



POLITECNICO
MILANO 1863



Advances in Deep Learning with Applications in Text and Image Processing

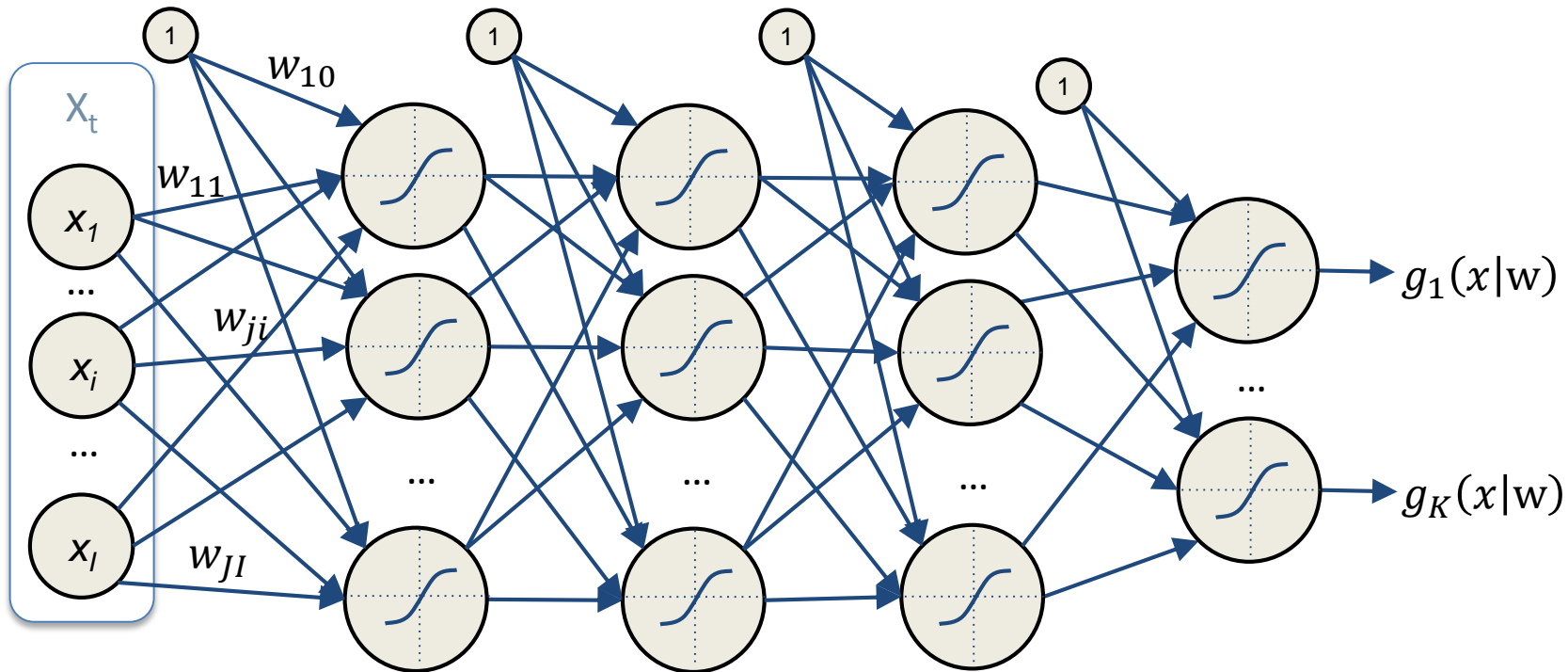
- Recurrent Neural Networks -

Prof. Matteo Matteucci – *matteo.matteucci@polimi.it*

Department of Electronics, Information and Bioengineering
Artificial Intelligence and Robotics Lab - Politecnico di Milano

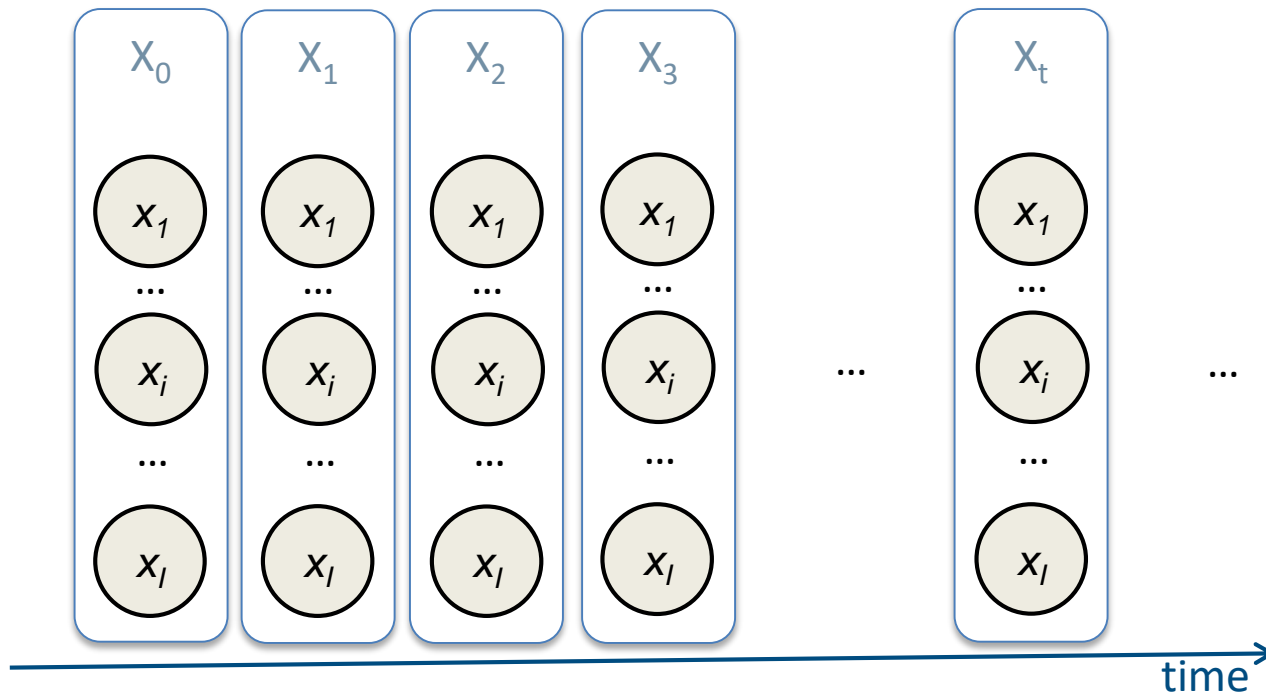
Sequence Modeling

So far we have considered only «static» datasets



Sequence Modeling

So far we have considered only «static» datasets



Sequence Modeling

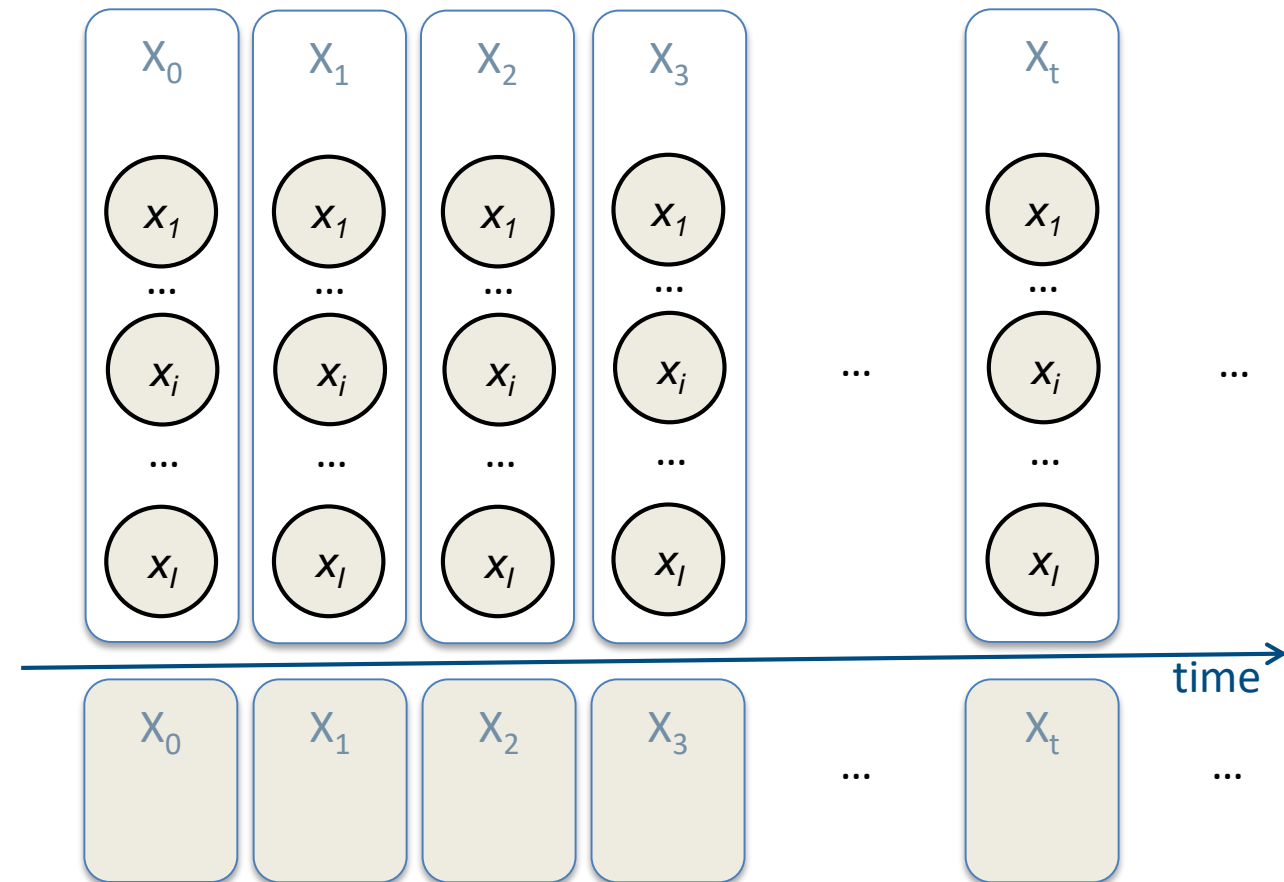
Different ways to deal with «dynamic» data:

Memoryless models:

- Autoregressive models
- Feedforward neural networks

Models with memory:

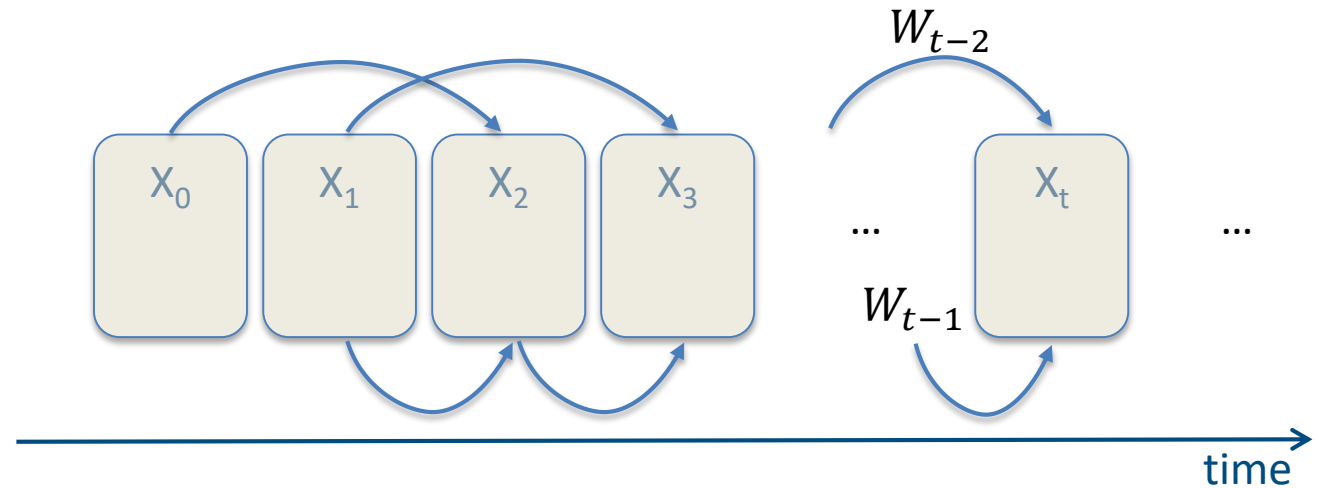
- Linear dynamical systems
- Hidden Markov models
- Recurrent Neural Networks
- ...



Memoryless Models for Sequences

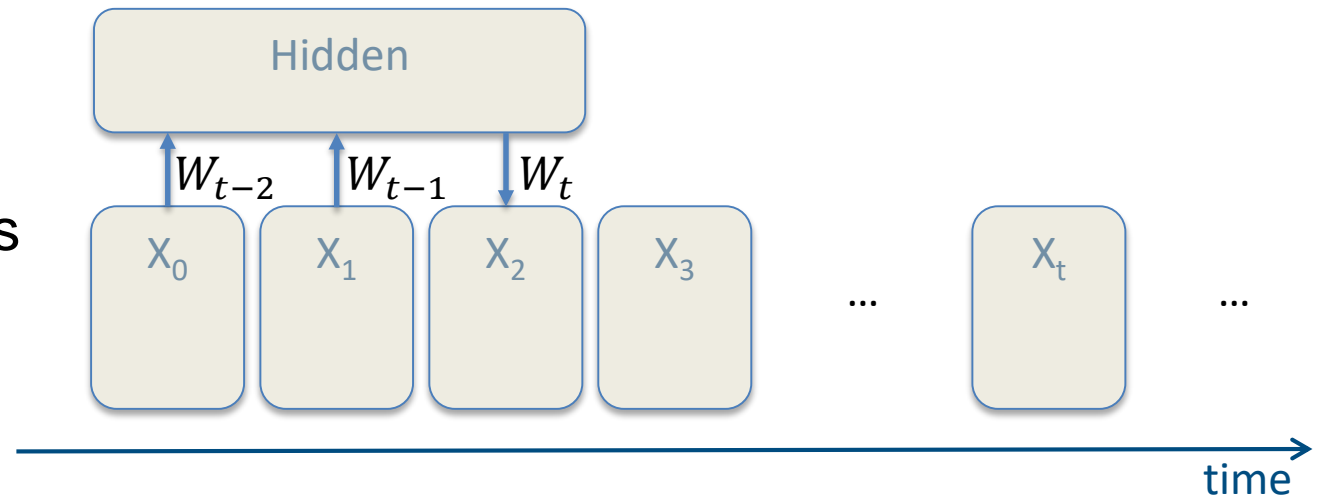
Autoregressive models

- Predict the next input from previous ones using «delay taps»



Feed forward neural networks

- Generalize autoregressive models using non linear hidden layers

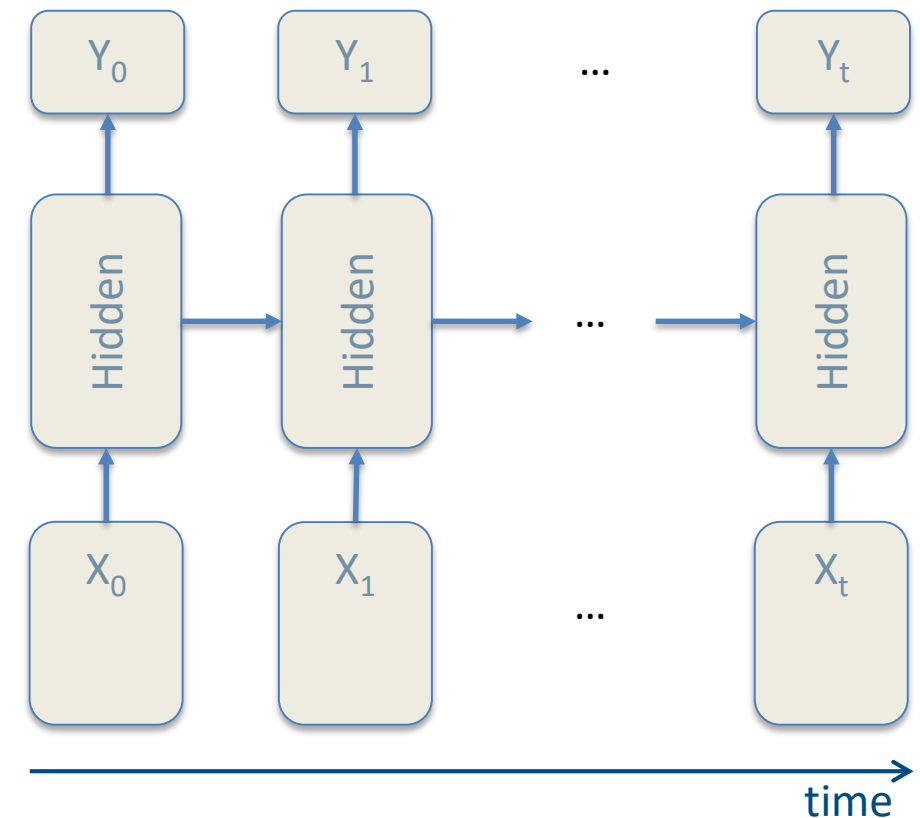


Generative models with a real-valued hidden state which cannot be observed directly

- The hidden state has some dynamics possibly affected by noise and produces the output
- To compute the output has to infer hidden state
- Input are treated as driving inputs

In linear dynamical systems this becomes:

- State continuous with Gaussian uncertainty
- Transformations are assumed to be linear
- State can be estimated using *Kalman filtering*

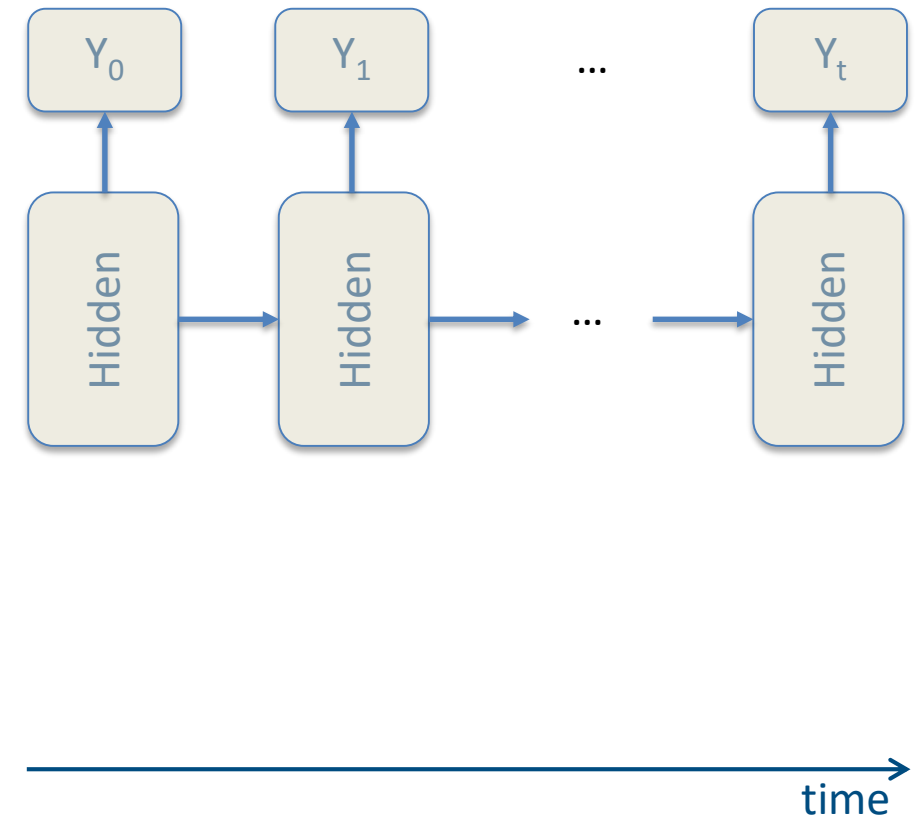


Generative models with a real-valued hidden state which cannot be observed directly

- The hidden state has some dynamics possibly affected by noise and produces the output
- To compute the output has to infer hidden state
- Input are treated as driving inputs

In hidden Markov models this becomes:

- State assumed to be discrete, state transitions are stochastic (transition matrix)
- Output is a stochastic function of hidden states
- State can be estimated via *Viterbi algorithm*.



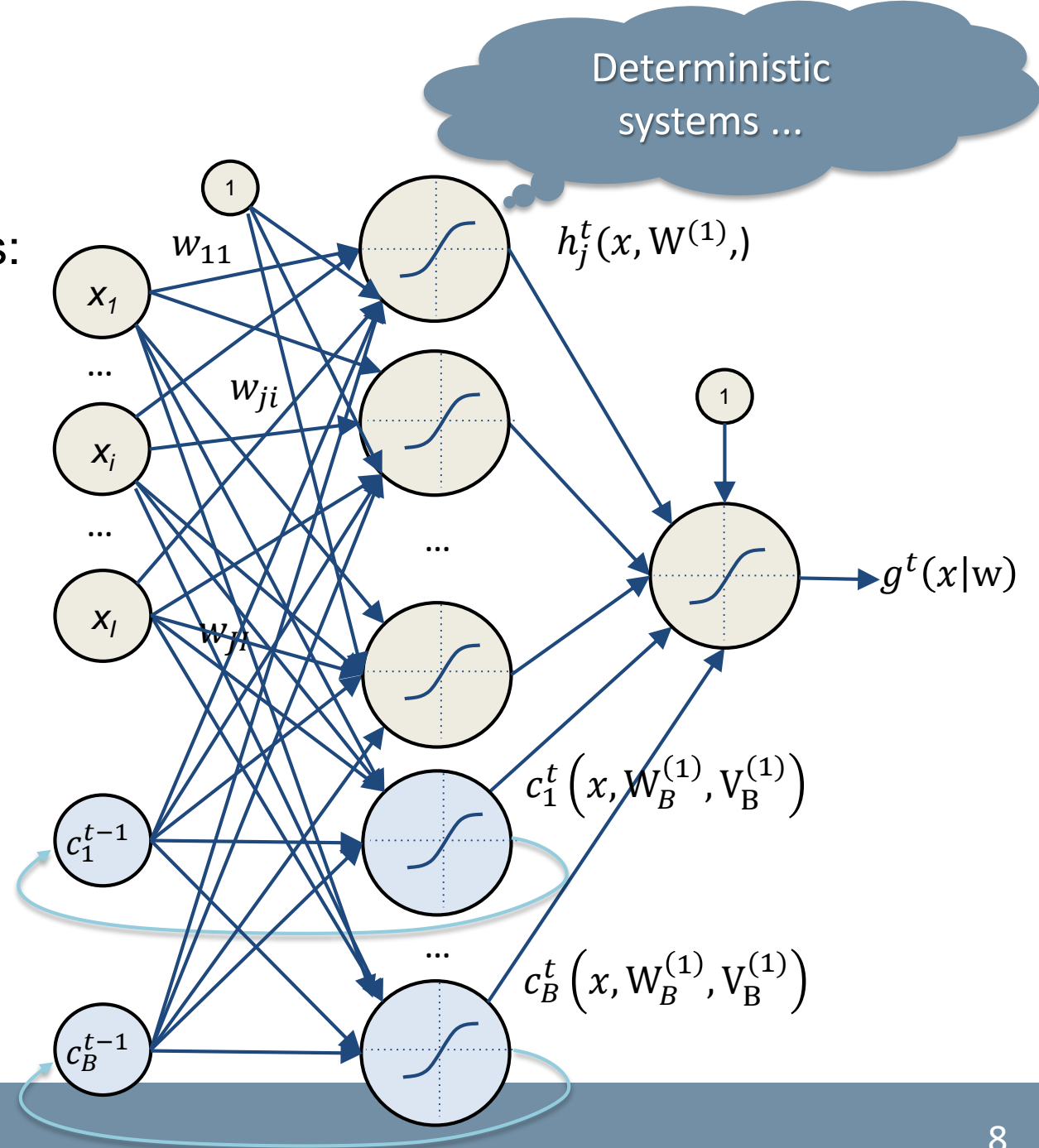
Recurrent Neural networks

Introduce memory by recurrent connections:

- Distributed hidden state allows to store a information efficiently
- Non-linear dynamics allows complex hidden state updates

“With enough neurons and time, RNNs can compute anything that can be computed by a computer.”

(Computation Beyond the Turing Limit
Hava T. Siegelmann, 1995)



Recurrent Neural networks

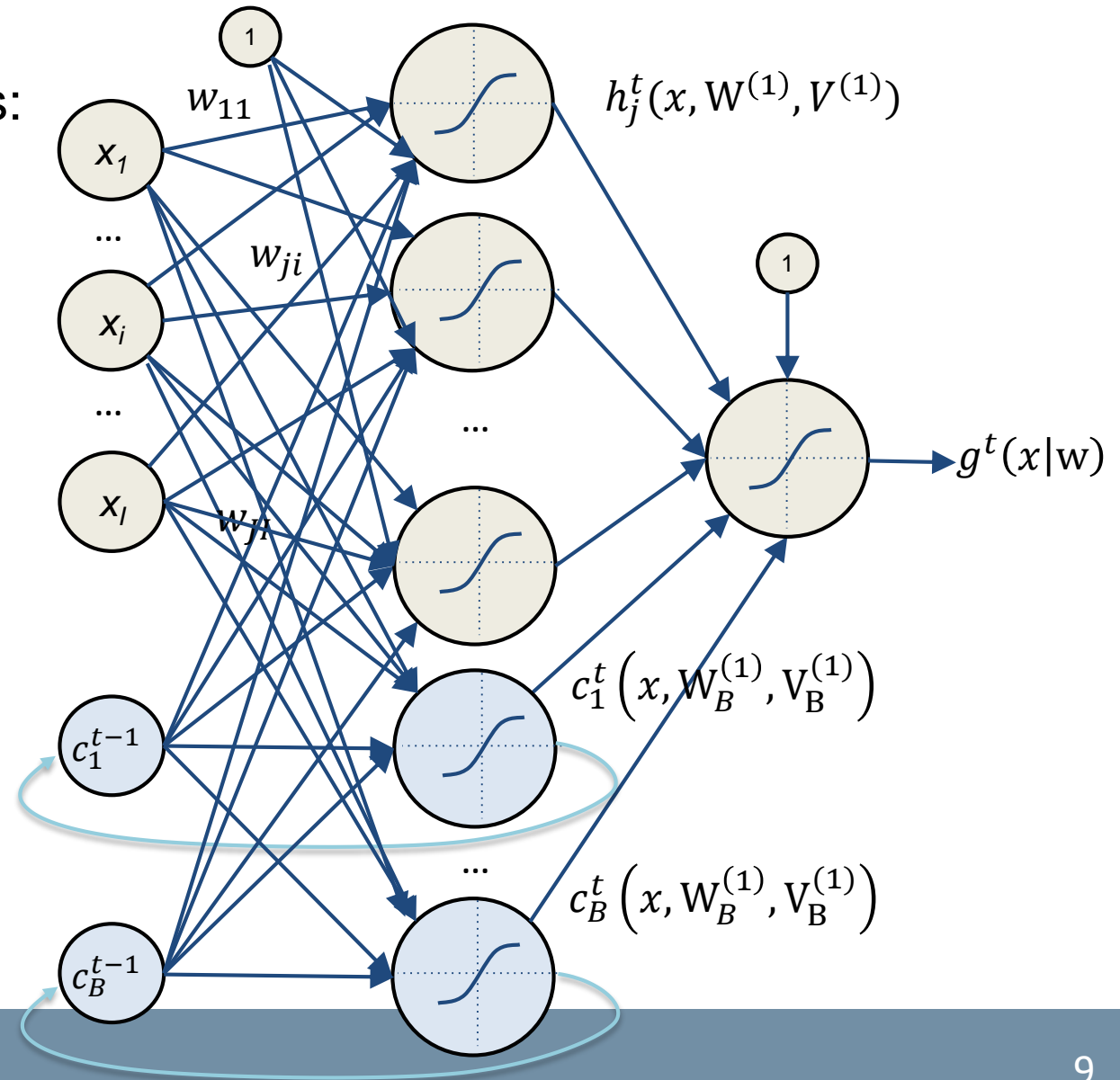
Introduce memory by recurrent connections:

- Distributed hidden state allows to store a information efficiently
- Non-linear dynamics allows complex hidden state updates

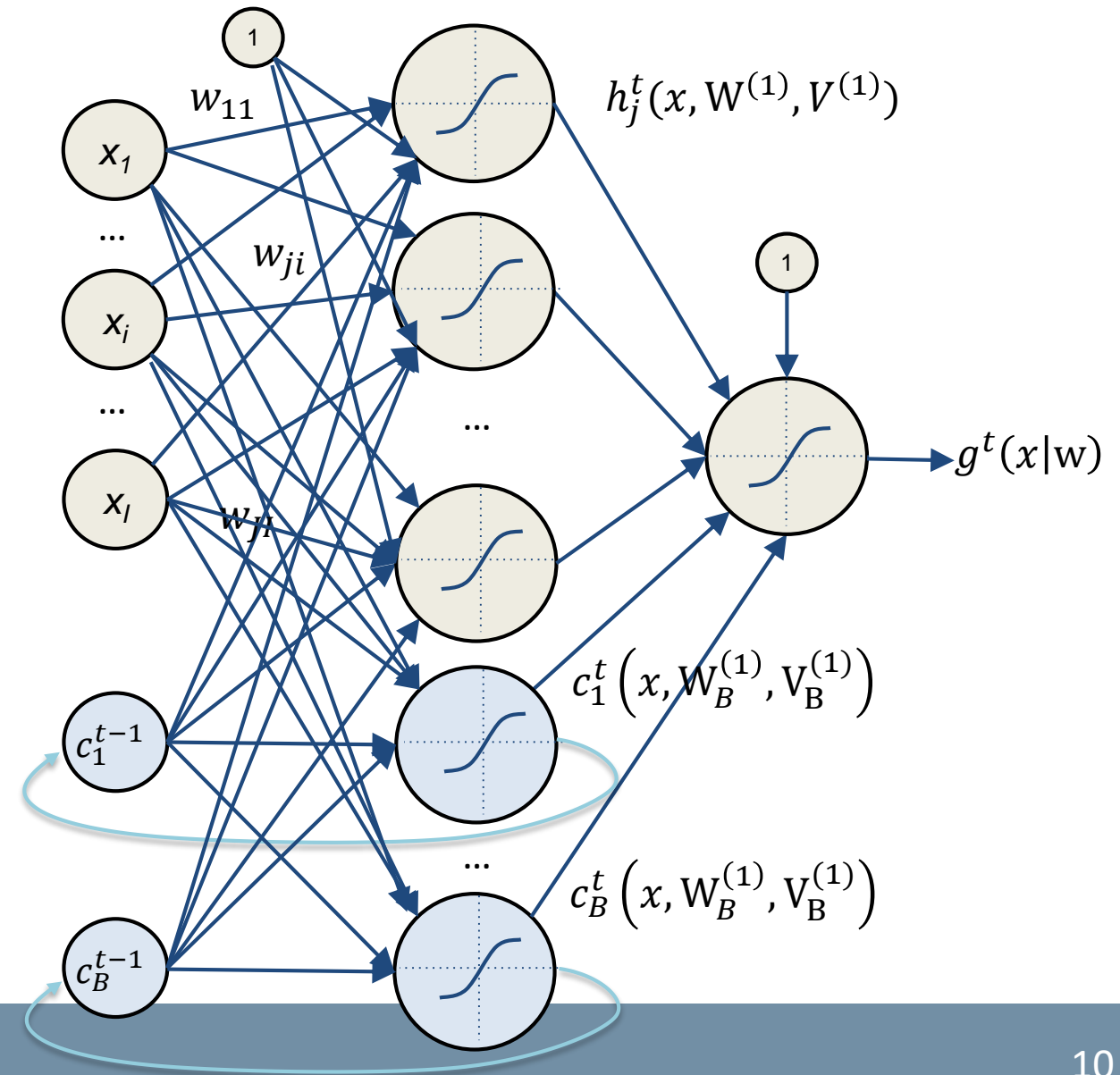
$$g^t(x_n|w) = g \left(\sum_{j=0}^J w_{1j}^{(2)} \cdot h_j^t(\cdot) + \sum_{b=0}^B v_{1b}^{(2)} \cdot c_b^t(\cdot) \right)$$

$$h_j^t(\cdot) = h_j^t \left(\sum_{i=0}^J w_{ji}^{(1)} \cdot x_{i,n} + \sum_{b=0}^B v_{jb}^{(1)} \cdot c_b^{t-1} \right)$$

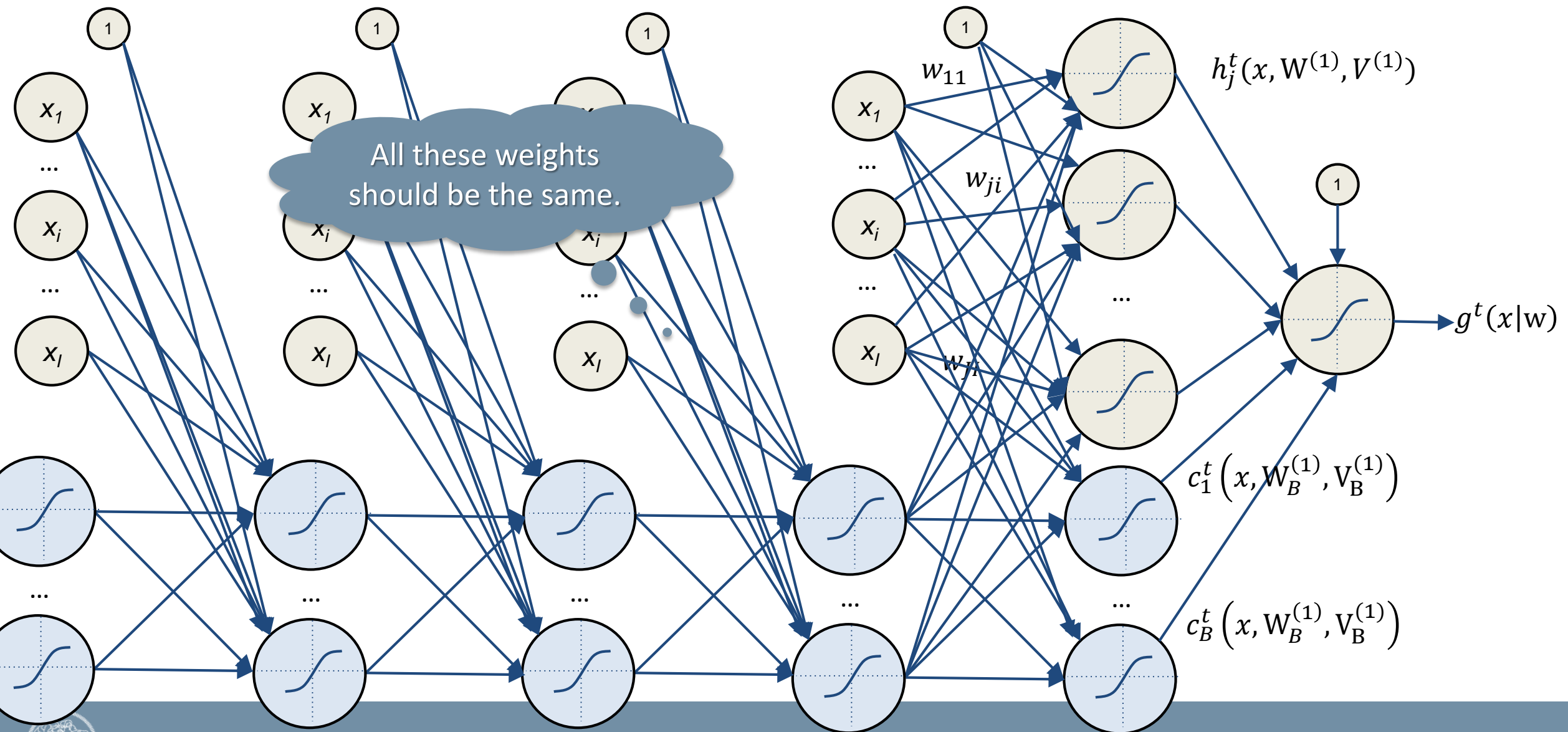
$$c_b^t(\cdot) = c_b^t \left(\sum_{i=0}^J v_{bi}^{(1)} \cdot x_{i,n} + \sum_{b'=0}^B v_{bb'}^{(1)} \cdot c_{b'}^{t-1} \right)$$



Backpropagation Through Time



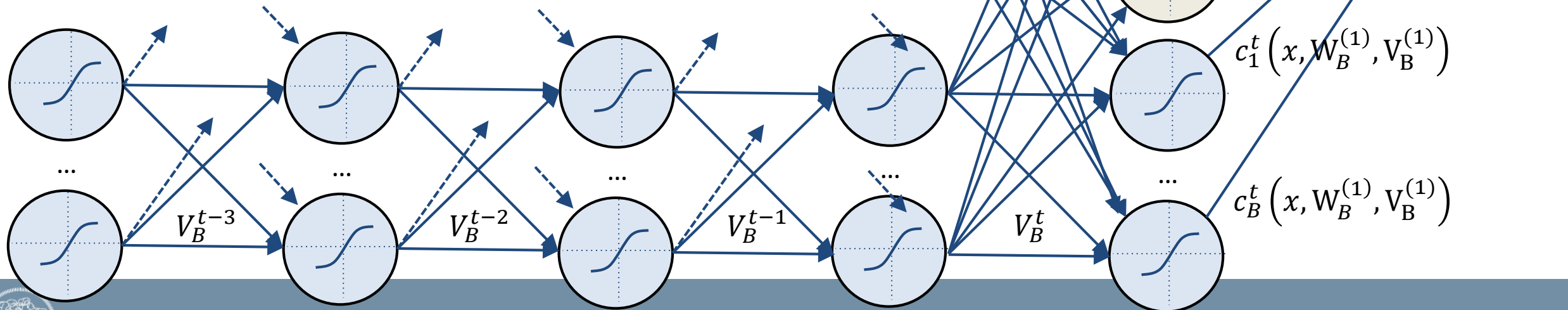
Backpropagation Through Time



Backpropagation Through Time

- Perform network unroll for U steps
- Initialize V, V_B replicas to be the same
- Compute gradients and update replicas with the average of their gradients

$$V = V - \eta \cdot \frac{1}{U} \sum_0^{U-1} V^{t-u} \quad V_B = V_B - \eta \cdot \frac{1}{U} \sum_0^{U-1} V_B^{t-u}$$



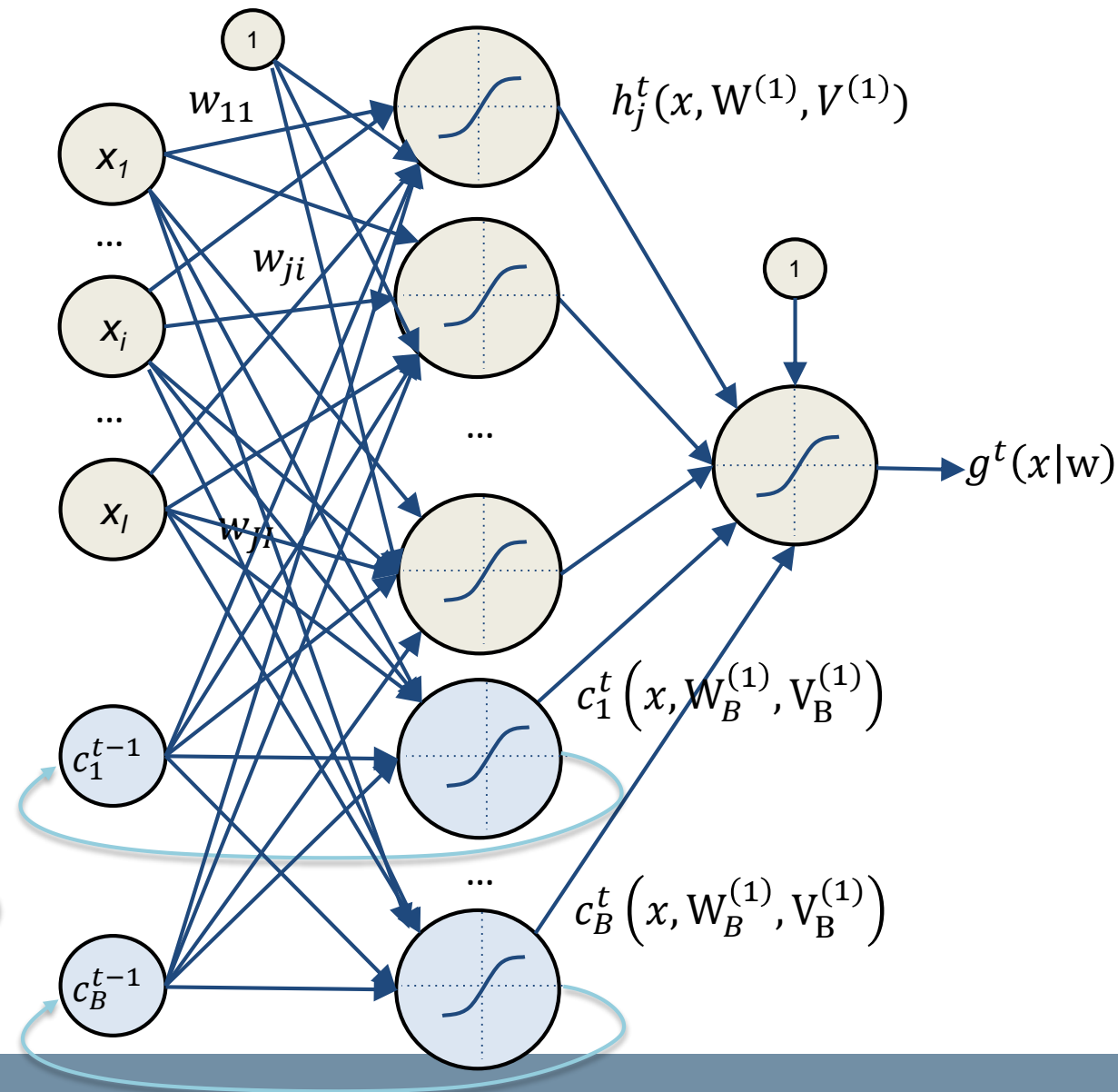
How much should we go back in time?

Sometime the output might be related to some input happened quite long before

Jane walked into the room. John walked in too.
It was late in the day. Jane said hi to <???

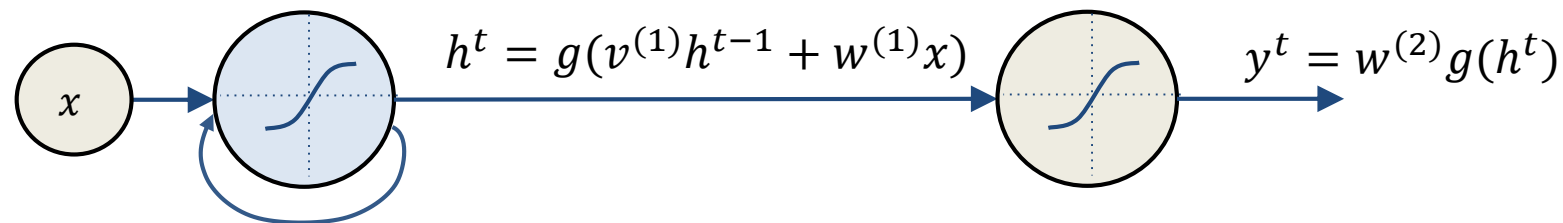
However backpropagation through time was not able to train recurrent neural networks significantly back in time ...

Was due to not being able to backprop through many layers ...



How much can we go back in time?

To better understand why it was not working let consider a simplified case:



Backpropagation over an entire sequence is computed as

$$\frac{\partial E}{\partial w} = \sum_{t=1}^S \frac{\partial E^t}{\partial w} \rightarrow \frac{\partial E^t}{\partial w} = \sum_{t=1}^t \frac{\partial E^t}{\partial y^t} \frac{\partial y^t}{\partial h^t} \frac{\partial h^t}{\partial h^k} \frac{\partial h^k}{\partial w} \rightarrow \frac{\partial h^t}{\partial h^k} = \prod_{i=k+1}^t \frac{\partial h_i}{\partial h_{i-1}} = \prod_{i=k+1}^t v^{(1)} g'(h^{i-1})$$

If we consider the norm of these terms

$$\left\| \frac{\partial h_i}{\partial h_{i-1}} \right\| = \|v^{(1)}\| \|g'(h^{i-1})\| \rightarrow \left\| \frac{\partial h^t}{\partial h^k} \right\| \leq (\gamma_v \cdot \gamma_{g'})^{t-k}$$

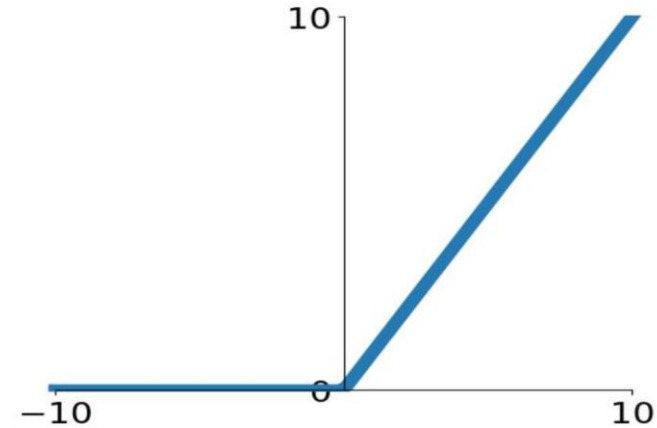
If $(\gamma_v \gamma_{g'}) < 1$ this converges to 0 ...

With Sigmoids and Tanh we have vanishing gradients

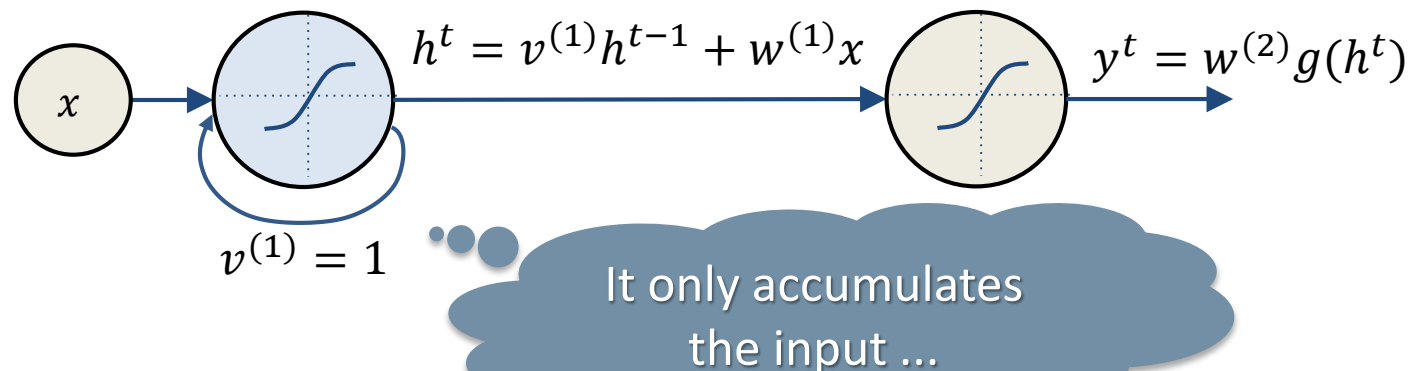
Dealing with Vanishing Gradient

Force all gradients to be either 0 or 1

$$g(a) = \text{ReLu}(a) = \max(0, a)$$
$$g'(a) = 1_{a>0}$$



Build Recurrent Neural Networks using small modules that are designed to remember values for a long time.



Long-Short Term Memories

Hochreiter & Schmidhuber (1997) solved the problem of vanishing gradient designing a memory cell using logistic and linear units with multiplicative interactions:

- Information gets into the cell whenever its “*write*” gate is on.
- The information stays in the cell so long as its “*keep*” gate is on.
- Information is read from the cell by turning on its “*read*” gate.

Can backpropagate through this since the loop has fixed weight.

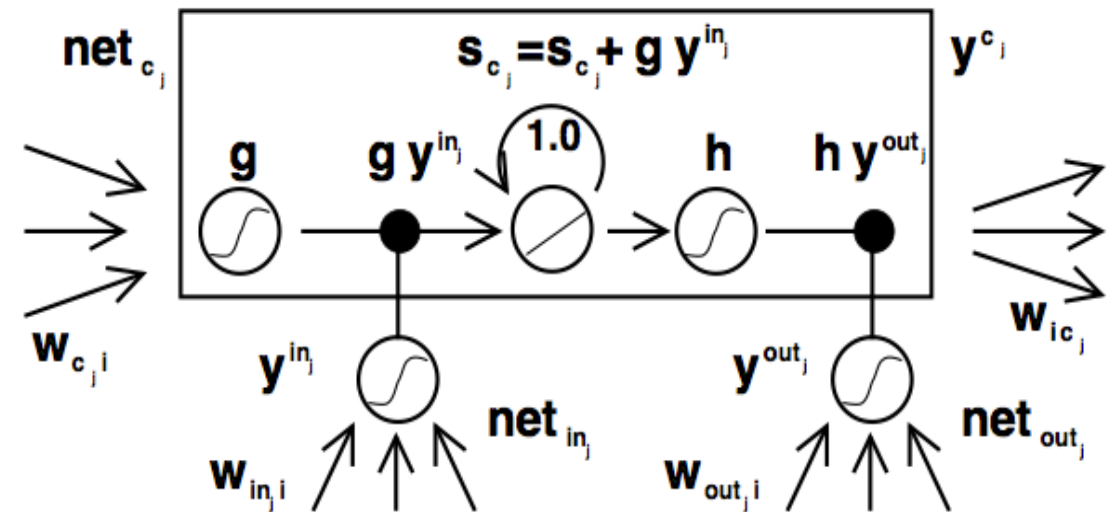
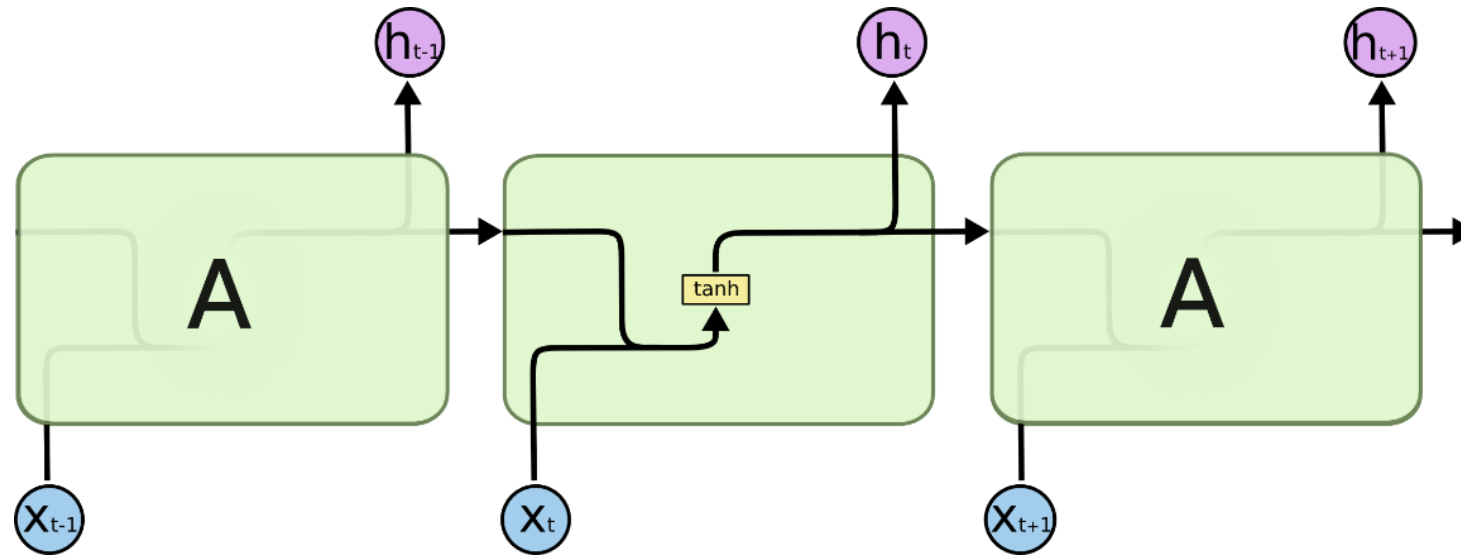


Figure 1: Architecture of memory cell c_j (the box) and its gate units in_j, out_j . The self-recurrent connection (with weight 1.0) indicates feedback with a delay of 1 time step. It builds the basis of the “constant error carousel” CEC. The gate units open and close access to CEC. See text and appendix A.1 for details.

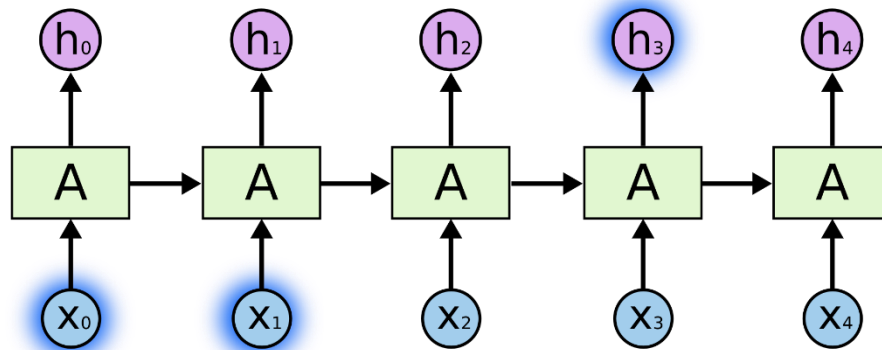
RNN vs. LSTM

RNN



$$a_h^t = \sum_{i=1}^I w_{ih} x_i^t + \sum_{h'=1}^H w_{h'h} b_{h'}^{t-1}$$

$$b_h^t = \theta_h(a_h^t)$$



Neural Network Layer

Pointwise Operation

Vector Transfer

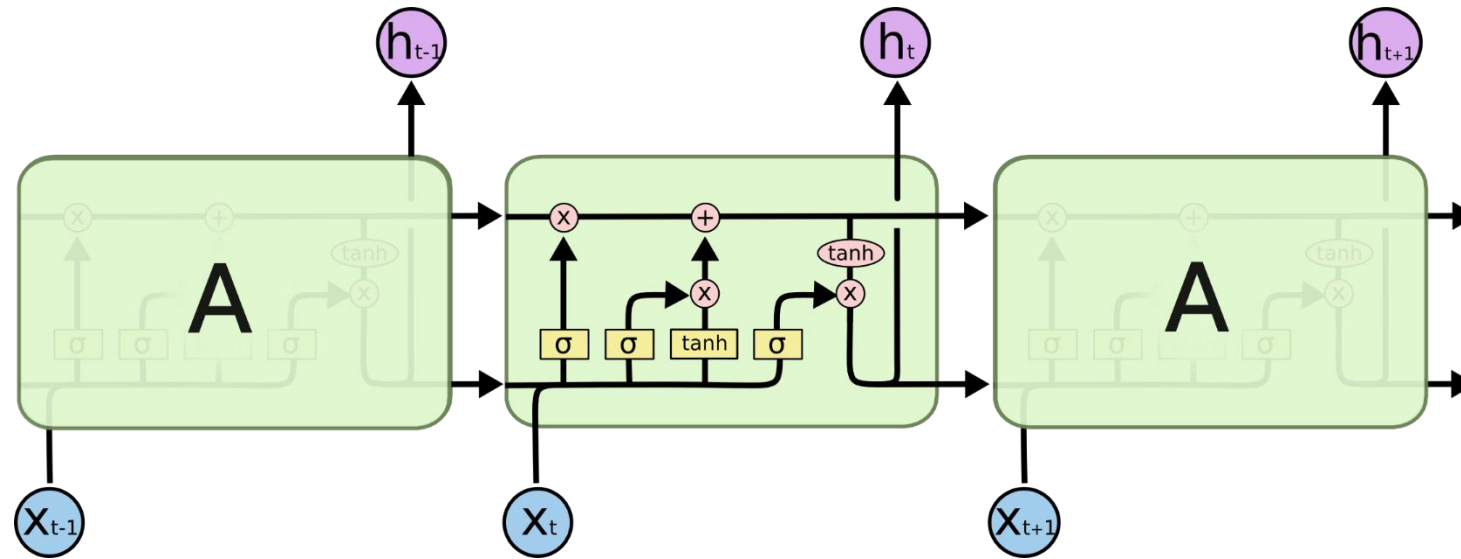
Concatenate

Copy



Long Short Term Memory

LSTM



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

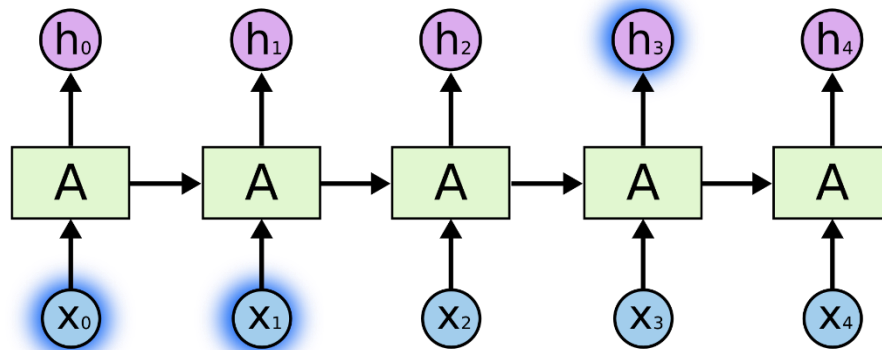
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

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

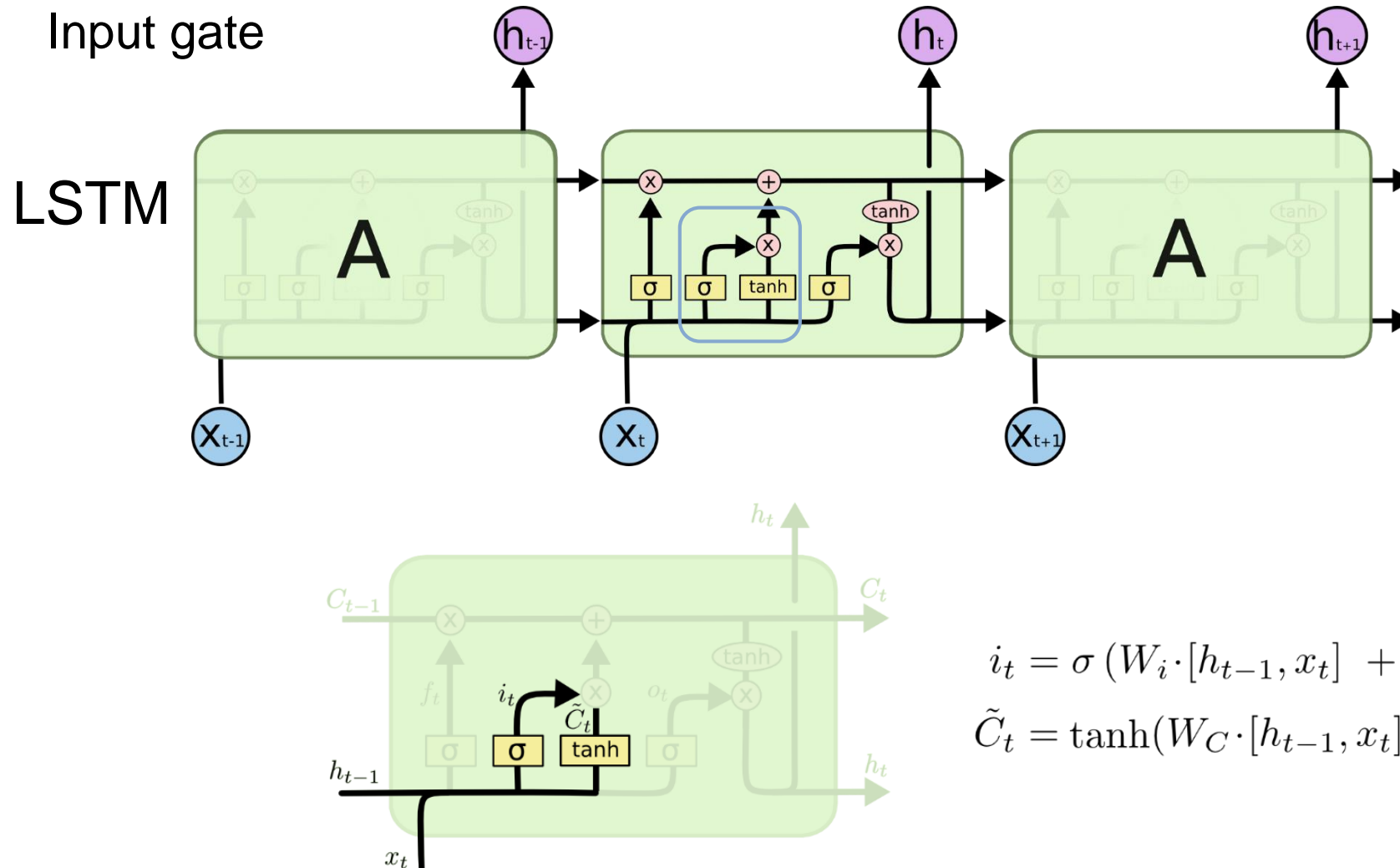
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

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

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



Long Short Term Memory



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

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

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

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

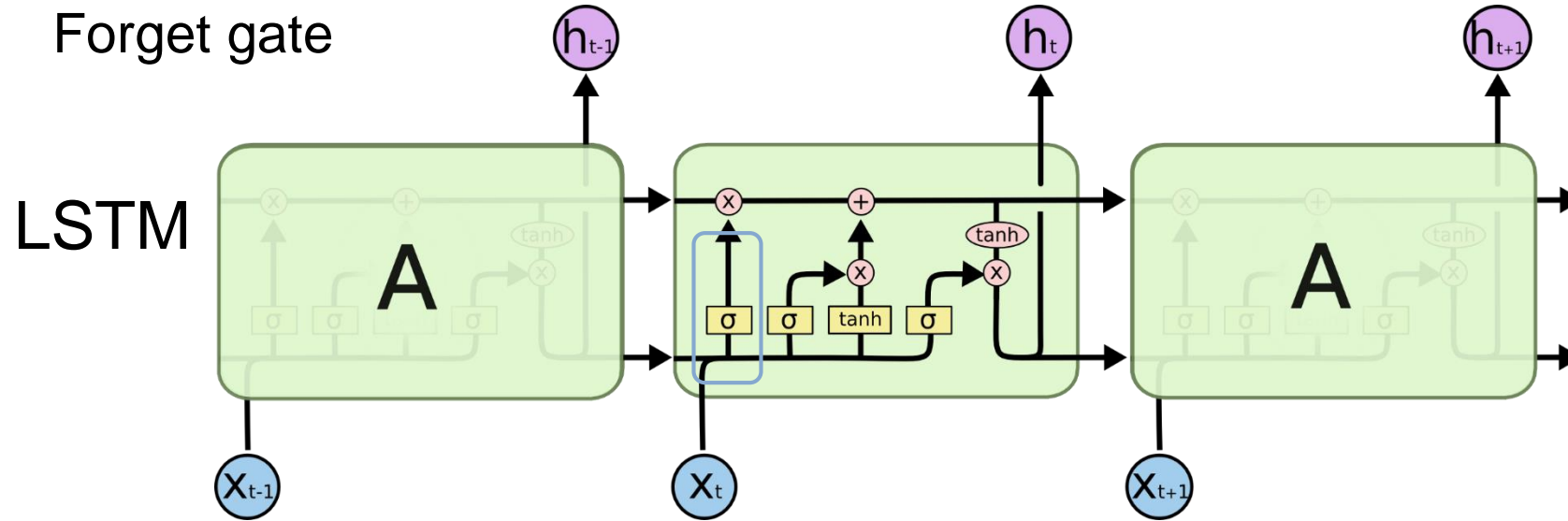
$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

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

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

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

Long Short Term Memory



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

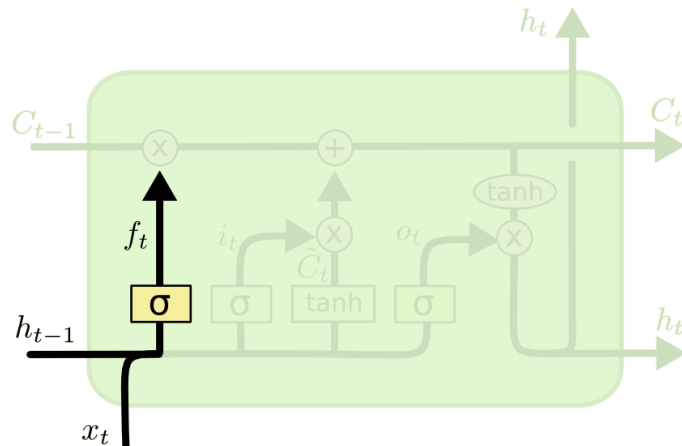
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

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

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

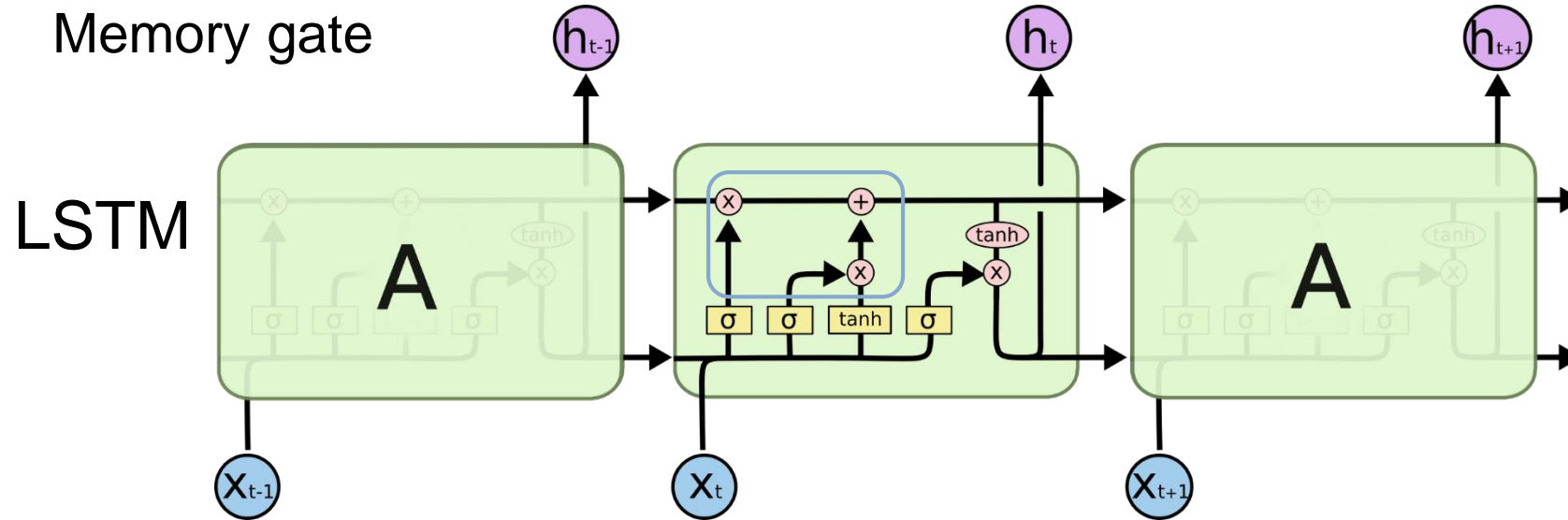
$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

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



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

Long Short Term Memory



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

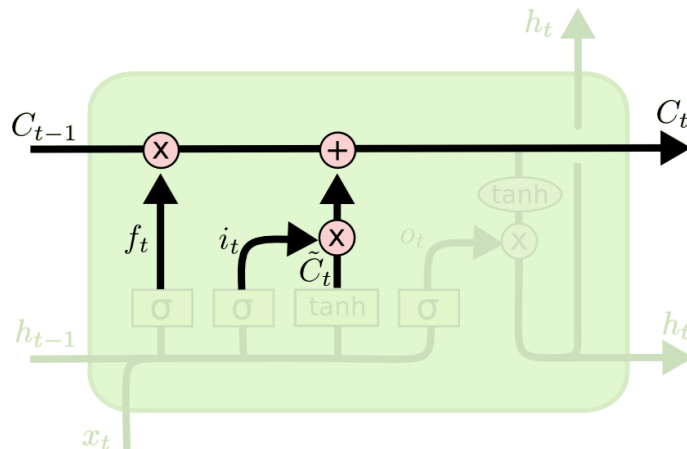
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

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

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

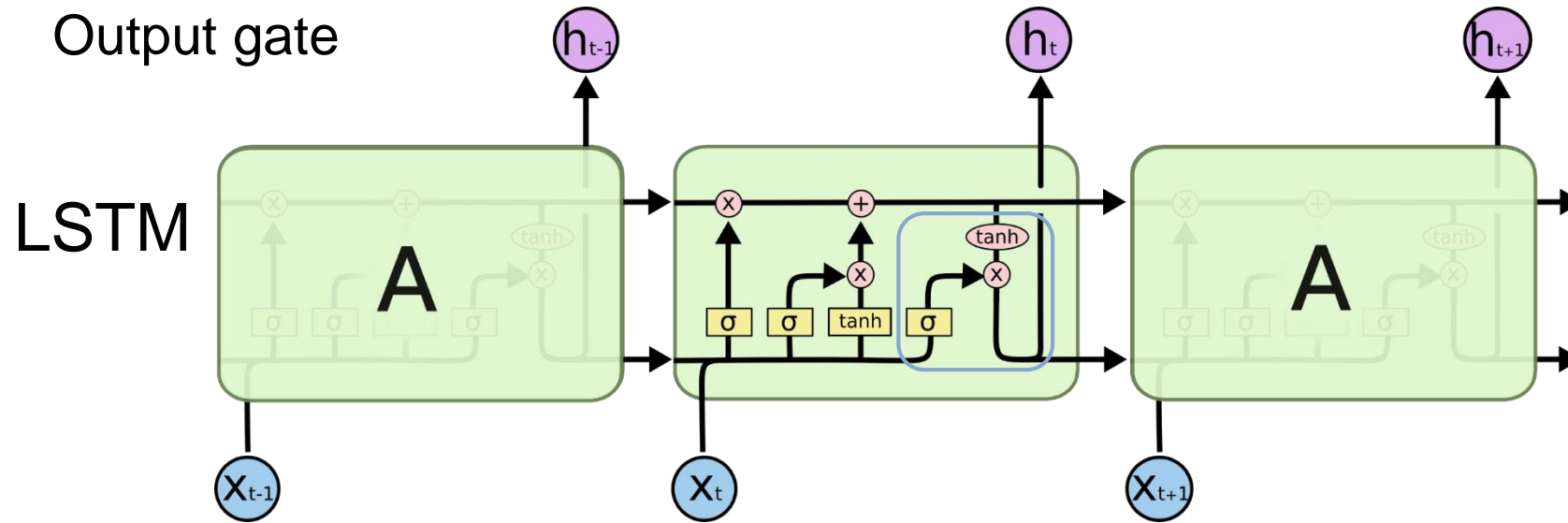
$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

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



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

Long Short Term Memory



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

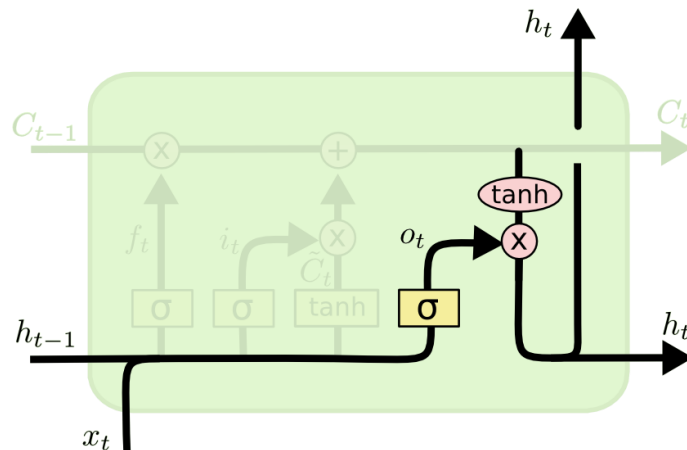
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

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

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

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

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

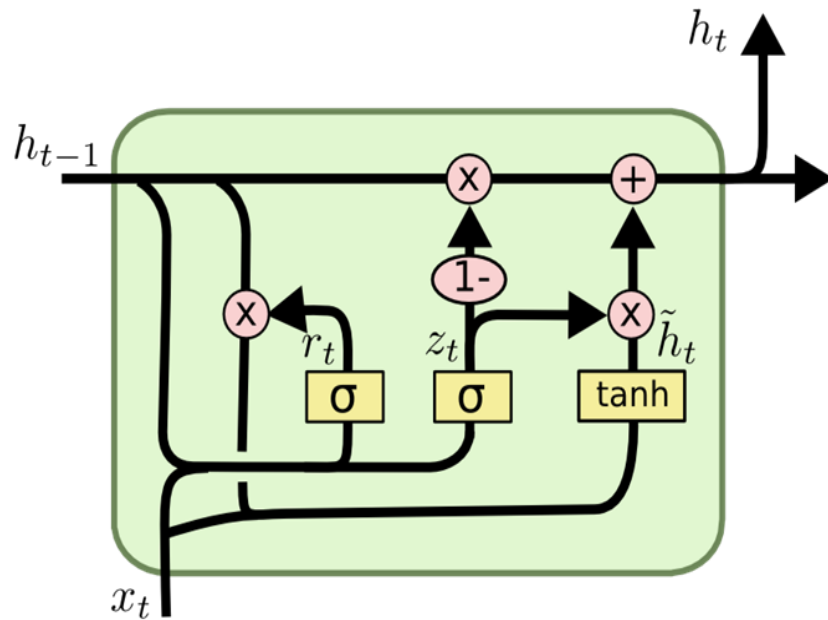


$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

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

Gated Recurrent Unit (GRU)

It combines the forget and input gates into a single “update gate.” It also merges the cell state and hidden state, and makes some other changes.



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

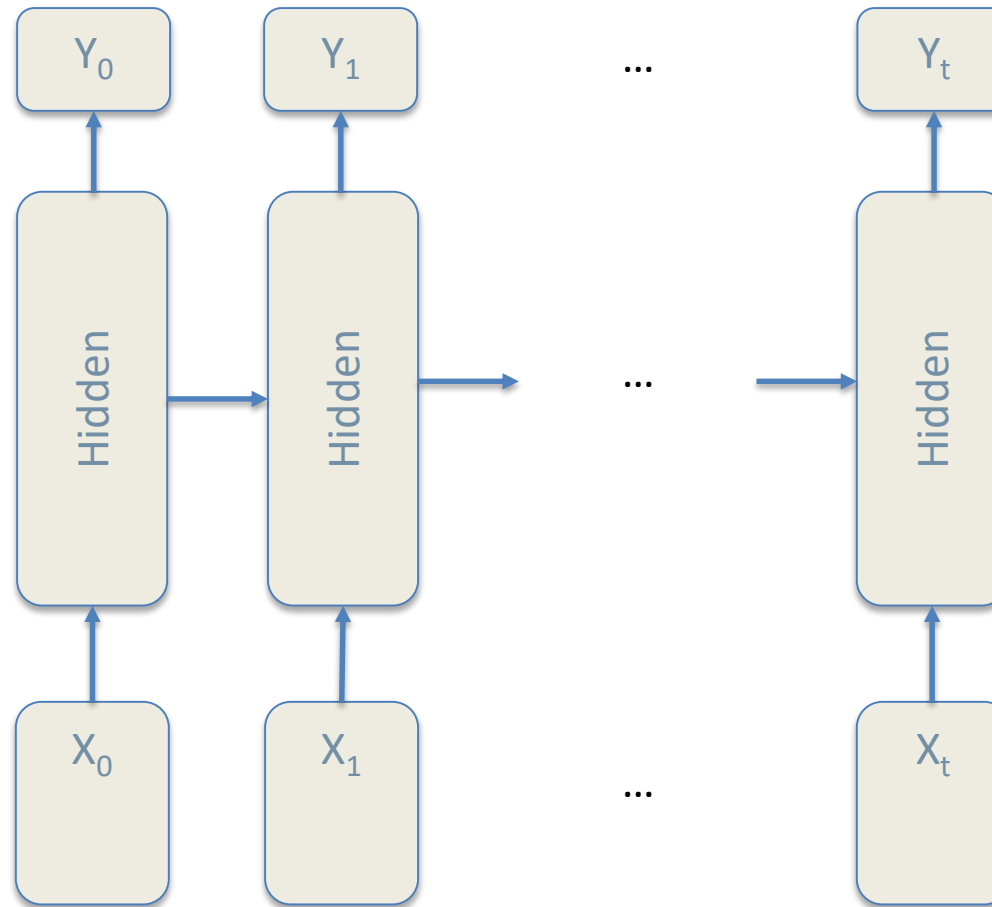
$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

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

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

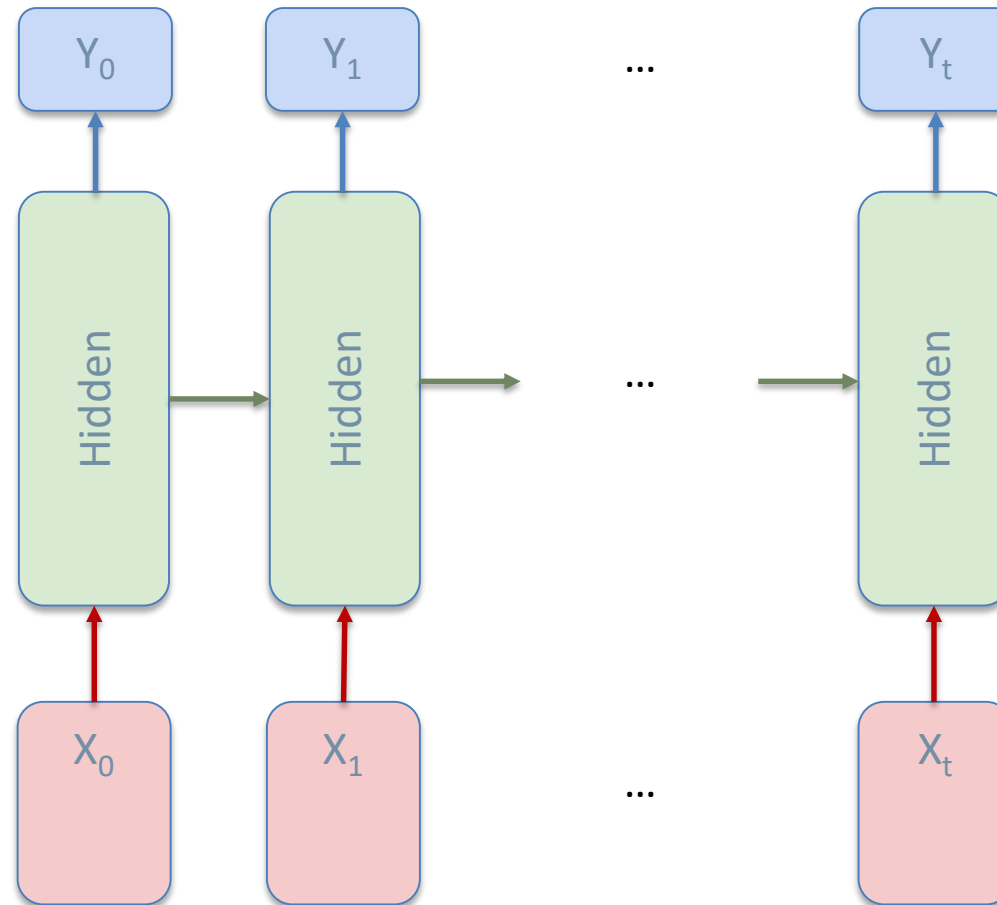
LSTM Networks

You can build a computation graph with continuous transformations.



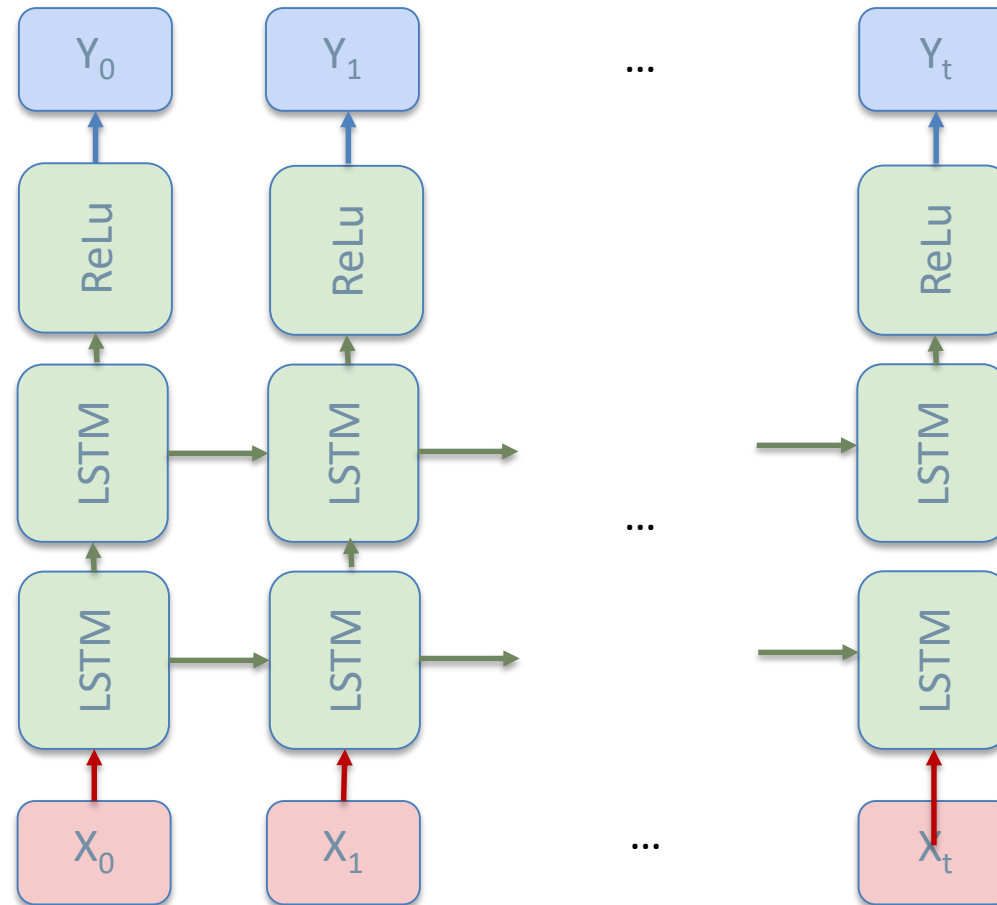
LSTM Networks

You can build a computation graph with continuous transformations.



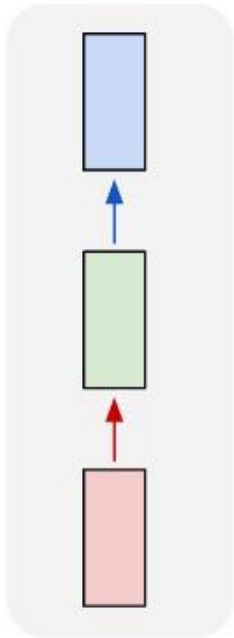
LSTM Networks

You can build a computation graph with continuous transformations.



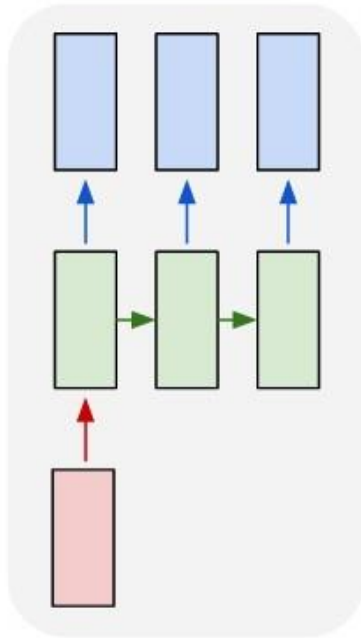
Sequential Data Problems

one to one



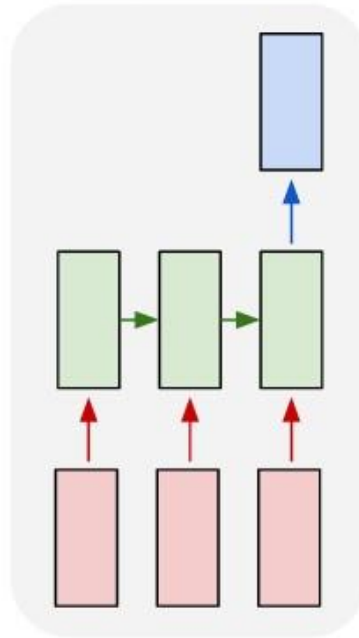
Fixed-sized input to fixed-sized output (e.g. image classification)

one to many



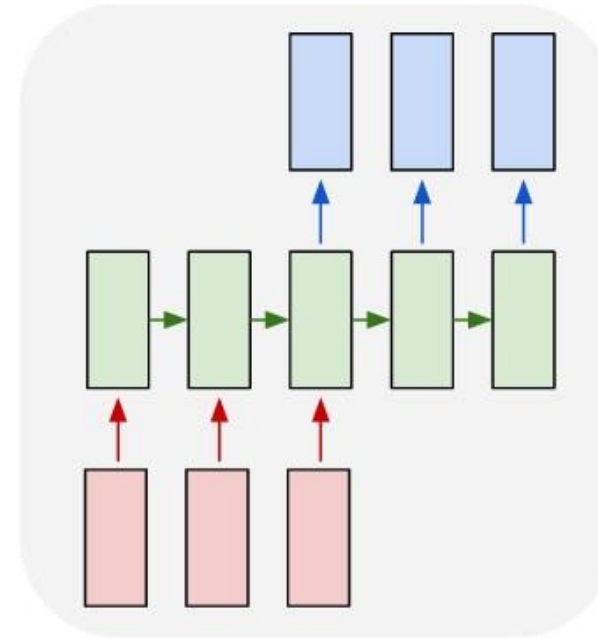
Sequence output (e.g. image captioning takes an image and outputs a sentence of words).

many to one



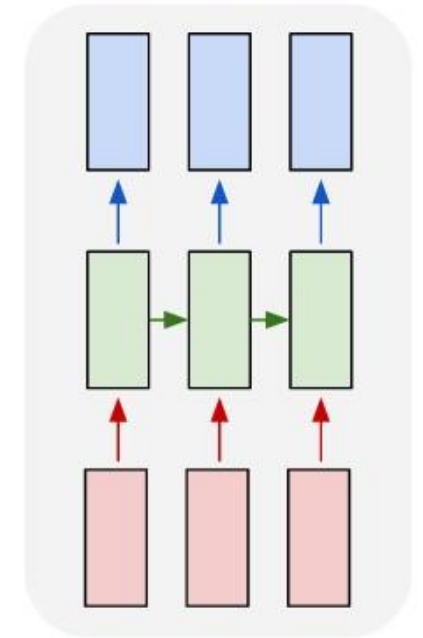
Sequence input (e.g. sentiment analysis where a given sentence is classified as expressing positive or negative sentiment).

many to many



Sequence input and sequence output (e.g. Machine Translation: an RNN reads a sentence in English and then outputs a sentence in French)

many to many



Synced sequence input and output (e.g. video classification where we wish to label each frame of the video)

Credits: Andrej Karpathy

Acknowledgements

This slides are highly based on material taken from:

- Geoffrey Hinton
- Hugo Larochelle
- Andrej Karpathy
- Nando De Freitas
- Chris Olah

You can find more details on the original slides

The amazing images on LSTM cell are taken from Chris Holo's blog:

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

