

Conjunctive Query Answering in the Description Logic \mathcal{EL} using a Relational Database System

Abstract

Conjunctive queries (CQ) are a fundamental mechanism for accessing description logic (DL) knowledge bases. We study CQ answering in (extensions of) the DL \mathcal{EL} , a popular choice for the formulation of large-scale ontologies that underlies the designated OWL2-EL profile of OWL2. Our main contribution is a novel approach to CQ answering that enables the use of conventional relational database systems as the basis for query execution. We evaluate our approach using the IBM DB2 system, with encouraging results.

1 Introduction

One of the main applications of ontologies in computer science is accessing data: an ontology formalizes conceptual information about data stored in one or multiple data sources. This information is then used to derive answers to queries over such sources. This general setup plays a central role in, e.g., querying the Semantic Web, ontology-based information integration, and peer-to-peer data management. In all these applications Description Logics (DLs) are popular ontology languages and conjunctive queries (CQs) are used as a fundamental querying tool. Hence, efficient and scalable approaches to CQ answering over DL ontologies are of vital importance.

Calvanese et al. have argued that, in the short run, true scalability of conjunctive query answering over DL ontologies can only be achieved by making use of standard relational database management systems (RDBMSs) [Calvanese et al., 2007b]. Alas, this is not a simple task as RDBMSs are unaware of TBoxes (the DL mechanism for storing conceptual information) and adopt the closed-world semantics. In contrast ABoxes (the DL mechanism for storing data) and the associated ontologies employ the open-world semantics. Existing approaches to overcome these differences have serious limitations: The approach of [Calvanese et al., 2007b], for example, applies only to DLs with data complexity of CQ answering in LOGSPACE; for many DLs, however, this problem is complete for PTIME or co-NP. In particular, the above approach cannot be directly used for DLs that include qualified existential restrictions, which play an important role in many

ontologies. This limitation is also shared by the rule-based approach presented in [Wu et al., 2008].

In this paper, we present a novel approach to using RDBMSs for CQ answering over DL ontologies that, in particular, accommodates qualified existential restrictions. We apply it to the \mathcal{EL} family of DLs, whose members are widely used as ontology languages for large-scale biomedical ontologies such as SNOMED CT and (early versions of) NCI [Spackman, 2000; Sioutos et al., 2006]. Our main result shows that CQ answering in $\mathcal{ELH}_{\perp}^{dr}$, the extension of basic \mathcal{EL} with the bottom concept, role inclusions, and domain and range restrictions, can be realized using an RDBMS. This result is of particular relevance as $\mathcal{ELH}_{\perp}^{dr}$ can be viewed as the core of the designated OWL-EL profile of the upcoming OWL Version 2 ontology language. We evaluate our approach using the IBM DB2 RDBMS and show that it scales to TBoxes with more than 50.000 axioms yielding answer times ranging from a fraction of a second to a few seconds.

The central idea of our approach is to incorporate the consequences of the TBox \mathcal{T} into the relational instance corresponding to the given ABox \mathcal{A} . Formally, we say that a DL admits *extended FO reducibility* if it is possible to effectively rewrite (i) \mathcal{A} and \mathcal{T} into an FO structure (independently of q) and (ii) q and (possibly) \mathcal{T} into an FO query q^* (independently of \mathcal{A}) such that answers are preserved, i.e., the answer to q^* over the above FO structure under the closed-world semantics is the same as the answer to q in $\mathcal{K} = (\mathcal{T}, \mathcal{A})$. The connection to RDBMSs then relies on a well-known equivalence between FO structures and relational databases, and FO queries and SQL queries. The notion of *extended FO reducibility* generalizes the notion of *first-order (FO) reducibility* [Calvanese et al., 2007b] where the TBox is incorporated into the query q rather than into the ABox \mathcal{A} while the ABox itself is used as a relational instance without any modification. Our approach has two major benefits:

1. It applies to DLs for which data complexity of CQ answering is PTIME-complete, such as $\mathcal{ELH}_{\perp}^{dr}$.
2. For $\mathcal{ELH}_{\perp}^{dr}$, both rewriting steps can be carried out in polynomial time and produce only a polynomial blowup.

Moreover, the query rewriting only depends on the input CQ and the role inclusions in \mathcal{T} (which are usually only few), but not on \mathcal{T} 's concept inclusions (usually very many). In contrast and to the best of our knowledge all existing approaches

used to define non-extended FO reducibility incur a blowup during query rewriting that is exponential in the number of concept inclusions.

In addition, we analyze the limitations of our approach and show that DLs with data complexity exceeding PTIME cannot enjoy *polynomial* extended FO reducibility, where the expansion of the ABox due to rewriting (but not necessarily of query rewriting) is bounded by a polynomial. We believe that this is a significant result as extended FO reducibility involving an exponential blowup of the data does not seem to be practical. In particular, the result implies that expressive DLs such as \mathcal{ALC} cannot enjoy polynomial FO reducibility. The same holds even for \mathcal{EL} enriched with negated ABox assertions and negated query atoms. For the latter case, we sketch an approach to query answering that involves only a polynomial blowup of the ABox, but is incomplete (in a way precisely characterized by an alternative semantics).

Due to space limitations, detailed proofs are given in the technical report [Anonymous, 2009].

2 Preliminaries

In \mathcal{ELH}_\perp^{dr} , concepts are built according to the syntax rule

$$C ::= A \mid \top \mid \perp \mid C \sqcap D \mid \exists r.C$$

where, here and in the remaining paper, A ranges over *concept names* taken from a countably infinite set N_C , r ranges over *role names* taken from a countably infinite set N_R , and C, D range over concepts. A *TBox* is a finite set of *concept inclusions* $C \sqsubseteq D$, *role inclusions* $r \sqsubseteq s$, *domain restrictions* $\text{dom}(r) \sqsubseteq C$, and *range restrictions* $\text{ran}(r) \sqsubseteq C$. An *ABox* is a finite set of *concept assertions* $A(a)$ and *role assertions* $r(a, b)$, where a, b range over a countably infinite set N_I of *individual names*. A *knowledge base* is a pair $(\mathcal{T}, \mathcal{A})$ with \mathcal{T} a TBox and \mathcal{A} an ABox.

The semantics of concepts, inclusions, restrictions, and assertions is as given as usual in terms of an *interpretation* $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$, we refer to [Baader *et al.*, 2008] for details. An interpretation \mathcal{I} is a *model* of a TBox \mathcal{T} (ABox \mathcal{A}) if it satisfies all inclusions and restrictions in \mathcal{T} (assertions in \mathcal{A}). It is a model of a knowledge base $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ if it is a model of \mathcal{T} and \mathcal{A} . A knowledge base that has a model is called *consistent*. For a concept inclusion, role inclusion, or assertion α , we write $\mathcal{K} \models \alpha$ if α is satisfied in all models of \mathcal{K} .

Let N_V be a countably infinite set of *variables*. Together, the sets N_V (of variables) and N_I (of individual names) form the set N_T of *terms*. A *first-order query* is a first-order formula $\varphi(\vec{v})$ built from N_T and the unary and binary predicates from N_C and N_R with $\vec{v} = v_1, \dots, v_k$ a vector of variables in N_V containing all free variables of φ . Variables in \vec{v} are called the *answer variables* of q . We call q *k-ary* if there are k answer variables. A *conjunctive query* is a first-order query q of the form $\exists \vec{u}.\psi(\vec{u}, \vec{v})$, where ψ is a conjunction of *concept atoms* $A(t)$ and *role atoms* $r(t, t')$. The variables in \vec{u} are called *quantified variables* of q . We use $\text{var}(q)$ to denote the set of all variables in \vec{u} and \vec{v} , $\text{qvar}(q)$ for the set of quantified variables, $\text{avar}(q)$ for the set of answer variables, and $\text{term}(q)$ for the terms in q . Slightly abusing notation, we write $\alpha \in q$ if the concept or role atom α occurs in q .

For \mathcal{I} an interpretation, $q = \varphi(\vec{v})$ a k -ary first-order query, and $a_1, \dots, a_k \in N_I$, we write $\mathcal{I} \models q[a_1, \dots, a_k]$ if \mathcal{I} satisfies q with v_i assigned to $a_i^{\mathcal{I}}$, $1 \leq i \leq k$. A *certain answer* for a k -ary conjunctive query q and a knowledge base \mathcal{K} is a tuple (a_1, \dots, a_k) such that a_1, \dots, a_k occur in \mathcal{K} and $\mathcal{I} \models q[a_1, \dots, a_k]$ for each model \mathcal{I} of \mathcal{K} . We use $\text{cert}(q, \mathcal{K})$ to denote the set of all certain answers for q and \mathcal{K} . This defines the querying problem studied in this paper: to compute $\text{cert}(q, \mathcal{K})$ given an \mathcal{ELH}_\perp^{dr} knowledge base \mathcal{K} and a CQ q .

Note that CQ answering generalizes *instance checking*, the problem of deciding whether $\mathcal{K} \models C(a)$, for C a \mathcal{ELH}_\perp^{dr} -concept: any such C can be easily unfolded into a conjunctive query q_C such that $a \in \text{cert}(q_C, \mathcal{K})$ iff $\mathcal{K} \models C(a)$.

In the remainder of this paper we use the *unique name assumption*: $a^{\mathcal{I}} \neq b^{\mathcal{I}}$ for all interpretations \mathcal{I} and all $a, b \in N_I$ with $a \neq b$. It is not hard to see that this has no impact on certain answers. We also assume, w.l.o.g., that (i) queries contain only individual names that occur in the KB against which they are asked, (ii) TBoxes do not contain domain restrictions, (iii) TBoxes contain exactly one range restriction per role name, (iv) if $\mathcal{K} \models r \sqsubseteq s$ and $\text{ran}(r) \sqsubseteq C$, $\text{ran}(s) \sqsubseteq D$ are in \mathcal{T} , then $C \sqsubseteq_{\mathcal{T}} D$, and (v) there are no $r, s \in N_R$ with $r \neq s$, $\mathcal{K} \models r \sqsubseteq s$, and $\mathcal{K} \models s \sqsubseteq r$. Justifications for these assumptions are given in [Anonymous, 2009].

3 ABox Rewriting / Canonical Models

The rewriting of the ABox consists of an extension of the ABox to a canonical model of the knowledge base. For the remainder of this section we fix an \mathcal{ELH}_\perp^{dr} -KB $\mathcal{K} = (\mathcal{T}, \mathcal{A})$. We use $\text{sub}(\mathcal{T})$ to denote the set of all subconcepts of concepts that occur in \mathcal{T} , $\text{rol}(\mathcal{T})$ for the set of role names that occur in \mathcal{T} , and $\text{Ind}(\mathcal{A})$ for the set of individual names that occur in \mathcal{A} . We also use $\text{ran}_{\mathcal{T}}(r)$ to denote the (unique) concept C with $\text{ran}(r) \sqsubseteq C \in \mathcal{T}$, and set

$$\begin{aligned} \text{ran}(\mathcal{T}) &:= \{\text{ran}_{\mathcal{T}}(r) \mid r \in \text{rol}(\mathcal{T})\} \\ N_{\text{aux}} &:= \{x_{C,D} \mid C \in \text{ran}(\mathcal{T}) \text{ and } D \in \text{sub}(\mathcal{T})\}, \end{aligned}$$

assuming $N_I \cap N_{\text{aux}} = \emptyset$. The canonical model $\mathcal{I}_{\mathcal{K}}$ of \mathcal{K} is defined in Figure 1. It is standard to show the following (similar to proofs in [Baader *et al.*, 2005b; Lutz and Wolter, 2007]):

Proposition 1. *If \mathcal{K} is consistent, then $\mathcal{I}_{\mathcal{K}}$ is a model of \mathcal{K} .*

Note that the cardinality of $\Delta^{\mathcal{I}_{\mathcal{K}}}$ is only quadratic in the size of \mathcal{K} , and linear if \mathcal{T} does not contain range restrictions. The model $\mathcal{I}_{\mathcal{K}}$ can be computed in polynomial time since subsumption and instance checking in \mathcal{ELH}_\perp^{dr} can be decided in poly-time [Baader *et al.*, 2005b]. In fact, there is an easy rule-based procedure for computing the canonical model and the rules can be implemented as standard database operations. Consistency of \mathcal{K} can also be checked in polynomial time. $\mathcal{I}_{\mathcal{K}}$ can be used for instance checking: $\mathcal{K} \models C(a)$ iff $\mathcal{I}_{\mathcal{K}} \models C(a)$, for \mathcal{K} consistent, C an \mathcal{ELH}_\perp^{dr} -concept, and $a \in \text{Ind}(\mathcal{A})$ (this can be shown by modifying similar proofs in [Baader *et al.*, 2005b]). Unfortunately, an analogous statement for conjunctive queries, namely $(a_1, \dots, a_k) \in \text{cert}(q, \mathcal{K})$ iff $\mathcal{I}_{\mathcal{K}} \models q[a_1, \dots, a_k]$, does not hold for two reasons. The first reason is that $\mathcal{I}_{\mathcal{K}}$ may contain unnecessary elements:

$$\begin{aligned}
\Delta^{\mathcal{I}_{\mathcal{K}}} &:= \text{Ind}(\mathcal{A}) \uplus \text{NI}_{\text{aux}} \text{ and } a^{\mathcal{I}_{\mathcal{K}}} := a \text{ for all } a \in \text{Ind}(\mathcal{K}) \\
A^{\mathcal{I}_{\mathcal{K}}} &:= \{a \in \text{ind}(\mathcal{A}) \mid \mathcal{K} \models A(a)\} \cup \{x_{C,D} \in \text{NI}_{\text{aux}} \mid \mathcal{K} \models C \sqcap D \sqsubseteq A\} \\
r^{\mathcal{I}_{\mathcal{K}}} &:= \{(a,b) \in \text{ind}(\mathcal{A}) \times \text{ind}(\mathcal{A}) \mid s(a,b) \in \mathcal{A} \text{ and } \mathcal{K} \models s \sqsubseteq r\} \cup \\
&\quad \{(a, x_{C,D}) \in \text{ind}(\mathcal{A}) \times \text{NI}_{\text{aux}} \mid \mathcal{K} \models \exists s.D(a), \text{ran}_{\mathcal{T}}(s) = C, \text{ and } \mathcal{K} \models s \sqsubseteq r\} \cup \\
&\quad \{(x_{C,D}, x_{C',D'}) \in \text{NI}_{\text{aux}} \times \text{NI}_{\text{aux}} \mid \mathcal{K} \models C \sqcap D \sqsubseteq \exists r.D', \text{ran}_{\mathcal{T}}(s) = C', \text{ and } \mathcal{K} \models s \sqsubseteq r\}
\end{aligned}$$

Figure 1: The canonical model $\mathcal{I}_{\mathcal{K}}$.

Example 2. Take $\mathcal{K}_1 = (\mathcal{T}_1, \mathcal{A}_1)$ with $\mathcal{T}_1 = \{A \sqsubseteq A\}$ and $\mathcal{A}_1 = \{B(a)\}$, and $q_1 = \exists u.(B(v) \wedge A(u))$. Then $x_{\top, A} \in A^{\mathcal{I}_{\mathcal{K}_1}}$ and so $\mathcal{I}_{\mathcal{K}_1} \models q_1[a]$, but clearly $a \notin \text{cert}(q_1, \mathcal{K}_1)$.

This deficiency of $\mathcal{I}_{\mathcal{K}}$ is easily repaired by restricting it to elements reachable from some $a^{\mathcal{I}_{\mathcal{K}}}$ with $a \in \text{Ind}(\mathcal{A})$. Formally, we define $\text{Ind}(\mathcal{A})^{\mathcal{I}} = \{a^{\mathcal{I}} \mid a \in \text{Ind}(\mathcal{A})\}$ for an interpretation \mathcal{I} . A path in \mathcal{I} is a finite sequence $d_0 r_1 d_1 \cdots r_n d_n$, $n \geq 0$, where $d_0 \in \text{Ind}(\mathcal{A})^{\mathcal{I}}$ and $(d_i, d_{i+1}) \in r_{i+1}^{\mathcal{I}}$ for all $i < n$ (note that the paths depend on the interpretation of $\text{Ind}(\mathcal{A})$ in \mathcal{I}). We denote by $\text{paths}_{\mathcal{A}}(\mathcal{I})$ the set of all paths in \mathcal{I} and by $\text{tail}(p)$ the last element d_n in p whenever $p \in \text{paths}_{\mathcal{A}}(\mathcal{I})$.

Now let $\mathcal{I}_{\mathcal{K}}^r$ denote the restriction of $\mathcal{I}_{\mathcal{K}}$ to those $d \in \Delta^{\mathcal{I}_{\mathcal{K}}}$ for which there exists a $p \in \text{paths}_{\mathcal{A}}(\mathcal{I}_{\mathcal{K}})$ such that $d = \text{tail}(p)$. Then $\mathcal{I}_{\mathcal{K}}^r$ provides the correct certain answers to the query q_1 from Example 2. A much more severe deficiency of $\mathcal{I}_{\mathcal{K}}$ (and $\mathcal{I}_{\mathcal{K}}^r$) is exemplified by the following:

Example 3. Take $\mathcal{K}_2 = (\mathcal{T}_2, \mathcal{A}_2)$ with $\mathcal{T}_2 = \{A \sqsubseteq \exists r.B \sqcap \exists s.B\}$ and $\mathcal{A}_2 = \{A(a)\}$, and $q_2 = \exists u.(r(v, u) \wedge s(v, u))$. Then $(a, x_{\top, B}) \in r^{\mathcal{I}_{\mathcal{K}_2}}$ and $(a, x_{\top, B}) \in s^{\mathcal{I}_{\mathcal{K}_2}}$ and therefore $\mathcal{I}_{\mathcal{K}_2} \models q_2[a]$, but clearly $a \notin \text{cert}(q_2, \mathcal{K}_2)$.

In principle, this problem can be overcome by replacing $\mathcal{I}_{\mathcal{K}}^r$ with its unraveling into a less constrained, tree-like model. In the following, we introduce unraveling as a general operation on models. Let \mathcal{R} be the set of role inclusions in \mathcal{K} . The $(\mathcal{A}, \mathcal{R})$ -unraveling \mathcal{J} of \mathcal{I} is defined as follows:

$$\begin{aligned}
\Delta^{\mathcal{J}} &:= \text{paths}_{\mathcal{A}}(\mathcal{I}) \text{ and } a^{\mathcal{J}} := a^{\mathcal{I}} \text{ for all } a \in \text{Ind}(\mathcal{K}) \\
A^{\mathcal{J}} &:= \{p \mid \text{tail}(p) \in A^{\mathcal{I}}\} \\
r^{\mathcal{J}} &:= \{(d, e) \mid d, e \in \text{Ind}^{\mathcal{I}} \wedge (d, e) \in r^{\mathcal{I}}\} \cup \\
&\quad \{(p, p \cdot se) \mid p, p \cdot se \in \Delta^{\mathcal{J}} \text{ and } \mathcal{R} \models s \sqsubseteq r\}
\end{aligned}$$

where “ \cdot ” denotes concatenation. Denote by $\mathcal{U}_{\mathcal{K}}$ the $(\mathcal{A}, \mathcal{R})$ -unraveling of $\mathcal{I}_{\mathcal{K}}^r$. Notice that the construction of $\mathcal{U}_{\mathcal{K}}$ from $\mathcal{I}_{\mathcal{K}}^r$ does not depend on the concept inclusions in \mathcal{T} , but only on \mathcal{R} . The following result is proved similarly to the analogous result for the DL \mathcal{ELI}^f in [Krisnadhi and Lutz, 2007].

Proposition 4. *If \mathcal{K} is consistent, then for all k -ary conjunctive queries q and all $a_1, \dots, a_k \in \text{Ind}(\mathcal{A})$, we have $(a_1, \dots, a_k) \in \text{cert}(q, \mathcal{K})$ iff $\mathcal{U}_{\mathcal{K}} \models q[a_1, \dots, a_k]$.*

By Proposition 4, $\mathcal{U}_{\mathcal{K}}$ gives the correct answers to conjunctive queries, but in contrast to $\mathcal{I}_{\mathcal{K}}^r$, it is typically infinite. Thus, we do not use it as a target for ABox rewriting and work with $\mathcal{I}_{\mathcal{K}}^r$. To overcome the problem indicated in Example 3, we use query rewriting.

4 Query Rewriting

Our aim is to rewrite the original CQ q into an FO query $q_{\mathcal{R}}^*$ such that $\mathcal{U}_{\mathcal{K}} \models q[a_1, \dots, a_k]$ iff $\mathcal{I}_{\mathcal{K}}^r \models q_{\mathcal{R}}^*[a_1, \dots, a_k]$ for all

$a_1, \dots, a_k \in \text{Ind}(\mathcal{A})$. By Proposition 4, we obtain the correct answers by using $\mathcal{I}_{\mathcal{K}}^r$ (the rewriting of \mathcal{A}) as a relational database instance and replacing q with $q_{\mathcal{R}}^*$. We now formulate the main result of this paper.

Theorem 5. *For every finite set of role inclusions \mathcal{R} and k -ary CQ q , one can construct in polynomial time a k -ary FO query $q_{\mathcal{R}}^*$ such that for all $\mathcal{ELH}_{\perp}^{dr}$ -KBS $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ with \mathcal{R} the set of role inclusions in \mathcal{T} and all $a_1, \dots, a_k \in \text{Ind}(\mathcal{A})$, the following holds:*

$$\mathcal{I}_{\mathcal{K}}^r \models q_{\mathcal{R}}^*[a_1, \dots, a_k] \text{ iff } \mathcal{U}_{\mathcal{K}} \models q[a_1, \dots, a_k].$$

In the remainder of this section, we show how to construct $q_{\mathcal{R}}^*$. The query $q_{\mathcal{R}}^*$ contains one additional unary predicate $\text{Aux}(x)$; we assume that Aux is always interpreted as $\Delta^{\mathcal{I}_{\mathcal{K}}^r} \setminus \text{Ind}(\mathcal{A})^{\mathcal{I}_{\mathcal{K}}^r}$ in $\mathcal{I}_{\mathcal{K}}^r$. Fix a finite set \mathcal{R} of role inclusions and a k -ary conjunctive query q . To construct $q_{\mathcal{R}}^*$, we use several auxiliary definitions. Let \sim_q denote the smallest relation on $\text{term}(q)$ that includes the identity relation, is transitive, and satisfies the following closure condition:

$$(*) \text{ if } r_1(s, t), r_2(s', t') \in q \text{ with } t \sim_q t', \text{ then } s \sim_q s'.$$

The relation \sim_q is central to our rewriting procedure. To see this, let $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ such that the role inclusions of \mathcal{T} coincide with \mathcal{R} . Intuitively, $\mathcal{U}_{\mathcal{K}}$ is produced from $\mathcal{I}_{\mathcal{K}}^r$ by keeping the $\text{Ind}(\mathcal{A})$ -part intact and relaxing the Aux -part into a collection of trees. To understand $(*)$, first assume $t = t'$. Then $(*)$ describes a non-tree situation in the query since $t = t'$ has two predecessors s and s' . Therefore, any match of the query in $\mathcal{I}_{\mathcal{K}}^r$ that maps t to the Aux -part should map s and s' to the same element; otherwise such a match cannot be reproduced in $\mathcal{U}_{\mathcal{K}}$. The case where $t \sim_q t'$ instead of $t = t'$ can be understood inductively.

It is not hard to verify that \sim_q is an equivalence relation and can be computed in time polynomial in the size of q . For $t \in \text{term}(q)$, we use $[t]$ to denote the equivalence class of t w.r.t. \sim_q and define, for any equivalence class ζ of \sim_q , the sets:

$$\begin{aligned}
\text{pre}(\zeta) &:= \{t \mid r(t, t') \in q \text{ for some } r \in \text{NR} \text{ and } t' \in \zeta\}, \text{ and} \\
\text{in}(\zeta) &:= \{r \mid r(t, t') \in q \text{ for some } t \in \text{term}(q) \text{ and } t' \in \zeta\}.
\end{aligned}$$

For $R \subseteq \text{NR}$ and $r \in \text{NR}$, r is called an *implicant* of R if $\mathcal{R} \models r \sqsubseteq s$ for all $s \in R$. It is called a *prime implicant* if, additionally, $\mathcal{R} \not\models r \sqsubseteq r'$ for all implicants r' of R with $r \neq r'$. By assumption (v) in Section 2, there is a prime implicant for any set $R \subseteq \text{NR}$ for which there is an implicant. We define:

- Fork_{\sim} is the set of pairs $(\text{pre}(\zeta), \zeta)$ with $\text{pre}(\zeta)$ of cardinality at least two;

- Fork_{\neq} is the set of variables $v \in \text{qvar}(q)$ such that there is no implicant of $\text{in}(\{v\})$;
- $\text{Fork}_{\mathcal{H}}$ is the set of pairs (I, ζ) such that $\text{pre}(\zeta) \neq \emptyset$, there is a prime implicant of $\text{in}(\zeta)$ that is not contained in $\text{in}(\zeta)$, and I is the set of all prime implicants of $\text{in}(\zeta)$;
- Cyc is the set of variables $v \in \text{qvar}(q)$ such that there are $r_0(t_0, t'_0), \dots, r_m(t_m, t'_m), \dots, r_n(t_n, t'_n) \in q$, $n, m \geq 0$, with $v \sim_q t_i$ for some $i \leq n$, $t'_i \sim_q t_{i+1}$ for all $i < n$, and $t'_n \sim_q t_m$.

It is not hard to see that $\text{Fork}_{=}$, Fork_{\neq} , $\text{Fork}_{\mathcal{H}}$, and Cyc can be computed in time polynomial in the size of q . For each equivalence class ζ of \sim_q , choose a representative $t_\zeta \in \zeta$ and if $\text{pre}(\zeta) \neq \emptyset$, choose a $t_\zeta^{\text{pre}} \in \text{pre}(\zeta)$. For $q = \exists \vec{u}. \psi$, the rewritten query $q_{\mathcal{R}}^*$ is now defined as $\exists \vec{u}. (\psi \wedge \varphi_1 \wedge \varphi_2 \wedge \varphi_3)$, where φ_1 , φ_2 , and φ_3 are as follows:

$$\begin{aligned} \varphi_1 &:= \bigwedge_{v \in \text{avar}(q) \cup \text{Fork}_{\neq} \cup \text{Cyc}} \neg \text{Aux}(v) \\ \varphi_2 &:= \bigwedge_{(\{t_1, \dots, t_k\}, \zeta) \in \text{Fork}_{=}} (\text{Aux}(t_\zeta) \rightarrow \bigwedge_{1 \leq i < k} t_i = t_{i+1}) \\ \varphi_3 &:= \bigwedge_{(I, \zeta) \in \text{Fork}_{\mathcal{H}}} (\text{Aux}(t_\zeta) \rightarrow \bigvee_{r \in I} r(t_\zeta^{\text{pre}}, t_\zeta)) \end{aligned}$$

We can show that $q_{\mathcal{R}}^*$ is as required [Anonymous, 2009]. The following examples illustrate the definition of $q_{\mathcal{R}}^*$. Note that φ_3 is trivial for $\mathcal{R} = \emptyset$.

(1) Let $\mathcal{R} = \emptyset$ and consider $q = \exists u. \psi$ with $\psi = r(v, u) \wedge r(v', u)$. This query illustrates the role of $\text{Fork}_{=}$. \sim_q consists of the equivalence classes $\{v, v'\}$ and $\{u\}$. We have $\text{pre}(\{u\}) = \{v, v'\}$ and $\text{in}(\{u\}) = \{r\}$. Hence $\text{Fork}_{=} = \{(\{v, v'\}, \{u\})\}$ and $\text{Fork}_{\neq} = \text{Fork}_{\mathcal{H}} = \text{Cyc} = \emptyset$. We obtain $q_{\mathcal{R}}^* = \exists u. (\psi \wedge \neg \text{Aux}(v) \wedge \neg \text{Aux}(v') \wedge (\text{Aux}(u) \rightarrow v = v'))$.

(2) Let $\mathcal{R} = \emptyset$ and consider $q = \exists u. (r(v, u) \wedge s(u, u))$. This query illustrates the role of Cyc . We have $u \in \text{Cyc}$ and so

$$q_{\mathcal{R}}^* = \exists u. (r(v, u) \wedge s(u, u) \wedge \neg \text{Aux}(v) \wedge \neg \text{Aux}(u)).$$

(3) Let $\mathcal{R} = \emptyset$ and consider $q = \exists u. (r(v, u) \wedge s(v, u))$ from Example 3. This query illustrates the role of Fork_{\neq} . We have $\text{in}(\{u\}) = \{r, s\}$. There is no implicant of $\text{in}(\{u\})$ in $\text{in}(\{u\})$ and thus $u \in \text{Fork}_{\neq}$. We obtain

$$q_{\mathcal{R}}^* = \exists u. (r(v, u) \wedge s(v, u) \wedge \neg \text{Aux}(v) \wedge \neg \text{Aux}(u)).$$

For $\mathcal{R} = \{s \sqsubseteq r\}$ and the same query q , s is an implicant of $\text{in}(\{u\})$ in $\text{in}(\{u\})$. Thus $u \notin \text{Fork}_{\neq}$. Observe that $\text{Fork}_{\mathcal{H}} = \emptyset$ as the prime implicant s of $\text{in}(\{u\})$ is contained in $\text{in}(\{u\})$. We obtain

$$q_{\mathcal{R}}^* = \exists u. (r(v, u) \wedge s(v, u) \wedge \neg \text{Aux}(v)).$$

Finally, assume $\mathcal{R} = \{r_0 \sqsubseteq r, r_0 \sqsubseteq s\}$. Again $u \notin \text{Fork}_{\neq}$, but now the prime implicant r_0 of $\text{in}(\{u\})$ is not contained in $\text{in}(\{u\})$. Thus, $\text{Fork}_{\mathcal{H}} = \{(\{r_0\}, \{u\})\}$ and we obtain

$$q_{\mathcal{R}}^* = \exists u. (r(v, u) \wedge s(v, u) \wedge \neg \text{Aux}(v) \wedge (\text{Aux}(u) \rightarrow r_0(v, u))).$$

(4) For queries $q_C = \exists \vec{v}. \psi$ that result from the unfolding of an \mathcal{EL} -concept C or have no quantified variables almost no query rewriting is needed: in both cases we have

$$q_{\mathcal{R}}^* = \exists \vec{v}. (\psi \wedge \bigwedge_{v \in \text{avar}(q)} \neg \text{Aux}(v)).$$

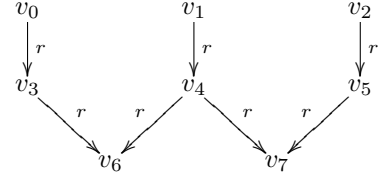


Figure 2: Query for Example (5).

(5) Let $\mathcal{R} = \emptyset$ and let $q = \exists v_0, \dots, v_7. \psi$ be the query shown in Figure 2, where all variables are quantified. Then \sim_q consists of the equivalence classes $\{v_0, v_1, v_2\}$, $\{v_3, v_4, v_5\}$, $\{v_6\}$, and $\{v_7\}$. Assume that the chosen representative for $\{v_3, v_4, v_5\}$ is v_3 . Then, we have

$$\begin{aligned} q_{\mathcal{R}}^* &= \exists v_0, \dots, v_7. (\psi \wedge \\ &\quad \text{Aux}(v_6) \rightarrow (v_3 = v_4) \wedge \\ &\quad \text{Aux}(v_7) \rightarrow (v_4 = v_5) \wedge \\ &\quad \text{Aux}(v_3) \rightarrow ((v_0 = v_1) \wedge (v_1 = v_2))) \end{aligned}$$

In the following we comment on properties of $q_{\mathcal{R}}^*$ and on consequences of our approach regarding the computational complexity of query answering in $\mathcal{ELH}_{\perp}^{\text{dr}}$. Firstly, the size of $q_{\mathcal{R}}^*$ is bounded by $\mathcal{O}(nm)$, where n is the size of q and m the size of \mathcal{R} . To see this, note that the number of conjuncts in φ_1 is bounded by the number of variables in q . For φ_2 , let $\text{Fork}_{=} = \{(T_1, \zeta_1), \dots, (T_\ell, \zeta_\ell)\}$. Then $|T_1| + \dots + |T_\ell|$ is bounded by the number of role atoms in q , and thus the size of φ_2 is $\mathcal{O}(n)$. Finally, the number of conjuncts in φ_3 is bounded by the number of quantified variables in q and each conjunct has at most m disjuncts.

Secondly, since ψ is a conjunct of the body of $q_{\mathcal{R}}^*$ it is readily checked that $q_{\mathcal{R}}^*$ is domain independent, and thus can be expressed as an SQL query. Moreover it is of the form $\exists \vec{u}. \psi$ with ψ quantifier-free. Thus, we obtain that the combined complexity of deciding whether $\mathcal{K} \models q[a_1, \dots, a_k]$ is in NP as it is sufficient to construct (in polynomial time) $\mathcal{I}_{\mathcal{K}}^r$, construct (in polynomial time) $q_{\mathcal{R}}^*$ and, finally, check whether $\mathcal{I}_{\mathcal{K}}^r \models \exists \vec{u}. \psi[a_1, \dots, a_k]$ using the obvious NP algorithm for model checking this class of formulas. Since an NP lower bound is trivial, we obtain NP-completeness. See [Krötzsch *et al.*, 2007; Rosati, 2007] for similar results regarding other variants of \mathcal{EL} .

5 Implementation and Experiments

To validate the value of the proposed approach we have conducted a series of experiments based on the NCI thesaurus (version 08.08d), which is a well-known ontology from the bio-medical domain. We have extracted an \mathcal{EL} -TBox that contains approximately 65 thousand (65K) primitive concept names, 70 primitive roles, and over 70K concept inclusions and concept definitions. The auxiliary part of canonical models, which is independent of the ABox, consists of 702K concept assertions and 171K role assertions (tuples). It was computed using the CEL system [Baader *et al.*, 2006] and loaded into the database. To add the tuples required to (a) complete \mathcal{A} to the ABox part of $\mathcal{I}_{\mathcal{K}}^r$ and (b) link the ABox part of $\mathcal{I}_{\mathcal{K}}^r$ with its auxiliary part, we have used iterated querying.

Number of assertions in ABox of \mathcal{K}										
Concept	100K	100K	100K	200K	200K	200K	400K	800K	1.6M	
Role	25K	50K	75K	40K	65K	90K	360K	1.5M	5.8M	
Number of assertions in $\mathcal{I}_{\mathcal{K}}^r$										
Concept	440K	440K	441K	683K	683K	684K	1.3M	2.6M	5.1M	
Role	197K	237K	273K	323K	371K	414K	986K	2.7M	8.2M	
Query Execution Time in seconds										
Q1 (2c1r)	0.19	0.19	0.20	0.23	0.25	0.24	0.27	0.46	0.59	
Q2 (3c2r)	0.23	0.22	0.23	0.52	0.25	0.56	0.33	0.42	0.69	
Q3 (3c2r)	0.25	0.27	0.26	0.31	0.31	0.31	0.42	0.86	1.13	
Q4 (4c3r)	0.24	0.23	0.23	0.25	0.26	0.25	0.31	0.42	1.44	
Q5 (5c5r)	0.36	0.36	0.30	0.60	0.34	0.45	2.24	7.93	128 ¹	

Figure 3: Summary of Experimental Results.

As relational systems are not optimized to handle tens of thousands of relatively small relations, one for each concept and role name in the ontology, we have used only two relations to represent $\mathcal{I}_{\mathcal{K}}^r$:

```
acbox(conceptid, indid) and
arbox(roleid, domain-indid, range-indid),
```

where `conceptid` and `roleid` are numerical identifiers for concept names and roles names, `indid`, `domain-indid`, and `range-indid` are numerical identifiers for individuals from $N_i \cup N_{aux}$, `acbox` represents concept memberships, and `arbox` role memberships. Indexes were generated on the attributes `conceptid`, `indid` and on `roleid`, `domain-indid` and `roleid`, `range-indid`, using B+trees. We distinguish individuals from N_i and N_{aux} by positive and negative ids, and thus need not store `Aux` as a relation. As an example for this representation take the query $q = Nerve(x) \wedge \neg Aux(x)$, which translates into the SQL statement

```
select indid from acbox
where conceptid=141723 and indid > 0
```

Data sets (ABox instances) for our experiments were generated randomly. When generating concept assertions, we have focused on most specific concept names, i.e., concept names without any subsumees in the TBox. The generation of role assertions was guided by the domain and range restrictions in NCI. The numbers of concept and role assertions in the initial and rewritten ABoxes used in our experiments are reported in Figure 3 in thousands/millions. The data was stored and the test queries were executed using IBM DB2 DBMS (version 9.5.0 running on SUN Fire-280R server with two UltraSPARC III 1.2GHz CPUs, 4GB memory, and 1TB storage under Solaris 5.10). Figure 3 summarizes the running times (in seconds) for each of the test queries for varying sizes of data; for each query we list the number of concept and role atoms in parentheses; the structure of the queries ranges from simple chains and star queries to queries with cycles in their bodies (we show only a few representative samples here).

The experimental results can be interpreted as follows. Firstly, the rewriting of moderately-sized ABoxes into the canonical model $\mathcal{I}_{\mathcal{K}}^r$ is well within the storage and query capabilities of existing relational technology. Secondly, query

¹This time is solely due to a large size of the result (60M tuples).

performance scales well even when the naive physical design described above is used. Thirdly, query rewriting does not increase the query processing times; the performance of rewritten queries is almost identical to the performance of the original queries over the completed ABox.

6 Limitations / Negation

Recall from Section 1 that a DL enjoys *polynomial* extended FO reducibility if it has extended FO reducibility such that the blowup of ABox rewriting (but not necessarily of query rewriting) is at most polynomial. Since ABoxes in realistic applications are large, extended FO reducibility that is not polynomial in this sense does not seem to be of much use. The following result gives a fundamental limitation of polynomial extended FO reducibility. Note that *ground* CQ answering, where a CQ may contain only individual names but not variables, is the decisional variant of CQ answering.

Theorem 6. *If the data complexity of ground CQ answering in a DL \mathcal{L} is not in PTIME, then \mathcal{L} does not enjoy polynomial extended FO reducibility.*

Proof. We show the contrapositive. Assume that \mathcal{L} enjoys polynomial extended FO reducibility, i.e., there are effectively computable mappings δ that takes each \mathcal{EL} -TBox $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ to a first-order structure $\delta(\mathcal{K})$ and γ that takes each pair (q, \mathcal{T}) , with q a k -ary conjunctive query and \mathcal{T} an \mathcal{EL} -TBox, to a first-order formula $\gamma(q, \mathcal{T})$ with k free variables such that δ can be computed in polynomial time, the size of $\delta(\mathcal{K})$ is polynomial in the size of \mathcal{K} for all \mathcal{K} , and the following condition holds: for all extended \mathcal{EL} -KBs $\mathcal{K} = (\mathcal{T}, \mathcal{A})$, all k -ary conjunctive queries q , and all tuples $(a_1, \dots, a_k) \in N_i$, $\mathcal{K} \models q[a_1, \dots, a_k]$ iff $\delta(\mathcal{K}) \models \gamma(q, \mathcal{T})[a_1, \dots, a_k]$.

Then the data complexity of ground CQ answering in \mathcal{L} is in PTIME: to check whether $\mathcal{K} \models q$ with q ground, we can compute $\delta(\mathcal{K})$ in polynomial time and $\gamma(q, \mathcal{T})$ in constant time (as the size of q and \mathcal{T} is constant), and then use FO model checking to decide whether $\delta(\mathcal{K}) \models \gamma(q, \mathcal{T})$. The latter can be done in LOGSPACE. \square

For expressive DLs such as \mathcal{ALC} and \mathcal{SHIQ} , the data complexity of ground CQ answering is co-NP-hard, thus Theorem 6 implies that these DLs do not enjoy polynomial extended FO reducibility unless PTIME = co-NP.

In many applications, it is natural to admit also negated concept assertions in the ABox and to extend conjunctive queries with negated concept atoms in order to query the resulting ABoxes, see e.g. [Patel *et al.*, 2007]. Let us call ABoxes, queries, and KBs of this form *literal*. The result stated in Theorem 6 also applies to literal KBs and literal queries. Since the data complexity of literal CQ answering over literal \mathcal{EL} -KBs (where \mathcal{EL} is $\mathcal{ELH}_{\perp}^{dr}$ without \perp , role hierarchies, and domain and range restrictions) is co-NP-hard [Schaerf, 1993], it follows that even this mild extension of \mathcal{EL} does not enjoy polynomial extended FO reducibility.

However, there is still a (pragmatic, yet formal) way to add negation to our approach. In the following, we sketch an incomplete approach to literal CQ answering over literal

$\mathcal{ELH}_{\perp}^{dr}$ -KBs. Its incompleteness can be precisely characterized in terms of an epistemic semantics for negated concept atoms in literal CQs, inspired by [Calvanese *et al.*, 2007a].

As a preliminary, we assume *standard names*, i.e., interpretation domains $\Delta^{\mathcal{I}}$ have to be a subset of the (countably infinite) set N_I of individual names and $a^{\mathcal{I}} = a$ for all \mathcal{I} and all $a \in N_I$. For a literal $\mathcal{ELH}_{\perp}^{dr}$ -KB \mathcal{K} , an interpretation \mathcal{I} , and a literal CQ $q = \exists \vec{u}. \psi(\vec{u}, \vec{v})$ with $\vec{v} = v_1, \dots, v_k$, we set $\mathcal{I} \models_e q[a_1, \dots, a_k]$ iff there exists a variable assignment π with $\pi(v_i) = a_i$ for $1 \leq i \leq k$ and

- $\mathcal{I} \models_{\pi} \alpha$ for all positive atoms α in ψ ,
- $\mathcal{K} \models \neg A(a)$ for all $\neg A(a)$ in ψ with $a \in N_I$, and
- $\mathcal{K} \models \neg A(\pi(v))$ for all $\neg A(v)$ in ψ with $v \in N_V$.

Now set $\mathcal{K} \models_e q[a_1, \dots, a_k]$ iff $\mathcal{I} \models_e q[a_1, \dots, a_k]$ for all models \mathcal{I} of \mathcal{K} with standard names. Thus, answers to negated atoms do not depend on a concrete interpretation \mathcal{I} , but only on deducibility. Epistemic semantics is sound: every answer is also an answer under the standard semantics (but not vice versa); it is also conservative: it yields the same answers as standard semantics when either the KB or the query does not contain negation. An example for which this epistemic semantics is different from the standard semantics is given by $\mathcal{A} = \{A'(a), \neg A(a)\}$, $\mathcal{T} = \{A' \sqsubseteq \exists r. \top, \exists r. B \sqsubseteq A\}$, and $q = \exists u. r(a, u) \wedge \neg B(u)$. Then $\mathcal{K} \models q$ but $\mathcal{K} \not\models_e q$.

We give a simple (and poly-time) reduction of epistemic literal CQ answering over literal $\mathcal{ELH}_{\perp}^{dr}$ -KBs to standard CQ answering over $\mathcal{ELH}_{\perp}^{dr}$ -KBs. Via the rewritings presented in the main part of this paper, the reduction enables the use of RDBMSs also for the case of $\mathcal{ELH}_{\perp}^{dr}$ with negation (under the epistemic semantics).

Given a consistent literal $\mathcal{ELH}_{\perp}^{dr}$ -KB $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ and a literal query q , replace every literal $\neg A$ in q with a fresh concept name \bar{A} . Additionally, add $\top \sqsubseteq \bar{A}$ to \mathcal{T} whenever $\mathcal{T} \models A \sqsubseteq \perp$, and $\bar{A}(a)$ to \mathcal{A} for all $a \in \text{Ind}(\mathcal{A})$ with $\mathcal{K} \models \neg A(a)$, and then remove all negated assertions from \mathcal{A} . Call the result $\mathcal{K}' = (\mathcal{T}', \mathcal{A}')$ and q' .

Theorem 7. *Let \mathcal{K} be a consistent literal $\mathcal{ELH}_{\perp}^{dr}$ -KB. Then \mathcal{K}' can be computed in polynomial time and $\mathcal{K} \models_e q[a_1, \dots, a_k]$ iff $\mathcal{K}' \models q'[a_1, \dots, a_k]$ for all literal CQs q and all $a_1, \dots, a_k \in \text{Ind}(\mathcal{A})$.*

7 Conclusion

We have proposed a novel approach to CQ answering in DLs using RDBMSs. In particular, our approach paves the way to using RDBMSs for CQ answering in the \mathcal{EL} family of DLs. Our experiments exhibit a promising performance even without a sophisticated physical design. One drawback of our approach is the blowup of the data, which is polynomial but still considerable on large data sets. As future work, it might be interesting to reduce this blowup by incorporating the TBox partly into the data and partly into the query. We will also develop effective approaches to update the canonical model/auxiliary data when assertions are added to or deleted from the ABox.

From a theoretical perspective, it would be interesting to further investigate the capabilities of our approach. For example, one could extend it to $\mathcal{ELH}_{\perp}^{dr}$ with transitive roles and try to adapt our technique to the DL-Lite family of DLs, thus providing an alternative to the approach in [Calvanese *et al.*, 2007b].

References

- [Anonymous, 2009] Anonymous. Conjunctive query answering in the DL \mathcal{EL} using an RDBMS. Anonymized Technical Report, 2009. Available from <http://ijcai09.001webs.com/>
- [Baader *et al.*, 2005b] F. Baader, S. Brandt, and C. Lutz. Pushing the \mathcal{EL} envelope. In *Proc. of IJCAI'05*, pages 364–369. Professional Book Center, 2005.
- [Baader *et al.*, 2006] F. Baader, C. Lutz, and B. Suntisrivaraporn. CEL—a polynomial time reasoner for life science ontologies. In *Proc. of IJCAR'06*, volume 4130 of *LNCS*, pages 287–291. Springer, 2006.
- [Baader *et al.*, 2008] F. Baader, S. Brandt, and C. Lutz. Pushing the \mathcal{EL} envelope further. In *In Proc. of OWLED'08*, 2008.
- [Calvanese *et al.*, 2007a] D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati. EQL-Lite: Effective first-order query processing in DLs. In *Proc. of IJCAI'07*, pages 274–279. AAAI press, 2007.
- [Calvanese *et al.*, 2007b] D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati. Tractable reasoning and efficient query answering in DLs: The DL-Lite family. *J. of Automated Reasoning*, 39(3):385–429, 2007.
- [Krisnadhi and Lutz, 2007] A. Krisnadhi and C. Lutz. Data complexity in the \mathcal{EL} family of DLs. In *Proc. of LPAR'07*, volume 4790 of *LNCS*, pages 333–347. Springer, 2007.
- [Krötzsch *et al.*, 2007] M. Krötzsch, S. Rudolph, and P. Hitzler. Conjunctive queries for a tractable fragment of OWL 1.1. In *Proc. of ISWC2007*, volume 4825 of *LNCS*, pages 310–323. Springer, 2007.
- [Lutz and Wolter, 2007] C. Lutz and F. Wolter. Conservative extensions in the lightweight description logic \mathcal{EL} . In *Proc. of CADE-21*, volume 4603 of *LNAI*, pages 84–99. Springer, 2007.
- [Patel *et al.*, 2007] C. Patel, J. J. Cimino, J. Dolby, A. Fokoue, A. Kalyanpur, A. Kershenbaum, L. Ma, E. Schonberg, and K. Srinivas. Matching patient records to clinical trials using ontologies. In *Proceedings of ISWC'07*, volume 4825 of *LNCS*, pages 816–829. Springer, 2007.
- [Rosati, 2007] Riccardo Rosati. On conjunctive query answering in \mathcal{EL} . In *Proc. of DL2007*, volume 250 of *CEUR-WS*, 2007.
- [Schaerf, 1993] A. Schaerf. On the complexity of the instance checking problem in concept languages with existential quantification. *J. of Intelligent Information Systems*, 2:265–278, 1993.
- [Sioutos *et al.*, 2006] N. Sioutos, S. de Coronado, M.W. Haber, F.W. Hartel, W.L. Shaiu, and L.W. Wright. NCI thesaurus: a semantic model integrating cancer-related clinical and molecular information. *J. of Biomedical Informatics*, 40(1):30–43, 2006.
- [Spackman, 2000] K.A. Spackman. Managing clinical terminology hierarchies using algorithmic calculation of subsumption: Experience with SNOMED-RT. *J. of the American Medical Informatics Association*, 2000.
- [Wu *et al.*, 2008] Z. Wu, G. Eadon, S. Das, E. Chong, V. Kolovski, M. Annamalai, and J. Srinivasan. Implementing an inference engine for RDFS/OWL constructs and user-defined rules in oracle. In *Proc. of ICDE'08*, pages 1239–1248. IEEE, 2008.