# Interfaces: Accessing Biographical Data and Metadata

## Matthias Reinert, Bernhard Ebneth

Historical Commission at the Bavarian Academy of Sciences and Humanities, Munich;
Head of department "*Deutsche Biographie*": Malte Rehbein (University of Passau)
reinert@hk.badw.de, ebneth@ndb.badw.de

### Abstract

Based on the principles of interaction and cooperation in the field of biographical lexicography, the paper outlines the design and setup of interfaces. Interfaces enable access to both individual and sets of biographical data entries. This paper briefly describes the design of metadata mapping to semantic formats (RDF) and graphs, when metadata is provided by the historical biographical information system *Deutsche Biographie*. It demonstrates exemplary usages of the APIs in the historical sciences.

**Keywords:** biographical metadata, interfaces, semantic metadata mapping, graph databases

## 1 Lexicography as interaction, cooperation, and exchange

Recent biographical dictionaries have emerged out of scientific debate and exchange on the international and interregional levels and in cross-disciplinary discourse. The majority of current biographical dictionaries have moved content – everything from indices to complete articles – to digital media or were set up to be born digital.[1]

In addition, digital media has made it easy to reference ongoing historiographical research, critical editions, documentary projects, library catalogs, and archival sources.

The common tasks of all lexicographers consist of extending, revising, and correcting the biographical database. Therefore, collaboration through the sharing and linking of information forms the backbone of future biographical dictionaries.

The digital transformation is confronted with the challenge of maintaining formal and scientific standards that have evolved in the past 100 years of lexicography. Current biographical dictionaries share a common set of elements or modules, e.g., names, dates, places of birth and death, family background, biographical description, and references to works, archival resources, secondary literature, portraits and authorship. These could probably be described as a common web ontology.

The similar structure of biographical articles in dictionaries allows similar strategies of visualization – ranging from static genealogical trees to dynamic relations to persons, institutions and places.

Authority files appear to be key to the interfaces of exchange and to ontology patterns and modes of visualization. They were conceived in the bibliographical field and proved valuable in interlinking both persons and concepts.(Ebneth and Reinert, 2018)

## 2 Data Resources.

The *Deutsche Biographie* is a joint endeavor of the Historical Commission at the Bavarian Academy of Sciences and Humanities and the Bavarian State Library in Munich.

Since the late 1990s, several measures have been taken to digitize two series of biographical dictionaries and to establish a website providing access to them.

The latest efforts have targeted cultural institutions, in order to enlarge the number of notable individuals in the corpus (Reinert et al., 2015).

### Components

The aggregated database consists of parts differing in provenance, density, and granularity of information:

- Digitized biographical texts of the Allgemeine Deutsche Biographie (ADB, 55 vols. and index vol. published 1875–1912) and Neue Deutsche Biographie (NDB, 26 [A–Vocke] of 28 vols. published since 1953) provide information on individuals, exact dates of birth and death, places of birth, death, and burial, as well as partially encoded information on entities related to the individual. The full text of both series have been digitized and structured in XML (Reinert, 2010).

- The index of persons and families mentioned in ADB and NDB was compiled manually and provides names, years of birth and death, a hierarchy of professions, and references to pages in the printed volumes. Almost all index entries are aligned with the German authority file (Gemeinsame Normdatei GND),[2] and further information can be derived from this resource (Hockerts, 2008), (Ebneth, 2009), (Busch and Jordan, 2011).

---

[1]In Europe there are only 9 national dictionaries accessible freely online: Historisches Lexikon der Schweiz – Dizionario storico della Svizzera – Dictionnaire historique de la Suisse (HLS), Neue Deutsche Biographie (NDB), Österreichisches Biographisches Lexikon 1815–1950 (ÖBL) mit 2. überarbeiteter Auflage ab 1815, Slovenska biografija (SBL), Nationaal Biografisch Woordenboek (NBW), Dizionario Biografico degli Italiani (DBI), Biographie Nationale de Belgique (BNB), Nouvelle Biographie Nationale numérisées de Belgique (NBNB), Kansallisbiografia (Finland), Norsk Biografisk Leksikon (NBL), Internetowy Polski Słownik Biograficzny (iPSB). Two others offer paid access: the Oxford Dictionary of National Biography and the Dictionary of Irish Biography.

[2]A complete dump is available under CC0 `http://www.dnb.de/lds` (Pfeifer, 2015).

- The main work base of the editorial office helps to curate the last two volumes that are to appear in print. A subset of this work base that includes all entries on deceased persons with a GND-identifier is merged with the *Deutsche Biographie*. The dataset includes differences in name spellings, pseudonyms, dates and places of birth and death. The places are rarely harmonized and aligned. While all entries are provided with a GND-identifier, further information from this resource can be derived (Hockerts, 2012).
- The last component is a set of persons and families – amounting to about 600,000 entries – provided by websites of cooperating partners[3]. The data for these entries are imported from the German national authority file GND (Ebneth, 2012), (Kraus et al., 2014).

**Data enhancement**

Metadata enhancement is crucial for extended search options. The main approach is to detect entities and determine the linguistic class they belong to, then identify entities against a given database, and finally detect relations and sentiments expressed with regard to them (Jurafsky and Martin, 2016).

The *Deutsche Biographie*' approach relies heavily on authority files. The index of persons was completely equipped with GND-entries and -identifiers in cooperation with the Bavarian State Library / Munich Digitization Center (Hockerts, 2008), (Busch and Jordan, 2011). This close cooperation makes the *Deutsche Biographie* a reference for biographical entries in the GND too.

With the appearance of each volume in print, the index is enlarged and enhanced in term of GND-metadata.

As soon as a new volume is published in print form, the previous volume will be transformed to deeply structured XML and put online.

**Entity detection**

The structure of the articles converted from PDF to XML first covers the main parts of the article, namely the headline, the genealogy, the life summary, and the technical parts listing awards, works, sources, secondary literature, and portraits.

It then deals with entities, like personal names occurring in verb phrases ("interpersonal relations") (Stotz and Reinert, 2013), (Stotz et al., 2015). The strategy of our choice, named "Local Grammars," is described in (Gross, 1997), (Geierhos, 2007), (Geierhos, 2010) and relies on dictionaries as described by (Guenthner and Maier, 1994), (Guenthner and Maier-Meyer, 1996). We set up different dictionaries of partial and complete forms of well known entities (first names, surnames, names of places, disciplines, fields of study, relevant adjectives and noun phrases).

---

[3]The partnering institutions provide individual personal information on a selection of individuals who are representative for their collecting focus and documentary field. A complete list can be found at https://www.deutsche-biographie.de/partner

The "Local Grammars" at hand are capable of describing institutional bodies (enterprises, educational institution, theatre and music groups, political parties), administrative geographical regions (populated places, regions, countries, religious territories) and the proper names of certain individual places (monasteries, churches, castles).

Finally biographical accounts refer to highly specific events (s. (Rospocher et al., 2016)). Of these only the proper names of the most prominent events like congresses, wars, and peace treaties have been directly detected with certain "Local Grammars."



Figure 1: The web version of Albert Einstein's article, https://www.deutsche-biographie.de/sfz68290.html, written by Max von Laue, first appeared in NDB 4 (1959). The RDF version linked on the left to https://www.deutsche-biographie.de/downloadRDF?url=sfz68290.rdf.

**Entity linking/disambiguation**

To enhance this database, different strategies were applied to link named entities to the internal database of indexed names. One obvious strategy relied on index entries and calculated scores for identified personal names on a given page in a given section (genealogy or biography), depending on the length of the named entity and any given birth or death dates as compared with those accounted for in the index of names for that page (Reinert et al., 2015).

The second strategy drew on professional descriptors that preceded the named entity and compared them with the most well-known index entry, which assumed that celebrity implied being named or connected with other famous personalities who had also been portrayed in the biographical dictionary.

Another strategy consisted of detecting places of birth, death and burial and linking them to authority files. These place names were prominent in the headlines of each biography of an individual person and taggable by regular expressions. Twelve thousand distinct names were uniquely detectable.
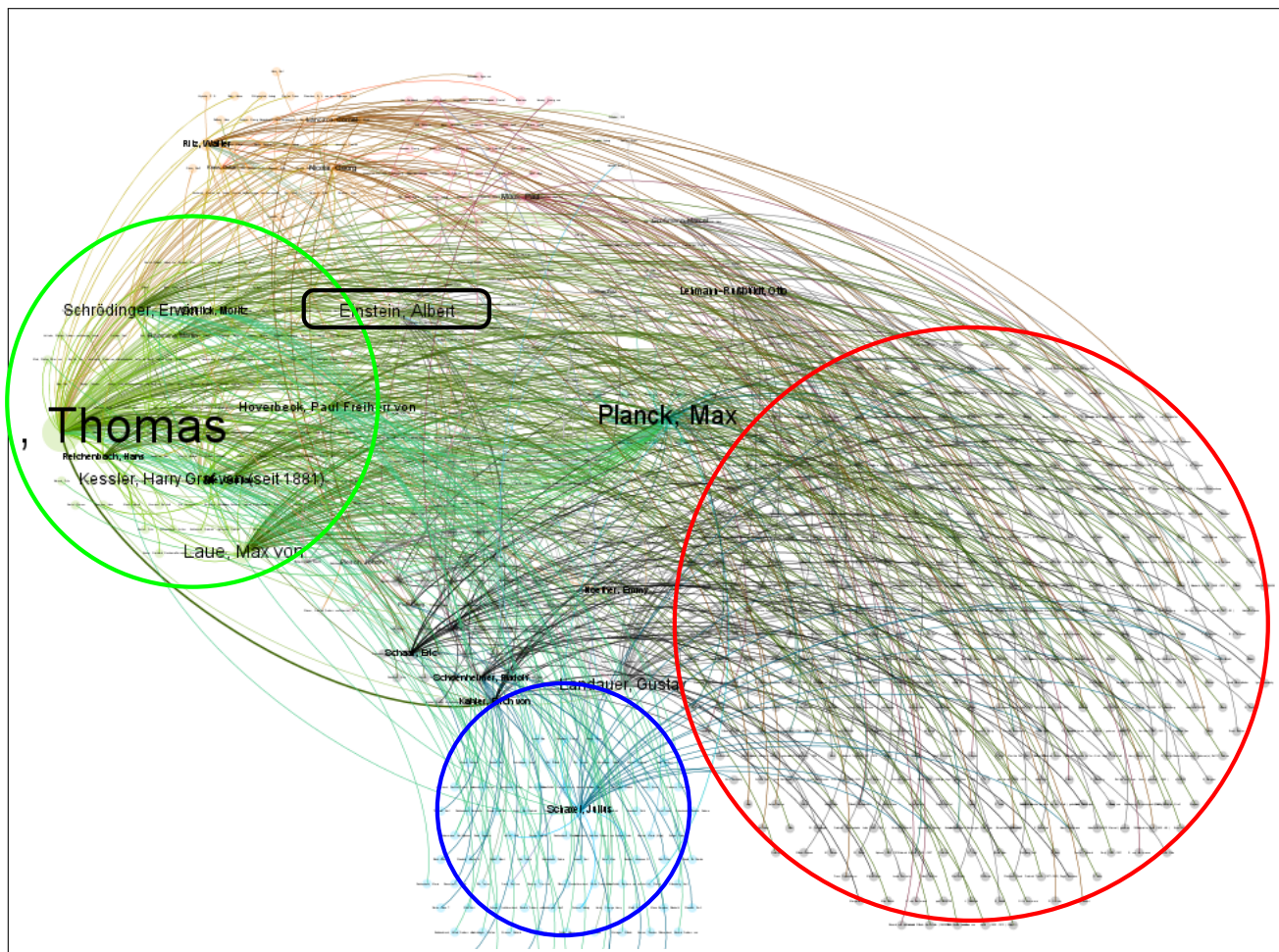
Figure 2: Albert Einstein's network up to a distance of 2 with interconnections, circle pack layout (Mike Bostock, in Gephi) coloured by religion/denomination. Red/right circle contains unidentified names, green/left circle represents (well connected) protestants, blue/center bottom circle represents Jewish and the upper middle (not circled) contains persons changing/not given denomination (like Einstein). Data `https://data.deutsche-biographie.de/graph-open/?id=sfz68290&depth=2` [3.11.2017].

By accessing the web services of OpenStreetMap (OSM)[4] and filtering the results by type of place and continent, about one third could be identified directly in OSM.

The majority of remaining places occurred only once in the corpus and had to be aligned to OSM manually. Although OSM does not provide unique IDs, the coordinates and some additional information were reusable.

With the help of the geographical coordinates, about half of these entries could be identified / reverse-located in Geonames[5]. Geonames provides unique and stable identifiers, but the qualification of the type of the place differs. Even though over 36,000 GND-entries for place- or regionlike administrative territories had been aligned to Geonames, the data was incoherent. The GND referred to administrative division in Geonames, while we were using populated places.

**Information extraction – summary**

Up to now, the following have been recognized and identified (linked to index database) in the *Deutsche Biographie*:

- 340,000 personal names tagged in the corpus of ADB and NDB, of which about 145,000 are identified as database entries;
- 380,000 place names, of which 12,500 are places of study, 1,900 "places of worship" (religious buildings); over 163,000 place names were identified as an instance of a given place in the database;
- 103,000 relations found in 21,330 biographies (NDB), of which 6,200 are teacher-student relations, 1,500 friendship predicates, about 1,000 predecessor/successor predicates, and another 3,000 membership relations and about 3,100 leadership relations relating to institutions or organizations;
- 12,000 organizational names;
- 8,000 time expressions;
- 11,500 mentions of discipline/fields of study.

---

[4] `http://nominatim.openstreetmap.org`.
[5] `http://www.geonames.org/`.

The efforts in geo-locating places led to new search options. Two start pages were introduced that provided a mapsearch and a zoomable geographical distribution of places mentioned in the *Deutsche Biographie*.

## 3 Interfaces

The notion of what an interface is ranges from intuitive user interaction surface to machine readable bit-stream. In this paper, interface is understood as a programmable data-delivering web application.

**Type, purpose, coverage of the APIs**

In *Deutsche Biographie* we mapped metadata on the entry level. This means that the RDF is directly accessible from the web-version of an individual article (s. fig. 1). Another way of mapping relations to a graph format is described below (s. sect. 3).

The APIs provided by `http://data.deutsche-biographie.de` split up into two groups: the Beacon-interface / Beacon-Aggregator and the Solr-based search index cover about 750,000 aggregated entries in the *Deutsche Biographie* whereas the RDF-based SPARQL-endpoint and the Neo4J-based graph-database cover only persons and families mentioned in the biographical dictionaries (about 100,000). The number of entries linked to others with an explicitly named relation shrinks to about 23,000, due to the limitations of named entity identification mentioned above.

The Beacon aggregation provided here offers a list of links to resources for a given GND-Identifier (s. fig. 3).



Figure 3: The first lines of the Beacon aggregation for Albert Einstein GND-ID: 118529579) `http://data.deutsche-biographie.de/rest/bd/gnd/118529579/alle_de.`

**Beacons and their aggregation**

The Beacon concept examined here was conceived by volunteers and Wikipedia-enthusiasts during a conference in Munich in 2010.[6] There is an informal description[7] that allows for aggregation of Beacons of different origin. Most Beacon files are announced at the given Wikipedia website, others are hosted by the Historical Commission or at Findbuch.de[8].

```
 #FORMAT: BEACON
#PREFIX: http://d-nb.info/gnd/
#TARGET: http://www.deutsche-biograph
ie.de/pnd{ID}.html#ndbcontent
118643525
118500015
...
```

Figure 4: First lines of Beacon file.

**Solr**

The Solr[9] index is configured for combined searches, auto-complete and faceted searches for names, places, and professions. Relations between personal entries are stored in an abbreviated serialized form (s. fig. 5). The decision to

```
...<arr name="bez"><!-- sfz@type@name -->
<str>sfz112378@Verwandt@Einstein, Maria
(Maja)</str>
<str>sfz68291@Verwandt@Einstein,
Alfred</str>
<str>sfz58393@Verwandt@Maric, Mileva</str>
<str>sfz107772@Verwandt@Maric,
Mileva</str></arr> ...
```

Figure 5: Serialised relations in index entries `http://data.deutsche-biographie.de/beta/solr-open/?q=defgnd:118529579.`

open the Solr index was supported by similar activities in the library sector,[10] reflecting the need for easily accessible APIs.

**RDF**

Different approaches were proposed in the field of semantic data modelling for biographies:

- a property-based approach is one in which persons are entitled with atomic properties that may be organized in classes (cf. GND approach with (Litz et al., 2012),
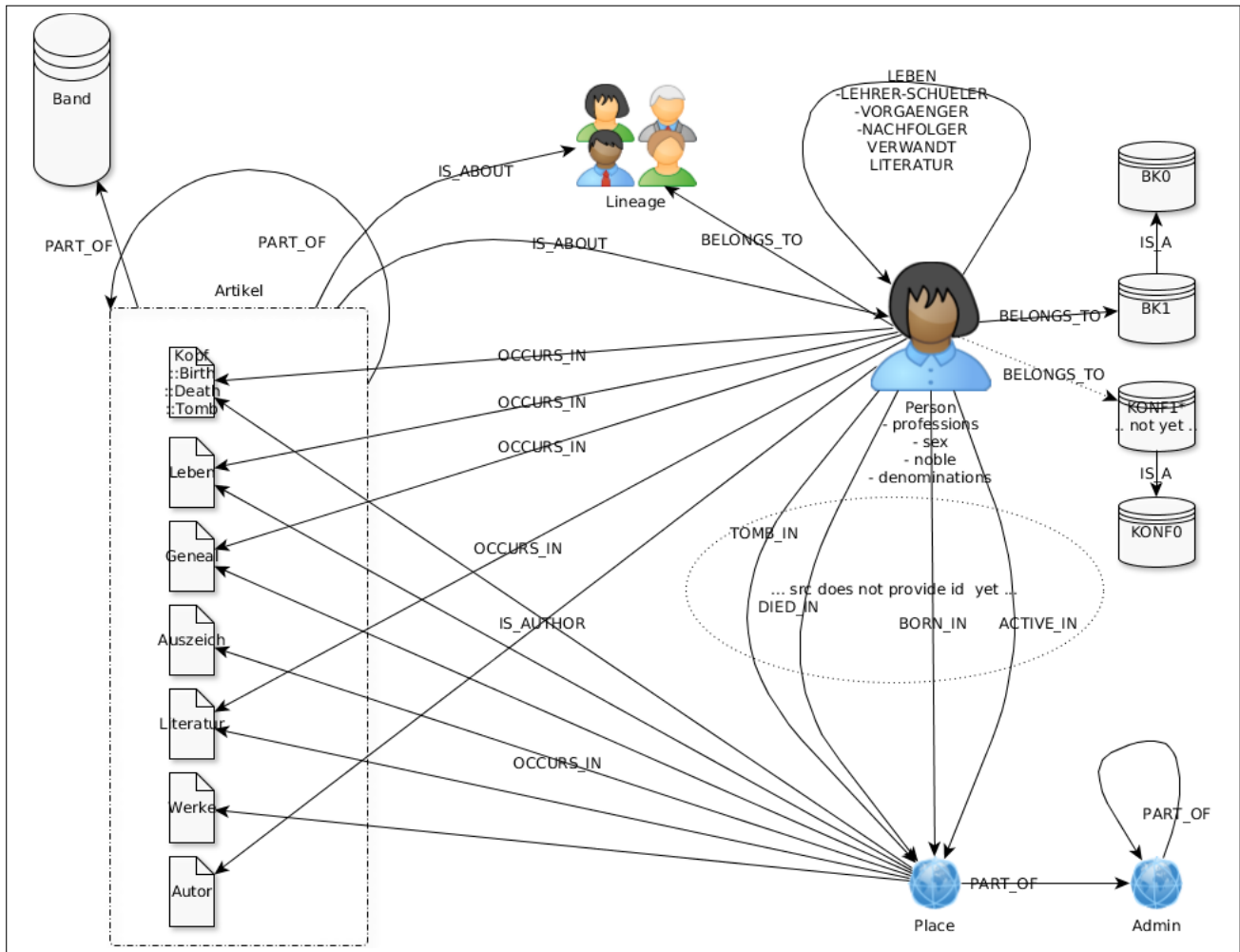
Figure 6: Model of NDB-Graph stored in Neo4J. Drawn with yEd.

the "factoids" proposed by (Michele Pasin and John Bradley, 2013) and the "aspects" deployed at the Personendatenrepositorium (Walkowski, 2009).

- event-based ontologies represent a person's life as a set of timespans or events of varying duration. Examples for this approach are the Erlangen CRM[11], (Trame et al., 2013) and the "life-spans" deployed at the Catalogus professorum Lipsiensis s. (Riechert et al., 2010).

The initial data-model of *Deutsche Biographie* was property based. It draws heavily on the concept of the GND-elementset[12] and initial work of (Litz et al., 2012), see (Brümmer, 2011). We publish dumps of all generated RDF-triples periodically and serve a Sparql-endpoint on request.

**Graph APIs**

The API provides access to a graph database based on Neo4J[13]. The graph is modelled according to the index of persons and the structure of the printed volumes and its articles. These hierarchically ordered units (`<is_part_of>`) of text refer to entries in the person and place databases by means of edges (`<occurs_in, is_author>`). Persons may refer to a lineage (a group of family-like related individuals) `<belongs_to>` or directly to others by named edges (`<lehrer, vorgaenger, nachfolger>`). A person is connected to places (`<born_in, died_in, tomb_in, active_in>`) and to a profession classification (`<belongs_to>`).

A second option makes it possible to export an ego-centered graph, based on the detected relations to identified personal names in the Solr index.[14]

The example presented in fig. 2 shows a data sample for

Albert Einstein and all individuals in his network up to a distance of 2, together with their cross-relations amongst each other. The dataset has easily been imported into Gephi[15], colored by an indexed category and visualized using a circle packing layout (Mike Bostock).

## 4 Summary

Recent incoming e-mail requests have showed a growing interest in accessing the data. There were genuine research questions (e.g., the network of Martin Luther, investigating German street names), integration tasks (RDF query of metadata for identified individual entries) and bridging tasks (querying personal names in a back-end for an archival CMS).

A proposal has been launched to raise funding for a web-based research laboratory that would work with the metadata encoded and aggregated in the *Deutsche Biographie*.

## 5 Acknowledgements

## 6 References

Martin Brümmer. 2011. Realisierung eines RDF-Interfaces für die Neue Deutsche Biographie. In Sören Auer, Johannes Schmidt, and Thomas Reichert, editors, *SKIL 2011 – Studentenkonferenz Informatik Leipzig 2011, Leipziger Beiträge zur Informatik Band XXVII, Leipziger Informatik-Verbund (LIV)*, pages 31–42.

Thomas Busch and Stefan Jordan. 2011. Vernetzte Lebensläufe. Der Einsatz von Normdatenbanken zur Verlinkung biographischer und bibliographischer Angebote im Internet. *Geschichte in Wissenschaft und Unterricht*, (11/12):684–691.

Bernhard Ebneth and Matthias Reinert. 2018. Potentiale der Deutschen Biographie (www.deutsche-biographie.de) als historisch-biographisches Informationssystem. In Ágoston Zénó Bernád, Christine Gruber, and Maximilian Kaiser, editors, *Europa baut auf Biographien. Aspekte, Bausteine, Normen und Standards für eine europäische Biographik.*, pages 283–295, Wien. new academic press.

Bernhard Ebneth. 2009. Vom digitalen Namenregister zum europäischen Biographie-Portal im Internet. In Martina Schattkowsky and Frank Metasch, editors, *Biografische Lexika im Internet*, volume 15 of *Bausteine aus dem Institut für Sächsische Geschichte u. Volkskunde*, pages 13–44. Dresden.

Bernhard Ebneth. 2012. Aktueller Stand der Genealogien in der Neuen Deutschen Biographie – Arbeit mit der Online-Version. In *64. Deutscher Genealogentag*, Augsburg.

Bernhard Ebneth. 2015. Auf dem Weg zu einem Historisch-biographischen Informationssystem. Datenintegration und Einsatz von Normdaten am Beispiel der Deutschen Biographie und des Biographie-Portals. *Jahrbuch für Universitätsgeschichte*, pages 261–290.

Michaela Geierhos. 2007. *Grammatik der Menschenbezeichner in biographischen Kontexten*. Arbeiten zur Informations- und Sprachverarbeitung. Band 2. München.

Michaela Geierhos. 2010. *BiographIE. Klassifikation und Extraktion karrierespezifischer Informationen*. Number 05 in Linguistic Resources for Natural Language Processing. Lincom, München.

Maurice Gross. 1997. The Construction of Local Grammars. In E. Roche and Y. Schabès, editors, *Finite-State Language Processing*, pages 329–354. Cambridge, Mass.

Franz Guenthner and Petra Maier-Meyer. 1996. Das CISLEX-Wörterbuchsystem. In H. Feldweg and E. W. Hinrichs, editors, *Lexikon und Text*, pages 69–82. Tübingen.

Franz Guenthner and Petra Maier, editors. 1994. *Das CISLEX Wörterbuchsystem*. München.

Hans Günter Hockerts. 2008. Vom nationalen Denkmal zum biographischen Portal: ADB und NDB. *Akademie Aktuell*, 33(2):19–22.

Hans Günter Hockerts. 2012. Zertifiziertes biographisches Wissen im Netz. Forschungsnahe Informationsinfrastruktur. Die „Deutsche Biographie" auf dem Weg zum zentralen historisch-biographischen Informationssystem für den deutschsprachigen Raum. *Akademie Aktuell*, 37(4):34–36.

Dan Jurafsky and James H. Martin. 2016. *Speech and Language Processing*.

Hans-Christof Kraus, Marco Jorio, Martina Schattkowsky, Bernhard Ebneth, Matthias Reinert, Thierry Declerck, Christine Gruber, and Eva Wandl-Vogt. 2014. Sektion: Vernetzung von historisch-biographischen Lexika und Fachportalen im Linked (Open) Data Framework. In *1. Jahrestagung der Digital Humanities im deutschsprachigen Raum (DHd 2014), Universität Passau, 25.-28.3.2014*.

Berenike Litz, Aenne Löhden, Jan Hannemann, and Lars Svensson. 2012. AgRelOn – An Agent Relationship Ontology. In Juan Manuel Dodero, Manuel Palomo-Duarte, and Pythagoras Karampiperis, editors, *Research Conference on Metadata and Semantic Research*, volume 6, pages 202–212. https://de.slideshare.net/larsgsvensson/agrelon-an-agent-relationship-ontology.

Michele Pasin and John Bradley. 2013. Factoid-Based Prosopography and Computer Ontologies: Towards an Integrated Approach. *Literary and Linguistic Computing*.

Barbara Pfeifer. 2015. Über Zweck und Nutzen der

---

Gemeinsamen Normdatei (GND). *Jahrbuch für Universitätsgeschichte*, 16:251–259.

Matthias Reinert, Maximilian Schrott, and Bernhard Ebneth. 2015. From Biographies to Data Curation-The Making of Www. Deutsche-Biographie. De. In *Biographical Data in a Digital World (BD)*, pages 13–19, Amsterdam.

Matthias Reinert. 2010. Biographisches Wissen auf einen Klick. *Akademie Aktuell*, 35(4):44–46.

Thomas Riechert, Ulf Morgenstern, Sören Auer, Sebastian Tramp, and Michael Martin. 2010. The Catalogus Professorum Lipsiensis – Semantics-Based Collaboration and Exploration for Historians. In *Proceedings of the 9th International Semantic Web Conference (ISWC2010)*, Shanghai/China.

Marco Rospocher, Marieke van Erp, Piek Vossen, Antske Fokkens, Itziar Aldabe, German Rigau, Aitor Soroa, Thomas Ploeger, and Tessel Bogaard. 2016. Building Event-Centric Knowledge Graphs from News. *Web Semantics: Science, Services and Agents on the World Wide Web*, 37-38:132–151, March.

Sophia Stotz and Matthias Reinert. 2013. Detecting and Encoding Interpersonal Relations with Unitex/Local Grammars. Université Paris Est-Marne-la-Vallée.

Sophia Stotz, Valentina Stuß, Matthias Reinert, and Maximilian Schrott. 2015. Interpersonal Relations in Biographical Dictionaries. A Case Study. In *Biographical Data in a Digital World (BD)*, pages 74–80.

Johannes Trame, Carsten Keßler, and Werner Kuhn. 2013. Linked Data and Time — Modeling Researcher Life Lines by Events. In *Proceedings of the 11th International Conference on Spatial Information Theory - Volume 8116*, COSIT 2013, pages 205–223, New York, NY, USA. Springer-Verlag New York, Inc.

Niels-Oliver Walkowski. 2009. Zur Problematik der Strukturierung und Abbildung von Personendaten in digitalen Systemen. In *Workshop Personendateien der Arbeitsgruppe „Elektronisches Publizieren" der Union der deutschen Akademien der Wissenschaften*.