



# DD2437 – Artificial Neural Networks and Deep Architectures (annda)

## Lecture 7: Temporal processing with ANNs: feedforward vs recurrent networks

Pawel Herman

Computational Science and Technology (CST)  
KTH Royal Institute of Technology

February 2018

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- ESN and LSTM

## Lecture overview

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time (BPTT)
- Echo state networks (ESNs)
- Long short-term memory model (LSTM)

- **Temporal processing with feedforward NNs**
- Recurrent architectures for sequence modelling
- Backpropagation through time
- ESN and LSTM

# Temporal aspects

- Time is an essential component of the description of many phenomena, observations, data structures
- Omnipotence of sequences – ordering of entities
  - numerical codes
  - language and speech
  - motor behaviour
  - signals, time series: sensor readings, market prices, biological recordings etc.
- Discrete vs continuous time
- Implicit vs explicit representation

- **Temporal processing with feedforward NNs**
- Recurrent architectures for sequence modelling
- Backpropagation through time
- ESN and LSTM

# Static MLP for handling dynamics

The use of a static MLP to account for temporal dimension

- short-term memory function
- nonlinear regression capabilities

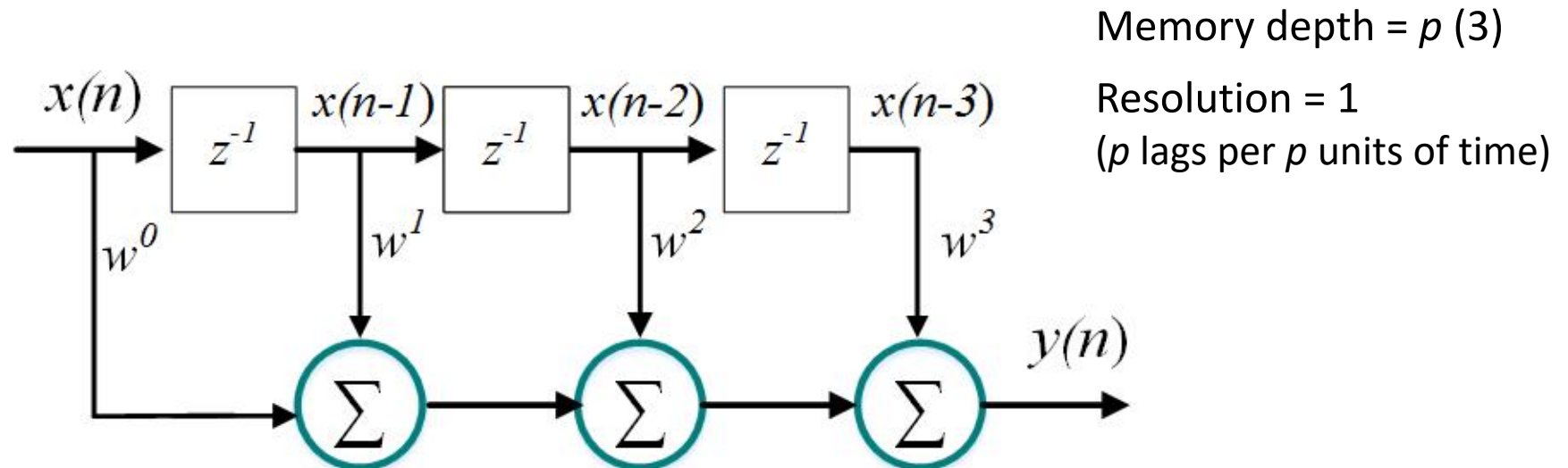
- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- ESN and LSTM

# Static MLP for handling dynamics

The use of a static MLP to account for temporal dimension

- short-term memory function
- nonlinear regression capabilities

Tapped delay line memory



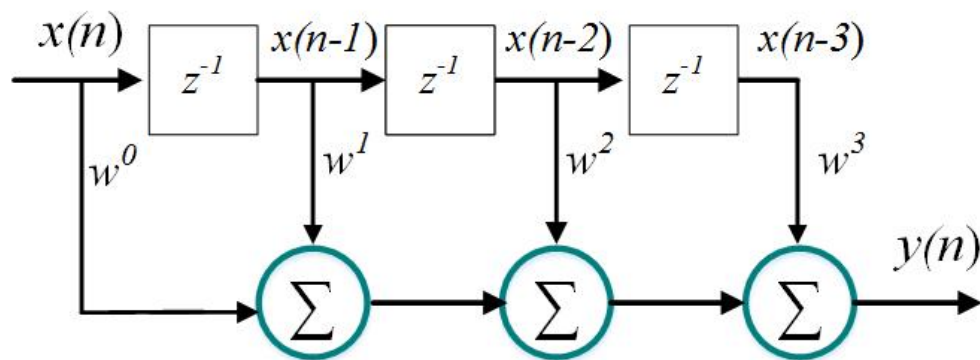
- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- ESN and LSTM

# Static MLP for handling dynamics

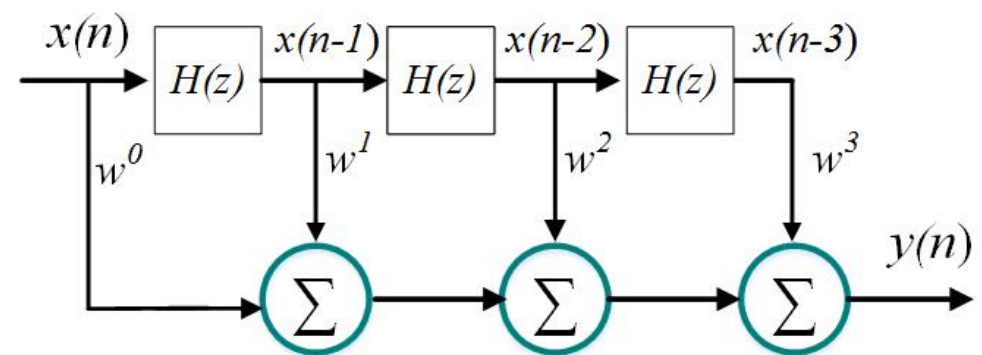
The use of a static MLP to account for temporal dimension

- short-term memory function
- nonlinear regression capabilities

Tapped delay line memory

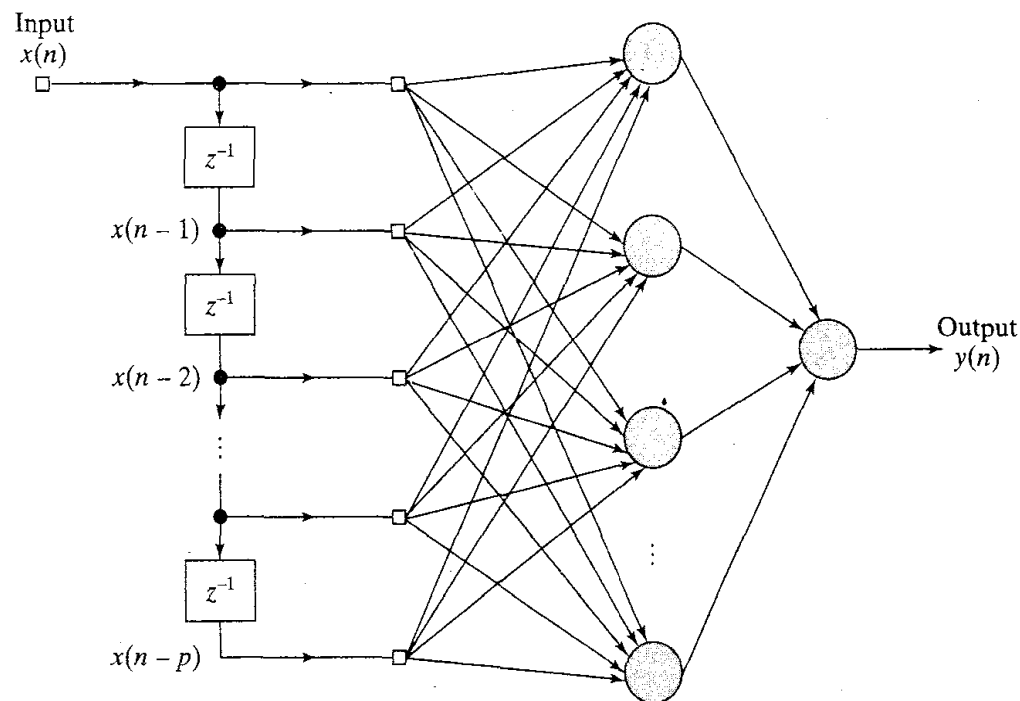


Generalized tapped delay line memory



- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- ESN and LSTM

## Learning approach to TLFN



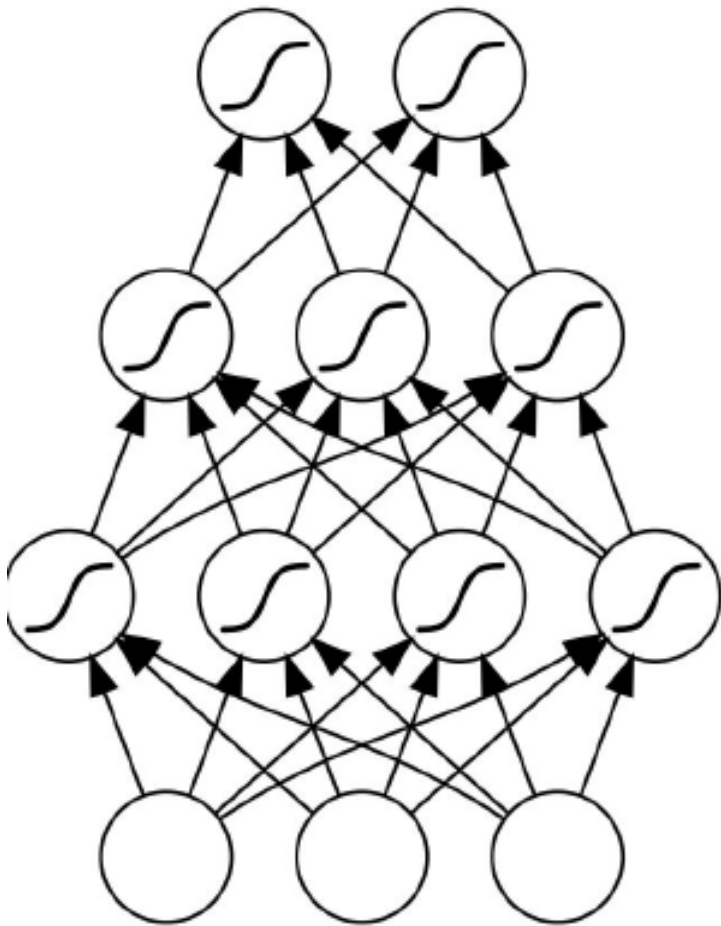
Backprop can be used with relatively simple *focused TLFNs*.

A general principle to unfold the network: form a large “static” network, and apply backprop.

Haykin, 1999

- Temporal processing with feedforward NNs
- **Recurrent architectures for sequence modelling**
- Backpropagation through time
- ESN and LSTM

# Recurrent neural networks (RNNs)

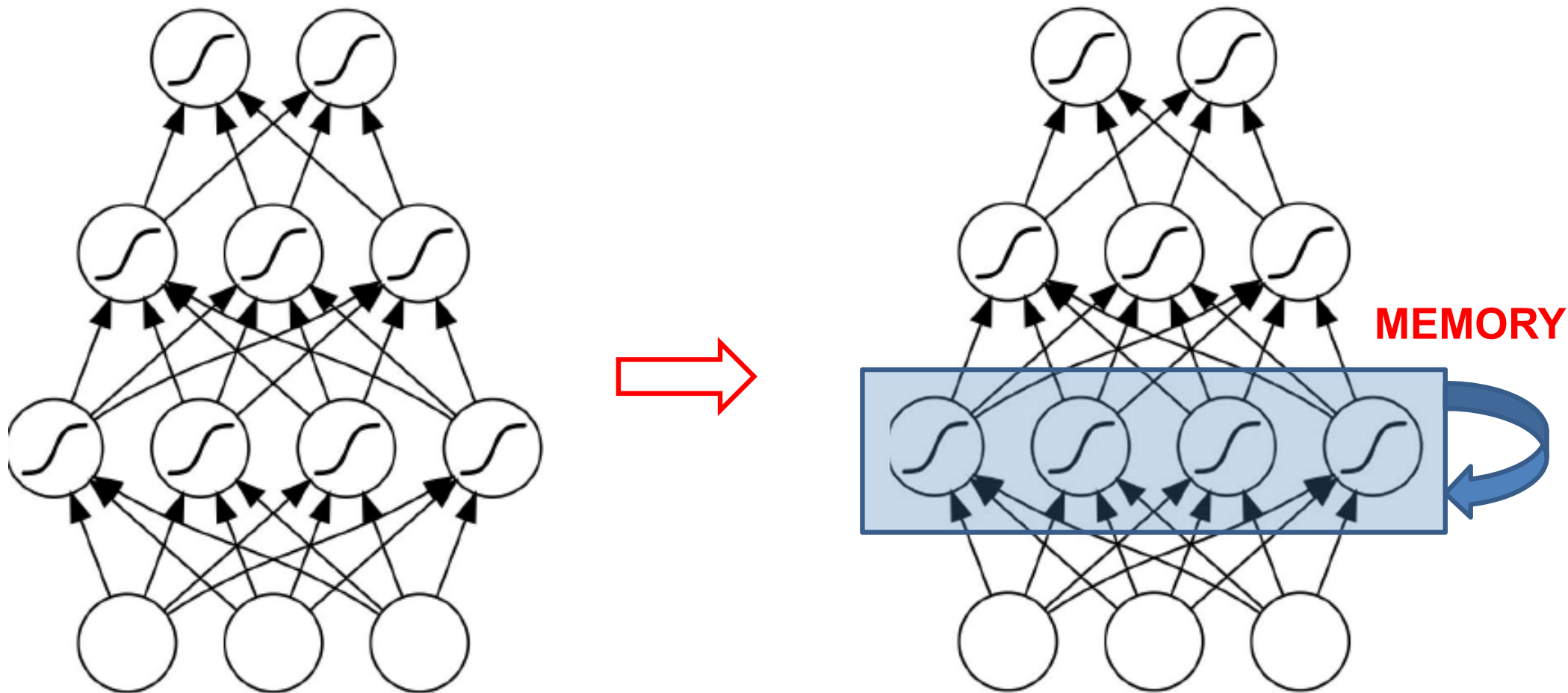


From MLP to ....



- Temporal processing with feedforward NNs
- **Recurrent architectures for sequence modelling**
- Backpropagation through time
- ESN and LSTM

# Recurrent neural networks (RNNs)

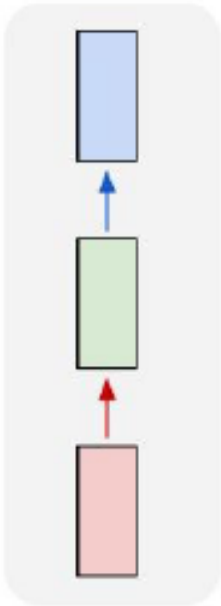


From MLP to RNN

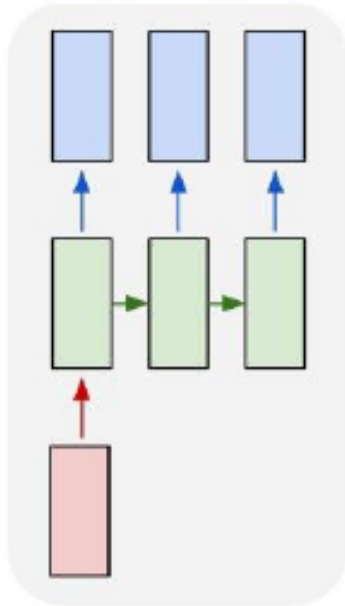
- Temporal processing with feedforward NNs
- **Recurrent architectures for sequence modelling**
- Backpropagation through time
- ESN and LSTM

# Repertoire of recurrent architectures

one to one

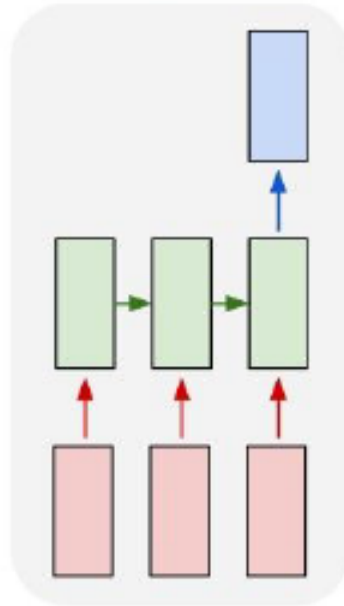


one to many



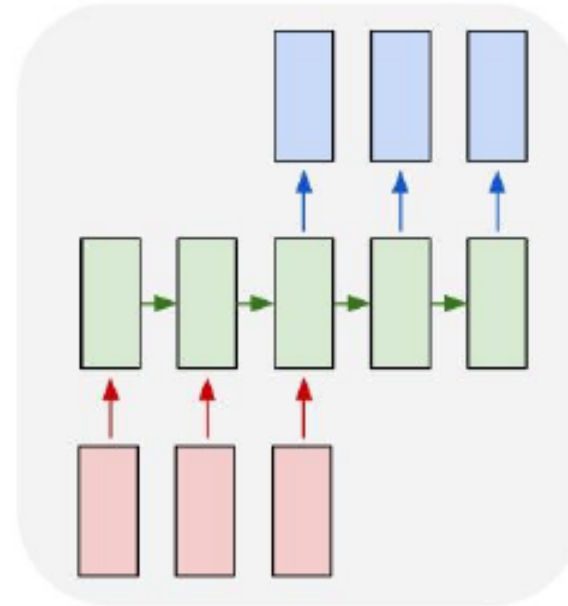
e.g. sequence generation

many to one



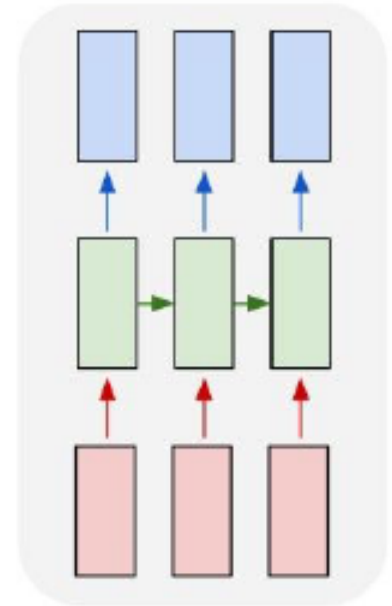
e.g. sentiment analysis

many to many



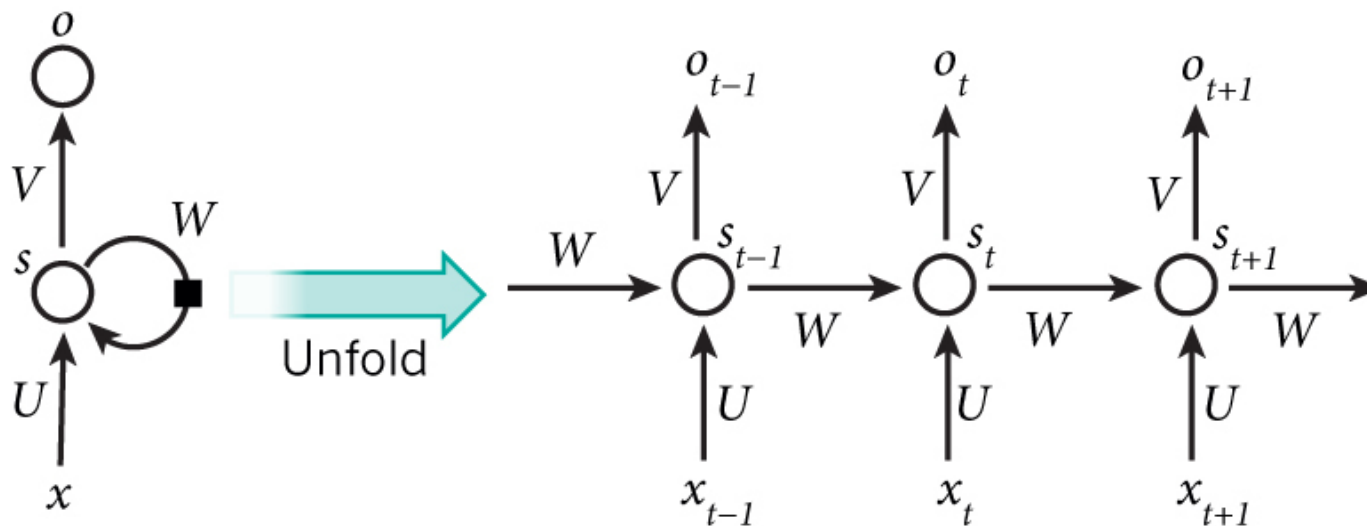
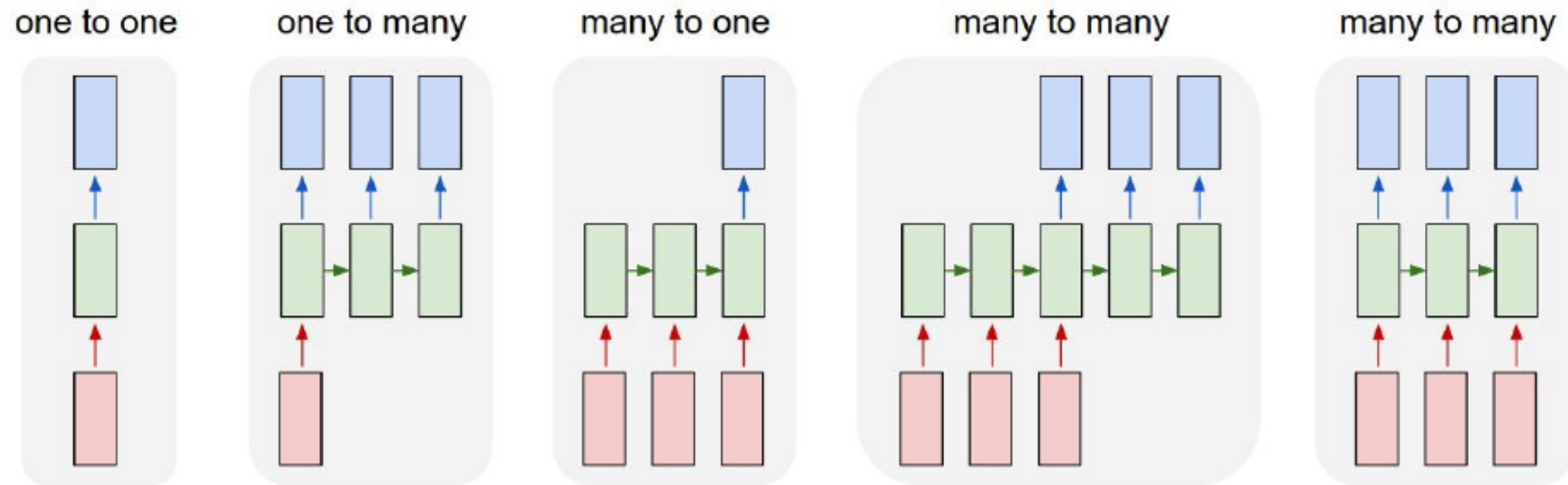
e.g. machine translation

many to many



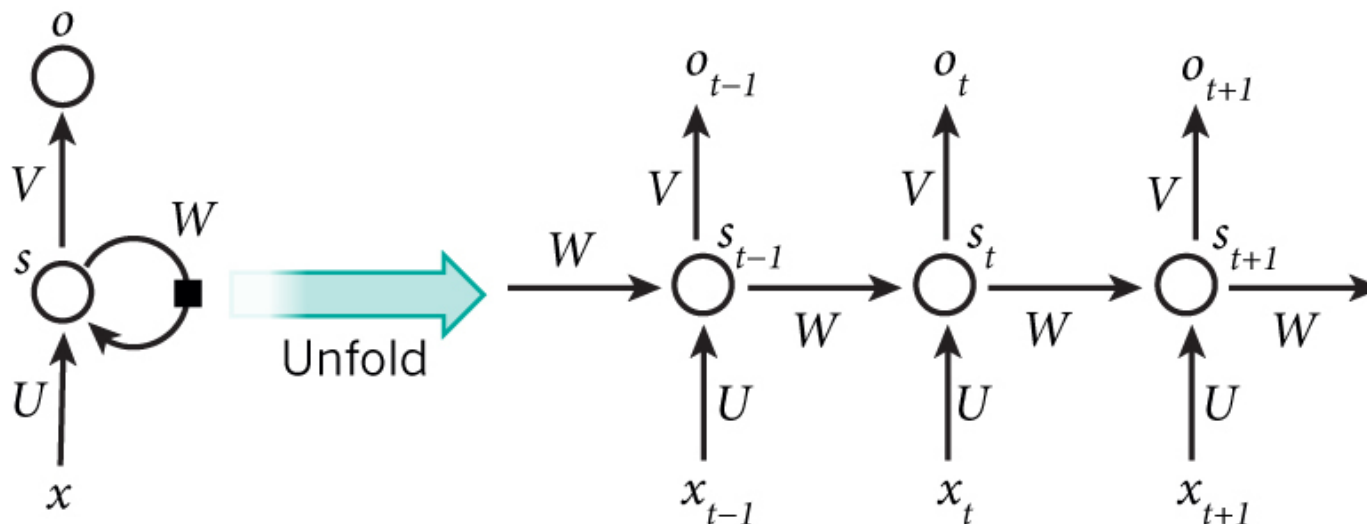
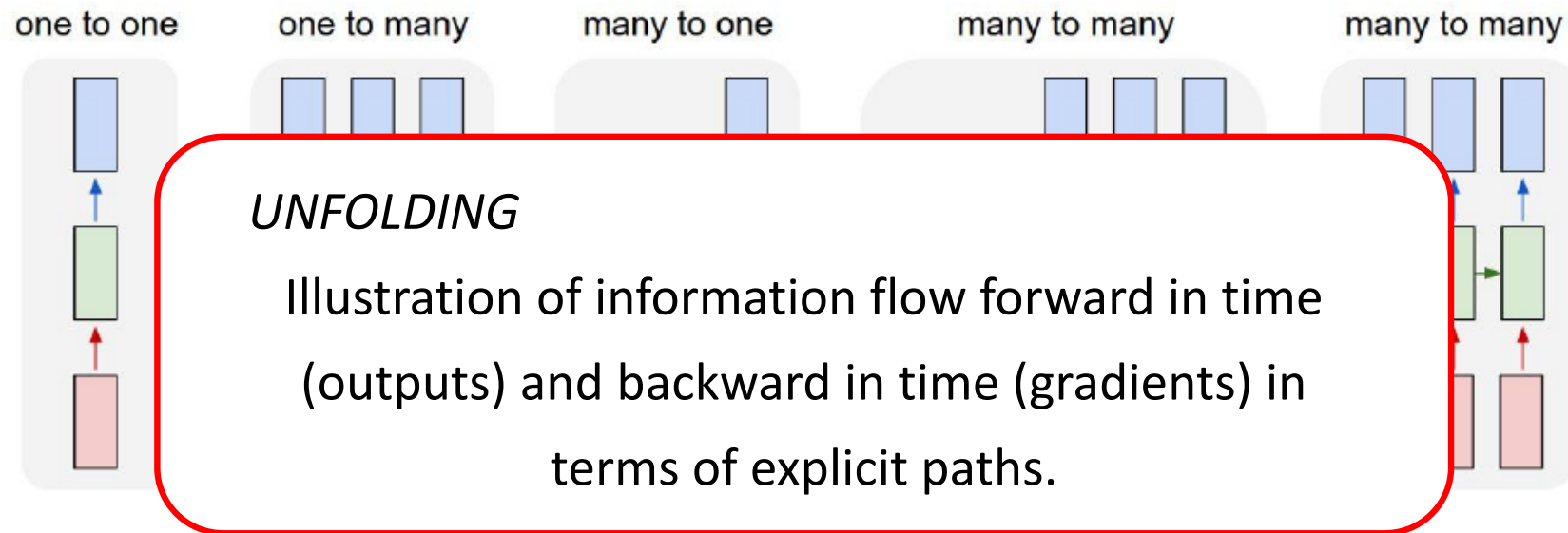
- Temporal processing with feedforward NNs
- **Recurrent architectures for sequence modelling**
- Backpropagation through time
- ESN and LSTM

# Repertoire of recurrent architectures, unfolding



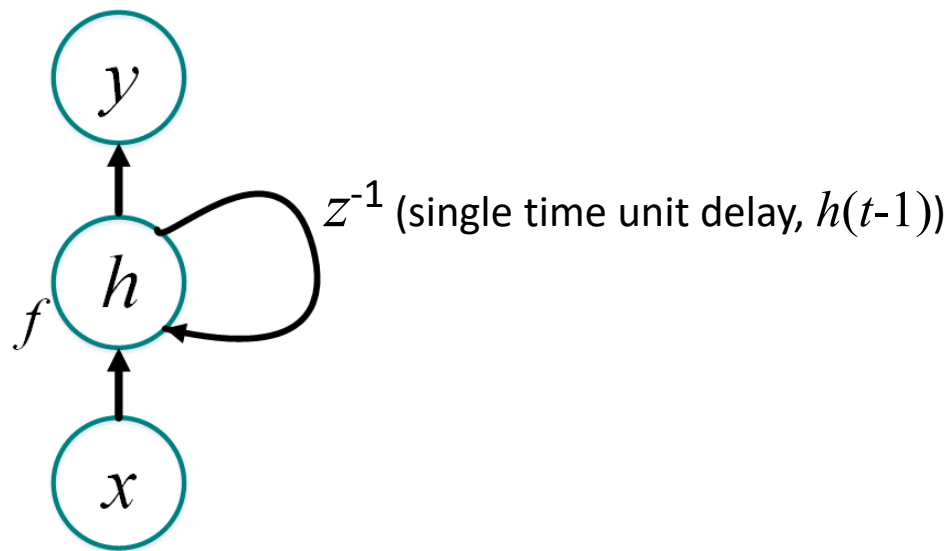
- Temporal processing with feedforward NNs
- **Recurrent architectures for sequence modelling**
- Backpropagation through time
- ESN and LSTM

# Repertoire of recurrent architectures, unfolding



- Temporal processing with feedforward NNs
- **Recurrent architectures for sequence modelling**
- Backpropagation through time
- ESN and LSTM

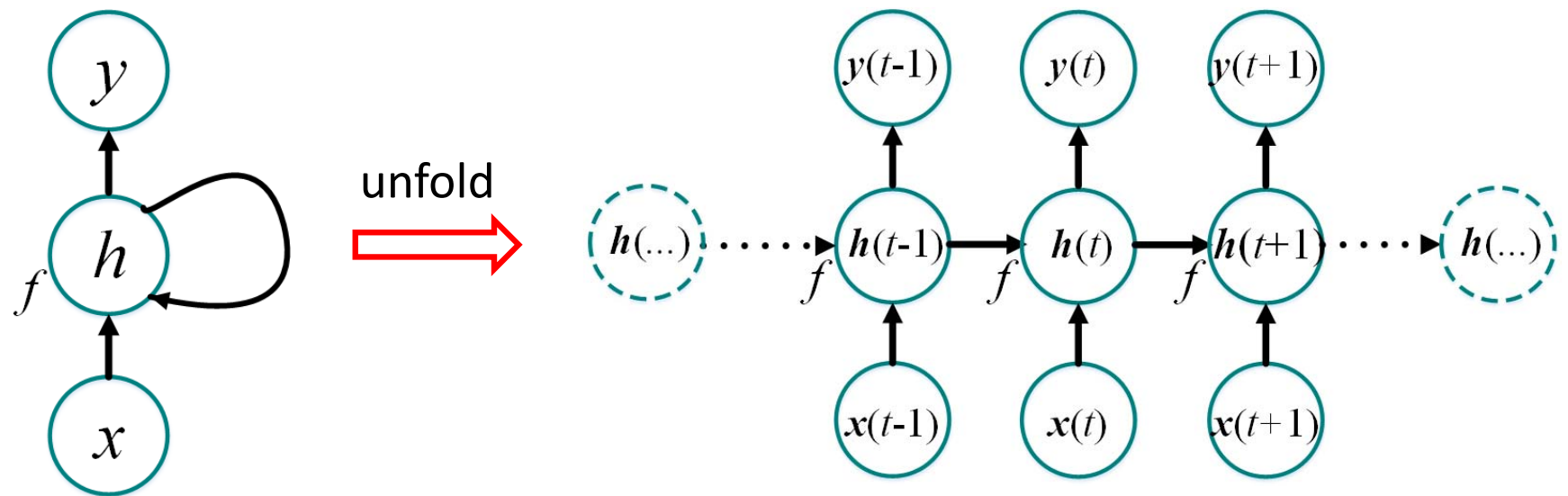
# Fundamental, vanilla RNN unit



$$\mathbf{h}(t) = f(\mathbf{h}(t-1), \mathbf{x}(t), \boldsymbol{\theta})$$

- Temporal processing with feedforward NNs
- **Recurrent architectures for sequence modelling**
- Backpropagation through time
- ESN and LSTM

# Fundamental, vanilla RNN – unfolded

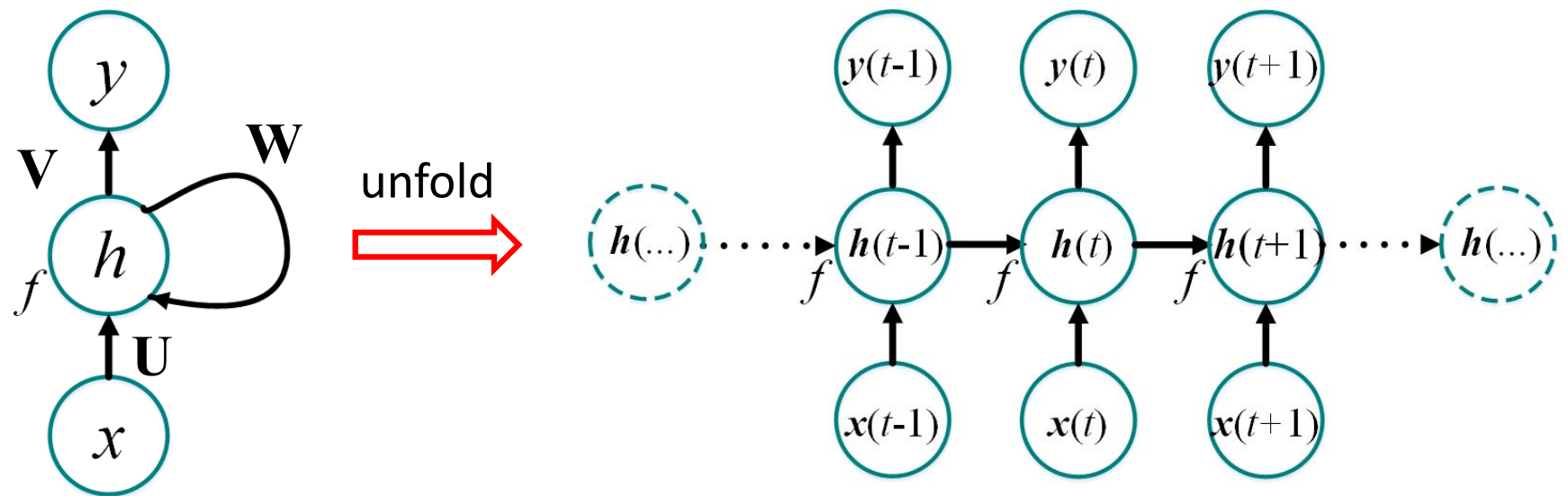


$$h(t) = f(h(t-1), x(t), \theta)$$

In a canonical form it allows for modelling sequences of varying length (though there are problems of technical nature).

- Temporal processing with feedforward NNs
- **Recurrent architectures for sequence modelling**
- Backpropagation through time
- ESN and LSTM

# Recurrent connections between hidden units



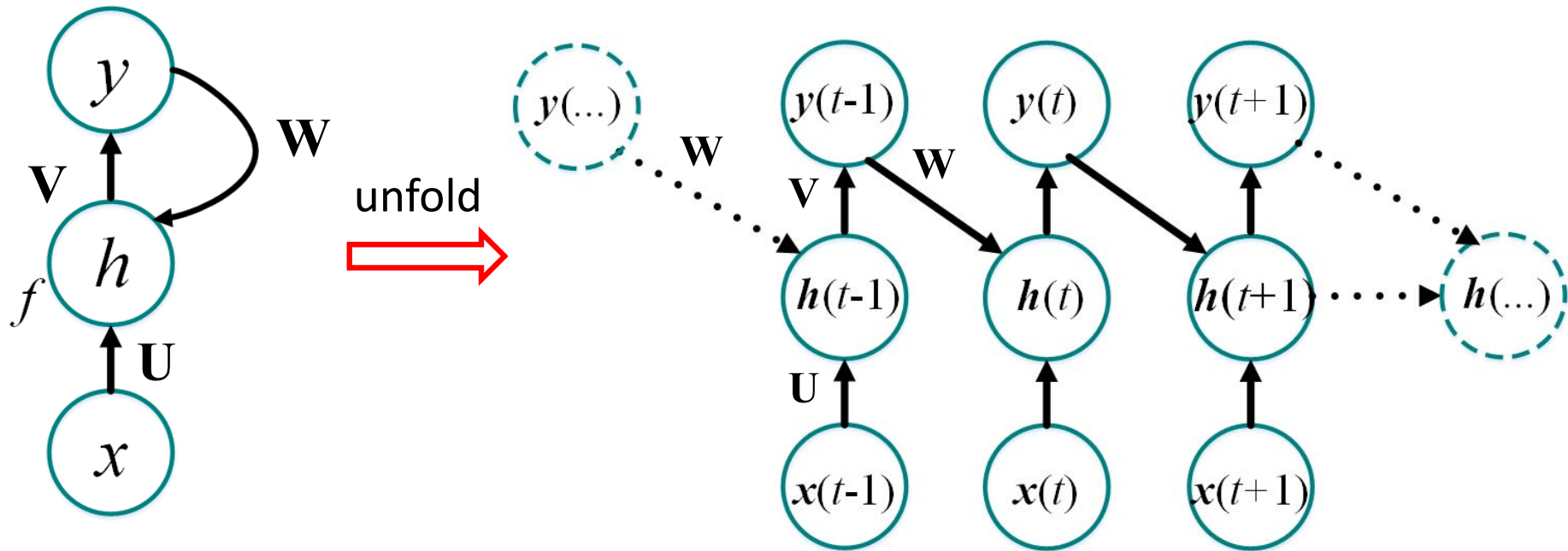
$$h(t) = f(\mathbf{W}h(t-1) + \mathbf{U}x(t) + bias)$$

$$y(t) = \mathbf{V}h(t) + bias$$

state-space  
description

- Temporal processing with feedforward NNs
- **Recurrent architectures for sequence modelling**
- Backpropagation through time
- ESN and LSTM

# Recurrent connection from output to hidden units



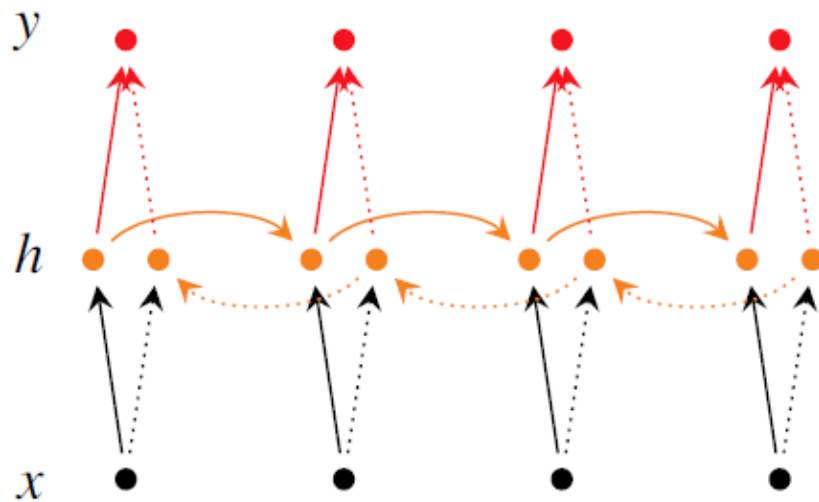
$$\mathbf{h}(t) = f(\mathbf{W}\mathbf{y}(t-1) + \mathbf{U}\mathbf{x}(t) + bias)$$

$$\mathbf{y}(t) = \mathbf{V}\mathbf{h}(t) + bias$$



- Temporal processing with feedforward NNs
- **Recurrent architectures for sequence modelling**
- Backpropagation through time
- ESN and LSTM

## .... bidirectional neural networks



$$\vec{h}_t = f(\vec{W}x_t + \vec{V}\vec{h}_{t-1} + \vec{b})$$

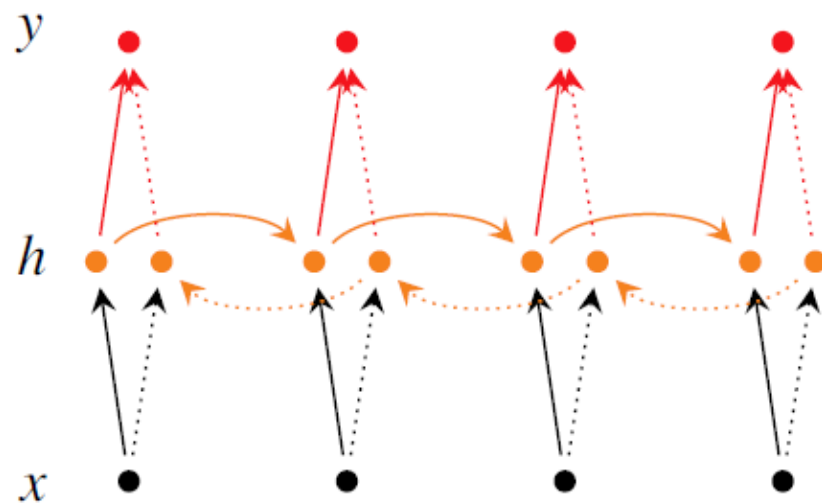
$$\overleftarrow{h}_t = f(\overleftarrow{W}x_t + \overleftarrow{V}\overleftarrow{h}_{t+1} + \overleftarrow{b})$$

$$y_t = g(U[\vec{h}_t; \overleftarrow{h}_t] + c)$$

To incorporate information from words or phonemes both preceding and following, e.g. where there is clear dependency of phonemes/words on the following neighbouring phonemes/words

- Temporal processing with feedforward NNs
- **Recurrent architectures for sequence modelling**
- Backpropagation through time
- ESN and LSTM

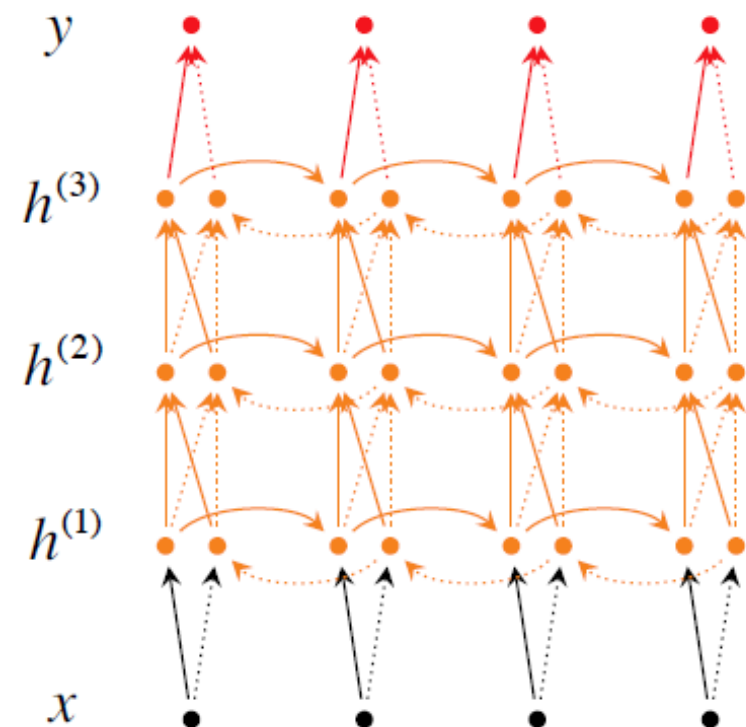
# Shallow vs deep bidirectional neural networks



$$\vec{h}_t = f(\vec{W}x_t + \vec{V}\vec{h}_{t-1} + \vec{b})$$

$$\overleftarrow{h}_t = f(\overleftarrow{W}x_t + \overleftarrow{V}\overleftarrow{h}_{t+1} + \overleftarrow{b})$$

$$y_t = g(U[\vec{h}_t; \overleftarrow{h}_t] + c)$$



$$\vec{h}_t^{(i)} = f(\vec{W}^{(i)} \vec{h}_t^{(i-1)} + \vec{V}^{(i)} \vec{h}_{t-1}^{(i)} + \vec{b}^{(i)})$$

$$\overleftarrow{h}_t^{(i)} = f(\overleftarrow{W}^{(i)} \overleftarrow{h}_t^{(i-1)} + \overleftarrow{V}^{(i)} \overleftarrow{h}_{t+1}^{(i)} + \overleftarrow{b}^{(i)})$$

$$y_t = g(U[\vec{h}_t^{(L)}; \overleftarrow{h}_t^{(L)}] + c)$$

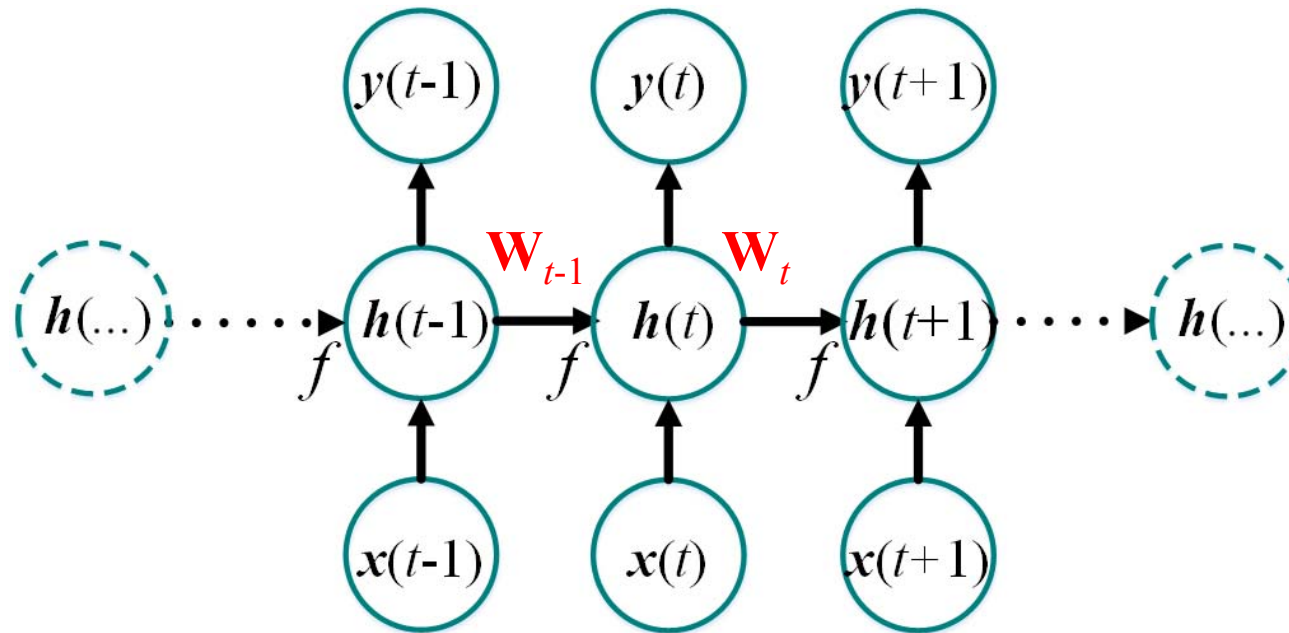
- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- **Backpropagation through time**
- ESN and LSTM

# Learning algorithms for RNNs

- Epochwise vs continuous training
  - “epoch” corresponds to a data sample – a sequence
  - RNN activity reset between epochs (in epochwise training)
- Backpropagation through time (unfolding into an MLP) can be both applied epochwise and in a continuous fashion

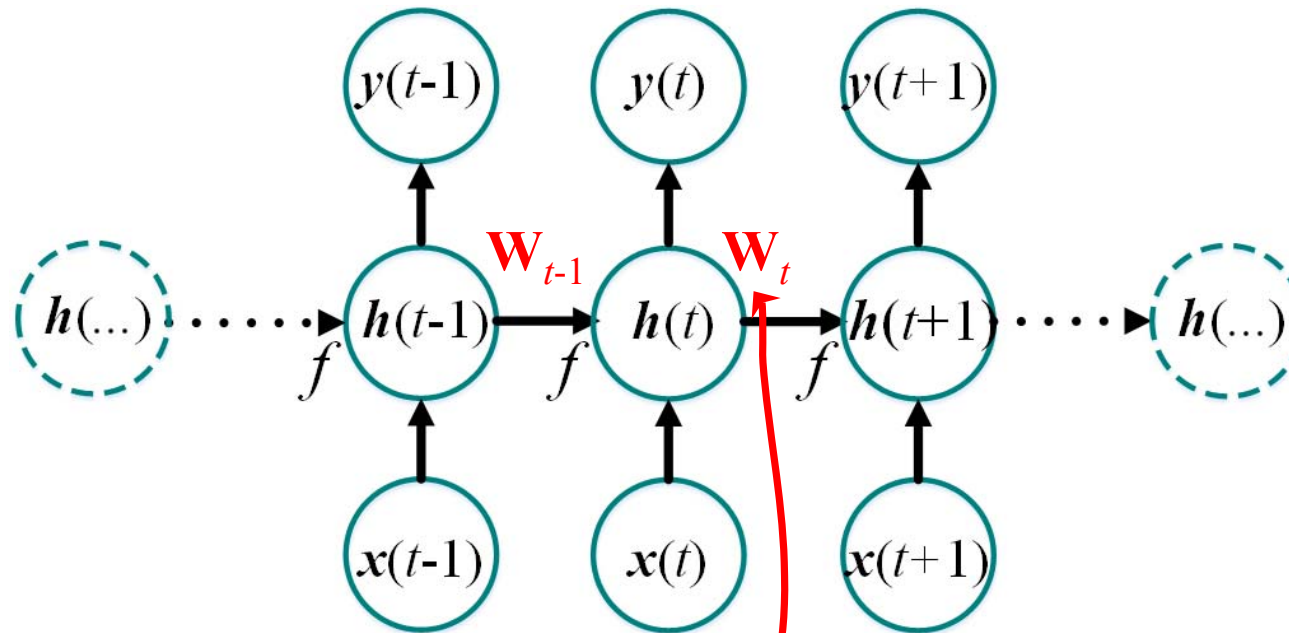
- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- **Backpropagation through time**
- ESN and LSTM

# Backpropagation through time



- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- **Backpropagation through time**
- ESN and LSTM

# Backpropagation through time



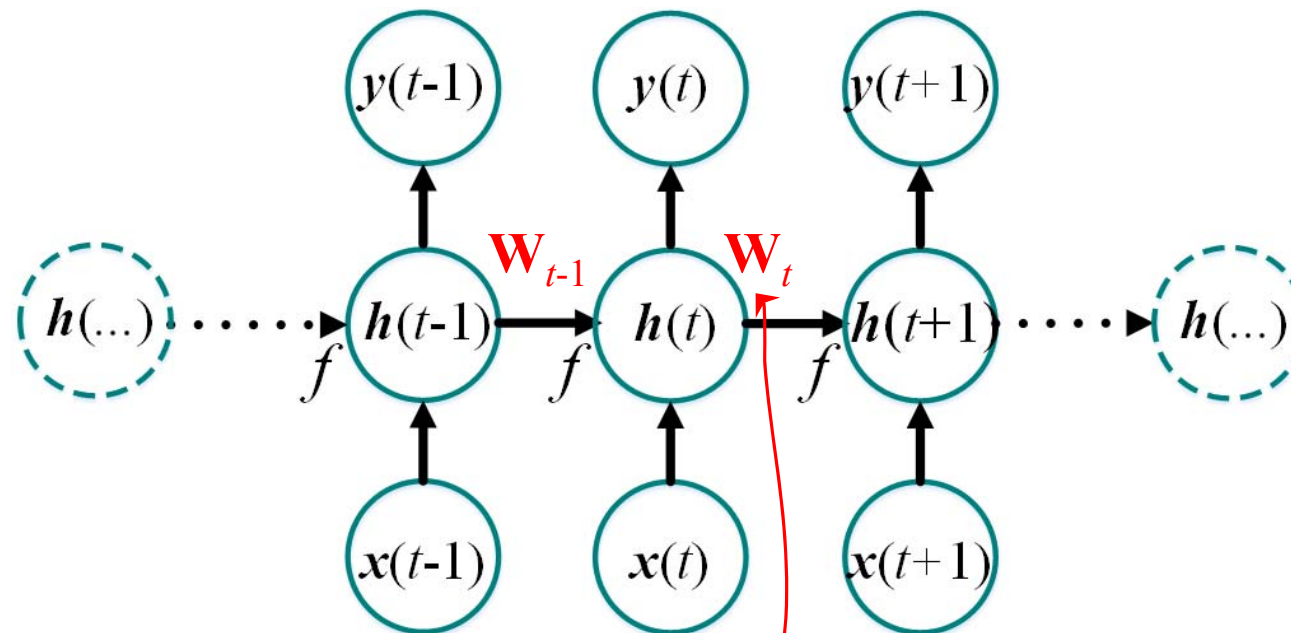
shared duplicated weights

Essentially,

$\mathbf{W} = \mathbf{W}_t$ ,  $\mathbf{W} = \mathbf{W}_{t-1}$  etc.

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- **Backpropagation through time**
- ESN and LSTM

# Backpropagation through time



shared duplicated weights

Essentially,

$\mathbf{W} = \mathbf{W}_t$ ,  $\mathbf{W} = \mathbf{W}_{t-1}$  etc.

$$\text{So: } \frac{\partial E}{\partial \mathbf{W}} = \frac{\partial E}{\partial \mathbf{W}_t} + \frac{\partial E}{\partial \mathbf{W}_{t-1}} + \dots$$

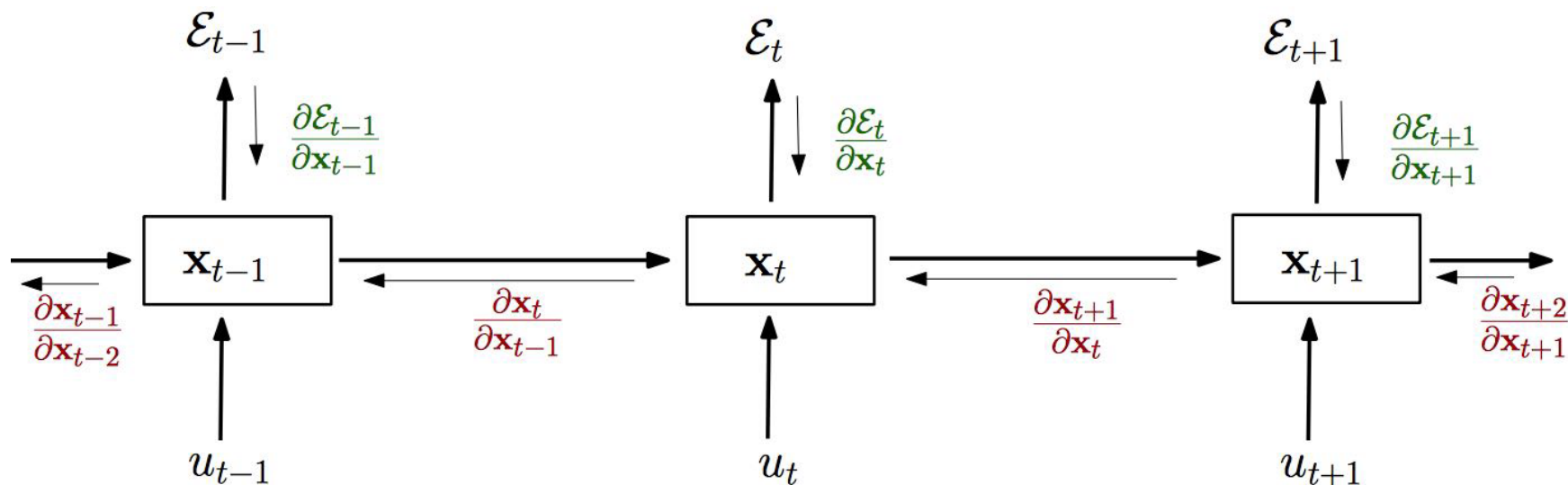
#time steps in a training sample

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- **Backpropagation through time**
- ESN and LSTM

# Backpropagation through time

For each pair of input-output sequences, the error is defined and  $n$  refers to the index over the duration of sequences (not the number of samples):

$$E = \sum_{n=1}^T (d(n) - y(n))^2 = \sum_{n=1}^T \varepsilon_n^2$$



Please note: According to our earlier notation,  $\mathbf{x}$  should be  $\mathbf{h}$  and  $u$  should be  $x$ .

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- **Backpropagation through time**
- ESN and LSTM

# Backpropagation through time

For each pair of input-output sequences, the error is defined and  $n$  refers to the index over the duration of sequences (not the number of samples):

$$E = \sum_{n=1}^T (d(n) - y(n))^2 = \sum_{n=1}^T \varepsilon_n^2$$

We can also think of this training algorithm in the time domain (*Hinton, 2013*):

- FORWARD PASS: a stack of the activities of all the units at each time step.
- BACKWARD PASS: activities are peeled off the stack to compute the error derivatives at each time step.
- THEN we add together the derivatives at all the different times for each weight.



- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- **Backpropagation through time**
- ESN and LSTM

# Backpropagation through time

For each pair of input-output sequences, the error is defined and  $n$  refers to the index over the duration of sequences (not the number of samples):

$$E = \sum_{n=1}^T (d(n) - y(n))^2 = \sum_{n=1}^T \varepsilon_n^2$$

Comments:

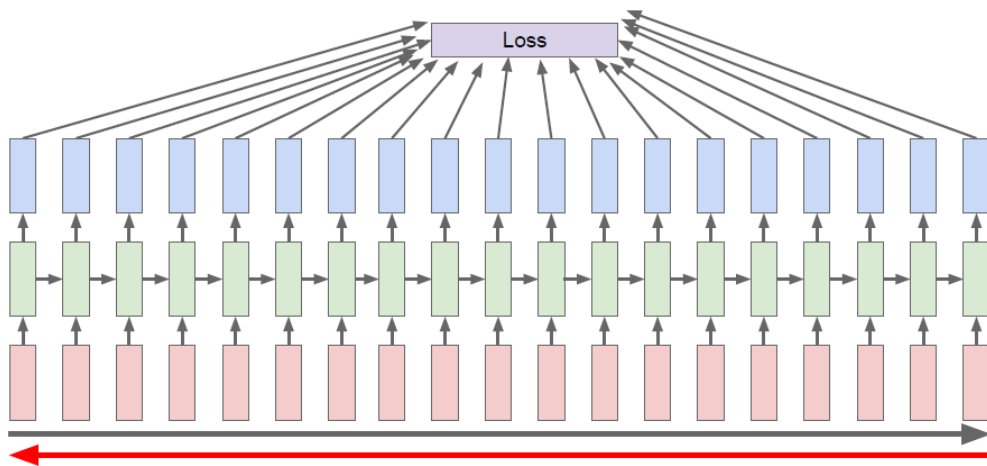
- epoch refers here to a single pair of input-output sequences
- if there are backprojections from the output to hidden layer, teacher output  $d(n)$  can be used in the computation of activations in layer  $n+1$  in the forward pass
  - teacher forcing is likely to speed the convergence
  - if it is exploited for the trained network however it may exhibit instability
- BPTT does not scale too well and may not converge even to the local minimum
- BPTT may result in instability

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- **Backpropagation through time**
- ESN and LSTM

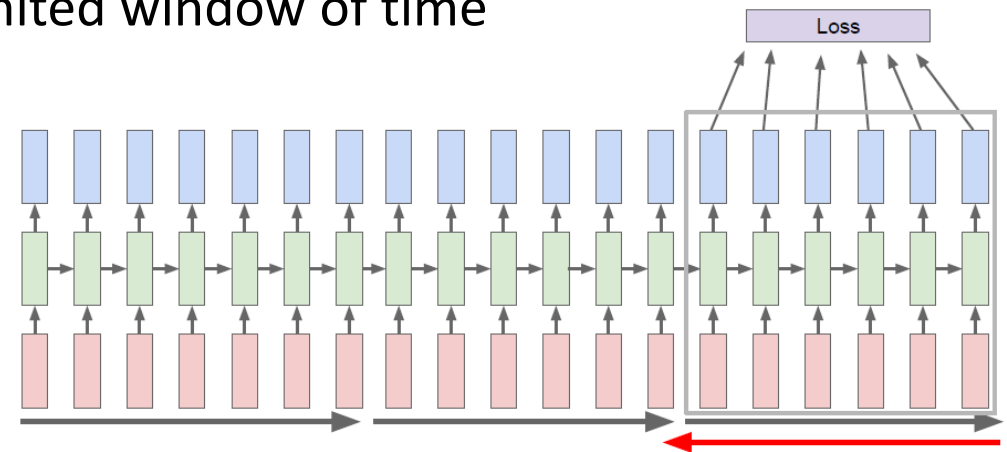
# Backpropagation through time – online with truncation

## Truncated BPTT for “online” learning

- if there are no batches and the network is supposed to work continuously, the number of recursive steps for weight updates has to be finite
- beyond the truncation there is no memory effect



with *truncation* backpropagate loss within a limited window of time



*For more details, please see Haykin and Goodfellow et al.*

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- **Backpropagation through time**
- ESN and LSTM

## Other learning algorithms for RNNs

- Real-time recurrent learning (Williams & Zipser, 1989)
  - the exact gradients are calculated and the synaptic updates are made at every step during network's processing stage
  - very high computational cost
- Kalman filters (*optimal filtering*)
  - solid theory formulated in the state-space concepts
  - rather than instantaneously estimate gradients, the network state is recursively estimated based on input data
  - the network has to be linearised
- Extended Kalman filter (*decoupled*)

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- **Backpropagation through time**
- ESN and LSTM

# Long-term dependencies problem

Vanishing (more rarely exploding) gradients for recurrent networks

$$\mathbf{h}(t) = \mathbf{W}^T \mathbf{h}(t-1) \Rightarrow \mathbf{h}(t) = (\mathbf{W}^t)^T \mathbf{h}(0)$$

$$\mathbf{W} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T \Rightarrow \mathbf{h}(t) = \mathbf{Q}^T \mathbf{\Lambda}^t \mathbf{Q} \mathbf{h}(0)$$



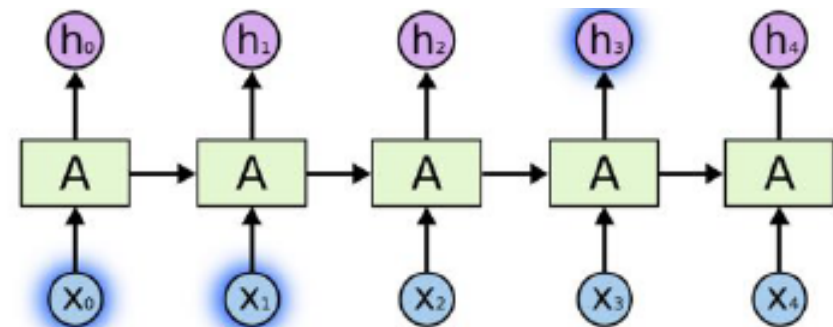
Eigenvalues are usually less than 1

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- **Backpropagation through time**
- ESN and LSTM

# Long-term dependencies problem

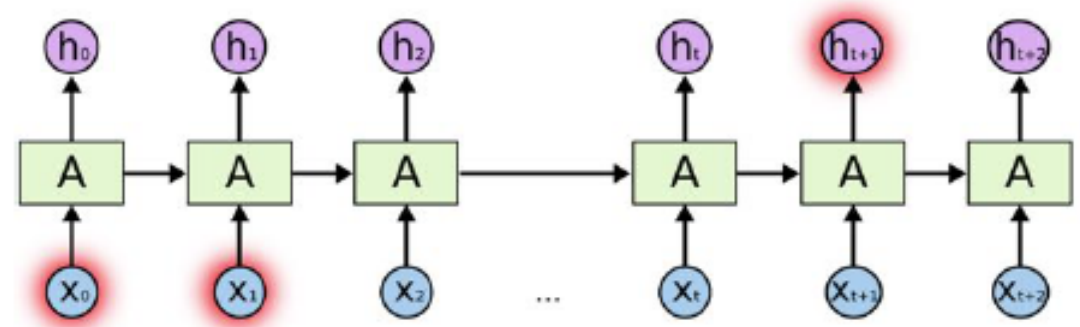
Handling (predicting) close words:

“The *clouds* are in the *sky*.”



Problem with the need for wider context:

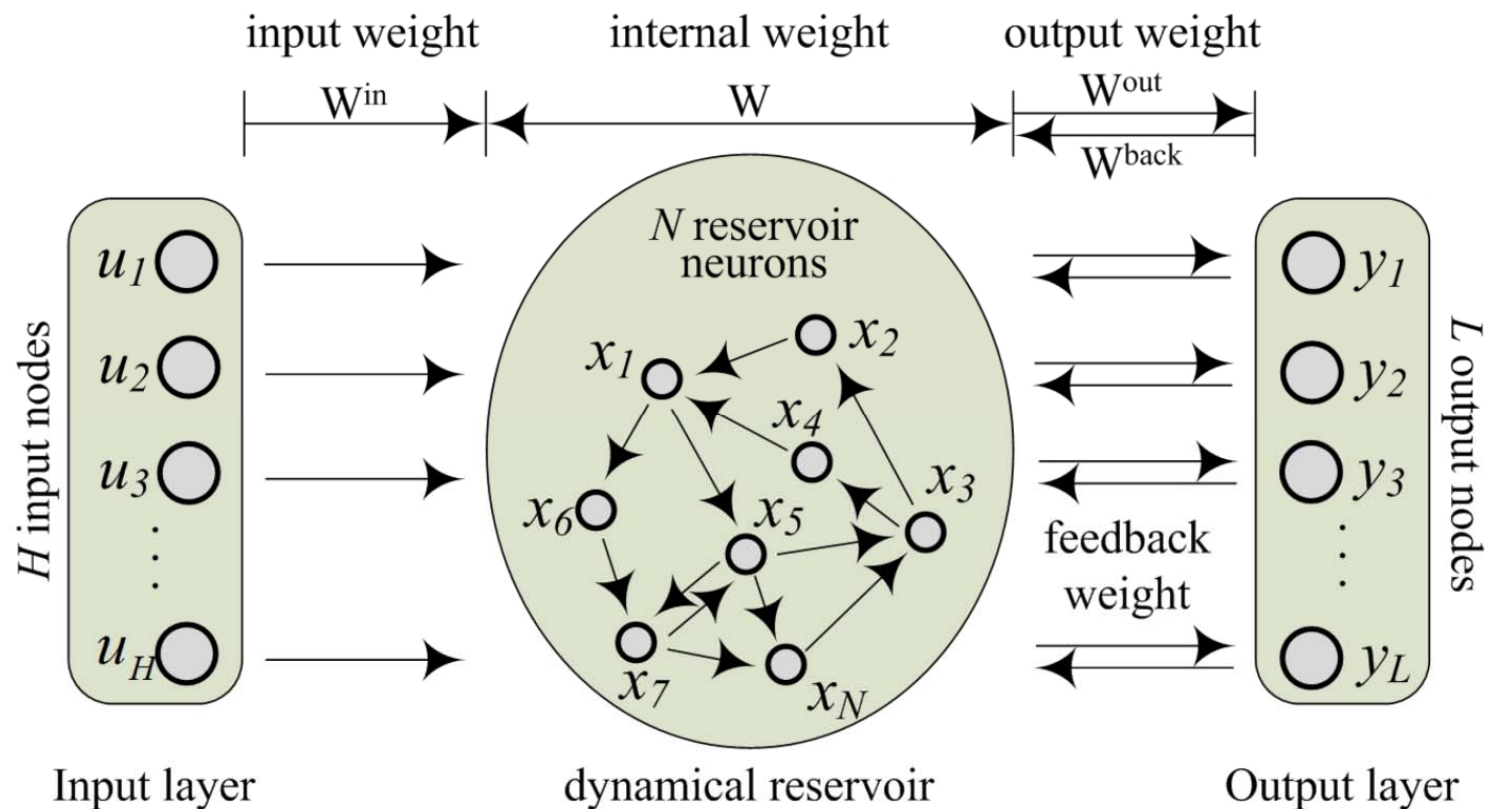
“I grew up in *France*... I speak fluent *French*.”



- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

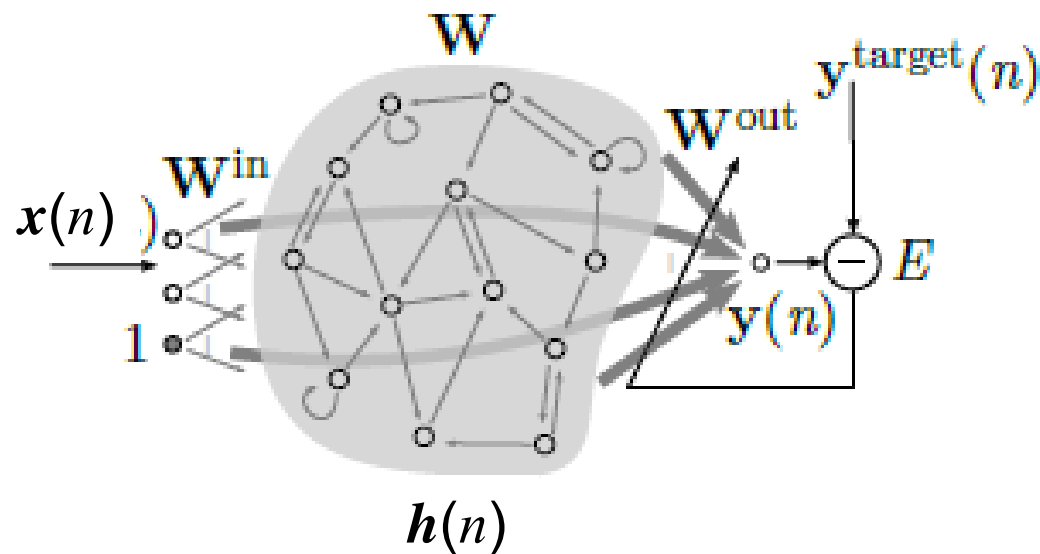
# Reservoir computing

## Echo state network (non-spiking version of liquid state machine)



- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Reservoir computing – fundamentals



$$y(n) = \mathbf{W}^{out} [1; \mathbf{h}(n); \mathbf{x}(n)]$$

biases

$$\tilde{\mathbf{h}}(n) = \tanh(\mathbf{W}_{in} [1; \mathbf{x}(n)] + \mathbf{W} \mathbf{h}(n-1))$$

$$\mathbf{h}(n) = (1 - \alpha) \mathbf{h}(n-1) + \alpha \tilde{\mathbf{h}}(n)$$

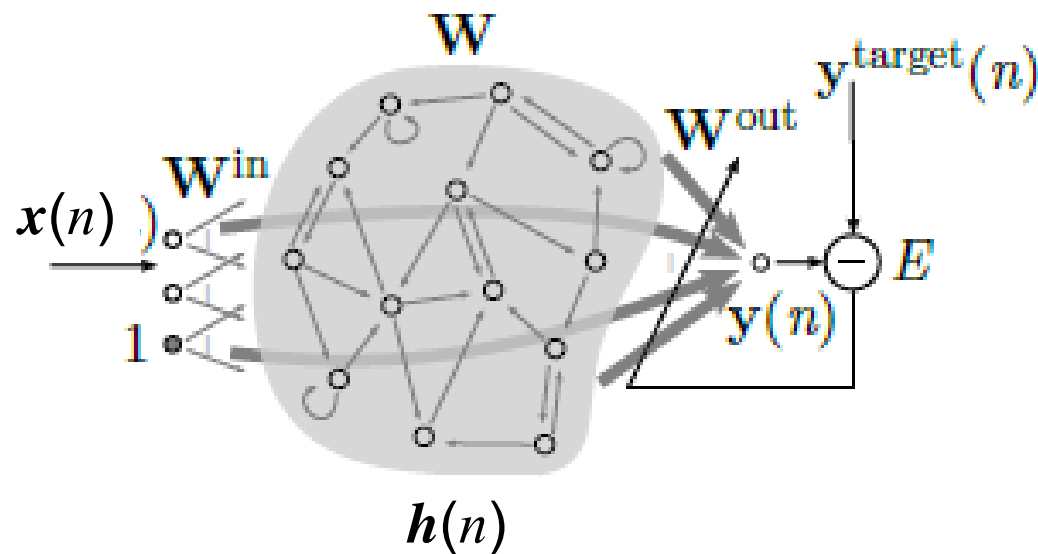
$$\mathbf{W} \in \mathbb{R}^{N_x \times N_x}$$

$$\mathbf{W}_{in} \in \mathbb{R}^{N_x \times (1 + N_u)}$$

$$\mathbf{W}^{out} \in \mathbb{R}^{N_y \times (1 + N_h + N_x)}$$

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Reservoir computing – fundamentals



$$y(n) = \mathbf{W}^{out} [1; \mathbf{h}(n); x(n)]$$

biases

$$\tilde{\mathbf{h}}(n) = \tanh(\mathbf{W}_{in} [1; x(n)] + \mathbf{W} \mathbf{h}(n-1) + \mathbf{W}^{fb} y(n-1))$$

$$\mathbf{h}(n) = (1 - \alpha) \mathbf{h}(n-1) + \alpha \tilde{\mathbf{h}}(n)$$

extra *feedback* projections from the output to the hidden layer

$$\mathbf{W} \in \mathbb{R}^{N_x \times N_x}$$

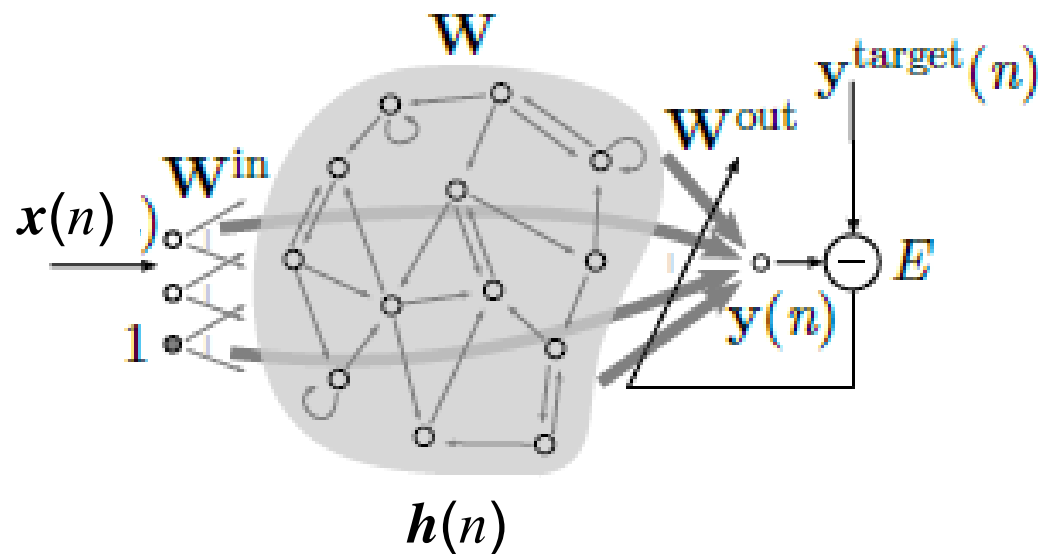
$$\mathbf{W}_{in} \in \mathbb{R}^{N_x \times (1 + N_u)}$$

$$\mathbf{W}^{out} \in \mathbb{R}^{N_y \times (1 + N_h + N_x)}$$



- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Reservoir computing – fundamentals



$$y(n) = \mathbf{W}^{out} [1; \mathbf{h}(n); x(n)]$$

biases

$$\tilde{\mathbf{h}}(n) = \tanh(\mathbf{W}_{in} [1; x(n)] + \mathbf{W} \mathbf{h}(n-1) + \mathbf{W}^{fb} y(n-1))$$

$$\mathbf{h}(n) = (1 - \alpha) \mathbf{h}(n-1) + \alpha \tilde{\mathbf{h}}(n)$$



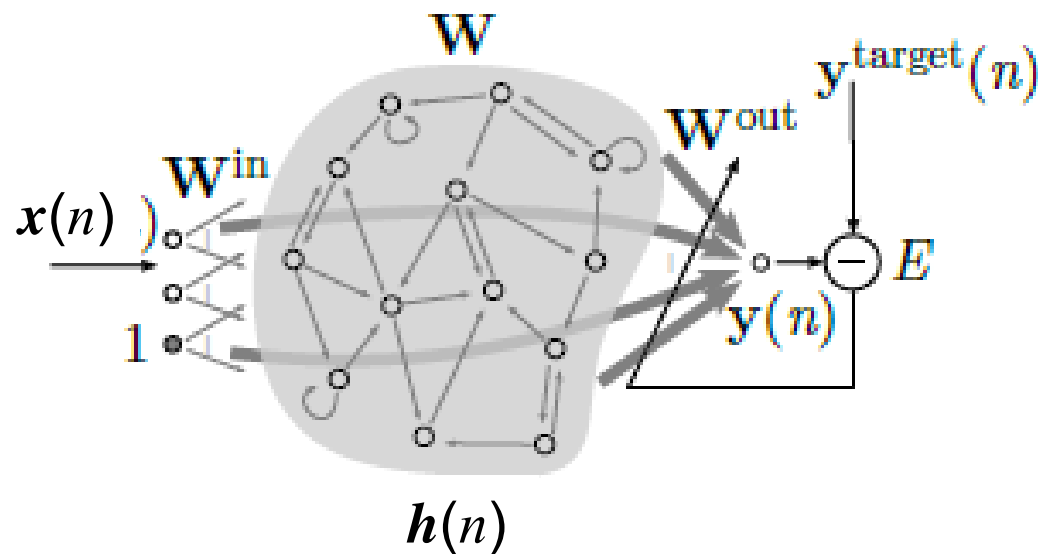
extra *feedback* projections from the output to the hidden layer

Could be used for providing feedback for training – *teacher forcing*:

$$y(n) = y^{target}(n)$$

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Reservoir computing – fundamentals



$$y(n) = \mathbf{W}^{out} [1; \mathbf{h}(n); \mathbf{x}(n)]$$

biases

$$\tilde{\mathbf{h}}(n) = \tanh(\mathbf{W}_{in} [1; \mathbf{x}(n)] + \mathbf{W} \mathbf{h}(n-1) + \mathbf{W}^{fb} \mathbf{y}(n-1))$$

$$\mathbf{h}(n) = (1 - \alpha) \mathbf{h}(n-1) + \alpha \tilde{\mathbf{h}}(n)$$

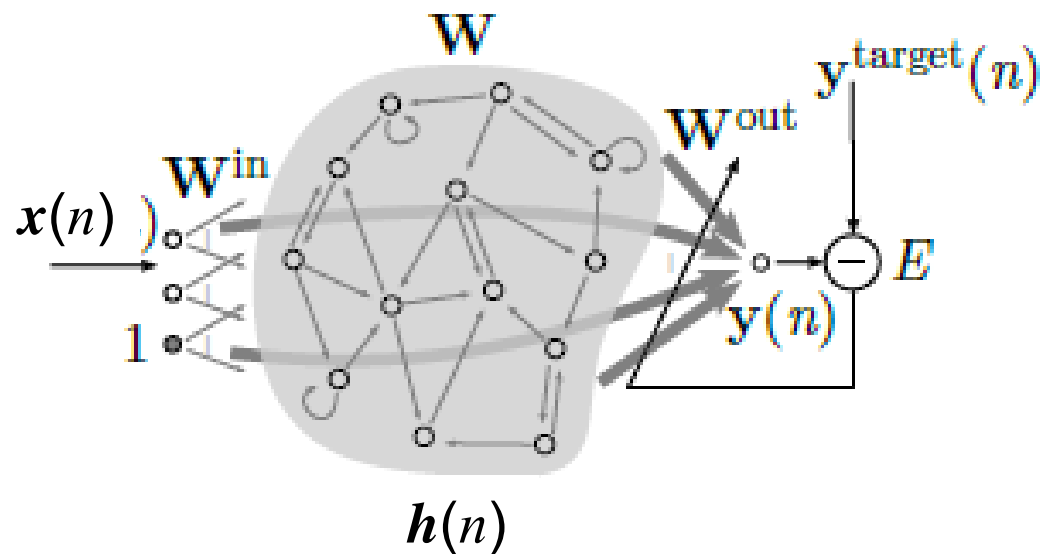
HOWEVER, in on-line learning the use of  $y(n)$  for feedback is preferred for the stability.

extra *feedback* projections from the output to the hidden layer

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Reservoir properties

**Reservoir serves as a memory (temporal context) and a nonlinear expansion**



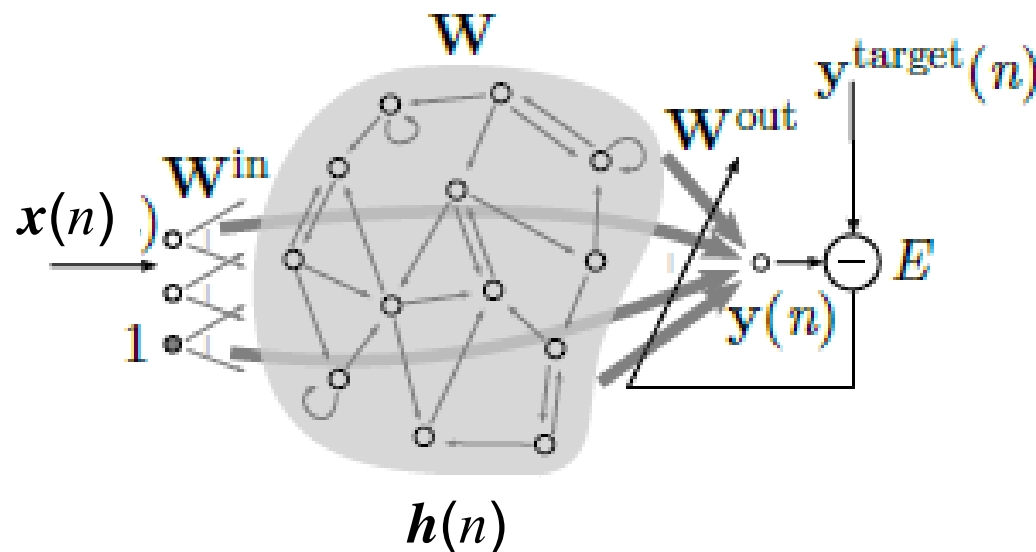
Key parameters:

- size ( $N_x$ ) (the bigger the better, even in the order of 10k)
- sparsity (2-20%) and the distribution of nonzero elements (uniform distribution)
- spectral radius,  $\rho$  (less than 1 to be near the edge of stability)

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Reservoir properties

**Reservoir serves as a memory (temporal context) and a nonlinear expansion**



Key parameters:

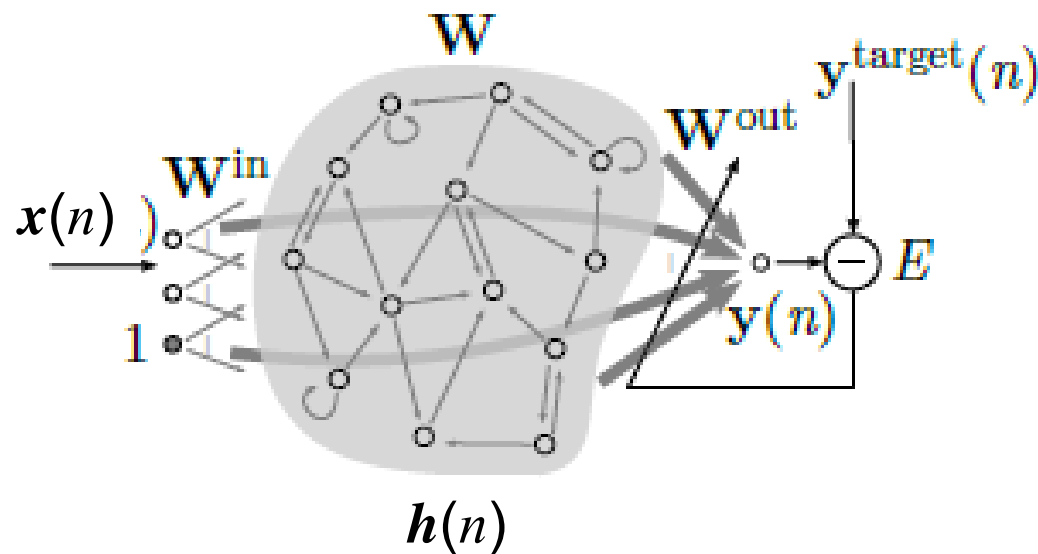
- size ( $N_x$ ) (the bigger the better, even in the order of 10k)
- sparsity (2-20%) and the distribution of nonzero elements (uniform distribution)
- spectral radius,  $\rho$  (less than 1 to be near the edge of stability)

For nonlinear networks, the system can still be contractive and stable for  $\rho > 1$ .  
So, commonly in practice,  $\rho \approx 3$ .

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Reservoir properties

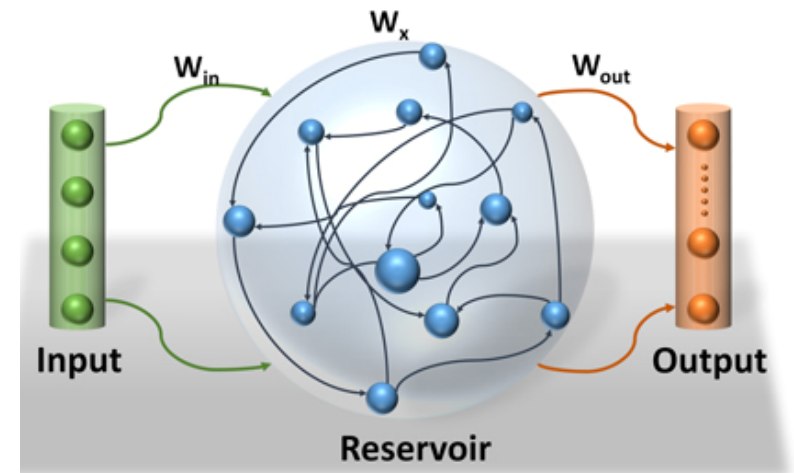
**Reservoir serves as a memory (temporal context) and a nonlinear expansion**



Key parameters:

- size ( $N_x$ )
- sparsity (2-20%)
- spectral radius,  $\rho$

Sparse, random and fixed connections in the reservoir...



....with “leaky” units

Leak modulated by  $\alpha$

$$h(t) = \alpha h(t-1) + (1 - \alpha)x(t)$$

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Training readouts

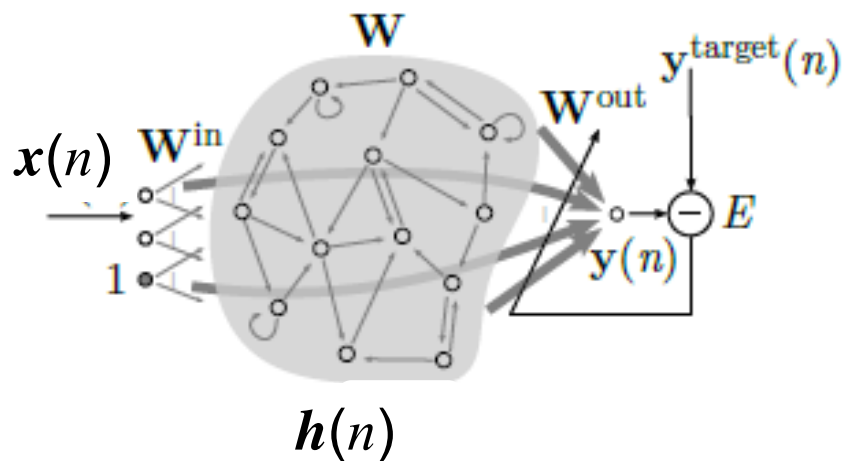
- One-shot learning with least mean square (LMS) error approaches
  - the design matrix is usually overdetermined -> ridge regression (regularization):

$$\mathbf{Y}^{target} = \mathbf{W}^{out} \mathbf{H} \rightarrow \underline{\mathbf{W}^{out} = \mathbf{Y}^{target} \mathbf{H}^T (\mathbf{H} \mathbf{H}^T + \beta \mathbf{I})^{-1}}$$

- be careful with extremely large values in  $\mathbf{W}$  as they may indicate problems with stability
- Online learning with, for example, *recursive* LMS
- The need to deal with initial transients, esp. for long sequences
- Possible use of teacher forcing (forcing output data through backprojections)

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Reservoir computing – overall recipe



## Recipe

1. Generation of the dynamic reservoir ( $\mathbf{W}_{in}$ ,  $\mathbf{W}$ )
2. Application of inputs,  $x(n)$ , and collecting the corresponding activation states,  $\mathbf{h}(n)$ .
3. Computation of the linear output weights from the reservoir with a linear regression approach (MSE error to be minimised).
4. Use of the RNN for new data – predictions.

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Reservoir computing

- set spectral radius  $\rho$ , large value means slow forgetting, e.g. storage of long time scales
- input scaling, small input means reservoir nodes operate on linear regime, large means binary operation
- output feedback weights if autonomous pattern generation needed (attractor property)
- connectivity structure: small world, distribution of loop lengths, fixed number of inputs to each node
- propagation delay on a fraction of connections



- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Long short-term memory (LSTM)

## Main motivation behind LSTMs

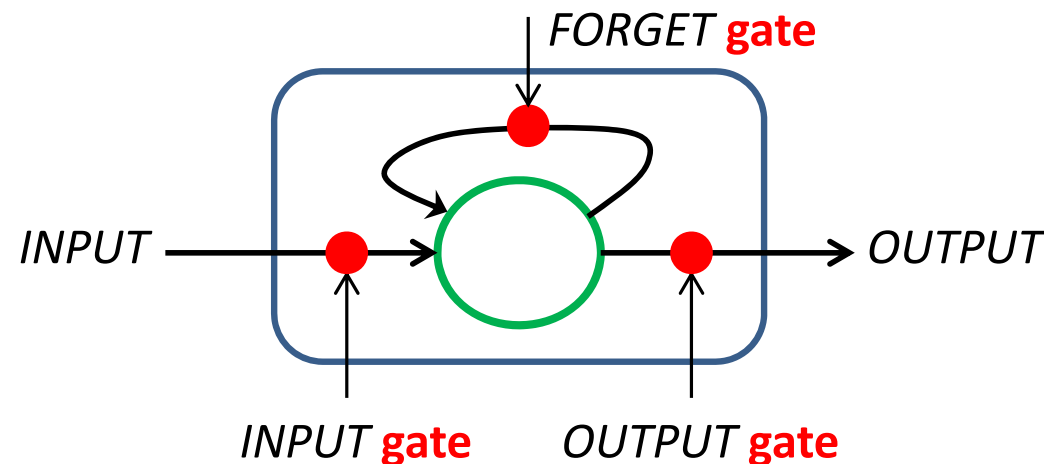
- vanishing gradients when using backprop through time for RNNs
- poor capacity to handle long-term dependencies

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Long short-term memory (LSTM)

## Main motivation behind LSTMs

- vanishing gradients when using backprop through time for RNNs
- poor capacity to handle long-term dependencies



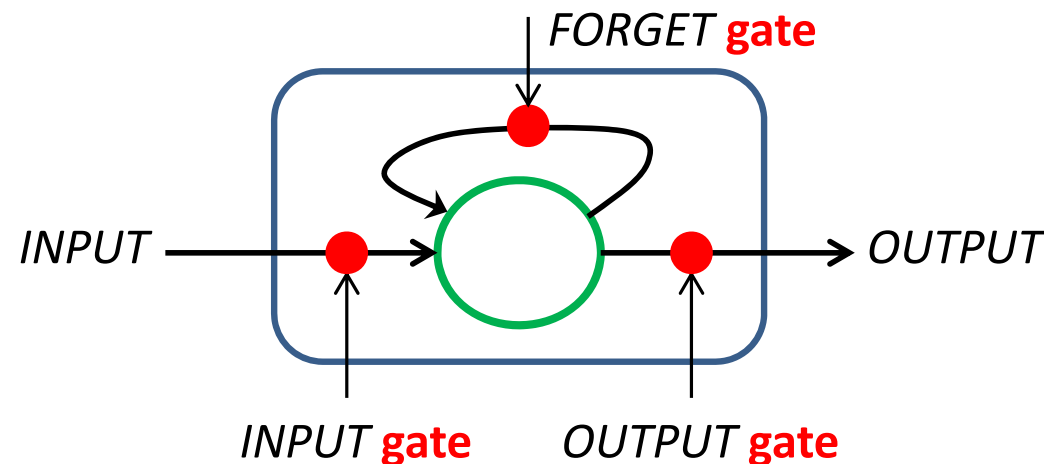
The key idea is to have a “memory cell” capable of keeping the state over time

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Long short-term memory (LSTM)

## Main motivation behind LSTMs

- vanishing gradients when using backprop through time for RNNs
- poor capacity to handle long-term dependencies



The key idea is to have a “memory cell” capable of keeping the state over time

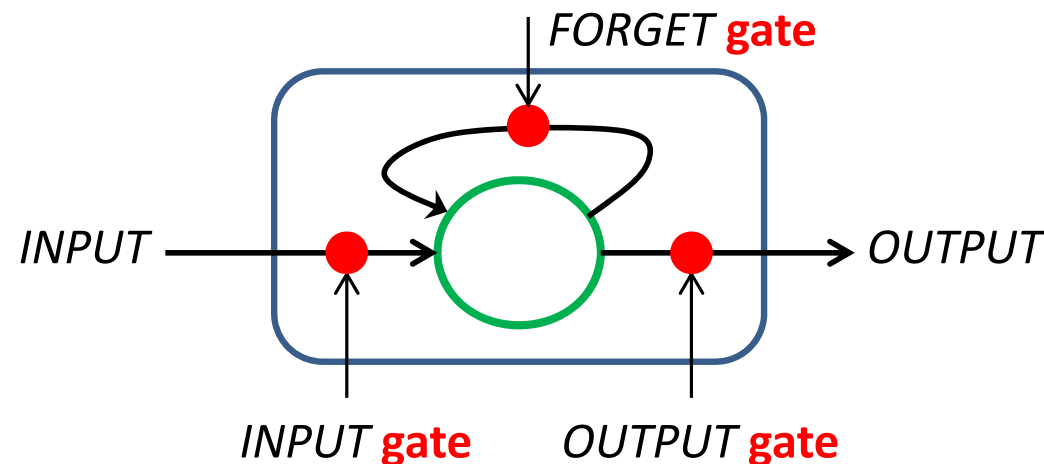
forget gate acts like a “leak” with tuneable gain  $h(t) = \alpha h(t-1) + (1 - \alpha)x(t)$

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Long short-term memory (LSTM)

## Main motivation behind LSTMs

- vanishing gradients when using backprop through time for RNNs
- poor capacity to handle long-term dependencies

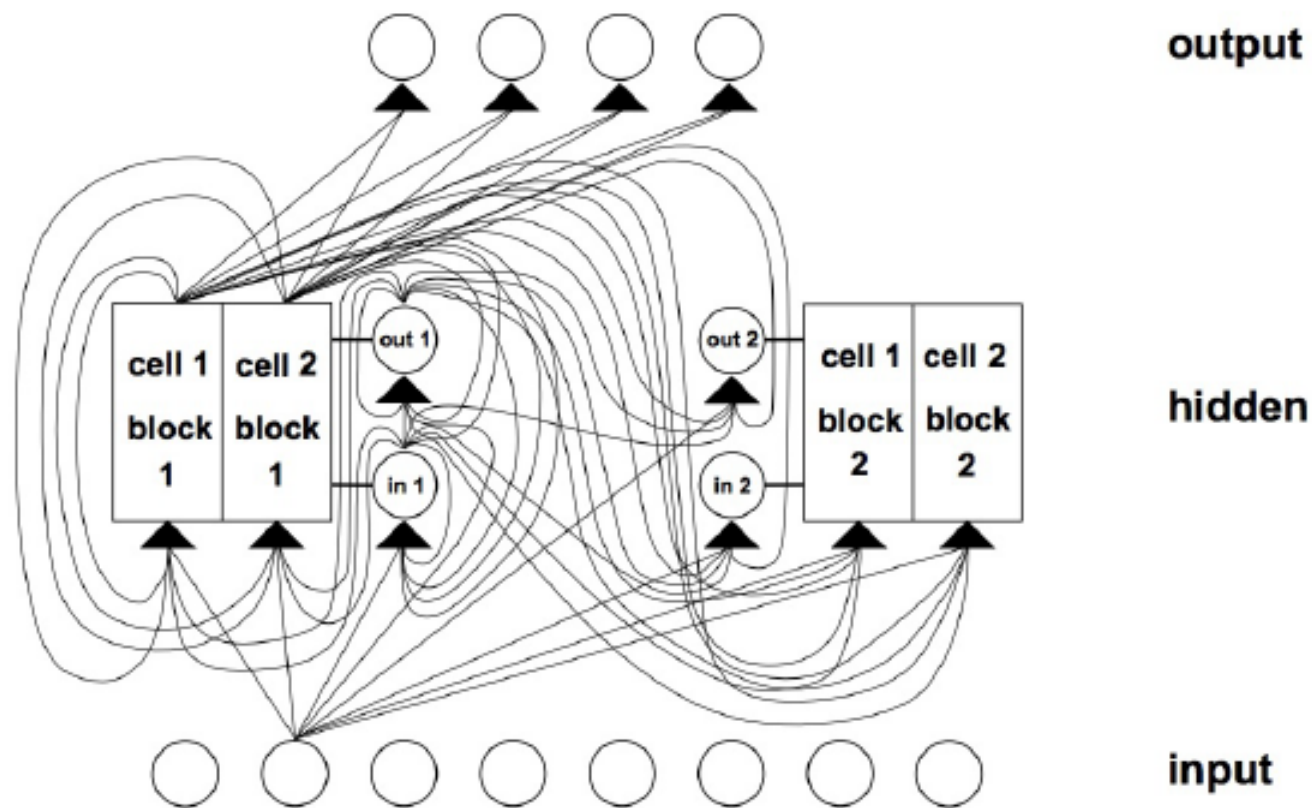


The key idea is to have a “memory cell” capable of keeping the state over time

- explicit memory – **cell state vector** (on top of the *hidden* state)
- regulatory mechanism for information flow in and out – **gating unit**

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

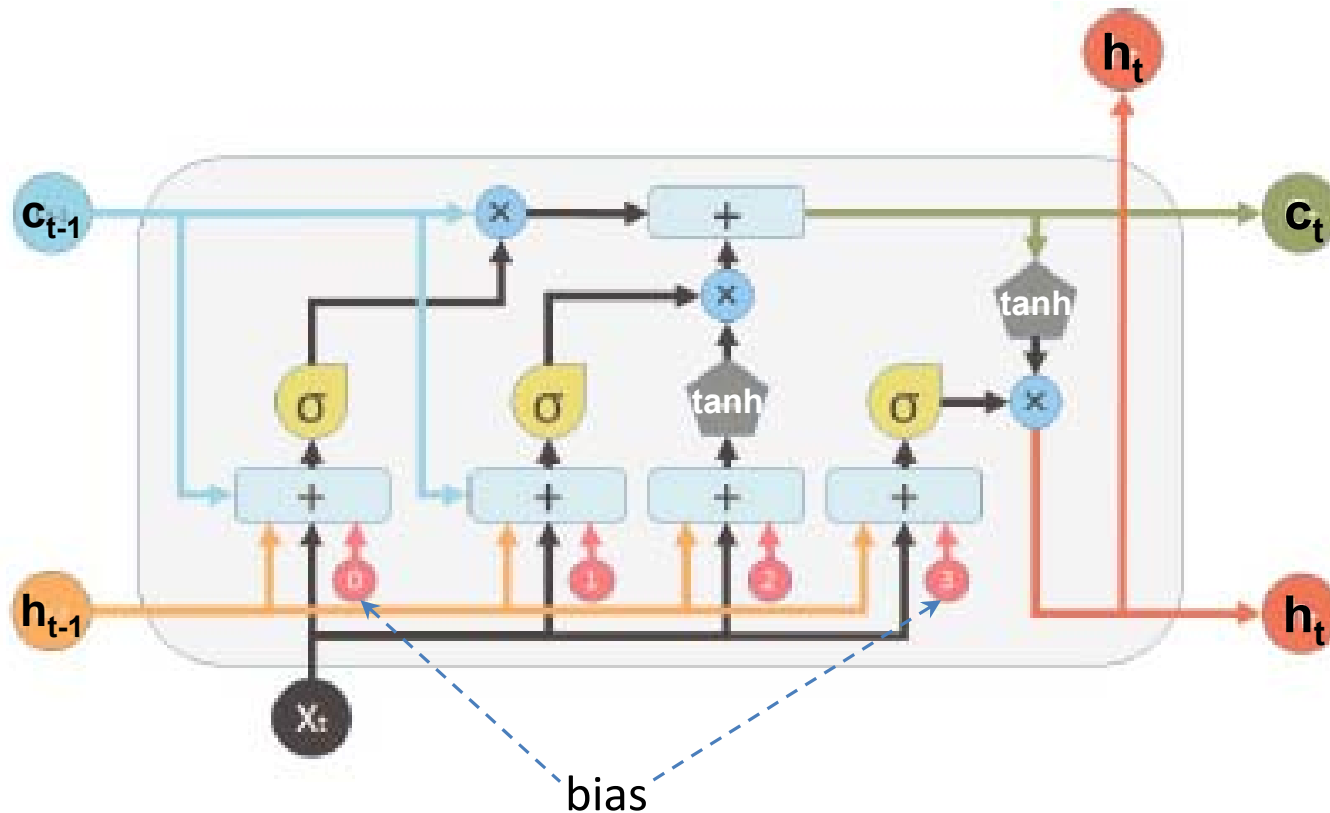
# Long short-term memory (LSTM) network



Hochreiter S, Schmidhuber J. Long short-term memory.  
*Neural computation*. 1997 Nov 15;9(8):1735-80.

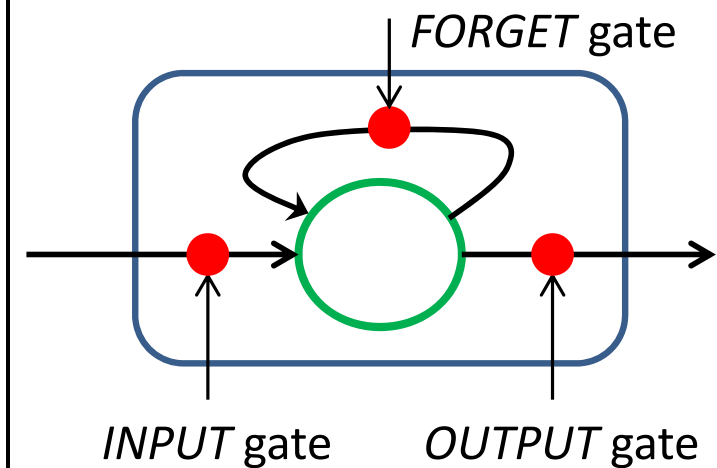
- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Long short-term memory (LSTM) cell



- $x_t$  input
- $c_{t-1}$  memory from previous cell
- $h_{t-1}$  output from previous cell

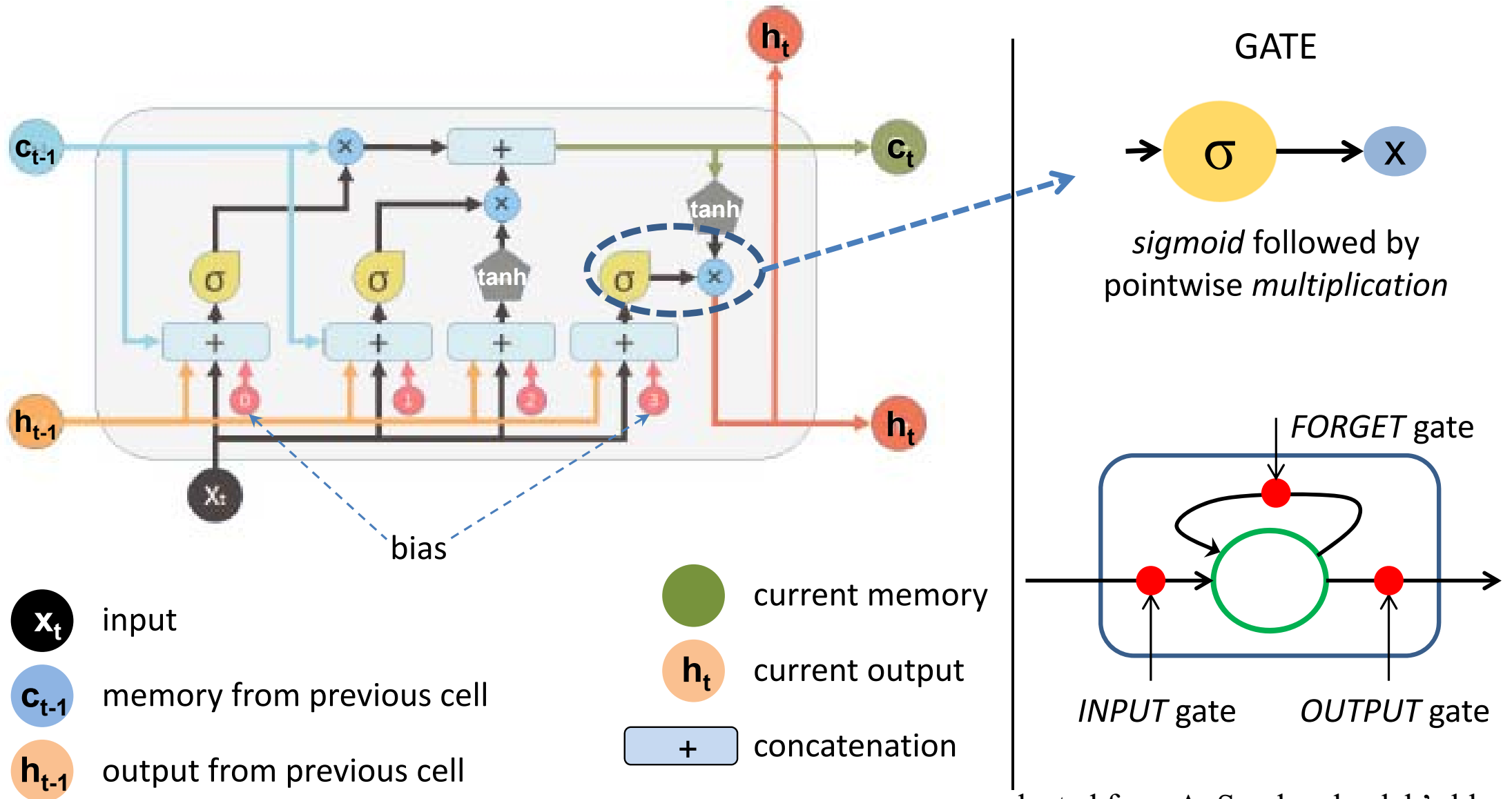
- current memory
- $h_t$  current output
- +
- concatenation



adapted from A. Sood and colah's blog

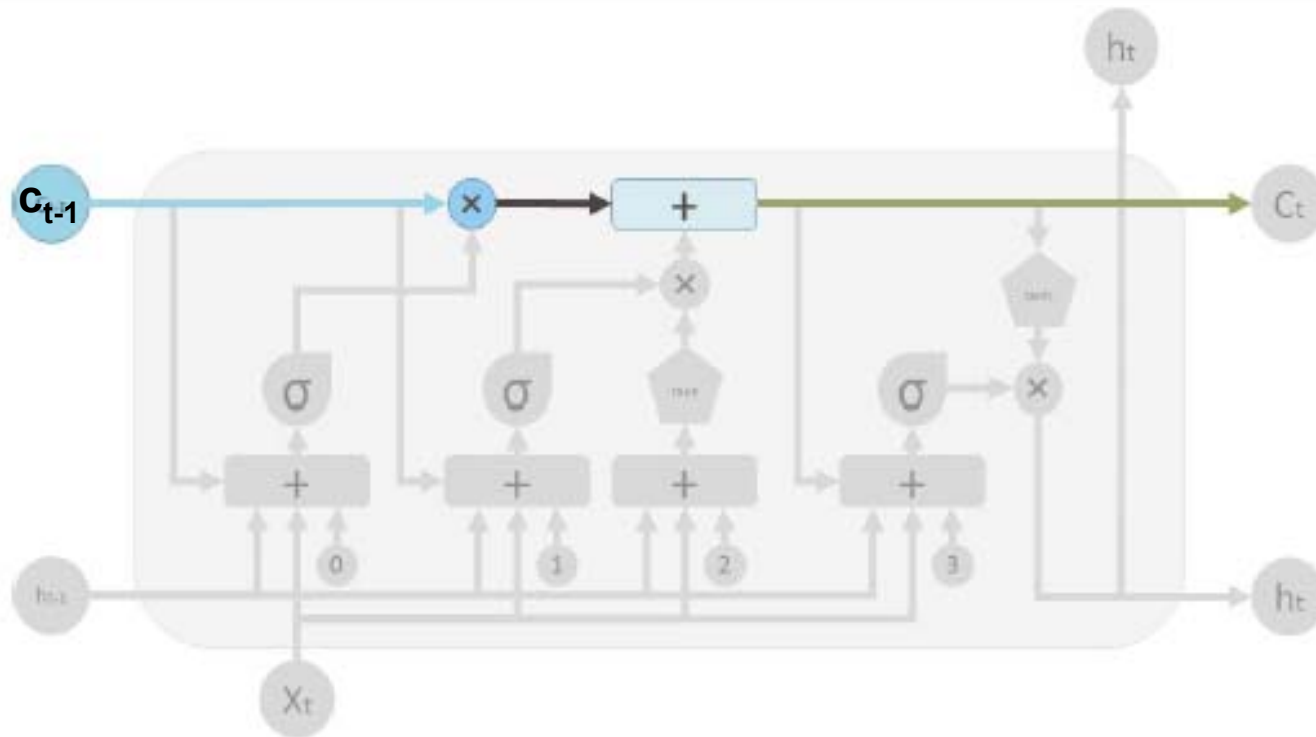
- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Long short-term memory (LSTM) cell



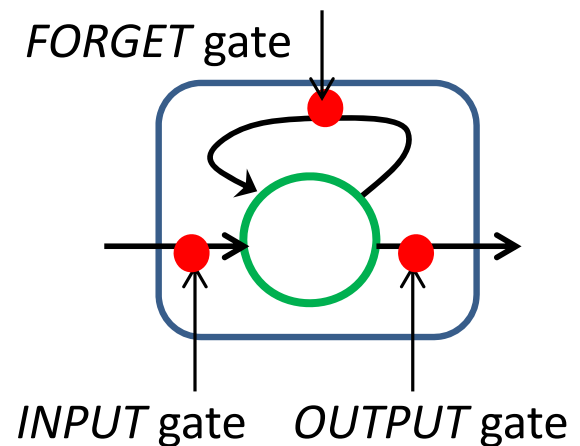
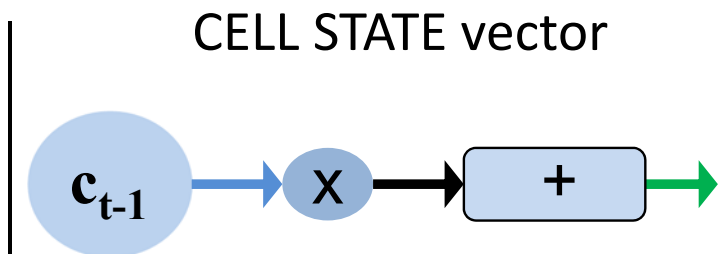
- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Memory cell components – cell state vector



## Cell state vector

- represents memory
- it is changing as a result of new information (input gate) and forgetting (forget gate)

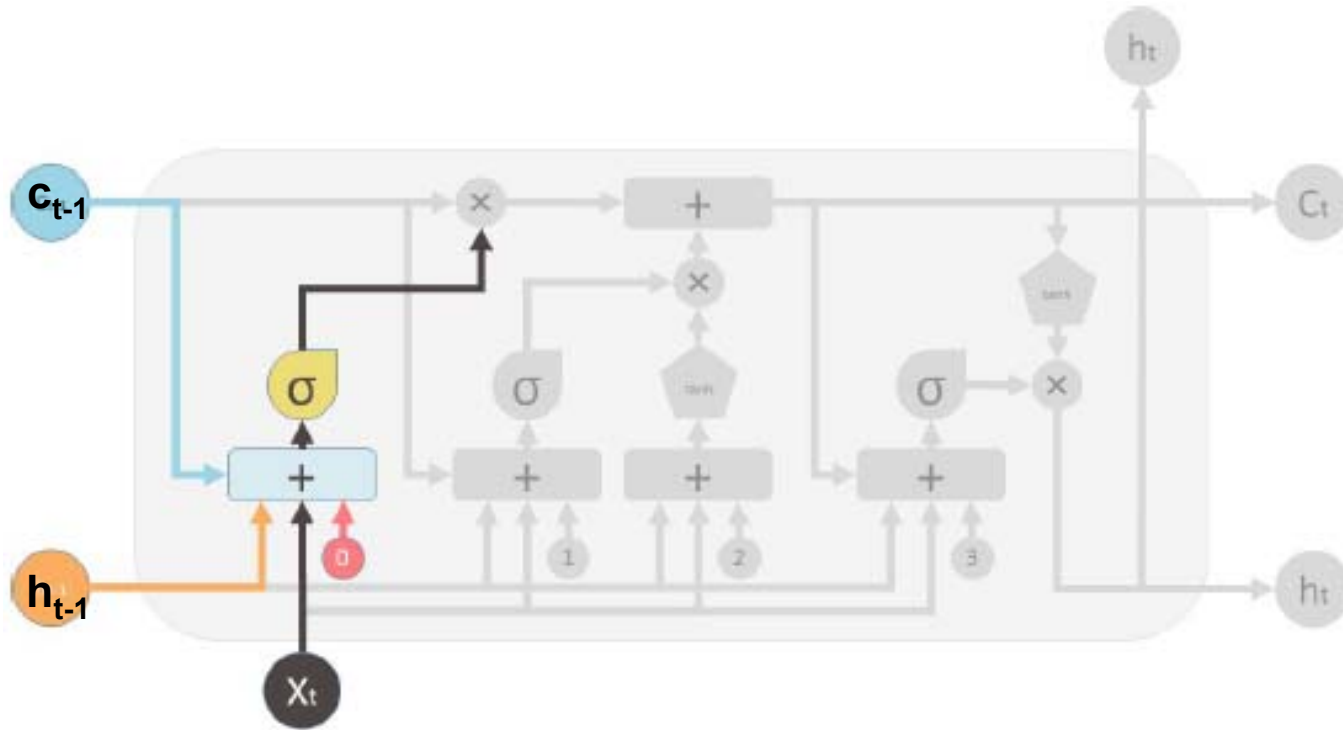


adapted from A. Sood and colah's blog



- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

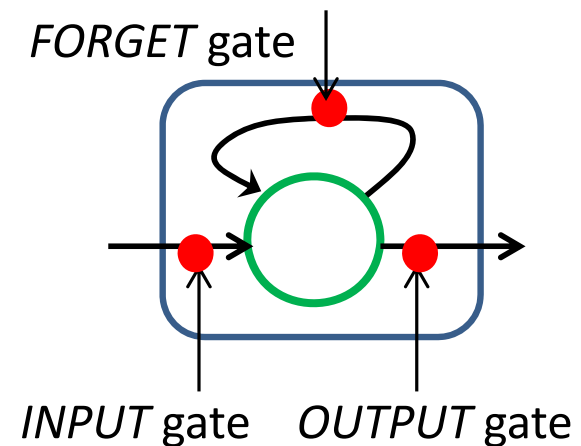
# Memory cell components – forget gate



## Forget gate

- controls what information should be forgotten (removed from memory)

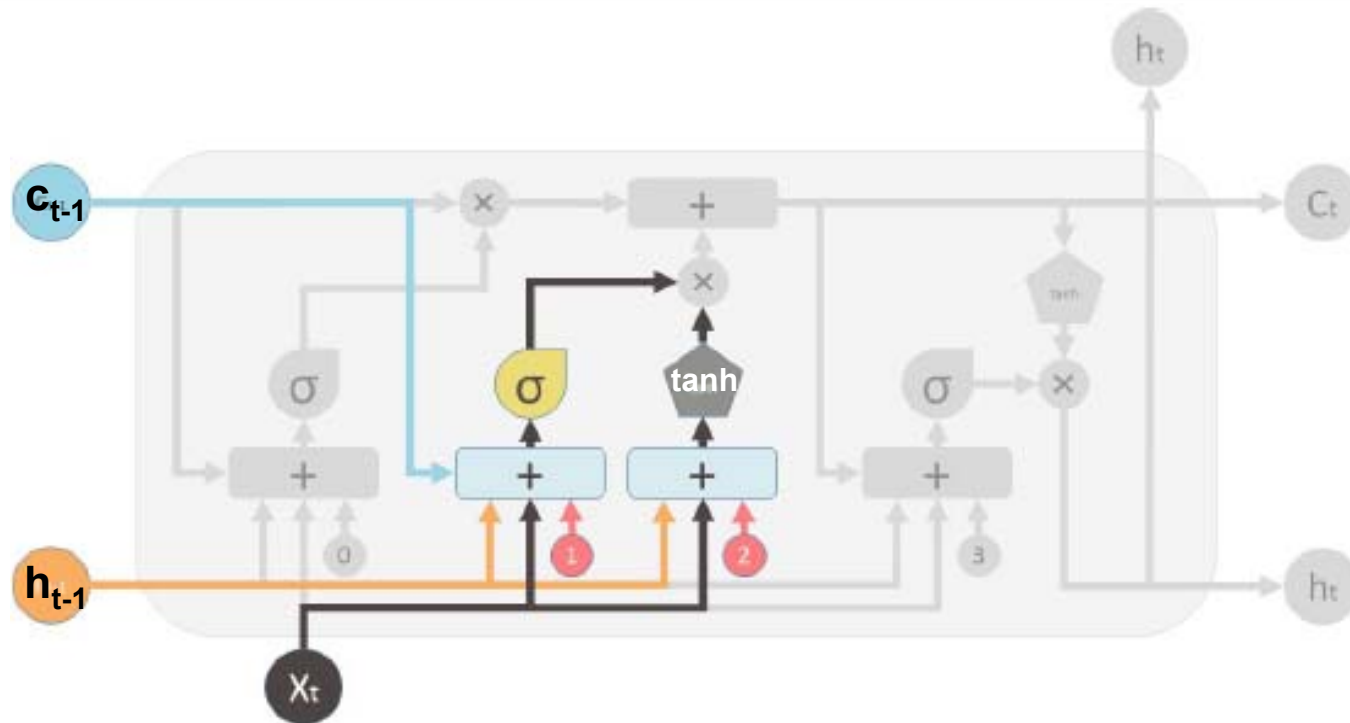
$$f(t) = \sigma(\mathbf{W}_f [h(t-1), x(t)] + b_f)$$



adapted from A. Sood and colah's blog

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Memory cell components – input gate

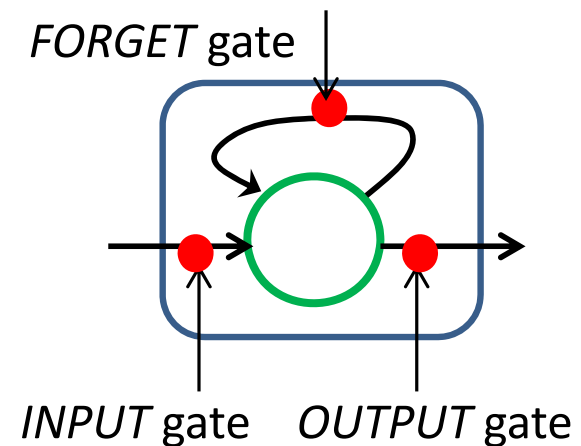


## Input gate

- controls what information should be added to cell state from the current input

$$i(t) = \sigma(\mathbf{W}_i[h(t-1), x(t)] + b_i)$$

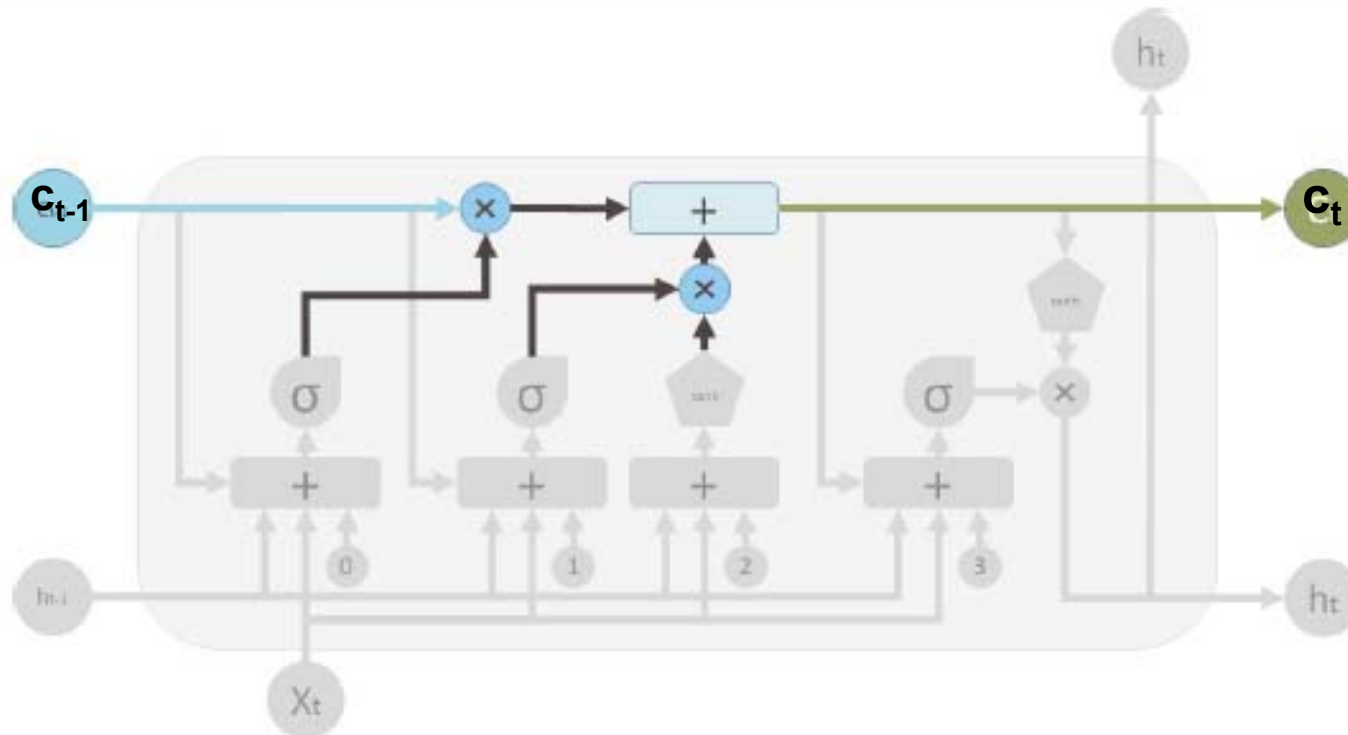
$$\tilde{C}(t) = \tanh(\mathbf{W}_c[h(t-1), x(t)] + b_c)$$



adapted from A. Sood and colah's blog

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

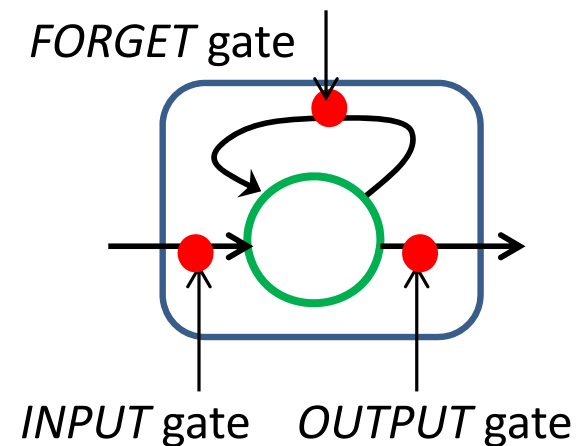
# Memory update



## Updating memory

- cell state vector aggregates old memory, gated by forget gate, and a new memory, filtered by the input gate

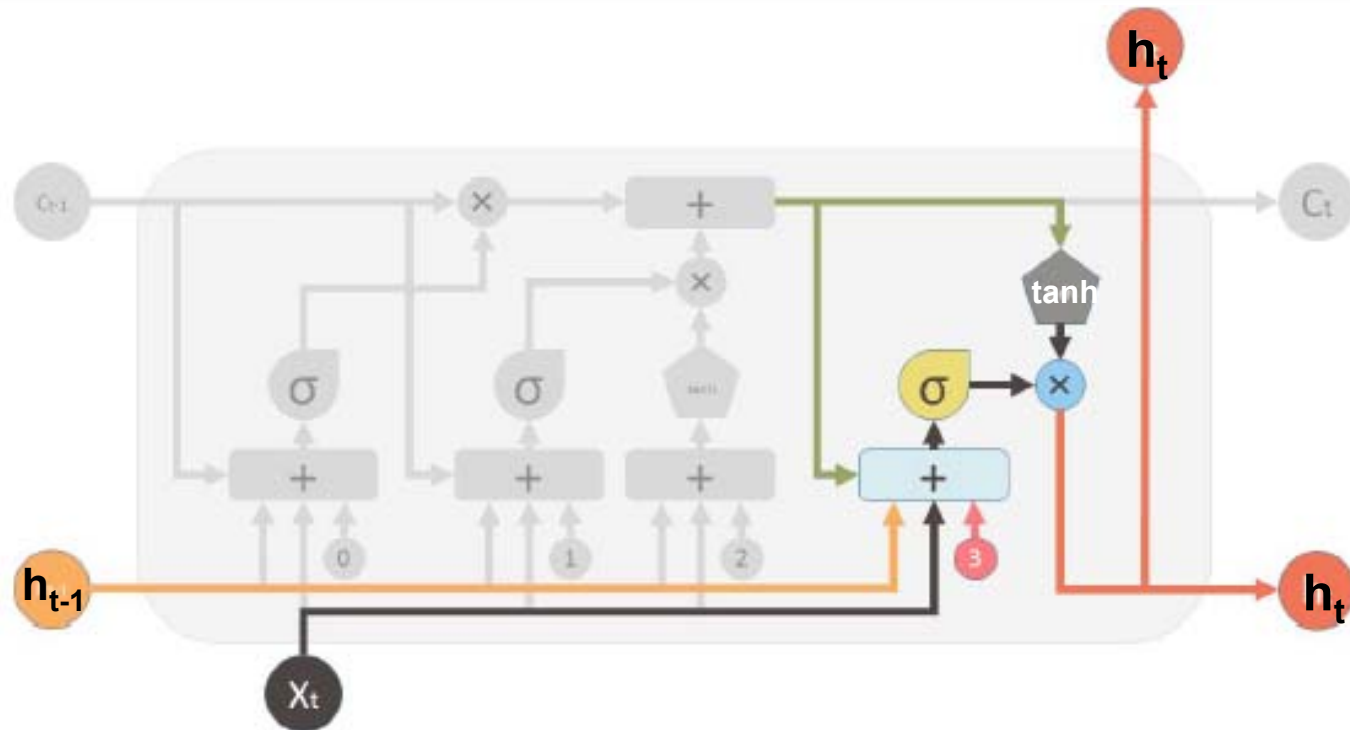
$$C(t) = f(t) \odot C(t-1) + i(t) \odot \tilde{C}(t)$$



adapted from A. Sood and colah's blog

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Memory cell components – output gate

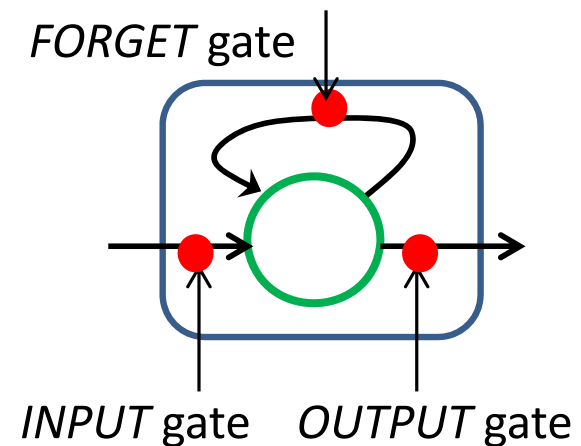


## Output gate

- controls what information is sent to the output

$$o(t) = \sigma(\mathbf{W}_o[h(t-1), x(t)] + b_o)$$

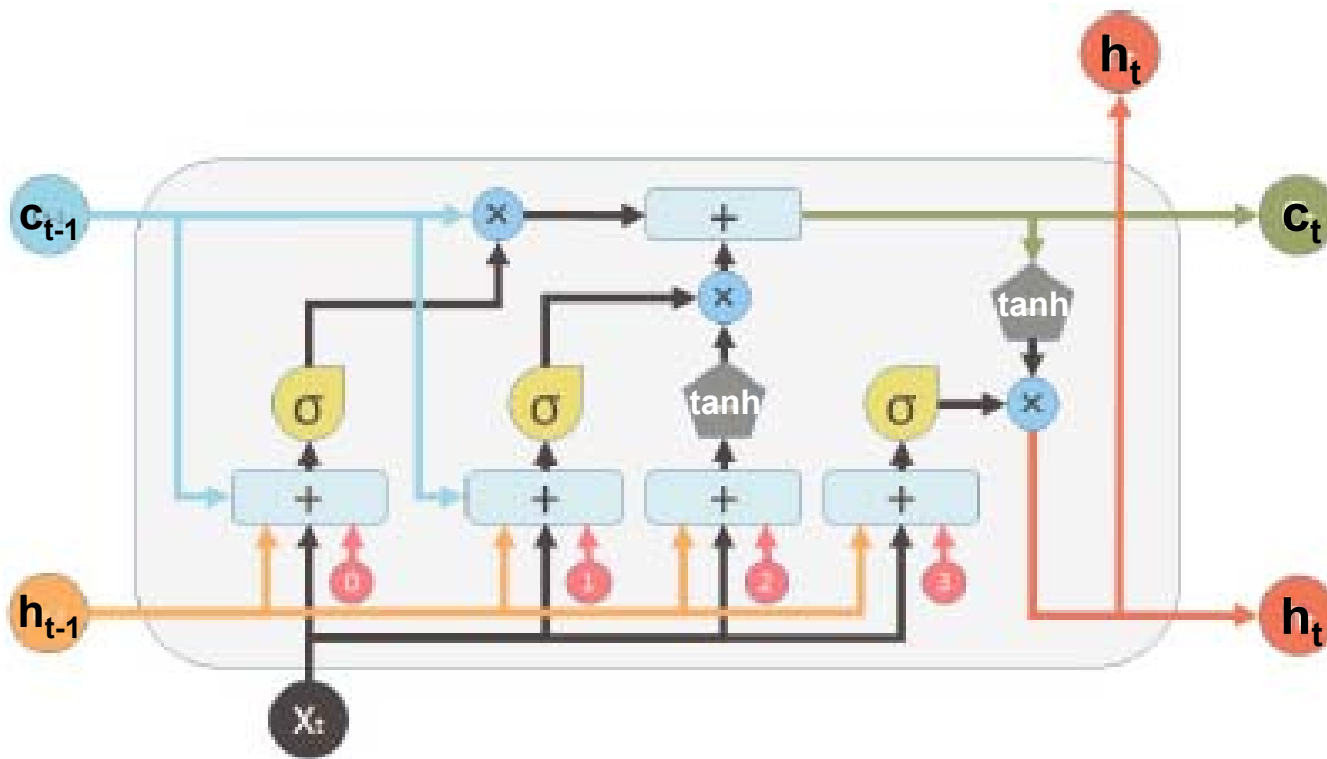
$$h(t) = o(t) \odot \tanh(C(t))$$



adapted from A. Sood and colah's blog

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# LSTM memory cell summary



$$f(t) = \sigma(\mathbf{W}_f[h(t-1), x(t)] + b_f)$$

$$i(t) = \sigma(\mathbf{W}_i[h(t-1), x(t)] + b_i)$$

$$o(t) = \sigma(\mathbf{W}_o[h(t-1), x(t)] + b_o)$$

$$\tilde{C}(t) = \tanh(\mathbf{W}_c[h(t-1), x(t)] + b_c)$$

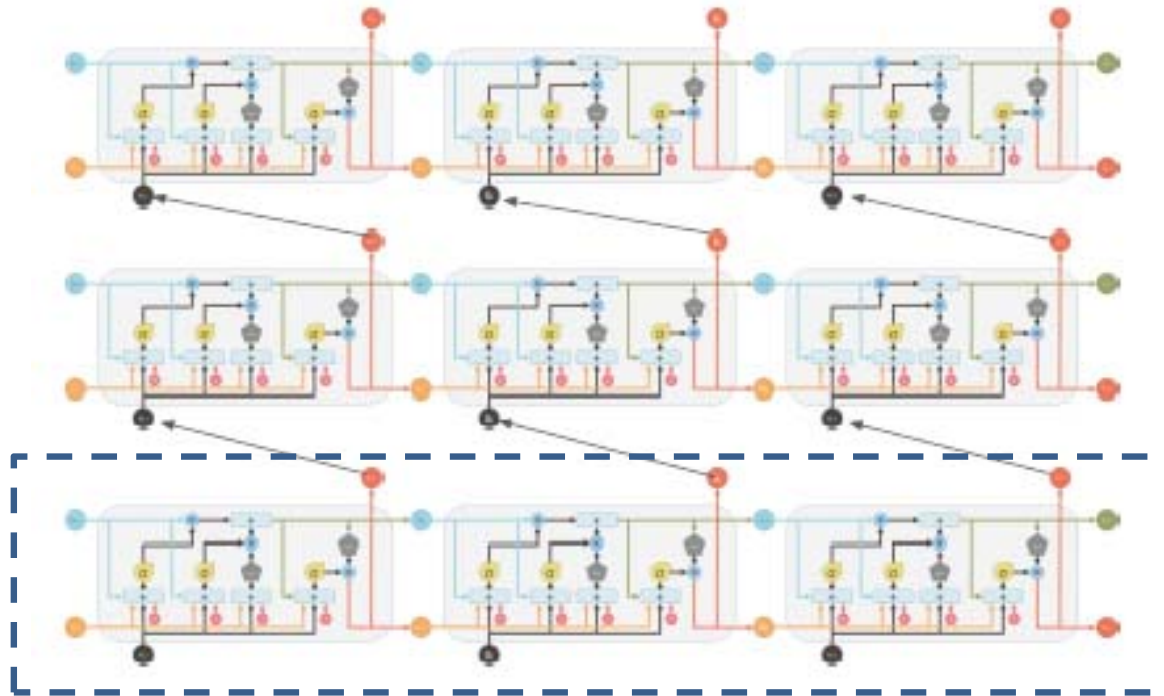
$$C(t) = f(t) \odot C(t-1) + i(t) \odot \tilde{C}(t)$$

$$h(t) = o(t) \odot \tanh(C(t))$$

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Deep LSTM

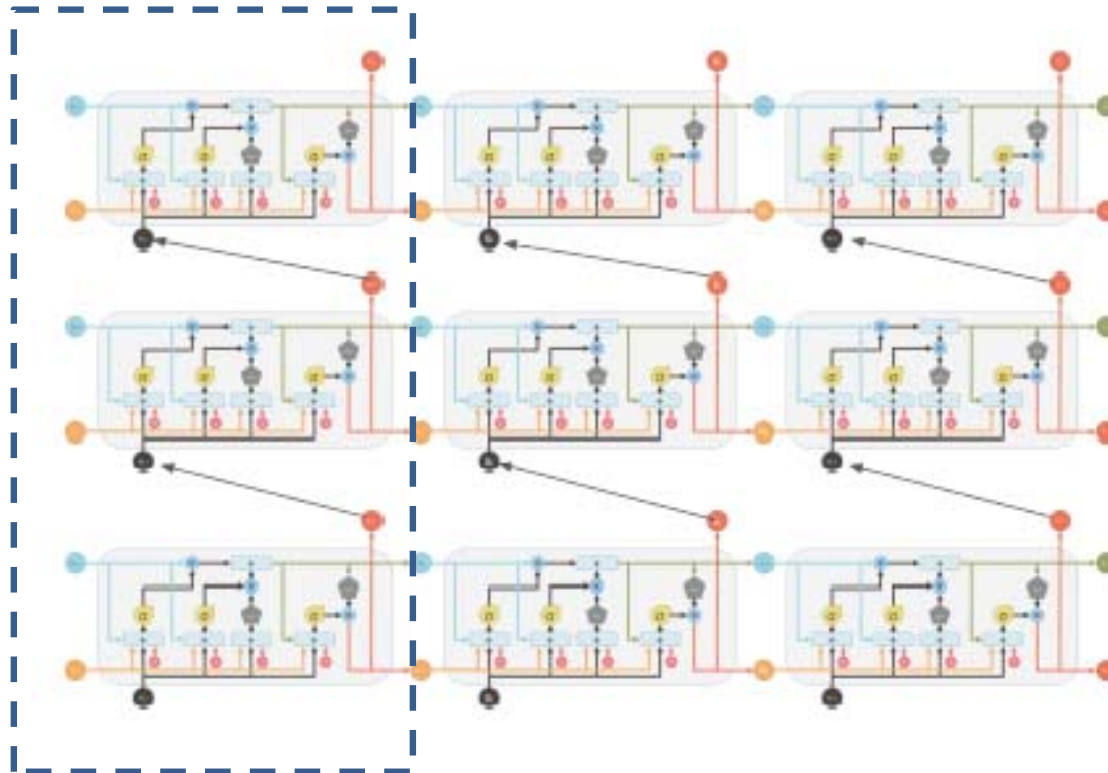
temporal unfolding



- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Deep LSTM

deep stacking  
output sequence of  
one layer constitutes  
the input sequence  
to another layer



## Why do we go deep?

- has the potential to perform better at handling temporal information at wide varying scales
- requires however many more parameters to be learnt

adapted from A. Sood

- Temporal processing with feedforward NNs
- Recurrent architectures for sequence modelling
- Backpropagation through time
- **ESN and LSTM**

# Video example of text understanding with LSTM

<https://youtu.be/mLxsbWAYlpw>