# SQL Nested Queries in SPARQL

Renzo Angles[1] and Claudio Gutierrez[2]

[1] Department of Computer Science, Universidad de Talca
[2] Department of Computer Science, Universidad de Chile

**Abstract.** SPARQL currently does not include any form of nested queries. In this paper we present a proposal to incorporate nested queries into SPARQL along the design philosophy of SQL nested queries. We present rewriting algorithms and show that all the proposed nested queries can be simulated by a subset of them.

## 1 Introduction

One of the most powerful features of a query language is the nesting of queries (or subqueries), that is the possibility of writing in a single expression a query which uses the output of other queries, i.e., a query can have an arbitrary number of subqueries nested within it. The current W3C recommendation of SPARQL [13] does not include any form of nesting, although it has been considered as an issue by the RDF Data Access Working Group, and have been gradually incorporated into SPARQL engines like ARQ and Virtuoso. In July 2009, the SPARQL Working Group has started up and is currently working, on defining potential extensions to SPARQL [9]; among them the notion of subqueries. To the best of our knowledge, there are not formal syntax and semantics defined for these types of queries in SPARQL.

For SPARQL in particular, the incorporation of nested queries has several motivations. One of the most important is the *reuse of queries.* Once a query is executed, the user may presumably direct the output to some storage medium, assign an IRI to it and then run a query against that extract. Also, with the right interface, the user might be able to just cut-and-paste a query that was debugged into ad hoc queries. Hence, this feature would allow to to building queries incrementally, from separately debugged pieces. Another important motivation, as SPARQL was thought to work on a distributed environment like the Web, is the notion of *distributed queries.* A SPARQL user will have access to a vast and time-varying input RDF graphs, containing a lot of data that is not of interest to the user. Hence the query computation can be distributed and only relevant data used to compute the final query. For example, AllegroGraph supports queries with distributed databases. A third and importante motivation is *query rewriting.* Complex queries can be structured in a way more intuitive and understandable for the user. We provide several examples of these uses.

It is well known that SQL has a long standing and user-tested facilities for nesting queries. Based on the fact of the similarities between SQL and SPARQL

[1,4], it seems natural to investigate how the SQL nesting features can be incorporated into SPARQL. In this paper, we address this challenge.

The findings and contributions of the paper are the following:

1. Define a proposal to incorporate nested queries to SPARQL along the design philosophy of SQL, by presenting the syntax and a formal semantics completely compatibles to current one, and discuss design features.
2. Show, via illustrative examples, that all SQL facilities are also relevant in SPARQL, and present a set of equivalences and rewriting rules among them.
3. Prove that all classical SQL nesting operators (IN, SOME, ALL, etc.) can be reduced into two types of simple and natural forms of nesting: ASK queries inside FILTER and SELECT queries as graph patterns.

This paper is organized as follows: Section 2 presents the syntax and the semantics of the extension of SPARQL with nested queries. Section 3 presents examples of nested queries. Section 4 presents the algorithms to rewrite among nested queries. Finally, Section 5 presents some conclusions.

## 1.1 Related Work

Query nesting in SPARQL has been considered as an issue by the RDF Data Access Working Group. It was raised on July 2004 and was denominated *cascadedQueries*[1]: "capability for multiple queries in the same request, cascading the bindings (or graph) returned from one query into another".

Currently, the W3C Working Group of SPARQL is working on new features for the language [9]. Among them, we found the notion of *Subqueries* (to nest a query within another query). The working draft of SPARQL 1.1 [7] introduces, as possible type of subquery, the inclusion of SELECT queries as graph patterns.

Regarding real-life practice, some implementations of SPARQL provide extensions that include support for some types of nested queries. ARQ, the query engine for Jena, supports a type of nested SELECT which uses aggregate functions[2]. Virtuoso has included some extensions[3] related to nested queries. Among them: an embedded select query in the place of a triple pattern, and filter conditions of the form "exists (<scalar subquery>)".

None of mentioned proposals and/or implementations present systematic covering, analysis of these extensions, nor a formal semantics for them. It is important to note that all of them introduce nested queries as a new pattern of the form (`SELECT ....`). This approach, from a theoretical point of view, combines the construction of solutions with filtering properties, a decision design which bring several problems. For example, it is not clear what the semantics of the evaluation pattern of the form (`SELECT ..?X.. ) AND (..?X..`), where a variable `?X` occurs free in the SELECT should be. Clearly, one of the interesting features of nested queries are correlated variables.

---

[1] http://www.w3.org/2001/sw/DataAccess/issues#cascadedQueries

[2] http://jena.sourceforge.net/ARQ/sub-select.html

[3] http://www.w3.org/2009/sparql/wiki/Extensions_Proposed_By_OpenLink#Nested_Queries

An alternative proposal is to incorporate nesting as a filtering device. For example, Polleres [12] suggested that boolean SPARQL queries (i.e., queries having ASK query form) can be safely allowed within filter constraints, although this work does not further develop this extension. The inspiration of our work is the design philosophy of SQL nested queries, which restrict the nested queries to a form of filtering in the WHERE condition. This approach allows to have a clean semantics for correlated variables, and permits to modularly extend the language. In this paper we investigate both approaches, i.e., SELECT queries as graph patterns and nested queries in filter constraints.

## 1.2 Nested queries in SQL

SQL is a paradigmatic example of the power of nested queries. In fact, without this capability the power of SQL would be severely restricted [8].

SQL allows nesting of query blocks in FROM and WHERE clauses, in any level of nesting. In the FROM clause, a subquery is *imported* and used to conform the set of relations to be queried. A subquery in the WHERE clause can be either aggregate (it returns a single value due to an aggregate operator) or non-aggregate (it returns either a set of values or empty).

Let $Q_A$ be an aggregate query, $Q_S$ and $Q'_S$ be non-aggregate queries, and $\theta$ be a scalar comparison operator ($<, \leq, >, \geq =, \neq$). A selection predicate containing nested queries can be of the form:

(1) $\langle value \mid attribute \rangle \ \theta \ (Q_A)$ (*value-set-comparison predicate*).
(2) $\langle value \mid attribute \rangle \ \theta$ SOME | ALL($Q_S$) (*quantified predicate*).
(3) EXISTS | NOT-EXISTS($Q_S$) (*existential predicate*).
(4) $\langle value \mid attribute \rangle$ IN | NOT-IN ($Q_S$) (*set-membership predicate*).
(5) ($Q_S$) CONTAINS | DOES-NOT-CONTAIN($Q'_S$) (*set-inclusion predicate*).
(6) ($Q_S$) = ($Q'_S$) (*set-equality predicate*).
(7) ($Q_S$) $\neq$ ($Q'_S$) (*set-inequality predicate*).

For example, consider the relations EMPLOYEES(EMP#,NAME,SAL,DEPT) and DEPARTMENTS(DEPT#,NAME,LOCATION). The following SQL expression, present the nested query $Q_1$ where $Q_2$ and $Q_3$ are aggregate and non-aggregate subqueries respectively.

| | |
|---|---|
| SELECT E.NAME FROM EMPLOYEES E | $(Q_1)$ |
| WHERE E.SAL > | |
|     (SELECT AVG(F.SAL) FROM EMPLOYEES F | $(Q_2)$ |
|      WHERE F.DEPT IN | |
|        (SELECT D.DEPT# FROM DEPARTMENTS D | $(Q_3)$ |
|        WHERE D.LOCATION = 'DENVER') ); | |

## 2 Syntax and Semantics of nested queries

In this section we extend the syntax and semantics of SPARQL to support nested queries. This extension, is based on the ideas and syntax of nested queries in SQL.

Considering that aggregate operators for SPARQL have not been defined yet, we are only considering non-aggregate queries (i.e., SELECT and ASK queries) as nested queries. Hence, we have not included value-set-comparison predicates of SQL as defined before. Additionally, we do not consider the orthogonal functionality of composition in SQL which is to have a query in the FROM. In SPARQL this would correspond to the discussion of how to nest queries in the FROM and FROM NAMED clauses. Schenk [14] proposes the use of views as parts of a dataset, that is, the inclusion of CONSTRUCT queries in FROM clauses. We do not address this topic in this paper.

In addition to SQL operators, we allow SELECT queries as graph patterns. As we will see in Section 4, this type of nesting is required for query rewriting.

The definition of this extension follows the formalization presented in [11].

## 2.1 Preliminaries: The RDF Model and RDF Datasets

Assume there are pairwise disjoint infinite sets $I$, $B$, $L$ (IRIs, blank nodes, and RDF literals respectively). We denote by $T$ the union $I \cup B \cup L$ (RDF *terms*). A tuple $(v_1, v_2, v_3) \in (I \cup B) \times I \times T$ is called an *RDF triple*, where $v_1$ is the *subject*, $v_2$ the *predicate*, and $v_3$ the *object*. An *RDF Graph* [10] (just graph from now on) is a set of RDF triples. Given a graph $G$, we denote by $\mathrm{term}(G)$ the set of elements of $T$ appearing in $G$ and by $\mathrm{blank}(G)$ the set of blank nodes in $G$. If $G$ is referred to by an IRI $u$, then $\mathrm{graph}(u)$ returns the graph available in $u$, i.e, $G = \mathrm{graph}(u)$.

We define two operations on two graphs $G_1$ and $G_2$. The *union* of graphs, denoted $G_1 \cup G_2$, is the set theoretical union of their sets of triples. The *merge* of graphs, denoted $G_1 + G_2$, is the graph $G_1 \cup G_2'$ where $G_2'$ is the graph obtained from $G_2$ by renaming its blank nodes to avoid clashes with those in $G_1$.

An *RDF dataset* is a set $D = \{G_0, \langle u_1, G_1 \rangle, \ldots, \langle u_n, G_n \rangle\}$ where each $G_i$ is a graph and each $u_j$ is an IRI. $G_0$ is called the *default graph* and each pair $\langle u_i, G_i \rangle$ is called a *named graph*. Every dataset satisfies that: (i) it always contains one default graph, (ii) there may be no named graphs, (iii) each $u_j$ is distinct, and (iv) $\mathrm{blank}(G_i) \cap \mathrm{blank}(G_j) = \emptyset$ for $i \neq j$. Given $D$, we denote by $\mathrm{term}(D)$ the set of terms occurring in the graphs of $D$. The default graph of $D$ is denoted $\mathrm{dg}(D)$. For a named graph $\langle u_i, G_i \rangle$ define $\mathrm{name}(G_i)_D = u_i$ and $\mathrm{graph}(u_i)_D = G_i$; otherwise $\mathrm{name}(G_i)_D = \emptyset$ and $\mathrm{graph}(u_i)_D = \emptyset$. We denote by $\mathrm{names}(D)$ the set of IRIs $\{u_1, \ldots, u_n\}$. Although $\mathrm{name}(G_0) = \varnothing$, we sometimes will use $g_0$ when referring to $G_0$. Finally, the *active graph* of $D$ is the graph $G_i$ used for querying the dataset.

## 2.2 Syntax of nested queries

Assume the existence of an infinite set $V$ of variables disjoint from $T$. Let $\mathrm{var}(\alpha)$ the function which returns the set of variables occurring in the structure $\alpha$.

A *triple pattern* is a tuple in $(T \cup V) \times (I \cup V) \times (T \cup V)$. A *nested query* is a tuple $(R, F, P)$[4] where $R$ is a result query form, $F$ is a set –possibly empty– of dataset clauses, and $P$ is a graph pattern. Next we define each component.

(1) If $W \subset V$ is a set of variables and $H$ is a set of triple patterns (called a *graph template*) then the expressions SELECT $W$, CONSTRUCT $H$, and ASK are *result query forms*.
(2) If $u \in I$ and $Q_C$ is a query of the form (CONSTRUCT $H, F, P$), then the expressions FROM $u$ and FROM NAMED $u$ are *dataset clauses*.
(3) A *filter constraint* is defined recursively as follows:
  − If $?X, ?Y \in V$ and $v \in I \cup L$ then $?X = v$, $?X = ?Y$, and bound($?X$) are (*atomic*) filter constraints[5].
  − If $u \in T$, $\theta$ is a scalar comparison operator $(=, \neq, <, <=, >, >=)$, and $Q_{?X}$ is a query of the form (SELECT $?X, F, P$), then the expressions $(u \ \theta \ \text{SOME}(Q_{?X}))$, $(u \ \theta \ \text{ALL}(Q_{?X}))$ and $(u \ \text{IN} \ (Q_{?X}))$ are filter constraints.
  − If $Q_A$ is a query of the form (ASK, $F, P$), then the expression EXISTS($Q_A$) is a filter constraint.
  − If $Q_1$ and $Q_2$ are SELECT queries, then the expressions $(Q_1 \ \text{CONTAINS} \ Q_2)$ and $(Q_1 = Q_2)$ are filter constraints.
  − If $C_1$ and $C_2$ are filter constraints, then $(\neg C_1)$, $(C_1 \wedge C_2)$, and $(C_1 \vee C_2)$ are (*complex*) filter constraints.
(4) A *graph pattern* is defined recursively as follows:
  − A triple pattern is a graph pattern.
  − If $P_1$ and $P_2$ are graph patterns then the expressions $(P_1 \ \text{AND} \ P_2)$, $(P_1 \ \text{OPT} \ P_2)$, $(P_1 \ \text{UNION} \ P_2)$, and $(P_1 \ \text{MINUS} \ P_2)$ are graph patterns.[6]
  − If $P$ is a graph pattern and $u \in I \cup V$ then the expression $(u \ \text{GRAPH} \ P)$ is a graph pattern.
  − If $P$ is a graph pattern and $C$ is a filter constraint then the expression $(P \ \text{FILTER} \ C)$ is a graph pattern.
  − If $P$ is a SELECT query then the expression $(Q_S)$ is a graph pattern.

Note that corresponding opposite operators of nesting can be represented by using the negation of filter constraints, i.e., NOT-IN can be expressed as $(\neg(v \ \text{IN} \ Q_S))$, NOT-EXISTS as $(\neg \text{EXISTS}(Q_A))$, DOES-NOT-CONTAIN as $(\neg(Q_1 \ \text{CONTAINS} \ Q_2))$, and set-inequality as $(\neg(Q_1 = Q_2))$.

**Definition 1 (Nested and flat queries).** *Let $Q = (R, F, P)$ be a query. A query $Q'$ is* nested *in $Q$ if and only if $Q'$ occurs in the graph pattern $P$, i.e., when $Q'$ is nested in $P$. In such case, $Q$ is known as the* outer query *and $Q'$ is known as the* inner query. *If $Q$ does not contain nested queries then $Q$ is called a* flat *query.*

Note that the language supports queries with any level of nesting.

---

[4] In this paper we do not consider the solution modifiers defined in [13].

[5] For a complete list of atomic filter constraints see the SPARQL specification [13]

[6] The MINUS operator is not defined in the SPARQL specification, however it can be simulated by a combination of the OPT and FILTER operators [1].

### 2.3 Semantics of nested queries

A *mapping* $\mu$ is a partial function $\mu : V \to T$. The domain of $\mu$, $\mathrm{dom}(\mu)$, is the subset of V where $\mu$ is defined. The *empty mapping* $\mu_0$ is a mapping such that $\mathrm{dom}(\mu_0) = \emptyset$. Given a triple pattern $t$ and a mapping $\mu$ such that $\mathrm{var}(t) \subseteq \mathrm{dom}(\mu)$, $\mu(t)$ is the triple obtained by replacing the variables in $t$ according to $\mu$. Abusing notation, for a query $Q$, we denote by $\mu(Q)$ the query resulting from replacing variables in $Q$ according to $\mu$.

Two mappings $\mu_1$ and $\mu_2$ are *compatible* when for all $?X \in \mathrm{dom}(\mu_1) \cap \mathrm{dom}(\mu_2)$ it satisfies that $\mu_1(?X) = \mu_2(?X)$, i.e., when $\mu_1 \cup \mu_2$ is also a mapping. The operations of *join, union, difference* and *left outer-join* between two sets of mappings $\Omega_1$ and $\Omega_2$ are defined as follows:

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

The *answer* for a query $Q = (R, F, P)$, denoted $\mathrm{ans}(Q)$, is a function which returns: (i) a set of mappings when $R$ is a SELECT query; (ii) an RDF graph when $R$ is a CONSTRUCT query; and (iii) a boolean value (*true / false*) when $R$ is an ASK query. We will use this informal definition of $\mathrm{ans}(\cdot)$ to define the semantics for the components of a query.

(1) *Semantics of result query forms.* Let $\mu$ be a mapping and $R$ be a result query form. The result of $R$ given $\mu$, denoted $\mathrm{result}(R, \mu)$, is defined as follows:
   - If $R$ is SELECT $W$ then $\mathrm{result}(R, \mu)$ is the restriction of $\mu$ to $W$, that is the mapping denoted $\mu_{|W}$ such that $\mathrm{dom}(\mu_{|W}) = \mathrm{dom}(\mu) \cap W$ and $\mu_{|W}(?X) = \mu(?X)$ for every $?X \in \mathrm{dom}(\mu_{|W})$.
   - If $R$ is CONSTRUCT $H$ then $\mathrm{result}(R, \mu)$ is the set of RDF triples (i.e. an RDF graph) $\{\mu(t) \mid t \in H \text{ and } \mu(t) \subset (I \cup B) \times I \times T\}$.
   - If $R$ is ASK then $\mathrm{result}(R, \mu)$ is *false* if $\mu = \emptyset$ and *true* otherwise.
(2) *Semantics of dataset clauses.* Let $F$ be a set of dataset clauses. The dataset resulting from $F$, denoted $\mathrm{dataset}(F)$, contains:
   (i) a default graph consisting of the merge of the graphs referred in clauses FROM $u$. If there is no FROM $u$, then the default graph is an empty graph $G_0 = \emptyset$; and
   (ii) a named graph $\langle u, \mathrm{graph}(u) \rangle$ for each dataset clause "FROM NAMED $u$".
(3) *Semantics of filter constraints.* Let $\mu$ be a mapping and $C$ be a filter constraint. We say that $\mu$ satisfies $C$, denoted $\mu \models C$, if:
   - $C$ is $?X = v$, $?X \in \mathrm{dom}(\mu)$, and $\mu(?X) = v$;
   - $C$ is $?X = ?Y$, $?X \in \mathrm{dom}(\mu)$, $?Y \in \mathrm{dom}(\mu)$, and $\mu(?X) = \mu(?Y)$;
   - $C$ is $\mathrm{bound}(?X)$ and $?X \in \mathrm{dom}(\mu)$;
   - $C$ is $(\neg C_1)$ and it is not the case that $\mu \models C_1$;
   - $C$ is $(C_1 \vee C_2)$ and $\mu \models C_1$ or $\mu \models C_2$;
   - $C$ is $(C_1 \wedge C_2)$, $\mu \models C_1$ and $\mu \models C_2$.

- $C$ is $(u\ \theta\ \mathrm{SOME}(Q_{?X}))$ and there exists a mapping $\mu' \in \mathrm{ans}(\mu(Q_{?X}))$ satisfying that either $u\ \theta\ \mu'(?X)$ when $u \in I \cup L$ or $\mu(u)\ \theta\ \mu'(?X)$ when $u \in V$.
- $C$ is $(u\ \theta\ \mathrm{ALL}(Q_{?X}))$ and for every mapping $\mu' \in \mathrm{ans}(\mu(Q_{?X}))$ it holds that either $u\ \theta\ \mu'(?X)$ when $u \in I \cup L$ or $\mu(u)\ \theta\ \mu'(?X)$ when $u \in V$.
- $C$ is $(u\ \mathrm{IN}\ (Q_{?X}))$ and there exists a mapping $\mu' \in \mathrm{ans}(\mu(Q_{?X}))$ satisfying that either $u\ \theta\ \mu'(?X)$ when $u \in I \cup L$ or $\mu(u)\ \theta\ \mu'(?X)$ when $u \in V$.
- $C$ is $\mathrm{EXISTS}(Q_A)$ and $\mathrm{ans}(\mu(Q_A))$ is *true*.
- $C$ is $(Q_1\ \mathrm{CONTAINS}\ Q_2)$ and for every mapping $\mu_2 \in \mathrm{ans}(\mu(Q_2))$ there is a mapping $\mu_1 \in \mathrm{ans}(\mu(Q_1))$ such that $\mu_1$ and $\mu_2$ are compatible.
- $C$ is $(Q_1 = Q_2)$ and it holds that both $(\mu(Q_1)\ \mathrm{CONTAINS}\ \mu(Q_2))$ and $(\mu(Q_2)\ \mathrm{CONTAINS}\ \mu(Q_1))$ are true.

(4) *Semantics of graph patterns.* The *evaluation* of a graph pattern $P$ over a dataset $D$ with active graph $G$, denoted $[\![\cdot]\!]_G^D$, is defined recursively as follows:
- $P$ is a triple pattern then $[\![P]\!]_G^D = \{\mu \mid \mathrm{dom}(\mu) = \mathrm{var}(P)$ and $\mu(P) \subseteq G\}$
- $[\![(P_1\ \mathrm{AND}\ P_2)]\!]_G^D = [\![P_1]\!]_G^D \bowtie [\![P_2]\!]_G^D$.
- $[\![(P_1\ \mathrm{OPT}\ P_2)]\!]_G^D = [\![P_1]\!]_G^D \mathbin{⟕} [\![P_2]\!]_G^D$.
- $[\![(P_1\ \mathrm{UNION}\ P_2)]\!]_G^D = [\![P_1]\!]_G^D \cup [\![P_2]\!]_G^D$.
- $[\![(P_1\ \mathrm{MINUS}\ P_2)]\!]_G^D = [\![P_1]\!]_G^D \setminus [\![P_2]\!]_G^D$.
- If $u \in I$ then $[\![(u\ \mathrm{GRAPH}\ P_1)]\!]_G^D = [\![P_1]\!]_{\mathrm{graph}(u)_D}^D$.
- If $?X \in V$ and $\mu_{?X \to v}$ is a mapping such that $\mathrm{dom}(\mu) = \{?X\}$ and $\mu(?X) = v$, then
  $[\![(?X\ \mathrm{GRAPH}\ P_1)]\!]_G^D = \bigcup_{v\ \in\ \mathrm{names}(D)}([\![P_1]\!]_{\mathrm{graph}(v)_D}^D \bowtie \{\mu_{?X \to v}\})$.
- $[\![(P_1\ \mathrm{FILTER}\ C)]\!]_G^D = \{\mu \mid \mu \in [\![P_1]\!]_G^D$ and $\mu \models C\}$
- If $P$ is a SELECT query $Q_S$ then $[\![P]\!]_G^D = \mathrm{ans}(Q_S)$.

**Definition 2 (Answer for a query).** *Let $Q = (R, F, P)$ be a query, $D$ be the dataset obtained from $F$, and $G$ be the default graph of $D$. The* answer *to $Q$, denoted $\mathrm{ans}(Q)$, is defined as follows:*

- *if $R$ is SELECT $W$ then $\mathrm{ans}(Q) = \{\mathrm{result}(R, \mu) \mid \mu \in [\![P]\!]_G^D\}$.*
- *if $R$ is CONSTRUCT $H$ and $\mathrm{blank}(H)$ is the set of blank nodes appearing in $H$, then $\mathrm{ans}(Q) = \{\beta_i(\mathrm{result}(R, \mu_i)) \mid \mu_i \in [\![P]\!]_G^D\}$ where $\beta_i : \mathrm{blank}(H) \to (B \setminus \mathrm{blank}(H)$ is a blank renaming function satisfying that for each pair of mappings $\mu_j, \mu_k \in [\![P]\!]_G^D$, $\mathrm{range}(\beta_j) \cap \mathrm{range}(\beta_k) = \emptyset$.*
- *if $R$ is ASK then $\mathrm{ans}(Q) = $ false when $[\![P]\!]_G^D = \emptyset$ (i.e., there exists no mapping $\mu \in [\![P]\!]_G^D$) and $\mathrm{ans}(Q) = $ true otherwise.*

Note that, similar to SQL, the proposed language supports correlated queries, that is, variables from an outer query block can be accessed inside a nested query block [2]. Such variables are called *correlated variables*. Moreover, the semantics for nested queries defined in this paper follows the *nested iteration method* [8], i.e., the results of the inner query may have to be retrieved once for each result of the outer query.

# 3 Examples of nested queries in SPARQL

Let $G_1, G_2$ be two RDF graphs identified by IRIs `foaf` and `bib` respectively. $G_1$ contains personal information using the FOAF vocabulary [7]. $G_2$ contains bibliographic information using the bibTex Vocabulary [8]. Consider the following examples of nested queries.

*Example 1.* The oldest people.

```
SELECT ?Per1 FROM foaf
WHERE ((?Per1 foaf:age ?Age1)
          FILTER (¬(?Age1 < SOME ( SELECT ?Age2 FROM foaf
                                        WHERE (?Per2 foaf:age ?Age2)))))
```

*Example 2.* The youngest people.

```
SELECT ?Per1 FROM foaf
WHERE ((?Per1 foaf:age ?Age1)
          FILTER ( ?Age1 ≤ ALL ( SELECT ?Age2 FROM foaf
                                        WHERE (?Per2 foaf:age ?Age2))))
```

*Example 3.* Mails of people being part of at least one group.

```
SELECT ?Mail FROM foaf
WHERE ((?Per foaf:mbox ?Mail)
          FILTER (?Per IN ( SELECT ?Mem FROM foaf
                                  WHERE (?Mem foaf:member ?Group))))
```

*Example 4.* Mails of people having at least one publication.

```
SELECT ?Mail FROM foaf
WHERE ((?Per foaf:mbox ?Mail)
          FILTER ( EXISTS (ASK FROM bib
                                  WHERE (?Art bib:has-author ?Per)))))
```

*Example 5.* Titles of articles whose authors work in the same group.

```
SELECT ?Tit FROM foaf FROM bib
WHERE (((((?Art bib:title ?Tit) AND (?Art bib:has-author ?Aut1)) AND
  (?Aut1 foaf:member ?Group)) FILTER
    ((SELECT ?Mem FROM foaf WHERE (?Mem foaf:member ?Group))
     CONTAINS
    (SELECT ?Aut2 FROM bib WHERE (?Art bib:has-author ?Aut2))))
```

*Example 6.* Mails of people having no publications.

```
SELECT ?Mail FROM foaf WHERE
((?Per foaf:mbox ?Mail) MINUS (SELECT ?Per FROM bib
                                        WHERE (?Art bib:has-author ?Per)))
```

---

[7] http://xmlns.com/foaf/spec/

[8] http://zeitkunst.org/bibtex/0.1/

From the above examples we can say that:

– IN expressions are less expressive than SOME expressions because the former are restricted to equality of values, whereas the latter allows all scalar comparison operators.
– Nested queries with SOME / ALL operators without correlated variables are better for query composition, i.e., simple and direct copy/paste of queries.
– The use of EXISTS is not adequate for distributed queries because it needs correlated variables to make sense. This helps the user to express complex queries but makes the evaluation harder (because the application of the nested iteration method).
– The inclusion of SELECT queries as graph patterns allows operations among queries (which is not supported in SPARQL). This approach of nesting is also good for query composition.

## 4 Equivalences among nested queries

In this section we present transformations among nested queries. We will show that all types of nested queries can be simulated by the EXISTS operator and SELECT queries as graph patterns. Several equivalences presented in this section are well know in SQL [3].

### 4.1 Normalization

In order to simplify the transformations, we will avoid complex filter constraints, i.e., expressions of the form $C_1 \wedge C_2$ and $C_1 \vee C_2$ where $C_1$ and $C_2$ are filter constraints. This assumption is supported by the following lemma.

**Lemma 1.** *Every graph pattern having complex filter constraints can be transformed in a graph pattern without complex filter constraints* [9].

*Proof.* Let $P$ be a graph pattern and $C$, $C_1$, $C_2$ be filter constraints. The lemma is supported by the following equivalences:

$$(P \, \text{FILTER}(C_1 \wedge C_2)) \equiv ((P \, \text{FILTER} \, C_1) \, \text{FILTER} \, C_2) \tag{1}$$

$$(P \, \text{FILTER}(C_1 \vee C_2)) \equiv ((P \, \text{FILTER} \, C_1) \, \text{UNION}(P \, \text{FILTER} \, C_2)) \tag{2}$$

$$(P \, \text{FILTER}(\neg C)) \equiv (P \, \text{MINUS}(P \, \text{FILTER} \, C)) \tag{3}$$

Is not hard to see that the equivalences holds.

---

[9] Lemma 1 is true under set semantics. The inclusion of bag semantics, as defined for SPARQL, introduces complexity issues which are not discussed here.

## 4.2 Transformations

Consider the following definition of query equivalence.

**Definition 3 (Equivalence of queries).** *Two graph patterns $P_1$ and $P_2$ are equivalent, denoted $P_1 \equiv P_2$, if and only if $[\![P_1]\!]_G^D = [\![P_2]\!]_G^D$ for every RDF dataset $D$ with active graph $G$. Additionally, given two queries $Q = (R, F, P)$ and $Q' = (R, F, P')$, we say that $Q$ and $Q'$ are equivalent, denoted $Q \equiv Q'$, if and only if $P \equiv P'$.*

Next we will define transformations among several types of nested queries. Based on transformations defined in Section 4.1, we assume that queries do not contain complex filter constraints.

**Proposition 1 (Transforming IN queries).** *Let $P$ be a pattern of the form $(P_1 \text{ FILTER}(u \text{ IN } \{Q_2\}))$ where $u \in T$. Then, $P$ is equivalent to expression:*

$$(P_1 \text{ FILTER}(u = \text{SOME}\{Q_2\})) \tag{4}$$

**Proposition 2 (Transforming SOME queries).** *Let $P$ be a pattern of the form $(P_1 \text{ FILTER}(u \; \theta \; \text{SOME}(Q_2)))$ where $u \in T$, $Q_2 = (\text{SELECT } ?X_2, F_2, P_2)$, and $\hat{\theta}$ is the inverse operator to $\theta$. Then, $P$ is equivalent to the following expressions:*

$$(P_1 \text{ FILTER}(\neg(u \; \hat{\theta} \; \text{ALL}(Q_2)))) \tag{5}$$

$$(P_1 \text{ FILTER EXISTS}(\text{ASK}, F_2, (P_2 \text{ FILTER}(u \; \theta \; ?X_2)))) \tag{6}$$

**Proposition 3 (Transforming ALL queries).** *Let $P$ be a pattern of the form $(P_1 \text{ FILTER}(u \; \theta \; \text{ALL}(Q_2)))$ where $u \in T$, $Q_2 = (\text{SELECT } ?X_2, F_2, P_2)$, and $\hat{\theta}$ is the inverse operator to $\theta$. Then, $P$ is equivalent to the following expressions:*

$$(P_1 \text{ FILTER}(\neg(u \; \hat{\theta} \; \text{SOME}(Q_2)))) \tag{7}$$

$$(P_1 \text{ FILTER}(\neg \text{EXISTS}(\text{ASK}, F_2, (P_2 \text{ FILTER}(u \; \hat{\theta} \; ?X_2))))) \tag{8}$$

For example, the following queries show the application of transformations (7) and (8) to the query of Example 2.

*Example 7.* The youngest people (using the SOME operator).

```
SELECT ?Per1 FROM foaf
WHERE ((?Per1 foaf:age ?Age1)
          FILTER ( ¬(?Age1 > SOME (SELECT ?Age2 FROM foaf
                                    WHERE (?Per1 foaf:age ?Age2)))))
```

*Example 8.* The youngest people (using the EXISTS operator).

```
SELECT ?Per1 FROM foaf
WHERE ((?Per1 foaf:age ?Age1)
          FILTER (¬ EXISTS (ASK FROM foaf
                              WHERE ((?Per2 foaf:age ?Age2)
                                      FILTER (?Age1 > ?Age2)))))))
```

**Proposition 4 (Transforming** CONTAINS **queries).** *Let $P_1$ be a graph pattern, and $Q_1, Q_2$ be queries. The following equivalence holds*

$$(P_1 \text{ FILTER}(Q_1 \text{ CONTAINS } Q_2)) \equiv$$
$$(P_1 \text{ FILTER}(\neg \text{EXISTS}(\text{ASK}((Q_2) \text{ MINUS}(Q_1))))) \quad (9)$$

**Proposition 5 (Transforming equality of nested queries).** *Let $P_1$ be a graph pattern, and $Q_1, Q_2$ be queries. The following equivalence holds*

$$(P_1 \text{ FILTER}(Q_1 = Q_2)) \equiv$$
$$(P_1 \text{ FILTER}((\neg \text{EXISTS}(\text{ASK}((Q_1) \text{ MINUS}(Q_2))))$$
$$\wedge (\neg \text{EXISTS}(\text{ASK}((Q_2) \text{ MINUS}(Q_1)))))) \quad (10)$$

From the transformations defined above we can present the following result.

**Theorem 1.** *Filter constraints using* SOME*,* ALL*,* IN*,* CONTAINS *and set-equality can be simulated by using* EXISTS *and* SELECT *queries as graph patterns.*

## 5 Conclusions

We have studied how to extend SPARQL to support nesting along the design philosophy of SQL. We showed that there is a simple syntax and semantics for such extensions in SPARQL.

We have shown that incorporating ASK queries in FILTERS (through the EXISTS operator) and SELECT queries as graph patterns, gives the full power and flexibility of SQL nesting, allowing additionally to extend the semantics of SPARQL in a clean and modular form.

The proposal presented here permits a simple and direct implementation of the full power of nesting as known in the relational world (and hence by known translation results) in the SPARQL world.

**Future work.** An interesting problem studied in the database literature is the efficient implementation of nested queries, where a well known approach is the development of algorithms which transform nested queries into equivalent non-nested queries which can be processed more efficiently by query-processing subsystems [8,5]. On this line, most results are concentrated on aggregate subqueries; optimization of non-aggregate subqueries has some limitations, specially for queries with multiple subqueries and null values [2].

Although decorrelation often results in cheaper non-nested plans, decorrelation is not always applicable, and even if applicable may not be the best choice in all situations since decorrelation carries a materialization overhead [15,6]. In this direction, the issue of efficient methods of processing nested queries is one of the main problems to be addressed in future works on this topic.

# References

1. R. Angles and C. Gutierrez. The Expressive Power of SPARQL. In *Proceedings of the 7th International Semantic Web Conference (ISWC)*, number 5318 in LNCS, pages 114–129, 2008.
2. B. Cao and A. Badia. A nested relational approach to processing sql subqueries. In *Proceedings of the 2005 ACM SIGMOD international conference on Management of data*, pages 191–202, New York, NY, USA, 2005. ACM Press.
3. S. Ceri and G. Gottlob. Translating sql into relational algebra: optimization, semantics, and equivalence of sql queries. *IEEE Transactions on Software Engineering*, 11(4):324–345, 1985.
4. R. Cyganiak. A relational algebra for SPARQL. Technical Report HPL-2005-170, HP Labs, 2005.
5. R. A. Ganski and H. K. T. Wong. Optimization of nested sql queries revisited. In *Proceedings of the 1987 ACM SIGMOD international conference on Management of data*, pages 23–33, New York, NY, USA, 1987. ACM Press.
6. R. Guravannavar, H. S. Ramanujam, and S. Sudarshan. Optimizing nested queries with parameter sort orders. In *Proc. of the 31st Int. Conf. on Very large Data Bases (VLDB)*, pages 481–492. VLDB Endowment, 2005.
7. S. Harris and A. Seaborne. SPARQL 1.1 Query. W3C Working Draft. http://www.w3.org/TR/2009/WD-sparql11-query-20091022/, October 22 2009.
8. W. Kim. On optimizing an sql-like nested query. *ACM Transactions on Database Systems (TODS)*, 7(3):443–469, 1982.
9. K. Kjernsmo and A. Passant. SPARQL New Features and Rationale. W3C Working Draft. http://www.w3.org/TR/2009/WD-sparql-features-20090702/, July 2 2009.
10. G. Klyne and J. Carroll. Resource Description Framework (RDF) Concepts and Abstract Syntax. http://www.w3.org/TR/2004/REC-115-concepts-20040210/, February 2004.
11. J. Pérez, M. Arenas, and C. Gutierrez. Semantics and Complexity of SPARQL. In *Proceedings of the 5th International Semantic Web Conference (ISWC)*, number 4273 in LNCS, pages 30–43. Springer-Verlag, 2006.
12. A. Polleres. From SPARQL to Rules (and back). In *Proceedings of the 16th International World Wide Web Conference (WWW)*, pages 787–796. ACM, 2007.
13. E. Prud'hommeaux and A. Seaborne. SPARQL Query Language for RDF. W3C Recommendation 15 January. http://www.w3.org/TR/2008/REC-115-sparql-query-20080115/, 2008.
14. S. Schenk. A SPARQL Semantics Based on Datalog. In *30th Annual German Conference on Advances in Artificial Intelligence (KI)*, volume 4667 of *LNCS*, pages 160–174. Springer, 2007.
15. P. Seshadri, H. Pirahesh, and T. Y. C. Leung. Complex query decorrelation. In *Proc. of the 12th Int. Conf. on Data Engineering (ICDE)*, pages 450–458. IEEE Computer Society, 1996.