



# MidTerm Report

## Self-Service Data Provisioning Through Semantic Enrichment of Data

Ahmad Assaf

EURECOM-Multimedia Communications  
Institut Mines-Télécom  
April 21st, 2014

**Supervisors:**  
Raphaël Troncy  
Aline Senart

**EURECOM**  
**SAP**

## 1. Introduction

Enterprises use a wide range of heterogeneous information systems in their business activities such as Enterprise Resource Planning (ERP), Customer Relationships Management (CRM) and Supply Chain Management (SCM) systems. An enterprise distributed IT landscape contains multiple systems using different technologies and data standards [1]. In addition to this heterogeneity, the amount of information in enterprise databases and on-line data stores expands exponentially each year. Enterprise Big Data isn't big in volume only, but in the associated file formats. The information is also often stored often in unstructured and unknown formats.

Data integration is the problem of combining data residing at different sources, and providing the user with a unified view of these data [2]. In large enterprises, it is a time and resource costly task. Various approaches have been introduced to solve this integration challenge. These approaches were primarily based on XML as the data representation syntax, Web Services to provide the data exchange protocols and Service Oriented Architecture (SOA) as a holistic approach for distributed systems architecture and communication [3][4]. However, it was found that these technologies are no sufficient to solve the integration problems in large enterprises. Recently, ontology-based data integration approaches have been suggested where ontologies are used to describe the data, queries and mappings between them [5]. A slightly different approach is the use of the Linked Data paradigm [6] for integrating enterprise data. Enterprises like Google and Microsoft are not only using the Linked Data integration paradigm for their information systems, but are also aiming at building enterprise knowledge bases ( like the Google Freebase<sup>1</sup>) that will act as a crystallization point for their structured data.

Linked Open Data (LOD) movement has gained lots of momentum in the last years. From 12 datasets cataloged in 2007, the Linked Open Data has grown to almost 300 datasets containing almost 32 billion triples [7]. Data is being published by both public and private sectors and covers a diverse set of domains from life sciences to military. This success lies in the cooperation between data publishers and consumers. Users are empowered to find, share and combine information in their applications easily.

Despite the legal issues surrounding Linked Data licenses [8], it is still considered a gold mine for organizations who are trying to leverage external data sources in order to produce more informed business decisions [9]. In [10] the authors see the potential economic effect unfolding in education, transportation, consumer products, electricity, oil and gas, health care and consumer finance. They estimate the potential annual value enabled by Open Data in these domains to be 3 trillion US Dollars across seven domains.

Data becomes more useful when it is open, widely available and in shareable formats, and when advanced computing and analysis can yield from it. The quality and amount of structured knowledge available make it now feasible for companies to mine this huge amount of public data and integrate it in their next-generation enterprise information management systems. Analyzing this new type of data within the context of existing enterprise data should bring them new or more accurate business insights and allow better recognition of sales and market opportunities [11].

Business Intelligence (BI) has always been about creating new insight for business by converting data into meaning that can be shared between people to drive change in the organization. One key aspect of creating meaning is driving a common shared understanding of information also known as Semantics.

Classic BI and even the newer Agile Visualization tools focus much of their selling features on attractive and unique visualizations, but preparing data for those visualizations still remains the far more challenging task in most BI projects large and small. self-service data provisioning aims at tackling this problem by providing intuitive datasets discovery, acquisition and integration techniques intuitively to the end user.

---

<sup>1</sup><http://freebase.com>

In this thesis, we aim at creating a framework that will enable self-service data provisioning in the enterprise. Our goal is to provide a mechanism that annotates and profiles tabular data and provide better dataset descriptions. Furthermore, we aim at providing a complete data quality metric and aggregate publicly available datasets descriptions so that people can search and browse through content

## 2. Challenges

In this thesis, we aim at creating a framework that leverages Semantic Web technologies in order to enrich enterprise data in general and Business Intelligence data in particular in order to facilitate self-service data provisioning. We investigate the following research challenges:

- **Dataset Integration and Enrichment:** The enterprise heterogeneous data sources raise tremendous challenges. They have inherently different file formats, access protocols or query languages. They possess their own data model with different ways of representing and storing the data. Data across these sources may be noisy (e.g. duplicate or inconsistent), uncertain or be semantically similar yet different [12]. Integration and provision of a unified view for these heterogeneous and complex data structures therefore require powerful tools to map and organize the data.
  - \* Attaching metadata and Semantic information to instances can be tricky. An entity is usually not associated with a single generic type in the knowledge base, but rather with a set of specific types which can be relevant or not given the context. The challenging task is finding the most relevant entity type within a given context.
  - \* Entities play a key role in knowledge bases in general and in the Web of Data in particular. Entities are generally described with a lot of properties, this is the case for DBpedia. It is, however, difficult to assess which ones are more “important” than others for particular tasks such data augmentation and visualizing the key facts of an entity.
  - \* Social Networks are not just gathering Internet users into groups of common interests, they are also helping people follow breaking news, contribute to online debates or learn from others. They are transforming Web usage in terms of users’ initial entry point, search, browsing and purchasing behavior. Integrating information from these Social Networks can be tricky due to the vast amount of data available which makes hard to spot what is relevant in a timely manner.
- **Dataset Discovery:** Even though popular datasets like DBpedia<sup>2</sup> and Freebase are well known and widely used, there are other hidden useful datasets not being used. Indeed these datasets may be useful for specialized domains, however without proper registry of topics, it is difficult for users to find them [13].
- **Dataset Quality Control:** Linked Data consists of structured information supported by models, ontologies and vocabularies and contains query endpoints and links. This makes data quality assurance a challenge. Despite the fact that Linked Open Data quality is a trending and highly demanded topic, very few efforts are currently trying to standardize, track and formalize frameworks to issue scores or certificates that will help data consumers in their integration tasks.

## 3. Proposal

Linked Open Data datasets are described using either the Vocabulary of Interlinked Datasets (VOID) [14] or the Data Catalog Vocabulary (DCAT) [15]. With these standards, discovery and usage of linked datasets can be performed both effectively and efficiently. In our framework, we plan to use DCAT as the common standard for homogenizing description metadata of datasets indexed by our crawler. This

---

<sup>2</sup><http://dbpedia.org>

choice came from the fact that the Open Data Support<sup>3</sup> is promoting the DCAT-AP (and consequently DCAT) as the standard for describing datasets and catalogs in Europe. In order to enable self-service data provisioning, we envisage building a framework (see Figure 1) that will be able to provide detailed DCAT descriptions for internal and external data sources. The framework is able to provide the following services:

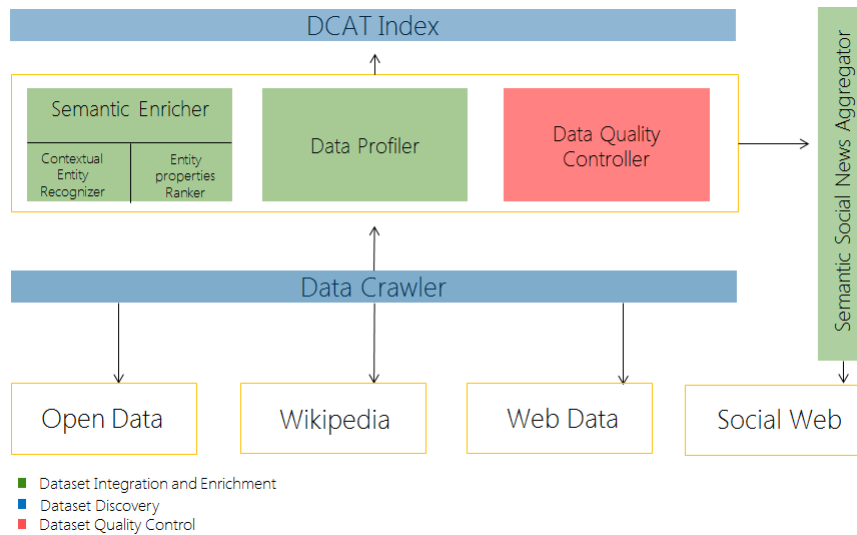


Figure 1. Overall Architecture

- **Data Acquisition:** Be able to sample data from the various structured data sources like Wikipedia tables, Open data portals, etc. While these are rich sources of information, sometimes live information streamed from the Social Networks is needed. As a result, we crawl Social Networks in order to aggregate semantically related information and connect it with the right resources.
- **Data Preparation:** This includes data profiling and validation, de-duplicating and enhancing relevant data sets with metadata. Profiling is used to examine data to understand its content, structure and data quality dependencies. The types of profiling tasks include:
  - \* Examining column data and getting statistical information such as min, max, average, median, null percentage, value distribution, pattern distribution.
  - \* Dependency tasks: Finds the values in one or more dependent columns that rely on values in a primary column
  - \* Redundancy tasks: Determine the degree of overlapping data values or duplication between two sets of columns
  - \* Uniqueness tasks: Returns the count and percentage of rows that contain non-unique data, for the set of column(s) selected.
  - \* Content type: Content type profiling provides suggested meaning based on the entities data in the columns.
  - \* Quality checks: An important aspect that we have to take into consideration while describing a dataset is its quality. For that, an objective Linked Data quality assessment framework should be created in order to issue quality profiles that extend the DCAT vocabulary. The framework helps on one hand data owners to rate the quality of their dataset and get some hints on possible

<sup>3</sup><http://opendatasupport.eu>

improvements, and on the other hand data consumers to choose their data sources from a ranked set.

- **Dataset Classification:** Classify and organize datasets based on the input from the previous tasks.

### 3.1. Contributions on Dataset Integration and Enrichment

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

- Building RUBIX which is a framework enabling mash-up of potentially noisy enterprise and external data.
- RUBIX improves instance and schema matching by adding Semantic metadata to data at the instance level.
- RUBIX improves data integration techniques by enabling clean representation of the data regardless of the languages used, existence of abbreviations, synonyms and typos.
- Build a Social crawler that queries several Social endpoints and aggregates news based on a set of defined keywords.
- Build a common Semantic model to represent Web documents.
- Building a Semantic Social News Aggregating service (SNARC) [16] that aggregates relevant Social news with regards to a Web document.
- Reverse engineering of Google Knowledge Graph Panel in order to define the top properties of a selected entity.
- High traction from the BI organization to integrate it as a service into various offerings.
- We presented RUBIX at the First International Workshop on Open Data [17][18].
- SNARC has won the first place at the AI Mashup Challenge<sup>4</sup> at ESWC13.

### 3.2. Contributions on Dataset Quality Control

Concerning our contributions on Linked Data quality assessment, we have achieved the following tasks:

- We identified five principle classes to describe the quality of a particular linked dataset. For each class, we list the principles that are involved at all stages of the data management process.
- We have presented our Data quality principles at the Sixth IEEE International Conference on Semantic Computing [19].
- We have surveyed the landscape of Linked Data quality assessment frameworks.
- We have surveyed the landscape of Linked Data quality assessment tools.
- We have refined the five principles in [19] towards a more objective framework.
- We have evaluated the surveyed tools with regards to the suggested framework.

## 4. Dataset Integration and Enrichment

Tagging and attaching metadata is often seen as additional work for data publishers with few paybacks. Moreover, different data creators use different terminologies which means that the same object maybe be represented using different metadata descriptions [20]. Presenting and enterprise taxonomies requires a considerable amount of time and effort, at least in the initial creation steps [4].

---

<sup>4</sup><http://aimashup.org/aimashup13>

#### 4.1. What are the Important Properties of 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 when augmenting extra columns into an existing dataset or when annotating instances with semantic tags.

Our approach consists in: (i) reverse engineering the Google Knowledge Panel by extracting the properties that Google considers as sufficiently important to show (Section ??), and (ii) analyzing users' preferences by conducting a user survey and comparing the results (Section ??).

We have shown that it is possible to reveal what are the "important" properties of entities in a large knowledge base by reverse engineering the choices made by Google when creating knowledge graph panels and by comparing with users preferences obtained from a user survey. Our motivation is to represent this choice explicitly, using the Fresnel vocabulary, so that any application could just read this configuration file for deciding which properties of an entity is worth to visualize. We are aware that this knowledge is highly dynamic, the Google Knowledge Graph panel differing from countries and varying along the time. We have provided the code that enables to perform new calculation at run time and we aim to study the temporal evolution of what are important properties on a longer period. This knowledge which has been captured will be made available shortly in a SPARQL endpoint. We are also investigating the use of Mechanical Turk to perform a larger survey for the complete set of DBpedia classes.

#### 4.2. RUBIX to Enhance Schema Matching

RUBIX is our approach to bootstrap the process of attaching meta information to data objects. We leverage DBpedia and Freebase as knowledge bases for our annotation process. In RUBIX we assign a vector of Semantic types for every object at the instance level. For example, Orange will be represented by a vector of rich types that contain (**Organization**, **Organism Classification**, **Place** ...). Currently, we rely on Freebase API to identify these rich types, but we have already started the effort to build an entity type ranking tool inspired by [21].

#### 4.3. RUBIX to Enhance Schema Matching

In the past, some work has tried to improve existing data schema [22] but literature mainly covers automatic or semi-automatic labeling of anonymous data sets through Web extraction. Examples include [23] that automatically labels news articles with a tree structure analysis or [24] that defines heuristics based on distance and alignment of a data value and its label. These approaches are however restricting label candidates to Web content from which the data was extracted. [25] goes a step further by launching speculative queries to standard Web search engines to enlarge the set of potential candidate labels. More recently, [26] applies machine learning techniques to respectively annotate table rows as entities, columns as their types and pairs of columns as relationships, referring to the YAGO ontology. The work presented aims however at leveraging such annotations to assist semantic search queries construction and not at improving schema matching. With the emergence of the Semantic Web, new work in the area has tried to exploit Linked Data repositories. The authors of [27] present techniques to automatically infer a semantic model on tabular data by getting top candidates from Wikitology [28] and classifying them with the Google page ranking algorithm. Since the authors' goal is to export the resulting table data as Linked Data and not to improve schema matching, some columns can be labeled incorrectly, and acronyms and languages are not well handled [27]. In the Helix project [29], a tagging mechanism is used to add semantic information on tabular data. A sample of instances values for each column is taken and a set of tags with scores are gathered from online sources such as Freebase. Tags are then correlated to infer annotations for the column. The mechanism is quite similar to ours but the resulting tags for the column are independent

of the existing column name and sampling might not always provide a representative population of the instance values.

In RUBIX we have implemented several matching algorithms (Cosine Similarity, Pearson Product-Moment Correlation Coefficient (PPMCC) and Spearman’s Rank Correlation Coefficient) in order to calculate similarity between data based on their rich types population. Our preliminary evaluation shows that for datasets where mappings were relevant yet not proposed, our framework provides higher quality matching results. Additionally, the number of matches discovered is increased when Linked Data is used in most datasets.

#### *4.4. Dataset Annotation and Domain Identification*

The increasing diversity of the datasets makes it difficult to annotate them with a fixed number of pre-defined tags. Moreover, manually entered tags are subjective and may not capture the essence and breadth of the dataset [13]. RUBIX can be used to Annotate datasets with meta information based on examining the data instances themselves. In addition to that, we can enhance our approach by applying techniques similar to [13] in order to identify topical domains with fine-grained classification.

#### *4.5. Semantic Social News Aggregation (SNARC)*

The Internet has created a paradigm shift in how we consume and disseminate information. Data nowadays is spread over heterogeneous silos of archived and live data. People willingly share data on social media by posting news, views, presentations, pictures and videos. SNARC is a service that uses semantic web technology and combines services available on the web to aggregate social news. SNARC brings live and archived information to the user that is directly related to his active page. The key advantage is an instantaneous access to complementary information without the need to dig for it. Information appears when it is relevant enabling the user to focus on what is really important.

Crawling data from these heterogeneous platforms implies the studying of related API specifications which differ in terms of use, privacy policy and described methods. Harvesting content spread over multiple platforms is a challenging task that has to ensure an easy and flexible way for integrating several social Web APIs. Tools such as API Blender [30] or Media Finder [31] provide such interface with the aim to save developers’ efforts for learning each API specification. However, for in SNARC we needed services that are not available in the mentioned services. As a result, we have implemented our own Social crawler for Twitter<sup>5</sup>, Google+<sup>6</sup>, YouTube<sup>7</sup>, Vimeo<sup>8</sup>, Slideshare<sup>9</sup>, Stack Exchange<sup>10</sup> and the Web.

## **5. Dataset Quality Control**

We are entering an era where open is the new default. Governments, universities, organizations and even individuals are publicly publishing huge amounts of open data. This openness should be accompanied with a certain level of trust or guarantees about the quality of data. To our knowledge, only one certificate is available to data publishers to assess the quality level of their datasets, the ODI certificate<sup>11</sup>.

---

<sup>5</sup><http://www.twitter.com>

<sup>6</sup><http://plus.google.com>

<sup>7</sup><http://www.youtube.com>

<sup>8</sup><http://www.vimeo.com>

<sup>9</sup><http://www.slideshare.com>

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

<sup>11</sup><https://certificates.theodi.org/>

This certificate provides a description of the published data quality in plain English. It aspires to act as a mark of approval that helps publishers understand how to publish good open data and users how to use it. It gives publishers the ability to provide assurance and support on their data while encouraging further improvements through an ascending scale.

ODI comes as an online and free questionnaire for data publishers focusing on certain characteristics about their data. The questions are classified into the following categories: general information (about dataset, publisher and type of release), legal information (e.g., rights to publish), licensing, privacy (e.g., whether individuals can be identified), practical information (e.g., how to reach the data), quality, reliability, technical information (e.g., format and type of data) and social information (e.g., contacts, communities, etc.). Based on the information provided by the data publisher, a certificate is created with one of four different ratings.

Although ODI is a great initiative, the issued certificates are self-certified. ODI does not verify or review submissions but retains the right to revoke a certificate at any time. The dynamicity of Linked Data makes it also very difficult to update the certificates manually, especially when these changes are frequent and affect multiple categories. There is clearly a need for automatic certification which can be supplemented with some manual input for categories that cannot be processed by machines.

The emerging critical need for large, distributed, heterogeneous, and complex structured datasets identified the necessity to establish industry cooperation between vendors of RDF and Graph database technologies in developing, endorsing, and publishing reliable and insightful benchmark results. The Linked Data Benchmark Council (LDBC)<sup>12</sup> aims to bridge the gap between the industry and the new trending stack of semantic technologies and their vendors.

LDBC more specifically aims at developing new benchmarks that will lead to significant progress in scalability, storage, indexing and query optimization techniques to become the de facto standard for publishing performance results. LDBC is promising initiative, but it is still work in progress with the final report expected on the first quarter of 2015.

In addition to the initiatives mentioned above, there exists a number of data quality frameworks and tools that are either standalone or implemented as modules in data integration tools.

LODGR<sub>efine</sub><sup>13</sup> is the Open Refine<sup>14</sup> of Linked Data. It does not act as a quality assessment tool, but it is powerful in cleaning and refining raw instance data. LODGR<sub>efine</sub> can help detect duplicates, empty values, spot inconsistencies, extract Named Entities, discover patterns and more. LODGR<sub>efine</sub> helps in improving the quality of the dataset by improving the quality of the data at the instance level.

PROLOD [32] is also not a quality assessment tool. It is a Linked Data profiling tool that provides clustering and labeling capabilities, schema discovery and statistics about data types and patterns. The statistics are about properties distribution, link-to-literal ratio, number of entities and RDF triples, average properties per entity and average error. PROLOD had been tested with DBpedia but the authors plan to improve its scalability to larger datasets.

Sieve [33] is framework for expressing quality assessment and fusion methods. It is implemented as a component of the Linked Data Integration Framework (LDIF)<sup>15</sup>. Sieve leverages the LDIF provenance metadata as quality indicators to produce quality assessment scores. However, despite its nice features, it is only targeted to perform data fusion based on user-configurable conflict resolution tasks. Moreover, since

---

<sup>12</sup><http://ldbc.eu/>

<sup>13</sup><http://code.zemanta.com/sparkica/>

<sup>14</sup><http://openrefine.org/>

<sup>15</sup><http://ldif.wb3g.de/>



Sieve main input is provenance metadata, it is only limited to domains that can provide such metadata associated with their data.

Quality Assessment of Data Sources (Flemming's Data Quality Assessment Tool)<sup>16</sup> calculates data quality scores based on manual user input. The user should assign weights to the predefined quality metrics and answer a series of questions regarding the dataset. These include, for example, the use of obsolete classes and properties by defining the number of described entities that are assigned disjoint classes, the usage of stable URIs and whether the publisher provides a mailing list for the dataset. The main disadvantage for using this tool is the manual intervention which requires deep knowledge in the dataset examined. Moreover, the tool lacks support for several quality concerns like completeness or consistency.

SWIQA [34] is composed of three layers: data acquisition, query and ontology layers. It uses query templates based on the SPARQL Inferencing Notation (SPIN)<sup>17</sup> to express quality requirements. The queries are built to compute weighted and unweighted quality scores. At the end of the assessment, it uses vocabulary elements to annotate important values of properties and classes, assigning inferred quality scores to ontology elements and classifying the identified data quality problems.

Despite all the recent efforts in providing frameworks and tools for data quality in Linked Open Data, there is still no framework for the objective assessment of such quality taking into account all aspects of Linked Open Data.

In our previous work [19] we have identified 24 different Linked Data quality attributes. We have refined these attributes into a condensed framework of 13 objective attributes. Since these attributes are rather abstract, we should rely on quality indicators that reflect data quality [35]. In this paper, we transform the quality indicators presented as a set of questions in [19] into more concrete quality indicator metrics. We extend them with the the objective quality indicators listed in the systematic review done in [36].

Table(1) lists the refined attributes along with their quality indicators.

---

<sup>16</sup><http://linkeddata.informatik.hu-berlin.de/LDSrcAss/datenquelle.php>

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

Table 1: Objective Assessment Framework for Linked Data Quality

Quality Attribute	Quality Category	ID	Quality Indicator
Completeness	Entity Level	QI.1	Covers of all the attributes needed for a given task [33]
		QI.2	Has complete language coverage [37]
		QI.3	Existence of documentation properties [38][37]
	Dataset Level	QI.4	Existence of all the necessary objects for a given task [33]
		QI.5	Existence of supporting structured metadata [39]
		QI.6	Supports multiple serializations [36]
		QI.7	Includes the correct MIME-type for the content [39]
		QI.8	Contains appropriate volume of data for a particular task [36]
		QI.9	Has different queryable endpoints to access the data [36]
		QI.10	Checked against syntactic errors [39]
		QI.11	Usage of datasets description vocabularies
		QI.12	Existence of descriptions about its size and categorization
	Links Level	QI.13	Existence of complete dereferenceable in-bound and out-bound links [39][37][40]
		QI.14	Existence of supporting linkage metadata [39]
	Model Level	QI.15	Covers the complete set of values [37]
		QI.16	Absence of disconnected graph clusters [37]
		QI.17	Absence of omitted top concept [39]
		QI.18	Absence of unidirectional related concepts [39]
		QI.19	Existence of supporting metadata about the kind and number of used vocabularies [36]
Availability	Dataset Level	QI.20	Existence of an RDF dump that can be downloaded by users [35][39]
		QI.21	Existence of queryable endpoints that respond to direct queries
Correctness	Entity Level	QI.22	Absence of missing or empty labels [41][37]
		QI.23	Absence of incorrect data type for typed literals [39][41]
		QI.24	Absence of omitted or invalid languages tags [42][37]
		QI.25	Absence of terms without any associative or hierarchical relationships [43]
	Links Level	QI.26	Existence of content related to the subject of the RDF triple [42][41]
		QI.27	Absence of syntactic errors [44]
	Model Level	QI.28	Contains marked top concepts [37]
QI.29		Absence of broader concepts for top concepts [37]	
Conciseness	Entity Level	QI.30	Absence of redundant attributes [33]
		QI.31	Existence of short URIs [36]
	Dataset Level	QI.32	Absence of redundant objects [33]
		QI.33	Follows the HTTP URI scheme [45][44]
Security	Dataset Level	QI.34	Uses login credentials to restrict access [36]
		QI.35	Uses SSL or SSH to provide access to their dataset [36]
		QI.36	Grants access to specific users [36]
Freshness	Entity Level	QI.37	Existence of timestamps that can keep track of its modifications [46]
Licensing	Dataset Level	QI.38	Existence of machine readable license information [45]
		QI.39	Existence of human readable license information [45]
		QI.40	Specifies permissions, copyrights and attributions [36]
Continued on next page			

Table 1 Objective Assessment Framework for Linked Data Quality

Quality Attribute	Quality Category	ID	Quality Indicator
Comprehensibility	Dataset Level	QI.41	Existence of at least one exemplary URI [36]
		QI.42	Existence of at least one exemplary SPARQL query [36]
		QI.43	Existence of regular expression pattern that matches the URIs of a dataset [36]
		QI.44	Existence of a list of used vocabularies
		QI.45	Existence of a mailing list or message board [35]
Consistency	Entity Level	QI.46	Existence of consistent preferred labels per language tag [47][37]
		QI.47	Absence of overlapping labels
		QI.48	Absence of disjoint labels [37]
		QI.49	Absence of extra white spaces in labels [42]
		QI.50	Existence of only one value of skos:prefLabel without a language tag [37][42]
	Dataset Level	QI.51	Absence of conflicting information [33]
	Model Level	QI.52	Absence of atypical use of collections, containers and reification [39]
		QI.53	Absence of overlapping usage of owl:sameAs and owl:differentFrom [39]
		QI.54	Absence of overlapping usage of owl:AllDifferent and owl:distinctMembers [39]
		QI.55	Absence of asserted members of owl:Nothing [39]
		QI.56	Absence of membership violations for disjoint classes [39]
Coherence	Model Level	QI.57	Absence of misplaced or deprecated classes or properties [39]
		QI.58	Absence of misused owl:DataTypeProperty or owl:ObjectProperty [39]
		QI.59	Absence of relation and mappings clashes [42]
		QI.60	Absence or minimal usage of blank nodes [45]
		QI.61	Absence of invalid inverse-functional values [39]
		QI.62	Absence of cyclic hierarchical relations [48][42][37]
		QI.63	Absence of undefined classes and properties usage [39]
		QI.64	Absence of solely transitive related concepts [37]
		QI.65	Absence of redefinitions of existing vocabularies [39]
		QI.66	Absence of valueless associative relations [37]
		QI.67	Absence of incomplete literals with datatype range [39]
Efficiency	Dataset Level	QI.68	Absence of slash-URIs [36]
		QI.69	Acceptable delay between the request and its response [49]
		QI.70	Low Latency HTTP requests [36]
		QI.71	Ability to scale [36]
Accuracy	Dataset Level	QI.72	Absence of outliers [36]
		QI.73	Absence of attributes that do not contain useful values for data entries [36]
Provenance	Entity Level	QI.74	Ability to construct decision networks informed by provenance graphs [50]
	Dataset Level	QI.75	Existence of metadata that describes its authoritative information [46]
		QI.76	Reliability and Trustworthiness of the publisher [46]
		QI.77	Usage of a provenance vocabulary
		QI.78	Usage of digital signatures [36]
	Model Level	QI.79	Trustworthiness of RDF statements [51]

As a result, we have identified the need for a complete quality framework that can assess all the quality indicators. Most of the tools were designed with limited coverage to certain aspects, for example, ontology and vocabulary checkers focus mainly on the coherence, completeness and correctness at the modeling level. Flemming’s tool covers several attributes like the completeness, correctness, conciseness, security, licensing and comprehensibility, but it falls short in measuring the consistency, coherence, efficiency and provenance.

We have also noticed the lack of tools to measure certain quality indicators like the dataset’s security. Moreover, to our knowledge, there are no tools that can measure all the provenance quality indicators (except for Flemming’s tool that is able to check for the use of digital signatures) although the literature covers several approaches to achieve that [51][46][52].

## 6. Conclusions and Future work

We have presented in this document our main contributions in some issues around Enriching Enterprise Data Towards self-service Data Provisioning. We have first focused on the aspect of data profiling through RUBIX. We plan to extend RUBIX to be able to work with DBpedia leveraging the planned entity type ranking module. We also plan to include data mining techniques to profile numerical data and provide statistical insights about data distribution.

Regarding the Linked Data quality module, we plan to develop a comprehensive objective Linked Data quality evaluation tool. The tool will be able to automatically measure the various quality indicators listed in this paper, introduce a scoring function with different weights for the various quality attributes and issue a quality certificate.

Regarding the Social integration, we would like to test SNARC on business web application, check if our annotations can be used to successfully query and attach relevant social snippets to the data.

We also plan to build our Linked Data crawler that will be responsible for the data acquisition phase which is the entry point for the work done so far. We also plan to investigate possible extensions to the current data description vocabularies to allow more comprehensive datasets categorization.

## References

- [1] Nandana Mihindukulasooriya, Raul Garcia-Castro, and Miguel Esteban Gutiérrez. Linked data platform as a novel approach for enterprise application integration. In *COLD*, 2013.
- [2] Maurizio Lenzerini. Data integration: A theoretical perspective. In *Proceedings of the Twenty-first ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*, PODS ’02, pages 233–246, New York, NY, USA, 2002. ACM.
- [3] Philipp Frischmuth, Sren Auer, Sebastian Tramp, Jrg Unbehauen, Kai Holzweigg, and Car-Martin Marquardt. Towards linked data based enterprise information integration. In *Proceedings of the Workshop on Semantic Web Enterprise Adoption and Best Practice (WASABI) 2013*, 2013.
- [4] Philipp Frischmuth, Jakub Klmeck, Sren Auer, Sebastian Tramp, Jrg Unbehauen, Kai Holzweigg, and Carl-Martin Marquardt. Linked data in enterprise information integration. 2012.
- [5] H. Wache, T. Vgele, U. Visser, H. Stuckenschmidt, G. Schuster, H. Neumann, and S. Hbner. Ontology-based integration of information - a survey of existing approaches. pages 108–117, 2001.
- [6] C. Bizer, T. Heath, and T. Berners-Lee. Linked data - the story so far. *Int. J. Semantic Web Inf. Syst.*, 5(3):122, 2009.
- [7] Chris Bizer, Anja Jentzsch, and Richard Cyganiak. State of the lod cloud, 2011.
- [8] Krzysztof Janowicz Prateek Jain, Pascal Hitzler and Chitra Venkatramani. Theres no money in linked data. 2013.
- [9] D Boyd and Kate Crawford. Six provocations for big data. *A Decade in Internet Time: Symposium on the Dynamics of the Internet and Society*, pages 1–17, 2011.
- [10] Diana Farrell Steve Van Kuiken Peter Groves James Manyika, Michael Chui and Elizabeth Almasi Doshi. Open data: Unlocking innovation and performance with liquid information. *McKinsey Business Technology Office*, 2013.
- [11] Rebecca Shockley Michael S. Hopkins Steve LaValle, Eric Lesser and Nina Kruschwitz. Big data, analytics and the path from insights to value. *MIT Sloan Management Review*, 2011.

- [12] G. Sudha; Shenoy Sangeetha N avitha, C.; Sadasivam. Ontology based semantic integration of heterogeneous databases. *European Journal of Scientific Research*;11/13/2011,, Vol. 64(Issue 1):p115, 2011.
- [13] Pascal Hitzler Amit Sheth Sarasi Lalithsena, Prateek Jain. Automatic domain identification for linked open data. *IEEE/WIC/ACM International Conference on Web Intelligence*, 2013.
- [14] Richard Cyganiak, Jun Zhao, Keith Alexander, and Michael Hausenblas. Describing linked datasets with the VoID vocabulary. W3C note, W3C, March 2011. <http://www.w3.org/TR/2011/NOTE-void-20110303/>.
- [15] Fadi Maali and John Erickson. Data catalog vocabulary (DCAT). Last call WD, W3C, August 2013. <http://www.w3.org/TR/2013/WD-vocab-dcat-20130801/>.
- [16] Ahmad Assaf, Aline Senart, and Raphaël Troncy. Snarc - an approach for aggregating and recommending contextualized social content. In *ESWC (Satellite Events)*, pages 319–326, 2013.
- [17] Ahmad Assaf, Eldad Louw, Aline Senart, Corentin Follenfant, Raphaël Troncy, and David Trastour. Rubix: a framework for improving data integration with linked data. In *Proceedings of the First International Workshop on Open Data, WOD '12*, pages 13–21, New York, NY, USA, 2012. ACM.
- [18] Ahmad Assaf, Eldad Louw, Aline Senart, Corentin Follenfant, Raphaël Troncy, and David Trastour. Improving schema matching with linked data. *CoRR*, abs/1205.2691, 2012.
- [19] Ahmad Assaf and Aline Senart. Data quality principles in the semantic web. *CoRR*, abs/1305.4054, 2013.
- [20] G. W. Furnas, T. K. Landauer, L. M. Gomez, and S. T. Dumais. The vocabulary problem in human-system communication. *COMMUNICATIONS OF THE ACM*, 30(11):964–971, 1987.
- [21] Alberto Tonon, Michele Catasta, Gianluca Demartini, Philippe Cudré-Mauroux, and Karl Aberer. Trank: Ranking entity types using the web of data. In *ISWC*, 2013.
- [22] Renée J. Miller and Periklis Andritsos. Schema discovery. *IEEE Data Eng. Bull.*, 26(3):40–45, 2003.
- [23] Davi de Castro Reis, Paulo B. Golgher, Altigran S. da Silva, and Alberto H. F. Laender. Automatic web news extraction using tree edit distance. In Stuart I. Feldman, Mike Uretsky, Marc Najork, and Craig E. Wills, editors, *Proceedings of the Thirteenth International World Wide Web Conference*, pages 502–601, New York, NY, May 2004. ACM Press.
- [24] Jiying Wang and Frederick H Lochovsky. Data Extraction and Label Assignment for Web Databases. *The World Wide Web Conference*, pages 187–196, 2003.
- [25] Altigran Soares da Silva, Denilson Barbosa, João M. B. Cavalcanti, and Marco A. S. Sevalho. Labeling data extracted from the web. In *OTM Conferences (1)*, pages 1099–1116, 2007.
- [26] Girija Limaye, Sunita Sarawagi, and Soumen Chakrabarti. Annotating and searching web tables using entities, types and relationships. *Proc. VLDB Endow.*, 3(1-2):1338–1347, September 2010.
- [27] Zareen Syed, Tim Finin, Varish Mulwad, and Anupam Joshi. Exploiting a Web of Semantic Data for Interpreting Tables. In *Proceedings of the Second Web Science Conference*, April 2010.
- [28] Tim Finin, Zareen Syed, James Mayfield, Paul McNamee, and Christine D. Piatko. Using wikilogy for cross-document entity coreference resolution. In *AAAI Spring Symposium: Learning by Reading and Learning to Read*, pages 29–35. AAAI, 2009.
- [29] Oktie Hassanzadeh, Songyun Duan, Achille Fokoue, Anastasios Kementsietsidis, Kavitha Srinivas, and Michael J. Ward. Helix: online enterprise data analytics. In *WWW (Companion Volume)*, pages 225–228. ACM, 2011.
- [30] Georges Gouriten and Pierre Senellart. Api blender: A uniform interface to social platform apis. *CoRR*, abs/1301.2086, 2013.
- [31] Hyunmo Kang and Ben Shneiderman. Mediafinder: an interface for dynamic personal media management with semantic regions. In Gilbert Cockton and Panu Korhonen, editors, *CHI Extended Abstracts*, pages 764–765. ACM, 2003.
- [32] Christophe Bohm, Felix Naumann, Ziawasch Abedjan, Fenz Dandy, Toni Grutze, Daniel Hefenbrock, Matthias Pohl, and David Sonnabend. Proling Linked Open Data with ProLOD.PDF. *ICDE 2010*, 2010.
- [33] PN Mendes, Hannes Mühleisen, and Christian Bizer. Sieve: linked data quality assessment and fusion. *LWDM2012 - Proceedings of the 2012 Joint EDBT*, 2012.
- [34] C Fürber and M Hepp. SWIQAA Semantic Web information quality assessment framework. *ECIS 2011*, 2011.
- [35] A Flemming. Quality characteristics of linked data publishing datasources, 2010.
- [36] Conceptual Framework, Amrapali Zaveri, Anisa Rula, Andrea Maurino, Ricardo Pietrobon, and Jens Lehmann. Quality Assessment Methodologies for Linked Open Data. *Under review, Semantic Web Journal*, 1:1–5, 2012.
- [37] Christian Mader, Bernhard Haslhofer, and Antoine Isaac. Finding quality issues in SKOS vocabularies. *Theory and Practice of Digital Libraries*, 2012.
- [38] A. Miles and S. Bechhofer. Skos simple knowledge organization system reference. w3c recommendation 18 august 2009., 2009.
- [39] Aidan Hogan, Andreas Harth, Alexandre Passant, Stefan Decker, and Axel Polleres. Weaving the pedantic web. *LDOW 2010*, 2010.
- [40] Christophe Guéret, Paul T. Groth, Claus Stadler, and Jens Lehmann. Assessing linked data mappings using network measures. In *Proceedings of the 9th Extended Semantic Web Conference*, volume 7295 of *Lecture Notes in Computer Science*, pages 87–102. Springer, 2012.
- [41] Maribel Acosta, Amrapali Zaveri, Elena Simperl, and Dimitris Kontokostas. Crowdsourcing Linked Data quality assessment. *ISWC 2013*, 2013.
- [42] Osmo Suominen and Eero Hyvönen. Improving the quality of skos vocabularies with skosify. In *Proceedings of the 18th*

- international conference on Knowledge Engineering and Knowledge Management, EKAW'12, pages 383–397, Berlin, Heidelberg, 2012. Springer-Verlag.
- [43] Henry Living. Review of: Hedden, heather. the accidental taxonomist medford, nj: Information today, inc., 2010. *Inf. Res.*, 15(2), 2010.
  - [44] Osma Suominen and Christian Mader. Assessing and Improving the Quality of SKOS Vocabularies. *Journal on Data Semantics*, June 2013.
  - [45] Aidan Hogan, JürRen Umbrich, Andreas Harth, Richard Cyganiak, Axel Polleres, and Stefan Decker. An empirical survey of linked data conformance. *Web Semant.*, 14:14–44, July 2012.
  - [46] Giorgos Flouris, Yannis Roussakis, and M Poveda-Villalón. Using provenance for quality assessment and repair in linked open data. pages 1–12, 2012.
  - [47] Antoine Isaac and Ed Summers. Skos simple knowledge organization system primer. World Wide Web Consortium, Working Draft WD-skos-primer-20080829, August 2008.
  - [48] Dagobert Soergel. Thesauri and ontologies in digital libraries. In Mary Marlino, Tamara Sumner, and Frank M. Shipman III, editors, *JCDL*, page 421. ACM, 2005.
  - [49] Christian Bizer. *Quality-Driven Information Filtering in the Context of Web-Based Information Systems*. PhD thesis, March 2007.
  - [50] Matthew Gamble. Quality, Trust, and Utility of Scientific Data on the Web: Towards a Joint Model.pdf. *WebSci'11*, 2011.
  - [51] Olaf Hartig and Jun Zhao. Using web data provenance for quality assessment. In *In: Proc. of the Workshop on Semantic Web and Provenance Management at ISWC*, 2009.
  - [52] Andreas Harth, Sheila Kinsella, and Stefan Decker. Using naming authority to rank data and ontologies for web search. *ISWC 2009*, 2, 2009.
  - [53] Besiki Stvilia, Les Gasser, Michael B. Twidale, and Linda C. Smith. A framework for information quality assessment. *Journal of the American Society for Information Science and Technology*, 58(12):1720–1733, October 2007.
  - [54] Mamoru Ohtai, Kouji Kozaki and Riichiro Mizoguchi. A Quality Assurance Framework for Ontology Construction and Renement.pdf. *Web Intelligence Conference (AWIC2011)*, 2011.
  - [55] RY Wang and DM Strong. Beyond accuracy: What data quality means to data consumers. *Journal of management information systems*, 1996.
  - [56] Allen Moulton, S Madnick, and M Siegel. Cross-organizational data quality and semantic integrity: learning and reasoning about data semantics with context interchange mediation. *MIT Sloan*, III(1):1–4, 2001.
  - [57] Barbara Pernici and Monica Scannapieco. Data quality in web information systems. *Journal on Data Semantics I*, pages 48–68, 2003.
  - [58] Anisa Rula. DC proposal: towards linked data assessment and linking temporal facts. *ISWC 2011*, 2011.
  - [59] Nickolai Toupikov, J Umbrich, and Renaud Delbru. DING! Dataset ranking using formal descriptions. *WWW09*, 2009.
  - [60] M D'Aquin. Formally measuring agreement and disagreement in ontologies. *K-CAP 09*, 2009.
  - [61] Renaud Delbru, Nickolai Toupikov, and Michele Catasta. Hierarchical link analysis for ranking web data. *The Semantic Web: Research and Applications*, 2010.
  - [62] Stuart Madnick and Hongwei Zhu. Improving data quality through effective use of data semantics. *Data & Knowledge Engineering*, (October):1–19, 2006.
  - [63] Astrid Duque-ramos, Jesualdo Tomás Fernández-breis, Robert Stevens, and Nathalie Aussenac-gilles. OQuaRE : A SQuaRE-based Approach for Evaluating the Quality of Ontologies. *Journal of Research and Practice in Information Technology, Software Engineering and Semantic Web Technologies*, 43(2), 2011.
  - [64] Christian Mader and Bernhard Haslhofer. Perception and Relevance of Quality Issues in Web Vocabularies. *I-SEMANTICS 2013*, 2013.
  - [65] Renaud Delbru. Sindice at SemSearch 2010. *WWW10*, 2010.
  - [66] Stephen Wood. The Equation between Semantics and Data Quality. *Other Conferences*, 2010.
  - [67] Christian Fürber and Martin Hepp. Towards a vocabulary for data quality management in semantic web architectures. *Proceedings of the 1st International Workshop on Linked Web Data Management - LWDM '11*, page 1, 2011.
  - [68] Dimitris Kontokostas, Amrapali Zaveri, S Auer, and J Lehmann. TripleCheckMate: A Tool for Crowdsourcing the Quality Assessment of Linked Data. *4th Conference on Knowledge Engineering and Semantic Web*, pages 1–8, 2013.
  - [69] Graeme Shanks and B Corbitt. Understanding data quality: Social and cultural aspects. *Proceedings of the 10th Australasian Conference on . . .*, (1998):785–797, 1999.
  - [70] Christian Fürber and Martin Hepp. Using semantic web resources for data quality management. *Knowledge Engineering and Management by the . . .*, pages 1–15, 2010.
  - [71] Christian Fürber and Martin Hepp. Using SPARQL and SPIN for Data Quality Management on the Semantic Web. *Business Information Systems*, (1):1–12, 2010.
  - [72] Li Ding, Rong Pan, Tim Finin, and Anupam Joshi. Finding and ranking knowledge on the semantic web. *ISWC 2005*, (November), 2005.
  - [73] Beverly K. Kahn, Diane M. Strong, and Richard Y. Wang. Information quality benchmarks: product and service performance. *Communications of the ACM*, 45(4ve):184–192, April 2002.
  - [74] PY Vandenbussche, CB Aranda, Aidan Hogan, and J Umbrich. Monitoring the Status of SPARQL Endpoints. *ISWC*

2013, 1380(3130617):3–6, 2013.

- [75] Samir Tartir, I Budak Arpinar, Michael Moore, Amit P Sheth, and Boanerges Aleman-meza. OntoQA : Metric-Based Ontology Quality Analysis University of Georgia. *IEEE Workshop on Knowledge Acquisition from Distributed, Autonomous, Semantically Heterogeneous Data and Knowledge Sources*, 2005.
- [76] C Buil-Aranda and Aidan Hogan. SPARQL Web-Querying Infrastructure: Ready for Action? *International . . .*, 2013.
- [77] L Ding, Tim Finin, A Joshi, R Pan, and RS Cost. Swoogle: A semantic web search and metadata engine. *CIKM04*, 2004.
- [78] Christian Fürber and Martin Hepp. Towards a vocabulary for data quality management in semantic web architectures. *Proceedings of the 1st International Workshop on Linked Web Data Management - LWDM '11*, page 1, 2011.
- [79] Li Ding, Tim Finin, and Yun Peng. Tracking rdf graph provenance using rdf molecules. *ISWC 2005*, pages 1–2, 2005.
- [80] Christian Fürber and Martin Hepp. Using SPARQL and SPIN for data quality management on the Semantic Web. *Business Information Systems*, (1):1–12, 2010.
- [81] JLG Sánchez, Roberto García, and JM Brunetti. Using SWET-QUM to Compare the Quality in Use of Semantic Web Exploration Tools. *Journal of Universal Computer Science*, 19(8):1025–1045, 2013.
- [82] Joseph. M. Juran and A. Blanton Godfrey. *Juran's quality handbook*. Juran's quality handbook, 5e. McGraw Hill, 1999.
- [83] Page Lawrence, Brin Sergey, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford University, 1998.
- [84] Berners-Lee Tim. Linked data. Technical report, W3C, July 2006. <http://www.w3.org/DesignIssues/LinkedData.html>.
- [85] Rakesh Agrawal and Ramakrishnan Srikant. Fast algorithms for mining association rules. In *Proc. of 20th Intl. Conf. on VLDB*, pages 487–499, 1994.
- [86] Yuanguai Lei, Victoria Uren, and Enrico Motta. A framework for evaluating semantic metadata. In *Proceedings of the 4th international conference on Knowledge capture*, K-CAP '07, pages 135–142, New York, NY, USA, 2007. ACM.
- [87] Jon M. Kleinberg. Authoritative sources in a hyperlinked environment. *J. ACM*, 46(5):604–632, September 1999.
- [88] Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual web search engine. In *Proceedings of the seventh international conference on World Wide Web 7*, WWW7, pages 107–117, Amsterdam, The Netherlands, The Netherlands, 1998. Elsevier Science Publishers B. V.
- [89] Soumen Chakrabarti, Byron E. Dom, S. Ravi Kumar, Prabhakar Raghavan, Sridhar Rajagopalan, Andrew Tomkins, David Gibson, and Jon Kleinberg. Mining the web's link structure, 1999.
- [90] Aidan Hogan, Andreas Harth, and Stefan Decker. Reconrank: A scalable ranking method for semantic web data with context. In *In 2nd Workshop on Scalable Semantic Web Knowledge Base Systems*, 2006.
- [91] Eyal Oren, Renaud Delbru, Michele Catasta, Richard Cyganiak, and Giovanni Tummarello. Sindice.com: A document-oriented lookup index for open linked data. *International Journal of Metadata, Semantics and Ontologies*, 3:2008.
- [92] Patricia Harpring. *Introduction to Controlled Vocabularies: Terminology for Art, Architecture, and Other Cultural Works*. Getty Research Institute, Los Angeles, 2010.
- [93] Mara Poveda-Villaln, MariCarmen Surez-Figueroa, and Asuncin Gmez-Prez. Validating ontologies with oops! In Annette Teije, Johanna Vlker, Siegfried Handschuh, Heiner Stuckenschmidt, Mathieu dAcquin, Andriy Nikolov, Nathalie Aussenac-Gilles, and Nathalie Hernandez, editors, *Knowledge Engineering and Knowledge Management*, volume 7603 of *Lecture Notes in Computer Science*, pages 267–281. Springer Berlin Heidelberg, 2012.
- [94] Christian Bizer, Jens Lehmann, Georgi Kobilarov, Sren Auer, Christian Becker, Richard Cyganiak, and Sebastian Hellmann. DBpedia - a crystallization point for the web of data. *Web Semantics: Science, Services and Agents on the World Wide Web*, 7(3):154–165, September 2009.
- [95] Evren Sirin, Michael Smith, and Evan Wallace. Opening, closing worlds - on integrity constraints. In Catherine Dolbear, Alan Ruttenberg, and Ulrike Sattler, editors, *OWLED*, volume 432 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2008.
- [96] Jiao Tao, Li Ding, and Deborah L. McGuinness. Instance data evaluation for semantic web-based knowledge management systems. In *HICSS*, pages 1–10. IEEE Computer Society, 2009.
- [97] Diego Berrueta, Sergio Fernandez, and Iván Frade. Cooking http content negotiation with vapour. In *In Proceedings of 4th workshop on Scripting for the Semantic Web 2008 (SFSW2008)*. co-located with *ESWC2008*, 2008.
- [98] Bernhard Haslhofer and Niko Popitsch. Dsnotify: Detecting and fixing broken links in linked data sets. In *8th International Workshop on Web Semantics (WebS &#8217;09)*, co-located with *DEXA 2009*, Berlin, Heidelberg, August 2009. Springer.
- [99] Jens Bleiholder and Felix Naumann. Data fusion. *ACM Comput. Surv.*, 41(1):1:1–1:41, January 2009.
- [100] N.I.S. Organization and National Information Standards Organization (U.S.). *Guidelines for the Construction, Format, and Management of Monolingual Controlled Vocabularies*. National information standards series. NISO Press, 2005.
- [101] Saeedeh Shekarpour and S.D. Katebi. Modeling and evaluation of trust with an extension in semantic web. *Web Semantics: Science, Services and Agents on the World Wide Web*, 8(1), 2010.
- [102] Khalid Belhajjame, James Cheney, David Corsar, Daniel Garijo, Stian Soiland-Reyes, Stephan Zednik, and Jun Zhao. Prov-o: The prov ontology. Technical report, 2012.
- [103] Ernesto Diaz-Aviles, Lucas Drumond, Lars Schmidt-Thieme, and Wolfgang Nejdl. Real-time top-n recommendation in social streams. In *Proceedings of the sixth ACM conference on Recommender systems - RecSys '12*, page 59, New York, New York, USA, 2012. ACM Press.

- [104] Valentina Zanardi and L Capra. Social ranking: uncovering relevant content using tag-based recommender systems. *ACM conference on Recommender systems*, 2008.
- [105] D Preotiuc-Pietro and S Samangooei. Trendminer: An architecture for real time analysis of social media text. *AAAI Publications, Sixth International AAAI Conference on Weblogs and Social Media*, pages 4–7, 2012.
- [106] Iván Cantador and Alejandro Bellogín. Semantic contextualisation of social tag-based profiles and item recommendations. *E-Commerce and Web ...*, (2), 2011.
- [107] Eytan Bakshy, Itamar Rosenn, Cameron Marlow, and Lada Adamic. The role of social networks in information diffusion. In *Proceedings of the 21st international conference on World Wide Web - WWW '12*, page 519, New York, New York, USA, 2012. ACM Press.
- [108] Marco Pennacchiotti and Siva Gurumurthy. Investigating topic models for social media user recommendation. In *WWW11*, page 101, New York, New York, USA, 2011. ACM Press.
- [109] Houda Khrouf, Ghislain Atemezing, and Giuseppe Rizzo. Aggregating Social Media for Enhancing Conference Experiences. *Proceedings of the 1st Int. Workshop on Real-Time Analysis and Mining of Social Streams*, 2012.
- [110] 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.
- [111] Ben Shneiderman. The Eyes Have It: A Task by Data Type Taxonomy for Information Visualizations. In *IEEE Symposium on Visual Languages*, pages 336–343, 1996.
- [112] Mike Bergman. Deconstructing the Google Knowledge Graph.  
<http://www.mkbergman.com/1009/deconstructing-the-google-knowledge-graph>.
- [113] 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.