

A Comparison of Query Rewriting Techniques for DL-Lite

Héctor Pérez-Urbina, Boris Motik, and Ian Horrocks

Computing Laboratory
University of Oxford
Oxford, UK

{hector.perez-urbina,boris.motik,ian.horrocks}@comlab.ox.ac.uk

1 Introduction

Answering queries under constraints is the problem of computing the answers to a query over an incomplete database w.r.t. a set of constraints [16]. Since the database is only partially specified, the task is to compute the so-called *certain answers*: the tuples that satisfy the query in every database instance that conforms to the partial specification and satisfies the constraints. Answering queries under constraints lies at the core of many problems, such as information integration [12], data exchange [9], and data warehousing [17].

There exist constraint languages for which the problem of query answering under constraints can be solved via query rewriting: given a query Q over an incomplete database, defined by a set of constraints C and a (partial) instance I , one computes a query a so-called rewriting Q' of Q w.r.t. C such that the answers to Q over C and I , and the answers to Q' over I coincide for every instance I . Thus, the set of constraints C is no longer needed to answer Q —in order to answer Q over an incomplete database consisting of C and any I , we can simply answer Q' over I . Description Logics (DLs) [3] can be viewed as very expressive but decidable first-order fragments, which makes them natural candidates for constraint languages. DLs are a family of knowledge representation formalisms that represent a given domain in terms of concepts (unary predicates), roles (binary predicates), and individuals (constants). A DL Knowledge Base (KB) consists of a *terminological* component called the TBox, and an *assertional* component called the ABox. In analogy to incomplete databases, the TBox can be seen as a conceptual schema containing the set of constraints C , and the ABox as some partial instance I of the schema. The use of DLs as constraint languages has already proven to be useful in a variety of scenarios such as Information Integration [7] and the Semantic Web [10].

Query rewriting under DL constraints for conjunctive queries has been considered by Rosati [15] for the \mathcal{EL} family of languages [2], and by Calvanese et al. for the DL-Lite family of languages [6]. The rewriting algorithm proposed by Calvanese et al., which we will refer to as the CGLLR algorithm, takes as input a conjunctive query Q and a DL-Lite_R TBox \mathcal{T} and computes a union of conjunctive queries that is a rewriting of Q w.r.t. \mathcal{T} . The CGLLR algorithm

lies at the core of reasoners like QUONTO¹ and OwlGres². The aforementioned approaches are closely related; however, they have been developed for two different DL families. Motivated by the prospect of developing a unified algorithm, in our previous work [13] we described a novel rewriting algorithm for the DL \mathcal{ELHI} [2]. The algorithm takes as input a conjunctive query Q and an \mathcal{ELHI} TBox \mathcal{T} , and produces a datalog query Q' that is a rewriting of Q w.r.t. \mathcal{T} . We have shown that our algorithm is worst-case optimal for various sublanguages of \mathcal{ELHI} : if \mathcal{T} is in \mathcal{ELHI} , \mathcal{ELH} or \mathcal{EL} , then Q' is a datalog query; if \mathcal{T} is in DL-Lite⁺ [13], then Q' consists of a union of conjunctive queries and a linear datalog query; and, if \mathcal{T} is in DL-Lite_R or DL-Lite_{core}, then Q' is a union of conjunctive queries. Hence, our approach not only handles the aforementioned sublanguages of \mathcal{ELHI} , but it is worst-case optimal for all of them, which makes it a generalization and extension of the approaches of [6] and [15].

An important difference between our algorithm and the CGLLR algorithm is the way in which existential quantification is handled. In our algorithm, functional terms are used to keep track of role successors. In contrast, the CGLLR algorithm deals with existential quantification by relying on a so-called reduction step, based on the standard notion of factoring [5]. Since our algorithm does not need to derive the queries produced by the reduction step, we conjecture that it will often produce smaller rewritings than the CGLLR algorithm. To test our conjecture we implemented both algorithms in their *original* form—that is, we did not implement any optimization to reduce the size of the rewritings, such as query containment [1], or forward/backward subsumption [5]. We then compared the two implementations using a benchmark suite containing different DL-Lite_R ontologies with TBoxes of varying complexity, and a set of test queries derived from applications that use these ontologies.

Our results provide insight into the strengths and weaknesses of the two approaches. They clearly indicate that the reduction step of the CGLLR algorithm can lead to a significant increase in the size of rewritings. The implementation of our algorithm—called REQUIEM (REsolution-based QUery rewriting for Expressive Models)—not only produced significantly smaller rewritings, but was also significantly faster in most cases.

2 Preliminaries

For P an *atomic role*, a *basic role* has the form P or P^- . For A and *atomic concept*, and R a *basic role*, an *antecedent concept* has the form A , or $\exists R$; and a *consequent concept* has the form A , $\exists R$, or $\exists R.A$. A DL-Lite_R TBox is a set of *inclusion assertions* or inclusion axioms of the form $A \sqsubseteq C$ or $R_1 \sqsubseteq R_2$, where A is an antecedent concept, C is a consequent concept, and R_1, R_2 are basic roles. The former inclusion assertion is called a *concept inclusion*, and the former, a *role inclusion*. An ABox is a set of *membership assertions* of the form $A(a)$ or $P(a, b)$, where A is an atomic concept, P is an atomic role, and a, b are constants.

¹ <http://www.dis.uniroma1.it/~quonto/>

² <http://pellet.owldl.com/owlgres/>

A DL-Lite_R *Knowledge Base* (KB) \mathcal{K} is a tuple $\langle \mathcal{T}, \mathcal{A} \rangle$, where \mathcal{T} is a DL-Lite_R TBox, and \mathcal{A} is an ABox. DL-Lite_{core} is obtained from DL-Lite_R by disallowing role inclusions. The semantics of a KB is defined as usual [3].

We use the well-known definitions of constants, variables, function symbols, terms, and atoms of first-order logic [8]. A *Horn clause* C is an expression of the form $D_0 \leftarrow D_1 \wedge \dots \wedge D_n$, where each D_i is an atom. The atom D_0 is called the *head*, and the set $\{D_1, \dots, D_n\}$ is called the *body*. Variables that occur in the body at most once and do not occur in the head are called *unbound* variables and may be denoted with the symbol $_$; all other variables are called *bound* variables. A Horn clause is *safe* if every variable occurring in the head also occurs in the body. A *conjunctive query* over a DL KB \mathcal{K} is a safe Horn clause whose head predicate does not occur in \mathcal{K} , and whose body is a set of atoms whose predicates are concept and role names occurring in \mathcal{K} . A *union of conjunctive queries* over \mathcal{K} is a set of conjunctive queries over \mathcal{K} with the same head up to variable renaming [4]. A tuple of constants \vec{a} is an *answer* to a union of conjunctive queries Q over \mathcal{K} if and only if $\mathcal{K} \cup Q \models Q_P(\vec{a})$, where Q_P is the head predicate of Q , and Q is considered to be a set of universally quantified implications with the usual first-order semantics [8]. The set of all answers or *certain answers* of Q over \mathcal{K} is denoted by $\text{ans}(Q, \mathcal{K})$. Given a conjunctive query Q and a TBox \mathcal{T} , a query Q' is said to be a *rewriting* of Q w.r.t. \mathcal{T} if $\text{ans}(Q, \mathcal{T} \cup \mathcal{A}) = \text{ans}(Q', \mathcal{A})$ for every ABox \mathcal{A} .

3 Query Rewriting Techniques

Axioms-as-rules Rewriting The CGLLR algorithm takes as input a conjunctive query Q and a DL-Lite_R TBox \mathcal{T} , and produces a union of conjunctive queries that is a rewriting of Q w.r.t. \mathcal{T} . The algorithm uses the axioms of \mathcal{T} as rewriting rules and applies them to the body atoms of Q . An axiom $\alpha \in \mathcal{T}$ is applicable to an atom $A(x)$ if α has A in its right-hand side; and α is applicable to an atom $P(x, y)$ if either (i) $y = _$ and the right-hand side of α is $\exists P$; (ii) $x = _$ and the right-hand side of α is $\exists P^-$; or (iii) α is a role inclusion and its right-hand side is either P or P^- . Let α be an axiom that is applicable to an atom D , we define the function $\text{ref}(D, \alpha)$ as follows:

- if $D = A(x)$, then we have that (i) if $\alpha = B \sqsubseteq A$, then $\text{ref}(D, \alpha) = B(x)$; (ii) if $\alpha = \exists P \sqsubseteq A$, then $\text{ref}(D, \alpha) = P(x, _)$; and (iii) if $\alpha = \exists P^- \sqsubseteq A$, then $\text{ref}(D, \alpha) = P(_, x)$.
- If $D = P(x, _)$, then we have that (i) if $\alpha = A \sqsubseteq \exists P$, then $\text{ref}(D, \alpha) = A(x)$; (ii) if $\alpha = \exists S \sqsubseteq \exists P$, then $\text{ref}(D, \alpha) = S(x, _)$; and (iii) if $\alpha = \exists S^- \sqsubseteq \exists P$, then $\text{ref}(D, \alpha) = S(_, x)$.
- If $D = P(_, x)$, then we have that (i) if $\alpha = A \sqsubseteq \exists P^-$, then $\text{ref}(D, \alpha) = A(x)$; (ii) if $\alpha = \exists S \sqsubseteq \exists P^-$, then $\text{ref}(D, \alpha) = S(x, _)$; and (iii) if $\alpha = \exists S^- \sqsubseteq \exists P^-$, then $\text{ref}(D, \alpha) = S(_, x)$.
- If $D = P(x, y)$, then we have that (i) if either $\alpha = S \sqsubseteq P$ or $\alpha = S^- \sqsubseteq P^-$, then $\text{ref}(D, \alpha) = S(x, y)$; and (ii) if either $\alpha = S \sqsubseteq P^-$ or $\alpha = S^- \sqsubseteq P$, then $\text{ref}(D, \alpha) = S(y, x)$.

```

Input: Conjunctive query  $Q$ , DL-LiteR TBox  $\mathcal{T}$ 
 $R = \{Q\}$ ;
repeat
  foreach query  $Q' \in R$  do
    (reformulation) foreach atom  $D$  in  $Q'$  do
      foreach axiom  $\alpha \in \mathcal{T}$  do
        if  $\alpha$  is applicable to  $D$  then
           $R = R \cup \{Q'[D/\text{ref}(D, \alpha)]\}$ ;
        end
      end
    (reduction) forall atoms  $D_1, D_2$  in  $Q'$  do
      if  $D_1$  and  $D_2$  unify then
         $\sigma = \text{MGU}(D_1, D_2)$ ;
         $R = R \cup \{\lambda(Q'\sigma)\}$ ;
      end
    end
  end
until no query unique up to variable renaming can be added to  $R$ ;
return  $R$ ;

```

Algorithm 1: The CGLLR algorithm

The CGLLR algorithm is shown in Algorithm 1. The expression $Q[D/D']$ denotes the conjunctive query obtained from Q by replacing the body atom D with a new atom D' . The function MGU takes as input two atoms and returns their most general unifier [4]. The function λ takes as input a conjunctive query Q and returns a new conjunctive query obtained by replacing each occurrence of an unbound variable in Q with the symbol $_.$. As can be seen, starting with the original query Q , the algorithm keeps producing queries in two steps until no other query unique up to variable renaming can be produced. In the *reformulation step* the algorithm rewrites the body atoms of a given query Q' by using applicable TBox axioms as rewriting rules, generating a new query for every atom reformulation. Then, in the *reduction step* the algorithm produces a new query $\lambda(Q'\sigma)$ for each pair of body atoms of Q' that unify. The second step is necessary to achieve completeness since an axiom that is not applicable to any body atom of Q' may be applicable to some body atom of $\lambda(Q'\sigma)$.

As a minor remark, we point out that the CGLLR algorithm does not handle qualified existential quantification natively—that is, there is no rewriting rule for axioms of the form $A \sqsubseteq \exists R.B$. Instead, w.l.o.g. every axiom of such a form occurring in \mathcal{T} is replaced with the set of axioms $\{A \sqsubseteq \exists P_1, \exists P_1^- \sqsubseteq B, P_1 \sqsubseteq R\}$, where P_1 is a new atomic role not occurring in \mathcal{T} .

Resolution-based Approach Our algorithm takes as input a conjunctive query Q and an \mathcal{ELHI} TBox \mathcal{T} , and produces a datalog query that is a rewriting of Q w.r.t. \mathcal{T} . We have shown that if \mathcal{T} is in DL-Lite_R, then the rewriting is a union of conjunctive queries. The algorithm first transforms Q and \mathcal{T} into

Input: Conjunctive query Q , DL-Lite_R TBox \mathcal{T}
 $R = \Xi(\mathcal{T}) \cup \{Q\}$;
repeat
 | (saturation) **forall** clauses C_1, C_2 in R **do**
 | | $R = R \cup \text{resolve}(C_1, C_2)$;
 | **end**
until no query unique up to variable renaming can be added to R ;
return $\{C \mid C \in \text{unfold}(\text{ff}(R)), \text{ and } C \text{ has the same head predicate as } Q\}$;
Algorithm 2: Our resolution-based algorithm

clauses, and then computes the rewriting by using a resolution-based calculus to derive new clauses from the initial set.

We present our algorithm in Algorithm 2. We omit the parts of the original algorithm that are irrelevant for DL-Lite_R. The expression $\Xi(\mathcal{T})$ denotes the set of clauses obtained from \mathcal{T} according to Table 1. The function `resolve` takes two clauses C_1 and C_2 , and returns a set containing every clause C_R that can be obtained by combining the atoms of C_1 and C_2 according to the *inference templates* shown in Table 2. A template of the form $\frac{P_1 \quad P_2}{R}$ denotes that given two clauses C_1, C_2 , if C_1 is of the form of P_1 , and C_2 is of the form of P_2 , then $\text{resolve}(C_1, C_2)$ contains all possible clauses of the form of R that can be constructed with C_1 and C_2 ; otherwise, $\text{resolve}(C_1, C_2) = \emptyset$. The function `ff` takes a set of clauses N and returns the set of function-free clauses of N . The function `unfold` takes a set of clauses N , and returns the set obtained by unfolding every clause in N ; for example, if $N = \{Q_P(x) \leftarrow A(x), A(x) \leftarrow B(x)\}$, then we have that $\text{unfold}(N) = N \cup \{Q_P(x) \leftarrow B(x)\}$, which results by unfolding $A(x) \leftarrow B(x)$ into $Q_P(x) \leftarrow A(x)$. A formal description of the unfolding step can be found in [13]. As can be seen, starting with the set of clauses $\Xi(\mathcal{T}) \cup \{Q\}$, the algorithm keeps producing clauses in the *saturation step* until no other clause unique up to variable renaming can be produced; the resulting function-free clauses are then unfolded; finally, every clause that does not have the same head predicate as Q is dropped.

As a minor remark, we point out that our algorithm does not handle axioms of the form $\exists R_1 \sqsubseteq \exists R_2$ natively. Instead, w.l.o.g. every axiom of such a form occurring in \mathcal{T} is replaced with the set of axioms $\{\exists R_1 \sqsubseteq A_1, A_1 \sqsubseteq \exists R_2\}$, where A_1 is a new atomic concept not occurring in \mathcal{T} .

The Main Difference The presented techniques mainly differ in their handling of *existential quantification*: while our algorithm deals with axioms containing existential quantifiers on the right-hand side by introducing functional terms, the CGLLR algorithm does so by restricting the applicability of such axioms, and relying on its reduction step. We explore this difference by means of an example (taken from [6]). Consider a DL-Lite_R TBox \mathcal{T} containing the following axioms:

$$\text{Professor} \sqsubseteq \exists \text{teaches}, \tag{1}$$

$$\exists \text{teaches}^- \sqsubseteq \text{Student}. \tag{2}$$

Table 1. Translating \mathcal{T} into a set of clauses $\Xi(\mathcal{T})$

DL-Lite _R clause	DL-Lite _R axiom
$B(x) \leftarrow A(x)$	$A \sqsubseteq B$
$P(x, f(x)) \leftarrow A(x)$	$A \sqsubseteq \exists P.B$
$B(f(x)) \leftarrow A(x)$	
$P(f(x), x) \leftarrow A(x)$	$A \sqsubseteq \exists P^-.B$
$B(f(x)) \leftarrow A(x)$	
$A(x) \leftarrow P(x, y)$	$\exists P \sqsubseteq A$
$A(x) \leftarrow P(y, x)$	$\exists P^- \sqsubseteq A$
$S(x, y) \leftarrow P(x, y)$	$P \sqsubseteq S, P^- \sqsubseteq S^-$
$S(x, y) \leftarrow P(y, x)$	$P \sqsubseteq S^-, P^- \sqsubseteq S$

Note 1. Each axiom of the form $A \sqsubseteq \exists R.B$ is uniquely associated with a function symbol f .

The TBox \mathcal{T} states that a professor teaches at least someone, and that someone that is taught is a student. Consider the query

$$Q(x) \leftarrow \text{teaches}(x, y) \wedge \text{Student}(y). \quad (3)$$

We first analyze the execution of the CGLLR algorithm. In the first iteration, the axiom (1) is not applicable to $\text{teaches}(x, y)$ because $\text{teaches}(x, y)$ has more than one bound variable. The reason why the applicability of (1) has to be restricted in this case is that the CGLLR algorithm does not keep track of information about role successors. In fact, naively allowing axioms of the form of (1) to be applicable in cases like this, would result in the loss of soundness. To illustrate this point, suppose that (1) were applicable to $\text{teaches}(x, y)$ and $\text{ref}(\text{teaches}(x, y), (1)) = \text{Professor}(x)$. Clearly, the algorithm would then obtain

$$Q(x) \leftarrow \text{Professor}(x) \wedge \text{Student}(y). \quad (4)$$

Note that the relation between x and y is lost—that is, the fact that the individual represented by y must be a teaches-successor of the individual represented by x is not captured by query (4). In particular, if \mathcal{T} were to contain axiom (1) only, then query (4) would clearly produce wrong answers.

Although the applicability of (1) is restricted, the axiom (2) is applicable to $\text{Student}(y)$ in (3), producing

$$Q(x) \leftarrow \text{teaches}(x, y) \wedge \text{teaches}(-, y). \quad (5)$$

In the next iteration, neither (1) nor (2) are applicable to any body atom of (5), so no query is added in the reformulation step. In the reduction step, however, the algorithm produces

$$Q(x) \leftarrow \text{teaches}(x, -), \quad (6)$$

Table 2. Inference templates for the function resolve

$$\frac{C(x) \leftarrow B(x) \quad B(f(x)) \leftarrow A(x)}{C(f(x)) \leftarrow A(x)}$$

$$\frac{B(x) \leftarrow P(x, y) \quad P(x, f(x)) \leftarrow A(x)}{B(x) \leftarrow A(x)} \quad \frac{B(x) \leftarrow P(x, y) \quad P(f(x), x) \leftarrow A(x)}{B(f(x)) \leftarrow A(x)}$$

$$\frac{B(x) \leftarrow P(y, x) \quad P(x, f(x)) \leftarrow A(x)}{B(f(x)) \leftarrow A(x)} \quad \frac{B(x) \leftarrow P(y, x) \quad P(f(x), x) \leftarrow A(x)}{B(x) \leftarrow A(x)}$$

$$\frac{S(x, y) \leftarrow P(x, y) \quad P(x, f(x)) \leftarrow A(x)}{S(x, f(x)) \leftarrow A(x)} \quad \frac{S(x, y) \leftarrow P(x, y) \quad P(f(x), x) \leftarrow A(x)}{S(f(x), x) \leftarrow A(x)}$$

$$\frac{S(x, y) \leftarrow P(y, x) \quad P(x, f(x)) \leftarrow A(x)}{S(f(x), x) \leftarrow A(x)} \quad \frac{S(x, y) \leftarrow P(y, x) \quad P(f(x), x) \leftarrow A(x)}{S(x, f(x)) \leftarrow A(x)}$$

$$\frac{Q_P(\vec{u}) \leftarrow B(t) \wedge \bigwedge D_i(\vec{t}_i) \quad B(f(x)) \leftarrow A(x)}{Q_P(\vec{u})\sigma \leftarrow A(x)\sigma \wedge \bigwedge D_i(\vec{t}_i)\sigma}$$

where $\sigma = \text{MGU}(B(t), B(f(x)))$, and $B(t)$ is deepest in its clause.

$$\frac{Q_P(\vec{u}) \leftarrow P(s, t) \wedge \bigwedge D_i(\vec{t}_i) \quad P(x, f(x)) \leftarrow A(x)}{Q_P(\vec{u})\sigma \leftarrow A(x)\sigma \wedge \bigwedge D_i(\vec{t}_i)\sigma}$$

where $\sigma = \text{MGU}(P(s, t), P(x, f(x)))$, and $P(s, t)$ is deepest in its clause.

$$\frac{Q_P(\vec{u}) \leftarrow P(s, t) \wedge \bigwedge D_i(\vec{t}_i) \quad P(f(x), x) \leftarrow A(x)}{Q_P(\vec{u})\sigma \leftarrow A(x)\sigma \wedge \bigwedge D_i(\vec{t}_i)\sigma}$$

where $\sigma = \text{MGU}(P(s, t), P(f(x), x))$, and $P(s, t)$ is deepest in its clause.

by unifying the body atoms of (5). In the following iteration, the axiom (1) is applicable to the only body atom of (6), producing

$$Q(x) \leftarrow \text{Professor}(x). \tag{7}$$

Note that without the reduction step, that caused query (6) to be produced, the algorithm would not have produced query (7). It can easily be verified that no more queries unique up to variable renaming can be produced; thus, the algorithm returns $\{(3), (5), (6), (7)\}$.

We now analyze the execution of our algorithm. The axioms (1) and (2) are translated into the following clauses:

$$\text{teaches}(x, f(x)) \leftarrow \text{Professor}(x), \tag{8}$$

$$\text{Student}(x) \leftarrow \text{teaches}(y, x). \tag{9}$$

In the saturation step the algorithm produces

$$\text{Student}(f(x)) \leftarrow \text{Professor}(x), \quad (\text{resolve}((8), (9))) \quad (10)$$

$$Q(x) \leftarrow \text{Professor}(x) \wedge \text{Student}(f(x)), \quad (\text{resolve}((3), (8))) \quad (11)$$

$$Q(x) \leftarrow \text{teaches}(x, f(x)) \wedge \text{Professor}(x), \quad (\text{resolve}((3), (10))) \quad (12)$$

$$Q(x) \leftarrow \text{Professor}(x). \quad (\text{resolve}((8), (12))) \quad (13)$$

Note the difference between queries (4) and (11). Since the function symbol f is uniquely associated with clause (8), unlike query (4), query (11) captures the fact that the individual represented by $f(x)$ must be a teaches-successor of the individual represented by x . It can easily be verified that no other clause is produced in the first step.

Clearly, $\text{ff}(R) = \{(3), (9), (13)\}$; therefore, we have that $\text{unfold}(\text{ff}(R))$ consists of $\text{ff}(R)$ and the query

$$Q(x) \leftarrow \text{teaches}(x, y) \wedge \text{teaches}(z, y), \quad (14)$$

which results from unfolding (9). Hence, since the clause (9) does not have the same head predicate as query (3), our algorithm returns $\{(3), (13), (14)\}$. Note that the queries (5) and (14), as well as (7) and (13), are equivalent up to variable renaming; therefore, the rewriting produced by our algorithm is contained in that produced by the CGLLR algorithm.

As can be seen in the example, the algorithms computed their rewritings in different ways. In particular, in the case of query (7)/(13), on the one hand, the CGLLR algorithm used its reduction step to derive (6) from (5); and then it derived (7) from (6) using its reformulation step; and on the other hand, it was possible for our algorithm to derive (13) from (8) and (12); (12) from (10) and (3); and (10) from (8) and (9), by syntactically keeping track of role successors using functional terms.

The use of functional terms is what mainly distinguishes it from the CGLLR algorithm. This distinctive feature makes our approach more goal-oriented, in the sense that it does not need to derive the queries produced by the reduction step of the CGLLR algorithm in order to be complete. Moreover, we have shown that every clause containing functional terms produced in the saturation step of our algorithm can be safely dropped. This is the main reason why we conjecture that our algorithm will often produce smaller rewritings than the CGLLR algorithm.

4 Evaluation

In order to test the conjecture that our algorithm produces smaller rewritings than the CGLLR algorithm, we implemented both algorithms. The implementation of our algorithm is called **REQUIEM** (REsolution-based QUery rewriting for Expressive Models). We refer to our implementation of the CGLLR algorithm as **C**. Both implementations are available at **REQUIEM**'s web site.³

³ <http://www.comlab.ox.ac.uk/projects/REQUIEM/>

As stated in the introduction, our objective is to compare the algorithms in their original forms; therefore, no optimization was implemented for either algorithm to reduce the size of the rewritings. The main goal of the evaluation is to compare the algorithms w.r.t. the size of the rewritings they produce. Since a rewriting containing fewer queries than another is not necessarily less complex, we consider the size of a rewriting as being the number of symbols needed to represent it. All tests were performed on a desktop computer with a 2.59 GHz Intel processor, 1.87 GB of RAM, running Windows XP. We used Java 1.6.0 Update 7 with a maximum heap size of 256 MB.

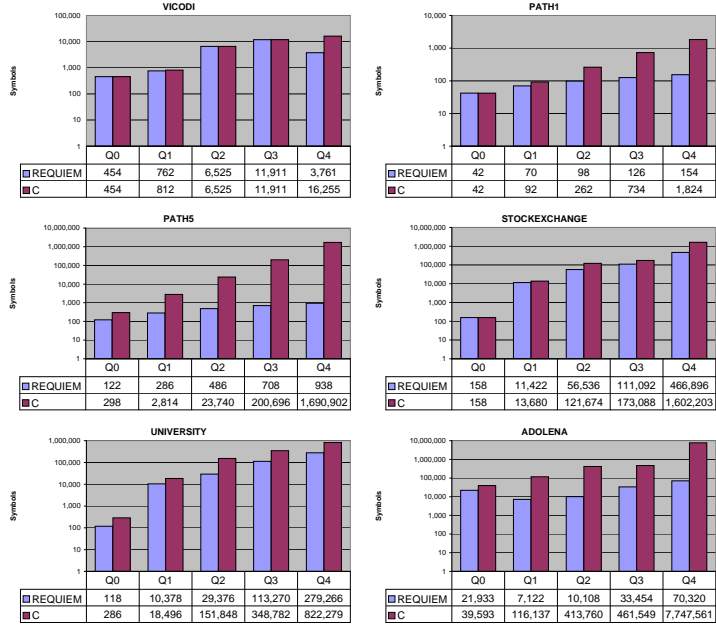
Test Ontologies and Queries We selected test DL-Lite_R ontologies with varying numbers of axioms containing existential quantification. We considered ontologies that were developed in the context of real projects. The queries we used were selected based on canonical examples used in the corresponding project (see Appendix A).

VICODI is an ontology developed in the homonymous EU-funded project,⁴ that captures information about European history that has been contextualized w.r.t. location, time and subject. STOCKEXCHANGE is an ontology developed for ontology-based data access [14], that represents the domain of financial institutions of the European Union. UNIVERSITY is a DL-Lite_R version of LUBM,⁵ a benchmark developed at Lehigh University for testing the performance of ontology management and reasoning systems; the ontology describes the organizational structure of universities. ADOLENA (Abilities and Disabilities OntoLogY for ENhancing Accessibility) is an ontology developed for ontology-based data access for the South African National Accessibility Portal [11]; the ontology describes abilities/disabilities, and devices. We decided to include two synthetic ontologies in our tests in order to have a controlled scenario to help us understand the impact of the reduction step of the CGLLR algorithm. PATH1 and PATH5 model information about graphs: nodes are represented by individuals, and vertices are ABox assertions of the form $\text{edge}(a, b)$. The ontology PATH5 contains concepts representing paths of length 1–5, while PATH1 contains a concept representing paths of length 1 only. All the ontologies and queries are available at REQUIEM’s web site.

Results The results of our tests are shown in Figures 1 and 2. Figure 1 shows the size of the rewritings, and Figure 2 shows the time it took to compute them. We decided to present the running time to give an indication of likely performance; we point out, however, that no special care was taken to obtain time efficient implementations. As can be seen in Figure 1, the rewritings produced by REQUIEM were never larger than those produced by C and were often significantly smaller. The most extreme case was the fifth query over ADOLENA, in which REQUIEM produced 624 queries, as opposed to the more than 78,500 produced by C. We compared the actual rewritings in order to verify whether for every query Q_R produced by REQUIEM, there was an equivalent query up to variable renaming Q_P produced by C. It turned out that this was the case for all

⁴ <http://www.vicodi.org/>

⁵ <http://swat.cse.lehigh.edu/projects/lubm/>

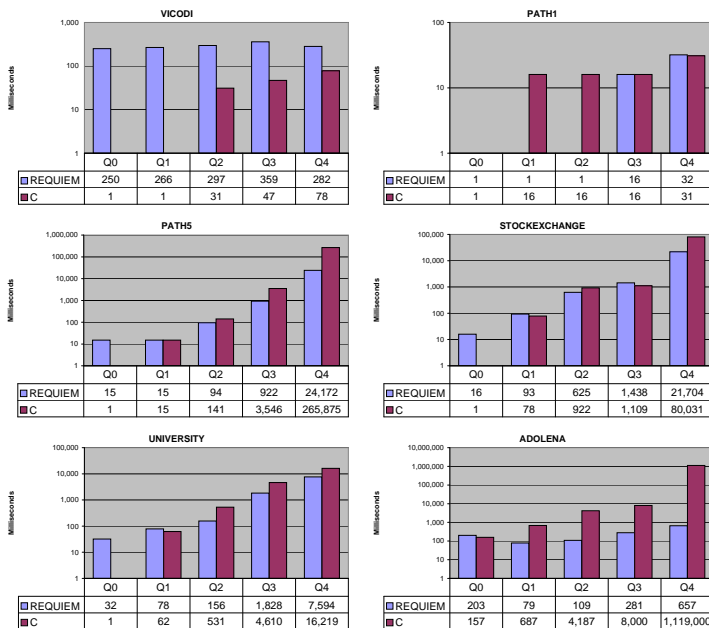
Fig. 1. Size of the rewritings

of our tests. A closer inspection of the rewritings produced by C revealed that most of them contained a considerable number of queries produced in the reduction step of the algorithm. This can be clearly seen in the tests involving PATH1 and PATH5. The queries we used with these two ontologies are of the form $Q_n(x_0) \leftarrow \text{edge}(x_0, x_1) \wedge \dots \wedge \text{edge}(x_n, x_{n+1})$, for $0 \leq n \leq 4$. Clearly, every pair of body atoms in queries of this form unify; therefore, for every Q' with m body atoms, C will produce $\binom{m}{2}$ new queries in the reduction step. The most extreme case for these ontologies was the fifth query over PATH5, in which REQUIEM produced 16 queries, as opposed to the more than 25,000 produced by C. Note that for every query Q'' produced from a query Q' in the reduction step, there is a substitution σ such that $Q'\sigma = Q''$. This implies that Q' subsumes Q'' [8]; therefore, every such query Q'' can be dropped from the rewriting without sacrificing completeness. Identifying such queries, however, is not straightforward since the CGLLR algorithm does not keep track of which queries were produced in the reduction step.

The results shown in Figure 2 suggest that REQUIEM will be significantly faster than C in case the queries are relatively complex and/or the ontologies contain a relatively large number of existential quantification axioms. The most extreme case was again the fifth query over ADOLENA, in which REQUIEM took a little over half a second, as opposed to the almost 20 minutes taken by C.

Finally, in order to determine the level of redundancy in the rewritings produced by our algorithm, we optimized the implementation of the algorithm in

Fig. 2. Running time



REQUIEM by using forward subsumption, and performing an a posteriori query containment check on the rewritings. We will refer to the optimized version as REQUIEM-Full. Note that, as explained earlier, in the case of the CGLLR algorithm for every query produced in the reduction step, there is a previously produced query that subsumes it; therefore, forward subsumption would get rid of every such query on the fly, which would compromise completeness. Moreover, given the size of the rewritings, a straightforward optimization based on an a posteriori query containment check may be impractical. Surprisingly, REQUIEM-Full produced the same rewritings as REQUIEM for 90% of the queries (but obviously took more time). The cases for which REQUIEM-Full produced fewer queries than REQUIEM were the first, second, and fourth queries over ADOLENA, where REQUIEM-Full obtained 27, 50, and 224, respectively; and REQUIEM produced 402, 103, and 492, respectively. This suggests that REQUIEM alone can be used effectively in many practical scenarios.

5 Future Work

We plan to develop optimization techniques for REQUIEM in order to reduce or eliminate recursion when possible. In order to evaluate the practicality of our query rewriting technique we plan to use REQUIEM in an ontology-based data access system using existing deductive database technology.

References

1. S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
2. F. Baader, S. Brandt, and C. Lutz. Pushing the EL Envelope. In *International Joint Conferences on Artificial Intelligence (IJCAI-05)*, 2005.
3. F. Baader and W. Nutt. *Basic Description Logics*, chapter 2, pages 47–100. Cambridge University Press, 2003.
4. F. Baader and W. Snyder. Unification Theory. In A. Robinson and A. Voronkov, editors, *Handbook of Automated Reasoning*, volume I, chapter 8, pages 445–532. Elsevier Science, 2001.
5. L. Bachmair and H. Ganzinger. Resolution Theorem Proving. In A. Robinson and A. Voronkov, editors, *Handbook of Automated Reasoning*, volume 1, chapter 2, pages 19–100. North Holland, 2001.
6. 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. of Automated Reasoning*, 2007.
7. D. Calvanese, G. D. Giacomo, M. Lenzerini, D. Nardi, and R. Rosati. Description Logic Framework for Information Integration. In *Principles of Knowledge Representation and Reasoning*, pages 2–13, 1998.
8. C.-L. Chang and R. C.-T. Lee. *Symbolic Logic and Mechanical Theorem Proving*. Academic Press, Inc., Orlando, FL, USA, 1997.
9. R. Fagin, P. G. Kolaitis, R. J. Miller, and L. Popa. Data Exchange: Semantics and Query Answering. In D. Calvanese, M. Lenzerini, and R. Motwani, editors, *ICDT*, volume 2572 of *Lecture Notes in Computer Science*, pages 207–224. Springer, 2003.
10. J. Hefflin and J. Hendler. A portrait of the semantic web in action. *IEEE Intelligent Systems*, 16(2):54–59, 2001.
11. C. M. Keet, R. Alberts, A. Gerber, and G. Chimamiwa. Enhancing web portals with ontology-based data access: The case study of south africa’s accessibility portal for people with disabilities. In *OWLED*, 2008.
12. M. Lenzerini. Data Integration: a theoretical perspective. In *PODS ’02: Proceedings of the twenty-first ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 233–246, New York, NY, USA, 2002. ACM Press.
13. H. Pérez-Urbina, B. Motik, and I. Horrocks. Rewriting Conjunctive Queries under Description Logic Constraints. In *Proceedings of the International Workshop on Logics in Databases*, May 2008.
14. M. Rodriguez-Muro, L. Lubyte, and D. Calvanese. Realizing ontology based data access: A plug-in for protg. In *Proc. of the Workshop on Information Integration Methods, Architectures, and Systems (IIMAS 2008)*, pages 286–289. IEEE Computer Society Press, 2008.
15. R. Rosati. On conjunctive query answering in EL. In *Proceedings of the 2007 International Workshop on Description Logics (DL2007)*, CEUR-WS, 2007.
16. R. van der Meyden. Logical Approaches to Incomplete Information: A Survey. In J. Chomicki and G. Saake, editors, *Logics for Databases and Information Systems*, pages 307–356. Kluwer, 1998.
17. J. Widom. Research Problems in Data Warehousing. In *4th International Conference on Information and Knowledge Management*, pages 25–30, Baltimore, Maryland, 1995.

A Queries

– VICODI

$$\begin{aligned}
Q_0(x_0) &\leftarrow \text{Location}(x_0) \\
Q_1(x_0, x_1) &\leftarrow \text{MilitaryPerson}(x_0) \wedge \text{hasRole}(x_1, x_0) \wedge \text{related}(x_0, x_2) \\
Q_2(x_0, x_1) &\leftarrow \text{TimeDependantRelation}(x_0) \wedge \text{hasRelationMember}(x_0, x_1) \wedge \\
&\quad \text{Event}(x_1) \\
Q_3(x_0, x_1) &\leftarrow \text{Object}(x_0) \wedge \text{hasRole}(x_0, x_1) \wedge \text{Symbol}(x_1) \\
Q_4(x_0) &\leftarrow \text{Individual}(x_0) \wedge \text{hasRole}(x_0, x_1) \wedge \text{Scientist}(x_1) \wedge \text{hasRole}(x_0, x_2) \wedge \\
&\quad \text{Discoverer}(x_2) \wedge \text{hasRole}(x_0, x_3) \wedge \text{Inventor}(x_3)
\end{aligned}$$

– PATH1/PATH5

$$Q_n(x_0) \leftarrow \text{edge}(x_0, x_1) \wedge \dots \wedge \text{edge}(x_n, x_{n+1}), \text{ for } 0 \leq n \leq 4.$$

– STOCKEXCHANGE

$$\begin{aligned}
Q_0(x_0) &\leftarrow \text{StockExchangeMember}(x_0) \\
Q_1(x_0, x_1) &\leftarrow \text{Person}(x_0) \wedge \text{hasStock}(x_0, x_1) \wedge \text{Stock}(x_1) \\
Q_2(x_0, x_1, x_2) &\leftarrow \text{FinancialInstrument}(x_0) \wedge \text{belongsToCompany}(x_0, x_1) \wedge \\
&\quad \text{Company}(x_1) \wedge \text{hasStock}(x_1, x_2) \wedge \text{Stock}(x_2) \\
Q_3(x_0, x_1, x_2) &\leftarrow \text{Person}(x_0) \wedge \text{hasStock}(x_0, x_1) \wedge \text{Stock}(x_1) \wedge \\
&\quad \text{isListedIn}(x_1, x_2) \wedge \text{StockExchangeList}(x_2) \\
Q_4(x_0, x_1, x_2, x_3) &\leftarrow \text{FinancialInstrument}(x_0) \wedge \text{belongsToCompany}(x_0, x_1) \wedge \\
&\quad \text{Company}(x_1) \wedge \text{hasStock}(x_1, x_2) \wedge \text{Stock}(x_2) \wedge \\
&\quad \text{isListedIn}(x_1, x_3) \wedge \text{StockExchangeList}(x_3)
\end{aligned}$$

– UNIVERSITY

$$\begin{aligned}
Q_0(x_0) &\leftarrow \text{worksFor}(x_0, x_1) \wedge \text{affiliatedOrganizationOf}(x_1, x_2) \\
Q_1(x_0, x_1) &\leftarrow \text{Person}(x_0) \wedge \text{teacherOf}(x_0, x_1) \wedge \text{Course}(x_1) \\
Q_2(x_0, x_1, x_2) &\leftarrow \text{Student}(x_0) \wedge \text{advisor}(x_0, x_1) \wedge \text{FacultyStaff}(x_1) \wedge \\
&\quad \text{takesCourse}(x_0, x_2) \wedge \text{teacherOf}(x_1, x_2) \wedge \text{Course}(x_2) \\
Q_3(x_0, x_1) &\leftarrow \text{Person}(x_0) \wedge \text{worksFor}(x_0, x_1) \wedge \text{Organization}(x_1) \\
Q_4(x_0) &\leftarrow \text{Person}(x_0) \wedge \text{worksFor}(x_0, x_1) \wedge \text{University}(x_1) \wedge \\
&\quad \text{hasAlumnus}(x_1, x_0)
\end{aligned}$$

– ADOLENA

$$\begin{aligned}
Q_0(x_0) &\leftarrow \text{Device}(x_0) \wedge \text{assistsWith}(x_0, x_1) \\
Q_1(x_0) &\leftarrow \text{Device}(x_0) \wedge \text{assistsWith}(x_0, x_1) \wedge \text{UpperLimbMobility}(x_1) \\
Q_2(x_0) &\leftarrow \text{Device}(x_0) \wedge \text{assistsWith}(x_0, x_1) \wedge \text{Hear}(x_1) \wedge \text{affects}(x_2, x_1) \wedge \\
&\quad \text{Autism}(x_2) \\
Q_3(x_0) &\leftarrow \text{Device}(x_0) \wedge \text{assistsWith}(x_0, x_1) \wedge \text{PhysicalAbility}(x_1) \\
Q_4(x_0) &\leftarrow \text{Device}(x_0) \wedge \text{assistsWith}(x_0, x_1) \wedge \text{PhysicalAbility}(x_1) \wedge \\
&\quad \text{affects}(x_2, x_1) \wedge \text{Quadriplegia}(x_2)
\end{aligned}$$