

# Integration, Quality Assurance and Usage of Geospatial Data using Semantic Tools Integration, Bewertung und Nutzung heterogener Datenquellen mittels semantischer Werkzeuge

Timo Homburg<sup>1</sup>, Claire Prudhomme<sup>1</sup>, Frank Boochs<sup>1</sup>,  
Ana-Maria Roxin<sup>2</sup>, Christophe Cruz<sup>2</sup>

{timo.homburg, claire.prudhomme,  
frank.boochs}  
@hs-mainz.de

{ana-maria.roxin,  
christophe.cruz}  
@u-bourgogne.fr

<sup>1</sup>Mainz University Of Applied Sciences  
Lucy-Hillebrand-Straße 2  
55128 Mainz, Germany

<sup>2</sup>Université de Bourgogne  
9 avenue Alain Savary  
21000 Dijon, France

**Abstract** In this article we want to present an integrational approach of geospatial data into the semantic web in the context of the Semantic GIS project. We first highlight the purpose and advantages of the integration and interpretation of data into the semantic web and further on describe the process of data acquisition, data interpretation, quality assurance and provenance and how to access the so integrated data. We continue to highlight the advantages of this integration method by presenting two fields of application of our research project: The evaluation of map data and the improvement of disaster management. We conclude the article by giving prospects of future work in our project.

In diesem Artikel stellen wir unsere Forschung in der Integration von Geodaten in einen Semantic Web Kontext in unserem Projekt Semantic GIS vor. Zunächst möchten wir den Zweck und die Vorteile einer Integration und Interpretation von Daten in das Semantic Web beleuchten und anschließend unseren Integrationprozess bestehend aus Datengewinnung, automatischer Interpretation, Qualitätssicherung und Provenance sowie den Datenzugriff erklären. Um die Anwendung unserer Forschung zu demonstrieren gehen wir auf zwei Anwendungsfälle in unserem Projekt ein: Die Bewertung von OpenStreetMap Daten und die Verbesserung des Katastrophenschutzes mittels semantischem Reasoning. Wir schließen den Artikel mit einem Fazit, sowie einem kurzen Ausblick auf zukünftige Forschung.

**Keywords:** Geospatial data, Linked data, Natural Language processing, Ontology, R2RML, SDI, Semantic web, Semantification, Data Quality, Provenance

**Keywords:** Geodaten, Linked data, Sprachverarbeitung natürlicher Sprachen, Ontologie, R2RML, SDI, Semantisches Web, Semantifikation, Datenqualität, Provenance

## 1 Introduction

Integration of heterogeneous datasets is a persisting problem in geographical computer science. Many classical GIS approaches exist making use of relational databases to achieve a tailor-made integration of geospatial data according to the needs of the current task. In the SemGIS project we are aiming at integrating heterogeneous geodatasets into a semantic web environment to take advantage of the flexibility of semantic data structures and to access a variety of related datasets that are already available in the semantic web. We intend to use the so-called geospatial knowledge base in the application field of disaster management in order to predict, mitigate or simplify decision making in an event of a disaster. As in our project we are possibly facing a large number of heterogeneous geodatasets of which we often do not know the origin nor intention nor the author and therefore lack an appropriate domain expert to help us understand data fields, we as non-domain experts would be left with a manual integration approach of said data. Dataset descriptions, if available, are often in natural language only which may give us hints but are hard to process in general and contain often hard to resolve ambiguities. However, despite mentioned obstacles we believe that a at least rudimentary classification and interlinking of our given data sets by means of the data values and data descriptions, is feasible. In addition, depending on the data source, data quality metrics as well as provenance information can be added to the to-be-imported data sets and change the way the data is treated not only for the geospatial community but also for the semantic web community. In this article we want to describe our approach to automatically find, process, analyse, interlink and quality-assure geospatial data sets on the web in the context of our project.

## 2 State Of The Art

The geospatial web provides several standards to distribute geospatial data. Since several years it is possible to publish geospatial data with the help of OGC webservice and to categorize said data using OGC catalogue web services (?). Despite this fact the access to geospatial data is very limited because it is not possible to search for geospatial data by means of their features and semantics and to make queries over geospatial data in the web on a large scale. This is due to several persisting problems in the publication of geospatial data:

- The scope of data is not semantically accessible using machines
- Geospatial data is not thematically clustered in the web of data
- Lack of a dedicated search engine for geospatial data
- Geospatial data is hard to index because of its heterogeneity

- Several non-quality annotated data sets depicting the same geometry and/or meaning and varying features might be found in the geospatial web
- Publications in the form of Map APIs like OpenStreetMap are if semantically interpretable only to a certain extent and to a limited amount of knowledge domains

We can conclude that there is no geospatial search engine nor a unified query interface being able to assess semantically interpreted and quality-assured geospatial webservices on a large scale. In addition many geospatial resources on the web are not even published as webservices or map APIs but in a variety of different formats and/or APIs (e.g. GeoTIFF, KML, GeoJSON, CovJSON) which in many cases are poorly documented and have often neither a quality assessment nor sufficient metadata about its source of origin.

## 2.1 Data Aquisition

To find thematic data in the geospatial web, traditional approaches are to find an appropriate CSW service which lists appropriate datasources that correspond to the description of the metadata or to its keywords. Recently, approaches to discover and link the geospatial web of data provided by services according to their keyword and service descriptions has been conducted (4). However, an automated analysis of features and their values was not provided in this work. In our work we would like to overcome this limitation or at least to test how far this limitation can be overcome in the geospatial domain.

## 2.2 Semantic Interpretability and Accessibility

Related work in the interpretation of geospatial data has been conducted for OpenStreetMap by the LinkedGeoData Project (8). Concepts of tags of OpenStreetMap data have been automatically generated and a virtual GeoSPARQL(5) interface accessing OSM data in realtime has been created. GeoSPARQL itself as an OGC standard provides us with a standardized method to access semantically interpreted geospatial data, so that in theory, the foundations to create a unified endpoint to search for features and types of geospatial data have been in place since its introduction in 2012.

## 2.3 Provenance and Data Quality

Provenance and Data Quality is of concern to the Semantic Web community as seen in (3) because the Semantic Web typically lacks such information appended to its knowledge bases. Semantic Web data are typically published without a rich provenance hierarchy and from various institutions without a record of trust and a history of how and in which quality the data has been gathered. However, concepts of which parameters to consider and which provenance information to gather can be found in repetitive standardized ontologies among others by W3C. ( ? ? ? ) In geospatial research data quality is defined in various standards such

as INSPIRE and in the respective literature (6; 7) which are useful in their respective applications. By integrating various kinds of heterogenous data, we intend to use as many quality criteria from those standards as possible to give endusers the possibility to choose among the criteria they deem fitting best.

### 3 SemGIS Project

In this section we describe how we implement the aforementioned steps of interpretation in the SemGIS project.

#### 3.1 Data Aquisition and discovery

To discover potential geospatial data in the SemGIS project we plan to develop a webcrawler similar to (4) which uses search engines and given addresses of geoportals to crawl for the data we would like to integrate. Along with the results we also allow lists of resources provided by users. We are prepared to integrate the following resources:

- Web services (WFS, WCS, WMS, CSW, SOS)
- Spatial Databases (PostGIS, Oracle Spatial)
- OGC data formats (SHP,KML,GML and dialects,GeoJSON,GeoTIFF)
- Using Geoparsing to make sense of geotagged webpages and/or web APIs

Once an appropriate resource has been found, the metadata found along with the resource has to be gathered as well if it is existent. For OGC webservices we can usually rely on given metadata standards of the webservice definition itself. Metadata for files might be stored along with them or on the surrounding homepage.

#### 3.2 Data Interpretation

Once a list of suitable geospatial resources has been gathered from the web or has been provided by the enduser, the data sets need to be interpreted to link them to Semantic Web concepts. Every aforementioned data source can be seen as one or many relational database tables, which we depict generically as shown in table 1.

ID	the_geom	Feature1	Feature 2	... FeatureN
Example.1	POINT(..)	123	"ExampleString"	3.4

Table 1: Example File "example" represented as a database table

When interpreting data from a spatial database, foreign keys aka. relations between existing database tables can be considered and extracted using established

methods like R2RML (?). Geodatasets in the form of files represent single relational database tables of which it is typically unknown which relations to other data sets exist. It is on this premise that we employ interpretation algorithms to create such relations to Semantic Web concepts. The goal of such efforts is to produce a so-called local ontology (figure 1) of each resource with induced links to other Semantic Web resources and concepts.

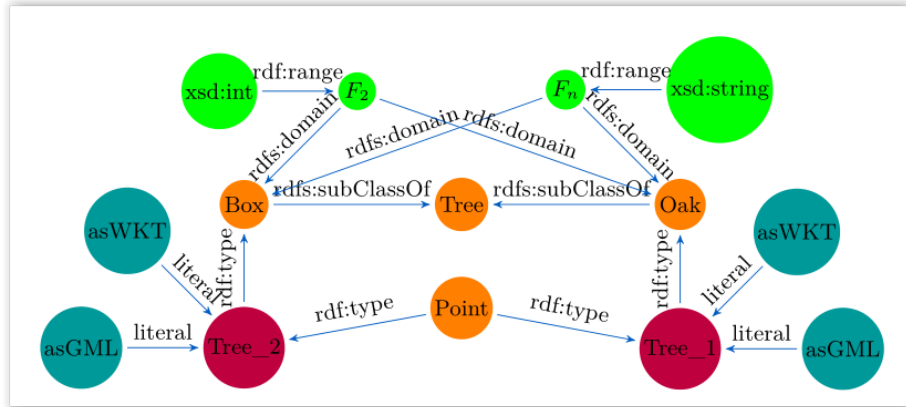


Figure 1: Local Ontology Example

**3.2.1 Concept Matching Process** The information we would like to extract from a single database table data set includes

- At least one concept for the whole data set
- At best one or more concepts per column of the database table
- At best several additional attributes using additional knowledge sources e.g. geocoding information

We can extract at least one concept for the whole data set by analyzing its filename/database table name or by using a reference data set of geometries e.g. OpenStreetMap/LinkedGeodata to find classifications of geometries in the vicinity of the geometries provided in the data set. Fitting concepts for columns can be found by either analyzing the columns' title and/or the values of the column using Natural Language Processing algorithms. To do that we rely on BabelNet (?), a multilingual network, partially connected to the Semantic Web and on the labels of ontologies we would like to link to e.g. DBPedia or Wikidata (?). Using a rudimentary detection of the language the dataset is written in we can analyze column titles and values for existing terms that we can subsequently match in the mentioned ontologies. Having sufficiently many values of a similar classification allows a generalization of a concept description for the respective

column. We are therefore able to detect at best one concept per column and if modelled the instance of each value present in the corresponding dataset. By doing further linguistic analysis we are able to detect the role of each column, which might be:

- A foreign key corresponding to an ObjectProperty in the Semantic Web
- A DataTypeProperty corresponding to a simple value
- An AnnotationProperty corresponding to metadata annotated to an instance or class in our knowledge base

Using this additional information we further interpret and distinguish columns by the following categories:

- Address columns: Columns that represent components of an address matchable with traditional geocoding
- SubClass columns: Columns including nouns that represent a subcategorization of the database tables content
- Object Property Columns: Columns including adjectives that represent a categorization of a relation or an attribute of the data set
- Common regular expression columns: Columns that can be associated by executing a common regular expression on its values (e.g. email addresses, UUIDs)
- Label and comment columns: Columns that represent a description of one row(=instance of the data set)
- Unit columns: Columns containing numbers which have an identifiable unit and/or concept when analyzing the column description

We are currently not able to analyze remaining columns so that they remain in the system as associated values in its primitive types (e.g. double, integer). The end-user is still able to access them, but the semantic meaning could not be determined automatically and is therefore not accessible if not corrected by a human being. Our goal is to improve the automated concept detection in further iterations of our software.

### 3.3 Quality and Provenance

Quality and Provenance are important metadata which can enhance the value of a data set for its daily usage. The existence of quality and provenance parameters in the geospatial web and in the semantic web is often not standardized and not common. We therefore propose to extract and generate such parameters in all data sets we integrate into a common knowledge base.

**3.3.1 Provenance** Provenance parameters can tell us information about who was when providing which data using which process of data preparation and which original data source or measurements. Usually provenance information can be found along within accompanying metadata or on the homepage/service page where a particular data set has been published. Provenance information can be modeled using the Prov-O ontology defined by W3C.(?) Examples of such provenance information publication are as follows:

- Provenance information of a shapefile: Creator, GPS measurement device, Measurement method, date of creation, date of modification etc.
- Provenance information of the publishing institution: Name, Email phone number etc.
- Provenance information of the publishing service: Domain, name, contact data, maintainer etc.

**3.3.2 Data Quality** The notion of data quality can be extended to various data quality dimensions. One definition of data quality could be

Data quality is the degree to which data fulfills requirements.

Which requirements are important for the data we are working with depends on its usecase. Every domain of knowledge can depend on various quality criteria. However, we can analyze as many quality criteria on our data as it is possible to prepare users to take qualified decisions about which data to use for their specific usecase. Examples of data quality parameters include:

- Positional accuracy of the geometry (with reference to a goldstandard)
- Geocodability of the data set
- Semantic Interpretability of the concept
- Completeness
- Open License/Cost of access
- Quality of Service

We are hereby focusing on known data quality concepts from the Semantic Web, GIS research as well as data quality provided by the knowledge domains we are connected to through features.

**3.3.3 Evaluating Provenance and Data Quality** When combining provenance information and quality parameters, datasets from specific resources can be associated with specific values of data quality. This allows not only to rank specific data sets but also to highlight data providers that are trustworthy because they have proven to provide data with a consistent data quality. If in doubt a reasoning system or the end user can take advantage of this information to choose the most trustworthy data set among several possible data sets for the fulfilment of his use case.

### 3.4 Data Access and Reasoning

To access data we have imported using the process described in the previous sections we rely on a GeoSPARQL (5) endpoint which allows us to use Egenhofer calculations in the semantic web. In addition we developed an extended vocabulary allowing us to use various PostGIS functions like geometry constructors to be used in GeoSPARQL. Queries that are often used or that lead to results that should be reused in a later stage of the development are standardized in so-called

reasoning rules in languages like SWRL (?) or SPIN (?). At this stage of the project we are at the point of developing reasoning strategies together with our projects partners. Therefore first realworld applications of reasoning are yet to be implemented in our research. To highlight a possible case of reasoning we refer to an example from (2) in which we highlighted the inference of nearest hospitals to a to-be-evacuated school as an example of automated reasoning in a disaster management case.

## 4 Applications

The SemGIS project is aimed at applications in disaster management and energy. However correct application cases also require trustworthy and correct map data which needs to be ensured while executing the use case calculations or beforehand.

### 4.1 Evaluation of OpenStreetMap Data

The largest repository of open geodata in the world is OpenStreetMap. It is used by various people around the world for many different purposes and is created by a vast amount of editors. To our knowledge a comprehensive analysis of the quality and provenance of OpenStreetMap data in Germany has not been undertaken yet. Therefore in our project, we would like to evaluate OpenStreetMap data by comparing them to the goldstandard provided by the German national authorities for cartography and geodesy. By semantically interpreting and by extracting and adding provenance as well as quality information we can compare German official data to OpenStreetMap data in as many aspects as needed. We can highlight conflicts in a separate layer on top of OpenStreetMap, evaluate which parts of OpenStreetMap are good enough to serve for which usecase we are aware of and can give hints to the OpenStreetMap community in which way to improve OpenStreetMap in the future. A preliminary development of this approach can be seen in figure 2 While doing so we also create a huge amount of quality-annotated data in the Semantic Web. This data serves the Semantic Web community which becomes increasingly interested in geospatial topics as well. In the context of disaster management we ensure that the resources we use to do flood simulations are correct to the extent we need them, so that predictions of flood and the consequences thereof are accurate. Lastly comparing datasets of different quality helps us to consolidate different features that are present in the different datasets. By knowing quality and provenance requirements of the enduser the system can end up with a merged dataset of high quality and the maximum amount of features possible.

### 4.2 Multi-Agent Natural Disaster Simulations

Disaster Management consists of various steps that can be highlighted in the so-called disaster management cycle (1). During an event of disaster for example



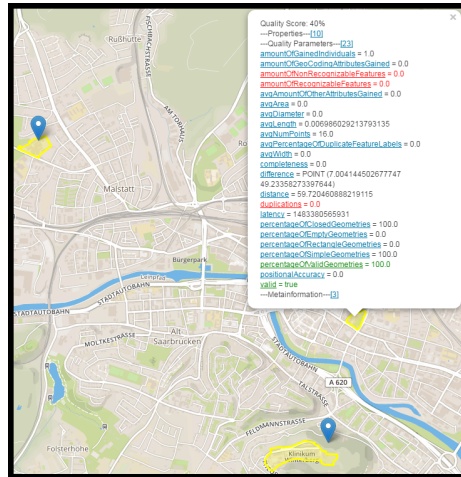


Figure 2: Preliminary Quality Comparison Screen of OpenStreetMap data vs. OpenData

a flood, various actors need to cooperate in order to prevent further damages, evacuate people and rescue endangered areas. The efficiency of these activities depends on many elements which need to be prepared. Three of these elements are the resources needed for this activities, activity planning and common data set shared by all actors. Each of these elements have an impact in the disaster management response. The activity planning improves the organisation and allows for knowing what you need to do according to the situation, and thus, to act quickly. Resources are key elements for the activity. If the resources are not enough to achieve the activities goal, the activity may be slowed down or even fail. The coordination between different actors becomes simplified when they are able to work with a common data set. A main problem in this field is that all actors do not have same rights and the same access of data. The identification of a data sets which could be used as a common data set for all actors (even if some of them can have more information) would be a good point for the coordination of the response activities. In order to assess these three elements, our project has aiming to simulate agents corresponding to real persons acting in a disaster event according to a rule-based system using gathered and interpreted data as described above our project. The simulation has aiming to support the preparation of disaster management response in assessing activity planning, resources and data sets.

## 5 Conclusion

Working towards a unified endpoint for semantically interpreted and quality assured geospatial data is a profitable approach for both the geospatial web of data as well as for the Semantic Web. In this article we have shown our efforts

on how to approach this goal and the progress we have achieved on the way. We have also shown how provenance and data quality parameters can be used in our system in the future to evaluate and append other sources of open data like OpenStreetMap or to act as a beneficial knowledge base for disaster management optimizations using Multi-Agent simulations. Our future work will continue on said use cases with our project partners and to investigate on how our concepts will help to improve the workflows of the several actors in disaster management.

## **6 Acknowledgements**

The SemGIS project was funded by the German Federal Ministry of Education and Research under Project Reference: 03FH032IX4.

## References

- Coppola, D.P.: Introduction to International Disaster Management. Elsevier (2011) [4.2](#)
- Homburg, T., Prudhomme, C., Würriehausen, F., Karmacharya, A., Boochs, F., Roxin, A., Cruz, C.: Interpreting heterogeneous geospatial data using semantic web technologies. In: International Conference on Computational Science and Its Applications. pp. 240–255. Springer (2016) [3.4](#)
- Mendes, P.N., Mühleisen, H., Bizer, C.: Sieve: Linked data quality assessment and fusion. In: Proceedings of the 2012 Joint EDBT/ICDT Workshops. pp. 116–123. EDBT-ICDT '12, ACM, New York, NY, USA (2012), <http://doi.acm.org/10.1145/2320765.2320803> [2.3](#)
- Pellicer, F.J.L.: Semantic linkage of the invisible geospatial web. Ph.D. thesis, Universidad de Zaragoza (2011) [2.1](#), [3.1](#)
- Perry, M., Herring, J.: Ogc geosparql-a geographic query language for rdf data. OGC Implementation Standard. Sept (2012) [2.2](#), [3.4](#)
- Redman, T.C.: Data quality: the field guide. Digital press (2001) [2.3](#)
- Shi, W., Fisher, P., Goodchild, M.F.: Spatial data quality. CRC Press (2003) [2.3](#)
- Stadler, C., Lehmann, J., Höffner, K., Auer, S.: Linkedgeodata: A core for a web of spatial open data. Semantic Web 3(4), 333–354 (2012) [2.2](#)