

SPARQL Formalization

Marcelo Arenas, Claudio Gutierrez, Jorge Pérez

Department of Computer Science
Pontificia Universidad Católica de Chile
Universidad de Chile

Center for Web Research
<http://www.cwr.cl>

SPARQL: A simple RDF query language

```
SELECT ?Name ?Email
WHERE
{
  ?X :name ?Name
  ?X :email ?Email
}
```

- ▶ The *semantics* of simple SPARQL queries is easy to understand, at least intuitively.

“Give me the name and email of the resources in the datasource”

But things can become more complex...

Interesting features of pattern matching on graphs

- ▶ Grouping
- ▶ Optional parts
- ▶ Nesting
- ▶ Union of patterns
- ▶ Filtering
- ▶

```
{ { P1
  P2 }
  OPTIONAL { P5 } }

{ P3
  P4 }
  OPTIONAL { P7 } }
  OPTIONAL { P8 } }
}
UNION
{ P9 }
  FILTER ( R ) }
```

A formal semantics for SPARQL is needed.

A formal approach would be beneficial

- ▶ Clarifying corner cases
- ▶ Helping in the implementation process
- ▶ Providing sound foundations

We will see:

- ▶ A formal compositional semantics based on **[PAG06: Semantics and Complexity of SPARQL]**
- ▶ This formalization is the starting point of the official semantics of the SPARQL language by the W3C.

Outline

Motivation

Basic Syntax

Semantics

Datasets

Query result forms

Dealing with bnodes

Dealing with duplicates

First of all, a simplified algebraic syntax

- ▶ Triple patterns: RDF triple + variables (no bnodes for now)

$$(?X, \text{name}, ?Name)$$

- ▶ The base case for the algebra is a set of triple patterns

$$\{t_1, t_2, \dots, t_k\}.$$

This is called **basic graph pattern** (BGP).

Example

$$\{ (?X, \text{name}, ?Name), (?X, \text{email}, ?Email) \}$$

First of all, a simplified algebraic syntax (cont.)

- ▶ We consider initially three basic operators:

AND, **UNION**, **OPT**.

- ▶ We will use them to construct graph pattern expressions from basic graph patterns.
- ▶ A SPARQL graph pattern:

$$(((\{t_1, t_2\} \text{ AND } t_3) \text{ OPT } \{t_4, t_5\}) \text{ AND } (t_6 \text{ UNION } \{t_7, t_8\}))$$

it is a **full parenthesized expression**

- ▶ Full parenthesized expressions give us **explicit** precedence/association.

Mappings: building block for the semantics

Definition

A mapping is a **partial function** from variables to RDF terms.

Given a mapping μ and a basic graph pattern P :

- ▶ $\text{dom}(\mu)$: the domain of μ .
- ▶ $\mu(P)$: the set obtained from P replacing the variables according to μ

Example

$$\mu = \{?X \rightarrow R_1, ?Y \rightarrow R_2, ?Name \rightarrow \text{john}, ?Email \rightarrow \text{J@ed.ex}\}$$

$$P = \{(?X, \text{name}, ?Name), (?X, \text{email}, ?Email)\}$$

$$\mu(P) = \{(R_1, \text{name}, \text{john}), (R_1, \text{email}, \text{J@ed.ex})\}$$

The semantics of basic graph pattern

Definition

The evaluation of the BGP P over a graph G , denoted by $[[P]]_G$, is the set of all mappings μ such that:

- ▶ $\text{dom}(\mu)$ is exactly the set of variables occurring in P
- ▶ $\mu(P) \subseteq G$

Example

G
(R_1 , name, john)
(R_1 , email, J@ed.ex)
(R_2 , name, paul)

$\llbracket \{ (?X, \text{name}, ?Y) \} \rrbracket_G$

$$\left\{ \begin{array}{l} \mu_1 = \{ ?X \rightarrow R_1, ?Y \rightarrow \text{john} \} \\ \mu_2 = \{ ?X \rightarrow R_2, ?Y \rightarrow \text{paul} \} \end{array} \right\} \begin{array}{l} \mu_1 \\ \mu_2 \end{array}$$

?X	?Y
R_1	john
R_2	paul

$\llbracket \{ (?X, \text{name}, ?Y), (?X, \text{email}, ?Z) \} \rrbracket_G$

$$\left\{ \mu = \{ ?X \rightarrow R_1, ?Y \rightarrow \text{john}, ?Z \rightarrow \text{J@ed.ex} \} \right\}$$

$$\mu$$

?X	?Y	?Z
R_1	john	J@ed.ex

Example

G
(R_1 , name, john)
(R_1 , email, J@ed.ex)
(R_2 , name, paul)

$\llbracket \{ (R_1, \text{webPage}, ?W) \} \rrbracket_G$
 $\{ \}$

$\llbracket \{ (R_3, \text{name}, \text{ringo}) \} \rrbracket_G$
 $\{ \}$

$\llbracket \{ (R_2, \text{name}, \text{paul}) \} \rrbracket_G$
 $\{ \mu_\emptyset = \{ \} \}$

$\llbracket \{ \} \rrbracket_G$
 $\{ \mu_\emptyset = \{ \} \}$

Compatible mappings: mappings that can be merged.

Definition

The mappings μ_1 , μ_2 are **compatibles** iff they **agree** in their **shared variables**:

- ▶ $\mu_1(?X) = \mu_2(?X)$ for every $?X \in \text{dom}(\mu_1) \cap \text{dom}(\mu_2)$.

$\mu_1 \cup \mu_2$ is also a mapping.

Example

	?X	?Y	?U	?V
μ_1	R_1	john		
μ_2	R_1		J@edu.ex	
μ_3			P@edu.ex	R_2
$\mu_1 \cup \mu_2$	R_1	john	J@edu.ex	
$\mu_1 \cup \mu_3$	R_1	john	P@edu.ex	R_2

$\mu_\emptyset = \{ \}$ is compatible with every mapping.

Sets of mappings and operations

Let M_1 and M_2 be sets of mappings:

Definition

Join: $M_1 \bowtie M_2$

- ▶ $\{\mu_1 \cup \mu_2 \mid \mu_1 \in M_1, \mu_2 \in M_2, \text{ and } \mu_1, \mu_2 \text{ are compatibles}\}$
- ▶ extending mappings in M_1 with compatible mappings in M_2

will be used to define **AND**

Definition

Union: $M_1 \cup M_2$

- ▶ $\{\mu \mid \mu \in M_1 \text{ or } \mu \in M_2\}$
- ▶ mappings in M_1 plus mappings in M_2 (the usual set union)

will be used to define **UNION**

Sets of mappings and operations

Definition

Difference: $M_1 \setminus M_2$

- ▶ $\{\mu \in M_1 \mid \text{for all } \mu' \in M_2, \mu \text{ and } \mu' \text{ are not compatibles}\}$
- ▶ mappings in M_1 that cannot be extended with mappings in M_2

Definition

Left outer join: $M_1 \bowtie M_2 = (M_1 \bowtie M_2) \cup (M_1 \setminus M_2)$

- ▶ extension of mappings in M_1 with compatible mappings in M_2
- ▶ plus the mappings in M_1 that cannot be extended.

will be used to define **OPT**

Definition

Given a graph G the evaluation of a pattern is recursively defined

- ▶ $[[(P_1 \text{ AND } P_2)]]_G = [[P_1]]_G \bowtie [[P_2]]_G$
- ▶ $[[(P_1 \text{ UNION } P_2)]]_G = [[P_1]]_G \cup [[P_2]]_G$
- ▶ $[[(P_1 \text{ OPT } P_2)]]_G = [[P_1]]_G \bowtie [[P_2]]_G$

the base case is the evaluation of a BGP.

Example (AND)

G : $(R_1, \text{name, john})$ $(R_2, \text{name, paul})$ $(R_3, \text{name, ringo})$
 $(R_1, \text{email, J@ed.ex})$ $(R_3, \text{email, R@ed.ex})$
 $(R_3, \text{webPage, www.ringo.com})$

$\llbracket \{ (?X, \text{name, ?N}) \} \text{ AND } \{ (?X, \text{email, ?E}) \} \rrbracket_G$

$\llbracket \{ (?X, \text{name, ?N}) \} \rrbracket_G \bowtie \llbracket \{ (?Y, \text{email, ?E}) \} \rrbracket_G$

	<table border="1"><thead><tr><th>?X</th><th>?N</th></tr></thead><tbody><tr><td>R_1</td><td>john</td></tr><tr><td>R_2</td><td>paul</td></tr><tr><td>R_3</td><td>ringo</td></tr></tbody></table>	?X	?N	R_1	john	R_2	paul	R_3	ringo	\bowtie	<table><tbody><tr><td>μ_4</td><td><table border="1"><thead><tr><th>?X</th><th>?E</th></tr></thead><tbody><tr><td>R_1</td><td>J@ed.ex</td></tr><tr><td>R_3</td><td>R@ed.ex</td></tr></tbody></table></td></tr><tr><td>μ_5</td><td></td></tr></tbody></table>	μ_4	<table border="1"><thead><tr><th>?X</th><th>?E</th></tr></thead><tbody><tr><td>R_1</td><td>J@ed.ex</td></tr><tr><td>R_3</td><td>R@ed.ex</td></tr></tbody></table>	?X	?E	R_1	J@ed.ex	R_3	R@ed.ex	μ_5	
?X	?N																				
R_1	john																				
R_2	paul																				
R_3	ringo																				
μ_4	<table border="1"><thead><tr><th>?X</th><th>?E</th></tr></thead><tbody><tr><td>R_1</td><td>J@ed.ex</td></tr><tr><td>R_3</td><td>R@ed.ex</td></tr></tbody></table>	?X	?E	R_1	J@ed.ex	R_3	R@ed.ex														
?X	?E																				
R_1	J@ed.ex																				
R_3	R@ed.ex																				
μ_5																					

	<table border="1"><thead><tr><th>?X</th><th>?N</th><th>?E</th></tr></thead><tbody><tr><td>R_1</td><td>john</td><td>J@ed.ex</td></tr><tr><td>R_3</td><td>ringo</td><td>R@ed.ex</td></tr></tbody></table>	?X	?N	?E	R_1	john	J@ed.ex	R_3	ringo	R@ed.ex
?X	?N	?E								
R_1	john	J@ed.ex								
R_3	ringo	R@ed.ex								
$\mu_1 \cup \mu_4$										
$\mu_3 \cup \mu_5$										

Example (UNION)

G : $(R_1, \text{ name, john})$ $(R_2, \text{ name, paul})$ $(R_3, \text{ name, ringo})$
 $(R_1, \text{ email, J@ed.ex})$ $(R_3, \text{ email, R@ed.ex})$
 $(R_3, \text{ webPage, www.ringo.com})$

$\llbracket \{(\text{?X}, \text{ email}, \text{ ?Info})\} \text{ UNION } \{(\text{?X}, \text{ webPage}, \text{ ?Info})\} \rrbracket_G$

$\llbracket \{(\text{?X}, \text{ email}, \text{ ?Info})\} \rrbracket_G \cup \llbracket \{(\text{?X}, \text{ webPage}, \text{ ?Info})\} \rrbracket_G$

μ_1	?X	?Info
	R_1	J@ed.ex
μ_2	R_3	R@ed.ex

 \cup

μ_3	?X	?Info
	R_3	www.ringo.com

	?X	?Info
μ_1	R_1	J@ed.ex
μ_2	R_3	R@ed.ex
μ_3	R_3	www.ringo.com

Boolean filter expressions (value constraints)

In filter expressions we consider

- ▶ the equality $=$ among variables and RDF terms
- ▶ a unary predicate **bound**
- ▶ boolean combinations (\wedge , \vee , \neg)

A mapping μ **satisfies**

- ▶ $?X = c$ if $\mu(?X) = c$
- ▶ $?X = ?Y$ if $\mu(?X) = \mu(?Y)$
- ▶ **bound**($?X$) if μ is defined in $?X$, i.e. $?X \in \text{dom}(\mu)$

Satisfaction of value constraints

- ▶ If P is a graph pattern and R is a value constraint then $(P \text{ FILTER } R)$ is also a graph pattern.

Definition

Given a graph G

- ▶ $[[(P \text{ FILTER } R)]]_G = \{ \mu \in [[P]]_G \mid \mu \text{ satisfies } R \}$
i.e. mappings in the evaluation of P that **satisfy** R .

Example (FILTER)

G : $(R_1, \text{name}, \text{john})$ $(R_2, \text{name}, \text{paul})$ $(R_3, \text{name}, \text{ringo})$
 $(R_1, \text{email}, \text{J@ed.ex})$ $(R_3, \text{email}, \text{R@ed.ex})$
 $(R_3, \text{webPage}, \text{www.ringo.com})$

$\llbracket (\{ (\{ (?X, \text{name}, ?N) \} \text{OPT} \{ (?X, \text{email}, ?E) \} \text{FILTER} \neg \text{bound}(?E)) \} \rrbracket_G$

	$?X$	$?N$	$?E$	
$\mu_1 \cup \mu_4$	R_1	john	J@ed.ex	$\neg \text{bound}(?E)$
$\mu_3 \cup \mu_5$	R_3	ringo	R@ed.ex	
μ_2	R_2	paul		

	$?X$	$?N$
μ_2	R_2	paul

FILTER: differences with the official specification

- ▶ We restrict to the case in which all variables in R are mentioned in P .
- ▶ This restriction is not imposed in the official specification by W3C.
- ▶ The semantics without the restriction does not modify the expressive power of the language.

- ▶ One of the interesting features of SPARQL is that a query may retrieve data from different sources.

Definition

A SPARQL **dataset** is a set

$$\mathcal{D} = \{G_0, \langle u_1, G_1 \rangle, \langle u_2, G_2 \rangle, \dots, \langle u_n, G_n \rangle\}$$

- ▶ G_0 is the default graph, $\langle u_i, G_i \rangle$ are named graphs
- ▶ $\text{name}(\mathcal{D}) = \{u_1, u_2, \dots, u_n\}$
- ▶ $d_{\mathcal{D}}$ is a function such $d_{\mathcal{D}}(u_i) = G_i$.

The GRAPH operator

if u is an IRI, $?X$ is a variable and P is a graph pattern, then

- ▶ $(u \text{ GRAPH } P)$ is a graph pattern
- ▶ $(?X \text{ GRAPH } P)$ is a graph pattern

GRAPH will permit us to dynamically change the graph against which our pattern is evaluated.

Definition

Given a dataset \mathcal{D} and a graph pattern P

$$\llbracket (u \text{ GRAPH } P) \rrbracket_G = \llbracket P \rrbracket_{d_{\mathcal{D}}(u)}$$

$$\llbracket (?X \text{ GRAPH } P) \rrbracket_G = \bigcup_{u \in \text{name}(\mathcal{D})} \left(\llbracket P \rrbracket_{d_{\mathcal{D}}(u)} \bowtie \{ \{ ?X \rightarrow u \} \} \right)$$

Definition

The evaluation of a general pattern P against a dataset \mathcal{D} , denoted by $\llbracket P \rrbracket_{\mathcal{D}}$, is the set $\llbracket P \rrbracket_{G_0}$ where G_0 is the default graph in \mathcal{D} .

Example (GRAPH)

\mathcal{D}

G_0 :
 $\langle \text{tb}, G_1: (R_1, \text{name}, \text{john}) \quad (R_2, \text{name}, \text{paul}) \rangle$
 $(R_1, \text{email}, \text{J@ed.ex})$
 $\langle \text{trs}, G_2: (R_4, \text{name}, \text{mick}) \quad (R_5, \text{name}, \text{keith}) \rangle$
 $(R_4, \text{email}, \text{M@ed.ex}) \quad (R_5, \text{email}, \text{K@ed.ex}) \rangle$

$[[(\text{trs GRAPH } \{(?X, \text{name}, ?N)\})]]_{\mathcal{D}}$

$[[(\text{trs GRAPH } \{(?X, \text{name}, ?N)\})]]_{G_0}$

$[[\{(?X, \text{name}, ?N)\}]]_{G_2}$

	$?X$	$?N$
μ_1	R_4	mick
μ_2	R_5	keith

Example (GRAPH)

\mathcal{D}

G_0 :
 $\langle \text{tb}, G_1: (R_1, \text{name}, \text{john}) \quad (R_2, \text{name}, \text{paul}) \rangle$
 $(R_1, \text{email}, \text{J@ed.ex})$
 $\langle \text{trs}, G_2: (R_4, \text{name}, \text{mick}) \quad (R_5, \text{name}, \text{keith}) \rangle$
 $(R_4, \text{email}, \text{M@ed.ex}) \quad (R_5, \text{email}, \text{K@ed.ex}) \rangle$

$\llbracket (?G \text{ GRAPH } \{(?X, \text{name}, ?N)\}) \rrbracket_{\mathcal{D}}$

$\llbracket \{(?X, \text{name}, ?N)\} \rrbracket_{G_1} \bowtie \{\{?G \rightarrow \text{tb}\}\} \cup$
 $\llbracket \{(?X, \text{name}, ?N)\} \rrbracket_{G_2} \bowtie \{\{?G \rightarrow \text{trs}\}\}$

μ_1

?X	?N
R_1	john
R_2	paul

 $\bowtie \{\{?G \rightarrow \text{tb}\}\} \cup$
 μ_3

?X	?N
R_4	mick
R_5	keith

 $\bowtie \{\{?G \rightarrow \text{trs}\}\}$

?G	?X	?N
tb	R_1	john
tb	R_2	paul
trs	R_4	mick
trs	R_5	keith

SELECT

- ▶ Up to this point we have concentrated in the **body** of a SPARQL query, i.e. in the graph pattern matching expression.
- ▶ A query can also process the values of the variables. The most simple processing operation is the **selection** of some variables appearing in the query.

Definition

- ▶ A SELECT query is a tuple (W, P) where P is a graph pattern and W is a set of variable.
- ▶ The answer of a SELECT query against a dataset \mathcal{D} is

$$\{\mu|_W \mid \mu \in \llbracket P \rrbracket_{\mathcal{D}}\}$$

where $\mu|_W$ is the restriction of μ to domain W .

CONSTRUCT

- ▶ A query can also output an RDF graph.
- ▶ The construction of the output graph is based on a **template**.
- ▶ A **template** is a set of triple patterns possibly with bnodes.

Example

$$T_1 = \{(?X, \text{name}, ?Y), (?X, \text{info}, ?I), (?X, \text{addr}, B)\}$$

with B a bnode

Definition

- ▶ A CONSTRUCT query is a tuple (T, P) where P is a graph pattern and T is a template.

Definition

The answer of a CONSTRUCT query (T, P) against a dataset \mathcal{D} is obtained by

- ▶ for every $\mu \in \llbracket P \rrbracket_{\mathcal{D}}$ create a template T_{μ} with **fresh bnodes**
- ▶ take the union of $\mu(T_{\mu})$ for every $\mu \in \llbracket P \rrbracket_{\mathcal{D}}$
- ▶ discard the **not valid** RDF triples
 - ▶ some variables have not been instantiated.
 - ▶ bnodes in predicate positions

Blank nodes in graph patterns

- ▶ We allow now bnodes in triple patterns.
- ▶ Bnodes act as existentials **scoped to the basic graph pattern**.

Definition

The evaluation of the BGP P with bnodes over the graph G denoted $\llbracket P \rrbracket_G$, is the set of all mappings μ such that:

- ▶ $\text{dom}(\mu)$ is exactly the set of variables occurring in P ,
- ▶ there exists a function θ from bnodes of P to G such that

$$\mu(\theta(P)) \subseteq G.$$

- ▶ A natural extension of BGPs without bnodes.
- ▶ The algebra remains the same.

Bag/Multiset semantics

- ▶ In a **bag**, a mapping can have cardinality greater than one.
- ▶ Every mapping μ in a bag M is annotated with an integer $c_M(\mu)$ that represents its cardinality ($c_M(\mu) = 0$ if $\mu \notin M$).
- ▶ Operations between sets of mappings can be extended to bags maintaining duplicates:

Definition

$$\mu \in M = M_1 \bowtie M_2, \quad c_M(\mu) = \sum_{\mu = \mu_1 \cup \mu_2} c_{M_1}(\mu_1) \cdot c_{M_2}(\mu_2),$$

$$\mu \in M = M_1 \cup M_2, \quad c_M(\mu) = c_{M_1}(\mu) + c_{M_2}(\mu),$$

$$\mu \in M = M_1 \setminus M_2, \quad c_M(\mu) = c_{M_1}(\mu).$$

- ▶ Intuition: we simply do not discard duplicates.

- ▶ R. Cyganiak, *A Relational Algebra for SPARQL*. Tech Report HP Laboratories, HPL-2005-170.
- ▶ E. Prud'hommeaux, A. Seaborne, *SPARQL Query Language for RDF*. W3C Working Draft, 2007.
- ▶ J. Pérez, M. Arenas, C. Gutierrez, *Semantics and Complexity of SPARQL*. In *Int. Semantic Web Conference 2006*.
- ▶ J. Pérez, M. Arenas, C. Gutierrez, *Semantics of SPARQL*. Tech Report Universidad de Chile 2006, TR/DCC-2006-17.