

Outline



1. Motivation

- 2. Research Question
- 3. Content Extraction
 - Requirements
 - Implementation
 - Processing Demo

4. Hierarchy Extraction

- Requirements
- Implementation
- Processing Demo
- 5. Architecture & Additional Technology
- 6. Evaluation
 - Content Extraction
 - Hierarchy Extraction
- 7. Conclusion & Future Work

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

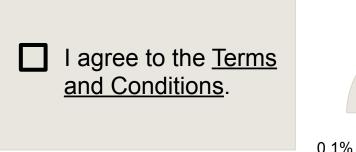


Figure 1.1: Motivation.



[Journal of Legal Studies] Bakos, Y., Marotta-Wurgler, F., Trossen, D.: Does Anyone Read the Fine Print? Consumer Attention to Standard-Form Contracts.

Outline



- 1. Motivation
- 2. Research Question
- 3. Content Extraction
 - Requirements
 - Implementation
 - Processing Demo
- 4. Hierarchy Extraction
 - Requirements
 - Implementation
 - Processing Demo
- 5. Architecture & Additional Technology
- 6. Evaluation
 - Content Extraction
 - Hierarchy Extraction
- 7. Conclusion & Future Work

Research Question



- RQ1: What are special requirements for the extraction of contracts in comparison with regular (i.e. news articles, blogposts) content?
- RQ2: How to extract the relevant parts from contracts in raw HTML?
- RQ3: How to extract the structure and the hierarchy of paragraphs, (sub-)titles, and related clauses?

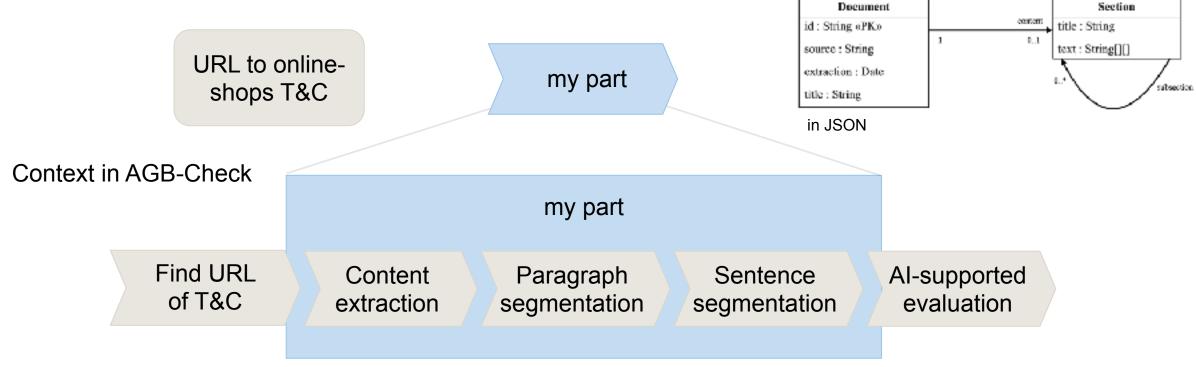


Figure 2.1: Context in AGB-Check.

Research Question



Example

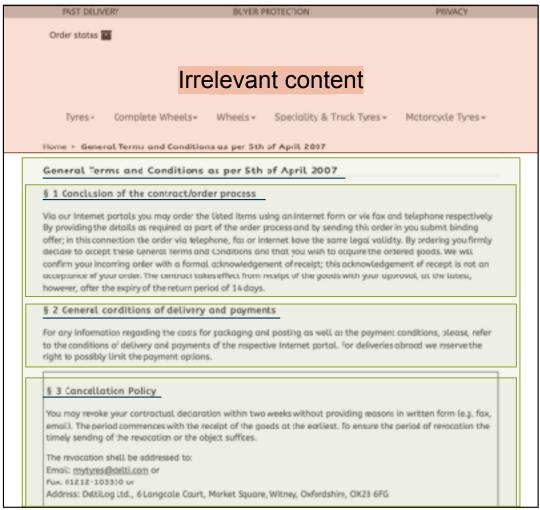


Figure 2.2: T&Cs of www.tyres-guru.co.uk, 210905.

```
"April",
"content": [
                                                       "1980",
                                                       ")",
    "subsections": [
                                                       "is",
        "subsections": [],
                                                       "excluded",
        "text": [
            "Via",
            "our",
            "Internet",
                                                  "title": "\u00a7 10
            "portals",
                                      Applicable law"
            "you",
            "may",
            "order",
                                              ],
            "the",
                                              "text": [
            "listed"
                                               []
            "items",
            "using",
                                              "title": "General Terms and
            "an",
            "Internet",
                                      Conditions as per 5th of April
            "form",
                                      2007"
            "or",
            "via",
            "fax",
            "and",
                                         "extractionDate": [
            "telephone",
                                           816109,
            "respectively",
                                           39,
                                           31.
                                           19,
                                           5,
            "providing",
                                           9,
            "the",
                                           2021
            "details",
            "as",
            "required",
                                         "id": 4803567503487537073,
            "as",
                                         "source": "https://www.tyres-
            "part",
                                      guru.co.uk/AGBs.html",
            "of",
            "the",
                                         "title": "General Terms and
            "order",
                                      Conditions as per 5th of April
            "process",
                                      2007 - Tyre Guru"
            "and",
            "by"
```

Code 2.1: Excerpt of parsed JSON, start & end.

Outline



- 1. Motivation
- 2. Research Question

3. Content Extraction

- Requirements
- Implementation
- Processing Demo

4. Hierarchy Extraction

- Requirements
- Implementation
- Processing Demo
- 5. Architecture & Additional Technology
- 6. Evaluation
 - Content Extraction
 - Hierarchy Extraction

7. Conclusion & Future Work

Content Extraction – Requirements



Evaluation of Existing Solutions

boilerpipe (Canola, LargestContent & Article), trafilatura, jusText

Methodology of Existing Solutions:

- Shallow Text Features to classify text as either Content or NotContent
- identification of largest container
- relative position of DOM-node

Problems:

- only extracts one of two large sections
- addresses, phone numbers, etc.

Manual Review

Potential Problems:

withdrawal forms (different style)

Properties of the Content:

- main content fills largest part of the page
- same style for the main content
- content is not interrupted (continuous)



Identify Most Common Style (MCS)

Find Lowest Common Ancestor holding a minimum of 85% of the characters in the MCS

Fallback Solution

find the largest subsequence of direct body-children holding the MCS

Rationale:

- uninterrupted continuous content
 - DOM-node or series of subsequent DOM-nodes form main content
- main content text usually shares a common visual style (and depth)
 - may also include different styles for headlines or withdrawal forms
- use of CMS (Content Management System)
 - main-content often grouped in a container (i.e. a <div>)



Identify Most Common Style (MCS)

Find Lowest Common Ancestor holding a minimum of 85% of the characters in the MCS

Fallback Solution

find the largest subsequence of direct body-children holding the MCS

- 1. count number of characters for each style
 - naive style and short text exclusion
- 2. determine the style with the highest amount of characters (MCS)

- style approximations
 - naive style: <tag>\${<attributes>}
 - rendered style (see Figure 3.1)
- additional rules
 - short text exclusion: only account for nodes containing more than 3 words

Font

isUnderlined : Bool

weight: Int

fontSize : Int

fontFamily: String

color : String

Figure 3.1: UML representation of fonts used to represent rendered styles.



Identify Most Common Style (MCS)

Find Lowest Common Ancestor
holding a minimum of 85% of the characters in the MCS

Fallback Solution

find the largest subsequence of direct body-children holding the MCS

- 1. search for the lowest node (highest depth) holding at least 85% of the characters of the MCS
- 2. if such a node is available, it is the main content node; if not, the fallback solution is used

- threshold 85%:
 - a higher threshold would extract too much content or trigger the fallback in more cases, worsening the overall content extraction performance
 - a lower threshold would extract to little content, i.e. a node holding only a part of the main content in more cases, worsening the overall content extraction performance

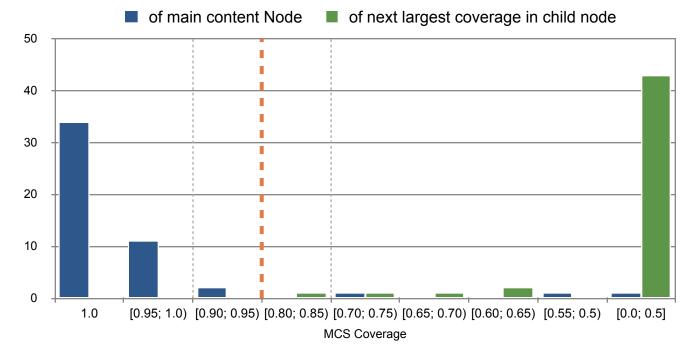


Figure 3.2: Amount of nodes covering a specific amount of the MCS characters*.

^{*} Only existing ranges are included in the diagram



Identify Most Common Style (MCS)

Find Lowest Common Ancestor holding a minimum of 85% of the characters in the MCS

Fallback Solution

find the largest subsequence of direct body-children holding the MCS

- 1. no node covering at holding at least 85% of MCS characters found by LowestCommonAncestor Extractor
- 2. Find largest subsequence of direct body-children holding the MCS
- assumption:
 - there is no container holding the main content
 - main content spread over multiple direct body-children
- risk:
 - there is a container holding the main content but it is made of less than 85% of the MCS characters

Content Extraction – Processing Demo



Raw HTMI

```
<html>
   Shead?
      <title>Terms and Conditions of Deno-Shop</title>
   -chead>
   Shody -
      64592
         <hi>ivelcome to the Beno-Shop*/hi>
         Sbd2NavigationS/bd2Sbd2ProductaS/td2
                   SaleMoout Us
             4/4192
      <4172
         <h8>Terms and Conditions</h8>
         <bs/>def>1. Loren Tpsur
         dolor sit amot, numecutetuer adipinging elit. Aenean commodo
         ligula eget dolor. Aemean massa. Cam sociis matoque penatibus
         et magnis dis parturient montes, mascetur ridiculus mas.
         <h6>1.1 Danec quans/h6>.
         felis, ultricies mec, pellentesque eu, pretium quis, sen.
         Mulla consequat massa quis emin. Dones pede justo, fringilla
         vel, ChDaliquet necC/bD, vulputate eget, arcu.C/pD
         <h6>1.2 In emim justo, rhoncus</h6>
         spout, imperdict a, wemenatis witae, justo. Mullam distun-
         felis en pede mollis pretium. Integer tincidunt. Cras
         dapibus. Vivamus elementum semper misi. Aemean vulputate
         eleifend tellus.
         <h8>2. lenean loc</h5>
         ligula, porttitor eu, consequat vitae, eleifend ac, enim.
         Aliquam lorem ante, dapilma im, viverra quis, fengiat a,
         tellus. Phasellus viverra nulla ut metus varius laoreet.
         Quisque rutrum. Acness imperdict.
      Thanks for visiting Deno-Shop
   Sybody>
</href>
```

Code 3.1: Raw HTML for the Demo-Shop used in the processing demo.

Welcome to the Demo-Shop

Navigation Products Sale About Us

Terms and Conditions

1. Lorem Ipsum

dolor sit amet, consectetuer adipiscing elit. Aenean commodo ligula eget dolor. Aenean massa, Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.

1.1 Dones quam

felis, ultricies nec, pellentesque eu, pretium quis, sem. Nulla consequat massa quis enim. Donce pede justo, fringilla vel, aliquet nec, vulputate eget, arcu-

1.2 In coim justo, rhoncus

ut, imperdiet a, venenatis vitae, justo. Nullam dictum felis eu pede mollis pretium. Integer tincidunt. Cras dapibus. Vivamus elementum semper nisi. Aenean vulputate eleifend tellus.

2. Acmean Ico

ligula, porttitor cu, consequat vitae, cleifend ac, cnim. Aliquam lorem ante, dapibus in, viverra quis, feugiat a, tellus. Phasellus viverra nulla ut metus varius lacreet. Quisque rutrum. Aenean imperdiet.

Thanks for visiting Demo-Shop

Figure 3.3: Rendered HTML of Code 3.1.

Content Extraction – Processing Demo



MCS, naive style and short text exclusion

Style	Characters	Characters discarded (less than 4 words)
h1\${}	24	0
td\${}	30	30
h3\${}	20	20
h5\${}	27	27
p\${}	742	0
h6\${}	40	14
b\${}	11	11

Table 3.1: Naive styles and number of relevant characters.

Welcome to the Demo-Shop



Terms and Conditions

1. Lorem Ipsum

dolor sit amet, consectetuer adipiscing elit. Aenean commodo ligula eget dolor. Aenean massa. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.

1.1 Dones quam

felis, ultricies nec, pellentesque eu, pretium quis, sem. Nulla consequat massa quis enim. Donce pede justo, fringilla vel, aliquet nec, vulputate eget, arcu-

1.2 In coim justo, rhoncus

ut, imperdiet a, venenatis vitae, justo. Nullam dictum felis eu pede mollis pretium. Integer tincidunt. Cras dapibus. Vivamus elementum semper nisi. Aenean vulputate eleifend tellus.

2. Acmean Ico

ligula, porttitor cu, consequat vitae, eleifend ac, enim. Aliquam lorem ante, dapibus in, viverra quis, feugiat a, tellus. Phasellus viverra nulla ut metus varius lacreet. Quisque rutrum, Aenean imperdiet.

Thanks for visiting Demo-Shop

Figure 3.4: Rendered HTML of Code 3.1.

Content Extraction – Processing Demo

Lowest Common Ancestor (threshold = 95%*)

Node (XPath)	Coverage	Depth
/html/body/div[2]	0.9664	1
/html/body/div[2]/p[1]	0.2355	2
/html/body/div[2]/p[2]	0.1984	2
/html/body/div[2]/p[3]	0.2529	2
/html/body/div[2]/p[4]	0.2796	2
/html/body/p	0.0336	1

Table 3.2: Nodes, coverage and depth.

Welcome to the Demo-Shop Irrelevant content

Navigation Products Sale About Us

Terms and Conditions

1. Lorem Ipsum

dolor sit amet, consectetuer adipiscing elit. Aenean commodo ligula eget dolor. Aenean massa. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.

1.1 Dones quam

felis, ultricies nec, pellentesque eu, pretium quis, sem. Nulla consequat massa quis enim. Donce pede justo, fringilla vel, aliquet nec, vulputate eget, arcu-

1.2 In coim justo, rhoncus

ut, imperdiet a, venenatis vitae, justo. Nullam dictum felis eu pede mollis pretium. Integer tincidunt. Cras dapibus. Vivamus elementum semper nisi. Aenean vulputate eleifend tellus.

2. Acmean Ico

ligula, porttitor cu, consequat vitae, eleifend ac, enim. Aliquam lorem ante, dapibus in, viverra quis, feugiat a, tellus. Phasellus viverra nulla ut metus varius lacreet. Quisque rutrum, Aenean imperdiet,

Thanks for visiting Demo-Shop

Irrelevant content

Figure 3.5: Rendered Demo-Shop with colored *content* and *not-content* areas.

^{*} threshold set to 95% for demonstration purposes due to the short document

Outline



- 1. Motivation
- 2. Research Question
- 3. Content Extraction
 - Requirements
 - Implementation
 - Processing Demo

4. Hierarchy Extraction

- Requirements
- Implementation
- Processing Demo
- 5. Architecture & Additional Technology
- 6. Evaluation
 - Content Extraction
 - Hierarchy Extraction
- 7. Conclusion & Future Work

Hierarchy Extraction – Requirements



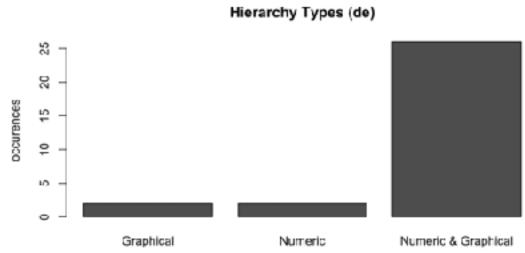


Figure 4.1: Hierarchy types in the German requirements sample.

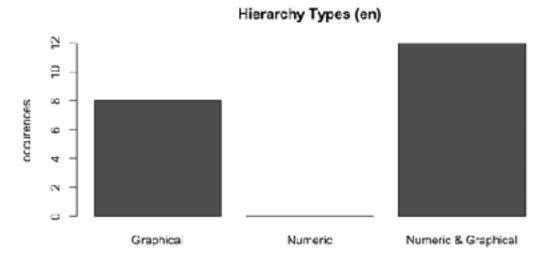


Figure 4.2: Hierarchy types in the English requirements sample.

Manual Review

Properties of the hierarchy representation:

- enumeration pattern (alphabetic, Arabic, Latin/Roman)
- headlines in a more prominent style, following paragraphs associated with the headline
- assumption: sections are not interrupted exception: lists

Potential Problems:

table of contents

[VLDB]Manabe, T., Tajima, K.: Extracting logical hierarchical structure of HTML documents based on headings.

[github]HEPS; [github]TC-Detection-Corpus



Block Formation

Visual Based Hierarchy Extraction

Enumeration Based Hierarchy Extraction Adjustments for Lists (<1i>)

Rationale:

- headlines with more prominent styles preface their associated section(s)
 - assumption: sections are not interrupted
- enumeration patterns of different styles
 - Ignore enumeration patterns in the table of contents
- lists need special treatment
 - may violate the assumption of uninterrupted sections
 - each <1i> forms a section



Block Formation

Visual Based Hierarchy Extraction

Enumeration Based Hierarchy Extraction Adjustments for Lists (<1i>>)

- 1. group all text elements and include line breaks forced by HTML tags
- 2. create blocks from text sequences between line breaks
 - the block's style is determined by the majority of characters; characters in <a> are not considered
- 3 enumerations for each block are detected
 - <1i> blocks receive a special numeration style (List)

Privacy Policy

We do not disclose buyers' information to third parties other than when order details are processed as part. of the order fulfilment process, e.g. courier companies. In this case, the third party will not disclose any of the details to any other third party. To read our full GDPR policy, please click **here**.

Cookies are used on this shopping site to keep track of the contents of your shopping cart, to store delivery addresses if the address book is used and to store your details if you select the 'Remember Me' Option. They are also used after you have logged on as part of that process. You can turn off cookies within your browser. by going to 'Tools | Internet Options | Privacy' and selecting to block cookies. If you turn off cookies, you will be unable to place orders or benefit from the other features that use cookies. Data collected by this site is used to:

a. Take and fulfill customer orders

b. Administer and enhance the site and service.

Figure 4.3: Example of block formation, www.telescopehouse.com.



Block Formation

Visual Based Hierarchy Extraction

Enumeration Based Hierarchy Extraction Adjustments for Lists (<1i>)

- 1. determine the MCS in the main content
 - rendered style and short text exclusion
- 2. identify the next headline style
 - headline styles are more prominent than the MCS
- 3. select all headlines of the current headline style and their associated content
 - associated content according to assumption: sections are not interrupted
- 4. continue within each of the content sections (from 2.)



Block Formation

Visual Based **Hierarchy Extraction**

Enumeration Based Hierarchy Extraction Adjustments for Lists (<1i>>)

- enumerations are detected using a regular expression
 - ▶\s[\(§]?(([IVXLivxl]{1,7})|([0-9]{1,2})|[a-zA-Z])([\.\-,:](([IVXLivxl]{1,7})|([0-9] $\{1,2\}$) | [a-zA-Z]) | $*[-:\.)$ | ?\s

A: adjust hierarchy of existing visual based hierarchy using enumerations

- 1. determine the next enumeration style (i.e. Latin/Roman, alphabetic, Arabic)
- 2. validate numeration style
 - numeral: continuous numeration within enumeration pattern, at least 2 occurrences
 - table of contents: short blocks (less than 10 words) need to be followed by a content block
- 3. apply adjustment and continue in the subsections recursively

B: add hierarchy within existing non-headline content blocks using enumerations



Block Formation

Visual Based Hierarchy Extraction

Enumeration Based Hierarchy Extraction Adjustments for Lists (<1i>)

- lists need special treatment
 - may violate the assumption of uninterrupted sections
- 1. check for list enumeration style within each content section
- 2. adjust hierarchy according to list
 - content before and after the list on the same hierarchy level is associated with the same parent as the list content



Block Formation

ID	Text	Font- Size	Bold	Under- lined
1	Terms and Conditions	18px		X
2	1. Lorem Ipsum	13px		X
3	dolor sit amet, consectetuer adipiscing elit. Aenean commodo ligula eget dolor. Aenean massa. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.	16px	X	X
4	1.1 Donec quam	10px 16px	✓	X
5	felis, ultricies nec, pellentesque eu, pretium quis, sem. Nulla consequat massa quis enim. Donec pede justo, fringilla vel, aliquet nec, vulputate eget, arcu.	16px	X	X
6	1.2 In enim justo, rhoncus	10px		X
7	ut, imperdiet a, venenatis vitae, justo. Nullam dictum felis eu pede mollis pretium. Integer tincidunt. Cras dapibus. Vivamus elementum semper nisi. Aenean vulputate eleifend tellus.	16px	X	X
8	2. Aenean leo	13px	✓	X
9	ligula, porttitor eu, consequat vitae, eleifend ac, enim. Aliquam lorem ante, dapibus in, viverra quis, feugiat a, tellus. Phasellus viverra nulla ut metus varius laoreet. Quisque rutrum. Aenean imperdiet.	16px	X	X

Table 4.1: List of all Blocks (incl. rendered style).

Welcome to the Demo-Shop Irrelevant content

Navigation Products Sale About Us

Terms and Conditions

1. Lorem Ipsum

dolor sit amet, consectetuer adipiscing elit. Aenean commodo ligula eget dolor. Aenean massa. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.

1.1 Dones quam

felis, ultricies nec, pellentesque eu, pretium quis, sem. Nulla consequat massa quis enim. Donce pede justo, fringilla vel, aliquet nec, vulputate eget, arcu.

1.2 In coim justo, rhoncus

ut, imperdiet a, venenatis vitae, justo. Nullam dictum felis eu pede mollis pretium. Integer tincidunt. Cras dapibus. Vivamus elementum semper nisi. Aenean vulputate eleifend tellus.

2. Acmean Ico

ligula, porttitor cu, consequat vitae, eleifend ac, enim. Aliquam lorem ante, dapibus in, viverra quis, feugiat a, tellus. Phasellus viverra nulla ut metus varius lacreet. Quisque rutrum, Aenean imperdiet,

Thanks for visiting Demo-Shop

Irrelevant content

Figure 4.4: Rendered Demo-Shop with colored *content* and *not-content* areas.

Visual Based

Determine MCS (rendered & short text exclusion):

Style	Characters	Characters discarded (less than 4 words)
18px, bold	20	20
13px, bold	27	27
16px	724	0
10px, bold	40	14

Table 4.2: MCS in relevant content.

Welcome to the Demo-Shon Irrelevant content

Navigation Products Sale About Us

Terms and Conditions

1. Lorem Ipsum

dolor sit amet, consectetuer adipiscing elit. Aenean commodo ligula eget dolor. Aenean massa. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.

1.1 Dones quam

felis, ultricies nec, pellentesque eu, pretium quis, sem. Nulla consequat massa quis enim. Donce pede justo, fringilla vel, aliquet nec, vulputate eget, arcu.

1.2 In coim justo, rhoncus

ut, imperdiet a, venenatis vitae, justo. Nullam dictum felis eu pede mollis pretium. Integer tincidunt. Cras dapibus. Vivamus elementum semper nisi. Aenean vulputate eleifend tellus.

2. Acmean Ico

ligula, porttitor cu, consequat vitae, eleifend ac, enim. Aliquam lorem ante, dapibus in, viverra quis, feugiat a, tellus. Phasellus viverra nulla ut metus varius lacreet. Quisque rutrum, Aenean imperdiet,

Thanks for visiting Demo-Shop

Irrelevant content

Figure 4.4: Rendered Demo-Shop with colored *content* and *not-content* areas.



Visual Based I

ID	Text	Font- Size	Bold	Under- lined
1	Terms and Conditions	18px	✓	X
2	1. Lorem Ipsum	13px	✓	X
3	dolor sit amet, consectetuer adipiscing elit. Aenean commodo ligula eget dolor. Aenean massa. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.	16px	X	X
4	1.1 Donec quam	10px	1	X
5	felis, ultricies nec, pellentesque eu, pretium quis, sem. Nulla consequat massa quis enim. Donec pede justo, fringilla vel, aliquet nec, vulputate eget, arcu.	16px	X	X
6	1.2 In enim justo, rhoncus	10px	✓	X
7	ut, imperdiet a, venenatis vitae, justo. Nullam dictum felis eu pede mollis pretium. Integer tincidunt. Cras dapibus. Vivamus elementum semper nisi. Aenean vulputate eleifend tellus.	16px	X	X
8	2. Aenean leo	13px	√	X
9	ligula, porttitor eu, consequat vitae, eleifend ac, enim. Aliquam lorem ante, dapibus in, viverra quis, feugiat a, tellus. Phasellus viverra nulla ut metus varius laoreet. Quisque rutrum. Aenean imperdiet.	16px	X	X

Welcome to the Demo-Shop Irrelevant content

Navigation Products Sale About Us

Terms and Conditions

1. Lorem Ipsum

dolor sit amet, consectetuer adipiscing elit. Aenean commodo ligula eget dolor. Aenean massa. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.

1.1 Dones quam.

felis, ultricies nec, pellentesque eu, pretium quis, sem. Nulla consequat massa quis enim. Donce pede justo, fringilla vel, aliquet nec, vulputate eget, arcu-

1.2 In enim justo, rhoncus

ut, imperdiet a, venenatis vitae, justo. Nullam dictum felis eu pede mollis pretium. Integer tincidunt. Cras dapibus. Vivamus elementum semper nisi. Aenean vulputate eleifend tellus.

2. Acmean Ico

ligula, portitior cu, consequat vitae, eleifend ac, enim. Aliquam lorem ante, dapibus in, viverra quis, feugiat a, tellus. Phasellus viverra nulla ut metus varius lacreet. Quisque rutrum. Aenean imperdiet,

Thanks for visiting Demo-Shop

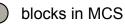
Irrelevant content

Figure 4.5: Rendered Demo-Shop with colored content and not-content areas and visual based hierarchy extraction, step 1.

Table 4.3: Visual based hierarchy extraction, step 1.



associated blocks (to preceding headline)





Visual Based II

ID	Text	Font- Size	Bold	Under- lined
1	Terms and Conditions	18px	✓	X
2	1. Lorem Ipsum	13px	√	X
3	dolor sit amet, consectetuer adipiscing elit. Aenean commodo ligula eget dolor. Aenean massa. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.	16px	X	X
4	1.1 Donec quam	10px	✓	X
5	felis, ultricies nec, pellentesque eu, pretium quis, sem. Nulla consequat massa quis enim. Donec pede justo, fringilla vel, aliquet nec, vulputate eget, arcu.	16px	X	X
6	1.2 In enim justo, rhoncus	10px	✓	X
7	ut, imperdiet a, venenatis vitae, justo. Nullam dictum felis eu pede mollis pretium. Integer tincidunt. Cras dapibus. Vivamus elementum semper nisi. Aenean vulputate eleifend tellus.	16px	X	X
8	2. Aenean leo	13px	✓	X
9	ligula, porttitor eu, consequat vitae, eleifend ac, enim. Aliquam lorem ante, dapibus in, viverra quis, feugiat a, tellus. Phasellus viverra nulla ut metus varius laoreet. Quisque rutrum. Aenean imperdiet.	16px	X	X

Welcome to the Demo-Shop Irrelevant content

Navigation Products Sale About Us

Terms and Conditions

1. Lorem Ipsum

dolor sit amet, consectetuer adipiscing elit. Aenean commodo ligula eget dolor. Aenean massa. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.

1.1 Donee quam

felis, ultricies nec, pellentesque eu, pretium quis, sem. Nulla consequat massa quis enim. Donce pede justo, fringilla vel, aliquet nec, vulputate eget, arcu-

1.2 In enim justo, rhoncus

ut, imperdiet a, venenatis vitae, justo. Nullam dictum felis eu pede mollis pretium. Integer tincidunt. Cras dapibus. Vivamus elementum semper nisi. Aenean vulputate eleifend tellus.

2. Acmean Ico

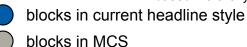
ligula, portitior cu, consequat vitae, eleifend ac, enim. Aliquam lorem ante, dapibus in, viverra quis, feugiat a, tellus. Phasellus viverra nulla ut metus varius lacreet. Quisque rutrum. Aenean imperdiet.

Thanks for visiting Demo-Shop

Irrelevant content

Figure 4.6: Rendered Demo-Shop with colored content and not-content areas and visual based hierarchy extraction, step 2.

Table 4.4: Visual based hierarchy extraction, step 2.



associated blocks (to preceding headline)



Visual Based III

ID	Text	Font- Size	Bold	Under- lined
1	Terms and Conditions	18px	✓	X
2	1. Lorem Ipsum	13px	√	X
3	dolor sit amet, consectetuer adipiscing elit. Aenean commodo ligula eget dolor. Aenean massa. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.	16px	X	X
4	1.1 Donec quam	10px	√	X
5	felis, ultricies nec, pellentesque eu, pretium quis, sem. Nulla consequat massa quis enim. Donec pede justo, fringilla vel, aliquet nec, vulputate eget, arcu.	16px	X	X
6	1.2 In enim justo, rhoncus	10px	✓	X
7	ut, imperdiet a, venenatis vitae, justo. Nullam dictum felis eu pede mollis pretium. Integer tincidunt. Cras dapibus. Vivamus elementum semper nisi. Aenean vulputate eleifend tellus.	16px	X	X
8	2. Aenean leo	13px	✓	X
9	ligula, porttitor eu, consequat vitae, eleifend ac, enim. Aliquam lorem ante, dapibus in, viverra quis, feugiat a, tellus. Phasellus viverra nulla ut metus varius laoreet. Quisque rutrum. Aenean imperdiet.	16px	X	X

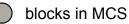
Welcome to the Demo-Shop Irrelevant content Navigation Products Sale About Us Terms and Conditions 1. Lorem Ipsum dolor sit amet, consectetuer adipiscing elit. Aenean commodo ligula eget dolor. Aenean massa. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. 1.1 Dones quanfelis, ultricies nec, pellentesque eu, pretium quis, sem. Nulla consequat massa quis enim. Donce pede justo, fringilla vel, aliquet nec, vulputate eget, arcu-1.2 In enim justo, rhoneus ut, imperdiet a, venenatis vitae, justo. Nullam dictum felis eu pede mollis pretium. Integer tincidunt. Cras dapibus. Vivamus elementum semper nisi. Aenean vulputate eleifend tellus. 2. Acmean Ico ligula, portitior cu, consequat vitae, eleifend ac, enim. Aliquam lorem ante, dapibus in, viverra quis, feugiat a, tellus. Phasellus viverra nulla ut metus varius lacreet. Quisque rutrum. Aenean imperdiet. Thanks for visiting Demo-Shop Irrelevant content

Figure 4.7: Rendered Demo-Shop with colored content and not-content areas and visual

Table 4.5: Visual based hierarchy extraction, step 3; done.

based hierarchy extraction, step 3; done. blocks in current headline style

associated blocks (to preceding headline)





Enumeration Based

ID	Enumeration	Pattern
1		
2	[1]	[numeric]
3		
4	[1, 1]	[numeric, numeric]
5		
6	[1, 2]	[numeric, numeric]
7		
8	[2]	[numeric]
9		

Table 4.6: Detected enumerations & patterns.

Enumeration valid ⇒ no adjustment needed No <1i> in html \Rightarrow no adjustment needed

Welcome to the Demo-Shop Irrelevant content Navigation Products Sale About Us Terms and Conditions 1. Lorem Ipsum dolor sit amet, consectetuer adipiscing elit. Aenean commodo ligula eget dolor, Aenean massa. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. 1.1 Dones quam felis, ultricies nec, pellentesque eu, pretium quis, sem. Nulla consequat massa quis enim. Donce pede justo, fringilla vel, aliquet nec, vulputate eget, arcu-1.2 In enim justo, rhonens ut, imperdiet a, venenatis vitae, justo. Nullam dictum felis eu pede mollis pretium. Integer tincidunt. Cras dapibus. Vivamus elementum semper nisi. Aenean vulputate eleifend tellus. 2. Acmean Ico ligula, portitior cu, consequat vitae, eleifend ac, enim. Aliquam lorem ante, dapibus in, viverra quis, feugiat a, tellus. Phasellus viverra nulla ut metus varius lacreet. Quisque rutrum. Aenean imperdiet. Thanks for visiting Demo-Shop Irrelevant content

Figure 4.7: Rendered Demo-Shop with colored content and not-content areas, visual based hierarchy extraction and enumeration based hierarchy extraction incl. adjustments for lists.

Outline



- 1. Motivation
- 2. Research Question
- 3. Content Extraction
 - Requirements
 - Implementation
 - Processing Demo

4. Hierarchy Extraction

- Requirements
- Implementation
- Processing Demo

5. Architecture & Additional Technology

- 6. Evaluation
 - Content Extraction
 - Hierarchy Extraction

7. Conclusion & Future Work

Architecture & Additional Technology



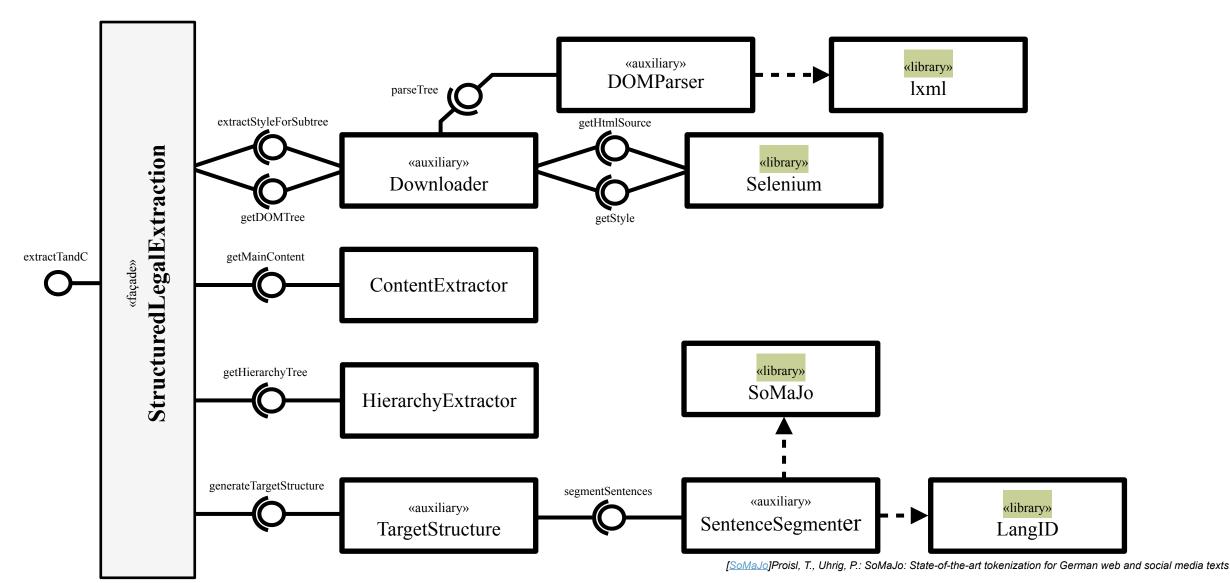


Figure 5.1: UML component diagram of the structuredlegalextraction library.

[aithub]SoMaJo; [aithub]Selenium; [aithub]LangID; [aithub]lxml

Architecture & Additional Technology



Web Page Download

HTML Parser

Sentence Segmentation and Tokenising

Language Determination

- download full HTMI
- extract CSS style
- full XPath support
- Java-Script support (opt.)

- full XPath support
- generate XPath from elements
- applicable in the domain of legal documents

- support German and English language

«library» Selenium «library» lxml

«library» SoMaJo

«library» LangID

Architecture & Additional Technology



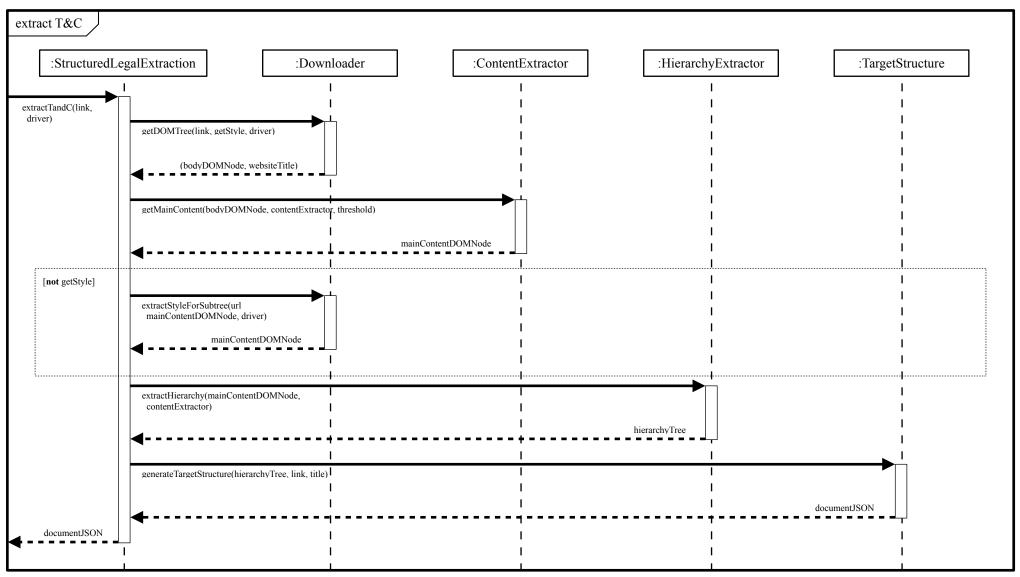


Figure 5.2: UML sequence diagram of the structuredlegalextraction libraries extractTandC method.

[github]SoMaJo; [github]Selenium; [github]LangID; [github]Ixml

Outline



- 1. Motivation
- 2. Research Question
- 3. Content Extraction
 - Requirements
 - Implementation
 - Processing Demo

4. Hierarchy Extraction

- Requirements
- Implementation
- Processing Demo
- 5. Architecture & Additional Technology
- 6. Evaluation
 - Content Extraction
 - Hierarchy Extraction

7. Conclusion & Future Work

Evaluation – Content Extraction



Two samples used during evaluation:

- Requirements sample used to derive requirements (30 de, 20 en)*
- Test sample to validate results (30 de, 20 en)**

Extractor	too late	too early	correct
LowestCommonAncestor	1	3	45
Extractor			
Boilerpipe ArticleExtractor	16	4	24
Boilerpipe	29	2	13
LargestContententExtractor			
Boilerpipe CanolaExtractor	7	15	22
JusText	6	11	25
Trafilatura	5	2	37

Extractor	too late	too early	correct
LowestCommonAncestor	3	2	44
Extractor			
Boilerpipe ArticleExtractor	23	13	8
Boilerpipe	32	1	11
LargestContententExtractor			
Boilerpipe CanolaExtractor	3	27	14
JusText	8	15	19
Trafilatura	5	4	35

Table 6.2: End, requirements sample.

Extractor	too late	too early	correct
LowestCommonAncestor Extractor - START	3	0	46
LowestCommonAncestor Extractor - END	0	2	47

Table 6.3: Start & End, test sample.

Performance in the center is not shown as LowestCommonAncestor Extractor extracts continuous content.

^{*} LowestCommonAncestor Extractor processed 49 of 50 web pages, the other extractors processed 44 of 50 web pages from the requirements sample.

** LowestCommonAncestor Extractor processed 49 of 50 web pages from the test sample.

[[]github]boilerpipe (Canola, LargestContent & Article); [github]trafilatura; [github]jusText; [github]TC-Detection-Corpus

Evaluation – Content Extraction

Sources of errors:

- 1. Threshold too high \Rightarrow no lowest common ancestor ⇒ fallback algorithm is triggered
- 2. No container surrounding the main content \Rightarrow no lowest common ancestor ⇒ fallback algorithm is triggered
- 3. Different tags with the same style \Rightarrow whenever one tag is held in its own container \Rightarrow only part of the content extracted
 - Example: part of T&C in <div> followed by numerous <1i> surrounded by <u1>/<o1>, see Figure 6.1

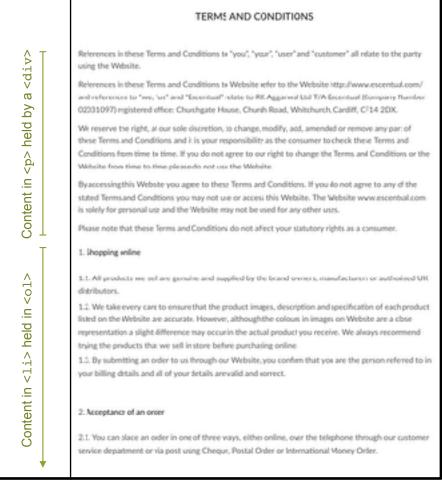


Figure 6.1: Terms and Conditions of escentual.com, 210903

Evaluation – Hierarchy Extraction



Error scores:

0: each section with correct parent, correct content, and correct title

• 0.4: wrong parent

• 0.5: wrong content

• 0.1: wrong title

Error Share	Amount
0	12
(0; 0.05]	12
(0.05; 0.1]	1
(0.1; 0.15]	2
(0.15; 0.2]	1
(0.2; 0.3]	1
(0.3; 0.5]	1
(0.5; 1]	0
Failed*	0

Table 6.4: Error scores of hierarchy extraction for the German requirements sample.

Error Share	Amount
0	5
(0; 0.05]	4
(0.05; 0.1]	1
(0.1; 0.15]	2
(0.15; 0.2]	0
(0.2; 0.3]	0
(0.3; 0.5]	1
(0.5; 1]	0
Failed*	6

Table 6.5: Error scores of hierarchy extraction for the English requirements sample.

Error Share	Amount
0	8
(0; 0.05]	11
(0.05; 0.1]	1
(0.1; 0.15]	1
(0.15; 0.2]	1
(0.2; 0.3]	1
(0.3; 0.5]	5
(0.5; 1]	0
Failed*	2

Table 6.6: Error scores of hierarchy extraction for the German test sample.**

Error Share	Amount
0	9
(0; 0.05]	3
(0.05; 0.1]	2
(0.1; 0.15]	1
(0.15; 0.2]	2
(0.2; 0.3]	2
[0.3; 0.5]	1
(0.5; 1]	0
Failed*	0

Table 6.7: Error scores of hierarchy extraction for the English test sample.

 $x \in [a;b) \mid x \ge x \land x < b$

 $error share = \frac{sum of total error scores}{sum of total nodes extracted}$

^{*} Failed content extraction.

^{**} One page was no longer available during the evaluation of the hierarchy extraction.

Evaluation – Hierarchy Extraction



Sources of errors:

- 1. Use of bold/large text, i.e. violation of assumptions about headline styles ⇒ regular content is identified as headline (see Figure 6.2)
- 2. Wrong enumeration \Rightarrow sequence of numeration is not valid ⇒ no hierarchy detected
- 3. Interrupted sections, i.e. violation of the assumption "Sections are not interrupted" ⇒ text is assigned to the wrong headline



Figure 6.2: Terms and Conditions of steber.de, 210903

- 4. Use of tables ⇒ each cell is its own block ⇒ many blocks in different styles/many numbers in the table ⇒ wrong visual/enumeration based adjustments
- 5. Failed style extraction \Rightarrow standard style is set \Rightarrow wrong visual based separation

Outline



- 1. Motivation
- 2. Research Question
- 3. Content Extraction
 - Requirements
 - Implementation
 - Processing Demo

4. Hierarchy Extraction

- Requirements
- Implementation
- Processing Demo
- 5. Architecture & Additional Technology
- 6. Evaluation
 - Content Extraction
 - Hierarchy Extraction

7. Conclusion & Future Work

Conclusion & Future Work



RQ1: What are special requirements for the extraction of contracts in comparison with regular (i.e. news articles, blogposts) content?

- continuous content
- addresses and phone numbers are part of the relevant content (no hint for imprint/footer)

RQ2: How to extract the relevant parts from contracts in raw HTML?

LowestCommonAncestor Extractor

RQ3: How to extract the structure and the hierarchy of paragraphs, (sub-)titles, and related clauses?

- rule based approach ignoring structural information except for <1i> and line breaks (blocks)
- enumeration patterns and visual styles

Future Work:

- evaluate LowestCommonAncestor Extractor in other domains
- address the sources of error in hierarchy extraction and improve time efficiency of rendered style extraction

