

# From SPARQL to Rules (and back) \*

Axel Polleres  
 Universidad Rey Juan Carlos  
 Tulipán s/n, 28933 Móstoles, Madrid, Spain  
 axel@polleres.net

## ABSTRACT

As the data and ontology layers of the Semantic Web stack have achieved a certain level of maturity in standard recommendations such as RDF and OWL, the current focus lies on two related aspects. On the one hand, the definition of a suitable query language for RDF, SPARQL, is close to recommendation status within the W3C. The establishment of the rules layer on top of the existing stack on the other hand marks the next step to be taken, where languages with their roots in Logic Programming and Deductive Databases are receiving considerable attention. The purpose of this paper is threefold. First, we discuss the formal semantics of SPARQL extending recent results in several ways. Second, we provide translations from SPARQL to Datalog with negation as failure. Third, we propose some useful and easy to implement extensions of SPARQL, based on this translation. As it turns out, the combination serves for direct implementations of SPARQL on top of existing rules engines as well as a basis for more general rules and query languages on top of RDF.

## Categories and Subject Descriptors

H.2.3 [Languages]: Query Languages; H.3.5 [Online Information Services]: Web-based services

## General Terms

Languages, Standardization

## Keywords

SPARQL, Datalog, Rules

## 1. INTRODUCTION

After the data and ontology layers of the Semantic Web stack have achieved a certain level of maturity in standard recommendations such as RDF and OWL, the query and the rules layers seem to be the next building-blocks to be finalized. For the first part, SPARQL [18], W3C's proposed query language, seems to be close to recommendation, though the Data Access working group is still struggling

with defining aspects such as a formal semantics or layering on top of OWL and RDFS. As for the second part, the RIF working group <sup>1</sup>, who is responsible for the rules layer, is just producing first concrete results. Besides aspects like business rules exchange or reactive rules, deductive rules languages on top of RDF and OWL are of special interest to the RIF group. One such deductive rules language is Datalog, which has been successfully applied in areas such as deductive databases and thus might be viewed as a query language itself. Let us briefly recap our starting points:

**Datalog and SQL.** Analogies between Datalog and relational query languages such as SQL are well-known and -studied. Both formalisms cover UCQ (unions of conjunctive queries), where Datalog adds recursion, particularly unrestricted recursion involving nonmonotonic negation (aka unstratified negation as failure). Still, SQL is often viewed to be more powerful in several respects. On the one hand, the lack of recursion has been partly solved in the standard's 1999 version [20]. On the other hand, aggregates or external function calls are missing in pure Datalog. However, also developments on the Datalog side are evolving and with recent extensions of Datalog towards Answer Set Programming (ASP) – a logic programming paradigm extending and building on top of Datalog – lots of these issues have been solved, for instance by defining a declarative semantics for aggregates [9], external predicates [8].

**The Semantic Web rules layer.** Remarkably, logic programming dialects such as Datalog with nonmonotonic negation which are covered by Answer Set Programming are often viewed as a natural basis for the Semantic Web rules layer [7]. Current ASP systems offer extensions for retrieving RDF data and querying OWL knowledge bases from the Web [8]. Particular concerns in the Semantic Web community exist with respect to adding rules including nonmonotonic negation [3] which involve a form of closed world reasoning on top of RDF and OWL which both adopt an open world assumption. Recent proposals for solving this issue suggest a “safe” use of negation as failure over finite contexts only for the Web, also called *scoped negation* [17].

**The Semantic Web query layer – SPARQL.** Since we base our considerations in this paper on the assumption that similar correspondences as between SQL and Datalog can be established for SPARQL, we have to observe that SPARQL inherits a lot from SQL, but there also remain substantial differences: On the one hand, SPARQL does not deal with nested queries or recursion, a detail which is indeed surpris-

\*An extended technical report of this article is available at <http://www.polleres.net/publications/>.

Copyright is held by the International World Wide Web Conference Committee (IW3C2). Distribution of these papers is limited to classroom use, and personal use by others.

WWW 2007, May 8–12, 2007, Banff, Alberta, Canada.  
 ACM 978-1-59593-654-7/07/0005.

<sup>1</sup><http://www.w3.org/2005/rules/wg>

ing by the fact that SPARQL is a graph query language on RDF where, typical recursive queries such as transitive closure of a property might seem very useful. Likewise, aggregation (such as count, average, etc.) of object values in RDF triples which might appear useful have not yet been included in the current standard. On the other hand, subtleties like blank nodes (aka bNodes), or optional graph patterns, which are similar but (as we will see) different to outer joins in SQL or relational algebra, are not straightforwardly translatable to Datalog.

The goal of this paper is to shed light on the actual relation between declarative rules languages such as Datalog and SPARQL, and by this also provide valuable input for the currently ongoing discussions on the Semantic Web rules layer, in particular its integration with SPARQL, taking the likely direction into account that LP style rules languages will play a significant role in this context.

Although the SPARQL specification does not seem 100% stable at the current point, just having taken a step back from candidate recommendation to working draft, we think that it is not too early for this exercise, as we will gain valuable insights and positive side effects by our investigation. More precisely, the contributions of the present work are:

- We refine and extend a recent proposal to formalize the semantics of SPARQL from Pérez et al. [16], presenting three variants, namely c-joining, s-joining and b-joining semantics where the latter coincides with [16], and can thus be considered normative. We further discuss how aspects such compositionality, or idempotency of joins are treated in these semantics.
- Based on the three semantic variants, we provide translations from a large fragment of SPARQL queries to Datalog, which give rise to implementations of SPARQL on top of existing engines.
- We provide some straightforward extensions of SPARQL such as a set difference operator MINUS, and nesting of ASK queries in FILTER expressions.
- Finally, we discuss an extension towards recursion by allowing bNode-free-CONSTRUCT queries as part of the query dataset, which may be viewed as a lightweight, recursive rule language on top of of RDF.

The remainder of this paper is structured as follows: In Sec. 2 we first overview SPARQL, discuss some issues in the language (Sec. 2.1) and then define its formal semantics (Sec. 2.2). After introducing a general form of Datalog with negation as failure under the answer set semantics in Sec. 3, we proceed with the translations of SPARQL to Datalog in Sec. 4. We finally discuss the above-mentioned language extensions in Sec. 5, before we conclude in Sec. 6.

## 2. RDF AND SPARQL

In examples, we will subsequently refer to the two RDF graphs in Fig. 1 which give some information about *Bob* and *Alice*. Such information is common in FOAF files which are gaining popularity to describe personal data. Similarities with existing examples in [18] are on purpose. We assume the two RDF graphs given in TURTLE [2] notation and accessible via the IRIs [ex.org/bob](http://ex.org/bob) and [alice.org](http://alice.org)<sup>2</sup>

<sup>2</sup>For reasons of legibility and conciseness, we omit the leading 'http://' or other schema identifiers in IRIs.

We assume the pairwise disjoint, infinite sets  $I$ ,  $B$ ,  $L$  and  $Var$ , which denote IRIs, Blank nodes, RDF literals, and variables respectively. In this paper, an *RDF Graph* is then a finite set, of triples from  $I \cup B \cup L \times I \times I \cup B \cup L$ ,<sup>3</sup> dereferenceable by an IRI. A SPARQL *query* is a quadruple  $Q = (V, P, DS, SM)$ , where  $V$  is a result form,  $P$  is a graph pattern,  $DS$  is a dataset, and  $SM$  is a set of solution modifiers. We refer to [18] for syntactical details and will explain these in the following as far as necessary. In this paper, we will ignore solution modifiers mostly, thus we will usually write queries as triples  $Q = (V, P, DS)$ , and will use the syntax for graph patterns introduced below.

**Result Forms.** Since we will, to a large extent, restrict ourselves to SELECT queries, it is sufficient for our purposes to describe result forms by sets variables. Other result forms will be discussed in Sec. 5. For instance, let  $Q = (V, P, DS)$  denote the query from Fig. 1, then  $V = \{?X, ?Y\}$ . Query results in SPARQL are given by partial, i.e. possibly incomplete, substitutions of variables in  $V$  by RDF terms. In traditional relational query languages, such incompleteness is usually expressed using null values. Using such null values we will write solutions as tuples where the order of columns is determined by *lexicographically ordering* the variables in  $V$ . Given a set of variables  $V$ , let  $\overline{V}$  denote the tuple obtained from lexicographically ordering  $V$ .

The query from Fig. 1 with result form  $\overline{V} = (?X, ?Y)$  then has solution tuples  $(\text{"Bob"}, \_ : a)$ ,  $(\text{"Alice"}, \text{alice.org\#me})$ ,  $(\text{"Bob"}, \_ : c)$ . We write substitutions in square brackets, so these tuples correspond to the substitutions  $[?X \rightarrow \text{"Bob"}, ?Y \rightarrow \_ : a]$ ,  $[?X \rightarrow \text{"Alice"}, ?Y \rightarrow \text{alice.org\#me}]$ , and  $[?X \rightarrow \text{"Bob"}, ?Y \rightarrow \_ : c]$ , respectively.

**Graph Patterns.** We follow the recursive definition of graph patterns  $P$  from [16]:

- a tuple  $(s, p, o)$  is a graph pattern where  $s, o \in I \cup L \cup Var$  and  $p \in I \cup Var$ .<sup>4</sup>
- if  $P$  and  $P'$  are graph patterns then  $(P \text{ AND } P')$ ,  $(P \text{ OPT } P')$ ,  $(P \text{ UNION } P')$ ,  $(P \text{ MINUS } P')$  are graph patterns.<sup>5</sup>
- if  $P$  is a graph pattern and  $i \in I \cup Var$ , then  $(\text{GRAPH } i \text{ } P)$  is a graph pattern.
- if  $P$  is a graph pattern and  $R$  is a filter expression then  $(P \text{ FILTER } R)$  is a graph pattern.

For any pattern  $P$ , we denote by  $vars(P)$  the set of all variables occurring in  $P$ . As *atomic filter expression*, SPARQL allows the unary predicates BOUND, isBLANK, isIRI, isLITERAL, binary equality predicates '=' for literals, and other features such as comparison operators, data type conversion

<sup>3</sup>Following SPARQL, we are slightly more general than the original RDF specification in that we allow literals in subject positions.

<sup>4</sup>We do not consider bNodes in patterns as these can be semantically equivalently replaced by variables in graph patterns [6].

<sup>5</sup>Note that AND and MINUS are not designated keywords in SPARQL, but we use them here for reasons of readability and in order to keep with the operator style definition of [16]. MINUS is syntactically not present at all, but we will suggest a syntax extension for this particular keyword in Sec. 5.

<pre># Graph: ex.org/bob @prefix foaf: &lt;http://xmlns.com/foaf/0.1/&gt; . @prefix bob: &lt;ex.org/bob#&gt; .  &lt;ex.org/bob&gt; foaf:maker _:a. _:a a foaf:Person ; foaf:name "Bob";     foaf:knows _:b.  _:b a foaf:Person ; foaf:nick "Alice". &lt;alice.org/&gt; foaf:maker _:b</pre>	<pre># Graph: alice.org @prefix foaf: &lt;http://xmlns.com/foaf/0.1/&gt; . @prefix alice: &lt;alice.org#&gt; .  alice:me a foaf:Person ; foaf:name "Alice" ;     foaf:knows _:c.  _:c a foaf:Person ; foaf:name "Bob" ;     foaf:nick "Bobby".</pre>
---	--

PREFIX foaf: <http://xmlns.com/foaf/0.1/>

```
SELECT ?Y ?X
FROM <alice.org>
FROM <ex.org/bob>
WHERE { ?Y foaf:name ?X . }
```

?X	?Y
"Bob"	_:a
"Bob"	_:c
"Alice"	alice.org#me

Figure 1: Two RDF graphs in TURTLE notation and a simple SPARQL query.

and string functions which we omit here, see [18, Sec. 11.3] for details. *Complex filter expressions* can be built using the connectives ' $\neg$ ', ' $\wedge$ ', ' $\vee$ '.

**Datasets.** The dataset  $DS = (G, \{(g_1, G_1), \dots, (g_k, G_k)\})$  of a SPARQL query is defined by a default graph  $G$  plus a set of named graphs, i.e. pairs of IRIs and corresponding graphs. Without loss of generality (there are other ways to define the dataset such as in a SPARQL protocol query), we assume  $G$  given as the merge of the graphs denoted by the IRIs given in a set of FROM and FROM NAMED clauses. For instance, the query from Fig. 1 refers to the dataset which consists of the default graph obtained from merging `alice.org`  $\sqcup$  `ex.org/bob` plus an empty set of named graphs.

The relation between names and graphs in SPARQL is defined solely in terms of that the IRI defines a resource which is represented by the respective graph. In this paper, we assume that the IRIs represent indeed network-accessible resources where the respective RDF-graphs can be retrieved from. This view has also been taken e.g. in [17]. Particularly, this treatment is not to be confused with so-called named graphs in the sense of [4]. We thus identify each IRI with the RDF graph available at this IRI and each set of IRIs with the graph merge [13] of the respective IRIs. This allows us to identify the dataset by a pair of sets of IRIs  $DS = (G, G_n)$  with  $G = \{d_1, \dots, d_n\}$  and  $G_n = \{g_1, \dots, g_k\}$  denoting the (merged) default graph and the set of named graphs, respectively. Hence, the following set of clauses

```
FROM <ex.org/bob>
FROM NAMED <alice.org>
```

defines the dataset  $DS = (\{ex.org/bob\}, \{alice.org\})$ .

## 2.1 Assumptions and Issues

In this section we will discuss some important issues about the current specification, and how we will deal with them here.

First, note that the default graph if specified by name in a FROM clause is not counted among the named graphs automatically [18, section 8, definition 1]. An unbound variable in the GRAPH directive, means any of the named graphs in  $DS$ , but does NOT necessarily include the default graph.

EXAMPLE 1. *This issue becomes obvious in the following*

*query with dataset  $DS = (\{ex.org/bob\}, \emptyset)$  which has an empty solution set.*

```
SELECT ?N WHERE { ?G foaf:maker ?M .
    GRAPH ?G { ?X foaf:name ?N } }
```

We will sometimes find the following assumption convenient to avoid such arguably unintuitive effects:

**Definition 1. (Dataset closedness assumption)** Given a dataset  $DS = (G, G_n)$ ,  $G_n$  implicitly contains (i) all graphs mentioned in  $G$  and (ii) all IRIs mentioned explicitly in the graphs corresponding to  $G$ .

Under this assumption, the previous query has both ("Alice") and ("Bob") in its solution set.

Some more remarks are in place concerning FILTER expressions. According to the SPARQL specification "Graph pattern matching creates bindings of variables [where] it is possible to further restrict solutions by constraining the allowable bindings of variables to RDF Terms [with FILTER expressions]." However, it is not clearly specified how to deal with filter constraints referring to variables which do not appear in simple graph patterns. In this paper, for graph patterns of the form  $(P \text{ FILTER } R)$  we tacitly assume *safe filter expressions*, i.e. that all variables used in a filter expression  $R$  also appear in the corresponding pattern  $P$ . This corresponds with the notion of safety in Datalog (see Sec.3), where the built-in predicates (which obviously correspond to filter predicates) do not suffice to safe unbound variables.

Moreover, the specification defines errors to avoid mistyped comparisons, or evaluation of built-in functions over unbound values, i.e. "any potential solution that causes an error condition in a constraint will not form part of the final results, but does not cause the query to fail." These errors propagate over the whole FILTER expression, also over negation, as shown by the following example.

EXAMPLE 2. *Assuming the dataset does not contain triples for the foaf : dummy property, the example query*

```
SELECT ?X
WHERE { { ?X a foaf:Person .
    OPTIONAL { ?X foaf:dummy ?Y . } }
    FILTER (  $\neg$ (isLITERAL (?Y)) ) }
```

*would discard any solution for ?X, since the unbound value for ?Y causes an error in the isLITERAL expression and thus the whole FILTER expression returns an error.*

We will take special care for these errors, when defining the semantics of FILTER expressions later on.

## 2.2 Formal Semantics of SPARQL

The semantics of SPARQL is still not formally defined in its current version. This lack of formal semantics has been tackled by a recent proposal of Pérez et al. [16]. We will base on this proposal, but suggest three variants thereof, namely (a) *bravely joining*, (b) *cautiously-joining*, and (c) *strictly-joining* semantics. Particularly, our definitions vary from [16] in the way we define joining unbound variables. Moreover, we will refine their notion of FILTER satisfaction in order to deal with error propagation properly.

We denote by  $T_{\text{null}}$  the union  $I \cup B \cup L \cup \{\text{null}\}$ , where **null** is a dedicated constant denoting the unknown value not appearing in any of  $I$ ,  $B$ , or  $L$ , how it is commonly introduced when defining outer joins in relational algebra.

A substitution  $\theta$  from  $Var$  to  $T_{\text{null}}$  is a partial function  $\theta : Var \rightarrow T_{\text{null}}$ . We write substitutions in postfix notation: For a triple pattern  $t = (s, p, o)$  we denote by  $t\theta$  the triple  $(s\theta, p\theta, o\theta)$  obtained by applying the substitution to all variables in  $t$ . The *domain* of  $\theta$ ,  $\text{dom}(\theta)$ , is the subset of  $Var$  where  $\theta$  is defined. For a substitution  $\theta$  and a set of variables  $D \subseteq Var$  we define the substitution  $\theta^D$  with domain  $D$  as follows:

$$x\theta^D = \begin{cases} x\theta & \text{if } x \in \text{dom}(\theta) \cap D \\ \text{null} & \text{if } x \in D \setminus \text{dom}(\theta) \end{cases}$$

Let  $\theta_1$  and  $\theta_2$  be substitutions, then  $\theta_1 \cup \theta_2$  is the substitution obtained as follows:

$$x(\theta_1 \cup \theta_2) = \begin{cases} x\theta_1 & \text{if } x\theta_1 \text{ defined and } x\theta_2 \text{ undefined} \\ \text{else: } x\theta_1 & \text{if } x\theta_1 \text{ defined and } x\theta_2 = \text{null} \\ \text{else: } x\theta_2 & \text{if } x\theta_2 \text{ defined} \\ \text{else: undefined} \end{cases}$$

Thus, in the union of two substitutions defined values in one take precedence over null values the other substitution. For instance, given the substitutions  $\theta_1 = [?X \rightarrow \text{"Alice"}, ?Y \rightarrow \text{"a"}, ?Z \rightarrow \text{null}]$  and  $\theta_2 = [?U \rightarrow \text{"Bob"}, ?X \rightarrow \text{"Alice"}, ?Y \rightarrow \text{null}]$  we get:  $\theta_1 \cup \theta_2 = [?U \rightarrow \text{"Bob"}, ?X \rightarrow \text{"Alice"}, ?Y \rightarrow \text{"a"}, ?Z \rightarrow \text{null}]$

Now, as opposed to [16], we define three notions of compatibility between substitutions:

- Two substitutions  $\theta_1$  and  $\theta_2$  are *bravely compatible* (*b-compatible*) when for all  $x \in \text{dom}(\theta_1) \cap \text{dom}(\theta_2)$  either  $x\theta_1 = \text{null}$  or  $x\theta_2 = \text{null}$  or  $x\theta_1 = x\theta_2$  holds. i.e., when  $\theta_1 \cup \theta_2$  is a substitution over  $\text{dom}(\theta_1) \cup \text{dom}(\theta_2)$ .
- Two substitutions  $\theta_1$  and  $\theta_2$  are *cautiously compatible* (*c-compatible*) when they are b-compatible and for all  $x \in \text{dom}(\theta_1) \cap \text{dom}(\theta_2)$  it holds that  $x\theta_1 = x\theta_2$ .
- Two substitutions  $\theta_1$  and  $\theta_2$  are *strictly compatible* (*s-compatible*) when they are c-compatible and for all  $x \in \text{dom}(\theta_1) \cap \text{dom}(\theta_2)$  it holds that  $x(\theta_1 \cup \theta_2) \neq \text{null}$ .

Analogously to [16] we define join, union, difference, and outer join between two sets of substitutions  $\Omega_1$  and  $\Omega_2$  over domains  $D_1$  and  $D_2$ , respectively, all except union parameterized by  $x \in \{b, c, s\}$ :

$$\begin{aligned} \Omega_1 \bowtie_x \Omega_2 &= \{\theta_1 \cup \theta_2 \mid \theta_1 \in \Omega_1, \theta_2 \in \Omega_2, \text{ are } x\text{-compatible}\} \\ \Omega_1 \cup \Omega_2 &= \{\theta \mid \exists \theta_1 \in \Omega_1 \text{ with } \theta = \theta_1^{D_1 \cup D_2} \text{ or} \\ &\quad \exists \theta_2 \in \Omega_2 \text{ with } \theta = \theta_2^{D_1 \cup D_2}\} \\ \Omega_1 -_x \Omega_2 &= \{\theta \in \Omega_1 \mid \forall \theta_2 \in \Omega_2, \theta \text{ and } \theta_2 \text{ not } x\text{-compatible}\} \\ \Omega_1 \sqsupset_x \Omega_2 &= (\Omega_1 \bowtie_x \Omega_2) \cup (\Omega_1 -_x \Omega_2) \end{aligned}$$

The semantics of a graph pattern  $P$  over dataset  $DS = (G, G_n)$ , can now be defined recursively by the evaluation function returning sets of substitutions.

*Definition 2. (Evaluation, extends [16, Def. 2])* Let  $t = (s, p, o)$  be a triple pattern,  $P, P_1, P_2$  graph patterns,  $DS = (G, G_n)$  a dataset, and  $i \in G_n$ , and  $v \in Var$ , then the  $x$ -joining evaluation  $[[\cdot]]_{DS}^x$  is defined as follows:

$$\begin{aligned} [[t]]_{DS}^x &= \{\theta \mid \text{dom}(\theta) = \text{vars}(P) \text{ and } t\theta \in G\} \\ [[P_1 \text{ AND } P_2]]_{DS}^x &= [[P_1]]_{DS}^x \bowtie_x [[P_2]]_{DS}^x \\ [[P_1 \text{ UNION } P_2]]_{DS}^x &= [[P_1]]_{DS}^x \cup [[P_2]]_{DS}^x \\ [[P_1 \text{ MINUS } P_2]]_{DS}^x &= [[P_1]]_{DS}^x -_x [[P_2]]_{DS}^x \\ [[P_1 \text{ OPT } P_2]]_{DS}^x &= [[P_1]]_{DS}^x \sqsupset_x [[P_2]]_{DS}^x \\ [[\text{GRAPH } i \ P]]_{DS}^x &= [[P]]_{DS}^{i, \emptyset} \\ [[\text{GRAPH } v \ P]]_{DS}^x &= \{\theta \cup [v \rightarrow g] \mid g \in G_n, \theta \in [[P[v \rightarrow g]]]_{(g, \emptyset)}^x\} \\ [[P \text{ FILTER } R]]_{DS}^x &= \{\theta \in [[P]]_{DS}^x \mid R\theta = \top\} \end{aligned}$$

Let  $R$  be a FILTER expression,  $u, v \in Var$ ,  $c \in I \cup B \cup L$ . The valuation of  $R$  on substitution  $\theta$ , written  $R\theta$  takes one of the three values  $\{\top, \perp, \varepsilon\}$ <sup>6</sup> and is defined as follows.  $R\theta = \top$ , if:

- (1)  $R = \text{BOUND}(v)$  with  $v \in \text{dom}(\theta) \wedge v\theta \neq \text{null}$ ;
- (2)  $R = \text{isBLANK}(v)$  with  $v \in \text{dom}(\theta) \wedge v\theta \in B$ ;
- (3)  $R = \text{isIRI}(v)$  with  $v \in \text{dom}(\theta) \wedge v\theta \in I$ ;
- (4)  $R = \text{isLITERAL}(v)$  with  $v \in \text{dom}(\theta) \wedge v\theta \in L$ ;
- (5)  $R = (v = c)$  with  $v \in \text{dom}(\theta) \wedge v\theta = c$ ;
- (6)  $R = (u = v)$  with  $u, v \in \text{dom}(\theta) \wedge u\theta = v\theta \wedge u\theta \neq \text{null}$ ;
- (7)  $R = (\neg R_1)$  with  $R_1\theta = \perp$ ;
- (8)  $R = (R_1 \vee R_2)$  with  $R_1\theta = \top \vee R_2\theta = \top$ ;
- (9)  $R = (R_1 \wedge R_2)$  with  $R_1\theta = \top \wedge R_2\theta = \top$ .

$R\theta = \varepsilon$ , if:

- (1)  $R = \text{isBLANK}(v), R = \text{isIRI}(v), R = \text{isLITERAL}(v)$ , or  $R = (v = c)$  with  $v \notin \text{dom}(\theta) \vee v\theta = \text{null}$ ;
- (2)  $R = (u = v)$  with  $u \notin \text{dom}(\theta) \vee u\theta = \text{null} \vee v \notin \text{dom}(\theta) \vee v\theta = \text{null}$ ;
- (3)  $R = (\neg R_1)$  and  $R_1\theta = \varepsilon$ ;
- (4)  $R = (R_1 \vee R_2)$  and  $(R_1\theta \neq \top \wedge R_2\theta \neq \top) \wedge (R_1\theta = \varepsilon \vee R_2\theta = \varepsilon)$ ;
- (5)  $R = (R_1 \wedge R_2)$  and  $R_1\theta = \varepsilon \vee R_2\theta = \varepsilon$ .

$R\theta = \perp$  otherwise.

We will now exemplify the three different semantics defined above, namely bravely joining (b-joining), cautiously joining (c-joining), and strictly-joining (s-joining) semantics. When taking a closer look to the AND and MINUS operators, one will realize that all three semantics take a slightly differing view only when joining null. Indeed, the AND operator behaves as the traditional natural join operator  $\bowtie$  in relational algebra, when no null values are involved.

Take for instance,  $DS = (\{\text{ex.org/bob, alice.org}\}, \emptyset)$  and  $P = ((?X, \text{name}, ?Name) \text{ AND } (?X, \text{knows}, ?Friend))$ . When viewing each solution set as a relational table with variables denoting attribute names, we can write:

?X	?Name		?X	?Friend
alice.org#me	"Bob"	$\bowtie$	alice.org#me	alice.org#me
alice.org#me	"Alice"		alice.org#me	alice.org#me
alice.org#me	"Bob"		alice.org#me	alice.org#me

  

?X	?Name	?Friend
alice.org#me	"Bob"	alice.org#me
alice.org#me	"Alice"	alice.org#me

Differences between the three semantics appear when joining over null-bound variables, as shown in the next example.

<sup>6</sup>  $\top$  stands for "true",  $\perp$  stands for "false" and  $\varepsilon$  stands for errors, see [18, Sec. 11.3] and Example 2 for details.

EXAMPLE 3. Let  $DS$  be as before and assume the following query which might be considered a naive attempt to ask for pairs of persons  $?X1$ ,  $?X2$  who share the same name and nickname where both, name and nickname are optional:

$P = ((?X1, a, \text{Person}) \text{ OPT } (?X1, \text{name}, ?N)) \text{ AND } ((?X2, a, \text{Person}) \text{ OPT } (?X2, \text{nick}, ?N))$

Again, we consider the tabular view of the resulting join:

$?X1$	$?N$	$?X2$	$?N$
$\_ : a$	"Bob"	$\_ : a$	null
$\_ : b$	null	$\_ : b$	"Alice"
$\_ : c$	"Bob"	$\_ : c$	"Bobby"
alice.org#me	"Alice"	alice.org#me	null

Now, let us see what happens when we evaluate the join  $\bowtie_x$  with respect to the different semantics. The following result table lists in the last column which tuples belong to the result of b-, c- and s-join, respectively.

$?X1$	$?N$	$X2$	
$\_ : a$	"Bob"	$\_ : a$	b
$\_ : a$	"Bob"	alice.org#me	b
$\_ : b$	null	$\_ : a$	b, c
$\_ : b$	"Alice"	$\_ : b$	b
$\_ : b$	"Bobby"	$\_ : c$	b
$\_ : b$	null	alice.org#me	b, c
$\_ : c$	"Bob"	$\_ : a$	b
$\_ : c$	"Bob"	alice.org#me	b
alice.org#me	"Alice"	$\_ : a$	b
alice.org#me	"Alice"	$\_ : b$	b, c, s
alice.org#me	"Alice"	alice.org#me	b

Leaving aside the question whether the query formulation was intuitively broken, we remark that only the s-join would have the expected result. At the very least we might argue, that the liberal behavior of b-joins might be considered surprising in some cases. The c-joining semantics acts a bit more cautious in between the two, treating null values as normal values, only unifiable with other null values.

Compared to how joins over incomplete relations are treated in common relational database systems, the s-joining semantics might be considered the intuitive behavior. Another interesting divergence (which would rather suggest to adopt the c-joining semantics) shows up when we consider a simple idempotent join.

EXAMPLE 4. Let us consider the following single triple dataset  $DS = (\{(alice.org\#me, a, \text{Person})\}, \emptyset)$  and the following simple query pattern:

$P = ((?X, a, \text{Person}) \text{ UNION } (?Y, a, \text{Person}))$

Clearly, this pattern, has the solution set

$[[P]]_{DS}^x = \{(alice.org\#me, null), (null, alice.org\#me)\}$

under all three semantics. Surprisingly,  $P' = (P \text{ AND } P)$  has different solution sets for the different semantics. First,  $[[P']]_{DS}^c = [[P]]_{DS}^c$ , but  $[[P']]_{DS}^s = \emptyset$ , since null values are not compatible under the s-joining semantics. Finally,

$[[P']]_{DS}^b = \{(alice.org\#me, null), (null, alice.org\#me), (alice.org\#me, alice.org\#me)\}$

As shown by this example, under the reasonable assumption, that the join operator is idempotent, i.e.,  $(P \bowtie P) \equiv P$ , only the c-joining semantics behaves correctly.

However, the brave b-joining behavior is advocated by the current SPARQL document, and we might also think of examples where this obviously makes a lot of sense. Especially, when considering no explicit joins, but the implicit joins within the OPT operator:

EXAMPLE 5. Let  $DS = (\{ex.org/bob, alice.org\}, \emptyset)$  and assume a slight variant of a query from [5] which asks for persons and some names for these persons, where preferably the foaf : name is taken, and, if not specified, foaf : nick.

$P = (((?X, a, \text{Person}) \text{ OPT } (?X, \text{name}, ?XNAME)) \text{ OPT } (?X, \text{nick}, ?XNAME))$

Only  $[[P]]_{DS}^b$  contains the expected solution  $(\_ : b, "Alice")$  for the bNode  $\_ : b$ .

All three semantics may be considered as variations of the original definitions in [16], for which the authors proved complexity results and various desirable features, such as semantics-preserving normal form transformations and compositionality. The following proposition shows that all these results carry over to the normative b-joining semantics:

PROPOSITION 1. Given a dataset  $DS$  and a pattern  $P$  which does not contain GRAPH patterns, the solutions of  $[[P]]_{DS}$  as in [16] and  $[[P]]_{DS}^b$  are in 1-to-1 correspondence.

PROOF. Given  $DS$  and  $P$  each substitution  $\theta$  obtained by evaluation  $[[P]]_{DS}^b$  can be reduced to a substitution  $\theta'$  obtained from the evaluation  $[[P]]_{DS}$  in [16] by dropping all mappings of the form  $v \rightarrow \text{null}$  from  $\theta$ . Likewise, each substitution  $\theta'$  obtained from  $[[P]]_{DS}$  can be extended to a substitution  $\theta = \theta'^{\text{vars}(P)}$  for  $[[P]]_{DS}^b$ .  $\square$

Following the definitions from the SPARQL specification and [16], the b-joining semantics is the only admissible definition. There are still advantages for gradually defining alternatives towards traditional treatment of joins involving nulls. On the one hand, as we have seen in the examples above, the brave view on joining unbound variables might have partly surprising results, on the other hand, as we will see, the c- and s-joining semantics allow for a more efficient implementation in terms of Datalog rules.

Let us now take a closer look on some properties of the three defined semantics.

**Compositionality and Equivalences.** As shown in [16], some implementations have a non-compositional semantics, leading to undesired effects such as non-commutativity of the join operator, etc. A semantics is called *compositional* if for each  $P'$  sub-pattern of  $P$  the result of evaluating  $P'$  can be used to evaluate  $P$ . Obviously, all three the c-, s- and b-joining semantics defined here retain this property, since all three semantics are defined recursively, and independent of the evaluation order of the sub-patterns.

The following proposition summarizes equivalences which hold for all three semantics, showing some interesting additions to the results of Pérez et al.

PROPOSITION 2 (EXTENDS [16, PROP. 1]). The following equivalences hold or do not hold in the different semantics as indicated after each law:

- (1)  $\text{AND}$ ,  $\text{UNION}$  are associative and commutative. (b, c, s)
- (2)  $(P_1 \text{ AND } (P_2 \text{ UNION } P_3)) \equiv ((P_1 \text{ AND } P_2) \text{ UNION } (P_1 \text{ AND } P_3)).$  (b)
- (3)  $(P_1 \text{ OPT } (P_2 \text{ UNION } P_3)) \equiv ((P_1 \text{ OPT } P_2) \text{ UNION } (P_1 \text{ OPT } P_3)).$  (b)
- (4)  $((P_1 \text{ UNION } P_2) \text{ OPT } P_3) \equiv ((P_1 \text{ OPT } P_3) \text{ UNION } (P_2 \text{ OPT } P_3)).$  (b)
- (5)  $((P_1 \text{ UNION } P_2) \text{ FILTER } R) \equiv ((P_1 \text{ FILTER } R) \text{ UNION } (P_2 \text{ FILTER } R)).$  (b, c, s)
- (6)  $\text{AND}$  is idempotent, i.e.  $(P \text{ AND } P) \equiv P.$  (c)

PROOF SKETCH.. (1-5) for the b-joining semantics are proven in [16], (1): for c-joining and s-joining follows straight from the definitions. (2)-(4): the substitution sets  $[[P_1]]^{c,s} = \{[?X \rightarrow a, ?Y \rightarrow b]\}$ ,  $[[P_2]]^{c,s} = \{[?X \rightarrow a, ?Z \rightarrow c]\}$ ,  $[[P_3]]^{c,s} = \{[?Y \rightarrow b, ?Z \rightarrow c]\}$  provide counterexamples for c-joining and s-joining semantics for all three equivalences (2)-(4). (5): The semantics of FILTER expressions and UNION is exactly the same for all three semantics, thus, the result for the b-joining semantics carries over to all three semantics. (6): follows from the observations in Example 4.  $\square$

Ideally, we would like to identify a subclass of programs, where the three semantics coincide. Obviously, this is the case for any query involving neither UNION nor OPT operators. Pérez et al. [16] define a bigger class of programs, including “well-behaving” optional patterns:

**Definition 3.** ([16, Def. 4]) A UNION-free graph pattern  $P$  is *well-designed* if for every occurrence of a sub-pattern  $P' = (P_1 \text{ OPT } P_2)$  of  $P$  and for every variable  $v$  occurring in  $P$ , the following condition holds: if  $v$  occurs both in  $P_2$  and outside  $P'$  then it also occurs in  $P_1$ .

As may be easily verified by the reader, neither Example 3 nor Example 5, which are both UNION-free, satisfy the well-designedness condition. Since in the general case the equivalences for Prop. 2 do not hold, we also need to consider nested UNION patterns as a potential source for null bindings which might affect join results. We extend the notion of well-designedness, which directly leads us to another correspondence in the subsequent proposition.

**Definition 4.** A graph pattern  $P$  is *well-designed* if the condition from Def. 3 holds and for every occurrence of a sub-pattern  $P' = (P_1 \text{ UNION } P_2)$  of  $P$  and for every variable  $v$  occurring in  $P'$ , the following condition holds: if  $v$  occurs outside  $P'$  then it occurs in both  $P_1$  and  $P_2$ .

**PROPOSITION 3.** *On well-designed graph patterns the c-, s-, and b-joining semantics coincide.*

PROOF SKETCH.. Follows directly from the observation that all variables which are re-used outside  $P'$  must be bound to a value unequal to null in  $P'$  due to well-designedness, and thus cannot generate null bindings which might carry over to joins.  $\square$

Likewise, we can identify “dangerous” variables in graph patterns, which might cause semantic differences:

**Definition 5.** Let  $P'$  a sub-pattern of  $P$  of either the form  $P' = (P_1 \text{ OPT } P_2)$  or  $P' = (P_1 \text{ UNION } P_2)$ . Any variable  $v$  in  $P'$  which violates the well-designedness-condition is called *possibly-null-binding* in  $P$ .

Note that, so far we have only defined the semantics in terms of a pattern  $P$  and dataset  $DS$ , but not yet taken the result form  $V$  of query  $Q = (V, P, DS)$  into account.

We now define *solution tuples* that were informally introduced in Sec. 2. Recall that by  $\bar{V}$  we denote the tuple obtained from lexicographically ordering a set of variables in  $V$ . The notion  $\bar{V}[V' \rightarrow \text{null}]$  means that, after ordering  $V$  all variables from a subset  $V' \subseteq V$  are replaced by null.

**Definition 6.** (*Solution Tuples*) Let  $Q = (V, P, DS)$  be a SPARQL query, and  $\theta$  a substitution in  $[[P]]_{DS}^x$ , then we call the tuple  $\bar{V}[(V \setminus \text{vars}(P)) \rightarrow \text{null}]\theta$  a solution tuple of  $Q$  with respect to the  $x$ -joining semantics.

Let us remark at this point, that as for the discussion of intuitivity of the different join semantics discussed in Examples 3-5, we did not yet consider combinations of different join semantics, e.g. using b-joins for OPT and c-joins for AND patterns. We leave this for further work.

### 3. DATALOG AND ANSWER SETS

In this paper we will use a very general form of Datalog commonly referred to as Answer Set Programming (ASP), i.e. function-free logic programming (LP) under the answer set semantics [1, 11]. ASP is widely proposed as a useful tool for various problem solving tasks in e.g. Knowledge Representation and Deductive databases. ASP extends Datalog with useful features such as negation as failure, disjunction in rule heads, aggregates [9], external predicates [8], etc. <sup>7</sup>

Let  $Pred$ ,  $Const$ ,  $Var$ ,  $exPr$  be sets of predicate, constant, variable symbols, and external predicate names, respectively. Note that we assume all these sets except  $Pred$  and  $Const$  (which may overlap), to be disjoint. In accordance with common notation in LP and the notation for external predicates from [7] we will in the following assume that  $Const$  and  $Pred$  comprise sets of numeric constants, string constants beginning with a lower case letter, or “” quoted strings, and strings of the form  $\langle \text{quoted-string} \rangle^{\sim} \langle \text{IRI} \rangle$ ,  $\langle \text{quoted-string} \rangle @ \langle \text{valid-lang-tag} \rangle$ ,  $Var$  is the set of string constants beginning with an upper case letter. Given  $p \in Pred$  an *atom* is defined as  $p(t_1, \dots, t_n)$ , where  $n$  is called the arity of  $p$  and  $t_1, \dots, t_n \in Const \cup Var$ .

Moreover, we define a fixed set of external predicates  $exPr = \{rdf, isBLANK, isIRI, isLITERAL, =, !=\}$ . All external predicates have a fixed semantics and fixed arities, distinguishing *input* and *output terms*. The atoms  $isBLANK[c](val)$ ,  $isIRI[c](val)$ ,  $isLITERAL[c](val)$  test the input term  $c \in Const \cup Var$  (in square brackets) for being valid string representations of Blank nodes, IRI References or RDF literals, returning an output value  $val \in \{t, f, e\}$ , representing truth, falsity or an error, following the semantics defined in [18, Sec. 11.3]. For the *rdf* predicate we write atoms as  $rdf[i](s, p, o)$  to denote that  $i \in Const \cup Var$  is an input term, whereas  $s, p, o \in Const \cup Var$  are output terms which may be bound by the external predicate. The external atom  $rdf[i](s, p, o)$  is true if  $(s, p, o)$  is an RDF triple entailed by the RDF graph which is accessibly at IRI  $i$ . For the moment, we consider simple RDF entailment [13] only. Finally, we write comparison atoms  $t_1 = t_2$  and  $t_1 != t_2$  in infix notation with  $t_1, t_2 \in Const \cup Var$  and the obvious semantics of (lexicographic or numeric) (in)equality. Here, for  $=$  either  $t_1$  or  $t_2$  is an output term, but at least one is an input term, and for  $!=$  both  $t_1$  and  $t_2$  are input terms.

**Definition 7.** Finally, a *rule* is of the form

$$h :- b_1, \dots, b_m, \text{not } b_{m+1}, \dots, \text{not } b_n. \quad (1)$$

where  $h$  and  $b_i$  ( $1 \leq i \leq n$ ) are atoms,  $b_k$  ( $1 \leq k \leq m$ ) are either atoms or external atoms, and **not** is the symbol for negation as failure.

We use  $H(r)$  to denote the head atom  $h$  and  $B(r)$  to denote the set of all body literals  $B^+(r) \cup B^-(r)$  of  $r$ , where  $B^+(r) = \{b_1, \dots, b_m\}$  and  $B^-(r) = \{b_{m+1}, \dots, b_n\}$ .

<sup>7</sup>We consider ASP, more precisely a simplified version of ASP with so-called HEX-programs [8] here, since it is up to date the most general extension of Datalog.

The notion of input and output terms in external atoms described above denotes the binding pattern. More precisely, we assume the following condition which extends the standard notion of safety (cf. [21]) in Datalog with negation: Each variable appearing in a rule must appear in  $B^+(r)$  in an atom or as an output term of an external atom.

*Definition 8.* A (logic) program  $\Pi$  is defined as a set of safe rules  $r$  of the form (1).

The *Herbrand base* of a program  $\Pi$ , denoted  $HB_\Pi$ , is the set of all possible ground versions of atoms and external atoms occurring in  $\Pi$  obtained by replacing variables with constants from  $Const$ , where we define for our purposes by  $Const$  the union of the set of all constants appearing in  $\Pi$  as well as the literals, IRIs, and distinct constants for each blank node occurring in each RDF graph identified<sup>8</sup> by one of the IRIs in the (recursively defined) set  $I$ , where  $I$  is defined by the recursive closure of all IRIs appearing in  $\Pi$  and all RDF graphs identified by IRIs in  $I$ .<sup>9</sup> As long as we assume that the Web is finite the grounding of a rule  $r$ ,  $ground(r)$ , is defined by replacing each variable with the possible elements of  $HB_\Pi$ , and the grounding of program  $\Pi$  is  $ground(\Pi) = \bigcup_{r \in \Pi} ground(r)$ .

An *interpretation relative to  $\Pi$*  is any subset  $\mathcal{I} \subseteq HB_\Pi$  containing only atoms. We say that  $\mathcal{I}$  is a *model* of atom  $a \in HB_\Pi$ , denoted  $\mathcal{I} \models a$ , if  $a \in \mathcal{I}$ . With every external predicate name  $lg \in exPr$  with arity  $n$  we associate an  $(n+1)$ -ary Boolean function  $f_{lg}$  (called *oracle function*) assigning each tuple  $(\mathcal{I}, t_1, \dots, t_n)$  either 0 or 1.<sup>10</sup> We say that  $\mathcal{I} \subseteq HB_\Pi$  is a *model* of a ground external atom  $a = g[t_1, \dots, t_m](t_{m+1}, \dots, t_n)$ , denoted  $\mathcal{I} \models a$ , if  $f_{lg}(\mathcal{I}, t_1, \dots, t_n) = 1$ .

The semantics we use here generalizes the answer-set semantics [11]<sup>11</sup>, and is defined using the *FLP-reduct* [9], which is more elegant than the traditional *GL-reduct* [11] of stable model semantics and ensures minimality of answer sets also in presence of external atoms.

Let  $r$  be a ground rule. We define (i)  $\mathcal{I} \models B(r)$  iff  $\mathcal{I} \models a$  for all  $a \in B^+(r)$  and  $\mathcal{I} \not\models a$  for all  $a \in B^-(r)$ , and (ii)  $\mathcal{I} \models r$  iff  $\mathcal{I} \models H(r)$  whenever  $\mathcal{I} \models B(r)$ . We say that  $\mathcal{I}$  is a *model* of a program  $\Pi$ , denoted  $\mathcal{I} \models \Pi$ , iff  $\mathcal{I} \models r$  for all  $r \in ground(\Pi)$ .

The *FLP-reduct* [9] of  $\Pi$  with respect to  $\mathcal{I} \subseteq HB_\Pi$ , denoted  $\Pi^\mathcal{I}$ , is the set of all  $r \in ground(\Pi)$  such that  $\mathcal{I} \models B(r)$ .  $\mathcal{I}$  is an *answer set* of  $\Pi$  iff  $\mathcal{I}$  is a minimal model of  $\Pi^\mathcal{I}$ .

We did not consider further extensions common to many ASP dialects here, namely disjunctive rule heads, strong negation [11]. We note that for non-recursive programs, i.e. where the predicate dependency graph is acyclic, the answer set is unique. For the pure translation which we will give in Sec. 4 where we will produce such non-recursive programs from SPARQL queries, we could equally take other seman-

<sup>8</sup>By “identified” we mean here that IRIs denote network accessible resources which correspond to RDF graphs.

<sup>9</sup>We assume the number of accessible IRIs finite.

<sup>10</sup>The notion of an oracle function reflects the intuition that external predicates compute (sets of) outputs for a particular input, depending on the interpretation. The dependence on the interpretation is necessary for instance for defining the semantics of external predicates querying OWL [8] or computing aggregate functions.

<sup>11</sup>In fact, we use slightly simplified definitions from [7] for HEX-programs, with the sole difference that we restrict ourselves to a fixed set of external predicates.

tics such as the well-founded [10] semantics into account, which coincides with ASP on non-recursive programs.

## 4. FROM SPARQL TO DATALOG

We are now ready to define a translation from SPARQL to Datalog which can serve straightforwardly to implement SPARQL within existing rules engines. We start with a translation for c-joining semantics, which we will extend thereafter towards s-joining and b-joining semantics.

*Translation  $\Pi_Q^c$ .* Let  $Q = (V, P, DS)$ , where  $DS = (G, G_n)$  as defined above. We translate this query to a logic program  $\Pi_Q^c$  defined as follows.

$$\begin{aligned} \Pi_Q^c = & \{ \text{triple}(S, P, O, \text{default}) \text{ :- rdf}[d](S, P, O). \mid d \in G \} \\ & \cup \{ \text{triple}(S, P, O, g) \text{ :- rdf}[g](S, P, O). \mid g \in G_n \} \\ & \cup \tau(V, P, \text{default}, 1) \end{aligned}$$

The first two rules serve to import the relevant RDF triples from the dataset into a 4-ary predicate **triple**. Under the dataset closedness assumption (see Def. 1) we may replace the second rule set, which imports the named graphs, by:

$$\text{triple}(S, P, O, G) \text{ :- rdf}[G](S, P, O), HU(G), \text{isIRI}(G).$$

Here, the predicate  $HU$  stands for “Herbrand universe”, where we use this name a bit sloppily, with the intention to cover all the relevant part of  $\mathcal{C}$ , recursively importing all possible IRIs in order to emulate the dataset closedness assumption.  $HU$ , can be computed recursively over the input triples, i.e.

$$\begin{aligned} HU(X) & \text{ :- triple}(X, P, O, D). \quad HU(X) \text{ :- triple}(S, X, O, D). \\ HU(X) & \text{ :- triple}(S, P, X, D). \quad HU(X) \text{ :- triple}(S, P, O, X). \end{aligned}$$

The remaining program  $\tau(V, P, \text{default}, 1)$  represents the actual query translation, where  $\tau$  is defined recursively as shown in Fig. 2.

By  $LT(\cdot)$  we mean the set of rules resulting from disassembling complex **FILTER** expressions (involving  $\neg, \wedge, \vee$ ) according to the rewriting defined by Lloyd and Topor [15] where we have to obey the semantics for errors, following Definition 2. In a nutshell, the rewriting  $LT - \text{rewrite}(\cdot)$  proceeds as follows: Complex filters involving  $\neg$  are transformed into negation normal form. Conjunctions of filter expressions are simply disassembled to conjunctions of body literals, disjunctions are handled by splitting the respective rule for both alternatives in the standard way. The resulting rules involve possibly negated atomic filter expressions in the bodies. Here,  $BOUND(v)$  is translated to  $v = \text{null}$ ,  $\neg BOUND(v)$  to  $v \neq \text{null}$ .  $\text{isBLANK}(v)$ ,  $\text{isIRI}(v)$ ,  $\text{isLITERAL}(v)$  and their negated forms are replaced by their corresponding external atoms (see Sec. 3)  $\text{isBLANK}[v](t)$  or  $\text{isBLANK}[v](f)$ , etc., respectively.

The resulting program  $\Pi_Q^c$  implements the c-joining semantics in the following sense:

**PROPOSITION 4** (SOUNDNESS AND COMPLETENESS OF  $\Pi_Q^c$ ).  
For each atom of the form  $\text{answer}_1(\vec{s}, \text{default})$  in the unique answer set  $M$  of  $\Pi_Q^c$ ,  $\vec{s}$  is a solution tuple of  $Q$  with respect to the c-joining semantics, and all solution tuples of  $Q$  are represented by the extension of predicate  $\text{answer}_1$  in  $M$ .

Without giving a proof, we remark that the result follows if we convince ourselves that  $\tau(V, P, D, i)$  emulates exactly

$$\tau(V, (s, p, o), D, i) = \text{answer}_1(\bar{V}, D) \text{ :- triple}(s, p, o, D). \quad (1)$$

$$\tau(V, (P' \text{ AND } P''), D, i) = \tau(\text{vars}(P'), P', D, 2*i) \cup \tau(\text{vars}(P''), P'', D, 2*i+1) \cup \text{answer}_1(\bar{V}, D) \text{ :- answer}_{2*i}(\text{vars}(P'), D), \text{answer}_{2*i+1}(\text{vars}(P''), D). \quad (2)$$

$$\tau(V, (P' \text{ UNION } P''), D, i) = \tau(\text{vars}(P'), P', D, 2*i) \cup \tau(\text{vars}(P''), P'', D, 2*i+1) \cup \text{answer}_1(\overline{V[(V \setminus \text{vars}(P')) \rightarrow \text{null}]}, D) \text{ :- answer}_{2*i}(\text{vars}(P'), D). \quad (3)$$

$$\text{answer}_1(\overline{V[(V \setminus \text{vars}(P'')) \rightarrow \text{null}]}, D) \text{ :- answer}_{2*i+1}(\text{vars}(P''), D). \quad (4)$$

$$\tau(V, (P' \text{ MINUS } P''), D, i) = \tau(\text{vars}(P'), P', D, 2*i) \cup \tau(\text{vars}(P''), P'', D, 2*i+1) \cup \text{answer}_1(\overline{V[(V \setminus \text{vars}(P')) \rightarrow \text{null}]}, D) \text{ :- answer}_{2*i}(\text{vars}(P'), D), \quad (5)$$

$$\text{not answer}_{2*i}'(\text{vars}(P') \cap \text{vars}(P''), D), \text{answer}_{2*i+1}'(\text{vars}(P'') \cap \text{vars}(P'), D). \quad (6)$$

$$\tau(V, (P' \text{ OPT } P''), D, i) = \tau(V, (P' \text{ AND } P''), D, i) \cup \tau(V, (P' \text{ MINUS } P''), D, i)$$

$$\tau(V, (P \text{ FILTER } R), D, i) = \tau(\text{vars}(P), P, D, 2*i) \cup \text{LT}(\text{answer}_1(\bar{V}, D) \text{ :- answer}_{2*i}(\text{vars}(P), D), R). \quad (7)$$

$$\tau(V, (\text{GRAPH } g \text{ } P), D, i) = \tau(V, P, g, i) \text{ for } g \in V \cup I$$

$$\text{answer}_1(\bar{V}, D) \text{ :- answer}_1(\bar{V}, g), \text{isIRI}(g), \text{not } g = \text{default}. \quad (8)$$

Alternate rules replacing (5)+(6):

$$\text{answer}_1(\overline{V[(V \setminus \text{vars}(P')) \rightarrow \text{null}]}, D) \text{ :- answer}_{2*i}(\text{vars}(P'), D), \text{not answer}_{2*i}'(\text{vars}(P'), D) \quad (5')$$

$$\text{answer}_{2*i}'(\text{vars}(P'), D) \text{ :- answer}_{2*i}(\text{vars}(P'), D), \text{answer}_{2*i+1}(\text{vars}(P''), D). \quad (6')$$

Figure 2: Translation  $\Pi_Q^c$  from SPARQL queries semantics to Datalog.

the recursive definition of  $[[P]]_{DS}^s$ . Moreover, together with Proposition 3, we obtain soundness and completeness of  $\Pi_Q$  for b-joining and s-joining semantics as well for well-designed query patterns.

**COROLLARY 1.** For  $Q = (V, P, DS)$ , if  $P$  is well-designed, then the extension of predicate  $\text{answer}_1$  in the unique answer set  $M$  of  $\Pi_Q^c$  represents all and only the solution tuples for  $Q$  with respect to the  $x$ -joining semantics, for  $x \in \{b, c, s\}$ .

Now, in order to obtain a proper translation for arbitrary patterns, we obviously need to focus our attention on the possibly-null-binding variables within the query pattern  $P$ . Let  $vnull(P)$  denote the possibly-null-binding variables in a (sub)pattern  $P$ . We need to consider all rules in Fig. 2 which involve  $x$ -joins, i.e. the rules of the forms (2), (5) and (6). Since rules (5) and (6) do not make this join explicit, we will replace them by the equivalent rules (5') and (6') for  $\Pi_Q^s$  and  $\Pi_Q^b$ . The “extensions” to s-joining and b-joining semantics can be achieved by rewriting the rules (2) and (6'). The idea is to rename variables and add proper FILTER expressions to these rules in order to realize the b-joining and s-joining behavior for the variables in  $V_N = vnull(P) \cap \text{vars}(P') \cap \text{vars}(P'')$ .

**Translation  $\Pi_Q^s$ .** The s-joining behavior can be achieved by adding FILTER expressions

$$R^s = ( \bigwedge_{v \in V_N} \text{BOUND}(v) )$$

to the rule bodies of (2) and (6'). The resulting rules are again subject to the *LT*-rewriting as discussed above for the rules of the form (7). This is sufficient to filter out any joins involving null values, thus achieving s-joining semantics, and we denote the program rewritten that way as  $\Pi_Q^s$ .

**Translation  $\Pi_Q^b$ .** Obviously, b-joining semantics is more tricky to achieve, since we now have to relax the allowed joins in order to allow null bindings to join with *any* other value. We will again achieve this result by modifying rules (2) and (6') where we first do some variable renaming and then add respective FILTER expressions to these rules.

**Step 1.** We rename each variable  $v \in V_N$  in the respective rule bodies to  $v'$  or  $v''$ , respectively, in order to disambiguate the occurrences originally from sub-pattern  $P'$  or  $P''$ , respectively. That is, for each rule (2) or (6'), we rewrite the body to:

$$\text{answer}_{2*i}(\text{vars}(P')[V_N \rightarrow V'_N], D),$$

$$\text{answer}_{2*i+1}(\text{vars}(P'')[V_N \rightarrow V''_N], D).$$

**Step 2.** We now add the following FILTER expressions  $R_{(2)}^b$  and  $R_{(6')}^b$ , respectively, to the resulting rule bodies which “emulate” the relaxed b-compatibility:

$$R_{(2)}^b = \bigwedge_{v \in V_N} ( ((v = v') \wedge (v' = v'')) \vee ((v = v') \wedge \neg \text{BOUND}(v'')) \vee ((v = v'') \wedge \neg \text{BOUND}(v')) )$$

$$R_{(6')}^b = \bigwedge_{v \in V_N} ( ((v = v') \wedge (v' = v'')) \vee ((v = v') \wedge \neg \text{BOUND}(v'')) \vee ((v = v'') \wedge \neg \text{BOUND}(v')) )$$

The rewritten rules are again subject to the *LT* rewriting. Note that, strictly speaking the filter expression introduced here does not fulfill the assumption of safe filter expressions, since it creates new bindings for the variable  $v$ . However, these can safely be allowed here, since the translation only creates valid input/output term bindings for the external Datalog predicate '='. The subtle difference between  $R_{(2)}^b$  and  $R_{(6')}^b$  lies in the fact that  $R_{(2)}^b$  preferably “carries over” bound values from  $v'$  or  $v''$  to  $v$  whereas  $R_{(6')}^b$  always takes the value of  $v'$ . The effect of this becomes obvious in the translation of Example 5 which we leave as an exercise to



the reader. We note that the potential exponential (with respect to  $|V_N|$ ) blowup of the program size by unfolding the filter expressions into negation normal form during the *LT* rewriting<sup>12</sup> is not surprising, given the negative complexity results in [16].

In total, we obtain a program which  $\Pi_Q^b$  which reflects the normative b-joining semantics. Consequently, we get sound and complete query translations for all three semantics:

**COROLLARY 2** (SOUNDNESS AND COMPLETENESS OF  $\Pi_Q^x$ ). *Given an arbitrary graph pattern  $P$ , the extension of predicate  $\text{answer}_1$  in the unique answer set  $M$  of  $\Pi_Q^x$  represents all and only the solution tuples for  $Q = (V, P, DS)$  with respect to the  $x$ -joining semantics, for  $x \in \{b, c, s\}$ .*

In the following, we will drop the superscript  $x$  in  $\Pi_Q$  implicitly refer to the normative b-joining translation/semantics.

## 5. POSSIBLE EXTENSIONS

As it turns out, the embedding of SPARQL in the rules world opens a wide range of possibilities for combinations. In this section, we will first discuss some straightforward extensions of SPARQL which come practically for free with the translation to Datalog provided before. We will then discuss the use of SPARQL itself as a simple RDF rules language<sup>13</sup> which allows to combine RDF fact bases with implicitly specified further facts and discuss the semantics thereof briefly. We conclude this section with revisiting the open issue of entailment regimes covering RDFS or OWL semantics in SPARQL.

### 5.1 Additional Language Features

**Set Difference.** As mentioned before, set difference is not present in the current SPARQL specification syntactically, though hidden, and would need to be emulated via a combination of **OPTIONAL** and **FILTER** constructs. As we defined the **MINUS** operator here in a completely modular fashion, it could be added straightforwardly without affecting the semantics definition.

**Nested queries.** Nested queries are a distinct feature of SQL not present in SPARQL. We suggest a simple, but useful form of nested queries to be added: Boolean queries  $Q_{\text{ASK}} = (\emptyset, P_{\text{ASK}}, DS_{\text{ASK}})$  with an empty result form (denoted by the keyword **ASK**) can be safely allowed within **FILTER** expressions as an easy extension fully compatible with our translation. Given query  $Q = (V, P, DS)$ , with sub-pattern  $(P_1 \text{ FILTER } (ASK Q_{\text{ASK}}))$  we can modularly translate such subqueries by extending  $\Pi_Q$  with  $\Pi_{Q'}$  where  $Q' = (\text{vars}(P_1) \cap \text{vars}(P_{\text{ASK}}), P_{\text{ASK}}, DS_{\text{ASK}})$ . Moreover, we have to rename predicate names  $\text{answer}_i$  to  $\text{answer}_i^{Q'}$  in  $\Pi_{Q'}$ . Some additional considerations are necessary in order to combine this within arbitrary complex filter expressions, and we probably need to impose well-designedness for variables shared between  $P$  and  $P_{\text{ASK}}$  similar to Def. 4. We leave more details as future work.

<sup>12</sup>Lloyd and Topor can avoid this potential exponential blowup by introducing new auxiliary predicates. However, we cannot do the same trick, mainly for reasons of preserving safety of external predicates as defined in Sec. 3.

<sup>13</sup>Thus, the "... (and back)" in the title of this paper!

## 5.2 Result Forms and Solution Modifiers

We have covered only **SELECT** queries so far. As shown in the previous section, we can consider **ASK** queries equally. A limited form of the **CONSTRUCT** result form, which allows to construct new triples could be emulated in our approach as well. Namely, we can allow queries of the form

$$Q_C = (\text{CONSTRUCT } P_C, P, DS)$$

where  $P_C$  is a graph pattern consisting only of bNode-free triple patterns. We can model these by adding a rule

$$\text{triple}(s, p, o, C) \text{ :- } \text{answer}_1(\overline{\text{vars}(P_C)}, \text{default}). \quad (2)$$

to  $\Pi_Q$  for each triple  $(s, p, o)$  in  $P_C$ . The result graph is then naturally represented in the answer set of the program extended that way in the extension of the predicate **triple**.

### 5.3 SPARQL as a Rules Language

As it turns out with the extensions defined in the previous subsections, SPARQL itself may be viewed as an expressive rules language on top of RDF. **CONSTRUCT** statements have an obvious similarity with view definitions in SQL, and thus may be seen as rules themselves.

Intuitively, in the translation of **CONSTRUCT** we "stored" the new triples in a new triple outside the dataset  $DS$ . We can imagine a similar construction in order to define the semantics of queries over datasets mixing such **CONSTRUCT** statements with RDF data in the same turtle file.

Let us assume such a mixed file containing **CONSTRUCT** rules and RDF triples web-accessible at IRI  $g$ , and a query  $Q = (V, P, DS)$ , with  $DS = (G, G_n)$ . The semantics of a query over a dataset containing  $g$  may then be defined by recursively adding  $\Pi_{Q_C}$  to  $\Pi_Q$  for any **CONSTRUCT** query  $Q_C$  in  $g$  plus the rules (2) above with their head changed to  $\text{triple}(s, p, o, g)$ . We further need to add a rule

$$\text{triple}(s, p, o, \text{default}) \text{ :- } \text{triple}(s, p, o, g).$$

for each  $g \in G$ , in order not to omit any of the implicit triples defined by such "CONSTRUCT rules". Analogously to the considerations for nested **ASK** queries, we need to rename the  $\text{answer}_i$  predicates and  $\text{default}$  constants in every subprogram  $\Pi_{Q_C}$  defined this way.

Naturally, the resulting programs possibly involve recursion, and, even worse, recursion over negation as failure. Fortunately, the general answer set semantics, which we use, can cope with this. For some important aspects on the semantics of such distributed rules and facts bases, we refer to [17], where we also outline an alternative semantics based on the well-founded semantics. A more in-depth investigation of the complexity and other semantic features of such a combination is on our agenda.

### 5.4 Revisiting Entailment Regimes

The current SPARQL specification does not treat entailment regimes beyond RDF simple entailment. Strictly speaking, even RDF entailment is already problematic as a basis for SPARQL query evaluation; a simple query pattern like  $P = (?X, \text{rdf:type}, \text{rdf:Property})$  would have infinitely many solutions even on the empty (sic!) dataset by matching the infinitely many axiomatic triples in the RDF(S) semantics.

Finite rule sets which approximate the RDF(S) semantics in terms of positive Datalog rules [17] have been im-

plemented in systems like TRIPLE<sup>14</sup> or JENA<sup>15</sup>. Similarly, fragments and extensions of OWL [12, 3, 14] definable in terms of Datalog rule bases have been proposed in the literature. Such rule bases can be parametrically combined with our translations, implementing what one might call RDFS<sup>-</sup> or OWL<sup>-</sup> entailment at least. It remains to be seen whether the SPARQL working group will define such reduced entailment regimes.

More complex issues arise when combining a nonmonotonic query language like SPARQL with ontologies in OWL. An embedding of SPARQL into a nonmonotonic rules language might provide valuable insights here, since it opens up a whole body of work done on combinations of such languages with ontologies [7, 19].

## 6. CONCLUSIONS & OUTLOOK

In this paper, we presented three possible semantics for SPARQL based on [16] which differ mainly in their treatment of joins and their translations to Datalog rules. We discussed intuitive behavior of these different joins in several examples. As it turned out, the s-joining semantics which is close to traditional treatment of joins over incomplete relations and the c-joining semantics are nicely embeddable into Datalog. The b-joining semantics which reflects the normative behavior as described by the current SPARQL specification is most difficult to translate. We also suggested some extension of SPARQL, based on this translation. Further, we hope to have contributed to clarifying the relationships between the Query, Rules and Ontology layers of the Semantic Web architecture with the present work.

A prototype of the presented translation has been implemented on top of the dlhex system, a flexible framework for developing extensions for the declarative Logic Programming Engine DLV<sup>16</sup>. The prototype is available as a plugin at <http://con.fusion.at/dlhex/>. The web-page also provides an online interface for evaluation, where the reader can check translation results for various example queries, which we had to omit here for space reasons. We currently implemented the c-joining and b-joining semantics and we plan to gradually extend the prototype towards the features mentioned in Sec. 5, in order to query mixed RDF+SPARQL rule and fact bases. Implementation of further extensions, such as the integration of aggregates typical for database query language, and recently defined for recursive Datalog programs in a declarative way compatible with the answer set semantics [9], are on our agenda. We are currently not aware of any other engine implementing the full semantics defined in [16].

## 7. ACKNOWLEDGMENTS

Special thanks go to Jos de Bruijn and Reto Krummenacher for discussions on earlier versions of this document, to Bijan Parsia, Jorge Pérez, and Andy Seaborne for valuable email-discussions, to Roman Schindlauer for his help on prototype implementation on top of dlhex, and to the anonymous reviewers for various useful comments. This work is partially supported by the Spanish MEC under the project TIC-2003-9001 and by the EC funded projects TripCom (FP6-027324) and KnowledgeWeb (IST 507482).

## 8. REFERENCES

- [1] C. Baral. *Knowledge Representation, Reasoning and Declarative Problem Solving*. Cambr. Univ. Press, 2003.
- [2] D. Beckett. Turtle - Terse RDF Triple Language. Tech. Report, 4 Apr. 2006.
- [3] J. de Bruijn, A. Polleres, R. Lara, D. Fensel. OWL DL vs. OWL Flight: Conceptual modeling and reasoning for the semantic web. In *Proc. WWW-2005*, 2005.
- [4] J. Carroll, C. Bizer, P. Hayes, P. Stickler. Named graphs. *Journal of Web Semantics*, 3(4), 2005.
- [5] R. Cyganiak. A relational algebra for sparql. Tech. Report HPL-2005-170, HP Labs, Sept. 2005.
- [6] J. de Bruijn, E. Franconi, S. Tessaris. Logical reconstruction of normative RDF. *OWL: Experiences and Directions Workshop (OWLED-2005)*, 2005.
- [7] T. Eiter, G. Ianni, A. Polleres, R. Schindlauer, H. Tompits. Reasoning with rules and ontologies. *Reasoning Web 2006*, 2006. Springer
- [8] T. Eiter, G. Ianni, R. Schindlauer, H. Tompits. A Uniform Integration of Higher-Order Reasoning and External Evaluations in Answer Set Programming. *Int'l Joint Conf. on Art. Intelligence (IJCAI)*, 2005.
- [9] W. Faber, N. Leone, G. Pfeifer. Recursive aggregates in disjunctive logic programs: Semantics and complexity. *Proc. of the 9th European Conf. on Art. Intelligence (JELIA 2004)*, 2004. Springer.
- [10] A. V. Gelder, K. Ross, J. Schlipf. Unfounded sets and well-founded semantics for general logic programs. *7<sup>th</sup> ACM Symp. on Principles of Database Systems*, 1988.
- [11] M. Gelfond, V. Lifschitz. Classical Negation in Logic Programs and Disjunctive Databases. *New Generation Computing*, 9:365–385, 1991.
- [12] B. N. Groszof, I. Horrocks, R. Volz, S. Decker. Description logic programs: Combining logic programs with description logics. *Proc. WWW-2003*, 2003.
- [13] P. Hayes. RDF semantics. *W3C Recommendation*, 10 Feb. 2004. <http://www.w3.org/TR/rdf-mt/>
- [14] H. J. ter Horst. Completeness, decidability and complexity of entailment for RDF Schema and a semantic extension involving the OWL vocabulary. *Journal of Web Semantics*, 3(2), July 2005.
- [15] J. W. Lloyd, R. W. Topor. Making prolog more expressive. *Journal of Logic Programming*, 1(3):225–240, 1984.
- [16] J. Pérez, M. Arenas, C. Gutierrez. Semantics and complexity of SPARQL. *The Semantic Web – ISWC 2006*, 2006. Springer.
- [17] A. Polleres, C. Feier, A. Harth. Rules with contextually scoped negation. *Proc. 3<sup>rd</sup> European Semantic Web Conf. (ESWC2006)*, 2006. Springer.
- [18] E. Prud'hommeaux, A. S. (ed.). SPARQL Query Language for RDF, *W3C Working Draft*, 4 Oct. 2006. <http://www.w3.org/TR/rdf-sparql-query/>
- [19] R. Rosati. Reasoning with Rules and Ontologies. *Reasoning Web 2006*, 2006. Springer.
- [20] SQL-99. Information Technology - Database Language SQL- Part 3: Call Level Interface (SQL/CLI). Technical Report INCITS/ISO/IEC 9075-3, INCITS/ISO/IEC, Oct. 1999. Standard specification.
- [21] J. D. Ullman. *Principles of Database and Knowledge Base Systems*. Computer Science Press, 1989.

<sup>14</sup><http://triple.semanticweb.org/>

<sup>15</sup><http://jena.sourceforge.net/>

<sup>16</sup><http://www.dlvsystem.com/>