

# A Context-Based Semantics for SPARQL Property Paths over the Web

Olaf Hartig<sup>1</sup> and Giuseppe Pirrò<sup>2</sup>

<sup>1</sup> University of Waterloo, Canada  
ohartig@uwaterloo.ca

<sup>2</sup> Institute for High Performance Computing and Networking, ICAR-CNR, Rende, Italy  
pirro@icar.cnr.it

**Abstract** As of today, there exists no standard language for querying Linked Data *on the Web*, where navigation across distributed data sources is a key feature. A natural candidate seems to be SPARQL, which recently has been enhanced with navigational capabilities thanks to the introduction of *property paths* (PPs). However, the semantics of SPARQL restricts the scope of navigation via PPs to *single* RDF graphs. This restriction limits the applicability of PPs on the Web. To fill this gap, in this paper we provide formal foundations for evaluating PPs on the Web, thus contributing to the definition of a query language for Linked Data. In particular, we introduce a query semantics for PPs that couples navigation at the data level with navigation on the Web graph. Given this semantics we find that for some PP-based SPARQL queries a complete evaluation on the Web is not feasible. To enable systems to identify queries that can be evaluated completely, we establish a decidable syntactic property of such queries.

## 1 Introduction

The increasing trend in sharing and interlinking pieces of structured data on the World Wide Web (WWW) is evolving the classical Web—which is focused on hypertext documents and syntactic links among them—into a Web of Linked Data. The Linked Data principles [4] present an approach to extend the scope of Uniform Resource Identifiers (URIs) to new types of resources (e.g., people, places) and represent their descriptions and interlinks by using the Resource Description Framework (RDF) [16] as standard data format. RDF adopts a graph-based data model, which can be queried upon by using the SPARQL query language [12]. When it comes to Linked Data on the WWW, the common way to provide query-based access is via SPARQL endpoints, that is, services that usually answer SPARQL queries over a single dataset. Recently, the original core of SPARQL has been extended with features supporting query federation; it is now possible, within a single query, to target multiple endpoints (via the `SERVICE` operator). However, such an extension is not enough to cope with an unbounded and a priori unknown space of data sources such as the WWW. Moreover, not all Linked Data on the WWW is accessible via SPARQL endpoints. Hence, as of today, there exists no standard query language for Linked Data on the WWW, although SPARQL is clearly a candidate.

While earlier research on using SPARQL for Linked Data is limited to fragments of the first version of the language [5,13,14,25], the more recent version 1.1 introduces a

feature that is particularly interesting in the context of queries over a graph-like environment such as Linked Data on the WWW. This feature is called *property paths* (PPs) and equips SPARQL with navigational capabilities [12]. However, the standard definition of PPs is limited to single, centralized RDF graphs and, thus, not directly applicable to Linked Data that is distributed over the WWW. Therefore, toward the definition of a language for accessing Linked Data live on the WWW, the following questions emerge naturally: “*How can PPs be defined over the WWW?*” and “*What are the implications of such a definition?*” Answering these questions is the broad objective of this paper. To this end, we make the following main contributions:

1. We formalize a query semantics for PP-based SPARQL queries that are meant to be evaluated over Linked Data on the WWW. This semantics is *context-based*; it intertwines Web graph navigation with navigation at the level of data.
2. We study the feasibility of evaluating queries under this semantics. We assume that query engines do not have complete information about the queried Web of Linked Data (as it is the case for the WWW). Our study shows that there exist cases in which query evaluation under the context-based semantics is not feasible.
3. We provide a decidable syntactic property of queries for which an evaluation under the context-based semantics is feasible.

The remainder of the paper is organized as follows. Section 2 provides an overview on related work. Section 3 introduces the formal framework for this paper, including a data model that captures a notion of Linked Data. In Section 4 we focus on PPs, independently from other SPARQL operators. In Section 5 we broaden our view to study PP-based SPARQL graph patterns; we characterize a class of *Web-safe* patterns and prove their feasibility. Finally, in Section 6 we conclude and sketch future work.

## 2 Related Work

The idea of querying the WWW as a database is not new (see Florescu et al.’s survey [11]). Perhaps the most notable early works in this context are by Konopnicki and Shmueli [18], Abiteboul and Vianu [1], and Mendelzon et al. [20], all of which tackled the problem of evaluating SQL-like queries on the traditional hypertext Web. While such queries included navigational features, the focus was on retrieving specific Web pages, particular attributes of specific pages, or content within them.

From a graph-oriented perspective, languages for the *navigation and specification* of vertices in graphs have a long tradition (see Wood’s survey [26]). In the RDF world, extensions of SPARQL such as PPARQL [2], nSPARQL [21], and SPARQLer [17] introduced navigational features since those were missing in the first version of SPARQL. Only recently, with the addition of *property paths* (PPs) in version 1.1 [12], SPARQL has been enhanced officially with such features. The final definition of PPs has been influenced by research that studied the computational complexity of an early draft version of PPs [3,19], and there also already exists a proposal to extend PPs with more expressive power [9]. However, the main assumption of all these navigational extensions of SPARQL is to work on a single, centralized RDF graph. Our departure point is different: *We aim at defining semantics of SPARQL queries (including property paths)*

over *Linked Data on the WWW*, which involves dealing with two graphs of different types; namely, an RDF graph that is distributed over documents on the WWW and the Web graph of how these documents are interlinked with each other.

To express queries over Linked Data on the WWW, two main strands of research can be identified. The first studies how to extend the scope of SPARQL queries to the WWW, with existing work focusing on basic graph patterns [5,13,25] or a more expressive fragment that includes AND, OPT, UNION and FILTER [14]. The second strand focuses on navigational languages such as NautiLOD [8,10]. These two strands have different departure points. The former employs navigation over the WWW to collect data for answering a given SPARQL query; here navigation is a means to discover query-relevant data. The latter provides explicit navigational features and uses querying capabilities to filter data sources of interest; here navigation (not querying) is the main focus. The context-based query semantics proposed in this paper combines both approaches. We believe that the outcome of this research can be a starting point toward the definition of a language for querying and navigating over Linked Data on the WWW.

### 3 Formal Framework

This section provides a formal framework for studying semantics of PPs over Linked Data. We first recall the definition of PPs as per the SPARQL standard [12]. Thereafter, we introduce a data model that captures the notion of Linked Data on the WWW.

#### 3.1 Preliminaries

Assume four pairwise disjoint, countably infinite sets  $\mathcal{I}$  (IRIs),  $\mathcal{B}$  (blank nodes),  $\mathcal{L}$  (literals), and  $\mathcal{V}$  (variables). An *RDF triple* (or simply *triple*) is a tuple from the set  $\mathcal{T} = (\mathcal{I} \cup \mathcal{B}) \times \mathcal{I} \times (\mathcal{I} \cup \mathcal{B} \cup \mathcal{L})$ . For any triple  $t \in \mathcal{T}$  we write  $\text{iris}(t)$  to denote the set of IRIs in that triple. A set of triples is called an *RDF graph*.

A *property path pattern* (or *PP pattern* for short) is a tuple  $P = \langle \alpha, \text{path}, \beta \rangle$  such that  $\alpha, \beta \in (\mathcal{I} \cup \mathcal{L} \cup \mathcal{V})$  and  $\text{path}$  is a *property path expression* (*PP expression*) defined by the following grammar (where  $u, u_1, \dots, u_n \in \mathcal{I}$ ):

$$\text{path} = u \mid !(u_1 \mid \dots \mid u_n) \mid \wedge \text{path} \mid \text{path}/\text{path} \mid (\text{path} \mid \text{path}) \mid (\text{path})^*$$

Note that the SPARQL standard introduces additional types of PP expressions [12]. Since these are merely syntactic sugar (they are defined in terms of expressions covered by the grammar given above), we ignore them in this paper. As another slight deviation from the standard, we do not permit blank nodes in PP patterns (i.e.,  $\alpha, \beta \notin \mathcal{B}$ ). However, standard PP patterns with blank nodes can be simulated using fresh variables.

**Example 1.** An example of a PP pattern is  $\langle \text{Tim}, (\text{knows})^*/\text{name}, ?n \rangle$ , which retrieves the names of persons that can be reached from Tim by an arbitrarily long path of knows relationships (which includes Tim). Another example are the two PP patterns  $\langle ?p, \text{knows}, \text{Tim} \rangle$  and  $\langle \text{Tim}, \wedge \text{knows}, ?p \rangle$ , both of which retrieve persons that know Tim.

The (standard) query semantics of PP patterns is defined by an evaluation function that returns multisets of *solution mappings* where a solution mapping  $\mu$  is a partial function

|   |  |
|---|--|
| <b>Function</b> $\text{ALP1}(\gamma, \text{path}, G)$<br><b>Input:</b> $\gamma \in (\mathcal{I} \cup \mathcal{B} \cup \mathcal{L})$ ,<br>path is a PP expression,<br>$G$ is an RDF graph. | <b>Function</b> $\text{ALP2}(\gamma, \text{path}, \text{Visited}, G)$<br><b>Input:</b> $\gamma \in (\mathcal{I} \cup \mathcal{B} \cup \mathcal{L})$ , path is a PP expression,<br>$\text{Visited} \subseteq (\mathcal{I} \cup \mathcal{B} \cup \mathcal{L})$ , $G$ is an RDF graph.                                    |
| 1: $\text{Visited} := \emptyset$<br>2: $\text{ALP2}(\gamma, \text{path}, \text{Visited}, G)$<br>3: <b>return</b> $\text{Visited}$   | 4: <b>if</b> $\gamma \notin \text{Visited}$ <b>then</b><br>5:   add $\gamma$ to $\text{Visited}$<br>6: <b>for all</b> $\mu \in \llbracket \langle ?x, \text{path}, ?y \rangle \rrbracket_G$ s.t. $\mu(?x) = \gamma$ <b>do</b><br>7: $\text{ALP2}(\mu(?y), \text{path}, \text{Visited}, G)$ // $?x, ?y \in \mathcal{V}$ |

**Figure 1. Auxiliary functions for defining the semantics of PP expressions of the form  $\text{path}^*$ .**

$\mu : \mathcal{V} \rightarrow (\mathcal{I} \cup \mathcal{B} \cup \mathcal{L})$ . Given a solution mapping  $\mu$  and a PP pattern  $P$ , we write  $\mu[P]$  to denote the PP pattern obtained by replacing the variables in  $P$  according to  $\mu$  (unbound variables must not be replaced). Two solution mappings, say  $\mu_1$  and  $\mu_2$ , are *compatible*, denoted by  $\mu_1 \sim \mu_2$ , if  $\mu_1(?v) = \mu_2(?v)$  for all variables  $?v \in (\text{dom}(\mu_1) \cap \text{dom}(\mu_2))$ .

We represent a *multiset* of solution mappings by a pair  $M = \langle \Omega, \text{card} \rangle$  where  $\Omega$  is the underlying set (of solution mappings) and  $\text{card} : \Omega \rightarrow \{1, 2, \dots\}$  is the corresponding *cardinality function*. By abusing notation slightly, we write  $\mu \in M$  for all  $\mu \in \Omega$ . Furthermore, we introduce a family of special (parameterized) cardinality functions that shall simplify the definition of any multiset whose solution mappings all have a cardinality of 1. That is, for any set of solution mappings  $\Omega$ , let  $\text{card}1^{(\Omega)} : \Omega \rightarrow \{1, 2, \dots\}$  be the *constant-1 cardinality function* that is defined by  $\text{card}1^{(\Omega)}(\mu) = 1$  for all  $\mu \in \Omega$ .

To define the aforementioned evaluation function we also need to introduce several SPARQL algebra operators. Let  $M_1 = \langle \Omega_1, \text{card}_1 \rangle$  and  $M_2 = \langle \Omega_2, \text{card}_2 \rangle$  be multisets of solution mappings and let  $V \subseteq \mathcal{V}$  be a finite set of variables. Then:

$M_1 \sqcup M_2 = \langle \Omega, \text{card} \rangle$  where  $\Omega = \Omega_1 \cup \Omega_2$  and (i)  $\text{card}(\mu) = \text{card}_1(\mu)$  for all solution mappings  $\mu \in \Omega \setminus \Omega_2$ , (ii)  $\text{card}(\mu) = \text{card}_2(\mu)$  for all  $\mu \in \Omega \setminus \Omega_1$ , and (iii)  $\text{card}(\mu) = \text{card}_1(\mu) + \text{card}_2(\mu)$  for all  $\mu \in \Omega_1 \cap \Omega_2$ .

$M_1 \bowtie M_2 = \langle \Omega, \text{card} \rangle$  where  $\Omega = \{ \mu_1 \cup \mu_2 \mid (\mu_1, \mu_2) \in \Omega_1 \times \Omega_2 \text{ and } \mu_1 \sim \mu_2 \}$  and, for every  $\mu \in \Omega$ ,  $\text{card}(\mu) = \sum_{(\mu_1, \mu_2) \in \Omega_1 \times \Omega_2 \text{ s.t. } \mu = \mu_1 \cup \mu_2} \text{card}(\mu_1) \cdot \text{card}(\mu_2)$ .

$M_1 \setminus M_2 = \langle \Omega, \text{card} \rangle$  where  $\Omega = \{ \mu_1 \in \Omega_1 \mid \mu_1 \not\sim \mu_2 \text{ for all } \mu_2 \in \Omega_2 \}$  and, for every  $\mu \in \Omega$ ,  $\text{card}(\mu) = \text{card}_1(\mu)$ .

$\pi_V(M_1) = \langle \Omega, \text{card} \rangle$  where  $\Omega = \{ \mu \mid \exists \mu' \in \Omega_1 : \mu \sim \mu' \text{ and } \text{dom}(\mu) = V \cap \text{dom}(\mu') \}$  and, for every  $\mu \in \Omega$ ,  $\text{card}(\mu) = \sum_{\mu' \in \Omega_1 \text{ s.t. } \mu \sim \mu'} \text{card}_1(\mu')$ .

In addition to these algebra operators, the SPARQL standard introduces auxiliary functions to define the semantics of PP patterns of the form  $\langle \alpha, \text{path}^*, \beta \rangle$ . Figure 1 provides these functions—which we call ALP1 and ALP2—adapted to our formalism.<sup>3</sup>

We are now ready to define the standard query semantics of PP patterns.

**Definition 1.** *The evaluation of a PP pattern  $P$  over an RDF graph  $G$ , denoted by  $\llbracket P \rrbracket_G$ , is a multiset of solution mappings  $\langle \Omega, \text{card} \rangle$  that is defined recursively as given in Figure 2 where  $\alpha, \beta \in (\mathcal{I} \cup \mathcal{L} \cup \mathcal{V})$ ,  $x_L, x_R \in (\mathcal{I} \cup \mathcal{L})$ ,  $?v_L, ?v_R \in \mathcal{V}$ ,  $u, u_1, \dots, u_n \in \mathcal{I}$ ,  $?v \in \mathcal{V}$  is a fresh variable, and  $\mu_\emptyset$  denotes the empty solution mapping ( $\text{dom}(\mu_\emptyset) = \emptyset$ ).*

<sup>3</sup> Variable  $?x$  in line 6 is necessary since PP patterns in our formalism do not have blank nodes.

$$\begin{aligned}
\llbracket \langle \alpha, u, \beta \rangle \rrbracket_G &= \left\langle \{ \mu \mid \text{dom}(\mu) = (\{\alpha, \beta\} \cap \mathcal{V}) \text{ and } \mu[\langle \alpha, u, \beta \rangle] \in G \}, \text{card1}^{(\Omega)} \right\rangle \\
\llbracket \langle \alpha, !(u_1 \mid \dots \mid u_n), \beta \rangle \rrbracket_G &= \left\langle \{ \mu \mid \text{dom}(\mu) = (\{\alpha, \beta\} \cap \mathcal{V}) \text{ and} \right. \\
&\quad \left. \exists \mu[\langle \alpha, u, \beta \rangle] \in G : u \in (\mathcal{I} \setminus \{u_1, \dots, u_n\}) \}, \text{card1}^{(\Omega)} \right\rangle \\
\llbracket \langle \alpha, \hat{\text{path}}, \beta \rangle \rrbracket_G &= \llbracket \langle \beta, \text{path}, \alpha \rangle \rrbracket_G \\
\llbracket \langle \alpha, \text{path}_1 / \text{path}_2, \beta \rangle \rrbracket_G &= \pi_{\{\alpha, \beta\} \cap \mathcal{V}} \left( \llbracket \langle \alpha, \text{path}_1, ?v \rangle \rrbracket_G \bowtie \llbracket \langle ?v, \text{path}_2, \beta \rangle \rrbracket_G \right) \\
\llbracket \langle \alpha, (\text{path}_1 \mid \text{path}_2), \beta \rangle \rrbracket_G &= \llbracket \langle \alpha, \text{path}_1, \beta \rangle \rrbracket_G \sqcup \llbracket \langle \alpha, \text{path}_2, \beta \rangle \rrbracket_G \\
\llbracket \langle x_L, (\text{path})^*, ?v_R \rangle \rrbracket_G &= \left\langle \{ \mu \mid \text{dom}(\mu) = \{?v_R\} \text{ and } \mu(?v_R) \in \text{ALP1}(x_L, \text{path}, G) \}, \text{card1}^{(\Omega)} \right\rangle \\
\llbracket \langle ?v_L, (\text{path})^*, ?v_R \rangle \rrbracket_G &= \left\langle \{ \mu \mid \text{dom}(\mu) = \{?v_L, ?v_R\} \text{ and } \mu(?v_L) \in \text{terms}(G) \text{ and} \right. \\
&\quad \left. \mu(?v_R) \in \text{ALP1}(\mu(?v_L), \text{path}, G) \}, \text{card1}^{(\Omega)} \right\rangle \\
\llbracket \langle ?v_L, (\text{path})^*, x_R \rangle \rrbracket_G &= \llbracket \langle x_R, (\hat{\text{path}})^*, ?v_L \rangle \rrbracket_G \\
\llbracket \langle x_L, (\text{path})^*, x_R \rangle \rrbracket_G &= \left\langle \begin{cases} \{\mu_\emptyset\} & \text{if } \exists \mu \in \llbracket \langle x_L, (\text{path})^*, ?v \rangle \rrbracket_G : \mu(?v) = x_R, \\ \emptyset & \text{else} \end{cases}, \text{card1}^{(\Omega)} \right\rangle
\end{aligned}$$

**Figure 2. SPARQL 1.1 W3C property paths semantics.**

### 3.2 Data Model

The standard SPARQL evaluation function for PP patterns (cf. Section 3.1) defines the expected result of the evaluation of a pattern over a single RDF graph. Since the WWW is not an RDF graph, the standard definition is insufficient as a formal foundation for evaluating PP patterns over Linked Data on the WWW. To provide a suitable definition we need a data model that captures the notion of a Web of Linked Data. To this end, we adopt the data model proposed in our earlier work [14]. Here, a *Web of Linked Data (WoLD)* is a tuple  $W = \langle D, \text{data}, \text{adoc} \rangle$  consisting of (i) a set  $D$  of so called *Linked Data documents (documents)*, (ii) a mapping  $\text{data} : D \rightarrow 2^{\mathcal{T}}$  that maps each document to a finite set of RDF triples (representing the data that can be obtained from the document), and (iii) a partial mapping  $\text{adoc} : \mathcal{I} \rightarrow D$  that maps (some) IRIs to a document and, thus, captures a IRI-based retrieval of documents. In this paper we assume that the set of documents  $D$  in any WoLD  $W = \langle D, \text{data}, \text{adoc} \rangle$  is finite, in which case we say  $W$  is *finite* (for a discussion of infiniteness refer to our earlier work [14]).

A few other concepts are needed for the subsequent discussion. For any two documents  $d, d' \in D$  in a WoLD  $W = \langle D, \text{data}, \text{adoc} \rangle$ , document  $d$  has a *data link* to  $d'$  if the data of  $d$  mentions an IRI  $u \in \mathcal{I}$  (i.e., there exists a triple  $\langle s, p, o \rangle \in \text{data}(d)$  with  $u \in \{s, p, o\}$ ) that can be used to retrieve  $d'$  (i.e.,  $\text{adoc}(u) = d'$ ). Such data links establish the *link graph* of the WoLD  $W$ , that is, a directed graph  $\langle D, E \rangle$  in which the edges  $E$  are all pairs  $\langle d, d' \rangle \in D \times D$  for which  $d$  has a data link to  $d'$ . Note that this graph, as well as the tuple  $\langle D, \text{data}, \text{adoc} \rangle$  typically are not available directly to systems that aim to compute queries over the Web captured by  $W$ . For instance, the complete domain of the partial mapping  $\text{adoc}$  (i.e., all IRIs that can be used to retrieve some document) is unknown to such systems and can only be disclosed partially (by trying to look up IRIs). Also note that the link graph of a WoLD is a different type of graph than the RDF “graph” whose triples are distributed over the documents in the WoLD.

## 4 Web-aware Query Semantics for Property Paths

We are now ready to introduce our framework, which does not deal with syntactic aspects of PPs but aims at defining query semantics that provide a formal foundation for using PP patterns as queries over a WoLD (and, thus, over Linked Data on the WWW).

### 4.1 Full-Web Query Semantics

As a first approach we may assume a full-Web query semantics that is based on the standard evaluation function (as introduced in Section 3.1) and defines an expected query result for any PP pattern in terms of *all data* on the queried WoLD. Formally:

**Definition 2.** *Let  $P$  be a PP pattern, let  $W = \langle D, data, adoc \rangle$  be a WoLD, and let  $G^*$  be an RDF graph such that  $G^* = \bigcup_{d \in D} data(d)$ , then the evaluation of  $P$  over  $W$  under full-Web semantics, denoted by  $\llbracket P \rrbracket_W^{\text{fw}}$ , is defined by  $\llbracket P \rrbracket_W^{\text{fw}} = \llbracket P \rrbracket_{G^*}$ .*

We emphasize that the full-Web query semantics is mostly of theoretical interest. In practice, that is, for a WoLD  $W$  that represents the “real” WWW (as it runs on the Internet), there cannot exist a system that guarantees to compute the given evaluation function  $\llbracket \cdot \rrbracket_W^{\text{fw}}$  over  $W$  using an algorithm that both terminates and returns complete query results. In earlier work, we showed such a limitation for evaluating other types of SPARQL graph patterns—including triple patterns—under a corresponding full-Web query semantics defined for these patterns [14]. This result readily carries over to the full-Web query semantics for PP patterns because any PP pattern  $P = \langle \alpha, path, \beta \rangle$  with PP expression `path` being an IRI  $u \in \mathcal{I}$  is, in fact, a triple pattern  $\langle \alpha, u, \beta \rangle$ . Informally, we explain this negative result by the fact that the three structures  $D$ ,  $data$ , and  $adoc$  that capture the queried Web formally, are not available in practice. Consequently, to enumerate the set of all triples on the Web (i.e., the RDF graph  $G^*$  in Definition 2), a query execution system would have to enumerate all documents (the set  $D$ ); given that such a system has limited access to mapping  $adoc$  (in particular,  $\text{dom}(adoc)$ )—the set of all IRIs whose lookup retrieves a document—is, at best, partially known, the only guarantee to discover all documents is to look up any possible (HTTP-scheme) IRI. Since these are infinitely many [7], the enumeration process cannot terminate.

### 4.2 Context-Based Query Semantics

Given the limited practical applicability of full-Web query semantics for PPs, we propose an alternative query semantics that interprets PP patterns as a language for navigation over Linked Data on the Web (i.e., along the lines of earlier navigational languages for Linked Data such as NautiLOD [8]). We refer to this semantics as *context-based*.

The main idea behind this query semantics is to restrict the scope of searching for any next triple of a potentially matching path to specific data within specific documents on the queried WoLD. As a basis for formalizing these restrictions we introduce the notion of a *context selector*. Informally, for each IRI that can be used to retrieve a document, the context selector returns a specific subset of the data within that document; this subset contains only those RDF triples that have the given IRI as their subject (such

a set of triples resembles Harth and Speiser’s notion of subject authoritative triples [13]). Formally, for any WoLD  $W = \langle D, data, adoc \rangle$ , the context selector of  $W$  is a function  $C^W: \mathcal{I} \cup \mathcal{B} \cup \mathcal{L} \cup \mathcal{V} \rightarrow 2^{\mathcal{T}}$  that, for each  $\gamma \in (\mathcal{I} \cup \mathcal{B} \cup \mathcal{L} \cup \mathcal{V})$ , is defined as follows:<sup>4</sup>

$$C^W(\gamma) = \begin{cases} \{ \langle s, p, o \rangle \in data(adoc(\gamma)) \mid \gamma = s \} & \text{if } \gamma \in \mathcal{I} \text{ and } \gamma \in \text{dom}(adoc), \\ \emptyset & \text{otherwise.} \end{cases}$$

Informally, we explain how a context selector restricts the scope of PP patterns over a WoLD as follows. Suppose a sequence of triples  $\langle s_1, p_1, o_1 \rangle, \dots, \langle s_k, p_k, o_k \rangle$  presents a path that already matches a sub-expression of a given PP expression. Under the previously defined full-Web query semantics (cf. Section 4.1), the next triple for such a path can be searched for in an arbitrary document in the queried WoLD  $W$ . By contrast, under the context-based query semantics, the next triple has to be searched for only in  $C^W(o_k)$ . Given these preliminaries, we now define context-based semantics:

**Definition 3.** *Let  $P$  be a PP pattern and let  $W = \langle D, data, adoc \rangle$  be a WoLD. The evaluation of  $P$  over  $W$  under context-based semantics, denoted by  $\llbracket P \rrbracket_W^{ctx}$ , returns a multiset of solution mappings  $\langle \Omega, card \rangle$  defined recursively as given in Figure 3, where  $u, \dots, u_n \in \mathcal{I}$ ;  $x_L, x_R \in (\mathcal{I} \cup \mathcal{L})$ ;  $?v_L, ?v_R \in \mathcal{V}$ ;  $\mu_\emptyset$  is the empty solution mapping (i.e.,  $\text{dom}(\mu_\emptyset) = \emptyset$ ); function ALPW1 is given in Figure 4; and  $?v \in \mathcal{V}$  is a fresh variable.*

There are three points worth mentioning w.r.t. Definition 3: First, note how the context selector restricts the data that has to be searched to find matching triples (e.g., consider the first line in Figure 3). Second, we emphasize that context-based query semantics is defined such that it resembles the standard semantics of PP patterns as close as possible (cf. Section 3.1). Therefore, for the part of our definition that covers PP patterns of the form  $\langle \alpha, \text{path}^*, \beta \rangle$ , we also use auxiliary functions—ALPW1 and ALPW2 (cf. Figure 4). These functions evaluate the sub-expression `path` recursively over the queried WoLD (instead of using a fixed RDF graph as done in the standard semantics in Figure 1). Third, the two base cases with a variable in the subject position (i.e., the third and the sixth line in Figure 3) require an enumeration of all IRIs. Such a requirement is necessary to preserve consistency with the standard semantics, as well as to preserve commutativity of operators that can be defined on top of PP patterns (such as the `AND` operator in SPARQL; cf. Section 5). However, due to this requirement there exist PP patterns whose (complete) evaluation under context-based semantics is infeasible when querying the WWW. The following example describes such a case.

**Example 2.** *Consider the PP pattern  $P_{E2} = \langle ?v, \text{knows}, \text{Tim} \rangle$ , which asks for the IRIs of people that know Tim. Under context-based semantics, any IRI  $u'$  can be used to generate a correct solution mapping for the pattern as long as a lookup of that IRI results in retrieving a document whose data includes the triple  $\langle u', \text{knows}, \text{Tim} \rangle$ . While, for any WoLD that is finite, there exists only a finite number of such IRIs, determining these IRIs and guaranteeing completeness requires to enumerate the infinite set of all IRIs and to check each of them (unless one knows the complete—and finite—subset of*

<sup>4</sup> To simplify the following formalization of context-based semantics, context selectors are defined not only over IRIs, but also over blank nodes, literals, and variables.

$$\begin{aligned}
\llbracket \langle u_L, p, \beta \rangle \rrbracket_W^{\text{ctx}} &= \langle \{ \mu \mid \text{dom}(\mu) = (\{\beta\} \cap \mathcal{V}) \text{ and } \mu[\langle u_L, p, \beta \rangle] \in C^W(u_L) \}, \text{card1}^{(\Omega)} \rangle \\
\llbracket \langle l_L, p, \beta \rangle \rrbracket_W^{\text{ctx}} &= \langle \emptyset, \text{card1}^{(\emptyset)} \rangle \\
\llbracket \langle ?v_L, p, \beta \rangle \rrbracket_W^{\text{ctx}} &= \langle \{ \mu \mid \text{dom}(\mu) = (\{?v_L, \beta\} \cap \mathcal{V}) \text{ and } \\
&\quad \mu[\langle ?v_L, p, \beta \rangle] \in \bigcup_{u \in \mathcal{I}} C^W(u) \}, \text{card1}^{(\Omega)} \rangle \\
\llbracket \langle u_L, !(u_1 \mid \dots \mid u_n), \beta \rangle \rrbracket_W^{\text{ctx}} &= \langle \{ \mu \mid \text{dom}(\mu) = (\{\beta\} \cap \mathcal{V}) \text{ and } \\
&\quad \exists \mu[\langle u_L, p, \beta \rangle] \in C^W(u_L) : p \notin \{u_1, \dots, u_n\} \}, \text{card1}^{(\Omega)} \rangle \\
\llbracket \langle l_L, !(u_1 \mid \dots \mid u_n), \beta \rangle \rrbracket_W^{\text{ctx}} &= \langle \emptyset, \text{card1}^{(\emptyset)} \rangle \\
\llbracket \langle ?v_L, !(u_1 \mid \dots \mid u_n), \beta \rangle \rrbracket_W^{\text{ctx}} &= \langle \{ \mu \mid \text{dom}(\mu) = (\{?v_L, \beta\} \cap \mathcal{V}) \text{ and } \\
&\quad \exists \mu[\langle ?v_L, p, \beta \rangle] \in \bigcup_{u \in \mathcal{I}} C^W(u) : p \notin \{u_1, \dots, u_n\} \}, \text{card1}^{(\Omega)} \rangle \\
\llbracket \langle \alpha, \wedge \text{path}, \beta \rangle \rrbracket_W^{\text{ctx}} &= \llbracket \langle \beta, \text{path}, \alpha \rangle \rrbracket_W^{\text{ctx}} \\
\llbracket \langle \alpha, \text{path}_1 / \text{path}_2, \beta \rangle \rrbracket_W^{\text{ctx}} &= \pi_{\{\alpha, \beta\} \cap \mathcal{V}} \left( \llbracket \langle \alpha, \text{path}_1, ?v \rangle \rrbracket_W^{\text{ctx}} \bowtie \llbracket \langle ?v, \text{path}_2, \beta \rangle \rrbracket_W^{\text{ctx}} \right) \\
\llbracket \langle \alpha, \text{path}_1 \mid \text{path}_2, \beta \rangle \rrbracket_W^{\text{ctx}} &= \llbracket \langle \alpha, \text{path}_1, \beta \rangle \rrbracket_W^{\text{ctx}} \sqcup \llbracket \langle \alpha, \text{path}_2, \beta \rangle \rrbracket_W^{\text{ctx}} \\
\llbracket \langle x_L, (\text{path})^*, ?v_R \rangle \rrbracket_W^{\text{ctx}} &= \langle \{ \mu \mid \text{dom}(\mu) = \{?v_R\} \text{ and } \mu(?v_R) \in \text{ALPW1}(x_L, \text{path}, W) \}, \text{card1}^{(\Omega)} \rangle \\
\llbracket \langle ?v_L, (\text{path})^*, ?v_R \rangle \rrbracket_W^{\text{ctx}} &= \langle \{ \mu \mid \text{dom}(\mu) = \{?v_L, ?v_R\} \text{ and } \mu(?v_L) \in \text{terms}(W) \text{ and } \\
&\quad \mu(?v_R) \in \text{ALWP1}(\mu(?v_L), \text{path}, W) \}, \text{card1}^{(\Omega)} \rangle \\
\llbracket \langle ?v_L, (\text{path})^*, x_R \rangle \rrbracket_W^{\text{ctx}} &= \llbracket \langle x_R, (\wedge \text{path})^*, ?v_L \rangle \rrbracket_W^{\text{ctx}} \\
\llbracket \langle x_L, (\text{path})^*, x_R \rangle \rrbracket_W^{\text{ctx}} &= \langle \begin{cases} \{\mu_\emptyset\} & \text{if } \exists \mu \in \llbracket \langle x_L, (\text{path})^*, ?v \rangle \rrbracket_W^{\text{ctx}} : \mu(?v) = x_R, \\ \emptyset & \text{else} \end{cases}, \text{card1}^{(\Omega)} \rangle
\end{aligned}$$

**Figure 3. Context-based query semantics for SPARQL property paths over the Web.**

*all IRIs that can be used to retrieve some document, which, due to the infiniteness of possible HTTP IRIs, cannot be achieved for the WWW).*

It is not difficult to see that the issue illustrated in the example exists for any triple pattern that has a variable in the subject position. On the other hand, triple patterns whose subject is an IRI do not have this issue. However, having an IRI in the subject position is not a sufficient condition in general. For instance, the PP pattern  $\langle \text{Tim}, \wedge \text{knows}, ?v \rangle$  has the same issue as the pattern in Example 2 (in fact, both patterns are semantically equivalent under context-based semantics). A question that arises is whether there exists a property of PP patterns that can be used to distinguish between patterns that do not have this issue (i.e., evaluating them over any WoLD is feasible) and those that do. We shall discuss this question for the more general case of PP-based SPARQL queries.

## 5 SPARQL with Property Paths on the Web

After considering PP patterns in separation, we now turn to a more expressive fragment of SPARQL that embeds PP patterns as the basic building block and uses additional operators on top. We define the resulting PP-based SPARQL queries, discuss the fea-



|   |   |
|---|---|
| <p><b>Function</b> <math>\text{ALPW1}(\gamma, \text{path}, W)</math><br/> <b>Input:</b> <math>\gamma \in (\mathcal{I} \cup \mathcal{B} \cup \mathcal{L})</math>,<br/> path is a PP expression,<br/> <math>W</math> is a WoLD.<br/> 1: <math>Visited := \emptyset</math><br/> 2: <math>\text{ALPW2}(\gamma, \text{path}, Visited, W)</math><br/> 3: <b>return</b> <math>Visited</math></p> | <p><b>Function</b> <math>\text{ALPW2}(\gamma, \text{path}, Visited, W)</math><br/> <b>Input:</b> <math>\gamma \in (\mathcal{I} \cup \mathcal{B} \cup \mathcal{L})</math>, path is a PP expression,<br/> <math>Visited \subseteq (\mathcal{I} \cup \mathcal{B} \cup \mathcal{L})</math>, <math>W</math> is a WoLD.<br/> 4: <b>if</b> <math>\gamma \notin Visited</math> <b>then</b><br/> 5:   add <math>\gamma</math> to <math>Visited</math><br/> 6:   <b>for all</b> <math>\mu \in \llbracket \langle ?x, \text{path}, ?y \rangle \rrbracket_W^{\text{ctx}}</math> s.t. <math>\mu(?x) = \gamma</math> <b>do</b><br/> 7:     <math>\text{ALPW2}(\mu(?y), \text{path}, Visited, W)</math> // <math>?x, ?y \in \mathcal{V}</math></p> |
|---|---|

**Figure 4. Auxiliary functions used for defining context-based query semantics.**

sibility of evaluating these queries over the Web, and introduce a syntactic property to identify queries for which an evaluation under context-based semantics is feasible.

### 5.1 Definition

By using the algebraic syntax of SPARQL [22], we define a *graph pattern* recursively as follows: (i) Any PP pattern  $\langle \alpha, \text{path}, \beta \rangle$  is a graph pattern; and (ii) if  $P_1$  and  $P_2$  are graph patterns, then  $(P_1 \text{ AND } P_2)$ ,  $(P_1 \text{ UNION } P_2)$ , and  $(P_1 \text{ OPT } P_2)$  are graph patterns.<sup>5</sup> For any graph pattern  $P$ , we write  $\forall(P)$  to denote the set of *all variables* in  $P$ .

By using PP patterns as the basic building block of graph patterns, we can readily carry over our context-based semantics to graph patterns: For any graph pattern  $P$  and any WoLD  $W$ , the *evaluation* of  $P$  over  $W$  under context-based semantics is a multiset of solution mappings, denoted by  $\llbracket P \rrbracket_W^{\text{ctx}}$ , that is defined recursively as follows:<sup>6</sup>

- If  $P$  is a PP pattern, then  $\llbracket P \rrbracket_W^{\text{ctx}}$  is defined in Definition 3.
- If  $P$  is  $(P_1 \text{ AND } P_2)$ , then  $\llbracket P \rrbracket_W^{\text{ctx}} = \llbracket P_1 \rrbracket_W^{\text{ctx}} \bowtie \llbracket P_2 \rrbracket_W^{\text{ctx}}$ .
- If  $P$  is  $(P_1 \text{ UNION } P_2)$ , then  $\llbracket P \rrbracket_W^{\text{ctx}} = \llbracket P_1 \rrbracket_W^{\text{ctx}} \sqcup \llbracket P_2 \rrbracket_W^{\text{ctx}}$ .
- If  $P$  is  $(P_1 \text{ OPT } P_2)$ , then  $\llbracket P \rrbracket_W^{\text{ctx}} = (\llbracket P_1 \rrbracket_W^{\text{ctx}} \bowtie \llbracket P_2 \rrbracket_W^{\text{ctx}}) \sqcup (\llbracket P_1 \rrbracket_W^{\text{ctx}} \setminus \llbracket P_2 \rrbracket_W^{\text{ctx}})$ .

### 5.2 Discussion

Given a query semantics for evaluating PP-based graph patterns over a WoLD, we now discuss the feasibility of such evaluation. To this end, we introduce the notion of *Web-safeness* of graph patterns. Informally, graph patterns are Web-safe if evaluating them completely under context-based semantics is possible. Formally:

**Definition 4.** A graph pattern  $P$  is Web-safe if there exists an algorithm that, for any finite WoLD  $W = \langle D, \text{data}, \text{adoc} \rangle$ , computes  $\llbracket P \rrbracket_W^{\text{ctx}}$  by looking up only a finite number of IRIs without assuming direct access to the sets  $D$  and  $\text{dom}(\text{adoc})$ .

<sup>5</sup> For this paper we leave out other types of SPARQL graph patterns such as filters. Adding them is an exercise that would not have any significant implication on the following discussion.

<sup>6</sup> Note that the definition uses the algebra operators introduced in Section 3.1.

**Example 3.** Consider graph pattern  $P_{E3} = (\langle \text{Bob, knows, } ?v \rangle \text{ AND } \langle ?v, \text{ knows, Tim} \rangle)$ . The right sub-pattern  $P_{E2} = \langle ?v, \text{ knows, Tim} \rangle$  is not Web-safe because evaluating it completely over the WWW is not feasible under context-based semantics (cf. Example 2). However, the larger pattern  $P_{E3}$  is Web-safe; it can be evaluated completely under context-based semantics. For instance, a possible algorithm may first evaluate the left sub-pattern, which is feasible because it requires the lookup of a single IRI only (the IRI Bob). Thereafter, the evaluation of the right sub-pattern  $P_{E2}$  can be reduced to looking up a finite number of IRIs only, namely the IRIs bound to variable  $?v$  in solution mappings obtained for the left sub-pattern. Although any other IRI  $u^*$  might also be used to discover matching triples for  $P_{E2}$ , each of these triples has IRI  $u^*$  as its subject (which is a consequence of restricting retrieved data based on the context selector introduced in Section 4.2). Therefore, the solution mappings resulting from such matching triples cannot be compatible with any solution for the left sub-pattern and, thus, do not satisfy the join condition established by the semantics of AND in pattern  $P_{E3}$ .

The example illustrates that some graph patterns are Web-safe even if some of their sub-patterns are not. Consequently, we are interested in a *decidable* property that enables to identify Web-safe patterns, including those whose sub-patterns are not Web-safe.

Buil-Aranda et al. study a similar problem in the context of SPARQL federation where graph patterns of the form  $P_S = (\text{SERVICE } ?v P)$  are allowed [6]. Here, variable  $?v$  ranges over a possibly large set of IRIs, each of which represents the address of a (remote) SPARQL service that needs to be called to assemble the complete result of  $P_S$ . However, many service calls may be avoided if  $P_S$  is embedded in a larger graph pattern that allows for an evaluation during which  $?v$  can be bound before evaluating  $P_S$ . To tackle this problem, Buil-Aranda et al. introduce a notion of *strong boundedness* of variables in graph patterns and use it to show a notion of safeness for the evaluation of patterns like  $P_S$  within larger graph patterns. The set of *strongly bound variables* in a graph pattern  $P$ , denoted by  $\text{SBV}(P)$ , is defined recursively as follows:

- If  $P$  is a PP pattern, then  $\text{SBV}(P) = \mathbb{V}(P)$  (recall that  $\mathbb{V}(P)$  are all variables in  $P$ ).
- If  $P$  is of the form  $(P_1 \text{ AND } P_2)$ , then  $\text{SBV}(P) = \text{SBV}(P_1) \cup \text{SBV}(P_2)$ .
- If  $P$  is of the form  $(P_1 \text{ UNION } P_2)$ , then  $\text{SBV}(P) = \text{SBV}(P_1) \cap \text{SBV}(P_2)$ .
- If  $P$  is of the form  $(P_1 \text{ OPT } P_2)$ , then  $\text{SBV}(P) = \text{SBV}(P_1)$ .

The idea behind the notion of strongly bound variables has already been used in earlier work (e.g., “*certain variables*” [23], “*output variables*” [24]), and it is tempting to adopt it for our problem. However, we note that one cannot identify Web-safe graph patterns by using strong boundedness in a manner similar to its use in Buil-Aranda et al.’s work alone. For instance, consider graph pattern  $P_{E3}$  from Example 3. We know that (i)  $P_{E3}$  is Web-safe and that (ii)  $\mathbb{V}(P_{E3}) = \{?v\}$  and also  $\text{SBV}(P_{E3}) = \{?v\}$ . Then, one might hypothesize that for every graph pattern  $P$ , if  $\text{SBV}(P) = \mathbb{V}(P)$ , then  $P$  is Web-safe. However, the PP pattern  $P_{E2} = \langle ?v, \text{ knows, Tim} \rangle$  disproves such a hypothesis because, even if  $\text{SBV}(P_{E2}) = \mathbb{V}(P_{E2})$ , pattern  $P_{E2}$  is not Web-safe (cf. Example 2).

We conjecture the following reason why strong boundedness cannot be used directly for our problem. For complex patterns (i.e., patterns that are not PP patterns), the sets of strongly bound variables of all sub-patterns are defined *independent* from each other, whereas the algorithm outlined in Example 3 leverages a specific relationship between

sub-patterns. More precisely, the algorithm leverages the fact that the same variable that is the subject of the right sub-pattern is also the object of the left sub-pattern.

Based on this observation, we introduce the notion of *conditionally Web-bounded variables*, the definition of which, for complex graph patterns, is based on specific relationships between sub-patterns. This notion shall turn out to be suitable for our case.

**Definition 5.** *The conditionally Web-bounded variables of a graph pattern  $P$  w.r.t. a set of variables  $X$  is the subset  $CBV(P | X) \subseteq V(P)$  that is defined recursively as follows:*

| <i>If <math>P</math> is:</i>   | <i>then <math>CBV(P   X)</math> is:</i>                          |
|--|--|
| 1) $\langle \alpha, u, \beta \rangle$ or $\langle \alpha, !(u_1   \dots   u_n), \beta \rangle$ such that $\alpha \in (\mathcal{I} \cup \mathcal{L})$ or $\alpha \in X$   | $V(P)$   |
| 2) $\langle \alpha, u, \beta \rangle$ or $\langle \alpha, !(u_1   \dots   u_n), \beta \rangle$ such that $\alpha \notin (\mathcal{I} \cup \mathcal{L})$ and $\alpha \notin X$  | $\emptyset$  |
| 3) $\langle \alpha, (\text{path})^*, \beta \rangle$ s.t. $\alpha \in \mathcal{V}$ and $\beta \notin \mathcal{V}$   | $CBV(\langle \beta, (\wedge \text{path})^*, \alpha \rangle   X)$ |
| 4) $\langle \alpha, (\text{path})^*, \beta \rangle$ s.t. (i) $\alpha \notin \mathcal{V}$ or $\beta \in \mathcal{V}$ , and (ii) for any two variables $?x, ?y \in \mathcal{V}$ it holds that $CBV(\langle ?x, \text{path}, ?y \rangle   \{?x\}) = \{?x, ?y\}$     | $CBV(\langle \alpha, \text{path}, \beta \rangle   X)$            |
| 5) $\langle \alpha, (\text{path})^*, \beta \rangle$ such that none of the above  | $\emptyset$  |
| 6) $\langle \alpha, \wedge \text{path}, \beta \rangle$ with $P' = \langle \beta, \text{path}, \alpha \rangle$  | $CBV(P'   X)$  |
| 7) $\langle \alpha, (\text{path}_1   \text{path}_2), \beta \rangle$ with $P' = ((\alpha, \text{path}_1, \beta) \text{ UNION } \langle \alpha, \text{path}_2, \beta \rangle)$   | $CBV(P'   X)$  |
| 8) $\langle \alpha, \text{path}_1 / \text{path}_2, \beta \rangle$ s.t. for any $?v \in \mathcal{V} \setminus (X \cup \{\alpha, \beta\})$ , $?v \in CBV(P'   X)$ where $P' = ((\alpha, \text{path}_1, ?v) \text{ AND } \langle ?v, \text{path}_2, \beta \rangle)$ | $CBV(P'   X) \setminus \{?v\}$                                   |
| 9) $\langle \alpha, \text{path}_1 / \text{path}_2, \beta \rangle$ such that none of the above  | $\emptyset$  |
| 10) $(P_1 \text{ AND } P_2)$ s.t. $CBV(P_1   X) = V(P_1)$ and $CBV(P_2   X) = V(P_2)$  | $V(P)$   |
| 11) $(P_1 \text{ AND } P_2)$ s.t. $CBV(P_1   X) = V(P_1)$ and $CBV(P_2   X \cup SBV(P_1)) = V(P_2)$  | $V(P)$   |
| 12) $(P_1 \text{ AND } P_2)$ s.t. $CBV(P_2   X) = V(P_2)$ and $CBV(P_1   X \cup SBV(P_2)) = V(P_1)$  | $V(P)$   |
| 13) $(P_1 \text{ AND } P_2)$ such that none of the above   | $\emptyset$  |
| 14) $(P_1 \text{ UNION } P_2)$   | $CBV(P_1   X) \cap CBV(P_2   X)$                                 |
| 15) $(P_1 \text{ OPT } P_2)$ s.t. $CBV(P_1   X) = V(P_1)$ and $CBV(P_2   X) = V(P_2)$  | $V(P)$   |
| 16) $(P_1 \text{ OPT } P_2)$ s.t. $CBV(P_1   X) = V(P_1)$ and $CBV(P_2   X \cup SBV(P_1)) = V(P_2)$  | $V(P)$   |
| 17) $(P_1 \text{ OPT } P_2)$ such that none of the above   | $\emptyset$  |

**Example 4.** For the PP pattern  $P_{E2} = \langle ?v, \text{knows}, \text{Tim} \rangle$ —which is not Web-safe (as discussed in Example 2)—if we use the set  $\{?v\}$  as condition, then, by line 1 in Definition 5, it holds that  $CBV(P_{E2} | \{?v\}) = \{?v\}$ . However, if we use the empty set instead, we obtain  $CBV(P_{E2} | \emptyset) = \emptyset$  (cf. line 2 in Definition 5).

While for the non-Web-safe pattern  $P_{E2}$  we thus observe  $CBV(P_{E2} | \emptyset) \neq V(P_{E2})$ , for graph pattern  $P_{E3} = (\langle \text{Bob}, \text{knows}, ?v \rangle \text{ AND } \langle ?v, \text{knows}, \text{Tim} \rangle)$ —which is Web-safe (cf. Example 3)—we have  $CBV(P_{E3} | \emptyset) = V(P_{E3})$ . The fact that  $CBV(P_{E3} | \emptyset) = \{?v\}$  follows from (i)  $CBV(\langle \text{Bob}, \text{knows}, ?v \rangle | \emptyset) = \{?v\}$ , (ii)  $SBV(\langle \text{Bob}, \text{knows}, ?v \rangle) = \{?v\}$ , (iii)  $CBV(\langle ?v, \text{knows}, \text{Tim} \rangle | \{?v\}) = \{?v\}$ , and (iv) line 11 in Definition 5.

The example seems to suggest that, if *all* variables of a graph pattern are conditionally Web-bounded w.r.t. the empty set of variables, then the graph pattern is Web-safe. The following result verifies this hypothesis.

**Theorem 1.** *A graph pattern  $P$  is Web-safe if  $CBV(P | \emptyset) = V(P)$ .*

**Note 1.** *Due to the recursive nature of Definition 5, the condition  $CBV(P | \emptyset) = V(P)$  (as used in Theorem 1) is decidable for any graph pattern  $P$ .*

We prove Theorem 1 based on an algorithm that evaluates graph patterns recursively by passing (intermediate) solution mappings to recursive calls. To capture the desired results of each recursive call formally, we introduce a special evaluation function for a graph pattern  $P$  over a WoLD  $W$  that takes a solution mapping  $\mu$  as input and returns only the solutions for  $P$  over  $W$  that are compatible with  $\mu$ .

**Definition 6.** Let  $P$  be a graph pattern, let  $W$  be a WoLD, and let  $\langle \Omega, \text{card} \rangle = \llbracket P \rrbracket_W^{\text{ctx}}$ . Given a solution mapping  $\mu$ , the  $\mu$ -restricted evaluation of  $P$  over  $W$  under context-based semantics, denoted by  $\llbracket P \mid \mu \rrbracket_W^{\text{ctx}}$ , is the multiset of solution mappings  $\langle \Omega', \text{card}' \rangle$  with  $\Omega' = \{ \mu' \in \Omega \mid \mu' \sim \mu \}$  and  $\text{card}'(\mu') = \text{card}(\mu')$  for all  $\mu' \in \Omega'$ .

The following lemma shows the existence of the aforementioned recursive algorithm.

**Lemma 1.** Let  $P$  be a graph pattern and let  $\mu_{\text{in}}$  be a solution mapping. If it holds that  $\text{CBV}(P \mid \text{dom}(\mu_{\text{in}})) = \vee(P)$ , there exists an algorithm that, for any finite WoLD  $W$ , computes  $\llbracket P \mid \mu_{\text{in}} \rrbracket_W^{\text{ctx}}$  by looking up a finite number of IRIs only.

Before providing the proof of the lemma (and of Theorem 1), we point out two important properties of Definition 6. First, it is easily seen that, for any graph pattern  $P$  and WoLD  $W$ ,  $\llbracket P \mid \mu_{\emptyset} \rrbracket_W^{\text{ctx}} = \llbracket P \rrbracket_W^{\text{ctx}}$ , where  $\mu_{\emptyset}$  is the empty solution mapping (i.e.,  $\text{dom}(\mu_{\emptyset}) = \emptyset$ ). Consequently, given an algorithm, say  $A$ , that has the properties of the algorithm described by Lemma 1, a trivial algorithm that can be used to prove Theorem 1 may simply call algorithm  $A$  with the empty solution mapping and return the result of this call (we shall elaborate more on this approach in the proof of Theorem 1 below). Second, for any PP pattern  $\langle \alpha, \text{path}, \beta \rangle$  and WoLD  $W$ , if  $\alpha$  is a variable and  $\text{path}$  is a base PP expression (i.e., one of the first two cases in the grammar in Section 3.1), then  $\llbracket P \mid \mu \rrbracket_W^{\text{ctx}}$  is empty for every solution mapping  $\mu$  that binds (variable)  $\alpha$  to a literal or a blank node. Formally, we show the latter as follows.

**Lemma 2.** Let  $P$  be a PP pattern of the form  $\langle ?v, u, \beta \rangle$  or  $\langle ?v, !(u_1 \mid \dots \mid u_n), \beta \rangle$  with  $?v \in \mathcal{V}$  and  $u, u_1, \dots, u_n \in \mathcal{I}$ , and let  $\mu$  be a solution mapping. If  $?v \in \text{dom}(\mu)$  and  $\mu(?v) \in (\mathcal{B} \cup \mathcal{L})$ , then, for any WoLD  $W$ ,  $\llbracket P \mid \mu \rrbracket_W^{\text{ctx}}$  is the empty multiset.

*Proof (Lemma 2).* Recall that, for any IRI  $u$  and any WoLD  $W$ , context  $C^W(u)$  contains only triples that have IRI  $u$  as their subject. As a consequence, for any WoLD  $W$ , every solution mapping  $\mu' \in \llbracket P \rrbracket_W^{\text{ctx}}$  binds variable  $?v$  to some IRI (and never to a literal or blank node); i.e.,  $\mu'(?v) \in \mathcal{I}$ . Therefore, if  $?v \in \text{dom}(\mu)$  and  $\mu(?v) \in (\mathcal{B} \cup \mathcal{L})$ , then  $\mu$  cannot be compatible with any  $\mu' \in \llbracket P \rrbracket_W^{\text{ctx}}$  and, thus,  $\llbracket P \mid \mu \rrbracket_W^{\text{ctx}}$  is empty.  $\square$

We use Lemma 2 to prove Lemma 1 as follows.

*Proof idea (Lemma 1).* We prove the lemma by induction on the possible structure of graph pattern  $P$ . For the proof, we provide Algorithm 1 and show that this (recursive) algorithm has the desired properties for any possible graph pattern (i.e., any case of the induction, including the base case). Due to space limitations, in this paper we only present a fragment of the algorithm and highlight essential properties thereof. The given fragment covers the base case (lines 1-11) and one pivotal case of the induction step, namely, graph patterns of the form  $(P_1 \text{ AND } P_2)$  (lines 57-72). The complete version of the algorithm and the full proof can be found in an extended version of this paper [15].

For the base case, Algorithm 1 looks up at most one IRI (cf. lines 2-5). The crux of showing that the returned result is sound and complete is Lemma 2 and the fact that the only possible *context* in which a triple  $\langle s, p, o \rangle$  with  $s \in \mathcal{I}$  can be found is  $C^W(s)$ .

For PP patterns of the form  $(P_1 \text{ AND } P_2)$  consider lines 57-72. By using Definition 5, we show  $\text{CBV}(P_i \mid \text{dom}(\mu_{\text{in}})) = \vee(P_i)$  and  $\text{CBV}(P_j \mid \text{dom}(\mu_{\text{in}}) \cup \text{dom}(\mu)) = \vee(P_j)$

---

**Algorithm 1**  $EvalCtxBased(P, \mu_{in})$ , which computes  $\llbracket P \mid \mu_{in} \rrbracket_W^{ctx}$ .

---

```

1: if  $P$  is of the form  $\langle \alpha, u, \beta \rangle$  or  $P$  is of the form  $\langle \alpha, !(u_1 \mid \dots \mid u_n), \beta \rangle$  then
2:   if  $\alpha \in \mathcal{I}$  then  $u' := \alpha$ 
3:   else if  $\alpha \in \mathcal{V}$  and  $\alpha \in \text{dom}(\mu_{in})$  and  $\mu_{in}(\alpha) \in \mathcal{I}$  then  $u' := \mu_{in}(\alpha)$ 
4:   else  $u' := \text{null}$ 

5:   if  $u'$  is an IRI and looking it up results in retrieving a document, say  $d$  then
6:      $G :=$  the set of triples in  $d$  (use a fresh set of blank node identifiers when parsing  $d$ )
7:      $G' := \{ \langle s, p, o \rangle \in G \mid s = u' \}$ 
8:      $\langle \Omega, card \rangle := \llbracket P \rrbracket_{G'}$  ( $\llbracket P \rrbracket_{G'}$  can be computed by using any algorithm that
       implements the standard SPARQL evaluation function)
9:     return a new multiset  $\langle \Omega', card' \rangle$  with  $\Omega' = \{ \mu' \in \Omega \mid \mu' \sim \mu_{in} \}$  and
        $card'(\mu') = card(\mu')$  for all  $\mu' \in \Omega'$ 

10:  else
11:    return a new empty multiset  $\langle \Omega, card \rangle$  with  $\Omega = \emptyset$  and  $\text{dom}(card) = \emptyset$ 

...

57: else if  $P$  is of the form  $(P_1 \text{ AND } P_2)$  then
58:   if  $\text{CBV}(P_1 \mid \text{dom}(\mu_{in})) = \text{V}(P_1)$  then  $i := 1; j := 2$  else  $i := 2; j := 1$ 
59:   Create a new empty multiset  $M = \langle \Omega, card \rangle$  with  $\Omega = \emptyset$  and  $\text{dom}(card) = \emptyset$ 
60:    $\langle \Omega^{P_i}, card^{P_i} \rangle := EvalCtxBased(P_i, \mu_{in})$ 
61:   for all  $\mu \in \Omega^{P_i}$  do
62:      $\langle \Omega^\mu, card^\mu \rangle := EvalCtxBased(P_j, \mu_{in} \cup \mu)$ 
63:     for all  $\mu' \in \Omega^\mu$  do
64:        $\mu^* := \mu \cup \mu'$ 
65:        $k := card^{P_i}(\mu) \cdot card^\mu(\mu')$ 
66:       if  $\mu^* \in \Omega$  then
67:          $old := card(\mu^*)$ 
68:         Adjust  $card$  such that  $card(\mu^*) = k + old$ 
69:       else
70:         Adjust  $card$  such that  $card(\mu^*) = k$ 
71:         Add  $\mu^*$  to  $\Omega$ 
72:   return  $M$ 

```

---

for all  $\mu \in \Omega^{P_i}$ . Therefore, by induction, all recursive calls (lines 60 and 62) look up a finite number of IRIs and return correct results; i.e.,  $\langle \Omega^{P_i}, card^{P_i} \rangle = \llbracket P_i \mid \mu_{in} \rrbracket_W^{ctx}$  and  $\langle \Omega^\mu, card^\mu \rangle = \llbracket P_j \mid \mu_{in} \cup \mu \rrbracket_W^{ctx}$  for all  $\mu \in \Omega^{P_i}$ . Then, since each  $\mu \in \Omega^{P_i}$  is compatible with all  $\mu' \in \Omega^\mu$  and all processed solution mappings are compatible with  $\mu_{in}$ , it is easily verified that the computed result is  $\llbracket (P_1 \text{ AND } P_2) \mid \mu_{in} \rrbracket_W^{ctx}$ .  $\square$

We are now ready to prove Theorem 1, for which we use Lemma 1, or more precisely the algorithm that we introduce in the proof of the lemma.

*Proof (Theorem 1).* Let  $P$  be a graph pattern s.t.  $\text{CBV}(P \mid \emptyset) = \text{V}(P)$ . Then, given the empty solution mapping  $\mu_\emptyset$  with  $\text{dom}(\mu_\emptyset) = \emptyset$ , we have  $\text{CBV}(P \mid \text{dom}(\mu_\emptyset)) = \text{V}(P)$ . Therefore, by our proof of Lemma 1 we know that, for any finite WoLD  $W$ , Algorithm 1 computes  $\llbracket P \mid \mu_\emptyset \rrbracket_W^{ctx}$  by looking up a finite number of IRIs. We also know that the empty solution mapping is compatible with any solution mapping. Consequently, by Definition 6,  $\llbracket P \mid \mu_\emptyset \rrbracket_W^{ctx} = \llbracket P \rrbracket_W^{ctx}$  for any WoLD  $W$ . Hence, by passing the empty solu-

tion mapping to it, Algorithm 1 can be used to compute  $\llbracket P \rrbracket_W^{ct,x}$  for any finite WoLD  $W$ , and during this computation the algorithm looks up a finite number of IRIs only.  $\square$

While the condition in Theorem 1 is sufficient to identify Web-safe graph patterns, the question that remains is whether it is a necessary condition (in which case it could be used to decide Web-safeness of *all* graph patterns). Unfortunately, the answer is no.

**Example 5.** Consider the graph pattern  $P = (P_1 \text{ UNION } P_2)$  with  $P_1 = \langle u_1, p_1, ?x \rangle$  and  $P_2 = \langle u_2, p_2, ?y \rangle$ . We note that  $CBV(P_1 | \emptyset) = \{?x\}$  and  $CBV(P_2 | \emptyset) = \{?y\}$ , and, thus,  $CBV(P | \emptyset) = \emptyset$ . Hence, the pattern does not satisfy the condition in Theorem 1. Nonetheless, it is easy to see that there exists a (sound and complete) algorithm that, for any WoLD  $W$ , computes  $\llbracket P \rrbracket_W^{ct,x}$  by looking up a finite number of IRIs only. For instance, such an algorithm, say  $A$ , may first use two other algorithms that compute  $\llbracket P_1 \rrbracket_W^{ct,x}$  and  $\llbracket P_2 \rrbracket_W^{ct,x}$  by looking up a finite number of IRIs, respectively. Such algorithms exist by Theorem 1, because  $CBV(P_1 | \emptyset) = \mathcal{V}(P_1)$  and  $CBV(P_2 | \emptyset) = \mathcal{V}(P_2)$ . Finally, algorithm  $A$  can generate the (sound and complete) query result  $\llbracket P \rrbracket_W^{ct,x}$  by computing the multiset union  $\llbracket P_1 \rrbracket_W^{ct,x} \sqcup \llbracket P_2 \rrbracket_W^{ct,x}$ , which requires no additional IRI lookups.

**Remark 1.** The example illustrates that “only if” cannot be shown in Theorem 1. It remains an open question whether there exists an alternative condition for Web-safeness that is both sufficient and necessary (and decidable).

## 6 Concluding Remarks and Future Work

This paper studies the problem of extending the scope of SPARQL property paths to query Linked Data that is distributed on the WWW. We have proposed a context-based query semantics and analyzed its peculiarities. Our perhaps most interesting finding is that there exist queries whose evaluation over the WWW is not feasible. We studied this aspect and introduced a decidable syntactic property for identifying feasible queries.

We believe that the presented work provides valuable input to a wider discussion about defining a language for accessing Linked Data on the WWW. In this context, there are several directions for future research such as the following three. First, studying a more expressive navigational core for property paths over the Web; e.g., along the lines of other navigational languages such as nSPARQL [21] or NautiLOD [8]. Second, investigating relationships between navigational queries and SPARQL federation. Third, while the aim of this paper was to introduce a formal foundation for answering SPARQL queries with PPs over Linked Data on the WWW, an investigation of how systems may implement efficiently the machinery developed in this paper is certainly interesting.

## References

1. Abiteboul, S., Vianu, V.: Queries and Computation on the Web. Theor. Comput. Sci. 239(2), 231–255 (2000)
2. Alkhateeb, F., Baget, J.F., Euzenat, J.: Extending SPARQL with Regular Expression Patterns (for querying RDF). J. Web Sem. 7(2), 57–73 (2009)

3. Arenas, M., Conca, S., Pérez, J.: Counting Beyond a Yottabyte, or how SPARQL 1.1 Property Paths will Prevent Adoption of the Standard. In: Proceedings of the 21st International Conference on World Wide Web (2012)
4. Berners-Lee, T.: Design issues: Linked Data. Online (Jul 2006)
5. Bouquet, P., Ghidini, C., Serafini, L.: Querying The Web Of Data: A Formal Approach. In: Proceedings of the 4th Asian Semantic Web Conference (2009)
6. Buil-Aranda, C., Arenas, M., Corcho, O., Polleres, A.: Federating Queries in SPARQL1.1: Syntax, Semantics and Evaluation. *Journal on Web Semantics* 18(1), 1–17 (2013)
7. Fielding, R., Gettys, J., Mogul, J.C., Frystyk, H., Masinter, L., Leach, P.J., Berners-Lee, T.: Hypertext Transfer Protocol – HTTP/1.1. RFC 2616 (Jun 1999)
8. Fionda, V., Gutierrez, C., Pirrò, G.: Semantic Navigation on the Web of Data: Specification of Routes, Web Fragments and Actions. In: Proceedings of the 21st International Conference on the World Wide Web (2012)
9. Fionda, V., Pirrò, G., Consens, M.: Extended Property Paths: Writing More SPARQL Queries in a Succinct Way. In: Proceedings of the 28th AAAI Conference on Artificial Intelligence (AAAI) (2015)
10. Fionda, V., Pirrò, G., Gutierrez, C.: NautiLOD: A Formal Language for the Web of Data Graph. *ACM Trans. Web* 9(1) (Jan 2015)
11. Florescu, D., Levy, A., Mendelzon, A.: Database Techniques for the World-Wide Web: A Survey. *SIGMOD Rec.* 27, 59–74 (1998)
12. Harris, S., Seaborne, A.: SPARQL 1.1 Query Language. W3C Recommendation (2013)
13. Harth, A., Speiser, S.: On Completeness Classes for Query Evaluation on Linked Data. In: Proceedings of the 26th AAAI Conference (2012)
14. Hartig, O.: SPARQL for a Web of Linked Data: Semantics and Computability. In: Proceedings of the 9th Extended Semantic Web Conference (2012)
15. Hartig, O., Pirrò, G.: A Context-Based Semantics for SPARQL Property Paths over the Web (Extended Version). CoRR abs/1503.04831 (2015), <http://arxiv.org/abs/1503.04831>
16. Klyne, G., Carroll, J.J.: Resource Description Framework (RDF): Concepts and Abstract Syntax (2006)
17. Kochut, K.J., Janik, M.: SPARQLeR: Extended SPARQL for Semantic Association Discovery. In: *The Semantic Web: Research and Applications*, pp. 145–159. Springer (2007)
18. Konopnicki, D., Shmueli, O.: Information Gathering in the World-Wide Web: The W3QL Query Language and the W3QS System. *ACM Transactions on Database Systems* 23(4), 369–410 (Dec 1998)
19. Loseman, K., Martens, W.: The Complexity of Evaluating Path Expressions in SPARQL. In: Proceedings of the 31st ACM Symposium on Principles of Database Systems (2012)
20. Mendelzon, A.O., Mihaila, G.A., Milo, T.: Querying the World Wide Web. In: 1 (ed.) *Int. J. on Digital Libraries*, vol. 1, pp. 54–97 (1997)
21. Pérez, J., Arenas, M., Gutierrez, C.: nSPARQL: A Navigational Language for RDF. *Journal on Web Semantics* 8(4), 255–270 (2010)
22. Pérez, J., Arenas, M., Gutierrez, C.: Semantics and Complexity of SPARQL. *ACM Transactions on Database Systems (TODS)* 34(3) (2009)
23. Schmidt, M., Meier, M., Lausen, G.: Foundations of SPARQL Query Optimization. In: Proceedings of the 13th International Conference on Database Theory (2010)
24. Toman, D., Weddell, G.E.: *Fundamentals of Physical Design and Query Compilation*. Synthesis Lectures on Data Management, Morgan & Claypool Publishers (2011)
25. Umbrich, J., Hogan, A., Polleres, A., Decker, S.: Link Traversal Querying for a diverse Web of Data. *Semantic Web Journal* (2014)
26. Wood, P.T.: Query Languages for Graph Databases. *SIGMOD Rec.* 41(1) (2012)