



Article

# Towards Interactive Analytics over RDF Graphs

Maria-Evangelia Papadaki<sup>1,2</sup> , Nicolas Spyratos<sup>3</sup> and Yannis Tzitzikas<sup>1,2</sup> \*<sup>1</sup> Institute of Computer Science, FORTH-ICS, 70013 Heraklion, Greece<sup>2</sup> Department of Computer Science, University of Crete, 70013 Heraklion, Greece<sup>3</sup> Université de Paris-Sud, Laboratoire de Recherche en Informatique, France

\* Correspondence: tzitzik@ics.forth.gr;

Version January 25, 2021 submitted to Algorithms

**Abstract:** The continuous accumulation of multi-dimensional data and the development of Semantic Web and Linked Data published in RDF bring new requirements for data analytics tools. Such tools should take into account the special features of RDF graphs, exploit the semantics of RDF and support flexible aggregate queries. In this paper, we present an approach for applying analytics to RDF data based on a high-level functional query language, called HIFUN. According to that language, each analytical query is considered as a well-formed expression of a functional algebra and its definition is independent of the nature and structure of the data. In this paper we investigate how HIFUN can be used for easing the formulation of analytic queries over RDF data. We detail the applicability of HIFUN over RDF, as well as the transformations of data that may be required, we introduce the translation rules of HIFUN queries to SPARQL and we describe a first implementation of the proposed model.

**Keywords:** Analytics; RDF; Linked data

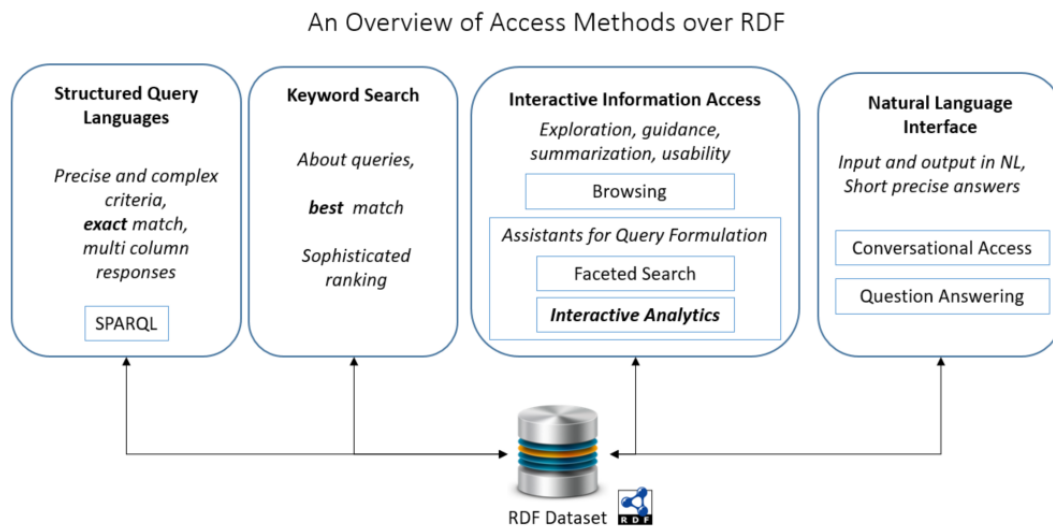
## 1. Introduction

The amount of data available on the Web today is increasing rapidly due to successful initiatives, such as the Linked Open Data movement<sup>1</sup>. More and more data sources are being exported or produced using the Resource Description Framework<sup>2</sup> (or RDF, for short) standardized by the W3C. There are thousands of published RDF datasets (see [1] for a recent survey), including cross-domain knowledge bases (KBs) (e.g., DBpedia [2] and Wikidata [3]), domain specific repositories (e.g., DrugBank [4], GRSF [5], ORKG [6], WarSampo [7], and recently COVID-19 related datasets [8–10] as well as Markup data through schema.org.

Figure 1 shows the general picture of access services over RDF. Apart from Structured Query Languages, we have Keyword Search systems over RDF (like [11]) that allow users to search for information using the familiar method they use for Web searching. We can also identify the category *Interactive Information Access* that refers to access methods that are beyond the simple “query-and-response” interaction, i.e. methods that offer more interaction options to the user and also exploit the *interaction session*. In this category, there are methods for RDF Browsing, methods for Faceted Search over RDF [12], as well as methods for Assistive (SPARQL) Query Building (e.g. [13]). Our work falls in this category, specifically we aim at providing an interactive method for analytics over RDF. Finally, in the category *natural language interfaces* there are methods for question answering, dialogue systems, and conversational interfaces.

<sup>1</sup> <http://lod-cloud.net/>

<sup>2</sup> <https://www.w3.org/RDF/>



**Figure 1.** An Overview of the Access Methods over RDF

31 As regards structured query languages, RDF data are mainly queried through structured query  
 32 languages, i.e. SPARQL<sup>3</sup>, which is the standard query language for RDF data. SPARQL supports  
 33 complex querying using regular path expressions, grouping, aggregation, etc., but the application  
 34 of analytics to RDF data and especially to large RDF graphs is not so straightforward. The structure  
 35 of such graphs tends to be complex due to several factors: (i) different resources may have different  
 36 sets of properties, (ii) properties can be multi-valued (i.e. there can be triples where the subject and  
 37 predicate are the same but the objects are different) and (iii) resources may or may not have types. On  
 38 the other hand, the regular methods of analytics are not capable of analyzing RDF graphs effectively,  
 39 as they (i) focus on relational data, (ii) can only work with a single homogeneous data set, (iii) neither  
 40 support multiple central concepts, nor RDF semantics, (iv) do not offer flexible choices of dimension,  
 41 measure, and aggregation and (v) demand deep knowledge of specific query languages depending on  
 42 data's structure.

43 In view of the above challenges, there is a need for a simple conceptual model able to guide data  
 44 analysis over one or more linked data sets, that demands no programming skill. Motivated from  
 45 this need, we are investigating an approach based on a high-level query language, called HIFUN  
 46 [14], for applying analytics to RDF graphs. We study how that language can be applied to RDF data  
 47 by clarifying how the concept of analysis context can be defined, what kind of transformations are  
 48 required and how HIFUN queries are translated to SPARQL. Note that with the translation approach  
 49 that we focus on, we can apply analytics to RDF sources, without having to transform the RDF data to  
 50 relational ones, nor to copy them.

51 The idea was first introduced in [15]. The current paper is an extended and enriched version of  
 52 that work presenting (a) a more complete related work, (b) a detailed analysis of the applicability of  
 53 HIFUN over RDF data, (c) the detailed algorithm for translating HIFUN queries over RDF data, (d)  
 54 the first implementation of an algorithm that makes that translation.

55 The remainder of the paper is organized as follows: Section 2 discusses the requirements for  
 56 analyzing RDF data and the research that has been conducted in that area. Section 3 introduces the  
 57 related background knowledge. Section 4 focuses on how HIFUN can be used as an interface to  
 58 RDF data. Section 5 investigates whether HIFUN can be applied to RDF data. Section 6 details the  
 59 translation algorithm, Section 7 discusses interactivity issues, and finally section 8 concludes this paper  
 60 and suggests directions for future research.

<sup>3</sup> <https://www.w3.org/TR/rdf-sparql-query/>

## 2. Requirements and Related Work

In this section, we describe the requirements of analyzing semantic warehouses and we survey the related work that has been conducted in the area of RDF analytics.

### 2.1. Requirements

In decision-support systems, in order to extract useful information from the data of an application, it is necessary to analyze large amounts of data accumulated over time - typically over a period of several months. This data is usually stored in a so-called "data warehouse" and analysed along various dimensions and at various levels in each dimension [16]. The end users of such warehouses are mainly analysts and decision-makers, who invariably ask for data aggregations (e.g. total sales by branch).

During last decade the development of Semantic Web data has led to the emergence of semantic warehouses; specific-domain warehouses [17,18] or general-purpose knowledge bases (e.g. DBpedia and WikiData<sup>4</sup>). Thus, it would be useful if the data of these warehouses could be analyzed to extract valuable information (e.g. identify patterns, predict values, discover data correlations), check the quality of semantic integration activities (e.g. for measuring the commonalities between several data sets [19,20]) or monitor the content and the quality of them (e.g. by evaluating the completeness or the representativeness of its data etc.).

However, the analysis of such warehouses introduces several challenges [21]. The data heterogeneity, its lack of a strict structure, its rich semantics and the possibility of incomplete data sources significantly complicates their analysis. For example, although one can reasonably expect to find information about a particular concept, they can not always find specific information for all the instances of it (e.g. the opening hours or closing days of all branch stores). Moreover, data warehouses follow a star schema and thus, the facts can be analyzed based on certain dimensions and measures, the choice of which is made at the data warehouse design time (e.g. if "branch" and "product" have been defined as dimensions, then aggregations over them are not allowed; one cannot find "the number of branches established in 2020", since "branch" is a dimension and relational data cubes do not allow aggregating over dimensions). In addition, different concepts (e.g. "branches", "products", "people") can be analyzed, only if each of them is modeled by a different schema and stored in a distinct data warehouse. Finally, even though such warehouses host data published along with a schema (which can facilitate the understanding of data), the structure of it tends to be complex. Note also, that the end-users, who are usually non-specialists, are unable to read the schema and formulate the queries necessary for their data analysis. Thus, it would be useful if apart from native RDF data, one could analyze and deduce further knowledge (inference) from RDF schemas, too (e.g. ask for all the relationships linking products to other entities).

Therefore, there is a need for applying analytics to any kind of RDF graph, - not only to multidimensional data expressed in RDF, but also to domain-specific or general-purpose semantic data; a way, that will be applicable to several RDF data sets, as well as to any data source. In general, we need an analytical tool that will allow the user to select the desired data set(s) or desired parts thereof, define the features (s)he is interested in at the query time, formulate an analytic query without having any programming background knowledge and will display the results in the form of tables, plots or any other kind of visualization for intuitive exploration.

### 2.2. Related Work

Statistical data is published as linked data in order to be combined with data sets that are published in disparate sources on the Web. Such data should have been modeled as data cubes describing data

---

<sup>4</sup> <https://www.wikidata.org>

104 in a multi-dimensional fashion. Towards this end, the RDF data cube vocabulary<sup>5</sup> (QB) is employed.  
105 This vocabulary provides a means to publish such data on the web using the W3C RDF standard. It  
106 consists of three main components: i) the *measures*, which are the observed values of primary interest,  
107 ii) the *dimensions*, which are the value keys that identify the measure and iii) the *attributes*, which are  
108 the metadata. However, even though that vocabulary can be used for structuring and publishing  
109 multi-dimensional data, it cannot be used for applying analytics over it. In view of this limitation,  
110 several approaches have been proposed.

111 These approaches can be divided into two major groups: (i) those assuming that the  
112 multidimensional data (MD) i.e. data related to more than two dimensions, has already been  
113 represented in the RDF format and (ii) those that do not. Our approach, as well as the works in [22–24]  
114 are related to the first group. On the other hand, the work in [25] considers that the multidimensional  
115 data has been stored as non-RDF data sets. In particular, it declares that the data cubes are retrieved  
116 from a relational database with SQL queries and then get triplicated.

117 The representation of MD data in RDF can further be organized in two categories: i) those that  
118 are based on specialized RDF vocabularies [23,26] and ii) those that implicitly define a data cube over  
119 existing RDF graphs<sup>6</sup> [24,27,28]. Even though the second category is promising, it cannot guarantee  
120 that the cubes on RDF graphs will be multi-dimensional compliant [23]. Additionally, to the best of  
121 our knowledge, the existing approaches support only homogeneous graphs [27], and thus they cannot  
122 handle any multi-valued attributes (e.g., a person being both “Greek” and “French”), nor semantics.

123 The existing methods can also be classified into i) those that require programming knowledge  
124 for analyzing the data and ii) those that do not deal with lower-level technicalities. The work in [29]  
125 presents a system for analytics over (large) graphs. It achieves efficient query answering by dividing  
126 the graph into partitions. However, in contrast to our work, the user should have some programming  
127 knowledge, since it is necessary to write a few lines of code to submit the query. The work in [30]  
128 presents a method for applying statistical calculations on numerical linked data. It stores the data  
129 in arrays and performs the calculations on the arrays’ values. Nevertheless, contrary to our work, it  
130 requires deep knowledge of SPARQL for formulating the queries.

131 In order to overcome one’s difficulty in background programming knowledge, high-level  
132 languages have been developed for data analysis, too. However, there has not been much activity in  
133 introducing high-level languages suitable for analytics on RDF data. While general-purpose languages,  
134 such as PIG Latin [31] and HiveQL [32] can be used, they are not tailored to address the peculiarities  
135 of the RDF data model. Even though, [33,34] present high-level query languages enabling OLAP  
136 querying of an extended format of data cubes [23], they are only applicable to data already represented  
137 and published using a corresponding vocabulary. As a consequence, they fall short in addressing a  
138 wide variety of analytical possibilities in non-statistical RDF data sources. In addition, [31] proposes  
139 a high-level language that supports semantics. However, it is targeted at processing structured relational  
140 data, limiting its use for semi-structured data such as RDF. Further, it provides only a finite set of  
141 primitives that is inadequate for the efficient expression of complex analytical queries.

142 A survey that is worth mentioning is [35], which introduces warehouse-style RDF analytics. There  
143 are similarities with our approach, since each analytical schema node corresponds to an RDF class,  
144 while each edge corresponds to an RDF property. Nonetheless, since the facts are encoded as unary  
145 patterns, they are limited to vertices instead of arbitrary sub-graphs (e.g. paths). Other related work  
146 includes [24] that focuses on how to reuse the materialized result of a given RDF analytical query  
147 (cube) in order to compute the answer to another cube, as well as recent systems for analytics over  
148 RDF like Spade [36] that suggests to users aggregates that are visually interesting.

---

<sup>5</sup> <https://www.w3.org/TR/vocab-data-cube/>

<sup>6</sup> <https://team.inria.fr/oak/projects/warg/>

149 In brief, in contrast to the aforementioned works, in this paper we focus on developing a  
150 user-friendly interface, where the user will be able to apply analytics to RDF data without dealing with  
151 lower-level technicalities of SPARQL. Indeed HIFUN is more simple for formulating analytic queries.  
152 We focus on the support of analytics over any RDF Data (not only over data expressed according  
153 to RDF Data Cube), and we focus on a query translation approach, i.e. an approach that does not  
154 require transforming or transferring the existing data; instead it can be directly applied over a SPARQL  
155 endpoint. Furthermore, the query translation approach allows exploiting the RDF Schema semantics  
156 that is supported by SPARQL, i.e. the inferred RDF triples are taken into account in the evaluation of  
157 the analytic queries.

### 158 3. Background

#### 159 3.1. Principles of Resource Description Framework (RDF)

160 **Resource Description Framework (RDF)** The Resource Description Framework (RDF) [37,38] is  
161 a graph-based data model for linked data interchanging on the web. It uses triples i.e. statements of the  
162 form *subject – predicate – object*, where the *subject* corresponds to an entity (e.g. a branch, a product,  
163 etc.), the *predicate* to a characteristic of the entity (e.g. name of branch) and the *object* to the value of  
164 the predicate for the specific subject (e.g. “*branch<sub>1</sub>*”). The triples are used for relating Uniform Resource  
165 Identifiers (URIs) or anonymous resources (blank nodes) with other URIs, blank nodes or constants  
166 (Literals). Formally, a *triple* is considered to be any element of  $T = (U \cup B) \times (U) \times (U \cup B \cup L)$ , where  
167  $U, B$  and  $L$  denote the sets of URIs, blank nodes and literals, respectively. Any finite subset of  $T$   
168 constitute an *RDF graph* (or *RDF data set*).

169 **RDF Schema.** RDF Schema<sup>7</sup> (RDFS) is a special vocabulary which comprises a set of classes with  
170 certain properties using the RDF extensible knowledge representation data model. Its intention is to  
171 structure RDF resources, since even though RDF uses URIs to uniquely identify resources, it lacks  
172 semantic expressiveness. It uses classes to indicate where a resource belongs, as well as properties to  
173 build relationships between the entities of a class and to model constraints. A class  $C$  is defined by a  
174 triple of the form  $\langle C \text{ rdf:type rdfs:Class} \rangle$  using the predefined class “*rdfs:Class*” and the predefined  
175 property “*rdf:type*”. For example, the triple  $\langle \text{ex:Product rdf:type rdfs:Class} \rangle$  indicates that “*Product*”  
176 is a class, while the triple  $\langle \text{ex:product1 rdf:type ex:Product} \rangle$  that individual “*product1*” is an instance  
177 of class *Product*. A property can be defined by stating that it is an instance of the predefined class  
178 “*rdf:Property*”. Optionally, properties can be declared to apply to certain instances of classes by  
179 defining their domain and range using the predicates “*rdfs:domain*” and “*rdfs:range*”, respectively. For  
180 example, the triples  $\langle \text{ex:hasProduct rdf:type rdf:Property} \rangle$ ,  $\langle \text{ex:hasProduct rdfs:domain ex:Branch} \rangle$ ,  
181  $\langle \text{ex:hasProduct rdfs:range ex:Product} \rangle$ , indicate that the domain of the property “*hasProduct*” is  
182 the class “*Branch*” and its range the class “*Product*”. RDFS is also used for defining hierarchical  
183 relationships among classes and properties. The predefined property “*rdfs:subclassOf*” is used as a  
184 predicate in a statement to declare that a class is a specialization of another more general class, while  
185 the specialization relationship between two properties is described using the predefined property  
186 “*rdfs:subPropertyOf*”. For example, the triple  $\langle \text{ex:Branch rdfs:subClassOf ex:Store} \rangle$  denotes that  
187 the class “*Branch*” is sub-class of “*Store*”, while the triple  $\langle \text{ex:hasDrinkProduct rdf:subPropertyOf} \text{ex:hasProduct} \rangle$   
188 that the property “*hasDrinkProduct*” is sub-property of “*hasProduct*”. Moreover,  
189 RDFS offers inference functionality<sup>8</sup> as additional information (i.e. discovery of new relationships  
190 between resources) about the data it receives. For example, if  $\langle \text{ex:Coca-Cola rdf:type ex:Drink} \rangle$  and  
191  $\langle \text{ex:Drink rdf:type ex:Product} \rangle$ , then it can be deduced that “*ex:Coca-Cola rdf:type ex:Product*”.

---

<sup>7</sup> [https://en.wikipedia.org/wiki/RDF\\_Schema](https://en.wikipedia.org/wiki/RDF_Schema)

<sup>8</sup> <https://www.w3.org/standards/semanticweb/inference>

192 We shall use the example of Fig. 2 as our running example throughout the paper. It is an RDF  
 193 Graph containing information about invoices and related information about them. Each invoice has a  
 194 URI, e.g. the invoice with URI `ex:ID4`. That invoice participates to the following five triples:  
 195 `ex:ID4 rdf:type ex:Invoice` .  
 196 `ex:ID4 ex:hasDate "2019-05-09"` .  
 197 `ex:ID4 ex:takesPlaceAt ex:branch3` .  
 198 `ex:ID4 ex:delivers ex:product4` .  
 199 `ex:ID4 ex:inQuantity "400"` .  
 200 meaning that the type of “`ex:ID4`” is Invoice, it took place in “2019-05-09” at “`branch3`”, and delivered  
 201 400 items of `ex:product4`. Since data is in RDF each product has a URI and in this particular example  
 202 we can see that the brand of “`product4`” is “Hermes” and that the founder of that brand is “Manousos”,  
 203 who is both Greek and French.

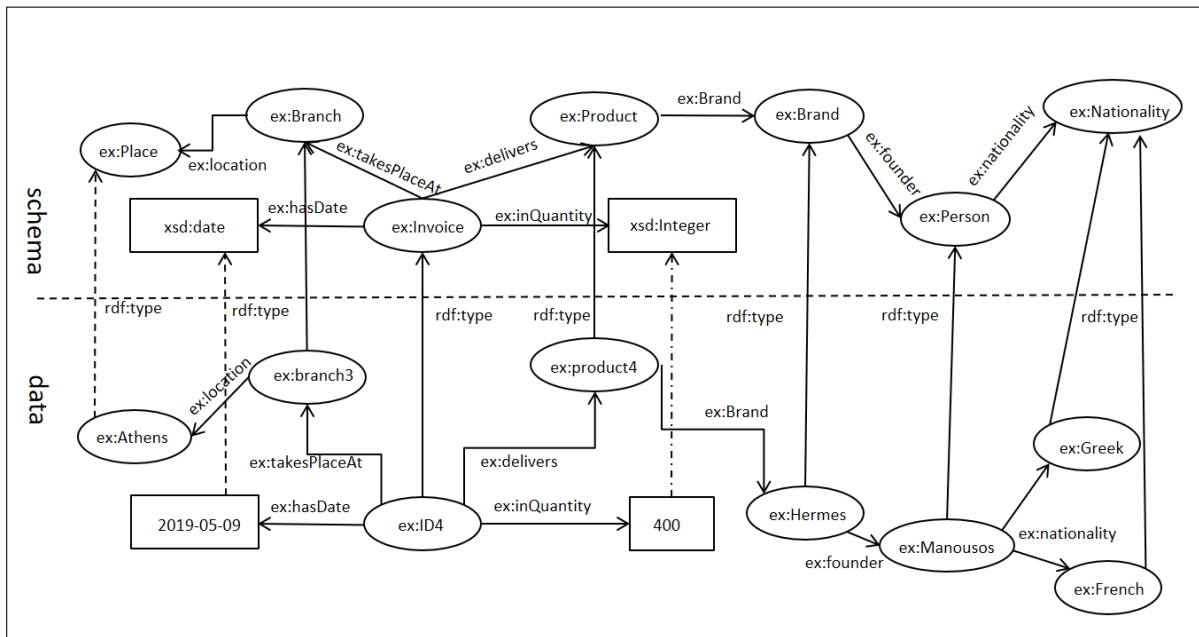


Figure 2. Running example

### 204 3.2. HIFUN - A High Level Functional Query Language for Big Data Analytics

205 HIFUN [14] is a high-level functional query language for defining analytic queries over big data  
 206 sets, independently of how these queries are evaluated. It can be applied over a data set that is  
 207 structured or unstructured, homogeneous or heterogeneous, centrally stored or distributed.

208 *Data set Assumptions.* To apply that language over a data set  $D$ , two assumptions should hold. The  
 209 data set should i) consist of *uniquely identified data items*, and ii) have a set of *attributes* each of which is  
 210 viewed as a *function* associating each data item of  $D$  with a value, in some set of values. For example, if  
 211 the data set  $D$  is a set of all delivery invoices over a year in a distribution center (e.g. Walmart) which  
 212 delivers products of various types in several branches, then the attribute “product type” (denoted  
 213 as  $pt$ ) is seen as a function  $pt : D \rightarrow String$  such that, for each invoice  $i$ ,  $pt(i)$  is the type of product  
 214 delivered according to the invoice  $i$ .

215 **Definition 1** (Analysis Context). Let  $D$  be a data set and  $A$  be the set of all attributes ( $a_1, \dots, a_k$ ) of  $D$ .  
 216 An *analysis context* over  $D$  is any set of attributes from  $A$ , and  $D$  is considered the *origin* (or *root*) of  
 217 that context. □

218 Roughly speaking, an analysis context is an acyclic directed labeled graph whose nodes and  
 219 arrows satisfy the following conditions:

- 220 1. one or more roots (i.e. nodes with no entering arrows representing the objects of an application)  
 221 may exist  
 222 2. at least one path from a root to every other node (i.e. attributes of the objects) exists  
 223 3. all arrow labels are distinct  
 224 4. each node is associated with a non-empty set of values

225 The number of roots of an analysis context indicates the number of data sets it is related to. While one  
 226 root means that data analysis concerns a single data set, the existence of two or more roots means that  
 227 data analysis relates to two or more different data sets, possibly sharing one or more attributes.

228 Fig. 3 shows our running example, expressed as a context. From a syntactic point of view, the  
 229 edges of it can be seen as triples of the form (source, label, target).

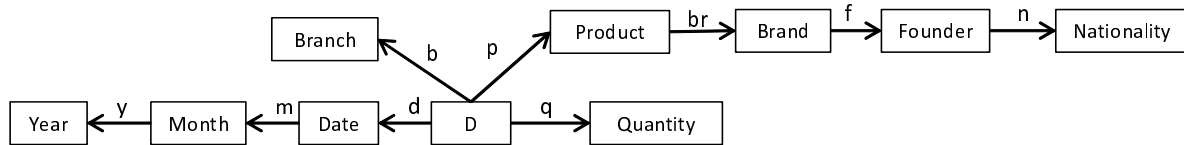


Figure 3. Running example expressed as a HIFUN context

230

231 *Direct & Derived Attributes.* The attributes of a context are divided into two groups, the *direct* and the  
 232 *derived*. The first group contains the attributes with origin  $D$ : these are the attributes whose values are  
 233 given. The second group contains the attributes whose origins are different than  $D$  and whose values  
 234 are computed based on the values of the direct attributes. For example, in Fig. 3 the attributes  $d$ ,  $b$ ,  $p$   
 235 and  $q$  are direct as their values appear on the delivery invoice  $D$ , whereas  $m$  and  $y$  are derived, since  
 236 their values can be computed from those of the attribute  $d$  (e.g. from the date 26/06/2019 one can  
 237 derive the month 06 and the year 2019).

238 **Definition 2** (HIFUN Analytic Query). A query in HIFUN is defined as an ordered triple  $Q = (g, m, op)$   
 239 such that  $g$  and  $m$  are attributes of the data set  $D$  having a common source and  $op$  is an aggregate  
 240 operation (or *reduction operation*) applicable on  $m$ -values. The first component of the triple is called  
 241 *grouping function*, the second *measuring function* (or the *measure*) and the third *aggregate operation* (or  
 242 *reduction operation*). □

243 Roughly speaking, an analytical query  $Q$  is a path expression over an analysis context  $C$ ; a well  
 244 formed expression whose operands are arrows from  $C$  and whose operators are those of the functional  
 245 algebra. It is formulated using paths starting at the root and is evaluated in a three-step process, as  
 246 follows: i) items with the same  $g$ -value  $g_i$  are grouped, ii) in each group of items created, the  $m$ -value  
 247 of each item in the group is extracted from  $D$  and iii) the  $m$ -values obtained in each group are aggregated  
 248 to obtain a single value  $v_i$ . Actually, the aggregate value  $v_i$  is the answer of  $Q$  on  $g_i$ . This means that a  
 249 query is a triple of functions and its answer  $AnsQ$  is a function, too.

#### 250 4. Using HIFUN as an Interface to RDF Dataset

251 There are several ways in which HIFUN can be used, such as for studying rewriting of analytic  
 252 queries in the abstract [14] or for defining an approach to data exploration [39]. In this paper, we use  
 253 HIFUN as a user-friendly interface for defining analytic queries over RDF data sets. To understand the  
 254 proposed approach, consider a data source  $S$  with query language  $L$  (e.g.  $S$  could be a relational data  
 255 set and  $L$  the SQL language). In order to use HIFUN as a user interface for  $S$ , we need to (a) define an  
 256 analysis context, that is a subset  $D$  of  $S$  to be analyzed and some attributes of  $D$  that are relevant for  
 257 the analysis and (b) define a mapping of HIFUN queries to queries in  $L$ .

258 Defining a subset  $D$  of  $S$  can be done using a query of  $L$  and defining  $D$  to be its answer (i.e.  $D$  is  
 259 defined as a view of  $S$ ); and similarly, the attributes that are relevant to the analysis can be defined

260 based on attributes of  $D$  already present in  $S$ . However, defining a mapping of HIFUN queries to  
261 queries in  $L$  might be a tedious task. In [39] such mappings have been defined from HIFUN queries to  
262 SQL queries and from HIFUN queries to MapReduce jobs.

263 The main objective of this paper is to define a user-friendly interface allowing users to perform  
264 analysis of RDF data sets. To this end, we use the HIFUN language as the interface. In other words, we  
265 consider the case where the data set  $S$  mentioned above is a set of RDF triples and its language  $L$  is the  
266 SPARQL language. Our main contributions are: (a) the proposal of tools for defining a HIFUN context  
267 from the RDF data set  $S$  and (b) defining a mapping from HIFUN queries to SPARQL queries. With  
268 these tools at hand, a user of the HIFUN interface can define an analysis context of interest over  $S$  and  
269 issue analytic queries using the HIFUN language. Each such query is then translated by the interface  
270 to a SPARQL query, which in turn is evaluated over the RDF triples of  $D$  and the answer is returned to  
271 the user.

## 272 5. Applicability of HIFUN over RDF

273 In Section 5.1 we discuss the prerequisites for applying HIFUN over RDF, and then (in Section  
274 5.2) we describe two methods for applying HIFUN over RDF: over the original data (in Section 5.3),  
275 and after transforming the original data (in Section 5.4).

### 276 5.1. Prerequisites for Applying HIFUN over RDF Data

277 Two assumptions should hold to apply HIFUN over a data set  $D$ , (i) the unique identification of  
278 its data items and (ii) the functionality of its attributes.

279 **RDF data.** The first assumption, the unique identification of the data items, is satisfied by the  
280 RDF data, since each resource is identified by a distinct URI. Consequently,  $D$  can be any subset of  
281 the set of all the available URIs. The second assumption, the functionality of attributes, is partially  
282 satisfied by the RDF properties. The functional (i.e. `owl:FunctionalProperty`) or the effectively functional  
283 properties (i.e. even if they are not declared as functional, they are single-valued for the resources in  
284 the data set  $D$ ) have only one unique value for each instance. However, there are also properties in  
285 RDF with (i) no or (ii) multiple values; a non-value property implies that a value may not exist (or it is  
286 unknown even if it exists) or it is incomplete, while a multi-valued infers that a property has more than  
287 one values for the same resource. Such cases require transforming the original data before applying  
288 HIFUN to it. These transformations can be made using the operators that will be described in §5.4.

289 **RDF Schema.** Each resource of an RDF schema is identified by a distinct URI; therefore, its data  
290 items are uniquely identified. However, a property (e.g. `rdf:type`, `rdfs:subClassOf` etc.) may appear  
291 more than once by relating different classes or classifying concepts in more than one classes (i.e. a  
292 class might be sub-class of several super-classes). Nevertheless, these relationships are considered  
293 distinct since they have different domain and/or range. Therefore, HIFUN allows analytics not only  
294 over the data, but over RDF schema(s) as well. Inference is supported, too, since it refers to automatic  
295 procedures that generate new relationships based on a set of rules; a process that is independent of  
296 HIFUN.

### 297 5.2. Methods to Apply HIFUN over RDF

298 We can identify two main methods for applying HIFUN over RDF:

- 299 I: Defining an Analysis Context over the Original RDF Data. Here the user selects some properties,  
300 all satisfying the aforementioned assumptions. This is discussed in Section 5.3.
- 301 II: Defining an Analysis Context after Transforming the Original RDF Data. Here the user transforms  
302 parts of the RDF graph in a way that satisfies the aforementioned assumptions. This is discussed  
303 in Section 5.4.



304 **5.3. I. Defining an Analysis Context over the Original RDF Data**

305 **Definition 3** (Analysis Context). An analysis context  $C$  over RDF data is defined as a set of resources  
306  $R$  to be analyzed along with a set of properties  $p_1, p_2, \dots, p_n$  that are relevant for the analysis.  $\square$

307 As the *root* of an analysis context in RDF can be selected any class (i.e. set of resources) of an RDF  
308 graph and as attributes any properties of that graph. For example, any of the classes “ex:Invoice”,  
309 “ex:Branch”, “ex:Product”, “ex:Brand”, “ex:Person”, “ex:Nationality” of Fig. 2 can be selected as  
310 the root of the context, while any of the properties “ex:hasDate”, “ex:takesPlaceAt”, “ex:delivers”,  
311 “ex:inQuantity”, “ex:Brand”, “ex:founder”, “ex:nationality” as the attributes of it.

312 **5.4. II: Defining an Analysis Context after Transforming the Original RDF Data**

313 A few feature operators that could be used for transforming the original (RDF) data to be in  
314 compliance with the assumptions of HIFUN are indicated in Table 1. That table lists the nine most  
315 frequent *Linked Data-based Feature Creation Operators* (for short *FCOs*), as defined in [40], and they have  
316 been re-grouped according to our requirements.  $\mathcal{T}$  denotes a set of triples,  $P$  a set of properties and  
317  $p, p_1, p_2$  properties. In detail,

- 318 •  $fc_{01}$  suits to the normal case and it can be exploited to confirm that all the properties are  
319 functional e.g. the date that each product was delivered, the branch where each invoice took  
320 place. The value can be numerical or categorical.
- 321 •  $fc_{02}$  and  $fc_{03}$  relate to issues that concern missing and multi-valued properties and can be used  
322 for turning properties with empty values into integers.
- 323 •  $fc_{04}$  can be used for converting a multi-valued property to a set of single-valued features, e.g.  
324 one boolean feature for each nationality, that a founder may have.
- 325 •  $fc_{05}$  and  $fc_{06}$  concern the degree of an entity and can be used to find the set of triples that  
326 contains a specific entity, defining its importance.
- 327 •  $fc_{07}$  to  $fc_{09}$  investigate paths in an RDF graph, e.g. whether at least one founder of a brand is  
328 “French”. It can be used for specifying a path (i.e. a sequence of properties  $p_1, p_2, \dots, p_n$  etc.) and  
329 treat it as an individual property  $p$ .

**Table 1.** Feature Creation Operators

id	Operator defining $f_i$	Type	$f_i(e)$
<b>Plain selection of one property</b>			
1	p.value	num/categ	$f_i(e) = \{v \mid (e, p, v) \in \mathcal{T}\}$
<b>For missing values and multi-valued properties</b>			
2	p.exists	boolean	$f_i(e) = 1$ if $(e, p, o)$ or $(o, p, e) \in \mathcal{T}$ , otherwise $f_i(e) = 0$
3	p.count	int	$f_i(e) =  \{v \mid (e, p, v) \in \mathcal{T}\} $
<b>For multi-valued properties</b>			
4	p.values.AsFeatures	boolean	for each $v \in \{v \mid (e, p, v) \in \mathcal{T}\}$ we get the feature $f_{iv}(e) = 1$ if $(e, p, v)$ or $(v, p, e) \in \mathcal{T}$ , otherwise $f_{iv}(e) = 0$
<b>General ones</b>			
5	degree	double	$f_i(e) =  \{(s, p, o) \in \mathcal{T} \mid s = e \text{ or } o = e\} $
6	average degree	double	$f_i(e) = \frac{ \text{triples}(C) }{ C }$ s.t. $C = \{c \mid (e, p, c) \in \mathcal{T}\}$ and $\text{triples}(C) = \{(s, p, o) \in \mathcal{T} \mid s \in C \text{ or } o \in C\}$
<b>Indicative extensions for paths</b>			
7	p1.p2.exists	boolean	$f_i(e) = 1$ if $\exists o2$ s.t. $\{(e, p1, o1), (o1, p2, o2)\} \subseteq \mathcal{T}$
8	p1.p2.count	int	$f_i(e) =  \{o2 \mid (e, p1, o1), (o1, p2, o2) \in \mathcal{T}\} $
9	p1.p2.value.maxFreq	num/categ	$f_i(e) = \text{most frequent } o2 \text{ in } \{o2 \mid (e, p1, o1), (o1, p2, o2) \in \mathcal{T}\}$

330 These features can be used for deriving a new RDF dataset that will be analyzed with HIFUN.  
331 This transformation can be done by using SPARQL CONSTRUCT queries: Suppose that the pair  $(R, F)$

332 expresses a context, where  $R$  is a set of resources and  $F$  the set of the features, the objects in  $R$  have.  
 333 Then, the resources  $R$ , as well as the features  $F$ , can be defined by the triple patterns i.e. “ $?s ?p ?o$ ”  
 334 in the *CONSTRUCT* clause of a SPARQL query, i.e. the bindings of “ $s$ ” (or “ $o$ ”) can correspond to  
 335 the resources, whereas the bindings of “ $p$ ” to the set of features. Alternatively, these features can be  
 336 defined by queries, but instead of constructing the triples, the definition of the features can be included  
 337 in the analytic queries in the form of nested queries (subqueries); in general any query translation  
 338 method for virtual integration [41] can be used. A concrete example will be given in Section 6.6.

339 Finally we should mention that the above list is by no means complete; the list of feature operators  
 340 can be expanded to cover the requirements that arise.

## 341 6. Translation of HIFUN Queries to SPARQL

342 Here we focus on how to translate a HIFUN query, over an analysis context over RDF (case I), to  
 343 a SPARQL query. Roughly, the grouping function will eventually yield variable(s) in the GROUP BY  
 344 clause, the measuring function will yield at least one variable in the WHERE clause, and the aggregate  
 345 operation corresponds to the appropriate aggregate SPARQL function in the SELECT clause (over the  
 346 measuring variable). We explain the translation method gradually using examples, assuming the  
 347 running example of Figure 2.

### 348 6.1. Simple Queries

349 Suppose that we would like to find the total quantities of products delivered to each branch.  
 350 This query would be expressed in HIFUN as  $(takesPlaceAt, inQuantity, SUM)$  and in SPARQL as (for  
 351 reasons of brevity we assume the namespace prefix “ $ex$ ” for each property of the following queries):

```
352
353 SELECT ?x2 SUM(?x3)
354 WHERE {
355   ?x1 ex:takesPlaceAt ?x2 .
356   ?x1 ex:inQuantity ?x3 .
357 }
358 GROUP BY ?x2
```

359 Therefore a HIFUN query  $(g, m, op)$  is translated to SPARQL as follows: the function  $g$  is translated  
 360 to a triple pattern  $?x1 \ g \ ?x2$  (in the WHERE clause) and the variable  $?x2$  is added to the SELECT  
 361 clause, and in the GROUP BY clause. The function  $m$  is translated to a triple pattern  $?x1 \ m \ ?x_N$  in the  
 362 WHERE clause, where  $x_N$  denotes a new variable. Finally the function  $op$  is translated to a  $op(right(m))$   
 363 in the SELECT clause, where  $right(m)$  refers to the “right” variable of the triple pattern derived by the  
 364 translation of  $m$ , i.e. to  $?x_N$ , in our example  $SUM(?x3)$ .  
 365

### 366 6.2. Attribute-Restricted Queries.

367 Suppose that we would like to find the total quantities of products, delivered to one particular  
 368 branch, say `branch1`. This query would be expressed in HIFUN as  
 369  $(takesPlaceAt/branch1, inQuantity, SUM)$  and in SPARQL as:

```
370
371 SELECT ?x2 SUM(?x3)
372 WHERE {
373   ?x1 ex:takesPlaceAt ?x2 .
374   ?x1 ex:inQuantity ?x3 .
375   ?x1 ex:takesPlaceAt branch1 .
376 }
377 GROUP BY ?x2
378
```

Therefore for translating the HIFUN query  $(g/v, m, op)$  we translate the restriction  $v$  by adding in the WHERE clause the triple pattern  $?x1 g v$ . Note that, here the restriction value refers to a URI. If that value had been represented with a literal, then a FILTER statement would have to be added in the WHERE clause. For example, consider the following example (where in this case the restriction is applied to the measuring function): Suppose that we would like to find the total quantities of products, delivered to each branch by considering only those invoices with quantity greater than or equal to 1. This query would be expressed in HIFUN as  $(takesPlaceAt, inQuantity/_{\geq 1}, SUM)$  and in SPARQL as:

```

387
388 SELECT ?x2 SUM(?x3)
389 WHERE {
390 ?x1 ex:takesPlaceAt ?x2 .
391 ?x1 ex:inQuantity ?x3 .
392 FILTER(?x3 ≥ xsd:integer("1")) .
393 }
394 GROUP BY ?x2
395

```

Consequently, a literal-attribute restriction in a HIFUN query  $(g, m/cond, op)$  would be translated by adding in the WHERE clause the following constraint  $FILTER(right(m) cond)$ .

### 6.3. Results-Restricted Queries.

Suppose that we would like to find the total quantities of products, delivered to each branch, but only for branches with total quantity greater than 1000. This query would be expressed in HIFUN as,  $(takesPlaceAt, inQuantity, SUM/_{>1000})$  and in SPARQL as:

```

402
403 SELECT ?x2 SUM(?x3)
404 WHERE {
405 ?x1 ex:takesPlaceAt ?x2 .
406 ?x1 ex:inQuantity ?x3 .
407 }
408 GROUP BY ?x2
409 HAVING (SUM(?x3) > 1000)

```

Therefore for translating the HIFUN query  $(g, m, op/cond)$  we translate the restriction  $cond$  by adding a HAVING clause with the following constraint  $HAVING Right(m) cond$ .

### 6.4. Complex Grouping Queries.

A grouping (as well as a measuring) function in HIFUN can be *more complex* using the following operations on functions, as defined in [14]: *composition* ( $\circ$ ) and *pairing* ( $\otimes$ ). These operations form the so called functional algebra [16] and they are well known, elementary operations.

#### 6.4.1. Composition

Suppose that we ask for the total quantities of products delivered by brand. This query would be expressed in HIFUN as  $(brand \circ delivers, inQuantity, SUM)$ , and in SPARQL as:

```

419
420 SELECT ?x3 SUM(?x4)
421 WHERE {
422 ?x1 ex:delivers ?x2 .
423 ?x2 ex:brand ?x3 .
424 ?x1 ex:inQuantity ?x4 .

```

```

425 }
426 GROUP BY ?x3
427
428

```

Therefore, a HIFUN query  $(f_k \circ \dots \circ f_2 \circ f_1, m, op)$  would be translated as follows. At first note that if instead of the composition we had one function  $f_1$ , then it would be interpreted as a single query, i.e. we would add the triple pattern  $?x1 f_1 ?x2$  to the WHERE clause and the variable  $right(f_1)$  would be added to the SELECT and to the GROUP BY clause. If we had the composition of two functions  $(f_2 \circ f_1)$ , then we would add the triple patterns  $?x1 f_1 right(f_1)$  and  $right(f_1) f_2 ?xf2r$  (where  $xf2r$  is a brand new variable) to the WHERE clause and the variable  $right(f_2)$  to the SELECT and to the GROUP BY clause.

Now suppose that the composition function comprises  $k$  functions,  $(f_k \circ \dots \circ f_2 \circ f_1)$ , which would be translated to the triple patterns  $?x1 f_1 right(f_1), right(f_1) f_2 right(f_2), \dots, right(f_{k-1}) f_k right(f_k)$  to the WHERE clause and the variable  $right(f_k)$  to the SELECT and to the GROUP BY clause. If we added one more function to the composition (reaching to  $k+1$  functions), i.e.  $(f_{k+1} \circ f_k \circ \dots \circ f_2 \circ f_1)$ , we would have to add the triple pattern  $right(f_k) f_{k+1} ?xnew$  (where  $?xnew$  is a brand new variable) to the WHERE clause and to replace the variable  $right(f_k)$  with the  $right(f_{k+1})$  in the SELECT and in the GROUP BY clause.

Now we shall provide an example of composition with a *derived* attribute. Suppose that we ask for the total quantities of products delivered by month. This query would be expressed in HIFUN as  $(month \circ date, inQuantity, SUM)$  and in SPARQL as:

```

446
447
448 SELECT month(?x2) SUM(?x3)
449 WHERE {
450 ?x1 ex:hasDate ?x2 .
451 ?x1 ex:inQuantity ?x3 .
452 }
453 GROUP BY month(?x2)
454
455

```

Therefore, a HIFUN query  $(f \circ g, m, op)$ , where the attribute  $f$  derives from  $g$ , would be translated by adding to the WHERE clause the triple pattern  $?x1 g ?xf2r$  (where  $?xf2r$  is a brand new variable). Then,  $f$  would be derived from  $right(g)$  by adding to the SELECT and to the GROUP BY clauses a SPARQL build-in function i.e.  $f(right(g))$  (in our example  $month(?x2)$ ); this function would extract the value  $f$  from that of  $right(g)$ .

#### 6.4.2. Pairing.

Suppose that we would like to find the total quantities delivered by branch and product. This query would be expressed in HIFUN as  $((takesPlaceAt \otimes delivers), inQuantity, SUM)$  and in SPARQL as:

```

465
466 SELECT ?x2 ?x4 SUM(?x3)
467 WHERE {
468 ?x1 ex:takesPlaceAt ?x2 .
469 ?x1 ex:inQuantity ?x3 .
470 ?x1 ex:delivers ?x4.
471 }
472 GROUP BY ?x2 ?x4

```

473

474 Therefore, a HIFUN query  $(f_k \otimes \dots, \otimes f_2 \otimes f_1, m, op)$  would be translated as follows: we would add (i)  
 475 the triple patterns  $?x_1 f_1 \text{ right}(f_1), ?x_1 f_2 \text{ right}(f_2), \dots, ?x_1 f_k \text{ right}(f_k)$  to the WHERE clause and (iii) the  
 476 variables  $\text{right}(f_1), \text{right}(f_2), \dots, \text{right}(f_k)$  to the SELECT and to the GROUP BY clauses. In other words,  
 477 we would join the pairing functions i.e.  $f_1, f_2, \dots, f_k$  on their shared variable i.e.  $?x_1$ .

#### 478 6.5. The Full Algorithm for Translating a HIFUN Query to a SPARQL Query

479 Let us now describe the full translation algorithm. Let  $q = (gE/rg, mE/rm, opE/ro)$  be a HIFUN  
 480 query where

- 481 •  $gE$  is the grouping expression,
- 482 •  $mE$  is the measuring expression,
- 483 •  $opE$  is the operation expression, and
- 484 •  $rg$  is a restriction on the grouping expression,
- 485 •  $rm$  is a restriction on the measuring expression, and
- 486 •  $ro$  is a restriction on the operation expression.

487 The final query is constructed by concatenating strings as shown below:

```
488 Q = "SELECT " + retVars(gE) + " " + opE(mE) + "\n"
489 + "WHERE {" + "\n"
490 + triplePatterns(gE) + "\n"
491 + triplePatterns(mE) + "\n"
492 + "}" + "\n"
493 + "GROUP BY " + retVars(gE) + "\n"
494 + "HAVING " + restr(Qans)
```

495

496 The way the variables  $\text{retVars}(gE), \text{opE}(mE), \text{triplePatterns}(gE), \text{triplePatterns}(mE), \text{retVars}(gE),$   
 497 and  $\text{restr}(Q_{ans})$  are constructed, is described next

498 1. We start the translation with the grouping expression  $gE$  by creating the string format of the  
 499 triple patterns in which the terms  $g_i$  of  $gE$  participates,  $\text{triplePatterns}(gE) += ?x_i g_i \text{ right}(g_i)$ , as  
 500 described in §6.1 and in §6.4.

501 If  $gE$  contains any restriction  $rg$  we supplementarily create the string format of the triple pattern  
 502 expressing that constraint:

- 503 1.1. if  $rg$  refers to a URI, then  $\text{triplePatterns}(gE) += ?x_i g_i rg$ ,
- 504 1.2. if  $rg$  is represented with a LITERAL, then  $\text{triplePatterns}(gE) += \text{FILTER}(\text{right}(g_i) rg)$ , as  
 505 described in §6.2.

506 2. We proceed with the translation of the measuring expression  $mE$  by creating the string format of  
 507 the triple patterns in which the terms  $m_i$  of  $mE$  participates,  $\text{triplePatterns}(mE) += x_i m_i \text{ right}(m_i)$ .  
 508 Since this expression can also be complex, the translation is made as described in §6.1 and in §6.4.

509 If  $mE$  contains any restriction  $rm$  we supplementarily create the string format of the triple pattern  
 510 expressing that constraint:

- 511 2.1. if  $rm$  refers to a URI, then  $\text{triplePatterns}(mE) += ?x_i m_i rm$ ,
- 512 2.2. if  $rm$  is represented with a LITERAL, then  $\text{triplePatterns}(mE) += \text{FILTER}(\text{right}(m_i) rm)$ , as  
 513 described in §6.2.

514 3. Following, we create the string format of the returned variables,  $\text{retVars}(gE) += \text{right}(g_i)$  as  
 515 described in §6.1 and in 6.4.

516 4. At last, we translate the aggregate expression  $opE$  by creating the string format of the operation  
 517  $op$  applied over the values of  $mE$ , i.e.,  $\text{opE}(mE) = \text{op}(\text{right}(m_i))$ , as described in §6.1.

518 5. Optionally, if any restrictions  $re$  are applied to the final answers  $Q_{ans}$ , then we create the string  
519 format of the condition expressing these constraints,  $restr(Q_{ans}) = right(m_i)$   $re$  as described in §6.3.

520 To give an example suppose that we would like to find the total quantities by branch and brand  
521 only for the month of January, by considering only (a) the invoices with quantity greater than or equal  
522 to 2, and (b) the branches with total quantity greater than 1000. This query would be expressed in  
523 HIFUN as  
524  $(takesPlaceAt \otimes (brand \circ delivers)) /_{month=01}, inQuantity_{>=2}, SUM /_{>1000}$  and (by following the above  
525 translation process) in SPARQL as:

```
526
527 SELECT ?x2 ?x5 SUM(?x3)
528 WHERE {
529   ?x1 ex:takesPlaceAt ?x2 .
530   ?x1 ex:inQuantity ?x3 .
531   ?x1 ex:delivers ?x4 .
532   ?x4 ex:brand ?x5 .
533   ?x1 ex:hasDate ?x6 .
534   FILTER((MONTH(?x6) = 01) && (?x3 >= xsd:integer("2")))
535 }
536 GROUP BY ?x2 ?x5
537 HAVING (SUM(?x3) > 1000)
538
539
```

540 The pseudocode of the algorithm that defines the variables  $retVars(gE)$ ,  $opE(mE)$ ,  
541  $triplePatterns(gE)$ ,  $triplePatterns(mE)$ ,  $retVars(gE)$ , and  $restr(Q_{ans})$ , for the simple case, where the  
542 HIFUN query *does not contain compositions and pairings*, is given in Algorithm 1.

---

#### Algorithm 1 Algorithm for computing the components of the translated query for the **Simple Case**

---

**Require:** A HIFUN query  $q = (g/rg, m/rm, op/ro)$

**Ensure:**  $retVars(g)$ ,  $op(m)$ ,  $triplePatterns(g)$ ,  $triplePatterns(m)$ ,  $retVars(g)$ , and  $restr(Q_{ans})$

```

1: right(g) ← newVariable()
2: triplePatterns(g).concat(?x1 g right(g))           ▷ Grouping function
3: if rg <> ε then                                   ▷ if rg is not empty
4:   if rg.endsWithURI() then                       ▷ Group restriction involving a URI
5:     triplePatterns(g).concat(?x1 g rg)           ▷ The restriction is expressed as a triple pattern
6:   else                                           ▷ Group restriction involving a literal
7:     triplePatterns(g).concat(FILTER(right(g) rg)) ▷ The restriction is expressed as a filter
8:   end if
9: end if
10: right(m) ← newVariable()
11: triplePatterns(m).concat(?x1 m right(m))        ▷ Measuring function
12: if rm <> ε then                                   ▷ if rm is not empty
13:   if rm.containsURI() then                       ▷ Measuring restriction involving a URI
14:     triplePatterns(m).concat(right(m) m rm)     ▷ The restriction is expressed as a triple pattern
15:   else                                           ▷ Measuring restriction involving a literal
16:     triplePatterns(m).concat(FILTER(right(m) rm)) ▷ The restriction is expressed as a filter
17:   end if
18: end if
19: retVars(g).concat(right(g))                     ▷ for the SELECT and the GROUP BY clauses
20: op(m) = op(right(m))                             ▷ for the SELECT clause
21: restr(Qans) = right(m) ro                       ▷ for the HAVING clause

```

---

543 However, in the general case we may have *compositions* in grouping, measuring and restrictions.  
544 They way compositions are translated is described in Algorithm 2. The extension of the composition  
545 algorithm that supports also *derived* attributes is given in Alg. 3. Note that all predefined functions of  
546 SPARQL with one parameter can be used straightforwardly as derived attributes.

547 In addition, Alg. 2 shows how *pairing* is translated. If  $gE$  involves *both* pairing and  
 548 composition it will have the form  $gc_1 \otimes \dots \otimes gc_k$  where each  $gc_i$  can be an individual function  
 549 or a composition of functions. Such expressions are translated as follows:  $\text{translate}(gE) =$   
 550  $\text{translatePairing}(\text{translateComposition}(gc_1) \otimes \dots \otimes \text{translateComposition}(gc_k))$ . The exact steps for  
 551 translating pairing of compositions are given in Alg. 2-PairingAndComposition.

---

**Algorithm 2** Auxiliary algorithms for compositions and pairings
 

---

```

1: procedure COMPOSITION( $fc = f_k \circ \dots \circ f_1$ )  $\triangleright$  returns the triplePatterns and the retVars for a composition
2:    $tp \leftarrow ""$ ;
3:   for  $i \in 1..k$  do
4:      $\text{right}(f_i) \leftarrow \text{newVariable}()$ 
5:     if  $k=1$  then
6:        $tp.\text{concat}(\text{?x1 } f_1 \text{ right}(f_1))$ 
7:     else
8:        $tp.\text{concat}(\text{right}(f_{i-1}) \text{ fi } \text{right}(f_i))$ 
9:     end if
10:  end for
11:  return  $tp, \text{right}(f_k)$ 
12: end procedure

1: procedure PAIRING( $fp = f_1 \otimes \dots \otimes f_k$ )  $\triangleright$  returns the triplePatterns and the retVars for a pairing expression
2:    $tp \leftarrow ""$ ;  $\text{retVars} \leftarrow ""$ 
3:   for  $i \in 1..k$  do
4:      $\text{right}(f_i) \leftarrow \text{newVariable}()$ 
5:      $tp.\text{concat}(\text{?x1 } f_i \text{ right}(f_i))$ 
6:      $\text{retVars}.\text{concat}(\text{right}(f_i))$ 
7:   end for
8:   return  $tp, \text{retVars}$ 
9: end procedure

1: procedure PAIRINGANDCOMPOSITION( $fp = gc_1 \otimes \dots \otimes gc_k$ )  $\triangleright$  Pairing over Compositions
2:    $tp \leftarrow ""$ ;  $\text{retVars} \leftarrow ""$ 
3:   for  $i \in 1..k$  do
4:      $tp.\text{concat}(\text{Composition}(gc_i).tp)$ 
5:      $\text{retVars}.\text{concat}(\text{Composition}(gc_i).\text{retVars})$ 
6:   end for
7:   return  $tp, \text{retVars}$ 
8: end procedure

```

---

**Algorithm 3** Algorithm for composition if derived attributes are involved
 

---

```

1: procedure COMPOSITIONSUPPORTINGDERIVED( $fc = f_k \circ \dots \circ f_1$ )
2:    $tp \leftarrow ""$ ;  $\text{retVars} \leftarrow ""$ ;
3:   for  $i \in 1..k$  do
4:     if  $f_i$  is not a derived attribute then
5:        $\text{right}(f_i) \leftarrow \text{newVariable}()$ 
6:       if  $i=1$  then
7:          $tp.\text{concat}(\text{?x1 } f_i \text{ right}(f_i))$ 
8:       else
9:          $tp.\text{concat}(\text{right}(f_{i-1}) \text{ fi } \text{right}(f_i))$ 
10:      end if
11:       $\text{retVars} \leftarrow \text{right}(f_i)$ 
12:    else
13:       $\text{retVars} \leftarrow f_i(\text{retVars})$   $\triangleright$  Is a derived attribute
14:    end if  $\triangleright$  No triple pattern will be produced
15:  end for
16:  return  $tp, \text{retVars}$ 
17: end procedure

```

---

552 Algorithm 4 is the algorithm for the general case, where compositions can occur in the restrictions  
 553 too. Notice that now  $rg$  is not necessarily a single URI or literal, but a path expression that  
 554 ends with a URI or literal. In this scenario, instead of  $(\text{takesPlaceAt}/\text{"branch1"}, \text{inQuantity}, \text{SUM})$   
 555 we write  $(\text{takesPlaceAt}/\text{takesPlaceAt}=\text{branch1}, \text{inQuantity}, \text{SUM})$ , and we now support expressions of  
 556 the form  $(\text{takesPlaceAt}/\text{location}\text{takesPlaceAt}=\text{ex.Athens}, \text{inQuantity}, \text{SUM})$ . The path of the restriction

557 is not necessarily the same with that of the grouping, e.g. we can get the sum of quantities  
 558 grouped by brand, only for those branches that are located in Athens by the following query  
 559 ( $brand \circ delivers/location \circ takesPlaceAt = Athens, inQuantity, sum$ ). Such expressions are supported also for  
 560 the restriction  $rm$  of the measuring function.

---

**Algorithm 4** Algorithm for computing the components of the translated query for the **General** case
 

---

**Require:** A HIFUN query  $q = (gE/rg, mE/rm, opE/ro)$   
**Ensure:**  $retVars(gE)$ ,  $opE(mE)$ ,  $triplePatterns(gE)$ ,  $triplePatterns(mE)$ ,  $retVars(gE)$ , and  $restr(Q_{ans})$

```

1: triplePatterns(gE) = PairingAndComposition(gE).tp           ▷ triplePatterns for Grouping
2: retVars(gE) = PairingAndComposition(gE).retVars           ▷ retVars for Grouping
3: if rg <> ε then                                           ▷ if rg is not empty
4:   triplePatterns(gE).concat(Composition(rg.functions).tp) ▷ For supporting restrictions by compositions
5:   t ← last(triplePatterns(gE))                            ▷ the last triple pattern of triplePatterns(gE)
6:   rgSuffix ← rg.LastWord                                  ▷ The last token after the function composition
7:   if isURI(rgSuffix) then                                  ▷ Group restriction involving a URI
8:     Replace the right variable of t with rgSuffix
9:   else                                                     ▷ Group restriction involving a literal
10:    Add "FILTER right(t) rg.op rgSuffix"
11:  end if
12: end if
13: triplePatterns(mE) = Composition(mE).tp                 ▷ triplePatterns for Measuring
14: restr(Qans) = Composition(mE).retVars ro                ▷ for the HAVING clause
15: if rm <> ε then                                           ▷ if rm is not empty
16:   triplePatterns(mE).concat(Composition(rm.functions).tp) ▷ For supporting restrictions by compositions
17:   t ← last(triplePatterns(mE))
18:   rmSuffix ← rm.LastWord                                  ▷ The last token after the function composition
19:   if isURI(rmSuffix) then                                  ▷ Measuring restriction involving a URI ▷ Measuring restriction involving a URI
20:     Replace the right variable t with rmSuffix
21:   else                                                     ▷ Measuring restriction involving a literal
22:     Add "FILTER right(t) rm.op rmSuffix"
23:   end if
24: end if
25: opE(mE) = opE(Composition(mE).retVars)                 ▷ for the SELECT clause

```

---

561 The above translation process can also support cases where the set of resources  $R$  of an analysis  
 562 context is defined by one unary query, a query that returns a set of URIs. If  $q_{str}(R)$  denotes the string  
 563 of that query, then  $q_{str}(R)$  can be used as the starting point of the above (query translation) process, i.e.  
 564 the triple patterns of  $q_{str}(R)$  will be the first to be added to the triples patterns of our query  $Q$ ; the only  
 565 constraint is that  $q_{str}(R)$  should have one variable  $?x1$  in the select clause.

566 In general, the translation algorithm works for all possible HIFUN queries.

#### 567 6.6. Cases where the Prerequisites of HIFUN are not Satisfied

568 Let us now discuss the problem of translation for the case where the requirements for applying  
 569 HIFUN (as described in Section 5.1) do not hold, as well as when features (as described in Section 5.4)  
 570 are required.

571 Consider that the running example of Figure 2 contained a property "ex:birthYear" with domain  
 572 the class "ex:Person" and range an integer-typed literal. Now suppose that we would like to compute  
 573 a *single number* being the average birth year of the founders of the products that were sold. Here  
 574 the problem of *incomplete information* arises, i.e. the dataset may not contain information about the  
 575 founders of all products, let alone their birth year, therefore the path  $delivers.brand.founder.birthYear$   
 576 will not be "applicable" to several invoices. If we formulate such a query in HIFUN, it will be translated  
 577 to a SPARQL query that will be evaluated successfully, however the results will not be complete, i.e.  
 578 only the invoices for which the path  $delivers.brand.founder.birthYear$  exists will be considered.

579 Moreover we may encounter the problem of *multiple values*, i.e. when a brand was founded by  
 580 more than one person. In that case, even if our dataset contained complete information, the path  
 581  $delivers.brand.founder.birthYear$  would not be functional. If we formulate such a query in HIFUN, it  
 582 will be translated to a SPARQL query that will be evaluated successfully, however the results will not



583 be accurate, i.e. all paths will be taken into account; valuating more than one birth year per product,  
 584 not only one birth year per product.

585 If we wanted to associate each product with *only one* birth year (before taking the average over all  
 586 products), then in case of multiple founders, we could define a *feature* that computes the average birth  
 587 year of each individual product. Having this feature enables the subsequent formulation of a HIFUN  
 588 query, that would compute the accurate answer. There are several possible methods to define such  
 589 features using the query language (as mentioned in Section 5.4). In our example, the average birth  
 590 year of the founder(s) for each individual product can be computed by one query that yields a variable  
 591 for that feature:

```
592
593  $qf =$ 
594 SELECT ?x2 AVG(?x5) as productFoundBirthYearAvg
595 WHERE {
596   ?x1 ex:delivers ?x2 .
597   ?x2 ex:brand ?x3 .
598   ?x3 ex:founder ?x4 .
599   ?x4 ex:birthYear ?x5 .
600 }
601 GROUP BY ?x2
602
```

603 To compute accurately that we want to (i.e. average birth year of the founders of the products  
 604 that were sold), we can exploit the notion of SPARQL *subqueries* (that was mentioned in the last part of  
 605 Section 5.4) to “embed” the aforementioned query  $qf$ .

606 Note that subqueries are a way to embed SPARQL queries inside other queries to allow the  
 607 expression of requests that are not possible otherwise. Subqueries are evaluated first and then the  
 608 outer query is applied to their results Only variables projected out of the subquery (i.e., appearing in  
 609 its SELECT clause) are visible to the outer query. Therefore the features can be expressed as subqueries  
 610 that can be placed in the WHERE clause.

611 In our case, the HIFUN query that computes the average birth year of the founders of  
 612 the products would have the following form:  $Q = (\epsilon, product.productFoundBirthYearAvg, AVG)$ ,  
 613 where  $productFoundBirthYearAvg$  is a feature (according to Table 1 we would write  
 614  $productFoundBirthYearAvg = delivers.brand.founder.birthYear.avg$ ); note that  $\epsilon$  denotes the empty  
 615 grouping function, since in our example we do not want to group the results, just to apply  $AVG$  to  
 616 the entire set. To compute this feature we can use a subquery that returns two variables, one for the  
 617 objects and one for the corresponding feature value, say  $vf1$  and  $vf2$  respectively, while ensuring that  
 618  $vf1$  should be the same with the variable used in the outer query for these objects. In our example, the  
 619 corresponding SPARQL query that computes the sought answer (where the subquery provides two  
 620 variables), is the following:

```
SELECT AVG(?productFoundBirthYearAvg)
WHERE {
  ?x1 ex:delivers ?x2 .
  {
    SELECT ?x2 (AVG(?x5) AS ?productFoundBirthYearAvg)
    WHERE {
      ?x1 ex:delivers ?x2 .
      ?x2 ex:brand ?x3 .
      ?x3 ex:founder ?x4 .
      ?x4 ex:birthYear ?x5 .
    }
  }
  GROUP BY ?x2
}
```

}

621 This example showcased how we can compute accurate results if the paths that are involved in  
622 the HIFUN query are not functional.

### 623 6.7. Analytics and RDF Schema Semantics

624 The translation of SPARQL allows leveraging the RDF Schema semantics, specifically the RDF  
625 Schema-related inference that is supported by SPARQL. To give a simple indicative example, consider  
626 two properties `directorOf` and `worksAt` such that `(directorOf, rdfs:subPropertyOf, worksAt)` and  
627 suppose the following data:

```
628 (p1,directorOf,brand1)
629 (p2,worksAt,brand1)
630 (p3,worksAt,brand1)
631 (p4,worksAt,brand1)
632 (p1,livesAt,Athens)
633 (p2,livesAt,Rhodes)
634 (p3,livesAt,Corfu),
635 (p4,livesAt,Corfu)
```

636 If we want to compute the locations where the persons related to `brand1` work at, and how  
637 many live at each place, we could use the following HIFUN query (`workAt,livesAt,COUNT`). If the  
638 inference of SPARQL is enabled, then the translated query will return

```
639 Athens,1
640 Rhodes,1
641 Corfu,2
```

642 The key point is that the location of the director will be considered, since it is inferred that  
643 `(p1,worksAt,brand1)`. Such inference would not be possible if we translated the data to the relational  
644 model.

645 The ability to leverage the inference rules of RDF Schema in the context of analytic queries is  
646 especially important for datasets which are described with ontologies that contain high number of  
647 `subClassOf` and `subPropertyOf` relationships for achieving semantic interoperability across various  
648 datasets, like those in the cultural and historic domain [7,42].

## 649 7. On Interactivity

650 As described in the introductory section, and illustrated in Figure 1, our ultimate objective is to  
651 provide a user friendly method for interactive analytics over any RDF graph. This requires:

652  $\mathcal{S}_{\textcircled{1}ctx}$  An interactive method for specifying an Analysis Context. Recall that (according to Def. 3), an  
653 analysis context  $C$  over RDF data is defined as a set of resources  $R$  to be analyzed along with a  
654 set of properties  $p_1, p_2, \dots, p_n$  that are relevant for the analysis.

655  $\mathcal{S}_{\textcircled{2}q}$  An interactive method for formulating the desired HIFUN query  $q = (g, m, op)$ , i.e. a method to  
656 select  $g$ ,  $m$ , and  $op$ .

657  $\mathcal{S}_{\textcircled{3}tr}$  A method to translate the HIFUN query to SPARQL.

658  $\mathcal{S}_{\textcircled{4}vis}$  A method for visualizing the results of the SPARQL query that is derived by translating  $q$ .

659 As regards  $\mathcal{S}_{\textcircled{1}ctx}$ , the analysis context can be specified over the original data (Section 5.3), or  
660 after a transformation (Section 5.4). In the former case, the user does not have to do anything: all  
661 resources of the RDF dataset and all properties can be exploited as an analysis context. In the latter  
662 case the intended transformation can be made by tools like `LODSynthesisML` [40] that allows the user to

663 interactively select the desired features. The output of this step can be either (a) a csv file or an RDF  
 664 file in RDF Data Cube format that contains the materialization of the defined features, or (b) a set of  
 665 feature specifications each being a pair (feature name, SPARQL query) which will be exploited in the  
 666 translation process as discussed in the last part of Section 5.4 and in Section 6.6.

667 As regards  $\mathcal{S}_{\text{tr}}$ , and assuming that an analysis context has been specified, we have developed  
 668 a tool, called  $\text{HIFUN}_{\text{RDF}}$ , where user is asked to select the functions of a HIFUN query i.e. (i) the  
 669 grouping function, (ii) the measuring function, (iii) the aggregate operation, and optionally, (iv) set  
 670 restrictions to the grouping, the measuring functions or to the final results. The above are accomplished  
 671 by *interactively selecting* the desired properties, i.e. the user does not have to know the SPARQL syntax.  
 672 Currently, this tool loads data represented in the RDF Data Cube format and stores them in a triplestore,  
 673 however this period we are extending this tool for supporting any RDF file (not only files in RDF Data  
 674 Cube format).

675 As regards  $\mathcal{S}_{\text{tr}}$ , we have implemented the translation method described in Section 6. Note  
 676 that the cost of the translation of a HIFUN query to SPARQL is negligible (the translation has linear  
 677 complexity with respect to the size of the HIFUN string; it does not depend on the size of the data).  
 678 The translated to SPARQL query is then executed on the triple store OpenLink Virtuoso<sup>9</sup> (where the  
 679 input data has been uploaded) and the returned results are displayed and saved in a “.csv” file in the  
 680 form of -  $var_1, var_2, \dots, var_i, \text{TOTALS}$ .

681 As regards  $\mathcal{S}_{\text{vis}}$ , we have embedded in  $\text{HIFUN}_{\text{RDF}}$  the jfreechart library<sup>10</sup> for reading the results  
 682 of the SPARQL query and preparing various charts including line, bar, pie and 3D pie charts. Three  
 683 screenshots are shown in Figure 4 that visualize the results of the query “total quantities by branch”  
 684 over a synthetic RDF dataset that uses the same schema with our running example.

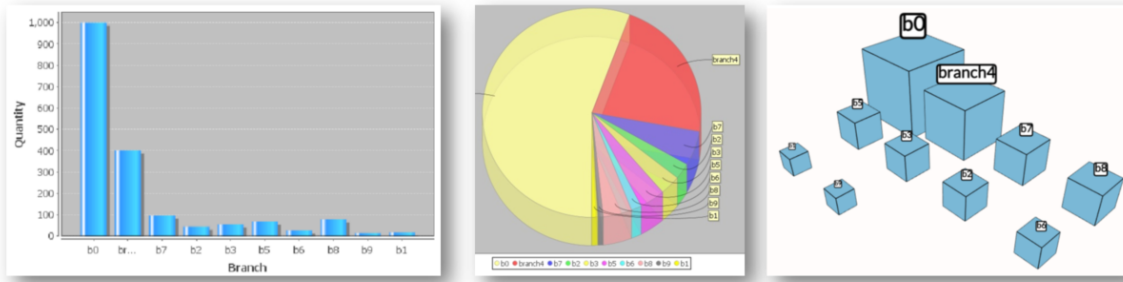


Figure 4. Indicative visualizations produced by  $\text{HIFUN}_{\text{RDF}}$

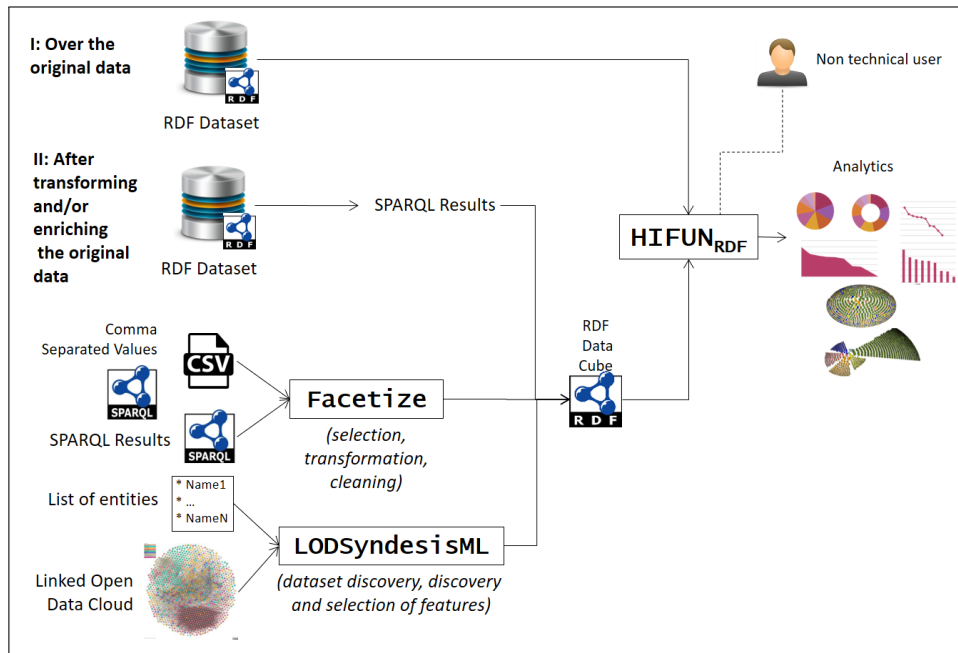
685 The indicative flows that are possible with this implementation, are illustrated in Figure 5. The  
 686 figure includes a workflow where the process starts with selecting the data set the user wants to  
 687 analyze using tools like Facetize [43] and LODSyndesisML [40], suitable for discovering, cleaning,  
 688 organizing data in hierarchies etc. Then, this data is converted into the RDF Data Cube format and  
 689 is loaded in  $\text{HIFUN}_{\text{RDF}}$ .  $\text{HIFUN}_{\text{RDF}}$  is not yet public, a public version will be released after tackling the  
 690 extensions described next in Section 7.1.

### 691 7.1. Future Work

692 Currently we are working on a second implementation that will support, *in a single system*, and  
 693 *interactively*, all steps of the process for  $\mathcal{S}_{\text{ctx}}$  to  $\mathcal{S}_{\text{vis}}$ . The objective is to provide a unified interface  
 694 that will enable the user to (i) select the RDF file or triple store (s)he wants to analyze, (ii) specify  
 695 and change the analysis context on the fly (and provide the capability to define features as described  
 696 in Section 5.4), and (iii) formulate an analytical query by defining its components by selecting the

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

<sup>10</sup> <https://www.jfree.org/jfreechart>



**Figure 5.** A few indicative workflows involved in  $HIFUN_{RDF}$ .

697 applicable properties. Then, the system will translate the query and finally visualize the results in  
 698 tabular form as well as in various other forms (like those described in [44]) including 3D (by extending  
 699 the work of [45]<sup>11</sup>) allowing the user to explore them, intuitively. For steps (ii) and (iii) we plan to  
 700 investigate extending the core model for exploratory search over RDF that is described in [12] (for  
 701 exploiting its capabilities of forming conditions that involve paths and complex expressions), with  
 702 actions for supporting  $S_{\mathcal{D}_{ctx}}$  and  $S_{\mathcal{Q}_q}$  (that are not supported by that model). That will enable the  
 703 user to specify interactively and with simple clicks complex restrictions that may concern the set of  
 704 resources of the analysis context, and/or the various restrictions of the analytic query (of grouping or  
 705 measuring functions, or of the final results).

706 Finally, we should mention that an interactive system for analytics that uses HIFUN as the  
 707 intermediate layer, enables the formulation of analytic queries over data sources with different data  
 708 models and query languages, since there are already mappings of HIFUN to SQL (for relational  
 709 sources), and to MapReduce using SPARK [46], thus our work enables applying that model over RDF  
 710 data sources too.

## 711 8. Concluding Remarks

712 In this paper we elaborated on the general problem of providing interactive analytics over RDF  
 713 data. We motivated this direction, we described the main requirements and challenges, and we  
 714 discussed the work that has been done in this area. Subsequently we investigated whether HIFUN, a  
 715 functional query language for analytics, can be used as a means for formulating analytic queries over  
 716 RDF data in a flexible manner. To this end, we analyzed the applicability of HIFUN over RDF, we  
 717 described various methods that can be used to apply HIFUN over such data, and then we focused  
 718 on the problem of translating HIFUN queries to SPARQL queries, starting from simple queries and  
 719 ending up to complex queries. We discussed what happens when HIFUN cannot be applied, and  
 720 how features can be employed to tackle properties and paths that are not functional. The presented  
 721 query translation approach does not require transforming or moving the data, and can leverage the

<sup>11</sup> <http://www.ics.forth.gr/isl/3DLod/>

722 inference that is provided by SPARQL. Subsequently we described how this approach can be exploited  
723 in an interactive context and we described our first implementation. Overall, we have shown that it is  
724 possible to provide a method for formulating complex analytic queries over RDF that is based on a  
725 few and basic concepts, those of HIFUN. In the future, we plan to implement a fully interactive system  
726 that will support all steps of the query formulation and visualization process and elaborate further on  
727 the interactivity of the approach.

## 728 References

- 729 1. Mountantonakis, M.; Tzitzikas, Y. Large-scale Semantic Integration of Linked Data: A Survey. *ACM*  
730 *Computing Surveys (CSUR)* **2019**, *52*, 103.
- 731 2. Bizer, C.; Lehmann, J.; Kobilarov, G.; Auer, S.; Becker, C.; Cyganiak, R.; Hellmann, S. DBpedia-A  
732 crystallization point for the Web of Data. *Journal of web semantics* **2009**, *7*, 154–165.
- 733 3. Vrandečić, D.; Krötzsch, M. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*  
734 **2014**, *57*, 78–85.
- 735 4. Wishart, D.S.; Feunang, Y.D.; Guo, A.C.; Lo, E.J.; Marcu, A.; Grant, J.R.; Sajed, T.; Johnson, D.; Li, C.;  
736 Sayeeda, Z.; others. DrugBank 5.0: a major update to the DrugBank database for 2018. *Nucleic acids research*  
737 **2018**, *46*, D1074–D1082.
- 738 5. Tzitzikas, Y.; Marketakis, Y.; Minadakis, N.; Mountantonakis, M.; Candela, L.; Mangiacrapa, F.; others.  
739 Methods and Tools for Supporting the Integration of Stocks and Fisheries. Chapter in *Information and*  
740 *Communication Technologies in Modern Agricultural Development*, Springer, 2019. Springer, 2019.
- 741 6. Jaradeh, M.Y.; Oelen, A.; Farfar, K.E.; Prinz, M.; D'Souza, J.; Kismihók, G.; Stocker, M.; Auer, S.  
742 Open Research Knowledge Graph: Next Generation Infrastructure for Semantic Scholarly Knowledge.  
743 Proceedings of the 10th International Conference on Knowledge Capture, 2019, pp. 243–246.
- 744 7. Koho, M.; Ikkala, E.; Leskinen, P.; Tamper, M.; Tuominen, J.; Hyvönen, E. WarSampo Knowledge Graph:  
745 Finland in the Second World War as Linked Open Data. *Semantic Web – Interoperability, Usability, Applicability*  
746 **2020**. In press, doi:10.3233/SW-200392.
- 747 8. Dimitrov, D.; Baran, E.; Fafalios, P.; Yu, R.; Zhu, X.; Zloch, M.; Dietze, S. TweetsCOVID19—A Knowledge Base  
748 of Semantically Annotated Tweets about the COVID-19 Pandemic. 29th ACM International Conference on  
749 Information and Knowledge Management (CIKM 2020), 2020.
- 750 9. COVID-19 Open Research Dataset (CORD-19). 2020.
- 751 10. R. Gazzotti, F. Michel, F.G. CORD-19 Named Entities Knowledge Graph (CORD19-NEKG). 2020.
- 752 11. Nikas, C.; Kadilierakis, G.; Fafalios, P.; Tzitzikas, Y. Keyword Search over RDF: Is a Single Perspective  
753 Enough? *Big Data and Cognitive Computing* **2020**, *4*, 22.
- 754 12. Tzitzikas, Y.; Manolis, N.; Papadacos, P. Faceted exploration of RDF/S datasets: a survey. *Journal of*  
755 *Intelligent Information Systems* **2017**, *48*, 329–364.
- 756 13. Kritsotakis, V.; Roussakis, Y.; Patkos, T.; Theodoridou, M. Assistive Query Building for Semantic Data.  
757 SEMANTICS Posters&Demos, 2018.
- 758 14. Spyratos, N.; Sugibuchi, T. HIFUN—a high level functional query language for big data analytics. *Journal of*  
759 *Intelligent Information Systems* **2018**, *51*.
- 760 15. Papadaki, M.E.; Tzitzikas, Y.; Spyratos, N. Analytics over RDF Graphs. International Workshop on  
761 Information Search, Integration, and Personalization, 2019.
- 762 16. Spyratos, N. A functional model for data analysis. International Conference on Flexible Query Answering  
763 Systems, 2006.
- 764 17. Tzitzikas, Y.; Allocca, C.; Bekiari, C.; Marketakis, Y.; Fafalios, P.; Doerr, M.; Minadakis, N.; Patkos, T.;  
765 Candela, L. Integrating heterogeneous and distributed information about marine species through a top  
766 level ontology. Research conference on metadata and semantic research, 2013.
- 767 18. Isaac, A.; Haslhofer, B. Europeana linked open data—data. europeana. eu. *Semantic Web* **2013**, *4*.
- 768 19. Mountantonakis, M.; Tzitzikas, Y. On measuring the lattice of commonalities among several linked datasets.  
769 *Proceedings of the VLDB Endowment* **2016**, *9*.
- 770 20. Mountantonakis, M.; Tzitzikas, Y. Scalable Methods for Measuring the Connectivity and Quality of Large  
771 Numbers of Linked Datasets. *Journal of Data and Information Quality (JDIQ)* **2018**, *9*.
- 772 21. Roatis, A. Analysing RDF Data: A Realm of New Possibilities. *ERCIM News* **2014**, *2014*.

- 773 22. Kämpgen, B.; O’Riain, S.; Harth, A. Interacting with statistical linked data via OLAP operations. *Extended*  
774 *Semantic Web Conference*, 2012.
- 775 23. Etcheverry, L.; Vaisman, A.A. QB4OLAP: a new vocabulary for OLAP cubes on the semantic web.  
776 *Proceedings of the Third International Conference on Consuming Linked Data*, 2012.
- 777 24. Azirani, E.A.; Goasdoué, F.; Manolescu, I.; Roatis, A. Efficient OLAP operations for RDF analytics. 2015  
778 31st IEEE International Conference on Data Engineering Workshops. IEEE, 2015, pp. 71–76.
- 779 25. Ruback, L.; Pesce, M.; Manso, S.; Ortiga, S.; Salas, P.E.R.; Casanova, M.A. A mediator for statistical linked  
780 data. *Proceedings of the 28th Annual ACM Symposium on Applied Computing*, 2013.
- 781 26. Etcheverry, L.; Vaisman, A.A. Enhancing OLAP analysis with web cubes. *Extended Semantic Web*  
782 *Conference*, 2012.
- 783 27. Zhao, P.; Li, X.; Xin, D.; Han, J. Graph cube: on warehousing and OLAP multidimensional networks.  
784 *Proceedings of the 2011 ACM SIGMOD International Conference on Management of data*, 2011.
- 785 28. Benatallah, B.; Motahari-Nezhad, H.R.; others. Scalable graph-based OLAP analytics over process  
786 execution data. *Distributed and Parallel Databases* **2016**, 34.
- 787 29. Wang, K.; Xu, G.; Su, Z.; Liu, Y.D. GraphQ: Graph Query Processing with Abstraction Refinement—Scalable  
788 and Programmable Analytics over Very Large Graphs on a Single {PC}. 2015 Annual Technical Conference  
789 (15), 2015.
- 790 30. Zapolko, B.; Mathiak, B. Performing statistical methods on linked data. *International conference on dublin*  
791 *core and metadata applications*, 2011.
- 792 31. Olston, C.; Reed, B.; Srivastava, U.; Kumar, R.; Tomkins, A. Pig latin: a not-so-foreign language for data  
793 processing. *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, 2008.
- 794 32. Thusoo, A.; Sarma, J.S.; Jain, N.; Shao, Z.; Chakka, P.; Zhang, N.; Antony, S.; Liu, H.; Murthy, R. Hive-a  
795 petabyte scale data warehouse using hadoop. 2010 IEEE 26th international conference on data engineering  
796 (ICDE 2010), 2010.
- 797 33. Etcheverry, L.; Vaisman, A.A. Querying Semantic Web Data Cubes. *Proc. Alberto Mendelzon Int.*  
798 *Workshop on Foundations of Data Management*, 2016.
- 799 34. Etcheverry, L.; Vaisman, A.A. Efficient Analytical Queries on Semantic Web Data Cubes. *Journal on Data*  
800 *Semantics* **2017**, 6.
- 801 35. Colazzo, D.; Goasdoué, F.; Manolescu, I.; Roatis, A. RDF analytics: lenses over semantic graphs.  
802 *Proceedings of the 23rd international conference on World wide web*, 2014.
- 803 36. Diao, Y.; Guzewicz, P.; Manolescu, I.; Mazuran, M. Spade: a modular framework for analytical exploration  
804 of RDF graphs. 2019.
- 805 37. Antoniou, G.; Van Harmelen, F. *A semantic web primer*; MIT press, 2004.
- 806 38. Mountantonakis, M.; Tzitzikas, Y. LODsyndesis: Global Scale Knowledge Services. *Heritage* **2018**, 1.
- 807 39. Spyratos, N.; Sugibuchi, T. Data Exploration in the HIFUN Language. *International Conference on Flexible*  
808 *Query Answering Systems*, 2019.
- 809 40. Mountantonakis, M.; Tzitzikas, Y. How linked data can aid machine learning-based tasks. *International*  
810 *Conference on Theory and Practice of Digital Libraries*, 2017.
- 811 41. Mami, M.N.; Graux, D.; Thakkar, H.; Scerri, S.; Auer, S.; Lehmann, J. The query translation landscape: a  
812 survey. *arXiv preprint arXiv:1910.03118* **2019**.
- 813 42. Fafalios, P.; Petrakis, C.; Samaritakis, G.; Doerr, K.; Tzitzikas, Y.; Doerr, M. FastCat: Collaborative Data  
814 Entry and Curation for Semantic Interoperability in Digital Humanities". *ACM Journal on Computing and*  
815 *Cultural Heritage* **2021**. accepted for publication.
- 816 43. Kokolaki, A.; Tzitzikas, Y. Facetize: An Interactive Tool for Cleaning and Transforming Datasets for  
817 Facilitating Exploratory Search. *arXiv preprint arXiv:1812.10734* **2018**.
- 818 44. Andrienko, G.; Andrienko, N.; Drucker, S.; Fekete, J.D.; Fisher, D.; Idreos, S.; Kraska, T.; Li, G.; Ma, K.L.;  
819 Mackinlay, J.; others. Big Data Visualization and Analytics: Future Research Challenges and Emerging  
820 Applications. *BigVis 2020: Big Data Visual Exploration and Analytics*, 2020.
- 821 45. Papadaki, M.E.; Papadakis, P.; Mountantonakis, M.; Tzitzikas, Y. An Interactive 3D Visualization for the  
822 LOD Cloud. *EDBT/ICDT Workshops*, 2018.
- 823 46. Zervoudakis, P.; Kondylakis, H.; Plexousakis, D.; Spyratos, N. Incremental Evaluation of Continuous  
824 Analytic Queries in HIFUN. *International Workshop on Information Search, Integration, and*  
825 *Personalization*. Springer, 2019, pp. 53–67.

826 **Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional  
827 affiliations.

828 © 2021 by the authors. Submitted to *Algorithms* for possible open access publication  
829 under the terms and conditions of the Creative Commons Attribution (CC BY) license  
830 (<http://creativecommons.org/licenses/by/4.0/>).