

Anonymization of German Legal Texts

Bachelor Thesis - Kickoff

Tom Schamberger, Kickoff 18th of June 2019

Chair of Software Engineering for Business Information Systems (sebis)
Faculty of Informatics
Technische Universität München
www.matthes.in.tum.de

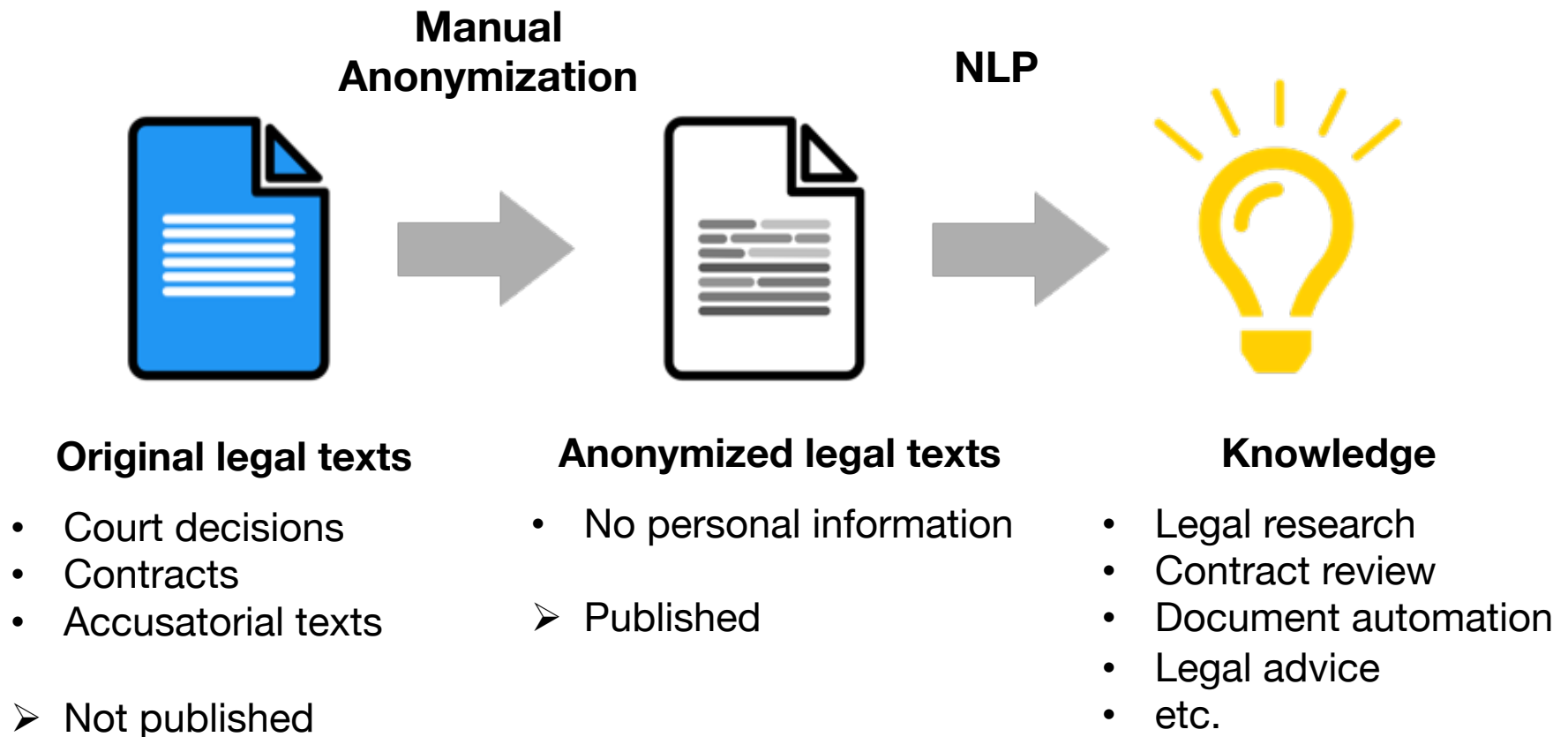
Motivation

Problem Statement

Approach

Research Questions

Schedule



Motivation

Problem Statement

Approach

Research Questions

Schedule

Example for court decisions:

Anonymization

6. Den Angeboten auf der Webseite ... liegt das Konzept der Produktdetailseite zugrunde. Dabei wird für jedes über die ...-Plattform angebotene Produkt jeweils nur eine Produktdetailseite angezeigt; jedes Produkt enthält eine spezifische ...-Produktidentifikationsnummer (...) zugewiesen.

- **Expensive manual** anonymization process
 - Leads to rare publications of legal texts
 - Results in few data sets
- Development of an **automated anonymization process**
- Only anonymized data sets available
 - Anonymization training **without non-anonymized data sets**

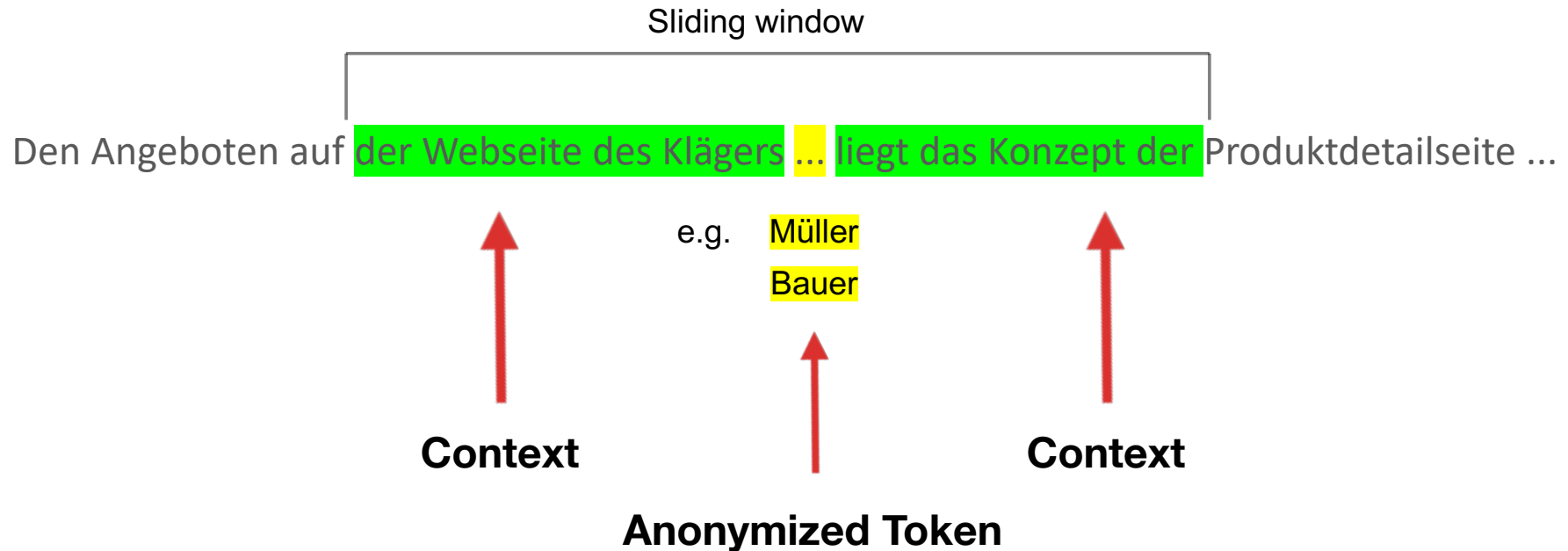
Motivation

Problem Statement

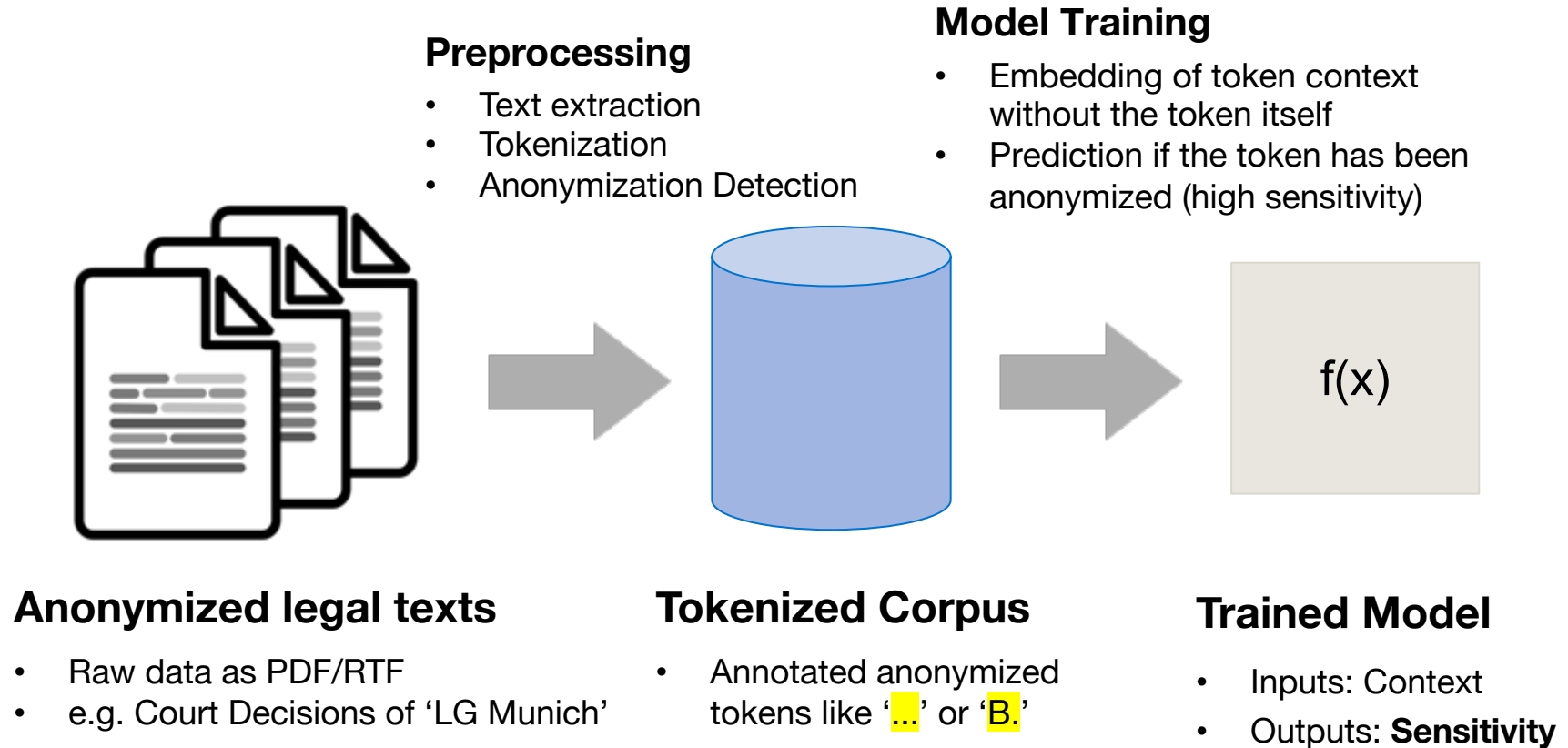
Approach

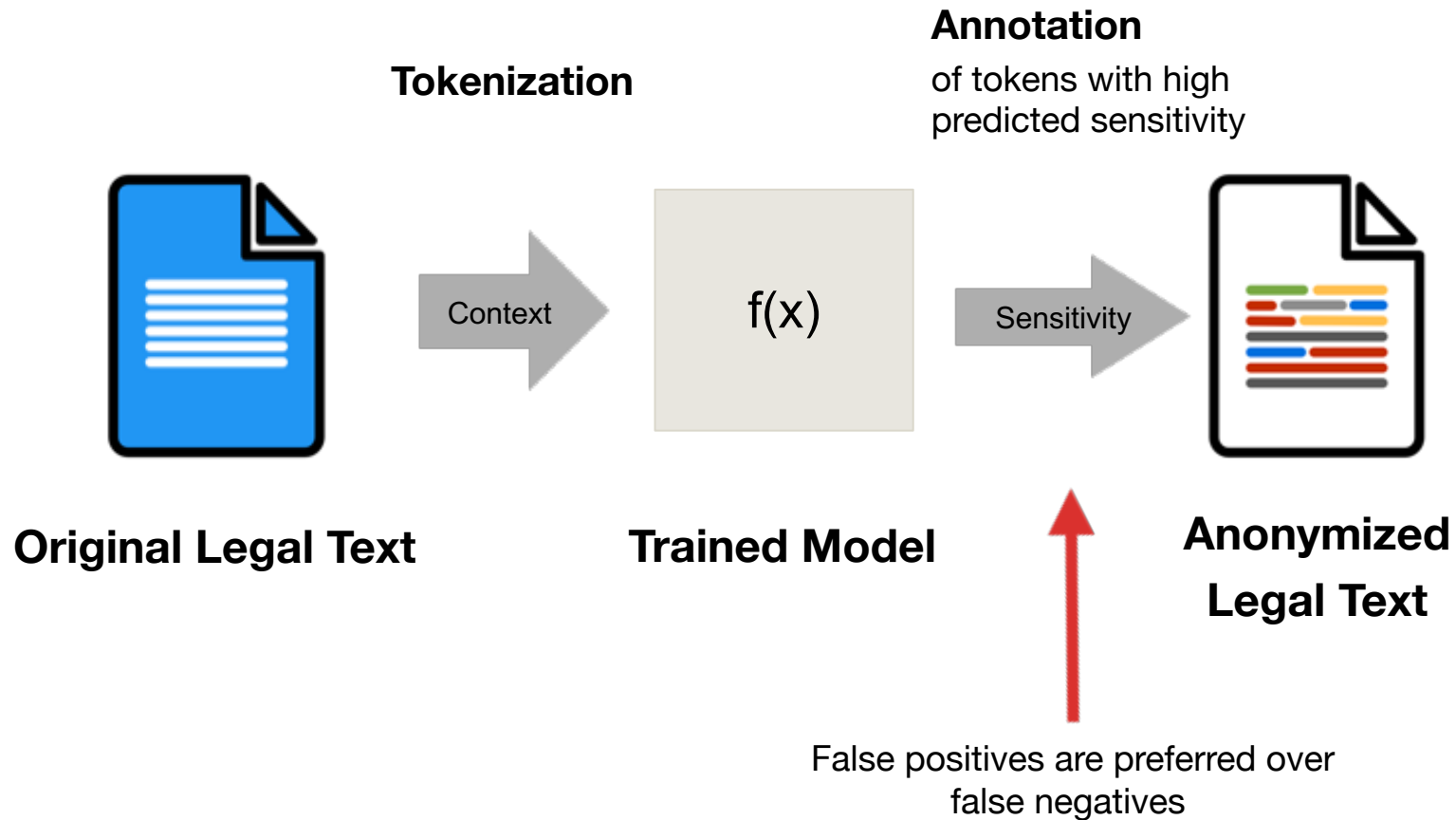
Research Questions

Schedule



- Statement: **Sensitivity** of token **depends only on context**, not the actual content itself
- Anonymized token in legal texts are annotated (e.g. by "...") and can be detected
- Anonymization model may be **trained using anonymized data**





1. Literature research

- existing suitable NLP models
- existing frameworks and implementations

2. NLP model training on anonymized court decisions

- Fetching and preprocessing of [court decision documents](#)
- Model training and hyperparameter fitting

3. Evaluation of the model

- using [court decision examples](#)
- using rewritten anonymized legal texts (e.g. using randomized NE replacement)

4. Application

- Implementation of the anonymization process using trained NLP model

Motivation

Problem Statement

Approach

Research Questions

Schedule

Is the textual context of a token within legal texts enough to predict the sensitivity of this token?

- Classical anonymization solutions in other domains additionally use NER (named entity recognition)
 - Improves stability
 - But there is no non-anonymized data
- Additional data, which may be taken into account:
 - the **commonness** of words in legal texts
 - e.g. words like *mit* or *Kläger* probably have a low sensitivity
 - **document metadata** such as title, court, etc.

How can placeholders be detected in anonymized legal documents?

- Examples for **possible anonymization patterns** for the word 'Schamberger':
 - ...
 - S...
 - S...
 - "S"
 - "S..."
 - s
 - S

Which machine learning approaches be used to automate anonymization using only anonymized data?

- Examples for contextual embeddings: **BERT**, ELMO, GPU, CMU
- Alternative classical embedding: Word2vec, GloVe
- Possible architectures: Convolution, RNN, Attention
- Possible sliding window sizes:
 - **Fixed number of words**
 - Full sentences
 - Full paragraphs

Motivation

Problem Statement

Approach

Research Questions

Schedule

Schedule

	June	July	August	September	October
Literature Research	Start – 2nd of August				
Implementation		1st of July – 13th of September			
Model Training		15th of July – 4th of October			
Evaluation			5th of August – 4th of October		
Writing			5th of August – End		
Review				16th of Sep. – End	



Tom Schamberger

Technische Universität München
Faculty of Informatics
Chair of Software Engineering for
Business Information Systems

Boltzmannstraße 3
85748 Garching bei München

Tel +49.89.289.17132

Fax +49.89.289.17136

matthes@in.tum.de
www.matthes.in.tum.de

