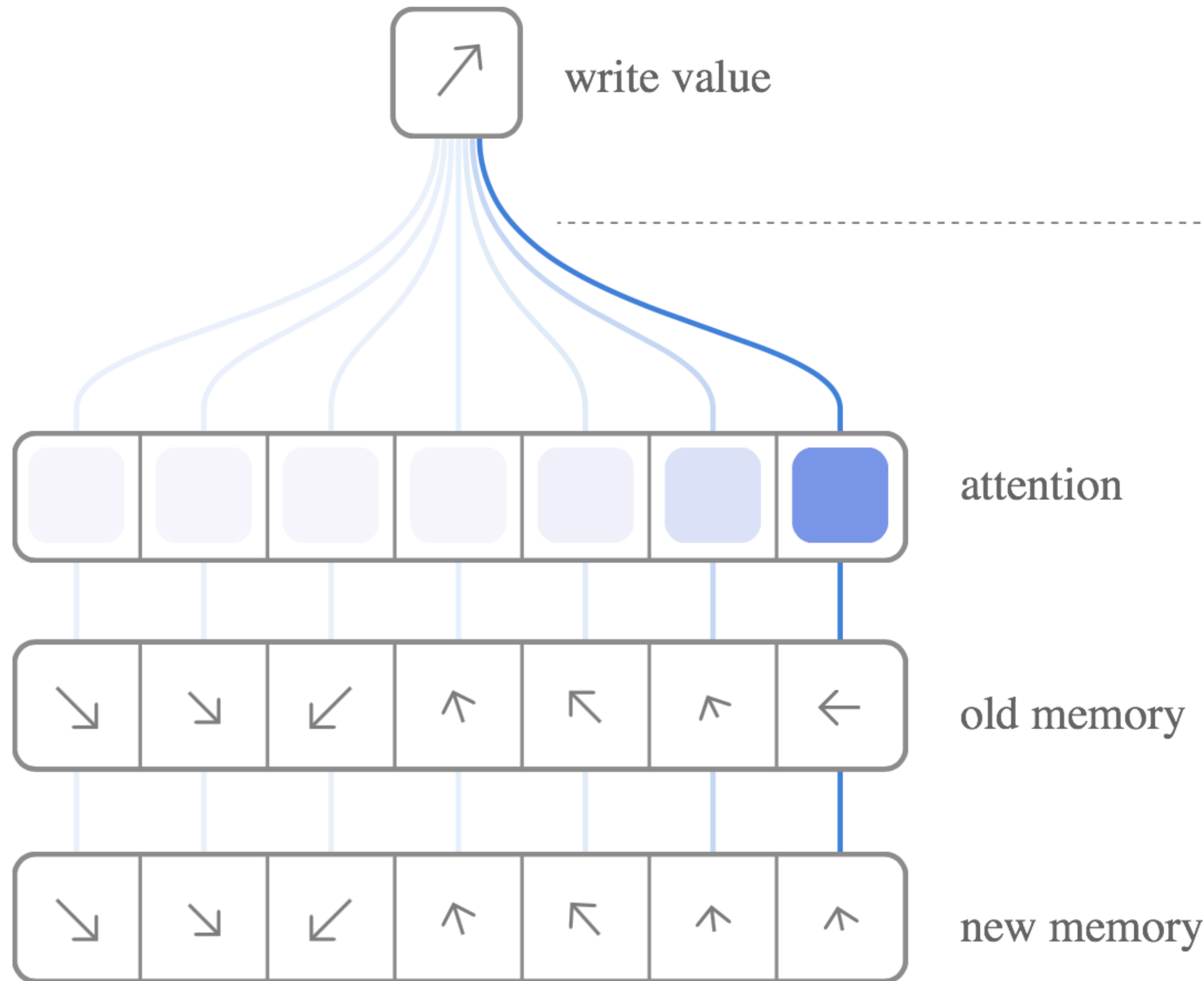
# Neural Turing Machines

Memory is an array of vectors.



Network A  
writes and reads  
from this memory  
each step.

# Neural Turing Machines



Instead of writing to one location, we write everywhere, just to different extents.

The RNN gives an **attention distribution**, describing how much we should change each memory position towards the write value.

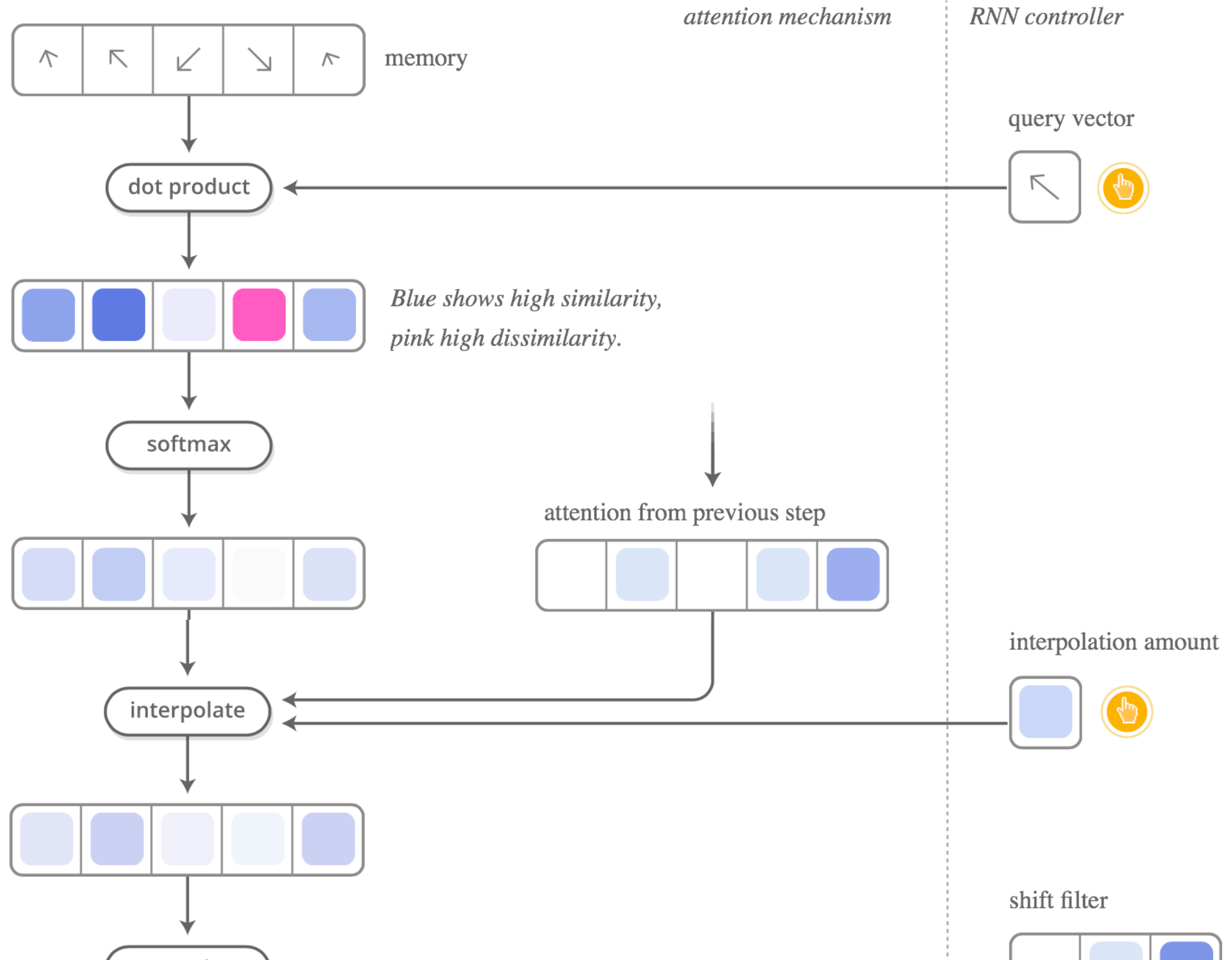
$$M_i \leftarrow a_i w + (1 - a_i) M_i$$

First, the controller gives a query vector and each memory entry is scored for similarity with the query.

The scores are then converted into a distribution using softmax.

Next, we interpolate the attention from the previous time step.

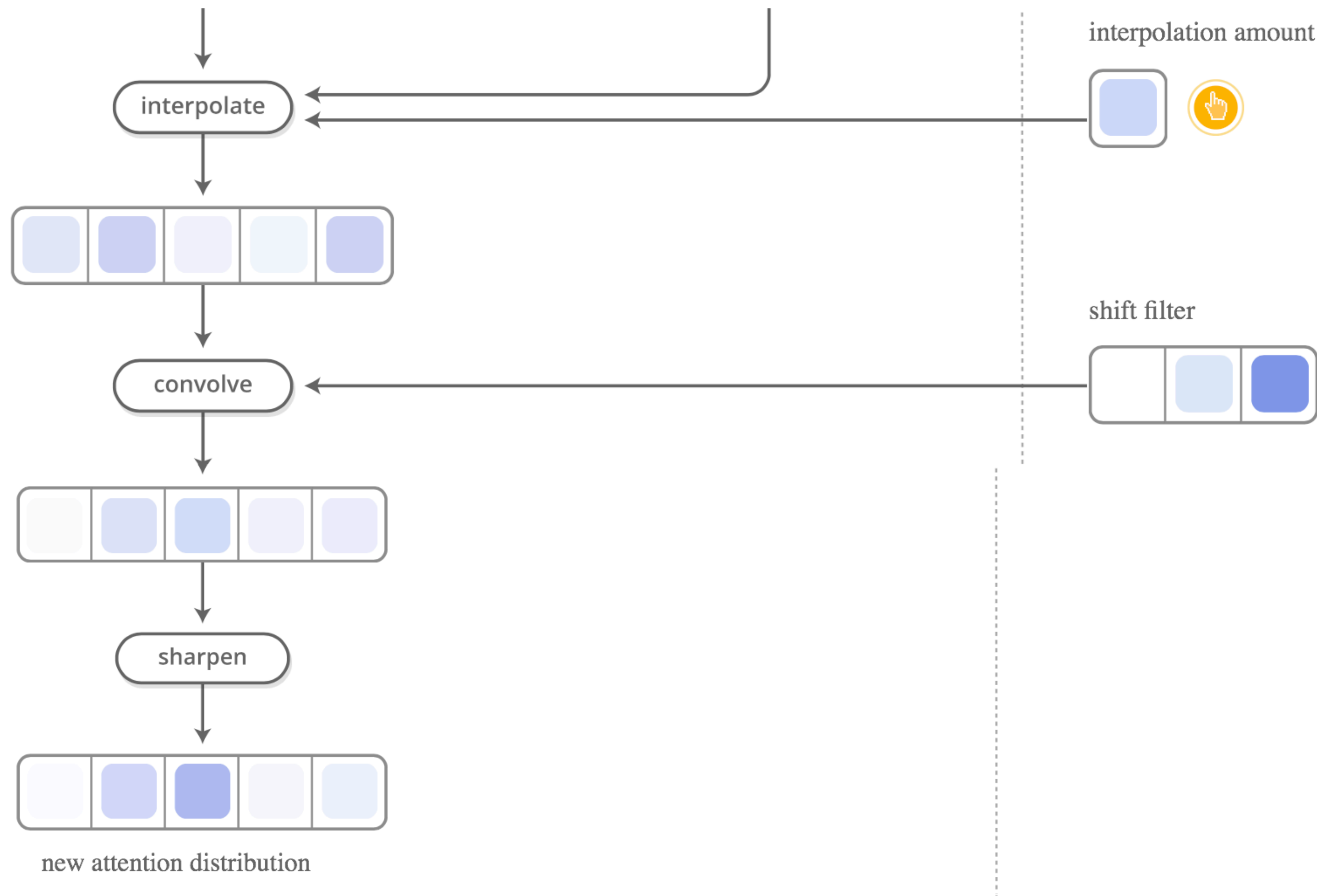
We convolve the attention with a shift filter—this allows the



Next, we interpolate the attention from the previous time step.

We convolve the attention with a shift filter—this allows the controller to move its focus.

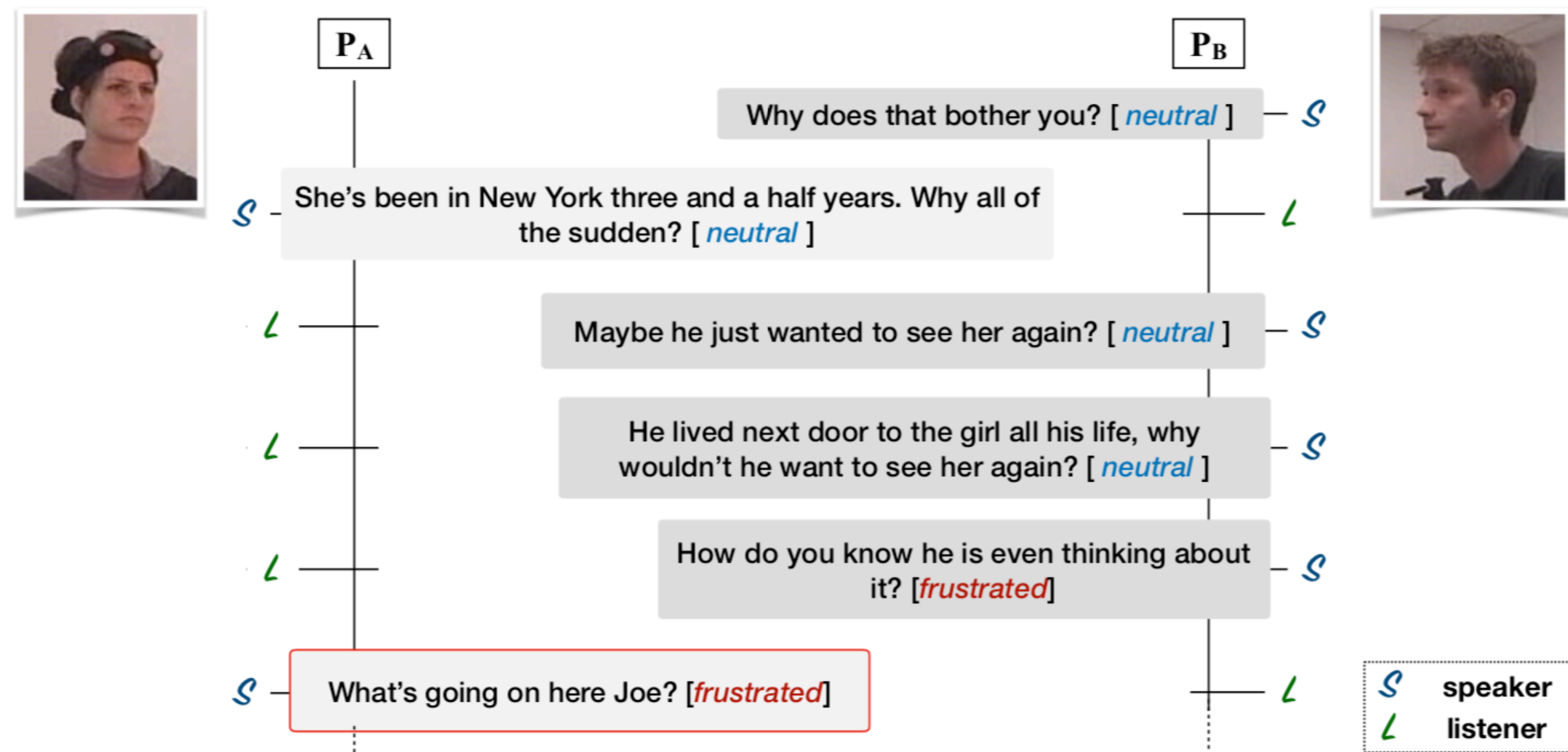
Finally, we sharpen the attention distribution. This final attention distribution is fed to the read or write operation.



# DialogueRNN: An Attentive RNN for Emotion Detection in Conversations

AAAI 2019

Navonil Majumder, Soujanya Poria, Devamanyu Hazarika, Rada Mihalcea, Alexander Gelbukh, Erik Cambria

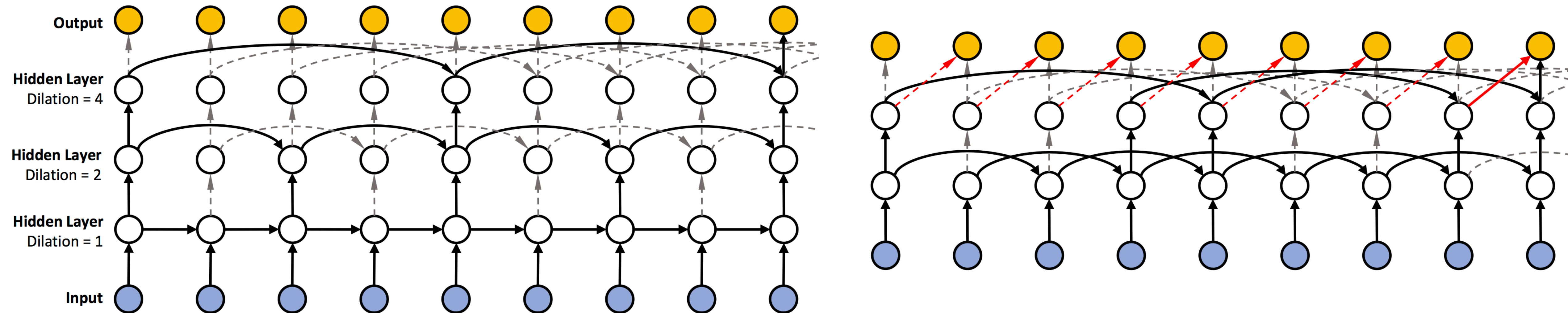




# Dilated Recurrent Neural Networks

NeurIPS 2017

Shiyu Chang, Yang Zhang, Wei Han, Mo Yu, Xiaoxiao Guo, Wei Tan, Xiaodong Cui,  
Michael Witbrock, Mark Hasegawa-Johnson, Thomas S. Huang



<http://papers.nips.cc/paper/6613-dilated-recurrent-neural-networks.pdf>

Next week!