

Deep Learning (CS324)

6. Recurrent Neural Networks*

Jianguo Zhang
SUSTech

*Based on <http://www.deeplearningbook.org> chapter 10 + sources mentioned in slides

Learning over sequences

- A lot of data naturally comes in the form of sequences

Videos



Source: <https://www.eltima.com/images/upload/elmedia/articles/frame/sample.jpg>

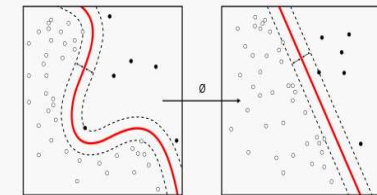
Text (and speech)

A **recurrent neural network (RNN)** is a class of [artificial neural networks](#) where connections between nodes form a [directed graph](#) along a temporal sequence. This allows it to exhibit temporal dynamic behavior. Unlike [feedforward neural networks](#), RNNs can use their internal state (memory) to process sequences of inputs. This makes them applicable to tasks such as unsegmented, connected [handwriting recognition](#)^[1] or [speech recognition](#).^{[2][3]}

The term "recurrent neural network" is used indiscriminately to refer to two broad classes of networks with a similar general structure, where one is [finite impulse](#) and the other is [infinite impulse](#). Both classes of networks exhibit temporal [dynamic behavior](#).^[4] A finite impulse recurrent network is a [directed acyclic graph](#) that can be unrolled and replaced with a strictly feedforward neural network, while an infinite impulse recurrent network is a [directed cyclic graph](#) that can not be unrolled.

Both finite impulse and infinite impulse recurrent networks can have additional stored state, and the storage can be under direct control by the neural network. The storage can also be replaced by another network or graph, if that incorporates time delays or has feedback loops. Such controlled states are referred to as gated state or gated memory, and are part of [long short-term memory networks \(LSTMs\)](#) and [gated recurrent units](#). This is also called Feedback Neural Network.

Machine learning and data mining



Problems

[\[show\]](#)

Supervised learning

(classification • regression)

[\[show\]](#)

Clustering

[\[show\]](#)

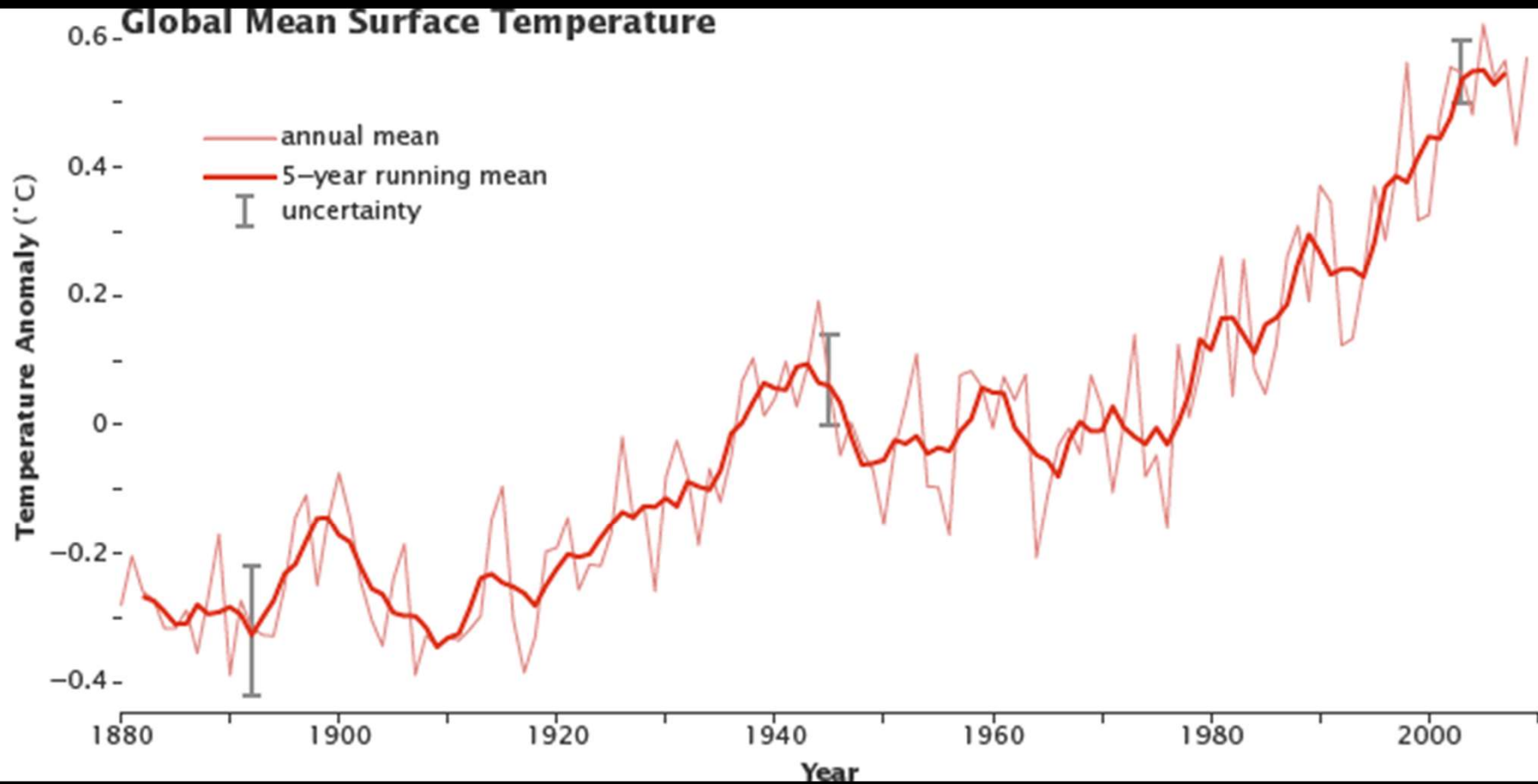
Dimensionality reduction

[\[show\]](#)

Structured prediction

[\[show\]](#)

Climate measurements



Music

Ode To Joy An die Freude

Transcribed for Clarinet and Piano

L. van Beethoven

Allegro assai vivace alla Marcia ♩. = 84

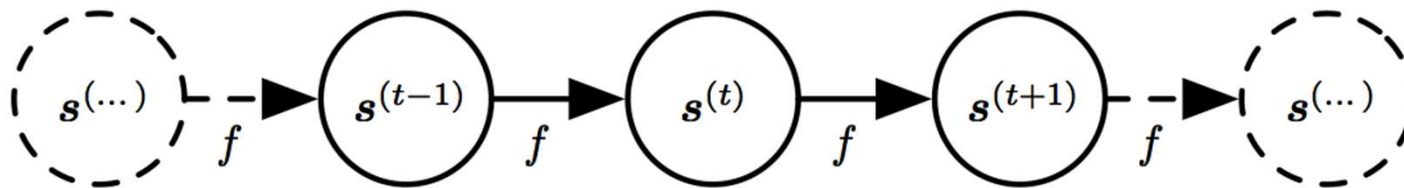
Clarinet in Bb

Piano

The musical score is for the first 8 measures of the 'Ode to Joy' movement. The Clarinet in Bb part is mostly rests. The Piano part features a melody in the right hand and a bass line in the left hand. Dynamics include *p*, *più p*, and *pp*. There are also markings for 'Leo.' and asterisks.

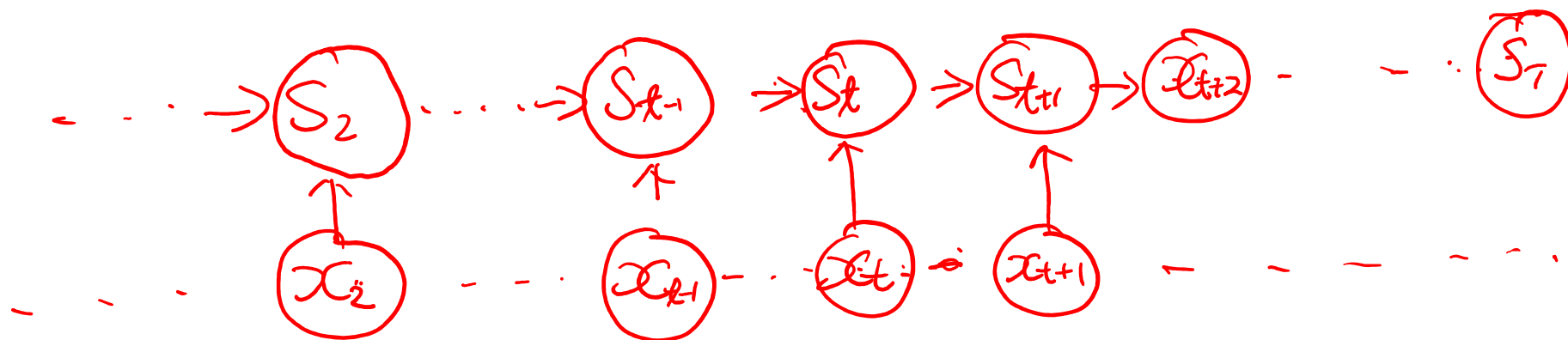
Learning over sequences

- As with images, we need to take structure (normally time temporal dependencies) into account, but note that in sequences **future and past** are not treated “symmetrically” (but could be, e.g., bidirectional ...)
- Also, we can have very **long temporal dependencies** and **inputs/outputs** of arbitrary length, even **infinite length**



Learning over sequences

- Transition matrix: $S_t = f(S_{t-1})$
 - 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3
 - 3
- However, it can't model the sequence below
(Example (noisy):
- 1.1, 2.2, 3.1, 1.3, 2.2, 3.3, 1.4, 2.4, 3.3, 1.1, 2.0,
3.1, 1.2, 2.1, 3.0, 1.1, 2.3, 3.4,
 - we need to use $S_t = f(S_{t-1}, X_t)$



Learning over sequences

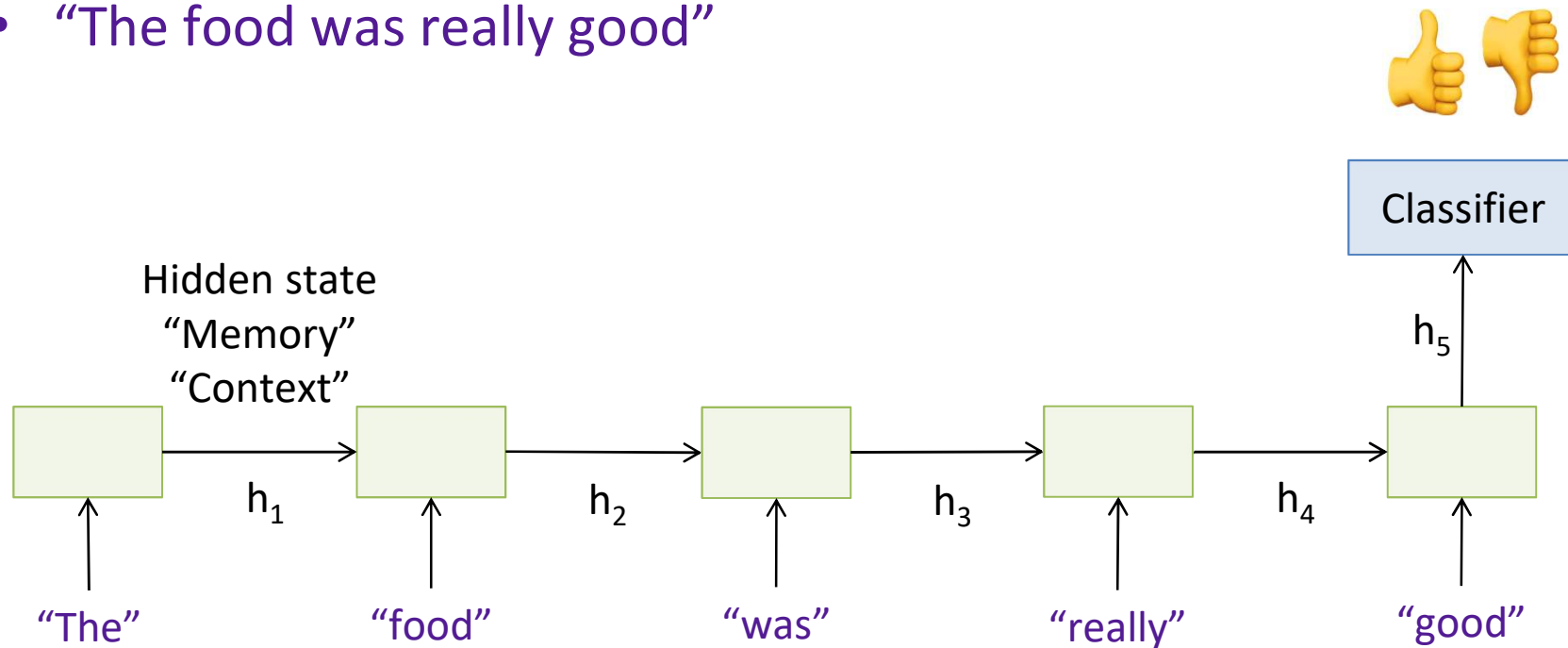
- Consider the word 'mountain'
- Given the string 'mountai', there is strong probability that the next character will be 'n'
- This is precisely what we want our neural network to be able to model
- More in general, we need to somehow **model context and memory**

Text classification: Examples

- **Sentiment classification:** classify a restaurant or movie or product review as positive or negative
 - “The food was really good”
 - “The vacuum cleaner broke within two weeks”
 - “The movie had slow parts, but overall was worth watching”
- What feature representation or predictor structure can we use for this problem?

Sentiment classification

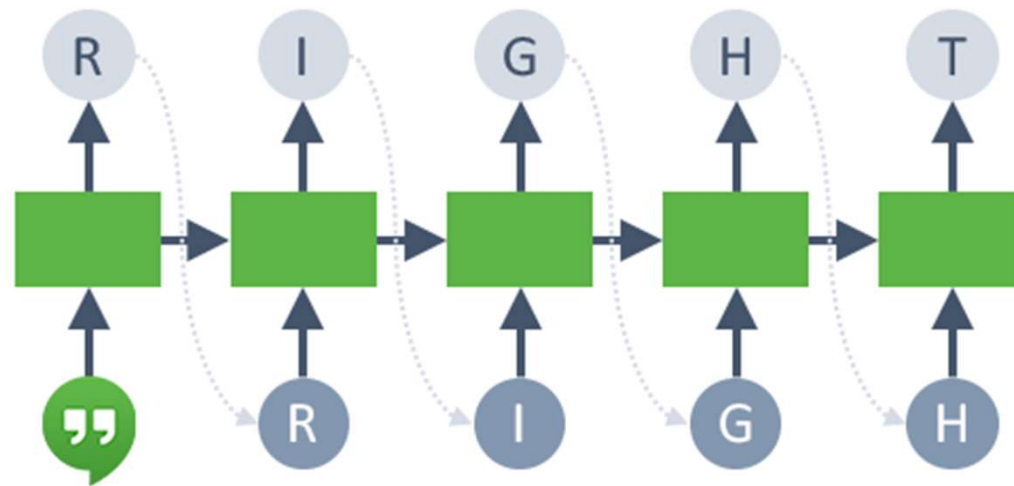
- “The food was really good”



Recurrent Neural Network (RNN)

Language Modeling

- Character RNN



Input representation

- In the following we assume our inputs are either characters (e.g., the network's task is to predict the next one) or words (e.g., the network's task is to predict the next word)
- In general, we will have a sequence of symbols in input that can be drawn from a larger vocabulary

Input representation

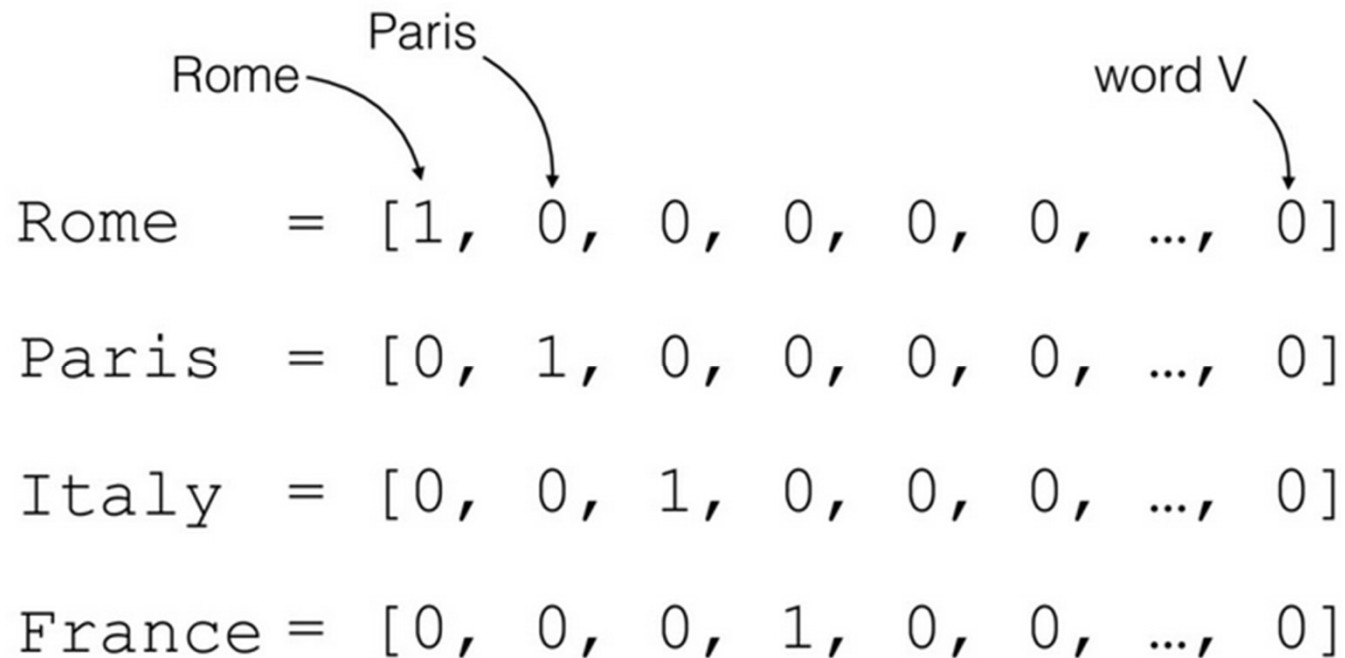
- We start with a one-hot vector representation

The diagram illustrates one-hot vector representations for four words: Rome, Paris, Italy, and France. Each word is represented by a vector of binary values (0s and 1s). The vectors are arranged in a grid. Arrows indicate the mapping from the word labels to the corresponding elements in the vectors: Rome points to the first element (1), Paris points to the second element (1), Italy points to the third element (1), and France points to the fourth element (1). The label 'word V' points to the last element (0) in each vector.

	Rome	Paris						word V
Rome	=	[1,	0,	0,	0,	0,	0,	..., 0]
Paris	=	[0,	1,	0,	0,	0,	0,	..., 0]
Italy	=	[0,	0,	1,	0,	0,	0,	..., 0]
France	=	[0,	0,	0,	1,	0,	0,	..., 0]

Input representation

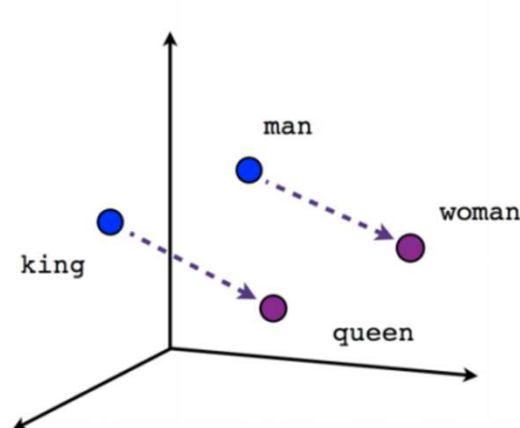
- We start with a one-hot vector representation



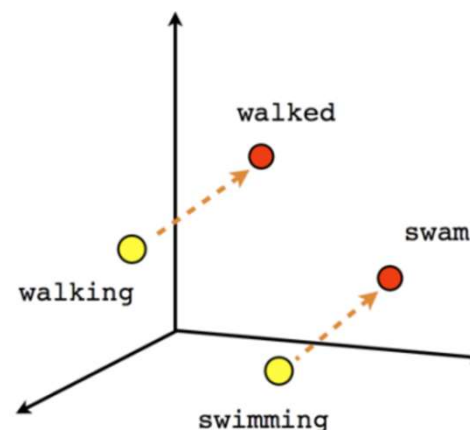
Note: we can do the same if our inputs are characters and one-hot encode them

Input representation

- We start with a one-hot vector representation
- If we're dealing with words, after the one-hot vector, we apply an embedding (e.g., Word2Vec)

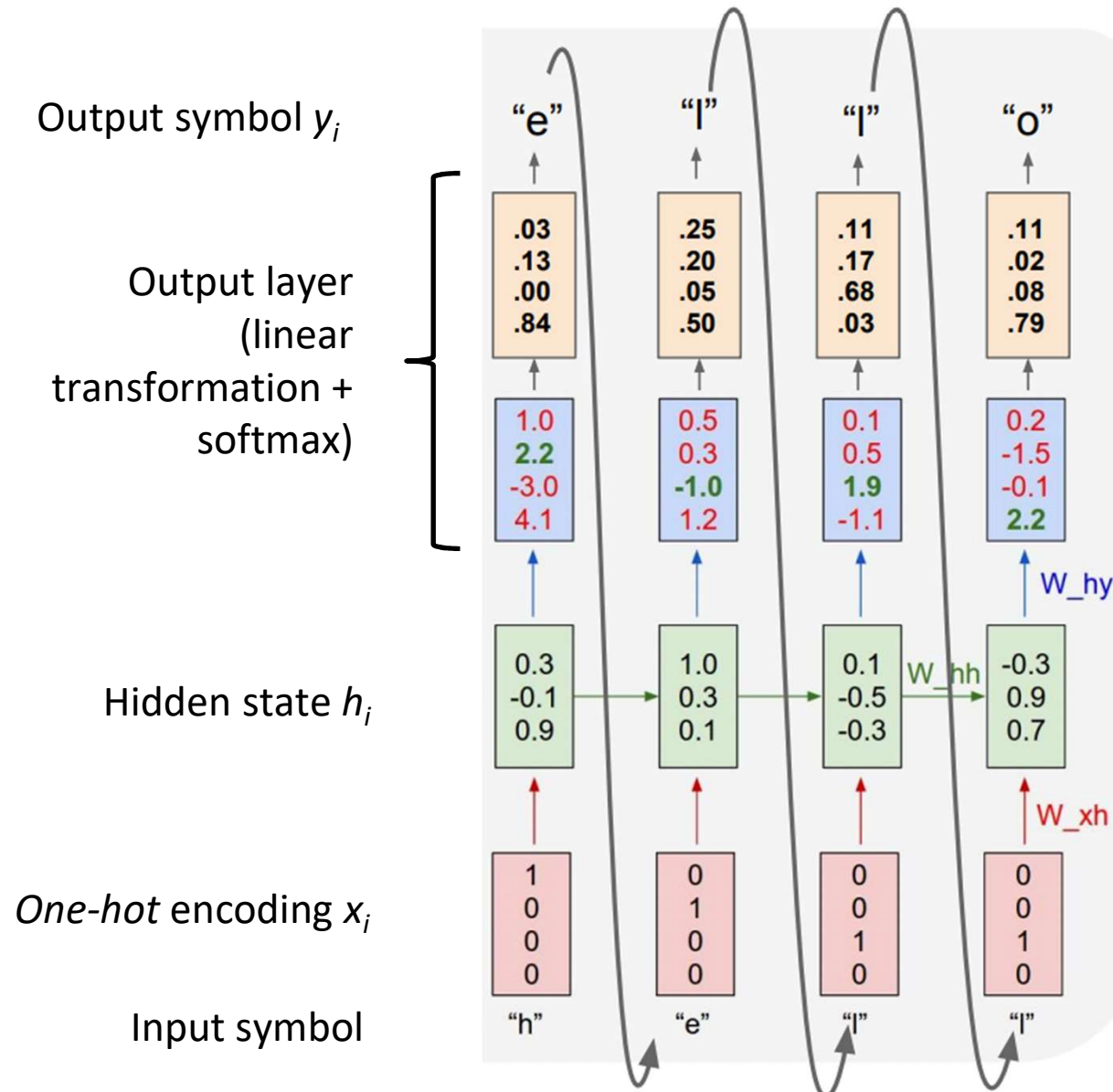


Male-Female



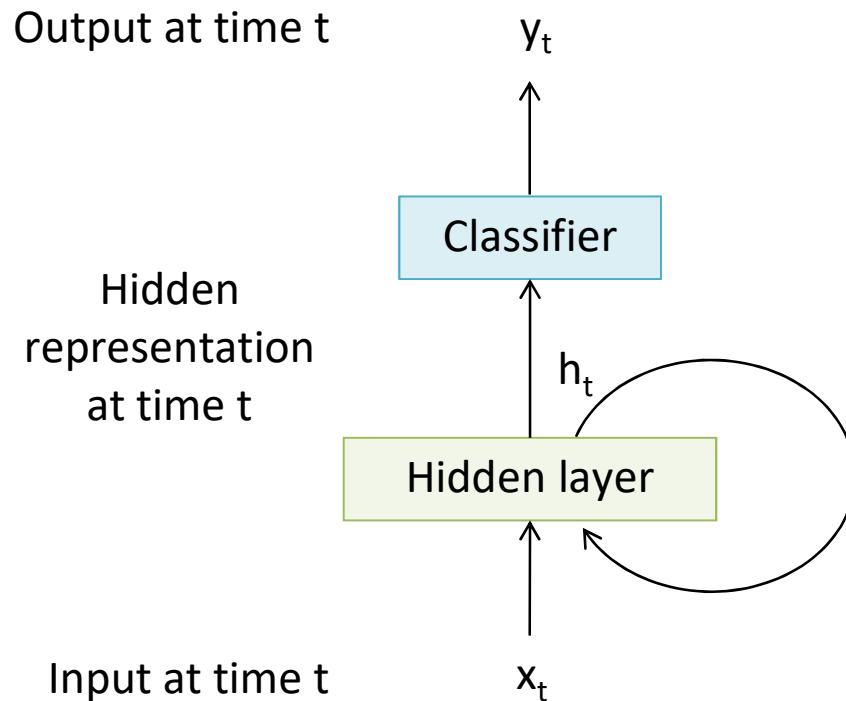
Verb tense

Character RNN



$$\begin{aligned}
 & p(y_1, y_2, \dots, y_n) \\
 &= \prod_{i=1}^n p(y_i | y_1, \dots, y_{i-1}) \\
 &\approx \prod_{i=1}^n P_W(y_i | h_i)
 \end{aligned}$$

Recurrent Neural Network (RNN)

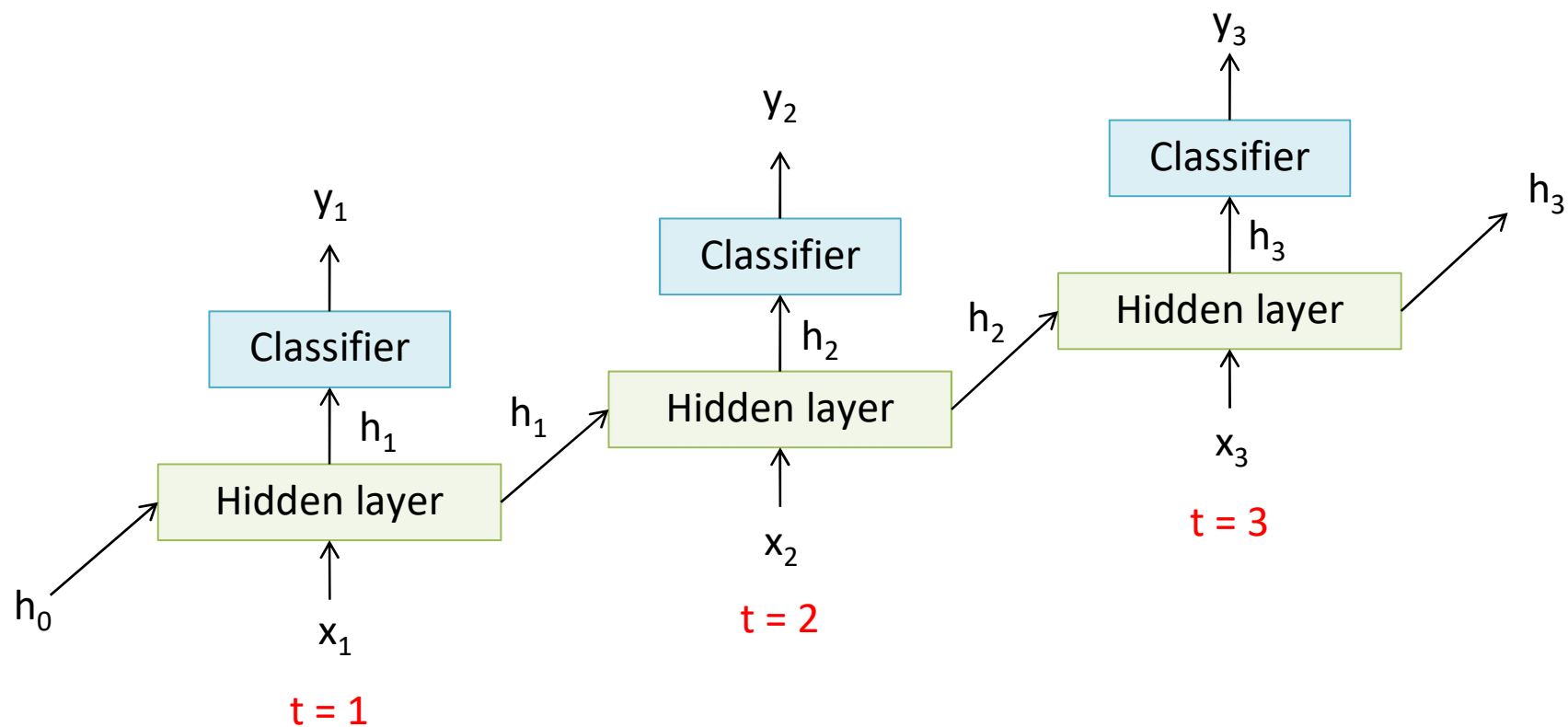


Recurrence:

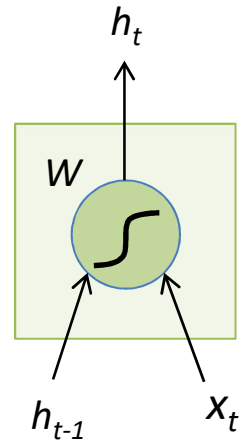
$$h_t = f_W(x_t, h_{t-1})$$

new state function of W input at time t old state

Unrolling the RNN

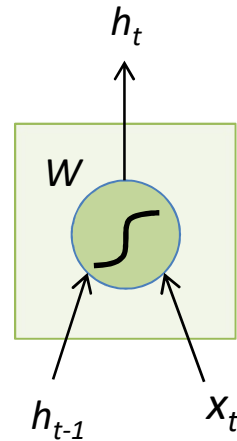


Vanilla RNN Cell

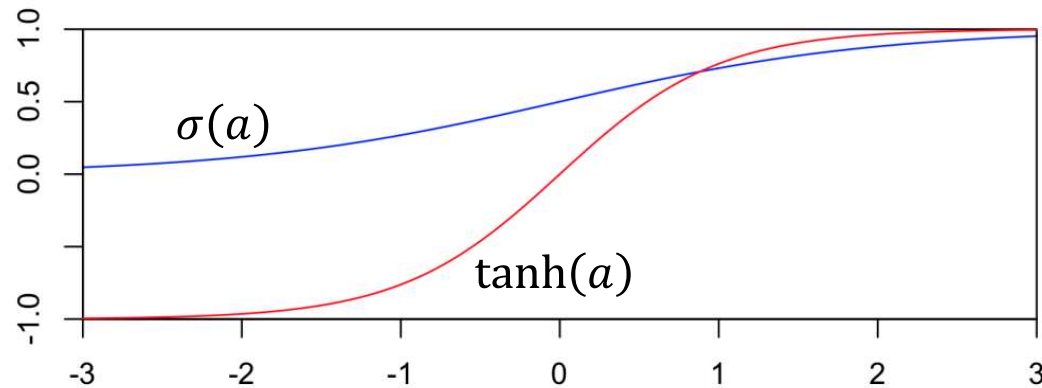


$$\begin{aligned} h_t &= f_W(x_t, h_{t-1}) \\ &= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} \end{aligned}$$

Vanilla RNN Cell

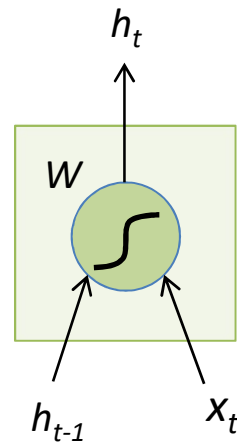


$$\begin{aligned} h_t &= f_W(x_t, h_{t-1}) \\ &= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} \end{aligned}$$

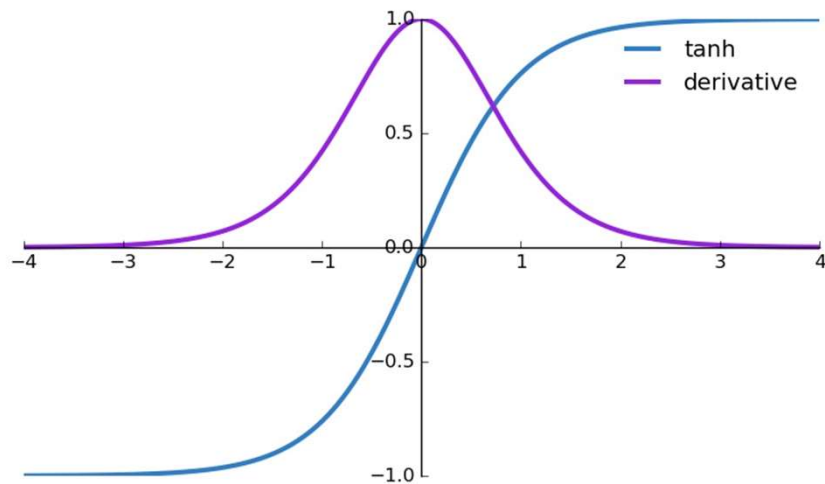


$$\begin{aligned} \tanh(a) &= \frac{e^a - e^{-a}}{e^a + e^{-a}} \\ &= 2\sigma(2a) - 1 \end{aligned}$$

Vanilla RNN Cell

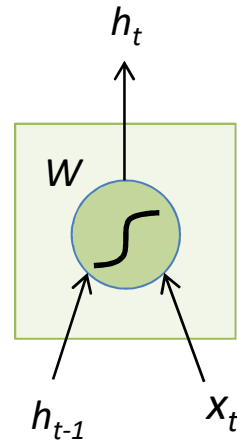


$$h_t = f_W(x_t, h_{t-1})$$
$$= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

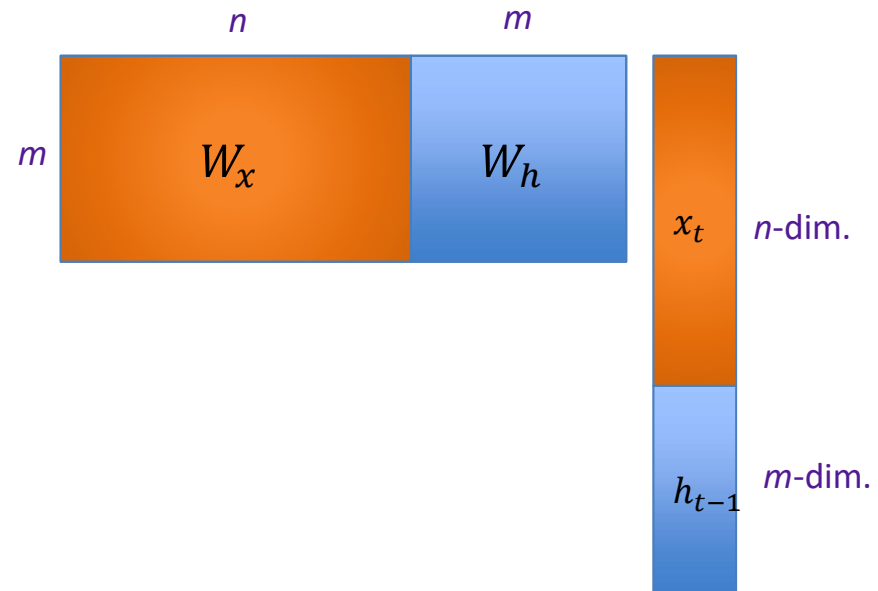


$$\frac{d}{da} \tanh(a) = 1 - \tanh^2(a)$$

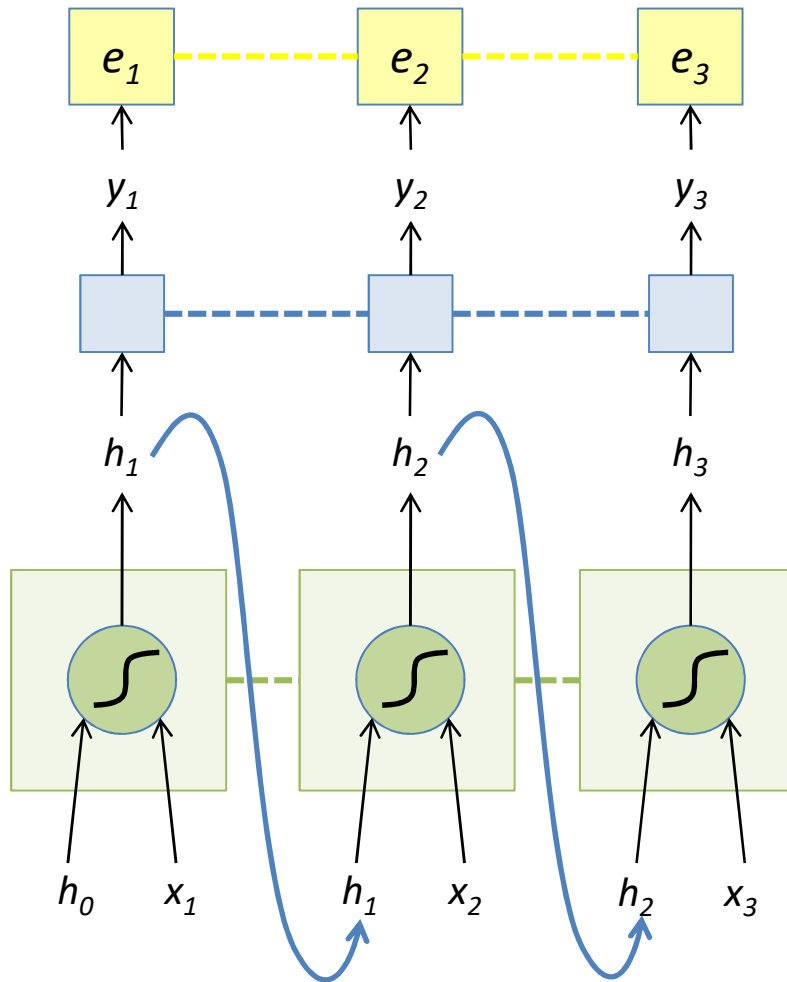
Vanilla RNN Cell



$$\begin{aligned} h_t &= f_W(x_t, h_{t-1}) \\ &= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} \\ &= \tanh(W_x x_t + W_h h_{t-1}) \end{aligned}$$



RNN Forward Pass



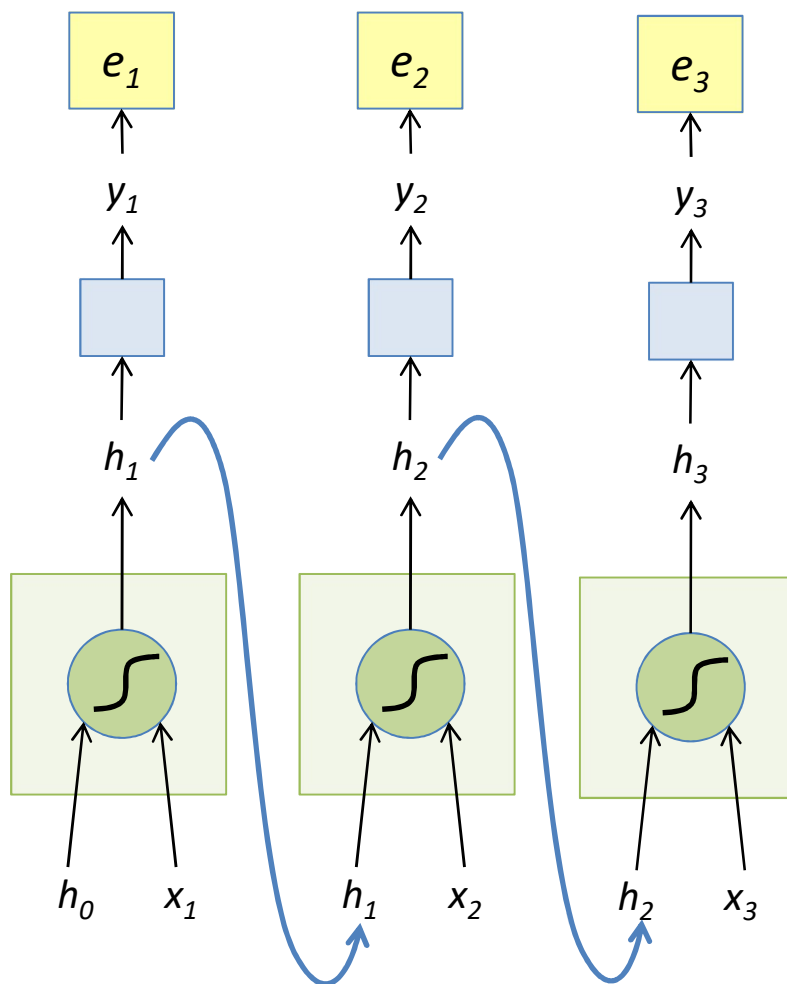
$$e_t = -\log(y_t(GT_t))$$

$$y_t = \text{softmax}(W_y h_t)$$

$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

----- shared weights

Unfolded RNN Forward Pass

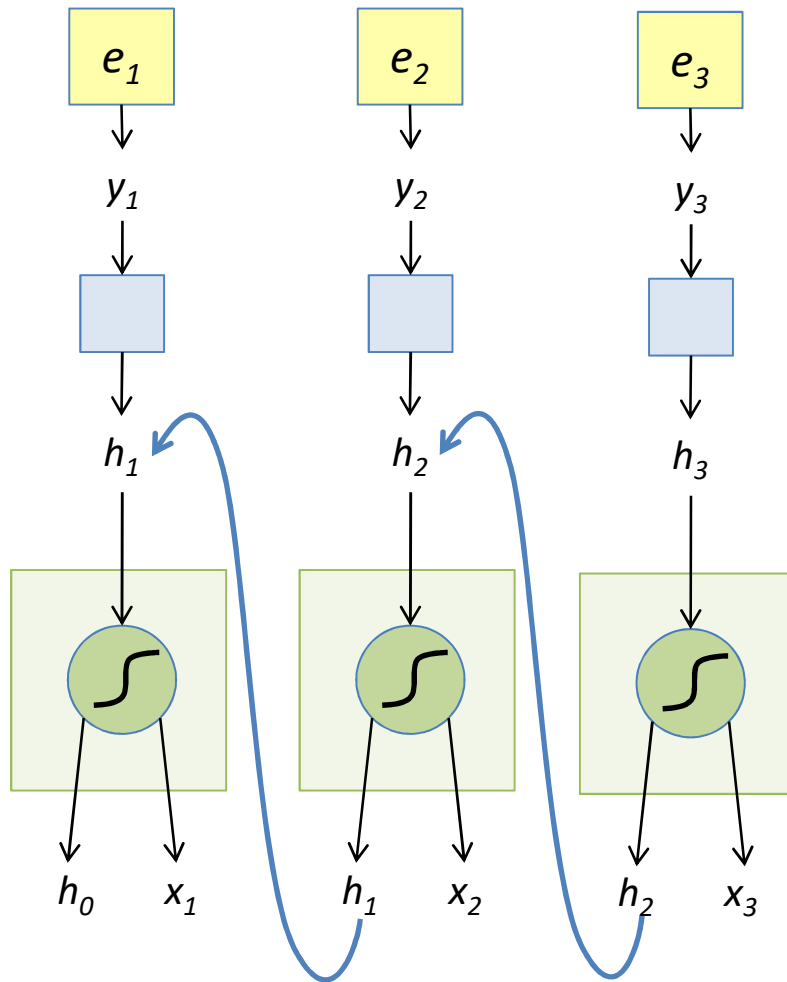


$$e_t = -\log(y_t(GT_t))$$

$$y_t = \text{softmax}(W_y h_t)$$

$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

Unfolded RNN Backward Pass



$$e_t = -\log(y_t(GT_t))$$

$$y_t = \text{softmax}(W_y h_t)$$

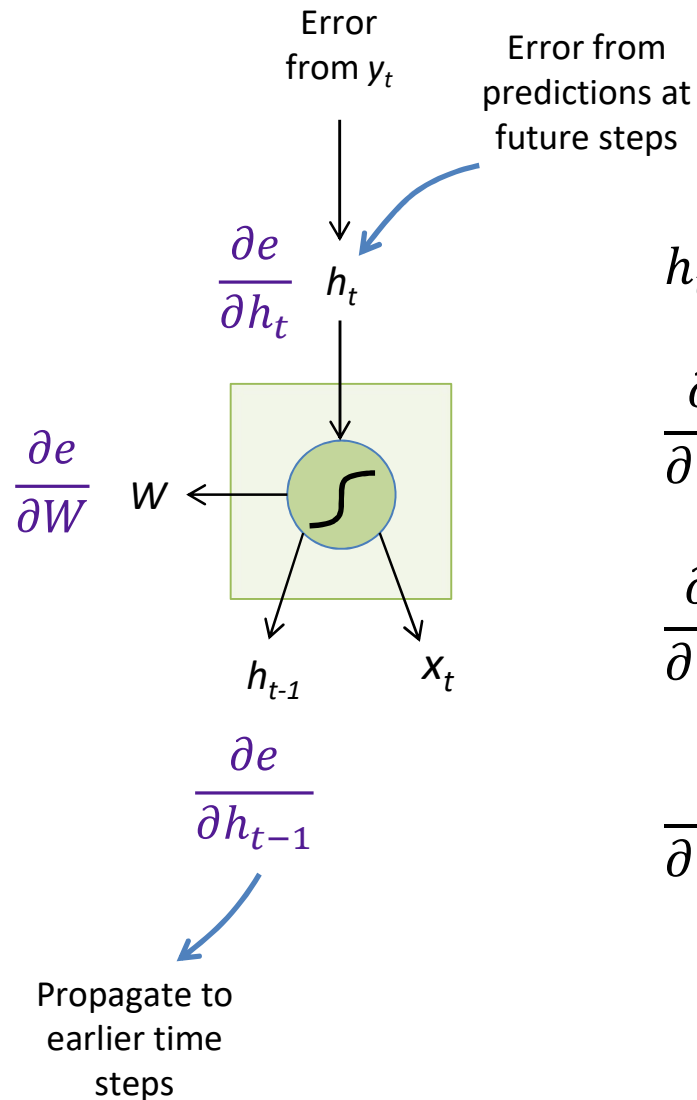
$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

Backpropagation Through Time (BPTT)

- Most common method used to train RNNs
- The unfolded network (used during forward pass) is treated as one big feed-forward network that accepts the whole time series as input
- The weight updates are computed for each copy in the unfolded network, then summed (or averaged) and applied to the RNN weights
- In practice, *truncated* BPTT is used: run the RNN forward k_1 time steps, propagate backward for k_2 time steps

<https://machinelearningmastery.com/gentle-introduction-backpropagation-time/>
http://www.cs.utoronto.ca/~ilya/pubs/ilya_sutskever_phd_thesis.pdf

RNN Backward Pass



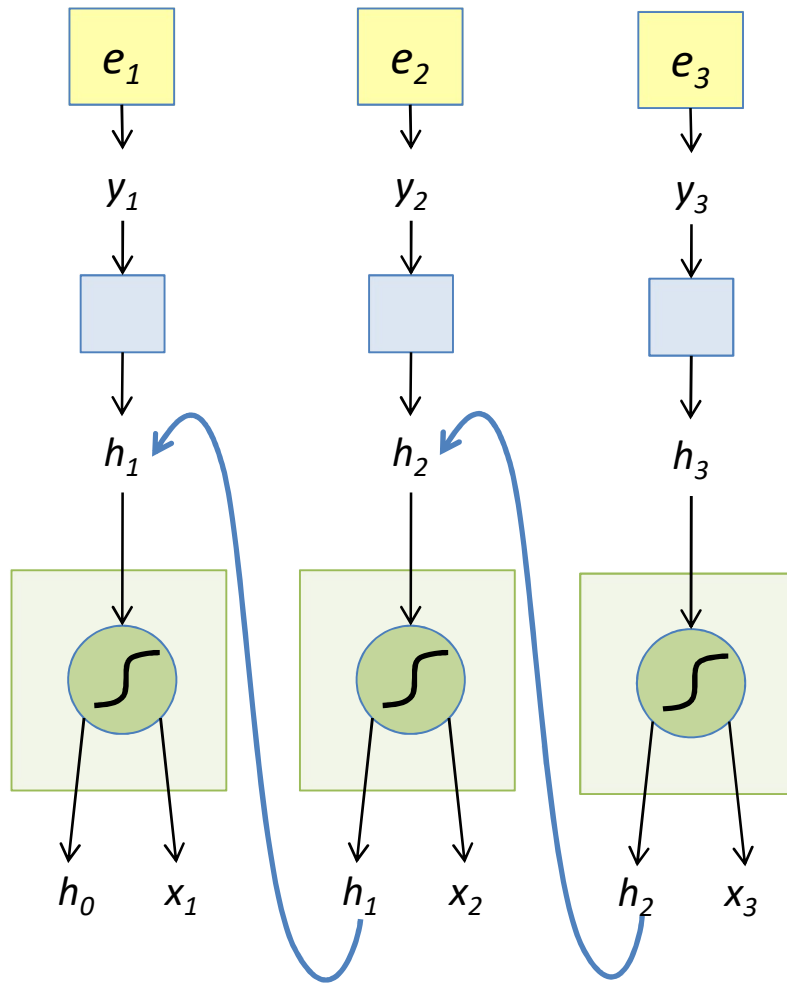
$$h_t = \tanh(W_x x_t + W_h h_{t-1})$$

$$\frac{\partial e}{\partial W_h} = \frac{\partial e}{\partial h_t} \odot (1 - \tanh^2(W_x x_t + W_h h_{t-1})) h_{t-1}^T$$

$$\frac{\partial e}{\partial W_x} = \frac{\partial e}{\partial h_t} \odot (1 - \tanh^2(W_x x_t + W_h h_{t-1})) x_t^T$$

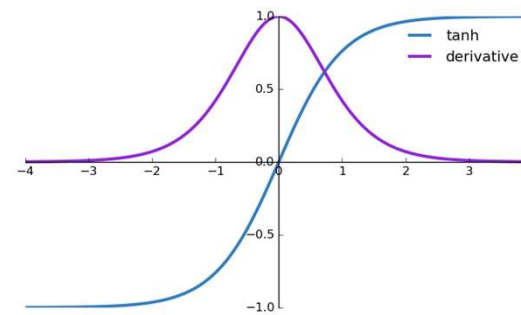
$$\frac{\partial e}{\partial h_{t-1}} = W_h^T (1 - \tanh^2(W_x x_t + W_h h_{t-1})) \odot \frac{\partial e}{\partial h_t}$$

RNN Backward Pass



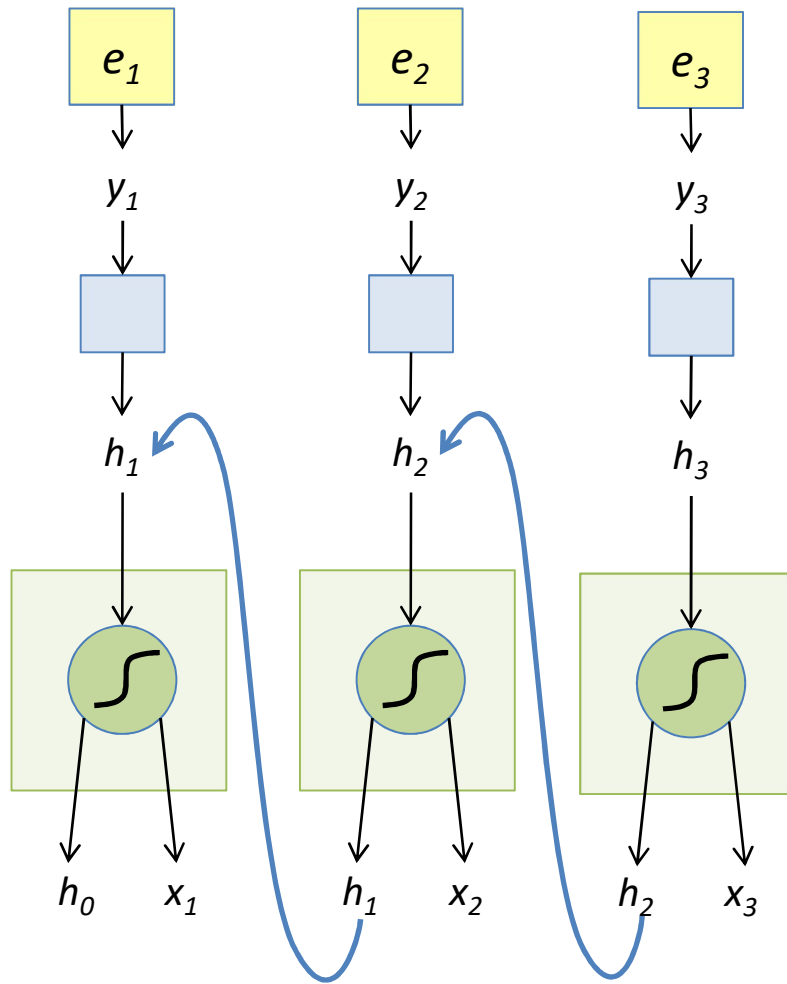
$$\frac{\partial e}{\partial h_{t-1}} = W_h^T (1 - \tanh^2(W_x x_t + W_h h_{t-1})) \odot \frac{\partial e}{\partial h_t}$$

Large tanh activations will give small gradients



Consider $\frac{\partial e_n}{\partial h_k}$ for $k \ll n$

RNN Backward Pass



$$\frac{\partial e}{\partial h_{t-1}} = W_h^T (1 - \tanh^2(W_x x_t + W_h h_{t-1})) \odot \frac{\partial e}{\partial h_t}$$

Gradients will vanish if
largest singular value of
 W_h is less than 1

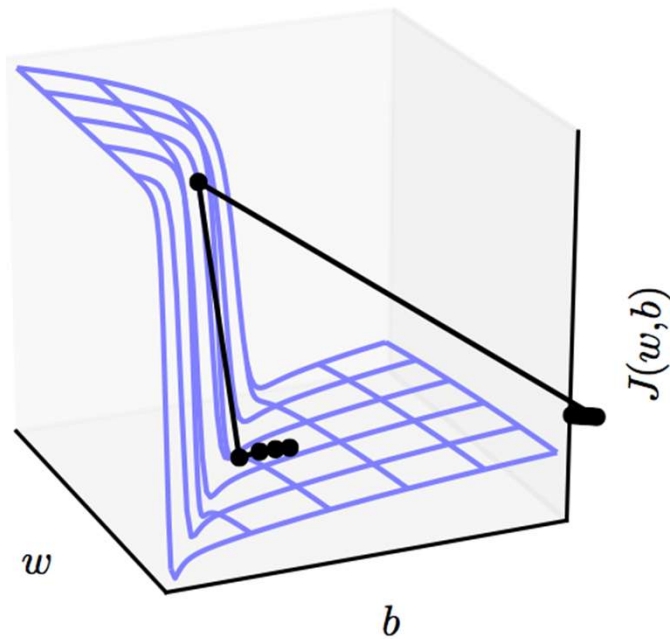
Consider $\frac{\partial e_n}{\partial h_k}$ for $k \ll n$

Backpropagation through time

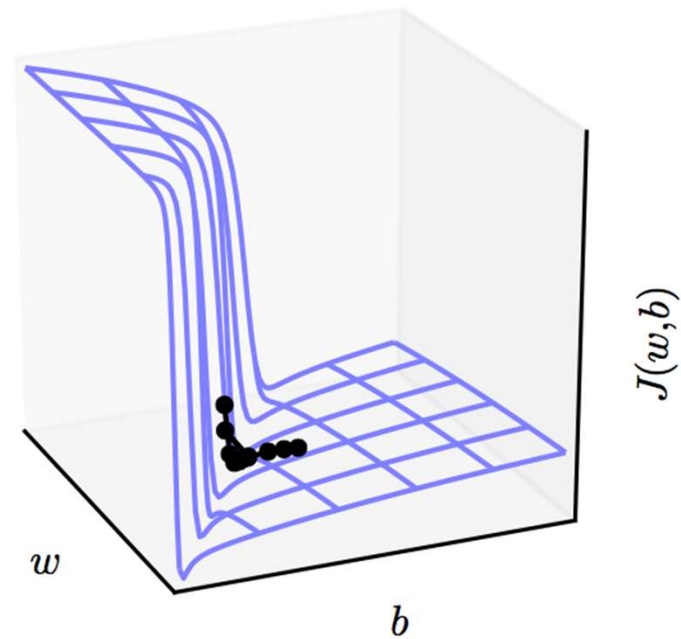
- Backpropagation is done similarly to MLP/CNN
- However the layers now correspond to different steps/times, so we call it **backpropagation through time** (section 10.2.2. in the book)
- As with very deep networks we need to multiply together many different gradients
- For long sequences, this is a problem
 - Vanishing/exploding gradients

Gradient clipping

Without clipping



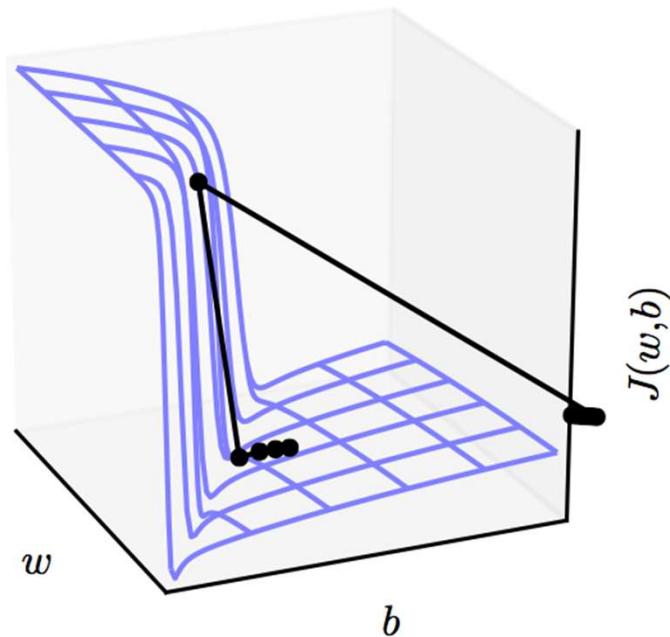
With clipping



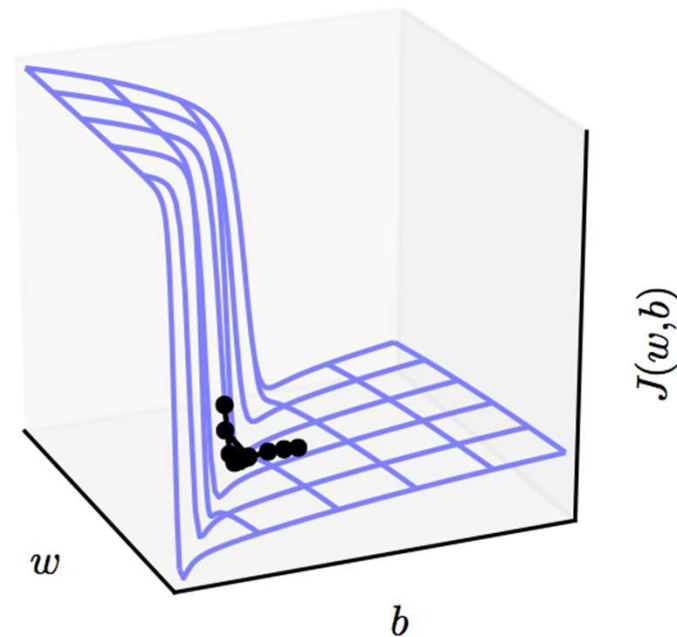
Exploding gradients are easily solved by rescaling or clipping the gradient

Gradient clipping

Without clipping



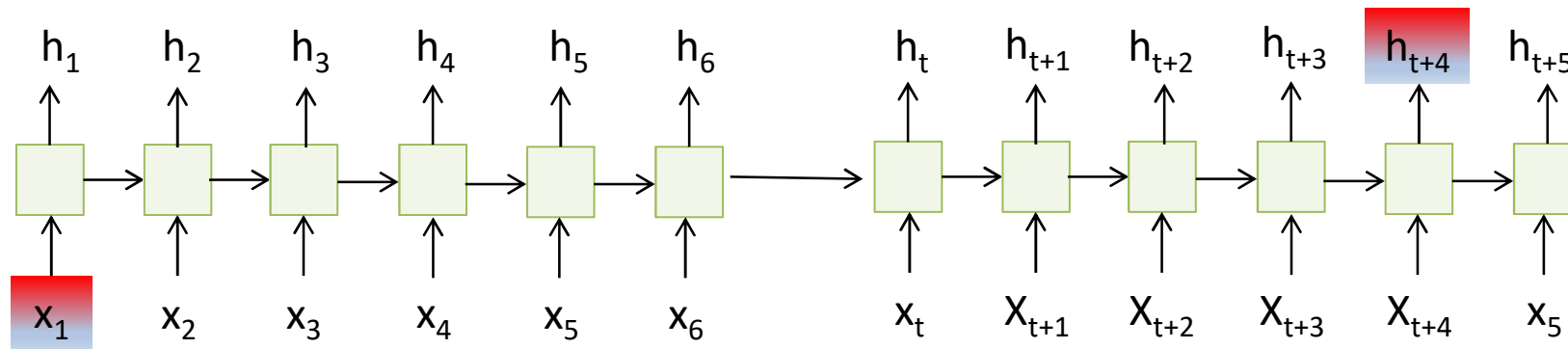
With clipping



But what about vanishing gradients?
There is no gradient to begin with, what should we clip/rescale?
Actually, there's even another problem...

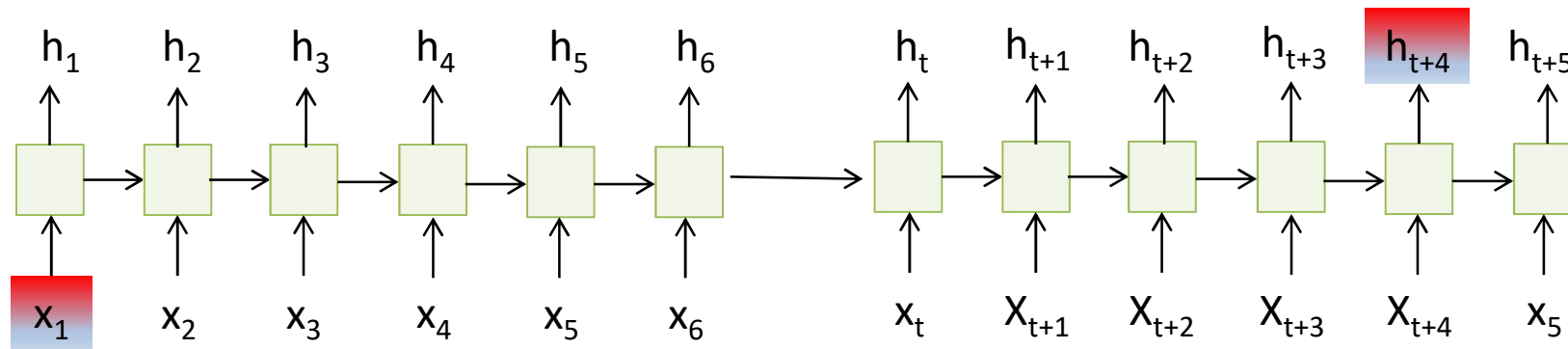
Long-term dependencies

Simplified version of RNN (omitting y)



In theory RNNs should be able to capture long-term dependencies
but in practice they're not!

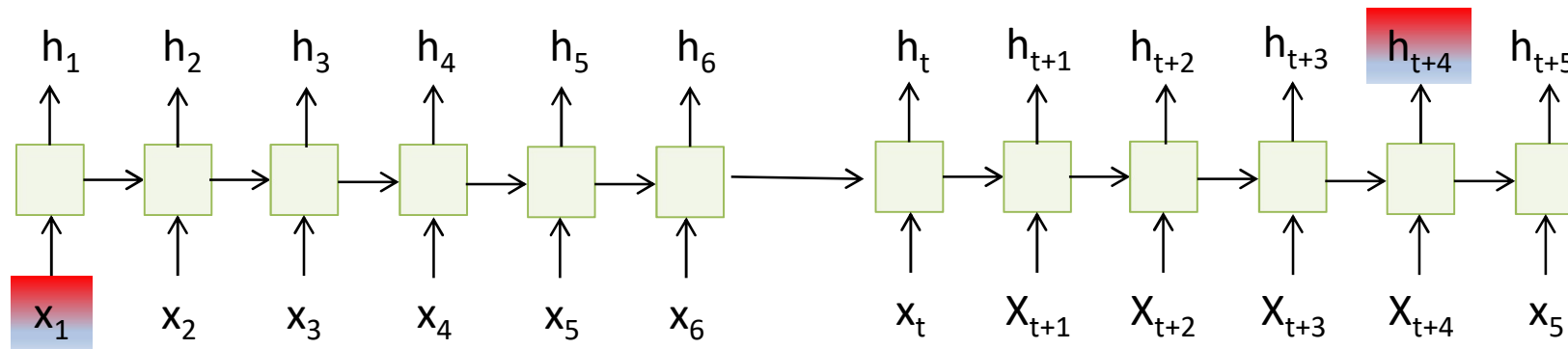
Long-term dependencies



$$h^{(t)} = W^{\top} h^{(t-1)}$$

Considering a very simply RNN described by this recurrence relation

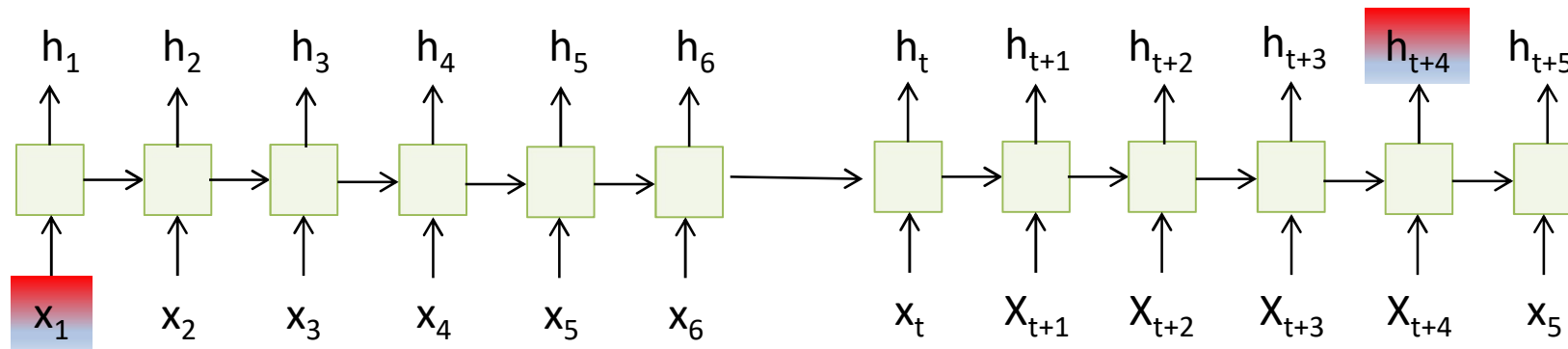
Long-term dependencies



$$\mathbf{h}^{(t)} = (\mathbf{W}^t)^\top \mathbf{h}^{(0)}$$

This can be simplified as shown above

Long-term dependencies

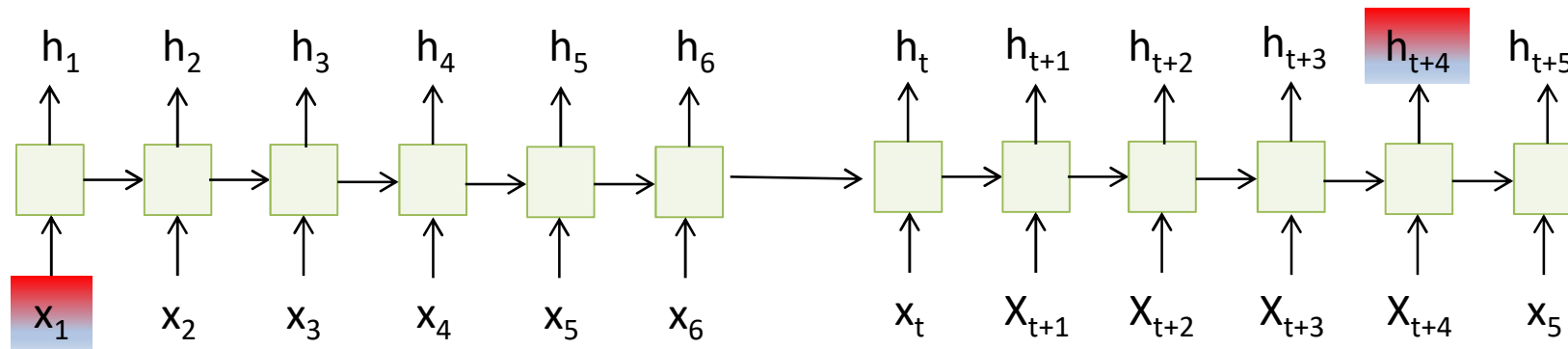


$$h^{(t)} = Q^\top \Lambda^t Q h^{(0)}$$

$$W = Q \Lambda Q^\top$$

Given the matrix decomposition of W , we get the formula above

Long-term dependencies

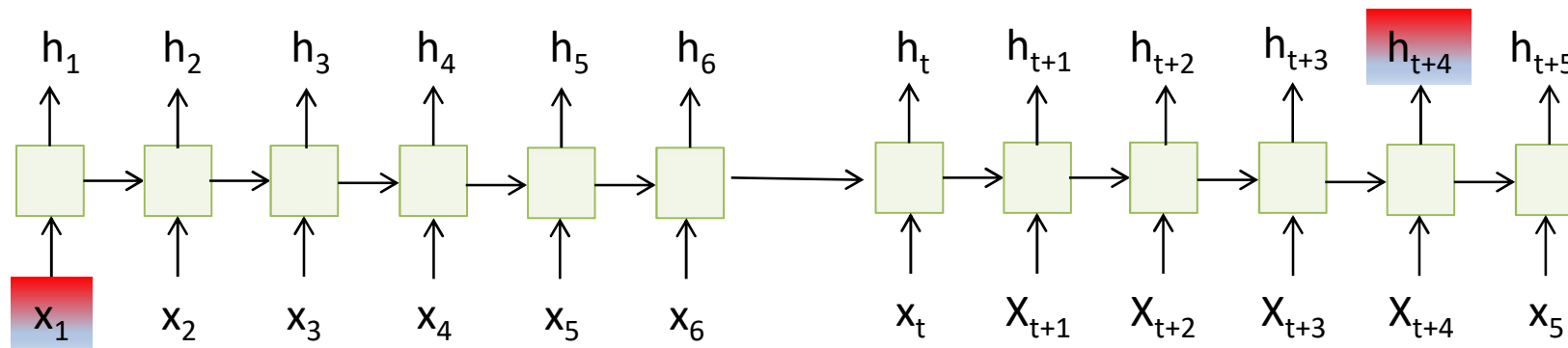


$$h^{(t)} = Q^\top \Lambda^t Q h^{(0)}$$

$$W = Q \Lambda Q^\top$$

eigenvalue with magnitude < 1 vanish, those with magnitude > 1 explode

Long-term dependencies

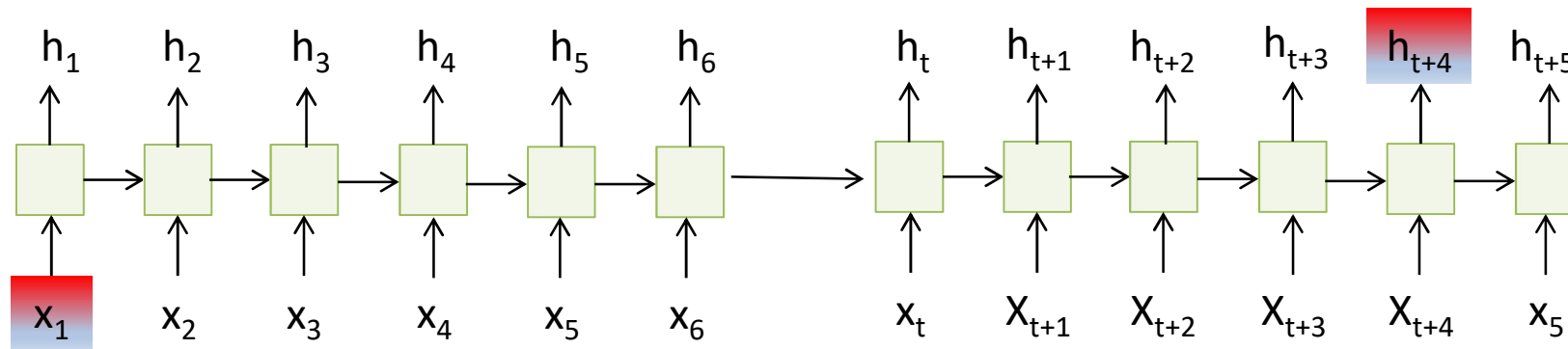


$$h^{(t)} = Q^T \Lambda^t Q h^{(0)}$$

$$W = Q \Lambda Q^T$$

Any component of $h(0)$ not aligned with max eigenvector of W is eventually discarded

Long-term dependencies

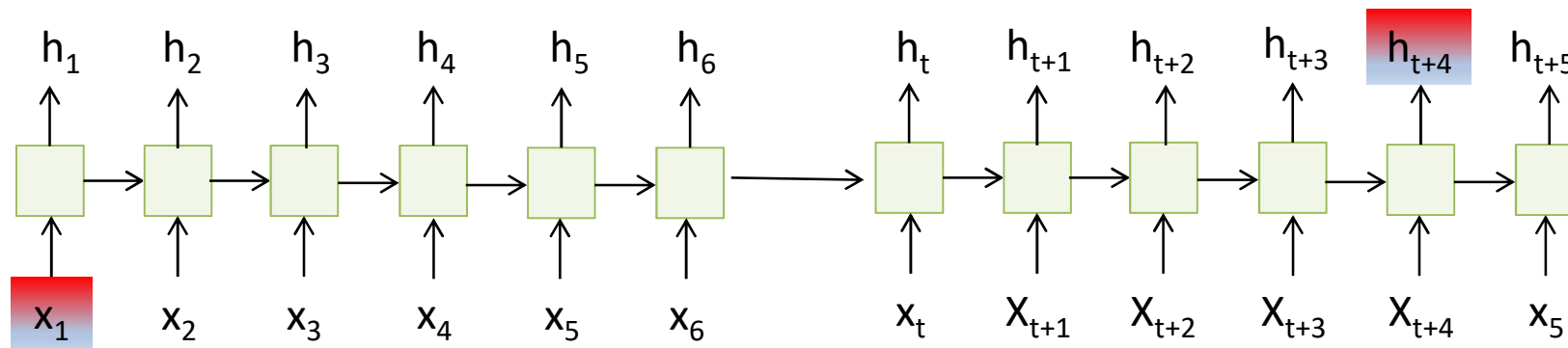


$$h^{(t)} = Q^\top \Lambda^t Q h^{(0)}$$

$$W = Q \Lambda Q^\top$$

Note that this problem is particular of RNNs. In MLPs W is not shared by all layers.

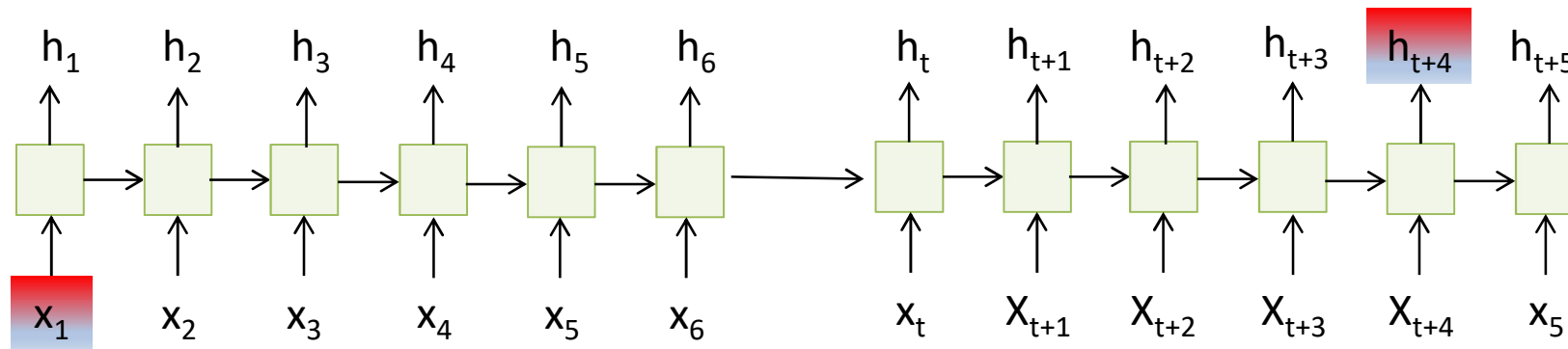
Long-term dependencies



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

Unfortunately it turns out that we can't avoid regions of the parameters space with exploding/vanishing gradients. It's there we need to go to learn long-term dependencies...

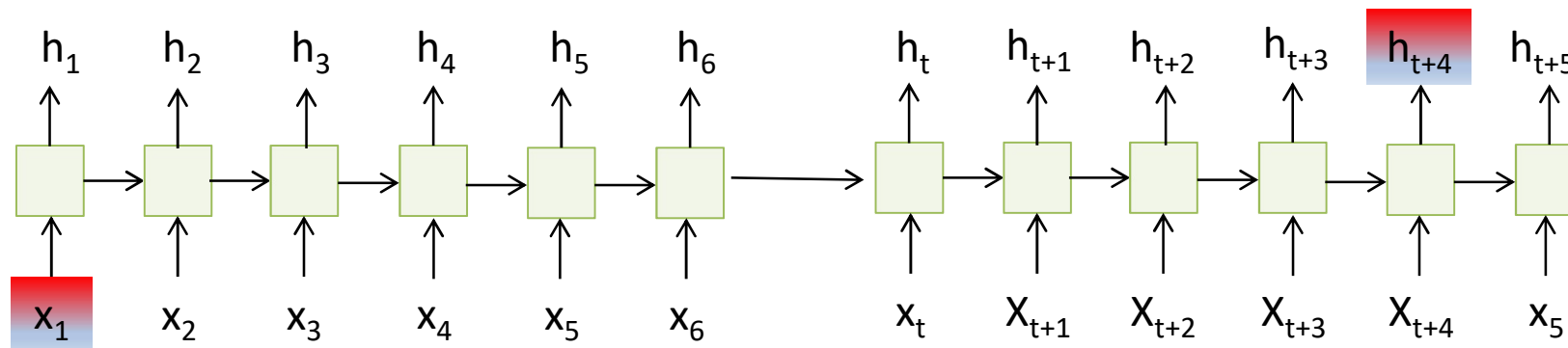
Long-term dependencies



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

There the training doesn't stop but becomes unsustainably slow.
In practice, we can't train RNNs for sequences longer than ~10

Long-term dependencies

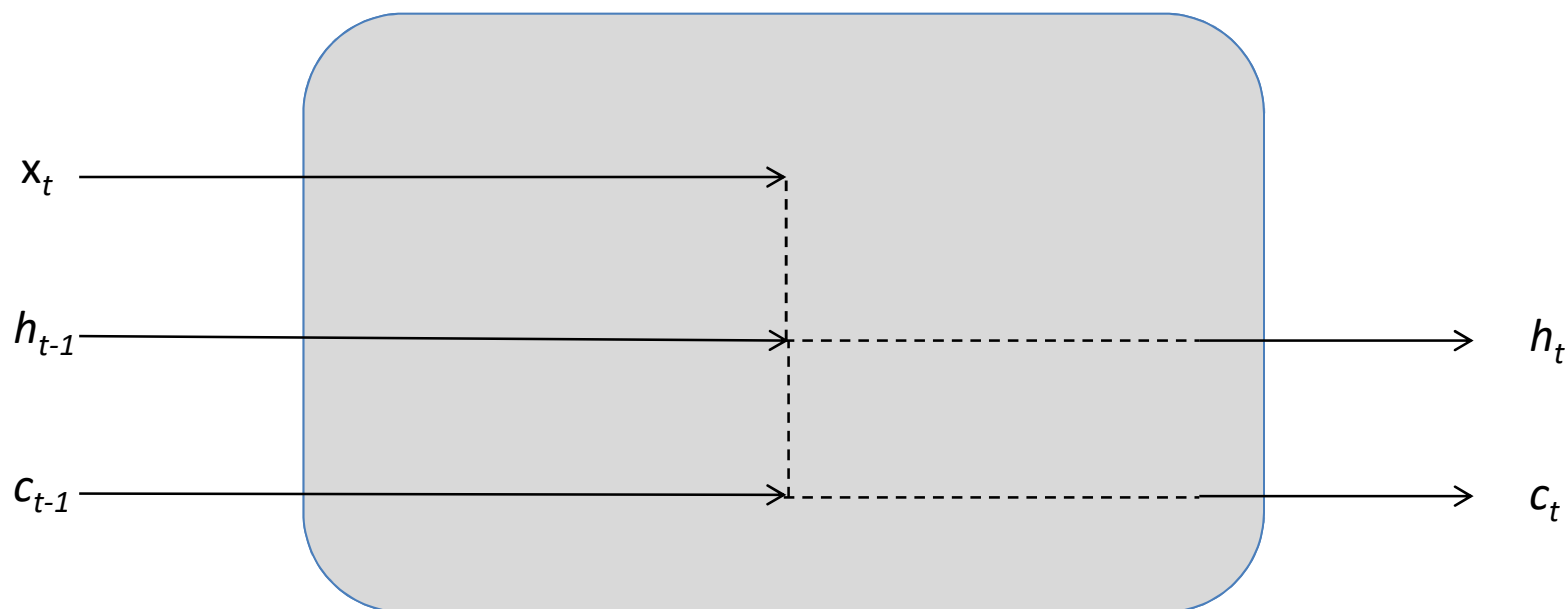


$$h^{(t)} = Q^T \Lambda^t Q h^{(0)}$$

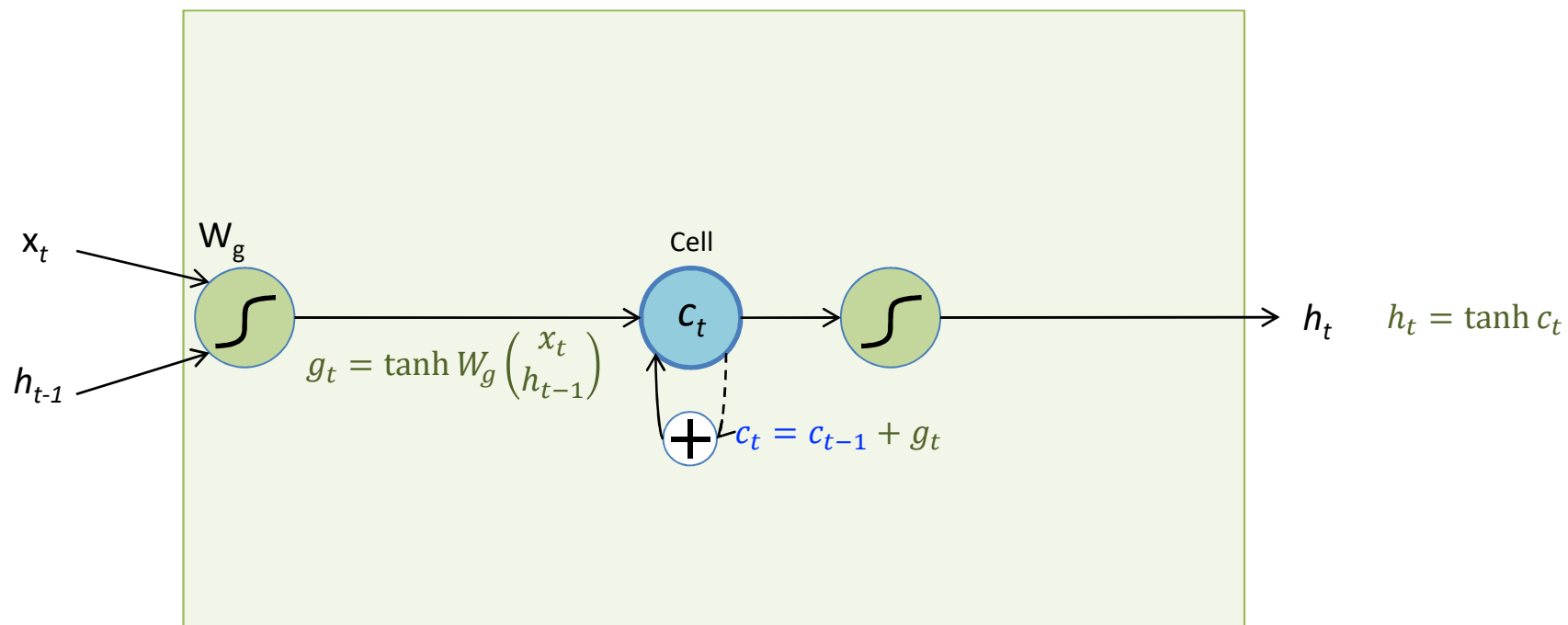
But there may be another way...

Long Short-Term Memory (LSTM)

- Add a *memory cell* that is not subject to matrix multiplication or squishing, thereby avoiding gradient decay

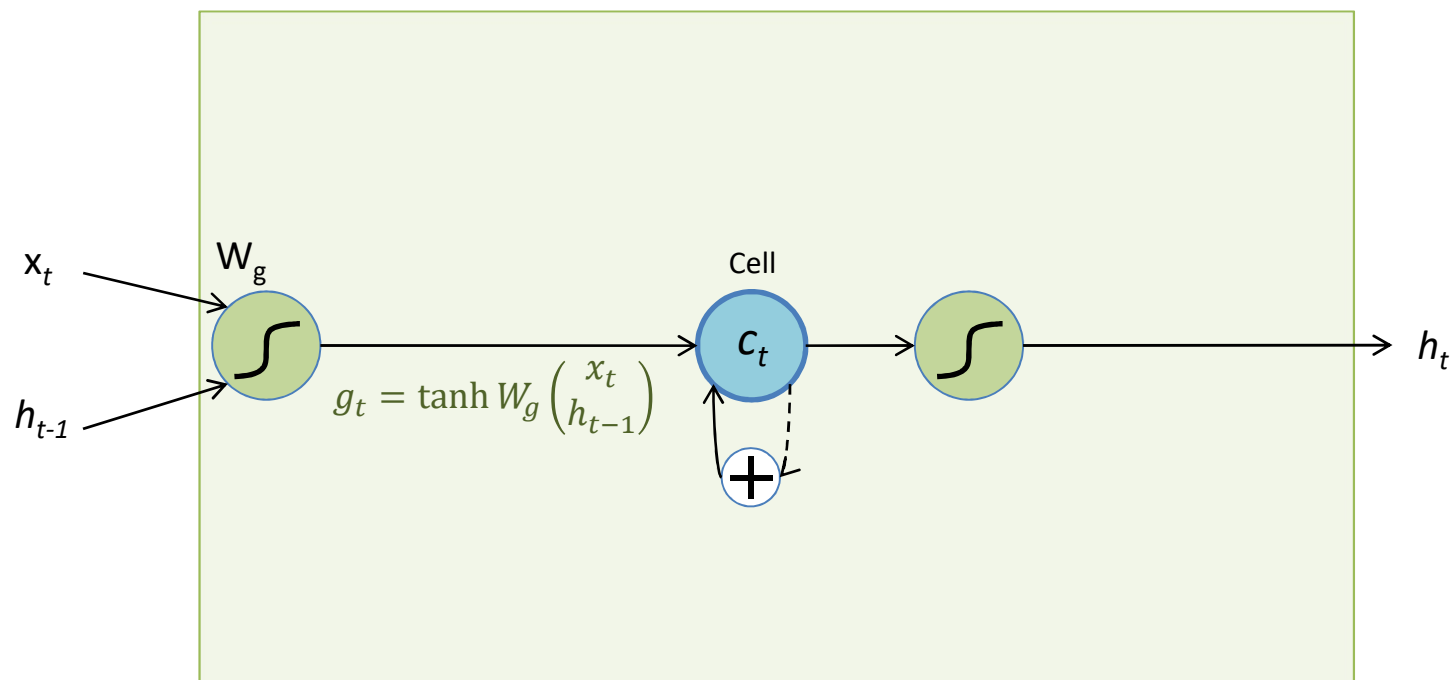


The LSTM Cell

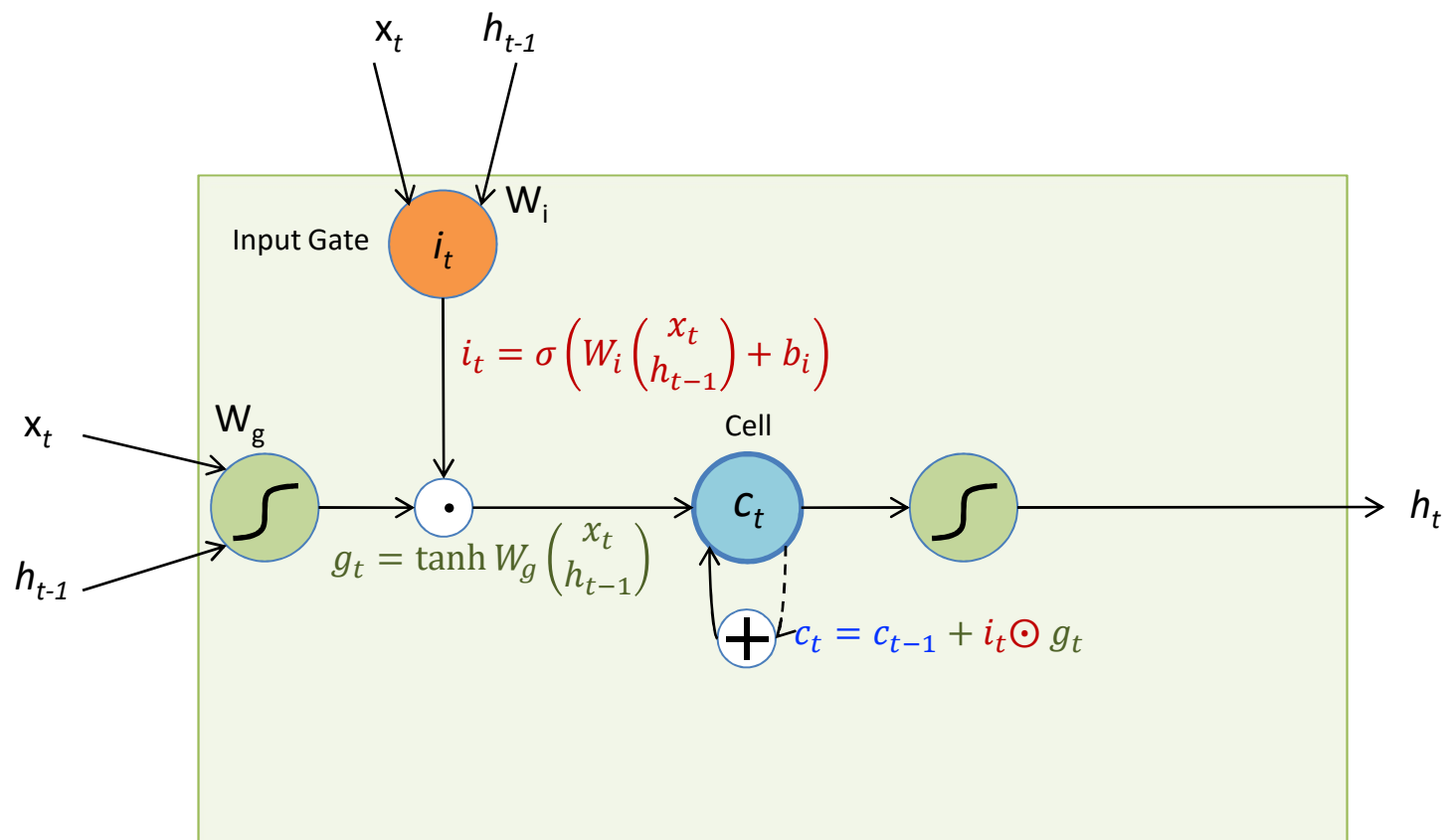


* Dashed line indicates time-lag

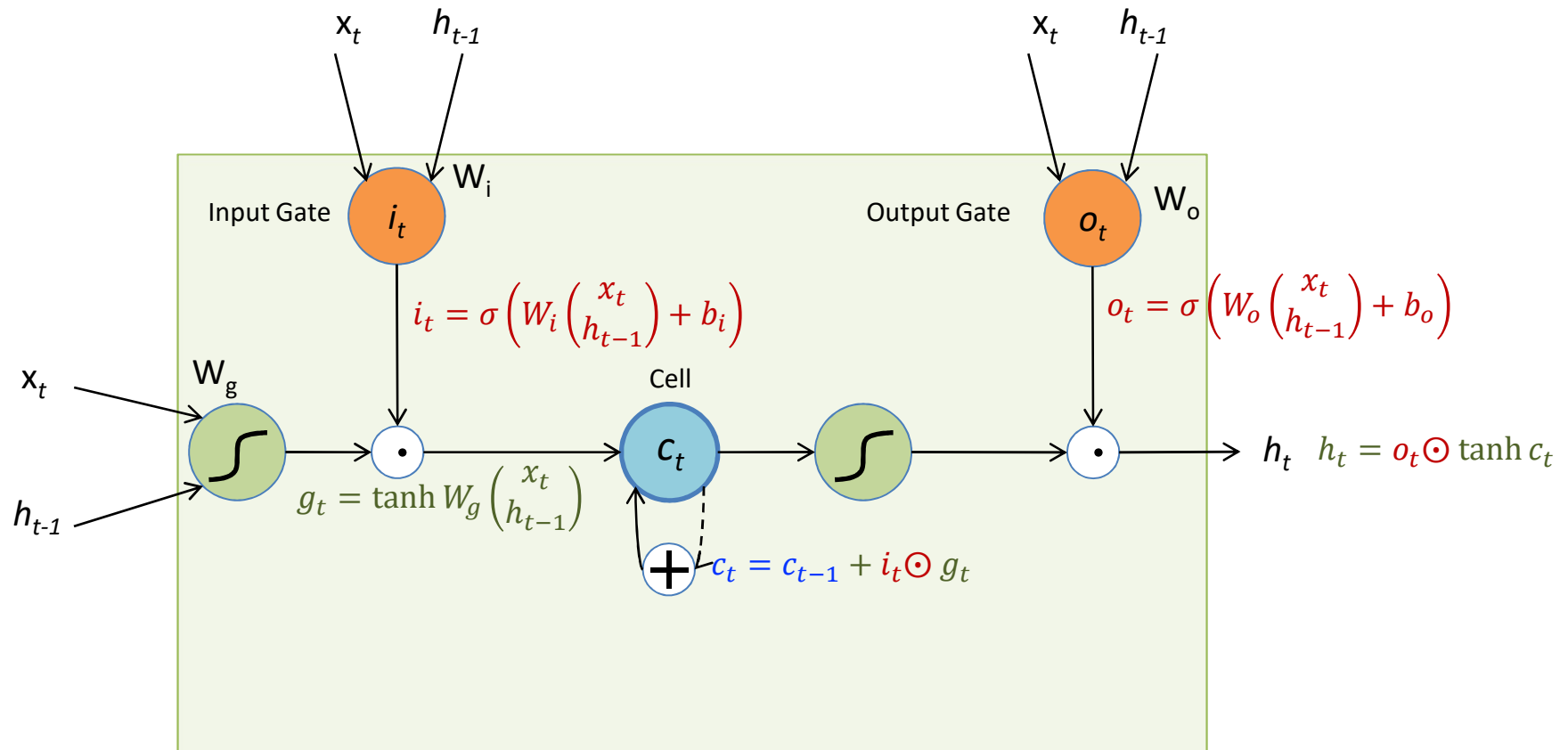
The LSTM Cell



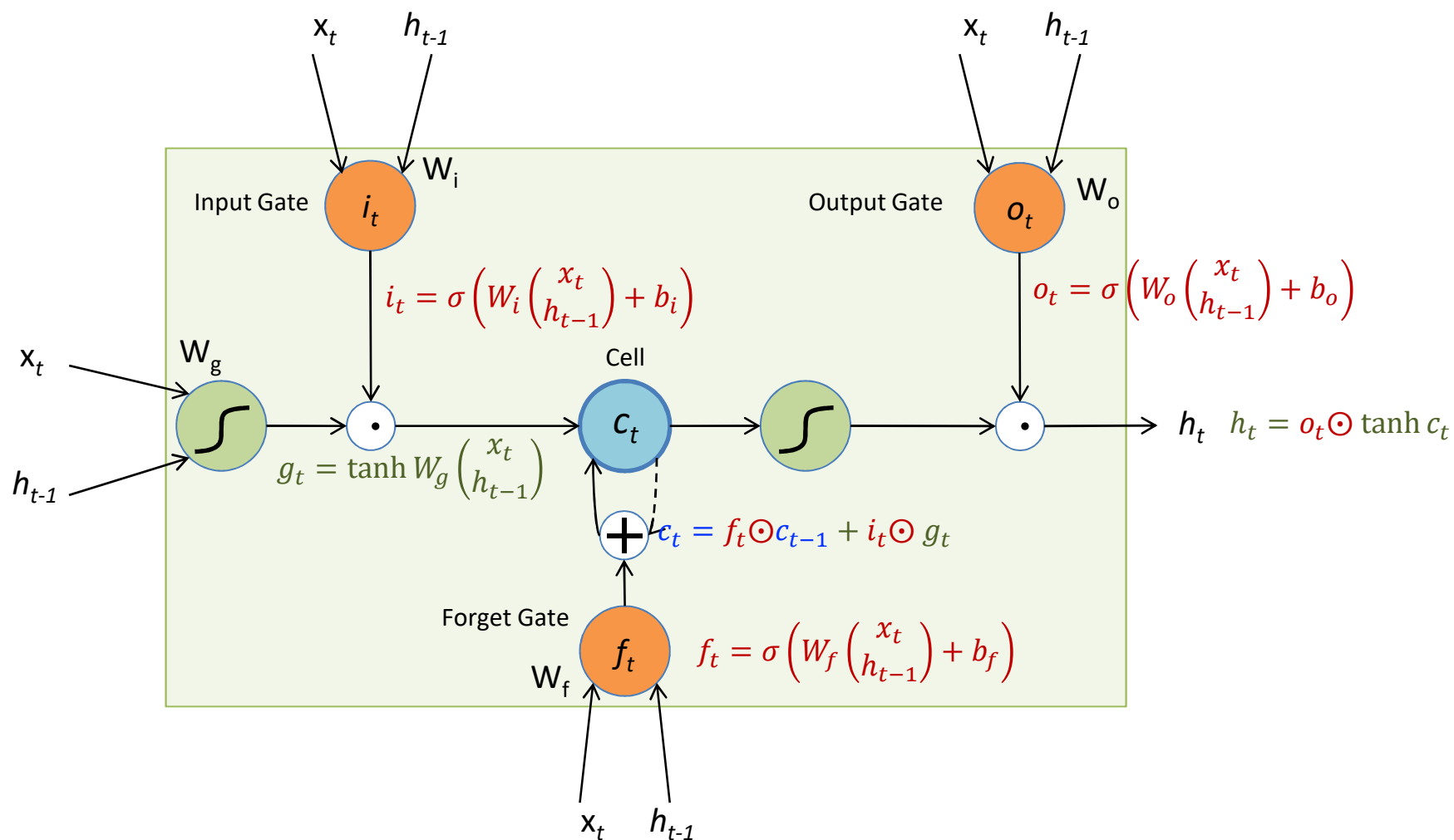
The LSTM Cell



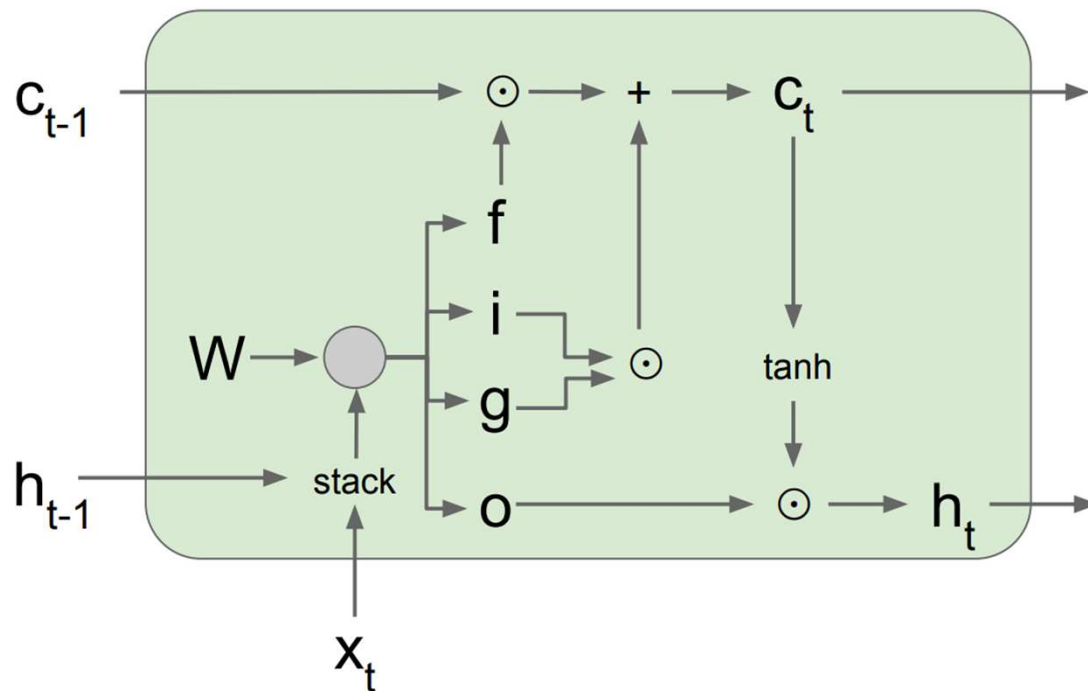
The LSTM Cell



The LSTM Cell



LSTM Forward Pass Summary

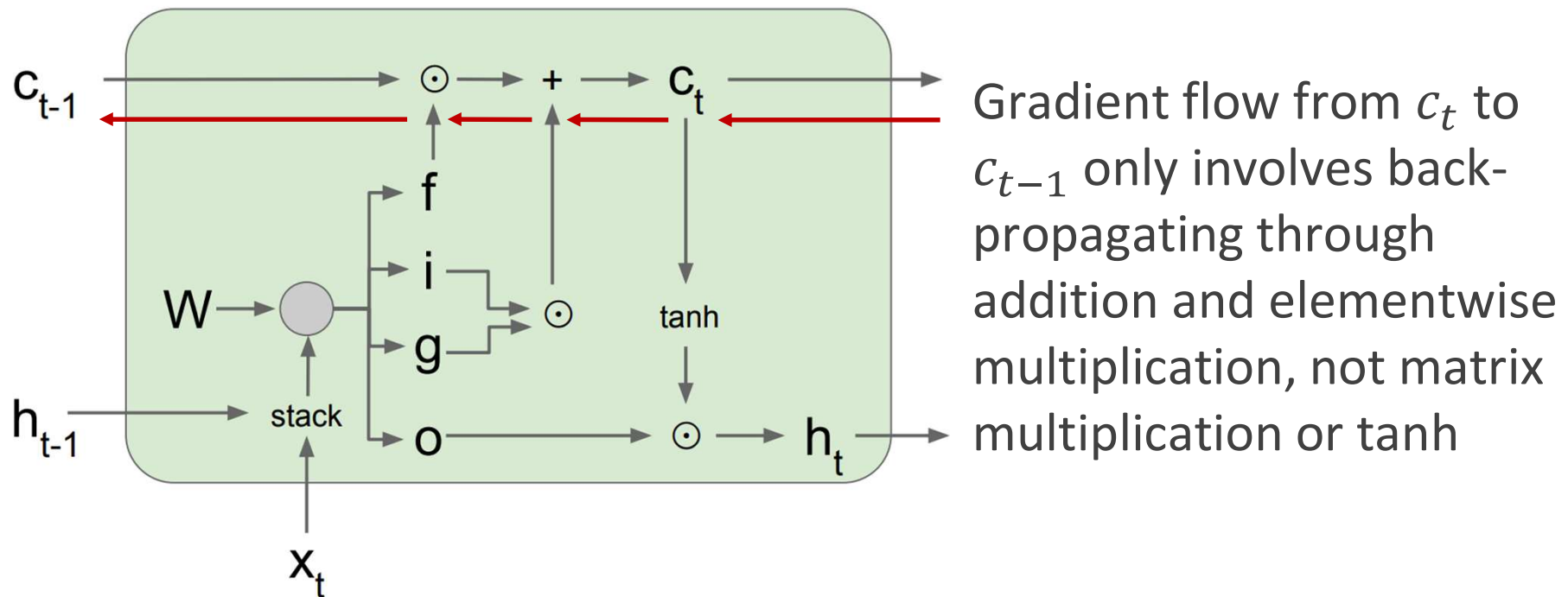


$$\begin{pmatrix} g_t \\ i_t \\ f_t \\ o_t \end{pmatrix} = \begin{pmatrix} \tanh \\ \sigma \\ \sigma \\ \sigma \end{pmatrix} \begin{pmatrix} W_g \\ W_i \\ W_f \\ W_o \end{pmatrix} \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

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

$$h_t = o_t \odot \tanh c_t$$

LSTM Backward Pass



For complete details: [Illustrated LSTM Forward and Backward Pass](#)

Figure source

- Next Week:
 - Continue: RNN – Gated Recurrent Unit
 - different applications and architecture of RNN
 - Manifold learning and Autoencoders