

# *Deep Learning*

*A course about theory & practice*



## Recurrent Neural Networks

Marco Piastra

# *The Basic Architecture*

# Recurrent Neural Networks

- **Feed-forward neural network**

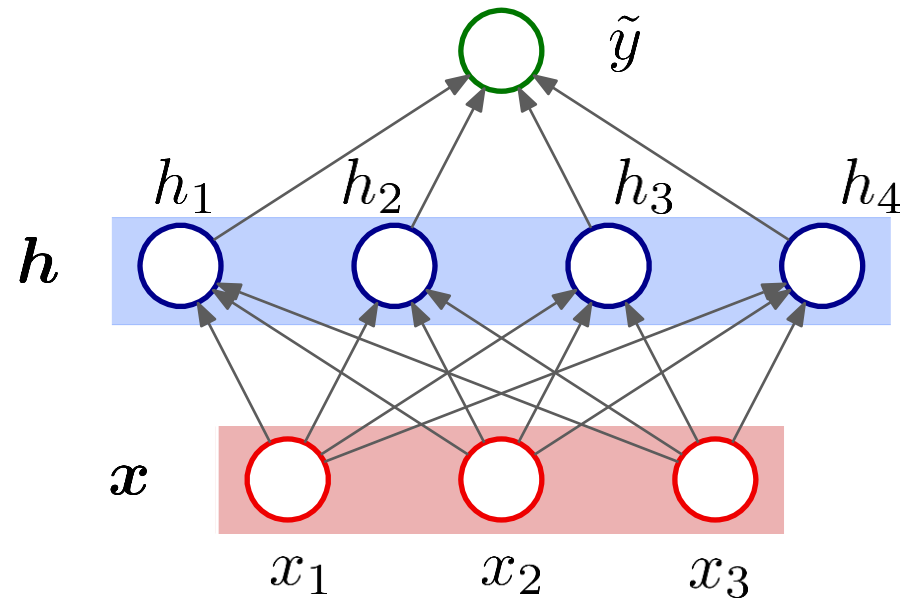
$$\tilde{y} = w \cdot h + b$$

where:

$$h := g(Wx + b)$$

hidden

input



# Recurrent Neural Networks

- **Feed-forward neural network**

$$\tilde{y} = \mathbf{w} \cdot \mathbf{h} + b$$

where:

$$\mathbf{h} := g(\mathbf{W}\mathbf{x} + \mathbf{b})$$

hidden                      input

- **Recurrent Neural Network**

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

where:

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

hidden at t                      input at t                      hidden at t-1

*The idea is to make the network output depend on the past 'history'*

# Recurrent Neural Networks

## Recurrent Neural Network

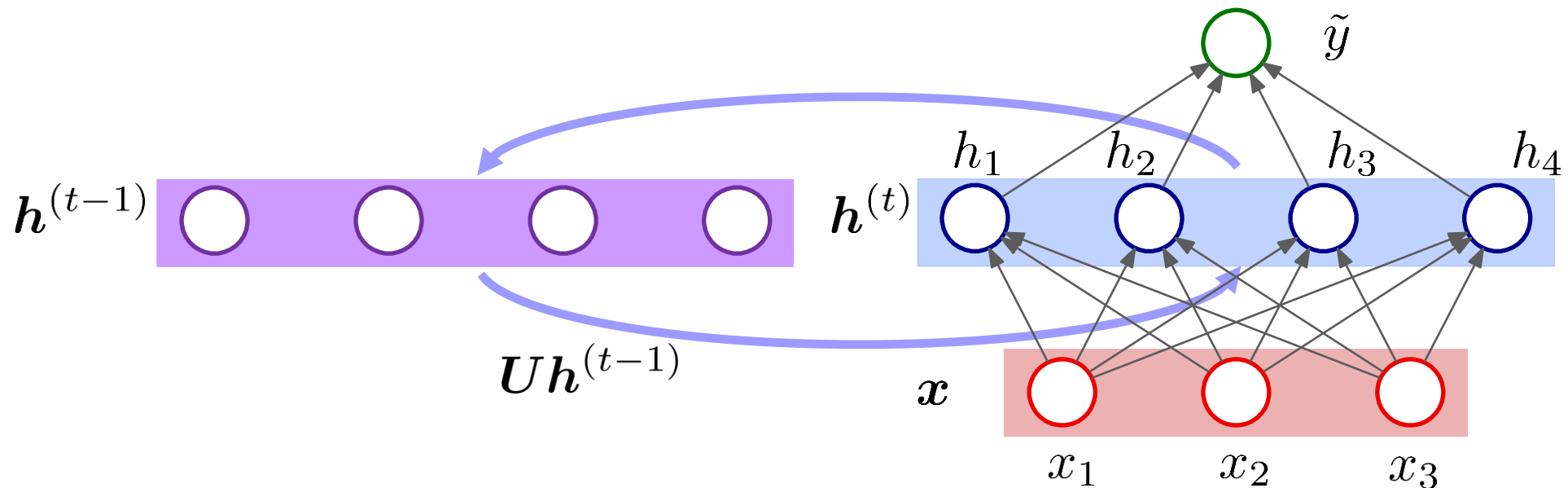
$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

where:

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

hidden at t                  input at t                  hidden at t-1

The idea is to make the network output depend on the past 'history'



# Recurrent Neural Networks

- **Recurrent Neural Network**

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

where:

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

hidden at t      input at t      hidden at t-1

# Recurrent Neural Networks

## Recurrent Neural Network

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

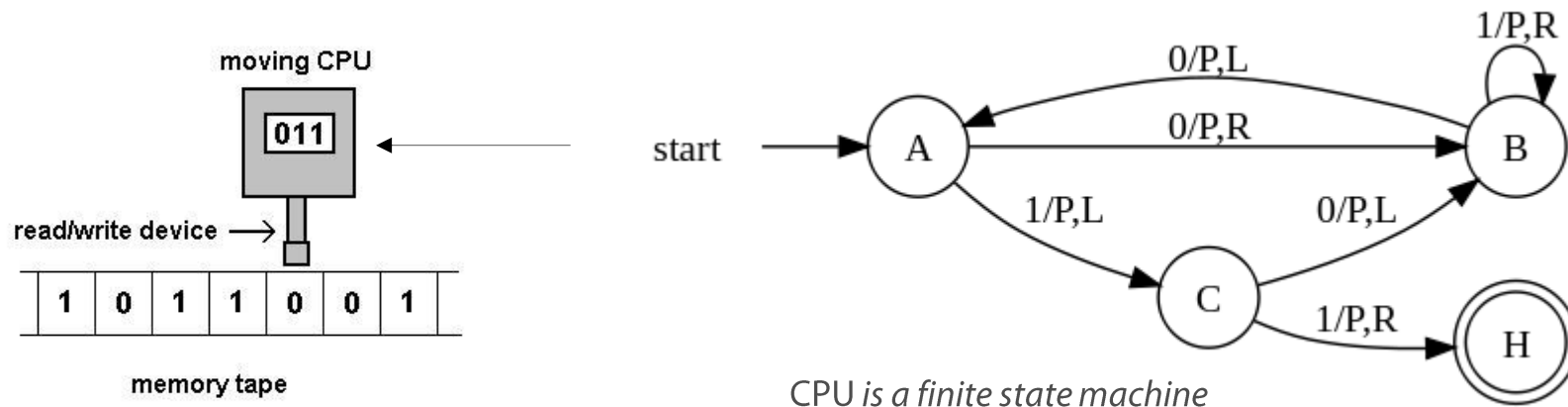
where:

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

hidden at t                  input at t                  hidden at t-1

## Computational power of RNNs (Siegelmann & Sontag, 1992)

“RNNs can simulate any Turing machine”



They are *much harder* to **train** than FF neural networks

There is no known general method to do this

# Recurrent Neural Networks

## Recurrent Neural Network

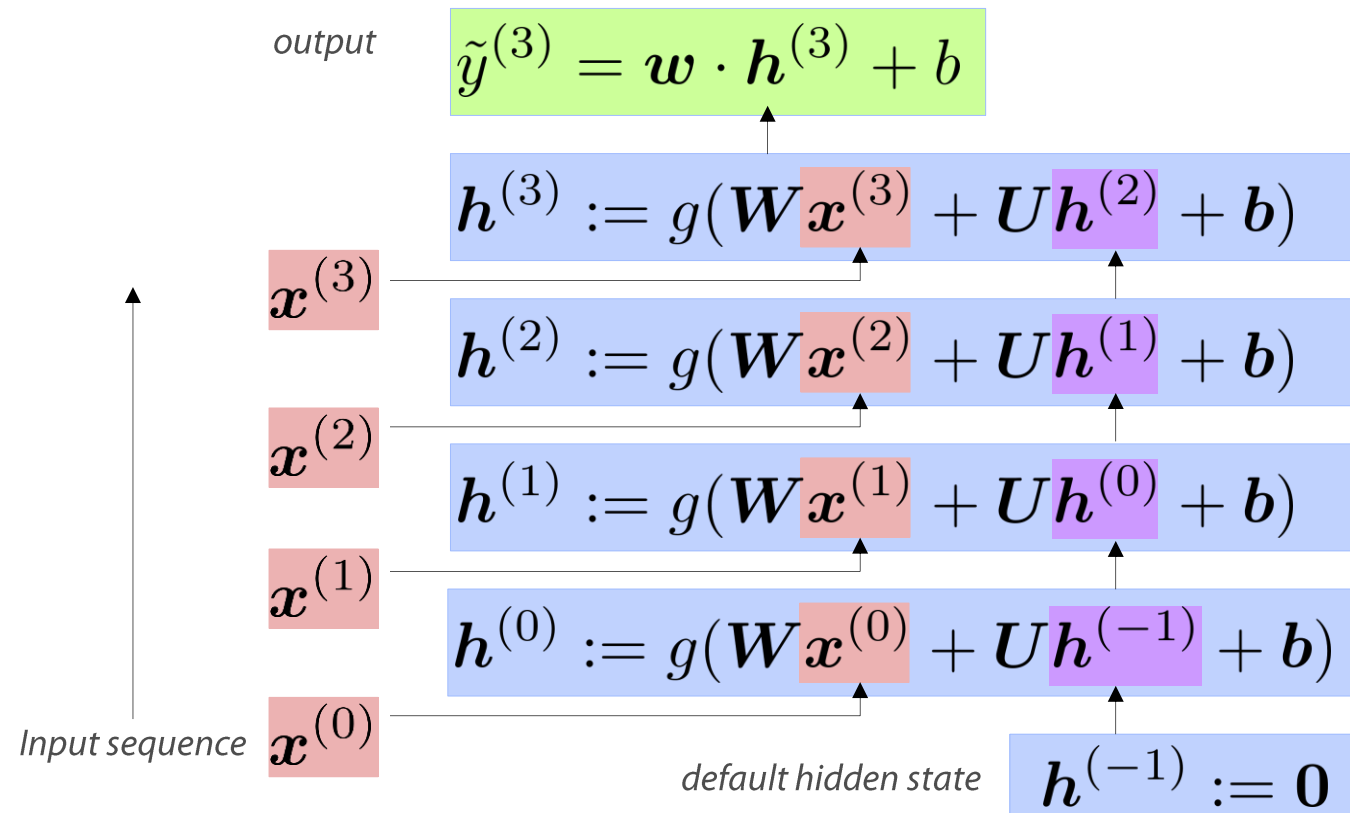
$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

where:

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

hidden at t      input at t      hidden at t-1

## Temporal Unfolding



# Recurrent Neural Networks

## Recurrent Neural Network

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

where:

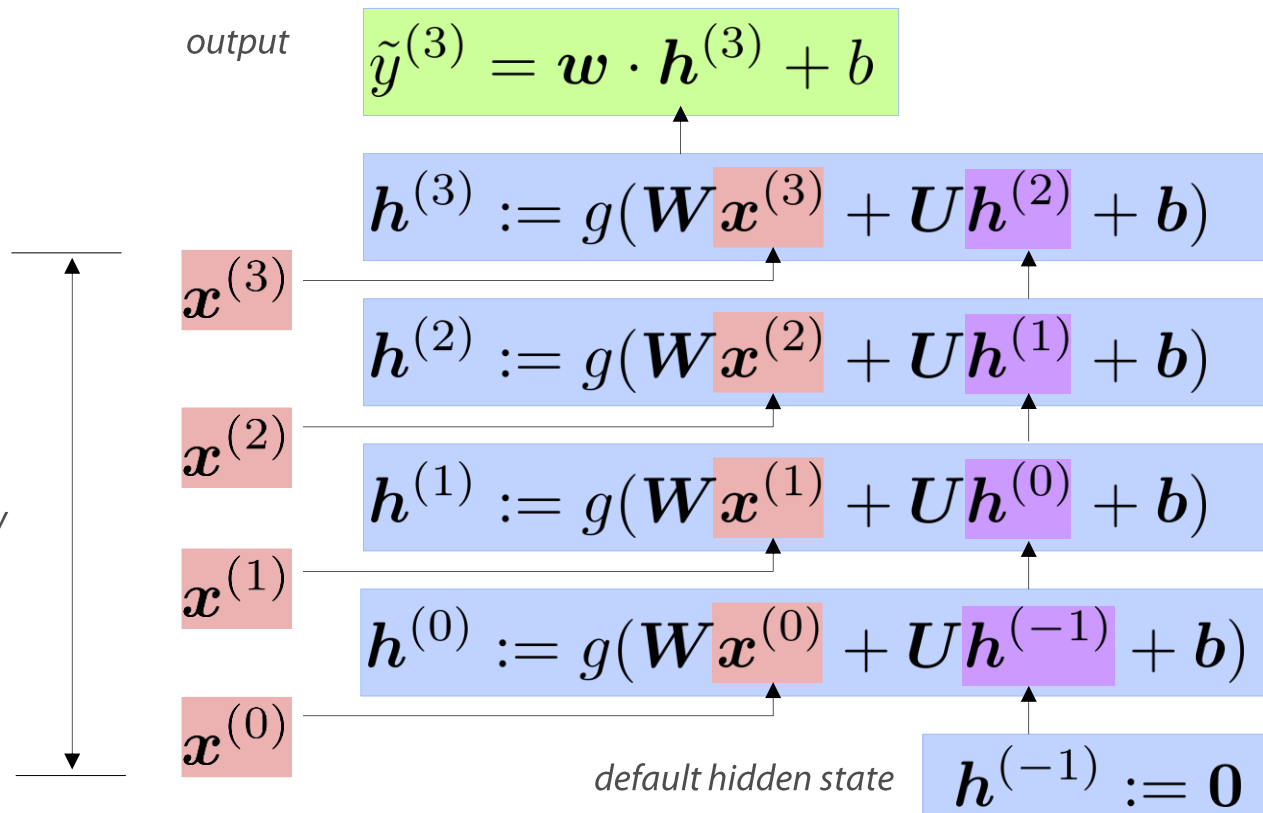
$$\mathbf{h}^{(t)} := g(\mathbf{W}\mathbf{x}^{(t)} + \mathbf{U}\mathbf{h}^{(t-1)} + \mathbf{b})$$

hidden at t      input at t      hidden at t-1

## Temporal Unfolding

This looks very similar to a deep feed-forward neural network ...

Context window



# Recurrent Neural Networks

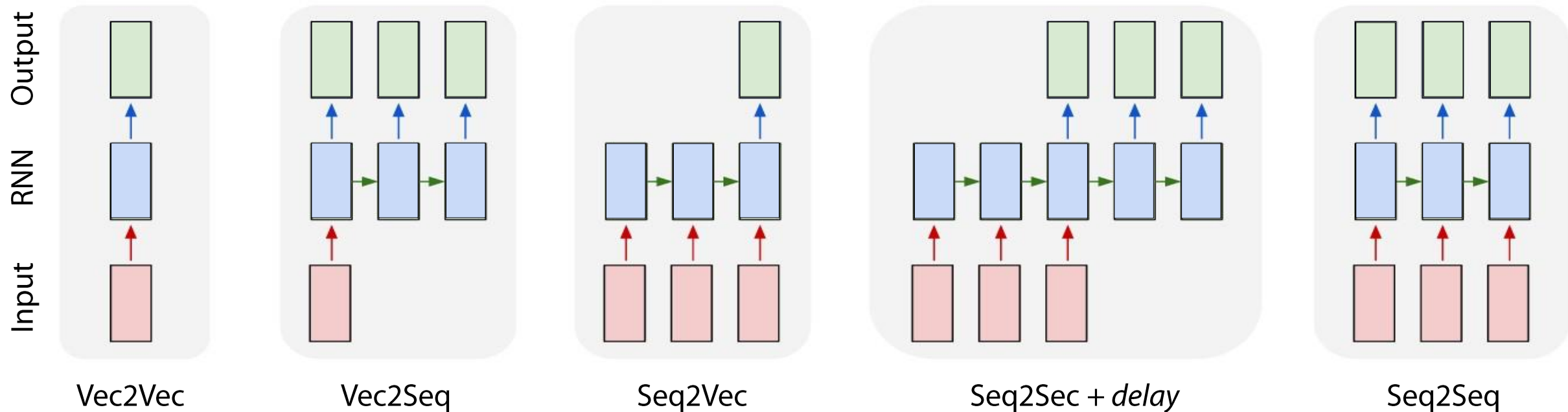
## Recurrent Neural Network

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

where:

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

## Input-Output Modes



# Recurrent Neural Networks

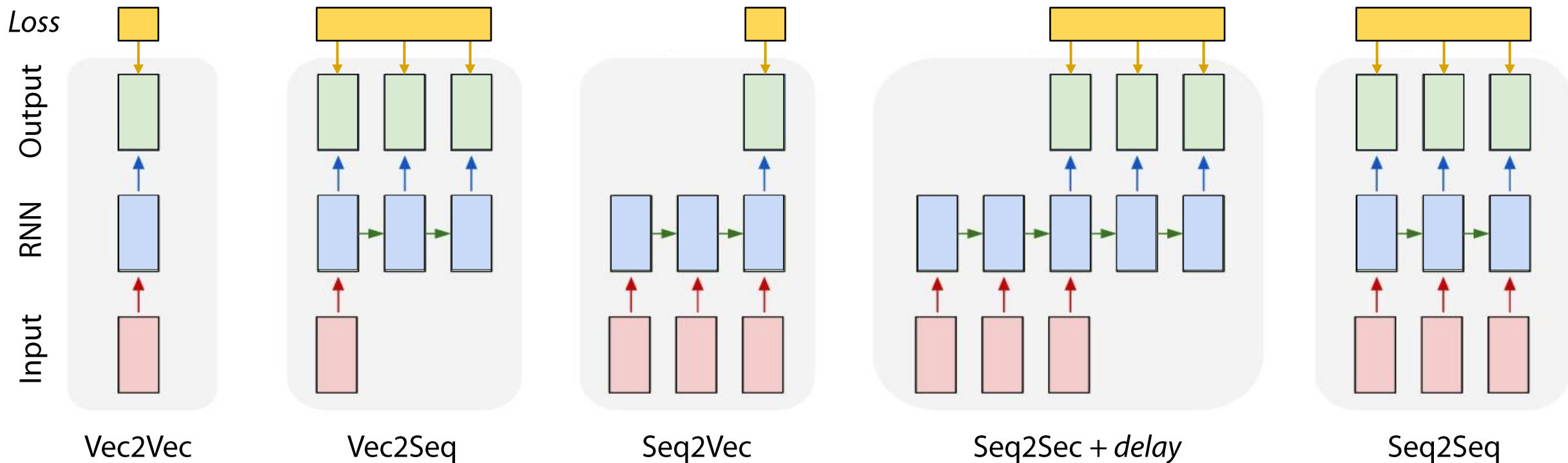
## Recurrent Neural Network

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

where:

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

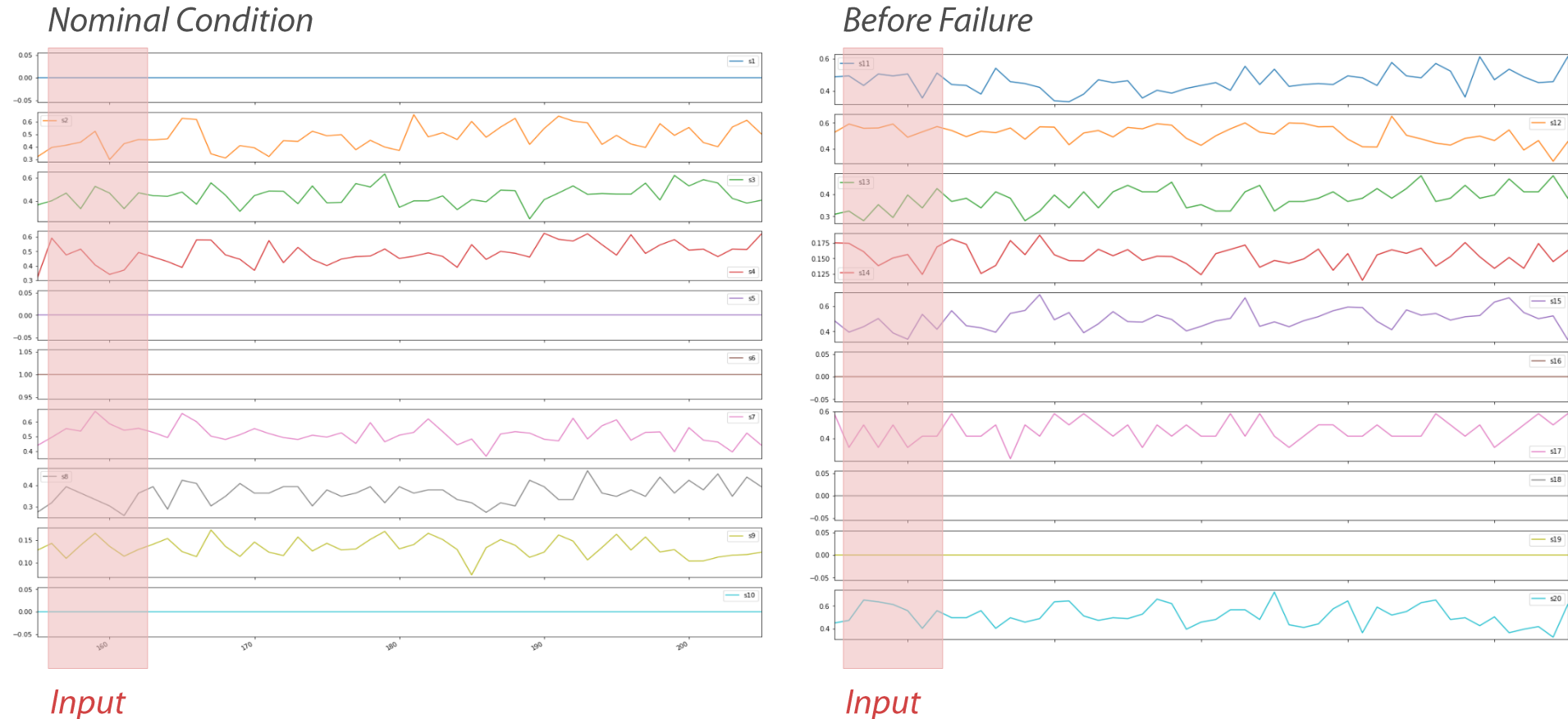
## Input-Output Modes



# *RNN Applications*

# Predictive Maintenance

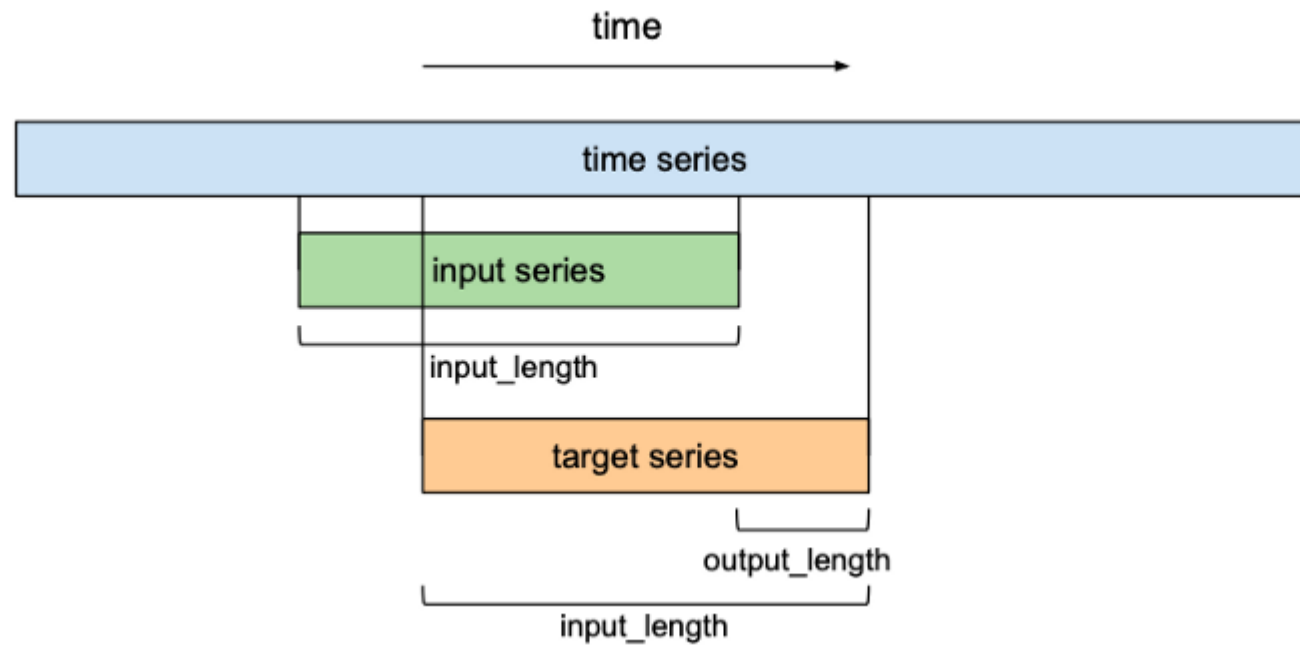
## ■ Detecting failure conditions from sensor readings



Training and Prediction occurs by using a sliding window of sensor readings as input

# Time Series Forecasting

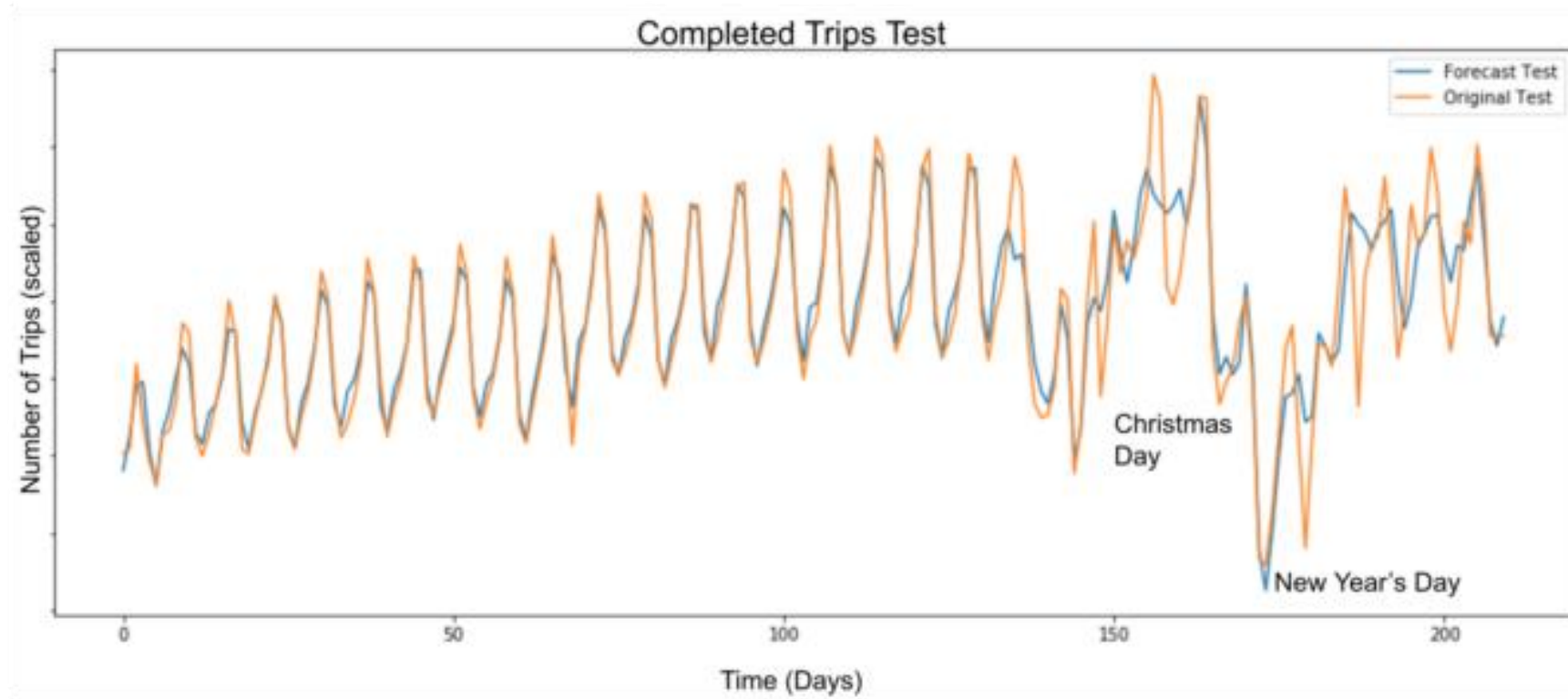
*Forecasting time series ahead of time*



[image from <https://medium.com/unit8-machine-learning-publication/temporal-convolutional-networks-and-forecasting-5ce1b6e97ce4>]

# Time Series Forecasting

*Detecting anomalies as differences from forecasted and actual*



[image from <https://eng.uber.com/neural-networks/>]

# *Long-Short Term Memory (LSTM)*

# LSTM Cell

- **Long-Short Term Memory** (Hochreiter & Schmidhuber, 1995)

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

$$\mathbf{h}^{(t)} := \mathbf{o}^{(t)} \odot \text{tanh}(\mathbf{c}^{(t)})$$

*elementwise product*

$$\mathbf{c}^{(t)} := \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \mathbf{c}_{in}^{(t)}$$

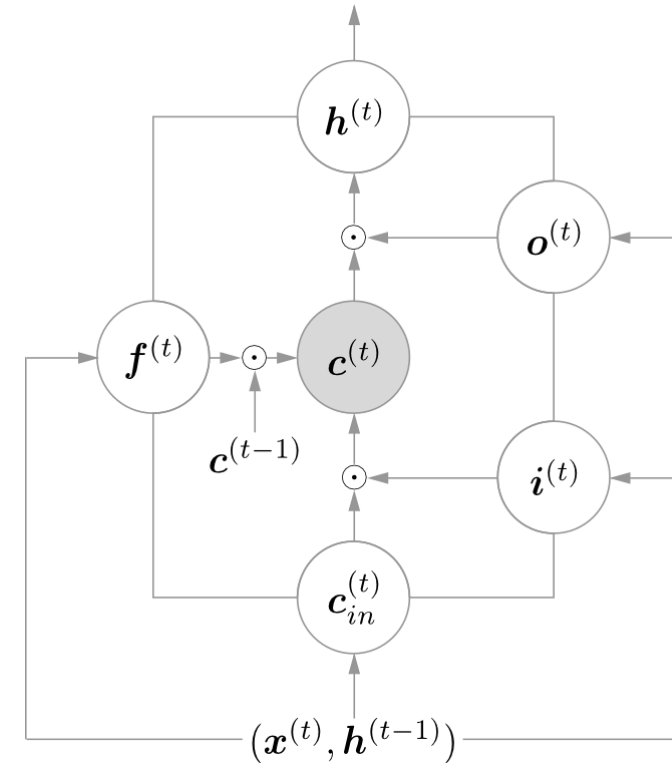
$$\mathbf{o}^{(t)} := \sigma(\mathbf{W}_o \mathbf{x}^{(t)} + \mathbf{U}_o \mathbf{h}^{(t-1)} + \mathbf{b}_o)$$

*sigmoid function*

$$\mathbf{f}^{(t)} := \sigma(\mathbf{W}_f \mathbf{x}^{(t)} + \mathbf{U}_f \mathbf{h}^{(t-1)} + \mathbf{b}_f)$$

$$\mathbf{i}^{(t)} := \sigma(\mathbf{W}_i \mathbf{x}^{(t)} + \mathbf{U}_i \mathbf{h}^{(t-1)} + \mathbf{b}_i)$$

$$\mathbf{c}_{in}^{(t)} := \tanh(\mathbf{W}_c \mathbf{x}^{(t)} + \mathbf{U}_c \mathbf{h}^{(t-1)} + \mathbf{b}_c)$$



# LSTM Cell

- **Long-Short Term Memory** (Hochreiter & Schmidhuber, 1995)

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

$$\mathbf{h}^{(t)} := \mathbf{o}^{(t)} \odot \tanh(\mathbf{c}^{(t)})$$

$$\mathbf{c}^{(t)} := \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \mathbf{c}_{in}^{(t)}$$

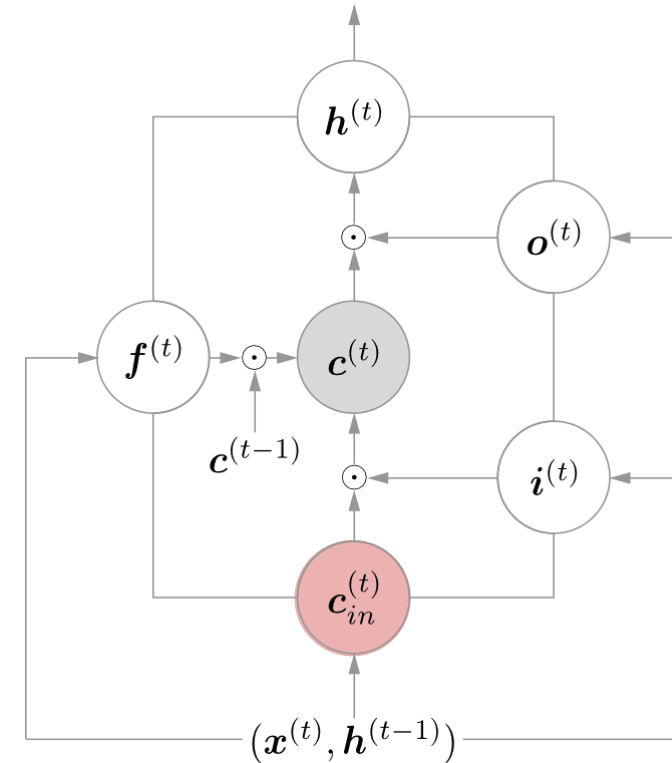
$$\mathbf{o}^{(t)} := \sigma(\mathbf{W}_o \mathbf{x}^{(t)} + \mathbf{U}_o \mathbf{h}^{(t-1)} + \mathbf{b}_o)$$

$$\mathbf{f}^{(t)} := \sigma(\mathbf{W}_f \mathbf{x}^{(t)} + \mathbf{U}_f \mathbf{h}^{(t-1)} + \mathbf{b}_f)$$

$$\mathbf{i}^{(t)} := \sigma(\mathbf{W}_i \mathbf{x}^{(t)} + \mathbf{U}_i \mathbf{h}^{(t-1)} + \mathbf{b}_i)$$

Combined input

$$\mathbf{c}_{in}^{(t)} := \tanh(\mathbf{W}_c \mathbf{x}^{(t)} + \mathbf{U}_c \mathbf{h}^{(t-1)} + \mathbf{b}_c)$$



# LSTM Cell

- **Long-Short Term Memory** (Hochreiter & Schmidhuber, 1995)

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

$$\mathbf{h}^{(t)} := \mathbf{o}^{(t)} \odot \tanh(\mathbf{c}^{(t)})$$

$$\mathbf{c}^{(t)} := \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \mathbf{c}_{in}^{(t)}$$

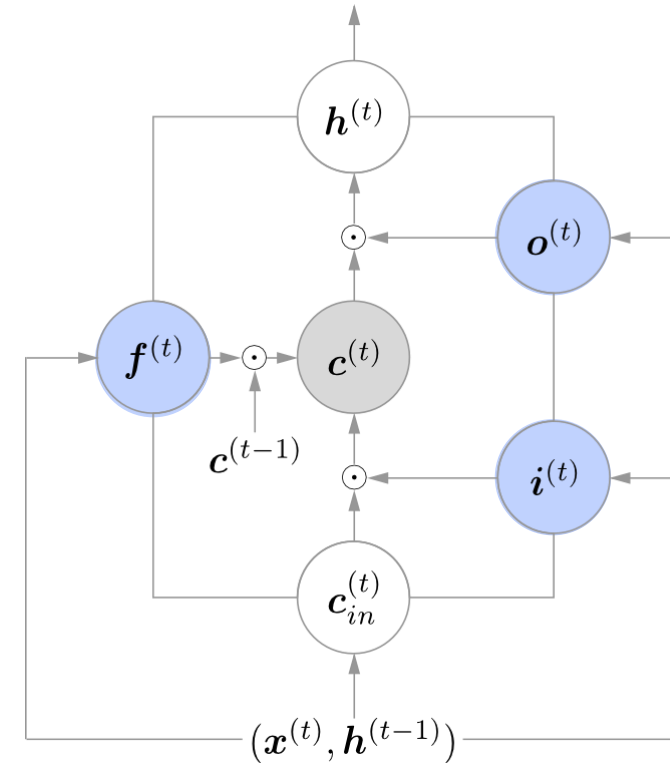
Gating values

$$\mathbf{o}^{(t)} := \sigma(\mathbf{W}_o \mathbf{x}^{(t)} + \mathbf{U}_o \mathbf{h}^{(t-1)} + \mathbf{b}_o) \quad \text{output}$$

$$\mathbf{f}^{(t)} := \sigma(\mathbf{W}_f \mathbf{x}^{(t)} + \mathbf{U}_f \mathbf{h}^{(t-1)} + \mathbf{b}_f) \quad \text{forget}$$

$$\mathbf{i}^{(t)} := \sigma(\mathbf{W}_i \mathbf{x}^{(t)} + \mathbf{U}_i \mathbf{h}^{(t-1)} + \mathbf{b}_i) \quad \text{input}$$

$$\mathbf{c}_{in}^{(t)} := \tanh(\mathbf{W}_c \mathbf{x}^{(t)} + \mathbf{U}_c \mathbf{h}^{(t-1)} + \mathbf{b}_c)$$



# LSTM Cell

- **Long-Short Term Memory** (Hochreiter & Schmidhuber, 1995)

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

*Applying gates*

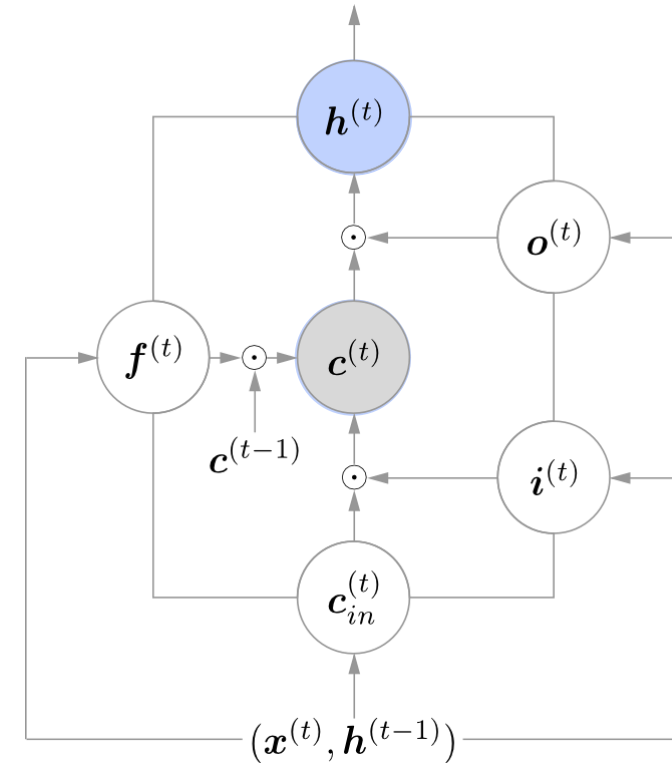
$$\mathbf{h}^{(t)} := \mathbf{o}^{(t)} \odot \tanh(\mathbf{c}^{(t)}) \quad \text{hidden}$$
$$\mathbf{c}^{(t)} := \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \mathbf{c}_{in}^{(t)} \quad \text{memory}$$

$$\mathbf{o}^{(t)} := \sigma(\mathbf{W}_o \mathbf{x}^{(t)} + \mathbf{U}_o \mathbf{h}^{(t-1)} + \mathbf{b}_o)$$

$$\mathbf{f}^{(t)} := \sigma(\mathbf{W}_f \mathbf{x}^{(t)} + \mathbf{U}_f \mathbf{h}^{(t-1)} + \mathbf{b}_f)$$

$$\mathbf{i}^{(t)} := \sigma(\mathbf{W}_i \mathbf{x}^{(t)} + \mathbf{U}_i \mathbf{h}^{(t-1)} + \mathbf{b}_i)$$

$$\mathbf{c}_{in}^{(t)} := \tanh(\mathbf{W}_c \mathbf{x}^{(t)} + \mathbf{U}_c \mathbf{h}^{(t-1)} + \mathbf{b}_c)$$



# LSTM Cell

- **Long-Short Term Memory** (Hochreiter & Schmidhuber, 1995)

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b \quad \text{Cell output}$$

$$\mathbf{h}^{(t)} := \mathbf{o}^{(t)} \odot \tanh(\mathbf{c}^{(t)})$$

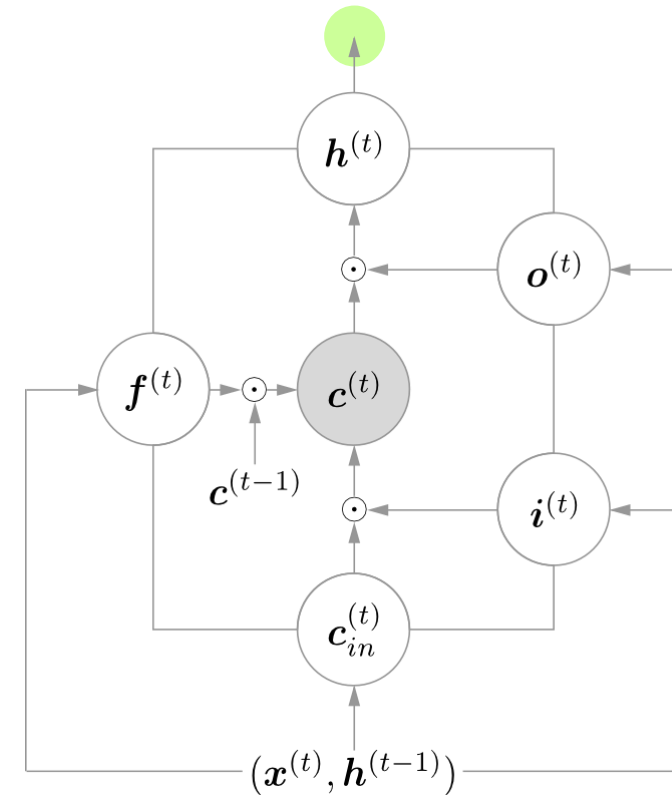
$$\mathbf{c}^{(t)} := \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \mathbf{c}_{in}^{(t)}$$

$$\mathbf{o}^{(t)} := \sigma(\mathbf{W}_o \mathbf{x}^{(t)} + \mathbf{U}_o \mathbf{h}^{(t-1)} + \mathbf{b}_o)$$

$$\mathbf{f}^{(t)} := \sigma(\mathbf{W}_f \mathbf{x}^{(t)} + \mathbf{U}_f \mathbf{h}^{(t-1)} + \mathbf{b}_f)$$

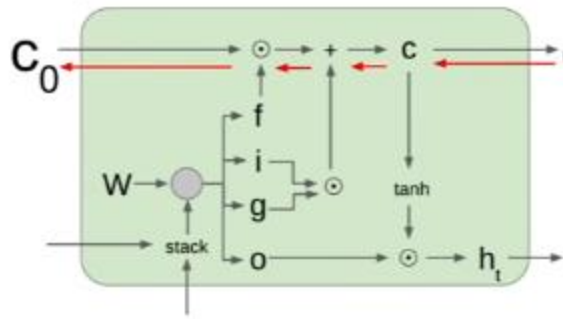
$$\mathbf{i}^{(t)} := \sigma(\mathbf{W}_i \mathbf{x}^{(t)} + \mathbf{U}_i \mathbf{h}^{(t-1)} + \mathbf{b}_i)$$

$$\mathbf{c}_{in}^{(t)} := \tanh(\mathbf{W}_c \mathbf{x}^{(t)} + \mathbf{U}_c \mathbf{h}^{(t-1)} + \mathbf{b}_c)$$

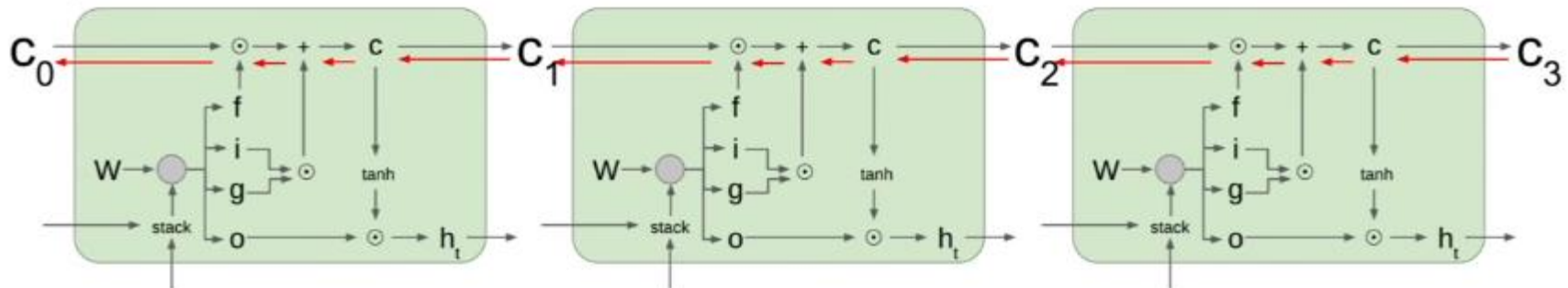


# LSTM Training

## Temporal Unfolding

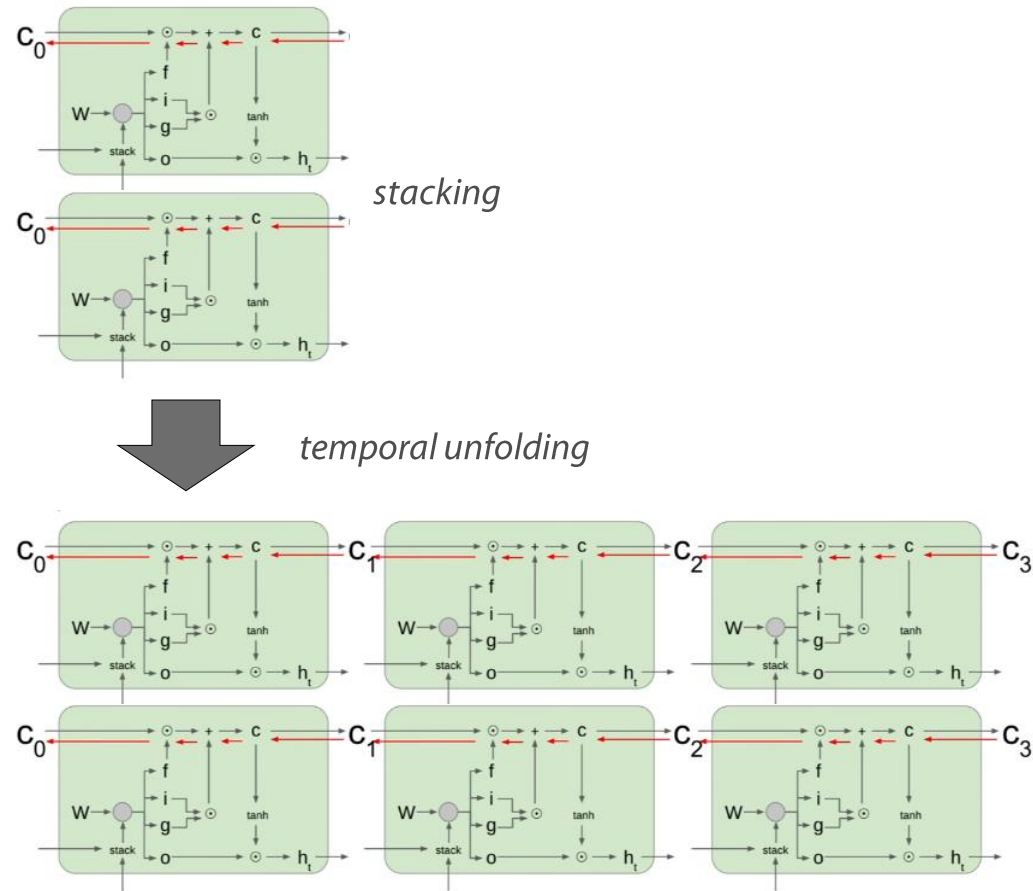


temporal unfolding



# LSTM Training

## Stacking and Temporal Unfolding



# GRU

- **Gated Recurrent Unit** (Kyunghyun Cho et al., 1995)

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

$$\mathbf{h}^{(t)} := \underbrace{(1 - \mathbf{z}^{(t)}) \odot \mathbf{h}^{(t-1)}}_{\text{(exponential moving average)}} + \mathbf{z}^{(t)} \odot \hat{\mathbf{h}}^{(t)}$$

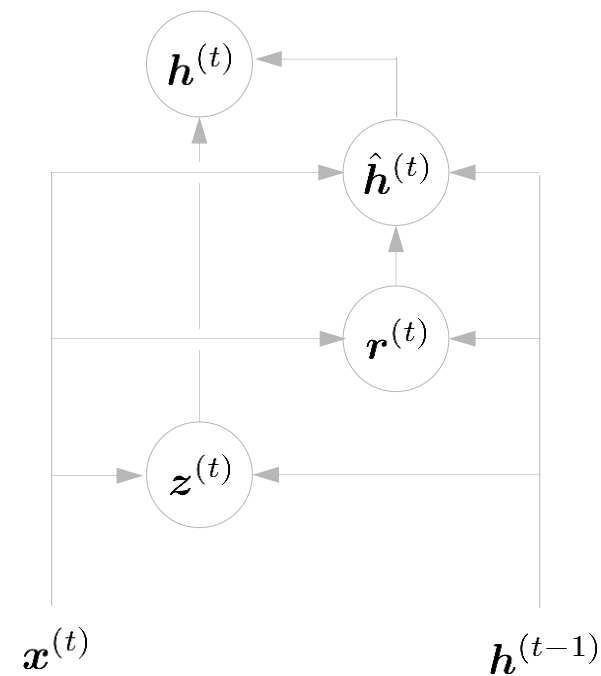
$$\hat{\mathbf{h}}^{(t)} := \tanh(\mathbf{W}_h \mathbf{x}^{(t)} + \mathbf{U}_h (\mathbf{r}^{(t)} \odot \mathbf{h}^{(t-1)})) + \mathbf{b}_h$$

$$\mathbf{r}^{(t)} := \sigma \left( \mathbf{W}_r \mathbf{x}^{(t)} + \mathbf{U}_r \mathbf{h}^{(t-1)} + \mathbf{b}_r \right)$$

*sigmoid function*

$$\mathbf{z}^{(t)} := \sigma(\mathbf{W}_z \mathbf{x}^{(t)} + \mathbf{U}_z \mathbf{h}^{(t-1)} + \mathbf{b}_z)$$

Simpler structure, no internal memory



# GRU

- **Gated Recurrent Unit** (Kyunghyun Cho et al., 1995)

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

$$\mathbf{h}^{(t)} := \underbrace{(1 - \mathbf{z}^{(t)}) \odot \mathbf{h}^{(t-1)}}_{\text{(exponential moving average)}} + \mathbf{z}^{(t)} \odot \hat{\mathbf{h}}^{(t)}$$

$$\hat{\mathbf{h}}^{(t)} := \tanh(\mathbf{W}_h \mathbf{x}^{(t)} + \mathbf{U}_h (\mathbf{r}^{(t)} \odot \mathbf{h}^{(t-1)})) + \mathbf{b}_h$$

Gating values

$$\mathbf{r}^{(t)} := \sigma(\mathbf{W}_r \mathbf{x}^{(t)} + \mathbf{U}_r \mathbf{h}^{(t-1)} + \mathbf{b}_r) \quad \text{reset}$$

$$\mathbf{z}^{(t)} := \sigma(\mathbf{W}_z \mathbf{x}^{(t)} + \mathbf{U}_z \mathbf{h}^{(t-1)} + \mathbf{b}_z) \quad \text{update}$$

Simpler structure, no internal memory

