

# Grundlagen der Künstlichen Intelligenz

## Wintersemester 25/26



### Vorlesung 10

Rekurrente Neuronale Netze

Anwendung: Natural Language Processing

12. Januar 2026

T.T.-Prof. Dr. Peer Nowack  
Prof. Dr. Pascal Friederich

# Wir werden immer mal wieder Umfragen machen...

Sie können sich bereits einloggen. Nutzen Sie das KIT Wi-Fi bei schlechtem Empfang.



1

Go to [wooclap.com](https://wooclap.com)

2

Enter the event code in the top banner

Event code

**HIPFAX**

 You cannot vote anymore



Welche der folgenden Beispiele sind geordnete Sequenzdaten (alle richtigen Antworten ...)



1

Messdaten von Lufttemperatur in Karlsruhe

72%

73 ✓

2

Videos

78%

79 ✓

3

Fotos von handgeschriebenen Postleitzahlen

14%

14 

4

Tonspur eines Podcasts

90%

91 ✓

5

Der Text eines Zeitungsartikels

82%

83 ✓

6

Eine Zeitreihe für ein Regressionsproblem in dem Samples zur Kreuzvalidierung gemischt wurden.

26%

26 



90%

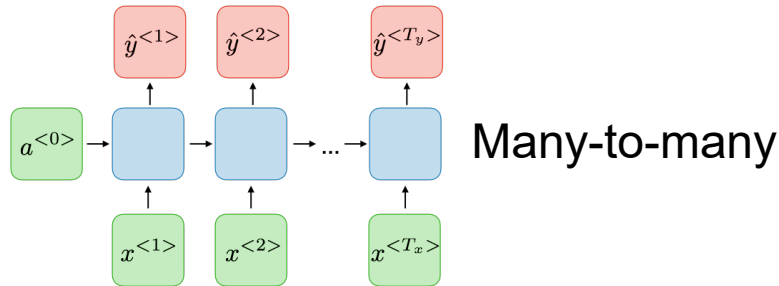
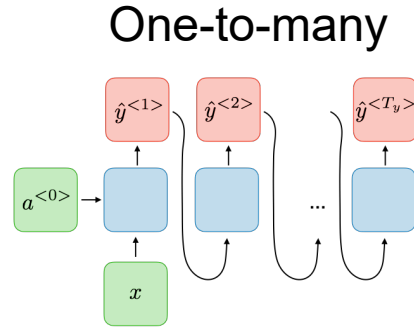
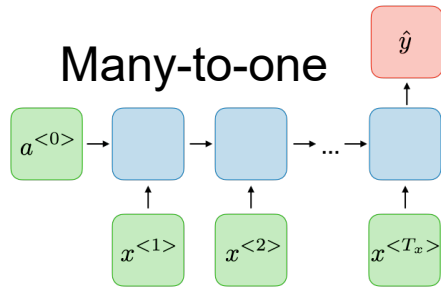


65% correct

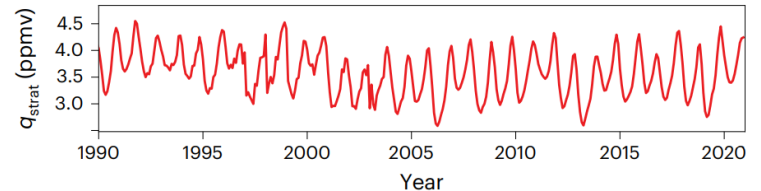
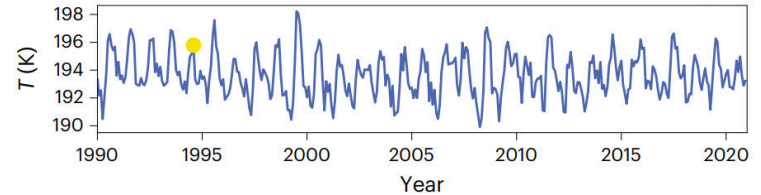
101 / 121



# Sequenzdaten (Textdaten, Zeitreihen von Messdaten, ...)



The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .



# KI-Landkarte

## Künstliche Intelligenz

### Modellierung und Schlussfolgerung

Variablen VL12 Inferenz

Logik VL11 Wissensrepräsentation

Zustände VL13 MDPs  
Suche

Reflex

### Anwendungen

Robotik VL14

Computer Vision VL9

Natürliche Sprache VL10

### Lernen

Optimierung und Generalisierung VL6

Vorhersage VL4 VL5 Neuronale Netze

Modellierung VL3 Unsupervised VL7  
Supervised VL8

### Historie und Philosophie

VL1 Geschichte

Personen

KI und Gesellschaft

Kritische Aspekte

### Mathematik

VL2

Lineare Algebra

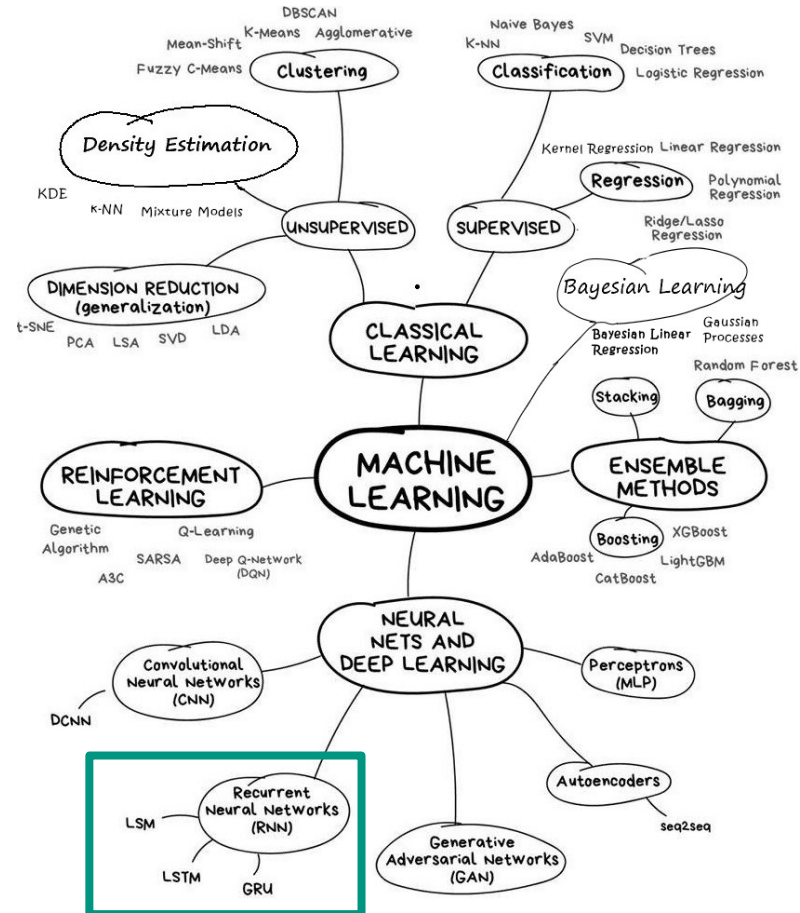
Statistik

Logik

Numerik

Analysis

# Der ML-Zoo...



# Die heutigen Vorlesungsfolien sind zu einem großen Teil basierend auf der Gastvorlesung von

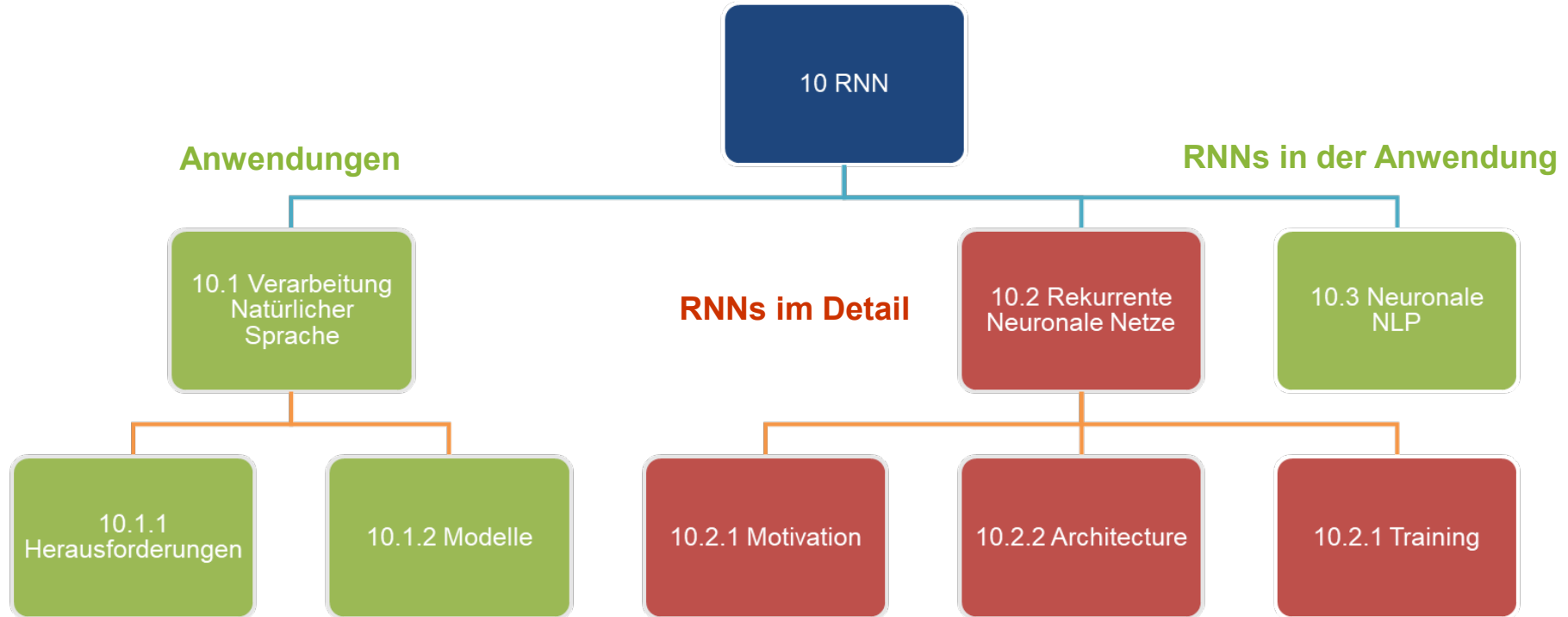


**Professor Jan Niehues**

Lehrstuhl KI für Sprachtechnologien

Institut für Anthropomatik

# Überblick



# Lernziele der heutigen Vorlesung

Nach dieser Vorlesung werden Sie in der Lage sein:

- Anwendungsbeispiele von **KI zur Text-/Sprachverarbeitung und Generierung** zu beschreiben und zu kontrastieren.
- Wichtige **Herausforderungen** in der Verarbeitung natürlicher Sprache zu benennen und aktuelle Lösungsstrategien zu erklären.
- **Recurrent Neural Networks (RNNs)** als eine bedeutende KI-Methode zu formulieren, die Architektur zu begründen, und deren Stärken und Schwächen zu bewerten.
- Das **Trainieren von RNNs mittels Backpropagation in der Zeit** zu erklären, sowie Limitierungen wie das **Vanishing/Exploding Gradients Problem** zu beschreiben.

# Verarbeitung natürlicher Sprache

## Was versteht man unter der „Verarbeitung natürlicher Sprache“?

### ChatGPT 5.2:

Was ist die "Verarbeitung natürlicher Sprache"?

Verarbeitung natürlicher Sprache (auf Englisch *Natural Language Processing*, kurz **NLP**) ist ein Teilgebiet der Künstlichen Intelligenz und der Informatik, das sich damit beschäftigt, wie Computer menschliche Sprache verstehen, verarbeiten und erzeugen können.

#### Einfach erklärt

NLP sorgt dafür, dass Maschinen mit Sprache arbeiten können – so wie Menschen es tun: lesen, schreiben, zuhören, sprechen und Bedeutung erfassen.

#### Typische Aufgaben der NLP

- **Text verstehen:** Worum geht es in einem Text? (z. B. Themen, Schlüsselbegriffe)
- **Spracherkennung:** Gesprochene Sprache in Text umwandeln
- **Maschinelle Übersetzung:** Texte von einer Sprache in eine andere übersetzen
- **Sentimentanalyse:** Erkennen, ob ein Text positiv, negativ oder neutral ist
- **Textgenerierung:** Automatisches Schreiben von Antworten, Zusammenfassungen oder Artikeln
- **Chatbots & Sprachassistenten:** Natürliche Dialoge mit Menschen führen

#### Wo begegnet dir NLP im Alltag?

- Sprachassistenten wie Siri oder Alexa
- Automatische Übersetzungen (z. B. Online-Übersetzer)
- Autokorrektur und Textvorschläge auf dem Smartphone
- Spam-Filter für E-Mails
- Suchmaschinen




# Automatische Zusammenfassung

**GIZMODO** + FOLLOW

Eric Limer  
Filed to: SMARTWATCHES Monday 4:31pm 175,377 1 ★

## The Best Smartwatches That Aren't the Apple Watch



## Five things the Pebble Time can do that the Apple Watch can't

**Summary:** The new Apple Watch isn't the only smartwatch to consider and if you own an iPhone then you should consider what the Pebble Time offers. Matthew lists five things to consider.

By Matthew Miller for The Mobile Gadgeteer | March 12, 2015 -- 14:25 GMT (07:25 PDT)

[Follow @palsolo](#) 8,013 followers [Get the ZDNet Microsoft newsletter now](#)

Comments 5 [Share on Facebook](#) 1 [Tweet](#) 81 [Share](#) 6 more +



## Apple Watch Has Big Drawbacks Interface, Reviews Say

Not overly positive reactions so far.

Chris Bling - Staff Reporter  
03/11/15 @11:11am in Tech 3.8K

50 [streetcred](#) [11](#) [twitter](#) [17](#) [facebook](#) [send via email](#) [share](#)



Apple's highly anticipated Apple Watch — a product developed behind a shroud of PR control and secrecy — is finally ready for prime time. And reviews of the Apple Watch are pouring in. But a number of the initial impressions are not great.

The Apple Watch has drawbacks. There are other smartwatches that offer more capabilities.

# Maschinelle Übersetzung

**"Il est impossible aux journalistes de rentrer dans les régions tibétaines"**

Bruno Philip, correspondant du "Monde" en Chine, estime que les journalistes de l'AFP qui ont été expulsés de la province tibétaine du Qinghai "n'étaient pas dans l'illégalité".

**Les faits** Le dalaï-lama dénonce l'"enfer" imposé au Tibet depuis sa fuite, en 1959

**Vidéo** Anniversaire de la rébellion tibétaine en Chine sur une vidéo





**"It is impossible for journalists to enter Tibetan areas"**

Philip Bruno, correspondent for "World" in China, said that journalists of the AFP who have been deported from the Tibetan province of Qinghai "were not illegal."

**Facts** The Dalai Lama denounces the "hell" imposed since he fled Tibet in 1959




**Video** Anniversary of the Tibetan rebellion: China on guard




Google   

[All](#) [Images](#) [Shopping](#) [Apps](#) [Videos](#) [More](#) [Search tools](#)

About 20,800,000 results (0.54 seconds)

Spanish   English 

**buenas noches** Edit **Goodnight**

 3 more translations

[Open in Google Translate](#)

# Persönliche Assistenten



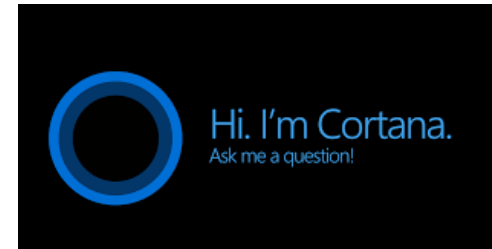
Listening..



Siri



amazon alexa



# Question answering

 amazon alexa



"Alexa, who was President when Barack Obama was nine?"

"Alexa, how's my commute?"

"Alexa, what's the weather?"

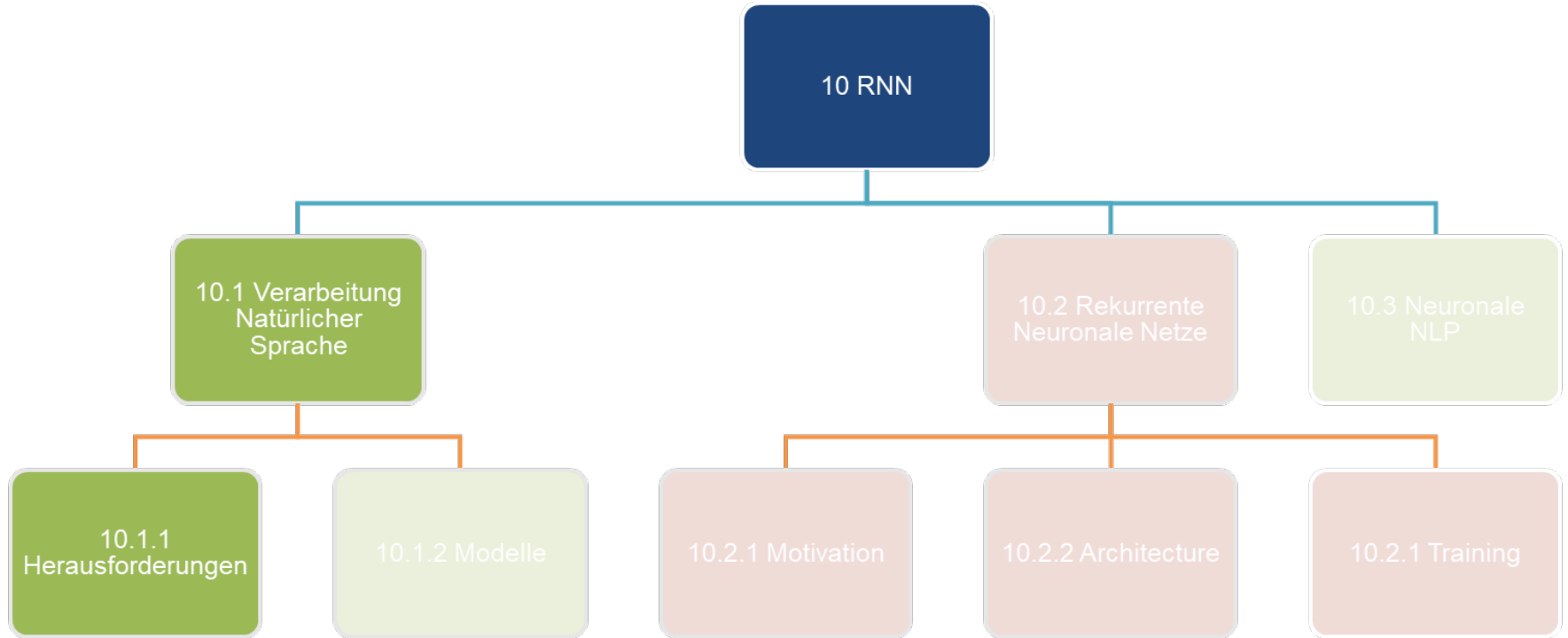
"Alexa, did the 49ers win?"



# NLP, formal

- Algorithmen zum **Verarbeiten, Analysieren und Verstehen** von Texten in natürlicher Sprache
  - **Struktur** verstehen
  - **Bedeutung** verstehen
- Schnittmenge zu vielen Bereichen
  - Informatik
    - Künstliche Intelligenz
  - Linguistik

# Überblick



# Verarbeitung natürlicher Sprache

ChatGPT 5.2 ▾

Was sind Herausforderungen bei der Verarbeitung natürlicher Sprache?

Bei der **Verarbeitung natürlicher Sprache (NLP)** gibt es mehrere zentrale Herausforderungen:

- **Mehrdeutigkeit:** Wörter und Sätze können je nach Kontext unterschiedliche Bedeutungen haben („Bank“ als Geldinstitut oder Sitzgelegenheit).
- **Kontextverständnis:** Die Bedeutung eines Satzes hängt oft vom vorherigen Text, vom Weltwissen oder von impliziten Annahmen ab.
- **Grammatik und Sprachvielfalt:** Sprachen haben komplexe Grammatikregeln, Dialekte, Umgangssprache und regionale Varianten.
- **Ironie und Sarkasmus:** Diese sind für Menschen oft leicht, für Computer aber sehr schwer zu erkennen.
- **Fehlerhafte oder unvollständige Sprache:** Tippfehler, Abkürzungen oder gesprochene Sprache mit Pausen und Füllwörtern erschweren die Verarbeitung.
- **Kulturelle und emotionale Nuancen:** Gefühle, Höflichkeit oder kulturelle Anspielungen sind schwer eindeutig zu interpretieren.
- **Datenabhängigkeit:** Gute NLP-Modelle benötigen große, qualitativ hochwertige und oft sprachspezifische Datensätze.

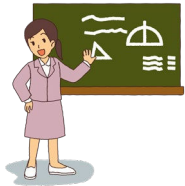
Diese Herausforderungen machen NLP zu einem komplexen, aber sehr aktiven Forschungsgebiet.

# Herausforderungen



# Herausforderungen

- Mehrdeutigkeiten



Vorzeichen



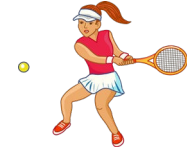
sign



omen



I am happy to introduce one of the best tennis players in the world.



# Mehrdeutigkeit

Find as many meanings as you can of this sentence:

I made her duck

- I cooked waterfowl for her benefit (to eat)
- I cooked waterfowl belonging to her
- I created the (plaster?) waterfowl she owns
- I caused her to quickly lower her head or body
- I waved my magic wand and turned her into undifferentiated waterfowl

# Mehrdeutigkeiten: Phonetik

**Aye mate, her duck**

I mate or duck

I'm eight or duck

Eye maid; her duck

I maid her duck

I'm aid her duck

I mate her duck

I'm ate her duck

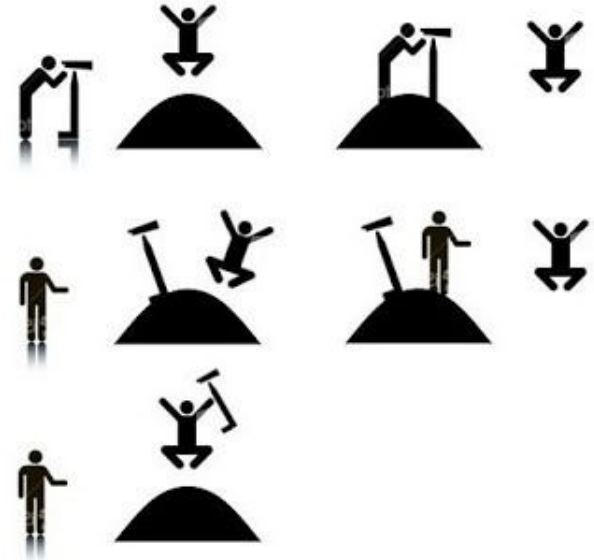
I'm ate or duck

I mate or duck



# Mehrdeutigkeit - Struktur

- “I saw the man on the hill with a telescope.”
  - I saw the man. The man was on the hill. I was using a telescope.
  - I saw the man. I was on the hill. I was using a telescope.
  - I saw the man. The man was on the hill. The hill had a telescope.
  - I saw the man. I was on the hill. The hill had a telescope.
  - I saw the man. The man was on the hill. I saw him using a telescope.



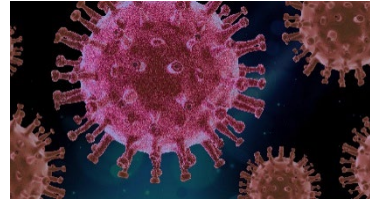
# Co-reference

„If the baby does not like the milk, boil it“



# Herausforderungen

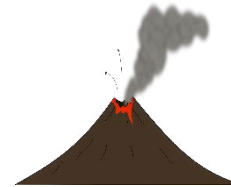
- Mehrdeutigkeiten
- Neue Wörter



COVID



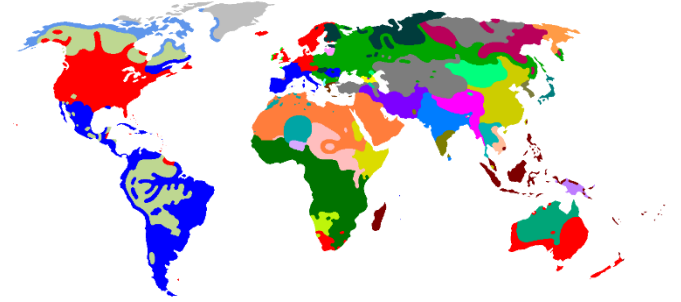
Brexit



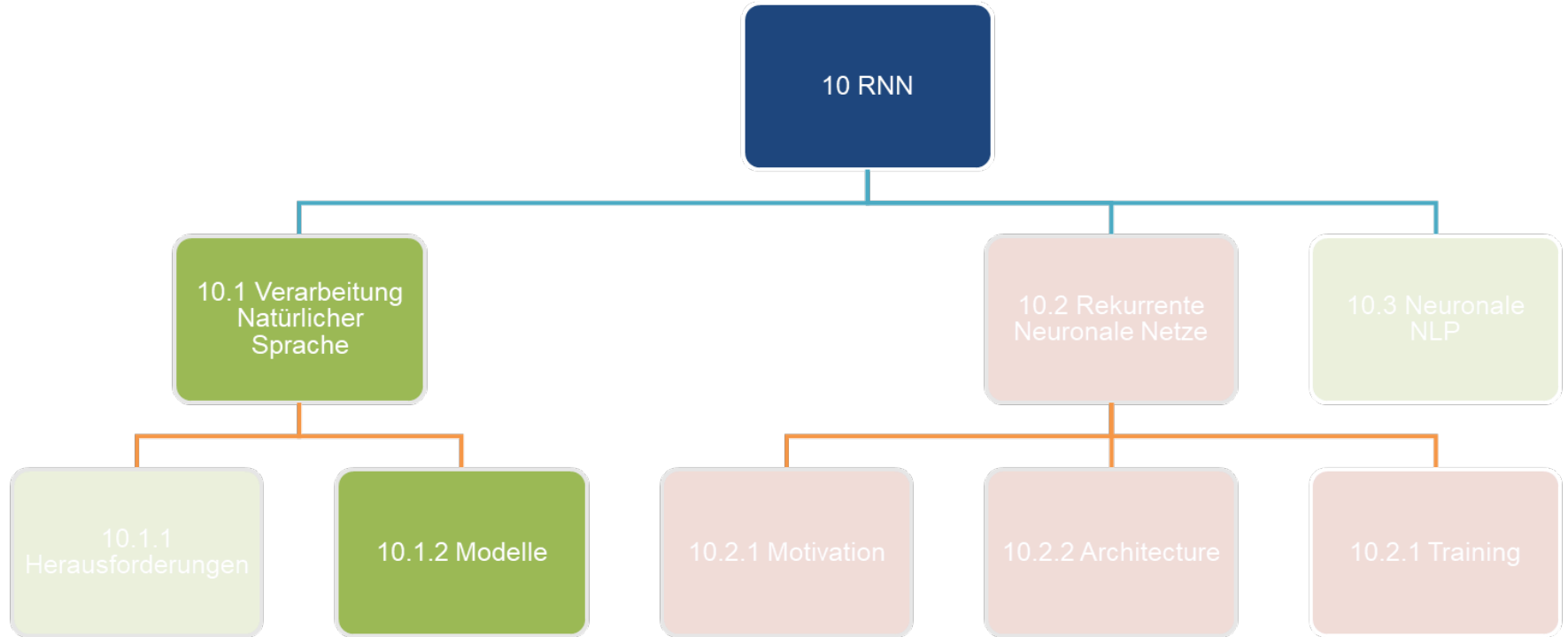
Eyjafjallajökull

# Herausforderungen

- Mehrdeutigkeiten
- Neue Wörter
- Diversität
  - 6000-7000 Sprachen
  - Modalitäten



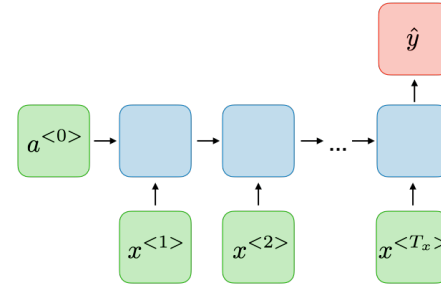
# Überblick



# Modelltypen

Sequenzklassifikation:

- Eingabe: Sequenz mit variabler Länge
- Ausgabe: Klasse



Beispiel: Sentimentanalyse

*„Ich bin mit dem Produkt äußerst zufrieden. Die Qualität ist hervorragend und die Lieferung erfolgte schneller als erwartet. Auch der Kundenservice war freundlich und sehr hilfsbereit.“*



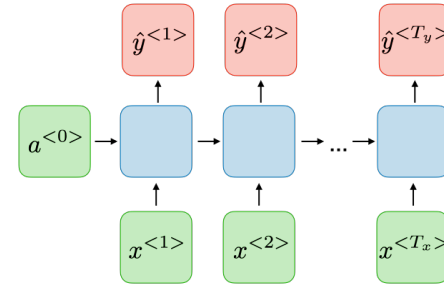
*„Man könnte meinen, die Qualität sei hochwertig, zumindest wenn man sich an den Beschreibungen orientiert. In der Praxis wirkt das Produkt jedoch alles andere als überzeugend, und die versprochene Zuverlässigkeit habe ich persönlich nicht wirklich erlebt. Wer einen reibungslosen Ablauf und hilfreichen Support erwartet, wird vermutlich ähnliche Erfahrungen machen.“*



# Modelltypen

Sequenzannotation:

- Eingang: Sequenz mit variabler Länge
- Ausgabe: Labelsequenz mit gleicher Länge



Beispiel: Erkennung von Entitäten (Named Entity Recognition) in

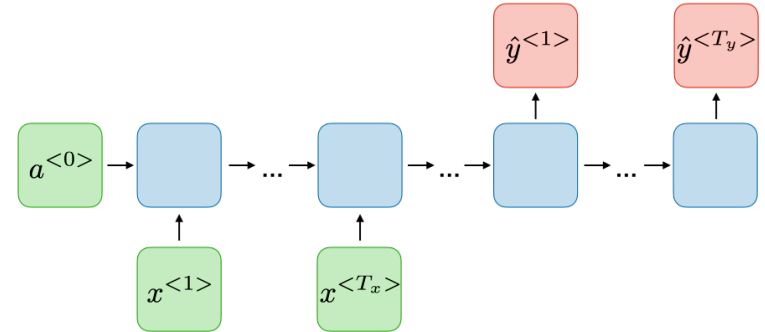
„Ich fliege mit American Airlines nach New York“

ich	fliege	mit	American	Airlines	nach	New	York

# Modelltypen

## Sequenz-zu-Sequenz-Modell:

- Eingabe: Sequenz mit variabler Länge
- Ausgabe: Sequenz mit variabler Länge



## Beispiel: Automatische Textzusammenfassung

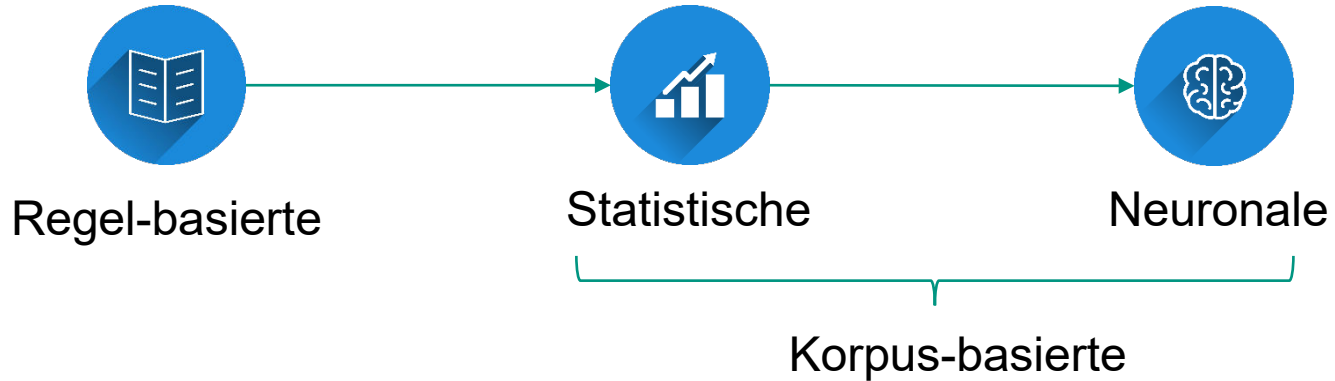
### Good quality summary output

**S:** a man charged with the murder last year of a british backpacker confessed to the slaying on the night he was charged with her killing , according to police evidence presented at a court hearing tuesday . ian douglas previte , ## , is charged with murdering caroline stuttle , ## , of yorkshire , england

**T:** man charged with british backpacker 's death confessed to crime police officer claims

**O:** man charged with murdering british backpacker confessed to murder

# Ansätze



**Experten erstellen Regeln**

**Lernen aus Daten**

# Überblick



# Kurze Pause: Lehrevaluation (Vorlesung – Nowack)



Vielen herzlichen Dank für Ihre Zeit  
und Ihr konstruktives Feedback!

↙  
Sachlich, spezifisch, begründet,  
lösungsorientiert:

- Was war gut und warum?
- Was könnte verbessert werden und wie?
- Ihre Kommentare und Vorschläge sind mir sehr willkommen!

<https://onlineumfrage.kit.edu/evasys/online.php?p=W28CQ>

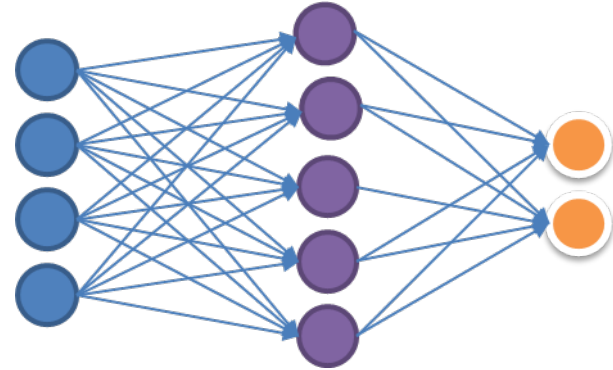
# Frühere Vorlesungen

## Mehrschichtiges Perzeptron (MLP)

- Annäherung von Funktionen: Zuordnen von Eingaben  $\mathbf{X}$  zu Ausgaben  $\mathbf{Y}$  mit Parametern  $\theta$

$$\mathbf{Y} = f(\mathbf{X}; \theta)$$

- Feste Anzahl von Eingabedimensionen  
→ ein Problem bei der Verarbeitung von Sequenzen unterschiedlicher Länge



# Recurrent Neural Networks (RNNs)

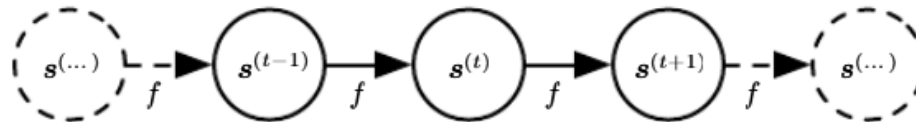
Eine Familie von neuronalen Netzwerken die spezialisiert darauf sind Sequenzdaten  $x_1, \dots, x_\tau$  von oftmals variabler Länge zu prozessieren.

Dafür wird wieder die Idee von „Parameter sharing“ zentral sein.

Notation:  $x_t$  beschreibt Vektoren, die die Sequenz bilden, mit dem Zeitschrittindex  $t$ . Dieser kann, muss sich aber nicht, auf eine wahre Zeitachse beziehen: „Die Position in der Sequenz.“

Zentral: Die Idee einer rekursiven/rekurrenten Berechnung eines dynamischen Systems mit Zustandsvektor  $s_t$

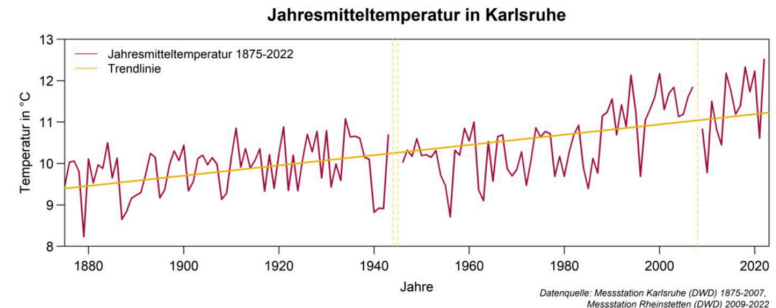
$$s_t = f(s_{t-1}; \theta)$$



# Sequenzen

Viele verschiedene Arten von Sequenzdaten

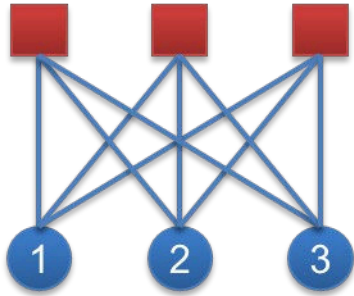
- Text
- Audio
- Videos
- Sensordaten
- Wetterdaten
- Finanzdaten (Aktienkurse)
- ...



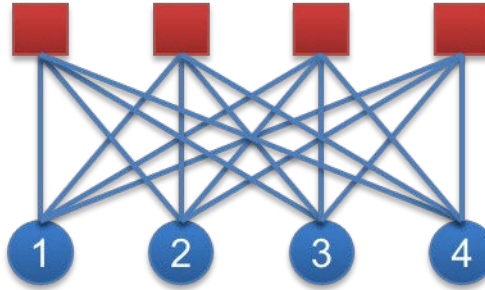
Stadt Karlsruhe

# Sequenzen

Herausforderung: Eingaben mit variabler Länge



I go home



My friend goes home

MLP: lerne alle Verbindungen über alle möglichen Positionen im Satz

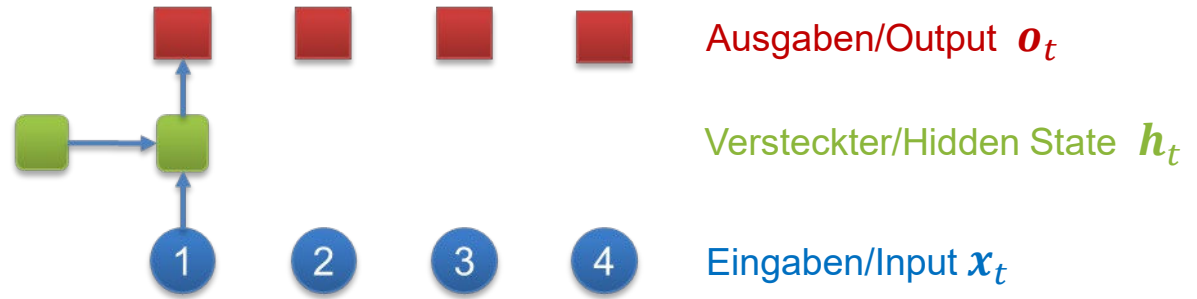
Herausforderung: Information unabhängig von der Position im Satz extrahieren:

*„I went to Nepal in 2009.“*

*„In 2009, I went to Nepal.“*

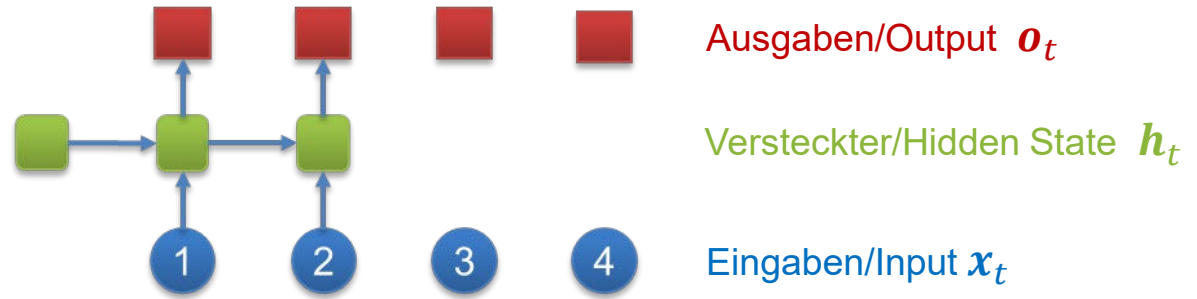
# Neuronale Netze für Sequenzen

Idee: Speichern des Verlaufs in einem Zustand (Vektor)



# Neuronale Netze für Sequenzen

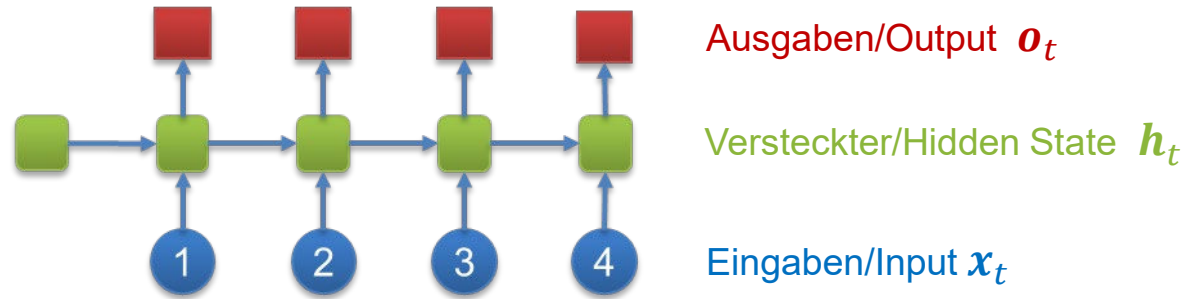
Idee: Speichern des Verlaufs in einem Zustand (Vektor) – Reihenfolge wichtig!



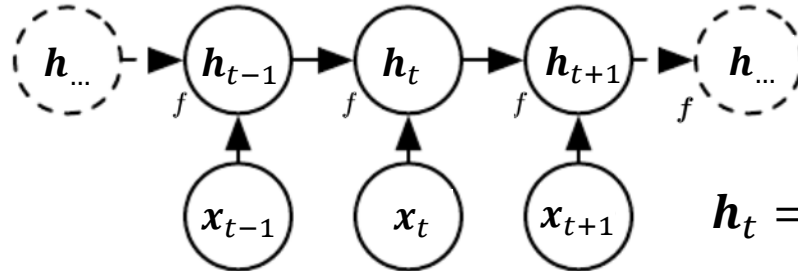
“Mary kills John” ist anders als “John kills Mary”

# Neuronale Netze für Sequenzen

Idee: Speichern des Verlaufs in einem Zustand (Vektor)  $h_t$

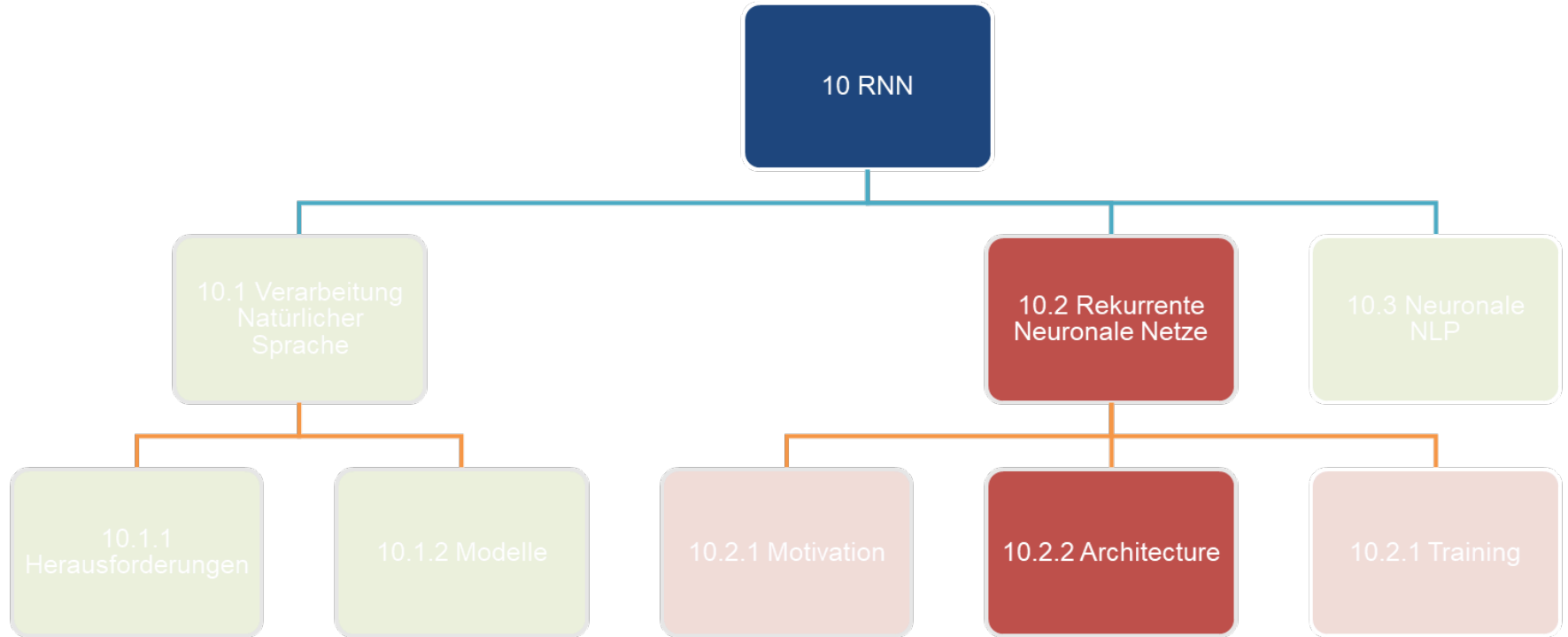


$$s_t = f(s_{t-1}; \theta)$$



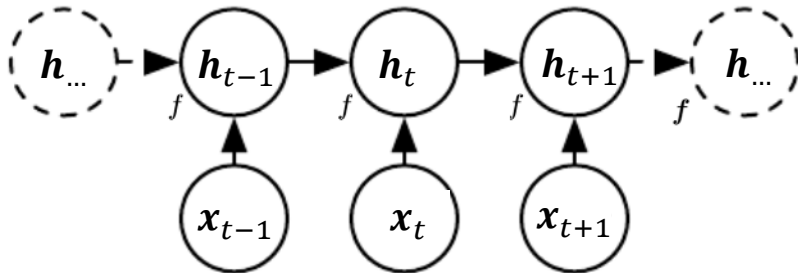
$$h_t = f(h_{t-1}, x_t; \theta)$$

# Überblick



# Was ist die Architektur solcher RNNs?

- Effektiv könnte jede neuronale Funktion  $f$  die einer rekurrenten Beschreibung wie z.B.  $\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t; \boldsymbol{\theta})$  folgt als RNN betrachtet werden...
- Daher viele Optionen!
- Es gibt aber z.B. Abweichungen darin wie der hidden state basierend auf den Eingaben und Ausgaben definiert wird.



Gemeinsame Vorteile:

- Unabhängig von Sequenzlänge
- Können die gleiche Übergangsfunktion  $f$  für jeden Zeitschritt nutzen

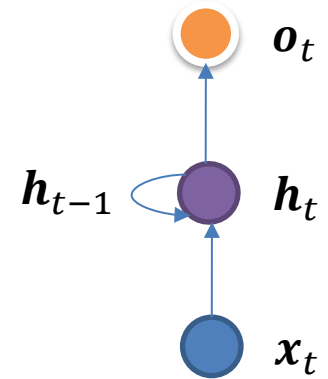
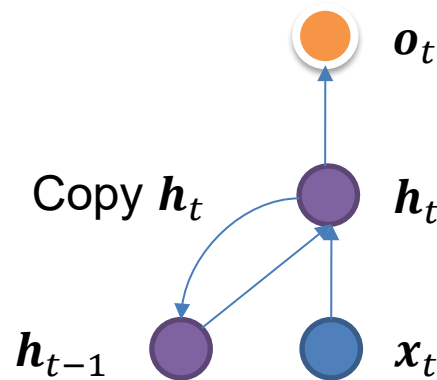
# Elman-Netzwerke

Architektur:

- Vorheriger hidden state wird zur Eingabe für den nächsten Zeitschritt

Implementierung:

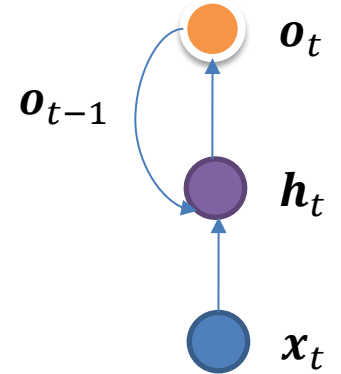
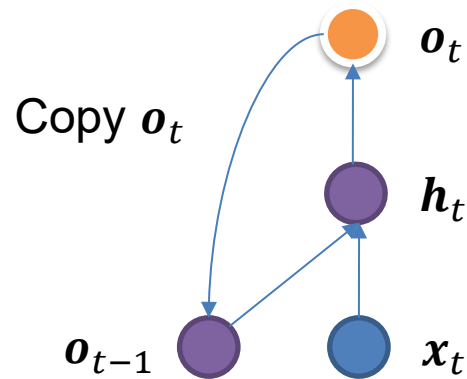
- Kopiere den hidden state des vorherigen Zeitschritts



# Jordan-Netzwerke

Gleicher Kopiermechanismus

- Aber für die Ausgabeschicht





# Wie werden RNNs in der Praxis implementiert?

Berechnung der versteckten Schicht:  $\mathbf{h}_t = f(\mathbf{x}_t, \mathbf{h}_{t-1}, \boldsymbol{\theta})$

Viele verschiedene Implementierungen/Arten von RNNs, z.B.:

- **Einfache rekurrente neuronale Netzwerke**

Rumelhart et al. 1986

- **Long Short-Term Memory (LSTM)**

Hochreiter und Schmidhuber 1997; adressiert Vanishing Gradients Problem

- **Gated Recurrent Unit (GRU)**

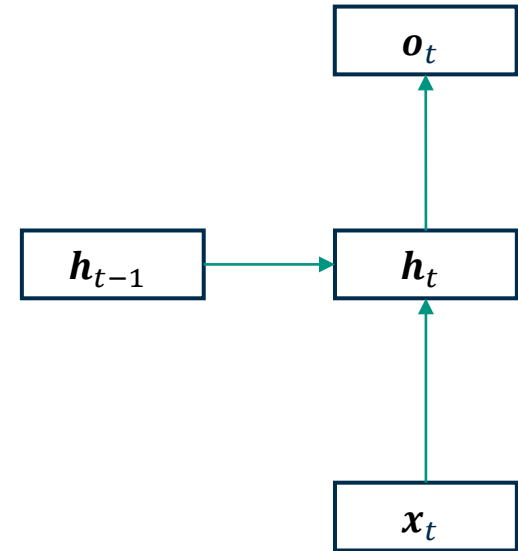
Cho et al. 2014; recheneffizientere Alternative

# Einfache rekurrente neuronale Netze

$$\mathbf{h}_t = f_{\text{act}}(\mathbf{W}^h \mathbf{h}_{t-1} + \mathbf{W}^x \mathbf{x}_t + \mathbf{b})$$

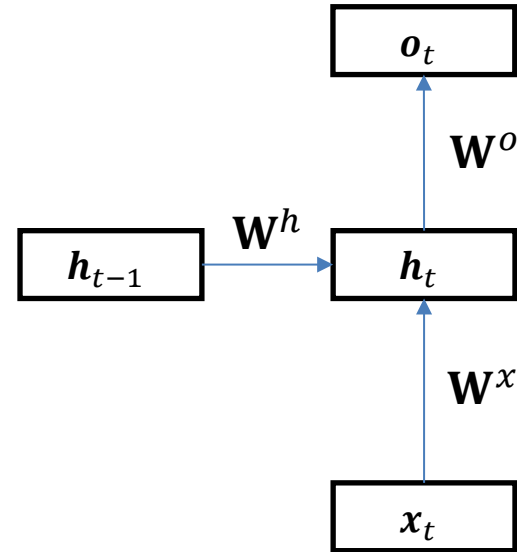
$f_{\text{act}}$ : eine Aktivierungsfunktion (Sigmoid, Tanh, ReLU)

Innerhalb der Klammer: eine lineare Kombination zwischen dem Speicher  $\mathbf{h}_{t-1}$  und der Eingabe  $\mathbf{x}_t$

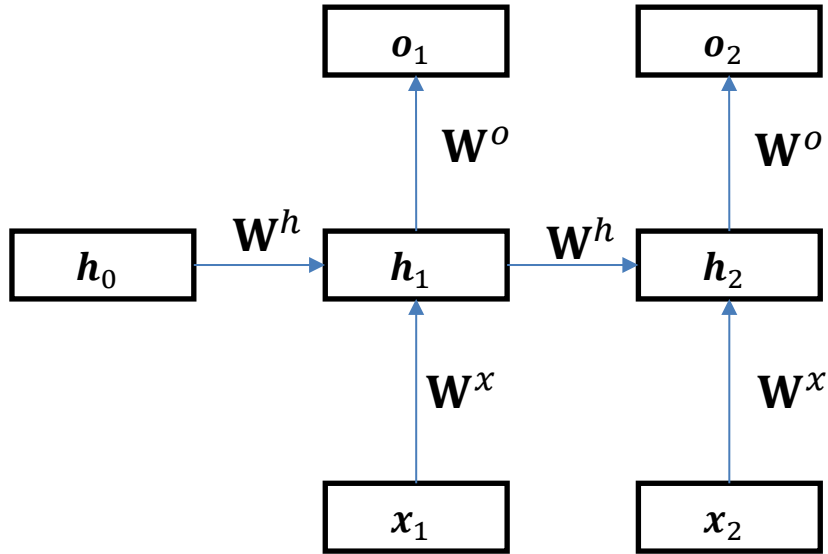


# Ausrollen/Roll-out

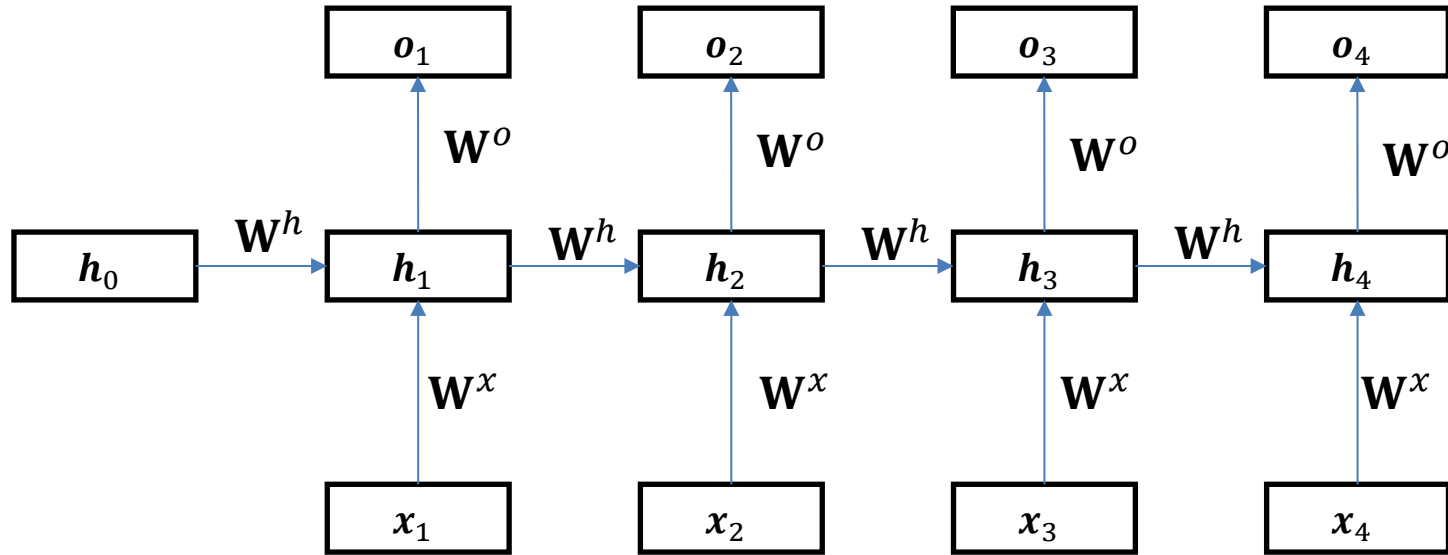
- Ausrollen des Netzwerks über die Zeit  
→ **geteilte Gewichte!**
- D.h. die gleichen drei Gewichtsmatrizen  $W^h$ ,  $W^x$  und (für die Ausgabe)  $W^o$  werden zu jedem Zeitschritt angewendet
- Damit haben wir ein Netzwerk, welches **unabhängig von der Sequenzlänge** eingesetzt werden kann
- Einfaches RNN: Feedforward-Netzwerk für jede neue Eingabe  $x_t$ , zudem basierend auf dem vorherigen hidden state  $h_{t-1}$ .



# Ausrollen/Roll-out



# Ausrollen/Roll-out

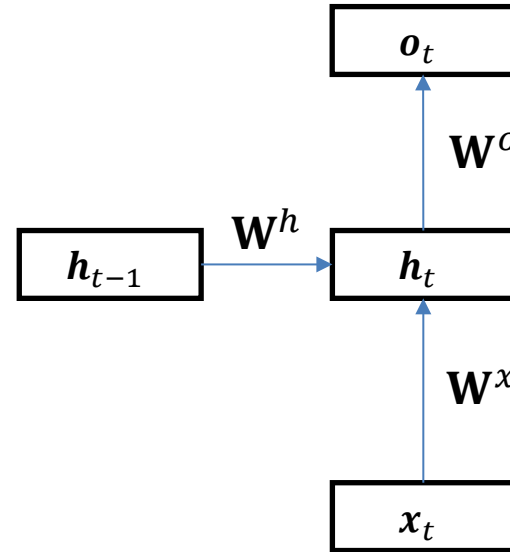


# Beispiel

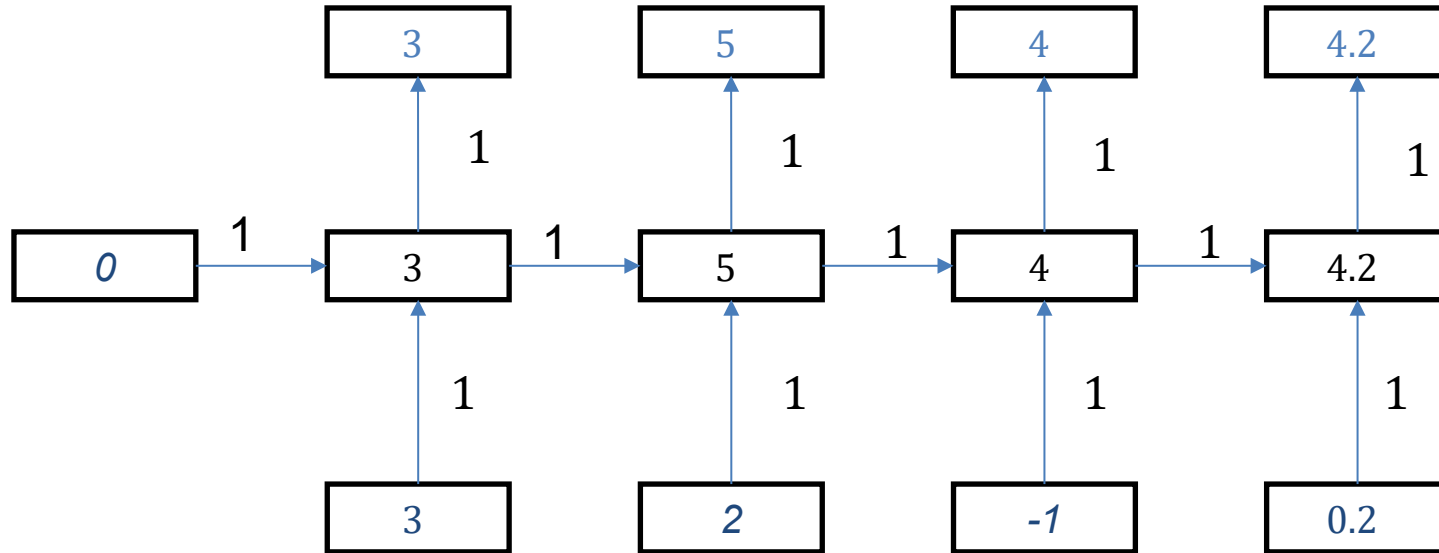
$$\mathbf{o}_t = \mathbf{W}^o \mathbf{h}_t$$

$$\mathbf{h}_t = \mathbf{W}^h \mathbf{h}_{t-1} + \mathbf{W}^x \mathbf{x}_t$$

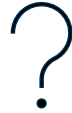
Frage: Wann wird das Output einfach zur Summe der Eingabesequenz?



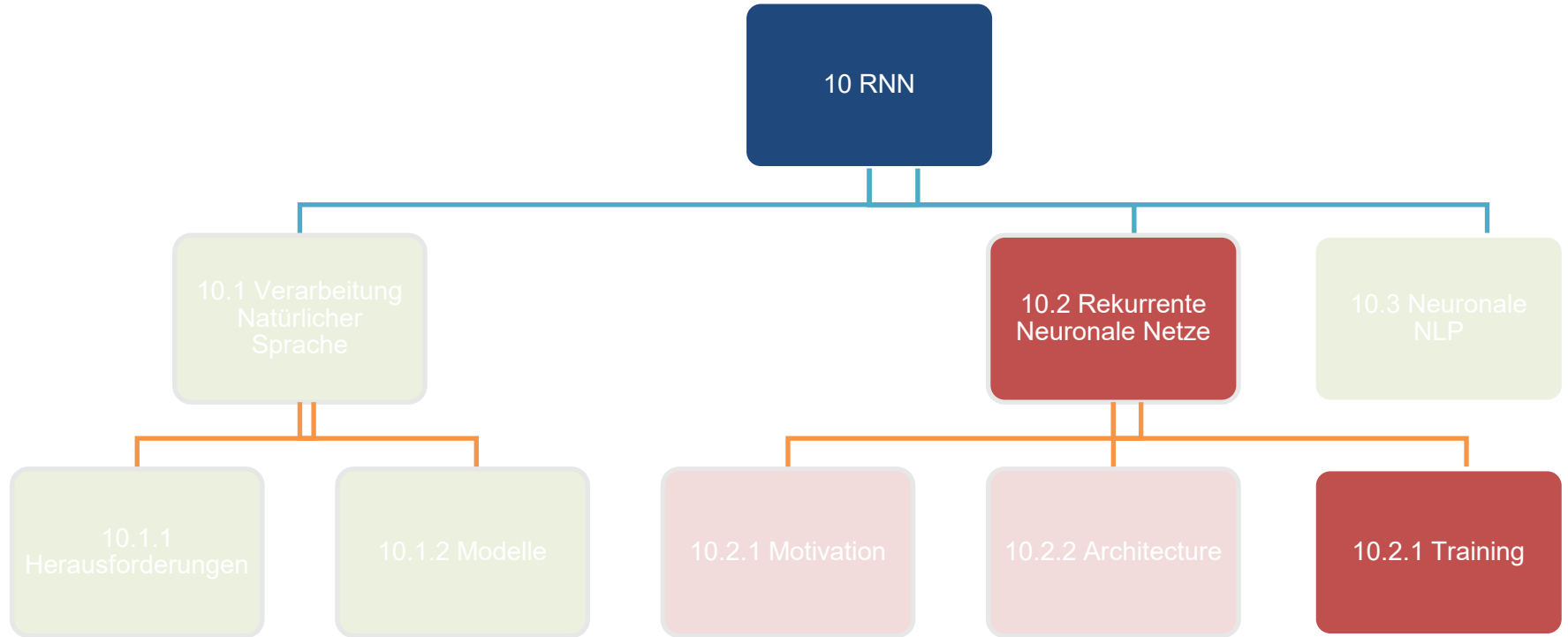
# Beispiel



# Fragen?



# Überblick



# Training? → Backpropagation durch die Zeit

Werbos (1990) Backpropagation through time: what it does and how to do it

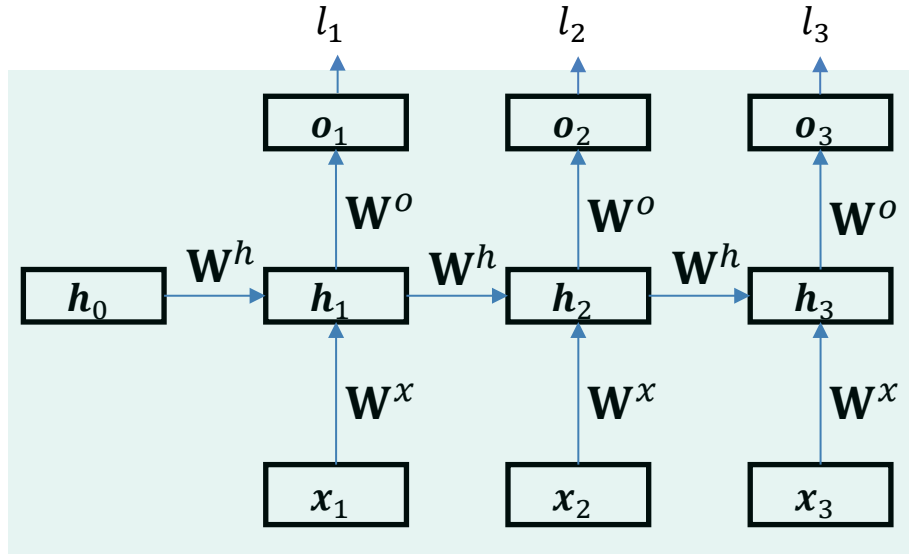
- **Backpropagation through time (BPTT)**
  - Für neuronale Netzwerke in denen der versteckte Zustand  $h_t$  rekurrent propagiert wird.
  - Gegeben: geordnete (→ Kausalität) Sequenz von Input-Output Paaren  $(x_t, y_t)$  sowie Anfangszustand (oft Nullvektor)
- Anfang: **Forward-Pass** um Aktivierungen zu berechnen.
- Dann: **Backpropagation** um die Gradienten der Loss-Funktion bezüglich allen Netzwerkparametern zu finden; **über die verschiedenen Zeitschritte hinweg**.
- Je nach RNN-Architektur: berechne Loss abhängig von Vorhersagen zu jedem Zeitschritt, oder nur am Ende, oder ...
- Hier: wegen **parameter sharing**, werden die gleichen Gewichtsmatrizen  $W^x, W^h, W^o$  zu jedem Zeitschritt genutzt. Daher werden wir auch die **Updates der Gewichte über alle Zeitschritte hinweg mitteln/summieren/gewichten**.

# Backpropagation durch die Zeit

$$\mathbf{h}_t = f_{\text{act}}(\mathbf{W}^x \mathbf{x}_t + \mathbf{W}^h \mathbf{h}_{t-1} + \mathbf{b}_h)$$

$$\mathbf{o}_t = f_{\text{out}}(\mathbf{W}^o \mathbf{h}_t + \mathbf{b}_o)$$

$$L = \frac{1}{\tau} \sum_{t=1}^{\tau} l_t(\mathbf{o}_t, \mathbf{y}_t)$$



## Für Gradient Descent

Gradienten bezüglich den drei Gewichtsmatrizen an jedem  $t$  mitteln:

$$\mathbf{W}^o = \mathbf{W}^o - \alpha \frac{\partial L}{\partial \mathbf{W}^o}$$

$$\frac{\partial L}{\partial \mathbf{W}^o} = \frac{1}{\tau} \sum_{t=1}^{\tau} \frac{\partial l_t}{\partial \mathbf{W}^o}$$

$$\mathbf{W}^h = \mathbf{W}^h - \alpha \frac{\partial L}{\partial \mathbf{W}^h}$$

$$\frac{\partial L}{\partial \mathbf{W}^h} = \frac{1}{\tau} \sum_{t=1}^{\tau} \frac{\partial l_t}{\partial \mathbf{W}^h}$$

$$\mathbf{W}^x = \mathbf{W}^x - \alpha \frac{\partial L}{\partial \mathbf{W}^x}$$

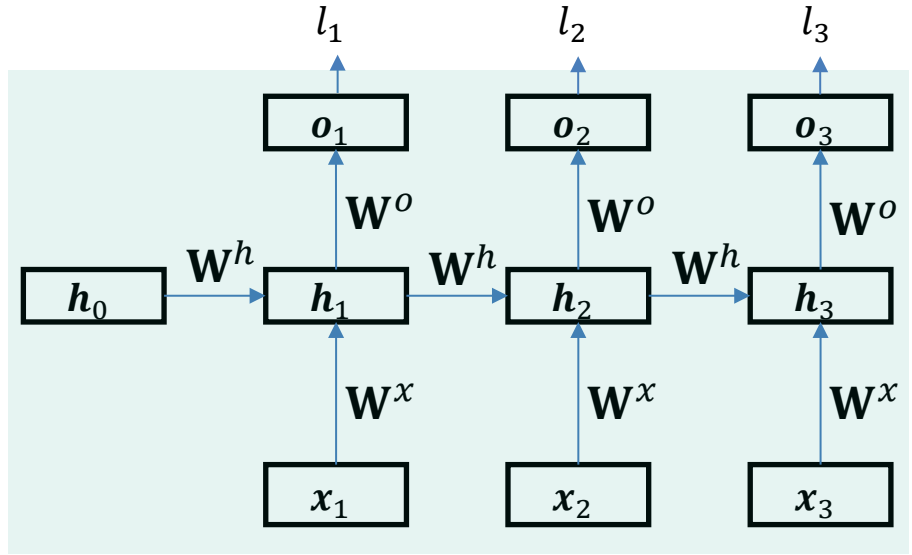
$$\frac{\partial L}{\partial \mathbf{W}^x} = \frac{1}{\tau} \sum_{t=1}^{\tau} \frac{\partial l_t}{\partial \mathbf{W}^x}$$

# Backpropagation durch die Zeit

$$\mathbf{h}_t = f_{\text{act}}(\mathbf{W}^x \mathbf{x}_t + \mathbf{W}^h \mathbf{h}_{t-1} + \mathbf{b}_h)$$

$$\mathbf{o}_t = f_{\text{out}}(\mathbf{W}^o \mathbf{h}_t + \mathbf{b}_o)$$

$$L = \frac{1}{\tau} \sum_{t=1}^{\tau} l_t(\mathbf{o}_t, \mathbf{y}_t)$$



Nutze den ausgerollten Graphen, um alle Abhängigkeiten zu sehen, die in der Berechnung der Gradienten berücksichtigt werden müssen.

Zum Beispiel für  $\mathbf{h}_3$  sind die direkten Abhängigkeiten:  $\mathbf{W}^x, \mathbf{W}^h, \mathbf{h}_2, \mathbf{x}_3$

Ableitungen des Losses hinsichtlich der drei Gewichtsmatrizen:  $\frac{\partial L}{\partial \mathbf{W}^x}, \frac{\partial L}{\partial \mathbf{W}^h}, \frac{\partial L}{\partial \mathbf{W}^o}$

# Backpropagation durch die Zeit

$$\mathbf{h}_t = \mathbf{W}^x \mathbf{x}_t + \mathbf{W}^h \mathbf{h}_{t-1}$$

$$\mathbf{o}_t = \mathbf{W}^o \mathbf{h}_t$$

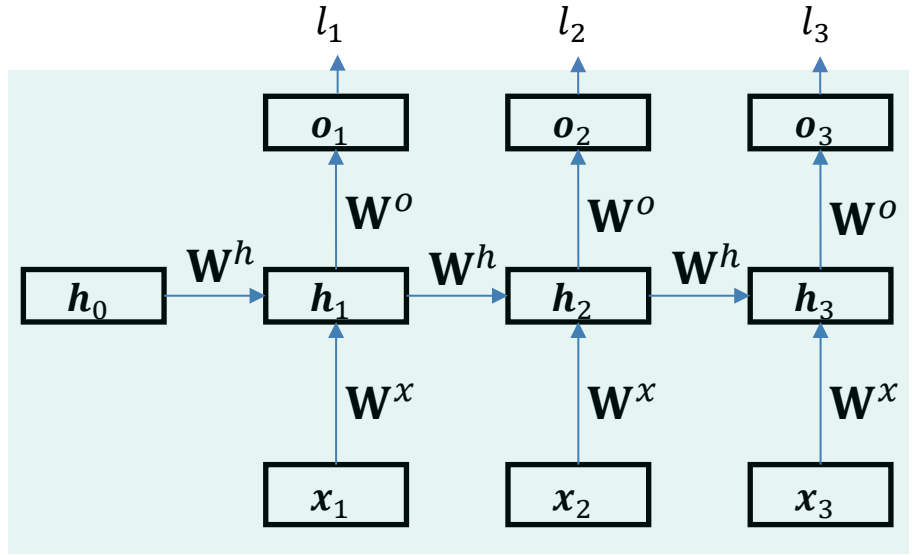
$$L = \frac{1}{\tau} \sum_{t=1}^{\tau} l_t(\mathbf{o}_t, \mathbf{y}_t)$$

Fehler hinsichtlich Vorhersage:

$$\mathbf{g}_t = \frac{\partial l_t}{\partial \mathbf{o}_t} \text{ und somit } \frac{\partial L}{\partial \mathbf{o}_t} = \frac{1}{\tau} \mathbf{g}_t$$

## Annahmen:

- $f_{\text{act}}, f_{\text{out}}$  sind Einheitsfunktionen
- keine  $\mathbf{b}$  Terme
- $\mathbf{h}_0$  wird als konstant, z.B. Nullvektor angenommen
- $\mathbf{o}_3$  ist das letzte Output

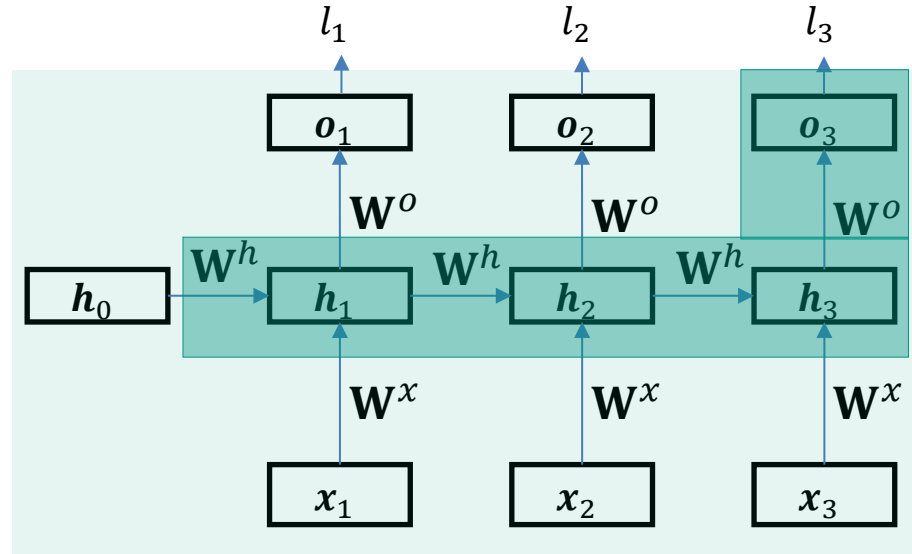


# Backpropagation durch die Zeit

$$\mathbf{h}_t = \mathbf{W}^x \mathbf{x}_t + \mathbf{W}^h \mathbf{h}_{t-1}$$

$$\mathbf{o}_t = \mathbf{W}^o \mathbf{h}_t$$

$$L = \frac{1}{\tau} \sum_{t=1}^{\tau} l_t(\mathbf{o}_t, \mathbf{y}_t)$$



Gradienten hinsichtlich  $\mathbf{W}^o$  und  $\mathbf{W}^h$ :

$$\frac{\partial l_3}{\partial \mathbf{W}^o} = \frac{\partial l_3}{\partial \mathbf{o}_3} \frac{\partial \mathbf{o}_3}{\partial \mathbf{W}^o} = \frac{\partial l_3}{\partial \mathbf{o}_3} \mathbf{h}_3^T$$

$\mathbf{h}_0$  Nullvektor

$$\frac{\partial l_3}{\partial \mathbf{W}^h} = \frac{\partial l_3}{\partial \mathbf{o}_3} \frac{\partial \mathbf{o}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{W}^h} + \frac{\partial l_3}{\partial \mathbf{o}_3} \frac{\partial \mathbf{o}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{W}^h} + \frac{\partial l_3}{\partial \mathbf{o}_3} \frac{\partial \mathbf{o}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{h}_1} \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}^h} + \frac{\partial l_3}{\partial \mathbf{o}_3} \frac{\partial \mathbf{o}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{h}_1} \frac{\partial \mathbf{h}_1}{\partial \mathbf{h}_0} \frac{\partial \mathbf{h}_0}{\partial \mathbf{W}^h}$$

$\mathbf{W}^o$   $\mathbf{h}_2^T$

$$= \sum_{t=1}^3 \frac{\partial l_3}{\partial \mathbf{o}_3} \frac{\partial \mathbf{o}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}^h}$$

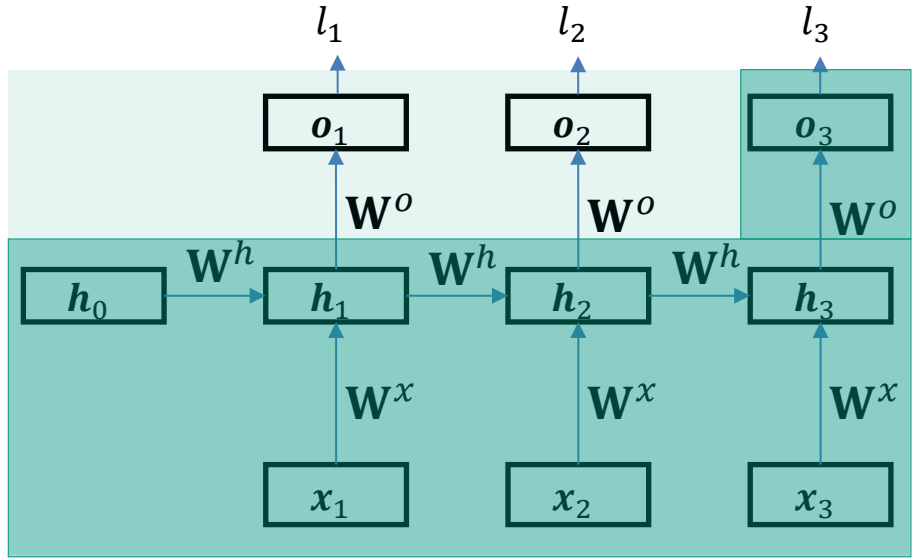
Nutze Kausalität in der Zeit um immer von fixen Quantitäten auszugehen!

# Backpropagation durch die Zeit

$$\mathbf{h}_t = \mathbf{W}^x \mathbf{x}_t + \mathbf{W}^h \mathbf{h}_{t-1}$$

$$\mathbf{o}_t = \mathbf{W}^o \mathbf{h}_t$$

$$L = \frac{1}{\tau} \sum_{t=1}^{\tau} l_t(\mathbf{o}_t, \mathbf{y}_t)$$



Gradienten hinsichtlich  $\mathbf{W}^x$  (läuft wieder über  $\mathbf{h}_t$ , daher sehr ähnlich)

$$\frac{\partial l_3}{\partial \mathbf{W}^x} = \frac{\partial l_3}{\partial \mathbf{o}_3} \frac{\partial \mathbf{o}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{W}^x} + \frac{\partial l_3}{\partial \mathbf{o}_3} \frac{\partial \mathbf{o}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{W}^x} + \frac{\partial l_3}{\partial \mathbf{o}_3} \frac{\partial \mathbf{o}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{h}_1} \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}^x}$$

$$\begin{matrix} \nearrow \mathbf{W}^o \\ \nearrow \mathbf{x}_3^T \end{matrix} = \sum_{t=1}^3 \frac{\partial l_3}{\partial \mathbf{o}_3} \frac{\partial \mathbf{o}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}^x}$$

Analog für andere Outputs, bias-Terme, eventuell  $\mathbf{h}_0$  Parameter...

# Probleme...

- Vanishing gradient problem / Exploding gradient problem  
→ durch wiederholte Multiplikation verschwinden/explodieren Gradienten

Illustrative Annahme:  $\mathbf{h}_t = \mathbf{W}^T \mathbf{h}_{t-1}$

Dann ist die Rekurrenz effektiv auf einfache Art und Weise beschrieben durch

$$\mathbf{h}_t = (\mathbf{W}^t)^T \mathbf{h}_0$$

Falls  $\mathbf{W}$  eine Eigenwertzerlegung  $\mathbf{W} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$  zulässt, erhalten wir

$$\mathbf{h}_t = \mathbf{V}^T \mathbf{\Lambda}^t \mathbf{V} \mathbf{h}_0$$

|Eigenwerte| < 1 „verschwinden“ und „explodieren“ andernfalls...

- Lösungsansätze: vor allem angepasste Architekturen, z.B. das LSTM

# Fragen?



You cannot vote anymore



Gegeben sei ein rekurrentes neuronales Netz mit geteilten Gewichten  $W^h$ , einer Ausgabeloss  $l_\tau$  am letzten Zeitschritt  $\tau$  und Hidden States  $h_t$ . Welche der folgenden Gleichungen beschreibt korrekt den Gradienten des Losses bezüglich der rekurrenten Gewichte  $W^h$ ?



Click on the projected screen to start the question

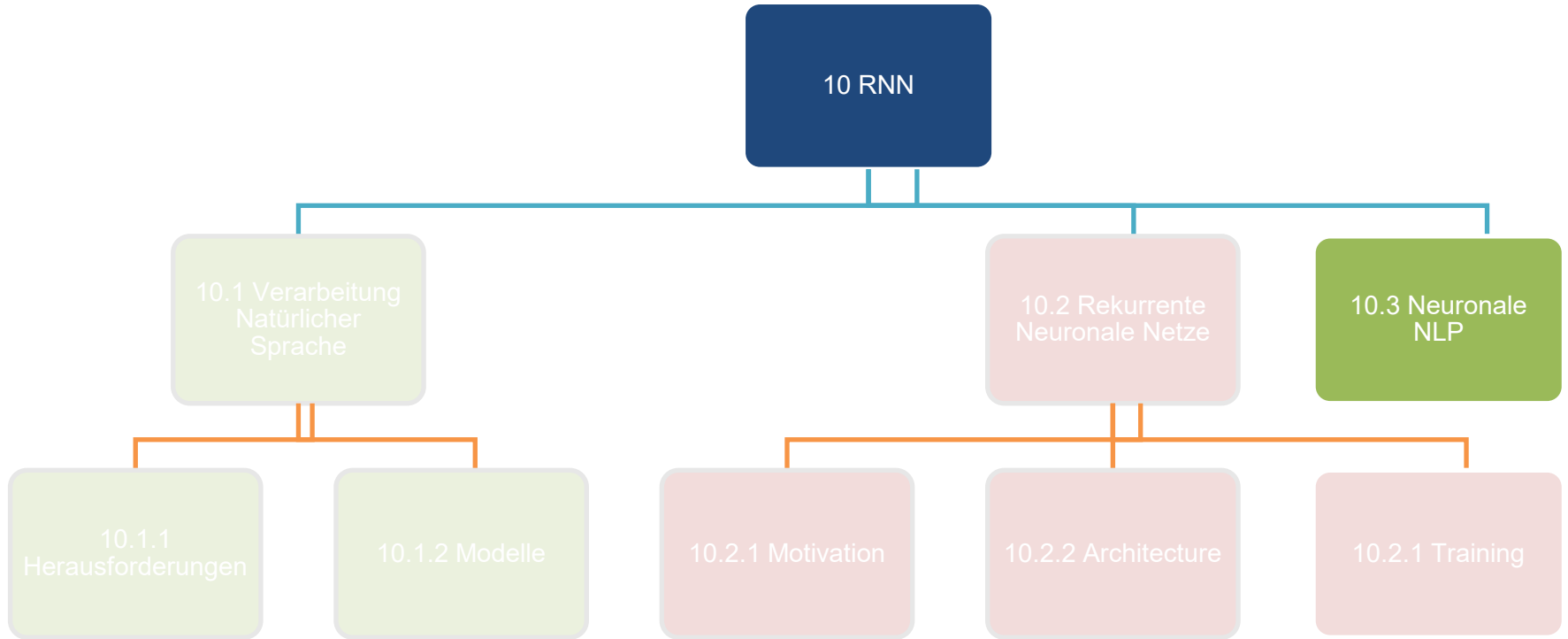
1  $\frac{\partial l_\tau}{\partial W^h} = \frac{\partial l_\tau}{\partial o_\tau} \frac{\partial o_\tau}{\partial h_\tau} \frac{\partial h_\tau}{\partial W^h}$  8% 6

2  $\frac{\partial l_\tau}{\partial W^h} = \sum_{t=1}^{\tau} \frac{\partial l_\tau}{\partial o_\tau} \frac{\partial o_\tau}{\partial h_\tau} \frac{\partial h_\tau}{\partial h_t} \frac{\partial h_t}{\partial W^h}$  53% 40 ✓

3  $\frac{\partial l_\tau}{\partial W^h} = \sum_{t=1}^{\tau} \frac{\partial l_\tau}{\partial h_t} \frac{\partial h_t}{\partial W^h}$  39% 29



# Überblick



# Modelltypen

Sequenzklassifikation:

- Eingabe: Sequenz mit variabler Länge
- Ausgabe: Klasse

Beispiel: Sentimentanalyse

*„Ich bin mit dem Produkt äußerst zufrieden. Die Qualität ist hervorragend und die Lieferung erfolgte schneller als erwartet. Auch der Kundenservice war freundlich und sehr hilfsbereit.“*



*„Man könnte meinen, die Qualität sei hochwertig, zumindest wenn man sich an den Beschreibungen orientiert. In der Praxis wirkt das Produkt jedoch alles andere als überzeugend, und die versprochene Zuverlässigkeit habe ich persönlich nicht wirklich erlebt. Wer einen reibungslosen Ablauf und hilfreichen Support erwartet, wird vermutlich ähnliche Erfahrungen machen.“*



# Grundlegendes Deep Learning Modell für NLP

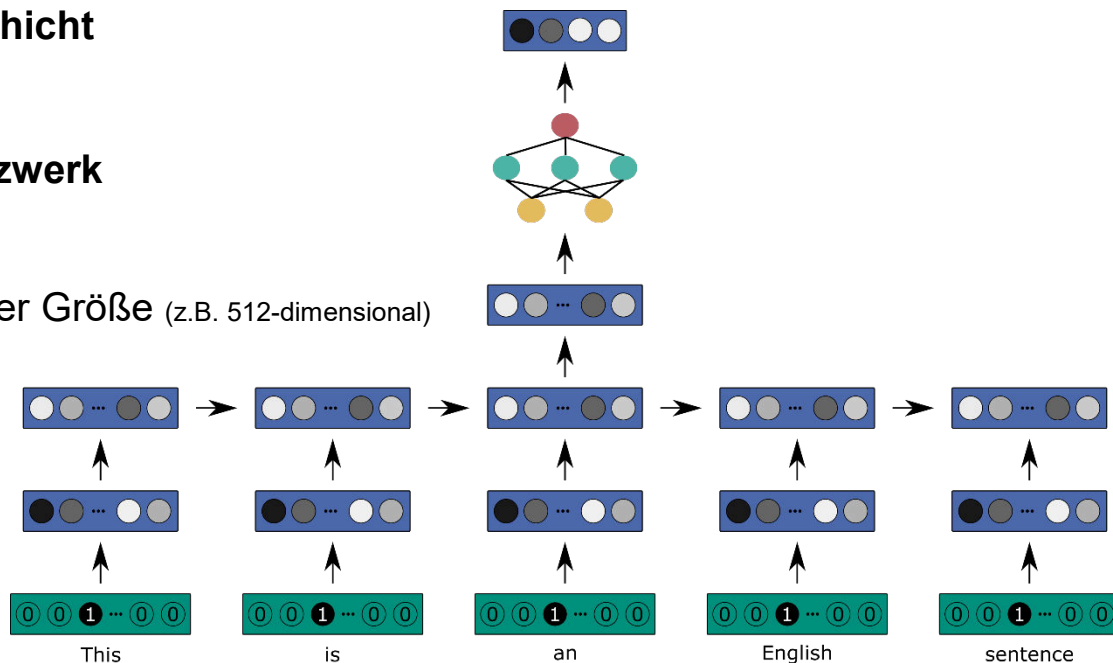
Klassifizierungsschicht

Feed-Forward-Netzwerk

Repräsentation fester Größe (z.B. 512-dimensional)

Encoder

Satz als Eingabe



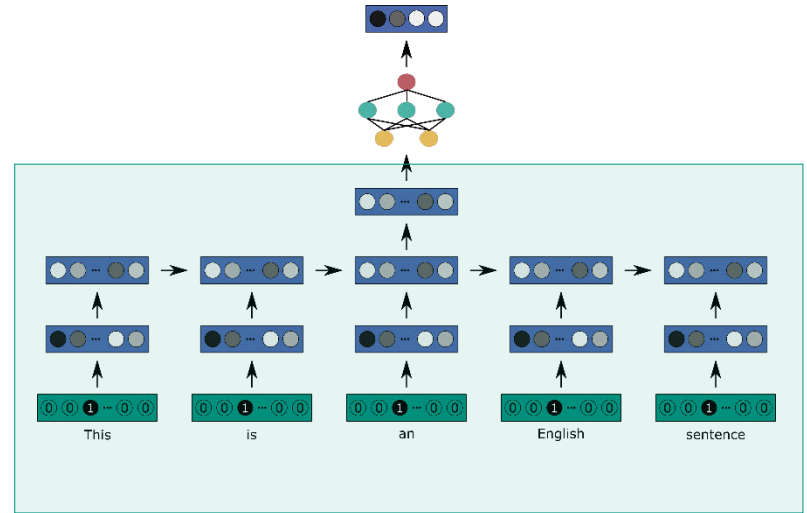
# Grundlegendes Deep Learning Modell für NLP

## Encoder:

Hier: eine Funktion zur Darstellung einer Wortfolge als Vektor fester Größe.

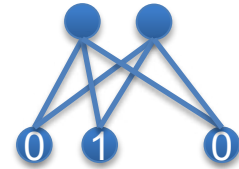
Teile eines solchen Textencoders:

- Word-Embeddings
- Sequenzschicht
- Aggregationsschicht



# Word embeddings

Erst **Zuordnung** von einem **Wort zu** einem kontinuierlichen **Vektor**



- Z.B. **one-hot-encoding Vektoren** von Wörtern/Tokens, bei dem das  $i$ -te Wort im Vokabular mit einem 1-bit in der  $i$ -ten Position und einer 0 in allen anderen Positionen encodiert ist. Besser: Zuordnung mittels  $V$  Integer IDs.
- Sehr ineffizient und kann auch Ähnlichkeiten zwischen Wörtern nicht darstellen.
- Lernen **niedrigdimensionale Vektoren** die jedes Wort repräsentieren (ähnliche Worte haben dann ähnliche Vektoren).
- Haben eine Sequenz von Wörtern in eine Sequenz von Vektoren gemapped...

## Herausforderungen:

- Die meisten Parameter
- Festes Vokabular

“aardvark” =  $[-0.7, +0.2, -3.2, \dots]$

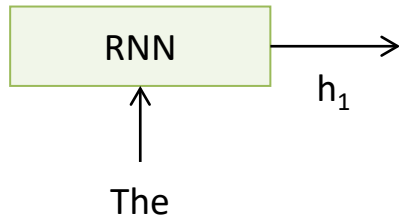
“abacus” =  $[+0.5, +0.9, -1.3, \dots]$

...

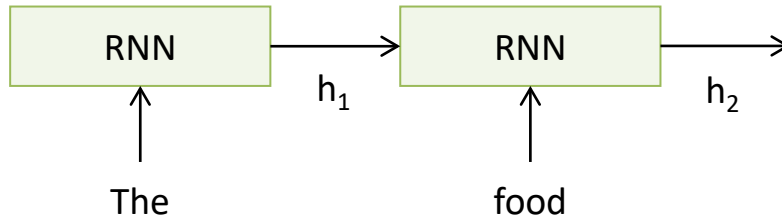
“zyzzyva” =  $[-0.1, +0.8, -0.4, \dots]$ .

Russel & Norvig, 4<sup>th</sup> edition, 2021

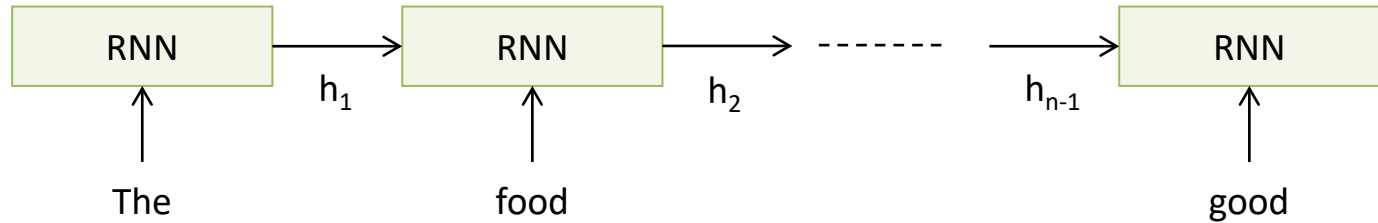
# Sequenzschicht (z.B. unser RNN)



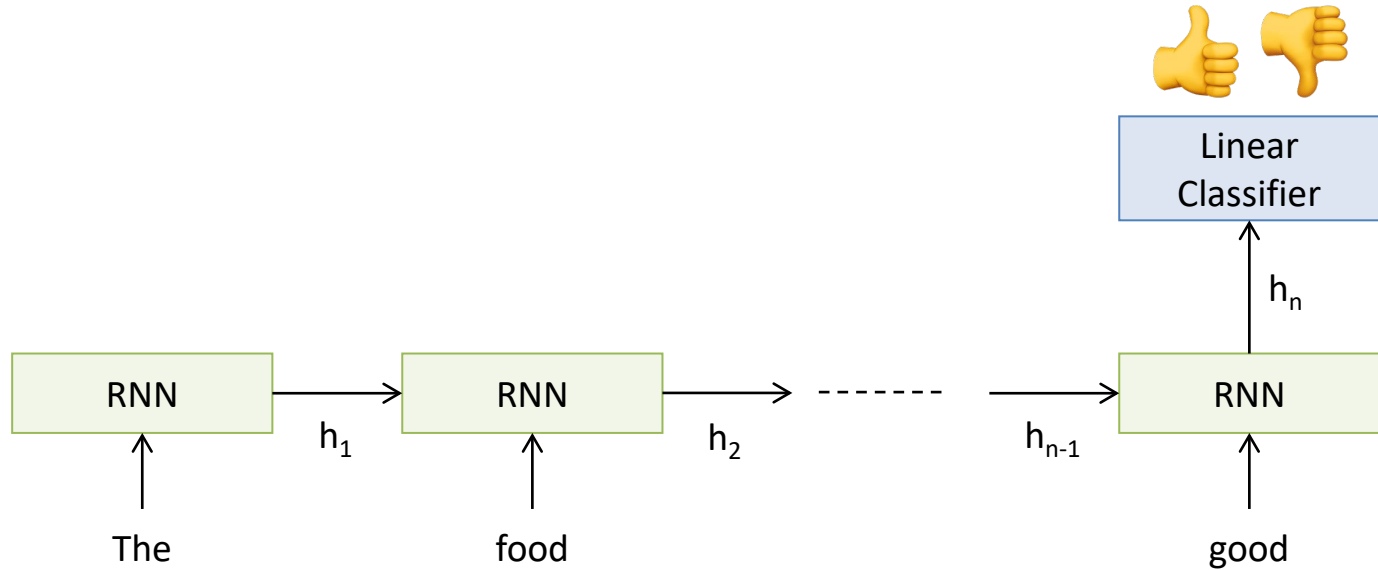
# Sequenzschicht (z.B. unser RNN)



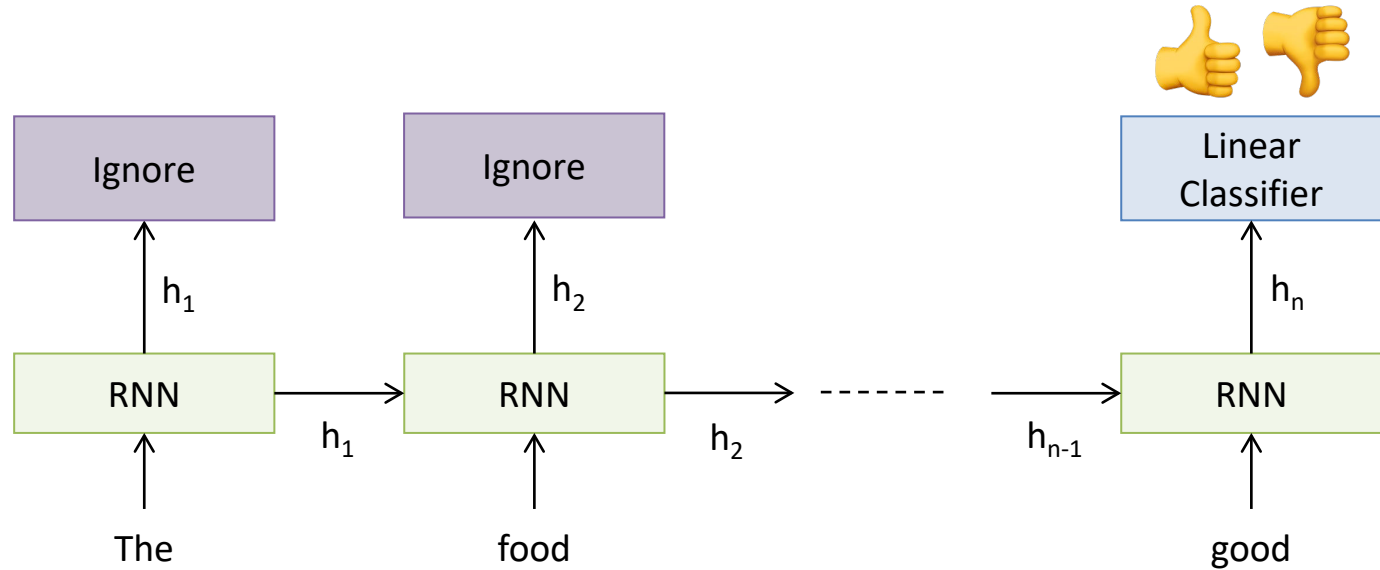
# Sequenzschicht (z.B. unser RNN)



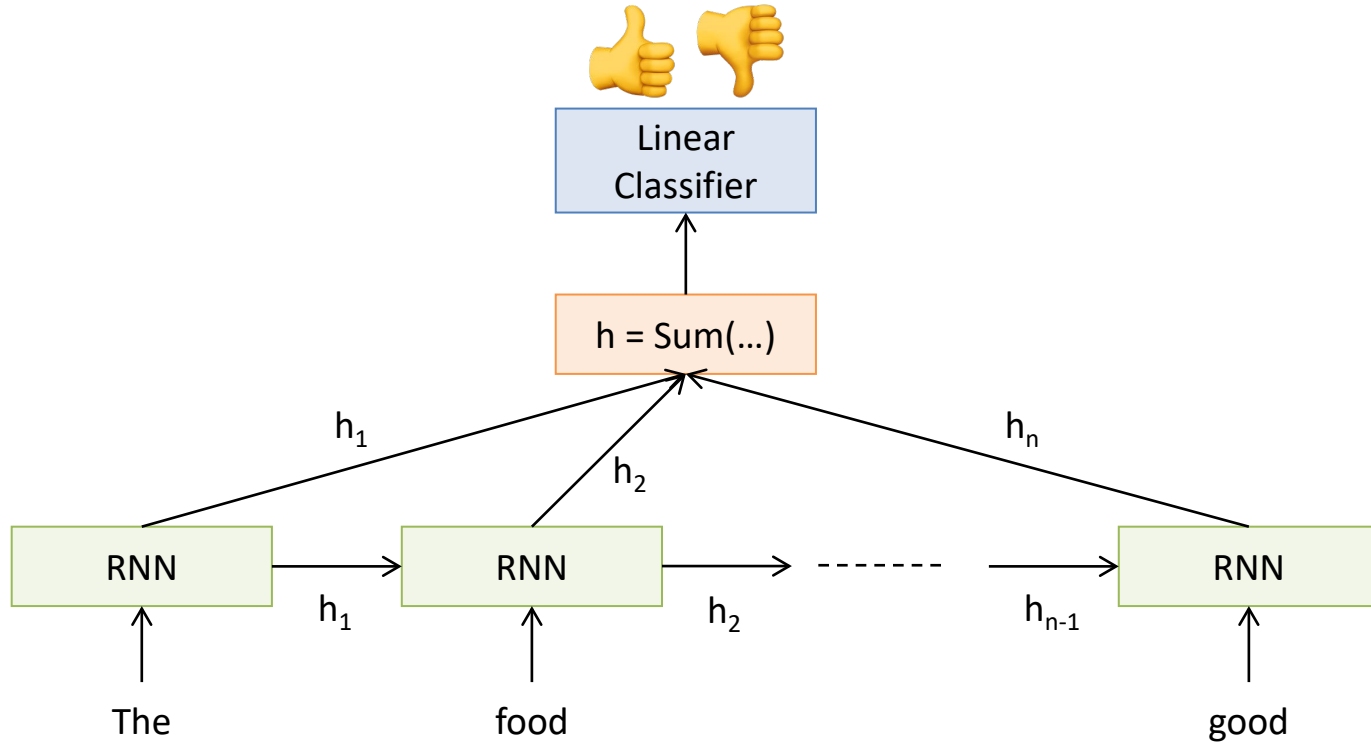
# Aggregationsschicht



# Aggregationsschicht

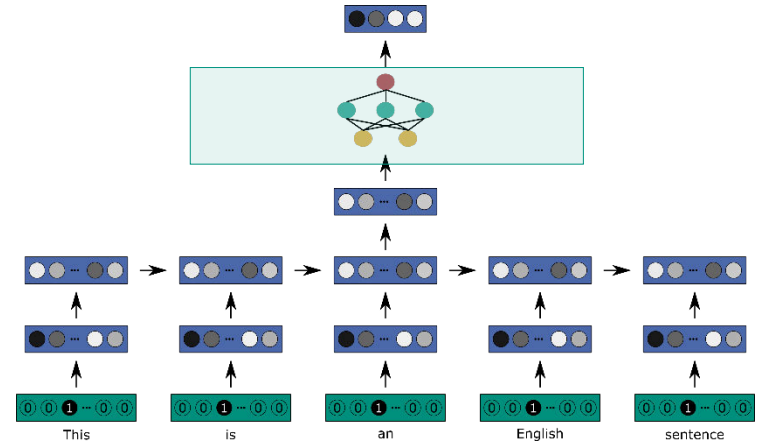


# Aggregationsschicht



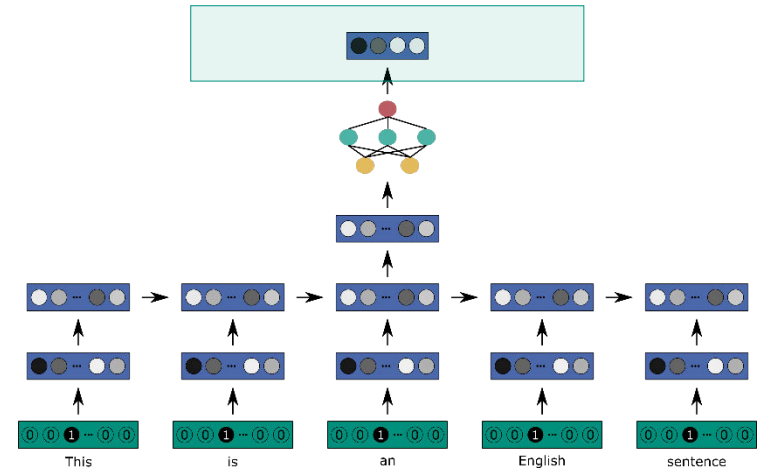
# Grundlegendes Deep Learning Modell für NLP

- Encoder
- Neuronale Netzwerkschichten
  - Bessere Darstellung des Satzinhalts
  - Neuronales Feedforward-Netzwerk
  - Eingabe: Darstellung des Satzes in fester Größe
  - Ausgabe: Darstellung mit fester Größe

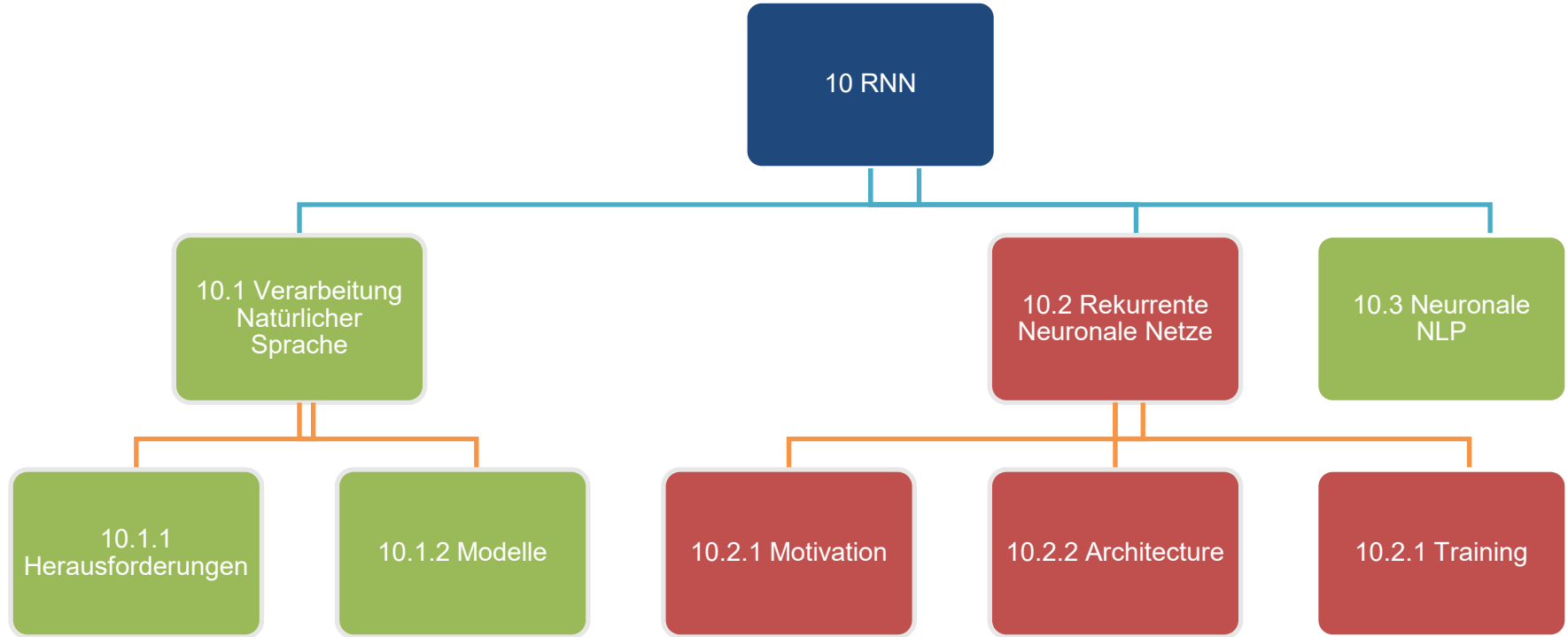


# Grundlegendes Deep Learning Modell für NLP

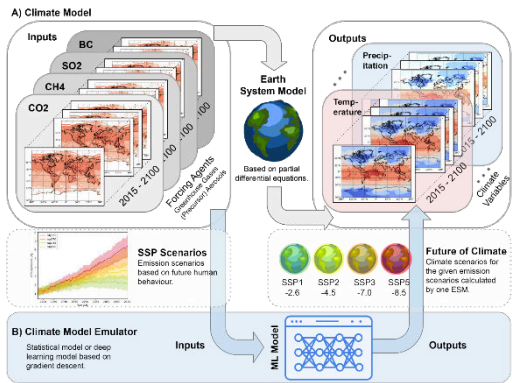
- Encoder
- Neuronale Netzwerkschichten
- Klassifizierungsschicht
  - Softmax-Schicht:
    - Ein Neuron pro Klasse
    - Aktivierung = Wahrscheinlichkeit dieser Klasse
    - $$o_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$
  - Loss Funktion: Cross-entropy Loss



# Überblick

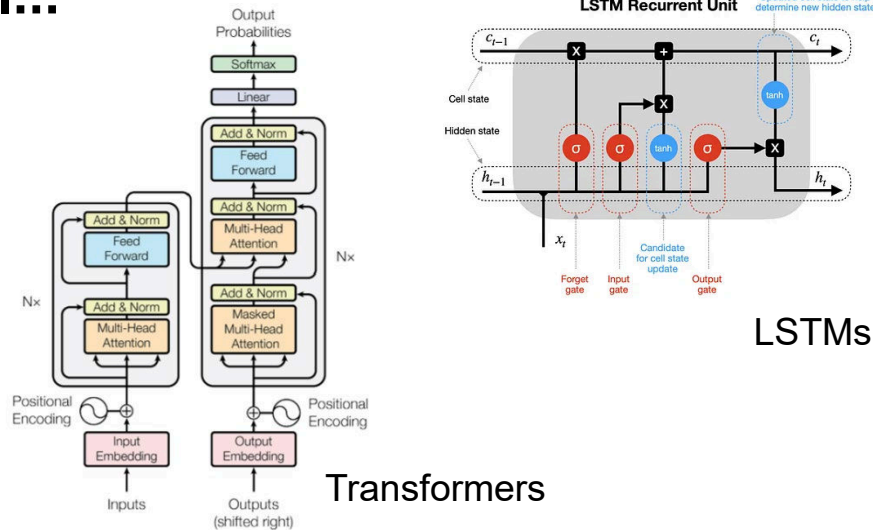
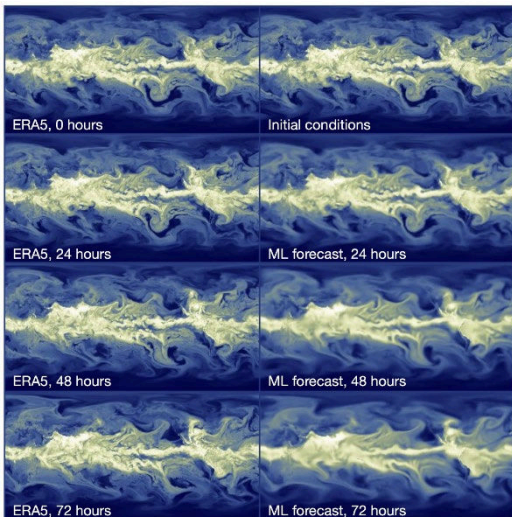


# Ausblick auf unser Mastermodul...



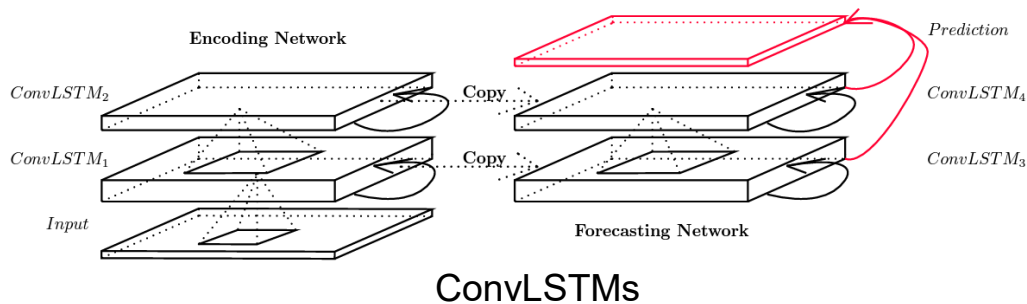
*Climate model emulation*

*Weather forecasting*



LSTMs

Transformers



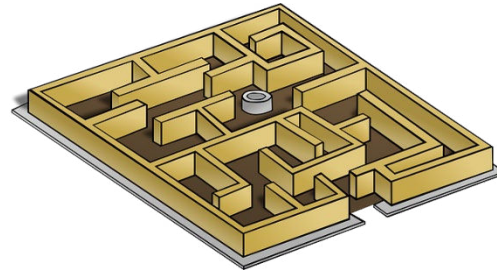
# Zusammenfassung



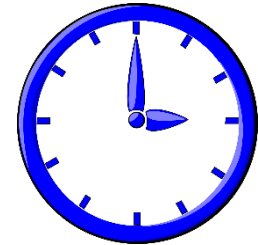
Verarbeitung natürlicher Sprache



Recurrent Neural Network



Herausforderungen



Backpropagation through time

# Lehrevaluation

Vielen herzlichen Dank für Ihre Zeit und Ihr konstruktives Feedback!

## Übung

[https://onlineumfrage.kit.edu/evasys/  
online.php?p=KPQEH](https://onlineumfrage.kit.edu/evasys/online.php?p=KPQEH)



## Vorlesung Friederich

[https://onlineumfrage.kit.edu/evasys/  
online.php?p=Y7FPN](https://onlineumfrage.kit.edu/evasys/online.php?p=Y7FPN)

