# STADS

# Machine Learning Course
# Workshop VI: Natural Language Processing

In **Kooperation** mit

**Prof. Dr. Martin Schlather**
Lehrstuhl für Stochastik und ihre Anwendungen

**Prof. Dr. Leif Döring**
Lehrstuhl für Stochastik

# Course structure

**Part I: Basics of Machine Learning**

04.10 – Introduction to machine learning

11.10 – Neural networks

**Part II: Computer Vision (CV)**

18.10 – Convolutional Neural Networks

25.10 – Practical application of CNNs

**08.11** – CV Challenge presentation

**Part III: Natural Language Processing (NLP)**

**15.11** – Recurrent Neural Networks

22.11 – Practical application of RNNs

**29.11** – NLP Challenge Presentation
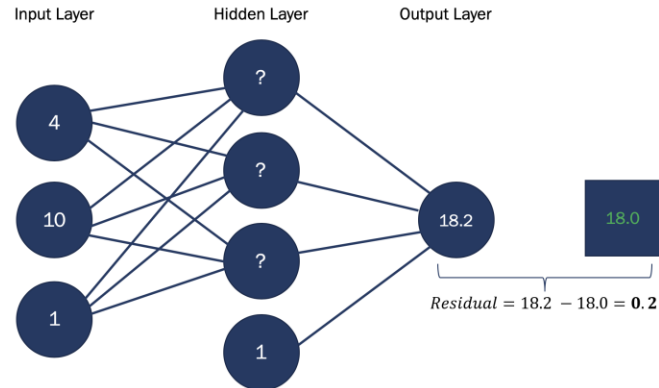
# Disclaimer

This presentation is based on the course on Sequence Models from DeepLearning.ai, as well as the Deep Learning course by Prof. Gemulla. These are highly recommended for in-depth knowledge.

# What happened so far...

# 43% of global companies are already using NLP to increase efficiency[1]



**Revenues from the natural language processing (NLP) market worldwide from 2017 to 2025**

Market in million U.S. dollars

| Year | Value |
|------|-------|
| 2017 | 3,185.7 |
| 2018 | 5,075 |
| 2019* | 8,211.5 |
| 2020* | 12,399.3 |
| 2021* | 17,578.4 |
| 2022* | 23,999.1 |
| 2023* | 30,356.6 |
| 2024* | 37,330.7 |
| 2025* | 43,289.9 |

Source
Tractica
© Statista 2022

Additional Information:
Worldwide; 2017 to 2018

statista

# What actually is Natural Language Processing?

| Natural Language Processing |
| --- |
| Natural language processing aims to develop machines that can understand text or voice input and respond with their own text or voice output - in the same way humans do.[1] |

| Use cases |
| --- |
| 1. Speech recognition<br>2. Sentiment Analysis<br>3. Named Entity Recognition (NEM)<br>4. Chat bots |

| Building blocks |
| --- |
| Classic machine learning or deep learning  **+**  LSTMs, Attention, Embeddings |

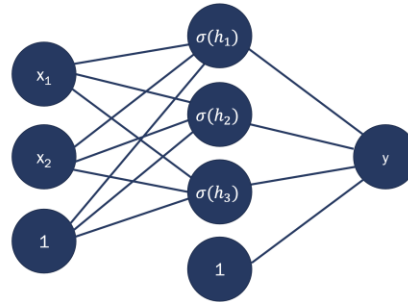1 https://www.ibm.com/cloud/learn/natural-language-processing

# Challenges of Natural Language Processing

## Variable Input-/Output lengths

- Typically, the number of words in sentences varies

- Therefore, classic neural networks with fixed inputs and outputs are unsuitable

## Sequence is important

*"The food in the cafeteria is good and not too expensive."* 🙂

*"The food in the cafeteria is not good and too expensive."* 🙁

## Long-term word dependencies

- Words at the beginning of a long sentence can refer to words at the end of the sentence.

    "The dinosaurs, which became extinct 66 million years ago, were present on every continent."

## Bidirectional dependencies

- The meaning of a word depends not only on the context before the word, but also the words following it.

    „Die Bank hat ein großes Gebäude in Frankfurt."

STADS | Students' Association for Data Analytics & Statistics

# The approach: Recurrent Neural Networks

# Important and exciting applications of RNNs

| | Input | Output |
|---|---|---|
| • Named entity recognition | Tomorrow Nicki Minaj plays in Mannheim $\longrightarrow$ | Tomorrow [Nicki Minaj]$_{Person}$ plays in [Mannehim]$_{Ort}$ |
| • Text classification | The movie was surprisingly bad. $\longrightarrow$ | ★☆☆☆☆ |
| • Machine translation | Das Schiffsverkehr ist eingestellt. $\longrightarrow$ | Shipping traffic is suspended. |
| • Text completion | Write a slogan for a dating app. $\longrightarrow$ | Find your match today! |
| • Speech recognition | | $\longrightarrow$ Alexa, what's the weather? |
| • Image generation | a teddy bear on a skateboard in times square $\longrightarrow$ | |
| • Video analysis | | $\longrightarrow$ Dog catching frisbee |

# Agenda

Introduction

**Theory of RNNs using NLP as an example**

Programming project

STADS | Students' Association for Data Analytics & Statistics

# How can we encode language?

A "dictionary" allows us to represent words as unit vectors (one-hot encoding):

Dictionary

| Aachen | ← | 1 |
| ... | | |
| ... | | |
| play | ← | 7984 |
| ... | | |
| Spotify | ← | 8120 |
| ... | | |
| ... | | |
| Zyzzyva | ← | 14000 |

Hey Siri, play Ed Sheeran on Spotify !

$$\ldots \quad \ldots \quad \begin{pmatrix} 0 \\ \ldots \\ \ldots \\ 1 \\ \ldots \\ 0 \\ \ldots \\ \ldots \\ 0 \end{pmatrix} \quad \ldots \quad \begin{pmatrix} 0 \\ \ldots \\ \ldots \\ 0 \\ \ldots \\ 1 \\ \ldots \\ \ldots \\ 0 \end{pmatrix} \quad \ldots$$

7984

8120

# Example: Named entity recognition

Alexa, play Ed Sheeran on Apple Music!

$\boldsymbol{x}$:
$\quad\quad\quad\quad x^{<1>} \quad x^{<2>} \quad\quad\quad \bullet\bullet\bullet\bullet \quad\quad\quad x^{<7>}$

$$\begin{pmatrix} 0 \\ 1 \\ \cdots \\ 0 \\ \cdots \\ 0 \\ \cdots \\ \cdots \\ 0 \end{pmatrix} \quad \begin{pmatrix} 0 \\ \cdots \\ \cdots \\ 0 \\ \cdots \\ 1 \\ \cdots \\ \cdots \\ 0 \end{pmatrix} \quad\quad\quad\quad\quad\quad \begin{pmatrix} 0 \\ \cdots \\ \cdots \\ 1 \\ \cdots \\ 0 \\ \cdots \\ \cdots \\ 0 \end{pmatrix}$$

$\quad\quad\quad\quad 1 \quad\quad\quad 0 \quad\quad 1 \quad\quad 1 \quad\quad 0 \quad\quad 0 \quad\quad 0$

$\boldsymbol{y}$:
$\quad\quad\quad\quad y^{<1>} \quad y^{<2>} \quad\quad\quad \bullet\bullet\bullet\bullet \quad\quad\quad y^{<7>}$
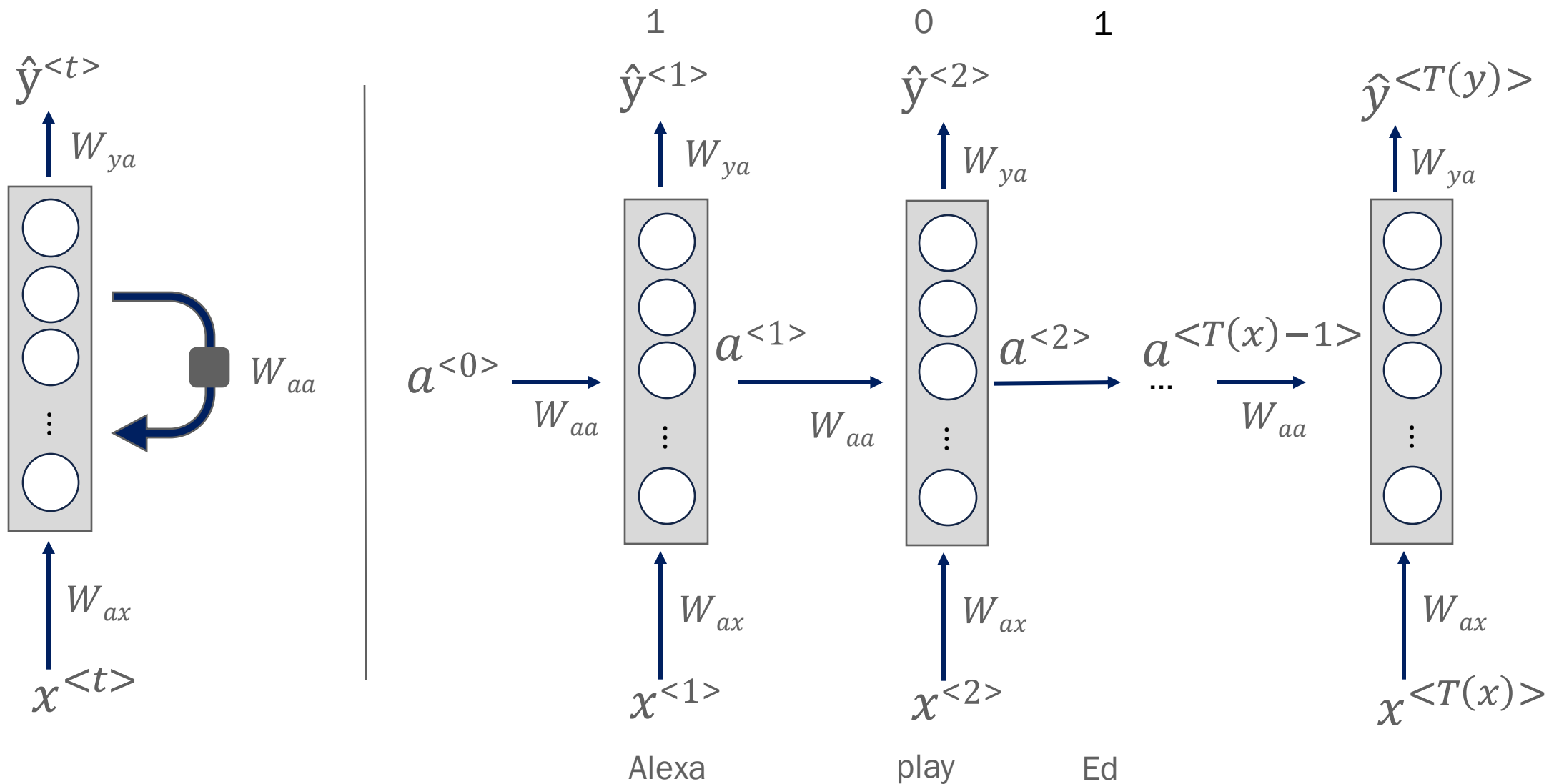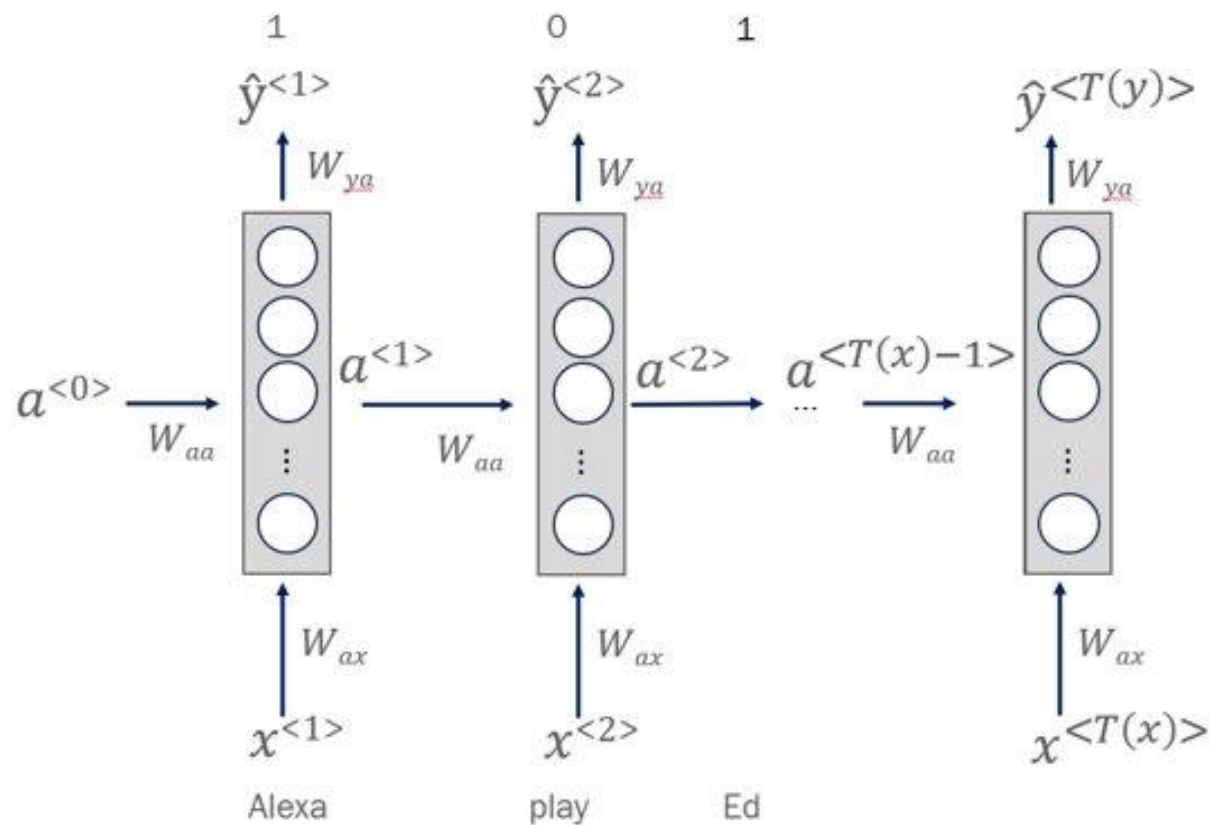
# Architecture of RNNs

# Building RNNs



$$a^{<0>} = 0$$

$$a^{<1>} = g_1(W_{aa} a^{<0>} + W_{ax} x^{<1>} + b_a)$$

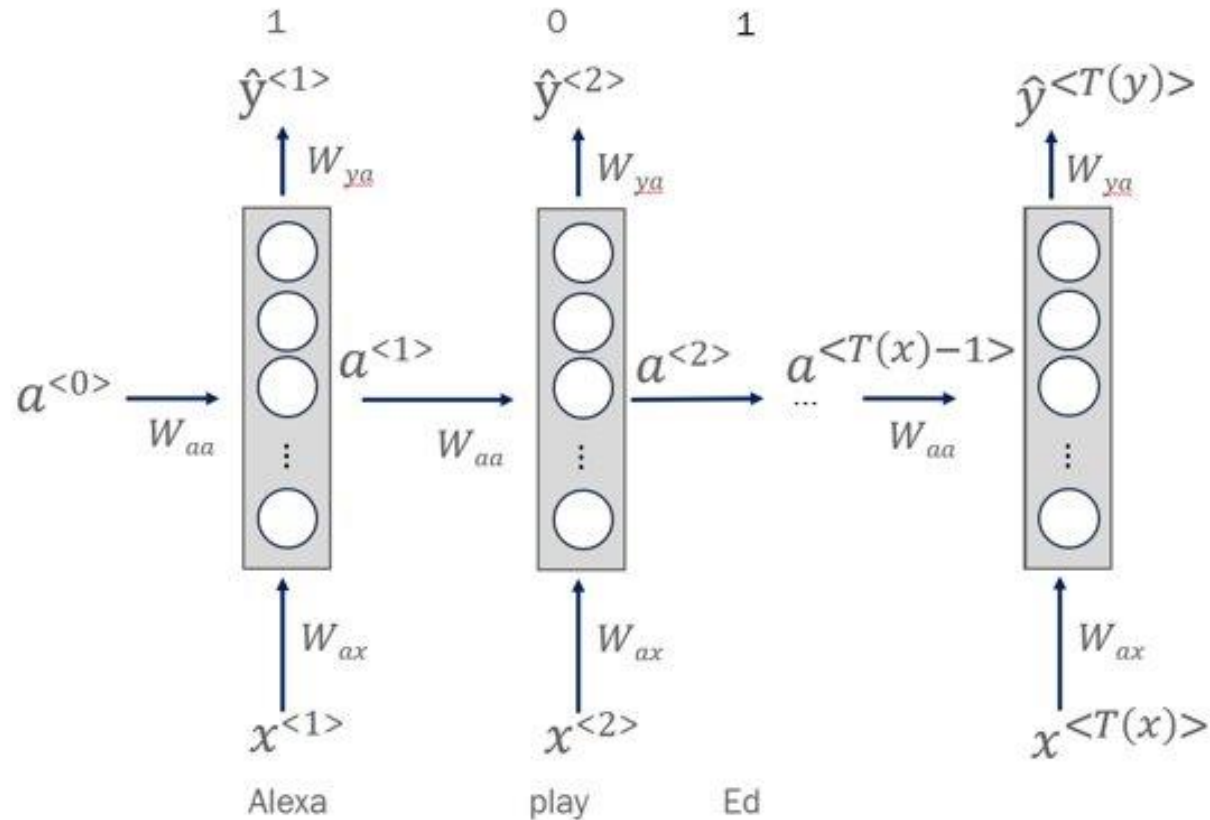$$\hat{y}^{<1>} = g_2(W_{ya} a^{<1>} + b_y)$$

$$\vdots$$

$$a^{<t>} = g(Waa_a^{<t-1>} + W_{ax} x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(W_{ya} a^{<t>} + b_y)$$

# Building RNNs (simplified Notation)



$$W_a = [W_{aa} \mid W_{ax}]$$

$$[a^{t-1} \, x^{<t>}] = \begin{pmatrix} a^{<t-1>} \\ x^{<t>} \end{pmatrix}$$

$$a^{<t>} = g(W_a[a^{<t-1>}, x^{<t>}] + b_a)$$

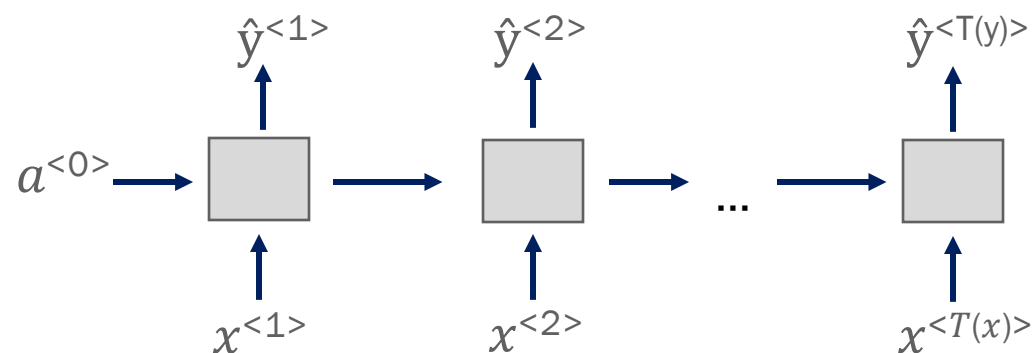$$\hat{y}^{<t>} = g(W_{ya} \, a^{<t>} + b_y)$$

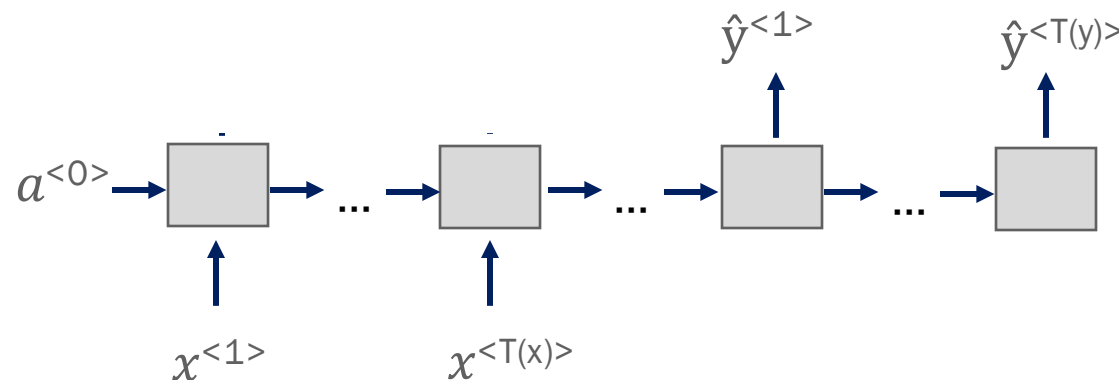# Example of RNN architecture

## Many-to-many $(T(x) = T(y))$
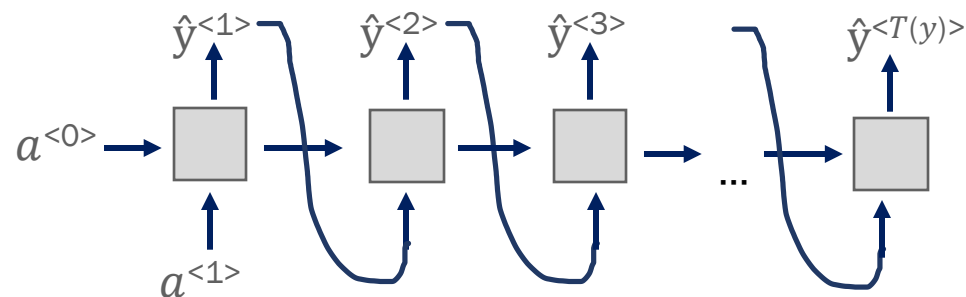
Named-Entity-Recognition, Speech-to-Text

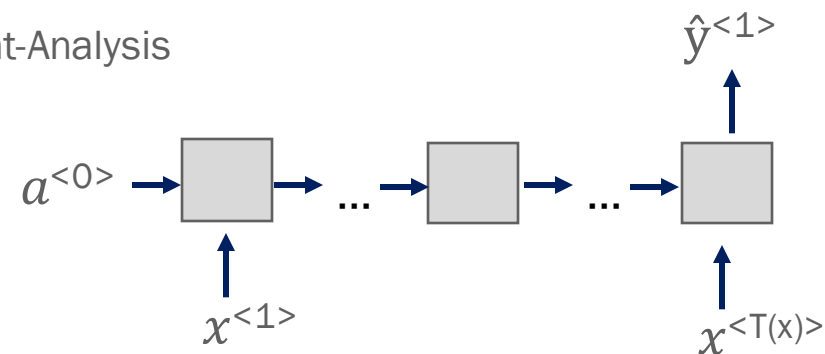## Many-to-many $(T(x) \neq T(y))$

Text-translation

## One-to-many

Music generation, Language modelling
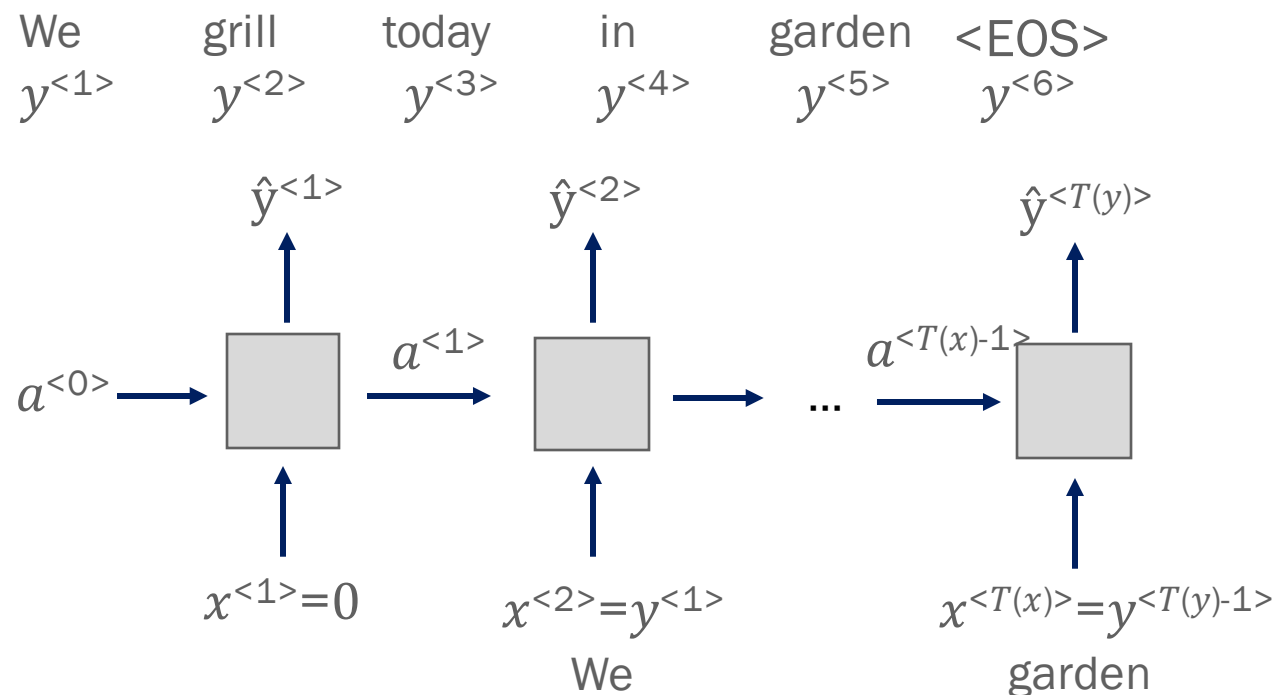
## Many-to-one

Sentiment-Analysis

# Language Models

## Distribution across sequences of words

- Predicting the next word/letter in a text
- Basis for many other NLP applications (text classification, question answering, speech recognition), examples: BERT, GPT-3,

Training data: large dataset of texts (e.g. Wikipedia articles).

| We | grill | today | in | garden | <EOS> |
|---|---|---|---|---|---|
| $y^{<1>}$ | $y^{<2>}$ | $y^{<3>}$ | $y^{<4>}$ | $y^{<5>}$ | $y^{<6>}$ |



$a^{<0>} \rightarrow$ ... $a^{<1>}$ ... $a^{<T(x)-1>}$

$\hat{y}^{<1>}$   $\hat{y}^{<2>}$   $\hat{y}^{<T(y)>}$

$x^{<1>}=0$   $x^{<2>}=y^{<1>}$   $x^{<T(x)>}=y^{<T(y)-1>}$

We   garden

- $\hat{y}^{<1>}$ - Softmax over dictionary $\rightarrow$ modeled $(p(aachen), \dots, p(spielen), \dots p(wir), \dots)$
- $\hat{y}^{<2>} \rightarrow p(\cdot|wir)$
- $\hat{y}^{<3>} \rightarrow p(\cdot|\text{"wir grillen"})$

- $p(y^{<1>}, \cdots, y^{<T(y)>}) = p(y^{<1>})p(y^{<2>}|y^{<1>}) * \cdots * p(y^{<T(y)>}| y^{<1>}, \cdots, y^{<T(y)-1>})$

- Loss:

$$\mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -\sum_i y_i^{<t>} \log(\hat{y}_i^{<t>})$$

$$\mathcal{L} = \sum_t \mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>}).$$

# Generating new Sequences



- Generated words are recursively used as the next input

- Sampling the words according to softmax output (probability vectors) using vocabulary

# Vanishing Gradients

$$\frac{\partial \mathcal{L}}{\partial W} = \sum_{t=1}^{T} \frac{\partial \mathcal{L}_t}{\partial W} = \sum_{t=1}^{T} \frac{\partial \mathcal{L}_t}{\partial y_t} \cdot \frac{\partial y_t}{\partial a_t} \cdot \frac{\partial a_t}{\partial a_{t-1}} \cdot \ldots \cdot \frac{\partial a_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial W} = \sum_{t=1}^{T} \frac{\partial \mathcal{L}_t}{\partial y_t} \cdot \frac{\partial y_t}{\partial a_t} \cdot \prod_{k=1}^{t-1} \frac{\partial a_{k+1}}{\partial a_k} \cdot \frac{\partial a_1}{\partial W}$$
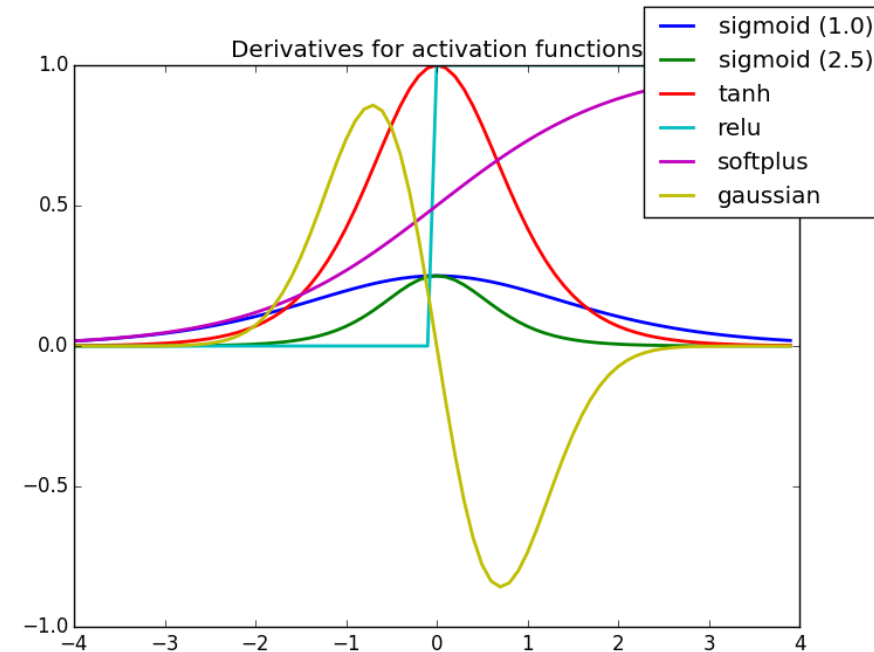
$$a_t = g(W_{aa} \cdot a^{t-1} + W_{ax} \cdot x^t + b_a)$$

$$\frac{\partial a_t}{\partial a_{t-1}} = g'(W_{aa} \cdot a_{t-1} + W_{ax} \cdot x_t + b_a) \cdot \frac{\partial}{\partial a_{t-1}}[W_{aa} \cdot a_{t-1} + W_{ax} \cdot x_t]$$

$$= g'(\ldots) \cdot W_{aa}$$

$$< 1$$



Derivatives for activation functions

Legend: sigmoid (1.0), sigmoid (2.5), tanh, relu, softplus, gaussian

Bildquelle (23.10.22):
https://miro.medium.com/max/1400/1*n1HFBpwv21FCAzGjmWt1sg.png

**Problem:** Since the gradient decreases with the number of words, long-term word dependencies are not taken into account when updating the parameters! But this is necessary for language models.
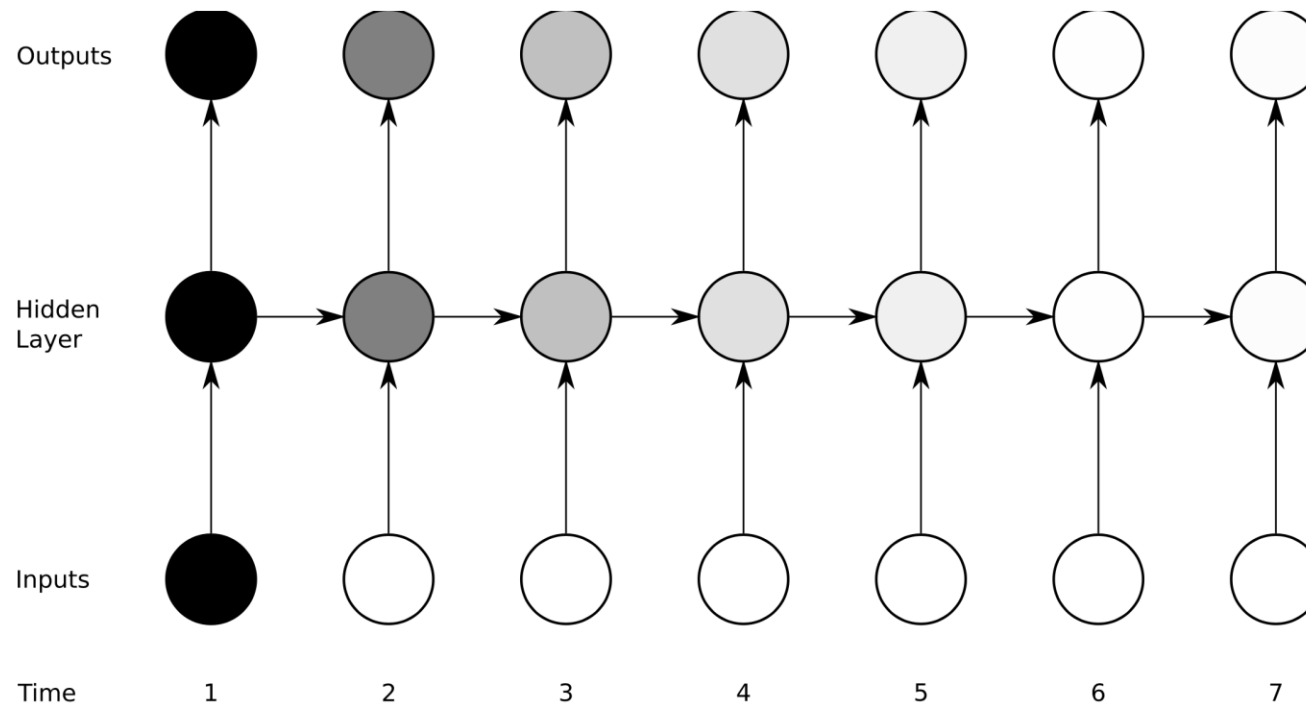
# Vanishing Gradients



Fig. 10.   Illustration of the vanishing gradient problem. The diagram represents a recurrent network unrolled in time. The units are shaded according to how sensitive they are to the input at time 1 (where black is high and white is low). As can be seen, the influence of the first input decays exponentially over time.
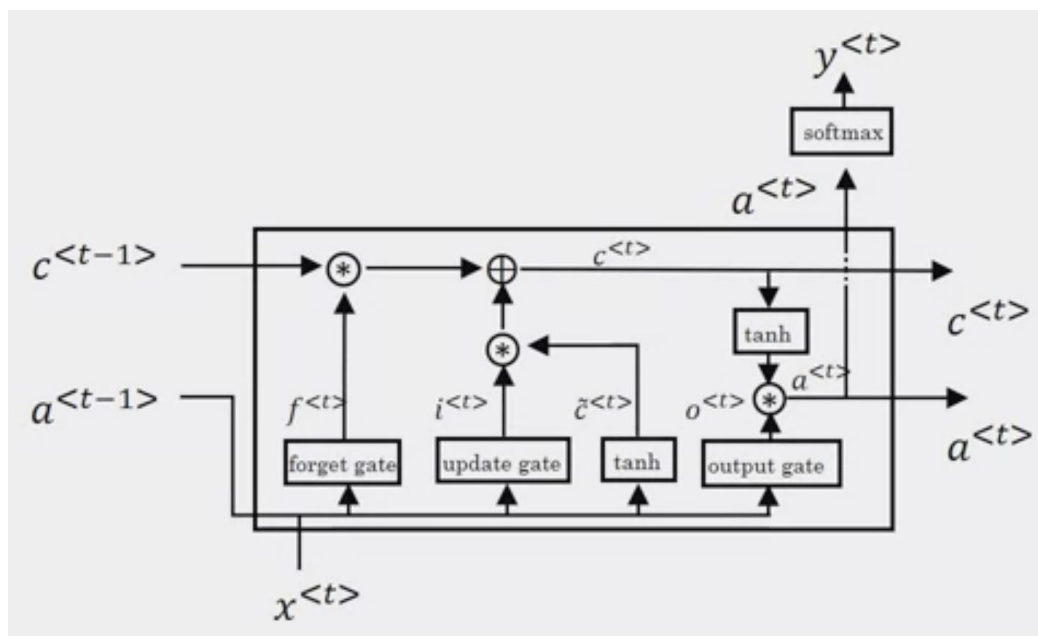
LSTMs can solve this problem!

# Long short-term memory cells (LSTM)

- Solution to the Vanishing Gradient Problem
- Introduction of a "memory cell" c<t> (same dimension as activation a<t>)
- "Gating mechanism" decides when memory cell is updated



$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

update $\longrightarrow \quad \Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$

forget $\longrightarrow \quad \Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$

output $\longrightarrow \quad \Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \tanh c^{<t>}$$

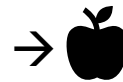*: Element-by-element multiplication of the vectors

Intuition: Sigmoid function σ to a small area either ≈ 0 or ≈ 1

$\longrightarrow$ Memory cell remains almost constant when $\Gamma_u \approx 0$ and $\Gamma_f \approx 1$

$\longrightarrow$ preserves sensitivity of the outputs to inputs from much earlier time steps
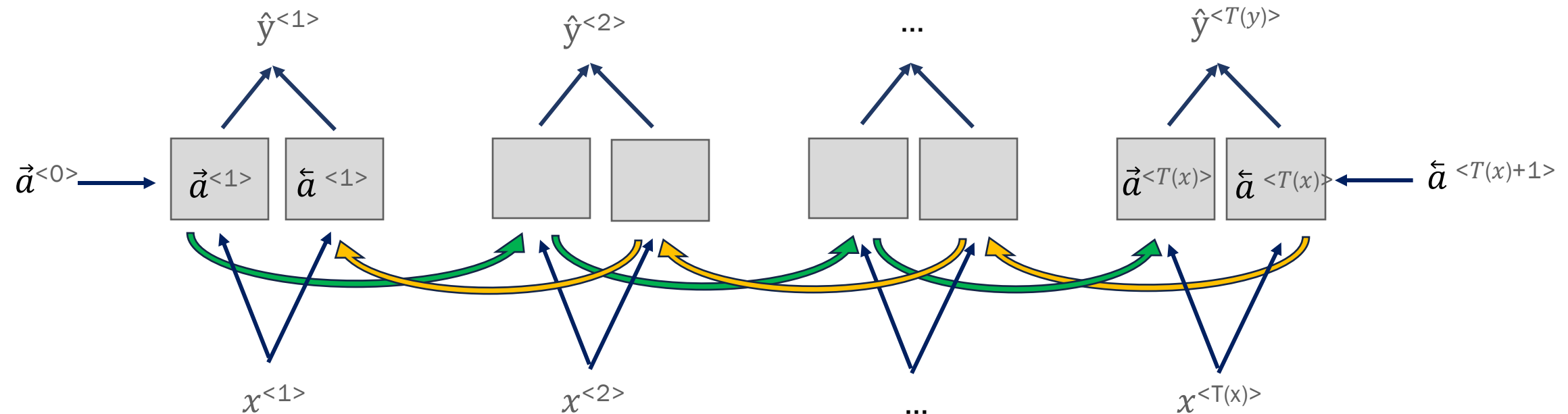
# Bidirectional RNN's

An apple a day is, keeps the doctor away. → 🍎

Apple devices use a lot of Machine Learning. → Brand/company name $\hat{y}^{<t>} = g(W_{ya}[\vec{a}^{<t>}, \overleftarrow{a}^{<t>}] + b_y)$



- Hidden units are also connected in the backward direction → enables use of information from the future
- Not suitable for real-time applications

# Word Embeddings

so far: Representation of every word from a dictionary of size V as a one-hot vector

$$\text{Frau} \triangleq \begin{pmatrix} 0 \\ \dots \\ 0 \\ 1 \\ 0 \\ \dots \\ \dots \\ \dots \\ 0 \end{pmatrix} = e_{26923}$$

| | |
|---|---|
| Mann | $\triangleq e_{51034}$ |
| Königin | $\triangleq e_{34263}$ |
| König | $\triangleq e_{34003}$ |
| super | $\triangleq e_{86831}$ |
| toll | $\triangleq e_{89239}$ |

Solution: feature vectors / word embeddings

- Representation of words using 300 dimensions in vector $emb_j$
- Embedding Matrix $E \in \mathbb{R}^{300 \times V}$
- $emb_i = E \cdot e_i$
- Embedding Layer implements efficient mapping
- Enables Transfer Learning

Problem:
- $\left\| e_i - e_j \right\|_2 = 1 \ \forall \ i \neq j$
- $\dim(e_i) = V$, typically $V >> 10000$ (word2vec $V \approx 3M$)

We want:

$e_{\text{Mann}} - e_{\text{Frau}} \approx e_{\text{König}} - e_{\text{Königin}}$

$e_{\text{super}} \approx e_{\text{toll}}$

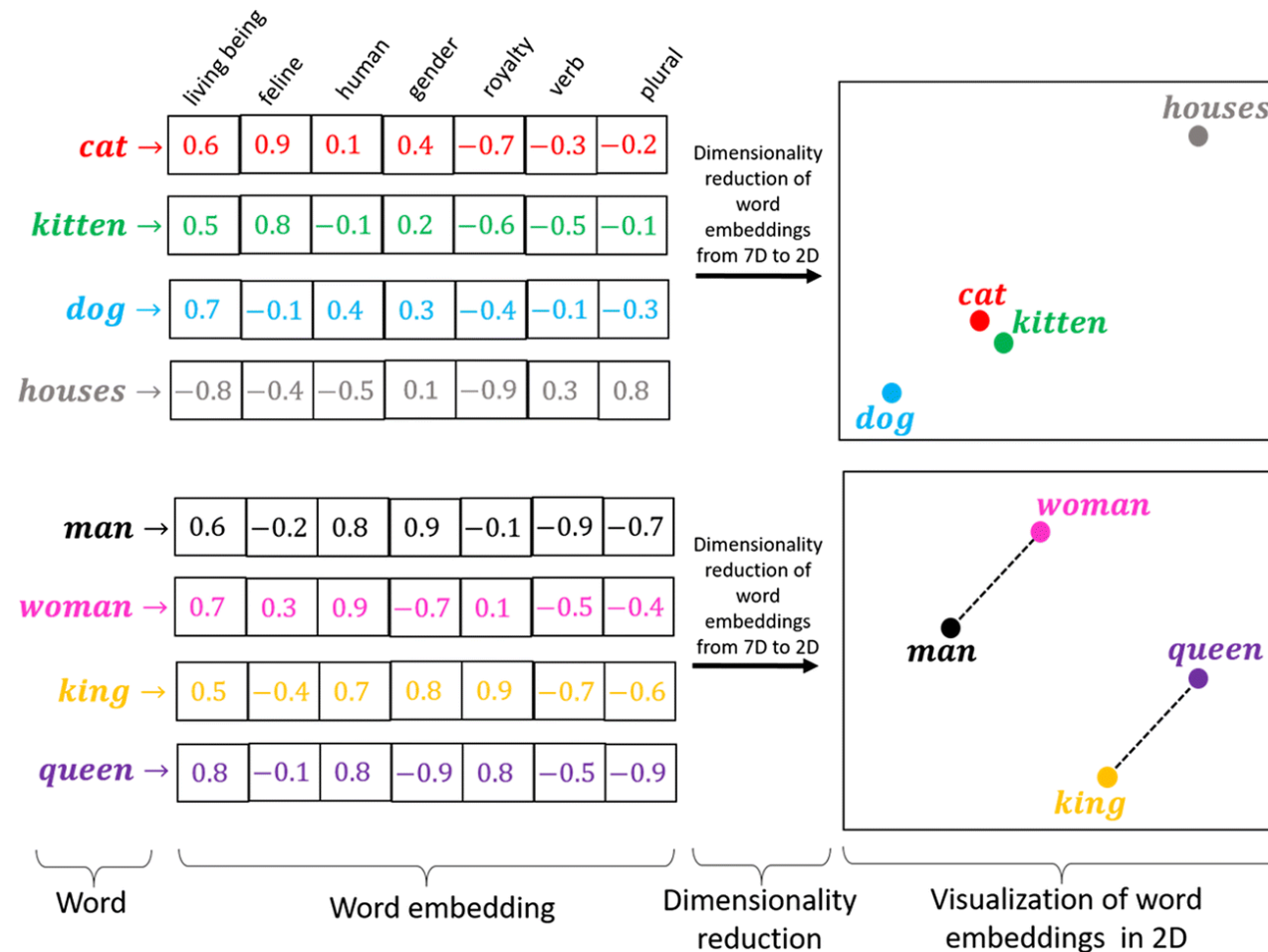Classification

"The film was awesome"
"The film was super"
NER
"Thomas is a dentist"
"Lisa is an optalmologist"

# Word Embeddings – Visualisation



Word embeddings fulfill the conditions such as:

$$Emb_{king}=argmax_{x\in emb}(x, emb_{man} - emb_{woman} + emb_{king}$$

Transfer learning:
- Using pre-trained word vectors on huge datasets as an embedding layer
- e.g. Word2Vec, GLoVe, ELMo
- Generalization regarding words that possibly did not appear in the small training data set

t-sne algorithm to plot high-dimensional data in 2 or 3D

# Lerning Word Embeddings

Skip-grams

| I | cook | Lasagna | with | fresh | tomatos. |
|---|------|---------|------|-------|----------|
| $e_{3849}$ | $e_{4123}$ | $e_{4476}$ | $e_{5134}$ | $e_{2490}$ | $e_{8123}$ |

Sample context/target pairs:

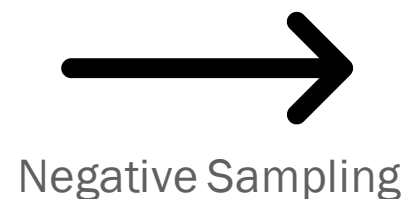| context | target |
|---------|--------|
| Lasagna | cook |
| Lasagna | fresh |
| Lasagna | tomatos |

Example: context=„Lasagna"  target=„fresh"

Model:

$$e_{context} \rightarrow E \rightarrow emb_{context} \rightarrow softmax(V) \rightarrow \hat{y} \in \mathbb{R}^V$$
$$\text{loss: } \mathcal{L}(\hat{y}, e_{target})$$

Problem: softmax over V classes is computationally very heavy.

From classification over V
to V binary classifications

| context | Wort | target |
|---------|------|--------|
| Lasagna | cook | 1 |
| Lasagna | house | 0 |
| Lasagna | queen | 0 |
| Lasagna | car | 0 |

neg. samples

Negative Sampling

$$e_{context} \rightarrow E \rightarrow emb_{context}$$

→ Sigmoid("cooks")
→ Sigmoid("house")
…
→ Sigmoid("car")

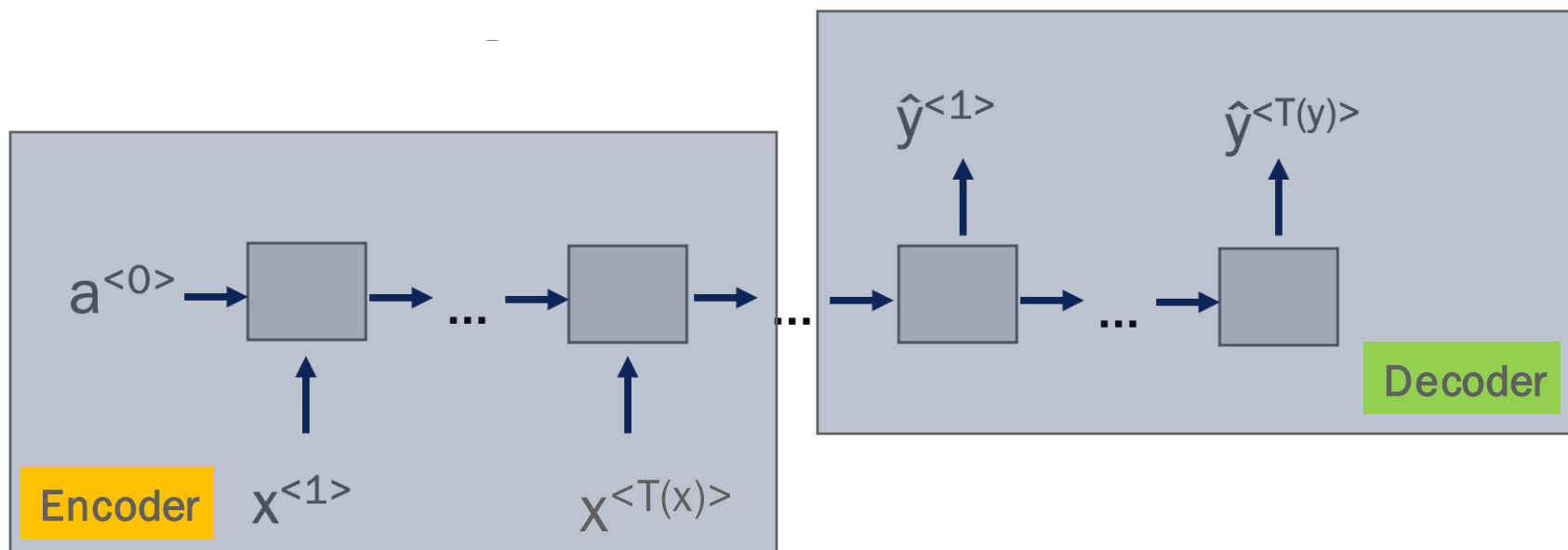# Sequence to Sequence Models

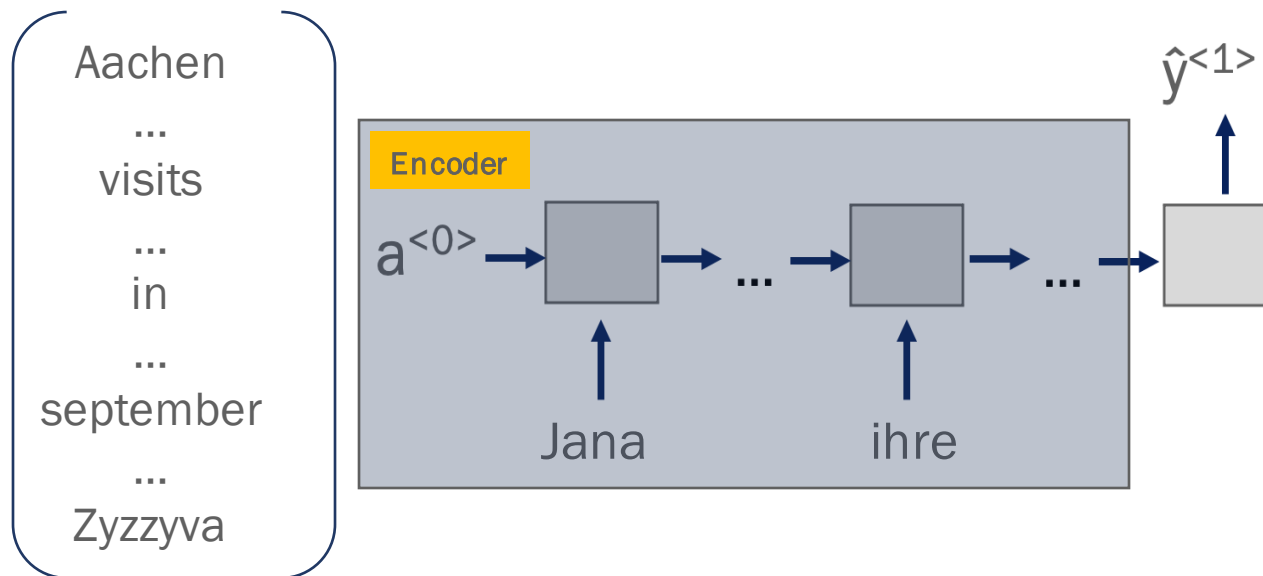**Input:** "Jana visits her parents in Aachen."

**Output:** "Jana besucht ihre Eltern in Aachen."

# Beam Search (B=3)

## 1. Step

Aachen
...
visits
...
in
...
september
...
Zyzzyva



$\hat{y}^{<1>}$

Encoder

$a^{<0>}$

...

...

Jana

ihre

Determine conditional probabilities of
1. all vocabulary words
   $p(\hat{y}^{<1>} \mid x)$

2. Choose **B=3** probability outputs.

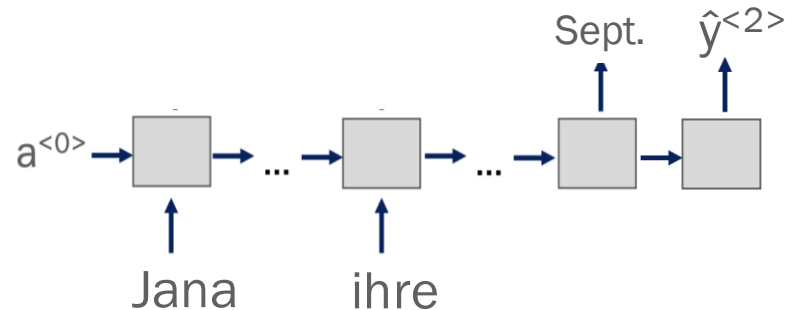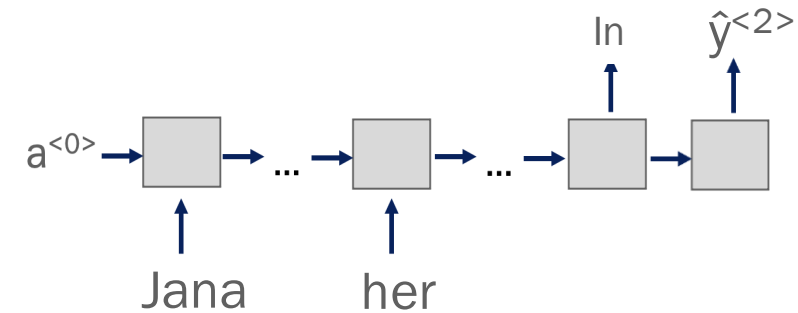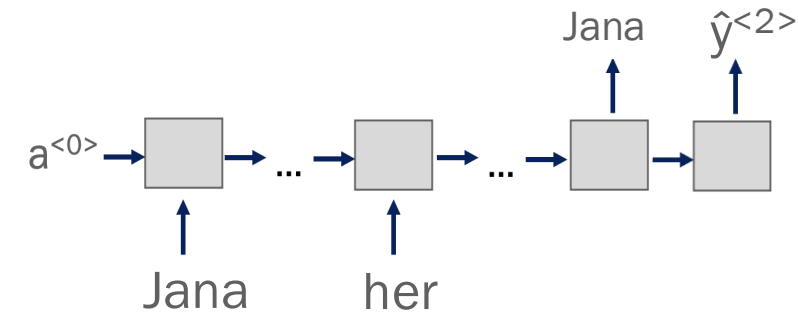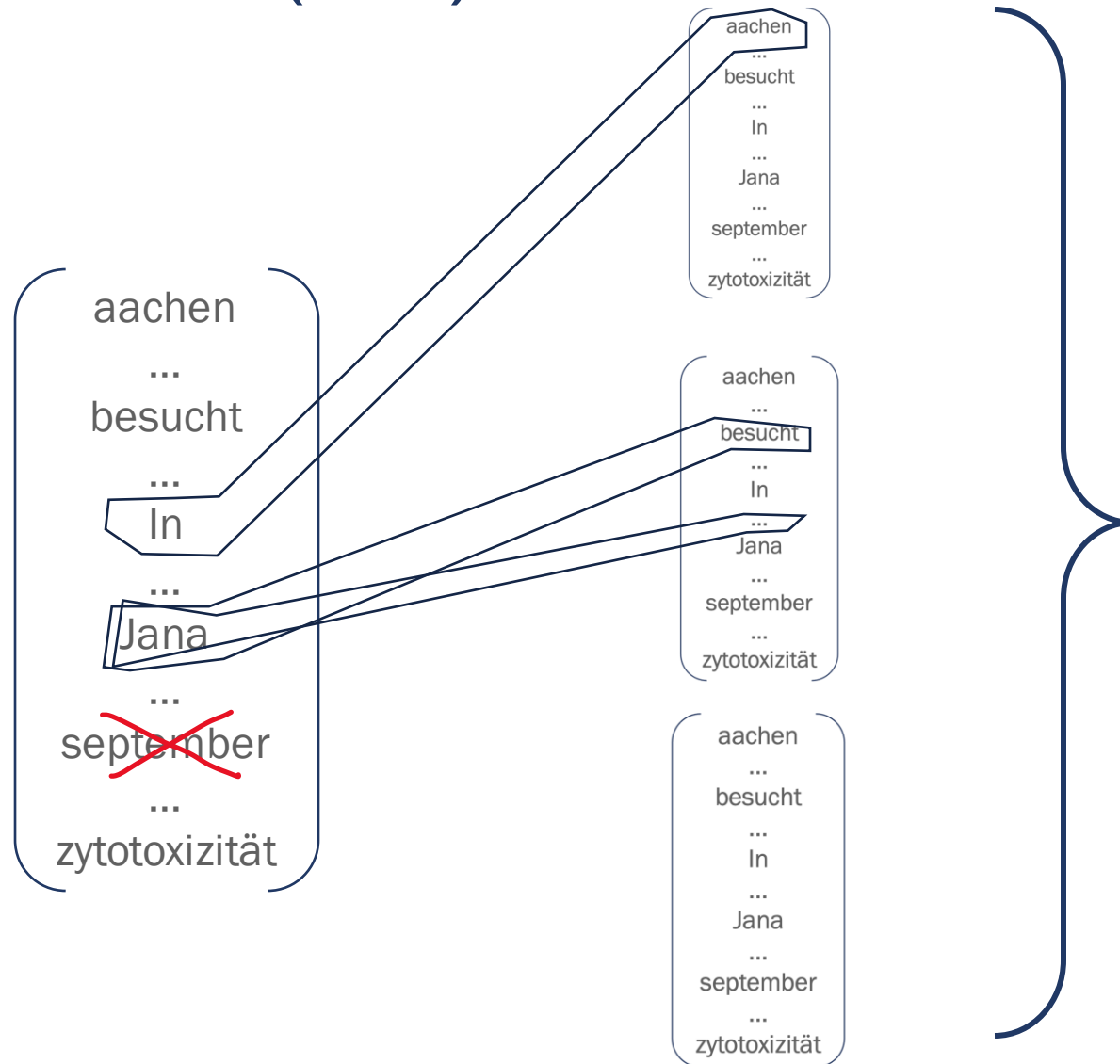3. In this case, e.g. Jana, in, September

STADS | Students' Association for Data Analytics & Statistics

# Beam Search (B=3)

## 2. Step

Aachen

…

visits

…

in

…

september

…

Zyzzyva



Jana $\hat{y}^{<2>}$

$a^{<0>}$

… …

Jana her

In $\hat{y}^{<2>}$

$a^{<0>}$

… …

Jana her

Sept. $\hat{y}^{<2>}$

$a^{<0>}$

… …

Jana ihre

Determine conditional probabilities in

1. each case
$$p(\hat{y}^{<2>} \mid x, \hat{y}^{<1>} = Jana),$$
$$p(\hat{y}^{<2>} \mid x, \hat{y}^{<1>} = In),$$
$$p(\hat{y}^{<2>} \mid x, \hat{y}^{<1>} = September)$$

2. Calculate
$$p(\hat{y}^{<1>}, \hat{y}^{<2>} \mid x)$$
$$= p(\hat{y}^{<1>} \mid x) \cdot p(\hat{y}^{<2>} \mid x, \hat{y}^{<1>})$$
and choose $(\hat{y}^{<1>}, \hat{y}^{<2>})$
with the highest probability

3. We are again given 3 options with which we can continue in the next step

STADS | Students' Association for Data Analytics & Statistics

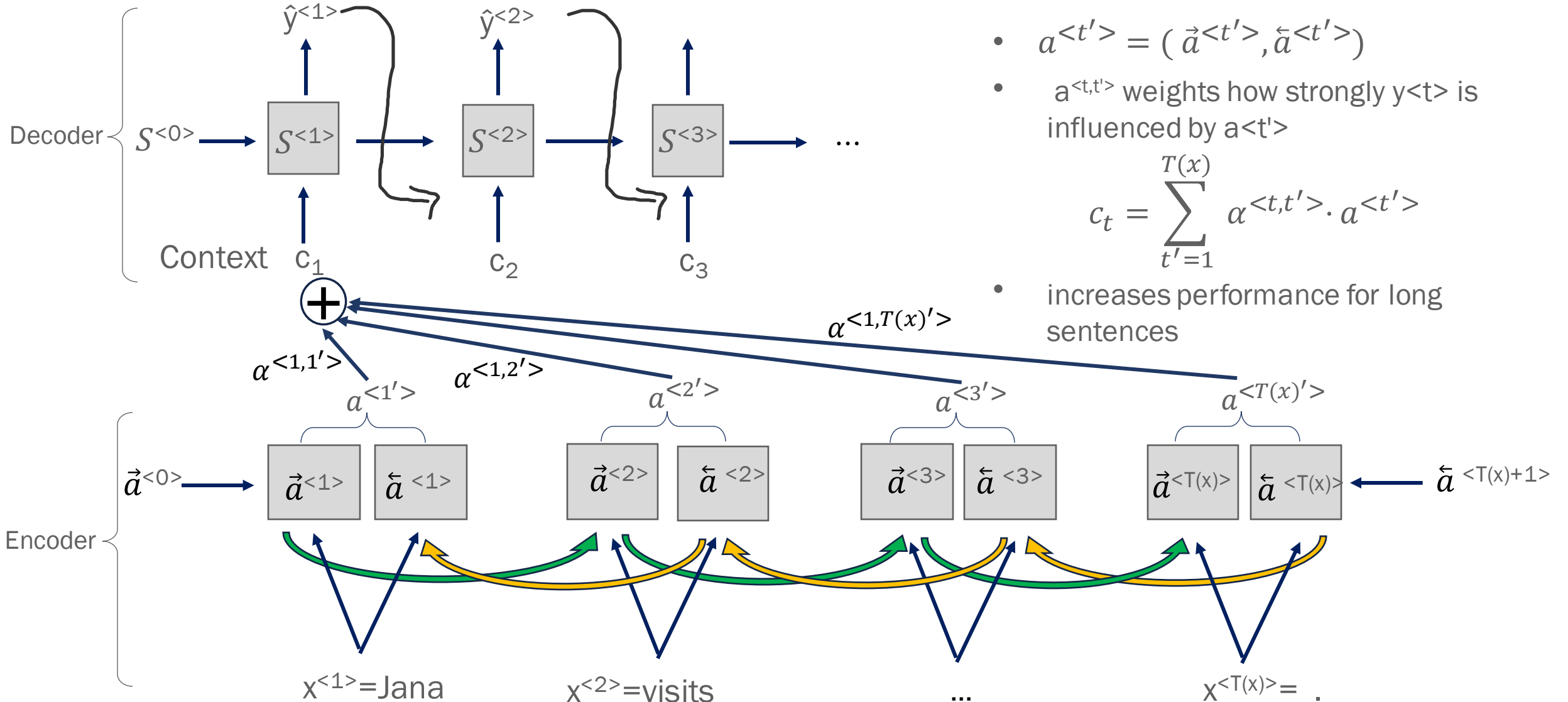# Beam Search (B=3)

## 2. Step



In this way we iteratively obtain the conditional probability we are looking for

$$p(\ \hat{y}^{<1>},\ \hat{y}^{<2>}, \dots, \hat{y}^{<T(y)>}\big|\ x)$$

# Attention Models – which Inputs are relevant?



- $a^{<t'>} = (\vec{a}^{<t'>}, \overleftarrow{a}^{<t'>})$

- $a^{<t,t'>}$ weights how strongly y<t> is influenced by a<t'>

$$c_t = \sum_{t'=1}^{T(x)} \alpha^{<t,t'>} \cdot a^{<t'>}$$

- increases performance for long sentences

Decoder

$S^{<0>}$ $S^{<1>}$ $S^{<2>}$ $S^{<3>}$ ...

$\hat{y}^{<1>}$ $\hat{y}^{<2>}$

Context $c_1$ $c_2$ $c_3$

$\alpha^{<1,T(x)'>}$

$\alpha^{<1,1'>}$ $\alpha^{<1,2'>}$

$a^{<1'>}$ $a^{<2'>}$ $a^{<3'>}$ $a^{<T(x)'>}$

Encoder

$\vec{a}^{<0>}$ $\vec{a}^{<1>}$ $\overleftarrow{a}^{<1>}$ $\vec{a}^{<2>}$ $\overleftarrow{a}^{<2>}$ $\vec{a}^{<3>}$ $\overleftarrow{a}^{<3>}$ $\vec{a}^{<T(x)>}$ $\overleftarrow{a}^{<T(x)>}$ $\overleftarrow{a}^{<T(x)+1>}$

x<1>=Jana x<2>=visits ... x<T(x)>= .

# Attention Models – which Inputs are relevant?

- $a^{<t'>} = (\vec{a}^{<t'>}, \overleftarrow{a}^{<t'>})$

- $a^{<t,t'>}$ weights how strongly y<t> is influenced by a<t'>

$$c_t = \sum_{t'=1}^{T(x)} \alpha^{<t,t'>} \cdot a^{<t'>}$$

- increases performance for long sentences

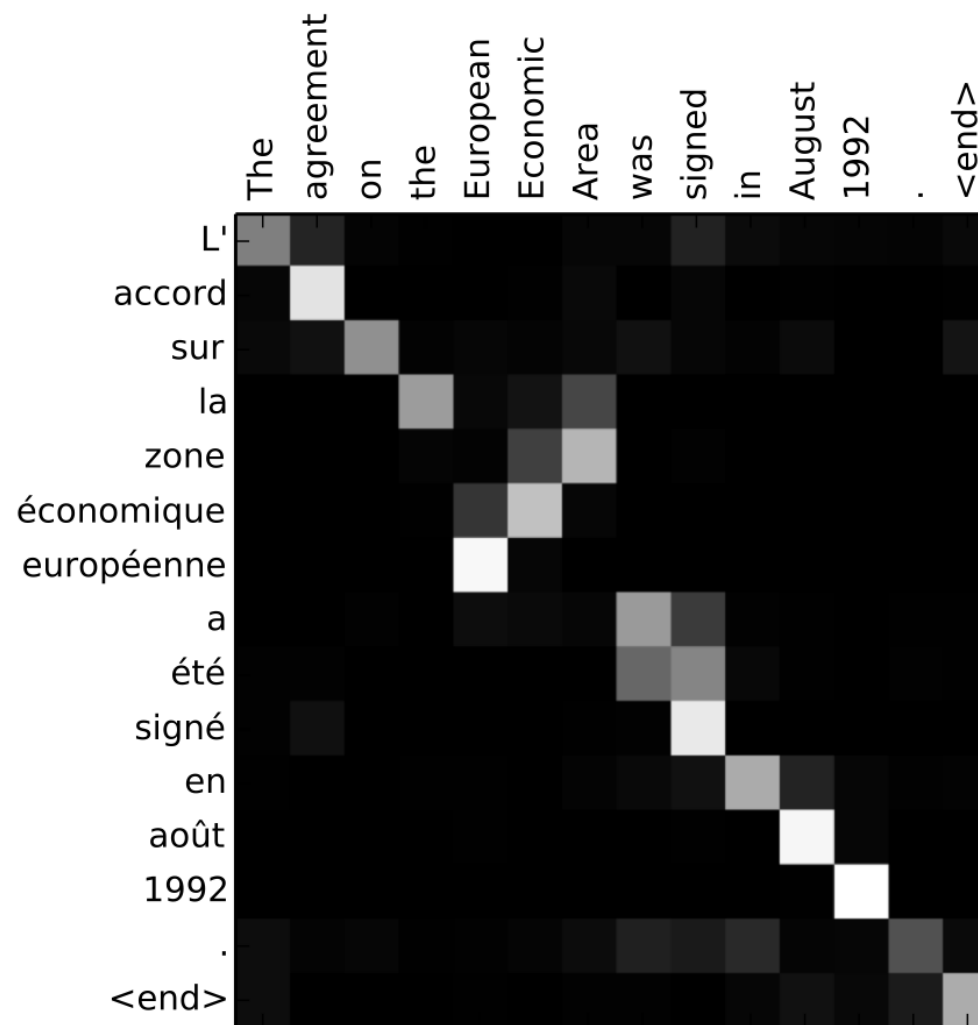# Attention Modelle: Berechnung der $\alpha^{<t,t'>}$

$$\alpha^{<t,t'>} = \frac{\exp\left(e^{<t,t'>}\right)}{\sum_{t'=1}^{T(x)} \exp(e^{<t,t'>})} \quad \text{(Softmax)}$$

$s^{<t-1>}$

$a^{<t'>}$

$e^{<t,t'>}$

→ α's and thus the context depend on the activations of the encoder as well as the activations/hidden states of the decoder at the previous time step

# Attention for Image Captioning



Figure 3. Examples of attending to the correct object (*white* indicates the attended regions, *underlines* indicated the corresponding word)
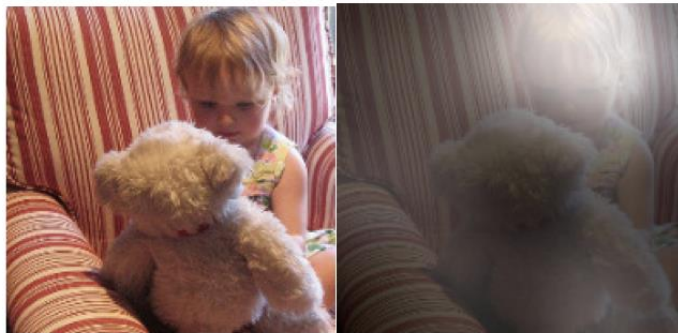
A woman is throwing a <u>frisbee</u> in a park.

A <u>dog</u> is standing on a hardwood floor.

A <u>stop</u> sign is on a road with a mountain in the background.

A little <u>girl</u> sitting on a bed with a teddy bear.

A group of <u>people</u> sitting on a boat in the water.

A giraffe standing in a forest with <u>trees</u> in the background.

# Agenda

**STADS**