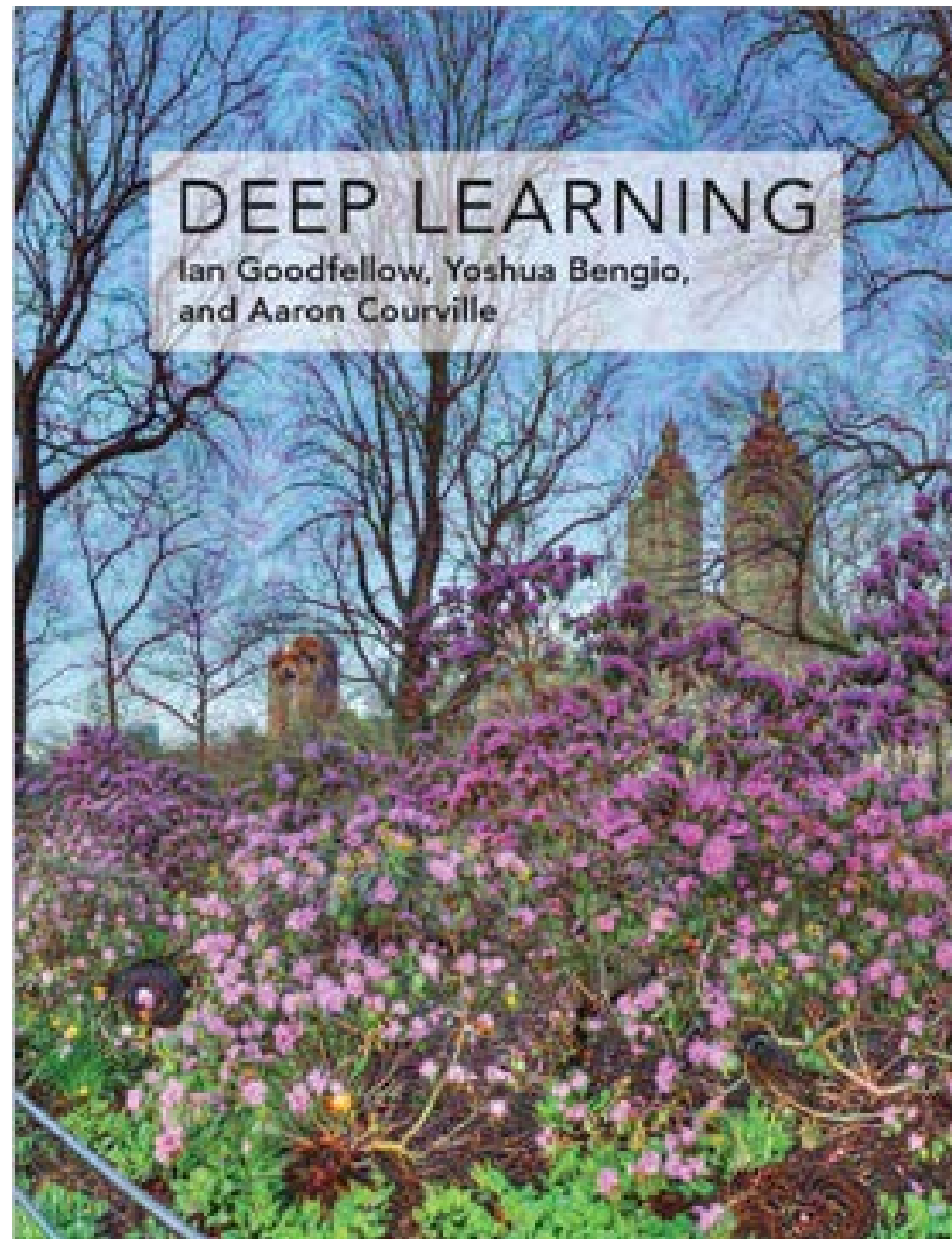**IF3270 Pembelajaran Mesin**

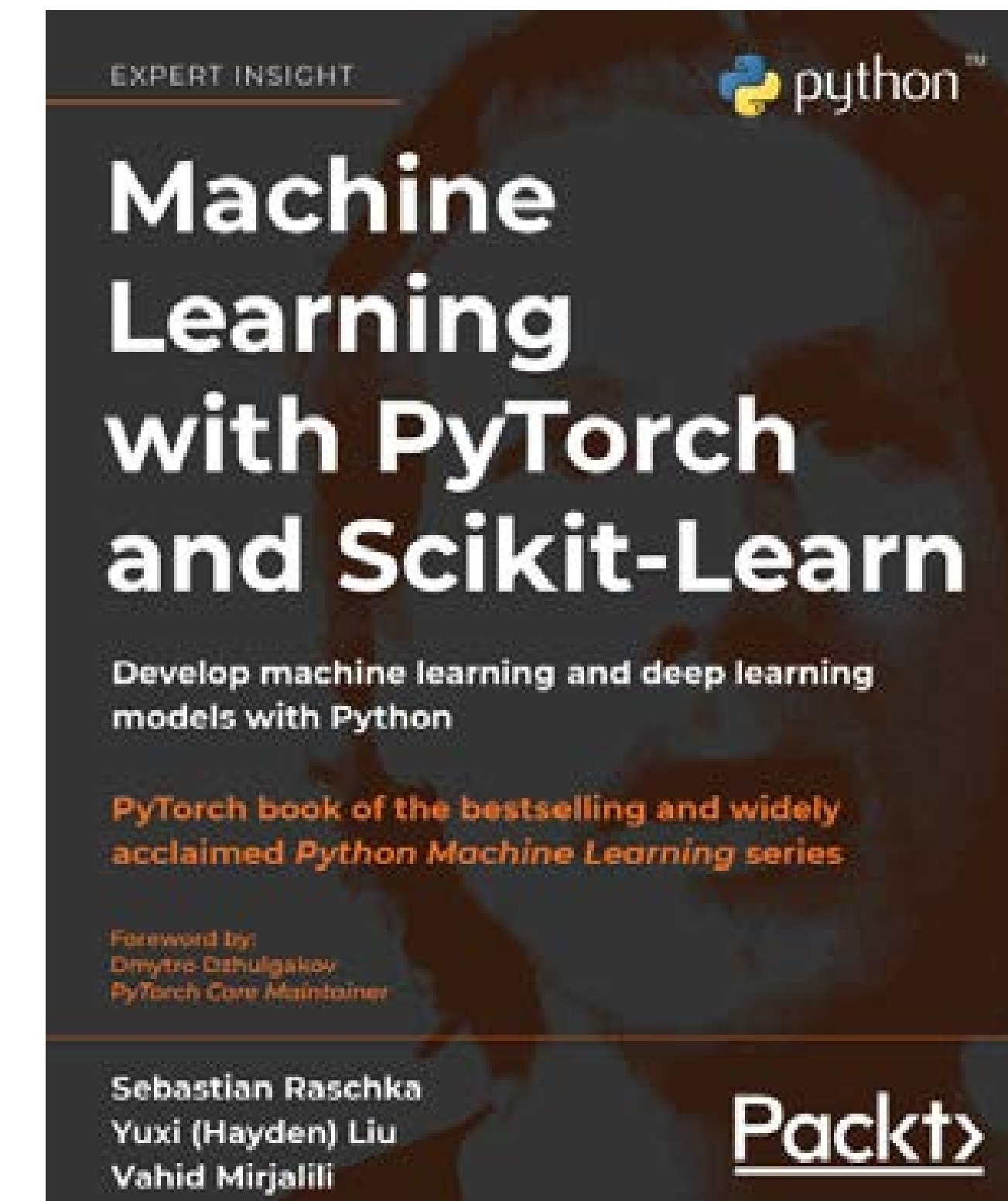# Recurrent Neural Network (RNN)

Tim Pengajar IF3270

# Review

- Machine learning (ML) overview

- Ensemble Methods → Supervised Learning

- Supervised Learning: Perceptron
- Supervised Learning: ANN: Feed Forward Neural Network
- Supervised Learning: ANN: Convolutional Neural Network
- Supervised Learning: ANN: Recurrent Neural Network → Today
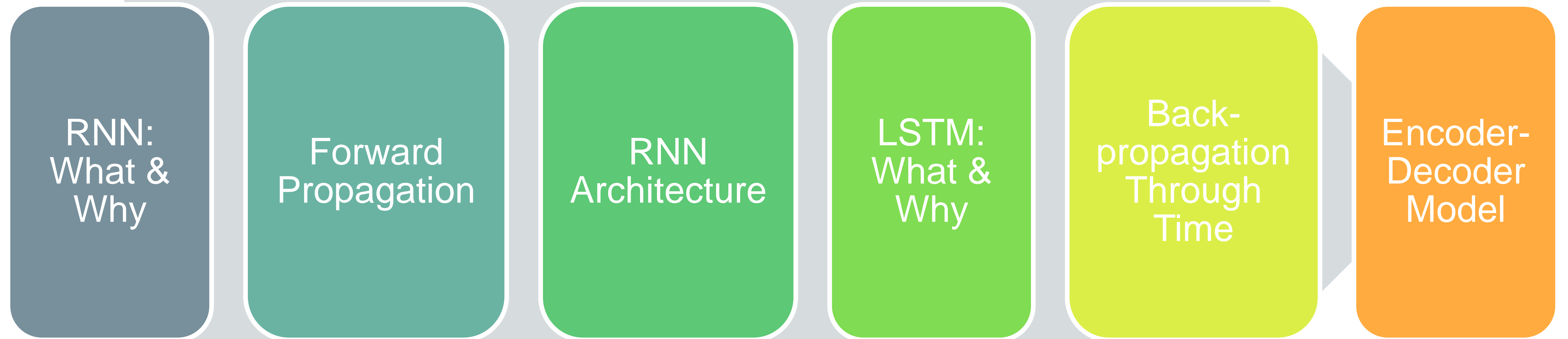
# References



Deep learning. I Goodfellow,
Y Bengio, A Courville, Y
Bengio. MIT press 1 (2), 2016
(Chapter 10)



Raschka, et.al., Machine Learning
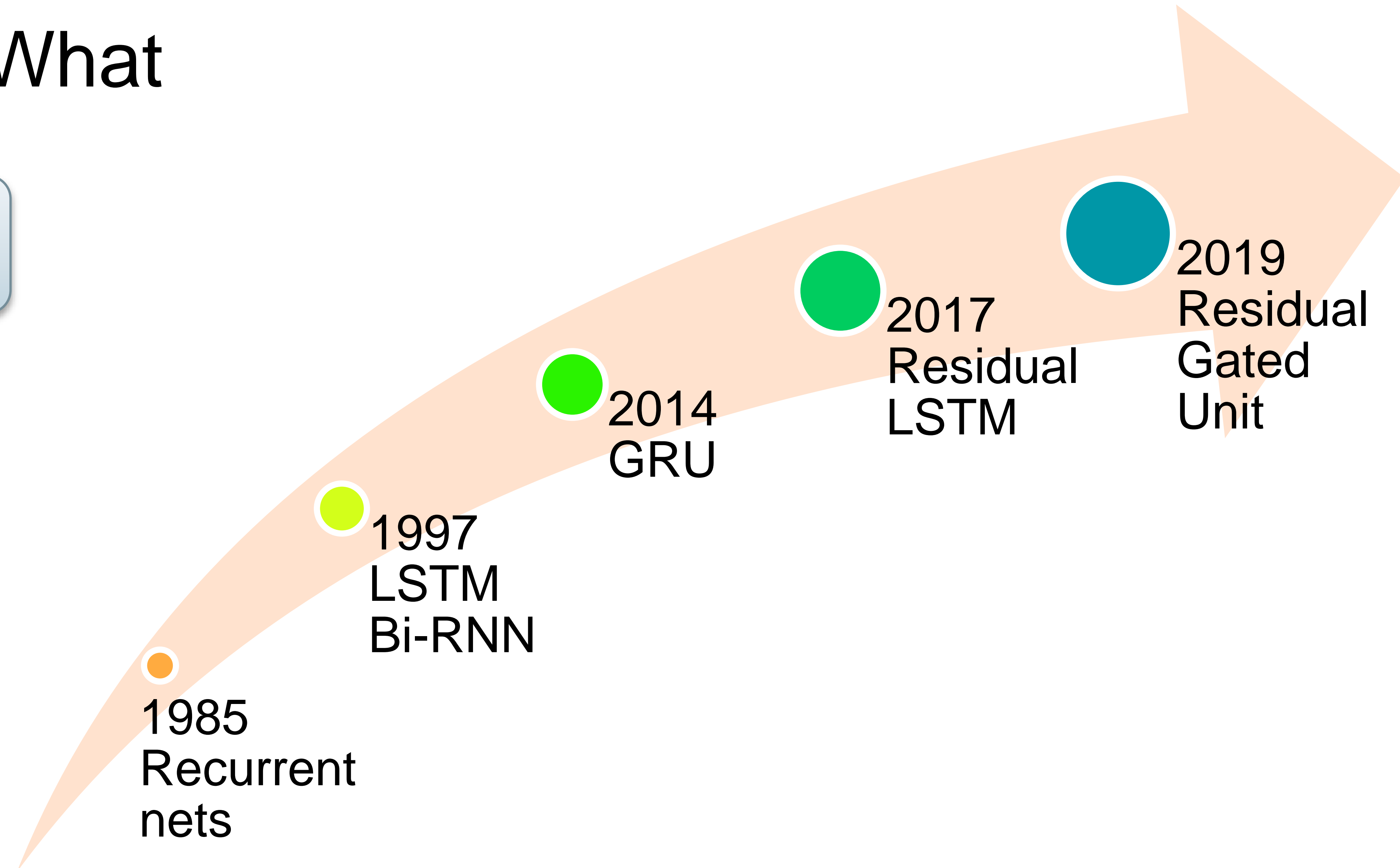with Pytorch and Scikit-Learn,
Packt Publishing Ltd., 2022
(Chapter 15)

# Outline

RNN: What & Why

Forward Propagation

RNN Architecture

LSTM: What & Why

Back-propagation Through Time

Encoder-Decoder Model

# RNN: What & Why

# Recurrent NN: What

ANN with forward
and backward link

2019
Residual
Gated
Unit

2017
Residual
LSTM

2014
GRU

1997
LSTM
Bi-RNN

1985
Recurrent
nets

Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1985). *Learning internal representations by error propagation* (No. ICS-8506). California Univ San Diego La Jolla Inst for Cognitive Science.
Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, *9*(8), 1735-1780.
Schuster, M., & Paliwal, K. K. (1997). Bidirectional recurrent neural networks. *IEEE transactions on Signal Processing*, *45*(11), 2673-2681.
Chung, J., Gulcehre, C., Cho, K., & Bengio, Y. (2014). Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*.
Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. In *Advances in neural information processing systems* (pp. 3104-3112).
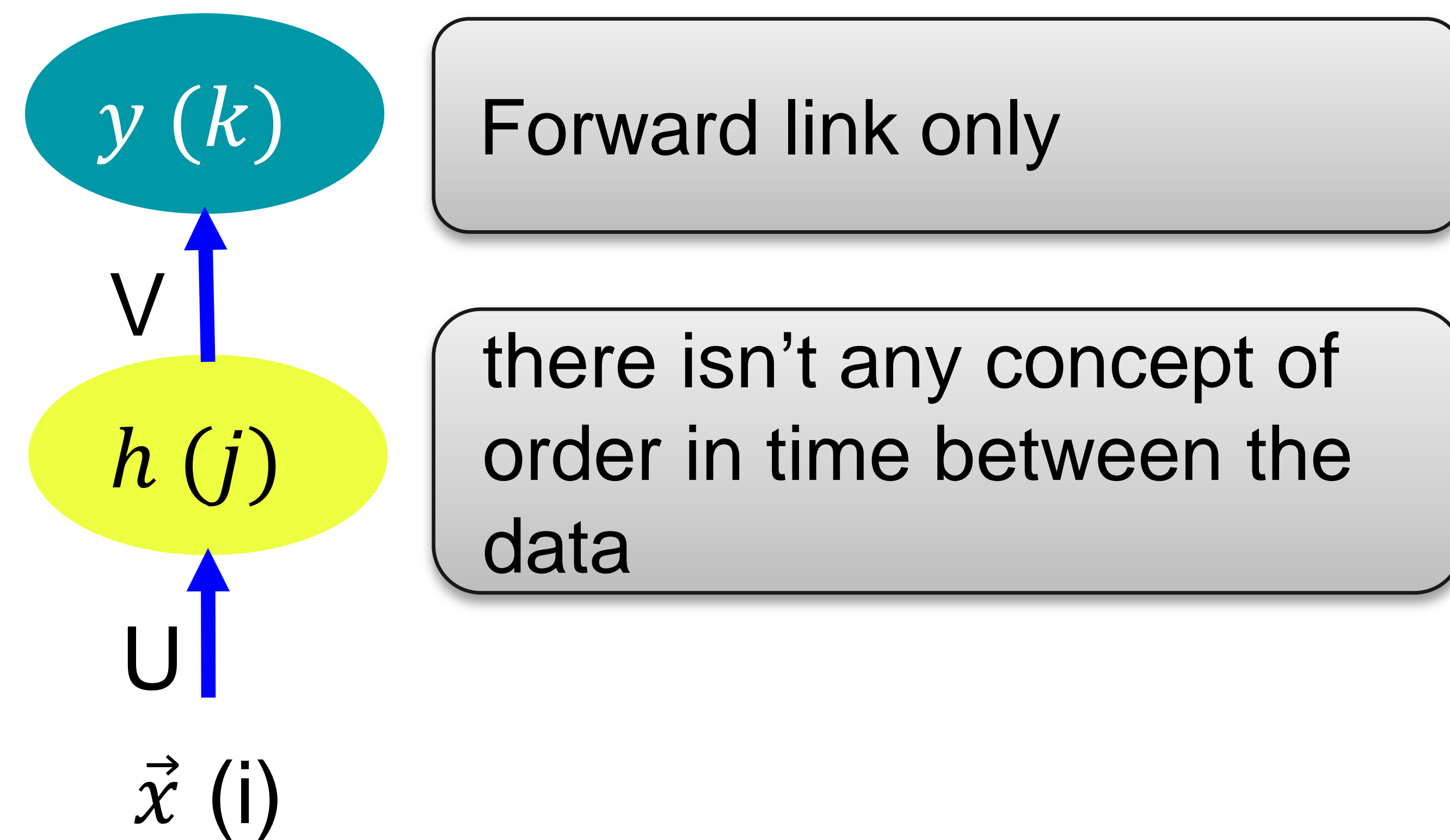Kim, J., El-Khamy, M., & Lee, J. (2017). Residual LSTM: Design of a deep recurrent architecture for distant speech recognition. *arXiv preprint arXiv:1701.03360*.
Luo, H., Li, T., Liu, B., & Zhang, J. (2019). DOER: Dual cross-shared RNN for aspect term-polarity co-extraction. *arXiv preprint arXiv:1906.01794*.
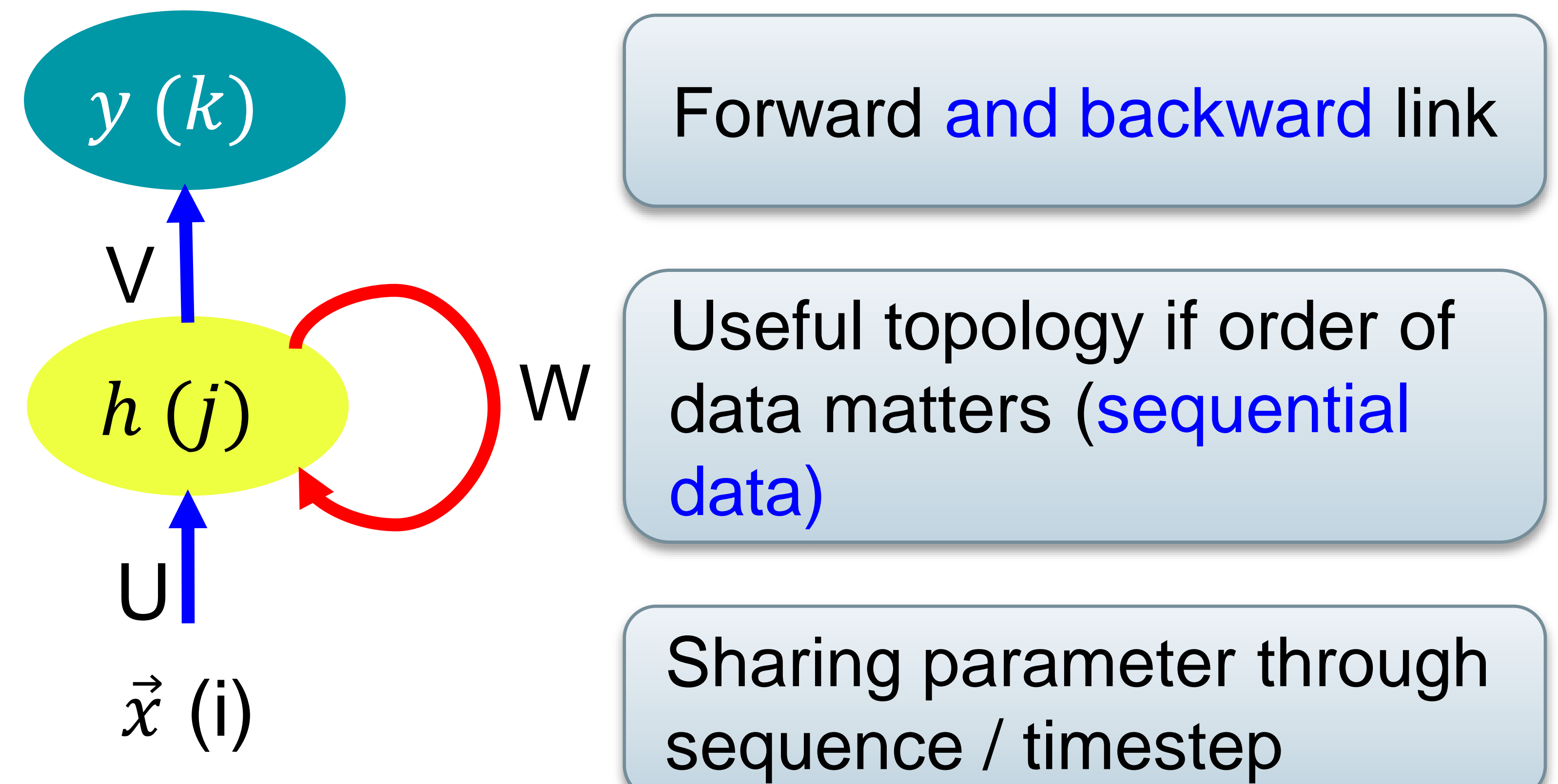
# Feedforward vs Feedback Neural Network



**Feedforward NN**

$y(k)$

Forward link only

V

$h(j)$

there isn't any concept of order in time between the data

U

$\vec{x}$ (i)

**Feedback (Recurrent) NN**

$y(k)$

Forward and backward link

V

$h(j)$   W

Useful topology if order of data matters (sequential data)

U

$\vec{x}$ (i)

Sharing parameter through sequence / timestep
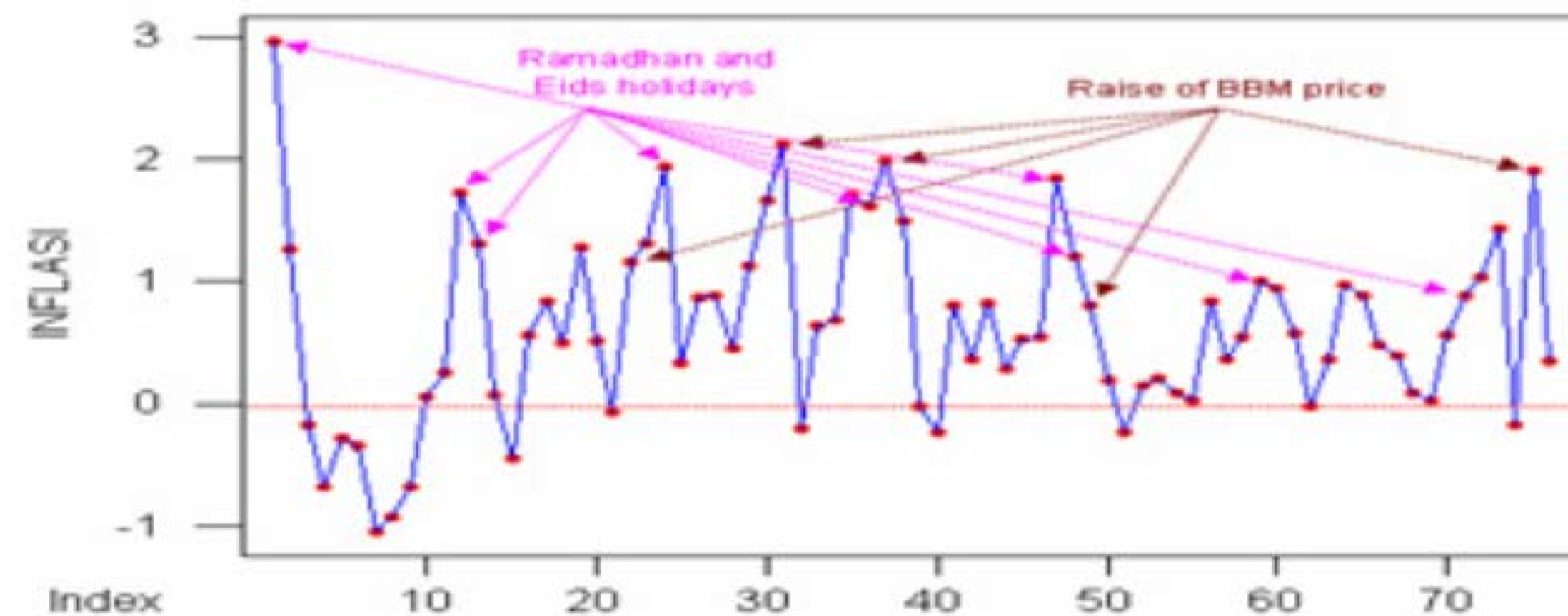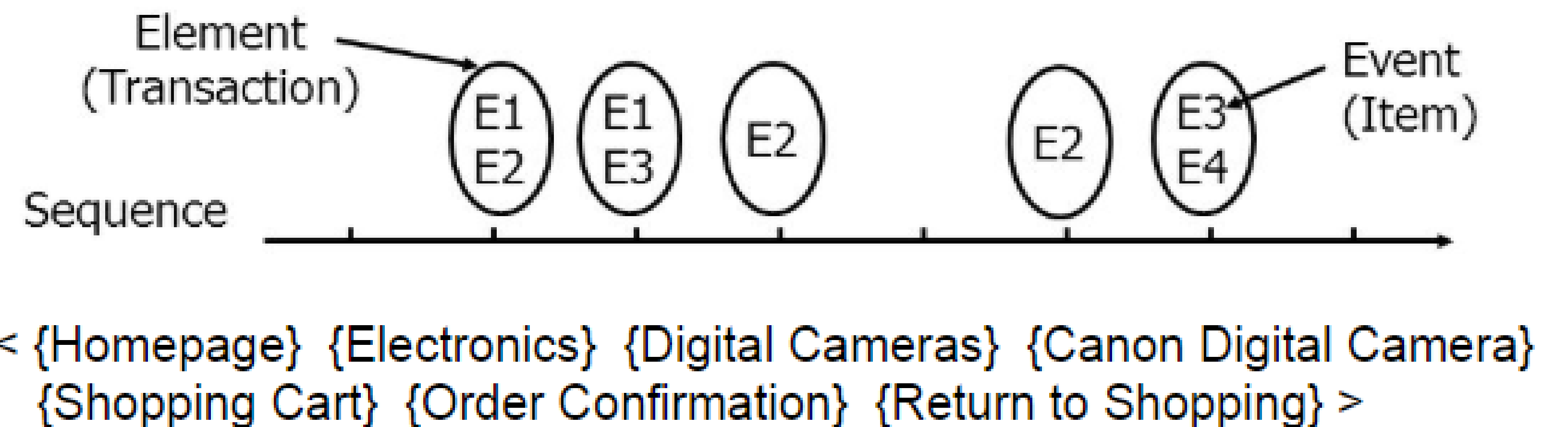
# Why RNN ? Order of data matters



**Time-series data**
consist of long sequences of data, recorded at **equal time intervals**

**Symbolic sequence data**
consist of long sequences of event or nominal data, which typically are not observed at equal time intervals.

**Natural language data**
text, speech

**Other sequence data**
video

< {Homepage} {Electronics} {Digital Cameras} {Canon Digital Camera} {Shopping Cart} {Order Confirmation} {Return to Shopping} >
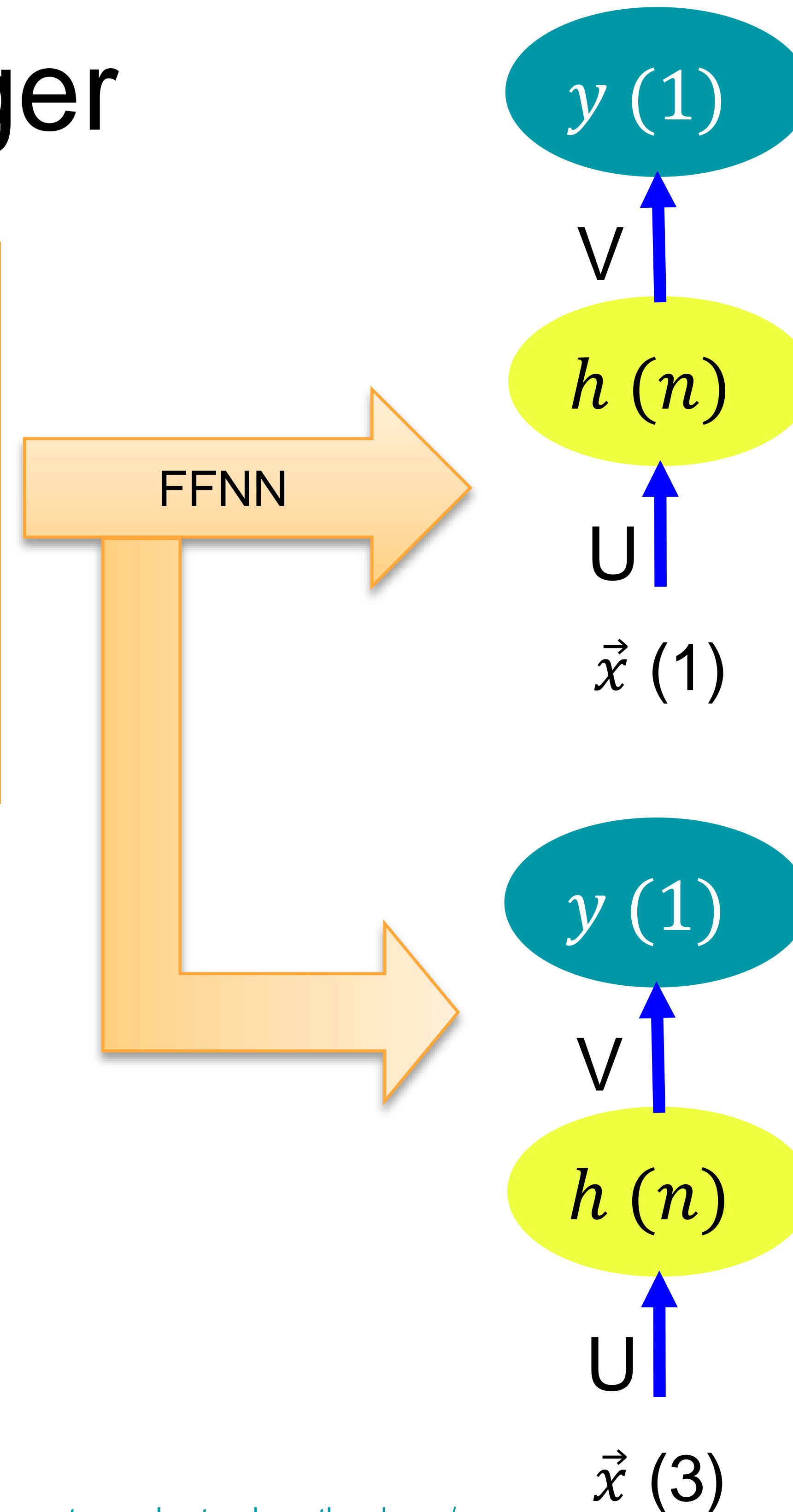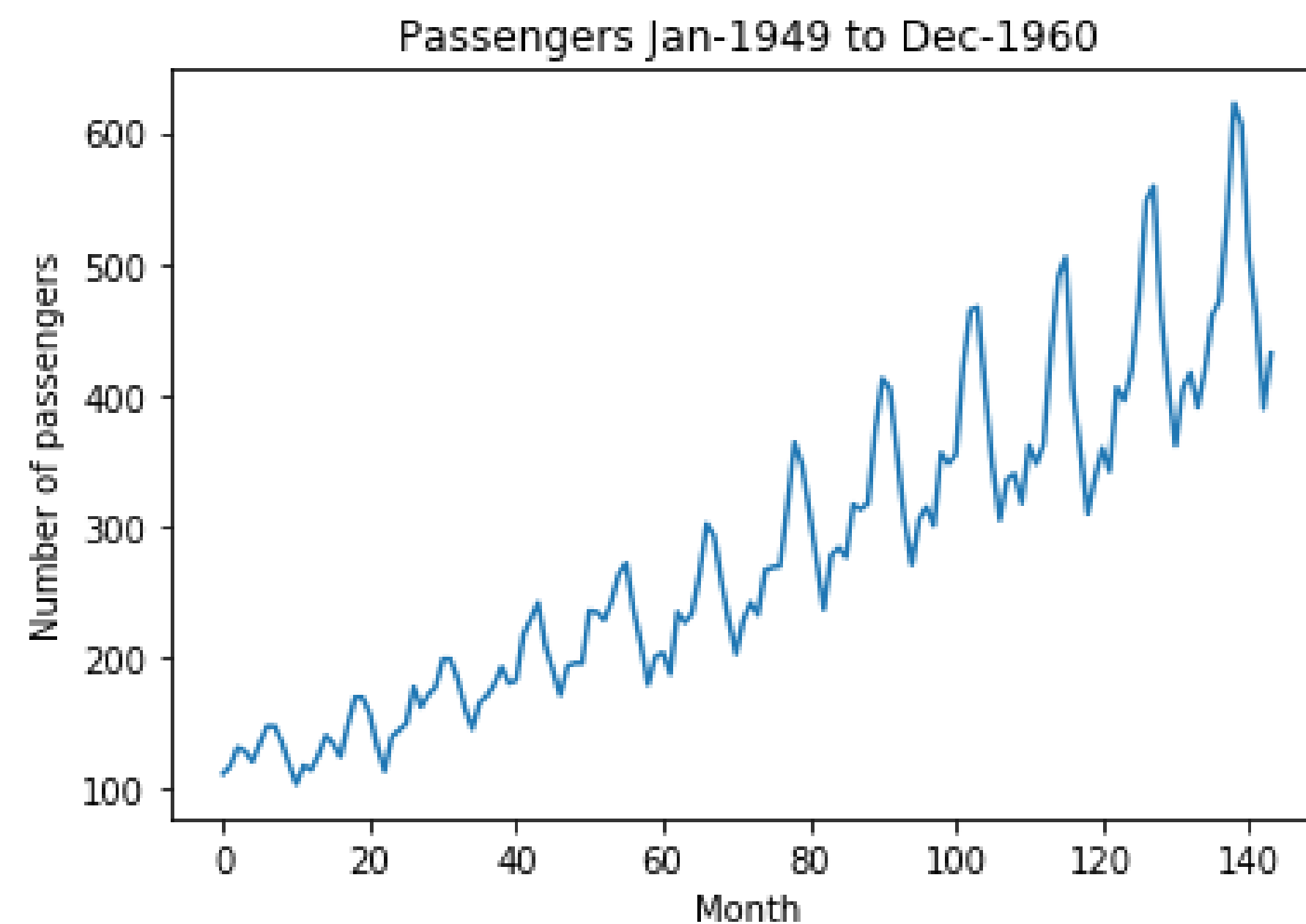
Hidayat, Y., Sutijo, B., Bon, A. T., & Supian, S. (2016). Indonesian financial data modeling and forecasting by using econometrics time series and neural network. *Global Journal of Pure and Applied Mathematics, 12*(4), 3745-3757.
Tan dkk. (2004): https://slideplayer.com/slide/778153/

# Airline Passenger

"Month","Passengers"
"1949-01",112
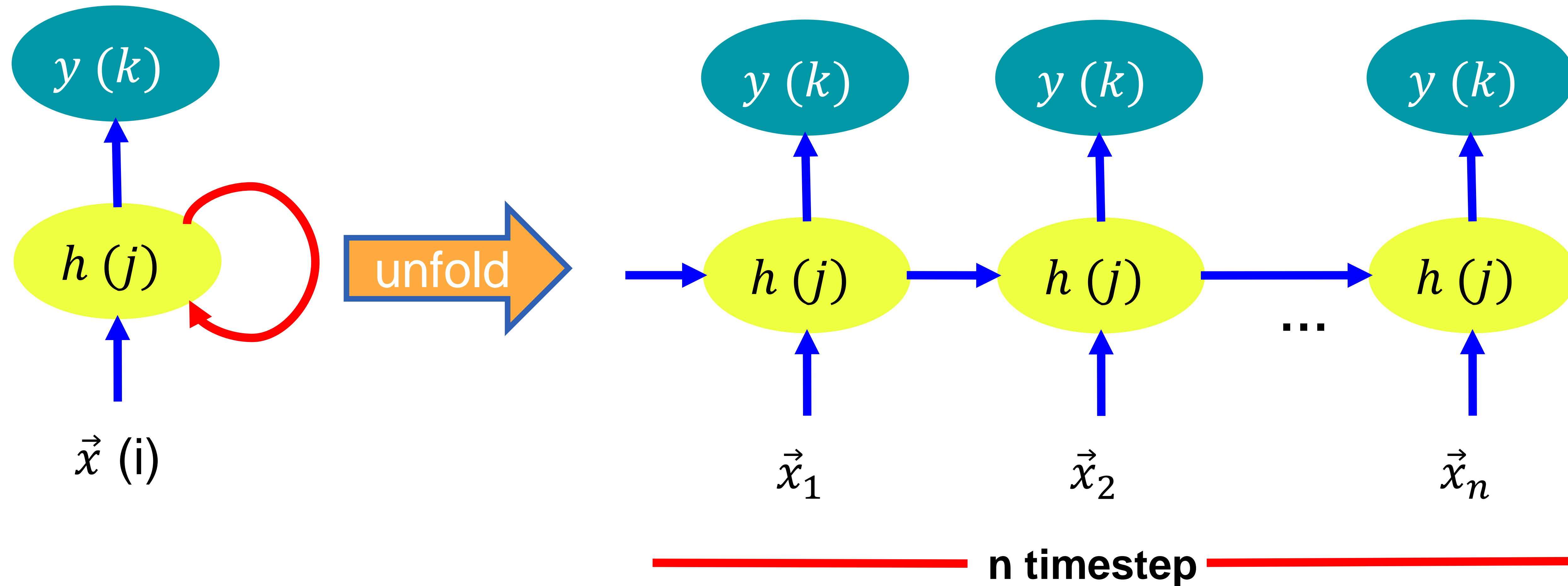"1949-02",118
"1949-03",132
"1949-04",129
"1949-05",121


Passengers Jan-1949 to Dec-1960

FFNN

$y$ (1)

V

$h$ (n)

U

$\vec{x}$ (1)

1 Feature Dataset:

| X=t | Y=(t+1) |
| --- | --- |
| 112 | 118 |
| 118 | 132 |
| 132 | 129 |
| 129 | 121 |
| 121 | 135 |

$y$ (1)

V

$h$ (n)

U

$\vec{x}$ (3)

3 Feature Dataset:

| X1=t-2 | X2=t-1 | X3=t | Y=(t+1) |
| --- | --- | --- | --- |
| 112 | 118 | 132 | 129 |
| 118 | 132 | 129 | 121 |
| 132 | 129 | 121 | 135 |
| 129 | 121 | 135 | 148 |
| 121 | 135 | 148 | 148 |

https://machinelearningmastery.com/time-series-prediction-lstm-recurrent-neural-networks-python-keras/

# RNN: Neuron Dependent



A RNN unit has loops in them that allow information to be carried across neurons while reading in input. RNN cannot rely on the input alone and must use its recurrent connection to keep track of the context to achieve this task.

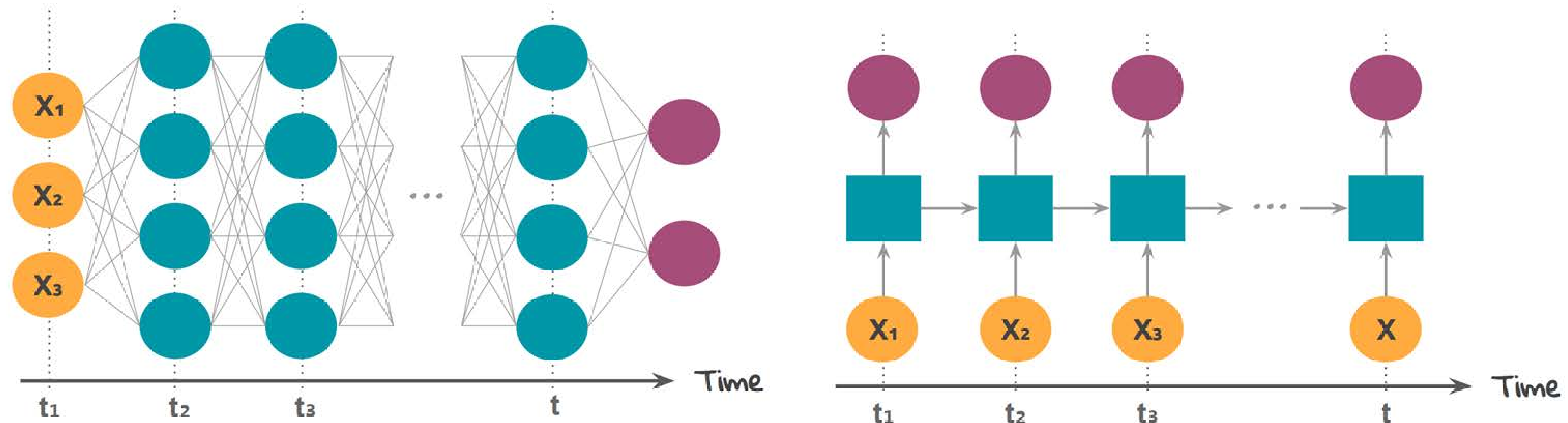http://colah.github.io/posts/2015-08-Understanding-LSTMs/

# RNN: Parameter Sharing



$$h_t = f(Ux_t + (Wh_{t-1} + b_{xh}))$$
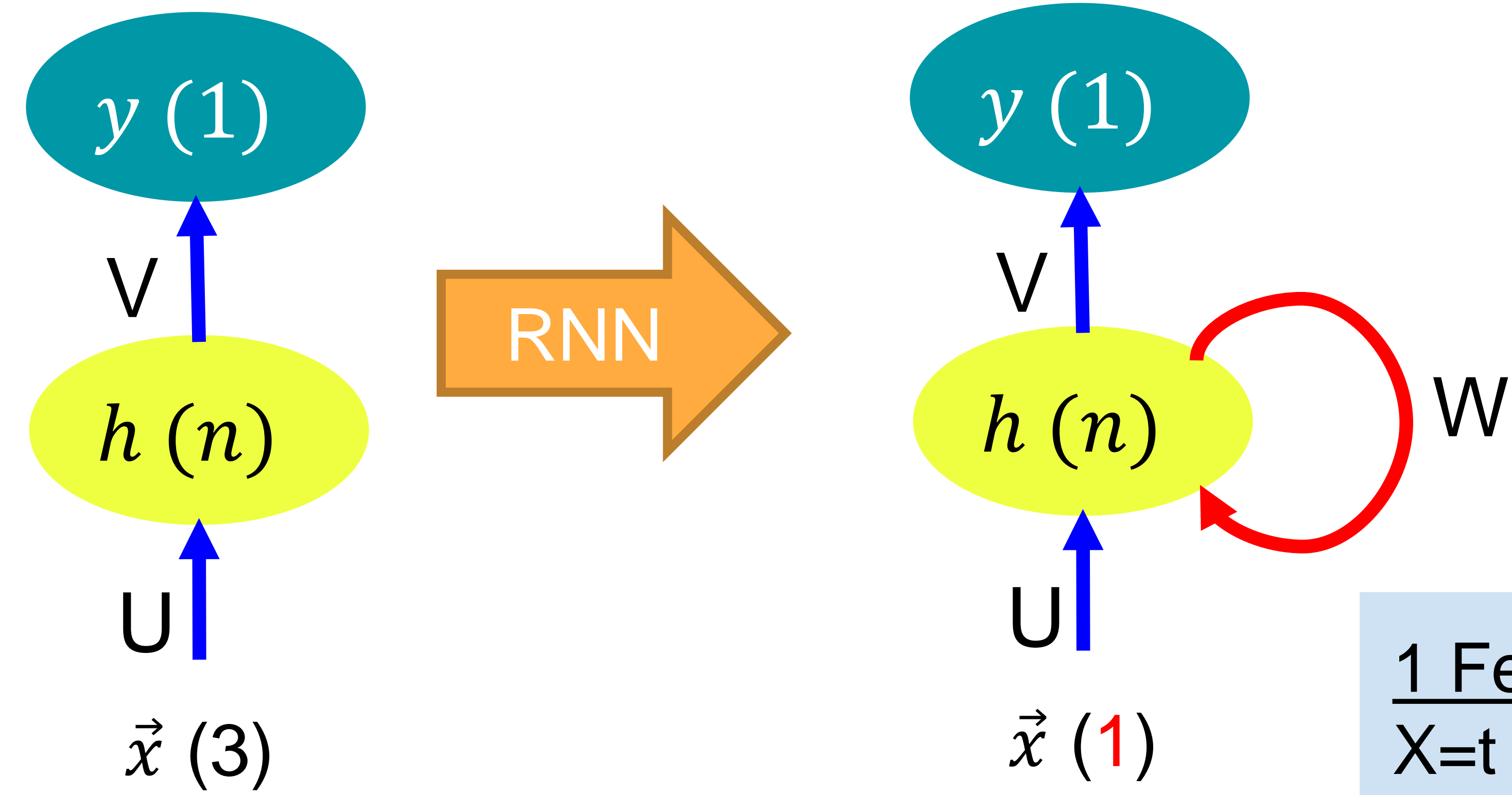$$y_t = f(Vh_t + b_{hy})$$

- $x_t$ : feature vector (i features) at step t
- $h_t$ : hidden state.
- $y_t$ : output at step t.
- f: activation function

11

# FFNN vs RNN: Sequential Data



- FFNN: there <span style="color:red">isn't</span> any concept of order in time between the data
- RNN: there is order in time between the data. We will input **X1** first and then input **X2** to the result of **X1** computation. So in the same way, **X3** is computed with the result from **X2** computation stage.

https://towardsdatascience.com/the-most-intuitive-and-easiest-guide-for-recurrent-neural-network-873c29da73c7

# FFNN vs RNN

$y\ (1)$

V

$h\ (n)$

U

$\vec{x}\ (3)$

RNN

$y\ (1)$

V

$h\ (n)$   W

U

$\vec{x}\ (\textcolor{red}{1})$

3 timestep

3 Feature Dataset:

| X1 | X2 | X3 | Y=(t+1) |
|---|---|---|---|
| 112 | 118 | 132 | 129 |
| 118 | 132 | 129 | 121 |
| 132 | 129 | 121 | 135 |
| 129 | 121 | 135 | 148 |
| 121 | 135 | 148 | 148 |

1 Feature Dataset:

| X=t | Y=(t+1) |
|---|---|
| 112 | 118 |
| 118 | 132 |
| 132 | 129 |
| 129 | 121 |
| 121 | 135 |

# Summary

RNN: Feedback NN

Why RNN ?
Order of data matters
(Sequence data)

Neuron dependent

Parameter sharing

Forward Propagation

# Forward Propagation
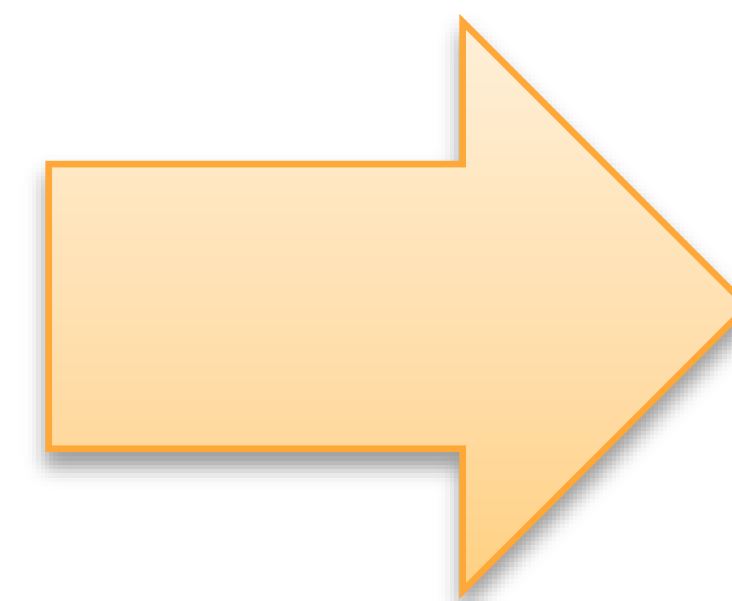
# Classification for Sequence Data

Sequence: ABCCD…

Browsing history: {Homepage}{Electronics}{Camera}{Camera}{ShoppingCart}…
Language model of character: hello…
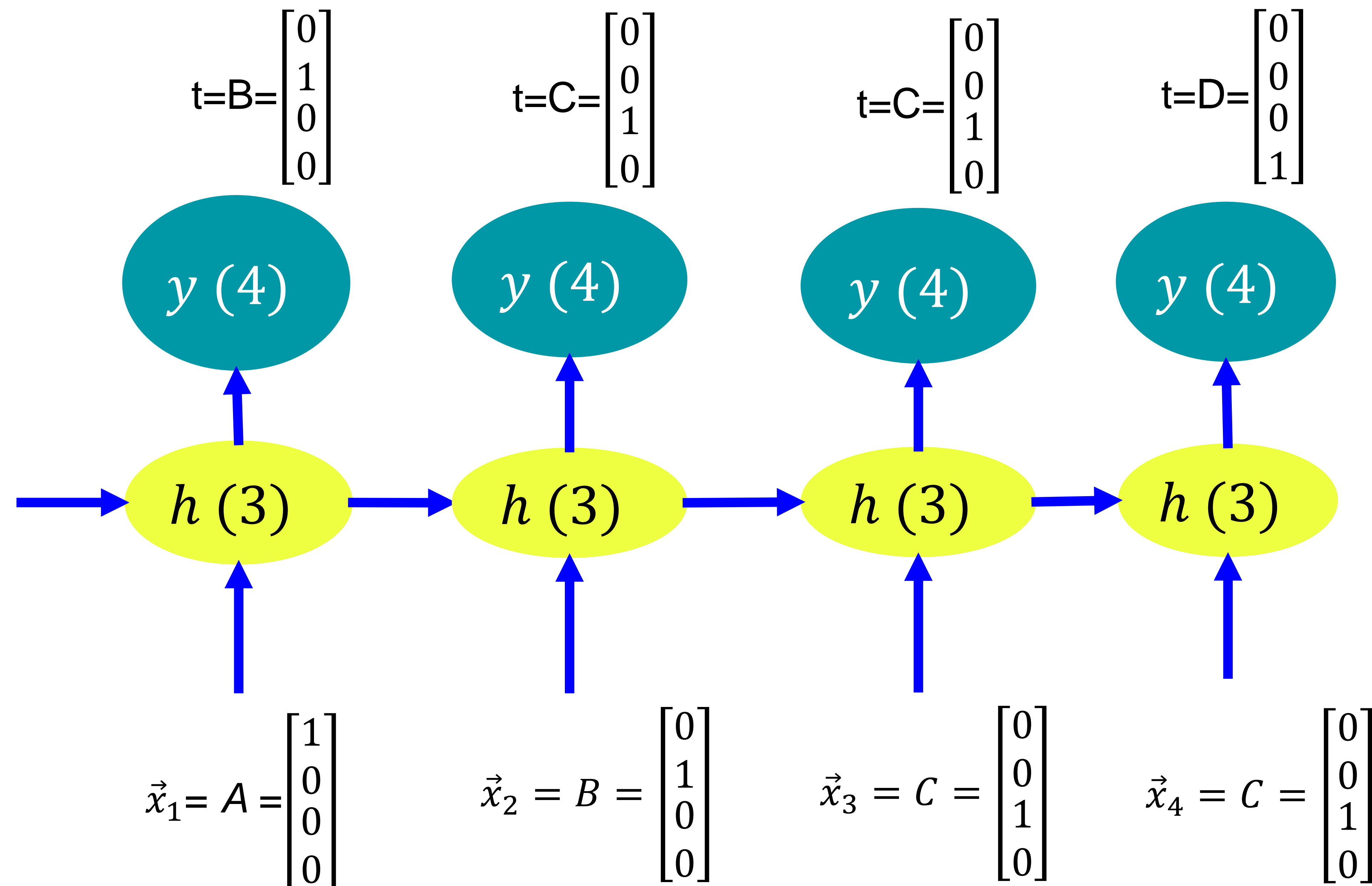Language model of words: aku sedang pura pura tertawa …

Vector representation:

$$A = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, B = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, C = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, D = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

$$ABCCD \ldots = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \ldots$$

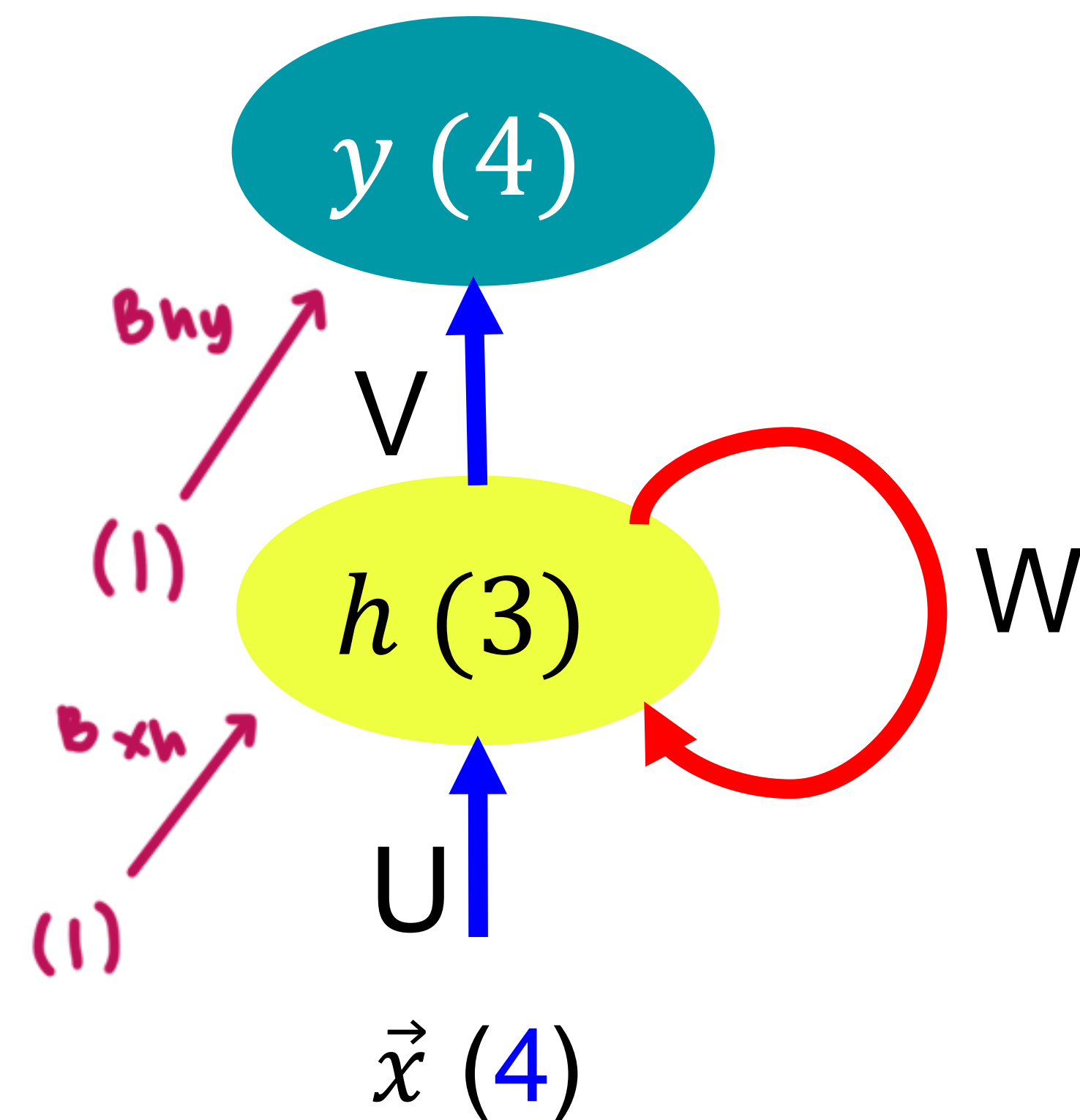# Dataset Construction: 4 timestep



Sequence: ABCCD…

| A1 | A2 | A3 | A4 | Class |
|----|----|----|----|-------|
| 1 | 0 | 0 | 0 | B |
| 0 | 1 | 0 | 0 | C |
| 0 | 0 | 1 | 0 | C |
| 0 | 0 | 1 | 0 | D |
| … | | | | |

# Sequence Classification: RNN



- U: matrix 3x4 (hidden neurons x input dimension)
- V: matrix 4x3 (output neurons x hidden neurons)
- W: matrix 3x3 (hidden neurons x hidden neurons)
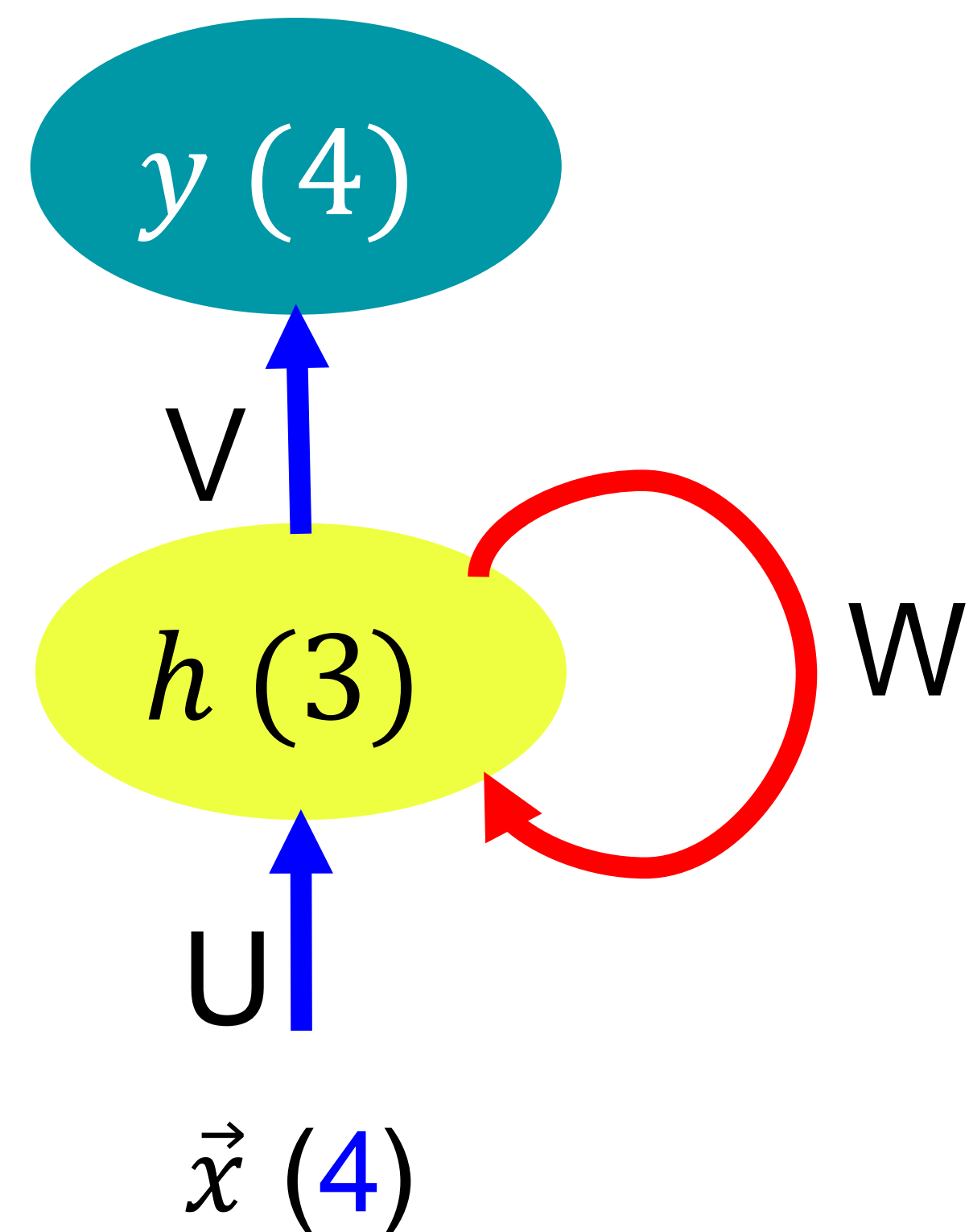- Bias$_{xh}$: matrix 3x1
- Bias$_{hy}$: matrix 4x1

$$h_t = tanh(Ux_t + (Wh_{t-1} + b_{xh}))$$
$$y_t = softmax(Vh_t + b_{hy})$$

$$h_t = \tanh(U \cdot x_t + (W \cdot h_{t-1} + b_{xh}))$$
$$y_t = \text{softmax}(V \cdot h_t + b_{hy})$$

# Weight Initialization: Example (Random)



| | U | | | | W | | | b_xh | ht-1 |
|---|---|---|---|---|---|---|---|---|---|
| 0.100 | 0.150 | 0.200 | 0.300 | 0.500 | 0.500 | 0.500 | | 0.100 | 0 |
| 0.150 | 0.200 | 0.300 | 0.100 | 0.500 | 0.500 | 0.500 | | 0.100 | 0 |
| 0.200 | 0.300 | 0.100 | 0.150 | 0.500 | 0.500 | 0.500 | | 0.100 | 0 |

| | V | | b_hy |
|---|---|---|---|
| 0.100 | 0.200 | 0.300 | 0.100 |
| 0.200 | 0.300 | 0.100 | 0.100 |
| 0.300 | 0.100 | 0.200 | 0.100 |
| 0.100 | 0.100 | 0.100 | 0.100 |

$$h_1 = \begin{bmatrix} 0.1 & 0.15 & 0.2 & 0.3 \\ 0.15 & 0.2 & 0.3 & 0.1 \\ 0.2 & 0.3 & 0.1 & 0.15 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$+ \left( \begin{bmatrix} 0.5 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0.5 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \end{bmatrix} \right)$$

$$= \begin{bmatrix} 0.1 \\ 0.15 \\ 0.2 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \end{bmatrix} = \begin{bmatrix} 0.2 \\ 0.25 \\ 0.3 \end{bmatrix}$$

$$\tanh\left( \begin{bmatrix} 0.2 \\ 0.25 \\ 0.3 \end{bmatrix} \right) = \begin{bmatrix} 0.197 \\ 0.245 \\ 0.291 \end{bmatrix}$$

19

# Forward Propagation: 1 instance, 4 timestep



$$y_1 = \begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.2 & 0.3 & 0.1 \\ 0.3 & 0.1 & 0.2 \\ 0.1 & 0.1 & 0.1 \end{bmatrix} \begin{bmatrix} 0.197 \\ 0.245 \\ 0.291 \end{bmatrix}$$

$$+ \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \end{bmatrix} = \begin{bmatrix} 0.256 \\ 0.242 \\ 0.2418 \\ 0.1733 \end{bmatrix}$$

$$y_1 = softmax\left(\begin{bmatrix} 0.256 \\ 0.242 \\ 0.2418 \\ 0.1733 \end{bmatrix}\right)$$

$$= \begin{bmatrix} 1.29175/5.0283 \\ 1.27379/5.0283 \\ 1.27354/5.0283 \\ 1.18922/5.0283 \end{bmatrix} = \begin{bmatrix} 0.257 \\ 0.253 \\ 0.253 \\ 0.236 \end{bmatrix}$$

$$h_2 = \begin{bmatrix} 0.1 & 0.15 & 0.2 & 0.3 \\ 0.15 & 0.2 & 0.3 & 0.1 \\ 0.2 & 0.3 & 0.1 & 0.15 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$+ \left(\begin{bmatrix} 0.5 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0.5 \end{bmatrix} \begin{bmatrix} 0.197 \\ 0.245 \\ 0.291 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \end{bmatrix}\right)$$

$$= \begin{bmatrix} 0.15 \\ 0.2 \\ 0.3 \end{bmatrix} + \begin{bmatrix} 0.3665 \\ 0.3665 \\ 0.3665 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \end{bmatrix} = \begin{bmatrix} 0.6165 \\ 0.6665 \\ 0.7665 \end{bmatrix}$$

$$tanh\left(\begin{bmatrix} 0.6165 \\ 0.6665 \\ 0.7665 \end{bmatrix}\right) = \begin{bmatrix} 0.549 \\ 0.583 \\ 0.645 \end{bmatrix}$$

20

# Computing $h_t$ and $y_t$ : Timestep t1 and t2

$$h_t = tanh(Ux_t + (Wh_{t-1} + b_{xh}))$$
$$y_t = softmax(Vh_t + b_{hy})$$

$t1=<\begin{bmatrix}1\\0\\0\\0\\0\end{bmatrix}, \begin{bmatrix}0\\1\\0\\0\\0\end{bmatrix}>$

| Uxt | Wht-1+bxh | net_ht | ht |
|---|---|---|---|
| 0.100 | 0.100 | 0.200 | 0.197 |
| 0.150 | 0.100 | 0.250 | 0.245 |
| 0.200 | 0.100 | 0.300 | 0.291 |

| Vht+bhy | exp(Vht+bhy) | yt |
|---|---|---|
| 0.256 | 1.292 | 0.257 |
| 0.242 | 1.274 | 0.253 |
| 0.242 | 1.274 | 0.253 |
| 0.173 | 1.189 | 0.236 |

$t2=<\begin{bmatrix}0\\1\\0\\0\\0\end{bmatrix}, \begin{bmatrix}0\\0\\1\\0\\0\end{bmatrix}>$

| Uxt | Wht-1+bxh | net_ht | ht |
|---|---|---|---|
| 0.150 | 0.467 | 0.617 | 0.549 |
| 0.200 | 0.467 | 0.667 | 0.583 |
| 0.300 | 0.467 | 0.767 | 0.645 |

| Vht+bhy | exp(Vht+bhy) | yt |
|---|---|---|
| 0.465 | 1.592 | 0.263 |
| 0.449 | 1.567 | 0.259 |
| 0.452 | 1.571 | 0.260 |
| 0.278 | 1.320 | 0.218 |

# Computing $h_t$ and $y_t$ : Timestep t3 and t4

$$h_t = tanh(Ux_t + (Wh_{t-1} + b_{xh}))$$
$$y_t = softmax(Vh_t + b_{hy})$$

$t3 = \left\langle \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \right\rangle$

| Uxt | Wht-1+bxh | net_ht | ht |
|---|---|---|---|
| 0.200 | 0.988 | 1.188 | 0.830 |
| 0.300 | 0.988 | 1.288 | 0.859 |
| 0.100 | 0.988 | 1.088 | 0.796 |

| Vht+bhy | exp(Vht+bhy) | yt |
|---|---|---|
| 0.594 | 1.811 | 0.299 |
| 0.603 | 1.828 | 0.302 |
| 0.594 | 1.812 | 0.299 |
| 0.349 | 1.417 | 0.234 |

$t4 = \left\langle \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \right\rangle$

| Uxt | Wht-1+bxh | net_ht | ht |
|---|---|---|---|
| 0.200 | 1.343 | 1.543 | 0.913 |
| 0.300 | 1.343 | 1.643 | 0.928 |
| 0.100 | 1.343 | 1.443 | 0.894 |

| Vht+bhy | exp(Vht+bhy) | yt |
|---|---|---|
| 0.645 | 1.906 | 0.315 |
| 0.650 | 1.916 | 0.317 |
| 0.645 | 1.907 | 0.315 |
| 0.373 | 1.453 | 0.240 |

# Forward Propagation: 1 instance, 4 timestep



$$\begin{bmatrix} 0.257 \\ 0.253 \\ 0.253 \\ 0.236 \end{bmatrix} \quad t=\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0.263 \\ 0.259 \\ 0.260 \\ 0.218 \end{bmatrix} \quad t=\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0.299 \\ 0.302 \\ 0.299 \\ 0.234 \end{bmatrix} \quad t=\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0.315 \\ 0.317 \\ 0.315 \\ 0.240 \end{bmatrix} \quad t=\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} 0.197 \\ 0.245 \\ 0.291 \end{bmatrix}$$

$$\begin{bmatrix} 0.549 \\ 0.583 \\ 0.645 \end{bmatrix}$$

$$\begin{bmatrix} 0.830 \\ 0.859 \\ 0.796 \end{bmatrix}$$

$$\begin{bmatrix} 0.913 \\ 0.928 \\ 0.894 \end{bmatrix}$$

$$\vec{x}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\vec{x}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\vec{x}_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\vec{x}_4 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

# Language Model of Character



- t=1 (input "h"): output="o", target="e"
- t=2 (input "e"): output="o", target="l"
- t=3 (input "l"): output="l", target="l"
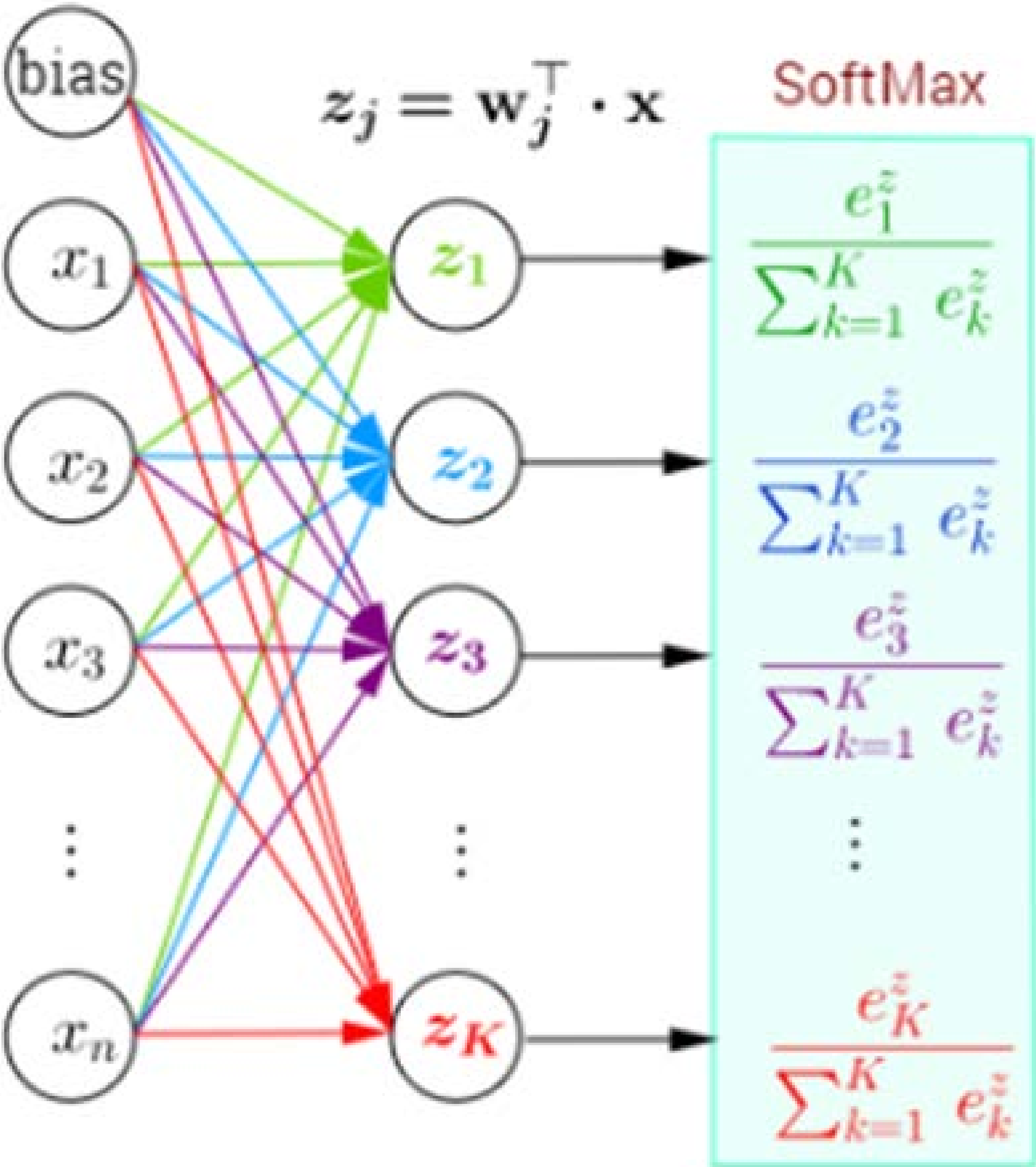- t=4 (input "l"): output="o", target="o"

http://karpathy.github.io/2015/05/21/rnn-effectiveness/

# Output with Softmax

"e"

| 1.0 |
| 2.2 |
| -3.0 |
| 4.1 |

Target: 'e'
Output: 'o'

| 0.3 |
| -0.1 |
| 0.9 |

| 1 |
| 0 |
| 0 |
| 0 |

"h"

| z | exp(z) | P(z) |
|---|--------|------|
| 1 | 2,718 | 0 |
| 2,2 | 9,025 | 0,1 |
| -3 | 0,05 | 0 |
| 4,1 | 60,34 | 0,8 |
|  | 72,13 | 1 |

bias

$z_j = \mathbf{w}_j^\top \cdot \mathbf{x}$

SoftMax

$x_1$    $z_1$

$x_2$    $z_2$

$x_3$    $z_3$

$x_n$    $z_K$

$$\frac{e_1^z}{\sum_{k=1}^{K} e_k^z}$$

$$\frac{e_2^z}{\sum_{k=1}^{K} e_k^z}$$

$$\frac{e_3^z}{\sum_{k=1}^{K} e_k^z}$$

$$\frac{e_K^z}{\sum_{k=1}^{K} e_k^z}$$

# Summary

| Sequence Classification | Dataset Construction | Computing $h_t$ and $y_t$ |
|---|---|---|

**RNN Architecture** →

# RNN Architecture

# General Architecture



$y\ (m)$

V

$hh$ (k)

$W_{hh}$

$U_{h...hh}$

...

$W_{h...}$

$U_{h1h...}$

$h1\ (j)$

$W_{h1}$

$U_{xh1}$

$\vec{x}\ (i)$

Return sequence = True/False

$h\ (j)$

$h\ (j)$

...

$h\ (j)$

$\vec{x}_1$

$\vec{x}_2$

$\vec{x}_n$

**n timestep**

# Architecture



one to many — 1 feature vector input

many to one — kelas output

many to many

many to many — harus baca semua sequence input baru sequence output digenerate

| fixed-sized input vector xt | RNN state st | fixed-sized output vector ot |

One to many: image captioning
Many to one: text classification
Many to many: machine translation, video frame classification, POS tagging

http://karpathy.github.io/2015/05/21/rnn-effectiveness/

29

# One to Many: Image Captioning

image → text
sequence of token

feature map (vektor)

image → CNN → RNN → text
one

ENCODER

A — person — riding — ... — <END>

thought vector

<START>

DECODER

CNN *Encoder* (Inception) - RNN *Decoder* (LSTM) (Vinyals dkk., 2014)

# Many to One: Text Classification



many to one

timestep
terakhir

$W^H$    $W^H$    $W^H$

$h_{t-1}$    $h_t$    $h_{t+1}$     ● ● ●     Binary Softmax Classifier     [ 0.09, .91 ]

$W^X_{t-1}$    $W^X_t$    $W^X_{t+1}$

$x_{t-1}$    $x_t$    $x_{t+1}$

The    movie    was

Max Sequence Length

https://www.oreilly.com/learning/perform-sentiment-analysis-with-lstms-using-tensorflow

31

# Many to Many: Sequence Tagging



Input is a sequence of words, and output is the sequence of POS tag for each word.

# Many to Many: Machine Translation

- Machine Translation: input is a sequence of words in source language (e.g. German). Output is a sequence of words in target language (e.g. English).
- A key difference is that our output only starts after we have seen the complete input, because the first word of our translated sentences may require information captured from the complete input sequence.
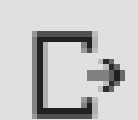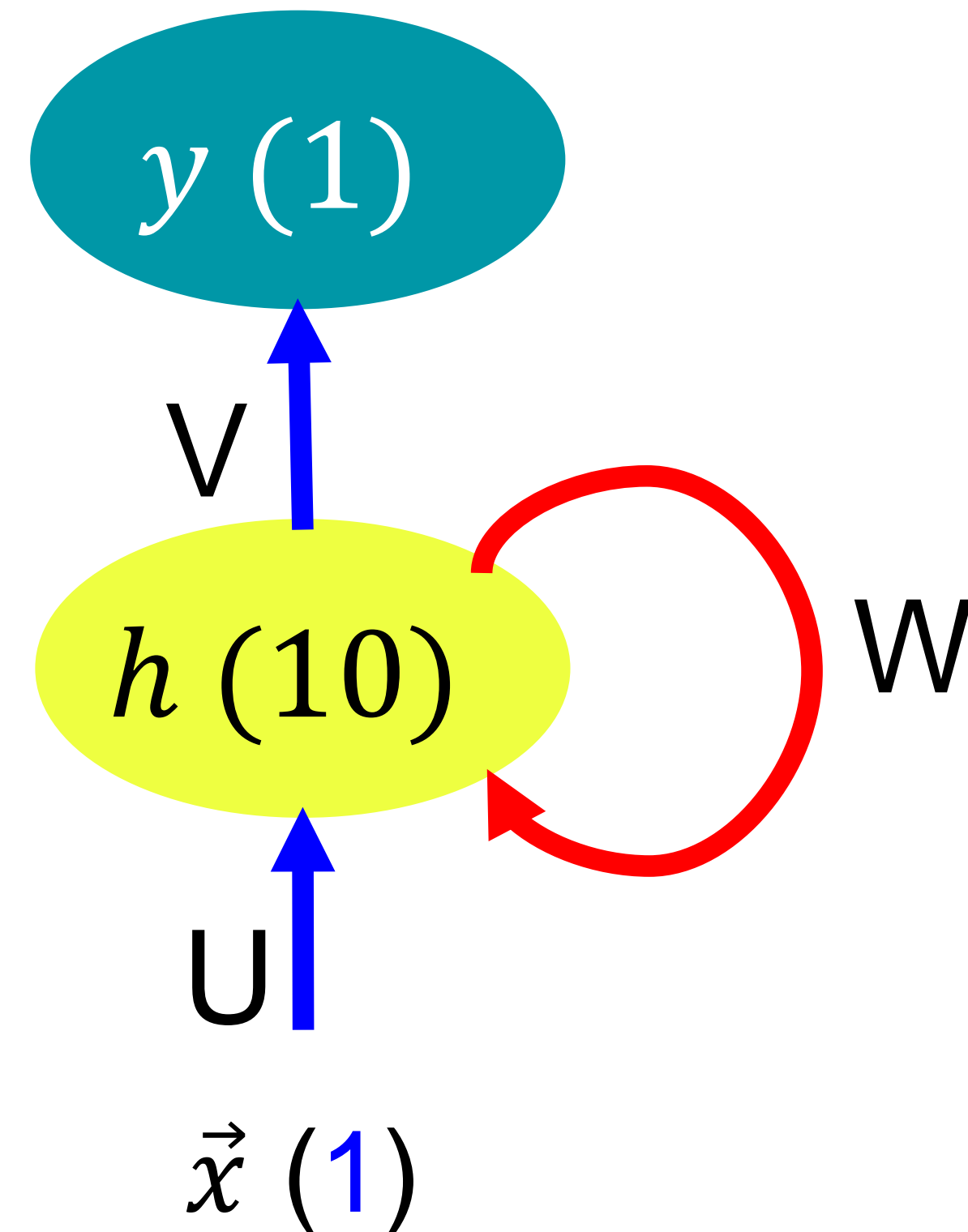
http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/

33

# Implementing RNN on Keras: Many to One



# predict amazon stock closing prices, RNN 50 timestep

```python
from keras import Sequential
from keras.layers import SimpleRNN, Dense
model = Sequential()
model.add(SimpleRNN(10, input_shape=(50,1)))
#simple recurrent layer, 10 neurons & process
50x1 sequences
model.add(Dense(1,activation='linear')) #linear
output because this is a regression problem
```

# Number of Parameter



$y$ (1)

V

$h$ (10)    W

U

$\vec{x}$ (1)

U : 10 × (1+1) = 20
W : 10 × 10 = 100     } 131
V : 1 × (10+1) = 11

Model: "sequential"

| Layer (type) | Output Shape | Param # |
|---|---|---|
| simple_rnn (SimpleRNN) | (None, 10) | 120 |
| dense (Dense) | (None, 1) | 11 |

Total params: 131
Trainable params: 131
Non-trainable params: 0

Total parameter = (1+10+1)*10+(10+1)*1=131

Simple RNN:
U: matrix hidden neurons x (input dimension + 1)
W: matrix hidden neurons x hidden neurons
V: matrix output neurons x (hidden neurons+1)

# Number of Parameter: Example 2

```python
model = Sequential() #initialize model
model.add(SimpleRNN(64, input_shape=(50,1), return_sequences=True))#64 neurons
model.add(SimpleRNN(32, return_sequences=True))#32 neurons
model.add(SimpleRNN(16)) #16 neurons
model.add(Dense(8,activation='tanh'))
model.add(Dense(1,activation='linear'))
```

1) $U = 64 \times (1+1) = 128$
$W = 64 \times 64 = 4096$ } 4224

2) $U = 32 \times (64+1) = 2080$
$W = 32 \times 32 = 1024$ } 3104

3) $U = 16 \times (32+1) = 528$
$W = 16 \times 16 = 256$ } 784

4) $V = 8 \times (16+1) = 136$

5) $V = 1 \times (8+1) = 9$

Model: "sequential_1"

| Layer (type) | Output Shape | Param # |
|---|---|---|
| simple_rnn_1 (SimpleRNN) | (None, 50, 64) | 4224 |
| simple_rnn_2 (SimpleRNN) | (None, 50, 32) | 3104 |
| simple_rnn_3 (SimpleRNN) | (None, 16) | 784 |
| dense_1 (Dense) | (None, 8) | 136 |
| dense_2 (Dense) | (None, 1) | 9 |

Total params: 8,257
Trainable params: 8,257
Non-trainable params: 0

= (1+64+1)*64=4224

= (64+32+1)*32=3104

= (32+16+1)*16=784
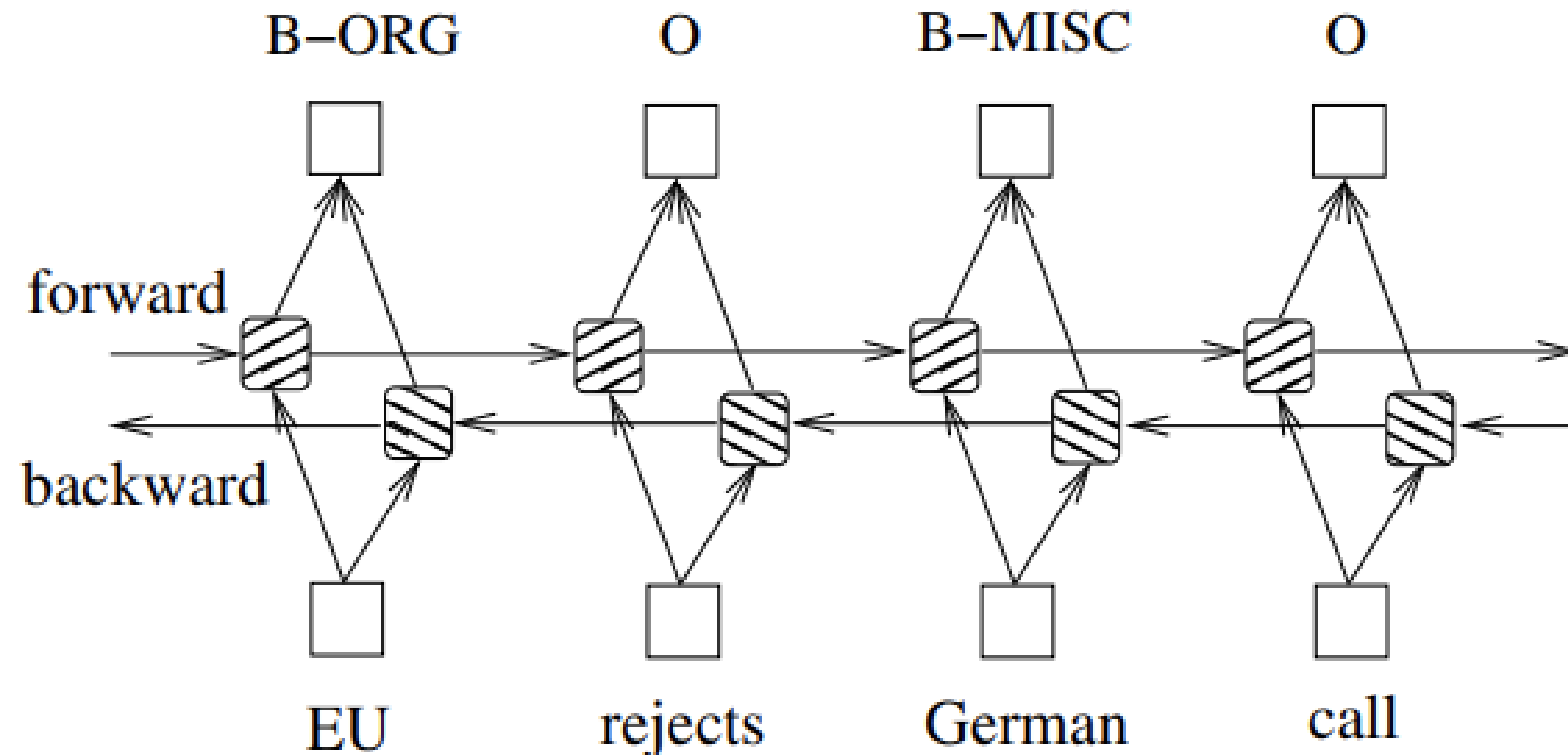
= (16+1)*8=136

= (8+1)*1=9

Total parameter = 8257

# Bidirectional RNNs

In many applications we want to output a prediction of y (t) which may depend on the whole input sequence. E.g. co-articulation in speech recognition, right neighbors in POS tagging, etc. **Bidirectional RNNs** combine an RNN that moves forward through time beginning from the start of the sequence with another RNN that moves backward through time beginning from the end of the sequence.



22nya jd input utk
output layer

https://www.cs.toronto.edu/~tingwuwang/rnn_tutorial.pdf

37

# Bidirectional RNNs for Information Extraction



https://www.depends-on-the-definition.com/sequence-tagging-lstm-crf/

38

# Summary

| Architecture: 1-to-n, n-to-1, n-to-n | Number of parameter RNN | Bidirectional RNN |
|---|---|---|

LSTM →

# Thank you