

Towards an Open Platform for Legal Information

Malte Ostendorff
Open Justice e.V.
mo@openlegaldata.io

Till Blume
Open Justice e.V.
tb@openlegaldata.io

Saskia Ostendorff
Open Justice e.V.
so@openlegaldata.io

ABSTRACT

Recent advances in the area of legal information systems have led to a variety of applications that promise support in processing and accessing legal documents. Unfortunately, these applications have various limitations, e. g., regarding scope or extensibility. Furthermore, we do not observe a trend towards open access in digital libraries in the legal domain as we observe in other domains, e. g., economics of computer science. To improve open access in the legal domain, we present our approach for an open source platform to transparently process and access Legal Open Data. This enables the sustainable development of legal applications by offering a single technology stack. Moreover, the approach facilitates the development and deployment of new technologies. As proof of concept, we implemented six technologies and generated metadata for more than 250,000 German laws and court decisions. Thus, we can provide users of our platform not only access to legal documents, but also the contained information.

CCS CONCEPTS

• **Information systems** → **Open source software**; **Digital libraries and archives**; *Information retrieval*; • **Applied computing** → **Law**; *Document searching*.

KEYWORDS

Legal information system, Open data, Open source, Legal data

1 INTRODUCTION

The importance of automatically processing legal documents is rising. Recent advances in research offer a portfolio of technologies to process legal documents, e. g., extracting, aggregating, and linking information from text. Hence, mostly commercial tools and platforms have emerged that promise support in processing and accessing legal documents. Unfortunately, there are various limitations when using these tools. For instance, they are country-specific, lack transparency and extensibility (closed source), or do not provide access to the raw data. These criteria are essential for legal data analysis [4], and the development of innovative technologies, e. g., visual query interfaces¹. Data analysis and visualization are of great benefit when interpreting the information, e. g., to investigate the mutual dependencies between statutes and the temporal evolution of law. In our opinion, democratizing the access to these tools and providing the data is fundamental when one is interested in facilitating access to justice and innovation in the legal domain. A key element to achieve these goals is open data, that can reduce integration costs, improve transparency, and harness the innovation of others [13].

In this paper, we present our approach for an open source platform to transparently process Legal Open Data by flexibly combining state-of-the-art technologies. Our approach enables the sustainable development of legal data processing tools by offering a single technology stack. The platform empowers others to quickly develop and deploy new technologies. As proof of concept, we implemented six technologies in an open processing pipeline, processed more than 250,000 laws and court decisions, and made them available on our Open Legal Data Platform. Our source code and our generated data is publicly available².

Below, we briefly discuss representative related projects and platforms. Subsequently, we present our approach as well as the implemented technologies.

2 RELATED WORK

CourtListener³ is a service for the United States, which is developed by the non-profit Free Law Project. CourtListener's goal is "to provide free, public, and permanent access to primary legal materials on the Internet for educational, charitable, and scientific purposes to the benefit of the general public and the public interest" [9]. CourtListener seeks to collect and freely distribute historical and current United States court opinions on state and federal level. However, other international jurisdictions are not in the scope of the project. Similarly, the Caselaw Access Project⁴ by Harvard Law Library aims to make all published U.S. court decisions freely available. The Finnish government developed the web service Finlex⁵, which provide laws and related legal documents as XML documents. In 2014, Frosterus et al. [5] improved Finlex in several ways, e. g., by transposing the XML documents to RDF documents following the Linked Open Data principles. They demonstrate the usefulness of Linked Open Data for content producers, application developers, and data analysts. OpenLaws⁶ is an open access platform for European legal information [7]. OpenLaws is built on top of open source software, but it does not provide access to the data. In summary, there are various projects scattered across the world that collect and publish legal documents. However, there is no single project that is open source, makes data openly accessible, and is not focused on a single country only.

3 THE OPEN LEGAL DATA PLATFORM

Our approach to develop a single legal technology platform is illustrated in Fig. 1. We provide the basic technology stack that legal engineers can build upon to develop new technologies. The developed technologies can be flexibly combined to provide country specific-platforms, e. g., for Germany. Developed technologies can

¹<https://www.vizlaw.de>, accessed: May 28, 2020

²<http://www.openlegaldata.io/>, accessed: May 28, 2020

³<http://www.courtlistener.com>, accessed: May 28, 2020

⁴<https://case.law>, accessed: May 28, 2020

⁵<http://www.finlex.fi>, accessed: May 28, 2020

⁶<http://www.openlaws.eu>, accessed: May 28, 2020

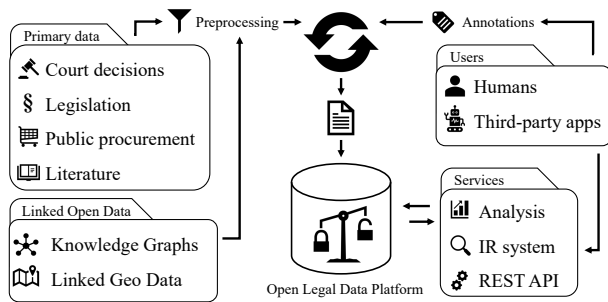


Figure 1: The Open Legal Data Platform vision: Harvesting legal documents from official sources, preprocessing of raw documents and enriching the documents with data from additional sources. The platform content is accessible through an API to facilitate analysis, information retrieval, and third-party apps.

be included in the global technology stack to make them accessible to developers around the world. Moreover, legal tech developers benefit from the tools and the data provided by the platforms via our REST API. Finally, researchers benefit from easy access and the latest technologies integrated into a single platform ecosystem to conduct analyses.

Although there are country-specific differences in the legal systems, we firmly believe that the development of a legal technology platform can be tackled with a single technology stack. Akoma-Ntoso [11] is one prominent example for an XML schema that aims to standardize legal documents on an international level. By integrating standards, like Akoma-Ntoso, the platform comparability for different countries can be facilitated. For any Open Legal Data platform, we need to access primary and secondary data sources (Fig. 1 left), process the data to generate additional information (Fig. 1 center), and provide services to users (Fig. 1 right).

For Open Legal Data, we access primary source like government services, and secondary sources liked Linked Open Data. The preprocessing system is designed to flexibly handle different types of documents, such as legislation or literature. The processing pipeline allows to enrich documents with additional information, e. g., automatically extracted references or manually created text annotations from domain experts. Finally, the data is made available to the public via information retrieval systems or REST APIs. As a foundation, we use the Django framework⁷. The user interfaces can be translated into different languages and adapted to specific information needs using Django’s template system. Django’s “app system” enables easy integration of new modules and the re-use of existing apps.

In the following, we describe in more detail technologies to access, process and provide legal documents. All technologies are integrated or are currently being integrated into our Open Legal Data technology stack.

3.1 Primary data sources & Linked Open Data

Finding and harvesting legal information is a challenging task due to several reasons. Accessing data directly from courts is time-consuming and expensive. Accessing data from sources on the Web induces quality issues. In the following, we present three alternative sources of legal information and our approach to include them.

Courts and Governmental data. We collaborate with courts to obtain decisions directly. Accessing data directly from courts has the highest level of trustworthiness. However, it is very time consuming since courts do not always make decisions publicly accessible. Furthermore, decisions obtained directly from courts are rarely in machine-readable formats and are not free of charge. Thus, information needs to be extracted from, e. g., from purchased PDF files.

Crawling trusted websites. Crawling trusted websites can significantly improve the amount of data. In Germany, there exists a small set of trusted websites. The German Federal Ministry of Justice and Consumer Protection (BMJV) operates websites with the latest version of federal legislation⁸ and decisions from federal courts⁹. In Germany, state-level legislation and decisions are not available on a central web service. Each state needs to be handled separately. On European level, the service EUR-LEX¹⁰ is the main data source. Additionally, we crawl legal blogs which have been shown to provide information for legal opinion mining [2]. Having different data sources requires the harmonization of the harvested data to avoid duplicates. To de-duplicate court decisions, we use the European Case Law Identifier [14].

Linked Open Data. According to the Open Data Monitor¹¹, 45% of all Open Data is currently provided in (semi-)structured and thus machine-readable format. Furthermore, the European Commission has identified the strong need to “opening up by default all scientific data” and to store and maintain it in the European Open Science Cloud¹². The European Union Open Data Portal¹³ serves as a single point of access to Open Data produced by EU institutions and bodies. In addition to major data portals, there exists a variety of small data providers. These data providers either provide data directly as RDF or embed Microformats in their websites. To automatically find and evaluate small data providers, we extend an existing pipeline to integrate Linked Open Data [1].

3.2 Open Processing Pipeline

Information in legal documents is hidden in the text and needs to be extracted to produce legal data. With this legal data, we can, e. g., implement question answering systems and structured text search. Furthermore, legal data can be linked to different (external) data sources to provide, e. g., background information or geo-locations.

We provide technologies to minimize the effort for tasks that can be applied in a semi-automatic setting. More specifically, we are

⁸<http://gesetze-im-internet.de>, accessed: May 28, 2020

⁹<http://rechtsprechung-im-internet.de>, accessed: May 28, 2020

¹⁰<http://www.eur-lex.europa.eu>, accessed: May 28, 2020

¹¹<http://www.opendatamonitor.eu/>, accessed: May 28, 2020

¹²http://www.europa.eu/rapid/press-release_IP-16-1408_en, accessed: May 28, 2020

¹³<https://www.europeandataportal.eu/en>, accessed: May 28, 2020

⁷<http://www.djangoproject.com>, accessed: May 28, 2020

18. It must therefore be held, in the light of the case-law referred to in paragraph 11 of the present position, that the General Court did not err in law in finding, in paragraphs 62 to 64 of the order under appeal, that it was possible for M. Group to withdraw its appeal of 21 August 2007 before the Board of Appeal and, in paragraph 66 of that order, that such a withdrawal meant, as a result, that the Board of Appeal was no longer required to rule on the incidental submissions presented by C.

Figure 2: Legal NER: Named entities like involved parties (red), organisations (green) or dates (blue) are automatically extracted for text documents like court decisions.

interested in the tasks of (1) reference extraction, (2) entity extraction and linking, (3) keyword and title generation, (4) information retrieval, (5) and visualizing networks. The individual components are combined in our open processing pipeline, whereby the term “open” refers to the fact that each component is integrated either as Python package or as external service over the API. With this approach, the components act as building blocks and can also be used in other projects.

Reference Extraction. Reference extraction to legislation and judicial decisions is of great interest [3]. A network analysis build on top of citation can reveal decisions with great influence [10]. We extract citations with a hybrid approach that combines rule-based methods with learning-based methods¹⁴.

Named Entity Recognition & Entity Linking. Extracting named entities such as locations, courts, dates, and times is a well-known information extraction task. We implemented this task based on the SpaCy framework¹⁵ and trained a German NER model based on the dataset provided by Leitner et al. [8]. Extracted entities need to be disambiguated to provide further cross-connections between documents as well as to external data sources. For example, we link mentions of locations such as cities or states with open geo-information systems like Linked Geo Data¹⁶ using the Nominatim service¹⁷.

Keyword and Title Generation. Keywords accurately describing the content of documents are of great value for legal information systems. To generate keywords, we implement a module that combines rule-based methods with statistical methods. In Germany, court decision identifiers follow strict rules that allow determining the court and the general domain. Thus, we can extract the general domain of court decision by parsing the identifiers. Furthermore, referenced laws indicate a predominant legal domain, e. g., civil law. Finally, statistical methods like TF-IDF in combination with the-sauri (CF-IDF) allow to generate accurate keywords by analyzing the content (full-text) [6]. We apply CF-IDF to generate keywords from legal texts. Overall, we obtain keywords by parsing the identifier, analyzing the references to laws, and by analyzing the content

of the court decision. Finally, we combine these keywords to create human-readable titles.

3.3 Services

We equip users and developers with the necessary interfaces to efficiently interact with the documents and data.

Information Retrieval System. Finding topically related court decisions is a crucial task for legal professionals. We implemented a full-text search based on Elasticsearch. Furthermore, we developed a text- and citation-based recommender system that assists users in finding relevant information [12]. To facilitate research in this area, the platform provides an open interface such that novel methods can be evaluated with real users in A/B test experiments.

Working Environment. Typical users access legal information systems with a particular purpose in mind. While data analysis is an important task, more frequent users are interested in documents regarding a specific topic. Finding these legal documents is, however, only the first step in a more complex workflow. The working environment integrates the hypothesis framework¹⁸. Essential pieces of information and extracted entities are highlighted automatically. Furthermore, users can interact with the documents by highlighting their own key phrases or making notes.

Visualization. In Fig. 3, we present an excerpt of our citation network generated from German court decisions. As one can see, federal courts like the “Bundesverfassungsgericht”, the “Bundesgerichtshof”, and the “Bundesverwaltungsgericht” take a central role in the network. Furthermore, we observe that the “Bundesverwaltungsgericht” cites the statute book “Verwaltungsgerichtsordnung” (VwGO) the most. Our preliminary analysis indicates that many citations to the VwGO are to cite reasons for rejections of revisions. Moreover, the “Bundesgerichtshof” cites the statute book “Bürgerliche Gesetzbuch” (BGB) and the statute book “Zivilprozessordnung” (ZPO) the most. This indicates that most decisions the highest court in Germany are in the civil procedure. Citing the ZPO can also indicate a rejection of revision, e. g., ZPO Äg 561. Additionally, we observe a high amount of citations from the “Bundesverfassungsgericht” (BVerfG) to the statute book “Gesetz über das Bundesverfassungsgericht” (BVerfGG), but also to the “Strafgesetzbuch” (StGB) and the “Strafprozessordnung” (StPO). The BVerfGG contains statutes explicitly regulating the BVerfG. Citations to the StGB and the StPO indicate that the court decision is in the criminal procedure.

Overall, we made some interesting observations in the data by visualizing it, which motivates us to conduct an extensive analysis of our citation network. However, our analysis has several limitations with regard to the data. For one, the high amount of court decisions on federal level more likely reflects a more open publishing attitude of federal courts rather than an actual higher workload. Thus, it is important to address the publishing policy of courts and raise awareness of the benefits of open access in the justice domain.

¹⁴<https://github.com/openlegaldata/legal-reference-extraction>, accessed: May 28, 2020

¹⁵<http://www.spacy.io>, accessed: May 28, 2020

¹⁶<http://linkedgeodata.org/>, accessed: May 28, 2020

¹⁷<https://wiki.openstreetmap.org/wiki/Nominatim>, accessed: May 28, 2020

¹⁸<http://www.hypothes.is>, accessed: May 28, 2020

¹⁹<https://openlegaldata.io/assets/pdfs/citation-network.pdf>, accessed: May 28, 2020

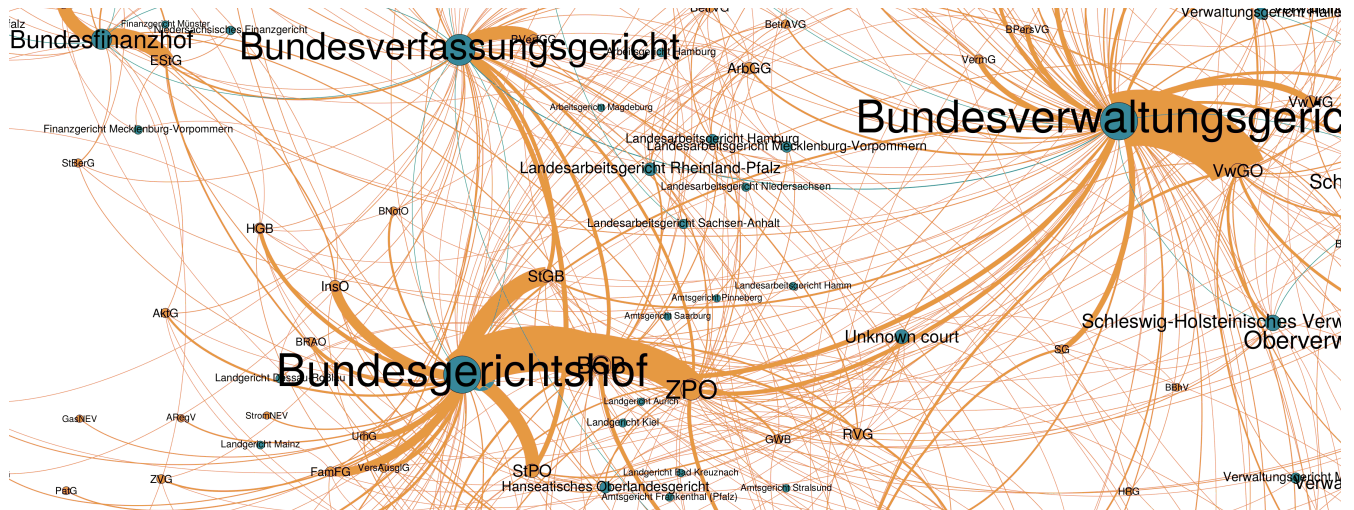


Figure 3: Excerpt of our citation network of German court decisions¹⁹. Blue vertices represent courts and orange vertices represent statutes books. The size of the vertices indicates the number of court decisions or number of statutes available. Blue edges visualize citations from a court decision to another court decision (of the respected court) and orange edges represent citations from court decisions to statute of the respected statute book. The width of the edges indicates the amount of citations.

4 CONCLUSION AND OUTLOOK

In this paper, we presented our approach to a single technology stack that allows accessing, processing, and providing legal information. We demonstrated that it is feasible to implement a variety of technologies in a single processing pipeline. We described in detail technologies that are implemented or currently being implemented in our open source project Open Legal Data. Furthermore, we published our first dataset of German court decisions²⁰. Based on this dataset, legal engineers developed the visual query interface VizLaw and were awarded the first place in the Berlin Legal Tech Hackathon 2019²¹. In conclusion, we see the Open Legal Data Platform as an important first step towards openness in the legal domain that will ultimately enable more collaboration among researchers and improve access to justice for the general public. In this context, we consider the MediaWiki software as role model, that powers all Wikipedias and helped to make encyclopedic knowledge available, and envision to achieve something comparable but for the legal domain. Making data technically open and accessible is only the beginning. In the future, we will also focus on Legal Design in order to make the data and information useful and usable not only for professional users, but also understandable and actionable for lay people.

ACKNOWLEDGMENTS

The research presented in this article is funded by the German Federal Ministry of Education and Research (BMBF) through the project Software-Sprint (grant no. 01IS16021).

REFERENCES

- [1] T. Blume and A. Scherp. 2019. FLUID: A Meta Model to Flexibly Define Schema-level Indices for the Web of Data. *arXiv preprint* (2019). arXiv:1908.01528
- [2] J. G. Conrad and F. Schilder. 2007. Opinion Mining in Legal Blogs. In *ICAIL*. ACM, 231–236.
- [3] E. De Maat, R. Winkels, and T. Van Engers. 2006. Automated Detection of Reference Structures in Law. *Frontiers in Artificial Intelligence and Applications* (2006), 41–50.
- [4] A. M. Fleckner and C. Coupette. 2018. Quantitative Rechtswissenschaft. *JuristenZeitung* 73, 8 (2018), 379.
- [5] M. Frosterus, J. Tuominen, and E. Hyvönen. 2014. Facilitating Re-use of Legal Data in Applications - Finnish Law as a Linked Open Data Service. *Legal Knowledge and Information Systems: JURIX 2014: The 27th Annual Conf.* (2014), 115–124.
- [6] L. Galke, F. Mai, A. Schelten, D. Brunsch, and A. Scherp. 2017. Using Titles vs. Full-text as Source for Automated Semantic Document Annotation. In *K-CAP*. ACM, 20:1–20:4.
- [7] T. J. Lampoltshammer, A. Guadamuz, C. Wass, and T. Heistracher. 2016. Open-laws.eu: Open justice in Europe through open access to legal information. *Achieving Open Justice through Citizen Participation and Transparency* (2016), 173–190.
- [8] E. Leitner, G. Rehm, and J. Moreno-Schneider. 2019. Fine-Grained Named Entity Recognition in Legal Documents. In *SEMANTECS*. 272–287.
- [9] M. Lissner. 2010. *CourtListener.com: A platform for researching and staying abreast of the latest in the law*. Master Thesis. University of California, Berkeley.
- [10] T. Neale. 2013. Citation Analysis of Canadian Case Law. *Journal of Open Access to Law* 1, 1 (2013), 1–51.
- [11] M. Palmirani and F. Vitali. 2011. Akoma-Netso for Legal Documents. In *Legislative XML for the Semantic Web*. Springer Netherlands, 75–100.
- [12] M. Schwarzer, C. Breiting, M. Schubotz, N. Meuschke, and B. Gipp. 2017. Citolytics: A Link-based Recommender System for Wikipedia. In *RecSys*. ACM, ACM Press, 360–361.
- [13] N. Shadbolt. 2011. Open For Business. *Think Quarterly: Data (UK)* (2011), 44–49.
- [14] M. Van Opijnen. 2011. European case law identifier: Indispensable asset for legal information retrieval. *Frontiers in Artificial Intelligence and Applications* (2011).

²⁰<https://openlegaldatal.io/research/2019/02/19/court-decision-dataset.html>, accessed: May 28, 2020

²¹<http://www.berlinlegal.tech>, accessed: May 28, 2020