

Natural Language Processing

Salar Mohtaj | DFKI

Natural Language Processing

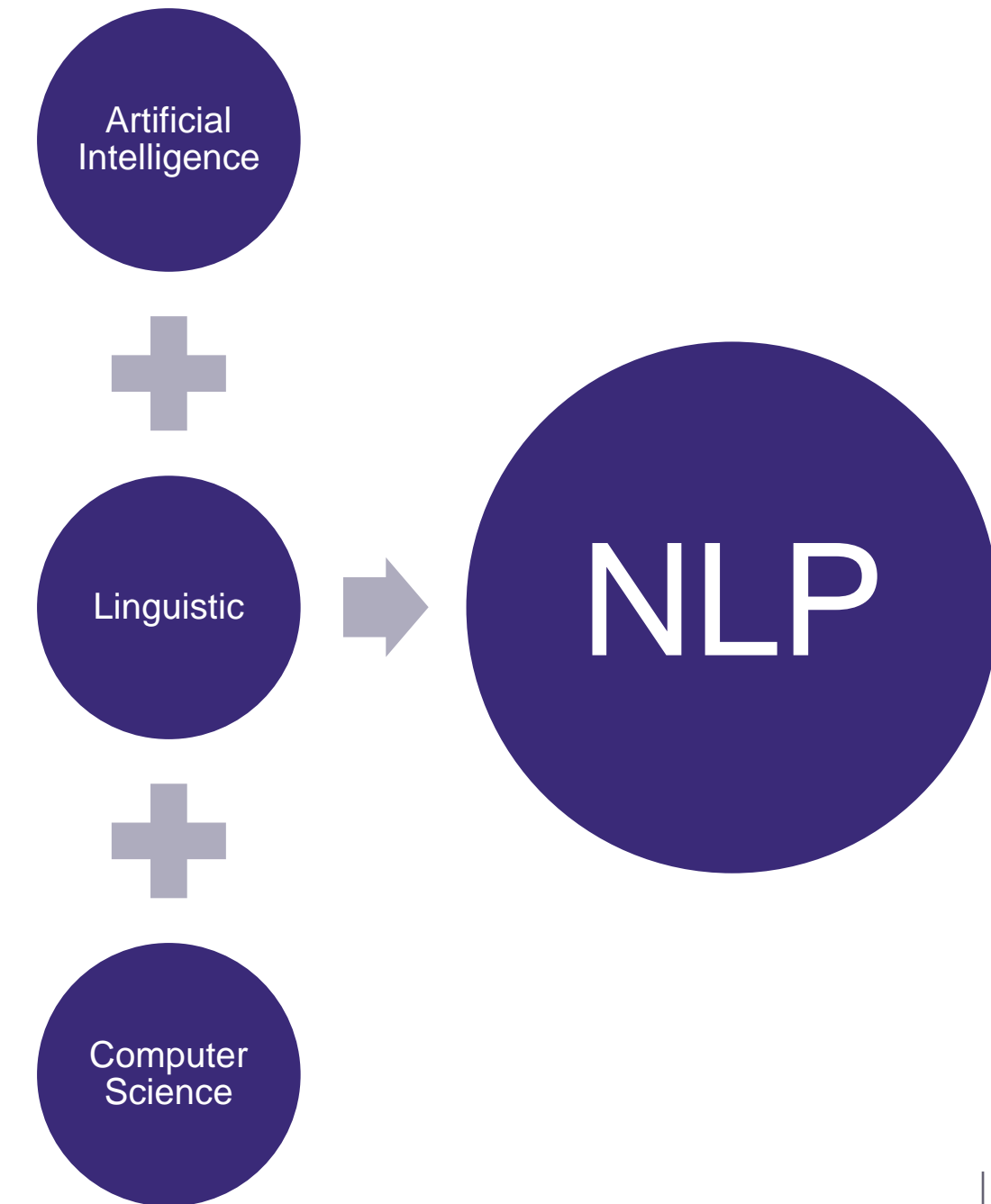
- Introduction to Natural Language Processing
- Everyday NLP applications
- Main NLP tasks
- Main approaches in NLP
- Who is this course for
- The Course structure
- NLP terminology

Natural Language Processing

- Introduction to Natural Language Processing
- Everyday NLP applications
- Main NLP tasks
- Main approaches in NLP
- Who is this course for
- The Course structure
- NLP terminology

Introduction to Natural Language Processing

- Natural language processing (NLP) is a branch of ***artificial intelligence*** that helps computers to understand, interpret and generate human language
- ***Natural language processing*** helps computers communicate with humans in their own language
- Most NLP techniques rely on machine learning to derive ***meaning*** from ***human languages***



Why is it difficult?

The hammer hit the glass and it broke!



Why is it difficult?

I saw someone on the hill with a telescope!



Images from www.thedailychain.com and www.storyblocks.com

Example from www.examples.yourdictionary.com

Why is it difficult?

- Ambiguity in language
 - The rules that dictate the passing of information using natural languages are not easy for computers to understand
 - Sarcastic remark

That's just what I needed today!

Why is it difficult?

- Ambiguity in language
 - The rules that dictate the passing of information using natural languages are not easy for computers to understand
 - Sarcastic remark
 - Multi meaning words

She will **park** the car so we can walk in the **park**.

The committee **chair** sat in the center **chair**.

Why is it difficult?

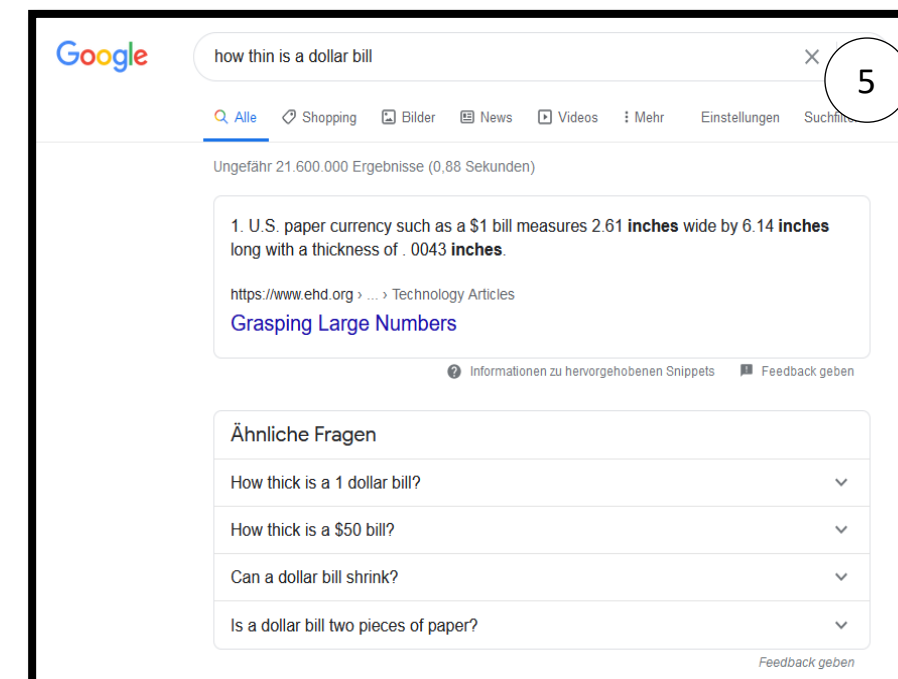
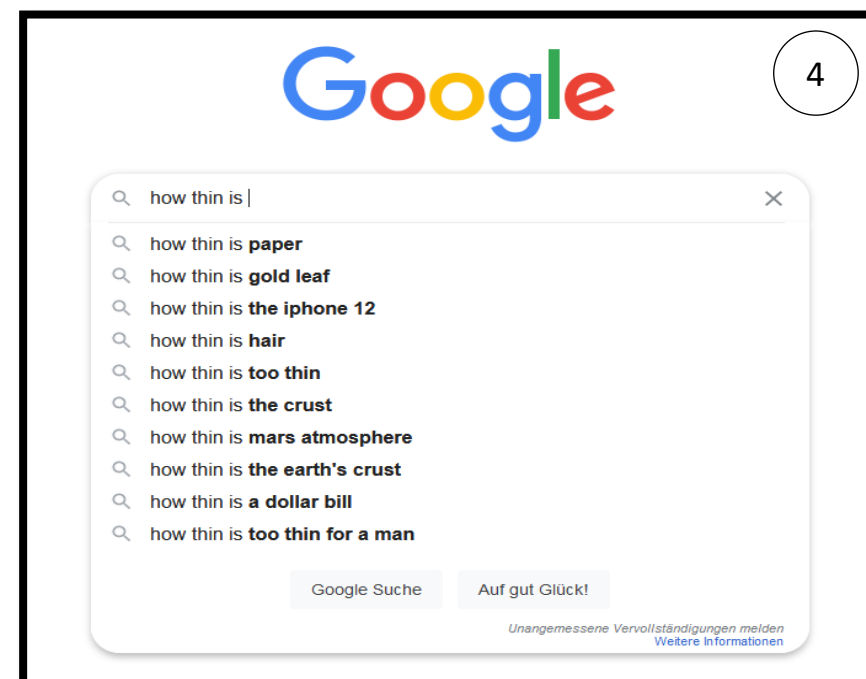
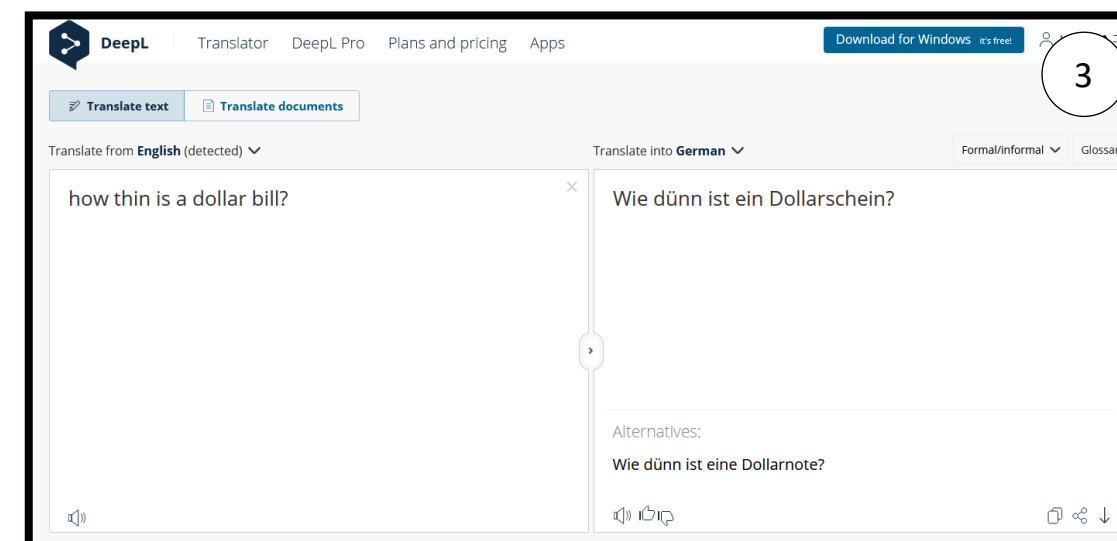
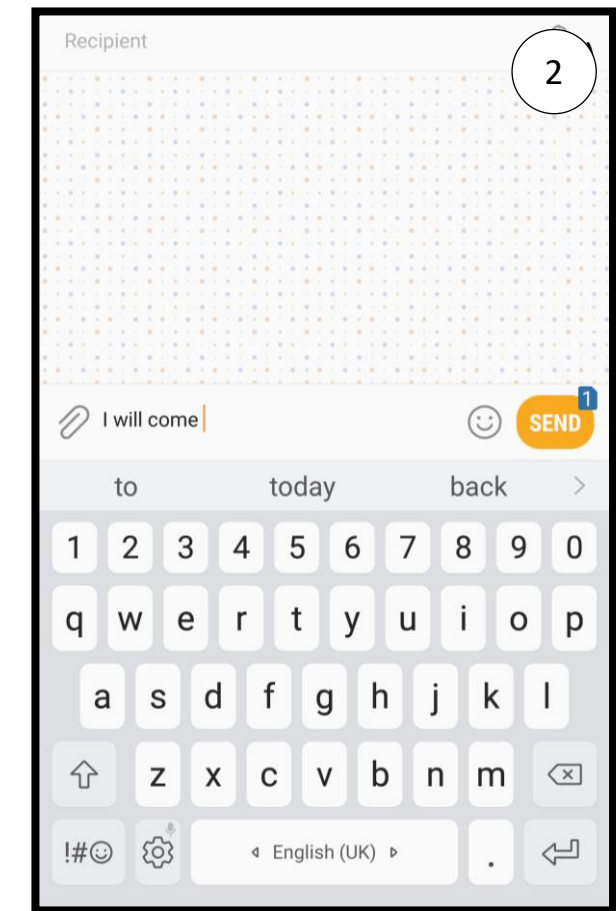
- Ambiguity in language
 - The rules that dictate the passing of information using natural languages are not easy for computers to understand
 - Sarcastic remark
 - Multi meaning words
- The lexicon of a language is usually enormous
 - Oxford dictionary has **273,000** headwords; **171,476** of them being in current use
 - An average person has a vocabulary range of about **20,000** to **35,000**

Natural Language Processing

- Introduction to Natural Language Processing
- **Everyday NLP applications**
- Main NLP tasks
- Main approaches in NLP
- Who is this course for
- The Course structure
- NLP terminology

Everyday NLP applications

- Email filters (spam detection) ⁽¹⁾
- Faster typing ⁽²⁾
- Language translation ⁽³⁾
- Question answering ⁽⁴⁾⁺⁽⁵⁾
- Smart assistant devices ⁽⁶⁾

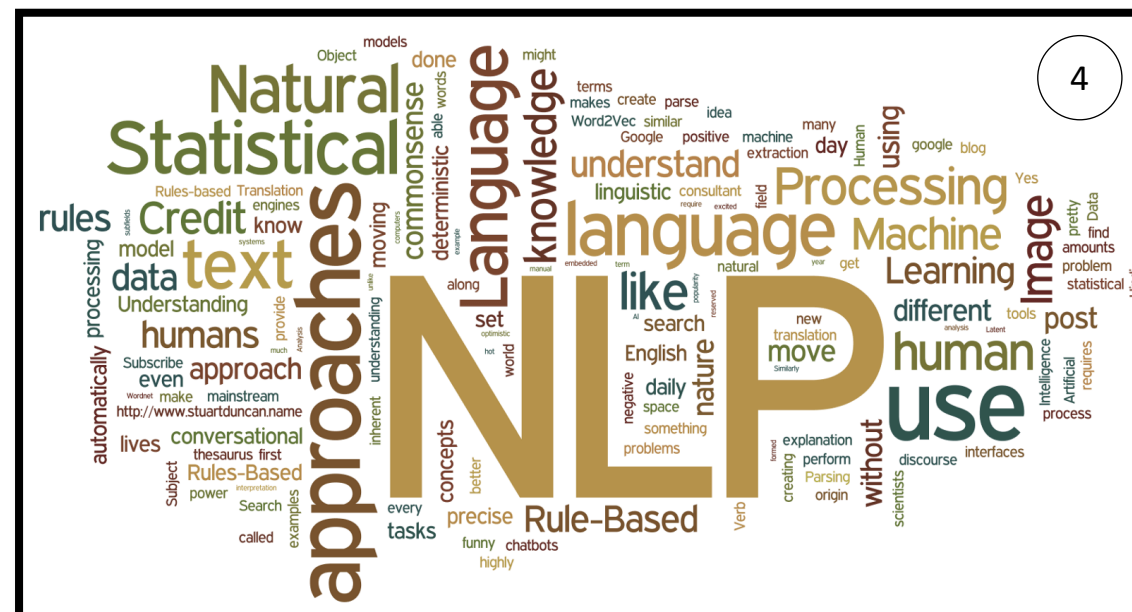
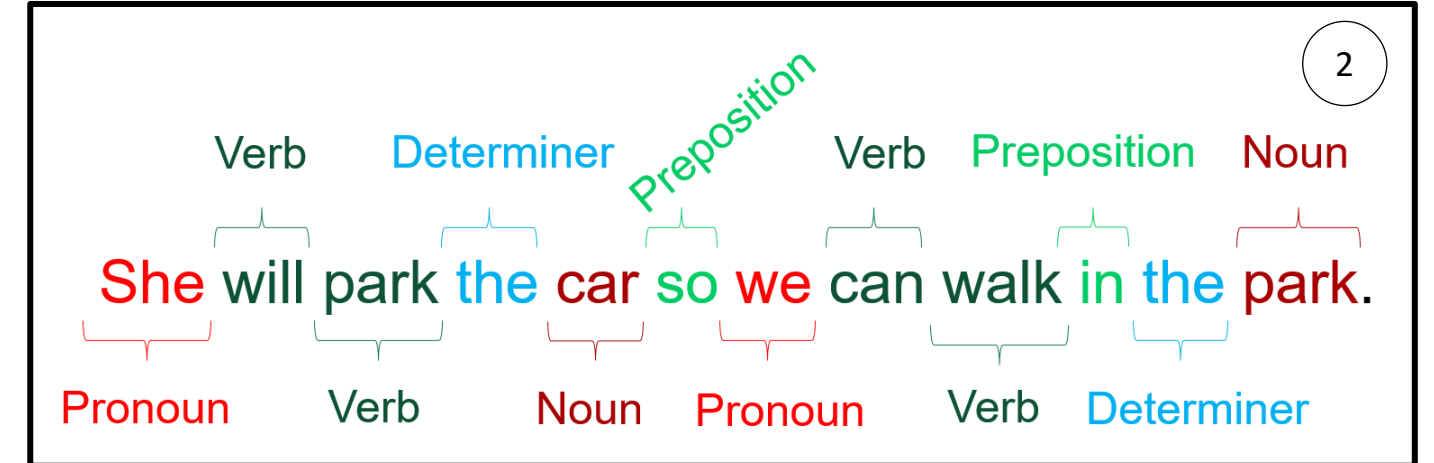


Natural Language Processing

- Introduction to Natural Language Processing
- Everyday NLP applications
- **Main NLP tasks**
- Main approaches in NLP
- Who is this course for
- The Course structure
- NLP terminology

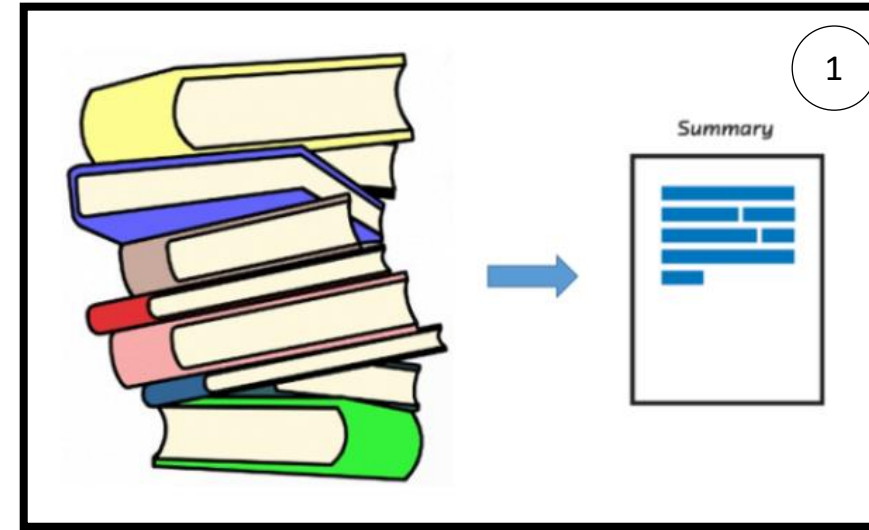
The main NLP tasks

- Text classification (1)
- Parts of speech tagging (2)
- Sentiment analysis (3)
- Keyword extraction (4)
- Text similarity (5)



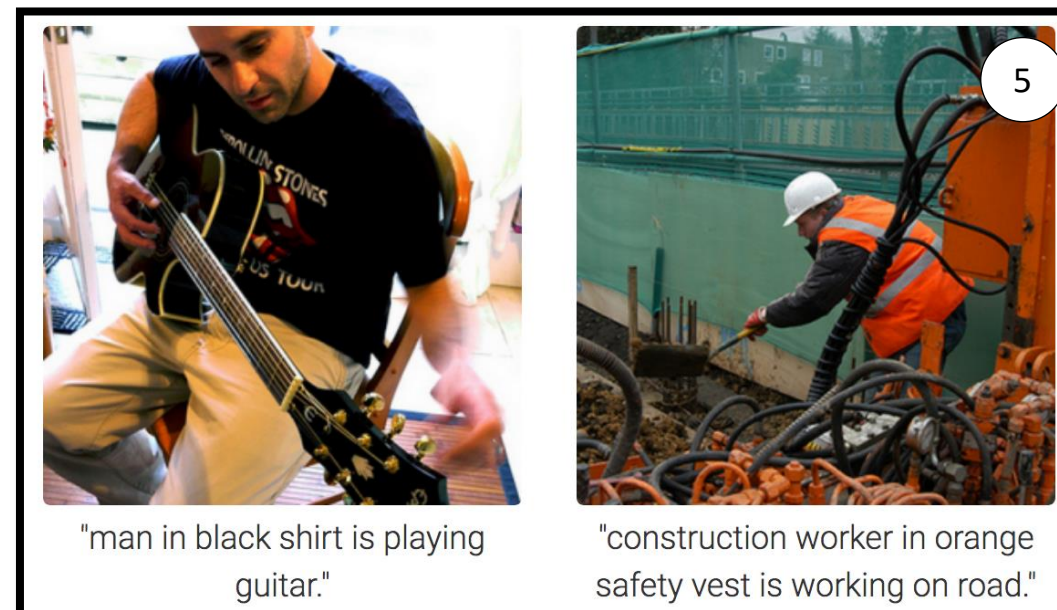
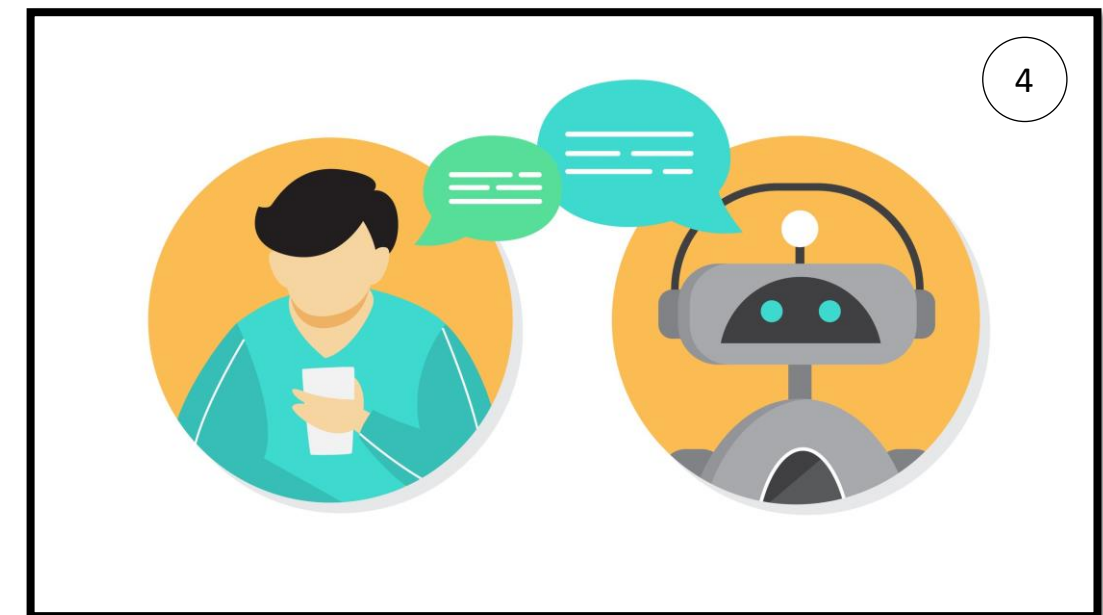
The main NLP tasks

- Text summarization ⁽¹⁾
- Named entity recognition (NER) ⁽²⁾
- Machine translation ⁽³⁾
- Question answering ⁽⁴⁾
- Image captioning ⁽⁵⁾



When **Sebastian Thrun** PERSON started at **Google** ORG in **2007** DATE, few people outside of the company took him seriously. "I can tell you very senior CEOs of major **American** NORP car companies would shake my hand and turn away because I wasn't worth talking to," said **Thrun** PERSON, now the co-founder and CEO of online higher education startup Udacity, in an interview with **Recode** ORG **earlier this week** DATE.

A little **less than a decade later** DATE, dozens of self-driving startups have cropped up while automakers around the world clamor, wallet in hand, to secure their place in the fast-moving world of fully automated transportation.



Natural Language Processing

- Introduction to Natural Language Processing
- Everyday NLP applications
- Main NLP tasks
- **Main approaches in NLP**
- Who is this course for
- The Course structure
- NLP terminology

Main Approaches in NLP

- Rule based approaches
- Classical machine learning
- Deep learning

Main Approaches in NLP

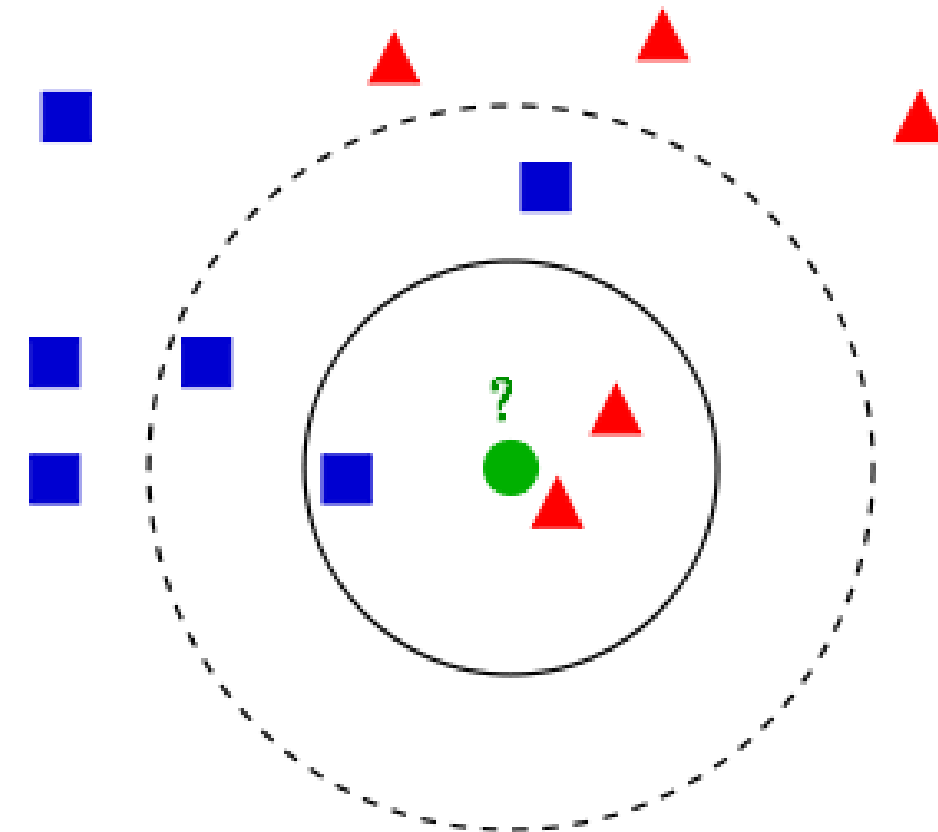
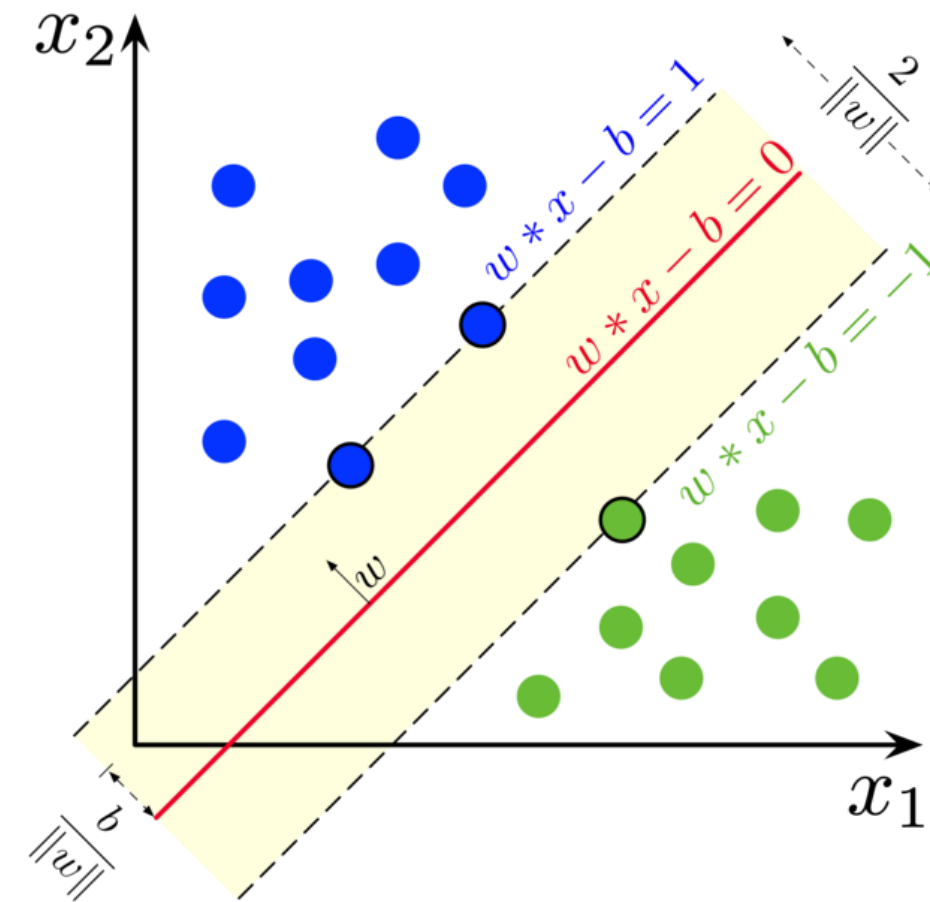
- Rule based approaches
 - Lack of enough accuracy

The film was **good** not **bad**.



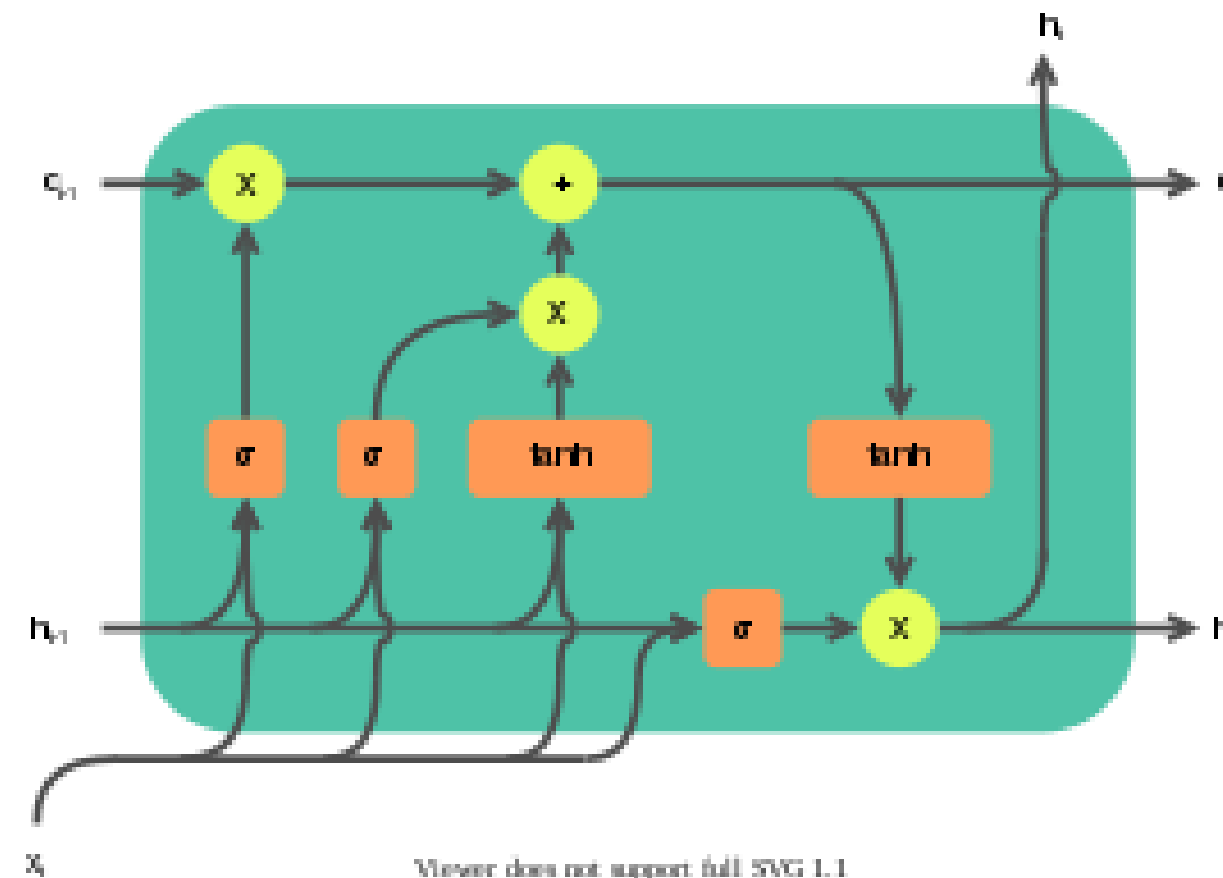
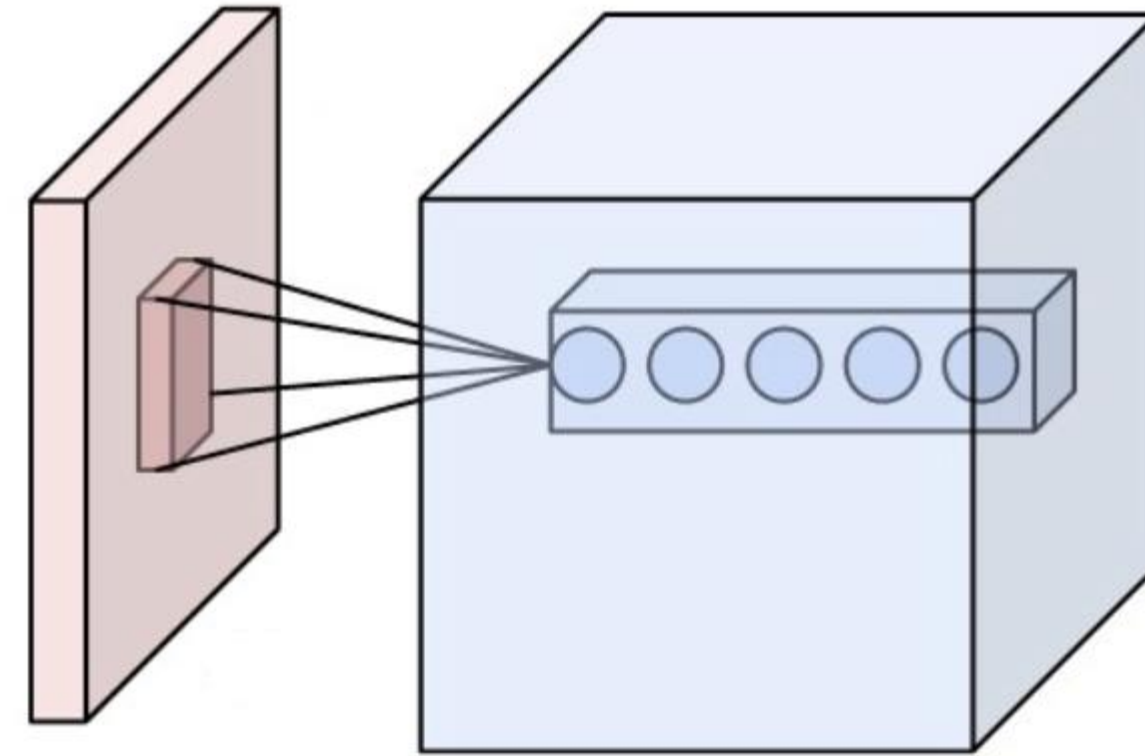
Main Approaches in NLP

- Rule based approaches
 - Lack of enough accuracy
- Classical machine learning
 - Training data
 - Feature engineering
 - Training a model



Main Approaches in NLP

- Rule based approaches
 - Lack of enough accuracy
- Classical machine learning
 - Training data
 - Feature engineering
 - Training a model
- Deep learning
 - More training data
 - Feature engineering is skipped
 - Training a model



Natural Language Processing

- Introduction to Natural Language Processing
- Everyday NLP applications
- Main NLP tasks
- Main approaches in NLP
- **Who is this course for**
- The Course structure
- NLP terminology

Who is this course for?

- Those who
 - don't want use NLP models as a black box
 - Review the state-of-the-art approaches
 - want to gain intuition for problem solving
 - you're asked to develop your own translation tool, what is the best approach
 - have a prior background on machine learning and deep learning



Supervised machine learning

Hidden layer

Backpropagation

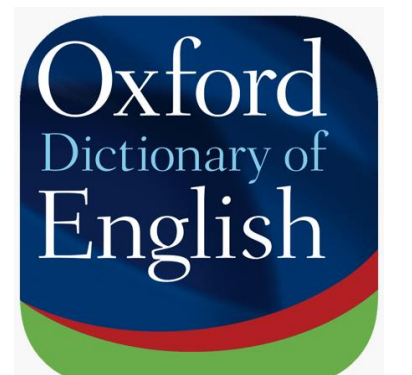
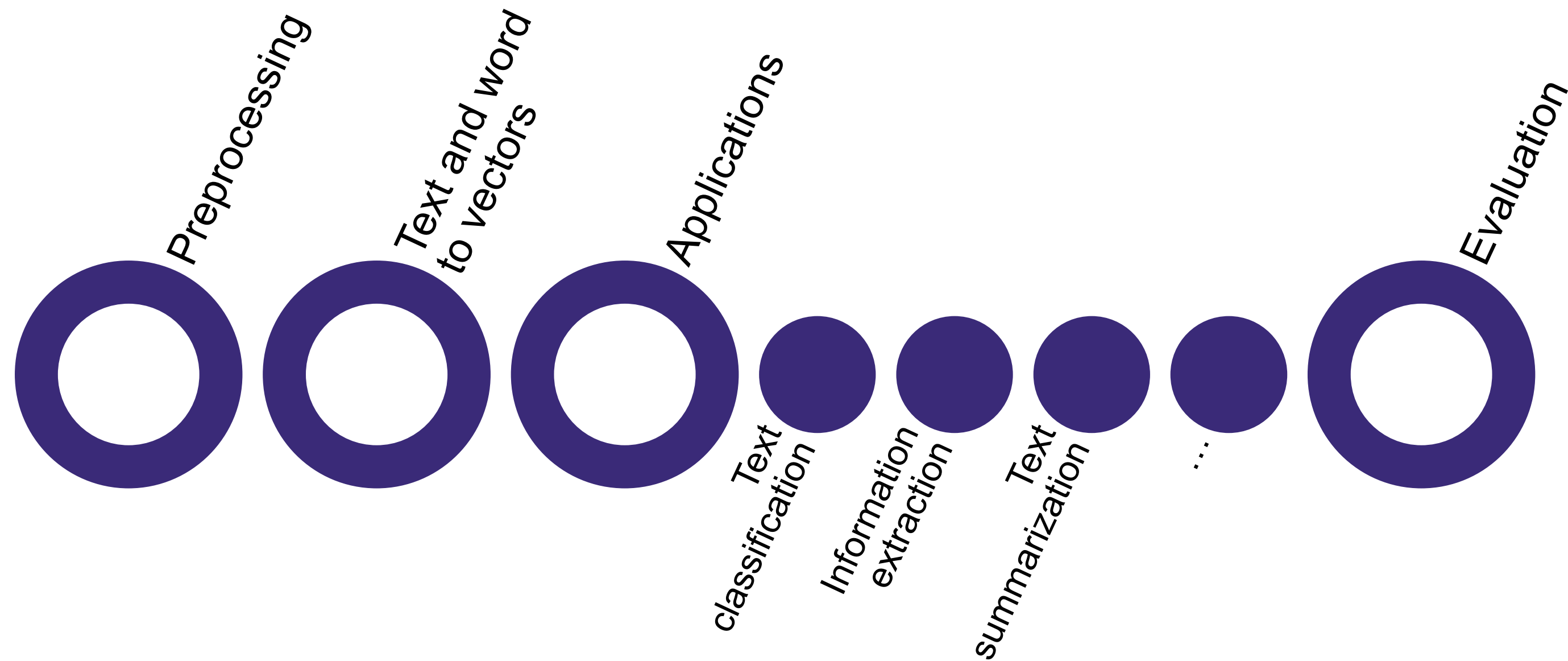
softmax

LSTM

Natural Language Processing

- Introduction to Natural Language Processing
- Everyday NLP applications
- Main NLP tasks
- Main approaches in NLP
- Who is this course for
- **The Course structure**
- NLP terminology

The Course structure

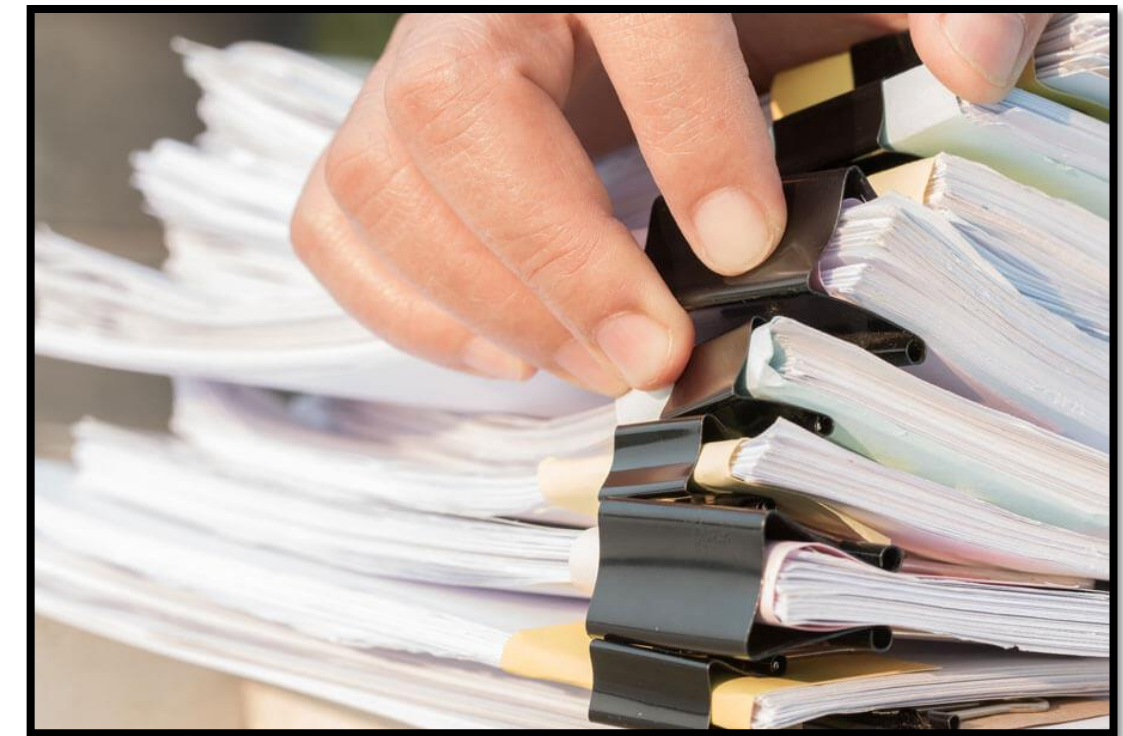


Natural Language Processing

- Introduction to Natural Language Processing
- Everyday NLP applications
- Main NLP tasks
- Main approaches in NLP
- Who is this course for
- The Course structure
- **NLP terminology**

NLP terminology

- Corpus (Plural: corpora)
 - A collection of text, usually contains several documents
 - Wikipedia articles
 - Collection of movies reviews
 - Internet comments
 - Collection of tweets
 - Corpora can be in a single language or multiple languages



NLP terminology

- Document
 - Document refers to a body of text in a corpus
 - A tweet in a twitter corpus
 - An email in a collection of emails
- Stop word
 - usually refers to the most common words in a language
 - Words like “**the**”, “**and**”, “**a**”, “**an**”, “**in**”

NLP terminology

- Vocabulary
 - The set of unique words used in the text corpus
 - Set of unique words which are used in all Wikipedia articles
- Out of Vocabulary (OOV)
 - Words that have not seen during the train, but in the test
 - We will encounter out of vocabulary terms when using our model for inference

Thank you!

„KI-Campus – Die Lernplattform für Künstliche Intelligenz“ ist ein Projekt von



www.ki-campus.org