

Défis en Intelligence Artificielle

Défi 3 : L'IA pour l'analyse et la prévision de séries temporelles

Souhaib BEN TAIEB



December 19, 2019

Kaggle competition

- <https://www.kaggle.com/c/hands-on-ai-umons-2019/overview>
- <https://www.kaggle.com/bsouhaib/starter>

Deadlines:

- December 19, 11:55pm: At least one Kaggle submission needs to have been made.
- January 8, 11:55pm: The Kaggle competition closes.
- January 22, 11:55pm: Upload to Moodle your project report and code, one per group.

Contents

I	Reccurent Neural Networks	5
II	Convolutional Neural Networks	33
III	Hyperparameter tuning	46

Why deep learning models?

- Deep learning model has been shown perform well in many scenarios
 - 2014 Global Energy Forecasting Competition ([link](#))
 - 2016 CIF International Time Series Competition ([link](#))
 - 2017 Web Traffic Time Series Forecasting ([link](#))
 - 2018 Corporación Favorita Grocery Sales Forecasting ([link](#))
 - 2018 M4-Competition ([link](#))
 - Flexible and expressive
 - Easily inject exogenous features into the model
 - Learn from large time series datasets
-

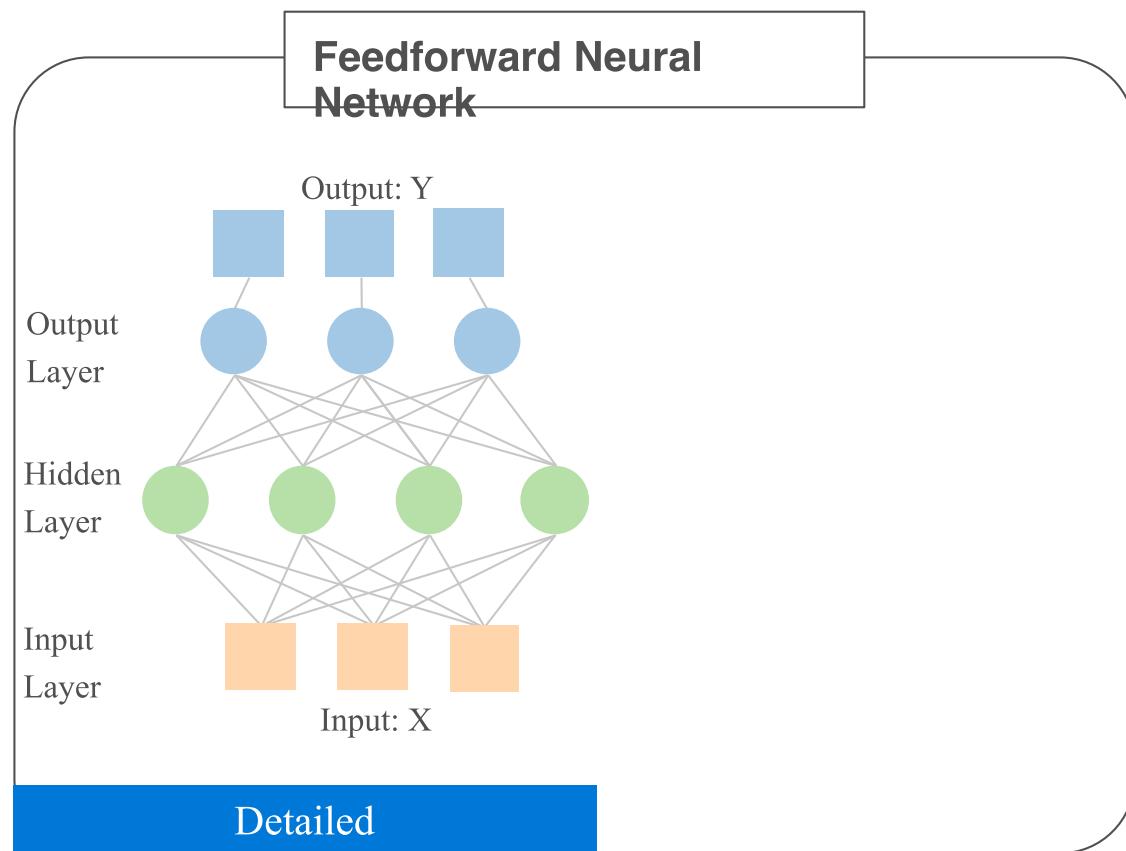
Part I

Recurrent Neural Networks

Agenda

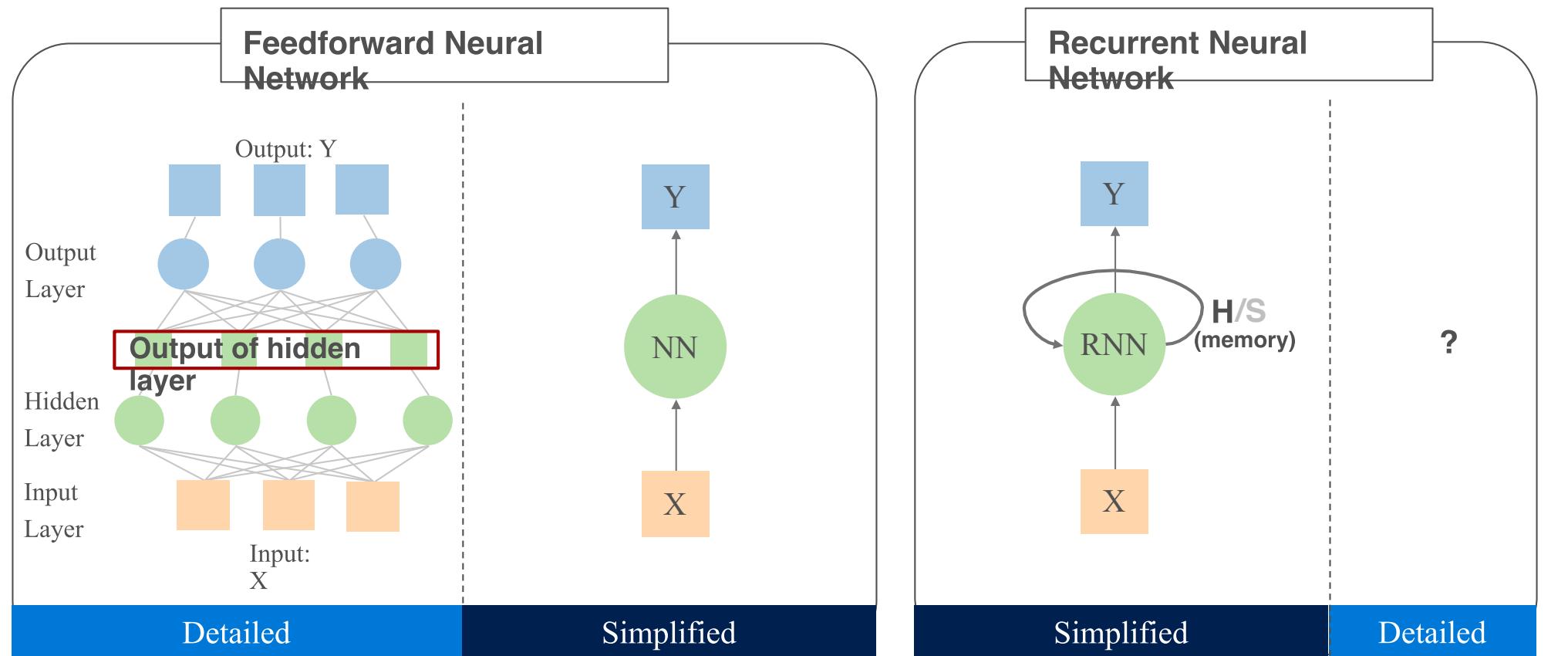
- What are RNNs?
 - How RNNs are trained: Backpropagation through time (BPTT)
 - Vanilla RNN and its gradient problems
 - Other RNN units
 - GRU
 - LSTM
 - RNN stacking
 - RNN for one step time series forecasting
 - Encode-decoder RNN for multi-step time series forecasting
-

What are RNNs?



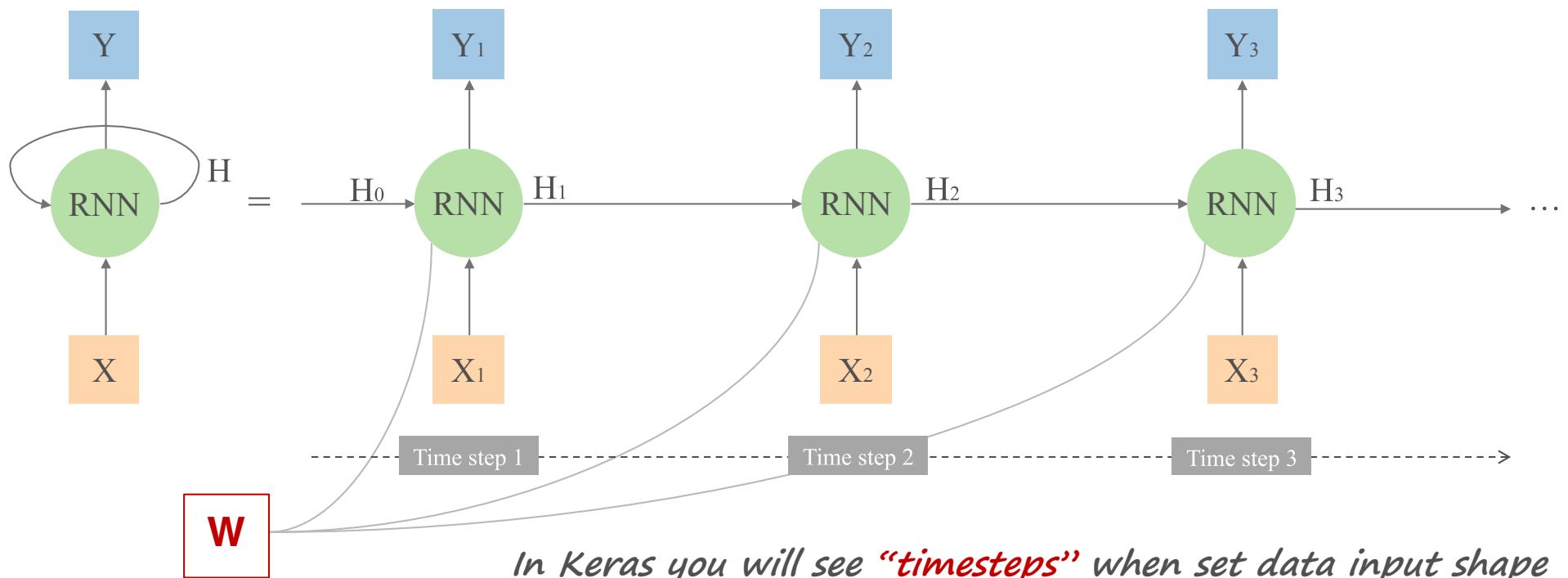
What are RNNs?

RNN has internal hidden state which can be fed back to network

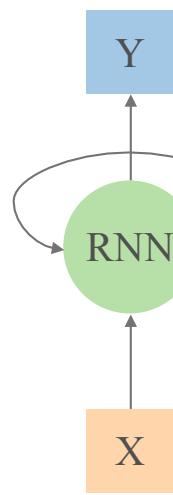


Unrolled RNN

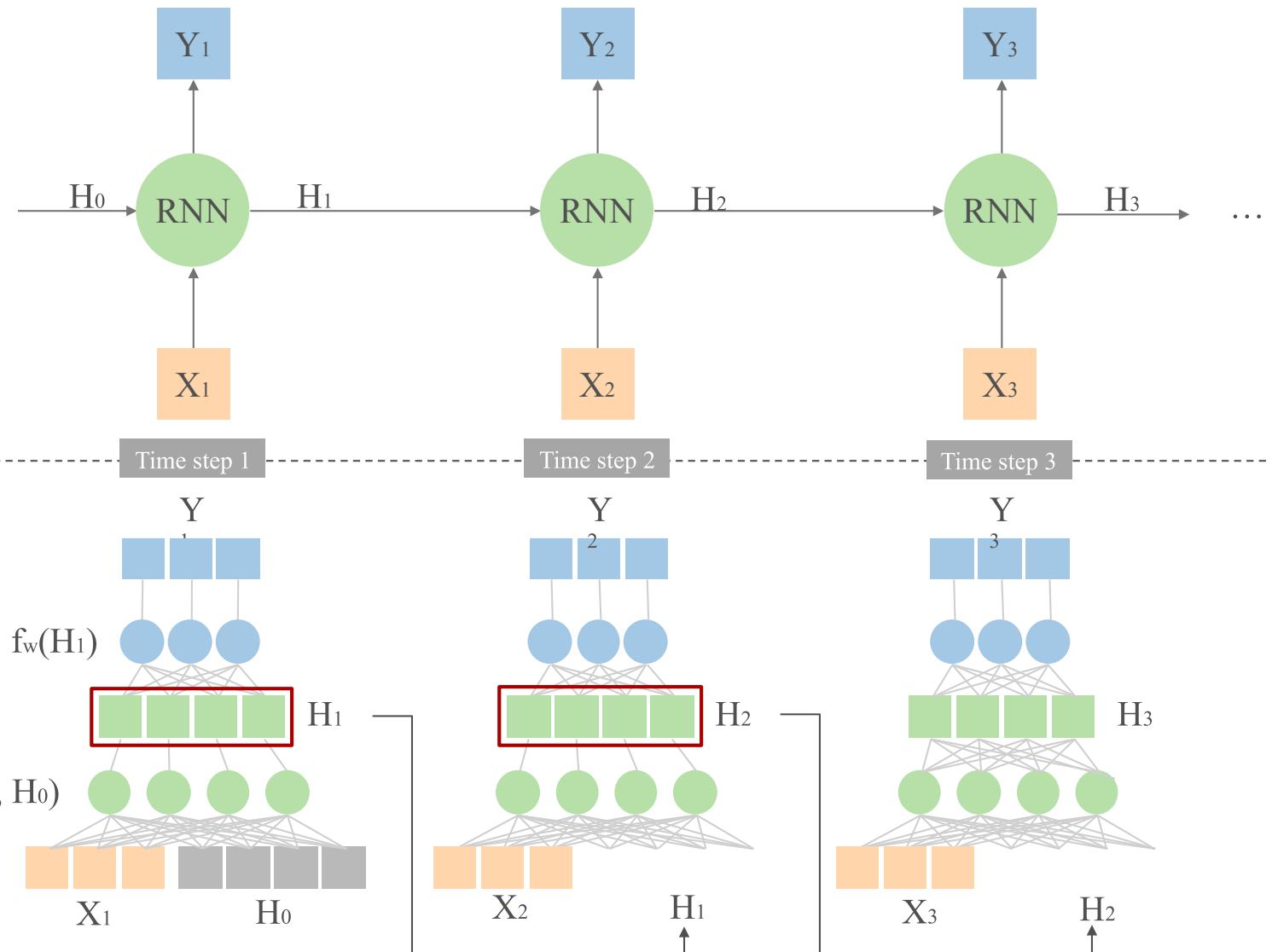
The same weight and bias shared across all the steps



Unrolled RNN



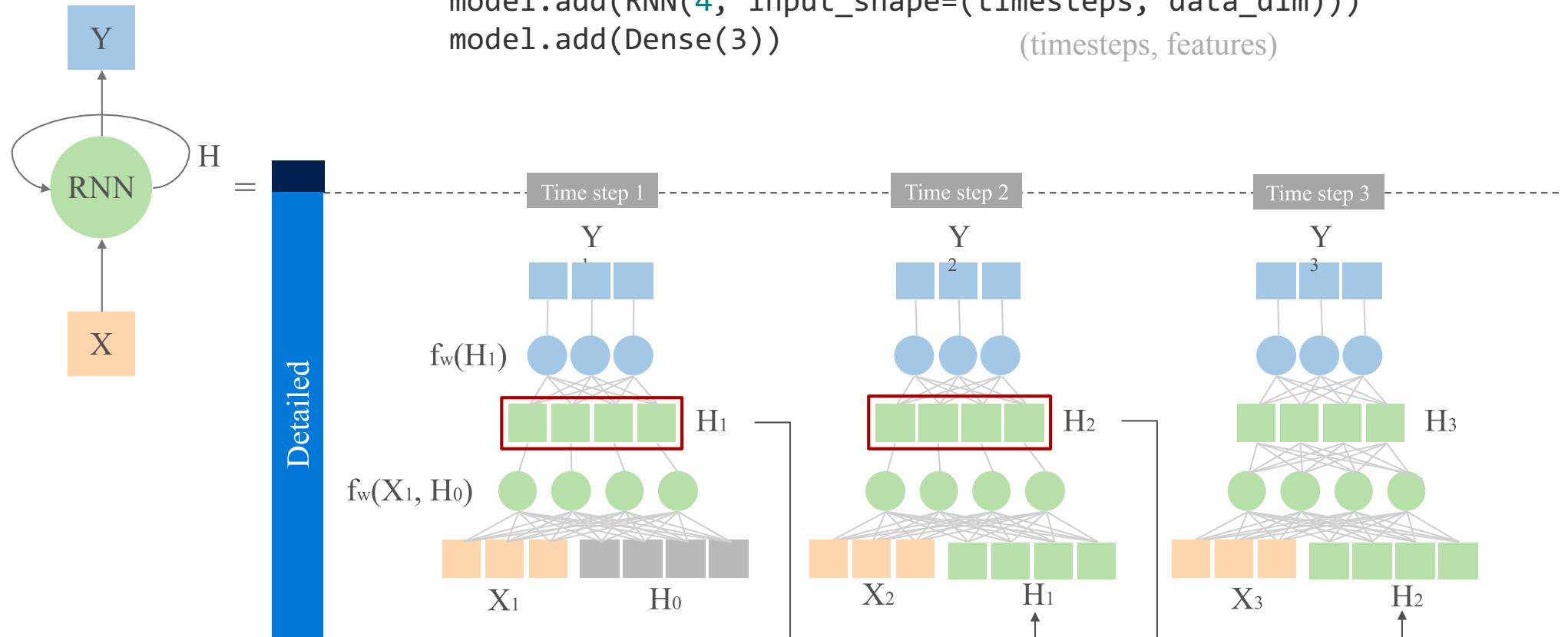
= **Simplified** **Detailed**



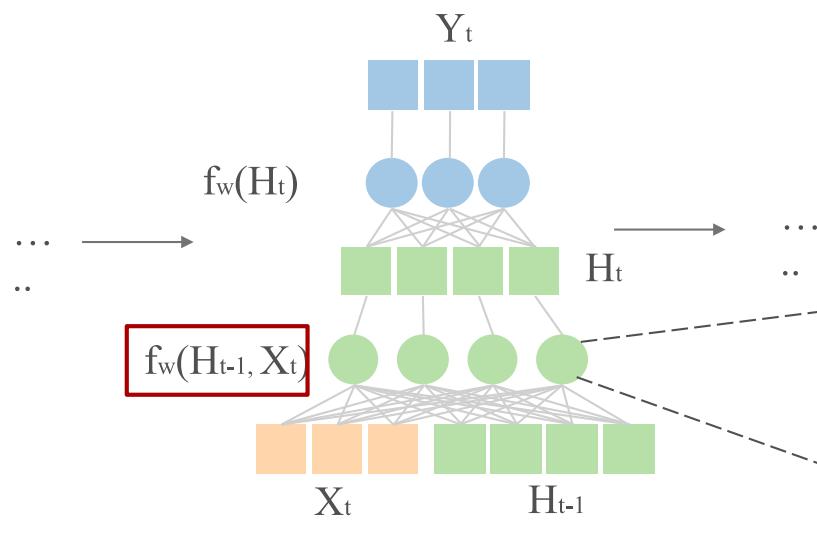
Unrolled RNN

*In Keras, the parameter “units” is dimensions of hidden state.
(Think of it as feedforward neural network number of units in hidden layer.)*

```
model = Sequential()  
model.add(RNN(4, input_shape=(timesteps, data_dim)))  
model.add(Dense(3))  
                                         (timesteps, features)
```



Vanilla RNN



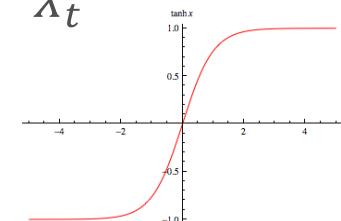
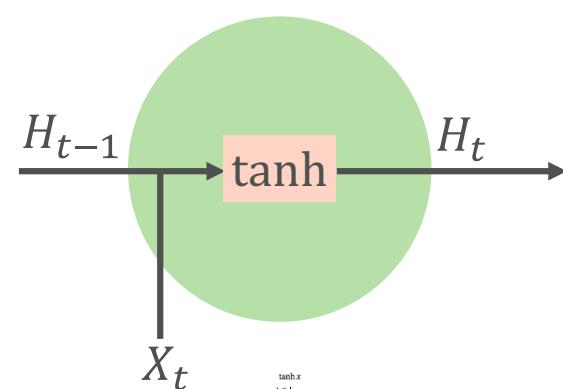
$$H_t = f_W(H_{t-1}, X_t)$$

some function with
parameter W

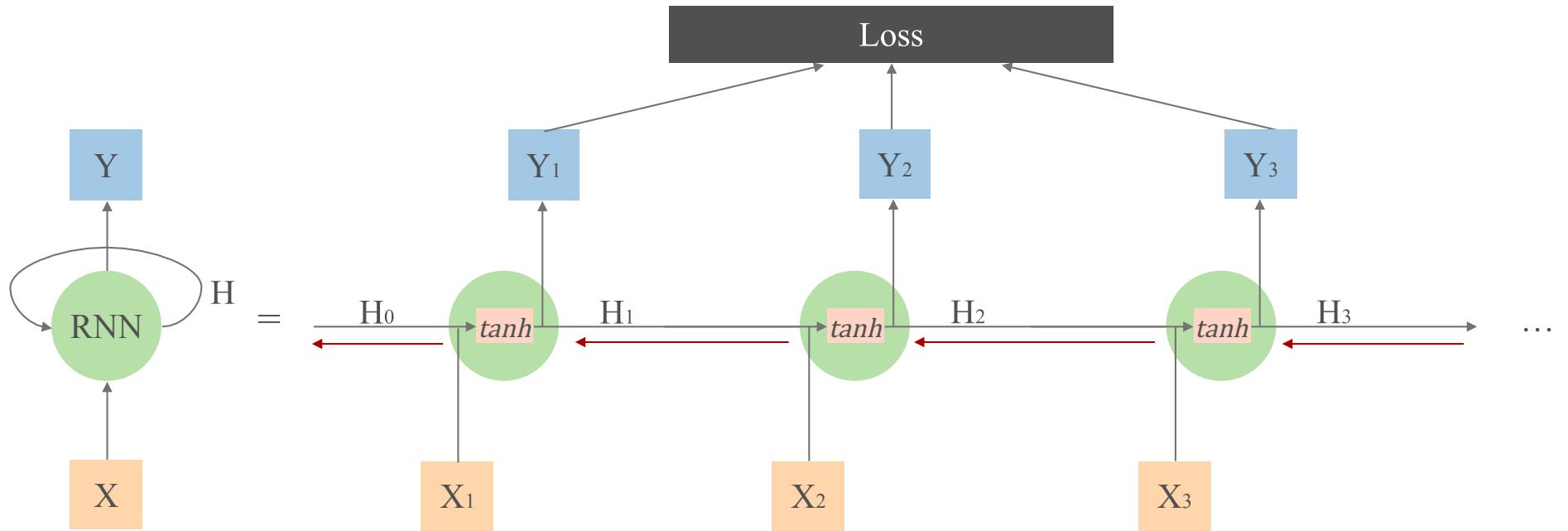
|
 new state
 |
 old state
 |
 Input vector at
 some time step

Vanilla RNN:

$$\begin{aligned} H_t &= \tanh(W_h H_{t-1} + W_x X_t) \\ &= \tanh(\mathbf{W} \cdot [H_{t-1}, X_t]) \end{aligned}$$



Vanilla RNN BPTT

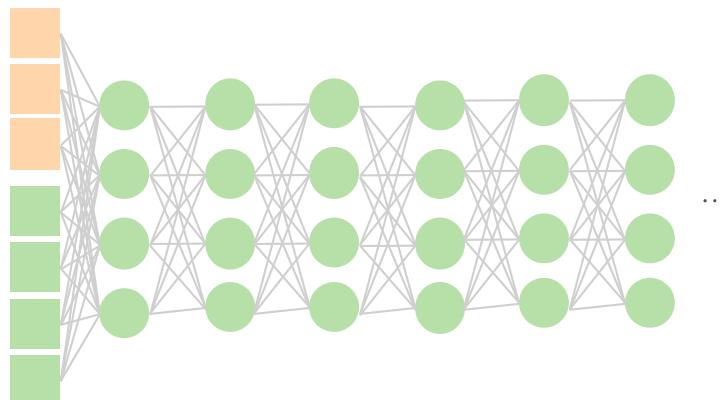


Computing gradient of h_0 involves repeated tanh and many factors of W

Vanilla RNN Gradient Problems

Computing gradient of h_0 involves repeated tanh and many factors of W which causes:

- **Exploding gradient** (e.g. $5*5*5*5*5*5*.....$)
- **Vanishing gradients** (e.g. $0.7*0.7*0.7*0.7*0.7*0.7*.....$)



100 time steps is similar to 100 layers feedforward neural net

Exploding Gradient

- Exploding gradients are obvious. Your gradients will become NaN (not a number) and your program will crash
- Solution: Gradient clipping
 - Clip the gradient when it goes higher than a threshold

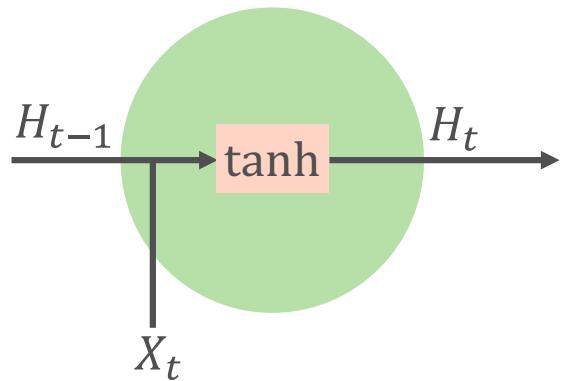
Vanishing Gradient

- Vanishing gradients are more problematic because it's not obvious when they occur or how to deal with them
- Solutions:
 - Change activation function to ReLU
 - Proper initialization
 - Regularization
 - Change architecture to LSTM or GRU

Gated Recurrent Unit (GRU)

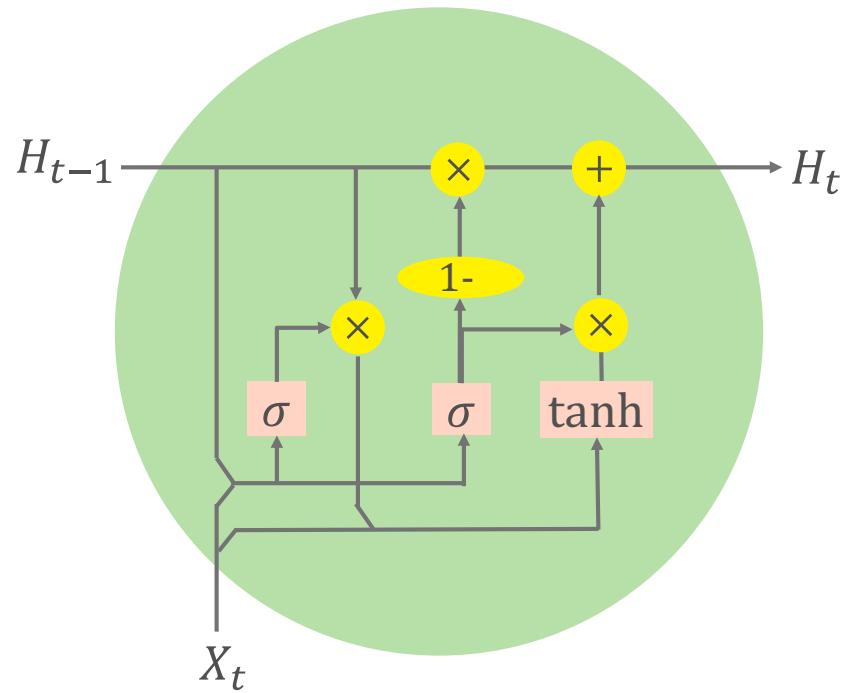
Vanilla RNN:

$$H_t = \underline{\tanh(\mathbf{W} \cdot [H_{t-1}, X_t])}$$

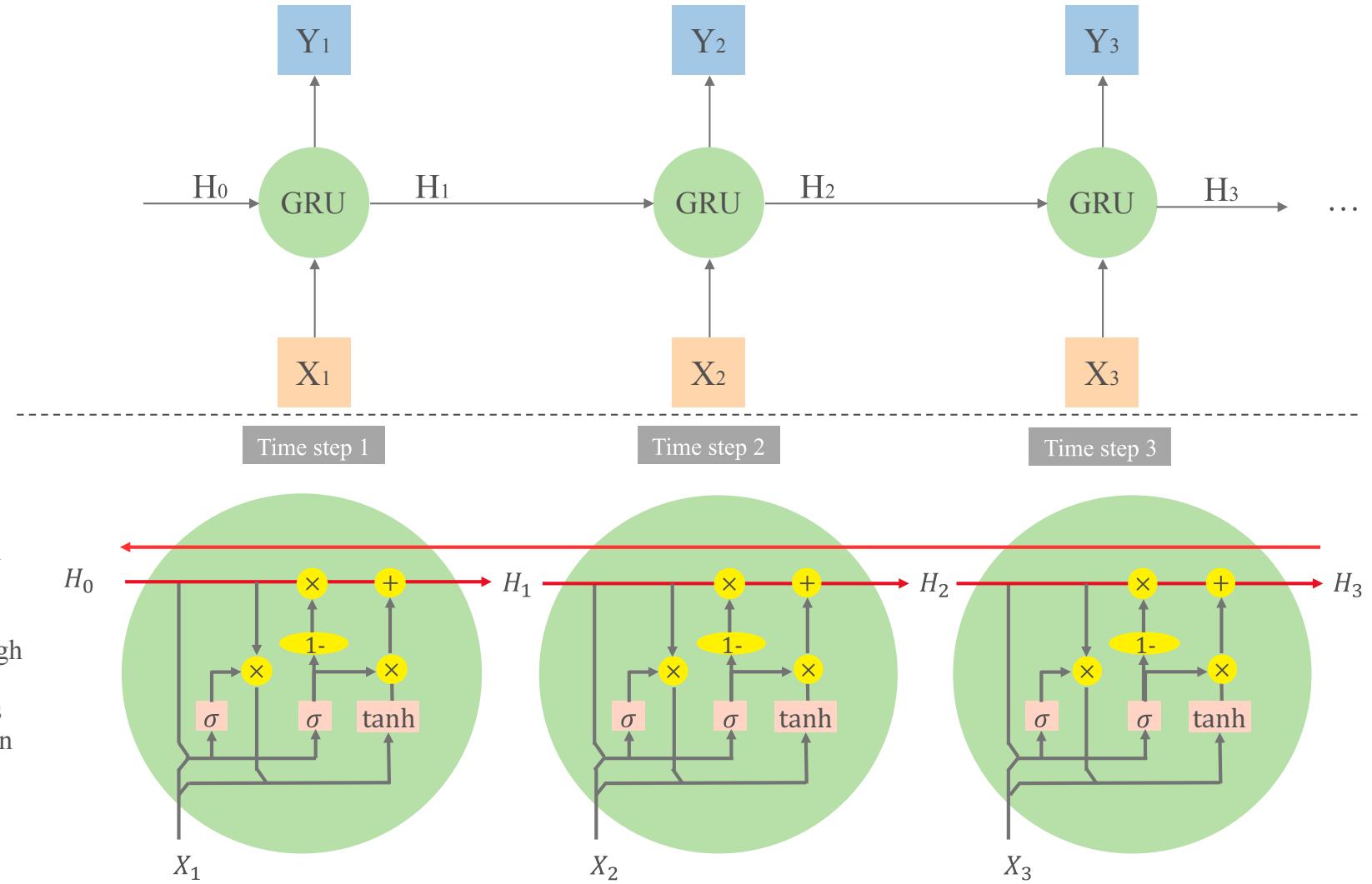


GRU:

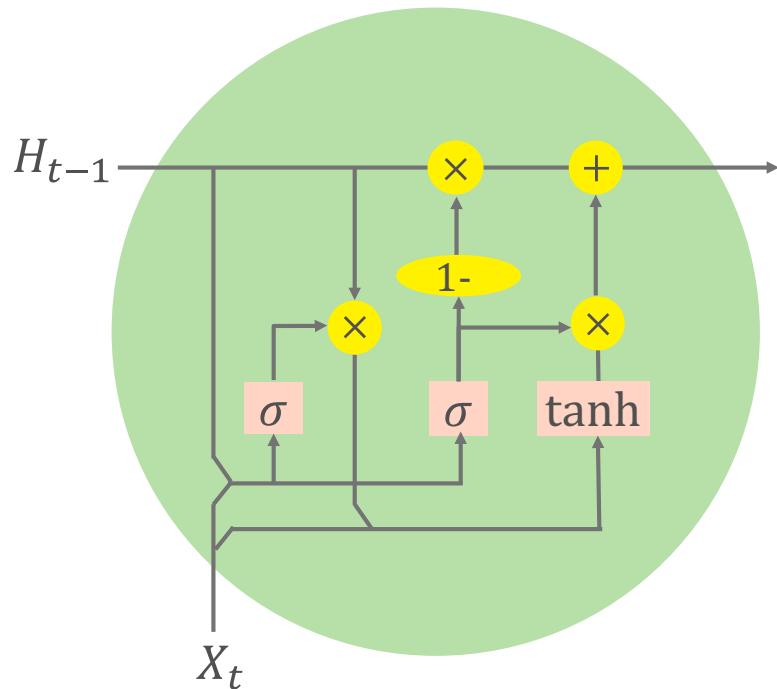
$$H_t = \underline{(1 - z_t) \times H_{t-1} + z_t \times \tilde{H}_t}$$



GRU



Gated Recurrent Unit (GRU)



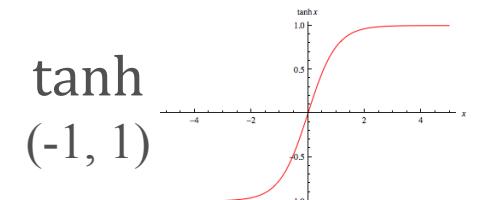
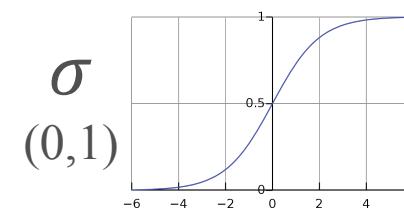
GRU [[Learning phrase representations using rnn encoderdecoder for statistical machine translation, Cho et al. 2014](#)]

$$\text{Hidden state: } H_t = (1 - z_t) \times H_{t-1} + z_t \times \tilde{H}_t$$

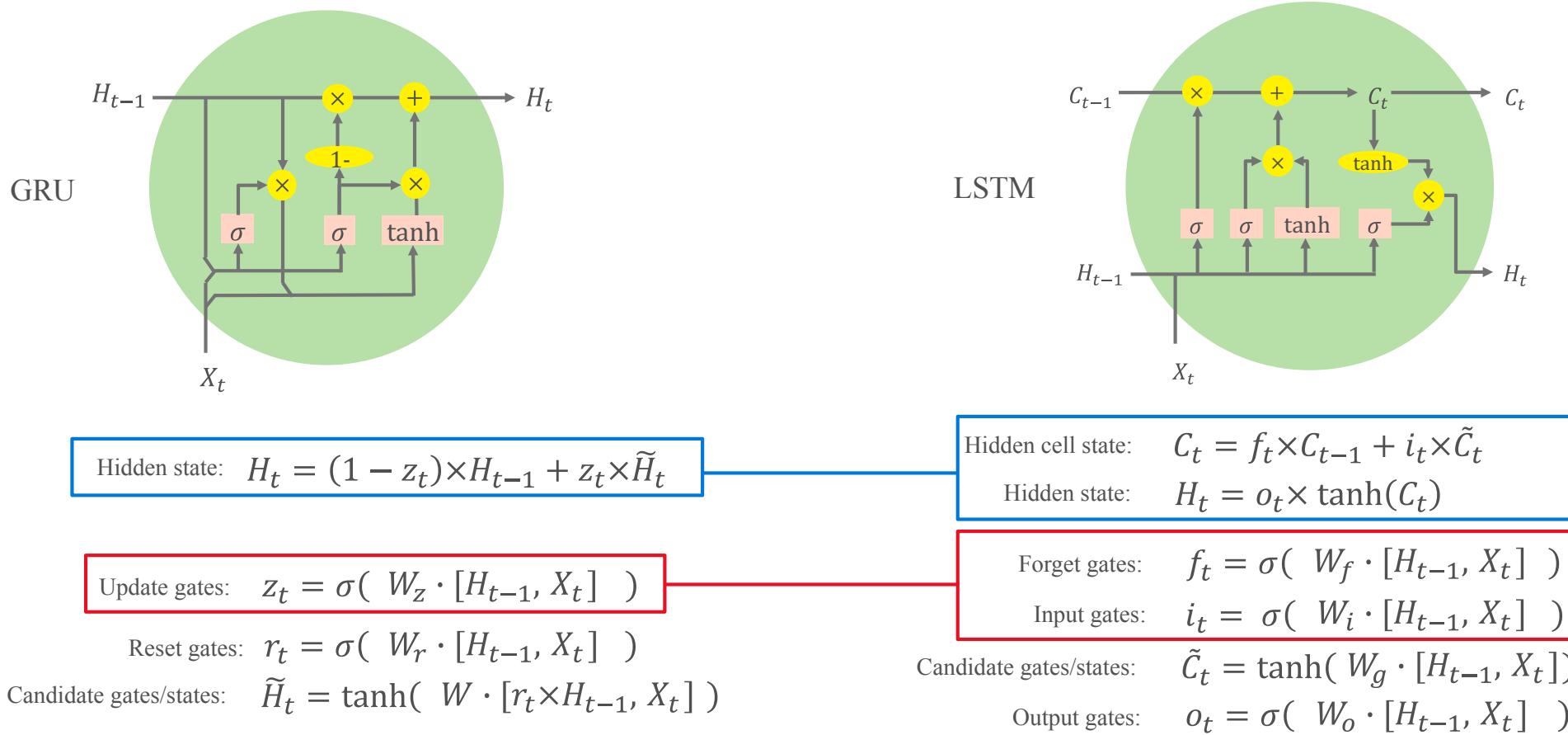
$$\text{Update gates: } z_t = \sigma(W_z \cdot [H_{t-1}, X_t])$$

$$\text{Candidate gates/states: } \tilde{H}_t = \tanh(W \cdot [r_t \times H_{t-1}, X_t])$$

$$\text{Reset gates: } r_t = \sigma(W_r \cdot [H_{t-1}, X_t])$$



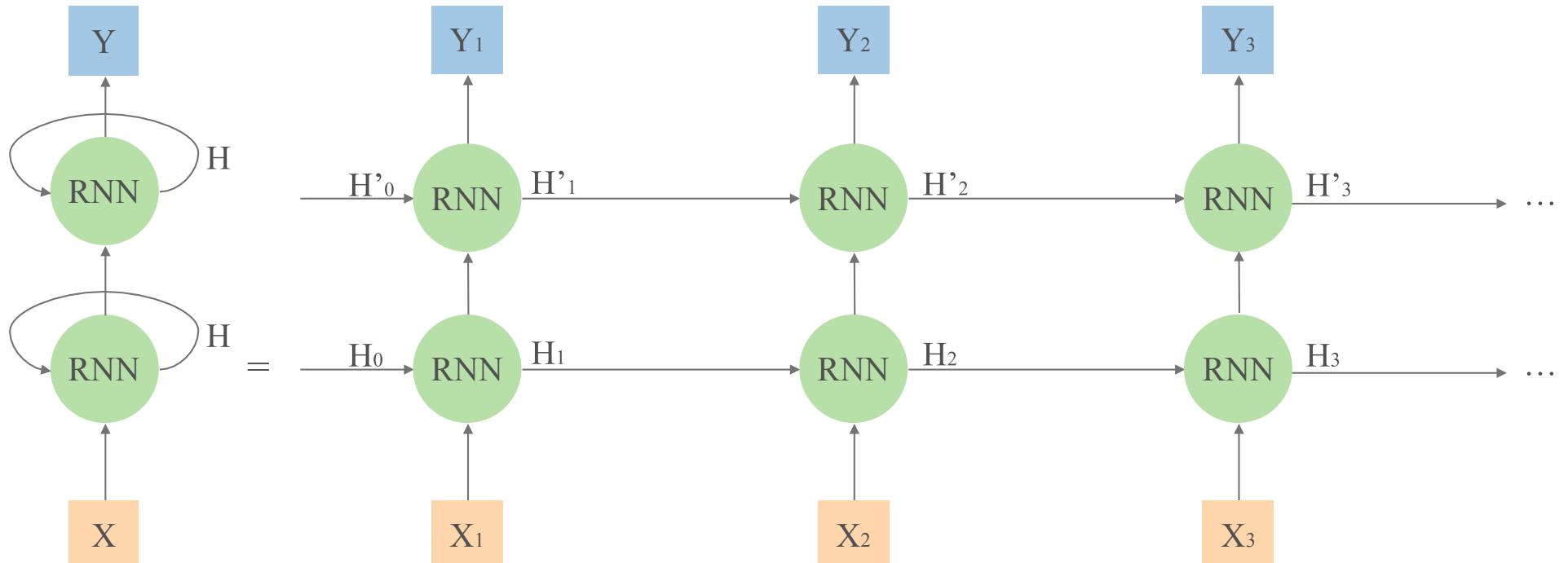
GRU vs LSTM (Long Short Term Memory)



[LSTM: A Search Space Odyssey, Greff et al., 2015](#) [Long Short-Term Memory, Hochreiter & Schmidhuber \(1997\)](#)

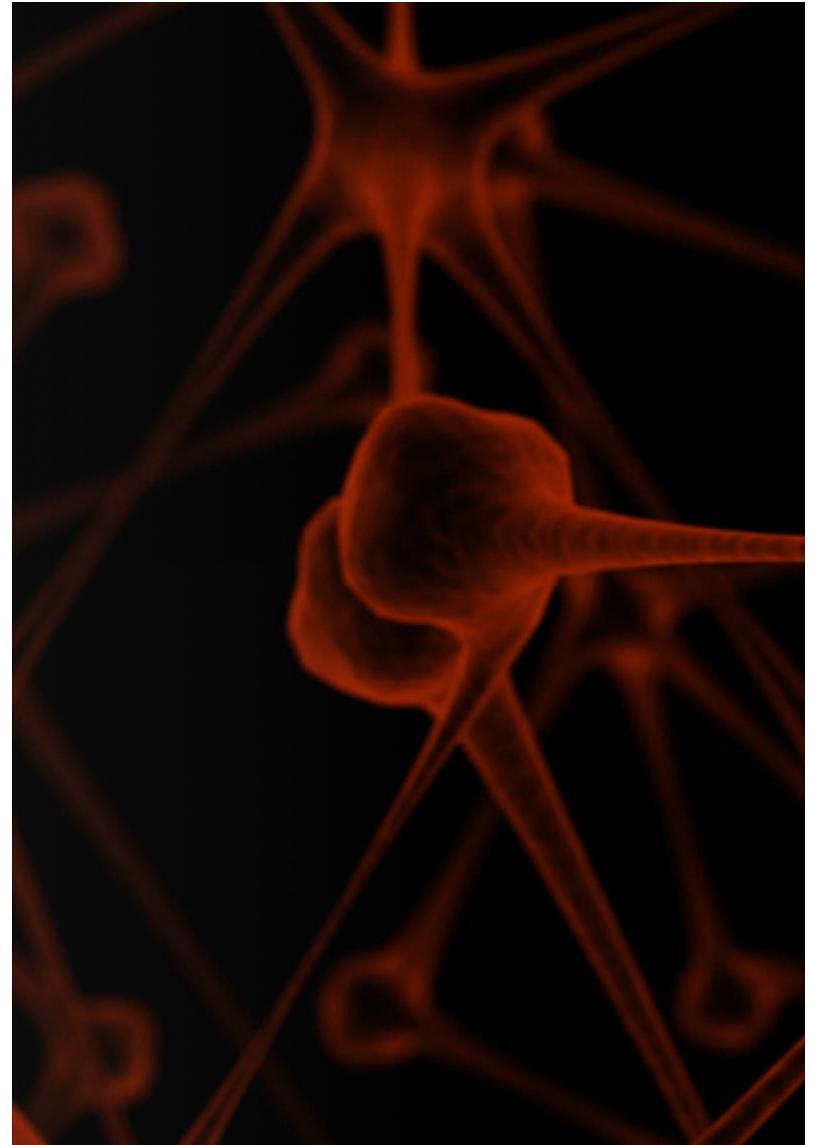
RNN Stacking

To learn more complex relationships, we can go deep by stacking the cells.



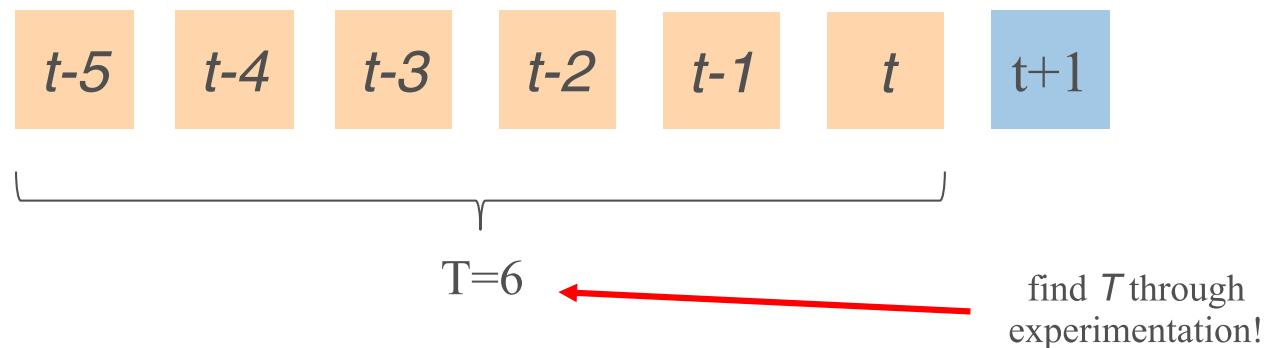
```
model = Sequential()  
model.add(RNN(4, return_sequences=True, input_shape=(timesteps, data_dim)))  
model.add(RNN(4))
```

RNN for one step time
series forecasting

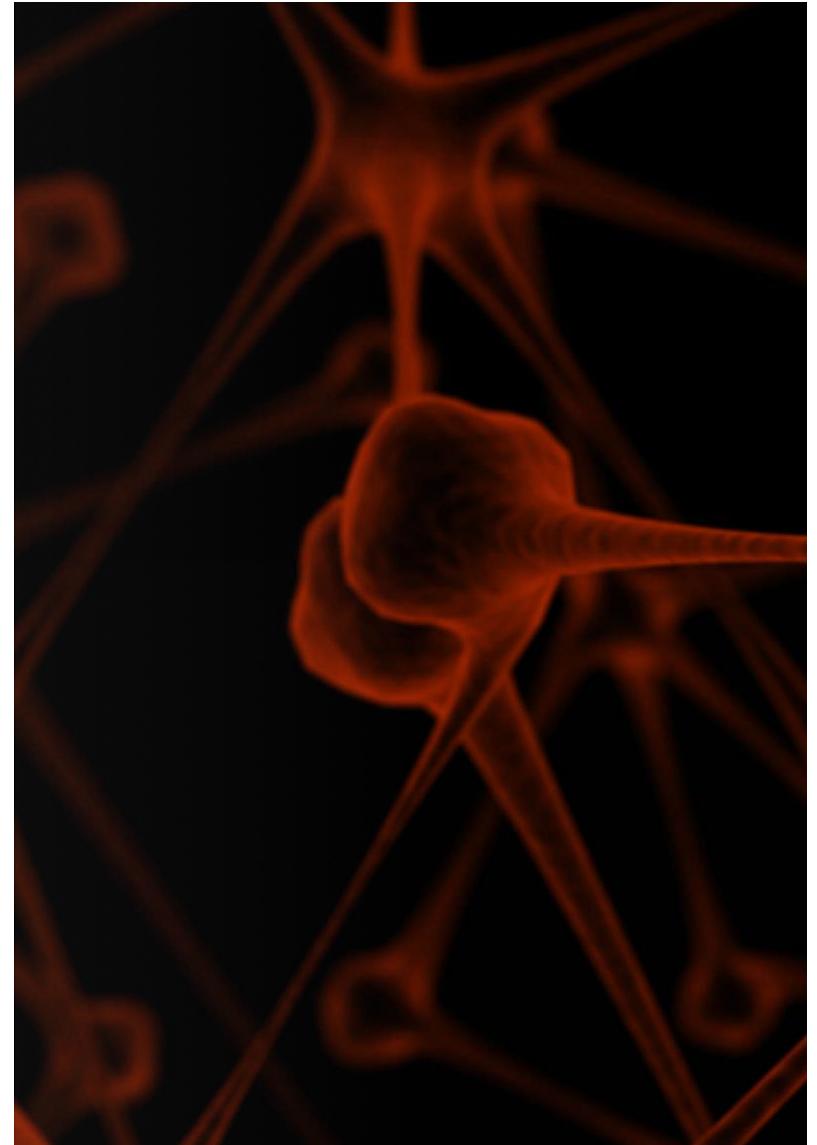


One-step forecast

- Assuming we are at time t ...
- ... predict the value at time $t+1$...
- ... conditional on the previous T values of the time series



RNN for multi-step time series forecasting

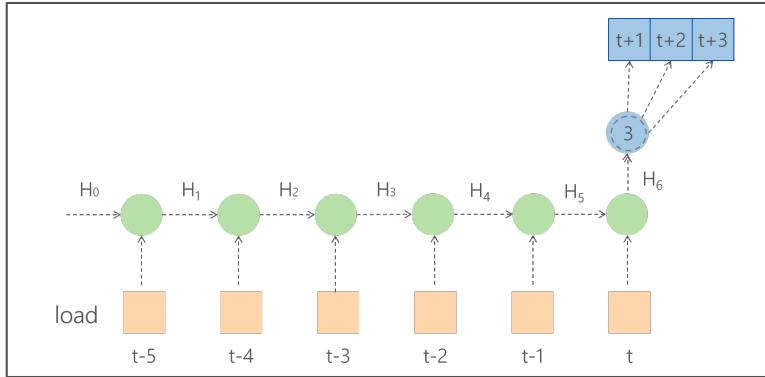


Multi-step forecast

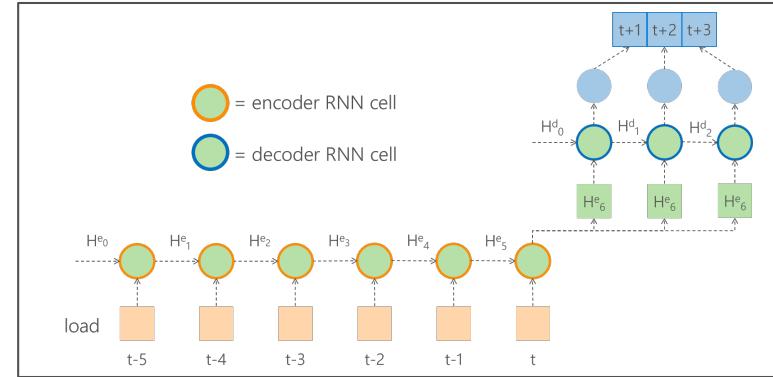
- Assuming we are at time t ...
- ... predict the values at times $(t+1, \dots, t+HORIZON)$...
- ... conditional on the previous T values of the time series



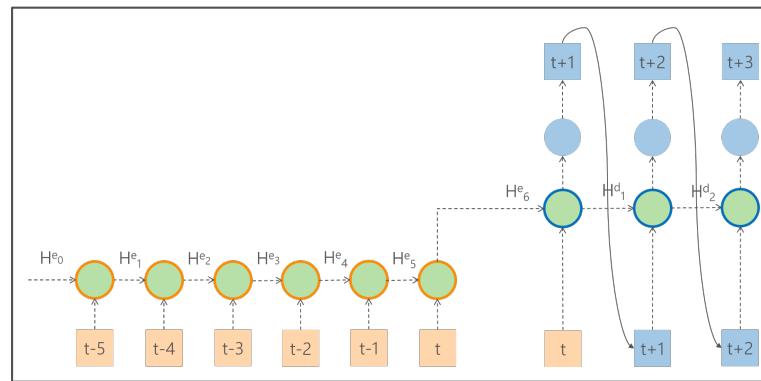
Multi-step forecast



Vector output

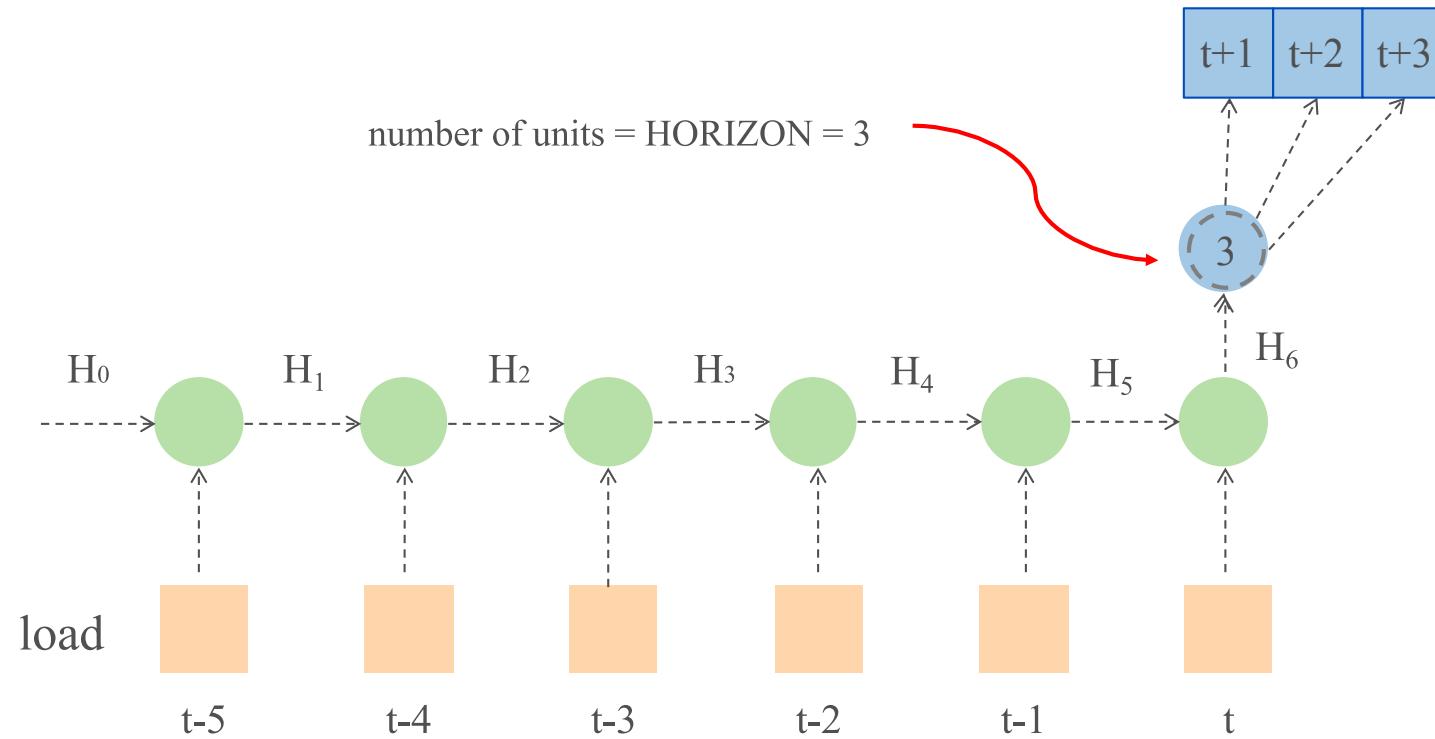


Simple encoder-decoder



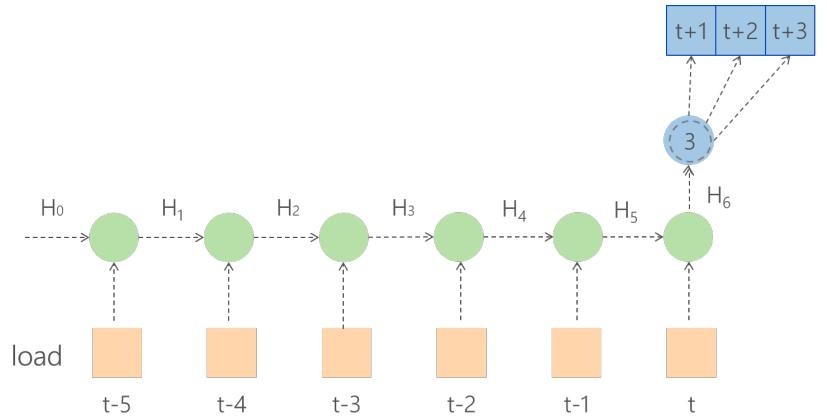
Recursive encoder-decoder

Vector output approach



Vector output approach

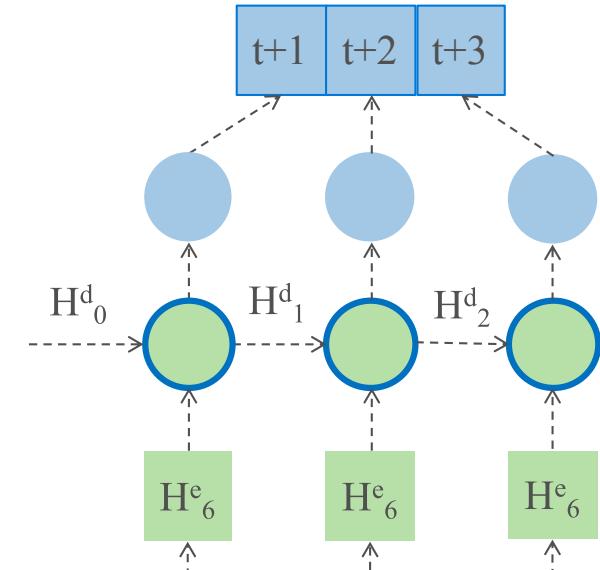
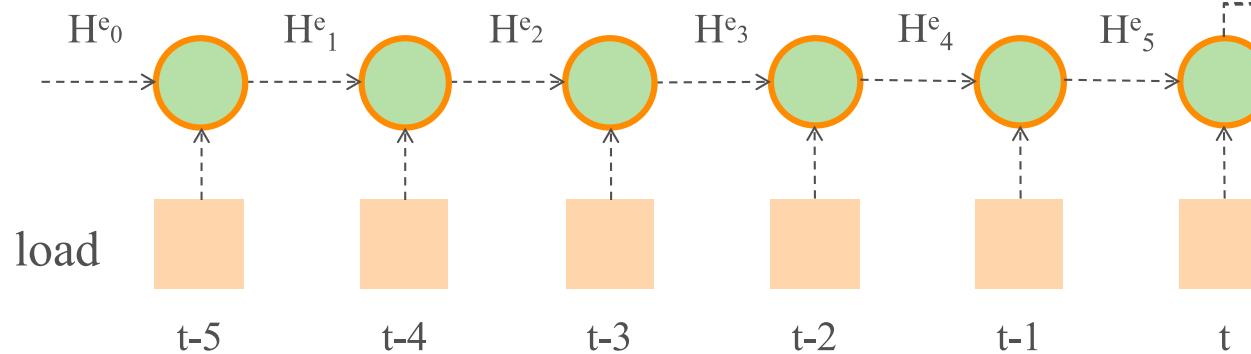
- 👍 Simplest to implement
- 👍 Fastest to train
- 👎 The forecasts are not dynamic



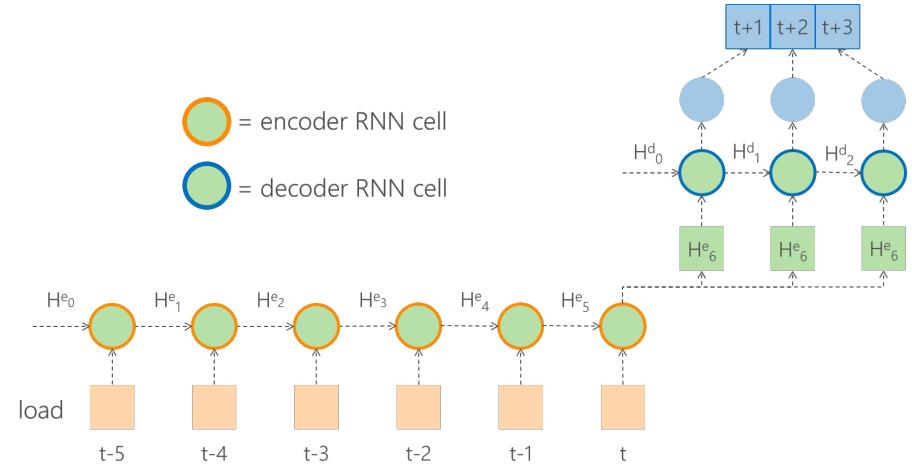
Simple encoder-decoder

 = encoder RNN cell

 = decoder RNN cell



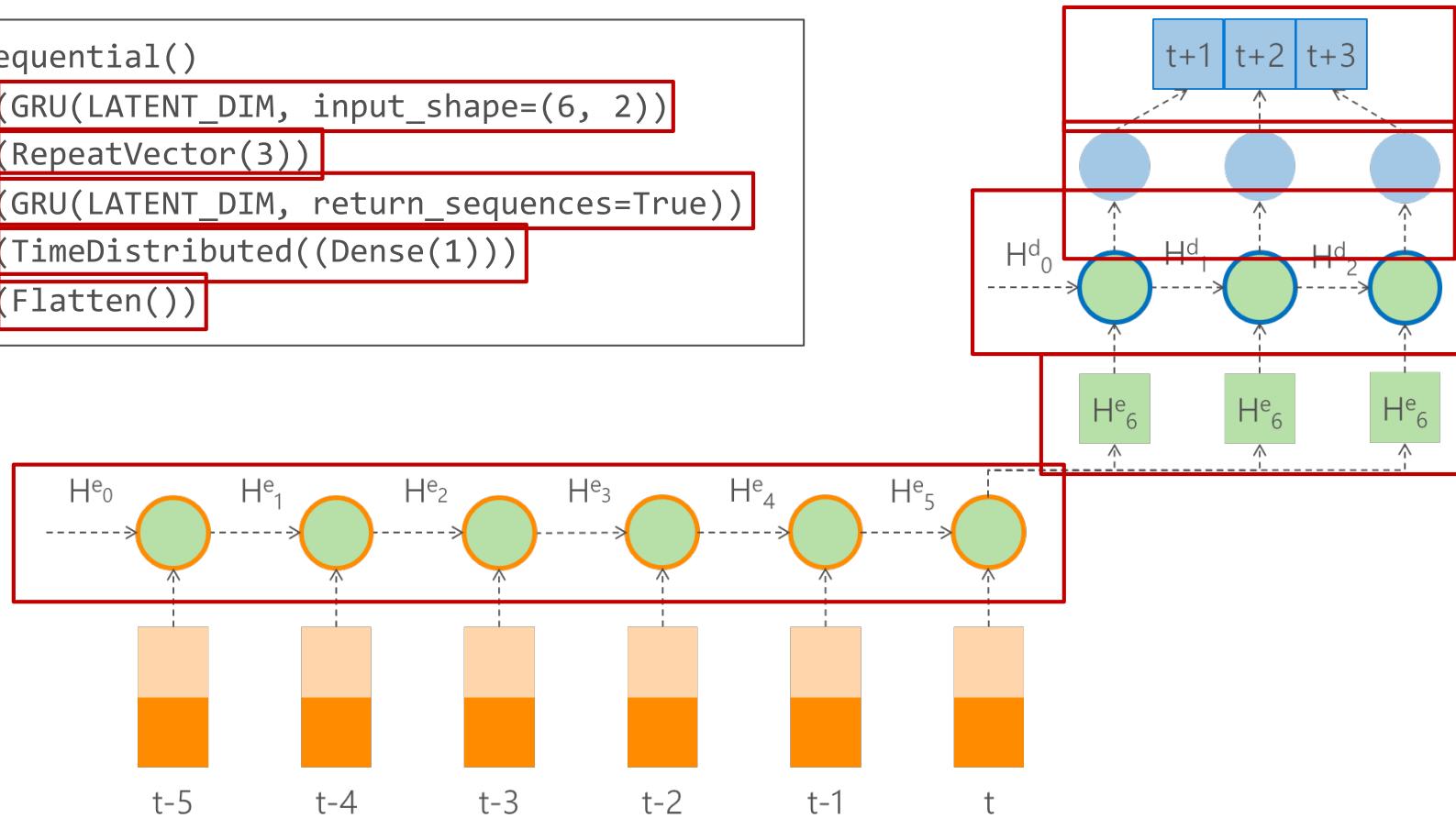
Simple encoder-decoder



- 👍 Fairly simple to implement
- 👍 Tries to capture dependencies between forecasted time steps through decoder hidden state
- 👎 Slower to train with stacked RNN layers

Simple encoder-decoder

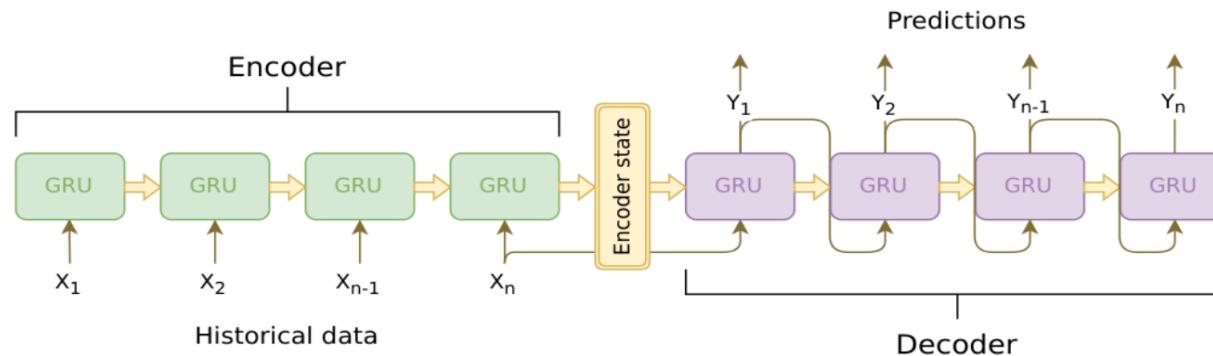
```
model = Sequential()  
model.add(GRU(LATENT_DIM, input_shape=(6, 2)))  
model.add(RepeatVector(3))  
model.add(GRU(LATENT_DIM, return_sequences=True))  
model.add(TimeDistributed(Dense(1)))  
model.add(Flatten())
```



Web traffic forecasting competition

Forecast traffic of 145K Wikipedia pages

Winning solution: github.com/Arturus/kaggle-web-traffic

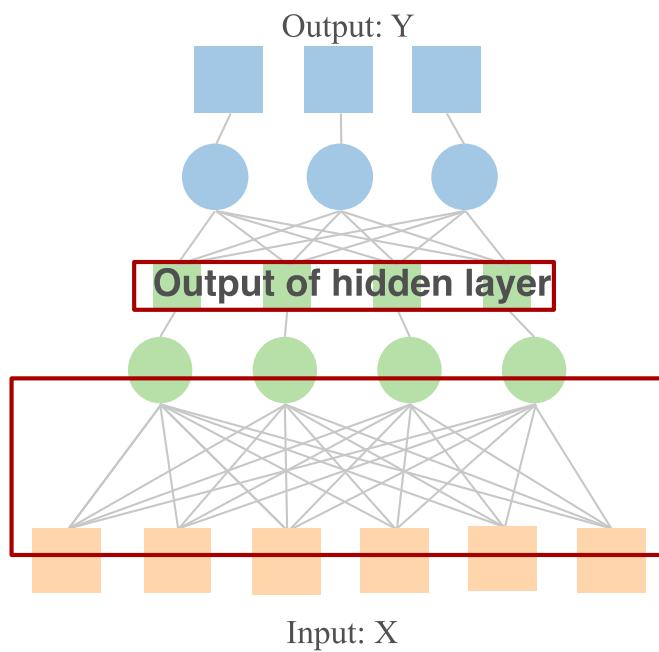


- Feature engineering: transform univariate to multivariate time series by adding seasonality features:
$$x_t \rightarrow (x_t, x_{t-\text{quarter}}, x_{t-\text{year}})$$
- COCOB optimizer that doesn't require learning rate tuning and converges considerably faster, custom implementation
- Sophisticated technique for building ensemble of RNN models, mainly for reducing model variance

Part II

Convolutional Neural Networks

Fully connected vs convolution layer



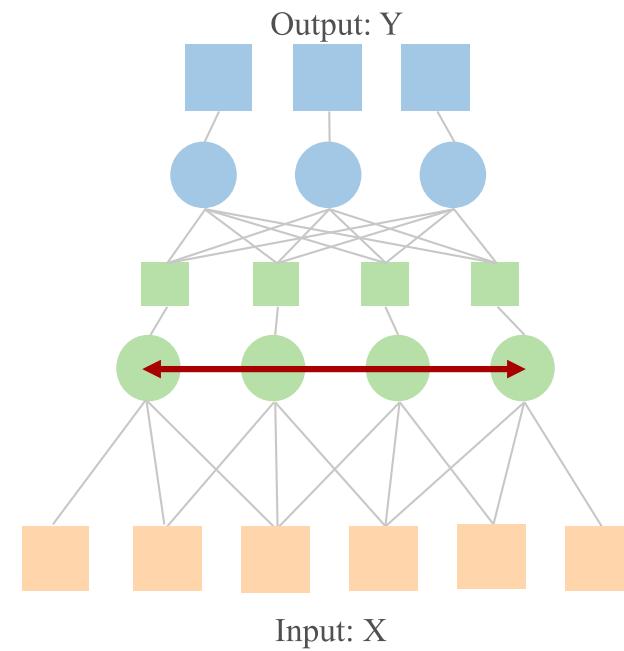
Fully connected:

- units in hidden layer are connected to every unit in previous layer

Fully connected vs convolution layer

Convolution layer:

- units in hidden layer operate on a field of the input
- weights are shared across input



1D Convolutions

- Apply a 1D filter to all elements of a 1D input vector
- Result is the sum of the element-wise product

$$\begin{array}{c} \boxed{5} \boxed{3} \boxed{2} \boxed{7} \boxed{1} \boxed{6} \\ \times \quad \boxed{1} \boxed{0} \boxed{-1} \\ = \quad \boxed{3} \boxed{-4} \boxed{1} \boxed{1} \end{array}$$

$5 \times 1 + 3 \times 0 + 2 \times -1 = 3$

Input: 1 x 6 Filter: 1 x 3 Output: 1 x 4

1D Convolutions

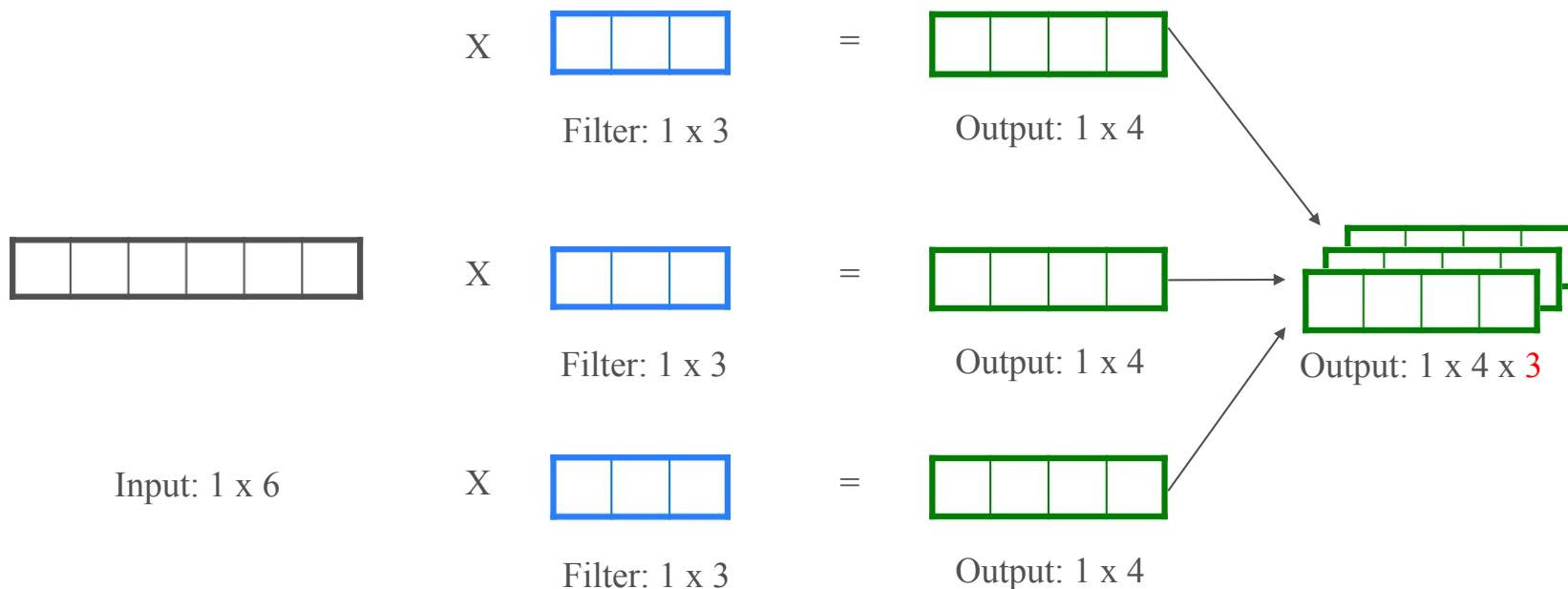
- Filters are trained to detect features in the input sequence
- Features can be detected regardless of where they appear in the input sequence

$$\begin{array}{|c|c|c|c|c|c|} \hline 5 & 3 & 2 & 7 & 1 & 6 \\ \hline \end{array} \quad X \quad \begin{array}{|c|c|c|} \hline \mathbf{w_1} & \mathbf{w_2} & \mathbf{w_3} \\ \hline \end{array} \quad = \quad \begin{array}{|c|c|c|c|} \hline 3 & 4 & 1 & 1 \\ \hline \end{array}$$

Input: 1 x 6 Filter: 1 x 3 Output: 1 x 4

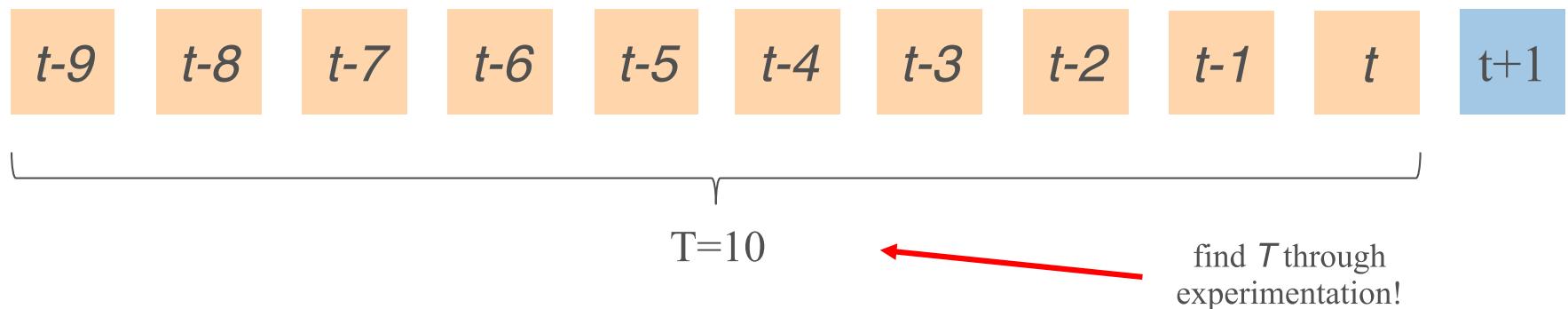
1D Convolutions

- Apply multiple filters to the input data to detect multiple features



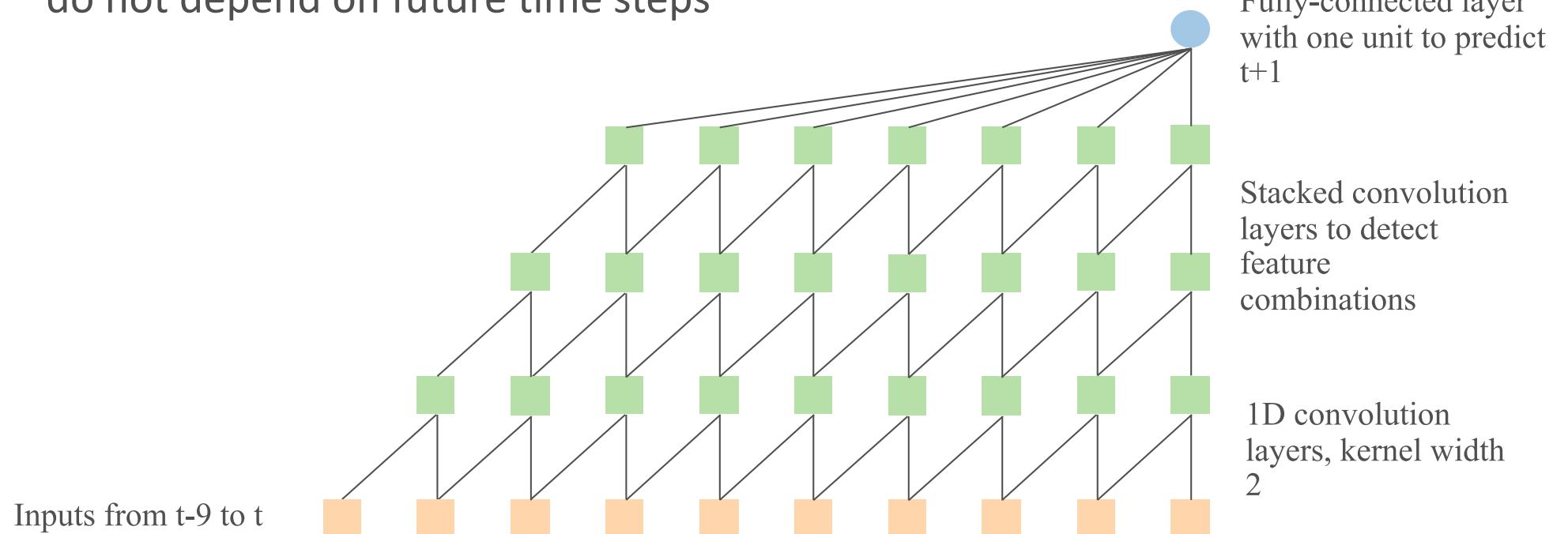
Application to forecasting

- Assuming we are at time t ...
- ... predict the value at time $t+1$...
- ... conditional on the previous T values of the time series



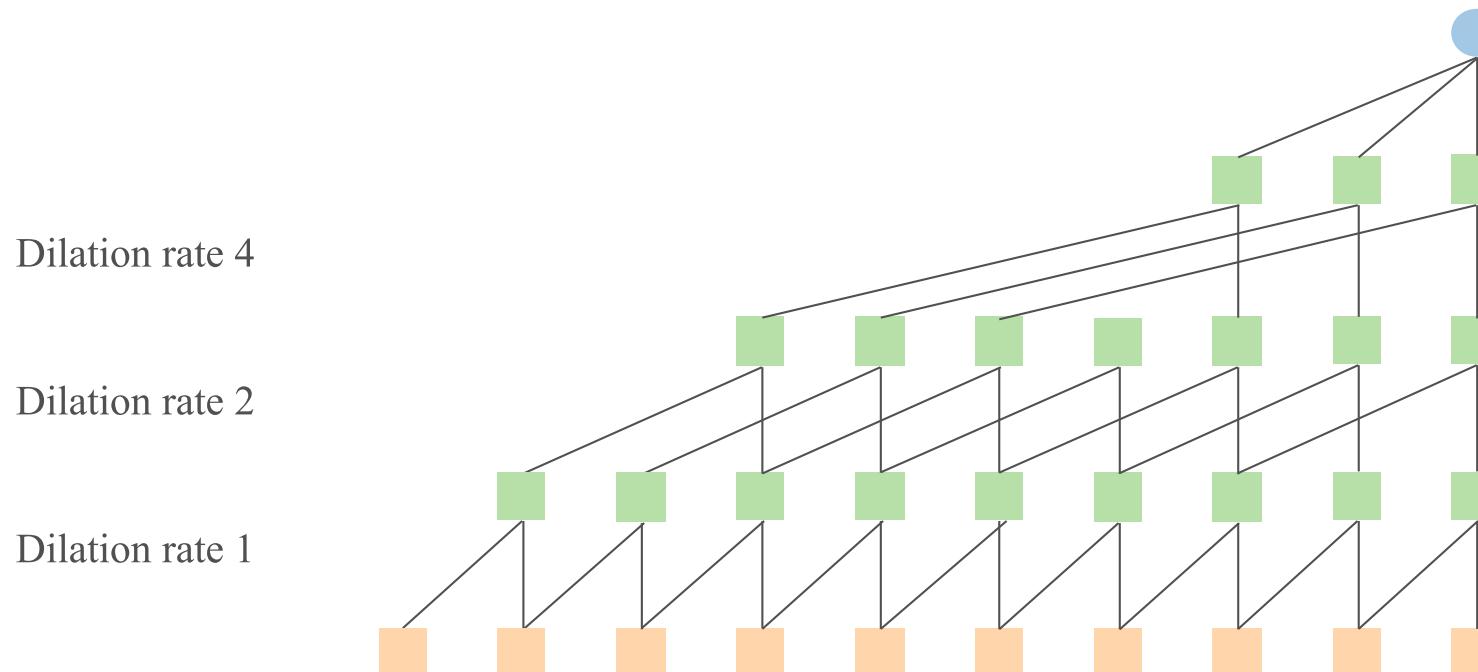
CNNs for forecasting

- Causal convolutions: the output at each time step do not depend on future time steps



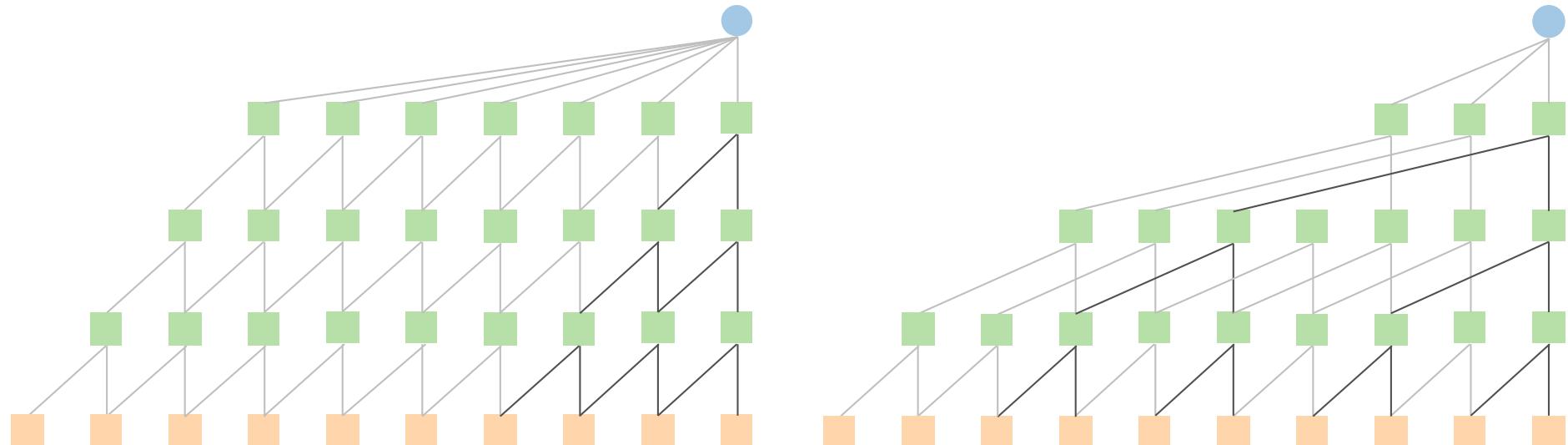
Dilated convolutions

- Skip outputs from previous layers



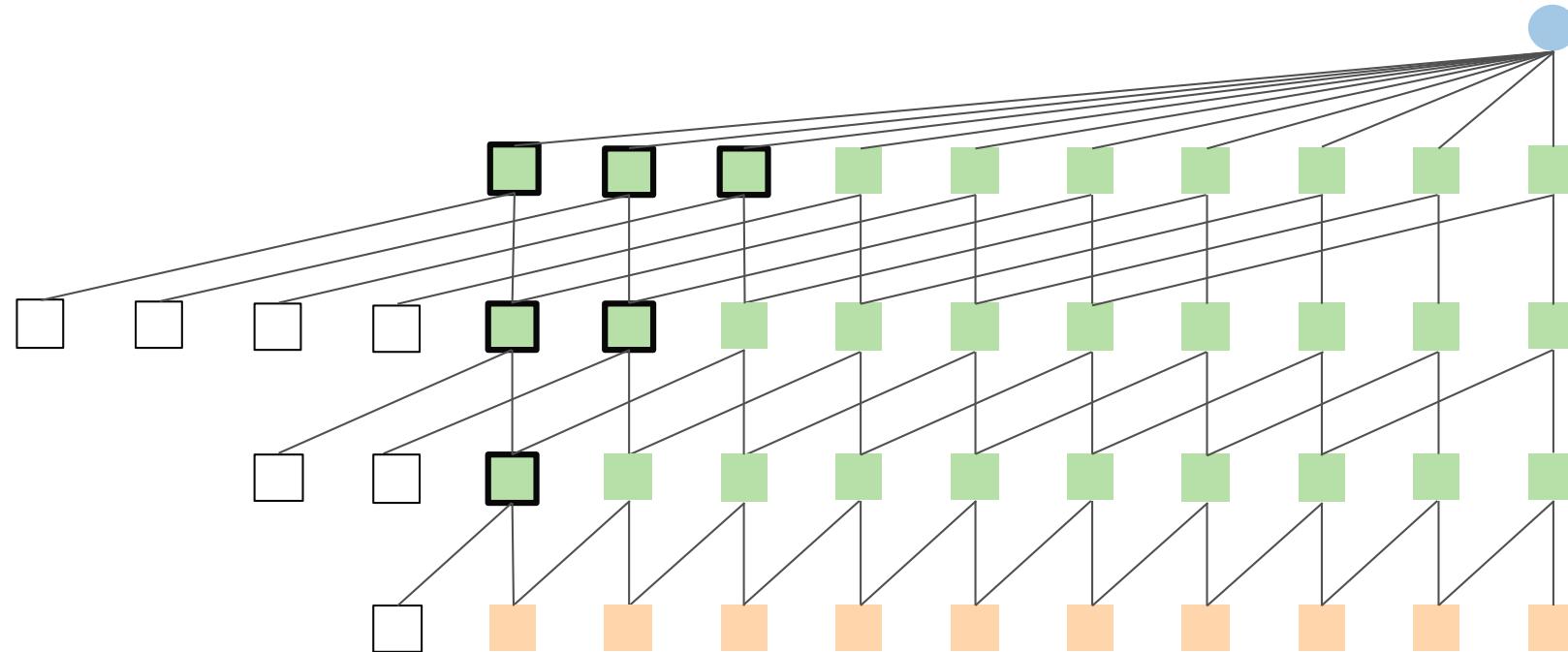
Why dilated convolutions?

- Normal convolutions = more connections = more weights to train
- Reach information from more distant values in the time series



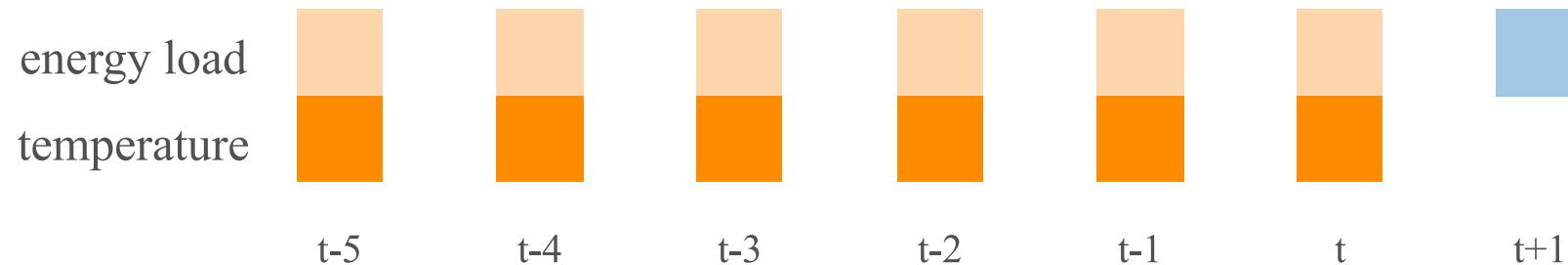
Causal padding

- Padding is necessary to preserve output dimension
- Allows all filter weights to apply to all inputs



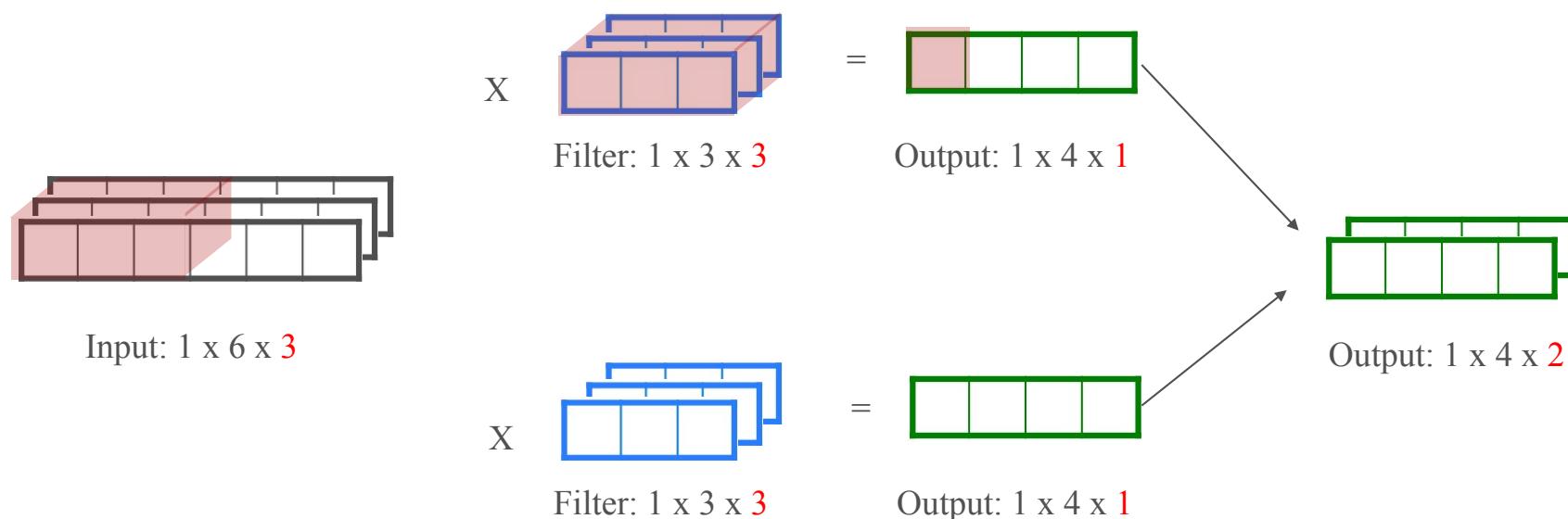
Multivariate time series

- Multiple time series in the input sequence



Multivariate time series

- Example: for 3 input time series use 3 channels

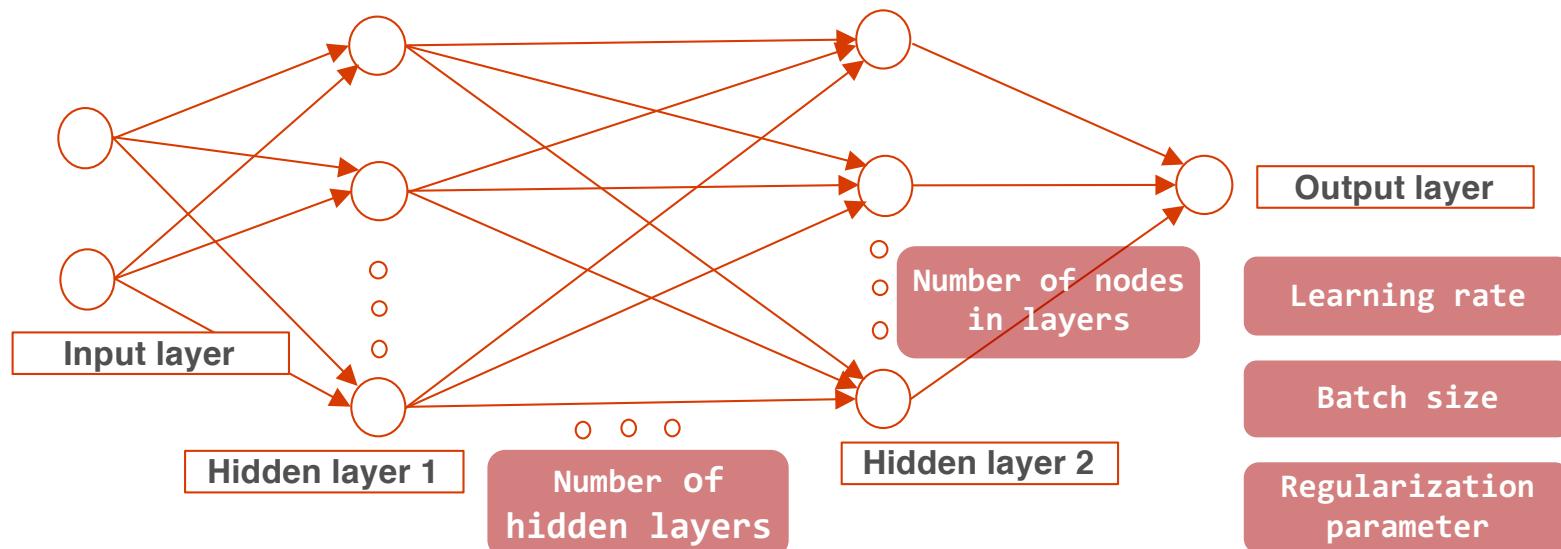


Part III

Hyperparameter tuning

What are hyperparameters?

- Adjustable parameters that govern model training
 - **Architecture**: number of layers, number of cells in each layer, connections
 - **Optimization**: learning rate, mini-batch size, stopping criterion, initialization
 - **Loss function**: weight of regularization term
- Chosen prior to training, stay constant during training, e.g.



Hyper-parameter tuning

- Search across various hyperparameter configurations
- Find the configuration that results in best performance

Challenges

- Huge search space to explore
 - Sparsity of good configurations
 - Expensive evaluation
 - Limited time and resources
-

Summary

- Basics and traditional forecasting
 - Time series features (trend, seasonality, cycle)
 - Autocorrelation
 - Stationarity
 - Transformations and statistical tests
 - Time series decomposition
 - Time series models (SARIMAX)
- Regression-based forecasting
 - Forecasting strategies (recursive, direct, multi-output)
 - Time series cross-validation
- (Deep) neural network architectures for forecasting
 - Convolutional neural networks
 - Recurrent neural networks (RNN, LSTM, GRU)
 - One-step and multi-step ahead forecasting

Acknowledgments

These slides are heavily based on the tutorial “**Deep Learning For Time Series Forecasting**” by Yijing Chen, Angus Taylor, Vanja Paunic, Dmitry Pechyoni at Microsoft. We would like to thank Yijing Chen for sharing the materials with us.