

# Querying for Provenance, Trust, Uncertainty and other Meta Knowledge in RDF<sup>1</sup>

Renata Dividino<sup>a</sup>, Sergej Sizov<sup>a</sup>, Steffen Staab<sup>a</sup>, Bernhard Schueler<sup>a</sup>

<sup>a</sup>*ISWeb - Information Systems and Semantic Web  
University of Koblenz-Landau, Germany*

---

## Abstract

The Semantic Web is based on accessing and reusing RDF data from many different sources, which one may assign different levels of authority and credibility. Existing Semantic Web query languages, like SPARQL, have targeted the retrieval, combination and reuse of facts, but have so far ignored all aspects of meta knowledge, such as origins, authorship, recency or certainty of data.

In this paper, we present an original, generic, formalized and implemented approach for managing many dimensions of meta knowledge, like source, authorship, certainty and others. The approach re-uses existing RDF modeling possibilities in order to represent meta knowledge. Then, it extends SPARQL query processing in such a way that given a SPARQL query for data, one may request meta knowledge without modifying the query proper. Thus, our approach achieves highly flexible and automatically coordinated querying for data and meta knowledge, while completely separating the two areas of concern.

*Key words:* Semantic Web, Meta Knowledge, SPARQL, RDF

---

## 1. Introduction

Integrating and re-using Semantic Web data becomes more and more fruitful and worthwhile in order to answer questions and deliver results. Typically, engines like Swoogle provide points of access for RDF data, crawlers may fetch relevant RDF data, and query languages like SPARQL with their corresponding query engines allow for selecting and re-using data in the appropriate format. With the arrival of more and more data in the Semantic Web and more sophisticated processing through query and reasoning engines, one now,

however, encounters challenging questions linked to meta knowledge about the data like:

- Where is this data from?
- Who provided the data?
- When was this data provided?
- Was the provider certain about the truth of this data?
- Was the data believed by others, too?

For instance, when querying the Semantic Web with the help of SPARQL for the affiliation of a person named “James Hendler”, one finds (at least) two answers, i.e. “University of Maryland” and “Rensselaer Polytechnic Institute”. Without further indication as to where, by whom, when, etc. such information was given, it is impossible to decide which of the two affiliations is still valid.

The problem might be remedied in several ways. First, an ideosyncratic solution by the search engine, such as returning the corresponding RDF

---

*Email addresses:* dividino@uni-koblenz.de (Renata Dividino), sizov@uni-koblenz.de (Sergej Sizov), staab@uni-koblenz.de (Steffen Staab).

<sup>1</sup> This proposal is a completely revised and extended version of [22]. Major changes include proofs showing that the meta knowledge evaluation of SPARQL queries is equivalent to the standard SPARQL evaluation, a discussion of soundness and completeness, as well as an extended empirical evaluation section.

files or links to sources of knowledge extraction (say `<http://www.cs.umd.edu/survey.pdf>` and `<http://www.rpi.edu/report.doc>`), might help in this special case. However, an ideosyncratic solution may not be appropriate in a second case in which the ‘when’ was more relevant than the ‘where’ or in a third case where such a piece of information had to be aggregated from several resources. Second, the person or system requesting the meta knowledge might manually extend the SPARQL query formalizing the request for the affiliation in order to return the where, the who and the when. Such a modification will, however, be very tedious, as it will include a number of additional optional statements, and expressing it manually will be error prone. Also, it will not help in delivering meta knowledge that arises from joining several statements, e.g. meta knowledge about uncertainty that was based on several meta knowledge statements with different values of uncertainty. Therefore, querying Semantic Web data requires a principled, generic approach to the treatment of meta knowledge that is able to adapt to many dimensions of meta knowledge and that is open to accommodate to new dimensions when the need arises. Such a principled, original framework is given in this paper. We start to explain our approach with a discussion of important design choices in section 2. We model meta knowledge in existing RDF structures by embedding a slightly more expressive language, which we call RDF<sup>+</sup>, into RDF. We define the abstract syntax of RDF<sup>+</sup>, its semantics and its embedding in RDF in Section 4. In Section 5, we extend the SPARQL syntax and semantics to work on data and meta knowledge of RDF<sup>+</sup>. The extension allows the user to extend a given conventional SPARQL query by a keyword for meta knowledge triggering the construction of meta knowledge by the query processor. Section 6 summarizes the overall use and processing of SPARQL queries with meta knowledge. Section 7 and Section 8 report on initial graceful results for meta knowledge processing from a theoretic point of view. Finally, Section 9 provides pointers to the prototype implementation of the system and initial experimental testing.

## 2. Scenario

In our scenario, the sample user aims to explore the knowledge and meta knowledge by querying knowledge extracted from Web pages of Computer Science departments. Here, we assume that he aims to find experts in the domain of Semantic Web and their affiliations. Example 2.1 shows the relevant facts that may

have been obtained from websites of different universities.

### Example 2.1 Relevant facts obtained from different websites

Site	JamesHendler’s research topic is SemanticWeb.
A	JamesHendler is affiliated with RensselaerPI
Site	JamesHendler’s research topic Robotics.
B	JamesHendler affiliated with UnivMaryland. RudiStuder’s research topic SemanticWeb. RudiStuder affiliated with Univ. Karlsruhe

In addition, the user wants to exploit meta knowledge for obtaining results with best certainty and for analyzing contradicting answers (e.g. different affiliations for the same person “James Hendler” in Example 2.1). An example of meta knowledge associated to the extracted facts is presented in Example 2.2 and shows that the facts have been extracted from different-sources, at different timepoints, and with different degrees of extraction confidence.

### Example 2.2 Associated meta knowledge of facts presented in Example 2.1

Site	source = www.rpi.edu/report.doc
A	certainty degree = 0.9 timestamp = 5/5/2007
Site	source = www.cs.umd.edu/survey.pdf
B	certainty degree = 0.6 timestamp = 6/6/2001

In the next section, we discuss design choices for representing meta knowledge taking into consideration the existing Semantic Web representation structures.

## 3. Design choices

This section summarizes and shortly motivates the design choices for our meta knowledge framework.

**Reification.** Establishing relationships between knowledge and meta knowledge requires appropriate reification mechanisms for supporting statements about statements. Our general objective is to execute queries on original data (i.e. without meta knowledge) directly without complex transformations. For compliance with existing applications that access the repository in a common way (e.g. using SPARQL queries), we do

not modify existing user data. This requirement does not allow us to use mechanisms like RDF reification, which decompose existing triples and fully change the representational model. In our framework described in section 4, we adopt the notion of Named RDF Graphs for meta knowledge representation [5,6].

**Storage mechanisms.** Following the overall philosophy of RDF, we do not separate meta knowledge from “normal” user knowledge in the repository. Following this paradigm, a user or developer has unlimited access to all contents of the triple store and can manipulate meta knowledge directly. In other words, the user can directly access meta knowledge (e.g. using suitable SPARQL queries). Beyond explicitly designed queries for meta knowledge access, in Section 5 we describe the extension of SPARQL that allows us to access meta knowledge about the result set automatically without user intervention.

**Dimensions of Meta Knowledge.** An important point for the application design is the definition of relevant meta knowledge properties and their suitable interpretation for arbitrarily complex query patterns. In general, these properties are application dependent and must be carefully chosen by the system administrator. In our scenario (Section 2 and Section 6) we discuss common and widely used properties, such as timestamp, source, and (un)certainty, and show ways of defining and utilizing them in our framework.

**Syntax extensions.** Seamlessly integrated access to meta knowledge requires corresponding extensions of existing querying mechanisms. These can be realized at different levels, for instance at the level of query languages (e.g. SPARQL) or at the level of application-specific interfaces (e.g. Sesame API). In Section 5 we describe our SPARQL extension for constructing query results with associated meta knowledge. It is system-independent and not related to some particular implementation of the RDF repository. Furthermore, it fully supports the existing SPARQL syntax and semantics. Compliance with existing established standards makes the integration with existing applications and interfaces substantially easier.

#### 4. Syntax and Semantics for RDF with Meta Knowledge

In order to adapt our sample application scenario to our framework, we formalize it using the notion of Named RDF Graphs [5,6]. We assume that the user utilizes knowledge which has been initially extracted from Web pages of Computer Science departments and

stored in form of RDF triples in his personal ‘active space’ [21], backed by a local RDF repository. Example 4.1 shows the relevant facts already presented in Example 2.1 in RDF triple language TriG [1] with Named Graphs in a simplified form that abstracts from (default) namespaces.

#### Example 4.1 Facts of Example 2.1 represented in TriG with Named Graphs

---

```
G1 { JamesHendler researchTopic SemanticWeb.
      JamesHendler affiliatedWith RensselaerPI}
G2 { JamesHendler researchTopic Robotics.
      JamesHendler affiliatedWith UnivMaryland.
      RudiStuder researchTopic SemanticWeb.
      RudiStuder affiliatedWith UnivKarlsruhe}
```

---

Likewise, the meta knowledge associated to the extracted knowledge presented in Example 2.2 is stored into the same RDF repository using the notion of Named RDF Graphs.

#### Example 4.2 Meta knowledge of Example 2.2 represented in TriG with Named Graphs

---

```
G3 { G1 mk:source <www.rpi.edu/report.doc>.
      G1 mk:certainty "0.9".
      G1 mk:timestamp "5/5/2007" }
G4 { G2 mk:source <www.cs.umd.edu/survey.pdf>.
      G2 mk:certainty "0.6".
      G2 mk:timestamp "6/6/2001" }
```

---

In the examples above we have used existing RDF modeling possibilities in order to represent meta knowledge. However we must distinguish the notation of RDF with only *implicit* notation of meta knowledge, but no semantic consequences specifically due to this meta knowledge, from a formally extended model of RDF with *explicit* notation of meta knowledge.

In the course of representing and reasoning with meta knowledge we embed a language with meta knowledge reasoning, i.e. RDF<sup>+</sup>, in a language without such specific facilities, i.e. in RDF. This embedding implies that we may consider an RDF snippet in its literal sense *and* we may possibly interpret it as making a meta knowledge statement. Embedding meta knowledge in RDF is not the most expressive means to deal with all needs of meta knowledge processing, but it retains upward compatibility with existing usage of the language and corresponding tools and methods, which is a major concern for Semantic Web approaches. The following def-

inition of RDF<sup>+</sup> helps us to draw this line very clearly and concisely. First we briefly describe existing foundations of RDF in Section 4.1, then the abstract syntax for this embedded language, RDF<sup>+</sup>, is given in Section 4.2 and its semantics in Section 4.3. Last we show how to embed RDF<sup>+</sup> in RDF with named graphs.

#### 4.1. Basic Definitions

RDF is a graph based knowledge representation language. The nodes in a graph are URIs, blank nodes (a kind of existentially quantified variables) or literals. Arcs between the nodes, labeled with URIs, represent their relationships. In the following definitions, we simplify the RDF graph model in order to come up with a more concise formal characterization like [17].

##### Definition 4.1 (RDF Terms, Triples, and Variables)

Let  $U$  be the set of URIs,  $L$  the set of RDF Literals and  $B$  the set of Blank Nodes as defined in [14].  $U$ ,  $L$  and  $B$  are pairwise disjoint. Let  $R = U \cup L \cup B$ . A statement is an RDF triple in  $R \times U \times R$ . If  $S = (s, p, o)$  is a statement,  $s$  is called the subject,  $p$  the predicate and  $o$  the object of  $S$ . We denote the union  $U \cup L \cup B$  by  $T$  (RDF terms). Assume additionally the existence of an infinite set  $V$  of variables disjoint from sets above.

##### Definition 4.2 (RDF Graph)

An RDF graph  $G$  is a set of statements. For every two RDF graphs  $G_1$  and  $G_2$  the sets of blank nodes used in  $G_1$  and in  $G_2$  are disjoint.

Named graphs [5,6] offer means to group a set of statements in a graph and to refer to this graph using a URI. This way information about the graph can be expressed in RDF using its name as subject or object:

##### Definition 4.3 (RDF Named Graph)

A named graph is a pair  $(u, G)$  of a URI  $u$ , called name, and an RDF graph  $G$ , called the extension.

Finally, RDF datasets can be defined as followed:

##### Definition 4.4 (RDF Dataset)

An RDF dataset  $D$  is a set  $\{(G_0), (u_1, G_1), \dots, (u_n, G_n)\}$ , denoted by  $name(D)$ , where each  $G_i$  is a graph and each  $u_i$  a URI.  $G_0$  is called the default graph of  $D$ , and each pair  $(u_i, G_i)$  is a named graph; define  $name(G_i)_D = u_i$  and  $gr(u_i) = G_i$ .

#### 4.2. An Abstract Syntax for RDF<sup>+</sup>

The abstract syntax of RDF<sup>+</sup> is based on the same building blocks as RDF:

- $U$  are Uniform Resource Identifiers (URIs).
- $L$  are all RDF literals.
- $G \subseteq U$  is the set of graph names.
- $P \subseteq U$  is the set of properties.

In addition, we must be able to refer to statements directly without use of reification. For this purpose, we exploit the notion of (internal) unique statement identifiers:

- $\Theta$  is a set of statement identifiers, which is disjoint from  $U$  and  $L$ .

Now, we may define RDF<sup>+</sup> literal statements that are placed in named graphs and have, in addition to RDF, a globally unique statement identity.

##### Definition 4.5 (RDF<sup>+</sup> Literal Statements)

The set of all RDF<sup>+</sup> literal statements,  $\mathfrak{S}$ , is defined as quintuples by:

$$\mathfrak{S} := \{(g, s, p, o, \theta) \mid g \in G, s \in U, p \in P, o \in U \cup L, \theta \in \Theta\}.$$

Thereby,  $\theta$  and  $(g, s, p, o)$  are keys such that there exists a bijection  $f_1$  with  $f_1(g, s, p, o) = \theta \wedge f_1^{-1}(\theta) = (g, s, p, o)$ . Moreover, we define the overloaded function  $f_5$  to return the complete quintuple given either  $\theta$  or  $(g, s, p, o)$ , i.e.  $f_5(\theta) := (g, s, p, o, \theta) =: f_5(g, s, p, o)$ , when  $f_1(g, s, p, o) = \theta$ .

The reader may note that we assume that  $f_1$  is fixed and given before any statement is defined. Furthermore, this definition of literal statements and the rest of this paper abstracts from RDF blank nodes in order to keep the formalization more concise. However, there is no conceptual problem in extending our treatment to blank nodes, too. The two statements of Graph  $G_1$  of Example 4.1 may now be represented in RDF<sup>+</sup> in the following way:

##### Example 4.3 Knowledge statements in RDF<sup>+</sup>

---


$$\mathfrak{S} \ni K \ni \{$$

$$(G_1, \text{JamesHendler}, \text{researchTopic}, \text{SemanticWeb}, \theta_1),$$

$$(G_1, \text{JamesHendler}, \text{affiliatedWith}, \text{RensselaerPI}, \theta_2)\}$$


---

Thereby, the exact form of statement identifiers in  $\Theta$  is up to the implementation, as they are only used for internal processing.

Having represented the literal interpretation of RDF statements in  $\text{RDF}^+$ , we may now address the representation of selected RDF statements as  $\text{RDF}^+$  meta knowledge. This is done using a structure of  $\text{RDF}^+$  meta knowledge statements,  $\mathfrak{M}$ , that is separate from the set of  $\text{RDF}^+$  literal statements:

**Definition 4.6 (RDF<sup>+</sup> Meta Knowledge Statements)**

Let  $\Gamma \subseteq P$  be the set of meta knowledge properties. Let  $\Omega_\gamma$ , with  $\gamma \in \Gamma$ , be sets providing possible value ranges for the meta knowledge property  $\gamma \in \Gamma$ . Then, the set of all  $\text{RDF}^+$  meta knowledge statements,  $\mathfrak{M}$ , is defined by:  $\mathfrak{M} := \{(\theta, \gamma, \omega) \mid \theta \in \Theta, \gamma \in \Gamma, \omega \in \Omega_\gamma\}$ .

The following example illustrates the target representation of the first two meta knowledge statements of graph  $G3$  from Example 4.3.

**Example 4.4** Meta knowledge statements in  $\text{RDF}^+$

---

$\mathfrak{M} \supseteq M \supseteq \{$   
 $(\theta_1, \text{mk:source}, \langle \text{www.rpi.edu/report.doc} \rangle).$   
 $(\theta_1, \text{mk:certainty}, "0.9")\}$

---

Together we may now define a  $\text{RDF}^+$  dataset.

**Definition 4.7 (RDF<sup>+</sup> Dataset)**

A  $\text{RDF}^+$  dataset  $D^+$  of literal statements and associated meta knowledge statements is a pair  $D^+ = (K, M)$  referring to a set of literal statements  $K \subseteq \mathfrak{S}$  and a set of meta knowledge statements  $M \subseteq \mathfrak{M}$ .

A (partial) example for such a dataset given by the pair  $(K, M)$  with definitions for  $K \subseteq \mathfrak{S}$  and  $M \subseteq \mathfrak{M}$  has been given in Example 4.3 and Example 4.4, respectively.

4.3. Semantics for  $\text{RDF}^+$

We now have an abstract syntax for representing RDF triples like “*JamesHendler researchTopic SemanticWeb*” as part of  $G1$  and meta knowledge statements like “*the source of the statement that James Hendler’s research topic is Semantic Web is found in the document <www.rpi.edu/report.doc>*”. However, such an abstract syntax may remain remarkably ambiguous if it cannot be linked to a formal semantics. Assume two meta knowledge statements:

$(\theta_1, \text{mk:source}, \langle \text{www.rpi.edu/draftReport.doc} \rangle)$   
 $(\theta_1, \text{mk:source}, \langle \text{www.rpi.edu/finalReport.doc} \rangle)$

For the same literal statement identified by  $\theta_1$ , the question may arise whether this means a disjunction, i.e. one of the two documents has provided the fact, or a conjunction, i.e. both documents have provided the fact, or a collective reading, i.e. the two documents together gave rise to the fact, or whether this situation constitutes invalid meta knowledge. In order to prevent such ambiguities we introduce a generic semantic framework for meta knowledge in  $\text{RDF}^+$ . However, the framework must also be able to reproduce the literal interpretations found in RDF. For the latter purpose, we first define a ‘standard’ model for a  $\text{RDF}^+$  dataset.

**Definition 4.8 (Standard Interpretation and Model)**

A standard interpretation  $I_s : \mathfrak{S} \rightarrow \{\top, \perp\}$  for a dataset  $D^+ = (K, M)$  assigns truth values to all statements<sup>2</sup> in  $K$ . A standard interpretation for  $K$  is a standard model for  $K$  if and only if it makes all statements in  $K$  become true. This is denoted by  $I_s \models_s (K, M)$ .

For instance, any standard model  $I_s$  for  $(K, M)$  in Example 4.3 would include  $(G1, \text{JamesHendler}, \text{researchTopic}, \text{SemanticWeb}, \theta_1)$  in its set of literal statements evaluating to  $\top$ . In order to address the level of meta knowledge we foresee an additional model layer that provides a different interpretation to each meta knowledge property.

**Definition 4.9 ( $\Gamma$ -Interpretation and Model)**

A  $\Gamma$ -interpretation  $I_\gamma : \mathfrak{S} \rightarrow \Omega_\gamma$  for a property  $\gamma \in \Gamma$  is a partial function mapping statements into the allowed value range of  $\gamma$ .

A  $\Gamma$ -interpretation  $I_\gamma$  is a  $\Gamma$ -model for  $(K, M)$  if and only if for all meta knowledge statements  $(\theta, \gamma, \omega) \in M$  where  $f_1(\theta) = (g, s, p, o)$  the value of the interpretation equals to  $\omega$ , i.e.  $I_\gamma((g, s, p, o, \theta)) = \omega$ . This is denoted by  $I_\gamma \models_\gamma (K, M)$ .

As an example, consider the literal statement  $(G1, \text{JamesHendler}, \text{researchTopic}, \text{SemanticWeb}, \theta_1)$  from Example 4.3, and the meta knowledge statement  $(\theta_1, \text{certainty}, 0.9)$  from Example 4.4. A  $\Gamma$ -interpretation  $I_{\text{certainty}}$ , that is a  $\Gamma$ -model, would map the literal statement onto  $0.9$ .

The ‘standard’ interpretation and  $\Gamma$ -interpretation may now be combined to define what an overall, un-

<sup>2</sup> Note that because  $f_1$  is fixed there are no two tuples  $(g, s, p, o, \theta_1), (g, s, p, o, \theta_2)$ , where  $\theta_1 \neq \theta_2$ . This implies that the standard interpretation is independent of the identifiers  $\theta_1, \theta_2$ .

ambiguous model is:

**Definition 4.10 (Meta Knowledge Interpretation)**

A meta knowledge interpretation  $\mathfrak{S}$  is a set including a standard interpretation  $I_s$  and the  $\Gamma$ -interpretations  $I_\gamma$  for all meta knowledge properties  $\gamma \in \Gamma$ .

**Definition 4.11 (Meta Knowledge Model)**

A meta knowledge interpretation  $\mathfrak{S}$  is a model for a RDF<sup>+</sup> dataset  $D^+ = (K, M)$  if and only if all its interpretations  $I \in \mathfrak{S}$  are a standard model or  $\Gamma$ -models for  $(K, M)$ . This is denoted by  $\mathfrak{S} \models (K, M)$ .

4.4. Mapping between RDF and RDF<sup>+</sup>

The mapping between RDF and RDF<sup>+</sup> needs to be defined in two directions. First, one must be able to map from RDF, as given in the examples from Section 2, to RDF<sup>+</sup>. Second, one must be able to map from RDF<sup>+</sup> to RDF. Because RDF<sup>+</sup> is more fine-grained than RDF the first direction will be easy. For the second a compromise on the granularity of the representation has to be made.

4.4.1. From RDF to RDF<sup>+</sup>

The examples of Section 2 reify groups of statements, i.e. the ones found in  $G1$  and  $G2$ , in order to associate meta knowledge, such as given in  $G3$  and  $G4$ . In order to allow for an interpretation of the meta knowledge as defined in the preceding section, we map RDF into RDF<sup>+</sup>. For all RDF statements, including statements in graphs  $G1$  and  $G2$  of Example 4.1, the mapping performed is close to an identity mapping. One only needs to add statement identifiers. The result for  $G1$  in RDF<sup>+</sup> is:

**Example 4.5 Representing graph  $G1$  in RDF<sup>+</sup>**

---

$K \supseteq \{$   
 $(G1, \text{JamesHendler}, \text{researchTopic}, \text{SemanticWeb}, \theta_1),$   
 $(G1, \text{JamesHendler}, \text{affiliatedWith}, \text{RensselaerPI}, \theta_2)\}$   
with  
 $\theta_1 := f_1(G1, \text{JamesHendler}, \text{researchTopic}, \text{SemanticWeb})$   
and  
 $\theta_2 := f_1(G1, \text{JamesHendler}, \text{affiliatedWith}, \text{RensselaerPI})$

---

The same mapping – close to the identity mapping – is performed for meta knowledge statements like statements of graph  $G3$ , resulting in their representation as literal statements:

**Example 4.6 Representing graph  $G3$  in RDF<sup>+</sup>**

---

$K \supseteq \{$   
 $(G3, G1, \text{mk:source}, \langle \text{www.rpi.edu/report.doc} \rangle, \theta_3),$   
 $(G3, G1, \text{mk:certainty}, \langle \text{"0.9"} \rangle, \theta_4),$   
 $\dots \}$

---

Note that this step is necessary in order to achieve upward and – limited – downward compatibility between RDF<sup>+</sup> and RDF. The interpretation of statements, like the ones found in  $G3$ , also require an interpretation as meta knowledge. This is achieved by mapping RDF statements with designated properties from  $\Gamma$  like  $\text{mk:source}$  and  $\text{mk:certainty}$  to the additional meta knowledge layer:

**Example 4.7 Meta Knowledge representation of statements in  $G3$**

---

$M \supseteq \{$   
 $(\theta_1, \text{mk:certainty}, \langle \text{"0.9"} \rangle),$   
 $(\theta_1, \text{mk:source}, \langle \text{www.rpi.edu/report.doc} \rangle),$   
 $\dots \}$

---

The mapping of predicates of these meta knowledge statements from RDF to RDF<sup>+</sup> is obvious, they are mapped to itself. Objects are mapped to the corresponding elements of the value ranges  $\Omega_\gamma$ . For the subjects, however, there arise modeling choices. For instance, if  $\text{mk:certainty}$  were interpreted using probability theory, one might assign a distributive or a collective reading. In the distributive reading, each fact in  $G1$  receives the probability value of 0.9 and, eventually, the distributive reading will assign a joint probability of close to 0 for a large number of  $n$  stochastically independent facts, i.e. the joint probability  $0.9^n$ . In the collective reading, the collection of facts in  $G1$  as a whole will receive the probability value 0.9. Therefore, the collective reading will assign an individual certainty close to 1 for each individual fact, when the number of facts is high and each fact is independent from the others, i.e. the individual probability would be  $\sqrt[n]{0.9}$ . A priori, none of the two (and more) modeling choices is better than the other, but they constitute different modeling targets.

The mapping from RDF to RDF<sup>+</sup> for the *distributive reading* of a meta property  $\gamma$  is easy to achieve.

**Definition 4.12 (Distributive Embedding)**

Given an RDF statement “ $G \{s p o\}$ ” and an RDF meta knowledge statement “ $H \{G \gamma \omega\}$ ”, a distributive embedding of RDF<sup>+</sup> in RDF adds the meta knowledge statement

$\{(\theta, \gamma, \omega) \mid \theta = f_1(G, s, p, o) \wedge f_5(\theta) \in K\}$   
to  $M$ .

This means that such a meta knowledge statement is applied individually to all statements in the graph to which it refers in RDF, as indicated in the example above. For certain  $\gamma$  there might be several RDF meta knowledge statements  $H \{G \gamma \omega_i\}$  which attach different values  $\omega_i$  to a graph  $G$  via a single meta knowledge property  $\gamma$ . In that case set-valued ranges have to be elements of  $\Omega_\gamma$  in order to be consistent with Definition 4.9.

#### 4.4.2. From RDF<sup>+</sup> to RDF

The serialization of RDF<sup>+</sup> data in an RDF knowledge base is straightforward. Each quintuple  $(g, s, p, o, \theta)$  is realized as a corresponding triple in a named graph and the tuple identifier  $\theta$  is discarded.

**Example 4.8** *Serialization of some RDF<sup>+</sup> statements in G5 in RDF*

---

```
(G5, JamesHendler, researchTopic, SemanticWeb,  $\theta$ )
```

-is mapped to-

---

```
G5 {JamesHendler researchTopic SemanticWeb}
```

---

For meta knowledge statements the situation is more challenging, because literal statements with different statement identifiers may belong to only one named graph. Their corresponding meta knowledge statements may differ, but the realization of the meta knowledge statements in RDF does not allow for retaining these fine-grained distinctions – unless one chooses to change the modeling approach drastically, e.g. by assigning each literal statement to a named graph of its own, which is undesirable (cf. discussion in Section 3).

We have preferred to pursue a more conventional modeling strategy for RDF with named graphs. Therefore, we weaken the association between meta knowledge statements and their corresponding literal statements when mapping to RDF, i.e. we group sets of meta knowledge property values into one complex value.

#### Definition 4.13 (Grouped Meta Knowledge)

Given an RDF<sup>+</sup> dataset  $(K, M)$ , RDF meta knowledge is generated by grouping RDF<sup>+</sup> meta knowledge statements as follows:

Add the triple  $(g \gamma \omega')$  to the RDF graph  $g' := \text{hashGraph}(g)$  for each

$$\omega' := \omega_1 \vee_\gamma \dots \vee_\gamma \omega_n,$$

where  $(\theta, \gamma, \omega_i) \in M \wedge (g, S, P, O, \theta) \in K$ . Further, *hashGraph* is a function mapping existing graph names onto graph names suitable for associating meta knowledge and  $\vee_\gamma$  is an operation defined on  $\Omega_\gamma$ .

If  $\omega'$  is set-valued then a set of triples is added to  $g'$  in order to represent  $\omega'$ . The suitability of *hashGraph* may be application specific. A general strategy may map graph names  $g$  to graph names prefixed by `<http://metaknowledge.semanticweb.org>` in a deterministic manner. Operations on meta knowledge properties are discussed in section 5.2. In the following example the grouping of meta knowledge values is illustrated.

#### Example 4.9 Group of meta knowledge values

---

```
K:= {
(G5, JamesHendler, researchTopic, SemanticWeb,  $\theta_1$ ).
(G5, JamesHendler, affiliatedWith, UnivMaryland,  $\theta_2$ )}
M:= {
( $\theta_1$ , mk:source, <www.rpi.edu/report.doc>).
( $\theta_2$ , mk:source, <www.cs.umd.edu/survey.pdf>)}
```

-is mapped to-

---

```
G5 { JamesHendler researchTopic SemanticWeb,
      JamesHendler affiliatedWith UnivMaryland}
G6 { G5 mk:source <www.rpi.edu/report.doc>,
      <www.cs.umd.edu/survey.pdf>}
```

---

In Example 4.9, the resulting grouped value is the set consisting of the two documents `<report.doc>` and `<survey.pdf>` which is represented by two triples. For specific meta knowledge properties, an additional function may be necessary to provide a mechanism for representing grouped values in an appropriate RDF data structure.

## 5. SPARQL for RDF and Meta Knowledge

For querying RDF repositories with meta knowledge awareness, we exploit the capabilities of the SPARQL query language. In this section we first introduce a small extension to standard SPARQL syntax [20] and then define how SPARQL can be applied to an RDF<sup>+</sup> knowledge base. The objective of our considerations is the derivation of meta knowledge about query results.

### 5.1. SPARQL Syntax Revisited

In our scenario, we assume that the user aims to explore the knowledge and meta knowledge using the RDF query language SPARQL. The following

SPARQL query shown in Example 5.1 enables the user to find experts in the domain of Semantic Web and their affiliations evaluated on the RDF dataset presented in Example 4.1.

**Example 5.1** *SPARQL query to be evaluated on the RDF knowledge base in Example 4.1*

---

```

CONSTRUCT {?x worksAt ?z}
FROM G1
FROM G2
WHERE { GRAPH ?g {?x affiliatedWith ?z .
                  ?x researchTopic 'SemanticWeb'}}

```

---

When using SPARQL to query RDF<sup>+</sup> we introduce one additional expression “WITH META *MetaList*”. This expression includes the named graphs specified in *MetaList* for treatment as meta knowledge. This statement is optional. When it is present the SPARQL processor may digest the RDF<sup>+</sup> meta knowledge statements derivable from the RDF named graphs appearing in the *MetaList*. The SPARQL processor will then use this meta knowledge to compute and output all the meta knowledge statements derivable by successful matches of RDF<sup>+</sup> literal statements with the WHERE pattern.

According to SPARQL semantics, FROM *g* expressions replicate all RDF triples of *g* into the default triple space of the query. When the same triple appears in multiple source graphs (say  $(s, p, o)$  appears in both *G1* and *G2*, and the query contains a clause FROM *G1* FROM *G2*), its meta knowledge may become ambiguous. As an intermediate step of our query processing, this meta knowledge is aggregated using RDF<sup>+</sup> interpretations as introduced in Section 4.3. We note that this intermediate step is not necessary for queries with WITH NAMED *g* clauses, which are also fully compatible with our framework.

Thus, SPARQL queries on RDF<sup>+</sup> have one of the two following overall forms:

**Definition 5.1 (SPARQL SELECT Query)**

*The structure of a SPARQL SELECT query has the following form:*

```

SELECT SelectExpression
(WITH META MetaList)?
(FROM GraphName)+
WHERE P

```

**Definition 5.2 (SPARQL CONSTRUCT Query)**

*The structure of a SPARQL CONSTRUCT query has the following form:*

```

CONSTRUCT ConstructExpression
(WITH META MetaList)?
(FROM GraphName)+
WHERE P

```

In these definitions, *P* refers to a graph pattern that explains how RDF<sup>+</sup> literal statements from named graphs specified using FROM statements are matched. Matches bind variables that are used for providing results according to the *SelectExpression* or the *ConstructExpression*.

5.2. SPARQL Semantics Revisited

In this subsection we define the semantics of SPARQL queries evaluated on an RDF<sup>+</sup> theory. For our definitions we use two building blocks: algebraic semantics of SPARQL [15,17] and the *how-provenance* calculated via annotated relations (cf. [10]).

The algebraic semantics of SPARQL queries are given based on set-theoretic operations for sets of variable assignments. Such a set of assignments may be assigned information about the so called *how-provenance* [10], i.e. the assignments may be annotated with formulae describing the individual derivation tree used to assign the variables. The how-provenance annotation may be represented by a function  $\Phi : (U \cup L)^{|V|} \rightarrow \mathfrak{F}$ , where  $(U \cup L)^{|V|}$  is the set of all tuples of the length  $|V|$  over the domain  $U \cup L$  and  $\mathfrak{F}$  is the set of formulae annotating variable assignments. The set of formulae  $\mathfrak{F}$  is given by all Boolean formulas constructed over the set of literal statements  $\mathfrak{S}$  and including a bottom element  $\perp$  and a top element  $\top$ . The formulae constitute an algebra  $(\mathfrak{F}, \wedge, \vee, \neg, \perp, \top)$ . The special element  $\perp$  is used as annotation of variable assignments which are not in the result set. The special element  $\top$  may be omitted, but it allows for simplification of complex formulas.

The following definition shows how a set of variable bindings is generalized to SPARQL queries of arbitrary complexity by a recursive definition of simultaneous query evaluation and computation of the annotations. The first step in evaluating a graph pattern is to find matches for the triple pattern contained in the query. Because the RDF<sup>+</sup> dataset  $D^+$  consists of quintuples, we need to adapt the SPARQL evaluation procedures. The statement identifiers do not need to be matched, as they depend functionally on graph name, subject, predicate and object. Therefore, we consider that the matching of triple patterns  $P = (\alpha, \beta, \delta)$ , given a graph name  $\lambda \in U$ , is always defined by the following SPARQL grammar (GRAPH  $\lambda P$ ). As a simplification of our formalization we assume that the keyword GRAPH together with a URI or a graph variable is used in any

given SPARQL query. If it is not used, we may expand a given SPARQL query to include it.

**Definition 5.3 (Basic Graph Pattern Matching)**

Given a basic graph pattern  $P = (\alpha\beta\delta)$ , a graph name  $\lambda \in U$  and a mapping  $\mu$  such that  $\text{var}(P) \subseteq \text{dom}(\mu)$ . Let  $D^+$  be an RDF<sup>+</sup> dataset,  $G = \text{gr}(\lambda)$  a RDF graph in  $D^+$ , and  $\mu(P)$  is the set of triples obtained by replacing the variables in the triples of  $P$  according to  $\mu$ . The (meta) evaluation of  $P$  over  $G$ , denoted by  $\llbracket P \rrbracket_G^{D^+}$  is defined as the set of annotated mappings  $\Phi$ ,  $\mu \subseteq \text{dom}(\Phi)$  :

$$\begin{aligned} \llbracket P \rrbracket_G^{D^+} = \{ & \mu \circ (\Phi(\mu)) \mid \mu : V \rightarrow U \cup L \wedge \\ & \Phi : (U \cup L)^{|V|} \rightarrow \tilde{\mathcal{F}} \wedge \\ & \text{dom}(\mu) = \text{var}(P) \wedge \\ & \exists \theta (\text{name}(G) \circ \mu(P) \circ (\theta) \in K_G) \} \end{aligned}$$

where  $K_G := \sigma_{\#1=\text{name}(G)}(K)$ , i.e. the selection of statements in  $K$  belonging to graph  $G$  and  $\circ$  denotes vector concatenation.

If  $\mu \circ (\Phi(\mu)) \in \llbracket P \rrbracket_G^{D^+}$ , we say that  $\mu \circ (\Phi(\mu))$  is a solution for  $P$  in  $G$ .

$$\Phi(\mu) = \begin{cases} \theta & \text{if } \mu(P) = (g, s, p, o) \wedge \\ & (g, s, p, o, \theta) \in K_G \wedge \\ & f_1(g, s, p, o) = \theta, \\ \perp & \text{else} \end{cases}$$

Assume, for the next SPARQL queries example, the following RDF<sup>+</sup> knowledge base  $D^+$ :

**Example 5.2** RDF<sup>+</sup> dataset  $D^+ = (K, M)$

---

$K = \{$   
 (G1, JamesHendler, researchTopic, SemanticWeb,  $\theta_1$ ),  
 (G1, JamesHendler, affiliatedWith, RensselaerPI,  $\theta_2$ ),  
 (G2, JamesHendler, researchTopic, Robotics,  $\theta_3$ ),  
 (G2, JamesHendler, affiliatedWith, UnivMaryland,  $\theta_4$ ),  
 (G2, RudiStuder, researchTopic, SemanticWeb,  $\theta_5$ ),  
 (G2, RudiStuder, affiliatedWith, UnivKarlsruhe,  $\theta_6$ )  
 $\}$

For corresponding  $M$ , see Example 5.14

---

Let the SPARQL query in Example 5.3 be evaluated on the dataset  $D^+$ :

**Example 5.3** SPARQL query to be evaluated on the knowledge base  $D^+$

---

SELECT ?g ?x ?y

---

```
FROM G1
FROM G2
WHERE {
  GRAPH ?g {?x researchTopic ?y} }
```

---

For the query of Example 5.3, we may find the following variable assignments in Example 5.4 using standard SPARQL processing and we may indicate, which atomic formulae, i.e. RDF<sup>+</sup> quintuples in this simple example, led to these variable assignments. This indication is given by the statement identifiers representing their statements.

**Example 5.4** Variable assignments for query of Example 5.3

?g	?x	?y	$\tilde{\mathcal{F}}$
G1	JamesHendler	SemanticWeb	$\theta_1$
G2	JamesHendler	Robotics	$\theta_3$
G2	RudiStuder	SemanticWeb	$\theta_5$

Basic pattern matching is not directly applicable, if an expression 'GRAPH  $\lambda$ ' appears outside a complex triple pattern. In such a case, we first need to distribute the expression 'GRAPH  $\lambda$ ' appropriately to atomic triple patterns in order to prescribe atomic SPARQL expressions accessible by basic pattern matching. Because named graphs cannot be nested, this distribution is always possible and unambiguous.

In the following we use the function  $\text{quads}(P)$  to denote the query resulting from this transformation. In example 5.5 this transformation is demonstrated on a conjunction of two triple patterns.

**Example 5.5** Distributing the expression 'GRAPH  $\lambda$ ' in a complex triple pattern with the function  $\text{quads}(P)$

---

$P_1 =$   
 GRAPH ?src {  
 { ?x researchTopic ?y .}  
 { ?x affiliatedWith ?z .}  
 $\text{quads}(P_1) =$   
 GRAPH ?src { ?x researchTopic ?y .}  
 GRAPH ?src { ?x affiliatedWith ?z .}

---

Now we define the evaluation of complex graph patterns by operations on sets of variable assignments similar to [15,17].

**Definition 5.4 (Complex Graph Pattern Matching)**

Given the basic triple graph patterns  $P, P_1, P_2$  and

a graph name  $\lambda \in U$ . Let  $D^+$  be an RDF<sup>+</sup> dataset,  $G = gr(\lambda)$  a RDF graph in  $D^+$ . The meta evaluation of  $P$ ,  $P_1$ , and  $P_2$  over  $G$ , denoted by  $\llbracket \cdot \rrbracket_G^{D^+}$ , is defined recursively as follow:

- $\llbracket P \rrbracket_G^{D^+}$  is given by definition 5.3,
- $\llbracket P_1 \text{ AND } P_2 \rrbracket_G^{D^+} = \llbracket P_1 \rrbracket_G^{D^+} \bowtie \llbracket P_2 \rrbracket_G^{D^+}$ ,
- $\llbracket P_1 \text{ OPT } P_2 \rrbracket_G^{D^+} = \llbracket P_1 \rrbracket_G^{D^+} \bowtie \llbracket P_2 \rrbracket_G^{D^+}$ ,
- $\llbracket P_1 \text{ UNION } P_2 \rrbracket_G^{D^+} = \llbracket P_1 \rrbracket_G^{D^+} \cup \llbracket P_2 \rrbracket_G^{D^+}$ ,
- $\llbracket P_1 \text{ FILTER } C \rrbracket_G^{D^+} = \sigma_c(\llbracket P_1 \rrbracket_G^{D^+})$ .

The definition uses the operation AND. In standard SPARQL the operation AND is denoted by the absence of an operator. Like [15,17] we still use the explicit term AND in order to facilitate referencing to this operator.

The recursion in the SPARQL query evaluation defined in Definition 5.3 is indeed identical to [15,17] (the work of [15,17] on the SPARQL semantics is briefly summarized in Appendix A). Only the basic pattern matching has been changed slightly. Basic pattern matching now considers the basic pattern triple with named graph matching, and it annotates variable assignments  $\Phi(\mu)$  from basic matches with atomic statements from  $\mathfrak{S}$  and variable assignments from complex matches with Boolean formulae  $F \in \mathfrak{F}$  over these statements.

The following definition specifies how complex formulae from  $\mathfrak{F}$ , which are used to construct formulas representing the how-provenance and serve as annotations for results of matching complex graph pattern, are derived using algebraic operations on sets of annotated variables bindings.

### Definition 5.5 (Algebra of Annotated Relations)

Let  $\Phi$ ,  $\Phi_1$  and  $\Phi_2$  be sets of annotated variable assignments. We define  $\bowtie$ ,  $\bowtie$ ,  $\cup$ ,  $\setminus$ ,  $\sigma$ , and  $\psi$  via operations on the annotations of the assignments (two binary operations  $\wedge$  and  $\vee$ , a unary operation  $\neg$ ) as following:

- $(\Phi_1 \bowtie \Phi_2)(\mu) = \Phi_1(\mu_1) \wedge \Phi_2(\mu_2)$ , where  $\forall x \in \text{dom}(\mu_1) \cap \text{dom}(\mu_2) : \mu_1(x) = \mu_2(x)$  and where  $\mu_1$  and  $\mu_2$  are compatible and  $\mu = \mu_1 \cup \mu_2$ <sup>3</sup>,
- $(\Phi_1 \cup \Phi_2)(\mu) = \Phi_1(\mu) \vee \Phi_2(\mu)$ ,  $\mu \in \Phi_1$  or  $\mu \in \Phi_2$
- $(\Phi_1 \setminus \Phi_2)(\mu) = \Phi_1(\mu) \wedge \neg \left( \bigvee_{\mu_i, \Phi_2(\mu_i) \neq \perp} \Phi_2(\mu_i) \right)$ , where  $\forall x \in \text{dom}(\mu_i) \cap \text{dom}(\mu) : \mu_i(x) = \mu(x)$ .
- $(\Phi_1 \bowtie \Phi_2)(\mu) = (\Phi_1 \bowtie \Phi_2)(\mu) \vee (\Phi_1 \setminus \Phi_2)(\mu)$ .
- $(\sigma_c(\Phi(\mu))) = \Phi(\mu) \wedge f_c(\mu)$ , where  $f_c(\mu)$  denotes a function mapping  $\mu$  to either  $\top$  or  $\perp$  according to the

<sup>3</sup> Two mappings  $\mu_1$  and  $\mu_2$  are compatible when for all  $x \in \text{dom}(\mu_1) \cap \text{dom}(\mu_2)$ , it is the case that  $\mu_1(x) = \mu_2(x)$ , then  $\mu_1 \cup \mu_2$  is also a mapping.

condition c.

- $(\psi_x(\Phi(\mu))) = \bigvee_{\mu \subseteq \nu, \text{dom}(\mu)=X, \Phi(\nu) \neq \perp} \Phi(\nu)$ .

Furthermore, following the algebraic definition of how-provenance formula, we present the set of laws for annotated relations.

### Definition 5.6 (Laws of Annotated Relations)

Let  $\Phi_1$ ,  $\Phi_2$  and  $\Phi_3$  be sets of annotated variable assignments. We define the laws for annotation relations via operations on the annotations of the assignments as following:

- *Idempotency:*  
 $\Phi_1(\mu_1) \wedge \Phi_1(\mu_1) = \Phi_1(\mu_1)$   
 $\Phi_1(\mu_1) \vee \Phi_1(\mu_1) = \Phi_1(\mu_1)$
- *Commutativity:*  
 $\Phi_1(\mu_1) \wedge \Phi_2(\mu_2) = \Phi_2(\mu_2) \wedge \Phi_1(\mu_1)$   
 $\Phi_1(\mu_1) \vee \Phi_2(\mu_2) = \Phi_2(\mu_2) \vee \Phi_1(\mu_1)$
- *Associativity:*  
 $\Phi_1(\mu_1) \wedge (\Phi_2(\mu_2) \wedge \Phi_3(\mu_3)) = (\Phi_1(\mu_1) \wedge \Phi_2(\mu_2)) \wedge \Phi_3(\mu_3)$   
 $\Phi_1(\mu_1) \vee (\Phi_2(\mu_2) \vee \Phi_3(\mu_3)) = (\Phi_1(\mu_1) \vee \Phi_2(\mu_2)) \vee \Phi_3(\mu_3)$
- *Distributivity:*  
 $\Phi_1(\mu_1) \wedge (\Phi_2(\mu_2) \vee \Phi_3(\mu_3)) = (\Phi_1(\mu_1) \wedge \Phi_2(\mu_2)) \wedge (\Phi_1(\mu_1) \wedge \Phi_3(\mu_3))$   
 $\Phi_1(\mu_1) \vee (\Phi_2(\mu_2) \wedge \Phi_3(\mu_3)) = (\Phi_1(\mu_1) \vee \Phi_2(\mu_2)) \vee (\Phi_1(\mu_1) \vee \Phi_3(\mu_3))$
- *Complement:*  
 $\Phi_1(\mu_1) \wedge \neg \Phi_1(\mu_1) = \perp$   
 $\Phi_1(\mu_1) \vee \neg \Phi_1(\mu_1) = \top$
- *Absorption:*  
 $\Phi_1(\mu_1) \wedge (\Phi_1(\mu_1) \vee \Phi_2(\mu_2)) = \Phi_1(\mu_1)$   
 $\Phi_1(\mu_1) \vee (\Phi_1(\mu_1) \wedge \Phi_2(\mu_2)) = \Phi_1(\mu_1)$
- *De Morgan's:*  
 $\neg(\Phi_1(\mu_1) \wedge \Phi_2(\mu_2)) = \neg \Phi_1(\mu_1) \vee \neg \Phi_2(\mu_2)$   
 $\neg(\Phi_1(\mu_1) \vee \Phi_2(\mu_2)) = \neg \Phi_1(\mu_1) \wedge \neg \Phi_2(\mu_2)$

In order to show the evaluation of a SPARQL query, we consider the query from Example 5.6 to be evaluated on the knowledge base  $D^+$ .

### Example 5.6 SPARQL query to be evaluated on the knowledge base $D^+$

```
SELECT ?h1 ?h2 ?x ?y
FROM G1
FROM G2
WHERE {
  {GRAPH ?h1 {?x affiliatedWith ?y}} AND
  {GRAPH ?h2 {?x researchTopic 'SemanticWeb'}}
```

FILTER {?x='JamesHendler'}

Let  $P$  be the graph pattern contained in the WHERE clause of the query. Then the evaluation of  $P$  is defined by an algebraic expression:

$$\begin{aligned}
\llbracket P \rrbracket_G^{D^+} &= \llbracket \{P_1 \text{ AND } P_2\} \rrbracket_G^{D^+} \\
&= \text{FILTER } \{?x = \text{JamesHendler}\} \llbracket P_1 \text{ AND } P_2 \rrbracket_G^{D^+} \\
&= \sigma_{?x=\text{JamesHendler}}(\llbracket P_1 \text{ AND } P_2 \rrbracket_G^{D^+}) \\
&= \sigma_{?x=\text{JamesHendler}}(\llbracket P_1 \rrbracket_G^{D^+} \bowtie \llbracket P_2 \rrbracket_G^{D^+}) \\
&= \sigma_{?x=\text{JamesHendler}}(\Phi_1 \bowtie \Phi_2)
\end{aligned}$$

where  $\Phi_1$  and  $\Phi_2$  are relations representing variable assignments and their annotations. In this example and in the preceding definition we have used algebraic operations on sets of annotated bindings. In order to evaluate the expression  $\sigma_{?x=\text{JamesHendler}}(\Phi_1 \bowtie \Phi_2)$  we need to determine  $\Phi_1$  and  $\Phi_2$  using Definition 5.3. The intermediate result is shown in Example 5.7.

**Example 5.7** Determining  $\Phi_1$  and  $\Phi_2$  - intermediate result of the expression  $\sigma_{?x=\text{JamesHendler}}(\Phi_1 \bowtie \Phi_2)$

$$\Phi_1 =$$

?h1	?x	?y	$\Xi_1$
G1	JamesHendler	RensselaerPI	$\theta_2$
G2	JamesHendler	UnivMaryland	$\theta_4$
G2	RudiStuder	UnivKarlsruhe	$\theta_6$

$$\Phi_2 =$$

?h2	?x	$\Xi_2$
G1	JamesHendler	$\theta_1$
G2	RudiStuder	$\theta_5$

To evaluate the conjunction of two patterns, the operation  $\bowtie$  is applied, the result is shown in Example 5.8. The annotation  $\theta_1 \wedge \theta_2$  of the first row represents that this assignment has been derived from the conjunction of the two literal statements  $\theta_1$  and  $\theta_2$  (see Example 5.2).

**Example 5.8** Determining  $\Phi_1 \bowtie \Phi_2$  - intermediate result of the expression  $\sigma_{?x=\text{JamesHendler}}(\Phi_1 \bowtie \Phi_2)$

$$\Phi_1 \bowtie \Phi_2 =$$

?h1	?h2	?x	?y	$\Xi_3$
G1	G1	JamesHendler	RensselaerPI	$\theta_1 \wedge \theta_2$
G1	G2	JamesHendler	UnivMaryland	$\theta_1 \wedge \theta_4$
G2	G2	RudiStuder	UnivKarlsruhe	$\theta_5 \wedge \theta_6$

Application of the  $\sigma$ -operation to the intermediate results gives the annotated relation shown in Example 5.9.

**Example 5.9** Result of  $\sigma_{?x=\text{JamesHendler}}(\Phi_1 \bowtie \Phi_2)$

$\sigma_{?x=\text{JamesHendler}}(\Phi_1 \bowtie \Phi_2) =$

?h1	?h2	?x	?y	$\Xi_4$
G1	G1	JamesHendler	RensselaerPI	$(\theta_1 \wedge \theta_2) \wedge \top$
G1	G2	JamesHendler	UnivMaryland	$(\theta_1 \wedge \theta_4) \wedge \top$
G2	G2	RudiStuder	UnivKarlsruhe	$\perp$

The annotations  $\Phi(\mu)$  can now be used to assign truth values for each tuple binding  $\mu$ .  $I_s$  (cf. Definition 4.8) assigns truth values to all atomic statements  $s_i \in K \subseteq \Xi$ . We extend the interpretation  $I_s$  to capture all the Boolean formulae over statements  $\Xi$ .

**Definition 5.7 (Standard Interpretation of Formulae)**

Let  $F, F_1, F_2 \in \mathfrak{F}$  be Boolean formulae over  $\Xi$ , let  $F_a \in \Xi$  be an atomic formula. We define the standard interpretation of formulae  $I_s^f$  as follows:

- $I_s^f(F_a) := I_s(F_a)$ ;
- $I_s^f(\neg F) := \perp$  if  $I_s^f(F) = \top$ ;  $I_s^f(\neg F) := \top$  if  $I_s^f(F) = \perp$ ;
- $I_s^f(F_1 \wedge F_2)$  is  $\top$  if  $I_s^f(F_1) = I_s^f(F_2) = \top$ , otherwise  $\perp$
- $I_s^f(F_1 \vee F_2)$  is  $\top$  if  $I_s^f(F_1) = \top$  or  $I_s^f(F_2) = \top$ , otherwise  $\perp$ .

For instance,  $I_s^f$  returns  $\top$  for the assignment shown in the first row of  $\Phi_1 \bowtie \Phi_2$  from Example 5.8, because the statements  $\theta_1$  and  $\theta_2$  are in the knowledge base.

Basically, the standard interpretation of formulae  $I_s^f$  assigns truth values for all variable bindings that are in the knowledge base, as well as capture all the Boolean formulae over these statements bindings. For producing the result set, we dismiss all variable bindings with  $\perp$  value assigned.

Analogously to  $I_s^f$ , we can extend a  $\Gamma$ -interpretation  $I_\gamma$  over RDF<sup>+</sup> statements to a  $\Gamma$ -interpretation  $I_\gamma^f$  over formulae. Similar to definition the standard interpretation of formulae, the  $\Gamma$ -interpretation assigns values which represent meta knowledge properties for all variable bindings which are not in a set of bindings annotated with  $\perp$ . Thus we do not consider meta knowledge interpretations of such bindings over formulae. Remember that meta knowledge interpretations allow for only one  $\omega$  per  $\theta \in \Theta$  and  $\gamma \in \Gamma$  (cf. Definition 4.9). In order to make use of the how-provenance represented by the annotations we require that for each meta knowl-

edge property  $\gamma$  an algebra  $(\Omega_\gamma, \wedge_\gamma, \vee_\gamma, \neg_\gamma, \top_\gamma, \perp_\gamma)$  with three operations  $\wedge_\gamma, \vee_\gamma, \neg_\gamma$  and two special elements  $\top_\gamma, \perp_\gamma \in \Omega_\gamma$  is defined. The definition of the algebras can be supplied by a modeler according to the intended semantics of the different meta knowledge properties.

**Definition 5.8 ( $\Gamma$ -Interpretation of Formulae)**

Let  $F, F_1, F_2 \in \mathfrak{F}$  be Boolean formulae over  $\mathfrak{S}$ , let  $F_a \in \mathfrak{S}$  be an atomic formula. We define the interpretation  $I_\gamma^f$  as follows:

- $I_\gamma^f(F_a) := I_\gamma(F_a)$ ;
- $I_\gamma^f(\neg F)$  is  $\neg_\gamma I_\gamma^f(F)$ ;
- $I_\gamma^f(F_1 \wedge F_2)$  is  $I_\gamma^f(F_1) \wedge_\gamma I_\gamma^f(F_2)$ ;
- $I_\gamma^f(F_1 \vee F_2)$  is  $I_\gamma^f(F_1) \vee_\gamma I_\gamma^f(F_2)$ ;

The definition of knowledge dimensions can be supplied by the administrator of the knowledge base, according to the intended semantics of particular meta knowledge properties. For illustration we consider in Example 5.10 the definition of fuzzy logic operations to calculate a possibility measure on variable assignments, operations defined on timestamps which calculate the time of the last modification, and set operations defined for source documents that construct the combined 'how-provenance'.

**Example 5.10 Fuzzy logic operations for meta knowledge properties**

---


$$\begin{aligned} \Omega_{certainty} &= [0, 1] \\ I_{certainty}^f(x_1 \wedge x_2) &= \min(I_{certainty}^f(x_1), I_{certainty}^f(x_2)) \\ I_{certainty}^f(x_1 \vee x_2) &= \max(I_{certainty}^f(x_1), I_{certainty}^f(x_2)) \\ I_{certainty}^f(\neg x_1) &= 1 - I_{certainty}^f(x_1) \\ \Omega_{time} &= [0, \infty) \\ I_{time}^f(x_1 \wedge x_2) &= \max(I_{time}^f(x_1), I_{time}^f(x_2)) \\ I_{time}^f(x_1 \vee x_2) &= \min(I_{time}^f(x_1), I_{time}^f(x_2)) \\ I_{time}^f(\neg x_1) &= 0 \\ \Omega_{source} &= 2^D, \text{ D the set of document URIs} \\ I_{source}^f(x_1 \wedge x_2) &= I_{source}^f(x_1) \cap I_{source}^f(x_2) \\ I_{source}^f(x_1 \vee x_2) &= I_{source}^f(x_1) \cup I_{source}^f(x_2) \\ I_{source}^f(\neg x_1) &= \{\} \end{aligned}$$


---

As discussed earlier, the graph pattern queries evaluation with meta knowledge algorithm proposed in this paper follows a compositional semantics, i.e. for each  $P'$  sub-pattern of  $P$  the result of evaluating  $P'$  can be used to evaluate  $P$ . In order to solve this

problem algorithmically, we have implemented, analogously to [15,16], a depth-first strategy algorithm for evaluating SPARQL query patterns with meta knowledge. This algorithm, denoted by  $Eval^+$ , has been implemented so that it emulates exactly the recursive definition of  $\llbracket \cdot \rrbracket$  for well-designed pattern as defined in [19]. Formally, given an RDF<sup>+</sup> dataset  $D^+$ , we define the evaluation of pattern  $P$  with the set of mappings  $\Omega$ , denoted by  $Eval_{D^+}^+(P, \Omega)$  as follows:

---


$$\begin{aligned} Eval_{D^+}^+(P, \Omega) &= \\ &\text{if } \Omega = \emptyset \text{ then return } (\emptyset) \\ &\text{if } P \text{ is a triple pattern } t \text{ then return } (\Omega \bowtie \llbracket t \rrbracket_{D^+}) \\ &\text{if } P = (P1 \text{ AND } P2) \text{ then return} \\ &\quad Eval_{D^+}^+(P2, Eval_{D^+}^+(P1, \Omega)) \\ &\text{if } P = (P1 \text{ OPT } P2) \text{ then return} \\ &\quad Eval_{D^+}^+(P1, \Omega) \bowtie Eval_{D^+}^+(P2, Eval_{D^+}^+(P1, \Omega)) \\ &\text{if } P = (P1 \text{ FILTER } R) \text{ then return} \\ &\quad \{ \Phi(\mu) \in Eval_{D^+}^+(P1, \Omega) \wedge f_R(\mu) \} \\ &\quad \text{where } f_R(\mu) \text{ denotes a function mapping } \mu \text{ to either} \\ &\quad \top \text{ or } \perp \text{ according to the condition } R \end{aligned}$$


---

Then, the evaluation of  $P$  against a dataset  $D^+$ , which we denote simply by  $Eval_{D^+}^+(P)$ , is  $Eval_{D^+}^+(P, \Phi(\mu)_\emptyset)$ , where  $\Phi(\mu)_\emptyset$  is the mapping with empty domain.

5.3. SPARQL Query forms

As mentioned before, in this paper we are only interested in standard SPARQL SELECT and CONSTRUCT query forms. The evaluation of SPARQL queries on RDF<sup>+</sup> data differs in that meta knowledge is attached to the results.

The evaluation of SELECT queries on an RDF<sup>+</sup> dataset is based on  $\psi_X(\llbracket P \rrbracket_G^{D^+})$  (see Definition 5.5), where  $X$  denotes the set of variables specified in the *Select-Expression* (see Definition 5.1). If  $X$  forms a proper subset of the variables used in the graph pattern then the annotations of all bindings  $\nu$  are grouped. This grouping is analogous to the generation of grouped meta knowledge described in Definition 4.13. As an example consider the query shown in Example 5.11, which is a slight modification of the query from Example 5.6, applied to the data shown in Example 5.2. For the result see Example 5.12. In contrast to Example 5.8 there is only one row for *JamesHendler*.

**Example 5.11 SPARQL query to be evaluated on knowledge base  $D^+$**

---

```
SELECT ?x
WITH META G3, G4
FROM G1
FROM G2
WHERE {
```

{GRAPH ?h1 {?x affiliatedWith ?y}} AND  
 {GRAPH ?h2 {?x researchTopic 'SemanticWeb'}}}

---

The result of a SELECT query is a set of extended bindings. Such an extended binding contains values for the specified variables and values for each meta knowledge property  $\gamma \in \Gamma$  which can be regarded as additional variables.

**Example 5.12** Determining  $\psi_{\{?x\}}(\Phi_1 \bowtie \Phi_2)$  - intermediate result of the expression defined in Example 5.11

?x	$\Xi_5$
JamesHendler	$(\theta_1 \wedge \theta_2) \vee (\theta_1 \wedge \theta_4)$
RudiStuder	$\theta_5 \wedge \theta_6$

For each binding  $\mu$  of query from the Example 5.11, the variables  $\gamma$  are bound to  $I_\gamma^f(\psi_X(\mathbb{I}P\mathbb{I}_G^{D^+})(\mu_i))$  (see Example 5.13). For this result the meta knowledge from Example 5.14 has been used. For instance  $I_{certainty}^f((\theta_1 \wedge \theta_2) \vee (\theta_1 \wedge \theta_4)) = 0.9$ . If no meta knowledge statement  $(\theta, \gamma, \omega)$  exists for a particular RDF<sup>+</sup> literal statement  $f_5(\theta)$  and a particular meta knowledge property  $\gamma$  then  $\perp_\gamma$  serves as default value. For the result of a SELECT query all bindings from  $\psi_X(\mathbb{I}P\mathbb{I}_G^{D^+})$  are extended in this way.

**Example 5.13** Variable bindings for  $I_\gamma^f(\psi_X(\mathbb{I}P\mathbb{I}_G^{D^+})(\mu_i))$  - results of the expression defined in Example 5.11

?x	certainty	time
JamesHendler	"0.9"	"5/5/2007"
RudiStuder	"0.6"	"6/6/2001"

**Example 5.14** Meta knowledge statements,  $D^+ = (K, M)$

$M = \{$   
 $(\theta_1, \text{mk:certainty}, "0.9")$   
 $(\theta_1, \text{mk:time}, "5/5/2007")$   
 $(\theta_2, \text{mk:certainty}, "0.9")$   
 $(\theta_2, \text{mk:time}, "5/5/2007")$   
 $(\theta_3, \text{mk:certainty}, "0.6")$   
 $(\theta_3, \text{mk:time}, "6/6/2001")$   
 $(\theta_4, \text{mk:certainty}, "0.6")$   
 $(\theta_4, \text{mk:time}, "6/6/2001")$   
 $(\theta_5, \text{mk:certainty}, "0.6")$   
 $(\theta_5, \text{mk:time}, "6/6/2001")$   
 $(\theta_6, \text{mk:certainty}, "0.6")$   
 $(\theta_6, \text{mk:time}, "6/6/2001")\}$

---

Analogously to standard evaluation, the evaluation of a CONSTRUCT query on an RDF<sup>+</sup> dataset results in a single RDF<sup>+</sup> graph which is built using the graph template specified in the *ConstructExpression* (see Definition 5.2). This is in line with the fact that the graph template consists of a conjunction of triple patterns and a named graph thus quadruple patterns cannot be stated.<sup>4</sup>

Similar to the evaluation of SELECT queries, the evaluation of CONSTRUCT queries is based on  $\psi_X(\mathbb{I}P\mathbb{I}_G^{D^+})$ , where  $X$  denotes the set of variables specified in the *ConstructExpression*. The RDF<sup>+</sup> graph is constructed as described in the following:

Let  $t_j$  be the triple pattern  $j$  specified in the *ConstructExpression*,  $P$  denote the graph pattern specified in the WHERE-clause,  $(s_{i,j}, p_{i,j}, o_{i,j})$  denote the triple obtained by replacing the variables in  $t_j$  according to a mapping  $\mu_i$  and  $\hat{g}$  denote a new graph name. Then, for each binding  $\mu_i \in \psi_X(\mathbb{I}P\mathbb{I}_G^{D^+})$  and for each  $t_j$  the quintuple  $(\hat{g}, s_{i,j}, p_{i,j}, o_{i,j}, \theta_{i,j})$  is added to  $\Xi$ , where  $\theta_{i,j}$  is the statement identifier  $f_1(\hat{g}, s_{i,j}, p_{i,j}, o_{i,j})$ . Further  $(\theta_{i,j}, \gamma, \omega_{i,j})$  is added to  $M$ , where  $\omega_{i,j} = I_\gamma^f(\psi_X(\mathbb{I}P\mathbb{I}_G^{D^+})(\mu_i))$ .

Each new quintuple inherits the meta knowledge properties  $\gamma$  associated with the binding which has been used to create that quintuple. The value of  $\omega_{i,j}$  is determined by applying  $I_\gamma^f$  to the formula which annotates the binding. Note that since  $\psi_X(\mathbb{I}P\mathbb{I}_G^{D^+})$  and the interpretations  $I_\gamma^f$  are functions and further the graph template in *ConstructExpression* is a set of triples the meta knowledge properties  $(\theta_{i,j}, \gamma, \omega_{i,j})$  are unique for a given  $\theta_{i,j}$ .

As an example for a CONSTRUCT statement consider Example 5.15. Meta knowledge for some of the RDF<sup>+</sup> statements presented in Example 5.2 is specified in Example 5.14. For graph pattern  $P$  contained in this query the result of  $\psi_X(\mathbb{I}P\mathbb{I}_G^{D^+})$  is identical to the annotated relation shown in Example 5.8 except for the first two columns. Based on the single triple pattern (?x worksAt ?y) contained in the graph template and the two bindings contained in  $\psi_X(\mathbb{I}P\mathbb{I}_G^{D^+})$  two quintuples are constructed and added to the RDF<sup>+</sup> literal statements  $K_{res}$  as shown in Example 5.16.  $M_{res}$  contains the corresponding meta knowledge statements resulting from  $I_\gamma^f(\psi_X(\mathbb{I}P\mathbb{I}_G^{D^+})(\mu_i))$ .

<sup>4</sup> Standard SPARQL does not allow for giving this graph a name. In order to associate meta knowledge, multiple named graphs as outputs are convenient. In order to remain standard compliant, the SPARQL engine may however also return data and meta knowledge in two different batches distinguished by some implementation-specific mechanism.

### Example 5.15 SPARQL query CONSTRUCT

```
CONSTRUCT {?x worksAt ?y}
WITH META G3, G4
FROM G1
FROM G2
WHERE {
  (GRAPH ?h1 {?x affiliatedWith ?y}) AND
  (GRAPH ?h2 {?x researchTopic 'SemanticWeb'})}
```

### Example 5.16 Variable binding for query of the Example 5.15

```
Kres = {
(Gnew1, JamesHendler, worksAt, RensselaerPI,  $\theta_{new1}$ )
(Gnew2, JamesHendler, worksAt, UnivMaryland,  $\theta_{new2}$ )
(Gnew3, RudiStuder, worksAt, UnivKarlsruhe,  $\theta_{new3}$ )
Mres = {
( $\theta_{new1}$ , mk:certainty, "0.9")
( $\theta_{new1}$ , mk:time, "5/5/2007")
( $\theta_{new2}$ , mk:certainty, "0.6")
( $\theta_{new2}$ , mk:time, "6/6/2001")
( $\theta_{new3}$ , mk:certainty, "0.6")
( $\theta_{new3}$ , mk:time, "6/6/2001")}
```

## 5.4. Advanced Query Forms

The SPARQL query forms presented so far can be seen as a straightforward adaptation of querying capabilities to a knowledge base with associated meta knowledge. Beyond this, the presented framework can potentially be extended for more complex query forms with meta knowledge awareness.

*Queries with conditions on meta knowledge*. A possible interesting extension of our framework is the support for additional selection and/or ranking conditions on associated meta knowledge attributes. From the conceptual point of view, this functionality can be seen as an extension of the WHERE clause. In our example introduced in Section 2, the query about affiliations of researchers (Example 5.1) may also require that only certain facts (say with associated certainty values greater 0.8) shall be used. This restriction would result in exclusion of triples from graph  $G_2$  (with insufficient certainty of 0.6) from evaluation. Consequently, the result set would only contain the statement *JamesHendlerworksAtRensselaerPI* which is completely constructed using triples from remaining graph  $G_1$ , with its high certainty of 0.9. This idea is reflected in [4], which gives examples of querying RDF through SPARQL queries with meta knowledge constraints,

using the notion of Named Graphs. Such functionality can be realized in our framework in a straightforward manner, using *nested* queries (or querying with meta knowledge from *views*). Since nested queries and result ranking are expected to become part of the SPARQL2 specification in near future, the presented extension can be seen as a promising upcoming task.

We note that integration of nested queries with our approach would result in a substantially more flexible and powerful solution than the approach from [4]. In particular, our framework potentially allows to express filtering and/or ranking conditions on meta knowledge of *query results*. Conceptually, this functionality can be seen as an extension of the HAVING clause (in the common SQL like sense). For instance, for our sample query from Example 5.1 we may require the minimum overall certainty of each result to be greater 0.8. The resulting filter condition is then applied to query results from Example 5.16, which contain complex interpretations of meta knowledge for possible variable bindings.

*Nested meta knowledge*. The *nested* meta knowledge (i.e. meta knowledge about meta knowledge) can potentially be expressed through RDF capabilities. Following our example from Section 2, in the nested setting we would be able to specify conditions regarding recency or reliability of associated meta knowledge. For instance, we may require that only query results with high certainty of their extraction sources shall be returned. This kind of queries was not yet a concern in our framework presented so far. However, the required adaptation is not too difficult. As discussed in Section 3, there is no conceptual separation between knowledge and meta knowledge in our repository. For this reason, complex queries may treat meta knowledge as knowledge and obtain its meta annotation (i.e. nested meta knowledge) in a quite straightforward manner. Likewise, advanced querying mechanisms with nested meta knowledge support can easily be supported. As a consequence from the preceding discussion, we note that advanced conditions with nested meta knowledge can also be specified both for repository contents (i.e. kind of advanced WHERE clause) and for results of complex queries (i.e. kind of advanced HAVING clause).

## 6. Tasks and Benefits

This section summarizes the discussed steps of meta knowledge representation and utilization for the sample scenario that was introduced in Section 2.

### 6.1. Tasks for the administrator

In order to represent and utilize meta knowledge, the system administrator has to make some design choices. In particular, the application-specific meta knowledge properties must be defined. In our sample scenario, we consider three meta knowledge properties: source, certainty, and timestamp. In the next step, the administrator defines the intended semantics of these properties in order to facilitate query processing with complex expressions and pattern combinations. For evaluating graph pattern expressions, we use the definition from Section 5.2, and we assume that the corresponding definitions for meta knowledge properties are defined, e.g. the fuzzy logic operations presented Example in 5.10.

Finally, data and available associated meta knowledge are represented in RDF using named graphs [5,6], and imported into our RDF<sup>+</sup>-based repository.

### 6.2. Processing performed by the System

We assume that the system imports the small sample dataset introduced in Section 2. The knowledge base is transformed into the RDF<sup>+</sup> quintuples shown in Example 5.2 as discussed in Section 4. Associated meta knowledge is transformed into further RDF<sup>+</sup> literal statements and RDF<sup>+</sup> meta knowledge statements. For the properties *mk:time* and *mk:certainty* an example is presented in Example 5.14.

Following our sample scenario, the query from Example 5.1 can be reformulated as the query from Example 5.15 which retrieves names of Semantic Web experts together with their affiliations and the associated meta knowledge. Internally, the query processor evaluates this query using graph patterns as discussed in Section 5.2. If  $P$  denotes the graph pattern from this query then all matches for all variables in  $P$  are given by the evaluation  $\llbracket P \rrbracket$ . The resulting set of annotated variable assignments is shown in Example 5.8. It contains possible variable assignments, and the how-provenance ( $\mathfrak{E}_3$ ) that explains how these source statements have been used.

By combining this information with definitions for meta knowledge properties and available meta knowledge statements, the query processor constructs the result shown in Example 5.16. This result is then serialized in RDF.

### 6.3. Benefits for the user/developer

The user or application developer can access the knowledge stored in the RDF<sup>+</sup>-based repository in different ways. On one hand, the repository does not change the existing SPARQL semantics and thus fully supports common SPARQL queries. This is an important advantage for compatibility with existing applications and interfaces. On the other hand the repository supports the SPARQL with metaknowledge (Section 5.2). Thus, the user obtains additional access to valuable meta knowledge that can be used for relevance ranking, conflict resolution, or other applications in connection with retrieved knowledge.

In our application scenario, the user may realize that the query answer is potentially contradictive (James Hendler is affiliated with Rensselaer PI and University of Maryland). By inspecting the associated meta knowledge, he would realize that the second fact was generated by mistake. In fact, it is based on outdated information (knowledge from the document *survey.pdf* with timestamp 6/6/2001) that was wrongly combined with knowledge from a more recent source (namely document *report.doc* with timestamp 5/5/2007). It turns out that the affiliation of James Hendler has actually changed from University Maryland to Rensselaer PI, and the erroneous tuple can be safely excluded from further processing.

## 7. Equivalences

The main goal of this section is to show the equivalence holding between the standard evaluation of SPARQL queries, denoted by  $\llbracket \cdot \rrbracket$  (see Appendix A), and the evaluation with meta knowledge, presented in this paper denoted by  $\lllbracket \cdot \rrlbracket$ . For this purpose, we first define the meta knowledge projection operator  $\Psi$  as following:

### Definition 7.1 (Meta Knowledge Projection)

Let  $\Phi$  be an annotated relation representing the set of tuples of variable assignments and their annotations, then

$$\Psi(\Phi) = \Pi_{\text{var}(P)}(\sigma_{(I_s \rightarrow \top)}(\Phi))$$

where  $\text{var}(P)$  is the set of variables occurring in  $P$ , and  $I_s \rightarrow \top$  denotes the standard interpretation  $I_s$  of an annotation formula evaluated to  $\top$ .

Basically, the meta knowledge projection operator se-

lects all the statements corresponding the variables bindings of the evaluation and ignores the meta knowledge properties assigned to the statements in the knowledge base. As discussed earlier the standard interpretation  $I_s$  of an annotation formula evaluates to  $\top$  if and only if the statements in the knowledge base are elements of the variable assignments obtained via evaluation, i.e. the corresponding variable assignment is in the set of variable assignments obtained via standard SPARQL evaluation extended with identifiers (see Definition 5.3). Based on this, the next proposition defines the equivalence holding between the evaluation with meta knowledge and the standard SPARQL evaluation.

**Proposition 7.1** *Given a graph pattern expression  $P$ , we say that the meta projection on the results of evaluation with meta knowledge of the graph pattern expression  $P$  over a RDF<sup>+</sup> dataset  $D^+$  is equivalent to the standard evaluation of  $P$  over the RDF dataset  $D$ , denoted by  $\Psi[\llbracket P \rrbracket^{D^+}] = \llbracket P \rrbracket^D$ .*

**Proof:** (sketch) For sets of variable assignments obtained via basic pattern, this follows immediately from Definition 5.3 and the application of operator  $\Psi$  over  $\Phi$ , so that for each substitution  $\Phi(\mu)$  obtained by evaluation over a dataset  $D^+$ ,  $\llbracket P \rrbracket^{D^+}$ , can be reduced to a substitution  $\mu$  obtained from the evaluation  $P$  over a RDF dataset  $D$ ,  $\llbracket P \rrbracket^D$ , defined in [15,17] by selecting all  $\mu$  where  $I_s$  of an annotation formula evaluates to  $\top$  and dropping all remaining substitutions of  $\theta$ . For complex graph patterns, the evaluation defined (see Definition 5.3) is indeed identical to [15,17].■

For relations annotated with Boolean formulae and standard relational algebra a proof for Proposition 7.1 can be found in [9].

Now, we can define the equivalence holding between equivalence graph pattern expressions with respect to the evaluation with meta knowledge, i.e. show that the evaluation of equivalent queries are annotated with equivalent meta knowledge properties. By equivalence we mean that the interpretation of formulae used to annotate bindings contained in the results of the equivalent queries are equivalent.

We use for the algebra of annotation formulae as syntactic objects  $(\mathfrak{F}, \wedge, \vee, \neg, \perp, \top)$  (see Section 5). The elements of  $\mathfrak{F}$  are Boolean formulae built over the set of all literal RDF<sup>+</sup> statements  $\mathfrak{S}$ . Instead of the statements themselves we use their identifiers. In the tagged tuple model which we adopt here a relation is represented by a function mapping all tuples of a domain to an annotation. Since we use  $\mathfrak{F}$  only as a placeholder

we could just as well directly annotate variable bindings with elements of different meta knowledge algebras. Therefore, first we need to define restrictions on the algebra for meta knowledge interpretation for each meta knowledge property  $\gamma$ ,  $(\Omega_\gamma, \wedge_\gamma, \vee_\gamma, \neg_\gamma, \top_\gamma, \perp_\gamma)$  as defined in Section 5, so that the equivalence holds. As a motivation for the reader, before presenting the restrictions for the meta knowledge algebra, in the next example we show one case where the equivalence holds and another where the equivalence does not hold when the meta knowledge properties are defined using the fuzzy logic operations as presented in Example 5.10. Let  $P = (P_1 \text{ AND } P_2)$  and  $P' = (P_2 \text{ AND } P_1)$  be equivalent complex graph pattern expressions. The annotation results for  $P$  and  $P'$  consist of  $(F_1 \wedge F_2)$ , and  $(F_2 \wedge F_1)$ , respectively, where  $F_1$  and  $F_2$  denote annotation Boolean formulae. We are going to evaluate  $P$  and  $P'$  with the meta knowledge property *time*. The  $\Gamma$ -interpretation for  $I_{time}$  for the  $(F_1 \wedge F_2)$  formulae is  $I_{time}^f(F_1 \wedge F_2) = \max(I_{time}^f(F_1) \wedge_{time} I_{time}^f(F_2)) = \max(I_{time}^f(F_1), I_{time}^f(F_2))$ . Likewise, the  $\Gamma$ -interpretation for  $I_{time}$  for the  $(F_2 \wedge F_1)$  formulae is  $I_{time}^f(F_2 \wedge F_1) = \max(I_{time}^f(F_2), I_{time}^f(F_1))$ . As shown in this example the meta knowledge evaluation for  $P$  and  $P'$  leads to the same result, thus it is equivalent. Now consider the equivalent graph patterns  $P$  and  $\neg(\neg(P))$ , and the respective annotation results  $F$  and  $\neg(\neg(F))$ , where  $F$  denote an annotation Boolean formula. The  $\Gamma$ -interpretation for  $I_{time}$  for the  $(\neg(\neg(F)))$  formulae is bounded to 0 which is not the case for the formulae  $F$ . Consequently the annotation results  $F$  and  $\neg(\neg(F))$  for the equivalent graph pattern expressions  $P$  and  $\neg(\neg(P))$  are not equivalent.

In the following we summarize some of the restrictions on the meta knowledge algebra. Let  $P_1$ ,  $P_2$  and  $P_3$  be graph pattern expressions and the built-in condition  $P$ :

- (i)  $\llbracket P_1 \text{ AND } P_2 \rrbracket_G^{D^+} = \llbracket P_2 \text{ AND } P_1 \rrbracket_G^{D^+}$ , and  $\llbracket P_1 \text{ UNION } P_2 \rrbracket_G^{D^+} = \llbracket P_2 \text{ UNION } P_1 \rrbracket_G^{D^+}$

Since annotations for results of graph pattern of the form  $P_1 \text{ AND } P_2$  are of the form  $F_1 \wedge F_2$ , where  $F_1$  and  $F_2$  denote annotation formulae, and annotations for results of graph pattern of the form  $P_1 \text{ UNION } P_2$  are of the form  $F_1 \vee F_2$  we require that:

$\wedge$  and  $\vee$  are associative and commutative.

for annotation formulae in order to guarantee that results of equivalent queries are associated with equivalent meta knowledge.

$$(ii) \quad \llbracket (P1 \text{ AND } (P2 \text{ UNION } P3)) \rrbracket_G^{D^+} = \llbracket ((P1 \text{ AND } P2) \text{ UNION } (P1 \text{ AND } P3)) \rrbracket_G^{D^+}$$

In order to satisfy this equivalence, the annotation formulae and algebras used to interpret these formulae need to satisfy the equivalence:

$$I_\gamma^f(F_1 \wedge (F_2 \vee F_3)) \equiv I_\gamma^f((F_1 \wedge F_2) \vee (F_1 \wedge F_3))$$

$$(iii) \quad \llbracket ((P1 \text{ UNION } P2) \text{ FILTER } R) \rrbracket_G^{D^+} = \llbracket ((P1 \text{ FILTER } R) \text{ UNION } (P2 \text{ FILTER } R)) \rrbracket_G^{D^+},$$

and

$$\llbracket (((P1 \text{ FILTER } R1) \text{ FILTER } R2)) \rrbracket_G^{D^+} = \llbracket (((P1 \text{ FILTER } (R1 \wedge R2))) \rrbracket_G^{D^+}.$$

Now we look at equivalences in the context of FILTER expressions. Equivalence requires:

$$I_\gamma^f((F_1 \vee F_2) \wedge r) \equiv I_\gamma^f((F_1 \wedge r)) \vee I_\gamma^f((F_2 \wedge r))$$

which follows from (i) and (ii).

Analogous:

$$I_\gamma^f((F_1 \wedge r_1) \wedge r_2) \equiv I_\gamma^f(F_1 \wedge r_3)$$

where  $r_3 = r_1 \wedge r_2$ .

$$(iv) \quad \llbracket (P1 \text{ AND } P1) \rrbracket_G^{D^+} = \llbracket P1 \rrbracket_G^{D^+}$$

Equivalence requires idempotency:

$$I_\gamma^f(F \wedge F) \equiv I_\gamma^f(F)$$

Based on the restrictions presented above, we see evidences for requiring the meta knowledge algebra to form such a commutative semiring [10] since the laws of commutative semirings are forced by certain expected identities in the meta knowledge algebra. For this purpose, we define in the next theorem that the meta knowledge algebra has to take form of a commutative ring algebraic structure in order to hold the equivalences presented above for meta knowledge interpretations. With the commutative ring approach, the extendability of our approach increases due to the low restriction level on the meta knowledge algebra required for the equivalences.

**Theorem 7.1** *Let  $P1$  and  $P2$  be equivalent graph pattern expressions using AND, OR, and FILTER operator and  $D^+$  a RDF<sup>+</sup> dataset. We have that  $\llbracket P1 \rrbracket_G^{D^+} = \llbracket P2 \rrbracket_G^{D^+}$  if and only if all meta knowledge properties form a commutative semiring.*

**Proof:**(sketch) This proof is an immediate consequence of the laws of commutative semirings defined in [10]. Suppose that the meta knowledge algebra  $(\Omega_\gamma, \wedge_\gamma, \vee_\gamma, \top_\gamma, \perp_\gamma)$  is not a commutative semiring with the following identities (see Definition 5.5, and Definition 5.6): *union* is associative, commutative

and has identity; *and* is associative, commutative and distributive over *union*; *projections* and *selections* commute with each other as well as with *unions* and *joins*. However, according to [10] these identities hold for meta knowledge algebra if and only if  $(\Omega_\gamma, \wedge_\gamma, \vee_\gamma, \top_\gamma, \perp_\gamma)$  is a commutative semiring, which contradicts our assumption. Hence the meta knowledge algebra is a commutative semiring, i.e., algebraic structure  $(\Omega_\gamma, \wedge_\gamma, \vee_\gamma, \top_\gamma, \perp_\gamma)$  such that  $(\Omega_\gamma, \wedge_\gamma, \top_\gamma)$  and  $(\Omega_\gamma, \vee_\gamma, \perp_\gamma)$  are commutative monoids,  $\wedge$  is distributive over  $\vee$  and  $\forall a, 0, a \wedge 0 = a \wedge 0 = 0$  ■

To illustrate an instance of this theorem, consider the meta knowledge formula  $F_1 \wedge F_2$  and  $F_2 \wedge F_1$  for the graph pattern expressions  $P = (P_1 \text{ AND } P_2)$  and  $P' = (P_2 \text{ AND } P_1)$ . Evaluating them in any of the meta knowledge dimensions, we get indeed the same value from the interpretation.

## 8. Complexity of Eval<sup>+</sup>

In this section we analyze how the construction of the annotations influences the complexity of the decision problem related to SPARQL. The decision problem associated with the standard evaluation of a SPARQL query can be stated as following [15]: *Given an RDF dataset  $D$ , a graph pattern  $P$  and a mapping  $\mu$ , determine whether  $\mu$  is in the result of  $P$  applied to  $D$ .* For this decision problem, which we denote by  $Eval_D(P)$ , an analysis of the complexity is presented in [15,16]. In the context of RDF<sup>+</sup> datasets and annotated variable assignments we have a slightly different decision problem: *Given an RDF<sup>+</sup> dataset  $D^+$ , an RDF<sup>+</sup> graph pattern  $P$ , a variable assignment  $\mu$  and an annotation  $\alpha$  determine whether  $\alpha$  is the correct annotation of  $\mu$ .* We denote this problem by  $Eval_{D^+}^+(P)$ . An annotation is correct iff the formula is *equivalent* (in the logical sense) to the formula obtained by evaluation as defined in Section 5.

With the following theorems we show the time complexity and space complexity of  $Eval^+$  for patterns that do not use the OPTIONAL operator. Note that the RDF<sup>+</sup> dataset  $D^+$  is built over a set of RDF literal statements and meta knowledge statements of the RDF dataset  $D$  (see Section 4.4). For  $D^+ = (K, M)$ , we assume that the size of  $|D^+| = |K| \cdot |M| \leq 1/2 |D|^2$ , and consequently the time complexity of building  $D^+$  is the order of  $O(|D|^2)$ .

The RDF counterparts of both theorems have been established by [15,16]. Like [15,16] we restrict to graph patterns which do not contain blank nodes. In

the first theorem we consider graph patterns which use only AND and FILTER operations. Bindings obtained by such patterns are annotated with formulae which do not contain any other operator besides  $\wedge$  according to Definition 5.4.

**Theorem 8.1** *Eval<sup>+</sup> can be solved in time  $O(|P| \cdot |D^+|)$  for graph pattern expressions constructed by using only AND and FILTER operators and for annotation formulas using only the operation  $\wedge$ .*

**Proof:** The proof consists of two parts: In the first part we construct a correct annotation  $\hat{\alpha}$  for  $\mu$  and in the second we check whether  $\alpha$  and  $\hat{\alpha}$  are equivalent. In the following let  $p_i$  denote pattern  $i$  from  $P$  and  $\mu(p_i)$  denote the triple obtained by replacing the variables in the pattern  $p_i$  according to  $\mu$ .

In order to construct  $\hat{\alpha}$  we start by evaluating  $\llbracket p_i \rrbracket^{D^+}$  for all pattern using Definition 5.3. In order to achieve this  $D^+$  needs to be searched for all  $\mu(p_i)$ . This can be performed in  $O(|P| \cdot |D^+|)$ . Then, we construct the algebraic expression  $\Phi$  by evaluating  $\llbracket P \rrbracket^{D^+}$  using definition 5.4. This can be performed by traversing  $P$ . The correct annotation  $\hat{\alpha}$  of  $\mu$  in the result of evaluating on  $D^+$  is defined as  $\Phi(\mu)$  (see Definition 5.5). We can construct  $\hat{\alpha} = \Phi(\mu)$  in a depth-first traversal of  $\Phi$  as following:

Let  $\Phi_1, \Phi_2$  be algebraic expressions part of  $\Phi$ . For each join operation in  $\Phi$  the annotation of  $(\Phi_1 \bowtie \Phi_2)(\mu)$  is given by  $\Phi_1(\mu) \wedge \Phi_2(\mu)$ . For each select operation  $\sigma_c(\Phi)(\mu)$  is given by  $\Phi(\mu)$  if condition  $c$  is fulfilled otherwise it is given by  $\perp$ . Since traversing and  $\Phi$  each has time complexity  $O(|P|)$  the evaluation of  $\hat{\alpha} = \Phi(\mu)$  remains in  $O(|P| \cdot |D^+|)$ .

Now we determine whether for a given annotation  $\alpha$  holds  $\alpha \equiv \hat{\alpha}$ . We transform both formulas into a normal form in  $O(|P| \cdot \log |P|)$  using associativity, commutativity and idempotency of  $\wedge$ . First, we remove all brackets then establish an order among atomic formulas (identifiers and  $\top, \perp$ ) and finally remove duplicates. For  $\alpha$  and  $\hat{\alpha}$  normalized this way syntactic equality bi-implies logical equivalence. Assuming  $|P| < |D^+|$  the overall complexity remains in  $O(|P| \cdot |D^+|)$ . ■

**Theorem 8.2** *Eval<sup>+</sup> is NP-complete for graph pattern expressions constructed by using only AND, FILTER and UNION operators.*

**Proof:** The proof consists of two parts: In the first part we show that *Eval<sup>+</sup>* is contained in NP and in the second we establish NP-hardness of *Eval<sup>+</sup>*. As above, let  $p_i$  denote triple pattern  $i$  from  $P$  and  $\mu(p_i)$  denote the triple obtained by replacing the variables in the pattern  $p_i$  according to  $\mu$ .

The first part consists of two steps as well: first we construct a correct annotation  $\hat{\alpha}$  for  $\mu$  and then we check whether  $\alpha$  and  $\hat{\alpha}$  are equivalent. In order to construct  $\hat{\alpha}$  we start by evaluating  $\llbracket p_i \rrbracket^{D^+}$  for all pattern using definition 5.3. In order to achieve this  $D^+$  needs to be searched for all  $\mu(p_i)$ . Then, we construct the algebraic expression  $\Phi$  by evaluating  $\llbracket P \rrbracket^{D^+}$  using Definition 5.4. This can be performed by traversing  $P$ . The correct annotation  $\hat{\alpha}$  of  $\mu$  in the result of evaluating  $P$  on  $D^+$  is defined as  $\Phi(\mu)$ , see definition 5.5. We can construct  $\hat{\alpha} = \Phi(\mu)$  in a depth-first traversal of  $\Phi$  as following:

Let  $\Phi_1, \Phi_2$  be algebraic expressions part of  $\Phi$ . For each join operation  $(\Phi_1 \bowtie \Phi_2)(\mu)$  is given by  $\Phi_1(\mu) \wedge \Phi_2(\mu)$ . For each union operation  $(\Phi_1 \cup \Phi_2)(\mu)$  is given by  $\Phi_1(\mu) \vee \Phi_2(\mu)$ . For each select operation  $\sigma_c(\Phi)(\mu)$  is given by  $\Phi(\mu)$  if condition  $c$  is fulfilled otherwise it is given by  $\perp$ .

Time complexity for evaluating  $\Phi(\mu)$  is  $O(|P| \cdot |D^+|)$ . The evaluation of  $\alpha \equiv \hat{\alpha}$  is more difficult if UNION operations are contained in  $P$ . But it is subsumed by checking equivalence of Boolean formulae which is a NP-complete problem. Thus, the decision problem *Eval<sup>+</sup>* is contained in NP.

We can deduce that *Eval<sup>+</sup>* is an NP-complete problem if *Eval*, which has been shown to be NP-complete [15,16], can be reduced to it. *Eval* can be reduced to *Eval<sup>+</sup>* if the evaluation defined in Section 5 results in an annotation  $\alpha \equiv \perp$  exactly for such bindings which are not in the result of standard SPARQL evaluation. For graph patterns that do not contain the OPTIONAL operator this can be shown by induction over the structure of algebraic expressions using Definition 5.4 and its standard SPARQL counterpart from [15]. ■

## 9. Implementation

The framework *metaK* described in this paper has been implemented in Java language, using libraries from the Sesame Project<sup>5</sup>. The prototype of *metaK* is available as an open source implementation at <http://isweb.uni-koblenz.de/Research/MetaKnowledge>.

From the practical point of view, the framework can easily be customized for new meta knowledge aspects. Dimension-specific interpretations of meta knowledge properties are implemented separately, as small Java classes. For defining a new meta knowledge aspect, the corresponding interpretation function must be implemented with respect to framework interfaces.

<sup>5</sup> Sesame Project: <http://www.openrdf.org>

In order to evaluate the overhead produced by the evaluation of meta knowledge properties for results of SPARQL queries we carried out two experiments based on the well-known LUBM benchmark [11]. Our main aim is to find out whether the evaluation of SPARQL queries remains feasible if meta knowledge is provided within the query results, i.e. conduct a family of experiments to validate the theoretical results of Section 8.

A key question is how to separate the additional effort for the evaluation of provenance and meta knowledge from standard SPARQL processing. Triples describing meta knowledge receive an *additional* meta knowledge interpretation according to section 4. At the same time they are also treated as ordinary RDF triples (and thus can be queried using standard SPARQL). Thus, if we add meta knowledge to a knowledge base this increases its overall size which also increases the workload for standard query processing. In order to account for this we compare query evaluations performed on the same knowledge base which includes meta knowledge. Our implementation is built on top of Sesame<sup>6</sup> 2.0 (beta 6) using query rewriting. The triple store is used to store both, knowledge and meta knowledge. We expect that a native implementation of the meta knowledge framework can achieve an increased performance.

*Query evaluation.* For the evaluation of SPARQL on RDF<sup>+</sup> we defined the results of SELECT queries to be set-valued (see Section 5.2). For standard SPARQL [20], however, a SELECT query returns a *solution sequence* which may contain duplicate elements. All 14 queries of the LUBM benchmark are SELECT queries. To allow for a consistent comparison of extended evaluation on one side and standard evaluation on the other we added the keyword DISTINCT to the queries for standard evaluation. This tells the query processor to eliminate duplicates. Since the evaluation of individual meta knowledge properties can be arbitrarily complex we compare the following three kinds of query evaluation:

- Standard evaluation with additional duplicate elimination (SD), as performed by Sesame,
- Evaluation of provenance formulae for each query result (PF) and
- Evaluation of four basic meta knowledge properties (M4), namely *agent*, *confidence*, *creation time* and *source*, see Example 5.10 (we evaluate *agent* analogously to *source*).

<sup>6</sup> Sesame: open source framework for storage, inferencing and querying of RDF data ([www.openrdf.org](http://www.openrdf.org))

Only the last type of query processing actually makes use of the additional triples.

*Data.* We added artificial meta knowledge to the LUBM data. Amount and granularity of the additional meta knowledge are key properties of the resulting dataset.

We created two datasets containing a different percentage of meta knowledge triples. The first dataset (MK29) was created based on LUBM OWL data for ten Universities. We added the meta knowledge by putting groups of ten consecutive original triples into one named graph and associating random values of the meta knowledge properties *agent*, *confidence*, *creation time* and *source* with this graph name, see Example 9.1. As a consequence the dataset contains 29 percent of meta knowledge which might be a reasonable scenario in a real world scenario. The resulting dataset consists of 1.8 million triples of which 1.3 million triples were created by the LUBM generator to which we added 0.5 million meta knowledge triples. It contains 0.3 million (additional) graph URIs.

**Example 9.1** *Original triples of dataset MK29 in one named graph associated with meta knowledge properties*

---

```

<graph>
  <uri>http://www.x-media-project.org/
    ontologies/someGraph#0-0_30 </uri>
  <triple>
    <uri>http://www.Department0.University0.edu/FullProfessor1 </uri>
    <uri>http://www.lehigh.edu/%7Ezhp2/2004/0401/univ-bench.owl#doctoralDegreeFrom </uri>
    <uri>http://www.University882.edu </uri>
  </triple>
  <triple>
    <uri>http://www.University882.edu </uri>
    <uri>http://www.w3.org/1999/02/22-rdf-syntax-ns#type </uri>
    <uri>http://www.lehigh.edu/%7Ezhp2/2004/0401/univ-bench.owl#University </uri>
  </triple>
  ...
<graph>
  <uri>http://www.metaknowledge.semanticweb.org0-0_30 </uri>
  <triple>
    <uri>http://www.x-media-project.org/ontologies/someGraph#0-0_30 </uri>
    <uri>http://www.x-media.org/ontologies/metaknow#source </uri>
    <uri>http://www.x-media-project.com/ex#source30 </uri>
  </triple>
</triple>

```

```

<uri>http://www.x-media-project.org/
  ontologies/someGraph#0-0_30 </uri>
<uri>http://www.x-media.org/ontologies/
  metaknow#confidence_degree </uri>
<typedLiteral datatype='http://www.w3.org/
  2001/XMLSchema#double'>0.7</typedLiteral>
</triple>
...

```

The data set MK400 was created based on LUBM OWL data for three Universities and we assigned each triple to a different graph and the four meta knowledge properties are associated with it. This way the knowledge base contains four times as many meta knowledge triples as knowledge triples. The resulting dataset consists of 1.7 million triples of which 0.35 million are original LUBM data and 1.4 million are additional meta-data. Note that the overall sizes of the two datasets are comparable. As indicated above the overall size of the knowledge base influences the workload for query processing. In fact, main memory consumption appeared to be a key factor. By choosing similar sizes for the two datasets we reduce the influence of factors related to the size of the knowledge base and concentrate on how the different evaluations influence the processing time.

	SD	PF	M4
MK29	313	894	1382
MK400	196	104	425

Table 1  
Random query sequence experiments: average processing time (ms)

*Random Query Sequence.* We conducted two experiments with each of the two datasets. One experiment aims at simulating behavior of a query engine in a real life scenario: first a dataset is loaded and then a sequence of queries is submitted to the query engine. We measured the average processing time of these queries. The query sequence was created based on 8 of the 14 queries from the LUBM benchmark. Query 2 was not included since we were not able to obtain results for this query with this version of Sesame, this dataset and an average machine. Five more queries are discarded since they require OWL inferencing or hierarchy information to obtain complete results and Sesame 2 was not able to obtain bindings using plain SPARQL processing. Since no meta knowledge needs to be calculated if the result set is empty using these queries would bias the evaluation in favor of the meta knowledge processing. The authors of the benchmark identified three main characteristics of queries with respect to plain SPARQL processing: *input size*, *selectivity* and *complexity*. The remaining 8 queries still cover different settings for these

features. Experiments with single queries will be presented below.

The query sequence consisted of a random shuffle of 20 copies of each of the remaining queries. For each query in the sequence we measured the time which elapses during issuing the query, obtaining the result and traversing the result sequentially. This measure is similar to the *query response time* defined in [11]. The only difference is that in [11] each query is performed ten times after the knowledge base has been loaded to measure the caching performance of the query engine. For each run we determined the average of the execution times of the queries in the sequence in question. The results are summarized in Table 1.

Evaluation of provenance formulae (PF) given dataset MK29 almost tripled the average query execution time. On average, the evaluation took about half a second longer than standard evaluation (SD). The evaluation of four meta knowledge properties (M4) adds again half a second to the evaluation of provenance. The overall overhead to obtain meta knowledge is about a second given a non-trivial dataset. We consider this to be an indication for the feasibility of our approach. Since these numbers are average values the question remains whether reasonable processing times are achieved for all individual queries as well. This will be analyzed below.

For the dataset MK400 the computations were faster for all three kinds of evaluations. This can be explained by the fact that this dataset contains a smaller number of knowledge triples and therefore the result sets for some of the queries contain fewer bindings as we will see below. Surprisingly, evaluation of provenance formulas needed less time than standard evaluation. A possible explanation is optimization performed by the query processor of Sesame 2. Since our implementation uses query rewriting different queries are evaluated for the different kinds of evaluations. Optimization techniques might be easier to apply to some of them. Why does the evaluation of MK29 not show similar characteristics? Here a possible explanation is the larger number of bindings involved in the evaluation. Main memory consumption was quite large in both cases. Possibly there was no space left for caching of (intermediary) results given dataset MK29. The calculation of the four meta knowledge properties did result in an increase in processing time as expected.

*Single LUBM Queries.* We also measured processing times for single queries. As stated above we measured the time which elapsed during issuing the query, obtaining the result and traversing the result sequentially. In

Query		Q1	Q4	Q5	Q6	Q7	Q8	Q10	Q14	Av.
# bindings		5	10	411	24019	3	1874	4	75547	
processing time (ms)	SD	127	163	211	357	135	1619	127	977	465
	PF	324	347	386	995	337	691	325	2320	716
	M4	340	375	426	2187	346	878	326	10501	1922

Table 2

Processing time of query evaluation with and without provenance and meta knowledge for individual queries and the MK29 dataset. Processing times are average values of 10 runs each.

Query		Q1	Q4	Q5	Q6	Q7	Q8	Q10	Q14	Av.
# bindings		5	10	411	6390	3	1874	4	19868	
processing time (ms)	SD	137	189	204	217	131	1596	124	379	372
	PF	367	350	367	517	316	676	300	769	458
	M4	346	359	440	859	321	1083	306	1953	708

Table 3

Processing time of query evaluation with and without provenance and meta knowledge for individual queries and the MK400 dataset. Processing times are average values of 10 runs each.

contrast to the previous experiment and to the definition of the *query response time* from [11] the application was restarted before each single query execution. That way each query was evaluated against a newly loaded knowledge base since we want to measure the effort it takes to evaluate the queries and not the caching strategy of the query engine. This procedure was repeated ten times for each query and each method of query evaluation. The average values from these runs are summarized in tables 2 and 3. The standard deviations estimated from these ten runs are less than or equal 10 percent for all queries and methods evaluated on the two datasets.

On average the calculation of provenance formulas (PF) increases processing time by factor 1.5. In absolute numbers the average increase is about 0.2 seconds. The largest increase (1.3 seconds for the MK29 dataset) can be observed for query 14 which also gives the largest result set. We attribute this to the main memory consumption of our implementation. For query 8 there even is a decrease in processing time. As for the case of a random query sequence evaluated on MK400 a possible explanation are optimizations of the query processor. The additional querying for meta knowledge might guide the optimization of the query execution.

Calculation of the four basic provenance properties (MD4) causes an average increase of factor 4.1 (1.5 seconds) for the MK29 dataset and factor 1.9 (0.3 seconds) for the MK400 dataset compared to standard evaluation. We ascribe the larger increase for the runs based on dataset MK29 to the larger number of results for queries 6 and 14 and the non-optimized memory con-

sumption of our implementation. These results are in line with the results from the experiments using a random query sequence shown in table 1. For query 8 the processing time decreases for query evaluation including meta knowledge as well which might be explained by query optimization of the underlying triple store as stated above.

*Discussion.* At first we consider the experiments with individual queries. From the estimated standard deviations of the 10 runs for each query, kind of evaluation and dataset we conclude that the measurements are reliable enough to draw two general conclusions: On the one hand we can observe a noticeable increase of processing time using our implementation and on the other hand the amount of this increase can be described by a small linear factor.

The results of the experiments using a random sequence of queries indicate that similar results also hold for real world scenarios. Here caching can be applied by a query processor. A few times evaluations of the same query which was repeated directly one after another in the sequence. The resulting average processing times are smaller but still comparable to the average values for the evaluations of single queries. A key insight is that the overall processing times remain feasible for evaluations on a dataset of up to 1.8 million triples and results of up to 75,000 bindings.

If we compare the processing times obtained for the two different datasets we might expect that for dataset MK400 processing times increase by a larger factor if meta knowledge is involved. MK400 contains a larger number of meta knowledge triples which need to be processed in order to evaluate meta knowledge – especially compared to the number of knowledge triples. However, at least with our implementation, the dominant characteristic of query evaluation appears to be the size of the result set.

## 10. Related Work

The importance of better understanding the ways by which the result came about is fundamental to many Semantic Web applications and scenarios. The specification of the Semantic Web proof layer was discussed in [13,18,12].

Meta knowledge has been initially studied as an extension for relational databases (i.e. data management systems based on relational algebra), probabilistic databases [9], and later adopted for RDF knowledge bases (i.e. for semantics of SPARQL query language).

In the area of database systems, meta knowledge is often represented using an extension of the relational data model, coined *annotated relations*. Its purpose is primarily the description of data origins (provenance) and the process by which it arrived as a query answer [7,2,3,8]. Basically they define custom (possibly different) interpretations for algebraic operations of boolean formulas (built on tuple identifiers as Boolean variables) for particular dimensions of meta knowledge (e.g. agent, timestamp, source, certainty) to obtain the  $m$ -dimensional record (i.e. query result). Indeed, a Boolean expression of the result set build from tuple identifiers does not only carry the information *which* triples have contributed to a variable assignment (why-provenance) but also *how* they contributed (how-provenance). Basically, our methodology follows the same idea and adopts the notion of provenance semirings which was introduced, as a generalization framework, in [10]. In contrast to other approaches discussed in [10] (based on the notion of annotated tuples in a relational schema), our solution is customized for RDF graphs (i.e. annotated RDF quadruples). The same holds for the corresponding query language (basically, SPARQL vs. SQL) and its semantics. An important conceptual difference to the relational model is the natural ability of RDF/SPARQL repositories for result serialization and thus seamless exchanging and utilization of knowledge and meta knowledge from our framework across multiple Semantic Web nodes without additional schema integration efforts.

In the Semantic Web field, meta knowledge has been recently considered in applications for assessing the trustworthiness of information [4,8]. In particular, meta knowledge shows how statements are organized among agents on the Web. Our approach is focused on an RDF language model and provides fine-grained meta knowledge management for retrieval queries with SPARQL that is not directly comparable with proof traces for OWL reasoning.

The use of Named Graphs for meta knowledge management and querying RDF with meta knowledge constraints is introduced in [4]. As discussed in Section 5.4, this functionality can be realized in our framework in a straightforward manner by the means of *nested* queries (or querying from *views*), which are expected features of the upcoming SPARQL2 query language. Moreover, our framework potentially allows to express filtering and/or ranking conditions on meta knowledge of *query results* and thus is more expressive and flexible than the solution presented in [4].

## 11. Conclusion and Future Work

In this paper, we presented an original, generic, formalized and implemented approach for the management of many dimensions of meta knowledge, like source, authorship, certainty, and others, for RDF repositories. Our method re-uses existing RDF modeling possibilities in order to represent meta knowledge. Then, it extends SPARQL query processing in such a way that given a SPARQL query for data, one may request meta knowledge without modifying the query proper. We achieve highly flexible and automatically coordinated querying for data and meta knowledge, while completely separating the two areas of concern. Our approach remains compatible to existing standards and query languages and can be easily integrated with existing applications and interfaces.

In the future, we will investigate the meta knowledge support for OWL-based knowledge bases with advanced reasoning capabilities. Due to the substantially higher complexity of inferencing and retrieval algorithms (e.g. reasoning in OWL-DL vs. RDF querying with SPARQL) and the distributed nature of knowledge sources in the Semantic Web, the notion of meta knowledge will require further, non-trivial justification. Our long-term objective is the generic, efficient and effective infrastructure for meta knowledge management as an integral part of the proof layer of the Semantic Web.

*Acknowledgements* This work was supported by the X-Media project ([www.x-media-project.org](http://www.x-media-project.org)) funded by the European Commission under EC grant number IST-FP6-026978 and by the project WeKnowIt ([www.weknowit.eu](http://www.weknowit.eu)) funded by the European Commission under EC grant number FP7-215453.

## References

- [1] C. Bizer, R. Cyganiak, The TriG Syntax, 2007.  
URL <http://sites.wiwi.fu-berlin.de/suhl/bizer/TriG/Spec/TriG-20070730/>
- [2] P. Buneman, S. Khanna, W. C. Tan, Data Provenance: Some Basic Issues, 20th Conference on Foundations of Software Technology and Theoretical Computer Science (FSTTCS), New Delhi, India (2000) 87–93.
- [3] P. Buneman, S. Khanna, W. C. Tan, Why and Where: A Characterization of Data Provenance, Proc. of ICDT (2001) 316–330.
- [4] J. J. Carroll, C. Bizer, P. Hayes, P. Stickler, Named graphs, provenance and trust, in: WWW '05: Proceedings of the 14th international conference on World Wide Web, ACM, New York, NY, USA, 2005.

- [5] J. J. Carroll, C. Bizer, P. J. Hayes, P. Stickler, Named graphs, *J. Web Sem.* 3 (4) (2005) 247–267.
- [6] J. J. Carroll, P. Stickler, TriX: RDF triples in XML, in: *Proceedings of the Extreme Markup Languages 2004*, Montreal, Canada, 2004.
- [7] Y. Cui, J. Widom, Practical Lineage Tracing in Data Warehouses, *Proc. of ICDE (2000)* 367–378.
- [8] L. Ding, P. Kolari, T. Finin, A. Joshi, Y. Peng, Y. Yesha, On Homeland Security and the Semantic Web: A Provenance and Trust Aware Inference Framework, in: *Proceedings of the AAAI Spring Symposium on AI Technologies for Homeland Security*, 2005.
- [9] N. Fuhr, T. Rölleke, A probabilistic relational algebra for the integration of information retrieval and database systems, *ACM Transactions on Information Systems* 15 (1) (1997) 32–66.
- [10] T. J. Green, G. Karvounarakis, V. Tannen, Provenance Semirings, in: *PODS*, 2007.
- [11] Y. Guo, Z. Pan, J. Heflin, LUBM: A Benchmark for OWL Knowledge Base Systems, *Journal of Web Semantics* 3 (2) (2005) 158–182.
- [12] D. McGuinness, P. Pinheiro da Silva, Explaining Answers from the Semantic Web: the Inference Web Approach, *J. Web Sem.* 1 (4) (2004) 397–413.
- [13] W. Murdock, D. McGuinness, P. Pinheiro da Silva, C. Welty, D. Ferrucci, Explaining Conclusions from Diverse Knowledge Sources, *International Semantic Web Conference (ISWC)*, Athens, USA (2006) 861–872.
- [14] P. Hayes (ed.), *Rdf semantics*, <http://www.w3.org/TR/rdf-mt/> (2004).
- [15] J. Perez, M. Arenas, C. Gutierrez, Semantics and Complexity of SPARQL, in: *Proc. of ISWC*, 2006.
- [16] J. Perez, M. Arenas, C. Gutierrez, Semantics and Complexity of SPARQL, *arXiv:cs/0605124v1 [cs.DB]* (May 2006).  
URL <http://arxiv.org/abs/cs.DB/0605124>
- [17] J. Perez, M. Arenas, C. Gutierrez, Semantics of SPARQL, *Tech. Rep. TR/DCC-2006-17*, Universidad de Chile (October 2006).
- [18] P. Pinheiro da Silva, D. McGuinness, R. Fikes, A Proof Markup Language for Semantic Web services, *Inf. Syst.* 31 (4-5) (2006) 381–395.
- [19] A. Polleres, From sparql to rules (and back), in: *WWW '07: Proceedings of the 16th international conference on World Wide Web*, ACM Press, New York, NY, USA, 2007.  
URL <http://dx.doi.org/10.1145/1242572.1242679>
- [20] E. Prud'hommeaux, A. Seaborne, SPARQL query language for RDF, Working draft, W3C, <http://www.w3.org/TR/rdf-sparql-query/> (March 2007).
- [21] M. Schraefel, N. Shadbolt, N. Gibbins, S. Harris, H. Glaser, CS AKTive Space: Representing Computer Science in the Semantic Web, *Proc. of WWW (2004)* 384–392.
- [22] B. Schueler, S. Sizov, S. Staab, D. T. Tran, Querying for meta knowledge, in: *WWW '08: Proceeding of the 17th international conference on World Wide Web*, ACM, New York, NY, USA, 2008.

## Appendix A. SPARQL Semantics

SPARQL is a query language for RDF based on graph pattern matching, which is defined in [20]. Query results in SPARQL are given by partial substitutions (mapping) of the query variables by RDF terms. The following definitions describe existing foundations of SPARQL introduced in [15,17].

### Definition A.1 (Triple and Basic Graph Patterns)

A tuple  $t$  in  $(U \cup L \cup V) \times (U \cup V) \times (U \cup L \cup V)$  is a triple pattern. A Basic Graph Pattern is a finite set of triple patterns. Given a triple pattern  $t$ ,  $\text{var}(t)$  is a set of variables occurring in  $t$ . Similarly, given a basic graph pattern  $P$ ,  $\text{var}(P) = \bigcup_{t \in P} \text{var}(t)$ , i.e.  $\text{var}(P)$  is the set of variables occurring in  $P$ .

### Definition A.2 (Mapping)

A mapping  $\mu$  from  $V$  to  $U \cup L$  is a partial function  $\mu : V \rightarrow U \cup L$ . The domain of  $\mu$ ,  $\text{dom}(\mu)$ , is a subset of  $V$  where  $\mu$  is defined. The empty mapping  $\mu_\emptyset$  is a mapping such that  $\text{dom}(\mu_\emptyset) = \emptyset$  (i.e.  $\mu_\emptyset = \emptyset$ ). Two mappings  $\mu_1$  and  $\mu_2$  are compatible when for all  $x \in \text{dom}(\mu_1) \cap \text{dom}(\mu_2)$ , it is the case that  $\mu_1(x) = \mu_2(x)$ , then  $\mu_1 \cup \mu_2$  is also a mapping.

Let  $\Omega_1$  and  $\Omega_2$  be set of mappings; join, union, the difference, and left outer-join between  $\Omega_1$  and  $\Omega_2$  are defined as:

- $\Omega_1 \bowtie \Omega_2 = \{\mu_1 \cup \mu_2 \mid \mu_1 \in \Omega_1, \mu_2 \in \Omega_2 \text{ are compatible mappings}\}$ ,
- $\Omega_1 \cup \Omega_2 = \{\mu \mid \mu \in \Omega_1 \text{ or } \mu \in \Omega_2\}$ ,
- $\Omega_1 \setminus \Omega_2 = \{\mu \in \Omega_1 \mid \text{for all } \mu' \in \Omega_2, \mu \text{ and } \mu' \text{ are not compatible}\}$ ,
- $\Omega_1 \bowtie \Omega_2 = (\Omega_1 \bowtie \Omega_2) \cup (\Omega_1 \setminus \Omega_2)$ .

### Definition A.3 (Basic Graph Pattern and Mappings)

Given a basic graph pattern  $P$  and a mapping  $\mu$  such that  $\text{var}(P) \subseteq \text{dom}(\mu)$ , we have that  $\mu(P) = \bigcup_{t \in P} \{\mu(t)\}$ , i.e.  $\mu(P)$  is the set of triples obtained by replacing the variables in the triples of  $P$  according to  $\mu$ .

### Definition A.4 (Basic Graph Pattern Matching)

Let  $G$  be an RDF graph over a RDF dataset  $D$ , and  $P$  a basic graph pattern. The evaluation of  $P$  over  $G$ , denoted by  $\llbracket P \rrbracket_G^D$  is defined as the set of mappings:

$$\llbracket P \rrbracket_G^D = \{\mu : V \rightarrow U \cup L \mid \text{dom}(\mu) = \text{var}(P) \text{ and } \mu(P) \subseteq G\}$$

If  $\mu \in \llbracket P \rrbracket_G^D$ , we say that  $\mu$  is a solution for  $P$  in  $G$ .

### Definition A.5 (Complex Graph Pattern Matching)

Let  $G$  be an RDF graph over a RDF dataset  $D$ , and  $P, P_1$  and  $P_2$  basic graph patterns. The evaluation of  $P, P_1$ , and  $P_2$  over  $G$  is defined recursively as follow:

- $\llbracket P \rrbracket_G^D$  is given by definition A.4,
- $\llbracket P_1 \text{ AND } P_2 \rrbracket_G^D = \llbracket P_1 \rrbracket_G^D \bowtie \llbracket P_2 \rrbracket_G^D$ ,
- $\llbracket P_1 \text{ OPT } P_2 \rrbracket_G^D = \llbracket P_1 \rrbracket_G^D \bowtie \llbracket P_2 \rrbracket_G^D$ ,
- $\llbracket P_1 \text{ UNION } P_2 \rrbracket_G^D = \llbracket P_1 \rrbracket_G^D \cup \llbracket P_2 \rrbracket_G^D$ ,
- $\llbracket P_1 \text{ FILTER } C \rrbracket_G^D = \sigma_c(\llbracket P_1 \rrbracket_G^D)$ ,
- $\llbracket u \text{ GRAPH } P \rrbracket_G^D = \llbracket P \rrbracket_{gr(u)_D}^D$
- $\llbracket ?X \text{ GRAPH } P \rrbracket_G^D = \bigcup_{v \in \text{names}(D)} (\llbracket P \rrbracket_{gr(u)_D}^D \Rightarrow \mu_{?X \rightarrow v})$

Finally, [15] defines the notion of equivalence for graph patterns.

### Definition A.6 (Equivalence of Graph Patterns)

Two graph pattern expressions  $P_1$  and  $P_2$  are equivalent, denoted by  $P_1 \equiv P_2$ , if  $\llbracket P_1 \rrbracket_G^D = \llbracket P_2 \rrbracket_G^D$  for every graph  $G$  and RDF dataset  $D$ .

**Proposition A.1** Let  $P_1, P_2$  and  $P_3$  be graph pattern expressions and  $R$  a built-in condition then:

- AND and UNION are associative and commutative.
- $(P_1 \text{ AND } (P_2 \text{ UNION } P_3)) \equiv ((P_1 \text{ AND } P_2) \text{ UNION } (P_1 \text{ AND } P_3))$
- $(P_1 \text{ OPT } (P_2 \text{ UNION } P_3)) \equiv ((P_1 \text{ OPT } P_2) \text{ UNION } (P_1 \text{ OPT } P_3))$ .
- $((P_1 \text{ UNION } P_2) \text{ OPT } P_3) \equiv ((P_1 \text{ OPT } P_3) \text{ UNION } (P_2 \text{ OPT } P_3))$ .
- $((P_1 \text{ UNION } P_2) \text{ FILTER } R) \equiv ((P_1 \text{ FILTER } R) \text{ UNION } (P_2 \text{ FILTER } R))$ .