

# Deep Learning and Applications

DSA 5204 • Lecture 5  
Dr Low Yi Rui (Aaron)  
Department of Mathematics



**NUS**  
National University  
of Singapore

# Last Time

**We introduced CNN to handle image data. The main idea is to exploit the spatial structure of images**

- Extract features using the convolution operation
- Weight sharing and sparse connectivity
- Equivariance and invariance

**Today we will discuss another class of neural network to handle temporal data: recurrent neural networks**



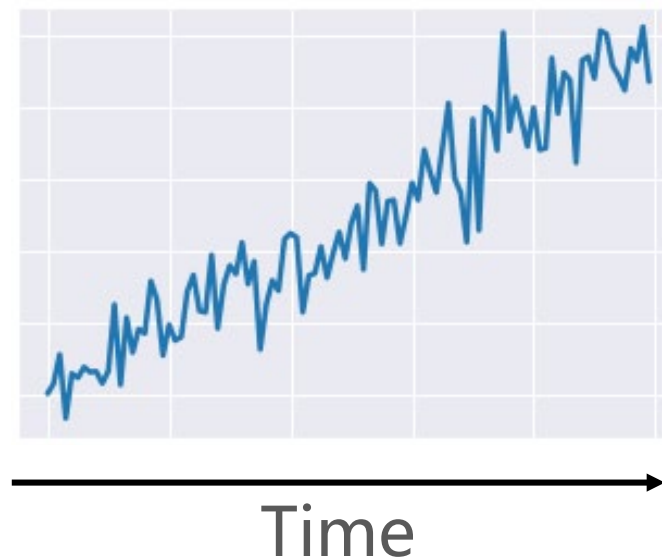
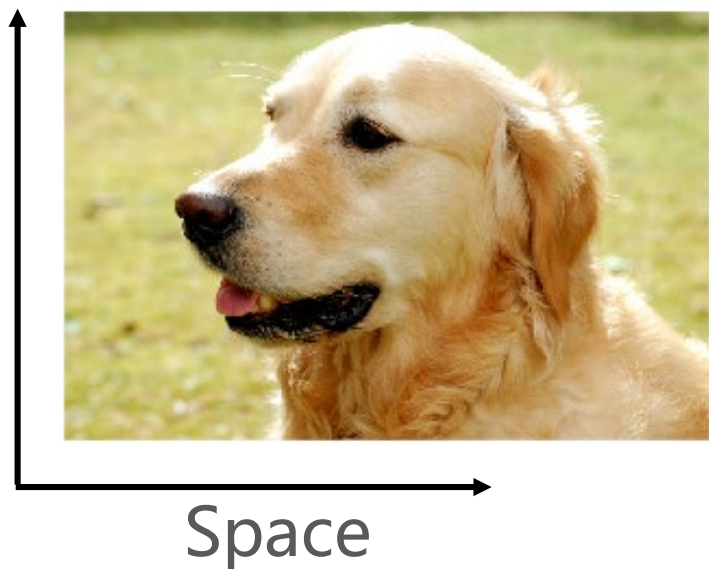
# Time series modelling

# What is a time series?



**Like images, a time series is another type of data with some local structure.**

**Instead of spatial structure in images, we consider temporal structures.**





**We may represent a time series as either**

- A discrete sequence

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \dots$$

- A continuous sequence

$$\{\mathbf{x}^{(t)} : t \geq 0\}$$

**Each slice of the time series,  $\mathbf{x}^{(t)}$  is a vector, or something that can be represented as a vector**

# Examples of Time Series Data

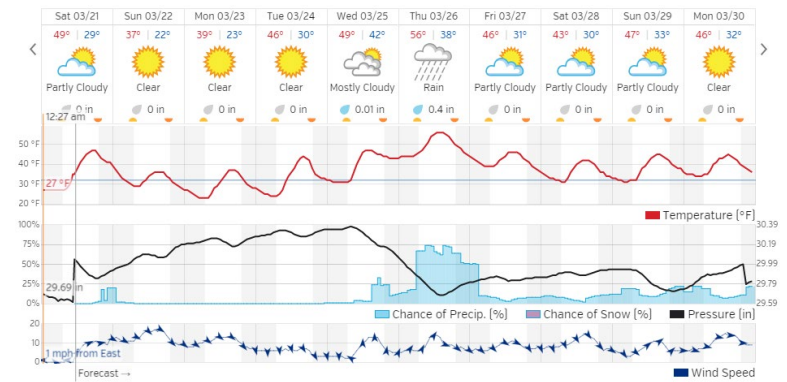
The honorable favour that yo<sup>r</sup> I<sup>st</sup> had afforded me in allowing me the liberty of myne own chamber, hath given me leave so much to respect and love my self, that now I am desirous to be well. And therefore for health and pleasure of yo<sup>r</sup> I<sup>st</sup> I desire to be so much more to shew my better that will be by yo<sup>r</sup> I<sup>st</sup> in our more own house, and surely so I may leave some provision and satisfaction, yo<sup>r</sup> I<sup>st</sup> shall work yo<sup>r</sup> I<sup>st</sup> grant me liberty to take the space about the house. The whole world is a strange improvement to us, which I am bound yo<sup>r</sup> I<sup>st</sup> right, but this favour may lengthen and better my life, so I desire to shew, only to help to reduce by my disease, and desire to do yo<sup>r</sup> I<sup>st</sup> service, my life is paid. All mighty god dwell ever in yo<sup>r</sup> I<sup>st</sup> heart and fill yo<sup>r</sup> I<sup>st</sup> good desires and grant them.

— J. B. L.

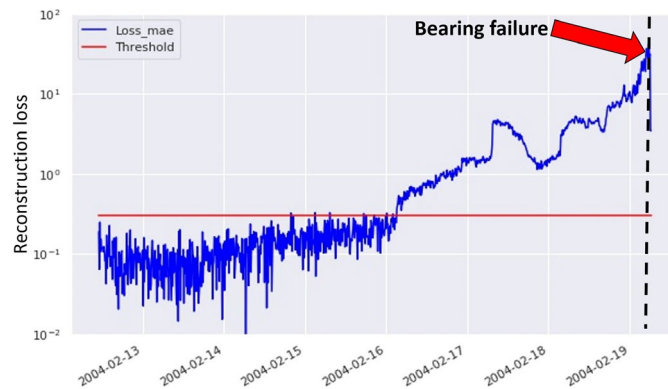
Text



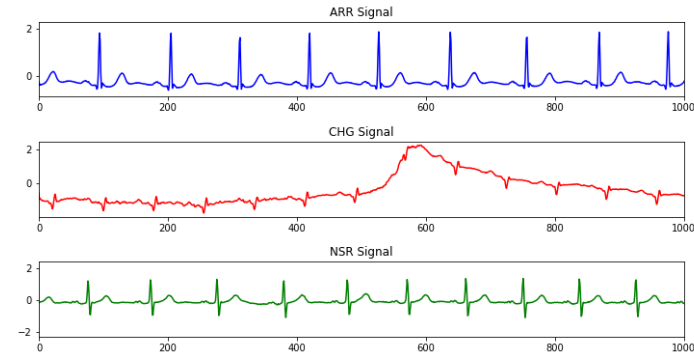
Stock Price



Weather



Machine Health



ECG Signals

# Task I: Sequence Prediction

Given a historical time series data:

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)}$$

Task is to predict

- The next data point  $\mathbf{x}^{(t+1)}$
- A sequence of next data points  $\mathbf{x}^{(t+1)}, \mathbf{x}^{(t+2)}, \dots, \mathbf{x}^{(t+\tau)}$

**Example applications**

- Stock price prediction, weather forecasting

# Task II: Sequence Classification/Regression

Given a time series data:

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)}$$

Task is to predict a label  $y$  of this sequence, which can be either discrete (classification) or continuous (regression)

## Example Applications

- Credit card transaction fraud detection
- Heart arrhythmia detection



# Task III: Sequence-to-sequence Modelling

Given a time series data:

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots \mathbf{x}^{(t)}$$

Task is to predict a corresponding time series

$$\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots \mathbf{y}^{(t)}$$

## Example Applications

- Machine translation
- Continuous health monitoring using wearable devices

# Task IV: Sequence Generation



**Given a seed sequence**

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(\tau)}$$

**Task is to generate a longer sequence starting from this, i.e.**

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(\tau)}, \mathbf{x}^{(\tau+1)}, \mathbf{x}^{(\tau+2)}, \dots$$

**according to some distribution**

**Example applications**

- Writing poems
- Composing music

# The Supervised Learning Problem

Recall that in supervised learning, we define the inputs, outputs and the target function that maps the former to the latter.

For time series modelling, we can define these similarly

- Inputs:  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \dots, \mathbf{x}^{(\tau)}$
- Outputs:  $\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \mathbf{y}^{(3)}, \dots, \mathbf{y}^{(\tau)}$
- Target:  $\{F_t^*\}$  with
$$\mathbf{y}^{(t)} = F_t^*(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(\tau)})$$
- Goal: Learn some  $\{\hat{F}_t\}$  to approximate  $\{F_t^*\}$

# Examples

- **Task I: Sequence Prediction**

$$F_t^*(x^{(1)}, x^{(2)}, \dots, x^{(\tau)}) = F_t^*(x^{(1)}, x^{(2)}, \dots, x^{(t)}) = x^{(t+1)}$$

- **Task II: Sequence Classification/Regression**

$$F_\tau^*(x^{(1)}, x^{(2)}, \dots, x^{(\tau)}) = y^{(\tau)} = y$$

- **Task III: Sequence-to-sequence**

$$F_t^*(x^{(1)}, x^{(2)}, \dots, x^{(\tau)}) = y^{(t)}, \quad t = 1, 2, \dots$$

or

$$F_t^*(x^{(1)}, x^{(2)}, \dots, x^{(t)}) = y^{(t)}, \quad t = 1, 2, \dots$$



# Recurrent Neural Networks

# Sharing Parameters in Time



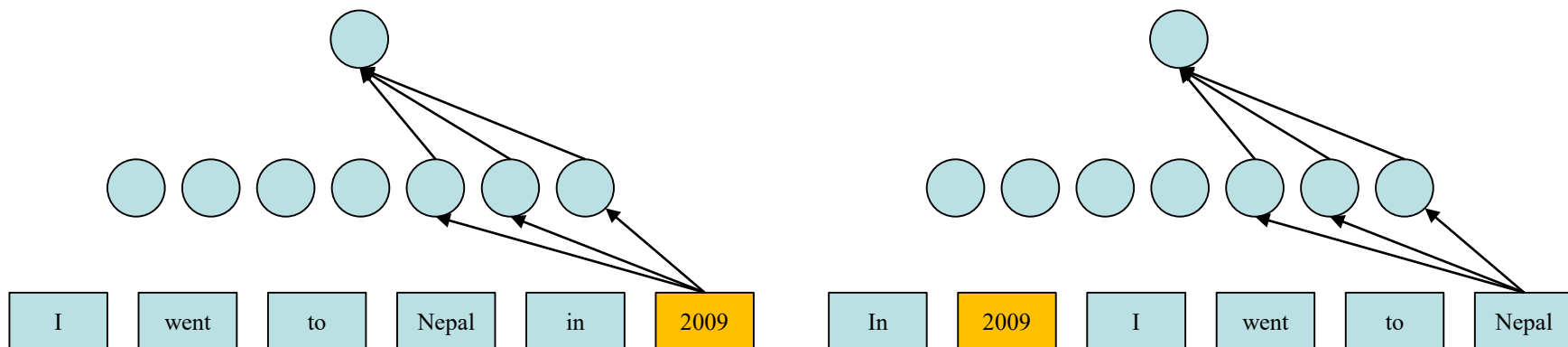
Suppose we want to extract information from two sentences

*"I went to Nepal in 2009"*

*"In 2009 I went to Nepal"*

We want a model to behave similarly for these inputs.

Suppose we use a FCNN for this...



# Dynamical System



A natural way to share parameters in time is to define a dynamical system depending on some parameter:

$$s^{(t+1)} = f(s^{(t)}; \theta)$$

- The vector  $s^{(t)}$  is the state of the dynamical system at time  $t$
- $\theta$  is a vector of parameters
- This is different from “feed-forward” dynamics

$$h^{(t+1)} = f(h^{(t)}; \theta^{(t)})$$

i.e. the parameter  $\theta$  is shared in time.

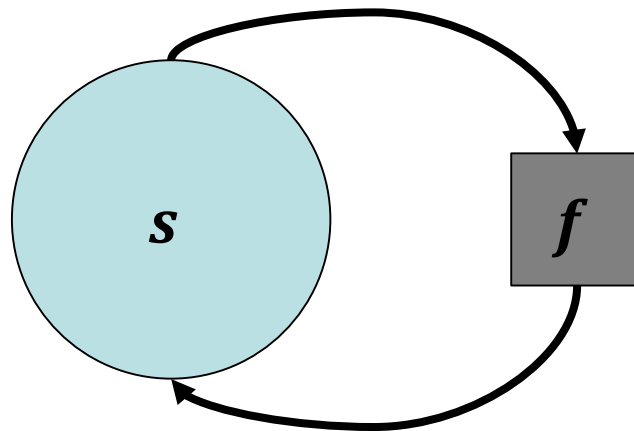
# Computational Graph Representation



We can represent the dynamics

$$s^{(t+1)} = f(s^{(t)}; \theta)$$

using the usual computational graph approach.



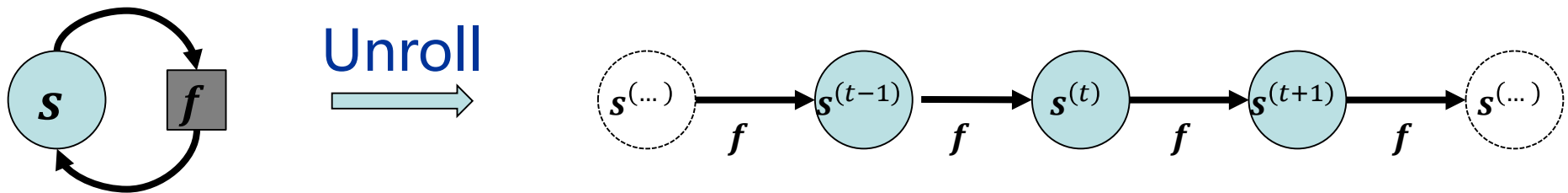


# Unrolling the Dynamics



Of course, we can also unroll the dynamics

$$\begin{aligned} s^{(t)} &= f(s^{(t-1)}; \theta) = f(f(s^{(t-2)}; \theta); \theta) \\ &= \dots = f(f(\dots f(s^{(0)}; \theta) \dots; \theta); \theta) \end{aligned}$$

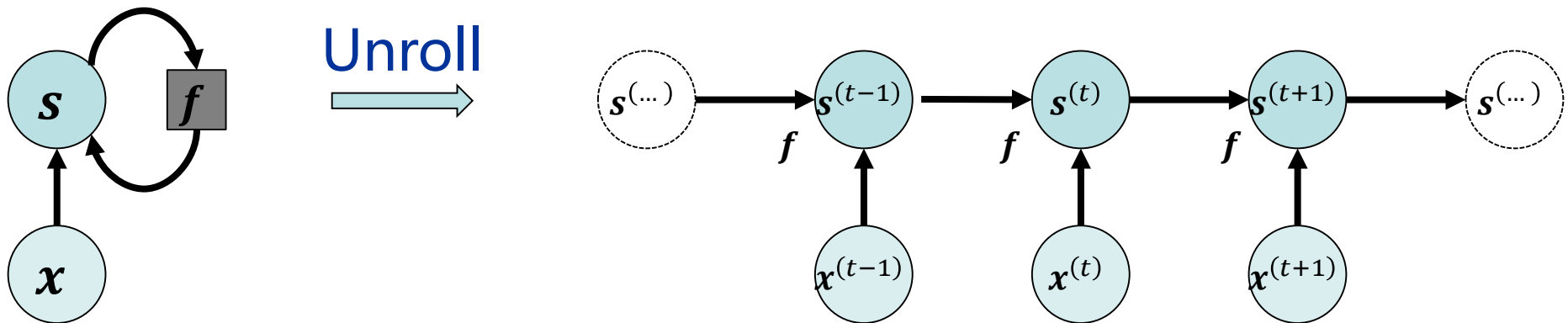


# Dynamical Systems with Inputs



Now, we consider a dynamical system with inputs or external forcing

$$s^{(t)} = f(s^{(t-1)}, x^{(t)}, \theta)$$



# Recurrent Neural Networks

A decorative network diagram in the top right corner, consisting of a series of interconnected nodes (small circles) and edges (thin lines), forming a complex web-like structure.

The basic architecture of recurrent neural networks uses a forced, hidden dynamical system as a basic hypothesis space

$$\begin{aligned}h^{(t)} &= f(h^{(t-1)}, x^{(t)}, \theta) \\ \hat{y}^{(t)} &= g(h^{(t)}, \phi)\end{aligned}$$

This implicitly parameterizes

$$\hat{F}_t(x^{(1)}, x^{(2)}, \dots, x^{(t)}; \theta, \phi) = \hat{y}^{(t)}$$

# Simple RNN



To be explicit, we consider the following simple RNN structure:

$$\mathbf{h}^{(t)} = \sigma_r(W\mathbf{h}^{(t-1)} + U\mathbf{x}^{(t)} + \mathbf{b})$$

$$\hat{\mathbf{y}}^{(t)} = \mathbf{o}^{(t)} = \sigma_o(V\mathbf{h}^{(t)} + \mathbf{c})$$

- The trainable parameters are  $(\theta, \phi) = (W, U, \mathbf{b}, V, \mathbf{c})$
- The recurrent activation  $\sigma_r$  is usually taken to be **tanh**
- The output activation function  $\sigma_o$  depends on application
- This is the Elman variant of RNN. The Jordan variant replaces  $\mathbf{h}^{t-1}$  by  $\mathbf{o}^{t-1}$  in the first equation

# Loss Functions

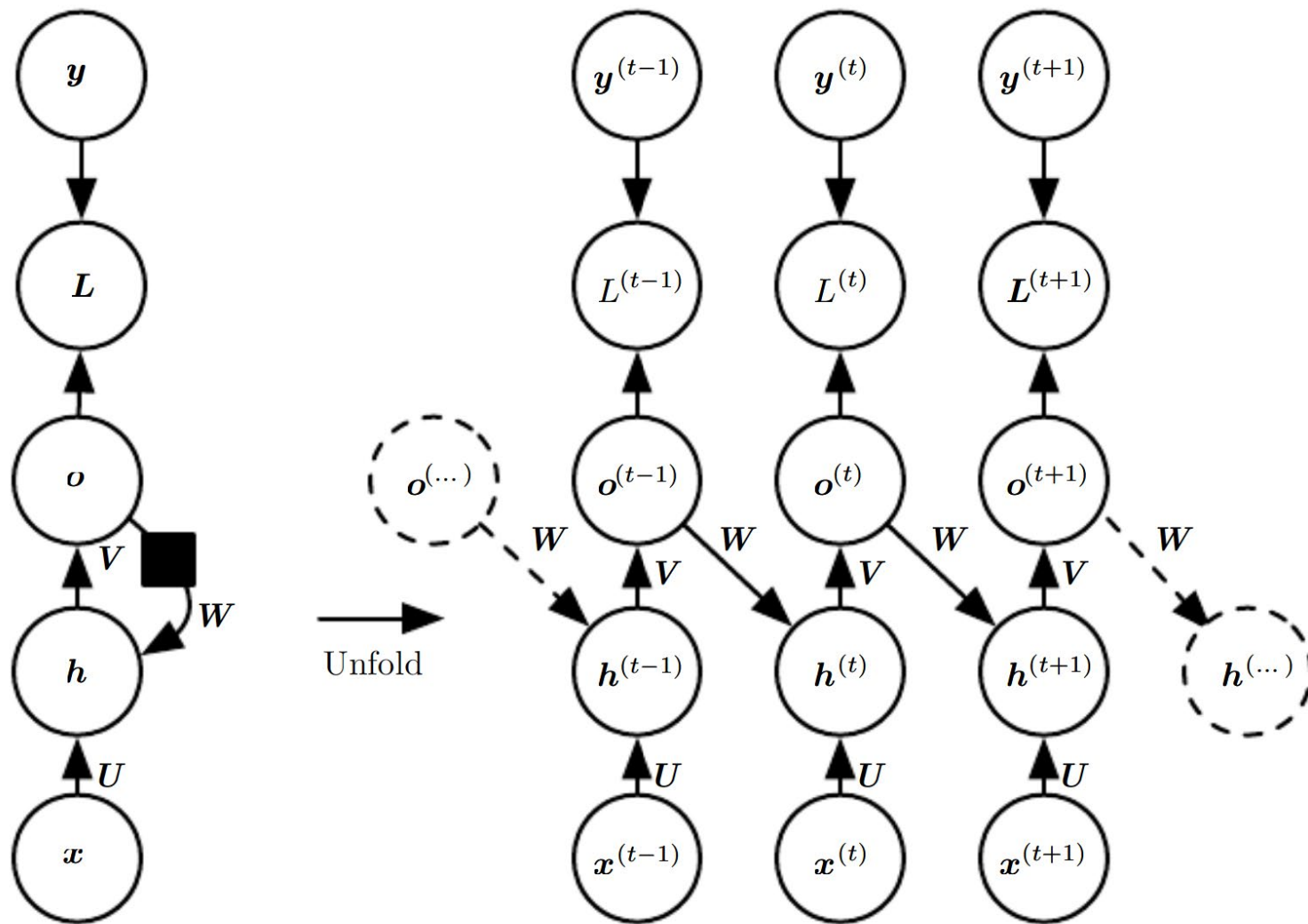
For single-prediction tasks (Tasks I and II), the loss function may be defined only at the end of the sequence, e.g.

$$L(\mathbf{y}, \hat{\mathbf{y}}^{(\tau)}) = \frac{1}{2} \|\mathbf{y} - \hat{\mathbf{y}}^{(\tau)}\|^2$$

For sequence-prediction tasks (Task III), we may take the sum

$$\sum_{t=1}^{\tau} L(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = \frac{1}{2} \sum_{t=1}^{\tau} \|\mathbf{y}^{(t)} - \hat{\mathbf{y}}^{(t)}\|^2$$

# Computational Graphs (Jordan Variant)



# Why a hidden dynamical system?

Consider instead the naive approach of just using FCNN to map each input to each output

$$\hat{y}^{(t)} = \text{FCNN}(x^{(t)})$$

## What is wrong?

- The prediction at time  $t$  only depends on  $x^{(t)}$
- This cannot model systems with memory, no matter how complex FCNN is

# Example: Hidden States

Consider generating from  $\{x^{(t)}, t \geq 1\}$  the outputs

$$y^{(t)} = x^{(t)} + x^{(t-1)} + x^{(t-2)}, \quad t \geq 1$$

where we define  $x^{(0)} = x^{(-1)} = 0$ .

Then, it is obvious that we cannot predict the value of  $y^{(t)}$  from just  $x^{(t)}$

To account for memory, we can form the linear model

$$\hat{y}^{(t)} = \sum_{s=1}^t a^{(s)} x^{(s)}$$

which can learn our system with  $a^{(t)} = a^{(t-1)} = a^{(t-2)} = 1$  and  $a^{(s)} = 0$  for  $s \leq t - 3$





**However, there are issues with this approach. The model**

$$\hat{y}^{(t)} = \sum_{s=1}^t a^{(s)} x^{(s)}$$

**operates on variable-length inputs. Also, it is hard to generalize to non-linear models.**

**Alternative approach: using hidden states**

$$\begin{aligned} h_1^{(t)} &= x^{(t)} \\ h_2^{(t)} &= h_1^{(t-1)} \\ h_3^{(t)} &= h_2^{(t-1)} \\ \hat{y}^{(t)} &= h_1^{(t)} + h_2^{(t)} + h_3^{(t)} \end{aligned}$$

**Note: this is simply a linear RNN!**

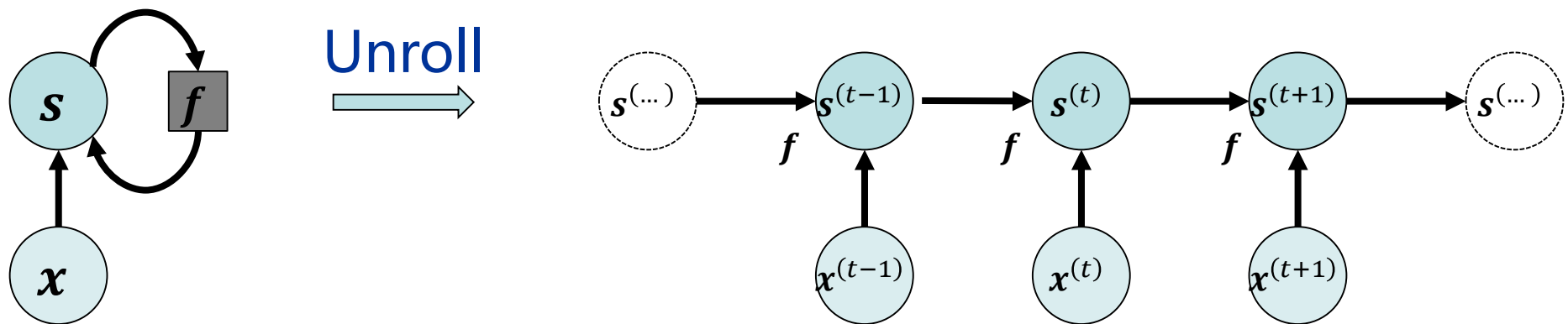


# Training Recurrent Neural Networks

# Unrolling Computational Graph



Once we unroll the computational graph, then RNN is just a parameter-tied feed-forward NN, thus we can apply back-propagation algorithm



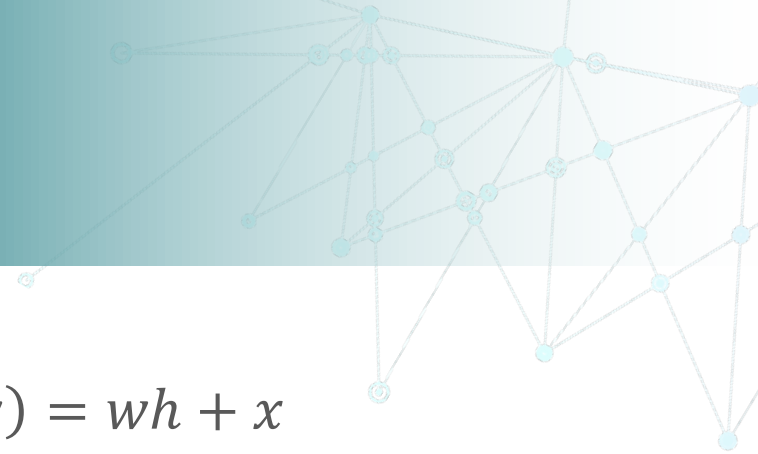
# Example: Regression on Linear RNN

Let us take a linear RNN in 1D with no bias

$$h^{(t)} = wh^{(t-1)} + x^{(t)}, \quad h^{(0)} = 0$$

$$\hat{y}^{(t)} = h^{(t)}$$

and terminal loss  $L(y, \hat{y}^{(\tau)})$ .



**Observe we can write**

$$h^{(t)} = H(h^{(t-1)}, x^{(t)}; w), \quad H(h, x, w) = wh + x$$

**Then, we can write**

$$\frac{dL}{dw} = \sum_{t=1}^{\tau} \frac{dL}{dh^{(t)}} \frac{\partial H}{\partial w}(h^{(t-1)}, x^{(t)}, w) = \sum_{t=1}^{\tau} \frac{dL}{dh^{(t)}} h^{(t-1)}$$

**Then, it remains to compute**

$$p^{(t)} := \frac{dL}{dh^{(t)}}$$

**which is easily shown to obey the backward recursion**

$$p^{(t-1)} = wp^{(t)}, \quad p^{(\tau)} = \frac{\partial L}{\partial \hat{y}}(y, \hat{y}^{(\tau)})$$

**This gives**

$$\frac{dL}{dw} = \frac{\partial L}{\partial \hat{y}}(y, \hat{y}^{(\tau)}) \sum_{t=1}^{\tau} w^{\tau-t} h^{(t-1)}$$

# Gradient Explosion and Vanishing

A decorative network diagram in the top right corner, consisting of a series of interconnected nodes (small circles) and edges (thin lines), forming a complex web-like structure.

**A big problem that plagues RNNs is the difficulty to train them.**

**One is often faced with a dilemma for large  $\tau$ :**

- Gradient explosion: for some parameter values, the gradients computed using the RNN diverges
- Gradient vanishing: for other parameter values, the gradients computed using the RNN gives vanishing weight to faraway inputs

# Example: Gradient Explosion/Vanishing

Let us go back to our previous example, where we computed

$$\frac{dL}{dw} = \sum_{t=1}^{\tau} w^{\tau-t} h^{(t-1)}$$

With some simplifications, we can show that

$$\frac{dL}{dw} = \sum_{t=1}^{\tau} (\tau - t) w^{\tau-t-1} x^{(t)}$$

- Now, if  $w > 1$ , gradients explode for large  $\tau$
- If  $w < 1$ , gradients do not explode, but attaches vanishing weights to earlier inputs,  
i.e.  $(\tau - 1)w^{\tau-2}x^{(1)} \ll 1$  for large  $\tau$



# Gated Recurrent Neural Networks



# Gating and long-term dependencies

A decorative network diagram in the top right corner, consisting of numerous light blue nodes connected by thin, light blue lines, forming a complex web-like structure.

**One of the most effective ways to overcome the inability for RNNs to learn long-term dependencies is their gated extensions**

## The key ideas

- Design RNN variants where gradients neither vanish nor explode
- Achieve this by the construction of “trainable gates”, which controls the flow of information
  - **Gates control the accumulation of information**
  - **When accumulated information is no-longer required, we just forget them**

# Long Short Term Memory (LSTM)

A decorative network diagram in the top right corner, consisting of a series of interconnected nodes and lines, resembling a neural network or a complex graph structure.

LSTM [Hochreiter and Schmidhuber, 1997] is one of the most commonly used gated RNNs.

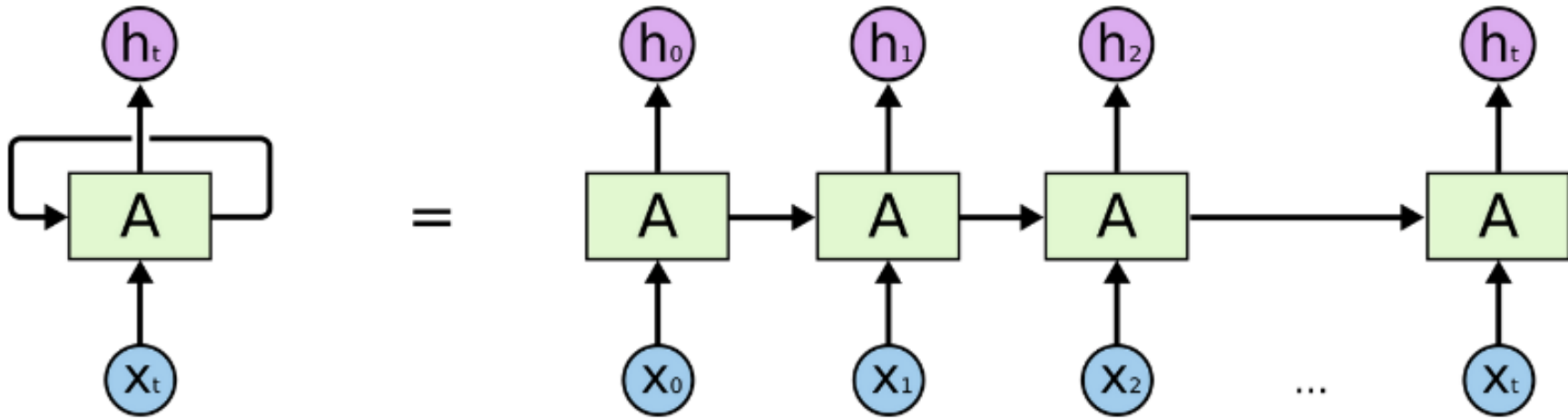
We will use a series of illustrations taken from <https://colah.github.io/posts/2015-08-Understanding-LSTMs/> to illustrate the main innovations of LSTM, as well as introducing the abstraction of RNN cells

# RNN Cells

RNNs operate by feeding the input and previous hidden states to a function to give the next hidden state, e.g.

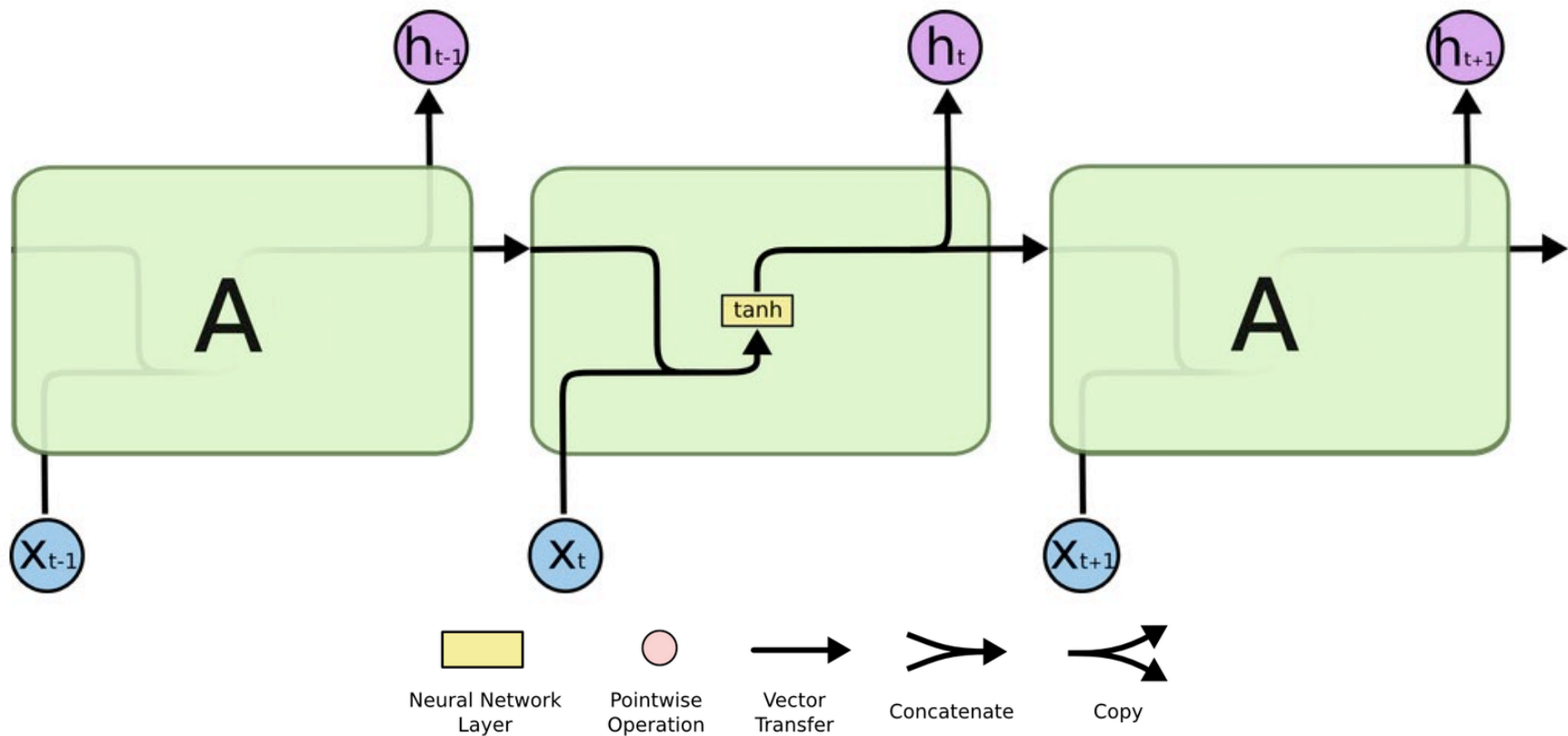
$$\mathbf{h}^{(t)} = \tanh(\mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)} + \mathbf{b})$$

We can abstract this operation as a black-box **cell A**

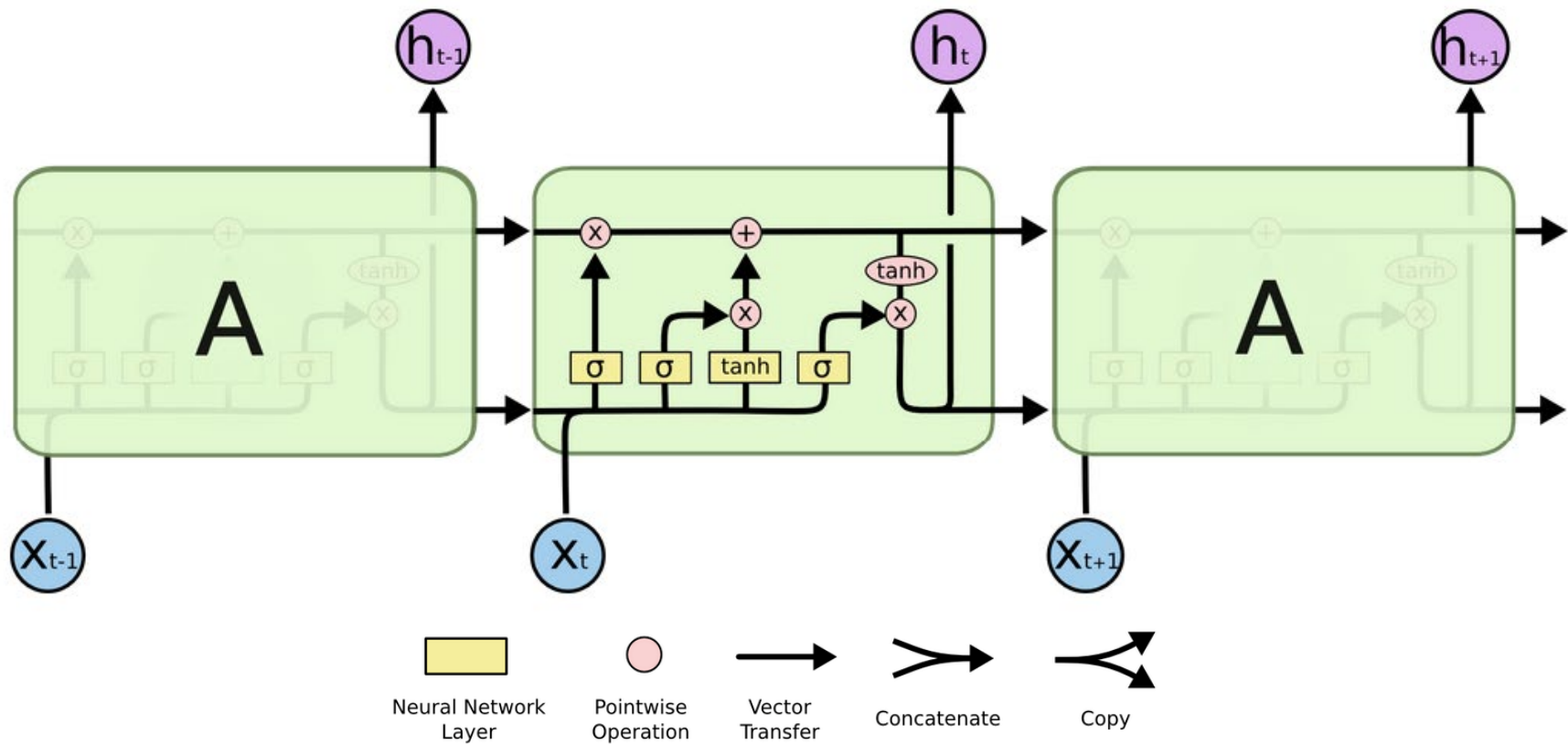


# Inside a Simple RNN Cell

$$\mathbf{h}^{(t)} = \tanh(W\mathbf{h}^{(t-1)} + U\mathbf{x}^{(t)} + \mathbf{b})$$

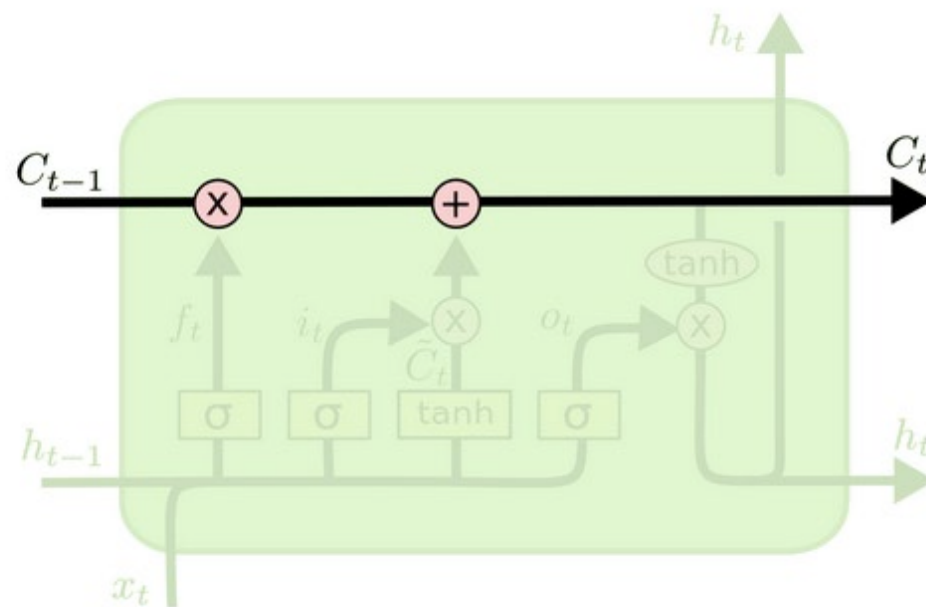


# Inside a LSTM Cell



# The Cell State

The cell state  $C_t$  is another hidden variable that is designed to flow through time with minimal interruptions



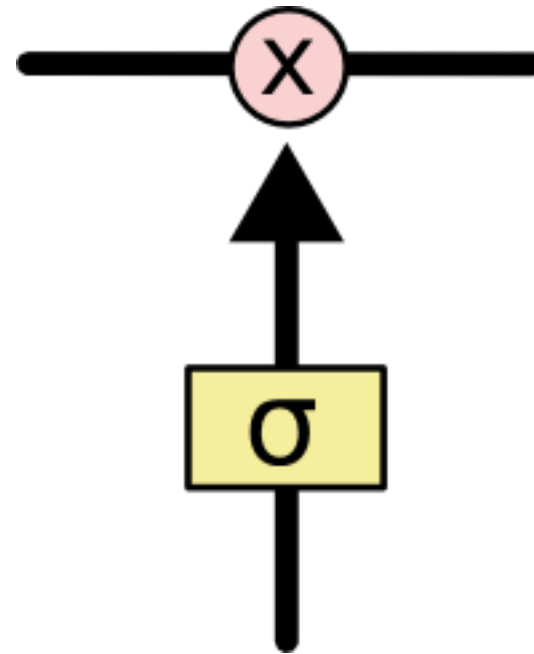
# Gates

The LSTM interacts with the cell state using gates, there are just layers with sigmoid outputs

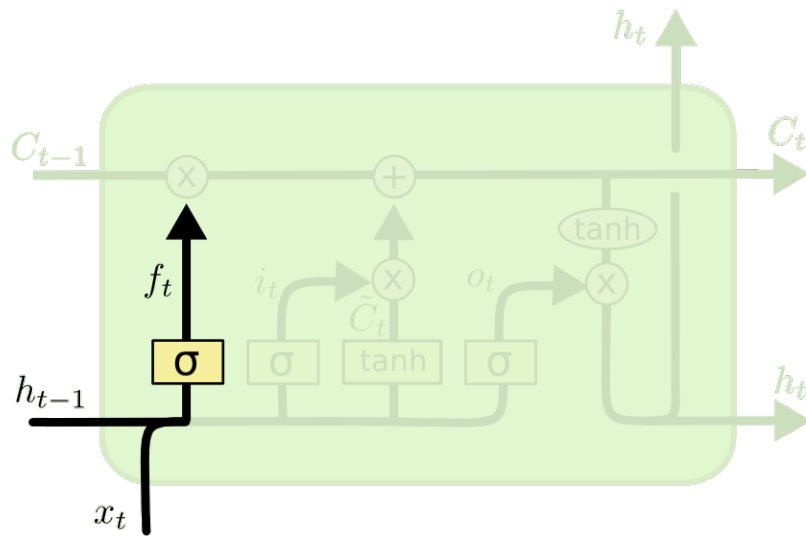
$$u \mapsto \sigma(Wu + b),$$

$$\sigma(z) = \frac{1}{1 + e^{-z}} \in (0,1)$$

Multiplying the cell state with the above makes a modification to it



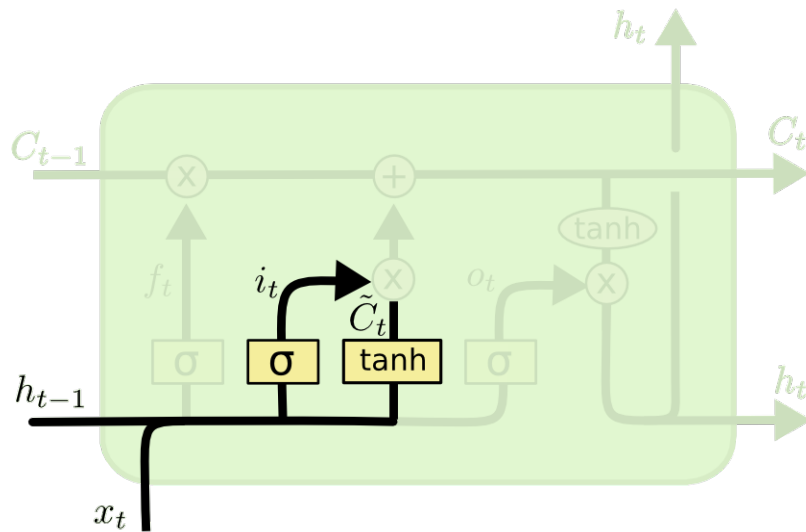
# LSTM Step 1



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



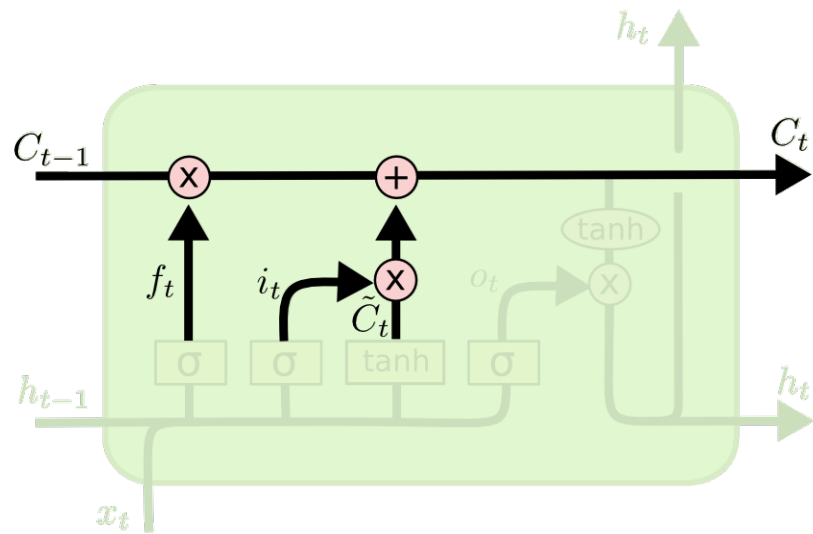
# LSTM Step 2



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

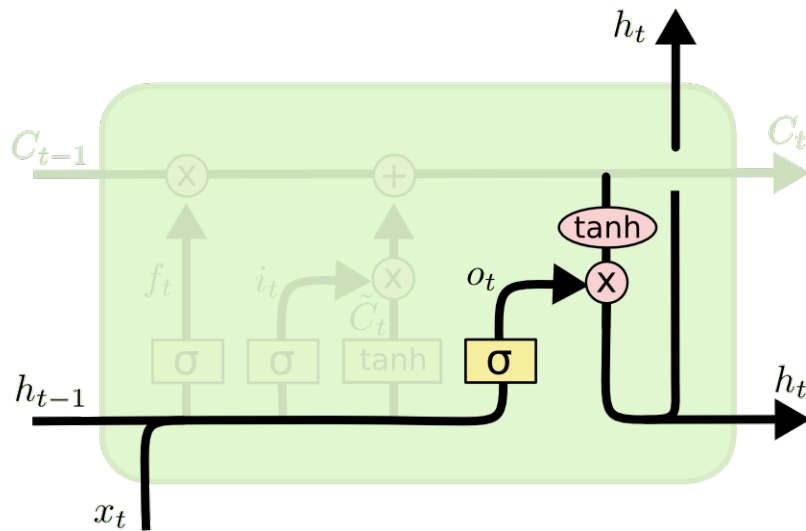
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

# LSTM Step 3



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

# LSTM Step 4



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

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

# Summary: LSTM Model (In our Notation)

$$\begin{aligned}f^{(t)} &= \sigma(W_f \mathbf{h}^{(t-1)} + U_f \mathbf{x}^{(t)} + \mathbf{b}_f) \\i^{(t)} &= \sigma(W_i \mathbf{h}^{(t-1)} + U_i \mathbf{x}^{(t)} + \mathbf{b}_i) \\o^{(t)} &= \sigma(W_o \mathbf{h}^{(t-1)} + U_o \mathbf{x}^{(t)} + \mathbf{b}_o) \\\mathbf{C}^{(t)} &= f_t \cdot \mathbf{C}_{t-1} + i_t \cdot \sigma_c(W_c \mathbf{h}^{(t-1)} + U_c \mathbf{x}^{(t)} + \mathbf{b}_c) \\\mathbf{h}^{(t)} &= o^{(t)} \cdot \sigma_h(\mathbf{C}^{(t)})\end{aligned}$$

**Gates:**  $f^{(t)}, i^{(t)}, o^{(t)}$

**Cell state:**  $\mathbf{C}^{(t)}$

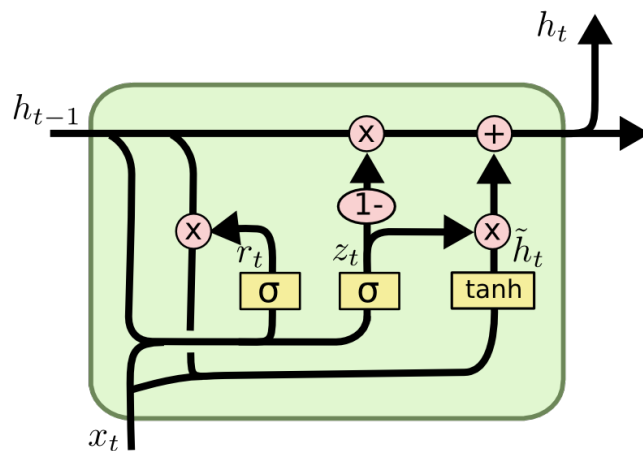
**Hidden/RNN state:**  $\mathbf{h}^{(t)}$

# Variants of the LSTM architecture



## There are many variants

- Peephole connections
- Coupled input/forget gates
- Gated Recurrent Unit (GRU)



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

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

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

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



# Deep RNNs and Other Architectures

# Deep Recurrent Neural Networks

A decorative network diagram in the top right corner, featuring a series of interconnected nodes (represented by small blue circles) and edges (represented by thin grey lines), forming a complex web-like structure.

**One may observe that the RNN, even in its simplest form, is already deep in the time direction**

**However, each time-step relied on a single layer of FCNN-like structure**

**Deep RNNs generalize this by using a deep NN for the recurrent step**

# The Basic Architecture

## Shallow/Simple RNN:

$$\begin{aligned}\mathbf{h}^{(t)} &= \sigma_r(W\mathbf{h}^{(t-1)} + U\mathbf{x}^{(t)} + \mathbf{b}) \\ \hat{\mathbf{y}}^{(t)} &= \sigma_o(V\mathbf{h}^{(t)} + \mathbf{c})\end{aligned}$$

## Deep RNN (One variant):

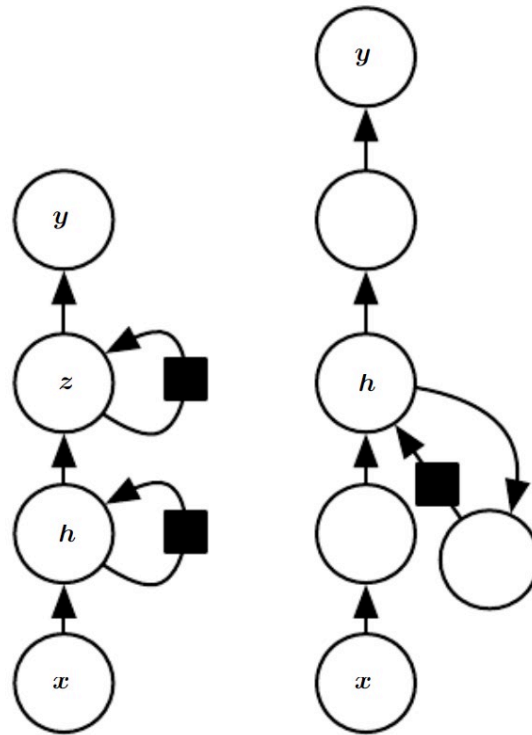
$$\begin{aligned}\mathbf{h}^{(t)} &= \sigma_r(W\mathbf{h}^{(t-1)} + U\mathbf{x}^{(t)} + \mathbf{b}) \\ \mathbf{z}^{(t)} &= \sigma_r(W\mathbf{z}^{(t-1)} + U\mathbf{h}^{(t)} + \mathbf{b}) \\ \hat{\mathbf{y}}^{(t)} &= \sigma_o(V\mathbf{z}^{(t)} + \mathbf{c})\end{aligned}$$



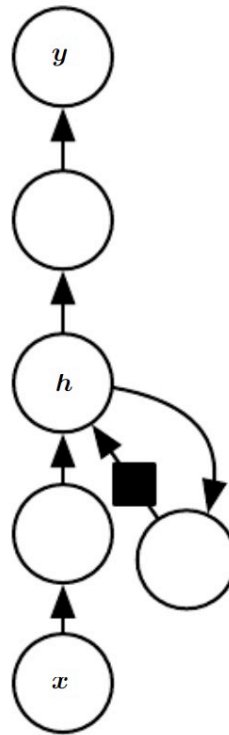
# Other Ways to obtain Deep RNN models



(a) Using Stacked Hidden/Cell States



(b) Using MLPs for Connections



# Other Important Architectures For Sequence Modelling



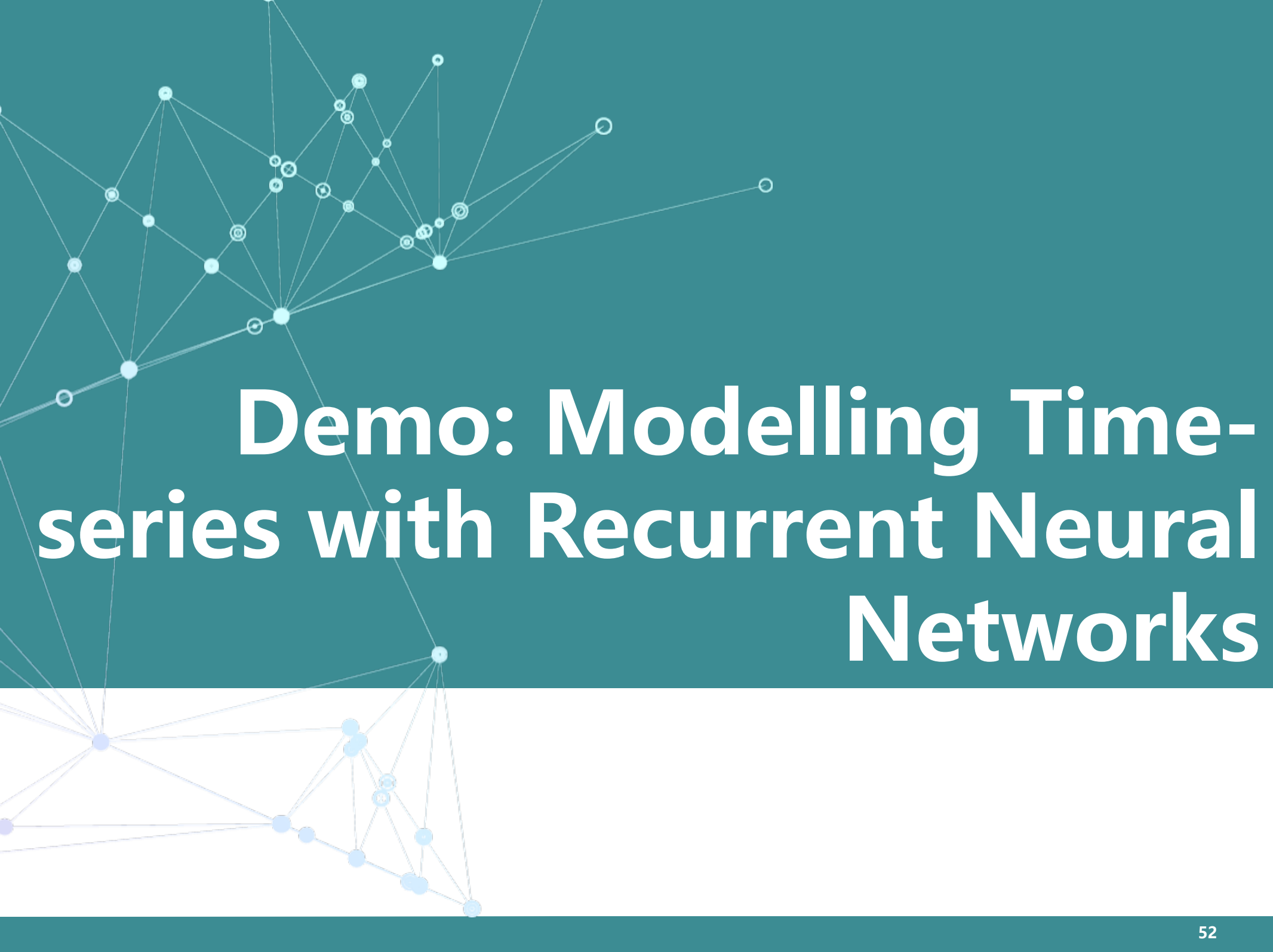
There are a number of other architectures for time-series modelling that are useful to know

Look at future, e.g translation.  
Non-eg price prediction

- Bi-directional RNNs
- Sequence-to-sequence autoencoders
- Attention mechanisms and the transformer network

## More information:

- Deep learning book, Chapter 10
- Graves, A. (2012). Supervised Sequence Labelling with Recurrent Neural Networks . Studies in Computational Intelligence.



# Demo: Modelling Time-series with Recurrent Neural Networks

# Summary

A decorative network graph in the top right corner, consisting of numerous small blue circular nodes connected by thin, light blue lines, forming a complex web-like structure.

- **We introduced recurrent neural networks as another type of parameter-sharing network for handling time series data**
- **Primary idea: incorporate memory into a hidden, forced dynamical system**
- **Extensions:**
  - Gated RNNs
  - Deep RNNs