

# An Evaluation Framework for LD Solutions

a case study at TU Wien

BACHELORARBEIT

zur Erlangung des akademischen Grades

**Bachelor of Science**

im Rahmen des Studiums

**Software und Information Engineering**

eingereicht von

**Lukas Baronyai**

Matrikelnummer 01326526

an der Fakultät für Informatik

der Technischen Universität Wien

Betreuung: Prof. Dr. Stefan Biffl

Mitwirkung: MSc., PhD Marta Sabou

Wien, 11. September 2017

---

Lukas Baronyai

---

Stefan Biffl



# An Evaluation Framework for LD Solutions

a case study at TU Wien

BACHELOR'S THESIS

submitted in partial fulfillment of the requirements for the degree of

**Bachelor of Science**

in

**Software and Information Engineering**

by

**Lukas Baronyai**

Registration Number 01326526

to the Faculty of Informatics

at the TU Wien

Advisor: Prof. Dr. Stefan Biffl

Assistance: MSc., PhD Marta Sabou

Vienna, 11<sup>th</sup> September, 2017

---

Lukas Baronyai

---

Stefan Biffl



# Erklärung zur Verfassung der Arbeit

Lukas Baronyai  
Huttengasse 51/10, 1160 Wien

Hiermit erkläre ich, dass ich diese Arbeit selbständig verfasst habe, dass ich die verwendeten Quellen und Hilfsmittel vollständig angegeben habe und dass ich die Stellen der Arbeit – einschließlich Tabellen, Karten und Abbildungen –, die anderen Werken oder dem Internet im Wortlaut oder dem Sinn nach entnommen sind, auf jeden Fall unter Angabe der Quelle als Entlehnung kenntlich gemacht habe.

Wien, 11. September 2017

---

Lukas Baronyai



# Abstract

Along with the rising popularity of the Semantic Web and Linked Open Data, there are many tools and frameworks present which support building LD applications. It is challenging to choose an approach among the many available options without a guide of comparison. This is especially true for TU Wien, since the university has to fulfill the needs of very different stakeholders at the same time while trying to integrate a LD solution into existing structures and workflow. This paper aims to give an overview of various LD tools and frameworks and compare them among each other to give stakeholders at TU Wien a guideline for a future LD project at the university.

In order to do that, an evaluation framework based on four families of comparison criteria was designed. It was then validated in two ways: First, it was applied to compare 8 LD solutions (Euclid project, LUCERO, Linked Data book, D2RQ Platform, Information Workbench, LDIF, Eclipse RDF4J and Apache Jena), which were found by conducting a literature study. Second, based on this comparison, a suitable solution was recommended to the use case of TU Wien. Our use of the framework in the TU Wien context, indicated that it makes the selection of an LD candidate well-informed and much faster.





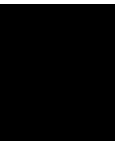
# Contents

<b>Abstract</b>	<b>vii</b>
<b>Contents</b>	<b>ix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Research Question . . . . .	1
1.2 Methodology . . . . .	2
1.3 Structure Of This Thesis . . . . .	2
<b>2 State Of The Art</b>	<b>3</b>
2.1 Benchmarking Database Representations of RDF/S Stores . . . . .	3
2.2 Ontology Alignment Evaluation Initiative (OAEL) . . . . .	3
2.3 Quality Assessment for Linked Data: A Survey . . . . .	4
2.4 An Evaluation Framework For Spreadsheet Tools Comparison . . . . .	4
2.5 The Berlin SPARQL Benchmark . . . . .	5
2.6 Benchmarking Fulltext Search Performance of RDF Stores . . . . .	6
<b>3 Methodology (RQ1)</b>	<b>7</b>
3.1 About The Difficulty Of Comparing Solutions . . . . .	7
3.1.1 The Term "Framework" . . . . .	7
3.1.2 Defining The Limits . . . . .	8
3.2 Methodology Of The Literature Study . . . . .	9
3.3 A Classification System . . . . .	9
<b>4 Overview Of LD Solutions (RQ1)</b>	<b>11</b>
4.1 Architectures Of Frameworks . . . . .	11
4.1.1 Euclid Project . . . . .	11
4.1.2 LUCERO . . . . .	13
4.1.3 Linked Data book . . . . .	15
4.2 Frameworks . . . . .	17
4.2.1 D2RQ Platform . . . . .	17
4.2.2 Information Workbench . . . . .	18
4.2.3 LDIF - Linked Data Integration Framework . . . . .	20
4.2.4 Eclipse RDF4J (formerly Sesame) . . . . .	22
	ix

4.2.5	Apache Jena . . . . .	23
4.3	Excluded Tools And Projects . . . . .	24
4.3.1	LD-Patterns . . . . .	24
4.3.2	LOD2 Stack . . . . .	24
4.3.3	LODUM . . . . .	24
4.3.4	Synth and SHDM . . . . .	24
<b>5</b>	<b>An Evaluation Framework for LD solutions (RQ2)</b>	<b>27</b>
5.1	Criteria from Previous Study . . . . .	27
5.2	Usability Criteria . . . . .	28
5.3	Data Formats . . . . .	28
5.4	Linked Data Publishing Checklist . . . . .	28
5.5	Evaluation framework . . . . .	29
<b>6</b>	<b>Comparison Of LD Solutions Based On The Evaluation Framework (RQ3)</b>	<b>31</b>
6.1	Classification . . . . .	32
6.2	Criteria Group 1: Criteria from Previous Study . . . . .	32
6.2.1	Summary Group 1 . . . . .	34
6.3	Criteria Group 2: Usability . . . . .	34
6.3.1	Summary Group 2 . . . . .	35
6.4	Criteria Group 3: Data Formats . . . . .	35
6.4.1	Summary Group 3 . . . . .	37
6.5	Criteria Group 4: Linked Data Publishing Checklist . . . . .	37
6.5.1	Summary Group 4 . . . . .	38
6.6	Summary . . . . .	39
<b>7</b>	<b>A Case Study Of Applying The Evaluation Framework At TU Wien (RQ4)</b>	<b>41</b>
7.1	General Considerations . . . . .	41
7.2	Situations & Requirements . . . . .	42
7.2.1	Stakeholders . . . . .	42
7.2.2	Situations . . . . .	42
7.3	Proposed Solutions . . . . .	43
7.3.1	Scenario 1: Specialized Single Solution . . . . .	43
7.3.2	Scenario 2: Function-rich Platform . . . . .	44
7.3.3	Scenario 3: Complete Controlled Platform . . . . .	44
7.4	Summary . . . . .	44
<b>8</b>	<b>Conclusion And Future Work</b>	<b>47</b>
8.1	Conclusion . . . . .	47
8.1.1	Existing LD solutions . . . . .	48
8.1.2	The Evaluation Framework For LD Solutions . . . . .	48
8.1.3	Application Of the Proposed Evaluation Framework . . . . .	48

8.1.4 Solution For TU Wien . . . . .	49
8.2 Future Work . . . . .	49
<b>9 Appendix</b>	<b>51</b>
List of Figures	55
List of Tables	55
References To Refereed Scientific Work	57
References To Non-Refereed Work	59
References To Websites	61





# Introduction

The Semantic Web is getting more and more popular and with it the need for ways to publish Linked (Open) Data. In order to cover these needs, a lot of different tools, frameworks and solutions were developed, some covering the whole process of publishing LD, others covering only steps in the process of it. These led to a vast diversity but also a big number of options. In order to develop a LD project, the responsible stakeholders now have to choose among these many options, which is tedious and time-consuming process]. There is a clear lack for a comparison framework of existing LD solutions. .

## 1.1 Research Question

Aim of this thesis is to compare existing and common LD solutions to give responsible stakeholders at TU Wien a decision guidance for choosing one. The concrete research question is as follows:

**RQ:** *What is a suitable evaluation framework for comparing LD solutions?*

1. **RQ1:** What are the most popular LD Solutions?
2. **RQ2:** What are criteria to compare LD solutions? How can an Evaluation Framework for comparing LD solutions look like?
3. **RQ3:** Is the evaluation framework suitable for comparing LD solutions?
4. **RQ4:** How can the Evaluation Framework and the results of the comparison help as guideline at TU Wien?

The conducted work is a follow-up work of a previous study (see [10]) done by the author and will build up on it.

.

### 1.2 Methodology

The work was done as a literature study. First a discussion of the term "framework" was done in order to achieve the research question. Then a range of existing LD projects and application was investigated, to extract used technologies from them. From this, candidates were retrieved and classified. After excluding and filtering some of the candidates, they were compared in four criteria groups: Criteria from the above mentioned study, usability, data formats and the Linked Data Publishing Checklist.

### 1.3 Structure Of This Thesis

This thesis starts with a state of the art section in which similar comparison will be described. The next chapter 3 is the definition of the used methodology, where a discussion of the term "framework" will be presented for this thesis (section 3.1), the process of the literature study (section 3.2) and the classification (section 3.3) system will be defined.

In chapter 4 the found solution will be reviewed and described as well as the excluded candidates (section 4.3). In chapter 5 the criteria as well as their according scala will be defined. The final comparison will be done in chapter 6, the summary of it can be found in section 6.6.

The last research question, RQ4, will be answered in chapter 7, investigating, which of the proposed solutions may be suitable for which situation at TU Wien.

The last chapter, 8, describes the overall summary and future work.

---

<sup>0</sup>**Note:** The term *solutions* is used here on purpose to abstract terms like *framework*, *tool* or *all-in-one solution*. Using only the term *framework* would lead to problems and cut out other options. For a more detailed discussion see section 3.1

## State Of The Art

Since there was no similar comparison about LD frameworks found, this section will cover general comparisons of frameworks in the field of semantic web.

### 2.1 Benchmarking Database Representations of RDF/S Stores

Theoharis et. al. conducted a benchmark in their paper *Benchmarking Database Representations of RDF/S Stores* ([11]) in order to compare different kind of RDF stores. They examined three popular database representations of RDF/S schemata and data: schema-aware, with explicit (ISA) or implicit (NOISA) storage of subsumption relationships, schema-oblivious using (ID) or not (URI) identifiers to represent resources and hybrid of the schema-aware and schema-oblivious representations. Further, they benchmarked two common approaches for evaluating taxonomic queries : on-the-fly (ISA, NOISA, Hybrid) and by precomputing the transitive closure of subsumption relationships (MatView, URI, ID).

"The main conclusion drawn is that the evaluation of taxonomic queries is most efficient over RDF/S stores utilizing the Hybrid and MatView representations. Of the rest, schema-aware representations (ISA, NOISA) exhibit overall better performance than URI, which is superior to that of ID, which exhibits the overall worst performance." ([11])

### 2.2 Ontology Alignment Evaluation Initiative (OAEI)

Since 2004 the initiative organizes every year evaluations of various ontology matching systems. The official result are published on their website <sup>1</sup>. There are also a various

---

<sup>1</sup><http://oaei.ontologymatching.org/>

kind of papers available, analyzing the outcome of the evaluations, the last ones are from 2016 ([12]), 2015 ([13]) and 2014 ([14]). The results of 2017 are not available at the point of writing this thesis.

In order to evaluate a benchmark, the initiative provides a set of test cases, ranging from real world cases to generated cases. Participants register their tools themselves, which then will be run by the organizers with both blind and published datasets. The evaluation then uses the SEALS infrastructure since 2011 and is since 2006 usually reported at the Ontology matching workshop of the International Semantic Web Conference.

### 2.3 Quality Assessment for Linked Data: A Survey

Aim of Zaeri et. al in their survey [1] was a systematic review of approaches for assessing the quality of LD, since they observed "widely varying data quality ranging from extensively curated datasets to crowdsourced and extracted data of relatively low quality". Their goal was to obtain an understanding of the differences of existing approaches. They gathered these approaches and analyzed them qualitatively, providing a list of 18 quality dimensions (grouped by accessibility, intrinsic, contextual and representational) and 69 metrics. Furthermore, they analysed 30 core approaches and 12 tools using a set of attributes.

As result, they observed that "this research area is still in its infancy and can benefit from the possible re-use of research from mature, related domains. Additionally, in most of the surveyed literature, the metrics were often not explicitly defined or did not consist of precise statistical measures."

### 2.4 An Evaluation Framework For Spreadsheet Tools Comparison

In the paper *Towards Evaluation and Comparison of Tools for Ontology Population from Spreadsheet Data* by Kovalenko et al ([15]) an evaluation framework was proposed by the authors in order to facilitate tools comparison, mainly for Ontology Population from Spreadsheet Data but also applicable for other comparisons. They analyzed different types of end users (Semantic Web Experts, Domain Experts and Software Developers) and their needs as well as performed a literature analysis on Tools for Ontology Population from Spreadsheet Data. The result are the following criteria:

- **(C1) General Information** with *Maturity*, *License* and *Type* as sub criteria
- **(C2) Usability** with *GUI* and *Required User Knowledge* as sub criteria
- **(C3) Input format**
- **(C4) Output format**



- (C5) **Mappings** with *Internal Representation*, *External Representation* and *Storage* as sub criteria
- (C6) **Expressiveness**
- (C7) **Multi-user support**
- (C8) **Required software**
- (C9) **Additional features**

This criteria can be used not only for comparison but also for classification of tools as well as for a search for an appropriate tool for a specific task/problem. In order to prove the usefulness, the authors performed a qualitative analysis and comparison of a representative set of seven tools (**RDF123**, **XLWrap**, **Mapping Master**, **Populous**, **Anzo for Excel**, **TopBraid**, **Google Refine**).

## 2.5 The Berlin SPARQL Benchmark

In the article *The Berlin SPARQL Benchmark* (BSBM) by Bizer and Schultz ([2]) the authors introduced a benchmark for comparing performance of native RDF stores with the performances of SPARQL-to-SQL rewriters across architectures. The benchmark query simulates the search and navigation pattern of a consumer looking for a product, representing an e-commerce use case. The article presents, next to the discussion about the benchmark itself, results of the benchmark including for popular RDF stores (**Sesame**, **Virtuoso**, **Jena TDB** and **Jena SDB**) as well as two SPARQL-to-SQL rewriters (**D2R Server** and **Virtuoso RDF Views**) and relational database management systems (**MySQL** and **Virtuoso RDBMS**). An excerpt of the result can be seen in figure 2.1.

Table 6. Load times for different stores and dataset sizes (in [day:]hh:min:sec)

	1M	25M	100M
<b>Sesame</b>	00:02:59	12:17:05	3:06:27:35
<b>Jena TDB</b>	00:00:49	00:16:53	01:34:14
<b>Jena SDB</b>	00:02:09	04:04:38	1:14:53:08
<b>Virtuoso TS</b>	00:00:23	00:39:24	07:56:47
<b>Virtuoso RV</b>	00:00:34	00:17:15	01:03:53
<b>D2R Server</b>	00:00:06	00:02:03	00:11:45
<b>MySQL</b>	00:00:06	00:02:03	00:11:45
<b>Virtuoso SQL</b>	00:00:34	00:17:15	01:03:53

Table 7. Speed-up factors between the runtime of the second query mix and the average runtime of a query mix in steady-state

	1M	25M	100M
<b>Sesame</b>	15,61	3,98	0,75
<b>Jena TDB</b>	3,03	0,52	0,00
<b>Jena SDB</b>	0,97	2,68	0,64
<b>Virtuoso TS</b>	0,47	26,14	46,65
<b>Virtuoso RV</b>	0,15	1,98	100,09
<b>D2R Server</b>	0,67	0,03	0,04
<b>MySQL</b>	26,30	17,37	8,49
<b>Virtuoso SQL</b>	1,03	13,58	247,20

Figure 2.1: Excerpt of the Berlin SPARQL Benchmark results

## 2.6 Benchmarking Fulltext Search Performance of RDF Stores

Another benchmark in the area of semantic web (next to many others) is the benchmark proposed by Minack et al in their paper *Benchmarking Fulltext Search Performance of RDF Stores* ([3]). They extended the existing LUBM benchmark with "synthetic scalable fulltext data and corresponding queries for fulltext-related query performance evaluation". The benchmark was then applied to four popular RDF stores (**Jena**, **Sesame2**, **Virtuoso**, and **YARS**). An excerpt of the result can be seen in figure 2.2.

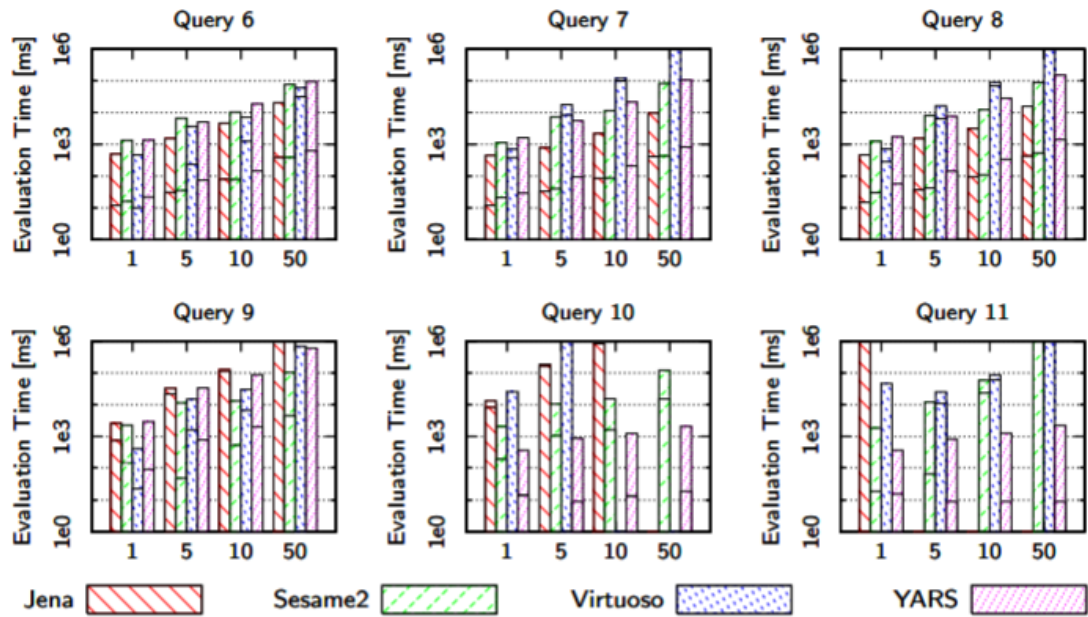


Fig. 4. Semantic IR queries combining fulltext search with structured queries

Figure 2.2: Excerpt of the Berlin SPARQL Benchmark results

## Methodology (RQ1)

### 3.1 About The Difficulty Of Comparing Solutions

One could expect, that this thesis may compare "frameworks". But there comes problems with this term. What actually is a "framework"? Is the term well-defined enough to clearly refer to a solution?

#### 3.1.1 The Term "Framework"

In order to use the term a first step must be to define *what* a framework actually is, since it is a very generic term. One way to define it could be the definition by Roberts and Johnson, [4]:

Frameworks are reusable designs of all part of a software described by a set of abstract classes and the way instances of those collaborate

Another way could be the explanation by Riehle in his PhD thesis, [16]:

Frameworks model a specific domain or an important aspect thereof. They represent the domain as an abstract design, consisting of abstract classes (or interfaces). The abstract design is more than a set of classes, because it defines how instances of the classes are allowed to collaborate with each other at runtime. Effectively, it acts as a skeleton, or a scaffolding, that determines how framework objects relate to each other.

A framework comes with reusable implementations in the form of abstract and concrete class implementations. Abstract implementations are abstract classes that implement parts of a framework abstraction (as expressed by

an abstract class or interface), but leave crucial implementation decisions to subclasses. [...]

Both of them refer frameworks as tools for coding, used when writing own applications. One of the most classical examples might be the Spring Framework in the Java world. In the mentioned projects, Apache Jena and RDF4J mostly apply on this definition.

But the problem is here, that the term is not always used and understood in this way, LDIF and the Silk frameworks define themselves as such, but providing in fact a set of tools without necessarily needing coding to work with them (except configuration files). Others may see tools like the Information Workbench or D2RQ as a framework for publishing.

On a higher level, the architectures proposed in section 4.1 might be seen as a high-level or meta-framework. And since the proposed tools in the other sections (partially) are using these architectures, one could argue, that they therefore are also frameworks.

It can be seen, that is actually problematic to use this term in the context of this thesis since it is too broad and not well-defined. In order to solve this issue, the term "solution" will be used to abstract the term and referring either to "tool", "framework" or "all-in-one solution".

#### 3.1.2 Defining The Limits

Next to the general problem about the term "framework", another problem is to set the borders of the examined topic. As mentioned in the introduction this thesis aims to compare LD solutions, the goal is to cover the whole process of publishing Linked Data, from the bottom persistence layer of accessing existing data, transforming data formats (e.g. relational to RDF), over cleaning and interlinking the data, over storing them in a triple store, up to making them available over an interface like SPARQL.

But there are not many tools/frameworks covering the whole process and supporting different data formats (e.g. relational data and CSV) at the same time. There are some tools like D2RQ only focusing on specific data formats, but providing the full stack, some tools like LDIF only focusing on a specific part of the process, without e.g. providing capabilities for SPARQL endpoints.

The best way is maybe using a stack of different tools to cover the whole workflow, combining them like Silk is integrated in LDIF. Or using the generic architecture, coding an own application and using partially the proposed tools.

But covering different areas, it is difficult to actually compare them. How to compare a persistence framework with a GUI framework? In order to handle this problem, a classification system will be introduced in section 3.3 and the criteria introduced in section 5 will take this difficulty into account.

## 3.2 Methodology Of The Literature Study

In order to answer RQ1, a literature study was conducted, but since the scope is difficult to define as described in section 3.1, it was not a pure study. As the aim of this thesis is to compare *common* frameworks and *best practices*, it would be not sufficient to review every possible paper about a LD framework or tool, therefore another approach was chosen: deriving candidates from projects. In order to do that, the following process was used:

1. Identify & find a LD application/project, ignoring the success of it
2. Find public documentation and/or scientific work of it
3. Analyse used technology, add as candidate if appropriate and if not disadvised
4. Classify candidates (see 3.3)
5. Analyse reference work for possible input for 1.)
6. Analyse reference work of tool/framework at its documentation

Using this approach led to a variety of candidates, which will be listed in section 4. The candidates from section 4.3 were mostly excluded because of step 2.), which ensured a better base for the following comparison.

## 3.3 A Classification System

Resulting from 3.1 different classifications were introduced to find classification-based criteria and to balance out the vast variation of the results. The classifications can be seen in table 3.1.

The Classifications are based on the idea, that a majority of applications are using in one way or another a variation or parts of the three layer architecture style, with components responsible for either UI, Business or Data Access. This does *not* necessarily mean, that they use the full concepts of this architecture or even implementing this style. It is only assumed that a component have a responsibility mappable to one of the layers. Accordingly it is assumed, that a solution can be associated with one of these responsibilities.

It is arguable, if the differentiation between "Full-Stack" and labeling a solution with the three layer class is necessary. The additional "Full-Stack" class was added to emphasize the "All-In-One" approach of such a solution, meaning that all components are provided, no further components are need. This also means, that the included components of the different responsibilities are either harmonized to each other or do not differentiate between these responsibilities. On the other side labeling a tool with the three layer classes, does not implicit this and can also mean, that the support of each of this layer can be optional.

<b>Class</b>	<b>Detail</b>
<b>Architecture</b>	A general architecture without concrete technology. A solution of this class can be used in combination of any other class.
<b>Full-Stack</b>	A solution which covers the whole stack and therefore does not need another component. An "All-In-One" Solution
<b>Presentation layer</b>	A solution which only covers the presentation or UI layer and therefore depends on other component. Managing how LD can be accessed from outside and how the data are exposed.
<b>Business Layer</b>	A solution which only covers the business layer and therefore depends on other components. Managing how LD are processed.
<b>Data Access Layer</b>	A solution which only covers the data access layer and therefore depends on other components. Managing how LD are stored and accessed by the application.

**Table 3.1:** Overview of the Classification

# Overview Of LD Solutions (RQ1)

In order to compare solutions an understanding of existing solutions is necessary. This section will look at existing solutions, what kind of solutions they are, which of them can be used for this thesis and which must be excluded. Furthermore, this section aims to understand how solutions look like and will examine the architecture of them.

## 4.1 Architectures Of Frameworks

In this subsection the thesis will look into three proposed models how solutions (and/or implemented LD-applications) should look like. There are many other existing architectures and ongoing projects exposing data as Linked (Open) Data, this thesis will use the following as representation of them.

### 4.1.1 Euclid Project



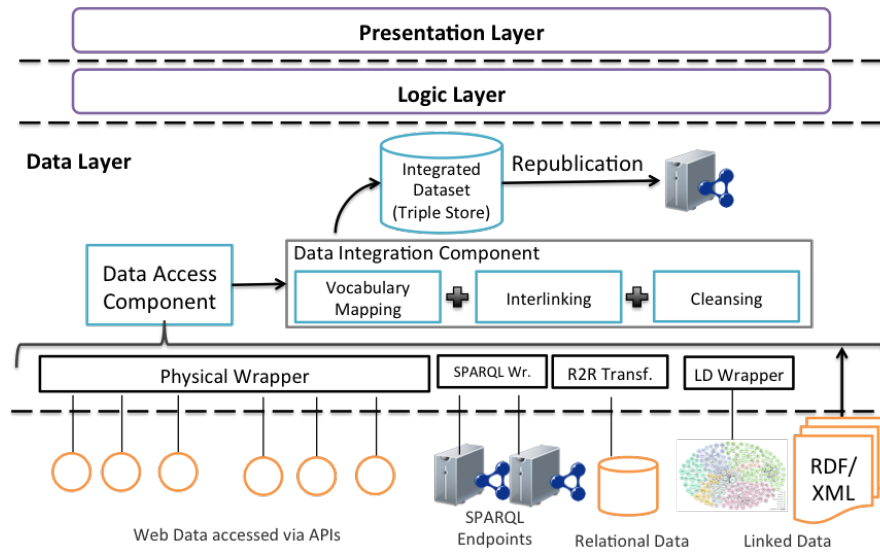
The EUCLID project <sup>1</sup> (EdUcational Curriculum for the usage of Linked Data) was founded under the *Seventh Framework Program of Research and Technological Development*, a funding program of the European Union/European Commission for 2007-2013 <sup>2,3</sup>.

---

<sup>1</sup> [27]

<sup>2</sup> [28]

<sup>3</sup> EUCLID in the CORDIS database: [http://cordis.europa.eu/project/rcn/103709\\_en.html](http://cordis.europa.eu/project/rcn/103709_en.html)



**Figure 4.1:** General EUCLID architecture

Aim of the project was (and still is) to gather existing knowledge and expertise of *"researchers, technology enthusiasts and early adopters in various European Member States"* and provide that accumulated as educational resources to enable the full benefit of L(O)D for European businesses. The project built upon a consortium experienced in *"over 20 LD projects with over 40 companies and public offices in more than 10 countries"* [29]

The outcome of this project is a a range of learning materials, fragmented into modules, and eLearning distribution channels. Overall there are six modules:

1. **Introduction and Application Scenarios** The introduction provides the knowledge to understand, *what* Linked Data are, the main principles, the standards and the required technologies. Further, an overview how to publish and to consume the data is given.
2. **Querying Linked Data** This chapter mainly describes SPARQL and how to use it for querying and updating.
3. **Providing Linked Data** This module deals with the production and exposure of Linked data, using the tools as R2RML (for relational databases), Open Refine (for spreadsheets), GATECloud (for natural language) and Silk (for interlinkage between datasets, see section 4.2.3 for details about this tool).
4. **Interaction with Linked Data** The projects describes in this chapter, how to explore Linked Data, using visualization tools, semantic browsers and applications, introducing search options like faceted search, concept-based search and hybrid search.



5. **Creating Linked Data Applications** This module describes how to build a Linked Data Application, which technologies to use and how to integrate common Web APIs.
6. **Scaling up** Finally this chapter examines the main issues of scalability regarding Linked Open Data and describes the relationship to Big Data.

For this thesis module 3 and 5<sup>4</sup> are the most interesting. Module 3 describes some useful technologies for various steps on the way of exposing LD, but module 5 introduce a high level architecture and some patterns, how a LD application might look like (see [30] for details). In detail, they provide a three-tier architecture (see figure 4.1 and three architecture patterns).

The architecture is very generic and consists of the classic three tiers: presentation, logic and data, each independent to the overlaying tier. Since the presentation and logic layer does not concern the actual publishing of the data, the data layer is the interesting one here. The layer consists of the *Data Access Component*, which represents the access to different data types like relational data or other Web APIs and transforms the data to RDF, the *Data Integration Component*, which does the vocabulary mapping and interlinking for the cleansing in order to e.g. identify and fix ambiguities in resource names, and finally the *Triple Store*, holding the integrated dataset for exposing it to the web.

The mentioned patterns to use for implementations are:

- **Crawling pattern** Used for loading the data in advance and storing them in a triple store, increasing the efficiency of data access. In exchange, the data might not be up to date when accessed.
- **On-The-Fly Dereferencing Pattern** Meaning that the URIs are dereferenced when the application need to access the data. This pattern provides up to date data but for the cost of performance when dereferencing many URIs.
- **(Federated) Query Pattern** Describing the use of complex queries on a fix set of data sources, enabling to work with current data directly retrieved from the sources. The pattern offers an access up-to-date data with adequate response time in specific situations but for the cost of the complex problem to find optimal queries.

#### 4.1.2 LUCERO

The LUCERO project ("Linking University Content for Education and Research Online")<sup>5</sup> was a project at the Open University, aiming to "*scope, prototype, pilot and evaluate*

---

<sup>4</sup> [30]

<sup>5</sup>The code is available in the Google Code Archive: <https://code.google.com/archive/p/lucero-project/wikis/StepByStepDocumentation.wiki>

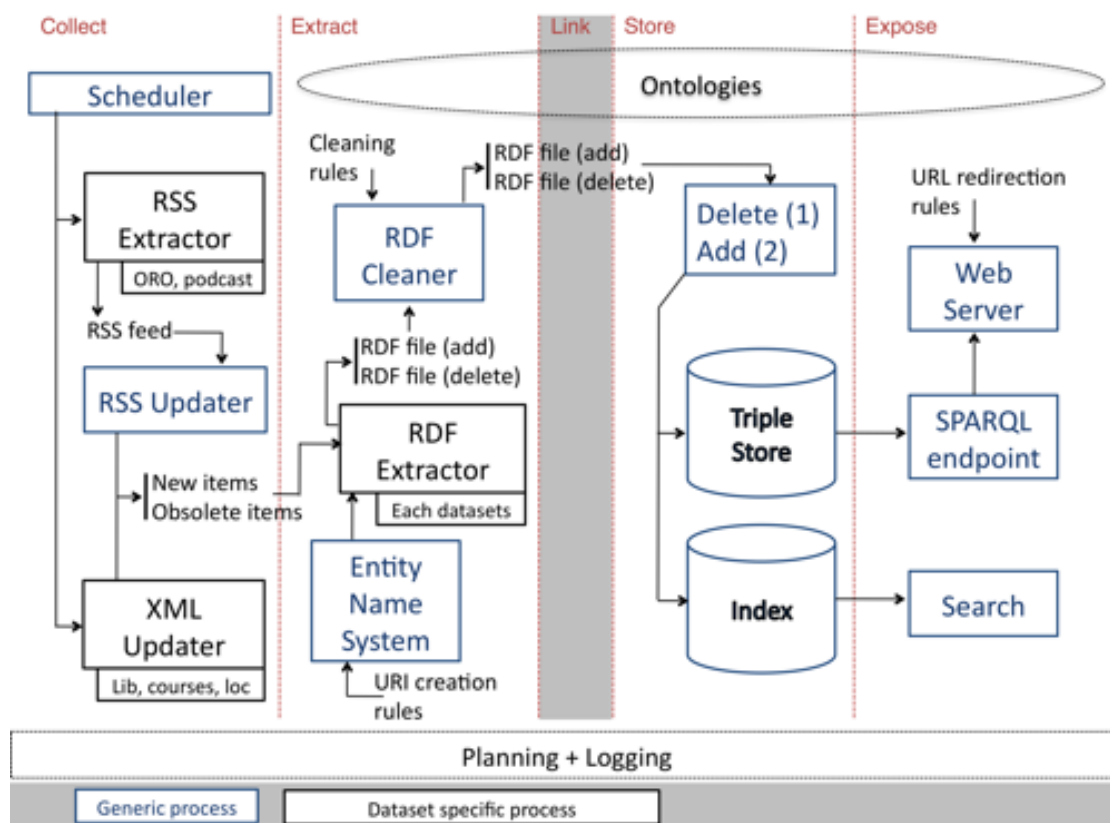


Figure 4.2: LUCERO work flow &amp; architecture

reusable, cost-effective solutions relying on the linked data principles and technologies for exposing and connecting educational and research content". It was founded for one year by the JISC Information Environment 2011 Program under the call Deposit of research outputs and Exposing digital content for education and research. [31]

The projects connected with other organizations through LinkedUniversities.org<sup>6</sup> to gather common issues and practices. The outcome was the first university linked data platform, <http://data.open.ac.uk/>, with a lot of impact on The Open University and the education community.

Looking at the architecture in figure 4.2 comparing to the Euclid architecture seen in the previous section, there are quite a lot of similarities. Both have components for accessing different kinds of data, here called *Extractors*, for cleaning the data, here called *Cleaner*, and a Triple Store, holding the data available. The lanes "Collect", "Extract", "Link" and "Store" can be seen as the data layer from the classic three-tier architecture, the "Expose" lane as the logic and presentation layer.

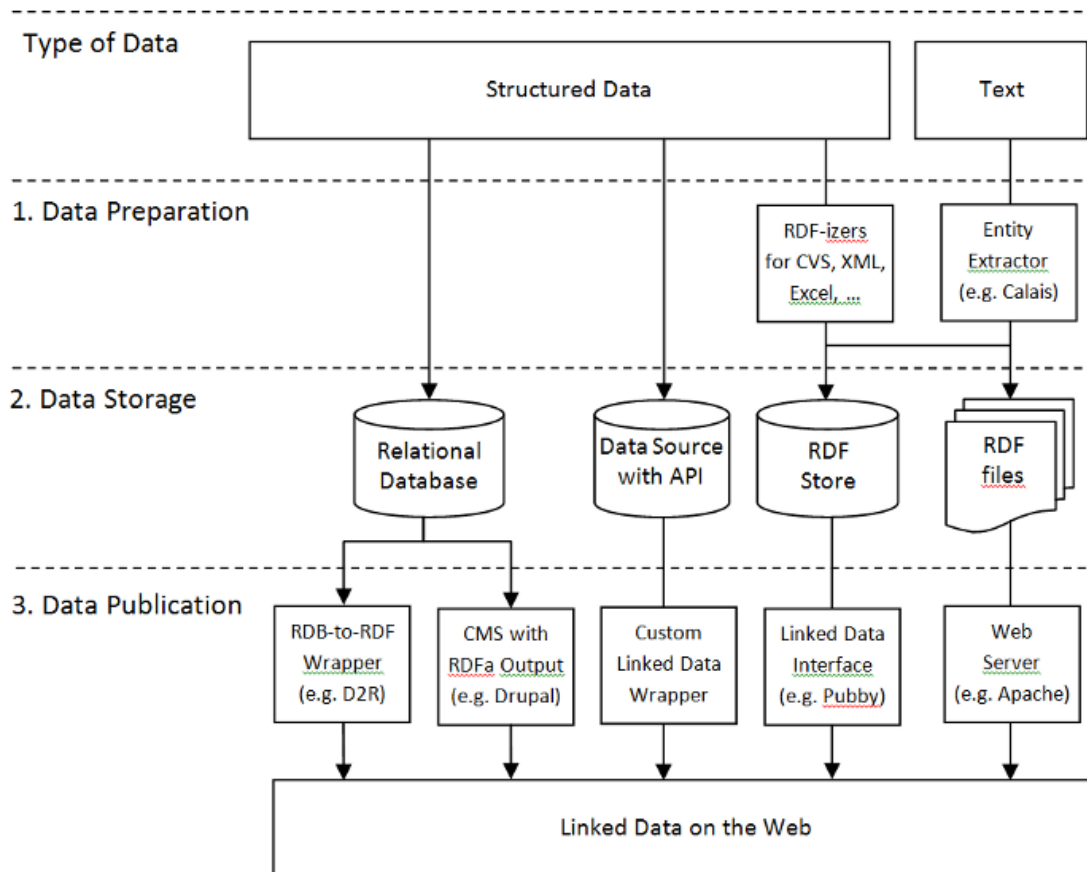
<sup>6</sup><http://linkeduniversities.org/>

Both using the crawling pattern to extract, map and store the data in a Linked Data format instead of transforming them for every request.

## TABLOID

One of the outcomes next to the LOD application itself was the Tabloid ("Toolkit ABout Linked Open Institutional Data"), *"a toolkit intended to help institutions and developers to both publish and consume linked data"*. It contains work-flows, documentations, examples and tools [32] trying to address different roles such as managers, developers and users. Tabloid tries to help people to understand LD, what can be done with it and gives advice on a technical perspective, how to publish and consume LD, providing at the same time a detailed and generic way.

### 4.1.3 Linked Data book



**Figure 4.3:** Linked Data Publishing Options and Workflow according to the LD book

Another big effort among many others of describing LD in general, how to publish and consume them and how to implement applications, was done by the book *"Linked Data:*

*Evolving the Web into a Global Data Space* by Heath and Bizer [5], which received a lot of attention.

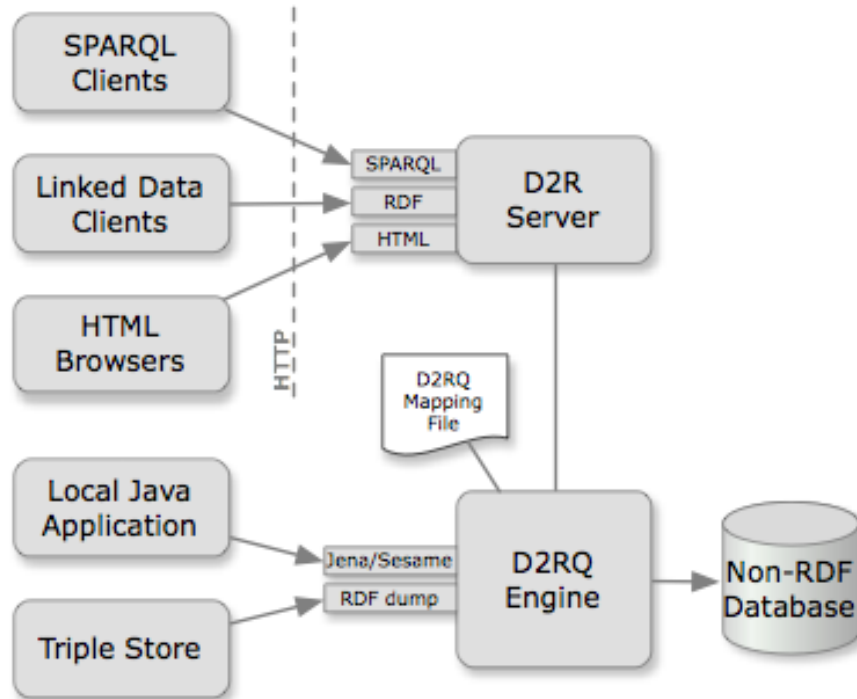
The book aims in general to give a basic understanding of LD and describing publication and consumption of LD. They providing advices and best practices, including architectures approaches, identifying the right set of URIs and vocabulary and much more. They also described an architecture, to be seen in figure 4.3

Next to patterns they also provide a general workflow for LD publishing, see figure 4.3. But comparing to the introduced architectures in the previous sections, the workflow has a different approach: instead of holding the data in a Triple Store, the workflow access and transforms the raw data on-the-fly for every request.

Next to this workflow, the book also provides various "recipes" for publishing LD and one of them is also to hold the data in a triple store as shown by Euclid and LUCERO. Furthermore the book provides a guide for the D2R-Server, which will be described in section 4.2.1.

## 4.2 Frameworks

### 4.2.1 D2RQ Platform



**Figure 4.4:** D2R Server architecture

**NOTE:** The last update on the D2RQ platform was in 2012 (version 0.8.1) and on the D2R Server in 2009 (version 0.7)

The D2RQ platform <sup>7</sup> was introduced by the Free University of Berlin and provides a database-to-RDF mapping. It is licensed under the terms of the GNU General Public License.

To map a relational database the platform provides a declarative mapping language, expressed in RDF, which is then be used to provide access to the database in the following, read-only, ways: [17]

- **RDF dumps**
- **RDF APIs**
- **SPARQL endpoint** (D2R Server)

<sup>7</sup><http://d2rq.org>

- **Linked Data**
- **HTML view** (D2R Server)

For an overview of the framework structure see figure 4.4.

### D2R Server

Part of the platform is the D2R Server <sup>8</sup>, which provides the public access to the platform over SPARQL and HTML, publishing it to the semantic web. More concrete, the server provides a dereferencing interface, for HTTP request dereferencing, and a SPARQL interface.

The server uses the mentioned **On-The-Fly Dereferencing Pattern** and does not provide a triple store, therefore it may be not has as good performance than tools with a triple store, although the team made a great effort to improve it.

Part of the server is also a tool which generates automatically a corresponding mapping and RDF vocabulary for an existing table structure, using table names as class names and column names as property names. The generated mapping file can then be customized. [18]

The following applications are examples using D2R-Server:

- DBLP Bibliography (University of Hannover) <sup>9</sup>
- DBtune (University of London) <sup>10</sup>
- Database of the Nobel Prize <sup>11</sup>

### 4.2.2 Information Workbench

The Information Workbench <sup>12</sup> is a high customizable tool to support the building of Linked Data applications, from basic data integration up to rich UI and visualization. The tool is developed by fluidOps and is published as Community Edition free available and under an Open Source License with a limited selection of capabilities and only for non-productive use (educational use, testing, development). The enterprise edition is also available but not for free.

---

<sup>8</sup><http://d2rq.org/d2r-server>

<sup>9</sup><http://dblp.uni-trier.de/>

<sup>10</sup><http://dbtune.org/>

<sup>11</sup><http://data.nobelprize.org/>

<sup>12</sup>[https://www.fluidops.com/en/products/information\\_workbench/](https://www.fluidops.com/en/products/information_workbench/)

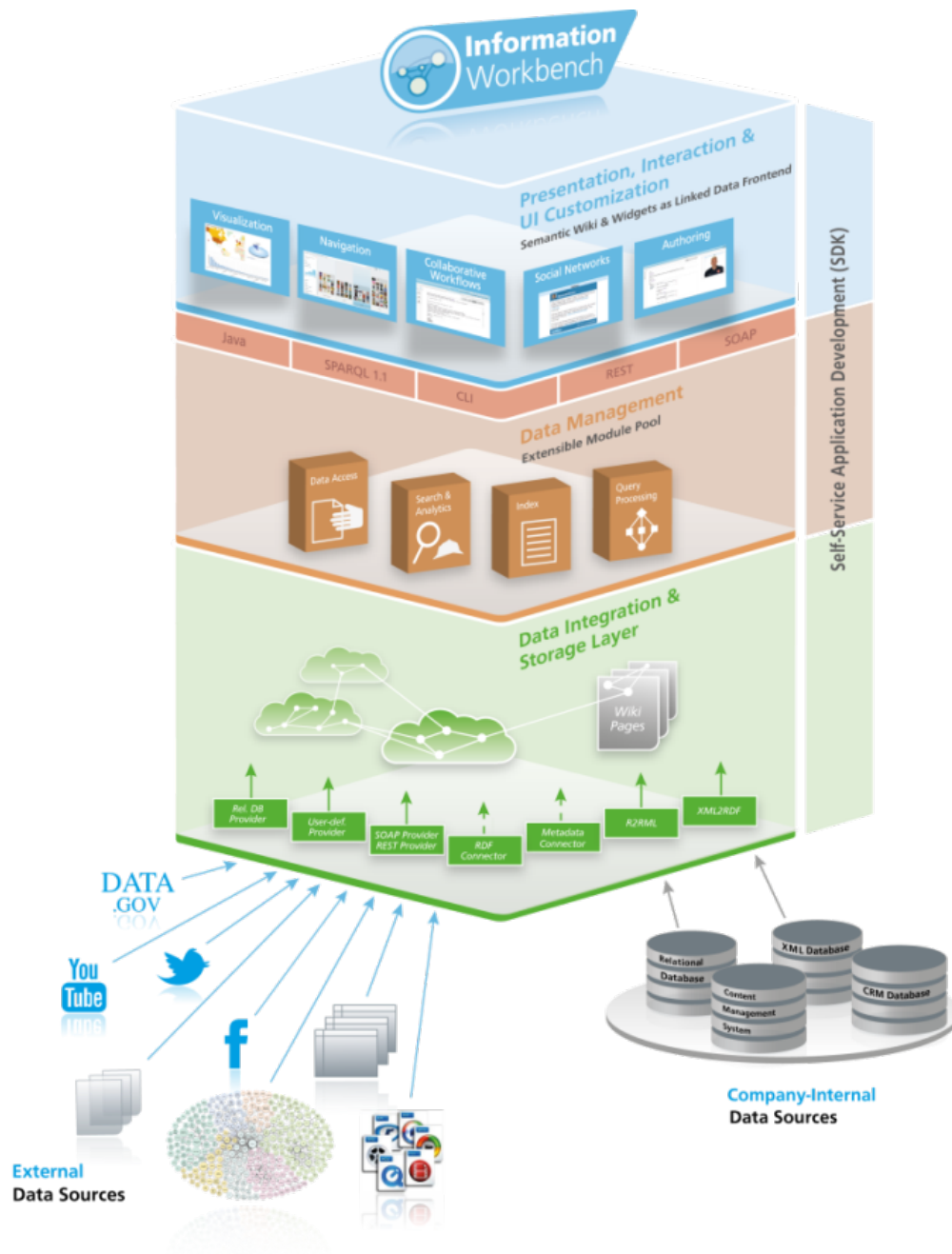


Figure 4.5: Architecture of the Information Workbench

The workbench consists of four layers (see figure 4.5 for an overview): [19] [20]

- **Persistence** Using so-called *providers*, the layers offer capabilities to integrate and convert data from different data source and stores them in a central triple store. Alternatively it also supports virtualized integration of local and public Linked Data sources using a *federation layer*.
- **Platform** On top of the persistence layer the core Platform layer a selection of modules and functionalities covering generic needs of Linked Data applications, the most important are a *Semantic Wiki & Widget Engine*, an *User Management & Access Control*, a *Search & Analytics Engine* and a *Workflow Engine*.
- **SDK** To support customized applications the workbench provides a SDK (Solution Development Kit) for developers to build domain specific applications, including *extensible data providers*, *data management facilities*, modified *ontologies*, *templates*, *widgets* and different APIs for extensive *system configuration*, *rules* and *workflow*.
- **Solution** On top of all layer stands the final solution, the application itself, which is either directly deployed through a RESTful API or over a zipped file for other installation approaches.

The resulting application is again customizable by widgets and different views, enabling data exploration and visualization.

### 4.2.3 LDIF - Linked Data Integration Framework

LDIF<sup>13</sup> was developed by the University of Mannheim and is published under the terms of the Apache Software License. It is implemented in Scala and aims to translate "*heterogeneous Linked Data from the Web into a clean, local target representation while keeping track of data provenance.*".

From a component perspective, LDIF consists of pluggable modules and a runtime environment, managing the data flows between them. The modules are: [21] [22]

- **Data Access Modules & Scheduler** For accessing the data to transform, LDIF provides several ways to import them. These import jobs are managed by a scheduler, which frequently fills a local cache. The module supports Triple/Quade Dump (for RDF/XML, N-Triples, N-Quads and Turtle formats), Crawler (using LDSpider<sup>14</sup>) and SPARQL imports.
- **Data Translation** For translating Web data using different vocabularies into a single target vocabulary, LDIF uses the R2R Mapping Language<sup>15</sup>.

---

<sup>13</sup><http://ldif.wbssg.de/>

<sup>14</sup><https://github.com/ldspider/ldspider>

<sup>15</sup><http://wifo5-03.informatik.uni-mannheim.de/bizer/r2r/spec/>



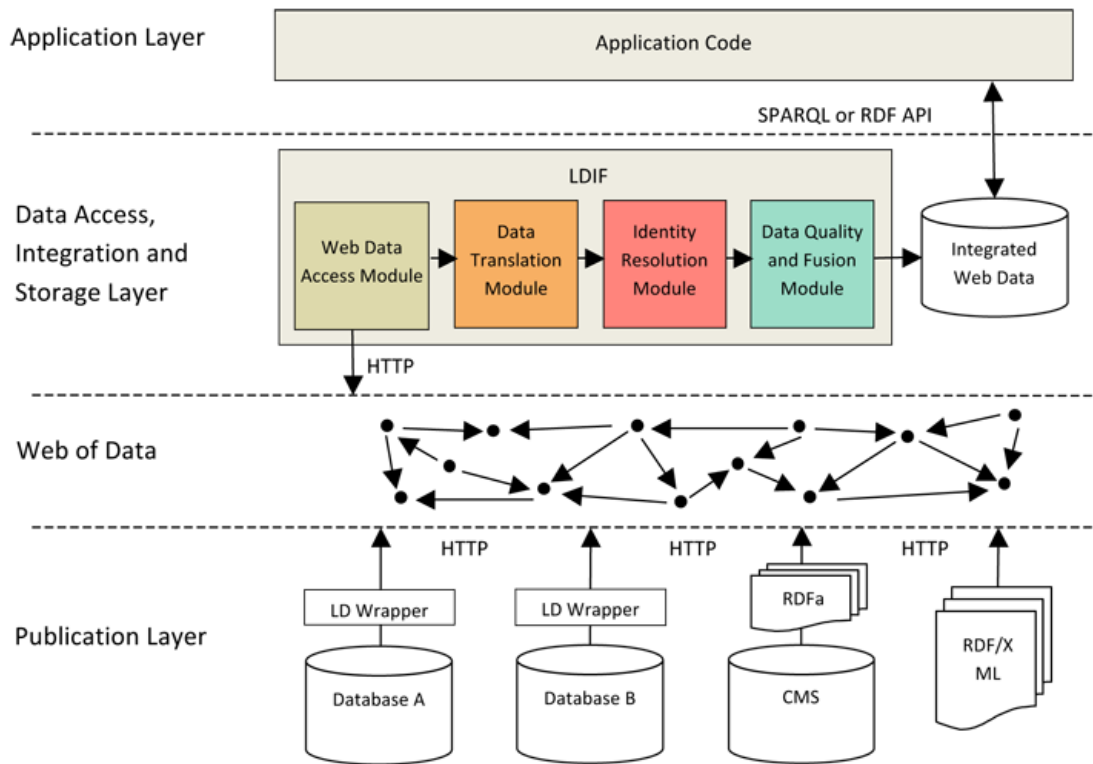


Figure 4.6: LDIF in the context of a LD application

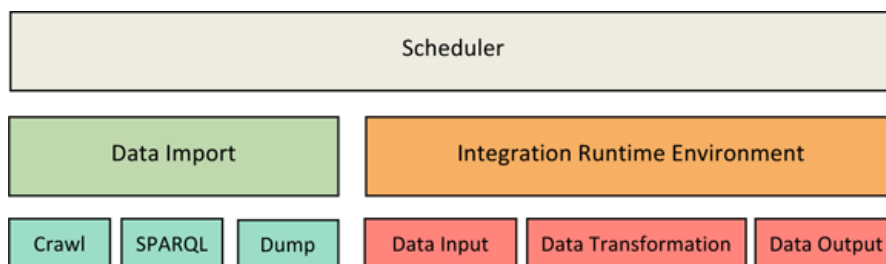


Figure 4.7: Components of LDIF

- **Identity Resolution** To find different URIs in different data pointing to the same entity, LDIF employs the Silk Link Discovery Framework with the Silk - Link Specification Language (Silk-LSL).
- **Data Quality Assessment and Fusion** For quality assessment, LDIF uses the Sieve Data Quality Assessment and Data Fusion Framework <sup>16</sup>.
- **Data Output** In the final step, LDIF write the cleaned data together with the provenance information in a single N-Quads file or without the meta-information in a N- Triples file.
- **Runtime Environment** As mentioned, the runtime environment manage the data flow between each module, providing an in-memory (fast, but limited scalable), a RDF store (using Apache Jena TDB and SPARQL queries, better scalable for the price of performance) and an Hadoop version.

### Silk

Silk <sup>17</sup> is "*an open source framework for integrating heterogeneous data sources.*" using the declarative Silk - Link Specification Language (Silk-LSL). It generates RDF links between data sets by custom link specifications. There are three different variations: [23]

- **Silk Single Machine** generates RDF links between two data items on a single machine.
- **Silk MapReduce** is for big scale datasets, using Hadoop and distributes to multiple machines.
- **Silk Server** is intended be used as an identity resolution component of as Linked Data consuming application. It provides a REST interface and runs as an HTTP server.

For details about the Link Discovery Framework see [6] and [24], for the server version consult [23].

### 4.2.4 Eclipse RDF4J (formerly Sesame)

Eclipse RDF4J <sup>18</sup> (formerly known as Sesame) is *a powerful Java framework for processing and handling RDF data. This includes creating, parsing, scalable storage, reasoning and querying with RDF and Linked Data. It offers an easy-to-use API that can be connected to all leading RDF database solutions.* It can be used as an embedded part of an application or as a stand-alone server.

---

<sup>16</sup><http://wifo5-03.informatik.uni-mannheim.de/bizer/sieve/>

<sup>17</sup><http://silkframework.org/>

<sup>18</sup><http://rdf4j.org/>

Originally developed as Sesame by Aduna as part of the "On-To-Knowledge" project (1999-2002), it was official forked into RDF4J. It is licensed under a BSD-style license.

The framework comes with many components, like Alibaba, an API for mapping Java classes onto ontologies. The RDF database API is unlikely similar solutions, it consists of stackable interfaces for adding functionality. Next to the intern abstract storage engine (SAIL, Storage and Inference Layer), many other triplestores are supported, like Ontotext GraphDB, Mulgara, and AllegroGraph

#### 4.2.5 Apache Jena

Apache Jena <sup>19</sup> is *a free and open source Java framework for building Semantic Web and Linked Data applications*. It was originally developed by HP Laboratories and now maintained by the Apache Software Foundation and is licensed under the Apache License 2.0.

The framework provides an API to extract data from and write to RDF, supporting relational databases, RDF/XML, Turtle and Notation 3. In contrast to RDF4J it also supports OWL.

More concrete, Jena can be used to manipulate RDF data, storing them in a triple store and publish it as a SPARQL access point. This HTTP interface is called *Fuseki*, which is in fact a sub-project of Jena and can be also run as stand-alone server using the Jetty web server.

---

<sup>19</sup><https://jena.apache.org/>

### 4.3 Excluded Tools And Projects

#### 4.3.1 LD-Patterns

The Linked Data Patterns book by Dodds and Davis (see [7]) tried to give an overview of existing design pattern regarding LD. But they don't give concrete architectures or architecture relating information, so this thesis will not use its content. But it is suggested, that this design pattern catalogue is used additionally when creating an application.

#### 4.3.2 LOD2 Stack

The LOD2 stack, introduced by Auer et. al., *is an integrated distribution of aligned tools which support the whole life cycle of Linked Data from extraction, authoring/creation via enrichment, interlinking, fusing to maintenance.* [25] For this thesis the proposed stack of technology was too generic to compare it with other frameworks and the website of the project <sup>20</sup> was at point of writing this thesis offline, therefore it was excluded of this thesis.

#### 4.3.3 LODUM

Another interesting project is the LODUM project (Linked Open Data University of Münster), the Open Data initiative of the university, hosted at the Institute for Geoinformatics' Semantic Interoperability Lab (MUSIL). The project team has co- initiated both LinkedUniversities.org and LinkedScience.org.

It was excluded for this thesis because the project don't provide public documentation of their architecture or any other part of their technical details <http://lodum.de/>

#### 4.3.4 Synth and SHDM

Synth <sup>21</sup> is a development environment for building SHDM <sup>22</sup> (Semantic Hypermedia Design Method) modeled applications, providing a set of modules, receiving SHDM generated models. Synth comes with a web browser GUI for adding and editing these models. A conceptual view of the architecture can be seen in figure 4.8, where the dashed boxes are modules and the whites boxes insides the module components. [8] The authors de Souza Bomfim and Schwabe describe in two papers, how a Linked Data application can be build with the environment: [8] and [26].

Since their description is very abstract and there are no further documentations of the tool, it was excluded for this thesis.

---

<sup>20</sup><http://stack.linkeddata.org/lod2//>

<sup>21</sup><http://www.tecweb.inf.puc-rio.br/synth>

<sup>22</sup>[https://www.w3.org/2005/Incubator/model-based-ui/wiki/SHDM\\_-\\_Semantic\\_Hypermedia\\_Design\\_Method](https://www.w3.org/2005/Incubator/model-based-ui/wiki/SHDM_-_Semantic_Hypermedia_Design_Method)

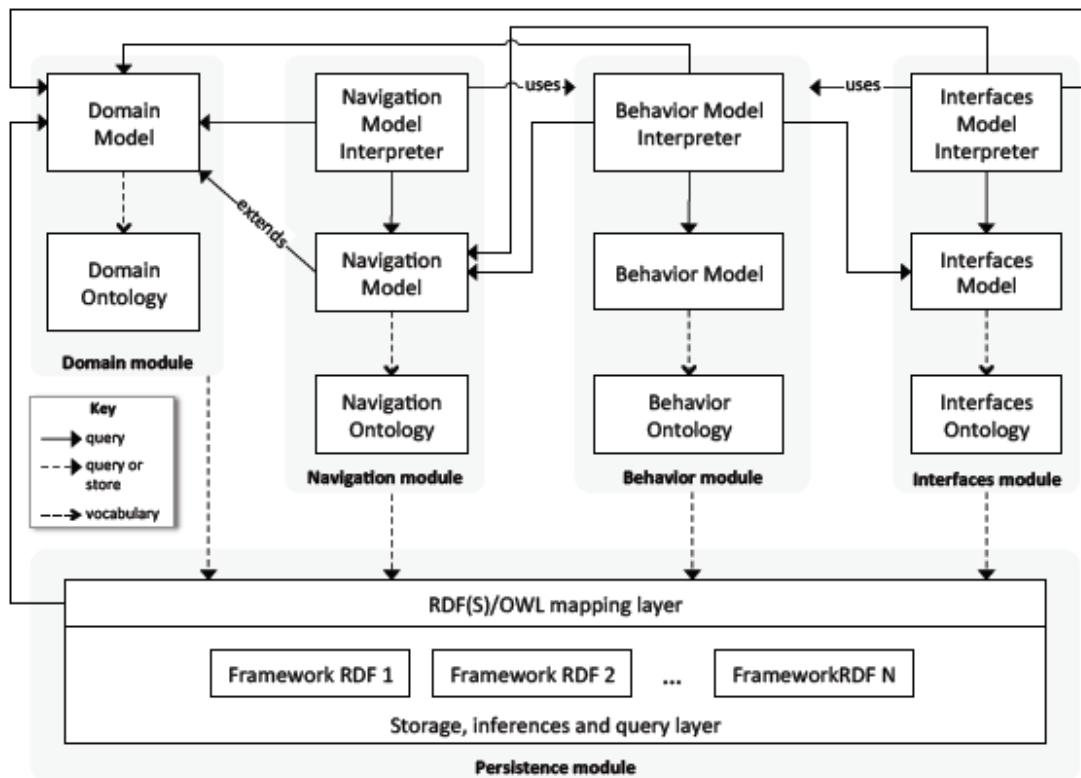


Figure 4.8: Concept of the Synth Architecture



# An Evaluation Framework for LD solutions (RQ2)

Since there are a wide variety of solutions, it is hard to find a set of criteria which can be applied to all in the same way in order to compare them. Therefore this thesis will use 4 different criteria groups to ensure higher cover of them and allowing at the same time, that criteria might not be applicable for some solutions.

## 5.1 Criteria from Previous Study

Criteria	Scala	Explanation
Maintainability	+/-	How much effort needs the maintenance? (less/much)
Data Freshness	yes/ no	Can it deal with new data?
Flexibility	yes/ no	Can it deal with heterogeneous and/or legacy data? Can it deal with changes in the ontology?

**Table 5.1:** Criteria group 1

From a previous study [10] conducted by the author, users at TU Wien are expecting and requesting from a LD application: **clear data ownership**, management of **data freshness** and **data quality** and **maintenance**. Since data ownership is a concern of organization and cannot be clarified by a tool, data freshness, data quality and maintainability are introduced as criteria. Since data quality is a very generic term, it will be used as criteria category. Additionally a concern of the stakeholders from the paper was how to deal with legacy data, therefore flexibility is included to Data quality.

The resulting criteria can be seen in table 5.1.

## 5.2 Usability Criteria

Criteria	Scala	Explanation
Effectiveness	+/-	How well do the users achieve their goals using the system? (good/bad)
Efficiency	+/-	What resources are consumed in order to achieve their goals? (less/much)
Satisfaction	+/-	How do the users feel about their use of the system? (good/bad)
Security	+/-	How well is security ensured? (good/bad)
Learnability	+/-	How much time is needed to learn the system? (less/much)

**Table 5.2:** Criteria Group 2: Usability

In every software application and especially solutions designed for end-users, usability is a huge and very important point these days. There are many definitions and measurements, ISO and other proposed models, trying to classify and define usability. This thesis will use an enhanced ISO model, proposed by Abran et.al. [9]. Since the goal of this thesis is not a complete usability analysis of the tools/frameworks and the variety of the chosen tools is too wide, the analysis of these aspect will be more of a general type. Abran et.al. are proposing various measurements for the different categories of their model, the interested reader might use them for a detailed analysis. For this thesis, these measurements will only be used as a guideline to estimate an assessment for the tools.

The model can be seen in table 5.2

## 5.3 Data Formats

An important aspect for the final decision for one of the tools can be the supported data format. It might be an (external) requirement or resulting from the fact of existing data. Since this highly depends on the context and the use case, this criteria will not be rated in any way, this thesis will only line out the supported data formats.

## 5.4 Linked Data Publishing Checklist

Since the whole thesis is about Linked Data, it is important to analyse not only the solutions itself but also the resulting LDs. In order to do that, the Linked Data Publishing Checklist by Heath et.al. [5] will be used. Another alternative could be the the LD definition itself by Tim Berners-Lee [33], but this thesis will presume, that a LD tool will produce valid LD. As one can expect from a checklist, the rating for this criteria will be only fulfilled/not fulfilled.



Criteria
Q1: Does your data set links to other data sets?
Q2: Do you provide provenance metadata?
Q3: Do you provide licensing metadata?
Q4: Do you use terms from widely deployed vocabularies?
Q5: Are the URIs of proprietary vocabulary terms dereferenceable?
Q6: Do you map proprietary vocabulary terms to other vocabularies?
Q7: Do you provide data set-level metadata?
Q8: Do you refer to additional access methods?

**Table 5.3:** Criteria Group 4: Linked Data Publishing Checklist

The checklist can be seen in table 5.3

## 5.5 Evaluation framework

Since this thesis is designed to be used as a help for decisions, the aim can not be to find the "best" solution, therefore there will not be such a thing like an end result. E.g. for some situation a solution with bad usability can be more appropriate because of the supported data formats than a solution with a higher overall rating.

Criteria group	Criteria
Group 1	Maintainability
	Data Freshness
	Flexibility
Group 2: Usability	Effectiveness
	Efficiency
	Satisfaction
	Security
	Learnability
Group 3: Data formats	Data formats
Group 4: LD Publishing Checklist	LD Publishing Checklist

**Table 5.4:** Complete Criteria Catalogue



## Comparison Of LD Solutions Based On The Evaluation Framework (RQ3)

In this section the comparison itself will be done. In order to do that, first in section 6.1 the classification introduced in section 3.3 will be applied to the found solution. Then in section 6 the criteria defined in section 5 will be applied for each of them, divided in the defined groups. The summary in section 6.6 will then give an overview of the done comparison.

An overview of the solutions to compare can be seen as recap in table 6.1, in order to simplify the following tables, each of the solutions is given an ID to refer.

ID	Framework
1	Euclid Project
2	LUCERO
3	Linked Data book
4	D2RQ Platform
5	Information Workbench
6	Linked Data Integration Framework
7	Eclipse RDF4J
8	Apache Jena

**Table 6.1:** Overview of the solutions

## 6.1 Classification

ID	Arch- itecture	Full- Stack	Present- ation layer	Business Layer	Data Access Layer	Note
1	x			x	x	
2	x			x	x	
3	x			x	x	
4			x	x	x	D2R Server includes HTML view and SPARQL endpoints
5		x				
6				x	x	
7				x	x	
8			x	x	x	Provides optionally SPARQL endpoints and stand-alone server with Jetty

Table 6.2: Classification

## 6.2 Criteria Group 1: Criteria from Previous Study

ID	Maintain- ability	Data Freshness	Flexibility	Note
1	+	yes*	yes*	*depending on implementation
2	?*	yes	yes*	*further investigations required
3	+	yes*	yes*	*depending on implementation
4	+	yes	yes*	*depending on configuration
5	?*	yes*	yes*	*depending on configuration
6	+	yes	yes*	
7	~*	?*	yes*	*depending on implementation
8	~*	?*	yes*	*depending on implementation

Table 6.3: Comparison Criteria Group 1

### 1: Euclid Project

Since the Euclid Project is very generic, the Data Freshness and Flexibility totally depends on how the concepts are implemented in an application. If the concepts are correctly implemented, Maintainability is probably well supported, since the three-layered architecture is widely recognized for it. For Data Freshness and Flexibility, the

implementation needs to be done focused on it, in order to achieve it, since the project does it not explicitly.

## **2: LUCERO**

The project supports Data Freshness by its Updaters and Schedulers, and can react to ontology changes by its Entity Name System. But it is unclear, how well it supports Maintainability, further investigations are required, which exceed the scope of this thesis.

## **3: Linked Data book**

Similar to the Euclid Project, the Data Freshness and Flexibility is possible well supported, but depends on the implementation. It is also not an explicit focus. Maintainability again is inherit of the concept.

## **4: D2RQ Platform**

The D2RQ Mapping Language can flexibly handle heterogeneous data as well as legacy data, the quality of Data Freshness depends on the configuration but is itself inherit of the concept.

## **5: Information Workbench**

Since the Workbench is designed to easily integrate different data sources, it can handle legacy data and ontology changes relatively good, depending on custom configurations. It might be a challenge to find the correct configurations, but the provided documentation is well done.

## **6: LDIF**

LDIF is explicit designed to handle various data sources and vocabularies, mapping them to an uniform vocabulary. The framework further provides scheduler to keep the data as fresh as needed (running hourly, daily etc.)

## **7: Eclipse RDF4J**

All three criteria are highly depending on the implementation of the framework and therefore how maintainability, freshness and flexibility are handled. The framework itself does only interact with a repository (read/write), which can potential be filled by another application which handles freshness and flexibility. Maintainability could be a bit tricky in this situation, if the other application is out of control of the RDF4J implementation.

## **8: Apache Jena**

Jena is very similar to RDF4J, all statements there are valid for Jena too.

### 6.2.1 Summary Group 1

Most of the solutions are supporting all data freshness and flexibility although most of them depending on the concrete implementation or configuration of the application. Only LDIF have a clear support of especially data freshness and flexibility. Maintainability is a critical point for most of the solutions.

## 6.3 Criteria Group 2: Usability

ID	Effectiveness	Efficiency	Satisfaction	Security	Learnability	Note
1	N.A.	N.A.	N.A.	-	+	Not applicable, good documentation
2	N.A.	N.A.	N.A.	-	-	No UI, no Security, bad documentation
3	N.A.	N.A.	N.A.	-	+	Not applicable, good documentation
4	+	+	+	-	+	Good documentation and UI, no security
5	+	+	+	+	+	* "easy to learn, hard to master"
6	N.A.	N.A.	N.A.	-	+	Not applicable, good documentation
7	N.A.	N.A.	N.A.	~	+	Not applicable, good documentation
8	N.A.	N.A.	N.A.	+	+	Not applicable, good documentation

**Table 6.4:** Comparison Criteria Group 2: Usability

### 1: Euclid Project

Since the Usability totally depends on the implemented application, no general statements can be given here. Security is not explicit. The documentation is very good.

### 2: LUCERO

The project does not support an explicit UI, therefore Usability is hard to evaluate. It also does not have any integrated security. The documentation is old and possible outdated, therefore learnability of the system is relatively bad.

### 3: Linked Data book

Again, similar to Euclid, the Usability strongly depends on the implementation and is therefore not applicable. Security is not explicit.

#### **4: D2RQ Platform**

The D2RQ server provides a basic but good readable UI. The documentation of the whole platform is very good. Security is not part of the project.

#### **5: Information Workbench**

The workbench has a very well designed UI, supporting the user to achieve their tasks relatively simple. The system is very good readable, enabling the user to explore the functionality and fulfilling their task. The according documentation is well done for "standard" users, but might be not enough for advanced development. Security is per default activated.

#### **6: LDIF**

LDIF does not have an explicit UI, therefore there can be no assessment of it done. Security is not part of the LDIF concept. The documentation is well done, but covers only simple topics, which could be problematic when running into problems while using the framework.

#### **7: Eclipse RDF4J**

The framework provides a simple UI for the RDF4J server and workbench for managing the repositories, but using RDF4 means writing Java Code. Therefore UI can not be topic of examination here. The RDF4J repositories have simple user management, but are handled per default by plain text cookies in the browser. The documentation is detailed and in good condition.

#### **8: Apache Jena**

Again, all statements from RDF4J are valid for Jena too. But in contrast to RDF4J Jena uses Apache Shiro for security. The server for Jena is called Fuseki.

#### **6.3.1 Summary Group 2**

Most of the solutions does not support an UI, except D2RQ and the information workbench. Security is only explicit mentioned in the information workbench and Apache Jena, RDF4J's default security is very basic. All solutions except LUCERO provide a good documentation.

### **6.4 Criteria Group 3: Data Formats**

#### **1: Euclid Project**

The Euclid Project does theoretically supports every kind of formats, since it requires the developer to write a consumer for each data source, additionally enhanced by a wrapper

Solution	Data formats	Note
1	Potential every possible format	Every format needs its own wrapper/consumer
2	Default RSS & XML	Additional data formats need custom extractors
3	Recipes for RDF, XML, HTML, relational databases, Wrapper	
4	Only relational databases	
5	Table-base (CSV, excel, groovy, JDBC, rest, TSV, SPARQL etc), tree-based (XML, JSON, etc), RDF	Each data source needs a configured data provider which provides R2RML (tabular data sources) and XML mappings to RDF (see)
6	N-Quads dumps, RDF/XML, N-Triples, Turtle dumps, dereferenced URIs, SPARQL	Using LDSpider for URI crawl import
7	Only RDF	
8	Only RDF & OWL	

**Table 6.5:** Comparison Criteria Group 3: Data formats

to transform the data to RDF.

## 2: LUCERO

Per default, the project supports only RSS and XML, but it provides also mechanism to integrate additional data extractors to the system, to enable addition data formats.

## 3: Linked Data book

The book provides recipes for RDF, XML, HTML, relational databases and wrapper for existing applications and Web APIs. Additional data formats can probably integrated by adopting the recipes.

## 4: D2RQ Platform

The scope of the platform are only the integration of relational databases.

## 5: Information Workbench

The workbench does support a wide range of data format. For each data source, a data provider has to be configured. For table- and tree-based data formats are mappings available to transform the data to RDF.



## 6: LDIF

LDIF provides four types of import jobs: Quad (import N-Quads dumps), Triple (import RDF/XML, N-Triples or Turtle dumps), Crawl (import by dereferencing URIs as RDF data, using the LDSpider Web Crawling Framework) and SPARQL Import Job (import by querying a SPARQL endpoint)

## 7: Eclipse RDF4J

The framework is mainly meant for working with data from a RDF repository and not explicit for putting/creating data in it in the first place.

## 8: Apache Jena

Jena is also developed with the focus of accessing data rather than creating the data store, therefore only RDF and additionally OWL are supported.

### 6.4.1 Summary Group 3

The solutions offer a wide range of data types in total, the Euclid Project and the LD book provide recipes for various data types, the Information Workbench does this out-of-the-box. D2RQ, RDF4J and Jena are specialized solutions, focusing only on single data types.

## 6.5 Criteria Group 4: Linked Data Publishing Checklist

ID	1	2	3	4	5	6	7	8
Q1: Does your data set links to other data sets?	x	x	x	x	x	x	?	?
Q2: Do you provide provenance metadata?	x	?	x	x	x	x	x	x
Q3: Do you provide licensing metadata?	x	?	x	x	x	x	x	x
Q4: Do you use terms from widely deployed vocabularies?	x	x	x	x	x	x	?	?
Q5: Are the URIs of proprietary vocabulary terms dereferenceable?	x	x	x	x	x	x	?	?
Q6: Do you map proprietary vocabulary terms to other vocabularies?	x	x	x	x	x	x	-	-
Q7: Do you provide data set-level metadata?	x	?	x	x	x	x	x	x
Q8: Do you refer to additional access methods?	x	-	x	x	x	x	?	?

**Table 6.6:** Comparison Criteria Group 4: Linked Data Publishing Checklist

## 1: Euclid Project

Since of its generic nature, possible every point of the checklist can be fulfilled depending on the implementation. But the architecture does not require any of the points.

### **2: LUCERO**

LUCERO supports interlinking of data and through its Entity Name System different vocabularies. Since its bad documentation, it is unclear, if any metadata are provided (assumable not).

### **3: Linked Data book**

Since Heath et al. propose the checklist in this book, it is also implicit and explicit integrated into the described solutions.

### **4: D2RQ Platform**

The D2RQ server provides comprehensive support for metadata, easy customizable by templates. The D2RQ Mapping language is very powerful in handling and mapping various kinds of vocabulary.

### **5: Information Workbench**

The workbench does support metadata, licensing metadata can be provided by custom implementations. Interlinking is as well supported as different kinds of vocabularies.

### **6: LDIF**

The framework implicit provides provenance next to the triple store, links to other data sets and maps proprietary vocabulary terms. The LD book explicit mentions LDIF as good example, therefore it can easily be assumed, that the checklist is fulfilled.

### **7: Eclipse RDF4J**

Again, same as the Euclid Project, nearly every point of the checklist could be fulfilled by an implementation, especially the points about meta data. The schema/ vocabulary/structure can be independent of the application and managed by another one, therefore it is out of the scope of RDF4J.

### **8: Apache Jena**

For Jena are the same arguments than for RDF4J valid: meta data are depended on the implementation, schema/vocabulary/structure can be out of the scope.

#### **6.5.1 Summary Group 4**

The checklist is overall well supported although it, again, is depending on the implementation of some of the solutions like Euclid, RDF4J or Jena. The last two are not managing the schema/vocabulary/structure by themselves, therefore the checklist is only partially applicable.

## 6.6 Summary

Overall the concept of having multiple criteria worked, since some solution could not be evaluated in some categories like Usability. But in the overall view, the criteria helped to find a standardized way to describe the solutions and compare them. In the following section, the evaluation of each solution will be summarized to have a better overview.

### 1: Euclid Project

Since the Euclid Project does only provide a generic architecture, most of the criteria are not directly applicable and are depending on the concrete implementation. Correctly applied it does however supports directly or indirectly maintainability, data freshness, flexibility, various kinds of data formats and the complete LD checklist. The architecture has no focus, neither explicit nor implicit, on security, therefore it needs to be done additionally if required. The documentation of the architecture is outstanding and very well done. Comparing to the other solutions, Euclid provides structures how other solutions or custom implementation can interact with each other. It can also be used as blueprint when combining other solutions.

### 2: LUCERO

LUCERO is on one hand overall bad documented and outdated. On the other hand, it does support data freshness and flexibility per default due its mechanisms and can support various data formats by custom extractors. But due the first facts, it is not recommended any more to use LUCERO in a real application.

### 3: Linked Data book

The summary for the LD book architecture is similar structured to the one of the Euclid project: good documentation but overall very generic. It can support maintainability, data freshness and flexibility if well implemented. The documentation is good and does include various recipes of data formats. A possible use case for it is the same as for Euclid: as blue print for combining other solutions.

### 4: D2RQ Platform

The D2RQ is as mentioned a specialized tool for relational databases. According to it, it has only limited data format supports. But due its simple structure and focused usage, data freshness, maintainability and flexibility can be assured. D2RQ includes a good UI and documentation. This solution is ideal for a very specific use case, requiring only a relational database to publish as L(O)D, but solving it as an all-in-one solution without the need of further tool integration.

### 5: Information Workbench

The Information Workbench is an all-in-one solution, rich with functionality and with a wide range of supported data formats. It can be used as a full stack solution with UI, integrating different data repositories. The documentation is well done but simple written, resulting maybe into problems handling more complex problems. The Workbench is ideal for an use case involving multiple data sources. For smaller use cases, it could be an overloaded solution.

### 6: LDIF

LDIF is due its concept on the one side a very flexible framework for handling different data sources with a wide range of vocabulary which also provides mechanism for keeping the retrieved data fresh. On the other hand the specialization leads to a specific focus resulting in a limit number of supported data formats. An use case for LDIF can be managing different data sources on base of RDF, SPARQL or something similar. Since it is not an all-in-one solution, it can/should be used in combination with another tool responsible for exposing the data to the web.

### 7: Eclipse RDF4J

RDF4J is completely different to the other investigated solutions, since it is a Java framework, requiring the developer to implement the given APIs. It is designed to handle a given repository and work with its data and/or expose them to the web. If an use case requires to publish data *not* in RDF format, it could be good idea to use RDF4J in combination with a tool like LDIF.

### 8: Apache Jena

As already mentioned Jena is similar to Sesame, but with the addition, that it provides support for OWL. It is also designed to handle a given repository and might be used in combination with something like LDIF.

# A Case Study Of Applying The Evaluation Framework At TU Wien (RQ4)

In this section, the thesis tries to answer research question 4: What a LD solution for TU Wien might look like.

## 7.1 General Considerations

In order to answer the question, it is important to give a context, since it is difficult respectively impossible to give a general, generic solution. In the following sections, three possible scenarios will be proposed, which can be a possible use case at TU Wien. For each scenario the requirement will be defined and a solution assigned.

From the previous study [10] conducted by the author the following requirements are important for the TU stakeholders:

- **Maintainable:**  
A LD application at TU Wien has to be easily maintainable. Citing stakeholders from the previous study, it "simply has to work without having to care about continuously".
- **Fresh Data:**  
It was important for the stakeholders, that the provided data are up-to-date.
- **Legacy Data:**  
A possible solutions for TU Wien has necessarily handle legacy data since this are the kind of data, that are candidates for publishing.

These requirements are therefore valid for each following scenario, meaning that criteria group 1 is higher weighted than the other groups. A solution, that does not full-fill one of this criteria, will not be take account of at all.

### 7.2 Situations & Requirements

As mentioned, to answer the research question, the context needs to be specified. This will be done by analyzing the stakeholders (already done by the previous study). These findings will then be combined to scenarios.

#### 7.2.1 Stakeholders

**Researcher** are interested in an easy access to the publications of the university and publications in general as well as an easy access to resources for their research subject. They need to *use* the data without further work in order to focus on their actual daily work. That also applies if the research is teaching at the university: in order to use an application, either by accessing it or providing data to it, it would be rejected by the stakeholders, if it means significant more work. But since they are used to technical work, the inhibition threshold to use a LD application is lower than for the administration staff.

**Administration staff** is more a subject of data source than of data user since they are holding a lot of interesting data. But on the other hand, the inhibition threshold is higher than for the researcher since the technical experience is at average lower. Therefore the argument of an "easy-to-use" of the previous stakeholder group is even more valid here.

**Students** are mainly data consumer and therefore more concerned with the frontend of a possible application. This group can build application that are consuming the published data.

#### 7.2.2 Situations

In the following sections will describe possible scenarios at TU Wien, how and which data are going to be published. There are two base scenarios: publishing a single data source and building a platform with various data sources. The second one will be additional split.

##### Situation 1: Specialized Single Solution

Only one or a small number of data source(s) needs to be published. The application has to be specific done for this data set. The project is small scaled and can be done quickly. A possible data source could be the publication database. This scenario can be used as a demonstration, how to *consume* Linked Open Data, but it is not advisable to use this approach if any other data source might be published in the future, since this scenario is *not* scalable

##### Requirement

- small scale
- small number of data sets
- specialized solution
- simple

### **Scenario 2: Function-rich Platform**

A comprehensive platform is needed, which covers a variety of data of different formats, including relational data as well as semi-structured data like XML. The project is medium scaled, the platform shall be able to handle data TU wide and is not only for internal purpose but also for external usage. Additional, the platform has to offer meta data, documentation and performance data. The solution has to be easy to use and implement. Licensing or Open Source does not play a role.

#### **Requirement**

- various kinds of data sets (in number and formats)
- medium scale
- platform with additional features
- easy to use/implement

### **Scenario 3: Complete Controlled Platform**

Similar to scenario 2, it is necessary to build platform combining different kinds of data sources. But in contrast, full control over the platform and the technology stack is necessary and no licensing is wanted, an Open Source Solution is necessary. In exchange, the requirement "easy to use/implement" is relaxed, but still relevant.

#### **Requirement**

- various kinds of data sets (in number and formats)
- medium scale
- platform with additional features
- no licensing
- Open Source

## **7.3 Proposed Solutions**

### **7.3.1 Scenario 1: Specialized Single Solution**

There are two solution thinkable:

**Option 1: D2RQ** can be a solution for this situation. It provides a simple mapping of a given data set and a method of publishing with the D2R Server. The project can be fast implemented and be ready. Other data sets may be added. Since D2RQ only handles relational data, this approach is not suitable for semi- structured data.

**Option 2: A combination of Jena/RDF4J with Euclid and LDIF** can be used to publish the given data. The Euclid architecture (assuring maintainability) can be used as a blueprint to implement an application using either the Apache Jena or Eclipse RDF4J framework. If a mapping is needed, LDIF can be used additionally. This solution is far more extensive than using D2RQ, but on the other hand more flexible and customizable. The scope can be controlled by defining the limits, specially by economizing the UI, but be aware that this approach can result in a lesser usability and the scope might be too big for this situation.

### 7.3.2 Scenario 2: Function-rich Platform

For this scenario is the **Information Workbench** suitable. The tool can integrate different kinds of data sources, combine the access points to a platform and providing additional features over modules. The implementation can be done relatively simple and fast. However, it needs a licensing if used outside an educational scope.

### 7.3.3 Scenario 3: Complete Controlled Platform

For this scenario no all-in-one solution of the proposed tools is suitable, instead the platform has to be developed by using a framework like **Jena or RDF4J** in combination with the **Euclid Project architecture**, similar to the solution proposed situation 1. If a mapping to an unified vocabulary is needed, **LDIF** can also be used. This project will probably consume more time and will have a bigger scope than the solution of situation 2. But on the other hand, this approach offers more customization options.

## 7.4 Summary

Taking a deeper look at the proposed situations and solutions, it can be seen, that it is most likely to use solutions like Jena or RDF4J in combination with a meta architecture like Euclid and utilization tools like LDIF. Tools like D2RQ are limiting and the Information Workbench needs a licensing. Looking at the popularity charts of sites like DB-Engines <sup>1</sup>, it currently seems like Jena is quite more popular.

This proposed solution is of course not complete, there are similar tools to Jena and RDF4J, like MarkLogic <sup>2</sup> or Virtuoso <sup>3</sup>. But it can be stated, that this *kind* of solution

---

<sup>1</sup>[https://db-engines.com/en/ranking\\_trend/system/Jena\protect\kern+.2777em\relaxRDF4J](https://db-engines.com/en/ranking_trend/system/Jena\protect\kern+.2777em\relaxRDF4J)

<sup>2</sup><http://www.marklogic.com/>

<sup>3</sup><https://virtuoso.openlinksw.com/>



Scenario	Requirement	Solution
<b>Specialised Single Solution</b>	<ul style="list-style-type: none"> <li>-) small scale</li> <li>-) small number of data sets</li> <li>-) specialized solution</li> <li>-) simple</li> </ul>	<ul style="list-style-type: none"> <li>-) Option 1: D2RQ</li> <li>-) Option 2: Jena + Euclid (+ LDIF)</li> </ul>
<b>Function-rich Platform</b>	<ul style="list-style-type: none"> <li>-) various kinds of data sets (in number and formats)</li> <li>-) medium scale</li> <li>-) platform with additional features</li> <li>-) easy to use/implement</li> </ul>	Information Workbench
<b>Complete Controlled Platform</b>	<ul style="list-style-type: none"> <li>-) various kinds of data sets (in number and formats)</li> <li>-) medium scale</li> <li>-) platform with additional features</li> <li>-) no licensing</li> <li>-) Open Source</li> </ul>	Jena/RDF4J + Euclid (+ LDIF)

**Table 7.1:** Summary of Scenarios & proposed solutions for TU Wien

is the most suitable of the proposed *kinds* of solutions.



# Conclusion And Future Work

The overall goal of this thesis was to compare common LD solutions and give TU Wien a guideline for choosing and developing such an application. The concrete research question was:

**RQ:** *How do common LD solutions compare against each?*

1. **RQ1:** What are existing LD solutions?
2. **RQ2:** What are criteria to compare solutions?
3. **RQ3:** How do they compare against each other?
4. **RQ4:** What can be a solution for TU Wien?

The first part of this thesis investigated the first research question by conducting a literature study and discussing the term "framework" and why it was not used instead of the term "solution" (see subsection 3.1). The methodology can be seen in chapter 3, the results in section 4.

In the second part (section 5, a set of criteria was developed in order to compare the found solutions and answer research question 2. These criteria were then used to investigate the found solutions under the aspects of them. The results can be found in section 6.

In the final part, the found solutions were examined on the usefulness for TU Wien, in section 7.

## 8.1 Conclusion

In the following subsections each research question will be addressed and revisited :

### 8.1.1 Existing LD solutions

Using the method of a literature review, eight solution were found for this work:

- Euclid Project
- LUCERO
- Linked Data Book
- D2RQ Platform
- LDIF
- Eclipse RDF4J
- Apache Jena

Of the found candidates a few were discarded since they were not (public) documented or too generic.

### 8.1.2 The Evaluation Framework For LD Solutions

Since the found solutions showed a high variation, it was necessary to find a set of criteria which can be partially not applicable while still meaningful in their entirety. Therefore four criteria group were defined: *Criteria from a Previous Study*, *Usability*, *Data formats and the Linked Data Checklist*. Since the aim of this thesis is not to find the "best" solution, the criteria were not weighted and had only a scala as assessment instead of e.g. a point system.

### 8.1.3 Application Of the Proposed Evaluation Framework

As expected, some criteria (especially Usability) were not applicable for some solutions, but the concept of multiple criteria groups worked out. The results (an overview of the full results can be seen in the Appendix) in short are:

- The Euclid Project was too generic to give a definite evaluation for the most criteria, since most of it is depending on an actual implementation. Correctly applied it does however supports directly or indirectly maintainability, data freshness, flexibility, various kinds of data formats and the complete LD checklist.
- LUCERO was found outdated and bad documented, it is not recommended to use it (any more).
- Similar to Euclid, the Linked Data Book was too generic to find explicit assessments for the criteria, too much is depending on the actual implementation.
- The D2RQ platform has (due its nature as specialized tool) only limited support for data formats, but received a good evaluation in the other criteria groups.

- As an all-in-one solution the Information Workbench has good results in all groups. It has to be noted, that this tool requires licensing when used outside an educational scope.
- LDIF is due its concept on the one side a very flexible framework for handling different data sources with a wide range of vocabulary. On the other hand the specialization leads to a specific focus resulting in a limit number of supported data formats. It is recommended to use this solution in specialized situation or in combination with other solutions like Jena or RDF4J.
- Eclipse RDF4J and Apache Jena both are very different to the other solutions, they are Java frameworks, requiring to write code against their API. For both the results of the comparison highly depend on the actual implementation.

#### 8.1.4 Solution For TU Wien

For the TU Wien a set of situations were developed in order to find a context for which a solution can be recommended. It was found, that the "ideal" solution, based on the given situations, requirements and stakeholders, would be an self- developed platform, using either Jena or RDF4J (or a similar tool) while using the Euclid architecture as a blueprint and LDIF if a mapping is necessary.

## 8.2 Future Work

It is recommended to extend the list of possible solutions and applying the given set of criteria in order to find a suitable solution for TU Wien. The found combination of Jena/RDF4J with Euclid and LDIF might be a way of developing a platform, but it is not necessary to use this *exact* combination. Nevertheless it should be keep in mind, that this way is more extensive than using e.g. the Information Workbench. But as an advantage, the full platform can be controlled.



# CHAPTER 9

## Appendix

ID	Arch- itecture	Full- Stack	Present- ation layer	Business Layer	Data Access Layer	Note
Euclid Project	x			x	x	
LUCERO	x			x	x	
Linked Data book	x			x	x	
D2RQ Platform			x	x	x	D2R Server includes HTML view and SPARQL endpoints
Information Workbench		x				
LDIF				x	x	
Eclipse RDF4J				x	x	
Apache Jena			x	x	x	Provides optionally SPARQL endpoints and stand-alone server with Jetty

**Table 9.1:** Classification

ID	Maintain-ability	Data Freshness	Flexibility	Effective-ness	Efficiency	Satis-faction	Security	Learn-ability
Euclid Project	+	yes	yes	N.A.	N.A.	N.A.	-	+
LUCERO	?	yes	yes	N.A.	N.A.	N.A.	-	-
Linked Data book	+	yes	yes	N.A.	N.A.	N.A.	-	+
D2RQ Platform	+	yes	yes	+	+	+	-	+
Information Workbench	?	yes	yes	+	+	+	+	+
LDIF	+	yes	yes	N.A.	N.A.	N.A.	-	+
Eclipse RDF4J	~	?	yes	N.A.	N.A.	N.A.	~	+
Apache Jena	~	?	yes	N.A.	N.A.	N.A.	+	+

Table 9.2: Comparison group 1 &amp; 2



ID	Data formats	Q1	Q2	Q3	Q4	Q5	Q6	Q7	Q8
<b>Euclid Project</b>	Potential every possible format	x	x	x	x	x	x	x	x
<b>LUCERO</b>	Default RSS & XML	x	?	x	x	x	x	x	x
<b>Linked Data book</b>	Recipes for RDF, XML, HTML, relational databases, Wrapper	x	?	x	x	x	x	x	x
<b>D2RQ Platform</b>	Only relational databases	x	x	x	x	x	x	?	?
<b>Information Workbench</b>	Table-based (csv, excel, groovy,jdbc, rest, tsv, sparql etc), tree-based (xml, json, etc), RDF	x	x	x	x	x	x	?	?
<b>LDIF</b>	N-Quads dumps, RDF/XML,N-Triples, Turtle dumps, dereferenced URIs, SPARQL	x	x	x	x	x	x	-	-
<b>Eclipse RDF4J</b>	Only RDF	x	?	x	x	x	x	x	x
<b>Apache Jena</b>	Only RDF & OWL	x	-	x	x	x	x	?	?

**Table 9.3:** Comparison group 3 & 4

Scenario	Requirement	Solution
<b>Specialised Single Solution</b>	<ul style="list-style-type: none"><li>-) small scale</li><li>-) small number of data sets</li><li>-) specialised solution</li><li>-) simple</li></ul>	<ul style="list-style-type: none"><li>-) Option 1: D2RQ</li><li>-) Option 2: Jena + Euclid (+ LDIF)</li></ul>
<b>Function-rich Platform</b>	<ul style="list-style-type: none"><li>-) various kinds of data sets (in number and formats)</li><li>-) medium scale</li><li>-) platform with additional features</li><li>-) easy to use/implement</li></ul>	Information Workbench
<b>Complete Controlled Platform</b>	<ul style="list-style-type: none"><li>-) various kinds of data sets (in number and formats)</li><li>-) medium scale</li><li>-) platform with additional features</li><li>-) no licensing</li><li>-) Open Source</li></ul>	Jena/RDF4J + Euclid (+ LDIF)

**Table 9.4:** Summary of Scenarios & proposed solutions for TU Wien

# List of Figures

2.1	Excerpt of the Berlin SPARQL Benchmark results . . . . .	5
2.2	Excerpt of the Berlin SPARQL Benchmark results . . . . .	6
4.1	General EUCLID architecture . . . . .	12
4.2	LUCERO work flow & architecture . . . . .	14
4.3	Linked Data Publishing Options and Workflow according to the LD book	15
4.4	D2R Server architecture . . . . .	17
4.5	Architecture of the Information Workbench . . . . .	19
4.6	LDIF in the context of a LD application . . . . .	21
4.7	Components of LDIF . . . . .	21
4.8	Concept of the Synth Architecture . . . . .	25

# List of Tables

3.1	Overview of the Classification . . . . .	10
5.1	Criteria group 1 . . . . .	27
5.2	Criteria Group 2: Usability . . . . .	28
5.3	Criteria Group 4: Linked Data Publishing Checklist . . . . .	29
5.4	Complete Criteria Catalogue . . . . .	29
6.1	Overview of the solutions . . . . .	31
6.2	Classification . . . . .	32
6.3	Comparison Criteria Group 1 . . . . .	32
6.4	Comparison Criteria Group 2: Usability . . . . .	34
6.5	Comparison Criteria Group 3: Data formats . . . . .	36
		55

6.6	Comparison Criteria Group 4: Linked Data Publishing Checklist . . . . .	37
7.1	Summary of Scenarios & proposed solutions for TU Wien . . . . .	45
9.1	Classification . . . . .	51
9.2	Comparison group 1 & 2 . . . . .	52
9.3	Comparison group 3 & 4 . . . . .	53
9.4	Summary of Scenarios & proposed solutions for TU Wien . . . . .	54

# References To Refereed Scientific Work

- [1] A. Zaveri, A. Rula, A. Maurino, R. Pietrobon, J. Lehmann, and S. Auer, „Quality assessment for linked data: A survey“, *Semantic Web*, vol. 7, no. 1, pp. 63–93, 2016.
- [2] C. Bizer and A. Schultz, *The Berlin SPARQL Benchmark*. 2009.
- [3] E. Minack, W. Siberski, and W. Nejdl, „Benchmarking fulltext search performance of RDF stores“, *The Semantic Web: Research and Applications*, pp. 81–95, 2009.
- [4] D. Roberts and R. Johnson, „Evolving frameworks“, *Pattern languages of program design*, vol. 3, 1996.
- [5] T. Heath and C. Bizer, „Linked data: Evolving the web into a global data space“, *Synthesis lectures on the semantic web: theory and technology*, vol. 1, no. 1, pp. 1–136, 2011.
- [6] J. Volz, C. Bizer, M. Gaedke, and G. Kobilarov, „Silk-A Link Discovery Framework for the Web of Data.“, *LDOW*, vol. 538, 2009.
- [7] L. Dodds and I. Davis, „Linked data patterns“, *Online: <http://patterns.dataincubator.org/book>*, 2011.
- [8] M. H. de Souza Bomfim and D. Schwabe, „Synth-Linked Data Application Implementation Environment“,
- [9] A. Abran, A. Khelifi, W. Suryn, and A. Seffah, „Usability meanings and interpretations in ISO standards“, *Software Quality Journal*, vol. 11, no. 4, pp. 325–338, 2003.



# References To Non-Refereed Work

- [10] L. Baronyai, K. Haller, and S. Gamerith, „Linked Open Data at the Vienna University of Technology – A case study about research data“, 2016.
- [11] Y. Theoharis, V. Christophides, and G. Karvounarakis, „Benchmarking database representations of RDF/S stores“, in *International Semantic Web Conference*, Springer, vol. 3729, 2005, pp. 685–701.
- [12] M. Achichi, M. Cheatham, Z. Dragisic, J. Euzenat, D. Faria, A. Ferrara, G. Flouris, I. Fundulaki, I. Harrow, V. Ivanova, *et al.*, „Results of the ontology alignment evaluation initiative 2016“, in *CEUR workshop proceedings*, RWTH, vol. 1766, 2016, pp. 73–129.
- [13] M. Cheatham, Z. Dragisic, J. Euzenat, D. Faria, A. Ferrara, G. Flouris, I. Fundulaki, R. Granada, V. Ivanova, E. Jiménez-Ruiz, *et al.*, „Results of the ontology alignment evaluation initiative 2015“, in *10th ISWC workshop on ontology matching (OM)*, No commercial editor., 2015, pp. 60–115.
- [14] Z. Dragisic, K. Eckert, J. Euzenat, D. Faria, A. Ferrara, R. Granada, V. Ivanova, E. Jiménez-Ruiz, A. O. Kempf, P. Lambrix, *et al.*, „Results of the ontology alignment evaluation initiative 2014“, in *Proceedings of the 9th International Conference on Ontology Matching-Volume 1317*, CEUR-WS. org, 2014, pp. 61–104.
- [15] O. Kovalenko, E. Serral, and S. Biffl, „Towards evaluation and comparison of tools for ontology population from spreadsheet data“, in *Proceedings of the 9th International Conference on Semantic Systems*, ACM, 2013, pp. 57–64.
- [16] D. Riehle, „Framework design“, PhD thesis, 2000.
- [17] C. Bizer and R. Cyganiak, „D2rq-lessons learned“, in *W3C Workshop on RDF Access to Relational Databases*, 2007, p. 35. [Online]. Available: <https://www.w3.org/2007/03/RdfRDB/papers/d2rq-positionpaper/>.
- [18] —, „D2r server-publishing relational databases on the semantic web“, in *Poster at the 5th international semantic web conference*, vol. 175, 2006.

- [19] P. Haase, M. Schmidt, and A. Schwarte, „The information workbench as a self-service platform for linked data applications“, in *Proceedings of the Second International Conference on Consuming Linked Data-Volume 782*, CEUR-WS. org, 2011, pp. 119–124.
- [20] A. Gossena, P. Haase, C. Hüttera, M. Meiera, A. Nikolova, C. Pinkela, M. Schmidta, A. Schwarte, and J. Tramea, „The Information Workbench—A Platform for Linked Data Applications“,
- [21] A. Schultz, A. Matteini, R. Isele, C. Bizer, and C. Becker, „Ldif-linked data integration framework“, in *Proceedings of the Second International Conference on Consuming Linked Data-Volume 782*, CEUR-WS. org, 2011, pp. 125–130.
- [22] A. Schultz, A. Matteini, R. Isele, P. N. Mendes, C. Bizer, and C. Becker, „LDIF-a framework for large-scale Linked Data integration“, in *21st International World Wide Web Conference (WWW 2012), Developers Track, Lyon, France*, 2012.
- [23] R. Isele, A. Jentzsch, and C. Bizer, „Silk server-adding missing links while consuming linked data“, in *Proceedings of the First International Conference on Consuming Linked Data-Volume 665*, CEUR-WS. org, 2010, pp. 85–96.
- [24] A. Jentzsch, R. Isele, and C. Bizer, „Silk-generating rdf links while publishing or consuming linked data“, in *Proceedings of the 2010 International Conference on Posters & Demonstrations Track-Volume 658*, CEUR-WS. org, 2010, pp. 53–56.
- [25] S. Auer, L. Bühmann, C. Dirschl, O. Erling, M. Hausenblas, R. Isele, J. Lehmann, M. Martin, P. N. Mendes, B. Van Nuffelen, *et al.*, „Managing the life-cycle of linked data with the LOD2 stack“, in *International semantic Web conference*, Springer, 2012, pp. 1–16.
- [26] M. H. de Souza Bomfim and D. Schwabe, „Design and implementation of linked data applications using SHDM and synth“, in *International Conference on Web Engineering*, Springer, 2011, pp. 121–136.



# References To Websites

- [27] EUCLID, *EUCLID — EdUcational Curriculum for the usage of LInked Data*, [Online; accessed 5-September-2016], 2012-2014. [Online]. Available: <http://euclid-project.eu/index.html>.
- [28] E. Union, *Framework Programmes for Research and Technological Development*, [Online; accessed 8-September-2016], 2012-2014. [Online]. Available: [https://ec.europa.eu/research/fp7/index\\_en.cfm](https://ec.europa.eu/research/fp7/index_en.cfm).
- [29] EUCLID, *About Euclid*, [Online; accessed 5-September-2016], 2012-2014. [Online]. Available: <http://euclid-project.eu/about/project-description.html>.
- [30] —, *EUCLID — Chapter 5: Building Linked Data Applications*, [Online; accessed 5-September-2016], 2012-2014. [Online]. Available: <http://euclid-project.eu/modules/chapter5.html>.
- [31] M. d’Aquin, F. Zablith, E. Motta, O. Stephens, S. Brown, S. Elahi, and R. Nurse, *The LUCERO project – About*, [Online; accessed 15-September-2016], 11 Jun 2010. [Online]. Available: <http://lucero-project.info/lb/about/index.html>.
- [32] —, *The LUCERO project – Tabloid*, [Online; accessed 15-September-2016], 1Jun 2011. [Online]. Available: <http://lucero-project.info/lb/tabloid/index.html>.
- [33] T. Berners-Lee, *Linked data, 2006*, 2006. [Online]. Available: <https://www.w3.org/DesignIssues/LinkedData.html>.