



# Learning Commonalities in RDF and SPARQL

Sara El Hassad, François Goasdoué, Hélène Jaudoin

## ► To cite this version:

Sara El Hassad, François Goasdoué, Hélène Jaudoin. Learning Commonalities in RDF and SPARQL. [Research Report] Université Rennes 1. 2016. hal-01386237v3

**HAL Id: hal-01386237**

**<https://inria.hal.science/hal-01386237v3>**

Submitted on 9 Mar 2017 (v3), last revised 23 May 2017 (v4)

**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.



Distributed under a Creative Commons Attribution 4.0 International License

# Learning Commonalities in RDF and SPARQL

v3 of the Research Report made available online in early september 2016 (v1).

The main change in v2 was the addition of Section 5. In v3, some results and discussions were added to the core Sections 3 and 4; also details were added to the description of the related work in Section 7.

Sara El Hassad  
IRISA, Univ. Rennes 1  
Lannion, France  
sara.el-hassad@irisa.fr

François Goasdoué  
IRISA, Univ. Rennes 1  
Lannion, France  
fg@irisa.fr

Hélène Jaudoin  
IRISA, Univ. Rennes 1  
Lannion, France  
helene.jaudoin@irisa.fr

## ABSTRACT

Finding commonalities between descriptions of data or knowledge is a fundamental task in Machine Learning. The formal notion characterizing precisely such commonalities is known as *least general generalization* of descriptions and was introduced by G. Plotkin in the early 70's, in First Order Logic.

Identifying least general generalizations has a large scope of database applications ranging from query optimization (e.g., to share commonalities between queries in view selection or multi-query optimization) to recommendation in social networks (e.g., to establish connections between users based on their commonalities between profiles or searches).

To the best of our knowledge, this is the first work that revisits the notion of least general generalizations in the entire Resource Description Framework (RDF) and popular conjunctive fragment of SPARQL, a.k.a. Basic Graph Pattern (BGP) queries. Our contributions include the definition and the computation of least general generalizations in these two settings, which amounts to finding the largest set of commonalities between incomplete databases and conjunctive queries, under deductive constraints. We also provide an experimental assessment of our technical contributions.

## 1. INTRODUCTION

Finding the commonalities between descriptions of data or knowledge is a foundational reasoning problem of Machine Learning, which was formalized in the early 70's as computing a *least general generalization* (lgg) of such descriptions [40, 41]. Since the early 90's, this problem has also received consideration in the Knowledge Representation field, where least general generalization were rebaptized *least common subsumers* [20], in Description Logics, e.g., [20, 11, 12, 51] and in Conceptual Graphs [19].

Interestingly, this problem has a large scope of applications when transposed in a database setting, i.e., when descriptions are expressed in a given data model or query lan-

guage. In *Query Optimization*, lggs of subsets of a query workload  $\mathcal{W}$  may be used to identify candidate views or potential query sharing, with which  $\mathcal{W}$  can be efficiently evaluated, in the View Selection problem or in the Multi-Query Optimization problem, respectively. In *Recommendation*, in particular in a social context, lggs of sets of user descriptions (i.e., profiles) may help recommending users other users, or to form a community, when they share enough interests. Also, lggs of sets of queries issued by distinct users may help recommending users to each other, if what they ask for is enough related. In *Exploration*, lggs of datasets may be used to classify/categorize them w.r.t. their common information, to identify common social graph patterns between organizations (e.g., criminal ones), or to help identifying new potential links between datasets in the Linked Data Cloud.

**Contributions.** In this work, we revisit the problem of computing least general generalizations in the RDF data model and its SPARQL query language, the prominent Semantic Web standards of the W3C. This is a necessary first step towards using lggs within central database problems. More precisely, our contributions to this problem are:

1. We define and study the problem of computing an lgg of RDF graphs in the *entire RDF standard*: we do not restrict RDF graphs in any way, i.e., neither their structure nor their semantics defined upon RDF entailment (inference). The recent work [22, 21] brings a limited solution to the problem. It allows finding the commonalities between *single entities* extracted from RDF graphs (e.g., users in a social network), *ignoring* RDF entailment. In contrast, we further aim at considering the problem in all its generality, i.e., finding the commonalities between *general* RDF graphs, hence modeling *multiple interrelated entities* (e.g., social networks of users), accurately w.r.t. their standard semantics.
2. We define and study the problem of computing an lgg of *general* SPARQL conjunctive queries, a.k.a. Basic Graph Pattern Queries (BGPQs), while the literature only considers unary tree-shaped conjunctive queries (UTCQ) [33]. Further, when available, we devise how to use background knowledge, formalized as ontological constraints modeling the application domain, in order to compute much more pregnant lggs.
3. We provide algorithms for our solutions to these two problems. In particular, since RDF graphs may be

RDF statement	Triple	Relational assertion
Class assertion	$(s, \text{rdf:type}, o)$	$o(s)$
Property assertion	$(s, p, o)$ with $p \neq \text{rdf:type}$	$p(s, o)$

RDFS statement	Triple	Relational constraint
Subclass	$(s, \text{rdfs:subClassOf}, o)$	$s \subseteq o$
Subproperty	$(s, \text{rdfs:subPropertyOf}, o)$	$s \subseteq o$
Domain typing	$(s, \text{rdfs:domain}, o)$	$\Pi_{\text{domain}}(s) \subseteq o$
Range typing	$(s, \text{rdfs:range}, o)$	$\Pi_{\text{range}}(s) \subseteq o$

Table 1: RDF (top) & RDFS (bottom) statements.

large, we provide algorithms that allow computing **lggs** of *small-to-huge general RDF graphs* (i.e., that fit either in memory, in data management systems or in MapReduce clusters) w.r.t. *any set of entailment rules* from the RDF standard.

4. We provide experiments using a synthetic (LUBM) dataset, which show to which extent **lggs** of BGPQs are much more precise when background knowledge is available.

The paper is organized as follows. In Section 2, we introduce the RDF data model and its SPARQL query language. Then, we study the problem of computing an **lgg** of RDF graphs in Section 3, and of BGPQs in Section 4. In Section 5, we present algorithms for the problems studied in the two preceding Sections. We report on the experiments we conducted in Section 6. Finally, we present related works in Section 7 and we conclude and draw some perspectives in Section 8.

## 2. PRELIMINARIES

We introduce the *RDF data model* in Section 2.1. Then, in Section 2.2, we present the SPARQL conjunctive queries, a.k.a. *Basic Graph Pattern queries (BGPQs)*.

### 2.1 Resource Description Framework (RDF)

#### 2.1.1 RDF graphs

The RDF data model allows specifying *RDF graphs*. An RDF graph is a set of *triples* of the form  $(s, p, o)$ . A triple states that its *subject*  $s$  has the *property*  $p$ , the value of which is the *object*  $o$ . Triples are built using three pairwise disjoint sets: a set  $\mathcal{U}$  of *uniform resources identifiers (URIs)*, a set  $\mathcal{L}$  of *literals* (constants), and a set  $\mathcal{B}$  of *blank nodes* allowing to support *incomplete information*. Blank nodes are local identifiers for missing values within an RDF graph (unknown URIs or literals); they are the RDF counterparts of labelled nulls or existential variables in incomplete relational databases [30, 10], as shown in [27]. *Well-formed triples*, as per the RDF specification [47], belong to  $(\mathcal{U} \cup \mathcal{B}) \times \mathcal{U} \times (\mathcal{U} \cup \mathcal{L} \cup \mathcal{B})$ ; we only consider such triples hereafter.

**Notations.** We use  $s, p, o$  in triples as placeholders. We note  $\text{Val}(\mathcal{G})$  the set of *values* occurring in an RDF graph  $\mathcal{G}$ , i.e., the URIs, literals and blank nodes; we note  $\text{B1}(\mathcal{G})$  the set of blank nodes occurring in  $\mathcal{G}$ . A blank node is written  $b$  possibly with a subscript, and a literal is a string between quotes. For instance, the triples  $(b, \text{hasTitle}, \text{"LGG in RDF"})$  and  $(b, \text{hasContactAuthor}, b_1)$  state that *something* ( $b$ ) *entitled "LGG in RDF" has somebody* ( $b_1$ ) *as contact author*.

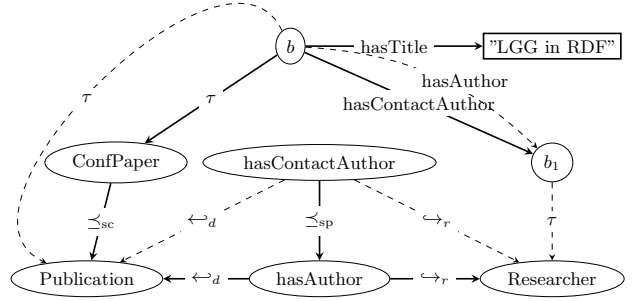


Figure 1: Sample RDF graph  $\mathcal{G}$ .

A triple models an assertion, either for a *class* (unary relation) or for a *property* (binary relation). Table 1 (top) shows the use of triples to state such assertions, as well as the relational assertions to which they correspond. The RDF standard [47] provides built-in classes and properties, as URIs within the **rdf** and **rdfs** pre-defined namespaces, e.g., **rdf:type** which can be used to state that the above  $b$  is a conference paper with the triple  $(b, \text{rdf:type}, \text{ConfPaper})$ .

#### 2.1.2 Adding ontological knowledge to RDF graphs

An essential feature of RDF is the possibility to enhance the descriptions in RDF graphs by declaring *ontological constraints* between the classes and properties they use. This is achieved with *RDF Schema (RDFS)* statements, which are triples using particular built-in properties. Table 1 (bottom) lists the allowed constraints, the triples to state them, as well as the *deductive* relational constraints to which they correspond; *domain* and *range* denote respectively the first and second attribute of every property. For example, the triple  $(\text{ConfPaper}, \text{rdfs:subClassOf}, \text{Publication})$  states that *conference papers are publications*, the triple  $(\text{hasContactAuthor}, \text{rdfs:subPropertyOf}, \text{hasAuthor})$  states that *having a contact author is having an author*, the triple  $(\text{hasAuthor}, \text{rdfs:domain}, \text{Publication})$  states that *only publications have authors* while the triple  $(\text{hasAuthor}, \text{rdfs:range}, \text{Researcher})$  states that *only researchers are authors of something*.

**Notations.** From now, for conciseness, we use the following shorthands for built-in properties:  $\tau$  for **rdf:type**,  $\preceq_{sc}$  for **rdfs:subClassOf**,  $\preceq_{sp}$  for **rdfs:subPropertyOf**,  $\leftrightarrow_d$  for **rdfs:domain**, and  $\leftrightarrow_r$  for **rdfs:range**.

Figure 1 displays the usual representation of the RDF graph  $\mathcal{G}$  made of the seven above-mentioned triples, which are called the *explicit triples* of  $\mathcal{G}$ . A triple  $(s, p, o)$  corresponds to the  $p$ -labeled directed edge from the  $s$  node to the  $o$  node. Explicit triples are shown as solid edges, while the *implicit ones*, which are derived using ontological constraints (see below), are shown as dashed edges.

Importantly, the semantics of ontological constraints differ from the typical relational database setting, where they are interpreted under the *closed world assumption (CWA)*, i.e., as *integrity constraints*. In RDF, they are interpreted under the *open world assumption (OWA)*, i.e., as *deductive constraints*. Under OWA, *all* the triples that can be deduced from the constraints, called *implicit triples*, are assumed to be part of the RDF graph even though they are not explicitly present in it. For instance, consider the RDF graph  $\mathcal{G}$  in Figure 1. The ontological constraint  $(\text{hasContactAuthor}, \text{rdfs:subPropertyOf}, \text{hasAuthor})$  is violated under CWA, as

Rule name [48]	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>	$(p_1, \preceq_{sp}, p_2), (p_2, \preceq_{sp}, p_3) \rightarrow (p_1, \preceq_{sp}, p_3)$
<b>rdfs7</b>	$(p_1, \preceq_{sp}, p_2), (s, p_1, o) \rightarrow (s, p_2, o)$
<b>rdfs9</b>	$(s, \preceq_{sc}, o), (s_1, \tau, s) \rightarrow (s_1, \tau, o)$
<b>rdfs11</b>	$(s, \preceq_{sc}, o), (o, \preceq_{sc}, o_1) \rightarrow (s, \preceq_{sc}, o_1)$
<b>ext1</b>	$(p, \leftarrow_d, o), (o, \preceq_{sc}, o_1) \rightarrow (p, \leftarrow_d, o_1)$
<b>ext2</b>	$(p, \leftarrow_r, o), (o, \preceq_{sc}, o_1) \rightarrow (p, \leftarrow_r, o_1)$
<b>ext3</b>	$(p, \preceq_{sp}, p_1), (p_1, \leftarrow_d, o) \rightarrow (p, \leftarrow_d, o)$
<b>ext4</b>	$(p, \preceq_{sp}, p_1), (p_1, \leftarrow_r, o) \rightarrow (p, \leftarrow_r, o)$

Table 2: Sample set of RDF entailment rules.

the triple  $(b, \text{hasContactAuthor}, b_1)$  is explicit in  $\mathcal{G}$  while  $(b, \text{hasAuthor}, b_1)$  is not. By contrast, under OWA, the same constraint allows deriving the “missing” triple  $(b, \text{hasAuthor}, b_1)$ , making it implicit within  $\mathcal{G}$ . Further, this implicit triple together with the explicit ontological constraint  $(\text{hasAuthor}, \text{rdfs:range}, \text{Researcher})$  yields the implicit triple  $(b_1, \text{rdf:type}, \text{Researcher})$ .

### 2.1.3 Deriving the implicit triples of an RDF graph

The RDF standard defines a set of *entailment rules* in order to derive automatically *all* the triples that are implicit to an RDF graph. Table 2 shows the strict subset of these rules that we will use to illustrate important notions as well as our contributions in the next sections; importantly, our contributions hold for the *whole* set of entailment rules of the RDF standard, and any subset of thereof. The rules in Table 2 concern the derivation of implicit triples using ontological constraints (i.e., *RDFS statements*). They encode the *propagation* of assertions through constraints (**rdfs2**, **rdfs3**, **rdfs7**, **rdfs9**), the *transitivity* of the  $\preceq_{sp}$  and  $\preceq_{sc}$  constraints (**rdfs5**, **rdfs11**), the *complementation* of domains or ranges through  $\preceq_{sc}$  (**ext1**, **ext2**), and the *inheritance* of domains and of ranges through  $\preceq_{sp}$  (**ext3**, **ext4**).

The *saturation* (a.k.a. *closure*) of an RDF graph  $\mathcal{G}$  w.r.t. a set  $\mathcal{R}$  of RDF entailment rules, is the RDF graph  $\mathcal{G}^\infty$  obtained by adding to  $\mathcal{G}$  all the implicit triples that can be derived from  $\mathcal{G}$  using  $\mathcal{R}$ . Roughly speaking, the saturation  $\mathcal{G}^\infty$  *materializes* the semantics of  $\mathcal{G}$ . It corresponds to the fixpoint obtained by repeatedly applying the rules in  $\mathcal{R}$  to  $\mathcal{G}$  in a forward-chaining fashion, a.k.a. *chasing*  $\mathcal{G}$  with  $\mathcal{R}$  in the database parlance. In RDF, the saturation is *always* finite and unique (up to blank node renaming), and does not contain implicit triples [47, 48].

The saturation of the RDF graph  $\mathcal{G}$  shown in Figure 1 corresponds to the RDF graph  $\mathcal{G}^\infty$  in which all the  $\mathcal{G}$  implicit triples have been made explicit. It is worth noting how, starting from  $\mathcal{G}$ , applying RDF entailment rules *mechanizes* the construction of  $\mathcal{G}^\infty$ . For instance, recall the reasoning sketched in Section 2.1.2 for deriving the triple  $(b_1, \text{rdf:type}, \text{Researcher})$ . This is automated by the following sequence of applications of RDF entailment rules:  $(\text{hasContactAuthor}, \text{rdfs:subPropertyOf}, \text{hasAuthor})$  and  $(b, \text{hasContactAuthor}, b_1)$  trigger the **rdfs7** rule which adds  $(b, \text{hasAuthor}, b_1)$  to the RDF graph. In turn, this new triple together with  $(\text{hasAuthor}, \text{rdfs:range}, \text{Researcher})$  triggers the **rdfs3** rule which adds  $(b_1, \text{rdf:type}, \text{Researcher})$ .

### 2.1.4 Comparing RDF graphs

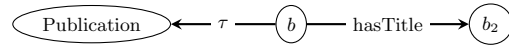


Figure 2: Sample RDF graph  $\mathcal{G}'$ .

The RDF standard defines a generalization/specialization relationship between two RDF graphs, called *entailment between graphs*. Roughly speaking, an RDF graph  $\mathcal{G}$  is more specific than another RDF graph  $\mathcal{G}'$ , or equivalently  $\mathcal{G}'$  is more general than  $\mathcal{G}$ , whenever there is an embedding of  $\mathcal{G}'$  into the *saturation* of  $\mathcal{G}$ , i.e., the complete set of triples that  $\mathcal{G}$  models.

More formally, given any subset  $\mathcal{R}$  of RDF entailment rules, an RDF graph  $\mathcal{G}$  *entails* an RDF graph  $\mathcal{G}'$ , denoted  $\mathcal{G} \models_{\mathcal{R}} \mathcal{G}'$ , iff there exists an homomorphism  $\phi$  from  $\text{B1}(\mathcal{G}')$  to  $\text{Val}(\mathcal{G}^\infty)$  such that  $[\mathcal{G}']_\phi \subseteq \mathcal{G}^\infty$ , where  $[\mathcal{G}']_\phi$  is the RDF graph obtained from  $\mathcal{G}'$  by replacing every blank node  $b$  by its image  $\phi(b)$ .

Figure 2 shows an RDF graph  $\mathcal{G}'$  that is entailed by the RDF graph  $\mathcal{G}$  in Figure 1 w.r.t. the RDF entailment rules displayed in Table 2. In particular,  $\mathcal{G} \models_{\mathcal{R}} \mathcal{G}'$  holds for the homomorphism  $\phi$  such that:  $\phi(b) = b$  and  $\phi(b_2) = \text{"LGG in RDF"}$ . By contrast, when  $\mathcal{R}$  is empty, this is not the case (i.e.,  $\mathcal{G} \not\models_{\mathcal{R}} \mathcal{G}'$ ), as the dashed edges in  $\mathcal{G}$  are not materialized by saturation, hence the  $\mathcal{G}'$  triple  $(b, \tau, \text{Publication})$  cannot have an image in  $\mathcal{G}$  through some homomorphism.

**Notations.** When relevant to the discussion, we designate by  $\mathcal{G} \models_{\mathcal{R}}^{\phi} \mathcal{G}'$  the fact that the entailment  $\mathcal{G} \models_{\mathcal{R}} \mathcal{G}'$  holds due to the graph homomorphism  $\phi$ . Also, when RDF entailment rules are disregarded, i.e.,  $\mathcal{R} = \emptyset$ , we note the entailment relation  $\models$  (without indicating the rule set under consideration). Observe that, in this particular case, the entailment between RDF graphs collapses to the database notion of *containment* between two incomplete instances, each stored into a **Triple**(**s**, **p**, **o**) relation.

Importantly, some remarkable properties follow directly from the definition of entailment between two RDF graphs:

1.  $\mathcal{G}$  and  $\mathcal{G}^\infty$  are equivalent, noted  $\mathcal{G} \equiv_{\mathcal{R}} \mathcal{G}^\infty$ , since clearly  $\mathcal{G} \models_{\mathcal{R}} \mathcal{G}^\infty$  and  $\mathcal{G}^\infty \models_{\mathcal{R}} \mathcal{G}$  hold,
2.  $\mathcal{G} \models_{\mathcal{R}} \mathcal{G}'$  iff  $\mathcal{G}^\infty \models \mathcal{G}'$  holds.

In particular, the second above property points out how entailment reduces to standard database containment once saturation is computed.

## 2.2 SPARQL conjunctive queries

### 2.2.1 Basic graph pattern queries

The well-established conjunctive fragment of SPARQL queries, a.k.a. *basic graph pattern queries* (*BGPQs*), is the counterpart of the relational select-project-join queries; it is the most widely used subset of SPARQL queries in real-world applications [37].

A *Basic Graph Pattern* (*BGP*) is a set of *triple patterns*, or simply triples by a slight abuse of language. They generalize RDF triples by allowing the use of variables. Given a set  $\mathcal{V}$  of variables, pairwise disjoint with  $\mathcal{U}$ ,  $\mathcal{L}$  and  $\mathcal{B}$ , triple patterns belong to:  $(\mathcal{V} \cup \mathcal{U} \cup \mathcal{B}) \times (\mathcal{V} \cup \mathcal{U}) \times (\mathcal{V} \cup \mathcal{U} \cup \mathcal{L} \cup \mathcal{B})$ .

**Notations.** We adopt the usual conjunctive query notation  $q(\bar{x}) \leftarrow t_1, \dots, t_\alpha$ , where  $\{t_1, \dots, t_\alpha\}$  is a BGP; the query head variables  $\bar{x}$  are called *answer variables*, and form a subset of the variables occurring in  $t_1, \dots, t_\alpha$ ; for boolean

queries,  $\bar{x}$  is empty. The head of  $q$ , noted  $head(q)$ , is  $q(\bar{x})$ , and the body of  $q$ , denoted  $body(q)$ , is the BGP  $\{t_1, \dots, t_\alpha\}$ . We use  $x$  and  $y$  in queries, possibly with subscripts, for answer and non-answer variables respectively. Finally, we denote by  $\text{VarBl}(q)$  the set of variables *and* blank nodes occurring in the query  $q$ ; we note  $\text{Val}(q)$  the set of all its values, i.e., URIs, blank nodes, literals and variables.

### 2.2.2 Entailing and answering queries

Two important notions characterize how an RDF graph contributes to a query.

**Query entailment.** The weaker notion indicates whether or not an RDF graph holds some answer(s) to a query. It generalizes entailment between RDF graphs, to account for the presence of variables in the query body, for establishing whether a graph entails a query, i.e., whether the query embeds in that graph.

Formally, given a BGPQ  $q$ , an RDF graph  $\mathcal{G}$  and a set  $\mathcal{R}$  of RDF entailment rules,  $\mathcal{G}$  *entails*  $q$ , denoted  $\mathcal{G} \models_{\mathcal{R}} q$ , iff  $\mathcal{G} \models_{\mathcal{R}} body(q)$  holds, i.e., there exists a homomorphism  $\phi$  from  $\text{VarBl}(q)$  to  $\text{Val}(\mathcal{G}^\infty)$  such that  $[body(q)]_\phi \subseteq \mathcal{G}^\infty$ .

The RDF graph  $\mathcal{G}$  in Figure 1 entails the query  $q(x_1, x_2) \leftarrow (x_1, \tau, x_2)$  asking for all the resources and their classes, for instance because of the homomorphism  $\phi$  such that  $\phi(x_1) = b$  and  $\phi(x_2) = \text{ConfPaper}$ . Observe that this entailment holds for any subset of RDF entailment rules, since the above  $\phi$  already holds for  $\mathcal{R} = \emptyset$ , i.e., considering only the explicit triples in Figure 1.

**Notations.** Similarly to RDF graph entailment, we denote by  $\mathcal{G} \models_{\mathcal{R}}^\phi q$  that the entailment  $\mathcal{G} \models_{\mathcal{R}} q$  holds due to the homomorphism  $\phi$ .

**Query answering.** The stronger notion characterizing how an RDF graph contributes to a query, which builds on the preceding one, identifies *all* the query answers that the graph holds. Formally, assuming the head of a BGPQ  $q$  is  $q(\bar{x})$ , 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  $\phi(x_i)$ . Observe that when  $\mathcal{R}$  is empty, BGPQ query answering coincide with *standard relational conjunctive query evaluation*.

The answer set to the above query  $q(x_1, x_2) \leftarrow (x_1, \tau, x_2)$  against the RDF graph  $\mathcal{G}$  in Figure 1 is:

- $\{(b, \text{ConfPaper}), (b, \text{Publication}), (b_1, \text{Researcher})\}$  for  $\mathcal{R}$  the set of entailment rules shown in Table 2, i.e., considering the explicit *and* implicit triples in Figure 1,
- $\{(b, \text{ConfPaper})\}$  for  $\mathcal{R} = \emptyset$ , i.e., considering only the explicit triples in Figure 1.

Finally, remark that query entailment and query answering *treat the blank nodes in a query exactly as non-answer variables* [49].

### 2.2.3 Comparing queries

Similarly to RDF graphs, queries can be compared through the generalization/specialization relationship of *entailment between queries*. This relationship is the counterpart of that *query containment* in the relational database setting.

Let  $q, q'$  be two BGPQs with *same* arity, whose heads are  $q(\bar{x})$  and  $q'(\bar{x}')$ , and  $\mathcal{R}$  the set of RDF entailment rules under consideration.  $q$  *entails*  $q'$ , denoted  $q \models_{\mathcal{R}} q'$ , iff  $body(q) \models_{\mathcal{R}}^\phi body(q')$  with  $(\bar{x}')_\phi = \bar{x}$  holds.

Here,  $body(q) \models_{\mathcal{R}}^\phi body(q')$  is the adaptation of the above-mentioned entailment relationships between RDF graphs to the fact that the query bodies may feature variables, i.e.,  $\phi$  is a homomorphism from  $\text{VarBl}(q')$  to  $\text{Val}(q)$  such that  $[body(q')]_\phi \subseteq body(q)^\infty$ ; the saturation of a BGP, here  $body(q)^\infty$ , is the generalization of RDF graph saturation that treats variables as blank nodes, since they both equivalently model unknown information within BGPs.

For instance, the query  $q_1(x) \leftarrow (x, \tau, \text{ConfPaper})$ ,  $(x, \text{hasContactAuthor}, y)$  entails the query  $q_2(x) \leftarrow (x, \tau, y)$  with  $\phi(x) = x$ ,  $\phi(y) = \text{ConfPaper}$  and any entailment rule set. However, it is worth noting that, *counterintuitively*,  $q_1$  does not entail the query  $q_3(x) \leftarrow (x, \tau, \text{Publication})$ ,  $(x, \text{hasAuthor}, y)$ . The reason is that query entailment does not consider extra ontological constraints, which would be needed by RDF entailment rules to find such elaborate entailment (here, that conference papers are publications, and that having a contact author is having an author). Extending standard entailment between BGPQs to take into account *extra ontological constraints* modeling background knowledge, is one contribution in this work (Section 4).

## 3. FINDING COMMONALITIES BETWEEN RDF GRAPHS

In Section 3.1, we first define the largest set of commonalities between RDF graphs as a particular RDF graph representing their *least general generalization* (**lgg** for short). Then, we devise a technique for computing such an **lgg** in Section 3.2.

### 3.1 Defining the lgg of RDF graphs

A *least general generalization* of  $n$  descriptions  $d_1, \dots, d_n$  is a most specific description  $d$  generalizing every  $d_{1 \leq i \leq n}$  for some generalization/specialization relation between descriptions [40, 41]. In our RDF setting, we use RDF graphs as descriptions and entailment between RDF graphs as relation for generalization/specialization:

**DEFINITION 1** (**lgg OF RDF GRAPHS**). *Let  $\mathcal{G}_1, \dots, \mathcal{G}_n$  be RDF graphs and  $\mathcal{R}$  a set of RDF entailment rules.*

- A generalization of  $\mathcal{G}_1, \dots, \mathcal{G}_n$  is an RDF graph  $\mathcal{G}_g$  such that  $\mathcal{G}_i \models_{\mathcal{R}} \mathcal{G}_g$  holds for  $1 \leq i \leq n$ .
- A least general generalization (**lgg**) of  $\mathcal{G}_1, \dots, \mathcal{G}_n$  is a generalization  $\mathcal{G}_{\text{lgg}}$  of  $\mathcal{G}_1, \dots, \mathcal{G}_n$  such that for any other generalization  $\mathcal{G}_g$  of  $\mathcal{G}_1, \dots, \mathcal{G}_n$ ,  $\mathcal{G}_{\text{lgg}} \models_{\mathcal{R}} \mathcal{G}_g$  holds.

Importantly, in the RDF setting, the following holds:

**THEOREM 1.** *An lgg of RDF graphs always exists, and is unique up to entailment.*

**PROOF.** An **lgg** of RDF graphs always exists, since we can always construct a (possibly empty) RDF graph that is the **lgg** of RDF graphs, in particular *the cover graph of the RDF graphs* devised in Section 3.2.

Also, an **lgg** of RDF graphs is *unique* up to entailment (since  $\mathcal{G}_{\text{lgg}} \models_{\mathcal{R}} \mathcal{G}_g$  holds for any  $\mathcal{G}_g$  in Definition 1). Indeed, if it were that RDF graphs have *multiple lgg's incomparable* w.r.t. entailment, say  $\text{lgg}_1, \dots, \text{lgg}_m$ , their merge<sup>1</sup>  $\text{lgg}_1 \uplus \dots \uplus \text{lgg}_m$  would be a *single strictly more specific lgg*, a contradiction.  $\square$

<sup>1</sup>The merge of RDF graphs, performed with the RDF spe-

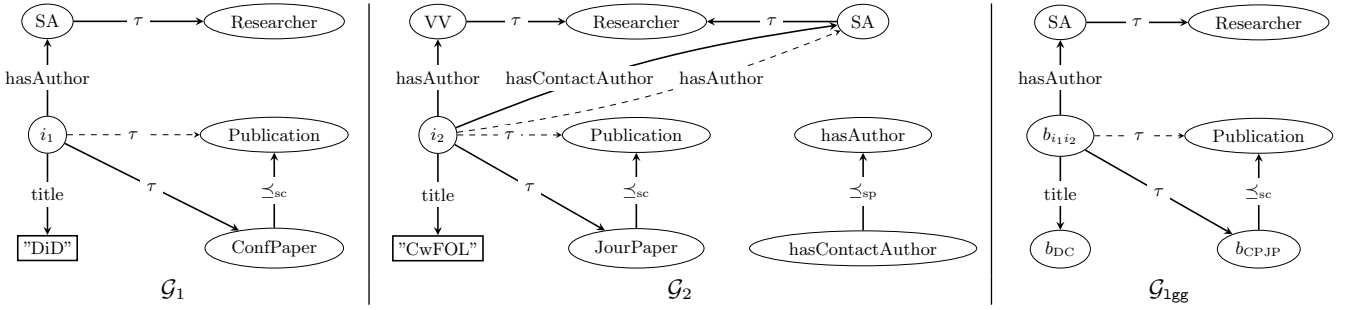


Figure 3: Sample RDF graphs  $\mathcal{G}_1$ ,  $\mathcal{G}_2$  and  $\mathcal{G}_{1gg}$ , with  $\mathcal{G}_{1gg}$  the minimal **lgg** of  $\mathcal{G}_1$  and  $\mathcal{G}_2$ ; their implicit triples (i.e., derived by the rules in Table 2) are shown as dashed edges.

Figure 3 displays two RDF graphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$ , as well as their minimal **lgg** (with lowest number of triples) when we consider the RDF entailment rules shown in Table 2:  $\mathcal{G}_{1gg}$ .  $\mathcal{G}_1$  describes a conference paper  $i_1$  with title “Disaggregations in Databases” and author Serge Abiteboul, who is a researcher; also conference papers are publications.  $\mathcal{G}_2$  describes a journal paper  $i_2$  with title “Computing with First-Order Logic”, contact author Serge Abiteboul and author Victor Vianu, who are researchers; moreover, journal papers are publications and having a contact author is having an author.  $\mathcal{G}_{1gg}$  states that their common information comprises the existence of a resource ( $b_{i_1i_2}$ ) having some type ( $b_{C(ont)P(aper)J(our)P(aper)}$ ), which is a particular case of publication, with some title ( $b_{D(iD)C(wFOL)}$ ) and author Serge Abiteboul, who is a researcher.

Though unique up to entailment (i.e., semantically unique), an **lgg** may have many syntactical forms due to *redundant* triples. Such triples can be either explicit ones that could have been left implicit if the set of RDF entailment rules under consideration allows deriving them from the remaining triples (e.g., materializing the only  $\mathcal{G}_{1gg}$  implicit triple in Figure 3 would make it redundant if we consider the RDF entailment rules in Table 2) or triples generalizing others without needing RDF entailment rules, i.e., in a purely relational fashion (e.g., adding the triple ( $b, \text{hasAuthor}, b'$ ) to  $\mathcal{G}_{1gg}$  in Figure 3 would be redundant w.r.t. ( $b_{i_1i_2}, \text{hasAuthor}, SA$ )). Also, an **lgg** may have *several minimal* syntactical variants obtained by pruning out redundant triples. For example, think of a minimal **lgg** comprising the triples ( $A, \preceq_{sc}, B$ ), ( $B, \preceq_{sc}, A$ ) and ( $b, \tau, A$ ), i.e., there exists an instance of the class  $A$ , which is equivalent to class  $B$ . Clearly, an equivalent and minimal variant of this **lgg** is the RDF graph comprising the triples ( $A, \preceq_{sc}, B$ ), ( $B, \preceq_{sc}, A$ ) and ( $b, \tau, B$ ).

Importantly, the above discussion is not specific to **lggs** of RDF graphs, since any RDF graph may feature redundancy. The detection and elimination of RDF graph redundancy has been studied in the literature, e.g., [34, 38, 39], hence we focus in this work on finding *some* **lgg** of RDF graphs.

The proposition below shows that an **lgg** of  $n$  RDF graphs, with  $n \geq 3$ , can be defined (hence computed) as a sequence of  $n - 1$  **lggs** of *two* RDF graphs. Intuitively, assuming

cific merge operator  $\oplus$ , is an RDF graph comprising the union of the input RDF graph *after renaming* their blank nodes with fresh ones, so that these RDF graphs do not share (thus join on) such values. Indeed, blank nodes are used to characterize the incompleteness of an RDF graph, hence are *local* to it (Section 2).

that  $\ell_{k \geq 2}$  is an operator computing an **lgg** of  $k$  input RDF graphs, the next proposition establishes that:

$$\begin{aligned} \ell_3(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3) &\equiv_{\mathcal{R}} \ell_2(\ell_2(\mathcal{G}_1, \mathcal{G}_2), \mathcal{G}_3) \\ \dots &\dots \\ \ell_n(\mathcal{G}_1, \dots, \mathcal{G}_n) &\equiv_{\mathcal{R}} \ell_2(\ell_{n-1}(\mathcal{G}_1, \dots, \mathcal{G}_{n-1}), \mathcal{G}_n) \\ &\equiv_{\mathcal{R}} \ell_2(\ell_2(\dots \ell_2(\ell_2(\mathcal{G}_1, \mathcal{G}_2), \mathcal{G}_3) \dots, \mathcal{G}_{n-1}), \mathcal{G}_n) \end{aligned}$$

**PROPOSITION 1.** *Let  $\mathcal{G}_1, \dots, \mathcal{G}_{n \geq 3}$  be  $n$  RDF graphs and  $\mathcal{R}$  a set of RDF entailment rules.  $\mathcal{G}_{1gg}$  is an **lgg** of  $\mathcal{G}_1, \dots, \mathcal{G}_n$  iff  $\mathcal{G}_{1gg}$  is an **lgg** of an **lgg** of  $\mathcal{G}_1, \dots, \mathcal{G}_{n-1}$  and  $\mathcal{G}_n$ .*

**PROOF.** The proof relies on the next lemma.

**LEMMA 1.** *Let  $\mathcal{G}_1, \dots, \mathcal{G}_n$  be RDF graphs and  $\mathcal{R}$  a set of RDF entailment rules. If  $\mathcal{G}_{1gg}^1$  is an **lgg** of  $\mathcal{G}_1, \dots, \mathcal{G}_{k < n}$  and  $\mathcal{G}_{1gg}^2$  is an **lgg** of  $\mathcal{G}_1, \dots, \mathcal{G}_n$ , then  $\mathcal{G}_{1gg}^1 \models_{\mathcal{R}} \mathcal{G}_{1gg}^2$  holds.*

Let us show that the above lemma holds.

Suppose that  $\mathcal{G}_{1gg}^1$  is an **lgg** of  $\mathcal{G}_1, \dots, \mathcal{G}_{k < n}$ , i.e., by Definition 1: (i)  $\mathcal{G}_{1 \leq i \leq k} \models_{\mathcal{R}} \mathcal{G}_{1gg}^1$  holds and (ii) for any RDF graph  $\mathcal{G}$  such that  $\mathcal{G}_{1 \leq i \leq k} \models_{\mathcal{R}} \mathcal{G}$  holds,  $\mathcal{G}_{1gg}^1 \models_{\mathcal{R}} \mathcal{G}$  holds.

Suppose also that  $\mathcal{G}_{1gg}^2$  is an **lgg** of  $\mathcal{G}_1, \dots, \mathcal{G}_n$ , i.e., by Definition 1: (i)  $\mathcal{G}_{1 \leq i \leq n} \models_{\mathcal{R}} \mathcal{G}_{1gg}^2$  holds and (ii) for any RDF graph  $\mathcal{G}'$  such that  $\mathcal{G}_{1 \leq i \leq n} \models_{\mathcal{R}} \mathcal{G}'$  holds,  $\mathcal{G}_{1gg}^2 \models_{\mathcal{R}} \mathcal{G}'$  holds.

Clearly,  $\mathcal{G}_{1gg}^2$  is a possible value for  $\mathcal{G}$  above, because  $\mathcal{G}_{1 \leq i \leq n} \models_{\mathcal{R}} \mathcal{G}_{1gg}^2$  implies  $\mathcal{G}_{1 \leq i \leq k < n} \models_{\mathcal{R}} \mathcal{G}_{1gg}^2$ , hence  $\mathcal{G}_{1gg}^1 \models_{\mathcal{R}} \mathcal{G}_{1gg}^2$  must hold.  $\triangleleft$

Now, let us prove Proposition 1 using the above lemma.

( $\Rightarrow$ ) Assume that  $\mathcal{G}_{1gg}$  is an **lgg** of  $\mathcal{G}_1, \dots, \mathcal{G}_n$  and let us show that it is also an **lgg** of some **lgg** of  $\mathcal{G}_1, \dots, \mathcal{G}_{n-1}$  and  $\mathcal{G}_n$ . By Definition 1, because  $\mathcal{G}_{1gg}$  is an **lgg** of  $\mathcal{G}_1, \dots, \mathcal{G}_n$ , we have: (i)  $\mathcal{G}_{1 \leq i \leq n} \models_{\mathcal{R}} \mathcal{G}_{1gg}$  holds and (ii) for any RDF graph  $\mathcal{G}$  such that  $\mathcal{G}_{1 \leq i \leq n} \models_{\mathcal{R}} \mathcal{G}$  holds,  $\mathcal{G}_{1gg} \models_{\mathcal{R}} \mathcal{G}$  holds.

Let  $\mathcal{G}_{1gg}'$  be some **lgg** of  $\mathcal{G}_1, \dots, \mathcal{G}_{n-1}$ . From (i) and the above lemma, (\*)  $\mathcal{G}_{1gg}' \models_{\mathcal{R}} \mathcal{G}_{1gg}$  and  $\mathcal{G}_n \models_{\mathcal{R}} \mathcal{G}_{1gg}$  holds.

Moreover, for any RDF graph  $\mathcal{G}$ , if it were that  $\mathcal{G}_{1gg}' \models_{\mathcal{R}} \mathcal{G}$  and  $\mathcal{G}_n \models_{\mathcal{R}} \mathcal{G}$  and  $\mathcal{G}_{1gg} \not\models_{\mathcal{R}} \mathcal{G}$ , then  $\mathcal{G}_{1 \leq i \leq n} \models_{\mathcal{R}} \mathcal{G}$  would hold and contradict the fact that  $\mathcal{G}_{1gg}$  is an **lgg** of  $\mathcal{G}_1, \dots, \mathcal{G}_n$ . Therefore, (\*\*) for any RDF graph  $\mathcal{G}$ , if  $\mathcal{G}_{1gg}' \models_{\mathcal{R}} \mathcal{G}$  and  $\mathcal{G}_n \models_{\mathcal{R}} \mathcal{G}$  holds, then  $\mathcal{G}_{1gg} \models_{\mathcal{R}} \mathcal{G}$  holds.

From Definition 1, (\*) and (\*\*) we get that  $\mathcal{G}_{1gg}$  is an **lgg** of an **lgg** of  $\mathcal{G}_1, \dots, \mathcal{G}_{n-1}$  and  $\mathcal{G}_n$ .

( $\Leftarrow$ ) Assume that  $\mathcal{G}_{1gg}$  is an **lgg** of an **lgg**  $\mathcal{G}'$  of  $\mathcal{G}_1, \dots, \mathcal{G}_{n-1}$  and  $\mathcal{G}_n$ , and let us show that it is also an **lgg** of  $\mathcal{G}_1, \dots, \mathcal{G}_n$ .

By Definition 1, (i)  $\mathcal{G}' \models_{\mathcal{R}} \mathcal{G}_{1gg}$  and  $\mathcal{G}_n \models_{\mathcal{R}} \mathcal{G}_{1gg}$  hold, hence  $\mathcal{G}_{1 \leq i \leq n-1} \models_{\mathcal{R}} \mathcal{G}_{1gg}$  and  $\mathcal{G}_n \models_{\mathcal{R}} \mathcal{G}_{1gg}$  hold, i.e., (\*)  $\mathcal{G}_{1 \leq i \leq n} \models_{\mathcal{R}} \mathcal{G}_{1gg}$  holds.

Moreover, (ii) for any RDF graph  $\mathcal{G}$  such that  $\mathcal{G}' \models_{\mathcal{R}} \mathcal{G}$  and  $\mathcal{G}_n \models_{\mathcal{R}} \mathcal{G}$  hold,  $\mathcal{G}_{\text{lgg}} \models_{\mathcal{R}} \mathcal{G}$  holds. By Definition 1,  $\mathcal{G}_{1 \leq i \leq n-1} \models_{\mathcal{R}} \mathcal{G}'$  holds, therefore we get: (\*\*') for any RDF graph  $\mathcal{G}$  such that  $\mathcal{G}_{1 \leq i \leq n} \models_{\mathcal{R}} \mathcal{G}$  holds,  $\mathcal{G}_{\text{lgg}} \models_{\mathcal{R}} \mathcal{G}$  holds.

From Definition 1, (\*\*) and (\*\*') we get that  $\mathcal{G}_{\text{lgg}}$  is an **lgg** of  $\mathcal{G}_1, \dots, \mathcal{G}_n$ .  $\square$

Based on the above result, without loss of generality, we focus in the next section on the following problem:

**PROBLEM 1.** *Given two RDF graphs  $\mathcal{G}_1, \mathcal{G}_2$  and a set  $\mathcal{R}$  of RDF entailment rules, we want to compute some **lgg** of  $\mathcal{G}_1$  and  $\mathcal{G}_2$ .*

### 3.2 Computing an lgg of RDF graphs

We first devise the *cover graph* of two RDF graphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$  (to be defined shortly, Definition 2 below), which is central to our technique for computing an **lgg** of  $\mathcal{G}_1$  and  $\mathcal{G}_2$ . We indeed show (Theorem 2) that this particular RDF graph corresponds to an **lgg** of  $\mathcal{G}_1$  and  $\mathcal{G}_2$  when considering their explicit triples *only*, i.e., ignoring RDF entailment rules. Then, we show the main result of this section (Theorem 3): an **lgg** of  $\mathcal{G}_1$  and  $\mathcal{G}_2$ , for *any set*  $\mathcal{R}$  of RDF entailment rules, is the cover graph of their saturations w.r.t.  $\mathcal{R}$ . We also provide the worst-case size of cover graph-based **lggs**, as well as the worst-case time to compute them.

**DEFINITION 2 (COVER GRAPH).** *The cover graph  $\mathcal{G}$  of two RDF graph  $\mathcal{G}_1$  and  $\mathcal{G}_2$  is the RDF graph such that for every property  $p$  in both  $\mathcal{G}_1$  and  $\mathcal{G}_2$ :*

$$(t_1, p, t_2) \in \mathcal{G}_1 \text{ and } (t_3, p, t_4) \in \mathcal{G}_2 \text{ iff } (t_5, p, t_6) \in \mathcal{G}$$

*with  $t_5 = t_1$  if  $t_1 = t_3$  and  $t_1 \in \mathcal{U} \cup \mathcal{L}$ , else  $t_5$  is the blank node  $b_{t_1 t_3}$ , and, similarly  $t_6 = t_2$  if  $t_2 = t_4$  and  $t_2 \in \mathcal{U} \cup \mathcal{L}$ , else  $t_6$  is the blank node  $b_{t_2 t_4}$ .*

The cover graph is more general than  $\mathcal{G}_1$  and  $\mathcal{G}_2$  as each of its triple  $(t_5, p, t_6)$  is a *least general anti-unifier* of a triple  $(t_1, p, t_2)$  from  $\mathcal{G}_1$  and a triple  $(t_3, p, t_4)$  from  $\mathcal{G}_2$ . The notion of *least general anti-unifier* [40, 41, 44] is dual to the well-known notion of *most general unifier* [43, 44]. Observe that  $\mathcal{G}$ 's triples result from anti-unifications of  $\mathcal{G}_1$  and  $\mathcal{G}_2$  triples with *same* property URI. Indeed, anti-unifying triples of the form  $(s_1, p, o_1)$  and  $(s_2, p', o_2)$ , with  $p \neq p'$ , would lead to a non-well-formed triples of the form  $(s, b_{pp'}, o)$  (recall that property values *must* be URIs in RDF graphs), where  $b_{pp'}$  is the blank node required to generalize the distinct values  $p$  and  $p'$ .

Further, the cover graph is an **lgg** for the explicit triples in  $\mathcal{G}_1$  and those in  $\mathcal{G}_2$  since, intuitively, we capture their *common structures* by consistently naming, across all the anti-unifications begetting  $\mathcal{G}$ , the blank nodes used to generalize pairs of distinct subject values or of object values: each time the distinct values  $t$  from  $\mathcal{G}_1$  and  $t'$  from  $\mathcal{G}_2$  are generalized by a blank node while anti-unifying two triples, it is *always* by the same blank node  $b_{t, t'}$  in  $\mathcal{G}$ . This way, we establish *joins* between  $\mathcal{G}$  triples, which reflect the common join structure on  $t$  within  $\mathcal{G}_1$  and on  $t'$  within  $\mathcal{G}_2$ . For instance in Figure 3, the *explicit* triples  $(i_1, \tau, \text{ConfPaper})$ ,  $(\text{ConfPaper}, \preceq_{\text{sc}}, \text{Publication})$ ,  $(i_1, \text{title}, \text{"DiD"})$  in  $\mathcal{G}_1$ , and  $(i_2, \tau, \text{JourPaper})$ ,  $(\text{JourPaper}, \preceq_{\text{sc}}, \text{Publication})$ ,  $(i_2, \text{title}, \text{"CwFOL"})$  in  $\mathcal{G}_2$ , lead to the triples  $(b_{i_1 i_2}, \tau, b_{\text{CPJP}})$ ,  $(b_{\text{CPJP}}, \preceq_{\text{sc}}, \text{Publication})$ ,  $(b_{i_1 i_2}, \text{title}, b_{\text{DC}})$  in the cover graph of  $\mathcal{G}_1$  and  $\mathcal{G}_2$  shown

in Figure 4 (left). The first above-mentioned  $\mathcal{G}$  triple results from anti-unifying  $i_1$  and  $i_2$  into  $b_{i_1 i_2}$ , and,  $\text{ConfPaper}$  and  $\text{JourPaper}$  into  $b_{\text{CPJP}}$ . The second results from anti-unifying *again*  $\text{ConfPaper}$  and  $\text{JourPaper}$  into  $b_{\text{CPJP}}$ , and,  $\text{Publication}$  and  $\text{Publication}$  into  $\text{Publication}$  (as a constant is its own least general generalization). Finally, the third results from anti-unifying *again*  $i_1$  and  $i_2$  into  $b_{i_1 i_2}$ , and, "DiD" and "CwFOL" into  $b_{\text{DC}}$ . By reusing consistently the same blank node name  $b_{i_1 i_2}$  for each anti-unification of the constants  $i_1$  and  $i_2$  (resp.  $b_{\text{CPJP}}$  for  $\text{ConfPaper}$  and  $\text{JourPaper}$ ), the cover graph triples join on  $b_{i_1 i_2}$  (resp.  $b_{\text{CPJP}}$ ) in order to reflect that, in  $\mathcal{G}_1$  and in  $\mathcal{G}_2$ , there exists a particular case of publication ( $i_1$  in  $\mathcal{G}_1$  and  $i_2$  in  $\mathcal{G}_2$ ) with some title ("DiD" in  $\mathcal{G}_1$  and "CwFOL" in  $\mathcal{G}_2$ ).

The next theorem formalizes the above discussion by stating that the cover graph of two RDF graphs is an **lgg** of them, *just in case of an empty set of RDF entailment rules*.

**THEOREM 2.** *The cover graph  $\mathcal{G}$  of the RDF graphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$  is an **lgg** of them for the empty set  $\mathcal{R}$  of RDF entailment rules (i.e.,  $\mathcal{R} = \emptyset$ ).*

**PROOF.** The proof can be directly derived from that of the more general Theorem 3 (see below), noting that in the particular case of Theorem 2  $\mathcal{R} = \emptyset$  holds, hence  $\mathcal{G}_1^\infty = \mathcal{G}_1$  and  $\mathcal{G}_2^\infty = \mathcal{G}_2$  holds.  $\square$

We provide below worst-case bounds for the time to compute a cover graph and for its size; these bounds are met when *all the triples of the two input graphs use the same property URI* (i.e., every pair of  $\mathcal{G}_1$  and  $\mathcal{G}_2$  triples begets a  $\mathcal{G}$  triple).

**PROPOSITION 2.** *The cover graph of two RDF graphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$  can be computed in  $O(|\mathcal{G}_1| \times |\mathcal{G}_2|)$ ; its size is bounded by  $|\mathcal{G}_1| \times |\mathcal{G}_2|$ .*

The main theorem below generalizes Theorem 2 in order to take into account any set of entailment rules from the RDF standard. It states that it is sufficient to compute the cover of the saturations of the input RDF graphs, instead of the input RDF graphs themselves.

**THEOREM 3.** *Let  $\mathcal{G}_1$  and  $\mathcal{G}_2$  be two RDF graphs, and  $\mathcal{R}$  a set of RDF entailment rules. The cover graph  $\mathcal{G}$  of  $\mathcal{G}_1^\infty$  and  $\mathcal{G}_2^\infty$  is an **lgg** of  $\mathcal{G}_1$  and  $\mathcal{G}_2$ .*

**PROOF.** We start by showing that  $\mathcal{G}$  is a generalization of  $\mathcal{G}_1$  and  $\mathcal{G}_2$ , i.e.,  $\mathcal{G}_1 \models_{\mathcal{R}} \mathcal{G}$  and that  $\mathcal{G}_2 \models_{\mathcal{R}} \mathcal{G}$ . Consider the RDF graph  $\mathcal{G}'$  obtained from  $\mathcal{G}$  by replacing every blank node  $b_{v_1 v_2}$  by the value  $v_1$ . Clearly,  $\mathcal{G}' \models_{\mathcal{R}} \mathcal{G}$  and, by construction of  $\mathcal{G}$ ,  $\mathcal{G}' = \mathcal{G}_1$ , i.e.,  $\mathcal{G}_1 \models_{\mathcal{R}} \mathcal{G}$ . Showing  $\mathcal{G}_2 \models_{\mathcal{R}} \mathcal{G}$  is done in a similar way. Hence,  $\mathcal{G}$  is a generalization of  $\mathcal{G}_1$  and  $\mathcal{G}_2$ .

Let us show now that  $\mathcal{G}$  is an **lgg** of  $\mathcal{G}_1$  and  $\mathcal{G}_2$ . Let  $\mathcal{G}_{\text{lgg}}$  be *any* **lgg** of  $\mathcal{G}_1$  and  $\mathcal{G}_2$ . To prove our claim, we need to show that  $\mathcal{G} \models_{\mathcal{R}} \mathcal{G}_{\text{lgg}}$  holds. Since  $\mathcal{G}_{\text{lgg}}$  is an **lgg** of  $\mathcal{G}_1$  and  $\mathcal{G}_2$ , there exist two homomorphisms  $\phi_1$  and  $\phi_2$  from the blank nodes in  $\mathcal{G}_{\text{lgg}}$  to the values in  $\mathcal{G}_1^\infty$  and in  $\mathcal{G}_2^\infty$  respectively, such that  $[\mathcal{G}_{\text{lgg}}]_{\phi_1} \subseteq \mathcal{G}_1^\infty$  and  $[\mathcal{G}_{\text{lgg}}]_{\phi_2} \subseteq \mathcal{G}_2^\infty$ . Consider the graph  $\mathcal{G}'_{\text{lgg}}$  obtained from  $\mathcal{G}_{\text{lgg}}$  by replacing every blank node  $b$  by  $b_{v_1 v_2}$  where  $v_1 = \phi_1(b)$  and  $v_2 = \phi_2(b)$ . Clearly,  $\mathcal{G}_{\text{lgg}} \equiv_{\mathcal{R}} \mathcal{G}'_{\text{lgg}}$  holds. Let us show that  $\mathcal{G} \models_{\mathcal{R}} \mathcal{G}'_{\text{lgg}}$  holds, i.e., there exists a homomorphism  $\phi$  from the blank nodes in  $\mathcal{G}'_{\text{lgg}}$  to

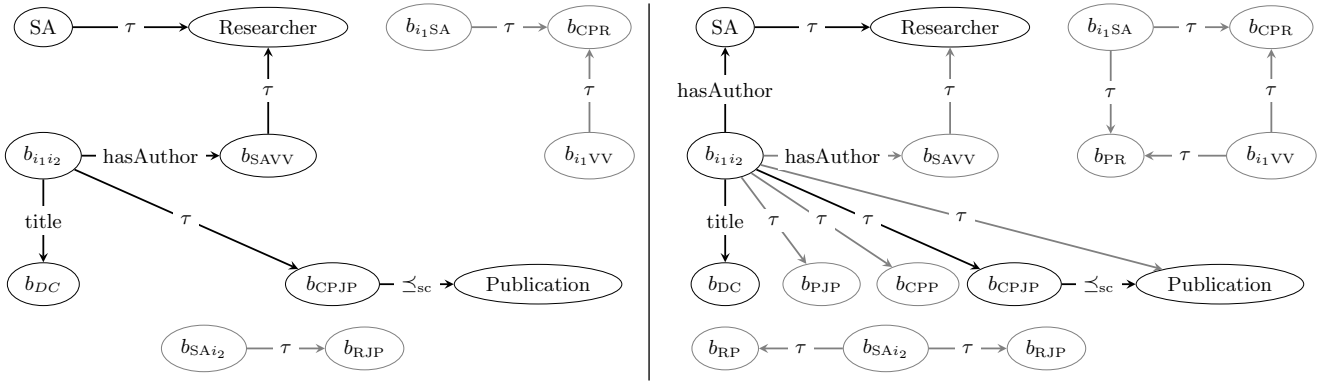


Figure 4: Cover graphs of  $\mathcal{G}_1$  and  $\mathcal{G}_2$  in Figure 3 (left) and of their saturations w.r.t. the entailment rules in Table 2 (right). Triples shown in grey are redundant w.r.t. those shown in black: they are part of the graph, while implicit to the triples shown in black.

the terms in  $\mathcal{G}^\infty$  such that  $[\mathcal{G}'_{1\text{gg}}]_\phi \subseteq \mathcal{G}^\infty$ . By construction of  $\mathcal{G}$ , it is easy to see that the above holds when  $\phi$  maps  $b_{v_1v_2}$  either to  $v_1$  if  $v_1 = v_2$  and  $v_1 \in \mathcal{U} \cup \mathcal{L}$ , or to itself otherwise. It therefore follows that  $\mathcal{G} \models_{\mathcal{R}} \mathcal{G}'_{1\text{gg}}$ , hence  $\mathcal{G} \models_{\mathcal{R}} \mathcal{G}_{1\text{gg}}$ .  $\square$

As an immediate consequence of the above results, we get the following worst-case bounds for the time to compute a cover graph-based **lgg** of two RDF graphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$ , and for its size. Here, we assume given the saturation  $\mathcal{G}_1^\infty$  and  $\mathcal{G}_2^\infty$ , as the times to compute them and their sizes depend on the RDF entailment rules at hand.

**COROLLARY 1.** *An **lgg** of two RDF graphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$  can be computed in  $O(|\mathcal{G}_1^\infty| \times |\mathcal{G}_2^\infty|)$  and its size is bounded by  $|\mathcal{G}_1^\infty| \times |\mathcal{G}_2^\infty|$ .*

Remark that computing naively the cover graph-based **lgg** of  $n$  RDF graphs of size  $M$  based on Proposition 1 may lead to an **lgg** of size  $M^n$ , in the unlikely worst-case where all the triples of all the RDF graphs use the same property URI. However, removing the redundant triples from the intermediate and final cover graph-based **lgg**s limits their size to at most  $M$ .

Figure 4 (right) displays an **lgg** of the RDF graphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$  in Figure 3 w.r.t. the entailment rules shown in Table 2. In contrast to Figure 4 (left), which shows an **lgg** of the same RDF graphs when RDF entailment rules are ignored, we further learn that Serge Abiteboul is an author of some particular publication ( $i_1$  in  $\mathcal{G}_1$  and  $i_2$  in  $\mathcal{G}_2$ ). Moreover, removing the redundant triples (the grey ones) yields precisely the **lgg**  $\mathcal{G}_{1\text{gg}}$  of  $\mathcal{G}_1$  and  $\mathcal{G}_2$  shown in Figure 3.

## 4. FINDING COMMONALITIES BETWEEN QUERIES

We consider now the problem of finding the largest set of commonalities between queries. Adapting the results for RDF graphs from the preceding section to BGPQs is rather direct (recall that a query body is a BGP, i.e., an RDF graph in which variables can be used in subject, property and object positions of triples), but of limited interest as the next example shows. Consider the two BGPQs  $q_1$  and  $q_2$  in Figure 6. The first asks for the conference papers having some contact author, while the second asks for the journal papers having some author. Clearly, a least general generalization

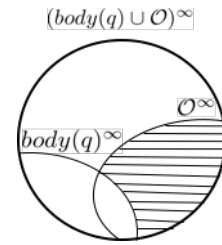


Figure 5: Characterization of the body of a saturated BGPQ  $q$  w.r.t. a set  $O$  of RDFS constraints

of  $q_1$  and  $q_2$  is the *very* general BGPQ  $q_{1\text{gg}}$  shown in Figure 6 that asks for the resources having some type. Observe that  $q_{1\text{gg}}$  is an **lgg** for *any* subset of the RDF entailment rules shown in Table 2, as there is *no* ontological constraint in the queries. However, by considering the extra ontological constraints displayed in Figure 7 that hold in the scientific publication domain, i.e., the context in which the queries are asked, a more pregnant **lgg** would be a BGPQ asking for *publications having some researcher as author*, since (i) having a contact authors is having an author, (ii) only publications have authors, (iii) only researchers are authors, and (iv) conference (resp. journal) papers are publications.

We therefore devise in Section 4.1 a notion of **lgg** of BGPQs in the presence of ontological constraints, which as exemplified above helps finding more interesting commonalities by taking into account the background knowledge of an application domain. Then, we provide a technique for computing such an **lgg** of BGPQs in Section 4.2.

### 4.1 Defining the **lgg** of BGPQs

Recall that the notion of **lgg** of  $n$  descriptions  $d_1, \dots, d_n$  is a most specific description  $d$  generalizing  $d_1, \dots, d_n$  for some generalization/specialization relation. Here, we adapt this notion by using BGPQs as descriptions and revisiting entailment between BGPQs (Section 2) in order to endow this RDF relation between queries with background knowledge.

We first devise in Section 4.1.1 *entailment between BGPQs w.r.t. ontological constraints* as a generalization of the standard entailment between BGPQs. Then, we define the **lgg** of BGPQs w.r.t. *ontological constraints* in Section 4.1.2.



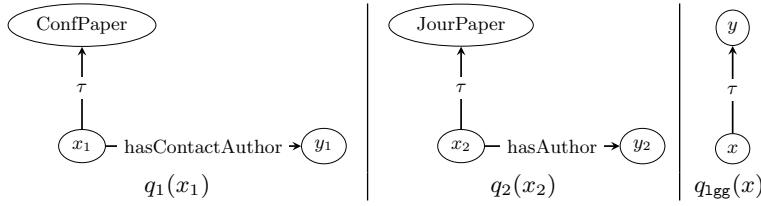


Figure 6: Sample BGPQs  $q_1, q_2$  and their minimal  $\text{lgg } q_{\text{lgg}}$ .

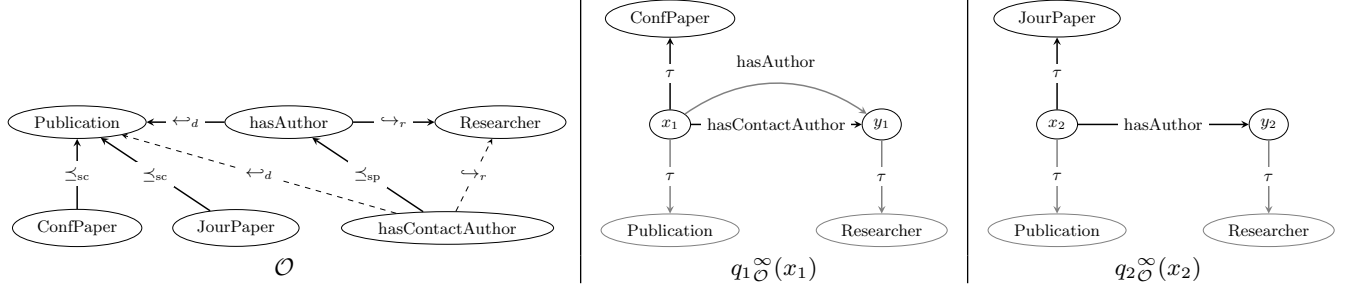


Figure 7: Sample set of ontological constraints  $\mathcal{O}$ ; saturations of the BGPQs  $q_1$  and  $q_2$  in Figure 6 w.r.t.  $\mathcal{O}$ . Triples shown in grey are redundant w.r.t. those shown in black.

#### 4.1.1 Entailment between BGPQs w.r.t. constraints

We start by generalizing the saturation of a query w.r.t. extra ontological constraints. The saturated query 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.

**DEFINITION 3 (SATURATION W.R.T. 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.

The above definition precisely characterizes  $\text{body}(q_{\mathcal{O}}^{\infty})$  as the non-hatched subset of  $(\mathcal{O} \cup \text{body}(q))^{\infty}$  shown in Figure 5. Also, remark that this definition coincides with the standard one of saturation (Section 2) when the set  $\mathcal{O}$  of ontological constraints is empty.

Figure 7 shows a set  $\mathcal{O}$  of ontological constraints, as well as the saturations of the BGPQs  $q_1$  and  $q_2$  in Figure 6 w.r.t.  $\mathcal{O}$ , named  $q_{1\mathcal{O}}^{\infty}$  and  $q_{2\mathcal{O}}^{\infty}$  respectively.

The theorem below states the fundamental properties for a query and its saturation w.r.t. ontological constraints: they are *equivalent* for the central RDF reasoning tasks of query entailment and query answering.

**THEOREM 4.** Let  $\mathcal{R}$  be a set of RDF entailment rules,  $\mathcal{O}$  a set of RDFS statements, and  $q$  a BGPQ whose saturation

<sup>2</sup>Blank nodes in  $q$  can be replaced by fresh non-answer variables, while preserving equivalence (Section 2). This assumption follows again from the semantics of blank nodes, which are used to model unknown information *within* an RDF graph or *within* a BGPQ.

w.r.t.  $\mathcal{O}$  is  $q_{\mathcal{O}}^{\infty}$ . For any RDF graph  $\mathcal{G}$  whose set of RDFS statements is  $\mathcal{O}$ ,

1.  $\mathcal{G} \models_{\mathcal{R}} q$  holds iff  $\mathcal{G} \models_{\mathcal{R}} q_{\mathcal{O}}^{\infty}$  holds,
2.  $q(\mathcal{G}^{\infty}) = q_{\mathcal{O}}^{\infty}(\mathcal{G}^{\infty})$  holds.

**PROOF.** Let us first show item 1). Observe that  $\mathcal{G} \models_{\mathcal{R}} q$  iff  $\mathcal{G} \models_{\mathcal{R}} q_{\mathcal{O}}^{\infty}$  holds iff  $\mathcal{G}^{\infty} \models q$  iff  $\mathcal{G}^{\infty} \models q_{\mathcal{O}}^{\infty}$  holds (recall Section 2). If  $\mathcal{G}^{\infty} \models^{\phi} q_{\mathcal{O}}^{\infty}$  holds for some homomorphism  $\phi$ , then  $\mathcal{G}^{\infty} \models^{\phi} q$  holds also since  $\text{body}(q) \subseteq \text{body}(q_{\mathcal{O}}^{\infty})$ . Now, if  $\mathcal{G}^{\infty} \models^{\phi} q$  holds for some homomorphism  $\phi$ , consider the subset  $[\text{body}(q)]_{\phi}$  of  $\mathcal{G}^{\infty}$ . Since  $\mathcal{O} \subseteq \mathcal{G}$ , and by definition of the saturation of  $q$  w.r.t.  $\mathcal{O}$ ,  $\phi$  can obviously be extended to a homomorphism  $\phi'$  such that  $[\text{body}(q_{\mathcal{O}}^{\infty})]_{\phi'} \subseteq \mathcal{G}^{\infty}$ , which maps  $\text{body}(q_{\mathcal{O}}^{\infty})$  to the maximal subset of  $(\mathcal{O} \cup [\text{body}(q)]_{\phi})^{\infty} \subseteq \mathcal{G}^{\infty}$ , such that for any of its subset  $\mathcal{G}'$ : if  $\mathcal{O} \models_{\mathcal{R}} \mathcal{G}'$  then  $[\text{body}(q)]_{\phi} \models_{\mathcal{R}} \mathcal{G}'$ .

Item 2) directly follows from: (i) the fact that, above, item 1) holds for any homomorphism  $\phi$  and (ii)  $q$  and  $q_{\mathcal{O}}^{\infty}$  have, by definition, the same answer variables.  $\square$

Based on the above result, we are now ready to extend the entailment relation between queries, i.e., the standard RDF generalization/specialization relation between queries, to that of *entailment between queries w.r.t. ontological constraints*.

**DEFINITION 4 (ENTAILMENT W.R.T. 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  $q_{\mathcal{O}}^{\infty} \models q'$  holds.

The above definition coincides with the standard one of entailment between BGPQs (Section 2) when the set  $\mathcal{O}$  of ontological constraints is empty.

The BGPQ  $q(x) \leftarrow (x, \tau, \text{Publication}), (x, \text{hasAuthor}, y)$  is *not* entailed by the queries  $q_1$  and  $q_2$  shown in Figure 6, while it is entailed by these queries w.r.t. the set  $\mathcal{O}$  of ontological constraints displayed in Figure 7:  $q_{1\mathcal{O}}^{\infty} \models^{p_1} q$  (resp.

$q_2^\infty \models^{\phi_2} q$  holds for  $\phi_1(x) = x_1$  and  $\phi_1(y) = y_1$  (resp.  $\phi_2(x) = x_2$  and  $\phi_2(y) = y_2$ ).

The next theorem establishes the fundamental properties for a query 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.

**THEOREM 5.** *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  $\mathcal{O}$ ,*

1. *if  $\mathcal{G} \models_{\mathcal{R}} q$  holds then  $\mathcal{G} \models_{\mathcal{R}} q'$  holds,*
2.  *$q(\mathcal{G}^\infty) \subseteq q'(\mathcal{G}^\infty)$  holds.*

**PROOF.** Let us first show item 1). By Theorem 4, item 1),  $\mathcal{G} \models_{\mathcal{R}} q$  holds iff  $\mathcal{G} \models_{\mathcal{R}} q^\infty$  holds. Consider some  $\phi$  such that  $\mathcal{G} \models_{\mathcal{R}}^\phi q^\infty$  holds. Also, by Definition 4, consider some  $\phi'$  such that  $q^\infty \models^{\phi'} q'$ . By composing  $\phi$  and  $\phi'$  into  $\psi$ , we get  $\mathcal{G} \models_{\mathcal{R}}^\psi q'$ , i.e.,  $\mathcal{G} \models_{\mathcal{R}} q'$ .

Item 2) directly follows from: (i) the fact that, above, item 1) holds for any pair of homomorphisms  $\phi, \phi'$  such that  $\mathcal{G} \models_{\mathcal{R}}^\phi q^\infty$  and  $q^\infty \models^{\phi'} q'$  hold, and (ii)  $q$  and  $q^\infty$  have, by definition, the same answer variables.  $\square$

Finally, it is worth noting that our notion of entailment between BGPQs w.r.t. ontological constraints is the RDF counterpart to that of *query containment under deductive constraints* in a database setting [18].

#### 4.1.2 $\mathbf{l_{gg}}$ of BGPQs w.r.t. constraints

With the above new RDF generalization/specialization relation between queries endowed with background knowledge, we are now ready to define the  $\mathbf{l_{gg}}$  of queries w.r.t. ontological constraints.

**DEFINITION 5** ( $\mathbf{l_{gg}}$  OF BGPQs). *Let  $q_1, \dots, q_n$  be BGPQs with the same arity,  $\mathcal{R}$  a set of RDF entailment rules and  $\mathcal{O}$  a set of RDFS statements.*

- *A generalization of  $q_1, \dots, q_n$  w.r.t.  $\mathcal{O}$  is a BGPQ  $q_g$  such that  $q_i \models_{\mathcal{R}, \mathcal{O}} q_g$  holds for  $1 \leq i \leq n$ .*
- *A least general generalization of  $q_1, \dots, q_n$  w.r.t.  $\mathcal{O}$  is a generalization  $q_{\mathbf{l_{gg}}}$  of  $q_1, \dots, q_n$  w.r.t.  $\mathcal{O}$  such that for any other generalization  $q_g$  of  $q_1, \dots, q_n$  w.r.t.  $\mathcal{O}$ ,  $q_{\mathbf{l_{gg}}} \models_{\mathcal{R}, \mathcal{O}} q_g$  holds.*

When  $\mathcal{O}$  is empty, the above definition is the direct adaptation of that of  $\mathbf{l_{gg}}$  of RDF graphs to the case of BGPQs, using standard entailment between BGPQs.

Importantly, in our SPARQL setting, the following holds:

**THEOREM 6.** *An  $\mathbf{l_{gg}}$  of BGPQs w.r.t. RDFS statements may not exist; when it exists, it is unique up to entailment w.r.t. the considered RDFS statements.*

**PROOF.** An  $\mathbf{l_{gg}}$  of BGPQs may not exist. For instance, consider  $q_1(x_1) \leftarrow (x_1, \text{hasAuthor}, y_1)$  asking for the resources having some author, and  $q_2(x_2) \leftarrow (y_2, \text{hasAuthor}, x_2)$  asking for the authors of some resource. Clearly, there is no BGPQ that generalizes them both.

Also, an  $\mathbf{l_{gg}}$  of queries is unique up to entailment w.r.t. ontological constraints, i.e., is semantically unique (as  $q_{\mathbf{l_{gg}}} \models_{\mathcal{R}, \mathcal{O}}$

$q_g$  holds for any  $q_g$ ). Indeed, if it were that queries have *multiple  $\mathbf{l_{gg}}$ s* incomparable w.r.t. entailment, say the BGPQs  $\mathbf{l_{gg}}_1(\bar{x}), \dots, \mathbf{l_{gg}}_m(\bar{x})$ , the BGPQ that merges<sup>3</sup> their bodies  $q_{\mathbf{l_{gg}}}(\bar{x}) \leftarrow \text{body}(\mathbf{l_{gg}}_1) \uplus \dots \uplus \text{body}(\mathbf{l_{gg}}_m)$  would be a *single strictly more specific  $\mathbf{l_{gg}}$* , a contradiction.  $\square$

Also, queries may have many syntactically different (though equivalent)  $\mathbf{l_{gg}}$ s due to redundant triples. However, differently from the case of RDF graphs, redundant triples are only those generalizing others in a *purely relation fashion*, as the semantics of BGPQs is *not* its saturation (recall Section 2). This implies that a query has a unique *minimal  $\mathbf{l_{gg}}$* , which may be obtained using standard database minimization techniques for relational conjunctive queries. Importantly, redundancy of triples is not specific to  $\mathbf{l_{gg}}$ s of BGPQs, since any query may feature redundancy. We therefore focus in this work on finding *some  $\mathbf{l_{gg}}$*  of BGPQs.

The proposition below is the counterpart of Proposition 1 for RDF graphs. It shows that an  $\mathbf{l_{gg}}$  of  $n$  BGPQs, with  $n \geq 3$ , can be defined (hence computed) as a sequence of  $n - 1$   $\mathbf{l_{gg}}$ s of *two* BGPQs. That is, assuming that  $\ell_{k \geq 2}$  is an operator computing an  $\mathbf{l_{gg}}$  of  $k$  input BGPQs, the next proposition establishes that:

$$\begin{aligned} \ell_3(q_1, q_2, q_3) &\equiv_{\mathcal{R}, \mathcal{O}} \ell_2(\ell_2(q_1, q_2), q_3) \\ \dots &\dots \\ \ell_n(q_1, \dots, q_n) &\equiv_{\mathcal{R}, \mathcal{O}} \ell_2(\ell_{n-1}(q_1, \dots, q_{n-1}), q_n) \\ &\equiv_{\mathcal{R}, \mathcal{O}} \ell_2(\ell_2(\dots \ell_2(\ell_2(q_1, q_2), q_3) \dots, q_{n-1}), q_n) \end{aligned}$$

**PROPOSITION 3.** *Let  $q_1, \dots, q_{n \geq 3}$  be  $n$  BGPQs,  $\mathcal{O}$  a set of RDFS statements and  $\mathcal{R}$  a set of RDF entailment rules.  $q_{\mathbf{l_{gg}}}$  is an  $\mathbf{l_{gg}}$  of  $q_1, \dots, q_n$  w.r.t.  $\mathcal{O}$  iff  $q_{\mathbf{l_{gg}}}$  is an  $\mathbf{l_{gg}}$  w.r.t.  $\mathcal{O}$  of an  $\mathbf{l_{gg}}$  of  $q_1, \dots, q_{n-1}$  and  $q_n$  w.r.t.  $\mathcal{O}$ .*

**PROOF.** The proof relies on the next lemma.

**LEMMA 2.** *Let  $q_1, \dots, q_n$  be BGPQs,  $\mathcal{O}$  a set of RDFS statements and  $\mathcal{R}$  a set of RDF entailment rules. If  $q_{\mathbf{l_{gg}}}^1$  is an  $\mathbf{l_{gg}}$  of  $q_1, \dots, q_{k < n}$  w.r.t.  $\mathcal{O}$  and  $q_{\mathbf{l_{gg}}}^2$  is an  $\mathbf{l_{gg}}$  of  $q_1, \dots, q_n$  w.r.t.  $\mathcal{O}$ , then  $q_{\mathbf{l_{gg}}}^1 \models_{\mathcal{R}, \mathcal{O}} q_{\mathbf{l_{gg}}}^2$  holds.*

Let us show that the above lemma holds.

Suppose that  $q_{\mathbf{l_{gg}}}^1$  is an  $\mathbf{l_{gg}}$  of  $q_1, \dots, q_{k < n}$  w.r.t.  $\mathcal{O}$ , i.e., by Definition 5: (i)  $q_{1 \leq i \leq k} \models_{\mathcal{R}, \mathcal{O}} q_{\mathbf{l_{gg}}}^1$  holds and (ii) for any BGPQ  $q$  such that  $q_{1 \leq i \leq k} \models_{\mathcal{R}, \mathcal{O}} q$  holds,  $q_{\mathbf{l_{gg}}}^1 \models_{\mathcal{R}, \mathcal{O}} q$  holds.

Suppose also that  $q_{\mathbf{l_{gg}}}^2$  is an  $\mathbf{l_{gg}}$  of  $q_1, \dots, q_n$  w.r.t.  $\mathcal{O}$ , i.e., by Definition 5: (i)  $q_{1 \leq i \leq n} \models_{\mathcal{R}, \mathcal{O}} q_{\mathbf{l_{gg}}}^2$  holds and (ii) for any BGPQ  $q'$  such that  $q_{1 \leq i \leq n} \models_{\mathcal{R}, \mathcal{O}} q'$  holds,  $q_{\mathbf{l_{gg}}}^2 \models_{\mathcal{R}, \mathcal{O}} q'$  holds.

Clearly,  $q_{\mathbf{l_{gg}}}^2$  is a possible value for  $q$  above, because  $q_{1 \leq i \leq n} \models_{\mathcal{R}, \mathcal{O}} q_{\mathbf{l_{gg}}}^2$  implies  $q_{1 \leq i \leq k < n} \models_{\mathcal{R}, \mathcal{O}} q_{\mathbf{l_{gg}}}^2$ , hence  $q_{\mathbf{l_{gg}}}^1 \models_{\mathcal{R}, \mathcal{O}} q_{\mathbf{l_{gg}}}^2$  must hold.  $\triangleleft$

Now, let us prove Proposition 3 using the above lemma.

( $\Rightarrow$ ) Assume that  $q_{\mathbf{l_{gg}}}$  is an  $\mathbf{l_{gg}}$  of  $q_1, \dots, q_n$  w.r.t.  $\mathcal{O}$  and let us show that it is also an  $\mathbf{l_{gg}}$  w.r.t.  $\mathcal{O}$  of some  $\mathbf{l_{gg}}$  of  $q_1, \dots, q_{n-1}$  and  $q_n$  w.r.t.  $\mathcal{O}$ . By Definition 5, because  $q_{\mathbf{l_{gg}}}$  is an  $\mathbf{l_{gg}}$  of  $q_1, \dots, q_n$  w.r.t.  $\mathcal{O}$ , we have: (i)  $q_{1 \leq i \leq n} \models_{\mathcal{R}, \mathcal{O}} q_{\mathbf{l_{gg}}}$

<sup>3</sup>The *merge* operator of RDF graphs  $\uplus$  is straightforwardly extended to union BGPQ bodies (i.e., BGPs) after renaming their blank nodes and non-answer variables with fresh ones, so that these BGPs do not share (thus join on) them. Blank nodes and non-answer variables in a BGPQ, which have the *same* semantics (recall Section 2), are used to characterize the answer variables of this BGPQ, hence are *local* to it.

holds and (ii) for any BGPQ  $q$  such that  $q_{1 \leq i \leq n} \models_{\mathcal{R}, \mathcal{O}} q$  holds,  $q_{1\text{gg}} \models_{\mathcal{R}, \mathcal{O}} q$  holds.

Let  $q'_{1\text{gg}}$  be an **lgg** of  $q_1, \dots, q_{n-1}$  w.r.t.  $\mathcal{O}$ . From (i) and the above lemma, (\*)  $q'_{1\text{gg}} \models_{\mathcal{R}, \mathcal{O}} q_{1\text{gg}}$  and  $q_n \models_{\mathcal{R}, \mathcal{O}} q_{1\text{gg}}$  hold. Moreover, for any BGPQ  $q$ , if it were that  $q'_{1\text{gg}} \models_{\mathcal{R}, \mathcal{O}} q$  and  $q_n \models_{\mathcal{R}, \mathcal{O}} q$  and  $q_{1\text{gg}} \not\models_{\mathcal{R}, \mathcal{O}} q$ , then  $q_{1 \leq i \leq n} \models_{\mathcal{R}, \mathcal{O}} q$  would hold and contradict the fact that  $q_{1\text{gg}}$  is an **lgg** of  $q_1, \dots, q_n$  w.r.t.  $\mathcal{O}$ . Therefore, (\*\*) for any BGPQ  $q$ , if  $q'_{1\text{gg}} \models_{\mathcal{R}, \mathcal{O}} q$  and  $q_n \models_{\mathcal{R}, \mathcal{O}} q$  holds, then  $q_{1\text{gg}} \models_{\mathcal{R}, \mathcal{O}} q$  holds.

From Definition 5, (\*) and (\*\*) we get that  $q_{1\text{gg}}$  is an **lgg** w.r.t.  $\mathcal{O}$  of an **lgg** of  $q_1, \dots, q_{n-1}$  and  $q_n$  w.r.t.  $\mathcal{O}$ .

( $\Leftarrow$ ) Assume that  $q_{1\text{gg}}$  is an **lgg** w.r.t.  $\mathcal{O}$  of an **lgg**  $q'$  of  $q_1, \dots, q_{n-1}$  and  $q_n$  w.r.t.  $\mathcal{O}$ , and let us show that it is also an **lgg** of  $q_1, \dots, q_n$  w.r.t.  $\mathcal{O}$ .

By Definition 5, (i)  $q' \models_{\mathcal{R}, \mathcal{O}} q_{1\text{gg}}$  and  $q_n \models_{\mathcal{R}, \mathcal{O}} q_{1\text{gg}}$  hold, hence  $q_{1 \leq i \leq n-1} \models_{\mathcal{R}, \mathcal{O}} q_{1\text{gg}}$  and  $q_n \models_{\mathcal{R}, \mathcal{O}} q_{1\text{gg}}$  hold, i.e., (\*)  $q_{1 \leq i \leq n} \models_{\mathcal{R}, \mathcal{O}} q_{1\text{gg}}$  holds.

Moreover, (ii) for any BGPQ  $q$  such that  $q' \models_{\mathcal{R}, \mathcal{O}} q$  and  $q_n \models_{\mathcal{R}, \mathcal{O}} q$  hold,  $q_{1\text{gg}} \models_{\mathcal{R}, \mathcal{O}} q$  holds. By Definition 5,  $q_{1 \leq i \leq n-1} \models_{\mathcal{R}, \mathcal{O}} q'$  holds, therefore we get: (\*\*') for any BGPQ  $q$  such that  $q_{1 \leq i \leq n} \models_{\mathcal{R}, \mathcal{O}} q$  holds,  $q_{1\text{gg}} \models_{\mathcal{R}, \mathcal{O}} q$  holds.

From Definition 5, (\*) and (\*\*') we get that  $q_{1\text{gg}}$  is an **lgg** of  $q_1, \dots, q_n$  w.r.t.  $\mathcal{O}$ .  $\square$

Based on the above result, without loss of generality, we focus in the next section on the following problem:

**PROBLEM 2.** *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, we want to compute some **lgg** of  $q_1$  and  $q_2$  w.r.t.  $\mathcal{O}$ .*

## 4.2 Computing an lgg of BGPQs

We first devise the notion of *cover query* of two BGPQs  $q_1$  and  $q_2$  (to be defined shortly, Definition 6 below) and we show (Theorem 7) that it corresponds to an **lgg** of  $q_1$  and  $q_2$  just in case extra ontological constraints and RDF entailment are ignored. Then, we show (Theorem 8) that their enhanced **lgg** as defined in Definition 5, i.e., which does consider extra ontological constraints and RDF entailment, can be obtained as the cover query of the saturations of  $q_1$  and of  $q_2$  with the ontological constraints and RDF entailment rules at hand (Definition 3). We also provide the size of these cover query-based **lggs**, as well as the time to compute them.

**DEFINITION 6 (COVER QUERY).** *Let  $q_1, q_2$  be two BGPQs with the same arity  $n$ .*

*If there exists the BGPQ  $q$  such that*

- *$(t_1, t_2, t_3) \in \text{body}(q_1)$  and  $(t_4, t_5, t_6) \in \text{body}(q_2)$  iff  $(t_7, t_8, t_9) \in \text{body}(q)$  with, for  $1 \leq i \leq 3$ ,  $t_{i+6} = t_i$  if  $t_i = t_{i+3}$  and  $t_i \in \mathcal{U} \cup \mathcal{L}$ , otherwise  $t_{i+6}$  is the variable  $v_{t_i t_{i+3}}$*
- *$\text{head}(q) = q(v_{x_1 x_2}, \dots, v_{x_1^n x_2^n})$  iff  $\text{head}(q_1) = q(x_1^1, \dots, x_1^n)$  and  $\text{head}(q_2) = q(x_2^1, \dots, x_2^n)$*

*then  $q$  is the cover query of  $q_1$  and  $q_2$ .*

The cover query body is defined (first item above) as a *slight generalization of the cover graph of two RDF graphs* (Definition 1 and Theorem 2), so that it corresponds to an **lgg** of the bodies of  $q_1$  and  $q_2$  when both extra ontological constraints and RDF entailment are ignored. Recall that the bodies of the BGPQs are BGPs, which generalize

RDF graphs by allowing variables in triples (Section 2). In particular, unknown values are allowed in the property position of BGP triples (using variables), whereas this was not possible in RDF graph triples in which the property values must be URIs. To generalize the results of Theorem 2 from RDF graphs to BGPs, (i) *variables* are used instead of blank nodes to generalize distinct values in the least general anti-unifications of triples begetting the cover query body (as they have the same semantics within a BGP, but only variables can be used at subject, object and property positions in triples), (ii) the cover query body comprises the least general anti-unifications of *every*  $q_1$  triple with *every*  $q_2$  triple instead of just those pairs of triples with same property URI (as now distinct property values can be generalized using variables), and (iii) again, across these anti-unifications characterizing the body of  $q$ , we consistently name the variables generalizing same pair of the distinct  $q_1$  and  $q_2$  values (so as to capture the common structure between the bodies of  $q_1$  and  $q_2$ ).

The cover query head is simply defined (second item in Definition 6) as the least general anti-unification of the  $q_1$ 's head and  $q_2$ 's head, in order to model the least general set of result tuples generalizing those of  $q_1$  and of  $q_2$ .

Figure 8 (left) shows the cover query  $q^c$  of the BGPQs  $q_1$  and  $q_2$  displayed in Figure 6. The head  $q^c(v_{x_1 x_2})$  of  $q^c$  results from anti-unifying  $q_1(x_2)$  and  $q_2(x_2)$ ; its body consists of the triple  $(v_{x_1 x_2}, \tau, v_{\text{CPJP}})$ , which results from anti-unifying  $(x_1, \tau, \text{ConfPaper})$  in  $q_1$  and  $(x_2, \tau, \text{JourPaper})$  in  $q_2$ , of the triple  $(v_{x_1 x_2}, v_{\text{hCA}\tau}, v_{y_1 \text{JP}})$ , which results from anti-unifying  $(x_1, \text{hasContactAuthor}, y_1)$  in  $q_1$  and  $(x_2, \tau, \text{JourPaper})$  in  $q_2$  etc.

Importantly, a cover query of two BGPQs may not exist (If ... then ... in Definition 6). For instance, recall the BGPQs  $q_1$  and  $q_2$  previously introduced in subsection 4.1.2, with which we pointed out that an **lgg** may not exist. Their cover query does not exist either, since Definition 6 leads to  $q(v_{x_1 x_2}) \leftarrow (v_{x_1 y_2}, \text{hasAuthor}, v_{y_1 x_2})$  which is *not* a BGPQ ( $v_{x_1 x_2}$  does not appear in  $q$ 's body).

The next theorem formalizes the above discussion:

**THEOREM 7.** *Given two BGPQs  $q_1, q_2$  with the same arity and empty sets  $\mathcal{R}$  of RDF entailment rules and  $\mathcal{O}$  of RDFS statements:*

1. *the cover query of  $q_1$  and  $q_2$  exists iff an **lgg** of  $q_1$  and  $q_2$  exists;*
2. *the cover query of  $q_1$  and  $q_2$  is an **lgg** of  $q_1$  and  $q_2$ .*

**PROOF.** The proof can be directly derived from that of the more general Theorem 8 (see below), noting that in the particular case of Theorem 7  $\mathcal{R} = \emptyset$  and  $\mathcal{O} = \emptyset$  hold, hence  $q_1^\infty = q_1$  and  $q_2^\infty = q_2$  hold.  $\square$

In Figure 8 (left), the cover query  $q^c$  of the BGPQs  $q_1$  and  $q_2$  displayed in Figure 6 is therefore an **lgg** of them when both extra ontological constraints and RDF entailment are ignored (Theorem 7). Remark that  $q^c$  is equivalent to the naïve minimal **lgg** of  $q_1$  and  $q_2$  shown Figure 6:  $q_{1\text{gg}}$ .

We provide below the worst-case time to compute a cover query, and its size.

**PROPOSITION 4.** *The cover query of two BGPQs  $q_1$  and  $q_2$  can be computed in  $O(|\text{body}(q_1)| \times |\text{body}(q_2)|)$ ; its size is  $|\text{body}(q_1)| \times |\text{body}(q_2)|$ .*

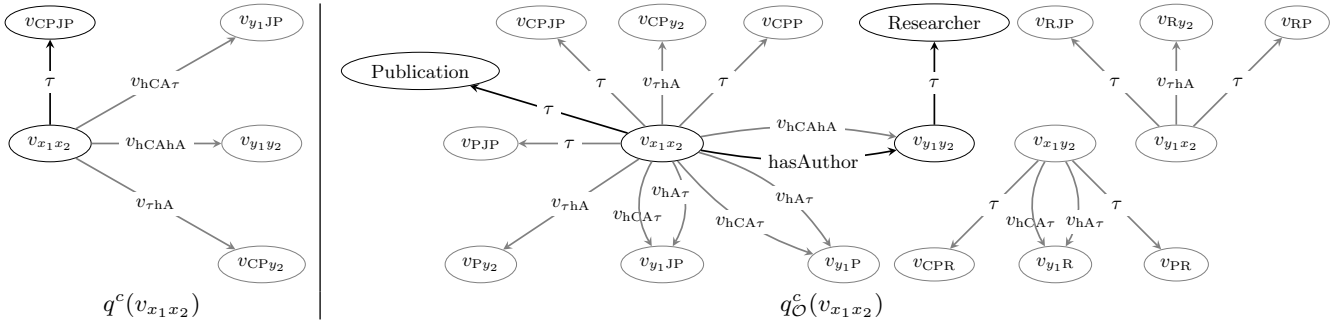


Figure 8: Cover queries of the BGPQs  $q_1$  and  $q_2$  in Figure 6 (left), and of their saturations w.r.t.  $\mathcal{O}$  in Figure 7 (right). Triples shown in grey are redundant w.r.t. those shown in black.

**THEOREM 8.** *Given a set  $\mathcal{R}$  of RDF entailment rules, a set  $\mathcal{O}$  of RDFS statements and two BGPQs  $q_1, q_2$  with the same arity,*

PROOF. We start by showing that the cover query  $q$ , when it exists, is a generalization w.r.t.  $\mathcal{O}$  of  $q_1$  and  $q_2$ , i.e.,  $q_1 \models_{\mathcal{R}, \mathcal{O}} q$  and that  $q_2 \models_{\mathcal{R}, \mathcal{O}} q$ . Clearly, once proved, this entails the fact that if  $q$  exists, then an **lgg** of  $q_1$  and  $q_2$  w.r.t.  $\mathcal{O}$  exists.

Therefore, the cover query  $q$  is a generalization of  $q_1$  and  $q_2$  w.r.t.  $\mathcal{O}$ , thus an  $\text{lgg}$  of  $q_1$  and  $q_2$  w.r.t.  $\mathcal{O}$  also exists.

To prove our claim, we need to show that  $q \models_{\mathcal{R}, \mathcal{O}} q_{1\text{gg}}$  holds. Since  $q_{1\text{gg}}$  is an **lgg** of  $q_1$  and  $q_2$  w.r.t.  $\mathcal{O}$ , there exist two homomorphisms  $\phi_1$  and  $\phi_2$  from the variables in  $q_{1\text{gg}}$  to the values (variables and constants) in  $q_1^\mathcal{O}$  and in  $q_2^\mathcal{O}$  respectively, such that  $[body(q_{1\text{gg}})]_{\phi_1} \subseteq body(q_1^\mathcal{O})$  and  $\phi_1(\bar{x}_{1\text{gg}}) = \bar{x}_1$ , and similarly,  $[body(q_{1\text{gg}})]_{\phi_2} \subseteq body(q_2^\mathcal{O})$  and  $\phi_2(\bar{x}_{1\text{gg}}) = \bar{x}_2$ . Consider the BGPQ  $q'_{1\text{gg}}$  obtained from  $q_{1\text{gg}}$  by replacing every variable  $v$  by  $v_{t_1 t_2}$  where  $t_1 = \phi_1(v)$  and  $t_2 = \phi_2(v)$ . Clearly,  $q_{1\text{gg}} \equiv_{\mathcal{R}, \mathcal{O}} q'_{1\text{gg}}$  holds. Assume that the head of  $q'_{1\text{gg}}$  is  $q'_{1\text{gg}}(\bar{x}'_{1\text{gg}})$ ; by construction it has the same head as  $q$ . Let us show that  $q \models_{\mathcal{R}, \mathcal{O}} q'_{1\text{gg}}$  holds, i.e., there exists a homomorphism  $\phi$  from the variables in  $q'_{1\text{gg}}$  to values in  $q^\mathcal{O}$  such that  $[body(q'_{1\text{gg}})]_\phi \subseteq body(q^\mathcal{O})$  and  $\phi(\bar{x}'_{1\text{gg}}) = \bar{x}$ . By Definition of  $\models$ , it is easy to see that the above holds when  $\phi$  maps  $v_{t_1 t_2}$  either to  $t_1$  if  $t_1 = t_2$  and  $t_1 \in \mathcal{U} \cup \mathcal{L}$ , or

As an immediate consequence of the above results, we get the following worst-case time to compute an  $\text{lgg}$  of two BGPQs  $q_1$  and  $q_2$ , and its size. We assume given the saturation  $q_1^\infty$  and  $q_2^\infty$  w.r.t. the sets  $\mathcal{O}$  of ontological constraints and  $\mathcal{R}$  of RDF entailment rules under consideration, as the times to compute  $q_1^\infty$  and  $q_2^\infty$ , and their sizes, depend on the particular  $\mathcal{O}$  and  $\mathcal{R}$  at hand.

Figure 8 (right) displays the cover query of the BGPQs  $q_1 \widetilde{\mathcal{O}}$  and  $q_2 \widetilde{\mathcal{O}}$  shown in Figure 7. It is therefore (Theorem 8) an **lgg** of the BGPQs  $q_1$  and  $q_2$  shown in Figure 6 w.r.t. the ontological constraints shown in Figure 7, using the RDF entailment rules shown in Table 2.

## 5. ALGORITHMS

---

**Algorithm 1** Least general anti-unification: **lgau**

---

**In:**  $T_1 = (t_1, \dots, t_n)$  and  $T_2 = (t'_1, \dots, t'_n)$ , boolean *bnodes*  
**Out:** least general anti-unification  $T$  of  $T_1$  and  $T_2$

```
1: for  $i = 1$  to  $n$  do  $\triangleright$  for each pair of  $i$ th attributes
2:   if  $t_i = t'_i$  and  $t_i \in \mathcal{U} \cup \mathcal{L}$  then
3:      $t_i^T \leftarrow t_i$   $\triangleright$  generalization of a constant by itself
4:   else if bnodes then  $\triangleright$  cover graph case
5:      $t_i^T \leftarrow b_{t_i, t'_i}$   $\triangleright$  generalization by a blank node
6:   else  $\triangleright$  cover query case
7:      $t_i^T \leftarrow v_{t_i, t'_i}$   $\triangleright$  generalization by a variable
8:   end if
9: end for
10: return  $(t_1^T, \dots, t_n^T)$   $\triangleright ()$  when  $n = 0$ 
```

---

---

**Algorithm 2** Cover graph of two RDF graphs: **lgg4g**

---

**In:** RDF graphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$   
**Out:**  $\mathcal{G}$  is the cover graph of  $\mathcal{G}_1$  and  $\mathcal{G}_2$

```
1:  $\mathcal{G} \leftarrow \emptyset$ 
2: for all  $T_1 = (s_1, p_1, o_1) \in \mathcal{G}_1$  do
3:   for all  $T_2 = (s_2, p_2, o_2) \in \mathcal{G}_2$  with  $p_1 = p_2$  do
4:      $\mathcal{G} \leftarrow \mathcal{G} \cup \{\text{lgau}(T_1, T_2, \text{true})\}$ 
5:   end for
6: end for
7: return  $\mathcal{G}$ 
```

---

## 5.1 Least general anti-unification

Algorithm 1, called **lgau**, computes a least general anti-unifier  $(t_1^T, \dots, t_n^T)$  of two  $n$ -ary tuples of RDF values  $(t_1, \dots, t_n)$  and  $(t'_1, \dots, t'_n)$ , made of constants (URIs and literals), blank nodes and variables. This is achieved by setting the  $i^{\text{th}}$  value  $t_i^T$  of the output tuple to the least general generalization of the values found at the  $i^{\text{th}}$  positions of the two input tuples:  $t_i$  and  $t'_i$ . Recall that a pair of a same constant is generalized by that constant itself, while in all other cases the generalization is either a blank node in case of RDF graph triple values (Section 3.2) or a variable in case of query head answer variables or of query body triple pattern values (Section 4.2). Generalizing pairs of different values with blank nodes or variables is controlled with the Boolean *bnodes* parameter. Crucially, these generated blank nodes or variables adopt the consistent naming scheme devised in Sections 3.2 and 4.2 that allows us preserving the common input RDF graphs (resp. BGPs) structure across the anti-unifications of triples (resp. triple patterns).

## 5.2 lgg of RDF graphs

Following Definition 2, Algorithm 2, called **lgg4g**, computes the cover graph  $\mathcal{G}$  of two input RDF graphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$ :  $\mathcal{G}$  comprises the least general anti-unification of every pair of  $\mathcal{G}_1$  and  $\mathcal{G}_2$  triples with *same* property. Therefore, given two RDF graphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$ , a call **lgg4g**( $\mathcal{G}_1, \mathcal{G}_2$ ) produces the cover graph-based **lgg** of  $\mathcal{G}_1$  and  $\mathcal{G}_2$  ignoring RDF entailment (Theorem 2), while a call **lgg4g**( $\mathcal{G}_1^\infty, \mathcal{G}_2^\infty$ ) produces the cover graph-based **lgg** of  $\mathcal{G}_1$  and  $\mathcal{G}_2$  taking into account the set of RDF entailment rules at hand (Theorem 3). In the latter case, the input RDF graphs can be saturated using standard algorithms implemented in RDF data management systems, like Jena [4] and Virtuoso [9].

Importantly, **lgg4g** assumes that the input RDF graphs, as well as their output cover graph, fit in memory. Checking

---

**Algorithm 3** Cover graph of two RDF graphs: **lgg4g-dms**

---

**In:** cursor  $c_1$  on RDF graph  $\mathcal{G}_1$ , cursor  $c_2$  on RDF graph  $\mathcal{G}_2$ , write access to an empty RDF graph  $\mathcal{G}$ , integer  $n$   
**Out:**  $\mathcal{G}$  is the cover graph of  $\mathcal{G}_1$  and  $\mathcal{G}_2$

```
1:  $c_1.\text{init}()$   $\triangleright c_1$  at beginning of  $\mathcal{G}_1$ 
2: while  $B_1 = c_1.\text{next}(n)$  do  $\triangleright$  fetch next  $n$   $\mathcal{G}_1$  triples
3:    $c_2.\text{init}()$   $\triangleright c_2$  at beginning of  $\mathcal{G}_2$ 
4:   while  $B_2 = c_2.\text{next}(n)$  do  $\triangleright$  fetch next  $n$   $\mathcal{G}_2$  triples
5:     for all  $T_1 = (s_1, p_1, o_1) \in B_1$  do
6:       for all  $T_2 = (s_2, p_2, o_2) \in B_2$  with  $p_1 = p_2$  do
7:         insert lgau( $T_1, T_2, \text{true}$ ) into  $\mathcal{G}$ 
8:       end for
9:     end for
10:   end while
11: end while
```

---

whether this is the case for the input RDF graphs under consideration can be done as follows.

The size of the input RDF graphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$  can be computed with the following SPARQL queries counting how many triples each of them holds: **SELECT count(\*) as ?size FROM  $\mathcal{G}_i$  with  $i \in [1, 2]$** . Recall that the worst-case size of the output cover graph is  $|\mathcal{G}| = |\mathcal{G}_1| \times |\mathcal{G}_2|$  in the unlikely case where *all* the  $\mathcal{G}_1$  and  $\mathcal{G}_2$  triples use the same property (Property 2 and Corollary 1).

The precise size of the output cover graph  $\mathcal{G}$  can be computed, *without* computing  $\mathcal{G}$ , with SPARQL queries. First, we calculate for each input RDF graph  $\mathcal{G}_i$ , with  $i \in [1, 2]$ , how many triples it holds per distinct property  $p$ :

$$S_{\mathcal{G}_i} = \{(p, n_i) \mid |\{(s, p, o) \in \mathcal{G}_i\}| = n_i\}$$

This can be computed with the SPARQL query: **SELECT ?p count(\*) as ?n<sub>i</sub> FROM  $\mathcal{G}_i$  WHERE  $\{(s, ?p, ?o)\}$  GROUP BY ?p**. Then, since every  $\mathcal{G}_1$  triple with property  $p$  anti-unifies with every  $\mathcal{G}_2$  triple with same property  $p$ , in order to beget  $\mathcal{G}$ , the size of  $\mathcal{G}$  is:

$$|\mathcal{G}| = \sum_{(p, n_1) \in S_{\mathcal{G}_1}, (p, n_2) \in S_{\mathcal{G}_2}} n_1 \times n_2$$

Overall,  $|\mathcal{G}|$  can be computed with the SPARQL query: **SELECT SUM(?n<sub>1</sub>\*?n<sub>2</sub>) as ?size WHERE  $\{\{S_{\mathcal{G}_1}\}\{S_{\mathcal{G}_2}\}\}$**  with  $S_{\mathcal{G}_1}$  and  $S_{\mathcal{G}_2}$  denoting the above-mentioned SPARQL queries computing these two sets, which join on their common answer variable  $?p$ .

When the input RDF graphs or their output cover graph cannot fit in memory, we propose variants of **lgg4g** that either assume that RDF graphs are stored in data management systems (DMSs, in short) or in a MapReduce cluster.

### 5.2.1 Handling large RDF graphs using DMSs

Algorithm 3, called **lgg4g-dms**, is an adaptation of **lgg4g**, which assumes that the input RDF graphs (already saturated if needed) and their cover graph are all stored in one or several DMSs. It further assumes that the system(s) storing the input RDF graphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$  feature(s) the well-known database mechanism of *cursor* [24, 42]. This is for instance the case for RDF graphs stored in relational servers like DB2 [1], MySQL [5], Oracle [6] and PostgreSQL [7], or in RDF servers like Jena-TDB [4] and Virtuoso [9]. Roughly speaking, a cursor is a pointer or iterator on tuples hold in a DMS, e.g., stored a relation or computed as the results to a query, that can be used to access these tuples. In particular,

a cursor can be used by an application to iteratively traverse all the tuples by fetching  $n$  of them at a time.

**lgg4g-dms** uses cursors to proceed similarly to **lgg4g** (remark that lines 5-9 in Algorithm 3 are almost the same as lines 2-6 in Algorithm 2) on pairs of  $n$ -triples subsets of  $\mathcal{G}_1$  and of  $\mathcal{G}_2$ , instead of on the whole RDF graphs themselves. It follows that the worst-case number of triples kept in memory by **lgg4g-dms** is  $M = (2 \times n) + 1$  at line 7 (i.e.,  $B_1$ ,  $B_2$  and the anti-unification triple output by **lgau**), with:

$$3 \leq M \leq |\mathcal{G}_1| + |\mathcal{G}_2| + 1.$$

The above lower bound is met for  $n$  set to 1, while the upper one is met for  $n$  set to  $\max(|\mathcal{G}_1|, |\mathcal{G}_2|)$ . Importantly, **lgg4g-dms** allows *choosing* the value of  $n$  in order to reflect the memory devoted to handling triples. For instance, if one wants to use 4GB of RAM for triples, assuming that any triple fits in less one 1KB (this value is much less when using dictionary encoding), the value of  $n$  can be set to 2M.

This clearly contrasts with the worst-case number of triples kept in memory by **lgg4g**:  $M = |\mathcal{G}_1| + |\mathcal{G}_2| + |\mathcal{G}|$  at line 7, with:

$$|\mathcal{G}_1| + |\mathcal{G}_2| \leq M \leq |\mathcal{G}_1| + |\mathcal{G}_2| + (|\mathcal{G}_1| \times |\mathcal{G}_2|).$$

The above lower bound is met when  $\mathcal{G}_1$  and  $\mathcal{G}_2$  have no property in common in their triples (i.e.,  $|\mathcal{G}| = 0$ ), while the upper one is met in the unlikely case where  $\mathcal{G}_1$  and  $\mathcal{G}_2$  use a same property in all their triples (i.e.,  $|\mathcal{G}| = |\mathcal{G}_1| \times |\mathcal{G}_2|$ ).

### 5.2.2 Handling huge RDF graphs using MapReduce

Algorithm 4, called **lgg4g-mr**, is a MapReduce (MR) variant of **lgg4g**. MR is a popular massively parallel programming framework [23], implemented by many large-scale data processing systems, like Hadoop [2] and Spark [8], which orchestrate clusters of compute nodes.

A MR program is organized in successive *jobs*, each of which comprises a *Map task* followed by a *Reduce task*. The Map task consists in reading some input data from the distributed file system<sup>4</sup> of the cluster, so as to partition the data into  $\langle \mathbf{k}, \mathbf{v} \rangle$  key-value pairs. Importantly, an MR engine transparently processes the Map task by running *Mapper processes* in parallel on cluster nodes, each process taking care of partitioning a portion of the input data by applying a **Map**(key: file, value: data unit) function on every data unit of a given input file. Key-value pairs thus produced are shuffled across the network, so that *all* pairs with *same* key  $\langle \mathbf{k}, \mathbf{v}_1 \rangle \dots \langle \mathbf{k}, \mathbf{v}_n \rangle$  are shipped to a same compute node. The Reduce task then consists in running *Reducer processes* in parallel, for every distinct key  $\mathbf{k}$  received by every compute node. Each process takes care of the set  $\mathcal{V}$  of values  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  emitted with key  $\mathbf{k}$ , by applying a **Reduce**(key:  $\mathbf{k}$ , values:  $\mathcal{V}$ ) function, and writing its results in a file. The result of an MR job, comprises the data, stored in a distributed fashion, in all the files output by Reducers.

In **lgg4g-mr**, the **Map** function applies to every  $(\mathbf{s}_i, \mathbf{p}_i, \mathbf{o}_i)$  triple of the input RDF graph  $\mathcal{G}_i$  stored in file  $G_i$ , and produces the corresponding key-value pair  $\langle \mathbf{p}_i, (G_i, (\mathbf{s}_i, \mathbf{p}_i, \mathbf{o}_i)) \rangle$ , for  $i \in [1, 2]$ . Hence, all the  $\mathcal{G}_1$  and  $\mathcal{G}_2$  triples with a same key/property  $p$  are shipped to the same cluster node. Then, similarly to **lgg4g** at lines 2-6, the **Reduce** functions process, on each node, the set  $\mathcal{V}$  of values emitted for every received

<sup>4</sup>For simplicity, we assume that input and output data of an MR job is stored on disk, like in Hadoop, while it may also reside in in-memory shared data structures, like in Spark.

---

#### Algorithm 4 Cover graph of two RDF graphs: **lgg4g-mr**

---

**In:** file  $G_1$  for RDF graph  $\mathcal{G}_1$ , file  $G_2$  for RDF graph  $\mathcal{G}_2$   
**Out:**  $\mathcal{G}$  is the cover graph of  $\mathcal{G}_1$  and  $\mathcal{G}_2$ , stored in  $G$ -\* files  
**Map**(key: file  $G_i$ , value: triple  $T_i = (\mathbf{s}_i, \mathbf{p}_i, \mathbf{o}_i)$ )  
1: **emit**( $\langle \mathbf{p}_i, (G_i, T_i) \rangle$ )  
**Reduce**(key:  $p$ , values: set  $\mathcal{V}$  of values emitted for key  $p$ )  
1:  $f \leftarrow \text{open}(G-p)$   
2: **for all**  $(G_1, T_1 = (\mathbf{s}_1, p, \mathbf{o}_1)) \in \mathcal{V}$  **do**  
3:   **for all**  $(G_2, T_2 = (\mathbf{s}_2, p, \mathbf{o}_2)) \in \mathcal{V}$  **do**  
4:      $f.write(\text{lgau}(T_1, T_2, \text{true}))$   
5:   **end for**  
6: **end for**  
7: **close**( $f$ )

---



---

#### Algorithm 5 Cover query of two BGPQs: **lgg4q**

---

**In:** BGPQs  $q_1$  and  $q_2$   
**Out:** Cover query  $q$  of  $q_1$  and  $q_2$  if it exists, else  $\perp$   
1:  $body(q) \leftarrow \emptyset$   
2: **for all**  $T_1 = (\mathbf{s}_1, \mathbf{p}_1, \mathbf{o}_1) \in \mathcal{G}_1$  **do**  
3:   **for all**  $T_2 = (\mathbf{s}_2, \mathbf{p}_2, \mathbf{o}_2) \in \mathcal{G}_2$  **do**  
4:      $body(q) \leftarrow body(q) \cup \{\text{lgau}(T_1, T_2, \text{false})\}$   
5:   **end for**  
6: **end for**  
7:  $head(q) \leftarrow \text{lgau}(head(q_1), head(q_2), \text{false})$   
8: **if**  $q$  is a well-formed BGPQ **then**  
9:   **return**  $q$   
10: **else**  
11:   **return**  $\perp$    ▷ error: a cover query does not exist  
12: **end if**

---

key  $p$ . The least general anti-unification triples obtained at line 4 are stored in the output file  $G-p$ . At the end of the MR job, the **lgg**  $\mathcal{G}$  of  $\mathcal{G}_1$  and  $\mathcal{G}_2$  is stored in the  $G$ -\* files of the distributed file system, where  $*$  denotes any key/property  $p$ .

In **lgg4g-mr**, the **Map** function holds at most a single  $\mathcal{G}_1$  or  $\mathcal{G}_2$  triple in memory. In contrast, the worst-case number of triples handled by the **Reduce** function for a given key  $p$  is:  $M = |\mathcal{G}_1| + |\mathcal{G}_2| + 1$  at line 4. This upper bound is met in the unlikely case where  $\mathcal{G}_1$  and  $\mathcal{G}_2$  use the same property  $p$  in all their triples. Similarly to **lgg4g-dms**, this upper bound can set to  $M = (2 \times n) + 1$ , with  $3 \leq M \leq |\mathcal{G}_1| + |\mathcal{G}_2| + 1$ , by first splitting the input RDF graphs in  $k_i$  files of  $n$   $\mathcal{G}_i$  triples (files  $G_i^1, \dots, G_i^{k_i}$ ), and then by processing every pair of  $n$ -triples of  $\mathcal{G}_1$  and  $\mathcal{G}_2$  files with an MR job (i.e., with  $k_1 \times k_2$  jobs), instead of a single MR job for the entire two input RDF graphs.

Finally, to take into account RDF entailment, the input RDF graphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$  can be saturated before being stored in the MR cluster, by using standard (centralized) techniques, or within the MR cluster, by using MR-based saturation techniques [46, 45]. Also, it is worth noting that RDF graphs, thus **lggs** of RDF graphs, stored in an MR cluster can be queried with MR-based SPARQL engines [29, 32, 36, 25].

### 5.3 lggs of BGPQs

Following Definition 6, Algorithm 5, called **lgg4q**, builds and returns either the cover query  $q$  of two input BGPQs  $q_1$  and  $q_2$  if it exists, or the  $\perp$  symbol otherwise.

At lines 2-6, **lgg4q** computes the body of  $q$ , which com-

---

**Algorithm 6** BGPQ saturation w.r.t. constraints: **sat4q**

---

**In:** BGPQ  $q$ , set  $\mathcal{R}$  of RDF entailment rules and set  $\mathcal{O}$  of RDFS statements  
**Out:**  $q^\infty$  is the saturation of  $q$  w.r.t.  $\mathcal{R}$  and  $\mathcal{O}$   
1:  $head(q^\infty) \leftarrow head(q)$   
2:  $\mathcal{O}_s \leftarrow \mathcal{O}^\infty$   
3:  $bq_s \leftarrow body(q)^\infty$   
4: rename  $bq_s$  blank nodes shared with  $\mathcal{O}_s \triangleright$  they are local  
5:  $body(q^\infty) \leftarrow bq_s \cup ((bq_s \cup \mathcal{O}_s)^\infty \setminus \mathcal{O}_s)$   
6: **return**  $q^\infty$

---

prises a least general anti-unifier of every pair of  $q_1$  body and  $q_2$  body triples (first item in Definition 6). Then, at line 7, it computes the head of  $q$  as the least general anti-unifier of the heads of  $q_1$  and  $q_2$  (second item in Definition 6). Finally, **lgg4q** checks whether  $q$  is a well-formed BGPQ, i.e., if its answer variables occur in its body triples, otherwise there is no cover query of  $q_1$  and  $q_2$ , and also no **lgg** of  $q_1$  and  $q_2$  (recall first item of Theorems 7 and 8).

Therefore, given two BGPQs  $q_1$  and  $q_2$ , a call **lgg4q**( $q_1, q_2$ ) produces the cover query-based **lgg** of  $q_1$  and  $q_2$  ignoring RDF entailment and extra ontological constraints, whenever it exists (Theorem 7), while a call **lgg4q**( $q_1^\infty, q_2^\infty$ ) produces the cover query-based **lgg** of  $q_1$  and  $q_2$  taking into account the set of RDF entailment rules at hand, as well as a set of extra ontological constraints, whenever it exists (Theorem 8). In the latter case, the saturated input BGPQs can be computed using Algorithm 6, called **sat4q**.

Following Definition 3, **sat4q** computes the saturation of a query  $q$  w.r.t. a set  $\mathcal{R}$  of RDF entailment rules and a set  $\mathcal{O}$  of RDFS constraints as a BGPQ  $q^\infty$  with same head as  $q$  (line 1), and whose body is computed based on its characterization displayed in Figure 5:  $body(q^\infty) = body(q)^\infty \cup ((body(q) \cup \mathcal{O})^\infty \setminus \mathcal{O}^\infty)$  (\*). Further, since  $(body(q) \cup \mathcal{O})^\infty$  is equivalent to  $(body(q)^\infty \cup \mathcal{O}^\infty)^\infty$  (because  $body(q)$  and  $\mathcal{O}$  are equivalent w.r.t.  $\mathcal{R}$  to  $body(q)^\infty$  and  $\mathcal{O}^\infty$  respectively),  $body(q^\infty) = body(q)^\infty \cup ((body(q)^\infty \cup \mathcal{O}^\infty)^\infty \setminus \mathcal{O}^\infty)$  (\*\*) holds. **sat4q** uses (\*\*) instead of (\*) since, *once*  $body(q)^\infty$  and  $\mathcal{O}^\infty$  are computed (lines 2 and 3), clearly  $(body(q)^\infty \cup \mathcal{O}^\infty)^\infty$  is faster to compute than  $(body(q) \cup \mathcal{O})^\infty$ , thus (\*\*) is faster to compute (line 5) than (\*).

## 6. EXPERIMENTS

We provide an experimental assessment of our technical contribution for computing **lggs** of BGPQs, which goes beyond that for RDF graphs (recall that BGPQ bodies generalize RDF graphs). In Section 6.1, we describe our experimental settings. Then, in Section 6.2, we study the saturation of a BGPQ w.r.t. ontological constraints. Finally, in Section 6.3, we study the computation of **lggs** with or without considering constraints, notably the gain in precision when constraints are considered.

### 6.1 Settings

**Software.** We implemented our framework for computing **lggs** of RDF graphs or of BGPQs in Java 1.8 [3], on top of the Jena 3.0.1 RDF reasoner [4] and a PostgreSQL 9.3.11 server [7]; all used with default settings.

We use Jena to compute the saturation of an RDF graph (against which queries must be evaluated to obtained their complete answer sets, recall Section 2), as well as the satu-

ration of queries w.r.t. extra ontological constraints that we devised in the preceding section.

PostgreSQL is used to evaluate *SQLized* BGPQs against a saturated RDF graph. This RDF graph, which comprises both data and constraints (i.e., RDF and RDFS statements), is stored in a **Triple**(**s**,**p**,**o**) PostgreSQL table, indexed by all permutations of the **s**, **p**, **o** columns, leading to a total of 6 indexes. This indexing choice is inspired by [35, 50], to give PostgreSQL efficient query evaluation opportunities.

There exist alternative data layouts like having one two-columns table **p**(**s**,**o**) per distinct property **p** in the RDF graph, which stores all the subject/object tuples of triples with property **p** [16], or the elaborate layout of [14] used in DB2 RDF. However, adopting another data layout than the **Triple** table would have no impact on our experiments, as the reported times are not related to query evaluation.

**RDF entailment rules.** We used the rules shown in Table 2, which allow reasoning about RDFS ontological constraints.

**Dataset.** We conducted experiments using LUBM [28]. We generated an RDF graph comprising 884k triples (including 242 ontological constraints, i.e., RDFS statements) before saturation and 1.08M triples after. This RDF graph is used to compare the quality of **lggs** of BGPQs w.r.t. their precision, when ontological are considered or not.

**Queries.** We borrowed from [15] a set of 9 BGPQs. Table 3 (top) displays the characteristics of these queries, which have a variety of structural aspects and numbers of answers. The queries can be found in the Appendix.

**Hardware.** We used an Intel Xeon (X5550) 2.67GHz machine with 32GB RAM, using Ubuntu 14.04.3 LTS (64bits).

In the sequel, all measured times are averaged over 5 warm runs and are in milliseconds.

### 6.2 BGPQ saturation w.r.t. constraints

Following Definition 3, we compute the saturation of a BGPQ  $q$  w.r.t. a set  $\mathcal{O}$  of ontological constraints as follows. We build a saturated query  $q^\infty$  with same output variables as  $q$  and whose body comprises the triples in  $(body(q) \cup \mathcal{O})^\infty$ , except those in  $\mathcal{O}^\infty$  that are not in  $body(q)^\infty$ , i.e.,  $body(q^\infty) = (body(q) \cup \mathcal{O})^\infty \setminus (\mathcal{O}^\infty \setminus body(q)^\infty)$ .

Table 3 (bottom) shows the size of our saturated test queries and the time to compute them. Enriching our test queries using the 242 LUBM constraints significantly augments their size: from  $\times 2$  for  $Q_{10}$  up to  $\times 5$  for  $Q_{22}$ . The query saturation time is always fast: a very few tens of milliseconds for all our test queries.

### 6.3 lggs of BGPQs w.r.t. constraints

Table 4 (lines 1 and 3) shows that cover query-based **lggs** of test queries are always computed fast whether or not we consider the LUBM constraints: 1 to 4ms when they are ignored, and 5 to 10ms when they are considered. In the latter case, overall, it takes between 50 and 59ms to compute an **lgg** in the worst case (i.e., when the two saturated test queries are computed *in sequence* before computing their cover query).

Further, observe that the cover query-based **lgg** of two BGPQs when ignoring extra ontological constraints, is *entailed* by that of these same queries when extra constraints are considered. Indeed, by construction, these **lggs** have the same answer variables and, clearly, the body of the for-

Queries:	$Q_{01}$	$Q_{03}$	$Q_{05}$	$Q_{07}$	$Q_{08}$	$Q_{10}$	$Q_{12}$	$Q_{22}$	$Q_{25}$
shape	star	star	graph	graph	star	tree	graph	star	star
size (number of triples)	3	3	4	3	2	4	6	2	3
number of constants/distinct variables/answer variables	5/2/2	5/2/2	5/3/3	4/2/2	3/2/2	6/3/3	7/4/3	2/3/2	5/2/2
cardinality (number of answers)	123	41	869	0	269	41 751	79	14 252	16
size of the saturation w.r.t. constraints	8	9	11	7	6	8	14	10	8
time to compute the saturation w.r.t. constraints	24	23	23	21	22	22	21	27	22

Table 3: Characteristics of our test queries (top) and of their saturations w.r.t. LUBM constraints (bottom); times are in ms.

Pairs of queries:	$Q_{01}Q_{07}$	$Q_{03}Q_{07}$	$Q_{01}Q_{22}$	$Q_{03}Q_{22}$	$Q_{08}Q_{22}$	$Q_{01}Q_{25}$	$Q_{22}Q_{25}$	$Q_{05}Q_{10}$	$Q_{05}Q_{12}$
time to compute the <b>lgg</b>	1	2	1	1	1	3	2	2	4
cardinality (number of answers) of the <b>lgg</b>	30 963	1 048 360	1 048 360	74 643	1 048 360	253 443	1 048 360	6 937 472	34 852
time to compute the <b>lgg</b> w.r.t. constraints	5	6	6	7	10	6	6	6	8
cardinality of the <b>lgg</b> w.r.t. constraints	14 285	14 285	56 358	56 358	185 124	19 053	56 358	520 365	57 774

Table 4: Characteristics of cover query-based **lggs** of test queries, w or w/o using the LUBM constraints; times are in ms.

mer is included in that of the latter (as it is built using the saturated queries and not the queries themselves). This observation allows comparing their numbers of answers in order to *quantify* to which extent **lggs** ignoring extra ontological constraints are unnecessarily more general than **lggs** which take them into account, i.e., to which extent **lggs** considering extra ontological constraints are more *precise* than those ignoring such constraints. Table 4 (lines 2 and 4) shows that ignoring extra constraints significantly increases the number of answers for some **lggs**, from a small  $\times 1.32$  for  $Q_{03}Q_{22}$  up to a striking  $\times 73.39$  for  $Q_{03}Q_{07}$ , with a significant average of  $\times 16.45$ . **lggs** with same number of answers in Table 4 were found either equivalent or almost equivalent so that their semantic difference is not visible on the LUBM dataset, e.g., the **lggs** of  $Q_{22}Q_{25}$ ,  $Q_{08}Q_{22}$  and  $Q_{01}Q_{22}$  with 1 048 360 answers (i.e., ignoring ontological constraints) are asking for all the (distinct) subject/object pairs in triples in the LUBM RDF graph; the **lggs** of  $Q_{22}Q_{25}$ ,  $Q_{03}Q_{22}$  with 56 358 answers (i.e., considering ontological constraints) ask respectively for faculties with some degree from some university and for employees with some doctoral degree from some university: they have the same number of answers in our LUBM RDF graph because all faculties have some doctoral degree and all employees with some doctoral degree are faculties.

## 7. RELATED WORK

The problem of computing an **lgg** was introduced in the early 70's by G. Plotkin [40, 41] to generalize First Order Logic clauses w.r.t.  $\theta$ -subsumption, a non-standard logical implication typical of Machine Learning. This problem has also been investigated in Knowledge Representation, for formalisms whose expressivity overlaps with our RDF setting, notably Description Logics (DLs) [31, 12, 51] and Conceptual Graphs (CGs) [19]. Finally, recently, this problem has started receiving attention in the Semantic Web field [22, 21, 33].

In DLs, computing an **lgg** of concepts (formulae) has been studied for  $\mathcal{EL}$  and extensions thereof [31, 12, 51]. The  $\mathcal{EL}$  setting translates into *particular tree-shaped* RDF graphs, which may feature RDFS subclass and domain constraints, and for which RDF entailment is limited to the use of these

two constraints only<sup>5</sup>. In these equivalent RDF and  $\mathcal{EL}$  fragments, the  $\mathcal{EL}$  technique that computes an **lgg** of  $\mathcal{EL}$  concepts, which is an  $\mathcal{EL}$  concept, provides *only a (non least general) generalization* of their corresponding tree-shaped RDF graphs w.r.t. the problem we study: the (minimal) cover graph-based **lgg** of tree-shaped RDF graphs is clearly a forest-shaped RDF graph in general. However, it can be shown that our technique for general RDF graph can be used to compute the  $\mathcal{EL}$  **lgg** of  $\mathcal{EL}$  concepts as shown in Example 1 below. Roughly speaking, given the  $\mathcal{EL}$  concepts  $C_1, \dots, C_n$  and their corresponding tree-shaped RDF graphs  $\mathcal{G}(C_1, b_1), \dots, \mathcal{G}(C_n, b_n)$ , the  $\mathcal{EL}$  **lgg** is that corresponding to the tree-shaped RDF graph rooted in  $b_{b_1, \dots, b_n}$  within the forest-shaped cover graph of  $\mathcal{G}(C_1, b_1), \dots, \mathcal{G}(C_n, b_n)$ .

In CGs, the so-called *simple CGs with unary and binary relations* correspond to *particular* RDF graphs (e.g., a property URI in a triple cannot be the subject or object of another triple, a class - URI or blank node - in a  $\tau$  triple cannot be the subject of another  $\tau$  triple nor the subject or object of another non- $\tau$  triple, etc), which may feature the four RDFS constraints, and for which RDF entailment is limited to the use of these RDFS constraints only [13]. In these equivalent RDF and CG fragments, we may interchangeably compute **lggs** with the CG technique in [19] or ours.

In RDF, computing an **lgg** has been studied for *particular* RDF graphs, called *r-graphs, ignoring RDF entailment* [22, 21]. An *r-graph* is an *extracted subgraph* of an RDF graph  $\mathcal{G}$ , rooted in the  $\mathcal{G}$  value  $r$  and comprising the  $\mathcal{G}$  triples *reachable from  $r$  through directed paths of length at most  $n$* . Such a rooted and directed *r-graph* can be defined recursively as  $\mathcal{S}(\mathcal{G}, r, n)$ , with:

$$\mathcal{S}(\mathcal{G}, v, 0) = \emptyset$$

$$\mathcal{S}(\mathcal{G}, v, n) = \bigcup_{(v, p, v') \in \mathcal{G}} \{(v, p, v')\} \cup \mathcal{S}(\mathcal{G}, v', n-1) \cup \mathcal{S}(\mathcal{G}, p, n-1)$$

<sup>5</sup>An  $\mathcal{EL}$  concept  $C$  recursively translates into the RDF graph rooted in the blank node  $b_r$  returned by the call  $\mathcal{G}(C, b_r)$ , with:  $\mathcal{G}(\top, b) = \emptyset$  for the universal  $\mathcal{EL}$  concept  $\top$ ,  $\mathcal{G}(A, b) = \{(b, \tau, A)\}$  for an atomic  $\mathcal{EL}$  concept  $A$ ,  $\mathcal{G}(\exists r.C, b) = \{(b, r, b')\} \cup \mathcal{G}(C, b')$ , with  $b'$  a fresh blank node, for an  $\mathcal{EL}$  existential restriction  $\exists r.C$ , and  $\mathcal{G}(C_1 \sqcap C_2, b) = \mathcal{G}(C_1, b) \cup \mathcal{G}(C_2, b)$  for an  $\mathcal{EL}$  conjunction  $C_1 \sqcap C_2$ ; the  $\mathcal{EL}$  constraints  $A_1 \sqsubseteq A_2$  and  $\exists r.\top \sqsubseteq A$  correspond to  $(A_1, \preceq_{sc}, A_2)$  and  $(r, \leftarrow_d, A)$  resp.



Intuitively, this purely structural definition of  $r$ -graph attempts carrying  $\mathcal{G}$ 's knowledge about  $r$ . **lggs** of  $r$ -graphs allow finding the commonalities between *single root entities*, while with general RDF graphs we further allow finding the commonalities between *sets of multiple interrelated entities*. The technique for computing an **lgg** of two  $r$ -graphs exploits their rooted and directed structure: it starts from their respective root and traverses them simultaneously considering triples reachable through directed paths of increasing size, while incrementally constructing an  $r$ -graph **lgg**. In contrast, the general RDF graphs we consider are unstructured; our technique blindly traverses the input RDF graphs to anti-unify their triples with same property, and captures their common structure across these anti-unifications thanks to the consistent naming scheme we devised for the blank nodes they generate. Further, when we ignore RDF entailment, computing **lggs** of  $r$ -graphs or of RDF graphs have the same worst-case time complexity ( $O(|\mathcal{G}_1| \times |\mathcal{G}_2|)$ , with  $\mathcal{G}_1, \mathcal{G}_2$  the input RDF graphs). This is for instance the case for the following star-shaped RDF graphs, which are trees,  $r$ -graphs and RDF graphs,  $\mathcal{G}_1 = \{(r_1, p, s_1^1), \dots, (r_1, p, s_1^m)\}$  and  $\mathcal{G}_2 = \{(r_2, p, s_2^1), \dots, (r_2, p, s_2^n)\}$ , the  $r$ -graph and cover graph-based **lgg** of which is built at some point by both techniques is:  $\mathcal{G}_{\text{lgg}} = \{(b_{r_1 r_2}, p, b_{s_1^1 s_2^1}), \dots, (b_{r_1 r_2}, p, b_{s_1^m s_2^n})\}$ . Also, as noted in [21], the computed  $r$ -graphs **lggs** are *only (non least general) generalizations* of  $r$ -graphs w.r.t. the standard semantics of RDF graphs defined upon RDF entailment. Finally, the  $r$ -graph technique that computes an **lgg** of  $r$ -graphs, which is an  $r$ -graph, gives *only a (non least general) generalization* of them w.r.t. the problem we study: the (minimal) cover graph-based **lgg** of  $r$ -graphs is clearly a general RDF graph as shown in Example 1 below.

In SPARQL, computing an **lgg** has been considered for *unary tree-shaped conjunctive queries (UTCQ)* [33]; a UTCQ **lgg** is computed by a simultaneous root-to-leaves traversal of the input queries. UTCQs are tree-shaped RDF graphs, when variables are viewed as blank nodes, for which RDF entailment is ignored [21]. The UTCQ technique that computes an **lgg** of UTCQs, which is a UTCQ, yields *only a (non least general) generalization* of their corresponding tree-shaped RDF graphs w.r.t. the problem we study: the (minimal) cover graph-based **lgg** of tree-shaped RDF graphs is clearly a forest-shaped RDF graph. However, similarly as for the  $\mathcal{EL}$  description logic above, it can be shown that our technique for general RDF graph can be used to compute the UTCQ **lgg** of UTCQs as shown in Example 1 below.

**EXAMPLE 1.** *Let us consider the two tree-shaped RDF graphs below that may correspond to  $\mathcal{EL}$  concepts,  $r$ -graphs or UTCQs:*

$\mathcal{G}_1 = \{(b_1, p_e, b_{11}), (b_1, p_s, b_{12})\}$  and  
 $\mathcal{G}_2 = \{(b_2, p_e, b_{21}), (b_{21}, p_s, b_{211})\}$ .

*Their RDF graph **lgg** is the general forest-shaped RDF graph  $\mathcal{G} = \{(b_{b_1 b_2}, p_e, b_{b_{11} b_{21}}), (b_{b_1 b_{21}}, p_s, b_{b_{12} b_{211}})\}$  while their  $\mathcal{EL}$ ,  $r$ -graph, as well as UTCQ **lgg** is the strictly more general tree-shaped RDF graph:  $\mathcal{G}' = \{(b_{b_1 b_2}, p_e, b_{b_{11} b_{21}})\}$ .*

## 8. CONCLUSION AND PERSPECTIVES

We have revisited the Machine Learning problem of computing a *least general generalization (lgg)* of some descriptions in the setting of RDF and SPARQL. Our contributions significantly extend the state of the art by considering the entire RDF standard and popular conjunctive fragment of

SPARQL, i.e., BGPQs. In particular, we neither restrict RDF graphs and BGPQs nor RDF entailment in any way, while closely related works only consider *rooted*-RDF graphs and *unary tree* BGPQs (i.e., describing a *single* root resource or answer variable, respectively), and, further, they *completely ignore* RDF entailment rules (i.e., their particular RDF graphs and BGPQs are simply compared in a *purely relational fashion*). Moreover, in the case of BGPQs, we also endowed with background knowledge the standard RDF entailment relation between queries. As our experiments showed, taking into account the ontological constraints describing the application domain in which queries are posed may enhance the precision of **lggs**.

As a short-term perspective, we want to study heuristics in order to efficiently prune out as much as possible redundant triples, while computing **lggs**. Indeed, as for instance Figures 4 and 8 show, our cover graph/query technique produces many redundant triples. This would allow having more readable **lggs**, as well as reducing the a posteriori elimination effort of redundant triples with standard technique from the literature. Moreover, in the case of BGPQs, removing redundant triples may significantly improves their evaluation time.

Computing **lggs** has many possible applications in databases (recall Section 1); in particular, we plan to apply our results for BGPQs to the optimization problem of view selection in RDF. Our idea is to select as views to materialize a set  $\mathcal{V}$  of **lggs** of queries from a workload, such that (i) the workload queries can be (partially or totally) rewritten using  $\mathcal{V}$  and (ii)  $\mathcal{V}$  minimizes a combination of query processing, view storage, and view maintenance costs. This would provide an alternative to [26], which applies to the so-called RDF database fragment of RDF (restricting RDF entailment) and BGPQs, and which is based on recursively decomposing the workload queries into subqueries that may be materialized.

Finally, we also want to investigate the problem of computing **lggs**, and apply it to view selection, in the setting of the DL-Lite $\mathcal{R}$  description logic [17], which underpins OWL2 QL, the other W3C's Semantic Web standard.

## Acknowledgment

This work has been partially funded by Lannion-Tregor Communauté and Région Bretagne (PAWS project).

## 9. REFERENCES

- [1] DB2. [www.ibm.com/analytics/us/en/technology/db2](http://www.ibm.com/analytics/us/en/technology/db2).
- [2] Hadoop. [hadoop.apache.org](http://hadoop.apache.org).
- [3] Java. [www.oracle.com/java](http://www.oracle.com/java).
- [4] Jena. [jena.apache.org](http://jena.apache.org).
- [5] MySQL. [www.mysql.com](http://www.mysql.com).
- [6] Oracle. [www.oracle.com/database](http://www.oracle.com/database).
- [7] PostgreSQL. [www.postgresql.org](http://www.postgresql.org).
- [8] Spark. [spark.apache.org](http://spark.apache.org).
- [9] Virtuoso. [virtuoso.openlinksw.com](http://virtuoso.openlinksw.com).
- [10] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [11] F. Baader, R. Kiisters, and R. Molitor. Computing least common subsumers in description logics with existential restrictions. In *IJCAI*, 1999.

- [12] F. Baader, B. Sertkaya, and A.-Y. Turhan. Computing the least common subsumer w.r.t. a background terminology. *Journal of Applied Logic*, 5(3), 2007.
- [13] J. Baget, M. Croitoru, A. Gutierrez, M. Leclère, and M. Mugnier. Translations between RDF(S) and conceptual graphs. In *ICCS*, 2010.
- [14] M. A. Bornea, J. Dolby, A. Kementsietsidis, K. Srinivas, P. Dantressangle, O. Udreă, and B. Bhattacharjee. Building an efficient RDF store over a relational database. In *SIGMOD*, 2013.
- [15] D. Bursztyn, F. Goasdoué, and I. Manolescu. Optimizing reformulation-based query answering in RDF. In *EDBT*, 2015.
- [16] D. Bursztyn, F. Goasdoué, and I. Manolescu. Teaching an RDBMS about ontological constraints. *PVLDB*, 9(12), 2016.
- [17] D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati. Tractable reasoning and efficient query answering in description logics: The *DL-Lite* family. *J. Autom. Reasoning*, 39(3), 2007.
- [18] D. Calvanese, G. De Giacomo, and M. Lenzerini. On the decidability of query containment under constraints. In *PODS*, 1998.
- [19] M. Chein and M. Mugnier. *Graph-based Knowledge Representation - Computational Foundations of Conceptual Graphs*. Springer, 2009.
- [20] W. W. Cohen, A. Borgida, and H. Hirsh. Computing least common subsumers in description logics. In *AAAI*, 1992.
- [21] S. Colucci, F. Donini, S. Giannini, and E. D. Sciascio. Defining and computing least common subsumers in RDF. *J. Web Semantics*, 39(0), 2016.
- [22] S. Colucci, F. M. Donini, and E. D. Sciascio. Common subsumers in RDF. In *AI\*IA*, 2013.
- [23] J. Dean and S. Ghemawat. Mapreduce: Simplified data processing on large clusters. In *OSDI*, 2004.
- [24] H. Garcia-Molina, J. D. Ullman, and J. Widom. *Database systems - the complete book*. Pearson Education, 2009.
- [25] F. Goasdoué, Z. Kaoudi, I. Manolescu, J. Quiané-Ruiz, and S. Zampetakis. Cliquesquare: Flat plans for massively parallel RDF queries. In *ICDE*, 2015.
- [26] F. Goasdoué, K. Karanasos, J. Leblay, and I. Manolescu. View selection in semantic web databases. *PVLDB*, 5(2), 2011.
- [27] F. Goasdoué, I. Manolescu, and A. Roatis. Efficient query answering against dynamic RDF databases. In *EDBT*, 2013.
- [28] Y. Guo, Z. Pan, and J. Heflin. LUBM: A benchmark for OWL knowledge base systems. *J. Web Sem.*, 3(2-3), Oct. 2005.
- [29] M. F. Husain, J. P. McGlothlin, M. M. Masud, L. R. Khan, and B. M. Thuraisingham. Heuristics-based query processing for large RDF graphs using cloud computing. *IEEE TKDE*, 23(9), 2011.
- [30] T. Imielinski and W. Lipski. Incomplete information in relational databases. *JACM*, 31(4), 1984.
- [31] R. Küsters. *Non-Standard Inferences in Description Logics*, volume 2100 of *LNCS*. Springer, 2001.
- [32] K. Lee and L. Liu. Scaling queries over big RDF graphs with semantic hash partitioning. *PVLDB*, 6(14), 2013.
- [33] J. Lehmann and L. Bühmann. Autosparql: Let users query your knowledge base. In *ESWC*, 2011.
- [34] M. Meier. Towards rule-based minimization of RDF graphs under constraints. In *RR*, 2008.
- [35] T. Neumann and G. Weikum. x-rdf-3x: Fast querying, high update rates, and consistency for RDF databases. *PVLDB*, 3(1), 2010.
- [36] N. Papailiou, D. Tsoumakos, I. Konstantinou, P. Karras, and N. Koziris. H<sub>2</sub>rdf+: an efficient data management system for big RDF graphs. In *SIGMOD*, 2014.
- [37] F. Picalausa, Y. Luo, G. H. Fletcher, J. Hidders, and S. Vansummeren. A structural approach to indexing triples. In *ESWC*, 2012.
- [38] R. Pichler, A. Polleres, S. Skritek, and S. Woltran. Redundancy elimination on RDF graphs in the presence of rules, constraints, and queries. In *RR*, 2010.
- [39] R. Pichler, A. Polleres, S. Skritek, and S. Woltran. Complexity of redundancy detection on RDF graphs in the presence of rules, constraints, and queries. *Semantic Web*, 4(4), 2013.
- [40] G. D. Plotkin. A note on inductive generalization. *Machine Intelligence*, 5, 1970.
- [41] G. D. Plotkin. A further note on inductive generalization. *Machine Intelligence*, 6, 1971.
- [42] R. Ramakrishnan and J. Gehrke. *Database management systems*. McGraw-Hill, 2003.
- [43] J. A. Robinson. A machine-oriented logic based on the resolution principle. *J. ACM*, 12(1), Jan. 1965.
- [44] J. A. Robinson and A. Voronkov, editors. *Handbook of Automated Reasoning*. Elsevier and MIT Press, 2001.
- [45] J. Urbani, S. Kotoulas, J. Maassen, F. van Harmelen, and H. E. Bal. Webpie: A web-scale parallel inference engine using mapreduce. *J. Web Sem.*, 10, 2012.
- [46] J. Urbani, S. Kotoulas, E. Oren, and F. van Harmelen. Scalable distributed reasoning using mapreduce. In *ISWC*, 2009.
- [47] Resource description framework 1.1. <https://www.w3.org/TR/rdf11-concepts>.
- [48] RDF 1.1 semantics. <https://www.w3.org/TR/rdf11-mt/>.
- [49] SPARQL protocol and RDF query language. <http://www.w3.org/TR/rdf-sparql-query>.
- [50] C. Weiss, P. Karras, and A. Bernstein. Hexastore: sextuple indexing for semantic web data management. *PVLDB*, 1(1), 2008.
- [51] B. Zarri  and A. Turhan. Most specific generalizations w.r.t. general EL-TBoxes. In *IJCAI*, 2013.

## 10. APPENDIX

<p>Q01(?X, ?Y) : –</p> <p>?X "http://www.w3.org/1999/02/22-rdf-syntax-ns#type" "http://swat.cse.lehigh.edu/onto/univ-bench.owl#Employee",</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#worksFor" "http://www.Department0.University0.edu",</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#degreeFrom" ?Y</p>
<p>Q03(?X, ?Y) : –</p> <p>?X "http://www.w3.org/1999/02/22-rdf-syntax-ns#type" "http://swat.cse.lehigh.edu/onto/univ-bench.owl#Employee",</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#worksFor" "http://www.Department0.University0.edu",</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#doctoralDegreeFrom" ?Y</p>
<p>Q05(?X, ?Y, ?Z) : –</p> <p>?X "http://www.w3.org/1999/02/22-rdf-syntax-ns#type" "http://swat.cse.lehigh.edu/onto/univ-bench.owl#Student",</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#advisor" ?Y,</p> <p>?Y "http://swat.cse.lehigh.edu/onto/univ-bench.owl#teacherOf" ?Z,</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#takesCourse" ?Z,</p>
<p>Q07(?X, ?Y) : –</p> <p>?X "http://www.w3.org/1999/02/22-rdf-syntax-ns#type" "http://swat.cse.lehigh.edu/onto/univ-bench.owl#Faculty",</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#degreeFrom" ?Y,</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#memberOf" ?Y</p>
<p>Q08(?X, ?Y) : –</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#degreeFrom" ?Y,</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#memberOf" "http://www.Department0.University0.edu",</p>
<p>Q10(?W, ?X, ?Y) : –</p> <p>?X "http://www.w3.org/1999/02/22-rdf-syntax-ns#type" "http://swat.cse.lehigh.edu/onto/univ-bench.owl#GraduateStudent",</p> <p>?Y "http://www.w3.org/1999/02/22-rdf-syntax-ns#type" "http://swat.cse.lehigh.edu/onto/univ-bench.owl#Faculty",</p> <p>?W "http://swat.cse.lehigh.edu/onto/univ-bench.owl#publicationAuthor" ?X,</p> <p>?W "http://swat.cse.lehigh.edu/onto/univ-bench.owl#publicationAuthor" ?Y</p>
<p>Q12(?W, ?X, ?Y) : –</p> <p>?X "http://www.w3.org/1999/02/22-rdf-syntax-ns#type" "http://swat.cse.lehigh.edu/onto/univ-bench.owl#GraduateStudent",</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#advisor" ?Y,</p> <p>?Y "http://swat.cse.lehigh.edu/onto/univ-bench.owl#teacherOf" ?Z,</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#takesCourse" ?Z,</p> <p>?W "http://swat.cse.lehigh.edu/onto/univ-bench.owl#publicationAuthor" ?X,</p> <p>?W "http://swat.cse.lehigh.edu/onto/univ-bench.owl#publicationAuthor" ?Y</p>
<p>Q22(?X, ?Y) : –</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#doctoralDegreeFrom" ?Z,</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#teacherOf" ?Y</p>
<p>Q25(?X, ?Y) : –</p> <p>?X "http://www.w3.org/1999/02/22-rdf-syntax-ns#type" "http://swat.cse.lehigh.edu/onto/univ-bench.owl#Faculty",</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#degreeFrom" "http://www.University532.edu",</p> <p>?X "http://swat.cse.lehigh.edu/onto/univ-bench.owl#memberOf" ?Y</p>

Figure 9: LUBM queries