



2015-ENST-0022



EDITE - ED 130

## Doctorat ParisTech

### THÈSE

pour obtenir le grade de docteur délivré par

**TELECOM ParisTech**

**Spécialité « Science de l'information de la communication »**

*présentée et soutenue publiquement par*

**Ghislain Auguste ATEMEZING**

le 10 avril 2015

### Publishing and Consuming Geo-Spatial and Government Data on the Semantic Web

Directeur de thèse : **Raphaël TRONCY**

#### Jury

M. Sören AUER, Professeur, IAIS, Université de Bonn, Allemagne	Rapporteur
Mme. Chantal REYNAUD, Professeur, INRIA Saclay, Université de Paris XI, France	Rapporteur
M. Roberto GARCIA, Maître de Conférences associé, Universitat de Lleida, Espagne	Examinateur
M. Andreas HARTH, Maître de Conférences, Institut AIFB, Karlsruhe, Allemagne	Examinateur
Mme. Elisabeth METAIS, Professeur, CNAM - Equipe ISID, France	Examinateur
M. Sébastien MUSTIERE, HDR, COGIT-IGN, France	Invité

T  
H  
È  
S  
E

*To all those who helped me  
to make this dream coming true.*

*“Things don’t have to change  
the world to be important”.*

-Steve Jobs -



# Acknowledgements

This thesis is the result of part of my work carried out in the context of two Projects: the DATALIFT Project (ANR-10-CORD-009), and the APPS4EUROPE Project. There are many people I want to thank since they kindly supported me in so many different ways for the successful completion of this thesis.

I am indebted to my thesis advisor Dr. Raphaël Troncy for giving me the opportunity for a PhD. at EURECOM / Telecom ParisTech. Throughout my PhD he provided helpful ideas and encouraging support in this exciting domain of Semantic Web .

I would like to thank my committee members, the reviewers Prof. Sören AUER and Pr. Chantal REYNAUD, and furthermore the examiners Pr. Elisabeth METAIS, Dr. Andreas HARTH and Dr. Roberto GARCIA for their precious time, shared positive insight and guidance. A special thanks to Dr. Sébastien MUSTIERE for accepting our invitation.

I would like to express my deepest appreciation for all the colleagues and partners of Datalift project for all the exchanges during the duration of the project. I would like to thank specially Bernard, Pierre-Yves, Nathalie, Laurent. It was a pleasure to work and exchange with them.

My sincere thanks are due to all the members of the W3C Government Linked Data Working Group, specially to Bernadette Hyland, Boris Villazón-Terrazas, Richard Cyganiak, Dave Reynolds and Phil Archer. We had many intensive moments of exchange and collaborations.

I would like also to show my sincerest gratitude to all my mentors at OEG/UPM (Madrid) for providing me first skills on ontology engineering: Pr. Asuncion Gómez and Dr. Oscar Corcho. Their enthusiasm and dedication to research on this field were truly inspiring.

My warmest thanks to my colleagues who supported me during my Ph.D. Precisely, I would like to thank Houda, Giuseppe, José Luis, Vuk and Ahmad. Also, I thank all those working at EURECOM, they made my stay at EURECOM very pleasant. Thanks also to Jodi Schneider for her insightful comments to this document.

I would like to express my deepest thanks to my wife, Verónica Álvarez Aguilar, who gave me all her patience and comprehension since the beginning of this adventure. I owe my family my profound gratitude because they always supported and believed in me: specially my parents Genevieve Djifack and Prosper Tabondjou; and all my brothers and sisters: Stephanie, Judith, Arlette, Mathias, Alvine, Edwige, Yannick. Lastly, special thanks to my friends for their unwavering friendship, moral and infinite support.



# Abstract

Over the past few years, the domain of Open Data has received an increasing attention from public administrations. Potential benefits for citizens include more transparency in the decision making, a better governance and a virtuous development of a digital eco-system that would create added-value apps when processing and analyzing open data. However, opening up and publishing data is not enough to create this data value chain due to a number of challenges such as the heterogeneity of formats (XML, CSV, Excel, PDF, Shape Files), the variety of access methods (API, database, dump) and the lack of nomenclature that would enable to better re-use and interconnect datasets. In this thesis, we explore how semantic web technologies can be used to tackle the research problems related to the integration and consumption of geo-spatial data.

This thesis applies the Linked Data principles in the domain of geographic information (a key domain for open government). In particular, we address three key challenging problems in the publishing workflow of geospatial open data publication and consumption, with real world use cases from the French National mapping agency (IGN): (1) How to efficiently represent and store geospatial data on the Web to ensure interoperable applications? (2) What are the best options for a user to interact with semantic content using visualizations? (3) What are the mechanisms that support preserving structured data of a high quality on the Web?

Our contributions are thus break down into three parts with applications in the geographical domain. We propose and model four vocabularies for representing coordinate reference systems (CRS), topographic entities and their geometries. These ontologies extend existing vocabularies and add two additional advantages: an explicit use of CRS identified by URIs for geometry, and the ability to describe structured geometries in RDF. We have contributed to the development of the Datalift platform that aims to support lay users in the process of "lifting" raw data into RDF. We have published the French authoritative database GEOFLA using this tool and we provide a systematic evaluation of the performances of the most used endpoints when dealing with spatial queries.

Regarding the consumption of linked data, after reviewing different categories of visualization (generic and Linked Data specific), we propose a vocabulary for Describing Visualization Applications (DVIA). We formalize and implement a novel workflow for visualizing datasets with the LDVizWiz tool: a Linked Data Visualization Wizard.

The last part of the thesis describes contributions to the Linked Open Vocabularies (LOV) catalogue: it shows how LOV can be used with an ontology modeling methodology (e.g. the NeOn methodology) to improve reuse of vocabularies. We propose an heuristic to align vocabularies and a ranking of vocabularies based on information content (IC) metrics. Finally, the thesis provides answers on how to check the license compatibility between vocabularies and datasets in the publication

workflow. Through this thesis, we demonstrate the benefits of using semantic technologies and W3C standards to improve the discovery, interlinking and visualization of geo-spatial government data for their publication on the Web.

# Contents

Abstract . . . . .	iii
Contents . . . . .	vii
List of Figures . . . . .	xi
List of Tables . . . . .	xiii
List of Listings . . . . .	xv
List of Publications . . . . .	xvii
Acronyms . . . . .	xx
<b>1 Introduction</b>	<b>1</b>
1.1 Context . . . . .	1
1.1.1 Geographic Information (GI) . . . . .	3
1.1.2 Geographic Data in the Linked Data Cloud . . . . .	4
1.2 Research Questions . . . . .	7
1.3 Contributions . . . . .	10
1.3.1 Modeling, Publishing and Querying geodata . . . . .	10
1.3.2 Visualization Tools in Linked Government Data . . . . .	11
1.3.3 Contributions to Standards . . . . .	12
1.4 Thesis Outline . . . . .	13
<b>I Modeling, Interconnecting and Generating Geodata on the Web</b>	<b>17</b>
<b>2 Geospatial Data on the Web</b>	<b>19</b>
2.1 Geographic Information . . . . .	20
2.1.1 Specificity . . . . .	21
2.1.2 Data Formats and Serialization . . . . .	22
2.2 A REST Service for Converting Geo Data . . . . .	24
2.2.1 Background . . . . .	24
2.2.2 Why a REST service is needed . . . . .	25
2.2.3 API Access and Implementation . . . . .	26
2.3 Current Modeling Approach . . . . .	30
2.3.1 Status of Vocabularies Usage for Geospatial Data . . . . .	30
2.3.2 Geospatial Vocabularies . . . . .	34
2.3.3 Georeferencing Data on the Web . . . . .	35
2.4 Vocabularies for Geometries and Feature Types . . . . .	39
2.4.1 Motivation . . . . .	40
2.4.2 A vocabulary for geometries . . . . .	40
2.4.3 A vocabulary for Topographic Entities . . . . .	42
2.4.4 Discussions . . . . .	45
2.5 Summary . . . . .	46

<b>3 Publishing, Interlinking and Querying Geodata</b>	<b>47</b>
3.1 Current Representation of Geodata on the Web . . . . .	48
3.1.1 DBpedia Modeling . . . . .	48
3.1.2 Geonames Modeling . . . . .	49
3.1.3 LinkedGeoData Modeling . . . . .	49
3.1.4 Discussion . . . . .	49
3.2 Existing Tools for Converting Geospatial Data . . . . .	50
3.2.1 Geometry2RDF . . . . .	50
3.2.2 TripleGeo . . . . .	51
3.2.3 shp2GeoSPARQL . . . . .	51
3.2.4 GeomRDF: Datalift tool for Converting Geodata . . . . .	51
3.2.5 Limitations of existing tools . . . . .	53
3.3 Interlinking Geospatial Vocabulary and Data . . . . .	53
3.3.1 Criteria for interlinking geospatial Data . . . . .	54
3.3.2 Functions for Comparing Geometries . . . . .	55
3.3.3 GeOnto Alignment Process Scenario . . . . .	57
3.4 Survey on Triple Stores and Workbench . . . . .	57
3.4.1 Generic Triple Stores . . . . .	58
3.4.2 Geospatial Triple Stores . . . . .	59
3.4.3 Assessing Triple Stores . . . . .	60
3.4.4 Workbench for Geospatial Data . . . . .	61
3.5 Datalift: A tool for Managing Linked (Geo)Data Publishing Workflow	63
3.5.1 Functionalities of the Datalift platform . . . . .	63
3.6 Publishing French Authoritative Datasets . . . . .	67
3.6.1 Publishing French Administrative Units (GEOFRA) . . . . .	68
3.6.2 Publishing French Gazetteer . . . . .	69
3.6.3 Publishing Addresses of OSM-France in RDF . . . . .	70
3.6.4 Status of French LOD cloud (FrLOD) . . . . .	73
3.7 Evaluation of Spatial Queries . . . . .	74
3.7.1 Querying LinkedGeodata . . . . .	75
3.7.2 Querying FactForge (OWLIM) . . . . .	76
3.7.3 Querying Structured Geometries . . . . .	77
3.8 Summary . . . . .	77
<b>II Generating Visualizations for Linked Data</b>	<b>79</b>
<b>4 Analyzing and Describing Visualization Tools and Applications</b>	<b>81</b>
4.1 Survey on Visualization Tools . . . . .	82
4.1.1 Tools for visualizing Structured Data . . . . .	82
4.1.2 Tools for visualizing RDF Data . . . . .	84
4.1.3 Discussion . . . . .	87
4.2 Describing Applications on the Web . . . . .	90
4.2.1 Motivation . . . . .	90

4.2.2	Catalogs of Applications . . . . .	90
4.3	Describing and Modeling Applications . . . . .	92
4.3.1	Typology of Applications . . . . .	92
4.3.2	Reusable applications . . . . .	93
4.3.3	A vocabulary for Describing VIualization Applications . . . . .	94
4.4	Summary . . . . .	96
<b>5</b>	<b>Creating and Generating Visual Applications</b>	<b>97</b>
5.1	Wizard for Visualizations: Theoretical foundations . . . . .	98
5.1.1	Dataset Analysis . . . . .	98
5.1.2	Mapping Datatype, Views and Vocabularies . . . . .	99
5.2	LDVizWiz: a Linked Data Visualization Wizard . . . . .	100
5.2.1	Workflow . . . . .	100
5.2.2	Implementation and Evaluation . . . . .	104
5.3	Finding Important Properties for an Entity . . . . .	107
5.3.1	Reverse Engineering the Google KG Panel . . . . .	107
5.3.2	Evaluation . . . . .	108
5.4	GeoRDFviz: Map visualization of Geodata Endpoints . . . . .	110
5.4.1	Back-End Description . . . . .	110
5.4.2	Front-End Interface . . . . .	112
5.5	An application consuming event datasets: Confomaton . . . . .	113
5.5.1	Background and Motivation . . . . .	113
5.5.2	Collecting and Modeling Data in Confomation . . . . .	115
5.5.3	Graphical User Interface . . . . .	119
5.6	Application consuming statistical datasets . . . . .	120
5.6.1	Dataset Modeling . . . . .	120
5.6.2	Interconnection . . . . .	123
5.6.3	User Interface . . . . .	123
5.7	Improving the discovery of applications contests in Open Data Events	124
5.7.1	Background . . . . .	124
5.7.2	Modeling Approaches . . . . .	124
5.7.3	Implementation and Application . . . . .	127
5.7.4	Discussion . . . . .	132
5.8	Summary . . . . .	132
<b>III</b>	<b>Contribution to Standards</b>	<b>135</b>
<b>6</b>	<b>Contributions to Linked Open Vocabularies (LOV)</b>	<b>137</b>
6.1	Catalog of Vocabularies . . . . .	139
6.2	Linked Open Vocabulary (LOV) and Vocabularies . . . . .	144
6.2.1	Linked Open Vocabularies . . . . .	145
6.2.2	LOV vs NeOn Methodology . . . . .	147
6.2.3	Prefixes harmonization . . . . .	149

6.3	Vocabulary Ranking Metrics . . . . .	158
6.3.1	Information Content Metrics . . . . .	159
6.3.2	Information Content in Linked Open Vocabularies . . . . .	159
6.3.3	Ranking Vocabularies using The Information Content . . . . .	160
6.4	Datalift Module for Selecting Vocabularies . . . . .	165
6.5	Summary . . . . .	166
<b>7</b>	<b>License Compatibility Checking</b>	<b>169</b>
7.1	Background . . . . .	169
7.2	Statistics about licensed vocabularies . . . . .	170
7.3	Related work about licenses in the Web of Data . . . . .	171
7.4	The LIVE Framework . . . . .	172
7.4.1	Licensing information from vocabularies and datasets. . . . .	173
7.4.2	Licenses compatibility verification. . . . .	174
7.4.3	Perspectives . . . . .	176
7.5	Summary . . . . .	176
<b>8</b>	<b>Conclusions and Future Perspectives</b>	<b>177</b>
8.1	Conclusions . . . . .	177
8.2	Future Perspectives . . . . .	179
8.2.1	Opportunities and Challenges for IGN-France . . . . .	179
8.2.2	Generic Visualizations on Linked Data . . . . .	181
8.2.3	Vocabularies and LOV . . . . .	182
	<b>Bibliography</b>	<b>185</b>
<b>A</b>	<b>Installation instructions for the JavaScript plugin</b>	<b>205</b>
A.1	Installing and configuring the REST-interface . . . . .	205
A.2	Installing and configuring the Admin-interface . . . . .	205
A.3	Installing and configuring the Event-website . . . . .	206
<b>B</b>	<b>Code source of vocabularies</b>	<b>207</b>
B.1	Vocabulary for geometry . . . . .	207
B.2	Vocabulary for CRS . . . . .	212
B.3	Vocabulary for French Administrative units . . . . .	223
B.4	Vocabulary for Visualization applications . . . . .	233
	<b>IV Résumé en Français</b>	<b>239</b>

# List of Figures

1.1	The LOD cloud as of May, 2007 . . . . .	2
1.2	RDF representation of the resource of Nice in DBpedia and Geonames . . . . .	5
1.3	Linking Open Data cloud diagram 2011, by Anja Jentzsch and Richard Cyganiak. <a href="http://lod-cloud.net/">http://lod-cloud.net/</a> . The highlighted portion corresponds to the geospatial datasets. . . . .	7
1.4	Linking Open Data cloud diagram 2014, by Max Schmachtenberg, Christian Bizer, Anja Jentzsch and Richard Cyganiak. <a href="http://lod-cloud.net/">http://lod-cloud.net/</a> . . . . .	8
2.1	Different geometries for representing a feature. . . . .	20
2.2	Vector representations of entities in one of the address database, BD ADRESSE®, produced by IGN-France. . . . .	21
2.3	The latitude and longitude angles represent the 2D geographic coordinate system. . . . .	24
2.4	The User Interface of the Geo Converter. . . . .	28
2.5	Results of conversion from WGS 84 to Lambert 93. Note: DD=Decimal Degree. . . . .	29
2.6	Results of conversion from WGS 84 to WGS 84 UTM. DD=Decimal Degree. . . . .	29
2.7	Illustration of the topological relations by Max Egenhofer. . . . .	32
2.8	Coordinate Reference Systems used in France . . . . .	37
2.9	High level classes of the <code>ignf</code> , <code>geom</code> and <code>topo</code> vocabularies; relationships between them and mappings with external vocabularies. . . . .	43
3.1	Generic architecture of tools for converting raw geospatial data into RDF. . . . .	53
3.2	Features of a geographic object adapted from [1]. . . . .	55
3.3	The most used distances on geometrical primitives: (a) the Euclidean distance and (b) the Hausdorff distance. . . . .	56
3.4	Architecture of the Geoknow Stack. . . . .	62
3.5	Lifting process of raw data source into RDF using Datalift Platform . . . . .	64
3.6	Architecture of Datalift platform. . . . .	65
3.7	French LOD cloud diagram based on the different datasets published in 4-5 stars. . . . .	75
4.1	Sample application of analyzers and visualizers in a LDVM pipeline. . . . .	87
4.2	Sample description of a Web application at the Open Data Service . . . . .	91
4.3	Conceptual Model of the DVIA vocabulary . . . . .	95
5.1	Big picture and architecture of the Linked Data visualization wizard. . . . .	101

5.2	Categories detected and visualization generated by the Linked Data visualization wizard in the case of EventMedia endpoint service. . . . .	105
5.3	Google Knowledge Graph Reverse Engineering Process. . . . .	108
5.4	Screenshot of the user interface. The circles with numbers highlight the different elements : (1) list of endpoints, (2) number of resources available in the map area, (3) A zoom to a given element and (4) description of the selected resource. . . . .	113
5.5	<i>Confomaton</i> general architecture. . . . .	115
5.6	Example of data modeled in <i>Confomaton</i> re-using multiple vocabularies	118
5.7	Steps for searching high schools in Antibes, France in a radius of 4000 meters. . . . .	125
5.8	The RDF triples before changes. Here we state that the application SBA Gems has won an award in the event Apps for Entrepreneurs.	126
5.9	The RDF triples after changes. . . . .	127
5.10	Universal JavaScript components . . . . .	129
5.11	Screenshots of the admin interface. On the left is the event listing page and on the right is the form for creating an event. . . . .	130
6.1	Datalift life cycle for publishing Linked Data. . . . .	138
6.2	Four steps to follow for publishing Linked Data by Governments and Institutions. . . . .	139
6.3	Evolution of the vocabularies inserted into LOV from 2011 to 2014. .	146
6.4	Equivalent classes and properties between foaf and dcterms . . . . .	149
6.5	Translations example for foaf:Person . . . . .	149
6.6	Meeting points between LOV and the NeOn methodology, derived from [2]. . . . .	150
6.7	Evolution of the number of prefix-namespace pairs registered in prefix.cc and LOV . . . . .	152
6.8	Matching data properties with ontology predicates in Data2Ontology module . . . . .	167
7.1	Licenses distribution in the LOV licensed vocabularies. . . . .	171
7.2	LIVE framework architecture. . . . .	173
7.3	LIVE tool user interface and sample results . . . . .	174
7.4	Licenses compatibility module. . . . .	175

# List of Tables

1.1	Format of an address in France with the example of EURECOM. . . . .	4
2.1	Statistics on the usage of the four main geographic vocabularies in the LOD cache. . . . .	30
2.2	Requirements and implementations for vocabulary definitions in GeoSPARQL.	31
2.3	Comparison of some geo-vocabularies with respect to the GeoSPARQL requirements. . . . .	33
2.4	URI schemes and conventions used for vocabularies and resources. . . . .	36
2.5	List of concept schemes used in the topographic ontology. . . . .	44
3.1	Geodata by provider and their different access type, either API, Web service or SPARQL endpoint. . . . .	48
3.2	Results of the alignment process between GeOnto and existing vocabularies/taxonomies. . . . .	58
3.3	Survey of some generic popular triple stores. . . . .	60
3.4	Triple stores survey with respect to geometry types supported and geospatial functions implemented. . . . .	61
3.5	Comparison of Datalift with the GeoKnow Stack . . . . .	67
3.6	Evaluation results in the interlinking process. . . . .	69
3.7	Interlinking results using the Hausdorff metric of LIMES tool between LinkedGeoData and toponyms in the French Gazetteer . . . . .	70
3.8	Initial mappings of Bano2RDF with LGD amenities resources respectively in Paris, Marseille and Lyon. The links are obtained using LIMES tool with a threshold of .97 using the Hausdorff distance. . . . .	74
3.9	Overview of the content of our contribution to the French LOD Cloud.	75
3.10	Results of the public buildings 10 km around EURECOM from LinkedGeoData endpoint. <code>1gdr</code> represents the base URI for <a href="http://linkedgeodata.org/triplify/">http://linkedgeodata.org/triplify/</a> . . . . .	76
4.1	Survey of tools used for creating visualizations on the Web. . . . .	89
4.2	Gathering reusable information from the openspending in Greece application . . . . .	93
4.3	Some innovative applications buit over Open Government Datasets . . . . .	94
5.1	A taxonomy of information visualizations for consuming Linked Datasets with suitable vocabulary space and visual elements. . . . .	100
5.2	Classification of the endpoints according to the datatype detected with our SPARQL generic queries . . . . .	106
5.3	Categories detected in some <i>dbpedia</i> endpoint domains, where “1” is the presence and “0” the absence of the given type of category. . . . .	106

5.4	Agreement on properties between users and the Knowledge Graph Panel . . . . .	109
5.5	Metadata provided by the Dog Food Server for the ISWC 2011 conference. . . . .	114
5.6	Media services used during the ISWC 2011 conference . . . . .	115
6.1	Summary of the best practices to publish Linked Data on the Web adapted from [3] . . . . .	142
6.2	Catalogs of vocabularies with respectively the number of the ontologies, the presence of a search feature, the catalog category and whether it is maintained or not. . . . .	143
6.3	Comparison of LOV, with respect to Swoogle, Watson and Falcons, based on part of the framework defined by D'Aquin and Noy in [4].	147
6.4	Type of issues encountered for vocabulary clashes . . . . .	154
6.5	LOV and prefix.cc conflicts resolution leading to contact vocabularies editors for negotiation. We provide the prefix, the URI in LOV and the action undertaken. . . . .	155
6.6	Experiments looking for stable results of finding vocabularies in prefix.cc. . . . .	157
6.7	Analysis of the URIs with no classes and no properties while using the LOV-Bot API . . . . .	158
6.8	Top 15 vocabularies according to their PIC. All the prefixes used for the vocabularies are the ones used by LOV . . . . .	161
6.9	Ranking of Top 20 terms (classes and properties) according to their IC value . . . . .	162
6.10	Sample of vocabularies with terms deprecated in LOV . . . . .	164
6.11	Comparing ranking position when using PIC in LOV with respect to prefix.cc and vocab.cc . . . . .	165
7.1	Evaluation of the LIVE framework. . . . .	175

# Listings

2.1	Sample output of the batch converter . . . . .	27
2.2	SPAQRL Query for creating sameAs links between data modeled with <i>ngeo</i> and <i>geom</i> vocabularies . . . . .	41
2.3	Definition in Turtle of the axiom defining a POINT. . . . .	41
3.1	Sample of structured geometry of the city of Nice. . . . .	52
3.2	The metadata information used in the BANO2RDF dataset. . . . .	71
3.3	This LIMES script is used to interconnect the Bano2RDF dataset located in the region of Marseille with buildings in LinkedGeoData. .	72
3.4	Query on LinkedGeodata endpoint to find all public buildings 10 km around Eurecom building in SophiaTech. . . . .	75
4.1	Snapshot in Turtle of the description of Event Media Live Application	96
5.1	Generic query to detect geo data from a SPARQL endpoint. . . . .	102
5.2	Generic query to detect time data from a SPARQL endpoint, using <i>time</i> , <i>dbpedia-owl</i> , <i>intervals</i> vocabularies. . . . .	102
5.3	Generic query to detect person categories from a SPARQL endpoint, using <i>foaf</i> , <i>dbpedia-owl</i> , <i>vcard</i> vocabularies. . . . .	102
5.4	Generic query to detect ORG data from a SPARQL endpoint. . . . .	103
5.5	Generic query to detect event data from a SPARQL endpoint, using <i>lode</i> , <i>event</i> , <i>dbpedia-owl</i> vocabularies. . . . .	103
5.6	Generic query to detect SKOS data from a SPARQL endpoint, using <i>skos</i> vocabulary. . . . .	103
5.7	Excerpt of a Fresnel lens in Turtle . . . . .	110
5.8	SPARQL query runs against each endpoint containing geodata to retrieve a random location. . . . .	110
5.9	Sample output of the JSON dataset used in the GeoRDFViz application.	111
5.10	Sample output of the media collector showing Google+ and Flickr results using #iswc2011 as the query term. . . . .	116
5.11	Sample output describing a resource. . . . .	117
5.12	Example configuration file of the <i>Confomaton</i> API, specifying event properties access. . . . .	119
5.13	Snapshot in Turtle for the school ID=0750676C, also at <a href="http://semantics.eurecom.fr/datalift/PerfectSchool/#school/0750676c">http:// semantics.eurecom.fr/datalift/PerfectSchool/#school/0750676c</a>	122
6.1	SPARQL query asking for all the equivalent classes and properties between the vocabularies foaf and dcterms. . . . .	148
6.2	SPARQL query for all the labels defined for the terms containing person. . . . .	149
6.3	Query to find namespaces in LOV not in prefix.cc . . . . .	152
6.4	Query to find disagreements LOV and prefix.cc . . . . .	154
6.5	Sample output of a response of the Check API . . . . .	156
6.6	SPARQL query for computing the occurrence of a class . . . . .	160

---

6.7	SPARQL query for computing the occurrence of a property . . . . .	160
B.1	Formal definition of the geometry vocabulary in turtle. The current version is deployed at <a href="http://data.ign.fr/def/geometrie">http://data.ign.fr/def/geometrie</a> . . . . .	207
B.2	Formal definition of the CRS vocabulary in turtle. The current version is deployed at <a href="http://data.ign.fr/def/ignf">http://data.ign.fr/def/ignf</a> . . . . .	212
B.3	Formal definition of the French administrative units in turtle. The current version is deployed at <a href="http://data.ign.fr/def/geofla">http://data.ign.fr/def/geofla</a> . .	223
B.4	Formal definition of DVIA vocabulary in turtle. The current version is deployed at <a href="http://purl.org/ontology/dvia">http://purl.org/ontology/dvia</a> . . . . .	233

# List of Publications

## Book Chapter

- Ghislain Atemezing and Raphaël Troncy: **Multimedia metadata.** In Encyclopedia of Social Network Analysis and Mining, Springer Verlag, 2014, ISBN: 978-1461461692

## Journal

1. Pierre-Yves Vandenbussche, Ghislain Atemezing, Maria Poveda, Bernard Vatant: **Linked Open Vocabularies (LOV): a gateway to reusable semantic vocabularies on the Web.** Semantic Web Journal, *under review*, 2015.
2. Ghislain Atemezing, Raphaël Troncy: **Modeling visualization tools and applications on the Web.** Semantic Web Journal, *under review*, 2015.
3. Ghislain Atemezing, Oscar Corcho, Daniel Garijo, José Mora, María Poveda-Villalón, Pablo Rozas, Daniel Vila-Suero and Boris Villazón-Terrazas: **Transforming meteorological data into linked data.** In Semantic Web journal, Special Issue on Linked Dataset descriptions, 2012. IOS Press.
4. Mari Carmen Suárez-Figueroa, Ghislain Atemezing and Oscar Corcho: **The landscape of multimedia ontologies in the last decade.** In Multimedia Tools and Applications, Vol. 55, Num. 3, December 2011.

## Conferences

1. Guido Governatori, Ho-Pun Lam, Antonino Rotolo, Serena Villata, Ghislain Atemezing and Fabien Gandon: **LIVE: a Tool for Checking Licenses Compatibility between Vocabularies and Data.** In 13th International Semantic Web Conference (ISWC 2014), Demo Track, October 2014, Riva del Garda, Italy.
2. Ghislain Atemezing and Raphaël Troncy: **Information content based ranking metric for linked open vocabularies.** In 10th International Conference on Semantic Systems (SEMANTICS 2014), September 2014, Leipzig, Germany.
3. Ahmad Assaf, Ghislain Atemezing, Raphaël Troncy and Elena Cabrio: **What are the important properties of an entity? Comparing users and knowledge graph point of view.** In 11th Extended Semantic Web Conference (ESWC 2014), Demo Track, May 2014, Heraklion, Crete.

4. François Scharffe, Ghislain Atemezing, Raphaël Troncy, Fabien Gandon, Serena Villata, Bénédicte Bucher, Fayçal Hamdi, Laurent Bihanic, Gabriel Képéklia, Franck Cotton, Jérôme Euzenat, Zhengjie Fan, Pierre-Yves Vandenbussche and Bernard Vatant: **Enabling linked-data publication with the datalift platform.** In AAAI 26th Conference on Artificial Intelligence, W10:Semantic Cities, July 2012, Toronto, Canada.
5. Ghislain Atemezing and Raphaël Troncy: **Vers une meilleure interopérabilité des données géographiques françaises sur le Web de données.** In 23èmes Journées Francophones d'Ingénierie des Connaissances (IC 2012), June 2012, Paris, France.
6. Houda Khrouf, Ghislain Atemezing, Thomas Steiner, Giuseppe Rizzo and Raphaël Troncy: **Confomaton: A conference enhancer with social media from the cloud.** In 9th Extended Semantic Web Conference (ESWC 2012), Demo Track, May 2012, Heraklion, Crete.

## Workshops

1. Ghislain Atemezing and Raphaël Troncy: **Towards a linked-data based visualization wizard.** In 5th International Workshop on Consuming Linked Data (COLD 2014), October 2014, Riva del Garda, Italy .
2. Guido Governatori, Ho-Pun Lam, Antonino Rotolo, Serena Villata, Ghislain Atemezing and Fabien Gandon: **Checking licenses compatibility between vocabularies and data.** In 5th International Workshop on Consuming Linked Data (COLD 2014), October 2014, Riva del Garda, Italy.
3. Ghislain Atemezing, Nathalie Abadie, Raphaël Troncy and Bénédicte Bucher: **Publishing Reference Geodata on the Web: Opportunities and Challenges for IGN France.** In 6th International Workshop on the Foundations, Technologies and Applications of the Geospatial Web (TERRA COGNITA 2014), October 2014, Riva del Garda, Italy.
4. Raphaël Troncy, Ghislain Atemezing, Nathalie Abadie and Cao-Vien Lam: **Modeling geometry and reference systems on the web of data.** In W3C Workshop on Linking Geospatial Data, March 2014, London, UK.
5. Ghislain Atemezing, Bernard Vatant, Raphaël Troncy and Pierre-Yves Vandenbussche: **Harmonizing services for LOD vocabularies: a case study.** In Workshop on Semantic Web Enterprise Adoption and Best Practice (WASABI 2013), October, 2013, Sydney, Australia.
6. Abdelfettah Feliachi, Nathalie Abadie, Hamdi Fayçal and Ghislain Atemezing: **Interlinking and visualizing linked open data with geospatial reference data.** In 9th International Workshop on Ontology Matching (OM 2014), October 2013, Sydney, Australia.

7. Ghislain Atemezing, Fabien Gandon, Gabriel Kepeklian, François Scharffe, Raphaël Troncy and Serena Villata: **When publishing linked data requires more than just using a tool.** In W3C Workshop on Open Data on the Web, April 2013, London, UK.
8. Ghislain Atemezing and Raphaël Troncy: **Towards interoperable visualization applications over linked data.** In 2nd European Data Forum (EDF 2013), April 2013, Dublin, Ireland.
9. Houda Khrouf, Ghislain Atemezing, Giuseppe Rizzo, Raphaël Troncy and Thomas Steiner: **Aggregating social media for enhancing conference experiences.** In 1st International Workshop on Real-Time Analysis and Mining of Social Streams (RAMSS 2012), June 2012, Dublin, Ireland.

## W3C Documents

1. Bernadette Hyland, Ghislain Atemezing and Boris Villazón-Terrazas: **Best Practices for Publishing Linked Data.** W3C Working Group Note published on January 9, 2014. URL: <http://www.w3.org/TR/ld-bp/>
2. Bernadette Hyland, Atemezing, Ghislain, Michael Pendleton, Biplav Srivastava: **Linked Data Glossary.** W3C Working Group Note published on June 27, 2013. URL: [www.w3.org/TR/ld-glossary/](http://www.w3.org/TR/ld-glossary/)



# Acronyms

Here are the main acronyms used in this document. The meaning of an acronym is usually indicated once, when it first appears in the text.

apps4x	Apps for X Co-creation Event Vocabulary
CRS	Coordinate Reference System
DCAT	Data Catalog Vocabulary
DOM	Data Object Model
DSRM	Data State Reference Model
DVIA	The Data VIualization Application Vocabulary
FrLOD	French Linked Open Data
IC	Information Content
LD	Linked Data
LDA	Linked Data API
LDVM	Linked Data visualization Model
LOD	Linked Open Data
LOV	Linked Open Vocabulary
GDAL	Geospatial Data Abstraction Library
GLD	Government Linked Data
GI	Geographic Information
GIS	Geographic Information System
GKP	Google Knowledge Panel
GPS	Global Positioning System
odapps	Open Data Applications Vocabulary
OGC	Open Geospatial Consortium
PIC	Partitioned Information Content
RDF	Resource Description Framework
RDFS	Resource Description Framework Schema
SDI	Spatial Data Information
SFA	Simple Features Access
SPARQL	SPARQL Protocol and RDF Query Language
UC	Use Case
UI	User Interface
URI	Uniform Resource Identifier
UPI	Universal Part Identifier
W3C	World Wide Web Consortium
WGS 84	World Geodetic System 1984



# CHAPTER 1

# Introduction

---

*“The Web as I envisaged it, we have not seen it yet.*

*The future is still so much bigger than the past.”*

Tim Berners-Lee

## 1.1 Context

The Web is currently in a transition phase from Web of documents to Web of Data. New devices and new ways to use them have emerged and changed user interaction with machines. The ubiquity of the Web also creates an unseen abundance of information. Data is flowing onto the Web, created by users, generated by sensors, and stored in ever-growing data farms. Geographic data is widely present on the Web as they are used to locate points of interest. At the same time, with the emergence of Open Data, many governments and local authorities are moving from legacy data stored in their databases to structured data on the Web. Structured data is already present in the many databases, metadata attached to media, and in the millions of spreadsheets created every day across the world.

The Web of Linked Data, unlike the Web of hypertext documents, is constructed with documents on the Web with links between arbitrary things described by RDF. Tim Berners-Lee [5] identifies 4 principles known as “Linked Data principles” as follows:

1. Use URIs as names for things
2. Use HTTP URIs so that people can look up those names.
3. When someone looks up a URI, provide useful information, using the standards (RDF\*, SPARQL)
4. Include links to other URIs. so that they can discover more things.

Hence, the URIs are used to identify any kind of object or concept on the Web via the HTTP protocol. Linked Data is continuously evolving, started in 2007 with a dozen of datasets (cf. Figure 1.1) to a large data space with thousands of datasets on different topics. From 2011 (see Figure 1.3)[6] to 2014, the number of datasets have tripled, with a significant growth of nearly 271% [7]. The new version of the Linked Open Data Cloud in April 2014 contains 1014 linked datasets which are connected by 2909 linksets, as depicted in Figure 1.4<sup>1</sup>. In order to enable Linked

---

<sup>1</sup> A more Web friendly version can be accessed at <http://data.dws.informatik.uni-mannheim.de/lodcloud/2014/>

Data applications to discover datasets and to ease the integration of data from multiple sources, Linked Data publishers should comply with the following set of best practices for publishing datasets on the Web [8]:

- **Data selection:** The dataset should be selected based on its potential relevance to be reused in an open format accessible somewhere on the Web.
- **Vocabulary usage:** The publishers should use terms from widely-used vocabularies in order to ease the interpretation of their data. If data providers use their own vocabularies, the terms of such proprietary vocabularies should be mapped to existing terms in “popular” vocabularies.
- **Linking:** By setting RDF links, data providers connect their datasets into a single global data graph which can be navigated by applications. Thus enables the discovery of additional data by following RDF links.
- **Dereferencable URIs:** By using HTTP URIs as identifiers for each resource, agents can easily look-up at the resources and “dereference” a URI in order to have access to the full representation identified by that URI. This helps build a network of URIs on the Web which enables navigation on different graphs.
- **Metadata provision and machine access to data:** Provide various ways for search engines and other automated processes to access data using standard Web mechanisms.

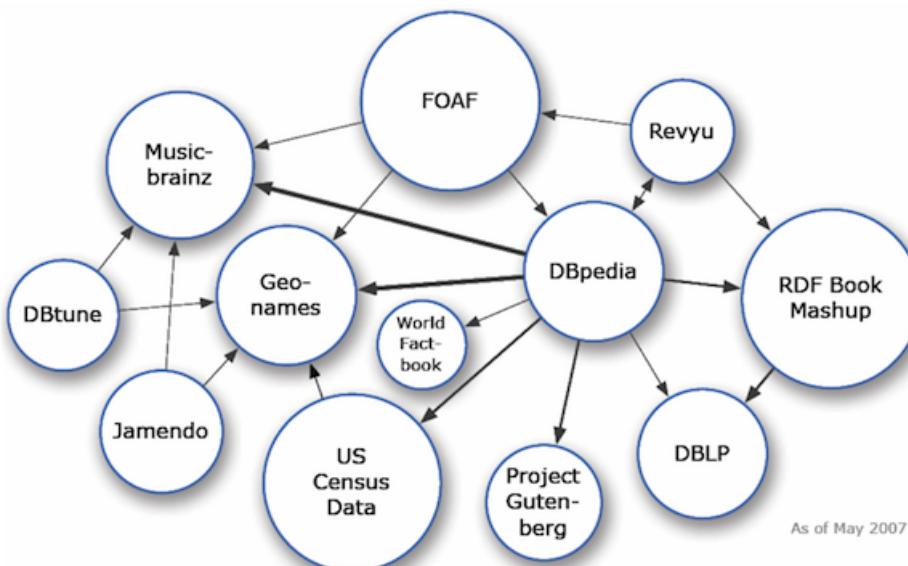


Figure 1.1: The LOD cloud as of May, 2007

However, geospatial data on the Web still lacks of interoperability for a better integration due to these three main factors:

1. Vendor specific geometry support, such as Google Maps API, Yahoo Geo Technologies, etc.
2. Different vocabularies, such as W3C Basic Geo<sup>2</sup>, NeoGeo<sup>3</sup>, GeoSPARQL [9], GML XMLLiteral<sup>4</sup> and vendor-specific.
3. Different spatial reference systems, such as Lambert93, WGS84, British National Grid, etc.

The potential of Linked Data for easing the access to government data is increasingly understood, as many countries such as UK, France, Australia, U.S.A. are releasing a significant amount of governmental and public sector data made accessible on the Web. Those organizations and public bodies are moving from legacy data stored in their databases to structured data on the Web. Structured data is already present in many databases, metadata attached to medias, and in the millions of spreadsheets created everyday across the world. However, the full potential of reusing such datasets will happen if they are published on the Web as Linked Data for more discovery and common understanding with appropriate semantically annotations.

### 1.1.1 Geographic Information (GI)

Geographic information (GI) surrounds us in our daily life. GI is commonly used to answer questions containing a *where* component. For example in “where shall Emmy study next year?”, the answer can show the different forms of expressing the location:

- In the answer: “Emmy will study in Nice”, Nice is used to identify the global area (locality) of where Emmy will study. In GI, this is an indirect reference. A map can display the geometry, and in the case of ambiguity, more context should be given, such as the country or the region.
- In the answer “Emmy will study in Sophia-Antipolis, Nice”; there is a relation “part-Of” (mereology) between Sophia-Antipolis and Nice, with more precision about the location of where Emmy will study.
- In the answer, “Emmy will study at Eurecom, Sophia-Antipolis, Nice”, Eurecom is the building “within” (topology)
- In the answer, “Emmy will study at 450, route des Chappes, 06160 Biot, Nice, France”, the full address is provided with standardized form of addressing, including a post code or zip code.

Table 1.1 shows the standard convention used to describe address in France. In [10], the authors summarized clearly the different forms of GI:

---

<sup>2</sup><http://www.w3.org/2003/01/geo/>

<sup>3</sup><http://geovocab.org/doc/neogeo/>

<sup>4</sup><http://www.opengeospatial.org/standards/gml>

Format	Example
<i>Addressee (Natural person/Organization)</i>	EURECOM
<i>More detailed description of addressee (optional)</i>	M. Frank Bender
<i>Housenumber + Streetname</i>	450 route des Chappes
<i>Postal code + uppercase town</i>	06410 Biot
<i>Country (if other than France)</i>	

Table 1.1: Format of an address in France with the example of EURECOM.

- Geometry: Geometries for GI can use raster or vector representations. In the former, GI is represented as an array of pixels (or cells), with each cell corresponding to a value and the position of each pixel corresponding to an area in the Earth. The latter is usually represented as a vector data of points, lines and polygons. This representation of GI is the mostly used by geographic datasets in the Web of Data
- Topology and Mereology: They are used to model GI topologically and mereologically. Topology expresses the spatial relationships that exist between different objects, and mereology describes the whole-part relationships.
- Textual representations, where GI is described as text or directions. Classifications, taxonomies and ontologies are often used to classify objects or geographic features. In this thesis, we use ontologies to describe French features in the French mapping agency.

In traditional GI, maps are the most obvious manner to visualize locations. A map is a two-dimensional representation of the landscape that can be used to navigate, visualize different aspects of a feature (street-view, satellite view, etc.).

### 1.1.2 Geographic Data in the Linked Data Cloud

The Linked Data Cloud originating from the Linking Open Data project classifies the data sets by domains, highlighting the diversity of data sets present in the Web of Data. The last classification in 2014 reveals 8 domains: government, publications, life sciences, user-generated content, cross-domain, media, geographic and social Web. Geography information is often used to connect information from varied topical domains. This appears in the Web of Data, where the Geonames<sup>5</sup> dataset serves as a hub for other datasets that have some geospatial component. Geonames is an open-license geographical database that publishes Linked Data about 8 million locations.

Figure 1.2 depicts a graph model of the city of NICE in France, starting from DBpedia [11] dataset and linked to Geonames with the relation “owl:sameAs”. In Geonames, the resource is linked back to DBpedia with a different relation, which

---

<sup>5</sup>url`http://www.geonames.org/`

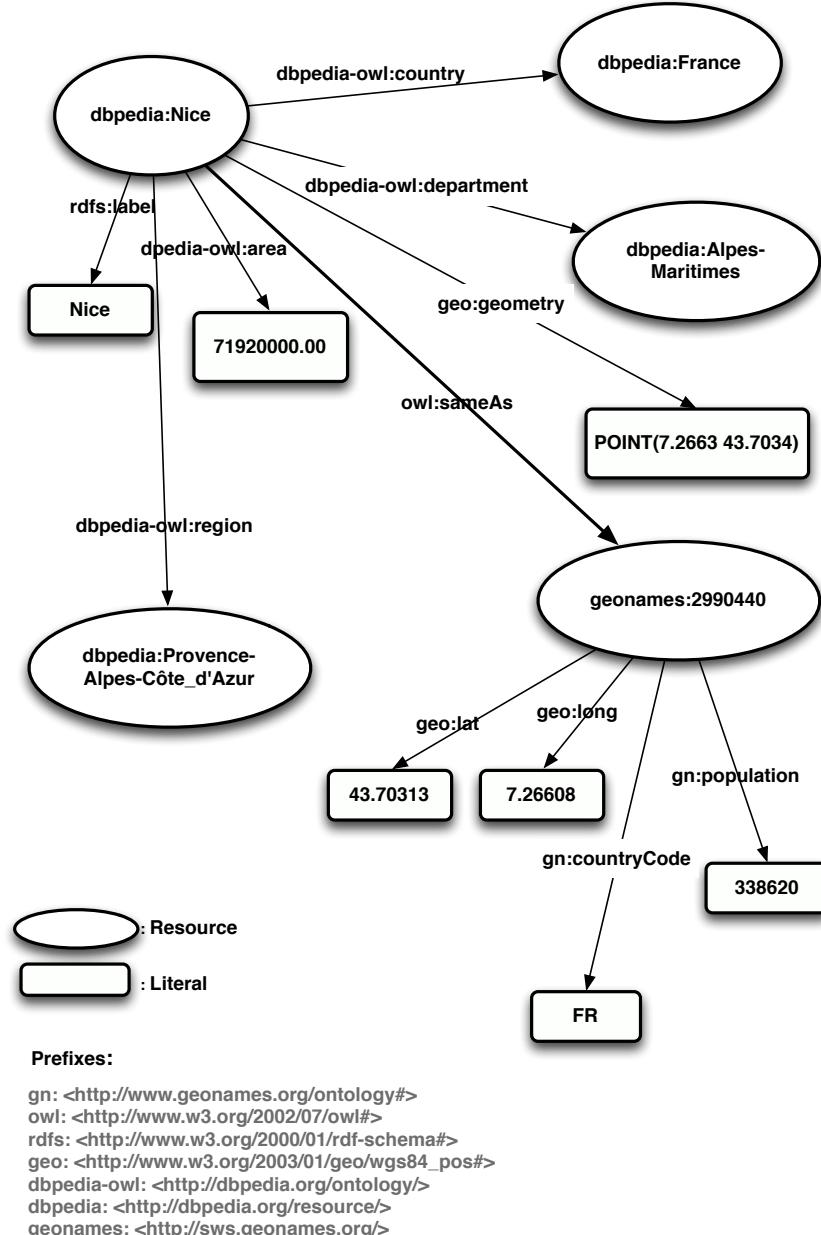


Figure 1.2: RDF representation of the resource of Nice in DBpedia and Geonames

is “`rdfs:seeAlso`”. These different connectors are useful to link different datasets and to achieve the fourth principle of Linked Data as set out by Berners-Lee [5]: “*include links to other URIs so that they can discover more things*”.

Another significant dataset in geospatial data is LinkedGeoData [12], a Linked Data conversion of data from the OpenStreetMap project, which provided information about more than 350 million spatial features. Wherever possible, locations in Geonames and LinkedGeoData are interlinked with corresponding locations in DBpedia,

ensuring there is a core of interlinked data about geographical locations.

Ordnance Survey (the national mapping agency of Great Britain) is one of the first mapping agency to publish Linked Data describing the administrative areas within the Great Britain [13], in efforts related to the [data.gov.uk](#) initiative. Since then, other initiative have begun to publish Linked Data, such as GeoLinkedData(es) [14], Linked Sensor Data, US Census data.

One of the key benefit of Linked Data is the use of the RDF model to manipulate data on the Web, interconnect it with other data and consume it in a variety of applications. In the process of getting those datasets effectively published, there are some barriers that prevent publishers to embrace the movement, such as the following:

- RDF and its different serializations is difficult to understand and use in practice compare to CSV or JSON.
- Given a dataset, choosing a suitable vocabulary to model the data is a big challenge.
- There are few easy-to-use tools that guide publishers in their process of publishing their dataset without being specialists in the different technologies: such as SPARQL, server configurations for dereferencing URIs, etc.
- The tools for converting data into RDF are either more domain-specific, or difficult to configure for non experts.

Nevertheless, the resulting “Web of Data” has started being populated in different domains, particularly with geospatial data, as proved by the efforts in [15, 12, 14, 16]. Those efforts and initiatives follow the vision of the *Semantic Geospatial Web* promotes by Max Egenhofer in [17] challenging GIS researchers to contribute to the Semantic Web effort by creating geospatial ontologies, query languages and processing techniques adapted to geospatial information on the Web.

The recent emergence of Linked Data radically changes the way structured data is being considered. By giving standard formats for the publication and interconnection of structured data, linked data transforms the Web into a giant database. While making data available on the Web, there is a need to build meaningful applications to show the value of all the data so that users can easily explore it, and derive new insights from it. As many information visualization (InfoVis) tools are already present in InfoVis community<sup>6</sup>, their easy adoption and usage for displaying structured data raise new challenges. Those challenges are the following:

1. How to specify and define semantic Web applications in terms of tools, widgets that can easily visualize RDF datasets?
2. How to mine efficiently heterogeneous structured data to derive patterns for automatically recommending adequate visualization tool to help users building innovative applications in an affordable time.

---

<sup>6</sup>[http://en.wikipedia.org/wiki/Information\\_visualization](http://en.wikipedia.org/wiki/Information_visualization)

3. How could we bridge the gap between traditional infoVis tools and Semantic Web technologies to built easily applications on top of datasets published as L(O)D?
4. How to represent and share visualizations built with datasets already present in different Open Data portals, such as <http://data.gouv.fr> and <http://data.gov.uk>.

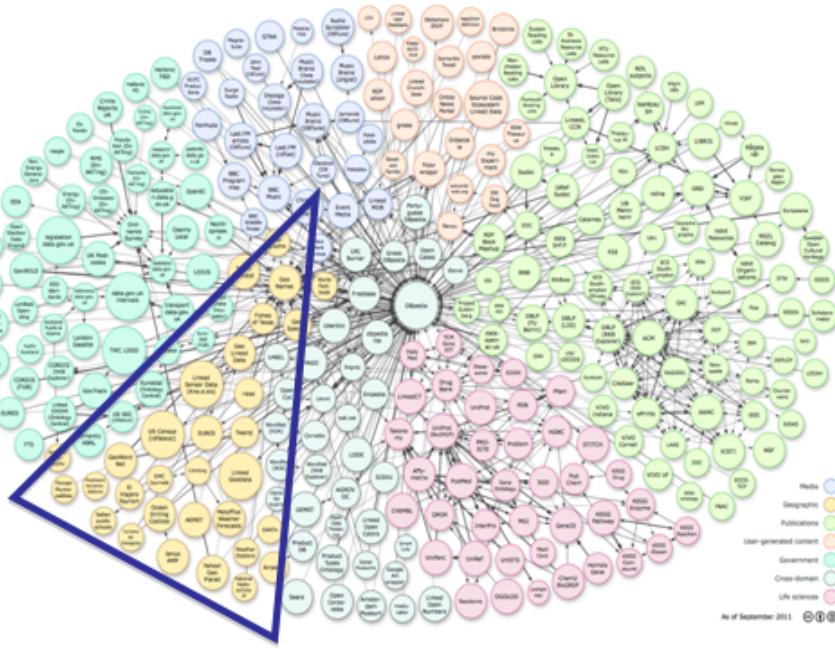


Figure 1.3: Linking Open Data cloud diagram 2011, by Anja Jentzsch and Richard Cyganiak. <http://lod-cloud.net/>. The highlighted portion corresponds to the geospatial datasets.

## 1.2 Research Questions

The ubiquity of the Web is creating an unseen abundance of information. Data is flowing onto the Web, created by users, generated by sensors, and stored in ever growing data farms. Geographic data is widely present on the Web as they are used for location of Point of Interest. At the same time, many organizations are moving from legacy data stored in their databases to structured data on the Web. Structured data is already present in many databases, metadata attached to medias, and in the millions of spreadsheets created everyday across the world. Many Linked Open Datasets have geospatial components, but still not having a common ways to describe features, spatial objects or geometries. Given the following three use-cases to express how challenging is to integrate geographic data from different datasets to obtain relevant answers:

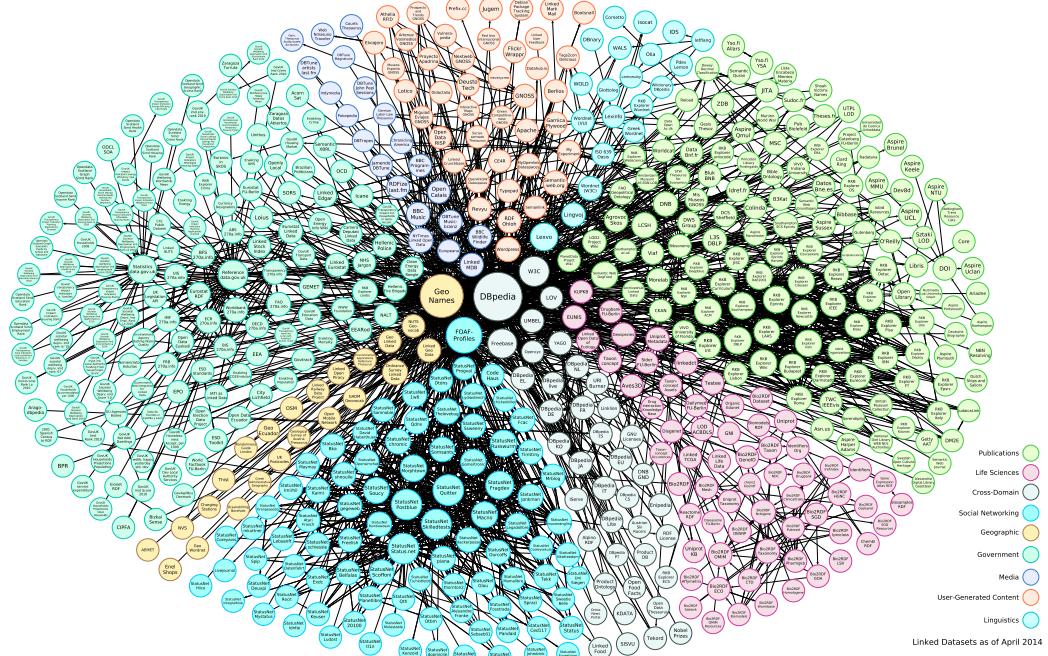


Figure 1.4: Linking Open Data cloud diagram 2014, by Max Schmachtenberg, Christian Bizer, Anja Jentzsch and Richard Cyganiak. <http://lod-cloud.net/>

- UC1: What DBpedia Historic Buildings are within walking distance from my current location ?
- UC2: What departments are located inside the bounding box composed of the Eurecom location and the Eiffel Tower?
- UC3: Give me the centroids of the administrative units in France in the projection Lambert93?

The aforementioned use-cases take into account “Concepts” (e.g: Historic Building, Department) that are defined differently depending on the provider of the dataset (e.g: DBpedia, OpenStreetMap, IGN). Further, the aforementioned use-cases implicitly make use of some specific topological functions widely used in the GIS applications, such as “within”, “inside”, “bounding box”. For example, Use Case 2 also mentioned Lambert93, a specific projection of data in France. Our aim is to contribute to actual efforts in representing geographic objects to leverage the barrier of integration of geospatial data both by the publishers and the lay users consuming the data.

Our main concern is to tackle the problems within the workflow of publication in two directions, more likely to happen at the beginning and the end:

- (i) Geographic Information on the Web of Data: as an application of the life-cycle of publishing geodata.

- (ii) Visualization tools for building innovative applications consuming structured data: as for leveraging the process of creating applications on-top of semantic data to highlight some relevant knowledge to the users.

In [18], the authors mentioned many research topics in the area of linked geospatial data. In this thesis, we address the following challenges:

1. *Vocabularies*: How do we model geospatial information on the Web? How do we assess geospatial ontologies? How should complex geometry be serialized?
2. *Query languages*: How do we write efficient queries that target geospatial Web? How do we store and index geodata in RDF ?
3. *Datasets*: How do we extract and convert geodata to expose it on the Web? What are the best practices for representing complex geometries on the Web? How can we integrate fully compatibility of Coordinate Reference Systems (CRSs) on datasets?
4. *Publication*: How can we develop scalable frameworks for covering the workflow of publishing geodata? What are the appropriate triple stores for handling geodata? What are the metrics to use for interconnecting different geodata resources on the Web?
5. *Applications and user interfaces*: How do we generate visualizations of linked geospatial data? What are appropriate high-level APIs that ease the development of user interfaces for geospatial data? Can we rely on existing map platforms such as Google Maps, Bing Maps or OpenStreetMap?

In this thesis, we tackle the issues and challenges both from publishers and users point of view. Publishers and users need pragmatic solutions that help them to choose a vocabulary, find a tool to convert ShapeFiles according to well-known vocabularies, then generate and publish the data following some best practices.

After the publication of the dataset using Linked Data paradigm, the logical question is: “What next”? Publishers and users need to be able to understand their dataset while developers have to build application out of the datasets. There is a huge amount of structured data on the Web, with the expectation to have more and more data to be published as LOD. This huge amount of “structured big data” is far to reach non-expert users so far, due to the complexity of different technologies used to model and query the data. Thus, it is important to create visualizations on top of the generated datasets to explore, analyze and showcase some of the potential of linked data to non-expert users. In this process, some challenging research questions have to be addressed, such as the following:

- How to find suitable visualizations according to the datasets without showing the complexity of SPARQL queries?

- What are the important properties to visualize for entities, depending on the domain and the users' expectations?
- How to bridge the gap between existing traditional tools of Information Visualization, mostly using CSV/XLS, JSON or proprietary formats to easily integrate RDF models as input?
- How to make interoperable applications built on top of Government Open Data catalogues? How to reuse existing applications?

While trying to answer to the aforementioned challenges, we report the state of the art and related approaches dealing with visualizations for Linked Data.

## 1.3 Contributions

In this section, we provide the main contributions of this thesis organized in three main parts of our thesis: model and publication of geospatial data, visualizations of data and applications on the Web and contribution to standards.

### 1.3.1 Modeling, Publishing and Querying geodata

The geolocation is crucial for many applications for both human and software agents. More and more data is opened and interlinked using the Linked Data principles, and it is worth modeling geographic data efficiently by reusing as much as possible from existing ontologies or vocabularies that describe both the geospatial features and their shapes. In the first part of our work, we survey different modeling approaches used by the Geographic Information System (GIS) and the Linked Open Data (LOD) communities. Our aim is to contribute to the actual efforts in representing geographic objects with attributes such as location, points of interest (POI), and addresses on the Web of Data. We focus on the French territory and we provide examples of representative vocabularies that can be used for describing geographic objects. We propose some alignments between various vocabularies (DBpedia, GeoNames, Schema.org, LinkedGeoData, Foursquare, etc.) in order to enable interoperability while interconnecting French geodata with other datasets.

Regarding this aspect of our research, we have achieved the following tasks:

1. We have proposed an ontology describing features and points of interest for the French territory, by reusing an existing taxonomy (GeOnto) and aligning it to other related vocabularies in the geolocation domain (Section 2.4.1).
2. We have studied how to extend the existing vocabularies for geographic domain to take into account efficient modeling of complex geometries. By doing so, we tackle the complex geometry representation issues in the Web of Data, describing the state of implementations of geospatial functions in triple stores and comparing them to the new GeoSPARQL standard (Section 2.3.1.1). We

finally make some recommendations and advocate for the reuse of more structured vocabularies for publishing topographic entities to better address the IGN-France requirements. (Section 2.3.1.2).

3. We have made a comparative study of triple stores, comparing their capability to store spatial information and their implementation of topological functions with respect to the ones existing in Open Geospatial Consortium<sup>7</sup> standards (Section 3.4).
4. We have designed and implemented vocabularies for describing complex geometries with different coordinate systems, with direct application to the French administrative units (Section 2.4).
5. We have interlinked French authoritative geodata resources with existing geospatial datasets on the Web, such as LinkedGeodata, GADM, NUTS and Geonames (Section 3.6.1.3).
6. We have contributed to the creation of a French LOD Cloud by publishing 8 datasets representing 340 million triples as LOD covering the French territory (Section 3.6.4).

Consuming data on the Web through visualizations is as challenging as applications have to deal with the graph structure of RDF, the underlying semantics of the dataset and the user interaction to easily understand what the data is about. In the following section, we present our contributions on visualization.

### 1.3.2 Visualization Tools in Linked Government Data

We first review some innovative applications that have been developed on top of datasets released as Open Data by governments (UK, USA, France) and local authorities. We have then derived and proposed 8 use cases (UCs) that can be developed to consume data from the different main providers in the French level: INSEE, DILA, IGN, FING, etc. We mention that the most interesting use cases are the ones which show the added value of having interconnected datasets. These UCs, developed and deployed, can be useful to show the benefits of Linked Data in a variety of domains such as education, tourism, cultural heritage, civil administrations, judicial courts, medicine, etc.

Regarding tools used for visualization, we have identified and classified them in two categories, providing for each of them relevant examples: (i)-tools that operate over RDF data, (ii) and tools that operate over other structured formats. We then provide some basic criteria for assessing a given visualization tool, with some weight attached to each of the criterion.

Moreover, regarding visualizations on top of datasets, we have contributed as follows:

---

<sup>7</sup><http://www.opengeospatial.org/>

1. We have built an application of the French first-round elections in 2012 using data from the <http://data.gouv.fr> and other public institutions. The application available at <http://www.eurecom.fr/~atemezin/DemoElection/> was built with the Exhibit Framework [19]; it aims to showcase the reconciliation of heterogeneous datasets: political results in CSV, unemployment rate, data of candidates, departments of France and more data from DBpedia. The user can filter by candidate image, unemployment rate and department to see the scores, with more enriched information about the department.
2. We have implemented a generic tool for exploring geodata on a map, based on the automatic detecting of SPARQL endpoints in the LOD cloud containing geospatial datasets (Section 5.4).
3. We have implemented an application consuming geodata and statistics combining multiple datasets in education from <http://data.gouv.fr> portal (Section 5.6).
4. we have built an application for conference events (Section 5.5) with their associated media reconciled from many social platforms (Instagram, Twitter, etc.).
5. We have built a vocabulary for structuring applications on the Web of Data. The vocabulary can be used for discovering visual tools or charts used to build applications (Section 4.3.3).
6. We have implemented a generic plugin for annotating applications developed for contests to be included in any web page, leveraging the generation of structured content out of webpages using the vocabulary.
7. We have implemented a wizard that analyses an RDF dataset and recommend visualization based on predefined categories, using generic SPARQL queries for easing the exploration of datasets published as LOD.

### 1.3.3 Contributions to Standards

We contributed to the W3C Government Linked Data Working Group (GLD WG)<sup>8</sup> activity from July 2011 until December 2013. The objective of the Working Group was to “*provide standards and other information which help governments around the world publish their data as effective and usable Linked Data using Semantic Web technologies*”.

We contributed to three task forces as follows:

- Task Force #1 aims to create a linked data community directory<sup>9</sup> and to maintain it on-line. The directory covers deployments, vendors, contractors,

---

<sup>8</sup><http://www.w3.org/2011/gld/>

<sup>9</sup><http://dir.w3.org>

end-user applications. We contributed to define the requirements and providing data for the French organizations in the directory.

- Task Force #2 aims at providing best Practices for publishing Linked Data by producing recommendations regarding vocabulary selection, URI construction, a Linked Data Cookbook, versioning, stability and provenance. Here, we have prepared a checklist to help government to select and re-use vocabularies in their project. We have also proposed our vision of the Linked Open Data Life cycle, with best practices for creating URIs. We served as editor for the Linked Data Glossary [20] published as a W3C Note document, apart from contributing in many sections of the document “Best Practices for Publishing Linked Data” [3].
- Task Force #3 goal was to provide relevant vocabularies to be used by governments or local authorities in their process of exposing their data. We have participated actively in the discussions on the different vocabularies published as recommendations by the W3C such as the Data Cube [21], the ORG vocabulary [22] and the DCAT [23] vocabulary.

Regarding the use of standard vocabularies we have contributed on:

- Proposing a method to harmonize prefixes on the Web of Data with two services: Linked Open Vocabularies (LOV)<sup>10</sup> and prefix.cc<sup>11</sup>. The former is currently a maintained hub of curated vocabularies on the Web, while the latter is a focal point for developer to register and look-up prefixes for their resources or ontologies. The approach proposed can be extended to any catalogue of vocabulary as long as the vocabularies fulfill the requirements to be inserted into LOV catalogue (Section 6.2.3).
- Designing and implementing a new method for ranking vocabularies based on the Information Content (IC) and Partitioned Information Content (PIC) metrics (Section 6.3).
- We have developed a tool that determined in real-time whether different licenses present in the dataset and vocabularies are either compatible or not (Section 7.4).

## 1.4 Thesis Outline

The work presented within this thesis is composed of three major parts. The remainder of the thesis proceeds as follows:

In the first part of this thesis, we focus on the various models and vocabularies for representing geography and geometry. We survey the state of triple stores and

---

<sup>10</sup><http://lov.okfn.org/dataset/lov/>

<sup>11</sup><http://prefix.cc>

describe particular problems such as coordinate systems, and highlight our contributions: new vocabularies, an online converter between CRS, etc. We also describe how geographic datasets can then be converted into RDF using the Datalift process to be published on the Web. We then show how those datasets can be interlinked (possibly trying different instance matching tools) and conclude with a thorough analysis of those alignments in the case of the French mapping agency (IGN-France) datasets. More specifically:

**Chapter 2** describes the current limitations of geodata on the Web and the different vocabularies we propose for geometries, coordinate reference systems and feature types. We also propose some best practices for publishing geodata on the Web. **Chapter 3** focuses on tools for publishing and querying geodata, their differences and applications. We describe the Datalift platform, an open source platform to “lift” raw data sources to semantic interlinked data sources. After comparing Datalift with the Geoknow stack, we apply it in the process of publishing French Administrative Units and French Gazetteer datasets. We then presents the status of the *French LOD(FrLOD)* cloud and some sample of queries over structured geometries published within the <http://data.ign.fr> endpoint.

In the second part of the thesis, we cover three main issues regarding how to present RDF to end-users. First, we give a state of the art review of existing tools and solutions for visual representation and exploration of RDF (Visualbox, LODSpeAKr, Map4RDF, Linked Data Visualization Model, etc.) Then, we present our contribution: the wizard for visualizations including the vocabulary for describing visualizations, the prototype itself, etc. Third, we present two applications applied to events and statistics to showcase the consumption of interlinked datasets in a new fashion. Then, we present a mechanism for extracting and reusing application in Open Data events. Finally, we provide some insights on revealing the “important” properties of entities for visualization by analyzing the Google Knowledge Panel (GKP) and evaluating with users preferences. This part is divided in two chapters:

**Chapter 4** provides a survey on visualization tools and applications, with their limitations. We also describe the status of the applications on the Web and provide a classification of so-called “Linked Data Applications”. In **Chapter 5**, we present our contribution on new approaches to generate visualizations and applications. We first propose a novel approach for category-based visualizations. We then show an application for geographic domain. Two applications related to events and statistics are also described. Finally, we propose how to improve the discovery of applications in Open Data events, through a model and a universal plugin for annotating Web pages with RDF.

In the last part of the thesis in **Chapter 6**, we describe various contributions to the Linked Open Vocabularies (catalog description, vocabulary publications, APIs and endpoints): vocabulary prefix harmonization, vocabulary ranking metrics using information content. In **Chapter 7** we present some insights on checking license compatibility between vocabularies and datasets with the defeasible deontic logic by creating an automatic tool for licenses checking for data on the Web.

In **Chapter 8**, we conclude by highlighting some limitations and suggest new re-

search directions.



## Part I

# Modeling, Interconnecting and Generating Geodata on the Web



# CHAPTER 2

# Geospatial Data on the Web

---

*“The Semantic Geospatial Web will be a significant advancement in the meaningful use of spatial information.[17]*

Max J. Egenhofer

## Introduction

The increasing number of initiatives for sharing geographic information on the Web of Data has significantly contribute to the interconnection of many data sets exposed as RDF based on the Linked Data principles. Many domains are represented in the Web of Data (media, events, academic publications, libraries, cultural heritage, life science, government data, etc.) while DBpedia is the most used dataset for interconnection. For many published datasets , geospatial information is required for rendering data on a map. In the current state of the art, different approaches and vocabularies are used to represent the “features” and their geometric shape although the POINT is the most common representation making use of the latitude/longitude properties defined in the W3C Geo vocabulary. Other geometries from the OpenGIS standard (POLYGON, LINESTRING, etc.) are more rarely exploited (e.g. LinkedGeoData, GeoLinkedData) while fine-grained geometry representations are often required.

In France, the National Geographic Institute (IGN) has started to publish data in RDF, as illustrated by the recent experimental LOD service <http://data.ign.fr>. IGN maintains large databases composed of different types of geographic entities, such as buildings, topographic information, occupied zones, etc. Using existing taxonomies and publishing them on the Web would ease the integration, retrieval and maintenance of French geospatial objects. Moreover, adding semantics to the current data on the Web can not only resolve ambiguities between datasets, but also will enable answering more complex queries than current GIS systems can handle, such as: *“show all buildings used as tribunal courts in the 7th Arrondissement of Paris”*. Another use case is the possibility to reason over parts of a structure: *“show the points where the river Seine touches a boundary of a district in Paris that contain an area of interest (AOI)”*.

In this chapter, we first describe the notion of geographic information, with its specificity and diversity of formats (Section 2.1). Then, we survey the models used on the Web to model existing geodata, by pointing some limitations (Section 2.3.1 and Section 2.3. We then move to distinguish two levels of georeferencing data (direct and indirect), and the importance of CRSs in the interpretation of geodata

(Section 2.3.3). The contributions start with a REST service for converting geodata in Section 2.2, followed by vocabularies developed for handling both geometries and features on the Web (Section 2.4). The chapter closes with a brief summary.

## 2.1 Geographic Information

Geographical phenomena require two descriptors to represent the real world; *what is represented*, and *where it is*. as reported by authors in [24]. For the former, concepts such as “town”, “school”, “river”, are used to recognize the phenomena and described in terms of well-established “objects” or “entities”. At the same time, the type of concepts used to describe a phenomena vary from one scale of resolution to another, depending on the perception of the human observation of the world. The *where it is* of the phenomena may be defined in terms of a geometrically exact or a relative location. The *what is represented* uses local or world coordinate systems - local or internationally accepted projections that uses geometrical coordinates of latitude and longitude - defined using a standard system of spheroids, projections, and coordinates [24]. Figure 2.1 depicts different geometries for representing features. Two approaches are generally used to represent geographical primitives in GIS: vector and raster approaches. In the vector data model, a space is represented by a geometry that describes the location and implicitly the shape, and information attributes such as the name, nature, length, or surface. In general, the geometry of a geographic object can be described using three primitives:

- **Point** represented in terms of XY coordinates.
- **Line** represented by a sets of XY coordinate pairs that define a connected path through space.
- **Polygon** represented in terms of the XY coordinates of its boundary, or in terms of the set of XY coordinates that are enclosed by such a boundary.

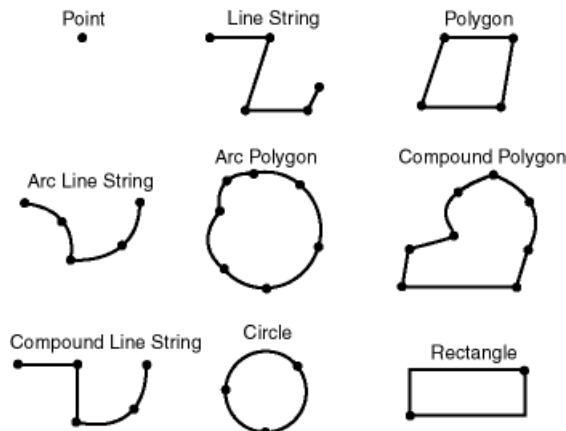


Figure 2.1: Different geometries for representing a feature.

Figure 2.1 shows different representations of an address in the BD ADRESS® database, including points, buildings, path and additional address location.

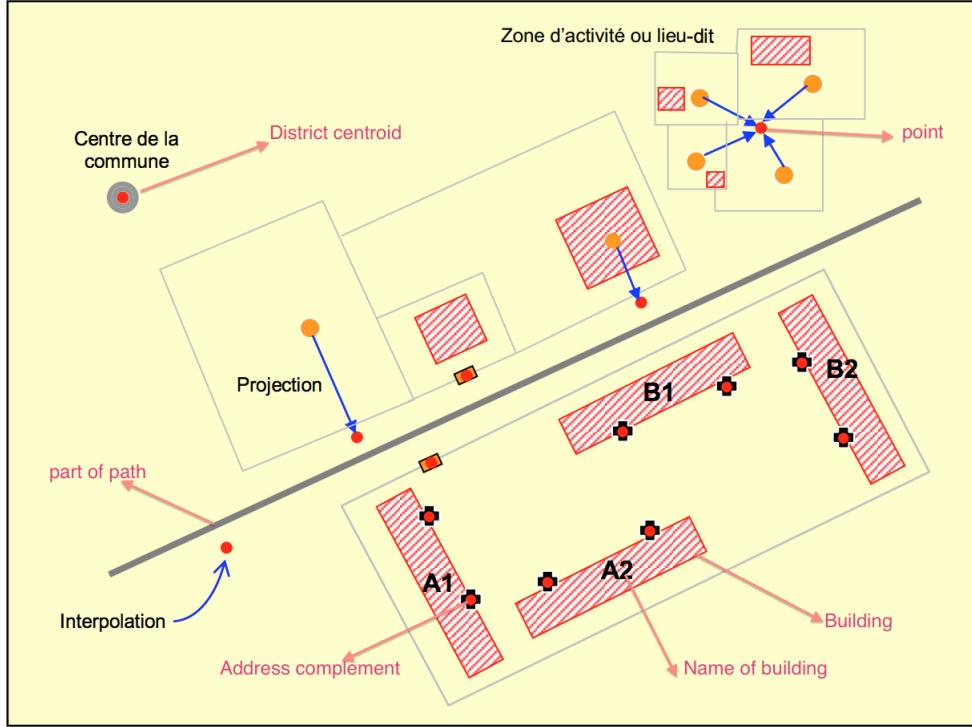


Figure 2.2: Vector representations of entities in one of the address database, BD ADRESSE®, produced by IGN-France.

### 2.1.1 Specificity

Depending on the level of the spatial resolution, a phenomenon reveals more or less detail, according to the intended use. So, increasing the level of resolution might reveal internal structure. For example, in the case of a town: sub-districts, suburbs, streets, houses, lamp-posts, traffic signs; which can be important for some purposes and not for others. The level of details also influences the representation of a given entity, and thus provides different views of the same entity in a GIS. A town could be represented by a point at the continental level of resolution but as a polygon entity at the regional level. A road at national level is adequately represented by a line; at the street level it becomes an area of paving.

Moreover, geographical data are stored in different GISs according to two different resolutions: the metric resolution and the “decametric” resolution [1]. This makes differences in the volume of the data, the scale used to gather and view them, and most importantly their relevance. Geodata providers explicitly provide details on the scale and the purpose of their datasets. Let's take the case of BDCARTO®, BDTOPO® and BD ADRESSE® of IGN-France. BD CARTO® represents the

elements using a decametric precision, and contains many themes: roads, administrative units, etc. It helps to manage data from 1 : 50 000 to 1 : 200 000 resolution. BD TOPO is used for producing maps at 1: 25 000 scale, and represents a 3D model of the territory and the amenities with addresses. BD ADRESSE is used for a precised location using the postal code in 1: 25 000 scale.

The aforementioned different resolutions, scalable issues represent constraints on abstraction, representation at different scales and requirements of geographical data create some issues both for users and producers. Those issues motivate the need for different relations between GIS datasets. The Web is a good medium to represent a real world phenomena with a unique Uniform Resource Identifier (URI) and associated semantics for referencing, interlinking and tracking the evolution of geodata over time.

### 2.1.2 Data Formats and Serialization

Diverse formats are used to store and exchange geodata in traditional GIS. Some of them are proprietary or closed formats, other are standards defined by the (Open Geospatial Consortium (OGC). We list below some of the most used formats:

**ESRI Shapefile:** A shapefile stores non topological geometry and attribute information for the spatial features in a dataset. The geometry for a feature is stored as a shape comprising a set of vector coordinates [25]. Concerning this format, four files are used for processing: the three mandatory dBASE file (.dbf), index file (.shx) and main file (.shp), plus the metadata file (.prj) that describes the CRS used by the dataset. The schema and the types of the thematic and geometric attribute are extracted from the main and dBase files.

**Geospatial DBMS:** A database that is optimized to store and query geospatial data, such as Oracle Spatial, PostGIS (a spatial extension to PostgreSQL), Spatial-Lite (a spatial extension to SQLite).

**GML:** An OGC encoding specification standard for geodata in XML that enables the storage, transport, processing and transformation of geographic information. GML serves as a modeling language for geographic systems as well as an open interchange format for geographic transactions on the Internet. As with most XML-based grammars, there are two parts to the grammar: the schema that describes the document and the instance document that contains the actual data. A GML document is described using a GML Schema. GML is also an ISO standard (ISO 19136:2007).

**Well-Known-Text (WKT):** Well-known text (WKT) is a text markup language for representing vector geometry objects on a map, spatial reference systems of spatial objects and transformations between spatial reference systems. The formats

were originally defined by OGC and described in the specifications for geographic information, simple feature access [26].

**GeoJSON:** GeoJSON is a format for encoding a variety of geographic data structures [27]. A GeoJSON object may represent a geometry, a feature, or a collection of features. It also supports the following geometry types: Point, LineString, Polygon, MultiPoint, MultiLineString, MultiPolygon, and GeometryCollection. Features in GeoJSON contain a geometry object and additional properties, and a feature collection represents a list of features. GeoJSON is the data format of choice for developers, which is widely implemented and supported by many tool chains. The default (and strongly recommended) Coordinate Reference System is WGS84 [28], but alternative systems can be specified. The recommended nomenclature for CRS systems is to use OGC URNs, for example urn:ogc:def:crs:OGC::CRS84 (for WGS84). EPSG identifiers, originally from the European Petroleum Survey Group and now maintained by the International Association of Oil and Gas Producers (OGP) can also be used. Alternatively, the parameters for a CRS can be linked to by URL. The example below represents in GeoJSON the location (point) of “Eurecom”:

```
{
  "type": "Feature",
  "geometry": {
    "type": "Point",
    "coordinates": [43.614151, 7.071414]
  },
  "properties": {
    "name": "EURECOM"
  }
}
```

**KML:** KML (Keyhole Markup Language)<sup>1</sup> is a language for the visualization of geographic information in 2D (on a map) or 3D (on a globe), including annotation of maps and images. Hence it can be used to specify layers for use in creating maps in a GIS system. A Placemark is one of the most commonly used features in Google Earth. It marks a position on the Earth’s surface, using a yellow pushpin as the icon. The simplest Placemark includes only a <Point> element, which specifies the location of the Placemark. A Placemark object contains the following elements:

- A *name* that is used as the label for the Placemark
- A *description* that appears in the “balloon” attached to the Placemark
- A *Point* that specifies the position of the Placemark on the Earth’s surface (longitude, latitude, and optional altitude)

---

<sup>1</sup><http://www.opengeospatial.org/standards/kml>

KML can be used to carry GML content, and GML can be “styled” to KML for the purposes of presentation. KML instances may be transformed loosely to GML, however roughly 90% of GML’s structures (such as metadata, coordinate reference systems, horizontal and vertical datums, etc.) cannot be transformed to KML. KML has very limited support for metadata as recommended by ISO 19115 [29]. The CRS in use is implicit and unique.

## 2.2 A REST Service for Converting Geo Data

### 2.2.1 Background

The Earth is shaped like a flattened sphere. This shape is called an ellipsoid. A datum is a model of the earth that is used in mapping. The datum consists of a series of numbers that define the shape and size of the ellipsoid and its orientation in space. A datum is chosen to give the best possible fit to the true shape of the Earth. There are a large number of datum in use. Many of them are optimized for its use in one particular part of the world. An example is the Geodetic 1949 datum that has been used in New Zealand. Another example used globally familiar to GPS users, is the WGS-84 datum.

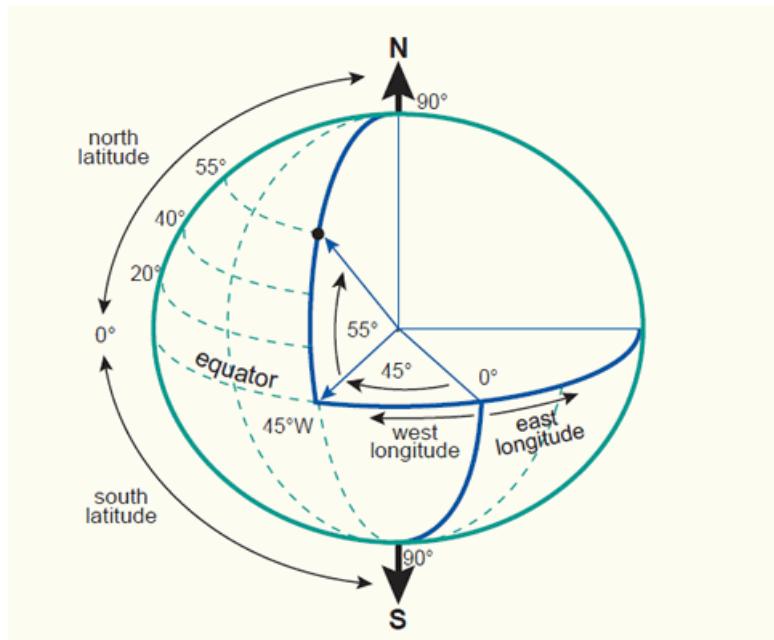


Figure 2.3: The latitude and longitude angles represent the 2D geographic coordinate system.

A point (location) is referenced by its longitude and latitude values. Longitude and latitude are angles measured from the earth’s center to a point on the earth’s surface. The angles often are measured in degrees (or in grads). It is important to keep in mind that latitude and longitude are always specified in terms of a datum.

The latitude and longitude of one current position are different for different datums. For example, the “*Théâtre National de Nice*” in Nice, France is at  $43.700594^{\circ}N$ ,  $7.277959^{\circ}E$  in WGS 84 coordinates and  $43.700570^{\circ}N$ ,  $4.941204^{\circ}E$  in NTF coordinates, an old France CRS. If the latter coordinates are used in WGS 84, they will point to a position which is approximately *295.17 kilometers* away from the theater. So when working with latitude/longitude coordinates and getting an error of a couple of hundred meters, it is most likely that the wrong datum is in used. Figure 2.2.1 shows an illustration of the 2D geographic coordinate system [30] with latitudes and longitudes given in angular degrees.

### 2.2.2 Why a REST service is needed

As we have seen, geodata interpretation relies on a coordinate reference system, and while the WGS84 CRS is the *de-facto* standard for GPS devices, many other CRS are in used. For example LAMBERT 93, RGM 04 or RGR 92 are respectively used for georeferencing points of interests in France continental, Mayotte or La Reunion. We have developed a REST service that is capable of transforming one dataset using a particular CRS into another one.

Previous closed-source software called Circé<sup>2</sup>, published by IGN-France, provides the abilities to convert coordinate between CRSs in France and WGS 84. It has two conversions modes: standard and grid. In each mode the user is required to input the source CRS values in order to convert. Based on what datum and projection method used, the number of required fields are different. Circé also has a batch converting mode which is done by supplying a file. The format of the content of the file is simple. The available formats are converted into *[location name] Lat/Lon Lon/Lat [Altitude/Height]*. There is an option to choose which format to use. But due to the fact that it is a closed source software, no one can use the software as a service for their system.

Another tool is a Web based application called the world coordinate converter at <http://twcc.free.fr>. This tool can convert between numerous CRSs. Unlike Circé which only supports France and WGS 84, this tool supports conversion of international CRSs and national CRSs. The result output from those previous two tools are the same. However, just like Circé, no API or any open service for the community to use their full power horse unless going to their website and use it like an application.

These tools are great as a standalone tool for end user, but not so great for the developer community. The developer can only use them to test their result. There is no way to use either fully functional algorithm except to develop it again.

**Our proposal:** The purpose of the REST Converter is to propose a web based service to perform conversion between various CRSs.

---

<sup>2</sup><http://geodesie.ign.fr/?p=53&page=circe>

The algorithms implemented are the ones described at <http://geodesie.ign.fr/index.php?page=algorithmes> and available within the standalone Circé software<sup>3</sup>. At the moment, the following features are implemented in the Geo Converter:

- from/to WGS 84 to/from WGS 84 UTM ;
- from/to WGS 84 to/from Lambert 93 and
- from/to WGS 84 UTM to/from Lambert 93

The API can also convert a file with space separated values. The API supports JSON as one of the output formats. The code of the REST service is available at <https://github.com/vienlam/Geo>.

### 2.2.3 API Access and Implementation

#### 2.2.3.1 Simple Converter

The API is working through URL like any RESTful service does. The syntax for the service is (note that {} is required and [] is optional): `http://{domainname}/eurecom.geo.rest/api/converter/{D1}[P1]{D2}[P2]?{Parameters}`

Where:

`domainname`: The service, currently hosted at Eurecom domain name.

`D1` and `D2`: Datum of source and target CRS respectively.

`P1` and `P2`: Projection of source and target CRS respectively.

`Parameters`: The required parameters as input to the service. See the table below for detail of required parameters for each converter. Parameters are provided as `p1=v1&p2=v1` where `p1` is the first parameter with `v1` is its value, etc.

An example:

```
eurecom.fr/eurecom.geo.rest/api/converter/WGS84RGF93Lambert93?lon=4.7021484375&lat=45.2130035559939
```

#### 2.2.3.2 Batch Converter

The API also support batch conversion from files input. The URL syntax is:

```
http://{domainname}/eurecom.geo.rest/api/converter/file/{D1}[P1]{D2}[P2]?{Parameters}
```

where `D1`, `P1`, `D2`, and `P2` are as before. The parameters are the source file and the encoding system. The order of the input coordinates in the file matters, the exact order is as follow: *longitude/x latitude/y [zone] [hemisphere]* The first value should be longitude, in the case of geographic coordinates, or x, in the case of planimetric coordinates. The second value is latitude or y. In case of UTM coordinates, the third value is the zone. The last value is hemisphere.

---

<sup>3</sup>[http://fr.wikipedia.org/wiki/Circe\\_\(logiciel\)](http://fr.wikipedia.org/wiki/Circe_(logiciel))

### 2.2.3.3 Result Format

Normally, the API will return the result in form of a normal string with the specify space character also as delimiter for each coordinate. In case of batch converter, each location will be in one line. For example:

```
eurecom.fr/eurecom.geo.rest/api/converter/WGS84RGF93Lambert93?lon=4.7021484375&lat=45.2130035559939
```

will return a string of:

```
833607.9336802219 6458515.660215093
```

The order of value respond to the order of coordinate as follow:

```
{longitude/x} {latitude/y} [zone] [hemisphere]
```

**JSON format result:** However, one can ask the API to return the result string in JSON format. To demand the API to do so, simply set the json parameter to the value 1 : json=1. For example:

```
eurecom.fr/eurecom.geo.rest/api/converter/WGS84RGF93Lambert93?lon=4.7021484375&lat=45.2130035559939&json=1
```

will return a string in JSON format (no line breaking):

```
{"y":6458515.660215093,"x":833607.9336802219}
```

JSON format result can also be specified for the batch converter. The command is as before. An example of JSON format results of batch converter from WGS84 to Lambert93:

```
1 {"point1":{"y":6543019.59988031,"x":882408.2999938729}, "point2":{"y":6544401.599880268,"x":881947.5999938913}, "point3":{"y":6581538.799879094,"x":849722.3999950405}, "point4":{"y":6561282.999879616,"x":917481.0999927416}, "point5":{"y":6561139.999879618,"x":917474.5999927415}}
```

Listing 2.1: Sample output of the batch converter

For better human readable result, a tool such as JSONLint<sup>4</sup> can be used for displaying the output of the JSON format.

---

<sup>4</sup> <http://jsonlint.com/>

### 2.2.3.4 User Interface

To access the User Interface (UI), the URL is: <http://{domainname}/eurecom.geo.rest/>. Figure 2.4 shows the landing page of the converter. The UI shows two systems, which are the source and the target system, to provide input. Which system is source or target depend on which button at the end is used. If the button *Convert S1-S2* is clicked, then the System 1 will be the source and System 2 is the target, and *vice versa* for button *Convert S2-S1*. The required inputs is just as mentioned before. First is the datum to use. Next is the projection method. If no project method is used then choose None. Next is the inputs of the source system. The map under the form will show where the coordinate in the source system point to through a marker. This marker can be used as an input to the source system as well. By dragging the marker, the system will update the input values of the source system, which is the system that have the last edited field or act as source system in last conversion, to correspond to the marker position. Furthermore, the "Use Marker's Position" button can be used to take marker's position and convert it to the system which the button resides in. This will not change the source system.

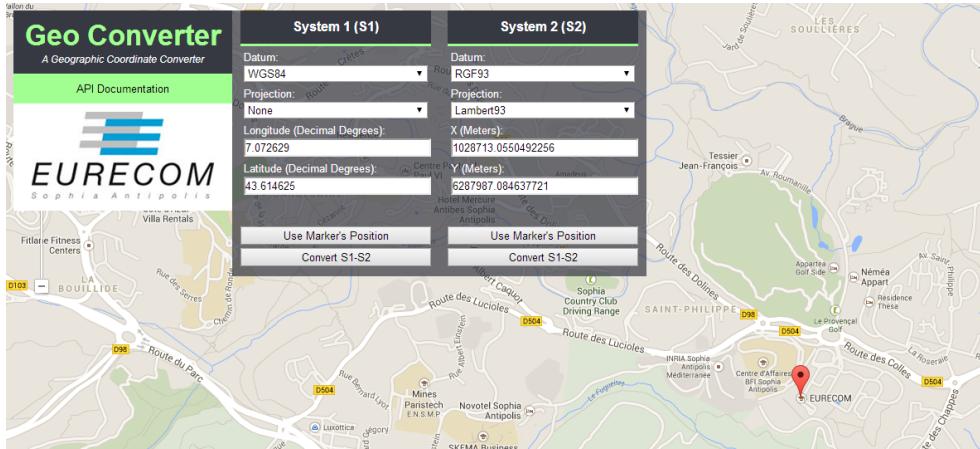


Figure 2.4: The User Interface of the Geo Converter.

### 2.2.3.5 Algorithm Evaluation

Figures Fig.2.5 and Fig.2.6 show the resulting conversion from our tool, Geo Converter, in comparison with Circé and twcc.free.fr (TWCC).

The results show that the results from Geo Converter are not much different from those from Circé and TWCC. This deviation can be tolerated, since when showing on the map, they are basically the same point, although the difference appears in the fifth decimal digit in the conversion from Lambert 93 to WGS 84. There is no need for evaluation of conversion between Lambert 93 to UTM, since it is needed an intermediate step of converting to WGS 84. The conversion from Lambert 93 to WGS 84 works well, as does the conversion from WGS 84 to WGS UTM.

Conversion from WGS 84 to Lambert 93					
Tools	Input		Output		
	Latitude (DD)	Longitude (DD)	X (m)	Y (m)	
Geo Converter	43.700594	7.277959	1044752.61	6298403.61	
Circé			1044752.20	6298404.17	
twcc.free.fr			1044752.61	6298403.61	
Conversion from Lambert 93 to WGS 84					
Tools	Input		Output		
	X (m)	Y (m)	Latitude (DD)	Longitude (DD)	
Geo Converter	1044752.61	6298403.61	43.700594	7.277959	
Circé			43.700589	7.277964	
twcc.free.fr			43.700594	7.277959	

Figure 2.5: Results of conversion from WGS 84 to Lambert 93. Note: DD=Decimal Degree.

Conversion from WGS 84 to WGS 84 UTM							
Tools	Input		Output				
	Latitude (DD)	Longitude (DD)	X (m)	Y (m)	Hem	Zone	
Geo Converter	43.700594	7.277959	361243.56	4840060.20	N	32	
Circé			361243.52	4840060.21	N	32	
twcc.free.fr			361243.52	4840060.21	N	32	
Conversion from WGS 84 UTM to WGS 84							
Tools	Input				Output		
	X (m)	Y (m)	Hem	Zone	Latitude (DD)	Longitude (DD)	
Geo Converter	361243.52	4840060.21	N	32	43.700586	7.277959	
Circé					43.700594	7.277959	
twcc.free.fr					43.700594	7.277959	

Figure 2.6: Results of conversion from WGS 84 to WGS 84 UTM. DD=Decimal Degree.

## 2.3 Current Modeling Approach

In this section, we review different approaches used to model geographical data on the Web, with their advantages and limitations. We first provide the status of vocabularies usage for geospatial data on the Web (Section 2.3.1), followed by a review of vocabularies for features and geometries, and finally the description of the GeoSPARQL specification and implementation. Further, section 2.4 details our contributions by implementing vocabularies for geometries and features.

### 2.3.1 Status of Vocabularies Usage for Geospatial Data

Publishing statistics concerning the actual usage of vocabularies on the LOD cloud<sup>5</sup> [31] provides not only an overview of best practices recommended by Tim Berners-Lee<sup>6</sup>, but also provides an overview of the vocabularies re-used in various datasets and domains. Concerning the geographic domain, the results show that W3C Geo<sup>7</sup> is the most widely used vocabulary, followed by the `spatialrelations`<sup>8</sup> ontology of Ordnance Survey (OS). At the same time, the analysis reveals that the property `geo:geometry` is used in 1,322,302,221 triples, exceeded only by the properties `rdf:type` (6,251,467,091 triples) and `rdfs:label` (1,586,115,316 triples). This shows the importance of geodata on the Web which contains almost 20% of triples of the LOD cloud datasets. Table 2.1 summarizes the results for four vocabularies (WGS84, OS spatial relation, Geonames ontology and OS admin geography) where the number of datasets using these vocabularies and the actual number of triples are computed<sup>9</sup>.

Ontologies	#Datasets	#Triples	SPARQL endpoint
W3C Geo	21	15 543 105	LOD cache
OS spatialrelations	10	9 412 167	OS dataset
Geonames ontology	5	8 272 905	LOD cache
UK administrative-geography	3	229 689	OS dataset

Table 2.1: Statistics on the usage of the four main geographic vocabularies in the LOD cache.

#### 2.3.1.1 GeoSPARQL Standard and specifications

OGC-Simple Features Access standard aims to support both representing and querying geospatial data on the Semantic Web. The standard [9] contains 30 requirements. It also defines a vocabulary for representing geospatial data in RDF and

<sup>5</sup><http://stats.lod2.eu>

<sup>6</sup><http://www.w3.org/DesignIssues/LinkedData.html>

<sup>7</sup>[http://www.w3.org/2003/01/geo/wgs84\\_pos](http://www.w3.org/2003/01/geo/wgs84_pos)

<sup>8</sup><http://data.ordnancesurvey.co.uk/ontology/spatialrelations>

<sup>9</sup>LOD cache should be understood as <http://lod.openlinksw.com/sparql/>. There are many more vocabularies used in the LOD cloud that contain also geographical information but that are never re-used.

provides an extension to the SPARQL query language for processing geospatial data. Moreover, GeoSPARQL defines functions that request or check properties of a geometry (e.g., `isSimple`, `isEmpty`, `Dimension`, `GeometryType`, `SRID`), function that test topological relations (e.g., `contains`), and functions that construct new geometries from existing ones (e.g., `buffer`). The proposed standard follows a modular design with five components:

- (i) A *core component* defining top-level RDFS/OWL classes for spatial objects;
- (ii) a *geometry component* defining RDFS data types for serializing geometry data, RDFS/OWL classes for geometry object types, geometry-related RDF properties, and non-topological spatial query functions for geometry objects;
- (iii) a *geometry topology component* defining topological query functions;
- (iv) a *topological vocabulary component* defining RDF properties for asserting topological relations between spatial objects; and
- (v) a *query rewrite component* defining rules for transforming a simple triple pattern that tests a topological relation between two features into an equivalent query involving concrete geometries and topological query functions.

Geo-Aspect	Requirement	Implementation Definition
Feature	Req 2	The Class <code>SpatialObject</code> should be defined & accepted
	Req 3	Defines Feature <code>rdfs:subClassOf SpatialObject</code>
	Req 4	Defines 8 Simple Features Object Properties(OP)
	Req 5	Defines 8 Egenhofer OP with domain and range
	Req 6	Defines 8 RCC OP with domain and range
Geometry	Req 7	Defines Geometry <code>rdfs:subClassOf SpatialObject</code>
	Req 8	Defines OP <code>hasGeometry</code> and <code>defaultGeometry</code>
	Req 9	Defines 6 Data Properties: e.g: <code>dimension</code> , <code>isEmpty</code> , etc.
Serialization	Req 10-13	<code>wktLiteral</code> definitions & URI encoding
	Req 14	Defines <code>asWKT</code> to retrieve WKT literal
	Req 15-17	GMLLiteral should be accepted
	Req 18	Defines <code>asGML</code> to retrieve GML literal

Table 2.2: Requirements and implementations for vocabulary definitions in GeoSPARQL.

Each of the components described above has associated requirements. Concerning the vocabulary requirements, Table 2.2 summarizes the seventeen requirements presented in the GeoSPARQL draft document. GeoSPARQL requires the implementation of thirty three functions extensions for geospatial, and provides also cor-

responding relationships between topological functions. The functions are listed below:

- 8 for RCC: equals (EQ), disconnected (DC), externally connected (EC), partially overlapping(PO), tangential proper part inverse(TPPI), tangential proper part (TPP), non-tangential proper part inverse (NTPPI), non-tangential proper part (NTPP)
- 8 for Simple Features: equals, disjoint, intersects, touches, within, contains, overlaps and cross.
- 8 for Egenhofer relations: equals, disjoint, meet, overlap, covers, coveredBy, inside and contains
- 9 for functions that generate new geometries (Non-topological functions)

### 2.3.1.2 Geospatial Vocabularies and Topological Functions

Based on the GeoSPARQL requirements, we were interested in comparing some geospatial vocabularies to see how far they take already into account topological functions and which are the standard they followed among OpenGIS Simple Features (SF), Region Connection Calculus (RCC) and Egenhofer relations (see Figure 2.3.1.2). We find that the NeoGeo (Spatial and Geometry) and OS Spatial vocabularies have integrated in their modeling partial or full aspects of topological functions as summarized in Table 2.3.

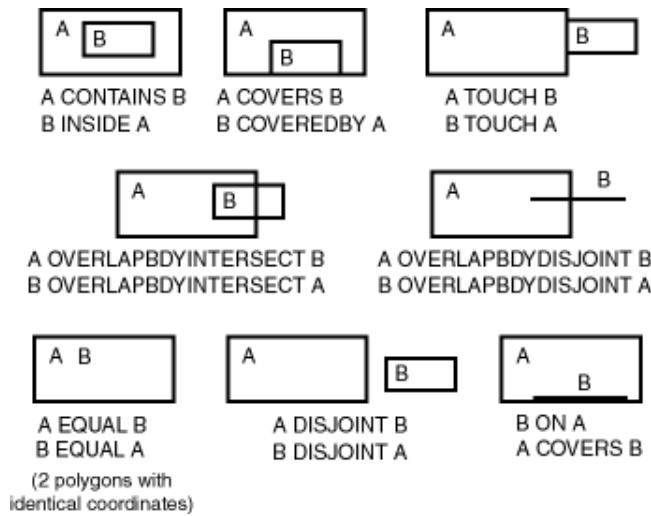


Figure 2.7: Illustration of the topological relations by Max Egenhofer.

As geodata has to be stored in triple stores with efficient geospatial indexing and querying capabilities, we also survey the current state of the art in supporting simple or complex geometries and topological functions compatible with SPARQL 1.1. Table 3.4 shows which triple stores can support part of the GeoSPARQL standard regarding serialization and spatial functions.

Geo-vocabulary	Topological Functions	GeoSPARQL Requirements	Standard Followed
Ordnance Survey Spatial	easting, northing, touches, within, contains	Part of Req 4	OpenGIS Simple Feature
Ordnance Survey Topography	contains, isContainedIn	Part of Req 4	OpenGIS Simple Feature
Place Ontology	in, overlaps, bounded_by	Subset of Req 4	N/A
NeoGeo Spatial	All RCC8 relations	Part of Req 3; Req 6	Region Connection Calculus (RCC)
NeoGeo Geometry	—	Req 10 to Req 14	N/A
FAO Geopolitical	isInGroup, hasBorderWith	—	—
OntoMedia Space	adjacent-below, adjacent-above, orbit-around, is_boundary-of, has-boundary	—	—

Table 2.3: Comparison of some geo-vocabularies with respect to the GeoSPARQL requirements.

### 2.3.2 Geospatial Vocabularies

#### 2.3.2.1 Vocabularies for Features

Modeling of features can be grouped into four categories depending on the structure of the data, the intended purpose of the data modeling, and the (re)-use of other resources.

1. One way to structure the features is to define high level codes (generally using a small finite set of codes) corresponding to specific types. Further, sub-types are attached to those codes in the classification. This approach is used in the Geonames ontology<sup>10</sup> for codes and classes (A, H, L, P, R, S, T, U, V), with each of the letter corresponding to a precise category (e.g.: A for administrative borders). Classes are then defined as `gn:featureClass` a `skos:ConceptScheme`, while codes are `gn:featureCode` a `skos:Concept`.
2. A second approach consists in defining a complete standalone ontology that does not reuse other vocabularies. A top level class is used under which a taxonomy is formed using the `rdfs:subClassOf` property. The LinkedGeoData ontology<sup>11</sup> follows this approach, where the 1294 classes are built around a nucleus of 16 high-level concepts which are: `Aerialway`, `Aeroway`, `Amenity`, `Barrier`, `Boundary`, `Highway`, `Historic`, `Landuse`, `Leisure`, `ManMade`, `Natural`, `Place`, `Power`, `Route`, `Tourism`, `Waterway`. The same approach is used for the French GeOnto ontology (Section 3.3), which defines two high-level classes `ArtificialTopographyEntity` and `NaturalTopographyEntity` with a total of 783 classes.
3. A third approach consists in defining several smaller ontologies, one for each sub-domain. An ontology network is built with a central ontology used to interconnect the different other ontologies. One obvious advantage of this approach is the modularity of the conceptualizing which should ease as much as possible the reuse of modular ontologies. The Ordnance Survey (OS) follows this approach providing ontologies for administrative regions<sup>12</sup>, for statistics decomposition<sup>13</sup> and for postal codes<sup>14</sup>. The `owl:imports` statements are used in the core ontology. Similarly, GeoLinkedData makes use of three different ontologies covering different domains.
4. A fourth approach consists in providing a *nearly flat list* of features or points of interest. This is the approach followed by popular Web APIs such as Foursquare types of venue<sup>15</sup> or Google Place categories<sup>16</sup>.

<sup>10</sup>[http://geonames.org/ontology/ontology\\_v3.0.rdf](http://geonames.org/ontology/ontology_v3.0.rdf)

<sup>11</sup><http://linkedgeodata.org/ontology>

<sup>12</sup><http://www.ordnancesurvey.co.uk/ontology/admingeo.owl>

<sup>13</sup><http://statistics.data.gov.uk/def/administrative-geography>

<sup>14</sup><http://www.ordnancesurvey.co.uk/ontology/postcode.owl>

<sup>15</sup><http://aboutfoursquare.com/foursquare-categories/>

<sup>16</sup>[https://developers.google.com/maps/documentation/places/supported\\_types](https://developers.google.com/maps/documentation/places/supported_types)

For this last approach, we have built an associated OWL vocabulary composed of alignments with other vocabularies.

### 2.3.2.2 Vocabularies for Geometry Shape

The geometry of a point of interest is also modeled in different ways. We include more details here to the survey started by Salas and Harth [16]:

- *Point representation*: the classical way to represent a location by providing the latitude and longitude in a given coordinate reference system (the most used on the Web is the WGS84 datum represented in RDF by the W3C Geo vocabulary). For example, Geonames defines the class `gn:Feature` a `skos:ConceptScheme` as a `SpatialThing` in the W3C Geo vocabulary.
- *Rectangle* (“bounding box”): which represents a location with two points or four segments making a geo-referenced rectangle. In this approach, the vocabulary provides more properties for each segment. The FAO Geopolitical ontology<sup>17</sup> uses this approach.
- *List of Points*: is a region represented by a collection of points, each of them being described by a unique RDF node identified by a lat/lon value. The `Node` class is used to connect one point of interest with its geometry representation. The POI are modeled either as `Node` or as `Waynode` (surfaces). This approach is followed by the LinkedGeoData [12].
- *Sequence of Points*: is represented by a group of RDF resources called a “curve” (similar to LineString of GML). The POI is connected to its geometry by the property `formedBy` and an attribute `order` to specify the position of each node in the sequence. This approach is used in GeoLinkedData [14].
- *Literals*: the vocabulary uses a predicate to include the GML representation of the geometry object, which is embedded in RDF as a literal. This approach is followed by Ordance Survey [13].
- *Structured representation*: the geometry shape is represented as a typed resource. In particular, polygons and lines are represented with an RDF collection of basic W3C Geo points. This approach is used by the NeoGeo vocabulary<sup>18</sup>.

### 2.3.3 Georeferencing Data on the Web

Georeferencing data either by direct or indirect spatial reference requires some reference datasets that can be used as the spatial frame for anchoring these thematic data. Especially, it requires data on both CRSs and named places, which must be published on the Web of Data.

---

<sup>17</sup><http://www.fao.org/countryprofiles/geoinfo/geopolitical/resource/>

<sup>18</sup><http://geovocab.org/doc/neogeo/>

### 2.3.3.1 Identifying and describing CRSs on the Web

In order to fulfill the need for CRS identification and description on the Web, OGC maintains a set of URIs for identifying the most commonly used CRS. While very useful, the main disadvantage of this proposal is that the URIs defined by OGC are not very intuitive for users who are not familiar with Spatial Reference System Identifiers defined by geographic information authorities like OGC or EPSG, such as “4326” (which actually refers to a WGS84 CRS defined by the EPSG). Moreover, many CRS commonly used locally, such as deprecated French projected CRS, are not available in that registry. In addition to OGC proposal, several registries have been proposed by the geographic information community for cataloguing existing CRSs. The EPSG Geodetic Parameter Registry<sup>19</sup> allows querying the Geodetic Parameter Dataset gathered by the EPSG. CRSs can be retrieved by name, by code, by type or by coverage area, and their characteristics are displayed on a HTML form. Unfortunately, there is no direct access to these data through dereferenceable URIs.

### 2.3.3.2 CRS requirements for the French Territory

As explained in Section 1.1, making explicit the CRS used in a given dataset is a very important issue when dealing with direct location data. This is especially important in the field of geographical information where different CRSs are commonly used due to technical or legal requirements. For INSPIRE Directive, CRS are considered as reference data used for linking thematic data [32], and must be described according to the ISO 19111 standard. To be consistent with Linked Data principles, CRS should be identified by URIs, like in OGC proposal. Moreover, as Linked Data users are not always familiar with CRS identifiers commonly used within the geographic information community, URIs used to identify CRS should use more intuitive names<sup>20</sup> by following some best practices [33]. Finally, consistently with our goal of contributing to better georeferencing of data on the French territory, we need access to the descriptions of all French CRSs, including some deprecated but still used CRSs like “Lambert 1”.

Prefix	URI
geofla	<a href="http://data.ign.fr/def/geofla#">http://data.ign.fr/def/geofla#</a>
geom	<a href="http://data.ign.fr/def/geometrie#">http://data.ign.fr/def/geometrie#</a>
ignf	<a href="http://data.ign.fr/def/ignf#">http://data.ign.fr/def/ignf#</a>
rgeofla	<a href="http://data.ign.fr/id/geofla/">http://data.ign.fr/id/geofla/</a>
topo	<a href="http://data.ign.fr/def/topo#">http://data.ign.fr/def/topo#</a>
rtopo	<a href="http://data.ign.fr/id/topo/">http://data.ign.fr/id/topo/</a>

Table 2.4: URI schemes and conventions used for vocabularies and resources.

<sup>19</sup><http://www.epsg-registry.org/>

<sup>20</sup><http://philarcher.org/diary/2013/uripersistence/>

The Information and Service System for European Coordinate Reference Systems<sup>21</sup> provides access to ISO 19111 standard-based descriptions of the main European CRSs but has the same limitation as the EPSG registry: access to the descriptions is not allowed by URI, but only through a cartographic interface. The <http://spatialReference.org> initiative aims at allowing the use of URI-based references to spatial reference systems, including some CRSs defined and maintained by IGN France. However, the proposed URL policy is not very intuitive. As an example, this URL identifies the projected system defined by IGN France, Lambert 93: <http://spatialreference.org/ref/sr-org/7527/>. Moreover, the definitions of some deprecated CRSs such as Lambert zone projected CRSs (which are still used in some datasets) seem to be referenced only for the authority EPSG and not for IGNF, which also maintains a registry of CRSs. ISO 19111 standard-based definitions of all CRSs defined and maintained by IGN France are published in an XML file<sup>22</sup>. References to equivalent definitions provided by the EPSG registry are explicitly stated with EPSG SRID. CRSs are identified by URIs using short names instead of numeric codes. For example, <http://registre.ign.fr/ign/IGNF/crs/NTFLAMB2E> is the URI designed for the “Lambert 2 étendu” projected system. Indeed “NTFLAMB2E” is used to identify the projected system “Lambert 2 étendu” which is based on NTF (New French Triangulation) geodetic reference system. Unfortunately, this registry is still in evolution and its URIs are not dereferenceable yet.

REGION	COORDINATE SYSTEM	ELLIPSOID	PROJECTION SYSTEM	ALTIMETRY SYSTEM
FRANCE METROPOLITAN	RGF93 RGM04 (ITRF2000)	IAG GRS 1980 IAG GR80	Lambert 93 and CC 9 Zones UTM 38 South	
MAYOTTE				SHOM 1953
GUYANE	RGFG95	IAG-GRS 1980	UTM 21 22 North	
MARTINIQUE	WGS84	IAG-GRS 1980	UTM 20 North	
GUADELOUPE	WGS84	IAG-GRS 1980	UTM 20 North	
LA RÉUNION	RGR92	IAG-GRS 1980	UTM 40 South	GGR 99
NOUVELLE-CALEDONIE	ITRF90	IAG-GRS 1980		
POLYNÉSIE	RGPF	IAG-GRS 1980	UTM 5, 6, 7 and 8 South	Tahiti IGN 1966
WALLIS ET FUTUNA	MOP87	International 1924		
SAINT-PIERRE ET MIQUELON	RGM01 (ITRF2000)	IAG GRS 1980	UTM 21 North	Danger 1950
ILE CLIPPERTON	Marine 1967	International	UTM 12 South	

Figure 2.8: Coordinate Reference Systems used in France

As no existing registry fulfilled all our requirements, we have designed a vocabu-

<sup>21</sup><http://www.crs-geo.eu>

<sup>22</sup><http://librairies.ign.fr/geoportail/resources/IGNF.xml>

lary<sup>23</sup>, inspired from the ISO 19111 schema for CRSs description. Then we have converted IGNF CRSs registry into RDF, and published this dataset on the Web with the Datalift platform<sup>24</sup>. Therefore, the description of the “NTF Lambert 2 étendu” projected CRS can be retrieved at this URL <http://data.ign.fr/id/ignf/crs/NTFLAMB2E>.

### 2.3.3.3 Direct georeferencing of data on the Web

Modeling direct location information such as coordinates or vector data geometries in RDF still poses some challenges. In [34], we have conducted a survey of the vocabularies used for representing geographical features from vocabularies of feature types to vocabularies for geometric primitives which provide ways for representing extents, shapes and boundaries of those features. Most of vocabularies dedicated to geometry representation reuse W3C Geo vocabulary which allows only WGS84 coordinates, such as NeoGeo<sup>25</sup>. With the rise of the Open Data movement, more and more publishers including governments and local authorities are releasing legacy data that are georeferenced using others CRSs. For example, IGN France releases data using different projected CRSs depending on the geographic extent of each dataset. In order to overcome this limitation on CRSs, the vocabulary designed by OGC GeoSPARQL standard does not reuse W3C Geo vocabulary but proposes another class, “Point”, instead. Geometries of geographical data represented in RDF with the GeoSPARQL vocabulary are represented by literals encoded consistently with other OGC standards. `gsp:wktLiteral` and `gsp:gmlLiteral` are thus respectively derived from Well-Known Text and GML encoding rules. In `wktLiteral` and `gmlLiteral`, the CRS used to define the coordinates of the point is identified by a dereferenceable URI which is explicitly stated at the beginning of the literal. This way of associating coordinate reference systems with geometries has the advantage of being consistent with Linked Data principles: each CRS is identified with a dereferenceable URI. The main drawback is that such literals cannot be easily queried with SPARQL, unless using regular expression-based filters. To overcome this limitation, we associate each geometry to the CRS used by its coordinates with the property `geom:crs` in the geometry vocabulary presented in Section 2.4.2.

### 2.3.3.4 Indirect georeferencing of data on the Web

**Location Vocabulary:** The Location Core Vocabulary<sup>26</sup> provides structure to describe a location in three different ways: by using a place name, a geometry or an address. The vocabulary is heavily based on the definition of ISO 19112 of a location, as “an identifiable geographic place”. A part from using simple string labels or names, the vocabulary provides a property to allow a location to be defined by a

---

<sup>23</sup><http://data.ign.fr/def/ignf>

<sup>24</sup>A service to lookup CRS in RDF can be found at <http://www.eurecom.fr/~atemezin/ignf-lookup/>

<sup>25</sup><http://geovocab.org/doc/neogeo/>

<sup>26</sup><http://www.w3.org/ns/locn>

URI, such as GeoNames or DBpedia URI. The geographic name used for a spatial object is consistent with the INSPIRE Data Specification on Geographical Names [32]. The Geometry Class denotes the notion of geometry at a conceptual level, and can be encoded in different formats including WKT, GML, KML, RDF+WKT/GML (GeoSPARQL), RDF (WGS84 lat/long, schema.org) and GeoHash URI references. In addition, the geometry property can be associated to either a literal (such as WKT, GML or KML) or a geometry class (e.g., `ogc:Geometry` and its subclasses, `geo:Point`, `schema:GeoCoordinates` and `schema:GeoShape`, a GeoHash URI reference). However, the CRS identifier of the geometry is either embedded in the literal (e.g., WKT, GML) or implicit in the more structured serialization (e.g., WGS84 lat/long), schema.org, GeoHash).

**Datasets using indirect georeferencing** Modeling indirect location information such as administrative units or named points of interest in RDF is preferably done by identifying such geographic features with URIs and describing them by their properties, so that they can be referenced by other datasets. This is the case in one of the most reused datasets of the Web of Data, namely Geonames<sup>27</sup>. However, so far there are very few authoritative datasets covering the French territory on the Web of Data. A simple example is the current resource for *Paris* in the French DBpedia<sup>28</sup>. The department's name associated to this resource is a literal named “Paris” and the different arrondissements composing the city are modeled as `skos:Concept` instead of `dbpedia-owl:Place`. Even Geonames data remain very limited, as French administrative units are provided as simple geometries (POINT). The “Official Geographic Code”<sup>29</sup> published by the French Statistical Institute (INSEE) is the most up-to-date and accurate dataset on French administrative units, but unfortunately it contains no geometrical description of their boundaries. The consequence here is therefore the lack of a baseline during the mapping process for application developers trying to consume specific data coming from France. Datasets describing administrative units, points of interest or postal addresses with their labels and geometries, and identifying these features with URIs could be beneficial not only for georeferencing other datasets, but also for interlinking datasets georeferenced by direct and indirect location information.

## 2.4 Vocabularies for Geometries and Feature Types

Direct georeferencing of data implies representing coordinates or geometries and associating them to a CRS. This requires vocabularies for geometries and CRSS. Further, indirect georeferencing of data implies associating them to other data on named places. Preferably, these data on named places should be also georeferenced by coordinates in order to serve as the basis for data linking between indirectly and

---

<sup>27</sup><http://sws.geonames.org/>

<sup>28</sup><http://fr.dbpedia.org/resource/Paris>

<sup>29</sup><http://rdf.insee.fr/sparql>

directly georeferenced datasets. In this section, we present the vocabularies that we have defined and reused for geographic data publishing. This requires reference geographic data on named places and therefore vocabularies for describing feature types and their properties.

### 2.4.1 Motivation

In [34], we already surveyed numerous vocabularies for representing geographical features and their geometries, either using a literal (e.g. by using the `wktLiteral` datatype) or a structured representation à la NeoGeo. We concluded the survey with some recommendations for geometry descriptions:

- the distinction of geometry versus feature and a property linking both classes (e.g. for attaching provenance information on how some points of a geometry have been collected),
- the ability to represent structured geometries (e.g. for performing simple spatial queries on the data, even when they are stored in a triple store that do not implement the GeoSPARQL standard),
- the integration of any coordinate reference system (e.g. for allowing projected coordinates for cartographic purposes).

In addition to these recommendations, we also suggest that the domain of the property used to link a feature to its geometry should be left empty in order to accept links between any type of resource and a geometry. This would be useful for example, to associate a person to the coordinates of their birthplace.

### 2.4.2 A vocabulary for geometries

On the current usage of georeferencing resources on the Web of Data, it is assumed that the coordinates should be in WGS84, and hence the definition of the point. However, publishers might have data in different CRSs according to the location. Thus, our proposal is to define a more generic class for a `POINT` with the benefit of choosing the CRS of the underlying data, as depicted in the Listing 2.3. The naming convention used for the `geom` vocabulary follows the terms used by the Simple Features vocabulary. The French translation of terms is based on the glossary of multilingual terminology of ISO/TC 211 available at <http://www.isotc211.org/Terminology.htm>.

**Axiom 1.** (*Geometry*): *A resource of type `geom:Geometry` should be associated to exactly one resource of type `ignf:CRS` via the property `geom:crs`.*

**Axiom 2.** (*Class Point*): **A `POINT` is a subclass of a `GEOMETRY`.**

Regarding alignments with some existing vocabularies, the class `geom:Geometry` is a subclass of both `sf:Geometry` and `ngeo:Geometry`. The class contains in addition

the property `geom:crs`. Moreover, it is possible to obtain equivalences between data modeled with `ngeo` vocabulary and `geom` vocabulary. The following SPARQL query make it possible:

```

1 1 CONSTRUCT {
2   []a geom:Point;
3   owl:sameAs ngeo:Point.
4
5
6 }WHERE {
7   []a geom:Geometry;
8   geom:Point;
9   geom:systCoord
10  <http://data.ign.fr/id/ignf/crs/WGS84G>.
11 }
```

Listing 2.2: SPAQRL Query for creating sameAs links between data modeled with `ngeo` and `geom` vocabularies

**Axiom 3.** (*Points*): *An instance of the class `geom:Point` is associated with exactly one instance of `ignf:CRS` via the property `geom:crs`. An instance of a `geom:Point` has exactly one coordinate X and exactly one coordinate Y. The coordinates are `xsd:double` and referred to the following properties:*

- *`geom:coordX` in an ellipsoidal CRS, refers to the longitude of a point and within a projected CRS and the value of false easting of a point.*
- *`geom:coordY` in an ellipsoidal CRS, refers to the latitude of a point and within a projected CRS and the value of false northing of a point.*

```

1 geom:Point a owl:Class;
2 rdfs:label "Point"@en, "Point"@fr;
3 rdfs:subClassOf geom:Geometry;
4 owl:equivalentClass
5 [a owl:Class ;
6 owl:intersectionOf
7 ([a owl:Restriction;
8   owl:onDataRange xsd:double;
9   owl:onProperty geom:coordY;
10  owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
11  [a owl:Restriction;
12  owl:onDataRange xsd:double;
13  owl:onProperty geom:coordX;
14  owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger])
15 ];
16 rdfs:subClassOf sf:Point.
```

Listing 2.3: Definition in Turtle of the axiom defining a POINT.

**Axiom 4.** (*PointsList*): *A `geom:PointsList` is a subclass of `rdf:List`. An instance of `geom:PointsList` is composed of only instances of type `geom:Point`.*

### Extending GeoSPARQL vocabulary

In order to fulfill these recommendations, we have developed a new vocabulary that re-uses and extends the existing vocabularies for representing geometries, namely:

- <http://www.opengis.net/ont/geosparql#> (prefix `gsp`): This vocabulary provides the basic concepts to represent geographical data such as `SpatialObject`, `Feature` or `Geometry`. A `Feature` is linked to a `Geometry` via the relation `gsp:hasGeometry`. The geometries are typed strings (`gsp:gmlLiteral` or `gsp:wktLiteral` corresponding respectively to the properties `gsp:asGML` and `gsp:asWKT`). The vocabulary contains also spatial functions.
- <http://www.opengis.net/ont/sf#> (prefix `sf`): This vocabulary is based on the OGC standard Simple Features Access [35]. The class `sf:Geometry` is a subclass of `gsp:Geometry`.

Reusing and extending GeoSPARQL Simple Features vocabulary with structured geometries à la NeoGeo enables us to represent geometries both with GeoSPARQL compliant literals and with structured geometries that can be handled easily with SPARQL. The extension for structured geometries consists in defining a subclass for each class from the `sf` vocabulary, and defining properties to associate its instances with a CRS and coordinates or other suitable geometric primitives. For example, the class `geom:Point` is a subclass of `sf:Point`. An instance of `geom:Point` is associated with exactly one instance of `ignf:CRS` via the property `geom:crs`, and it has exactly one coordinate X and exactly one coordinate Y. It can also have a Z coordinate. The coordinates are `xsd:double` and correspond to the properties `geom:coordX:`, `geom:coordY:` and `geom:coordZ:` respectively. Other complex geometries are also defined, such as Linestrings, LinearRings, Polygons or Multi-Polygons. Their definitions are based on the class `geom:Point`. As an example, an instance of `geom:Linestring` is defined as an instance of `geom:PointsList` which is an ordered `rdf>List` of instances of `geom:Point` designated by the property `geom:points`.

We have also defined a property `geom:geometry` with an empty domain. Thus, our proposal defines a more generic class for a `POINT` with the benefit of choosing the CRS of the underlying data. Figure 2.9 gives an overview of the relationships between the high level concepts with geometries, CRS and topographic features.

#### 2.4.3 A vocabulary for Topographic Entities

A topographic entity `topo:EntiteTopographique` is the class associated to a phenomenon with an associated location on the Earth<sup>30</sup>. `topo`<sup>31</sup> vocabulary contains 8 direct subclasses: Buildings and structures, Cableway transport line, Energy transport infrastructure, Inland hydrographic feature, Relief feature, Road transport

<sup>30</sup>The names used in the description of the classes or properties are `rdfs:label[@lang='en']` of the actual URIs in the models.

<sup>31</sup><http://data.ign.fr/def/topo>

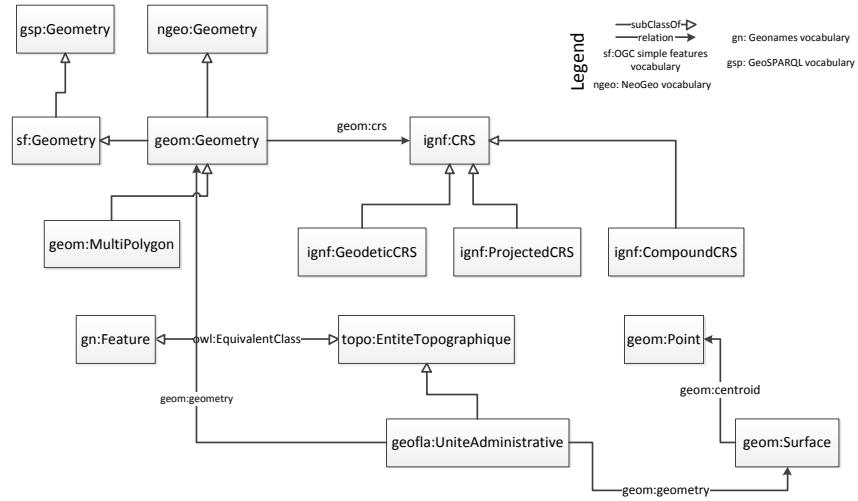


Figure 2.9: High level classes of the `ignf`, `geom` and `topo` vocabularies; relationships between them and mappings with external vocabularies.

feature, Transport by rail feature, Vegetation area and Working area of interest. Furthermore, we use a second level of classes to model direct subclasses of the previous ones. For example, a Building is specialized in 6 different classes: Cemetery, Sports ground, Structure, Tank and Taxiway. A building is then connected with the property `topo:typeDeBatiment` (type of building) to different ‘nature’ of buildings, which are modeled as SKOS concepts. SKOS is intensively used to easily group concepts into different schemes (using `skos:hasTopConcept`) and provide semantic relationships (e.g: `skos:broader`, `skos:narrowMatch`) among them. This gives flexibility in the model by defining few high level classes and restrictions for the properties. Table 2.5 gives a listing of the different SKOS namespaces defined to capture some high level concepts.

Indirect georeferencing of resources on the Web requires reference geographic data on named places and therefore vocabularies for describing feature types and their properties. Therefore, we have chosen to publish a reference dataset on administrative units called GEOFLA®, which is already available in GIS format under an Open Data license. We have also made tests of data conversion and interlinking with another largest dataset on French names places. We have produced and published two vocabularies to describe these datasets, to make sure that all concepts and properties needed would be available. In the GEOFLA® vocabulary, 5 classes have been defined: commune, canton, arrondissement, department and region. In the BD TOPO® vocabulary<sup>32</sup> 35 main classes have been defined. They represent the main types of geographic features represented in the BD TOPO® database.

<sup>32</sup><http://data.ign.fr/def/topo>

URL	Description
cdtopo:typedezai:liste	list of different areas of activities and interest
cdtopo:typedebatiment:liste	list of types of buildings
cdtopo:typedeterraindesport:liste	list of types of sports ground
cdtopo:typedeconstruction:liste	list of types of construction
cdtopo:typedereservoir:liste	list of types of tanks
cdtopo:typedevegetation:liste	list of types of vegetation
cdtopo:typederelief:liste	list of types of relief features
cdtopo:typedevoieferree:liste	list of types of railway track
cdtopo:typedetransportcable:liste	list of types of cableway transport
cdtopo:typedefranchissement:liste	list of types of crossing
cdtopo:typederoute:liste	list of types of road
cdtopo:typedereservoir:liste	list of types of waterhole
cdtopo:typedelaisse:liste	list of types of tide line

Table 2.5: List of concept schemes used in the topographic ontology.

In both vocabularies, properties have been defined based on the attributes of their related classes in the databases. The geographic feature types defined as values of attributes “nature” are modeled as instances of `skos:Concept`. SKOS is intensively used to easily group concepts into different schemes (using `skos:hasTopConcept`) and provide semantic relationships (e.g: `skos:broader`, `skos:narrowMatch`) among them. We also provide alignments with Geonames vocabulary, where `topo:Place` is subclass of `gn:S` and `owl:sameAs` linked concepts.<sup>33</sup>

All the classes are defined as subclasses of `topo:EntiteTopographique` which defines the representation of a real world entity associated to a location relative to the Earth, consistently with ISO TC 211 and OGC standards. The GEOFLA®’s application schema is composed of classes representing different types of french administrative units, namely communes, cantons, arrondissements and departments. In `geofla` vocabulary, we add a class `Region` from the instances of the class department via two attributes that precise to what region each instance of department belongs. Their properties are defined based on the attributes of their related classes in GEOFLA® database.

A Commune has an attribute called “*nature*” whose enumerated values precise whether the commune is the capital of a bigger administrative unit, modeled in the vocabulary by the ObjectProperty `geofla:statut` with range `skos:Concept` defined in this specific `skos:ConceptScheme`

`http://{BASE}/codes/geofla/typedecommune/liste` pointing to the different types of French administrative unit’s capital.

<sup>33</sup><https://github.com/gatemazing/ign-iswc2014/blob/master/vocabularies/mappingsGeonames.ttl>

#### 2.4.4 Discussions

The alignment of existing taxonomies for describing geodata enables interoperability of symbolic descriptions. The need for a better choice of geometric structure, typically the choice between literal versus structured representations depends on four criteria:

1. the coverage of all the complex geometries as they appear in the data;
2. a rapid mechanism for connecting “features” to their respective “geometry”;
3. the possibility to serialize geodata into traditional formats used in GIS applications (GML, KML, etc.) and
4. the choice of triple stores supporting as many functions as possible to perform quantitative reasoning on geodata.

It is clear that a trade-off should be made depending on the technological infrastructure (e.g: data storage capacity, further reasoning on specific points on a complex geometry). The following points helps understanding better some of the challenges:

- **Complex Geometry Coverage:** We have seen that on the Web of Data, there is little modeling of geodata with their correct shape represented as lines or polygons. However, some content providers (e.g. IGN) need to publish all types of geodata including complex geometries representing roads, rivers, administrative regions, etc. Two representations are suitable: the *OS Spatial* and *NeoGeo* ontologies (Table 2.2). A direct representation of the GeoSPARQL vocabulary is also suitable.
- **Features connected to Geometry:** In modeling geodata, we advocate a clear separation between the features and their geometry. This is consistent with the consensus obtained from the different GeoVocamps<sup>34</sup> and the outcome of this approach is expressed in the modeling design of NeoGeo. The top level classes `spatial:Feature` and `geom:Geometry` are connected with the property `geom:geometry`.
- **Literal versus structured Geometry:** Decomposing a LINE or a POLYGON results in an “explosion” in the size of the dataset and the creation of numerous blank nodes. However, sharing points between descriptions is a use case with many real-world applications. IGN has such use cases and the natural solution at this stage is to consider reusing the NeoGeo ontology . The choice of the triple store (e.g., Virtuoso vs Open Sahara) is not really an issue, as the IndexingSail<sup>35</sup> service could also be wrapped on top of Virtuoso to support full OpenGIS Simple Features functions<sup>36</sup>.

---

<sup>34</sup><http://www.vocamp.org>

<sup>35</sup><https://dev.opensahara.com/projects/useekm/wiki/IndexingSail>

<sup>36</sup><http://www.opengeospatial.org/standards/sfs>

## Publishing Structured Geometries from Geographic Data

The vocabulary for geometry reused by a geodata converter that takes traditional GIS data as input and outputs RDF data with geometries defined both with a `gsp:wktLiteral` and with a structured representation is compliant with the vocabularies presented in the previous sections. Geometries are automatically associated with the chosen CRS. This converter is implemented as a plugin of the Datalift platform (Section 3.2.4) and can be reused easily for geographic data publishing purposes. The result in the RDF data is a set of triples where each member of a complex geometry can be identified by a single URI with properties for latitude and longitude values.

## 2.5 Summary

In this chapter, we reviewed the formats and the different vocabularies used to model geospatial data on the Web, based on direct and indirect georeferencing. Afterwards, we identified some limitations the lack of an explicit CRS reference on the data. We then proposed a REST service for converting between different CRSs to help the publishers to be able to handle different projections of their datasets. Then, we proposed and implemented three vocabularies for geometries, CRSs and topographic entities (`geom`, `ingf`, `topo`). The vocabularies extend existing ones and integrate two additional advantages: an explicit use of CRS identified by URIs for geometry, and the ability to describe structured geometries in RDF. Some of our findings and model description are under discussion for standardization at the W3C, such as to extend GeoSPARQL to handle CRS, or the extension of the location and address vocabulary [36] to support different CRSs properties. In the next chapter, we will go through the process of the conversion and the publication of geospatial data on the Web.

## CHAPTER 3

# Publishing, Interlinking and Querying Geodata

---

*“If you’re a geospatial developer, AJAX is not a domestic cleaning product.  
If you’re a web dev, a polygon is not a dead parrot”<sup>1</sup>.*

Steve Peters  
(UK Government’s Department  
for Communities and Local Government)

## Introduction

So far, the Web of Data has taken advantage of geocoding technologies for publishing large amounts of data. For example, Geonames provides more than 10 million records (e.g. 5 million resources of the form <http://sws.geonames.org/10000/>) while LinkedGeoData has more than 60,356,364 triples. All the above mentioned data is diverse in nature and format, access point (SPARQL endpoint, Web service or API), the entities they represent and the vocabularies used for describing them. Table 3.1 summarizes for different providers the quantity of geodata available (resources, triples) and how the data can be accessed.

To publish geospatial data on the Web, conversion tools are required to transform native formats into RDF. Those tools are based on Open Geospatial Consortium (OGC) libraries for extracting features, and specific vocabularies as input to build the RDF. Once the dataset is converted, it has to be stored in an efficient way, interconnected to other datasets, and then consumed using SPARQL queries. Moreover, all those steps might be integrated in frameworks that ease the overall process of publishing geodata. In this chapter, we first discuss different representation on the Web of the 7<sup>th</sup> arrondissement of Paris in DBpedia, Geoname and LinkedGeoData (Section 3.1). Section 3.2 provides an overview of four tools for converting geospatial data into RDF, with a brief discussion on their limitations. Criteria to interlink geodata is then described in Section 3.3, followed by a survey on triple stores (Section 3.4). We then describe Datalift, a tool for managing the workflow of publishing raw data to RDF (Section 3.5), followed by our contributions of publishing French authoritative datasets using Datalift. An evaluation on time execution of built-on geo-functions of three triple stores is presented in Section 3.7, and a brief summary (Section 3.8) to conclude the chapter.

---

<sup>1</sup><http://www.w3.org/2014/03/lgd/report>

### 3.1 Current Representation of Geodata on the Web

In this section, we review the modeling of the 7<sup>th</sup> arrondissement of Paris, France in different geospatial datasets on the Web. The 7<sup>th</sup> arrondissement of Paris is one of the 20 arrondissements (administrative districts) of the capital city of France. It includes some of Paris's major tourist attractions such as the Eiffel Tower, some world famous museums (e.g.: musée d'Orsay) and contains a number of French national institutions, including numerous government offices<sup>2</sup>. We use the 7<sup>th</sup> arrondissement throughout this chapter to highlight the diversity of representations one can use for geographical entity. We assume that this district should be modeled as a POLYGON composed of a number of POINTs needed to “interpolate” its effective boundaries. We assume the use of the WGS84 geodetic system [28].

Provider	#Geodata	Data access
DBpedia	727,232 triples	SPARQL endpoint
Geonames	5,240,032 (feature).	API
LinkedGeoData	60,356,364 triples	SPARQL endpoint, Snorql
Foursquare	N/A	API
Freebase	8,5MB	RDF Freebase Service
Ordnance Survey(Cities)	6,295 triples	Talis API
GeoLinkedData.es	101,018 triples	SPARQL endpoint
Google Places	N/A	Google API
GADM	682,605 triples	Web Service
NUTS	316,238 triples	Web Service
IGN experimental	629,716 triples	SPARQL endpoint
LOD Greek	634 KTriples	SPARQL endpoint

Table 3.1: Geodata by provider and their different access type, either API, Web service or SPARQL endpoint.

#### 3.1.1 DBpedia Modeling

We provide below an excerpt of the DBpedia 3.8 description for the 7<sup>th</sup> arrondissement of Paris:

```
dbpedia:7th_arrondissement_of_Paris a gml:_Feature ;
  a <http://dbpedia.org/class/yago/1900SummerOlympicVenuEs>
  rdfs:label "7. arrondissementti (Pariisi)"@fi; (14 different languages)
  dbpprop:commune "Paris" ;
  dbpprop:departement dbpedia:Paris ;
  dbpprop:region dbpedia:Ile-de-France_(region) ;
  grs:point "48.85916666666667 2.3127777777777778" ;
  geo:geometry "POINT(2.31278 48.8592)" ;
  geo:lat "48.859165"^^xsd:float;
```

<sup>2</sup><http://sws.geonames.org/2988760>

```
geo:long "2.312778"^^xsd:float.
```

First, we observe that the type `gml:_Feature` and the property `grs:point` are not resolvable and are not associated to any OWL ontologies published providing descriptions and definitions. Second, the property `geo:geometry` used by DBpedia is not defined in the WGS84 vocabulary. For the geometry, the 7th arrondissement is a simple POINT defined by a latitude and a longitude.

### 3.1.2 Geonames Modeling

In Geonames, the 7<sup>th</sup> arrondissement is considered to be a 3<sup>rd</sup> order administrative division, represented by a POINT for the geometry model. The RDF description of this resource gives other information such as the alternate name in French, the country code and the number of inhabitants.

```
gnr:6618613 a gn:Feature ;
  gn:name "Paris 07";
  gn:alternateName "7eme arrondissement";
  gn:featureClass gn:A [
    a skos:ConceptScheme ;
    rdfs:comment "country, state, region ..."@en .
  ] ;
  gn:featureCode gn:A.ADM4 [
    a skos:Concept ;
    rdfs:comment "a subdivision of a 3rd order admin division"@en .
  ];
  gn:countryCode "FR";
  gn:population "57410";
  geo:lat "48.8565";
  geo:long "2.321".
```

### 3.1.3 LinkedGeoData Modeling

In LinkedGeoData dataset, the district is a `lgdo:Suburb` which is subClass of `lgdo:Place`. Its geometry is still modeled as a POINT and not as a complex geometry of type POLYGON as we could have expected for this type of spatial object.

```
lgd:node248177663 a lgdo:Suburb ;
  rdfs:label "7th Arrondissement"@en , "7e Arrondissement" ;
  lgdo:contributor lgd:user13442 ;
  lgdo:ref%3AINSEE 75107 ;
  lgdp:alt_name "VIIe Arrondissement" ;
  georss:point "48.8570281 2.3201953" ;
  geo:lat 48.8570281 ;
  geo:long 2.3201953 .
```

### 3.1.4 Discussion

These samples from DBpedia, Geonames and LinkedGeoData give an overview of at least three different views of the same reality, in this case the district of the

$7^{th}$  Arrondissement in Paris. Regarding the “symbolic representation”, two datasets opted for “Feature” (DBpedia and Geonames) while LGD classifies it as a “Suburb” or “Place”. They all represent the shape of the district as a POINT which is not very appropriate if we consider a query such as *show all monuments of international importance located within the 7<sup>th</sup> arrondissement*. To address this type of query and more complicated ones, there is a need for more advanced modeling as we describe in the next section.

## 3.2 Existing Tools for Converting Geospatial Data

To address the need for converting geodata into a graph model such as RDF, some tools for generating RDF data from legacy geospatial datasets have been proposed. The differences among the tools are based on four main factors:

1. **Input format:** The different types of formats accepted as input of the tool;
2. **Vocabulary:** The vocabulary used to handle the geometry shape of the spatial data during the RDF conversion step;
3. **CRS converter:** The presence or not of a CRS converter between different CRSs in the final output of the tool;
4. **Output:** The type of serialization in RDF, and more importantly the choice between structured geometry, OGC compliant (WKT, GML literals) or both for the geometry part of the features.

Almost all the conversion tools use the OGC libraries for parsing and extracting features from shape formats, such as GDAL (Geospatial Data Abstraction Library) and GeoTools. We know describe each of the four tools: Geometry2RDF, TripleGeo, shp2GeoSPARQL and GeomRDF in turn.

### 3.2.1 Geometry2RDF

Geometry2RDF [14] is a Java-based tool<sup>3</sup> that generates RDF triples from geometrical information, which can be available in GML or WKT. The tool takes as input any ESRI shapefiles, spatial DBMS (Oracle, PostgreSQL, MySQL, etc), transforms the data into GML (using GeoTools<sup>4</sup>) and then generates RDF (using Jena<sup>5</sup>) consistent with the NeoGeo vocabulary. The default CRS used for the geometry is WGS84. The architecture is modular enough to run as a standalone platform or as a library.

---

<sup>3</sup><https://github.com/boricles/geometry2rdf>

<sup>4</sup><http://www.geotools.org>

<sup>5</sup><https://jena.apache.org/>

### 3.2.2 TripleGeo

TripleGeo [37] is an Extract-Transform-Load (ETL)<sup>6</sup> tool derived from Geometry2RDF (Cf. Section 3.2.1) to transform a variety of geospatial databases and shapefiles (including KML and INSPIRE compliant files) into RDF triples. Triples can be exported according to the GeoSPARQL standard, the WGS84 vocabulary and the Virtuoso RDF vocabulary<sup>7</sup>. In addition TripleGeo allows on-the-fly reprojection between CRSs, e.g., transform geometries from GreekGrid87 (a local CRS) into WGS84 (used for GPS locations).

### 3.2.3 shp2GeoSPARQL

shp2GeoSPARQL [38] is also an extension of Geometry2DF which transforms Shapefiles into RDF in the cadastral domain using ISO 19152 [39] (Land Administration Domain Model) and GeoSPARQL. The geometries obtained from shp2GeoSPARQL are consistent with the GeoSPARQL geometry vocabulary.

### 3.2.4 GeomRDF: Datalift tool for Converting Geodata

GeomRDF [40] is a tool developed within the Datalift platform that transforms geospatial dataset from traditional GIS formats into RDF, and overcome the limitations of the existing tools mentioned in Section 3.2.5. GeoRDF is based on a vocabulary that reuses and extends GeoSPARQL and NeoGeo so that geometries can be defined in any CRS, and represented both as structured geometries and GeoSPARQL standard compliant. GEOMRDF is composed of three components:

- **input parsers:** This component extracts from the different input format all the features and their descriptions.
- **feature parsers:** This component extracts for all the features their properties depending on their type, either thematic (attributes and properties of a geographic entity) or geometric (the geometry associated to the entity). For example, in the case of a Multipolygon, the parser stores first all the Polygons composing the original Multipolygon. Then, it stores the exterior and the eventual LinearRings for each Polygon, as well as the points included in the LinearRings. Finally the coordinates of each point are also stored. At this stage, GeoRDF provides a “*CRS reprojection*” functionality, which consists of transforming on-the-fly between different CRSs, before storing the geometries.
- **RDF Builder:** This module is responsible for generating RDF triples according to the `geom` vocabulary for geometry, `geofla` and `topo` vocabularies for different topographic entities.

---

<sup>6</sup><https://github.com/GeoKnow/TripleGeo>

<sup>7</sup><http://docs.openlinksw.com/virtuoso/rdfsparqlgeospat.html>

GeomRDF is implemented as a module of the RDF publication platform Datalift. Moreover, it has been validated against the French Administrative Units dataset available at <http://data.ign.fr/id/sparql>. GeomRDF can also be used as a stand-alone library and can be accessed at <https://github.com/fhamdi/GeomRDF>. Listing 3.1 presents a snippet in Turtle [41] of GeomRDF for geometry of the city of Nice (France). It also contains the structured modeling of a MULTIPOLYGON as a set of POLYGON, containing POINTs in an LINEARING.

```

1 @prefix geom:<http://data.ign.fr/def/geometrie#> .
2 @prefix rgeofla:<http://data.ign.fr/id/geofla/commune/> .
3 @prefix gsp: <http://www.opengis.net/ont/geosparql#> .
4
5 rgeofla:Multipolygon_11130 a geom:MultiPolygon ;
6 geom:crs <http://data.ign.fr/id/ignf/crs/WGS84GDD> ;
7 geom:polygonMember _:polygon_11130_1 ;
8 geom:polygonMember _:polygon_11130_2 .
9
10 _:polygon_11130_1 a geom:Polygon ;
11 geom:crs <http://data.ign.fr/id/ignf/crs/WGS84GDD> ;
12 geom:exterior _:linearRing_11130_1 .
13
14 _:polygon_11130_2 a geom:Polygon ;
15 geom:crs <http://data.ign.fr/id/ignf/crs/WGS84GDD> ;
16 geom:exterior _:linearRing_11130_2 .
17
18
19 _:linearRing_11130_1 a geom:LinearRing ;
20 geom:crs <http://data.ign.fr/id/ignf/crs/WGS84GDD> ;
21 geom:points _:points_11130_1 ;
22 geom:firstAndLast _:point_11130_10 .
23
24
25 _:linearRing_11130_2 a geom:LinearRing ;
26 geom:crs <http://data.ign.fr/id/ignf/crs/WGS84GDD> ;
27 geom:points _:points_11130_2 ;
28 geom:firstAndLast _:point_11130_2 .
29
30
31 _:points_11130_1 a geom:PointsList; geom:firstAndLast _:point_11130_10 ;
32 rdf:rest _:point_11130_11 .
33 _:point_11130_11 a geom:PointsList; rdf:first _:point_11130_12 ;
34 rdf:rest _:point_11130_13 .
35 ....
36 _:point_11130_129 a geom:PointsList; rdf:first _:point_11130_130 ;
37 rdf:rest _:point_11130_1 .
38
39
40 _:point_11130_10 a geom:Point ;
41 geom:crs <http://data.ign.fr/id/ignf/crs/WGS84GDD> ;
42 geom:coordX "7.308745254207776"^^xsd:double ;
43 geom:coordY "43.69237084089203"^^xsd:double .
44
45 rgeofla:06088 a geofla:Commune ;
46 rdfs:label "NICE"@fr ;
47 geom:geometry rgeofla:Multipolygon_11130 ;
48 gsp:asWKT "<http://data.ign.fr/id/ignf/crs/WGS84GDD> MULTIPOLYGON
    (((7.308745254207776 43.69237084089203, 7.306051040744396 43.68445728916297, ....,
    7.308745254207776 43.69237084089203)))"^^gsp:wktLiteral .

```

---

Listing 3.1: Sample of structured geometry of the city of Nice.

### 3.2.5 Limitations of existing tools

Currently, the tools achieving the transformations of geospatial data into RDF still suffer from some limitations. On one hand, they are all compatible with the GeoSPARQL standard, which in turns has a handful of endpoints that implement all the requirements. Moreover, geometries are exposed as literals (e.g., wktLiteral, gmlLiteral) with CRS embedded in the literal. On the other hand, tools based on the NeoGeo vocabulary output structured geometries that can be easily handled by existing Triple Stores and SPARQL queries. However, NeoGeo only allows geometry in WGS84 CRS. The ideal scenario would be a tool that provides output consistent with both GeoSPARQL and structured geometries handling multiple CRSs. Figure 3.1 shows a generic architecture for tools to convert ESRI ShapeFiles into a different flavor of RDF. Tools differ mainly in the presence or not of the CRS reprojection, the vocabularies used for the RDF output, and the compatibility with GML and WKT representations.

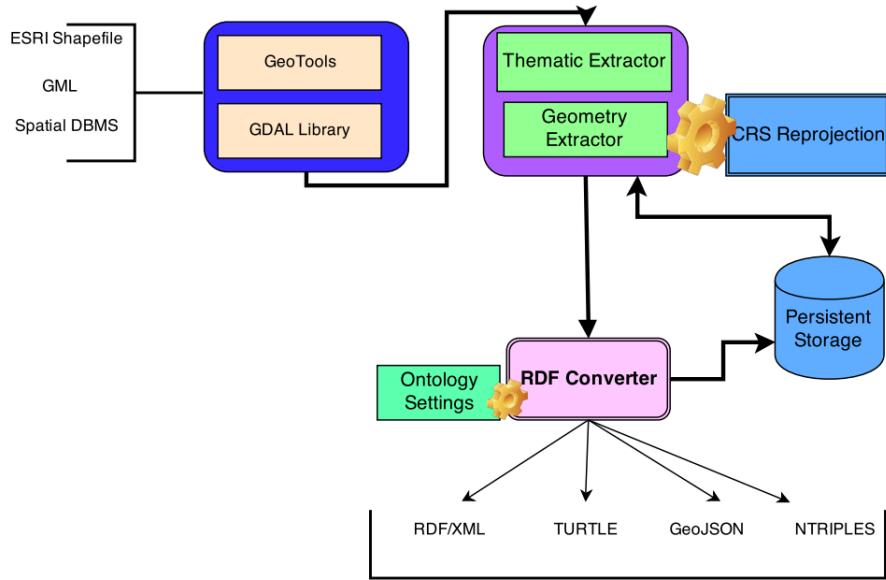


Figure 3.1: Generic architecture of tools for converting raw geospatial data into RDF.

## 3.3 Interlinking Geospatial Vocabulary and Data

We present in this section different criteria that can be used to interlink geospatial datasets in the context of Web of Data, as well as some distance measures adapted for

combining two geospatial datasets. We first provide results of aligning two existing ontologies of previous project GeOnto with five ontologies and two flat taxonomies containing geospatial. For more details on mapping techniques, the readers can refer to [42], and functions implemented in SILK [43, 44] and LIMES [45] interlinking tools. Among the set of properties usually used as data linking criteria, geolocation (addresses, postal code, latitude/longitude, etc) remains one of the most commonly used.

In this thesis, our linking approach is mainly based on geolocation properties comparison, although is by either combining well-known functions implemented in SILK and LIMES. The result of this interlinking process is a list of “*owl:sameAs*” links between entities of each datasets, with a certain score.

### 3.3.1 Criteria for interlinking geospatial Data

Data interlinking task is generally performed by comparing the property values of each resource of a source dataset with corresponding property of the resources described in the target datasets [46]. For geographic databases, data matching is also performed by comparing properties, and especially geometries (points, lines, polygons, etc.) that are used to represent the shape and the location of geographic features. This task is usually based on the distance measures chosen according to the type of the geometric primitives that must be compared [47, 1, 48, 49].

In [1], Olteanu mentions many different forms of interlinking geospatial datasets. We highlight below four important criteria to be taken into account in the process of interlinking two geospatial datasets:

- **Geometrical criteria:** This criterion is based on the geometry of the objects, which is very specific to geographic data. In general, geometry refers to the location and the implicit information regarding the shape (length, orientation, etc.).
- **Topological criteria:** Topology describes the relations of inclusion (e.g., part of, inside, etc) between objects and use the notion of neighborhood. Topological relationships are used for relations such as: the forest borders the road, two roads are connected, etc. Topological relationships are created from the geometry of the initial geographical objects.
- **Attributes criteria:** This criterion uses different attributes belonging to the geographical object, such as name, nature or number of ways. As illustrated in Figure 3.3.1, the attributes of a geographic object can be quantitative or qualitative. The nature of the geographical object is the most important attribute to use in the process of interlinking geospatial data [1]. This motivates the creation of domain ontologies to provide context in the process of interlinking, where domain-specific ontology alignment frameworks such as TaxoMap [50, 51] give good precision for a large topographic ontology.

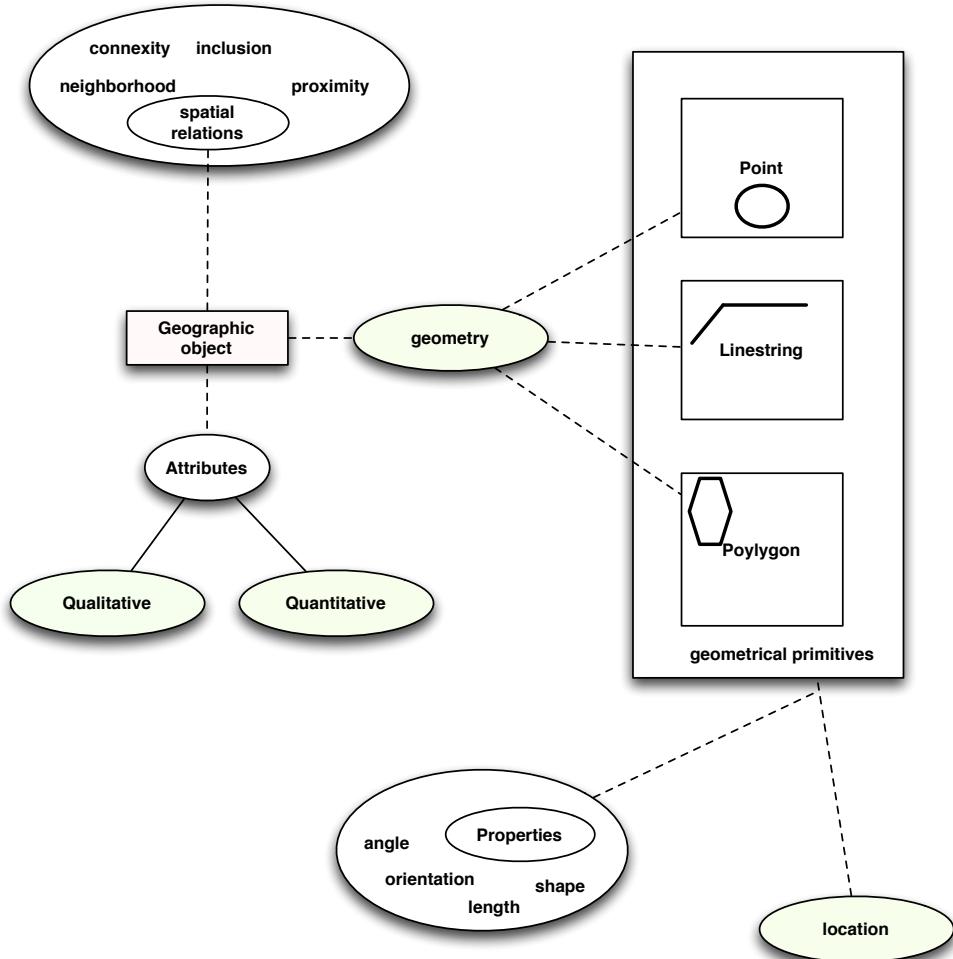


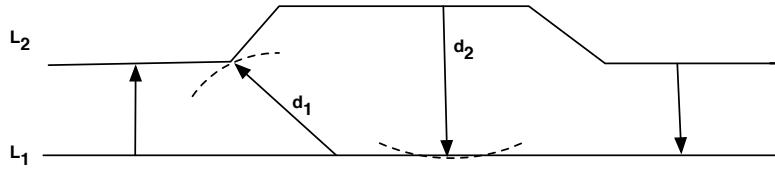
Figure 3.2: Features of a geographic object adapted from [1].

### 3.3.2 Functions for Comparing Geometries

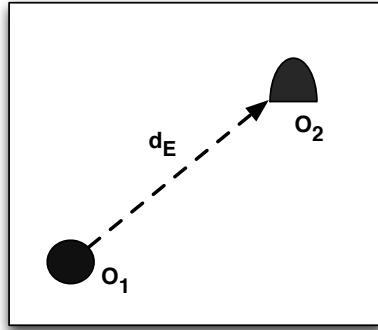
**Euclidean Distance:** For two objects represented with points, the principal distance used for measuring their position is the Euclidean distance (see Figure 3.3.2). Given two objects  $O_1 = (x_1, y_1)$  and  $O_2 = (x_2, y_2)$ , the euclidean distance  $d_E$  between the two objects is defined as follows:

$$d_E = \sqrt{((x_1 - x_2)^2 + ((y_1 - y_2)^2)} \quad (3.1)$$

**Hausdorff distance:** For complex geometries (polylines, linestrings, polygons, etc.) there are different measures in the literature, such as the mean distance, Fréchet distance and the Hausdorff distance [1]. However, at the time of writing this thesis, a Hausdorff distance measure is mostly used in the LIMES framework



(a) Hausdorff distance between two linestring



(b) Euclidean distance between two points

Figure 3.3: The most used distances on geometrical primitives: (a) the Euclidean distance and (b) the Hausdorff distance.

for complex geometries. Considering two LINESTRINGS  $L_1$  and  $L_2$ , the Hausdorff distance represents the maximal distance between the two lines according to the following equation:

$$d_H = \max(d_1, d_2) \quad (3.2)$$

where  $d_1$  and  $d_2$  are defined using the Euclidean distance  $d_E$  as follows:

$$d_1 = \max_{p_1 \in L_1} [\min_{p_2 \in L_2} [d_E(p_1, p_2)]] \quad (3.3)$$

$$d_2 = \max_{p_2 \in L_2} [\min_{p_1 \in L_1} [d_E(p_2, p_1)]] \quad (3.4)$$

Currently on the Web of Data, resources that actually refer to complex topographic features are generally described by very simple geolocation properties, such as a position defined by coordinates (longitude, latitude). On the other hand, geographic reference datasets provide more precise geometric information about geographic features. Interlinking thematic Linked Open Datasets with geographic reference datasets enable to take advantage of both information sources to link independent thematic datasets and create rich cartographic applications for data visualization [52].

### 3.3.3 GeOnto Alignment Process Scenario

IGN-France has previously developed two complementary vocabularies (GeOnto and bdtopo) which differ in their provenance but have the same scope, for describing geographic entities in the French territory. GeOnto is the product of a French research project<sup>8</sup> aiming at building and aligning heterogeneous ontologies in the geographic domain. A “light” version of the final ontology at <http://semantics.eurecom.fr/datalift/tc2012/vocabs/GeoOnto/> defines two high level classes in a total of 783 classes and 17 properties (12 Data Properties and 5 Object Properties). GeOnto contains labels in both French and English, without comments specified for the terms. The bdtopo ontology is derived from a geospatial database with the same name. It contains 237 classes and 51 properties (47 Data Properties / 4 Object Properties). All the labels and comments are in French.

The first step towards interoperability of French geographic features and the existing vocabularies is to align GeOnto to other vocabularies. We choose GeOnto because it covers a large number of categories and also has labels in English. We have performed the alignment with five OWL vocabularies (bdtopo, LGD, DBpedia, Schema.org and Geonames) and two flat taxonomies (Foursquare, Google Place). For the latter, we have transformed the flat list of types and categories into an OWL ontology. For each alignment performed, we only consider `owl:equivalentClass` axioms. We use the Silk tool [48] to compute the alignment using two metrics for string comparison: the *levenshteinDistance* and *jaro* distances. They work on the English labels except for the alignment with bdtopo where we use the French labels. We apply the average aggregation function on these metrics with an empirically derived threshold. However, for generating the final mapping file for vocabularies of small size, we manually validate and insert relations of type `rdfs:subClassOf`. The threshold to validate the results is set to 100% for links considered to be correct and greater than 40% for links to be verified. The alignment with Geonames is special, considering the property restriction used in the ontology for codes.

Table 3.2 summarizes the result of the alignment process between GeOnto and the existing vocabularies and taxonomies. All the resources of this work are available at <http://semantics.eurecom.fr/datalift/tc2012/>.

In general, we obtain good results with Silk, with precision beyond 80%: Google Place: 94%, LGD: 98%, DBpedia: 89%, Foursquare: 92%, Geonames: 87% and bdtopo: 92%. We obtained a precision of only 50% with schema.org due to numerous fine-grained categories that are badly aligned (e.g. `ign:Berge owl:equivalentClass schema:Park`).

## 3.4 Survey on Triple Stores and Workbench

A triple store is a back-end that has some form of persistent storage of RDF data and provides mechanisms to run SPARQL [53] queries against that data. The SPARQL

---

<sup>8</sup><http://geonto.lri.fr/Livrables.html>

Vocabulary	#Classes	#Aligned Classes
LGD	<code>owl:Class:1294</code>	178
DBpedia	<code>owl:Class:366</code>	42
Schema.org	<code>owl:Class:296</code>	52
Geonames	<code>skos:ConceptScheme:12</code> <code>skos:Concept:699</code>	— 287
Foursquare	359	46
Google Place	126	41
bdtopo	<code>owl:Class:237</code>	153

Table 3.2: Results of the alignment process between GeOnto and existing vocabularies/taxonomies.

support can either be built in as part of the main tool, or an add-on installed separately. In this section, we discuss the most used triple stores based on the list maintained by the sparqls tool [54] and we classify them in two group: (1) “generic triple stores” that handle any RDF triples and (2) “geospatial triple stores”, the ones designed to handle geographic data and implement topological functions.

### 3.4.1 Generic Triple Stores

We briefly describes some well-known triple stores that have proven to have passed the “10 Billion Statements”, that is they are able to load and handle more than 10 billion triples.

**Virtuoso:** Virtuoso is a triple store developed by OpenLink<sup>9</sup> and released both in open source and commercial versions. It implements part of SPARQL1.1 and stores billion of triples. The geospatial extension implements subset of SQL spatial (SQL/MM) [55] functions. A geometry stores in Virtuoso uses a special RDF typed literal `virtrdf:Geometry`<sup>10</sup>. Such a literal is automatically indexed in an R tree index [56]. Only geometry in WGS84 can be manipulated. Virtuoso is used as the back end for one of the most used dataset in the LOD cloud, DBPEDIA.<sup>11</sup>

**OWLIM:** OWLIM is a semantic repository developed by Ontotext<sup>12</sup>, which provides support and querying of two-dimensional point geometries modeled with the W3C Geo vocabulary. It implements spatial predicates represented as property functions. The four available operations are: Buffer, Distance, Nearby and Point-in-polygon. OWLIM is implemented as a storage and inference layer for Sesame, with custom spatial index.

<sup>9</sup><http://virtuoso.openlinksw.com/>

<sup>10</sup><http://www.openlinksw.com/schemas/virtrdf#>

<sup>11</sup><http://dbpedia.org/sparql>

<sup>12</sup><http://www.ontotext.com/owlim>

**AllegroGraph:** AllegroGraph<sup>13</sup> is another RDF store developed and maintained by Franz Inc. which stores geospatial data types as native data structures. Support is provided both for Cartesian coordinate systems and for spherical coordinate systems. Every datum in an AllegroGraph store is a universal part identifier (UPI), and for geospatial data, the added UPI type is :geospatial with type code +geospatial+. Geometries are assigned to geometric objects through the use of property <<http://franz.com/ns/allegrograph/3.0/geospatial/pos>>, and for querying a GEO operator is introduced to express geospatial query patterns in SPARQL. AllegroGraph provides some operations on geodata, such as Bounding Box, Distance and Buffer.

### 3.4.2 Geospatial Triple Stores

We discuss in this section two triple stores built for indexing and querying geospatial data, Parliament and Strabon. We acknowledge that there might be some other solutions (e.g., Oracle DBMS version 111g, release 2). More details can be found in [18, 57, 58]. In [58], the authors present a benchmark of Geospatial RDF stores in the context of the TELIOS project<sup>14</sup>, which uses real-world and synthetic data to test the offered functionality and the performance of some prominent geospatial RDF stores.

#### 3.4.2.1 Parliament

The RDF store Parliament<sup>15</sup> developed by BBN Technologies implements most of the functionality of GeoSPARQL [57]. Parliament has been extended to provide the first implementation of the GeoSPARQL standard, therefore geometries are represented using the WKT and GML serializations. Therefore, topological functions belonging in the OGC Simple Feature Access, Egenhofer and RCC8 families are exposed by Parliament [57]. Topological properties and non-topological functions are also available. In addition, multiple coordinate reference systems may be used. Parliament is implemented in C++, and Java Native Interface is used to couple with Jena, and includes a rule engine, serving as a means of inference. It registers the presence of a spatial object in an R-Tree.

#### 3.4.2.2 Strabon

Strabon<sup>16</sup> is a semantic spatio-temporal RDF store for stRDF<sup>17</sup> and stSPARQL [59]. Strabon extends the well-known RDF store Sesame<sup>18</sup>, allowing it to manage both thematic and spatial data expressed in stRDF. PostGIS is used as the relational backend of Strabon. stRDF uses OGC standards (the OGC SFA specification in

---

<sup>13</sup><http://franz.com/agraph/allegrograph/>

<sup>14</sup><http://geographica.di.uoa.gr>

<sup>15</sup><http://parliament.semwebcentral.org>

<sup>16</sup><http://www.strabon.di.uoa.gr>

<sup>17</sup><http://strdf.di.uoa.gr/ontology>

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

particular) for the representation of geospatial data. The datatypes strdf:WKT and strdf:GML are introduced to represent geometries serialized using the OGC standards WKT and GML. Strabon allows geometries to be expressed in any coordinate reference system defined by the EPSG (European Petroleum Survey Group) or OGC. In addition, the latest version of Strabon implements the GeoSPARQL Core, Geometry extension and Geometry topology extension components.

### 3.4.3 Assessing Triple Stores

There are many criteria that can be used to assess triple stores based on the requirements dataset publishers. Most importantly are the compatibility with the set of standards for RDF and SPARQL (e.g., SPARQL1.0, SPARQL1.1). Table 3.3 provides an overview comparing “generic” triple stores. Currently, none of the triple stores listed supports security mechanisms built natively. Regarding specific geospatial triple stores (or extensions of geospatial in RDF stores), more specific requirements are needed to index and provide native functions/operations over geometries. In Table 3.3, we compare triple stores based on the following features:

- the types of geometries supported,
- the coverage of spatial functions,
- the compliance with the GeoSPARQL standard,
- the extensions to existing vocabularies and SPARQL to manage geodata.

Table 3.3: Survey of some generic popular triple stores.

Feature	Virtuoso	OWLIM	AllegroGraph	4Store	Sesame	Fuseki
SPARQL1.0	Yes	Yes	Yes	Yes	Yes	Yes
SPARQL1.1	Partial	Yes	Partial	Partial	Partial	Yes
SPARQL Update	Non-std	Yes	Yes	Yes	Yes	Yes
Reasoning	Rules	Rules	Rules	Add-on	Partial	Rules
10 billion statements	Yes	Yes	Yes	No	No	No
Clustering	Yes	Yes	Yes	Yes	No	No
Open source	Yes	No	No	Yes	Yes	Yes

**Serialization and Triple stores:** We also advocate the use of properties that can provide compatibility with other formats (GML, KML, etc.). This choice can be triple store independent, as there could be ways to use content negotiation to reach the same result. In Table 3.4, `Open Sahara`, `Parliament` and `Virtuoso` are

WKT/GML-compliant with respectively 23 and 13 functions dealing with geodata. Moreover, the choice of the triple store (e.g., Virtuoso<sup>19</sup> vs Open Sahara) is not really an issue, as the IndexingSail<sup>20</sup> service could also be wrapped on-top of Virtuoso to support full OpenGIS Simple Features functions<sup>21</sup>.

Table 3.4: Triple stores survey with respect to geometry types supported and geospatial functions implemented.

Triplestore Geometry supported	WKT-Compliance Geospatial Functions	GML-Compliance Geo-vocabulary
Virtuoso Point	Yes SQL/MM (subset)	Yes W3C Geo + Typed Literal
AllegroGraph Point	- Buffer, Bounding Box, Distance	- “strip” mapping data
OWLIM-SE Point	N/A Distance, Buffer, Nearby, Within	N/A W3C Geo
Open Sahara Point, Line, Polygons	Yes OGC-SFA, Egenhofer, RCC-8	Yes Typed Literal
Parliament Point, Line, Polygons	Yes OGC-SFA, Egenhofer, RCC-8	Yes GeoSPARQL vocabulary
Strabon Point, Line, Polygons	Yes OGC-SFA, Egenhofer, RCC-8	Yes stRDF

### 3.4.4 Workbench for Geospatial Data

The GEOKNOW STACK is a workbench developed within the Geoknow project<sup>22</sup>, which aims to bring geospatial knowledge integration to the Linked Data, with reasoning on billion triples, data provenance, and adaptive authoring, exploration and curation of geospatial data. The GeoKnow Stack consists of eight tools integrated in six modules for Extraction, storage and querying, authoring, linking, enrichment and exploration. Figure 3.4.4 shows the architecture.

Below is a brief description of the main modules:

- **Extraction and Loading:** The module is in charge of loading/importing RDF datasets, extracting and converting legacy datasets using extractors/mappers such as TripleGeo (Cf. Section 3.2) and Sparqlify<sup>23</sup>. Sparqlify is a SPARQL-SQL rewriter that enables one to define RDF views on relational

<sup>19</sup>Here we used Virtuoso Open Edition, V6.xx

<sup>20</sup><https://dev.opensahara.com/projects/useekm/wiki/IndexingSail>

<sup>21</sup><http://www.opengeospatial.org/standards/sfs>

<sup>22</sup><http://geoknow.eu/>

<sup>23</sup><http://sparqlify.org/>

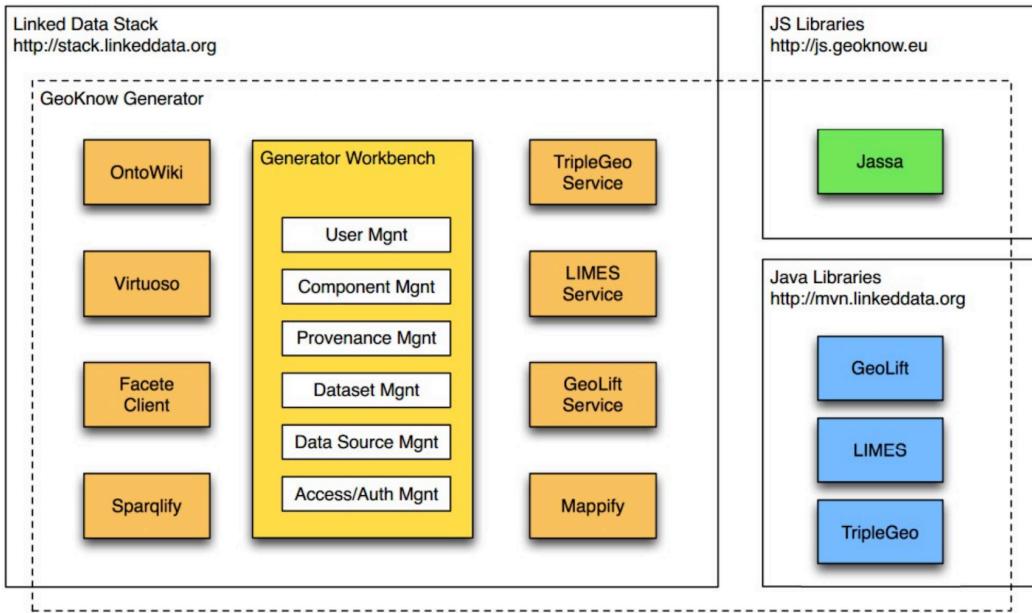


Figure 3.4: Architecture of the Geoknow Stack.

databases and query them with SPARQL. The module powers the Linked-Data Interface of the LinkedGeoData Server and provides access to billions of virtual triples from the OpenStreetMap database.

- **Storage and Querying:** The module is powered by Virtuoso 7.0 triple store for querying the datasets.
- **Authoring module** integrated the OntoWiki tool [60] that facilitates the visual presentation of a knowledge base as an information map, with different views on instance data. It enables intuitive authoring of semantic content, with an inline “*What You See is What You Get*” (WYSIWYG) editing mode for editing RDF content.
- **Linking and Fusion:** To achieve this module, LIMES is the tool integrated in the workbench. LIMES is an abbreviation of the LInk discovery framework for MEtric Spaces, a tool for interlinking resources on the Web of Data<sup>24</sup>. It implements time-efficient approaches for large-scale link discovery based on the characteristics of metric spaces [45]. It is easily configurable via a web interface or can also be downloaded as standalone tool for carrying out link discovery locally. LIMES can reduce the number of comparisons needed during the mapping process by several orders of magnitude [61]. The approaches implemented in LIMES include the original LIMES algorithm for edit distances, REEDED for weighted edit distances, HR3, HYPPO, and ORCHID.

<sup>24</sup><http://aksw.org/Projects/LIMES.html>

The algorithms implemented include the supervised, active and unsupervised versions of EAGLE, COALA and EUCLID [62].

- **Enrichment:** During the enrichment step, the tool GEOLIFT [63] is used to enrich geographic content using three techniques: linking, dereferencing and Natural Language Processing (NLP). The code is available at <https://github.com/AKSW/GeoLift/>
- **Exploration:** The workflow proposes two tools for this module: MAPPIFY and FACETE. The former is a tool for exploring (geographical) Linked Data datasets on the Web, while the latter allows to explore the specific slice of data named “facet” of a Linked Data endpoint in a graphical way, by defining set of constraints on properties of the database. Once the facet is defined, the information in the facet can be clicked-through in a tabular interface and visualized on a map. FACETE is available at <http://144.76.166.111/facete/>.

## 3.5 Datalift: A tool for Managing Linked (Geo)Data Publishing Workflow

In this section, we present the Datalift platform, as a tool to help for lifting raw data to RDF which integrate some existing tools to manage the workflow of publishing geodata on the Web. To the best of our knowledge, a related framework providing similar functionalities is the GEOKNOW STACK that we have described in section 3.4.4.

Datalift is an open source platform [64] helping to lift raw data sources or legacy data to semantic interlinked data sources. The ambition of DataLift is to act as a catalyst for the emergence of the Web of Data by providing a complete path from raw data to fully interlinked, identified, and qualified linked datasets. The Datalift platform supports the following stages in lifting the data:

1. Selection of ontologies for publishing data;
2. Conversion of data to the appropriate format (e.g., from CSV to RDF);
3. Interlinking of data with other data sources;
4. Publication of linked data ;
5. Access control and license management.

Figure 3.5 gives an overview of the different steps in lifting raw source data into RDF using different modules of Datalift.

### 3.5.1 Functionalities of the Datalift platform

The architecture of Datalift is modular. Several levels of abstraction allow decoupling between the different stages from raw data to semantic data. The dataset

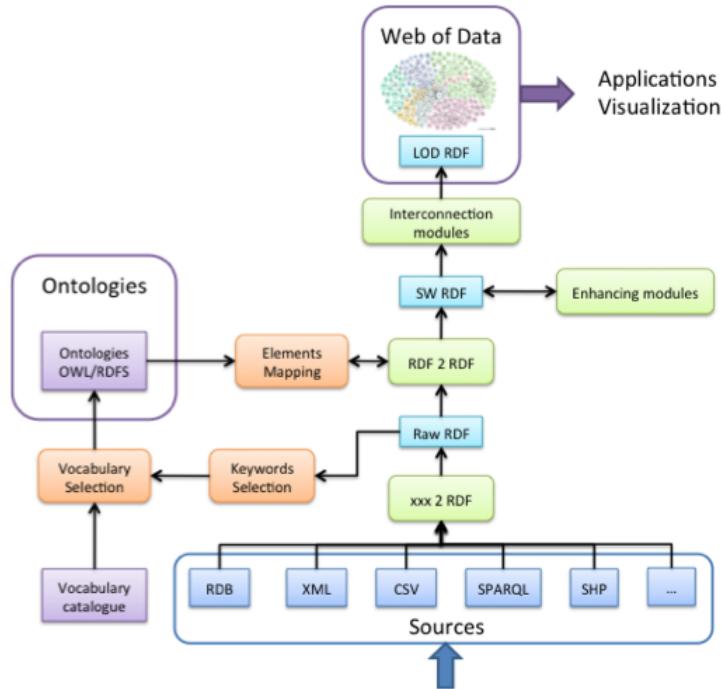


Figure 3.5: Lifting process of raw data source into RDF using Datalift Platform

selection allows us to identify the data to be published and migrate them to a first RDF version. The ontologies selection step asks the user to input a set of vocabularies', terms that will be used to describe the lifted data. Once the terms are selected, they can be mapped to the raw RDF and then converted to properly formatted RDF. The data is then published on the DataLift SPARQL endpoint. Finally, the process aims at providing links from the newly published data to other datasets already published as Linked Data on the Web. Figure 3.5 corresponds to a visual depiction of the workflow to convert raw data into “*structured*” RDF data. Figure 3.5.1 depicts the architecture of Datalift, consisting of different modules for:

1. **Dataset Selection** The first step of the data lifting process is to identify and access the datasets to be processed. A dataset is either a file or the result of a query to retrieve data from a datastore. The kinds of files currently considered are CSV, RDF, XML, GML and Shape files. Queries are SQL queries sent to an RDBMS or SPARQL queries on a triple store.
2. **Ontologies Selection:** The publisher of a dataset should be able to select the vocabularies that are the most suitable to describe the data, and the least possible terms should be created specifically for a dataset publication task. The Linked Open Vocabularies [65] (LOV) developed in Datalift provides easy access methods to this ecosystem of vocabularies, and in particular by making explicit the ways they link to each other and providing metrics on how they are used in the linked data cloud. LOV is integrated as a module in the DataLift

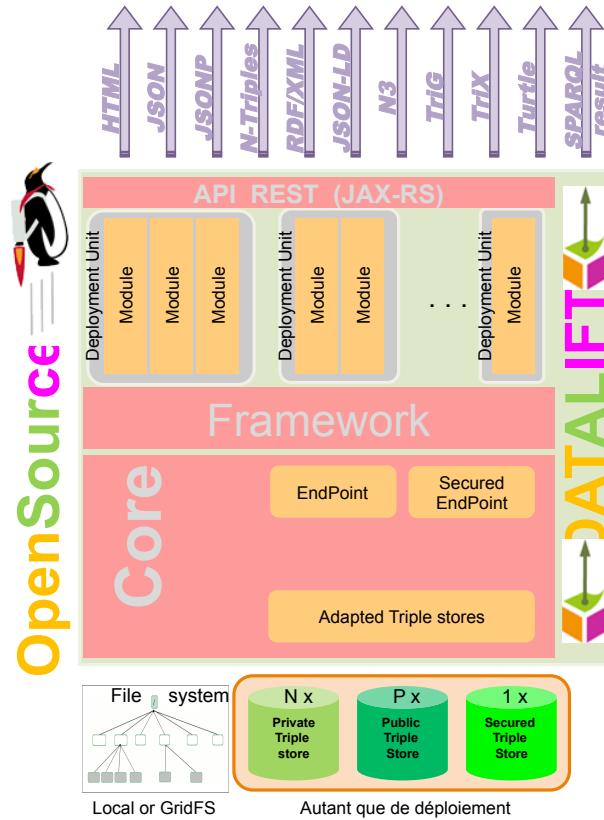


Figure 3.6: Architecture of DataLIFT platform.

platform to assist the ontology selection.

**3. Data Conversion:** Once URIs are created and a set of vocabulary terms able to represent the data is selected, it is time to convert the source dataset into a more precise RDF representation. Many tools exist to convert various structured data sources to RDF. The major source of structured data on the Web comes from spreadsheets, relational databases and XML files. Two steps are provided. First, a conversion from the source format to raw RDF is performed. Second, a conversion of the raw RDF into “well-formed” RDF using selected vocabularies is performed using SPARQL Construct queries. Most tools provide spreadsheet conversion to CSV, and CSV to RDF is straightforward, each line becoming a resource, and columns becoming RDF properties. The W3C RDB2RDF WG<sup>25</sup> proposes the Direct Mapping to automatically generate RDF from the tables but without using any vocabulary, and R2RML<sup>26</sup> to assign vocabulary terms to the database schema. In the case of XML, a generic XSLT transformation is performed to produce RDF from a wide range of XML documents. The DataLIFT platform provides a graphical interface to

<sup>25</sup><http://www.w3.org/2001/sw/rdb2rdf/>

<sup>26</sup><http://www.w3.org/TR/r2rml/>

help mapping the data to selected vocabulary terms.

4. **Data Protection:** The data protection is linked to a Java security framework Apache Shiro<sup>27</sup> for obtaining the information, i.e., username and password, about the user who is accessing the platform. The module<sup>28</sup> checks which data is targeted by the user's query and then verifies whether the user can access the requested data. This verification leads to three kinds of possible answers depending on the access privileges of the user: some of the requested data is returned, all the requested data is returned, or no data is returned. This means that the user's query is filtered in such a way that she is allowed to access only the data she is granted access to. The access policies are expressed using RDF and SPARQL 1.1 [53] Semantic Web languages thus provide a completely standard way of expressing and enforcing access control rules.
5. **Data Interlinking:** The interlinking step provides means to link datasets published through the Datalift platform with other datasets available on the Web of Data. Technically, the module helps to find equivalence links in the form of "owl:sameAs" relations. An analysis of the vocabulary terms used by the published data set and a potential data set to be interlinked is performed. When the vocabulary terms are different, the module checks if alignments between the terms used by the two data sets are available. Here the alignment server provided with the Alignment API<sup>29</sup> is used for that purpose. The correspondences are translated into SPARQL graph patterns and transformation functions are combined into a SILK script.
6. **Data Publication:** This module aims at publishing the data obtained from the previous steps to a triple store, either public or private. The providers can restrict which graphs can be accessible, they could decide whether to provide just a "Linked Data" or a "Linked Open Data". Datalift comes by default with Sesame triple store , but provides API for connecting to Allegrograph, OWLIM, and Virtuoso triple stores as well.

**Installation:** All the documentation for installing Datalift is available at [http://datalift.org/wiki/index.php/Platform\\_installation\\_\(english\)](http://datalift.org/wiki/index.php/Platform_installation_(english)). The latest version of the platform is available at <http://datalift.org/en/node/24>, for downloading and installation.

**Usage:** The data lifting workflow has several distinct steps. DataLift makes it possible to replay each step in producing different results for each step. To facilitate access to all the different treatments and their results, they are grouped as one project. The project gathers together the various sources used and the results of all treatments done. Each module has its own way to be used within the lifting process

---

<sup>27</sup><http://shiro.apache.org/>

<sup>28</sup><http://wimmics.inria.fr/projects/shi3ld/>

<sup>29</sup><http://alignapi.gforge.inria.fr/>

in DataLift. For more details, the readers are encouraged to read this resource at [http://datalift.org/wiki/index.php/How\\_to\\_use\\_the\\_Datalift\\_platform\\_to\\_publish\\_a\\_dataset\\_on\\_the\\_Web#The\\_lifting\\_project](http://datalift.org/wiki/index.php/How_to_use_the_Datalift_platform_to_publish_a_dataset_on_the_Web#The_lifting_project).

**Comparison between Datalift and Geoknow Stack** As described above, the Geoknow Stack shares similar goals and functionalities with Datalift. However, while the latter is cross-domain and generic, the former is more targeted for geospatial data. Moreover, Datalift implement connectors for many triple stores, while Geoknow Stack is powered by Virtuoso. Datalift is available in different platforms (Linux, Windows and Mac OS), while the current version of Geoknow can be installed only in Linux. Table 3.5 gives more features used to compare both frameworks, such as modules, platform for installation, access control and authoring modules.

Features	Geoknow Stack	Datalift
Scope	Geospatial data	Cross-domain/Generic
Triple Store	Virtuoso 7	Sesame, Virtuoso 6, AllegroGraph
Shape2RDF	Few models for geometry	More generic approach
Interlinking Tool	LIMES	SILK
Publication	Publish directly in a graph	Export dump data in CSV, Turtle, NTriples
Access Control	N/A	Security Access Module
Installation	Expert level	One-click
Platform	Linux	Multi-platform
Deployment Environment	N/A	IGN, INSEE
Provenance	Partially (user)	Tracking PROV for each project
Visualization	Facete, Mappify	Sgvizler, RDFViz
Access Control	N/A	Security Access Module
Authoring	Wiki integrated	N/A

Table 3.5: Comparison of Datalift with the GeoKnow Stack

## 3.6 Publishing French Authoritative Datasets

In this section, we provide our contribution on the process of publishing different datasets from the French National Mapping agency and OSM-France. We detailed the workflow starting from the conversion of source files, their modeling using the vocabularies implemented, the interconnection with relevant data sources. We conclude this section by showing the current status of the French LOD cloud according

to our contributions so far.

### 3.6.1 Publishing French Administrative Units (GEOFLA)

As a dataset dedicated to administrative units, GEOFLA is very likely to be reused by other datasets, either by reusing directly its URIs for georeferencing needs, or by reusing its description of administrative units-labels, properties and geometries- for interlinking purposes.

#### 3.6.1.1 Data conversion

GEOFLA is delivered as a set of 4 shapefiles that describe the boundaries and properties of administrative units of mainland France (for CRS reasons, overseas territories are delivered within different shapefiles): communes, cantons, arrondissements and départements. For the sake of our application, we have generated another shapefile describing regions by aggregating the geometries of the instances of departments based on their region's foreign key value. GEOFLA is updated every year. Publishing this data in RDF with unique identifiers on the Web will ease the interlinking with some existing datasets describing French boundaries in the wild. We follow a two-step conversion: we use the SHP2RDF module of Datalift to obtain a raw RDF from shapefiles, and the RDF2RDF module of Datalift using a set of SPARQL construct queries<sup>30</sup> for getting a refined RDF datasets using suitable vocabularies.

#### 3.6.1.2 URI design policy

One of the requirements to publish data is to have unique identifiers and stable URIs<sup>31</sup>. Since our legacy databases have unique IDs to refer to the objects, we had to make sure they were unique at Web level. Thus, the base scheme for vocabularies URIs is: <http://data.ign.fr/def/>. Besides, the base schema for identifying a real world resource uses <http://{BASE}/id/>. For example, IGN main buildings are located in the commune with the URI `rgeofla:commune/94067`, corresponding to Saint-Mandé, and `rgeofla:departement/94` corresponding to the department “Val de Marne” to which the commune belongs. We had to make choices based on the set of best practices related to URI design<sup>32</sup> which should guarantee both stable and human readable identifiers.

#### 3.6.1.3 Interlinking with existing GeoData

We interlinked our datasets with the NUTS, DBpedia FR<sup>33</sup> and GADM datasets. SILK [43] is used to interlink the departments in our dataset with departments in DBpedia FR, using labels and INSEE Code. We obtained 93 matches (all correct)

<sup>30</sup> <https://github.com/gatemezing/ign-iswc2014/tree/master/rdf2rdf>

<sup>31</sup> <http://www.w3.org/TR/ld-bp/>

<sup>32</sup> <http://www.w3.org/TR/ld-bp/#HTTP-URIS>

<sup>33</sup> <http://fr.dbpedia.org/>

while three resources are missing for the departments 07, 09 and 75<sup>34</sup>. The LIMES tool<sup>35</sup> is then used to perform the rest of the interlinking tasks [62] with the trigrams function based on the labels with restriction to France.

- GEOFLA-RDF with DBpedia FR: **23,252** links obtained. This results show the missing of nearly 13,435 communes not correctly typed in DBpedia FR as `Spatial Feature` or `Place`, or not having a French Wikipedia entry.
- GEOFLA-RDF with GADM (8,314,443 features): **70** links obtained: 10 communes, 51 departments and 9 regions. The property `gadm:in_country` is used to restrict the interlinking to France. E.g.: The city of Saint-Alban in Quebec is a commune in France.
- GEOFLA-RDF with NUTS (316,236 triples): Using a “naive” script with `trigrams` function on `geofla:Commune/rdfs:label` and `spatial:Feature/ramon:name` reveal two odd results located in Germany and Switzerland. The former named “*Celle*” and the latter being the *JURA*. In order to remove those odd effects, we add another restrictions based on `ramon:code` by filtering the ones located in France (136 features). The final matchings give a total of **105** correct links: 14 communes, 75 departments and 16 regions.

The results in Table 3.6 show good precision of the matching algorithm (score above 0.98) and a rather low recall value with DBpedia-FR (0.627). The few number of matched entities is likely due to the low coverage of French features in the target and external datasets used for the interlinking.

Datasets	Precision	Recall	F1-score
NUTS	0.98	1	0.90
GADM	1	0.86	0.92
DBpedia-FR	1	0.627	0.77

Table 3.6: Evaluation results in the interlinking process.

The SPARQL endpoint for the French Administrative dataset is available for querying at <http://data.ign.fr/id/sparql>.

### 3.6.2 Publishing French Gazetteer

In this section, we present some first tests of converting BDTOPO® into RDF and interlinking with LinkedGeoData using LIMES. The results confirm the need for geographic publishers to publish georeference data on the Web.

<sup>34</sup><https://github.com/gatemezing/ign-iswc2014/tree/master/interlinking/matched>

<sup>35</sup><https://github.com/AKSW/LIMES>.

**Data conversion, URIs and Interlinking:** Shapefiles are converted into RDF using the same two conversion process as for GEOFLA®. The URIs for each resource follow the pattern: `rtopo:CLASS/ID` for the feature, while `rtopo:geom/CLASS/ID` is used to reference the geometry of the resource. The gazetteer dataset in RDF is part of the BD TOPO® database consisting of 1,137,543 triples (103,413 features). We chose LinkedGeoData (LGD)<sup>36</sup> to perform the alignments using the main class `lgdo:Amenity`<sup>37</sup> (5,543,001 triples), as they are closed to the features contained in the gazetteer. We perform the interlinking on the geometries using the `hausdorff` metric of LIMES tool. A total of **654** alignments was obtained above the threshold (0.9). This relatively low number of hits can be explained by the coverage of French data in LGD, and the subset of BDTOPO® used for the interlinking. Table 3.7 provides details of the alignments with subclasses of Amenity.

LGD Class	#links matched
<code>lgdo:Shop</code>	252
<code>lgdo:TourismThing</code>	30
<code>lgdo:Craft</code>	3
<code>lgdo:AerowayThing</code>	37
<code>lgdo:AerialwayThing</code>	11
<code>lgdo:EmergencyThing</code>	56
<code>lgdo:HistoricThing</code>	257
<code>lgdo:MilitaryThing</code>	8

Table 3.7: Interlinking results using the Hausdorff metric of LIMES tool between LinkedGeoData and toponyms in the French Gazetteer

### 3.6.3 Publishing Addresses of OSM-France in RDF

OpenStreetMap France is working on providing the location addresses of France in different formats: CSV, ShapeFiles in a collaborative and open source fashion. The BANO project already contains 15 million indirect georeferencing locations. The geometries are POINTs and use WGS84 CRS. One of the requirements is to provide an RDF-ize version of the data, enriching the dataset with existing relevant ones. The goal of transforming BANO to RDF (BANO2RDF) consists of six basic requirements:

- ◆ **Requirement 1-Model:** Model the existing data according to an existing vocabulary, generic enough to cover the scope of the existing dataset.
- ◆ **Requirement 2-Provenance:** Add relevance metadata on the published dataset such as provenance, spatial coverage, licensing, authorship, etc.

<sup>36</sup><http://linkedgeodata.org/sparql>

<sup>37</sup><http://linkedgeodata.org/ontology/>

- ◆ **Requirement 3-Stable URIs:** Define a policy to provide and ensure stable URIs for identifying uniquely each address entity on the Web.
- ◆ **Requirement 4-Interconnection:** Find relevant dataset already published on the Web to which interconnect for better discoverability.
- ◆ **Requirement 5-Access to data:** Provide primarily a frequent dump of the dataset in RDF.
- ◆ **Requirement 6-GeoSPARQL interoperability:** Provide different representations of geometries of locations that are interoperable with GeoSPARQL standard.

Based on the above requirements, the Location Core Vocabulary [36] is used in an *ad-hoc* script<sup>38</sup> to convert CSV data into RDF. URIs used to identify objects are of the form: <http://id.osmfr.org/bano/INSEE-CODE+FANTOIR-CODE+Street-Number>. For example, the location corresponding to EURECOM in SophiaTech (*450, route des chappes, 06410 Biot, FRANCE*) is identified by the URI <http://id.osmfr.org/bano/060180238L-450>. Moreover, the property `locn:location` links to the French Statistic dataset<sup>39</sup> for communes in France. Metadata are inserted at the beginning of the dataset in RDF corresponding to a department, using vocabularies to model license, spatial coverage of the data (i.e., `dcat`, `dcterms`, `foaf`). The extract below represents the metadata for the department of “ALPES-MARITIMES” to which the commune of Biot belongs:

```

1 <http://www.openstreetmap.fr/bano/data/> a dcat:Catalog ;
2   dcterms:title "Donnees des adresses du projet BANO en RDF"@fr ;
3   dcterms:description "Le projet BANO en RDF de Base d'Adresses
4   Nationale
5   Ouverte par OpenStreetMap France."@fr ;
6   foaf:homepage <http://openstreetmap.fr/bano> ;
7   dcterms:language "fr" ;
8   dcterms:license <http://www.opendatacommons.org/licenses/odbl/>
9   ;
10  dcterms:publisher <http://www.openstreetmap.fr/> ;
11  dcterms:issued "2014-05-14"^^xsd:date ; # dataset issued
12  dcterms:modified "2014-08-21"^^xsd:date ; #last modification
13  dcterms:spatial <http://id.insee.fr/geo/departement/06>,
14  <http://id.insee.fr/geo/pays/france> .

```

Listing 3.2: The metadata information used in the BANO2RDF dataset.

Geometries are provided in three different representations: W3C WGS84, typed literal in WKT and geo URI; all using the property `loc:geometry`. Currently, the

<sup>38</sup><https://github.com/osm-fr/bano/blob/master/out/csv2ttl.py>

<sup>39</sup><http://id.insee.fr>

dataset is mainly available as dump files at <http://www.openstreetmap.fr/bano/>. However, an experimental endpoint has been set up to query at <http://eventmedia.eurecom.fr/sparql> with the named graph <http://data.osm.fr/bano/>, consisting of nearly 170 million triples.

**Interconnection with LinkedGeoData:** The first results of mappings with LinkedGeoData amenities using the LIMES tool against three big cities of France (Paris, Lyon and Marseille). The dataset is loaded in an endpoint, and the postal code is used to filter the relevant location address. Furthermore, the HAUSDORFF function is used on the geometry (polygons) of the target and source datasets, with a threshold of 0.97, i.e, we consider “*very close*” resources with a distance of 30 meter, based on the following equation:

$$\theta = \frac{1}{1+d}$$

$$d = \frac{1-\theta}{\theta}$$

Where  $\theta = 0.97$  (threshold) and  $d$ = distance in Kilometer (km).

Listing 3.3 shows the configuration used for finding corresponding buildings of LinkedGeoData in the Bano2RDF dataset.

```

1 <LIMES>
2
3 <PREFIX>
4   <NAMESPACE>http://geovocab.org/geometry#</NAMESPACE>
5   <LABEL>geom</LABEL>
6 </PREFIX>
7
8 <PREFIX>
9   <NAMESPACE>http://www.w3.org/ns/locn#</NAMESPACE>
10  <LABEL>locn</LABEL>
11 </PREFIX>
12
13 <PREFIX>
14   <NAMESPACE>http://www.opengis.net/ont/geosparql#</NAMESPACE>
15   <LABEL>geos</LABEL>
16 </PREFIX>
17
18 <PREFIX>
19   <NAMESPACE>http://linkedgeodata.org/ontology/</NAMESPACE>
20   <LABEL>lgdo</LABEL>
21 </PREFIX>
22 <PREFIX>
23   <NAMESPACE>http://www.geonames.org/ontology#</NAMESPACE>
24   <LABEL>gn</LABEL>
25 </PREFIX>
26 <SOURCE>
27   <ID>bano2RDF</ID>
28   <ENDPOINT>http://eventmedia.eurecom.fr/sparql</ENDPOINT>
29   <VAR>?x</VAR>
30   <PAGESIZE>5000</PAGESIZE>
31   <RESTRICTION>?x a locn:Address</RESTRICTION>
32   <RESTRICTION>?x locn:postalCode ?code</RESTRICTION>
```

```

33   <RESTRICTION>FILTER(regex(str(?code), '690'))</RESTRICTION>
34   <PROPERTY>locn:geometry/geos:asWKT RENAME polygon</PROPERTY>
35   <TYPE>SPARQL</TYPE>
36   </SOURCE>
37   <TARGET>
38   <ID>linkedgeodata</ID>
39   <ENDPOINT>http://linkedgeodata.org/sparql</ENDPOINT>
40   <VAR>?y</VAR>
41   <PAGESIZE>5000</PAGESIZE>
42   <RESTRICTION>?y a lgdo:Building</RESTRICTION>
43   <PROPERTY>geom:geometry/geos:asWKT RENAME polygon</PROPERTY>
44   <TYPE>SPARQL</TYPE>
45   </TARGET>
46
47   <METRIC>hausdorff(x.polygon, y.polygon)</METRIC>
48   <ACCEPTANCE>
49   <THRESHOLD>0.97</THRESHOLD>
50   <FILE>bano690xx-lgdBuilding_verynear.nt</FILE>
51   <RELATION>lgdo:near</RELATION>
52   </ACCEPTANCE>
53   <REVIEW>
54   <THRESHOLD>0.95</THRESHOLD>
55   <FILE>bano690xx-lgdBuilding_place_near.nt</FILE>
56   <RELATION>lgdo:near</RELATION>
57   </REVIEW>
58
59   <EXECUTION>Simple</EXECUTION>
60   <OUTPUT>TAB</OUTPUT>
61   </LIMES>

```

Listing 3.3: This LIMES script is used to interconnect the Bano2RDF dataset located in the region of Marseille with buildings in LinkedGeoData.

Table 3.8 shows the first results of the mappings. For each type of amenities in LinkedGeoData, it shows the number of matched resources, and the percentage of the result with respect to the distinct individuals of the category. For example, the mapping reveals 735 parkings located in Paris (postal code starting with 750). Since there are 250,516 resources of type “Parking” in LinkedGeoData, the percentage is 0.293%, i.e.,  $(735/250680)*100$ . The table also reveals the low presence of the addresses for French territory in LinkedGeoData. All the results of the mappings are available at <https://github.com/gatemezing/bano2rdf-matching>. Based on the results of the interlinking process shown in Table 3.8, shops and restaurants are the resources with most `owl:sameAs` links between the Bano2RDF and LGD datasets.

### 3.6.4 Status of French LOD cloud (FrLOD)

We summarize in this section our contributions for a *French LOD (FrLOD)* cloud, consisting of the different datasets published in RDF based on best practices for publishing data in 4-5 stars<sup>40</sup> on the Web. Datasets published belong to different domains: geographical, governmental, statistical and educational. A VOID [66] de-

---

<sup>40</sup><http://5stardata.info/>

LGData Amenities	Bano-750xx (248,052)		Bano-130xx (401,404)		Bano-690xx (89,061)	
	#matched	%	#matched	%	#matched	%
Building (22,283)	05	0.022	12	0.053	0	0
Parking (250,516)	735	0.293	625	0.24	210	0.083
Shop (778,680)	21,171	2.71	8,556	1.098	3,049	0.391
School (318,287)	883	0.277	411	0.129	197	0.061
PlaceOfWorship (357,445)	272	0.076	193	0.053	31	0.008
Restaurant (260,675)	13,567	5.204	2,654	1.018	1,882	0.721
PublicBuilding (26,735)	97	0.362	64	0.239	21	0.078
PostOffice (87,731)	971	1.106	555	0.632	173	0.197

Table 3.8: Initial mappings of Bano2RDF with LGD amenities resources respectively in Paris, Marseille and Lyon. The links are obtained using LIMES tool with a threshold of .97 using the Hausdorff distance.

scription of the FrLOD gives details of the different datasets published<sup>41</sup>, as well as links and applications/visualizations built on top of them. Figure 3.7 depicts a static diagram of the FrLOD, while in Table 3.9 gives an overview of the datasets, number of triples and domains. Altogether, the FrLOD represents 340 million RDF triples, which is nearly 10% of the DBpedia 2014 release<sup>42</sup>. It is expected to add the FrLOD datasets in the Linking Open Data Cloud at <http://datahub.io/group/lodcloud>

### 3.7 Evaluation of Spatial Queries

We illustrate in this section different queries making use of geospatial data and geometries functions implemented in three triple stores (Virtuoso, OWLIM, Sesame). Within the queries, we use the following points in WGS84, in the the form (longitude latitude): Eurecom, located at (7.0463, 43.6266) and Eiffel Tower, located at (2.2942 48.8628). We perform the time for retrieving the results as well as the expressivity of the types of SPARQL quieries available to users to get the results.

<sup>41</sup>An interactive version can be accessed at <http://www.eurecom.fr/~atemezin/work/frenchLOD.svg>

<sup>42</sup><http://wiki.dbpedia.org/Datasets>

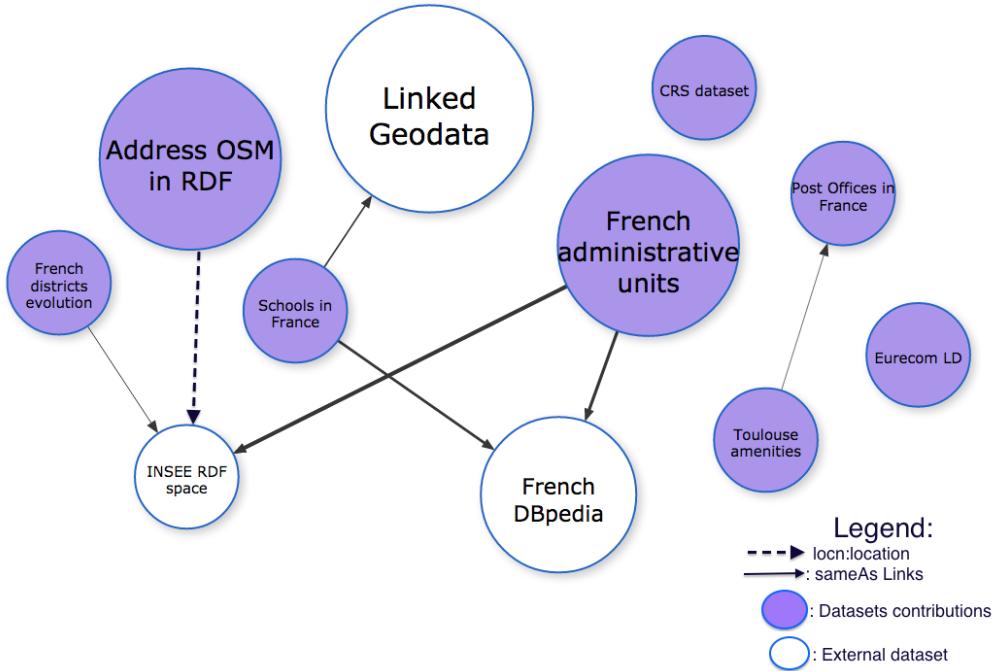


Figure 3.7: French LOD cloud diagram based on the different datasets published in 4-5 stars.

Dataset	#Triples	Domain
French Admin Unit	173,639,128	Geography
Schools in France	1,454,564	Statistics
Eurecom LD	34,651	Education
CRS dataset	23, 377	Government
Toulouse Amenities	963, 618	Government
Bano Project	163, 723, 230	Geography
French Post Offices	371, 381	Government
French District Evolution	150, 356	Geography

Table 3.9: Overview of the content of our contribution to the French LOD Cloud.

### 3.7.1 Querying LinkedGeodata

The query below finds public buildings 10 km around Eurecom with their corresponding distance from the LinkedGeodata endpoint. In this query, we use the functions of `st_distance`, `st_intersects` and `st_point`. Time used to answer this query is **341 ms**, monitored by using the time command on `CURL -g <query>`. The result is listed in Table 3.10.

```

1 Prefix lgd:<http://linkedgeodata.org/>
2 Prefix lgdo:<http://linkedgeodata.org/ontology/>
3
4 SELECT ?s ?name (bif:st_distance(bif:st_point(?long, ?lat), bif:st_point(7.0463, 43.6266)) as ?distance)
5 WHERE
6 {
7     ?s a lgdo:PublicBuilding ; rdfs:label ?name .
8     ?s geo:lat ?lat ; geo:long ?long .
9     filter (bif:st_intersects(bif:st_point(?long, ?lat), bif:st_point(7.0463, 43.6266), 10 ))
10 } ORDER BY ASC (?distance)

```

Listing 3.4: Query on LinkedGeodata endpoint to find all public buildings 10 km around Eurecom building in SophiaTech.

URI	Label	Distance
lgdr:node1645286605	Fondation Sophia Antipolis	1.43949
lgdr:node1274078114	Caisse d'allocations familiales	4.12131
lgdr:node1082377863	Caisse des écoles	5.06445
lgdr:node1082377868	Centre Administratif	5.14575
lgdr:node960122458	Trésorerie de Mougins	5.54149
lgdr:node471974902	Centre d'information jeunesse	7.07093
lgdr:node2209129663	Not available	7.59737

Table 3.10: Results of the public buildings 10 km around EURECOM from Linked-GeoData endpoint. `lgdr` represents the base URI for <http://linkedgeodata.org/triplify/>

The result time is faster when using built-in geospatial functions for Virtuoso. Additionally, it might take much more time (x100) and a resulting long query if just using standard SPARQL functions without geographic functions.

### 3.7.2 Querying FactForge (OWLIM)

The query below find in FactForge<sup>43</sup> the French departments within the Bounding Box of Eurecom and the Eiffel Tower. We use the function `omgeo:within`. The results<sup>44</sup> consist of 36 departments and the query takes **1.369 s** to be completed.

```

PREFIX omgeo: <http://www.ontotext.com/owlim/geo#>
PREFIX geo-ont: <http://www.geonames.org/ontology#>
SELECT DISTINCT ?link ?m
WHERE {
    ?link geo-ont:name ?m.
    ?link geo-ont:featureCode geo-ont:A.ADM2 .
    ?link geo-ont:parentCountry dbpedia:France .

```

<sup>43</sup><http://factforge.net/sparql>

<sup>44</sup><http://goo.gl/dLc45p>

```
?link geo-pos:lat ?lat2 .
?link geo-pos:long ?long2 .
?link omgeo:within( 43.6266 2.2942 48.8628 7.0463 ) }.
```

### 3.7.3 Querying Structured Geometries

We query the `http://data.ign.fr` endpoint that contains structured geometries of French departments, powered by Datalift, using the default triple store Sesame. The query finds all the departments containing in the Bounding Box formed by EURE-COM and Eiffel Tower. The results consisting of 94 departments is obtained after **23.496s** when launched against `data.ign.fr/id/sparql1`. Not all the SPARQL1.1 property paths<sup>45</sup> are fully implemented in Virtuoso endpoints.

```
SELECT DISTINCT ?name WHERE {
?dep a geofla:Departement .
?dep rdfs:label ?name .
?dep geom:geometry/geom:center ?ct .
?ct geom:coordX ?long ; geom:coordY ?lat .
?dep geom:geometry/geom:polygonMember/geom:exterior/geom:points ?pl .
?pl rdf:type geom:PointsList .
?pl (rdf:rest*/rdf:first)|geom:firstAndLast ?pm .
?pm rdf:type geom:Point .
?pm geom:coordX ?x .
?pm geom:coordY ?y .
FILTER ((?x > 2.2942 && ?x < 7.0463) && (?y > 43.6266 && ?y < 48.8628)) .
}
```

The result in section 3.7.3 suggests that it may be possible to reduce the time of querying geographic data by just using standard SPARQL queries without built-in geospatial functions.

## 3.8 Summary

In this chapter, we have surveyed tools for extracting and converting geospatial data in RDF. Then, we described GeomRDF, a tool developed within the Datalift project which goes beyond the state of the art by providing structured geometries and conforming to GeoSPARQL. Moreover, we have described the limitation of the existing data models by discussing some recommendations to publishers of geodata on storage aspects. Similarly, an extensive description of the Datalift tool used to publish data on the Web has been provided, with a special focus of our contribution to build the French LOD cloud with datasets in 4-5 stars according to Linked Data

---

<sup>45</sup><http://www.w3.org/TR/sparql11-query/#propertypaths>

principles. Finally, we have shown some real world use-cases of SPARQL queries making use of either structured geometries or built-in geospatial functions. Depending on the requirements of the users and the underlying datasets, the user can choose between the simplicity of SPARQL query, with limitations to the endpoint (for e.g., when using built-in geospatial functions of a given endpoint), or the expressivity of the `geom` vocabulary regardless the endpoint, with a more longest SPARQL query. The next chapter is about consuming datasets published on the Web, where we explore tools for visualizing data published as LOD, and later propose our approach which is a category-based visualization wizard.

## Part II

# Generating Visualizations for Linked Data



## CHAPTER 4

# Analyzing and Describing Visualization Tools and Applications

---

*“I think we have consensus, RDF is something you don’t show your end users.”<sup>1</sup>*

Phil Archer (W3C Data Activity Lead)

## Introduction

According to [67] the main goal of information visualization is to translate abstract information into a visual form that provides new insight about that information, in a clear and precise form. In the field of information visualization, data classification, either quantitative or categorical, is useful for visualization, and can be used to make the difference between tools. For example, hierarchical faceted metadata can be used to build a set of category hierarchies where each dimension is relevant to the collection for navigation. The resulting interface is known as faceted navigation, or guided navigation [68]. However, visualizing structured data in RDF by taking advantage of the underlying semantics is challenging both for the publishers and for the users. On the one hand, publishers need to build nice visualizations on top of their 4-5 star datasets (section 4.2). On the other hand, lay users shouldn’t need to understand the complexity of the Semantic Web stack in order to quickly get insights from the data. Thus, adapting visual tools for exploring RDF datasets can bridge the gap between the complexity of Semantic Web and simplicity offered by the field of information exploration. In this chapter, we survey tools for visualizing both structured data (section 4.1.1) and RDF data section 4.1.2). We then provide a classification of the tools for creating applications in the context of LOD (section 4.1.3), along with the way applications are described on the Web. Section 4.3 describes Linked Data applications, followed by the relevant information to describe applications built on top of government open datasets (section 4.3.2). The chapter ends with a brief summary.

---

<sup>1</sup><https://twitter.com/philarcher1/status/507856407127814145>

## 4.1 Survey on Visualization Tools

In this section, we describe also visualization tools that natively do not take as input RDF data for two reasons:

- those tools are relatively “popular” for analyzing data exposed by the government and agencies (most of them in XLS, CSV) as they quickly make it easy to the users to build chart maps and compare with other datasets. One widely application is in the data journalism where facts are analyzed by those tools without waiting for the semantic publication of the data
- Also these tools have many options for visualizing data and are not totally compatible with RDF resources for visualizing graphs.

### 4.1.1 Tools for visualizing Structured Data

In this section, we review tools that are used to visualize structured data and RDF data. The former categories can be extended to support RDF, while the latter tools are trying to cover visualization charts/graphs as much as possible. The authors in [69] propose an overview of all JavaScript visualization libraries and frameworks that are open source, supporting at least bar charts, line charts and pie charts. However, the focus is rather for tools to create dashboard in open data in general, and two of the tools their report are also included in this section, such as D3.js and Google Charts.

#### 4.1.1.1 Choosel

**Choosel** [70] is built on top of GWT and the Google App Engine (the backend can be modified to run on any servlet container). The client-side framework facilitates interaction with visualization components, which can be wrappers around third party components and toolkits such as the Simile Timeline, Protovis and FlexViz. Choosel can integrate components developed using different technologies such as Flash and JavaScript. It is possible to implement visualization components that are compatible with the Choosel visualization component API. These visualization components can then be used to take advantage of Choosel features such as management of view synchronization, management of selections, and support for hovering and details on demand.

#### 4.1.1.2 Many Eyes

**Many Eyes** [71] is a website that can be used to visualize data such as numbers, text and geographic information. It provides a range of visualizations including unusual ones such as “treemaps” and “phrase trees”. All the charts made in Many Eyes are interactive, so it is possible to change what data is shown and how it is displayed. Many Eyes is also an online community where users can create groups (such as “Ebola Crisis” or “Kobane War”) to organize, share and discuss data visualizations.

Users can also comment on visualizations made by others, which is a good way to improve their work. The authors claim that with Many Eyes, the users “*can build quick and easily visualizations from their own data, with the possibility to share them*”. Data input formats are XLS, Plain text and HTML. The output formats are PNG or embeddable. However, using Many Eyes makes your data and the visualizations created with it public. The license of Many Eyes is proprietary of IBM.

#### 4.1.1.3 D3.js

D3.js [72] is a JavaScript library for manipulating documents based on data. D3 uses HTML, SVG and CSS. D3 combines powerful visualization components, plugins<sup>2</sup> and a data-driven approach to Document Object Model (DOM) manipulation. D3 solves the problem of efficient manipulation of documents based on data. It avoids proprietary representation and affords flexibility, exposing the full capabilities of Web standards such as CSS3, HTML5 and SVG. D3 supports large datasets and dynamic behaviors for interaction and animation.

D3 is intended to gradually replace Protovis<sup>3</sup>, which is another tool to build custom visualizations in the browser, created by the same authors and which is no longer under active development. Although D3 is based on many of the same concepts as Protovis, it improves support for animation and interaction. The difference between D3 and Protovis is in the type of visualizations they enable and the method of implementation. While Protovis excels at concise, declarative representations of static scenes, D3 focuses on efficient transformations: scene changes. This makes animation, interaction, complex and dynamic visualizations much easier to implement in D3. Also, by adopting the browser’s native representations (HTML & SVG), D3 better integrates with other Web technologies, such as CSS3 and other developer tools.

#### 4.1.1.4 Google Visualization API

The Google Visualization API<sup>4</sup> establishes two common interfaces to expose data and visualize it on the Web: (1) to expose data on the Web and (2) to provide data to visualizations [73]. Because the Google Visualization API provides a platform that can be used to create, share and reuse visualizations written by the developer community at large, it provides a means to create reports and dashboards as well as the choice to analyze and display data through the wealth of available visualization applications. Many kinds of visualizations are available. Google Visualization accepts data in two different ways: a direct construction as well as a JSON literal object, instantiated via the object `google.visualization.DataTable`. In the latter, the structure of this JSON format is the convention that Google API data

---

<sup>2</sup><https://github.com/d3/d3-plugins>

<sup>3</sup><http://mbostock.github.com/protovis/>

<sup>4</sup><https://developers.google.com/chart/interactive/docs/reference>

sources are expected to return. So, a `google.visualization.DataTable` can be created using the results of an AJAX response. Thus, Google Visualization API can be used to and visualize RDF data. As long as the URL retrieved returns a Google Visualization JSON, an application can create a `DataTable` and send it to the visual construct by the `draw()` function. The results of a SPARQL query can be converted to the Google Visualization JSON using XSL like one used at RPI for data.gov<sup>5</sup>. A sample performing these steps is presented in the Tetherless World Constellation, named `SparqlProxy`<sup>6</sup>. It performs these steps for a client with a single HTTP request, by providing the URL of a SPARQL endpoint to be queried (using `service_uri`), a query (using `query` or `query-uri`), and a specification for return format as Google Visualization JSON (using `output=gvds`).

All the visualizations are based on the type of the columns/fields of the data. While this is normal for tabular data, it is not the case for data exploiting semantics. In Linked Data, vocabularies are used for modeling datasets in RDF, thus making it difficult to directly reuse those tools. There is a need to build more generic tools that exploits the semantics and reuse the visual tools aforementioned.

#### **4.1.2 Tools for visualizing RDF Data**

Regarding the tools for visualizing Linked Data, the paper [74] analyses in detail the current approaches used to browse and visualize Linked Data, by identifying requirements for two groups of users: tech-savvy and lay-users. As the authors extensively surveyed more generic Linked Data browsers, with text-based presentation and visualization options, they provide some recommendations according to the size of the data such as fine-grained analysis among others. However, they do not target their study on tools that can easily help building visual Semantic Web-based applications. By contrast, our approach is to study the tools used to build innovative applications for detecting the components that could be reused across different domain or scope.

##### **4.1.2.1 Linked Data API**

The Linked Data API (LDA) [75], provides a configurable way to access RDF data using simple RESTful URIs that are translated into queries to a SPARQL endpoint. The API layer is intended to be deployed as a proxy in front of a SPARQL endpoint to support:(i) the generation of documents (information resources) for publishing of Linked Data; (ii) the provision of sophisticated querying and data extraction features, without the need for end users to write SPARQL queries and (iii) delivery of multiple output formats from these APIs, including a simple serialization of RDF in JSON syntax.

ELDA<sup>7</sup> is a Java implementation of the LDA by the company Epimorphics. ELDA

---

<sup>5</sup><http://data-gov.tw.rpi.edu/ws/sparqlxml2googlejson.xsl>

<sup>6</sup><http://data-gov.tw.rpi.edu/ws/sparqlproxy.php>

<sup>7</sup><http://www.epimorphics.com/web/tools/elda.html>

comes with some pre-built samples and documentation, which allow developers to build the specification to leverage the connection between the back-end (data in the triple store) and the front-end (visualizations for the user). The API layer helps to associate URIs with processing logic that extract data from the SPARQL endpoint using one or more SPARQL queries and then serialize the results using the format requested by the client. A URI is used to identify a single resource whose properties are to be retrieved or to identify a set of resources.

#### 4.1.2.2 Sgvizler

Sgvizler [76] is a JavaScript library which renders the result of SPARQL SELECT queries into charts or HTML elements. The tool relies on queries against SPARQL endpoints using visualizations based on Google Visualization API, SPARQLer, Snorql<sup>8</sup> and Spark<sup>9</sup>. All the major chart types offered by the Google Visualization API are supported by Sgvizler. The user inputs a SPARQL query which is sent to a designated SPARQL endpoint. The endpoint must return the results back in XML format or SPARQL Query Results in JSON format. Sgvizler is able to parse the results into Google compatible JSON format and displays the result chart using the Google Visualization API or any customize visualization. Sgvizler needs, in addition to the Google Visualization API, the JavaScript framework jQuery to work. One of the drawback of Sgvizler that it is up to the user to test the query and embed it into the HTML page.

#### 4.1.2.3 Facete

Facete [77] is an exploration tool for (geographical) Linked Data datasets on the Web. Also called “Semmap”, the application allows the user to explore the specific slice of data named ‘facet’ of a Linked Data endpoint in a graphical way, available at <http://144.76.166.111/facete/>. The facet is created by defining a set of constraints on properties in the database. Once the facet is defined, the information in the facet can be clicked-through in a tabular interface and visualized on a map. The user can choose a SPARQL endpoint and graph for listing the content and visualizing the dataset. The application has three main views:

1. Selection: A tree-based structure of the dataset. It shows all items' properties and sub-properties.
2. Data: A tabular representation of the data in one facet. All properties that have been marked with an arrow symbol in the facet tree are shown as columns. The columns contain the property values for every item based on the selected filter criteria.
3. Geographical: A map view showing a representation of the elements with geo-coordinates in the facet.

---

<sup>8</sup><http://dbpedia.org/snorql/>

<sup>9</sup><http://code.google.com/p/rdf-spark>

#### 4.1.2.4 VisualBox

Visualbox<sup>10</sup> is another tool that aims at facilitating the creation of visualizations by providing an editor for SPARQL queries and different visual tools to visualize the data. Visualbox is derived from LODSPeakr [78] mainly based on the Model-View-Component (MVC) paradigm. A visualization is created in a Component consisting of one or more SPARQL queries (models), and usually one (but sometimes more) templates (Views). Visualbox is designed for users that have at least some basic knowledge of SPARQL and an understanding of RDF, and it runs the query on the server side. Visualbox uses Haanga<sup>11</sup>, a template engine that provides a syntax for creating templates by defining markers in a document (usually an HTML page) of the form variable that later will be compiled and replaced by values taken from a data source. One of the drawback of Visualbox is that it cannot be extended with custom visualization third-party filters cannot be used. Currently, it implements visualization filters for D3.js (5), Google Maps, Google Charts(6) and TimeKnots library (TimeLine with events)<sup>12</sup>.

#### 4.1.2.5 Payola

Payola [79] is a Web framework for analyzing and visualizing Linked Data, and enables users to build instances of Linked Data visualization Model (LDVM) pipelines [80]. LDVM is an adaptation of the Data State Reference Model (DSRM) proposed by Chi [81] applied to visualizing RDF and Linked Data. It extends DSRM with three additional concepts that are reusable software components:

- **Analyzers:** They take as input compatible datasets (hierarchical dataset, geocoordinates dataset, etc) and perform adapted SPARQL queries: .
- **Visualization transformers:** They can be any software component that transforms data between different formats or performs aggregations for better visualization. They are generally SPARQL CONSTRUCT queries, with input signatures corresponding to the FROM clauses and output data samples corresponding to the CONSTRUCT clauses.
- **Visualizers:** They consume RDF data and produce a visualization a user can interact with. They are visual tools libraries consuming data in RDF/JSON<sup>13</sup>.

A developer builds different instances of LDVM based on the datasets used in the analyzers and transformers. Figure 4.1.2.5 depicts a sample of a LDVM pipeline applied to two different datasets published as LOD.

---

<sup>10</sup><http://alangrafu.github.io/visualbox/>

<sup>11</sup><http://haanga.net>

<sup>12</sup><https://github.com/alangrafu/timeknots>

<sup>13</sup><https://dvcs.w3.org/hg/rdf/raw-file/default/rdf-json/index.html>

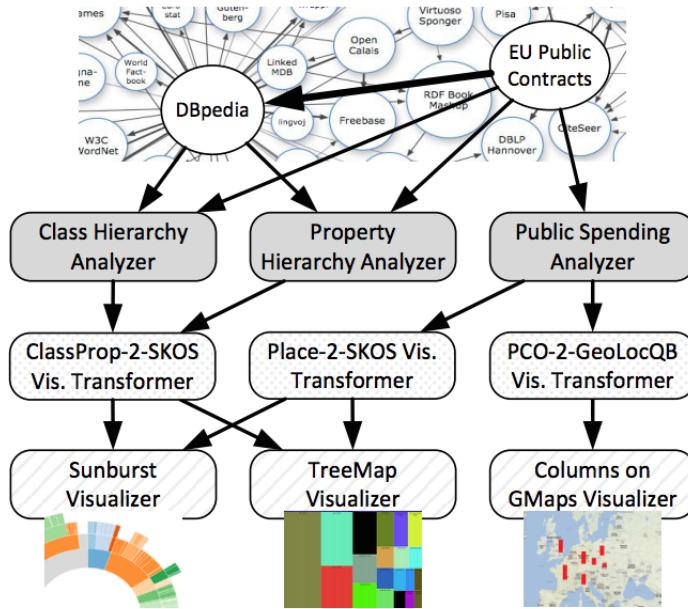


Figure 4.1: Sample application of analyzers and visualizers in a LDVM pipeline.

### 4.1.3 Discussion

There are currently many projects aiming at visualizing (RDF) Linked Data. A survey by Dadzie and Rowe [82] concluded with the fact that many visualization tools are not easy to use by lay users. In [83] a recent review of some visualizations tools that can be summarized as follows:

- *Vocabulary-based visualization tools*: these tools are built for specific vocabularies and they help in visualizing data modelled according to those vocabularies, such as CubeViz [84], FOAF explorer<sup>14</sup> and Map4RDF [85]. They aim at visualizing data modelled respectively with `dq`, `foaf` and `geo+scovo`.
- *Mashup tools*: they are used to create mashup visualizations with different widgets and some data analysis, such as DERI Pipes [86]. Mashup tools can be integrated into the LD wizard to combine different views of the data.
- *Generic RDF visualization tools*: they typically support data browsing and entity rendering. They can also be used to build applications. In this category, we can mention Graphity<sup>15</sup>, lodlive<sup>16</sup> and Balloon Synopsis<sup>17</sup>.

While above-mentioned tools are often extensible and support specific domain datasets, they suffer from the following drawbacks:

<sup>14</sup><http://foaf-visualizer.gnu.org.ua/>

<sup>15</sup><https://github.com/Graphity/graphity-browser>

<sup>16</sup><http://en.lodlive.it/>

<sup>17</sup><https://github.com/schlegel/balloon-synopsis>

- *They are not easy for lay users to set up and use.* Sometimes, users just need a visual summary of a dataset in order to start exploring the data. Our approach to this challenge is to provide such a lightweight JavaScript-based tool that supports a quick exploration task.
- *They do not make recommendation based on categories.* A tool similar to our approach is Facete<sup>18</sup>[77] which shows a tree-based structure of a dataset based on some properties of an endpoint more relevant for geodata. A tabular view enables slices of data to be visualized and a map view can be activated when there is geodata. Our approach aims to be more generic, offering more views (tabular, map, graph, charts, etc.) based on a systematic analysis of the high-level categories present in a dataset.

The outcome of this state of the art can then be used to assess different visual tools in the process of creating Web-based visualizations. Some criteria can be used for assessing visual tools, such as (i) usability, (ii) visualization capabilities, (iii) data accessibility, (iv) deployment and (v) extensibility. In [87], the readers can find more details on this survey. Table 4.1 gives an overview of the selected tools studied based on the following features:

- *Data Formats* for the format of data taken as input by the tool;
- *Data Access*, for the way to access the data from the tool, such as Web service, SPARQL endpoint, etc.
- *Language code*, the programming language used to develop the tool;
- *Type of views*, the different views potentially accessible when using the tool;
- *Imported libraries*, the external libraries available within the tool,
- *License* for the Intellectual Property rights of the tool,
- *SemWeb compliant*, whether the tool can be easily extended or is compliant with structured data; and

---

<sup>18</sup><http://cstadler.aksw.org/facete/>

Table 4.1: Survey of tools used for creating visualizations on the Web.

Tools Views	Data Formats Libraries	Data Access License	Language SemWeb App
Choosel Text/Map/Bar chart	XLS, CSV Time (Simile)/ProteoVis/Flexvis	API Open	GWT No
Fresnel Property/Labels	RDF Welkin/IsaViz/Haystack/CSS	- Open	RDF Yes
Spark Charts/Tabular	RDF-JSON -	SPARQL Open	PHP Yes
LDA -	RDF -	SPARQL Open	Java/PHP Yes
SemWeb Import Graph Node	RDF -	SPARQL CECILL-B	Netbeans Yes
Many Eyes Charts/Trees/Graphs/Maps	XLS/Text/HTML -	API IBM	Java/Flash No
D3.js Charts/Trees/Graphs/Maps	CSV/SVG, GeoJson Jquery/sizzle/colorbrewer	API Open	JavaScript Maybe
Facet Map, Facet view	RDF-JSON JQuery/ dynatree	SPARQL Open	JavaScript Yes
Sgvizler Map/Line chart, Timeline/Sparkline	RDF-JSON Google visualization API	SPARQL Open	JavaScript Yes
Visual Box Map/Charts/ Timeline/Graphs	RDF Google charts/TimeKnots/D3.js	SPARQL Open	PHP/Django Yes
Map4rdf Facet/Map	RDF-JSON OSM Layers, Google Maps	SPARQL Open	Java/GWT Yes
Exhibit Map/Tile/Thumbnail/Tabular/Timeline	JSON Exhibit -	Data dump Open	JavaScript Yes
Google Visualization API Charts/ Charts/Maps/Dashboard	JSON/CSV AJAX API	API Open	JavaScript Possible

## 4.2 Describing Applications on the Web

### 4.2.1 Motivation

As many initiatives on Linked Open Data are growing, tools and technologies are getting more and more mature to help consumers to leverage the lifting process of the data. At the same times, standardization bodies such as the W3C are helping by providing best practices to publish Open Government Data by using appropriate vocabularies, taking care of stability in the URIs policies, and making links to other datasets. For instance the Government Linked Data Working Group<sup>19</sup> has released best practices and vocabularies to help governments publishing their data using Semantic Web technologies. Having a look at different proposals of the Life Cycle of Government Linked Data, one of the last stage is “Publication” where the data is released according to the 4-5 star principles<sup>20</sup>, with a given access to a SPARQL endpoint. However, for a better understanding of the data, one of the next steps is usually to generate visualizations through intuitive visual tools(charts, graphs, etc.) that will benefit to citizens, data journalists and other public authorities. Currently, one way of creating new applications is to look around previous initiatives to see what type of applications exist already and to make something similar according to a given dataset and domain. Another approach is by organizing *contests* where the challenges are to mash up unexpected datasets with clear and beautiful visualizations. This approach is harder because developers also try figure out which tool and library is used for different applications. What if we describe applications according to the facets/views, datasets, visual tools used to build them? How are the types of information that can help create a vocabulary for annotating Web-based visualizations online?

### 4.2.2 Catalogs of Applications

We provide below two use cases of the current description of applications developed by datasets published on the Web. We expose also the limitation of the approach as they do not fully make use of semantics for improve discovery of the visual tools and datasets used to develop such applications.

#### 4.2.2.1 Open Data Service

The Open Data Service at the University of Southampton<sup>21</sup> has a register of all the applications developed using their datasets. A catalog of the Applications using the data is available at <http://id.southampton.ac.uk/dataset/apps>. Each application is described by giving three main categories of information:

1. The available distributions corresponding to the different formats HTML, RDF/XML and RDF/Turtle ;

---

<sup>19</sup><http://www.w3.org/2011/gld/>

<sup>20</sup><http://5stardata.info/>

<sup>21</sup><http://data.southampton.ac.uk/apps.html>

<a href="http://id.southampton.ac.uk/app/soton-map-amenities">http://id.southampton.ac.uk/app/soton-map-amenities</a>	
Searchable map for finding buildings, amenities and bus stops in and around University sites	
App type:	Web
Created:	6th March 2011
Created by:	Colin R. Williams emax Jarutas Pattanaphanchai
Uses:	Buildings and Places Southampton Bus Information Local Amenities Catering Teaching Room Features

Figure 4.2: Sample description of a Web application at the Open Data Service

2. Dataset information, which defines the type, the number of triples, license information, the publisher and the publication date.
3. The provenance, such as files used to generate the dataset for building the application, as well as the script itself.

Currently, some vocabularies are used to model the catalog, such as the DCAT vocabulary [23] and proprietary vocabularies. Each application is then described the type (Web, mobile Web, android, etc.), the authors, the date of creation and the datasets used to build the application. Figure 4.2.2.1 depicts the HTML view of a Web application for a searchable map for finding buildings within the University sites. This initiative seems to be not extended. Thus, there is a real need to have a common layer of semantics for describing such applications. This would benefit the interoperability and discovery of applications on the Web.

#### 4.2.2.2 RPI Applications

Another approach, taken by the researchers at the Rensselaer Polytech Institute<sup>22</sup>, is to put at the bottom of the static page of a demo/application showcasing the benefits of Open Data for [data.gov](http://data.gov)<sup>23</sup> some basic metadata (description, URL to dataset, author), and also a link to the SPARQL query used for generating the application. As this information is human-readable and can help, the main drawback is the lack of a machine-readable version, using semantics to discover and connect different demos and datasets with authors. A more vocabulary can leverage the issue by annotating

<sup>22</sup><http://data-gov.tw.rpi.edu>

<sup>23</sup><https://www.data.gov/>

such applications to help discovering and aggregating other similar applications in other Open Data initiatives.

## 4.3 Describing and Modeling Applications

According to [88], *Visualization* is “*the use of computer-supported, interactive visual representations to amplify cognition*”. So the unique object of visualization is to develop insights from collected data. That explains why each time a new dataset is released, users always expect some showcases to play with the underlying datasets. It is true that many public open initiatives use incentives actions like *challenges*, *datahack-day* or *contest*, etc. to find innovative applications that actually exhibit the benefits of datasets published. Visualizations play crucial role as they can easily find errors in a large collection, detect patterns in a dataset or help navigate through the dataset.

### 4.3.1 Typology of Applications

Jeni Tennison defines in her blog<sup>24</sup> three categories of applications using online data:

- (i) *data-specific applications*, which are constructed around particular data sets that are known to the developer of the application; hence the visualizations obtained are of data-specific applications. Examples are the famous applications of “*Where does my money go*” in Greece<sup>25</sup> or the UK<sup>26</sup>. Those applications are also called “*mashups*”.
- (ii) *vocabulary-specific applications*, which are constructed around particular vocabularies, wherever the data might be found that uses them. Examples here are the FaceBook Social Graph API<sup>27</sup> and IsaViz [89], among others.
- (iii) *generic applications*, which are constructed to navigate through any RDF that they find; e.g., Tabulator [90], OpenLink Data Explorer<sup>28</sup>.

Because most mash-ups are data-specific applications, it is important to know what information the dataset contains. This could be achieved by giving the meaning of some properties or classes of the vocabularies used to create the dataset. Hence, what the data publisher very often needs to do is to make sure that the data they publish is documented. However, what is used in practice, is to consider using an intuitive visualization self-descriptive to both show the added-value of the data and its documentation.

---

<sup>24</sup><http://www.jenitennison.com/blog/node/126>

<sup>25</sup><http://publicspending.medialab.ntua.gr/en/index.php>

<sup>26</sup><http://wheredoesmymoneygo.org/>

<sup>27</sup><https://developers.facebook.com/docs/plugins/>

<sup>28</sup><http://ode.openlinksw.com/>

Table 4.2: Gathering reusable information from the openspending in Greece application

Features	Value
Access Url	<a href="http://publicspending.medialab.ntua.gr/">http://publicspending.medialab.ntua.gr/</a>
Scope/Domain	Public spending, Government
Description	The application helps visualizing the most characteristic facts of the Greek public spending, interconnected to foreign expenditure and other data.
Supported Platform	Web
URL Policy	<a href="http://BASE/en/NAME-CHART.php">http://BASE/en/NAME-CHART.php</a> e.g., <a href="http://{BASE}/en/toppayersday.php">http://{BASE}/en/toppayersday.php</a>
Data Source	<a href="http://opendata.diavgeia.gov.fr">http://opendata.diavgeia.gov.fr</a> ; Greek Tax data (TAXIS)
Type of views	Bubble tree, column and bar charts
Visualization tools	HighchartsJS, Bubble TreeJS JqueryJS ; RaphaelJS
License	Open
Business Value	Not Commercial (Free)

### 4.3.2 Reusable applications

Many applications are built on top of datasets exposed in different open data governments initiatives. Generally, they are used to provide insight about the datasets and their usefulness. However, some of the applications could be generalized and reused if published adequately. Having some best practices in publishing applications on the Web could boost the interoperability between datasets and visual tools. To achieve this task, we first review some applications that have been developed on top of datasets opened by governments (UK, USA, France) and public local authorities. We made a random survey of thirteen (13) innovative applications [91] in various domains such as of security, health, finance, transportation, housing, city, foreign aid and education. Table 4.3 provides a summary of the surveyed applications; with names, types, countries and brief description.

The main template used in the survey gathered the following information:

- the name of the application;
- the scope or the target domain of the application;
- a small and concise description;
- the platform on which the application can be deployed and view;
- the policy used for creating the URL of the application;

- the legacy data used to build the application, and a mention of the process of the “lifting process” of the raw data to RDF if available;
- the different views available from the application;
- comments or relevant drawbacks;
- the license of the application.

Table 4.2 provides the information extracted from `openspending` in Greece using the aforementioned template. Such information can be published using a vocabulary to help discover all the applications built on top of public spending data across different platforms.

Table 4.3: Some innovative applications built over Open Government Datasets

<b>Application</b>	<b>Domain</b>	<b>Type</b>	<b>Country</b>
UK Crime	Crimes	Web	UK
UK Pharmacy	Health, Pharmacy	Mobile/Android	UK
Numberhood	Local area dynamics	iPhone/iPad	UK
BuSit London	Public Transportation	Web and mobile	UK
UK School Finder	Education	Web	UK
Where-can-I-Live	Homes, Transportation	Web	UK
Opendatacommunities	Local Government	Web	UK
FlyOnTime	Flights/airlines	Web	USA
White House Visitor Search	White House	Web	USA
US-USAID/UK-DFID	Foreign aid	Web	USA
Fourmisante	Medicine/Health-care	Web	France
MaVilleVueDuCiel	Local Government	Web	France
Home'n'Go	Housing	Web	France

#### 4.3.3 A vocabulary for Describing Visualization Applications

We have implemented a vocabulary, DVIA<sup>29</sup>, that aims to describe any applications developed to consume datasets in compliant with 4-5 stars rate, using visual tools to showcase the benefits of Linked Data. It reuses four existing vocabularies: Dublin Core<sup>30</sup>, dataset catalogue (DCAT)<sup>31</sup>, Dublin Core Metadata Initiative<sup>32</sup> and the

<sup>29</sup><http://bit.ly/Vb4L8k>

<sup>30</sup><http://purl.org/dc/terms/>

<sup>31</sup><http://www.w3.org/ns/dcat#>

<sup>32</sup><http://purl.org/dc/dcmitype>

Organization vocabulary <sup>33</sup>. It is composed of three main classes:

- **Application:** This class represents the application or the mashup developed for demoing or consuming data in LD fashion. It is a subclass of **dctype:Software**
- **Platform:** The platform where to host or use the application, could be on the Web (Firefox, Chrome, IE, etc..) or mobile (android, iOS, mobile ) or even

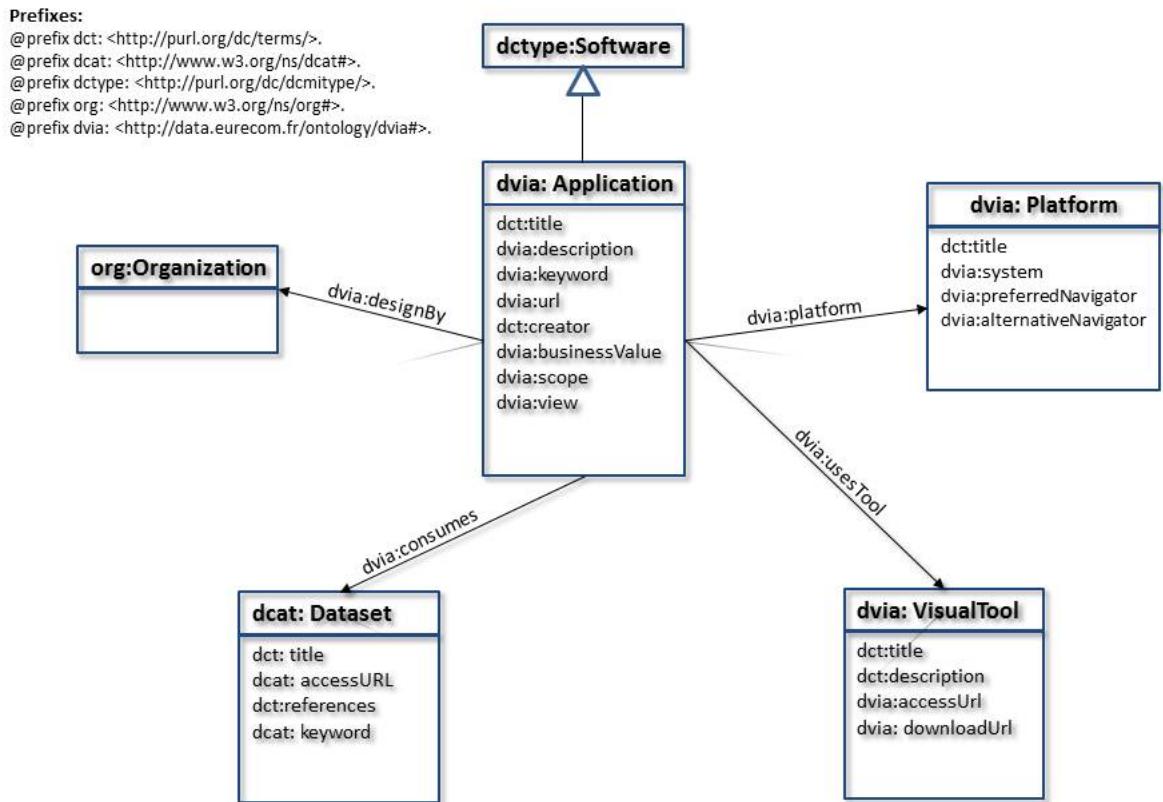


Figure 4.3: Conceptual Model of the DVIA vocabulary

The main classes and properties are depicted in Figure 4.3. The current version of the vocabulary in Turtle format can be found at <http://purl.org/ontology/dvia>. Listing 4.1 is a snapshot of the description of *EventMedia Live*<sup>34</sup> application in DVIA. It depicts apart from some metadata about the application (`dct:title`, `dct:name`, `dct:issued`, `dct:creator` and `dct:license`), the different visualization libraries integrated for building EventMedia Live (e.g. The Google API, Backbone, etc), as

<sup>33</sup><http://www.w3.org/ns/org#>

<sup>34</sup><http://challenge.semanticweb.org/2012/winners.html>

well as the operating systems where it is designed for, the different views/facets available in the application (map, charts, graphs, force-directed layout, etc.) and the heterogeneous datasets used to implement it.

```

1 visuapp:eventMedia01
2   a dvia:Application ;
3   dct:title "EventMedia Live"@en;
4   dvia:description "An application for reconciling Live events with media" ;
5   dvia:url <http://eventmedia.eurecom.fr> ;
6   dct:issued "2012-11-10"^^xsd:date ;
7   dvia:businessValue "not commercial" ;
8   dvia:keyword "events, media"^^xsd:string ;
9   dct:license <http://www.opendatacommons.org/licenses/pddl/1.0/> ;
10  dvia:platform [ a dvia:Platform ;
11    dct:title: "Desktop" ;
12    dvia:preferredNavigator "Google Chrome" ;
13    dvia:alternativeNavigator "FireFox" ;
14    dvia:system "Mac OS, Windows, Linux"^^string ] ;
15  dvia:usesTool [ a dvia:visualTool;
16    dct:title "Google visualization Tool" ;
17    dct:description "Google visualization API" ;
18    dvia:accessUrl <https://developers.google.com/chart/interactive/docs/reference> ;
19    dvia:downloadUrl <http://www.google.com/uds/modules/gviz/gviz-api.js/> ] ;
20  dvia:usesTool visuapp:visualTool02 ;
21  dvia:consumes [ a dcat:Dataset; dct:title "BBC dump" ] ;
22  dvia:consumes [ a dcat:Dataset; dct:title "last.fm scrapped dataset" ] ;
23  dvia:consumes [ a dcat:Dataset; dct:title "upcoming scrapped dataset" ] ;
24  dvia:consumes [ a dcat:Dataset; dct:title "eventful scrapped dataset" ] ;
25  dvia:consumes [ a dcat:Dataset; dct:title "Flickr scrapped dataset" ] ;
26  dvia:consumes [ a dcat:Dataset; dct:title "Music Brainz" ] ;
27  dvia:consumes [ a dcat:Dataset; dct:title "Foursquare Json file" ] ;
28  dvia:consumes [ a dcat:Dataset; dct:title "DBpedia" ] ;
29  dct:creator [foaf:mbox "khrouf@eurecom.fr"; foaf:name "Houda Khrouf"];
30  dct:creator [foaf:mbox "vuk@eurecom.fr"; foaf:name "Vuk Milicik"];
31  dct:creator [foaf:mbox "raphael.troncy@eurecom.fr"; foaf:name "Raphael Troncy"];
32  dvia:view "map, chart, graph, force-directed layout" ;
33 ...

```

Listing 4.1: Snapshot in Turtle of the description of Event Media Live Application

The full version of this sample description is available at <http://www.eurecom.fr/~atemezin/datalift/visumodel/eventMedia-sample.ttl>. The current version of the DVIA intends to be small enough to cover the concepts that are needed to reuse partial or full parts of applications.

## 4.4 Summary

In this chapter, we have described different tools used for visualizing data, structured and graph data. We have also discussed different types of applications currently built on top of government open data initiatives. The goal of this survey is to propose some new approaches of generating and discovering visualizations and applications on the Web of Data. We designed and implemented DVIA, a vocabulary that aims to model applications for more interoperability and discovery of applications and tool visualizations on the Web.

## CHAPTER 5

# Creating and Generating Visual Applications

---

*“A Semantic Web application is one whose schema is expected to change.”*  
(David Karger, MIT CSAIL)<sup>1</sup>

## Introduction

The objective of visualization is to develop insights from collected data. Moreover, according to Information Theory, human vision is the sense that has the largest bandwidth (100 Mbits/s), which makes it the best suited channel to convey information to the brain [92]. Based on the Visual Information Seeking Mantra of Shneiderman: *“overview first, zoom and filter, then details on demand”* [93], we advocate for more visual interactive representations of RDF graphs using SPARQL endpoints. At the same time, we use the term “Linked Data Visualization”, to refer to a *combination of charts, graphics, and other visual elements built on top of 4-5 star datasets accessible via a SPARQL endpoint*. Linked Data offers some great advantages for publishing government data. The approach makes it easy to publish information in a way that allows it to be combined with other sets of data. The benefits also arise from the semantics associated to things, common identifiers with resources, from the inherent extensibility of the RDF data model, and from the publication of data in a standard format. Linked Data is a great way of publishing information for diverse and distributed organizations, such as government [94]. Despite the presence of more and more datasets published as Linked Data, there is still a need to help end users to discover what (unknown) datasets describe by hiding the complexity of SPARQL queries from such users. The RDF model, its various serializations and the SPARQL query language are foreign to the majority of developers who understandably want to be able to use the tool chains that they are familiar with to access government data. Sometimes, publishing data purely as RDF, and providing access purely through SPARQL queries raises an unacceptable barrier onto the use of that data. Moreover, the task of identifying the key categories of datasets can help in selecting and matching the most suitable visualization types. In this chapter, we present our contribution on consuming datasets by generating applications target to lay-users. The remainder of this chapter is structured as follows. We first propose in section 5.1.1 some important categories that are worth

---

<sup>1</sup>A statement during his keynote at ESWC 2013 conference

visualizing and in section 5.1.2 a set of mapping views associated with vocabularies. In Section 5.2, we describe the implementation of a wizard that can work on top of any RDF dataset. We detail the results of an experiment where high level categories and associated visualizations have been generated from numerous SPARQL endpoints (Section 5.3.2). We do reverse engineer the GKP (Section 5.3) to look for the most important properties of an Entity. Then, we present two domain applications around events (Section 5.5) and statistics (Section 5.6). Finally, we discuss how to improve the discovery of applications developed in open Data events (Section 5.7) by proposing a vocabulary and a plugin to easily annotate their Web pages and generate RDF data.

## 5.1 Wizard for Visualizations: Theoretical foundations

### Background

With the growing adoption of the Linked Data principles, there is a real need to support data consumers in quickly exploring a dataset through visualizations. In order to involve more general Web users into the Semantic Web and Linked Data world, tools that reuse existing visualization libraries are needed for showing the key information about RDF datasets. Many datasets are published on top of SPARQL endpoints and yet, are not “visually” accessible. Thus, understanding the underlying graphs and consuming them require lay users to have some knowledge in writing queries.

In this section, we propose a first step towards making available a semi-automatic way for the production of possible visualization of linked data sets of high-level categories grouping objects that are worth viewing and we associate them with some very well known vocabularies. Then, we describe the implementation of a Linked Data Visualization Wizard and its main components. This wizard can be used to easily visualize slices of datasets based on generic types detected.

#### 5.1.1 Dataset Analysis

When developing an application, there are some “*important*” classes/categories, objects or datatypes that can be detected first to help to guide in the progress of creating a set of visualizations tied with those categories. We distinguish seven categories while acknowledging that this is not necessarily an exhaustive list:

- [Geographic information]: This category is for data generated and modeled using `geo:SpatialThing`, `dbpedia-owl:Place`, `schema:Place` or `gml:_Feature` classes.
- [Temporal information]: This category also includes datasets containing date, time (e.g: `xsd:dateTime`) and period or interval of time, using the OWL Time ontology.

- **[Event information]:** This category is for any action or activity occurring at some place at some time.
- **[Agent/Person information]:** This category is heavily influenced by the use of `foaf:Person` or `foaf:Agent`.
- **[Organization information]:** This category is related to organizations or companies data, with the use of the `org` vocabulary<sup>2</sup> or the `foaf:Organization` class.
- **[Statistical information]:** This category refers to statistical data generally modeled using the `data cube` vocabulary<sup>3</sup> or the SDMX model<sup>4</sup>.
- **[Knowledge classification]:** This category refers to dataset describing schemas, classifications or taxonomies using the `SKOS` vocabulary.

### 5.1.2 Mapping Datatype, Views and Vocabularies

The On-line Library of Information Visualization Environments (OLIVE)<sup>5</sup> is a web site describing eight categories of information visualization environments differentiated by data type and collected by students, following a visualization course given at Maryland College Park, mostly inspired from the work of Ben Shneiderman [93]. Based on the classification provided by OLIVE, we propose a set of mappings between those categories (excluding the workspace dimension), views that can be applied to this category and a suitable list of vocabularies from the Linked Open Vocabularies catalogue [65]<sup>6</sup> that correspond to those categories. Those vocabularies are easy to be found as there is a manual classification of vocabularies by the curators of the catalogue based on the content and scope of the terms and properties. According to the seven categories defined in Section 5.1.1, we have identified some of their corresponding one to one mapping with the set of vocabularies:

- **Geography** space, consisting of 21 vocabularies for features: `geo`, `gn`, `gf`, `om`, `geop`, `md`, `lgdo`, `loc`, `igeo`, `osadm`, `geod`, `ostop`, `place`, `geos`, `locn`, `coun`, `postcode`, `osr`, `geof`, `g50k` and `ad`.
- **Geometry** space, for vocabularies dealing with the geometries, mostly combined with the features, such as:
- **Time** space, consisting of 14 vocabularies, `cal`, `date`, `gts`, `interval`, `ncal`, `oh`, `te`, `thors`, `ti`, `time`, `tl`, `tm`, `tvc` and `tzont`.
- **Event** space, containing vocabularies such as `event`, `lode`, `music`, `sem`, `situ`, `sport`, `stories`, `theatre`, `tis` and `tisc`.

<sup>2</sup><http://www.w3.org/TR/vocab-org/>

<sup>3</sup><http://www.w3.org/TR/vocab-data-cube/>

<sup>4</sup><http://sdmx.org/>

<sup>5</sup><http://lte-projects.umd.edu/Olive/>

<sup>6</sup><http://lov.okfn.org/dataset/lov/>

- **Government** space, with 9 vocabularies (`cgov`, `ctorg`, `elec`, `few`, `gc`, `gd`, `oan`, `odd`, `parl`) and the `org` vocabulary belonging to the W3C recommendation vocabularies at <http://lov.okfn.org/dataset/lov/lov#W3C>.

Vocabularies used to provide metadata information, such as `rdfs`, `dcterms` or `dce` can be used in association with any of the visual elements to give a basic description of the resource of a given dimension. For example, popup information can be shown on a map view to display the relevant information of a geodata resource such as the label, the abstract or the description. Another application is to detect which visualization is best suited for geodata. Geodata belongs to a two-dimension visual representation. Geodata is usually displayed using geographical-based visualizations (map, geo charts, etc.) and it is often modeled by vocabularies in the space named **Geometry and Geography**<sup>7</sup> vocabularies in RDF datasets. Hence those vocabularies can be combined to detect whether a dataset contains geographic information, and thus determine whether to recommend a map view. Table 5.1 gives an overview of those mappings. The tabular representation is the “default” visual representation of RDF data and can be used by any vocabulary without restriction.

Dimension	Vocabulary Space	Visual element
Temporal	Time space	TimeLine
One dimension	any	Tabular, text
Two dimension	Geography space	Map view
	Geometry space	Maps view
Three dimension	Event space	Map + TimeLine
Multi dimension	<code>qb</code> , <code>sdmx-model</code> , <code>scovo</code>	Charts, graphs
Tree	<code>skos</code> , Government space	Treemap, Org view
Network	any vocab.	Graph, network map

Table 5.1: A taxonomy of information visualizations for consuming Linked Datasets with suitable vocabulary space and visual elements.

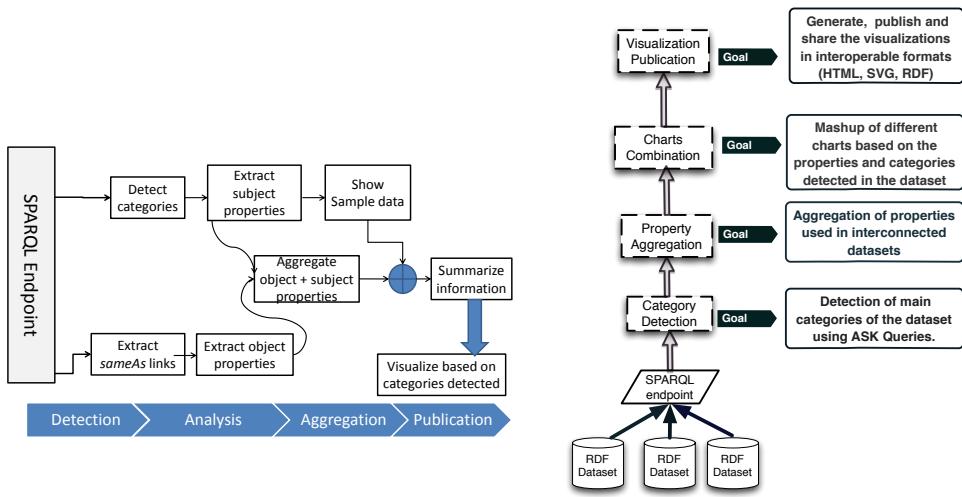
## 5.2 LDVizWiz: a Linked Data Visualization Wizard

### 5.2.1 Workflow

We propose a workflow composed of four principal steps for a Linked Data Visualization Wizard, as depicted in Figure 5.1. Our goal is to provide a tool that hides the complexity of SPARQL from lay users and at the same time, can be embedded in the existing Linked Data infrastructure and workflow. First, we proposed detecting the presence of data belonging to one of the seven categories (Table 5.1) by using generic SPARQL queries. More precisely, we perform ASK queries to test whether or not a particular query pattern has a solution. Second, we look at entities in a

<sup>7</sup>All the prefixes used for the vocabularies are the same used in LOV catalog.

dataset that have `owl:sameAs` links with external objects and we retrieve the properties associated to those objects. We argue that the objects that are interlinked with other datasets are of primary importance in a visualization. We show the results of this mining process to the user (the categories that have been detected, the properties going with the categories and the external domain). Based on this information, the user can make a personalized “mashup” by aggregating suitable visualization widgets. Some default visualizations are available according to the top categories detected. The last step is to publish the visualization and a report in RDF/XML, Turtle or N3.



(a) The workflow of the different modules interacting in the Linked Data visualization wizard. (b) High level functionalities of the Linked Data visualization wizard

Figure 5.1: Big picture and architecture of the Linked Data visualization wizard.

Let consider a graph  $\langle G, c \rangle$  to be  $G = \{(s, p, o) | p \in URI, s \in URI, o \in (URI \cup LIT)\}$  where  $URI$  is the set of URIs,  $LIT$  is the set of literals, and  $c$  is the context. We define  $L = \{V_1, V_2, \dots, V_n | V_i = P_i \cup T_i\}$  the list of vocabularies in LOV, with  $P_i$  and  $T_i$  respectively the properties and a vocabulary  $V_i$ . Let also  $D = \{D_1, D_2, \dots, D_m\}$  be the domains of vocabularies. We assume  $\Phi(L, D) = \{\forall V \in L, \exists D_k, \Phi(V) \in D_k\}$ . A domain of a vocabulary is considered to be the primary scope covered by the vocabulary. We define a generic function  $\Sigma : (G, c) \mapsto B$  to detect categories in a dataset as follows:  $\Sigma((G, c)) = \{B | (\exists(s, p, o) \in G : p \in V) \cup (\exists(s, rdf:type, o) \in G : o \in V)\}$  where  $B = \{True, False\}$ .

In the following sections, we describe each of the steps involved in the Linked Data Visualization Wizard in more detail.

### 5.2.1.1 Category Detection

The goal of the category detection task is to use SPARQL queries to detect the presence of some high level categories in the dataset. We perform ASK queries as implementation of the  $\Sigma$  function using standard vocabularies as defined in the Table 5.1. We start with six domains, namely: geographic information, person, organization, event, time and knowledge organization systems. We select popular vocabularies based on two existing catalogues: LOV [95] and prefix.cc<sup>8</sup>.

```

1 ASK WHERE {
2 {
3 ?x a ?o.
4 filter (?o= dbpedia-owl:Place ||
5 ?o=gml:_Feature ||
6 ?o=geo:SpatialFeature || ?o=gn:Feature ||
7 ?o=admingeo:CivilAdministrativeArea ||
8 ?o=spatial:Feature ||
9 ?o=vcard:Location)
10 }
11 UNION {
12 ?x ?p ?o. filter(?p=geo:lat || ?p=geo:long ||
13 ?p=georss:point || ?p=geo:geometry ||
14 geom:geometry)
15 }
16 }
```

Listing 5.1: Generic query to detect geo data from a SPARQL endpoint.

Listing 5.1 shows seven classes of different vocabularies are used, respectively for the namespaces `dbpedia-owl`, `geo`, `gn`, `admingeo`, `spatial` and `vcard`, with relevant classes to check the presence of geographic data.

```

1 ASK WHERE {?x a ?o. filter(?o=time:TemporalEntity ||
2 ?o=time:Instant ||
3 ?o=time:Interval || ?o=dbpedia-owl:TimePeriod ||
4 ?o=time:DateTimeInterval || ?o=intervals:CalendarInterval)
5 }
6 UNION{ ?x ?p ?o. filter(?p=time:duration ||
7 ?p=time:hasBeginning ||
8 ?p=time:inDateTime || ?p=time:hasDateTimeDescription
9 || ?p=time:hasEnd)}
```

Listing 5.2: Generic query to detect time data from a SPARQL endpoint, using `time`, `dbpedia-owl`, `intervals` vocabularies.

Listing 5.2 detects the presence of time information, while Listing 5.3, 5.4 and 5.5 detect persons, organizations and events respectively.

```

1 ASK WHERE {?x a ?o. filter(?o = foaf:Person ||
2 ?o=dbpedia-owl:Person ||
3 ?o=vcard:Individual) }
```

Listing 5.3: Generic query to detect person categories from a SPARQL endpoint, using `foaf`, `dbpedia-owl`, `vcard` vocabularies.

<sup>8</sup><http://prefix.cc>

```

1 PREFIX org:<http://www.w3.org/ns/org#>
2 PREFIX foaf: <http://xmlns.com/foaf/0.1/>
3 PREFIX dbpedia-owl: <http://dbpedia.org/ontology/>
4 ASK WHERE {?x a ?o . filter (?o=org:Organization ||
5 ?o=org:OrganizationalUnit ||
6 ?o=foaf:Organization ||
7 ?o=dbpedia-owl:Organisation)}
```

Listing 5.4: Generic query to detect ORG data from a SPARQL endpoint.

```

1 ASK WHERE{?x a ?o. filter (?o= lode:Event || ?o=event:Event ||
2 ?o=dbpedia-owl:Event)}
```

Listing 5.5: Generic query to detect event data from a SPARQL endpoint, using `lode`, `event`, `dbpedia-owl` vocabularies.

For detecting data organized as taxonomy, `skos` vocabulary is used along with the most used classes and properties as showed in Listing 5.6.

```

1 ASK WHERE {{?x a ?o. filter(?o=skos:Concept ||
2 ?o=skos:ConceptScheme || ?o=skos:Collection )}|
3 UNION{ ?x ?p ?o. filter(?p=skos:featureCode ||
4 ?p=skos:altLabel || ?p=skos:prefLabel || ?p=skos:relatedMatch)}}
```

Listing 5.6: Generic query to detect SKOS data from a SPARQL endpoint, using `skos` vocabulary.

### 5.2.1.2 Property Aggregation

We take the benefits of the `owl:sameAs` links between entities to have access to the properties of the entities in the external namespaces different from the origin dataset. This module also aggregates the properties found in the dataset with the ones found in the interlinked sets. This is based on the assumption that during the linkage process, external datasets not only help in not breaking the *follow-your-nose* principle, but also add more information to be viewed in visualization applications. As shown in the code below, at this stage, we have collected and aggregated external properties gathered from the enrichment process of the workflow.

```

1-LET Namespace(?s) = S and LET Namespace(?t) =T
2-SELECT owl:sameAs links
LET SEMTERM = list of ?s owl:sameAs ?t
WITH T != S
3-IN T, SELECT distinct properties used in dataset
4-AGGREGATE (3) with properties FROM S.
```

### 5.2.1.3 Visualization Generator

This module aims at recommending the appropriate visualizations based on the categories detected by the wizard. It might also help the user to make a report

summarizing the result of the mining process, and then use different visualization libraries to view the data. This module can be viewed as a recommender system because it derives visualizations based on the categories. The input to build each visualization is the corresponding SELECT query of each ASK queries used to detect the categories. Moreover, some adjustment are made to avoid blank nodes and to get the labels of the resources. The generator can be coupled with a mashup widget generator for some categories. For example, users could expect for event data, a combination of map view (where), a timeline (when) and facets based on the agents (who).

#### 5.2.1.4 Visualization Publisher

The publisher module aims at exporting the combined visualizations, along with the report of all the process of mining the dataset, in a format easy to share, either as HTML, SVG or in the different RDF syntax flavor. For the latter, apart from using metadata information (creator, issued date, license), we model the categories we have detected using the `dcterms:subject` property of a `dcat:Dataset`, the queries used (using the `prov:wasDerivedFrom` property), the sample resources for each category (using the `void:exampleResource` property) and the visualization generated (using the `dvia` and `chart`<sup>9</sup> vocabularies).

### 5.2.2 Implementation and Evaluation

In this section, we describe the experiments and report the evaluation on detecting categories on 444 endpoints. We then describe a prototype as a “proof-of-concept” of the proposal.

#### 5.2.2.1 Implementation

A first prototype, implemented with Javascript and the Bootstrap framework<sup>10</sup>, is available at <http://semantics.eurecom.fr/datalift/rdfViz/apps/>, as a proof of concept. We aim at providing a lightweight tool for lay users to quickly understand what the data is about and so that they get first visualizations based on categories detected in the datasets. We also reuse *sgvizler* [76] for generating charts according to the categories retrieved by the wizard. In the current implementation, the user can enter any SPARQL endpoint, and with a “click”, the user can receive the list of categories detected together with sample resources. In the second step, the wizard retrieves the properties from the objects and subjects part of `owl:sameAs` links. The last step shows different tabs with the summary of the previous steps, the visualizations available for each categories, and a report both in human and machine readable formats. Figure 5.2.2.1 depicts a sample visualization generated by the wizard for geo data and statistics data.

---

<sup>9</sup><http://data.lirmm.fr/ontologies/chart>

<sup>10</sup><http://getbootstrap.com/>

The system can be used in any tool consuming Linked Data in which the complexity of SPARQL analysis and visualizations of RDF datasets is hidden to the lay users, with the benefits of showing that information encoded in triples is not only "beautiful", but also useful in the sense of traditional wizard-based tools.

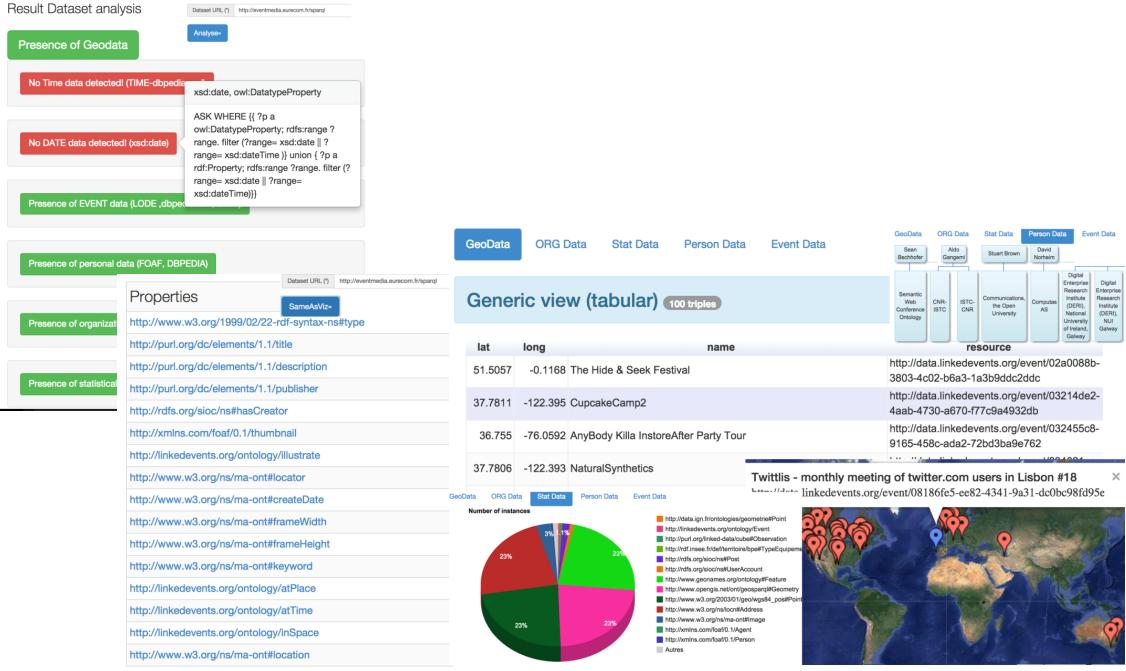


Figure 5.2: Categories detected and visualization generated by the Linked Data visualization wizard in the case of EventMedia endpoint service.

### 5.2.2.2 Experiment Set Up

We have evaluated our approach on the list of 444 endpoints referenced at <http://sparqles.okfn.org/> monitoring the availability, performance, interoperability and discoverability of SPARQL Endpoints registered in Datahub [96]. We have implemented a script in Python to speed up the process and obtain the results. Every ASK query for the different category is implemented in a separate function requesting a JSON response.

### 5.2.2.3 Evaluation

We evaluate our algorithm of detecting categories using SPARQL queries against endpoints retrieved from the LOD cloud. From the 444 endpoints used on the detection category module, 278 endpoints (62.61%) were able to give satisfactory (yes/no on one of the seven categories) results based on the queries. However, almost 37.38% of the endpoints were either down at the time of our experiments or

the response header was in XML instead of JSON (as set up in the script). This result shows that our proposal with the current implementation (not covering all the vocabularies in LOV) make use of most popular vocabularies reused in the Linked Data.

Category	number	Percentage
GEO DATA	97	21.84%
EVENT DATA	16	3.60%
TIME DATA	27	6.08%
SKOS DATA	2	0.45%
ORG DATA	48	10.81%
PERSON DATA	59	13.28%
STAT DATA	29	6.6%

Table 5.2: Classification of the endpoints according to the datatype detected with our SPARQL generic queries

This also implies a good coverage of the method that uses standard queries and yet can be extended. The full result of the detection module on the queried services is available at <http://cf.datawrapper.de/3FuiV/2/>, where for each column, the value 0 stands for *no presence* and 1 for the *presence* of the categories. As provided in Table 5.2, 21.84% of geo data was detected, 13.288% of person data, 10.81% of org data and 3.6% of SKOS data.

Endpoint	event	geo	org	person	skos	time
dbpedia.org	0	1	1	1	0	0
de.dbpedia.org	0	1	1	1	0	0
el.dbpedia.org	1	1	1	1	0	0
fr.dbpedia.org	1	1	1	1	0	1
ja.dbpedia.org	1	1	1	1	0	0
live.dbpedia.org	1	1	1	1	0	1
nl.dbpedia.org	1	1	1	1	0	0
pt.dbpedia.org	1	1	1	1	0	0

Table 5.3: Categories detected in some *dbpedia* endpoint domains, where “1” is the presence and “0” the absence of the given type of category.

Table 5.3 summarizes some findings for 8 DBpedia chapters endpoints where it’s easy to note the absence of SKOS data, and less presence of data modeled using **time** vocabulary. The Table also shows the differences in the standard vocabularies used to convert the Wikipedia data into RDF across different chapters.

## 5.3 Finding Important Properties for an Entity

In many knowledge bases, entities are described with numerous properties. However, not all properties have the same importance. Some properties are considered as keys for performing instance matching tasks while other properties are generally chosen for quickly providing a summary of the key facts attached to an entity. Our motivation is to provide a method enabling to select what properties should be used when depicting the summary of an entity, for example in a multimedia question answering system such as QakisMedia<sup>11</sup> or in a second screen application providing more information about a particular TV program<sup>12</sup>. Our approach consists in: (i) reverse engineering the Google Knowledge Panel by extracting the properties that Google considers as sufficiently important to show (Section 5.3.1), and (ii) analyzing users' preferences by conducting a user survey and comparing the results (Section 5.3.2). We finally show how we can explicitly represent this knowledge of preferred properties to attach to an entity using the Fresnel vocabulary [97].

### 5.3.1 Reverse Engineering the Google KG Panel

We aim at capturing the properties depicted in the Google Knowledge Panel (GKP) that are injected in search result pages [98] by using Web scraping. Web scraping is a technique for extracting data from traditional Web pages. We have developed a Node.js<sup>13</sup> application that queries all DBpedia concepts that have at least one instance which is `owl:sameAs` with a Freebase resource in order to increase the probability that the search engine result page (SERP) for this resource will contain a GKP. We assume in our experiments that the properties displayed for an entity are type and context dependent (country, time, query) which can affect the results. Moreover, we filter out generic concepts by excluding those who are direct subclasses of `owl:Thing` since they will trigger ambiguous queries. We obtained a list of 352 concepts<sup>14</sup>.

For each of these concepts, we retrieve  $n$  instances (in our experiment,  $n$  was equal to 100 random instances). For each of these instances, we issue a search query to Google containing the instance label. Google does not serve the GKP for all user agents and we had to mimic a browser behavior by setting the *User – Agent* to a particular browser. We use CSS selectors to check the existence of and to extract data from a GKP. An example of a query selector is `._om` (all elements with class name `_om`) which returns the property DOM element(s) for the concept described in the GKP. From our experiments, we found out that we do not always get a GKP in a SERP. If this happens, we try to disambiguate the instance by issuing a new query with the concept type attached. However, if no GKP was found

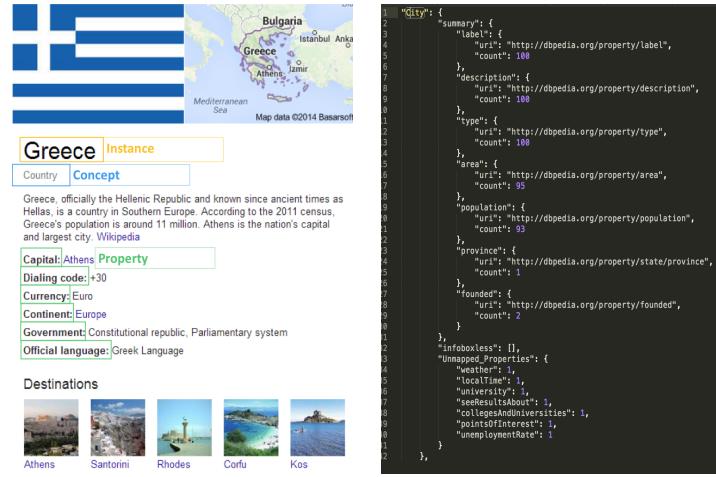
---

<sup>11</sup><http://qakis.org/>

<sup>12</sup><http://www.linkedtv.eu/demos/linkednews/>

<sup>13</sup><http://nodejs.org/>

<sup>14</sup>See also the SPARQL query at <http://goo.gl/EYuGm1>



(a) Google's Knowledge Panel for Greece (b) Results of the crawling for the class City

Figure 5.3: Google Knowledge Graph Reverse Engineering Process.

again, we capture that for manual inspection later on. Listing 1 gives the high level algorithm for extracting the GKP. The full implementation can be found at <https://github.com/ahmadassaf/KBE>. We finally observe that this experiment is only valid for the English Google.com search results since GKP varies according to top level names.

Figure 5.3 shows the GKP for Greece, and the results of crawling and mapping for the class `City`. In this particular case, each of the properties of the GKP are directly mapped with the corresponding property in DBpedia and the frequency. Hence, a `City` is likely to be described in the GKP with at least the following properties: label, description, type, and population.

### 5.3.2 Evaluation

We conducted a user survey in order to compare what users think should be the important properties to display for a particular entity and what the GKP shows.

**User survey.** We set up a survey<sup>15</sup> on February 25th, 2014 and for three weeks in order to collect the preferences of users in term of the properties they would like to be shown for a particular entity. We select only one representative entity for nine classes: `TennisPlayer`, `Museum`, `Politician`, `Company`, `Country`, `City`, `Film`, `SoccerClub` and `Book`. 152 participants have provided answers, 72% from academia, 20% coming from the industry and 8% having not declared their affiliation. 94% of the respondents have heard about the Semantic Web while 35% were not familiar

<sup>15</sup>The survey is at <http://eSrv.org?u=entityviz>

**Algorithm 1** Google Knowledge Panel reverse engineering algorithm

---

```

1: INITIALIZE equivalentClasses(DBpedia, Freebase) AS vectorClasses
2: Upload vectorClasses for querying processing
3: Set n AS number-of-instances-to-query
4: for each conceptType ∈ vectorClasses do
5:   SELECT n instances
6:   listInstances ← SELECT-SPARQL(conceptType, n)
7:   for each instance ∈ listInstances do
8:     CALL http://www.google.com/search?q=instance
9:     if knowledgePanel exists then
10:      SCRAP GOOGLE KNOWLEDGE PANEL
11:    else
12:      CALL http://www.google.com/search?q=instance+conceptType
13:      SCRAP GOOGLE KNOWLEDGE PANEL
14:    end if
15:    gkpProperties ← GetData(DOM, EXIST(GKP))
16:   end for
17:   COMPUTE occurrences for each prop ∈ gkpProperties
18: end for
19: RETURN gkpProperties

```

---

with specific visualization tools. The detailed results<sup>16</sup> show the ranking of the top properties for each entity. We only keep the properties having received at least 10% votes for comparing with the properties depicted in a KGP. We observe that users do not seem to be interested in the INSEE code identifying a French city while they expect to see the population or the points of interest of this city.

**Comparison with the Knowledge Graphs.** The results of the Google Knowledge Panel (GKP) extraction<sup>17</sup> clearly show a long tail distribution of the properties depicted by Google, with a top N properties (N being 4, 5 or 6 depending on the entity) counting for 98% of the properties shown for this type. We compare those properties with the ones revealed by the user study. Table 5.4 shows the agreement between the users and the choices made by Google in the GKP for the 9 classes. The highest agreement concerns the type **Museum** (66.97%) while the lowest one is for the **TennisPlayer** (20%) concept. We think properties for museums or books are more stable than for types such as person/agent which vary significantly. We acknowledge the fact that more than one instance should be tested in order to draw meaningful conclusion regarding what are the important properties for a type.

Classes	TennisPlayer	Museum	Politician	Company	Country	City	Film	SoccerClub	Book
Aggr.	20%	66.97%	50%	40%	60%	60%	60%	50%	60%

Table 5.4: Agreement on properties between users and the Knowledge Graph Panel

With this set of 9 concepts, we are covering 301,189 DBpedia entities that have an existence in Freebase, and for each of them, we can now empirically define the most important properties when there is an agreement between one of the biggest knowledge base (Google) and users preferences.

<sup>16</sup><https://github.com/ahmadassaf/KBE/blob/master/results/agreement-gkp-users.xls>

<sup>17</sup><https://github.com/ahmadassaf/KBE/blob/master/results/survey.json>

**Modeling the preferred properties with Fresnel.** Fresnel<sup>18</sup> is a presentation vocabulary for displaying RDF data. It specifies *what* information contained in an RDF graph should be presented with the core concept `fresnel:Lens`. We use the Fresnel and PROV-O ontologies<sup>19</sup> to explicitly represent what properties should be depicted when displaying an entity. This dataset can now be re-used as a configuration for any consuming application.

```

1 :tennisPlayerGKPDefaultLens rdf:type fresnel:Lens ;
2   fresnel:purpose fresnel:defaultLens ;
3   fresnel:classLensDomain dbpedia-owl:TennisPlayer ;
4   fresnel:group :tennisPlayerGroup ;
5   fresnel:showProperties (dbpedia-owl:abstract dbpedia-owl:birthDate
6     dbpedia-owl:birthPlace dbpprop:height dbpprop:weight
7     dbpprop:turnedpro dbpprop:siblings) ;
8   prov:wasDerivedFrom
9   <http://www.google.com/insidesearch/features/search/knowledge.html> .

```

Listing 5.7: Excerpt of a Fresnel lens in Turtle

## 5.4 GeoRDFviz: Map visualization of Geodata Endpoints

With the growing interest of publishing geolocation data according to linked data principles, many endpoints are provided without any visual interface to help navigating on top of the data. This lead to make it difficult for lay users to grab the essence of the data without learning some SPARQL queries. Visual depiction of a location in a map makes it easier to identifier a resource or a group of resources in a given bounding box. We adapted the LDVIZWiz in the geospatial datasets published in the Linked Open Data cloud. We implemented **GeoRDFviz**, which is a lightweight tool that help understanding the geodata resources in any endpoint with a more attractive way for non-experts users.

### 5.4.1 Back-End Description

GeoRDFviz is an application of the generic wizard for Linked Data presented in section 5.2. First, we take the result of the output of the analysis of endpoints, based on the geodata category as defined in listing 5.1. For each of the endpoints containing geodata, we trigger a SELECT query to take randomly one resource with latitude and longitude according to the `geo` vocabulary, as depicted in listing 5.8.

```

1 PREFIX geo:<http://www.w3.org/2003/01/geo/wgs84_pos#>
2 SELECT ?lat ?long ?name ?p
3 WHERE{
4   ?s geo:lat ?lat ; geo:long ?long.
5   ?p ?o ?s.
6   ?p rdfs:label ?name .
7   FILTER(!isblank(?p))
8 } LIMIT 1

```

<sup>18</sup><http://www.w3.org/2005/04/fresnel-info/>

<sup>19</sup><http://www.w3.org/TR/prov-o/>

---

Listing 5.8: SPARQL query runs against each endpoint containing geodata to retrieve a random location.

After we output the result of this process in a JSON file containing not only the result of the SPARQL query, but also the encoding uri for the DESCRIBE query of the resource. Listing 5.9 shows a sample of the dataset used for GeoRDFviz.

```

1  [
2  {
3      "endpoint": "http://healthdata.tw.rpi.edu/sparql",
4      "geoData": {
5          "items": [
6              {
7                  "describe": "http://healthdata.tw.rpi.edu/sparql?
8                      default-graph-uri=&query=DESCRIBE%3Chttp%3A%2F%2
9                      Fwww.w3.org%2FPeople%2FBerners-Lee%2Fcard%23i%3E
10                     &format=text%2Frdf%2Bn%3&timeout=0&debug=on",
11
12                  "lat": "42.361860",
13                  "long": "-71.091840",
14                  "name": "Tim Berners-Lee",
15                  "resource": "http://www.w3.org/People/Berners-Lee/
16
17                      card#i"
18              }
19          ]
20      }
21  },
22  {
23      "endpoint": "http://purl.org/twc/hub/sparql",
24      "geoData": {
25          "items": [
26              {
27                  "describe": "http://purl.org/twc/hub/sparql?default-
28                      graph-uri=&query=DESCRIBE%3Chttp%3A%2F%2Fdbpedia.
29                      org%2Fresource%2
30                      FAdministrative_Office_of_the_United_States_Courts%
31                      3E&format=text%2Frdf%2Bn%3&timeout=0&debug=on",
32
33                  "lat": "38.8951",
34                  "long": "-77.0367",
35                  "name": "Administrative Office of the United States
36
37                      Courts",
38                  "resource": "http://dbpedia.org/resource/
39
40                      Administrative_Office_of_the_United_States_Courts"
41              }
42          ]
43      }
44  }
45 ]
```

```

28         ]
29     }
30   },
31
32 {
33   "endpoint": "http://linguistic.linkeddata.es/sparql",
34   "geoData": [
35     "items": [
36       {
37         "describe": "http://linguistic.linkeddata.es/sparql?
38           default-graph-uri=&query=DESCRIBE%3Chttp%3A%2F%2
39             Flinguistic.linkeddata.es%2Fmlode%2Fresource%2
40               FVillage%2FBandiagara%3E&format=text%2Frdf%2Bn&
41                 timeout=0&debug=on",
42         "lat": "14.35",
43         "long": "-3.6167",
44         "name": "Bandiagara",
45         "resource": "http://linguistic.linkeddata.es/mlode/
46           resource/Village/Bandiagara"
47       }
48     ]
49   }
50 ]

```

Listing 5.9: Sample output of the JSON dataset used in the GeoRDFViz application.

#### 5.4.2 Front-End Interface

GeoRDFviz is built using generic queries in SPARQL and three visual actions: (i) zooming, (ii) filtering and (iii) describing of resources with geometry. GeoRDFviz prototype is available at <http://semantics.eurecom.fr/datalift/GeoRDFviz/>. The user first select the endpoint to visualize. GeoRDFviz picks up randomly a resource with its location contained in the endpoint and direct the user to that position in a map. From that point, the user can zoom-in to find list of resources around that area of the map, and progressively can reach a more specific resource. By clicking on the resource, a pop-up menu shows more connections of the resource modeled as SKOS concepts, obtained from the DESCRIBE query on the resource. The application is a client application implemented using Backbone.js library<sup>20</sup> and can be used to explore the geospatial LOD cloud by zooming, selecting and viewing resources.

---

<sup>20</sup><http://backbonejs.org/>

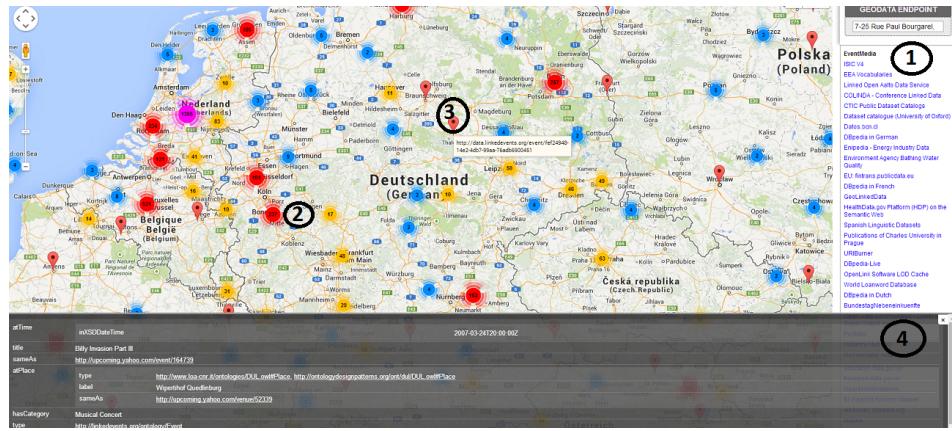


Figure 5.4: Screenshot of the user interface. The circles with numbers highlight the different elements : (1) list of endpoints, (2) number of resources available in the map area, (3) A zoom to a given element and (4) description of the selected resource.

## 5.5 An application consuming event datasets: Confomaton

In this section, our goal is to create a rich environment that enables users to navigate events and their various representative media such as photos, slides and tweets. A typical usage is to gather data about a scientific conference and investigate the added value of collecting scientific-related media. A non-trivial task in such applications is to connect structured data with extremely noisy content, especially in the case of a major conference.

### 5.5.1 Background and Motivation

We consider a scientific conference, the International Semantic Web Conference which took place in Bonn, Germany in November 2011. We estimate that it attracted more than 1,500 participants including all co-authors, people who have participated in the reviewing process, people who physically attended the conference or tried to follow it on social networks. The conference organizers publish a lot of structured data about the conference including the list of accepted papers, their authors and institutions, the detailed program composed of sub-events with the exact timetable and the location (rooms) of the talks. This data is modeled using the SWC ontology<sup>21</sup>, which was designed to describe academic events, and uses classes and properties from other ontologies such as FOAF (for people) and SWRC (BibTeX elements for the papers). The main conference (type `swc:ConferenceEvent`) is related to a set of sub-events, namely (`WorkshopEvent`, `TutorialEvent`, `SessionEvent`, `TalkEvent`

<sup>21</sup><http://data.semanticweb.org/ns/swc/ontology>

via the property `swc:isSuperEventOf`. Table 5.5 shows some statistics about the data provided by the DogFood server about the ISWC 2011 conference.

We first notice that the data is incomplete. The conference hosted 16 workshops in total, but the 75 papers are only associated with 8 of them while the 8 others did not have any papers according to the corpus. Furthermore, we find 133 papers that are not connected to any of the events via the predicate `swc:hasRelatedDocument`. Finally, some useful information is also missing such as the keynote speakers and the *Semantic Web Death Match* (panel) event. This lack of knowledge is also a motivation for our work: can we collect and analyse social network activities in order to complete the factual description of this event?

Main Event	Sub-event	Number of events	Papers	Authors
Conference Event	Workshop Event	16	75	185
	Tutorial Event	7	7	20
	Session Event	1	66	202
	Talk Event	93	93	275
	-	-	133	385
Total (distinct)		117	292	735

Table 5.5: Metadata provided by the Dog Food Server for the ISWC 2011 conference.

We collected social network data in real time during the six days of the conference using the main tags advertised by the organizers (`#iswc2011`, `#cold2011`, `#derive2011`, etc.). Table 5.6 shows some statistics about the different media services used by the attendees along with the number of items from a number of distinct users. As expected, Twitter is by far the most used service: we have been able to collect 3,390 tweets from 519 different users. A significant proportion of tweets contain hyperlinks that we have further analysed. Hence, we extracted 384 different websites indexed by so-called URL shorteners (such as Bit.ly) found in 1,464 tweets (43% of tweets). These links represent a rich source of media, as they point to various Web resources categories including blogs, slides, photos, publications and projects. For example, 25% of these links pointed to a PDF document, typically one of the conference papers but could also be related papers relevant for the followers of the conference. We also analyse those links to extract the various media services used by Twitter.

*Confomaton* is a Semantic Web application that produces and consumes Linked Data. The name *Confomaton* is a word play on the French term *Photomaton* (English photo booth) and *conference*. Just like a Photomaton illustrates the scene inside of the photo booth, the *Confomaton* illustrates an event such as a conference enriched with social media. The system is available at <http://semantics.eurecom.fr/confomaton/iswc2011>. It is composed of four main components (Figure 5.5):

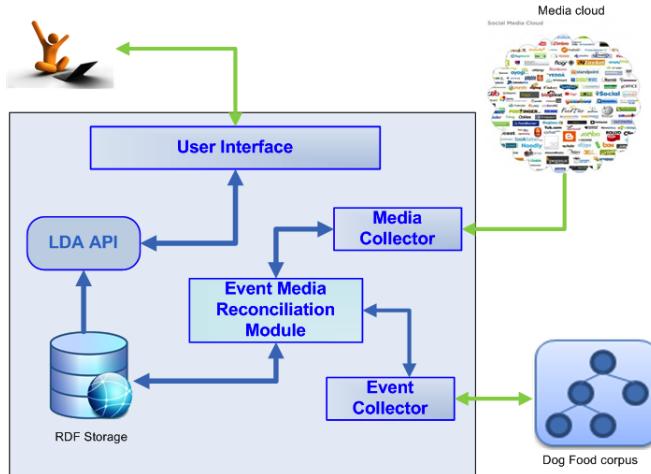
1. an Event Collector that extracts events descriptions such as the ones available in the Semantic Web Dog Food corpus;
2. a Media Collector that collects social media content and represents it in RDF

Media Service	Items	Users
Twitter	3390 tweets	519
pic.twitter	12 photos	6
yfrog	10 photos	9
Twitpic	10 photos	6
Flickr	47 photos	6
Google+	30 posts	26
Slideshare	25 slides	20

Table 5.6: Media services used during the ISWC 2011 conference

using various vocabularies;

3. a Reconciliation Module playing the role of associating social media with sub-events and external knowledge;
4. a User Interface powered by an instance of the Linked Data API as a logical layer connecting all the data in the triple store with the front-end visualizations.

Figure 5.5: *Confomaton* general architecture.

### 5.5.2 Collecting and Modeling Data in Confomation

**Media Collector: A Media Crawler** In the context of *Confomaton*, we have developed a *media collector* with the purpose of searching various social networks and media platforms for event-related media items such as photos, videos, and slides. We currently support 4 social networks (Google+, MySpace, Facebook, and Twitter) and 7 media platforms (Instagram, YouTube, Flickr, MobyPicture, img.ly, yfrog and Twitpic). Our approach is agnostic of media providers, as we offer a common alignment schema for all of them:

- **Media URI**, the deep link to the media item, e.g. `http://farm7.staticflickr.com/6059/6290784192_567346ba6a_o.jpg`
- **Type**, the type of the media item, e.g. “photo”
- **Story URI**, the URI of the micropost or story where the media item appeared, e.g. `http://www.flickr.com/photos/96628098@N00/6290784192/`
- **Message**, the concrete micropost or description text in raw format, e.g. “Laura. #iswc2011, #semanticweb, #bonn, #germany”
- **Clean**, the concrete cleaned micropost or description text with some characters (e.g. hash signs) removed, e.g. “Laura. iswc2011, semanticweb, bonn, germany”
- **User**, the URI of the author of the micropost, e.g. `http://www.flickr.com/photos/96628098@N00/`
- **Published**, the timestamp of when the micropost was authored, or the media item was uploaded, e.g. `2011-10-27T12:24:41Z`

```

1  {
2    "GooglePlus": [
3      {
4        "mediaurl": "http://software.ac.uk/sites/default/files/images/
5          content/Bonn.jpg",
6        "storyurl": "https://plus.google.com/107504842282779733854/
7          posts/6ucw1Udb5NT",
8        "message": {...}
9      },
10     "Flickr": [
11       {
12         "mediaurl": "http://farm7.staticflickr.com/6226/6290782640_e8
13           a1ffdcc2_o.jpg",
14         "storyurl": "http://www.flickr.com/photos/96628098@N00/629078
15           2640/",
16         "message": {...}
17       }
18     ]
19   }

```

Listing 5.10: Sample output of the media collector showing Google+ and Flickr results using `#iswc2011` as the query term.

In order to retrieve data from media providers, we use the particular media provider’s search Application Programming Interfaces (API) where available, and otherwise fall back to Web scraping the media provider’s website. In some cases, we initially use the search API, but then have to fall back to Web scraping in order to get more

details on the results, such as the **Media URI**, which is not exposed by all APIs. Of the media providers, Twitter plays a special role, as it can serve as a host for other media providers. For example, it is very common for tweets to contain links to media items hosted on external media providers such as Twitpic. Other media providers treat media items as first class objects, i.e. have dedicated object keys in their API results for media items, which is not in all cases true for Twitter. We handle this by searching for a list of URIs of known media providers in combination with the actual search term. To illustrate this, when searching for media items for the search term “ISWC 2012” on Twitter, we would actually search for “iswc 2012 AND (twitpic.com OR flic.kr)” in the background, whereas on all other media providers, the search term “iswc 2012” is sufficient. The media collector can be tested at <http://webmasterapp.net/social/>.

**Data Modeling of Confomaton** The Event Collector takes as input the Dog Food corpus described using the SWC ontology and converts all events into the LODE ontology<sup>22</sup>, a minimal model that encapsulates the most useful properties for describing events. We use the Room ontology<sup>23</sup> for describing the various rooms contained in the conference centre. An explicit relationship between an event and its representative media (photo, slide, micropost, etc.) is realised through the lode:illustrate property. For describing those media, we re-use two popular vocabularies: the W3C Ontology for Media Resources<sup>24</sup> for photos and videos, and SIOC<sup>25</sup> for tweets, status, posts and slides. Listing 5.11 below shows how a tweet is represented in *Confomaton*.

```

1 \begin{verbatim}
2 <http://data.linkedevents.org/tweet/af557cef-5d5b-49c6-a4c3-bc9c41ce1555>
3 a sioc:Post;
4 dcterms:created "2011-10-23T13:34:03+00:00";
5 sioc:content "@smeh Good luck for your presentation at #ssn2011
   ...";
6 sioc:hasCreator <http://www.twitter.com/BadmotorF>;
7 lode:illustrate <http://data.semanticweb.org/workshop/ssn/2011>;
8 gc:hashtag "#ssn2011";
9 owl:sameAs <http://twitter.com/BadmotorF/status/128071685235671040>.
```

Listing 5.11: Sample output describing a resource.

Figure 5.6 depicts how all these vocabularies are used together. The ISWC 2011 conference is illustrated by a photo shared on Flickr, has a sub-event, EvoDyn

<sup>22</sup><http://linkedevents.org/ontology/>

<sup>23</sup><http://vocab.deri.ie/rooms>

<sup>24</sup><http://www.w3.org/TR/mediaont-10/>

<sup>25</sup><http://rdfs.org/sioc/spec/>

2011 workshop, in which one of the tweets posted mentioned the recognised named entity Natasha Noy who is also a general chair of the conference. All the data in the *Confomaton* graph is in a public SPARQL endpoint available at <http://semantics.eurecom.fr/sparql>.

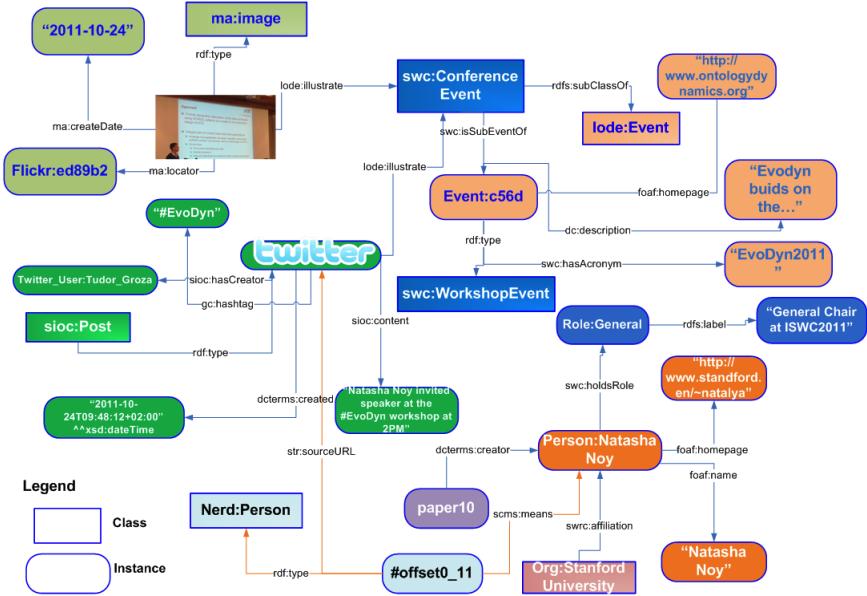


Figure 5.6: Example of data modeled in *Confomaton* re-using multiple vocabularies

**Event Media Reconciliation Module** The event media reconciliation module aims to align the incoming stream of social media with their appropriate events and to interlink some descriptions with general knowledge available in the LOD cloud<sup>26</sup> (e.g. people and institutions descriptions). Attaching social media to fine-grained events is a challenging problem. We tackle it by pre-processing the data with two successive filters in order to reduce the noise: one filter relies on keyword search applied to some fields such as title and tag, while the other one filters data based on temporal clues. The reconciliation is then ensured through a pre-configured mapping between a set of keywords and their associated events. This map enables us to associate media with the macro-events that people explicitly refer to in their posts or photos. For example, we connect all media items containing the tag #iswc2011 with the general ISWC 2011 conference. However, this method is not at all convenient for associating media items with sub-events. For instance, in the ISWC 2011 conference, there are 99 sub-events of type TalkEvent, which could be the presentation of a paper, a keynote speech or any other kind of talk. Social network users typically do not specify a particular tag for such events. We hence advocate the need for more advanced classifiers to associate media with sub-events. These classifiers can exploit a variety of parameters such as social network graphs and named entities extracted

<sup>26</sup><http://lod-cloud.net/>

from media content.

### 5.5.3 Graphical User Interface

The Graphical User Interface (GUI) of *Confomaton* is built around four perspectives characterising an event:

1. “Where does the event take place?”,
2. “What is the event about?”,
3. “When does the event take place?”, and finally
4. “Who are the attendees of the event?”.

In addition, the user interface offers full text search for these four dimensions. The *Confomaton* user interface is powered by the Linked Data API<sup>27</sup>, which provides a configurable way to access RDF data using simple RESTful URLs that are translated into queries to our SPARQL endpoint. More precisely, we use the Elda<sup>28</sup> implementation developed by Epimorphics. Elda comes with some pre-built samples and documentation, which allow to build specification to leverage the connection between the back-end (data in the triple store) and the front-end (visualizations for the user). The API layer helps to associate URIs with processing logic that first extracts data from the SPARQL endpoint using one or more SPARQL queries and then serializes the results using the format requested by the client. A URI is used to identify a single resource whose properties are to be retrieved or to identify a set of resources, either through the structure of the URI, or through the query parameters. Listing 5.12 shows an example of the configuration file of the *Confomaton* API specifying the event, media and tweet viewers followed by the events and media properties access.

```

1 <#MyAPI> a api:API ;
2   rdfs:label "Confomaton API"@en ;
3   api:endpoint <#event>,<#media>,<#tweet>,<#agent>,<#venue>,<#user>,<#eventbyid>,
4     <#mediabid>,<#tweetbyid>,<#agentbyid>,<#venuebyid>,<#userbyid>;
5   api:sparqlEndpoint <http://semantics.eurecom.fr/sparql>;
6   # specification of the event viewer (all properties appear in the json file)
7   spec:eventViewer a api:Viewer ;
8   api:property "title","description","space.lat","space.lon","time.datetime","inagent.label",...
9   <#eventbyid> a api:ItemEndpoint;
10  api:uriTemplate "/event/{id}";
11  api:itemTemplate "http://data.linkedevents.org/event/{id}";
```

Listing 5.12: Example configuration file of the *Confomaton* API, specifying event properties access.

On the left side of the main view, the user can select the main conference event or one of the sub-events as provided by the Dog Food metadata corpus. In the

<sup>27</sup><http://code.google.com/p/linked-data-api/wiki/Specification>

<sup>28</sup><http://code.google.com/p/elda>

centre, the default view is a map centered on where the event took place (e.g. Bonn, Germany) and the user is also encouraged to explore potential other types of events (concerts, exhibitions, sports, etc.) happening nearby, based on data provided by EventMedia [99]. The *What* tab is media-centred and allows to quickly see what illustrates a selected event (tweets, photos, slides). Zooming in an event triggers a popup window that contains the title and timetable of the event, the precise room location and a slideshow gallery of all the media items collected for this event. For the *When* tab, a timeline is provided in order to filter events according to a day time period. Finally, the *Who* tab aims to show all the participants of the conference. This is intrinsically bound to a social component, aiming not only to present relevant information about participants (their affiliations, homepages, or roles at the conference), but also the relationships between participants themselves and with events.

## 5.6 Application consuming statistical datasets

In this application, we describe the creation of an application consuming statistical dataset in the French context. The vocabulary used to model the data is compatible with the Data Cube vocabulary [21], useful for retrieving statistical data. The **Perfect School** application is intended to provide useful information on schools in France using semantic technologies, with RDF-ized data enriched with other datasets in the wild. The application and the vocabulary have successfully passed the integrity checker<sup>29</sup> of an implementation for the candidate recommendation of Data Cube vocabulary [21], a recommendation from the W3C.

### 5.6.1 Dataset Modeling

**Legacy Datasets** In order to build the application, we had to look at some relevant datasets in <http://data.gouv.fr>. The ones selected for building the application are the following:

- The file at <http://www.data.gouv.fr/DataSet/564055> in CSV format, containing a list of 67,201 schools (name, status, type), with geolocation positions in Lambert 93, for the academic year 2011-2012. The file contains the following attributes (in French):
  - code of the school (e.g. 0010002X),
  - official name of the school (e.g.: College Saint-Exupery),
  - principal name (e.g. COLLEGE),
  - patronymic name (e.g. Saint-Exupery),
  - status of the school (1 = open, 2 = to be closed and 3 = to be opened),
  - label of the type of school (e.g. 1= first degree, 3 = second degree).

---

<sup>29</sup>[http://www.w3.org/2011/gld/wiki/Data\\_Cube\\_Implementations](http://www.w3.org/2011/gld/wiki/Data_Cube_Implementations)

- A dataset at <http://www.data.gouv.fr/DataSet/572165> in .CSV format, giving results for professional schools, indicators), from one academic year (2011-2012). The file contains the following attributes:
  - name of school,
  - city code,
  - code of the school,
  - district where the school is located,
  - sector (PR= PRivate, PU= PUblic),
  - several observation measures with statistics on success rates, school versus academic/national rates, name of the academy it belongs to, as well as the department.
- A dataset at <http://www.data.gouv.fr/DataSet/572162> (.CSV) containing statistics for 2296 public high schools and indicators. It complements the statistics from INSEE.

**Ontology Modeling** We have reused some external ontologies for more interoperability:

- **aiiso**<sup>30</sup> for the type of school and codes of school.
- **geofla**<sup>31</sup> since the schools are considered as a topographic entities,
- **geom**<sup>32</sup> for representing the different geometries (points with latitude and longitude) in a given coordinate reference systems with **ignf** ontology at <http://data.ign.fr/def/ignf#>.
- **skos:Concept** for describing the 30 types of nature of schools.
- **qb:DimensionProperty** and **qb:MeasureProperty** [21] for modeling the dimensions and different indicators available for a given school.
- **geom**<sup>33</sup> for representing the different geometries (points with latitude and longitude) in a given coordinate reference systems with **ignf** ontology at <http://data.ign.fr/def/ignf#>.
- **skos:Concept** for describing the 30 types of nature of schools.
- **qb:DimensionProperty** and **qb:MeasureProperty** [21] for modeling the dimensions and different indicators available for a given school.

---

<sup>30</sup><http://vocab.org/aiiso/schema>

<sup>31</sup><http://data.ign.fr/def/geofla#>

<sup>32</sup><http://data.ign.fr/def/geometrie#>

<sup>33</sup><http://data.ign.fr/def/geometrie#>

The resulting vocabulary is available at <http://purl.org/ontology/dvia/eco>. We use the Datalift platform for transforming the different CSV files into RDF. The final data is available at <http://eventmedia.eurecom.fr/sparql> with the named graph <http://data.eurecom.fr/school>.

**URI Policies** We use the following patterns for the URIs of the vocabularies or the resources in our namespace.

- Vocabulary: <http://data.eurecom.fr/ontologies/{SECTOR}>.  
e.g: <<http://data.eurecom.fr/ontologies/eco>> for the ontology that we have developed.
- Resources: <http://data.eurecom.fr/id/{SECTOR}/{CLASS}>.  
e.g: <<http://data.eurecom.fr/id/school>> for the schools URIs,  
<<http://data.eurecom.fr/id/school/slice>> for qb:Slices.
- Taxonomies: we use SKOS for modeling concepts and codes as follows:  
<http://data.eurecom.fr/codes/{SECTOR}/{CONCEPT-TYPE}/{CODE}>. e.g. <<http://data.eurecom.fr/codes/eco/natureUAI>> for the collection of natures of the schools ;  
<<http://data.eurecom.fr/codes/eco/natureUAI/101>> for a particular concept with code 101 and label “École maternelle”.

Besides the aforementioned URI policy, each school in the user interface can be reached directly in the UI by using directly the following URI: <http://semantics.eurecom.fr/datalift/PerfectSchool/#school/{SCHOOL-CODE}/>, with {SCHOOL-CODE/} in lowercase.

*Example:* The school “Albert Camus” in the city “Le Mans” with the code school 0720800D can be viewed in the application directly at <http://semantics.eurecom.fr/datalift/PerfectSchool/#school/0720800d/>

#### Sample School Data in RDF

```

1
2 school:0750676c a aiiso:School, geofla:EntiteTopographique .
3 school:0750676c a ecole:Etablissement .
4 school:0750676c rdfs:label "LYCEE DORIAN (PROFESSIONNEL)"@fr ;
5 school:0750676c dcterms:title "Lycee polyvalent et lycee des metiers de la .."@fr ;
6   aiiso:code "0750676C" ;
7   ecole:denominationPrincipale "LPO LYCEE DES METIERS"@fr ;
8   ecole:patronyme "DORIAN"@fr ;
9   ecole:ville "PARIS 11"@fr ;
10  ecole:codeCommune "75111" ;
11  ecole:secteur "PU" ;
12  ecole:academie "PARIS"@fr ;
13  ecole:departement "PARIS"@fr ;
14  ecole:cycle "3"^^xsd:int ;
15  ecole:etat "1"^^xsd:int ;
16  ecole:nature bpenat:306 .
17
18 slice:0750676c ecole:etablissement school:0750676c .

```

```

19 school:0750676c geom:geometrie _:vb42480647 , _:vb42531960 .
20
21 _:vb42480647 a geom:Point;
22 geom:systCoord ignfr:wgs34g;
23 geom:coordX "48.85429801"^^xsd:double;
24 geom:coordY "2.39231163"^^xsd:double .
25
26 _:vb42531960 a geom:Point;
27 geom:systCoord ignfr:ntflamb35e ;
28 geom:coordX "655410.1"^^xsd:double;
29 geom:coordY "6861755.9"^^xsd:double .

```

Listing 5.13: Snapshot in Turtle for the school ID=0750676C, also at <http://semantics.eurecom.fr/datalift/PerfectSchool/#school/0750676c/>

### 5.6.2 Interconnection

For the interconnection process, we didn't use the current module. Instead of that we have used the SILK [43] platform as it is re-packaged in the workflow of Datalift. We believe the scripts used for this task can be easily reused within the Datalift platform. Two datasets were used for finding `owl:sameAs` links:

1. DBpedia French chapter<sup>34</sup>, as the scope of the application was limited to France. We have found only 7 match links with our schools datasets.
2. LinkedGeoData<sup>35</sup>, as the underlying data used comes from the community project Open Street Map (OSM). Here, we got a total of 601 matching links in the category of `lgdo:BuildingSchool`<sup>36</sup>.

### 5.6.3 User Interface

The target device for the application is Mobile phone, using principally two frameworks: JQuery mobile<sup>37</sup> and backbone JavaScript<sup>38</sup>. The application provides geolocation, search by city/district, graph charts for stats, table views of relevant results aggregated or group by some other aspects. **Perfect School Application** provides 3 main views:

1. **Search form:** The interface retrieves the location automatically ( à la Google Maps API fashion), and offers choices based on the School type: first degree / second degree. When choosing first degree, the user can further select one of primary school, elementary school or other. For the second degree, apart from looking for one of college, high school etc., the user can look for public or private schools. The search button launches the query behind the scenes for retrieving the collection of data matching the user's criteria.

<sup>34</sup><http://fr.dbpedia.org/sparql>

<sup>35</sup><http://linkedgeodata.org/sparql>

<sup>36</sup><http://linkedgeodata.org/ontology/BuildingSchool>

<sup>37</sup><http://jquerymobile.com/>

<sup>38</sup><http://backbonejs.org/>

2. **Results of searching:** The search action returns a collection of schools plotted in a map. A cursor on the left side helps users zoom to get more details about schools retrieved in a given region, or street. When selecting a given school, the name is displayed and with the possibility to see the route from the centroid of the result on the map.
3. **Description of the school:** This panel is divided in 3 different tabs: (a) General information (name, cycle, principal denomination, nature and patronym used); (b) Stats with all the different statistics in form of charts and graphs comparing the school with the others and (c) information DBpedia-FR<sup>39</sup> if available, obtained with the `owl:sameAs` links for enriching the original dataset with information such as founder, date of creation, web site, population, head of school etc.

Figure 5.7 shows the three different views on a running example when using the application in a mobile phone.

## 5.7 Improving the discovery of applications contests in Open Data Events

### 5.7.1 Background

One of the challenges in the Apps for Europe project and many open data projects in general is the discovery of existing applications using the open data. The discovery of existing applications and ideas is important in order to prevent people from reinventing the wheel, when they instead should put their effort in refining existing applications or developing completely new applications. This will hopefully lead to more diverse applications and also to higher quality applications, because existing applications can be enhanced by other people and organizations [100].

In order to promote the reuse and discovery of open data applications, Apps for Europe has created an RDF vocabulary<sup>40</sup> that can be used for describing open data events and also the applications that have been built on top of the open data. Furthermore, Apps for Europe has also developed a Wordpress plugin<sup>41</sup> that open data event organizers can install on their web pages, to make it simpler to populate the RDF data for the events and applications.

### 5.7.2 Modeling Approaches

**Modeling Events and Applications in RDF** One of the key factor in improving the discovery of open data events is to present the information in a structured machine readable format, so that applications (such as crawlers) can easily consume

---

<sup>39</sup><http://fr.dbpedia.org>

<sup>40</sup><https://github.com/mmlab/apps4eu-vocabulary>

<sup>41</sup><https://github.com/mmlab/AppsForX>

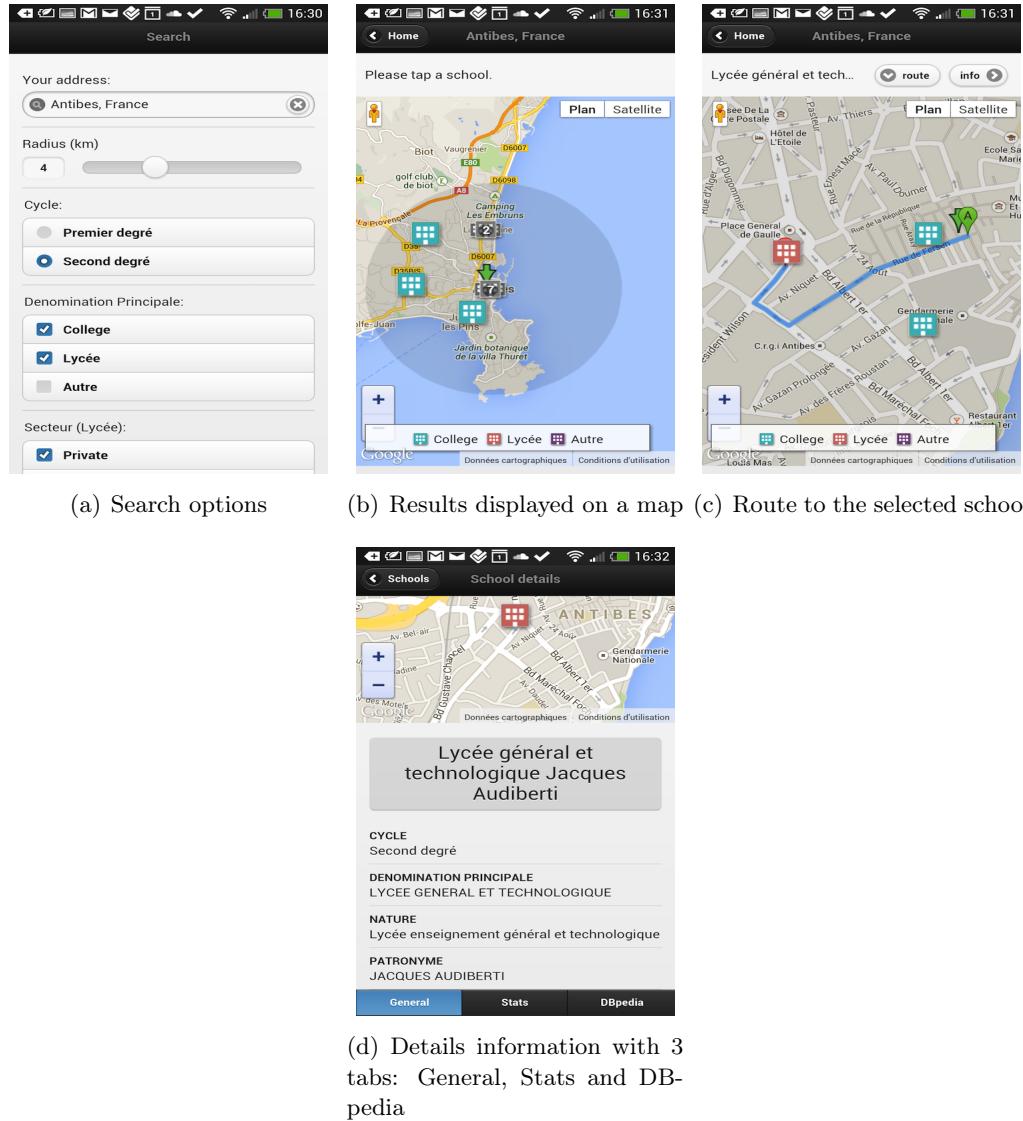


Figure 5.7: Steps for searching high schools in Antibes, France in a radius of 4000 meters.

the data. This follows the Semantic Web principles where the goal is to convert the current, unstructured documents (readable by humans) into structured documents. In order to present the event and application information in structured format, Apps for Europe has created an RDF vocabulary available at <https://github.com/mmlab/apps4eu-vocabulary> for modeling the events and applications. The quality of this vocabulary was tested by manually modeling past events and applications and looking for potential problems and improvements in the ontology. The problems and improvements made to the ontology is presented below.

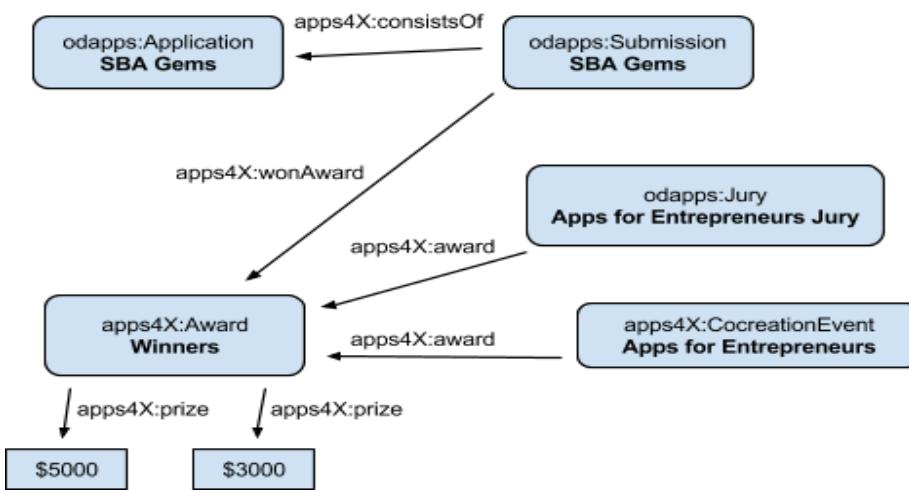


Figure 5.8: The RDF triples before changes. Here we state that the application SBA Gems has won an award in the event Apps for Entrepreneurs.

In 5.8, the model is complex, because it requires us to specify an extra class (`:Submission`) in order to use the `:wonAward`-property, and yet we are unable to specify which prize the application won (\$5000 or \$3000). In Section 5.9, the model is straightforward to specify that an application (SBA Gems) won a prize, because we don't need to use the `Submission` intermediate class. We are also able to specify which prize the application won (`:FirstPrize`). Note that now the jury is connected to a prize by specifying a `jury-attribute` for prizes (instead of having a `prize-attribute` for the jury i.e. direction of the arrow changed).

**Improving the model for specifying winners** Open data events often include competitions where the best applications are rewarded, and this information was modeled in the Apps for Europe RDF vocabulary as well. However, we found that the current model for this was unnecessarily complicated and therefore the model was simplified. Previously, winning apps required a separate “`Submission`”-class in order to state that they had won a competition. Furthermore, the vocabulary didn't allow stating which position (winner, second, third etc.) the winning application was awarded. These problems were fixed by introducing a new property “`wonPrize`” for

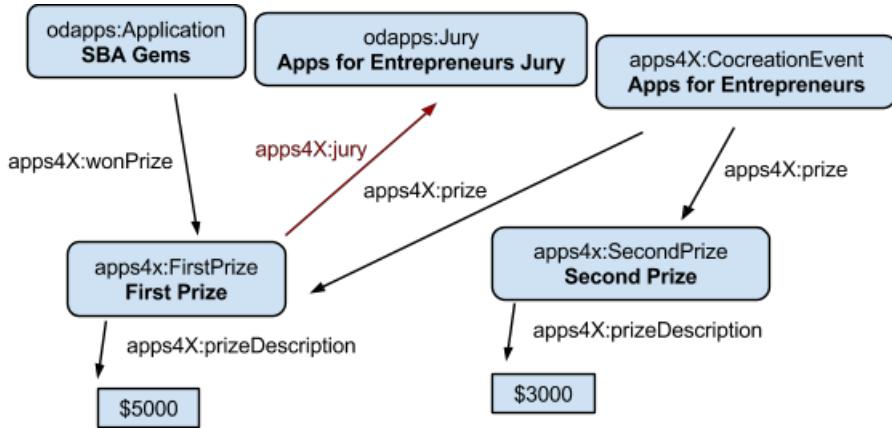


Figure 5.9: The RDF triples after changes.

applications, so that we can directly state that an application won a prize without needing to have an intermediary “Submission”-class, and by introducing the new classes “FirstPrize”, “SecondPrize”, and “ThirdPrize” that are all subclasses of the “Prize”-class. We concluded that it is satisfactory to be able to state positions one to three for the winning apps, and therefore we didn’t implement a more complicated vocabulary that would have enabled to state an arbitrary position (by having a position attribute for the prize for example).

Also, we add some other minor changes to the Apps4X vocabulary:

- Use the word “Prize” instead of “Award” (i.e. rename in all places)
- Normalize the values for juryRate and usersRate attributes (because otherwise they can’t be compared to each other). Alternatively reuse the review vocabulary from vocab.org<sup>42</sup> or schema.org<sup>43</sup>
- Introduce properties `connectedApp` (for Event) and `connectedEvent` (for Application), so that applications and events can be linked together.

### 5.7.3 Implementation and Application

In this section, we describe the implementation of the plugin based on the vocabulary described in the previous section, the RESTful API and the evaluation on creating a knowledge-base in RDF of 112 past open data events in Europe.

**Universal JavaScript plugin for RDF population** The inspiration for the JavaScript-based solution came from services like <https://muut.com/>, which offers embeddable commenting and discussion forums and it can be installed on any web page independent of the CMS. Since almost all people have JavaScript support

<sup>42</sup><http://vocab.org/review/terms.html>

<sup>43</sup><http://schema.org/Review>

enabled in their browsers<sup>44</sup> it enables us to implement a solution that can be installed on any web page and also works on almost any browser.

The most striking difference between the CMS-based plugin and the JavaScript plugin is the user base, since the JavaScript solution works on any Web page while the CMS plugin obviously only can be installed on a dedicated websites. Also the technical approach differs significantly, because in the CMS-based approach the server will do all the work for producing the HTML and RDFa markup for the events, while in the JavaScript approach it is the client browser that is responsible for this task.

**Technical description:** The JavaScript plugin consists of three distinct parts (Figure 5.10) that are technically independent of each other, but they can still be deployed on the same server if needed:

**RESTful API** for creating, editing and removing events and applications from the database. Both interfaces listed below (2 and 3) communicates and manipulates the data using the RESTful API.

**Admin interface** for event organizers. Event organizers manage their events and applications through this interface, which is provided as a “software as a service” (SaaS<sup>45</sup>).

**Embeddable script** that displays the event and application information both in human readable form and in computer readable form (RDFa format) on the event organizers’ Web page. The script fetches the event and application information using the RESTful API and then manipulates the document object model (DOM ) after page load and inject the event and application description into the DOM.

**RESTful API:** The API on server side was implemented in Node.js<sup>46</sup>, which is a platform for running applications written in JavaScript. It is built on Google Chrome’s JavaScript runtime and provides very good performance thanks to its asynchronous data manipulation model. Node.js has become very popular recently and there are huge number of third party extensions and frameworks written for it. Furthermore, since the programming language is JavaScript, many libraries and frameworks written for browsers will also function in Node.js. For more details on how to configure the plugin, see Appendix A.

The data is stored in MongoDB, which is the leading NoSQL database<sup>47</sup>. In addition, a Mongoose<sup>48</sup> object data mapper was used in order to simplify the data

---

<sup>44</sup>According to a study by Yahoo only about 1% have JavaScript disabled in their browsers. Source: <https://developer.yahoo.com/blogs/ydnfourblog/many-users-javascript-disabled-14121.html> [Referenced 2014-06-02]

<sup>45</sup>[http://en.wikipedia.org/wiki/Software\\_as\\_a\\_service](http://en.wikipedia.org/wiki/Software_as_a_service)

<sup>46</sup><http://nodejs.org/>

<sup>47</sup><http://www.mongodb.com/leading-nosql-database>

<sup>48</sup><http://mongoosejs.com/index.html>

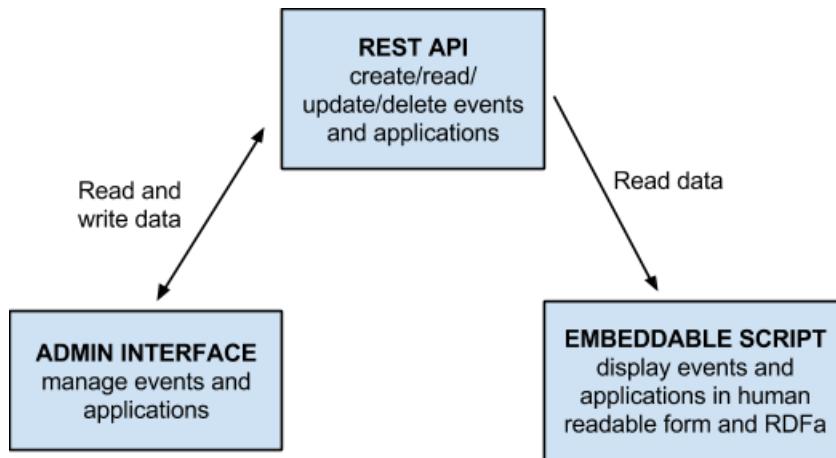


Figure 5.10: Universal JavaScript components

manipulation and to give structure to the database. Finally, node-restful was used to transform the Mongoose schema into a working REST API. Node-RESTful internally relies on Express<sup>49</sup>, which is the most popular web application framework for Node.js.

**Admin Interface:** The admin interface provides an intuitive and easy-to-use interface for event organizers to manage their events and the related applications. It is provided as software as a service and can be hosted and managed on server side. Event organizers will register to the service and after this they have access to all the features of the admin interface. Two screenshots of the service is presented in Figure 5.11.

The admin interface is implemented in HTML, CSS and JavaScript. Bootstrap<sup>50</sup> was used as a CSS framework. It simplifies the development of responsive websites and also provides clean aesthetics for the website. AngularJS<sup>51</sup> was used as a JavaScript framework.

**Embeddable script** The final component of the puzzle is the embeddable script, that is responsible of displaying the event and application information in human readable form on the event organizer's website. The event organizer will get the embed code from the admin interface, which she will then insert into the HTML code of her website. The embed code is just a reference to an external JavaScript file that will be loaded from the Apps for Europe server and then executed in the browser. A detailed description of how the script works is given below:

<sup>49</sup> [urlhttp://expressjs.com/](http://expressjs.com/)

<sup>50</sup> <http://getbootstrap.com/>

<sup>51</sup> <https://angularjs.org/>

Title	Date	Related Apps
Accessible App competition	2013-08-30 – 2013-08-30	0 apps
Developing Solutions day	2011-11-25 – 2011-11-25	2 apps
Appening	2011-03-18 – 2011-03-19	2 apps
HSL mobile competition - Developer Forum 16th March 2011	2011-03-16 – 2011-03-16	20 apps
Open concept phase of the competition Public Data In Play - join us!	2009-10-05 – 2009-11-15	6 apps

Figure 5.11: Screenshots of the admin interface. On the left is the event listing page and on the right is the form for creating an event.

- 
1. The event organizer will embed the code on his own web page.
  2. On page load, the <script>-tag is parsed by the browser and the actual script is loaded from the Apps for Europe server.
  3. The script is executed.
    - (a). The script will fetch the event or application information using an AJAX-request.
    - (b). After the server has responded with the event/application information, the information is displayed on the page in human readable form by modifying the DOM of the page. Furthermore, the same information is also presented in computer readable format by embedding RDFa in the DOM.
- 

**Creating the Knowledge-base for Past Events** The final task was to create the knowledge base for past events by using the Apps for Europe RDF vocabulary and generating data to feed the endpoint at <http://apps4europe.eurecom.fr/sparql>. The list of events was extracted from the Apps for Europe Google spreadsheet<sup>52</sup>.

The goal was to populate the events and applications, so that they would follow and fulfill the Apps for Europe vocabulary as well as possible. However, the format and information provided on the event web pages varied, and therefore we were usually only able to populate the following information for most of the events:

- Title
- Dates (start and end date)
- Description (free text)

<sup>52</sup><https://docs.google.com/spreadsheet/ccc?key=0AiXRLGASq8I0dDlfZURkWGpS0DBJaWotQUp3eGNwNGc&usp=sharing>

- Prizes (if the event included an application contest)
- Jury members (if the event included an application contest)
- Location (usually only country was known)

Similarly, the information provided for applications varied even more, and sometimes we were able only to populate the title, description and the homepage of the application. Because of these problems, the data populated into the triple store doesn't include all the information we would have hoped to have.

The original idea was to automate the triple store data population as far as possible using content scraping , but because of the heterogeneity of the web pages (in terms of both structure and data) a semi-automated approach was used. Two particular problems prevented a fully automatic approach from being implemented:

1. First, because of the heterogeneity in the structure of the web pages, it is for example extremely hard to "automatically" know which part of the page contains the event location or application homepage link.
2. Second, many links to the event pages were broken, and thus manual work was required in order to relocate the actual page.

In order to be as efficient as possible in the data population a semi-automatic approach was used. Here, the key idea was to populate event pages with only a few application entries (less than 10) manually using the JavaScript plugin, and for the rest of the event pages web scraping was used to populate the application entries (the scraping script had to be configured for each event website separately). More details of the script can be found in Appendix A.

The goal was to populate all the past events (112 in total). However, many of the websites for the past events had already disappeared or didn't include enough information for data population. Furthermore, the data population process turned out to take much more time than anticipated, and therefore we managed to populate only 28 events and in total 889 application entries. The average number of application entries per event was 34 and the median 22.

Visualizations and applications could be built on top of the knowledge base, but because of the small size of the current knowledge base, it is difficult to extract quantitative data from the knowledge base. Some ideas for visualizations are presented below:

- Map visualization of where past events have been organized
- Most popular categories/themes for applications
- Gallery of application screenshots (visual inspiration for developers).

The dataset is available below in different formats as dump for download in RDF/-Turtle and MongoDB:

- RDF in Turtle format:  
<https://www.dropbox.com/s/3075qsblxzau2fk/rdfInTurtleFormat.tar.gz>
- MongoDB dump: <https://www.dropbox.com/s/m2sr4na12v3yk07/apps4europe.tar.gz>

#### 5.7.4 Discussion

The embeddable script is non-optimal in the sense that the event or application information is loaded only after the actual page (where the script is embedded) has loaded. This causes some additional delay before the page is fully rendered, but the problem can be alleviated by showing loading indicators.

CSS styling for the events and applications is provided using Bootstrap. All CSS rules have been made specific to the container div-element of the plugin, in order to prevent the CSS from conflicting with the CSS of the page. In the future, another solution called shadow DOM could be used to prevent conflicts between different components of the page, but at the moment the browser support for shadow DOM<sup>53</sup> is not satisfactory. Since the event and application information is directly embedded on the page using DOM, the event organizer can add his own CSS styling to the plugin by overriding the provided CSS rules.

The RDF information about events and applications is populated only if event organizers actually use it on their websites, and therefore it is crucial to ensure that the plugin is appealing in the eyes of the event organizers. It would be especially useful to study the usability of the plugin and also discuss with the event organizers how well the Apps for Europe ontology fits with their data model. In addition to this, also the security of the JavaScript-plugin should be improved, although the information entered to the plugin is not very critical from a security perspective.

## 5.8 Summary

In this chapter we have presented an approach for creating visualizations on top of Linked Data based on Semantic Web technologies. We first defined seven categories of objects worth viewing in a dataset, and we propose to associate them with commonly used and domain vocabularies. We then present a description of the main components of a Linked Data Visualization Wizard. We describe a lightweight implementation in JavaScript as a *proof-of-concept* of our proposal, with the benefits to be usable on-line or being extensible. We advocate that such a tool can be easily integrated in any workflow/framework for publishing and linking data on the Web, such as Datalift or the GeoKnow Stack. Besides, we have performed experiments on GKP to look for important properties in entities, and evaluated against users' preferences. Then we presented two applications in the domain of statistics and events, consuming different datasets in RDF on real-world scenario. We discussed on how to improve applications developed for contests, by proposing a vocabulary

---

<sup>53</sup><http://caniuse.com/shadowdom>

and a tool for populating the model by using a universal plugin. Some past events have been already semi-automatically curated using both the vocabulary and the plugin.



## Part III

# Contribution to Standards



## CHAPTER 6

# Contributions to Linked Open Vocabularies (LOV)

---

*“Cool URIs don’t change”*

Tim Berners-Lee

## Introduction

In recent years, governments worldwide have mandated publication of open government content to the public Web for the purpose of facilitating open societies and to support governmental accountability and transparency initiatives. In order to realize the goals of open government initiatives, the W3C Government Linked Data Working Group have provided some guidance to aid in the access and re-use of open government data. Since Linked Data can provide a simple mechanism for combining data from multiple sources across the Web, it also addresses many objectives of open government transparency initiatives through the use of international Web standards for the publication, dissemination and reuse of structured data [3].

To publish data on the Web, stakeholders have to follow different tasks sometimes included in different life-cycle models. As described in [3], different proposals for life-cycles all share common activities, summarized in the need to specify, model and publish data in standard open Web formats. Below are details about four different life cycles.

- Hyland et al. [101] provide a six-step “cookbook” to model, create, publish, and announce government linked data. They highlight the role of the World Wide Web Consortium (W3C) which is currently driving specifications and best practices for the publication of governmental data. Hyland et al. lifecycle consists of the following activities: (1) Identify, (2) Model, (3) Name, (4) Describe, (5) Convert, (6) Publish, and (7) Maintain.
- According to Hausenblas et al. [102] existing data management approaches assume control over schema, data and data generation, which is not the case in the Web because it is an open, de-centralized environment. Based on several years of experience in Linked Data publishing and consumption over the past years, they identify involved parties and fundamental phases, which provide for a multitude of so called Linked Data life cycles that consist of the following steps: (1) data awareness, (2) modeling, (3) publishing, (4) discovery, (5) integration, and (6) use cases.

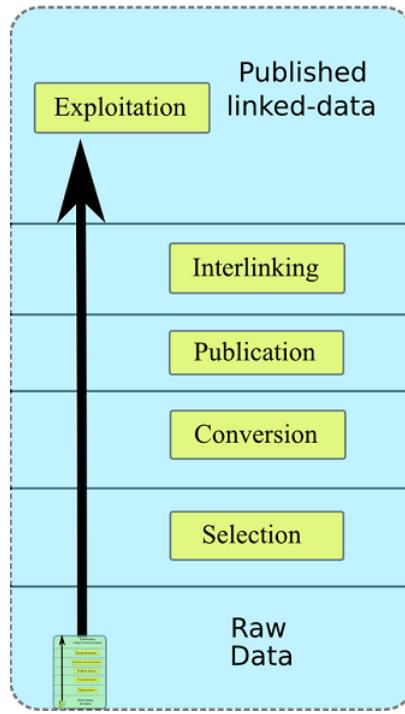


Figure 6.1: Datalift life cycle for publishing Linked Data.

- Villazón-Terrazas et al. propose in [103] a first step to formalize their experience gained from developing government Linked Data into a preliminary set of methodological guidelines for generating, publishing and exploiting Linked Government Data. Their life cycle consists of the following activities: (1) Specify, (2) Model, (3) Generate, (4) Publish, and (5) Exploit.
- In the Datalift vision [64], the process is divided into 3 principal phases: modeling the data, publishing and exploitation. Figure 6 depicts the different steps which are similar to the architecture implemented for the platform:
  1. Modeling the data consists of: (1) supporting the data selection, (2) identifying the relevant vocabulary, (3) defining a schema pattern for the URIs and (4) converting between formats
  2. Publishing the dataset consists of: (1) interconnecting with external datasets, (2) attaching provenance metadata information, (3) managing access rights to the dataset and (4) storing the data in a triple store
  3. Exploitation, such as visualizing the dataset, re-publishing and versioning.

Although the process of lifting raw data to interlinked data in Datalift seems to be linear, it can also be cyclic with a maintenance task in the life cycle. Table 6.1 summarizes the ten best practices to be taken into account to publish data as Linked

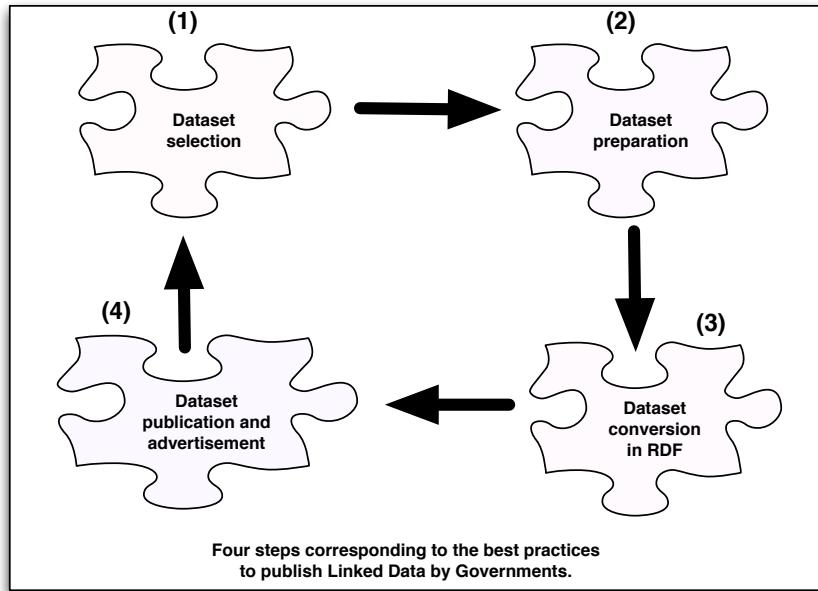


Figure 6.2: Four steps to follow for publishing Linked Data by Governments and Institutions.

Data on the Web. Figure 6 depicts the four main steps to consider when publishing Linked Data according to Table 6.1.

In this chapter, we mainly present our contributions around Linked Open Vocabularies (LOV). Section 6.2 describes the scope and the how LOV can be used with an ontology methodology (e.g., the NeOn methodology) to improve the search and the quality of the reused vocabularies. In section 6.2.3, we propose and implement an heuristic to align vocabularies on the Web of Data and a ranking of vocabularies based on information content (IC) metrics. Then, we describe a module of Datalift that reuses the LOV catalogue in the process of converting raw dataset for reusing terms already defined in other vocabularies (Section 6.4). Finally we conclude the chapter in Section 6.5.

## 6.1 Catalog of Vocabularies

A vocabulary is a collection of “terms” with more or less complex semantics used to model a domain [20]. Vocabularies can range from simple such as the widely used RDF Schema, FOAF and Dublin Core Metadata Element Set to complex vocabularies with thousands of terms, such as those used in healthcare to describe symptoms, diseases and treatments. Vocabularies play a very important role in Linked Data, specifically to help with data integration. Its can also overlap with Ontology in the context of Linked Data. Regarding the catalog of vocabularies, we refer the reader to [4] for a systematic survey of ontology libraries, but we give our own classification

of ontology repositories (Table 6.2). In particular, we distinguish six categories of catalogs:

1. *Catalogs of generic vocabularies/schemas* similar to the LOV catalog, but without any relations among the vocabularies. Example of catalogs falling in this category are vocab.org<sup>1</sup>, ontologi.es<sup>2</sup>, JoinUp Semantic Assets or the Open Metadata Registry.
2. *Catalogs of ontologies for a specific domain* such as biomedicine with the Bio-Portal<sup>3</sup>, geospatial ontologies with SOCoP+OOR<sup>4</sup>, Marine Metadata Interoperability and the SWEET ontologies<sup>5</sup>.
3. *Catalogs of ontologies from a project* such as the famous DAML repository of ontologies<sup>6</sup>.
4. *Catalogs of ontology Design Patterns (ODP)* focused on reusable patterns in ontology engineering.
5. *Catalogs of editors' ontologies* used to test some features of a tool and to keep track of the ontologies built by a tool, such as Web Protégé or TONES.
6. *Catalogs of ontologies maintained by a single organization* which often uses a platform such as Neologism<sup>7</sup> for publishing vocabularies.
7. *Vocabularies crawled by Semantic Web search engines* containing snapshots at the time of the crawls such as Watson<sup>8</sup>, Sindice<sup>9</sup>, Falcon-s<sup>10</sup> or Swoogle. For example, the NanJing Vocabulary Repository (NJVR) - a dump of Falcon-s ontologies, reported as of June 17th, 2014 2,996 vocabularies crawled from 261 pay-level domains.

We observe that the existing catalogs of vocabularies in the literature have some limitations compared with LOV. In terms of coverage, the number of vocabularies indexed by LOV is constantly growing and it is the only catalog, to the best of our knowledge, that provides all types of search criteria (metadata search, within/across ontologies search), both an API and a SPARQL endpoint access and that can be simultaneously classified as an “Application platform”, an ontology directory, and an ontology registry. Using the categories of ontology libraries defined in [4], LOV falls under two categories: “curated ontology directory” and “application platform”

---

<sup>1</sup><http://vocab.org/>

<sup>2</sup><http://ontologi.es/>

<sup>3</sup><http://bioportal.bioontology.org/>

<sup>4</sup><http://socop.oor.net/>

<sup>5</sup><http://sweet.jpl.nasa.gov/2.1/>

<sup>6</sup><http://daml.org/ontologies/>

<sup>7</sup><http://neologism.deri.ie>

<sup>8</sup><http://watson.kmi.open.ac.uk/>

<sup>9</sup><http://www.sindice.com>

<sup>10</sup><http://ws.nju.edu.cn/falcons/>

because the ontologies are curated manually with statistics automatically generated, and because it exposes its data via an API. Furthermore, LOV provides an answer to some of the issues mentioned in the survey reported in [4], such as “where has an ontology been used before?” or “is this ontology compatible with mine?”. In particular, LOV provides vocabulary usage statistics from the LOD Cloud datasets and it exposes vocabulary dependencies using the Vocabulary-of-A-Friend (**VOAF**) ontology.

`vocab.cc`<sup>11</sup> is a service which is similar to `prefix.cc` since it enables to look up and search for Linked Data vocabularies while providing more specific information about the usage of a particular class or property in the Billion Triple Challenge Dataset (BTCD). It also provides the ranking of those properties or classes. The authors mentioned that “common prefixes are resolved with data from `prefix.cc`”. Although they don’t give further details, this service is related to `prefix.cc`. `Triple-Checker`<sup>12</sup> is a web service based on `prefix.cc` which aims at finding typos and common errors in RDF data. It parses a given URI/URL and the output is divided into two sections: the namespaces and the term section. The former matches against `prefix.cc` to determine whether they are “common prefixes” and the latter provides the term definition.

---

<sup>11</sup><http://vocab.cc>

<sup>12</sup><https://github.com/cgutteridge/TripleChecker>

Step	Name	Description
STEP #1	PREPARE STAKEHOLDERS	Prepare stakeholders by explaining the process of creating and maintaining Linked Open Data
STEP #2	SELECT A DATASET	Select a dataset that provides benefit to others for reuse.
STEP #3	MODEL THE DATA	Modeling Linked Data involves representing data objects and how they are related in an application-independent way
STEP #4	SPECIFY AN APPROPRIATE LICENSE	Specify an appropriate open data license with a clear statement about the origin, ownership and terms related to the use of the published data.
STEP #5	GOOD URIs FOR LINKED DATA	Consider a good URI naming strategy and implementation plan, based on HTTP URIs. Consideration for naming objects, multilingual support, data change over time and persistence strategy are the building blocks for useful Linked Data.
STEP #6	USE STANDARD VOCABULARIES	Describe objects with previously defined vocabularies whenever possible. Extend standard vocabularies where necessary, and create vocabularies (only when required) that follow best practices whenever possible.
STEP #7	CONVERT DATA	Convert data to a Linked Data representation, typically done by script or other automated processes.
STEP #8	PROVIDE MACHINE ACCESS TO DATA	Provide various ways for search engines and other automated processes to access data using standard Web mechanisms.
STEP #9	ANNOUNCE NEW DATA SETS	Remember to announce new data sets on an authoritative domain. Importantly, remember that as a Linked Open Data publisher, an implicit social contract is in effect.
STEP #10	RECOGNIZE THE SOCIAL CONTRACT	Recognize your responsibility in maintaining data once it is published. Ensure that the dataset(s) remain available where your organization says it will be and is maintained over time.

Table 6.1: Summary of the best practices to publish Linked Data on the Web adapted from [3]

Catalog name	Number of vocabularies	Search Feature	Category	Vocabulary maintenance
vocab.org	19	No	Catalog of generic vocabularies	N/A
ontologi.es	39	No	-//-	N/A
Joinup Semantic Assets	112	Yes	-//-	Yes
Open Metadata Registry	308	Yes	-//-	Yes
BioPortal	355	Yes	Catalog of Domain vocabularies	Yes
SOCoP + OOR	40	Yes	-//-	Yes
Marine Metadata Interoperability	55	Yes	-//-	Yes
SWEET 2.2	200	No	-//-	N/A
DAML	282	No	-//-	No
ODPs	101	No	Catalog of ODPs	Yes
vocab.derie.ie	68	No	Catalog of Organizations	Yes
data.lirmm.fr ontologies	15	No	-//-	Yes
TONES	219	No	Catalog of editors' vocabularies	N/A
Web Protégé	69	No	-//-	Yes

Table 6.2: Catalogs of vocabularies with respectively the number of the ontologies, the presence of a search feature, the catalog category and whether it is maintained or not.

## 6.2 Linked Open Vocabulary (LOV) and Vocabularies

The Linked Open Vocabularies (LOV) initiative aims to bring more insights about published vocabularies in order to foster their reuse. Compared to other projects, LOV benefits from a community<sup>13</sup>

- to assess the quality (including documentation and metadata) and the reuse potential of a vocabulary before it is indexed. LOV currently contains 350+ reusable and well-documented vocabularies;
- to augment vocabularies with explicit information not originally defined in the RDF vocabulary. For example, only 55% of vocabularies have explicit metadata of at least one creator, contributor or editor. In LOV, we augmented this information leading to more than 85% of vocabularies with this information;
- to automatically extract the implicit relations between vocabularies using the Vocabulary Of Friend<sup>14</sup> (VOAF) ontology. These relations can be used as a new metric for ranking terms based on their popularity at the schema level;
- to consider vocabulary semantics in the result rankings: a literal value matched in the `rdfs:label` property has a higher score than for the `dcterms:comment` property.

The way vocabularies are considered in LOV is similar to the way datasets are considered in the LOD cloud [104]. Hence, while the Vocabulary of Interlinked Datasets (VoID) is used to describe relationships between datasets and their vocabularies [105], VOAF is used to describe the mutual relationships between vocabularies. VOAF itself reuses over popular vocabularies such as Dublin Core Terms (dcterms), Vocabulary Of Interlinked Datasets (VoID), Vocabulary for ANNotating vocabulary (vann) and the BIBliographic Ontology (bibo). The vocabulary also introduces new classes such as `voaf:Vocabulary` and `voaf:VocabularySpace`.

The LOV-Bot is the tool that automatically keeps up-to-date the relationships and the metadata about the vocabularies indexed in LOV, using the following steps:

- LOV-Bot checks daily for vocabulary updates (any difference in the vocabulary formal description fetched using content negotiation);
- LOV-Bot uses SPARQL constructs to detect relationships and metadata and creates explicit metadata descriptions in the LOV dataset;
- LOV-Bot annotations are then listed in a back-office administration dashboard for review. This manual procedure enables LOV curators to interact with vocabularies authors and the wider community to raise issues and make remarks or suggestions.

---

<sup>13</sup><https://plus.google.com/communities/108509791366293651606>

<sup>14</sup><http://lov.okfn.org/vocab/voaf/>

The LOV dataset is synchronized with the information presented in human readable format. The Linked Open Vocabularies initiative not only monitors the current state of the ecosystem, but also aims at storing and giving access to vocabulary history. To achieve this goal, the LOV database contains each different version of a vocabulary over the time since its first issue. For each version, a user can access the file and a log of modifications since the previous version.

### 6.2.1 Linked Open Vocabularies

LOV gathers and makes visible indicators not yet harvested before, such as interconnection between vocabularies, versioning history and maintenance policy, past and current referent (individual or organization) if any. The number of vocabularies indexed by LOV is constantly growing (390 as of January 2014) thanks to a community effort and it is the only catalog, to the best of our knowledge, that provides all types of search criteria (metadata search, within/across ontologies search), both an API and a SPARQL endpoint access. According to the categories of ontology libraries defined in [4], LOV falls under the category of “*curated ontology directory*” and “*application platform*”. Figure 6.3 depicts the evolution of vocabularies inserted into LOV catalogue by the curators.

The development of LOV has highlighted a number of interesting research issues such as “*where to find the best domain vocabulary to reuse?*”, and “*is it possible to create a curated catalogue of vocabularies that are links?*” Below we illustrate some of the LOV features useful for ontology search and reuse activities:

**Domain filtering.** Each vocabulary is inserted into LOV based on its domain and/or scope. This information is guided by the scope of the vocabulary, such as City, Science, Library, Metadata, Media, etc. This feature helps in disambiguating the results of the querying service and in classifying vocabularies.

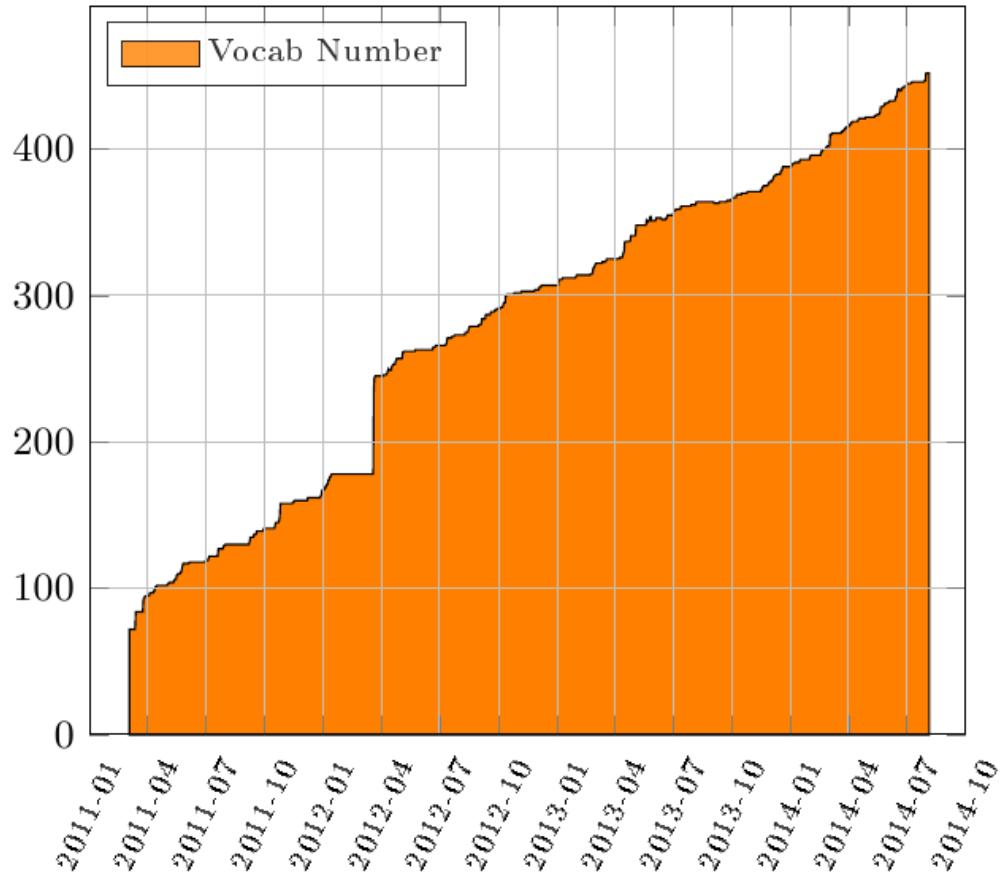
**Content aware Search.** If the searched term matches an `rdfs:label` it will have a higher score than if it matches `dcterms:comment`.

**Links between vocabularies.** One of the key features of LOV’ design is the explicit links between vocabularies,

**Scope of LOV.** The intended use is to promote and facilitate the reuse of vocabularies in the linked data ecosystem.

**Vocabulary Curation.** The collection of vocabularies is maintained by curators in charge of validating and inserting vocabularies in the LOV ecosystem, who double-check the new versions of the vocabulary and give some reviews. Each vocabulary is then automatically enriched with more information about the datasets using it and its relations to other vocabularies.

LOV focuses only on vocabularies (subpart of semantic documents of the web) submitted by any user, reviewed and validated by curators. In addition, LOV keeps



**Fig. 3.** Evolution of the number of vocabularies in LOV.

Figure 6.3: Evolution of the vocabularies inserted into LOV from 2011 to 2014.

track of different versions of the vocabularies in the server that can be retrieved for comparing the differences between along the time evolution. In contrast, Swoogle is designed to automatically discover Semantic Web Documents (SWDs), indexes their metadata and answers queries about it. Thus, the result of a search query retrieved any semantic document. For example, a query of the term *person* gives 16,438 results while in LOV, the term only appears in 134 vocabularies. Watson works similarly to Swoogle, crawling and indexing semantic document at a small scale, explicitly distinguishing for each document (resource), concepts, properties and individuals if available. While in Swoogle the ranking score is displayed, Watson shows the language of the resource and the size. Falcons is a keyword-based search system for concepts and objects on the Semantic Web, and is equipped with entity summarization for browsing. Falcons limits the search to ontologies and a recommendation feature is provided according to users' preferences. However, it does not provide any relationships between the related ontologies, nor any domain

classification of the vocabularies. Table 6.3 compares some key features of Swoogle, Watson, Falcons and LOV.

Feature	Swoogle	Watson	Falcons	LOV
Browsing ontologies	Yes	Yes	Yes	Yes
Scope	SWDs	SWDs	Concepts	ontologies
Metrics	Ranking	Ranking	Ranking	LOD popularity
Domain filtering	No	No	No	Yes
Comments and review	No	Yes	No	Only by curators
Ranking	Doc. based	Doc. based	Doc. based	Metric-based
Web service access	Yes	Yes	Yes	Yes
SPARQL endpoint	No	No	No	Yes
Read/Write	Read	Read & Write	Read	Read
Ontology directory	No	No	No	Yes
Application platform	No	No	No	Yes
Storage	Cache	-	-	Dump & endpoint
Interaction with Contributors	No	-	No	Yes

Table 6.3: Comparison of LOV, with respect to Swoogle, Watson and Falcons, based on part of the framework defined by D'Aquin and Noy in [4].

### 6.2.2 LOV vs NeOn Methodology

The NeOn Methodology is a scenario-based methodology that supports the collaborative aspects of ontology development and reuse, as well as the dynamic evolution of ontology networks in distributed environments. The key assets of the NeOn Methodology are [2]:

- A set of nine scenarios for building ontologies and ontology networks, emphasizing the reuse of ontological and non-ontological resources, the re-engineering and merging, and taking into account collaboration and dynamism.
- The NeOn Glossary of Processes and Activities, which identifies and defines the processes and activities carried out when ontology networks are collaboratively built by teams.
- Methodological guidelines for different processes and activities of the ontology network development process, such as the reuse and re-engineering of ontological and non-ontological resources, the ontology requirements specification, the ontology localization, the scheduling, etc.

LOV is a catalog and API that can fit well within the NeOn methodology for building vocabularies and ontologies. Based on the NeOn Methodology's glossary of activities for building ontologies, LOV is relevant in four activities:

**Ontology Search.** Main LOV's feature is the search of vocabulary terms. These vocabularies are categorized within LOV according to the domain they address. In this way, LOV contributes to ontology search by means of (a) keyword search and (b) domain browsing.

**Ontology Assessment.** LOV provides a score for each term retrieved by a keyword search. This score can be used during the assessment stage and includes a unique term statistical feature<sup>15</sup> which provides for each term registered in LOV the following information: (a) “LOV distribution” that represents the number of vocabularies in LOV that refer to a particular element; (b) “LOV popularity” that shows the number of other vocabulary elements that refers to a particular one; and (c) “LOD distribution” that refers to the number of datasets in LOD which use a particular vocabulary; and (d) “LOD popularity” that refers to the number of vocabulary element occurrences in the LOD.

**Ontology Mapping.** In LOV vocabularies rely on each other in seven different ways. These relationships are explicitly stated using VOAF vocabulary. This data could be useful to find alignments between ontologies, for example one user might be interested in finding equivalent classes for a given class or all the equivalent classes among two ontologies. Listing 6.1 shows the retrieved data when asking for all the equivalent classes and properties between the vocabularies `foaf` and `dcterms` by means of the related VOAF query<sup>16</sup>:

```

1  SELECT DISTINCT ?elem1 ?alignment ?elem2 {
2    {?elem1 <http://www.w3.org/2002/07/owl#equivalentClass> ?elem2}
3    UNION {?elem1 <http://www.w3.org/2002/07/owl#equivalentProperty> ?elem2}
4    UNION {?elem2 <http://www.w3.org/2002/07/owl#equivalentClass> ?elem1}
5    UNION {?elem2 <http://www.w3.org/2002/07/owl#equivalentProperty> ?elem1}
6    FILTER(!isBlank(?elem2))
7    FILTER(!isBlank(?elem1))
8    ?elem1 ?alignment ?elem2.
9    ?elem1 rdfs:isDefinedBy <http://xmlns.com/foaf/0.1/>.
10   ?elem2 rdfs:isDefinedBy <http://purl.org/dc/terms/>.
11 } ORDER BY ?alignment

```

Listing 6.1: SPARQL query asking for all the equivalent classes and properties between the vocabularies `foaf` and `dcterms`.

Figure 6.4 shows the alignments between `foaf` and `dcterms` vocabularies by mean of `owl:equivalentClass` and `owl:equivalentProperty`.

**Ontology Localization.** Labels in different languages are stored in the LOV endpoint for the ontology terms that provide such information. These annotations can be used when translating terms into different languages. This information could be extracted by querying the SPARQL endpoint as shown in Listing

<sup>15</sup><http://lov.okfn.org/dataset/lov/stats/>

<sup>16</sup><http://goo.gl/sTIGQ6>. Prefixes are omitted for readability purpose. You can find the correct namespace for a prefix in LOV.

elem1	alignment	elem2
<a href="http://xmlns.com/foaf/0.1/Agent">http://xmlns.com/foaf/0.1/Agent</a>	<a href="http://www.w3.org/2002/07/owl#equivalentClass">http://www.w3.org/2002/07/owl#equivalentClass</a>	<a href="http://purl.org/dc/terms/Agent">http://purl.org/dc/terms/Agent</a>
<a href="http://xmlns.com/foaf/0.1/maker">http://xmlns.com/foaf/0.1/maker</a>	<a href="http://www.w3.org/2002/07/owl#equivalentProperty">http://www.w3.org/2002/07/owl#equivalentProperty</a>	<a href="http://purl.org/dc/terms/creator">http://purl.org/dc/terms/creator</a>

Figure 6.4: Equivalent classes and properties between foaf and dcterms

"Person"	<a href="http://xmlns.com/foaf/0.1/Person">http://xmlns.com/foaf/0.1/Person</a>
"Persona"@es	<a href="http://xmlns.com/foaf/0.1/Person">http://xmlns.com/foaf/0.1/Person</a>
"Personne"@fr	<a href="http://xmlns.com/foaf/0.1/Person">http://xmlns.com/foaf/0.1/Person</a>
"Person"@en	<a href="http://xmlns.com/foaf/0.1/Person">http://xmlns.com/foaf/0.1/Person</a>

Figure 6.5: Translations example for foaf:Person

6.2<sup>17</sup> where all the labels defined for the terms that have at least one *rdfs:label* containing strictly “person”:

```

1  SELECT DISTINCT ?label2 ?element{
2    ?element rdfs:label ?label1 .
3    ?element rdfs:label ?label2 .
4    FILTER (?label1 != ?label2 ) .
5    FILTER(REGEX(STR(?label1), "person", "i")) .
6  }ORDER BY ?element

```

Listing 6.2: SPARQL query for all the labels defined for the terms containing person.

An excerpt of the query result is shown in Figure 6.5. Based on that result, ‘‘Persona’’@es and ‘‘Personne’’@fr could be used as Spanish and French translations for the English term ‘‘Person’’.

Figure 6.6 shows the activities LOV can support within the overall NeOn methodology activities workflow.

### 6.2.3 Prefixes harmonization

RDF vocabularies bring their meaning to linked data by defining classes and properties, and their formal semantics. Relying on W3C standards RDFS or OWL, those vocabularies are a fundamental layer in the architecture of the Semantic Web. Without the explicit semantics declared in vocabularies, linked data, even using RDF, would be just linked pieces of information where links have no meaning. Interoperability between data and datasets rely heavily on shared vocabularies, but given the distributed nature of the Web, vocabularies are published by independent parties and there is no centralized coordination of this publication, nor should it be. Various independent services have been developed in order to discover vocabularies

<sup>17</sup><http://goo.gl/JJCJ01>

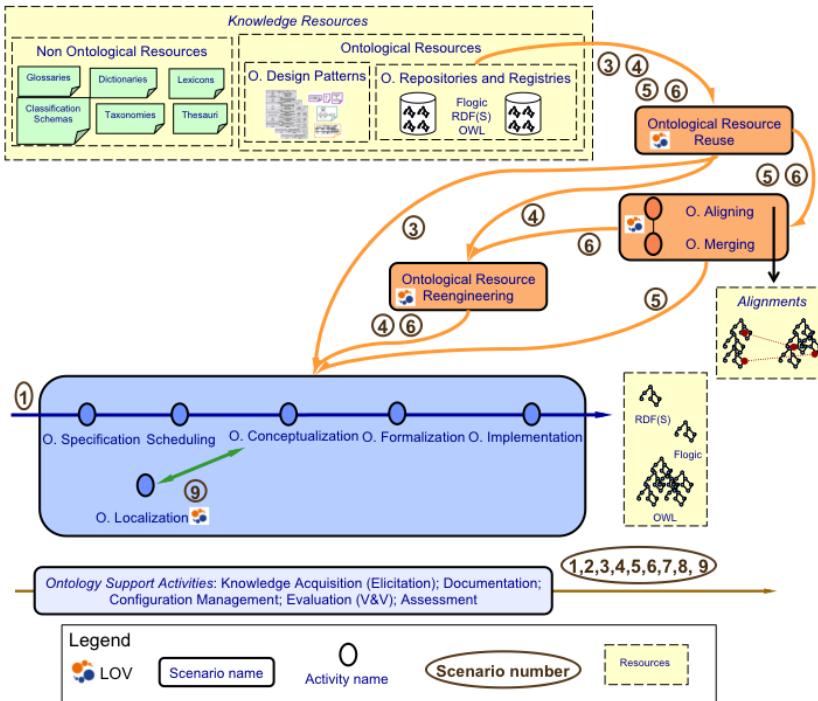


Figure 6.6: Meeting points between LOV and the NeOn methodology, derived from [2].

and provide information about them, and the community of data publishers and vocabulary managers have all interest in complementarity and coordination between such services. In this section, we focus on a specific aspect of vocabularies: their identification by namespaces and associated prefixes.

In the original XML syntax of RDF, prefixes are simply local shortcuts associated with XML namespaces using `xmlns` declarations. The usage of prefixes has been further extended to other syntaxes of RDF such as N3 and Turtle. Although a prefix to namespace association is syntactically limited to the local context of the file in which it is declared, common prefixes such as `rdf:`, `rdfs:`, `owl:`, `skos:`, `foaf:` and many more have become de facto standards. For example, RDFa 1.1 has a default profile made of 11 well-used vocabularies based on their general usage on the Semantic Web according to the crawl of Yahoo! and Sindice as of March 2013<sup>18</sup>. Similarly, the YASGUI SPARQL editor has a list of built-in prefix-namespace associations to ease the construction of SPARQL queries. However, this list of “standard” prefixes is open-ended. Interfaces such as SPARQL endpoints (e.g. Virtuoso) use a list of built-in prefixes declaration for more and more namespaces but the choice of entries in this list is all but transparent. Hence, the reason of a given namespace being or not in this list could be interpreted in many ways, a potential source of technical and social conflicts. Therefore, the notion has been slowly spreading, at least im-

<sup>18</sup><http://www.w3.org/2010/02/rdfa/profile/data/>

plicitly, that common prefixes could and indeed should have a global use, implying some kind of governance and good practices. More and more vocabularies explicitly recommend the prefix that should be used for their namespace, generally using a common if not written good practice to avoid frontal clashes by recommending a prefix not already used. But there is no global policy except implicit rules of fair use to avoid potential conflicts resulting from polysemy (different namespaces using or recommending the same prefix) or synonymy (different prefixes used for the same namespace).

A vocabulary publisher needs to have access to some services capable of monitoring the existing prefixes usage in order to stick to those rules. Moreover, we focus on two services providing such information on prefixes usage namely prefix.cc<sup>19</sup> and LOV (Linked Open Vocabularies) [106]. Both services provide associations between prefixes and namespaces but following a different logic. The prefix.cc service allows anybody to suggest a prefix to namespace association. It supports polysemy and synonymy, and has a very loose control on its crowd-sourced information. What it provides is more a measure of popularity of prefixes and namespaces than a way to put order in them. LOV has a much more strict policy forbidding polysemy and synonymy, enforced by a dedicated back-office database infrastructure, ensuring that each vocabulary in the LOV database is uniquely identified by a prefix, this unique identification allowing the usage of prefixes in various LOV publication URIs. This requirement leads sometimes to a situation where LOV uses prefixes different from the ones recommended by the vocabulary publishers.

### 6.2.3.1 Aligning LOV with Prefix.cc

In this section, we present how we perform the alignment between the two services LOV and prefix.cc. Figure 6.7 shows the evolution of the number of prefixes registered in these two services between April 2009 and July 2013. Our main goals are to align Qnames (prefix) to a unique URI in LOV and to make sure that all the vocabularies in LOV are actually inserted in prefix.cc.

We propose to perform SPARQL queries over all the files of prefix.cc at <http://prefix.cc/popular/all.file.vann> in the FROM clause and compare them to the content of the LOV SPARQL endpoint<sup>20</sup> via a SERVICE<sup>21</sup> call. The SERVICE keyword defined in the SPARQL 1.1 Query Language instructs a federated query processor to invoke a portion of a SPARQL query against a remote SPARQL endpoint [107]. Results are returned to the federated query processor and are combined with results from the rest of the query. To be more generic and standards-compliant, the queries could be run with the Jena ARQ command-line tool to produce a CSV or a JSON serialization that could be easily consumed either by the prefix.cc backend via phpMyAdmin or by the LOV backend.

<sup>19</sup>Service: <http://prefix.cc/>; Code: <https://github.com/cygri/prefix.cc>

<sup>20</sup><http://lov.okfn.org/endpoint/lov>

<sup>21</sup><http://www.w3.org/2009/sparql/docs/fed/service>

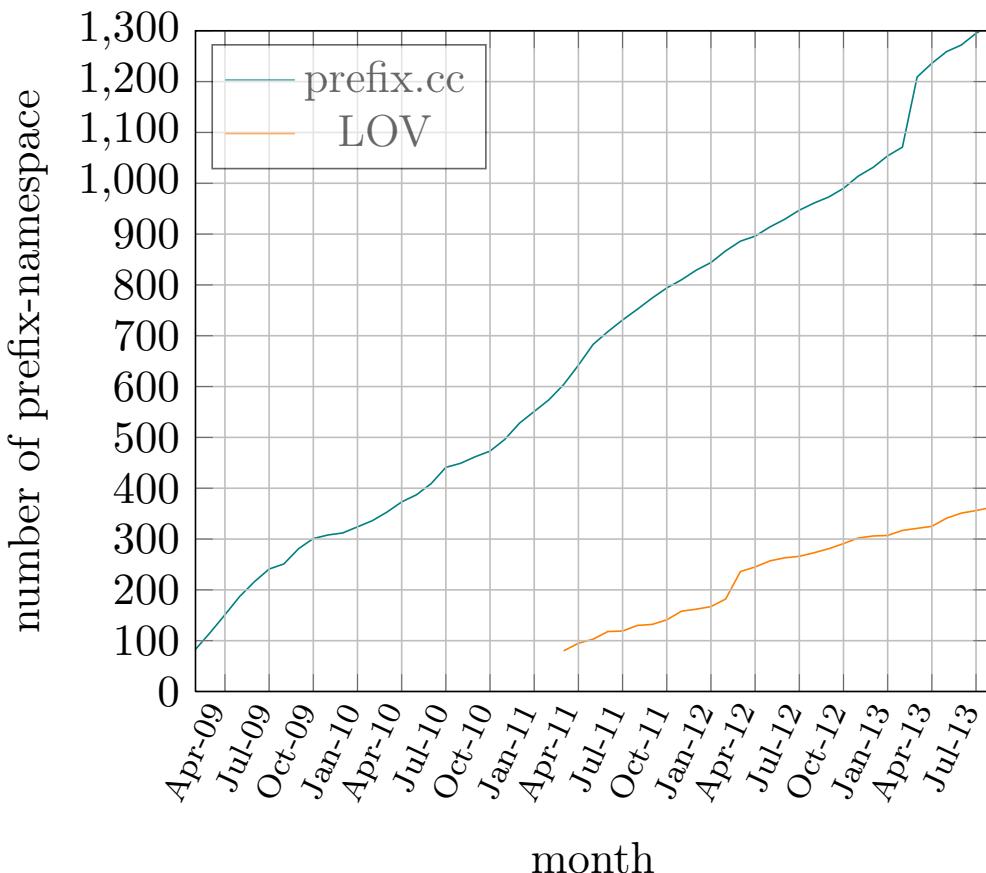


Figure 6.7: Evolution of the number of prefix-namespace pairs registered in prefix.cc and LOV

#### 6.2.3.2 First Task: prefixes in LOV not present in Prefix.cc

First, we compute  $\langle LOV \rangle \text{ INTERSECTS } \langle PREFIX.CC \rangle$  and  $\langle LOV \rangle \text{ MINUS } \{\langle LOV \rangle \text{ INTERSECTS } \langle PREFIX.CC \rangle\}$ . The following SPARQL query finds namespace URIs in LOV that do not exist in prefix.cc along with their LOV prefix.

```

1
2 PREFIX vann: <http://purl.org/vocab/vann/>
3 SELECT ?prefix ?lovURI
4 FROM <http://prefix.cc/popular/all.file.vann> {
5   SERVICE <http://lov.okfn.org/endpoint/lov> {
6     SELECT ?prefix ?lovURI {
7       []vann:preferredNamespacePrefix ?prefix;
8       vann:preferredNamespaceUri ?lovURI;
9     }

```

```

10   }
11 FILTER (NOT EXISTS { []vann:preferredNamespaceUri ?lovURI })
12 OPTIONAL {
13   []vann:preferredNamespacePrefix ?prefix;
14   vann:preferredNamespaceUri ?pccURI;
15 }
16 }
17 ORDER BY ?prefix

```

Listing 6.3: Query to find namespaces in LOV not in prefix.cc

The first results<sup>22</sup> show the following:  $card(LOV) \cap card(PREFIX.cc) = 188$ <sup>23</sup> and  $card(LOV) - card(PREFIX.cc) = 133$ <sup>24</sup> prefixes in LOV not yet registered in prefix.cc. At this point, a first batch of 80 prefixes/namespaces from LOV were safely imported in prefix.cc since there were no conflicts. For the remaining ones, they needed more in-depth analysis due to the conflicts.

#### 6.2.3.3 Second Task: Dealing with Conflicts between Prefix.cc and LOV

In the process of alignment, there were two types of conflicts and we provide appropriate actions and/or solutions accordingly:

- Clashes: cases where we have in both services the same prefix but different URIs;
- Disagreements on preferred namespace: cases where for the same URI, we found different prefixes.

**Clashes.** We performed a SPARQL query as above to identify clashes in vocabularies (30). In Table 6.4, we identify seven different types of issues to deal with, such as (i) real conflicts, (ii) URIs are 404, (iii) URIs are obsolete versions and (iv) two URIs redirecting to the same resource.

**Disagreements on namespace URIs.** The general idea is that if vocabulary editors have not included explicitly a `vann:preferredNamespacePrefix` in their description, the curators of LOV are free to change it and put whatever seems appropriate. At the same time, in prefix.cc, having multiple prefixes for the same namespace IRI is not a problem. However, we computed those prefixes in LOV that have different prefixes in prefix.cc. The following query retrieves the URIs falling in those disagreements:

<sup>22</sup>This query was performed in two weeks between March, 2nd and March, 20th 2013 and at this time,  $card(LOV) = 321$  vocabularies while  $card(Prefix.cc) = 925$

<sup>23</sup><http://www.eurecom.fr/~atemezin/iswc2013/experiments/firstAlignments/intersection-prefixLOV-02-03.csv>

<sup>24</sup><http://www.eurecom.fr/~atemezin/iswc2013/experiments/firstAlignments/inLovNotINPrefixcc-02-03.csv>

Type of issue	# Vocabularies	%
pccURI and lovURI redirect to same resource	8	26.67%
lovURI already in prefix.cc as secondary	7	23.3%
Real conflicts	6	20%
pccURI is 404	4	13.3%
pccURI is an obsolete version	3	10%
lovURI is 404	1	3.3%
lovURI is an obsolete version	1	3.3%

Table 6.4: Type of issues encountered for vocabulary clashes

```

1 PREFIX vann: <http://purl.org/vocab/vann/>
2 SELECT ?prefix ?lovURI ?prefixcc
3 FROM <http://prefix.cc/popular/all.file.vann> {
4   SERVICE <http://lov.okfn.org/endpoint/lov> {
5     SELECT ?prefix ?lovURI {
6       []vann:preferredNamespacePrefix ?prefix;
7       vann:preferredNamespaceUri ?lovURI;
8     }
9   }
10 FILTER (?pccURI = ?lovURI && ?prefix != ?prefixcc)
11 OPTIONAL {
12   []vann:preferredNamespacePrefix ?prefixcc;
13   vann:preferredNamespaceUri ?pccURI;
14 }
15 }
16 ORDER BY ?prefix

```

Listing 6.4: Query to find disagreements LOV and prefix.cc

From the results of this query (61 cases), we have three actions to perform:

- add the lovPrefix (prefix in LOV) in prefix.cc (e.g: adding `geod: http://vocab.lenka.no/geo-deling` to the existing `ngeo` in `pccPrefix`.)
- add more alternative URIs to the existing prefix in prefix.cc (e.g: adding `prov: http://purl.org/net/provenance/ns#` to the existing `hartigprov`, `prv` in `pccPrefix`)
- change a prefix in LOV<sup>25</sup> (e.g: lovPrefix `dc` for `http://purl.org/dc/terms` not in the list `{dcterm, dcq, dct, dcterms}` has been replaced by `dce` in LOV).
- No changes when the lovPrefix is contained in the set of prefixes of prefix.cc.

<sup>25</sup><http://www.eurecom.fr/~atemezin/iswc2013/material/action-sameUriDifferentPrefixes.pdf>

### 6.2.3.4 Social Aspects

Several vocabularies are maintained by a community of users. As part of the alignment process, we contacted the authors, creators or maintainers (if they exist) of vocabularies to involve them as well in the process of changing prefixes, and agree with them to fix some issues regarding their vocabularies. From the homepages of the vocabulary authors and editors collected in LOV, we connect to their social platform accounts such as LinkedIn, Google+ or Twitter. Table 6.5 summarizes some cases of real conflicts where the LOV curators have to find and contact the editors of the vocabularies for negotiation.

prefix	lovURI	Remark
sp	<a href="http://data.lirmm.fr/ontologies/sp#">http://data.lirmm.fr/ontologies/sp#</a>	contact editor at LIRMM ( <i>sp</i> $\Rightarrow$ <i>osp</i> )
scot	<a href="http://scot-project.net/scot/ns#">http://scot-project.net/scot/ns#</a>	contact editors at lovURI
media	<a href="http://purl.org/media#">http://purl.org/media#</a>	contact editors for negotiation
pro	<a href="http://purl.org/spar/pro/">http://purl.org/spar/pro/</a>	contact editors for negotiation
swp	<a href="http://www.w3.org/2004/03/trix/swp-1/">http://www.w3.org/2004/03/trix/swp-1/</a>	contact editors, fix on LOV side
wo	<a href="http://purl.org/ontology/wo/core#">http://purl.org/ontology/wo/core#</a>	contact editors
idemo	<a href="http://rdf.insee.fr/def/demo#">http://rdf.insee.fr/def/demo#</a>	to resolve with INSEE

Table 6.5: LOV and prefix.cc conflicts resolution leading to contact vocabularies editors for negotiation. We provide the prefix, the URI in LOV and the action undertaken.

### 6.2.3.5 Finding Vocabularies in Prefix.cc

We want to find out in prefix.cc, which of the couples (prefix, URI) could be potentially a vocabulary to be further assess to be included in the LOV catalog. To address this question, we first compute all the differences on prefix.cc NOT in LOV, i.e. *PREFIX.CC MINUS (LOV < INTERSECT > PREFIX.CC)*, performing using a SPARQL query. This results in 742 URIs to be checked<sup>26</sup>.

**LOV Check API** We have implemented an API<sup>27</sup> that allows a user to run the LOV-Bot over a distant vocabulary. It takes as parameter the vocabulary URI to process and the time out (integer) specified to stop the process. The result of this action is a set of 26 property-values from which we are interested in using only 8 of them, namely:

- **uri** (string); uri of the vocabulary.
- **namespace** (string) ; namespace of the vocabulary.
- **prefix** (string) ; prefix of the vocabulary

<sup>26</sup><http://www.eurecom.fr/~atemezin/iswc2013/experiments/input/notInLOV.json>

<sup>27</sup><http://lov.okfn.org/dataset/lov/apidoc/>

- **inLOV** (boolean) ; indicates if the vocabulary is already in the Linked Open Vocabularies ecosystem.
- **nbClasses** (int) ; Number of classes defined in the vocabulary namespace.
- **nbProperties** (int) ; Number of properties defined in the vocabulary namespace.
- **dateIssued** (string) ; Vocabulary date of issue.
- **title** (Taxonomy) ; List of titles with language information if available.

The code below gives the response of our algorithm for the vocabulary identified at <http://ns.aksw.org/Evolution/>.

```

1
2 {
3   "dateIssued": "None",
4   "inLOV": false,
5   "namespace": "http://www.agfa.com/w3c/2009/clinicalProcedure#",
6   "nbClasses": 47,
7   "nbProperties": 29,
8   "pccURI": "http://www.agfa.com/w3c/2009/clinicalProcedure",
9   "prefix": "clinproc",
10  "title": [
11    {
12      "dataType": null,
13      "language": "en",
14      "value": "Clinical Procedure"
15    }
16  ],
17  "uri": "http://www.agfa.com/w3c/2009/clinicalProcedure"
18 }
```

Listing 6.5: Sample output of a response of the Check API

#### 6.2.3.6 Experiments

We wrote a script calling the LOV Check API on the URIs in prefix.cc in order to determine the candidate vocabularies to be inserted in LOV using the algorithm in Algorithm 2. We ran four times the experiments (possibly due to some network instabilities) in order to determine from which results what should be assessed. Table 6.6 gives an overview of the number of URIs with respectively the attribute “inLOV=false”(TP), “inLOV=true”(FP) and the errors occurred (Null returned, http/proxy or time out reached by the API).

Regarding the experiments, **Experiment4** gives stable results with less network errors. Therefore, we focus on this experiment when reporting our findings and analysis. We found that 227 (43.48%) are vocabularies in the sense of LOV since they have at least one property or one class. 297 vocabularies (56.51%) might have some problems (or are even not vocabularies at all) as they have neither classes nor properties. Regarding the presence of prefixes, we found 140 (61.67%) of them. The 227

	TP(inLOV=false)	FP(inLOV=true)	Errors
Experiment1	525	44	173
Experiment2	403	26	313
Experiment3	351	28	363
Experiment4	522	44	176

Table 6.6: Experiments looking for stable results of finding vocabularies in prefix.cc.

vocabularies could all be inserted in the LOV catalog since they fulfill the current requirements of what is a “ LOV-able vocabulary”. In this list, we found vocabularies such as `rdf`, `rdfs`, `owl` that are used to build other vocabularies but are not yet integrated in the LOV catalog.

---

**Algorithm 2** finding vocabularies NOT in LOV from prefix.cc

---

```

1: Open notInLOV.json file containing the prefix.cc URIs NOT in LOV
2: initialize item as List
3: Initialize result as collection of item
4: for each pccURI ∈ notInLOV file do
5:   uri ← value of pccURI
6:   uriv ← construct-valid uri
7:   call LOV-Check API with parameter uriv
8:   try/catch HTTPError, URLError, IOError, ValueError
9:   while no error raised do
10:    initialize item to an empty List
11:    append pccURI, prefix, inLOV, namespace, title, dateIssued, nbClasses, nbProperties
      in item List
12:    append item to result
13:   end while
14: end for
15: RETURN output – result

```

---

From the list of URIs that were not LOV-able vocabularies, we wanted to do more analysis by checking the RDF files using the Triple-Checker tool. Our aim is to ensure that we did not leave out some candidate vocabularies or if there are other type of errors such as parsing errors. Table 6.7 provides results classified into 4 categories:

- General errors such as loading files or proxy errors: 78.30%

Total URIs	295	100%
Loading/404 errors	182	61.69%
Vocabularies	36	12.20%
Proxy errors	27	9.15%
50x, 40x errors	22	7.45%
Parsing errors	9	3.05%
Web Pages containers	9	3.05%
No triples found	8	2.71%
RDF data	2	0.67%

Table 6.7: Analysis of the URIs with no classes and no properties while using the LOV-Bot API

- Candidate LOV-able vocabularies: 12.20%
- Clearly not vocabularies (`nbClasses = nbProperties = 0`), typically instances, datasets, html pages: 6.45%
- Others (mainly parsing errors): 3.05%

### 6.3 Vocabulary Ranking Metrics

The linked data principles have gained significant momentum over the last few years as a best practice for sharing and publishing structured data on the Semantic Web [104]. Before being published, data is modeled and ontologies or vocabularies are one of the key elements of a dataset. Vocabularies are the artefact that bring semantics to raw data. One of the major barriers to the deployment of linked data is the difficulty for data publishers to determine which vocabularies should be used since developing new vocabularies has a cost. Catalogues of ontologies are therefore a useful resource for searching terms (classes and properties) defined in those vocabularies. The Linked Open Vocabulary (LOV) initiative [64] is playing a significant role in providing such services to users who can search within curated vocabularies, fostering ontologies reuse. LOV focuses only on vocabularies submitted by users, which are then reviewed and validated by curators. In addition, LOV computes dependencies between vocabularies, keeps track of different versions of them in order to enable their temporal evolution.

To the best of our knowledge, recommending vocabularies to reuse are limited to “popular” or “well-known” ones. We propose a metric combining different features such as how vocabularies are interlinked, or how they are used in real world datasets. This contribution originates also in the desire to bring the traditional concept of Information Content (IC) into the field of the semantic web applied to vocabularies. Many catalogs of ontologies already provide some ranking metrics based on some features. However, we are interested in applying the principles of IC on vocabularies

to investigate if such techniques can give more insights in ontology ranking and ontology usage (e.g in visualization applications).

This Section is organized as follows: Section 6.3.1 defines the theory of Information Content, and the features used for applying Partition Information Content to vocabularies. We present our experiments on the LOV catalogue in the Section 6.3.3.1. We discuss how this ranking metric can be used for vocabulary design and maintenance in Section 6.3.3.2. We compare our results with other rankings for vocabularies in Section 6.3.3.3.

### 6.3.1 Information Content Metrics

Based on probability theory, Information Content (IC) is computed as a measure of generated amount of surprise [108]. More common terms in a given corpus with higher chance of occurrence cause less surprise and accordingly carry less information, whereas infrequent ones are more informative. We reuse the notion of informativeness as the value of information associated with a given entity, where Information Content has a negative relation with its probability. The concept of Information Content can be used to rank each entity, term, or alphabet in the corpus. We apply the Partitioned Information Content to measure the informativeness of Linked Open Vocabularies as a semantic network of resources connected together using different range of relations, as described in [109]. Partitioned Information Content (PIC) is derived from the IC value using some weights. We empirically set those weights according to three features:

- (i) datasets using the vocabulary (*weight* = 2);
- (ii) *outlinks* from a vocabulary, i.e. whether a vocabulary reused other vocabularies (*weight* = 1);
- (iii) *inlinks* to a vocabulary, i.e. whether other vocabularies are reusing this vocabulary (*weight* = 3).

### 6.3.2 Information Content in Linked Open Vocabularies

This experiment aims at bringing the concept of informativeness in the field of terms semantically related as it is the case within semantic web ontologies. The ranking obtained can give additional information based on the Information Content theory to help reusing terms and detecting the ones that are less popular. This can then be used by applications consuming datasets described with these vocabularies. The equation (1) gives the formula for computing the IC value of a term (class or property):

$$IC(t) = -\log_2\left(\frac{\varphi(t)}{N}\right), \quad (6.1)$$

where  $N$  is set to be the maximum value corresponding to the term occurrence in the LOV aggregator (as of June 2014, this value is 3958, and it corresponds to the

popularity of the `skos:prefLabel` property); and  $\varphi(t)$  is the occurrence of the term (but not its popularity).

For computing  $\varphi(t)$ , we use two types of SPARQL queries depending on whether the term is a class (Listing 6.6) or a property (Listing 6.7 considers `owl:ObjectProperty`, `owl:DatatypeProperty` and `rdfs:Property`). Note that we do not yet take into account the `owl:equivalentClass` and `owl:equivalentProperty` axioms that may appear in some vocabularies.

```

1  SELECT (count(?uri1) as ?occ)
2  WHERE {
3    ?uri1 ?p %%classURI . }
```

Listing 6.6: SPARQL query for computing the occurrence of a class

```

1  SELECT (count(?uri1) as ?occ)
2  WHERE {
3    ?uri1 +objectURI+ ?uri2.
4    FILTER (?uri1 != ?uri2) }
```

Listing 6.7: SPARQL query for computing the occurrence of a property

### 6.3.3 Ranking Vocabularies using The Information Content

For computing the PIC value, we use the following formula:

$$PIC(f) = w_f \times \sum_{i=1}^n IC(t_i), \quad (6.2)$$

where  $w_f$  is the weight related to vocabulary  $f$ .

We consider very important that a vocabulary is being reused by other vocabularies and implemented within real world datasets. For example, the `foaf` ontology is weighted 6 because it reuses vocabularies (1), it has been used in some datasets (2) and it is being reused by other vocabularies (3). The `dul`<sup>28</sup> vocabulary is weighted 3 because it doesn't reuse any vocabulary although it is used by several other vocabularies.

#### 6.3.3.1 Experiments on Vocabularies

We use the LOV catalogue, and particularly the LOV aggregator<sup>29</sup> to look at the terms (classes and properties) to compute their Information Content (IC). LOV defines the *LOV Distribution* as the number of vocabularies in LOV that refer to

<sup>28</sup><http://www.ontologydesignpatterns.org/ont/dul/DUL.owl>

<sup>29</sup>[http://lov.okfn.org/endpoint/lov\\_aggregator](http://lov.okfn.org/endpoint/lov_aggregator)

a particular element and the *LOV popularity* as the number of other vocabulary elements that refers to a particular one. Based on the concept of Partitioned Information Content, we implement our ranking measure using the algorithm 3. We take the subset of classes and/or properties with LOV popularity and LOV distribution greater than one. The initial set of vocabularies in LOV is 366. After filtering the candidate terms, we came out with a set of 161 vocabularies (44% or 161 vocabularies) for computing their ranking.

The Table 6.8 gives the Top 15-ranking of the vocabularies according to the informativeness of the classes and properties used within the LOV ecosystem. As the function is proportional to the number of terms, we use a threshold of 22 terms in the vocabularies. For example, the PIC value of `dcterms` is higher than `foaf`'s because the former uses 53 terms (39 properties and 14 classes), while the latter only 35 terms (9 classes and 26 properties), although they both have the same weight value.

The Table 6.9 shows the Top 20 namespaces of vocabularies according to the informativeness of the classes and properties used within the LOV ecosystem, along with their Information Content Value.

Rank	Prefix	PIC score
1	<code>dcterms</code>	1724.844
2	<code>schema</code>	1588.700
3	<code>gr</code>	1261.101
4	<code>foaf</code>	1033.197
5	<code>bibo</code>	876.205
6	<code>time</code>	816.2020
7	<code>skos</code>	805.287
8	<code>dul</code>	797.328
9	<code>ptop</code>	773.167
10	<code>rdafrbr</code>	640.834
11	<code>vaem</code>	630.621
12	<code>ma-ont</code>	508,694
13	<code>prov</code>	497.524
14	<code>swrc</code>	437.394
15	<code>dce</code>	428.618

Table 6.8: Top 15 vocabularies according to their PIC. All the prefixes used for the vocabularies are the ones used by LOV

Rank	vocab term	IC value
1	skos:example	7.7806
2	dce:contributor	4.674
3	skos:scopeNote	4.365
4	dcterms:source	4.299
5	mads:code	3.922
6	mads:authoritativeLabel	3.922
7	vs:userdocs	3.847
8	dce:title	3.79
9	skos:hasTopConcept	3.4547
10	dce:description	2.758
11	dcterms:issued	2.553
12	dce:creator	2.518
13	skos:inScheme	2.202
14	skos:notation	1.924
15	dcterms:description	1.646
16	coll>List	0.761
17	vs:term_status	0.735
18	skos:definition	0.43
19	skos:prefLabel	0.009
20	foaf:Person	0

Table 6.9: Ranking of Top 20 terms (classes and properties) according to their IC value

**Algorithm 3** Vocabulary Ranking algorithm

---

```

1: REQUIRE Dump of lovaggregator file
2: Upload in a triple store for querying
3: Select subset of candidate vocabs LOVaggregatorendpoint
4: for term ∈ lovaggregator do
5:   if LOV distribution ≥ 1 then
6:     if LOV Popularity ≥ 1 then
7:       candidateterms ← append term
8:     end if
9:   end if
10: end for
11: for each term ∈ candidateterms do
12:   GROUP BY vocabulary namespace
13:   COMPUTE weight for each vocabulary
14: end for
15: INITIALIZE PICvector AS a vector
16: for each term ∈ candidateterms do
17:   while term ∈ vocabularySpace do
18:     ICterm ← function IC(term, vocabPrefix)
19:     ICvocab ←  $\sum ICterm$ 
20:   end while
21:   PICvocab ← weight(vocab) × ICvocab
22:   PICvector ← append (PICvocab)
23: ORDER PICvector
24: end for
25: RETURN PICvector

```

---

**6.3.3.2 Application of Information Content on Vocabularies**

We envision various applications using the ranking method based on the Information Content metric while designing semantic web applications, vocabulary life-cycle management or novel recommendation services. We make the following recommendations when using the PIC ranking method on vocabularies:

- Vocabularies on the Top PIC-ranking can be used in visualization applications, i.e. to be displayed to the user as much as possible.
- Terms with lower IC can be used in faceted browsing, and they seem appropriate for generating `sameAs` links during the interconnection and enrichment process. They might also be used for promoting the reuse of terms in vocabularies in general.
- The PIC-ranking could help the ontology designers to monitor and to assess the usage of some terms and lead to update the ontology accordingly. For example,

it can be useful in extending the use of the properties such as `vs:term_status` or `owl:deprecated`.

- Such a ranking can be used to rank organizations or publishers of vocabularies in a time period (e.g. annual) as a way to encourage good qualities vocabularies and/or datasets on the cloud.

The use of the information content on LOV vocabularies can be applied in the datasets interlinking task and visualization applications workflow. For interlinking datasets, this method can help detecting properties with a lower PIC which will be a candidate for the interlinking tool. The PIC score can further be used to track the vocabularies terms status (i.e. `vs:term_status`) or `owl:deprecated` properties by dataset maintainers. From the list of namespaces having deprecated terms (Table 6.10), we observe some correlations with the PIC rank for the vocabularies `dcat` (8), `vcard` (36), `gr` (6), `wl` (2), `pav` (1) and `bibo` (1)<sup>30</sup>. More precisely, the presence of `gr` and `bibo` provides evidence of such a correlation, while the presence of `dcat` and `card` can be explained by the fact that those two vocabularies are in a review process at W3C and subject to re-modeling respectively. Table 6.10 gives an overview of some namespaces with their deprecated terms.

### 6.3.3.3 Related Work and Discussion

In this section, we look at three other catalogues providing rankings for vocabularies: `vocab.cc`, `LODStats` and `prefix.cc`. `vocab.cc`<sup>31</sup> does not provide a ranking for vocabularies but rather proposes a rank for classes and properties. The proposed ranking presented in Table 6.11 is taken from the ranking of classes assuming the namespace is used only once per class.

prefix	#DeprecatedTerms	dcterms:modified
<code>vcard</code>	36	2013-09-25
<code>dcat</code>	8	2013-09-20
<code>gr</code>	6	2011-10-01
<code>wl</code>	2	2013-05-30
<code>pav</code>	1	2013-08-30
<code>bibo</code>	1	2009-11-04

Table 6.10: Sample of vocabularies with terms deprecated in LOV

The LODStats ranking is focused on covering the number of datasets reused in the linked open data cloud [31], which is partially taken into account in our approach. The evidence of that is the first three vocabularies used (`RDF`, `RDFS`, `OWL`) which are considered as the meta model for designing vocabularies. Those vocabularies are not included into the LOV catalog and they do not appear in our ranking. The

<sup>30</sup>As of June 2014, there are 60 terms deprecated in LOV with the query <http://bit.ly/1aqcDf3>

<sup>31</sup><http://vocab.cc/v/tco>

Rank	LOV-PIC	prefix.cc	vocab.cc	lodstats
1	dcterms	yago	intervals	rdf
2	schema	rdf	foaf	rdfs
3	gr	foaf	time	owl
4	foaf	dbp	qb	dcterms
5	bibo	dce	scovo	skos
6	time	owl	freebase	foaf
7	skos	rdfs	mo	dce
8	dul	dbo	owl	void
9	ptop	rss	metalex	geo
10	rdafrbr	skos	doap	aktors
11	vaem	gldp	prov	ro
12	ma-ont	geo	void	obo
13	prov	sc	frbr	app
14	swrc	fb	skos	repo
15	dce	gn	dcterms	time

Table 6.11: Comparing ranking position when using PIC in LOV with respect to prefix.cc and vocab.cc

relative stable position of **foaf** in the four columns of the table suggests that there are equal popular terms. In addition, two other vocabularies have “relative” similar ranking using PIC and LODStats: **skos** and **dcterms**. Regardless the metric used, a short list of the “most popular vocabularies” based on their presence in the Top-15 of the four catalogues is: **foaf**, **skos** followed by **dcterms**, **time**, **dce**, **prov**.

Closer to our work, Schaible *et al.* reported on an empirical study involving 75 linked data experts and practitioners assessing reuse strategies based on various ranking decisions [110]. The goal is to find objective criteria for choosing which vocabularies to reuse and how many can be combined. LODStats and LOV are used to obtain the number of datasets using a specific vocabulary while *vocab.cc* is used for getting the number of occurrence of a vocabulary term. We propose a different metric to rank existing vocabularies that can be furthermore added as a new feature in such a study. One drawback in the model is to use the same weight for two vocabularies with different number of datasets reused. This could be address in the future by using a “function based” weighting for datasets reused (e.g. inverse logarithm) for computing the PIC score.

## 6.4 Datalift Module for Selecting Vocabularies

Datalift platform comes with a module to map data objects and properties to ontology classes and predicates available in the LOV catalogue. Data2Ontology takes an input a “raw RDF”, that is a dataset that has been converted directly from legacy format to triples. The goal is to help to publishers reusing existing ontologies for

converting their dataset for easy discovery and interlinking. It consists of three main components assisting the publisher in selecting properties suitable for the dataset to be published.

- **LOV component:** This component is in charge to connect with the LOV catalogue to retrieve up-to-date ontologies using the LOV search API<sup>32</sup>.
- **Matching Workflow:** Data2Ontology offers to map the data to LOV by automatically proposing a list of best matches. The suggestions are based on the algorithm that:
  1. maximizes linguistic proximity between data properties and ontology predicates,
  2. maximizes the estimated quality of the target ontology using LOV evaluation criteria, and
  3. minimizes the number of candidates target ontologies to increase the overall semantic of the transformed dataset.
- **SPARQL Generator:** This module receives as input the desired mappings and creates the SPARQL CONSTRUCT query needed to implement the mapping. The query can further be modified before the execution to generate a new dataset in the lifting process with Datalift.

Figure 6.8 illustrates the process of matching the properties with ontologies. Depending on the data properties and object properties used to map the elements of the datasets, Data2Ontology can also automatically infer the class model. However, the user can still distribute the predicate matches among several interconnected classes, as well as to add new predicates and values. Once the properties are matched with the desired predicates and their classes, the resulting model can be viewed as a graph.

## 6.5 Summary

We have presented in this chapter our contribution to Linked Open Vocabularies, as part of implementations of the benefits of using LOV in ontology engineering (case of NeOn methodology), prefixes harmonization and alignment of vocabularies published on the Web, or on ranking vocabularies using the principles of Information Content. By applying this latter to Linked Open Vocabularies, we tried to use features that we consider “relevant” to be taken into account when comparing vocabularies (e.g: datasets reused, external vocabularies). We compare with other rankings that are mostly based on the “popularity” of vocabularies. This work can path the way for assessing vocabularies with applications in a more systemic approach for recommending classes/properties in ontology management, or in visualization applications to propose the most “*oh yeah?*” suitable property to be visualized for RDF entities when there is large a large number of properties.

---

<sup>32</sup><http://lov.okfn.org/dataset/lov/apidoc/#lov2search>

DataLift - Data2Ontology    Last LOV update : 28/10/14 | [Update](#) |

**1 Select**    **2 Match**    **3 Refine**    **4 Convert**

Source properties (csv-test (RDF #1))		Target property																											
<input type="text" value="date-de-publication-de-l-etablissement"/>		<a href="#">Match to LOV</a>	<a href="#">Match to project ontologies</a>																										
<b>Suggested matches</b> <table border="1"> <thead> <tr> <th>Source</th> <th>Target</th> </tr> </thead> <tbody> <tr> <td>categorie</td> <td>rdf:categories</td> </tr> <tr> <td>mention</td> <td>http://schema.org/mentions</td> </tr> <tr> <td>adresse</td> <td>http://swrc.ontoware.org/ontology#address</td> </tr> <tr> <td>telephone</td> <td>http://schema.org/telephone</td> </tr> <tr> <td>classement</td> <td>topo:numero</td> </tr> <tr> <td>commune</td> <td>topo:itVert</td> </tr> </tbody> </table>		Source	Target	categorie	rdf:categories	mention	http://schema.org/mentions	adresse	http://swrc.ontoware.org/ontology#address	telephone	http://schema.org/telephone	classement	topo:numero	commune	topo:itVert	date-de-publication-de-l-etablissement <b>Vocabulary space</b> All Results (3) <table border="1"> <thead> <tr> <th>Predicate</th> <th>Vocabulary Space</th> <th>Score (/100)</th> </tr> </thead> <tbody> <tr> <td>poste:complementAdresse</td> <td>Society</td> <td>29.04</td> </tr> <tr> <td>osp:ligneAdresse</td> <td>Society</td> <td>19.32</td> </tr> <tr> <td>osp:adresse</td> <td>Society</td> <td>18.60</td> </tr> </tbody> </table>		Predicate	Vocabulary Space	Score (/100)	poste:complementAdresse	Society	29.04	osp:ligneAdresse	Society	19.32	osp:adresse	Society	18.60
Source	Target																												
categorie	rdf:categories																												
mention	http://schema.org/mentions																												
adresse	http://swrc.ontoware.org/ontology#address																												
telephone	http://schema.org/telephone																												
classement	topo:numero																												
commune	topo:itVert																												
Predicate	Vocabulary Space	Score (/100)																											
poste:complementAdresse	Society	29.04																											
osp:ligneAdresse	Society	19.32																											
osp:adresse	Society	18.60																											
First < 1 > Last																													

Figure 6.8: Matching data properties with ontology predicates in Data2Ontology module



## CHAPTER 7

# License Compatibility Checking

---

*“What kind of Web do we want?”*

Tim Berners-Lee<sup>1</sup>

## Introduction

In this chapter, we report the state of license information attached to both datasets and ontologies on the Web. Given the diversity of the licenses, we analyze the LOD cloud and LOV to report vocabularies available to describe license information (Section 7.2). Section 7.3 presents related work on licenses on the Web, with emphasis on datasets. We then propose a framework to check compatibility between ontologies and datasets based on the defeasible theory (Section 7.4). We conclude in Section 7.5 with some perspectives for future work.

## 7.1 Background

The license of a dataset in the Web of Data can be specified within the data, or outside of it, for example in a separate document linking the data. In line with the Web of Data philosophy [111], licenses for such datasets should be specified in RDF, for instance through the Dublin Core vocabulary<sup>2</sup>. Despite such guidelines, still a lot of effort is needed to enhance the association of licenses to data on the Web, and to process licensed material in an automated way. The scenario becomes even more complex when another essential component in the Web of Data is taken into account: the vocabularies. Our goal is to support the data provider in assigning a license to her data, and verifying its compatibility with the licenses associated to the adopted vocabularies.

We answer this question by proposing an online framework called LIVE<sup>3</sup> (LICenses VERification) that exploits the formal approach to licenses composition proposed in [112] to verify the compatibility of a set of heterogeneous licenses. LIVE, after retrieving the licenses associated to the vocabularies used in the dataset under analysis, supports data providers in verifying whether the license assigned to the dataset is compatible with those of the vocabularies, and returns a warning when this is not the case.

---

<sup>1</sup><https://webwewant.org/>

<sup>2</sup><http://purl.org/dc/terms/license>

<sup>3</sup>The online tool is available at <http://www.eurecom.fr/~atemezin/licenseChecker/>

## 7.2 Statistics about licensed vocabularies

The first step to be addressed consists in analyzing how many vocabularies are licensed, and what is the distribution of such licenses. To achieve this goal, we started our analysis on the Linked Open Vocabularies repository. The LOV initiative stands as an observatory for the re-usable linked vocabularies ecosystem. The initiative goes beyond collecting and highlighting vocabulary metadata, and it now plays a major social role in promoting good practice and improving overall ecosystem quality of publishing vocabularies<sup>4</sup>. We crawled the LOV repository together with the LODstats one searching for licensed vocabularies. The results we obtained are as follows:

**Licensed vocabularies** : we have considered in total 419 vocabularies. The licensed vocabularies are 64 out of 419, eating about 16% of the the total number of considered vocabularies. The properties used to specify the licenses in the vocabularies are <http://creativecommons.org/ns#license> from the Creative Commons vocabulary<sup>5</sup> and <http://purl.org/dc/terms/license> from the Dublin Core vocabulary<sup>6</sup>. Even if the number of licensed vocabularies available on LOV and LODstats is rather low, analyzing the edits of the licensing information it is possible to note an increasing interest in providing further metadata about the data published on the Web of Data, and it holds also for vocabularies<sup>7</sup>. Another interesting result shows that 4 out of 13 vocabularies retrieved searching for the “top 20 most used entities in the LOD cloud<sup>8</sup>” has an explicit license associated. These vocabularies are Good Relations<sup>9</sup>, DBpedia, GeoNames<sup>10</sup>, and FRBR<sup>11</sup>.

**Licenses distribution** : the distribution of the licenses in the licensed vocabularies we retrieved is visualized in Figure 7.1. The most adopted license is Creative Commons Attribution (CC-BY) (30 out of 64 licensed vocabularies), and Creative Commons licenses in general represent the 85% of the licenses used for vocabularies licensing, confirming the trend shown for licensed datasets [113, 112]. Another popular license is ODC Public Domain Dedication and License (PDDL)<sup>12</sup> followed by the W3C Software Notice and License<sup>13</sup>.

We are aware that the obtained results are referred to the data available on LOV and LODstats, and that such repositories of vocabularies are not exhaustive. However,

<sup>4</sup>For more details about LOV, see <http://ercim-news.ercim.eu/en96/special/linked-open-vocabularies>

<sup>5</sup><http://creativecommons.org/ns>

<sup>6</sup><http://dublincore.org/documents/dcmi-terms/>

<sup>7</sup>In this work, we consider vocabularies as data. Other interpretation of the role of vocabularies are discussed in the conclusions.

<sup>8</sup><http://lod-cloud.net/>

<sup>9</sup><http://purl.org/goodrelations/v1#>

<sup>10</sup><http://www.geonames.org/ontology#>

<sup>11</sup><http://vocab.org/frbr/core.html#>

<sup>12</sup><http://opendatacommons.org/licenses/pddl/1.0/>

<sup>13</sup><http://bit.ly/W3C-license>

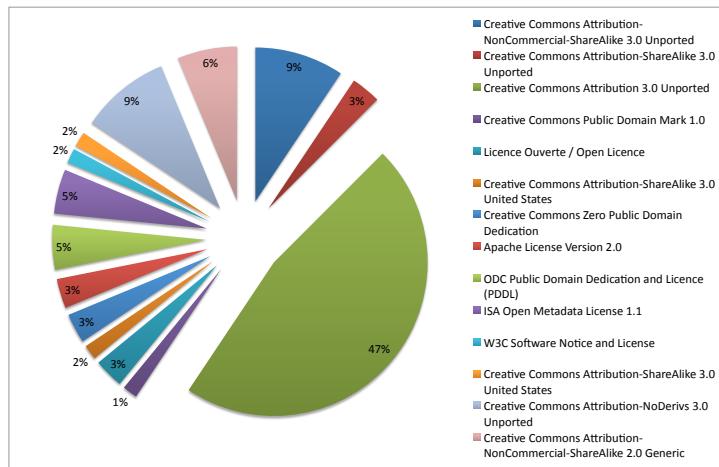


Figure 7.1: Licenses distribution in the LOV licensed vocabularies.

they provide us a reliable picture of what is the ongoing trend in licensing vocabularies<sup>14</sup>. This is particularly true in the case of LOV that furthermore supports a number of good practices for the publication of a vocabulary, among which the addition of the license associated to the vocabulary is highly encouraged.

### 7.3 Related work about licenses in the Web of Data

In the Web scenario, a number of works address the problem of representing and/or reasoning over licensing information. Iannella <sup>15</sup> presents the Open Digital Rights Language (ODRL) for expressing rights information over content, and Gangadharan et al. [114] further extend ODRL developing the ODRL-S language to implement the clauses of service licensing. Gangadharan et al. [115] address the issue of service license composition and compatibility analysis basing on ODRL-S. They specify a matchmaking algorithm which verifies whether two service licenses are compatible. In case of a positive answer, the services can be composed and the framework determines the license of the composite service. Nadah et al. [116] propose to assist licensors' work by providing them a generic way to instantiate licenses, independent from specific formats, and then they translate the license expressed in generic terms into more specific terms compliant with the specific standards used by distribution systems, i.e., ODRL and MPEG Rights Data Dictionaries. Truong et al. [117] address the issue of analyzing data contracts, based on ODRL-S again. Contract analysis leads to the definition of a contract composition where first the comparable contractual terms from the different data contracts are retrieved, and second an evaluation of the new contractual terms for the data mash-up is addressed. Krotzsch and Speiser [118] present a semantic framework for evaluating ShareAlike recur-

<sup>14</sup>We refer the reader interested into statistics about the distribution of licenses on the Web of Data to [113, 112].

<sup>15</sup><http://odrl.net/1.1/ODRL-11.pdf>

sive statements. In particular, they develop a general policy modelling language, then instantiated with OWL DL and Datalog, for supporting self-referential policies as expressed by CC. Gordon [119] presents a legal prototype for analyzing open source licenses compatibility using the Carneades argumentation system. Finally, Rodriguez-Doncel et al. [120, 121] discuss licenses patterns for Linked Data. In particular, they first analyze and discuss six rights expression languages, abstracting their commonalities and outlining their underlying pattern. Second, they propose the License Linked Data Resources pattern which provides a solution to describe existing licenses and rights expressions both for open and not open scenarios. All these works either propose new ways to model licenses information or new formal frameworks to deal with rights. In this paper, we do not address none of these issues, and we adopt the formal framework proposed in [112]. Also Pucella and Weissman [122] propose a logic to check whether the user's actions follow the licenses' specifications. However, as they do not deal with compatibility, do not provide a deontic account of licenses' conclusions, and their logic is not able to handle conflicting licenses, we choose and adapt the deontic logic of [112], which better suits our needs.

Up to our knowledge, the issue of licensed vocabularies has never been addressed. More precisely, no available framework exists dealing with such licenses and verifying in an automated way their potential compatibility with the license associated to datasets. The goal of supporting users in dealing with licensing information has been recently addressed by Cabrio et al. [123] with a different goal, i.e., supporting data publishers in creating RDF licenses representations from natural language texts.

## 7.4 The LIVE Framework

The LIVE framework is a Javascript application, combining HTML and Bootstrap. Hence, installation has no prerequisite. Since the tool is written in Javascript, the best way to monitor the execution time is with the `performance.now()` function. We use the 10 LOD datasets with the highest number of links towards other LOD datasets available at <http://lod-cloud.net/state/#links>. For each of the URLs in Datahub, we retrieve the VoID<sup>16</sup> file in Turtle format, and we use the `voidChecker` function<sup>17</sup> of the LIVE tool to retrieve the associated license, if any. The goal of the LIVE framework is to support data producers to assign a license to the data ensuring the consistency of such license with respect to the licenses assigned to the vocabularies she exploits in the dataset. The input of the LIVE framework (Figure 7.2) consists in the dataset (URI or VoID) whose license has to be verified. The framework is composed by two modules. The first module takes care of retrieving the vocabularies used in the dataset, and for each vocabulary, retrieves the associate license<sup>18</sup> (if any) querying the LOV repository. The module searches out also the license associated to the dataset itself. When all the licensing information of inter-

---

<sup>16</sup><http://www.w3.org/TR/void/>

<sup>17</sup><http://www.eurecom.fr/~atemezin/licenseChecker/voidChecker.html>

<sup>18</sup>Note that the LIVE framework relies on the dataset of machine-readable licenses (RDF, Turtle syntax) presented in [123].

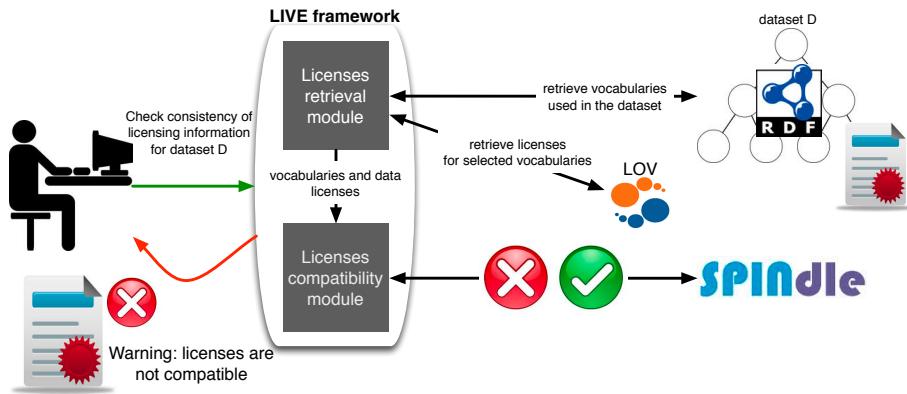


Figure 7.2: LIVE framework architecture.

est has been obtained, the module provides such set of licenses to the compatibility checking module. The second module takes as input the set of licenses (i.e., the licenses of the vocabularies used in the dataset as well as the license assigned to the dataset) to verify whether they are compatible with each others. The result returned by the module is a *yes/no* answer. In case of negative answer, the data provider is invited to change the license associated to the dataset and check back again with the LIVE framework whether further inconsistencies arise.

#### 7.4.1 Licensing information from vocabularies and datasets.

Two use-cases are taken into account: a SPARQL endpoint, or a VoID file in Turtle syntax. In the first use case, the tool retrieves the named graphs present in the repository, and then the user is asked to select the URI of the graph that needs to be checked. Having that information, a SPARQL query is triggered, looking for entities declared as `owl:Ontology`, `voaf:Vocabulary` or object of the `void:vocabulary` property. The final step is to look up the LOV catalogue to check whether they declare any license. There are two options for checking the license: *(i)* a “*strict checking*” where the `FILTER` clause contains exactly the namespace of the submitted vocabulary, or *(ii)* a “*domain checking*”, where only the domain of the vocabulary is used in the `FILTER` clause. This latter option is recommended in case only one vocabulary has to be checked for the license. In the second use case, the module parses a VoID file using a N3 parser for Javascript<sup>19</sup>, and then collects the declared vocabularies in the file, querying again LOV<sup>20</sup> to check their licensing information. When the URIs of the licenses associated to the vocabularies and the dataset are retrieved, the module retrieves the machine-readable description of the licenses in the dataset of licenses [123]. More specifically, such dataset is composed by 37 licenses, comprising all the licenses adopted to certify data in the Linked Data cloud (as

<sup>19</sup><https://github.com/RubenVerborgh/N3.js>

<sup>20</sup>Since LOV endpoint does not support the JSON format in the results, we have uploaded the data in [eventmedia.eurecom.fr/sparql](http://eventmedia.eurecom.fr/sparql).

all the Creative Commons licenses<sup>21</sup>), software licenses (as Mozilla Public License<sup>22</sup> and Microsoft License<sup>23</sup>), and additional licenses for other material on the Web (as the UK Open Government license, and the New Free Documentation License<sup>24</sup>). The dataset provides the licenses in RDF using the Turtle syntax, however Creative Commons licenses are also available in XML/RDF format on the CC website<sup>25</sup>. Figure 7.3 shows the user interface for querying a graph and sample results provided by LIVE tool.

### LIVE LicenseTool We help you detect the right licenses for your dataset

The screenshot shows the LIVE LicenseTool interface. At the top, there are two input fields: '1-Endpoint URL (\*)' containing 'http://eventmedia.eurecom.fr/sparql' and '2-Choose Graph URI (\*)' containing 'http://data.eurecom.fr/bpe'. Below these are two buttons: 'List Graphs»' and 'Check Vocabs»'. The 'Check Vocabs»' button is highlighted. To its right is a table titled 'Graphs Detected' with a green header bar containing 'Sample Results compatibility SPINdle»'. The table has columns: 'Graph URI', 'License A', 'License B', 'Compatibility', and 'Time (ms)'. The data in the table is as follows:

Graph URI	License A	License B	Compatibility	Time (ms)
<a href="http://www.w3.org/2010/09/owl#">http://www.w3.org/2010/09/owl#</a>	CC-BY	CC-BY	Yes	0
<a href="http://www.w3.org/2010/09/rdf#">http://www.w3.org/2010/09/rdf#</a>	CC-BY	OGL	Yes	7
<a href="http://www.w3.org/2002/03/xhtml1/DTD/xhtml1-strict.dtd">http://www.w3.org/2002/03/xhtml1/DTD/xhtml1-strict.dtd</a>	CC-BY	ODBL	Yes	6
<a href="http://www.w3.org/2002/03/xhtml1/DTD/xhtml1-transitional.dtd">http://www.w3.org/2002/03/xhtml1/DTD/xhtml1-transitional.dtd</a>	CC-BY	CCO	Yes	3
<a href="http://www.w3.org/2002/03/xhtml1/DTD/xhtml1-strict.dtd">http://www.w3.org/2002/03/xhtml1/DTD/xhtml1-strict.dtd</a>	CC-BY	CC-BY-SA	Yes	6
<a href="http://www.w3.org/2002/03/xhtml1/DTD/xhtml1-transitional.dtd">http://www.w3.org/2002/03/xhtml1/DTD/xhtml1-transitional.dtd</a>	PDDL	OGL	Yes	6
<a href="http://www.w3.org/2002/03/xhtml1/DTD/xhtml1-strict.dtd">http://www.w3.org/2002/03/xhtml1/DTD/xhtml1-strict.dtd</a>	PDDL	EUROSTAT	yes	9
<a href="http://www.w3.org/2002/03/xhtml1/DTD/xhtml1-transitional.dtd">http://www.w3.org/2002/03/xhtml1/DTD/xhtml1-transitional.dtd</a>	VocabURI	VocabLicense	DatasetLicense	Compatible?
<a href="http://rdf.insee.fr/de">http://rdf.insee.fr/de</a>	<a href="http://purl.org/vocommons/voaf">http://purl.org/vocommons/voaf</a>	<a href="http://creativecommons.org/licenses/by/3.0/">http://creativecommons.org/licenses/by/3.0/</a>	<a href="http://creativecommons.org/publicdomain/zero/1.0/">http://creativecommons.org/publicdomain/zero/1.0/</a>	Yes
<a href="http://rdf.insee.fr/de">http://rdf.insee.fr/de</a>	<a href="http://purl.org/vocommons/voaf">http://purl.org/vocommons/voaf</a>	<a href="http://creativecommons.org/licenses/by/3.0/">http://creativecommons.org/licenses/by/3.0/</a>	<a href="http://creativecommons.org/publicdomain/zero/1.0/">http://creativecommons.org/publicdomain/zero/1.0/</a>	Yes
<a href="http://rdf.insee.fr/de/geo">http://rdf.insee.fr/de/geo</a>	<a href="http://purl.org/vocommons/voaf">http://purl.org/vocommons/voaf</a>	<a href="http://creativecommons.org/licenses/by/3.0/">http://creativecommons.org/licenses/by/3.0/</a>	<a href="http://creativecommons.org/publicdomain/zero/1.0/">http://creativecommons.org/publicdomain/zero/1.0/</a>	Yes

Figure 7.3: LIVE tool user interface and sample results

#### 7.4.2 Licenses compatibility verification.

The logic proposed in [112] and the licenses compatibility verification process has been implemented using SPINdle [124] – a defeasible logic reasoner capable of inferencing defeasible theories with hundredth of thousand rules.

As depicted in Figure 7.4, after receiving queries from users, the selected licenses (represented using RDF) will be translated into the DFL formalism supported by SPINdle using the *RDF-Defeasible Theory Translator*. That is, each RDF-triple will be translated into a defeasible rule based on the subsumption relation between the *subject* and *object* of a RDF-triples. In our case, we can use the subject and object of the RDF-triples as the antecedent and head of a defeasible rule, respectively.

<sup>21</sup><http://creativecommons.org/licenses/>

<sup>22</sup><http://www.mozilla.org/MPL/2.0/>

<sup>23</sup><http://referencesource.microsoft.com/referencesourcelicensing.aspx>

<sup>24</sup><http://www.gnu.org/copyleft/fdl.html>

<sup>25</sup>For instance, Creative Commons Attribution 4.0 license is available at <http://creativecommons.org/licenses/by/4.0/rdf>

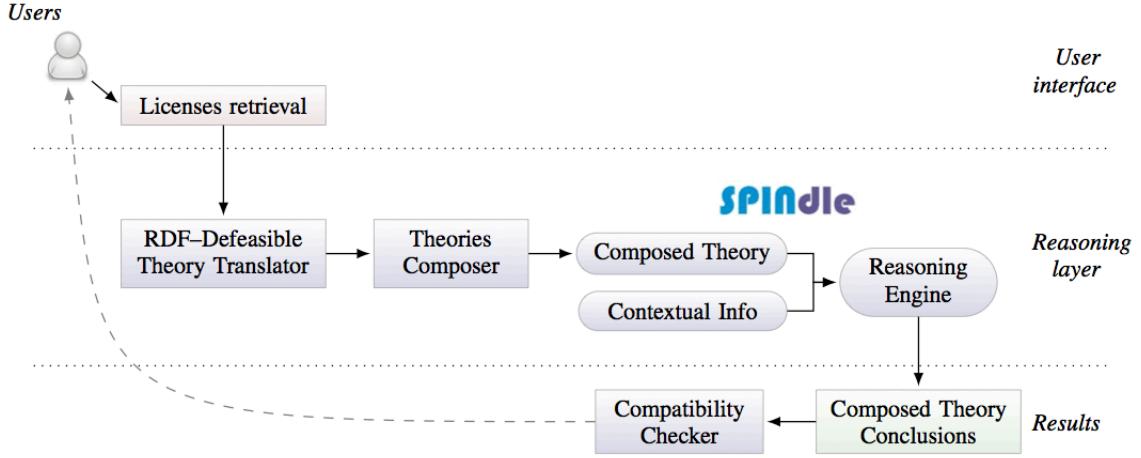


Figure 7.4: Licenses compatibility module.

Besides, the translator also supports direct import from the Web and processing of RDF data into SPINdle theories. The *RDF-Defeasible Theory Translator* will translate the RDF-licenses into the DFL formalism supported by SPINdle.

The translated defeasible theories will then be composed into a single defeasible theory based on the logic proposed in [112], using the *Theories Composer*. Afterwards, the composed theory, together with other contextual information (as defined by user), will be loaded into the SPINdle reasoner to perform a compatibility check before returning the results to the users.

We have evaluated the time performances of the LIVE framework in two steps (Table 7.1).

Dataset	LicRetrievl(ms)	#vocabularies	LicCompatibility(ms)	LIVE(ms)
rkb-explorer-dblp	4,499	1	0	4,499
rkb-explorer-southampton	14,693	1	0	14,693
rkb-explorer-eprints	3,220	1	0	3,220
rkb-explorer-acm	3,007	1	0	3,007
rkb-explorer-wiki	14,598	1	0	14,598
rkb-explorer-rae2001	3,343	1	0	3,343
rkb-explorer-citeseer	2,760	1	0	2,760
rkb-explorer-newcastle	3,354	1	0	3,354
rkb-explorer-kisti	4,094	5	6	4,100
270a.info	13,202	48	8	13,210

Table 7.1: Evaluation of the LIVE framework.

First, we evaluate the time performances of the licenses compatibility module: it

needs about 6ms to compute the compatibility of two licenses. Second, we evaluate time performances (Chrome v. 34) of the whole LIVE framework for the 10 LOD datasets with the highest number of links towards other LOD datasets, considering both the licenses retrieval module and the licenses compatibility one. The results show that LIVE provides the compatibility evaluation in less than 5 seconds for 7 of the selected datasets. Time performances of LIVE are mostly affected by the first module while the compatibility module does not produce a significant overhead. For instance, consider the Linked statistical Data-spaces<sup>26</sup>, a dataset where we retrieve the licensing information in both the dataset and the adopted vocabularies. In this case, LIVE retrieves 48 vocabularies in 13.20s , the license for the dataset is CC-BY, and the PDDL license is attached one of the vocabularies<sup>27</sup>. The time for verifying the compatibility is 8ms, leading to a total of 13.208s.

#### 7.4.3 Perspectives

In the present work, we consider vocabularies as data but this is not the only possible interpretation. For instance, we may see vocabularies as a kind of compiler, such that, after the creation of the dataset then the external vocabularies are no more used. In this case, what is a suitable way of defining a compatibility verification? We will investigate this issue as well as we will evaluate the usability of the online LIVE tool to subsequently improve the user interface.

### 7.5 Summary

In this chapter, we have presented an online tool to check the compatibility between datasets and vocabularies based on the RDF-defeasible of SPINdle. We have introduced the LIVE framework for licenses compatibility. The goal of the framework is to verify the compatibility of the licenses associated to the vocabularies exploited to create a RDF dataset and the license associated to the dataset itself. Several points have to be taken into account for future work.

---

<sup>26</sup><http://270a.info/>

<sup>27</sup><http://purl.org/linked-data/cube>

# CHAPTER 8

# Conclusions and Future Perspectives

---

*“How do we know that Semantic Web technologies were actually better here, as opposed to being what the developers found most familiar?”*  
(D. Karger)<sup>1</sup>

## 8.1 Conclusions

This thesis is focused on the challenges of publishing geodata on the Web and a more generic approach to visualize data as Linked Data target to lay-users. The former considers the diversity of different formats used to publish legacy geospatial data, the different projections (or Coordinate Reference Systems) and the representation of complex geometries. The latter approach is different from the state-of-the-art in visualizations where the complexity of SPARQL and RDF is not sufficiently hidden from the users. A deep analysis of the literature has revealed some limitations in the publication of geospatial data and visualization tools, namely:

- Limited of complex geometries exposed in structured representation, instead of literals.
- The absence of an explicit reference to CRSs in direct georeference data on the Web.
- Absence of visualization tool targeted to lay users to easily grasp the essence of the underlying data published as LOD.
- Many data silos for applications built and published on the Web, lost in many HTML pages.
- Few tools that provide an integrated environment for publishing raw data into Linked Data, from data modeling until the final step of storing the dataset in an endpoint.
- The difficulty for publishers to understand and check the compatibility of the licenses between vocabularies and datasets.

---

<sup>1</sup><http://goo.gl/hQQ3h5>

In this thesis, we have provided different vocabularies that all together support the publication of geodata integrating almost all the CRSs, extending the existing vocabularies. The vocabularies have been used to publish the French Administrative Units, with the data compatible with GeoSPARQL standards. Regarding the visualizations, after reviewing visual tools and existing applications on the Web, we have developed an ontology to better expose the data on the Web for better interoperability. We have also proposed a framework for automatically generating visualizations based on categories detected on datasets published as Linked Data, using predefined high level categories used in Information Visualization taxonomy and mapped with vocabularies.

### **Review of the Contributions**

This section reviews the main contributions of this thesis and the solutions to solved some of the open research problems in publishing and consuming data on the Semantic Web:

- We modelled and implemented of a vocabulary for geometry, topological entities and Coordinate Reference Systems (see Section 2.4).
- We have implemented of an API for converting data between different CRSs accessible on the Web (see Section 2.2).
- We have published different projections systems used in France with unique URIs to improve look up and integration in structured geometries on the Web (see Section 2.3.3.2).
- We have contributed in the development of the Datalift platform, an integrated environment to publish raw data on the Web (see Section 3.5).
- We have provided a comparison of triple stores for geodata against the geometries handled (literal or structured) to assess which one to use when publishing geospatial data (see Section 3.4.2).
- We have published French administrative units available as LOD available at <http://data.ign.fr> endpoint, based on the vocabularies developed and implemented. Moreover, we have provided interlinking with relevant existing geospatial datasets (see Section 3.6.1 and 3.6.2 ).
- We have published in RDF 15 millions of addresses from Open Street Map France using the location address vocabulary (see Section 3.6.3).
- We have contributed to the *French LOD (FrLOD)* cloud, with more datasets published using the Datalift platform, and covering the French territory (see Section 3.6.4).

- We surveyed and classified applications built on top of Open government portals, and proposed a vocabulary for semantically annotate and improve the discovery of applications contests in Open Data event (see Section 4.2 and Section 4.3).
- We have proposed a generic approach for automatically generating visualizations based on predefined categories (see Section 5.2).
- We have implemented and evaluated an approach for determining which properties are suitable to use for an entity, based on the Google Knowledge Panel (see Section 5.3).
- We have developed two innovative applications consuming events and statistical datasets (see Section 5.5 and Section 5.6).
- We have proposed a generic plugin tool that can improve the discovery of applications contests in Open Data events (see Section 5.7).
- We have also proposed an approach to harmonize prefixes used in different catalogues of vocabulary, with an evaluation based on Linked Open Vocabulary (see Section 6.2.3).
- We have developed new ranking metrics for vocabularies based on Information Content theories and applied in LOV (see Section 6.3).
- Finally, we have built a more efficient tool for checking license compatibility between vocabularies and datasets (see Section 7.4).

## 8.2 Future Perspectives

In this thesis, we have tackled some open research problems within the context of publishing and consuming open data on the Web but there are still open issues and challenges for future work. We mention some of the most important from our perspective, based on different aspects related to the workflow of publishing Linked Data, more specifically in the geospatial domain.

### 8.2.1 Opportunities and Challenges for IGN-France

The need for interoperable reference geographic data to share and combine georeferenced environmental spatial information is highlighted by the INSPIRE Directive. The INSPIRE Directive [32] aims to create a European Union (EU) spatial data infrastructure<sup>2</sup>. INSPIRE is based on a number of high level common principles, with some of them very closed to the key concepts of Semantic Web goals, and specifically the Linked Data principles. We provide below the correspondence of our contributions mapped to the five goals of INSPIRE:

---

<sup>2</sup><http://inspire.jrc.ec.europa.eu/index.cfm/pageid/48>

- **P1:** *Data should be collected only once and kept where it can be maintained most effectively.* The use of good and stable URI policies can help achieve this principle. IGN as a French geospatial dataset provider, is committed to accurate information, and so will be the URIs chosen for the experimental portal.
- **P2:** *It should be possible to combine seamless spatial information from different sources across Europe and share it with many users and applications.* This is more or less the goal of the interlinking tasks performed with other datasets on the wild. The models developed and well-documented can ease the conversion by other mapping agencies or institutions of their datasets.
- **P3:** *It should be possible for information collected at one level/scale to be shared with all levels/scales; detailed for thorough investigations, general for strategic purposes.* One of the drawback of the models proposed is that they don't currently admit many geometries attached to a feature. This will certainly be one of the extension foreseen for the models. However, the precise classifications for the features is a starter to fulfill this principle.
- **P4:** *Geographic information needed for good governance at all levels should be readily and transparently available.* Publishing `data.ign.fr` is one of the objective to have also data both in human and machine readable manner using semantic concepts and technologies.
- **P5:** *Easy to find what geographic information is available, how it can be used to meet a particular need, and under which conditions it can be acquired and used.* Publishing data on the web contribute *per se* in leveraging their discovery and integration. Moreover, an explicit license attached to datasets published help achieving this principle.

For geographic data producers, the benefit of publishing their data on the Web according to Linked Data (LD) principles is twofold:

1. First, their data are interoperable with other published datasets and they can be referenced by external resources and used as spatial reference data, which would not have been straightforward when published according to spatial data infrastructures (SDI) standards.
2. Second, the use of Semantic Web technologies can help addressing interoperability issues which are not solved yet by geographic information standards.

Moreover, the French national mapping agency (IGN) has different types of license policies for accessing data from their professional portal<sup>3</sup> (e.g., research purpose, commercial use, access on demand, etc.), with some of them not necessary “open” or free to access: (e.g., BD TOPO®). Although there is a clear understanding of

---

<sup>3</sup><http://professionnels.ign.fr/>

the benefits of publishing and interconnecting data on the Web, ongoing investigations on how to combine licenses on datasets are under consideration at IGN. Two solutions are under investigation:

1. Different license policies attached to given datasets: Here the attached license is given directly when published. So for example, if it is an open license, the endpoint is publicly available to be queried without any restriction.
2. The use of a security access mechanism on top of the datasets granting access according to a predetermined configuration list of named graphs, resources and operations allowed. This solution goes along with the work of Rotolo et al.[125], where even if there is an endpoint, a module for configuring the types of queries to perform and access policies have to be defined for subsets with special care to take into account compositions of licenses in the results.

According to Linked Data principles URIs should remain stable, even if administrative units change or disappear. This implies adapting the data vocabulary in order to handle data versioning and time scale evolution of the data. This issue will be addressed in our future work, as we are working on releasing a spatio-temporal dataset describing the evolution of communes since the French Revolution. Another issue deals with the automation of the whole publication process, from traditional geographic data to fully interconnected RDF data. The last issue deals with the use of multiple geometries for describing a geographic feature: geometries with different levels of detail, different CRS, different representational choices. This has been superficially addressed in our use case with the use of both polygons and points for representing respectively the surface and the centroid of departments, but should be further investigated for both query answering and map design purposes.

### 8.2.2 Generic Visualizations on Linked Data

We plan to use a more exhaustive set of vocabularies in our generic queries for detecting those categories, plugging into directly the wizard to the LOV catalogue. The aggregation properties can be extended to take other semantic relations (e.g: `skos:exactMatch`) into account. Additionally, we plan to make an evaluation of the prototype and compare it to related tools such as the ones aiming to build profiles of datasets. We also need to quantify when a category is “important” within a dataset. For example, is it enough for a dataset to be classified GEODATA with ten triples containing location? From which number of triples could the categories and hence the visualizations be assigned? These issues can further be investigated to find the best trade-off. Another drawback of our work on visualizations is the lack of user evaluation, with experiments to understand users’ needs, focusing more on the semantic aspects than just the exploration ones (webby-interface). A natural follow-up is use these evaluations and re-adapt the applications/visualizations based on the results.

### 8.2.3 Vocabularies and LOV

Work on the harmonization of prefixes can be extended in several directions. Sticking to the two services we have studied and already contributed to harmonize, the possible next steps would be to automate as far as possible the tasks that have been made semi-automatically so far:

- *i)* developing a unique interface for submitting namespaces and prefixes to both services;
- *ii)* bridging the LOV back-office and the prefix-cc database using both services API in order to publish a list of common recommended prefixes.

The latter goes beyond the limited framework of the two original services since such a list could be consolidated and endorsed by the main actors in vocabulary publication and management, and recommended for use in linked data applications. This could be picked up by the upcoming W3C Vocabulary Management Working Group as part of the new Data Activity<sup>4</sup>.

As per ranking vocabularies, we aim to take into account the equivalence axioms (between classes and properties) when computing the Information Content, and more generally, all sort of semantic relationships between terms. Also, we plan to compare our ranking model with other ranking approaches such as graph-based ones (e.g., pagerank). Another future direction is to investigate the dependency ranking between vocabularies, by focusing on a specific type of “inlinks” (i.e. extensions, generalizations) and study how they affect the information content (PIC) metrics.

We have made the assumption in this thesis that access to data was either by querying a SPARQL endpoint, or by browsing or by downloading the dumps. Recently, a new way of accessing the data on the Web is emerging: through triple pattern fragments<sup>5</sup>. Linked Data Fragments [126] aims at exploring endpoints with simple fragments to solve queries at the client side with server data. Servers can offer data at low processing cost in a way that enables client-side querying, thus, moving intelligence from the server to the client. One possible direction of study could be to use client-side concept for evaluating endpoints consuming only structured geometries versus literal for real-world applications. Finally, triple fragment concepts can be applied also to detect also patterns for visualization in different endpoints.

As the Linked Data grows, so will datasets and ontologies on geospatial data. Geodata publishers will release more often and frequently their data on the Web. There will be a need for more analytical tools, especially in data mining to provide feedback to the publishers with respect to triples usage and retrieval. Streaming

---

<sup>4</sup><http://www.w3.org/2013/05/odbp-charter.html>

<sup>5</sup><http://linkeddatafragments.org/>

geospatial data on the Web will require efficient implementations of spatial functions to be able to query on-the-fly data with temporal information. Thus, geo-temporal streaming data modeling, querying and analysis on the Web are likely to be the next challenges that Semantic Web technologies will have to solve.



# Bibliography

- [1] Ana-Maria Olteanu. *Fusion de connaissances imparfaites pour l'appariement de données géographiques*. PhD thesis, Universite Paris-Est, France, October 2008. xi, 21, 54, 55
- [2] Maria del Carmen. Suárez-Figueroa. *NeOn Methodology for Building Ontology Networks: Specification, Scheduling and Reuse*. PhD thesis, Universidad Politecnica de Madrid, Spain, June 2010. <http://oa.upm.es/3879/>. xii, 147, 150
- [3] Bernadette Hyland, Ghislain Atemezing, and Boris Villazon-Terrazas (eds). Best practices for publishing linked data. W3C Working Group Note, 2014. <http://www.w3.org/TR/ld-bp/>. xiv, 13, 137, 142
- [4] Mathieu DÁquin and Natasha F. Noy. Where to publish and find ontologies? a survey of ontology libraries. *Web Semantics: Science, Services and Agents on the World Wide Web*, 11(0):96–111, 2012. xiv, 139, 140, 141, 145, 147
- [5] Tim Berners-Lee. Design issues for linked data, 2006. <http://www.w3.org/DesignIssues/LinkedData.html>. 1, 5
- [6] A. Jentzsch, R. Cyganiak, and C. Bizer. State of the lod cloud (september 2011), 2011. <http://lod-cloud.net/state/>. 1
- [7] Schmachtenberg Max, Bizer Christian, and Paulheim Heiko. Adoption of the linked data best practices in different topical domains. In *Proceedings of the ISWC 2014, RDB Track (To appear)*, 2014. <http://dws.informatik.uni-mannheim.de/fileadmin/lehrstuehle/ki/pub/SchmachtenbergBizerPaulheim-AdoptionOfLinkedDataBestPractices.pdf>. 1
- [8] Tom Heath and Christian Bizer. *Linked Data: Evolving the Web into a Global Data Space*. Morgan & Claypool, 1st edition, 2011. 2
- [9] Matthew Perry and John Herring. OGC GeoSPARQL- A Geographic Query Language for RDF Data. In *OGC Implementation Standard, ref: OGC 11-052r4*, 2012. 3, 30
- [10] Glen Hart and Catherine Dolbear. *Linked Data: A Geographic Perspective*. CRC Press, 1 edition, January 2013. 3
- [11] Jens Lehmann, Christian Bizer, Georgi Kobilarov, Soren Auer, Christian Becker, Richard Cyganiak, and Sebastian Hellmann. Dbpedia - a crystallization point for the web of data. *Journal of Web Semantics: Science, Services and Agents on the World Wide Web*, 7(3):154–165, 2009. 4

- [12] Sören Auer, Jens Lehmann, and Sebastian Hellmann. LinkedGeoData - Adding a Spatial Dimension to the Web of Data. In *8<sup>th</sup> International Semantic Web Conference (ISWC'09)*, 2009. 5, 6, 35
- [13] John Goodwin, Catherine Dolbear, and Glen Hart. Geographical Linked Data: The Administrative Geography of Great Britain on the Semantic Web. *Transactions in GIS*, 12:19–30, 2008. 6, 35
- [14] Alexander de León, Luis M. Vilches, Boris Villazón-Terrazas, Freddy Priyatna, and Oscar Corcho. Geographical linked data: a Spanish use case. In *International Conference on Semantic Systems (I-SEMANTICS'10)*, Graz, Austria, 2010. 6, 35, 50
- [15] John Goodwin, Catherine Dolbear, and Glen Hart. Geographical linked data: The administrative geography of great britain on the semantic web. *Transactions in GIS*, 12:19–30, 2008. 6
- [16] Juan Salas and Andreas Harth. Finding spatial equivalences across multiple RDF datasets. In *4<sup>th</sup> International Terra Cognita Workshop*, pages 114–126, Bonn, Germany, 2011. 6, 35
- [17] Max J Egenhofer. Toward the semantic geospatial web. In *Proceedings of the 10th ACM International Symposium on Advances in geographic information systems*, pages 1–4. ACM, 2002. 6, 19
- [18] Manolis Koubarakis, Manos Karpathiotakis, Kostis Kyzirakos, Charalampos Nikolaou, and Michael Sioutis. Data models and query languages for linked geospatial data. In *Reasoning Web. Semantic Technologies for Advanced Query Answering - 8th International Summer School 2012, Vienna, Austria, September 3-8, 2012. Proceedings*, pages 290–328, 2012. [http://dx.doi.org/10.1007/978-3-642-33158-9\\_8](http://dx.doi.org/10.1007/978-3-642-33158-9_8). 9, 59
- [19] David F Huynh, David R Karger, and Robert C Miller. Exhibit: lightweight structured data publishing. In *Proceedings of the 16th international conference on World Wide Web*, pages 737–746. ACM, 2007. 12
- [20] Bernadette Hyland, Ghislain Atemezing, Michael Pendleton, and Biplav Srivastava (eds). Linked data glossary. W3C Working Group Note, 2013. <http://www.w3.org/TR/ld-glossary/>. 13, 139
- [21] Richard Cyganiak and Dave Reynolds (eds). The rdf data cube vocabulary. W3C Recommendation, 2014. <http://www.w3.org/TR/vocab-data-cube/>. 13, 120, 121
- [22] Dave Reynolds. The organization ontology. W3C Recommendation, 2014. <http://www.w3.org/TR/vocab-org/>. 13
- [23] Fadi Maali and John Erickson (eds). Data catalog vocabulary (dcat). W3C Recommendation, 2014. <http://www.w3.org/TR/vocab-dcat/>. 13, 91

- [24] Peter A Burrough, Rachael McDonnell, Peter A Burrough, and Rachael McDonnell. *Principles of geographical information systems*, volume 333. Oxford university press Oxford, 1998. 20
- [25] ESRI. Esri shapefile technical description. An ESRI White Paper, 1998. <http://www.esri.com/library/whitepapers/pdfs/shapefile.pdf>. 22
- [26] Open Geospatial Consortium Inc. Simple feature access - part 1: Common architecture. Technical report, OGC, 2011. <http://www.opengeospatial.org/standards/sfa>. 23
- [27] Howard Butler (Hobu Inc.), Martin Daly (Cadccorp), Allan Doyle (MIT), Sean Gillies (UNC-Chapel Hill), Tim Schaub (OpenGeo), and Christopher Schmidt (MetaCarta). The geojson format specification, 2008. <http://geojson.org/geojson-spec.html>. 23
- [28] Muneendra Kumar. World geodetic system 1984: A modern and accurate global reference frame. *Marine Geodesy*, 12(2):117–126, 1988. 23, 48
- [29] International Organization for Standardization. ISO 19115 Geographic Information - Metadata, 2003. 24
- [30] R.A. Knippers. Geometric aspects of mapping. Non-published educational notes, 2009. <http://kartoweb.itc.nl/geometrics/index.html>. 25
- [31] Jan Demter, Sören Auer, Michael Martin, and Jens Lehmann. LODStats – An Extensible Framework for High-performance Dataset Analytics. In *18<sup>th</sup> International Conference on Knowledge Engineering and Knowledge Management (EKAW)*, pages 353–362, Galway, Ireland, 2012. 30, 164
- [32] INSPIRE Thematic WG CRS and Geographical Grid Systems . Guidelines INSPIRE Specification on Coordinate Reference Systems , 2009. [http://inspire.ec.europa.eu/documents/Data\\_Specifications/INSPIRE\\_Specification\\_CRS\\_v3.0.pdf](http://inspire.ec.europa.eu/documents/Data_Specifications/INSPIRE_Specification_CRS_v3.0.pdf). 36, 39, 179
- [33] Phil Archer, Stijn Goedertier, and Nikolaos Loutas. Study on persistent URIs, with identification of best practices and recommendations on the topic for the MSs and the EC. Deliverable, 2012. 36
- [34] Auguste Ghislain Atemezing and Raphaël Troncy. Comparing Vocabularies for Representing Geographical Features and Their Geometry. In *5<sup>th</sup> International Terra Cognita Workshop*, Boston, USA, 2012. 38, 40
- [35] International Organization for Standardization . ISO 19125-1, Geographic information- Simple feature access - Part 1: Common architecture, 2004. 42
- [36] EU ISA Programme Core Vocabularies Working Group (Location Task Force). Isa programme location vocabulary. W3C document, 2013. <http://www.w3.org/ns/locn>. 46, 71

- [37] Kostas Patroumpas, Michalis Alexakis, Giorgos Giannopoulos, and Spiros Athanasiou. Triplegeo: an etl tool for transforming geospatial data into rdf triples. In K. Selçuk Candan, Sihem Amer-Yahia, Nicole Schweikardt, Vassilis Christophides, , and Vincent Leroy, editors, *EDBT/ICDT Workshops*, pages 275–278, 2014. 51
- [38] Jhonny Saavedra, Luis M. Vilches-Blázquez, and Alberto Boada. Cadastral data integration through linked data. In Granell (Eds): Connecting a Digital Europe through Location Huerta, Schade and Place., editors, *Proceedings of the AGILE'2014 International Conference on Geographic Information Science Castellón, June, 3-6, 2014*. <http://hdl.handle.net/10234/98742>. 51
- [39] International Organization for Standardization . ISO 19152, Geographic information: Land Administration Domain Model (LADM), 2012. [http://www.iso.org/iso/iso\\_catalogue/catalogue\\_tc/catalogue\\_detail.htm?csnumber=51206](http://www.iso.org/iso/iso_catalogue/catalogue_tc/catalogue_detail.htm?csnumber=51206). 51
- [40] Fayçal Hamdi, Nathalie Abadie, Bénédicte Bucher, and Abdelfettah Feliachi. Geomrdf: A fine-grained structured representation of geometry in the web. In *Proceedings of the 1st International Workshop on Geospatial Linked Data, 1 September, Leipzig, Germany*, 2014. 51
- [41] E Prud'hommeaux, G Carothers, D Beckett, and T Berners-Lee. Rdf 1.1 turtle—terse rdf triple language. *W3C Recommendation*, 2014. <http://www.w3.org/TR/turtle/>. 52
- [42] Jérôme Euzenat, Pavel Shvaiko, et al. *Ontology matching*, volume 18. Springer, 2007. 54
- [43] A. Jentzsch, R. Isele, and C. Bizer. Silk-generating rdf links while publishing or consuming linked data. In *Poster at the International Semantic Web Conference (ISWC2010), Shanghai*, 2010. 54, 68, 123
- [44] Robert Isele and Christian Bizer. Learning linkage rules using genetic programming. In *Proceedings of the Sixth International Workshop on Ontology Matching*, pages 13–24, 2011. 54
- [45] Axel-Cyrille Ngonga Ngomo and Sören Auer. Limes - a time-efficient approach for large-scale link discovery on the web of data. In *Proceedings of IJCAI*, 2011. 54, 62
- [46] François Scharffe, Jérôme Euzenat, et al. Méthodes et outils pour lier le web des données. In *Actes 17e conférence AFIA-AFRIF sur reconnaissance des formes et intelligence artificielle (RFIA)*, pages 678–685, 2010. 54
- [47] Sébastien Mustière and Thomas Devogele. Matching networks with different levels of detail. *GeoInformatica*, 12(4):435–453, 2008. 54

- [48] Julius Volz, Christian Bizer, Martin Gaedke, and Georgi Kobilarov. Discovering and Maintaining Links on the Web of Data. In *8<sup>th</sup> International Semantic Web Conference (ISWC)*, Washington DC, USA, 2009. 54, 57
- [49] Volker Walter and Dieter Fritsch. Matching spatial data sets: a statistical approach. *International Journal of Geographical Information Science*, 13(5):445–473, 1999. 54
- [50] Chantal Reynaud and Brigitte Safar. Techniques structurelles d’alignement pour portails web. *Revue RNTI W3, Fouille du Web*, 2007. 54
- [51] Fayçal Hamdi, Chantal Reynaud, and Brigitte Safar. Pattern-based mapping refinement. In *Knowledge Engineering and Management by the Masses*, pages 1–15. Springer, 2010. 54
- [52] Abdelfettah Feliachi, Nathalie Abadie, Fayçal Hamdi, and Ghislain Auguste Atemezing. Interlinking and visualizing linked open data with geospatial reference data. In *OM’13*, pages 237–238, 2013. 56
- [53] W3C SPARQL Working Group. Sparql query language for rdf. W3C Recommendation 21 March 2013, 2013. <http://www.w3.org/TR/sparql11-overview/>. 57, 66
- [54] Carlos Buil Aranda, Aidan Hogan, Jürgen Umbrich, and Pierre-Yves Vandenbussche. SPARQL web-querying infrastructure: Ready for action? In Harith Alani, Lalana Kagal, Achille Fokoue, Paul T. Groth, Chris Biemann, Josiane Xavier Parreira, Lora Aroyo, Natasha F. Noy, Chris Welty, and Krzysztof Janowicz, editors, *The Semantic Web - ISWC 2013 - 12th International Semantic Web Conference, Sydney, NSW, Australia, October 21-25, 2013, Proceedings, Part II*, volume 8219 of *Lecture Notes in Computer Science*, pages 277–293. Springer, 2013. 58
- [55] M Ashworth. Information technology, database languages, sql multimedia, and application packages, part 3: Spatial, iso. *IEC 13249*, 3, 1999. 58
- [56] Antonin Guttman. *R-trees: a dynamic index structure for spatial searching*, volume 14. ACM, 1984. 58
- [57] Robert Battle and Dave Kolas. Enabling the geospatial semantic web with parliament and geosparql. *Semantic Web*, 3(4):355–370, 2012. 59
- [58] George Garbis, Kostis Kyziros, and Manolis Koubarakis. Geographica: A benchmark for geospatial RDF stores (long version). In *The Semantic Web - ISWC 2013 - 12th International Semantic Web Conference, Sydney, NSW, Australia, October 21-25, 2013, Proceedings, Part II*, pages 343–359, 2013. [http://dx.doi.org/10.1007/978-3-642-41338-4\\_22](http://dx.doi.org/10.1007/978-3-642-41338-4_22). 59

- [59] Kostis Kyziarakos, Manos Karpathiotakis, and Manolis Koubarakis. Strabon: a semantic geospatial dbms. In *The Semantic Web–ISWC 2012*, pages 295–311. Springer, 2012. 59
- [60] Bert Van Nuffelen, Valentina Janev, Michael Martin, Vuk Mijovic, and Sebastian Tramp. Supporting the linked data life cycle using an integrated tool stack. In Sören Auer, Volha Bryl, and Sebastian Tramp, editors, *Linked Open Data – Creating Knowledge Out of Interlinked Data*, Lecture Notes in Computer Science, pages 108–129. Springer International Publishing, 2014. 62
- [61] Tommaso Soru and Axel-Cyrille Ngonga Ngomo. Rapid execution of weighted edit distances. In *Proceedings of the Ontology Matching Workshop*, 2013. 62
- [62] Axel-Cyrille Ngonga Ngomo. Orchid - reduction-ratio-optimal computation of geo-spatial distances for link discovery. In *Proceedings of ISWC 2013*, 2013. 63, 69
- [63] GeoKnow Project. Spatial mapping framework for enriching rdf datasets with geo-spatial information. Manual-Deliverable, 2014. <https://github.com/GeoKnow/GeoLift/blob/master/GeoLiftManual/GeoLiftManual.pdf>. 63
- [64] François Scharffe, Ghislain Atemezing, Raphaël Troncy, Fabien Gandon, Serena Villata, Bénédicte Bucher, Fayçal Hamdi, Laurent Bihanic, Gabriel Képéklia, Franck Cotton, Jérôme Euzenat, Zhengjie Fan, Pierre-Yves Vandenbussche, and Bernard Vatant. Enabling linked-data publication with the datalift platform. In *26th Conference on Artificial Intelligence (AAAI-12)*, 2012. 63, 138, 158
- [65] Pierre-Yves Vandenbussche, Bernard Vatant, and L. Rozat. Linked open vocabularies: an initiative for the web of data. In *QetR Workshop*, Chambery, France, 2011. 64, 99
- [66] Keith Alexander and Michael Hausenblas. Describing linked datasets: on the design and usage of void, the vocabulary of interlinked datasets. In *In Linked Data on the Web Workshop (LDOW 09), in conjunction with 18th International World Wide Web Conference (WWW 09)*. Citeseer, 2009. 73
- [67] Marti A. Hearst. *Search User Interfaces*. Cambridge University Press, New York, NY, USA, 1st edition, 2009. 81
- [68] Marti Hearst, Ame Elliott, Jennifer English, Rashmi Sinha, Kirsten Swearinged, and Ka-Ping Yee. Finding the flow in web site search, 2002. 81
- [69] Ulrich Atz, Tom Heath, Michael Heil, Jack Hardinges, and Jamie Fawcett. Best practice visualisation, dashboard and key figures report. WP2: Research studies and stakeholder analysis, 2014. <http://project.opendatamonitor.eu>. 82

- [70] Lars Grammel and Margaret-Anne Storey. Choosel: Web-based visualization construction and coordination for information visualization novices. Poster: IEEE Information Visualization Conference, 2010. 82
- [71] IBM Research. Many eyes, 2010. <http://www-958.ibm.com/software/data/cognos/maneyes/>. 82
- [72] Mike Bostock. Data-driven documents, 2012. <http://d3js.org/>. 83
- [73] Tetherless Constellation. How to use google visualization api. Wiki notes, 2012. [http://data-gov.tw.rpi.edu/wiki/How\\_to\\_use\\_Google\\_Visualization\\_API](http://data-gov.tw.rpi.edu/wiki/How_to_use_Google_Visualization_API). 83
- [74] Aba-Sah Dadzie and Matthew Rowe. Approaches to visualising linked data: A survey. *Semantic Web*, 2(2):89–124, 2011. 84
- [75] Dave Reynolds, Jeni Tennison, Leigh Dodds, and et. al. Linked data api, 2012. <https://code.google.com/p/linked-data-api/>. 84
- [76] G. Martin Skjæveland. Sgvizler: A javascript wrapper for easy visualization of sparql result sets. In *9<sup>th</sup> Extended Semantic Web Conference (ESWC'12)*, 2012. 85, 104
- [77] Claus Stadler, Michael Martin, and Sören Auer. Exploring the Web of Spatial Data with Facete. In *Companion proceedings of 23rd International World Wide Web Conference (WWW)*, pages 175–178, 2014. 85, 88
- [78] Alvaro Graves. Creation of visualizations based on linked data. In *Proceedings of the 3rd International Conference on Web Intelligence, Mining and Semantics*, WIMS '13, pages 41:1–41:12, New York, NY, USA, 2013. ACM. 86
- [79] Jakub Klímek, Jirí Helmich, and Martin Necaský. Payola: Collaborative linked data analysis and visualization framework. In *The Semantic Web: ESWC 2013 Satellite Events - ESWC 2013 Satellite Events, Montpellier, France, May 26-30, 2013, Revised Selected Papers*, pages 147–151, 2013. 86
- [80] Josep Maria Brunetti, Sören Auer, Roberto García, Jakub Klímek, and Martin Nečaský. Formal linked data visualization model. In *Proceedings of International Conference on Information Integration and Web-based Applications & Services*, IIWAS '13, pages 309:309–309:318, New York, NY, USA, 2013. ACM. 86
- [81] Ed H. Chi. A taxonomy of visualization techniques using the data state reference model. In *Proceedings of the IEEE Symposium on Information Visualization 2000*, INFOVIS '00, pages 69–, Washington, DC, USA, 2000. IEEE Computer Society. 86
- [82] Aba-Sah Dadzie and Matthew Rowe. Approaches to visualising linked data: A survey. *Semantic Web Journal*, 2(2):89–124, 2011. 87

- [83] J. Klimek, J. Helmich, and M. Neasky. Application of the Linked Data Visualization Model on Real World Data from the Czech LOD Cloud. In *6<sup>th</sup> International Workshop on the Linked Data on the Web (LDOW'14)*, 2014. 87
- [84] Percy E Salas, Michael Martin, Fernando Maia Da Mota, Karin Breitman, Sören Auer, and Marco A Casanova. Publishing statistical data on the web. In *Proceedings of 6<sup>th</sup> International IEEE Conference on Semantic Computing*, IEEE 2012. IEEE, 2012. 87
- [85] Alexander de Leon, Filip Wisniewski, Boris Villazón-Terrazas, and Oscar Corcho. Map4rdf - Faceted Brower for Geospatial Datasets. In *Using Open Data: policy modeling, citizen empowerment, data journalism (PMOD'12)*, 2012. 87
- [86] Danh Le Phuoc, Axel Polleres, Christian Morbidoni, Manfred Hauswirth, and Giovanni Tummarello. Rapid semantic web mashup development through semantic web pipes. In *18<sup>th</sup> International World Wide Web Conference (WWW'09)*, Madrid, Spain, 2009. 87
- [87] G. Atemezing and R. Troncy. Tools for visualization (v.1.2). Deliverables - D.6.2 of DataLift project, 2012. 88
- [88] Stuart K. Card, Jock D. Mackinlay, and Ben Shneiderman, editors. *Readings in information visualization: using vision to think*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1999. 92
- [89] Emmanuel Pietriga. Isaviz: A visual authoring tool for rdf, 2004. <http://www.w3.org/2001/11/IsaViz/>. 92
- [90] Tim Berners-lee, Yuhsin Chen, Lydia Chilton, Dan Connolly, Ruth Dhanaraj, James Hollenbach, Adam Lerer, and David Sheets. Tabulator: Exploring and analyzing linked data on the semantic web. In *In Proceedings of the 3rd International Semantic Web User Interaction Workshop*, 2006. 92
- [91] G. Atemezing and R. Troncy. Usage scenarii for applications (v.1.1). Deliverables- D.6.1 of DataLift project, 2012. 93
- [92] Colin Ware. *Information Visualization, Second Edition: Perception for Design*. Morgan Kaufmann Publishers Inc.; 2 edition (April 21, 2004), San Francisco, CA, USA, 2014. 97
- [93] B Shneiderman. The Eyes Have It: A Task by Data Type Taxonomy for Information Visualizations. In *IEEE Symposium on Visual Languages '96, IEEE, Los Alamos, CA (September 1996)*, pages 336–343, 1996. 97, 99
- [94] Jeni Tennison. Guest post: A developers' guide to the linked data apis, July 2010. <http://data.gov.uk/blog/guest-post-developers-guide-linked-data-apis-jeni-tennison>. 97

- [95] Pierre-Yves Vandenbussche and Bernard Vatant. Metadata Recommendations For Linked Open Vocabularies. OKFN, 2012. [http://lov.okfn.org/dataset/lov/Recommendations\\_Vocabulary\\_Design.pdf](http://lov.okfn.org/dataset/lov/Recommendations_Vocabulary_Design.pdf). 102
- [96] Carlos Buil-Aranda, Aidan Hogan, Jürgen Umbrich, and Pierre-Yves Vandenbussche. Sparql web-querying infrastructure: Ready for action? In *The Semantic Web–ISWC 2013*, pages 277–293. Springer, 2013. 105
- [97] Emmanuel Pietriga, Christian Bizer, David Karger, and Ryan Lee. Fresnel: A Browser-Independent Presentation Vocabulary for RDF. In *5<sup>th</sup> International Semantic Web Conference (ISWC'06)*, pages 158–171, 2006. 107
- [98] Mike Bergman. Deconstructing the Google Knowledge Graph.  
<http://www.mkbergman.com/1009/deconstructing-the-google-knowledge-graph.html>. 107
- [99] Raphaël Troncy, Bartosz Malocha, and André Fialho. Linking Events with Media. In *6<sup>th</sup> International Conference on Semantic Systems (I-SEMANTICS'10)*, Graz, Austria, 2010. 120
- [100] Pieter Colpert, Anastasia Dimou, Ghislain Auguste Atemezing, and Raphaël Troncy. Technical toolset for open data competitions apps for europe. Technical report, Ghent University, EURECOM, 2014. [http://www.appsforeurope.eu/sites/default/files/TechnicaltoolsetforopendatacompetitionsAppsforEurope\\_0.pdf](http://www.appsforeurope.eu/sites/default/files/TechnicaltoolsetforopendatacompetitionsAppsforEurope_0.pdf). 124
- [101] Bernadette Hyland and David Wood. The joy of data - cookbook for publishing linked government data on the web, 2011. [http://www.w3.org/2011/gld/wiki/Linked\\_Data\\_Cookbook](http://www.w3.org/2011/gld/wiki/Linked_Data_Cookbook). 137
- [102] Michael Hausenblas and Richard Cyganiak. Linked data life cycles, 2012. <http://linked-data-life-cycles.info/>. 137
- [103] Boris Villazón-Terrazas; et al. Methodological guidelines for publishing government linked data. [http://link.springer.com/chapter/10.1007/978-1-4614-1767-5\\_2](http://link.springer.com/chapter/10.1007/978-1-4614-1767-5_2). 138
- [104] Christian Bizer, Tom Heath, and Tim Berners-Lee. Linked Data - The Story So Far. *International Journal on Semantic Web and Information Systems*, 5:1–22, 2009. 144, 158
- [105] Keith Alexander, Richard Cyganiak, Michael Hausenblas, and Jun Zhao. Describing linked datasets. In *2<sup>nd</sup> Workshop on Linked Data on the Web (LDOW)*, Madrid, Spain, 2009. 144
- [106] Bernard Vatant and Pierre-Yves Vandenbussche. Catalogue de Vocabulaires. Datalift, D2.2, 2013. <http://datalift.org/en/node/18>. 151

- [107] Eric Prud'hommeaux and Carlos Buil-Aranda. SPARQL 1.1 Federated Query. W3C Recommendation, 2013. <http://www.w3.org/TR/sparql11-federated-query/>. 151
- [108] S. M. Ross. A First Course in Probability, 2002. 159
- [109] R. Meymandpour and J. G. Davis. Ranking Universities Using Linked Open Data. In *5<sup>th</sup> International Workshop on the Linked Data on the Web (LDOW'13)*, 2013. 159
- [110] Johann Schaible, Thomas Gottron, and Ansgar Scherp. Survey on Common Strategies of Vocabulary Reuse in Linked Open Data Modeling. In *11<sup>th</sup> Extended Semantic Web Conference (ESWC'14)*, pages 457–472, 2014. 165
- [111] Tom Heath and Christian Bizer. *Linked Data: Evolving the Web into a Global Data Space*. Morgan & Claypool, 2011. 169
- [112] Guido Governatori, Antonino Rotolo, Serena Villata, and Fabien Gandon. One license to compose them all - a deontic logic approach to data licensing on the web of data. In *International Semantic Web Conference (1)*, volume 8218 of *Lecture Notes in Computer Science*, pages 151–166. Springer, 2013. 169, 170, 171, 172, 174, 175
- [113] Víctor Rodríguez-Doncel, Asunción Gómez-Pérez, and Nandana Mihindukulasooriya. Rights declaration in linked data. In Hartig et al. [127]. 170, 171
- [114] G. R. Gangadharan, V. D'Andrea, R. Iannella, and M. Weiss. Odrl service licensing profile (ODRL-S). In *Proceedings of Virtual Goods*, 2007. 171
- [115] G. R. Gangadharan, Michael Weiss, Vincenzo D'Andrea, and Renato Iannella. Service license composition and compatibility analysis. In *Proceedings of ICSOC, LNCS 4749*, pages 257–269. Springer, 2007. 171
- [116] Nadia Nadah, Mélanie Dulong de Rosnay, and Bruno Bachimont. Licensing digital content with a generic ontology: escaping from the jungle of rights expression languages. In *Proceedings of ICAIL-the Eleventh International Conference on Artificial Intelligence and Law, Stanford Law School, Stanford, California, USA*, pages 65–69. ACM, 2007. 171
- [117] Hong Linh Truong, G. R. Gangadharan, Marco Comerio, Schahram Dustdar, and Flavio De Paoli. On analyzing and developing data contracts in cloud-based data marketplaces. In *Proceedings of APSCC, IEEE*, pages 174–181, 2011. 171
- [118] Markus Krötzsch and Sebastian Speiser. ShareAlike Your Data: Self-referential Usage Policies for the Semantic Web. In *Proceedings of ISWC, LNCS 7031*, pages 354–369. Springer, 2011. 171

- [119] Thomas F. Gordon. Analyzing open source license compatibility issues with Carneades. In *Proceedings of ICAIL*, pages 51–55. ACM, 2011. 172
- [120] V. Rodriguez-Doncel, M.C. Suarez-Figueroa, A. Gómez-Pérez, and M. Poveda Villalón. Licensing patterns for linked data. In *Proceedings of the 4th International Workshop on Ontology Patterns*, 2013. 172
- [121] V. Rodriguez-Doncel, M.C. Suarez Figueroa, A. Gómez-Pérez, and M. Poveda Villalón. License linked data resources pattern. In *Proceedings of the 4th International Workshop on Ontology Patterns*, 2013. 172
- [122] Riccardo Pucella and Vicky Weissman. A logic for reasoning about digital rights. In *Proceedings of CSFW*, pages 282–294. IEEE, 2002. 172
- [123] Elena Cabrio, Alessio Palmero Aprosio, and Serena Villata. These are your rights: A natural language processing approach to automated rdf licenses generation. In *ESWC2014, LNCS*, 2014. 172, 173
- [124] Ho-Pun Lam and Guido Governatori. The making of SPINdle. In *Proceedings of RuleML, LNCS 5858*, pages 315–322. Springer, 2009. 174
- [125] Antonino Rotolo, Serena Villata, and Fabien Gandon. A deontic logic semantics for licenses composition in the web of data. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Law*, pages 111–120. ACM, 2013. 181
- [126] Ruben Verborgh, Olaf Hartig, Ben De Meester, Gerald Haesendonck, Laurens De Vocht, Miel Vander Sande, Richard Cyganiak, Pieter Colpaert, Erik Manens, and Rik Van de Walle. Low-cost queryable linked data through triple pattern fragments. In *Proceedings of the 13th International Semantic Web Conference: Posters and Demos*, 2014. 182
- [127] Olaf Hartig, Juan Sequeda, Aidan Hogan, and Takahide Matsutsuka, editors. *Proceedings of the Fourth International Workshop on Consuming Linked Data, COLD 2013, Sydney, Australia, October 22, 2013*, volume 1034 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2013. 194
- [128] Bernardo Cuenca Grau, Ian Horrocks, Boris Motik, and Ulrike Sattler, editors. *Proceedings of the 22nd International Workshop on Description Logics (DL 2009), Oxford, UK, July 27-30, 2009*, volume 477 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2009. 200
- [129] Sadok Ben Yahia and Jean-Marc Petit, editors. *Extraction et gestion des connaissances (EGC'2010), Actes, 26 au 29 janvier 2010, Hammamet, Tunisie*, volume RNTI-E-19 of *Revue des Nouvelles Technologies de l'Information*. Cépaduès-Éditions, 2010. 201

- [130] Birte Glimm and David Huynh, editors. *Proceedings of the ISWC 2012 Posters & Demonstrations Track, Boston, USA, November 11-15, 2012*, volume 914 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2012. 202
- [131] P.Y. Vandenbussche, B. Vatant, and L. Rozat. Linked open vocabularies: an initiative for the web of data. In *QeR Workshop, Chambéry, France.*, 2011.
- [132] Guido Governatori, Antonino Rotolo, Serena Villata, and Fabien G. One license to compose them all a deontic logic approach to data licensing on the web of data, 2012.
- [133] Pierre-Yves Vandenbussche and Bernard Vatant. Metadata Recommendations For Linked Open Vocabularies. OKFN, 2012. [http://lov.okfn.org/dataset/lov/Recommendations\\_Vocabulary\\_Design.pdf](http://lov.okfn.org/dataset/lov/Recommendations_Vocabulary_Design.pdf).
- [134] Krzysztof Janowicz, Sven Schade, Arne Bröring, Carsten Kessler, Christoph Stasch, Patrick Maué, and Thorsten Diekhof. A transparent semantic enablement layer for the geospatial web. In *2nd International Terra Cognita Workshop*, 2009.
- [135] International Organization for Standardization (TC 211). ISO 19107: Geographic information - Spatial Schema., 2003.
- [136] International Organization for Standardization (TC 211). ISO 19109: Geographic information - Rules for application schema., 2005.
- [137] International Organization for Standardization (TC 211). ISO 19111: Geographic information - Spatial referencing by coordinates, 2007.
- [138] Bénédicte Bucher, Nathalie Abadie, and Ghislain Auguste Atemezing. Modélisation de connaissances spécifiques: information spatiale. Rapport de Recherche. Délivrable WP2.T.2.3- Projet Datalift, 2013.
- [139] Krzysztof Janowicz, Pascal Hitzler, Benjamin Adams, Dave Kolas, and Charles Vardeman II. Five stars of linked data vocabulary use. *Semantic Web Journal*, 2014. <http://geog.ucsb.edu/~jano/swj653.pdf>.
- [140] W3C Semantic Web Interest Group (SWIG). Wgs84 rdf geoposition vocabulary, 2004. <http://www.w3.org/2003/01/geo/>.
- [141] Tim Berners-lee, Yuhsin Chen, Lydia Chilton, Dan Connolly, Ruth Dhanaraj, James Hollenbach, Adam Lerer, and David Sheets. Tabulator: Exploring and analyzing linked data on the semantic web. In *In Proceedings of the 3rd International Semantic Web User Interaction Workshop*, 2006.
- [142] S. Harris and A Seaborne. SPARQL 1.1 Query Language, 2013. <http://www.w3.org/TR/sparql11-query/>.
- [143] Vladimir Geroimenko and Chaomei Chen. Visualising the semantic web, 2006.

- [144] Josep Maria Brunetti, Soren Auer, and Roberto Garcia. The linked data visualization model. In *In Proceedings of the 11th International Semantic Web Conference*, 2012.
- [145] Antoine Isaac, William Waites, Jeff Young, and Marcia Zeng. Library Linked Data Incubator Group: Datasets, Value Vocabularies, and Metadata Element Sets. W3C Incubator Group Report, 2011. <http://www.w3.org/2005/Incubator/ld/XGR-ld-vocabdataset/>.
- [146] Sean Bechhofer and Alistair Miles. SKOS Simple Knowledge Organization System Reference. W3C Recommendation, 2009. <http://www.w3.org/TR/skos-reference/>.
- [147] Bernadette Hyland, Ghislain Auguste Atemezing, Michael Pendleton, and Bipav Srivastava. Linked Data Glossary. W3C Working Group Note, 2013. <http://www.w3.org/TR/ld-glossary/>.
- [148] Tom Heath and Christian Bizer. *Linked Data: Evolving the Web into a Global Data Space: Theory and Technology*. Morgan & Claypool Publishers, 2011.
- [149] Thomas Steiner and Stefan Mirea. SEKI@home or Crowdsourcing an Open Knowledge Graph. In *1<sup>st</sup> International Workshop on Knowledge Extraction & Consolidation from Social Media (KECSM'12)*, Boston, USA, 2012.
- [150] Ashutosh Jadhav, Hemant Purohit, Pramod Ananthram, Ajith Ranabahu, Vinh Nguyen, Pablo Mendes, Alan Gary Smith, Michael Cooney, and Amit Sheth. Twitris 2.0 : Semantically Empowered System for Understanding Perceptions From Social Data. In *Semantic Web Challenge at the 9<sup>th</sup> International Semantic Web Conference (ISWC'10)*, Shanghai, China, 2010.
- [151] Hila Becker, Mor Naaman, and Luis Gravano. Learning similarity metrics for event identification in social media. In *3<sup>rd</sup> ACM International Conference on Web Search and Data Mining*, pages 291–300, New York, NY, USA, 2010.
- [152] Smitashree Choudhury and John G. Breslin. Extracting Semantic Entities and Events from Sports Tweets. In *Making Sense of Microposts (#MSM2011)*, 2011.
- [153] Ralph Grishman and Beth Sundheim. Message Understanding Conference-6: a brief history. In *16<sup>th</sup> International Conference on Computational linguistics (COLING'96)*, pages 466–471, Copenhagen, Denmark, 1996.
- [154] Xueliang Liu, Raphaël Troncy, and Benoît Huet. Using Social Media to Identify Events. In *3<sup>rd</sup> Workshop on Social Media (WSM'11)*, Scottsdale, Arizona, USA, 2011.
- [155] Pablo Mendes, Alexandre Passant, and Pavan Kapanipathi. TWARQL: Tapping into the Wisdom of the Crowd. In *6<sup>th</sup> International Conference on Semantic Systems (I-SEMANTICS'10)*, Graz, Austria, 2010.

- [156] Knud Möller, Tom Heath, Siegfried Handschuh, and John Domingue. Recipes for Semantic Web dog food - The ESWC and ISWC metadata projects. In *6<sup>th</sup> International Semantic Web Conference (ISWC'07)*, pages 802–815, Busan, Korea, 2007.
- [157] Giuseppe Rizzo and Raphaël Troncy. NERD: A Framework for Evaluating Named Entity Recognition Tools in the Web of Data. In *10<sup>th</sup> International Semantic Web Conference (ISWC'11), Demo Session*, pages 1–4, Bonn, Germany, 2011.
- [158] Matthew Rowe and Milan Stankovic. Aligning Tweets with Events: Automation via Semantics. *Semantic Web Journal*, 2011.
- [159] R. Shaw, R. Troncy, and L. Hardman. LODE: Linking Open Descriptions of Events. In *4<sup>th</sup> Asian Semantic Web Conference (ASWC'09)*, pages 153–167, Shanghai, China, 2009.
- [160] Milan Stankovic. Modeling Online Presence. In *1<sup>st</sup> Social Data on the Web Workshop (SDoW'08)*, 2008.
- [161] Katrin Weller, Evelyn Drge, and Cornelius Puschmann. Citation Analysis in Twitter: Approaches for Defining and Measuring Information Flows within Tweets during Scientific Conferences. In *Making Sense of Microposts (#MSM2011)*, pages 1–12, 2011.
- [162] Alan Ritter, Sam Clark, Mausam, and Oren Etzioni. Named entity recognition in tweets: An experimental study. In *EMNLP*, pages 1524–1534, 2011.
- [163] Guido Governatori, Ho-Pun Lam, Antonino Rotolo, Serena Villata, and Fabien Gandon. Heuristics for licenses composition. In *Proceedings of JURIX*, pages 77–86. IOS Press, 2013.
- [164] Kevin D. Ashley, editor. *Legal Knowledge and Information Systems - JURIX 2013: The Twenty-Sixth Annual Conference, December 11-13, 2013, University of Bologna, Italy*, volume 259 of *Frontiers in Artificial Intelligence and Applications*. IOS Press, 2013.
- [165] N. Cristianini and J. Shawe-Taylor. *An Introduction to Support Vector Machines and Other Kernel-based Learning Methods*. Cambridge University Press, 2010.
- [166] Paul T. Groth, Yolanda Gil, James Cheney, and Simon Miles. Requirements for provenance on the web. *IJDC*, 7(1):39–56, 2012.
- [167] Michael J. Maher, Andrew Rock, Grigoris Antoniou, David Billington, and Tristan Miller. Efficient defeasible reasoning systems. *International Journal of Artificial Intelligence Tools*, 10:483–501, 2001.

- [168] Peter Buneman, Sanjeev Khanna, and Wang Chiew Tan. Why and where: A characterization of data provenance. In Jan Van den Bussche and Victor Vianu, editors, *Proceedings of Database Theory - ICDT 2001, 8th International Conference, London, UK, January 4-6, 2001*, volume 1973 of *Lecture Notes in Computer Science*, pages 316–330. Springer, 2001.
- [169] Marco Comerio. *Web Service Contracts: Specification, Selection and Composition*. PhD thesis, University of Milano-Bicocca, 2009.
- [170] Guido Governatori. On the relationship between carneades and defeasible logic. In *Proceedings of ICAIL*, pages 31–40. ACM, 2011.
- [171] Kevin D. Ashley and Tom M. van Engers, editors. *The 13th International Conference on Artificial Intelligence and Law, Proceedings of the Conference, June 6-10, 2011, Pittsburgh, PA, USA*. ACM, 2011.
- [172] G. R. Gangadharan, Hong Linh Truong, Martin Treiber, Vincenzo D’Andrea, Schahram Dustdar, Renato Iannella, and Michael Weiss. Consumer-specified service license selection and composition. In *Proceedings of ICCBSS, IEEE*, pages 194–203, 2008.
- [173] *Seventh International Conference on Composition-Based Software Systems (ICCBSS 2008), February, 25-29, 2008, Madrid, Spain, Proceedings*. IEEE Computer Society, 2008.
- [174] Christian Bizer, Tom Heath, and Tim Berners-Lee. Linked data - the story so far. *Int. J. Semantic Web Inf. Syst.*, 5(3):1–22, 2009.
- [175] James J. Park, Christos Nikolaou, and Jiannong Cao, editors. *2011 IEEE Asia-Pacific Services Computing Conference, APSCC 2011, Jeju, Korea (South), December 12-15, 2011*. IEEE, 2011.
- [176] Bernd J. Krämer, Kwei-Jay Lin, and Priya Narasimhan, editors. *Service-Oriented Computing - ICSOC 2007, Fifth International Conference, Vienna, Austria, September 17-20, 2007, Proceedings*, volume 4749 of *Lecture Notes in Computer Science*. Springer, 2007.
- [177] Lora Aroyo, Chris Welty, Harith Alani, Jamie Taylor, Abraham Bernstein, Lalana Kagal, Natasha Fridman Noy, and Eva Blomqvist, editors. *The Semantic Web - ISWC 2011 - 10th International Semantic Web Conference, Bonn, Germany, October 23-27, 2011, Proceedings, Part I*, volume 7031 of *Lecture Notes in Computer Science*. Springer, 2011.
- [178] Cássia Trojahn dos Santos and Jérôme Euzenat. Consistency-driven argumentation for alignment agreement. In Pavel Shvaiko, Jérôme Euzenat, Fausto Giunchiglia, Heiner Stuckenschmidt, Ming Mao, and Isabel F. Cruz, editors,

- Proceedings of the 5th International Workshop on Ontology Matching (OM-2010), Shanghai, China, November 7, 2010*, volume 689 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2010.
- [179] Nicoletta Fornara and Marco Colombetti. Ontology and time evolution of obligations and prohibitions using semantic web technology. In Matteo Baldoni, Jamal Bentahar, M. Birna van Riemsdijk, and John Lloyd, editors, *Declarative Agent Languages and Technologies VII, 7th International Workshop, DALT 2009, Budapest, Hungary, May 11, 2009. Revised Selected and Invited Papers*, volume 5948 of *Lecture Notes in Computer Science*, pages 101–118. Springer, 2010.
  - [180] Moreau et al. *The Open Provenance Model Core Specification (v1.1)*, 2009. <http://eprints.ecs.soton.ac.uk/18332/1/opm.pdf>.
  - [181] Jun Zhao. *Guide to the Open Provenance Vocabulary*, 2010. url <http://purl.org/net/opmv/guide>.
  - [182] P. Miller, R. Styles, and T. Heath. Open data commons, a license for open data. In *LDOW*, 2008.
  - [183] Renato Iannella. *Open Digital Rights Language (ODRL)*, 2002. <http://odrl.net/1.1/ODRL-11.pdf>.
  - [184] Hal Abelson, Ben Adida, Mike Linksvayer, and Nathan Yergler. *ccREL: The Creative Commons Rights Expression Language*, 2008.
  - [185] *ODC Public Domain Dedication and License*, 2008. [http://download.opencontentlawyer.com/ODC\\_PDDL.pdf](http://download.opencontentlawyer.com/ODC_PDDL.pdf).
  - [186] Rui Zhang, Alessandro Artale, Fausto Giunchiglia, and Bruno Crispo. Using description logics in relation based access control. In Grau et al. [128].
  - [187] Owen Sacco and Alexandre Passant. A Privacy Preference Ontology (PPO) for Linked Data. In *Proceedings of the 4th Workshop about Linked Data on the Web (LDOW-2011)*, 2011.
  - [188] James Hollenbach, Joe Presbrey, and Tim Berners-Lee. Using RDF Metadata To Enable Access Control on the Social Semantic Web. In *Proceedings of the Workshop on Collaborative Construction, Management and Linking of Structured Knowledge (CK-2009)*, 2009.
  - [189] Fabian Abel, Juri Luca De Coi, Nicola Henze, Arne Wolf Koesling, Daniel Krause, and Daniel Olmedilla. Enabling advanced and context-dependent access control in rdf stores. In *Proceedings of the 6th International Semantic Web Conference (ISWC-2007)*, LNCS 4825, pages 1–14, 2007.

- [190] Karl Aberer, Key-Sun Choi, Natasha Fridman Noy, Dean Allemang, Kyung-Il Lee, Lyndon J. B. Nixon, Jennifer Golbeck, Peter Mika, Diana Maynard, Riichiro Mizoguchi, Guus Schreiber, and Philippe Cudré-Mauroux, editors. *The Semantic Web, 6th International Semantic Web Conference, 2nd Asian Semantic Web Conference, ISWC 2007 + ASWC 2007, Busan, Korea, November 11-15, 2007*, volume 4825 of *Lecture Notes in Computer Science*. Springer, 2007.
- [191] Hannes Muhleisen, Martin Kost, and Johann-Christoph Freytag. SWRL-based Access Policies for Linked Data. In *Proceedings of the 2nd Workshop on Trust and Privacy on the Social and Semantic Web (SPOT2010)*, 2010.
- [192] Stephanie Stroka, Sebastian Schaffert, and Tobias Burger. Access Control in the Social Semantic Web - Extending the idea of FOAF+SSL in KiWi. In *Proceedings of the 2nd Workshop on Trust and Privacy on the Social and Semantic Web (SPOT2010)*, 2010.
- [193] Christian Bizer, Tom Heath, and Tim Berners-lee. Linked data-the story so far. *International Journal on Semantic Web and Information Systems*, 5(3):1–22, 2009.
- [194] Juri Luca De Coi, Daniel Olmedilla, Sergej Zerr, Piero A. Bonatti, and Luigi Sauro. A trust management package for policy-driven protection & personalization of web content. In *Proceedings of the 9th IEEE International Workshop on Policies for Distributed Systems and Networks (POLICY 2008)*, 2-4 June 2008, Palisades, New York, USA, pages 228–230. IEEE Computer Society, 2008.
- [195] Olaf Hartig. Querying trust in rdf data with tsql. In *Proceedings of the 6th European Semantic Web Conference (ESWC-2009)*, LNCS 5554, pages 5–20, 2009.
- [196] Lora Aroyo, Paolo Traverso, Fabio Ciravegna, Philipp Cimiano, Tom Heath, Eero Hyvönen, Riichiro Mizoguchi, Eyal Oren, Marta Sabou, and Elena Paslaru Bontas Simperl, editors. *The Semantic Web: Research and Applications, 6th European Semantic Web Conference, ESWC 2009, Heraklion, Crete, Greece, May 31-June 4, 2009, Proceedings*, volume 5554 of *Lecture Notes in Computer Science*. Springer, 2009.
- [197] Jeremy Carroll, Christian Bizer, Patrick Hayes, and Patrick Stickler. Named graphs. *J. Web Sem.*, 3(4):247–267, 2005.
- [198] Michel Buffa, Catherine Faron-Zucker, and Anna Kolomoyskaya. Gestion sémantique des droits d'accès au contenu : l'ontologie AMO. In Yahia and Petit [129], pages 471–482.

- [199] Fausto Giunchiglia, Rui Zhang, and Bruno Crispo. Ontology driven community access control. In *Proceedings of the 1st Workshop on Trust and Privacy on the Social and Semantic Web (SPOT-2009)*, 2009.
- [200] John Breslin, Alexandre Passant, and Stefan Decker. *The Social Semantic Web*. Springer, Heidelberg, 2009.
- [201] Barbara Carminati, Elena Ferrari, Raymond Heatherly, Murat Kantarcioglu, and Bhavani M. Thuraisingham. Semantic web-based social network access control. *Computers & Security*, 30(2-3):108–115, 2011.
- [202] Paul Miller, Rob Styles, and Tom Heath. Open data commons, a license for open data. In *Proceedings of LDOW*, 2008.
- [203] Hal Abelson, Ben Adida, Mike Linksvayer, and Nathan Yergler. ccrel: The creative commons rights expression language. Technical report, Technical report, Creative Commons, 2008. <http://wiki.creativecommons.org/Image:CCrel-1.0.pdf>, 2008.
- [204] Paul Miller, Rob Styles, and Tom Heath. Open data commons, a license for open data. *LDOW*, 369, 2008. [http://download.opencontentlawyer.com/ODC\\_PDDL.pdf](http://download.opencontentlawyer.com/ODC_PDDL.pdf).
- [205] Raul Palma, Jens Hartmann, and Peter Haase. Omv-ontology metadata vocabulary for the semantic web. Technical report, Technical report, Universidad Politécnica de Madrid, University of Karlsruhe, 2008. Version 2.4. Available at <http://omv.ontoware.org>, 2008.
- [206] Richard Raysman, Edward A. Pisacreta, and Kenneth A. Adler. *Intellectual Property Licensing: Forms and Analysis*. Law Journal Press, 1999.
- [207] Serena Villata and Fabien Gandon. Towards licenses compatibility and composition in the web of data. In Glimm and Huynh [130].
- [208] Serena Villata and Fabien Gandon. Licenses compatibility and composition in the web of data. In *Proceedings of COLD*, CEUR Workshop Proceedings 905, 2012.
- [209] Juan Sequeda, Andreas Harth, and Olaf Hartig, editors. *Proceedings of the Third International Workshop on Consuming Linked Data, COLD 2012, Boston, MA, USA, November 12, 2012*, volume 905 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2012.
- [210] Guido Governatori, Antonino Rotolo, Serena Villata, and Fabien Gandon. One license to compose them all. In *The Semantic Web-ISWC 2013*, pages 151–166. Springer, 2013.

- [211] Richard Socher, John Bauer, Christopher D. Manning, and Andrew Y. Ng. Parsing with compositional vector grammars. In *ACL (1)*, pages 455–465, 2013.
- [212] Daniel Gerber and Axel-Cyrille Ngonga Ngomo. Extracting multilingual natural-language patterns for rdf predicates. In *EKAW*, pages 87–96, 2012.



## APPENDIX A

# Installation instructions for the JavaScript plugin

---

The system consists of three components, each having their own code repository:

- Admin-interface: This is the interface that event organizers use to create new events and to get the embed code for embedding the event on their own website. It is available at <https://github.com/EurecomApps4Eu/admin-interface>.
- REST-interface: This interface provides a RESTful service for events and applications, available at <https://github.com/EurecomApps4Eu/rest-interface>
- Embeddable script at <https://github.com/EurecomApps4Eu/event-website>. This repository contains the code that is embedded on the event organizer's own website. It will fetch the event and application information from the REST-interface, and then display the information directly on the event organizers webpage.

The three different components can be installed on different computers, if needed. REST-interface requires Node.js and MongoDB to be installed, whereas Admin-interface and Event-website produces static files (when running the build task) that can be hosted on any web server that is capable of serving static files.

## A.1 Installing and configuring the REST-interface

1. Clone the repository
2. Install dependencies by running the command "npm install"
3. Configure the system by editing "config.js"-file
4. Start the service by running "node app.js {PORT}", and replace {PORT} with the port you want to run the application in (typically port 80 for HTTP)

## A.2 Installing and configuring the Admin-interface

1. Clone the repository
2. Install Node dependencies by running the command "npm install"
3. Install Bower dependencies by running the command "bower install"  
(if bower is not installed run "npm install bower -g")

4. Configure the system by editing file "app/scripts/app.js".

Look for appSettings-part in the file.

5. Build the system with command "grunt build"

(if grunt is not installed, run "npm install grunt -g").

This will produce static HTML/CSS/JS-files that can be hosted on any web server.

### **A.3 Installing and configuring the Event-website**

1. Clone the repository

2. Install Node dependencies by running the command "npm install"

3. Install Bower dependencies by running the command "bower install"

(if bower is not installed run "npm install bower -g")

4. Configure the system by editing file "config.js"

5. Build static HTML/CSS/JS-files by running "grunt all"

## APPENDIX B

# Code source of vocabularies

## B.1 Vocabulary for geometry

```
1 @prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#>.
2 @prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>.
3 @prefix xsd: <http://www.w3.org/2001/XMLSchema#>.
4 @prefix sf: <http://www.opengis.net/ont/sf#>.
5 @prefix xml: <http://www.w3.org/XML/1998/namespace>.
6 @prefix owl: <http://www.w3.org/2002/07/owl#>.
7 @prefix geom: <http://data.ign.fr/def/geometrie#>.
8 @prefix foaf: <http://xmlns.com/foaf/0.1/> .
9 @prefix voaf: <http://purl.org/voccommons/voaf#> .
10 @prefix cc: <http://creativecommons.org/ns#> .
11 @prefix vann: <http://purl.org/vocab/vann/> .
12 @prefix dcterms: <http://purl.org/dc/terms/> .
13 @prefix gsp: <http://www.opengis.net/ont/geosparql#>.
14 @prefix ngeo: <http://geovocab.org/geometry#>.
15 @prefix ignf: <http://data.ign.fr/def/ignf#>.
16
17 # ----- Ontology metadata -----
18 <http://data.ign.fr/def/geometrie> a owl:Ontology;
19   dcterms:description "Ontology for geometries representing shpaes and location of spatial entities"@
20     en;
21   dcterms:title "Ontologie des primitives géométriques"@fr;
22   cc:license <http://www.data.gouv.fr/Licence–Ouverte–Open–Licence> ;
23   cc:license <http://creativecommons.org/licenses/by/2.0/> ;
24   cc:license <http://opendatacommons.org/licenses/by/> ;
25   dcterms:creator <http://recherche.ign.fr/labos/cogit/cv.php?prenom=Nathalie&nom=Abadie> ;
26   dcterms:creator <http://www.eurecom.fr/~atemezin/> ;
27   dcterms:contributor <http://www.eurecom.fr/~troncy/> ;
28   dcterms:contributor <http://data.semanticweb.org/person/bernard–vatant> ;
29   dcterms:contributor <http://recherche.ign.fr/labos/cogit/cv.php?prenom=Bénédicte&nom=
30     Bucher> ;
31   dcterms:issued "2013–06–11"^^xsd:date ;
32   dcterms:modified "2014–08–22"^^xsd:date ;
33   dcterms:publisher <http://fr.dbpedia.org/resource/Institut_national_de_l%27information_g%
34     C3%A9ographique_et_foresti%C3%A8re> ;
35   dcterms:rights "Copyright 2014, IGN" ;
36
37 # ----- Contributors -----
38 <http://data.semanticweb.org/person/bernard–vatant> a foaf:Person.
39 <http://www.eurecom.fr/~atemezin/> a foaf:Person.
40 <http://www.eurecom.fr/~troncy/> a foaf:Person.
41 <http://recherche.ign.fr/labos/cogit/cv.php?prenom=Nathalie&nom=Abadie> a foaf:Person.
42 <http://recherche.ign.fr/labos/cogit/cv.php?prenom=Bénédicte&nom=Bucher> a foaf:Person.
43
44 # ----- Classes definition -----
45
```

```

46 geom:Geometry a owl:Class;
47   rdfs:label "Géométrie"@fr, "Geometry"@en;
48   owl:equivalentClass [ a owl:Restriction;
49     owl:onClass ignf:CoordinatesSystem;
50     owl:onProperty geom:crs;
51     owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger];
52   rdfs:subClassOf ngeo:Geometry;
53   rdfs:subClassOf sf:Geometry.
54
55 geom:Point a owl:Class;
56   rdfs:label "Point"@en, "Point"@fr;
57   rdfs:subClassOf geom:Geometry;
58   rdfs:subClassOf sf:Point;
59   owl:equivalentClass [ a owl:Class ;
60     owl:intersectionOf (
61       [ a owl:Restriction;
62         owl:onDataRange xsd:double;
63         owl:onProperty geom:coordY;
64         owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger ]
65
66       [ a owl:Restriction;
67         owl:onDataRange xsd:double;
68         owl:onProperty geom:coordX;
69         owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger ]
70
71       [ a owl:Restriction;
72         owl:onDataRange xsd:double;
73         owl:onProperty geom:coordZ;
74         owl:maxQualifiedCardinality "1"^^xsd:nonNegativeInteger ]
75
76       [ a owl:Restriction;
77         owl:onDataRange xsd:double;
78         owl:onProperty geom:coordM;
79         owl:maxQualifiedCardinality "1"^^xsd:nonNegativeInteger ]
80     )
81   ].
82
83 geom:Curve a owl:Class;
84   rdfs:label "Courbe"@fr, "Curve"@en;
85   rdfs:subClassOf geom:Geometry;
86   rdfs:subClassOf sf:Curve.
87
88 geom:Surface a owl:Class;
89   rdfs:label "Surface"@en, "Surface"@fr;
90   rdfs:subClassOf geom:Geometry;
91   rdfs:subClassOf sf:Surface.
92
93 geom:Envelope a owl:Class;
94   rdfs:label "Envelope"@en, "Enveloppe"@fr;
95   owl:equivalentClass [ a owl:Class ;
96     owl:intersectionOf (
97       [ a owl:Restriction;
98         owl:onClass geom:Point;
99         owl:onProperty geom:upperCorner;
100        owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger ]
101
102       [ a owl:Restriction;
103         owl:onClass geom:Point;
104         owl:onProperty geom:lowerCorner;
105         owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger ]
106     )
107   )

```

```

108      ] ;
109      rdfs:subClassOf geom:Geometry.
110
111      geom:Polygon a owl:Class;
112      rdfs:label "Polygon"@en, "Polygone"@fr;
113      rdfs:subClassOf geom:Surface;
114      owl:equivalentClass [
115          a owl:Restriction;
116          owl:someValuesFrom geom:LinearRing;
117          owl:onProperty geom:exterior];
118      rdfs:subClassOf [ a owl:Class ;
119          owl:intersectionOf (
120              [ a owl:Restriction;
121                  owl:onClass geom:LinearRing;
122                  owl:onProperty geom:exterior;
123                  owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger ]
124
125              [ a owl:Restriction;
126                  owl:someValuesFrom geom:LinearRing;
127                  owl:onProperty geom:interior;
128              ]
129          )
130      ] ;
131      rdfs:subClassOf sf:Polygon.
132
133      geom:LineString a owl:Class;
134      rdfs:label "Line string"@en, "Polyligne"@fr;
135      rdfs:subClassOf geom:Curve;
136      owl:equivalentClass [
137          a owl:Restriction;
138          owl:someValuesFrom geom:PointsList;
139          owl:onProperty geom:points];
140      rdfs:subClassOf [
141          a owl:Restriction;
142          owl:onClass geom:PointsList;
143          owl:onProperty geom:points;
144          owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger];
145      rdfs:subClassOf sf:LineString.
146
147      geom:LinearRing a owl:Class;
148      rdfs:label "Anneau"@fr, "Linear ring"@en;
149      rdfs:subClassOf geom:LineString;
150      owl:equivalentClass [
151          a owl:Restriction;
152          owl:someValuesFrom [ a owl:Class;
153              owl:intersectionOf(
154                  geom:PointsList
155                  [ a owl:Restriction;
156                      owl:onClass geom:Point;
157                      owl:onProperty geom:firstAndLast;
158                      owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger; ])
159              ];
160          owl:onProperty geom:points;
161      ];
162      rdfs:subClassOf sf:LinearRing.
163
164      geom:Line a owl:Class;
165      rdfs:label "Ligne"@fr, "Line"@en;
166      rdfs:subClassOf geom:LineString;
167      owl:equivalentClass[
168          a owl:Restriction;
169          owl:someValuesFrom [ a owl:Class;

```

```

170     owl:intersectionOf(
171         geom:PointsList
172             [ a owl:Restriction;
173                 owl:someValuesFrom [ a owl:Class;
174                     owl:intersectionOf(
175                         geom:PointsList
176                             [ a owl:Restriction;
177                                 owl:onProperty rdf:rest;
178                                 owl:hasValue rdf:nil; ]]);
179                     owl:onProperty rdf:rest; ]
180             );
181         owl:onProperty geom:points; ];
182     rdfs:subClassOf sf:Line.
183
184 geom:GeometryCollection a owl:Class;
185     rdfs:label "Collection de géométries"@fr, "Geometry collection"@en;
186     rdfs:subClassOf geom:Geometry;
187     rdfs:subClassOf sf:GeometryCollection.
188
189 geom:MultiCurve a owl:Class;
190     rdfs:label "Multi curve"@en, "Multicourbe"@fr;
191     rdfs:subClassOf geom:GeometryCollection;
192     rdfs:subClassOf sf:MultiCurve.
193
194 geom:MultiPoint a owl:Class;
195     rdfs:label "Multi point"@en, "Multipoint"@fr;
196     owl:equivalentClass [
197         a owl:Restriction;
198             owl:someValuesFrom geom:Point;
199             owl:onProperty geom:pointMember];
200     rdfs:subClassOf geom:GeometryCollection;
201     rdfs:subClassOf sf:MultiPoint.
202
203 geom:MultiPolygon a owl:Class;
204     rdfs:label "Multi polygon"@en, "Multipolygone"@fr;
205     owl:equivalentClass [
206         a owl:Restriction;
207             owl:someValuesFrom geom:Polygon;
208             owl:onProperty geom:polygonMember];
209     rdfs:subClassOf geom:MultiSurface;
210     rdfs:subClassOf sf:MultiPolygon.
211
212 geom:MultiLineString a owl:Class;
213     rdfs:label "Multi line string"@en, "Multipolyline"@fr;
214     owl:equivalentClass [
215         a owl:Restriction;
216             owl:someValuesFrom geom:LineString;
217             owl:onProperty geom:lineStringMember];
218     rdfs:subClassOf geom:MultiCurve;
219     rdfs:subClassOf sf:MultiLineString.
220
221 geom:MultiSurface a owl:Class;
222     rdfs:label "Multi surface"@en, "Multisurface"@fr;
223     rdfs:subClassOf geom:GeometryCollection;
224     rdfs:subClassOf sf:MultiSurface.
225
226 geom:PointsList a owl:Class;
227     rdfs:label "List of points"@en, "Liste de points"@fr;
228     rdfs:subClassOf rdf:List;
229     rdfs:subClassOf [ a owl:Restriction;
230         owl:allValuesFrom geom:Point;
231         owl:onProperty rdf:first].

```

```

232
233 # ----- Object Properties definition -----
234
235 geom:firstAndLast a owl:ObjectProperty;
236   rdfs:domain geom:PointsList;
237   rdfs:label "first and last"@en, "premier et dernier"@fr;
238   rdfs:subPropertyOf rdf:first;
239   rdfs:range geom:Point.
240
241 geom:points a owl:ObjectProperty;
242   rdfs:domain geom:Curve;
243   rdfs:label "points"@en, "points"@fr;
244   rdfs:range geom:PointsList.
245
246 geom:crs a owl:ObjectProperty;
247   rdfs:domain geom:Geometry;
248   rdfs:label "coordinate reference system"@en, "système de coordonnées"@fr;
249   rdfs:range ignf:CRS.
250
251 geom:boundary a owl:ObjectProperty;
252   rdfs:domain geom:Polygon;
253   rdfs:label "frontière"@fr, "boundary"@en;
254   rdfs:range geom:LinearRing.
255
256 geom:interior a owl:ObjectProperty;
257   rdfs:domain geom:Polygon;
258   rdfs:label "intérieur"@fr, "interior"@en;
259   rdfs:range geom:LinearRing;
260   rdfs:subPropertyOf geom:boundary.
261
262 geom:exterior a owl:ObjectProperty;
263   rdfs:comment "Relie un polygone à un anneau décrivant le contour extérieur de sa surface."@fr;
264   rdfs:domain geom:Polygon;
265   rdfs:label "extérieur"@fr, "exterior"@en;
266   rdfs:range geom:LinearRing;
267   rdfs:subPropertyOf geom:boundary.
268
269 geom:pointMember a owl:ObjectProperty;
270   rdfs:domain geom:MultiPoint;
271   rdfs:label "point membre"@fr, "point member"@en;
272   rdfs:range geom:Point.
273
274 geom:lineStringMember a owl:ObjectProperty;
275   rdfs:domain geom:MultiLineString;
276   rdfs:label "polyligne membre"@fr, "line string member"@en;
277   rdfs:range geom:LineString.
278
279 geom:polygonMember a owl:ObjectProperty;
280   rdfs:domain geom:MultiPolygon;
281   rdfs:label "polygone membre"@fr, "polygon member"@en;
282   rdfs:range geom:Polygon.
283
284 geom:geometry a owl:ObjectProperty;
285   rdfs:label "a pour géométrie"@fr, "has geometry"@en;
286   rdfs:range geom:Geometry.
287
288 geom:centroid a owl:ObjectProperty;
289   rdfs:domain geom:Surface;
290   rdfs:label "centroid"@en;
291   rdfs:label "centroïde"@fr;
292   rdfs:range geom:Point.
293

```

```

294 geom:envelope a owl:ObjectProperty;
295   rdfs:domain geom:Geometry;
296   rdfs:label "envelope"@en;
297   rdfs:label "enveloppe"@fr;
298   rdfs:range geom:Envelope.
299
300 geom:upperCorner a owl:ObjectProperty;
301   rdfs:domain geom:Envelope;
302   rdfs:label "upper corner"@en;
303   rdfs:label "coin supérieur"@fr;
304   rdfs:range geom:Point.
305
306 geom:lowerCorner a owl:ObjectProperty;
307   rdfs:domain geom:Envelope;
308   rdfs:label "lower corner"@en;
309   rdfs:label "coin inférieur"@fr;
310   rdfs:range geom:Point.
311
312 # ----- Datatype Properties definition -----
313
314 geom:coordX a owl:DatatypeProperty;
315   rdfs:domain geom:Point;
316   rdfs:label "x"@fr;
317   rdfs:label "x"@en;
318   rdfs:range xsd:double.
319
320 geom:coordY a owl:DatatypeProperty;
321   rdfs:domain geom:Point;
322   rdfs:label "y"@fr;
323   rdfs:label "y"@en;
324   rdfs:range xsd:double.
325
326 geom:coordZ a owl:DatatypeProperty;
327   rdfs:domain geom:Point;
328   rdfs:label "z"@fr;
329   rdfs:label "z"@en;
330   rdfs:range xsd:double.
331
332 geom:coordM a owl:DatatypeProperty;
333   rdfs:domain geom:Point;
334   rdfs:label "m"@fr;
335   rdfs:label "m"@en;
336   rdfs:range xsd:double.
337
338 # ----- Instances definition -----

```

Listing B.1: Formal definition of the geometry vocabulary in turtle. The current version is deployed at <http://data.ign.fr/def/geometrie>

## B.2 Vocabulary for CRS

```

1 @prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#>.
2 @prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>.
3 @prefix xsd: <http://www.w3.org/2001/XMLSchema#>.
4 @prefix sf: <http://www.opengis.net/ont/sf#>.
5 @prefix xml: <http://www.w3.org/XML/1998/namespace>.
6 @prefix owl: <http://www.w3.org/2002/07/owl#>.
7 @prefix ignf: <http://data.ign.fr/def/ignf#>.

```

```

8  @prefix foaf: <http://xmlns.com/foaf/0.1/> .
9  @prefix voaf: <http://purl.org/vocommons/voaf#> .
10 @prefix cc: <http://creativecommons.org/ns#> .
11 @prefix vann: <http://purl.org/vocab/vann/> .
12 @prefix dcterms: <http://purl.org/dc/terms/> .
13 @prefix dc: <http://purl.org/dc/elements/1.1/> .
14 @prefix geom: <http://data.ign.fr/def/geometrie#>.
15 @prefix qudt: <http://qudt.org/schema/qudt#> .
16
17
18
19 #####--Metadata here ----#####
20 #####----Persons URIs here -----#####
21 #####----Classes here-----#####
22
23
24 <http://data.ign.fr/def/ignf> a owl:Ontology;
25 dcterms:description "Codes pour la description de systèmes de référence de coordonnées conforme
   ISO TC/211. Les traductions franÃ§aises des termes et leurs définitions sont pour la plupart
   issues du glossaire multilingue ISO/TC 211 disponible ici: http://www.isotc211.org/
   Terminology.htm"@fr;
26 dcterms:title "Ontologie des systèmes de coordonnées"@fr;
27 dcterms:creator <http://recherche.ign.fr/labos/cogit/cv.php?prenom=Nathalie&nom=Abadie> ;
28 dcterms:creator <http://www.eurecom.fr/~atemezin/> ;
29 dcterms:contributor <http://www.eurecom.fr/~troncy/> ;
30 dcterms:issued "2013-06-11"^^xsd:date ;
31 dcterms:modified "2014-04-09"^^xsd:date ;
32 dcterms:publisher <http://fr.dbpedia.org/resource/Institut_national_de_l%27information_g%
   C3%A9ographique_et_foresti%C3%A8re> ;
33 dcterms:rights "Copyright 2014, IGN" ;
34 vann:preferredNamespacePrefix "ignf" ;
35 vann:preferredNamespaceUri <http://data.ign.fr/def/ignf#> ;
36 cc:license <http://www.data.gouv.fr/Licence-Ouverte-Open-Licence> ;
37 cc:license <http://creativecommons.org/licenses/by/2.0/> ;
38 cc:license <http://opendatacommons.org/licenses/by/> ;
39 rdfs:seeAlso <http://librairies.ignf.fr/geoportail/resources/IGNF.xml> ;
40 owl:versionInfo "Version 1.0 – 2014-04-09" .

41
42
43 #####--Metadata here ----#####
44 #####----Persons URIs here -----#####
45 #####----Classes here-----#####
46
47 <http://www.eurecom.fr/~atemezin/> a foaf:Person.
48 <http://www.eurecom.fr/~troncy/> a foaf:Person.
49 <http://recherche.ign.fr/labos/cogit/cv.php?prenom=Nathalie&nom=Abadie> a foaf:Person.

50
51
52 #####----Classes here-----#####
53 #####----Classes here-----#####
54 #####----Classes here-----#####
55
56
57 ignf:CRS a owl:Class;
58 rdfs:label "Coordinates reference system"@en, "Système de référence de coordonnées"@fr ;
59 rdfs:comment "Métadonnées permettant de préciser, selon la dimension spatiale des coordonnées 1
   D, 2D ou 3D, les éléments de définition associés au jeu de coordonnées: le Système de Réfé
   rence Terrestre, l'ellipsoïde géodésique, le méridien origine, le type de coordonnées (carté
   siennes géocentriques, planes, géographiques,...), les unités dans lesquelles sont exprimées les
   coordonnées, la projection cartographique, le référentiel altimétrique.(http://geodesie.ign.fr/
   index.php?page=glossaire)"@fr;
60 rdfs:subClassOf [ a owl:Restriction;
```

```

61      owl:someValuesFrom ignf:Extent;
62      owl:onProperty ignf:domainOfValidity ];
63  rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
64
65  ignf:CoordinateSystem a owl:Class;
66  rdfs:label "Coordinate system"@en, "Système de coordonnées"@fr ;
67  rdfs:comment "Ensemble d'axes couvrant un espace de coordonnées, et de règles mathématiques
   permettant l'affectation de coordonnées à un point."@fr;
68  rdfs:subClassOf [ a owl:Restriction;
69      owl:onClass ignf:CoordinateSystemAxis;
70      owl:onProperty ignf:axis;
71      owl:minCardinality "1"^^xsd:nonNegativeInteger];
72  rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
73
74  ignf:SingleCRS a owl:Class;
75  rdfs:label "single CRS"@en, "SRC simple"@fr;
76  rdfs:comment "CRS composé d'un système de coordonnées et d'un référentiel."@fr;
77  rdfs:subClassOf ignf:CRS;
78  rdfs:subClassOf [ a owl:Class ;
79      owl:intersectionOf
80          ([ a owl:Restriction;
81              owl:onClass ignf:Datum;
82              owl:onProperty ignf:datum;
83              owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
84
85          [ a owl:Restriction;
86              owl:onClass ignf:CoordinateSystem;
87              owl:onProperty ignf:coordinateSystem;
88              owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger] )
89      ] ;
90  rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
91
92  ignf:GeodeticCRS a owl:Class, rdfs:Class;
93  rdfs:label "Geodetic CRS"@en, "Système de coordonnées de référence géodésique"@fr;
94  rdfs:comment "Système de référence de coordonnées associé à un référentiel géodésique."@fr;
95  rdfs:subClassOf ignf:SingleCRS ;
96  rdfs:subClassOf [ a owl:Restriction;
97      owl:allValuesFrom ignf:GeodeticDatum;
98      owl:onProperty ignf:datum];
99  rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
100
101 ignf:ProjectedCRS a owl:Class, rdfs:Class;
102 rdfs:label "Projected CRS"@en, "Système de coordonnées de référence projeté"@fr;
103 rdfs:comment "Système de coordonnées de référence dérivé par projection cartographique d'un système de coordonnées de référence bidimensionnel."@fr;
104 rdfs:subClassOf ignf:SingleCRS ;
105 rdfs:subClassOf [ a owl:Restriction;
106     owl:allValuesFrom ignf:GeodeticCRS;
107     owl:onProperty ignf:baseCRS];
108 rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
109
110 ignf:CompoundCRS a owl:Class, rdfs:Class;
111 rdfs:label "Compound CRS"@en, "Système de coordonnées de référence combiné"@fr;
112 rdfs:comment "Système de référence de coordonnées utilisant deux systèmes de référence de coordonnées simples."@fr;
113 rdfs:subClassOf ignf:CRS ;
114 rdfs:subClassOf [ a owl:Restriction;
115     owl:onClass ignf:SingleCRS;
116     owl:onProperty ignf:includesSingleCRS;
117     owl:minCardinality "2"^^xsd:nonNegativeInteger];
118 rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
119

```

```

120 ignf:VerticalCRS a owl:Class, rdfs:Class;
121   rdfs:label "Vertical CRS"@en, "Système de coordonnées de référence vertical"@fr;
122   rdfs:comment "Système de référence de coordonnées à une dimension, associé à un référentiel
123   vertical."@fr;
124   rdfs:subClassOf ignf:SingleCRS ;
125   rdfs:subClassOf [ a owl:Class ;
126     owl:intersectionOf
127     ([ a owl:Restriction;
128       owl:allValuesFrom ignf:VerticalDatum;
129       owl:onProperty ignf:datum]
130
131     [ a owl:Restriction;
132       owl:allValuesFrom ignf:VerticalCS;
133       owl:onProperty ignf:coordinateSystem] )
134   ] ;
135   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
136
137 ignf:Extent a owl:Class, rdfs:Class;
138   rdfs:label "Extent"@en, "Région de validité"@fr ;
139   rdfs:comment "Zone ou intervalle de temps dans lequel la référence est valide."@fr;
140   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
141
142 ignf:Datum a owl:Class, rdfs:Class;
143   rdfs:label "Datum"@en, "Référentiel"@fr;
144   rdfs:comment "Paramètre ou ensemble de paramètres définissant la position de l'origine, l'échelle
145   et l'orientation d'un système de coordonnées."@fr;
146   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
147
148 ignf:GeodeticDatum a owl:Class, rdfs:Class;
149   rdfs:label "Geodetic datum"@en, "Référentiel géodésique"@fr;
150   rdfs:comment "Référentiel décrivant la relation entre un système de coordonnées à deux ou trois
151   dimensions et la Terre."@fr;
152   rdfs:subClassOf ignf:Datum;
153   rdfs:subClassOf [ a owl:Class ;
154     owl:intersectionOf
155     ([ a owl:Restriction;
156       owl:onClass ignf:Ellipsoid;
157       owl:onProperty ignf:ellipsoid;
158       owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
159
160     [ a owl:Restriction;
161       owl:onClass ignf:PrimeMeridian;
162       owl:onProperty ignf:primeMeridian;
163       owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger] )
164   ] ;
165   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
166
167 ignf:VerticalDatum a owl:Class, rdfs:Class;
168   rdfs:label "Vertical datum"@en, "Référentiel vertical"@fr;
169   rdfs:comment "Référentiel décrivant la relation entre les hauteurs ou les profondeurs relatives à la
170   gravité et la Terre."@fr;
171   rdfs:subClassOf ignf:Datum;
172   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
173
174 ignf:VerticalCS a owl:Class, rdfs:Class;
175   rdfs:label "Vertical CS"@en, "Système de coordonnées vertical"@fr;
176   rdfs:comment "Système de coordonnée à une dimension utilisé pour les mesures de hauteur ou de
177   profondeur relatives à la gravité."@fr;
178   rdfs:subClassOf ignf:CoordinateSystem ;
179   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
180
181 ignf:CartesianCS a owl:Class, rdfs:Class;
```

```

177 rdfs:label "Cartesian CS"@en, "Système cartésien de coordonnées"@fr;
178 rdfs:comment "Système de coordonnées donnant la position des points relativement à axes
    perpendiculaires deux à deux."@fr;
179 rdfs:subClassOf ignf:CoordinateSystem ;
180 rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
181
182 ignf:EllipsoidalCS a owl:Class, rdfs:Class;
183     rdfs:label "Ellipsoidal CS"@en, "Système de coordonnées ellipsoïdal"@fr;
184     rdfs:comment "Système de coordonnées dans lequel la position est spécifiée par la latitude géodé
        sique, la longitude géodésique et (dans le cas tridimensionnel) la hauteur ellipsoïdale."@fr;
185     rdfs:subClassOf ignf:CoordinateSystem ;
186     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
187
188 ignf:CoordinateSystemAxis a owl:Class, rdfs:Class;
189     rdfs:label "Axis"@en, "Axe"@fr;
190     rdfs:comment "Axe par rapport auquel une coordonnées d'un point est spécifiée dans un système
        de coordonnées."@fr;
191     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
192
193 ignf:Ellipsoid a owl:Class, rdfs:Class;
194     rdfs:label "Ellipsoid"@en;
195     rdfs:label "Ellipsoïde"@fr;
196     rdfs:comment "Surface de révolution engendrée par une ellipse tournant autour de son petit axe, d
        éfinie par le rayon équatorial et un paramètre d'aplatissement, et sensiblement géocentrique.
        Note : Il s'agit d'un modèle mathématique du géoïde, c'est-à-dire de la Terre débarrassée de
        son relief. Il existe de nombreux ellipsoïdes géodésiques."@fr;
197     rdfs:subClassOf [ a owl:Class ;
198         owl:intersectionOf
199             ([ a owl:Restriction;
200                 owl:onClass qudt:QuantityValue;
201                 owl:onProperty ignf:semiMajorAxis;
202                 owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger
203             ]
204
205             [ a owl:Class ;
206                 owl:unionOf(
207                     [a owl:Restriction;
208                         owl:onClass qudt:QuantityValue;
209                         owl:onProperty ignf:semiMinorAxis;
210                         owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
211
212                     [a owl:Restriction;
213                         owl:onClass qudt:QuantityValue;
214                         owl:onProperty ignf:inverseFlattening;
215                         owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
216
217                     [a owl:Restriction;
218                         owl:onDataRange xsd:boolean;
219                         owl:onProperty ignf:isSphere;
220                         owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger])
221                 ]
222             );
223             rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
224
225 ignf:PrimeMeridian a owl:Class, rdfs:Class;
226     rdfs:label "Prime meridian"@en;
227     rdfs:label "Méridien origine"@fr;
228     rdfs:comment "Méridien à partir duquel les longitudes d'autres méridiens sont mesurées."@fr;
229     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
230
231 ignf:GeographicBoundingBox a owl:Class, rdfs:Class;
232     rdfs:label "Geographic bounding box"@en;

```

```

233 rdfs:label "Cadre englobant géographique"@fr;
234 rdfs:comment "Cadre délimitant une zone d'intérêt."@fr;
235 rdfs:subClassOf [ a owl:Class ;
236   owl:intersectionOf
237     ([ a owl:Restriction;
238       owl:onDataRange xsd:double;
239       owl:onProperty ignf:westBoundLongitude;
240       owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
241
242     [ a owl:Restriction;
243       owl:onDataRange xsd:double;
244       owl:onProperty ignf:eastBoundLongitude;
245       owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
246
247     [ a owl:Restriction;
248       owl:onDataRange xsd:double;
249       owl:onProperty ignf:southBoundLatitude;
250       owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
251
252     [ a owl:Restriction;
253       owl:onDataRange xsd:double;
254       owl:onProperty ignf:northBoundLatitude;
255       owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger] )
256   ] ;
257   rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
258
259 ignf:CoordinateOperation a owl:Class, rdfs:Class;
260   rdfs:label "Coordinate operation"@en;
261   rdfs:label "Opération sur les coordonnées"@fr;
262   rdfs:comment "Opération mathématique sur des coordonnées dans un SRC source en vue de les
263   convertir vers un SRC cible."@fr;
264   rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
265
266 ignf:SingleOperation a owl:Class, rdfs:Class;
267   rdfs:label "Single coordinate operation"@en;
268   rdfs:label "Opération simple sur les coordonnées"@fr;
269   rdfs:comment "Opération non concaténée sur des coordonnées."@fr;
270   rdfs:subClassOf ignf:CoordinateOperation;
271   rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
272
273 ignf:Transformation a owl:Class, rdfs:Class;
274   rdfs:label "Transformation"@en;
275   rdfs:label "Transformation"@fr;
276   rdfs:comment "Opération par laquelle des coordonnées en entrée et en sortie sont associées à diffé
277   rents référentiels."@fr;
278   rdfs:subClassOf ignf:SingleOperation;
279   rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
280
281 ignf:Conversion a owl:Class, rdfs:Class;
282   rdfs:label "Conversion"@en;
283   rdfs:label "Conversion"@fr;
284   rdfs:comment "Opération par laquelle les coordonnées en sortie sont associées au même réfé
285   rentiel que les données en entrée (ex: projection cartographique)."@fr;
286   rdfs:subClassOf ignf:SingleOperation;
287   rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
288
289 ignf:OperationMethod a owl:Class, rdfs:Class;
290   rdfs:label "Operation method"@en;
291   rdfs:label "Méthode"@fr;
292   rdfs:comment "Méthode utilisée pour faire une opération sur des coordonnées."@fr;
293   rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
294
295

```

```

292 ignf:OperationParameter a owl:Class, rdfs:Class;
293   rdfs:label "Operation parameter"@en;
294   rdfs:label "Paramètre"@fr;
295   rdfs:comment "Paramètre utilisé par une méthode pour faire une opération sur des coordonnées."@fr;
296   rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
297
298 ignf:OperationParameterValue a owl:Class, rdfs:Class;
299   rdfs:label "Operation parameter value"@en;
300   rdfs:label "Valeur de paramètre"@fr;
301   rdfs:comment "Valeur d'un paramètre utilisé par une méthode pour faire une opération sur des coordonnées."@fr;
302   rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
303
304 ######
305 ## Axioms here ##
306 ######
307
308
309 ignf:CompoundCRS owl:disjointWith ignf:GeodeticCRS .
310 ignf:VerticalCS owl:disjointWith ignf:VerticalCRS .
311 ignf:CartesianCS owl:disjointWith ignf:GeodeticCRS .
312 ignf:EllipsoidalCS owl:disjointWith ignf:GeodeticCRS .
313 ignf:PrimeMeridian owl:disjointWith ignf:Ellipsoid .
314 ignf:CRS owl:disjointWith ignf:CoordinateSystem .
315
316 ######
317 ### ---- Properties here -----###
318 ######
319
320
321 ignf:domainOfValidity a owl:ObjectProperty ;
322   rdfs:label "domain of validity"@en, "domaine de validité"@fr ;
323   rdfs:comment "Zone ou intervalle de temps dans lequel un SRC est valide."@fr;
324   rdfs:range ignf:Extent ;
325   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
326
327 ignf:includesSingleCRS a owl:ObjectProperty ;
328   rdfs:label "includes single CRS"@en, "inclus un simple SRC"@fr ;
329   rdfs:comment "Désigne les SRC simples qui composent un SRC composé."@fr;
330   rdfs:domain ignf:CompoundCRS ;
331   rdfs:range ignf:SingleCRS ;
332   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
333
334 ignf:baseCRS a owl:ObjectProperty ;
335   rdfs:label "base CRS"@en, "SRC de base"@fr ;
336   rdfs:comment "Désigne le SRC géodésique sur lequel repose un SRC projeté."@fr;
337   rdfs:domain ignf:ProjectedCRS ;
338   rdfs:range ignf:GeodeticCRS ;
339   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
340
341 ignf:coordinateSystem a owl:ObjectProperty ;
342   rdfs:label "uses coordinate system"@en, "utilise le système de coordonnées"@fr ;
343   rdfs:comment "Désigne le système de coordonnées utilisé par un SRC."@fr;
344   rdfs:domain ignf:SingleCRS ;
345   rdfs:range ignf:CoordinateSystem ;
346   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
347
348 ignf:ellipsoidalCS a owl:ObjectProperty ;
349   rdfs:label "uses ellipsoidal CS"@en, "utilise le système de coordonnées ellipsoïdal"@fr;
350   rdfs:comment "Désigne le système de coordonnées ellipsoïdal utilisé par un SRC géodésique."@fr;
351   rdfs:subPropertyOf ignf:coordinateSystem;

```

```

352    rdfs:domain ignf:GeodeticCRS ;
353    rdfs:range ignf:EllipsoidalCS ;
354    rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
355
356    ignf:cartesianCS a owl:ObjectProperty ;
357    rdfs:label "uses Cartesian CS"@en, "utilise le système cartésien de coordonnées"@fr;
358    rdfs:comment "Désigne le système de coordonnées cartésien utilisé par un SRC géodésique ou
359    projeté."@fr;
360    rdfs:subPropertyOf ignf:coordinateSystem;
361    rdfs:domain [ a owl:Class ;
362        owl:unionOf (ignf:GeodeticCRS ignf:ProjectedCRS )] ;
363    rdfs:range ignf:CartesianCS ;
364    rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
365
366    ignf:verticalCS a owl:ObjectProperty ;
367    rdfs:label "uses Vertical CS"@en, "utilise le système vertical de coordonnées"@fr;
368    rdfs:comment "Désigne le système de coordonnées vertical utilisé par un SRC vertical."@fr;
369    rdfs:subPropertyOf ignf:coordinateSystem;
370    rdfs:domain ignf:VerticalCRS ;
371    rdfs:range ignf:VerticalCS ;
372    rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
373
374    ignf:datum a owl:ObjectProperty ;
375    rdfs:label "uses datum"@en, "utilise le référentiel"@fr;
376    rdfs:comment "Désigne le référentiel utilisé par un SRC simple."@fr;
377    rdfs:domain ignf:SingleCRS;
378    rdfs:range ignf:Datum ;
379    rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
380
381    ignf:geodeticDatum a owl:ObjectProperty ;
382    rdfs:label "uses geodetic datum"@en, "utilise le référentiel géodésique"@fr;
383    rdfs:comment "Désigne le référentiel géodésique utilisé par un SRC géodésique."@fr;
384    rdfs:subPropertyOf ignf:datum;
385    rdfs:domain ignf:GeodeticCRS;
386    rdfs:range ignf:GeodeticDatum ;
387    rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
388
389    ignf:verticalDatum a owl:ObjectProperty ;
390    rdfs:label "uses vertical datum"@en, "utilise le référentiel vertical"@fr;
391    rdfs:comment "Désigne le référentiel vertical utilisé par un SRC vertical."@fr;
392    rdfs:subPropertyOf ignf:datum;
393    rdfs:domain ignf:VerticalCRS ;
394    rdfs:range ignf:VerticalDatum ;
395    rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
396
397    ignf:axis a owl:ObjectProperty ;
398    rdfs:label "uses Axis"@en, "utilise l'axe"@fr;
399    rdfs:comment "Désigne un axe utilisé par un système de coordonnées ellipsoïdal ou cartésien."@fr;
400    rdfs:domain ignf:CoordinateSystem ;
401    rdfs:range ignf:CoordinateSystemAxis ;
402    rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
403
404    ignf:primeMeridian a owl:ObjectProperty ;
405    rdfs:label "uses prime meridian"@en, "utilise le méridien origine"@fr;
406    rdfs:comment "Désigne le méridien origine d'un référentiel géodésique."@fr;
407    rdfs:domain ignf:GeodeticDatum ;
408    rdfs:range ignf:PrimeMeridian ;
409    rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
410
411    ignf:ellipsoid a owl:ObjectProperty ;
412    rdfs:label "uses ellipsoid"@en, "utilise l'ellipsoïde"@fr;
413    rdfs:comment "Désigne l'ellipsoïde utilisé par un référentiel géodésique."@fr;

```

```

413  rdfs:domain ignf:GeodeticDatum;
414  rdfs:range ignf:Ellipsoid ;
415  rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
416
417  ignf:geographicElement a owl:ObjectProperty ;
418  rdfs:label "geographic element"@en, "élément géographique"@fr;
419  rdfs:comment "Désigne le cadre englobant géographique d'une région de validité."@fr;
420  rdfs:domain ignf:Extent;
421  rdfs:range ignf:GeographicBoundingBox ;
422  rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
423
424  ignf:semiMajorAxis a owl:ObjectProperty;
425  rdfs:label "semi major axis"@en, "demi grand axe"@fr;
426  rdfs:comment "Désigne la longueur du demi grand axe d'un ellipsoïde."@fr;
427  rdfs:domain ignf:Ellipsoid;
428  rdfs:range [a owl:Class ;
429    owl:intersectionOf
430      (qudt:QuantityValue
431        [ a owl:Restriction;
432          owl:allValuesFrom qudt:LengthUnit;
433          owl:onProperty qudt:unit
434        ] );
435  rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
436
437  ignf:semiMinorAxis a owl:ObjectProperty;
438  rdfs:label "semi minor axis"@en, "demi petit axe"@fr;
439  rdfs:comment "Désigne la longueur du demi petit axe d'un ellipsoïde."@fr;
440  rdfs:domain ignf:Ellipsoid;
441  rdfs:range [a owl:Class ;
442    owl:intersectionOf
443      (qudt:QuantityValue
444        [ a owl:Restriction;
445          owl:allValuesFrom qudt:LengthUnit;
446          owl:onProperty qudt:unit
447        ] );
448  rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
449
450  ignf:inverseFlattening a owl:ObjectProperty;
451  rdfs:label "inverse flattening"@en, "aplatissement inverse"@fr;
452  rdfs:comment "Désigne la valeur d'aplatissement inverse d'un ellipsoïde, exprimée sous la forme d'un nombre ou d'un ratio (pourcentage, parties par million, etc.)."@fr;
453  rdfs:domain ignf:Ellipsoid;
454  rdfs:range [a owl:Class ;
455    owl:intersectionOf
456      (qudt:QuantityValue
457        [ a owl:Restriction;
458          owl:allValuesFrom qudt:CountingUnit;
459          owl:onProperty qudt:unit
460        ] );
461  rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
462
463  ignf:greenwichLongitude a owl:ObjectProperty;
464  rdfs:label "greenwich longitude"@en, "longitude par rapport àGreenwich"@fr;
465  rdfs:comment "Désigne la longitude par rapport au méridien de Greenwich."@fr;
466  rdfs:domain ignf:PrimeMeridian;
467  rdfs:range [a owl:Class ;
468    owl:intersectionOf
469      (qudt:QuantityValue
470        [ a owl:Restriction;
471          owl:allValuesFrom qudt:AngleUnit;
472          owl:onProperty qudt:unit
473        ] );

```

```

474 rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
475
476 ignf:usesMethod a owl:ObjectProperty ;
477   rdfs:label "uses method"@en, "utilise la méthode"@fr;
478   rdfs:comment "Désigne la méthode utilisée par une opération sur des coordonnées."@fr;
479   rdfs:domain ignf:CoordinateOperation;
480   rdfs:range ignf:OperationMethod ;
481   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
482
483 ignf:usesValue a owl:ObjectProperty ;
484   rdfs:label "uses value"@en, "utilise la valeur"@fr;
485   rdfs:comment "Désigne une valeur utilisée par une opération sur des coordonnées."@fr;
486   rdfs:domain ignf:CoordinateOperation;
487   rdfs:range ignf:OperationParameterValue ;
488   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
489
490 ignf:usesParameter a owl:ObjectProperty ;
491   rdfs:label "uses parameter"@en, "utilise le paramètre"@fr;
492   rdfs:comment "Désigne un paramètre utilisé par une méthode pour faire une opération sur des
493   coordonnées."@fr;
494   rdfs:domain ignf:OperationMethod;
495   rdfs:range ignf:OperationParameter;
496   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
497
498 ignf:parameter a owl:ObjectProperty ;
499   rdfs:label "parameter"@en, "paramètre"@fr;
500   rdfs:comment "Désigne le paramètre auquel est associée une valeur."@fr;
501   rdfs:domain ignf:OperationParameterValue;
502   rdfs:range ignf:OperationParameter;
503   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
504
505 ignf:sourceCRS a owl:ObjectProperty ;
506   rdfs:label "source CRS"@en, "SRC source"@fr;
507   rdfs:comment "Désigne le SRC des données en entrée d'une opération."@fr;
508   rdfs:domain ignf:CoordinateOperation;
509   rdfs:range ignf:CRS;
510   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
511
512 ignf:targetCRS a owl:ObjectProperty ;
513   rdfs:label "target CRS"@en, "SRC cible"@fr;
514   rdfs:comment "Désigne le SRC des données en sortie d'une opération."@fr;
515   rdfs:domain ignf:CoordinateOperation;
516   rdfs:range ignf:CRS;
517   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
518 #####
519 #####-- Dataproperties here ---#####
520 #####
521
522 ignf:westBoundLongitude a owl:DatatypeProperty;
523   rdfs:domain ignf:GeographicBoundingBox;
524   rdfs:label "west bound longitude"@en;
525   rdfs:label "longitude ouest minimale"@fr;
526   rdfs:comment "Désigne la longitude ouest minimale du cadre englobant."@fr;
527   rdfs:range xsd:double.
528
529 ignf:eastBoundLongitude a owl:DatatypeProperty;
530   rdfs:domain ignf:GeographicBoundingBox;
531   rdfs:label "east bound longitude"@en;
532   rdfs:label "longitude est maximale"@fr;
533   rdfs:comment "Désigne la longitude est maximale du cadre englobant."@fr;
534   rdfs:range xsd:double.

```

```

535 ignf:southBoundLatitude a owl:DatatypeProperty;
536   rdfs:domain ignf:GeographicBoundingBox;
537   rdfs:label "south bound latitude"@en;
538   rdfs:label "latitude sud minimale"@fr;
539   rdfs:comment "Désigne la latitude sud minimale du cadre englobant."@fr;
540   rdfs:range xsd:double.
541
542 ignf:northBoundLatitude a owl:DatatypeProperty;
543   rdfs:domain ignf:GeographicBoundingBox;
544   rdfs:label "north bound longitude"@en;
545   rdfs:label "latitude nord maximale"@fr;
546   rdfs:comment "Désigne la latitude nord maximale du cadre englobant."@fr;
547   rdfs:range xsd:double.
548
549 ignf:scope a owl:DatatypeProperty ;
550   rdfs:label "scope"@en, "portée"@fr ;
551   rdfs:comment "Désigne la portée de la référence."@fr;
552   rdfs:range xsd:string.
553
554 ignf:codeSpace a owl:DatatypeProperty;
555   rdfs:label "code space"@en, "espace de codage"@fr;
556   rdfs:comment "Désigne la règle ou l'autorité dont résulte la valeur de la référence."@fr;
557   rdfs:range xsd:string.
558
559 ignf:conversion a owl:DatatypeProperty;
560   rdfs:label "defined by conversion"@en, "défini par conversion"@fr;
561   rdfs:comment "Désigne la conversion utilisée pour définir un SRC projeté."@fr;
562   rdfs:domain ignf:ProjectedCRS;
563   rdfs:range xsd:anyURI.
564
565 ignf:axisAbbrev a owl:DatatypeProperty;
566   rdfs:label "axis abbreviation"@en, "abréviation utilisée pour désigner l'axe"@fr;
567   rdfs:comment "Désigne l'abréviation utilisée pour désigner l'axe."@fr;
568   rdfs:domain ignf:CoordinateSystemAxis;
569   rdfs:range xsd:string.
570
571 ignf:axisDirection a owl:DatatypeProperty;
572   rdfs:label "axis direction"@en, "direction de l'axe"@fr;
573   rdfs:comment "Désigne la direction de l'axe."@fr;
574   rdfs:domain ignf:CoordinateSystemAxis;
575   rdfs:range xsd:string.
576
577 ignf:isSphere a owl:DatatypeProperty;
578   rdfs:label "is sphere"@en, "est une sphère"@fr;
579   rdfs:comment "Indique si l'ellipsoïde est une sphère."@fr;
580   rdfs:domain ignf:Ellipsoid;
581   rdfs:range xsd:boolean.
582
583 ignf:epsgID a owl:DatatypeProperty;
584   rdfs:label "espg identifier"@en, "identifiant epsg"@fr;
585   rdfs:comment "Indique l'identifiant EPSG de la ressource."@fr;
586   rdfs:range xsd:string.
587
588 ignf:operationVersion a owl:DatatypeProperty;
589   rdfs:label "operation version"@en, "version de l'opération"@fr;
590   rdfs:comment "Indique la version d'une opération appliquée sur des coordonnées."@fr;
591   rdfs:domain ignf:CoordinateOperation;
592   rdfs:range xsd:string.
593
594 ignf:methodFormula a owl:DatatypeProperty;
595   rdfs:label "method formula"@en, "formule de la méthode"@fr;
596

```

```

597 rdfs:comment "Indique la formule utilisée par une méthode pour faire opération sur des coordonné
      es."@fr;
598 rdfs:domain ignf:OperationMethod;
599 rdfs:range xsd:string.
600
601 ignf:sourceDimension a owl:DatatypeProperty;
602 rdfs:label "source CRS dimension"@en, "dimension du SRC source"@fr;
603 rdfs:comment "Indique la dimension du SRC des données en entrée d'une opération."@fr;
604 rdfs:domain ignf:OperationMethod;
605 rdfs:range xsd:integer.
606
607 ignf:targetDimension a owl:DatatypeProperty;
608 rdfs:label "target CRS dimension"@en, "dimension du SRC cible"@fr;
609 rdfs:comment "Indique la dimension du SRC des données en sortie d'une opération."@fr;
610 rdfs:domain ignf:OperationMethod;
611 rdfs:range xsd:integer.

```

Listing B.2: Formal definition of the CRS vocabulary in turtle. The current version is deployed at <http://data.ign.fr/def/ignf>

## B.3 Vocabulary for French Administrative units

```

1 @prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#>.
2 @prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>.
3 @prefix xsd: <http://www.w3.org/2001/XMLSchema#>.
4 @prefix sf: <http://www.opengis.net/ont/sf#>.
5 @prefix xml: <http://www.w3.org/XML/1998/namespace>.
6 @prefix owl: <http://www.w3.org/2002/07/owl#>.
7 @prefix ignf: <http://data.ign.fr/def/ignf#>.
8 @prefix foaf: <http://xmlns.com/foaf/0.1/> .
9 @prefix voaf: <http://purl.org/vocommons/voaf#> .
10 @prefix cc: <http://creativecommons.org/ns#> .
11 @prefix vann: <http://purl.org/vocab/vann/> .
12 @prefix dcterms: <http://purl.org/dc/terms/> .
13 @prefix dc: <http://purl.org/dc/elements/1.1/> .
14 @prefix geom: <http://data.ign.fr/def/geometrie#>.
15 @prefix qudt: <http://qudt.org/schema/qudt#> .

16
17
18
19 #####--Metadata here--#####
20 #####--#
21 #####
22
23
24 <http://data.ign.fr/def/ignf> a owl:Ontology;
25   dcterms:description "Codes pour la description de systèmes de référence de coordonnées conforme
      ISO TC/211. Les traductions franÃ§aises des termes et leurs définitions sont pour la plupart
      issues du glossaire multilingue ISO/TC 211 disponible ici: http://www.isotc211.org/
      Terminology.htm"@fr;
26   dcterms:title "Ontologie des systèmes de coordonnées"@fr;
27   dcterms:creator <http://recherche.ign.fr/labos/cogit/cv.php?prenom=Nathalie&nom=Abadie> ;
28   dcterms:creator <http://www.eurecom.fr/~atemezin/> ;
29   dcterms:contributor <http://www.eurecom.fr/~troncy/> ;
30   dcterms:issued "2013-06-11"^^xsd:date ;
31   dcterms:modified "2014-04-09"^^xsd:date ;
32   dcterms:publisher <http://fr.dbpedia.org/resource/Institut\_national\_de\_l%27information\_g%C3%A9ographique\_et\_foresti%C3%A8re> ;

```

```

33 dcterms:rights "Copyright 2014, IGN" ;
34 vann:preferredNamespacePrefix "ignf" ;
35 vann:preferredNamespaceUri <http://data.ign.fr/def/ignf#> ;
36 cc:license <http://www.data.gouv.fr/Licence—Ouverte—Open—Licence> ;
37 cc:license <http://creativecommons.org/licenses/by/2.0/> ;
38 cc:license <http://opendatacommons.org/licenses/by/> ;
39 rdfs:seeAlso <http://librairies.ignf.fr/geoportail/resources/IGNF.xml> ;
40 owl:versionInfo "Version 1.0 – 2014–04–09" .
41
42
43 #####--Persons URIs here ----#####
44 #####---Classes here----#
45 #####----#
46
47 <http://www.eurecom.fr/~atemezin/> a foaf:Person.
48 <http://www.eurecom.fr/~troncy/> a foaf:Person.
49 <http://recherche.ign.fr/labos/cogit/cv.php?prenom=Nathalie&nom=Abadie> a foaf:Person.
50
51
52 #####----#
53 #####---#
54 #####----#
55
56
57 ignf:CRS a owl:Class;
58   rdfs:label "Coordinates reference system"@en, "Système de référence de coordonnées"@fr ;
59   rdfs:comment "Métadonnées permettant de préciser, selon la dimension spatiale des coordonnées 1
D, 2D ou 3D, les éléments de définition associés au jeu de coordonnées: le Système de Réfé
rence Terrestre, l'ellipsoïde géodésique, le méridien origine, le type de coordonnées (carté
siennes géocentriques, planes, géographiques,...), les unités dans lesquelles sont exprimées les
coordonnées, la projection cartographique, le référentiel altimétrique.(http://geodesie.ign.fr/
index.php?page=glossaire)"@fr;
60   rdfs:subClassOf [ a owl:Restriction;
61     owl:someValuesFrom ignf:Extent;
62     owl:onProperty ignf:domainOfValidity ];
63   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
64
65 ignf:CoordinateSystem a owl:Class;
66   rdfs:label "Coordinate system"@en, "Système de coordonnees"@fr ;
67   rdfs:comment "Ensemble d'axes couvrant un espace de coordonnées, et de règles mathématiques
permettant l'affectation de coordonnées à un point."@fr;
68   rdfs:subClassOf [ a owl:Restriction;
69     owl:onClass ignf:CoordinateSystemAxis;
70     owl:onProperty ignf:axis;
71     owl:minCardinality "1"^^xsd:nonNegativeInteger];
72   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
73
74 ignf:SingleCRS a owl:Class;
75   rdfs:label "single CRS"@en, "SRC simple"@fr;
76   rdfs:comment "CRS composé d'un système de coordonnées et d'un référentiel."@fr;
77   rdfs:subClassOf ignf:CRS;
78   rdfs:subClassOf [ a owl:Class ;
79     owl:intersectionOf
80       ([ a owl:Restriction;
81         owl:onClass ignf:Datum;
82         owl:onProperty ignf:datum;
83         owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
84
85       [ a owl:Restriction;
86         owl:onClass ignf:CoordinateSystem;
87         owl:onProperty ignf:coordinateSystem;
88         owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger] )

```

```

89      ] ;
90      rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
91
92      ignf:GeodeticCRS a owl:Class, rdfs:Class;
93      rdfs:label "Geodetic CRS"@en, "Système de coordonnées de référence géodésique"@fr;
94      rdfs:comment "Système de référence de coordonnées associé à un référentiel géodésique."@fr;
95      rdfs:subClassOf ignf:SingleCRS ;
96      rdfs:subClassOf [ a owl:Restriction;
97          owl:allValuesFrom ignf:GeodeticDatum;
98          owl:onProperty ignf:datum];
99      rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
100
101     ignf:ProjectedCRS a owl:Class, rdfs:Class;
102     rdfs:label "Projected CRS"@en, "Système de coordonnées de référence projeté"@fr;
103     rdfs:comment "Système de coordonnées de référence dérivé par projection cartographique d'un système de coordonnées de référence bidimensionnel."@fr;
104     rdfs:subClassOf ignf:SingleCRS ;
105     rdfs:subClassOf [ a owl:Restriction;
106         owl:allValuesFrom ignf:GeodeticCRS;
107         owl:onProperty ignf:baseCRS];
108     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
109
110     ignf:CompoundCRS a owl:Class, rdfs:Class;
111     rdfs:label "Compound CRS"@en, "Système de coordonnées de référence combiné"@fr;
112     rdfs:comment "Système de référence de coordonnées utilisant deux systèmes de référence de coordonnées simples."@fr;
113     rdfs:subClassOf ignf:CRS ;
114     rdfs:subClassOf [ a owl:Restriction;
115         owl:onClass ignf:SingleCRS;
116         owl:onProperty ignf:includesSingleCRS;
117         owl:minCardinality "2"^^xsd:nonNegativeInteger];
118     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
119
120     ignf:VerticalCRS a owl:Class, rdfs:Class;
121     rdfs:label "Vertical CRS"@en, "Système de coordonnées de référence vertical"@fr;
122     rdfs:comment "Système de référence de coordonnées à une dimension, associé à un référentiel vertical."@fr;
123     rdfs:subClassOf ignf:SingleCRS ;
124     rdfs:subClassOf [ a owl:Class ;
125         owl:intersectionOf
126             ([ a owl:Restriction;
127                 owl:allValuesFrom ignf:VerticalDatum;
128                 owl:onProperty ignf:datum]
129
130                 [ a owl:Restriction;
131                     owl:allValuesFrom ignf:VerticalCS;
132                     owl:onProperty ignf:coordinateSystem] )
133             ] ;
134     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
135
136     ignf:Extent a owl:Class, rdfs:Class;
137     rdfs:label "Extent"@en, "Région de validité"@fr ;
138     rdfs:comment "Zone ou intervalle de temps dans lequel la référence est valide."@fr;
139     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
140
141     ignf:Datum a owl:Class, rdfs:Class;
142     rdfs:label "Datum"@en, "Référentiel"@fr;
143     rdfs:comment "Paramètre ou ensemble de paramètres définissant la position de l'origine, l'échelle et l'orientation d'un système de coordonnées."@fr;
144     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
145
146     ignf:GeodeticDatum a owl:Class, rdfs:Class;

```

```

147 rdfs:label "Geodetic datum"@en, "Référentiel géodésique"@fr;
148 rdfs:comment "Référentiel décrivant la relation entre un système de coordonnées à deux ou trois
    dimensions et la Terre."@fr;
149 rdfs:subClassOf ignf:Datum;
150 rdfs:subClassOf [ a owl:Class ;
151     owl:intersectionOf
152     ( [ a owl:Restriction;
153         owl:onClass ignf:Ellipsoid;
154         owl:onProperty ignf:ellipsoid;
155         owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
156
157         [ a owl:Restriction;
158             owl:onClass ignf:PrimeMeridian;
159             owl:onProperty ignf:primeMeridian;
160             owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger] )
161     ] ;
162 rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
163
164 ignf:VerticalDatum a owl:Class, rdfs:Class;
165     rdfs:label "Vertical datum"@en, "Référentiel vertical"@fr;
166     rdfs:comment "Référentiel décrivant la relation entre les hauteurs ou les profondeurs relatives à la
        gravité et la Terre."@fr;
167     rdfs:subClassOf ignf:Datum;
168     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
169
170 ignf:VerticalCS a owl:Class, rdfs:Class;
171     rdfs:label "Vertical CS"@en, "Système de coordonnées vertical"@fr;
172     rdfs:comment "Système de coordonnée à une dimension utilisé pour les mesures de hauteur ou de
        profondeur relatives à la gravité."@fr;
173     rdfs:subClassOf ignf:CoordinateSystem ;
174     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
175
176 ignf:CartesianCS a owl:Class, rdfs:Class;
177     rdfs:label "Cartesian CS"@en, "Système cartésien de coordonnées"@fr;
178     rdfs:comment "Système de coordonnées donnant la position des points relativement à axes
        perpendiculaires deux à deux."@fr;
179     rdfs:subClassOf ignf:CoordinateSystem ;
180     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
181
182 ignf:EllipsoidalCS a owl:Class, rdfs:Class;
183     rdfs:label "Ellipsoidal CS"@en, "Système de coordonnées ellipsoïdal"@fr;
184     rdfs:comment "Système de coordonnées dans lequel la position est spécifiée par la latitude géodé
        sique, la longitude géodésique et (dans le cas tridimensionnel) la hauteur ellipsoïdale."@fr;
185     rdfs:subClassOf ignf:CoordinateSystem ;
186     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
187
188 ignf:CoordinateSystemAxis a owl:Class, rdfs:Class;
189     rdfs:label "Axis"@en, "Axe"@fr;
190     rdfs:comment "Axe par rapport auquel une coordonnées d'un point est spécifiée dans un système
        de coordonnées."@fr;
191     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
192
193 ignf:Ellipsoid a owl:Class, rdfs:Class;
194     rdfs:label "Ellipsoid"@en;
195     rdfs:label "Ellipsoïde"@fr;
196     rdfs:comment "Surface de révolution engendrée par une ellipse tournant autour de son petit axe, d
        éfinie par le rayon équatorial et un paramètre d'aplatissement, et sensiblement géocentrique.
        Note : Il s'agit d'un modèle mathématique du géoïde, c'est-à-dire de la Terre débarrassée de
        son relief. Il existe de nombreux ellipsoïdes géodésiques."@fr;
197     rdfs:subClassOf [ a owl:Class ;
198         owl:intersectionOf
199             ([ a owl:Restriction;

```

```

200      owl:onClass qudt:QuantityValue;
201      owl:onProperty ignf:semiMajorAxis;
202      owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger
203      ]
204
205      [ a owl:Class ;
206      owl:unionOf(
207          [a owl:Restriction;
208              owl:onClass qudt:QuantityValue;
209              owl:onProperty ignf:semiMinorAxis;
210              owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
211
212          [a owl:Restriction;
213              owl:onClass qudt:QuantityValue;
214              owl:onProperty ignf:inverseFlattening;
215              owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
216
217          [a owl:Restriction;
218              owl:onDataRange xsd:boolean;
219              owl:onProperty ignf:isSphere;
220              owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger])
221      ]
222      ] ;
223      rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
224
225 ignf:PrimeMeridian a owl:Class, rdfs:Class;
226     rdfs:label "Prime meridian"@en;
227     rdfs:label "Méridien origine"@fr;
228     rdfs:comment "Méridien à partir duquel les longitudes d'autres méridiens sont mesurées."@fr;
229     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
230
231 ignf:GeographicBoundingBox a owl:Class, rdfs:Class;
232     rdfs:label "Geographic bounding box"@en;
233     rdfs:label "Cadre englobant géographique"@fr;
234     rdfs:comment "Cadre délimitant une zone d'intérêt."@fr;
235     rdfs:subClassOf [ a owl:Class ;
236         owl:intersectionOf
237             ([ a owl:Restriction;
238                 owl:onDataRange xsd:double;
239                 owl:onProperty ignf:westBoundLongitude;
240                 owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
241
242             [ a owl:Restriction;
243                 owl:onDataRange xsd:double;
244                 owl:onProperty ignf:eastBoundLongitude;
245                 owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
246
247             [ a owl:Restriction;
248                 owl:onDataRange xsd:double;
249                 owl:onProperty ignf:southBoundLatitude;
250                 owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger]
251
252             [ a owl:Restriction;
253                 owl:onDataRange xsd:double;
254                 owl:onProperty ignf:northBoundLatitude;
255                 owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger] )
256         ]
257     rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
258
259 ignf:CoordinateOperation a owl:Class, rdfs:Class;
260     rdfs:label "Coordinate operation"@en;
261     rdfs:label "Opération sur les coordonnées"@fr;

```

```

262 rdfs:comment "Opération mathématique sur des coordonnées dans un SRC source en vue de les
263   convertir vers un SRC cible."@fr;
264 rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
265
265 ignf:SingleOperation a owl:Class, rdfs:Class;
266   rdfs:label "Single coordinate operation"@en;
267   rdfs:label "Opération simple sur les coordonnées"@fr;
268   rdfs:comment "Opération non concaténée sur des coordonnées."@fr;
269   rdfs:subClassOf ignf:CoordinateOperation;
270   rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
271
272 ignf:Transformation a owl:Class, rdfs:Class;
273   rdfs:label "Transformation"@en;
274   rdfs:label "Transformation"@fr;
275   rdfs:comment "Opération par laquelle des coordonnées en entrée et en sortie sont associées à diffé
276   rents référentiels."@fr;
276 rdfs:subClassOf ignf:SingleOperation;
277 rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
278
279 ignf:Conversion a owl:Class, rdfs:Class;
280   rdfs:label "Conversion"@en;
281   rdfs:label "Conversion"@fr;
282   rdfs:comment "Opération par laquelle les coordonnées en sortie sont associées au même réfé
283   rentiel que les données en entrée (ex: projection cartographique)."@fr;
283 rdfs:subClassOf ignf:SingleOperation;
284 rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
285
286 ignf:OperationMethod a owl:Class, rdfs:Class;
287   rdfs:label "Operation method"@en;
288   rdfs:label "Méthode"@fr;
289   rdfs:comment "Méthode utilisée pour faire une opération sur des coordonnées."@fr;
290   rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
291
292 ignf:OperationParameter a owl:Class, rdfs:Class;
293   rdfs:label "Operation parameter"@en;
294   rdfs:label "Paramètre"@fr;
295   rdfs:comment "Paramètre utilisé par une méthode pour faire une opération sur des coordonnées."@
296   fr;
296 rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
297
298 ignf:OperationParameterValue a owl:Class, rdfs:Class;
299   rdfs:label "Operation parameter value"@en;
300   rdfs:label "Valeur de paramètre"@fr;
301   rdfs:comment "Valeur d'un paramètre utilisé par une méthode pour faire une opération sur des
302   coordonnées."@fr;
302 rdfs:isDefinedBy <http://data.ign.fr/def/ignf>.
303
304 ######
305 ## Axioms here ##
306 ######
307
308
309 ignf:CompoundCRS owl:disjointWith ignf:GeodeticCRS .
310 ignf:VerticalCS owl:disjointWith ignf:VerticalCRS .
311 ignf:CartesianCS owl:disjointWith ignf:GeodeticCRS .
312 ignf:EllipsoidalCS owl:disjointWith ignf:GeodeticCRS .
313 ignf:PrimeMeridian owl:disjointWith ignf:Ellipsoid .
314 ignf:CRS owl:disjointWith ignf:CoordinateSystem .
315
316 ######
317 ### ----Properties here ----###
318 ######

```

```
319  
320  
321 ignf:domainOfValidity a owl:ObjectProperty ;  
322   rdfs:label "domain of validity"@en, "domaine de validité"@fr ;  
323   rdfs:comment "Zone ou intervalle de temps dans lequel un SRC est valide."@fr;  
324   rdfs:range ignf:Extent ;  
325   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .  
326  
327 ignf:includesSingleCRS a owl:ObjectProperty ;  
328   rdfs:label "includes single CRS"@en, "inclus un simple SRC"@fr ;  
329   rdfs:comment "Désigne les SRC simples qui composent un SRC composé."@fr;  
330   rdfs:domain ignf:CompoundCRS ;  
331   rdfs:range ignf:SingleCRS ;  
332   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .  
333  
334 ignf:baseCRS a owl:ObjectProperty ;  
335   rdfs:label "base CRS"@en, "SRC de base"@fr ;  
336   rdfs:comment "Désigne le SRC géodésique sur lequel repose un SRC projeté."@fr;  
337   rdfs:domain ignf:ProjectedCRS ;  
338   rdfs:range ignf:GeodeticCRS ;  
339   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .  
340  
341 ignf:coordinateSystem a owl:ObjectProperty ;  
342   rdfs:label "uses coordinate system"@en, "utilise le système de coordonnées"@fr ;  
343   rdfs:comment "Désigne le système de coordonnées utilisé par un SRC."@fr;  
344   rdfs:domain ignf:SingleCRS ;  
345   rdfs:range ignf:CoordinateSystem ;  
346   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .  
347  
348 ignf:ellipsoidalCS a owl:ObjectProperty ;  
349   rdfs:label "uses ellipsoidal CS"@en, "utilise le système de coordonnées ellipsoïdal"@fr;  
350   rdfs:comment "Désigne le système de coordonnées ellipsoïdal utilisé par un SRC géodésique."@fr;  
351   rdfs:subPropertyOf ignf:coordinateSystem;  
352   rdfs:domain ignf:GeodeticCRS ;  
353   rdfs:range ignf:EllipsoidalCS ;  
354   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .  
355  
356 ignf:cartesianCS a owl:ObjectProperty ;  
357   rdfs:label "uses Cartesian CS"@en, "utilise le système cartesien de coordonnées"@fr;  
358   rdfs:comment "Désigne le système de coordonnées cartésien utilisé par un SRC géodésique ou  
      projeté."@fr;  
359   rdfs:subPropertyOf ignf:coordinateSystem;  
360   rdfs:domain [ a owl:Class ;  
361     owl:unionOf (ignf:GeodeticCRS ignf:ProjectedCRS ) ] ;  
362   rdfs:range ignf:CartesianCS ;  
363   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .  
364  
365 ignf:verticalCS a owl:ObjectProperty ;  
366   rdfs:label "uses Vertical CS"@en, "utilise le système vertical de coordonnées"@fr;  
367   rdfs:comment "Désigne le système de coordonnées vertical utilisé par un SRC vertical."@fr;  
368   rdfs:subPropertyOf ignf:coordinateSystem;  
369   rdfs:domain ignf:VerticalCRS ;  
370   rdfs:range ignf:VerticalCS ;  
371   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .  
372  
373 ignf:datum a owl:ObjectProperty ;  
374   rdfs:label "uses datum"@en, "utilise le référentiel"@fr;  
375   rdfs:comment "Désigne le référentiel utilisé par un SRC simple."@fr;  
376   rdfs:domain ignf:SingleCRS;  
377   rdfs:range ignf:Datum ;  
378   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .  
379
```

```

380 ignf:geodeticDatum a owl:ObjectProperty ;
381   rdfs:label "uses geodetic datum"@en, "utilise le référentiel géodésique"@fr;
382   rdfs:comment "Désigne le référentiel géodésique utilisé par un SRC géodésique."@fr;
383   rdfs:subPropertyOf ignf:datum;
384   rdfs:domain ignf:GeodeticCRS;
385   rdfs:range ignf:GeodeticDatum ;
386   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
387
388 ignf:verticalDatum a owl:ObjectProperty ;
389   rdfs:label "uses vertical datum"@en, "utilise le référentiel vertical"@fr;
390   rdfs:comment "Désigne le référentiel vertical utilisé par un SRC vertical."@fr;
391   rdfs:subPropertyOf ignf:datum;
392   rdfs:domain ignf:VerticalCRS ;
393   rdfs:range ignf:VerticalDatum ;
394   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
395
396 ignf:axis a owl:ObjectProperty ;
397   rdfs:label "uses Axis"@en, "utilise l'axe"@fr;
398   rdfs:comment "Désigne un axe utilisé par un système de coordonnées ellipsoïdal ou cartésien."@fr;
399   rdfs:domain ignf:CoordinateSystem ;
400   rdfs:range ignf:CoordinateSystemAxis ;
401   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
402
403 ignf:primeMeridian a owl:ObjectProperty ;
404   rdfs:label "uses prime meridian"@en, "utilise le méridien origine"@fr;
405   rdfs:comment "Désigne le méridien origine d'un référentiel géodésique."@fr;
406   rdfs:domain ignf:GeodeticDatum ;
407   rdfs:range ignf:PrimeMeridian ;
408   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
409
410 ignf:ellipsoid a owl:ObjectProperty ;
411   rdfs:label "uses ellipsoid"@en, "utilise l'ellipsoïde"@fr;
412   rdfs:comment "Désigne l'ellipsoïde utilisé par un référentiel géodésique."@fr;
413   rdfs:domain ignf:GeodeticDatum;
414   rdfs:range ignf:Ellipsoid ;
415   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
416
417 ignf:geographicElement a owl:ObjectProperty ;
418   rdfs:label "geographic element"@en, "élément géographique"@fr;
419   rdfs:comment "Désigne le cadre englobant géographique d'une région de validité."@fr;
420   rdfs:domain ignf:Extent;
421   rdfs:range ignf:GeographicBoundingBox ;
422   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
423
424 ignf:semiMajorAxis a owl:ObjectProperty;
425   rdfs:label "semi major axis"@en, "demi grand axe"@fr;
426   rdfs:comment "Désigne la longueur du demi grand axe d'un ellipsoïde."@fr;
427   rdfs:domain ignf:Ellipsoid;
428   rdfs:range [a owl:Class ;
429     owl:intersectionOf
430       (qudt:QuantityValue
431         [ a owl:Restriction;
432           owl:allValuesFrom qudt:LengthUnit;
433           owl:onProperty qudt:unit
434         ] );
435   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
436
437 ignf:semiMinorAxis a owl:ObjectProperty;
438   rdfs:label "semi minor axis"@en, "demi petit axe"@fr;
439   rdfs:comment "Désigne la longueur du demi petit axe d'un ellipsoïde."@fr;
440   rdfs:domain ignf:Ellipsoid;
441   rdfs:range [a owl:Class ;

```

```

442    owl:intersectionOf
443        (qudt:QuantityValue
444            [ a owl:Restriction;
445                owl:allValuesFrom qudt:LengthUnit;
446                owl:onProperty qudt:unit
447                ] );
448    rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
449
450 ignf:inverseFlattening a owl:ObjectProperty;
451     rdfs:label "inverse flattening"@en, "aplatissement inverse"@fr;
452     rdfs:comment "Désigne la valeur d'aplatissement inverse d'un ellipsoïde, exprimée sous la forme d'un nombre ou d'un ratio (pourcentage, parties par million, etc.)."@fr;
453     rdfs:domain ignf:Ellipsoid;
454     rdfs:range [a owl:Class ;
455         owl:intersectionOf
456             (qudt:QuantityValue
457                 [ a owl:Restriction;
458                     owl:allValuesFrom qudt:CountingUnit;
459                     owl:onProperty qudt:unit
460                     ] );
461     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
462
463 ignf:greenwichLongitude a owl:ObjectProperty;
464     rdfs:label "greenwich longitude"@en, "longitude par rapport à Greenwich"@fr;
465     rdfs:comment "Désigne la longitude par rapport au méridien de Greenwich."@fr;
466     rdfs:domain ignf:PrimeMeridian;
467     rdfs:range [a owl:Class ;
468         owl:intersectionOf
469             (qudt:QuantityValue
470                 [ a owl:Restriction;
471                     owl:allValuesFrom qudt:AngleUnit;
472                     owl:onProperty qudt:unit
473                     ] );
474     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
475
476 ignf:usesMethod a owl:ObjectProperty ;
477     rdfs:label "uses method"@en, "utilise la méthode"@fr;
478     rdfs:comment "Désigne la méthode utilisée par une opération sur des coordonnées."@fr;
479     rdfs:domain ignf:CoordinateOperation;
480     rdfs:range ignf:OperationMethod ;
481     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
482
483 ignf:usesValue a owl:ObjectProperty ;
484     rdfs:label "uses value"@en, "utilise la valeur"@fr;
485     rdfs:comment "Désigne une valeur utilisée par une opération sur des coordonnées."@fr;
486     rdfs:domain ignf:CoordinateOperation;
487     rdfs:range ignf:OperationParameterValue ;
488     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
489
490 ignf:usesParameter a owl:ObjectProperty ;
491     rdfs:label "uses parameter"@en, "utilise le paramètre"@fr;
492     rdfs:comment "Désigne un paramètre utilisé par une méthode pour faire une opération sur des coordonnées."@fr;
493     rdfs:domain ignf:OperationMethod;
494     rdfs:range ignf:OperationParameter;
495     rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
496
497 ignf:parameter a owl:ObjectProperty ;
498     rdfs:label "parameter"@en, "paramètre"@fr;
499     rdfs:comment "Désigne le paramètre auquel est associée une valeur."@fr;
500     rdfs:domain ignf:OperationParameterValue;
501     rdfs:range ignf:OperationParameter;

```

```

502 rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
503
504 ignf:sourceCRS a owl:ObjectProperty ;
505   rdfs:label "source CRS"@en, "SRC source"@fr;
506   rdfs:comment "Désigne le SRC des données en entrée d'une opération."@fr;
507   rdfs:domain ignf:CoordinateOperation;
508   rdfs:range ignf:CRS;
509   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
510
511 ignf:targetCRS a owl:ObjectProperty ;
512   rdfs:label "target CRS"@en, "SRC cible"@fr;
513   rdfs:comment "Désigne le SRC des données en sortie d'une opération."@fr;
514   rdfs:domain ignf:CoordinateOperation;
515   rdfs:range ignf:CRS;
516   rdfs:isDefinedBy <http://data.ign.fr/def/ignf> .
517
518 #####--#
519 #####---Dataproperties here ---#####
520 ##########
521
522 ignf:westBoundLongitude a owl:DatatypeProperty;
523   rdfs:domain ignf:GeographicBoundingBox;
524   rdfs:label "west bound longitude"@en;
525   rdfs:label "longitude ouest minimale"@fr;
526   rdfs:comment "Désigne la longitude ouest minimale du cadre englobant."@fr;
527   rdfs:range xsd:double.
528
529 ignf:eastBoundLongitude a owl:DatatypeProperty;
530   rdfs:domain ignf:GeographicBoundingBox;
531   rdfs:label "east bound longitude"@en;
532   rdfs:label "longitude est maximale"@fr;
533   rdfs:comment "Désigne la longitude est maximale du cadre englobant."@fr;
534   rdfs:range xsd:double.
535
536 ignf:southBoundLatitude a owl:DatatypeProperty;
537   rdfs:domain ignf:GeographicBoundingBox;
538   rdfs:label "south bound latitude"@en;
539   rdfs:label "latitude sud minimale"@fr;
540   rdfs:comment "Désigne la latitude sud minimale du cadre englobant."@fr;
541   rdfs:range xsd:double.
542
543 ignf:northBoundLatitude a owl:DatatypeProperty;
544   rdfs:domain ignf:GeographicBoundingBox;
545   rdfs:label "north bound longitude"@en;
546   rdfs:label "latitude nord maximale"@fr;
547   rdfs:comment "Désigne la latitude nord maximale du cadre englobant."@fr;
548   rdfs:range xsd:double.
549
550 ignf:scope a owl:DatatypeProperty ;
551   rdfs:label "scope"@en, "portée"@fr ;
552   rdfs:comment "Désigne la portée de la référence."@fr;
553   rdfs:range xsd:string.
554
555 ignf:codeSpace a owl:DatatypeProperty;
556   rdfs:label "code space"@en, "espace de codage"@fr;
557   rdfs:comment "Désigne la règle ou l'autorité dont résulte la valeur de la référence."@fr;
558   rdfs:range xsd:string.
559
560 ignf:conversion a owl:DatatypeProperty;
561   rdfs:label "defined by conversion"@en, "défini par conversion"@fr;
562   rdfs:comment "Désigne la conversion utilisée pour définir un SRC projeté."@fr;
563   rdfs:domain ignf:ProjectedCRS;

```

```

564    rdfs:range xsd:anyURI.
565
566 ignf:axisAbbrev a owl:DatatypeProperty;
567   rdfs:label "axis abbreviation"@en, "abréviation utilisée pour désigner l'axe"@fr;
568   rdfs:comment "Désigne l'abréviation utilisée pour désigner l'axe."@fr;
569   rdfs:domain ignf:CoordinateSystemAxis;
570   rdfs:range xsd:string.
571
572 ignf:axisDirection a owl:DatatypeProperty;
573   rdfs:label "axis direction"@en, "direction de l'axe"@fr;
574   rdfs:comment "Désigne la direction de l'axe."@fr;
575   rdfs:domain ignf:CoordinateSystemAxis;
576   rdfs:range xsd:string.
577
578 ignf:isSphere a owl:DatatypeProperty;
579   rdfs:label "is sphere"@en, "est une sphère"@fr;
580   rdfs:comment "Indique si l'ellipsoïde est une sphère."@fr;
581   rdfs:domain ignf:Ellipsoid;
582   rdfs:range xsd:boolean.
583
584 ignf:epsgID a owl:DatatypeProperty;
585   rdfs:label "espg identifier"@en, "identifiant epsg"@fr;
586   rdfs:comment "Indique l'identifiant EPSG de la ressource."@fr;
587   rdfs:range xsd:string.
588
589 ignf:operationVersion a owl:DatatypeProperty;
590   rdfs:label "operation version"@en, "version de l'opération"@fr;
591   rdfs:comment "Indique la version d'une opération appliquée sur des coordonnées."@fr;
592   rdfs:domain ignf:CoordinateOperation;
593   rdfs:range xsd:string.
594
595 ignf:methodFormula a owl:DatatypeProperty;
596   rdfs:label "method formula"@en, "formule de la méthode"@fr;
597   rdfs:comment "Indique la formule utilisée par une méthode pour faire opération sur des coordonné
      es."@fr;
598   rdfs:domain ignf:OperationMethod;
599   rdfs:range xsd:string.
600
601 ignf:sourceDimension a owl:DatatypeProperty;
602   rdfs:label "source CRS dimension"@en, "dimension du SRC source"@fr;
603   rdfs:comment "Indique la dimension du SRC des données en entrée d'une opération."@fr;
604   rdfs:domain ignf:OperationMethod;
605   rdfs:range xsd:integer.
606
607 ignf:targetDimension a owl:DatatypeProperty;
608   rdfs:label "target CRS dimension"@en, "dimension du SRC cible"@fr;
609   rdfs:comment "Indique la dimension du SRC des données en sortie d'une opération."@fr;
610   rdfs:domain ignf:OperationMethod;
611   rdfs:range xsd:integer.

```

Listing B.3: Formal definition of the French administrative units in turtle. The current version is deployed at <http://data.ign.fr/def/geofla>

## B.4 Vocabulary for Visualization applications

1 @prefix dct: <<http://purl.org/dc/terms/>>.  
 2 @prefix dcat: <<http://www.w3.org/ns/dcat#>>.  
 3 @prefix foaf: <<http://xmlns.com/foaf/0.1/>>.

```

4  @prefix dctype: <http://purl.org/dc/dcmitype/>.
5  @prefix xsd: <http://www.w3.org/2001/XMLSchema#>.
6  @prefix org: <http://www.w3.org/ns/org#>.
7  @prefix dvia: <http://purl.org/ontology/dvia#>.
8  @prefix owl: <http://www.w3.org/2002/07/owl#>.
9  @prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#>.
10 @prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>.
11 @prefix skos: <http://www.w3.org/2004/02/skos/core#>.
12
13 #---Meta data goes here -----
14 <http://purl.org/ontology/dvia>
15   a owl:Ontology;
16   rdfs:comment "dvia is an RDF vocabulary designed to facilitate interoperability between
17     visualization applications published on the Web."@en;
18   dct:title "The visualization vocabulary for LOD applications"@en;
19   dct:description "dvia is an RDF vocabulary designed to facilitate interoperability between
20     visualization applications published on the Web."@en;
21   dct:issued "2012-07-25"^^xsd:date;
22   dct:modified "2012-10-31"^^xsd:date;
23   dct:modified "2012-11-27"^^xsd:date;
24   dct:modified "2013-02-01"^^xsd:date;
25   dct:modified "2013-10-09"^^xsd:date;
26   dct:title "Ontologie des applications de visualisation sur le web"@fr;
27   dct:creator [foaf:mbox "atemezin@eurecom.fr"; foaf:name "Ghislain Atemezing"];
28   dct:contributor [foaf:mbox "rtroncy@eurecom.fr"; foaf:name "Raphael Troncy"];
29   dct:license <http://www.opendatacommons.org/licenses/pddl/1.0/> .
30
31 #---Classes here -----
32 <http://purl.org/ontology/dvia#Application> a rdfs:Class,
33   owl:Class;
34   rdfs:comment "The application or the mashup developed for demo-ing or consuming data in LD
35     fashion"@en;
36   rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
37   rdfs:label "Application"@en;
38   rdfs:label "Application"@fr;
39   rdfs:subClassOf dctype:Software.
40
41 <http://purl.org/ontology/dvia#Platform> a rdfs:Class,
42   owl:Class;
43   rdfs:comment "The platform where to host or use the application, could be on the web (firefox,
44     chrome, IE, etc..) or mobile (android, etc..) or event desktop"@en;
45   rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
46   rdfs:label "Plate forme"@fr;
47   rdfs:label "Platform"@en.
48
49 <http://purl.org/ontology/dvia#VisualTool> a rdfs:Class,
50   owl:Class;
51   rdfs:comment "The tool or library used to build the application"@en;
52   rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
53   rdfs:label "visual Tool"@en;
54   rdfs:label "Outil de visualisation"@fr;
55   owl:disjointWith dvia:Application.
56
57 # Object properties here -----
58
59 <http://purl.org/ontology/dvia#author> a rdf:Property,
60   owl:ObjectProperty;
61   rdfs:comment "links to the authors of the application or the tools, libraries"@en;

```

```

62 rdfs:domain dvia:Application;
63 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
64 rdfs:label "author"@en;
65 rdfs:label "auteur"@fr;
66 rdfs:range foaf:Person.
67
68 <http://purl.org/ontology/dvia#platform> a rdf:Property,
69   owl:ObjectProperty;
70 rdfs:comment "This property links the application to a platform to actually use the application."@en;
71 rdfs:domain dvia:Application;
72 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
73 rdfs:label "platform"@en;
74 rdfs:label "plate forme"@fr;
75 rdfs:range dvia:Platform.
76
77 <http://purl.org/ontology/dvia#consumes> a rdf:Property,
78   owl:ObjectProperty;
79 rdfs:comment "links to the dataset used to make the application, and could be of different types or formats"@en;
80 rdfs:domain dvia:Application;
81 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
82 rdfs:label "consumes"@en;
83 rdfs:label "consomme"@fr;
84 rdfs:range dcat:Dataset.
85
86 <http://purl.org/ontology/dvia#designBy> a rdf:Property,
87   owl:ObjectProperty;
88 rdfs:comment "links to the organization which builds the application"@en;
89 rdfs:domain dvia:Application;
90 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
91 rdfs:label "design By"@en;
92 rdfs:label "implémente par"@fr;
93 rdfs:range org:Organization.
94
95 <http://purl.org/ontology/dvia#usesTool> a rdf:Property,
96   owl:ObjectProperty;
97 rdfs:comment "This property links to the tools or libraries used for the application"@en;
98 rdfs:domain dvia:Application;
99 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
100 rdfs:label "uses Tool"@en;
101 rdfs:label "utilise l'outil"@fr;
102 rdfs:range dvia:VisualTool.
103
104 <http://purl.org/ontology/dvia#hasLicense> a rdf:Property,
105   owl:ObjectProperty;
106 rdfs:comment "This property links to the license of the application"@en;
107 rdfs:domain dvia:Application;
108 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
109 rdfs:label "hasLicense"@en;
110 rdfs:label "a pour license"@fr;
111 rdfs:range dct:License.
112
113
114 # Data properties here -----
115
116 <http://purl.org/ontology/dvia#preferredNavigator> a rdf:Property,
117   owl:DatatypeProperty;
118 rdfs:comment "The name of the preferred navigator to be used by the application"@en;
119 rdfs:domain dvia:platform;
120 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
121 rdfs:label "preferred navigator"@en;

```

```

122 rdfs:label "navigateur préféré"@fr;
123 rdfs:range xsd:string.
124
125 <http://purl.org/ontology/dvia#alternativeNavigator> a rdf:Property,
126   owl:DatatypeProperty;
127 rdfs:comment "The name of the alternate navigator if applicable"@en;
128 rdfs:domain dvia:platform;
129 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
130 rdfs:label "alternative navigator"@en;
131 rdfs:label "navigateur alternatif"@fr;
132 rdfs:range xsd:string.
133
134 <http://purl.org/ontology/dvia#businessValue> a rdf:Property,
135   owl:DatatypeProperty;
136 rdfs:comment "The business value of the application; generally could be commercial or free. Also
137   depending on the license"@en;
138 rdfs:domain dvia:Application;
139 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
140 rdfs:label "businessValue"@en;
141 rdfs:label "valeur commerciale"@fr;
142 rdfs:range xsd:string.
143
144 <http://purl.org/ontology/dvia#url> a rdf:Property,
145   owl:DatatypeProperty;
146 rdfs:comment "the url of the application."@en;
147 rdfs:domain dvia:Application;
148 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
149 rdfs:label "url"@en;
150 rdfs:label "url"@fr;
151 rdfs:range xsd:anyURI.
152
153 <http://purl.org/ontology/dvia#keyword> a rdf:Property,
154   owl:DatatypeProperty;
155 rdfs:comment "keywords used for the application."@en;
156 rdfs:domain dvia:Application;
157 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
158 rdfs:label "keyword"@en;
159 rdfs:label "mot clé"@fr;
160 rdfs:range xsd:string.
161
162 <http://purl.org/ontology/dvia#scope> a rdf:Property,
163   owl:DatatypeProperty;
164 rdfs:comment "The scope or domain of the application."@en;
165 rdfs:domain dvia:Application;
166 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
167 rdfs:label "scope"@en;
168 rdfs:label "domain d'usage"@fr;
169 rdfs:range skos:Concept.
170
171 <http://purl.org/ontology/dvia#view> a rdf:Property,
172   owl:DatatypeProperty;
173 rdfs:comment "The types of view available in the application, such as maps, charts, graphs, etc."@
174   en;
175 rdfs:domain dvia:Application;
176 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
177 rdfs:label "view"@en;
178 rdfs:label "vue"@fr;
179 rdfs:range xsd:string.
180
181 <http://purl.org/ontology/dvia#system> a rdf:Property,
182   owl:DatatypeProperty;
183 rdfs:comment "The operating system where the application runs."@en;

```

```
182 rdfs:domain dvia:Platform;
183 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
184 rdfs:label "system"@en;
185 rdfs:label "système"@fr;
186 rdfs:range xsd:string.
187
188 <http://purl.org/ontology/dvia#datasetDescription> a rdf:Property,
189 owl:DatatypeProperty;
190 rdfs:comment "Property for a given descriptive informations of the datasets used for making the
application. It could be used when no more details are given on the datasets like URL,
formats, etc."@en;
191 rdfs:domain dvia:Application;
192 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
193 rdfs:label "description du jeu de données"@fr;
194 rdfs:label "dataset description"@en.
195
196
197 <http://purl.org/ontology/dvia#libUrl> a rdf:Property,
198 owl:DatatypeProperty;
199 rdfs:comment "The url to the page describing the library or the tool for visualization."@en;
200 rdfs:domain dvia:VisualTool;
201 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
202 rdfs:label "library url"@en;
203 rdfs:label "url de la librairie"@fr;
204 rdfs:range xsd:anyURI.
205
206 <http://purl.org/ontology/dvia#downloadUrl> a rdf:Property,
207 owl:DatatypeProperty;
208 rdfs:comment "The download url of the tool for visualization."@en;
209 rdfs:domain dvia:VisualTool;
210 rdfs:isDefinedBy <http://purl.org/ontology/dvia>;
211 rdfs:label "download url"@en;
212 rdfs:label "url téléchargement"@fr;
213 rdfs:range xsd:anyURI.
```

Listing B.4: Formal definition of DVIA vocabulary in turtle. The current version is deployed at <http://purl.org/ontology/dvia>.



## Part IV

# Résumé en Français





# Résumé

Au cours de ces dernières années, le domaine de l'Open Data a reçu une attention croissante de la part des administrations publiques qui veulent tirer avantage de la publication de données ouvertes sur le Web. Les bénéfices supposés de cette ouverture pour les citoyens font référence à une meilleure transparence dans les prises de décisions publiques, à une meilleure gouvernance ou encore au développement d'un éco-système numérique qui tirerait un profit économique des applications analysant ces données. Cependant, la réalité montre que la simple ouverture et la publication de données par les administrations ne sont pas suffisantes au regard des défis liés à la variété des formats (XML, CSV, Excel, PDF, Shape), des méthodes d'accès (API, base de données) et à l'absence de nomenclature qui permettrait une meilleure réutilisation et interconnexion avec d'autres jeux de données. Dans cette thèse, nous explorons comment l'utilisation des standards et des technologies du web sémantique peut aider à résoudre les problèmes causés par l'hétérogénéité et la diversité des formats de données et des structures de représentations dans le domaine géographique.

Cette thèse applique les principes des « données liées » dans le domaine de l'information géographique, un domaine clef pour les administrations publiques qui couvrent, par définition, un territoire. En particulier, nous traitons de trois aspects essentiels dans le workflow de traitement et de publication de données géo-spatiales et de leur consommation (visualisation), avec des scénarios d'utilisation issus de l'Institut Nationale de l'Information Géographique et Forestière (IGN) : (1) Comment représenter efficacement et stocker des données géospatiales sur le Web pour assurer des applications interopérables ? (2) Quelles sont les meilleures options pour un utilisateur pour interagir avec des données sémantiques interconnectées ? (3) Quels mécanismes peuvent être mis en place pour aider à la préservation des données structurées de haute qualité sur le Web ?

Nos contributions sont structurées en trois grandes parties correspondantes aux problématiques susmentionnées, avec des applications spécifiques dans le domaine géographique. Nous proposons et développons trois vocabulaires pour représenter des systèmes de coordonnées de référence (CRS), des entités topographiques et la géométrie associée à ces entités. Ces ontologies étendent d'autres vocabulaires existants et ajoutent deux avantages

supplémentaires : l'utilisation explicite de CRS identifiés par des URIs pour représenter la géométrie, et la capacité de décrire des géométries structurées en RDF. Nous avons ainsi publié la base de données GEOFLA, en contribuant et utilisant la plate-forme Datalift, un outil permettant de convertir et publier des données brutes en données liées. Nous avons également évalué de manière systématique la performance des points d'accès SPARQL pour traiter des requêtes spatiales.

Concernant la « consommation » de données RDF, après avoir examiné les différentes catégories des outils de visualisation (génériques et spécifiques à des jeux de données), nous proposons un vocabulaire pour décrire les applications de visualisation (DVIA). En outre, nous formalisons et mettons en œuvre un workflow pour visualiser des données sémantiques interconnectées à travers l'outil LDVizWiz, un assistant de visualisation générique de données liées sur le Web.

La dernière partie de la thèse décrit des contributions au catalogue des vocabulaires liées (LOV) ainsi qu'une proposition originale pour utiliser LOV avec une méthodologie de création d'ontologie telle que NeOn dans le but d'améliorer la réutilisation des vocabulaires. Nous proposons une heuristique pour aligner les vocabulaires et un classement de ces derniers en fonction de métriques liées au contenu de l'information des termes définis dans les vocabulaires. Enfin, la thèse apporte des réponses sur la façon de vérifier la compatibilité des licences entre les vocabulaires et les jeux de données dans le workflow de publication. Tout au long de la thèse, nous démontrons les avantages de l'utilisation des technologies sémantiques et des standards du W3C pour mieux découvrir, interconnecter et visualiser les données géospatiales gouvernementales pour leur publication sur le Web.

Our main concern is to tackle the problems within the workflow of publication in two directions, more likely to happen at the beginning and the end :

- (i) Geographic Information on the Web of Data : as an application of the life-cycle of publishing geodata.
- (ii) Visualization tools for building innovative applications consuming structured data : as for leveraging the process of creating applications on-top of semantic data to highlight some relevant knowledge to the users.

Notre but principal est de résoudre les problèmes liés à la chaîne de publication dans deux directions, plus susceptible de se produire en début et en fin de chaîne :

- (i) l'information géographique sur le Web de données : comme une application du cycle de vie de publication des données geo-spatiales.
- (ii) des outils de visualisation pour créer des applications innovantes utilisant les données structurées : comme pour tirer parti du processus de création des applications au-dessus des données sémantiques afin de mettre en exergue des connaissances pertinentes pour les utilisateurs.

# Questions de recherche

Dans cette thèse, nous proposons des solutions dans les défis liés à la publication des données géographique sur le Web des données, qui sont les suivantes :

1. *Vocabulaires* : Comment modéliser l'information géographique sur le Web ? Comment évaluer les ontologies du domaine géographique ? Comment sérialiser les géométries complexes dans un environnement comme le Web ?
2. *Languages de requêtes* : Comment pouvons-nous écrire des requêtes efficaces qui ciblent les données géospatiales sur le Web ? Comment pouvons-nous stocker et indexer les géodonnées en RDF ?
3. *Données* : Comment pouvons-nous extraire et convertir les géodonnées pour publication sur le Web ? Quelles sont les bonnes pratiques pour représenter des géométries complexes sur le Web ? Comment pouvons-nous intégrer pleinement la compatibilité des systèmes de coordonnées sur des jeux de données ?
4. *Publication* : Comment pouvons-nous développer des environnements qui passent à l'échelle pour couvrir le chaîne de publication des géodonnées ? Quels sont les triples stores appropriées pour le stockage des géodonnées ? Quelles sont les métriques à utiliser pour l'interconnexion de différentes ressources de géodonnées sur le Web ?
5. *Applications et interfaces utilisateurs* : Comment pouvons-nous générer des visualisations de données géospatiales liées entre elles ? Quels sont les API de haut niveau appropriées qui facilitent le développement d'interfaces utilisateur pour les données géospatiales ? Pouvons-nous réutiliser les outils de cartographie existants tels que Google Maps, Bing Maps ou OpenStreetMap ?

Dans cette thèse, nous abordons les enjeux de publication des données du point de vue tant par les éditeurs que des utilisateurs. Les éditeurs et les utilisateurs ont besoin de solutions pragmatiques qui les aident à choisir un vocabulaire, trouver un outil pour convertir des fichiers shape ShapeFiles selon des vocabulaires existants, puis transformer en RDF et publier les données suivant des bonnes pratiques.

Après la publication du jeux de données sur le Web, les éditeurs et les utilisateurs doivent comprendre ces données pendant que les développeurs doivent pouvoir créer des applications. Le Web commence à contenir de plus en plus de données structurées, qui ne sont pas toujours exploitées par les utilisateurs finaux, à cause de la complexité dans l'usage du modèle RDF et de son langage de requête, le SPARQL. Ainsi, il est important de créer des visualisations pour explorer, analyser et montrer les bénéfices du Linked Data aux non-experts. Dans ce processus, il existe des questions de recherches à résoudre telles que :

- Comment trouver visualisations adaptées selon les jeux de données tout en masquant la complexité du langage de requêtes SPARQL ?
- Quelles sont les propriétés importantes pour visualiser les ressources du Web, en fonction du domaine et des attentes des utilisateurs ?
- Comment combler le fossé entre les outils traditionnels existants de visualisation de l'information, la plupart du temps aux formats CSV/XLS, JSON ou formats propriétaires pour intégrer facilement le modèle de données RDF en entrée ?
- Comment développer des applications interopérables sur les catalogues de données gouvernementaux en Open Data ? Comment réutiliser les applications existantes ?

Tout en essayant de répondre aux défis ci-dessus mentionnés, nous exposons l'état de l'art et approches existantes dans le domaine des visualisations de données liées.

## Contributions

Les contributions de cette thèse sont organisées en trois parties principales : la modélisation et la publication des données géospatiales, la visualisation de données et des applications sur le Web et la contribution dans les standards.

### Modélisation et Publication des données géospatiales

La géolocalisation est cruciale pour de nombreuses applications tant pour agents humains que les logiciels. De plus en plus des masses de données sont ouvertes et interconnectées sur le Web. Une modélisation des données géographiques de manière efficace en réutilisant autant que possible des ontologies ou des vocabulaires existants qui décrivent à la fois les fonctionnalités géospatiales et leurs formes. Dans la première partie de notre travail, nous examinons différentes approches de modélisation utilisées dans les systèmes d'information géographique (SIG) et la communauté des données ouvertes (LOD). Notre objectif est de contribuer aux efforts réels dans la représentation des objets géographiques avec des attributs tels que l'emplacement, les points d'intérêt (POI), et les adresses sur le Web de données. Nous nous

concentrons sur le territoire français et nous fournissons des exemples de vocabulaires représentatifs qui peuvent être utilisés pour décrire les objets géographiques. Nous proposons quelques alignements entre différents vocabulaires (DBpedia, schema.org, LinkedGeoData, Foursquare, etc.) afin de permettre l'interopérabilité tout en interconnectant les géodonnées en France avec d'autres jeux de données.

Concernant cet aspect de notre recherche, nos contributions sont les suivantes :

1. Nous avons proposé et développé une ontologie décrivant les caractéristiques et les points d'intérêt pour le territoire français, en réutilisant une taxonomie existante (GeOnto) en l'alignant sur d'autres vocabulaires connexes dans le domaine de la géolocalisation.
2. Nous avons étudié comment étendre les vocabulaires existants dans le domaine géographique afin de prendre en compte une modélisation efficace des géométries complexes. Ce faisant, nous abordons les questions de représentation de géométrie complexe dans le Web de données, décrivant l'état de mise en oeuvre des fonctions géospatiales dans triples stores et une comparaison avec la nouvelle norme GeoSPARQL. Nous faisons enfin quelques recommandations et plaidons pour la réutilisation des vocabulaires plus structurées pour la publication d'entités topographiques pour mieux répondre aux exigences des données issues de IGN-France.
3. Nous avons fait une étude comparative des triples stores, comparant leur capacité de stockage des informations spatiales et leur implémentation des fonctions topologiques ra rapport à celles déjà existantes dans les normes de l'Open Geospatial Consortium (OGC)<sup>1</sup>.
4. Nous avons conçu et développé des vocabulaires pour décrire les géométries complexes avec différents systèmes de coordonnées, avec application directe aux unités administratives françaises.
5. Nous avons interconnecté des géodonnées du contexte français avec des jeux de données géospatiales existantes sur le Web, tels que LinkedGeo-data, GADM, NUTS et Geonames.
6. Nous avons contribué à la création du “nuage données” ( LOD Cloud) représentant la publication de 8 jeux de données, soit 340 millions de triplets couvrant le territoire français.

Consommer des données sur le Web grâce à des visualisations comporte autant de défis que les applications doivent de conformer à la structure du graphe RDF, la sémantique sous-jacente du jeu de données et de l'interaction homme-machine pour comprendre facilement de quoi traitent ces données. Dans la section suivante, nous présentons nos contributions sur la visualisation.

---

1. <http://www.opengeospatial.org/>

## Outils de visualisation des données gouvernementales liées

Nous étudions d'abord quelques applications innovantes qui ont été développées sur des jeux de données publiées en Open Data par les gouvernements (Royaume-Uni, USA, France) et des administrations locales. Nous avons ensuite dérivé et proposé 8 cas d'utilisation (scénarios) qui peuvent être développés pour consommer des données provenant des différents fournisseurs principaux en France : INSEE, DILA, IGN, FING, etc. Nous mentionnons que les cas d'utilisation les plus intéressants sont ceux qui montrent la valeur ajoutée des jeux de données interconnectées. Ces scénarios développés et déployés, peuvent être utiles pour montrer les avantages de données liées dans une variété de domaines tels que l'éducation, le tourisme, le patrimoine culturel, les administrations civiles, les tribunaux, la médecine, etc.

En ce qui concerne les outils utilisés pour la visualisation, nous avons identifié et classer en deux catégories, en fournissant pour chacun d'eux des exemples pertinents : (i) - des outils qui fonctionnent sur des données RDF, et (ii) des outils qui fonctionnent sur d'autres formats structurés. Nous proposons donc des critères de base pour évaluer un outil de visualisation de donnée en général, avec des poids attachés à chaque critère.

Nos contributions sur la visualisation sont les suivantes :

1. Nous avons construit une application des élections présidentielles du premier tour français en 2012 en utilisant les données de <http://data.gouv.fr> et d'autres institutions publiques. L'application disponible à <http://www.eurecom.fr/~atemezin/DemoElection/> a été construit avec l'outil Exhibit. Il vise à mettre en valeur l'intégration des jeux de données hétérogènes : les résultats politiques en CSV, le taux de chômage, les données des candidats, les informations des départements de France provenant des données DBpedia. L'utilisateur peut filtrer par image du candidat, le taux de chômage et par département pour voir les scores, avec des informations plus enrichies sur le département.
2. Nous avons mis en place un outil générique pour explorer les géodonnées sur une carte, en fonction de la détection automatique des données via de requêtes SPARQL dans le nuage LOD contenant des jeux de données géospatiales.
3. Nous avons développé une application consommatrice de données géospatiales et statistiques combinant plusieurs jeux de données dans l'éducation de provenant du portail <http://data.gouv.fr> .
4. Nous avons développé une application sur les événements dans une conférence avec leurs médias supports réconciliés provenant de nombreuses plates-formes sociales (Instagram, Twitter, etc.).
5. Nous avons implanté un vocabulaire pour structurer les applications sur le Web de données. Le vocabulaire peut être utilisé pour découvrir des outils visuels ou graphiques utilisés pour créer des applications.

6. Nous avons implémenté un plugin générique pour annoter des applications développées pour des hackathon pouvant être inclus dans une page web, permettant la génération de contenu structuré de pages Web en utilisant le vocabulaire développé.
7. Nous avons mis en place un assistant qui analyse un jeu de données RDF et recommande une visualisation basée sur des catégories pré-définies, en utilisant des requêtes SPARQL génériques pour faciliter l'exploration des jeux de données publiés sur le LOD.

## Contributions aux standards

We contributed to the W3C Government Linked Data Working Group (GLD WG)<sup>2</sup> activity from July 2011 until December 2013. The objective of the Working Group was to “provide standards and other information which help governments around the world publish their data as effective and usable Linked Data using Semantic Web technologies”.

Nous avons contribué aux activités du groupe du W3C sur les données gouvernementales liées de travail (GLD WG)<sup>3</sup> de juillet 2011 jusqu'à décembre 2013. L'objectif du Groupe de travail était de «fournir des normes et d'autres informations qui aident les gouvernements à travers le monde dans la publication de leurs données aussi efficace qu'utilisable à l'aide des technologies du Web sémantique ».

Nous avons contribué à trois groupes de travail, avec en particulier dans deux documents :

1. Un glossaire<sup>4</sup> pour la description des termes utilisés dans le domaine du Linked data pour les potentiels producteurs et consommateurs de données gouvernementales sur le Web
2. Un document sur les bonnes pratiques de publication des données gouvernementales sur le Web<sup>5</sup>

En ce qui concerne l'utilisation de vocabulaires standards, nous avons contribué à :

- Proposer une méthode pour harmoniser les préfixes sur le Web de données avec deux services : Linked Open vocabulaires (LOV)<sup>6</sup> prefix.cc<sup>7</sup>. le premier service est actuellement un catalogue à jour des vocabulaires utilisés sur le Web, tandis que le dernier est un service pour les développeurs pour choisir, valider et chercher des préfixes pour leurs ressources ou ontologies. L'approche proposée peut être étendue à tout

---

2. <http://www.w3.org/2011/gld/>  
 3. <http://www.w3.org/2011/gld/>  
 4. <http://www.w3.org/TR/ld-glossary/>  
 5. <http://www.w3.org/TR/ld-bp/>  
 6. <http://lov.okfn.org/dataset/lov/>  
 7. <http://prefix.cc>

le catalogue du vocabulaire tant que les vocabulaires remplissent les conditions pour être insérées dans le catalogue LOV.

- Concevoir et mettre en œuvre une nouvelle méthode de classement des vocabulaires sur la base des métriques du contenu de l'information et de l'information partitionnée.
- Nous avons développé un outil qui détermine en temps réel si les différentes licences présentes dans un jeu de données et les vocabulaires associés sont soit compatible ou non.

# Plan de la Thèse

Dans la première partie de cette thèse , nous nous concentrons sur l'étude des différents modèles et vocabulaires pour représenter la géographie et de la géométrie . Nous étudions les points d'accès aux données et décrivons les problèmes particuliers tels que les systèmes de coordonnées, et mettons en évidence nos contributions dans ce dommaine : création de nouveaux vocabulaires réutilisant les vocabulaires existants, implémentation d'un convertisseur en ligne entre des différents systèmes de coordonées, etc. Nous décrivons également comment des jeux de données géographiques peuvent ensuite être convertis en RDF en utilisant le processu d'élévation des données du projet Datalift afin de leur publication sur le Web. Nous montrons ensuite comment ces jeux de données peuvent être alignées entre elles et concluons par une analyse approfondie de ces alignements dans le cas des jeux de données de cartographie française fournis par l'Institut Géographique et Forestière (IGN -France ) .

Plus précisement :

**Le Chaptitre 1** décrit les limites actuelles de représentation des géo-données sur le Web et notre contribution sur les différents vocabulaires pour représenter les géométries, les systèmes de coordonnées de référence et les ressources topographiques. Nous proposons également des bonnes pratiques pour la publication des données géospatiales sur le Web. **Le Chapitre 2** met l'accent sur les outils de publication et des requêtes d'interrogation des géodonnées, leurs différences et leurs applications. Nous décrivons la plate-forme Datalift, une plate-forme ouverte servant de catalyseur des sources de données brutes vers des données sémantiques et interconnectés. Après avoir comparé Datalift avec Geoknow, nous l'appliquons dans le processus de publication d'unités administratives et le Gazetteer français. Nous présentons ensuite l'état du nuage français LOD (FrLOD) des données liées et des exemples de requêtes sur des géométries structurées publiées dans le point d'accès <http://data.ign.fr>.

Dans la deuxième partie de la thèse, nous couvrons trois principales questions relatives à la façon de présenter les données en RDF aux utilisateurs finaux. Tout d'abord, nous présentons l'état de l'art des outils et des solutions pour la représentation visuelle et l'exploration des données en RDF (Visualbox, LODSpeaKr, Map4RDF, le modèle “Linked Data Visualization”,

etc.). Ensuite, nous présentons notre contribution : un assistant pour faciliter les visualisations automatique des points d'accès aux données sur le Web, y compris le vocabulaire spécifiant les visualisations et le prototype implémenté. Par la suite, nous présentons deux applications dans les domaines événementiel et statistique pour mettre en exergue de manière innovante la réutilisation des jeux de données liés. Enfin, nous implementons un algorithme permettant de révéler les propriétés les plus "importantes" des entités des ressources pour leur visualisation en partant de la Base de Connaissance de Google (GKP) ainsi qu'une évaluation faite sur les préférences des utilisateurs.

Cette partie est divisée en deux chapitres :

Le **Chapitre 3** fournit une revue de littérature sur des outils de visualisation et les applications, avec leurs limites. Nous décrivons également l'état de l'art des applications sur le Web et proposons une classification des «Applications de données liées».

Dans le **Chapitre 4**, nous présentons notre contribution sur de nouvelles approches pour générer des visualisations et des applications. Nous proposons tout d'abord une nouvelle approche pour les visualisations basées sur des catégories. Nous montrons ensuite une application dans le domaine géographique. Deux applications liées aux événements et aux statistiques sont également décrites. Enfin, nous proposons la façon d'améliorer la découverte d'applications dans les événements Open Data grâce à un modèle et un plugin universel pour annoter des pages Web en RDFa.

Dans la dernière partie de la thèse dans le **Chapitre 5**, nous décrivons diverses contributions aux vocabulaires ouverts liés (description du catalogue, les publications des vocabulaires, des API et des points d'accès) : l'harmonisation des prefixes des vocabulaires, les métriques pour classer les vocabulaires en utilisant le contenu de l'information.

Dans le **Chapitre 6**, nous présentons quelques idées sur la vérification de compatibilité des licences entre les vocabulaires et les jeux de données en utilisant la logique déontique en créant un outil en ligne pour la détection automatique des licences sur les données du Web.

Dans le **Chapitre 7**, nous concluons en mettant en évidence certaines limites et perspectives pour de nouvelles directions de recherche.

# Partie I : Intégration des données geo-spatiales sur le Web

Cette partie est divisée en deux chapitres et consacrée à l'état de l'art sur les formats et les différents vocabulaires utilisés dans la littérature pour

## Chapitre I

Dans ce chapitre, nous faisons une revue de la littérature des formats et des différents vocabulaires utilisés pour modéliser des données géospatiales sur le Web, en distinguant deux types de géoréférencement : direct et indirect. Ensuite, nous identifions certaines limitations liés à l'absence d'une référence explicite du Système de coordonné géographique (SCG) dans les jeux de données actuellement publiés sur le Web. Nous proposons ensuite un service REST pour la conversion entre différents SCG pour aider les éditeurs à être capable de gérer différentes projections dans les jeux données. En outre, nous proposons et implémentons trois vocabulaires pour les géométries, les SCG et les entités topographiques qui sont en ligne aux adresses respectives à <http://data.ign.fr/def/geometrie>, <http://data.ign.fr/def/ignf> et <http://data.ign.fr/def/topo>. Les vocabulaires étendent ceux existants et intègrent deux avantages supplémentaires : un usage explicite de SCG identifiés par des URIs pour la géométrie et la capacité à décrire des géométries structurées en RDF. Certains de nos résultats et la description du modèle sont en cours de discussion pour la standardisation au W3C, comme par exemple étendre le standard GeoSPARQL pour intégrer de manière plus explicite les coordonnées géographiques.

## Chapitre II

Dans ce chapitre, nous présentons une étude des outils d'extraction et de conversion de données géospatiales en RDF. Ensuite, nous décrivons GeomRDF, un outil développé au sein du projet Datalift qui va au-delà de l'état

de l'art en fournissant des géométries structurées et conformes au standard GeoSPARQL. En outre, nous présentons les limites des modèles de données existants en suggérant des recommandations aux éditeurs de géodonnées sur les aspects de stockage de gros volumes de donnés. De même, une description détaillée de l'outil Datalift utilisé pour publier des données sur le Web est fournie, avec une attention particulière sur notre contribution à la construction du nuage des données du Linked Open Data sur des données du territoire français avec des jeux de données en 4-5 étoiles selon les principes de données liées. Enfin, nous montrons quelques cas d'utilisation du monde réel des requêtes SPARQL faisant usage tour à tour de la géométrie structurée ou des fonctions géospatiales intégrées dans le triple store. Selon les besoins des utilisateurs et les jeux de données sous-jacentes, l'utilisateur peut choisir entre la simplicité du langage de requête SPARQL, avec des limitations au niveau du triple store (par exemple, lors de l'usage des fonctions géospatiales intégrées), ou l'expressivité du vocabulaire que nous proposons (`geom`), comme critère dans le choix du triple store et du stockage des données géospatiales.

# Partie II : Visualisation des graphes de données sur le Web

## Chapitre III

Dans ce chapitre, nous décrivons les différents outils utilisés pour la visualisation des données structurées et des graphes. Nous discutons également de différents types d'applications actuellement construites basées sur des initiatives de données ouvertes du gouvernement en Open Data. Le but de cet état de l'art est de proposer de nouvelles approches de génération et des outils de visualisations et des applications sur le Web de données. Nous avons conçu et mis en œuvre un vocabulaire, DVIA, qui vise à modéliser des applications pour plus d'interopérabilité et de découverte d'applications et d'outil de visualisations sur le Web.

## Chapitre IV

Dans ce chapitre, nous avons présenté une approche pour créer des visualisations au-dessus de données liées basés sur les technologies du Web sémantique. Nous avons d'abord défini sept catégories des entités qui peuvent être associés à la visualisation des jeux de données, et nous proposons de les mapper à d'autres vocabulaires de domaine. Nous présentons ensuite une description des principales composantes d'un assistant (wizard) de visualisation dans le contexte de Linked Data. Nous décrivons une implémentation en JavaScript comme *preuve-de-concept* de notre proposition, avec les avantages d'être disponible en ligne et extensible. Nous pensons qu'un tel outil peut être facilement intégré dans une chaîne globale de publication données sur le Web, tels que Datalift ou GeoKnow. En outre, nous avons effectué des expériences sur le graphe de connaissances de Google pour détecter des propriétés importantes à visualiser dans des entités, et avons évalué en fonction aux préférences des utilisateurs. Ensuite, nous avons présenté deux applications dans le domaine des statistiques et des événements, consommant différents jeux de données en RDF sur le scénario des données du monde réel. Nous avons discuté de la façon d'améliorer les applications développées dans des

contextes de hackathon, en proposant un vocabulaire et un outil pour peupler le modèle en utilisant un plugin universel. Des exemples d'événements passés ont déjà été transformés de manière semi-automatique utilisant à la fois le vocabulaire et le plugin.

# Partie III : Contribution au catalogue des vocabulaires liés

## Chapitre V

Nous avons présenté dans ce chapitre notre contribution au catalogue des vocabulaires ouverts et liés, comme partie d'implémentations des avantages de l'utilisation de LOV dans la création et la gestion de l'ontologie (cas de la méthodologie de NeOn), de l'harmonisation des préfixes et l'alignement des vocabulaires publiés sur le Web, ou sur du ranking des vocabulaires en utilisant la théorie du contenu de l'information. En appliquant ce dernier aux vocabulaires, nous avons essayé d'utiliser les fonctionnalités que nous jugeons "pertinents" à prendre en compte lorsque l'on veut des vocabulaires (par exemple : les jeux de données réutilisés, les liens vers les vocabulaires externes). Nous comparons notre approche avec d'autres classements qui sont principalement basées sur la "popularité" des vocabulaires. Ce travail peut ouvrir la voie vers une évaluation des vocabulaires avec des applications dans une approche plus systémique de recommandation des classes ou propriétés dans la gestion de l'ontologie, ou dans des applications de visualisation afin de proposer la propriété la plus appropriée à visualiser dans les ressource RDF contenant un grand nombre de propriétés.

## Chapitre VI

Dans ce chapitre, nous avons présenté un outil en ligne pour vérifier la compatibilité entre des jeux de données et des vocabulaires basés sur la logique "RDF-defeasible" de SPINdle. Nous avons implémenté le framework LIVE pour tester la compatibilité des licences sur les données publiés sur le Web. Le but de ce framework est de vérifier la compatibilité des licences associées aux vocabulaires utilisées pour générer un jeu de données RDF et la licence associée au jeu de données final. Plusieurs aspects d'ordre plus juridique doivent être pris en compte pour les travaux futurs. Plus précisément, nous considérons que les vocabulaires comme des données à part entière, mais ce n'est pas la seule interprétation possible. Par exemple, nous

pouvons voir des vocabulaires comme une sorte de compilateur, de telle sorte qu'après la création du jeu de données, les vocabulaires externes ne sont plus utilisées. Dans ce cas, quel serait le moyen approprié pour définir un système de vérification de compatibilité? Comme travail futur, nous étudierons en profondeur cette question ainsi que nous ferons une évaluation sur la facilité d'utilisation de l'outil en ligne LIVE pour améliorer l'interface utilisateur.

# Conclusion et Perspectives

Cette thèse est consacrée aux défis de la publication des données géospatiales sur le Web et une approche plus générique de visualiser les données liées pour les utilisateurs. La première considère la diversité des différents formats utilisés pour publier les données géospatiales propriétaires, les différentes projections (ou systèmes de coordonnées de référence) et la représentation des géométries complexes. Cette dernière approche est différente de l'état-de-l'art dans les visualisations où la complexité du langage SPARQL et RDF est pas suffisamment cachée des utilisateurs. Une analyse approfondie de la littérature a révélé certaines limites dans la publication des données géospatiales et des outils de visualisation, à savoir :

- Une présence limitée des géométries complexes représentées de manière structurée, au lieu de littéraux.
- L'absence d'une référence explicite aux SCG dans les données au géo-référencement direct sur le Web.
- Absence d'outil de visualisation destiné aux utilisateurs permettant de comprendre facilement l'essence des données sous-jacentes publiées en LOD.
- Beaucoup de silos de données pour les applications publiées sur le Web, perdues dans de nombreuses pages HTML.
- Peu d'outils qui fournissent un environnement intégré pour la publication des données brutes en données liées, partant la modélisation de données jusqu'à l'étape finale de stockage du jeu de données dans un store RDF.
- La difficulté pour les éditeurs de données de comprendre et de vérifier la compatibilité des licences entre les vocabulaires et les jeux de données qu'ils réutilisent venant du LOD.

Dans cette thèse, nous avons fourni des vocabulaires qui aident à la modélisation et la publication des données géospatiales intégrant la quasi-totalité des SCG, qui étendent les vocabulaires existants. Les vocabulaires ont été utilisés pour publier les unités administratives françaises, avec les données compatibles au standard GeoSPARQL. En ce qui concerne les visualisations, après avoir examiné des outils visuels et les applications existantes sur le Web, nous avons développé une ontologie pour mieux exposer les données sur le Web pour une meilleure interopérabilité. Nous avons également pro-

posé un framework pour générer automatiquement des visualisations basées sur les catégories détectés sur des jeux de données liées et publiées, en utilisant les catégories prédéfinies de haut niveau utilisées dans la taxonomie de la visualisation de l'information ; celle-ci mappée avec les vocabulaires.

## Revue des contributions

Cette section examine les principales contributions de cette thèse et les solutions que nous avons apporté comme contributions dans le contexte de la publication des données géospatiales sur le Web. Nos contributions décrites tout au long de cette thèse sont les suivantes :

- Nous avons modélisé et implémenté un vocabulaire pour la géométrie, les entités topologiques et les systèmes de coordonnées géographiques.
- Nous avons mis en place une API pour convertir des données en ligne entre les différents systèmes de coordonnées accessibles sur le Web.
- Nous avons publié les différents systèmes de projections utilisés en France avec des URI uniques pour améliorer la recherche et l'intégration des géométries structurées sur le Web.
- Nous avons contribué à l'élaboration de la plate-forme Datalift, un environnement intégré de publication des données brutes de formats hétérogènes sur le Web.
- Nous avons fourni une comparaison des triples stores pour les géo-données en projectant les types de géométries nativement incorporées (littéral ou structurée) pour aider à la recommandation lors de la publication des données géospatiales.
- Nous avons publié les données sur les circonscriptions administratives françaises selon les bonnes pratiques du LOD accessible au <http://data.ign.fr> basées sur les vocabulaires que nous avons développé. En outre, nous avons fourni des alignements avec des jeux de données géospatiales pertinentes existantes, tel que Geonames, GADM, NUTS, INSEE, etc.
- Nous avons publié en RDF 15 millions d'adresses provenant d'Open Street Map France en utilisant le vocabulaire des adresses proposé par le W3C.
- Nous avons contribué à la création du nuage de données dans le contexte français (FrLOD), en utilisant la plateforme Datalift, ainsi que des alignements avec des jeux de données existantes. Ces données ont la principale caractéristique de couvrir la France.
- Nous avons revu la littérature et avons classifié les applications construites sur des portails des données ouvertes des gouvernements, et avons proposé un vocabulaire pour annoter sémantiquement et améliorer la recherche et l'extraction d'applications créées dans le cadre des hackathon sur les données en Open Data.
- Nous avons proposé une approche générique pour générer automatique-

ment des visualisations basées sur des catégories prédéfinies à l'aide des requêtes SPARQL.

- Nous avons implémenté et évalué une approche pour déterminer les propriétés qui conviennent le mieux à utiliser pour choisir une entité à visualiser, basé sur une approche similaire à celle mise en oeuvre dans panel de recherche de la base de connaissance Google.
- Nous avons développé deux applications innovantes consommant des données événementiels et statistiques mutualisées avec des données externes présentes dans le nuage des données liées.
- Nous avons proposé un plugin générique pour améliorer la découverte des applications construites dans les hackathons sur les données en Open Data.
- Nous avons également proposé une approche pour harmoniser les préfixes utilisés dans les différents catalogues de vocabulaire, avec une évaluation faite dans le cas des vocabulaires sur Linked Open Vocabulary (LOV).
- Nous avons développé de nouvelles mesures de ranking pour les vocabulaires basées sur le contenu des informations et appliquée dans LOV.
- Enfin, nous avons construit un outil plus efficace pour vérifier la compatibilité des licences entre les vocabulaires et les jeux de données.

## Perspectives

Dans cette thèse, nous avons abordé certains problèmes ouverts de recherche dans le cadre de la publication et la consommation (visualisation) des données ouvertes sur le Web, mais il reste encore des questions en suspens et des défis pour des travaux futurs. Nous mentionnons quelques-uns des plus importants dans la section suivante, basés sur différents aspects liés à la chaîne d'édition des données liées, plus spécifiquement dans le domaine géospatial.

## Opportunités et Défis pour IGN-France

Le besoin de données de référence géographiques interopérables pour partager et combiner des informations environnementales spatiales géoréférencées est mis en évidence par la directive européenne INSPIRE. La directive INSPIRE vise à créer dans l'espace de l'Union Européenne (UE) une infrastructure de données spatiales<sup>8</sup>. INSPIRE est basée sur un certain nombre de principes communs de haut niveau, avec certains d'entre eux très propres des principes clés appliqués dans les fondements du Web sémantique, et en particulier dans son implantations dans les données ouvertes et liées. Nous

---

8. <http://inspire.jrc.ec.europa.eu/index.cfm/pageid/48>

fournissons ci-dessous la correspondance de nos contributions ayant un lien avec les cinq objectifs de la directive INSPIRE :<sup>9</sup>

- **P1** : *Les données doivent être collectées une seule fois et conservées où elles peuvent être maintenues le plus efficacement possible.* L'utilisation de bonnes politiques et des URIs stables peuvent aider à atteindre ce principe. IGN comme un fournisseur de données géospatiales en France est commis à des informations exactes, ainsi que seront donc les URIs choisis et utilisés pour le portail sémantique.
- **P2** : *Il devrait être possible de combiner des informations spatiales transparente provenant de différentes sources à travers l'Europe et les partager avec de nombreux utilisateurs et applications.* Ce principe est plus ou moins l'objectif des tâches d'interconnexion avec d'autres jeux de données dans le Web. Il faudrait restreindre les domaines de recherche dans les données européennes. Les modèles développés et bien documentés peuvent faciliter la conversion des données par d'autres organisations ou institutions utilisant ou produisant les données cartographiques.
- **P3** : *Il devrait être possible pour l'information recueillie à un niveau d'échelle être partagée à tous les autres niveaux ou échelles ; détaillée pour des recherches approfondies et générale à des fins stratégiques.* Un des inconvénients des modèles proposés est qu'ils n'admettent pas actuellement de nombreuses géométries attachées à une entité géographique ou adresse. Ce sera certainement l'une des extensions prévue pour les vocabulaires développés. Cependant, une classification précise des entités géographiques et topologiques est un début pour remplir ce principe.
- **P4** : *L'information géographique nécessaire pour la bonne gouvernance à tous les niveaux devrait être facilement disponible et transparente.* La publication du portail <http://data.ign.fr> est l'un des objectifs afin d'avoir également des données à la fois processable par la machine et lisible par l'homme avec l'aide des concepts et des technologies sémantiques.
- **P5** : *Accès facile pour retrouver quelle l'information géographique est disponible, comment elle peut être utilisé pour répondre à un besoin particulier, et dans quelles conditions elle peut être acquise et utilisée.* La publication des données sur le Web contribue en soit à tirer profit de leur découverte et intégration . En outre, une licence explicite attachée à aux jeux de données publiées contribue à atteindre de ce principe.

Pour les fournisseurs de données géographiques, les avantages de publier leur donnée sur le Web selon les principes du Linked Data sont de deux sortes :

1. Tout d'abord, leurs données sont interopérables avec d'autres jeux de
9. La traduction des objectifs est faite par nos propres soins.

données publiées et peuvent être référencées par des ressources externes et utilisées comme des données à référence spatiale, ce qui n'auraient pas été le cas si elles étaient publiées selon les normes des infrastructures de données spatiales (SDI).

2. Deuxièmement, l'utilisation des technologies du Web sémantique peut aider à résoudre les problèmes d'interopérabilité qui ne sont pas encore résolus par les normes et standards actuels dans le domaine de l'information géographique.

En outre, l'agence nationale de cartographie française (IGN) dispose de différents types de politiques de licence pour accéder aux données à partir de leur portail professionnel<sup>10</sup> (par exemple pour des fins de recherche, l'utilisation commerciale, l'accès à la demande, etc.), avec certains accès pas nécessairement "ouvert" ou libre d'accès : (par exemple, BD TOPO®). Bien qu'il y ait une compréhension claire des avantages de la publication et de l'alignement des données sur le Web, les recherches à l'IGN sont en cours sur la manière de combiner les licences sur des jeux de données. Deux solutions sont à l'étude :

1. Les différentes politiques de licence attachées aux jeux de données : Ici, la licence attachée est donnée directement lors de la publication. Ainsi, par exemple, s'il s'agit d'une licence libre, le point d'accès SPARQL est publiquement disponible et peut être interrogé sans aucune restriction.
2. L'utilisation d'un mécanisme d'accès aux données donnant accès selon une liste de configuration prédéterminée de graphes dédiés, des ressources et des opérations autorisées. Cette solution va en droite ligne avec les propositions des chercheurs comme Rotolo et al. dans l'application de la logique déontique. Cette solution suggère que même s'il y a un point d'accès, un module de configuration des types de requêtes à réaliser et des politiques d'accès doivent être définis pour les sous-ensembles de données avec un soin particulier pour tenir compte des compositions de licences dans les résultats.

Selon les principes du Linked Data, les URIs devraient rester stables, même si les unités administratives changent ou disparaissent. Cela implique l'adaptation du vocabulaire de données afin de gérer les versions des données, l'évolution temporelle et la granularité des données. Cette question sera abordée dans nos travaux futurs, comme nous travaillons sur la publication du jeu de données spatio-temporelle décrivant l'évolution des communes depuis la Révolution française. Une autre question de recherche porte sur l'automatisation de l'ensemble du processus de publication, partant des données géographiques aux formats traditionnels (SHAPE, CSV, etc) pour arriver à des données RDF pleinement interconnectées .

---

10. <http://professionnels.ign.fr/>

La dernière question porte sur l'utilisation de plusieurs géométries pour décrire une entité géographique : des géométries avec différents niveaux de détail, avec différents CRS ainsi que des choix différents de représentation. Cela a été superficiellement abordées dans notre cas d'utilisation avec l'utilisation de deux polygones et de points pour représenter respectivement la surface et le centre de gravité des communes, mais doit être étudiée en profondeur pour proposer une solution intégrant à la fois les contraintes de requêtes et d'affichage d'informations sur des fonds de carte en fonction des besoins utilisateurs.

## Visualisations Génériques des données liées sur le Web

Nous prévoyons d'utiliser un ensemble plus exhaustif des vocabulaires dans nos requêtes génériques pour détecter les catégories, en prenant la liste des vocabulaires du catalogue LOV pour alimenter l'assistant. Les propriétés d'agrégation peuvent être étendues afin de prendre d'autres relations sémantiques (par exemple, prendre en compte SKOS:`exactMatch`). En outre, nous prévoyons de faire une évaluation du prototype et le comparer à des outils connexes comme celles permettant de produire des statistiques des jeux de données. Nous avons également besoin de quantifier quand une catégorie est «importante» dans un jeu de données. Par exemple, est-ce suffisant pour un jeu de données pour être classé dans la catégorie “GEODATA” avec juste dix triplets contenant des adresses ? À partir de quel nombre de triplets et donc quelle proportion pourrait-on utiliser les catégories, donc les librairies de visualisations associées ? Ces questions peuvent en outre être étudiées pour trouver le meilleur compromis entre le pourcentage de représentativité de certaines catégories et les librairies correspondantes. Un inconvénient de notre travail sur les visualisations est l'absence d'évaluation au niveau de l'utilisateur final, avec un protocole bien défini pour comprendre les besoins des utilisateurs, en se concentrant davantage sur les aspects sémantiques que sur ceux liés juste à l'exploration (interface web). Un travail futur naturel est d'utiliser ces évaluations et ré-adapter les applications/ visualisations basées sur les résultats d'une étude utilisateur.

## Vocabulaires et LOV

Les travaux sur l'harmonisation des préfixes peut être étendu dans plusieurs directions. En se limitant aux deux services que nous avons étudié et déjà contribué à harmoniser les prefixes, les prochaines étapes possibles seraient d'automatiser autant que possible les tâches qui ont été faites de manière semi-automatique à ce jour :

- *i)* le développement d'une interface unique pour soumettre les espaces de noms et les préfixes aux deux services ;
- *ii)* la couverture simultanée des prefixes dans LOV et prefix.cc pour

harmoniser les URIs des vocabulaires présents dans les deux services afin de ne confondre les utilisateurs et les développeurs dans le choix ou la gestion des espaces de noms, et ainsi proposer une liste recommandée namespaces-URIs des vocabulaires les plus importants.

Ce dernier aspect va au-delà du domaine d'application des deux services car une telle liste pourrait être consolidé et approuvé par les principaux acteurs de publication et de gestion du vocabulaire, et recommandé pour une utilisation dans les applications de données liées. Cela pourrait être pris en charge par le prochain groupe de travail du W3C<sup>11</sup> chargé de gérer les vocabulaires dans le cadre de la nouvelle activité de gestion des données sur le Web.

Pour le ranking des vocabulaires, nous souhaitons prendre en compte les axiomes d'équivalence (entre les classes et les propriétés) lors du calcul du contenu de l'information, et plus généralement, toutes sortes de relations sémantiques entre les termes. En outre, nous prévoyons de comparer notre modèle de classement avec d'autres approches telles que celles de classement fondées sur les graphes (par exemple, le PageRank). Une autre orientation future est de chercher la dépendance dans la position entre les vocabulaires, en se concentrant sur un type spécifique de "liens entrants" (à savoir les extensions, les généralisations) et d'étudier comment ils affectent les métriques que nous avons présenté dans cette thèse.

Nous avons fait l'hypothèse dans cette thèse que l'accès aux données se faisait soit en interrogeant un accès SPARQL, en parcourant le graphe par le principe du "follow-your-nose" ou en téléchargeant des dumps. Récemment, une nouvelle façon d'accéder aux données sur le Web est en train d'émerger : à travers des motifs de fragments triplets liés<sup>12</sup>. Ce concept vise à explorer les accès de données avec des fragments simples de données pour résoudre les requêtes côté client avec les données hébergées dans un serveur. Les serveurs peuvent servir des données à faible coût de traitement d'une manière qui permet l'interrogation côté client en déplaçant du même coup l'intelligence passant du serveur vers le client. Un travail futur pourrait être d'utiliser le concept côté client pour évaluer les points d'accès aux données contenant seulement les géométries structurées ou littéraux pour les applications du monde réel. Enfin, le concept de fragments de triplets peut également être appliqué pour détecter des patterns pour la visualisation des différents points d'accès de données.

Avec la croissance continue et soutenue de la publication des données ouvertes et liées sur le Web, il en sera aussi des jeux de données et des ontologies

---

11. <http://www.w3.org/2013/05/odbp-charter.html>

12. <http://linkeddatafragments.org/>

sur des données géospatiales. Les producteurs de données géographiques vont continuer de libérer de plus en plus fréquemment leurs données sur le Web. Cela va créer un besoin d'outils pour facilement créer des analyses, en particulier dans l'extraction et la fouille de gros volumes de données pour retro-alimenter les éditeurs quant à l'utilisation effective des triplets. La gestion des flux de données géospatiales sur le Web va demander des implémentations plus efficaces dans le domaine des données spatiales pour être en mesure d'interroger à la volée des données contenant aussi de l'information temporelle. Ainsi, la modélisation des flux de données en streaming géo-temporel, l'interrogation et l'analyse sur le Web sont susceptibles d'être les prochains défis que les technologies du Web sémantique devront faire face et résoudre.