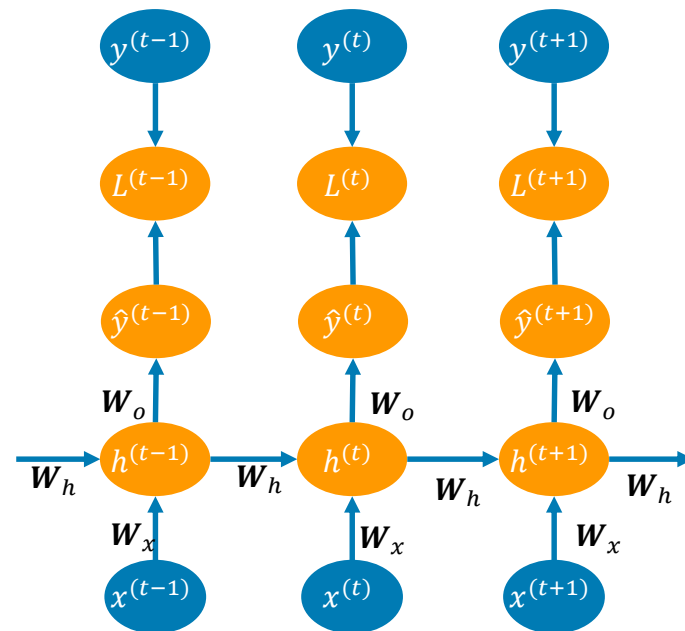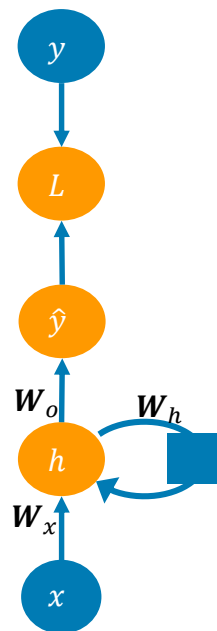# VL Deep Learning for Natural Language Processing

12. Recurrent Neural Networks

*Prof. Dr. Ralf Krestel*

*AG Information Profiling and Retrieval*

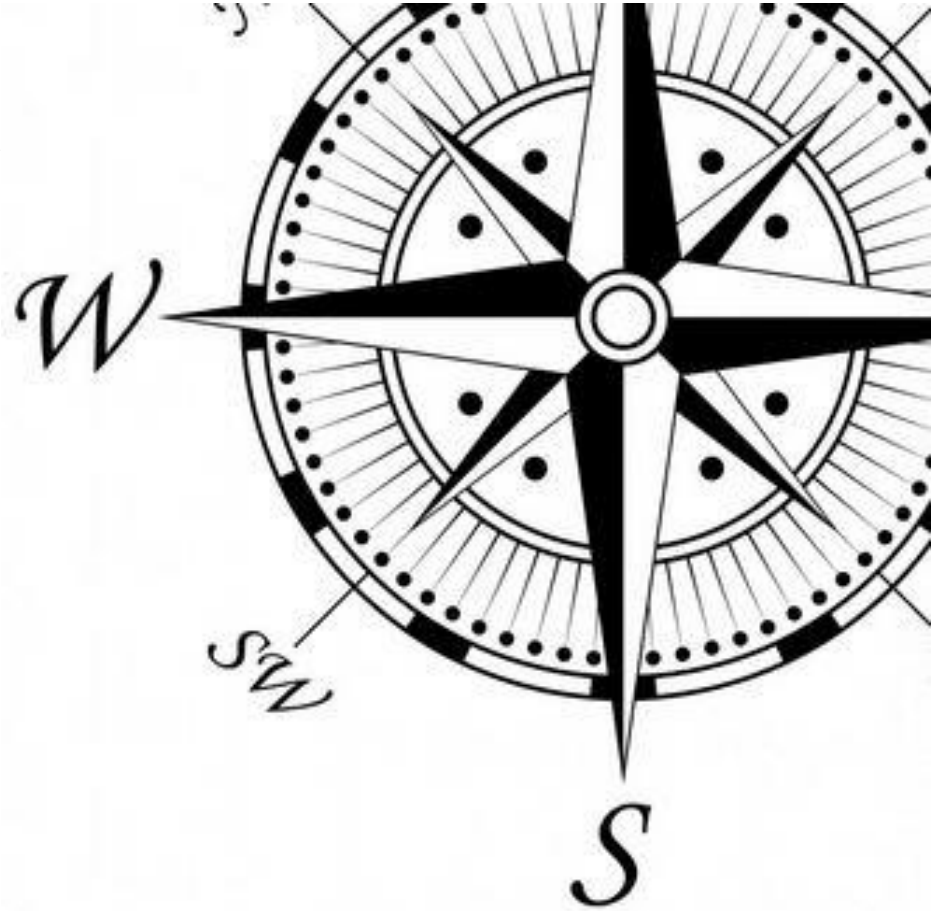# Recurrent Neural Networks

# Learning Goals for this Chapter

- Adapt task description so that RNNs can be used to solve the problem
- Understand BPTT
- Explain different kinds of RNNs and how they work
- Implement and evaluate a simple RNN model

- Relevant chapters:
  - P6.2
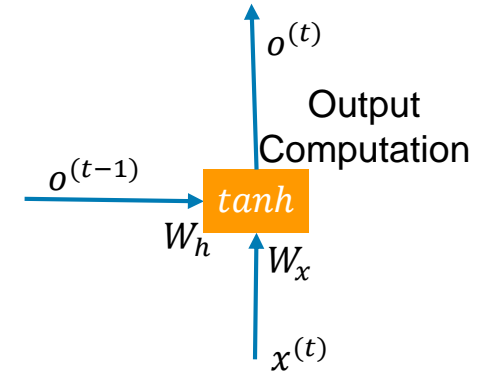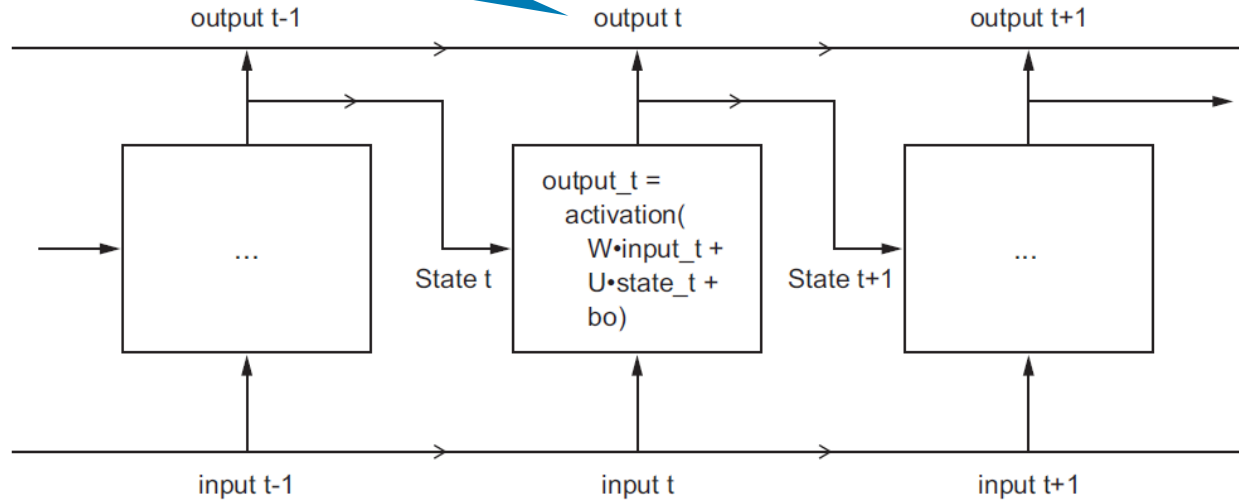  - S6 (2021) https://www.youtube.com/watch?v=0LixFSa7yts

# Topics Today

1. **Recurrent Neural Networks (RNN)**
2. Backprop Through Time (BPTT)
3. Different RNN Types
4. A Simple RNN

# Rolled-Out Layer

- Output of layer at timestep $t$
  - $h^{(t)} = \tanh(W_h h^{(t-1)} + W_x x^{(t)} + b_h)$

**Output layer above, e.g. softmax**



output t-1      output t      output t+1

output_t =
activation(
W•input_t +
U•state_t +
bo)

State t    State t+1

input t-1      input t      input t+1

$o^{(t)}$

Output Computation

$o^{(t-1)}$

$tanh$

$W_h$   $W_x$

$x^{(t)}$

# Forward Computation

- $L^{(t)}(\theta) = CE\left(\hat{\boldsymbol{y}}^{(t)}, \boldsymbol{y}^{(t)}\right) = -\sum_{j=1}^{|V|} y_j^{(t)} \log \hat{y}_j^{(t)}$

- Total loss is averaged:

$$L(\theta) = \frac{1}{T} \sum_{t=1}^{T} L^{(t)}(\theta)$$

# Forward Computation Implementation

```python
import numpy as np
timesteps = 100
input_features = 32
output_features = 64
inputs = np.random.random((timesteps, input_features))
h_t = np.zeros((output_features,))
W_x = np.random.random((output_features, input_features))
W_h = np.random.random((output_features, output_features))
b_h = np.random.random((output_features,))
successive_outputs = []
for input_t in inputs:
        output_t = np.tanh(np.dot(W_x, input_t) + np.dot(W_h, h_t) + b_h)
        successive_outputs.append(output_t)
        h_t = output_t
final_output_sequence = np.concatenate(successive_outputs, axis=0)
```

**Components of an input sequence**

**Here: Input is random noise**
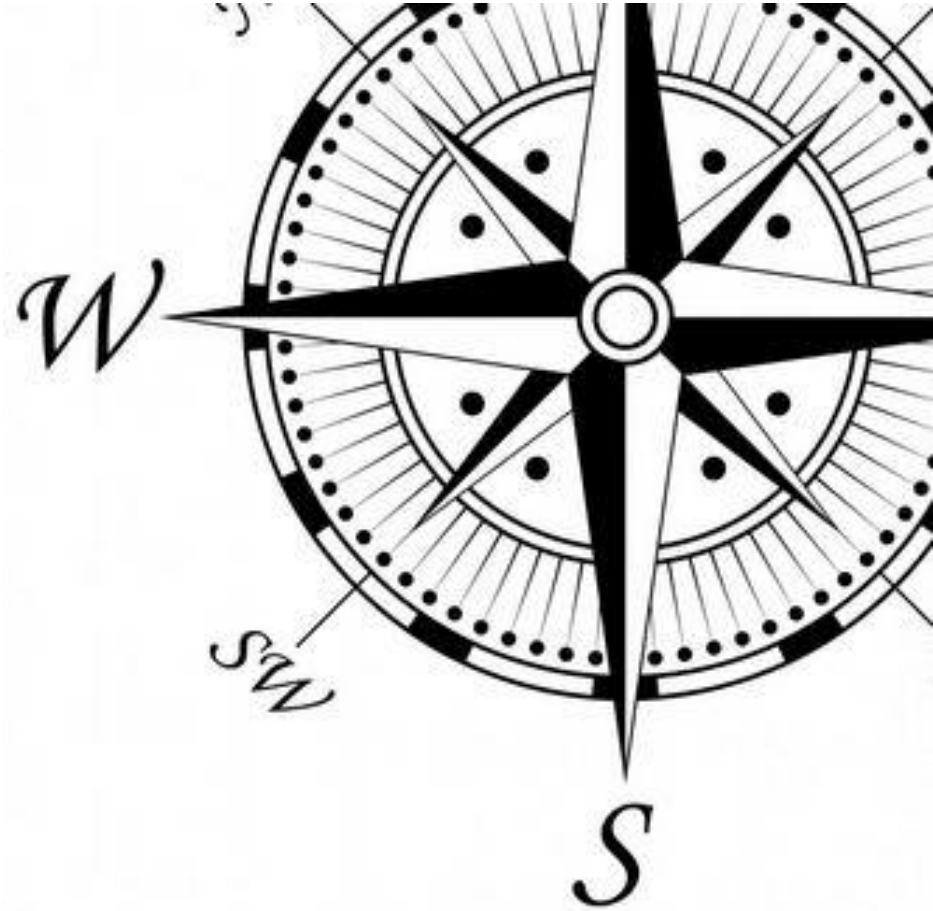
**Initial state = 0-vector**

**Also the weight matrices are randomly initialized**

**output_t, sufficient, since it contains information about the whole sequence**

**Output is a 2d-tensor of shape (timesteps, output_features)**

# Topics Today

1. Recurrent Neural Networks (RNN)
2. **Backprop Through Time (BPTT)**
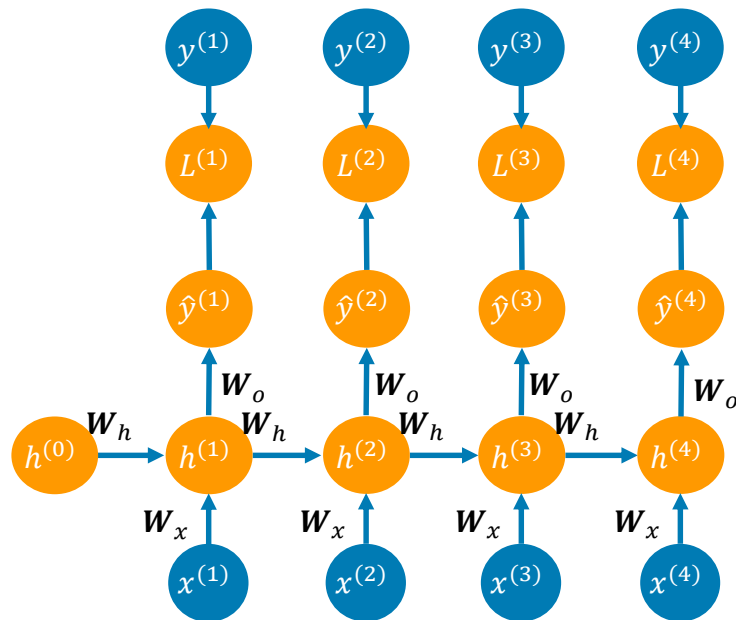3. Different RNN Types
4. A Simple RNN

# Backward Computation I

- **Backpropagation through time (BPTT)**
- Given: Multi-variable function $f(x, y)$ and two functions with one variable $x(t)$ and $y(t)$, then this is the multi-variable chain rule

$$\frac{d}{dt} f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

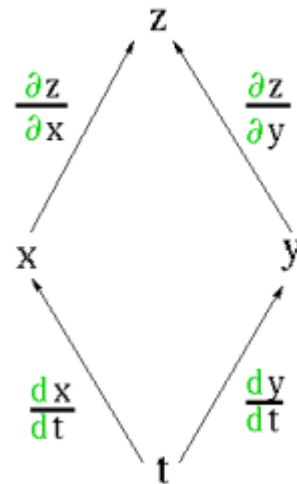# Multi-Variable Chain Rule: Example

- Let $z = x^2 y - y^2$ where $x$ and $y$ are parametrized as $x = t^2$ and $y = 2t$
- Then

$$\frac{dz}{dt} = \frac{\partial z}{\partial x}\frac{dx}{dt} + \frac{\partial z}{\partial y}\frac{dy}{dt}$$
$$= (2xy)(2t) + (x^2 - 2y)(2)$$
$$= (2t^2 \cdot 2t)(2t) + \left((t^2)^2 - 2(2t)\right)(2)$$
$$= 8t^4 + 2t^4 - 8t$$
$$= 10t^4 - 8t$$
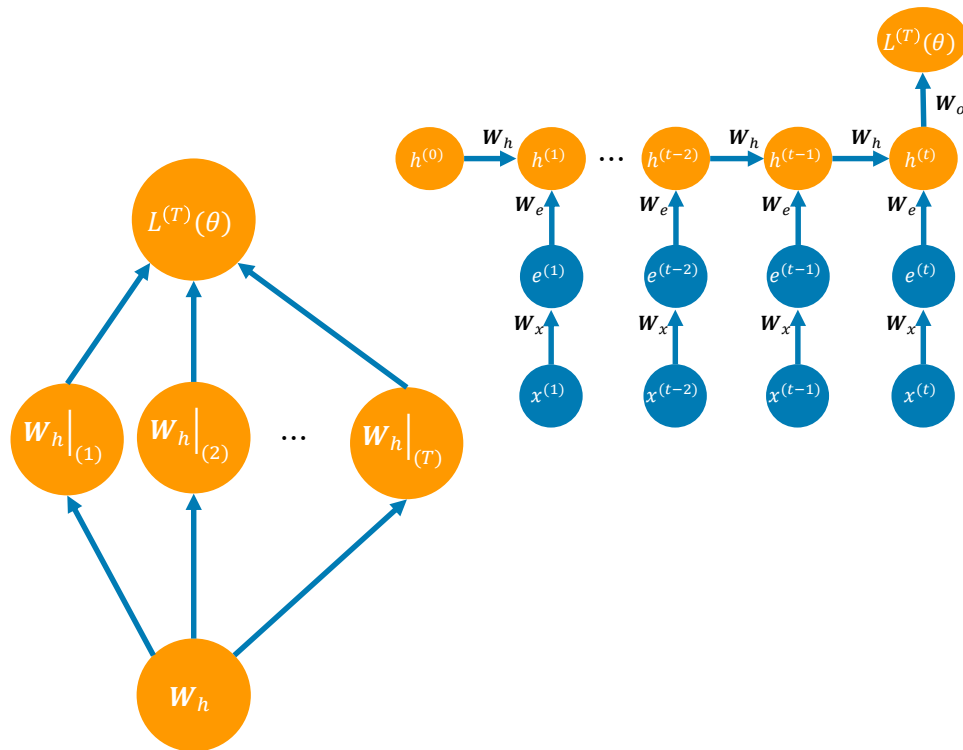
# Backward Computation II

$$\frac{d}{dt}f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x}\frac{dx}{dt} + \frac{\partial f}{\partial y}\frac{dy}{dt}$$

$$L(\theta) = \frac{1}{T}\sum_{t=1}^{T} L^{(t)}(\theta)$$

- Derivation of the loss fuction $L^{(t)}(\theta)$ with respect to repeating $\boldsymbol{W}_h$

$$\frac{\partial L^{(T)}}{\partial \boldsymbol{W}_h} = \sum_{t=1}^{T}\left(\frac{\partial L^{(T)}}{\partial \boldsymbol{W}_h}\bigg|_{(t)} \cdot \frac{\partial \boldsymbol{W}_h|_{(t)}}{\partial \boldsymbol{W}_h}\right)$$

$= 1$

$$= \sum_{t=1}^{T}\frac{\partial L^{(T)}}{\partial \boldsymbol{W}_h}\bigg|_{(t)}$$

# Back Propagation Through Time (BPTT) I



$$L^{(t)}(\theta) = CE\big(\hat{\boldsymbol{y}}^{(t)}, \boldsymbol{y}^{(t)}\big) = -\sum_{j=1}^{|V|} y_j^{(t)} \log \hat{y}_j^{(t)}$$

$$\hat{\boldsymbol{y}}^{(t)} = softmax\big(\boldsymbol{h}^{(t)} \boldsymbol{W}_o + \boldsymbol{b}_o\big)$$

$$\boldsymbol{z}_1^{(t)}$$

$$\boldsymbol{h}^{(t)} = sigmoid\big(\boldsymbol{h}^{(t-1)} \boldsymbol{W}_h + \boldsymbol{e}^{(t)} \boldsymbol{W}_e + \boldsymbol{b}_h\big)$$

$$\boldsymbol{z}_2^{(t)}$$

$$\boldsymbol{e}^{(t)} = \boldsymbol{x}^{(t)} \boldsymbol{W}_x$$

# Back Propagation Through Time (BPTT) II



$$L^{(t)}(\theta) = CE(\hat{\boldsymbol{y}}^{(t)}, \boldsymbol{y}^{(t)}) = -\sum_{j=1}^{|V|} y_j^{(t)} \log \hat{y}_j^{(t)}$$

$$\hat{\boldsymbol{y}}^{(t)} = softmax\big(\underbrace{\boldsymbol{h}^{(t)}\boldsymbol{W}_o + \boldsymbol{b}_o}_{\boldsymbol{z}_1^{(t)}}\big)$$

$$\boldsymbol{h}^{(t)} = sigmoid\big(\underbrace{\boldsymbol{h}^{(t-1)}\boldsymbol{W}_h + \boldsymbol{e}^{(t)}\boldsymbol{W}_e + \boldsymbol{b}_h}_{\boldsymbol{z}_2^{(t)}}\big)$$

$$\boldsymbol{e}^{(t)} = \boldsymbol{x}^{(t)}\boldsymbol{W}_x$$

**Error signal**

$$\boldsymbol{\delta}_1^{(t)} = \frac{\partial L^{(t)}}{\partial \boldsymbol{z}_1^{(t)}} = \hat{\boldsymbol{y}}^{(t)} - \boldsymbol{y}^{(t)}$$

$$\boldsymbol{\delta}_2^{(t)} = \frac{\partial L^{(t)}}{\partial \boldsymbol{z}_2^{(t)}} = \frac{\partial L^{(t)}}{\partial \boldsymbol{z}_1^{(t)}} \frac{\partial \boldsymbol{z}_1^{(t)}}{\partial \boldsymbol{h}^{(t)}} \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{z}_2^{(t)}} = \boldsymbol{\delta}_1^{(t)} \boldsymbol{W}_o^T \odot \sigma'(\boldsymbol{z}_2^{(t)})$$

# Gradients of an RNN I

$$\frac{\partial L^{(t)}}{\partial W_o} = \frac{\partial L^{(t)}}{\partial z_1^{(t)}} \frac{\partial z_1^{(t)}}{\partial W_o} = \boldsymbol{\delta}_1^{(t)} \otimes \boldsymbol{h}^{(t)}$$

$$\frac{\partial L^{(t)}}{\partial b_o} = \boldsymbol{\delta}_1^{(t)}$$

**Outer product**

$$\left.\frac{\partial L^{(t)}}{\partial W_e}\right|_{(t)} = \frac{\partial L^{(t)}}{\partial z_2^{(t)}} \frac{\partial z_2^{(t)}}{\partial W_o} = \boldsymbol{\delta}_2^{(t)} \otimes \boldsymbol{e}^{(t)}$$

$$\left.\frac{\partial L^{(t)}}{\partial W_h}\right|_{(t)} = \frac{\partial L^{(t)}}{\partial z_2^{(t)}} \frac{\partial z_2^{(t)}}{\partial W_o} = \boldsymbol{\delta}_2^{(t)} \otimes \boldsymbol{h}^{(t-1)}$$

$$\left.\frac{\partial L^{(t)}}{\partial b_h}\right|_{(t)} = \frac{\partial L^{(t)}}{\partial z_2^{(t)}} \frac{\partial z_2^{(t)}}{\partial b_h} = \boldsymbol{\delta}_2^{(t)}$$

$$\frac{\partial L^{(t)}}{\partial W_{x_{x(t)}}} = \frac{\partial L^{(t)}}{\partial z_2^{(t)}} \frac{\partial z_2^{(t)}}{\partial W_{x_{x(t)}}} = \boldsymbol{\delta}_2^{(t)} \otimes \boldsymbol{W}_e$$

$$\frac{\partial L^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} = \frac{\partial L^{(t)}}{\partial z_2^{(t-1)}} \frac{\partial z_2^{(t-1)}}{\partial \boldsymbol{h}^{(t-1)}} = \boldsymbol{\delta}_2^{(t-1)} \otimes \boldsymbol{W}_h$$

$$\boldsymbol{\delta}_1^{(t)} = \frac{\partial L^{(t)}}{\partial z_1^{(t)}} = \widehat{\boldsymbol{y}}^{(t)} - \mathbf{y}^{(t)}$$

**Hadamard product**

$$\boldsymbol{\delta}_2^{(t)} = \frac{\partial L^{(t)}}{\partial z_2^{(t)}} = \frac{\partial L^{(t)}}{\partial z_1^{(t)}} \frac{\partial z_1^{(t)}}{\partial \boldsymbol{h}^{(t)}} \frac{\partial \boldsymbol{h}^{(t)}}{\partial z_2^{(t)}} = \boldsymbol{\delta}_1^{(t)} \boldsymbol{W}_o^T \odot \sigma'(z_2^{(t)})$$

$$L^{(t)}(\theta) = CE\big(\widehat{\boldsymbol{y}}^{(t)}, \boldsymbol{y}^{(t)}\big) = -\sum_{j=1}^{|V|} y_j^{(t)} \log \widehat{y}_j^{(t)}$$

$$\widehat{\boldsymbol{y}}^{(t)} = softmax(\underbrace{\boldsymbol{h}^{(t)} \boldsymbol{W}_o + \boldsymbol{b}_o}_{\boldsymbol{z}_1^{(t)}})$$

$$\boldsymbol{h}^{(t)} = sigmoid\big(\underbrace{\boldsymbol{h}^{(t-1)} \boldsymbol{W}_h + \boldsymbol{e}^{(t)} \boldsymbol{W}_e + \boldsymbol{b}_h}_{\boldsymbol{z}_2^{(t)}}\big)$$

$$\boldsymbol{e}^{(t)} = \boldsymbol{x}^{(t)} \boldsymbol{W}_x$$

# Gradients of an RNN II

$$\frac{\partial L^{(t)}}{\partial W_e}\bigg|_{(t-1)} = (e^{(t-1)})^T (\delta_3^{(t-1)} \odot \sigma'(z_2^{(t-1)}))$$

$$\frac{\partial L^{(t)}}{\partial W_h}\bigg|_{(t-1)} = (h^{(t-2)})^T (\delta_3^{(t-1)} \odot \sigma'(z_2^{(t-1)}))$$

$$\frac{\partial L^{(t)}}{\partial b_h}\bigg|_{(t-1)} = \delta_3^{(t-1)} \odot \sigma'(z_2^{(t-1)})$$

$$\frac{\partial L^{(t)}}{\partial W_{x_{x(t-1)}}}\bigg|_{(t-1)} = \delta_3^{(t-1)} \odot \sigma'(z_2^{(t-1)}) W_e^T$$

$$\delta_1^{(t)} = \frac{\partial L^{(t)}}{\partial z_1^{(t)}} = \hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)}$$

$$\delta_3^{(t-1)} = \delta_1^{(t-1)} W_o^T$$

$$L^{(t)}(\theta) = CE(\hat{\mathbf{y}}^{(t)}, \mathbf{y}^{(t)}) = -\sum_{j=1}^{|V|} y_j^{(t)} log \hat{y}_j^{(t)}$$

$$\hat{\mathbf{y}}^{(t)} = softmax(\underbrace{h^{(t)} W_o + b_o}_{z_1^{(t)}})$$

$$h^{(t)} = sigmoid(\underbrace{h^{(t-1)} W_h + e^{(t)} W_e + b_h}_{z_2^{(t)}})$$

$$e^{(t)} = x^{(t)} W_x$$

# Tutorials/Blogposts/Videos…
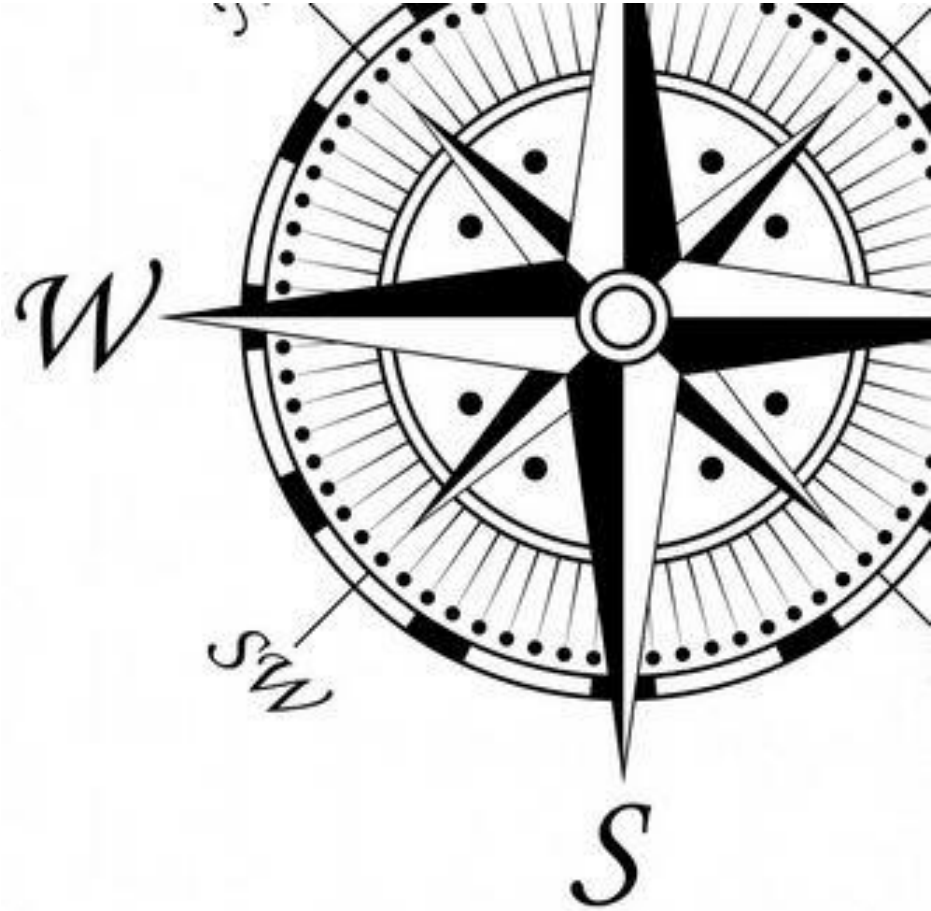
- https://github.com/go2carter/nn-learn/blob/master/grad-deriv-tex/rnn-grad-deriv.pdf
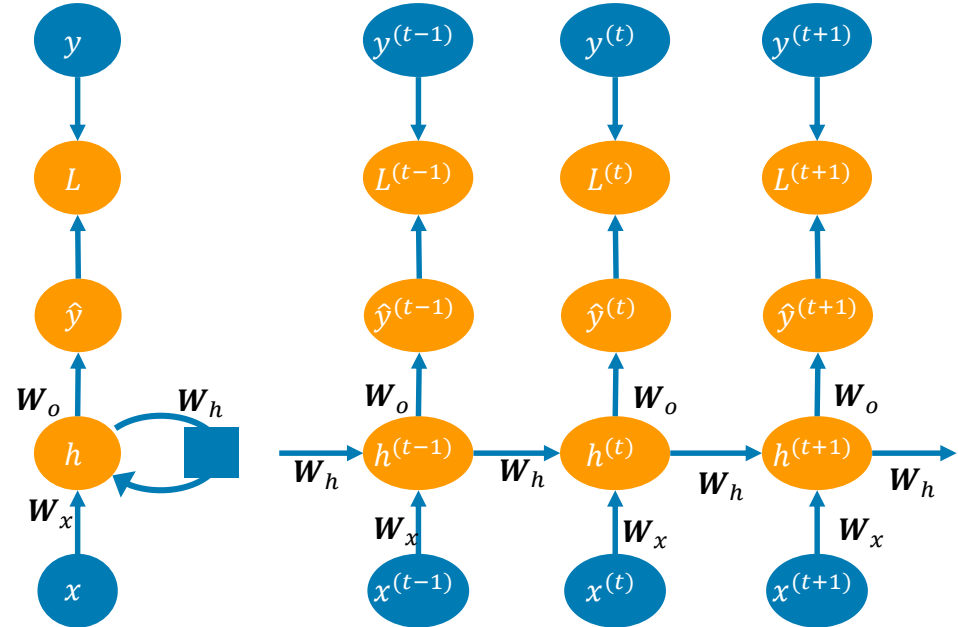- https://mmuratarat.github.io/2019-02-07/bptt-of-rnn
- https://arxiv.org/pdf/1610.02583.pdf
- https://www.youtube.com/watch?v=nFTQ7kHQWtc
- https://www.youtube.com/watch?v=q4mVeRLitsU

# Topics Today

1. Recurrent Neural Networks (RNN)
2. Backprop Through Time (BPTT)
3. **Different RNN Types**
4. A Simple RNN

ZBW CAU

Leibniz
Gemeinschaft

# Categorization

- RNNs can be categorized based on their input and output.
  - **Many-to-many**
    - Standard RNN
    - Sequence-to-Sequence
  - **Many-to-one**
    - Summary
  - **One-to-many**
    - Generative models
  - **(One-to-one)**
    - Standard NN
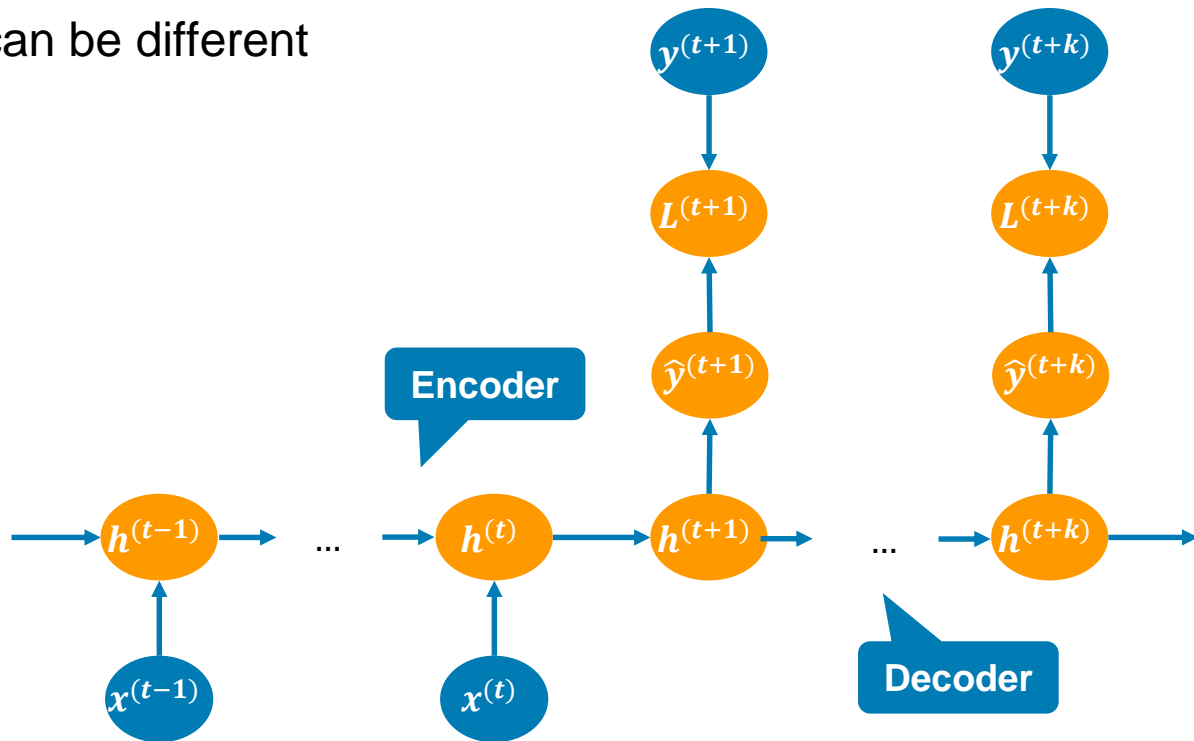
# Many-to-Many: Standard RNN

- RNNs are very powerful
  - Turing-complete

- Through the hidden layer, access to all information from the past

- Input length equals output length
  - At each time step $t$ there is an input $x^{(t)}$ and an output $y^{(t)}$

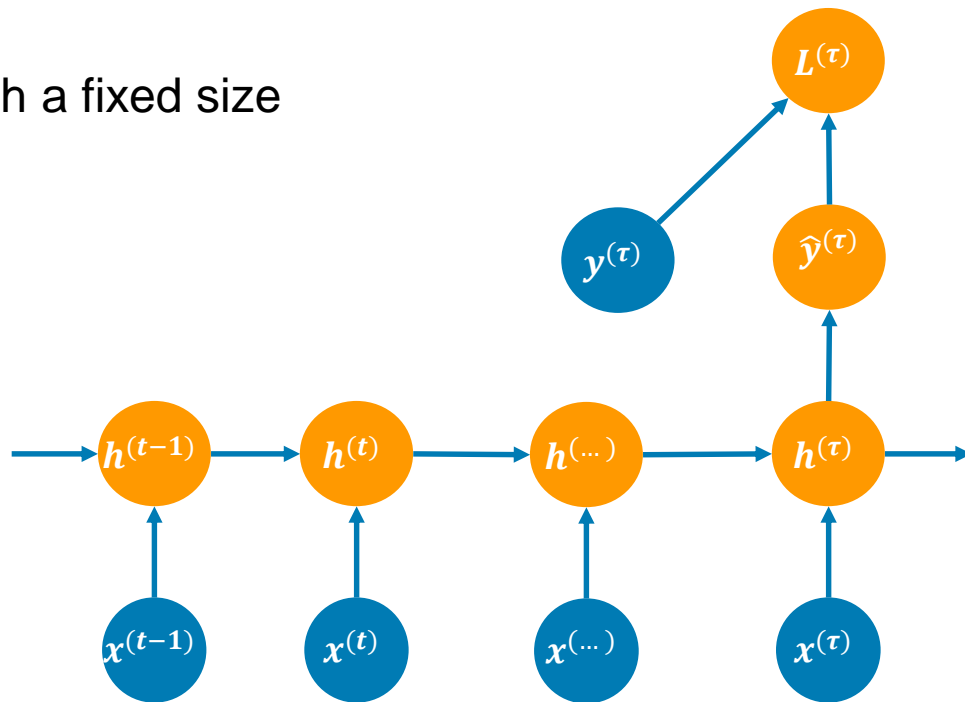# Many-to-Many: Sequence-to-Sequence

- Input and output length can be different
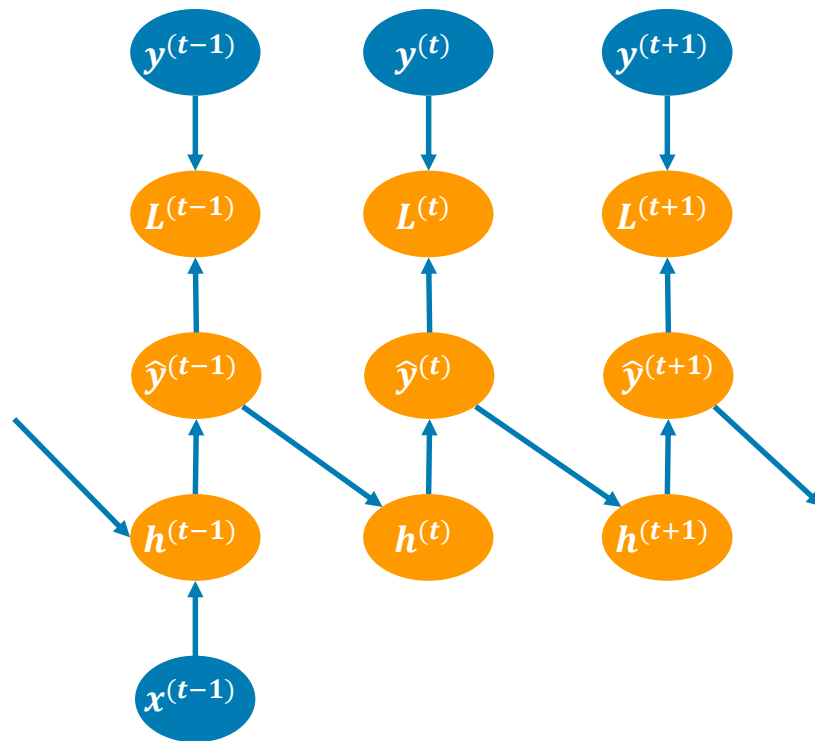- E.g. machine translation

# Many-to-One: Summary

- Can summarize a sequence
  - Representation of a sequence with a fixed size
- E.g. sentiment analyse

- As input for further layers
  - Output will be learnt by backprop

- Here with ist own target value

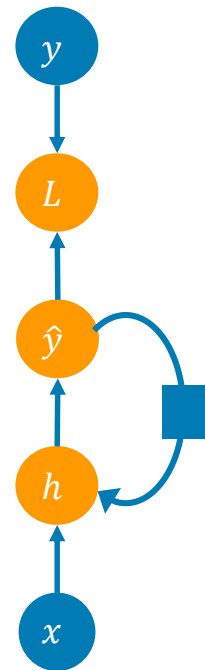# One-to-Many: Generative Models

- E.g. music generation
- Input could be, e.g., an integer, to decide on the style
  - Input could also be empty
    $x = \emptyset$

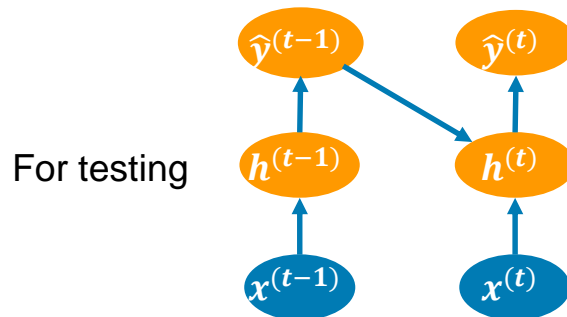# Many/One-to-Many: RNN Variants

- Less powerful
  - Only output is transmitted
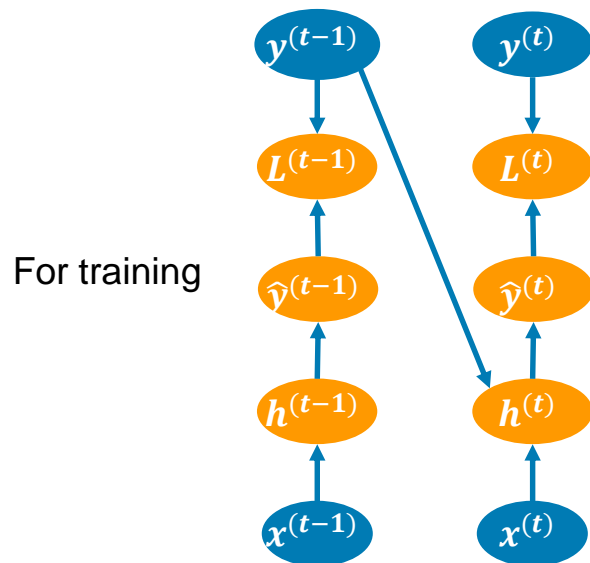
- Only indirect connection of previous hidden layer to current hidden layer
  - Via output layer
    - Typically less dimensions than hidden layer

- Potentially simpler to train
  - Parallelization

# Many/One-to-Many: Teacher Forcing

- Method for RNNs to learn from ground-truth
- When the model is deployed, the ground-truth is approximated by the output



For training

For testing

# Recurrent Layer
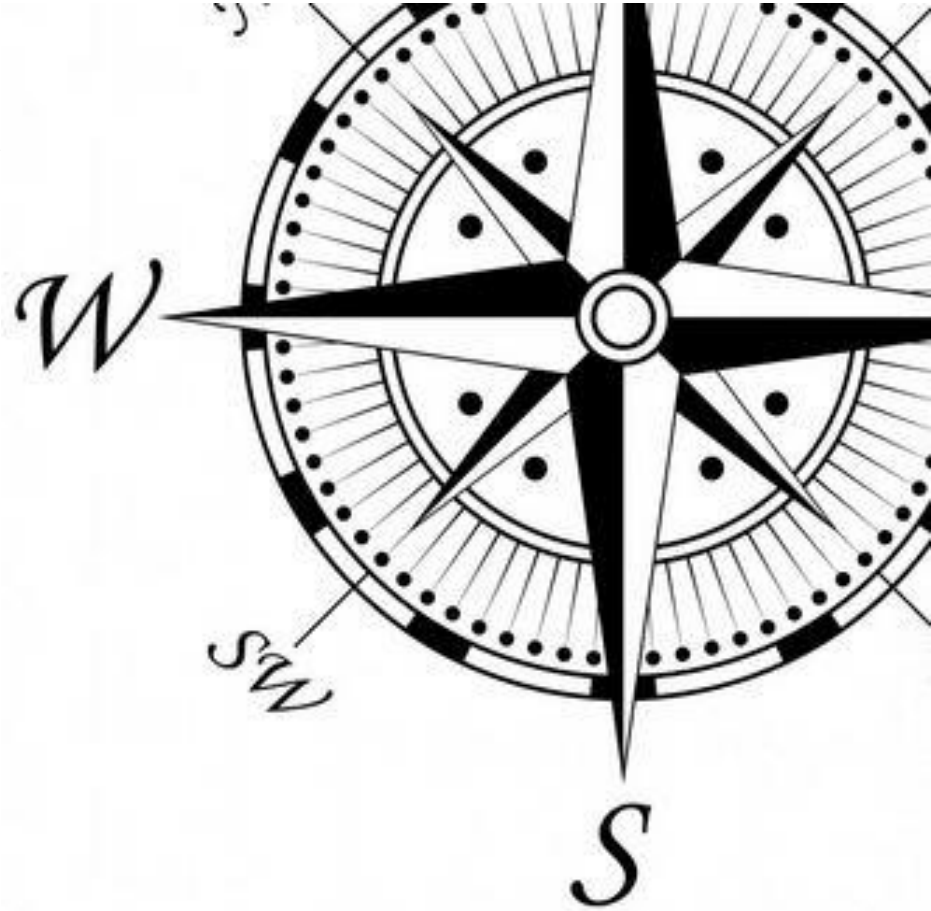
- Think about suitable scenarios for the presented RNNs variants.
- What are advantages and disadvantages of teacher forcing?

Start  5  4  3  2  1  End

# Topics Today

1. Recurrent Neural Networks (RNN)
2. Backprop Through Time (BPTT)
3. Different RNN Types
4. **A Simple RNN**

# Simple RNN in Keras

```python
from keras.models import Sequential
from keras.layers import Embedding, SimpleRNN
model = Sequential()
model.add(Embedding(10000, 32))
model.add(SimpleRNN(32, return_sequences=True))
model.summary()
```

```
Layer (type)                    Output Shape              Param #
=================================================================
embedding_22 (Embedding)        (None, None, 32)          320000

simplernn_10 (SimpleRNN)        (None, 32)                2080
=================================================================
Total params: 322,080
Trainable params: 322,080
Non-trainable params: 0
```

# Stacked Simple RNN in Keras

```
model = Sequential()
model.add(Embedding(10000, 32))
model.add(SimpleRNN(32, return_sequences=True))
model.add(SimpleRNN(32, return_sequences=True))
model.add(SimpleRNN(32, return_sequences=True))
model.add(SimpleRNN(32))
model.summary()
```

| Layer (type) | Output Shape | Param # |
|---|---|---|
| embedding_24 (Embedding) | (None, None, 32) | 320000 |
| simplernn_12 (SimpleRNN) | (None, None, 32) | 2080 |
| simplernn_13 (SimpleRNN) | (None, None, 32) | 2080 |
| simplernn_14 (SimpleRNN) | (None, None, 32) | 2080 |
| simplernn_15 (SimpleRNN) | (None, 32) | 2080 |

```
Total params: 328,320
Trainable params: 328,320
Non-trainable params: 0
```

# IMDB Example: Preprocessing

```python
from keras.datasets import imdb
from keras.preprocessing import sequence
max_features = 10000
maxlen = 500
batch_size = 32
print('Loading data...')
(input_train, y_train), (input_test, y_test) =
imdb.load_data(num_words=max_features)
print(len(input_train), 'train sequences')
print(len(input_test), 'test sequences')
print('Pad sequences (samples x time)')
input_train = sequence.pad_sequences(input_train, maxlen=maxlen)
input_test = sequence.pad_sequences(input_test, maxlen=maxlen)
print('input_train shape:', input_train.shape)
print('input_test shape:', input_test.shape)
```

Batch processing, i.e. multiple sequences simultaneously
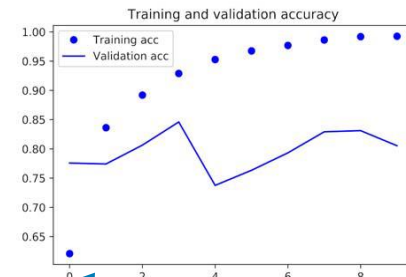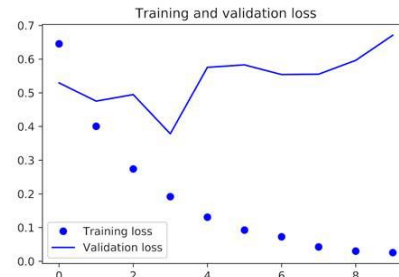
# IMDB Example: Model

```
from keras.layers import Dense
model = Sequential()
model.add(Embedding(max_features, 32))
model.add(SimpleRNN(32))
model.add(Dense(1, activation='sigmoid'))
model.compile(optimizer='rmsprop',
          loss='binary_crossentropy',
          metrics=['acc'])
history = model.fit(input_train, y_train,
            epochs=10,
            batch_size=batch_size,
            validation_split=0.2)
```

# IMDB Example: Validation

```python
import matplotlib.pyplot as plt
acc = history.history['acc']
val_acc = history.history['val_acc']
loss = history.history['loss']
val_loss = history.history['val_loss']
epochs = range(1, len(acc) + 1)


plt.plot(epochs, acc, 'bo', label='Training acc')
plt.plot(epochs, val_acc, 'b', label='Validation acc')
plt.title('Training and validation accuracy')
plt.legend()
plt.figure()
plt.plot(epochs, loss, 'bo', label='Training loss')
plt.plot(epochs, val_loss, 'b', label='Validation loss')
plt.title('Training and validation loss')
plt.legend()
plt.show()
```



Training and validation loss



Training and validation accuracy

**Test accuracy without RNN: 88%;**
**With RNN only 85%;**

# Exercise

- Implement the previous example
  - Stack RNN-layers
  - Play around with other hyperparameters

Start — 10 — 9 — 8 — 7 — 6 — 5 — 4 — 3 — 2 — 1 — End

# Learning Goals for this Chapter

- Adapt task description so that RNNs can be used to solve the problem
- Understand BPTT
- Explain different kinds of RNNs and how they work
- Implement and evaluate a simple RNN model

- Relevant chapters:
  - P6.2
  - S6 (2021) https://www.youtube.com/watch?v=0LixFSa7yts