

# Answering SPARQL Queries over Databases under OWL 2 QL Entailment Regime

Roman Kontchakov<sup>1</sup>, Martin Rezk<sup>2</sup>, Mariano Rodríguez-Muro<sup>3</sup>, Guohui Xiao<sup>2</sup>, and Michael Zakharyashev<sup>1</sup>

<sup>1</sup> Department of Computer Science and Information Systems,  
Birkbeck, University of London, U.K.

<sup>2</sup> Faculty of Computer Science, Free University of Bozen-Bolzano, Italy

<sup>3</sup> IBM T.J. Watson Research Center, Yorktown Heights, NY, USA

**Abstract.** We present an extension of the ontology-based data access platform *Ontop* that supports answering SPARQL queries under the OWL 2 QL direct semantics entailment regime for data instances stored in relational databases. On the theoretical side, we show how any input SPARQL query, OWL 2 QL ontology and R2RML mappings can be rewritten to an equivalent SQL query solely over the data. On the practical side, we present initial experimental results demonstrating that by applying the *Ontop* technologies—the tree-witness query rewriting,  $\mathcal{T}$ -mappings compiling R2RML mappings with ontology hierarchies, and  $\mathcal{T}$ -mapping optimisations using SQL expressivity and database integrity constraints—the system produces scalable SQL queries.

## 1 Introduction

Ontology-based data access and management (OBDA) is a popular paradigm of organising access to various types of data sources that has been developed since the mid 2000s [11,17,24]. In a nutshell, OBDA separates the user from the data sources (relational databases, triple stores, etc.) by means of an ontology which provides the user with a convenient query vocabulary, hides the structure of the data sources, and can enrich incomplete data with background knowledge. About a dozen OBDA systems have been implemented in both academia and industry; e.g., [27,30,24,4,23,15,12,8,20,22]. Most of them support conjunctive queries and the OWL 2 QL profile of OWL 2 as the ontology language (or its generalisations to existential datalog rules). Thus, the OBDA platform *Ontop* [29] was designed to query data instances stored in relational databases, with the vocabularies of the data and OWL 2 QL ontologies linked by means of global-as-view (GAV) mappings. Given a conjunctive query in the vocabulary of such an ontology, *Ontop* rewrites it to an SQL query in the vocabulary of the data, optimises the rewriting and delegates its evaluation to the database system.

One of the main aims behind the newly designed query language SPARQL 1.1—a W3C recommendation since 2013—has been to support various entailment regimes, which can be regarded as variants of OBDA. Thus, the OWL 2 direct semantics entailment regime allows SPARQL queries over OWL 2 DL ontologies and RDF graphs (which can be thought of as 3-column database tables). SPARQL queries are in many aspects more expressive than conjunctive queries as they offer more complex query

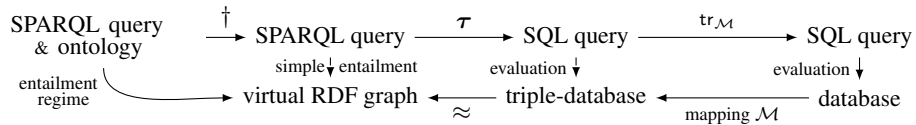
constructs and can retrieve not only domain elements but also class and property names using second-order variables. (Note, however, that SPARQL 1.1 does not cover all conjunctive queries.) OWL 2 DL is also vastly superior to OWL 2 QL, but this makes query answering under the OWL 2 direct semantics entailment regime intractable (CONP-hard for data complexity). For example, the query evaluation algorithm of [19] calls an OWL 2 DL reasoner for each possible assignment to the variables in a given query, and therefore cannot cope with large data instances.

In this paper, we investigate answering SPARQL queries under a less expressive entailment regime, which corresponds to OWL 2 QL, assuming that data is stored in relational databases. It is to be noted that the W3C specification<sup>1</sup> of SPARQL 1.1 defines entailment regimes for the profiles of OWL 2 by restricting the general definition to the profile constructs that can be used in the queries. However, in the case of OWL 2 QL, this generic approach leads to a sub-optimal, almost trivial query language, which is essentially subsumed by the OWL 2 RL entailment regime.

The first aim of this paper is to give an optimal definition of the OWL 2 QL direct semantics entailment regime and prove that—similarly to OBDA with OWL 2 QL and conjunctive queries—answering SPARQL queries under this regime is reducible to answering queries under *simple entailment*. More precisely, in Theorem 4 we construct a rewriting  $\cdot^\dagger$  of any given SPARQL query and ontology under the OWL 2 QL entailment regime to a SPARQL query that can be evaluated on any dataset directly.

In a typical *Ontop* scenario, data is stored in a relational database whose schema is linked to the vocabulary of the given OWL 2 QL ontology via a GAV mapping in the language R2RML. The mapping allows one to transform the relational data instance into an RDF representation, called the virtual RDF graph (which is not materialised in our scenario). The rewriting  $\cdot^\dagger$  constructs a SPARQL query over this virtual graph.

Our second aim is to show how such a SPARQL query can be translated to an equivalent SQL query over a relational representation of the virtual RDF graph as a 3-column table (translation  $\tau$  in Theorem 7). The third aim is to show that the resulting SQL query can be unfolded, using a given R2RML mapping  $\mathcal{M}$ , to an SQL query over the original database ( $\text{tr}_{\mathcal{M}}$  in Theorem 12), which is evaluated by the database system.



Unfortunately, each of these three transformations may involve an exponential blowup. We tackle this problem in *Ontop* using the following optimisation techniques. (i) The mapping is compiled with the ontology into a  $\mathcal{T}$ -mapping [29] and optimised by database dependencies (e.g., primary, candidate and foreign keys) and SQL disjunctions. (ii) The SPARQL-to-SQL translation is optimised using null join elimination (Theorem 8). (iii) The unfolding is optimised by eliminating joins with mismatching R2RML IRI templates, de-IRIing the join conditions (Section 3.3) and using database dependencies.

Our contributions (Theorems 4, 7, 8 and 12 and optimisations in Section 3.3) make *Ontop* the first system to support the W3C recommendations OWL 2 QL, R2RML, SPARQL and the OWL 2 QL direct semantics entailment regime; its architecture is out-

<sup>1</sup> <http://www.w3.org/TR/sparql11-entailment>

lined in Section 4. We evaluate the performance of *Ontop* using the LUBM Benchmark [16] extended with queries containing class and property variables, and compare it with two other systems that support the OWL 2 entailment regime by calling OWL DL reasoners (Section 5). Our experiments show that *Ontop* outperforms the reasoner-based systems for most of the queries over small datasets; over larger datasets the difference becomes dramatic, with *Ontop* demonstrating a solid performance even on 69 million triples in LUBM<sub>500</sub>. Finally, we note that, although *Ontop* was designed to work with existing relational databases, it is also applicable in the context of RDF triple stores, in which case approaches such as the one from [3] can be used to generate suitable relational schemas. Omitted proofs and evaluation details can be found in the full version at <http://www.dcs.bbk.ac.uk/~michael/ISWC-14-v2.pdf>.

## 2 SPARQL Queries under OWL 2 QL Entailment Regime

SPARQL is a W3C standard language designed to query RDF graphs. Its vocabulary contains four pairwise disjoint and countably infinite sets of symbols: **I** for *IRIs*, **B** for *blank nodes*, **L** for *RDF literals*, and **V** for *variables*. The elements of  $\mathbf{C} = \mathbf{I} \cup \mathbf{B} \cup \mathbf{L}$  are called *RDF terms*. A *triple pattern* is an element of  $(\mathbf{C} \cup \mathbf{V}) \times (\mathbf{I} \cup \mathbf{V}) \times (\mathbf{C} \cup \mathbf{V})$ . A *basic graph pattern (BGP)* is a finite set of triple patterns. Finally, a *graph pattern, P*, is an expression defined by the grammar

$$P ::= \text{BGP} \mid \text{FILTER}(P, F) \mid \text{BIND}(P, v, c) \mid \text{UNION}(P_1, P_2) \mid \\ \text{JOIN}(P_1, P_2) \mid \text{OPT}(P_1, P_2, F),$$

where  $F$ , a *filter*, is a formula constructed from atoms of the form  $\text{bound}(v)$ ,  $(v = c)$ ,  $(v = v')$ , for  $v, v' \in \mathbf{V}$ ,  $c \in \mathbf{C}$ , and possibly other built-in predicates using the logical connectives  $\wedge$  and  $\neg$ . The set of variables in  $P$  is denoted by  $\text{var}(P)$ .

A *SPARQL query* is a graph pattern  $P$  with a *solution modifier*, which specifies the *answer variables*—the variables in  $P$  whose values we are interested in—and the form of the output (we ignore other solution modifiers for simplicity). The values to variables are given by *solution mappings*, which are *partial maps*  $s: \mathbf{V} \rightarrow \mathbf{C}$  with (possibly empty) domain  $\text{dom}(s)$ . In this paper, we use the set-based (rather than bag-based, as in the specification) semantics for SPARQL. For sets  $S_1$  and  $S_2$  of solution mappings, a filter  $F$ , a variable  $v \in \mathbf{V}$  and a term  $c \in \mathbf{C}$ , let

- $\text{FILTER}(S, F) = \{s \in S \mid F^s = \top\}$ ;
- $\text{BIND}(S, v, c) = \{s \oplus \{v \mapsto c\} \mid s \in S\}$  (provided that  $v \notin \text{dom}(s)$ , for  $s \in S$ );
- $\text{UNION}(S_1, S_2) = \{s \mid s \in S_1 \text{ or } s \in S_2\}$ ;
- $\text{JOIN}(S_1, S_2) = \{s_1 \oplus s_2 \mid s_1 \in S_1 \text{ and } s_2 \in S_2 \text{ are compatible}\}$ ;
- $\text{OPT}(S_1, S_2, F) = \text{FILTER}(\text{JOIN}(S_1, S_2), F) \cup \{s_1 \in S_1 \mid \text{for all } s_2 \in S_2, \\ \text{either } s_1, s_2 \text{ are incompatible or } F^{s_1 \oplus s_2} \neq \top\}$ .

Here,  $s_1$  and  $s_2$  are *compatible* if  $s_1(v) = s_2(v)$ , for any  $v \in \text{dom}(s_1) \cap \text{dom}(s_2)$ , in which case  $s_1 \oplus s_2$  is a solution mapping with  $s_1 \oplus s_2: v \mapsto s_1(v)$ , for  $v \in \text{dom}(s_1)$ ,  $s_1 \oplus s_2: v \mapsto s_2(v)$ , for  $v \in \text{dom}(s_2)$ , and domain  $\text{dom}(s_1) \cup \text{dom}(s_2)$ . The *truth-value*  $F^s \in \{\top, \perp, \varepsilon\}$  of a filter  $F$  under a solution mapping  $s$  is defined inductively:

- $(\text{bound}(v))^s$  is  $\top$  if  $v \in \text{dom}(s)$  and  $\perp$  otherwise;
- $(v = c)^s = \varepsilon$  if  $v \notin \text{dom}(s)$ ; otherwise,  $(v = c)^s$  is the classical truth-value of the predicate  $s(v) = c$ ; similarly,  $(v = v')^s = \varepsilon$  if either  $v$  or  $v' \notin \text{dom}(s)$ ; otherwise,  $(v = v')^s$  is the classical truth-value of the predicate  $s(v) = s(v')$ ;
- $(\neg F)^s = \begin{cases} \varepsilon, & \text{if } F^s = \varepsilon, \\ \neg F^s, & \text{otherwise,} \end{cases}$  and  $(F_1 \wedge F_2)^s = \begin{cases} \perp, & \text{if } F_1^s = \perp \text{ or } F_2^s = \perp, \\ \top, & \text{if } F_1^s = F_2^s = \top, \\ \varepsilon, & \text{otherwise.} \end{cases}$

Finally, given an RDF graph  $G$ , the *answer to a graph pattern  $P$  over  $G$*  is the set  $\llbracket P \rrbracket_G$  of solution mappings defined by induction using the operations above and starting from the following base case: for a basic graph pattern  $B$ ,

$$\llbracket B \rrbracket_G = \{s : \text{var}(B) \rightarrow \mathbf{C} \mid s(B) \subseteq G\}, \quad (1)$$

where  $s(B)$  is the set of triples resulting from substituting each variable  $u$  in  $B$  by  $s(u)$ . This semantics is known as *simple entailment*.

*Remark 1.* The condition ‘ $F^{s_1 \oplus s_2}$  is not true’ in the definition of OPT is different from ‘ $F^{s_1 \oplus s_2}$  has an effective Boolean value of false’ given by the W3C specification:<sup>2</sup> the effective Boolean value can be undefined (type error) if a variable in  $F$  is not bound by  $s_1 \oplus s_2$ . As we shall see in Section 3.1, our reading corresponds to LEFT JOIN in SQL. (Note also that the informal explanation of OPT in the W3C specification is inconsistent with the definition of DIFF; see the full version for details.)

Under the *OWL 2 QL direct semantics entailment regime*, one can query an RDF graph  $G$  that consist of two parts: an *extensional* sub-graph  $\mathcal{A}$  representing the *data* as OWL 2 QL class and property assertions, and the *intensional* sub-graph  $\mathcal{T}$  representing the background *knowledge* as OWL 2 QL class and property axioms. We write  $(\mathcal{T}, \mathcal{A})$  in place of  $G$  to emphasise the partitioning. To illustrate, we give a simple example.

*Example 2.* Consider the following two axioms from the LUBM ontology  $(\mathcal{T}, \mathcal{A})$  (see Section 5), which are given here in the functional-style syntax (FSS):

SubClassOf(ub:UGStudent, ub:Student), SubClassOf(ub:GradStudent, ub:Student).

Under the entailment regime, we can write a query that retrieves all named *subclasses* of students in  $(\mathcal{T}, \mathcal{A})$  and all *instances* of each of these subclasses (cf.  $q'_9$  in Section 5):

SELECT  $?x ?C$  WHERE {  $?C$  rdfs:subClassOf ub:Student.  $?x$  rdf:type  $?C$ . }.

Here  $?C$  ranges over the class names (IRIs) in  $(\mathcal{T}, \mathcal{A})$  and  $?x$  over the IRIs of individuals. If, for example,  $\mathcal{A}$  consists of the two assertions on the left-hand side, then the answer to the query over  $(\mathcal{T}, \mathcal{A})$  is on the right-hand side:

$\mathcal{A}$	$?x$	$?C$
ClassAssertion(ub:UGStudent, ub:jim)	ub:jim	ub:UGStudent
ClassAssertion(ub:Student, ub:bob)	ub:jim	ub:Student
	ub:bob	ub:Student

<sup>2</sup> <http://www.w3.org/TR/sparql11-query/#sparqlAlgebra>

To formally define SPARQL queries that can be used under the OWL 2 QL direct semantics entailment regime, we assume that the set  $I$  of IRIs is partitioned into disjoint and countably infinite sets of *class names*  $I_C$ , *object property names*  $I_R$  and *individual names*  $I_I$ . Similarly, the variables  $V$  are also assumed to be a disjoint union of countably infinite sets  $V_C, V_R, V_I$ . Now, we define an *OWL 2 QL BGP* as a finite set of triple patterns representing OWL 2 QL axiom and assertion templates in the FSS such as:<sup>3</sup>

$$\begin{array}{ll} \text{SubClassOf}(\text{SubC}, \text{SuperC}), & \text{DisjointClasses}(\text{SubC}_1, \dots, \text{SubC}_n), \\ \text{ObjectPropertyDomain}(\text{OP}, \text{SuperC}), & \text{ObjectPropertyRange}(\text{OP}, \text{SuperC}), \\ \text{SubObjectPropertyOf}(\text{OP}, \text{OP}), & \text{DisjointObjectProperties}(\text{OP}_1, \dots, \text{OP}_n), \\ \text{ClassAssertion}(\text{SuperC}, I), & \text{ObjectPropertyAssertion}(\text{OP}, I, I), \end{array}$$

where  $I \in I_I \cup V_I$  and  $OP, \text{SubC}$  and  $\text{SuperC}$  are defined by the following grammar with  $C \in I_C \cup V_C$  and  $R \in I_R \cup V_R$ :

$$\begin{array}{ll} OP ::= R & | \text{ObjectInverseOf}(R), \\ \text{SubC} ::= C & | \text{ObjectSomeValuesFrom}(OP, \text{owl:Thing}), \\ \text{SuperC} ::= C & | \text{ObjectIntersectionOf}(\text{SuperC}_1, \dots, \text{SuperC}_n) \quad | \\ & \text{ObjectSomeValuesFrom}(OP, \text{SuperC}). \end{array}$$

*OWL 2 QL graph patterns* are constructed from OWL 2 QL BGPs using the SPARQL operators. Finally, an *OWL 2 QL query* is a pair  $(P, V)$ , where  $P$  is an OWL 2 QL graph pattern and  $V \subseteq \text{var}(P)$ . To define the answer to such a query  $(P, V)$  over an RDF graph  $(\mathcal{T}, \mathcal{A})$ , we fix a *finite vocabulary*  $I_{\mathcal{T}, \mathcal{A}} \subseteq I$  that includes all names (IRIs) in  $\mathcal{T}$  and  $\mathcal{A}$  as well as the required finite part of the OWL 2 RDF-based vocabulary (e.g.,  $\text{owl:Thing}$  but not the infinite number of the  $\text{rdf:}_n$ ). To ensure finiteness of the answers and proper typing of variables, in the following definition we only consider solution mappings  $s: \text{var}(P) \rightarrow I_{\mathcal{T}, \mathcal{A}}$  such that  $s^{-1}(I_\alpha) \subseteq V_\alpha$ , for  $\alpha \in \{C, R, I\}$ . For each BGP  $B$ , we define the *answer*  $\llbracket B \rrbracket_{\mathcal{T}, \mathcal{A}}$  to  $B$  over  $(\mathcal{T}, \mathcal{A})$  by taking

$$\llbracket B \rrbracket_{\mathcal{T}, \mathcal{A}} = \{s: \text{var}(B) \rightarrow I_{\mathcal{T}, \mathcal{A}} \mid (\mathcal{T}, \mathcal{A}) \models s(B)\},$$

where  $\models$  is the entailment relation given by the OWL 2 direct semantics. Starting from the  $\llbracket B \rrbracket_{\mathcal{T}, \mathcal{A}}$  and applying the SPARQL operators in  $P$ , we compute the set  $\llbracket P \rrbracket_{\mathcal{T}, \mathcal{A}}$  of *solution mappings*. The *answer to  $(P, V)$  over  $(\mathcal{T}, \mathcal{A})$*  is the restriction  $\llbracket P \rrbracket_{\mathcal{T}, \mathcal{A}}|_V$  of the solution mappings in  $\llbracket P \rrbracket_{\mathcal{T}, \mathcal{A}}$  to the variables in  $V$ .

*Example 3.* Suppose  $\mathcal{T}$  contains

$$\begin{array}{l} \text{SubClassOf}(\text{:A}, \text{ObjectSomeValuesFrom}(\text{:P}, \text{owl:Thing})), \\ \text{SubObjectPropertyOf}(\text{:P}, \text{:R}), \quad \text{SubObjectPropertyOf}(\text{:P}, \text{ObjectInverseOf}(\text{:S})). \end{array}$$

Consider the following OWL 2 QL BGP  $B$ :

$$\text{ClassAssertion}(\text{ObjectSomeValuesFrom}(\text{:R}, \text{ObjectSomeValuesFrom}(\text{:S}, \text{ObjectSomeValuesFrom}(\text{:T}, \text{owl:Thing}))), ?x).$$

<sup>3</sup> The official specification of legal queries under the OWL 2 QL entailment regime only allows  $\text{ClassAssertion}(C, I)$  rather than  $\text{ClassAssertion}(\text{SuperC}, I)$ , which makes the OWL 2 QL entailment regime trivial and essentially subsumed by the OWL 2 RL entailment regime.

Assuming that  $\mathcal{A} = \{\text{ClassAssertion}(:A, :a), \text{ObjectPropertyAssertion}(:T, :a, :b)\}$ , it is not hard to see that  $\llbracket B \rrbracket_{\mathcal{T}, \mathcal{A}} = \{?x \mapsto :a\}$ . Indeed, by the first assertion of  $\mathcal{A}$  and the first two axioms of  $\mathcal{T}$ , any model of  $(\mathcal{T}, \mathcal{A})$  contains a domain element  $w$  (not necessarily among the individuals in  $\mathcal{A}$ ) such that  $\text{ObjectPropertyAssertion}(:R, :a, w)$  holds. In addition, the third axiom of  $\mathcal{T}$  implies  $\text{ObjectPropertyAssertion}(:S, w, :a)$ , which together with the second assertion of  $\mathcal{A}$  mean that  $\{?x \mapsto :a\}$  is an answer.

The following theorem shows that answering OWL 2 QL queries under the direct semantics entailment regime can be reduced to answering OWL 2 QL queries under simple entailment, which are evaluated only on the extensional part of the RDF graph:

**Theorem 4.** *Given any intensional graph  $\mathcal{T}$  and OWL 2 QL query  $(P, V)$ , one can construct an OWL 2 QL query  $(P^\dagger, V)$  such that, for any extensional graph  $\mathcal{A}$  (in some fixed finite vocabulary),  $\llbracket P \rrbracket_{\mathcal{T}, \mathcal{A}}|_V = \llbracket P^\dagger \rrbracket_{\mathcal{A}}|_V$ .*

*Proof sketch.* By the definition of the entailment regime, it suffices to construct  $B^\dagger$ , for any BGP  $B$ ; the rewriting  $P^\dagger$  is obtained then by replacing each BGP  $B$  in  $P$  with  $B^\dagger$ . First, we instantiate the class and property variables in  $B$  by all possible class and property names in the given vocabulary and add the respective BIND operations. In each of the resulting BGPs, we remove the class and property axioms if they are entailed by  $\mathcal{T}$ ; otherwise we replace the BGP with an empty one. The obtained BGPs are (SPARQL representations of) conjunctive queries (with non-distinguished variables in complex concepts *SuperC* of the assertions  $\text{ClassAssertion}(\textit{SuperC}, I)$ ). The second step is to rewrite these conjunctive queries together with  $\mathcal{T}$  into unions of conjunctive queries (BGPs) that can be evaluated over any extensional graph  $\mathcal{A}$  [5,21]. (We emphasise that the SPARQL algebra operations, including difference and OPT, are applied to BGPs and do not interact with the two steps of our rewriting.)  $\square$

We illustrate the proof of Theorem 4 using the queries from Examples 2 and 3.

*Example 5.* The class variable  $?C$  in the query from Example 2 can be instantiated, using BIND, by all possible values from  $I_C \cap I_{\mathcal{T}, \mathcal{A}}$ , which gives the rewriting

```
SELECT ?x ?C WHERE {
  { ?x rdf:type ub:Student. BIND(ub:Student as ?C) } UNION
  { ?x rdf:type ub:GradStudent. BIND(ub:GradStudent as ?C) } UNION
  { ?x rdf:type ub:UGStudent. BIND(ub:UGStudent as ?C) } }.
```

The query from Example 3 is equivalent to a (tree-shaped) conjunctive query with three non-distinguished and one answer variable, which can be rewritten to

```
SELECT ?x WHERE { { ?x :R ?y. ?y :S ?z. ?z :T ?u. } UNION
  { ?x rdf:type :A. ?x :T ?u. } }.
```

### 3 Translating SPARQL under Simple Entailment to SQL

A number of translations of SPARQL queries (under simple entailment) to SQL queries have already been suggested in the literature; see, e.g., [9,13,7,32,27]. However, none

of them is suitable for our aims because they do not take into account the three-valued logic used in the OPTIONAL and BOUND constructs of the current SPARQL 1.1 (the semantics of OPTIONAL was not compositional in SPARQL 1.0). Note also that SPARQL has been translated to Datalog [25,2,26].

We begin by recapping the basics of relational algebra and SQL (see e.g., [1]). Let  $U$  be a finite (possibly empty) set of *attributes*. A *tuple over  $U$*  is a map  $t: U \rightarrow \Delta$ , where  $\Delta$  is the underlying domain, which always contains a distinguished element *null*. A ( $|U|$ -ary) *relation over  $U$*  is a finite set of tuples over  $U$  (again, we use the set-based rather than bag-based semantics). A *filter  $F$  over  $U$*  is a formula constructed from atoms  $isNull(U')$ ,  $(u = c)$  and  $(u = u')$ , where  $U' \subseteq U$ ,  $u, u' \in U$  and  $c \in \Delta$ , using the connectives  $\wedge$  and  $\neg$ . Let  $F$  be a filter with variables  $U$  and let  $t$  be a tuple over  $U$ . The *truth-value  $F^t \in \{\top, \perp, \varepsilon\}$  of  $F$  over  $t$*  is defined inductively:

- $(isNull(U'))^t$  is  $\top$  if  $t(u)$  is *null*, for all  $u \in U'$ , and  $\perp$  otherwise;
- $(u = c)^t = \varepsilon$  if  $t(u)$  is *null*; otherwise,  $(u = c)^t$  is the classical truth-value of the predicate  $t(u) = c$ ; similarly,  $(u = u')^t = \varepsilon$  if either  $t(u)$  or  $t(u')$  is *null*; otherwise,  $(u = u')^t$  is the classical truth-value of the predicate  $t(u) = t(u')$ ;
- $(\neg F)^t = \begin{cases} \varepsilon, & \text{if } F^t = \varepsilon, \\ \neg F^t, & \text{otherwise,} \end{cases}$  and  $(F_1 \wedge F_2)^t = \begin{cases} \perp, & \text{if } F_1^t = \perp \text{ or } F_2^t = \perp, \\ \top, & \text{if } F_1^t = F_2^t = \top, \\ \varepsilon, & \text{otherwise.} \end{cases}$

(Note that  $\neg$  and  $\wedge$  are interpreted in the same three-valued logic as in SPARQL.) We use standard relational algebra operations such as union, difference, projection, selection, renaming and natural (inner) join. Let  $R_i$  be a relation over  $U_i$ ,  $i = 1, 2$ .

- If  $U_1 = U_2$  then the standard  $R_1 \cup R_2$  and  $R_1 \setminus R_2$  are relations over  $U_1$ .
- If  $U \subseteq U_1$  then  $\pi_U R_1 = R_1|_U$  is a relation over  $U$ .
- If  $F$  is a filter over  $U_1$  then  $\sigma_F R_1 = \{t \in R_1 \mid F^t = \top\}$  is a relation over  $U_1$ .
- If  $v \notin U_1$  and  $u \in U_1$  then  $\rho_{v/u} R_1 = \{t_{v/u} \mid t \in R_1\}$ , where  $t_{v/u}: v \mapsto t(u)$  and  $t_{v/u}: u' \mapsto t(u')$ , for  $u' \in U_1 \setminus \{u\}$ , is a relation over  $(U_1 \setminus \{u\}) \cup \{v\}$ .
- $R_1 \bowtie R_2 = \{t_1 \oplus t_2 \mid t_1 \in R_1 \text{ and } t_2 \in R_2 \text{ are compatible}\}$  is a relation over  $U_1 \cup U_2$ . Here,  $t_1$  and  $t_2$  are *compatible* if  $t_1(u) = t_2(u) \neq \text{null}$ , for all  $u \in U_1 \cap U_2$ , in which case a tuple  $t_1 \oplus t_2$  over  $U_1 \cup U_2$  is defined by taking  $t_1 \oplus t_2: u \mapsto t_1(u)$ , for  $u \in U_1$ , and  $t_1 \oplus t_2: u \mapsto t_2(u)$ , for  $u \in U_2$  (note that if  $u$  is *null* in either of the tuples then they are incompatible).

To bridge the gap between partial functions (solution mappings) in SPARQL and total mappings (on attributes) in SQL, we require one more operation (expressible in SQL):

- If  $U \cap U_1 = \emptyset$  then the *padding*  $\mu_U R_1$  is  $R_1 \bowtie \text{null}^U$ , where  $\text{null}^U$  is the relation consisting of a single tuple  $t$  over  $U$  with  $t: u \mapsto \text{null}$ , for all  $u \in U$ .

By an *SQL query*,  $Q$ , we understand any expression constructed from relation symbols (each over a fixed set of attributes) and filters using the relational algebra operations given above (and complying with all restrictions on the structure). Suppose  $Q$  is an SQL query and  $D$  a data instance which, for any relation symbol in the schema under consideration, gives a concrete relation over the corresponding set of attributes. The

answer to  $Q$  over  $D$  is a relation  $\|Q\|_D$  defined inductively in the obvious way starting from the base case: for a relation symbol  $Q$ ,  $\|Q\|_D$  is the corresponding relation in  $D$ .

We now define a translation,  $\tau$ , which, given a graph pattern  $P$ , returns an SQL query  $\tau(P)$  with the same answers as  $P$ . More formally, for a set of variables  $V$ , let  $ext_V$  be a function transforming any solution mapping  $s$  with  $dom(s) \subseteq V$  to a tuple over  $V$  by padding it with *nulls*:

$$ext_V(s) = \{v \mapsto s(v) \mid v \in dom(s)\} \cup \{v \mapsto null \mid v \in V \setminus dom(s)\}.$$

The *relational answer to  $P$  over  $G$*  is  $\|P\|_G = \{ext_{var(P)}(s) \mid s \in \|P\|_G\}$ . The SQL query  $\tau(P)$  will be such that, for any RDF graph  $G$ , the relational answer to  $P$  over  $G$  coincides with the answer to  $\tau(P)$  over  $triple(G)$ , the database instance storing  $G$  as a ternary relation  $triple$  with the attributes  $subj, pred, obj$ . First, we define the translation of a SPARQL filter  $F$  by taking  $\tau(F)$  to be the SQL filter obtained by replacing each  $bound(v)$  with  $\neg isNull(v)$  (other built-in predicates can be handled similarly).

**Proposition 6.** *Let  $F$  be a SPARQL filter and let  $V$  be the set of variables in  $F$ . Then  $F^s = (\tau(F))^{ext_V(s)}$ , for any solution mapping  $s$  with  $dom(s) \subseteq V$ .*

The definition of  $\tau$  proceeds by induction on the construction of  $P$ . Note that we can always assume that graph patterns *under simple entailment* do not contain blank nodes because they can be replaced by fresh variables. It follows that a BGP  $\{tp_1, \dots, tp_n\}$  is equivalent to  $JOIN(\{tp_1\}, JOIN(\{tp_2\}, \dots))$ . So, for the basis of induction we set

$$\tau(\langle s, p, o \rangle) = \begin{cases} \pi_{\emptyset} \sigma_{(subj=s) \wedge (pred=p) \wedge (obj=o)} triple, & \text{if } s, p, o \in I \cup L, \\ \pi_s \rho_s / subj \sigma_{(pred=p) \wedge (obj=o)} triple, & \text{if } s \in V \text{ and } p, o \in I \cup L, \\ \pi_{s,o} \rho_s / subj \rho_o / obj \sigma_{pred=p} triple, & \text{if } s, o \in V, s \neq o, p \in I \cup L, \\ \pi_s \rho_s / subj \sigma_{(pred=p) \wedge (subj=obj)} triple, & \text{if } s, o \in V, s = o, p \in I \cup L, \\ \dots & \end{cases}$$

(the remaining cases are similar). Now, if  $P_1$  and  $P_2$  are graph patterns and  $F_1$  and  $F$  are filters containing only variables in  $var(P_1)$  and  $var(P_1) \cup var(P_2)$ , respectively, then we set  $U_i = var(P_i)$ ,  $i = 1, 2$ , and

$$\begin{aligned} \tau(\text{FILTER}(P_1, F_1)) &= \sigma_{\tau(F_1)} \tau(P_1), \\ \tau(\text{BIND}(P_1, v, c)) &= \tau(P_1) \bowtie \{v \mapsto c\}, \\ \tau(\text{UNION}(P_1, P_2)) &= \mu_{U_2 \setminus U_1} \tau(P_1) \cup \mu_{U_1 \setminus U_2} \tau(P_2), \\ \tau(\text{JOIN}(P_1, P_2)) &= \bigcup_{\substack{V_1, V_2 \subseteq U_1 \cap U_2 \\ V_1 \cap V_2 = \emptyset}} [(\pi_{U_1 \setminus V_1} \sigma_{isNull(V_1)} \tau(P_1)) \bowtie (\pi_{U_2 \setminus V_2} \sigma_{isNull(V_2)} \tau(P_2))], \\ \tau(\text{OPT}(P_1, P_2, F)) &= \tau(\text{FILTER}(\text{JOIN}(P_1, P_2), F)) \cup \\ &\quad \mu_{U_2 \setminus U_1} (\tau(P_1) \setminus \bigcup_{V_1 \subseteq U_1 \cap U_2} \mu_{V_1} \pi_{U_1 \setminus V_1} \tau(\text{FILTER}(\text{JOIN}(P_1^{V_1, U_1 \cap U_2}, P_2), F))), \end{aligned}$$

where  $P^{V,U} = \text{FILTER}(P, \bigwedge_{v \in V} \neg bound(v) \wedge \bigwedge_{v \in U \setminus V} bound(v))$ . It is readily seen that any  $\tau(P)$  is a valid SQL query and defines a relation over  $var(P)$ ; in particular,  $\tau(\text{JOIN}(P_1, P_2))$  is a relation over  $\bigcup_{i=1,2} (U_i \setminus V_i) = U_1 \cup U_2 = var(\text{JOIN}(P_1, P_2))$ .



**Theorem 7.** For any RDF graph  $G$  and any graph pattern  $P$ ,  $\|P\|_G = \|\tau(P)\|_{\text{triple}(G)}$ .

*Proof.* The proof is by induction on the structure of  $P$ . Here we only consider the induction step for  $P = \text{JOIN}(P_1, P_2)$ . Let  $U_i = \text{var}(P_i)$ ,  $i = 1, 2$ , and  $U = U_1 \cap U_2$ .

If  $t \in \|\text{JOIN}(P_1, P_2)\|_G$  then there is a solution mapping  $s \in \llbracket \text{JOIN}(P_1, P_2) \rrbracket_G$  with  $\text{ext}_{U_1 \cup U_2}(s) = t$ , and so there are  $s_i \in \llbracket P_i \rrbracket_G$  such that  $s_1$  and  $s_2$  are compatible and  $s_1 \oplus s_2 = s$ . Since,  $\text{ext}_{U_i}(s_i) \in \|P_i\|_G$ , by IH,  $\text{ext}_{U_i}(s_i) \in \|\tau(P_i)\|_{\text{triple}(G)}$ . Let  $V = \text{dom}(s_1) \cap \text{dom}(s_2)$ ,  $V_1 = U \setminus \text{dom}(s_1)$  and  $V_2 = \text{dom}(s_1) \cap (U \setminus \text{dom}(s_2))$ . Then  $V_1$ ,  $V_2$  and  $V$  are disjoint and partition  $U$ . Consider  $t_i = \text{ext}_{U_i \setminus V_i}(s_i)$ , and let  $Q_i = \pi_{U_i \setminus V_i}(\sigma_{\text{isNull}(V_i)} \tau(P_i))$ . By definition,  $t_i \in \|Q_i\|_{\text{triple}(G)}$ . Since  $s_1$  and  $s_2$  are compatible and  $V$  are the common non-null attributes of  $t_1$  and  $t_2$ , the tuples  $t_1$  and  $t_2$  are compatible and  $t_1 \oplus t_2 \in \|Q_1 \bowtie Q_2\|_{\text{triple}(G)}$ . It remains to observe that  $t = t_1 \oplus t_2$ .

If  $t \in \|\tau(\text{JOIN}(P_1, P_2))\|_{\text{triple}(G)}$  then there are disjoint  $V_1, V_2 \subseteq U$  and compatible tuples  $t_1$  and  $t_2$  with  $t = t_1 \oplus t_2$  and  $t_i \in \|\pi_{U_i \setminus V_i} \sigma_{\text{isNull}(V_i)} \tau(P_i)\|_{\text{triple}(G)}$ . Let  $s_i$  be the solution mapping such that  $t_i = \text{ext}_{U_i}(s_i)$ . Then  $s_1$  and  $s_2$  are compatible and  $\text{ext}_{U_i}(s_i) \in \|\tau(P_i)\|_{\text{triple}(G)}$ . By IH,  $\text{ext}_{U_i}(s_i) \in \|P_i\|_G$  and  $s_i \in \llbracket P_i \rrbracket_G$ . Therefore,  $s_1 \oplus s_2 \in \llbracket \text{JOIN}(P_1, P_2) \rrbracket_G$  and  $\text{ext}_{U_1 \cup U_2}(s_1 \oplus s_2) = t \in \|\text{JOIN}(P_1, P_2)\|_G$ .  $\square$

### 3.1 Optimising SPARQL JOIN and OPT

By definition,  $\tau(\text{JOIN}(P_1, P_2))$  is a union of exponentially many natural joins ( $\bowtie$ ). Observe, however, that for any BGP  $B = \{tp_1, \dots, tp_n\}$ , none of the attributes in the  $\tau(tp_i)$  can be *null*. So, we can drastically simplify the definition of  $\tau(B)$  by taking

$$\tau(\{tp_1, \dots, tp_n\}) = \tau(tp_1) \bowtie \dots \bowtie \tau(tp_n).$$

Moreover, this observation can be generalised. First, we identify the variables in graph patterns that are not necessarily bound in solution mappings:

$$\begin{aligned} \nu(B) &= \emptyset, & B \text{ is a BGP,} \\ \nu(\text{FILTER}(P_1, F)) &= \nu(P_1) \setminus \{v \mid \text{bound}(v) \text{ is a conjunct of } F\}, \\ \nu(\text{BIND}(P_1, v, c)) &= \nu(P_1), \\ \nu(\text{UNION}(P_1, P_2)) &= (\text{var}(P_1) \setminus \text{var}(P_2)) \cup (\text{var}(P_2) \setminus \text{var}(P_1)) \cup \nu(P_1) \cup \nu(P_2), \\ \nu(\text{JOIN}(P_1, P_2)) &= \nu(P_1) \cup \nu(P_2), \\ \nu(\text{OPT}(P_1, P_2, F)) &= \nu(P_1) \cup \text{var}(P_2). \end{aligned}$$

Thus, if a variable  $v$  in  $P$  does not belong to  $\nu(P)$ , then  $v \in \text{dom}(s)$ , for any solution mapping  $s \in \llbracket P \rrbracket_G$  and RDF graph  $G$  (but not the other way round). Now, we observe that the union in the definition of  $\tau(\text{JOIN}(P_1, P_2))$  can be taken over those subsets of  $\text{var}(P_1) \cap \text{var}(P_2)$  that only contain variables from  $\nu(P_1) \cup \nu(P_2)$ . This gives us:

**Theorem 8.** If  $\text{var}(P_1) \cap \text{var}(P_2) \cap (\nu(P_1) \cup \nu(P_2)) = \emptyset$  then we can define

$$\tau(\text{JOIN}(P_1, P_2)) = \tau(P_1) \bowtie \tau(P_2), \quad \tau(\text{OPT}(P_1, P_2, F)) = \tau(P_1) \bowtie_{\tau(F)} \tau(P_2),$$

where  $R_1 \bowtie_F R_2 = \sigma_F(R_1 \bowtie R_2) \cup \mu_{U_2 \setminus U_1}(R_1 \setminus \pi_{U_1}(\sigma_F(R_1 \bowtie R_2)))$ , for  $R_i$  over  $U_i$ .

(Note that the relational operation  $\bowtie_F$  corresponds to LEFT JOIN in SQL with the condition  $F$  placed in its ON clause.)

*Example 9.* Consider the following BGP  $B$  taken from the official SPARQL specification (‘find the names of people who do not know anyone’):

```
FILTER(OPT({ ?x foaf:givenName ?n }, { ?x foaf:knows ?w },  $\top$ ),  $\neg$ bound(?w)).
```

By Theorem 8,  $\tau(B)$  is defined as  $\sigma_{isNull(w)}(\pi_{x,n}Q_1 \bowtie \pi_{x,w}Q_2)$ , where  $Q_1$  and  $Q_2$  are  $\sigma_{pred=foaf:givenName} \rho_{x/subj} \rho_{n/obj}$  triple and  $\sigma_{pred=foaf:knows} \rho_{x/subj} \rho_{w/obj}$  triple, respectively (we note in passing that the projection on  $x$  is equivalent to  $\pi_x Q_1 \setminus \pi_x Q_2$ ).

### 3.2 R2RML Mappings

The SQL translation of a SPARQL query constructed above has to be evaluated over the ternary relation  $triple(G)$  representing the virtual RDF graph  $G$ . Our aim now is to transform it to an SQL query over the actual database, which is related to  $G$  by means of an R2RML mapping [10]. A variant of such a transformation has been suggested in [27]. Here we develop the idea first presented in [28]. We begin with a simple example.

*Example 10.* The following R2RML mapping (in the Turtle syntax) populates an object property `ub:UGDegreeFrom` from a relational table `students`, whose attributes `id` and `degreeuniid` identify graduate students and their universities:

```
_:m1 a rr:TripleMap;
      rr:logicalTable [ rr:sqlQuery "SELECT * FROM students WHERE stype=1" ];
      rr:subjectMap [ rr:template "/GradStudent{id}" ];
      rr:predicateObjectMap [ rr:predicate ub:UGDegreeFrom ;
                              rr:objectMap [ rr:template "/Uni{degreeuniid}" ] ]
```

More specifically, for each tuple in the query, an R2RML processor generates an RDF triple with the predicate `ub:UGDegreeFrom` and the subject and object constructed from attributes `id` and `degreeuniid`, respectively, using IRI templates.

Our aim now is as follows: given an R2RML mapping  $\mathcal{M}$ , we are going to define an SQL query  $tr_{\mathcal{M}}(triple)$  that constructs the relational representation  $triple(G_{D,\mathcal{M}})$  of the virtual RDF graph  $G_{D,\mathcal{M}}$  obtained by  $\mathcal{M}$  from any given data instance  $D$ . Without loss of generality and to simplify presentation, we assume that each triple map has

- one logical table (`rr:sqlQuery`),
- one subject map (`rr:subjectMap`), which does not have resource typing (`rr:class`),
- and one predicate-object map with one `rr:predicateMap` and one `rr:objectMap`.

This normal form can be achieved by introducing predicate-object maps with `rdf:type` and splitting any triple map into a number of triple maps with the same logical table and subject. We also assume that triple maps contain no referencing object maps (`rr:parentTriplesMap`, etc.) since they can be eliminated using joint SQL queries [10]. Finally, we assume that the term maps (i.e., subject, predicate and object maps) contain no constant shortcuts and are of the form `[rr:column  $v$ ]`, `[rr:constant  $c$ ]` or `[rr:template  $s$ ]`.

Given a triple map  $m$  with a logical table (SQL query)  $R$ , we construct a selection  $\sigma_{\neg isNull(v_1)} \cdots \sigma_{\neg isNull(v_k)} R$ , where  $v_1, \dots, v_k$  are the *referenced columns* of  $m$  (attributes of  $R$  in the term maps in  $m$ )—this is done to exclude tuples that contain *null* [10]. To construct  $tr_m$ , the selection filter is prefixed with projection  $\pi_{subj,pred,obj}$

and, for each of the three term maps, either with renaming (e.g., with  $\rho_{obj/v}$  if the object map is of the form  $[rr:column\ v]$ ) or with value creation (if the term map is of the form  $[rr:constant\ c]$  or  $[rr:template\ s]$ ; in the latter case, we use the built-in string concatenation function  $\|$ ). For instance, the mapping  $\_ :m1$  from Example 10 is converted to the SQL query

```
SELECT ('/GradStudent' \| id) AS subj, 'ub:UGDegreeFrom' AS pred,
      ('/Uni' \| degreeuniid) AS obj FROM students
WHERE (id IS NOT NULL) AND (degreeuniid IS NOT NULL) AND (stype=1).
```

Given an R2RML mapping  $\mathcal{M}$ , we set  $\text{tr}_{\mathcal{M}}(\text{triple}) = \bigcup_{m \in \mathcal{M}} \text{tr}_m$ .

**Proposition 11.** *For any R2RML mapping  $\mathcal{M}$  and data instance  $D$ ,  $t \in \|\text{tr}_{\mathcal{M}}(\text{triple})\|_D$  if and only if  $t \in \text{triple}(G_{D,\mathcal{M}})$ .*

Finally, given a graph pattern  $P$  and an R2RML mapping  $\mathcal{M}$ , we define  $\text{tr}_{\mathcal{M}}(\tau(P))$  to be the result of replacing every occurrence of the relation  $\text{triple}$  in the query  $\tau(P)$ , constructed in Section 3, with  $\text{tr}_{\mathcal{M}}(\text{triple})$ . By Theorem 7 and Proposition 11, we obtain:

**Theorem 12.** *For any graph pattern  $P$ , R2RML mapping  $\mathcal{M}$  and data instance  $D$ ,  $\|P\|_{G_{D,\mathcal{M}}} = \|\text{tr}_{\mathcal{M}}(\tau(P))\|_D$ .*

### 3.3 Optimising SQL Translation

The straightforward application of  $\text{tr}_{\mathcal{M}}$  to  $\tau(P)$  can result in a very complex SQL query. We now show that such queries can be optimised by the following techniques:

- choosing matching  $\text{tr}_m$  from  $\text{tr}_{\mathcal{M}}(\text{triple})$ , for each occurrence of  $\text{triple}$  in  $\tau(P)$ ;
- using the distributivity of  $\bowtie$  over  $\cup$  and removing sub-queries with *incompatible IRI templates* and *de-IRIing* join conditions;
- functional dependencies (e.g., primary keys) for self-join elimination [6,18,29,30].

To illustrate, suppose we are given a mapping  $\mathcal{M}$  containing  $\_ :m1$  from Example 10 and the following triple maps (which are a simplified version of those in Section 5):

```
\_ :m2 a rr:TripleMap;
      rr:logicalTable [ rr:sqlQuery "SELECT * FROM students WHERE stype=0" ];
      rr:subjectMap [ rr:template "/UGStudent{id}"; rr:class ub:Student ].
```

```
\_ :m3 a rr:TripleMap;
      rr:logicalTable [ rr:sqlQuery "SELECT * FROM students WHERE stype=1" ];
      rr:subjectMap [ rr:template "/GradStudent{id}"; rr:class ub:Student ].
```

which generate undergraduate and graduate students (both are instances of  $\text{ub:Student}$ , but their IRIs are constructed using different templates [16]). Consider the following query (a fragment of  $q_2^{obg}$  from Section 5):

```
SELECT ?x ?y WHERE { ?x rdf:type ub:Student. ?x ub:UGDegreeFrom ?y }.
```

The translation  $\tau$  of its BGP (after the SPARQL JOIN optimisation of Section 3.1) is

$$\begin{aligned} & (\pi_{x,\rho_x/subj}\sigma_{(pred=rdf:type) \wedge (obj=ub:Student)} \text{triple}) \bowtie \\ & (\pi_{x,y,\rho_x/subj,\rho_y/obj}\sigma_{pred=ub:UGDegreeFrom} \text{triple}) \end{aligned}$$

First, since *triple* always occurs in the scope of some selection operation  $\sigma_F$ , we can choose only those elements in  $\bigcup_{m \in \mathcal{M}} tr_m$  that have matching values of *pred* and/or *obj*. In our example, the first occurrence of *triple* is replaced by  $tr_{:m2} \cup tr_{:m3}$ , and the second one by  $tr_{:m1}$ . This results in the natural join of the following union, denoted A:

```
(SELECT DISTINCT '/UGStudent' || id AS x FROM students
 WHERE (id IS NOT NULL) AND (stype=0))
 UNION (SELECT DISTINCT '/GradStudent' || id AS x FROM students
        WHERE (id IS NOT NULL) AND (stype=1))
```

and of the following query, denoted B:

```
SELECT DISTINCT '/GradStudent' || id AS x, '/Uni' || degreeuniid AS y FROM students
 WHERE (id IS NOT NULL) AND (degreeuniid IS NOT NULL) AND (stype=1)
```

Second, observe that the IRI template in B is compatible only with the second component of A. Moreover, since the two compatible templates coincide, we can *de-IRI* the join, namely, replace the join over the constructed strings ( $A.x = B.x$ ) by the join over the numerical attributes ( $A.id = B.id$ ), which results in a more efficient query:

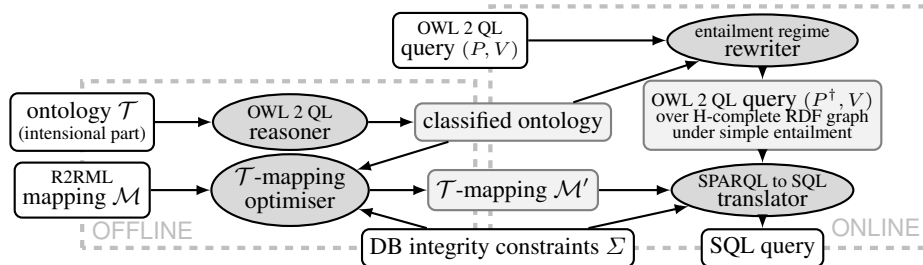
```
SELECT DISTINCT A.x, B.y FROM
 (SELECT id, '/GradStudent' || id AS x FROM students
  WHERE (id IS NOT NULL) AND (stype=1)) A
 JOIN
 (SELECT id, '/GradStudent' || id AS x, '/Uni' || degreeuniid AS y FROM students
  WHERE (id IS NOT NULL) AND (degreeuniid IS NOT NULL) AND (stype=1)) B
 ON A.id = B.id
```

Finally, by using self-join elimination and the fact that *id* and *stype* are the composite primary key in *students*, we obtain the query (without *DISTINCT* as *x* is unique)

```
SELECT '/GradStudent' || id AS x, '/Uni' || degreeuniid AS y FROM students
 WHERE (degreeuniid IS NOT NULL) AND (stype=1)
```

## 4 Putting it all Together

The techniques introduced above suggest the following architecture to support answering SPARQL queries under the OWL 2 QL entailment regime with data instances stored in a database. Suppose we are given an ontology with an intensional part  $\mathcal{T}$  and an extensional part stored in a database,  $D$ , over a schema  $\Sigma$ . Suppose also that the languages of  $\Sigma$  and  $\mathcal{T}$  are connected by an R2RML mapping  $\mathcal{M}$ . The process of answering a given OWL 2 QL query  $(P, V)$  involves two stages, off-line and on-line.



The *off-line* stage takes  $\mathcal{T}$ ,  $\mathcal{M}$  and  $\Sigma$  and proceeds via the following steps:

- 1 An OWL 2 QL reasoner is used to obtain a complete class / property hierarchy in  $\mathcal{T}$ .

② The composition  $\mathcal{M}^T$  of  $\mathcal{M}$  with the class and property hierarchy in  $\mathcal{T}$  is taken as an initial  $\mathcal{T}$ -mapping. Recall [29] that a mapping  $\mathcal{M}'$  is a  $\mathcal{T}$ -mapping over  $\Sigma$  if, for any data instance  $D$  satisfying  $\Sigma$ , the *virtual* (not materialised) RDF graph  $G_{D,\mathcal{M}'}$  obtained by applying  $\mathcal{M}'$  to  $D$  contains all class and property assertions  $\alpha$  with  $(\mathcal{T}, G_{D,\mathcal{M}'}) \models \alpha$ . As a result,  $G_{D,\mathcal{M}'}$  is complete with respect to the class and property hierarchy in  $\mathcal{T}$  (or H-complete), which allows us to avoid reasoning about class and property inclusions (in particular, inferences that involve property domains and ranges) at the query rewriting step ④ and drastically simplify rewritings (see [29] for details).

③ The initial  $\mathcal{T}$ -mapping  $\mathcal{M}^T$  is then optimised by (i) eliminating redundant triple maps detected by query containment with inclusion dependencies in  $\Sigma$ , (ii) eliminating redundant joins in logical tables using the functional dependencies in  $\Sigma$ , and (iii) merging sets of triple maps by means of interval expressions or disjunctions in logical tables (see [29] for details). Let  $\mathcal{M}'$  be the resulting  $\mathcal{T}$ -mapping over  $\Sigma$ .

The *on-line* stage takes an OWL 2 QL query  $(P, V)$  as an input and proceeds as follows:

④ The graph pattern  $P$  and  $\mathcal{T}$  are rewritten to the OWL 2 QL graph pattern  $P^\dagger$  over the H-complete virtual RDF graph  $G_{D,\mathcal{M}'}$  under simple entailment by applying the classified ontology of step ① to instantiate class and property variables and then using a query rewriting algorithm (e.g., the tree-witness rewriter of [29]); see Theorem 4.

⑤ The graph pattern  $P^\dagger$  is transformed to the SQL query  $\tau(P^\dagger)$  over the 3-column representation *triple* of the RDF graph (Theorem 7). Next, the query  $\tau(P^\dagger)$  is unfolded into the SQL query  $\text{tr}_{\mathcal{M}'}(\tau(P^\dagger))$  over the original database  $D$  (Theorem 12). The unfolded query is optimised using the techniques similar to the ones employed in step ③.

⑥ The optimised query is executed by the database.

As follows from Theorems 4, 7 and 12, the resulting query gives us all correct answers to the original OWL 2 QL query  $(P, V)$  over  $\mathcal{T}$  and  $D$  with the R2RML mapping  $\mathcal{M}$ .

## 5 Evaluation

The architecture described above has been implemented in the open-source OBDA system *Ontop*<sup>4</sup>. We evaluated its performance using the OWL 2 QL version of the Lehigh University Benchmark LUBM [16]. The ontology contains 43 classes, 32 object and data properties and 243 axioms. The benchmark also includes a data generator and a set of 14 queries  $q_1$ – $q_{14}$ . We added 7 queries with second-order variables ranging over class and property names:  $q'_4, q''_4, q'_9, q''_9$  derived from  $q_4$  and  $q_9$ , and  $q_2^{obg}, q_4^{obg}, q_{10}^{obg}$  taken from [19]. The LUBM data generator produces an OWL file with class and property assertions. To store the assertions in a database, we created a database schema with 11 relations and an R2RML mapping with 89 predicate-object maps. For instance, the information about undergraduate and graduate students (id, name, etc.) from Example 10 is collected in the relation *students*, where the attribute *stype* distinguishes between the types of students (*stype* is known as a discriminant column in databases); more details including primary and foreign keys and indexes are provided in the full version.

We experimented with the data instances  $\text{LUBM}_n$ ,  $n = 1, 9, 20, 50, 100, 200, 500$  (where  $n$  specifies the number of universities;  $\text{LUBM}_1$  and  $\text{LUBM}_9$  were used in [19]).

<sup>4</sup> <http://ontop.inf.unibz.it>

Q	LUBM <sub>1</sub>				LUBM <sub>9</sub>			LUBM <sub>100</sub>		LUBM <sub>200</sub>	LUBM <sub>500</sub>
	O	OB <sub>H</sub>	OB <sub>P</sub>	P	O	OB <sub>H</sub>	P	O	P	O	O
q <sub>1</sub>	2	8	29	1	3	97	1	3	1	3	2
q <sub>2</sub>	2	25	11 137	19	3	2 531	256	16	30 593	36	88
q <sub>3</sub>	1	6	86	9	2	78	158	2	2 087	63	12
q <sub>4</sub>	13	7	19	14	15	44	164	27	2 093	24	22
q <sub>5</sub>	16	12	4 451	10	22	98	158	32	2 182	28	23
q <sub>6</sub>	455	27	32	21	5 076	411	317	58 968	10 781	123 578	434 349
q <sub>7</sub>	5	21	34 005	10	6	429	157	8	2 171	8	9
q <sub>8</sub>	726	195	95 875	80	760	917	192	796	2 131	820	855
q <sub>9</sub>	60	972	168 978	78	668	189 126	857	7 466	12 125	15 227	44 598
q <sub>10</sub>	2	6	126	9	3	97	158	2	2 134	3	2
q <sub>11</sub>	4	5	58	10	6	43	160	11	2 093	18	44
q <sub>12</sub>	3	4	19	15	4	70	236	3	2 114	5	5
q <sub>13</sub>	6	4	67	8	7	40	157	14	2 657	38	58
q <sub>14</sub>	91	20	24	15	1 168	329	287	13 524	4 457	29 512	92 376
q <sub>4</sub> <sup>r</sup>	93	58	190	46	99	98	767	92	4 422	95	107
q <sub>4</sub> <sup>r</sup>	108	21	35	63	122	72	719	115	9 179	108	127
q <sub>9</sub> <sup>r</sup>	257	716	91 855	174	4 686	40 575	1 385	54 092	19 945	115 110	295 228
q <sub>9</sub> <sup>r</sup>	557	951	65 916	102	6 093	178 401	1 214	67 123	19 705	151 376	356 176
q <sub>2</sub> <sup>obg</sup>	150	30	57 141	29	9 992	520	348	39 477	5 411	79 351	206 061
q <sub>4</sub> <sup>obg</sup>	6	7	241	25	31	40	273	7	3 969	7	494
q <sub>10</sub> <sup>obg</sup>	641	760	31 269	253	6 998	149 191	2 258	163 308	17 929	174 362	459 669
start up	3.1s	13.6s	7.7s	3.6s	3.1s	80m33s	18s	3.1s	3m23s	3.1s	3.1s
data load	10s	n/a	n/a	n/a	15s	n/a	n/a	1m56s	n/a	3m35s	10m17s

**Table 1.** Start up time, data loading time (in s) and query execution time (in ms): O is *Ontop*, OB<sub>H</sub> and OB<sub>P</sub> are OWL-BGP with Hermit and Pellet, respectively, and P is standalone Pellet.

Here we only show the results for  $n = 1, 9, 100, 200, 500$  containing 103k, 1.2M, 14M, 28M and 69M triples, respectively; the complete table can be found in the full version. All the materials required for the experiments are available online<sup>5</sup>. We compared *Ontop* with two other systems, OWL-BGP r123 [19] and Pellet 2.3.1 [31] (Stardog and OWLIM are incomplete for the OWL 2 QL entailment regime). OWL-BGP requires an OWL 2 reasoner as a backend; as in [19], we employed Hermit 1.3.8 [14] and Pellet 2.3.1. The hardware was an HP Proliant Linux server with 144 cores @3.47GHz, 106GB of RAM and a 1TB 15k RPM HD. Each system used a single core and was given 20 GB of Java 7 heap memory. *Ontop* used MySQL 5.6 database engine.

The evaluation results are given in Table 1. OWL-BGP and Pellet used significantly more time to start up (last but one row) because they do not rely on query rewriting and require costly pre-computations. OWL-BGP failed to start on LUBM<sub>9</sub> with Pellet and on LUBM<sub>20</sub> with Hermit; Pellet ran out of memory after 10hrs loading LUBM<sub>200</sub>. For *Ontop*, the start up is the off-line stage described in Section 4; it does not include the time of loading the data into MySQL, which is specified in the last row of Table 1 (note that the data is loaded only once, not every time *Ontop* starts; moreover, this could be improved with CSV loading and delayed indexing rather than SQL dumps we used).

On queries  $q_1$ – $q_{14}$ , *Ontop* generally outperforms OWL-BGP and Pellet. Due to the optimisations, the SQL queries generated by *Ontop* are very simple, and MySQL is able to execute them efficiently. This is also the case for large datasets, where *Ontop* is able to maintain almost constant times for many of the queries. Notable exceptions are  $q_6$ ,  $q_8$  and  $q_{14}$  that return a very large number (hundreds of thousands) of results (low

<sup>5</sup> <https://github.com/ontop/iswc2014-benchmark>

selectivity). A closer inspection reveals that execution time is mostly spent on fetching the results from disk. On the queries with second-order variables, the picture is mixed. While indeed these queries are not the strongest point of *Ontop* at the moment, we see that in general the performance is good. Although Pellet outperforms *Ontop* on small datasets, only *Ontop* is able to provide answers for very large datasets. For second-order queries with high selectivity (e.g.,  $q'_4$  and  $q''_4$ ) and large datasets, the performance of *Ontop* is very good while the other systems fail to return answers.

## 6 Conclusions

In this paper, we gave both a theoretical background and a practical implementation of a procedure for answering SPARQL 1.1 queries under the OWL 2 QL direct semantics entailment regime in the scenario where data instances are stored in a relational database whose schema is connected to the language of the given OWL 2 QL ontology via an R2RML mapping. Our main contributions can be summarised as follows:

- We defined an entailment regime for SPARQL 1.1 corresponding to the OWL 2 QL profile of OWL 2 (which was specifically designed for ontology-based data access).
- We proved that answering SPARQL queries under this regime is reducible to answering SPARQL queries under simple entailment (where no reasoning is involved).
- We showed how to transform such SPARQL queries to equivalent SQL queries over an RDF representation of the data, and then unfold them, using R2RML mappings, into SQL queries over the original relational data.
- We developed optimisation techniques to substantially reduce the size and improve the quality of the resulting SQL queries.
- We implemented these rewriting and optimisation techniques in the OBDA system *Ontop*. Our initial experiments showed that *Ontop* generally outperforms reasoner-based systems, especially on large data instances.

Some aspects of SPARQL 1.1 (such as RDF types, property paths, aggregates) were not discussed here and are left for future work.

**Acknowledgements.** Our work was supported by EU project Optique. We thank S. Kollia-Ebri for help with the experiments, and I. Kollia and B. Glimm for discussions.

## References

1. Abiteboul, S., Hull, R., Vianu, V.: Foundations of Databases. Addison-Wesley (1995)
2. Angles, R., Gutierrez, C.: The expressive power of SPARQL. In: Proc. of ISWC. LNCS, vol. 5318, pp. 114–129. Springer (2008)
3. Bornea, M., Dolby, J., Kementsietsidis, A., Srinivas, K., Dantressangle, P., Udrea, O., Bhattacharjee, B.: Building an efficient RDF store over a relational database. In: Proc. of SIGMOD 2013, pp. 121–132. ACM (2013)
4. Calvanese, D., De Giacomo, G., Lembo, D., Lenzerini, M., Poggi, A., Rodriguez-Muro, M., Rosati, R., Ruzzi, M., Savo, D.F.: The MASTRO system for ontology-based data access. Semantic Web 2(1), 43–53 (2011)
5. Calvanese, D., De Giacomo, G., Lembo, D., Lenzerini, M., Rosati, R.: Tractable reasoning and efficient query answering in description logics: The *DL-Lite* family. J. of Automated Reasoning 39(3), 385–429 (2007)

6. Chakravarthy, U.S., Grant, J., Minker, J.: Logic-based approach to semantic query optimization. *ACM Transactions on Database Systems* 15(2), 162–207 (1990)
7. Chebotko, A., Lu, S., Fotouhi, F.: Semantics preserving SPARQL-to-SQL translation. *Data Knowl. Eng.* 68(10), 973–1000 (2009)
8. Chortaras, A., Trivela, D., Stamou, G.: Optimized query rewriting for OWL 2 QL. In: *Proc. of CADE-23. LNCS*, vol. 6803, pp. 192–206. Springer (2011)
9. Cyganiak, R.: A relational algebra for SPARQL. *Tech. Rep. HPL-2005-170, HP Labs* (2005)
10. Das, S., Sundara, S., Cyganiak, R.: R2RML: RDB to RDF Mapping Language (September 2012), <http://www.w3.org/TR/r2rml>
11. Dolby, J., Fokoue, A., Kalyanpur, A., Ma, L., Schonberg, E., Srinivas, K., Sun, X.: Scalable grounded conjunctive query evaluation over large and expressive knowledge bases. In: *Proc. of ISWC. LNCS*, vol. 5318, pp. 403–418. Springer (2008)
12. Eiter, T., Ortiz, M., Šimkus, M., Tran, T.K., Xiao, G.: Query rewriting for Horn-SHIQ plus rules. In: *Proc. of AAAI. AAAI Press* (2012)
13. Elliott, B., Cheng, E., Thomas-Ogbuji, C., Özsoyoglu, Z.M.: A complete translation from SPARQL into efficient SQL. In: *Proc. of IDEAS*. pp. 31–42. ACM (2009)
14. Glimm, B., Horrocks, I., Motik, B., Stoilos, G.: Optimising ontology classification. In: *Proc. of ISWC, part I. LNCS*, vol. 6496, pp. 225–240. Springer (2010)
15. Gottlob, G., Orsi, G., Pieris, A.: Ontological queries: Rewriting and optimization. In: *Proc. of ICDE*. pp. 2–13. IEEE Computer Society (2011)
16. Guo, Y., Pan, Z., Heflin, J.: LUBM: A benchmark for OWL knowledge base systems. *J. of Web Semantics* 3(2–3), 158–182 (2005)
17. Heymans, S. *et al.*: Ontology reasoning with large data repositories. In: *Ontology Management, Semantic Web, Semantic Web Services, and Business Applications*. Springer (2008)
18. King, J.J.: Query Optimization by Semantic Reasoning. Ph.D. thesis, Stanford, USA (1981)
19. Kollia, I., Glimm, B.: Optimizing SPARQL query answering over OWL ontologies. *J. of Artificial Intelligence Research* 48, 253–303 (2013)
20. König, M., Leclère, M., Mugnier, M.-L., Thomazo, M.: On the exploration of the query rewriting space with existential rules. In: *Proc. of RR. LNCS*. pp. 123–137. Springer (2013)
21. Kontchakov, R., Rodríguez-Muro, M., Zakharyashev, M.: Ontology-based data access with databases: A shortcourse. In: *Reasoning Web. LNCS*, vol. 8067, pp. 194–229. Springer (2013)
22. Lutz, C., Seylan, I., Toman, D., Wolter, F.: The combined approach to OBDA: Taming role hierarchies using filters. In: *Proc. of ISWC. LNCS*, vol. 8218, pp. 314–330. Springer (2013)
23. Pérez-Urbina, H., Rodríguez-Díaz, E., Grove, M., Konstantinidis, G., Sirin, E.: Evaluation of query rewriting approaches for OWL 2. In: *SSWS+HPCSW. CEUR-WS*, vol. 943 (2012)
24. Poggi, A., Lembo, D., Calvanese, D., De Giacomo, G., Lenzerini, M., Rosati, R.: Linking data to ontologies. *J. on Data Semantics X*, 133–173 (2008)
25. Polleres, A.: From SPARQL to rules (and back). In: *Proc. WWW*. pp. 787–796. ACM (2007)
26. Polleres, A., Wallner, J.P.: On the relation between SPARQL 1.1 and Answer Set Programming. *J. of Applied Non-Classical Logics* 23(1–2), 159–212 (2013)
27. Priyatna, F., Corcho, O., Sequeda, J.: Formalisation and experiences of R2RML-based SPARQL to SQL query translation using Morph. In: *Proc. of WWW*. pp. 479–490 (2014)
28. Rodríguez-Muro, M., Hardi, J., Calvanese, D.: Quest: Efficient SPARQL-to-SQL for RDF and OWL. In: *Proc. of the ISWC 2012 P&D Track*. vol. 914. CEUR-WS.org (2012)
29. Rodríguez-Muro, M., Kontchakov, R., Zakharyashev, M.: Ontology-based data access: Ontop of databases. In: *Proc. of ISWC. LNCS*, vol. 8218, pp. 558–573. Springer (2013)
30. Sequeda, J.F., Miranker, D.P.: Ultrawrap: SPARQL execution on relational data. *J. of Web Semantics* 22, 19–39 (2013)
31. Sirin, E., Parsia, B., Cuenca Grau, B., Kalyanpur, A., Katz, Y.: Pellet: A practical OWL-DL Reasoner. *J. of Web Semantics* 5(2), 51–53 (2007)
32. Zemke, F.: Converting SPARQL to SQL. *Tech. rep.*, Oracle Corp. (2006)



## A On the Semantics of SPARQL

*Remark 1.* The condition ‘ $F^{s_1 \oplus s_2}$  is not true’ in our definition of OPT is slightly different from ‘ $F^{s_1 \oplus s_2}$  has an effective Boolean value of false’ given by the W3C specification<sup>6</sup>. The two definitions do not necessarily coincide because the effective Boolean value can be undefined (type error) if, for example, a variable in  $F$  is not bound by  $s_1 \oplus s_2$ . As we see from Section 3.1, our reading corresponds to LEFT JOIN in SQL.

We also find the informal explanation of the semantics for OPT in the W3C recommendation inconsistent with the definition of DIFF, which forms the second component of the union in the definition of OPT. It suggests that  $\text{DIFF}(S_1, S_2, F)$  is equivalent to

$$\begin{aligned} \text{DIFF}'(S_1, S_2, F) &= \{s_1 \in S_1 \mid s_1 \text{ and } s_2 \text{ are incompatible, for all } s_2 \in S_2\} \\ &\cup \{s_1 \in S_1 \mid \text{there is } s_2 \in S_2 \text{ compatible with } s_1 \text{ such that } F^{s_1 \oplus s_2} = \perp\}. \end{aligned}$$

Observe that there may be  $s_2, s'_2 \in S_2$  that are both compatible with  $s_1$ ,  $F^{s_1 \oplus s_2} = \top$  and  $F^{s_1 \oplus s'_2} = \perp$ , in which case  $s_1 \in \text{DIFF}'(S_1, S_2, F) \setminus \text{DIFF}(S_1, S_2, F)$ .

## B Proof of Theorem 1

**Theorem 4.** *Given any intensional graph  $\mathcal{T}$  and OWL2 QL query  $(P, V)$ , one can construct an OWL2 QL query  $(P^\dagger, V)$  such that, for any extensional graph  $\mathcal{A}$  (in some fixed finite vocabulary),*

$$\llbracket P \rrbracket_{\mathcal{T}, \mathcal{A} | V} = \llbracket P^\dagger \rrbracket_{\mathcal{A} | V}.$$

*Proof.* By the definition of the entailment regime, it suffices to construct a rewriting  $B^\dagger$ , for any *basic* graph pattern  $B$ , such that  $\llbracket B \rrbracket_{\mathcal{T}, \mathcal{A}} = \llbracket B^\dagger \rrbracket_{\mathcal{A}}$ . A rewriting  $P^\dagger$  of  $P$  is obtained by replacing every BGP  $B$  in it with  $B^\dagger$ .

Take any BGP  $B$  that occurs in  $P$ . Let  $?c_1, \dots, ?c_m$  be all class variables in  $B$ , and let  $?r_1, \dots, ?r_n$  be all property variables in  $B$ . For any  $m$ -tuple  $\mathbf{C} = (C_1, \dots, C_m)$  of class names and  $n$ -tuple  $\mathbf{R} = (R_1, \dots, R_n)$  of property names in the given vocabulary, we take the result  $B'_{\mathbf{C}, \mathbf{R}}$  of replacing every  $?c_i$  in  $B$  with  $C_i$  and every  $?r_j$  in  $B$  with  $R_j$ . By definition,  $B'_{\mathbf{C}, \mathbf{R}}$  contains no class or property variables and contains only OWL2 QL class and property axioms (rather than templates) and assertions. For any class or property axiom in  $B'_{\mathbf{C}, \mathbf{R}}$ , we check whether it is entailed by (logically follows from)  $\mathcal{T}$ . If at least one of the axioms is not entailed by  $\mathcal{T}$ , we set  $B''_{\mathbf{C}, \mathbf{R}}$  to be the empty BGP; otherwise, let  $B''_{\mathbf{C}, \mathbf{R}}$  be the result of removing all class or property axioms from  $B'_{\mathbf{C}, \mathbf{R}}$ .

The constructed BGP  $B''_{\mathbf{C}, \mathbf{R}}$  contains only class and property assertions and can be regarded as (a SPARQL representations of) a conjunctive query. As is well-known [29], we can rewrite this query and  $\mathcal{T}$  into a union of conjunctive queries over  $\mathcal{A}$ , which can be represented as a union  $B^\dagger_{\mathbf{C}, \mathbf{R}}$  of BGPs. Now, we take  $B^\dagger$  to be the union (UNION) of

$$B^\dagger_{\mathbf{C}, \mathbf{R}}[?c_1 \mapsto C_1, \dots, ?c_m \mapsto C_m, ?r_1 \mapsto R_1, \dots, ?r_n \mapsto R_n],$$

<sup>6</sup> <http://www.w3.org/TR/sparql11-query/#sparqlAlgebra>

for all possible  $m$ -tuples  $C$  of class names and  $n$ -tuples  $R$  of property names in the given vocabulary such that  $B''_{C,R}$  is not empty (the empty union is, by definition, the empty BGP). Here

$$B[?v_1 \mapsto V_1, \dots, ?v_k \mapsto V_k] = \text{BIND}(\dots \text{BIND}(B, ?v_1, V_1), \dots, ?v_k, V_k).$$

It follows from the construction that  $\llbracket B \rrbracket_{\mathcal{T}, \mathcal{A}} = \llbracket B^\dagger \rrbracket_{\mathcal{A}}$ .  $\square$

We emphasise that  $\llbracket P^\dagger \rrbracket_{\mathcal{A}}$  is computed only on the extensional part of the RDF graph and under simple entailment.

## C Proof of Proposition 6 and Theorem 7

**Proposition 6.** *Let  $V$  be a set of variables and  $F$  a SPARQL filter expression with variables in  $V$ . For each solution mapping  $s$  with  $\text{dom}(s) \subseteq V$ , we have  $F^s = (\tau(F))^{ext_V(s)}$ .*

*Proof.* The proof is based on the observation that  $ext_V(s): v \mapsto null$  just in case  $v \notin \text{dom}(s)$ , for each  $v \in V$ . Thus, the claim holds for  $F = bound(v)$ . Also, due to this observation, the clauses in the definition of  $v = c$  and  $v = v'$  and the truth tables for  $\neg$  and  $\wedge$  coincide for SPARQL and SQL filter expressions.  $\square$

**Theorem 7.** *For any RDF graph  $G$  and any graph pattern  $P$ ,*

$$\|P\|_G = \|\tau(P)\|_{triple(G)}.$$

*Proof.* The proof is by induction on the structure of  $P$ .

For the basis of induction, let  $P$  be a triple pattern of the form  $\langle s, p, o \rangle$ . Since each of the components of the triple pattern is either a variable in  $\mathbf{V}$  or an RDF term in  $\mathbf{I} \cup \mathbf{L}$ , there are 15 possible cases (recall that we have no blank nodes):

$s, p, o \in \mathbf{I} \cup \mathbf{L}$	$s, p, o \in \mathbf{I} \cup \mathbf{L}$	$o \in \mathbf{V}, s, p \in \mathbf{I} \cup \mathbf{L}$
$s \in \mathbf{V}, p, o \in \mathbf{I} \cup \mathbf{L}$	$p \in \mathbf{V}, s, o \in \mathbf{I} \cup \mathbf{L}$	$p, o \in \mathbf{V}, p \neq o, s \in \mathbf{I} \cup \mathbf{L}$
$s, p \in \mathbf{V}, s \neq p, o \in \mathbf{I} \cup \mathbf{L}$	$s, o \in \mathbf{V}, s \neq o, p \in \mathbf{I} \cup \mathbf{L}$	$p, o \in \mathbf{V}, p = o, s \in \mathbf{I} \cup \mathbf{L}$
$s, p \in \mathbf{V}, s = p, o \in \mathbf{I} \cup \mathbf{L}$	$s, o \in \mathbf{V}, s = o, p \in \mathbf{I} \cup \mathbf{L}$	$p, o \in \mathbf{V}, p = o, s \in \mathbf{I} \cup \mathbf{L}$
$s, p, o \in \mathbf{V}, s \neq p, p \neq o, o \neq s$	$s, p, o \in \mathbf{V}, s \neq p, p \neq o, o \neq s$	$s, p, o \in \mathbf{V}, o = s \neq p$
$s, p, o \in \mathbf{V}, s = p \neq o$	$s, p, o \in \mathbf{V}, p = o \neq s$	
	$s, p, o \in \mathbf{V}, s = p = o$	

We consider just one case with  $s, p \in \mathbf{V}, s \neq p$  and  $o \in \mathbf{I} \cup \mathbf{L}$  and leave all other cases to the reader. By definition, we have  $\tau(\langle s, p, o \rangle) = \pi_{s,p} \rho_{s/subj} \rho_{p/pred} \sigma_{obj=o} triple$  and  $\text{var}(\langle s, p, o \rangle) = \{s, p\}$ . Then the following are equivalent:

- $\|\langle s, p, o \rangle\|_G$  contains tuple  $\{s \mapsto a, p \mapsto b\}$ ;
- $\llbracket \langle s, p, o \rangle \rrbracket_G$  contains solution mapping  $\{s \mapsto a, p \mapsto b\}$ ;
- $G$  contains triple  $(a, b, o)$ ;
- $triple(G)$  contains tuple  $\{subj \mapsto a, pred \mapsto b, obj \mapsto o\}$ ;
- $\|\pi_{s,p} \rho_{s/subj} \rho_{p/pred} \sigma_{obj=o} triple\|_{triple(G)}$  contains tuple  $\{s \mapsto a, p \mapsto b\}$ .

For the induction step, we consider the five cases of SPARQL algebra operations.

$P = \text{FILTER}(P_1, F)$ . Denote  $U = \text{var}(P)$  (recall that all variables of  $F$  are in  $U$ ).

- ( $\subseteq$ ) Let  $t \in \|\text{FILTER}(P_1, F)\|_G$ . Then there is  $s \in \llbracket \text{FILTER}(P_1, F) \rrbracket_G$  such that  $\text{ext}_U(s) = t$ . By definition, we have  $s \in \llbracket P_1 \rrbracket_G$  and  $F^s = \top$ . Then  $t \in \|\llbracket P_1 \rrbracket_G\|_G$ , whence, by the induction hypothesis,  $t \in \|\tau(P_1)\|_{\text{triple}(G)}$  and, by Proposition 6,  $(\tau(F))^t = \top$ . Thus,  $t \in \|\sigma_{\tau(F)}\tau(P_1)\|_{\text{triple}(G)}$ , which coincides with  $\|\tau(\text{FILTER}(P_1, F))\|_{\text{triple}(G)}$ .
- ( $\subseteq$ ) Let  $t \in \|\tau(\text{FILTER}(P_1, F))\|_{\text{triple}(G)}$ . By definition,  $t \in \|\tau(P_1)\|_{\text{triple}(G)}$  and  $(\tau(F))^t = \top$ . By the induction hypothesis,  $t \in \|\llbracket P_1 \rrbracket_G\|_G$ . Then there is a solution mapping  $s$  such that  $t = \text{ext}_U(s)$  and  $s \in \llbracket P_1 \rrbracket_G$ . By Proposition 6,  $F^s = \top$  and thus, we obtain  $s \in \llbracket \text{FILTER}(P_1, F) \rrbracket_G$  and  $t \in \|\text{FILTER}(P_1, F)\|_G$ .
- $P = \text{BIND}(P_1, v, c)$ . (Recall that  $v \notin \text{var}(P_1)$ .)
- ( $\subseteq$ ) Let  $t \in \|\text{BIND}(P_1, v, c)\|_G$ . Then there is  $s \in \llbracket \text{BIND}(P_1, v, c) \rrbracket_G$  such that  $\text{ext}_U(s) = t$ . By definition, we have  $s' \in \llbracket P_1 \rrbracket_G$  for  $s'$  that coincides with  $s$  on  $\text{dom}(s) \setminus \{v\}$  and is undefined on  $v$ . Then  $s' \in \|\llbracket P_1 \rrbracket_G\|_G$ , whence, by the induction hypothesis,  $s' \in \|\tau(P_1)\|_{\text{triple}(G)}$ . Thus,  $t \in \|\tau(P_1) \times \{v \mapsto c\}\|_{\text{triple}(G)}$ , which coincides with  $\|\tau(\text{BIND}(P_1, v, c))\|_{\text{triple}(G)}$ .
- ( $\subseteq$ ) Let  $t \in \|\tau(\text{BIND}(P_1, v, c))\|_{\text{triple}(G)}$ . By definition,  $t(v) = c$  and the restriction  $t'$  of  $t$  to  $\text{var}(P) \setminus \{v\}$  is in  $\|\tau(P_1)\|_{\text{triple}(G)}$ . By the induction hypothesis,  $t' \in \|\llbracket P_1 \rrbracket_G\|_G$ . So, there is a solution mapping  $s$  with  $t' = \text{ext}_U(s)$  and  $s \in \llbracket P_1 \rrbracket_G$ . Thus,  $t \in \llbracket \text{BIND}(P_1, v, c) \rrbracket_G$ .
- $P = \text{UNION}(P_1, P_2)$ . Denote  $U_i = \text{var}(P_i)$ , for  $i = 1, 2$ .
- ( $\subseteq$ ) Let  $t \in \|\text{UNION}(P_1, P_2)\|_G$ . Then there is  $s \in \llbracket \text{UNION}(P_1, P_2) \rrbracket_G$  such that  $t = \text{ext}_{U_1 \cup U_2}(s)$ . By definition, we have either  $s \in \llbracket P_1 \rrbracket_G$  or  $s \in \llbracket P_2 \rrbracket_G$  and so, either  $\text{ext}_{U_1}(s) \in \|\llbracket P_1 \rrbracket_G\|_G$  or  $\text{ext}_{U_2}(s) \in \|\llbracket P_2 \rrbracket_G\|_G$ . By the induction hypothesis, either  $\text{ext}_{U_1}(s) \in \|\tau(P_1)\|_{\text{triple}(G)}$  or  $\text{ext}_{U_2}(s) \in \|\tau(P_2)\|_{\text{triple}(G)}$ , which implies  $\text{ext}_{U_1 \cup U_2}(s) = t \in \|\tau(\text{UNION}(P_1, P_2))\|_{\text{triple}(G)}$ .
- ( $\supseteq$ ) Let  $t \in \|\tau(\text{UNION}(P_1, P_2))\|_{\text{triple}(G)}$ . Let  $s$  be such that  $t = \text{ext}_{U_1 \cup U_2}(s)$ . By definition, either  $\text{ext}_{U_1}(s) \in \|\tau(P_1)\|_{\text{triple}(G)}$  or  $\text{ext}_{U_2}(s) \in \|\tau(P_2)\|_{\text{triple}(G)}$ . By the induction hypothesis, either  $\text{ext}_{U_1}(s) \in \|\llbracket P_1 \rrbracket_G\|_G$  or  $\text{ext}_{U_2}(s) \in \|\llbracket P_2 \rrbracket_G\|_G$ , which implies  $s \in \llbracket P_1 \rrbracket_G$  or  $s \in \llbracket P_2 \rrbracket_G$ . Thus,  $s \in \llbracket \text{UNION}(P_1, P_2) \rrbracket_G$  and thus,  $t \in \|\text{UNION}(P_1, P_2)\|_G$ .
- $P = \text{JOIN}(P_1, P_2)$ . Let  $U_i = \text{var}(P_i)$ , for  $i = 1, 2$ , and  $U = U_1 \cap U_2$ .
- ( $\subseteq$ ) If  $t \in \|\text{JOIN}(P_1, P_2)\|_G$  then there is a solution mapping  $s \in \llbracket \text{JOIN}(P_1, P_2) \rrbracket_G$  with  $\text{ext}_{U_1 \cup U_2}(s) = t$ , and so there are compatible  $s_1$  and  $s_2$  with  $s_1 \oplus s_2 = s$  and  $s_i \in \llbracket P_i \rrbracket_G$ , for  $i = 1, 2$ . By definition,  $\text{ext}_{U_i}(s_i) \in \|\llbracket P_i \rrbracket_G\|_G$ , for  $i = 1, 2$ , whence, by the induction hypothesis,  $\text{ext}_{U_i}(s_i) \in \|\tau(P_i)\|_{\text{triple}(G)}$ . Suppose  $V = \text{dom}(s_1) \cap \text{dom}(s_2)$ ,  $V_1 = U \setminus \text{dom}(s_1)$  and  $V_2 = \text{dom}(s_1) \cap (U \setminus \text{dom}(s_2))$ . Then  $V_1, V_2$  and  $V$  are disjoint and partition  $U$  (note that  $V_1$  and  $V_2$  could be chosen differently since, in general,  $(U \setminus \text{dom}(s_1)) \cap (U \setminus \text{dom}(s_2)) \neq \emptyset$ ). Consider  $t_i = \text{ext}_{U_i \setminus V_i}(s_i)$ , and let  $Q_i = \pi_{U_i \setminus V_i}(\sigma_{\text{isNull}(V_i)}\tau(P_i))$ . By definition,  $t_i: v \mapsto \text{null}$ , for each  $v \in V_i$ , and  $t_i: v \mapsto \text{ext}_{U_i} s_i(v)$ , for each  $v \in U_i \setminus V_i$ . Therefore,  $t_i \in \|\llbracket Q_i \rrbracket_{\text{triple}(G)}\|_G$ . Since  $s_1$  and  $s_2$  are compatible and  $V$  are the common non-null attributes of  $t_1$  and  $t_2$ , the tuples  $t_1$  and  $t_2$  are compatible and  $t_1 \oplus t_2 \in \|\llbracket Q_1 \bowtie Q_2 \rrbracket_{\text{triple}(G)}\|_G$ . It remains to observe that  $t = t_1 \oplus t_2$ .

( $\supseteq$ ) If  $t \in \|\tau(\text{JOIN}(P_1, P_2))\|_{\text{triple}(G)}$  then there are disjoint  $V_1, V_2 \subseteq U$  and compatible tuples  $t_1$  and  $t_2$  with  $t = t_1 \oplus t_2$  and  $t_i \in \|\pi_{U_i \setminus V_i} \sigma_{\text{isNull}(V_i)} \tau(P_i)\|_{\text{triple}(G)}$ . Let  $s_i$  be the solution mapping such that  $t_i = \text{ext}_{U_i}(s_i)$ . Then  $s_1$  and  $s_2$  are compatible and  $\text{ext}_{U_i}(s_i) \in \|\tau(P_i)\|_{\text{triple}(G)}$ . By the induction hypothesis,  $\text{ext}_{U_i}(s_i) \in \|P_i\|_G$  and  $s_i \in \|P_i\|_G$ . Therefore,  $s_1 \oplus s_2 \in \|\text{JOIN}(P_1, P_2)\|_G$  and  $\text{ext}_{U_1 \cup U_2}(s_1 \oplus s_2) = t \in \|\text{JOIN}(P_1, P_2)\|_G$ .

$P = \text{OPT}(P_1, P_2, F)$ . Let  $U_i = \text{var}(P_i)$ , for  $i = 1, 2$ , and  $U = U_1 \cap U_2$ . Recall that we have assumed that the variables of  $F$  are in  $U_1 \cup U_2$ . Observe first that, for any  $V_1 \subseteq U$  and any solution mapping  $s$ , we have

$$s \in \|P_1\|_G \text{ and } \text{dom}(s) \cap U = U \setminus V_1 \quad \text{iff} \quad s \in \|P_1^{V_1, U}\|_G. \quad (*)$$

Now, we claim that, for any  $V_1 \subseteq U$  and any solution mapping  $s$ , the following two conditions are equivalent:

- $s \in \|P_1\|_G$ ,  $\text{dom}(s) \cap U = U \setminus V_1$  and there is  $s_2 \in \|P_2\|_G$  such that  $s$  and  $s_2$  are compatible and  $F^{s \oplus s_2} = \top$ ;
- $\text{ext}_{U_1}(s) \in \|\mu_{V_1} \pi_{U_1 \setminus V_1} \tau(\text{FILTER}(\text{JOIN}(P_1^{V_1, U}, P_2), F))\|_{\text{triple}(G)}$ .

Indeed, suppose  $s \in \|P_1\|_G$ ,  $\text{dom}(s) \cap U = U \setminus V_1$  and there is  $s_2 \in \|P_2\|_G$  such that  $s$  and  $s_2$  are compatible and  $F^{s \oplus s_2} = \top$ . By (\*), we have  $s \in \|P_1^{V_1, U}\|_G$ , whence  $s \oplus s_2 \in \|\text{FILTER}(\text{JOIN}(P_1^{V_1, U}, P_2), F)\|_G$ . By repeating the argument for FILTER and JOIN, we get  $\text{ext}_{U_1 \cup U_2}(s \oplus s_2) \in \|\tau(\text{FILTER}(\text{JOIN}(P_1^{V_1, U}, P_2), F))\|_{\text{triple}(G)}$ . Since  $\text{dom}(s) \cap U = U \setminus V_1$ , the projection of  $\text{ext}_{U_1 \cup U_2}(s \oplus s_2)$  on  $U_1 \setminus V_1$  coincides with the projection of  $\text{ext}_{U_1}(s)$ ; moreover, the extension of this projection by *nulls* to  $V_1$  coincides with  $\text{ext}_{U_1}(s)$ . Therefore, we obtain that  $\text{ext}_{U_1}(s)$  belongs to  $\|\mu_{V_1} \pi_{U_1 \setminus V_1} \tau(\text{FILTER}(\text{JOIN}(P_1^{V_1, U}, P_2), F))\|_{\text{triple}(G)}$ , as required.

Conversely, let  $\text{ext}_{U_1}(s) \in \|\mu_{V_1} \pi_{U_1 \setminus V_1} \tau(\text{FILTER}(\text{JOIN}(P_1^{V_1, U}, P_2), F))\|_{\text{triple}(G)}$ . Then  $\text{dom}(s) \cap V_1 = \emptyset$ , and there is a tuple  $t'$  over  $U_1 \cup U_2$  that belongs to  $\|\tau(\text{FILTER}(\text{JOIN}(P_1^{V_1, U}, P_2), F))\|_{\text{triple}(G)}$  and such that  $\text{ext}_{U_1}(s)$  coincides with  $t'$  on  $U_1 \setminus V_1$ . By repeating the argument for FILTER, we get  $s' \in \|\text{JOIN}(P_1^{V_1, U}, P_2)\|_G$  such that  $t' = \text{ext}_{U_1 \cup U_2}(s')$  and  $F^{s'} = \top$ . By definition, there are  $s_1 \in \|P_1^{V_1, U}\|_G$  and  $s_2 \in \|P_2\|_G$  such that  $s_1$  and  $s_2$  are compatible and  $s' = s_1 \oplus s_2$ . By (\*), we have  $s_1 \in \|P_1\|_G$  and  $\text{dom}(s_1) \cap U = U \setminus V_1$ . Since  $\text{ext}_{U_1}(s)$  coincides with  $t' = \text{ext}_{U_1 \cup U_2}(s_1 \oplus s_2)$  on  $U_1 \setminus V_1$ , it also coincides with  $s_1$  on  $U_1 \setminus V_1$  ( $= U_1 \setminus U \cup U \setminus V_1$ ). On the other hand, both  $s$  and  $s_1$  are undefined on  $V_1$ , whence  $s_1 = s$ .

The remainder of the proof for  $\text{OPT}(P_1, P_2, F)$  follows from the induction hypothesis and the arguments used in the cases of FILTER and JOIN.

This completes the proof of Theorem 7.  $\square$

*Example 9.* Consider the following query taken from the official SPARQL document<sup>7</sup> ('find the names of people who do not know anyone'):

```
SELECT ?name WHERE {
  ?x foaf:givenName ?name.
  OPTIONAL { ?x foaf:knows ?who }.
  FILTER (! BOUND (?who))
},
```

which is represented by the following graph pattern

$$\text{FILTER}(\text{OPT}(\{ ?x \text{ foaf:givenName } ?name \}, \{ ?x \text{ foaf:knows } ?who \}, \top), \neg \text{bound}(?who)). \quad (2)$$

For the translation of the OPT operator, we first require to translate

$$\text{JOIN}(\{ ?x \text{ foaf:givenName } ?name \}, \{ ?x \text{ foaf:knows } ?who \}),$$

which results in the following relational expression (we remove trivial projections and filters):

$$\begin{aligned} & \pi_{x,name} \sigma_{\text{pred}=\text{foaf:givenName}} \rho_{x/\text{subj}} \rho_{\text{name}/\text{obj}} \text{triple} \bowtie \\ & \pi_{x,who} \sigma_{\text{pred}=\text{foaf:knows}} \rho_{x/\text{subj}} \rho_{\text{who}/\text{obj}} \text{triple} \cup \\ & \mu_x (\pi_{\text{name}} \sigma_{\text{isNull}(x)} \pi_{x,name} \sigma_{\text{pred}=\text{foaf:givenName}} \rho_{x/\text{subj}} \rho_{\text{name}/\text{obj}} \text{triple} \bowtie \\ & \pi_{x,who} \sigma_{\text{pred}=\text{foaf:knows}} \rho_{x/\text{subj}} \rho_{\text{who}/\text{obj}} \text{triple}) \cup \\ & \mu_x (\pi_{x,name} \sigma_{\text{pred}=\text{foaf:givenName}} \rho_{x/\text{subj}} \rho_{\text{name}/\text{obj}} \text{triple} \bowtie \\ & \pi_{\text{who}} \sigma_{\text{isNull}(x)} \pi_{x,who} \sigma_{\text{pred}=\text{foaf:knows}} \rho_{x/\text{subj}} \rho_{\text{who}/\text{obj}} \text{triple}). \end{aligned}$$

Since *triple* contains no *nulls*, the above relational expression is clearly equivalent to (cf. Theorem 8):

$$Q = \pi_{x,name} \sigma_{\text{pred}=\text{foaf:givenName}} \rho_{x/\text{subj}} \rho_{\text{name}/\text{obj}} \text{triple} \bowtie \pi_{x,who} \sigma_{\text{pred}=\text{foaf:knows}} \rho_{x/\text{subj}} \rho_{\text{who}/\text{obj}} \text{triple}$$

Then  $\text{OPT}(\{ ?x \text{ foaf:givenName } ?name \}, \{ ?x \text{ foaf:knows } ?who \}, \top)$  is translated into the following relational expression:

$$Q \cup \mu_{\text{who}} (\pi_{x,name} \sigma_{\text{pred}=\text{foaf:givenName}} \rho_{x/\text{subj}} \rho_{\text{name}/\text{obj}} \text{triple} \setminus \pi_{x,name} Q).$$

It can be verified (cf. Theorem 8) that this expression is in fact equivalent to

$$\pi_{x,name} \sigma_{\text{pred}=\text{foaf:givenName}} \rho_{x/\text{subj}} \rho_{\text{name}/\text{obj}} \text{triple} \bowtie \pi_{x,who} \sigma_{\text{pred}=\text{foaf:knows}} \rho_{x/\text{subj}} \rho_{\text{who}/\text{obj}} \text{triple}.$$

<sup>7</sup> <http://www.w3.org/TR/sparql-features>

Finally, the filter expression in graph pattern (2) is translated into  $\neg \neg isNull(who)$ , that is  $isNull(who)$ , and the graph pattern itself to

$$\sigma_{isNull(who)} \left( \pi_{x, name} \sigma_{pred=foaf:givenName} \rho_{x/subj} \rho_{name/obj} \text{triple} \bowtie \pi_{x, who} \sigma_{pred=foaf:knows} \rho_{x/subj} \rho_{who/obj} \text{triple} \right).$$

whose projection onto  $\{x\}$  can also be expressed as follows:

$$\pi_x \rho_{x/subj} \sigma_{pred=foaf:givenName} \text{triple} \setminus \pi_x \rho_{x/subj} \sigma_{pred=foaf:knows} \text{triple}$$

(informally, find those individuals who do not know anyone).

## D Proof of Theorem 8

We begin by formalising the intuition behind the definition of  $\nu$ :

**Proposition 13.** *Let  $P$  be a graph pattern. Then, for any RDF graph  $G$  and any solution mapping  $s \in \llbracket P \rrbracket_G$ , we have  $var(P) \setminus \nu(P) \subseteq dom(s)$ .*

*Proof.* The proof is by induction on the structure of  $P$ . The basis of induction is immediate: all variables in any BGP  $P$  are in the domain of any  $s \in \llbracket P \rrbracket_G$ . For the induction step, we consider the cases of SPARQL algebra operations:

$P = \text{FILTER}(P_1, F)$ . If  $s \in \llbracket \text{FILTER}(P_1, F) \rrbracket_G$  then  $s \in \llbracket P_1 \rrbracket_G$  and so, by the induction hypothesis,  $var(P_1) \setminus \nu(P_1) \subseteq dom(s)$ , from which the claim follows.

$P = \text{BIND}(P_1, v, c)$ . If  $s \in \llbracket \text{BIND}(P_1, v, c) \rrbracket_G$  then  $v \in dom(s)$  and the restriction  $s'$  of  $s$  onto  $dom(s) \setminus \{v\}$  is in  $\llbracket P_1 \rrbracket_G$ . By the induction hypothesis,  $var(P_1) \setminus \nu(P_1) \subseteq dom(s')$ , whence the claim:  $(var(P_1) \cup \{v\}) \setminus \nu(P_1) \subseteq dom(s') \cup \{v\} = dom(s)$ .

$P = \text{UNION}(P_1, P_2)$ . If  $s \in \llbracket \text{UNION}(P_1, P_2) \rrbracket_G$  then either  $s \in \llbracket P_i \rrbracket_G$ , for  $i = 1$  or  $i = 2$ . By the induction hypothesis,  $var(P_i) \setminus \nu(P_i) \subseteq dom(s)$ , for  $i = 1$  or  $i = 2$ . Let  $v \in var(P_1)$  but  $v \notin (var(P_1) \setminus var(P_2)) \cup (var(P_2) \setminus var(P_1)) \cup \nu(P_1) \cup \nu(P_2)$ . It follows that  $v \in var(P_2)$  but  $v \notin \nu(P_1)$  and  $v \notin \nu(P_2)$ . So,  $v \in var(P_i) \setminus \nu(P_i)$ , for both  $i = 1$  and  $i = 2$ . By the mirror image argument, if  $v \in var(P_2)$  then  $v \in var(P_i) \setminus \nu(P_i)$ , for both  $i = 1$  and  $i = 2$ . Thus,  $v \in dom(s)$ .

$P = \text{JOIN}(P_1, P_2)$ . If  $s \in \llbracket \text{JOIN}(P_1, P_2) \rrbracket_G$  then there are  $s_1 \in \llbracket P_1 \rrbracket_G$  and  $s_2 \in \llbracket P_2 \rrbracket_G$  such that  $s_1$  and  $s_2$  are compatible and  $s = s_1 \oplus s_2$ . By the induction hypothesis,  $var(P_i) \setminus \nu(P_i) \subseteq dom(s_i)$ . Let  $v \in var(P_i)$ , for either  $i = 1$  or  $i = 2$  but  $v \notin \nu(P_1) \cup \nu(P_2)$ . Clearly,  $v \in dom(s_i)$  and so, in either case  $v \in dom(s)$ .

$P = \text{OPT}(P_1, P_2, F)$ . If  $s \in \llbracket \text{OPT}(P_1, P_2, F) \rrbracket_G$  then either there is  $s_1 \in \llbracket P_1 \rrbracket_G$  and  $s_2 \in \llbracket P_2 \rrbracket_G$  such that  $s_1$  and  $s_2$  are compatible,  $s = s_1 \oplus s_2$  and  $F^s = \top$  or  $s \in \llbracket P_1 \rrbracket_G$  and, for all  $s_2 \in \llbracket P_2 \rrbracket_G$  either  $s$  and  $s_2$  are incompatible or  $F^{s \oplus s_2} \neq \top$ . Let  $v \in var(P_1)$  but  $v \notin \nu(P_1) \cup var(P_2)$ . By the induction hypothesis, for the two options above, we have  $v \in dom(s_1) \subseteq dom(s)$  and  $v \in dom(s)$ , respectively. The choice  $v \in var(P_2)$  but  $v \notin \nu(P_1) \cup var(P_2)$  is impossible.

This completes the proof of Proposition 13.  $\square$

**Theorem 8.** *If  $\text{var}(P_1) \cap \text{var}(P_2) \cap (\nu(P_1) \cup \nu(P_2)) = \emptyset$  then we can define*

$$\begin{aligned}\tau(\text{JOIN}(P_1, P_2)) &= \tau(P_1) \bowtie \tau(P_2), \\ \tau(\text{OPT}(P_1, P_2, F)) &= \tau(P_1) \bowtie_{\tau(F)} \tau(P_2),\end{aligned}$$

where  $R_1 \bowtie_F R_2 = \sigma_F(R_1 \bowtie R_2) \cup \mu_{U_2 \setminus U_1}(R_1 \setminus \pi_{U_1}(\sigma_F(R_1 \bowtie R_2)))$ , for relations  $R_1$  and  $R_2$  over  $U_1$  and  $U_2$ , respectively.

*Proof.* We clearly have  $\|\tau(P_1) \bowtie \tau(P_2)\|_{\text{triple}(G)} \subseteq \|\tau(\text{JOIN}(P_1, P_2))\|_{\text{triple}(G)}$  because  $\tau(P_1) \bowtie \tau(P_2)$  is a component of the union in the definition of  $\tau(\text{JOIN}(P_1, P_2))$ , with  $V_1 = V_2 = \emptyset$ . For the converse inclusion, consider a component of the union in the definition of  $\tau(\text{JOIN}(P_1, P_2))$ :

$$\left( (\pi_{U_1 \setminus V_1}(\sigma_{\text{isNull}(V_1)} \tau(P_1))) \right) \bowtie \left( (\pi_{U_2 \setminus V_2}(\sigma_{\text{isNull}(V_2)} \tau(P_2))) \right).$$

By Proposition 13 and Theorem 7, if  $V_1 \cap \nu(P_1) \neq \emptyset$  then  $P_1$  contains a variable that is always bound and so,  $\|\sigma_{\text{isNull}(V_1)} \tau(P_1)\|_{\text{triple}(G)} = \emptyset$ , for any RDF graph  $G$ . Therefore, in this case the component is empty and can be removed from the union. If the condition in the theorem is satisfied then only  $V_1 = \emptyset$  and  $V_2 = \emptyset$  will give rise to a possibly non-empty component. Thus,  $\|\tau(\text{JOIN}(P_1, P_2))\|_{\text{triple}(G)} \subseteq \|\tau(P_1) \bowtie \tau(P_2)\|_{\text{triple}(G)}$ .

The claim for  $\text{OPT}(P_1, P_2, F)$  is then immediate from the claim for  $\text{JOIN}(P_1, P_2)$  and the definition of the left join relational operation.  $\square$

## E Proof of Proposition 11 and Theorem 12

**Proposition 11.** *For any R2RML mapping  $\mathcal{M}$  and data instance  $D$ ,  $t \in \|\text{tr}_{\mathcal{M}}(\text{triple})\|_D$  if and only if  $t \in \text{triple}(G_{D, \mathcal{M}})$ .*

*Proof.* ( $\Rightarrow$ ) Let  $t \in \|\text{tr}_{\mathcal{M}}(\text{triple})\|_D$  then there is  $m$  in  $\mathcal{M}$  such that  $t \in \|\text{tr}_m\|_D$ . That is, the logical table of  $m$  matches the selection of  $\text{tr}_m$  (minus the  $\neg\text{isNull}(v_i)$  operators) and the term maps (subject, predicate and object) of  $m$  match the *subj*, *pred*, and *obj* projections of  $\text{tr}_m$ . It follows that, by the procedure described in [10, Section 11],  $t$  is part of the *generated triples* of  $m$  and, therefore, belongs to  $\text{triple}(G_{D, \mathcal{M}})$ .

( $\Leftarrow$ ) If a triple  $t$  is produced by  $\mathcal{M}$ , then there is triple map  $m$  with a predicate object map  $po$  that produces it by the procedure in [10, Section 11]. If this is the case, the logical table of  $m$  returns a tuple  $s$ , for which the values of the referenced columns in the term maps of  $m$  are not *null* and that generates  $t$ . By construction,  $m$  gives rise to  $\text{tr}_m$  in  $\text{tr}_{\mathcal{M}}(\text{triple})$  whose selection is the logical table of  $m$  and its projection matches the term maps of  $m$ . Thus,  $s$  produces  $t$  in  $\|\text{tr}_m\|_D$  and so, in  $\|\text{tr}_{\mathcal{M}}(\text{triple})\|_D$ .

**Theorem 12.** *For any graph pattern  $P$ , R2RML mapping  $\mathcal{M}$  and data instance  $D$ ,*

$$\|P\|_{G_{D, \mathcal{M}}} = \|\text{tr}_{\mathcal{M}}(\tau(P))\|_D.$$

*Proof.* Follows from Theorem 7 and Proposition 11.  $\square$

Q	LUBM <sub>1</sub>		LUBM <sub>9</sub>		LUBM <sub>20</sub>		LUBM <sub>50</sub>		LUBM <sub>100</sub>		LUBM <sub>200</sub>		LUBM <sub>500</sub>	
	O	OB <sub>H</sub>	OB <sub>P</sub>	P	O	OB <sub>H</sub>	P	O	P	O	P	O	P	O
$q_1$	2	8	29	1	3	97	1	3	1	3	1	3	1	3
$q_2$	2	25	11137	19	3	2531	256	5	633	12	11312	16	30593	36
$q_3$	1	6	86	9	2	78	158	2	345	3	932	2	2087	63
$q_4$	13	7	19	14	15	44	164	15	347	22	943	27	2093	24
$q_5$	16	12	4451	10	22	98	158	16	343	21	945	32	2182	28
$q_6$	455	27	32	21	5076	411	317	11610	730	28674	4321	58968	10781	123578
$q_7$	5	21	34005	10	6	429	157	6	343	32	939	8	2171	8
$q_8$	726	195	95875	80	760	917	192	748	374	755	974	796	2131	820
$q_9$	60	972	168978	78	668	189126	857	1594	5093	3734	5175	7466	12125	15227
$q_{10}$	2	6	126	9	3	97	158	2	344	2	930	2	2134	3
$q_{11}$	4	5	58	10	6	43	160	6	349	8	939	11	2093	18
$q_{12}$	3	4	19	15	4	70	236	4	567	4	937	3	2114	5
$q_{13}$	6	4	67	8	7	40	157	9	348	11	935	14	2657	38
$q_{14}$	91	20	24	15	1168	329	287	2782	694	6444	1866	13524	4457	29512
$q'_1$	93	58	190	46	99	98	767	98	715	91	2672	92	4422	95
$q'_4$	108	21	35	63	122	72	719	117	3394	140	3994	115	9179	108
$q'_6$	257	716	91855	174	4686	40575	1385	11659	3415	26418	8932	54092	19945	115110
$q'_9$	557	951	65916	102	6093	178401	1214	14934	3038	34178	9296	67123	19705	151376
$q_2^{avg}$	150	30	57141	29	9992	520	348	8232	939	19486	2607	39477	5411	79351
$q_6^{avg}$	6	7	241	25	31	40	273	29	735	32	1699	7	3969	7
$q_{10}^{avg}$	641	760	31269	253	6998	149191	2258	17106	5440	43006	8642	163308	17929	174362
start up time	3.1s	13.6s	7.7s	3.6s	3.1s	80m33s	18s	3.1s	39.8s	3.1s	1m38s	3.1s	3m23s	3.1s
loading time	10s	n/a	n/a	n/a	15s	n/a	n/a	26s	n/a	n/a	n/a	1m56s	n/a	3m35s
														10m17s

**Table 2.** Start up time, data loading time (in s) and query execution time (in ms) for LUBM in *Ontop* (O), OWL-BGP +Hermit (OB<sub>H</sub>), OWL-BGP +Pellet (OB<sub>P</sub>) and Pellet (P); OWL-BGP +Pellet times out on LUBM<sub>9</sub>, OWL-BGP +Hermit on LUBM<sub>20</sub> and Pellet on LUBM<sub>200</sub>.



## F SQL Schema for LUBM

We evaluated the performance of *Ontop* using the OWL 2 QL version of the Lehigh University Benchmark LUBM [16]. We experimented with the data instances LUBM<sub>*n*</sub>, for  $n = 1, 9, 20, 50, 100, 200, 500$  (where  $n$  specifies the number of universities; note that LUBM<sub>1</sub> and LUBM<sub>9</sub> were used for experiments in [19]). The results of the experiments are shown in Table 2.

The LUBM data generator creates an OWL file with class and property assertions. As *Ontop* has been designed to work with relational databases, we had to store the class and property assertions in a database. The simplest solution to this issue would be to use a generic schema, as usually done when storing triples in RDBMS backends. The two most common examples of this are (a) a schema with a single relation containing three attributes: *subj*, *pred*, *obj*, or (b) a schema with one unary relation for each class and one binary relation for each property. Such generic schemas, however, do not allow for efficient SQL translations in *Ontop* (or any other SPARQL-to-SQL system) because, for example, they require multiple and expensive (self)-join operations (if multiple properties are needed for an individual) and cause exponential blowup (due to class and property hierarchies).

In order to obtain efficient SQL queries, it is necessary that the schema follows standard best practices in the DB schema design, for example, normalisation. Usually, a normalised schema is obtained through a top-down approach that starts with the design of a conceptual model and follows by a translation of the conceptual model into a relational model. In the case of LUBM, we were not able to do so because the data generator was fixed. Instead, we studied the specification of the generator, extracted a model out of it, and created a schema based on that model. In particular, we identified disjoint classes (such as `ub:UndergraduateStudent` and `ub:GraduateStudent`) and functional properties (such as `ub:name`). When creating the schema we used the following principles:

- Classes on the top levels of hierarchies (e.g., `ub:Student`) have their own relations (e.g., `students`).
- The class membership in hierarchies is encoded using discriminating attributes (e.g., instances of `ub:UndergraduateStudent` and `ub:GraduateStudent` are also stored in the relation `students` for their superclass `ub:Student`, but are distinguished by the value of the discriminating attribute `stype`).
- Each functional property is included in the relation for its domain (e.g., property `ub:name` is the attribute `name` in `students`).
- Each 1-to-N and N-to-N property together with its attributes has a separate relation (e.g., relation `takescourses` represents `ub:takesCourse`).

The resulting schema consists of eleven relations:

- `coauthors` (data about the authors of a publication),
- `courses` (data about courses and the teachers assigned to those courses),
- `departments` (university departments),
- `heads` (heads of departments),
- `publications` (publications and the main author of a publication),

- ra (research assistants),
- researchgroups (research groups),
- students (data about students including their degrees and supervisors),
- ta (data about teaching assistants and courses they teach),
- takescourses (data about students and courses they take),
- teachers (data about teachers and their departments).

Instead of storing complete URIs, which were generated automatically by the LUBM data generator following a certain pattern, we store their components separately. For example, URIs of instances of `ub:GraduateStudent` are of the form

`http://www.Department12.University54.edu/GraduateStudent22`

where 54 refers to the university, 12 to the department in the university, and 22 to the graduate student in that department. We extracted those IDs (54, 12, 22 for this instance) and stored them in the respective attributes of the relations in the database (`uniid`, `depid`, `id` of the relation `students`, respectively). The obtained attributes together with the class names uniquely identify individuals, and so they form *primary keys* of the respective relations (e.g., attributes `depid`, `uniid`, `stype`, `id` constitute the primary key for relation `students`,<sup>8</sup> note that the discriminating attribute `stype` encodes the class name). *Foreign keys* are defined when the relation contained attributes that referred to the IDs of entities stored in a different relation.

In addition to the indexes that are defined by default on primary keys, we added indexes on attributes that would likely participate in join operations between relations. These were, mostly, the attributes that store IDs of entities from different relations. In total, we defined 39 additional indexes out of which 7 are composite. Note that in DB tuning, it is usual that indexes are defined with respect to query workload. Since we did not proceed in this way, there are indexes that could be added to obtain a better performance for some of the queries we evaluated, and some indexes that could be removed since they are not relevant for the workload of the evaluation.

Finally, R2RML mappings were defined so that the (virtual) RDF graph entailed by the mappings would consist of all the triples that were initially used to populate the database.

To illustrate our rationale in more detail, we consider the case of `ub:Student` and its two disjoint subclasses, `ub:GraduateStudent` and `ub:UndergraduateStudent`. The class `ub:Student` is what is known as an abstract class in ER modelling, that is, a class that has no instances; only `ub:GraduateStudent` and `ub:UndergraduateStudent` have instances. In addition, each `ub:Student` instance has exactly one value for the properties `ub:name`, `ub:telephone`, `ub:degreeFrom`, `ub:emailAddress` and `ub:advisor` (note that these properties are identified as functional on the basis of the specification of the data generator rather than the ontology). So, we defined the relation `students` shown in Fig. 1. Indexes in this relation include the (composite) index of the primary key (`depid`, `uniid`, `stype`, `id`) as well as indexes on the attributes `degreeuniid`, `advisorstype`, `advisorid`. In addition, it contains a composite foreign key on the pair `depid`, `uniid` referring to the attributes `id`, `universityid` in the relation `departments`.

<sup>8</sup> In the simplified example in Section 3.3 we do not have `depid` and `uniid`.

```

CREATE TABLE 'students' (
  'depid' smallint(6) NOT NULL,
  'uniid' smallint(6) NOT NULL,
  'stype' tinyint(4) NOT NULL,
  'id' smallint(6) NOT NULL,
  'name' varchar(45) DEFAULT NULL,
  'degreeuniid' smallint(6) DEFAULT NULL,
  'email' varchar(255) DEFAULT NULL,
  'phone' varchar(255) DEFAULT NULL,
  'advisorstype' tinyint(4) DEFAULT NULL,
  'advisorid' smallint(6) DEFAULT NULL,
  PRIMARY KEY ('depid', 'uniid', 'stype', 'id'),
  INDEX 'idx_stud_1' ('degreeuniid'),
  INDEX 'idx_stud_2' ('advisorstype'),
  INDEX 'idx_stud_3' ('advisorid'),
  CONSTRAINT 'fk_students_1' FOREIGN KEY ('depid', 'uniid')
    REFERENCES 'departments' ('departmentid', 'universityid')
);

```

**Fig. 1.** Relations for the subclasses `ub:UndergraduateStudent` and `ub:GraduateStudent` of `ub:Student`.

The relation `students` is mapped to the LUBM classes and properties using R2RML triple maps such as the one presented in Fig. 2. The mapping generates all triples in which the subject is a URI of an undergraduate student (indicated by the value 0 of the discriminating attribute `stype`).

```

@prefix rr: <http://www.w3.org/ns/r2rml#> .
@prefix ub: <http://swat.cse.lehigh.edu/onto/univ-bench.owl#> .

[ a rr:TriplesMap ;
  rr:logicalTable [ a rr:R2RMLView ;
    rr:sqlQuery "select * from students where stype=0"
  ] ;
  rr:subjectMap [ a rr:TermMap, rr:SubjectMap ;
    rr:class ub:UndergraduateStudent ;
    rr:termType rr:IRI ;
    rr:template
      "http://www.Department{depid}.University{uniid}.edu/UndergraduateStudent{id}"
  ] ;
  rr:predicateObjectMap [ a rr:PredicateObjectMap ;
    rr:predicate ub:telephone ;
    rr:objectMap [ a rr:TermMap, rr:ObjectMap ;
      rr:termType rr:Literal ;
      rr:column "phone"
    ]
  ] ;
  rr:predicateObjectMap [ a rr:PredicateObjectMap ;
    rr:predicate ub:emailAddress ;
    rr:objectMap [ a rr:TermMap, rr:ObjectMap ;
      rr:termType rr:Literal ;
      rr:column "email"
    ]
  ] ;
  rr:predicateObjectMap [ a rr:PredicateObjectMap ;
    rr:predicate ub:memberOf ;
    rr:objectMap [ a rr:TermMap, rr:ObjectMap ;
      rr:termType rr:IRI ;
      rr:template
        "http://www.Department{depid}.University{uniid}.edu"
    ]
  ]
] .

```

**Fig. 2.** R2RML mapping for instances of `ub:UndergraduateStudent`.