# Recurrent Neural Networks
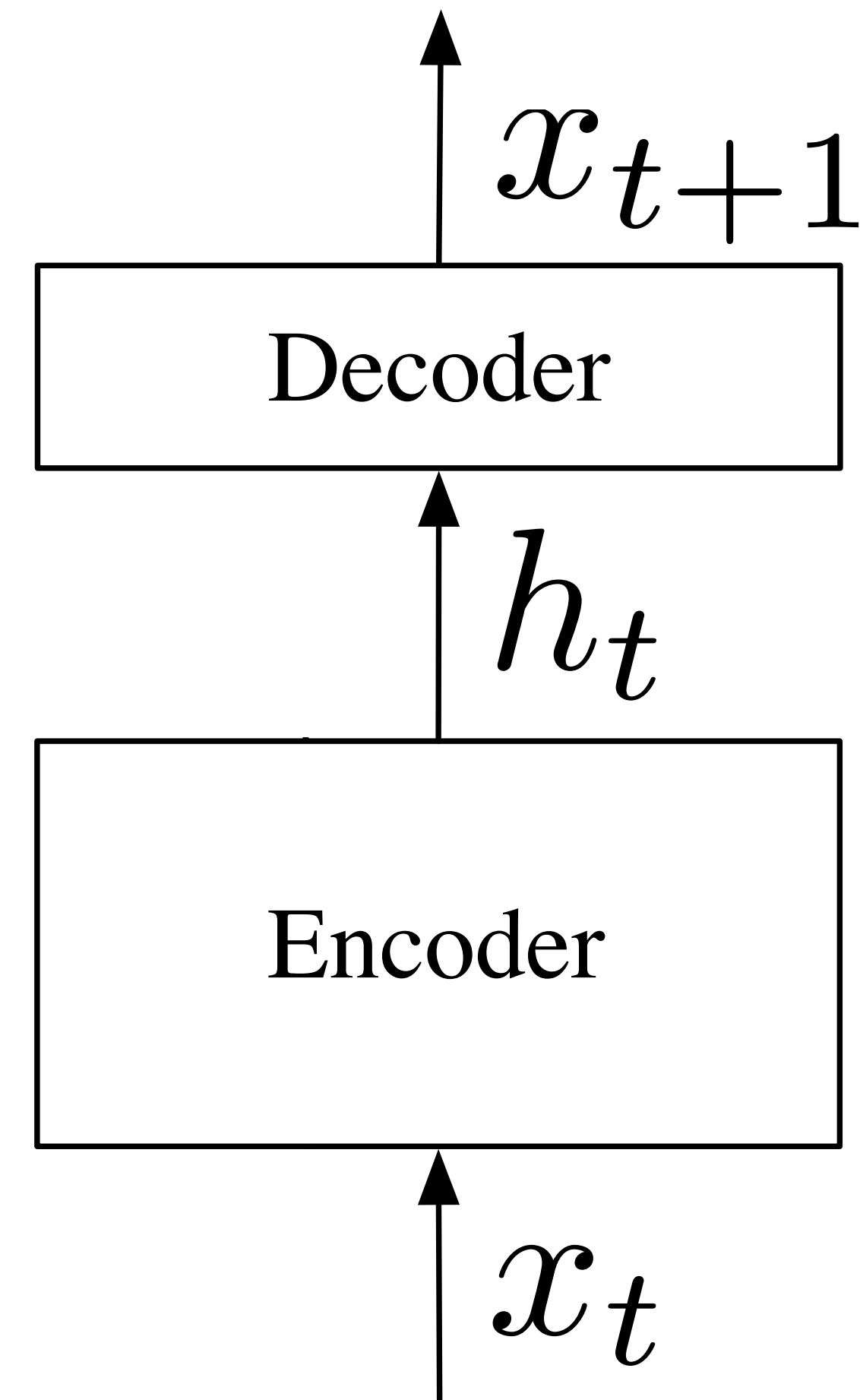
Elman Network
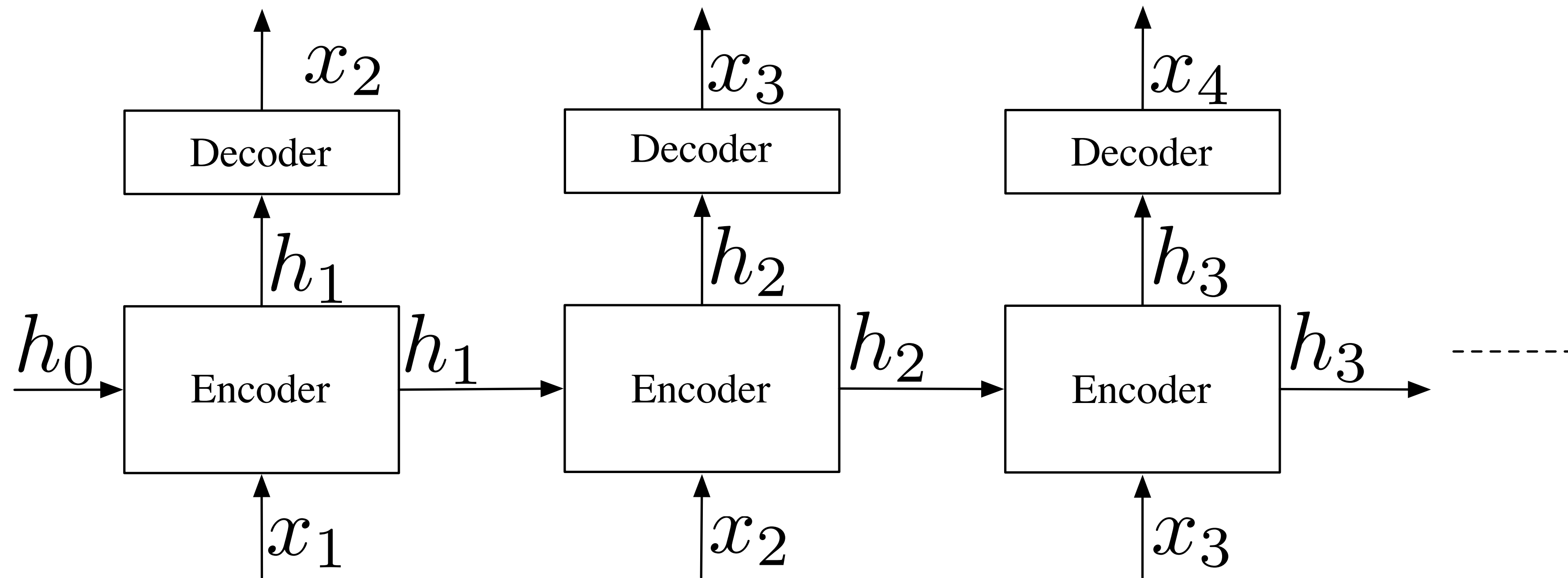
$$h_t = \sigma(W_1 x_t + W_2 h_{t-1})$$

$$x_{t+1} = \rho(W_3 h_t)$$

Can be viewed as a non-linear IIR Filter
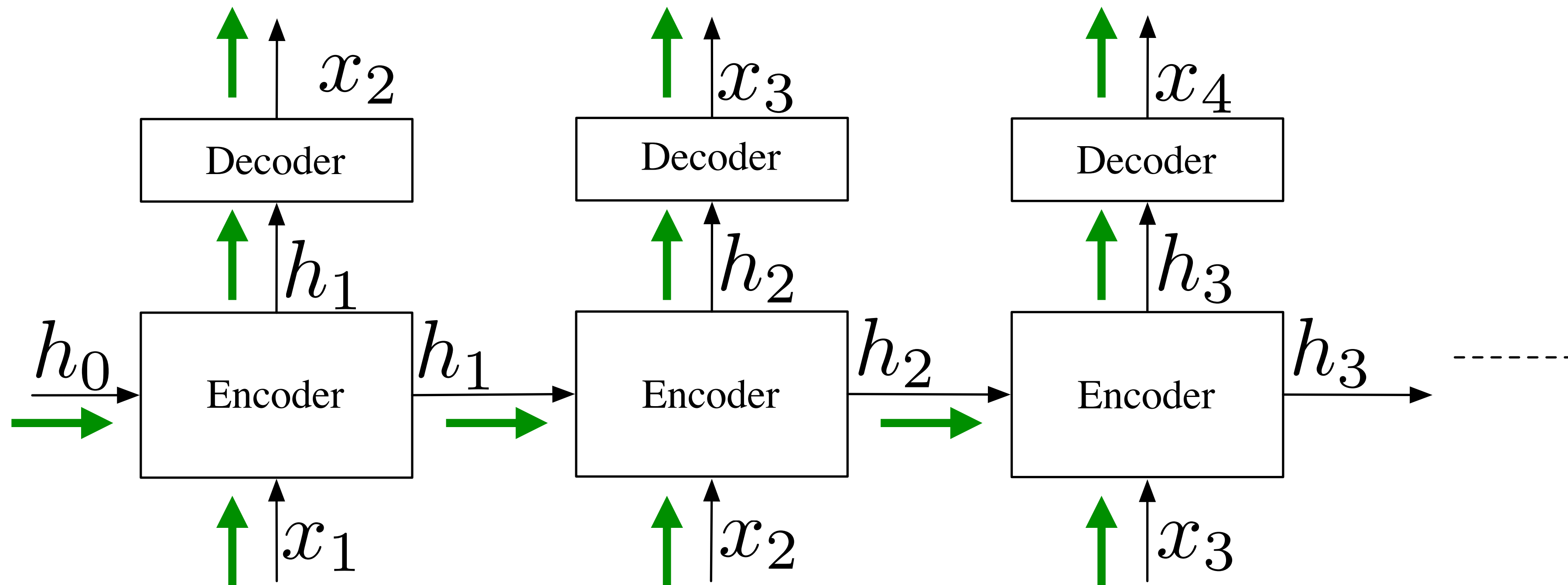
# Elman Network



$$h_t = \sigma(W_1 x_t + W_2 h_{t-1})$$
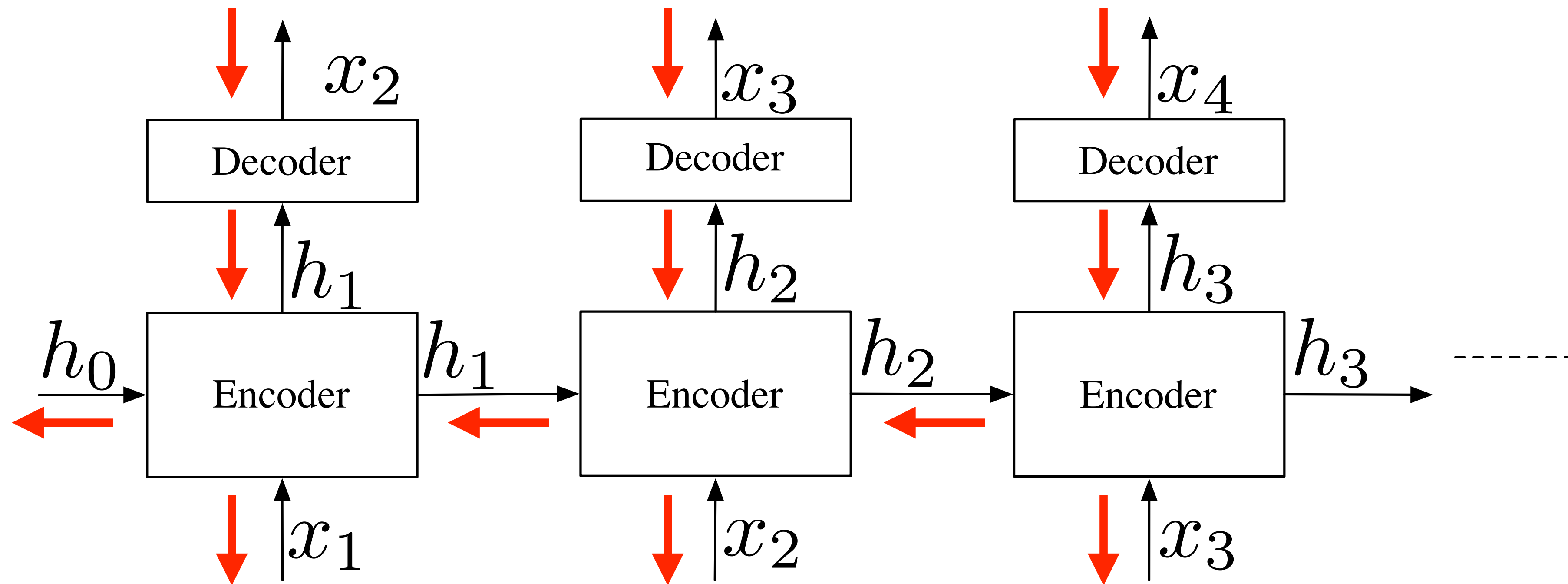
$$x_{t+1} = \rho(W_3 h_t)$$

# Training: Back Propagation Through Time



Forward Pass

Pass inputs through the unrolled network

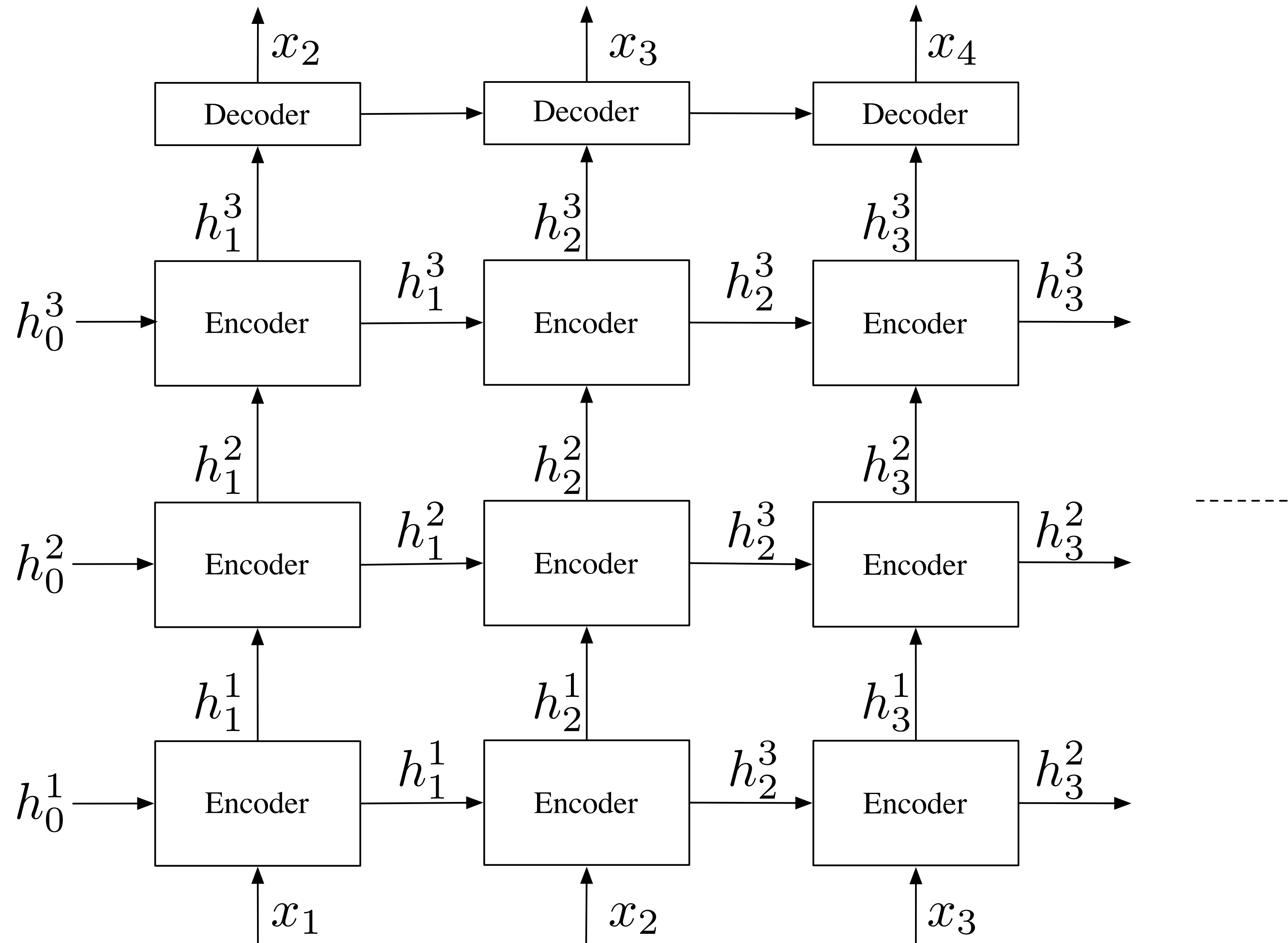compute the error at every time step

# Training: Back Propagation Through Time



Backward Pass

Compute the contribution of error at every time step

accumulate this gradient while going back in time

update the parameters

# Deep RNNs

# Shortcoming of Elman Nets

Exploding gradients

Vanishing gradients

Unable to capture long-term dependencies

Training is somewhat brittle

# Long Short-Term Memory (LSTM)
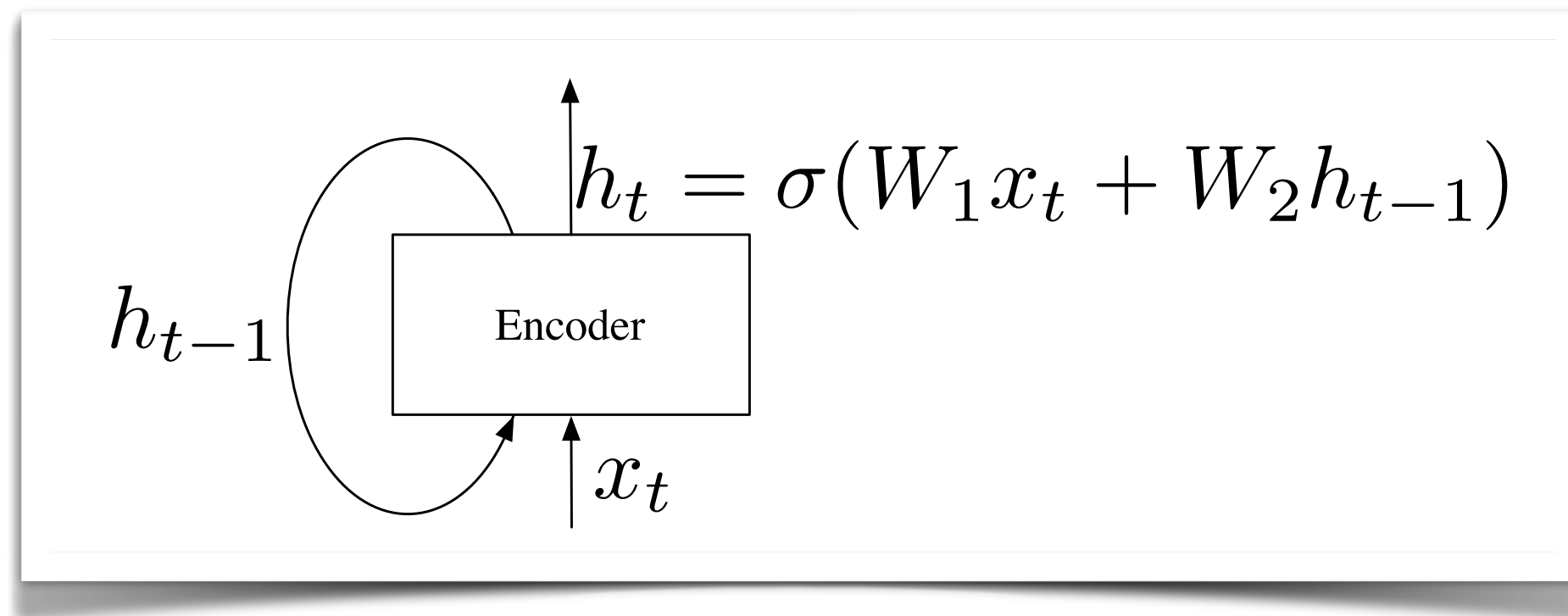
recently gained a lot of popularity

have explicit memory "cells" to store short-term activations

the presence of additional gates partly alleviates the vanishing gradient problem

multi-layer versions shown to work quite well on tasks which have "medium term" dependencies
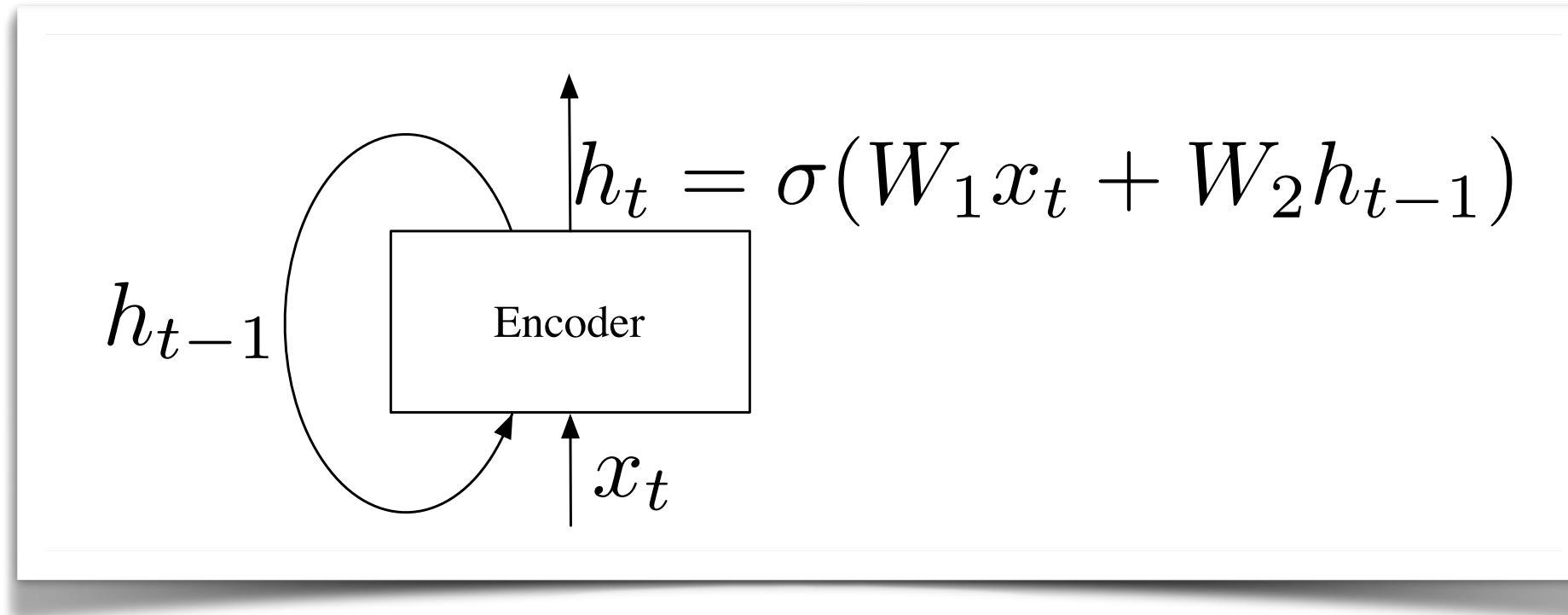
Hochreiter et.al., 1997: Long Short-Term Memory

# Long Short-Term Memory (LSTM)



$$h_t = \sigma(W_1 x_t + W_2 h_{t-1})$$

with $h_{t-1}$, Encoder, $x_t$

Hochreiter et.al., 1997: Long Short-Term Memory

# Long Short-Term Memory (LSTM)

$$h_t = \sigma(W_1 x_t + W_2 h_{t-1})$$

$h_{t-1}$   Encoder   $x_t$

$$h_t = c_t \cdot o_t$$

$$o_t$$

$$x_t$$

$$h_{t-1}$$

$$1.0 \quad \text{Cell} \quad c_t = c_{t-1} + g_t \cdot i_t$$

$$g_t = \sigma(W_1 x_t + W_2 h_{t-1})$$
$$i_t = \sigma(W_3 x_t + W_4 h_{t-1})$$
$$o_t = \sigma(W_5 x_t + W_6 h_{t-1})$$

$$i_t$$

$$x_t$$

$$h_{t-1}$$

$$g_t$$

$$h_{t-1} \quad x_t$$

Hochreiter et.al., 1997: Long Short-Term Memory

# Long Short-Term Memory (LSTM)

$$h_t = c_t \cdot o_t$$

$$c_t = f_t \cdot c_{t-1} + g_t \cdot i_t$$

Cell

$o_t$

$x_t$

$h_{t-1}$

$f_t$

$x_t$

$h_{t-1}$

$i_t$

$x_t$

$h_{t-1}$

$g_t$

$h_{t-1}$  $x_t$
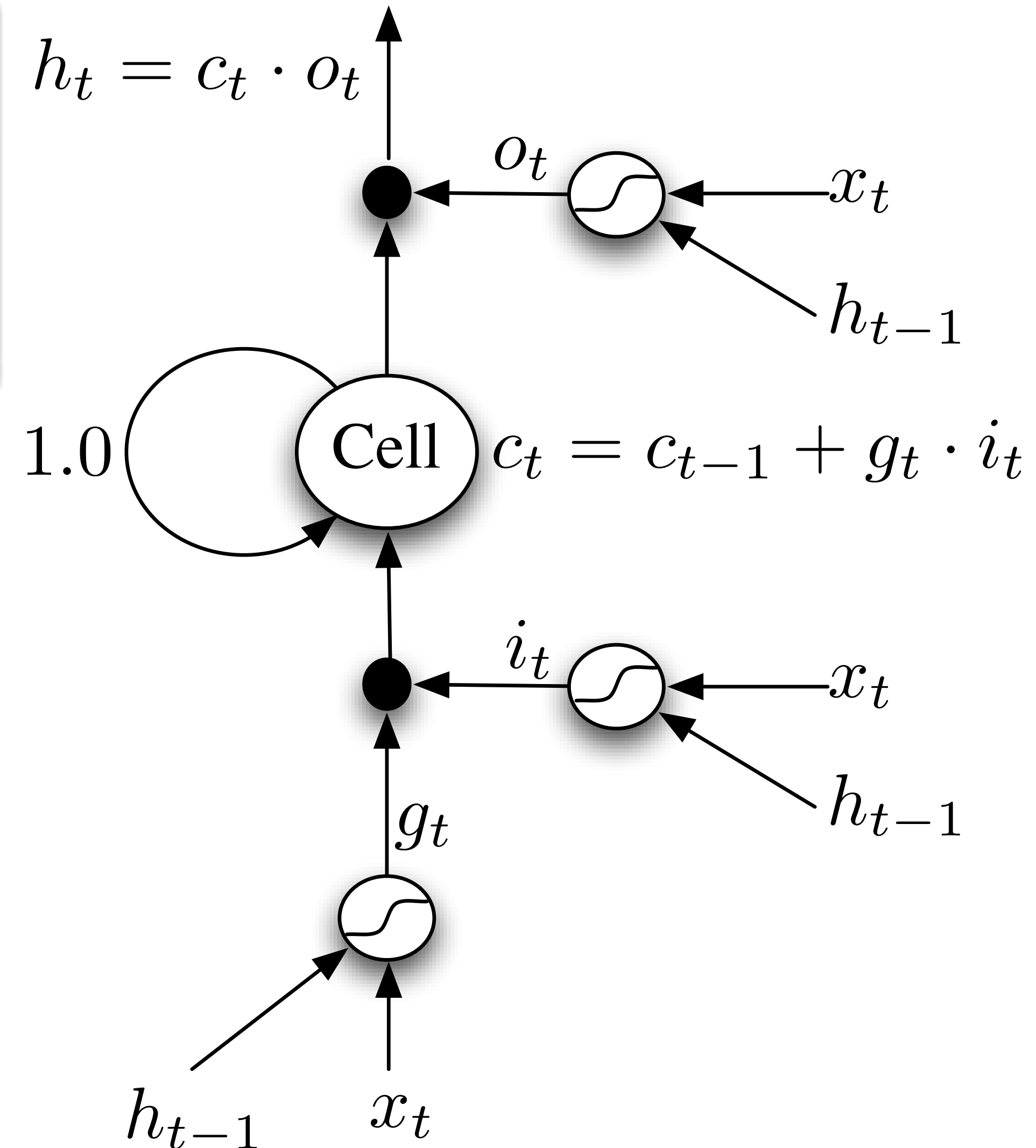
$$g_t = \sigma(W_1 x_t + W_2 h_{t-1})$$
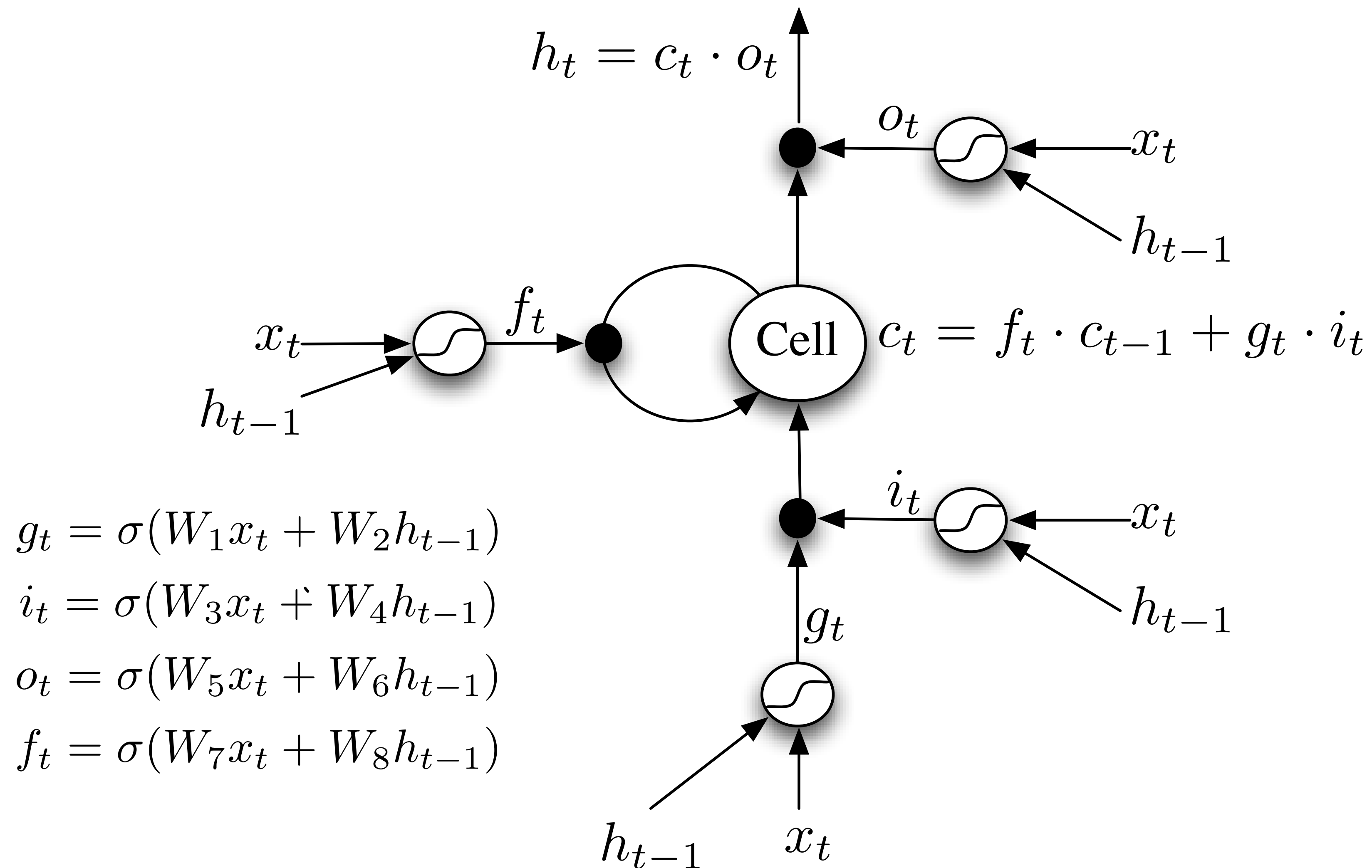
$$i_t = \sigma(W_3 x_t + W_4 h_{t-1})$$

$$o_t = \sigma(W_5 x_t + W_6 h_{t-1})$$

$$f_t = \sigma(W_7 x_t + W_8 h_{t-1})$$
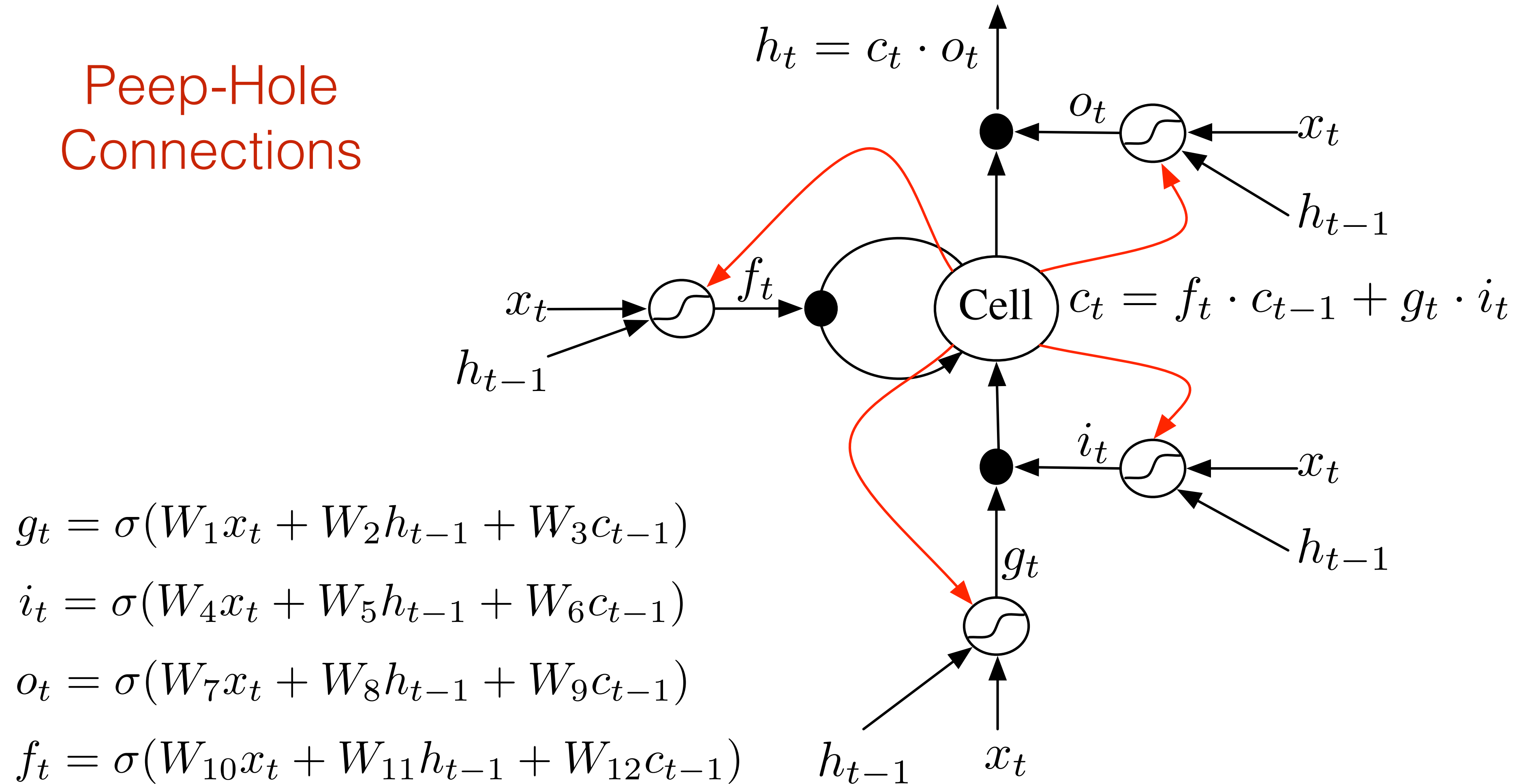
Hochreiter et.al., 1997: Long Short-Term Memory

# Long Short-Term Memory (LSTM)

$$h_t = c_t \cdot o_t$$

$$c_t = f_t \cdot c_{t-1} + g_t \cdot i_t$$

$$g_t = \sigma(W_1 x_t + W_2 h_{t-1} + W_3 c_{t-1})$$

$$i_t = \sigma(W_4 x_t + W_5 h_{t-1} + W_6 c_{t-1})$$

$$o_t = \sigma(W_7 x_t + W_8 h_{t-1} + W_9 c_{t-1})$$

$$f_t = \sigma(W_{10} x_t + W_{11} h_{t-1} + W_{12} c_{t-1})$$

Hochreiter et.al., 1997: Long Short-Term Memory

# Long Short-Term Memory (LSTM)



$$h_t = c_t \cdot o_t$$

$$c_t = f_t \cdot c_{t-1} + g_t \cdot i_t$$

**Encoder**

Hochreiter et.al., 1997: Long Short-Term Memory

# Long Short-Term Memory (LSTM)



$$h_t = c_t \cdot o_t$$

$$c_t = f_t \cdot c_{t-1} + g_t \cdot i_t$$

**Encoder**

Hochreiter et.al., 1997: Long Short-Term Memory