

# The SPARQL query language for RDF

Gilles Falquet

# SPARQL – based on graph patterns

Triple pattern a triple from

$(\text{RDF-Term} \cup \text{Var}) \times (\text{IRI} \cup \text{Var}) \times (\text{RDF-Term} \cup \text{Var})$

- RDF-Term : IRI or literal or blank
- Var : variable

Graph pattern

- a set of triple patterns

Same syntax as Turtle + variables

```
{ ?x ex:color "red" }
```

```
{ ?x ex:friend ?y . ?y rdf:type ex:Cat }
```

```
{ x:bob ex:member ?v . ?v rdf:type ex:Club . ?v ex:country ?w }
```

```
{ ?x ex:address _:adr . _:adr ex:city ex:Madrid }
```

# SPARQL Basic Graph Pattern Query

prefix definitions  
select *output variables*  
from *graph*  
where { *basic graph pattern* }

```
prefix foaf: <...>
select ?p
from <http://cui.unige.ch/g1>
where {
    ?p foaf:knows ?q . ?q foaf:familyName "Zep"
}
```

# Definition: Basic Graph Pattern Matching

Let BGP be a basic graph pattern and let  $G$  be an RDF graph.

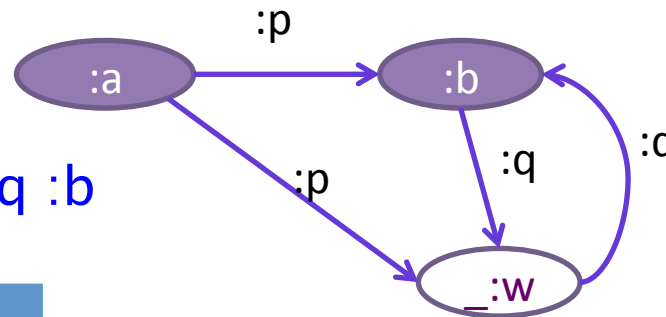
$\mu$  is a **solution** for BGP from  $G$  when

- there is a pattern instance mapping  $P$  such that  $P(\text{BGP})$  is a subgraph of  $G$ 
  - $P$  maps variables and blank nodes to RDF-terms
- and  $\mu$  is the restriction of  $P$  to the query variables in BGP.

# Example

On the graph

`:a :p :b.`   `:a :p _:w.`   `:b :q _:w.`   `_:w :q :b`



solutions of { `?x ?y :b` }

<code>?x</code>	<code>?y</code>
<code>:a</code>	<code>:p</code>
<code>_:w</code>	<code>:q</code>

solutions of { `?x :p ?y . ?y :q ?z` }

<code>?x</code>	<code>?y</code>	<code>?z</code>
<code>:a</code>	<code>:b</code>	<code>_:w</code>
<code>:a</code>	<code>_:w</code>	<code>:b</code>

solutions of { `?x :p _:h . _:h :q ?z` }

<code>?x</code>	<code>?z</code>
<code>:a</code>	<code>_:w</code>
<code>:a</code>	<code>:b</code>

# Simple Graph Patterns are not Enough

Need to express

- disjunctions (match this or that)
- optional parts in patterns (match if possible)
- negations (match this but not that)
- conditions on variable values (  $<$ ,  $>$ ,  $=$ , ...)
- multiple paths (path expressions) in patterns

Need to process the results

- combine the solution variables (+, -, ...)
- aggregation functions (sum, average, ...)
- ordering
- grouping

# Optional parts

$\{pattern_1 \text{ OPTIONAL } \{ pattern_2 \}\}$

Find solutions for  $\{ pattern_1 pattern_2 \}$  and for  $\{ pattern_1 \}$

In the solutions for  $pattern_1$  only, the variables that appear in  $pattern_2$  only are unbound.

# Example

On the graph

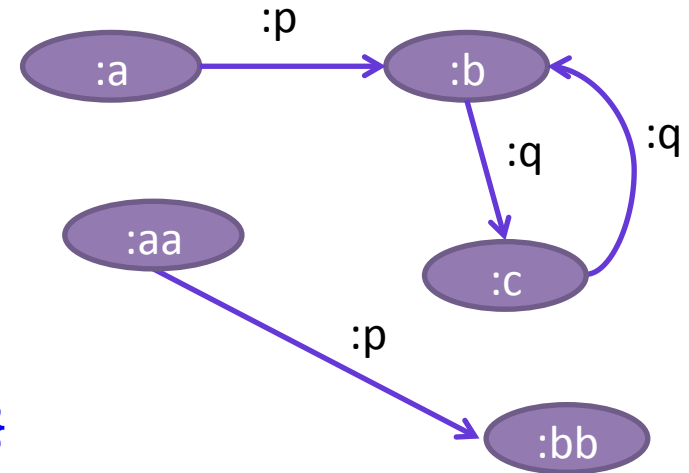
`:a :p :b. :aa :p :bb. :b :q :c.`

The solutions of

`{ ?x :p ?y OPTIONAL {?y :q ?z} }`

are

?x	?y	?z
:a	:b	:c
:aa	:bb	UNBOUND





# Union

To represent disjunctions

a solution to

*pattern1* UNION *pattern2*

is a solution to *pattern1* or to *pattern2* (or both)

# Example

"Find people who own a cat or a dog"

```
{?p a :Person. ?p :owns ?a. ?a a :Cat }
```

UNION

```
{?p a :Person. ?p :owns ?a. ?a a :Dog }
```

can be simplified by using group graph patterns

```
{?p a :Person. ?p :owns ?a. {{?a a :Cat} UNION {?a a :Dog}}}
```

# Filtering with a boolean expression

*{pattern FILTER( expression ) }*

retain only the solutions to *pattern* for which *expression* evaluates to true

## Example

```
{ ?x a :Car. ?x :price ?p. ?x :category ?c  
  FILTER(?p < 10000 && ?c != :sport) }
```

# Testing For the Absence of a Pattern

Data:

```
:alice rdf:type foaf:Person . :alice foaf:name "Alice" .  
:bob rdf:type foaf:Person .
```

Query:

```
SELECT ?person  
WHERE { ?person rdf:type foaf:Person .  
        FILTER NOT EXISTS { ?person foaf:name ?name } }
```

Query Result:

```
:bob
```

# Testing For the Presence of a Pattern

Query:

```
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
SELECT ?person
WHERE { ?person rdf:type foaf:Person .
        FILTER EXISTS { ?person foaf:name ?name } }
```

Query Result:

```
<http://example/alice>
```

# Removing Possible Solutions

MINUS evaluates both its arguments, then calculates solutions in the left-hand side that are not compatible with the solutions on the right-hand side.

```
:alice foaf:givenName "Alice" ; foaf:familyName "Smith" .  
:bob foaf:givenName "Bob" ; foaf:familyName "Jones" .  
:carol foaf:givenName "Carol" ; foaf:familyName "Smith" .
```

```
SELECT DISTINCT ?s  
WHERE { ?s ?p ?o . MINUS { ?s foaf:givenName "Bob" . } }
```

Results:

```
<http://example/carol>  
<http://example/alice>
```

# Relationship and differences between NOT EXISTS and MINUS

NOT EXISTS and MINUS represent two ways of thinking about negation

- one based on testing whether a pattern exists in the data, given the bindings already determined by the query pattern,
- one based on removing matches based on the evaluation of two patterns. In some cases they can produce different answers.

@prefix : <http://example/> .

:a :b :c .

```
SELECT * { ?s ?p ?o FILTER NOT EXISTS { ?x ?y ?z } }
```

No solutions because { ?x ?y ?z } matches given any ?s ?p ?o

```
SELECT * { ?s ?p ?o MINUS { ?x ?y ?z } }
```

There is no shared variable between the first part (?s ?p ?o) and the second (?x ?y ?z) so no bindings are eliminated.

## Results:

```
<http://example/a> <http://example/b><http://example/c>
```



# Property path

iri	
$\wedge$ elt	inverse path
elt / elt	sequence
elt   elt	alternative
elt*	repetition (0...n)
elt+	repetition (1...n)
elt?	option
$!iri$ or $!(iri_1   \dots   iri_n)$	negation
$!\wedge iri$ or $!(\wedge iri_1   \dots   \wedge iri_n)$	negation of the inverse
$!(iri_1   \dots   iri_j   \wedge iri_{j+1}   \dots   \wedge iri_n)$	

# Using property path to access lists

Recall that a list structure, written as

```
:france :flagColors (:blue :white :red) .
```

is an abbreviation for:

```
:france :flagColors [  
  rdf:type rdf:List ;  
  rdf:first :blue ;  
  rdf:rest [  
    rdf:type rdf:List ;  
    rdf:first :white ;  
    rdf:rest [  
      rdf:type rdf:List ;  
      rdf:first :red ;  
      rdf:rest rdf:nil ]]]
```

Some queries over such structures can only be solved by using property path expressions.<sup>1</sup>

### Find all the colors in the french flag

```
select ?c
where { :france :flagColors/rdf:rest*/rdf:first ?c }
```

### Find the last color of the french flag

```
select ?c
where { :france :flagColors/rest* ?last.
       ?last rdf:rest rdf:nil. ?last rdf:first ?c }
```

1. Or with entailment regimes that take into account transitive properties

# Accessing Trees

Example: a part is decomposed into subparts, sub-subparts, etc. linked through a `:partOf` property.

Display all the parts that belong to `:b`

```
select ?p where {?p :partOf* :b. }
```

If part `:a` belongs to part `:b`, display its part number

```
select ?pn where { :a :partOf* :b. :a :partNo ?pn }
```

What are the subparts of `:b` that have more than 10 subparts (at any level)

```
select ?p where {?p :partOf* :b.  
filter count(select ?q where {?q partOf* ?p}) > 10 }
```

# Entailment Regimes

Several entailment regimes for RDF, to infer new triples

basic RDF entailment

- e.g.  $(x \text{ p } y) \rightarrow (p \text{ rdf:type } \text{rdf:Property})$

RDFS entailment

- e.g.  $(x \text{ p } y), (p \text{ rdfs:range } z) \rightarrow (y \text{ rdf:type } z)$

OWL entailment

OWL direct semantics

...

SPARQL answers should take into account inferred triples

- (1) `ex:book1 rdf:type ex:Publication .`
- (2) `ex:book2 rdf:type ex:Article .`
- (3) `ex:Article rdfs:subClassOf ex:Publication .`
- (4) `ex:publishes rdfs:range ex:Publication .`
- (5) `ex:MITPress ex:publishes ex:book3 .`

**SELECT ?prop WHERE { ?prop rdf:type rdf:Property }**

Simple entailment  $\rightarrow \emptyset$

RDF entailment  $\rightarrow$

using RDF the entailment rule

$(x \ p \ y) \rightarrow (p \text{ rdf:type } \text{rdf:Property})$

<b>?prop</b>
rdf:type
rdfs:subClassOf
rdfs:range
ex:publishes

- (1) `ex:book1 rdf:type ex:Publication .`
- (2) `ex:book2 rdf:type ex:Article .`
- (3) `ex:Article rdfs:subClassOf ex:Publication .`
- (4) `ex:publishes rdfs:range ex:Publication .`
- (5) `ex:MITPress ex:publishes ex:book3 .`

`SELECT ?pub WHERE { ?pub rdf:type ex:Publication }`

with RDF entailment →

?pub
ex:book1

with RDFS entailment →

using the rules

$(x \text{ a } C), (C \text{ rdfs:subClassOf } D) \rightarrow (x \text{ a } D)$

$(x \text{ p } y), (y \text{ rdfs:range } C) \rightarrow (y \text{ a } D)$

?pub
ex:book1
ex:book2
ex:book3

# Blank nodes in graphs and results

ex:MITPress

ex:published ex:every ;

ex:published \_:2 .

ex:every ex:title "Everything" .

\_:2 ex:date "2011" .

- Blank nodes are local
- They have no URI
- They cannot be "exported" to the answer
- The answer mapping must "invent" blank nodes

select ?pub where {ex:MITPress ex:published ?pub}

Infinitely many possible answers ?

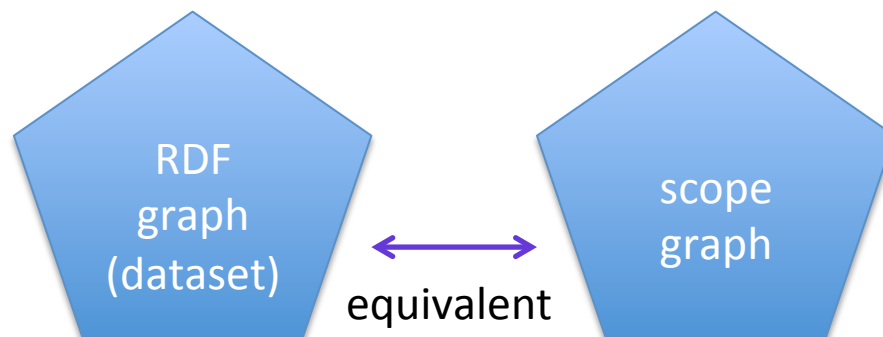
?pub	?pub	?pub
ex:every	ex:every	ex:every
_:cx323322	_:abc	_:u123



# Scoping Graph – to avoid infinite answers

- The SPARQL definition allows the solution mapping to bind a variable in a basic graph pattern, BGP, to a blank node in G.
- Since SPARQL treats blank node identifiers in a results format document as **scoped to the document**, they cannot be understood as identifying nodes in the active graph of the dataset.
- If DS is the dataset of a query, pattern solutions are therefore understood to be not from the active graph of DS itself, but from an RDF graph, called the **scoping graph**, which is graph-equivalent to the active graph of DS but shares no blank nodes with DS or with BGP.
- The same scoping graph is used for all solutions to a single query.

- Since RDF blank nodes allow infinitely many redundant solutions for many patterns, there can be infinitely many pattern solutions (obtained by replacing blank nodes by different blank nodes).
- It is necessary, therefore, to somehow delimit the solutions for a basic graph pattern.
- SPARQL uses the subgraph match criterion to determine the solutions of a basic graph pattern. There is one solution for each distinct pattern instance mapping from the basic graph pattern to a subset of the active graph.



# RDF Entailment

- A solution mapping  $\mu$  is a *possible solution for BGP from G under RDF entailment* if
  - $\text{dom}(\mu) = V(\text{BGP})$
  - and there is an RDF instance mapping  $\sigma$  from  $B(\text{BGP})$  to RDF-T such that  $\text{dom}(\sigma) = B(\text{BGP})$
  - and the pattern instance mapping  $P = (\mu, \sigma)$  is such that  $P(\text{BGP})$  are well-formed RDF triples that are RDF entailed by SG.
- A possible solution  $\mu$  is a *solution for BGP from SG under RDF entailment* if:
  - (C1) The RDF triples  $\text{sk}(P(\text{BGP}))$  are ground and RDF entailed by  $\text{sk}(SG)$ .
  - (C2) For each variable  $x$  in  $V(\text{BGP})$ ,  $\mu(x)$  occurs in SG or in  $\text{rdfV-Minus}$ .
- The multiplicity of  $\mu$  in the multiset of solutions is the maximal number of distinct RDF instance mappings  $\sigma$  that yield a pattern instance mapping  $P = (\mu, \sigma)$  for which  $\mu$  is a solution.

