



HAL
open science

WaRG: Warehousing RDF Graphs

Dario Colazzo, Tushar I. Ghosh, François Goasdoué, Ioana Manolescu,
Alexandra Roatis

► **To cite this version:**

Dario Colazzo, Tushar I. Ghosh, François Goasdoué, Ioana Manolescu, Alexandra Roatis. WaRG: Warehousing RDF Graphs. Bases de Données Avancées, Oct 2013, Nantes, France. hal-00868670

HAL Id: hal-00868670

<https://inria.hal.science/hal-00868670v1>

Submitted on 1 Oct 2013

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

WaRG: Warehousing RDF Graphs*

Dario Colazzo[‡] Tushar I. Ghosh^{§,*} François Goasdoué[†]
Ioana Manolescu^{§,*} Alexandra Roatis^{*,§}

[§]Inria Saclay, France ^{*}Université Paris-Sud, France [‡]Université Paris Dauphine, France [†]Université Rennes 1, France
firstname.lastname [§]@inria.fr ^{*}@lri.fr [‡]@dauphine.fr [†]@univ-rennes1.fr

ABSTRACT

We propose to demonstrate *WaRG*, a system for performing warehouse-style analytics on RDF graphs. To our knowledge, our framework is the first to keep the warehousing process purely in the RDF format and take advantage of the heterogeneity and semantics inherent to this model.

1. INTRODUCTION

Databases of *facts*, each characterized by multiple *dimensions*, whose values are recorded in *measures*, are at the core of *multidimensional data warehouses* (DWs in short) [13]. The facts can then be analyzed by means of *aggregating* the measures, e.g., “*what is the average sale price of item I every month in every store?*”

Data warehouses are typically built to *analyze (some aspects of) an enterprise’s business processes*. First, analysts choose the data sources of interest for a given class business questions. Next, they describe the facts, dimensions, and measures to be analyzed. Then, for each relevant business question, an *analytical query* (commonly known as a *cube*) is formulated, by (i) classifying facts along a set of dimensions and (ii) reporting the aggregated values of their measures. Data warehousing has attracted enormous interest, both from practitioners, e.g., [14] and from the research community, e.g., [8, 12]; warehousing tools are now part of major relational database servers. Relational data warehousing can be thus considered a pretty mature area.

Recent years witnessed a steady and important interest in Semantic Web data, represented within the W3C’s Resource Description Framework (or RDF, in short) [18]. The RDF model allows describing *resources*, by specifying the *values* of their *properties*. The RDF language is used to *export, share, and collaboratively author data*, in many settings. It serves as a *metadata language*, e.g., to describe cultural artifacts in large digital libraries, and to encode protein sequence data (e.g., the Uniprot data set). RDF is a natural target for representing heterogeneous facts contributed by millions of Wikipedia users, gathered within the DBpedia data source, as well as for the *Linked Open Data* effort, aiming at connecting and sharing collectively produced data and knowledge.

The current popularity of RDF raises interest in *models and tools for RDF data analytics*. Consider applications seeking to *harvest, aggregate and analyze user data* from various sources (such as social networks, blog posts, comments on public Web sites). The data is *heterogeneous* (facts about the user age, gender, region, endorsements)

and *graph-structured* (describes relationships between users, places, companies). It comes from multiple sources and may have attached *semantics*, based on ontologies for which RDF is an ideal format. Analyzing such data with warehouse-style tools enables understanding user behavior and needs, and proposing services better suited to the users.

Despite the perceived need, there is currently *no satisfactory conceptual and practical solution for large-scale RDF analytics*. Relational DW tools are not suitable, since loading RDF data in a relational analytical schema may lead to facts with unfilled or multiply-defined dimensions or measures; the latter does not comply with the relational multidimensional setting and DW tools. Further, to fully exploit RDF graphs, *the heterogeneity and rich semantics of RDF data should be preserved* through the warehouse processing chain and up to the analytical queries. In particular, RDF analytical queries should be allowed to *jointly query the schema and the data*, e.g., ask for *most frequently specified properties of a CollegeStudent*, or the *top-three most represented groups of Inhabitants*. Changes to the underlying database (e.g., a new subclass of *Inhabitant*) should not cause the warehouse schema to be re-designed; instead, the new resources (and their properties) should propagate smoothly to the analysis schema and cubes.

We propose to demonstrate WARG, a tool embodying our novel approach for *warehouse-style analytics on Semantic Web data*. WARG supports the first all-RDF formal model for data warehousing. Its main features are:

- a formal model for *analytical schemas*, capturing *the data of interest for a given analysis* of an RDF data set. Importantly, an analytical schema instance is an RDF graph itself, and as such, it preserves the heterogeneity, semantics, and flexible mix of schema and data present in the RDF model;
- *analytical queries*, RDF counter-parts of the relational analytical cubes, and supporting typical analytical operations (slice, dice, drill-down etc.)

In the sequel, Section 2 presents our analytical RDF model, the WARG platform architecture is discussed in Section 3 and we describe use scenarios in Section 4. Finally, we discuss related work in Section 5.

2. DATA MODEL

An **RDF graph** is a set of *triples* of the form **s p o**. A triple states that its *subject* **s** has the *property* **p**, and the value of that property is the *object* **o**. We only consider well-formed triples, **s** $\in U \cup B$, **p** $\in U$ and **o** $\in U \cup B \cup L$, where U is a set of URIs, L of literals (constants), and B are blank nodes (unknown URIs or literals), such that U , B and L are pairwise disjoint.

*This work was performed during the authors’ affiliation with Université Paris-Sud and Inria Saclay, France.

node n	$\lambda(n)$	$\delta(n)$
n_1	<i>person</i>	$q(x):- x \text{ rdf:type Person}$
n_2	<i>message</i>	$q(x):- y \text{ wrote } x,$ $x \text{ inBl } b, b \text{ rdf:type Blog}$
n_4	<i>product</i>	$q(x):- x \text{ rdf:type Product}$
n_5	<i>age</i>	$q(x):- y \text{ hasAge } x$
n_6	<i>pname</i>	$q(x):- y \text{ hasName } x$
n_{10}	<i>pType</i>	$q(x):- x \text{ rdfs:subClassOf Product}$

edge $n \rightarrow n'$	$\lambda(n \rightarrow n')$	$\delta(n \rightarrow n')$
$n_1 \rightarrow n_1$	<i>acquaintance</i>	$q(x, y):- x \text{ knows } y$
$n_1 \rightarrow n_4$	<i>purchase</i>	$q(x, y):- x \text{ bought } y$
$n_1 \rightarrow n_5$	<i>personAge</i>	$q(x, y):- x \text{ rdf:type Person},$ $x \text{ hasAge } y$
$n_2 \rightarrow n_3$	<i>messgInBlog</i>	$q(x, y):- x \text{ rdf:type Message},$ $x \text{ inBl } y$
$n_4 \rightarrow n_{10}$	<i>productType</i>	$q(x, y):- x \text{ rdf:type Product},$ $x \text{ rdf:type } y$

Table 1: Labels and queries of some nodes and edges of the analytical schema (AnS) shown in Figure 2.

EXAMPLE 3 (ANALYTICAL SCHEMA INSTANCE). Table 2 shows part of the instance of the analytical schema introduced in Example 2. For each triple, we indicate at right the node (or edge) of the AnS which has produced it.

$\mathcal{I}(\mathcal{S}, \mathcal{G}') =$	$\{user_1 \text{ rdf:type person}$	n_1
	$user_2 \text{ rdf:type person}$	n_1
	$user_3 \text{ rdf:type person}$	n_1
	$user_1 \text{ acquaintance user}_2$	$n_1 \rightarrow n_1$
	$user_1 \text{ acquaintance user}_3$	$n_1 \rightarrow n_1$
	$post_1 \text{ messgInBlog blog}_1$	$n_2 \rightarrow n_3$
	$product_1 \text{ rdf:type product}$	n_4
	$user_1 \text{ personAge "28"}$	$n_1 \rightarrow n_5$
	$user_1 \text{ personName "Bill"}$	$n_1 \rightarrow n_6$
	$Notebook \text{ rdf:type pType}$	n_{10}
	$SmartPhone \text{ rdf:type pType}$	n_{10}
	$product_1 \text{ pPrice "\$400"}$	$n_4 \rightarrow n_7 \dots\}$

Table 2: Partial instance of the AnS in Figure 2.

Crucial to our ability to handle RDF heterogeneity is *the disjunctive semantics of an AnS*, materialized by the two levels of \cup operators in Definition 2. Each node and each edge of an AnS populates \mathcal{I} through an independent query, and the resulting triples are simply combined through unions. This has two benefits: (i) significant flexibility when designing the AnS, and (ii) the ability to build a (heterogeneous data) warehouse on top of a heterogeneous RDF graph. Consider for instance the three users in the original graph \mathcal{G} (Figure 1) and their properties: $user_1$, $user_2$ and $user_3$ are part of the person class in our AnS instance \mathcal{I} (through n_1 's query), although $user_2$ and $user_3$ lacks a name. However, those user properties present in the original graph, are reflected by the AnS edges $n_1 \rightarrow n_2$, $n_1 \rightarrow n_8$ etc. Thus, the inherent heterogeneity of RDF graphs is accepted in the input and present in the output of an AnS.

2.2 Analytical queries

Data warehouse analysis summarizes facts according to relevant criteria into *cubes*, which analyze facts characterized by some *dimensions*, using a *measure*. We consider a set of dimensions d_1, d_2, \dots, d_n , such that each dimension d_i may range over the value set $\{d_i^1, \dots, d_i^{n_i}\}$; the Cartesian product of all dimensions $d_1 \times \dots \times d_n$ defines a multidimensional

space \mathcal{M} . To each tuple t in this multidimensional space \mathcal{M} corresponds a *subset* \mathcal{F}_t of the analyzed facts, having for each dimension d_i , $1 \leq i \leq n$, the value of t along d_i .

A *measure* is a set of values¹ characterizing each analyzed fact f . The facts in \mathcal{F}_t are summarized by the *cube cell* $\mathcal{M}[t]$ by the result of an *aggregation* function \oplus (e.g., count, sum, average, etc.) applied to the union of the measures of the \mathcal{F}_t facts: $\mathcal{M}[t] = \oplus(\bigcup_{f \in \mathcal{F}_t} v_f)$.

An *analytical query* consists of two (rooted) queries and an aggregation function. The first query, known as a *classifier* in traditional data warehouse settings, defines the *dimensions* d_1, d_2, \dots, d_n according to which the facts matching the query root will be analyzed. The second query defines the *measure* according to which these facts will be summarized. Finally, the aggregation function is used for *summarizing* the analyzed facts.

DEFINITION 3 (ANALYTICAL QUERY). Given an analytical schema $\mathcal{S} = \langle \mathcal{N}, \mathcal{E}, \lambda, \delta \rangle$, an analytical query (AnQ) rooted in the node $r \in \mathcal{N}$ is a triple

$$Q = \langle c(x, d_1, \dots, d_n), m(x, v), \oplus \rangle$$

where:

- $c(x, d_1, \dots, d_n)$ is a query rooted in the node r_c of its graph G_c . This query is called the classifier of x w.r.t. the n dimensions d_1, \dots, d_n .
- $m(x, v)$ is a query rooted in the node r_m of its graph G_m . This query is called the measure of x .
- \oplus is a function computing a value (a literal) from an input set of values. This function is called the aggregator for the measure of x w.r.t. its classifier.
- For every homomorphism h_c from the classifier to \mathcal{S} and every homomorphism h_m from the measure to \mathcal{S} , $h_c(r_c) = h_m(r_m) = r$ holds.

The last item above guarantees the “well-formedness” of the analytical query, that is: the facts for which we aggregate the measure, are indeed those classified along the desired dimensions. It is worth noticing that, from a practical viewpoint, this condition can be easily and naturally guaranteed by giving explicitly in the classifier and the measure either the type of the facts to analyze, using $x \text{ rdf:type } \lambda(r)$, or a property describing those facts, using $x \text{ } \lambda(r \rightarrow n) \circ$ with $r \rightarrow n \in \mathcal{E}$. As a result, since the labels are unique in an AnS (its labeling function is injective), every homomorphism from the classifier (respectively the measure) to the AnS does map the query’s root node labeled with x to the AnS’s node r .

Definition 3 can be easily extended to more than one measure, by defining a set of measure queries and an associate set of aggregation functions.

EXAMPLE 4 (ANALYTICAL QUERY). The next query asks for the number of blogs where the user posts, classified by the user age and city:

$$\langle c(x, a, c), m(x, y), \text{count} \rangle$$

where the classifier and measure queries are defined by:

$$c(x, a, c):- x \text{ personAge } a, x \text{ personCity } c$$

$$m(x, b):- x \text{ wroteMessg } o, o \text{ messgInBlog } b$$

¹It is a set and not a value, due to the structural heterogeneity of the AnS instance, which is also an RDF graph: each fact may have zero or more values for a given measure.

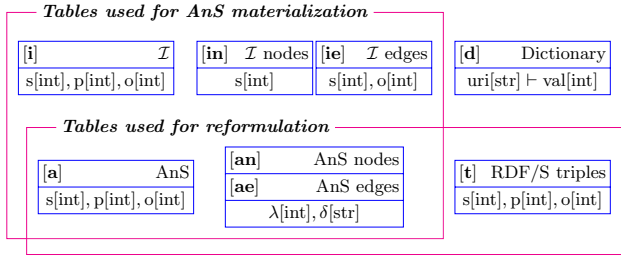


Figure 3: Data layout of the RDF warehouse.

The semantics of an analytical query is:

DEFINITION 4 (ANSWER SET OF AN ANQ). Let \mathcal{I} be the instance of an AnS with respect to some RDF graph. Let $Q = \langle c(x, d_1, \dots, d_n), m(x, v), \oplus \rangle$ be an AnQ against \mathcal{I} . The answer set of Q against \mathcal{I} , denoted $ans(Q, \mathcal{I})$, is:

$$ans(Q, \mathcal{I}) = \{ \langle d_1^j, \dots, d_n^j, \oplus(q^j(\mathcal{I})) \rangle \mid \langle x^j, d_1^j, \dots, d_n^j \rangle \in c(\mathcal{I}) \text{ and } q^j \text{ is defined as } q^j(v) := m(x^j, v) \}$$

assuming that the type of each value returned by $q^j(\mathcal{I})$ belongs (or can be converted by the SPARQL rules [19]) to the input type of the aggregator \oplus . Otherwise, the answer set is undefined.

The analytical query returns each tuple of dimension values from the answer of the classifier, together with the aggregated result of the measure query. The answer set of an AnQ can thus be represented as a cube of n dimensions, holding in each cube cell the corresponding aggregate measure.

EXAMPLE 5 (ANALYTICAL QUERY ANSWER). Consider the query in Example 4, over the analytical schema in Figure 2. Some triples from the instance of this analytical schema were shown in Table 2. The classifier query returns:

$$\langle user_1, 28, \text{"Madrid"} \rangle, \langle user_3, 35, \text{"NY"} \rangle$$

while the measure query returns:

$$\langle user_1, blog_1 \rangle, \langle user_1, blog_2 \rangle, \langle user_3, blog_2 \rangle$$

Aggregating the blogs among the classification dimensions leads to the AnQ answer:

$$\langle 28, \text{"Madrid"}, 2 \rangle, \langle 35, \text{"NY"}, 1 \rangle$$

3. ARCHITECTURE

We deploy our approach using `kdb+` v3.0 (64 bits) [1], an in-memory column database used in decision support analytics applications. `kdb+` provides arrays (tables), which can be manipulated through the `q` interpreted programming language. We store in `kdb+` the RDF graph G , the AnS definitions and AnS instance, when we choose to materialize it. We also translate BGP queries into `q` queries that `kdb+` interprets.

Figure 3 outlines our data layout. As customary in RDF databases, we *dictionary-encode* all URIs into integers; these mappings are stored within the `[d]` table. The graph saturation is stored in `[t]`. Analytical schema *definitions* are stored as follows. The `[a]` table stores the λ function encoded as triples $\lambda(n_1) \lambda(n_1 \rightarrow n_2) \lambda(n_2)$ for all the nodes n_1, n_2 and edges $n_1 \rightarrow n_2 \in \text{AnS}$. The `[an]` table stores the δ function for AnS nodes, and similarly `[ae]` stores the δ function for AnS edges. Finally, we use either the `[i]` table to store the AnS instance \mathcal{I} , or several `[in]` and `[ie]` tables when considering partitioned tables.

`kdb+` stores each column in a table independently; search by value inside a column is fast, supported by an in-memory index, and *self-joins* are also fast. Joins on different columns (e.g., `s=o`) are not advantaged.

The user-facing GUI is implemented in Java 1.6, using Prefuse [9], a Java-based framework for visualizing and interacting with data. The interface retrieves data from `kdb+` and stores it into an internal data structure provided by Prefuse, who then renders the graphs.

4. SCENARIO

The demo will rely on a subset of DBpedia.

Visualizing analytical schemas Figure 4 shows an example of a small analytical schema shown in our GUI. Our visualization tool allows selecting a node, which results in highlighting its connected nodes within the schema. This feature is useful when building an analytical query, since such a query must be a connected component from the schema.

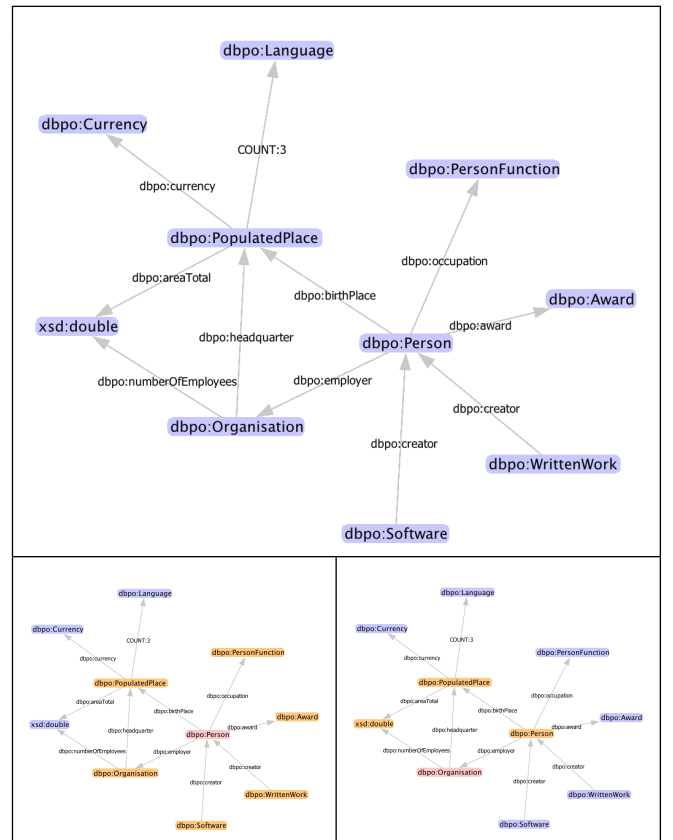


Figure 4: Analytical schema view.

Upon selecting a node additional information regarding the node itself will be displayed. Such information consists of the short (*namespace:term*), the full URI, and the query defining the node. Also if the node has underlying data, this information is displayed too. For instance when selecting the node `dbpo:Person` the user will see the information illustrated in Table 3.

The software is intended to handle large analytical schemas, possibly containing multiple edges between two nodes. In order to increase visibility, our visualization tool groups multiple edges connecting two nodes: a single edge labeled with a value counting all the distinct edges is displayed (see for

Short URI	dbpo:Person
Full URI	http://dbpedia.org/ontology/Person
Query	q(?x):- ?x rdf:type dbpo:Person
Is empty	FALSE
Data:	dbpr : Aristotle dbpr : Neil_Armstrong dbpr : Tim_Berners-Lee ...
Available edges	dbpo:language dbpo:officialLanguage dbpo:regionalLanguage
Short URI	dbpo:language
Full URI	http://dbpedia.org/ontology/language
Query	q(?x,?y):- ?x dbpo:language ?y
Is empty	FALSE
Data:	dbpr : France-dbpr : French_language dbpr : Germany-dbpr : German_language dbpr : United_Kingdom-dbpr : English_language ...

Table 3: Node (top) and edge (bottom) information.

example the edge between dbpo:PopulatedPlace and dbpo:Language in Figure 4). Selecting an edge in a graph leads to obtaining information about the edge itself, and the same holds for edges. When a multiple edge is selected, all the edges are displayed (Table 3).

Evaluating queries Users may build and evaluate analytical queries over a defined schema. Building an analytical query consists of defining each of its three components: classifier query, measure query and aggregation function. The classifier query is built starting from a selected node. The system then suggests (highlights) the nodes connected to it. As more nodes are selected, their neighbors are highlighted as potential options. Only nodes that have underlying data are considered at any time, because the intention is to build queries that can have answers. A selected node is interpreted as a variable in the query. The user may also add constants by selecting values from the node’s underlying data. The process goes on until the classifier is complete and the dimensions (distinguished variables) are selected. The measure query is built similarly. Finally, an aggregation function is chosen from a drop-down list of functions applicable to the chosen measure.

5. RELATED WORKS

Relational DWs have been thoroughly studied [13], and many efficient tools exist. In the context of the Web, “Web Data Warehouses” have been proposed [3, 16], however, their focus is on distribution, data exchange or enrichment, rather than RDF and analytics.

Some works [6, 20] propose new vocabularies for describing relational cubes in RDF; in contrast, WARG is a generic approach for defining and using pure RDF cubes. Data management techniques for RDF have been a hot topic, with recent works tackling storage [2, 17], query processing and updates [15], materialized views [7], and Map-Reduce based RDF processing [10, 11]; Oracle 11g provides a “Semantic Graph” extension etc. Ours is the first work to formalize RDF analytics and propose analytic schema and queries with well-defined semantics.

Recent works [5, 21] have focused on graph warehousing. While we share certain general principles, these approaches are not meant for *heterogeneous* graphs and thus cannot handle multi-valued attributes. More importantly, they do not consider the semantically rich RDF model.

The separation between grouping and aggregation present in our analytical queries is similar to the MD-join operator introduced in [4] for relational DWs.

Finally, the recent SPARQL 1.1 language [19] includes group-by and aggregation constructs closely inspired from SQL. While this bears similarities to our AnQs, our proposal is the first to outline a full-RDF data analytics chain, including a general framework for specifying analytical schemas which are RDF graphs themselves, and therefore carry all the rich features of the data model, most importantly structure heterogeneity and rich semantics.

6. REFERENCES

- [1] [kx] white paper. kx.com/papers/KdbPLUS.Whitepaper-2012-1205.pdf.
- [2] D. J. Abadi, A. Marcus, S. R. Madden, and K. Hollenbach. Scalable semantic web data management using vertical partitioning. In *VLDB*, 2007.
- [3] S. Abiteboul, E. Antoine, and J. Stoyanovich. Viewing the web as a distributed knowledge base. In *ICDE*, 2012.
- [4] M. Akinde, D. Chatziantoniou, T. Johnson, and S. Kim. The MD-join: An operator for complex OLAP. In *ICDE*, 2001.
- [5] D. Bleco and Y. Kotidis. Business intelligence on complex graph data. In *EDBT Workshops*, 2012.
- [6] L. Etcheverry and A. A. Vaisman. Enhancing OLAP analysis with web cubes. In *ESWC*, 2012.
- [7] F. Goasdoue, K. Karanasos, J. Leblay, and I. Manolescu. View selection in Semantic Web databases. *PVLDB*, 5(1), 2012.
- [8] V. Harinarayan, A. Rajaraman, and J. D. Ullman. Implementing data cubes efficiently. In *SIGMOD*, 1996.
- [9] J. Heer, S. K. Card, and J. Landay. Prefuse: A toolkit for interactive information visualization. In *ACM Human Factors in Computing Systems (CHI)*, 2005.
- [10] J. Huang, D. J. Abadi, and K. Ren. Scalable SPARQL Querying of Large RDF Graphs. *PVLDB*, 4(11), 2011.
- [11] M. Husain, J. McGlothlin, M. M. Masud, L. Khan, and B. M. Thuraisingham. Heuristics-Based Query Processing for Large RDF Graphs Using Cloud Computing. *IEEE Trans. on Knowl. and Data Eng.*, 2011.
- [12] M. Jarke, M. Lenzerini, Y. Vassiliou, and P. Vassiliadis. *Fundamentals of Data Warehouses*. Springer, 2001.
- [13] C. S. Jensen, T. B. Pedersen, and C. Thomsen. *Multidimensional Databases and Data Warehousing*. Synthesis Lectures on Data Management. Morgan & Claypool Publishers, 2010.
- [14] R. Kimball and M. Ross. *The Data Warehouse Toolkit*. Wiley, 2002. 2nd Edition.
- [15] T. Neumann and G. Weikum. x-RDF-3X: Fast querying, high update rates, and consistency for RDF databases. *PVLDB*, 3(1), 2010.
- [16] N. Preda, G. Kasneci, F. M. Suchanek, T. Neumann, W. Yuan, and G. Weikum. Active knowledge: dynamically enriching RDF knowledge bases by web services. In *SIGMOD*, 2010.
- [17] L. Sidirourgos, R. Goncalves, M. Kersten, N. Nes, and S. Manegold. Column-store support for RDF data management: not all swans are white. *PVLDB*, 1(2), 2008.
- [18] W3C. Resource description framework. <http://www.w3.org/RDF/>.
- [19] W3C. SPARQL 1.1 query language. <http://www.w3.org/TR/sparql11-query/>.
- [20] W3C. The RDF data cube vocabulary. <http://www.w3.org/TR/vocab-data-cube/>, 2012.
- [21] P. Zhao, X. Li, D. Xin, and J. Han. Graph cube: on warehousing and OLAP multidimensional networks. In *SIGMOD*, 2011.