

Humboldt-Universität zu Berlin
Philosophische Fakultät
Institut für Bibliotheks- und Informationswissenschaft

Modelling and Automated Retrieval of Provenance Relationships

Masterarbeit im Rahmen des Weiterbildenden Masterstudiengangs
Bibliotheks- und Informationswissenschaft im Fernstudium

vorgelegt von
Thomas Schneider

Gutachter:
Prof. Dr. Robert Jäschke
Christian Rüter

Erfurt, den 1. Juni 2023

Todo list

<input type="checkbox"/> Discuss <i>all</i> of Q1–Q9 ?	3
<input type="checkbox"/> Discuss SNA/HNA and general significance of relations; see intro report SoNAR AP 2 . . .	3
<input type="checkbox"/> Further thoughts on example queries:	3
<input type="checkbox"/> be more precise	4
<input type="checkbox"/> possibly adapt this paragraph later	4
<input type="checkbox"/> Some more comments on NEL and choice of data sources from [Menzel, Schnaitter, et al. 2021; Meiners 2022]? see Exposé	8
<input type="checkbox"/> Also mention ontology-based data integration? See https://en.wikipedia.org/wiki/Ontology-based_data_integration and references therein	12
<input type="checkbox"/> elaborate, depending on whether provenance plays a larger role in subsequent chapters	14
<input type="checkbox"/> omit this subsection?	14
<input type="checkbox"/> Discuss instances of Q1 , Q3–Q9 , and their differences (extend remarks from Section 1.1)?	17
<input type="checkbox"/> Explain this further? A case for ontologies!	18
<input type="checkbox"/> give examples of commonalities (OPAC) and differences (discovery vs. OPAC)?	18
<input type="checkbox"/> Unify typesetting of terms (i.e., Term instead of “Term”); unify (FRBR) relation names (e.g., has_creator vs. creator)	31
<input type="checkbox"/> Formulate decision problem, discuss (data) complexity & algorithms. (Reduction to FO/SQL?)	35
<input type="checkbox"/> Relate this “machinery” to the example queries. In particular:	35
<input type="checkbox"/> TODO: Discuss further extensions:	35
<input type="checkbox"/> list data and link to repository; paste and fill in data management plan)	52
<input type="checkbox"/> finish DMP	53

Contents

1	Introduction	1
1.1	Background	1
1.2	Aim and Research Question	3
1.3	Methods and Outline	4
2	Context	5
2.1	Network Analysis and the SoNAR Project	5
2.2	Linked Data and Data Integration	11
2.3	Provenance Indexing	14
3	Example Queries and Answers: A Case Study	15
3.1	Example Queries	15
3.2	Manual Answer Retrieval	16
3.3	The Quality of Query Answers	17
3.4	From Manual to Automated Query Answering	18
4	Analysis of Available Data Sources and Techniques	20
4.1	Standards for the Description of Bibliographic Resources	20
4.2	Data Formats and Communication Protocols	21
4.3	Data Sources	23
4.4	Conclusions	28
5	A General Model of Provenance Relationships	29
5.1	Basic Notions	29
5.2	Labelled Directed Graphs	30
5.3	Modelling Data Sources, Queries, and Answers	32
5.4	Decision Procedures	35
5.5	Discussion of the Modelling Decisions	35
5.6	Possible Extensions	35
6	Automated Retrieval of Provenance Relationships	37
7	Conclusion	38
	References	39
	Bibliography	39
	Web Resources	43
	Acronyms	50
A	Research Data Management Plan (in German)	52
B	Selbstständigkeitserklärung	54

List of Figures

4.1	FRBR entities and basic relations of Groups 1 and 2	21
5.1	A directed graph	30
5.2	A labelled directed graph that represents data concerning an exemplar of Copernicus' <i>De revolutionibus</i> and some of its owners	31
5.3	A graph representing example query Q2'	33
5.4	Positive (a,b) and negative (c, d) examples for query answers	34
5.5	An example homomorphism	35

1 Introduction

1.1 Background

Provenance research is concerned with the origins and ownership history of cultural objects. Its main objective is the reconstruction of “object biographies” in the historical context. Application areas of provenance research include the study of private and public collections, the detection of forgery, and the identification and restitution of loot. Since the release of the *Washington Conference Principles on Nazi-Confiscated Art* [Web1] in 1998, provenance research has received increased attention.¹

Regarding the holdings of university and research libraries, particularly interesting provenances are those related to the change of owners when a book copy is passed on or distributed [Hakelberg 2016, p. 2]. Provenances can be reconstructed by marks of ownership such as stamps, bookplates (ex libris), or handwritten signatures: with the help of these features, it is possible to retrace the “history” of a book item or the extent of library holdings that have been scattered in the meantime [Hakelberg 2016, p. 2].

In order to enable provenance research, libraries index the provenances of their historical holdings in their catalogues. A provenance entry contains an entity such as a person, corporate body, or collection, and a feature that indicates ownership. The identification and disambiguation of entities is achieved via links to records in authority files such as the “Gemeinsame Normdatei”, the Integrated Authority File of the German National Library (GND) [Web2]. The indexed provenance data makes it possible, for example, to query and reconstruct the items owned or held by a single person, to query the whereabouts of relevant items, or to retrace the distribution of all indexed exemplars derived from a given work.

Nowadays, provenance entries are recorded in electronic catalogues. German university and research libraries typically do not maintain their own individual catalogue; rather they are provided with a central catalogue by the library network in which they participate.² These central catalogues are equipped with an underlying database and standardised data formats for internal representation and data export. For example, the German library networks GBV and SWB maintain and use the common catalogue (database) K10plus, which internally uses the data format PICA [Web4] and allows for exports in the data formats MARC 21, MAB2, and Pica+ [Web5]. Despite the use of a uniform data format, there are several possibilities to record provenance entries. As Hakelberg [2016, §4] explains, libraries even within the same network often use diverse representations for the same type of provenance entry, and the differences are considerable: for example, some GBV libraries record their provenance entries in data fields on the bibliographic level, while others use the level on the exemplar level. These deviations lead to large differences in the presentation of the holdings in the online catalogue, which hinders the retrieval of relevant items and historical holdings.

Data about persons and corporate bodies are recorded in Germany- or worldwide authority files and further databases, such as the GND, the Library of Congress Name Authority File (LCNAF)

¹↑ This paragraph is a brief summary of the introductory chapter in Zuschlag’s Introduction to Provenance Research [2022, §1].

²↑ There are six library networks (*Bibliothekverbünde*) for scientific libraries in Germany [Web3].

[Web6], the *WorldCat* [Web7], databases of projects such as International Standard Name Identifier (ISNI) [Web8] and Virtual International Authority File (VIAF) [Web9], or in Wikidata [Web10]. These data sources usually support standardized bibliographic or generic data formats for data export. Data about a person contain, among others, the name and alternative name forms (which can be manifold in the case of, e.g., scholars of previous centuries), places of birth, death, and work, as well as relationships to corporate bodies and other persons (e.g., coauthors and students). The extent of a dataset of the same person can differ between data sources, which is witnessed, for example, by the entries on the scholar Georg Joachim Rheticus in *ISNI*, *VIAF*, *WorldCat*, *GND*, and Wikidata.³ Hence, the state of data on persons and corporate body is heterogeneous as well, and depending on the concrete individual, it may be necessary to consult several data sources and combine the obtained data.

Given this diversity and heterogeneity of the existing data sources, it is currently difficult to retrieve provenance relationships which require data that is distributed over several sources. This is in fact the case for a number of research questions that can arise in cultural, historical, or provenance research. The following list contains some examples for such questions, which were obtained in personal communication with researchers.⁴

- Q1** Who read work *X*, in which manifestation and in which year?
- Q2** Which exemplars⁵ of work *X* were passed from one of its owners to a student?
- Q3** What are the relationships between the recipients of manifestation *Y* of work *X*?
- Q4** Which exemplars from a collection *X* were passed on by its owner to a family member?
- Q5** Which exemplars from the holdings of library *X* were acquired from bookseller *Y* between 1933 and 1945?
- Q6** Who participated in the sale of collection *X*?
- Q7** Via which paths did exemplars from collection *X* enter library *Y*?
- Q8** Which libraries own the exemplars once owned by person *X*?
- Q9** Where did person *X* acquire exemplars and did they know the previous owners?

In these examples, variables *X*, *Y* are used as placeholders for arbitrary works, manifestations, collections, etc. Therefore, **Q1–Q9** are actually query *patterns*, each of which represents a set of possible queries that can be obtained by instantiating the variables with concrete objects. For the introductory purposes of this chapter, we continue to use the placeholders and refer to **Q1–Q9** simply as queries. We will get back this distinction in Chapter 3.

Query **Q1** addresses works as well as their manifestations (e.g., editions of the same work in various languages). Answering this query would help trace the reception of the same work over several eras. For example, Duchess Luise Dorothea of Saxe-Gotha-Altenburg read French editions of English works by John Milton and Alexander Pope.⁶ Obviously, there is a difference between the action “read” used in this query and the relation “owned” represented by provenance (entries). We neglect this difference for the moment and will get back to it in Chapter 3.

³↑ This can be verified by following, at the end of the Wikipedia page on Rheticus in the table “Authority Control” [Web11], the links to *ISNI*, *VIAF*, *WorldCat*, and *GND*, and inspecting the Wikidata dataset [Web12].

⁴↑ Dietrich Hakelberg (head of Dept. “Collection Development and Cataloguing”, Research Library Gotha of the University of Erfurt) [Web13], Michaela Scheibe (deputy head of Dept. “Manuscripts and Early Printed Books”), [Web14], and Joëlle Weis (head of Research Area “Digital Literary and Cultural Studies”, University of Trier), [Web15].

⁵↑ Conforming to the FRBR model [IFLA SG FRBR 2009], the precise wording should be “exemplars of manifestations of expressions of *X*” but, here and in the following, we omit the intermediate entities for brevity when no misunderstanding is expected.

⁶↑ See, for example, the provenance entries of the respective exemplars in the Research Library Gotha [Web16].

Queries **Q2** and **Q3** aim at highlighting the network that spans between the recipients of a work. For example, one of the two exemplars of Nicolaus Copernicus’s main work *De revolutionibus orbium coelestium* [Copernicus 1543] that are now held by the Gotha Research Library of the University of Erfurt have been owned, by several scholars from the circle around the author, some of which were in the teacher-student relation. This information can be concluded from the accounts of Gingerich [Gingerich 2002, p. 69] and Salatowsky and Lotze [Salatowsky and Lotze 2015, p. 142], but it can also be obtained by looking up the entry and its owners in the electronic catalogue of the library, following the links to the **GND** entries of the owners, and inspecting the relationships between the owners in **GND**; see Chapter 3 for a more detailed derivation.

An important difference between Queries **Q2** and **Q3** is the following. While **Q2** fixes a relation between persons (“collaborator” or “student”) and asks for works in whose context this relation has instances, **Q3** does not fix a particular relation but asks for the entire context of the work or a manifestation or exemplar thereof.

Query **Q4** is very similar in structure to **Q2**. Query **Q5** is important in the context of research on Nazi loot.

Discuss *all* of **Q1–Q9**?

Discuss SNA/HNA and general significance of relations; see intro report SoNAR AP 2

When attempting to answer queries such as **Q1–Q9**, it does not suffice to consult a single data source such as a catalogue, authority file, or knowledge base. Instead, it is necessary to consult several data sources and combine the information found. This process is highly laborious, given not just the number of data sources but also their diverse data formats. Therefore, automated support is essential. This is one of the reasons for Hakelberg to raise the following question [Hakelberg 2016, p. 46, translated from German]: “How can historical provenance relationships be formulated and represented in a machine-readable way?”

In order to implement suitable tools, it is necessary to analyse available data sources, data models, and data integration techniques, to develop an abstract model of data sources and possible queries, and to devise a method for obtaining answers in this abstract framework.

Further thoughts on example queries:

- While the difference between persons and institutions may be negligible from a modelling point of view, it is highly relevant concerning the amount of query answers and their handling.

1.2 Aim and Research Question

In this thesis, we pursue the goal of facilitating the automated retrieval of provenance relationships. More precisely, we want to develop a method for answering provenance queries that refer to bibliographic entities, people, and corporate bodies as well as the relationships between those. This method should consult standard data sources such as library catalogues, authority files, and knowledge bases. The method should furthermore be implementable as a retrieval system that supports the user in formulating their queries, answering them, and exploring the data that supports the query answers. In the long term, we envisage that such a retrieval system will support provenance research by prospectively retrieving potentially interesting constellations.

This goal leads to the following central research question for this thesis.

How can provenance relationships be modelled and automatically retrieved?

This question implies several subordinate questions:

- SQ1** *Which data sources are available for answering provenance queries?*
- SQ2** *Which techniques and tools are available for the integration of data from heterogeneous sources?*
- SQ3** *Based on the structure of the identified data sources, how data sources, queries, and answers be modelled in an abstract framework?*
- SQ4** *What is a suitable method for retrieving provenance relationships in that framework?*

1.3 Methods and Outline

In order to answer our research questions, we will proceed as follows.

In Chapter 2, we review existing work and connect it with our research questions. In Chapter 3, we revisit the exemplary queries from Section 1.1, demonstrate a manual attempt at answering them, and discuss the expected quality of query answers and difficulties with this manual process. This discussion is used as a point of reference for the subsequent considerations.

In order to answer Questions **SQ1** and **SQ2**, we review data sources, data models, and data integration techniques in Chapter 4. Based on the results of this analysis, we develop an abstract model of data sources, queries, and answers in Chapter 5. This mathematical model subsumes the previous examples while being vastly more general: it gives a formal description of how to build admissible queries, without restricting their contents (i.e., the specific names, attributes, concepts, and relations used) or their complexity. Using this model, we hope to provide a means for library science, which is the origin of the research question, to benefit from established methods from mathematics and computer science.

Based on our model and the insights from the preceding analysis, we develop in Chapter 6 a method for answering queries that can serve as the basis of a retrieval system as indicated in the previous section. Finally, we draw conclusions in Chapter 7.

be more
precise

possibly
adapt
this para-
graph
later

2 Context

This thesis aims at providing support for provenance research, which is part of historical research. As we have seen in Chapter 1, networks that represent people and their relationships play a central role. Therefore, we begin our exploration of the scientific context of our work with the topic of historical (social) network analysis. In Section 2.1, we collect insights from a recent research project aimed at providing a research infrastructure for historical network analysis and apply them to our research questions.

It has also become clear in Chapter 1 that the relevant data is distributed over a multitude of heterogeneous data sources. Therefore we review the literature on linked data, data integration, and data provenance in Section 2.2.

Finally, we give a brief overview of the state of provenance indexing with authority data in Section 2.3.

2.1 Network Analysis and the SoNAR Project

2.1.1 Network Analysis

According to Jansen [2003], the notion of a *network* is a central tool for the analysis of modern societies in sociology, political science, and economics. In these disciplines, networks of political actors, companies, or researchers, among many others, are the subject of study. Additionally, networks play an important role in organisational psychology, biology, and web science [Web17; Web18]. The following paragraph briefly summarises the main constituents of network analysis as described by Jansen [2003].

Network analysis defines networks and provides a statistical toolset for describing and analysing them. A network is a graph, which consists of nodes and edges. Nodes represent actors, events or objects (e.g., people, companies, institutions, countries), and edges represent relations between those (in the case of social networks, e.g., friendship, collaboration, or family relationships). Networks are defined via certain modelling decisions such as the commitment to a set of actors and relations that are of relevance, or the decision whether relations are directed or undirected and whether they are dichotomous (an edge is present or not) or weighted by values representing frequencies or extent. The tools provided by network analysis include various metrics that apply to single nodes, pairs of nodes, paths in the network, or the whole network; those often come in several variants. Examples are connectivity, in-/outdegree, density, size, cohesion, multiplexity, reachability, and path distance. Since a graph can conveniently be represented by its adjacency matrix, methods from matrix algebra are part of the toolset [cf. Jansen 2003, §§1.1, 3.3, 5.3].

In applications where social structures are the object of study, the term *social network analysis* (SNA) is used, and it predominantly refers to the *method* of investigating social interactions [Otte and Rousseau 2002].

In the context of historical research, the notion of *historical network analysis* (HNA) has emerged recently; it focusses on the reconstruction of historical networks [Menzel, M.-J. Bludau, et al. 2020]. A distinguishing feature of HNA is its retrospective view, i.e., it is used to analyse historical data

extracted from available sources [Fangerau et al. 2022]. Menzel et al. [2020] name “a lack of awareness with regard to the availability of suitable research data” as a limiting factor of HNA; in particular, a large amount of data is distributed over heterogeneous data sources.

2.1.2 SoNAR (IDH)

The project “SoNAR (IDH), Interfaces to Data for Historical Social Network Analysis and Research” [M. Bludau et al. 2020; Menzel, M.-J. Bludau, et al. 2020; Web19], which was funded by the German Research Foundation (DFG) from 2019 to 2021, developed approaches to building “an advanced research technology environment supporting Historical Network Analysis and related research” [Web19]. In the long-term vision, that environment is expected to integrate data from a variety of existing repositories, thus providing researchers with an extensive, standardised, and interregional infrastructure for answering research questions using methods from HNA. According to the project proposal and the final reports [Web20], the project participants undertook a systematic analysis of processing and managing the source data for the purposes of HNA, designed a model of a structured data analysis for HNA based on SNA methods, developed visualization approaches and interfaces, evaluated all components for a scientific usage, and developed a concept for implementation and operation.

One of the project’s four work packages (WPs) addresses the development of the research technology environment and its evaluation against real research questions. In the remainder of this subsection, we summarise the insights described in the final report on this WP [Fangerau et al. 2022] that are relevant for the research goals of this thesis.

The report emphasises relationships and their evolution over time as central aspects of historical research and acknowledges the increased importance of methods of HNA. In this context, the authors note that visualisation plays an important role as a means for restructuring information and thus contributes to the progression of knowledge. They also address the *hermeneutic circle* [Malpas and Gander 2015] as a general approach to answering historical research questions: this term refers to an iterative, circular work process where the initial question guides the work with the source material and is, in turn, readjusted based on the answers obtained.

In order to evaluate the research technology, the project participants developed two research scenarios with several research questions. The report names four concrete questions that were considered central. The nature of these questions is very general and “global” in the sense that they refer to a large part of a network: for example, they ask for the point in time when a scientific area became a separate discipline, or for the role of academic and familial links in the course of a given scientist’s career.

In this context, the report distinguishes two approaches to developing research questions in HNA. One of these is the explorative approach, where suitable questions are developed based on the available data. According to the authors, this approach includes serendipity, i.e., the hope that an inspection of the data helps identify unexpected phenomena that contribute to the shaping of the research question.

The report uses the term *data source* for original documents, media, and artefacts from which “network-compatible” (i.e., essentially relational) data can be obtained; data sources can be identified via *repositories* of various kinds, including library catalogues, archive portals, databases, and more. The report highlights the advantages of using data from authority files such as the GND, which are standardised, subject to quality control, and essential for disambiguating personal names. Moreover, authority files are freely accessible and provide information on the provenance of their data.

Concerning the repositories used, the authors report the general problem of missing or erroneous data, which leads to distorted answers to the questions, and which was not solved systematically in the project. In particular, the [GND](#) is focused on German library holdings and thus contains a disproportionately high amount of German-speaking persons. The data is biased towards people who have published at all, towards elite research, towards men (in particular among authors before the 1950s), and against certain disciplines, such as economy. The authors conclude that these restrictions significantly affect the answers to historical research questions. As possible remedies, they include further repositories, such as the German Union Catalogue of Serials (ZDB) [[Web21](#)], the German National Library (DNB) [[Web22](#)], the Kalliope Union Catalogue (KPE) [[Web23](#)], and, tentatively, other authority files (in particular, Wikidata and [VIAF](#)); we will get back to these in Section 4.3. However, the evaluation showed that even the sum of these repositories does not provide enough data for a differentiated temporal analysis; in particular, biographical data of the authors dominate, but those cover long time spans and do not provide sufficient evidence for or against dynamic relationships such as teacher-student relationships.

Concerning the selection of data, the report distinguishes between direct and indirect information on relationships. For example, if one is looking for support for the hypothesis of a social relationship between two persons *A* and *B*, then a family relationship between *A* and *B* explicit in the data supports the hypothesis directly, while matching biographical dates for *A* and *B* are an indirect indicator. In the latter case, the relationship explicit in the data (matching biographical dates) acts as *substitute information* (“Stellvertreterinformation”) for the relationship under consideration. The authors also address a special kind of indirect relationship, called intellectual, which mostly involves a third person or event, such as the co-citation of *A*’s and *B*’s texts by someone else.

A main constituent of the report is a catalogue of requirements to [HNA](#) and the research technology. This catalogue consists of 61 requirements that reflect the specific needs of researchers. The following of these requirements are relevant for our purposes (descriptions translated from the original German text and slightly rephrased and/or shortened):

- R003** Researcher wants all information on relationships contained in the metadata of the resources to be shown as derived social relationships in the data model.
- R004** Researcher wants all information on non-social relationships (see above) to be distinguished from direct social relationships.
- R005** Researcher wants to connect non-social and non-explicit social relationships with further conditions (e.g., overlapping lifespans).
- R006** Researcher wants to group people with comparable attributes (e.g., overlapping lifespans).
- R016** Researcher wants a list of people connected to a given person.
- R017** Researcher wants a list of corporate bodies connected to a given person.
- R018** Researcher wants a list of people connected to a given person via some corporate body.
- R029** Researcher wants, given a node, a description of the “potentially available attributes” and the provenance of the data.
- R030** [dito, but with edges instead of nodes]
- R043** [Researcher wants] a filter on attributes, e.g. biographical data, in order to restrict relationships to adulthood.
- R061** Researcher wants to know which kinds of relationships are contained in the dataset.

2.1.3 Graph Technologies for the Analysis of HNAs Using Heterogeneous Data Sources

In the paper of the same name, Menzel, M.-J. Bludau, et al. [2020] report on work in the research project “SoNAR (IDH), Interfaces to Data for Historical Social Network Analysis and Research” (SoNAR) [M. Bludau et al. 2020; Menzel, M.-J. Bludau, et al. 2020; Web19] project concerning the “creation and operation of research infrastructure for [HNA] based on heterogeneous data sources from cultural heritage institutions”. In particular, the authors present insights on modelling and transformation of heterogeneous data sources and on the design of visualisation for historical networks. In the remainder of this subsection, we summarise their insights on data selection and processing.

In the described approach, the authors integrate data from 6 heterogeneous repositories into a large uniform graph that is stored in a graph database and managed by the highly performant graph database management system Neo4j [Web24]. The original repositories comprise an authority file GND, federated library catalogues (DNB, ZDB, KPE), and portals offering electronic full texts of historical newspapers (ZEFYS [Web25], Exile Press [Web26]); altogether they contain some 30 million records. The data model is restricted to 9 entity types (e.g., Resource or PersonName) and 9 very generic relation types (e.g., RelationToPersonName, RelationToTopicTerm). There is also a distinction between explicit and implicit relations: the former are present in the data, and the latter have to be inferred automatically via certain “guidelines” and marked as such.

Since some of the above data sources include full texts, it is necessary to link names of people, corporate bodies, or geographical entities to entities in authority files. The use of the underlying technique, *named entity linking*, is discussed further by Menzel, Schnaitter, et al. [2021] in the context of SoNAR, and by Meiners [2022] independently of SoNAR.

Among the challenges associated specifically with processing data, Menzel, M.-J. Bludau, et al. [2020] name the sheer size of the combined graph, which causes performance issues, and the fact that the normalisation led to errors and inconsistencies. More generic challenges include the integration of domain knowledge and the adaptation of a more performant graph database engine (GraphDB), which requires extensive remodelling.

Some more comments on NEL and choice of data sources from [Menzel, Schnaitter, et al. 2021; Meiners 2022]? see Exposé

2.1.4 Insights Relevant for This Thesis

We now put the insights from the SoNAR project described in the previous two subsections into relation with the research goals pursued in this thesis. First, the analysis of historical networks is a broad concept, and SoNAR supports a setting that is much more general than just provenance research. Thus, the setting that we want to support in this thesis is a specific section of HNA and subsumed by the SoNAR setting.

The central role of relationships and of temporal aspects in historical research have to be reflected in the model and method that we are going to develop. Furthermore, the hermeneutic approach to historical research and the explorative approach to answering general research questions need to be supported by our method. In fact, the repeated exploration of data is exactly what we hope to support by enabling researchers to ask specific queries over a combination of data sources and use the answers to obtain “local” views of the whole network and thus inform the next steps in their research process.

Locality is an important feature that distinguishes our approach from the SoNAR setting: The exemplary research questions from the SoNAR report appear to be rather global in the sense that

they focus on wider areas in the network and that their analysis requires “global” techniques such as graph metrics and statistical methods. In contrast, we aim at answering more local queries that are part of a larger research question and which focus on a single item/actor and its neighbourhood in the network (e.g., owners of a book or students of a scholar; see the exemplary questions in Section 1.1). Answering such queries requires the retrieval of certain “candidates” from that neighbourhood which provide the necessary evidence and which then can be used to inform another iteration within the explorative approach. In the light of these considerations, we should not adopt the a-priori restriction to a fixed set of entity types or (generic) relation types; instead, it should be to use arbitrary concepts and relations from the data sources in queries and answers. Hence, we will focus on the described local perspective and leave the accommodation of global techniques for future work. In a similar spirit, we will not study visualisation; however, our envisaged retrieval system can be used *along with* visualisation tools and provide entry points to a large network.

The construction of a uniform data source (*graph*) that integrates all the available data from the heterogeneous repositories appears to be an integral part of SoNAR, according to Menzel et al.’s [2020] discussion. The obvious advantages include the interaction with a single database, the use of a single data model, autonomous hosting independently of the single repositories, and thus full control over the data. Furthermore, inconsistencies between the repositories can be resolved a-priori (at creation time), and the created database is independent on possible changes in the data models of the repositories. However, the static nature of the uniform database also has disadvantages: a data model has to be fixed upfront, and ultimately needs to be adapted when the models of the repositories change or new repositories are added; regular updates are necessary when the contents of the repositories changes; finally, as we have learnt above, a large graph database uses a large amount of memory and requires a highly performant database system.

As a result of this discussion, we want our approach to exploit the benefits of a dynamic solution. More precisely, while our abstract model will centre around a single data source (graph) that represents the combination of the distributed repositories, our method will abstain from constructing that graph explicitly and, instead, answer queries “in place” over the original repositories. This way, our approach will not depend on hosting capacities, and it will always have direct access to the current content of the repositories. On the downside, our approach will depend on external web services provided by the repositories, and it will be sensitive to changes in their data models. Our dynamic approach also requires that inconsistencies are resolved a posteriori, i.e., every time a query answer is retrieved. As a final advantage, the dynamic approach is flexible in the sense that it can be applied to a static integrated data source as well, thus benefiting from the advantages of the static approach.

A feature that we miss in the reports and publications of SoNAR is a rigorous definition of admissible queries that specifies exactly which queries are in the scope and which are not. We aim at providing such a rigorous definition for our setting which, at the same time, should be as general as possible. We consider this rigorous definition a main feature of our approach.

Data sources in the sense used in the SoNAR project report, i.e., original documents, media, and further artefacts are not part of our setting, as we do not consider the part of the research process that consults those data sources. Our work focuses on the level of what the SoNAR authors refer to as repositories, and we will relate, in Section 4.3, our choice of repositories to Menzel et al.’s [2020] selection. In this thesis, we continue to use the term “data source” for what SoNAR refers to as a repository, and we speak of *original data sources* to refer to original documents etc.

In order to allow researchers to refer to the original data sources, it becomes clear that *data provenance* plays a crucial role: the answer to a specific query in our approach needs to refer to the original data sources that provide evidence for the information contained in that answer. This

reference can point to original documents (SoNAR: data sources), and/or to data sets in repositories. For a more detailed discussion of data provenance, see Section 2.2.3.

The observed relevance of indirect or substitute information in the context of SoNAR seems important to keep in mind for our work as well: It has to be noted that not all information that one might be looking for is recorded in the data. For example, there is no record of who *read* a book, and ownership has to be used as a substitute for readership although evidence of ownership is only a necessary condition for readership and not a sufficient one—although one might argue that, in the case of, say, a scientist, ownership is very likely to indicate readership. Therefore, answers to questions about readership retrieved on the grounds of recorded ownership require further interpretation and investigation by the researcher who asked the question in the first place. The example from the SoNAR report that uses “having the same biographical dates” as a substitute for a social relationship is very similar, but in addition the fact that two persons have the same biographical dates is not a relationship that is given explicitly in the data but requires a certain amount of reasoning on several attributes of the entries for the two persons. This general observation has to be taken into account by our approach.

Finally, the insights on advantages and disadvantages of the repositories used in SoNAR will affect our discussion of data sources in Section 4.3.

2.1.5 Further Initiatives Providing Research Infrastructures

Menzel, M.-J. Bludau, et al. [2020] mention several projects that connect decentralised and heterogeneous bibliographic data sources and/or extract historical networks within the social sciences and humanities. We briefly summarise some of these projects.

DARIAH-DE [Web27] is an association that develops a digital research infrastructure for the arts and humanities following the example of large digital research infrastructures for the natural sciences. It comprises of 16 partner institutions from the academic sector and is part of the European network DARIAH-EU [Web28]. DARIAH-DE’s services for research include the development and hosting of services for data analysis and visualisation [Web29].

Culturegraph [Vorndran 2018; Web30] is a service offered by the DNB which aggregates bibliographic data from library networks in German-speaking countries, from the DNB itself, and from the British Library. According to Vorndran [2018], Culturegraph makes over 160 million data sets available “for data analyses, evaluation of connections and statistical analyses”. This is achieved, among other things, by enrichment via external data sources such as GND and ORCID [Web31], thus allowing for the identification and disambiguation of persons.

Menzel, M.-J. Bludau, et al. [2020] further note “an increase in joint research projects that are focused on the extraction of historical networks within the social sciences and the humanities.” These projects include *Six Degrees of Francis Beacon* [Warren et al. 2016; Web32], *histoGraph* [Novak et al. 2014; Web33], *Issues with Europe* [Web34], *APIS* [Gruber and Wandl-Vogt 2017; Web35], and *Gesellschaftliche Wissensproduktion in der Aufklärung* [Purschwitz 2018]. Judging from Menzel et al.’s [2020] summary, these projects put an emphasis on statistical methods, visualisation, and/or exploration.

2.2 Linked Data and Data Integration

2.2.1 Basic Concepts

Data Integration refers to the problem of making a set of autonomous and heterogeneous data sources uniformly accessible [Doan, Halevy, and Ives 2012, p.6]. The current landscape of data sources includes highly structured data represented and accessed using classical database techniques, as well as more loosely structured data accessible via (Semantic) Web techniques. Given this vast and diverse landscape, data integration is a complex problem, and various techniques have been developed for this purpose. The textbook by Doan, Halevy, and Ives [2012] provides a comprehensive introduction to data integration.

In the remainder of this section, we restrict our focus on those aspects of data integration that are most relevant in the context of bibliographic data sources on the Web.

Linked data [Web36; Domingue, Fensel, and J. A. Hendler 2011] refers to data published on the World Wide Web that is structured and connected with other data. Based on a small and uniform set of simple technologies the linked data paradigm provides applications with access to “a global, unbounded dataspace” [Domingue, Fensel, and J. A. Hendler 2011]. In the context of the Semantic Web [Berners-Lee, J. Hendler, and Lassila 2001; Marshall and Shipman 2003], linked data is one of several developments aimed at making data on the Web more machine-understandable.

The main technologies underlying linked data are the following.

URIs A Uniform Resource Identifier (URI) [Berners-Lee, Fielding, and Masinter 2005] is a string that is assigned to a physical or logical resource in order to identify that resource uniquely. URLs (Uniform Resource Locators) are a special case of URIs which additionally ensure that a resource can be located in a network (e.g., the internet) and that information can be retrieved from it. URIs are an integral part of RDF and the Web Ontology Language OWL (see below).

HTTP The Hypertext Transfer Protocol (HTTP) [Web37] is the fundamental protocol for data transfer and communication underlying the World Wide Web (WWW).

RDF The Resource Description Framework (RDF) [Web38] “is a framework for expressing information about resources. Resources can be anything, including documents, people, physical objects, and abstract concepts” [Web38]. Its main ingredients are *resources* and binary *properties* for linking resources. RDF statements are *triples* of the format subject–predicate–object (S–P–O) consisting of a resource, a property, and a resource (or literal). Resources and properties are described using **URIs**. RDF thus provides a simple and flexible data model that is particularly useful for the publication and linking of data on the Web; it has been a standard of the World Wide Web Consortium (W3C) since 2004 [Web39].

In the context of the Semantic Web, **RDF** and associated technologies are important tools for the definition and use of ontologies. For a definition of the term “ontology” in this setting, we cite Horrocks and Patel-Schneider [2011]: “A major feature of the Semantic Web is the ability to provide definitions for objects and types of objects that are accessible and manipulable from within the Semantic Web. In Computer Science, a collection of these sorts of definitions about a particular domain is called an ontology, although philosophers may (and probably will) have a different understanding of what constitutes an ontology.” More importantly for this thesis, ontologies can be used to represent knowledge about a specific domain or about the world in general, in an unambiguous and machine-readable way, using a formal language. What is more, *reasoning* mechanisms can be employed to derive implicit knowledge, i.e., knowledge that logically follows from the explicitly represented knowledge.

The following technologies associated with **RDF** and ontologies are of interest for the purposes of this thesis:

RDFS RDF Schema (RDFS) [Web40] is a specific **RDF** vocabulary which can be used for modelling and thus serves as a very basic ontology language. It provides terms for modelling, among others, classes, properties, and domain and range restrictions.

OWL **OWL** is the ontology language recommended by the W3C. It is based on knowledge representation languages from the description logic (DL) family [Baader et al. 2017]. DLs have a well-defined syntax and model-theoretic semantics, which makes them suitable as a formal ontology language in the sense mentioned above. The members of the DL family vary regarding their expressive power and, closely related, regarding the computational complexity of their basic reasoning problems. OWL 2, which is the current version of OWL, is based on an expressive DL where reasoning is still decidable and reasoning algorithms have been implemented in various inference machines that perform well on a wide range of real-world ontologies. Besides these knowledge representation languages and support for reasoning, OWL includes infrastructure for interaction with ontologies and interoperability within the Web, such as Internationalized Resource Identifiers (IRIs), **XML** schema datatypes, import mechanisms, and many more.

SPARQL The Simple Protocol And RDF Query Language (SPARQL) [Web41] is a W3C-recommended query language for the Semantic Web. More precisely, it is a language for querying **RDF** graphs via pattern matching, whose syntax is based on the ISO standard SQL [Web42] for managing and querying relational databases. However, SPARQL is “far more powerful than SQL, since it is designed for the *open, decentralized, and fluid* Web” [Della Valle and Ceri 2011]. SPARQL also provides a communication protocol for the interaction between clients and *endpoints*; answers are returned in **RDF** or **XML**. Additionally, SPARQL exploits inference mechanisms, i.e., answers may contain facts that are not explicitly stated in the original **RDF** graph but are obtained involving certain sets of inference rules, based on **OWL** or other standards.

Linked open data (LOD) is linked data that is also *open data*, i.e., accessible and shareable by everyone [Web43].

Also mention ontology-based data integration? See https://en.wikipedia.org/wiki/Ontology-based_data_integration and references therein

2.2.2 Linked Data in the Library Domain

The state of LOD in the library domain in 2013 is summarised by Pohl and Danowski [2013] in their introduction to the edited volume “(Open) Linked Data in Libraries” [Danowski and Pohl 2013] as follows (translated from German and rephrased). Libraries and related organisations began to experiment with linked data as early as 2008. This was followed by linked-data initiatives by important players such as OCLC, the Library of Congress (LoC) [Web44], and the **DNB**; in particular, the **LoC** started the initiative “Bibliographic Framework for the Digital Age” in 2011, declaring the renunciation of the library-specific standards **MARC 21** and **Z39.50**, and announcing the development of an infrastructure based on **RDF** as a basic data model. The advantages of LOD for the library domain include retrievability of, e.g., catalogue data by search engines, mechanisms for permanently linking, e.g., hit lists and hit views, the use of a much more flexible data model, interoperability ensured by the use of open web standards, and reusability via open licences. A prominent example of this development is the provision of authority data as LOD, which can be shared and reused on the Web: e.g., authority data on person names from the **GND** has been used by Wikipedia since 2005 in order to provide links to further reading in Wikipedia articles, see also [Hengel and Pfeifer 2005]. In addition to the **GND**, further providers have made authority data

available as LOD, among them [LoC](#) and [VIAF](#).

The state of L(O)D in libraries in 2013 was extensively surveyed in the special issue “Linked Data, Semantic Web and Libraries” of the *Journal of Library Metadata* [Bair 2013a]. In her preface to this issue, Bair [2013b] refers to the challenges to linked data implementations in general and in the library domain that were reported in previous work [Bizer, Heath, and Berners-Lee 2009; Byrne and Goddard 2010; [Web45](#); Alemu et al. 2012], and she emphasises that these challenges remain. The special issue contains a survey on the perception of linked data in libraries, as well as “seven case studies of the experimental efforts of LAM institutions to use linked data to increase access to their collections and user services, plus four others that aim to increase awareness of and educate on key topics” [Bair 2013b, p.76]. The challenges highlighted in these articles fall into five areas, among them data and schema mapping and interoperability, and data quality and trust.

Further work on L(O)D in the library domain is reported in the following publications.

Burrows et al. [2021] describe an LOD model that links three large manuscript databases in the *Mapping Manuscript Migrations (MMM)* project. MMM aims at “providing large-scale analysis and visualizations of the history and provenance of medieval and Renaissance manuscripts” [Burrows et al. 2021, p.3].

Lígia Triques, Ventura Amorim Gonçalves, and Albuquerque [2022] give an overview of the data collection and integration technology in the Digital Public Library of America (DPLA), with a focus on interoperability and challenges for integrated data access.

N. Freire et al. [2019] study metadata aggregation in the context of *Europeana* [Isaac, Clayphan, and Haslhofer 2012; Petras and Stiller 2017], a web portal of the European Union that provides unified access to more than 56 million digital objects from the collections of over 4000 cultural heritage institutions in Europe [[Web46](#)]. N. Freire et al.’s case study involves two Dutch institutions as data provider and aggregator, and aims at improved discoverability of the data. The main challenge was the transition from traditional data models to flexible semantic data models. The article presents the results of a requirement analysis, the workflow that was developed, and its implementation.

Ullah et al. [2018] review the current state of L(O)D in cataloguing in the context of the trending transfer of bibliographic metadata towards L(O)D by major libraries. Their review provides an extensive survey of the recent literature. The main findings include the observations that L(O)D is becoming “the mainstream trend in library cataloging especially in the major libraries and research projects of the world” [Ullah et al. 2018, p.47] and that bibliographic metadata is becoming increasingly meaningful and reusable of through the emergence of Linked Open Vocabularies.

Hauser [2014] gives an overview of the linked data service of the [DNB](#), which started in 2010 with the publication of [GND](#) authority data as linked data. Since 2012, the [DNB](#) has been using and maintaining its own [GND](#) ontology [[Web47](#)] in order to bridge the gap from the original MARC-based data model to a flexible, open data model. In particular, the [GND](#) ontology provides a vocabulary for describing entities such as persons, corporate bodies, and places, thus allowing to disambiguate and link these entities.

2.2.3 Metadata Provenance

In the data integration scenario, the origins of (meta)data are important: When querying the combination of several data sources, the retrieved answer should contain, for every fact, information that identifies the original data source or repository from which that fact was taken. This information is useful as a “justification” for the retrieved answer or as a pointer for further research; it is especially desirable when the consulted data sources contain conflicting information. In order to

delineate the origin of data about a book (or other item) from the provenance of that book itself, we adopt the notion of *metadata provenance* [Eckert 2012].

Metadata provenance on the Web has been studied from the beginnings of linked data [see, e.g. Hartig 2009; Moreau, Groth, et al. 2008; Moreau, J. Freire, et al. 2008]. In the context of linked data from cultural heritage institutions and especially the Europeana portal, metadata provenance is studied by Eckert [2013a; 2013b; 2012].

elaborate, depending on whether provenance plays a larger role in subsequent chapters

Data provenance has also received attention in the classical database domain [see, e.g., Doan, Halevy, and Ives 2012, §14].

2.2.4 Some Recent Developments

omit this subsection?

Boumechaal and Boufaïda [2023] developed a framework for transforming queries formulated in natural language (NL) to SPARQL queries in order to “query linked and heterogeneous semantic data on the web” [Boumechaal and Boufaïda 2023, p.1]. In contrast to most of the previous work, their approach focuses on complex queries involving negation, numbers, superlatives, and comparative adjectives. The aspect of translation from NL to SPARQL (or any other formal query language) is important in the context of historical research, where the use of a formal query language would be a barrier for most users. This aspect is, however, not in the scope of this thesis and should be considered in future work.

2.3 Provenance Indexing

In an earlier master’s thesis, Hakelberg [2016] studies the state of provenance indexing with authority data. We summarise his conclusions in the remainder of this paragraph. The GND provides authority data that is very suitable for recording provenances across institutions; this is ensured by the structured data model and the open and collaborative nature of the GND. The GND has thus become a central instrument for recording provenance information. Provenance indexing is laborious, and practices vary greatly between German libraries (see Section 1.1). Provenance records are useful only if they are retrievable across library networks. For this purpose, indexed data needs to be homogeneous and the usability of catalogues of library networks needs to be ensured. Among other things, an overview of the normed vocabulary used for provenance indexing should be provided to users as a search entry. Hakelberg also recommends the further development of the Thesaurus of Provenance Terms (T-PRO) [Web48], which is a controlled vocabulary for the specification of provenances via single descriptors or chains thereof. In particular, the T-PRO descriptors need to be assigned URIs in order to make them reusable as linked open data.

We draw the following insights from these observations. GND as well as library catalogues should be central data sources within our approach. Obstacles may be caused by the heterogeneous situation of indexing in library catalogues as well as the fact that T-PRO is not ready for LOD. The need for the presentation of the used vocabulary as a search entry should be taken into account.

3 Example Queries and Answers: A Case Study

As a first step towards delineating the type of queries that should be covered by our approach, we examine the queries used as motivating examples in Section 1.1 more closely. Those queries will serve as a point of reference for the analysis of the available data sources in Chapter 4, and they will be generalised by the abstract framework developed in Chapter 5. With that framework, we will provide a rigorous definition of the admissible queries that specifies their structure but does not impose any restriction on their content.

3.1 Example Queries

As already noted in Section 1.1, the queries introduced there are in fact *query patterns*, and we will make this distinction in the remainder of this section. We repeat the query patterns here in order to discuss them in more depth.

- Q1** Who read work *X*, in which manifestation and in which year?
- Q2** Which exemplars of work *X* were passed from one of its owners to a student?
- Q3** What are the relationships between the recipients of manifestation *Y* of work *X*?
- Q4** Which exemplars from a collection *X* were passed on by its owner to a family member?
- Q5** Which exemplars from the holdings of library *X* were acquired from bookseller *Y* between 1933 and 1945?
- Q6** Who participated in the sale of collection *X*?
- Q7** Via which paths did exemplars from collection *X* enter library *Y*?
- Q8** Which libraries own the exemplars once owned by person *X*?
- Q9** Where did person *X* acquire exemplars and did they know the previous owners?

These patterns can be instantiated by substituting each variable with concrete object, which yields a specific *query*. For example, if we substitute the variable *X* in **Q2** with the seminal work *De revolutionibus orbium coelestium* (short: *De revolutionibus*; English translation: *On the Revolutions of the Heavenly Spheres*) [Copernicus 1543] by the astronomer Nicolaus Copernicus (1473–1543), we obtain the following query.

- Q2'** Which exemplars of *De revolutionibus* were owned by some scientist who passed them on to a student?

It is important to note that these exemplary query patterns are not meant to be representative for the range of queries that provenance researchers are interested in asking. Finding a representative choice would require a systematic analysis of queries relevant to or useful for researchers. Such an analysis would need to comprise an extensive interview study based on very generic questions of a predominantly open-ended nature, requiring a labour-intensive evaluation which could easily fill a separate thesis. As indicated above, the general framework that we will develop is informed by the available data sources and designed to cover a wide range of possible queries. Hence, it is reasonable to assume that tools developed on its basis will be helpful for provenance researchers.

In subsequent work, when our method will hopefully have been implemented in a prototype retrieval system, the extent to which researchers' needs are served can be determined by means of a more focused user study with more specific questions, which in turn can inform possible extensions of the framework.

3.2 Manual Answer Retrieval

In order to demonstrate how a researcher could proceed (manually) when answering a query, we take Pattern **Q2** and fix work *W* to be the seminal work *De revolutionibus orbium coelestium* (short: *De revolutionibus*; English translation: *On the Revolutions of the Heavenly Spheres*) [Copernicus 1543] by the astronomer Nicolaus Copernicus (1473–1543). We thus consider the following concrete query.

Q2' Which exemplars of *De revolutionibus* were owned by some scientist who passed them on to a student?

When answering **Q2'**, an obvious way to proceed is the following: first, our researcher finds exemplars of *De revolutionibus* in online catalogues of libraries and library networks. For each such exemplar, they then inspect the provenance entries that name owners who were people (not corporate bodies). Finally, our researcher will have to find those names in databases such as authority files or Wikidata and, for each entry, explore the specified profession ("scientist") and relationships to other people ("student").

For example, the online catalogue (OPAC) of the Gotha Research Library of the University of Erfurt (Forschungsbibliothek Gotha) lists two printed exemplars of *De revolutionibus* [Web49]. One of those bears the signature Druck 4° 00466, and its provenance entries name the following previous owners [Web50]:

- Hieronymus Tilesius (1529–1566): autograph and date 1551
- NN: note, date 1553, name scraped out
- Johann Hommel (1518–1562), autograph
- Valentin Thau (1531–1575), note (greek proverb, possibly not denoting ownership)
- Ernest II, Duke of Saxe-Gotha-Altenburg (1745–1804): stamp/seal, initial
- Ducal Library, Gotha (a predecessor organisation of Gotha Research Library): stamp marking a duplicate
- Ernestine Gymnasium, Gotha: stamp
- Landesbibliothek Gotha: stamp

Our researcher can immediately decide that they can ignore the second entry (no name given) and the last three entries (corporate bodies). For the remaining four entries, our researcher follows the links given in the OPAC to the [GND](#). On inspection of these [GND](#) entries, it turns out that Ernest II was a regent and very probably not a scientist, and that the other three people—Tilesius, Hommel, and Thau—had professions such as theologian, mathematician, and astronomer, which qualifies them as scientists. Furthermore, Hommel's entry contains a reference to Thau via the relation "has student" (and Thau's entry contains the inverse reference to Hommel). From this reference, our researcher can conclude that two scientists in the teacher-student relation have both possessed the exemplar.

Unfortunately, the data available does not imply that Hommel passed the exemplar (directly) to

Thau; furthermore, the provenance entry for Thau raises doubts as to whether Thau was really an owner of the exemplar. Therefore, the retrieved data can only be regarded as a *candidate answer* which entails the hypothesis that the found exemplar was passed from Hommel to Thau. Our researcher can now engage in further research in order to verify that hypothesis.

Discuss instances of **Q1**, **Q3–Q9**, and their differences (extend remarks from Section 1.1)?

3.3 The Quality of Query Answers

Our simple example already illustrates that the quality of query answers strongly depends on the quality of the underlying data, regardless of whether answers are obtained manually or automatically. In particular, missing or spurious data will lead to missing or spurious answers. Before we begin to deal with automated query answering, we need to analyse the sources of missing or spurious data and their effects on the answers obtained.

For the sake of a principled discussion, we call the set of answers to a query obtained by some manual or automatic procedure an *answer set*, and we call the (set of) answers that the query has in reality the *true answers* and the *true answer set*. In the above example, the answer set obtained by the described manual process consists of the single answer “Druck 4° 00466”, if we assume that our researcher interprets the data retrieved generously (i.e., as an indication that Thau might have been an owner and might have received the book directly from Hommel) and draws additional conclusions (i.e., that every mathematician or astronomer is a scientist). The true answer set is not known and will most probably never be known; it is a term of rather philosophical nature.

Ideally, both answer sets should coincide. Incomplete data may cause the answer set to exclude some true answers, and spurious data may cause it to contain some answers that are not true. We call the respective answers *missing* and *spurious*. On the basis of our example, we can identify four distinct causes for data being incomplete or spurious, as discussed in the following. We use the word *term* as an umbrella term for concepts (such as “scientist”) and relations (such as “is student of”).

1. In the example, there is no information available in the [GND](#) on whether any of the persons involved is indeed a *scientist*. Instead, [GND](#) provides information on more specific professions, e.g., for Thau and Hommel (theologian, mathematician, and astronomer). The information that they were scientists has to be *derived* based on *general knowledge of the world*. Indeed, a search in the [GND](#) catalogue reveals that [GND](#) does have an entry for the subject “Wissenschaftler” (scientist) [[Web51](#)] which, however, is used only sporadically: its 46 subordinate terms do not include, for example, “Mathematiker”, “Astronom”, or “Theologe” (mathematician, astronomer, theologian) [[Web52](#)].

The same effect can occur with relations instead of concepts: if **Q2** were to ask for family members instead of collaborators or students but the data only supported more specific relations such as “has sister” or “has father”, then relationships using “has family member” would need to be derived. Unfortunately, the cataloguing rules for [GND](#) authority datasets do not require that relationships are recorded exhaustively nor using a normed vocabulary; see Section 4.3 for more details.

In summary, terms can be missing in the data sources because they have not been recorded or because they are implicit in more specific terms. If the query uses such missing terms, then the answer set is always empty unless the person or machine determining the answers makes derivations based on their knowledge of the world.

Clearly, it would not be advisable to attempt adding all implicit knowledge to the data because that would massively inflate the data, as most terms have several superordinate concepts or relations and, furthermore, implicit knowledge is not restricted to taxonomic knowledge.

Explain this further? A case for ontologies!

2. In the example, there is no information available on whether any of the owners of the exemplar *passed it on* to another owner. Similarly, attempts to answer instances of Query Pattern **Q1** will have to deal with the problem that there is no information available as to who *read* books.

More generally speaking, terms can be missing in the data sources because they are not recorded at all, as a consequence of either a general lack of evidence or a general design decision for the data source. The underlying reasons can be manifold: for example, relationships such as who actually *read* a book are very hard to confirm, or terms may not be part of the fixed vocabulary for a data field in a source. If the query uses such terms, then the answer set is always empty, as in Case 1.

3. In the example, it is possible that single owners of the exemplar or single relationships between owners have not been recorded because no evidence has been found yet. More generally, concept memberships or relationships can be missing sporadically, which may lead to missing answers.
4. In the example, if the provenance entry on Hommel is incorrect, then the answer obtained on the grounds of that entry is incorrect. More generally, spurious concept memberships or relationships may lead to spurious answers.

These cases reveal a striking qualitative difference concerning the effects of data incompleteness or spuriousness on the answer set: Cases 3 and 4 have effects on single answers only, i.e., some answers are missing or spurious. However, Cases 1 and 2 generally make *all* answers missing unless further provisions are made. Such “provisions” might include *reasoning* (e.g., deriving the implicit knowledge that Hommel was a scientist from the explicitly recorded fact that he was a mathematician) and *hypothesising* (e.g., assuming that Thau was an owner and received the book directly from Hommel). Reasoning can be supported using semantic technologies such as (domain-specific or top-level) ontologies and related infrastructure. A possible way to support hypothesising is to provide for the use of *substitute information*, which has been addressed in the context of the [SoNAR](#) project; see Section 2.1. That is, users could be allowed to declare terms occurring in the data that may act as substitutes for certain terms used in their queries.

In the light of these observations, it is appropriate to consider query answers as *candidates* that necessitate (and inspire) further research. For this purpose, spurious answers (in manageable numbers) are less harmful than missing answers. Consequently, it is important to find ways to avoid missing answers without generating too many spurious answers. In other words, it is desirable to obtain an answer set that is a slight *overapproximation* of the true answer set.

3.4 From Manual to Automated Query Answering

The manual process that we have described in Section 3.2 is cumbersome, laborious, and prone to errors and omissions for several reasons: The search in library catalogues for exemplars of works and their provenances requires expert skills. Catalogues with potential matches need to be selected manually, and each catalogue needs to be queried individually, using its own search functionality and syntax. The traversal through all retrieved exemplars and the pursuit of each potential

give examples of commonalities (OPAC) and differences (discov-

relevant provenance entry per exemplar multiplies the amount of manual work necessary. Finally, it is not clear what an effective and efficient way to “explore” relationships would be: while it is easy to find direct relationships such as “is student of” in the view for a person’s entry in databases such as [GND](#) or Wikidata, there are relationships that cannot be discovered easily by hand, e.g., “ P_1 and P_2 are students of the same scholar”.

These considerations suggest that query answering will strongly benefit from automated support, which can help reduce the amount of manual work, integrate heterogeneous data sources, incorporate background knowledge (e.g., every mathematician is a scientist), and discover relationships between entities that are not necessarily direct. We envisage a retrieval system that enables researchers to formulate queries and which computes answers consulting data sources selected by the user. The vocabulary used for formulating queries should be based on the vocabulary present in the data sources, but it should also include terms such as “scientist” or “read”, which are not recorded in the data, as discussed in Section 3.3. The retrieval system thus needs to implement ways to map those terms with the available vocabulary, as well as techniques for integrating data from heterogeneous sources. It is our general vision that the retrieval system will serve as an instrument for prospectively finding interesting candidates that inspire further research. In the remainder of this thesis, we want to lay the foundations and develop a precise method that can serve as a basis for a future implementation of a retrieval system.

4 Analysis of Available Data Sources and Techniques

In order to build a system for retrieving provenance relationships, data sources need to be selected. It is the aim of this chapter to provide information for this selection can be based later on. Furthermore, the abstract model of data sources, provenance queries, and answers that we will develop in the following chapter needs to incorporate the characteristics of the data models used in the appropriate data sources. For this purpose, we review data sources that contain the relevant metadata (Section 4.3), as well as bibliographic standards, data formats, and communication interfaces relevant for these sources (Sections 4.1 and 4.2). Given the observation from the previous chapters that answering provenance queries is likely to involve several data sources, we pay special attention to data transfer and integration techniques, including linked data. At the end of this chapter (Section 4.4), we draw conclusions from the review that inform our model and method in the following chapters.

Before we start the review, we need to define the central term *metadata*. For this purpose, we follow Hider and Harvey’s [2008] definition and reflection: The term *metadata* is used in multiple ways. The most general use is according to its literal meaning, “data about data”. This definition is not restricted to the library domain but includes bibliographic records for documents (which contain the primary *data*). A more specific variant is the use of “metadata” to refer to structured data describing digital objects, again independently of the library domain. Hider and Harvey also note that these two meanings can no longer be clearly distinguished in view of the progressing digital transformation. In this thesis, we adopt their decision to use the term *metadata* in the most general sense, “applying to all information resources” [Hider and Harvey 2008, p.13].

4.1 Standards for the Description of Bibliographic Resources

The following two standards are widely adhered to in bibliographic data sources. This discussion is a brief summary of Wiesenmüller and Horny’s introduction [2015, §§2, 3], unless indicated otherwise.

FRBR The Functional Requirements for Bibliographic Records (FRBR) [IFLA SG FRBR 2009] constitute a conceptual entity-relationship model developed by the International Federation of Library Associations and Institutions (IFLA) [Web53] in 1998 (and updated in 2009) in order to support users in finding, identifying, selecting, and accessing bibliographic items [Wiesenmüller and Horny 2015, p.17]. FRBR’s entities represent things that need to be represented in the data. These entities fall into three groups, of which we introduce the first two. Entities of Group 1 are Work, Expression, Manifestation, Item, and their common superordinate concept Endeavour; entities of Group 2 are Person, CorporateBody, and their superordinate concept ResponsibleEntity. Furthermore, FRBR contains relations that link Group-1 entities with each other and with Group-2 entities, respectively. These entities and the basic relations are shown in Figure 4.1.

Beside FRBR, the IFLA developed a conceptual entity-relationship model for authority data, the Functional Requirements for Authority Data (FRAD) [IFLA WG FRANAR 2008], as well as a continuation of FRBR, the Functional Requirements for Subject Authority Data (FRSAD) [IFLA WG FRSAR 2010].

FRBR serves as a basic principle of RDA, which we discuss next.

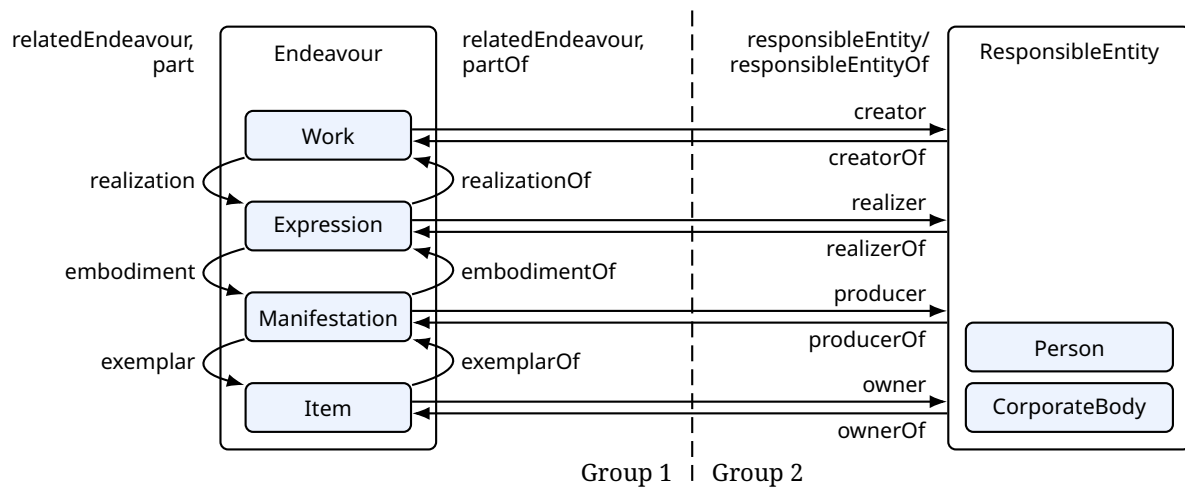


Figure 4.1: FRBR entities and basic relations of Groups 1 and 2 [following [Web54](#); [Web55](#)]

RDA Resource Description and Access (RDA) [[Web56](#)] is a standard for the cataloguing of publications in cultural heritage institutions, and particularly in libraries. RDA is implemented and used in the library systems of several countries, including the US, the UK, Canada, Australia, and the German-speaking countries [[Web57](#)].

RDA and the application guidelines for the German-speaking area stipulate a special practice for recording FRBR Group-1 entities and relating them to each other: Bibliographic records (“Titel-datensätze”) have a bibliographic level for describing a manifestation and an exemplar level for describing the related exemplars. The description of works and expressions is considered part of the description of a manifestation and thus recorded on the bibliographic level. However, it is possible to create and link an authority record for a work or expression containing the relevant description. In that case, the source of the information can be recorded as well [cf. Wiesenmüller and Horny 2015, §5.1].

For our approach, this observation means that catalogues of German libraries do not use a uniform way of linking, e.g., a work X with its manifestations. In particular, if there is no authority record for X , then X is represented in the data of its manifestations only implicitly (and possibly ambiguously).

4.2 Data Formats and Communication Protocols

In order to identify data sources in the following section, we first need to describe relevant standards and formats for the description and communication of (bibliographic) data. The following discussion is based on Hider and Harvey’s overview [2008, §10], unless indicated otherwise

4.2.1 Markup Languages

We briefly describe markup languages, which are repeatedly used in the technologies described in the following subsections. Markup languages are machine-readable languages used for the structuring and formatting of files.

HTML The most prominent example is the Hypertext Markup Language (HTML) [[Web58](#)], the core language of the [WWW](#). The markup features of HTML are also used to identify and describe certain elements of digital objects. In particular, HTML provides for meta tags, i.e., labels for

metadata. One of the restrictions of HTML is its fixed set of labels, whose meanings cannot be changed. Further standards were developed to overcome this restriction, such as SGML and XML.

XML The Extensible Markup Language (XML) [Web59] is both a file format and markup language for the storage and transmission of arbitrary data; it allows for a hierarchical structuring of data in a text file format that is readable by both humans and machines [Web60]. XML is the basis for many modern web developments. Most metadata schemas have standard expressions in XML, as we will learn below.

4.2.2 Data Communication Protocols

Z39.50 is a protocol for data communication between bibliographic databases. It is an **ISO** standard as well as an **ANSI/NISO** standard. Z39.50 is a set of rules developed specifically for translating search and retrieval commands between databases such as OPACs. Using Z39.50, it is possible to issue commands specific to a local database and obtain results from a remote database that uses a possibly different command set. This protocol can be implemented over the internet and similar networks. Z39.50 is widely used in the library domain for providing access to catalogues, for example, by the **LoC**. Still, it is not fully implemented by all library systems, and not all libraries use the most recent version. It is also not very widely spread outside the library community.

SRU Search/Retrieval via URL (SRU) [Web61] is a successor of **Z39.50** and an **LoC** standard. It has the same function, which is the standardisation of search and retrieval across online databases. SRU is based on **HTTP**, **XML**, and the Contextual Query Language (CQL) [Web62], a standard syntax for representing queries to information retrieval systems. Thanks to this technology, SRU is more dynamic than **Z39.50** and less restricted to the library domain.

OAI-PMH The Open Archives Initiative (OAI) [Web63] is a project devoted to interoperability standards in information agencies in general, including libraries. The OAI Protocol for Metadata Harvesting (OAI-PMH) [Web64] is a standard developed by the OAI and is widely used in the library domain, albeit not as widely as **Z39.50**.

4.2.3 Data Description Schemata

MARC Machine-Readable Cataloguing (MARC) [Web65] “is the main data communication standard in use in libraries today” [Hider and Harvey 2008, p.198]. It is a family of formats for the exchange of bibliographic and further data between library systems, which is being extensively used by libraries around the world. MARC was developed by the **LoC** in 1969 and has been updated several times. The MARC family comprises more than 20 dialects that have been developed as an official standard in several countries. Those are very similar to each other and provide very detailed record structures for cataloguing of bibliographic data. They all adhere to an international **ISO** standard. The development of MARC pursued the goals of labour and cost reduction, and of standardising the cataloguing process as well as data communication and transfer. MARC “allows flexibility in library information systems” [Hider and Harvey 2008, p.201] by providing an easy way of data exchange and allowing for better resource discovery (among others, in comprehensive union catalogues).

MARC 21, was developed in 1999 and is “the major version of MARC in use internationally” [Hider and Harvey 2008, p.205]. It consists of five formats for specific kinds of data. Two of these are for bibliographic data and authorities data, respectively.

MARC has been criticised for its inflexible structure, the fragmentation into several dialects, and the restricted compatibility with current computer technologies, among other things. Hider and

Harvey [2008, p.212] give more details and further references. There have been several attempts to redeem some of these disadvantages, among them the development of new standards such as MODS (see below) and of XML schemas based on MARC 21 such as MARCXML and Turbomarc.

The problem of standard proliferation was not restricted to the variety of MARC dialects but also occurred within non-MARC formats. Attempts to solve this problem have been made via translation and unification. One tool for unification is Dublin Core (DC); see below.

MAB The “Maschinelles Austauschformat für Bibliotheken” (MAB) [Web66], translating as “machine data exchange format for libraries”, is a legacy bibliographic data format that has been developed solely for the purpose of data exchange by the DNB. Although Hider and Harvey [2008, p.204] classify it as a dialect of MARC, it is in fact conceptually different from MARC, assigning exactly one bibliographic element to each data field and allowing a more flexible ordering of elements [Web67]. In 2013, the DNB completely abandoned the delivery of its bibliographic and authority data in MAB in favour of MARC 21 [Web66].

MODS The Metadata Object Description Schema (MODS) [Web68] is a standard developed by the LoC in order to overcome the problems with MARC mentioned above. MODS is based on XML and a subset of MARC fields; it can thus be used alongside MARC and as a “switching format between MARC and non-MARC schemas” [Hider and Harvey 2008, p.219]. The LoC also maintains an analogous standard for authority data, the Metadata Authority Description Schema (MADS) [Web69].

DC The Dublin Core Metadata Element Set, short: Dublin Core (DC) [Web70] is an international standard consisting of 15 essential metadata terms for describing digital or physical resources. It was formulated by the Dublin Core Metadata Initiative (DCMI) [Web71], a project of the US-American non-profit organisation ASIS&T, for the purpose of locating information resources on the WWW. The 15 core elements are tailored towards the “primary metadata needs” across domains [Hider and Harvey 2008, p.215], thus making it a flexible data model. Applications of DC include resource description, the combination of metadata vocabularies from different standards, and the provision of interoperability in the linked data context. DC is extended by the *DCMI Metadata Terms*.

RDF We have already introduced RDF in Section 2.2. In the context of the following description of data sources, RDF is a very flexible data schema and the core technology for providing metadata as linked data.

BIBFRAME The Bibliographic Framework (BIBFRAME) [Web72] is a data model for bibliographic description, which was designed to replace the MARC standards and use linked data principles to improve access to bibliographic metadata within and outside the library domain [Web73]. The most recent version 2.0 was released by the LoC in 2016 [Web72]. So far, only a handful of libraries are using BIBFRAME, most of them in test mode [Web73].

4.3 Data Sources

In this section, we provide information on data sources that provide metadata on objects and individuals relevant for provenance research, such as works, expressions, manifestations, exemplars, persons, ownership, and social relationships. We first give an overview of eligible data sources (Subsection 4.3.1) and then select a few particularly relevant ones, for which we provide a detailed analysis (Subsection 4.3.2). The choice of the latter was motivated by the insights from the SoNAR project (see Section 2.1) and Hakelberg’s [2016, §4] overview of the state of provenance indexing with authority data in the German-speaking area. This analysis will have to be extended whenever

the model and method that we are going to develop in the following chapters is implemented in a retrieval system in future work.

4.3.1 Collection of Data Sources

We have identified four categories of relevant data sources: library catalogues, special databases for provenance research, authority files, and knowledge bases. We next describe, for each of these categories, the relevant information that is contained in the respective data sources, and we list examples.

Library Catalogues Catalogues contain the information relevant for provenance research, such as

- bibliographic entities (works, manifestations, expressions, exemplars);
- relations between these entities (e.g., `isManifestationOf`);
- attributes for these entities (e.g., year of publication);
- people and corporate bodies and relations such as authorship;
- provenance entries (including, e.g., owners, the ownership relation, and further attributes such as the year of ownership);
- (implicitly) the current ownership;
- (ideally) references to entries in authority files for all these types of entities.

Examples of relevant library catalogues include

- catalogues of national libraries, which often aim at collecting literature exhaustively and which are most likely to make metadata available in interoperable formats, in particular: **the catalogues of DNB, LoC, and the British Library** [Web22; Web74; Web75];
- meta-catalogues that enable a federated search in several catalogues, in particular: **WorldCat** and the **Karlsruhe Virtual Catalogue (KVK)** [Web7; Web76];
- catalogues of (German) library networks, in particular: **K10plus, B3KAT**, the **union catalogues of hbz and hebis** [Web77; Web78; Web79; Web80];
- catalogues of single libraries, if not covered by any of the above;
- special catalogues, e.g., those used for historical research in the **SoNAR** project, in particular: **ZDB, KPE, ZEFYS, ExilePress** [Web21; Web23; Web25; Web26].

Authority Files Authority files contain information relevant for provenance research, such as

- entities such as persons, corporate bodies, places, works;
- relations such as social relations, family relations, professional relations;
- attributes for the entities such as their profession;

They are usually subject to a strict quality control. Examples include the **GND** [Web2] in the German-speaking area (which has been used extensively in **SoNAR**), the North-American **LCNAF** [Web6], and the international projects **ISNI** and **VIAF** [Web8; Web9].

Knowledge Bases According to the insights from the **SoNAR** project, (open) knowledge bases (KBs) can be useful when trying to overcome the problem with missing or unbalanced data in authority files such as the **GND** (see Section 2.1). KBs are usually edited independently of the library domain by a wider community of contributors, and their quality cannot be expected to meet the

standards of catalogues or authority files edited by library personnel. A prominent example of an open, cross-domain KB is **Wikidata** [Web10], which has tentatively been used in the context of the **SoNAR** project and which we have already identified in Section 1.1 as possible source of metadata for complementing authority data.

Cultural Heritage Databases This category contains generic federated portals of cultural heritage items as well as special databases created especially for the support of provenance research. Examples for these two kinds of databases are, respectively, **Europeana** [Web46] (see Section 2.2) and **Proveana**—the research database of the German Lost Art Foundation [Web81].

4.3.2 Analysis of Data Sources

We now provide a review of four data sources from the previous list, namely **DNB**, **K10plus**, **GND**, and **Wikidata**. For each data source, we collect the following information.

- *Scope*: thematic focus; coverage; number of records; standards for cataloguing
- *Technical infrastructure*: data format; data model; interfaces; support for linked data
- *Data quality* (from the **SoNAR** grant proposal [Schneider-Kempf et al. 2018, p. 19ff.]), if applicable: use of persistent identifiers and **URIs**; adherence to standards for, e.g., time and date specification
- *Features useful in the context of SoNAR* (see our discussion in Section 2.1.4), if applicable: recording of temporal attributes or relations; recording of data provenance
- Further features specific to the respective data source, if applicable

DNB The following information is taken from the **DNB**’s websites under the category “DNB Professional” [Web82; Web83; Web84].

The DNB adheres to a legal collection mandate, according to which the DNB collects “all texts, images and sound recordings published in Germany or in the German language, translated from German or relating to Germany that have been issued since 1913” [Web82]. This includes all physical publications, and, since 2006, electronic publications made available via the internet. The mandate commits the DNB to a complete and unbiased collection that includes, among others, “printed works compiled or published between 1933 and 1945 by German-speaking emigrants” [Web82].

The DNB catalogues its entire collection both descriptively and by subject, adhering to standards such as **RDA**, using authority data, and including persistent identifiers such as ISBN, ISSN, URN, and DOI. The cataloguing data feeds the German National Bibliography.

According to the 2021 annual report [German National Library 2022], the DNB’s holdings comprise 43.7 million physical or digitally accessible units, which are represented by almost 26 million records in the German National Bibliography.

Metadata can be obtained from the DNB freely (under the CC0 1.0 licence) via the interfaces **SRU** and **OAI-PMH**, which support **XML** serialisations of the data formats **MARC 21**, **RDF**, DNB Casual (an **XML**-based **DC** format), and **MODS**. Further data formats are available via an individual access to the DNB’s *Data Shop*. The DNB’s Linked Data Service provides open access to its bibliographic and authority data in **RDF** under the same license. Instructions on how to interact with these interfaces are given on the DNB’s webpage on metadata services [Web84]. There DNB website also reports on a project for the conversion of its data into the **BIBFRAME** format [cf. Web85]: for title records, the DNB OPAC provides a download button for the **BIBFRAME** format. The state of this project beyond 2014 is left unclear.

According to the detailed specifications on the [MARC](#) format by the DNB [[Web86](#); [Web87](#)], the recording of dates and times, languages, geographic area codes, and countries conforms to the respective [ISO](#) standards. Since 2015, the DNB has been using [MARC](#) field 883 for recording the metadata provenance for a selection of data fields [[Web88](#)]. Furthermore, the DNB's [MARC 21](#) schema provides for the cataloguing of the following temporal data concerning title (i.e., manifestation) records: production, publication, distribution, and manufacturing dates (including reproductions and reissues); notes on dates and times; chronological subject term (as full text); graduation year (for dissertations etc.).

K10plus [K10plus](#) is the joint catalogue of the German library networks [GBV](#) and [SWB](#). Together, these two networks comprise more than 1400 national, regional, academic, and public libraries [[Web89](#); [Web90](#)], of which 838 participate in [K10plus](#), according to the list of participating institutions in the [K10plus Wiki](#) [[Web91](#)]. Thus, [K10plus](#) comprises the data from the majority of German academic institutions [cf. [Web92](#)].

As of 31 December 2022, [K10plus](#) contains 80.8 million bibliographic records (“Titeldatensätze”) with 235.4 million ownership records [[Web93](#)]. Cataloguing adheres to the same standards and guidelines as in the DNB, including similar specifications for the recording of temporal data and the use of [ISO](#) standards for dates, times, etc. In addition, the structured variants of provenance indexing used in [GBV](#) and [SWB](#) (see below) enable the recording of ownership dates if available.

[K10plus](#) provides metadata freely (under the CC0 1.0 licence), mainly via the interfaces [Z39.50](#) and [SRU](#). Those support the data formats [MARC 21](#), (including its XML variants), [PICA+](#) and [PICA-XML](#) (a provider-specific internal data format loosely based on [MARC 21](#)), as well as [DC](#), [MODS](#), and the legacy format [MAB2](#). Detailed instructions on how to interact with these interfaces are given in the [K10plus Wiki](#) [[Web91](#)]. In addition, snapshots of [K10plus](#) data are provided regularly in [MARC-XML](#) and as linked open data ([RDF-XML](#)), but the information on these in the [K10plus Wiki](#) is incomplete [[Web94](#)].

[K10plus](#) uses the [GND](#) as the central authority file for disambiguating persons, corporate bodies, works, subject terms, and further entities (see below for more details on the [GND](#)). For this purpose, [K10plus](#) includes copies of all [GND](#) records, which are synchronised via an [OAI-PMH](#) interface within minutes after modification [cf. [BSZ](#) and [VZG 2022](#)]. [K10plus](#) furthermore uses authority data from the [Basisklassifikation \(BK\)](#) [[Facharbeitsgruppe Sacherschließung 1995](#)] and the [Regensburger Verbundklassifikation \(RVK\)](#) [[Web95](#)], cf. the [K10plus Wiki](#) [[Web96](#)].

Concerning the state of provenance indexing in the library networks [GBV](#) and [SWB](#), [Hakelberg \[2016, §4\]](#) reports the following. The catalogue systems of [GBV](#) and [SWB](#), which are united in [K10plus](#), support the recording and display of provenance features on the bibliographic as well as on the exemplar level of a record. However, the practice of provenance indexing differs greatly between the two networks. Since 2014, [GBV](#) has been advising the recording on the bibliographic level, in order to allow provenance research across libraries. The names of owners are linked to the local authority record that represents, and redirects to, the respective record in [GND](#). The physical provenance features are documented via chains of terms from [T-PRO](#) (see [Section 2.3](#)) [plus dates if available]. As of 2014, this standard was being tested and successively introduced in the participating libraries. Before 2014, provenance features were documented in free text in a data field for exemplar-specific comments on the exemplar level, using [T-PRO](#) descriptors and referring to local authority records that were *not* directly linked to [GND](#) records. In contrast, the [SWB](#) libraries have always been using a dedicated data field on the exemplar level to record provenance features. This field provides subfields for the structured recording of [T-PRO](#) terms [and dates]. However, most [SWB](#) libraries have abandoned the creation of exemplar records altogether and needed to resort to using isolated local systems for provenance indexing [cf. [Hakelberg 2016](#)].

In summary, provenance features are recorded within [K10plus](#) in a heterogeneous and incomplete

way, and developers of a retrieval system that uses K10plus have to be aware of this problem.

GND The “Gemeinsame Normdatei”, the Integrated Authority File of the German National Library (GND) [Web2], is operated by the DNB “in cooperation with many other libraries, libra[r]y networks and other cultural and academic institutions. At present, the GND contains around 9 million authority records for persons, corporate bodies, congresses, geographic entities, specialised terms and works; these are supplemented, updated and used frequently” [Web83]. More detailed statistics can be found in the DNB’s 2021 annual report [German National Library 2022, p.49]. GND adheres to the same cataloguing standards and offers the same technical infrastructure as the DNB catalogue; in particular, GND data is accessible under the CC0 1.0 license alongside national bibliographic data via the same interfaces and data formats, including several RDF serialisations.

The DNB’s MARC 21 schema [Web86] provides for the cataloguing of the following temporal data in authority records: biographical data (persons), date of publication (works), dates of existence (conferences, corporations), year of discovery, dates of activity (persons, corporations), years of effect of relationships between entities of different kinds. Examples for temporal data of the latter kind include the time interval during which a scholar was member of a university or during which a person was a student of another.

We have already learnt from the SoNAR project (Section 2.1) that the data in the GND is incomplete for several reasons. Another reason is certainly connected with the individualisation guidelines in the DNB’s guidelines for recording persons and families in the GND [Web97, part EH-P-16], which specify the additional information that needs to be provided in order to disambiguate a person or family. These guidelines define two groups of features that can be used to individualise (i.e., disambiguate) a person, such as biographical data or relationships to other persons, among many others. Depending on the level of the respective authority record (which indicates the state and quality of the record), one or only a few of the many features from these two groups have to be recorded. Therefore, there is no guarantee that, for example, relationships to other people or temporal data of the kind mentioned above is recorded. In addition, relationships to other people are recorded in an unnormed way, using one of four very generic codes (e.g., acquaintance, professional relationship, family relationship) and specifying the exact kind of relationship in free text. Similarly, professions may be recorded in free text or via a link to an authority record.

Wikidata Wikidata [Web10] is “a free, collaborative, multilingual, secondary database, collecting structured data to provide support for Wikipedia, Wikimedia Commons, the other wikis of the Wikimedia movement, and to anyone in the world” [Web98]. Hence, there is no thematic restriction on the content of Wikidata.

The main unit of information in Wikidata is the *item*. Items “are used to represent all the *things* in human knowledge, including topics, concepts, and objects” [Web99]. Currently, Wikidata contains more than 100 million items [Web100]. This number is two orders of magnitudes higher than the number of articles in Wikipedia because Wikidata stores structured information on *all* Wikimedia projects, including the large media file repository Wikimedia Commons [cf. Web99]. Data on each item is stored in *statements*, which have the same **S–P–O** form as **RDF**, here called *item–property–value*. Statements can have annotations [cf. Web101]. All items and properties have unique identifiers consisting of a Q or P, respectively, followed by a sequence of digits. These IDs are also part of the URL of the respective Wikidata page and of the item’s or property’s URI. For example, the item representing the Polish mathematician and astronomer Nicolaus Copernicus has the ID Q619, and the statements about this item are listed on the respective Wikidata page [cf. Web102]. The first statement in this list says that Copernicus is a human, using the property *instance_of* (P31) and the value human (item Q5).

Data is recorded in Wikidata collaboratively, i.e., “[d]ata is entered and maintained by Wikidata editors, who decide on the rules of content creation and management. Automated bots also enter

data into Wikidata” [Web98]. Values such as times and dates are recorded conforming to the respective ISO standard.

The Wikidata page on data access [Web100] lists numerous ways to access data, all under the free CC0 license. For retrieving individual entries via URIs, there is a linked data interface. For retrieving a relatively small set of entries that are not known in advance, SPARQL queries can be posed against the Wikidata Query Service. For retrieving larger sets of entries, Wikidata offers the Linked Data Fragments endpoint, which requires computational power on the client side. Furthermore, there are two APIs mostly for editing Wikidata, as well as a recent changes stream and complete exports (dumps) [cf. Web100].

Annotations of statements (see above) seem to be suitable for recording metadata provenance and temporal attributes of relationships. It remains to investigate the use of annotations systematically in order to determine the extent to which metadata provenance and temporal attributes are recorded.

4.4 Conclusions

We have collected 20 data sources and reviewed four of them in detail. These four have in common that they provide their metadata as linked data (among other formats), and the three bibliographic data sources adhere to the FRBR entity-relationship model. From this commonality, we conclude that entities and relationships—in mathematical terms, constants and unary and binary relations—should play a central role in the model of queries, answers, and data sources that we are going to develop next. This means that a graph-based type of model, which also features in the SoNAR project, is suitable. In case relations of higher arity turn out to be needed, for example when representing the provenance of an S-P-O statement, those can always be represented by several binary relations via reification [cf. Doan, Halevy, and Ives 2012, p.339f.].

The choice of the reviewed data sources was motivated by the insights from the SoNAR project and Hakelberg’s thesis [Hakelberg 2016]. For the abstract model to be developed next, the choice of concrete data sources is largely unimportant because the model should be independent on the contents or specific data models of the sources. The same holds for the general method for answering queries that we are going to develop based on our model. However, future implementations of our method will have to take the specifics of the data sources into account. For this reason, it will be necessary to review further data sources, compare them quantitatively with respect to the extent and quality of their data, and find structured ways for dealing with the problems concerning the data quality reported above, such as the ambiguous links to works and expressions in RDA-based catalogues, the heterogeneous and partly unstructured recording of provenance, and the bias and incompleteness of the GND. The latter problem might be alleviated by using Wikidata (sacrificing the strict quality control of the GND), and a promising candidate for an alternative provenance database is Proveana.

5 A General Model of Provenance Relationships

In this chapter, we develop a generic approach to modelling provenance relationships. More precisely, we need to model three central notions: queries that a user may want to ask, data sources that are to be consulted in order to answer a query, and answers given to a query. In order to obtain a generic approach, we aim at providing rigorous definitions for these central concepts, and we seek intensional rather than extensional definitions. In particular, those definitions should not depend on concrete example queries or data sources such as the ones discussed in Chapters 3 and 4; neither should they depend on concrete objects, concepts, or relations (such as “Copernicus”, “exemplar”, or “student”). Instead, we will develop an abstract model that formalises the notions of a query, data source, and answer. This model can then be instantiated with a multitude of concrete queries and data sources, and it constitutes the basis of the method for finding answers that we will develop in the next chapter.

As a basis for our abstract model, we choose standard concepts and techniques from graph theory. The concept of a graph is widely used in computer science and discrete mathematics; see standard textbooks [e.g., Diestel 2012]. Graphs and graph techniques are widely applied in various areas such as computer science, linguistics, physics and chemistry, social sciences, and biology [Web103]. They are also the foundation of social networks [Galety et al. 2022] and therefore highly relevant for historical research, e.g., in the SoNAR project, as we have seen in Section 2.1. Graphs are used to represent large knowledge bases [e.g., Ehrlinger and Wöß 2016] and are a fundamental ingredient of RDF. Furthermore, the basic definition of a graph is conceptually simple, and graph theory provides a plethora of well-understood concepts and algorithms. By utilising graph theory, our application scenario can benefit from these concepts, algorithms [Diestel 2012; Even 2012], and implementations [Web104; Web105].

The main idea of our abstract model is the following. Both queries and data sources are represented as graphs (typically a rather small graph for the query and a very large graph for the data source). Answers to the query are those parts of the graph that have the same structure as the query. In more formal words, answers are found using *pattern matching* techniques known from querying RDF graphs via SPARQL [Della Valle and Ceri 2011] and from database and graph theory [Abiteboul, Hull, and Vianu 1995; Diestel 2012].

In the following sections, we introduce our model by defining the basic terminology (Section 5.1), the specific notion of a graph (Sections 5.2, and the abstract notions of data sources, queries, and answers (Section 5.3). We also discuss decision procedures related to query answering (Section 5.4). Up to that point, the model remains conceptually simple and is based on elementary and self-contained mathematical definitions. We provide additional explanations and illustrations for the sake of readers with little or no background in mathematics. In order to obtain a more flexible and comprehensive model, we discuss the underlying modelling decisions in Section 5.5 and possible extensions in Section 5.6.

5.1 Basic Notions

We start by introducing the basic terminology that we are going to use in the following, which consists of the usual terms from conceptual modelling [Brodie, Mylopoulos, and Schmidt 1984].

Objects An *object* is a specific entity, for example, the work “Graph Theory” by Reinhard Diestel, its expression in English, its manifestation as the 4th edition by Springer, a specific copy, or the person Reinhard Diestel.

Concepts A *concept* is a class of objects, such as the entities Work, Expression, etc. from the FRBR model (see Section 4.1).

We determine that names for concepts and objects start with an uppercase letter.

Relations An *n*-ary *relation* is a set of *n*-tuples of objects; for example, the binary relation creator consists of pairs of objects including (Graph_Theory, Reinhard_Diestel) if we assume for the sake of simplicity that Graph_Theory is the unambiguous name of said work.

Constants and concepts are in fact nullary and unary relations, but it is more intuitive to use those separate terms. As we have already argued in Section 4.4, we will hardly need to deal with relations of arity beyond 2; therefore we will mostly discuss binary relations and refer to them as *relations* when no confusion may arise.

To distinguish relations from concepts and objects, we determine that their names start with a lowercase letter.

Converses of relations The *converse* of a relation *R* is the relation obtained from *R* by swapping the components in each pair, e.g., the converse creator_of of creator contains the pair (Reinhard_Diestel, Graph_Theory).

Instances An *instance* of a concept or relation is an object or a pair of objects from that set; e.g., Reinhard_Diestel is an instance of Person, and (Graph_Theory, Reinhard_Diestel) is an instance of creator.

Relationships Instances of relations are called *relationships*.

Literals A *literal* is a fixed value, such as a year, date, or identifier.

5.2 Labelled Directed Graphs

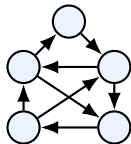


Figure 5.1: A directed graph

Graphs consist of nodes and edges. Edges link nodes and can be directed or undirected. Graphs are easy to visualise; nodes are represented as circles or rectangles, and edges as arrows (directed) or lines (undirected). Figure 5.1 shows an abstract example of a directed graph.

For our purposes, nodes represent objects or literals. Edges represent relations, which are directed by the above definition (as the order of components in a pair matters). For example, an edge representing the relation creator points *from* a person or corporate body *to* a work, whereas its converse creator_of points into the opposite direction. Therefore we use *directed* graphs. Symmetric relations, such as relatedEndeavour, can be represented via pairs of edges pointing in both directions.

Furthermore, we want to assign a unique name to each node of a graph and one or several labels to each node and each edge: The name of a node specifies the object that is represented by that node. The labels of a node specify the concepts of which that node is an instance. For example, a node representing the physicist Albert Einstein (object) may be labelled, among others, with the concepts Person, Scientist, and Physicist. The labels of an edge specify the relations of which the pair of nodes represented by that edge is an instance. For example, if a person P_1 has a student P_2 and also collaborates with P_2 , then this can be represented via an edge from P_1 to

P_2 with the label {has_student, collaborates_with} (and/or an edge from P_2 to P_1 with the label {is_student_of, collaborates_with}). These considerations lead us to a straightforward extension of the notion of a directed graph: a *labelled directed graph*.

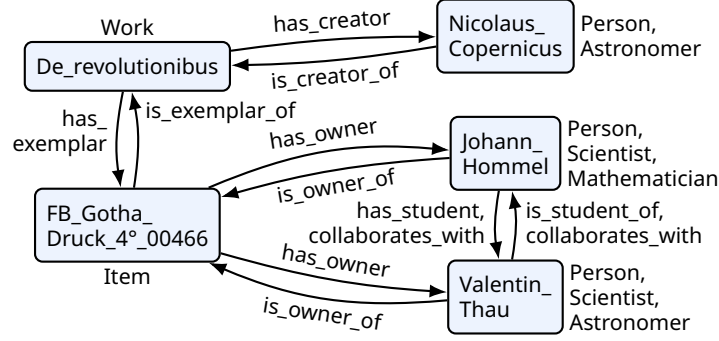


Figure 5.2: A labelled directed graph that represents data concerning an exemplar of Copernicus’ *De revolutionibus* and some of its owners

In order to visualise a labelled directed graph, node names are written into the respective node, and node and edge labels are written next to the node or edge. Multiple labels of the same node or edge are delimited with commas. An example is given in Figure 5.2. The shown graph represents a part of the data described in Section 3.2 concerning an exemplar of Copernicus’ *De revolutionibus* at the Gotha Research Library (*FB Gotha*). It contains a node for the work (labelled with the FRBR entity Work), a node for the exemplar (labelled with the FRBR entity Item), and nodes for the author and two of the owners (labelled with their professions according to their GND entries). For the sake of this example, the owners are additionally labelled with the profession Scientist, which is implicit in the real data. The ten edges represent relationships that are instances of FRBR relations and others between Work, Item, and Person. For the sake of simplicity, the graph deviates from the FRBR model [IFLA SG FRBR 2009] by omitting the FRBR entities “Expression” and “Manifestation” that should occur between the nodes labelled “Work” and “Item”.

Unify typetting of terms (i.e., Term instead of “Term”); unify (FRBR) relation names (e.g., has_creator vs. creator)

As we will see in the following, labelled directed graphs can be used in our setting to represent (combinations of) data sources as well as queries. They allow us to draw on standard notions from graph theory and query answering in order to define admissible query answers and to devise methods for obtaining those.

The above explanations can be cast into a rigorous mathematical definition, which uses sets to represent nodes, a binary relation over the set of nodes to represent edges, and functions over the nodes and edges to represent names and labels. In order for the range of those functions to be well-defined, the definition of a labelled directed graph is relative to a namespace which contains all the names of objects, concepts and relations that are relevant. The contents of this namespace is arbitrary and may consist, for example, of all the names found in the relevant data sources. The following definition introduces the notions of a namespace and a labelled directed graph.

Definition 1. Let $N = (N_O, N_C, N_R)$ be a *namespace* consisting of a set N_O of *object names*, a set N_C of *concept names*, and a set N_R of *relation names*. A *labelled directed graph* over N is a triple $G = (V, E, N, \mathcal{L})$ where

- V is a set, whose members are called *nodes*;^a
- $E \subseteq V \times V$ is a set of pairs of nodes, whose members are called *edges*;

- $\mathcal{N} : V \rightarrow \mathcal{N}_O$ is an injective function that assigns to each node a unique object (called the node's *name*);
- $\mathcal{L} : V \cup E \rightarrow \mathcal{N}_V \cup 2^{\mathcal{N}_R}$ is a function that assigns to each node a set of concept names (called the node's *labels*) and to each edge a non-empty set of relation names (called the edge's *labels*); we call \mathcal{L} a *labelling function*.

^a↑In classical graph theory, nodes are called *vertices*; thus the set of nodes of a graph is denoted by V . We adopt the denotation V for conformity and the more modern term “node” for brevity.

Definition 1 formalises the following commitments regarding names and labels. (1) Every node has a unique name, and no two nodes have the same name (the latter being ensured by injectivity). (2) A node can have an arbitrary number of labels, including no label (in case the node belongs to no concept). (3) An edge can have an arbitrary number of labels, but that number must not be zero – the effect of an edge having no labels can be achieved by simply omitting that edge.

In order to illustrate the components of Definition 1, we refer to the graph depicted in Figure 5.2: V consists of five nodes, and E of ten edges (each single arrow constitutes an edge since the direction matters). There are, among others, the following node names and labels:

- The node at the top left has the name `De_revolutionibus` and the single label `Work`.
- The node at the bottom right has the two labels `Person` and `Astronomer`.
- The edge from the node named `Johann_Hommel` to the node named `Valentin_Thau` has the two labels `has_student` and `collaborates_with`, and the edge pointing into the converse direction has the labels `is_student_of` and `collaborates_with`.

5.3 Modelling Data Sources, Queries, and Answers

We can now use our notion of a labelled directed graph to model data sources and queries, and to obtain a rigorous definition of a query answer.

5.3.1 Data Sources

Data sources correspond exactly to our notion of a graph.

Definition 2. A *data source over the namespace* $\mathcal{N} = (\mathcal{N}_O, \mathcal{N}_C, \mathcal{N}_R)$ is a labelled directed graph over \mathcal{N} .

In our model, we assume that there is always a *single* data source against which a query is posed and evaluated. When the model is applied to real-word queries and data sources, the abstract notion of a data source is instantiated by the union of all available concrete data sources (such as catalogues, authority files, knowledge bases), including mappings between them if applicable.

5.3.2 Queries

In order to model queries based on the same notion of a graph, we need to distinguish two special groups of nodes that act as placeholders (1) for the object(s) after which the query asks and (2) for further objects that are mentioned in the query without being named explicitly. For example, consider Query **Q2'** from Section 3.2:

Q2' Which exemplars of *De revolutionibus* were owned by some scientist who passed them on to a student?

To model this query, we do not only need a node representing the work *De revolutionibus*, but also a node representing an exemplar that satisfies the conditions stated in the query and whose name is asked for (Group 1), and two nodes representing the owner and their student (Group 2). Since these three individuals are not known, we need to use *variables* for naming them. Query **Q2'** can then be modelled by the graph shown in Figure 5.3.

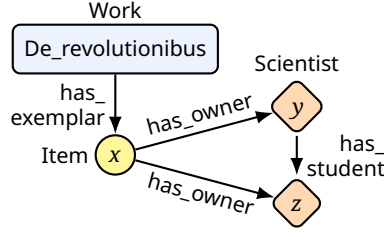


Figure 5.3: A graph representing example query **Q2'**

The nodes of this graph fall into three groups:

- (1) The node named De_revolutionibus represents that work;
- (2) Node x falls into Group 1 as explained above;
- (3) Nodes y and z fall into Group 2 as explained above.

Node names x, y, z are the variables mentioned above, and we call x the *answer variable* and y, z the *anonymous variables* of this query. From now on, we fix two sets VAR_{ANS} and VAR_{ANON} of *answer variables* and *anonymous variables*, respectively, and we assume that they both contain a countably infinite number of elements—which simply ensures that there is an unlimited supply of variables. We furthermore require that these two sets are disjoint with each other and with any set N_O of object names. Thus, according to Definition 2, graphs representing data sources cannot use any variables as node names.

In order to allow queries to use variables, the following definition of a query is now immediate.

Definition 3. A query over the namespace (N_O, N_C, N_R) is a labelled directed graph over $(N_O \uplus \text{VAR}_{\text{ANS}} \uplus \text{VAR}_{\text{ANON}}, N_C, N_R)$.

The operator “ \uplus ” used in this definition stands for *disjoint union*, i.e., “ $A \uplus B$ ” stands for the union of the *disjoint* sets A, B . The wording of the definition also ensures that we do not have to mention variables explicitly when specifying the namespace of a query.

5.3.3 Query Answers

Based on the representation of both data sources and queries as graphs, we can now define the notion of an answer to a query with respect to a data source. For this purpose, it is important to realise that, typically, a query is a small graph and a data source is a large graph, and that finding answers simply means finding small parts of the large graph that have the same structure as the small graph. For the example query given in Figure 5.3, this means that every subgraph of the data source consisting of four nodes with the same edges and labels should be an answer. Regarding the

previous example, the subgraphs given in Figure 5.4 (a, b) constitute answers, while the subgraphs in Figure 5.4 (c, d) do not: Subgraph (a) is identical to the query graph with the only exception that it contains proper objects instead of variables in the node names; Subgraph (b) is the same graph but extended with additional information; Subgraph (c) lacks the edge labelled `has_student` between the two owners; Subgraph (d) resembles the structure of the query but does not contain the required node named `De_revolutionibus`.

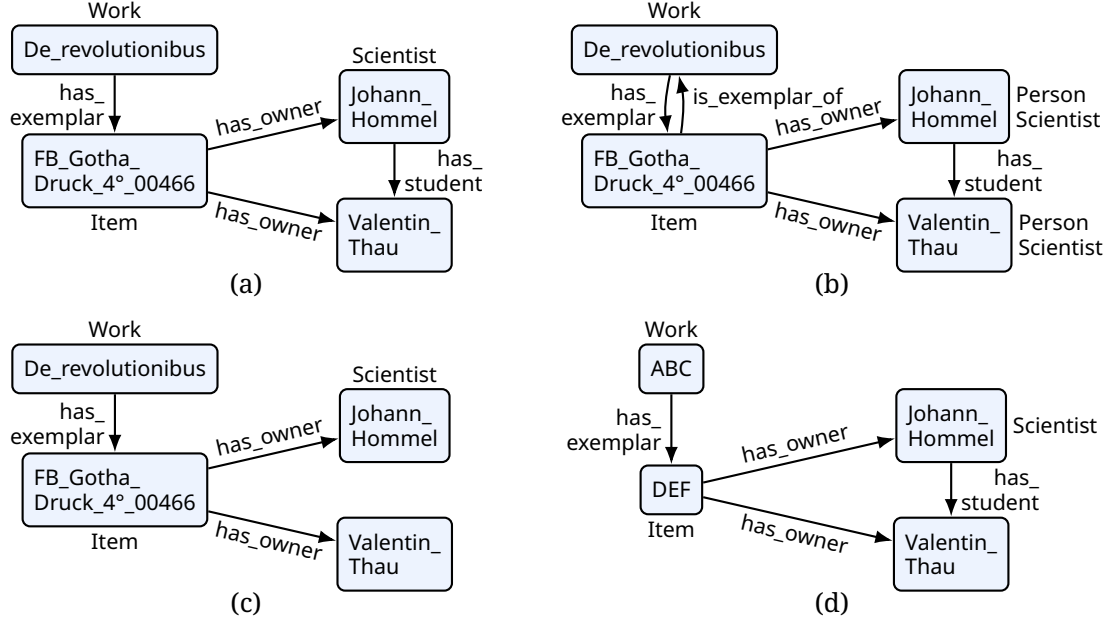


Figure 5.4: Positive (a, b) and negative (c, d) examples for query answers

In order to identify subgraphs of a given graph that have the same structure as another given graph, we use the notion of a homomorphism. A homomorphism is a function that maps some object to another while preserving the structure of the former. We therefore need to define a variant of homomorphisms that maps queries to data sources. This variant is given in the following.

Definition 4. Let $N = (N_O, N_C, N_R)$ be a namespace, $G = (V, E, \mathcal{N}, \mathcal{L})$ a query over N , and $G' = (V', E', \mathcal{N}', \mathcal{L}')$ a data source over N . A *homomorphism from G to G'* is a map $h : V \rightarrow V'$ that satisfies the following properties.

H1 $\mathcal{N}(v) = \mathcal{N}'(h(v))$ for every node $v \in V$ with $\mathcal{N}(v) \in N_O$.

H2 $\mathcal{L}(v) \subseteq \mathcal{L}'(h(v))$ for every node $v \in V$.

H3 $\mathcal{L}(v_1, v_2) \subseteq \mathcal{L}'(h(v_1), h(v_2))$ for every edge $(v_1, v_2) \in E$.

If h is a homomorphism from G to G' , we write $h : G \rightarrow G'$. If there is some homomorphism from G to G' , we write $G \lesssim G'$.

Property **H1** requires that a homomorphism maps each node in G that is named with an object to that node in G' which is named with the same object. Nodes named with variables in G can be mapped to arbitrary nodes in G' . Properties **H2** and **H3** require that homomorphisms preserve node and edge labels; more precisely, the image of a node (or edge) under h must have *at least* the same labels (and may have additional labels).

Figure 5.5 shows a homomorphism h (dashed lines) from the query depicted in Figure 5.3 to the graph from Figure 5.2.

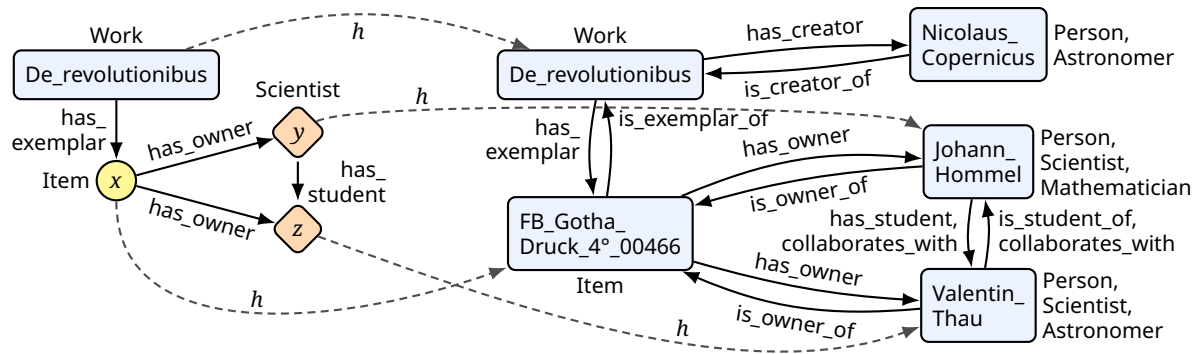


Figure 5.5: An example homomorphism

5.4 Decision Procedures

Formulate decision problem, discuss (data) complexity & algorithms. (Reduction to FO/SQL?)

5.5 Discussion of the Modelling Decisions

Relate this “machinery” to the example queries. In particular:

- Comment on Boolean queries if necessary.
- Discuss specific requirements for modelling **Q1** and **Q3**:
 - **Q1** seems to require answer variables representing sets (the owners) and an appropriate extension of the definition of a homomorphism;
 - the same holds for **Q3**; additionally the answer should include the relationships between the images of the answer variables (the relationships between the owners), i.e., some sort of spanned subgraph
- ↪ extensions needed; description or sketch of these and other extensions in the next section

5.6 Possible Extensions

TODO: Discuss further extensions:

- relations of arbitrary arity
- data provenance
- concrete values (“year of publication”)
- attributes on relationships:
 - sketch idea: e.g., provide year for relationship has_owner – example: “passed on to” requires descending year numbers *and* no successor with intermediate year number

- solution: quads instead of triples (in [RDF](#) speak); add attributes to the graph model? (Markus Krötzsch's work?)
- explain difficulties: more complex formal machinery (def. of graphs, queries, and matches); missing data, e.g.:
 - * From which year *to which year* did person *X* own item *Y*?
 - * Was person *Z* a student of person *X*'s *at the point in time when X passed the item on to Z*?
- discuss usefulness: false positives due to incomplete data as discussed in [Section 3.3](#) \leadsto manual inspection is necessary anyway; attributes may still help hide false answers
- Discuss [SoNAR](#) requirements R0xx:
 - Some requirements can be addressed directly with our queries, e.g., R006! (R016?!?)
 - R016 seems to require the “ \top -role” in queries!
 - Edge weights;
 - E.g. R029, 30, 61: explorative (for every node/edge the available attributes) vs. ??? (ask query knowing those attributes)
- Discuss further SoNAR insights, e.g., substitute relations

6 Automated Retrieval of Provenance Relationships

- develop method for answering queries in the model just defined
- online vs. offline scenario:
 - offline:

data source graph is generated explicitly (using data integration techniques) and updated in fixed intervals

at runtime: provenance query is formulated as a [SPARQL](#) query and posed against the graph
 - online:

data source graph is implicit; data sources are queried “on the fly” (this is the preferred scenario here; see delineation from SoNAR in §2.1)

data sources need to be defined and interactions with them programmed in advance

at runtime: decompose query into single [SPARQL](#) (sub)queries and pose them against several data sources, potentially iteratively; compose final answer from the partial answers \leadsto demonstrate this for example query/ies!

7 Conclusion

- get back to initial research question and its subordinate questions
- Acknowledgements? See comments.

References

Bibliography

- Abiteboul, Serge, Richard Hull, and Victor Vianu (1995). *Foundations of databases*. Addison-Wesley (cit. on p. 29).
- Alemu, Getaneh et al. (2012). “Linked Data for libraries: Benefits of a conceptual shift from library-specific record structures to RDF-based data models”. In: *New Library World* 113.11-12, pp. 549–570. DOI: <https://doi.org/10.1108/03074801211282920> (cit. on p. 13).
- Baader, Franz et al. (2017). *An Introduction to Description Logic*. Cambridge University Press. URL: <http://www.cambridge.org/de/academic/subjects/computer-science/knowledge-management-databases-and-data-mining/introduction-description-logic?format=PB#17zVGeWD2TZUeu6s.97> (cit. on p. 12).
- Bair, Sheila, ed. (2013a). *Journal of Library Metadata* 13.2–3. Special issue “Linked Data, Semantic Web and Libraries”. URL: <https://www.tandfonline.com/toc/wjlm20/13/2-3> (cit. on p. 13).
- (2013b). “Linked Data—The Right Time?” In: *Journal of Library Metadata* 13.2-3, pp. 75–79. DOI: [10.1080/19386389.2013.828527](https://doi.org/10.1080/19386389.2013.828527) (cit. on p. 13).
- Berners-Lee, Tim, Roy T. Fielding, and Larry M. Masinter (Jan. 2005). *Uniform Resource Identifier (URI): Generic Syntax*. RFC 3986. DOI: [10.17487/RFC3986](https://doi.org/10.17487/RFC3986) (cit. on pp. 11, 51).
- Berners-Lee, Tim, James Hendler, and Ora Lassila (2001). “The Semantic Web”. In: *Scientific American* 284.5, pp. 34–43. DOI: [10.1038/scientificamerican0501-34](https://doi.org/10.1038/scientificamerican0501-34) (cit. on p. 11).
- Bizer, Christian, Tom Heath, and Tim Berners-Lee (2009). “Linked Data – The Story So Far”. In: *International Journal on Semantic Web and Information Systems* 5.3, pp. 1–22 (cit. on p. 13).
- Bludau, Mark-Jan et al. (2020). “SoNAR (IDH): Datenschnittstellen für historische Netzwerkanalyse”. In: *7. Tagung Digital Humanities im deutschsprachigen Raum (DHd 2020)*. Ed. by Tara Andrews et al. DOI: [10.5281/zenodo.4621861](https://doi.org/10.5281/zenodo.4621861) (cit. on pp. 6, 8, 51).
- Boumechaal, Hasna and Zizette Boufaïda (2023). “Complex Queries for Querying Linked Data”. In: *Future Internet* 15.3. DOI: [10.3390/fi15030106](https://doi.org/10.3390/fi15030106) (cit. on p. 14).
- Brodie, M. L., J. Mylopoulos, and J. W. Schmidt, eds. (1984). *On Conceptual Modelling: Perspectives from Artificial Intelligence, Databases, and Programming Languages*. 1st. Topics in Information Systems. Springer (cit. on p. 29).
- BSZ and VZG (May 2022). *Normdaten*. German. Handbook for using and creating GND authority data in K10plus, linked in K10plus Wiki. URL: https://opus.k10plus.de/frontdoor/deliver/index/docId/410/file/K10plus_Normdaten.pdf (visited on 29/05/2023) (cit. on p. 26).
- Burrows, Toby et al. (2021). “A New Model for Manuscript Provenance Research: The Mapping Manuscript Migrations Project”. In: *Manuscript Studies* 6.1, pp. 131–144. URL: https://repository.upenn.edu/mss_sims/vol6/iss1/5/ (cit. on p. 13).
- Byrne, Gillian and Lisa Goddard (2010). “The Strongest Link: Libraries and Linked Data”. In: *D-Lib Magazine* 16.11/12. DOI: [doi:10.1045/november2010-byrne](https://doi.org/10.1045/november2010-byrne) (cit. on p. 13).

- Copernicus, Nicolaus (1543). *Nicolai Copernici Torinensis De revolutionibus orbium coelestium, Libri VI*. Latin. Ed. by Andreas Osiander. Nürnberg: Petreius, Johann, [6], 196 sheets. URL: <https://opac.uni-erfurt.de/LNG=EN/DB=1/PPNSET?PPN=567506266> (cit. on pp. 3, 15, 16).
- Danowski, Patrick and Adrian Pohl, eds. (2013). *(Open) Linked Data in Bibliotheken*. German. Berlin, Boston: De Gruyter Saur. DOI: [doi:10.1515/9783110278736](https://doi.org/10.1515/9783110278736) (cit. on p. 12).
- Della Valle, Emanuele and Stefano Ceri (2011). “Querying the Semantic Web: SPARQL”. In: *Handbook of Semantic Web Technologies*. Ed. by John Domingue, Dieter Fensel, and James A. Hendler. Berlin, Heidelberg: Springer. Chap. 8, pp. 299–363 (cit. on pp. 12, 29).
- Diestel, Reinhard (2012). *Graph Theory*. 4th. Vol. 173. Graduate texts in mathematics. Springer (cit. on p. 29).
- Doan, AnHai, Alon Y. Halevy, and Zachary G. Ives (2012). *Principles of Data Integration*. Morgan Kaufmann. DOI: <https://doi.org/10.1016/C2011-0-06130-6> (cit. on pp. 11, 14, 28).
- Domingue, John, Dieter Fensel, and James A. Hendler, eds. (2011). *Handbook of Semantic Web Technologies*. Berlin, Heidelberg: Springer. DOI: [10.1007/978-3-540-92913-0](https://doi.org/10.1007/978-3-540-92913-0) (cit. on p. 11).
- Eckert, Kai (2012). “Metadata Provenance in Europeana and the Semantic Web”. German. Master’s thesis. Humboldt-Universität zu Berlin. DOI: <http://dx.doi.org/10.18452/14172> (cit. on p. 14).
- (2013a). “Die Provenienz von Linked Data”. German. In: *(Open) Linked Data in Bibliotheken*. Ed. by Patrick Danowski and Adrian Pohl. Berlin, Boston: De Gruyter Saur, pp. 97–121. DOI: [doi:10.1515/9783110278736.97](https://doi.org/10.1515/9783110278736.97) (cit. on p. 14).
- (Sept. 2013b). “Provenance and Annotations for Linked Data”. In: *International Conference on Dublin Core and Metadata Applications*, pp. 9–18. URL: <https://dcpapers.dublincore.org/pub/s/article/view/3669> (cit. on p. 14).
- Ehrlinger, Lisa and Wolfram Wöß (2016). “Towards a Definition of Knowledge Graphs”. In: *Joint Proceedings of SEMANTiCS 2016 and SuCESS 2016, Posters and Demos Track*. Ed. by Michael Martin, Martí Cuquet, and Erwin Folmer. Vol. 1695. CEUR Workshop Proceedings. CEUR-WS.org. URL: <https://ceur-ws.org/Vol-1695/paper4.pdf> (cit. on p. 29).
- Even, Shimon (2012). *Graph algorithms*. Ed. by Guy Even and Richard M. Karp. 2nd. Literaturangaben. Cambridge [u.a.]: Cambridge Univ. Press. XII, 189 (cit. on p. 29).
- Facharbeitsgruppe Sacherschließung, ed. (1995). *Basisklassifikation für den Bibliotheksverbund Niedersachsen/Sachsen-Anhalt/Thüringen*. German. 2nd, revised edition. Göttingen. URL: <https://kxp.k10plus.de/DB=2.1/DB=2.1/PPNSET?PPN=192887122> (cit. on pp. 26, 50).
- Fangerau, Heiner et al. (Jan. 2022). *SoNAR AP 2*. German. Final report on Work Package 2 of the project. URL: <https://github.com/sonar-idh/reports/blob/main/AP2-UDK-Projektdokumentation.pdf> (cit. on p. 6).
- Freire, Nuno et al. (2019). “Aggregation of Linked Data in the Cultural Heritage Domain: A Case Study in the Europeana Network”. In: *Information* 10.8. DOI: [10.3390/info10080252](https://doi.org/10.3390/info10080252) (cit. on p. 13).
- Galety, Mohammad Gouse et al., eds. (2022). *Social network analysis: theory and applications*. Hoboken, NJ (cit. on p. 29).
- German National Library (June 2022). *Jahresbericht 2021*. German. Annual report 2021 by the DNB. URL: <https://d-nb.info/1257467816/34> (cit. on pp. 25, 27).
- Gingerich, Owen (2002). *An annotated census of Copernicus’ “De revolutionibus” (Nuremberg, 1543 and Basel, 1566)*. Studia Copernicana. Leiden: Brill. 402 pp. (cit. on p. 3).
- Gruber, Christine and Eveline Wandl-Vogt (2017). “Mapping historical networks: Building the new Austrian Prosopographical / Biographical Information System (APIS). Ein Überblick”. German. In: *Europa baut auf Biographien. Aspekte, Bausteine, Normen und Standards für eine europäische*

- Biographik*. Ed. by Ágoston Bernad, Christine Gruber, and Maximilian Kaiser. new academic press (cit. on p. 10).
- Hakelberg, Dietrich (2016). “Herkunft finden und vernetzen: Stand und Perspektiven der Provenienzerschließung mit Normdaten”. German. Master’s thesis. Humboldt-Universität zu Berlin (cit. on pp. 1, 3, 14, 23, 26, 28).
- Hartig, Olaf (2009). “Provenance Information in the Web of Data”. In: *Proceedings of the Second Workshop on Linked Data on the Web (LDOW 2009)* (cit. on p. 14).
- Hauser, Julia (2014). “Der Linked Data Service der Deutschen Nationalbibliothek”. German. In: *Dialog mit Bibliotheken* 26.1, pp. 38–42 (cit. on p. 13).
- Hengel, Christel and Barbara Pfeifer (2005). “Kooperation der Personennamendatei (PND) mit Wikipedia.” German. In: *Dialog mit Bibliotheken* 17.3, pp. 18–24 (cit. on p. 12).
- Hider, Philip and Ross Harvey (2008). *Organising Knowledge in a Global Society: Principles and Practice in Libraries and Information Centres*. Topics in Australasian Library and Information Studies 29. Wagga Wagga: Centre for Information Studies, Charles Stuart University. URL: <http://ebookcentral.proquest.com/lib/huberlin-ebooks/detail.action?docID=1639577> (cit. on pp. 20–23).
- Horrocks, Ian and Peter F. Patel-Schneider (2011). “KR and Reasoning on the Semantic Web: OWL”. In: *Handbook of Semantic Web Technologies*. Ed. by John Domingue, Dieter Fensel, and James A. Hendler. Berlin, Heidelberg: Springer, pp. 365–398. DOI: [10.1007/978-3-540-92913-0](https://doi.org/10.1007/978-3-540-92913-0) (cit. on p. 11).
- IFLA Study Group on the Functional Requirements for Bibliographic Records (Feb. 2009). *Functional Requirements for Bibliographic Records. Final report*. Ed. by International Federation of Library Associations and Institutions (IFLA). München. URL: <https://repository.ifla.org/handle/123456789/811> (cit. on pp. 2, 20, 31, 50).
- IFLA Working Group on Functional Requirements and Numbering of Authority Records (FRANAR) (Dec. 2008). *Functional Requirements for Authority Data. A Conceptual Model*. Ed. by International Federation of Library Associations and Institutions (IFLA). URL: http://www.ifla.org/files/assets/cataloguing/frad/frad_2013.pdf (cit. on pp. 20, 50).
- IFLA Working Group on Functional Requirements for Subject Authority Records (FRSAR) (June 2010). *Functional Requirements for Subject Authority Data (FRSAD). A Conceptual Model*. Ed. by International Federation of Library Associations and Institutions (IFLA). URL: <https://repository.ifla.org/handle/123456789/835> (cit. on pp. 20, 50).
- Isaac, Antoine, Robina Clayphan, and Bernhard Haslhofer (2012). “Europeana: Moving to Linked Open Data”. In: *Information Standards Quarterly* 24.2/3, p. 34. DOI: [10.3789/isqv24n2-3.2012.06](https://doi.org/10.3789/isqv24n2-3.2012.06) (cit. on p. 13).
- Jansen, Dorothea (2003). *Einführung in die Netzwerkanalyse. Grundlagen, Methoden, Forschungsbeispiele*. German. 2nd. UTB 2241. Opladen: Leske + Budrich. 312 pp. (cit. on p. 5).
- Lígia Triques, Maria, Paula Regina Ventura Amorim Gonçalves, and Ana Cristina de Albuquerque (2022). “Integration of cultural data from digital repositories: an overview of the DPLA Hubs.” In: *Revista Digital de Biblioteconomia e Ciência da Informação* 20, pp. 1–21. DOI: [10.20396/rdbci.v20i00.8666967](https://doi.org/10.20396/rdbci.v20i00.8666967) (cit. on p. 13).
- Malpas, Jeffrey Edward and Hans-Helmuth Gander, eds. (2015). *The Routledge companion to hermeneutics*. Routledge Philosophy Companions. London: Routledge. 778 p (cit. on p. 6).
- Marshall, Catherine C. and Frank M. Shipman (2003). “Which semantic web?” In: *UK Conference on Hypertext* (cit. on p. 11).

- Meiners, Ole (2022). “Evaluation von Named Entity Recognition-Modellen für historische deutschsprachige Texte am Beispiel frühneuzeitlicher Ego-Dokumente”. German. Bachelor’s thesis. Humboldt-Universität zu Berlin. 62 pp. (cit. on p. 8).
- Menzel, Sina, Mark-Jan Bludau, et al. (2020). “Graph Technologies for the Analysis of Historical Social Networks Using Heterogeneous Data Sources”. In: *Graph Technologies in the Humanities 2020 (GRAPH 2020)*. Vol. 3110. CEUR Workshop Proceedings, pp. 124–149 (cit. on pp. 5, 6, 8–10, 51).
- Menzel, Sina, Hannes Schnaitter, et al. (2021). “Named Entity Linking mit Wikidata und GND – Das Potenzial handkuratierter und strukturierter Datenquellen für die semantische Anreicherung von Volltexten”. In: *Qualität in der Inhaltserschließung*. Ed. by Michael Franke-Maier et al. Berlin, Boston: De Gruyter Saur, pp. 229–258. DOI: [doi:10.1515/9783110691597-012](https://doi.org/10.1515/9783110691597-012) (cit. on p. 8).
- Moreau, Luc, Juliana Freire, et al. (2008). “The Open Provenance Model: An Overview”. In: *Provenance and Annotation of Data and Processes*. Ed. by Juliana Freire, David Koop, and Luc Moreau. Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 323–326 (cit. on p. 14).
- Moreau, Luc, Paul Groth, et al. (Apr. 2008). “The Provenance of Electronic Data”. In: *Commun. ACM* 51.4, pp. 52–58. DOI: [10.1145/1330311.1330323](https://doi.org/10.1145/1330311.1330323) (cit. on p. 14).
- Novak, Jasminko et al. (2014). “HistoGraph – A Visualization Tool for Collaborative Analysis of Networks from Historical Social Multimedia Collections”. In: *18th International Conference on Information Visualisation*, pp. 241–250. DOI: [10.1109/IV.2014.47](https://doi.org/10.1109/IV.2014.47) (cit. on p. 10).
- Otte, Evelien and Ronald Rousseau (2002). “Social network analysis: a powerful strategy, also for the information sciences”. In: *Journal of Information Science* 28.6, pp. 441–453. DOI: [10.1177/016555150202800601](https://doi.org/10.1177/016555150202800601) (cit. on p. 5).
- Petras, Vivien and Juliane Stiller (2017). “A Decade of Evaluating Europeana – Constructs, Contexts, Methods & Criteria”. In: *Proceedings of the 21st International Conference on Theory and Practice of Digital Libraries (TPDL)*. Ed. by Jaap Kamps et al. Vol. 10450. Lecture Notes in Computer Science. Springer, pp. 233–245. DOI: [10.1007/978-3-319-67008-9_19](https://doi.org/10.1007/978-3-319-67008-9_19) (cit. on p. 13).
- Pohl, Adrian and Patrick Danowski (2013). “Linked Open Data in der Bibliothekswelt: Grundlagen und Überblick”. German. In: *(Open) Linked Data in Bibliotheken*. Ed. by Patrick Danowski and Adrian Pohl. Berlin, Boston: De Gruyter Saur, pp. 1–44. DOI: [doi:10.1515/9783110278736.1](https://doi.org/10.1515/9783110278736.1) (cit. on p. 12).
- Purschwitz, Anne (Dec. 2018). “Netzwerke des Wissens – Thematische und personelle Relationen innerhalb der halleschen Zeitungen und Zeitschriften der Aufklärungsepoche (1688–1818)”. German. In: *Journal of Historical Network Research* 2.1, pp. 109–142. URL: <http://jhnr.uni.lu/index.php/jhnr/article/view/47> (cit. on p. 10).
- Salatowsky, Sascha and Karl-Heinz Lotze, eds. (2015). *Himmelsspektakel. Ausstellung der Universitäts- und Forschungsbibliothek Erfurt/Gotha*. German. Veröffentlichung der Forschungsbibliothek Gotha 52. Gotha: University und Research Library Erfurt/Gotha. 231 pp. (cit. on p. 3).
- Schneider-Kempf, Barbara et al. (2018). *Beschreibung des Vorhabens – Projektanträge im Bereich “Wissenschaftliche Literaturversorgungs- und Informationssysteme” (LIS)*. Project proposal. URL: <https://github.com/sonar-idh/reports/blob/main/SoNAR-Projektantrag-short.pdf> (cit. on p. 25).
- Ullah, Irfan et al. (Dec. 2018). “An Overview of the Current State of Linked and Open Data in Cataloging”. In: *Information Technology and Libraries (Online)* 37.4, pp. 47–80. DOI: [10.6017/ita1.v37i4.10432](https://doi.org/10.6017/ita1.v37i4.10432) (cit. on p. 13).

- Vorndran, Angela (Dec. 2018). "Hervorholen, was in unseren Daten steckt! Mehrwerte durch Analysen großer Bibliotheksdatenbestände". German. In: *o-bib* 5.4, pp. 166–180. DOI: [10.5282/o-bib/2018H4S166-180](https://doi.org/10.5282/o-bib/2018H4S166-180) (cit. on p. 10).
- Warren, Christopher N. et al. (2016). "Six Degrees of Francis Bacon: A Statistical Method for Reconstructing Large Historical Social Networks". In: *Digital humanities quarterly* 10.3 (cit. on p. 10).
- Wiesenmüller, Heidrun and Silke Horny (2015). *Basiswissen RDA*. German. Berlin, München, Boston: De Gruyter Saur. DOI: [doi:10.1515/9783110311471](https://doi.org/10.1515/9783110311471) (cit. on pp. 20, 21).
- Zuschlag, Christoph (2022). *Einführung in die Provenienzforschung*. German. München: C.H.BECK (cit. on p. 1).

Web Resources

- [Web1] U.S. Department of State. *Washington Conference Principles on Nazi-Confiscated Art*. URL: <https://web.archive.org/web/20170426113213/https://www.state.gov/p/eur/rt/hlcst/270431.htm> (visited on 12/05/2023) (cit. on p. 1).
- [Web2] German National Library. *The Integrated Authority File (GND)*. URL: <https://www.dnb.de/EN/Professionell/Standardisierung/GND/gnd.html> (visited on 30/03/2023) (cit. on pp. 1, 24, 27, 50).
- [Web3] Wikipedia. *Bibliotheksverbund: Deutschland*. German. URL: <https://de.wikipedia.org/wiki/Bibliotheksverbund#Deutschland> (visited on 30/03/2023) (cit. on p. 1).
- [Web4] GBV Common Library Network. *PICA-Format*. German. URL: <https://format.gbv.de/pica> (visited on 30/03/2023) (cit. on p. 1).
- [Web5] Library Service Center Baden-Württemberg and GBV Common Library Network. *Exportformate*. German. URL: <https://wiki.k10plus.de/display/K10PLUS/Exportformate> (visited on 30/03/2023) (cit. on p. 1).
- [Web6] Library of Congress. *Library of Congress Name Authority File (LCNAF)*. URL: <https://id.loc.gov/authorities/names.html> (visited on 19/05/2023) (cit. on pp. 2, 24, 50).
- [Web7] OCLC, Inc. *WorldCat®*. URL: <https://www.worldcat.org/> (visited on 31/03/2023) (cit. on pp. 2, 24).
- [Web8] ISNI International Agency. *ISNI*. URL: <https://isni.org/> (visited on 31/03/2023) (cit. on pp. 2, 24, 50).
- [Web9] OCLC, Inc. *VIAF*. URL: <https://viaf.org/> (visited on 31/03/2023) (cit. on pp. 2, 24, 51).
- [Web10] Wikimedia Foundation. *Wikidata*. URL: https://www.wikidata.org/wiki/Wikidata:Main_Page (visited on 31/03/2023) (cit. on pp. 2, 25, 27).
- [Web11] Wikipedia. *Georg Joachim Rheticus*. URL: https://en.wikipedia.org/wiki/Georg_Joachim_Rheticus#External_links (visited on 31/03/2023) (cit. on p. 2).
- [Web12] Wikidata. *Georg Joachim Rheticus*. Data Set Q93588. URL: <https://www.wikidata.org/wiki/Q93588> (visited on 31/03/2023) (cit. on p. 2).
- [Web13] Research Library Gotha of the University of Erfurt. *Contact persons: Dr. Dietrich Hakelberg*. URL: <https://www.uni-erfurt.de/en/gotha-research-library/library/contact/contact-persons/dr-dietrich-hakelberg> (visited on 18/05/2023) (cit. on p. 2).
- [Web14] Stiftung Preußischer Kulturbesitz: Staatsbibliothek zu Berlin. *Manuscripts and Early Printed Books: Contact*. URL: <https://staatsbibliothek-berlin.de/en/about-the>

- library/departments/manuscripts-and-early-printed-books/contact (visited on 18/05/2023) (cit. on p. 2).
- [Web15] Kompetenzzentrum – Trier Center for Digital Humanities, University of Trier. *Dr Joëlle Weis: Head of Research Area*. URL: <https://tcdh.uni-trier.de/en/person/dr-joelle-weis> (visited on 18/05/2023) (cit. on p. 2).
- [Web16] University of Erfurt. *OPAC search result for French exemplars of English works owned by Luise Dorothea*. URL: [https://opac.uni-erfurt.de/DB=1/SET=5/TTL=1/CMD?ACT=SRCHA&IKT=1016&SRT=YOP&TRM=\(per+milton+or+pope\)+and+slw+provenienz:+luise+dorothea*](https://opac.uni-erfurt.de/DB=1/SET=5/TTL=1/CMD?ACT=SRCHA&IKT=1016&SRT=YOP&TRM=(per+milton+or+pope)+and+slw+provenienz:+luise+dorothea*) (visited on 16/05/2023) (cit. on p. 2).
- [Web17] Wikipedia. *Soziale Netzwerkanalyse*. German. URL: https://de.wikipedia.org/wiki/Soziale_Netzwerkanalyse (visited on 26/04/2023) (cit. on p. 5).
- [Web18] — *Network Science. Network Analysis*. URL: https://en.wikipedia.org/wiki/Network_science#Network_analysis (visited on 26/04/2023) (cit. on p. 5).
- [Web19] University of Applied Sciences Potsdam. *SoNAR (IDH) – Interfaces to Data for Historical Social Network Analysis and Research*. URL: <https://sonar.fh-potsdam.de/> (visited on 20/04/2023) (cit. on pp. 6, 8, 51).
- [Web20] SoNAR (IDH) research collective. *Reports*. URL: <https://github.com/sonar-idh/reports> (visited on 20/04/2023) (cit. on p. 6).
- [Web21] German National Library. *Zeitschriftendatenbank (German Union Catalogue of Serials)*. URL: <https://zdb-katalog.de> (visited on 05/05/2023) (cit. on pp. 7, 24, 51).
- [Web22] — *DNB Catalogue*. URL: <https://katalog.dnb.de/EN/home.html?v=plist> (visited on 31/03/2023) (cit. on pp. 7, 24, 50).
- [Web23] Prussian Cultural Heritage Foundation: Berlin State Library. *Kalliope Union Catalog*. URL: <https://kalliope-verbund.info/> (visited on 05/05/2023) (cit. on pp. 7, 24, 50).
- [Web24] Neo4j, Inc. *neo4j*. URL: <https://neo4j.com/> (visited on 05/05/2023) (cit. on p. 8).
- [Web25] Prussian Cultural Heritage Foundation: Berlin State Library. *ZEFYS – the newspaper information system of Berlin State Library*. URL: <https://zefys.staatsbibliothek-berlin.de/> (visited on 05/05/2023) (cit. on pp. 8, 24, 51).
- [Web26] German National Library. *Digital Exile Press*. URL: https://www.dnb.de/EN/Sammlungen/DEA/Exilpresse/exilpresse_node.html (visited on 05/05/2023) (cit. on pp. 8, 24).
- [Web27] Göttingen State and University Library. *DARIAH-DE*. URL: <https://de.dariah.eu/> (visited on 06/05/2023) (cit. on p. 10).
- [Web28] DARIAH ERIC. *DARIAH-EU*. URL: <https://www.dariah.eu/> (visited on 06/05/2023) (cit. on p. 10).
- [Web29] Wikipedia. *DARIAH-DE*. German. URL: <https://de.wikipedia.org/wiki/DARIAH-DE> (visited on 06/05/2023) (cit. on p. 10).
- [Web30] German National Library. *Culturegraph*. German. URL: Culturegraph.org (visited on 06/05/2023) (cit. on p. 10).
- [Web31] ORCID, Inc. *ORCID. Connecting research and researchers*. URL: <https://orcid.org/> (visited on 06/05/2023) (cit. on p. 10).
- [Web32] Carnegie Mellon University Libraries. *Six Degrees of Francis Bacon*. URL: <http://www.sixdegreesoffrancisbacon.com/> (visited on 06/05/2023) (cit. on p. 10).
- [Web33] CVCE Digital Humanities Lab. *histograph*. URL: <http://histograph.eu/> (visited on 06/05/2023) (cit. on p. 10).

- [Web34] University of Innsbruck. “Issues with Europe”. URL: <https://www.uibk.ac.at/projects/issues-with-europe/index.html.en> (visited on 06/05/2023) (cit. on p. 10).
- [Web35] Austrian Centre for Digital Humanities and Cultural Heritage. *APIS*. URL: <https://www.oeaw.ac.at/acdh/projects/completed-projects/apis> (visited on 06/05/2023) (cit. on p. 10).
- [Web36] W3C. *LinkedData*. URL: <https://www.w3.org/wiki/LinkedData> (visited on 07/05/2023) (cit. on p. 11).
- [Web37] IETF HTTP Working Group. *HTTP Documentation*. URL: <https://httpwg.org/specs/> (visited on 12/05/2023) (cit. on pp. 11, 50).
- [Web38] Schreiber, Guus and Yves Raimond, eds. *RDF 1.1 Primer*. W3C Working Group Note, 24 June 2014. URL: <https://www.w3.org/TR/rdf11-primer/> (visited on 13/05/2023) (cit. on pp. 11, 51).
- [Web39] World Wide Web Consortium. *Resource Description Framework (RDF)*. URL: <https://www.w3.org/RDF/> (visited on 08/04/2023) (cit. on p. 11).
- [Web40] Brickley, Dan and R.V. Guha, eds. *RDF Schema 1.1*. W3C Recommendation, 25 February 2014. URL: <https://www.w3.org/TR/rdf-schema/> (visited on 13/05/2023) (cit. on pp. 12, 51).
- [Web41] Harris, Steve and Andy Seaborne, eds. *SPARQL 1.1 Query Language*. W3C Recommendation, 21 March 2013. URL: <https://www.w3.org/TR/sparql11-query/> (visited on 13/05/2023) (cit. on pp. 12, 51).
- [Web42] International Organization for Standardization (ISO). *ISO/IEC 9075-1:2016. Information technology — Database languages — SQL — Part 1: Framework (SQL/Framework)*. URL: <https://www.iso.org/standard/63555.html> (visited on 13/05/2023) (cit. on p. 12).
- [Web43] Wikipedia. *Linked data*. URL: https://en.wikipedia.org/wiki/Linked_data (visited on 09/05/2023) (cit. on p. 12).
- [Web44] Library of Congress. *Library of Congress*. URL: <https://www.loc.gov/> (visited on 20/05/2023) (cit. on pp. 12, 50).
- [Web45] Gonzalez, Gloria. *Linked Open Data: A Beckoning Paradise*. Guest post in the Library of Congress Blog “The Signal – Digital Happenings at the Library of Congress”. URL: <https://blogs.loc.gov/thesignal/2011/06/linked-open-data-a-beckoning-paradise/> (visited on 14/05/2023) (cit. on p. 13).
- [Web46] Europeana Foundation. *Discover Europe’s digital cultural heritage*. URL: <https://www.europeana.eu> (visited on 14/05/2023) (cit. on pp. 13, 25).
- [Web47] German National Library. *GND Ontology*. URL: <https://d-nb.info/standards/elementset/gnd#> (visited on 14/05/2023) (cit. on p. 13).
- [Web48] GBV Common Library Network. *T-PRO Thesaurus der Provenienzbegriffe*. German. URL: https://provenienz.gbv.de/T-PRO_Thesaurus_der_Provenienzbegriffe (visited on 17/05/2023) (cit. on pp. 14, 51).
- [Web49] University of Erfurt. *OPAC search result for “De Revolutionibus”*. URL: https://opac.uni-erfurt.de/DB=1/CMD?ACT=SRCHA&IKT=1016&SRT=YOP&TRM=tit+de+revolutionibus+and+per+kopernikus+and+jah+15**+and+bbg+a* (visited on 31/03/2023) (cit. on p. 16).
- [Web50] — *OPAC entry for one exemplar of “De Revolutionibus”*. URL: <https://opac.uni-erfurt.de/LNG=EN/DB=1/XMLPRS=N/PPN?PPN=567506266> (visited on 31/03/2023) (cit. on p. 16).

- [Web51] German National Library. *Dataset for the subject term “Wissenschaftler” (scientist) in the GND catalogue*. German. URL: <https://portal.dnb.de/opac/opacPresentation?cqlMode=true&reset=true&referrerPosition=5&referrerResultId=wissenschaftler+and+mat=subjects&any&query=idn=040665674> (visited on 18/05/2023) (cit. on p. 17).
- [Web52] — *List of all subordinate terms for the subject term “Wissenschaftler” (scientist) in the GND catalogue*. German. URL: <https://portal.dnb.de/opac/simpleSearch?reset=true&cqlMode=true&query=ubRef=040665674&selectedCategory=any> (visited on 18/05/2023) (cit. on p. 17).
- [Web53] International Federation of Library Associations and Institutions. *IFLA*. URL: <https://www.ifla.org/> (visited on 23/05/2023) (cit. on pp. 20, 50).
- [Web54] Voss, Jakob. *File:FRBR-Group-1-entities-and-basic-relations.svg*. URL: <https://en.wikipedia.org/wiki/File:FRBR-Group-1-entities-and-basic-relations.svg> (visited on 27/05/2023) (cit. on p. 21).
- [Web55] — *File:FRBR-Group-2-entities-and-relations.svg*. URL: <https://en.wikipedia.org/wiki/File:FRBR-Group-2-entities-and-relations.svg> (visited on 27/05/2023) (cit. on p. 21).
- [Web56] German National Library. *RDA*. URL: https://www.dnb.de/EN/Professionell/Standardisierung/Standards/_content/rda.html?nn=147858 (visited on 21/05/2023) (cit. on pp. 21, 51).
- [Web57] Wikipedia. *Resource Description and Access*. German. URL: https://de.wikipedia.org/wiki/Resource_Description_and_Access (visited on 24/05/2023) (cit. on p. 21).
- [Web58] WHATWG community. *HTML*. URL: <https://html.spec.whatwg.org/multipage/> (visited on 27/05/2023) (cit. on pp. 21, 50).
- [Web59] World Wide Web Consortium. *Extensible Markup Language (XML)*. URL: <https://www.w3.org/XML/> (visited on 26/05/2023) (cit. on pp. 22, 51).
- [Web60] Wikipedia. *Extensible Markup Language*. URL: https://de.wikipedia.org/wiki/Extensible_Markup_Language (visited on 27/05/2023) (cit. on p. 22).
- [Web61] Library of Congress. *SRU – Search/Retrieval via URL*. URL: <https://www.loc.gov/standards/sru/> (visited on 21/05/2023) (cit. on pp. 22, 51).
- [Web62] — *CQL: The Contextual Query Language*. URL: <https://www.loc.gov/standards/sru/cql/index.html> (visited on 21/05/2023) (cit. on pp. 22, 50).
- [Web63] Open Archives Initiative. *Open Archives Initiative*. URL: <https://www.openarchives.org> (visited on 21/05/2023) (cit. on pp. 22, 51).
- [Web64] — *The Open Archives Initiative Protocol for Metadata Harvesting*. URL: <https://www.openarchives.org/OAI/openarchivesprotocol.html> (visited on 21/05/2023) (cit. on pp. 22, 51).
- [Web65] Library of Congress. *MARC Standards*. URL: <https://www.loc.gov/marc/> (visited on 27/05/2023) (cit. on pp. 22, 50).
- [Web66] German National Library. *MAB*. German. URL: https://web.archive.org/web/20181121033620/https://www.dnb.de/DE/Standardisierung/Formate/MAB/mab_node.html (visited on 27/05/2023) (cit. on pp. 23, 50).
- [Web67] Wikipedia. *Maschinelles Austauschformat für Bibliotheken*. URL: https://en.wikipedia.org/wiki/Maschinelles_Austauschformat_f%C3%BCr_Bibliotheken (visited on 27/05/2023) (cit. on p. 23).

- [Web68] Library of Congress. *Metadata Object Description Schema: MODS*. URL: <http://www.loc.gov/standards/mods/> (visited on 27/05/2023) (cit. on pp. 23, 51).
- [Web69] — *Metadata Authority Description Schema (MADS)*. URL: <http://www.loc.gov/standards/mads/> (visited on 27/05/2023) (cit. on pp. 23, 50).
- [Web70] Dublin Core™ Metadata Initiative (DCMI). *DCMI: Dublin Core™*. URL: <https://www.dublincore.org/specifications/dublin-core/> (visited on 25/05/2023) (cit. on pp. 23, 50).
- [Web71] — *DCMI: Home*. URL: <https://www.dublincore.org> (visited on 25/05/2023) (cit. on pp. 23, 50).
- [Web72] Library of Congress. *Overview of the BIBFRAME 2.0 Model*. URL: <https://www.loc.gov/bibframe/docs/bibframe2-model.html> (visited on 28/05/2023) (cit. on pp. 23, 50).
- [Web73] Wikipedia. *BIBFRAME*. URL: <https://en.wikipedia.org/wiki/BIBFRAME> (visited on 28/05/2023) (cit. on p. 23).
- [Web74] Library of Congress. *LC Catalog*. URL: <https://catalog.loc.gov/> (visited on 28/05/2023) (cit. on p. 24).
- [Web75] British Library. *Explore the British Library*. URL: https://explore.bl.uk/primo_library/libweb/action/search.do (visited on 28/05/2023) (cit. on p. 24).
- [Web76] Karlsruhe Institute of Technology (KIT). *KVK – Karlsruhe Virtual Catalog*. URL: <http://kvk.bibliothek.kit.edu/index.html?lang=en> (visited on 20/05/2023) (cit. on pp. 24, 50).
- [Web77] Library Service Center Baden-Württemberg and GBV Common Library Network. *K10plus – Kooperationsprojekt BSZ und GBV*. German. URL: <https://www.bszgbv.de/services/k10plus/> (visited on 30/03/2023) (cit. on pp. 24, 50).
- [Web78] Bibliotheksverbund Bayern. *B3Kat – kurz und bündig*. German. URL: <https://www.bib-bvb.de/web/b3kat> (visited on 20/05/2023) (cit. on pp. 24, 50).
- [Web79] Hochschulbibliothekszentrum des Landes Nordrhein-Westfalen. *Hochschulbibliothekszentrum NRW*. German. URL: <https://www.hbz-nrw.de/> (visited on 20/05/2023) (cit. on pp. 24, 50).
- [Web80] hebis-Verbundzentrale. *hebis*. German. URL: <https://www.hebis.de/> (visited on 20/05/2023) (cit. on pp. 24, 50).
- [Web81] German Lost Art Foundation. *Proveana – Provenance Research Database*. URL: <https://www.proveana.de/en/start> (visited on 20/05/2023) (cit. on p. 25).
- [Web82] German National Library. *Our collection mandate*. URL: https://www.dnb.de/EN/Professionell/Sammeln/sammeln_node.html (visited on 21/05/2023) (cit. on p. 25).
- [Web83] — *Cataloguing media works*. URL: https://www.dnb.de/EN/Professionell/Erschliessen/erschliessen_node.html (visited on 21/05/2023) (cit. on pp. 25, 27).
- [Web84] — *Metadata services*. URL: https://www.dnb.de/EN/Professionell/Metadaten/metadaten_node.html (visited on 21/05/2023) (cit. on p. 25).
- [Web85] — *BIBFRAME – Bibliographic Framework Initiative*. URL: https://www.dnb.de/EN/Professionell/ProjekteKooperationen/Projekte/BIBFRAME/bibframe_node.html (visited on 28/05/2023) (cit. on p. 25).
- [Web86] — *MARC 21*. URL: https://www.dnb.de/EN/Professionell/Metadaten/Exportformate/MARC21/marc21_node.html (visited on 22/05/2023) (cit. on pp. 26, 27).

- [Web87] German National Library. *MARCXML*. URL: https://www.dnb.de/EN/Professionelle11/Standardisierung/Standards/_content/marcxml.html (visited on 22/05/2023) (cit. on p. 26).
- [Web88] — *Metadatenherkunft in der DNB und in MARC 883*. German. URL: <https://wiki.dnb.de/display/ILTIS/Metadatenherkunft+in+der+DNB+und+in+MARC+883> (visited on 22/05/2023) (cit. on p. 26).
- [Web89] Library Service Center Baden-Württemberg and GBV Common Library Network. *About us*. German. URL: <https://www.bszgbv.de/organisation/projekt/> (visited on 23/05/2023) (cit. on p. 26).
- [Web90] GBV Common Library Network. *About the Head Office*. German. URL: https://www.gbv.de/informationen/Verbundzentrale/ueber_die_VZG (visited on 23/05/2023) (cit. on p. 26).
- [Web91] Library Service Center Baden-Württemberg and GBV Common Library Network. *Homepage of the K10plus Wiki*. German. URL: <https://wiki.k10plus.de/> (visited on 23/05/2023) (cit. on p. 26).
- [Web92] Library Service Centre Baden-Württemberg. *K10plus-Recherche*. German. URL: <https://www.bsz-bw.de/K10plus.html> (visited on 23/05/2023) (cit. on p. 26).
- [Web93] GBV Common Library Network. *Holdings statistics of the union catalogue*. German. URL: https://www.gbv.de/informationen/Verbundzentrale/Datenbankstatistik/Datenbankstatistik_1540 (visited on 23/05/2023) (cit. on p. 26).
- [Web94] Library Service Center Baden-Württemberg and GBV Common Library Network. *K10plus Wiki: Open Data*. German. URL: <https://wiki.k10plus.de/display/K10PLUS/Open+Data> (visited on 23/05/2023) (cit. on p. 26).
- [Web95] Universitätsbibliothek Regensburg. *RVK Online*. German. URL: <https://rvk.uni-regensburg.de/regensburger-verbundklassifikation-online> (visited on 29/05/2023) (cit. on pp. 26, 51).
- [Web96] Library Service Center Baden-Württemberg and GBV Common Library Network. *K10plus Wiki: Authority Data*. German. URL: <https://wiki.k10plus.de/display/K10PLUS/Normdaten> (visited on 29/05/2023) (cit. on p. 26).
- [Web97] German National Library. *Erfassungshilfen für Personen und Familien*. German. URL: <https://wiki.dnb.de/pages/viewpage.action?pageId=90411361> (visited on 29/05/2023) (cit. on p. 27).
- [Web98] Wikimedia Foundation. *Wikidata:Introduction*. URL: https://www.wikidata.org/wiki/Wikidata:Main_Page (visited on 30/05/2023) (cit. on pp. 27, 28).
- [Web99] — *Help:Items*. URL: <https://www.wikidata.org/wiki/Help:Items> (visited on 30/05/2023) (cit. on p. 27).
- [Web100] — *Wikidata:Data access*. URL: https://www.wikidata.org/wiki/Wikidata:Data_access (visited on 30/05/2023) (cit. on pp. 27, 28).
- [Web101] — *Help:Statements*. URL: <https://www.wikidata.org/wiki/Help:Statements> (visited on 30/05/2023) (cit. on p. 27).
- [Web102] — *Nicolaus Copernicus (Q619)*. URL: <https://www.wikidata.org/wiki/Q619> (visited on 30/05/2023) (cit. on p. 27).
- [Web103] Wikipedia. *Graph Theory: Applications*. URL: https://en.wikipedia.org/wiki/Graph_theory#Applications (visited on 06/04/2023) (cit. on p. 29).
- [Web104] Python Software Foundation. *Python Graph Libraries*. URL: <https://wiki.python.org/moin/PythonGraphLibraries> (visited on 15/04/2023) (cit. on p. 29).

- [Web105] Naveh, Barak and Contributors. *JGraphT: a Java library of graph theory data structures and algorithms*. URL: <https://jgrapht.org/> (visited on 08/04/2023) (cit. on p. 29).
- [Web106] American National Standards Institute, Inc. (ANSI). *American National Standards Institute – ANSI Home*. URL: <https://www.ansi.org/> (visited on 26/05/2023) (cit. on p. 50).
- [Web107] GBV Common Library Network. *Welcome*. German. URL: <https://en.gbv.de/> (visited on 28/05/2023) (cit. on p. 50).
- [Web108] International Organization for Standardization (ISO). *ISO – International Organization for Standardization*. URL: <https://www.iso.org/> (visited on 26/05/2023) (cit. on p. 50).
- [Web109] National Information Standards Organization. *National Information Standards Organization | NISO Website*. URL: <http://www.niso.org/> (visited on 26/05/2023) (cit. on p. 51).
- [Web110] Hitzler, Pascal et al., eds. *OWL 2 Web Ontology Language Primer (Second Edition)*. W3C Recommendation, 11 December 2012. URL: <https://www.w3.org/TR/owl2-primer/> (visited on 13/05/2023) (cit. on p. 51).
- [Web111] Library Service Centre Baden-Württemberg. *Home – BSZ*. German. URL: <https://www.bsz-bw.de/> (visited on 28/05/2023) (cit. on p. 51).
- [Web112] Library of Congress. *Z39.50 Maintenance Agency Page*. URL: <https://www.loc.gov/z3950/agency/> (visited on 26/05/2023) (cit. on p. 51).

Acronyms

ANSI	American National Standards Institute [Web106]
B3KAT	the union catalogue of the German library networks BVB and KOBV [Web78]
BIBFRAME	Bibliographic Framework [Web72]
BK	Basisklassifikation [Facharbeitsgruppe Sacherschließung 1995]
CQL	Contextual Query Language [Web62]
DC	The Dublin Core Metadata Element Set, short: Dublin Core [Web70]
DCMI	Dublin Core Metadata Initiative [Web71]
DNB	German National Library [Web22]
FRAD	Functional Requirements for Authority Data [IFLA WG FRANAR 2008]
FRBR	Functional Requirements for Bibliographic Records [IFLA SG FRBR 2009]
FRSAD	Functional Requirements for Subject Authority Data [IFLA WG FRSAR 2010]
GBV	Gemeinsamer Bibliotheksverbund, the joint library network of the German states Bremen, Hamburg, Mecklenburg-Vorpommern, Niedersachsen, Sachsen-Anhalt, Schleswig-Holstein, Thüringen, and the Foundation of Prussian Cultural Heritage [Web107]
GND	“Gemeinsame Normdatei”, the Integrated Authority File of the German National Library [Web2]
hbz	library network of the German states North Rhine-Westphalia and Rhineland-Palatinate [Web79]
hebis	library network of the state of Hesse and the region of Rhine Hesse (Germany) [Web80]
HNA	historical network analysis
HTML	Hypertext Markup Language [Web58]
HTTP	Hypertext Transfer Protocol [Web37]
IFLA	International Federation of Library Associations and Institutions [Web53]
ISNI	International Standard Name Identifier [Web8]
ISO	International Organization for Standardization [Web108]
KB	knowledge base
K10plus	the union catalogue of the German library networks GBV and SWB [Web77]
KPE	Kalliope Union Catalogue [Web23]
KVK	Karlsruhe Virtual Catalogue [Web76]
LoC	Library of Congress [Web44]
LCNAF	Library of Congress Name Authority File [Web6]
MAB	“Maschinelles Austauschformat für Bibliotheken” [Web66]
MADS	Metadata Authority Description Schema [Web69]
MARC	Machine-Readable Cataloguing [Web65]

MODS	Metadata Object Description Schema [Web68]
NISO	National Information Standards Organization of the USA [Web109]
OAI	Open Archives Initiative [Web63]
OAI-PMH	OAI Protocol for Metadata Harvesting [Web64]
OWL	the Web Ontology Language recommended by the W3C [Web110]
RDA	Resource Description and Access [Web56]
RDF	Resource Description Framework [Web38]
RDFS	RDF Schema [Web40]
RVK	Regensburger Verbundklassifikation [Web95]
SNA	social network analysis
SoNAR	research project “SoNAR (IDH), Interfaces to Data for Historical Social Network Analysis and Research” [M. Bludau et al. 2020 ; Menzel, M.-J. Bludau, et al. 2020 ; Web19]
SPARQL	Simple Protocol And RDF Query Language [Web41]
S-P-O	subject–predicate–object
SRU	Search/Retrieval via URL [Web61]
SWB	Südwestdeutscher Bibliotheksverbund, the joint library network of the German states Baden-Württemberg, Sachsen, Saarland, and further institutions [Web111]
T-PRO	Thesaurus of Provenance Terms [Web48]
URI	Uniform Resource Identifier [Berners-Lee, Fielding, and Masinter 2005]
VIAF	Virtual International Authority File [Web9]
WWW	World Wide Web
XML	Extensible Markup Language [Web59]
Z39.50	ANSI/NISO Standard “Information Retrieval: Application Service Definition and Protocol Specification” [Web112]
ZDB	German Union Catalogue of Serials [Web21]
ZEFYS	the newspaper information system of Berlin State Library [Web25]

Appendix A

Research Data Management Plan (in German)

list data and link to repository; paste and fill in data management plan)

see also https://www.dfg.de/foerderung/grundlagen_rahmenbedingungen/forschungsdaten/empfehlungen/index.html

A.1 Allgemein

A.1.1 Thema

A.1.1.1 Wie lautet die primäre Forschungsfrage der Abschlussarbeit?

Wie können Provenienzbeziehungen modelliert und maschinell gestützt aufgefunden werden?

A.1.1.2 Bitte geben Sie einige Schlagwörter zur Forschungsfrage bzw. Fragestellung an.

- DDC:⁷
 - 005.72 Datenaufbereitung und Datenrepräsentation
 - 006.332 Wissensrepräsentation
 - 020.0113 Computermodellierung in Bibliotheks- & Informationswissenschaften
- 2012 ACM Computing Classification System:⁸
 - Information systems / Information retrieval / Retrieval tasks and goals / Question answering
 - Information systems / Information retrieval / Retrieval tasks and goals / Information extraction

A.1.1.3 Welchen Regeln oder Richtlinien (HU) zum Umgang mit den in der Abschlussarbeit erhobenen Forschungsdaten folgen Sie für den DMP? Bitte referenzieren Sie diese hier inklusive Version bzw. Veröffentlichungsjahr.

Institut für Bibliotheks- und Informationswissenschaft: Leitlinie zum Umgang mit Forschungsdaten in Abschlussarbeiten. Beschlossen im Institutsrat des IBI am 08.12.2021, in Kraft getreten am 01.02.2022.⁹

⁷<https://deweysearchde.pansoft.de/webdeweysearch/mainClasses.html?catalogs=DNB>

⁸<https://dl.acm.org/ccs>

⁹<https://www.ibi.hu-berlin.de/de/studium/rundumdasstudium/fdm-fuer-studierende>

A.2 Inhaltliche Einordnung

NB: Bitte beschreiben Sie jeden Datensatztyp oder Datensammlung einzeln in dem jeweiligen Kapitel, wo sinnvoll.

A.2.1 Datensatz

A.2.1.1 Um welche Arten von Daten handelt es sich? Bitte in wenigen Zeilen kurz beschreiben.

Für die Literaturstudie und Analyse der Datenquellen wurden Arbeiten aus der Literatur und Webseiten gesammelt.

A.2.2 Datenursprung

A.2.2.1 Werden die Daten selbst erzeugt oder nachgenutzt?

Die Daten wurden selbst erhoben.

A.2.2.2 Wenn die Daten nachgenutzt werden, wer hat die Daten erzeugt? Bitte mit Angabe des Identifiers, falls vorhanden, z.B. DOI.

Trifft nicht zu; siehe oben.

A.2.3 Reproduzierbarkeit

A.2.3.1 Sind die Daten reproduzierbar, d. h. ließen sie sich, wenn sie verloren gingen, erneut erstellen oder erheben?

Die Daten sind nur bedingt reproduzierbar: Anhand des Literaturverzeichnisses der Arbeit lassen sich alle Quellen wiederfinden, aber bei den meisten Webseiten besteht die Gefahr, dass der Inhalt sich ändert oder die Webseite inaktiv wird. Letzteres kann auch bei denjenigen der wissenschaftlichen Arbeiten passieren, die nur online verfügbar sind.

A.2.4 Nachnutzung

A.2.4.1 Für welche Personen, Gruppen oder Institutionen könnte dieser Datensatz (für die Nachnutzung) von Interesse sein? Für welche Szenarien ist dies denkbar?

*** TBC ***

finish DMP



Name: Schneider Vorname: Thomas

Matr.Nr.: 624025

Eidesstattliche Erklärung zur

- ☐ **Hausarbeit ***
☐ **Bachelorarbeit ***
☒ **Masterarbeit ***
☐ **Abschlussarbeit im Bibliotheksreferendariat ***

* Die eingereichte PDF-Datei ist mit den Printexemplaren identisch.

Ich erkläre ausdrücklich, dass es sich bei der von mir eingereichten schriftlichen Arbeit mit dem Titel

Modelling and Automated Retrieval of Provenance Relationships

.....
.....

um eine von mir erstmalig, selbstständig und ohne fremde Hilfe verfasste Arbeit handelt.

Ich erkläre ausdrücklich, dass ich *sämtliche* in der oben genannten Arbeit verwendeten fremden Quellen, auch aus dem Internet (einschließlich Tabellen, Grafiken u. Ä.) als solche kenntlich gemacht habe. Insbesondere bestätige ich, dass ich ausnahmslos sowohl bei wörtlich übernommenen Aussagen bzw. unverändert übernommenen Tabellen, Grafiken u. Ä. (Zitaten) als auch bei in eigenen Worten wiedergegebenen Aussagen bzw. von mir abgewandelten Tabellen, Grafiken u. Ä. anderer Autorinnen und Autoren (Paraphrasen) die Quelle angegeben habe.

Mir ist bewusst, dass Verstöße gegen die Grundsätze der Selbstständigkeit als Täuschung betrachtet und entsprechend der Prüfungsordnung und/oder der Fächerübergreifenden Satzung zur Regelung von Zulassung, Studium und Prüfung (ZSP-HU) geahndet werden.

Datum

Unterschrift