# Artificial Neural Networks and Deep Learning

Week 5
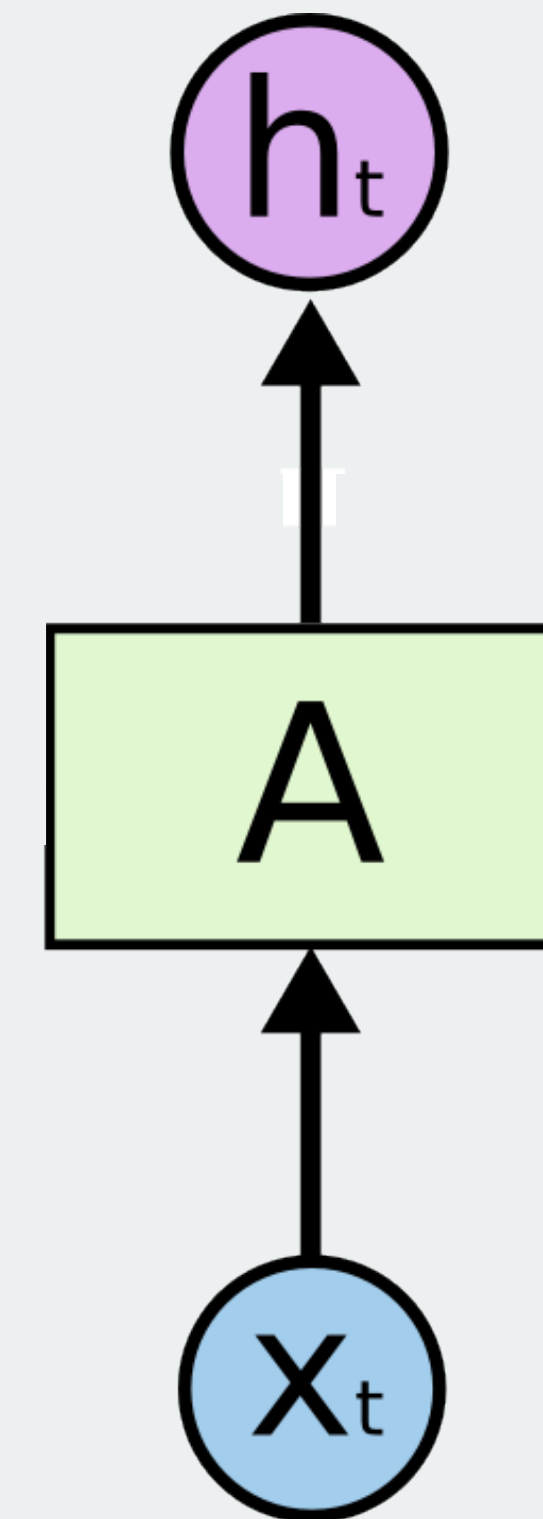
# Recurrent neural networks

# Recurrent Neural Networks

THE neural network architecture to use for sequential data

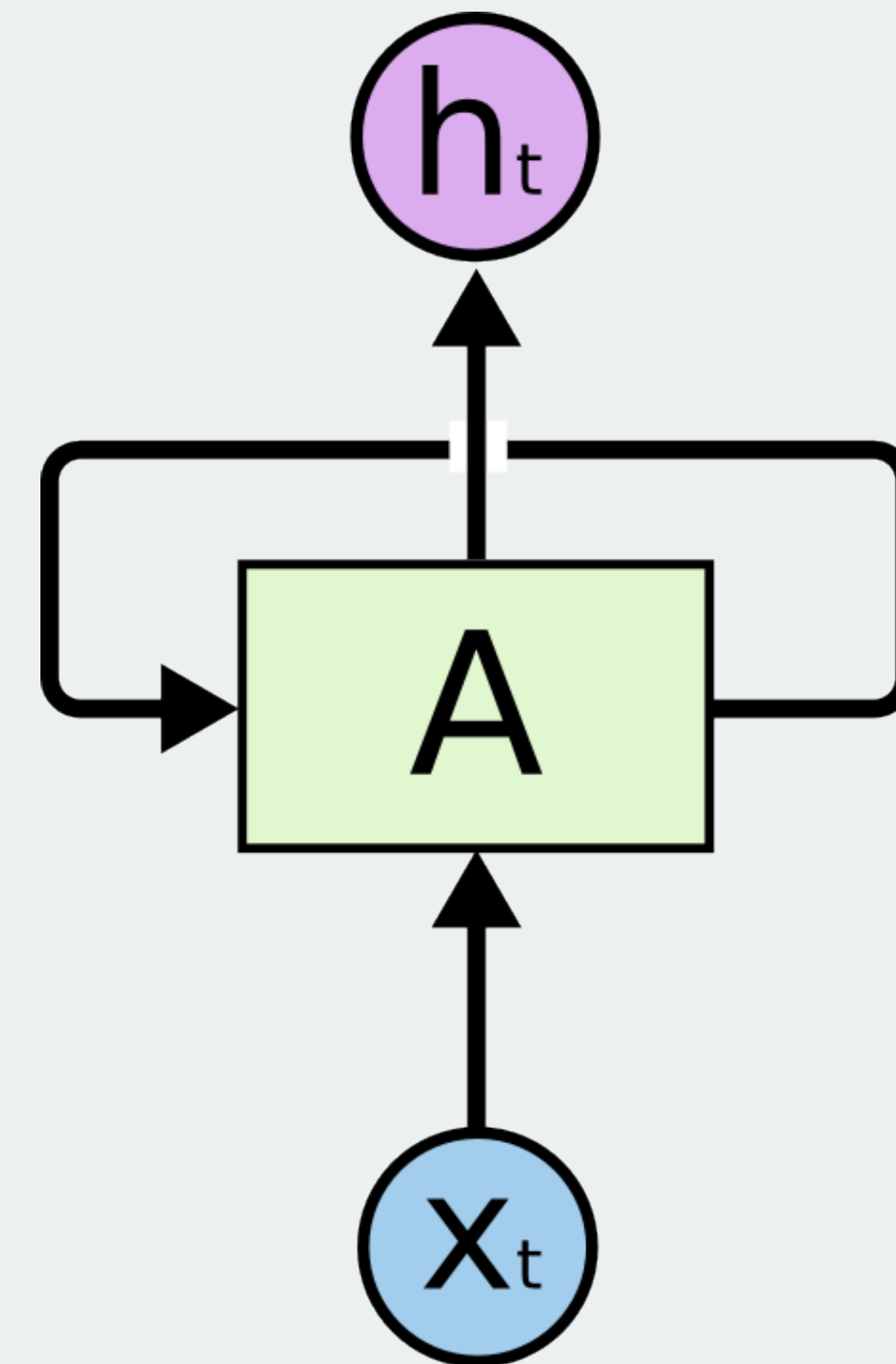# The **problem** with **all** feed forward neural networks

- Input and output must be of hard-assigned dimensions

- The network makes a fixed number of computations

- If input is sequence-like (video, sound, etc.) the network is ignorant to the order of samples

$h_t$

$A$

$X_t$

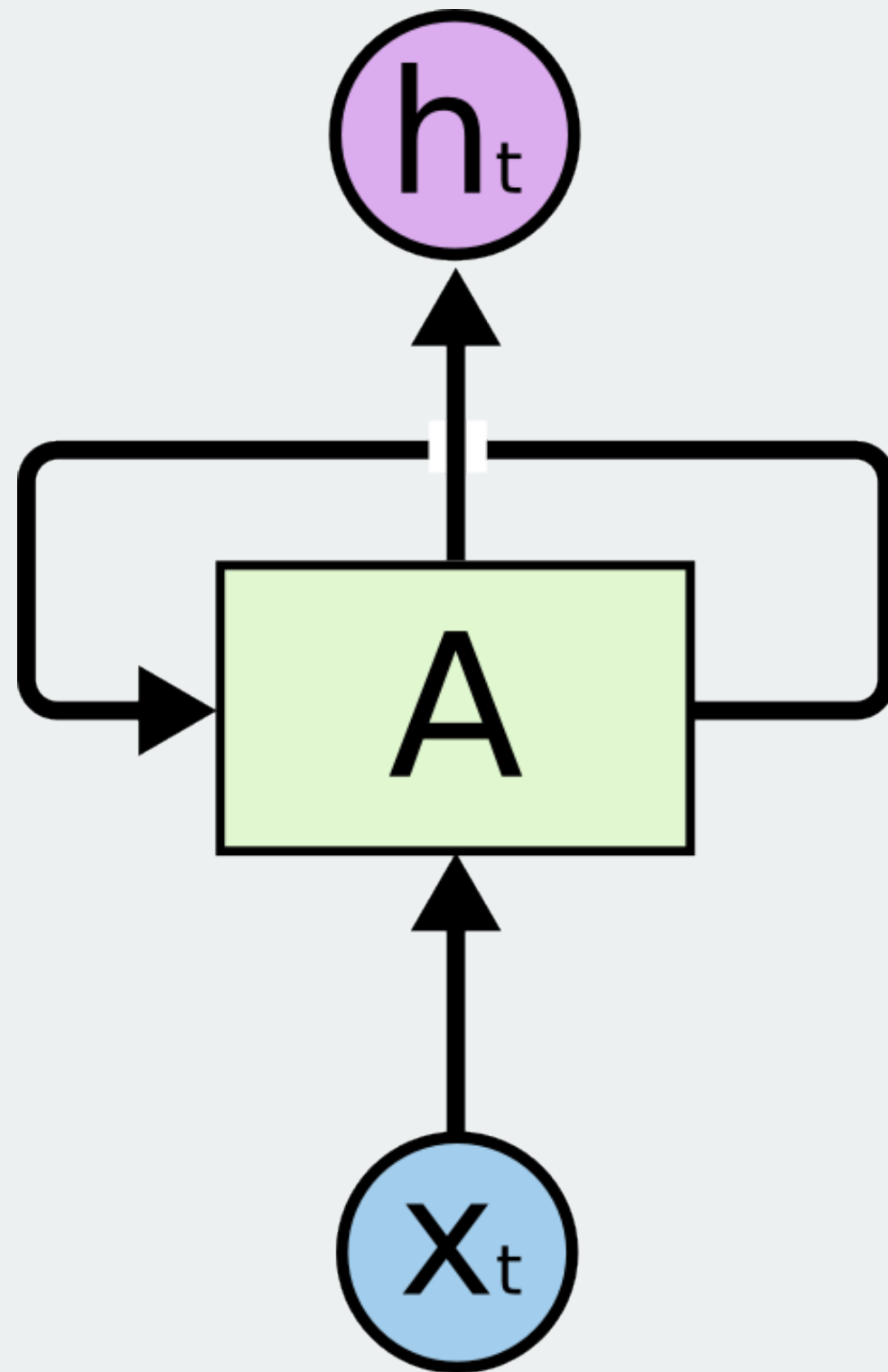# The **problem** with **all** feed forward neural networks
> **Solution:** *Recurrence!*

- Input and output must be of hard-assigned dimensions

- The network makes a fixed number of computations

- If input is sequence-like (video, sound, etc.) the network is ignorant to the order of samples
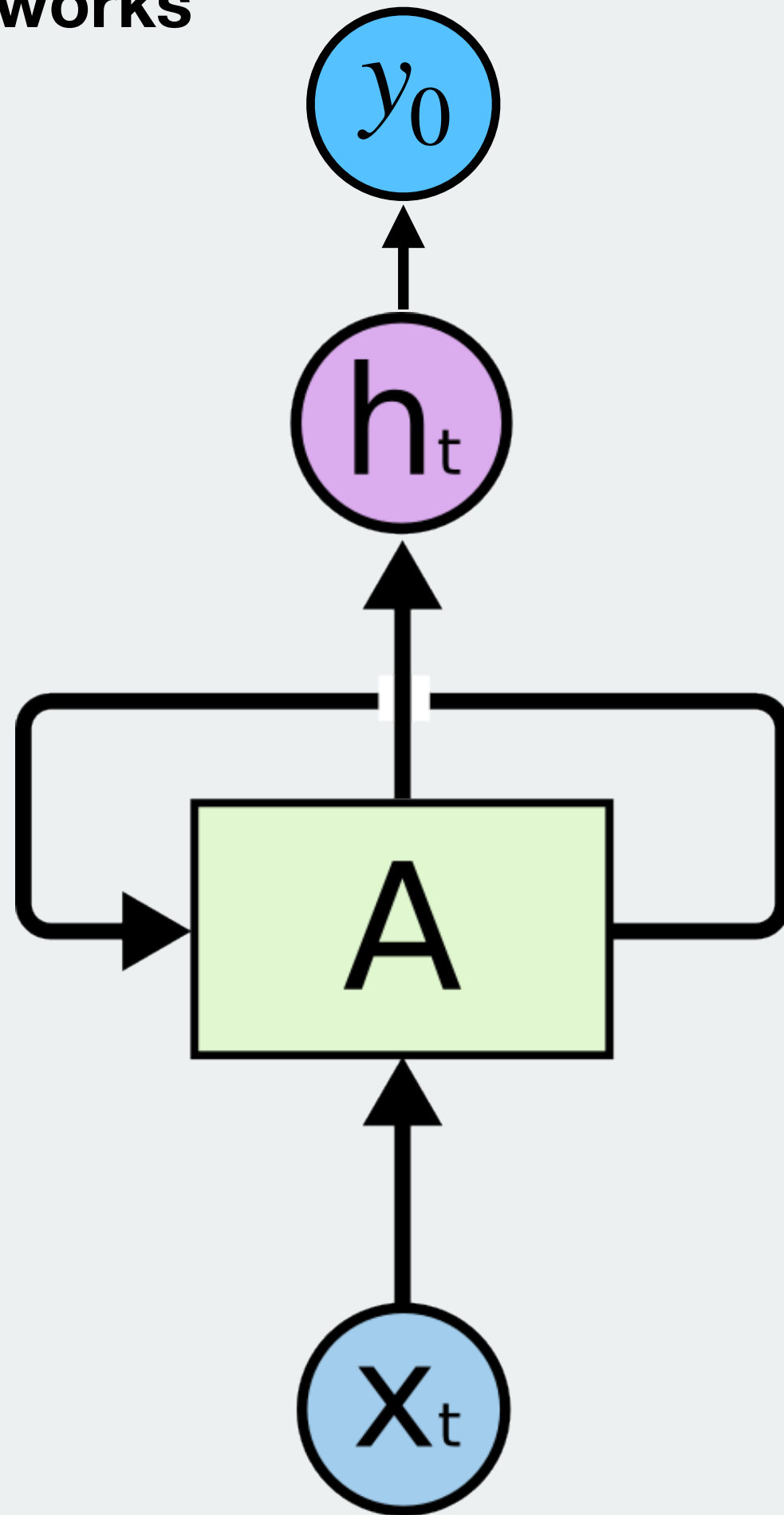
**Recurrent neural networks**

> Fundamental idea



$$h_t = f_W(h_{t-1}, x_t)$$

**Notice:** W is the same in each iteration

**Recurrent neural networks**
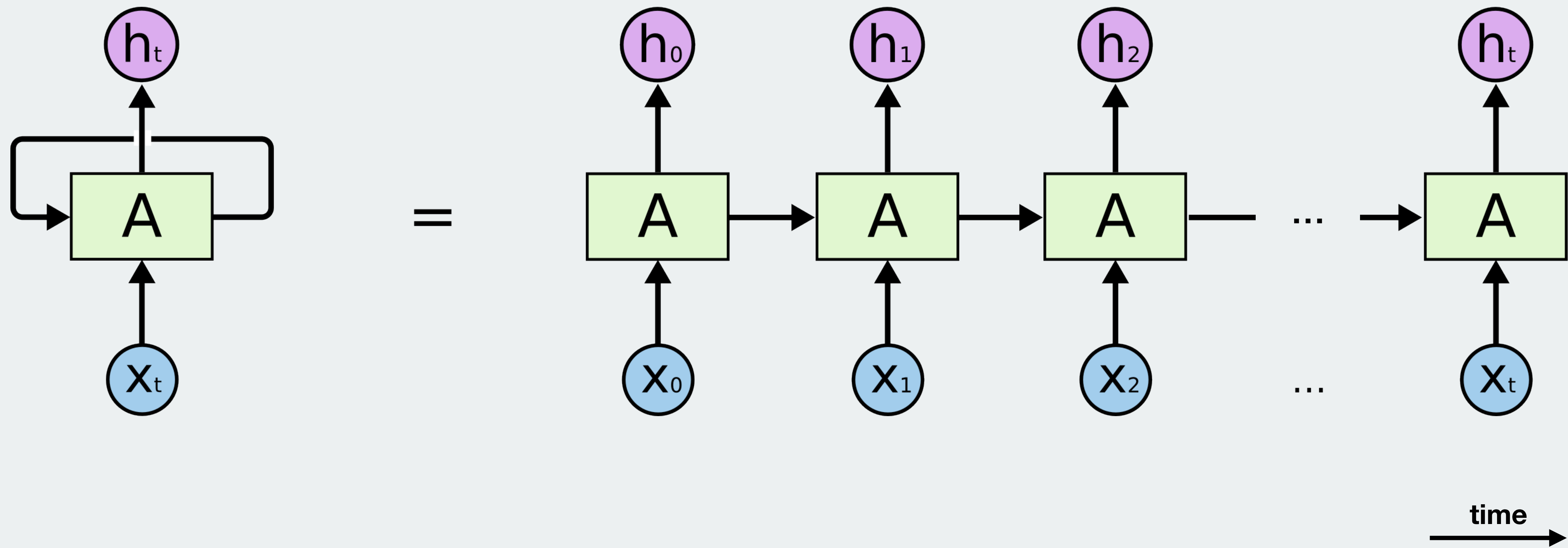> Fundamental idea

$y_0$

$h_t$

A

$x_t$

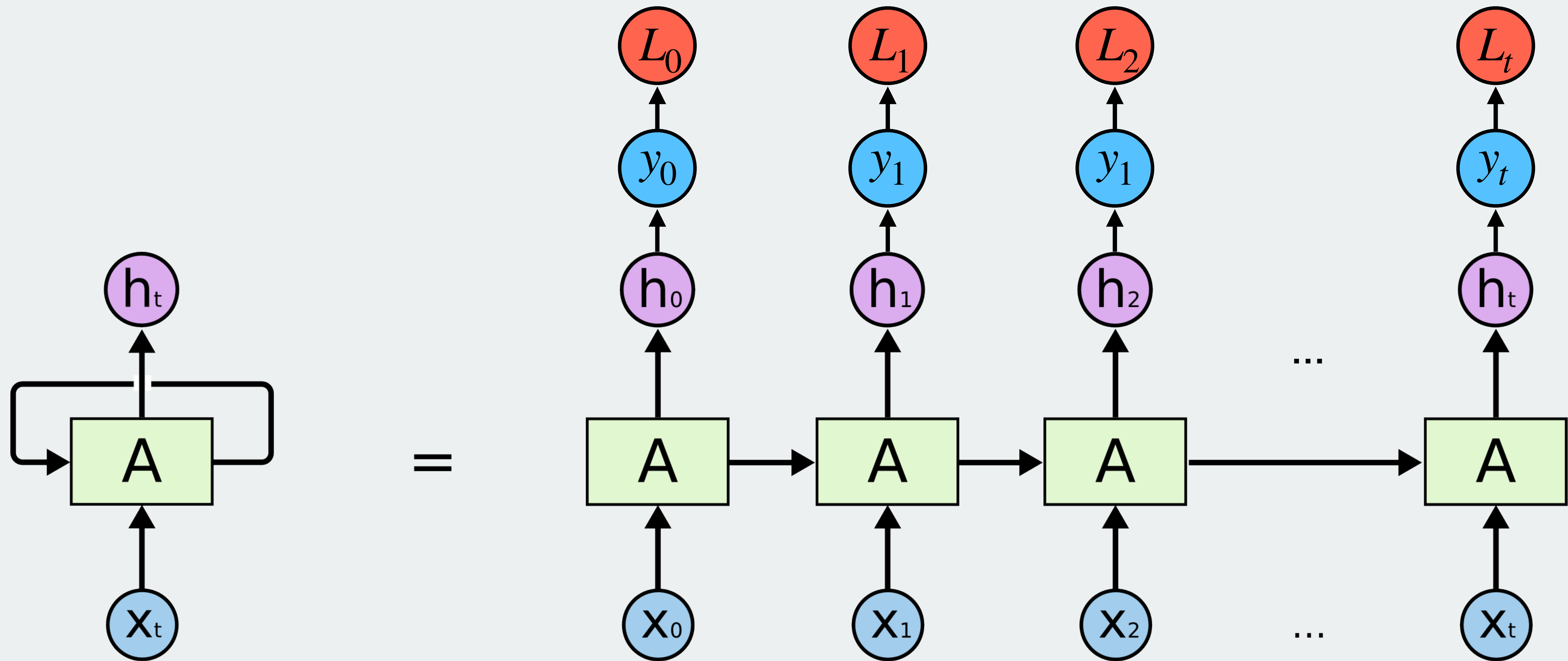$$h_t = f_W(h_{t-1}, x_t)$$

$$y_t = W_{hy}h_t$$

# Recurrent neural networks

> Unrolled in time

**Recurrent neural networks**

> Backpropagation
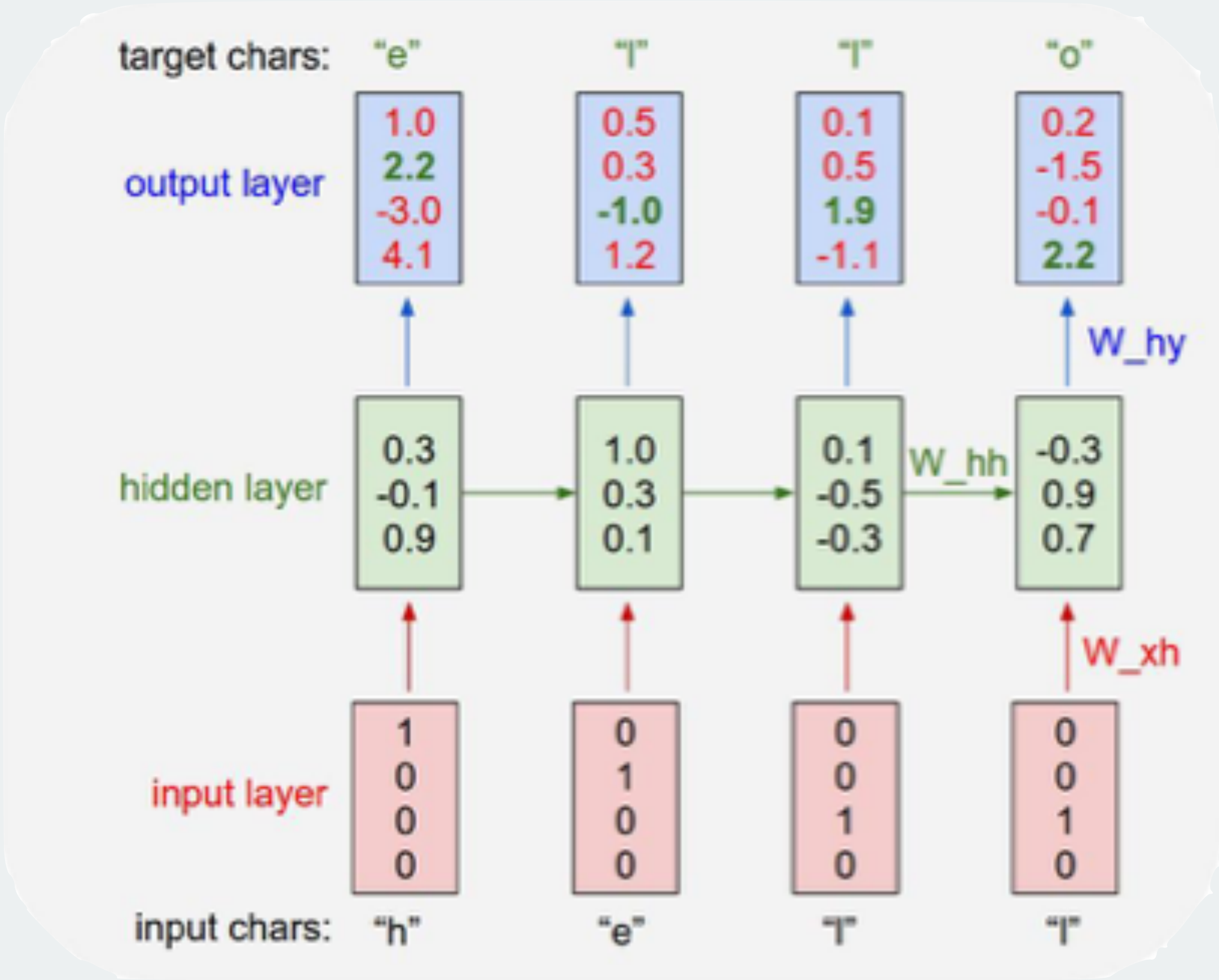
# Recurrent neural networks

> Ways to process sequential data



one to one     one to many     many to one     many to many     many to many
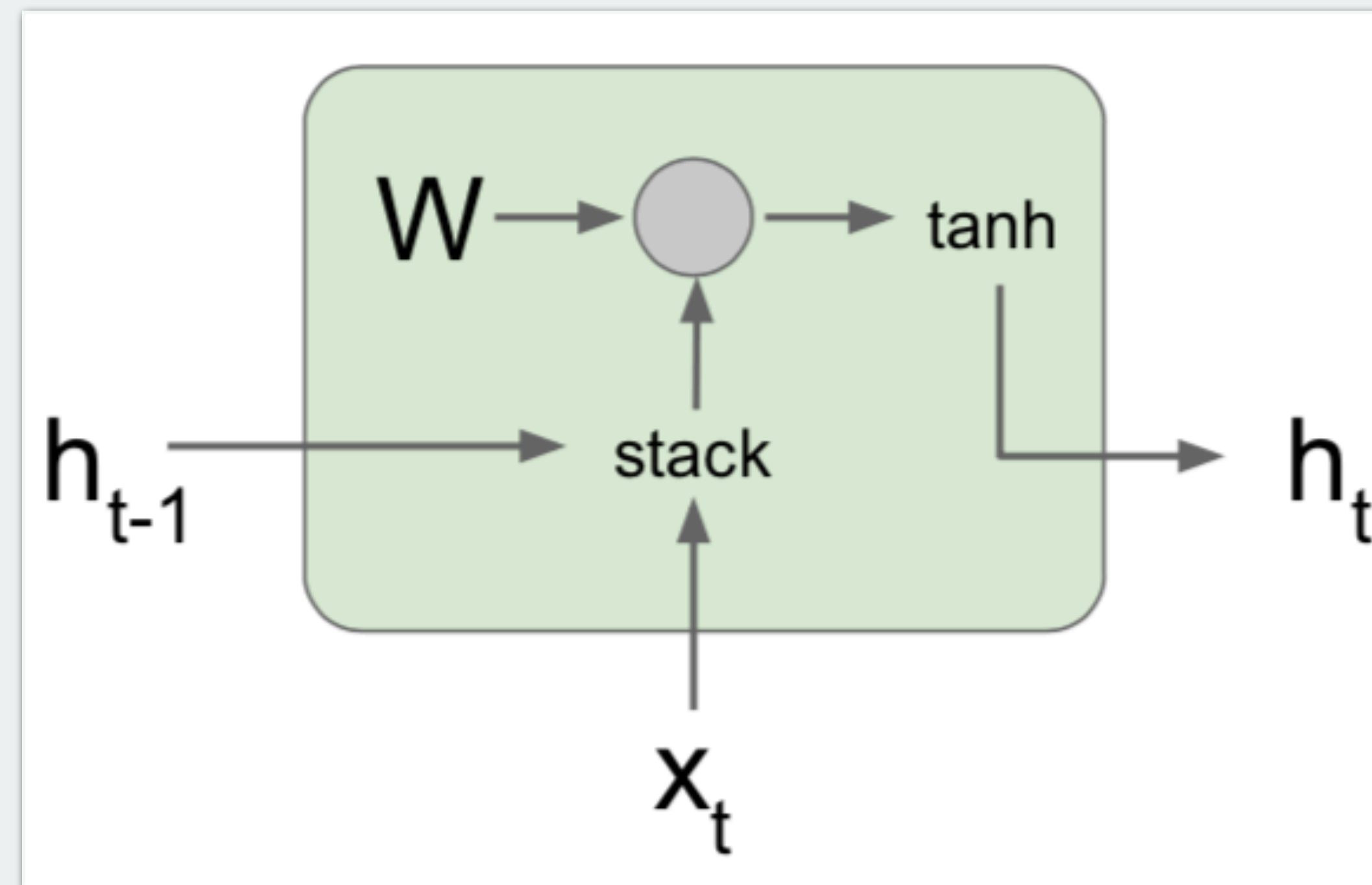
# Recurrent neural networks

> Example: predicting next character

training sequence**: "hello"**
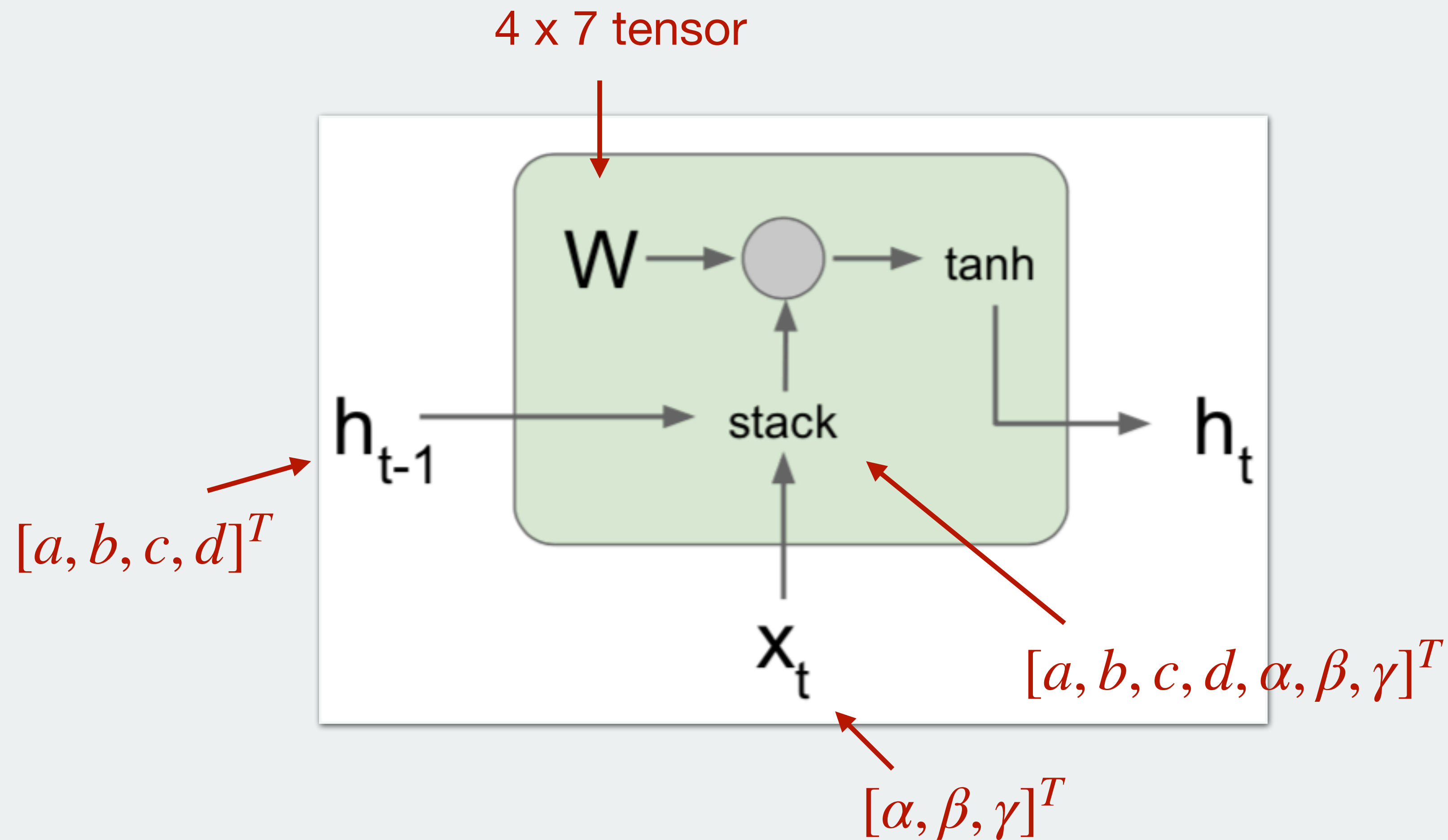
# Recurrent neural networks
> Vanilla RNN, architecture
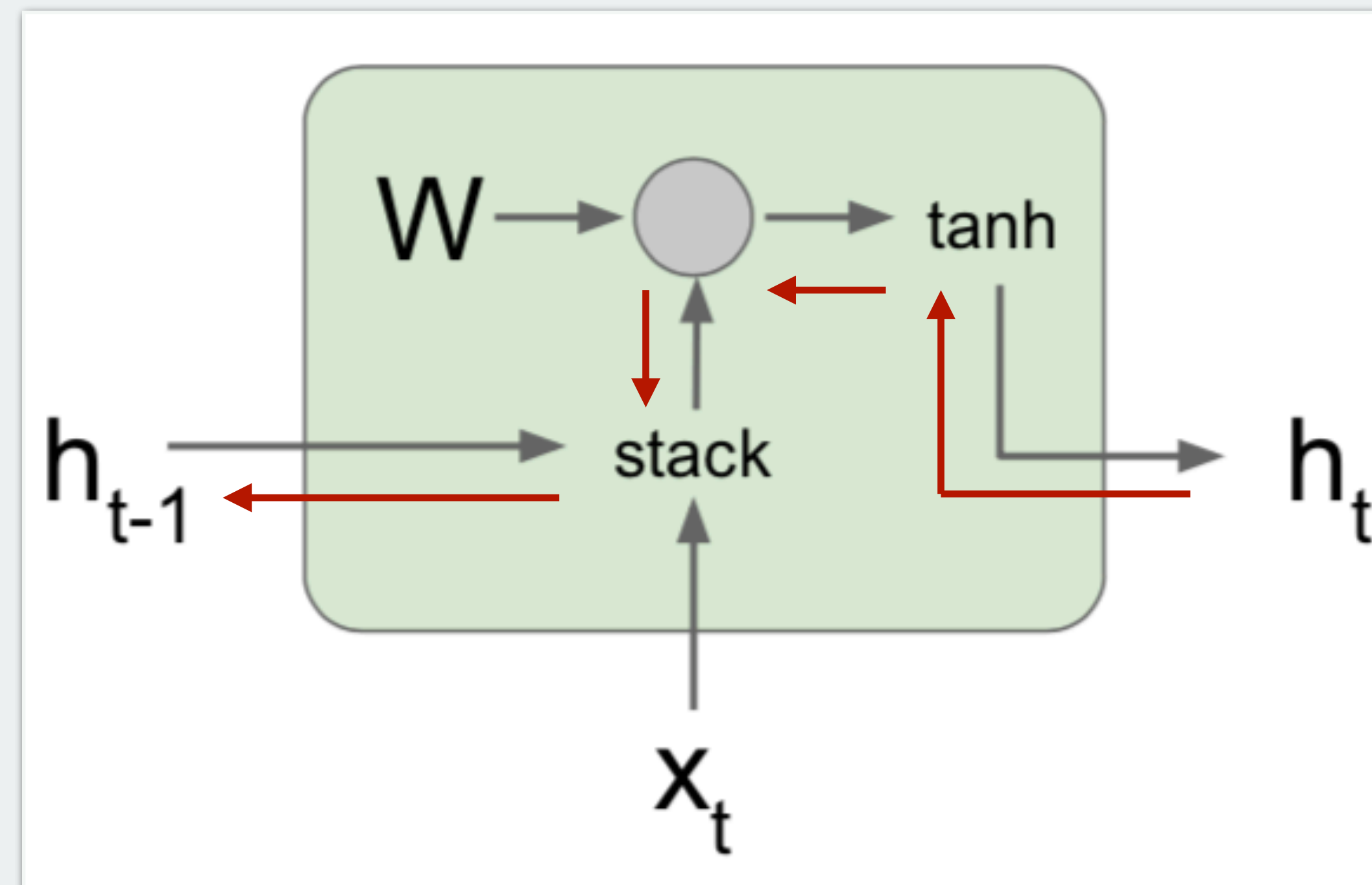
# Recurrent neural networks
> Vanilla RNN, architecture

4: because dotting (h, t) onto it should result in a new vector with 4 elements
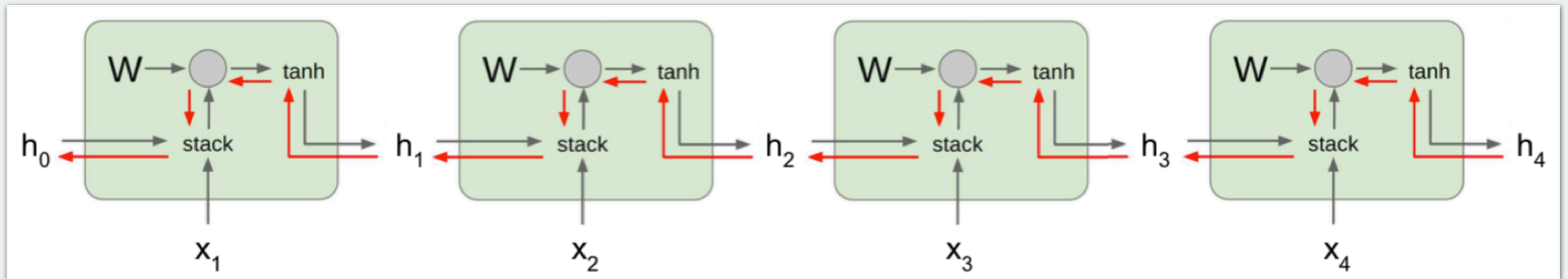7: because the (h, t) vector we are dotting onto it has 7 elements in it

4 x 7 tensor



$[a, b, c, d]^T$

$[a, b, c, d, \alpha, \beta, \gamma]^T$

$[\alpha, \beta, \gamma]^T$

http://cs231n.stanford.edu/syllabus.html

# Recurrent neural networks

> Vanilla RNN, backpropagation

# Recurrent neural networks

> Vanilla RNN, backpropagation



**Problem:** Repeated multiplications by **W** during backpropagation

**Leads to:** Exploding/vanishing gradients

**Solution:** Gradient clipping (solves exploding gradients), or **change architecture**

**Long Short Term Memory (LSTM)**

**Vanilla RNN**

$$h_t = \tanh\left(W \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix}\right)$$

$$\begin{bmatrix} i \\ f \\ o \\ g \end{bmatrix} = \begin{bmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{bmatrix} W \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

**Recurrent neural networks**

> Vanilla RNN vs. LSTM
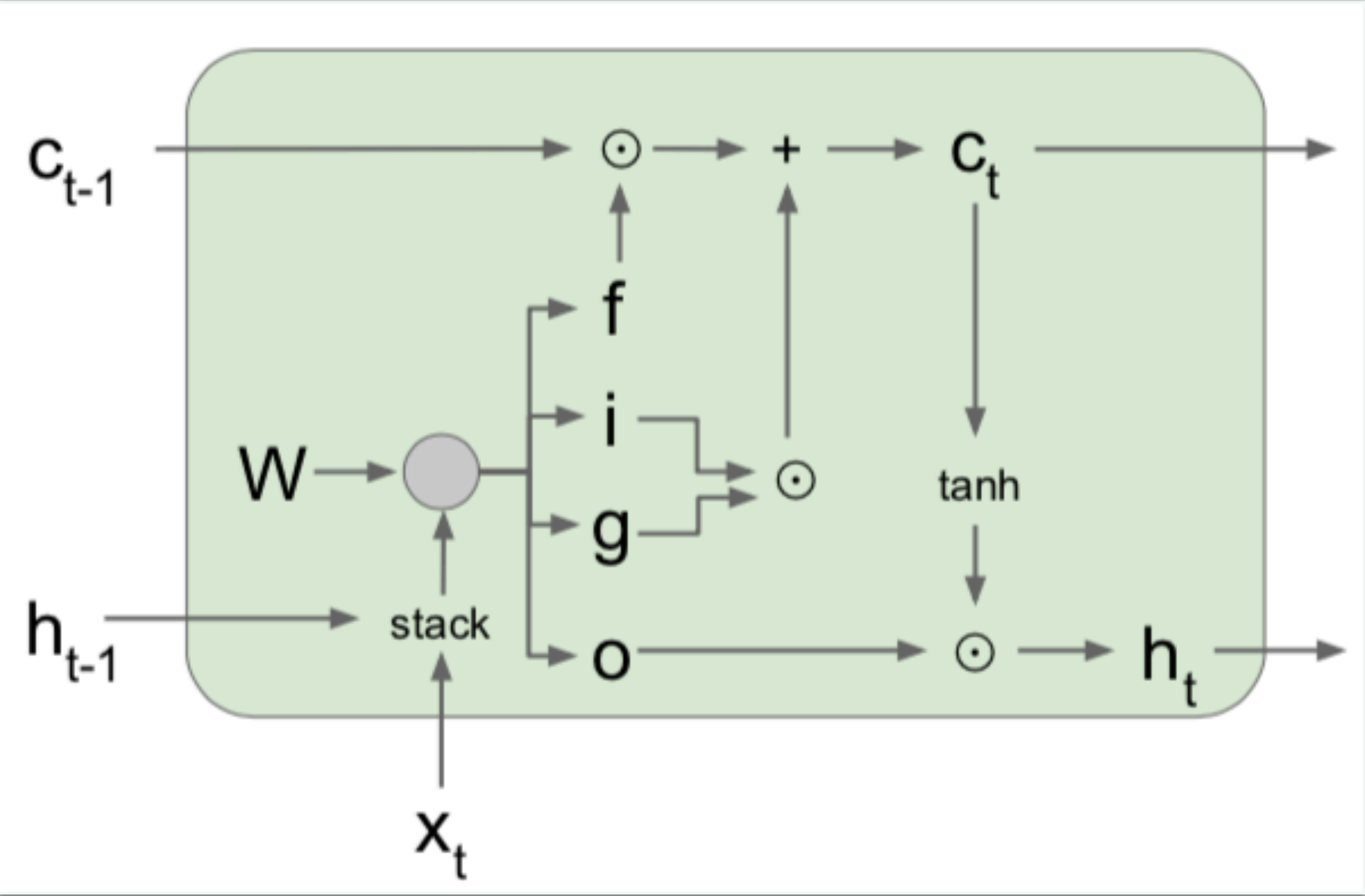
**Long Short Term Memory (LSTM)**



$$\begin{bmatrix} i \\ f \\ o \\ g \end{bmatrix} = \begin{bmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{bmatrix} W \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$
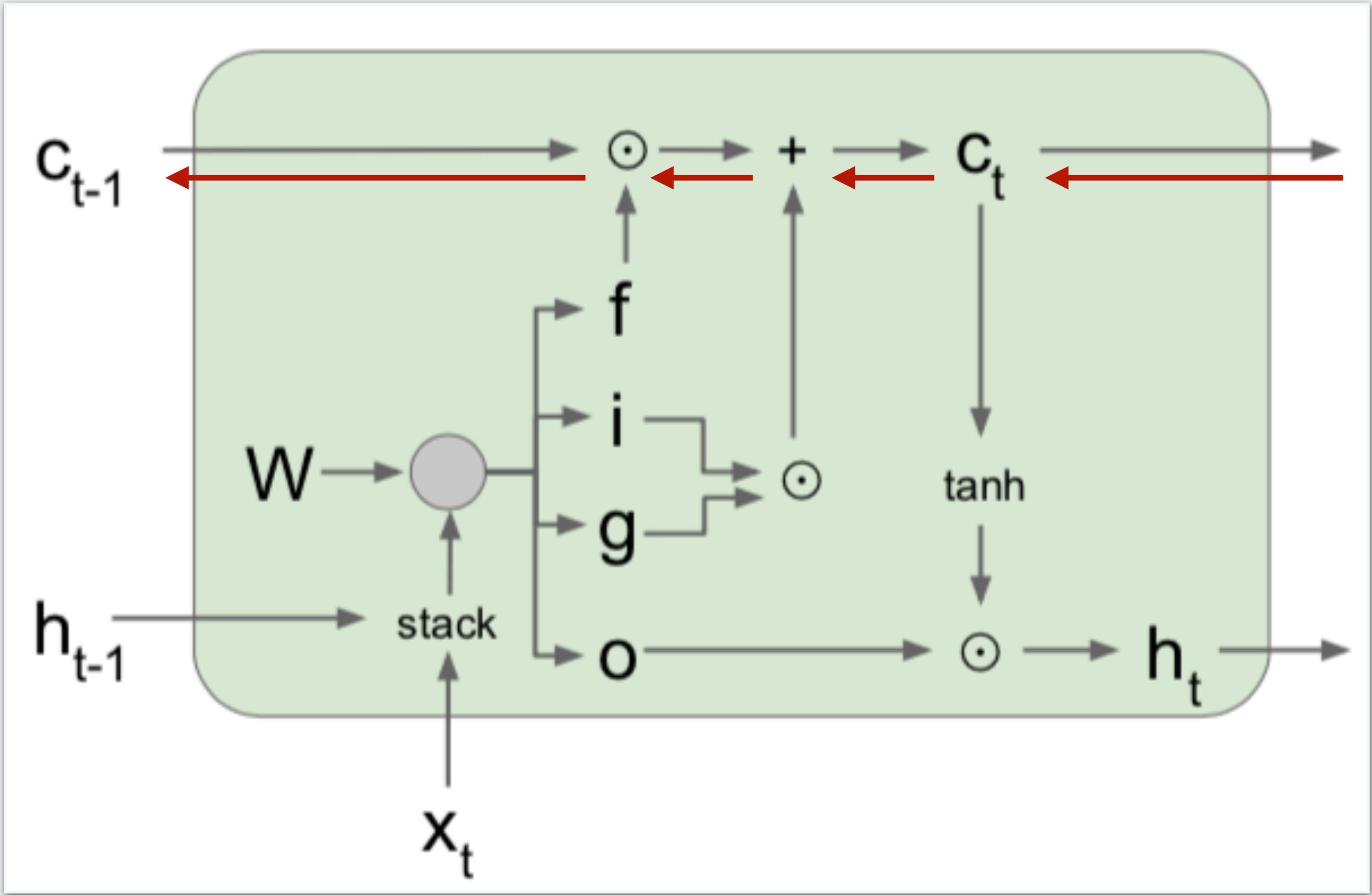
# Recurrent neural networks
> Vanilla RNN vs. LSTM

**"Gradient highway":** solves vanishing/exploding gradient problems

**Long Short Term Memory (LSTM)**



$$
\begin{bmatrix} i \\ f \\ o \\ g \end{bmatrix} = \begin{bmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{bmatrix} W \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix}
$$

$$
c_t = f \odot c_{t-1} + i \odot g
$$

$$
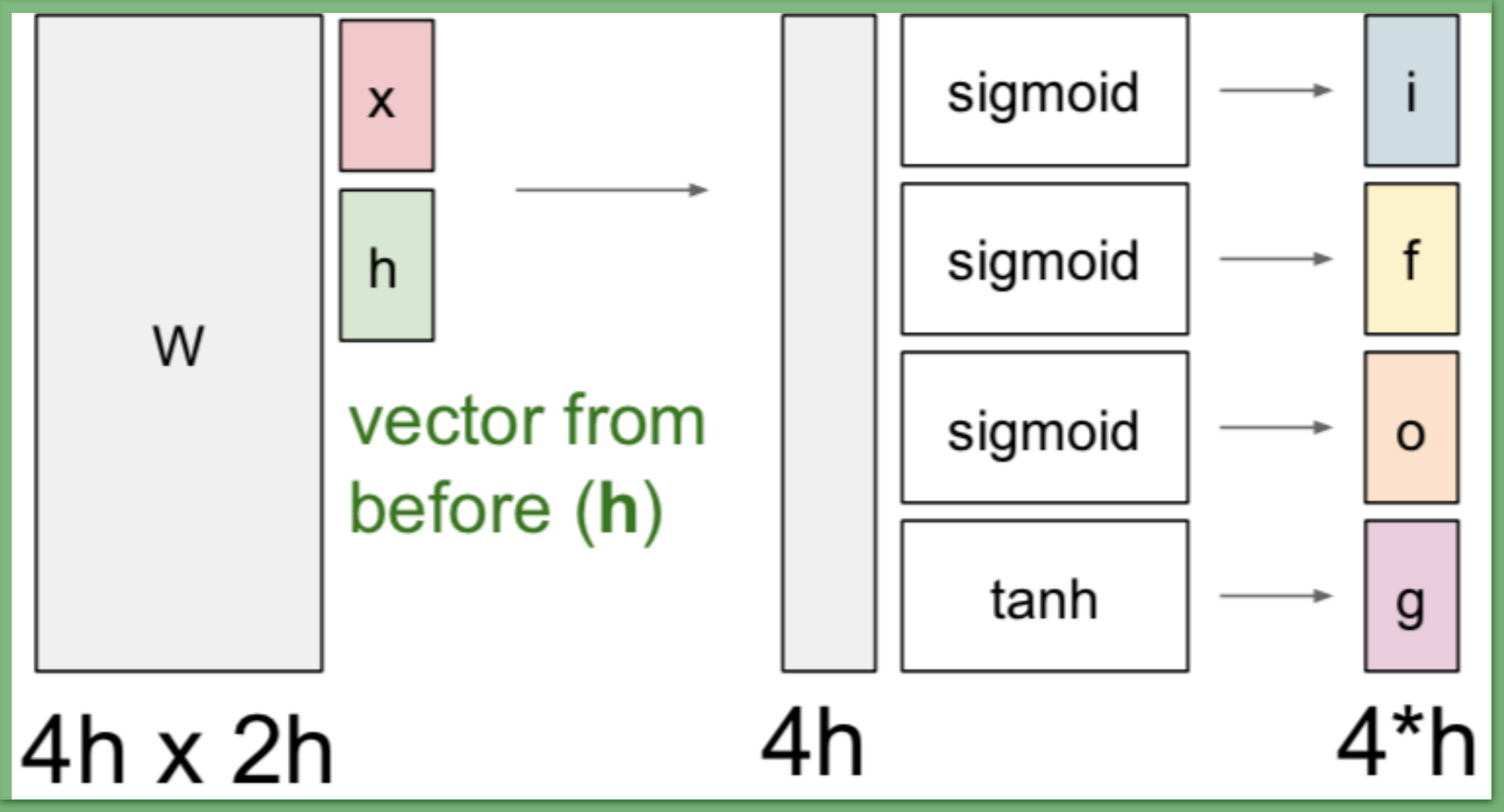h_t = o \odot \tanh(c_t)
$$

**Long Short Term Memory (LSTM)**



$$\begin{bmatrix} i \\ f \\ o \\ g \end{bmatrix} = \begin{bmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{bmatrix} W \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

http://cs231n.stanford.edu/syllabus.html