

Titel der Arbeit

Optionaler Untertitel der Arbeit

BACHELORARBEIT

zur Erlangung des akademischen Grades

Bachelor of Science

im Rahmen des Studiums

Software und Information Engineering

eingereicht von

Kevin Haller

Matrikelnummer 1325694

an der Fakultät für Informatik

der Technischen Universität Wien

Betreuung: Ao.Univ.Prof. Dipl.-Ing. Dr.techn. Mag.rer.soc.oec Stefan Biffl, Univ.Doz.

Mitwirkung: MSc PhD Reka Marta Sabou, Project Ass.

Wien, 1. August 2016

Kevin Haller

Stefan Biffl

Title

Optional Subtitle of the Thesis

BACHELOR'S THESIS

submitted in partial fulfillment of the requirements for the degree of

Bachelor of Science

in

Software and Information Engineering

by

Kevin Haller

Registration Number 1325694

to the Faculty of Informatics

at the TU Wien

Advisor: Ao.Univ.Prof. Dipl.-Ing. Dr.techn. Mag.rer.soc.oec Stefan Biffl, Univ.Doz.

Assistance: MSc PhD Reka Marta Sabou, Project Ass.

Vienna, 1st August, 2016

Kevin Haller

Stefan Biffl

Erklärung zur Verfassung der Arbeit

Kevin Haller
Address

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, 1. August 2016

Kevin Haller

Danksagung

TODO: Ihr Text hier.

Acknowledgements

TODO: Enter your text here.

Kurzfassung

TODO: Ihr Text hier.

Abstract

TODO: Enter your text here.

Contents

Kurzfassung	xi
Abstract	xiii
Contents	xv
1 Introduction	1
1.1 Motivation	1
1.2 Problem description	3
1.3 Structure of the thesis	4
2 Background	5
2.1 Web of Data	5
2.2 RDF	5
2.3 Ontology	5
2.4 SPARQL	5
3 Related Work	7
3.1 Geospatial ontologies	7
3.2 Linked Open Data based map applications	12
3.3 Indoor modelling	13
4 Solution	17
4.1 Data acquisition	17
4.2 Prototype of ontology	20
4.3 Architectural prototype	24
4.4 Map application	29
5 Discussion	33
5.1 Evaluation	33
5.2 Further work	33
5.3 Conclusion	33
List of Figures	35

List of Tables	35
List of Algorithms	37
Acronyms	39
Bibliography	41

Introduction

1.1 Motivation

Today's universities usually have to manage a significant amount of information and have to provide systems to handle common services like finding and booking courses for students. As the current time demands it, those services are often provided over web sites participating in the Web of Documents (see figure 1.1). The domain of a university is quite complex and as a consequence it is likely that different isolated information systems handling a particular part of the domain evolve; resulting in an environment, where information is distributed over multiple disconnected data silos that may have different formats and/or data owners. Such a situation prevents universities of fully exploiting their data [30].

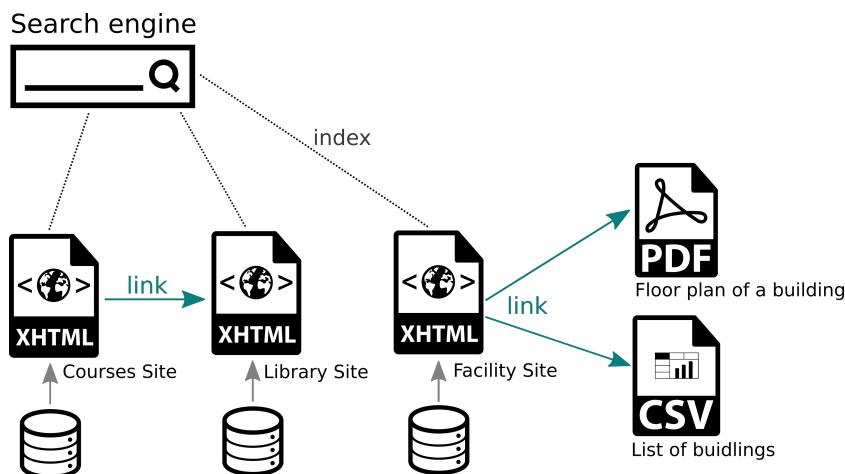


Figure 1.1: Illustration of the Web of Documents

1. INTRODUCTION

The same applies at the moment to the Vienna University of Technology (TU Vienna). In order to find a special location like a certain lecture hall and how it can be accessed without obstacles like stairways for persons with mobility-impairments, one has to search for information on different web sites, scan floor plans and eventually construct a convenient route to the location based on the gathered knowledge. For humans this procedure does not constitute a problem, but think of machines. The process of information retrieval through data mining or harvesting is quite difficult and/or time consuming. As a consequence application developers that may have innovative ideas, which would be a benefit for the information environment of the university, face a barrier that is hard to overcome. Linked Data (LD) is one way to transform this information published on multiple web sites into a university-wide data space (see figure 1.2).

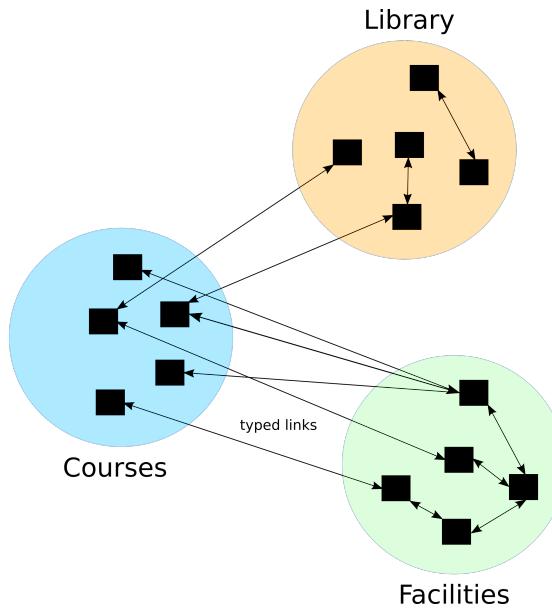


Figure 1.2: Illustration of the Web of Data

LD describes a method of publishing structured data so that it can be interlinked and become more useful. It builds upon standard Web technologies, but rather than using them to serve web pages for human readers, it extends them to share information in a way that can be read automatically by computers. This enables data from different sources to be connected and queried [6].

This thesis aims to show the potential that could evolve, if spatial data about TU Vienna is transformed into a machine-readable form such as LD, by proposing a prototype of a map application based on a subset of spatial data about TU Vienna. This subset is

intended to comply with the principles of LD [4] and to fulfil the submission criteria to the Linking Open Data cloud¹, which are as follows [8]:

- There must be resolvable `http://` (or `https://`) URIs.
- They must resolve, with or without content negotiation, to RDF data in one of the popular RDF formats (RDFa, RDF/XML, Turtle, N-Triples).
- The dataset must contain at least 1000 triples.
- The dataset must be connected via RDF links to a dataset that is already in the diagram. This means, either your dataset must use URIs from the other dataset, or vice versa. We arbitrarily require at least 50 links.
- Access of the entire dataset must be possible via RDF crawling, via an RDF dump, or via a SPARQL endpoint.

1.2 Problem description

As already mentioned does the TU Vienna distribute spatial information over different sources. Few of this spatial data is actually in a machine-readable format, so that it can be transformed into LD automatically. The **1st problem** to solve is therefore the transformation process of the unstructured part of the data into a machine-readable format. Section 4.1 is suggesting a solution for this problem.

The **2nd problem** is how the structured spatial data shall be described so that it can be easily used by application developers as well as integrated into the Web of Data. The developed ontology shall be compliant to the best practises for publishing ontologies [5]. Section 4.2 is suggesting a prototype of an ontology by taking already existing approaches (see section 3.1 and 3.3) into consideration.

After the spatial data was transformed into LD, a system must be designed to expose this data on the Web in order to qualify as Linked Open Data (LOD), representing the **3rd problem** to solve. The solution shall take best practises for publishing LD[13] and best practises for spatial data on the Web[24] into consideration. Section 4.3 is proposing an architectural prototype of such a system.

The **4th problem** is the development of a map application that is based on the resulting LOD. The application shall answer the following question: *"Give me all learning rooms that are nearby that are free in a given time range and are accessible without obstacles for person with mobility-impairments"* to show the potential of the data. Section 4.4 presents a solution to this problem.

¹<http://lod-cloud.net/>

1.3 Structure of the thesis

This chapter outlined the motivation for writing this thesis and gave a description of the problems for which a solution will be suggested. The rest of this thesis is structured as follows: Chapter 2 discusses the basic concepts, principles and technologies that build the foundation of LD and the Semantic Web. It is intended to be a brief introduction for readers that are not familiar with this topic. Chapter 3 provides a summary of past efforts in research and works related to the focus of this thesis. In the subsequent chapter 4 a solution for the given problems will be suggested and discussed. This includes the prototype of an ontology that models the problem domain as well as an architectural prototype of a system that manages a subset of spatial data about the TU Vienna in form of LD and provides a human- and machine-friendly interface to it. Finally, the implemented solution is evaluated and conclusions are drawn in the last chapter 5. It provides furthermore an outlook to future work and potential improvements.

CHAPTER 2

Background

2.1 Web of Data

TODO: Enter your text here.

2.2 RDF

Resource description framework (RDF) **TODO:** Enter your text here.

2.3 Ontology

TODO: Enter your text here.

2.3.1 RDFS

TODO: Enter your text here.

2.3.2 OWL

TODO: Enter your text here.

2.4 SPARQL

SPARQL protocol and RDF query language (SPARQL) **TODO:** Enter your text here.

CHAPTER 3

Related Work

In this chapter important efforts and research related to this thesis are outlined. Section 3.1 describes and evaluates geospatial ontologies that were suggested by [24] or discovered in the dataset of Linked Open Vocabulary (LOV) through a systematic search. Section 3.2 lists already existing map applications that are based on LOD and describes their capabilities. The final section 3.3 describes done research and designed ontologies for modelling indoor environments.

3.1 Geospatial ontologies

This section is going to outline common ontologies for describing geospatial data and to evaluate them from the perspective of the given problem description (see 1.2). The ontologies were either be suggested by [24] or discovered in the dataset of LOV by common spatial search terms like '*feature*', '*geometry*' or '*location*'. *Feature* and *Geometry* are widely used terms in these ontologies and describe two different concepts of geospatial science. A *feature* is simply a spatial entity that can be everything from a building with fixed position to a movable food truck as long as it has a spatial extent. *Geometry* as the name suggests is a certain geometric shape from points, lines to polygons. Geometric shapes can be used to describe the spatial extent of *features*.

The ontologies are visualized using the VOWL2 notation[16]. Classes are modelled as circles, properties in form of rectangles; data type properties have a green colour and object properties a blue one. Literals are presented in a yellow rectangle.

3.1.1 Basic Geo Vocabulary (WGS84)

Basic Geo (WGS84 long/lat) Vocabulary (GEO)[7] is a lightweight ontology published by W3C to describe the position of a spatial *feature* using the WGS84 geodetic reference

3. RELATED WORK

datum¹ with the properties longitude, latitude and attitude. The general class SpatialThing is the domain of all these properties. Figure 3.1 shows a visualization of this ontology.

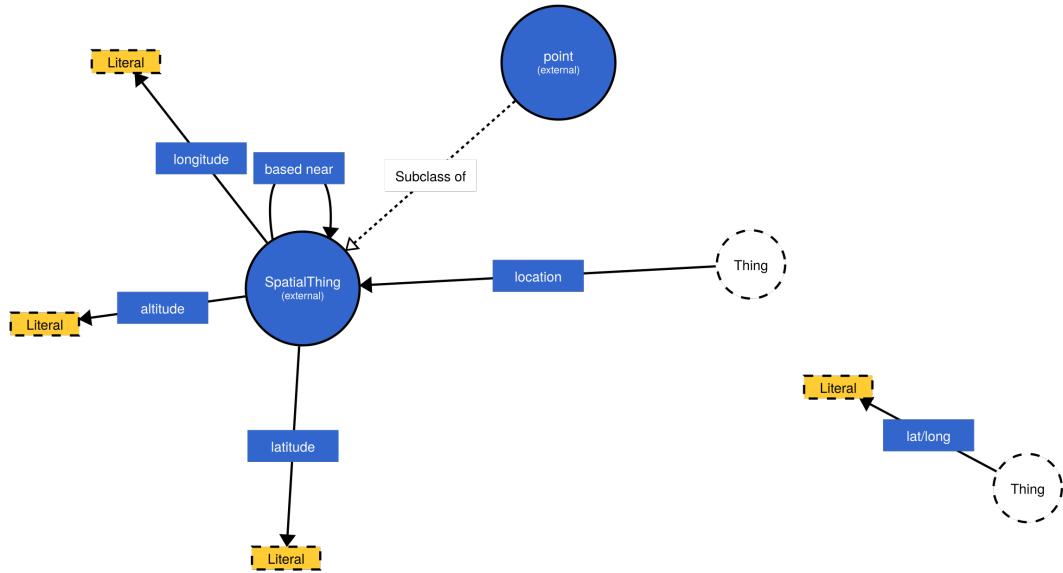


Figure 3.1: Visualization of GEO

Advantages are its simplicity and popularity as an analysis of the *Comprehensive Knowledge Archive* (CKAN, "The Data Hub"²) indicates. This analysis shows that GEO is used in **538** different datasets in the archive and is thereby the 4th most used ontology regarding usage in different datasets [21]. It is used by big players like DBpedia³, LinkedGeoData⁴ and GeoNames⁵. Listing 3.1 shows a data snippet of DBpedia describing the town square 'Karlsplatz' in Vienna with GEO.

Listing 3.1: Snippet of DBpedia

```

@prefix dbr: <http://dbpedia.org/resource/> .
@prefix geo: <http://www.w3.org/2003/01/geo/wgs84_pos#> .

dbr:Karlsplatz rdf:type geo:SpatialThing ;
  rdfs:label "Karlsplatz"@en , "Karlsplatz (Wien)"@de ;
  geo:lat "48.19916534423828125"^^xsd:float ;
  geo:long "16.370000839233398438"^^xsd:float .
  
```

¹http://en.wikipedia.org/wiki/World_Geodetic_System

²<http://thedatahub.org>

³<http://dbpedia.org>

⁴<http://linkedgeodata.org/About>

⁵<http://linkedgeodata.org/About>

GEO evolved into a standard for presenting the location of points of interest and as described by [13] in an article about best practises, standardized ontologies shall be used wherever possible to facilitate inclusion into the Web of Data.

However, this ontology is not sufficient for describing spatial data more complex than a single point in a coordinate reference system.

3.1.2 NeoGeo Geometry/Spatial Ontology

NeoGeo Geometry Ontology (NGEO)⁶ is an ontology that aims to provide a comprehensive descriptive power for modelling geographic regions, whereas NeoGeo Spatial Ontology (SPATIAL)⁷ aims to describe topological relationships between *features*. Both ontologies were designed to have a strict distinction between *features* and *geometries* [18].

NGEO

NGEO follows the principle of modelling geometries in pure RDF, which leads to a number of classes and properties for covering all shapes suggested by Simple Features Profile⁸. This includes points, lines and polygons. Shapes that go beyond a simple point are represented by a RDF collection of point instances, whereby point is a external class of GEO. Figure 3.2 shows a visualization of this ontology.

One advantage of this approach is that triple stores do not have to meet special requirements in order to enable the querying over geometric data like it is the case for GeoSPARQL. However, this shifts the burden of writing geometric queries to developers. A further problem is the high demand for resource identifiers or blank nodes (at least for each point of a shape), as the example of describing the geometry of Iceland shows⁹; this significantly increases the verbosity of the data without adding particular gains of expressivity to it[3]. NGEO has also a generic property asWKT for pointing to the Well-known text (WKT) serialization of a geometry, which has at the moment of writing the status 'deprecated'.

LinkedGeoData uses the *geometry* property of this ontology to point to the *geometry* of a *feature*, but the *asWKT* property of GeoSPARQL to represent the geometry.

SPATIAL

SPATIAL is an ontology that aims to provide properties to make topological relationships of *features* explicit, which else would only exist implicit in the geometric data. Included properties cover all topological relationships described by RCC-8 like does a *feature* overlap with another *feature*. One problem of making the relationships explicit is that these relationships have to be adopted if the *geometry* of *features* changes.

⁶<http://geovocab.org/geometry>

⁷<http://geovocab.org/spatial>

⁸<http://www.ogcnetwork.net/gml-sf>

⁹http://nuts.geovocab.org/id/IS_geometry.ttl

3. RELATED WORK

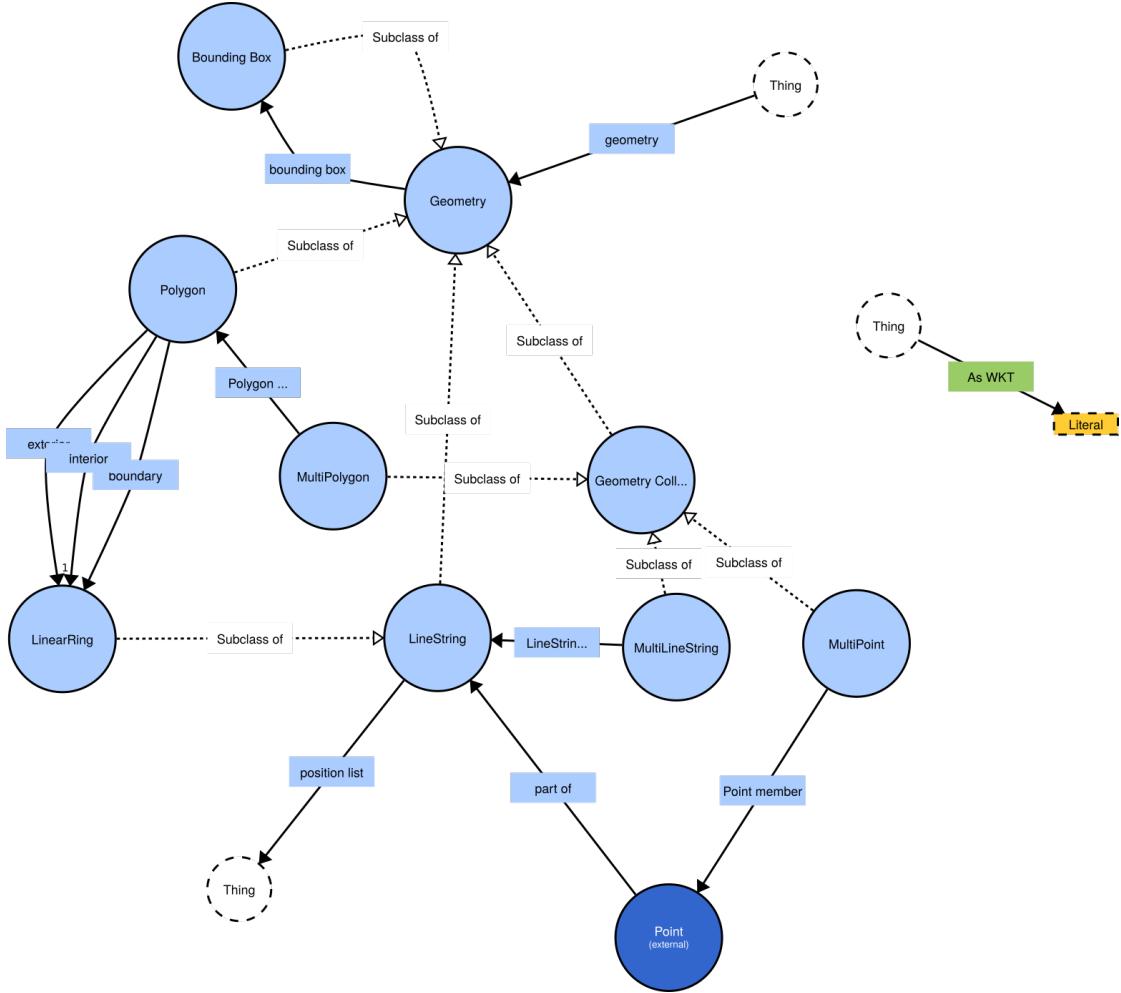


Figure 3.2: Visualization of NGEO

3.1.3 GeoSPARQL

GeoSPARQL defines a set of SPARQL extension functions, a set of Rule Interchange Format (RIF) rules, and a core ontology for geographic information[20]. Figure 3.3 shows a visualization of this core ontology. It has two major classes namely Feature to represent spatial *features* and Geometry to represent *geometries*. An instance of Geometry can be assigned to an instance of Feature with the property hasGeometry. It is based on the Simple Features model¹⁰ like NGEO. However, GeoSPARQL defines two properties that can be used to point to the serialization of the geometry; either asWKT for WKT or asGML for Geography Markup Language (GML). In contrast to NGEO, where *geometries* are represented in pure RDF.

¹⁰<http://www.ogcnetwork.net/gml-sf>

TODO: advantages, disadvantages

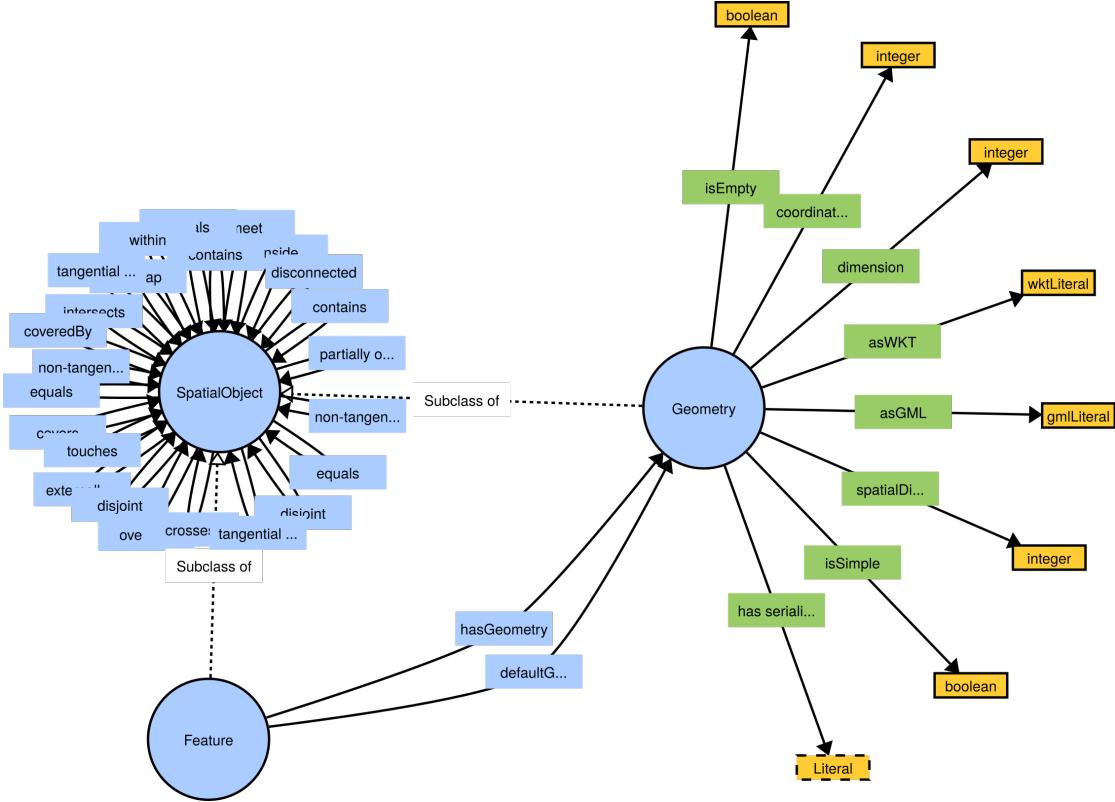


Figure 3.3: Visualization of GeoSPARQL

3.1.4 schema.org

schema.org¹¹ is an initiative that was started by Google, Bing and Yahoo! (subsequently joined by Yandex) to promote a common ontology for publishing structured data mark-up on web pages [10]. This machine-readable information can then be extracted from web pages by search engines to improve search results.

schema.org is a quite broad ontology covering a set of different domains from events, organizations to places. For this thesis especially the class **Place** is of interest describing entities with a fixed, physical extent; it has a set of more specific subclasses. This class has properties to assign an address to its instances (**address**) as well as to provide geo-information (**geo**). The value of the property **geo** can either be instances of the class **GeoCoordinates** or **GeoShapes**. **GeoCoordinates** represents a point in a georeference system similar to GEO. **GeoShapes** on the other hand provide the ability to describe more complex geometric shapes including circles, boxes, polygons and simple lines. All those shapes are expressed in textual form following an own specific pattern, in

¹¹<http://schema.org/>

3. RELATED WORK

contrast to GeoSparql, which declares the use of GML or WKT. However, such literals could be transformed into WKT by string operations (supported in SPARQL) except for the circle, which has no direct representation in WKT.

One advantage of this ontology is the visibility to search engines when exposed on web pages in one of the supported formats (Microdata¹², RDF in Attributes (RDFa)¹³ and JSON-LD¹⁴). As mentioned earlier, the embedded data can then be extracted and in fact major search engines crawl for it to enhance search results[11]. However, there are odd ascriptions of properties like that instances of the class Beach (subclass of Place) can have the property faxNumber; it must be considered that schema.org is meant for web masters to add structured metadata to their web sites [28] and not to describe a specific domain precisely.

3.1.5 ISA Programme Location Core Vocabulary

ISA Programme Location Core Vocabulary (LOCN) is an ontology that provides a set of properties and classes to describe the geometry, address and location of a *feature*. Figure 3.4 shows a visualization of this ontology. It provides a comprehensive set of properties to describe the address of a *feature*, but a minimal set for geometries and locations. For geometries this ontology suggests the use of external classes or a simple literal representing the serialization of the geometry in formats such as WKT and GML. Those suggested classes are *Geometry* of GeoSPARQL with its more specific subclasses and *GeoCoordinates* as well as *GeoShapes* from schema.org. In case of representing simple points also point of GEO is mentioned. [19] However, this quite broad range of possibilities does not make it easy for the consumer of such LD.

3.2 Linked Open Data based map applications

Multiple universities (e.g. Linked Universities¹⁵) exposed their public university data and map applications played out to be a common use case for this data.

University of Münster started an Open Data initiative named LODUM[14] from which a map application¹⁶ evolved that shows the location of buildings and which departments are located in the selected building, but there is no deeper insight into the buildings. However, there is a publication of the initiative that deals with the issue of indoor navigation [15]. Furthermore an indoor map ontology named LIMAP¹⁷ was designed.

University of Southampton is a further university that exposed spatial information about their campus as LOD. In contrast to the map of the LODUM initiative Southampton's

¹²<https://www.w3.org/TR/microdata/>

¹³<https://www.w3.org/TR/rdfa-syntax/>

¹⁴<http://json-ld.org/>

¹⁵<http://linkeduniversities.org>

¹⁶<http://app.uni-muenster.de/Karte/>

¹⁷<http://lodum.de/results/>

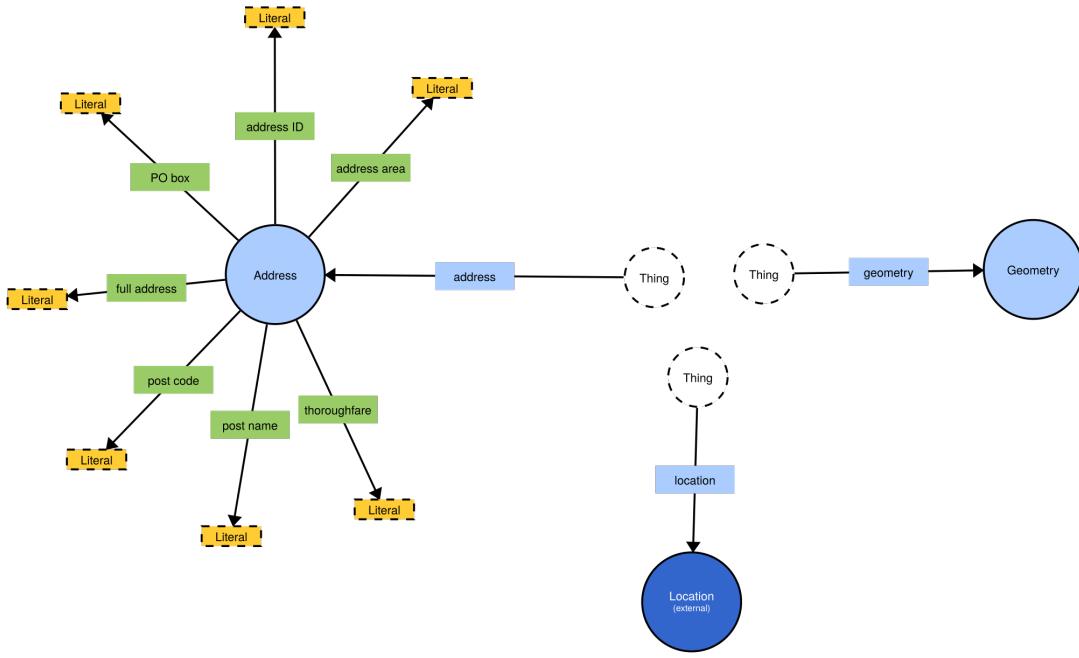


Figure 3.4: Visualization of LOCN

map application¹⁸ (see figure 3.5) gives insight into buildings and shows the location of certain rooms with the ability to switch the floor. It shows also computer rooms with their current estimated capacity and opening hours.

University of Oxford has also exposed spatial information in form of LOD on their Open Data portal and developed an application named *University Science Area Map*¹⁹. It highlights buildings, where departments of a certain field are located, on a map, but it gives no insight into the buildings.

3.3 Indoor modelling

Section 3.1 outlined ontologies to describe the spatial extent of *features* and topological relationships between them, but those ontologies are not intended for expressing the characteristics of *features* in an indoor environment; like for example that a *feature* is a learning room with certain opening hours or an elevator to get from one floor to another. This section discusses previous efforts and ontologies to model indoor environments and also how to navigate in them.

IndoorGML²⁰ is a Extensible Markup Language (XML) schema of OGC that aims to

¹⁸<http://maps.southampton.ac.uk/>

¹⁹<https://data.ox.ac.uk/explore/science-area/>

²⁰<http://indoorgml.net/>

3. RELATED WORK

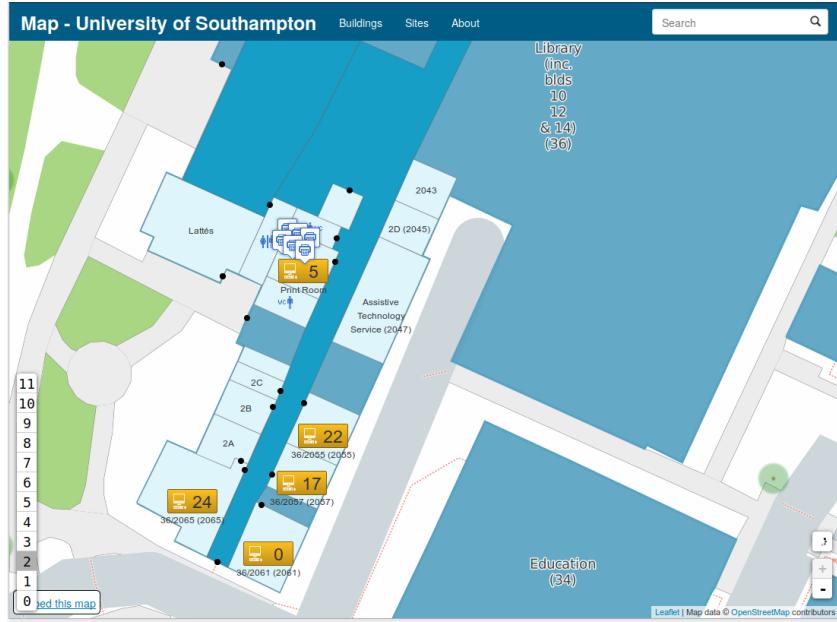


Figure 3.5: Map application of University of Southampton

provide a framework for representing and exchanging indoor spatial information, but no efforts has yet been made to transform it's model to the LD domain. OntoNav presents a semantic indoor navigation system that aims to follow a user-centric paradigm in which the capabilities, limitations and preferences of the user are taken into consideration for computing the best suiting path. A path is according to OntoNav a sequence of connected corridors, passages (stairways, ramps, elevators, etc.) and exits from one location to another[1]. ONALIN is a ontology that models the indoor environment similar to OntoNav, but it follows the approach of modeling buildings as network of hallways consisting of different points (start, decision and corner points); resulting in a finer granularity. ONALIN also takes the American Disability Act into consideration, which leads to properties for describing the height and number of steps of stairways or the height of sinks in a restroom [9]. Both ontologies are at the moment not locatable over the Web. iLoc is an ontology that has a minimal set of classes to describe the internal structure of a building from rooms, floors to vertical passages like elevators and stairways. Points of interest like entrances and landmarks play a special role in the navigation. A route section connects then two points and asserts that there is a walk-able path between them. This route section can have certain constraints like an access control measurement is required or steps have to be climbed [23].

3.3.1 Other ontologies

The ontologies mentioned in this section have been discovered in the repository of LOV by searching for terms such as '*room*', '*building*' and '*toilets*' or explored in datasets of

map applications based on LOD (see section 3.2).

Buildings and Rooms Vocabulary (ROOMS) provides a minimal set of classes and properties to describe the basic structure of a building. It has 6 classes Building, Floor, FloorSection, which is a named (identifiable) section of a floor, Room and Desk as well as Site, which describes an area of land like a campus.

CHAPTER 4

Solution

This chapter suggests solutions for the given problem description (see 1.2). Section 4.1 is dealing with the problem of transforming unstructured spatial data of the TU Vienna into a machine-readable format. The question of how this structured data shall be described is answered in section 4.2, where a prototype of an ontology is discussed. The subsequent section 4.3 is proposing an architectural prototype of a system for integrating and publishing LD consistent with the LD principles [4]. The final section 4.4 is presenting a map application to show the potential of the generated spatial LOD.

4.1 Data acquisition

Spatial data about the TU Vienna is distributed over multiple sources with different formats and data owners. In order to be eventually transformed into LD, this data must be extracted and prepared for the transformation. Figure 4.1 visualizes the different data sources.

The major source for building information is the web site of the facility management unit of the TU Vienna (GUT¹). This site lists all buildings and an enumeration of their floors in form of a (X)HTML table. It also contains links to the floor plans that are provided as Portable Document Format (PDF) files as well as to a PDF file for each building consisting of a table of all contained rooms with their intended function (e.g. office room, sanitary room). Section 4.1.3 deals with the problem of extracting and transforming tables from PDFs and web pages. The approach for processing floor plans is discussed in section 4.1.2.

The central information system TISS² provides a RESTful API that gives access to organizational information and the public address book. For this thesis, this API is useful

¹<http://www.gut.tuwien.ac.at/>

²<https://tiss.tuwien.ac.at/>

for linking persons and organizations to their offices. TISS has also a web page dedicated to presenting the event schedule of certain rooms like lecture halls, but this data is not available over the mentioned RESTful API; only in a human-readable form.

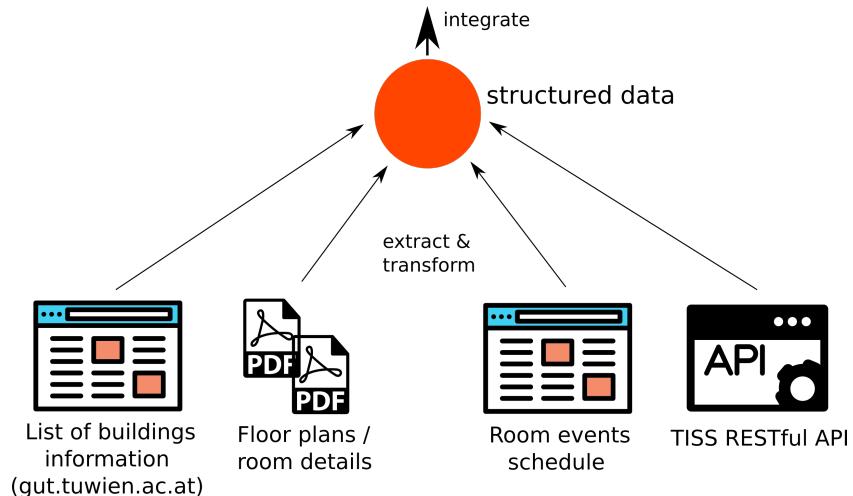


Figure 4.1: Visualization of the data sources

4.1.1 (X)HTML and PDF tables

As mentioned earlier, some of the data is presented to users only in form of (X)HTML tables and must therefore be extracted in order to be transformed into LD. Web scrapers/Web harvester are an useful tool to extract such tables from web pages manually or by using XPATH expressions and to store them in an intermediate representation. For this thesis a chrome extension named 'Web Scraper'³ was used. This intermediate representation was then exported into a Comma-separated values (CSV) file.

For tables in PDF files a different approach is required and the tool named 'Tabula'⁴ is qualified for this kind of problem. It can detect tables automatically and transforms them into an intermediate representation, but this does not always work appropriately. In this case, the user can support the program by selecting the relevant parts manually. This intermediate representation was then also exported into a CSV file.

After the information was extracted from web pages and PDF files and transformed into a structured format, it still was messy data that is not convenient for a further transformation into LD. In this case the tool OpenRefine⁵ is quite useful. It provides a set of abilities to manipulate messy, tabular data including an own expression language named GREL. LODRefine and RDF extension are extensions for OpenRefine, which enable the transformation of tabular data into RDF by mapping columns to classes and

³<http://webscraper.io/>

⁴<http://tabula.technology/>

⁵<http://openrefine.org/>

relationships between columns to properties. This mapping is then computed for each row.

4.1.2 Floor plans

Floor plans are ideal for describing the internal structure of a building, but they are mostly available in form of images and vector graphics. This representation is fine for humans, but not for machines; because the floor plan should not only be available as single image, but rather be the knowledge base for applications to find paths and locate rooms in the corresponding building.

Google Maps is a map application that also gives insight into certain buildings. It provides a tool⁶, where any user can submit floor plans and this floor plan then will be automatically transformed into an internal representation that can thenceforward be looked at in the application (if accepted). The problem is that the user cannot demand access to this internal representation of the plan, although the user may be the owner. Figure 4.3 shows the indoor plan of TU Vienna's main library in Google Maps.

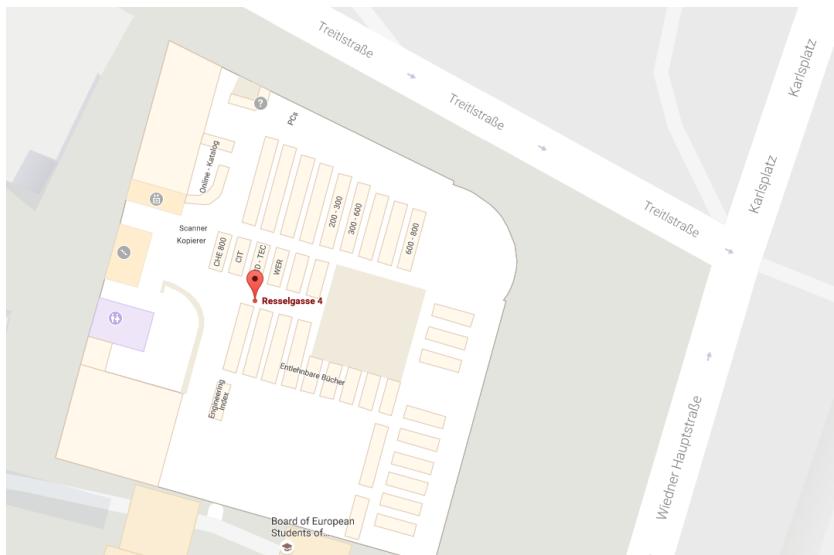


Figure 4.2: Google Maps: Indoor plan of TU Vienna's main library

Due to the lack of alternatives, the approach used for this thesis was to manually transform two floor plans of one building into shape files. The tool used for this approach was the Open Source software QGIS⁷. It has two plugins that are useful for the transformation, the 'Open Layers' plug-in to enable the use of OpenStreetMaps among others and the GDAL georeferencing plug-in to position a raster image on a map. Especially the last plug-in is important for positioning the floor plan to the correct position on the map, whereby

⁶<https://maps.google.com/floorplans/find>

⁷<http://www.qgis.org/>

the map should already contain the boundary of the building. After georeferencing the floor plan, a shape file of the rooms and other *features* of interest can be created. Each of the created shapes has a table of attributes with an unique ID that can be set by the user. The unique ID of each shape is an Internationalized Resource Identifier (IRI), so that the shape is ready to be transformed into LD. The resulting shape file can then be exported in multiple formats including GeoJSON⁸. For this thesis, a CSV file with the WKT serialization of each shape of the file and the corresponding ID per row showed to be the easiest way to integrate this data. Figure 4.3 shows a shape file of rooms on the basement floor of the informatics institute's building at the TU Vienna.



Figure 4.3: Room shapes file of the basement floor of the informatics institute's building

4.1.3 XML

The RESTful API of TISS provides results of methods in a (semi-)structured format like XML and JavaScript Object Notation (JSON), which is why no further preparations for the transformation are needed.

TODO: Describe ...

4.2 Prototype of ontology

In this section two ontologies will be proposed that are intended to cover the specific domain of spatial university data. Section 4.2.1 presents an ontology for describing the campus from buildings to indoor environments, whereas section 4.2.2 proposes an ontology for modelling indoor navigation. The design process took best practises for

⁸<http://geojson.org/>

publishing ontologies[5] into consideration. In order to enhance the interchangeability of the described data it is commonly understood to reuse terms of common ontologies, which is why chapter 3 outlined common geospatial ontologies as well as ontologies and approaches for modelling indoor environments. The conclusion of the evaluation was to use the GeoSPARQL ontology for describing the geometry of indoor features and buildings. Hence, the system architecture has to consider an additional prerequisite; the used triple store should support the GeoSPARQL extensions and for the purpose of better scalability also spatial indexing.

4.2.1 Spatial ontology

TODO: describe ...

4.2.2 Navigation ontology

The spatial ontology (see section 4.2.1) provides the necessary expressiveness to describe the indoor environment of buildings from rooms, corridors to entrances, but no assertions can be made about how these *features* are actually connected. Some conclusion can be drawn by looking at the *geometries* of *features*. In order to get for example the two rooms a door is connecting, the distance to the *geometry* of all rooms must be computed and the two rooms with the lowest distance by considering a reasonable threshold (in case one or both rooms are unknown) could be selected. Furthermore, would it be possible to compute the path to a room by considering the *geometries*, but this would still need information about how certain *features* are connected. A door for example could require an access control measure to pass it in one direction, but in the other direction it could be passed without constraints. A further problem to consider, if the path should be computed based on the *geometries*, is that not all paths may be accessible for all persons. A space could contain steps that are hard to overcome for persons with mobility impairments, or the space could be occupied by other obstacles. As an example, directly in front of or behind a door could be some steps that would force such persons to use another entrance.

For this thesis the approach of pre-defining routes based on *features* and points of interest was chosen for the purpose of simplification. However, the previous mentioned techniques could be used to compute certain parts of the route and humans can take over, when it became more complex. Another important question to address is the granularity. Should a route consist of a collection of connected rooms, horizontal passageways like corridors and vertical passageways like elevators or be more detailed with different paths inside of rooms and corridors. The problem of the first approach is that there could be some obstacles inside of rooms and corridors (e.g. steps in the middle of the corridor) that would be overseen. In case of the more detailed routes consisting of a collection of points of interest such obstacles can be taken into consideration, but it is consequently costlier than the first approach.

4. SOLUTION

Figure 4.4 shows the designed ontology for the indoor navigation. It has two major classes named `NavigationEntity` and `Constraint`. The class `NavigationEntity` has two subclasses namely `Route` and `Access`.

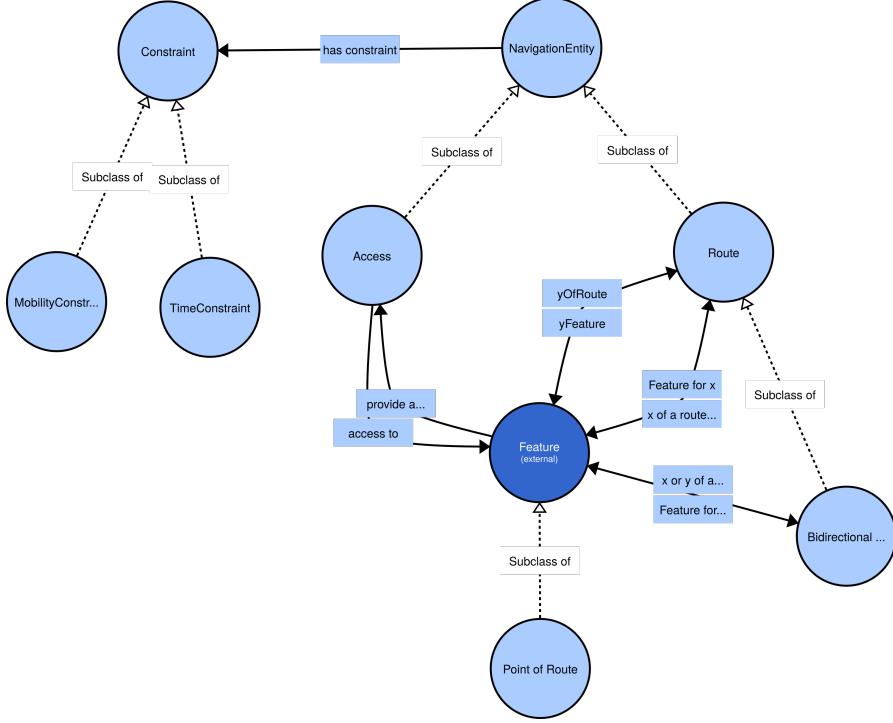


Figure 4.4: Visualization of the navigation ontology

Route is a path $R(x,y)$ between two spatial *features* x and y , whereby x and y are instances of the class `Feature` of GeoSPARQL. The path is walkable from x to y by considering the given constraints of this route. No assumptions about the other direction $R(y,x)$ are made. The x and y values of a route are assigned to it by the properties `xOfRoute` (inverse property `xFeature`) and `yOfRoute` (inverse property `yFeature`). A door that requires an access control measurement in one direction and not in the other is one example of a `Route`.

BiRoute is a subclass of `Route` and represents a path $R(x,y)$ between two spatial *features*, whereas the path between x and y is walkable in both directions considering the given constraints. The values x and y can be assigned to an instance of `BiRoute` by using the property `xyOfRoute` (inverse property `xyFeature`). This property is a subproperty of `xOfRoute` and `yOfRoute`; means that each value of `xyOfRoute` is automatically x and y of a `Route`. Consequently also $R(x,x)$ and $R(y,y)$ are valid paths, which must be filtered. A simple path inside of a corridor is an example for a `BiRoute`.

It would be nice to query for the path between two spatial *features* in pure SPARQL. The most recent version of the language (v1.1) introduced a new feature called property paths. These property paths can be used to recursively follow property links[12]. However, property paths lack expressiveness and queries that count the number of possible paths between two nodes or compute the shortest path are not expressible with SPARQL v1.1 [22][2]. Nonetheless, a SPARQL query could prepare the input for a path finding algorithm.

Listing 4.1: Prepare for path finding

```

SELECT ?xRouteFeature ?route ?yRouteFeature ?length WHERE
{
    values (?startFeature ?endFeature) {
        <http://finder.tuwien.ac.at/spatial/accessunit/id/HEEGMD>
        <http://finder.tuwien.ac.at/spatial/elevator/id/HFEG02C>
    }
    # Fetch the features of possible paths
    ?startFeature (navi:xOfRoute/navi:yFeature)+ ?routeFeature .
    ?xRouteFeature navi:xOfRoute ?route ;
        gsp:hasGeometry [ gsp:asWKT ?xRouteFeatureGeom ] .
    ?route navi:yFeature ?yRouteFeature .
    ?yRouteFeature (navi:xOfRoute/navi:yFeature)* ?endFeature ;
        gsp:hasGeometry [ gsp:asWKT ?yRouteFeatureGeom ] .
    FILTER(?xRouteFeature != ?yRouteFeature).
    # Computes the length of the path in meters.
    BIND(gspf:distance(?xRouteFeatureGeom, ?yRouteFeatureGeom, uom:metre)
        as ?length) .
    # Remove all routes that have certain constraints.
    FILTER NOT EXISTS {
        ?route navi:hasConstraint ex:building-id-H-opening-hours .
    }
}
    
```

Listing 4.1 shows a query returning a result set of all *features* (xRouteFeature) that are part of a possible path between two given *features* (startFeature, endFeature) and all the adjacent, noteworthy *features* (yRouteFeature) for each of them. An adjacent *feature* (yRouteFeature) is noteworthy, if there is a possible path to the target *feature* (endFeature) and the route R(xRouteFeature, yRouteFeature) has no undesirable constraints. Additionally, the length will be determined for each route in order to have a weight for the path finding algorithm. The length is simply the distance between the x- and y-*feature* of a route.

Access ...

TODO: finish ...

4.3 Architectural prototype

This section proposes a architectural prototype of a software system that manages LOD from integration to publication. Figure 4.5 shows an overview of this architecture, which is separated into three layers; the data management, service and presentation layer. The shown data sources (described in section 4.1) build the knowledge base of the system. Section 4.3.2 discusses the data management layer, which deals with the integration of those data sources and the maintenance of it. The linking of the local spatial LD with external sources like LinkedGeoData and GeoNames is described in section 4.3.3. After the local spatial data has been integrated into the system, the resulting LD can be published. Section 4.3.4 discusses different approaches for the publication with the goal to make the underlying LD easily accessible by machine and human agents. The prototype of a map application that make use of the provided services of the system is presented in the subsequent section 4.4.

In order to implement the suggested prototype the programming language Java was used and section 4.3.1 outlines the common frameworks for handling LD in this language. It also describes the reasoning, why RDF4J⁹ was chosen.

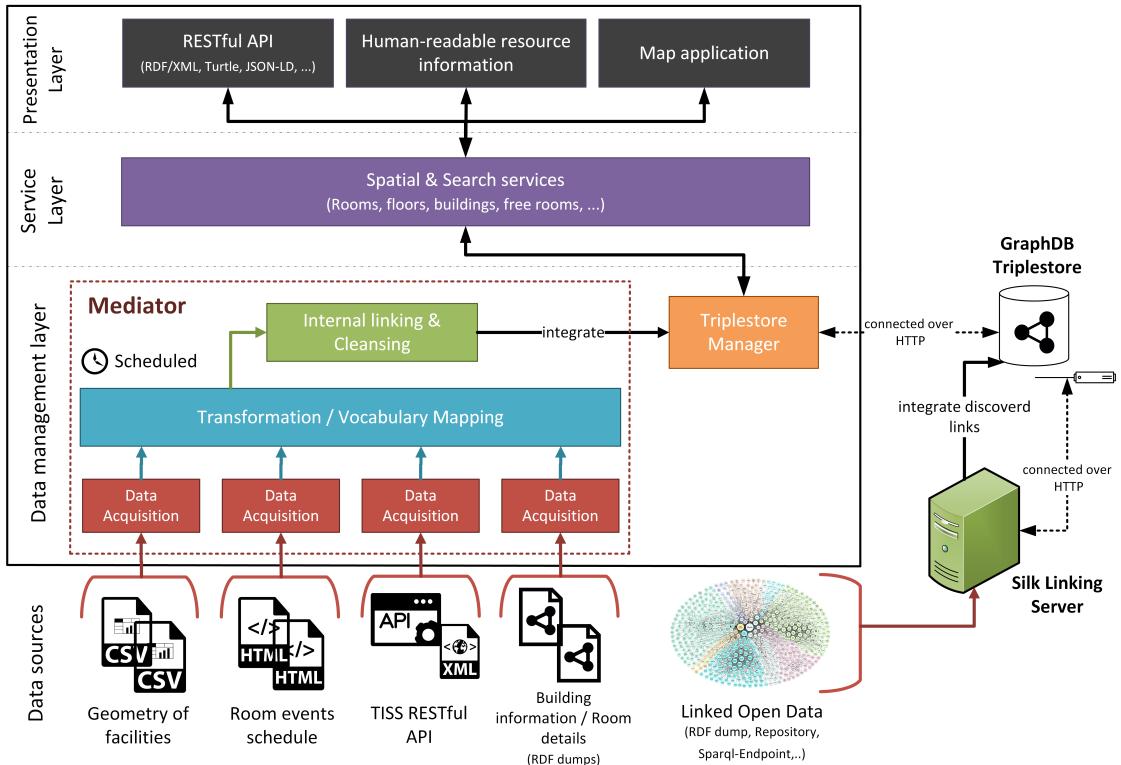


Figure 4.5: Overview of the system architecture

⁹<http://rdf4j.org/>

4.3.1 RDF Framework

Handling RDF data (as representation of LD) is an essential part of the proposed software system and a reliable framework for fulfilling this task is needed. Besides reliability the framework should also support most of the following requirements:

1. it should be freely available, favourably Open Source.
2. it should be possible to store data into either an embedded or external triple store.
3. it should support SPARQL to query data.
4. it should support common serialization formats like RDF/XML, Turtle, LD-JSON and N-Triples.
5. it should support spatial functionalities in SPARQL queries compliant to GeoSPARQL (see section 3.1). At least the minimal distance between two geometries must be computable.
6. it should be possible to reason over stored LD given a list of ontologies.

There are currently two common frameworks for handling RDF data in the Java environment namely Apache Jena¹⁰ and Sesame, which has recently been migrated to the project RDF4J¹¹. Both are Open Source, but none fulfils all of the mentioned requirements, why it is also necessary to take plug-ins and external triple stores into consideration.

Apache Jena is a framework that can be used to store, manipulate and query LD. It's architecture consists basically of three layers, the store API at the bottom, the inference API in the middle, the RDF API and SPARQL API respectively as interface to the application code [25]. In contrary to RDF4J, RDF data can be parsed into a model instance and thenceforward this model can be used to query and reason over it. RDF4J requires this RDF data to be at least stored into a repository embedded in memory in order be queried using SPARQL. One advantage of Jena over RDF4J is also the built-in reasoner that supports simple RDF Schema (RDFS) reasoning as well as multiple profiles of Web ontology language (OWL) [26]. Jena supports different types of storage from native triple store to a memory based one; also custom triple stores can be used, but in this case an adopter is required. In case of spatial functionalities Jena does support spatial relationships like near by, within box, intersect box or north/south/west/east of [27], but it is not compliant to the GeoSPARQL standard . Functions like computing the minimal distance between two geometries are not covered. There is also no support for GML serializations of geometries. As a consequence requirement **5** is not fulfilled. A way to circumvent this problem is to use an external triple store that provides these spatial functionalities, but as mentioned above this triple store must offer an adaptor for Jena, which limits the number of considerable triple stores.

¹⁰<https://jena.apache.org/>

¹¹<http://rdf4j.org/>

RDF4J is as well a framework to store, manipulate and query LD, but there are conceptual differences. As mentioned above in RDF4J actions are carried out against repositories and a model has less functionality than in Jena; it is simply a collection of statements (triples). RDF4J provides RDFS reasoning, but nothing more expressive than that. However, custom reasoners can be added like in Apache Jena too. An advantage of RDF4J over Jena is the simplicity of adding other storages besides the built-in ones (from native triple stores to in-memory repositories). It provides for example classes to use a SPARQL endpoint as repository. In fact there are multiple external triple stores that can be used with RDF4J over HTTP like GraphDB¹², Stardog¹³ or Blazegraph¹⁴. RDF4J does have a package named `rdf4j-queryalgebra-geosparql` to add spatial functionalities compliant with GeoSPARQL, but during testing the minimal distance function always returned null. It was not clear, if there was a bug in the framework or it was used incorrectly. This leads to a major drawback of RDF4J in contrast to Jena, it has a sparse documentation at the moment of writing.

Conclusion

The support of spatial functionalities in SPARQL queries like computing the minimal distance between two geometries is crucial for the proposed software system, not only for the navigation task, but also for spatial linking of resources with external sources (see section 4.3.3). In this regard both frameworks would be insufficient, but custom/external triple stores can help here. GraphDB is a triple store that is able to take care of reasoning (with support of RDFS as well as profiles of OWL). Additionally, it offers a GeoSPARQL plug-in that not only adds the corresponding spatial functionalities, but also spatial indexing resulting in a performance gain. This is the reason why RDF4J was chosen as framework, because it is easier to use such an external triple store, which then can also mitigate some of the drawbacks of RDF4J.

4.3.2 LD Management

At the bottom of the proposed, layered software system is the data management layer, which is responsible for integrating and maintaining LD. The first step is called mediation and it deals with the problem of integrating data from different data sources and updating them in a frequency that fits the characteristics of the dataset. The triple store manager controls the access to the used triple store. For this thesis, GraphDB was chosen as explained in the previous section 4.3.1, but the idea of the triple store manager is to hide this fact; also other storage solution that are compatible with RDF4J and fulfil the mentioned requirements can be used. All these aspects are explained in the following in more detail, starting with the designed naming scheme for resources.

¹²<http://ontotext.com/products/graphdb/>

¹³<http://docs.stardog.com/>

¹⁴<https://www.blazegraph.com/product/>

Naming scheme

In order to handle a larger amount of LD a consistent naming scheme for resources is inevitable. This also helps to simplify the design of a RESTful API (see section 4.3.4). As suggested by the first two principles of LD brought forward by Tim-Berners-Lee [4], unique IRIs should be used to name resources and be based on Hypertext Transfer Protocol (HTTP) so that people can look up those resources. A driving factor for exposing data as LD is the idea of lowering bars for innovative, motivated developers for coming up with useful applications using the exposed resources. Furthermore other datasets may link to those resources, which is why it must be considered that the designed IRIs are as stable and persistent as possible. However, this issue goes beyond the designed software system and has to include administrators of the infrastructure, where the system is going to be deployed. For this thesis the base IRI was chosen in advance to be <http://finder.tuwien.ac.at/>, but it is not hard-coded and can be changed before deployment.

The data that will be integrated into the system is predominately spatial data, but also organizational data and information about events are of interest for certain services. Information about events for instance can be used to determine which lecture room is available at a given time. All of this different types of data will be stored into different named graphs. Table 4.1 shows the chosen naming scheme for the named graphs. The catalog graph contains metadata about the different named graphs as suggested by best practises for data on the Web[29]. This metadata includes general information (like title, keywords, publisher, frequency of changes), license information, quality details and available forms of distribution. The data catalog vocabulary[17] was used to provide this information.

The suggested naming scheme of resources is then considering which named graph the resource belongs to. Table 4.2 shows the scheme for spatial resources. The IRIs are intended to be self-descriptive and it is assumed that each resource has an unique identifier. In case of rooms, buildings and floors it is easy, because the TU Vienna has already specified an unique identifier for those spatial entities and this is also reflected in the gathered data. Assigning an identifier to resources like addresses is more complicated, due to the lack of a predefined global identifier; this leads to the question how addresses that are considered to be equal, but gathered from different data sources, should be identified? One one hand a string consisting of a part that identifies the data source uniquely and a second part that identifies the resource inside of this data source uniquely can be used, whereby internal linking must be computed afterwards. The transformation process becomes easier with this approach, but consequently owl:sameAs statements will be produced during the linking, which could have an impact on the querying performance. The approach used for address resources in this thesis was a combination of concatenating key properties and hashing. This makes the transformation process more complicated, but prevents unnecessary owl:sameAs statements. The address identifier consists of the country code, the postal code, a MD5 hash of the street name and eventually the locator. The address '*Favoritenstraße 9-11, 1040 Vienna, Austria*' for example has the

4. SOLUTION

Name of graph	Description
<base>/spatial	Dataset containing statements about resources of the spatial domain.
<base>/organizational	Dataset containing statements about resources concerned with organizational information like persons, institutes and faculties.
<base>/event	Dataset containing statements about resources concerned with events like lectures.
<base>/catalog	Dataset containing statements describing the mentioned datasets (named graphs) like f.e. last update or frequency of changes. This dataset can be used as starting point to explore the datasets mentioned above and to inform yourself about available distributions.

Table 4.1: Named graphs used in the proposed software system

Naming of spatial resource	Description
<base>/spatial/building/id/<id>	Describes building with the given id.
<base>/spatial/building/id/<id>/address	Describes the default address of the building with the given id.
<base>/spatial/floor/id/<id>	Describes floor with the given id.
<base>/spatial/room/id/<id>	Describes room with the given id.
<base>/spatial/room/id/<id>/address	Describes the default address of the room with the given id.
...	

Table 4.2: Naming scheme for spatial resources

identifier AT1040-3ebd6229134d2d9c36a91657af166825-9-11.

Mediation

Data acquisition is the first step of the mediation process.

4.3.3 Linking

TODO: Description of strategies for linking with external sources and to enrich the data

(see figures 4.6 and 4.7)

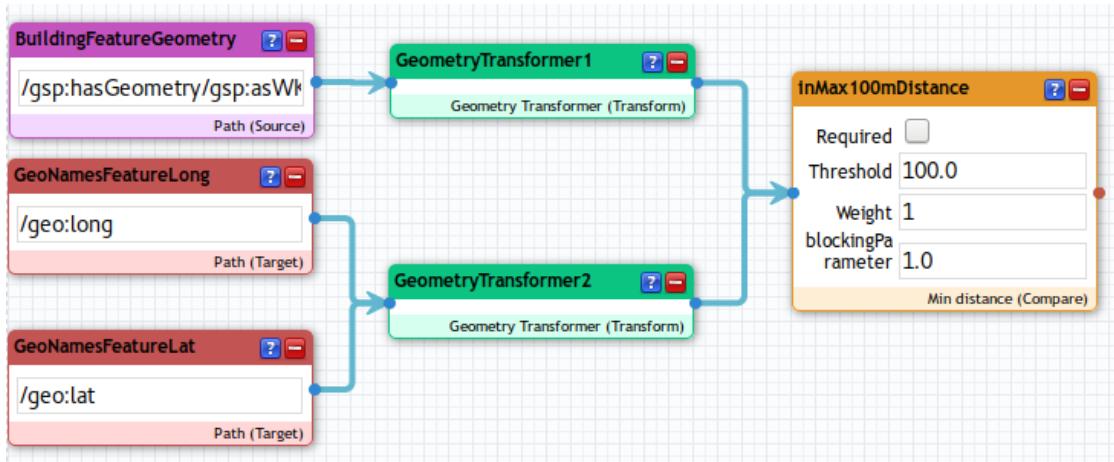


Figure 4.6: GeoNames features that are nearby (max. 100m)

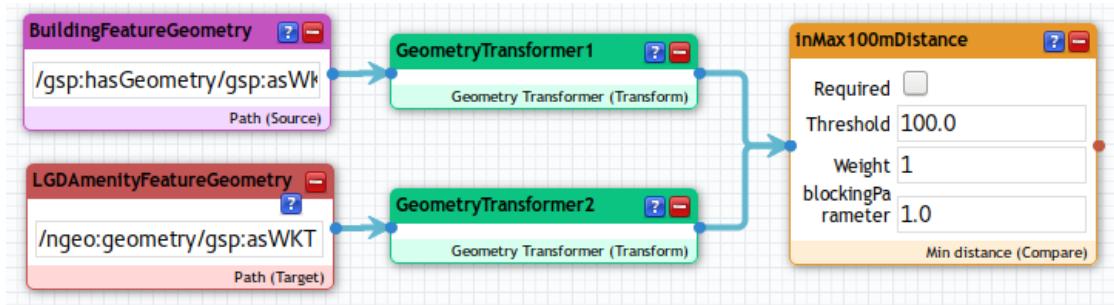


Figure 4.7: LinkedGeoData features that are nearby (max. 100m)

4.3.4 Linked Data publishing

TODO: Description of the data publishing and compliance to best practise

RESTful API

Content negotiation

4.4 Map application

TODO: A more comprehensive presentation of the user interface for humans to show the targeted use case.

4. SOLUTION

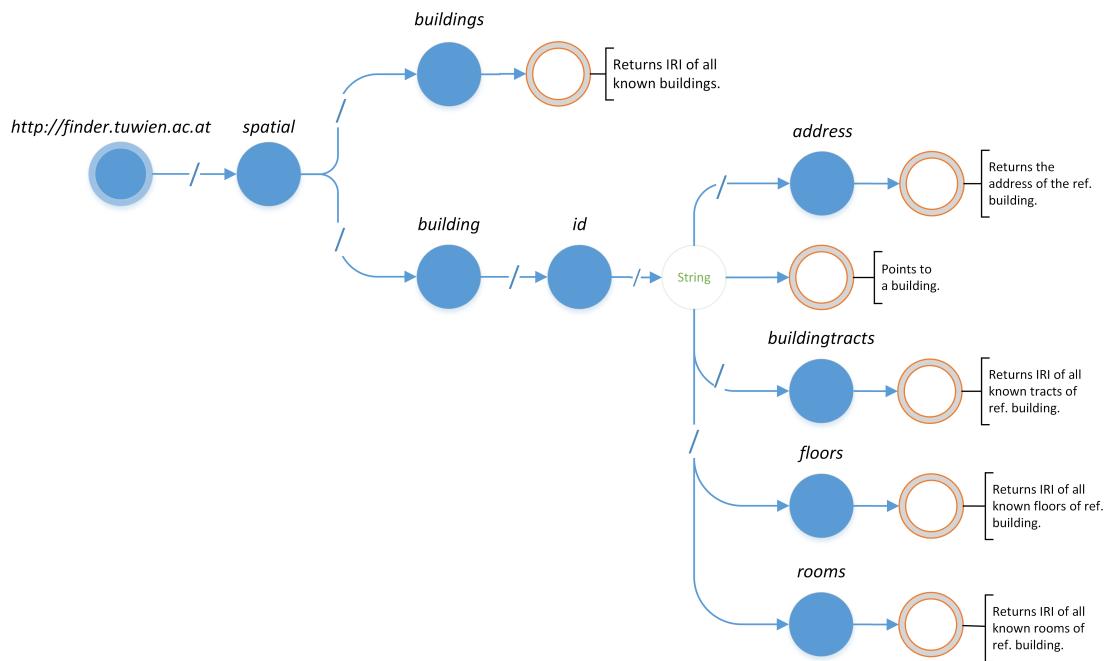


Figure 4.8: URI design for the services concerning buildings

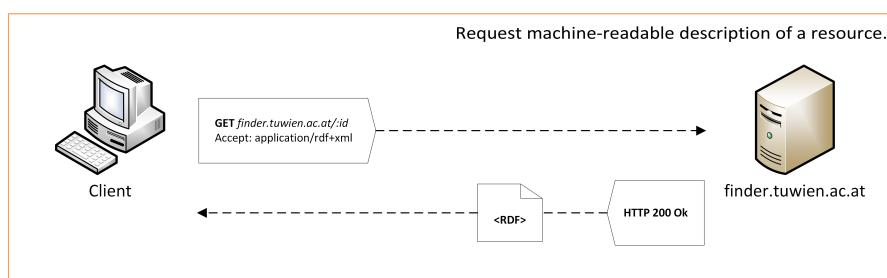


Figure 4.9: Requesting machine-readable representation of a given resource.

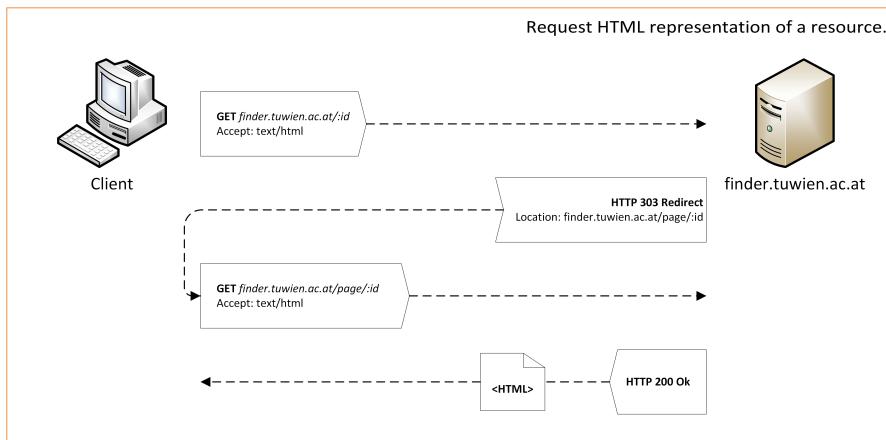


Figure 4.10: Requesting human-readable representation of a given resource.

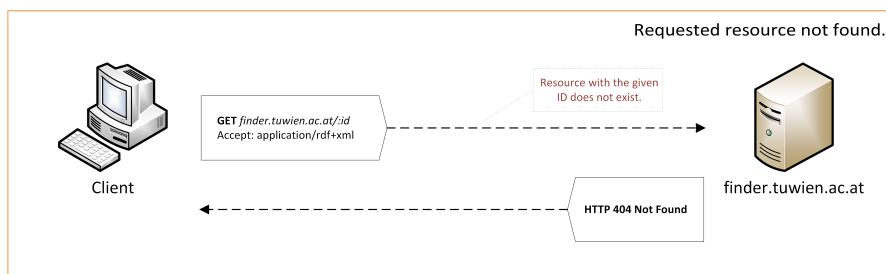


Figure 4.11: Requesting an unknown resource.

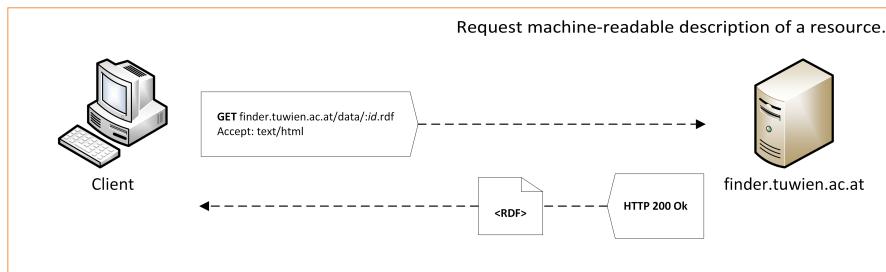


Figure 4.12: Forcing the return of a specific machine-readable representation.

CHAPTER 5

Discussion

5.1 Evaluation

TODO: Enter your text here.

5.2 Further work

TODO: Enter your text here.

5.3 Conclusion

TODO: Enter your text here.

List of Figures

1.1	Illustration of the Web of Documents	1
1.2	Illustration of the Web of Data	2
3.1	Visualization of GEO	8
3.2	Visualization of NGEO	10
3.3	Visualization of GeoSPARQL	11
3.4	Visualization of LOCN	13
3.5	Map application of University of Southampton	14
4.1	Visualization of the data sources	18
4.2	Google Maps: Indoor plan of TU Vienna's main library	19
4.3	Room shapes file of the basement floor of the informatics institute's building	20
4.4	Visualization of the navigation ontology	22
4.5	Overview of the system architecture	24
4.6	GeoNames features that are nearby (max. 100m)	29
4.7	LinkedGeoData features that are nearby (max. 100m)	29
4.8	URI design for the services concerning buildings	30
4.9	Requesting machine-readable representation of a given resource.	30
4.10	Requesting human-readable representation of a given resource.	31
4.11	Requesting an unknown resource.	31
4.12	Forcing the return of a specific machine-readable representation.	31

List of Tables

4.1	Named graphs used in the proposed software system	28
4.2	Naming scheme for spatial resources	28

List of Algorithms

Acronyms

- CSV** Comma-separated values. 18, 20
- GEO** Basic Geo (WGS84 long/lat) Vocabulary. 7–9, 11, 12, 27
- GML** Geography Markup Language. 10, 12
- IRI** Internationalized Resource Identifier. 20, 24
- LD** Linked Data. 2–4, 12, 14, 17, 18, 20
- LOCN** ISA Programme Location Core Vocabulary. 12, 13, 27
- LOD** Linked Open Data. 3, 7, 12, 13, 15, 17
- LOV** Linked Open Vocabulary. 7, 14
- NGEO** NeoGeo Geometry Ontology. 9, 10, 27
- PDF** Portable Document Format. 17, 18
- RDF** Resource description framework. 5, 9, 10, 12, 18, 23, 31
- RDFa** RDF in Attributes. 12
- RIF** Rule Interchange Format. 10
- ROOMS** Buildings and Rooms Vocabulary. 15
- SPARQL** SPARQL protocol and RDF query language. 5, 12, 22
- SPATIAL** NeoGeo Spatial Ontology. 9
- TU Vienna** Vienna University of Technology. 2–4, 17, 19, 20, 27
- WKT** Well-known text. 9, 10, 12, 20
- XML** Extensible Markup Language. 13, 20

Bibliography

- [1] Christos Anagnostopoulos, Vassileios Tsetsos, Panayotis Kikiras, and others. On-toNav: A semantic indoor navigation system. In *1st Workshop on Semantics in Mobile Environments (SME05), Ayia*. Citeseer, 2005.
- [2] Maurizio Atzori. Computing recursive SPARQL queries. In *Semantic Computing (ICSC), 2014 IEEE International Conference on*, pages 258–259. IEEE, 2014.
- [3] Robert Battle and Dave Kolas. Geosparql: enabling a geospatial semantic web. *Semantic Web Journal*, 3(4):355–370, 2011.
- [4] Tim Berners-Lee. Linked data-design issues (2006). *URL <http://www.w3.org/DesignIssues/LinkedData.html>*, 2009.
- [5] Diego Berraeta, Jon Phipps, Alistair Miles, Thomas Baker, and Ralph Swick. Best practice recipes for publishing RDF vocabularies. *Working draft, W3C*, 2008.
- [6] Christian Bizer, Tom Heath, and Tim Berners-Lee. Linked Data - The Story So Far. *International Journal on Semantic Web and Information Systems*, 2009.
- [7] D Brickley. Basic geo vocabulary. *W3C Semantic Web Interest Group. Available at: <http://www.w3.org/2003/01/geo>*, 2003. Accessed: 2016-09-26.
- [8] Richard Cyganiak and Anja Jentzsch. Linking open data cloud diagram. *LOD Community (<http://lod-cloud.net/>)*, 12, 2011.
- [9] Patrick M Dudas, Mahsa Ghafourian, and Hassan A Karimi. ONALIN: Ontology and algorithm for indoor routing. In *2009 Tenth International Conference on Mobile Data Management: Systems, Services and Middleware*, pages 720–725. IEEE, 2009.
- [10] Ramanathan Guha. Introducing schema.org: Search engines come together for a richer web. *Google Official Blog*, 2011.
- [11] Kevin Haas, Peter Mika, Paul Tarjan, and Roi Blanco. Enhanced results for web search. In *Proceedings of the 34th international ACM SIGIR conference on Research and development in Information Retrieval*, pages 725–734. ACM, 2011.

- [12] Steve Harris, Andy Seaborne, and Eric Prud'hommeaux. SPARQL 1.1 query language. *W3C Recommendation*, 21, 2013.
- [13] Bernadette Hyland, Ghislain Atemezing, and Boris Villazón-Terrazas. Best practices for publishing linked data. *W3C Working Group Note*, 2014.
- [14] Carsten Keßler and Tomi Kauppinen. Linked Open Data University of Münster–Infrastructure and Applications. In *Extended Semantic Web Conference*, pages 447–451. Springer, 2012.
- [15] Nemanja Kostic and Simon Scheider. Automated generation of indoor accessibility information for mobility-impaired individuals. In *AGILE 2015*, pages 235–252. Springer, 2015.
- [16] Steffen Lohmann, Stefan Negru, Florian Haag, and Thomas Ertl. VOWL 2: user-oriented visualization of ontologies. In *International Conference on Knowledge Engineering and Knowledge Management*, pages 266–281. Springer, 2014.
- [17] Fadi Maali, John Erickson, and Phil Archer. Data catalog vocabulary (DCAT). *W3C Recommendation*, 2014.
- [18] Barry Norton, Luis M. Vilches, Alexander De León, John Goodwin, Claus Stadler, Suchith Anand, and Dominic Harries. NeoGeo Vocabulary Specification - Madrid Edition. 2012.
- [19] Andrea Perego and Michael Lutz. ISA Programme Location Core Vocabulary. 2015.
- [20] Matthew Perry and John Herring. OGC GeoSPARQL-A geographic query language for RDF data. *OGC Implementation Standard. Sept*, 2012.
- [21] Research Group: AKWS. LODStats - Basic Statistics of the LOD cloud. Accessed: 2016-09-26.
- [22] Juan L Reutter, Adrián Soto, and Domagoj Vrgoč. Recursion in SPARQL. In *International Semantic Web Conference*, pages 19–35. Springer, 2015.
- [23] Barnabás Szász, Rita Fleiner, and András Micsik. iLOC–Building In-door Navigation Services using Linked Data. 2010.
- [24] Jeremy Tandy, Payam Barnaghi, and Linda van den Brink. Spatial Data on the Web Best Practices, January 2016. Accessed: 2016-09-26.
- [25] The Apache Software Foundation. Jena architecture overview. Accessed: 2016-11-12.
- [26] The Apache Software Foundation. Jena Ontology API. Accessed: 2016-11-12.
- [27] The Apache Software Foundation. Spatial searches with SPARQL. Accessed: 2016-11-13.

- [28] Csaba Veres and Eivind Elseth. Schema.org for the Semantic Web with MaDaME. In *I-SEMANTICS (Posters & Demos)*, pages 11–15. Citeseer, 2013.
- [29] W3C. Data on the Web Best Practices, August 2016. Accessed at: 2016-11-16.
- [30] Fouad Zablith, Mathieu d'Aquin, Stuart Brown, and Liam Green-Hughes. Consuming linked data within a large educational organization. In *Proceedings of the Second International Conference on Consuming Linked Data- Volume 782*, pages 85–96. CEUR-WS. org, 2011.