# Specification and Extraction of Semantic Patterns in German Laws based on Linguistics Features using Apache Ruta

Patrick Ruoff, June 27th, 2016

Software Engineering für betriebliche Informationssysteme (sebis)
Fakultät für Informatik
Technische Universität München

wwwmatthes.in.tum.de

# Agenda

**sebis**

- Processes of Legal Experts (Scientists and Lawyers) are…
  - ... time-intensive
  - ... knowledge-intensive
  - ... data-intensive

- Legal Data Science is becoming more and more attractive, because
  - ... process time and memory space are cheap
  - ... algorithms can process data fast and accurate

- In order to achieve highest accuracy, algorithms and data models need an adaption to the domain
  - German legal texts (laws, contracts, etc.)
  - Data model tailored to legislative data

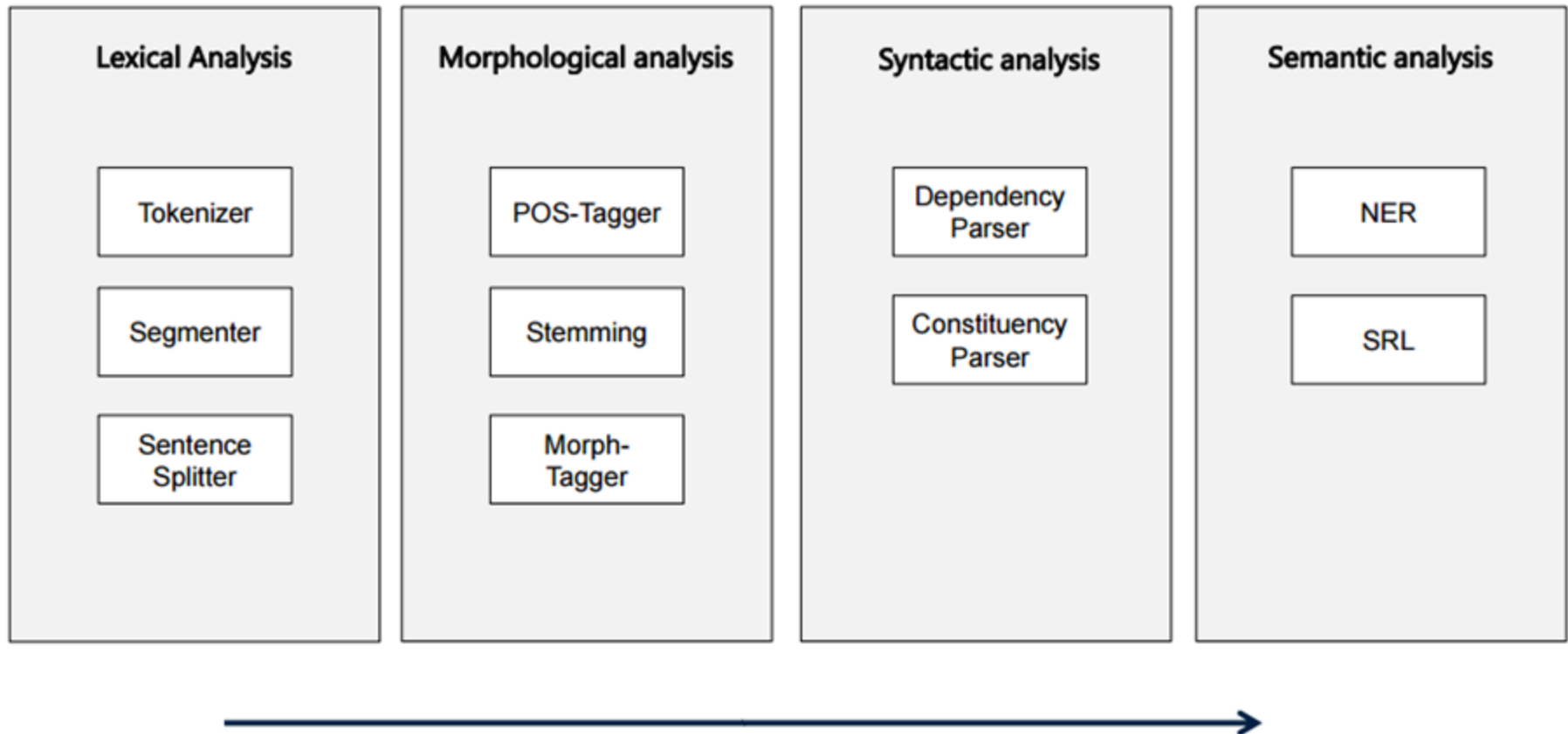    $\rightarrow$ There are still more Patterns that are not yet implemented

1. Maat and Winkels (2010)

   - Classification of norms regarding linguistic structures

   - Regular Expressions

   - Limitation: no consideration of linguistic properties, such as nouns, etc.

2. Bommarito and Katz (2014)

   - Analysis of semantic and structural properties

3. Grabmair et al. (2015)

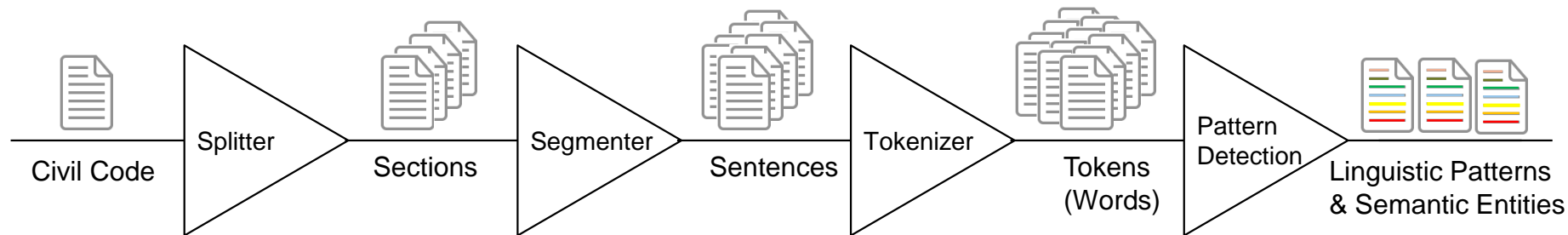   - Using Apache UIMA for legal text analysis (LUIMA)

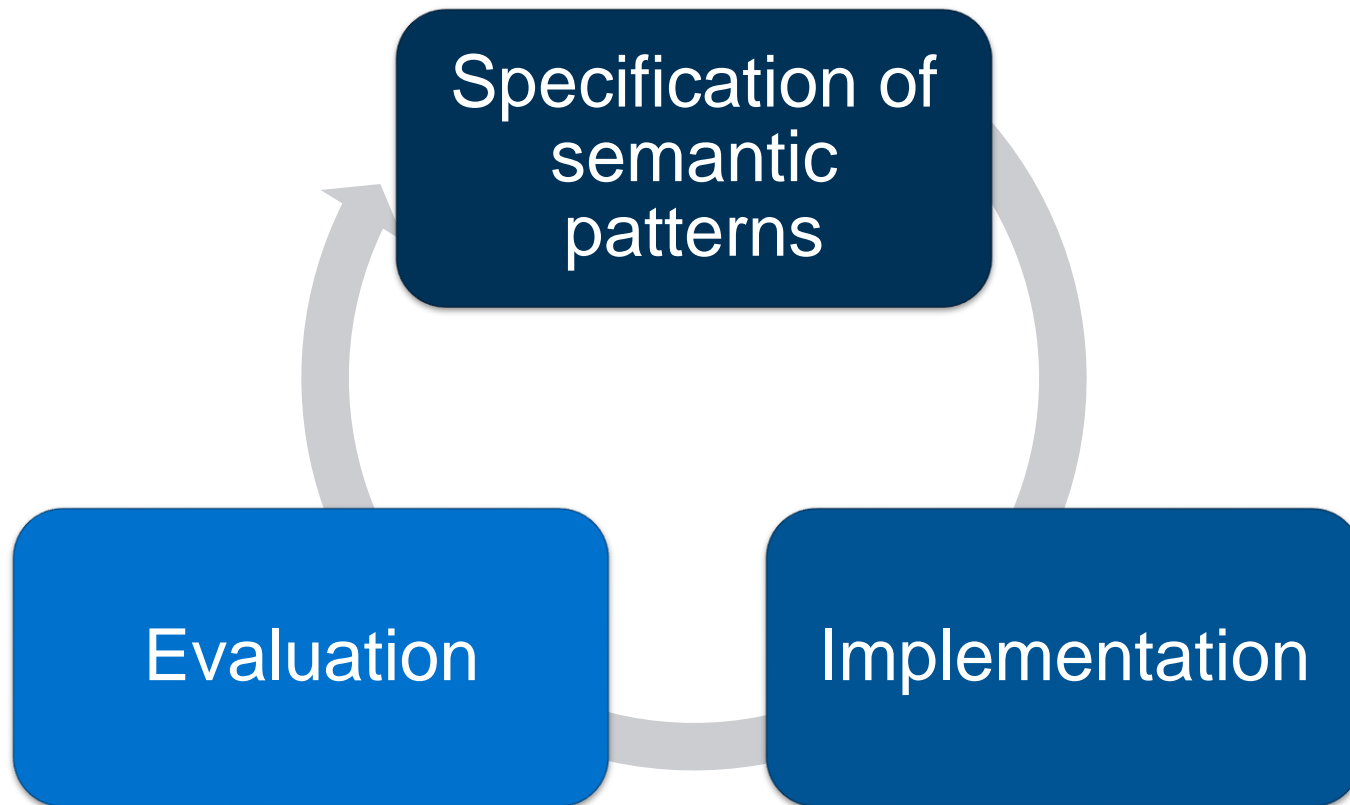| **Short summary** |
|---|
| - Data analysis is well established in legal informatics<br>- Adaption to domain is crucial to achieve highest accuracy …<br>    - ... data model<br>    - ... algorithms<br>- Reusable code and implementations to avoid re-inventions |

**UIMA (Unstructured Information Management Architecture)**

- A common software architecture for text mining/processing
  - Alternatives: GATE, NLTK, etc.
  - Base line for IBM Watson
- Pipes & Filters architecture



Civil Code → Splitter → Sections → Segmenter → Sentences → Tokenizer → Tokens (Words) → Pattern Detection → Linguistic Patterns & Semantic Entities

- Thread-safe (usage in a web application with multiple users/requests)
- Apache Ruta for complex pattern specification engine
  - Analogy: Jape grammar (GATE)

Specification of semantic patterns

Evaluation

Implementation

**Specification of semantic patterns**

➢ In Cooperation with Konrad Heßler from the Juristic Faculty of the LMU
➢ Introduction into juristic work needed

**Implementation**

➢ Ruta scripts
➢ Executed on Lexia

**Evaluation**

➢ Testing on German laws
➢ Comparison with handmade Classifications

# Code Example

PatrickO1.ruta

```
1  // Import types
2  IMPORT PACKAGE de.tudarmstadt.ukp.dkpro.core.api.lexmorph.type.pos FROM GeneratedDKProCoreTypes AS pos;
3
4  DECLARE LegalConcept;
5  ((pos.ADJ)[2,3] pos.N) {-> LegalConcept};
6
7  // Annotate the sentence being a legal concept as LegalConcept
8  DECLARE LegalConceptSentence;
9  Sentence{CONTAINS(Patrick01.LegalConcept) -> Patrick01.LegalConceptSentence};
```

sebis

# Roadmap

| June | July | August | September | October |
|------|------|--------|-----------|---------|

Literature Research

Concept

Specification

Evaluation

Implementation

Evaluation

Writing the thesis

# Thank you for your attention!

**Patrick Ruoff**

**sebis**

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

Boltzmannstraße 3
85748 Garching bei München

Tel    +49.89.289.17124
Fax   +49.89.289.17136

ga54kuc@mytum.de
wwwmatthes.in.tum.de

```
 1  // Basic linguistic vocabulary
 2  DECLARE ISDG;
 3  "im Sinne dieses Gesetzes" -> LDSache.ISDG;
 4  "im Sinne des Gesetzes" -> LDSache.ISDG;
 5
 6  DECLARE IST;
 7  "ist|sind" -> LD.IST;
 8
 9  DECLARE NEG;
10  "keine|kein|nicht" -> LD.NEG;
11
12  DECLARE LDIdentifier; // Declare the indicator for legal definitions
13  DECLARE LegalEntity; // Declare the legally defined entity
14  DECLARE LegalDefinition; // Declare the legal definition
15
16  // Definition of linguistic patterns and rules
17  // {{ADJ}} {{NOUN}} im Sinne dieses|des Gesetzes ist {{Phrase}}
18  ((pos.N? pos.N) {-> LD.LegalEntity} LD.ISDG) {-> LD.LDIdentifier};
19  ((pos.ADJ+ pos.N) {-> LD.LegalEntity} LD.ISDG) {-> LD.LDIdentifier};
20
21  // {{NOUN}} ist kein {{NOUN}}
22  (pos.N {-> LD.LegalEntity} LD.IST LD.NEG pos.N) {-> LD.LDIdentifier};
23  (pos.N{-PARTOF(LD.LegalEntity) -> LD.LegalEntity} LD.ISDG){->LD.LDIdentifier};
24
25  // Annotate the sentence being a legal definition as LegalDefinition
26  Sentence{CONTAINS(LD.LDIdentifier) -> LD.LegalDefinition};
27
28  // Remove temporary annotations
29  LD.IST {-> UNMARK(LD.IST)};
30  LD.NEG {-> UNMARK(LD.NEG)};
31  LD.ISDG {-> UNMARK(LD.ISDG)};
32  LD.LDIdentifier{-> UNMARK(LD.LDIdentifier)};
```

*Listing 1: Linguistic pattern descriptions (LD.ruta) for the semantic entity Legal Definition using Apache Ruta*