



# Reformulation-based query answering for RDF graphs with RDFS ontologies

Maxime Buron, François Goasdoué, Ioana Manolescu, Marie-Laure Mugnier

► **To cite this version:**

Maxime Buron, François Goasdoué, Ioana Manolescu, Marie-Laure Mugnier. Reformulation-based query answering for RDF graphs with RDFS ontologies. ESWC: European Semantic Web Conference, Jun 2019, Portoroz, Slovenia. hal-02051413

**HAL Id: hal-02051413**

**<https://hal.archives-ouvertes.fr/hal-02051413>**

Submitted on 12 Mar 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Reformulation-based query answering for RDF graphs with RDFS ontologies

Maxime Buron<sup>1,2</sup>, François Goasdoué<sup>3</sup>, Ioana Manolescu<sup>1,2</sup>, and Marie-Laure Mugnier<sup>4</sup>

<sup>1</sup> Inria Saclay, `firstname.lastname@inria.fr`

<sup>2</sup> LIX (UMR 7161, CNRS and Ecole polytechnique, France)

<sup>3</sup> Univ Rennes, CNRS, IRISA, `fg@irisa.fr`

<sup>4</sup> Univ. Montpellier, LIRMM, Inria, `mugnier@lirmm.fr`

**Abstract.** Query answering in RDF knowledge bases has traditionally been performed either through graph saturation, i.e., adding all implicit triples to the graph, or through query reformulation, i.e., modifying the query to look for the explicit triples entailing precisely what the original query asks for. The most expressive fragment of RDF for which Reformulation-based query answering exists is the so-called database fragment [12], in which implicit triples are restricted to those entailed using an RDFS ontology. Within this fragment, query answering was so far limited to the interrogation of data triples (non-RDFS ones); however, a powerful feature specific to RDF is the ability to query data and schema triples together. In this paper, we address the general query answering problem by reducing it, through a pre-query reformulation step, to that solved by the query reformulation technique of [12]. We also report on experiments demonstrating the low cost of our reformulation algorithm.

**Keywords:** Query answering · Query reformulation · RDF · RDFS

## 1 Introduction

RDF is the standard model for sharing data and knowledge bases. The rapid increase in number and size of RDF graphs makes efficient query answering on RDF quite a challenging task. *Reasoning* raises a performance challenge: query answering on an RDF graph no longer reduces to *evaluating* the query on the graph (by finding all the homomorphisms, or embeddings, of the query in the graph). Instead, it requires taking into account also the possible ontology (or knowledge) rules, which specify how different classes and properties of an RDF graph relate to each other, and may lead to query answers that evaluation alone cannot compute. Moreover, SPARQL, the standard query language of RDF, allows *querying the data and the ontology together*. This is a radical departure both from relational databases, and from Description Logics (DL)-style models for RDF data and queries.

For what concerns reasoning, two main methods have been explored: graph saturation, which injects the ontology knowledge into the graph, and query reformulation, which pushes it into the query. Saturation adds to the graph all the triples it entails through the ontology. Evaluating a query on a saturated graph

<b>RDF assertions</b>	Triple notation
Class assertion	$(\mathbf{s}, \tau, \mathbf{o})$
Property assertion	$(\mathbf{s}, \mathbf{p}, \mathbf{o})$ with $\mathbf{p} \notin \{\tau, \prec_{sc}, \prec_{sp}, \leftrightarrow_d, \leftrightarrow_r\}$
<b>RDFS constraints</b>	Triple notation
Subclass	$(\mathbf{s}, \prec_{sc}, \mathbf{o})$
Subproperty	$(\mathbf{s}, \prec_{sp}, \mathbf{o})$
Domain typing	$(\mathbf{s}, \leftrightarrow_d, \mathbf{o})$
Range typing	$(\mathbf{s}, \leftrightarrow_r, \mathbf{o})$

**Table 1.** RDF statements.

can be quite efficient; however, saturation takes time to compute, space to store, and needs to be updated when the data and/or ontology rules change. Reformulation leaves the graph unchanged and builds a reformulated query which, evaluated on the original graph, computes all the answers, including those that hold due to entailed triples. Each query reformulation method, thus, targets a certain ontology language and a query dialect. The most expressive RDF fragment for which sound and complete reformulation-based query answering exists is the so-called database fragment [12], in which RDF Schema (RDFS, in short) is used to describe the ontology, while queries only carry over the data triples. In this work, we present a novel *reformulation-based query answering* under *RDFS* ontologies for *Basic Graph Pattern (BGP) queries over both the data and the ontology*. This goes beyond the closest algorithm previously known [12] which is restricted to queries over the data only (not over the ontology). The algorithm we present here also goes beyond those of RDF platforms such as Jena, Virtuoso or Stardog, which we found experimentally to be incomplete when answering through reformulation queries over the data and the ontology of an RDF graph. Below, we recall some terminology (Section 2) and discuss the state of the art (Section 3). Then, Section 4 introduces our novel query reformulation algorithm, which we implemented in the platform used in [12, 9], leveraging an efficient relational database (RDBMS) engine for query answering. Our experiments (Section 5) demonstrate the practical interest of our reformulation approach.

## 2 Preliminaries

We present the basics of the RDF graph data model (Section 2.1), of RDF entailment used to make explicit the implicit information RDF graphs encode (Section 2.2), as well as how they can be queried using the widely-considered SPARQL Basic Graph Pattern queries (Section 2.3).

### 2.1 RDF Graph

We consider three pairwise disjoint sets of values:  $\mathcal{I}$  of IRIs (resource identifiers),  $\mathcal{L}$  of literals (constants) and  $\mathcal{B}$  of blank nodes modeling unknown IRIs or literals, a.k.a. to *labelled nulls* [4]. A *well-formed triple* belongs to  $(\mathcal{I} \cup \mathcal{B}) \times \mathcal{I} \times (\mathcal{L} \cup \mathcal{I} \cup \mathcal{B})$ , and an *RDF graph*  $G$  is a set of well-formed triples. A triple  $(\mathbf{s}, \mathbf{p}, \mathbf{o})$  states that its *subject*  $\mathbf{s}$  has the *property*  $\mathbf{p}$  with the *object* value  $\mathbf{o}$  [1]. We denote by  $\text{Val}(G)$  the set of all values (IRIs, blank nodes and literals) occurring in an RDF graph  $G$ , and by  $\text{Bl}(G)$  its set of blank nodes.

Within an RDF graph, triples model either factual *assertions* for unary relations called *classes* and binary relations called *properties*, or *RDFS ontological constraints* between classes and properties. The RDFS constraints are of four flavours: **subclass** constraints, **subproperty** constraints, typing of the **domain** (first attribute) or of the **range** (second attribute) of a property. The triple notations we adopt for RDF assertions and constraints are shown in Table 1. In a triple, we use  $_:b$  (possibly with indices) to denote blank nodes, and strings between quotes to denote literals.

We consider RDF graphs with *RDFS ontologies*, i.e., constraints of the four flavors above, excluding constraints which would alter the commonly-accepted RDFS semantics. For instance,  $(\leftrightarrow_d, \prec_{sp}, \leftrightarrow_r)$  is not allowed as it would make domain typing a particular case of range typing. Let  $\text{RDFS}(G)$  denote the RDFS constraints of a graph  $G$ . We define:

**Definition 1 (RDF graph with an RDFS ontology).** *An RDFS ontology (or ontology in short) is a set of RDFS constraints, whose subjects and objects are either IRIs (other than  $\prec_{sc}, \prec_{sp}, \leftrightarrow_d, \leftrightarrow_r, \tau$ ) or blank nodes. An RDF graph  $G$  with ontology  $O$  is such that:  $\text{RDFS}(G) = O$ .*

*Example 1 (Running example).* Consider the following RDF graph:

$$G_{\text{ex}} = \{(\text{:worksFor}, \leftrightarrow_d, \text{:Person}), (\text{:worksFor}, \leftrightarrow_r, \text{:Org}), (\text{:PubAdmin}, \prec_{sc}, \text{:Org}), (\text{:Comp}, \prec_{sc}, \text{:Org}), (:_b_C, \prec_{sc}, \text{:Comp}), (\text{:hiredBy}, \prec_{sp}, \text{:worksFor}), (\text{:ceoOf}, \prec_{sp}, \text{:worksFor}), (\text{:ceoOf}, \leftrightarrow_r, \text{:Comp}), (\text{:p}_1, \text{:ceoOf}, \text{:c}), (\text{:c}, \tau, \text{:}_b_C), (\text{:p}_2, \text{:hiredBy}, \text{:a}), (\text{:a}, \tau, \text{:PubAdmin})\}$$

The ontology of  $G_{\text{ex}}$ , i.e., the first eight triples, states that persons are working for organizations, some of which are public administrations or companies. Further, there exists a special kind of company (modeled by  $:_b_C$ ). Being hired by or being CEO of an organization are two ways of working for it; in the latter case, this organization is a company. The assertions of  $G_{\text{ex}}$ , i.e., the four remaining triples, states that  $\text{:p}_1$  is CEO of  $\text{:c}$ , which is a company of the special kind  $:_b_C$ , and  $\text{:p}_2$  is hired by the public administration  $\text{:a}$ .

A *homomorphism between RDF graphs* allows characterizing whether an RDF graph *simply entails* another, based on their explicit triples only:

**Definition 2 (RDF graph homomorphism).** *Let  $G$  and  $G'$  be two RDF graphs. A homomorphism from  $G$  to  $G'$  is a function  $\varphi$  from  $\text{Val}(G)$  to  $\text{Val}(G')$ , which is the identity on IRIs and literals, such that for any triple  $(s, p, o)$  in  $G$ , the triple  $(\varphi(s), \varphi(p), \varphi(o))$  is in  $G'$ .*

Note that, according to the previous definition, a  $G$  blank node can be mapped to any  $G'$  value. A graph  $G'$  simply entails a graph  $G$  if there is a homomorphism  $\varphi$  from  $G$  to  $G'$ , which we denote by  $G' \models^\varphi G$ .

## 2.2 RDF Entailment Rules

The semantics of an RDF graph consists of the explicit triples it contains, and of the implicit triples that can be derived from it using *RDF entailment rules*.

**Definition 3 (RDF entailment rule).** *An RDF entailment rule  $r$  has the form  $\text{body}(r) \rightarrow \text{head}(r)$ , where  $\text{body}(r)$  and  $\text{head}(r)$  are RDF graphs, respectively called *body* and *head* of the rule  $r$ .*

Rule [2]	Entailment rule
<b>rdfs2</b>	$(p, \leftarrow_d, o), (s_1, p, o_1) \rightarrow (s_1, \tau, o)$
<b>rdfs3</b>	$(p, \leftarrow_r, o), (s_1, p, o_1) \rightarrow (o_1, \tau, o)$
<b>rdfs5</b>	$(p_1, \prec_{sp}, p_2), (p_2, \prec_{sp}, p_3) \rightarrow (p_1, \prec_{sp}, p_3)$
<b>rdfs7</b>	$(p_1, \prec_{sp}, p_2), (s, p_1, o) \rightarrow (s, p_2, o)$
<b>rdfs9</b>	$(s, \prec_{sc}, o), (s_1, \tau, s) \rightarrow (s_1, \tau, o)$
<b>rdfs11</b>	$(s, \prec_{sc}, o), (o, \prec_{sc}, o_1) \rightarrow (s, \prec_{sc}, o_1)$
<b>ext1</b>	$(p, \leftarrow_d, o), (o, \prec_{sc}, o_1) \rightarrow (p, \leftarrow_d, o_1)$
<b>ext2</b>	$(p, \leftarrow_r, o), (o, \prec_{sc}, o_1) \rightarrow (p, \leftarrow_r, o_1)$
<b>ext3</b>	$(p, \prec_{sp}, p_1), (p_1, \leftarrow_d, o) \rightarrow (p, \leftarrow_d, o)$
<b>ext4</b>	$(p, \prec_{sp}, p_1), (p_1, \leftarrow_r, o) \rightarrow (p, \leftarrow_r, o)$

**Table 2.** RDFS entailment rules.

The standard RDF entailment rules are defined in [2]. In this work, we consider the rules shown in Table 2, which we call *RDFS entailment rules*; all values except the  $\tau, \prec_{sc}, \prec_{sp}, \leftarrow_d, \leftarrow_r$  properties are blank nodes. These rules are the most frequently used for RDFS entailment; they produce implicit triples by exploiting the RDFS ontological constraints of an RDF graph. For example, the rule **rdfs9**, which propagates values from subclasses to their superclasses, is defined by  $\text{body}(\mathbf{rdfs9}) = \{(s, \prec_{sc}, o), (s_1, \tau, s)\}$  and  $\text{head}(\mathbf{rdfs9}) = \{(s_1, \tau, o)\}$ . The *direct entailment* of an RDF graph  $G$  with a set of RDF entailment rules  $\mathcal{R}$ , denoted by  $C_{G, \mathcal{R}}$ , characterizes the set of implicit triples resulting from rule applications that use solely the explicit triples of  $G$ . It is defined as:

$$C_{G, \mathcal{R}} = \{\varphi(\text{head}(r)) \mid r \in \mathcal{R}, G \models^\varphi \text{body}(r)\}$$

For instance, the rule **rdfs9** applies to the graph  $G_{\text{ex}}$ :  $G_{\text{ex}} \models^\varphi \text{body}(\mathbf{rdfs9})$  through the homomorphism  $\varphi$  defined as  $\{s \mapsto :b_C, o \mapsto :Comp, s_1 \mapsto :c\}$ , hence allows deriving the implicit triple  $(:c, \tau, :Comp)$ .

The *saturation* of an RDF graph allows materializing its semantics, by iteratively augmenting it with the triples it entails using a set  $\mathcal{R}$  of RDF entailment rules, until reaching a fixpoint; this process is finite [2]. Formally:

**Definition 4 (RDF graph saturation).** *Let  $G$  be an RDF graph and  $\mathcal{R}$  a set of entailment rules. We recursively define a sequence  $(G_i^{\mathcal{R}})_{i \in \mathbb{N}}$  of RDF graphs as follows:*

- $G_0^{\mathcal{R}} = G$ , and
- $G_{i+1}^{\mathcal{R}} = G_i^{\mathcal{R}} \cup C_{G_i^{\mathcal{R}}, \mathcal{R}}$  for  $0 \leq i$ .

The saturation of  $G$  w.r.t.  $\mathcal{R}$ , denoted by  $G^{\mathcal{R}}$ , is  $G_n^{\mathcal{R}}$  for  $n$  the smallest integer such that  $G_n^{\mathcal{R}} = G_{n+1}^{\mathcal{R}}$ .

*Example 2.* The saturation of  $G_{\text{ex}}$  w.r.t. the set  $\mathcal{R}$  of RDFS entailment rules shown in Table 2 is attained after the following *two* saturation steps:

$$\begin{aligned} (G_{\text{ex}})_1^{\mathcal{R}} &= G_{\text{ex}} \cup \{(:b_C, \prec_{sc}, :Org), (:hiredBy, \leftarrow_d, :Person), (:hiredBy, \leftarrow_r, :Org), \\ &\quad (:ceoOf, \leftarrow_d, :Person), (:ceoOf, \leftarrow_r, :Org), \\ &\quad (:p_1, :worksFor, :c), (:c, \tau, :Comp), (:p_2, :worksFor, :a), (:a, \tau, :Org)\} \\ (G_{\text{ex}})_2^{\mathcal{R}} &= (G_{\text{ex}})_1^{\mathcal{R}} \cup \{(:p_1, \tau, :Person), (:p_2, \tau, :Person), (:c, \tau, :Org)\} \end{aligned}$$

Simple entailment between RDF graphs, which is based on their explicit triples only, generalizes to *entailment between RDF graphs w.r.t. a set of RDF entailment rules*, to also take into account their implicit triples.

A graph  $G$  entails a graph  $G'$  w.r.t. a set of rules  $\mathcal{R}$ , noted  $G \models_{\mathcal{R}}^\varphi G'$ , whenever

there is a homomorphism  $\varphi$  from  $G'$  to  $G^{\mathcal{R}}$ . Of course, simple entailment and entailment between RDF graphs coincide when  $\mathcal{R} = \emptyset$ . For simplicity, we will just write  $G \models_{\mathcal{R}} G'$  whenever  $\varphi$  is not needed for the discussion.

*In this work, unless otherwise specified,  $\mathcal{R}$  denotes the rules from Table 2.*

### 2.3 Basic Graph Pattern Queries

A popular fragment of the SPARQL query language consists of conjunctive queries, also known as basic graph pattern queries. Let  $\mathcal{V}$  be a set of variable symbols, disjoint from  $\mathcal{I} \cup \mathcal{B} \cup \mathcal{L}$ . A *basic graph pattern* (BGP) is a set of *triple patterns* (triples in short) belonging to  $(\mathcal{I} \cup \mathcal{B} \cup \mathcal{V}) \times (\mathcal{I} \cup \mathcal{V}) \times (\mathcal{I} \cup \mathcal{B} \cup \mathcal{L} \cup \mathcal{V})$ . For a BGP  $P$ , we denote by  $\text{Var}(P)$  the set of variables occurring in  $P$ , by  $\text{Bl}(P)$  its set of blank nodes, and by  $\text{Val}(P)$  its set of values (IRIs, blank nodes, literals and variables).

**Definition 5 (BGP query).** *A BGP query (BGPQ)  $q$  is of the form  $q(\bar{x}) \leftarrow P$ , where  $P$  is a BGP also denoted by  $\text{body}(q)$  and  $\bar{x} \subseteq \text{Var}(P)$  is the set of  $q$ 's answer variables. The arity of  $q$  is that of  $\bar{x}$ , i.e.,  $|\bar{x}|$ .*

**Partially instantiated BGPQs** generalize BGPQs and have been used for reformulation-based query answering [12]. Starting from a BGPQ  $q$ , partial instantiation replaces *some* variables and/or blank nodes with values from  $\mathcal{I} \cup \mathcal{L} \cup \mathcal{B}$ , as specified by a substitution  $\sigma$ ; the partially instantiated query is denoted  $q_{\sigma}$ . Observe that when  $\sigma = \emptyset$ ,  $q_{\sigma}$  coincides with  $q$ . Further, due to  $\sigma$ , and in contrast with standard BGPQs, some answer variables of  $q_{\sigma}$  can be bound:

*Example 3.* Consider the BGPQ asking for *who is working for which kind of company*:  $q(x, y) \leftarrow (x, \text{:worksFor}, z), (z, \tau, y), (y, \prec_{sc}, \text{:Comp})$ , and the substitution  $\sigma = \{x \mapsto \text{:p}_1\}$ . The partially instantiated BGPQ  $q_{\sigma}$  corresponds to  $q(\text{:p}_1, y) \leftarrow (\text{:p}_1, \text{:worksFor}, z), (z, \tau, y), (y, \prec_{sc}, \text{:Comp})$ .

The semantics of a (partially instantiated) BGPQ on an RDF graph is defined through homomorphisms from the query body to the saturation of the queried graph. The homomorphisms needed here are a straightforward extension of RDF graph homomorphisms (Definition 2) to also take variables into account.

**Definition 6 ((Non-standard) BGP to RDF graph homomorphism).** *A homomorphism from a BGP  $q$  to an RDF graph  $G$  is a function  $\varphi$  from  $\text{Val}(\text{body}(q))$  to  $\text{Val}(G)$  such that for any triple  $(s, p, o) \in \text{body}(q)$ , the triple  $(\varphi(s), \varphi(p), \varphi(o))$  is in  $G$ . For a standard homomorphism, as per the SPARQL recommendation,  $\varphi$  is the identity on IRIs and literals; for a non standard one,  $\varphi$  is the identity on IRIs, literals and on blank nodes.*

We distinguish *query evaluation*, whose result is just based on the explicit triples of the graph, i.e., on BGP to RDF graph homomorphisms, from *query answering* that also accounts for the implicit graph triples, i.e., based on both BGP to RDF graph homomorphisms *and* RDF entailment. In this paper, we use two flavors of query evaluation and of query answering, which differ in relying either on standard or on non-standard BGP to RDF graph homomorphisms.

**Definition 7 ((Non-standard) evaluation and answering).** *Let  $q_{\sigma}$  be a partially instantiated BGPQ  $q_{\sigma}$  obtained from a BGPQ  $q$  and a substitution  $\sigma$ .*

The standard answer set to  $q_\sigma$  on an RDF graph  $G$  w.r.t. a set  $\mathcal{R}$  of RDF entailment rules is:

$$q_\sigma(G, \mathcal{R}) = \{\varphi(\bar{x}_\sigma) \mid G \models_{\mathcal{R}}^{\varphi} \text{body}(q)_\sigma\}$$

where  $\bar{x}_\sigma$  and  $\text{body}(q)_\sigma$  denote the result of replacing the variables and blank nodes in  $\bar{x}$  and  $\text{body}(q)$ , respectively, according to  $\sigma$ .

If  $\bar{x} = \emptyset$ ,  $q_\sigma$  is a Boolean query, in which case  $q_\sigma$  is false when  $q_\sigma(G, \mathcal{R}) = \emptyset$  and true when  $q_\sigma(G, \mathcal{R}) = \{\langle \rangle\}$ , i.e., the answer to  $q_\sigma$  is an empty tuple.

We call  $q_\sigma(G, \emptyset)$  the standard evaluation of  $q_\sigma$  on  $G$ , written  $q_\sigma(G)$  for short, which solely amounts to standard BGP to RDF graph homomorphism finding.

The non-standard answer set, denoted  $\overbrace{q_\sigma(G, \mathcal{R})}$ , and non-standard evaluation  $\overbrace{q_\sigma(G)}$  of  $q_\sigma$  on  $G$  w.r.t.  $\mathcal{R}$  only differ from the standard ones by using non-standard BGP to RDF graph homomorphisms.

These notions and notations naturally extend to *unions* of BGPQs.

*Example 4.* Consider again the BGPQs from the preceding example. Their standard evaluations on  $G_{\text{ex}}$  are empty because  $G_{\text{ex}}$  has no explicit `:worksFor` assertion, while their standard answer sets on  $G_{\text{ex}}$  w.r.t.  $\mathcal{R}$  are  $\{\langle :p_1, :b_C \rangle\}$  because `:p1` being CEO of `:c`, `:p1` implicitly works for it, and `:c` is explicitly a company of the particular unknown type `:b_C`.

Consider now the BGPQ  $q(x) \leftarrow (x, :worksFor, y), (y, \tau, :b_C)$ . Under standard query answering, it asks for *who is working for some kind of organization* and its answer set is  $\{\langle :p_1 \rangle, \langle :p_2 \rangle\}$ ; by contrast, under non-standard query answering, it asks for *who is working for an organization of the particular unknown type `:b_C`* in  $G_{\text{ex}}$  and its answer set is just  $\{\langle :p_1 \rangle\}$ .

### 3 Prior related work

Two main techniques for answering BGPQs on RDF graphs have been investigated in the literature.

*Saturation-based query answering.* This technique directly follows from the definition of query answers in the W3C's SPARQL recommendations [3], recalled in Section 2.3 for BGPQs. Indeed, it trivially follows from Definition 7 that

$q(G, \mathcal{R}) = q(G^{\mathcal{R}})$  (resp.  $\overbrace{q(G, \mathcal{R})} = \overbrace{q(G^{\mathcal{R}})}$ ), i.e., query answering reduces to query evaluation on the *saturated* RDF graph.

Saturation-based query answering is typically fast, because it only requires query evaluation, which can be efficiently performed by a data management engine. However, saturation takes time to be computed, requires extra space to be stored, and must be recomputed or maintained (e.g., [8, 7, 12]) upon updates. Many RDF data management systems use saturation-based query answering. They either allow computing graph saturation, e.g., Jena and RDFox, or simply assume that RDF graphs have been saturated before being stored, e.g., DB2RDF.

*Reformulation-based query answering.* This technique also reduces query answering to query evaluation, however, the reasoning needed to ensure complete answers is performed on the query and not on the RDF graph. A given query

$q$ , asked on an RDF graph  $G$  w.r.t.  $\mathcal{R}$  is *reformulated* into a query  $q'$  such that  $q(G, \mathcal{R}) = q'(G)$  or  $q(G, \mathcal{R}) = \overbrace{q'(G)}$  holds. Standard or non-standard query evaluation is needed on the reformulated query, depending on the considered RDF fragment: when blank nodes are allowed in RDFS constraints, non-standard evaluation is used [12], while standard evaluation is sufficient otherwise [5, 11]. Different SPARQL dialects have been adopted for BGPQ reformulation in more limited settings than the one considered in this paper, i.e., the database fragment of RDF and unrestricted BGPQs. *Unions of BGPQs (UBGPQs in short)* have been used in [5, 11, 12]. However, these works are restricted to input BGPQs that must be matched on *RDF assertions* only. BGPQs aiming at interrogating solely the RDFS ontology, or the ontology *and* the assertions are not considered, even though such joint querying is a major novelty of RDF and SPARQL. The techniques adopt unions of BGPQs [5] or of *partially instantiated* BGPQs [11, 12], depending on whether variables can be used in class and property positions in queries, e.g., whether a query triple  $(x, \tau, z)$  or  $(x, y, z)$  is allowed. Reformulation-based query answering in the DL fragment of RDF, which is strictly contained in the database fragment of RDF, has been investigated for relational conjunctive queries [5, 10], while the slight extension thereof considered in [6, 11, 13, 17] has been investigated for one-triple BGPQs [13, 17], BGPQs [11], and SPARQL queries [6]. In [6], SPARQL queries are reformulated into *nested* SPARQL, allowing nested regular expressions in property position in query triples. These reformulations allow sound and complete query answering on restricted RDF graphs with RDFS ontologies: these graph *must not contain blank nodes*. While such nested reformulations are more compact, the queries we produce are more practical, since their evaluation can be delegated to any off-the-shelf RDBMS, or to an RDF engine such as RDF-3X [16] even if it is unaware of reasoning; further, we do not impose restrictions on RDF graphs.

In Section 4, we devise a reformulation-based query answering technique for the entire database fragment of RDF and unrestricted BGPQs.

Reformulation-based query answering is well-suited to frequently updated RDF graphs, because it uses the queried RDF graph at query time (and not its saturation). However, reformulated queries tend to be more complex than the original ones, thus costly to evaluate. To mitigate this, [9] provides an *optimized reformulation framework* whereas an incoming BGPQ is reformulated into a *join of unions of BGPQs (JUBGPQ in short)*. This approach being based on a database-style *cost model*, JUBGPQ reformulations are very efficiently evaluated.

Some available RDF data management systems use reformulation-based query answering but return incomplete answer sets in the RDF setting we consider<sup>5</sup>, e.g., AllegroGraph<sup>6</sup> and Stardog<sup>7</sup> miss answers because they cannot evaluate

---

<sup>5</sup>See discussion at

<https://team.inria.fr/cedar/rdfs-reasoning-experiments/>.

<sup>6</sup><https://franz.com/agraph/support/documentation/current/reasoner-tutorial.html>

<sup>7</sup>[https://www.stardog.com/docs/#\\_owl\\_rule\\_reasoning](https://www.stardog.com/docs/#_owl_rule_reasoning)



triples with a variable property on the schema, while Virtuoso<sup>8</sup> only exploits subclass and subproperty constraints, but not domain and range ones.

Finally, *Hybrid* approaches have also been studied, e.g., in [18], where some one-triple queries are chosen for materialization and reused during reformulation-based answering.

## 4 Extending query reformulation to queries over the ontology

We now present the main contribution of this paper: a reformulation-based query answering (QA) technique able to compute all answers to a BGPQ against *all* the explicit and implicit triples of an RDF graph, i.e., its RDF assertions *and* RDFS constraints, as per the SPARQL and RDF recommendations [3, 2].

The central idea is to *reduce* this full QA problem to an *assertion-level* QA, i.e., where the query is confined to just the explicit and implicit RDF assertions. To this aim, we divide query reformulation in two steps: the first reformulation step implements the reduction, while the second step relies on the reformulation technique of [12], which considers assertion-level QA.

### 4.1 Overview of our query reformulation technique

Let us first notice that the body of any BGPQ  $q$  can be divided into three disjoint subsets of triples  $(s, p, o)$ , according to the nature of term  $p$ : the set  $b_c$  of RDFS triples where  $p$  is a built-in RDFS property ( $\prec_{sc}$ ,  $\prec_{sp}$ ,  $\leftrightarrow_d$ ,  $\leftrightarrow_r$ ); the set  $b_a$  of assertion triples where  $p$  is  $\tau$  or a user-defined property; and the set  $b_v$  where  $p$  is a variable. We denote by  $q_c$ ,  $q_a$  and  $q_v$  the subqueries respectively associated with these bodies. If  $b_v$  is not empty,  $q$  can be reformulated as a union of BGPQs, say  $\mathcal{Q}$ , composed of all BGPQs that can be obtained from  $q$  by substituting some (possibly none) variables occurring in  $q_v$  with one of the four built-in RDFS properties. We assume this preprocessing step to simplify the explanations, even if in practice it may not be performed. Then, the answers to any BGPQ  $q' \in \mathcal{Q}$  can be computed in two steps:

1. compute the answers to the subquery  $q'_c$ , i.e., with body restricted to the RDFS triples; if  $q'_c$  has no answer, neither has  $q'$ . Otherwise, each answer to  $q'_c$  defines a (partial) instantiation  $\sigma$  of the variables in  $q'$ .
2. compute the assertion-level answers to each partially instantiated query  $(q'_{a,v})_\sigma$ , where  $q'_{a,v}$  is the subquery with body  $b'_a \cup b'_v$ , and return the union of all the obtained answers.

To summarize, Step 1 computes answers to RDFS triples, which allows one to produce a set of partially instantiated queries that no longer contain RDFS triples. Hence, these queries can then be answered using RDF assertions only, which is the purpose of Step 2. Our two-step query reformulation follows this decomposition. It furthermore considers a partition of the set  $\mathcal{R}$  of RDFS entailment rules (recall Table 2) into two subsets: the set of rules  $\mathcal{R}_c$  that produces *RDFS constraints* and the set of rules  $\mathcal{R}_a$  that produces *RDF assertions*:

<sup>8</sup><http://docs.openlinksw.com/virtuoso/rdfsparqlruleimpl>

- $\mathcal{R}_c = \{\text{rdfs5}, \text{rdfs11}, \text{ext1}, \text{ext2}, \text{ext3}, \text{ext4}\};$
- $\mathcal{R}_a = \{\text{rdfs2}, \text{rdfs3}, \text{rdfs7}, \text{rdfs9}\}.$

The reason of this decomposition is that query answering remains complete if, on the one hand, only  $\mathcal{R}_c$  is considered to answer queries made of RDFS triples (Step 1: for any graph  $G$ ,  $q'_c(G, \mathcal{R}) = q'_c(G, \mathcal{R}_c)$ ), and, on the other hand, only  $\mathcal{R}_a$  is considered to answer queries on RDF assertions only, as shown in [12].

Query reformulation does not directly work on the entailment rules as classical backward-chaining techniques would do. Instead, a set of so-called *reformulation rules* is specifically associated with  $\mathcal{R}_c$  (resp.  $\mathcal{R}_a$ ). We can now outline the two-step query reformulation algorithm:

**Step 1. Reformulation w.r.t.  $\mathcal{R}_c$ :** The input BGPQ  $q$  is first reformulated into a union  $\mathcal{Q}_c$  of partially instantiated BGPQs, using the set of reformulation rules associated with  $\mathcal{R}_c$  (see Figure 1). This reformulation step is sound and complete for query answering w.r.t.  $\mathcal{R}_c$ , i.e., for any graph  $G$ ,  $q(G, \mathcal{R}_c) = \overline{\mathcal{Q}_c(G)}$ ; furthermore, it preserves the answers with respect to the set  $\mathcal{R}$ , i.e.,  $q(G, \mathcal{R}_c \cup \mathcal{R}_a) = \overline{\mathcal{Q}_c(G, \mathcal{R}_c \cup \mathcal{R}_a)}$  (see Theorem 1 in Section 4.3).

**Step 2. Reformulation w.r.t.  $\mathcal{R}_a$ :** We recall that  $\mathcal{Q}_c$  consists of queries that do not contain RDFS triples. It is given as input to the query reformulation algorithm of [12], which relies on a set of reformulation rules associated with  $\mathcal{R}_a$  to output a union  $\mathcal{Q}_{c,a}$  of partially instantiated BGPQs. This reformulation step being sound and complete for query answering on the RDF assertions of an RDF graph, we obtain the soundness and completeness of the two-step reformulation, i.e.,  $q(G, \mathcal{R}_c \cup \mathcal{R}_a) = \overline{\mathcal{Q}_c(G, \mathcal{R}_a)} = \overline{\mathcal{Q}_{c,a}(G)}$  (see Theorem 2 in Section 4.3).

## 4.2 Reformulation rules associated with $\mathcal{R}_c$

We now detail reformulation rules associated with  $\mathcal{R}_c$ , see Figure 1. Each reformulation rule is of the form  $\frac{\text{input}}{\text{output}}$ , where the input is composed of a triple from a partially instantiated query  $q_\sigma$  and a triple from  $O$  and the output is a new query obtained from  $q_\sigma$  by instantiating a variable, removing the input triple, or replacing it by one or two triples. The notation *old triple/new triple(s)* means that *old triple* is replaced by *new triple(s)*. The specific case where *old triple* is simply removed is denoted by *old triple/–*. The notations for the triples themselves are the following:

- a bold character like **c**, **p**, **s** or **o** represents an IRI or a blank node
- a *v* character represents a variable of the query
- *s* and *o* characters represent either variables, IRIs or blank nodes, in subject and object positions respectively.

The four rules (1) substitute a variable in a property position by one of the four built-in RDFS properties. All the other rules take as input query triples of the form  $(s, \mathbf{p}, o)$ , where  $\mathbf{p}$  is a built-in RDFS property. Rule (2) simply removes from  $q_\sigma$  an (instantiated) input triple found in  $O$ .

Query triples with a domain ( $\leftrightarrow_d$ ) or range property ( $\leftrightarrow_r$ ) are processed by Rules (3)-(11). Given a triple  $(\mathbf{p}, \leftrightarrow, \mathbf{c})$  in  $O$  (where  $\leftrightarrow$  stands for  $\leftrightarrow_d$  or  $\leftrightarrow_r$ ),

$$\frac{(s, v, o) \in q_\sigma}{q_\sigma \cup \{v \rightarrow \prec_{sc}\}}, \frac{(s, v, o) \in q_\sigma}{q_\sigma \cup \{v \rightarrow \prec_{sp}\}}, \frac{(s, v, o) \in q_\sigma}{q_\sigma \cup \{v \rightarrow \leftrightarrow_d\}}, \frac{(s, v, o) \in q_\sigma}{q_\sigma \cup \{v \rightarrow \leftrightarrow_r\}} \quad (1)$$

$$\frac{(\mathbf{s}, \mathbf{p}, \mathbf{o}) \in q_\sigma, (\mathbf{s}, \mathbf{p}, \mathbf{o}) \in O}{q_\sigma[(\mathbf{s}, \mathbf{p}, \mathbf{o})/-]} \quad (2)$$

$$\frac{(v_1, \leftrightarrow, v_2) \in q_\sigma, (\mathbf{p}, \leftrightarrow, \mathbf{c}) \in O}{q_\sigma[(v_1, \leftrightarrow, v_2)/(v_1, \prec_{sp}, \mathbf{p}), (\mathbf{c}, \prec_{sc}, v_2)]} \quad (3)$$

$$\frac{(v_1, \leftrightarrow, v_2) \in q_\sigma, (\mathbf{p}, \leftrightarrow, \mathbf{c}) \in O}{q_\sigma \cup \{v_1 \rightarrow \mathbf{p}\}} \quad (4)$$

$$\frac{(v_1, \leftrightarrow, v_2) \in q_\sigma, (\mathbf{p}, \leftrightarrow, \mathbf{c}) \in O}{q_\sigma \cup \{v_2 \rightarrow \mathbf{c}\}} \quad (5)$$

$$\frac{(v, \leftrightarrow, \mathbf{c}) \in q_\sigma, (\mathbf{p}, \leftrightarrow, \mathbf{c}) \in O}{q_\sigma \cup \{v \rightarrow \mathbf{p}\}} \quad (6)$$

$$\frac{(\mathbf{p}, \leftrightarrow, v) \in q_\sigma, (\mathbf{p}, \leftrightarrow, \mathbf{c}) \in O}{q_\sigma \cup \{v \rightarrow \mathbf{c}\}} \quad (7)$$

$$\frac{(v, \leftrightarrow, \mathbf{c}) \in q_\sigma, (\mathbf{p}, \leftrightarrow, \mathbf{c}) \in O}{q_\sigma[(v, \leftrightarrow, \mathbf{c})/(v, \prec_{sp}, \mathbf{p})]} \quad (8)$$

$$\frac{(\mathbf{p}, \leftrightarrow, v) \in q_\sigma, (\mathbf{p}, \leftrightarrow, \mathbf{c}) \in O}{q_\sigma[(\mathbf{p}, \leftrightarrow, v)/(\mathbf{c}, \prec_{sc}, v)]} \quad (9)$$

$$\frac{(s, \leftrightarrow, \mathbf{c}_1) \in q_\sigma, (\mathbf{c}, \prec_{sc}, \mathbf{c}_1) \in O, \mathbf{c} \neq \mathbf{c}_1}{q_\sigma[(s, \leftrightarrow, \mathbf{c}_1)/(s, \leftrightarrow, \mathbf{c})]} \quad (10)$$

$$\frac{(\mathbf{p}, \leftrightarrow, o) \in q_\sigma, (\mathbf{p}, \prec_{sp}, \mathbf{p}_1) \in O, \mathbf{p} \neq \mathbf{p}_1}{q_\sigma[(\mathbf{p}, \leftrightarrow, o)/(\mathbf{p}_1, \leftrightarrow, o)]} \quad (11)$$

$$\frac{(v_1, \prec, v_2) \in q_\sigma, (\mathbf{c}_1, \prec, \mathbf{c}_2) \in O}{q_\sigma \cup \{v_1 \rightarrow \mathbf{c}_1\}} \quad (12)$$

$$\frac{(v, \prec, \mathbf{c}_2) \in q_\sigma, (\mathbf{c}_1, \prec, \mathbf{c}_2) \in O}{q_\sigma \cup \{v \rightarrow \mathbf{c}_1\}} \quad (13)$$

$$\frac{(\mathbf{c}_1, \prec, v) \in q_\sigma, (\mathbf{c}_1, \prec, \mathbf{c}_2) \in O}{q_\sigma \cup \{v \rightarrow \mathbf{c}_2\}} \quad (14)$$

$$\frac{(\mathbf{c}_1, \prec, o) \in q_\sigma, (\mathbf{c}_1, \prec, \mathbf{c}_2) \in O, \mathbf{c}_1 \neq \mathbf{c}_2}{q_\sigma[(\mathbf{c}_1, \prec, o)/(\mathbf{c}_2, \prec, o)]} \quad (15)$$

$$\frac{(s, \prec, \mathbf{c}_2) \in q_\sigma, (\mathbf{c}_1, \prec, \mathbf{c}_2) \in O, \mathbf{c}_1 \neq \mathbf{c}_2}{q_\sigma[(s, \prec, \mathbf{c}_2)/(s, \prec, \mathbf{c}_1)]} \quad (16)$$

**Fig. 1.** Reformulation rules for a partially instantiated query  $q_\sigma$  w.r.t. an RDFS ontology  $O$ . For compactness, we factorize similar rules, using the symbol  $\leftrightarrow$  to denote either  $\leftrightarrow_d$  or  $\leftrightarrow_r$ , and  $\prec$  to denote either  $\prec_{sc}$  or  $\prec_{sp}$ .

Rule (3) replaces a query triple of the form  $(v_1, \leftrightarrow, v_2)$  by two triples  $(v_1, \prec_{sp}, p)$  and  $(c, \prec_{sc}, v_2)$ . This rule relies on the fact that a triple  $(p', \leftrightarrow, c')$  belongs to the saturation of the RDF graph by  $\mathcal{R}_c$  if and only if  $p'$  is a subproperty of  $p$  (including  $p = p'$ ) and  $c$  is a subclass of  $c'$  (including  $c = c'$ ), see Lemma 1 in Section 4.3. However, we do not assume that the ontology ensures the reflexivity of the subclass and subproperty relations, hence Rules (4)-(7), whose sole purpose is to deal with the cases  $c = c'$  and  $p = p'$ . Should the ontology contain axiomatic triples ensuring the reflexivity of subclass and subproperty, these four rules would be useless. Note that a natural candidate rule to deal with the case where  $c \neq c'$  and  $p \neq p'$  would have been the following:

$$\frac{(p', \leftrightarrow, c') \in q_\sigma, (p, \leftrightarrow, c) \in O}{q_{\sigma[(p', \leftrightarrow, c')/(p', \prec_{sp}, p), (c, \prec_{sc}, c')]} \quad (17)$$

However, such a rule is flawed: it would blindly consider all triples  $(p, \leftrightarrow, c)$  from  $O$ , which causes a combinatorial explosion. Instead, we propose Rules (10) and (11), which use  $p'$  and  $c'$  as guides to replace  $(p', \leftrightarrow, c')$  by other domain / range triples based on the subproperty-chains from  $p'$  and the subclass-chains to  $c'$ . Query triples with a subclass ( $\prec_{sc}$ ) or subproperty ( $\prec_{sp}$ ) property are processed by Rules (12)-(16). Rules (12), (13), (14) instantiate a variable using an ontology triple of the form  $(c_1, \prec, c_2)$ . In Rule (12), which considers a query triple with two variables and instantiates one of these variables, we arbitrarily chose to instantiate the first variable. The two last rules allow to go up or down in the class and property hierarchies.

### 4.3 Reformulation algorithm associated with $\mathcal{R}_c$

The reformulation algorithm itself, denoted by  $\text{Reformulate}_c$ , is presented in Algorithm 1. The set of queries to be explored (named *toExplore*) initially contains  $q$ . Exploring a query consists of generating all new queries that can be obtained from it by applying a reformulation rule (lines 7–9). Newly generated queries are put in the set named *produced*. The algorithm proceeds in a breadth-first manner, exploring at each step the queries that have been generated at the previous step. When no new query can be generated at a step, the algorithm stops, otherwise the next step will explore the newly generated queries (line 11). Note the use of a set named *explored*, which contains all explored queries; the purpose of this set is to avoid infinite generation of the same queries when the subclass or subproperty hierarchy contains cycles (other than loops), otherwise it is useless. Importantly, not all explored queries are returned in the resulting set, but only those that no longer contain RDFS triples (lines 5–6). Indeed, on the one hand RDFS triples that contain variables are instantiated by the rules in all possible ways using the ontology, and, on the other hand, instantiated triples that belong to the ontology are removed (by Rule (2)). Finally, note that a variable  $v$  in a triple of the form  $(s, v, o)$  is replaced by a built-in RDFS property in some queries (by Rule (1)) and left unchanged in others as it may also be later mapped to a user-defined property in the RDF graph  $G$ .

A simple analysis of the reformulation rule behavior shows that the worst-case time complexity of algorithm  $\text{Reformulate}_c$  is polynomial in the size of  $O$  and

**Algorithm 1:** Reformulate<sub>c</sub>


---

**Input** : BGPQ  $q$  and ontology  $O$   
**Output**: the reformulation of  $q$  with the rules from Fig. 1

```

1 result ← ∅; toExplore ← {q}; explored ← ∅
2 while toExplore ≠ ∅ do
3   produced ← ∅
4   for each qσ ∈ toExplore do
5     if qσ does not contain any RDFS triple then
6       result ← result ∪ {qσ}
7     for each RDFS triple t in qσ do
8       for each q'σ obtained by applying a reformulation rule to t do
9         produced ← produced ∪ {q'σ}
10    explored ← explored ∪ {qσ}
11  toExplore ← produced \ explored
12 return result
```

---

simply exponential in the size of  $q$ . More precisely:

**Proposition 1.** *The algorithm Reformulate<sub>c</sub> runs in time  $\mathcal{O}(|Val(O)|^{6|q|})$ , where  $|q|$  is the number of triples in the body of  $q$ .*

The correctness of the algorithm relies on the following lemma, which characterizes the saturated graph  $G^{\mathcal{R}_c}$  from the triples of  $G$ . We call  $\prec_{sc}$ -chain (resp.  $\prec_{sp}$ -chain) from  $s$  to  $o$  a possibly empty sequence of triples  $(s_i, \prec_{sc}, o_i)$  (resp.  $(s_i, \prec_{sp}, o_i)$ ) with  $1 \leq i \leq n$ , such that  $s_1 = s$ ,  $o_n = o$  and, for  $i > 1$ ,  $s_i = o_{i-1}$ . Since we do not enforce the reflexivity of the subclass relation, a triple  $(c, \prec_{sc}, c)$  belongs to  $G^{\mathcal{R}}$  if and only if there is a non-empty  $\prec_{sc}$ -chain from  $c$  to  $c$  (which includes the case  $(c, \prec_{sc}, c) \in G$ ). The same holds for the subproperty relation.

**Lemma 1.** *Let  $G$  be an RDF graph. It holds that:*

- $(c, \prec_{sc}, c') \in G^{\mathcal{R}_c}$  iff  $G$  contains a non-empty  $\prec_{sc}$ -chain from  $c$  to  $c'$ ;
- $(p, \prec_{sp}, p') \in G^{\mathcal{R}_c}$  iff  $G$  contains a non-empty empty  $\prec_{sp}$ -chain from  $p$  to  $p'$ ;
- $(p', \leftrightarrow_d, c') \in G^{\mathcal{R}_c}$  iff  $G$  contains a triple  $(p, \leftrightarrow_d, c)$ , a (possibly empty)  $\prec_{sp}$ -chain from  $p'$  to  $p$  and a (possibly empty)  $\prec_{sc}$ -chain from  $c$  to  $c'$ . The case for  $(p', \leftrightarrow_r, c') \in G^{\mathcal{R}_c}$  is similar (replace  $\leftrightarrow_d$  by  $\leftrightarrow_r$  in the statement above).

Below, we assume *without loss of generality* that the input query does not contain blank nodes; if needed, these have been equivalently replaced by variables. Therefore, all blank nodes that occur in the output reformulation have been introduced by the reformulation rules, and specifically refer to unknown classes and properties they identify within the ontology at hand. This justifies the subsequent use of non-standard query evaluation and answering in the next theorems.

**Theorem 1.** *Let  $G$  be an RDF graph with ontology  $O$  and  $q$  be a BGP query without blank nodes. Let  $\mathcal{Q}_c$  be the output of Reformulate<sub>c</sub>( $q, O$ ). Then:*

$$q(G, \mathcal{R}_c) = \overbrace{q(G, \mathcal{R}_c)} = \overbrace{\mathcal{Q}_c(G)} \quad (18)$$

$$q(G, \mathcal{R}_c \cup \mathcal{R}_a) = \overbrace{q(G, \mathcal{R}_c \cup \mathcal{R}_a)} = \overbrace{\mathcal{Q}_c(G, \mathcal{R}_c \cup \mathcal{R}_a)} \quad (19)$$

*Example 5.* Consider the BGPQ asking for *how someone is related to some particular kind of company*:  $q(x, y) \leftarrow (x, y, z), (z, \tau, t), (y, \prec_{sp}, :worksFor), (t, \prec_{sc}, :Comp)$ . Its answer set on  $G_{ex}$  w.r.t.  $\mathcal{R}$ , which can be easily checked using  $(G_{ex})^{\mathcal{R}}$  provided in Section 2, is:  $q(G_{ex}, \mathcal{R}) = \{ \langle :p1, :ceoOf \rangle \}$ . The output of  $\text{Reformulate}_c(q, \text{RDFS}(G_{ex}))$  is:

$$\begin{aligned} \mathcal{Q}_c = \{ & q'(x, :ceoOf) \leftarrow (x, :ceoOf, z), (z, \tau, \_ :b_C), \\ & q''(x, :hiredBy) \leftarrow (x, :hiredBy, z), (z, \tau, \_ :b_C) \} \end{aligned}$$

where  $q'$  and  $q''$  are obtained by binding, using Rule (13),  $y$  to either  $:ceoOf$  or  $:hiredBy$ , and  $t$  to  $\_ :b_C$ . Further, these bindings have also produced the fully instantiated RDFS constraints  $(:ceoOf, \prec_{sp}, :worksFor)$  and  $(\_ :b_C, \prec_{sc}, :Comp)$  in  $q'$ , as well as  $(:hiredBy, \prec_{sp}, :worksFor)$  and  $(\_ :b_C, \prec_{sc}, :Comp)$  in  $q''$ , which have then been eliminated by Rule (2).

The *non-standard* answering of  $\mathcal{Q}_c$  on  $G_{ex}$  w.r.t.  $\mathcal{R}$ , i.e.,  $\overbrace{q'(G_{ex}, \mathcal{R})} \cup \overbrace{q''(G_{ex}, \mathcal{R})}$  provides the correct answer set  $\{ \langle :p1, :ceoOf \rangle \}$ , whose only tuple results from  $q'$ . Note that, using *standard* answering, the incorrect answer  $\langle :p2, :hiredBy \rangle$  would have also been obtained from  $q''$ , since under this semantics  $q''$  asks for *who is hired by an organization of some type* (this is the case of  $:p2$  who is hired by a public administration) and not *who is hired by an organization of the particular unknown type of company designated by  $\_ :b_C$  in  $G_{ex}$* .

We now rely on the query reformulation algorithm, from [12], say  $\text{Reformulate}_a$ , which takes as input a partially instantiated BGPQ  $q$  *without RDFS triples*, and a graph  $G$  and, using a set of reformulation rules associated with  $\mathcal{R}_a$ , outputs a reformulation  $\mathcal{Q}_a$  such that:  $q(G, \mathcal{R}_c \cup \mathcal{R}_a) = q(G, \mathcal{R}_a) = \overbrace{\mathcal{Q}_a(G)}$ . The adaptation of  $\text{Reformulate}_a$  to an input UBGPQ instead of a BGPQ is straightforward. Furthermore, we notice that the algorithm would consider potential blank nodes in the input query as if they were IRIs. Hence, denoting by  $\mathcal{Q}_{c,a}$  the output of  $\text{Reformulate}_a(\mathcal{Q}_c, O)$ , we obtain:

$$\overbrace{\mathcal{Q}_c(G, \mathcal{R}_c \cup \mathcal{R}_a)} = \overbrace{\mathcal{Q}_c(G, \mathcal{R}_a)} = \overbrace{\mathcal{Q}_{c,a}(G)} \quad (20)$$

Putting together (20) and statement (19) in Theorem 1, we can prove the correctness of the global reformulation algorithm:

**Theorem 2.** *Let  $G$  be an RDF graph and  $q$  be a BGPQ without blank nodes. Let  $\mathcal{Q}_{c,a}$  be the reformulation of  $q$  by the 2-step algorithm described in Section 4.1. Then:*

$$q(G, \mathcal{R}_c \cup \mathcal{R}_a) = \overbrace{\mathcal{Q}_{c,a}(G)}$$

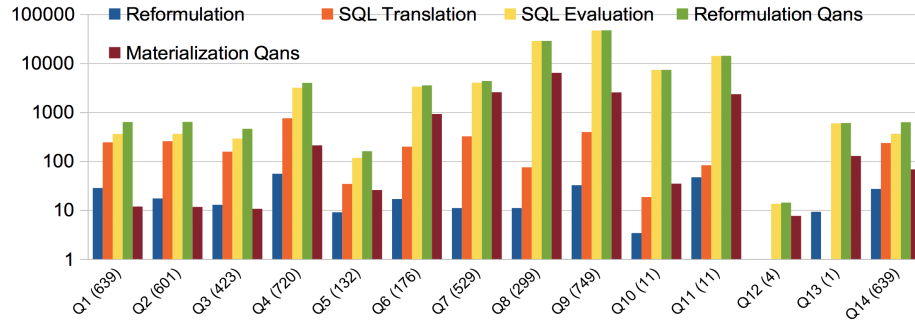


Fig. 2. Query answering times through reformulation and saturation.

## 5 Experimental evaluation

We have implemented our reformulation algorithm on top of OntoSQL (<https://ontosql.inria.fr>), a Java platform providing efficient RDF storage, saturation, and query evaluation on top of an RDBMS [9, 12]; we used Postgres v9.6. To save space, OntoSQL encodes IRIs and literals into integers, and a dictionary table which allows going from one to the other. It stores all resources of a certain type in a one-attribute table, all subject, object pairs for each data property in a table, and all schema triples in another table; the tables are indexed. Our server has a 2,7 GHz Intel Core i7 and 160 GB of RAM; it runs CentOS Linux 7.5.

We generated LUBM<sup>3</sup> data graphs [15] of 10M triples and restricted the ontology to RDFS, leading to 175 triples ( $123 \prec_{sc}$ ,  $5 \prec_{sp}$ ,  $25 \leftrightarrow_d$  and  $22 \hookrightarrow_r$ ). We devised 14 queries having from 3 to 7 triples; one has no result, while the others have a few dozen to three hundred thousand results. Each has 1 or 2 triples which match the ontology (and must be evaluated on it for correctness), including (but not limited to) the generic triple  $(x, y, z)$ , which appears 7 times overall in our workload. Some of our queries are not handled through reformulation by AllegroGraph and Stardog, nor by Virtuoso (recall Section 3).

Figure 2 shows for each query: the size of the UBG PQ reformulation (in parenthesis after the query name on the  $x$  axis), i.e., the number of BGPQs it contains; the reformulation time (with both  $\mathcal{R}_c$  and  $\mathcal{R}_a$ ); the time to translate the reformulation into SQL; the time to evaluate this SQL query; the total query answering time through reformulation, and (for comparison) through saturation. Note the logarithmic  $y$  axis. *Details of our experiments are available online*<sup>5</sup>. The reformulation time is very short (0.2 ms to 55 ms). Unsurprisingly, the time to convert the reformulation into SQL is closely correlated with the reformulation size. The overhead of our approach is quite negligible, given that the answering time through reformulation is very close to the SQL evaluation time.

As expected, saturation-based query answering is faster; however, saturating this graph took more than 1289 seconds, while the slowest query (Q9) took 46 seconds. As in [12], we compute for each query  $Q$  a *threshold*  $n_Q$  which is the smallest number of times we need to run  $Q$ , so that saturating  $G$  and running  $Q$   $n_Q$  times on  $G^{\mathcal{R}}$  is faster than  $n_Q$  runs of  $Q$  through reformulation; intuitively, *after  $n_Q$  runs of  $Q$ , the saturation cost amortizes*. For our queries,  $n_Q$  ranged from 29 (Q9) to 9648 (Q5), which shows that saturation costs take a while to

amortize. If the graph or the ontology change, requiring maintenance of the saturated graph, reformulation may be even more competitive.

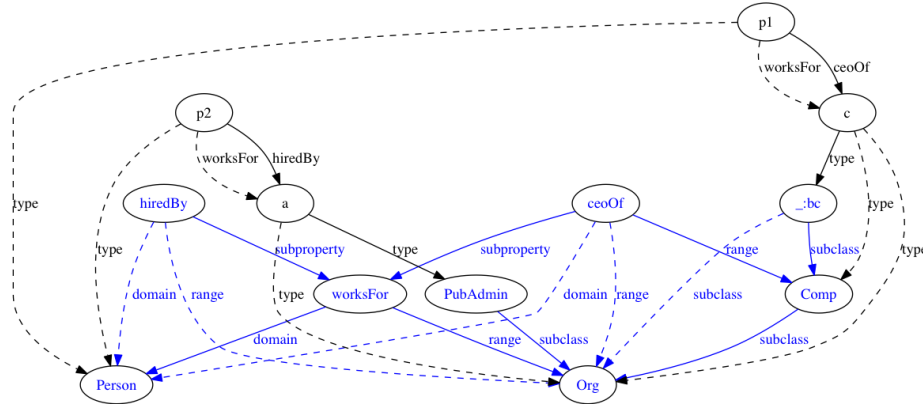
## 6 Conclusion

We have presented a novel reformulation-based query answering technique for RDF graphs with RDFS ontologies. Its novelty lies in its capacity to handle query triples over both the assertions and the ontology; such queries are not always handled correctly by existing RDF engines. In the future, we plan to integrate our reformulation technique in the cost-based optimized reformulation framework we introduced in [9] to improve its performance, and to an OBDA setting along the lines of [14].

**Acknowledgements:** This work is supported by the Inria Project Lab grant iCoda, a collaborative project between Inria and several major French media.

## 7 Appendix

This appendix provides the figure of the running example graph and proofs of the results claimed in the paper.



**Fig. 3.** Illustration of Example 1, where black edges represent assertion triples and blue edges represent RDFS triples. Plain edges are those contain in  $G_{\text{ex}}$  and dotted edges those added its saturation  $G_{\text{ex}}^{\mathcal{R}}$ .

### Proof of Proposition 1

*Proof.* We provide an upper bound the number of reformulations explored during the reformulation of a BGPQ by analyzing the producer-consumer dependencies among rules w.r.t. the form of the query. Given a query form  $Q$  and an ontology  $O$ , we denote by  $\#\text{explored}(Q, O)$  the number of explored reformulations during the execution of  $\text{Reformulate}_c(q, O)$  for a BGPQ  $q$  of form  $Q$ .



First, we notice that for the most general query form  $Q(\bar{x}) \leftarrow t_1, t_2, \dots, t_n$  (where the  $t_i$  are triples), it holds that:

$$\#\text{explored}(Q, O) \leq \prod_{1 \leq i \leq n} \#\text{explored}(Q_i, O)$$

where  $Q_i$  has the form  $Q_i(x_i) \leftarrow t_i$ , with  $x_i$  the list of variables in  $t_i$ . We now analyze the reformulations obtained for the different forms of queries composed of a single triple.

Let us consider the query form  $Q_0(\bar{x}) \leftarrow (s, v, o)$ , where  $s$  and  $o$  are values or variables and  $v$  is a variable. The rules in (1) are the only rules that can consume  $Q_0$ . They produce queries of the form  $Q_1(\bar{x}_1) \leftarrow (s, \leftarrow_d, o)$  or  $Q_2(\bar{x}_2) \leftarrow (s, \prec_{sc}, o)$  (and there are two similar cases with properties  $\hookrightarrow_r$  and  $\prec_{sp}$ ). Since no rule feeds rule (1), it holds that:

$$\#\text{explored}(Q_0, O) \leq 2(\#\text{explored}(Q_1, O) + \#\text{explored}(Q_2, O))$$

Queries of the form  $Q_2$  only feed rules from (12) to (16). These rules always produce queries of the form  $Q_2(\bar{x}_2) \leftarrow (s, \prec_{sc}, o)$ , where  $s$  and  $o$  are either values from  $\text{Val}(O)$  or variables. Moreover, there are at most 2 variables in  $\text{Var}(Q_2)$ , which can only be instantiated by values from  $\text{Val}(O)$  or variables.

In the end,

$$\#\text{explored}(Q_2, O) \leq (\#\text{Val}(O) + 1)^2 (\#\text{Val}(O) + 1)^2 = (\#\text{Val}(O) + 1)^4.$$

Concerning a query of the form  $Q_1$ , either rule (3) can be applied, then produced queries have the form  $Q_3(\bar{x}_3) \leftarrow (v_1, \prec_{sp}, \mathbf{p}), (\mathbf{c}, \prec_{sc}, v_2)$ , or a rule from (4) to (11) can be applied. In the later case, we observe that all further produced queries will have the form  $Q_4(\bar{x}_4) \leftarrow (s, \mathbf{p}, o)$  with  $\mathbf{p} \in \{\leftarrow_d, \prec_{sc}, \prec_{sp}\}$ ,  $s$  and  $o$  belonging to  $\text{Val}(O) \cup \text{Var}(Q_4)$ , and there is at most one variable among  $s$  and  $o$ .

So, we obtain by counting the number of possible values by position that

$$\#\text{explored}(Q_3, O) \leq ((\#\text{Val}(O) + 1)^2 \#\text{Val}(O))^2$$

and

$$\#\text{explored}(Q_4, O) \leq 3(\#\text{Val}(O) + 1)^4.$$

In the end,

$$\begin{aligned} \#\text{explored}(Q_1, O) &\leq (\#\text{Val}(O) + 1)^4 \#\text{Val}(O)^2 + 3(\#\text{Val}(O) + 1)^4 \\ &\leq (\#\text{Val}(O) + 1)^6 + 3(\#\text{Val}(O) + 1)^4 \end{aligned}$$

It follows that the maximal number of explored reformulations for an input query with a single triple is  $\mathcal{O}(\#\text{Val}(O)^6)$ , hence in  $\mathcal{O}(\#\text{Val}(O)^{6|q|})$  for a BGPQ  $q$ . Note that  $|q|$  is a rough upper bound, since one should only consider the number of triples in  $q$  with an RDFS property or a variable in property position.

**Proof of Lemma 1** The only entailment rule in  $\mathcal{R}_c$  that allows one to infer a new triple with property  $\prec_{sc}$  (respectively  $\prec_{sp}$ ) is the rule **rdfs11** (resp. **rdfs5**). Since this rule states the transitivity of the property  $\prec_{sc}$  (resp.  $\prec_{sp}$ ), it holds that  $(c, \prec_{sc}, c') \in G^{\mathcal{R}_c}$  iff  $G$  contains a non-empty  $\prec_{sc}$ -chain from  $c$  to  $c'$  (resp.  $(p, \prec_{sp}, p') \in G^{\mathcal{R}_c}$  iff  $G$  contains a non-empty empty  $\prec_{sp}$ -chain from  $p$  to  $p'$ ). Assume now that  $(p', \leftrightarrow_d, c') \in G^{\mathcal{R}_c}$ . The only entailment rules in  $\mathcal{R}_c$  that entail a triple with property  $\leftrightarrow_d$  are **ext1** and **ext3**. The body of these rules contain a triple with property  $\leftrightarrow_d$ , so there exists an entailment chain (of triples with  $\leftrightarrow_d$  property) of length  $l \geq 0$  starting from  $G$  and using only rules **ext1** and **ext3**. We prove by induction on  $l$  that  $G$  contains a triple  $(p, \leftrightarrow_d, c)$ , a (possibly empty)  $\prec_{sp}$ -chain from  $p'$  to  $p$  and a (possibly empty)  $\prec_{sc}$ -chain from  $c$  to  $c'$ .

- If  $l = 0$ , then  $(p', \leftrightarrow_d, c') \in G$  and there are an empty  $\prec_{sp}$ -chain from  $p'$  to  $p'$  and an empty  $\prec_{sc}$ -chain from  $c'$  to  $c'$ .
- Otherwise ( $l > 0$ ), the last rule applied in the chain is:
  - either **ext1**, so  $G^{\mathcal{R}_c}$  contains a triple  $(p', \leftrightarrow_d, c_1)$ , which results from an entailment chain of length  $l-1$  starting from  $G$  and using only rules **ext1** and **ext3**, and a triple  $(c_1, \prec_{sc}, c')$ . By induction hypothesis, we know that  $G$  contains a triple  $(p, \leftrightarrow_d, c)$ , a (possibly empty)  $\prec_{sp}$ -chain from  $p'$  to  $p$  and a (possibly empty)  $\prec_{sc}$ -chain from  $c$  to  $c_1$ . Moreover, by using the first point of the lemma (proved above),  $(c_1, \prec_{sc}, c') \in G^{\mathcal{R}_c}$  implies that  $G$  contains a non-empty  $\prec_{sc}$ -chain from  $c_1$  to  $c'$ . So, concatenating the two  $\prec_{sc}$ -chains, we obtain a  $\prec_{sc}$ -chain from  $c$  to  $c'$ . Hence,  $G$  contains a triple  $(p, \leftrightarrow_d, c)$ , a (possibly empty)  $\prec_{sp}$ -chain from  $p'$  to  $p$  and a  $\prec_{sc}$ -chain from  $c$  to  $c'$ .
  - or **ext3**, and the proof is similar to that for **ext1**, replacing  $\prec_{sc}$ -chains by  $\prec_{sp}$ -chains.

We have proven that  $(p', \leftrightarrow_d, c') \in G^{\mathcal{R}_c}$  implies that  $G$  contains a triple  $(p, \leftrightarrow_d, c)$ , a (possibly empty)  $\prec_{sp}$ -chain from  $p'$  to  $p$  and a (possibly empty)  $\prec_{sc}$ -chain from  $c$  to  $c'$ . The converse implication is straightforward: from the two first points of the lemma, we obtain  $(c, \prec_{sc}, c') \in G^{\mathcal{R}_c}$  and  $(p', \prec_{sp}, p) \in G^{\mathcal{R}_c}$ , then by one application of each entailment rule **ext1** and **ext3**, we obtain  $(p', \leftrightarrow_d, c') \in G^{\mathcal{R}_c}$ .

**Proof of Theorem 1** For the sake of readability, we assume in the following that  $G$  does not contain blank nodes. So, we do not need non-standard query evaluation. This assumption can be done without loss of generality. Indeed, we may define a one-to-one mapping  $f$  from the blank nodes of  $G$  to fresh IRIs, apply  $f$  to  $G$  before any processing, and apply the inverse mapping  $f^{-1}$  to the answer tuples obtained considering  $f(G)$  to get answers considering  $G$ .

With the above assumption, to prove statement 18, it remains to prove that  $q(G, \mathcal{R}_c) = \mathcal{Q}_c(G)$  holds. We first prove  $\mathcal{Q}_c(G) \subseteq q(G, \mathcal{R}_c)$  (soundness) then  $q(G, \mathcal{R}_c) \subseteq \mathcal{Q}_c(G)$  (completeness).

**(soundness)** We want to prove that for all  $q'_\sigma$ , reformulation of  $q$  in  $\mathcal{Q}_c$ , for all tuple  $t$  answer to  $q'_\sigma$  in  $G$ , there is  $G'$  obtained from  $G$  by application of some entailment rules to  $G$  such that  $t$  is an answer to  $q$  in  $G'$ . In other words, we want

to prove that  $q'_{\sigma'}(G, \emptyset) \subseteq q(G, \mathcal{R}_c)$ . Since  $q'_{\sigma'}(G, \emptyset) \subseteq q'_{\sigma'}(G, \mathcal{R}_c)$ , it is sufficient to prove that  $q'_{\sigma'}(G, \mathcal{R}_c) \subseteq q(G, \mathcal{R}_c)$ .

The proof is done by induction on the length  $l$  of a sequence of reformulation rules leading to  $q'_{\sigma'}$ , starting from  $O$  and  $q$ .

*Base step* For  $l = 0$ , we have  $q'_{\sigma'} = q$ , so  $q'_{\sigma'}(G, \mathcal{R}_c) \subseteq q(G, \mathcal{R}_c)$ .

*Inductive step* For  $l < \alpha$ , suppose that  $q'_{\sigma'}(G, \mathcal{R}_c) \subseteq q(G, \mathcal{R}_c)$  holds. Now at  $l = \alpha$ ,  $q'_{\sigma'}$  has been produced from  $q''_{\sigma''}$  by the application of a reformulation rule (i) and  $q''_{\sigma''}$  is a reformulation of  $q$ . So that sequence being of length  $< \alpha$ , we get  $q''_{\sigma''}(G) \subseteq q(G, \mathcal{R}_c)$  by induction hypothesis. We will show that  $q'_{\sigma'}(G, \mathcal{R}_c) \subseteq q''_{\sigma''}(G, \mathcal{R}_c)$ .

There are basically two cases:

- the reformulation rule (i) instantiates a variable of  $q''_{\sigma''}$  to generate  $q'_{\sigma'}$  i.e., rule (i) is one of the following (1), (4)-(7), (12)-(14). In this case,  $q'_{\sigma'}$  is contained in  $q''_{\sigma''}$ , so  $q'_{\sigma'}(G, \mathcal{R}_c) \subseteq q''_{\sigma''}(G, \mathcal{R}_c)$ .
- the reformulation rule (i) has the form  $\frac{t_1 \in q_{\sigma}, t_2 \in O}{q_{\sigma}[t_1/t_3]}$  that replaces a triple in  $q''_{\sigma''}$  by another one (or two for the rule (3)). Observe here that  $\sigma' = \sigma''$  holds. If  $\varphi(\bar{x}_{\sigma'}) \in q'_{\sigma'}(G, \mathcal{R}_c)$ , then  $\varphi(t_{3\sigma'}) \in G^{\mathcal{R}_c}$ . Furthermore, the reformulation rules ensure that  $\varphi(t_{3\sigma'}), t_2 \models_{\mathcal{R}_c} \varphi(t_{1\sigma''})$ . As a result,  $\varphi(t_{1\sigma''}) \in G^{\mathcal{R}_c}$ , and  $\varphi$  is a total assignment of the variables of  $q''_{\sigma''}$  such that  $\varphi(\bar{x}_{\sigma'}) = \varphi(\bar{x}_{\sigma''}) \in q''_{\sigma''}(G, \mathcal{R}_c)$ .

In each case, we get  $q'_{\sigma'}(G, \mathcal{R}_c) \subseteq q''_{\sigma''}(G, \mathcal{R}_c)$  which concludes the proof of  $q'_{\sigma'}(G, \emptyset) \subseteq q(G, \mathcal{R}_c)$ , so  $\mathcal{Q}_c(G) \subseteq q(G, \mathcal{R}_c)$ .

**(completeness)** We now show that  $q(G, \mathcal{R}_c) \subseteq \mathcal{Q}_c(G)$  with

$\mathcal{Q}_c(G) = \bigcup_{q'_{\sigma'} \in \text{Reformulate}_c(q, O)} q'_{\sigma'}(G, \emptyset)$ , i.e., for each answer tuple  $a \in q(G, \mathcal{R}_c)$ , there exists  $q'_{\sigma'} \in \mathcal{Q}_c$  a reformulation of  $q$  using  $O$  such that  $a \in q'_{\sigma'}(G, \emptyset)$ .

In the following, we will consider that  $\mathcal{Q}_c$  contains queries in which all the instantiated RDFS triples that belong to the ontology are kept; in other words, the triples removed by applications of rule (2) are restored in the resulting queries. This has no impact on the completeness of the algorithm, since the reformulations output in both versions have the same answers in  $G$ .

Let the query  $q$  be defined by  $q(\bar{x}) \leftarrow t_1, t_2, \dots, t_n$  with  $t_i$  being the body triples of  $q$ . An answer from  $q(G, \mathcal{R}_c)$  has the form  $\varphi(\bar{x})$ , where  $\varphi$  is a homomorphism from  $\text{body}(q)$  to  $G^{\mathcal{R}_c}$ . If for all triples  $t_i$  from the body of  $q$ ,  $\varphi(t_i)$  is not an RDFS triple, then  $\varphi(\text{body}(q)) \in G$  (because data triples are not entailed by  $\mathcal{R}_c$ ), so a valid reformulation of  $q$  is  $q$  itself, since  $q(G, \mathcal{R}_c) = q(G, \emptyset)$ . Otherwise, there exists a triple  $t_i$  from the body of  $q$  such that  $\varphi(t_i) \in O^{\mathcal{R}_c}$  and we will show that there exists  $q'_{\sigma'}$  a reformulation of  $q$  where only  $t_i$  has been replaced by a BGP  $P$  such that  $P \subseteq O$  and  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .

First case,  $\varphi(t_i) = (\mathbf{c}, \prec_{sc}, \mathbf{c}') \in G^{\mathcal{R}_c}$ ; according to Lemma 1, there is  $C = ((\mathbf{c}_i, \prec_{sc}, \mathbf{c}_{i+1}))_{1 \leq i < c}$  a  $\prec_{sc}$ -chain in  $G$  such that  $\mathbf{c} = \mathbf{c}_1$  and  $\mathbf{c}' = \mathbf{c}_c$ . The triple  $t_i$  can have one of the following forms:

- $(\mathbf{c}, \prec_{sc}, \mathbf{c}')$ , then we consider  $q'_{\sigma'}$  obtained from  $q$  by applying rule (15) for each triple of  $C$ ; finally  $(\mathbf{c}, \prec_{sc}, \mathbf{c}')$  is replaced by  $(\mathbf{c}_{c-1}, \prec_{sc}, \mathbf{c}') \in O$ . Since  $\sigma' = \emptyset$ ,  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .

- $(c, \prec_{sc}, v')$ , then we consider  $q'_{\sigma'}$  obtained from  $q$  by applying rule (15) for each triple of  $C$  then (14); finally  $(c, \prec_{sc}, v')$  is replaced by  $(c_{c-1}, \prec_{sc}, c') \in O$ . Since  $\sigma' = \{v' \mapsto c'\}$ ,  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .
- $(c, v_p, v)$ , then we consider  $q'_{\sigma'}$  obtained from  $q$  by applying rule (1) then (15) for each triple of  $C$  then (14); finally  $(c, v_p, v)$  is replaced by  $(c_{c-1}, \prec_{sc}, c') \in O$ . Since  $\sigma' = \{v_p \mapsto \prec_{sc}, v \mapsto c'\}$ ,  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .
- $(v, \prec_{sc}, c')$ , then we consider  $q'_{\sigma'}$  obtained from  $q$  by applying rule (16) for each triple of  $C$  in inverse order then (13); finally  $(v, \prec_{sc}, c')$  is replaced by  $(c, \prec_{sc}, c_2) \in O$ . Since  $\sigma' = \{v \mapsto c\}$ ,  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .
- $(v, v_p, c')$ , then we consider  $q'_{\sigma'}$  obtained from  $q$  by applying rule (1) then (16) for each triple of  $C$  in inverse order then (13); finally  $(v, v_p, c')$  is replaced by  $(c, \prec_{sc}, c_2) \in O$ . Since  $\sigma' = \{v_p \mapsto \prec_{sc}, v \mapsto c\}$ ,  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .
- $(v, \prec_{sc}, v')$ , then we consider  $q'_{\sigma'}$  obtained from  $q$  by applying rule (12) then (15) for each triple of  $C$  then (14); finally  $(v, \prec_{sc}, v')$  is replaced by  $(c_{c-1}, \prec_{sc}, c') \in O$ . Since  $\sigma' = \{v' \mapsto c', v \mapsto c\}$ ,  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .
- $(v, v_p, v')$ , then we consider  $q'_{\sigma'}$  obtained from  $q$  by applying rule (1) then (12) then (15) for each triple of  $C$  then (14); finally  $(v, v_p, v')$  is replaced by  $(c_{c-1}, \prec_{sc}, c') \in O$ . Since  $\sigma' = \{v_p \mapsto \prec_{sc}, v' \mapsto c', v \mapsto c\}$ ,  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .

Second case,  $\varphi(t_i) = (\mathbf{p}, \leftrightarrow_d, \mathbf{c}) \in G^{\mathcal{R}c}$ ; according to Lemma 1, there are three cases, depending on whether a chain is empty or not. We describe the case where none of the chains is empty, hence assuming that there exists  $P = ((\mathbf{p}_i, \prec_{sp}, \mathbf{p}_{i+1}))_{1 \leq i \leq p}$  a  $\prec_{sp}$ -chain in  $G$  from  $\mathbf{p}$  to  $\mathbf{p}'$  and  $(\mathbf{p}', \leftrightarrow_d, \mathbf{c}') \in G$  and there exists  $C = ((\mathbf{c}_i, \prec_{sc}, \mathbf{c}_{i+1}))_{1 \leq i \leq c}$  a  $\prec_{sc}$ -chain in  $G$  from  $\mathbf{c}'$  to  $\mathbf{c}$ . The other cases are handled similarly using also rules (4) and (5). The triple  $t_i$  can have the following forms:

- $(\mathbf{p}, \leftrightarrow_d, \mathbf{c})$ , then we consider  $q'_{\sigma'}$  obtained from  $q$  by applying rule (10) for each triple in  $C$  in inverse order then (11) for each triple in  $P$ ; finally  $(\mathbf{p}, \leftrightarrow_d, \mathbf{c})$  is replaced by  $(\mathbf{p}', \leftrightarrow_d, \mathbf{c}') \in O$ . Since  $\sigma' = \emptyset$ ,  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .
- $(\mathbf{p}, \leftrightarrow_d, v')$ , then we consider  $q'_{\sigma'}$  obtained from  $q$  by applying rule (11) for each triple in  $P$  then (9) then (15) for each triple in  $C$  then (14); finally  $(\mathbf{p}, \leftrightarrow_d, \mathbf{c})$  is replaced by  $(c_{c-1}, \prec_{sc}, \mathbf{c}) \in O$ . Since  $\sigma' = \{v' \mapsto \mathbf{c}\}$ ,  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .
- $(\mathbf{p}, v_p, v')$ , then we consider  $q'_{\sigma'}$  obtained from  $q$  by applying rule (1) then (11) for each triple in  $P$  then (9) then (15) for each triple in  $C$  then (14); finally  $(\mathbf{p}, v_p, v')$  is replaced by  $(c_{c-1}, \prec_{sc}, \mathbf{c}) \in O$ . Since  $\sigma' = \{v_p \mapsto \leftrightarrow_d, v' \mapsto \mathbf{c}\}$ ,  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .
- $(v, \leftrightarrow_d, \mathbf{c})$ , then we consider  $q'_{\sigma'}$  obtained from  $q$  by applying rule (10) for each triple in  $C$  in inverse order then (8) then (16) for each triple in  $P$  in inverse order then (13); finally  $(v, \leftrightarrow_d, \mathbf{c})$  is replaced by  $(\mathbf{p}, \prec_{sp}, \mathbf{p}_2) \in O$ . Since  $\sigma' = \{v \mapsto \mathbf{p}\}$ ,  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .
- $(v, v_p, \mathbf{c})$ , then we consider  $q'_{\sigma'}$  obtained from  $q$  by applying rule (1) then (10) for each triple in  $C$  in inverse order then (8) then (16) for each triple in  $P$  in inverse order then (13); finally  $(v, v_p, \mathbf{c})$  is replaced by  $(\mathbf{p}, \prec_{sp}, \mathbf{p}_2) \in O$ . Since  $\sigma' = \{v_p \mapsto \leftrightarrow_d, v \mapsto \mathbf{p}\}$ ,  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .
- $(v, \leftrightarrow_d, v')$ , then we consider  $q'_{\sigma'}$  is obtained from  $q$  by applying rule (3)

- then (16) for each triple in  $P$  inverse order then (13) then on the other triple, (15) for each triple in  $C$  then (14) ; finally  $(v, \leftrightarrow_d, v')$  is replaced by  $(\mathbf{p}, \prec_{sp}, \mathbf{p}_2), (\mathbf{c}_{c-1}, \prec_{sp}, \mathbf{c}) \in O$ . Since  $\sigma' = \{v \mapsto \mathbf{p}, v' \mapsto \mathbf{c}\}$ ,  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .
- $(v, v_p, v')$ , then we consider  $q'_{\sigma'}$ , obtained from  $q$  by applying rule (1) then (3) then (16) for each triple in  $P$  in inverse order then (13) then on the other triple, (15) for each triple in  $C$  then (14) ; finally  $(v, v_p, v')$  is replaced by  $(\mathbf{p}, \prec_{sp}, \mathbf{p}_2), (\mathbf{c}_{c-1}, \prec_{sp}, \mathbf{c}) \in O$ . Since  $\sigma' = \{v \mapsto \mathbf{p}, v' \mapsto \mathbf{c}\}$ ,  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ .

Hence, for each triple  $t_i$  in  $q$  such that  $\varphi(t_i) \in O^{\mathcal{R}_c}$ , there is  $q'_{\sigma'}$ , a reformulation of  $q$ , where only  $t_i$  has been replaced by a BGP  $P$  such that  $P \subseteq O$  and  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma'})$ . It follows that there is  $q''_{\sigma''}$ , a reformulation of  $q$ , in which all body triples of  $q$  mapped by  $\varphi$  to  $O^{\mathcal{R}_c}$  have been replaced by triples that belong to  $O$ , such that  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma''})$ . Since the other triples of  $q$  are necessarily mapped by  $\varphi$  to  $G$  (actually,  $G \setminus O$ ), we conclude that  $\varphi(\bar{x}) = \varphi(\bar{x}_{\sigma''})$  is an answer to  $q''_{\sigma''}$  in  $G$ . This concludes the proof of statement (18), which is the only part of Theorem 1 needed in the proof of Theorem 2. Statement (19) actually follows from the next lemma (Lemma 2).

## Proof of Theorem 2

**Lemma 2.** *For all RDF graph  $G$ , it holds that:*

$$G^{\mathcal{R}_a \cup \mathcal{R}_c} = (G^{\mathcal{R}_a})^{\mathcal{R}_c}$$

*Proof.* For one direction:  $(G^{\mathcal{R}_a})^{\mathcal{R}_c} \subseteq G^{\mathcal{R}_a \cup \mathcal{R}_c}$ . The proof is trivial.

For the converse direction  $G^{\mathcal{R}_a \cup \mathcal{R}_c} \subseteq (G^{\mathcal{R}_a})^{\mathcal{R}_c}$ . We take a triple  $t \in G^{\mathcal{R}_a \cup \mathcal{R}_c}$ , and differentiate two cases:

- either  $t$  is not an RDFS triple, then by applying Theorem 1 of [12],  $t \in G^{\mathcal{R}_a}$ . In other words, assertion rules suffice to entail all RDF assertions.
  - or  $t$  is an RDFS triple. Since the RDFS ontology  $O$  of  $G$  does not contain an RDFS property as subject or object, the entailment rule `rdfs7` does not entail RDFS triples. So,  $t \in O$  or  $t$  has been produced by a rule in  $\mathcal{R}_c$ . Moreover, all rules in  $\mathcal{R}_c$  have a body that contains only RDFS triples, so  $t \in O$  or  $t$  has been entailed from  $O$  using rules in  $\mathcal{R}_c$ , i.e.,  $t \in O^{\mathcal{R}_c}$ . We also know that  $O^{\mathcal{R}_c} \subseteq (G^{\mathcal{R}_a})^{\mathcal{R}_c}$ , so  $t \in (G^{\mathcal{R}_a})^{\mathcal{R}_c}$ .
- In both cases, we have proven that  $t \in (G^{\mathcal{R}_a})^{\mathcal{R}_c}$ .

*Proof (of the theorem).*

$$\begin{aligned}
q(G, \mathcal{R}_a \cup \mathcal{R}_c) &= q(G^{\mathcal{R}_a \cup \mathcal{R}_c}) \\
&= q((G^{\mathcal{R}_a})^{\mathcal{R}_c}) \quad \text{by Lemma 2} \\
&= q(G^{\mathcal{R}_a}, \mathcal{R}_c) \quad \text{by definition of query answering} \\
&= \overbrace{Q_c(G^{\mathcal{R}_a})} \quad \text{by Theorem 1, statement (18)} \\
&= \overbrace{Q_c(G, \mathcal{R}_a)} \quad \text{by definition of query answering} \\
&= \overbrace{Q_{c,a}(G)} \quad \text{by Theorem 6 of [12] since } Q_c \text{ is without RDFS triples}
\end{aligned}$$

## References

1. RDF 1.1 Concepts and Abstract Syntax, <https://www.w3.org/TR/rdf11-concepts/>
2. RDF 1.1 Semantics, <https://www.w3.org/TR/rdf11-nt/#rdfs-entailment>
3. SPARQL 1.1 Query Language, <https://www.w3.org/TR/sparql11-query/>
4. Abiteboul, S., Hull, R., Vianu, V.: Foundations of Databases. Addison-Wesley (1995)
5. Adjiman, P., Goasdoué, F., Rousset, M.C.: SomeRDFS in the semantic web. JODS **8** (2007)
6. Arenas, M., Gutierrez, C., Pérez, J.: Foundations of RDF databases. In: Reasoning Web (2009)
7. Bishop, B., Kiryakov, A., Ognyanoff, D., Peikov, I., Tashev, Z., Velkov, R.: OWLIM: A family of scalable semantic repositories. Semantic Web **2**(1) (2011)
8. Broekstra, J., Kampman, A.: Inferencing and truth maintenance in RDF schema. In: PSSS1 Workshop (2003), <http://ceur-ws.org/Vol-89/broekstra-et-al.pdf>
9. Bursztyn, D., Goasdoué, F., Manolescu, I.: Optimizing reformulation-based query answering in RDF. In: EDBT (2015)
10. Calvanese, D., Giacomo, G.D., Lembo, D., Lenzerini, M., Rosati, R.: Tractable reasoning and efficient query answering in description logics: The DL-Lite family. Journal of Automated Reasoning (JAR) **39**(3) (2007)
11. Goasdoué, F., Karanasos, K., Leblay, J., Manolescu, I.: View selection in semantic web databases. PVLDB **5**(2) (2011), <https://hal.inria.fr/inria-00625090v1>
12. Goasdoué, F., Manolescu, I., Roatis, A.: Efficient query answering against dynamic RDF databases. In: EDBT (2013), <https://hal.inria.fr/hal-00804503v2>
13. Kaoudi, Z., Miliaraki, I., Koubarakis, M.: RDFS reasoning and query answering on DHTs. In: ISWC (2008)
14. Lanti, D., Xiao, G., Calvanese, D.: Cost-driven ontology-based data access. In: ISWC (2017), [https://doi.org/10.1007/978-3-319-68288-4\\_27](https://doi.org/10.1007/978-3-319-68288-4_27)
15. Lutz, C., Seylan, I., Toman, D., Wolter, F.: The combined approach to OBDA: taming role hierarchies using filters. In: ISWC (2013)
16. Neumann, T., Weikum, G.: The RDF-3X engine for scalable management of RDF data. VLDB J. (2010)
17. Urbani, J., van Harmelen, F., Schlobach, S., Bal, H.: QueryPIE: Backward reasoning for OWL Horst over very large knowledge bases. In: ISWC (2011)
18. Urbani, J., Piro, R., van Harmelen, F., Bal, H.E.: Hybrid reasoning on OWL RL. Semantic Web **5**(6) (2014)