

## Towards Learning Commonalities in SPARQL

Sara El Hassad, François Goasdoué, H el ene Jaudoin

► **To cite this version:**

Sara El Hassad, Fran ois Goasdou e, H el ene Jaudoin. Towards Learning Commonalities in SPARQL. Extended Semantic Web Conference (ESWC), May 2017, Portoroz, Slovenia. <hal-01508720v2>

**HAL Id: hal-01508720**

**<https://hal.inria.fr/hal-01508720v2>**

Submitted on 26 Apr 2017

**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 destin ee au d ep ot et  a la diffusion de documents scientifiques de niveau recherche, publi es ou non,  emanant des  tablissements d'enseignement et de recherche fran ais ou  trangers, des laboratoires publics ou priv es.

# Towards Learning Commonalities in SPARQL

Sara El Hassad, François Goasdoué, and Hélène Jaudoin

IRISA, Univ. Rennes 1, Lannion, France  
{sara.el-hassad,fg,helene.jaudoin}@irisa.fr

**Abstract.** Finding the commonalities between descriptions of data or knowledge is a foundational reasoning problem of Machine Learning, which amounts to computing a *least general generalization* (**lgg**) of such descriptions. We revisit this old problem in the popular conjunctive fragment of SPARQL, a.k.a. Basic Graph Pattern Queries (BGPQs). In particular, we define this problem in all its generality by considering general BGPQs, while the literature considers unary tree-shaped BGPQs only. Further, when *ontological knowledge* is available as RDF Schema constraints, we take advantage of it to devise much more pregnant **lggs**.

**Keywords:** BGP queries, RDF, RDFS, least general generalization

## 1 Introduction

Finding commonalities between descriptions of data and knowledge is a fundamental Machine Learning problem. It was formalized in early 70's as computing a *least general generalization* (**lgg**) of First Order Logic formulae [4].

We revisit this old reasoning problem in the setting of SPARQL, the RDF query language by W3C, which may have varied theoretical and practical applications. For instance, an **lgg** of queries is a best upper approximation thereof by a single query in *knowledge approximation*, is the largest set of commonalities that may be recommended for view materialization or shared processing in *query optimization*, or may help recommending users to each other, especially in a social context, if what they ask for is enough related in *recommendation*, etc.

Our contribution is to carefully study and define a pregnant notion of **lgg** for the well-established conjunctive fragment of SPARQL, a.k.a. Basic Graph Pattern Queries (BGPQs). Our results significantly depart from the literature by considering *general* BGPQs, instead of *unary tree-shaped* BGPQs [1, 3], and crucially by taking advantage of *ontological knowledge* formalized as RDF Schema constraints, when available. Proofs for this paper's claims are delegated to [2].

## 2 Preliminaries

The RDF data model allows specifying *RDF graphs*, which are sets of *well-formed triples* from  $(\mathcal{U} \cup \mathcal{B}) \times \mathcal{U} \times (\mathcal{U} \cup \mathcal{B} \cup \mathcal{L})$  with  $\mathcal{U}$ ,  $\mathcal{B}$ ,  $\mathcal{L}$  pairwise disjoint sets of URIs, of blank nodes (unknown values) and of literals (constants) respectively [5]. A triple  $(s, p, o)$  states that its *subject*  $s$  has *property*  $p$  whose value is the *object*  $o$ . Importantly, the RDF standard provides built-in property URIs to state facts for classes (unary relations) and properties (binary relations), called *RDF statements*, and ontological constraints relating classes and

		Rule [6]	Entailment rule
		<b>rdfs2</b>	$(p, \leftrightarrow_d, o), (s_1, p, o_1) \rightarrow (s_1, \tau, o)$
		<b>rdfs3</b>	$(p, \leftrightarrow_r, o), (s_1, p, o_1) \rightarrow (o_1, \tau, o)$
<b>RDF statement</b>	Triple		
Class assertion	$(s, \text{rdf:type}, o)$	<b>rdfs5</b>	$(p_1, \preceq_{\text{sp}}, p_2), (p_2, \preceq_{\text{sp}}, p_3) \rightarrow (p_1, \preceq_{\text{sp}}, p_3)$
Property assertion	$(s, p, o)$ with $p \neq \text{rdf:type}$	<b>rdfs7</b>	$(p_1, \preceq_{\text{sp}}, p_2), (s, p_1, o) \rightarrow (s, p_2, o)$
		<b>rdfs9</b>	$(s, \preceq_{\text{sc}}, o), (s_1, \tau, s) \rightarrow (s_1, \tau, o)$
<b>RDFS statement</b>	Triple		
Subclass	$(s, \text{rdfs:subClassOf}, o)$	<b>rdfs11</b>	$(s, \preceq_{\text{sc}}, o), (o, \preceq_{\text{sc}}, o_1) \rightarrow (s, \preceq_{\text{sc}}, o_1)$
Subproperty	$(s, \text{rdfs:subPropertyOf}, o)$	<b>ext1</b>	$(p, \leftrightarrow_d, o), (o, \preceq_{\text{sc}}, o_1) \rightarrow (p, \leftrightarrow_d, o_1)$
Domain typing	$(s, \text{rdfs:domain}, o)$	<b>ext2</b>	$(p, \leftrightarrow_r, o), (o, \preceq_{\text{sc}}, o_1) \rightarrow (p, \leftrightarrow_r, o_1)$
Range typing	$(s, \text{rdfs:range}, o)$	<b>ext3</b>	$(p, \preceq_{\text{sp}}, p_1), (p_1, \leftrightarrow_d, o) \rightarrow (p, \leftrightarrow_d, o)$
		<b>ext4</b>	$(p, \preceq_{\text{sp}}, p_1), (p_1, \leftrightarrow_r, o) \rightarrow (p, \leftrightarrow_r, o)$

**Table 1.** RDF & RDFS statements.**Table 2.** Sample RDF entailment rules.

properties, called *RDF Schema (RDFS) statements*, as shown in Table 1. Hereafter, we use the shorthands  $\tau$ ,  $\preceq_{\text{sc}}$ ,  $\preceq_{\text{sp}}$ ,  $\leftrightarrow_d$  and  $\leftrightarrow_r$  for the built-in property URIs `rdf:type`, `rdfs:subClassOf`, `rdfs:subPropertyOf`, `rdfs:domain` and `rdfs:range` respectively. The semantics of an RDF graph  $\mathcal{G}$  is its *saturation* (a.k.a. *closure*), denoted  $\mathcal{G}^\infty$ , defined as the set of  $\mathcal{G}$  triples together with all the *implicit* triples that can be derived from them and entailment rules from the RDF standard. Table 2 shows some rules that use RDFS constraints to derive implicit facts and constraints.

The *Basic Graph Pattern Queries (BGPQs)* form the conjunctive (or select-project-join) fragment of SPARQL. A BGPQ is of the form  $q(\bar{x}) \leftarrow t_1, \dots, t_\alpha$ , where  $\{t_1, \dots, t_\alpha\}$  is a subset of  $(\mathcal{U} \cup \mathcal{B} \cup \mathcal{V}) \times (\mathcal{U} \cup \mathcal{V}) \times (\mathcal{U} \cup \mathcal{B} \cup \mathcal{L} \cup \mathcal{V})$  with  $\mathcal{V}$  a set of variables pairwise disjoint with  $\mathcal{U}, \mathcal{B}, \mathcal{L}$ , and  $\bar{x}$  is a subset of the variables occurring in  $t_1, \dots, t_\alpha$  called *answer variables*; for boolean queries,  $\bar{x}$  is empty. The head of  $q$  is  $head(q) = q(\bar{x})$  and the body of  $q$  is  $body(q) = \{t_1, \dots, t_\alpha\}$ .

Two standard reasoning tasks characterize how RDF graphs contribute to queries. *Query entailment* indicates if an RDF graph holds some answer(s) to a query. Given a BGPQ  $q$ , an RDF graph  $\mathcal{G}$  and a set  $\mathcal{R}$  of RDF entailment rules,  $\mathcal{G}$  *entails*  $q$ , noted  $\mathcal{G} \models_{\mathcal{R}} q$ , iff  $\mathcal{G} \models_{\mathcal{R}} body(q)$  holds, i.e., there exists a homomorphism  $\phi$  from  $q$ 's variables and blank nodes to  $\mathcal{G}^\infty$ 's values (URIs, literals and blank nodes) such that  $[body(q)]_\phi \subseteq \mathcal{G}^\infty$ . Importantly,  $\mathcal{G} \models_{\mathcal{R}} q$  holds iff  $\mathcal{G}^\infty \models_\emptyset q$  holds. We note  $\mathcal{G} \models_{\mathcal{R}}^\phi q$  the entailment  $\mathcal{G} \models_{\mathcal{R}} q$  due to the homomorphism  $\phi$ . *Query answering* identifies *all* the answers to a query that an RDF graph holds. Given a BGPQ  $q$  with head  $q(\bar{x})$ , the *answer set of  $q$  against  $\mathcal{G}$*  is  $q(\mathcal{G}) = \{(\bar{x})_\phi \mid \mathcal{G} \models_{\mathcal{R}}^\phi body(q)\}$  where we denote by  $(\bar{x})_\phi$  the tuple of  $\mathcal{G}^\infty$  values obtained by replacing every answer variable  $x_i \in \bar{x}$  by its image  $\phi(x_i)$ .

Finally, queries can be compared through the generalization/specialization relationship of *entailment between queries*, which is the obvious adaptation of query entailment to the presence of variables in queries. Given two BGPQs  $q, q'$  with *same* arity, whose heads are  $q(\bar{x})$  and  $q'(\bar{x}')$ , and a set  $\mathcal{R}$  of RDF entailment rules at hand,  $q$  *entails*  $q'$ , noted  $q \models_{\mathcal{R}} q'$ , iff  $body(q) \models_{\mathcal{R}}^\phi body(q')$  with  $(\bar{x}')_\phi = \bar{x}$ .

### 3 Least General Generalization of BGPQs

A *least general generalization* (**lgg**) of two<sup>1</sup> descriptions  $d_1, d_2$  is a most specific description  $d$  generalizing  $d_1, d_2$  for some generalization/specialization rela-

<sup>1</sup> This easily generalizes to **lggs** of  $n$  descriptions [2].

tion [4]. In our SPARQL setting, we use BGPQs as descriptions and entailment between BGPQs as generalization/specialization relation:

**Definition 1 (l<sub>gg</sub> of BGPQs).** Let  $q_1, q_2$  be two BGPQs with the same arity and  $\mathcal{R}$  a set of RDF entailment rules.

- A generalization of  $q_1, q_2$  is a BGPQ  $q_g$  such that  $q_1 \models_{\mathcal{R}} q_g$  and  $q_2 \models_{\mathcal{R}} q_g$ .
- A least general generalization of  $q_1, q_2$  is a generalization  $q_{\text{l<sub>gg of  $q_1, q_2$  such that for any other generalization  $q_g$  of  $q_1, q_2$ :  $q_{\text{l<sub>gg.</sub>$</sub>$

Unfortunately, this natural definition is of limited practical interest as exemplified next. Consider the BGPQs  $q_1$  and  $q_2$  in Figure 1, which respectively ask for the conference papers having some contact author, and for the journal papers having some author. Clearly, with the RDF entailment rules shown in Table 2, an l<sub>gg</sub> of  $q_1$  and  $q_2$  is the *very* general BGPQ  $q_{\text{l<sub>gg asking for *the resources having some type*. However, by considering the ontological constraints displayed in Figure 1 that hold in the scientific publication domain, i.e., the context in which the queries are asked, a more pregnant l<sub>gg</sub> would be  $q_{\text{l<sub>gg asking for *the publications having some researcher as author*, since (i) having a contact author is having an author, (ii) only publications have authors, (iii) only researchers are authors, and (iv) conference (resp. journal) papers are publications.</sub>$</sub>$

To devise such elaborate l<sub>ggs</sub> that rely on ontological knowledge, we revisit the notion of entailment between BGPQs in order to account for extra RDFS constraints. We first complement a BGPQ w.r.t. ontological knowledge:

**Definition 2 (BGPQ saturation w.r.t. RDFS constraints).** Let  $\mathcal{R}$  be a set of RDF entailment rules,  $\mathcal{O}$  a set of RDFS statements, and  $q$  a BGPQ the body of which, without loss of generality, does not contain blank nodes<sup>2</sup>. The saturation of  $q$  w.r.t.  $\mathcal{O}$ , denoted  $q_{\mathcal{O}}^{\infty}$ , is a BGPQ with the same answer variables as  $q$  and whose body, denoted  $\text{body}(q_{\mathcal{O}}^{\infty})$ , is the maximal subset of  $(\mathcal{O} \cup \text{body}(q))^{\infty}$  such that for any of its subset  $\mathcal{S}$ : if  $\mathcal{O} \models_{\mathcal{R}} \mathcal{S}$  holds then  $\text{body}(q) \models_{\mathcal{R}} \mathcal{S}$  holds.

Intuitively, the saturation of a BGPQ comprises all the triples in the saturation of its body augmented with the constraints, except those triples that only follow from the ontological constraints, i.e., which are not related to what the query is asking for. This corresponds to the *non-hatched* subset of  $(\mathcal{O} \cup \text{body}(q))^{\infty}$  shown in Figure 2. This Figure also displays the saturations  $q_{1\mathcal{O}}^{\infty}, q_{2\mathcal{O}}^{\infty}$  of the two BGPQs  $q_1, q_2$  w.r.t. the constraints  $\mathcal{O}$  shown in Figure 1. Importantly, we proved that a BGPQ and its saturation w.r.t. ontological constraints are *equivalent for the central RDF reasoning tasks of query entailment and query answering* [2]:

**Theorem 1.** Let  $\mathcal{R}$  be a set of RDF entailment rules,  $\mathcal{O}$  a set of RDFS statements, and  $q$  a BGPQ whose saturation w.r.t.  $\mathcal{O}$  is  $q_{\mathcal{O}}^{\infty}$ . For any RDF graph  $\mathcal{G}$  whose set of RDFS statements is  $\mathcal{O}$ , (i)  $\mathcal{G} \models_{\mathcal{R}} q$  holds iff  $\mathcal{G} \models_{\mathcal{R}} q_{\mathcal{O}}^{\infty}$  holds, and (ii)  $q(\mathcal{G}) = q_{\mathcal{O}}^{\infty}(\mathcal{G})$  holds.

Building on BGPQ saturation, we generalize entailment between BGPQs to:

<sup>2</sup> In SPARQL queries, blank nodes are equivalent to non-answer variables [7].

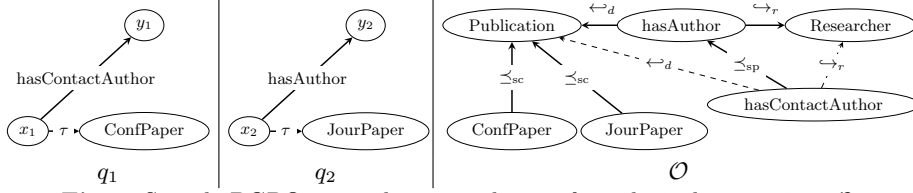


Fig. 1. Sample BGPQs  $q_1$  and  $q_2$ ; sample set of ontological constraints  $\mathcal{O}$ .

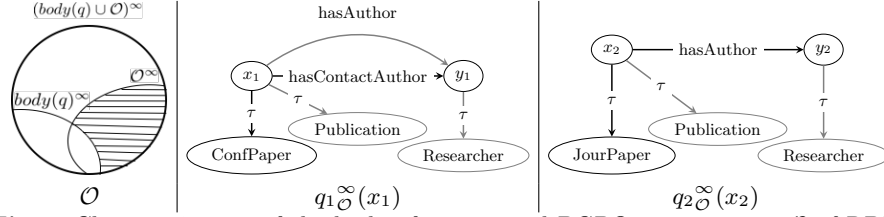


Fig. 2. Characterization of the body of a saturated BGPQ  $q$  w.r.t. a set  $\mathcal{O}$  of RDFS constraints; saturations of  $q_1$  and  $q_2$  w.r.t.  $\mathcal{O}$ , triples in grey are added by saturation.

**Definition 3 (Entailment between BGPQs w.r.t. RDFS constraints).** Given a set  $\mathcal{R}$  of RDF entailment rules, a set  $\mathcal{O}$  of RDFS statements, and two BGPQs  $q$  and  $q'$  with the same arity,  $q$  entails  $q'$  w.r.t.  $\mathcal{O}$ , denoted  $q \models_{\mathcal{R}, \mathcal{O}} q'$ , iff  $q_{\mathcal{O}}^{\infty} \models_{\emptyset} q'$  holds.

When  $\mathcal{O}$  is empty, the above definition coincides with standard entailment between BGPQs. Further, we proved fundamental properties for a BGPQ entailed by another w.r.t. ontological constraints: the former *generalizes* the latter for the central RDF reasoning tasks of query entailment and query answering [2]:

**Theorem 2.** Let  $\mathcal{R}$  be a set of RDF entailment rules,  $\mathcal{O}$  a set of RDFS statements, and two BGPQs  $q$  and  $q'$  such that  $q \models_{\mathcal{R}, \mathcal{O}} q'$ . For any RDF graph  $\mathcal{G}$  whose set of RDFS statements is  $\mathcal{O}$ , (i) if  $\mathcal{G} \models_{\mathcal{R}} q$  holds then  $\mathcal{G} \models_{\mathcal{R}} q'$  holds, and (ii)  $q(\mathcal{G}) \subseteq q'(\mathcal{G})$  holds.

With the above notion of entailment between BGPQs endowed with ontological knowledge, we revise the definition of **lgg** (Definition 1) in order to use  $\models_{\mathcal{R}, \mathcal{O}}$  instead of  $\models_{\mathcal{R}}$ . We therefore propose to investigate as next challenge:

*Problem 1.* Given two BGPQs  $q_1, q_2$  with same arity, a set  $\mathcal{O}$  of RDFS statements, and a set  $\mathcal{R}$  of RDF entailment rules, compute an **lgg** of  $q_1, q_2$  w.r.t.  $\mathcal{O}$ .

## References

1. Bühmann, L., Lehmann, J., Westphal, P.: DL-learner - a framework for inductive learning on the semantic web. *J. Web Semantics* 39 (2016)
2. El Hassad, S., Goasdoué, F., Jaudoin, H.: Learning commonalities in RDF and SPARQL (research report). <https://hal.inria.fr/hal-01386237> (2016)
3. Lehmann, J., Bühmann, L.: Autosparql: Let users query your knowledge base. In: *ESWC* (2011)
4. Plotkin, G.D.: A note on inductive generalization. *Machine Intelligence* 5 (1970)
5. Resource description framework 1.1. <https://www.w3.org/TR/rdf11-concepts>
6. RDF 1.1 semantics. <https://www.w3.org/TR/rdf11-mt/>
7. SPARQL. <http://www.w3.org/TR/rdf-sparql-query>