



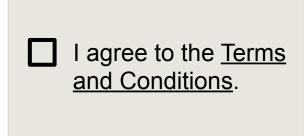
- 1. Motivation
- 2. Research Question
 - Context in AGB Check
- 3. Solution Approaches
- 4. Timeline

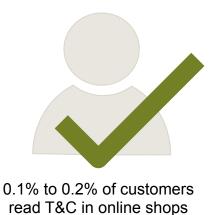
Motivation



Whom of you did ever read the <u>Terms and Conditions</u> when buying products online?

- no standard solutions for parts of the NLP pipe (structured extraction from websites)
- structured extraction is needed in legal documents to allow the drawing of conclusions later in the pipe
 - no existing solutions for structured & hierarchical extraction of T&C





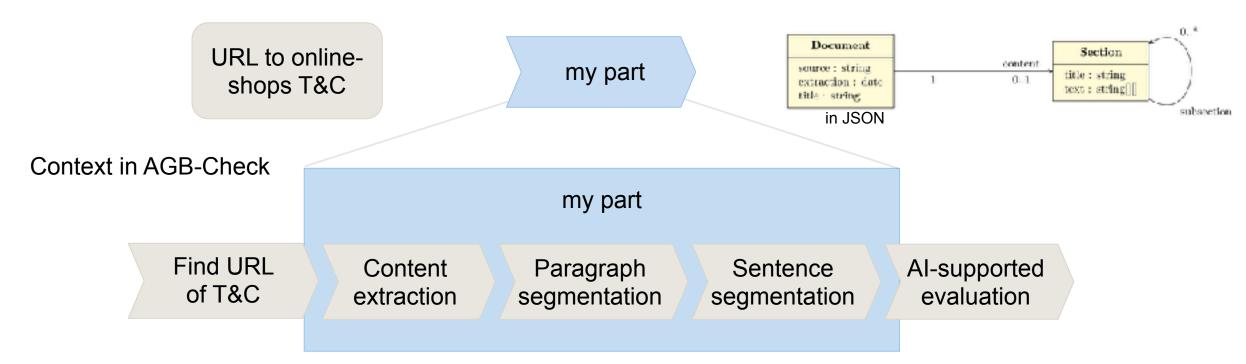


- 1. Motivation
- 2. Research Question
 - Context in AGB Check
- 3. Solution Approaches
- 4. Timeline

Research Question (Context in AGB-Check)



- RQ1: What are special requirements for the extraction of contracts in comparison with regular content?
- RQ2: How to extract the relevant parts from contracts in raw html?
- RQ3: How to handle certain characteristic elements (avoid classification as main content)?
- RQ4: How to extract structure and hierarchy of paragraphs, (sub-)titles and related clauses?
 - Which existing HTML and CSS properties provide useful information and can be used to structure the content?



For more information visit <u>AGB-Check</u> (http://www.matthes.in.tum.de)

Research Question (Context in AGB-Check)



Example



Terms and Conditions of Aldi UK, 210325

"subsections": ["App", "subsections": [], "text": ["our", "title": "" "Privacy" "Notice" "which", "subsections": [], "sets", "text": ["out", "the", "our", "terms", "on", "Content", "which", "Management", "we", "Policy", "process" "any", "personal", "which", "data", "governs", "we", "how", "collect", "we", "from", "you", "deal" "with", "or", "content", "that", "which", "you", "provide", "you", "to", "submit", "us", "to", "us",], "You", "and" "consent", "to", "such", "processing", "title": "" "by", "using", "our", "Site", "subsections": [],

excerpt of parsed T&C



- 1. Motivation
- 2. Research Question
 - Context in AGB Check
- 3. Solution Approaches
- 4. Timeline

Solution Approaches



Content

Paragraph segmentation

Sentence segmentation

Selenium with driver to read entire page source code (+style)

Removing irrelevant content (based on tags, amount of text)

trafilatura, jusText, boilerpipe Segmenting paragraphs (based on fontsize [style], containers, enumerations)

SoMaJo shows good results for sentence segmentation & tokenization

NLTK, spaCy

- exist. solutions
- approaches
- exist. libraries

[SoMaJo]Proisl, T., Uhrig, P.: SoMaJo: State-of-the-art tokenization for German web and social media texts

[github]boilerpipe (Canola & Article); [github]trafilatura; [github]jusText; [github]selenium; [github]NLTK; [github]spaCy

Solution Approaches



Paragraph Content Sentence extraction segmentation segmentation getFontAndStyle Hierarchy Sentence WebExtraction Detection Segmentation getHTMLTree getParagraphs getContent MainContent Paragraph Detection Detection segmentSentences



- 1. Motivation
- 2. Research Question
 - Context in AGB Check
- 3. Solution Approaches
- 4. Timeline

Timeline



