

Towards Integrated Information Extraction and Facetted Search Applications in Nephrology

Danilo Schmidt¹, Hans-Jürgen Profitlich², and Daniel Sonntag²

¹ Nephrology Department
Charité - Universitätsmedizin Berlin
10117 Berlin, Germany

² German Research Center for Artificial Intelligence (DFKI)
66123 Saarbrücken, Germany

Abstract. This work focusses on our first integration steps of complex and partly unstructured medical data into a clinical research database. Our main application is an integrated facetted search tool in nephrology based on automatic information extraction results from textual documents. We describe the details of our technical architecture which is based on open-source tools—to be replicated at other universities, research institutes, or hospitals.

1 Introduction

As medical records may cover a very long history of diseases (up to 30 years) and include a vast number of diagnoses, symptoms, results, medications, and laboratory values, we could highly benefit from advanced search capabilities in clinical information systems to allow for the retrieval of relevant data. However, medical information systems often suffer from good search capabilities for data which has many unstructured text parts. Therefore, concepts to implement knowledge based systems, based on textual information extraction in medicine, are in focus of many recent research initiatives [11].

In this paper, we propose a three stage process: (1) offline textual information extraction from medical records in transplant medicine; (2) the generation of interesting facetted search capabilities on the results of the previous stage; (3) the combination of the information extraction results with structured laboratory values (ongoing work). Such a facetted search application uses techniques for accessing information organised according to a facetted medical classification system, allowing users to explore a collection of diagnoses, symptoms, results, medications, and laboratory values by applying multiple filters. Thus, facetted search allows clinicians to analyse complex data sets along a medical and cognitive (reflective) chain of decision-making; in particular, facetted search applications allow physicians to identify groups of patients with similar attributes. This can provide valuable decision support, when physicians are confronted with situations where rare or complex diseases require a high degree of specialist knowledge to filter and interpret (unstructured) medical data.

2 Background and Related Work

The faceted search application is based on the nephrology database TBase[®]. The web-based electronic patient record TBase[®] has been implemented in a German kidney transplantation programme as a cooperation between the Nephrology of Charité Universitätsmedizin Berlin and the AI Lab of the Institute of Computer Sciences of the Humboldt University of Berlin [3,10]. Currently, TBase[®] automatically integrates essential laboratory data (9.9 million values), clinical pharmacology (237.000 prescribed medications), diagnostic findings from radiology, pathology and virology (146.000 findings), and administrative data from the SAP-system of the Charité (70.000 diagnoses, 25.000 hospitalisations). Two groups of use cases for the application of faceted search in the medical field, and nephrology, can be identified: first, the use in clinical research, and second, the implementation in the individual treatment as a decision support system in the clinical routine.

Sacco [7] describes an approach of a guided interactive diagnostic system based on dynamic taxonomies. Biron et al. [2] describe an information retrieval system for computerised patient records. We extend these approaches by a special multi-facet functionality. Our approach shows the following main advantages:

- In our faceted search application, the user may remove *any* restriction he or she may have made in previous steps. This allows for a much better navigation through the search space where related systems only allow the subsequent thinning [8,9].
- The ranking of facet values by cardinality supports the survey of remaining subsets.
- We base automatically generated facets (e.g., disease/symptom relationships and negations) on multi-term extraction and relation extraction, by employing state-of-the-art, high-precision textual information extraction modules.

Only recently, new text mining approaches on Web-based medical literature have been proposed. For extracting adverse drug events from text [6] or automatic symptom extraction from texts on rare diseases [4], for example. However, clinical information extraction from patient records is still underrepresented and underdeveloped in clinical settings. Earlier work includes evaluating context features for medical relation mining on medical abstracts; the identification of semantic relations, such as substance A treats disease B, remains a non-trivial task [13]. Recent work and comparative baseline experiments include temporal information extraction [5]. A special trend becomes apparent, the need for ontology modelling of medical terminology and corresponding information extraction results [12]. Because of enormous annotation costs, mainly unsupervised methods are being used [1]. In industry and in the context of reliable clinical relevance, however, very detailed (and labor-intensive) supervised rule-based approaches represent the state-of-the-art.³

³ Here, we use our research project partner's solution (Averbis), which is based on shallow text parsing, see <https://averbis.com/en/research/>

3 System Architecture

The annotated texts are transferred in XMI format⁴ and stored in a local database at DFKI (see figure 1). Important components are the Solr search platform, the information extraction module, and the faceted search and presentation user interface modules. Solr⁵ is an open source enterprise search platform

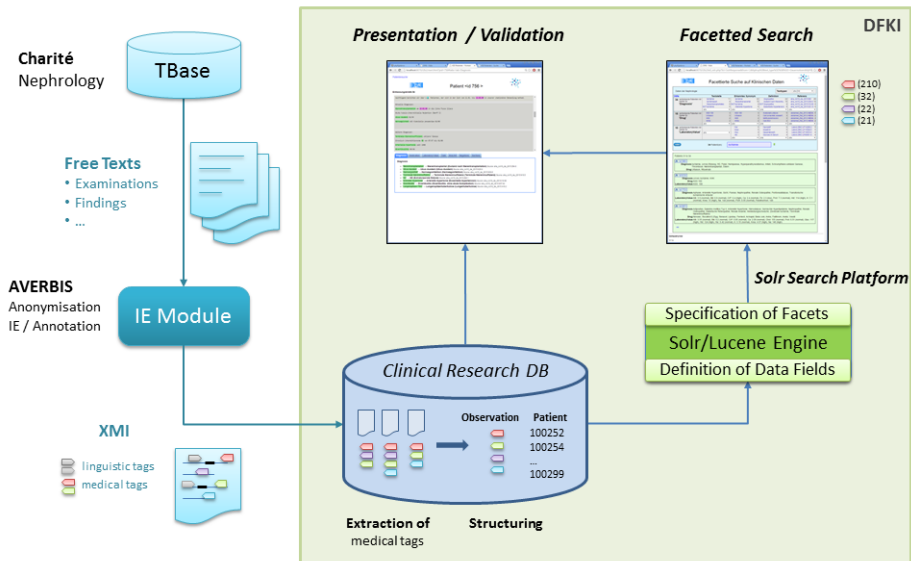


Fig. 1. System Architecture of the Integrated Information Extraction and Faceted Search Application

used in many large websites and applications and is one of the most popular enterprise search engines.⁶ Solr runs as a standalone full-text search server and uses the Lucene Java search library at its core for full-text indexing and (faceted) search. We chose the Solr system mainly because of some interesting features like faceted navigation, a query language that supports structured and textual search, the possibility for automatic result clustering based on Carrot2⁷, its scalability and extensibility through plug-ins, and its various APIs for input (text, xml, JSON, etc.) and output (JSON, XML, PHP, python, etc.).

The first step in our process is offline informative extraction. The text data for our system originate from the TBase[®] database of Charité Berlin containing medical information about nephrology patients. In the first phase we only

⁴ <http://www.omg.org/spec/XMI/>

⁵ <http://lucene.apache.org/solr/>

⁶ <http://db-engines.com/en/ranking/search+engine>

⁷ <http://project.carrot2.org/>

used about 5000 unstructured, free texts (no meta data or structured data of patients) of four types: "Befunde" (findings), "Untersuchungen" (visits), "Entlassungsbriefe" (clinical reports), and "Verläufe" (progress reports). These free texts are processed by the project partner Averbis, which anonymises the texts and adds annotations based on several medical reference systems and dictionaries (LOINC⁸, ICD10⁹, ABDAMED¹⁰).

A software module extracts the relevant medical tags and features and stores these in a database structure similar to the i2b2 star structure (in order to simplify the updates of the target system i2b2¹¹). The user interface to search and explore the annotated text database by using facets is built as a web service based on the Solr extension "solarium" for PHP systems.¹² This extension provides for an API to specify all parameters necessary to create complex Solr requests. The presentation/validation user interface of the system (see figure 2) consists of two parts: the upper part shows the original text with highlighted annotations, the lower part contains tabs listing the different relevant annotations. Clicking on an item in the lower part scrolls the text above to the corresponding

Untersuchungen 10000

Farbkodierte Dopplersonographie und Power-Doppler des **Nierentransplantates** vom **12.11.2001** :

Befund:

1. **postoperative Kontrolle** nach NTX.

Nierentransplantat im rechten Unterbauch, Volumen 100 **cm³**. Parenchymsaum nicht verschmälert und von regelrechter Echogenität mit gut abgrenzbaren Markkegeln. NBKS geschlossen, **kein** Konkrementnachweis.

Dopplersonographisch zeigt sich eine gute Durchblutung mit RI = 0,61 bis 0,7, gemessen an mehreren Segmentarterien. Im Bereich der arteriellen Anastomose **kein Nachweis** einer Stenose; vmax hier 80 **cm/s**. Im Nierenhilus **unauffälliges** arterielles Spektrum. Guter venöser Abstrom.

Ergebnis:

Morphologisch **unauffälliger Nierentransplantat** mit guter Durchblutung.
Kein Anhalt für Nierenarterienstenose.
Keine freie Flüssigkeit, **kein perirenales Hämatom**.

Diagnosis | Date | AhaUnit | Negations | Sections

Diagnosis

- **Nierentransplantates** → **Nierentransplantat (Zustand nach Nierentransplantation)** Source: aha_icd10_de_2015.Z94.0
- **postoperative Kontrolle** → **Postoperative Kontrolle (Nachuntersuchung nach chirurgischem Eingriff wegen anderer Krankheitszustände)** Source: aha_icd10_de_2015.Z09.0
- **Nierentransplantat** → **Nierentransplantat (Zustand nach Nierentransplantation)** Source: aha_icd10_de_2015.Z94.0
- **Nierentransplantat** (negated: 'unauffälliges' 6349) → **Nierentransplantat (Zustand nach Nierentransplantation)** Source: aha_icd10_de_2015.Z94.0
- **Nierenarterienstenose** (negated: 'Kein Anhalt für' 6410) → **Nierenarterienstenose (Nierenarteriosklerose)** Source: aha_icd10_de_2015.I70.1
- **perirenales Hämatom** (negated: 'kein' 6481) → **Perirenales Hämatom (Prellung und Hämatom der Niere)** Source: aha_icd10_de_2015.S37.01

Fig. 2. Presentation/Validation User Interface

⁸ <https://loinc.org/>

⁹ <http://www.icd-code.de/>

¹⁰ <http://www.wuv-gmbh.de/abdata-pharma-daten-service/datenangebot/abdamed/>

¹¹ <https://www.i2b2.org/about/intro.html>

¹² <http://www.solarium-project.org/>

position. The original XMI contents representing the complete original annotation information is shown in a pop-up window when a highlighted annotation in the text is clicked. Accordingly, this page serves two different purposes: (1) the presentation of the original text snippet found by the faceted search, and (2) the validation of the annotations.

4 Conclusion and Outlook

We demonstrated that new faceted search applications in the use case of transplant medicine in nephrology, based on open-source software tools and exchangeable information extraction modules, are feasible and a very suitable decision-support tool for the doctor: this type of a knowledge based system provides physicians with a practicable tool for the analysis of medical data and decision support for cohort selection. We developed a user interface for faceted search which is based on the Solr Engine. In the next project phase, we will extend the capabilities of the faceted search application, mainly including the following aspects: (1) integration of existing structural information about patients and treatments which includes numerical values, in relation to laboratory values or medications in particular; (2) extending the user interface by adding visual search and presentation techniques like "foamtree"¹³ to further facilitate the users exploration of the search space; (3) the integration of faceted search into special use cases moving towards individualised medicine [11].

Acknowledgements

This research is part of the project "clinical data intelligence" (KDI) which is founded by the Federal Ministry for Economic Affairs and Energy (BMWi).

References

1. Alicante, A.: Unsupervised entity and relation extraction from clinical records in italian. *Computers in Biology and Medicine* 72(1), 263–275 (2016)
2. Biron, P., Metzger, M., Pezet, C., Sebban, C., Barthuet, E., Durand, T.: An Information Retrieval System for Computerized Patient Records in the Context of a Daily Hospital Practice: the Example of the Leon Berard Cancer Center (France). *Applied Clinical Informatics* 5(1), 191–205 (2014)
3. Lindemann, G.: A web-based patient record for hospitals - the design of tbase2. In: Bruch, H.P. (ed.) *New Aspects of Hight Technology in Medicine: Hannover (Germany)*, pp. 409–414. Monduzzi Editore, International Proceedings Division (2000)
4. Métivier, J., Serrano, L., Charnois, T., Cuissart, B., Widlöcher, A.: Automatic symptom extraction from texts to enhance knowledge discovery on rare diseases. In: *Artificial Intelligence in Medicine - 15th Conference on Artificial Intelligence in Medicine, AIME 2015, Pavia, Italy, June 17-20, 2015. Proceedings.* pp. 249–254 (2015)

¹³ <https://carrotsearch.com/foamtree-overview>

5. Mkrtchyan, T., Sonntag, D.: Deep parsing at the CLEF2014 IE task. In: Working Notes for CLEF 2014 Conference, Sheffield, UK, September 15-18, 2014. pp. 138–146 (2014)
6. Odom, P., Bangera, V., Khot, T., Page, D., Natarajan, S.: Extracting adverse drug events from text using human advice. In: Artificial Intelligence in Medicine - 15th Conference on Artificial Intelligence in Medicine, AIME 2015, Pavia, Italy, June 17-20, 2015. Proceedings. pp. 195–204 (2015)
7. Sacco, G.: Guided interactive diagnostic systems. In: Computer-Based Medical Systems. pp. 117–122 (2005)
8. Sacco, G.: Dynamic taxonomies and guided searches. *Journal of the American Society for Information Science and Technology* 57(6), 792–796 (2006)
9. Sacco, G.: Dynamic taxonomies for intelligent information access. In: Khosrow-Pour, M. (ed.) *Encyclopedia of Information Science and Technology*, pp. 3883–3892. 3 edn. (2014)
10. Schröter, K.: Tbase2, a web-based electronic patient record. *Fundamenta Informaticae* 43(1-4), 343–353 (2000)
11. Sonntag, D., Tresp, V., Zillner, S., Cavallaro, A., Hammon, M., Reis, A., Fasching, A.P., Sedlmayr, M., Ganslandt, T., Prokosch, H.U., Budde, K., Schmidt, D., Hinrichs, C., Wittenberg, T., Daumke, P., Oppelt, G.P.: The clinical data intelligence project. *Informatik-Spektrum Journal* pp. 1–11 (2015)
12. Sonntag, D., Wennerberg, P., Buitelaar, P., Zillner, S.: Pillars of ontology treatment in the medical domain. *J. Cases on Inf. Techn.* 11(4), 47–73 (2009)
13. Vintar, S., Todorovski, L., Sonntag, D., Buitelaar, P.: Evaluating context features for medical relation mining. In: Proceedings of the ECML/PKDD Workshop on Data Mining and Text Mining for Bioinformatics (2003)