# Recap

- ► Word2Vec
  - ► Method to represent words in vector space
  - ► Requires large quantity of raw text
  - ► Pre-trained embeddings can be shared
  - ► Embeddings capture (some aspects of) lexical meaning

# Large Language Models
## Sprachverarbeitung (VL + Ü)

Nils Reiter

June 29, 2023

# Group Exercise

1. In which situations have you talked about ChatGPT (& co)?
2. For which tasks can it be put to use?
3. For which tasks *should* it not be used? Why not?

# Brief history of Computational Linguistics II

- ▶ 1984: First corpus-based commercial MT system                    Nagao (1984)
- ▶ 1992: Study programs established in Germany (Saarbrücken/Stuttgart)
- ▶ 2011: IBM Watson beats two humans in Jeopardy `YouTube` / Apples Siri launched
- ▶ 2013: Word embeddings (e.g., word2vec)                    Mikolov et al. (2013)
- ▶ 2017: Launch of the DeepL Translator (a Cologne-based company)
- ▶ 2018: Transformer models: BERT                    Devlin et al. (2019)
- ▶ 2022: ChatGPT `chat.openai.com`
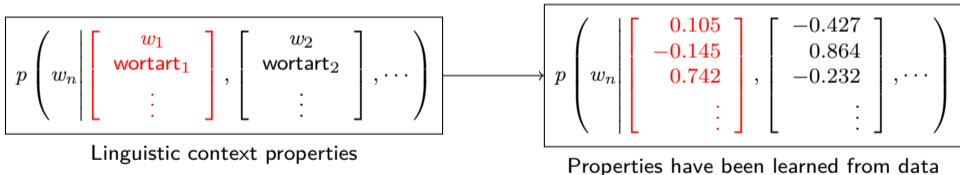  - ⚠ Yes, we need to talk about ChatGPT ⬇

# Large Language Models

- ▶ Term LLM used in contrast to classical language models
- ▶ Family of »transformer models«: BERT, GPT, …
  - ▶ BERT by Google, GPT-X by OpenAI
  - ▶ BERT model can be downloaded and used locally
- ▶ Huge amount of training data (e.g., the web)
- ▶ High computing costs
  - ▶ »Just how much does it cost to train a model? Two correct answers are ›depends‹ and ›a lot‹.«                                                        Sharir et al. (2020, 1)
  - ▶ BERT w/ 340 million parameters: $ 10k / $ 200k

# Key Idea 1: Learned Representation

- Classical ML: Instances are represented by their features
- Neural ML
    - Words/texts are represented by vectors
    - Vectors are learned representations
        - I.e., vectors are optimised for some task, usually filling gaps in texts

# Key Idea 1: Learned Representation

▶ Classical ML: Instances are represented by their features
▶ Neural ML
  ▶ Words/texts are represented by vectors
  ▶ Vectors are learned representations
    ▶ I.e., vectors are optimised for some task, usually filling gaps in texts

$$p\left(w_n \middle| \begin{bmatrix} w_1 \\ \text{wortart}_1 \\ \vdots \end{bmatrix}, \begin{bmatrix} w_2 \\ \text{wortart}_2 \\ \vdots \end{bmatrix}, \cdots \right)$$

Linguistic context properties

$$p\left(w_n \middle| \begin{bmatrix} 0.105 \\ -0.145 \\ 0.742 \\ \vdots \end{bmatrix}, \begin{bmatrix} -0.427 \\ 0.864 \\ -0.232 \\ \vdots \end{bmatrix}, \cdots \right)$$

Properties have been learned from data

# Key Idea 2: Not every token is equally important

- ▶ »Attention Is All You Need« <span>Vaswani et al. (2017)</span>
- ▶ Idea: During training, model learns which tokens are relevant to predict the output
  - ▶ Additional parameters to train …
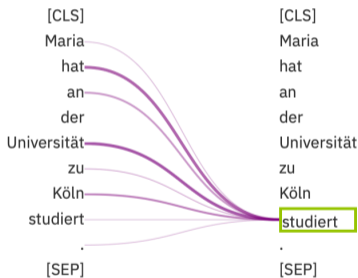


Figure: Attention given when predicting the word »studiert«. Screenshot taken from

https://huggingface.co/spaces/exbert-project/exbert

# Key Idea 3: Training Process Split into two Phases

▶ Traditionally (Naive Bayes, Decision tree, …), we train a model and are done
▶ Transformer architecture:
  ▶ Pre-Training: Model is trained on huge data set to do generic task
  ▶ Fine-Tuning: We continue training the model, but on the task we are actually interested in (!)

# Key Idea 3: Training Process Split into two Phases
BERT Training Tasks

Masked Language Modeling (MLM)

- ▶ Sentence-wise
- ▶ 15% of the tokens are »masked« by a special token
- ▶ Model predicts these, having access to all other tokens

# Key Idea 3: Training Process Split into two Phases
BERT Training Tasks

Masked Language Modeling (MLM)

▶ Sentence-wise

▶ 15% of the tokens are »masked« by a special token

▶ Model predicts these, having access to all other tokens

Next sentence prediction (NSP)

▶ Two (masked) sentences are concatenated

▶ Model has to predict wether second sentence follows on the first or not

# Key Idea 4: Scale Up

▶ With the transformer recipe, many parameters have simply been scaled up



Figure: Statistics about NLP models (Sharir et al., 2020; Wikipedia)

# Key Idea 4: Scale Up

▶ With the transformer recipe, many parameters have simply been scaled up



Figure: Statistics about NLP models (Sharir et al., 2020; Wikipedia)

# Key Idea 4: Scale Up

Large Numbers are Complicated

| Kurze Skala | | Lange Skala | | | Zehner-potenz | Vorsätze |
|---|---|---|---|---|---|---|
| Name | Systematik | Chuquet | mit -arde | Systematik | | |
| [Einheit] | Tausend$^{1-1}$ | [Einheit] | [Einheit] | Million$^0$ | $10^0$ | [Einheit] |
| Tausend | Tausend$^{1+0}$ | Tausend | Tausend | Million$^{½}$ | $10^3$ | Kilo |
| Million | Tausend$^{1+1}$ | Million | Million | Million$^1$ | $10^6$ | Mega |
| **Bil**lion | Tausend$^{1+2}$ | Tausend Millionen | Milli**arde** | Million$^{1½}$ | $10^9$ | Giga |
| **Tri**llion | Tausend$^{1+3}$ | Billion | **Bi**llion | Million$^2$ | $10^{12}$ | Tera |
| **Quadr**illion | Tausend$^{1+4}$ | Tausend Billionen | Billi**arde** | Million$^{2½}$ | $10^{15}$ | Peta |
| **Quint**illion | Tausend$^{1+5}$ | Trillion | **Tri**llion | Million$^3$ | $10^{18}$ | Exa |
| **Sext**illion | Tausend$^{1+6}$ | Tausend Trillionen | Trilli**arde** | Million$^{3½}$ | $10^{21}$ | Zetta |
| **Sept**illion | Tausend$^{1+7}$ | Quadrillion | **Quadr**illion | Million$^4$ | $10^{24}$ | Yotta |

# Key Ideas

- ▶ Input representation trained
- ▶ Attention to identify relevant tokens
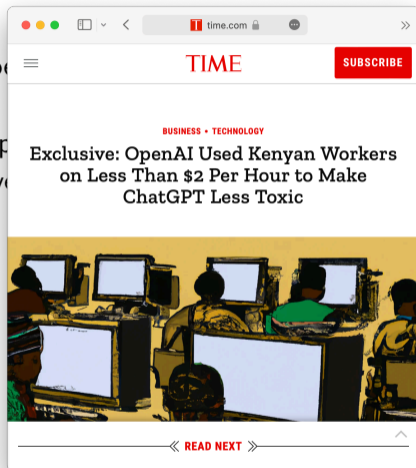- ▶ Two phases for training processes
- ▶ Scale up

# ChatGPT

- ▶ »OpenAI« is just a name – nothing this company does is ›open‹
  - ▶ I.e., we don't know many details
- ▶ Running ChatGPT is expensive (rumors: $ 100 000 per day)
  - ▶ Usually, running a service costs a fraction of the development/training cost

# ChatGPT

- ▶ »OpenAI« is just a name – nothing this company does is ›open‹
  - ▶ I.e., we don't know many details
- ▶ Running ChatGPT is expensive (rumors: $ 100 000 per day)
  - ▶ Usually, running a service costs a fraction of the development/training cost
- ▶ There are multiple ugly sides

# ChatGPT

- »OpenAI« is just a name – nothing this company do[es]
  - I.e., we don't know many details
- Running ChatGPT is expensive (rumors: $ 100 000 p[er]
  - Usually, running a service costs a fraction of the dev[...]
- There are multiple ugly sides



time.com

TIME

SUBSCRIBE

BUSINESS • TECHNOLOGY

**Exclusive: OpenAI Used Kenyan Workers on Less Than $2 Per Hour to Make ChatGPT Less Toxic**

≪ READ NEXT ≫

# ChatGPT

# ChatGPT

- ▶ »OpenAI« is just a name – nothing this company does is ›open‹
  - ▶ I.e., we don't know many details
- ▶ Running ChatGPT is expensive (rumors: $ 100 000 per day)
  - ▶ Usually, running a service costs a fraction of the development/training cost
- ▶ There are multiple ugly sides

## ChatGPT predicts probable next words

- ▶ There is no model of the world behind it
- ▶ There is no factual knowledge or reasoning about anything behind it
- ▶ No one is able to guarantee, that the produced text is factually correct

# Discussion

Do we need legal regulation of »AI«, and if so, what exactly?

Section 1

Summary

# Summary

► Summary