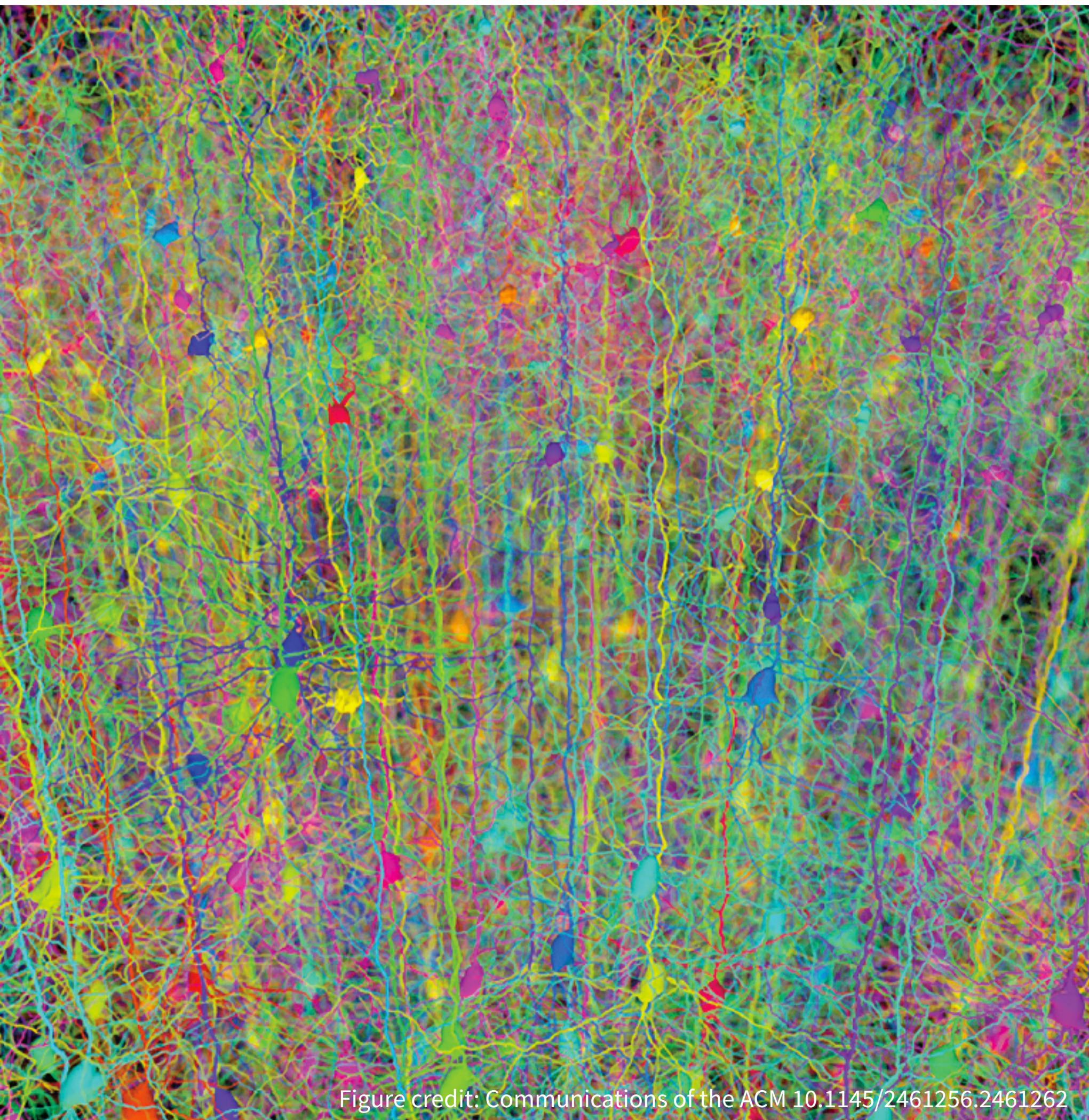


# Recurrent Neural Networks

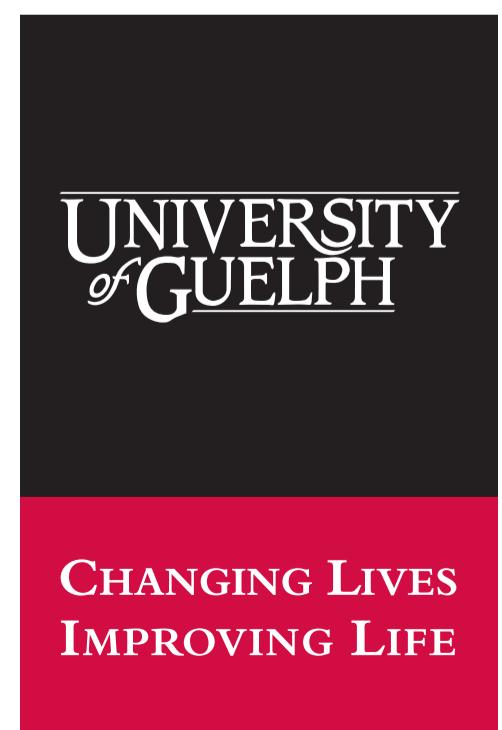


GRAHAM TAYLOR

VECTOR INSTITUTE

SCHOOL OF ENGINEERING  
UNIVERSITY OF GUELPH

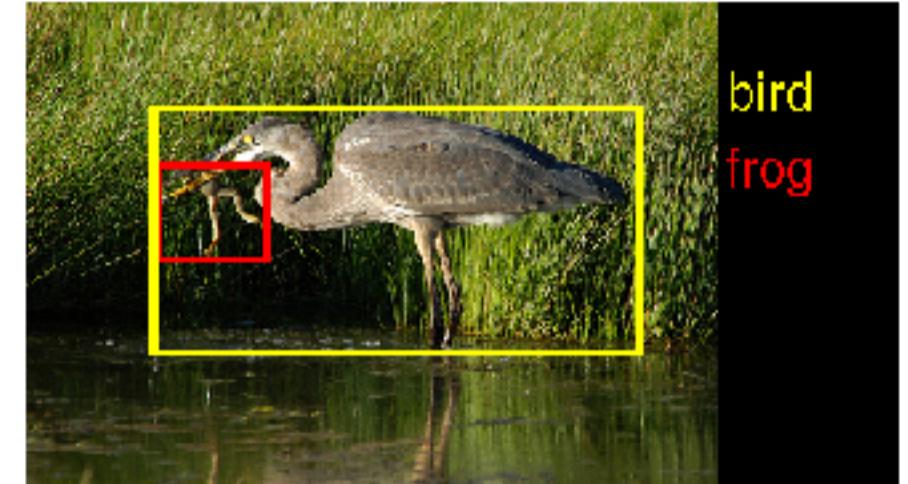
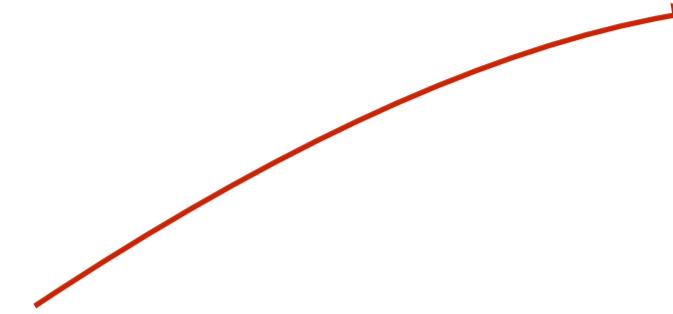
CANADIAN INSTITUTE  
FOR ADVANCED RESEARCH



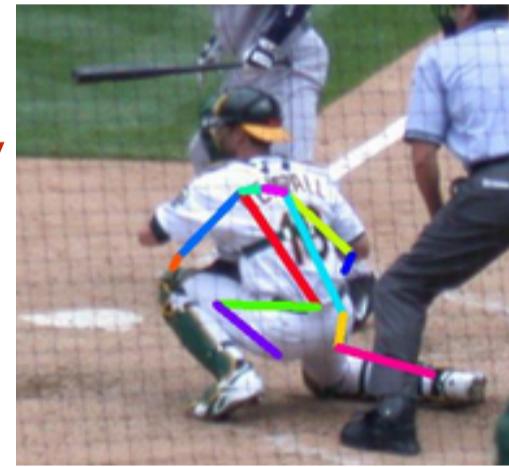
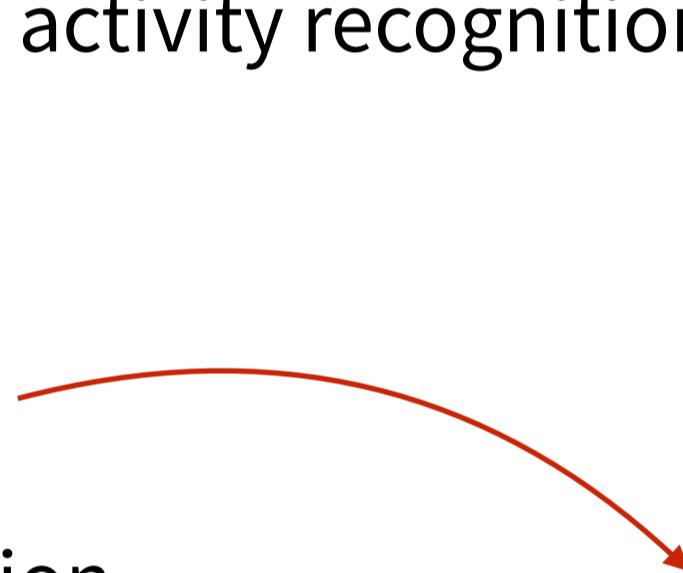
**CIFAR**  
CANADIAN  
INSTITUTE  
FOR  
ADVANCED  
RESEARCH

# Deep Learning Successes

- Vision
  - Object recognition and detection



- Speech
  - Speech recognition
  - Speaker authentication

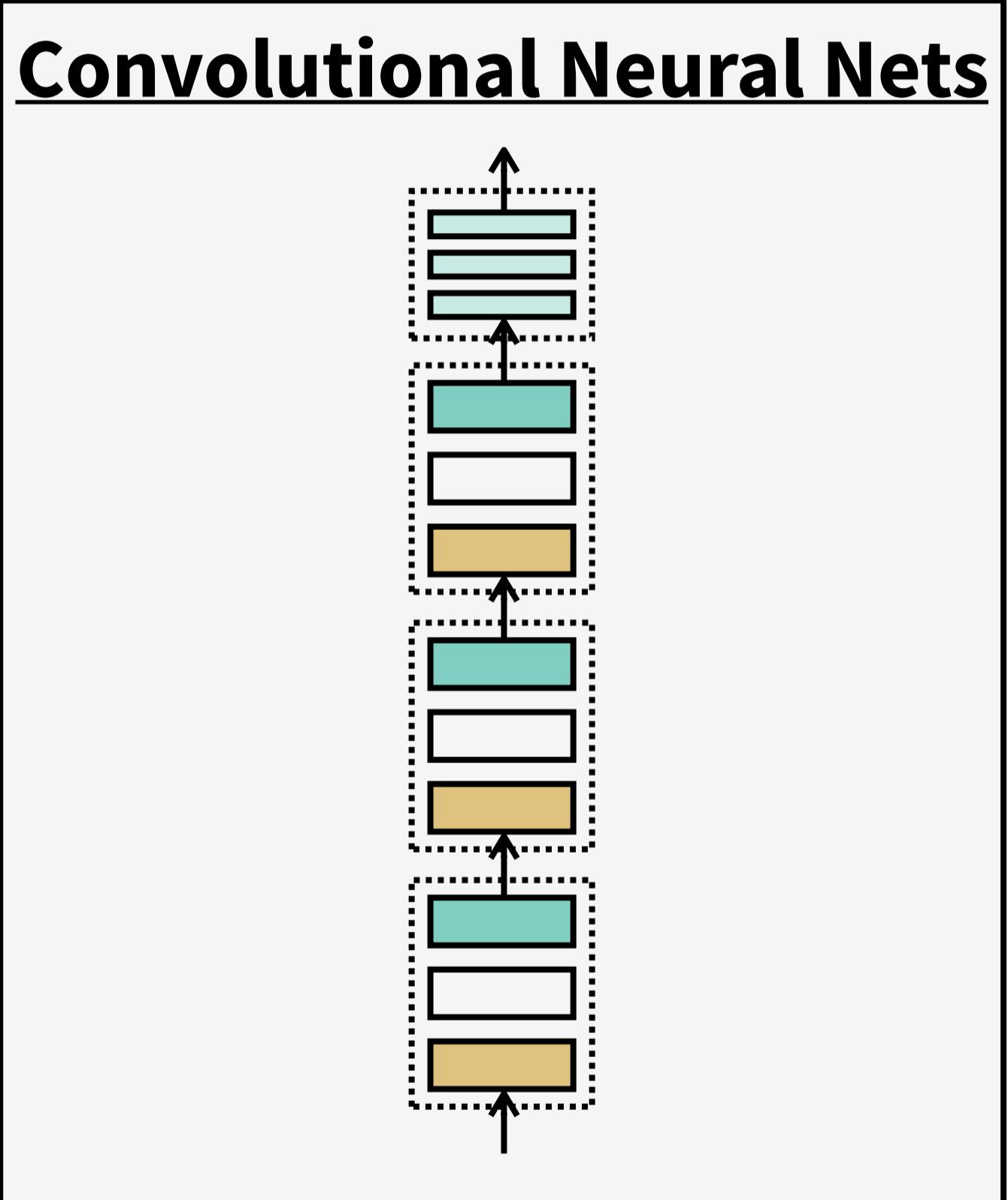


- Natural Language Processing
  - Question answering
  - Machine translation



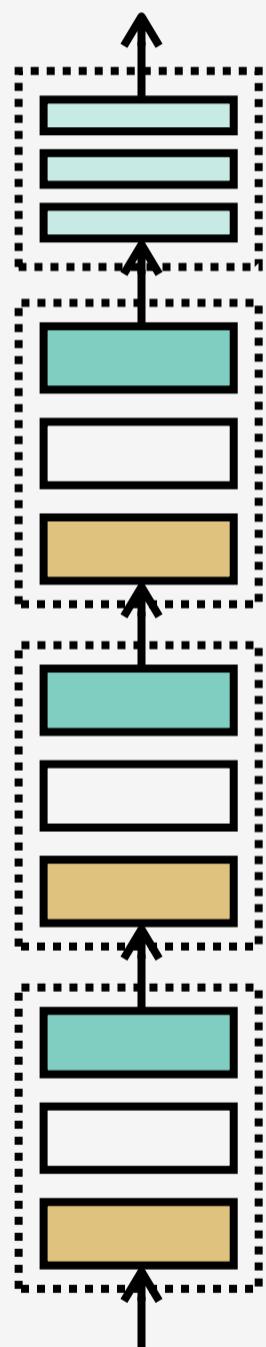
# Two “Hammers”

# Two “Hammers”

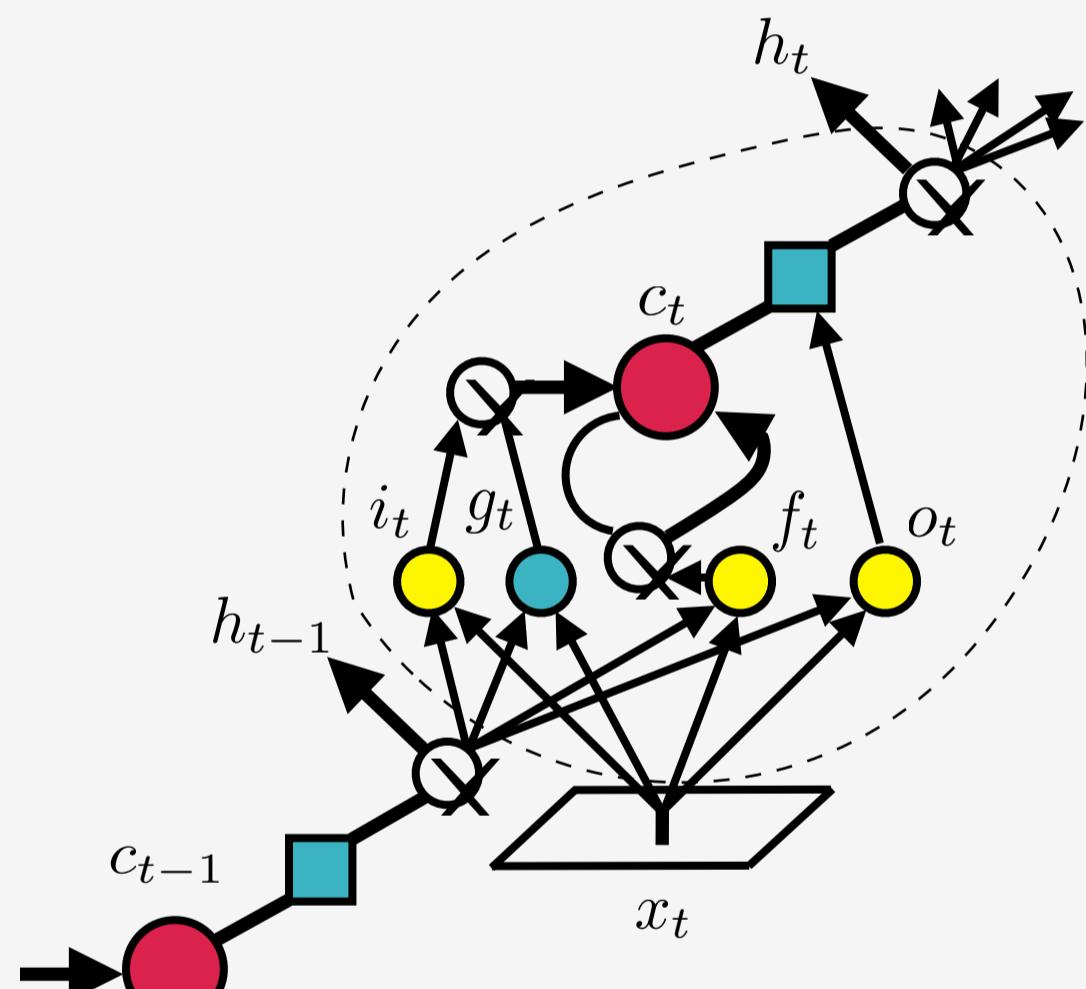


# Two “Hammers”

## Convolutional Neural Nets

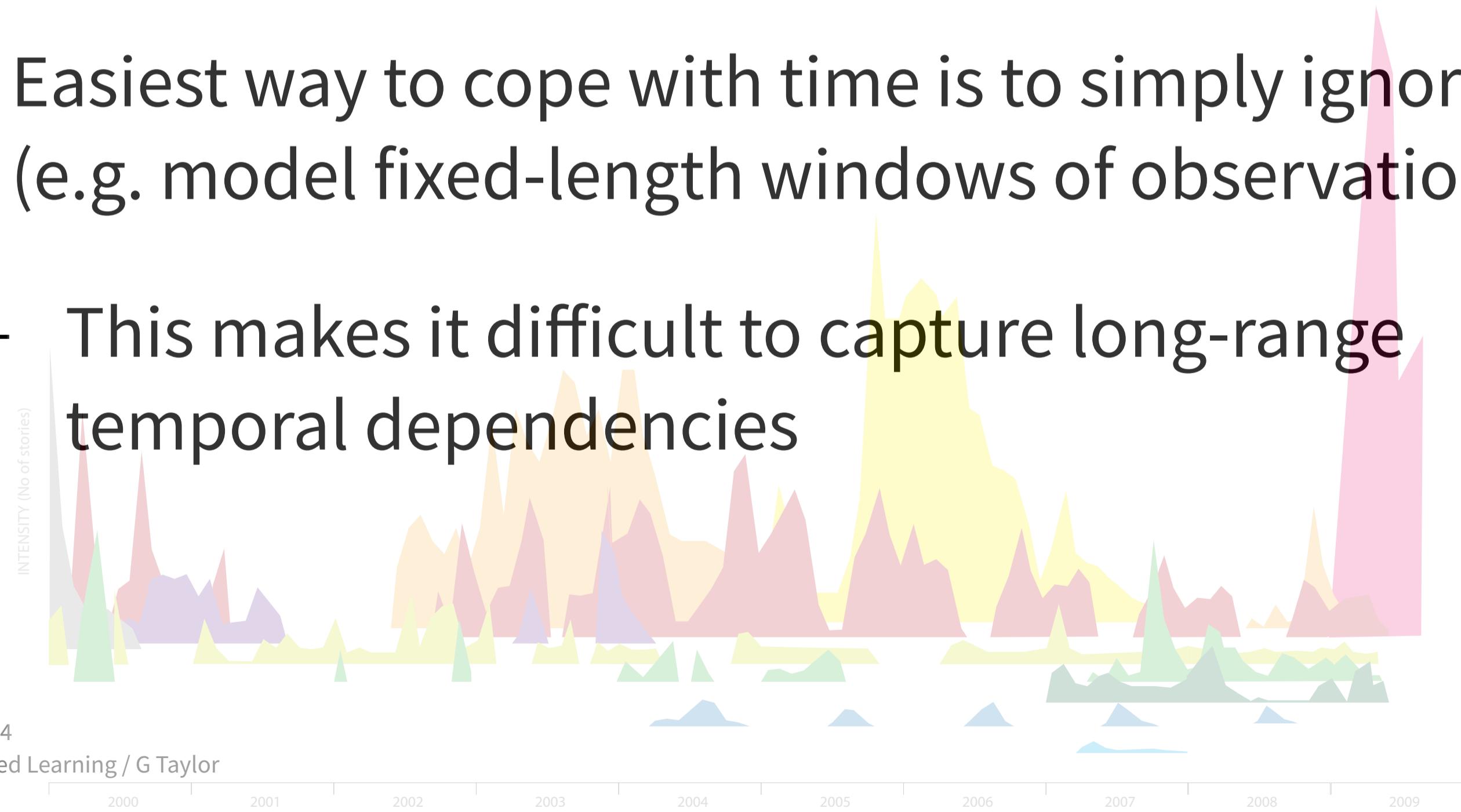


## Recurrent Neural Nets



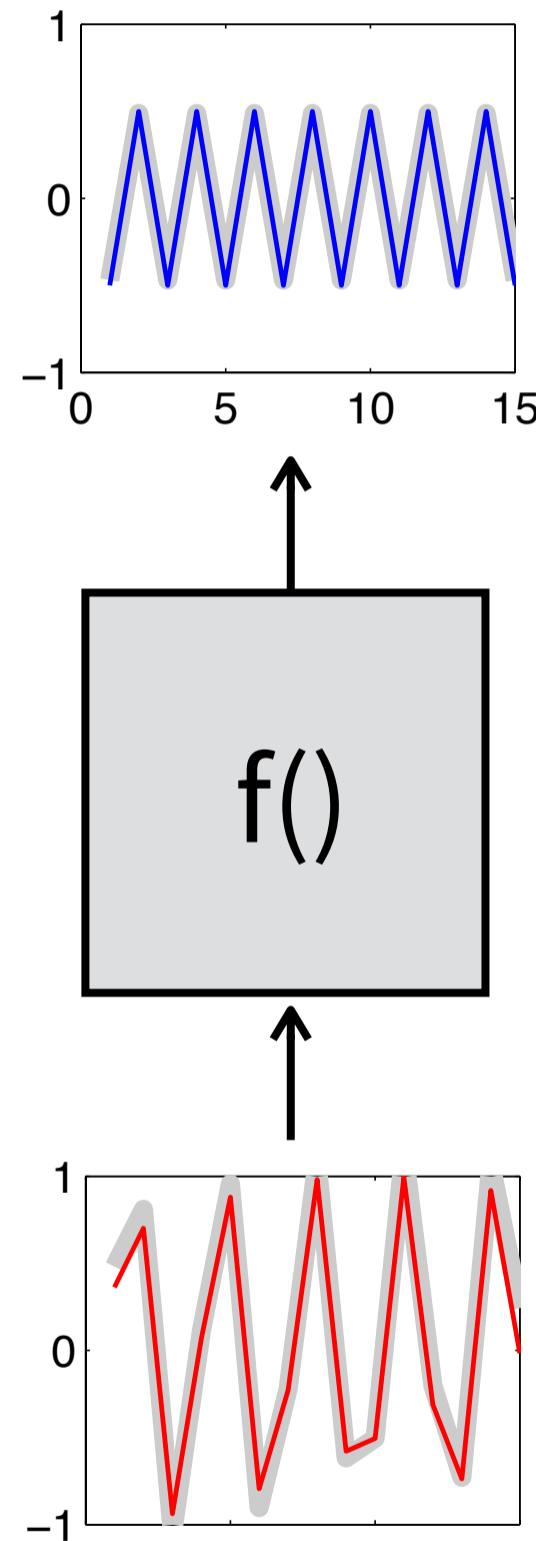
# Learning Sequence Representations

- Time is an integral part of many human behaviours (e.g. motion, reasoning, dialogue)
- Easiest way to cope with time is to simply ignore it (e.g. model fixed-length windows of observations)
  - This makes it difficult to capture long-range temporal dependencies

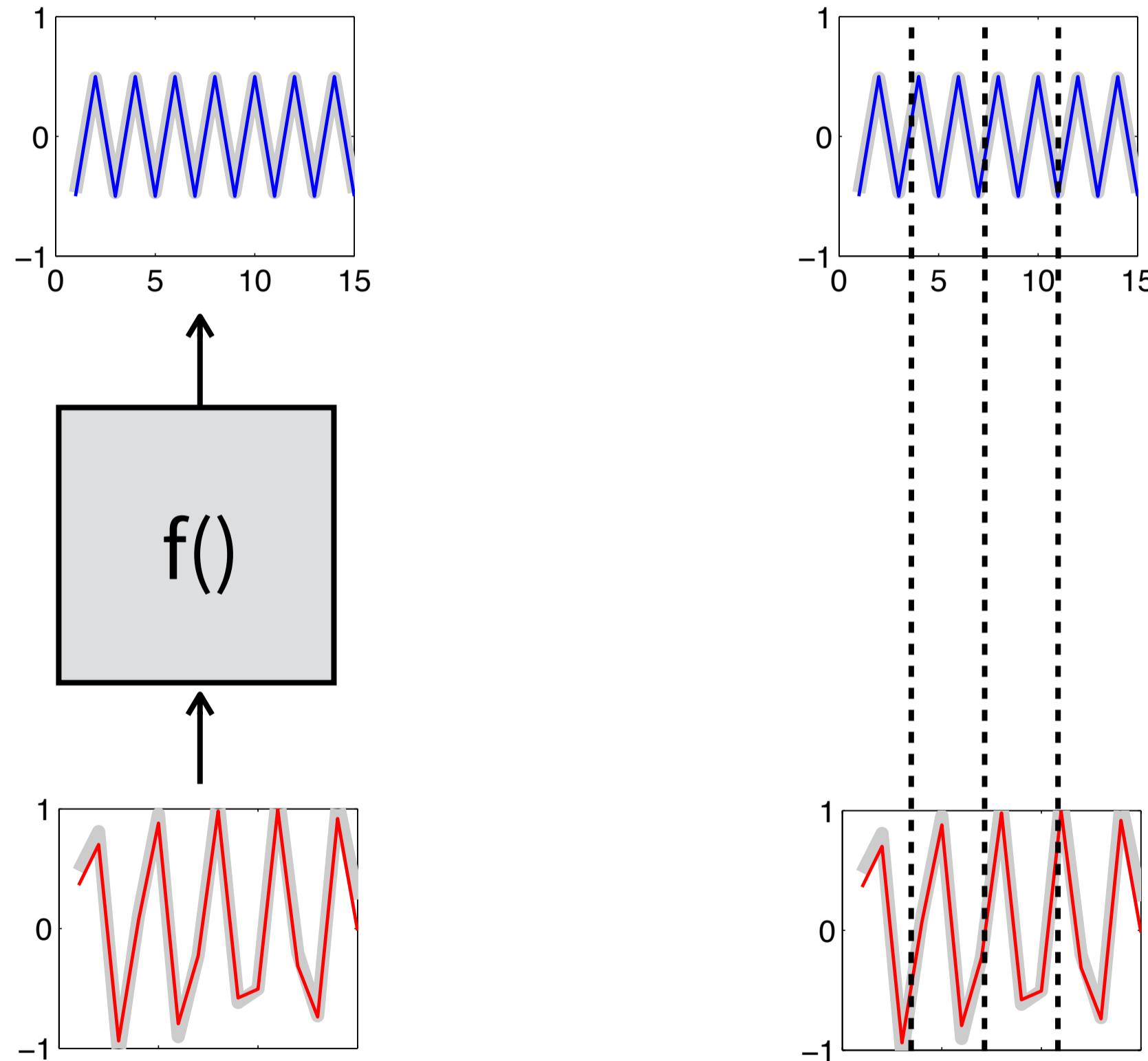


# Modeling fixed-length windows?

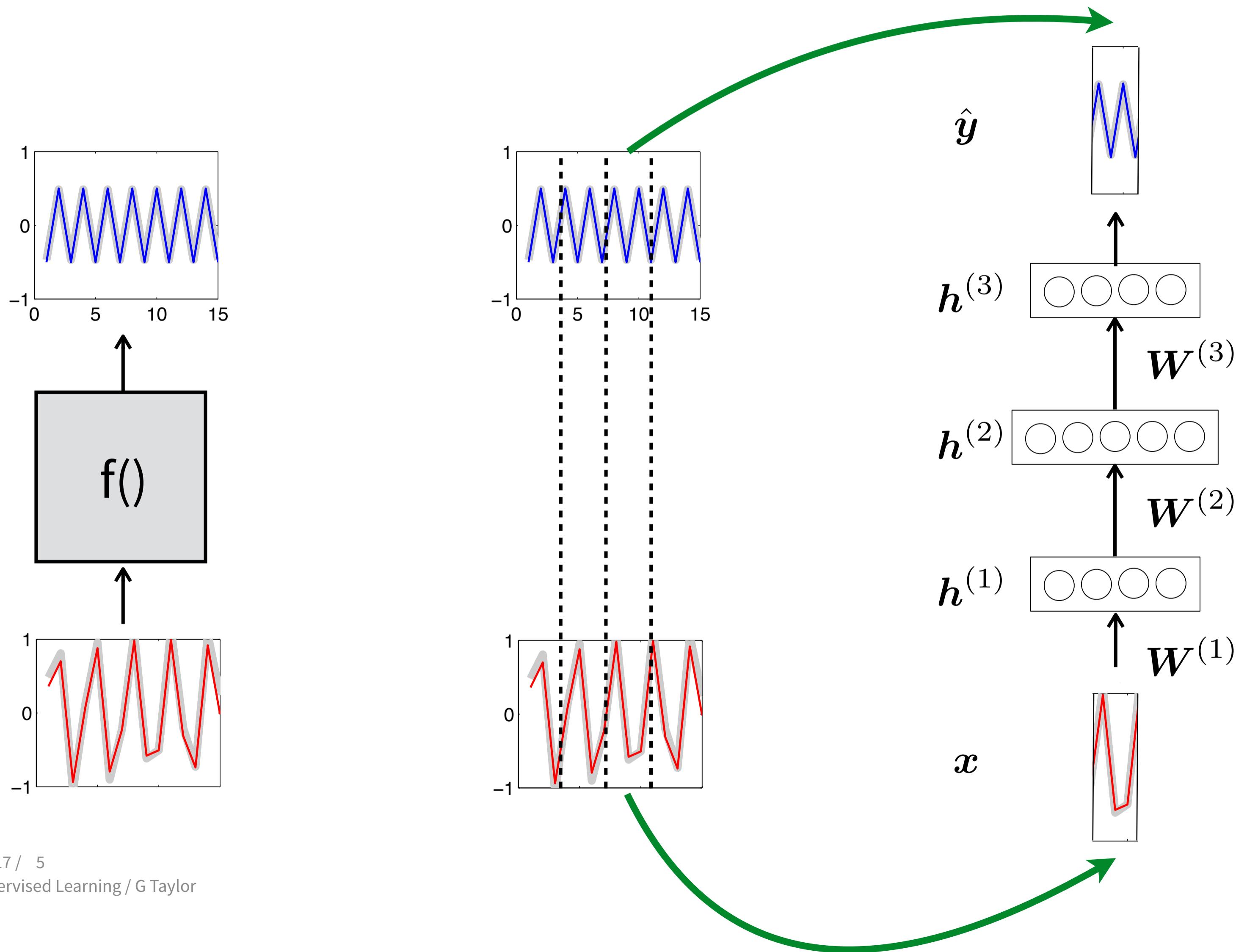
# Modeling fixed-length windows?



# Modeling fixed-length windows?

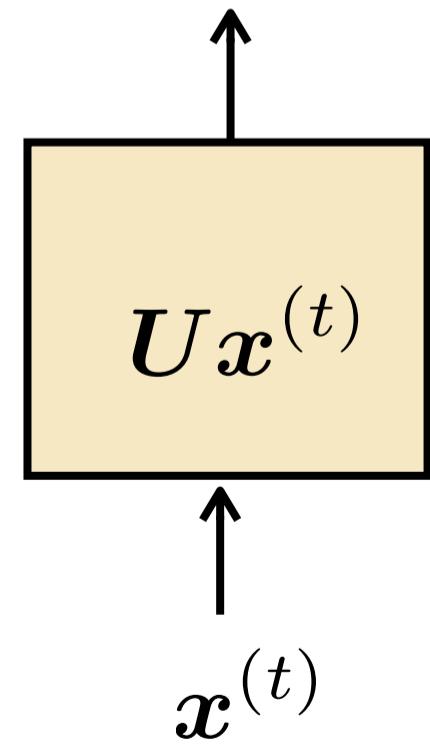


# Modeling fixed-length windows?

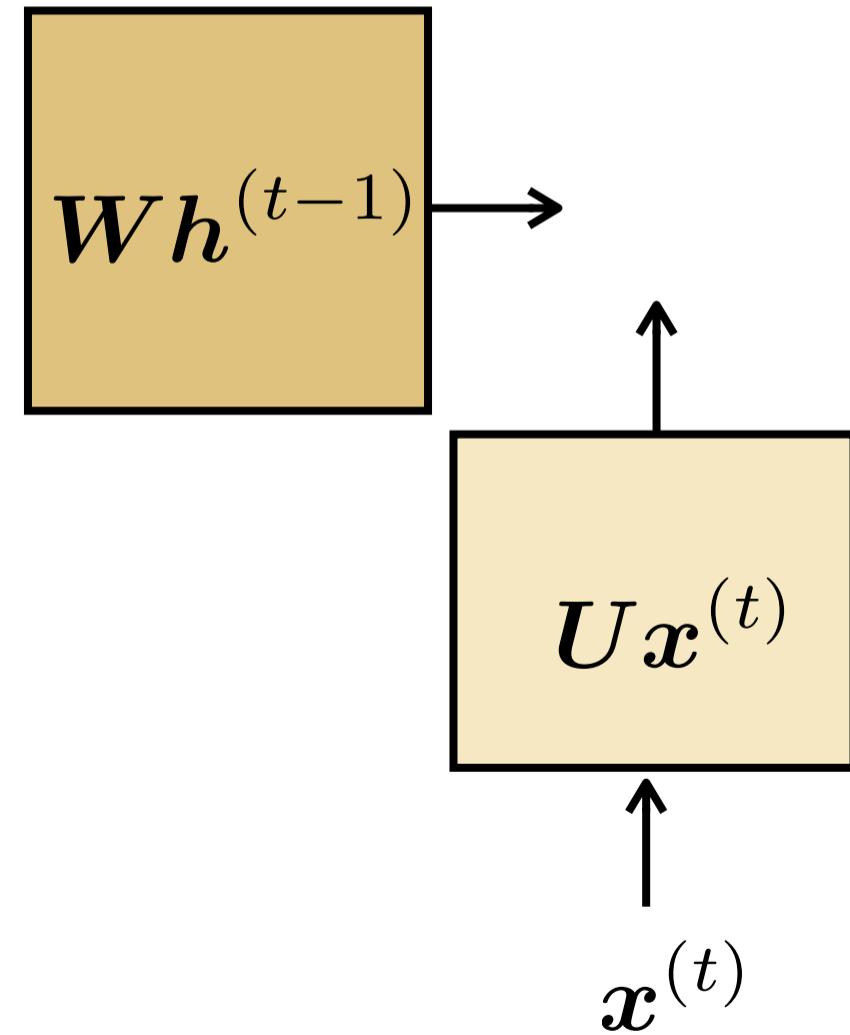


# Recurrent Neural Networks

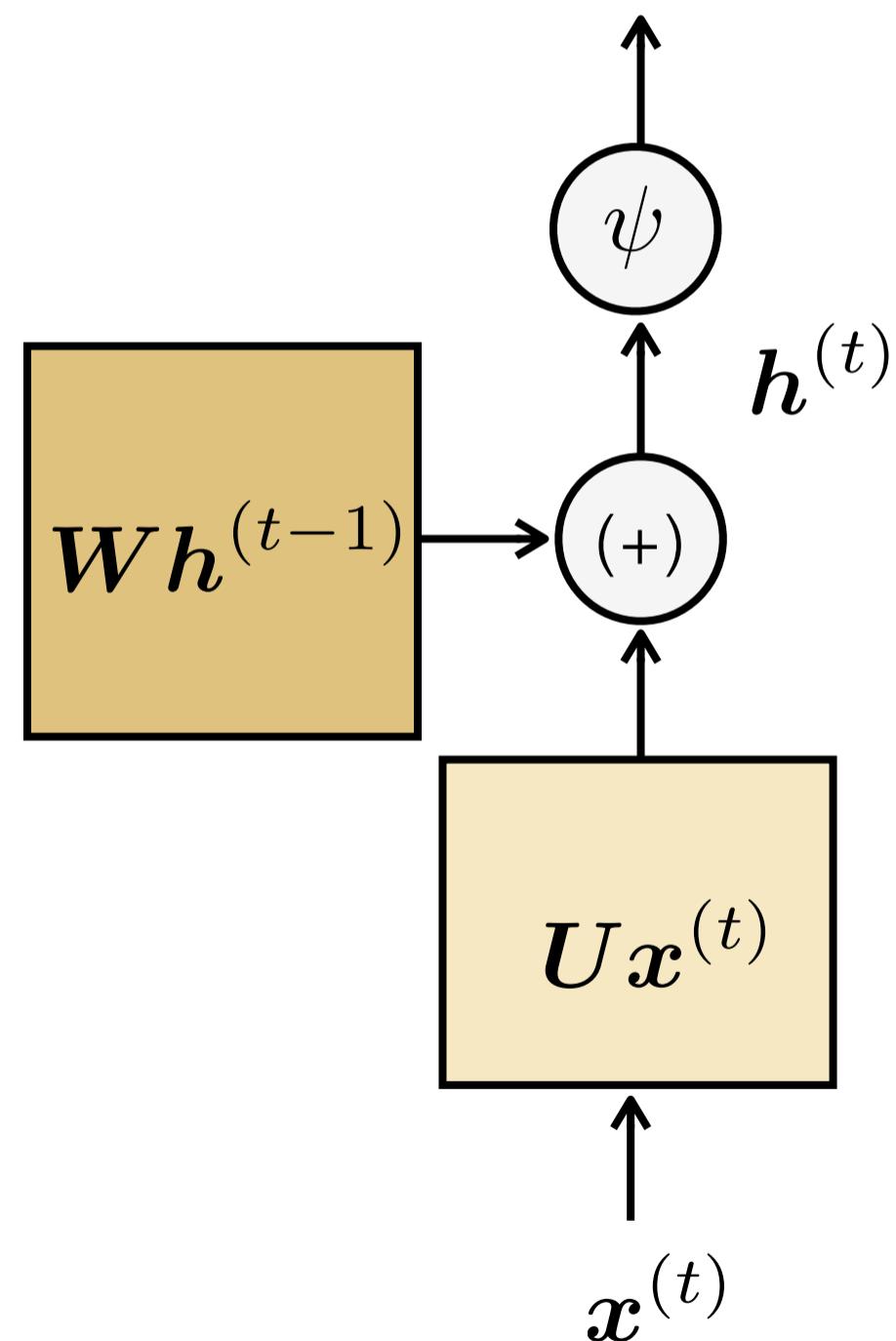
# Recurrent Neural Networks



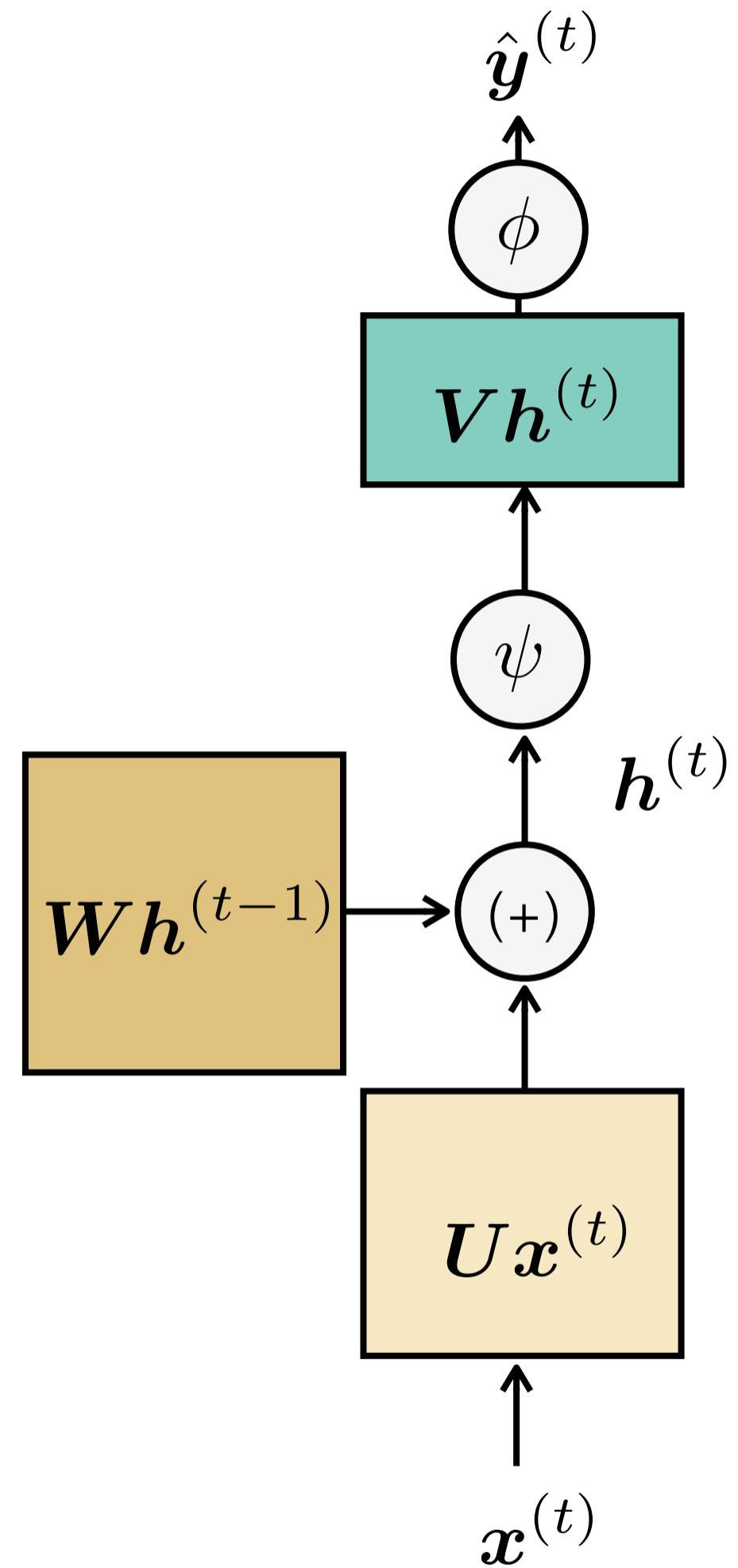
# Recurrent Neural Networks



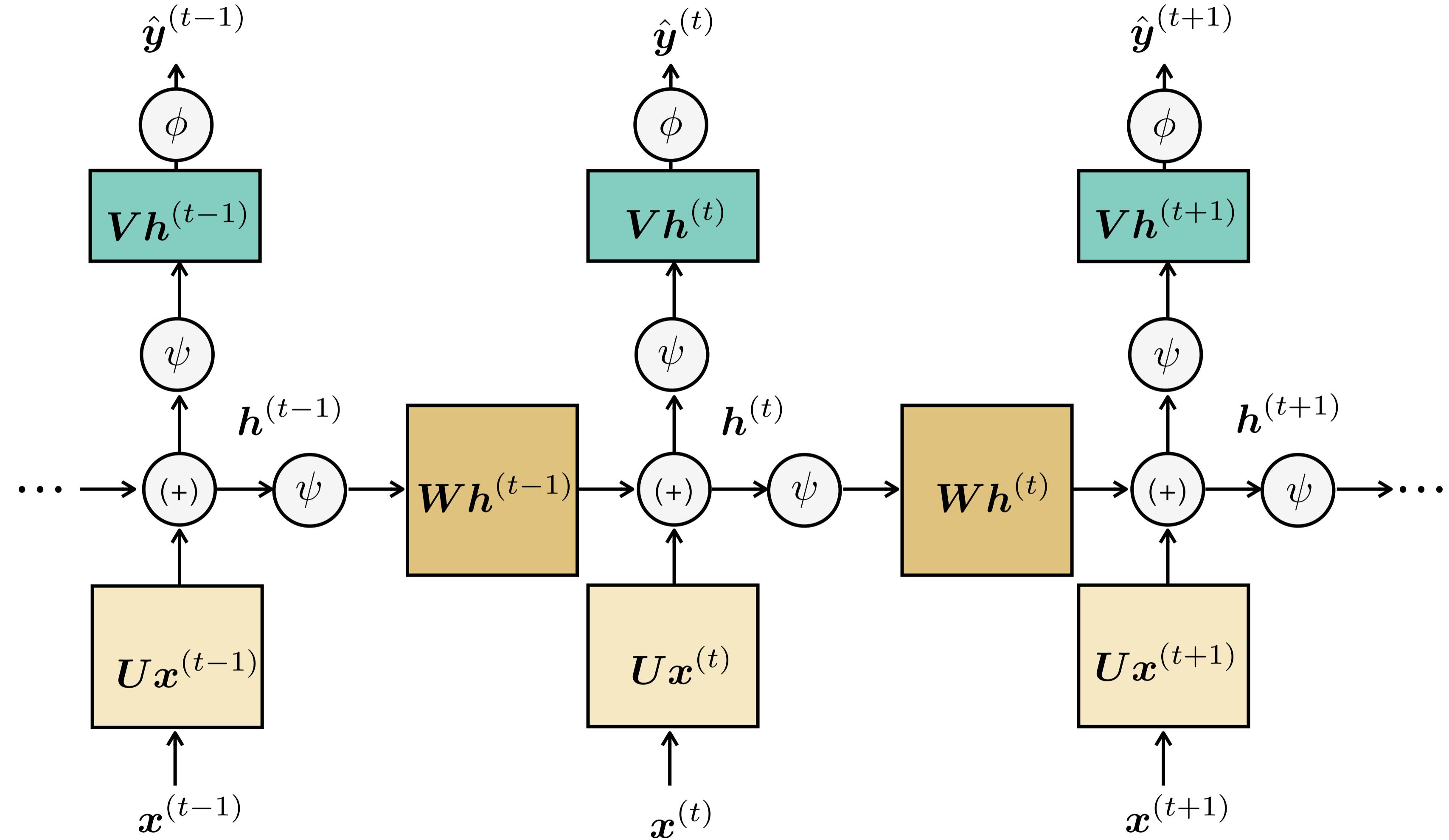
# Recurrent Neural Networks



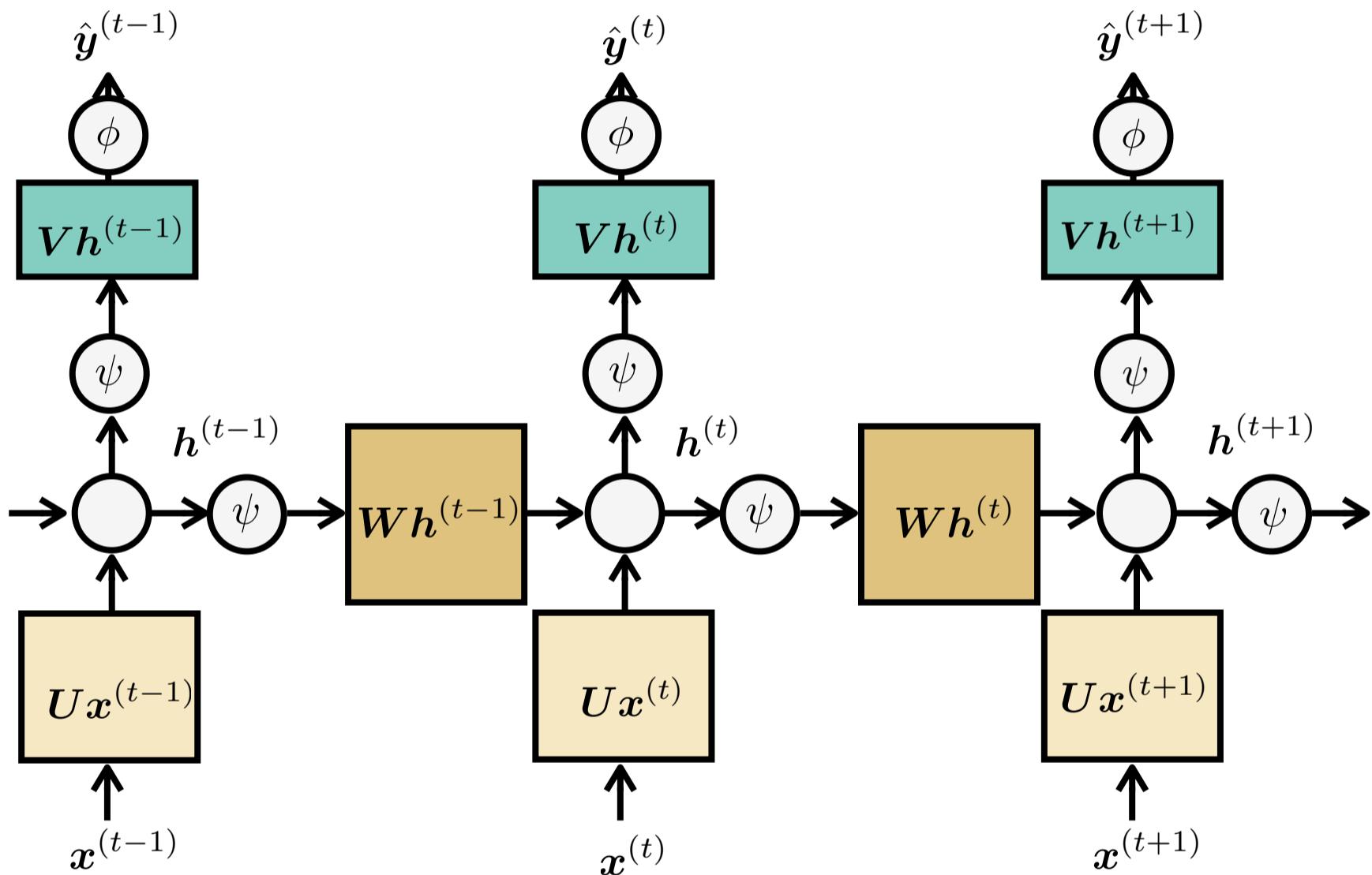
# Recurrent Neural Networks



# Recurrent Neural Networks



# RNN (Update)



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

# RNNs: The Promise

- Possibly high-dimensional, distributed, internal representation and nonlinear dynamics allow model, in theory, model complex time series
  - and capture **long-range dependencies!**
- Exact gradients can be computed exactly via **Backpropagation Through Time**

# Generality

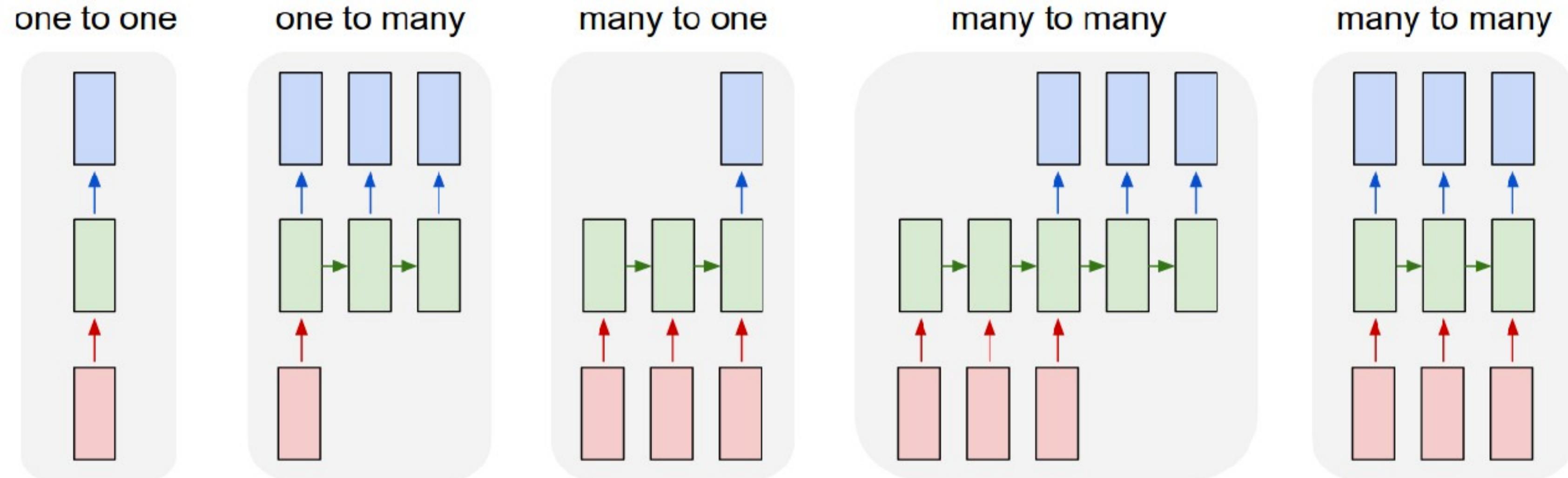
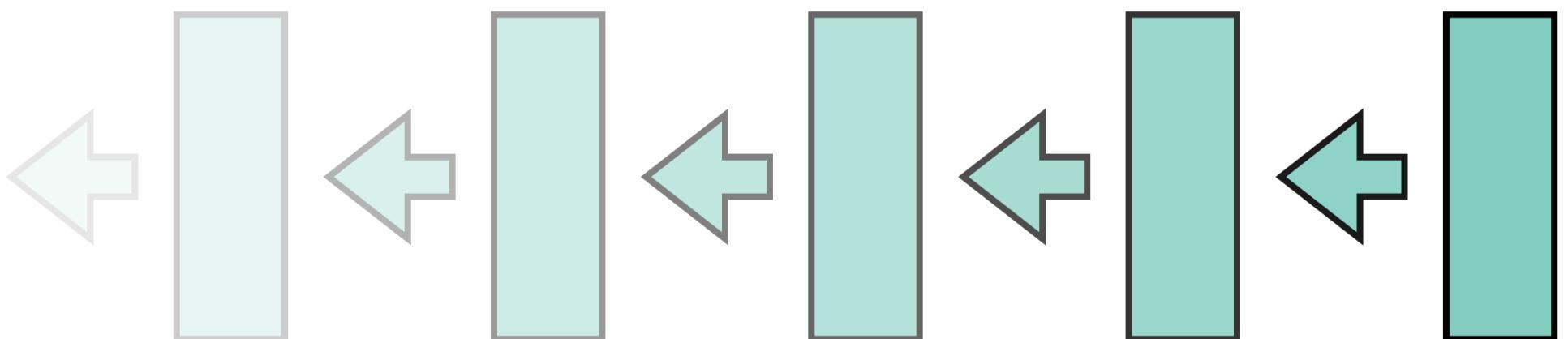


Figure: Andrej Karpathy  
<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

# RNNs: The Reality

- A powerful and flexible architecture - but?
  - Training RNNs via gradient descent **fails on simple problems** (typically memorization tasks)
  - Attributed to “vanishing” or “exploding” gradients

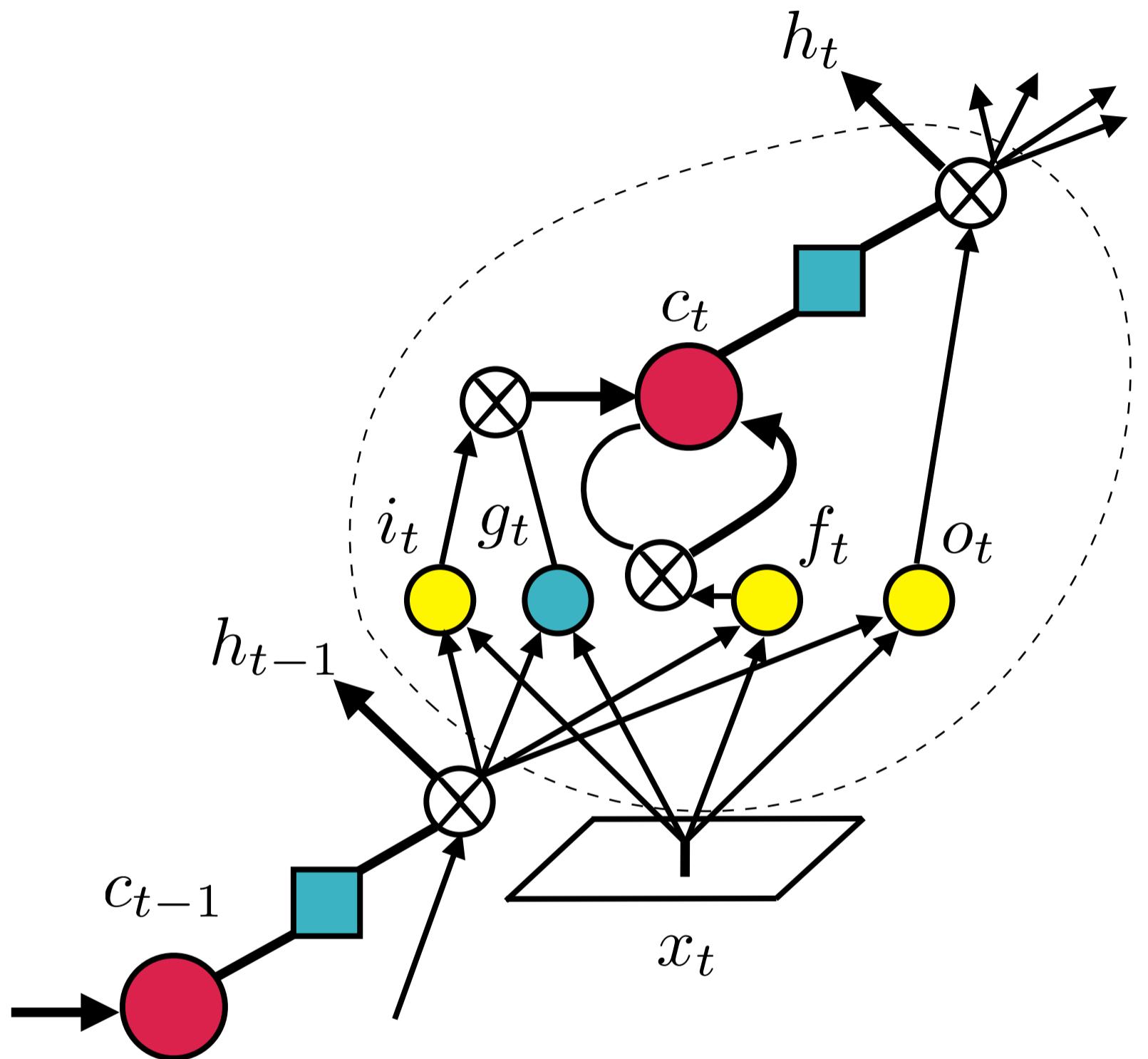


# Solutions: Architectural

- Gating 
- Long Short-term Memory (Hochreiter et al. 1997)
- Gated Recurrent Units (Chung et al. 2014)
- “Structured” depth (Pascanu et al. 2014)
- Rectifiers (Le, Jaitly and Hinton 2015) 
- Units operating at different timescales (Hihi and Bengio 1996, Koutnik et al. 2014, Chung et al. 2017)

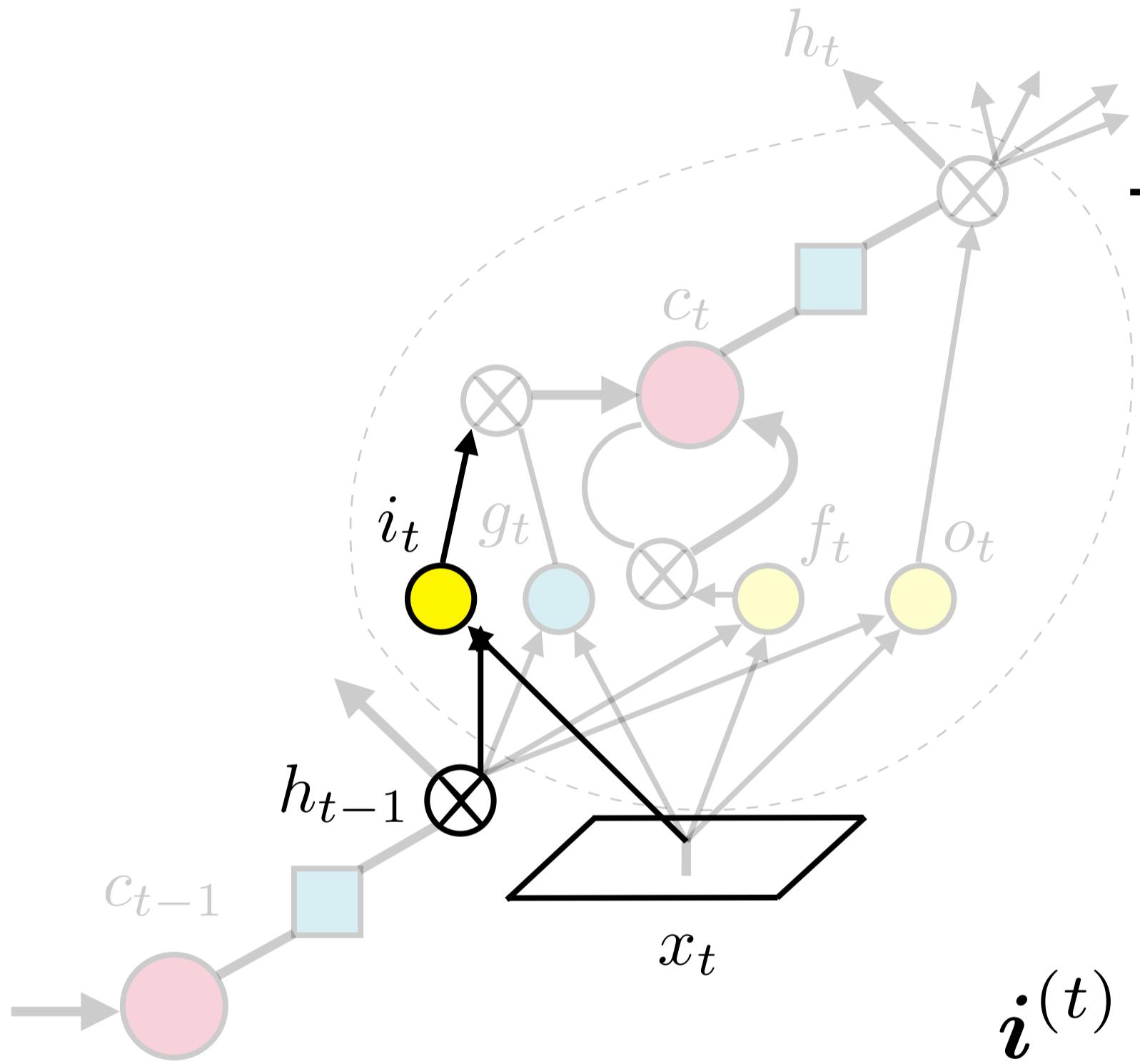
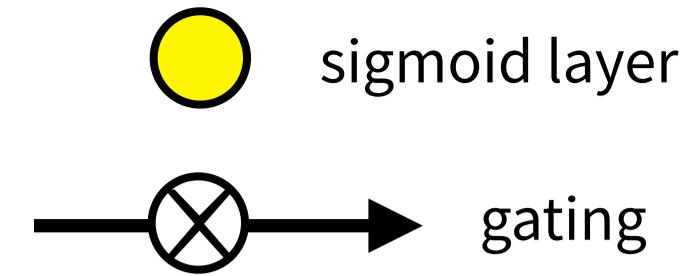
# LSTM

- (Yellow circle) sigmoid layer
- (Teal circle)  $tanh$  layer
- (Teal square) element-wise  $tanh$
- (Red circle) memory cell
- (Line with circle and cross) gating



RNN with memory cells  
designed to retain long-term information

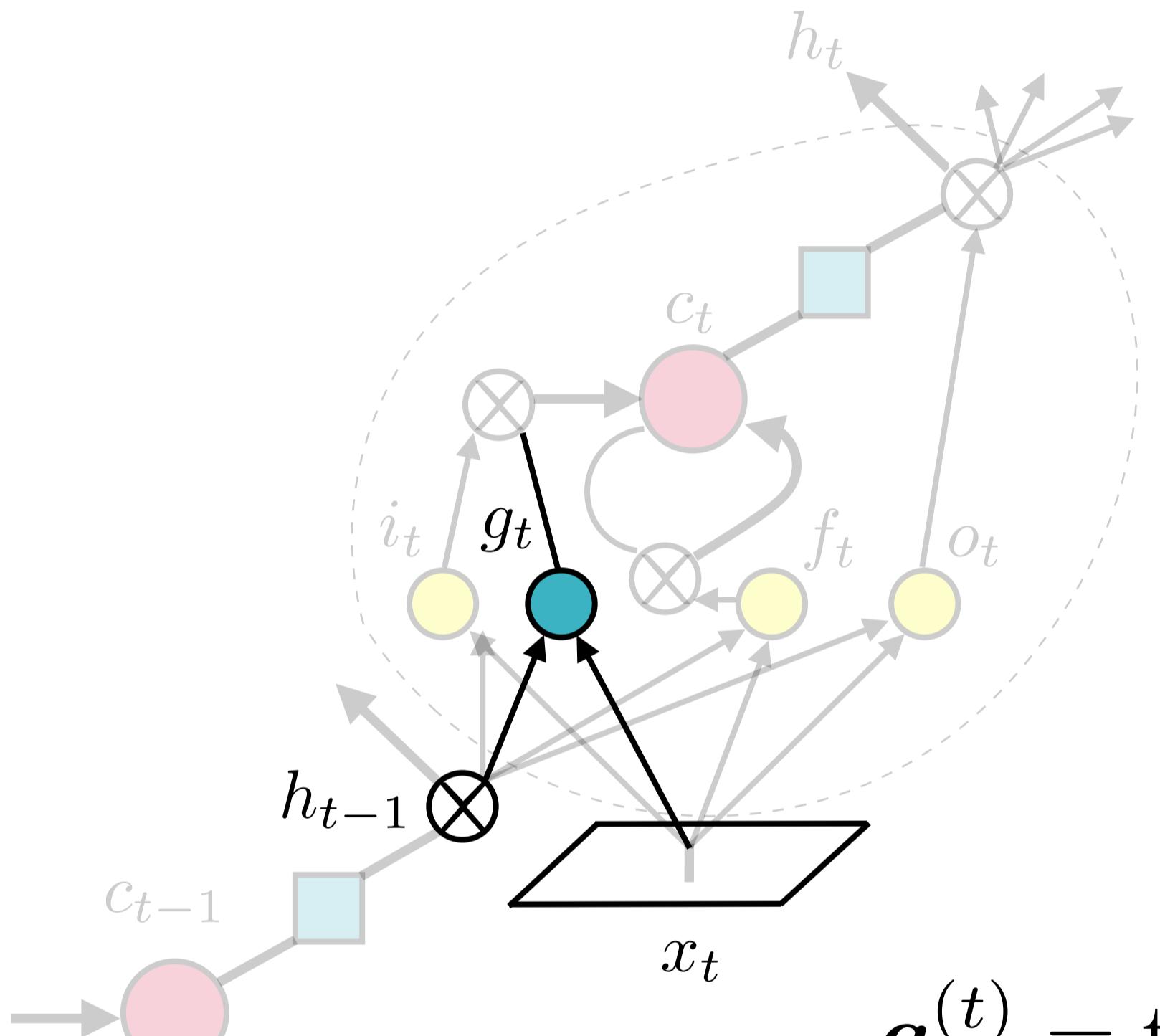
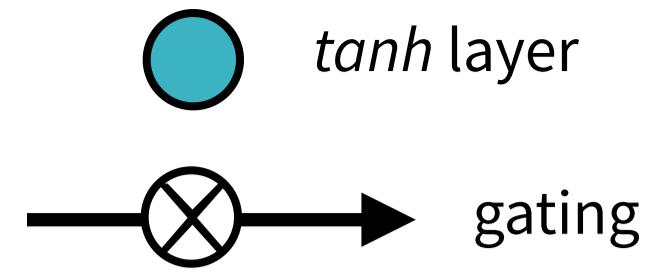
# LSTM - Input gate



The input gate regulates the flow of new information into the cell

$$i^{(t)} = \sigma(b^i + W^i [h^{(t-1)}, x^{(t)}])$$

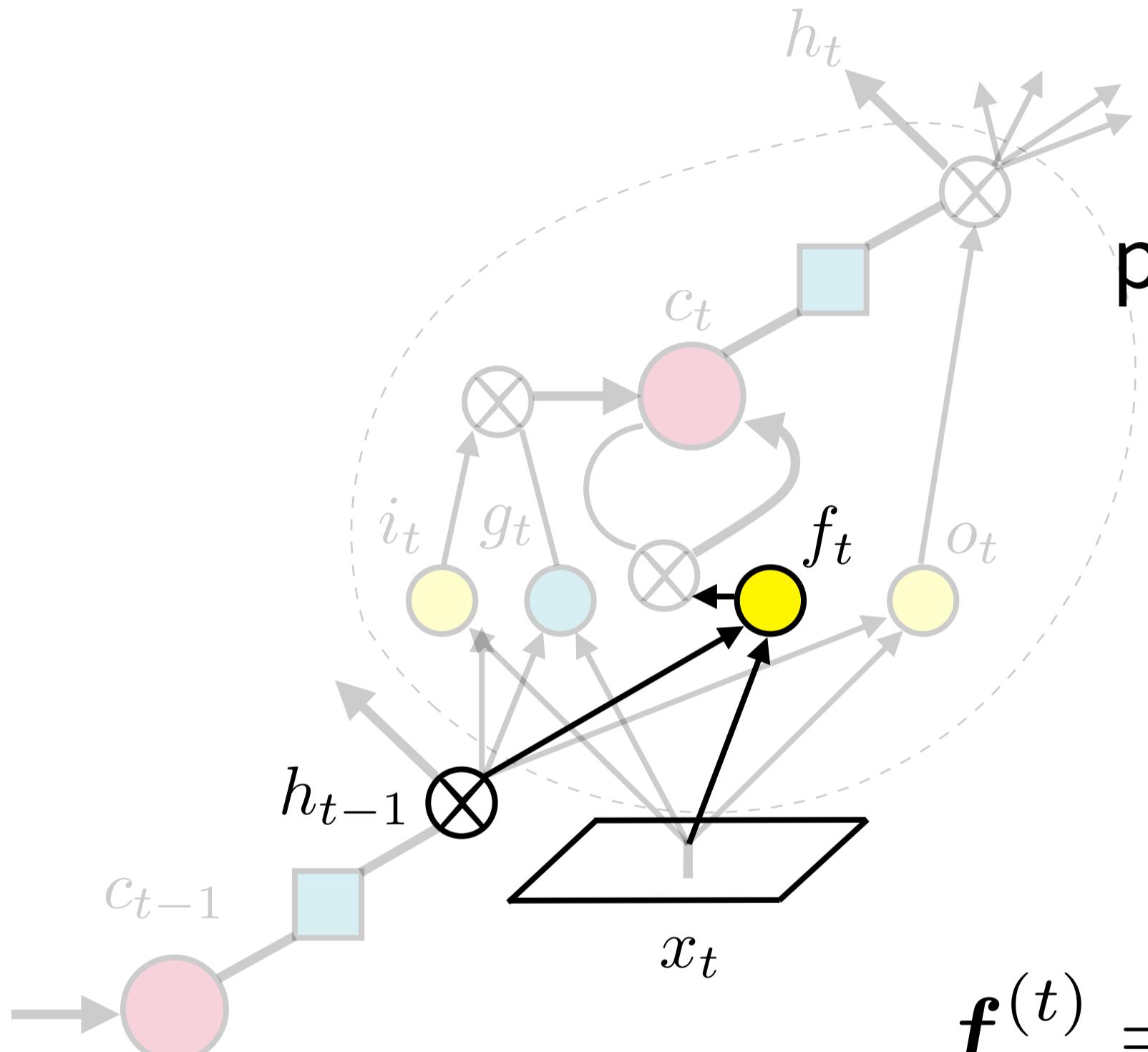
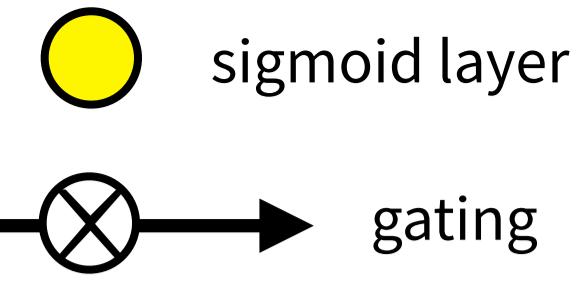
# LSTM - Content



The new information that flows into the cell is a function of the past hidden state and the input (similar to RNN)

$$g^{(t)} = \tanh \left( b^g + W^g [h^{(t-1)}, x^{(t)}] \right)$$

# LSTM - Forget gate

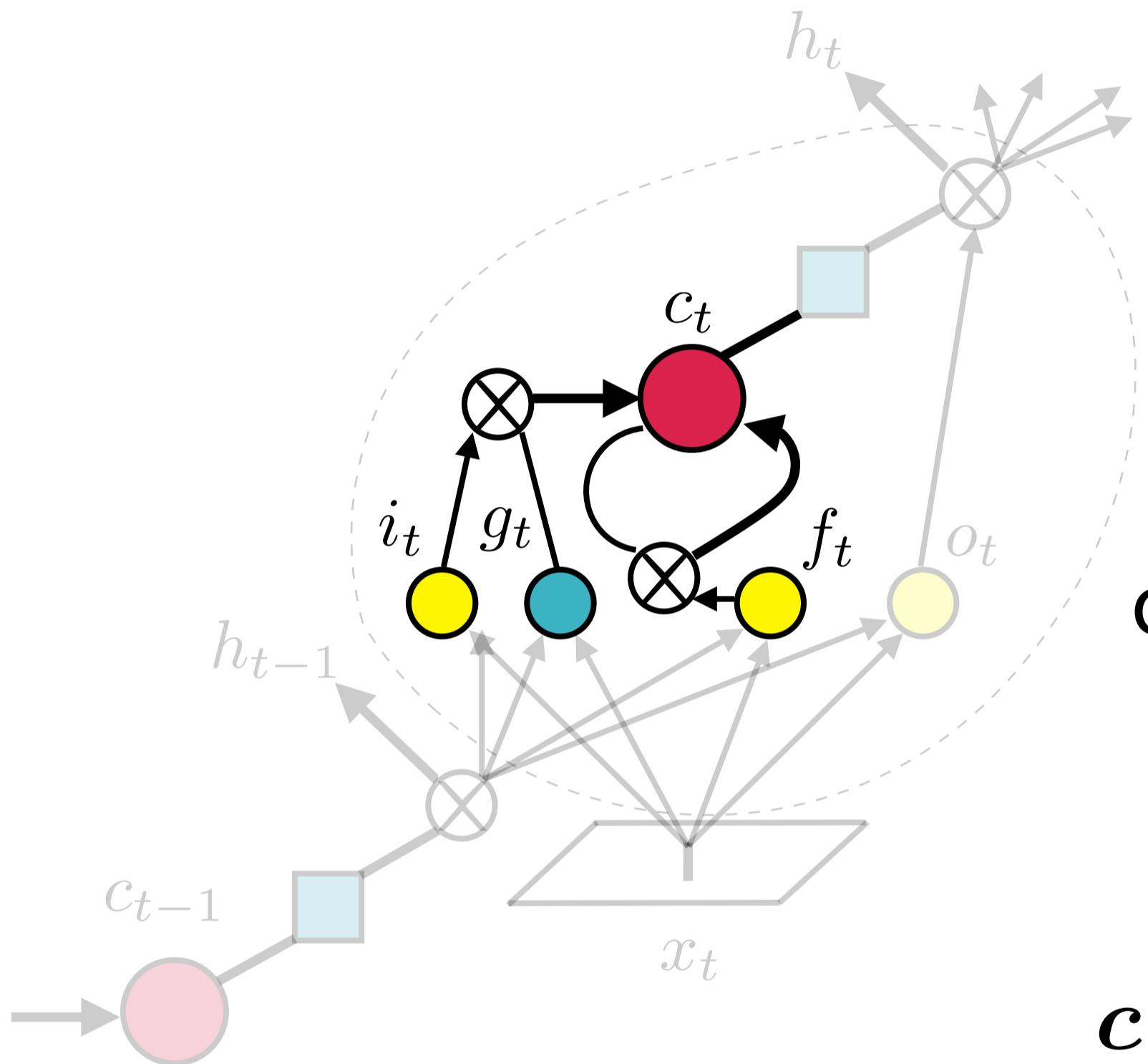


The forget gate regulates  
preservation of content in the cell

$$f^{(t)} = \sigma(b^f + W^f [h^{(t-1)}, x^{(t)}])$$

# LSTM - Cell update

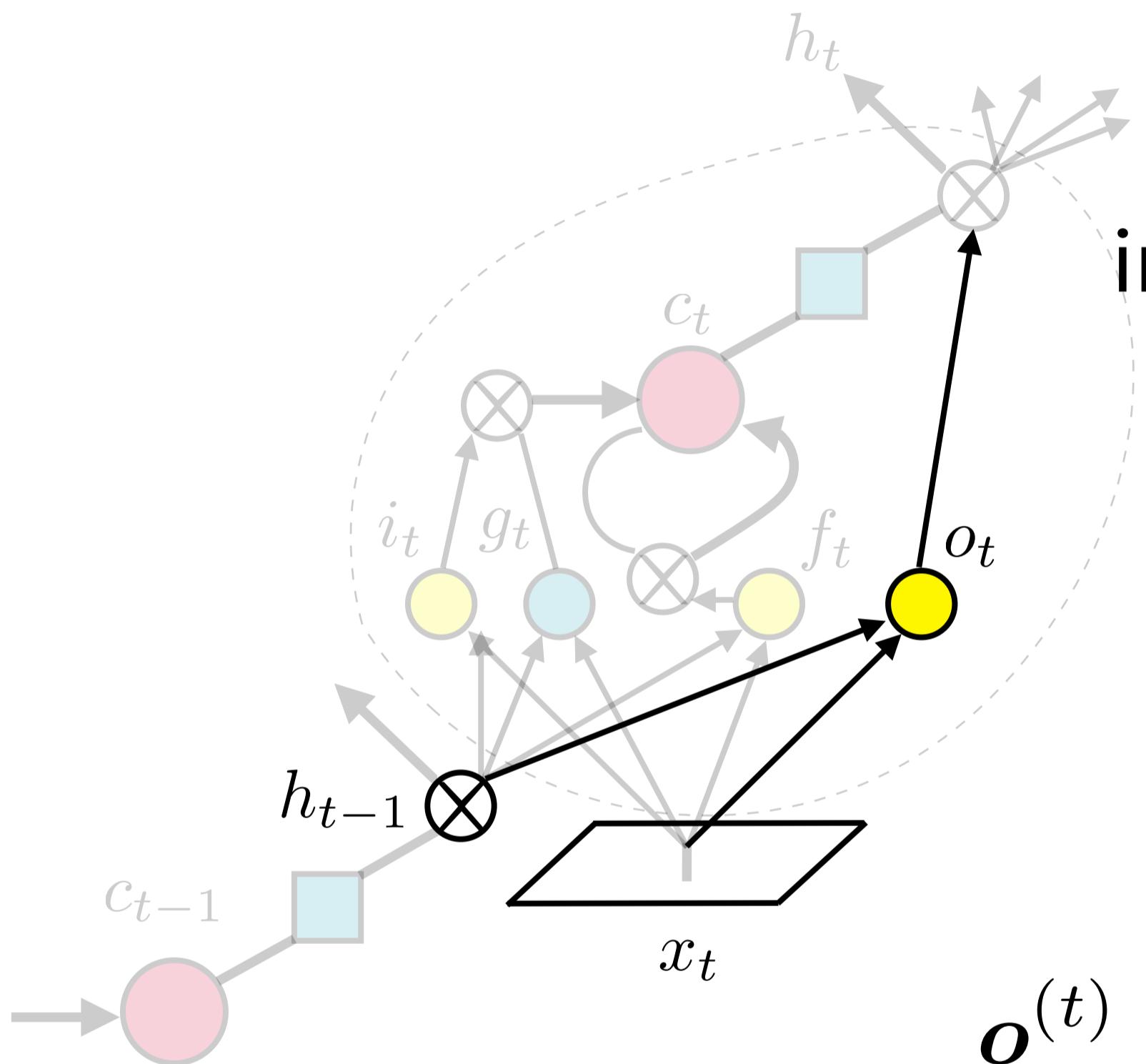
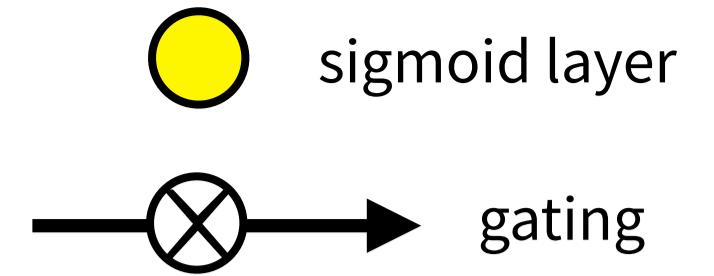
- sigmoid layer
- $\tanh$  layer
- memory cell
-  gating



The cell is updated with a combination of new and existing information

$$c^{(t)} = f^{(t)} \odot c^{(t-1)} + i^{(t)} \odot g^{(t)}$$

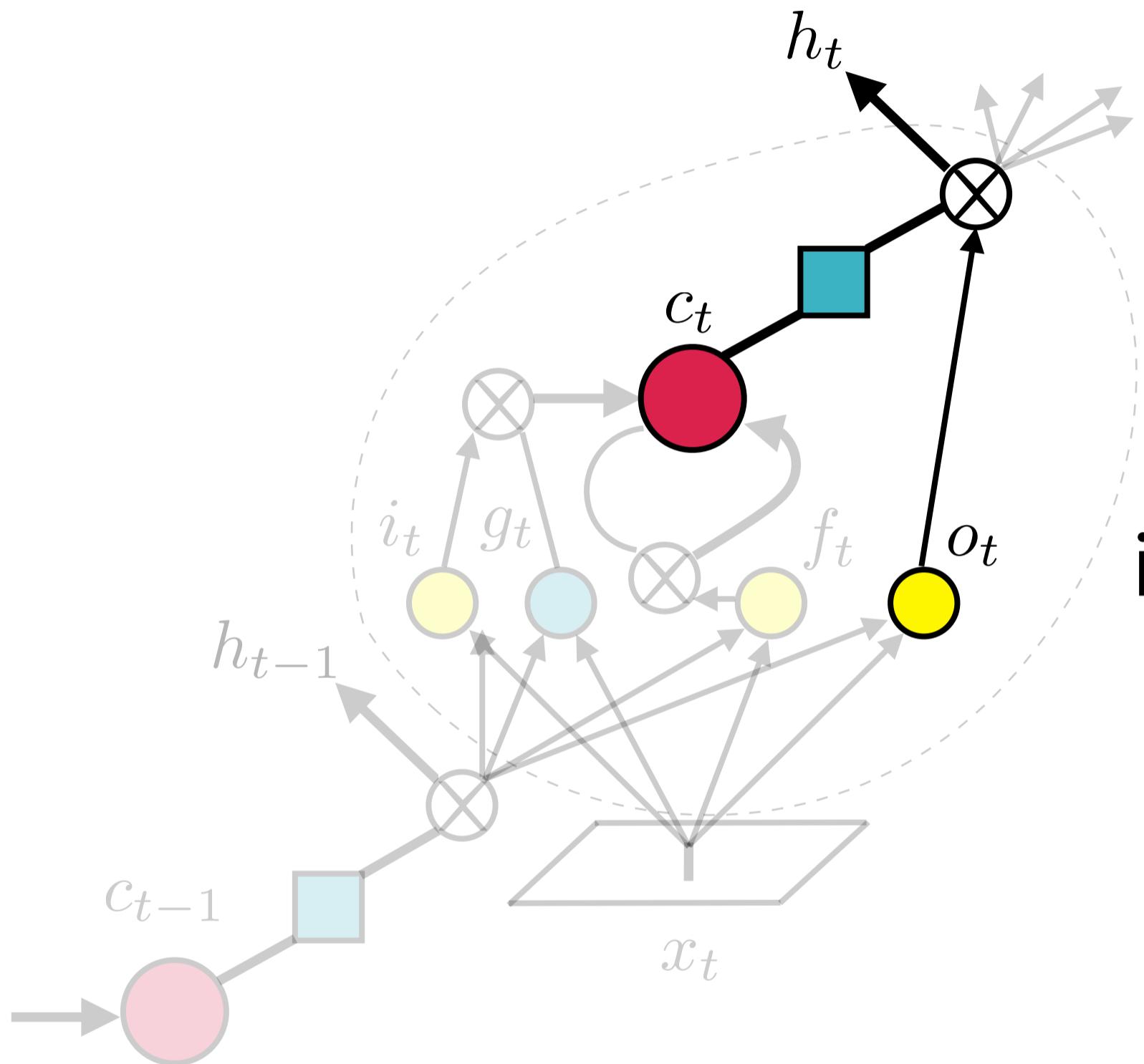
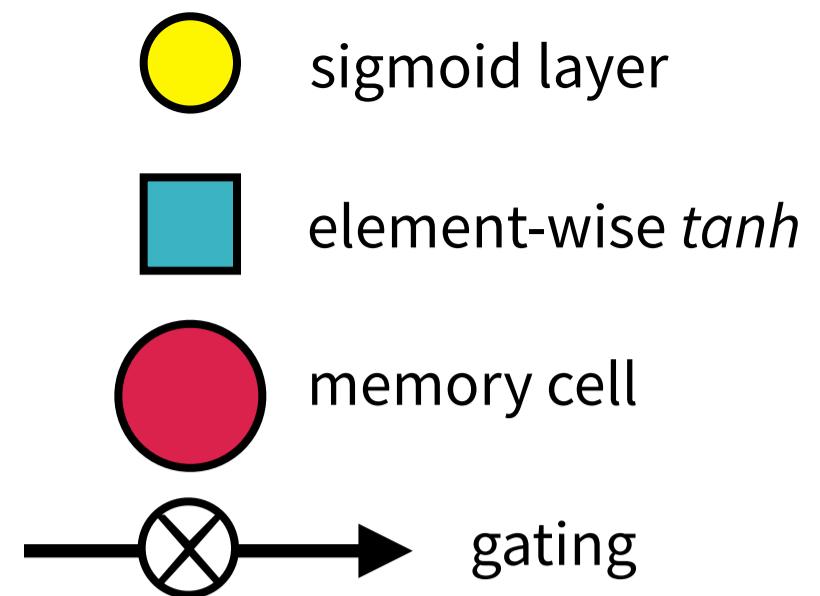
# LSTM - Output gate



The output gate allows information to flow out of the cell

$$o^{(t)} = \sigma \left( b^o + W^o [h^{(t-1)}, x^{(t)}] \right)$$

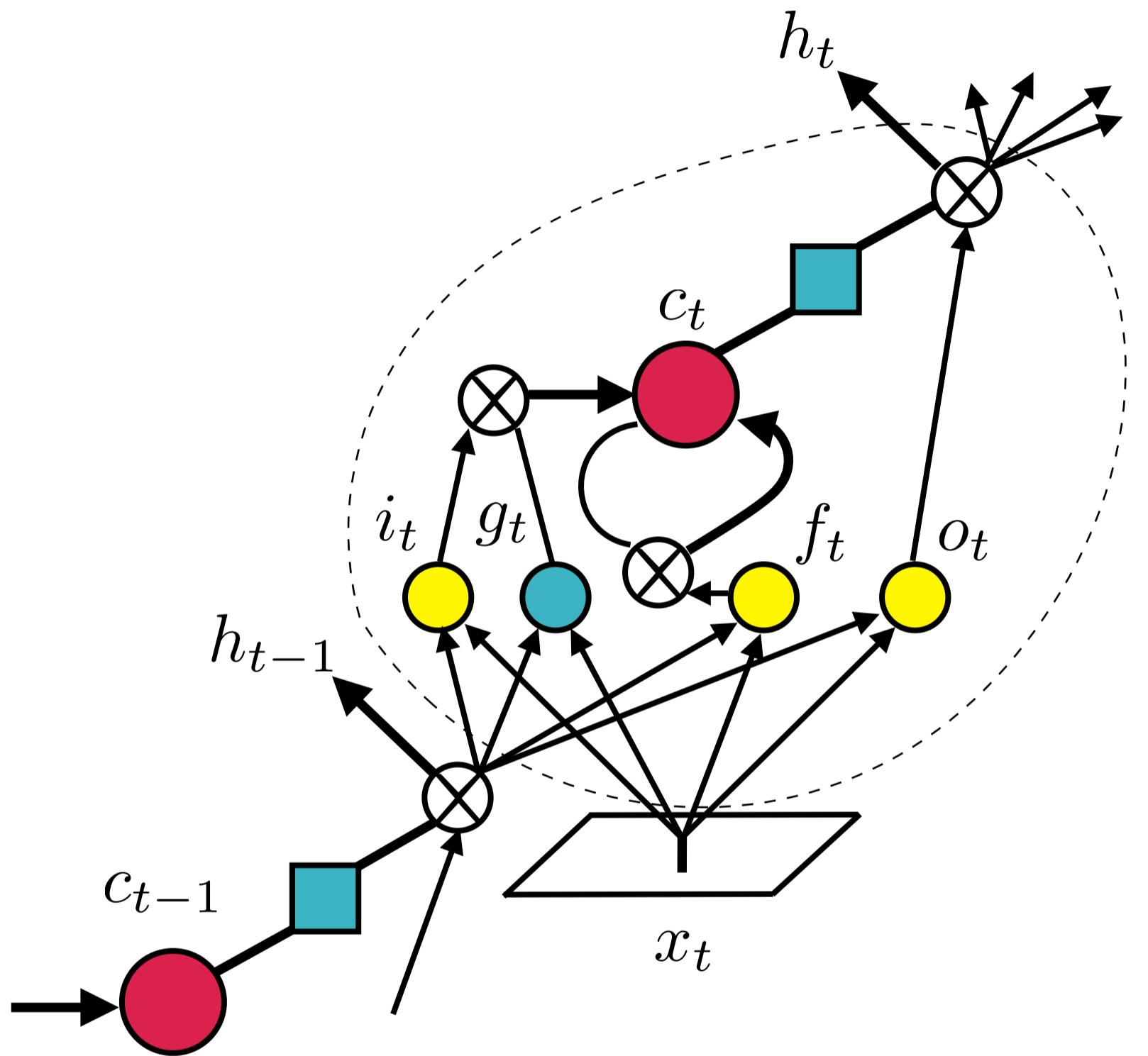
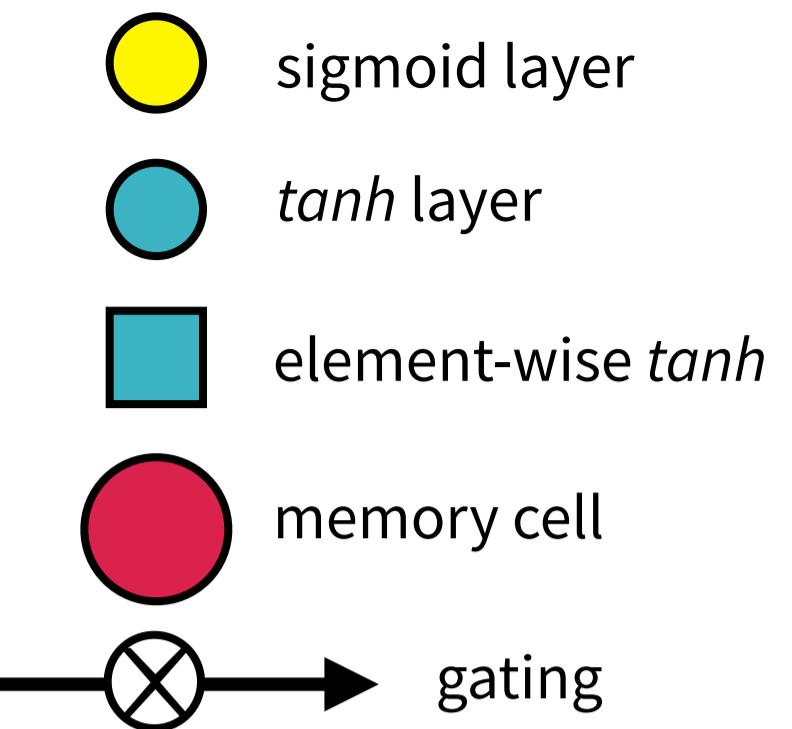
# LSTM - New state



The hidden state update  
is a nonlinear function of the cell  
regulated by the output gate

$$h^{(t)} = o^{(t)} \odot \tanh(c^{(t)})$$

# LSTM (recap)



$$\begin{aligned}
 i^{(t)} &= \sigma \left( \mathbf{b}^i + \mathbf{W}^i \left[ \mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 g^{(t)} &= \tanh \left( \mathbf{b}^g + \mathbf{W}^g \left[ \mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 f^{(t)} &= \sigma \left( \mathbf{b}^f + \mathbf{W}^f \left[ \mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{c}^{(t)} &= f^{(t)} \odot \mathbf{c}^{(t-1)} + i^{(t)} \odot g^{(t)} \\
 o^{(t)} &= \sigma \left( \mathbf{b}^o + \mathbf{W}^o \left[ \mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{h}^{(t)} &= o^{(t)} \odot \tanh(\mathbf{c}^{(t)})
 \end{aligned}$$

Recommended reading:

Chris Olah: Understanding LSTMs (blog)

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

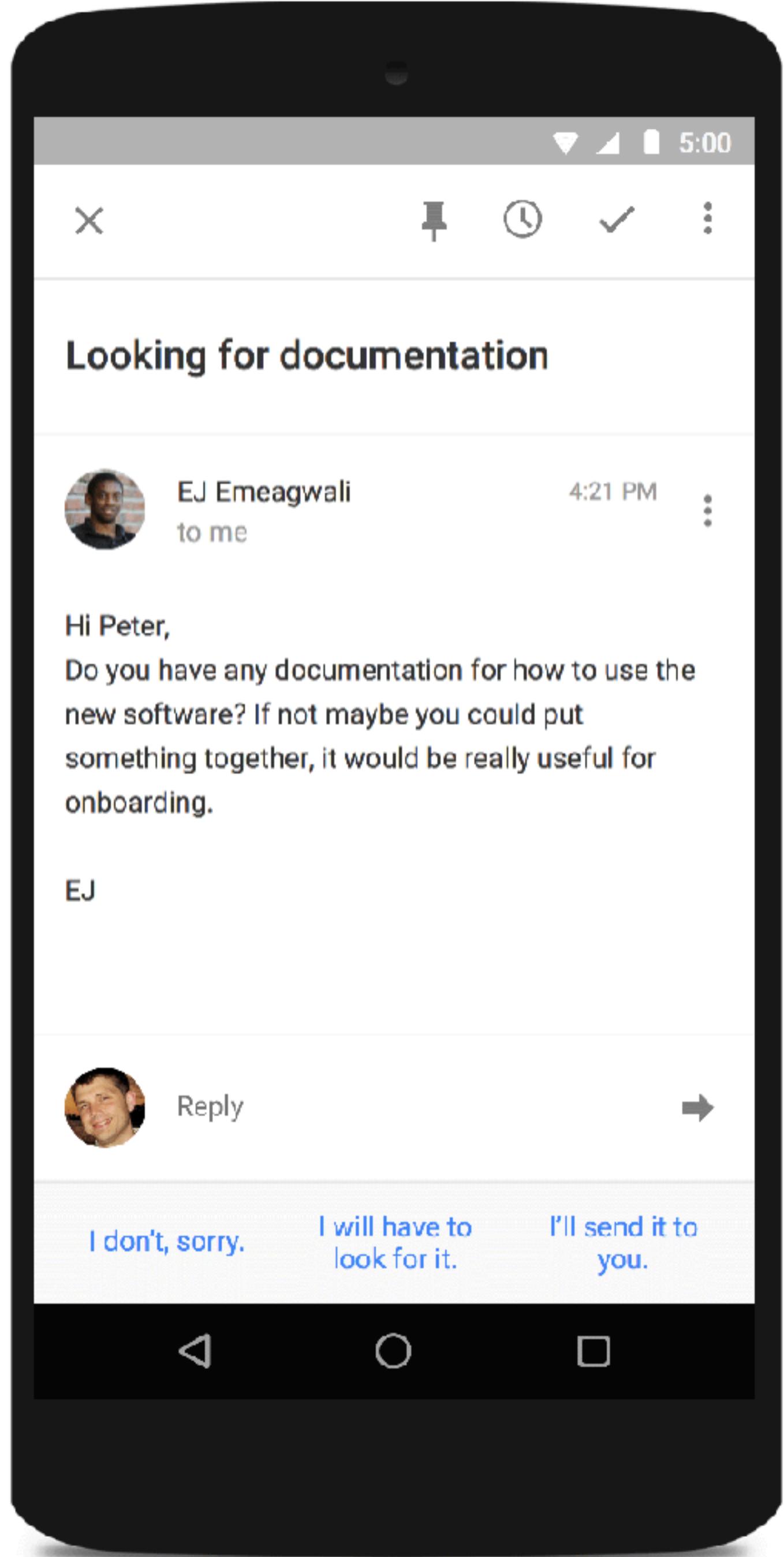


Figure: Google

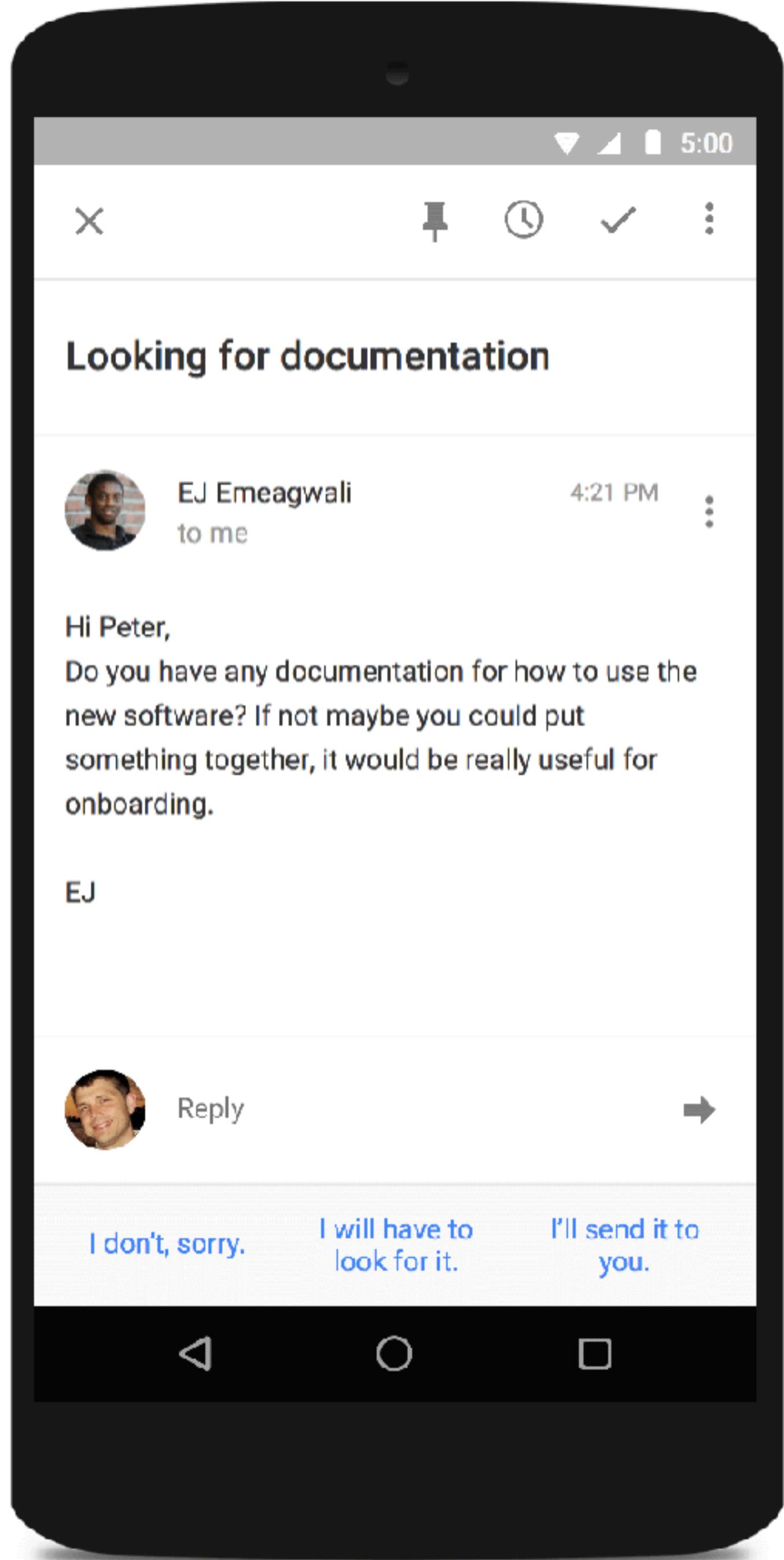
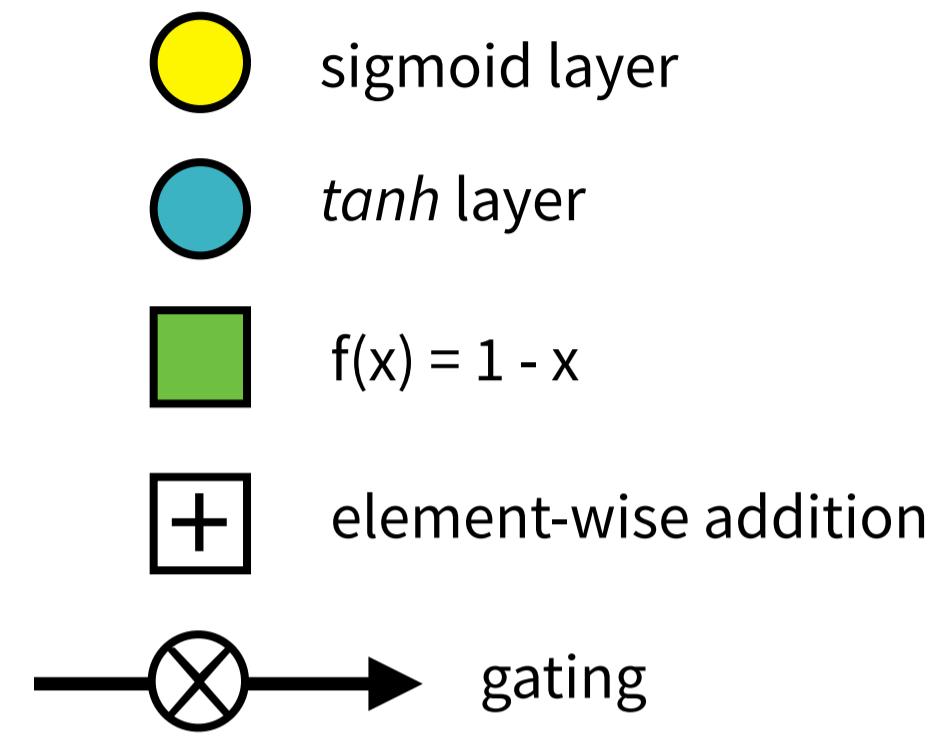
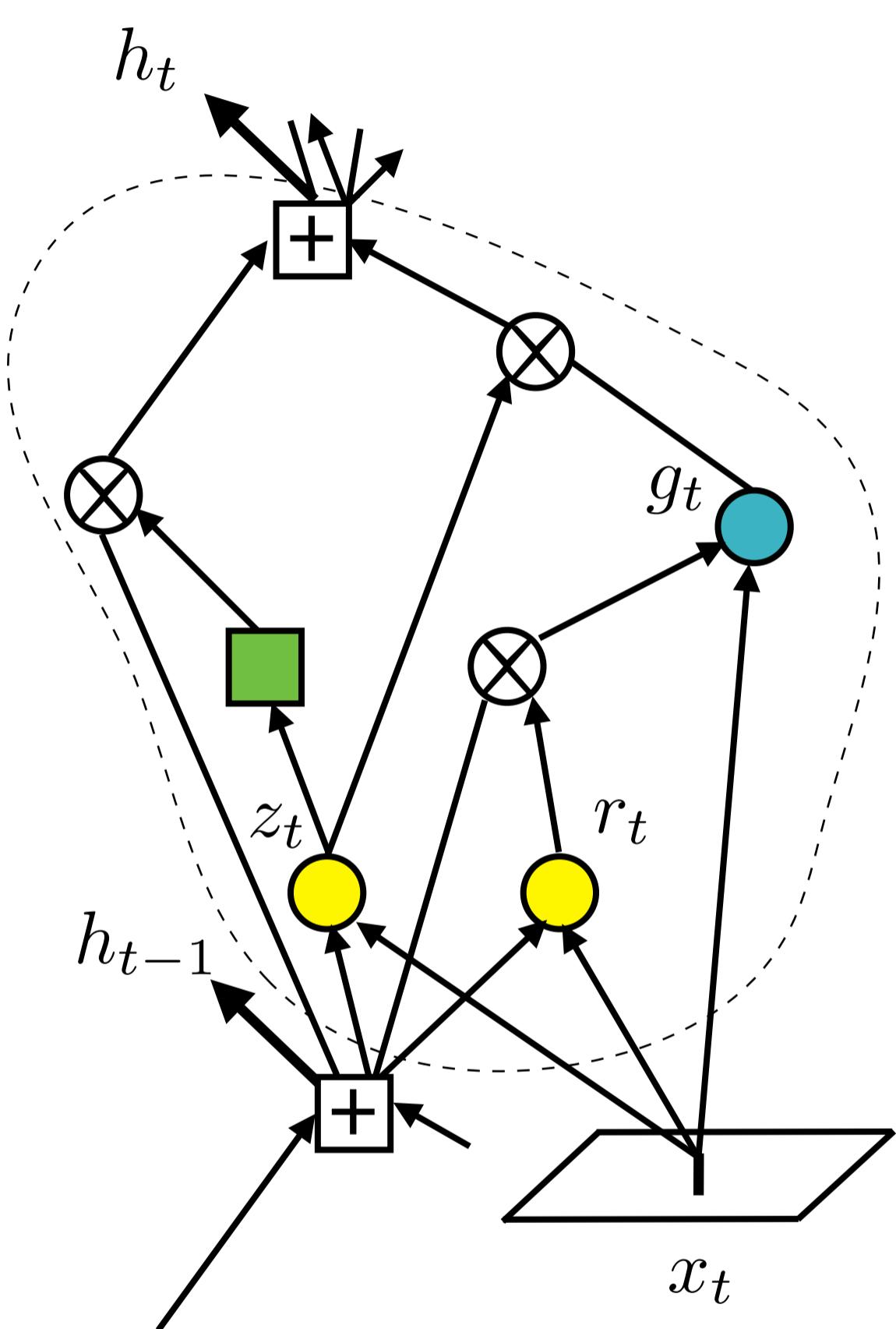


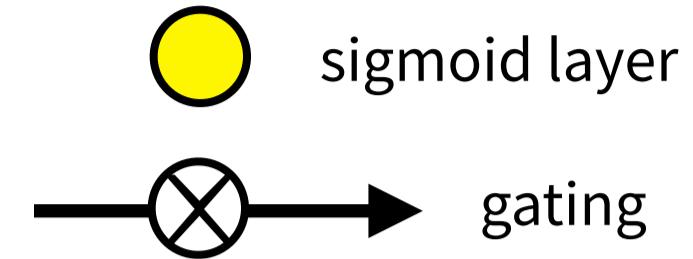
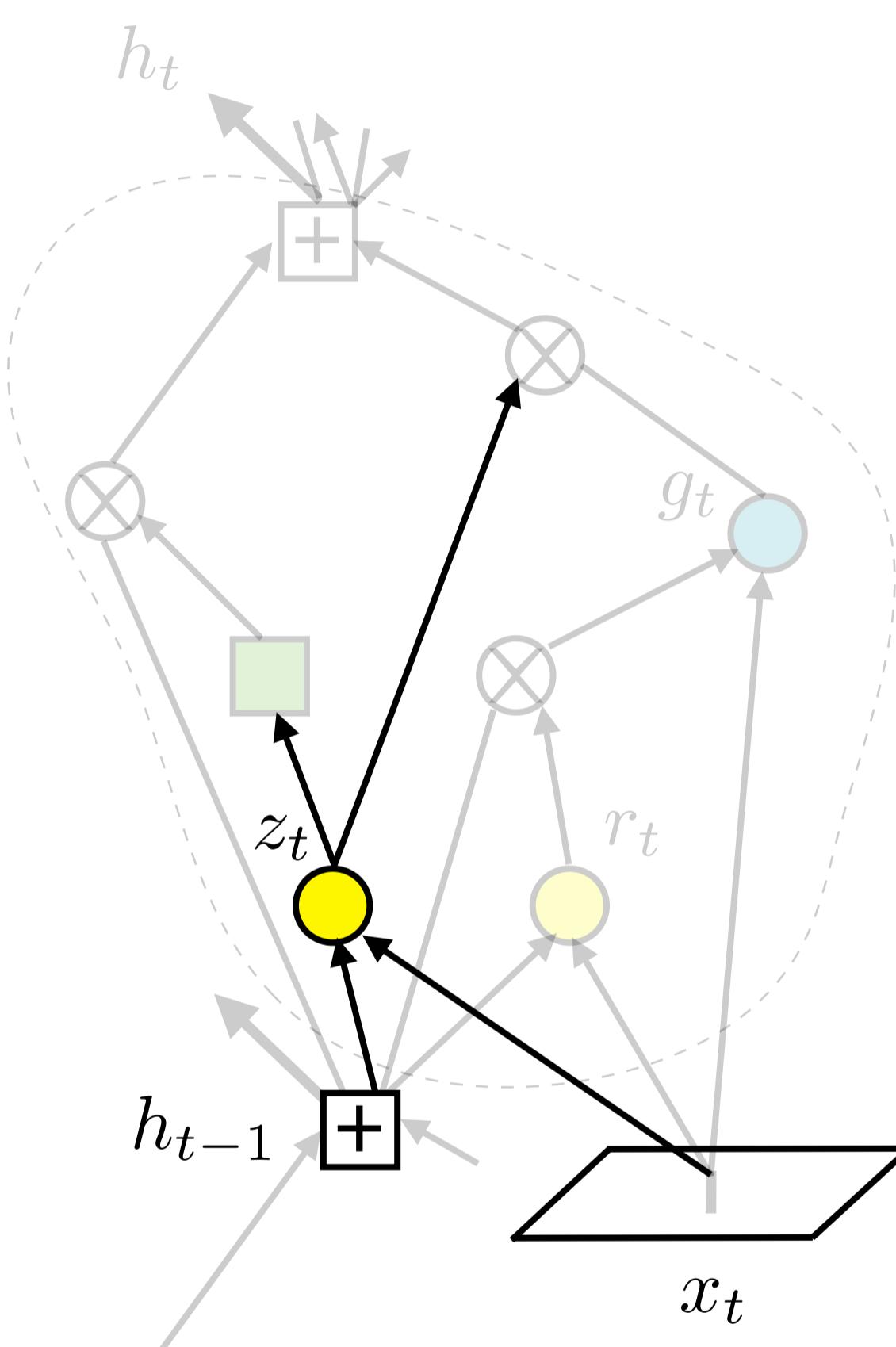
Figure: Google

# Gated Recurrent Units (GRU)



Combine forget and input gates  
into a single “update gate”  
Merges cell state and hidden state

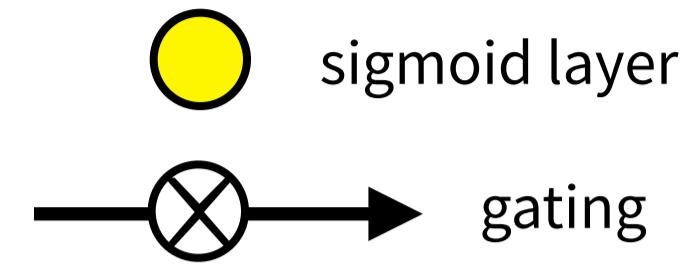
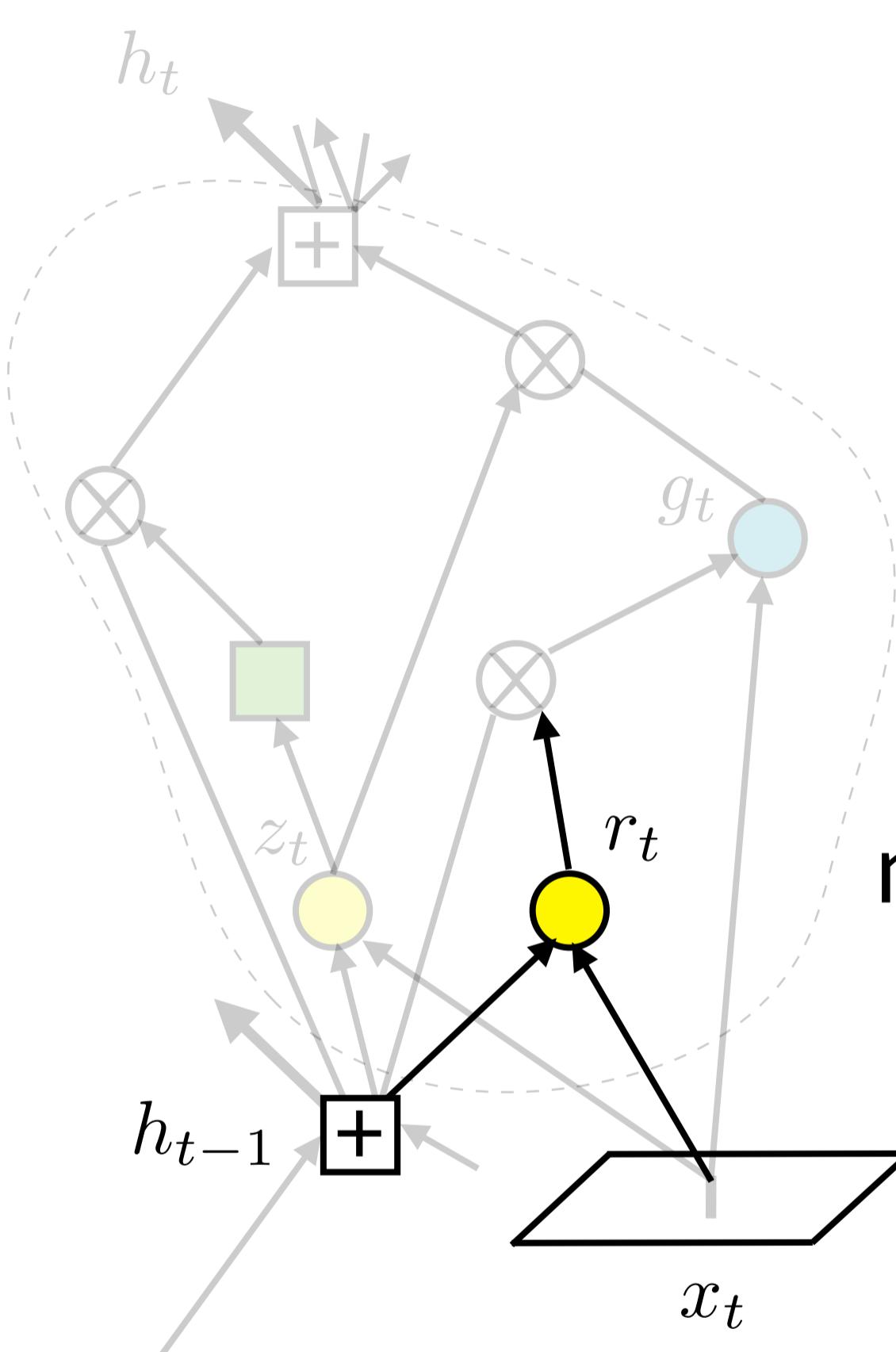
# GRU - Update gate



The update gate simultaneously controls the forgetting factor and the decision to update the state

$$z^{(t)} = \sigma(b^z + W^z [h^{(t-1)}, x^{(t)}])$$

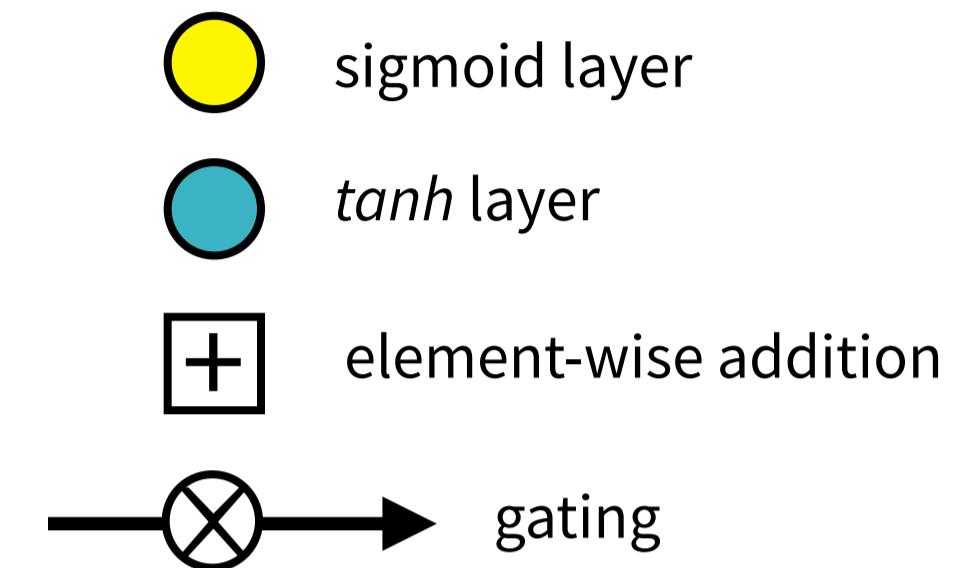
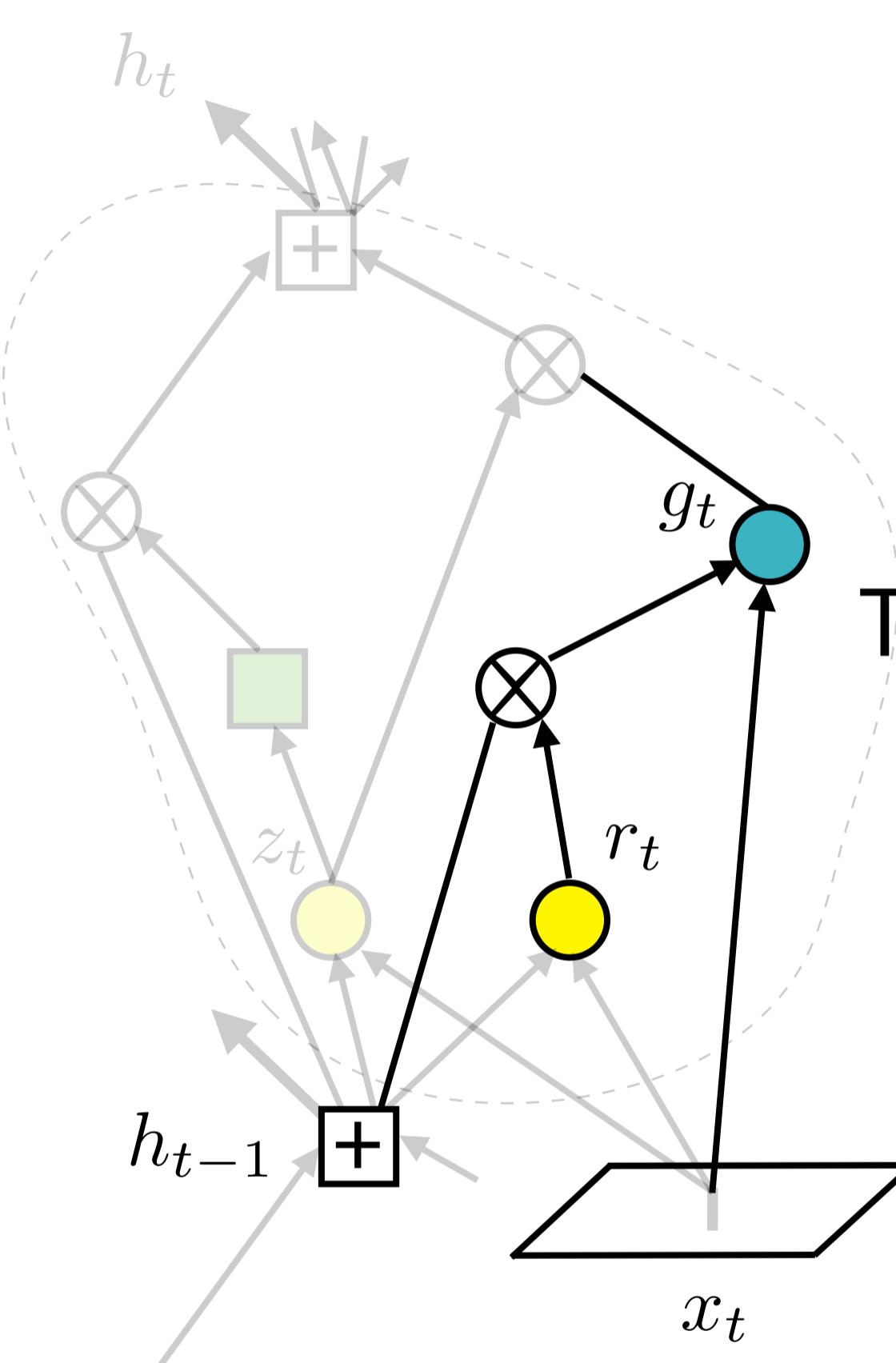
# GRU - Reset gate



The reset gate controls which part of the old state is used to compute the next state, introducing an additional nonlinear effect b/w past and present

$$r^{(t)} = \sigma \left( b^r + W^r [h^{(t-1)}, x^{(t)}] \right)$$

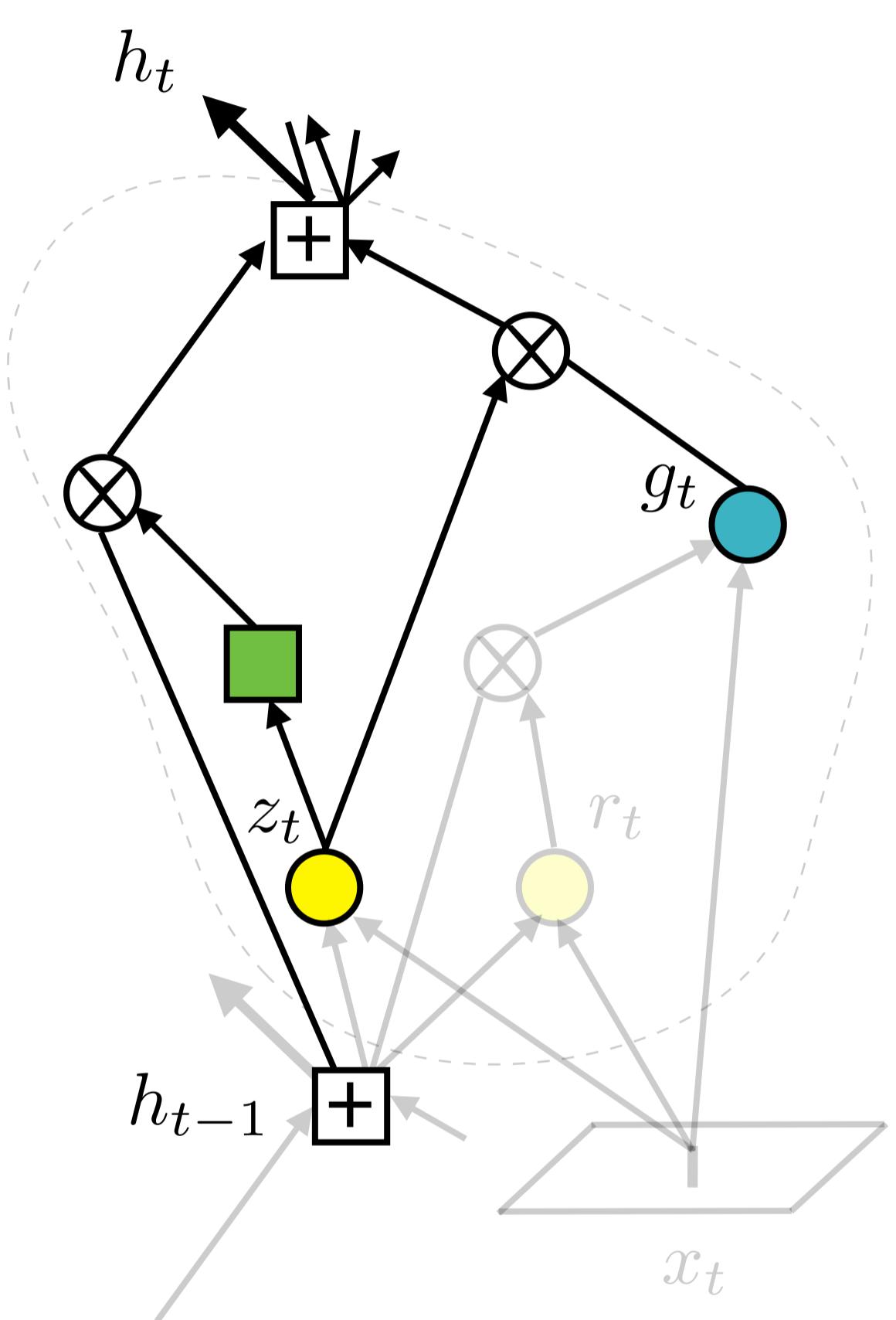
# GRU - Target state



The target state is a nonlinear function  
of the previous state and input

$$g^{(t)} = \tanh \left( b^g + W^g [r^{(t)} \odot h^{(t-1)}, x^{(t)}] \right)$$

# GRU - New state

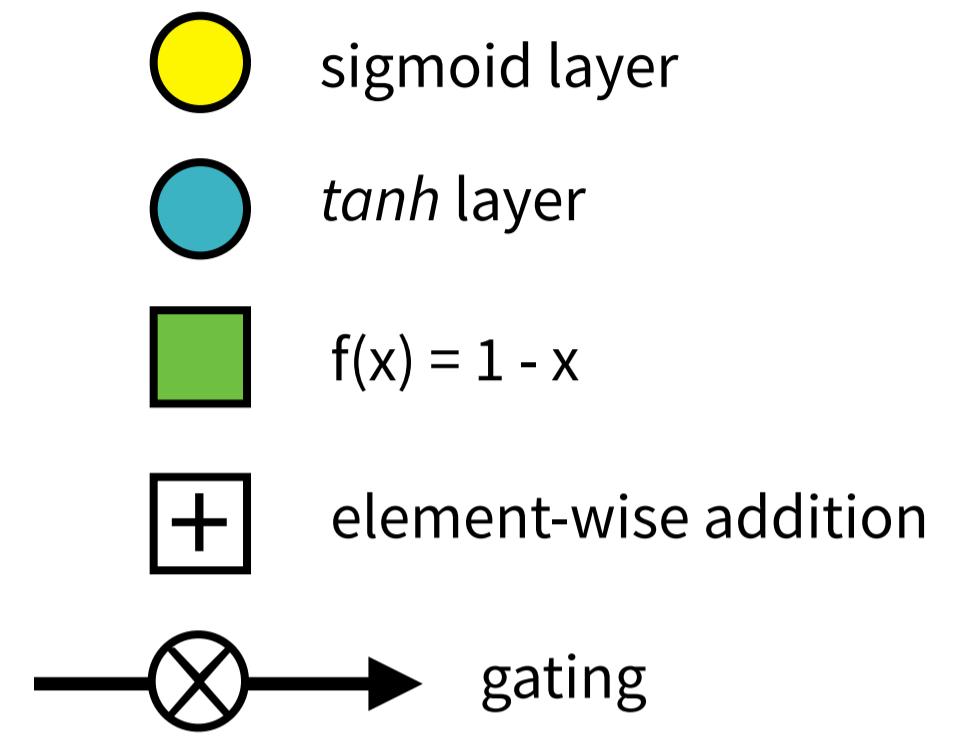
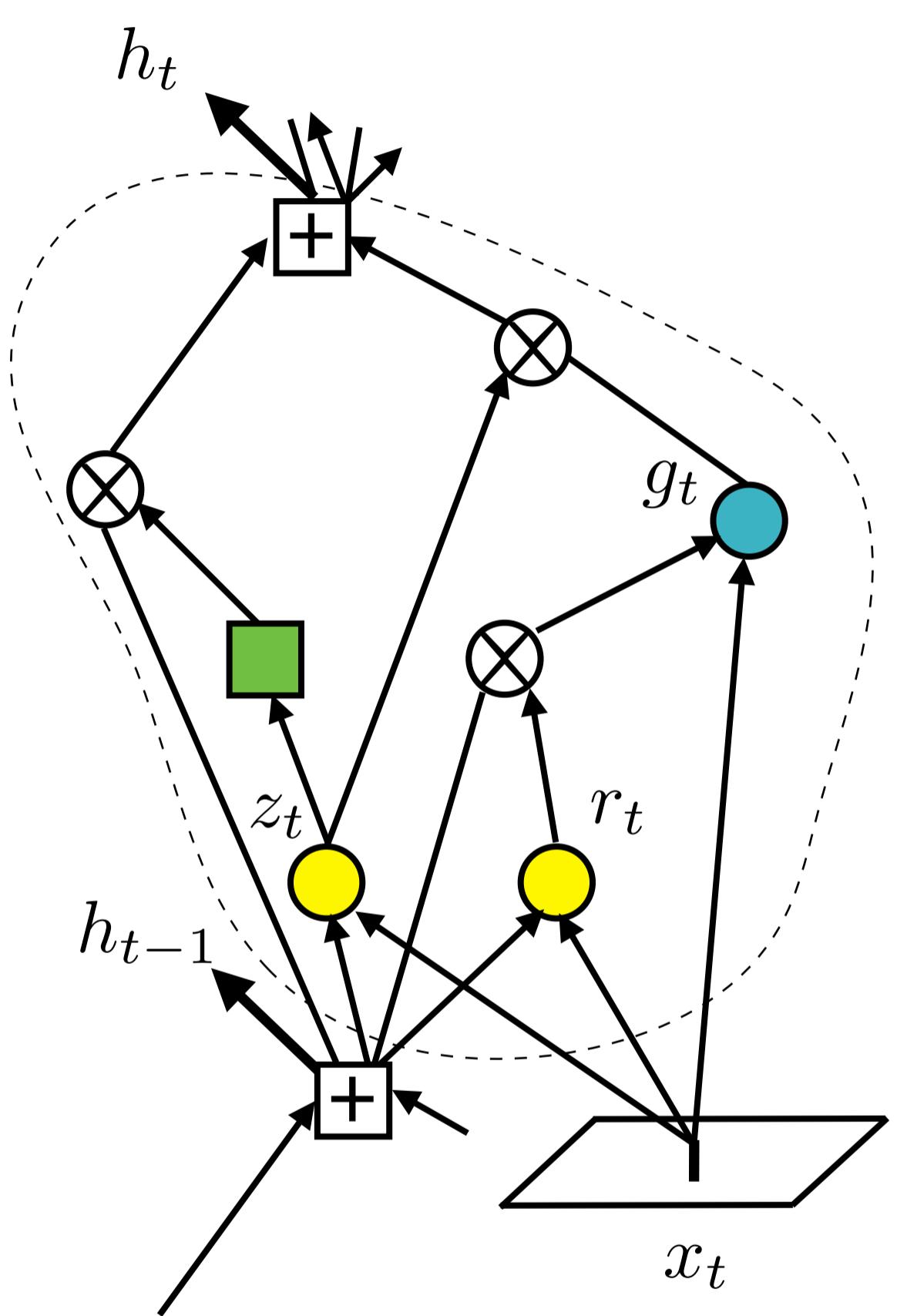


- (Yellow circle) sigmoid layer
- (Cyan circle)  $tanh$  layer
- (Green square)  $f(x) = 1 - x$
- (Plus sign in a box) element-wise addition
- (Crossed-out arrow) gating

The new state is a function of  
the target state and previous state

$$h^{(t)} = (1 - z^{(t)}) \odot h^{(t-1)} + z^{(t)} \odot g^{(t)}$$

# GRU (recap)

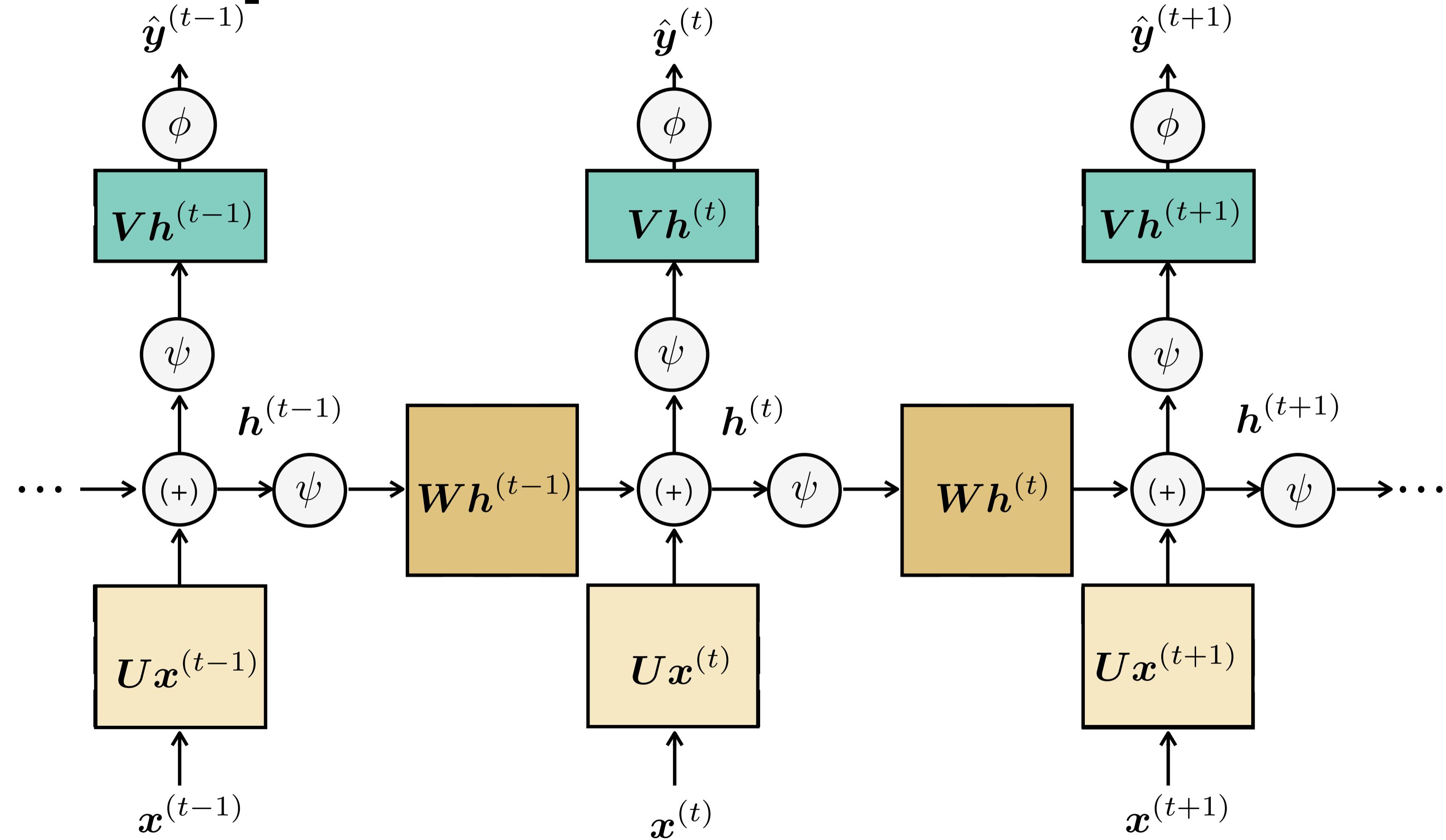


$$\begin{aligned}
 z^{(t)} &= \sigma \left( \mathbf{b}^z + \mathbf{W}^z \left[ h^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 r^{(t)} &= \sigma \left( \mathbf{b}^r + \mathbf{W}^r \left[ h^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 g^{(t)} &= \tanh \left( \mathbf{b}^g + \mathbf{W}^g \left[ r^{(t)} \odot h^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 h^{(t)} &= (1 - z^{(t)}) \odot h^{(t-1)} + z^{(t)} \odot g^{(t)}
 \end{aligned}$$

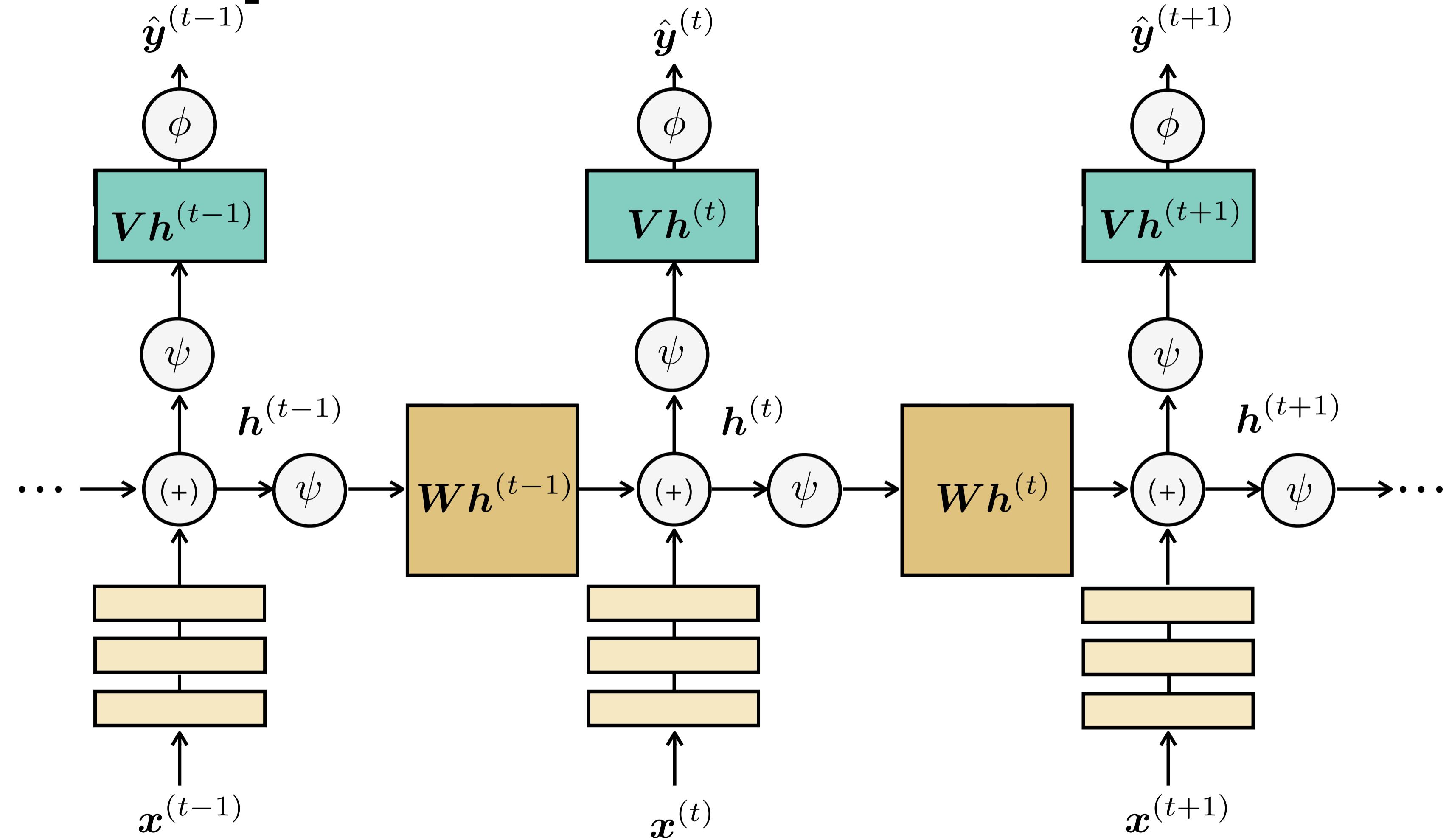
# LSTM: A Search Space Odyssey

- So many possibilities!
- Several investigations over architectural variations of the LSTM and GRU found that **no variant beat the “baseline”** over a wide variety of tasks (Greff et al., 2015), (Jozefowicz et al., 2015)
- Grief et al. (2015) found the **forget gate** and **output activation function** to be the LSTM’s most critical parts

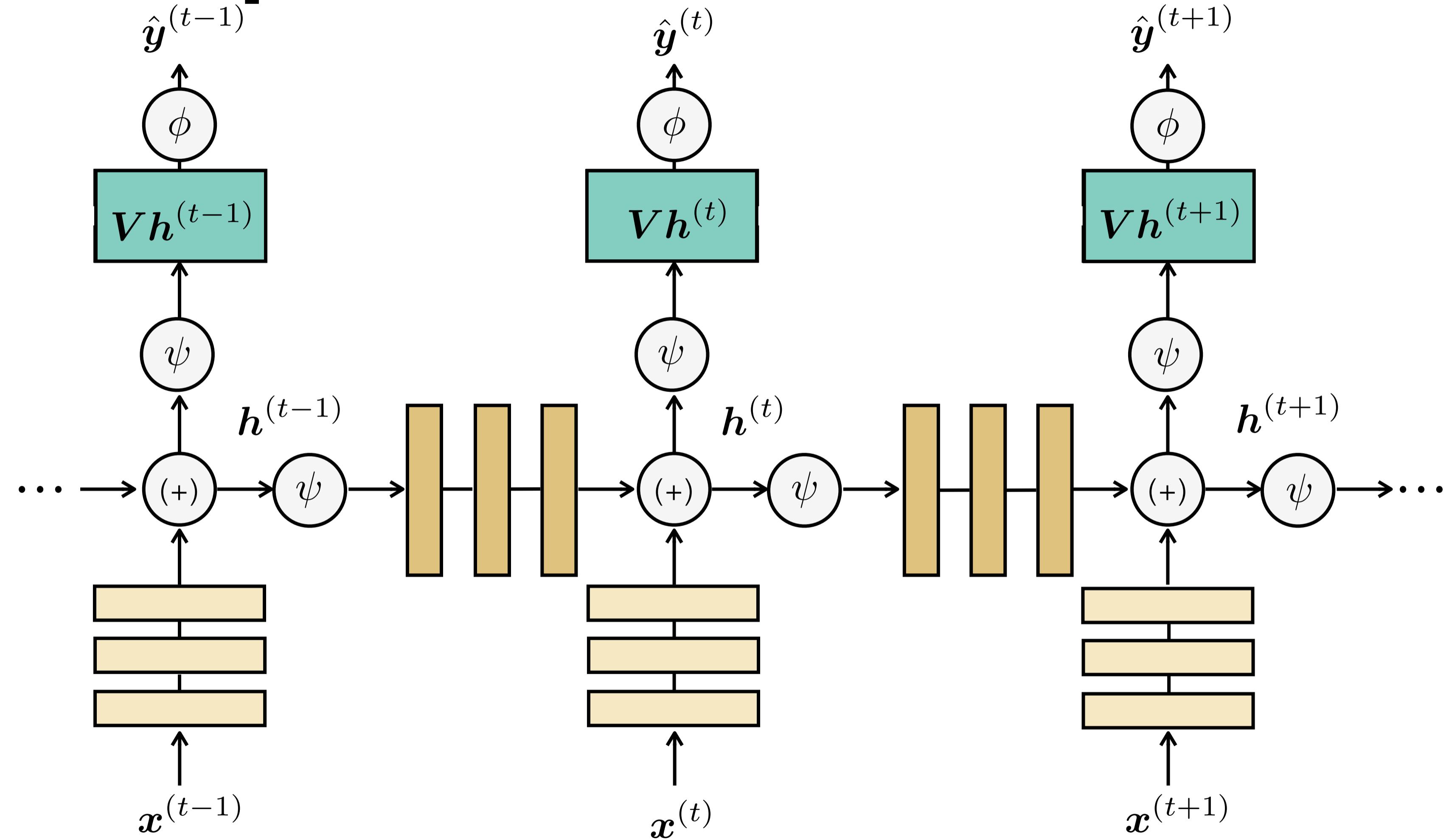
# Deep RNNs



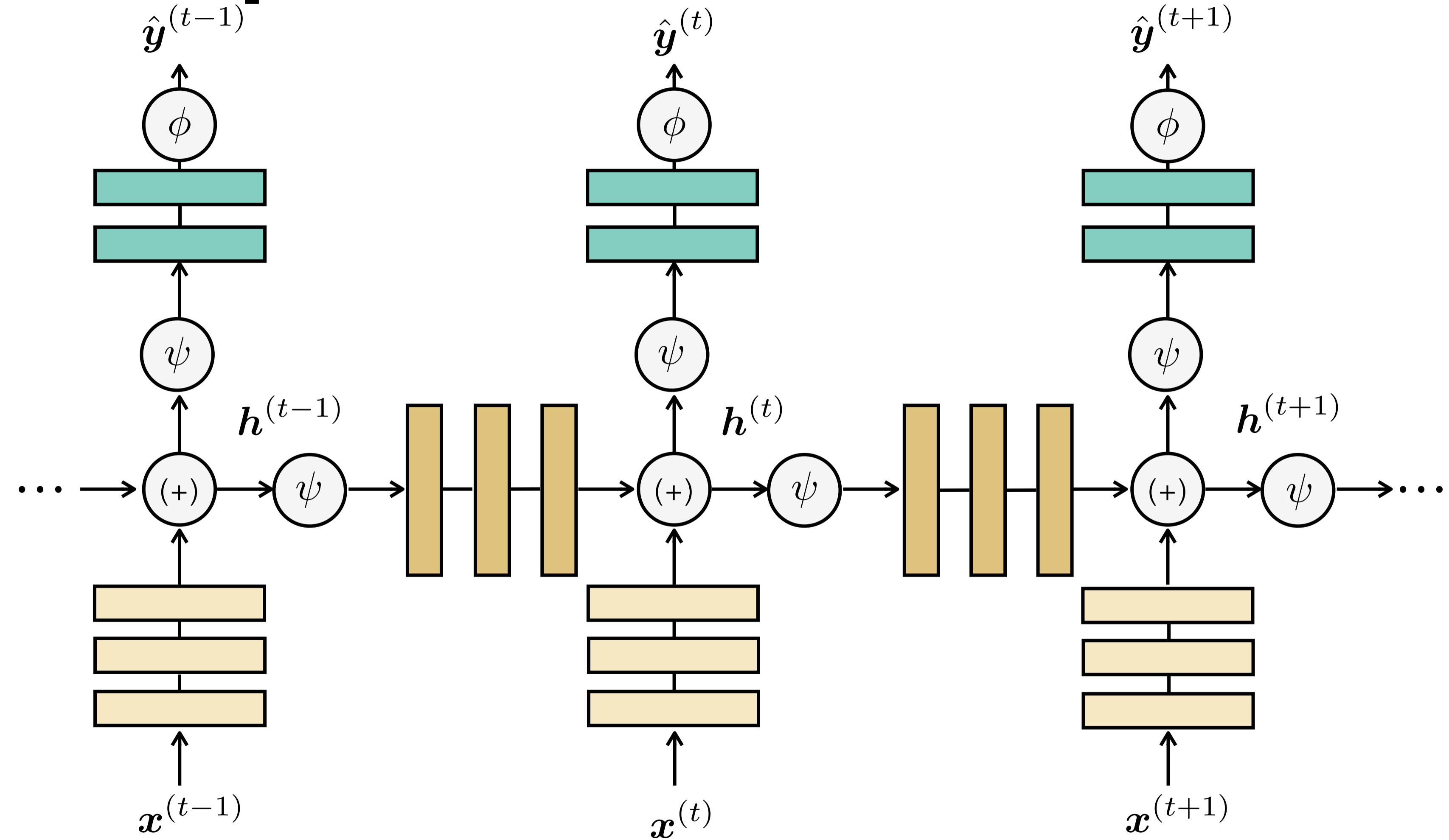
# Deep RNNs



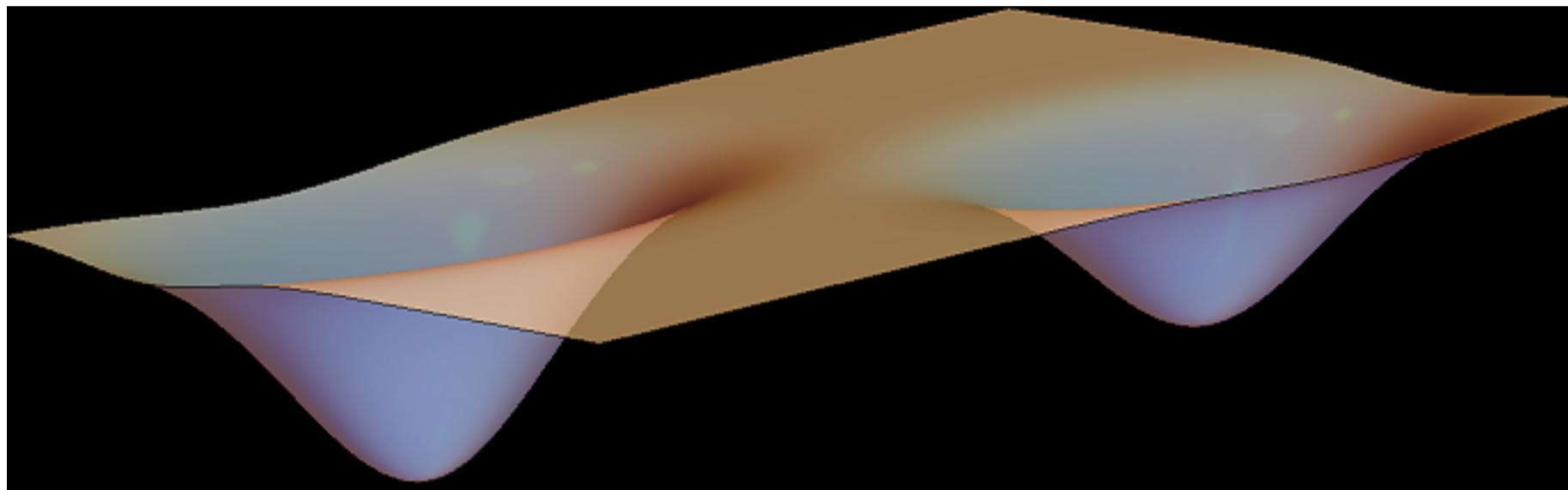
# Deep RNNs



# Deep RNNs



# Solutions: Optimization



- 2<sup>nd</sup>-order methods (Martens and Sutskever 2011)
- Reservoirs (Jaeger and Haas 2004)
- Gradient clipping/normalization (Pascanu et al. 2012, Mikolov 2012)
- Careful initialization (Sutskever et al. 2013)

# Summary (and Advice!)

- Keep the architecture **simple!**
  - GRUs, gradient clipping, careful initialization
- Use well-tested, well-supported **open-source tools:**
  - Tensor Flow, PyTorch, Caffe
- Don't forget the **flexibility** of input-output structure