

Deep Learning

A course about theory & practice



Deep Learning and Time Series

Marco Piastra

Recurrent Neural Networks

Recurrent Neural Networks

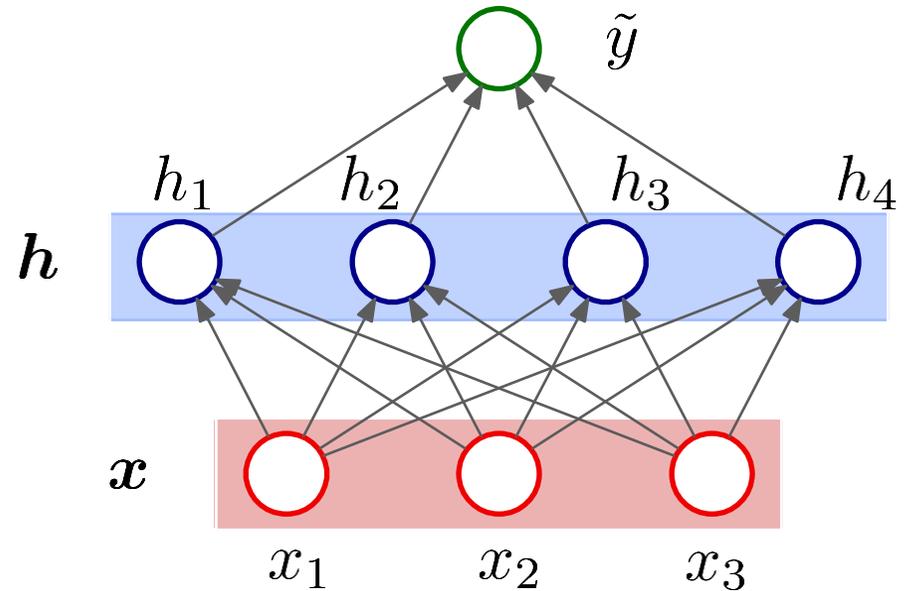
- **Feed-forward neural network**

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

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

Recurrent Neural Networks

Recurrent Neural Network

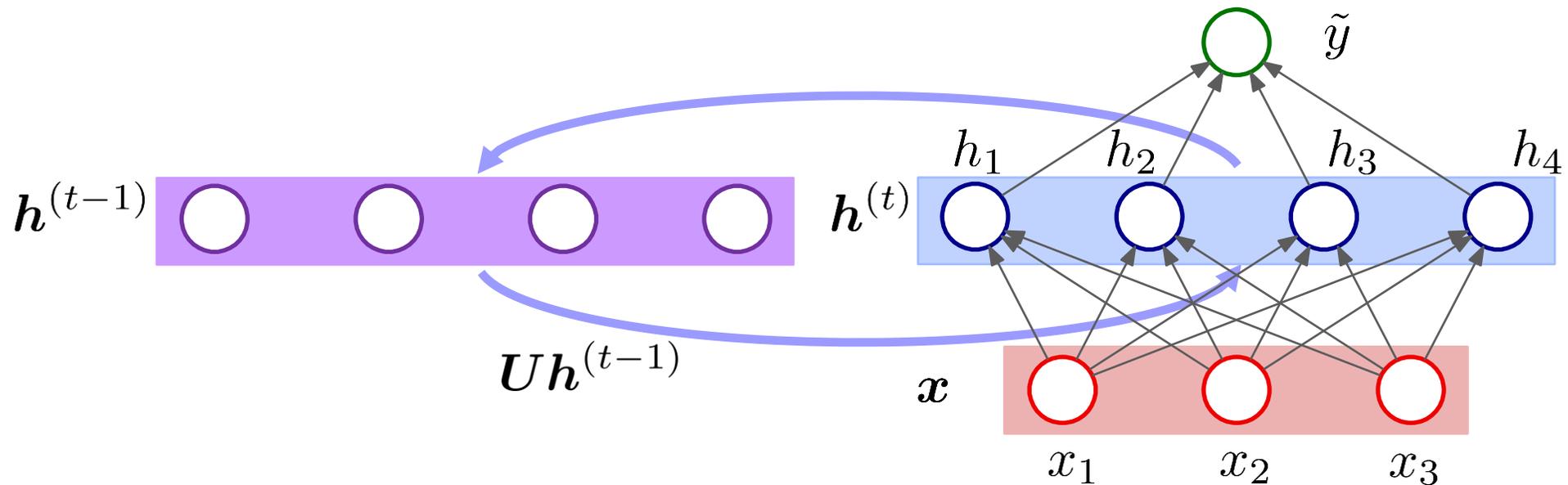
$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

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 basic 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

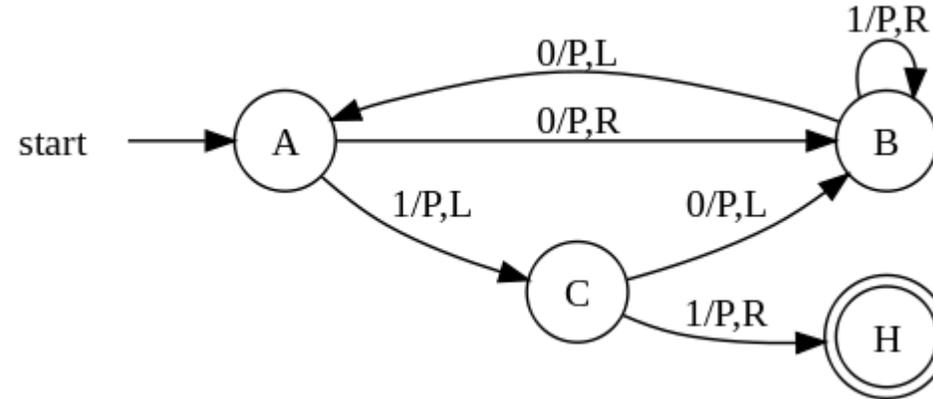
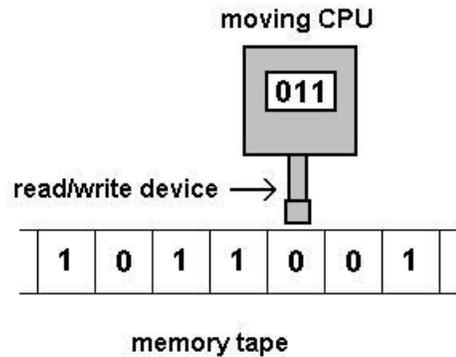
$$\mathbf{h}^{(t)} := g(\mathbf{W} \mathbf{x}^{(t)} + \mathbf{U} \mathbf{h}^{(t-1)} + \mathbf{b})$$

hidden at t input at t hidden at t-1

RNN are Turing Machines

- **Computational power of RNNs** (Siegelmann & Sontag, 1992)

“RNNs can simulate any Turing machine”



This means that they can compute anything a Turing Machine could

[image from https://en.wikipedia.org/wiki/Turing_machine]

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

■ General Properties

A recurrent neural network (RNN) is even more powerful than a FF neural network
It can approximate any Turing machine (i.e. a general theoretical model of computation)

It is much harder to train than a FF neural network

Although, with *temporal unfolding* (see above), *gradient descent* methods can be applied

Recurrent Neural Networks

Recurrent Neural Network

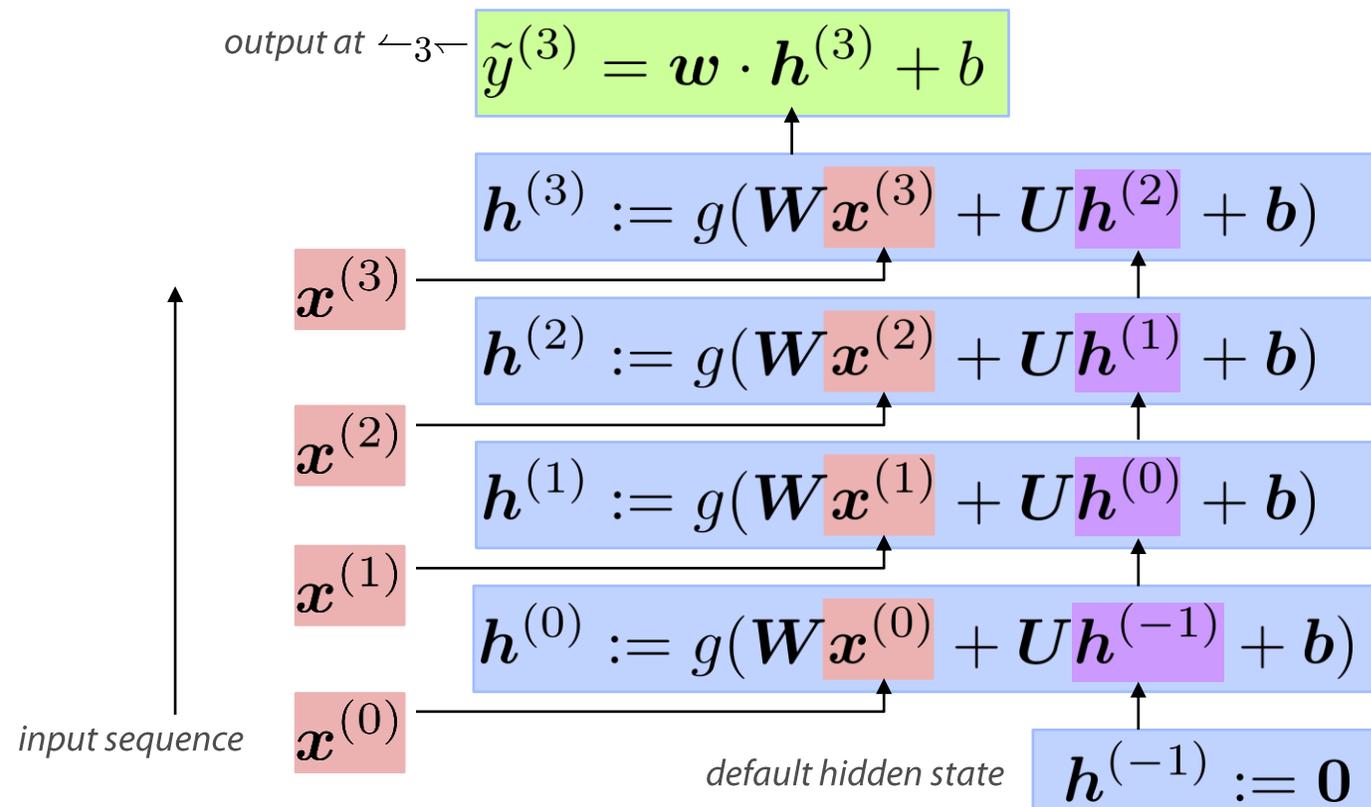
$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b$$

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

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

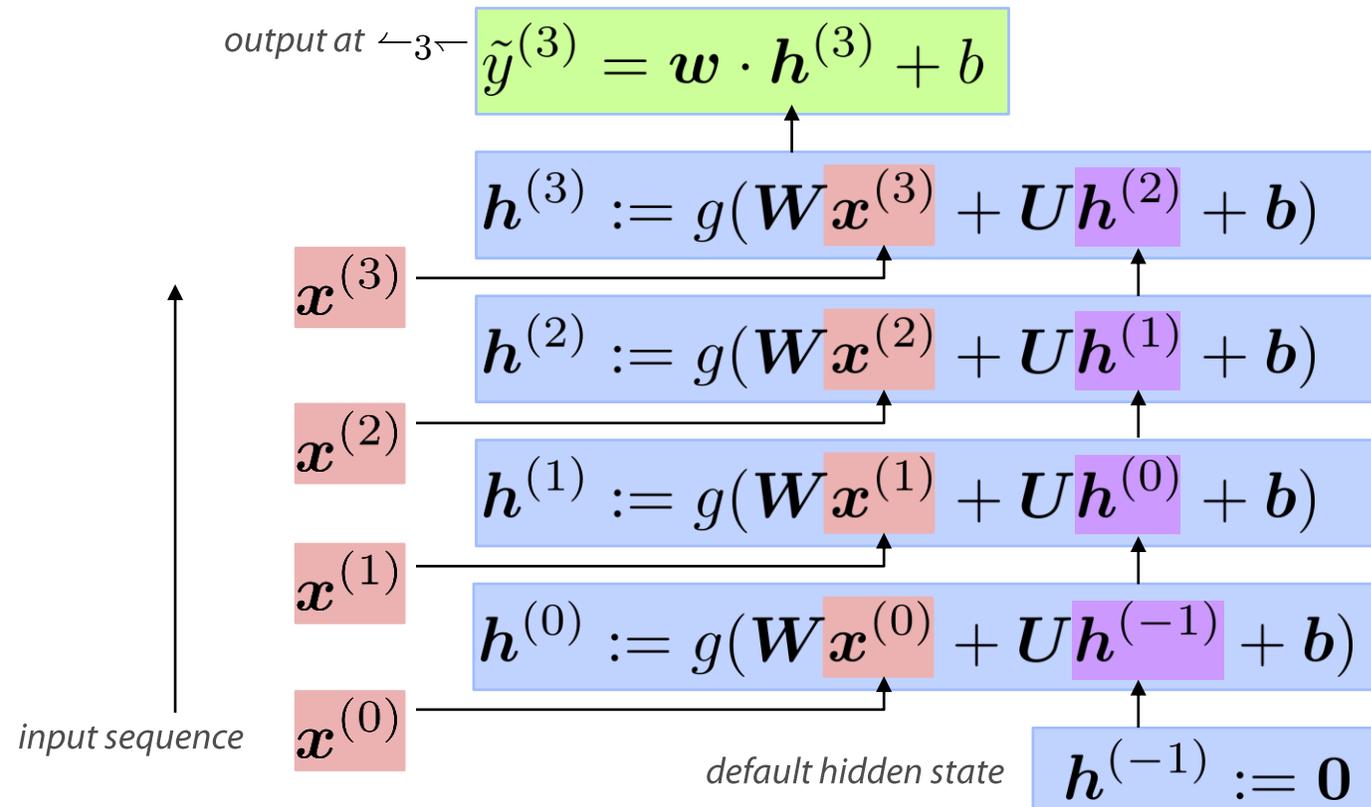
where

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

hidden at t input at t hidden at t-1

Temporal Unfolding

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

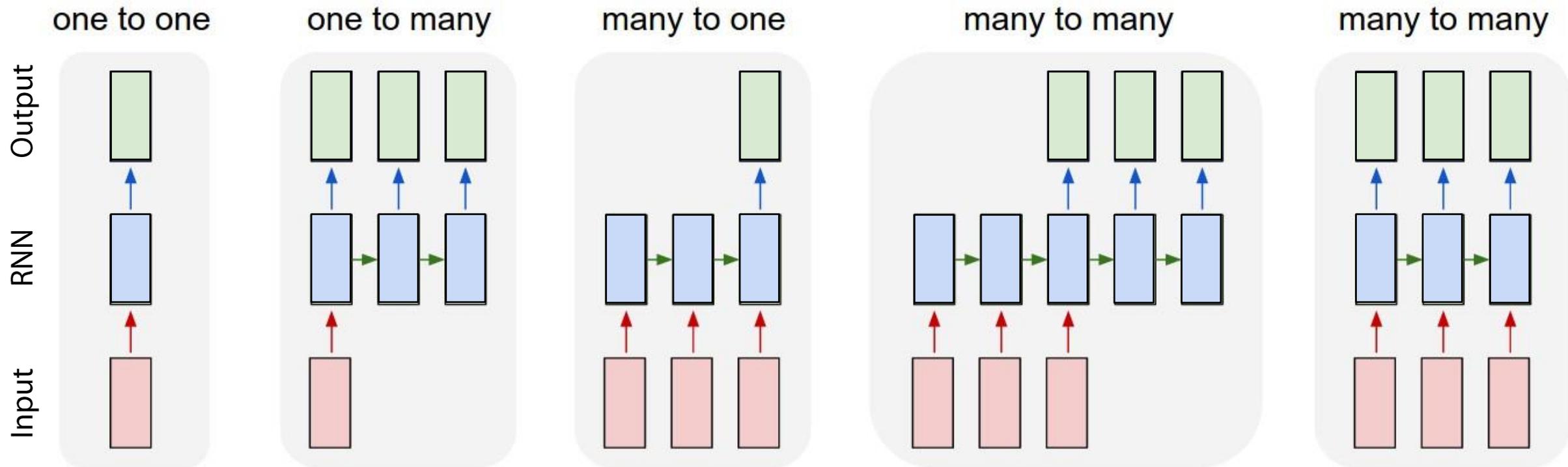


Recurrent Neural Networks

Recurrent Neural Network

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b \quad \text{where} \quad \mathbf{h}^{(t)} := g(\mathbf{W}\mathbf{x}^{(t)} + \mathbf{U}\mathbf{h}^{(t-1)} + \mathbf{b})$$

Input-Output Modes

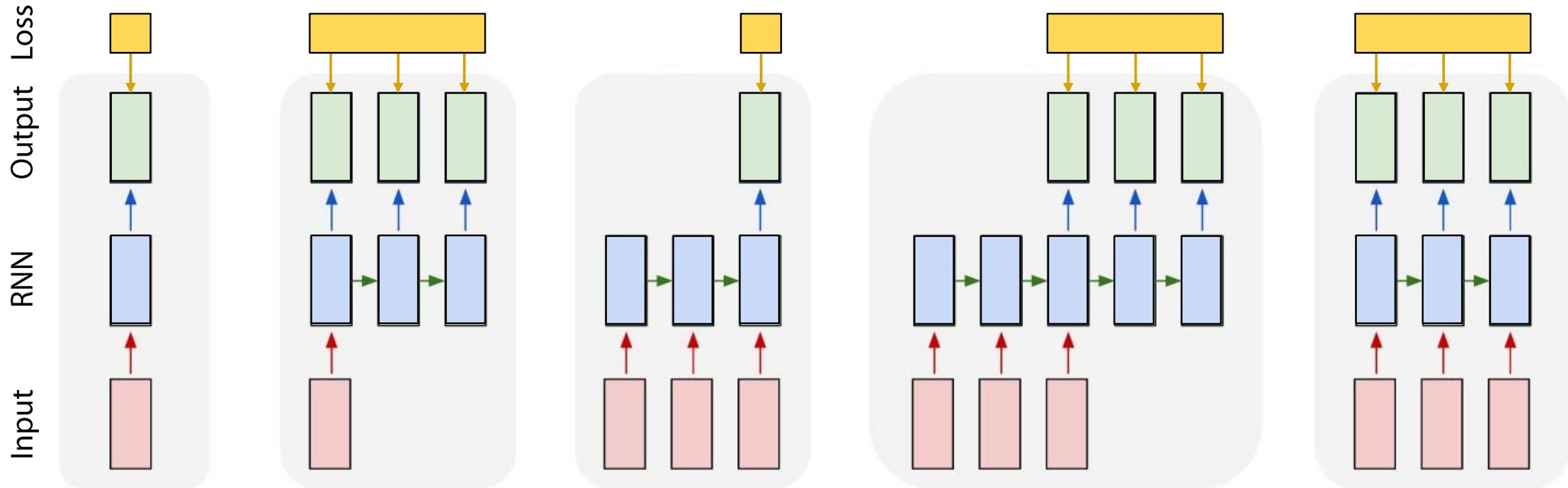


Recurrent Neural Networks

Recurrent Neural Network

$$\tilde{y}^{(t)} = \mathbf{w} \cdot \mathbf{h}^{(t)} + b \quad \text{where} \quad \mathbf{h}^{(t)} := g(\mathbf{W}\mathbf{x}^{(t)} + \mathbf{U}\mathbf{h}^{(t-1)} + \mathbf{b})$$

Input-Output Modes



Long-Short Term Memory (LSTM)

LSTM

- **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

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

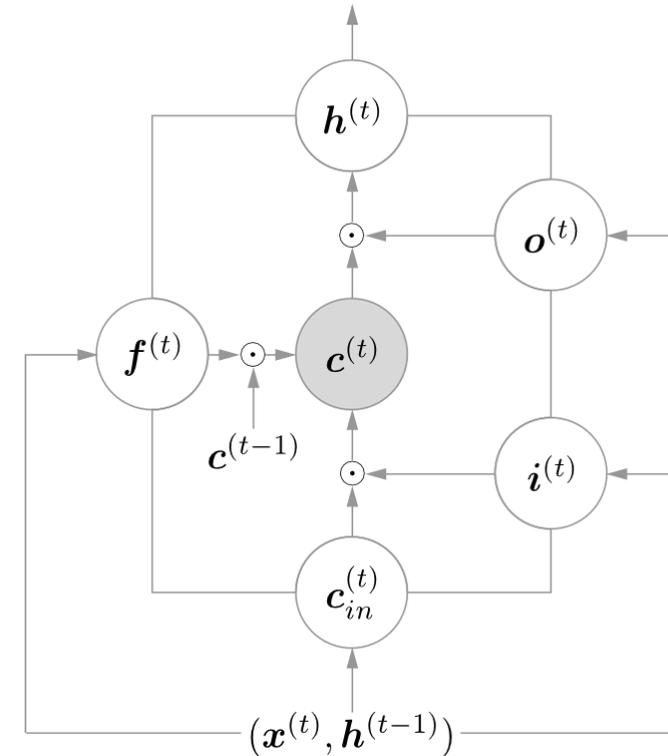
$$\mathbf{o}^{(t)} := \sigma(\mathbf{W}_o \mathbf{x}^{(t)} + \mathbf{U}_o \mathbf{h}^{(t-1)} + \mathbf{b}_o)$$

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

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

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

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

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

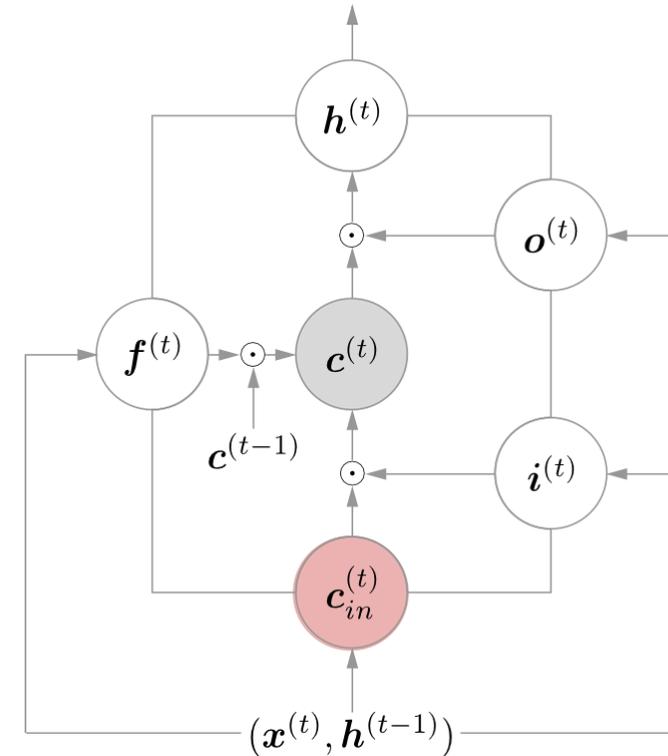
$$\mathbf{o}^{(t)} := \sigma(\mathbf{W}_o \mathbf{x}^{(t)} + \mathbf{U}_o \mathbf{h}^{(t-1)} + \mathbf{b}_o)$$

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

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

Combined input

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



LSTM

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

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

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

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

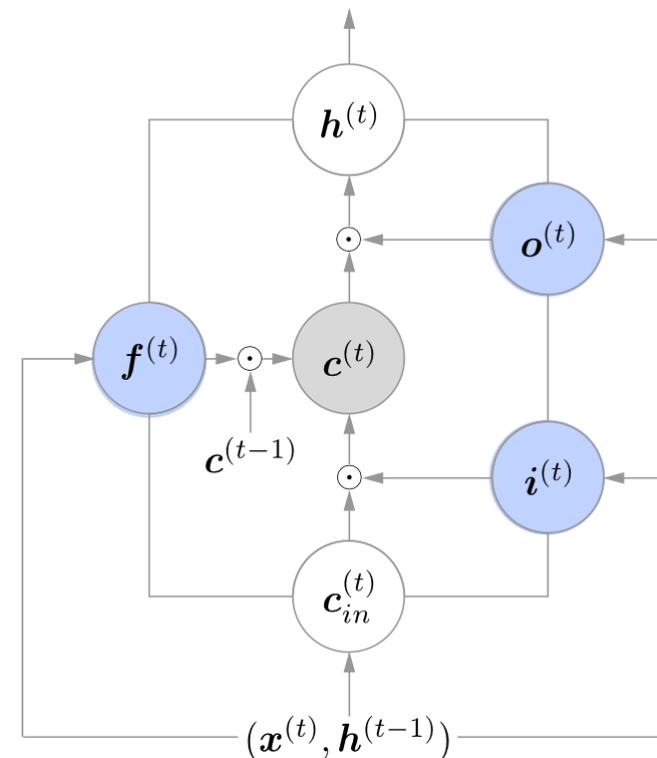
Gating values

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

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

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

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



LSTM

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

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

Applying gates

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

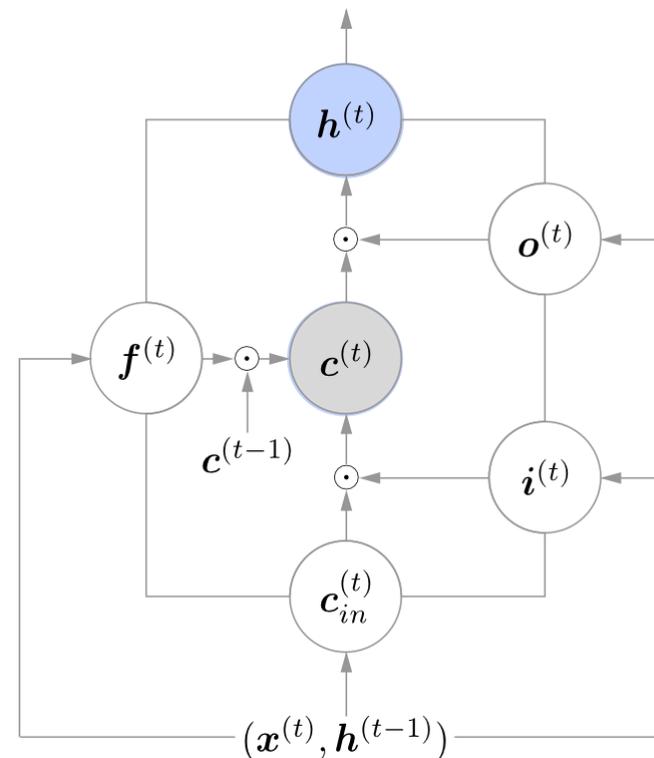
$$\mathbf{c}^{(t)} := \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \mathbf{c}_{in}^{(t)} \quad \text{memory}$$

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

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

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

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



LSTM

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

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

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

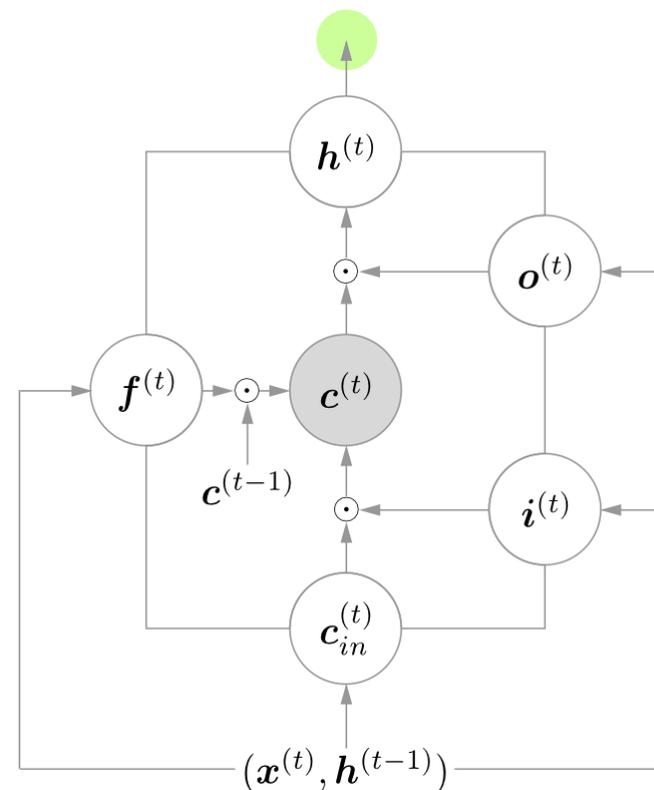
$$\mathbf{c}^{(t)} := \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \mathbf{c}_{in}^{(t)}$$

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

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

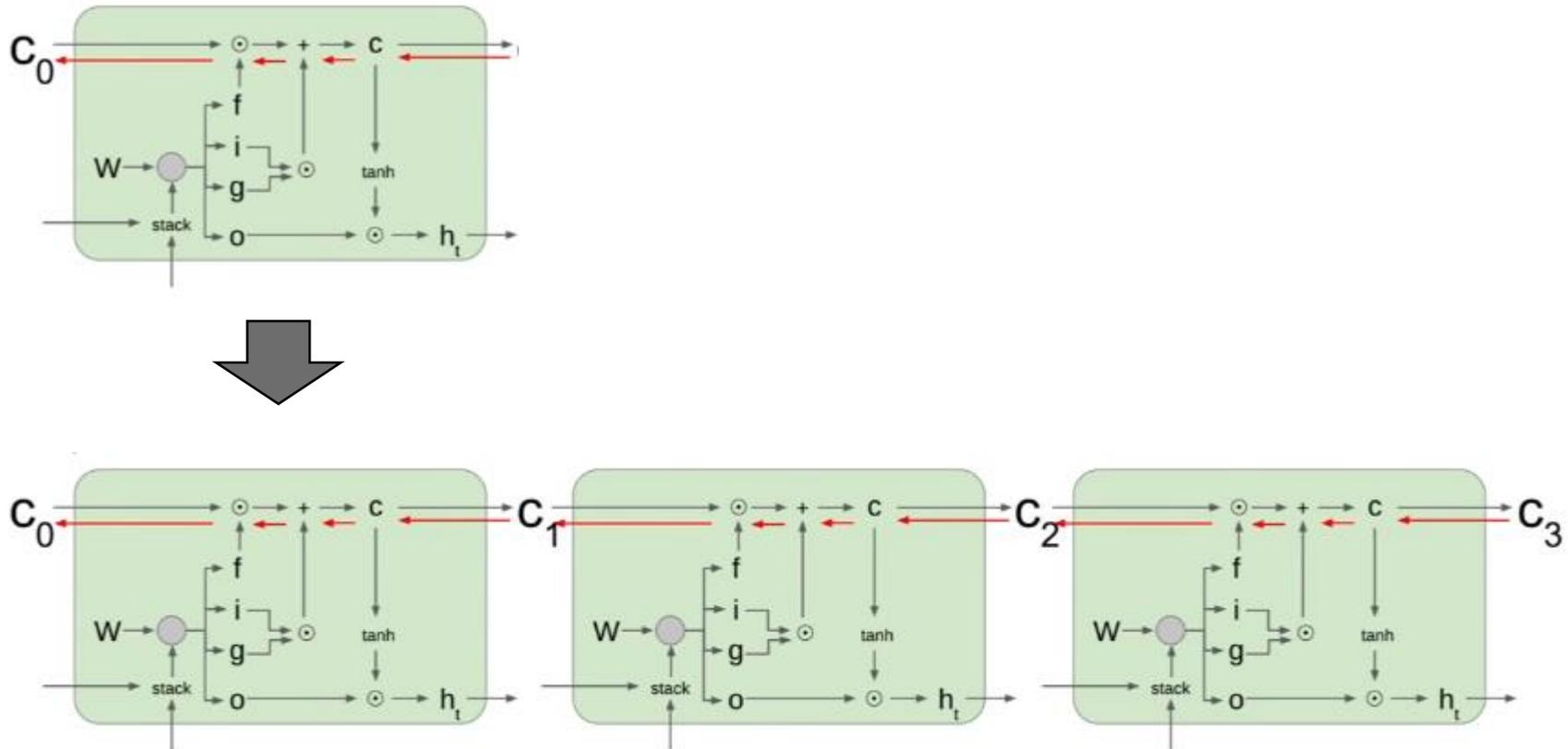
$$\mathbf{i}^{(t)} := \sigma(\mathbf{W}_i \mathbf{x}^{(t)} + \mathbf{U}_i \mathbf{h}^{(t-1)} + \mathbf{b}_i)$$

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



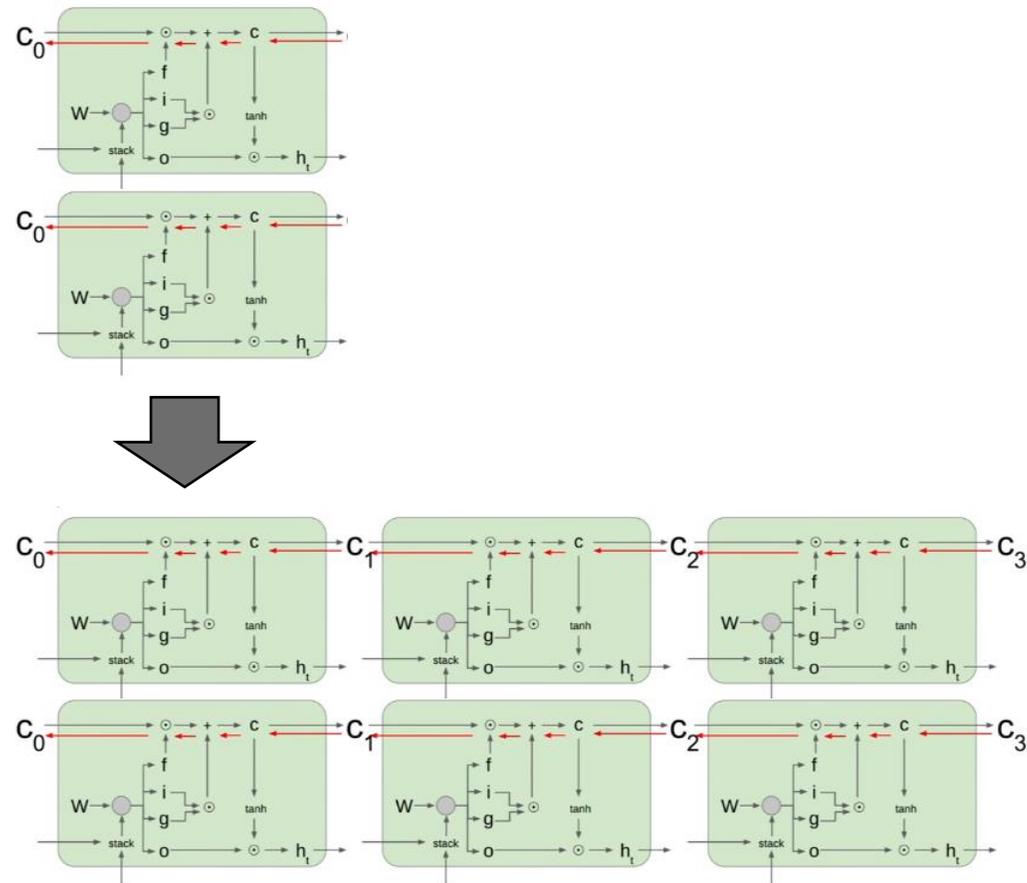
LSTM

Temporal Unfolding



LSTM

Stacking and Temporal Unfolding



GRU

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

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

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

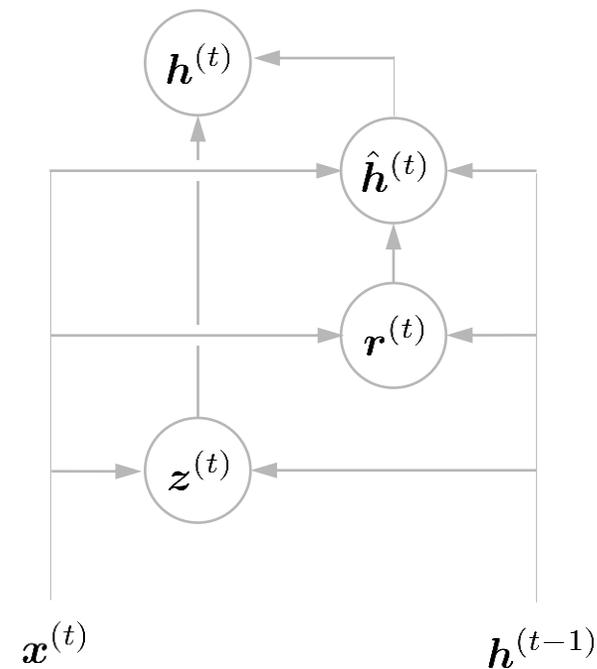
$$\hat{\mathbf{h}}^{(t)} := \tanh(\mathbf{W}_h \mathbf{x}^{(t)} + \mathbf{U}_h (\mathbf{r}^{(t)} \odot \mathbf{h}^{(t-1)})) + \mathbf{b}_h$$

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

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)

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

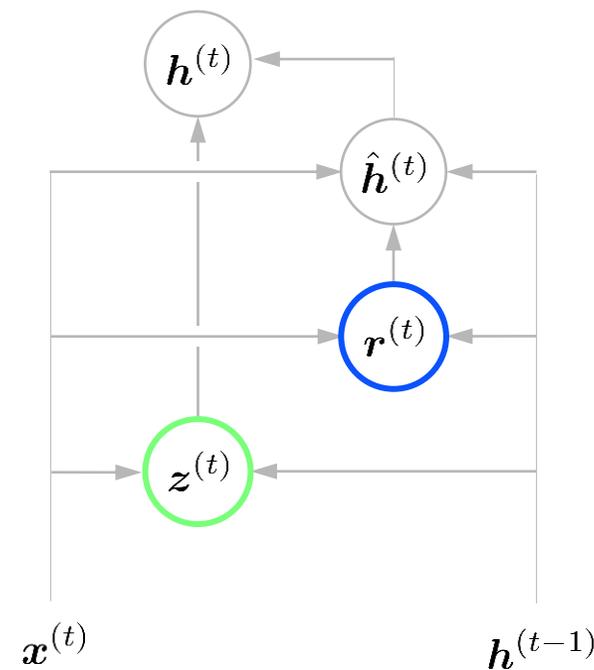
$$\mathbf{h}^{(t)} := \underbrace{(1 - \mathbf{z}^{(t)})}_{\text{(exponential moving average)}} \odot \mathbf{h}^{(t-1)} + \mathbf{z}^{(t)} \odot \hat{\mathbf{h}}^{(t)}$$

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

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

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

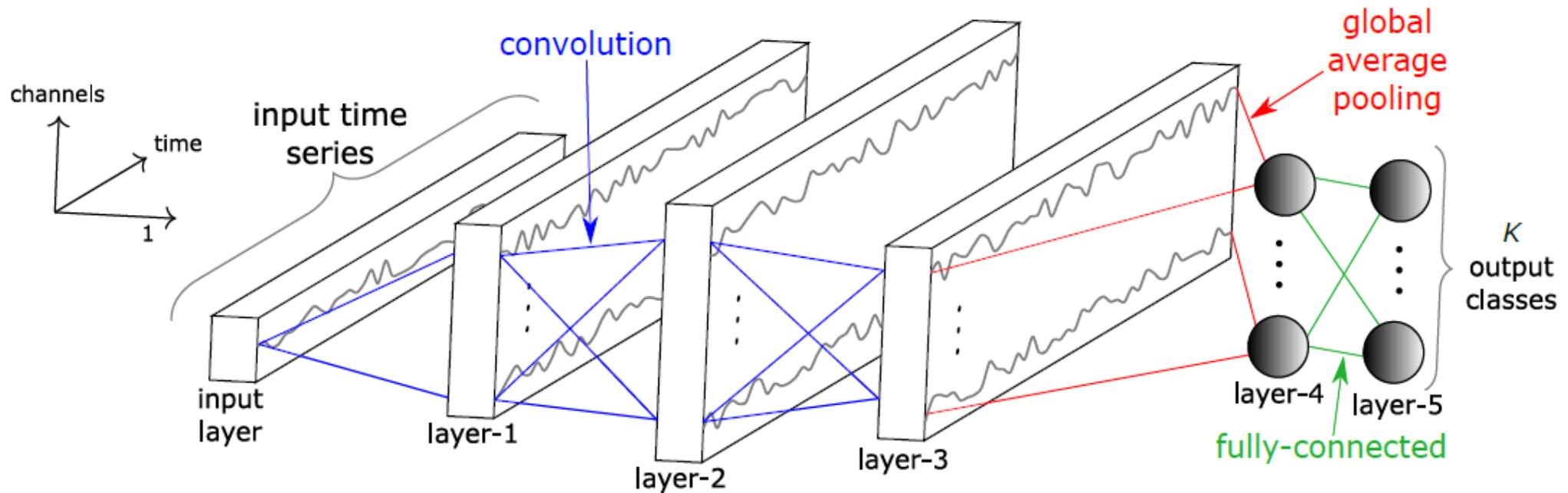
Simpler structure, no internal memory



Deep Convolutional Networks for Time Series Analysis

1D Convolution Over Time

Time windows are treated as 1D 'images'

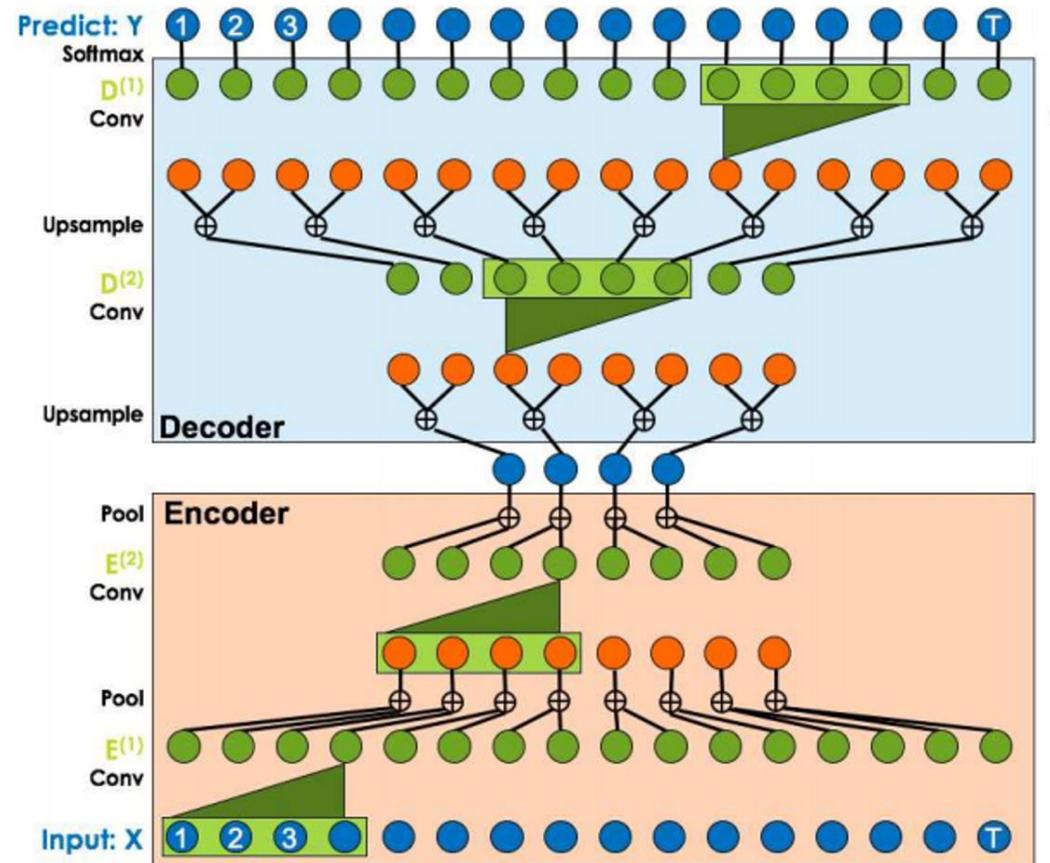


[image from <https://link.springer.com/article/10.1007%2Fs10618-019-00619-1>]

Temporal Convolution Networks

1D Convolution-Deconvolution in an autoencoder architecture (Lea et al., 2016)

Effective in segmenting actions and predicting time series

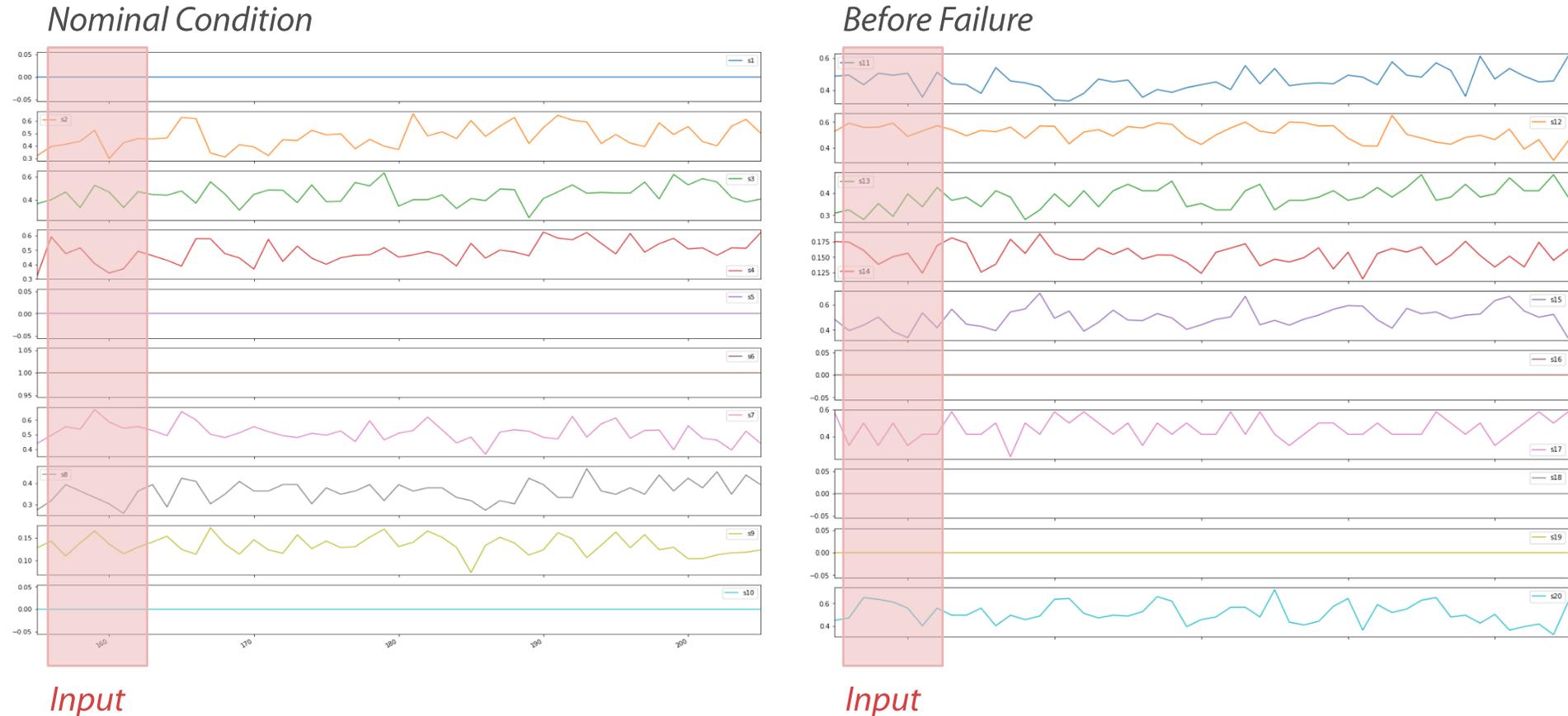


[image from <https://arxiv.org/abs/1611.05267>]

RNN applications

Predictive Maintenance

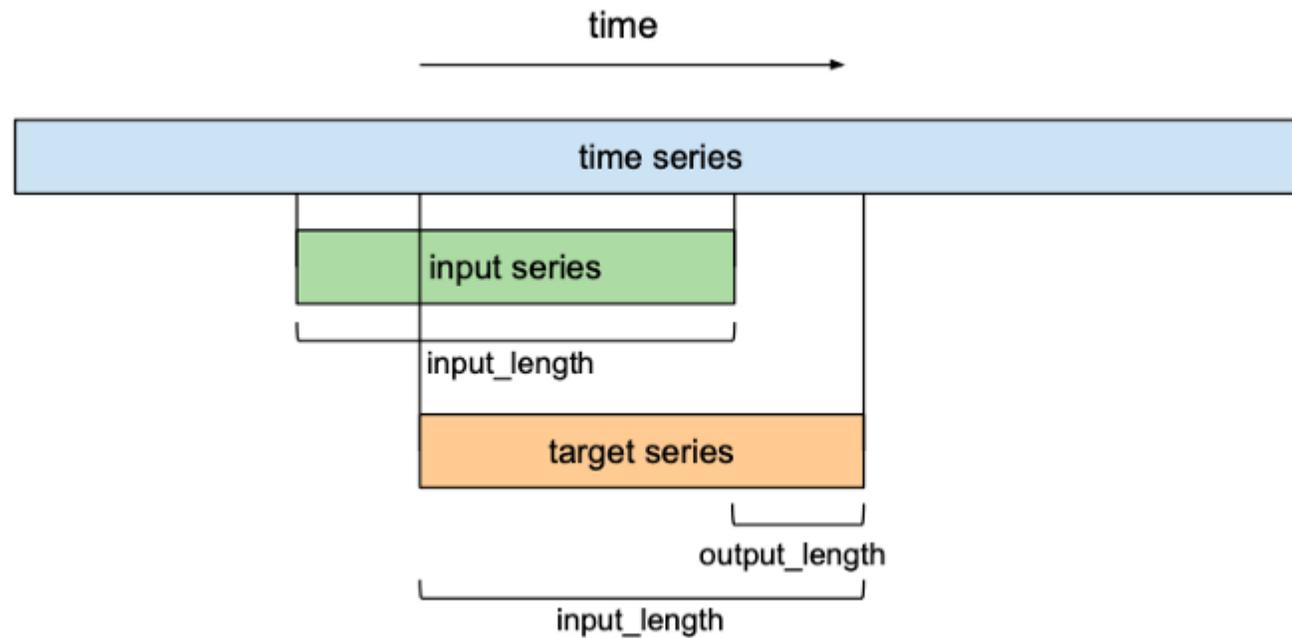
■ Detecting failure conditions from sensor readings



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

Time Series Forecasting

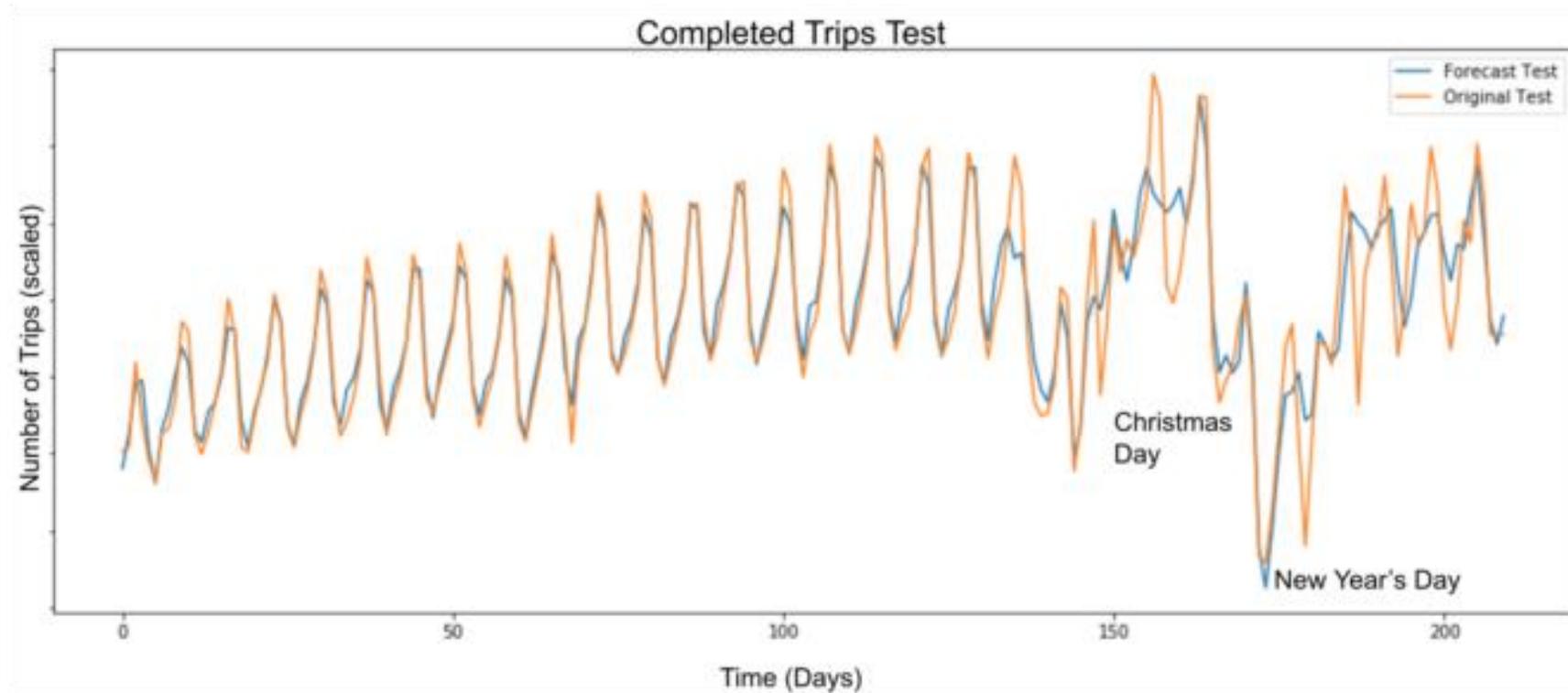
Forecasting time series ahead of time



[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/>]

An Aside: Embedded AI (at UniPV)

Credits are due to M. Musci and E. Torti

STM SensorTile for Wearable Devices

Microcontroller STM32 (ARM Cortex M4)

80MHz Clock Frequency

128KB SRAM

1MB Flash memory

Multiple Onboard Sensors

Microphone

3D Accelerometer + 3D Gyroscope

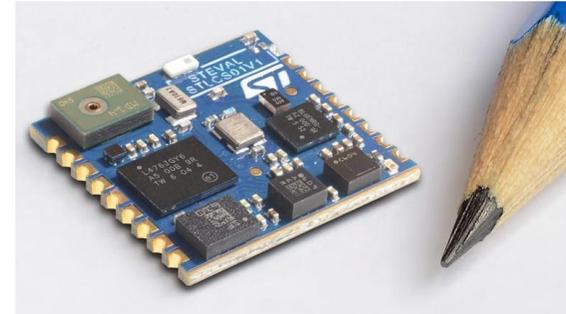
e-Compass, 3D Accelerometer, 3D Magnetic sensor

Barometer

Peripherals

100 mAh Li-Ion battery

Bluetooth Low Energy (BLE) radio module



Why Embedding?

Edge computing is a method of optimizing **cloud computing systems** by taking some portion of an application, its data, or services away from one or more central nodes (the "core") to the other logical extreme (the "edge") of the Internet which makes contact with the physical world or end users [Wikipedia]

- **A critical problem for (any) Intelligent Wireless Sensor**

Wireless transmission of raw sensor data to the cloud requires a substantial amount of power

All Internet-of-Things Low-Power Wide Area protocols (LPWA: Lora, Sigfox, NB-IoT) but also BLE, ZigBee, etc. are optimized for sparse and infrequent short messages

*Ideally, in an Intelligent Sensor, data processing must be performed **onboard**: short messages should be sent only when relevant events or state transitions occur*

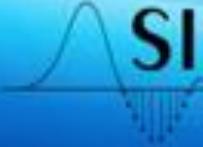
Human Fall (rehearsed)



F01

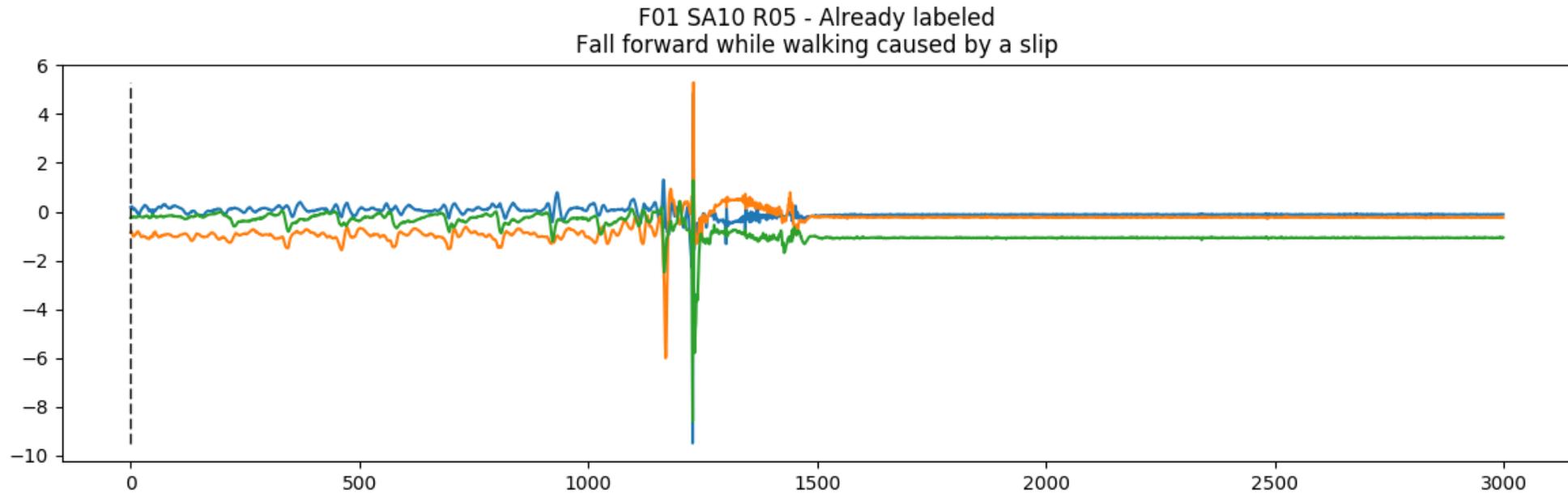
**Fall forward while walking
caused by a slip**


UNIVERSIDAD
DE ANTIOQUIA
1803


SISTEMIC
Sistemas Embebidos e
Inteligencia Computacional

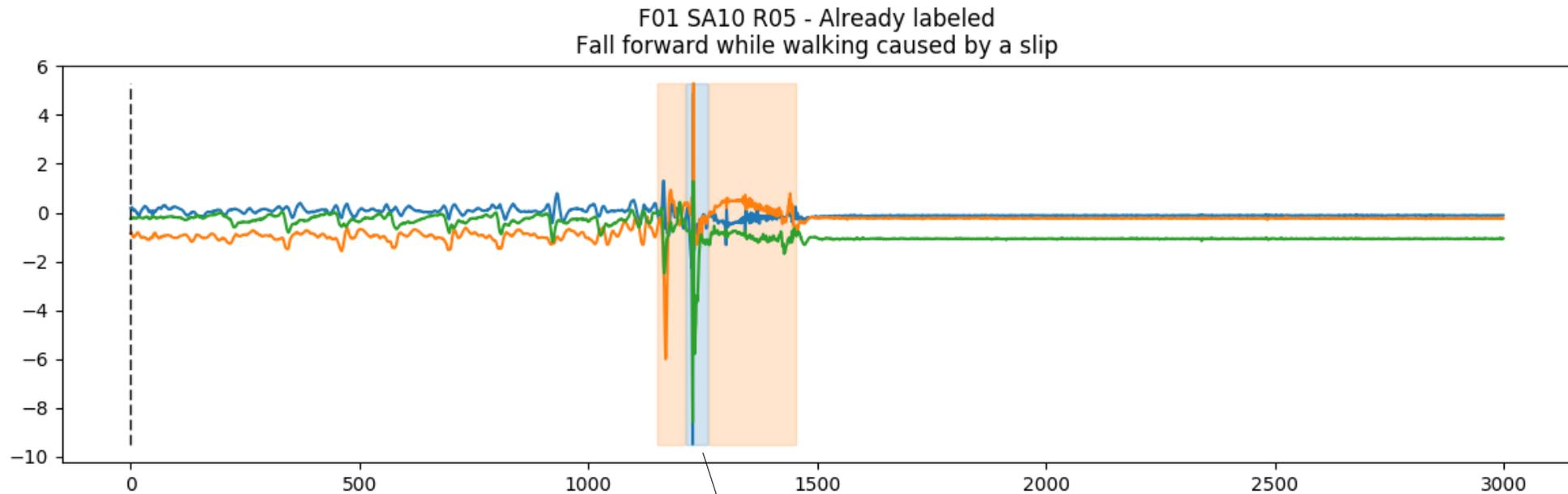
Human Fall (rehearsed)

- Accelerometers



Human Fall (rehearsed)

- Accelerometers (annotated)



the actual fall occurred here

Fall Detection

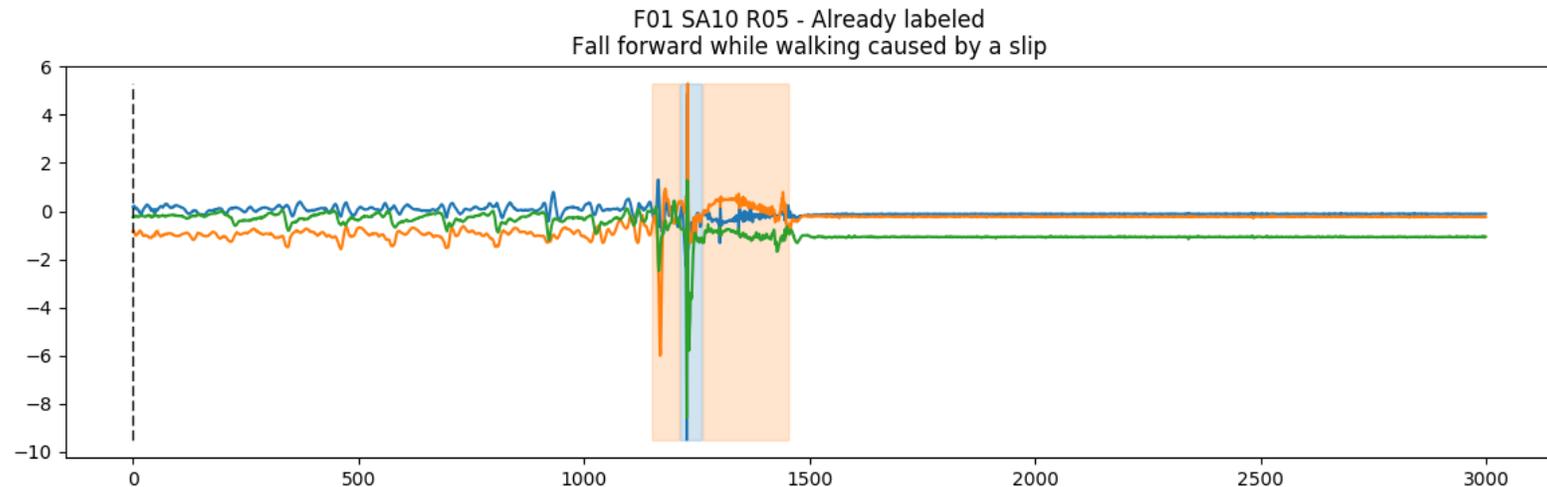
- *Detecting human falls from accelerometer signals is difficult*

The input is time-variant (*unlike a still image*)

Information is limited (*just three scalar values at each time instant*)

It is the 'history' (*i.e. the shape*) of signals that describes the event

False negatives are to be avoided, but even false positives....



Smart Sensors: a Case Study

■ Fall detection with wearable sensors (IPHSDM)

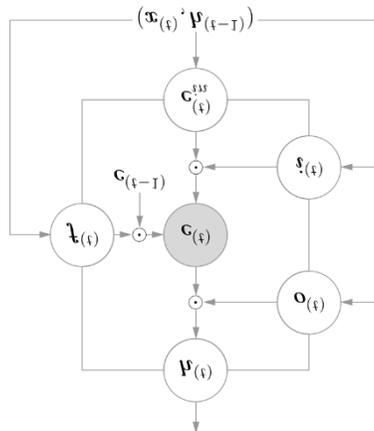
A project co-funded by Regione Lombardia



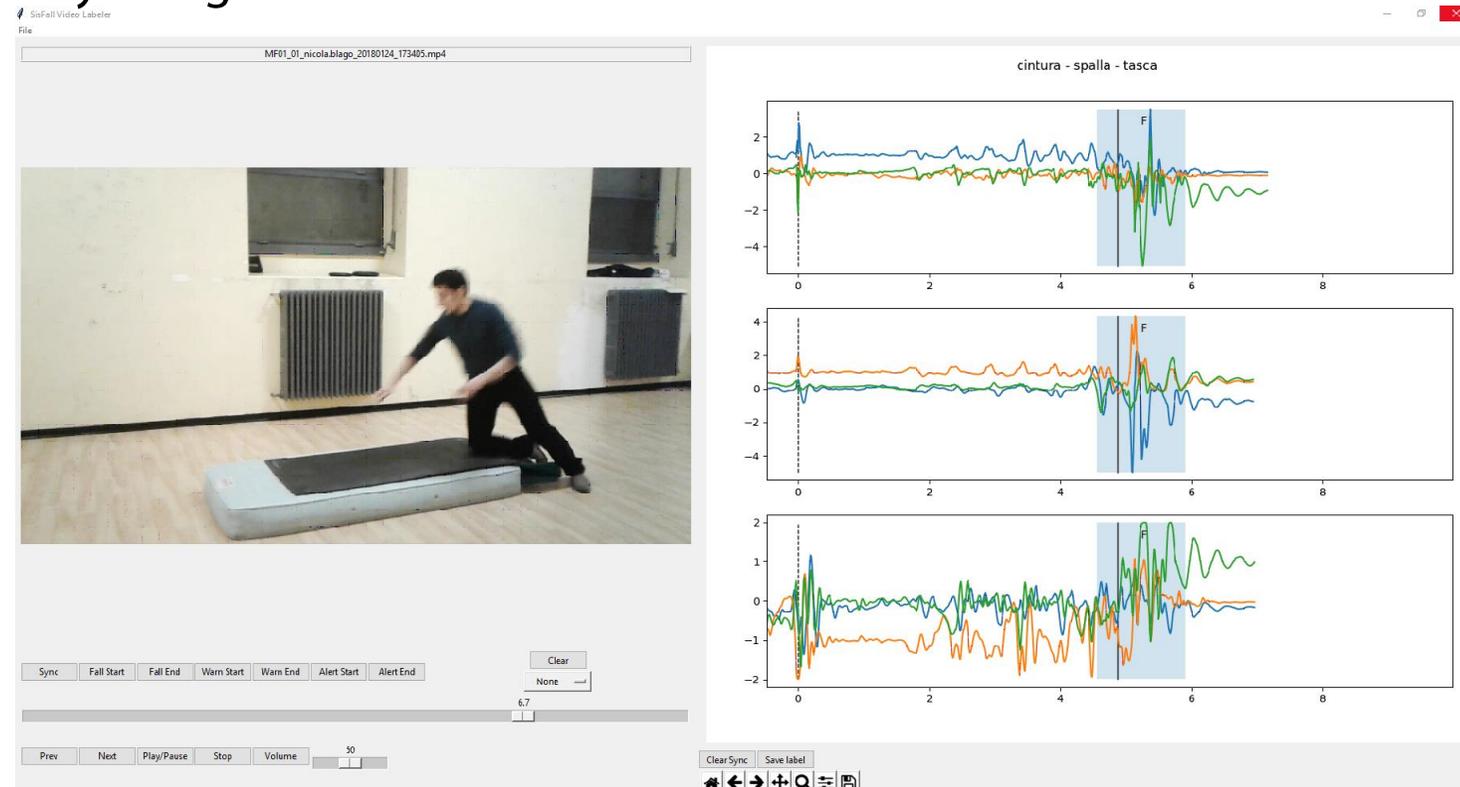
Event detection is performed by the smart sensors (SensorTile)
Messages are sent over BLE only to signal relevant events



Long-Short Term Memory (LSTM)

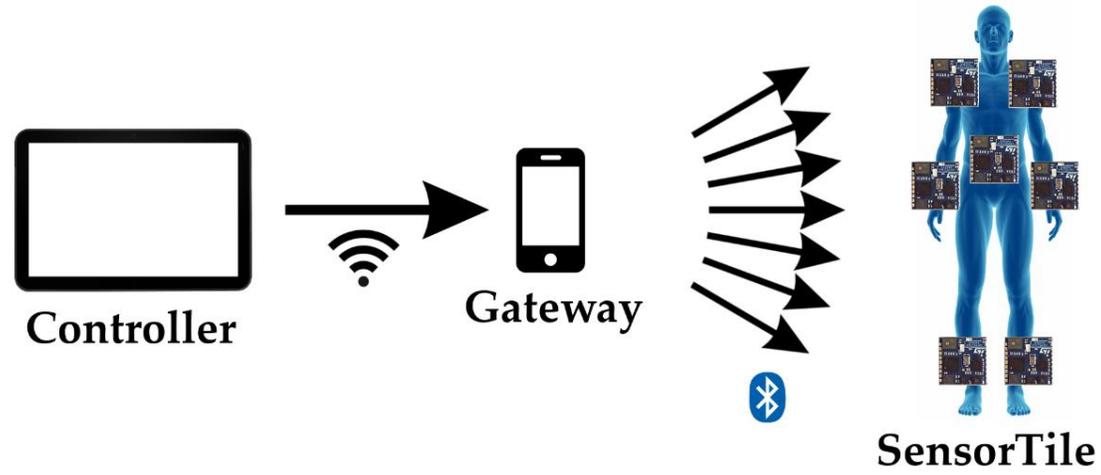


Embedded on Device



Creating a dataset *(for fall detection)*

Body Network



Body Network



Simulated Falls



Dataset Annotation

