

Publishing and Consuming Geo-spatial and Government Data on the Semantic Web

Ghislain Auguste Atemezing

Multimedia Department, Eurecom

Supervisor:
Dr. Raphaël Troncy, MM Department, Eurecom



FINANCIÉ PAR
ANR

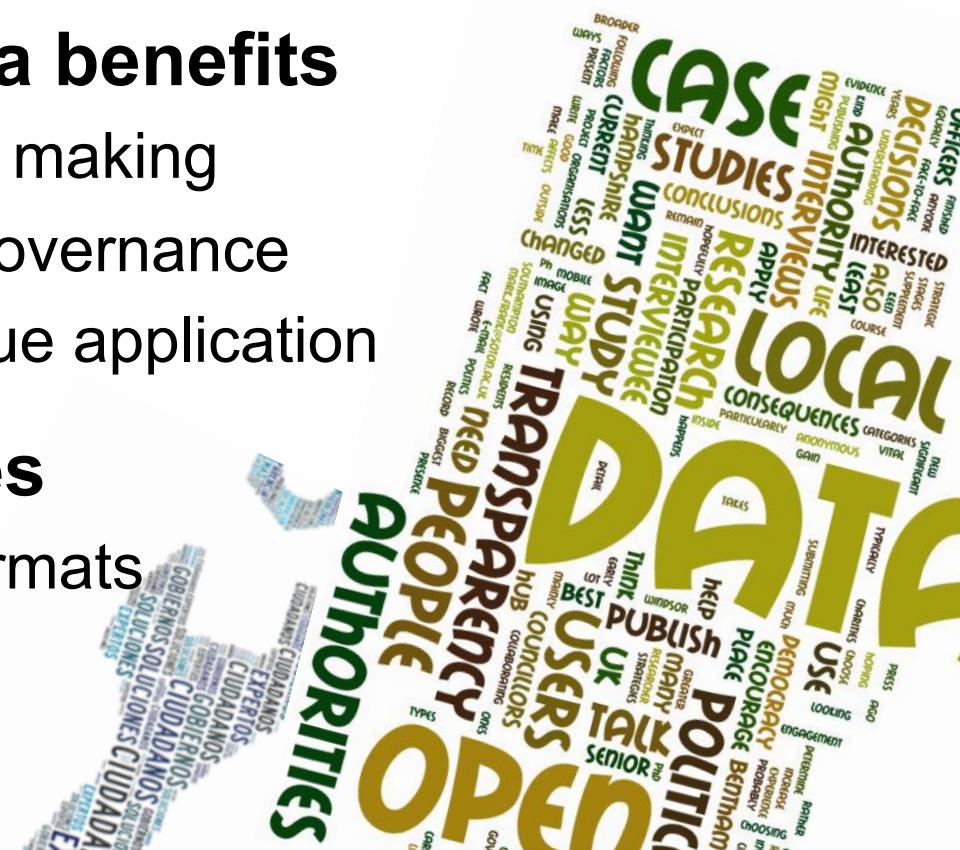
Thesis Context

▪ Open Government Data benefits

- Transparency in decision making
 - Better governance or e-governance
 - Eco-system of added value application

■ Barriers and challenges

- Heterogeneity of data formats
 - Variety of access method
 - Lack of nomenclature



In this thesis we explore how Semantic Web technologies can be used for better integration and consumption of geo-spatial data.

Outline

- **Introduction**
- **Research Questions**
- **Contributions**
 - Semantic publishing of geospatial data
 - Visualizations of Government Linked Data
 - Best Practices for metadata in vocabularies
- **Conclusion**
- **Publications**

Geospatial data: why it matters

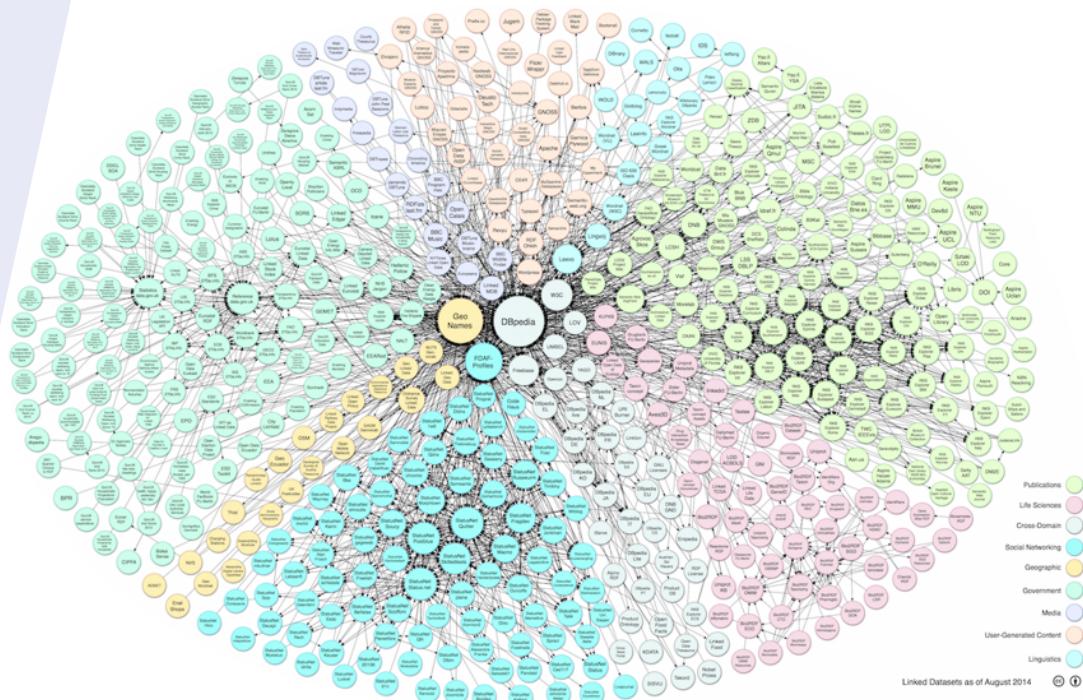
“80% of needs for decisions from public authorities have a geospatial component”.

(Philippe Grelot, IGN-France)

The collage includes:

- A green circular badge with a crown icon and the text "CHECK-IN HERE ON foursquare".
- A photograph of the Eiffel Tower with a small inset map showing its location.
- A Google Maps interface showing the Eiffel Tower with "Directions" and "Write a review" buttons.
- A text box containing search queries related to food in Nice: "where to eat in nice", "where to eat in niceville fl", "where to eat in nice france", and "where to eat in nice old town".
- A map of Paris with numbered arrondissements (1-20) and landmarks like Montmartre, Champs Elysées, and the Louvre.
- A TomTom GO satellite navigation device showing a route to New York City.
- A screenshot of a TripAdvisor search results page for "restaurants in Nice - French Riviera".
- A Wikipedia page for the "7th arrondissement of Paris", featuring a map and text about its history and landmarks.
- A screenshot of a mobile phone displaying a map of Paris with numerous red location pins.

GeoData on the LOD Cloud



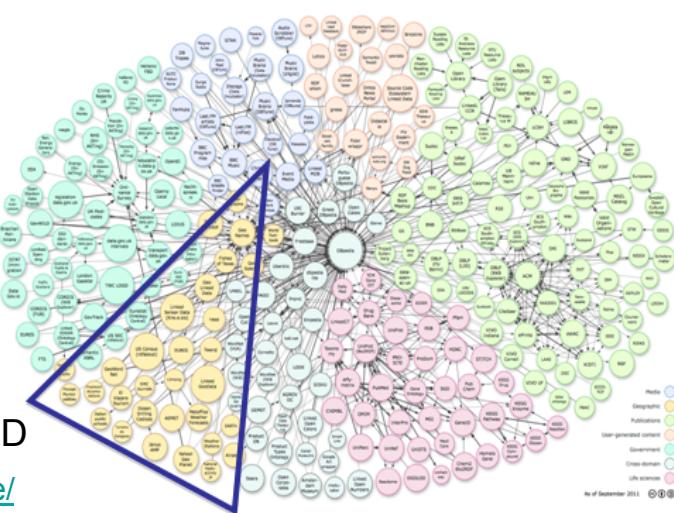
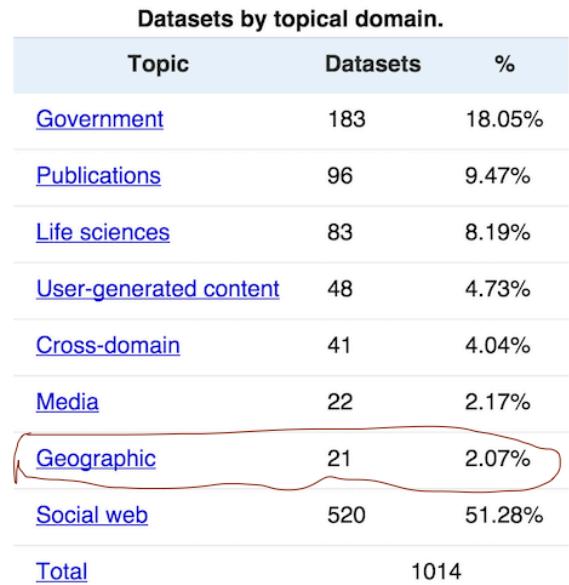
<http://lod-cloud.net/versions/2014-08-30/>

Linking Open Data cloud diagram 2014, by Max Schmachtenberg, Christian Bizer, Anja Jentzsch and Richard Cyganiak.

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

In 2011 19,43% → 31 geo-datasets in LOD

<http://lod-cloud.net/state/>



French IGN- Reference Geodata

« ..describes the French national territory and the occupation of its land, elaborates and updates perpetual inventory of the forest resources »

- ✓ Different databases with **overlapping** information:
BD ORTHO, BD PARCELLAIRE, POINT ADRESSE, BD ALTI 25m,
BD TOPO; etc.
- ✓ CRS: LAMBERT93 or RGF93

Q: "Give me all the bridges in a radius of 2km from the "Eiffel Tower"?"
A: Not straightforward

CATALOGUE

■ Données par type >>
RGE®

Photos aériennes
Cartographie
Topographie et foncier
3D
Adresse
Altimétrie
Administratif
Routier
Historique
Aéronautique
Données par échelle >>
Données raster >>
Données vecteur >>
Données DOM-TOM >>
Données internationales >>
Données gratuites >>
Données INSPIRE >>

IMPRIMER **ENVOYER** **PARTAGER**

Données par type

Bases de données du référentiel à grande échelle (RGE®)

L'État a confié à l'IGN le développement du référentiel à grande échelle (RGE®). Pour ce faire, il fait appel à ses moyens propres ainsi qu'à des partenariats avec des producteurs principalement de la sphère publique. Le RGE® est constitué de 4 composantes : orthophotographique, topographique et altimétrique, parcellaire, adresse.

BD ORTHO >> La composante orthophotographique du RGE®

BD ADRESSE® >> La composante adresse du RGE®

BD TOPO® >> La composante topographique du RGE®

RGE®ALTI >> La composante altimétrique du RGE®

BD PARCELLAIRE® >> La composante parcellaire du RGE®


HAUT DE PAGE

Photographies aériennes et orthophotographies

L'IGN propose une gamme complète d'orthophotographies (BD ORTHO®, ORTHO HR et BD ORTHO® Historique) ainsi que des prises de vues aériennes.

BD ORTHO® >> L'orthophotographie IGN de référence

ORTHO HR >> L'orthophotographie IGN à Haute Résolution

BD ORTHO Historique® >> L'orthophotographie historique de l'IGN

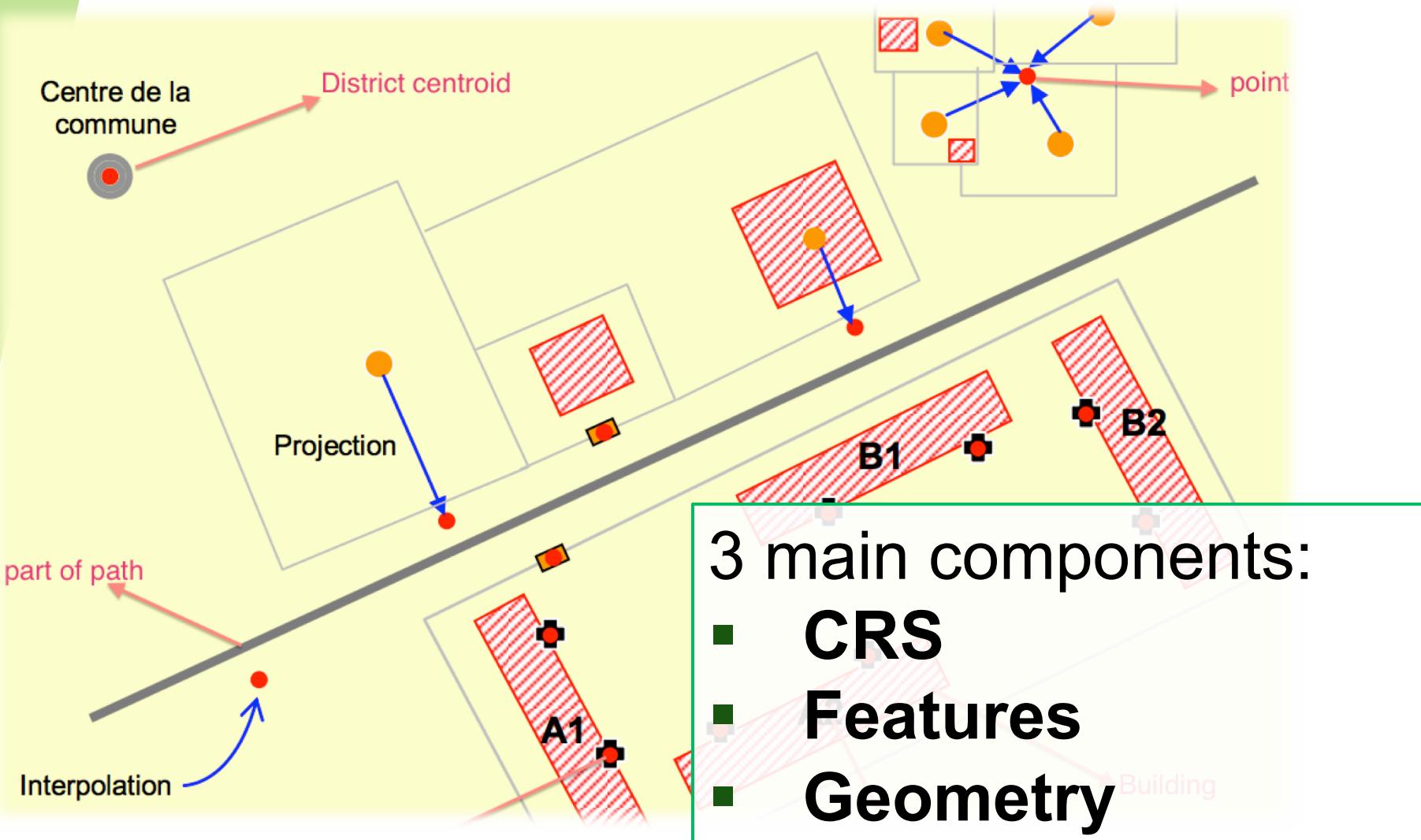

Research Questions

- **How to efficiently represent and store geospatial data on the Web to ensure interoperable applications?**
 - the publication of real datasets
 - the interlinking with other datasets having overlapping coverage
- **What are the best options for a user to discover, browse and interact with semantic content?**
 - Can we propose a generic model for visualizing semantic content?
- **What are the mechanisms to help preserving structured data of a high quality on the Web?**
 - How to improve reusing vocabularies for better interoperability
 - How to detect incompatibility between data and metadata

Part 1

Semantic Publication of Geo-spatial Data

Modeling geospatial objects



Coordinate Reference Systems (CRS)

- **A representation of the locations of geographic features within a common geographic framework**
- **Each CRS is defined by its measurement framework**
 - Geographic: spherical coordinates (unit: decimal degrees)
 - Planimetric: 2D planar surface (unit: meters) + a map projection
 - Additional measurement properties: ellipsoid, datum, standard parallels, central meridians, etc.
- **Situation:**
 - Several hundred Geographical Coordinate Systems (WGS84, ETRS89, etc.)
 - Several thousands Projected Coordinate Systems (UTM, Lambert93, etc.),
<http://resources.esri.com/help/9.3/arcgisserver/apis/rest/pcs.html>

CRS in France (Metropolitan + overseas)

- 10 CRSs coverage (mostly RGF93)
- 10 projections (mostly used Lambert93)
- 3 different ellipsoids to define the CRS

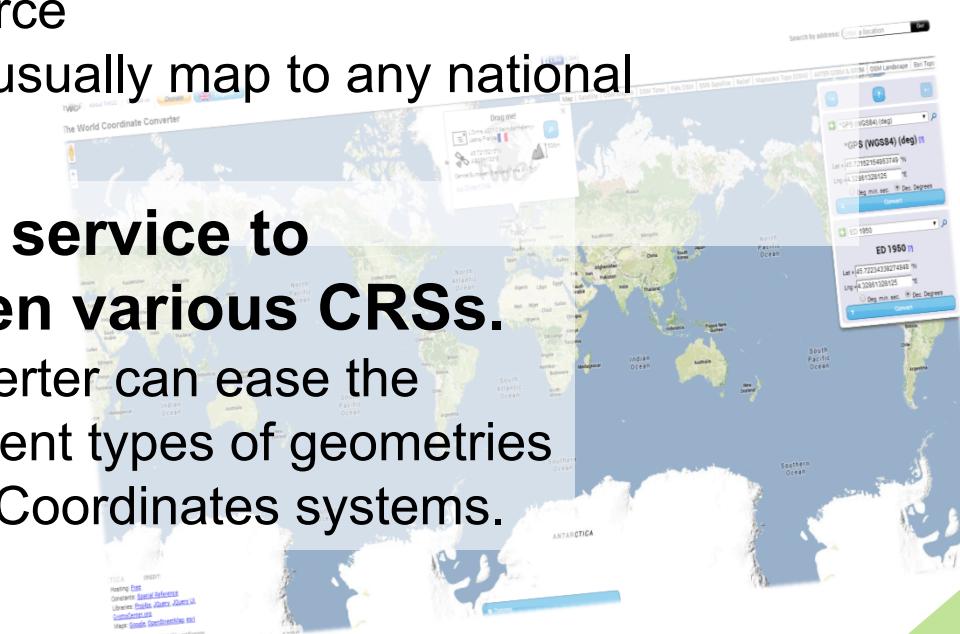
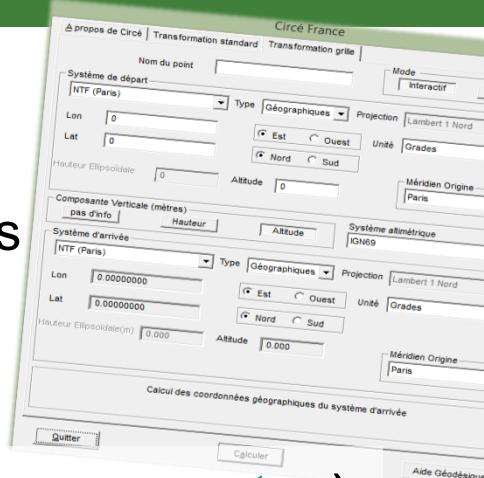
The Web only used
WGS84

Publishers need converters
for their geodata

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

CRS Converters and Limitations

- **Circé: A Converter Software from IGN**
 - Allows conversion between some France CRSs to WGS84
 - Closed Source, No open service
- **Existing Web tools (e.g., [world coordinate converter](#))**
 - No open service, closed source
 - Converting algorithms don't usually map to any national mapping authority.
- **Contribution: a Web based service to perform conversion between various CRSs.**
 - The integration of such a converter can ease the process of the publishing different types of geometries on the web regardless of their Coordinates systems.



Developing a REST converter service

- Many regional's published data involve with local CRS

International Communication

A medium

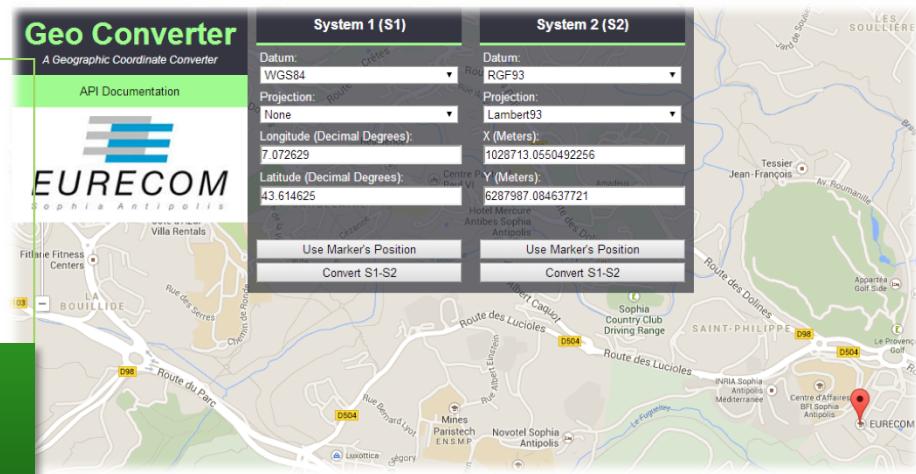
- Interpret the coordinates between these CRSs

- Open services for community
- RESTful Web Service

A Converter Service

- WGS 84 \leftrightarrow Lambert 93
- WGS 84 \leftrightarrow UTM

GOOD accuracy of the results



Vocabularies for Modeling Features

- **Authority list of terms (e.g. Foursquare)**
 - No semantics at all!
- **SKOS Categories (e.g. GeoNames)**
 - Classes are `skos:conceptScheme`, codes are `skos:Concept`
- **Domain specific ontologies (OrdS, GeoLD)**
 - Interconnected subdomain ontologies (transport, hydrography, etc.)
- **Data driven ontologies (LGD, GeOnto)**
 - Deeper taxonomy to structure the ontology
 - Many classes

Modeling Geometry: State of the Art

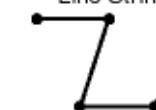
- **Point (lat/long)**

- WGS 84 vocabulary described by W3C

Point



Line String



Polygon



- **Rectangle (“bounding box”)**

- Geopolitical Vocabulary (FAO)

- **Points in a List**

- Sequence of points (LinkedGeoData)
- An object is “*formedBy*” a ListOfPoints (GeoLinkedData.es)

Arc Line String



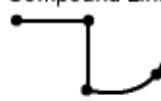
Arc Polygon



Compound Polygon



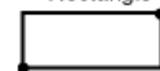
Compound Line String



Circle



Rectangle



- **Literals WKT datatype in RDF**

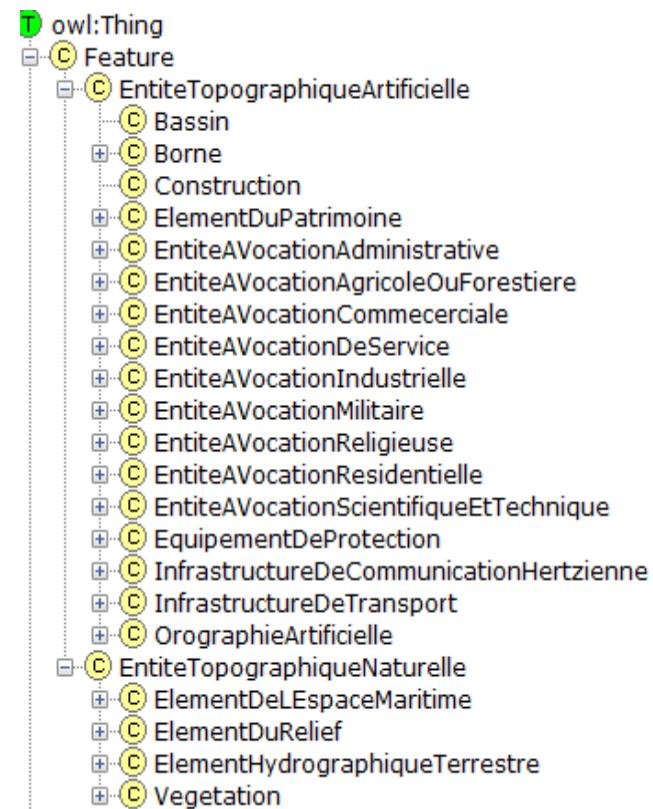
- Ordnance Survey (UK), GeoSPARQL embedding CRS in literals

- **More structured representation of complex geometry**

- NeoGeo Vocabulary (GeoVocamp), <http://geovocab.org/>

Reusing Existing Ontologies (GeOnto)

- **Ontology for geographic objects (POI)**
 - Output of a French (ANR) research project
 - Obtained from NLP tools
- **Classes in French**
 - rdfs:labels in FR & EN
 - No rdfs:comments
 - Few owl:ObjectProperty
 - **783 classes**
- **Overlap with other vocab**
 - Need for alignment

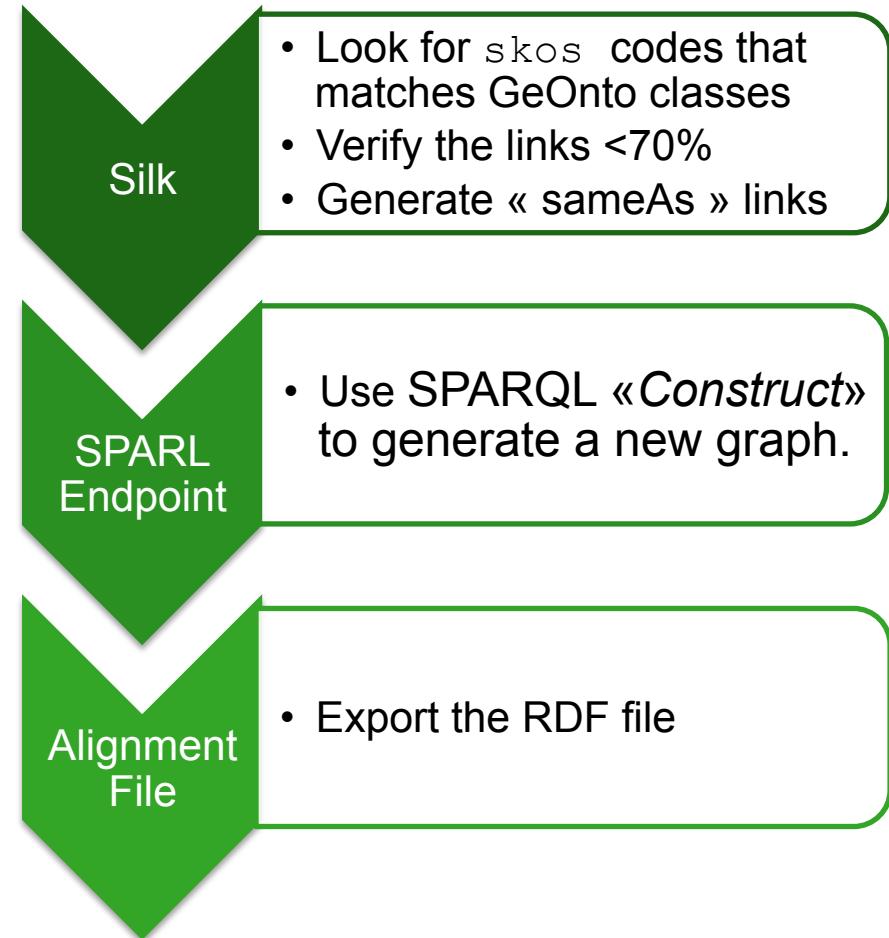


Aligning GeOnto with existing Ontologies

- **Alignment of GeOnto with 5 ontologies and 2 simple taxonomies**
 - LGD, DBpedia, Schema.org, GeoNames, bdtopo
 - Foursquare, Google Places
- **Goal: finding owl:equivalentClass**
 - Tool : Silk framework
 - Metrics : LevenshteinDistance, Jaro
 - Labels : @en des classes
 - Aggregation Function: Mean
- **Manual validation**
 - For « rdfs:subClassOf »
 - Specific alignments with GeoNames codes

Alignment Process (GeoNames)– Results

Vocab/ taxonomies	#Classes	#Classes aligned
Bdtopo	237	153 (64.65%)
GeoNames	699	287 (41.06%)
Google Place (*)	126	41 (32.54%)
Schema.org (*)	296	52 (17.57%)
LGD	1294	178 (13.76%)
Foursquare	359	46 (12.81%)
DBpedia	366	42 (11.48%)



- High precisions > 80%
- BUT $P(\text{Schema.org}) = 50\%$ → GeOnto entities are more specific to France

Proposal: CheckList for modeling geodata

- **Complex Geometry Coverage**
 - Need to publish more data with complex geometries
 - Reuse and extend suitable ontologies (NeoGeo, GeoSPARQL)
- **Features MUST be connected to Geometry**
 - Sometimes it may requires two namespaces
- **Serialization in other GIS formats**
 - Provide serialization in other GIS formats (GML, WKT, KML, etc.)
- **Structured Representation**
 - Use of structured representation for complex geometry
 - This covers some of the Use Cases at IGN

Proposal: Vocabulary for complex Geometry

- Extend and reuse existing vocabularies
- Available at <http://data.ignf.fr/def/geometrie>

```
@prefix ngeo: <http://geovocab.org/geometry#>.  
@prefix sf: <http://www.opengis.net/ont/sf#>.  
[...]
```

```
geom:Geometry a owl:Class;  
    rdfs:comment "Primitive géométrique non instanciable, racine de  
        l'ontologie des primitives géométriques. Une géométrie est  
        associée à un système de coordonnées et un seul."@fr;  
    rdfs:label "Géométrie"@fr, "Geometry"@en;  
    owl:equivalentClass [ a owl:Restriction;  
        owl:onClass ignf:CoordinatesSystem;  
        owl:onProperty geom:crs;  
        owl:qualifiedCardinality "1"^^xsd:nonNegativeInteger];  
    rdfs:subClassOf ngeo:Geometry;  
    rdfs:subClassOf sf:Geometry.
```

Ontologie des primitives géométriques

IRI: <http://data.ign.fr/def/geometrie>
Current version: Version 1.0 - 2014-08-22
Authors:
<http://recherche.ign.fr/labos/cogit/cv.php?prenom=Nathalie&nom=Abadie>
<http://www.eurecom.fr/~atemezini/>
Contributors:
<http://data.semanticweb.org/person/bernard-vatant>
<http://recherche.ign.fr/labos/cogit/cv.php?prenom=Bénédicte&nom=Bucher>
<http://www.eurecom.fr/~troncy/>
Publisher:
http://fr.dbpedia.org/resource/Institut_national_de_l%27information_g%C3%A9ographique
Other visualisation:
[Ontology source](#)

Copyright 2014, IGN

Associating CRS to Geometries

■ A vocabulary for describing CRS

- Subset of ISO 19111 model
- Available at <http://data.ignf.fr/def/ingf>

■ A dataset for French CRS

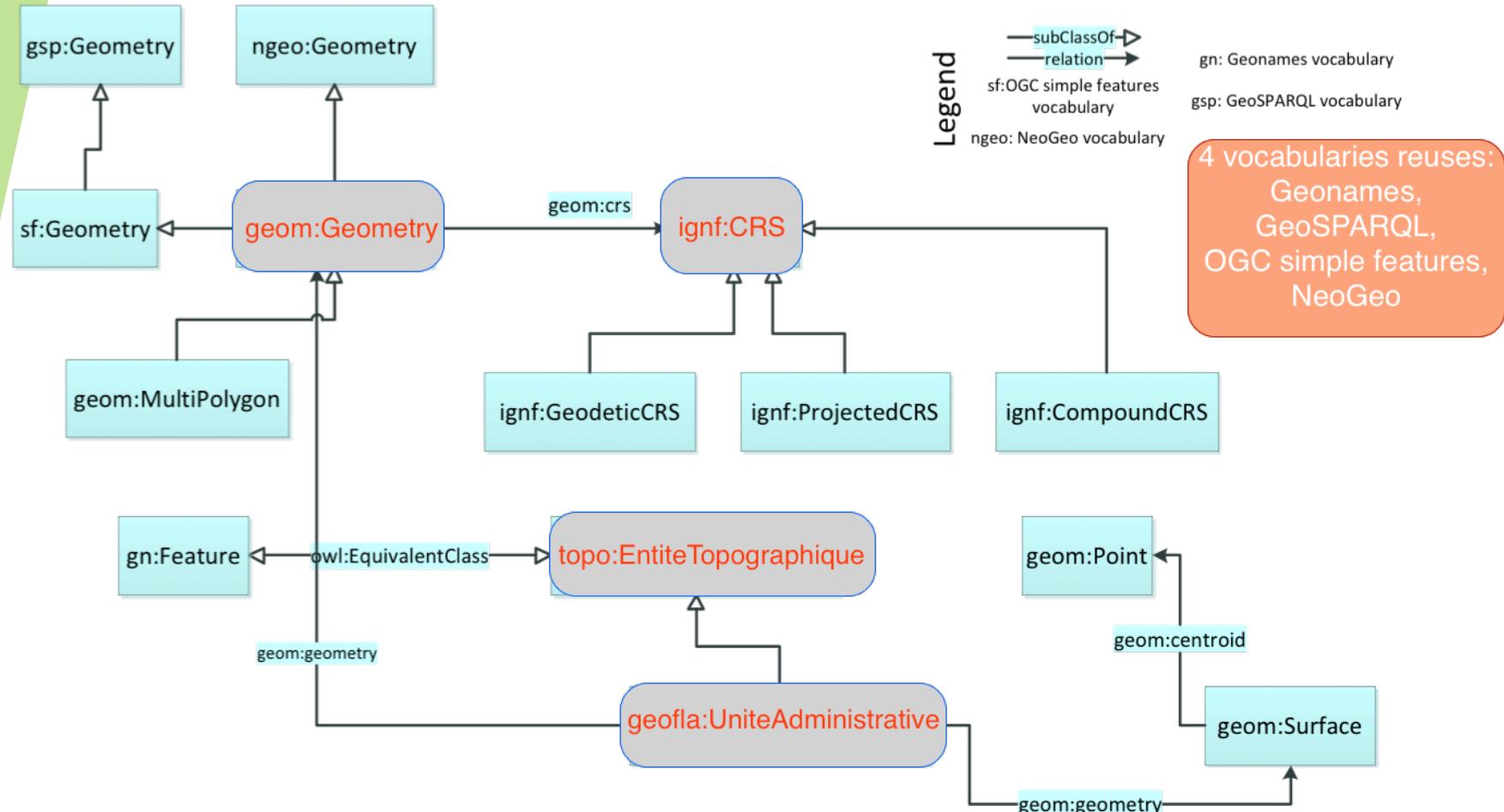
- Convert from XML data published by IGN France to RDF
- Eg: “*Lambert 2 étendu*” <http://data.ign.fr/id/ignf/crs/NTFLAMB2E>

CRS lookup

CRS:	ntfl
	http://data.ign.fr/id/ignf/crs/NTFLAMBGC
	http://data.ign.fr/id/ignf/crs/NTFLAMB1
	http://data.ign.fr/id/ignf/crs/NTFLAMB1C
	http://data.ign.fr/id/ignf/crs/NTFLAMB2
	http://data.ign.fr/id/ignf/crs/NTFLAMB2C
	http://data.ign.fr/id/ignf/crs/NTFLAMB2E
	http://data.ign.fr/id/ignf/crs/NTFLAMB3
	http://data.ign.fr/id/ignf/crs/NTFLAMB3C
	http://data.ign.fr/id/ignf/crs/NTFLAMB4

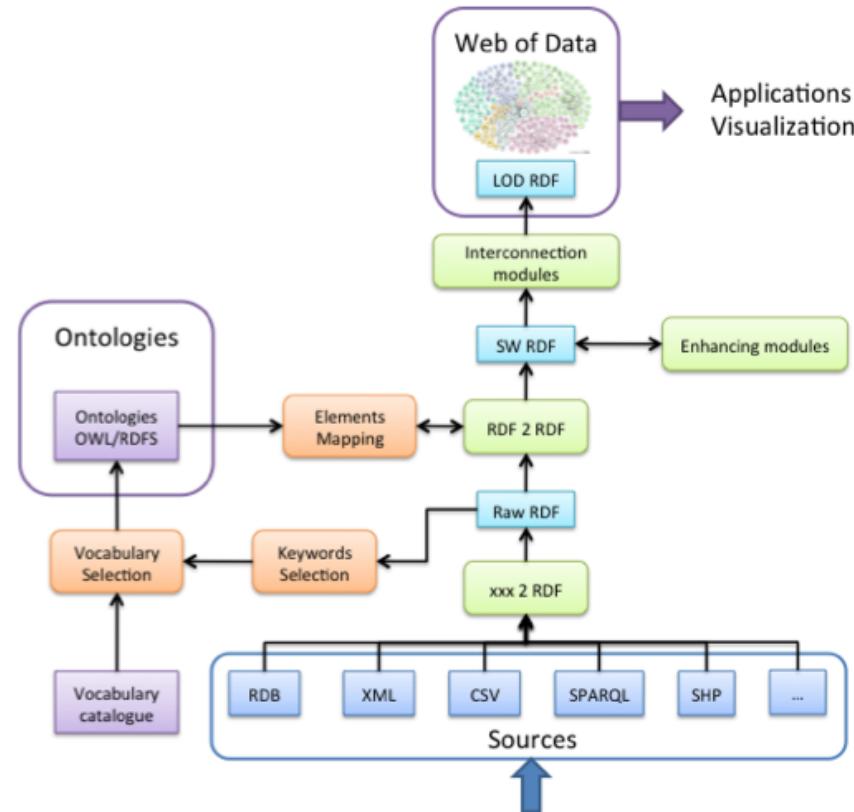
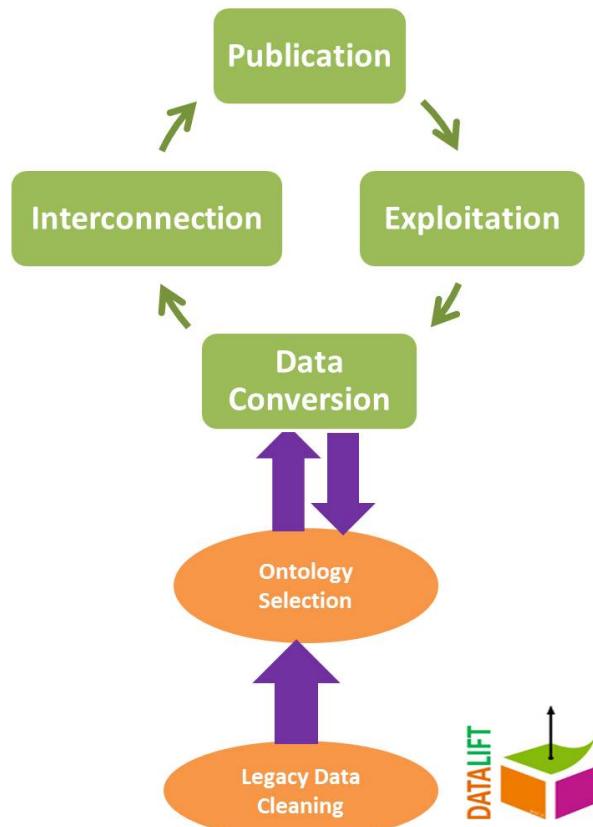
Description du nœud <http://data.ign.fr/id/ignf/crs/NTFLAMB2E>		
Rechercher	Texte	dans
	Toutes les colonnes	
Sujet	Prédicat	Objet
http://data.ign.fr/id/ignf/crs/NTFLAMB2E	http://www.w3.org/1999/02/22-rdf-syntax-ns#type	http://data.ign.fr/def/ingf#ProjectedCRS
	http://www.w3.org/2000/01/rdf-schema#label	"NTF Lambert II étendu"@fr
	http://data.ign.fr/def/ingf#scope	"NATIONALE, HISTORIQUE"@fr
	http://data.ign.fr/def/ingf#epsgID	"urn:ogc:def:crs:EPSG:6.11.2.2757"
	http://data.ign.fr/def/ingf#epsgID	"urn:ogc:def:crs:EPSG:3:20002120"
	http://data.ign.fr/def/ignf#domainOfValidity	http://data.ign.fr/id/ignf/extension/NTFLAMB2E
	http://data.ign.fr/def/ignf#cartesianCS	http://data.ign.fr/ignf/cartesiancs/TYP_CRG32
	http://data.ign.fr/def/ignf#baseCRS	http://data.ign.fr/ignf/crs/NTPCGRAD
	http://data.ign.fr/def/ignf#conversion	http://data.ign.fr/ignf/conversion/PRC012013
	http://data.ign.fr/def/ignf#includesSingleCRS	http://data.ign.fr/ignf/crs/NTFLAMB2E

Overview of the vocabularies and relationships



A workflow for publishing Linked(Geo)Data

- **DATALIFT: a set of modular tools for “lifting” raw data in RDF.**



Publishing French Reference Geodata in RDF

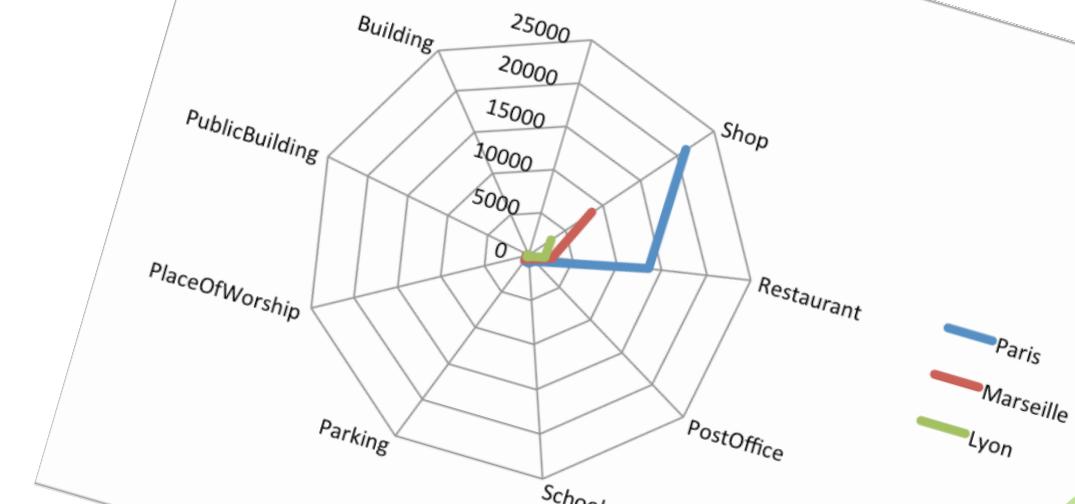
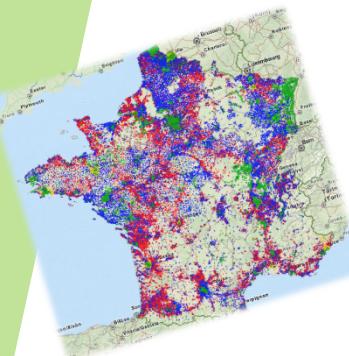
- **GEOFLA® DATASET ON FRENCH ADMINISTRATIVE UNITS IN RDF**
 - Features are of type <http://data.ign.fr/geofla>
- **FRENCH GAZETTEER DATASET**
 - Features are of type <http://data.ign.fr/topo>

Mapping results with external dataset			
DATASET GEOFLA DATASET		#mappings	Tool
	INSEE	37,020	SILK
	DBPEDIA-FR	23, 252 (communes) and 93 departments	LIMES
	NUTS	105 links (14 comm., 75 depts, 16 regions)	SILK
	GADM	70 links (10 comm., 51 depts, 9 regions)	SILK
FRENCH GAZETTEER	LINKED GEODATA	654 links (Igdo:Amenity)	LIMES

French Open Addresses in RDF Mappings

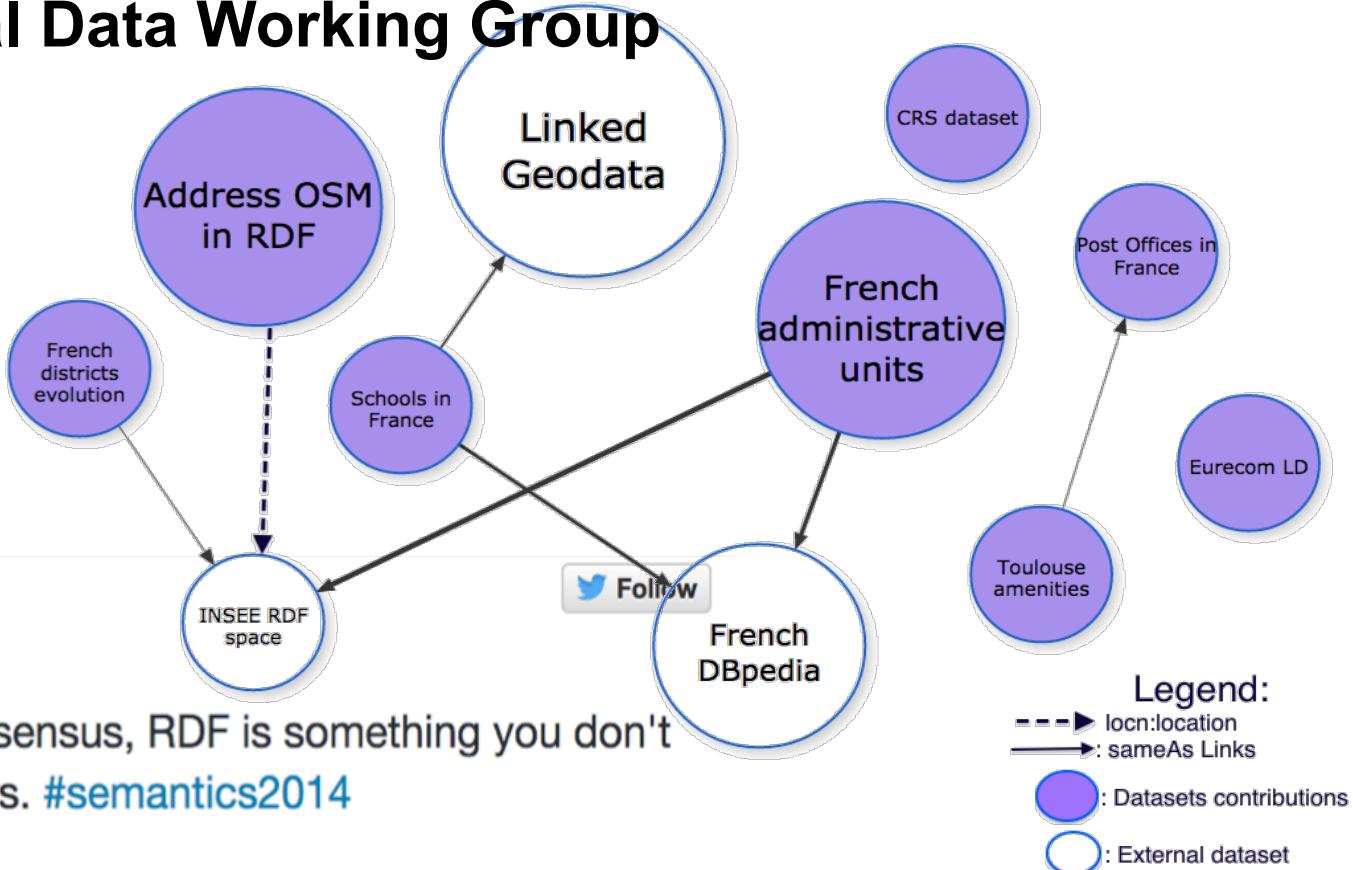
LGData Amenities	Paris (248052)		Marseille (401404)		Lyon (89061)	
	#matched	%matched	#matched	%matched	#matched	%matched
Shop (778680)	21171	2.71	8556	1.098	3049	0.391
Restaurant (260675)	13567	5.204	2654	1.018	1882	0.721
PostOffice (87731)	971	1.106	555	0.632	173	0.197
School (318287)	883	0.277	411	0.129	197	0.061
Parking (250516)	735	0.293	625	0.24	210	0.083
PlaceOfWorship (357445)	272	0.076	193	0.053	31	0.008
PublicBuilding (26735)	97	0.362	64	0.239	21	0.078
Building (22283)	5	0.022	12	0.053	0	0

BANO2RDF



Contribution to the LOD Cloud (FrLOD)

- 340 millions of triples in RDF (**10% of DBpedia 2014**)
- Part of our work is under reused by the W3C/OGC Spatial Data Working Group



Phil Archer
@philarcher1

I think we have consensus, RDF is something you don't show your end users. #semantics2014

1:42 PM - 5 Sep 2014

25 RETWEETS 14 FAVORITES



Why Visualizations Matters

*“Don’t ask what you can do for
the Semantic Web; ask what
The Semantic Web can do for you!”
(D. Karger, MIT CSAIL)*

*“If you use our Linked Data, please let
us know, or we might switch off!”
(Ordnance Survey)*

1. How to build **bridge** to fill the gap between traditional InfoVis tools and Semantic Web technologies
2. How can Semantic Web help in visualization?



Part 2

Generating Visualizations For Linked Data

Visualization Categories in Government Portals

- **Study of applications consuming Open Data**
 - 13 applications from UK (7), USA (3) and France (3)
 - Domain: education, health, transport, government, city, housing, criminality, foreign aid
- **Different dimensions**
 - Platform (web, mobile), data sources, which views are available (maps, charts, timeline, etc.)
 - URL policy for identifying data objects
 - Licenses for the application / for the data
 - Commercial / non-commercial

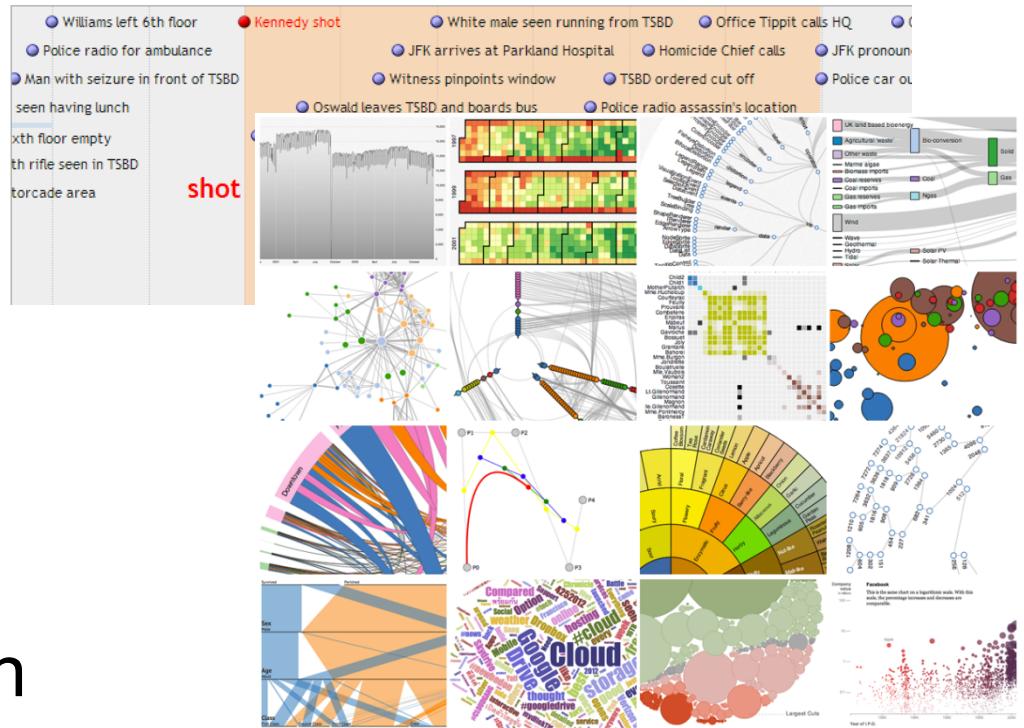
Relevant Features in Visualization Tools

- Data format given as input (csv, xml, shp, rdf, etc.)
- Data access (API, dump, etc.)
- Language code
- Type of view
- External Libraries
- License
- Metadata: author, organization

Timeline

Web Widget for Visualizing Temporal Data

With this widget, you can make beautiful interactive timelines like the one below. Try dragging it horizontally or using your mouse-wheel. Click on each event for more details.



Describing an Application: Opendatacom

Scope/Domain: Department for Communities and Local Government, datasets access

Description: visualize available datasets (finance, housing, deprivation, geography) by authorities or postcode.

On the dashboard, it provides graphs showing the national distribution of a district and how the values for this local authority compare with others in England.

Supported Platform: Web

URL Policy: <http://{domain}/id/{...}> with redirection to the corresponding document at: <http://{domain}/doc/{...}>.

Hampshire County Council is:

<http://opendatacommunities.org/id/county-council/hampshire>

Data Sources: 36 datasets from DCLG, Administrative Geography and Postcodes from Ordnance Survey.

Type of View: Graph, Map views.

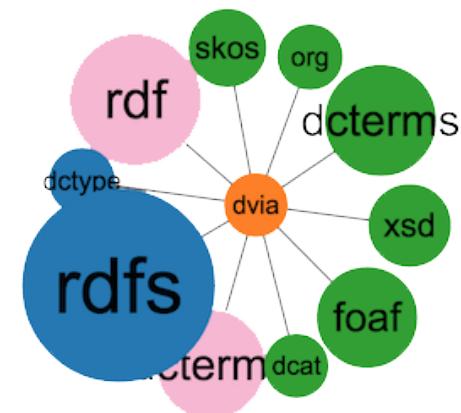
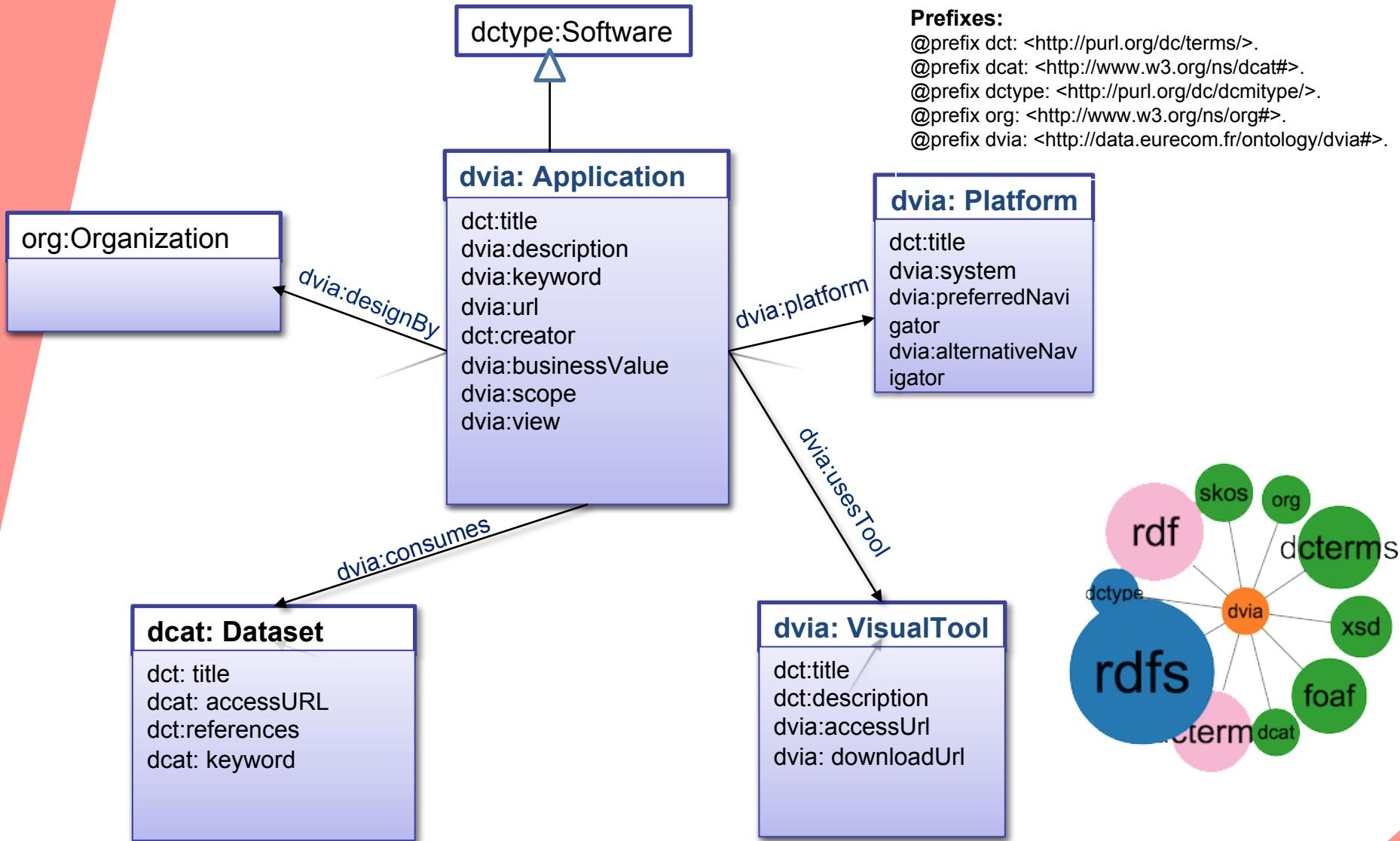
Visualization Tools: google visualization API, raphael.js

License: Open Government license [OGL]

Business Value: Non commercial



DVIA: A vocabulary to describe Applications



DVIA in Real World Datasets

- **4 applications re-using DVIA**
 - Use to populate 22 past events of hack-a-thons in Europe, with 889 applications from [Apps4Europe](#).
 - Implementation of a universal JavaScript plugin to embed RDFa in organizers events
 - An [extension for Wordpress](#) uses DVIA.



The screenshot shows a WordPress plugin demo at localhost. At the top, there's a navigation bar with links for Home, Plugin 3, Plugin2, and Universal Plugin Demo. Below the navigation, a user login form is displayed. To the right of the login form, a red dashed box highlights a "Test Event" section. This section contains a title "Test Event", a long descriptive text block, and a bulleted list under the heading "Bullets". To the far right, there's a sidebar with event details: Location (Somerset, UK), Starts (2014-06-15), Ends (2014-06-15), Organizers (Sample Organization), Sponsors (Sample Sponsor), Theme, and Territory.

localhost

Home Plugin 3 Plugin2 Universal Plugin Demo

Home

User login

Username * admin

Password * [Create new account](#) [Request new password](#)

Log in

Test Event

Lorem ipsum dolor sit amet, consectetur adipisicing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Bullets

- can
- be used

Location
Somerset, UK
Starts
2014-06-15
Ends
2014-06-15
Organizers
• Sample Organization
Sponsors
• Sample Sponsor
Theme
Territory

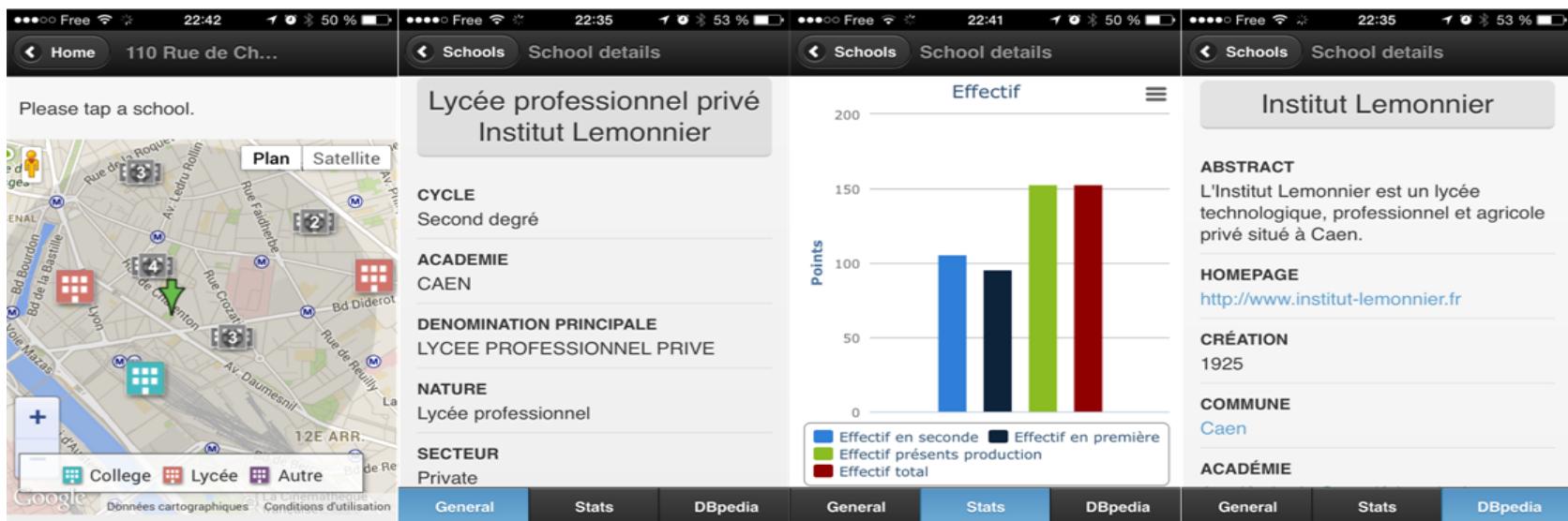
Developing Web Applications: from specific to generic approach

■ Scenario 1

- Known Datasets, Known vocabularies → Specific SPARQL queries
- Visualizations: dataset specific

■ Example

- Datasets on schools in France
- Vocabularies: geo vocab, data cube, geometrie,
- Application: [PerfectSchool](#)



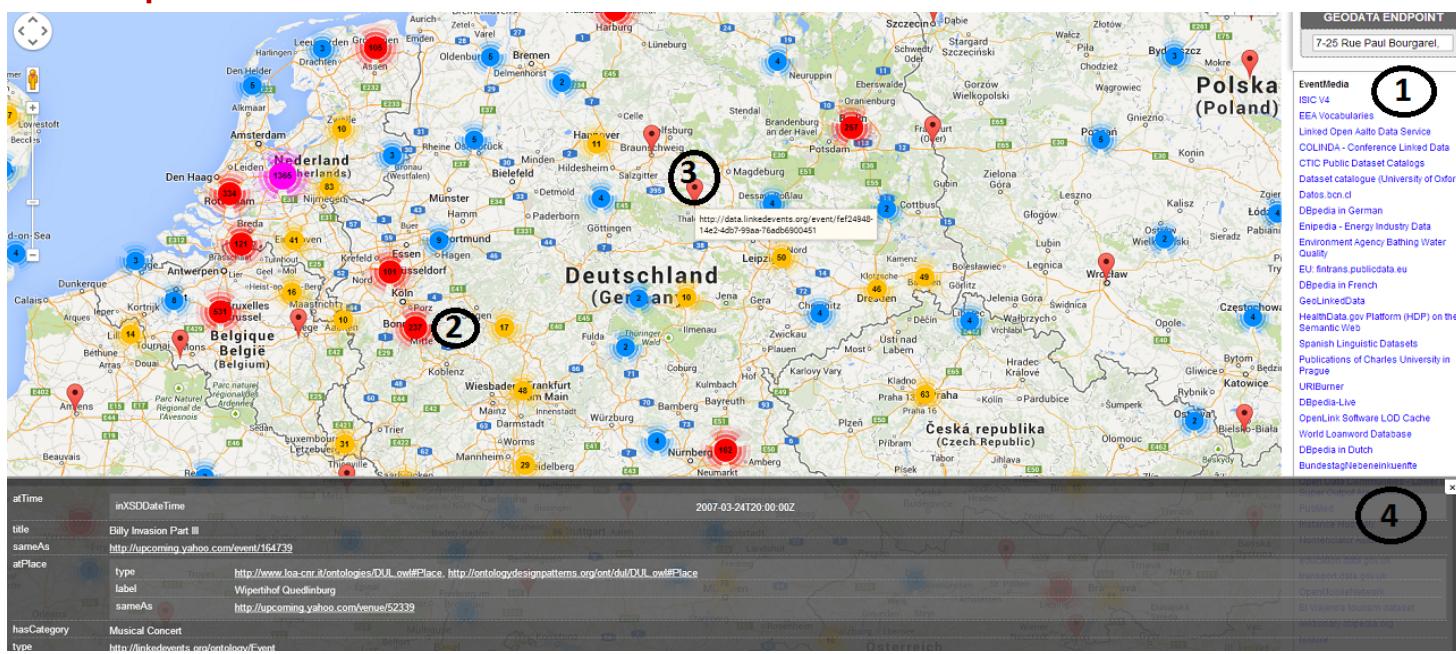
Developing Web Applications: from specific to generic approach

■ Scenario 2

- Unknown Datasets, Known domains, so domain-specific SPARQL queries
 - Visualizations: domain specific

■ Example

- Endpoints of geodatasets
 - Domain: geospatial
 - Application: [GeoRDFviz](#)



Developing Web Applications: from specific to generic approach

■ Scenario 3

- Unknown Datasets, Unknown domains, so generic SPARQL queries
- Visualizations: adapted to domains specific
- Any endpoints
- Multiple domains: geodata, statistics, persons, cross-domains, etc..
- Application: ?????

Related work on configuring Semantic Web widgets by data mapping. [1]

Application: Efficient search for Semantic News demonstrator in Cultural Heritage Dataset

Tool: [ClioPatria](#)



...but "*method not apply to create interfaces on top of arbitrary SPARQL endpoints*"

[1] Hildebrand, Michiel, and Jacco Van Ossenbruggen. "Configuring semantic web interfaces by data mapping." *Visual Interfaces to the Social and the Semantic Web (VISSW 2009)* 443 (2009): 96.

Our Proposal

Linked Data Vizualization Wizard (LDVizWiz)

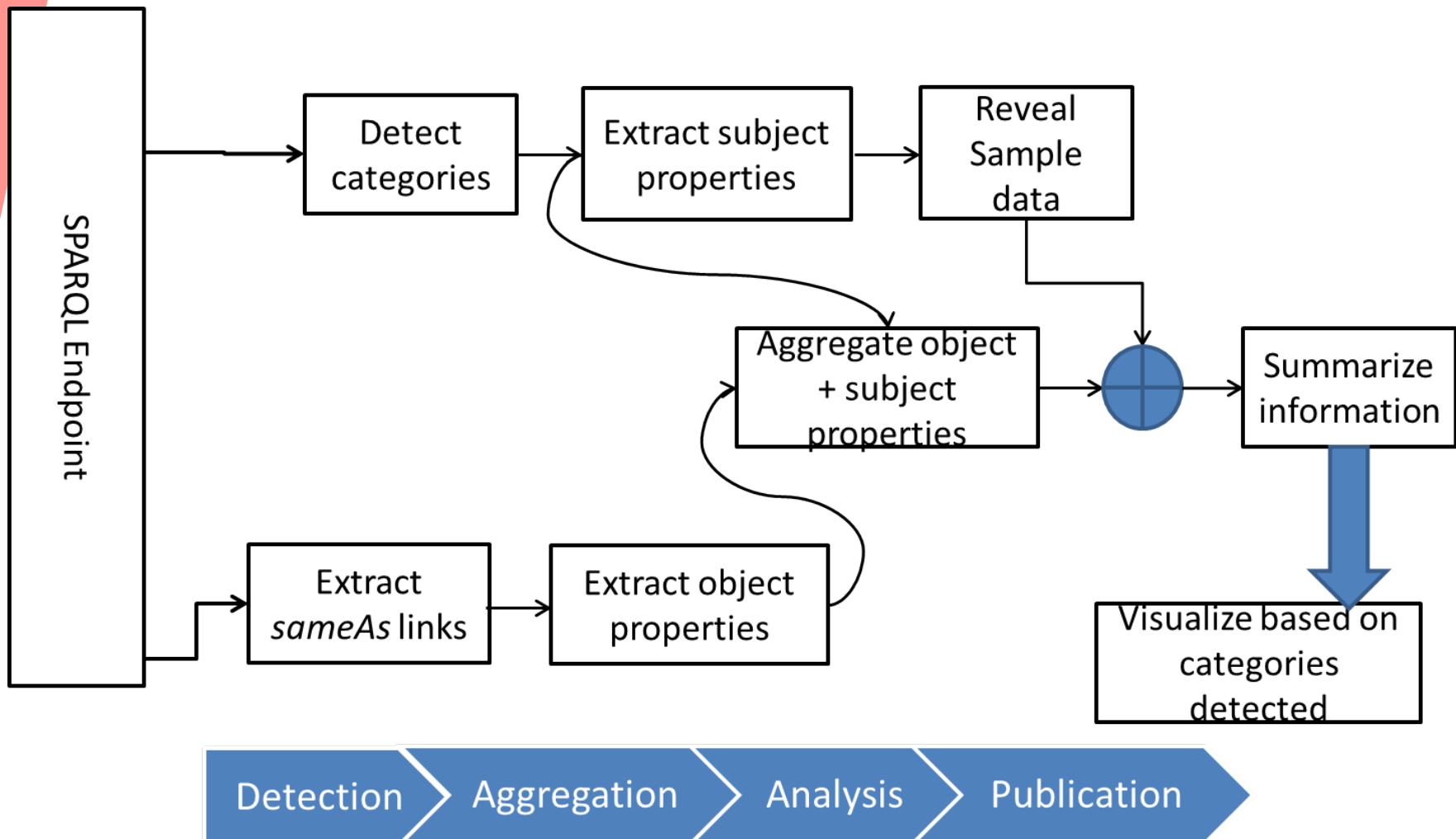
Requirements of LDVizWiz (LDViz-”Wise”)

- Predefined **categories** associated to visual elements
- Build on top of RDF standards
 - e.g., SPARQL queries ; Semantic Web technologies
- Reuse existing Visualization libraries
 - e.g., Google Maps, Google Charts, D3.js, etc.
- Reuse On-line Library of Information Visualization ([OLIVE](#))
- Target to non “RDF/SPARQL speakers”
- Input: Datasets published as LOD

Mapping Categories and vocabularies

- **Geographic information**
 - geo: vocab, schema:Place, etc.
- **Temporal information**
 - Time, interval ontologies
- **Event information**
 - lode, event, sport, etc..
- **Agent/Person**
 - foaf:Person/foaf:Agent
- **Organization information**
 - ORG vocabulary, foaf:Organization
- **Statistics information**
 - Data cube, SDMX model
- **Knowledge information**
 - Schemas, classifications using SKOS vocabulary

LDVizWiz Workflow



Step 1: Categories Detection

- **Detection of main categories in datasets**
 - ASK SPARQL queries on predefined categories
 - Uses well-known vocabularies in LOV
 - Condition the type of visual elements

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	qb, sdmx-model, scovo	Charts, graphs
Tree	skos, Government space	Treemap, Org view
Network	any vocab.	Graph, network map

Detection



Olive

A close-up photograph of a pile of green olives, showing their texture and color.

Experiment: Categories Detection

Category	Number	%
GEO DATA	97	21.84%
EVENT DATA	16	3.60%
TIME DATA	27	6.08%
SKOS DATA	02	0.45%
ORG DATA	48	10.81%
PERSON DATA	59	13.28%
STAT DATA	29	6.6%

- 444 endpoints (*) analyzed, 278 good answers (**62.61%**) using ASK queries.
- Few **taxonomies** in SKOS, **many GEO DATA**

■ Applications

- Automatic detection of endpoints categories
- More “*trustable*” than human tagging
- Map categories detected with “suitable” visual elements for the visualizations (e.g., TimeLine + maps for events data)



Detection

(*) All the endpoints retrieved from sparqles.org

Step 2: Properties Aggregation

- **Goal: Exploit the “connectors” between graphs**
- **“connectors” are used to enrich a given graph**
 - e.g., owl:sameAs ; rdfs:seeAlso ,
skos:exactMatch
- **Retrieve properties from external datasets**
 - So called “enriched properties”
- **Build candidate properties for visualization**
 - For pop-up menus
 - For facet browsing
 - For charts display



Step 3: Publication

■ Visualization Generator

- Recommend the visual elements based on categories
- Transform ASK queries to SELECT or CONSTRUCT queries for input to visual library.

■ Visualization Publisher

- Export the description of visualization in RDF
- Add metadata for the visualization (charts) and steps used to create it.
- e.g., dcat:Dataset, prov:wasDerivedFrom, chart vocabulary, void:ExampleResource.

Detection

Aggregation

Publication

Current Implementation

- Javascript light version as “*proof-of-concept*”
- Url: <http://semantics.eurecom.fr/datalift/rdfViz/apps/>

The screenshot shows a user interface for dataset analysis. At the top, there's a 'Dataset URL' input field with the value 'http://eventmedia.eurecom.fr/sparql' and a 'Analyse' button. Below this, several green and red status boxes provide information about the dataset:

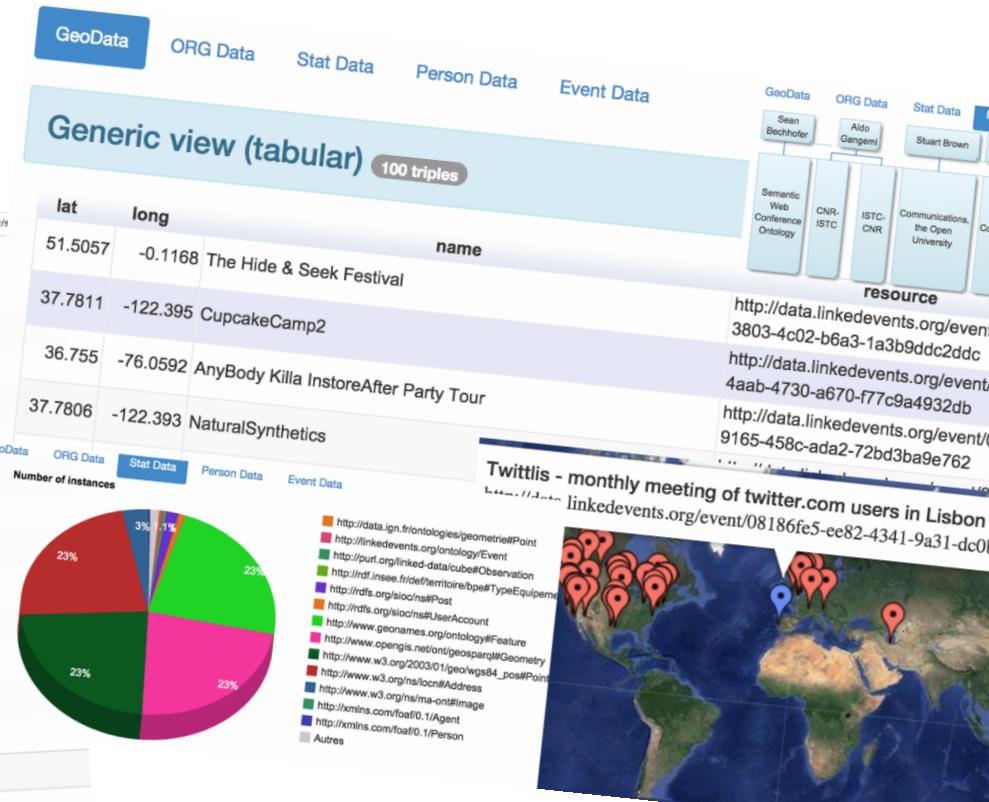
- Presence of Geodata
- No Time data detected! (TIME-dbpedia)
- No DATE data detected! (xsd:date)
- Presence of EVENT data (LODE ,dbpedia)
- Presence of personal data (FOAF, DBPEDIA)
- Presence of organization data (FOAF, ORG, dbpedia)
- Presence of statistical data(Datacube)

A central panel titled 'Properties' lists the following triples:

```
xsd:date, owl:DatatypeProperty
ASK WHERE { ?p a owl:DatatypeProperty+ rdfs:range ?r
range, filter (?r
range= xsd:date
range= xsd:dat
rdf:Property; rdf
range= xsd:dat
xsd:dateTime)}
```

Below this is a table with columns 'lat' and 'long'.

lat	long	name
51.5057	-0.1168	The Hide & Seek Festival
37.7811	-122.395	CupcakeCamp2
36.755	-76.0592	AnyBody Killa InstoreAfter Party Tour
37.7806	-122.393	NaturalSynthetics



Part 3

Best Practices for metadata in vocabularies

W3C Govn't LD Best Practices

- **10** best practices to help government worldwide to access and reuse their by taking benefit of Linked Data mechanism.
- **4** steps to rich **5★** ratings datasets in TimBL scale.



Best Practices for Publishing Linked Data

W3C Working Group Note 21 December 2013

This version:

<http://www.w3.org/TR/2013/NOTE-ld-bp-20131221/>

Latest published version:

<http://www.w3.org/TR/lد-bp/>

Latest editor's draft:

<https://dvcs.w3.org/hg/gld/raw-file/default/bp/index.html>

Previous version:

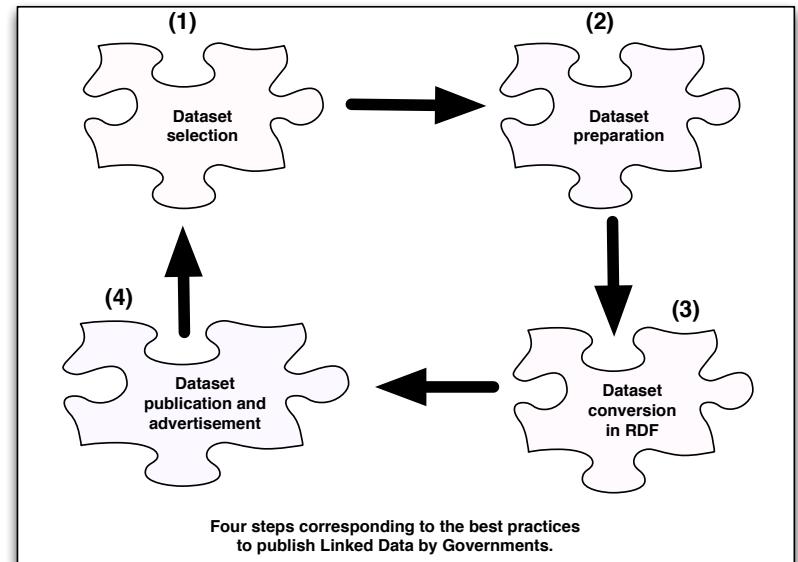
<http://www.w3.org/TR/2013/WD-lد-bp-20130625/>

Editors:

Bernadette Hyland, 3 Round Stones, Inc.

Ghislain Atemezing, EURECOM

Boris Villazón-Terrazas, iSOCO, Intelligent Software Components S.A.

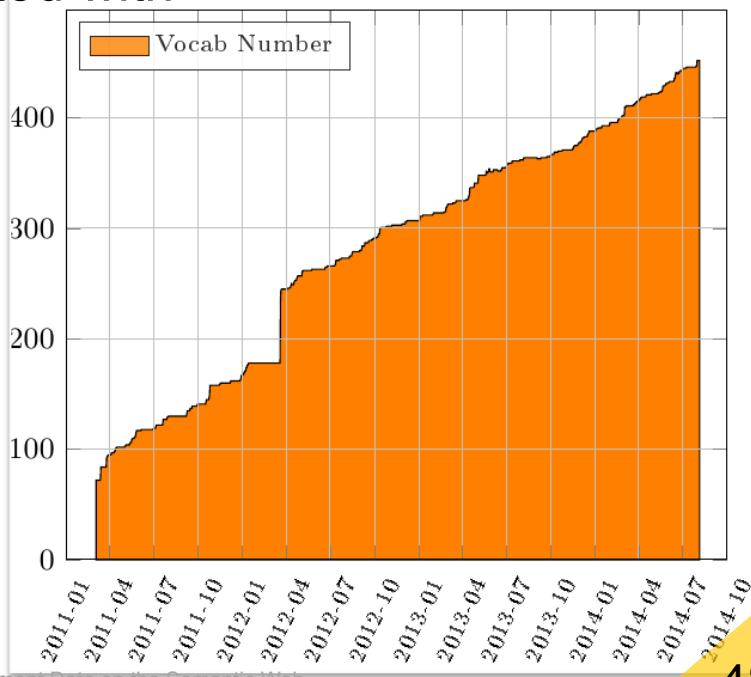


Copyright © 2012-2013 W3C® (MIT, ERCIM, Keio, Beihang). W3C [liability](#), [trademark](#) and [document use](#) rules apply.

■ A curated list of vocabularies

- More than 495 vocabularies
- Each of them described by **vocabulary-of-a-friend (voaf)**
- Provide a dump in N3 of the different versions of a vocabulary
- Quasi **linearity** of the growth, started with 75 vocabularies

The glitch in 2012 corresponds to the migration to OKFN hosting



Focus on vocabularies: Disambiguating Vocabulary Prefixes

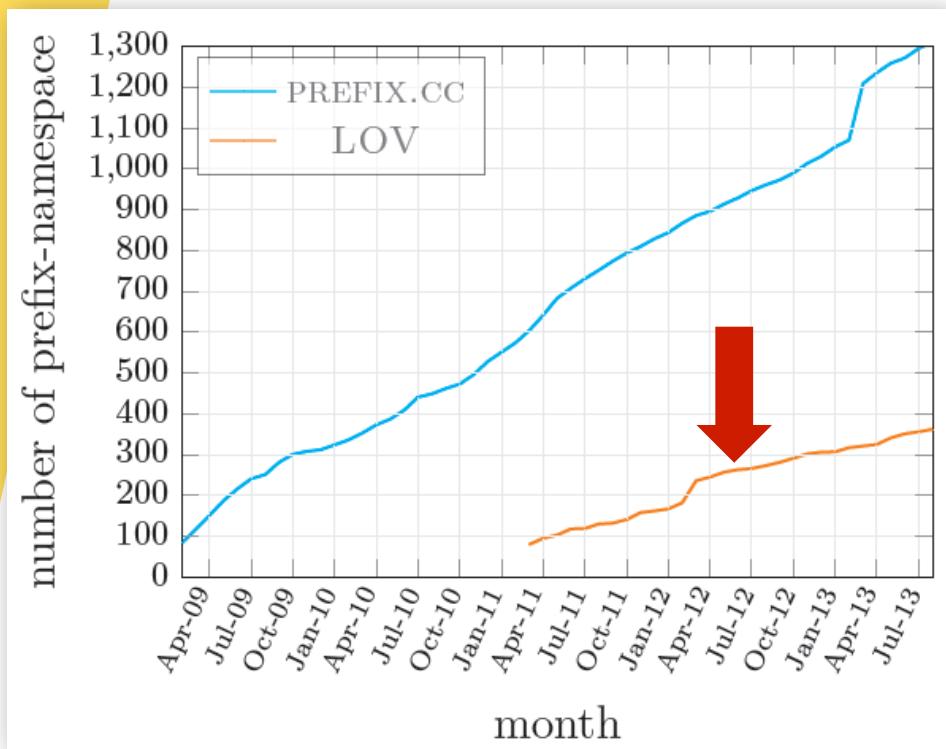
Goal:

**align services against Linked Open Vocabularies to
harmonize and manage vocabularies' namespaces**

Prefix	URI
SearchResults	http://www.zillow.com/static/xsd/SearchResults.xsd
UpdatedPropertyDetails	http://www.zillow.com/static/xsd/UpdatedPropertyDetails.xsd
a	http://www.w3.org/2005/Atom
aapi	http://rdf.alchemyapi.com/rdf/v1/s/aapi-schema#
activity	http://activitystrea.ms/spec/1.0/
address	http://schemas.talis.com/2005/address/schema#
admin	http://webns.net/mvc#/
aiiso	http://purl.org/Vocab/aiiso/schema#
amz	http://webservices.amazon.com/AWSECommerceService/2011-08-01
atom	http://atomowl.org/ontologies/atomrdf#
audio	http://purl.org/media/audio#
awol	http://bblfish.net/work/atom-owl/2006-06-06#
aws	http://soap.amazon.com/
b3s	http://b3s.openlinksw.com/
batch	http://schemas.google.com/gdata/batch
bestbuy	http://remix.bestbuy.com/
bibo	http://purl.org/ontology/bibo/
bif	bif:
book	http://purl.org/NET/book/vocab#
bookmark	http://www.w3.org/2002/01/bookmark#
bugzilla	http://www.openlinksw.com/schemas/bugzilla#
buzz	http://schemas.google.com/buzz/2010
c	http://www.w3.org/2002/12/cal/icaltzd#
category	http://dbpedia.org/resource/Category:
	http://www.crunchbase.com/

- **Global namespaces**
 - With good practices to recommend a prefix
 - Have a more transparent list of built-in prefixes
 - All the services understand each other with prefixes
 - Some de facto prefixes emerging: **rdfs:, foaf:, rdf:, owl:, skos:,**

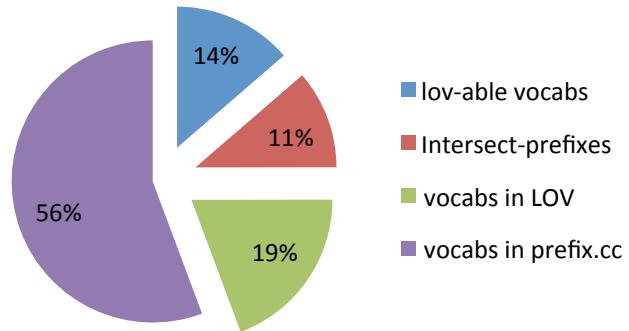
LOV vs PREFIX.CC-Alignment Findings



More than 200 prefixes in prefix.cc are vocabularies

Category	Number
lov-able vocab	227
Intersect-prefixes	188
vocab in LOV	321
vocab in prefix.cc	925

Findings during alignment process



Vocabulary Search and Ranking

Goal

Ranking vocabularies based on Information Content Metrics

- **Metrics**

- Information Content Metric (IC): **value of information** associated with a given entity
- Partition Information Content Metric (PIC)
- Proposed a ranking based on IC and PIC

Ranking Algorithm

Candidate terms selection in LOV

Grouping terms by namespace & weight assignment

Compute IC score

Compute PIC score

Output ranking

Require: Dump of *lovaggregator* file

```
1: Upload in a triple store for querying
2: Select subset of candidate LOVaggregatorendpoint 1 vocab
3: for term ∈ lovaggregator do
4:   if (LOVdistribution ≥ 1 )and (LOVPopularity ≥ 1) then
5:     candidateterms ← append term
6:   end if
7: end for
8: for each term ∈ candidateterms do 2
9:   GROUP BY vocabulary namespace
10:  COMPUTE weight for each vocabulary
11: end for
12: INITIALIZE PICvector AS a vector
13: for each term ∈ candidateterms do
14:   while term ∈ vocabularySpace do 3
15:     ICterm ← function IC(term, vocabPrefix)
16:     ICvocab ←  $\sum$  ICterm
17:   end while
18:   PICvocab ← weight(vocab) × ICvocab
19:   PICvector ← append (PICvocab)
20: ORDER PICvector
21: end for
22: return PICvector
```

Ranking Algorithm

- **dcterms:**
<http://purl.org/dc/terms/>
- **Candidate terms: 53 (39 properties + 14 classes)**
- **wf = 1+ 2+3 = 6**
- **PIC = 1724.844**

- **foaf:**
<http://xmlns.com/foaf/0.1/>
- **Candidate terms: 35 (26 properties + 9 classes)**
- **wf = 1+ 2+ 3 = 6**
- **PIC = 1033.197**

Publisher	Dublin Core Metadata Initiative
Class number	22
Property number	55
Homepage	http://dublincore.org/documents/dcmi-terms
See also	http://stats.lod2.eu/vocabularies/9
Represented by	format-dcterms
Has review	(2013-03-07) Bernard Vatant : Prefix restore (2014-03-14) Bernard Vatant : This vocabulary purl redirection is sometimes down, like at

Property number	62
Homepage	http://www.foaf-project.org/
See also	http://stats.lod2.eu/vocabularies/90
Represented by	format-foaf
Has review	(2013-06-04) Bernard Vatant : From the specification 2000. There is now a stable core of classes and their documentation to track implementation feed (2014-01-15) Bernard Vatant : Looking forward for

PIC(dcterms) > PIC(foaf)

Ranking Algorithm

- Top-15 terms (IC value)
- Top-15 vocabs (PIC value)

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

Rank	Prefix	PIC score
1	dcterms	1724.844
2	schema	1588.700
3	gr	1261.101
4	foaf	1033.197
5	bibo	876.205
6	time	816.2020
7	skos	805.287
8	dul	797.328
9	ptop	773.167
10	rdafrbr	640.834
11	vaem	630.621
12	ma-ont	508,694
13	prov	497.524
14	swrc	437.394
15	dce	428.618

Applications of the Ranking Metrics

- **Vocabulary life-cycle management**
 - Help assessing the use of terms and vocabulary updates
 - Monitoring the use of *owl:deprecated* or
http://www.w3.org/2003/06/sw-vocab-status/ns#:term_status
- **Semantic Web applications**
 - Vocabularies with **higher PIC** might be proposed to a user as much as possible, e.g. for choosing properties to display in a faceted browsing interface
- **Interlinking datasets**
 - Generate *sameAs* links between resources based on vocabularies terms with **lower IC** value

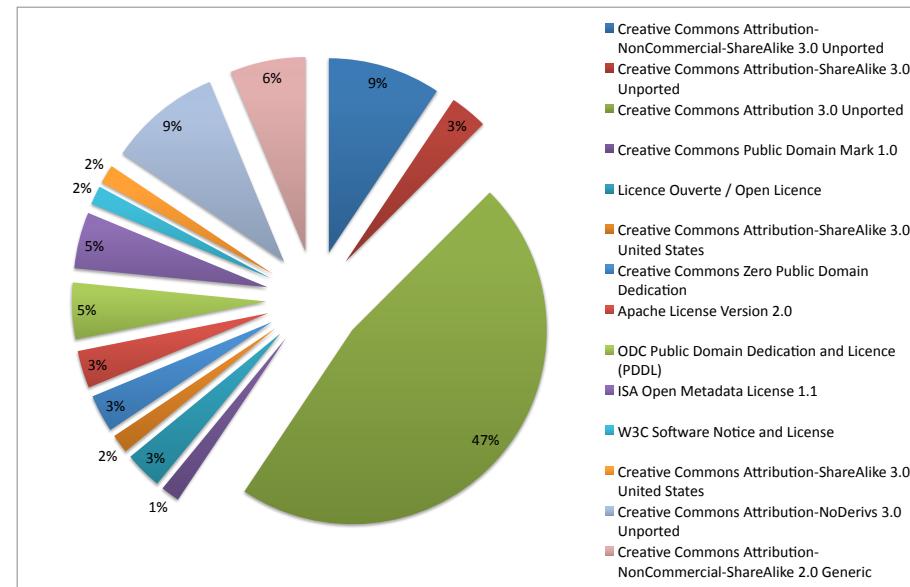
Licenses Compatibility

Goal

Reasoning on Licenses for checking compatibilities between vocabularies and datasets

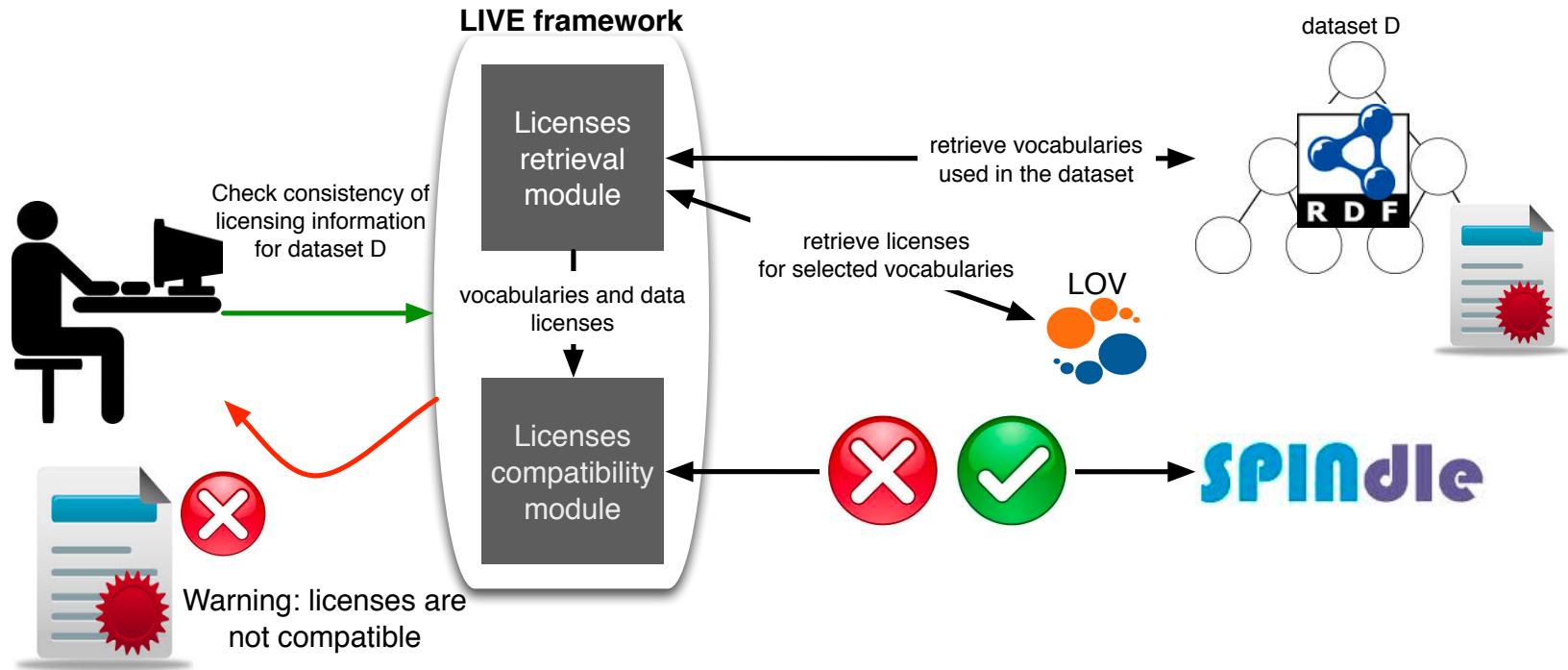
Approach:

- Licenses retrieval both for the dataset and for the vocabularies used in the dataset.
- Deontic logic [1] to compute compatibility using RDF representation of licenses



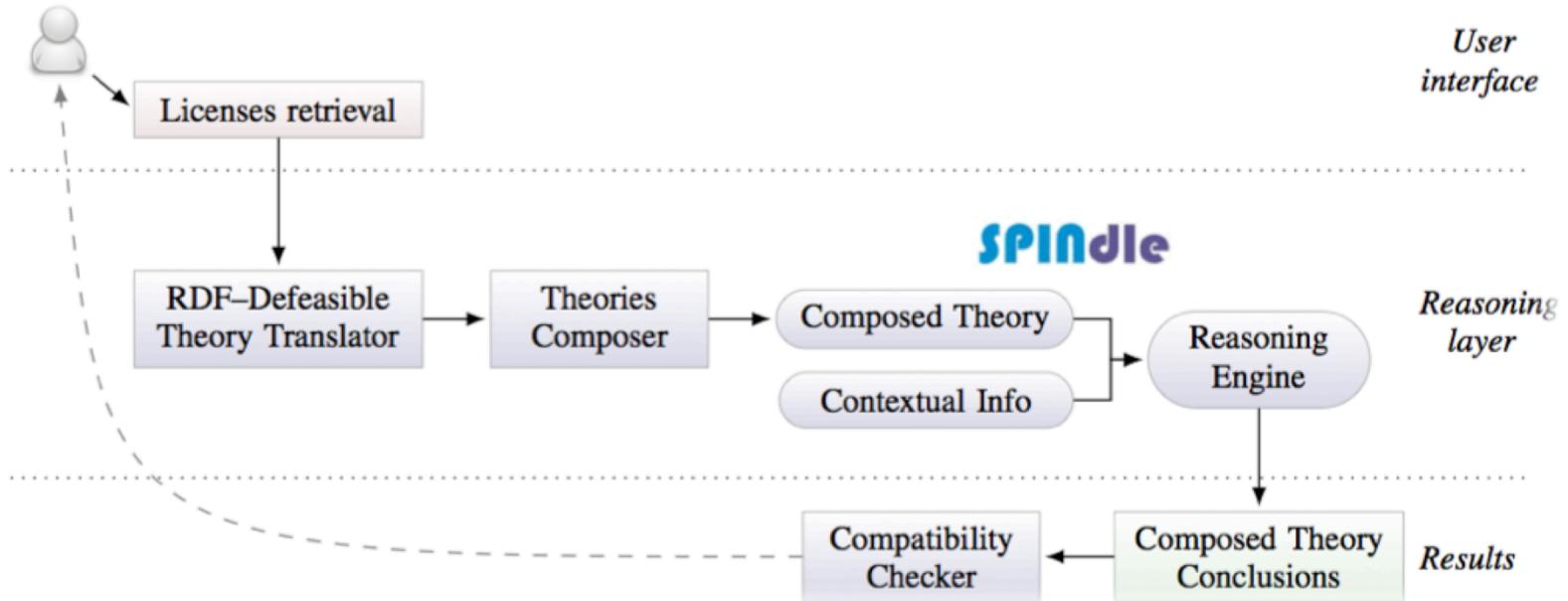
[1] Governatori, Guido, et al. "One License to Compose Them All." The Semantic Web–ISWC 2013. Springer Berlin Heidelberg, 2013. 151–166.

Our Proposal: LIVE Framework



<http://www.eurecom.fr/~atemezin/licenseChecker/>

The Logic



- **RDF licenses:** <http://purl.org/NET/rdflicense>
- **Logic of deontic rules:**
 - constructive account of basic deontic modalities (obligation, prohibition, permission)
 - compute the set of all conclusions for each license and then check whether incompatible conclusions are obtained.

Live Evaluation

Dataset	LicRetrieval(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

LIVE provides the compatibility in less than **5 seconds** for 7 datasets

LIVE retrieves 48 vocabularies in less than **14 seconds**

Conclusions

- **We have presented models for Geodata**
 - For better handling complex geometries
 - Supporting ***all*** CRS
 - For easy querying in SPARQL
- **We have proposed a generic visualization wizard**
 - Based on predefined categories
 - Targeted to lay-users
- **We have contributed in managing vocabulary metadata**
 - Disambiguating prefixes between vocabularies
 - Improving vocabulary search and ranking
 - Providing a license compatibility framework
 - Contributing to Best Practices standard

Future Work

- **Managing Data updates and versioning**
 - Versioning: integrate [Memento](#) protocol?
 - Spatio-temporal evolution
- **Multiple representation: need for metadata?**
 - Level of detail
 - Geometry modeling rules and reasoning
- **Tracking Provenance of geodata**
 - To ensure quality of published dataset
 - To ensure trust from application consumers

Future Work

▪ Visualizations

- Extend categories and vocabularies for detection
- Provide templates for generating “mash-ups” to combine domains, an mash-up widget generator
- Investigate the “*importance*” of a category in dataset
- Provide a user evaluation

▪ Metadata management

- Publish a list of common recommended prefixes
- Foster and support current effort towards a more sustainable governance of vocabularies.
- Compare (P)IC with other graph-based ranking (e.g. pagerank)
- Investigate the dependency ranking between vocabularies

Publications

- Pierre-Yves Vandenbussche, **G.A. Atemezing**, Maria Poveda, Bernard Vatant: Linked Open Vocabularies (LOV): a gateway to reusable semantic vocabularies on the Web. Semantic Web Journal, *under review*, 2015.
- **G.A. Atemezing**, Raphael Troncy: Modeling visualization tools and applications on the Web. Semantic Web Journal, *under review*, 2015.
- **G.A. Atemezing** et al.: Transforming meteorological data into linked data. In Semantic Web journal, Special Issue on Linked Dataset descriptions, 2012. IOS Press.
- Guido Governatori, Ho-Pun Lam, Antonino Rotolo, Serena Villata, **G.A. Atemezing** and Fabien Gandon: LIVE: a Tool for Checking Licenses Compatibility between Vocabularies and Data.(ISWC 2014, Demo Track)
- **G.A. Atemezing** and Raphael Troncy: Information content based ranking metric for linked open vocabularies. (SEMANTICS 2014)
- Ahmad Assaf, **G.A. Atemezing**, Raphael 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 (Demo Track, ESWC 2014)

Publications

- Francois Scharffe, **G.A Atemezing**, Raphael Troncy, Fabien Gandon, Serena Villata, Bénédicte Bucher, Faycal Hamdi, Laurent Bihanic, Gabriel Kepeklian, Franck Cotton, Jerome Euzenat, Zhengjie Fan, Pierre-Yves Vandenbussche and Bernard Vatant: Enabling linked-data publication with the datalift platform. (AAAI, W10:Semantic Cities, 2012)
- **G.A. Atemezing** and Raphael Troncy: Vers une meilleure interopérabilité des données géographiques françaises sur le Web de données. (IC 2012).
- Houda Khrouf, **G.A. Atemezing**, Thomas Steiner, Giuseppe Rizzo and Raphael Troncy: Confomaton: A conference enhancer with social media from the cloud. (ESWC 2012, Demo Track)
- Bernadette Hyland, **G.A. Atemezing** and Boris Villazón-Terrazas (editors): Best Practices for Publishing Linked Data. W3C Working Group Note published on January 9, 2014. URL: <http://www.w3.org/TR/ld-bp/>
- Bernadette Hyland, **G.A Atemezing**, Michael Pendleton, Biplav Srivastava (editors): Linked Data Glossary. W3C Working Group Note published on June 27, 2013. URL: www.w3.org/TR/ld-glossary/

Thank you for your attention!



The stars of the show: Miss Globe and Mr Cube, created by Frans Knibbe of Geodan

Credits layout

Mariella Sabatino: mll.sabatino@gmail.com