

Unleashing Semantics of Research Data

Florian Stegmaier¹, Christin Seifert¹, Roman Kern², Patrick Höfler²,
Sebastian Bayer¹, Michael Granitzer¹, Harald Kosch¹, Stefanie Lindstaedt²,
Belgin Mutlu², Vedran Sabol², Kai Schlegel¹, and Stefan Zwicklbauer¹

¹ University of Passau, Germany

² Know-Center, Graz, Austria

Abstract. Research depends to a large degree on the availability and quality of primary research data, i.e., data generated through experiments and evaluations. While the Web in general and Linked Data in particular provide a platform and the necessary technologies for sharing, managing and utilizing research data, an ecosystem supporting those tasks is still missing. The vision of the CODE project is the establishment of a sophisticated ecosystem for Linked Data. Here, the extraction of knowledge encapsulated in scientific research paper along with its public release as Linked Data serves as the major use case. Further, Visual Analytics approaches empower end users to analyse, integrate and organize data. During these tasks, specific Big Data issues are present.

Keywords: Linked Data, Natural Language Processing, Data Warehousing, Big Data.

1 Introduction

Within the last ten years, the Web reinvented itself over and over, which led from a more or less static and silo-based Web to an open Web of data, the so called Semantic Web¹. The main intention of the Semantic Web is to provide an open-access, machine-readable and semantic description of content mediated by ontologies. Following this, Linked Data [1] is the de-facto standard to publish and interlink distributed data sets in the Web. At its core, Linked Data defines a set of rules on how to expose data and leverages the combination of Semantic Web best practices, e.g., RDF² and SKOS³.

However, the Linked Data cloud is mostly restricted to academic purposes due to unreliability of services and a lack of quality estimations of the accessible data. The vision of the CODE project⁴ is to improve this situation by the creation of a web-based, commercially oriented ecosystem for the Linked Science cloud, which is the part of the Linked Data cloud focusing in research data. This ecosystem offers a value-creation chain to increase the interaction between all peers,

¹ <http://www.w3.org/standards/semanticweb/>

² <http://www.w3.org/RDF/>

³ <http://www.w3.org/2004/02/skos/>

⁴ <http://www.code-research.eu/>

e.g., data vendors or analysts. The integration of a marketplace leads on the one hand to crowd-sourced data processing and on the other hand to sustainability. By the help of provenance data central steps in the data lifecycle, e.g., creation, consumption and processing, along corresponding peers can be monitored enabling data quality estimations. Reliability in terms of retrieval will be ensured by the creation of dynamic views over certain Linked Data endpoints. The portions of data made available through those views can be queried with data warehousing functionalities serving as entry point for visual analytics applications.

The motivation behind the CODE project originated from obstacles of daily research work. When working on a specific research topic, the related work analysis is an crucial step. Unfortunately, this has to be done in a manual and time consuming way due to the following facts: First, experimental results and observations are (mostly) locked in PDF documents, which are out of the box unstructured and not efficiently searchable. Second, there exist a large amount of conferences, workshops, etc. leading to an tremendous amount of published research data. Without doubt, the creation of a comprehensive overview over ongoing research activities is a cumbersome task. Moreover, these issues can lead to a complete wrong interpretation of the evolution of a research topic. Specifically for research on ad-hoc information retrieval, Armstrong et al. [2] discovered in an analysis of research papers issued within a decade, that no significant progress has been achieved.

In contrast to scientific events, ongoing benchmarking initiatives such as the Transaction Processing Council (TPC) exist. The main output of the TPC is the specification and maintenance of high-impact benchmarks for the database technology with members from Oracle, Microsoft, Sybase etc. Obviously, the industries are interested in running these benchmarks to show their competitive abilities. The results of those runs are published on the TPC website⁵ as well as scientific workshops, e.g., Technology Conference on Performance Evaluation and Benchmarking (TPCTC)⁶. Unfortunately, it is currently very cumbersome to interact with the data to create comparisons or further visual analysis.

On the basis of this observations, the present issues could be improved by the use of the services established by the CODE project. Focusing on benchmarking initiatives, CODE technologies can be used to integrate the results of specific test runs, align them with extra information and therefore create an integrated TPC data warehouse to perform in depth analysis on the data, e.g., time series analysis.

This paper introduces the CODE project, along with its main processing steps. The main contributions are as follows:

- Data sources available in the research community will be described and the correlation to Big Data issues are given.
- The CODE project along with its main components is introduced with respect to the already defined Big Data processing pipeline.

⁵ <http://www.tpc.org/information/results.asp>

⁶ <http://www.tpc.org/tpctc/>

Table 1. Processable research data available in the CODE project

Type	Data Set Description	Data Characteristic
Research paper	PDF documents	Aggregated facts like tables, figures or textual patterns. Low volume, but high integration effort.
Primary research data	Evaluation data of research campaigns available in a spreadsheet like format (1. normal form) or via Web-APIs	Data generated by mostly automated means. Large volumes, low schema complexity.
Retrievable data	Linked Open Data endpoints	Semantically rich, interconnected data sets. Large volumes, hard to query (technically and from a usability point of view). Mostly background knowledge.
Embedded data	Microdata, Microformat, RDFa	Semantically rich, but distributed data. Less of interest.

The remainder of the paper is as follows: Section 2 highlights the data sources which can be processed by the CODE technologies. Here, an correlation to Big Data issues will be given. To get an understanding of the actual workflow, Section 3 proposes a processing pipeline, which is compliant to the overall definition of a *Big Data processing pipeline*. Finally, Section 4 concludes the paper and gives insights in the current achievements of the project.

2 Rediscovering Hidden Insights In Research

Research data is made available in various ways to the research society, e.g., stored in digital libraries or just linked to a specific website. Table 1 summarizes four data sources that are taken into account in the aforementioned usage scenario.

Research papers are a valuable source of state-of-the-art knowledge mostly stored in digital libraries reaching an amount of several Terabytes. Apart from the overall storage, the actual size of a single PDF document does not exceed a few Megabytes. The main task is to extract meaningful, encapsulated information such as facts, table of contents, figures and – most important – tables carrying the actual evaluation results. The present diversity of extracted data leads to a high integration effort for a unified storage. In contrast to that, primary research data is released in a more data centric form, such as table-based data. This kind of data is mostly issued by (periodic) evaluation campaigns or computing challenges. Famous examples are the CLEF initiative focusing on the promotion of research, innovation, and development of information retrieval

systems. The outcome of such activities is thousands of raw data points stored in Excel sheets. Here, the volume of the data is most likely very large but defined by a specific schema with less complexity than PDF documents. Both data sources share an unstructured nature due to missing semantics on the schema and data level. To overcome this issue, the two remaining data sources of Table 1 are utilized. In this light, Linked Open Data endpoints serve as source for retrievable data, such as DBpedia⁷ or PubMed⁸. On the one hand, these endpoints expose their data following the 5 star open data rule meaning the data is openly available, annotated with clear semantics and interconnected in the distributed Linked Open Data cloud. On the other hand, due to its distributed nature, efficient federated retrieval is a hard task. The last data source mentioned is embedded data meaning content of websites semantically annotated with microdata, microformat or RDFa. This information can be embedded table-based primary research data or auxiliary information, such as biographic data of a person.

As one can observe, there is a large amount of research data already available on the Web. The major drawback in this data landscape is the fact, that those are unconnected. Due to this fact, a comprehensive view is not possible, which leads to a loss of information. By the help of the CODE ecosystem, in particular by its data warehouse, this data gets connected and inference with respect to new knowledge is enabled.

Before considering the details of the knowledge extraction process, the correlation to the buzzword *Big Data* will be discussed. In todays research, the term Big Data⁹ [3] is often used as a fuzzy concept without clear defined semantics. The following dimensions, the “3Vs”, have to be mentioned when speaking of Big Data:

Volume is the most obvious characteristic for Big Data. Nearly every application domain produces an tremendous amount of data and is even increased by user interactions. This observation is also observable in terms of research data, when thinking of the amount of papers published with the corresponding monitored user interactions, such as citing.

Velocity makes it possible to state the production rate of the data. Huge data portions may be produced in real time in ongoing sensor systems, e.g., astronomy data, as an batch-like outcome of events, such as a conference or an evaluation campaign or single publications, such as white papers.

Variety takes the structure of the data itself into account. As already discussed, the data can be unstructured in silo-based PDF storage, semi-structured in Excel spreadsheets, or available in information retrieval systems.

Those three characteristics are commonly discussed by the community. It is clear, that Big Data at its core defines the data itself and the way it is processed and analyzed by corresponding pipelines. Linked Data on the other hand brings

⁷ <http://www.dbpedia.org/>

⁸ <http://pubmed.bio2rdf.org/>

⁹ <http://cra.org/ccc/docs/init/bigdatawhitepaper.pdf>

in the techniques to semantically interlink and publish this heterogeneous portions of data. A recent white paper [4] issued by Mitchell and Wilson extend those Vs with respect to a data centric way:

Value of the data is the key to real interpretation and knowledge generation by answering the question which interaction steps of a processing chain made the data portions really “worthy”.

The last characteristic can be directly aligned to the proposed approach. Here, crowd-sourced enabled data processing and analysis is combined with provenance chains to estimate the quality of the underlying data. By the help of Linked Data publishing techniques, the basis is given towards opening data silos for sophisticated interaction.

3 Big Data Pipeline Approach

When working with Big Data, Labrinidis and Jagadish [5] argue that “we lose track of the fact that there are multiple steps to the data analysis pipeline, whether the data are big or small”. The Big Data processing pipeline proposed by CODE in terms of knowledge extraction of research data is illustrated in Figure 1.

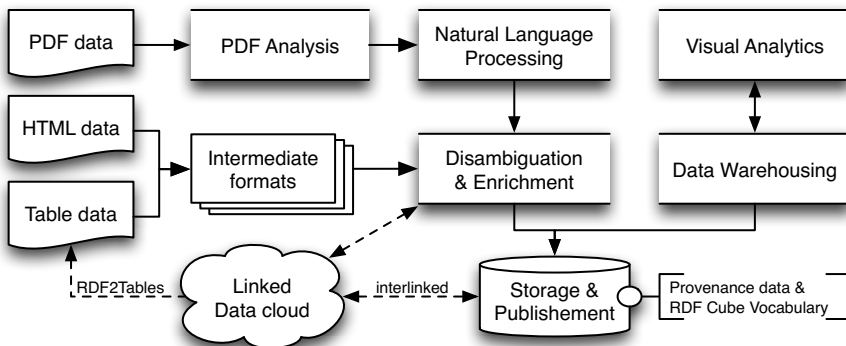


Fig. 1. Conceptual processing chain of knowledge creation and consumption

On the left hand side of Figure [5] the data sources introduced in Section 2 serve as an input for the conceptual processing chain. The data flow (continuous arrows) as well as dependencies (dashed arrows) are also plotted in the image. The central components are *PDF analysis*, *Natural Language Processing*, *Disambiguation & Enrichment*, *Data Warehousing* and *Visual Analytics* and will be discussed in the following.

3.1 PDF Analysis

Most of the research papers are stored in the PDF format. The quality of output of the PDF analysis thereby highly influences subsequent steps in the CODE processing chain. PDF is a page description language which allows low level control of the layout, but in this process the logical structure of the text is lost. For instance, text in multiple columns is often rendered across the columns, not adhering to the natural reading order. Especially tables are challenging because there is no annotation of logical tables defined in the PDF format. Still tables are assumed to contain lot of factual quantitative information. In general the challenges for PDF analysis can be summarised as:

- Text content extraction, extracting raw textual content (ignoring images and tables).
- Metadata extraction, e.g. extracting author names, titles, journal titles for scientific publications.
- Structure annotation, annotating document structure, e.g. for generating automatic table of contents.
- Block detection, detection of logical blocks like tables, abstracts.
- Table decomposition, extraction of table data according to its logical structure.

In recent years considerable research progress has been made with regard to these challenges. Text content extraction methods are able to extract text in human-reading order [6]. Metadata extraction already quite well extracts relevant metadata from scientific papers [7, 8]. Block detection has been approached [8], but especially the extraction of complex tables is in the focus of ongoing research [9, 10].

Despite the progress in the single steps, there is no general solution which can provide all information in the quality needed within the CODE project in sufficient quality. Thus, the task is to aggregate results from recent research on PDF analysis into the CODE prototype and adapt or refine existing approaches. Further, we expect manual post-processing to be necessary for achieving certain analysis results.

3.2 Natural Language Processing

Based upon the textual representation of a research article, the contained facts should be mined. Therefore techniques from the field of natural language processing are employed. As an initial step, named entities within the text are identified. Depending on the actual domain of the articles (biomedical domain, computational science, ...) the type of named entities varies.

Domain adaptation in the CODE project is foreseen to be transformed into a crowd-sourcing task. For example, in the computer science domain, where ontologies and annotated corpora are scarce, the users of the CODE platform themselves annotate the relevant concepts. Starting with the automatic detection of named entities, the relationship between those are identified in a second

step. This way the textual content is analysed and domain dependant, factual information is extracted and stored for later retrieval.

3.3 Disambiguation and Enrichment

Entity disambiguation is the task of identifying a real world entity for a given entity mentioning. In presence of a semantic knowledge base, disambiguation is the process of linking an entity to the specific entity in the knowledge base.

Within the CODE project, entity disambiguation is applied to identify and link scientific knowledge artefacts mentioned in scientific papers. Subsequently background information from the Linked Science cloud can be presented to the user while reading or writing scientific papers.

The challenges regarding entity disambiguation within the CODE project are the following: (i) variance and specificity of scientific domains: not only do scientific papers cover a wide variety of topics but each domain very in-depth; (ii) synonyms in Linked Data repositories, and (iii) evolving knowledge: topic changes in scientific papers and in Linked Data endpoints.

Disambiguation using general purpose knowledge bases (mostly Wikipedia) has been widely covered in research, e.g. [11–13]. While approaches for specific knowledge bases exist, e.g. [14] for biomedical domain, the applicability of the approaches to a combination of general and specific knowledge bases and the resulting challenges (scalability, synonyms) has to be investigated within the CODE project.

After disambiguation, the gathered information for an entity can be extended by knowledge available in the Linked Data cloud. This extra information will be validated by user feedback and then integrated into the knowledge base. This process yields to an automatic and intelligent Linked Data endpoint facing the following research tasks: (i) integration and usage of provenance data, (ii) ranking and similarity estimations of Linked Data repositories or RDF instances, and (iii) quality of service parameter (e.g., response time). This process is often termed Linked Data Sailing. Currently, there exist frameworks to calculate similarity between Linked Data endpoints, e.g., SILK [15], and Linked Data traversal frameworks, e.g., Gremlin¹⁰, which serves as a basis for further developments.

3.4 Storage and Publishing

The persistence layer of the CODE framework consists of a triple store, which has to offer certain abilities: (i) Linked Data compatible SPARQL endpoint and free text search capability, (ii) federated query execution, e.g., SPARQL 1.1 federated query¹¹, and (iii) caching strategies to ensure efficient retrieval. Those requirements are fulfilled by the Linked Media Framework [16], which has been selected for storage. For data modelling tasks, two W3C standardization efforts are in scope, which will be soon issued as official recommendations. The

¹⁰ <https://github.com/tinkerpop/gremlin/>

¹¹ <http://www.w3.org/TR/sparql11-federated-query/>

PROV-O¹² ontology will be used to express and interchange provenance data. Further, the W3C proposes the RDF Cube Vocabulary¹³ as foundation for data cubes, which are the foundation of data warehouses. Both vocabularies will be interconnected to ensure a sophisticated retrieval process.

3.5 Data Warehousing

As already mentioned, the basis for OLAP functionalities is the data cube. The data cube model is a collection of statistical data, called observations. All observations are defined by dimensions along with measures (covering the semantics) and attributes (qualify and interpret the observation). Well-known data warehousing retrieval functionalities would last from simple aggregation functions, such as *AVG*, up to high-level *roll up* or *drill down operators*. During retrieval the following functionality has to be ensured: (i) interconnection of RDF cubes, (ii) independence of dimensions, and (iii) high-level analytical retrieval in graph structures. Current research is dealing with the integration of RDF data into single data cubes [17, 18], but do not take an interconnection / federation into scope. Within the CODE framework, algorithms of relational data warehousing systems will be evaluated with respect to their applicability to graph structures. By the help of data cube interconnections complex analytical workflows can be created.

3.6 Visual Analytics

One important aspect of the CODE project is to make data available to end users in an easy-to-use way. This data might be already Linked Data as well as semantic data extracted from scientific PDFs. The goal is to build a web-based Visual Analytics interface for users who have no prior knowledge about semantic technologies. The main challenges regarding Visual Analytics in the scope of the CODE projects are:

- building an easy-to-use web-based interfaces for querying, filtering and exploring semantic data,
- developing semantic descriptions of Visual Analytics components to facilitate usage with semantic data, and
- building an easy-to-use web-based interfaces for creating visual analytic dashboards.

A query wizard is envisioned, with which users can search for relevant data, filter it according to their needs, and explore and incorporate related data. Once the relevant data is selected and presented to the user in tabular form, the Visualization Wizard helps them to generate charts based on the data in order to make it easier understandable, generate new insights, and communicate those insights in a visual way. One of the tools for visualizing the data will be MeisterLabs' web-based MindMeister mind mapping platform.

¹² <http://www.w3.org/TR/prov-o/>

¹³ <http://www.w3.org/TR/vocab-data-cube/>

4 Conclusion

In this paper the challenges of the CODE project have been outlined. Further, the connection and the relevance to Big Data topics has been argued. In the current phase of the project, prototypes for certain issues of the introduced pipeline have been developed¹⁴. Within the second year of the project, those will be integrated into a single platform. Periodic evaluations will be conducted to ensure the required functionality and usability of the prototypes.

Acknowledgement. The presented work was developed within the CODE project funded by the EU Seventh Framework Programme, grant agreement number 296150. The Know-Center is funded within the Austrian COMET Program under the auspices of the Austrian Ministry of Transport, Innovation and Technology, the Austrian Ministry of Economics, Family and Youth and by the State of Styria. COMET is managed by the Austrian Research Promotion Agency FFG.

References

1. Bizer, C., Heath, T., Berners-Lee, T.: Linked data – the story so far. *International Journal on Semantic Web and Information Systems* 5(3), 1–22 (2009)
2. Armstrong, T.G., Moffat, A., Webber, W., Zobel, J.: Improvements that don't add up: ad-hoc retrieval results since 1998. In: *Conference on Information and Knowledge Management*, pp. 601–610 (2009)
3. Dumbill, E.: What is big data? An introduction to the big data landscape. O'Reilly Strata (January 11, 2012), <http://strata.oreilly.com/2012/01/what-is-big-data.html>
4. Mitchell, I., Wilson, M.: *Linked Data - Connecting and exploiting Big Data*. White Paper (March 2012), [http://www.fujitsu.com/uk/Images/Linked-data-connecting-and-exploiting-big-data-\(v1.0\).pdf](http://www.fujitsu.com/uk/Images/Linked-data-connecting-and-exploiting-big-data-(v1.0).pdf)
5. Labrinidis, A., Jagadish, H.V.: Challenges and opportunities with big data. *PVLDB* 5(12), 2032–2033 (2012)
6. Hasan, I., Parapar, J., Barreiro, Á.: Improving the extraction of text in pdfs by simulating the human reading order. *Journal of Universal Computer Science* 18, 623–649 (2012), http://www.jucs.org/jucs_18_5/improving_the_extraction_of
7. Granitzer, M., Hristakeva, M., Knight, R., Jack, K., Kern, R.: A comparison of layout based bibliographic metadata extraction techniques. In: *Proceedings of the 2nd International Conference on Web Intelligence, Mining and Semantics, WIMS 2012*, pp. 19:1–19:8. ACM, New York (2012)
8. Kern, R., Jack, K., Hristakeva, M.: *TeamBeam - Meta-Data Extraction from Scientific Literature*. *D-Lib Magazine* 18 (July 2012)
9. Fang, J., Gao, L., Bai, K., Qiu, R., Tao, X., Tang, Z.: A table detection method for multipage pdf documents via visual separators and tabular structures. In: *2011 International Conference on Document Analysis and Recognition (ICDAR)*, pp. 779–783 (September 2011)

¹⁴ <http://www.code-research.eu/results>

10. Liu, Y., Bai, K., Gao, L.: An efficient pre-processing method to identify logical components from pdf documents. In: Huang, J.Z., Cao, L., Srivastava, J. (eds.) PAKDD 2011, Part I. LNCS (LNAI), vol. 6634, pp. 500–511. Springer, Heidelberg (2011)
11. Kataria, S.S., Kumar, K.S., Rastogi, R.R., Sen, P., Sengamedu, S.H.: Entity disambiguation with hierarchical topic models. In: Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2011, pp. 1037–1045. ACM, New York (2011)
12. Fader, A., Soderl, S., Etzioni, O.: Scaling wikipediabased named entity disambiguation to arbitrary web text. In: Proc. of WikiAI (2009)
13. Dredze, M., McNamee, P., Rao, D., Gerber, A., Finin, T.: Entity disambiguation for knowledge base population. In: Proceedings of the 23rd International Conference on Computational Linguistics, COLING 2010, Stroudsburg, PA, USA, pp. 277–285. Association for Computational Linguistics (2010)
14. Rebholz-Schuhmann, D., Kirsch, H., Gaudan, S., Arregui, M., Nenadic, G.: Annotation and disambiguation of semantic types in biomedical text: a cascaded approach to named entity recognition. In: Proceedings of the EACL Workshop on Multi-Dimensional Markup in NLP, Trento, Italy (2006)
15. Volz, J., Bizer, C., Gaedke, M., Kobilarov, G.: Discovering and maintaining links on the web of data. In: Bernstein, A., Karger, D.R., Heath, T., Feigenbaum, L., Maynard, D., Motta, E., Thirunarayan, K. (eds.) ISWC 2009. LNCS, vol. 5823, pp. 650–665. Springer, Heidelberg (2009)
16. Kurz, T., Schaffert, S., Bürger, T.: LMF – a framework for linked media. In: Proceedings of the Workshop on Multimedia on the Web Collocated to i-KNOW/i-SEMANTICS, pp. 1–4 (September 2011)
17. Kämpgen, B., Harth, A.: Transforming statistical linked data for use in olap systems. In: Proceedings of the 7th International Conference on Semantic Systems, I-Semantics 2011, New York, NY, USA, pp. 33–40. ACM (2011)
18. Zhao, P., Li, X., Xin, D., Han, J.: Graph cube: on warehousing and olap multidimensional networks. In: Proceedings of the International Conference on Management of Data, pp. 853–864 (2011)