

# Computing FO-Rewritings in $\mathcal{EL}$ in Practice: from Atomic to Conjunctive Queries

Peter Hansen and Carsten Lutz

University of Bremen, Germany  
{hansen, clu}@informatik.uni-bremen.de

**Abstract.** A prominent approach to implementing ontology-mediated queries (OMQs) is to rewrite into a first-order query, which is then executed using a conventional SQL database system. We consider the case where the ontology is formulated in the description logic  $\mathcal{EL}$  and the actual query is a conjunctive query and show that rewritings of such OMQs can be efficiently computed in practice, in a sound and complete way. Our approach combines a reduction with a decomposed backwards chaining algorithm for OMQs that are based on the simpler atomic queries, also illuminating the relationship between first-order rewritings of OMQs based on conjunctive and on atomic queries. Experiments with real-world ontologies show promising results.

## 1 Introduction

One of the most important tools in ontology-mediated querying is *query rewriting*: reformulate a given ontology-mediated query (OMQ) in an equivalence-preserving way in a query language that is supported by a database system used to store the data. Since SQL is the dominating query language in conventional database systems, rewriting into SQL and into first-order logic (FO) as its logical core has attracted particularly much attention [3, 4, 5, 6, 7, 10, 12, 15]. In fact, the DL-Lite family of description logics (DLs) was invented specifically with the aim to guarantee that FO-rewritings of OMQs (whose ontology is formulated in DL-Lite) always exist [1, 7], but is rather restricted in expressive power. For essentially all other DLs, there are OMQs which cannot be equivalently rewritten into an FO query. However, ontologies used in real-world applications tend to have a very simple structure and, consequently, one may hope that FO-rewritings of practically relevant OMQs exist in the majority of cases. This hope was confirmed in an experimental evaluation carried out in the context of the  $\mathcal{EL}$  family of description logics where less than 1% of the considered queries was found not to be FO-rewritable [12]; moreover, most of the negative cases seemed to be due to modeling mistakes in the ontology.

In this paper, we focus on the description logic  $\mathcal{EL}$ , which can be viewed as a logical core of the OWL EL profile of the OWL 2 ontology language [20]. We use  $(\mathcal{L}, \mathcal{Q})$  to denote the OMQ language that consists of all OMQs where the ontology is formulated in the description logic  $\mathcal{L}$  and the actual query is formulated in the

query language  $\mathcal{Q}$ . Important choices for  $\mathcal{Q}$  include *atomic queries* (AQs) and the much more expressive *conjunctive queries* (CQs). It has been shown in [6] that for OMQs from  $(\mathcal{EL}, \text{AQ})$ , it is EXPTIME-complete to decide FO-rewritability. Combining the techniques from [6] and the backwards chaining approach to query rewriting brought forward e.g. in [8, 15], a practical algorithm for computing FO-rewritings of OMQs from  $(\mathcal{EL}, \text{AQ})$  was then developed in [12]. This algorithm is based on a *decomposed version* of backwards chaining that implements a form of structure sharing. It was implemented in the *Grind* system and shown to perform very well in practice [12]. It is important to remark that the algorithm is *complete*, that is, it computes an FO-rewriting whenever there is one and reports failure otherwise.

The aim of this paper is to devise a way to efficiently compute FO-rewritings of OMQs from  $(\mathcal{EL}, \text{CQ})$ , and thus the challenge is to deal with conjunctive queries instead of only with atomic ones. Note that, as shown in [5], FO-rewritability in  $(\mathcal{EL}, \text{CQ})$  is still EXPTIME-complete. Our approach is to combine a reduction with the decomposed algorithm from [12], also illuminating the relationship between first-order rewritings of OMQs based on CQs and on AQs. It is worthwhile to point out that naive reductions of FO-rewritability in  $(\mathcal{EL}, \text{CQ})$  to FO-rewritability in  $(\mathcal{EL}, \text{AQ})$  fail. In particular, FO-rewritability of all AQs that occur in a CQ  $q$  are neither a sufficient nor a necessary condition for  $q$  to be FO-rewritable. As a simple example, consider the OMQ that consists of the ontology and query

$$\mathcal{O} = \{\exists r.A \sqsubseteq A, \exists s.\top \sqsubseteq A\} \quad \text{and} \quad q(x) = \exists y (A(x) \wedge s(x, y))$$

and which is FO-rewritable into  $\exists y s(x, y)$ , but the only AQ  $A(x)$  that occurs in  $q$  is not FO-rewritable in the presence of  $\mathcal{O}$ .<sup>1</sup> In fact, it is not clear how to attain a reduction of FO-rewritability in  $(\mathcal{EL}, \text{CQ})$  to FO-rewritability in  $(\mathcal{EL}, \text{AQ})$ , and even less so a polynomial time one. This leads us to considering mildly restricted forms of CQs and admitting reductions that make certain assumptions on the algorithm used to compute FO-rewritings in  $(\mathcal{EL}, \text{AQ})$ —all of them are satisfied by the decomposed backwards chaining algorithm implemented in Grind.

We first consider the class of *tree-quantified CQs* (*tqCQs*) in which the quantified parts of the CQ form a collection of directed trees. In this case, we indeed achieve a polynomial time reduction to FO-rewritability in  $(\mathcal{EL}, \text{AQ})$ . To also transfer actual FO-rewritings from the OMQ constructed in the reduction to the original OMQ, we make the assumption that the rewriting of the former takes the form of a UCQ (union of conjunctive queries) in which every CQ is tree-shaped and that, in a certain sense made precise in the paper, atoms are never introduced into the rewriting ‘without a reason’. Both conditions are very natural in the context of backwards chaining and satisfied by the decomposed algorithm.

We then move to *rooted CQs* (*rCQs*) in which every quantified variable must be reachable from some answer variable (in an undirected sense, in the query

<sup>1</sup> OMQs also allow to fix the signature (set of concept and role names) that can occur in the ABox. In this example, we do not assume any restriction on the ABox signature.

graph). We consider this a mild restriction and expect that almost all queries in practical applications will be rCQs. In the rCQ case, we do not achieve a ‘black box’ reduction. Instead, we assume that FO-rewritings of the constructed OMQs from  $(\mathcal{EL}, \text{AQ})$  are obtained from a certain straightforward backwards chaining algorithm or a refinement thereof as implemented in the Grind system. We then show how to combine the construction of (several) OMQs from  $(\mathcal{EL}, \text{AQ})$ , similar to those constructed in the tqCQ case, with a modification of the assumed algorithm to decide FO-rewritability in  $(\mathcal{EL}, \text{rCQ})$  and to construct actual rewritings. The approach involves exponential blowups, but only in parameters that we expect to be very small in practical cases and that, in particular, only depend on the actual query contained in the OMQ but not on the ontology.

We have implemented our approach in the Grind system and carried out experiments on five real-world ontologies with 10 hand-crafted CQs for each. The average runtimes are between 0.5 and 19 seconds (depending on the ontology), which we consider very reasonable given that we are dealing with a complex static analysis problem.

Proofs are deferred to the appendix, which is made available at <http://www.cs.uni-bremen.de/tdki/research/papers.html>.

**Related Work.** We directly build on our prior work in [12] as discussed above, and to a lesser degree also on [5,6]. The latter line of work has recently been picked up in the context of existential rules [3]. The distinguishing features of our work are that (1) our algorithms are sound, complete, and terminating, that is, they find an FO-rewriting if there is one and report failure otherwise, and (2) we rely on the decomposed calculus from [12] that implements structure sharing for constructing small rewritings and achieving practical feasibility. We are not aware of other work that combines features (1) and (2) and is applicable to OMQs based on  $\mathcal{EL}$ . In the context of the description logic DL-Lite, though, the construction of small rewritings has received a lot of attention, see e.g. [11, 13, 22, 23]. Producing small rewritings of OMQs whose ontology is a set of existential rules has been studied in [14], but there are no termination guarantees. Constructing small *Datalog*-rewritings of OMQs based on  $\mathcal{EL}$ , which are guaranteed to always exist, was studied e.g. in [9, 21, 25, 26]. A different approach to answering  $\mathcal{EL}$ -based OMQs using SQL databases is the combined approach where the consequences of the ontology are materialized in the data [18, 24].

## 2 Preliminaries

Let  $\mathbb{N}_C$ ,  $\mathbb{N}_R$ , and  $\mathbb{N}_I$  be countably infinite sets of *concept names*, *role names*, and *individual names*. An  $\mathcal{EL}$ -*concept* is formed according to the syntax rule

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

where  $A$  ranges over  $\mathbb{N}_C$  and  $r$  over  $\mathbb{N}_R$ . An  $\mathcal{EL}$ -*TBox*  $\mathcal{T}$  is a finite set of *concept inclusions*  $C \sqsubseteq D$ , with  $C$  and  $D$   $\mathcal{EL}$ -concepts. Throughout the paper, we use  $\mathcal{EL}$ -TBoxes as ontologies. An *ABox* is a finite set of *concept assertions*  $A(a)$  and

role assertions  $r(a, b)$  where  $a$  and  $b$  range over  $\mathbf{N}_I$ . We use  $\text{Ind}(\mathcal{A})$  to denote the set of individual names in the ABox  $\mathcal{A}$ . A *signature* is a set of concept and role names. When an ABox uses only symbols from a signature  $\Sigma$ , then we call it a  $\Sigma$ -ABox. To emphasize that a signature  $\Sigma$  is used to constrain the symbols admitted in ABoxes, we sometimes call  $\Sigma$  an *ABox signature*.

The semantics of concepts, TBoxes, and ABoxes is defined in the usual way, see [2]. We write  $\mathcal{T} \models C \sqsubseteq D$  if the concept inclusion  $C \sqsubseteq D$  is satisfied in every model of  $\mathcal{T}$ ; when  $\mathcal{T}$  is empty, we write  $\models C \sqsubseteq D$ . As usual in ontology-mediated querying, we make the *standard names assumption*, that is, an interpretation  $\mathcal{I}$  satisfies a concept assertion  $A(a)$  if  $a \in A^{\mathcal{I}}$  and a role assertion  $r(a, b)$  if  $(a, b) \in r^{\mathcal{I}}$ .

A *conjunctive query* (CQ) takes the form  $q(\mathbf{x}) = \exists \mathbf{y} \varphi(\mathbf{x}, \mathbf{y})$  with  $\mathbf{x}, \mathbf{y}$  tuples of variables and  $\varphi$  a conjunction of atoms of the form  $A(x)$  and  $r(x, y)$  that uses only variables from  $\text{var}(q) = \mathbf{x} \cup \mathbf{y}$ . The variables  $\mathbf{x}$  are the *answer variables* of  $q$ , denoted  $\text{avar}(q)$ , and the *arity* of  $q$  is the length of  $\mathbf{x}$ . Unless noted otherwise, we allow equality in CQs, but we assume w.l.o.g. that equality atoms contain only answer variables, and that when  $x = y$  is an equality atom in  $q$ , then  $y$  does not occur in any other atoms in  $q$ . Other occurrences of equality can be eliminated by identifying variables. An *atomic query* (AQ) is a conjunctive query of the form  $A(x)$ . A *union of conjunctive queries* (UCQ) is a disjunction of CQs that share the same answer variables.

An *ontology-mediated query* (OMQ) is a triple  $Q = (\mathcal{T}, \Sigma, q)$  where  $\mathcal{T}$  is a TBox,  $\Sigma$  an ABox signature, and  $q$  a CQ. We use  $(\mathcal{EL}, \text{AQ})$  to denote the set of OMQs where  $\mathcal{T}$  is an  $\mathcal{EL}$ -TBox and  $q$  is an AQ, and similarly for  $(\mathcal{EL}, \text{CQ})$  and so on. We do generally not allow equality in CQs that are part of an OMQ. Let  $Q = (\mathcal{T}, \Sigma, q)$  be an OMQ,  $\mathcal{A}$  a  $\Sigma$ -ABox and  $\mathbf{a} \subseteq \text{Ind}(\mathcal{A})$ . We write  $\mathcal{A} \models Q(\mathbf{a})$  if  $\mathcal{I} \models q(\mathbf{a})$  for all models  $\mathcal{I}$  of  $\mathcal{T}$  and  $\mathcal{A}$ . In this case,  $\mathbf{a}$  is a *certain answer* to  $Q$  on  $\mathcal{A}$ .

*Example 1.* Consider an example from the medical domain. The following ABox holds data about patients and diagnoses:

$$\mathcal{A} = \{\text{Person}(a), \text{hasDisease}(a, \text{oca}_1), \text{Albinism}(\text{oca}_1)\}$$

A TBox  $\mathcal{T}_1$  is used to make domain knowledge available:

$$\begin{aligned} \mathcal{T}_1 = \{ & \text{Albinism} \sqsubseteq \text{HereditaryDisease}, \\ & \text{Person} \sqcap \exists \text{hasDisease.HereditaryDisease} \sqsubseteq \text{GeneticRiskPatient} \} \end{aligned}$$

Let  $Q_1$  be the OMQ  $(\mathcal{T}_1, \Sigma_{\text{full}}, q_1(x))$ , where  $q_1(x) = \text{GeneticRiskPatient}(x)$ , and  $\Sigma_{\text{full}}$  contains all concept and role names. It can be verified that  $\mathcal{A} \models Q_1(a)$ .  $\dashv$

We do not distinguish between a CQ and the set of atoms in it and associate with each CQ  $q$  a directed graph  $G_q := (\text{var}(q), \{(x, y) \mid r(x, y) \in q\})$  (equality atoms are not reflected). A CQ  $q$  is *tree-shaped* if  $G_q$  is a directed tree and  $r(x, y), s(x, y) \in q$  implies  $r = s$ . A *tree CQ* (tCQ) is a tree-shaped CQ with the root the only answer variable and a *tree UCQ* (tUCQ) is a disjunction

of tree CQs. Every  $\mathcal{EL}$ -concept can be viewed as a tree-shaped CQ and vice versa; for example, the  $\mathcal{EL}$ -concept  $A \sqcap \exists r.(B \sqcap \exists s.A)$  corresponds to the CQ  $q(x) = \exists y, z A(x) \wedge r(x, y) \wedge B(y) \wedge s(y, z) \wedge A(z)$ . We will not always distinguish between the two representations and even mix them. We might thus write  $\exists r.q$  to denote an  $\mathcal{EL}$ -concept when  $q$  is a tree-shaped CQ; if  $q(x)$  is as in the example just given, then  $\exists r.q$  is the  $\mathcal{EL}$ -concept  $\exists r.(A \sqcap \exists r.(B \sqcap \exists s.A))$ . If convenient, we also view a CQ  $q$  as an ABox  $\mathcal{A}_q$  which is obtained from  $q$  by dropping equality atoms and then replacing each variable with an individual (not distinguishing answer variables from quantified variables). A *rooted CQ* (*rCQ*) is a CQ  $q$  such that in the undirected graph induced by  $G_q$ , every quantified variable is reachable from some answer variable. A *tree-quantified CQ* (*tqCQ*) is an rCQ  $q$  such that after removing all atoms  $r(x, y)$  with  $x, y \in \text{avar}(q)$ , we obtain a disjoint union of tCQs. We call these tCQs the *tCQs in  $q$* . For example,  $q(x_1, x_2) = \exists y_1, y_2 r(x_1, x_2) \wedge r(x_2, x_1) \wedge r(x_1, y_1) \wedge s(x_2, y_2)$  is a tqCQ and the tCQs in  $q$  are  $\exists y_1 r(x_1, y_1)$  and  $\exists y_2 s(x_2, y_2)$ ; by adding to  $q$  the atom  $r(y_1, y_2)$ , we obtain an rCQ that is not a tqCQ.

An OMQ  $Q = (\mathcal{T}, \Sigma, q)$  is *FO-rewritable* if there is a first-order (FO) formula  $\varphi$  such that  $\mathcal{A} \models Q(\mathbf{a})$  iff  $\mathcal{A} \models \varphi(\mathbf{a})$  for all  $\Sigma$ -ABoxes  $\mathcal{A}$ . In this case,  $\varphi$  is an *FO-rewriting of  $Q$* . When  $\varphi$  happens to be a UCQ, we speak of a *UCQ-rewriting* and likewise for other classes of queries. It is known that FO-rewritability coincides with UCQ-rewritability for OMQs from  $(\mathcal{EL}, \text{CQ})$  [4, 6]; note that equality is important here as, for example, the OMQ  $(\{B \sqsubseteq \exists r.A\}, \{B, r\}, q)$  with  $q(x, y) = \exists z(r(x, z) \wedge r(y, z) \wedge A(z))$  rewrites into the UCQ  $q \vee (B(x) \wedge x = y)$ , but not into an UCQ that does not use equality.

*Example 2.* We extend the TBox  $\mathcal{T}_1$  from Example 1 to additionally describe the hereditary nature of genetic defects:

$$\mathcal{T}_2 := \mathcal{T}_1 \cup \{\text{Person} \sqcap \exists \text{hasParent.GeneticRiskPatient} \sqsubseteq \text{GeneticRiskPatient}\}.$$

The OMQ  $Q'_1 = (\mathcal{T}_2, \Sigma_{\text{full}}, q_1(x))$  with  $q_1(x)$  as in Example 1, is not FO-rewritable, intuitively because it expresses unbounded reachability along the *hasParent* role. In contrast, consider the OMQ  $Q_2 = (\mathcal{T}_2, \Sigma_{\text{full}}, q_2(x))$  where  $q_2(x) = \exists y \text{GeneticRiskPatient}(x) \wedge \text{hasDisease}(x, y) \wedge \text{Albinism}(y)$ . Even though  $q_2$  is an extension of  $q_1$  with additional atoms,  $Q_2$  is FO-rewritable, with  $\varphi(x) = q_2(x) \vee (\exists y \text{Person}(x) \wedge \text{hasDisease}(x, y) \wedge \text{Albinism}(y))$  a concrete rewriting.  $\dashv$

We shall sometimes refer to the problem of (*query*) *containment* between two OMQs  $Q_1 = (\mathcal{T}_1, \Sigma, q_1)$  and  $Q_2 = (\mathcal{T}_2, \Sigma, q_2)$ ; we say  $Q_1$  is *contained in  $Q_2$*  if  $\mathcal{A} \models Q_1(\mathbf{a})$  implies  $\mathcal{A} \models Q_2(\mathbf{a})$  for all  $\Sigma$ -ABoxes  $\mathcal{A}$  and  $\mathbf{a} \subseteq \text{Ind}(\mathcal{A})$ . If both OMQs are from  $(\mathcal{EL}, \text{rCQ})$  and  $\mathcal{T}_1 = \mathcal{T}_2 = \mathcal{T}$ , then we denote this with  $q_1 \subseteq_{\mathcal{T}} q_2$ .

We now introduce two more technical notions that are central to the constructions in Section 4. Both notions have been used before in the context of ontology-mediated querying, see for example [16, 17]. They are illustrated in Example 3 below.

**Definition 1 (Fork rewriting).** Let  $q_0$  be a CQ. Obtaining a CQ  $q$  from  $q_0$  by fork elimination means to select two atoms  $r(x_0, y)$  and  $r(x_1, y)$  with  $y$  an

existentially quantified variable, then to replace every occurrence of  $x_{1-i}$  in  $q$  with  $x_i$ , where  $i \in \{0, 1\}$  is chosen such that  $x_i$  is an answer variable if any of  $x_0, x_1$  is an answer variable, and to finally add the atom  $x_i = x_{1-i}$  if  $x_{1-i}$  is an answer variable. When  $q$  can be obtained from  $q_0$  by repeated (but not necessarily exhaustive) fork elimination, then  $q$  is a fork rewriting of  $q_0$ .

For a CQ  $q$  and  $V \subseteq \text{var}(q)$ , we use  $q|_V$  to denote the restriction of  $q$  to the variables in  $V$ , that is,  $q|_V$  is the set of atoms in  $q$  that use only variables from  $V$ .

**Definition 2 (Splitting).** Let  $\mathcal{T}$  be an  $\mathcal{EL}$ -TBox,  $q$  a CQ, and  $\mathcal{A}$  an ABox. A splitting of  $q$  w.r.t.  $\mathcal{A}$  and  $\mathcal{T}$  is a tuple  $\Pi = \langle R, S_1, \dots, S_\ell, r_1, \dots, r_\ell, \mu, \nu \rangle$ , where  $R, S_1, \dots, S_\ell$  is a partitioning of  $\text{var}(q)$ ,  $r_1, \dots, r_\ell$  are role names,  $\mu : \{1, \dots, \ell\} \rightarrow R$  assigns to each set  $S_i$  a variable from  $R$ ,  $\nu : R \rightarrow \text{Ind}(\mathcal{A})$  assigns to each variable from  $R$  and individual name from  $\mathcal{A}$ , and the following conditions are satisfied:

1.  $\text{avar}(q) \subseteq R$  and  $x = y \in q$  implies  $\nu(x) = \nu(y)$ ;
2. if  $r(x, y) \in q$  with  $x, y \in R$ , then  $r(\nu(x), \nu(y)) \in \mathcal{A}$ ;
3.  $q|_{S_i}$  is tree-shaped and can thus be seen as an  $\mathcal{EL}$ -concept  $C_{q|_{S_i}}$ , for  $1 \leq i \leq \ell$ ;
4. if  $r(x, x') \in q$  then either (i)  $x, x'$  belong to the same set  $R, S_1, \dots, S_\ell$ , or (ii)  $x \in R$  and, for some  $i$ ,  $r = r_i$  and  $x'$  root of  $q|_{S_i}$ .

The following lemma illustrates the combined use and raison d'être of both fork rewritings and splittings. A proof is standard and omitted, see for example [17]. It does rely on the existence of *forest models* for ABoxes and  $\mathcal{EL}$ -TBoxes, that is, for every ABox  $\mathcal{A}$  and TBox  $\mathcal{T}$ , there is a model  $\mathcal{I}$  whose shape is that of  $\mathcal{A}$  with a directed (potentially infinite) tree attached to each individual.

**Lemma 1.** Let  $Q = (\mathcal{T}, \Sigma, q_0)$  be an OMQ from  $(\mathcal{EL}, CQ)$ ,  $\mathcal{A}$  a  $\Sigma$ -ABox, and  $\mathbf{a} \subseteq \text{Ind}(\mathcal{A})$ . Then  $\mathcal{A} \models Q(\mathbf{a})$  iff there exists a fork rewriting  $q$  of  $q_0$  and a splitting  $\langle R, S_1, \dots, S_\ell, r_1, \dots, r_\ell, \mu, \nu \rangle$  of  $q$  w.r.t.  $\mathcal{A}$  and  $\mathcal{T}$  such that the following conditions are satisfied:

1.  $\nu(\mathbf{x}) = \mathbf{a}$ ,  $\mathbf{x}$  the answer variables of  $q_0$ ;
2. if  $A(x) \in q$  and  $x \in R$ , then  $\mathcal{A}, \mathcal{T} \models A(\nu(x))$ ;
3.  $\mathcal{A}, \mathcal{T} \models \exists r_i. C_{q|_{S_i}}(\nu(\mu(i)))$  for  $1 \leq i \leq \ell$ .

*Example 3.* To illustrate the described notions, consider the following CQ.

$$q_3(x) = \exists y_1, y_2, z \text{ Person}(x) \wedge \\ \text{hasDisease}(x, y_1) \wedge \text{MelaminDeficiency}(y_1) \wedge \text{causedBy}(y_1, z) \wedge \\ \text{hasDisease}(x, y_2) \wedge \text{ImpairedVision}(y_2) \wedge \text{causedBy}(y_2, z) \wedge \\ \text{GeneDefect}(z)$$

It asks for persons suffering from two conditions connected with the same gene defect. Let the ABox  $\mathcal{A}$  consist only of the assertion  $\text{OCA1aPatient}(a)$ . We extend the TBox  $\mathcal{T}_2$  from Example 2, as follows:

$$\mathcal{T}_3 := \mathcal{T}_2 \cup \{ \text{OCA1aPatient} \sqsubseteq \text{Person} \sqcap \text{hasDisease.OCA1aAlbinism} \\ \text{OCA1aAlbinism} \sqsubseteq \text{ImpairedVision} \sqcap \text{MelaninDeficiency} \\ \text{OCA1aAlbinism} \sqsubseteq \exists \text{causedBy.GeneDefect} \}$$

Let  $Q = (\mathcal{T}_3, \Sigma_{\text{full}}, q_3(x))$ . It can be verified that  $\mathcal{A} \models Q(a)$ . By Lemma 1, this is witnessed by a fork rewriting and a splitting  $\Pi$ . The fork rewriting is

$$q'_3(x) = \exists y_1, z \text{ Person}(x) \wedge \\ \text{hasDisease}(x, y_1) \wedge \text{MelaminDeficiency}(y_1) \wedge \text{ImpairedVision}(y_1) \wedge \\ \text{causedBy}(y_1, z) \wedge \text{GeneDefect}(z)$$

The splitting  $\Pi = \langle R, S_1, r_1, \mu, \nu \rangle$  of  $q'_3$  wrt.  $\mathcal{A}$  and  $\mathcal{T}_3$  is defined by setting

$$R = \{x\}, S_1 = \{y_1, z\}, r_1 = \text{hasDisease}, \mu(1) = x, \nu = (x \mapsto a)$$

It can be verified that the conditions given in Lemma 1 are satisfied.  $\dashv$

### 3 Tree-quantified CQs

We reduce FO-rewritability in  $(\mathcal{EL}, \text{tqCQ})$  to FO-rewritability in  $(\mathcal{EL}, \text{AQ})$  and, making only very mild assumptions on the algorithm used for solving the latter problem, show that rewritings of the OMQs produced in the reduction can be transformed in a straightforward way into rewritings of the original OMQ. The mild assumptions are that the algorithm produces a tUCQ-rewriting and that, informally, when constructing the tCQs of the tUCQ-rewriting it never introduces atoms ‘without a reason’—this will be made precise later.

Let  $Q = (\mathcal{T}, \Sigma, q_0)$  be from  $(\mathcal{EL}, \text{tqCQ})$ . We can assume w.l.o.g. that  $q_0$  contains only answer variables: every tCQ in  $q$  with root  $x$  can be represented as an  $\mathcal{EL}$ -concept  $C$  and we can replace the tree with the atom  $A_C(x)$  (unless it has only a single node) and extend  $\mathcal{T}$  with  $C \sqsubseteq A_C$  where  $A_C$  is a fresh concept name that is not included in  $\Sigma$ . Clearly, the resulting OMQ is equivalent to the original one.

Let  $Q$  be an OMQ from  $(\mathcal{EL}, \text{tqCQ})$ . We show how to construct an OMQ  $Q' = (\mathcal{T}', \Sigma', q'_0)$  from  $(\mathcal{EL}, \text{AQ})$  with the announced properties; in particular,  $Q$  is FO-rewritable if and only if  $Q'$  is. Let  $\text{CN}(\mathcal{T})$  and  $\text{RN}(\mathcal{T})$  denote the set of concept names and role names that occur in  $\mathcal{T}$ , and let  $\text{sub}_L$  denote the set of concepts that occur on the left-hand side of a concept inclusion in  $\mathcal{T}$ , closed under subconcepts. Reserve a fresh concept name  $A^x$  for every  $A \in \text{CN}(\mathcal{T})$  and  $x \in \text{avar}(q_0)$ , and a fresh role name  $r^x$  for every  $r \in \text{RN}(\mathcal{T})$  and  $x \in \text{avar}(q_0)$ . Set

$$\Sigma' = \Sigma \cup \{A^x \mid A \in \text{CN}(\mathcal{T}) \cap \Sigma \text{ and } x \in \text{avar}(q_0)\} \\ \cup \{r^x \mid r \in \text{RN}(\mathcal{T}) \cap \Sigma \text{ and } x \in \text{avar}(q_0)\}.$$

Additionally reserve a concept name  $A_{\exists r.E}^x$  for every concept  $\exists r.E \in \text{sub}_L(\mathcal{T})$  and every  $x \in \text{avar}(q_0)$ . Define

$$\mathcal{T}' := \mathcal{T} \cup \{C_L^x \sqsubseteq D_R^x \mid x \in \text{var}(q_0) \text{ and } C \sqsubseteq D \in \mathcal{T}\} \\ \cup \{\exists r^x.C \sqsubseteq A_{\exists r.C}^x \mid x \in \text{var}(q_0) \text{ and } \exists r.C \in \text{sub}_L(\mathcal{T})\} \\ \cup \{C_L^y \sqsubseteq A_{\exists r.C}^x \mid r(x, y) \in q_0 \text{ and } \exists r.C \in \text{sub}_L(\mathcal{T})\} \\ \cup \left\{ \prod_{A(x) \in q_0} A^x \sqsubseteq N \right\}$$

where for a concept  $C = A_1 \sqcap \dots \sqcap A_n \sqcap \exists r_1.E_1 \sqcap \dots \sqcap \exists r_m.E_m$ , the concepts  $C_L^x$  and  $C_R^x$  are given by

$$\begin{aligned} C_L^x &= A_1^x \sqcap \dots \sqcap A_n^x \sqcap A_{\exists r_1.E_1}^x \sqcap \dots \sqcap A_{\exists r_m.E_m}^x \\ C_R^x &= A_1^x \sqcap \dots \sqcap A_n^x \sqcap \exists r_1^x.E_1 \sqcap \dots \sqcap \exists r_m^x.E_m \end{aligned}$$

Moreover, set  $q'_0 := N(x)$ .

*Example 4.* Consider the OMQ  $Q = (\mathcal{T}_1, \Sigma_{\text{full}}, q(x, y))$  with  $\mathcal{T}_1$  as in Example 1 and let  $q(x, y)$  the following tqCQ:<sup>2</sup>

$$\begin{aligned} q(x, y) &= \exists z \text{GeneticRiskPatient}(x) \wedge \text{hasDisease}(x, y) \wedge \\ &\quad \text{Disease}(y) \wedge \text{hasDisease}(x, z) \wedge \text{Albinism}(z) \end{aligned}$$

We first remove quantified variables: all atoms that contain the variable  $z$  are replaced by  $A_{\exists \text{hasDisease.Albinism}}(y)$ , and the TBox is extended with the inclusion  $\exists \text{hasDisease.Albinism} \sqsubseteq A_{\exists \text{hasDisease.Albinism}}$ . We then construct  $\mathcal{T}'_1$ , which we give here only partially. The final concept inclusion in  $\mathcal{T}'_1$  is

$$\text{GeneticRiskPatient}^x \sqcap \text{Disease}^y \sqcap A_{\exists \text{hasDisease.Albinism}}^x \sqsubseteq N,$$

representing the updated query without role atoms; for example, the concept name  $\text{Disease}^y$  stands for the atom  $\text{Disease}(y)$ . Among others,  $\mathcal{T}'_1$  contains the further concept inclusions

$$\begin{aligned} \exists \text{hasDisease}^x.\text{HereditaryDisease} &\sqsubseteq A_{\exists \text{hasDisease.HereditaryDisease}}^x \\ \text{HereditaryDisease}^y &\sqsubseteq A_{\exists \text{hasDisease.HereditaryDisease}}^x \end{aligned}$$

where, intuitively, the lower concept inclusion captures that case that the truth of the concept  $\exists \text{hasDisease.HereditaryDisease}$  is witnessed at  $y$  (the role atom  $\text{hasDisease}(x, y)$  from  $q$  is only implicit here) while the upper concept inclusion deals with other witnesses.  $\dashv$

Before proving that the constructed OMQ  $Q'$  behaves in the desired way, we give some preliminaries. It is known that, if an OMQ from  $(\mathcal{EL}, \text{AQ})$  has an FO-rewriting, then it has a tUCQ-rewriting, see for example [6, 12]. A tqCQ  $q$  is *conformant* if it satisfies the following properties:

1. if  $A(x)$  is a concept atom, then either  $A$  is of the form  $B^y$  and  $x$  is the answer variable or  $A$  is not of this form and  $x$  is a quantified variable;
2. if  $r(x, y)$  is a role atom, then either  $r$  is of the form  $s^z$  and  $x$  is the answer variable or  $r$  is not of this form and  $x$  is a quantified variable.

A *conformant tUCQ* is then defined in the expected way. The notion of conformance captures what we informally described as never introducing atoms into the rewriting ‘without a reason’. By the following lemma, FO-rewritability of the OMQs constructed in our reduction implies conformant tUCQ-rewritability, that is, there is indeed no reason to introduce any of the atoms that are forbidden in conformant rewritings.

<sup>2</sup> We only use here that  $\mathcal{T}_1$  contains the concept  $\exists \text{hasDisease.HereditaryDisease}$  on the left-hand side of a concept inclusion.



**Lemma 2.** *Let  $Q$  be from  $(\mathcal{EL}, \text{tqCQ})$  and  $Q'$  the OMQ constructed from  $Q$  as above. If  $Q'$  is FO-rewritable, then it is rewritable into a conformant tUCQ.*

When started on an OMQ produced by our reduction, the algorithms presented in [12] and implemented in the Grind system produce a conformant tUCQ-rewriting. Indeed, this can be expected of any reasonable algorithm based on backwards chaining. Let  $q'$  be a conformant tUCQ-rewriting of  $Q'$ . The *corresponding UCQ for  $Q$*  is the UCQ  $q$  obtained by taking each CQ from  $q'$ , replacing every atom  $A^x(x_0)$  with  $A(x)$  and every atom  $r^x(x_0, y)$  with  $r(x, y)$ , and adding all atoms  $r(x, y)$  from  $q_0$  such that both  $x$  and  $y$  are answer variables. The answer variables in  $q$  are those of  $q_0$ . Observe that  $q$  is a union of tqCQs.

**Proposition 1.**  *$Q$  is FO-rewritable iff  $Q'$  is FO-rewritable. Moreover, if  $q'$  is a conformant tUCQ-rewriting of  $Q'$  and  $q$  the corresponding UCQ for  $Q$ , then  $q$  is a rewriting of  $Q$ .*

The proof strategy is to establish the ‘moreover’ part and to additionally show how certain UCQ-rewritings of  $Q$  can be converted into UCQ-rewritings of  $Q'$ . More precisely, a CQ  $q$  is a *derivative* of  $q_0$  if it results from  $q_0$  by exchanging atoms  $A(x)$  for  $\mathcal{EL}$ -concepts  $C$ , seen as tree-shaped CQs rooted in  $x$ . We are going to prove the following lemma in Section 4.

**Lemma 3.** *If an OMQ  $(\mathcal{T}, \Sigma, q_0)$  from  $(\mathcal{EL}, \text{tqCQ})$  is FO-rewritable, then it has a UCQ-rewriting in which each CQ is a derivative of  $q_0$ .*

Let  $q$  be a UCQ in which every CQ is a derivative of  $q_0$ . Then the *corresponding UCQ for  $Q'$*  is the UCQ  $q'$  obtained by taking each CQ from  $q$ , replacing every atom  $A(x)$ ,  $x$  answer variable, with  $A^x(x_0)$ , every atom  $r(x, y)$ ,  $x$  answer variable and  $y$  quantified variable, with  $r^x(x_0, y)$ , and deleting all atoms  $r(x_1, x_2)$ ,  $x_1, x_2$  answer variables. The answer variable in  $q'$  is  $x_0$ . Note that  $q'$  is a tUCQ. To establish the ‘only if’ direction of Proposition 1, we show that when  $q$  is a UCQ-rewriting of  $Q$  in which every CQ is a derivative of the query  $q_0$ , then the corresponding UCQ for  $Q'$  is a rewriting of  $Q'$ .

## 4 Rooted CQs

We consider OMQs based on rCQs, a strict generalization of tqCQs. In this case, we are not going to achieve a ‘black box’ reduction, but rely on a concrete algorithm for solving FO-rewritability in  $(\mathcal{EL}, \text{AQ})$ . This algorithm is a straightforward and not necessarily terminating backwards chaining algorithm or a (potentially terminating) refinement thereof, as implemented in the Grind system. We show how to combine the construction of (several) OMQs from  $(\mathcal{EL}, \text{AQ})$  with a modification of the assumed algorithm to decide FO-rewritability in  $(\mathcal{EL}, \text{rCQ})$  and to construct actual rewritings.

We start with introducing the straightforward backwards chaining algorithm mentioned above which we refer to as  $\text{bc}_{\text{AQ}}$ . Central to  $\text{bc}_{\text{AQ}}$  is a backwards chaining step based on concept inclusions in the TBox used in the OMQ. Let  $C$

and  $D$  be  $\mathcal{EL}$ -concepts,  $E \sqsubseteq F$  a concept inclusion, and  $x \in \text{var}(C)$  (where  $C$  is viewed as a tree-shaped CQ). Then  $D$  is *obtained from  $C$  by applying  $E \sqsubseteq F$  at  $x$*  if  $D$  can be obtained from  $C$  by

- removing  $A(x)$  for all concept names  $A$  with  $\models F \sqsubseteq A$ ;
- removing  $r(x, y)$  and the tree-shaped CQ  $G$  rooted at  $y$  when  $\models F \sqsubseteq \exists r.G$ ;
- adding  $A(x)$  for all concept names  $A$  that occur in  $E$  as a top-level conjunct (that is, that are not nested inside existential restrictions);
- adding  $\exists r.G$  as a CQ with root  $x$ , for each  $\exists r.G$  that is a top-level conjunct of  $E$ .

Let  $C$  and  $D$  be  $\mathcal{EL}$ -concepts. We write  $D \prec C$  if  $D$  can be obtained from  $C$  by removing an existential restriction (not necessarily on top level, and potentially resulting in  $D = \top$  when  $C$  is of the form  $\exists r.E$ ). We use  $\prec^*$  to denote the reflexive and transitive closure of  $\prec$  and say that  $D$  is  *$\prec$ -minimal with  $\mathcal{T} \models D \sqsubseteq A_0$*  if  $\mathcal{T} \models D \sqsubseteq A_0$  and there is no  $D' \prec D$  with  $\mathcal{T} \models D' \sqsubseteq A_0$ .

Now we are in the position to describe algorithm  $\text{bc}_{\text{AQ}}$ . It maintains a set  $M$  of  $\mathcal{EL}$ -concepts that represent tCQs. Let  $Q = (\mathcal{T}, \Sigma, A_0)$  be from  $(\mathcal{EL}, \text{AQ})$ . Starting from the set  $M = \{A_0\}$ , it exhaustively performs the following steps:

1. find  $C \in M$ ,  $x \in \text{var}(C)$ , a concept inclusion  $E \sqsubseteq F \in \mathcal{T}$ , and  $D$ , such that  $D$  is obtained from  $C$  by applying  $E \sqsubseteq F$  at  $x$ ;
2. find  $D' \prec^* D$  that is  $\prec$ -minimal with  $\mathcal{T} \models D' \sqsubseteq A_0$ , and add  $D'$  to  $M$ .

Application of these steps might not terminate. We use  $\text{bc}_{\text{AQ}}(Q)$  to denote the potentially infinitary UCQ  $\bigvee M|_{\Sigma}$  where  $M$  is the set obtained in the limit and  $q|_{\Sigma}$  denotes the restriction of the UCQ  $q$  to those disjuncts that only use symbols from  $\Sigma$ . Note that, in Point 2, it is possible to find the desired  $D'$  in polynomial time since the subsumption ' $\mathcal{T} \models D' \sqsubseteq A_0$ ' can be decided in polynomial time. The following is standard to prove, see [12, 15] and Lemma 5 below for similar results.

**Lemma 4.** *Let  $Q$  be an OMQ from  $(\mathcal{EL}, \text{AQ})$ . If  $\text{bc}_{\text{AQ}}(Q)$  is finite, then it is a UCQ-rewriting of  $Q$ . Otherwise,  $Q$  is not FO-rewritable.*

*Example 5.* Consider the TBox

$$\mathcal{T} = \{\text{Person} \sqcap \exists \text{hasParent}.\text{GeneticRiskPatient} \sqsubseteq \text{GeneticRiskPatient}\}$$

and let  $Q = (\mathcal{T}, \Sigma, \text{GeneticRiskPatient}(x))$  with  $\Sigma = \{\text{Person}, \text{GeneticRiskPatient}\}$ . Note that the role name `hasParent` does not occur in  $\Sigma$ . Even though the set  $M$  generated by  $\text{bc}_{\text{AQ}}$  (in the limit of its non-terminating run) is infinite,  $\text{bc}_{\text{AQ}}(Q) = \text{GeneticRiskPatient}(x)$  is finite and a UCQ-rewriting of  $Q$ .  $\dashv$

The algorithm for deciding FO-rewritability in  $(\mathcal{EL}, \text{AQ})$  presented in [12] and underlying the Grind system can be seen as a refinement of  $\text{bc}_{\text{AQ}}$ . Indeed, that algorithm always terminates and returns  $\bigvee M|_{\Sigma}$  if that UCQ is finite and reports non-FO-rewritability otherwise. Moreover, the UCQ-rewriting is represented in a decomposed way and output as a non-recursive Datalog program for efficiency

and succinctness. For our purposes, the only important aspect is that, when started on an FO-rewritable OMQ, it computes (a non-recursive Datalog program that is equivalent to) the UCQ-rewriting  $\bigvee M|_{\Sigma}$ .

We next introduce a generalized version  $\text{bc}_{\text{AQ}}^+$  of  $\text{bc}_{\text{AQ}}$  that takes as input an OMQ  $Q = (\mathcal{T}, \Sigma, A_0)$  from  $(\mathcal{EL}, \text{AQ})$  and an additional  $\mathcal{EL}$ -TBox  $\mathcal{T}^{\text{min}}$ , such that termination and output of  $\text{bc}_{\text{AQ}}^+$  agrees with that of  $\text{bc}_{\text{AQ}}$  when the input satisfies  $\mathcal{T}^{\text{min}} = \mathcal{T}$ . Starting from  $M = \{A_0\}$ , algorithm  $\text{bc}_{\text{AQ}}^+$  exhaustively performs the following steps:

1. find  $C \in M$ ,  $x \in \text{var}(C)$ , a concept inclusion  $E \sqsubseteq F \in \mathcal{T}$ , and  $D$ , such that  $D$  is obtained from  $C$  by applying  $E \sqsubseteq F$  at  $x$ ;
2. find  $D' \prec^* D$  that is  $\prec$ -minimal with  $\mathcal{T}^{\text{min}} \models D' \sqsubseteq A_0$ , and add  $D'$  to  $M$ .

We use  $\text{bc}_{\text{AQ}}^+(Q, \mathcal{T}^{\text{min}})$  to denote the potentially infinitary UCQ  $\bigvee M|_{\Sigma}$ ,  $M$  obtained in the limit. Note that  $\text{bc}_{\text{AQ}}^+$  uses the TBox  $\mathcal{T}$  for backwards chaining and  $\mathcal{T}^{\text{min}}$  for minimization while  $\text{bc}_{\text{AQ}}$  uses  $\mathcal{T}$  for both purposes. The refined version of  $\text{bc}_{\text{AQ}}$  implemented in the Grind system can easily be adapted to behave like a terminating version of  $\text{bc}_{\text{AQ}}^+$ .

Our aim is to convert an OMQ  $Q = (\mathcal{T}, \Sigma, q_0)$  from  $(\mathcal{EL}, \text{rCQ})$  into a set of pairs  $(Q', \mathcal{T}^{\text{min}})$  with  $Q'$  an OMQ from  $(\mathcal{EL}, \text{AQ})$  and  $\mathcal{T}^{\text{min}}$  an  $\mathcal{EL}$ -TBox such that  $Q$  is FO-rewritable iff  $\text{bc}_{\text{AQ}}^+(Q', \mathcal{T}^{\text{min}})$  terminates for all pairs  $(Q', \mathcal{T}^{\text{min}})$  and, moreover, if this is the case, then the resulting UCQ-rewritings can straightforwardly be converted into a rewriting of  $Q$ .

Let  $Q = (\mathcal{T}, \Sigma, q_0)$ . We construct one pair  $(Q_{q_r}, \mathcal{T}_{q_r}^{\text{min}})$  for each fork rewriting  $q_r$  of  $q_0$ . We use  $\text{core}(q_r)$  to denote the minimal set  $V$  of variables that contains all answer variables in  $q_r$  and such that after removing all atoms  $r(x, y)$  with  $x, y \in V$ , we obtain a disjoint union of tree-shaped CQs. We call these CQs the *trees in  $q_r$* . Intuitively, we separate the tree-shaped parts of  $q_r$  from the cyclic part, the latter identified by  $\text{core}(q_r)$ . This is similar to the definition of tqCQs where, however, cycles cannot involve any quantified variables. In a forest model of an ABox and a TBox as mentioned before Lemma 1, the variables in  $\text{core}(q_r)$  must be mapped to the ABox part of the model (rather than to the trees attached to it). Now  $(Q_{q_r}, \mathcal{T}_{q_r}^{\text{min}})$  is defined by setting  $Q_{q_r} = (\mathcal{T}_{q_r}, \Sigma_{q_r}, N(x))$  and

$$\begin{aligned} \mathcal{T}_{q_r} = & \mathcal{T} \cup \{C_R^x \sqsubseteq D_R^x \mid x \in \text{core}(q_r), C \sqsubseteq D \in \mathcal{T}\} \\ & \cup \left\{ \bigcap_{C(x) \text{ a tree in } q_r} C_R^x \sqsubseteq N \right\} \end{aligned}$$

where  $C_R^x$  is defined as in Section 3, and  $\Sigma_{q_r}$  is the extension of  $\Sigma$  with all concept names  $A^x$  and role names  $r^x$  used in  $\mathcal{T}_{q_r}$  such that  $A, r \in \Sigma$ .

It remains to define  $\mathcal{T}_{q_r}^{\text{min}}$ , which is  $\mathcal{T}_{q_r}$  extended with one concept inclusion for each fork rewriting  $q$  of  $q_0$  and each splitting  $\Pi = \langle R, S_1, \dots, S_\ell, r_1, \dots, r_\ell, \mu, \nu \rangle$  of  $q$  w.r.t.  $\mathcal{A}_{q_r}$ , as follows. For each  $x \in \text{avar}(q_r)$ , the equality atoms in  $q_r$  give rise to an equivalence class  $[x]_{q_r}$  of answer variables, defined in the expected way. We only consider the splitting  $\Pi$  of  $q$  if it preserves answer variables modulo

equality, that is, if  $x \in \text{avar}(q)$ , then there is a  $y \in [x]_{q_r}$  such that  $\nu(x) = y$ . We then add the inclusion

$$\left( \prod_{\substack{A(x) \in q \\ \text{with } x \in R}} A^{\nu(x)} \right) \sqcap \left( \prod_{1 \leq i \leq \ell} \exists r_i^{\nu(\mu(i))}. C_{q|S_i} \right) \sqsubseteq N$$

It can be shown that, summing up over all fork rewritings and splittings, only polynomially many concepts  $\exists r_i^{\nu(\mu(i))}. C_{q|S_i}$  are introduced (this is similar to the proof of Lemma 6 in [17]). Note that we do not introduce fresh concept names of the form  $A_{\exists r.C}^x$  as in Section 3. This is not necessary here because of the use of fork rewritings and splittings in  $\mathcal{T}^{\min}$ .

*Example 6.* Consider query  $q_3$  from Example 3 and TBox  $\mathcal{T}_1$  from Example 1. Constructing  $\mathcal{T}_{q_3}$  (thus considering  $q_3$  as a fork rewriting of itself) would add concept inclusions like

$$\text{Person}^x \sqcap \exists \text{hasDisease}^x. \text{HereditaryDisease} \sqsubseteq \text{GeneticRiskPatient}^x$$

The final concept inclusion added is the following, listing concepts needed at  $x, y_1, y_2$ , and  $z$  that result in a match of  $q_3$ :

$$\text{Person}^x \sqcap \text{MelaminDeficiency}^{y_1} \sqcap \text{ImpairedVision}^{y_2} \sqcap \text{GeneDefect}^z \sqsubseteq N$$

When building the TBox  $\mathcal{T}_{q_3}^{\min}$ , it is necessary to look for matches of  $q_3$  by a splitting  $\Pi$  of a fork rewriting of  $q_3$  w.r.t.  $\mathcal{A}_{q_3}$  and  $\mathcal{T}_1$ . We consider here the splitting  $\Pi = \langle R, S_1, r_1, \mu, \nu \rangle$  of the fork rewriting  $q'_3$  of  $q_3$  given in Example 3, defined by setting

$$R = \{x\}, S_1 = \{y_1, z\}, r_1 = \text{hasDisease}, \mu(1) = x, \nu = (x \mapsto x)$$

For  $\Pi$ , the following concept inclusion is added to  $\mathcal{T}_{q_3}^{\min}$ :

$$\text{Person}^x \sqcap \exists \text{hasDisease}^x. (\text{MelaminDeficiency} \sqcap \text{ImpairedVision} \sqcap \text{causedBy.GeneDefect}) \sqsubseteq N \quad \dashv$$

It can be seen that when  $\text{bc}_{\text{AQ}}^+(Q_{q_r}, \mathcal{T}_{q_r}^{\min})$  is finite, then it is a conformant tUCQ in the sense of Section 3. Thus, we can also define a *corresponding UCQ*  $q$  for  $Q$  as in that section, that is,  $q$  is obtained by taking each CQ from  $q'$ , replacing every atom  $A^x(x_0)$  with  $A(x)$  and every atom  $r^x(x_0, y)$  with  $r(x, y)$ , and adding all atoms  $r(x, y)$  from  $q_r$  such that  $x, y \in \text{core}(q_r)$ . The answer variables in  $q$  are those of  $q_0$ .

**Proposition 2.** *Let  $Q = (\mathcal{T}, \Sigma, q_0)$  be an OMQ from  $(\mathcal{EL}, rCQ)$ . If  $\text{bc}_{\text{AQ}}^+(Q_{q_r}, \mathcal{T}_{q_r}^{\min})$  is finite for all fork rewritings  $q_r$  of  $q_0$ , then  $\bigvee_{q_r} \hat{q}_{q_r}$  is a UCQ-rewriting of  $Q$ , where  $\hat{q}_{q_r}$  is the UCQ for  $Q$  that corresponds to  $\text{bc}_{\text{AQ}}^+(Q_{q_r}, \mathcal{T}_{q_r}^{\min})$ . Otherwise,  $Q$  is not FO-rewritable.*

To prove Proposition 2, we introduce a backwards chaining algorithm  $\text{bc}_{\text{rCQ}}$  for computing UCQ-rewritings of OMQs from  $(\mathcal{E}\mathcal{L}, \text{rCQ})$  that we refer to as  $\text{bc}_{\text{rCQ}}$ . In a sense,  $\text{bc}_{\text{rCQ}}$  is the natural generalization of  $\text{bc}_{\text{AQ}}$  to rCQs. We then show a correspondence between the run of  $\text{bc}_{\text{rCQ}}$  on the input OMQ  $Q$  from  $(\mathcal{E}\mathcal{L}, \text{rCQ})$  and the runs of  $\text{bc}_{\text{AQ}}^+$  on the constructed inputs of the form  $(Q_{q_r}, \mathcal{T}_{q_r}^{\text{min}})$ .

On the way, we also provide the missing proof for Lemma 3, which in fact is a consequence of the correctness of  $\text{bc}_{\text{rCQ}}$  (states as Lemma 5 in the appendix) and the observation that, when  $Q = (\mathcal{T}, \Sigma, q_0)$  is from  $(\mathcal{E}\mathcal{L}, \text{tqCQ})$ , then  $\text{bc}_{\text{rCQ}}(Q)$  contains only derivatives of  $q_0$ . The latter is due to the definition of the  $\text{bc}_{\text{rCQ}}$  algorithm, which starts with a set of minimized fork rewritings of  $q_0$ , and the fact that the only fork rewriting of a tqCQ is the query itself.

There are two exponential blowups in the presented approach. First, the number of fork rewritings of  $q_0$  might be exponential in the size of  $q_0$ . We expect this not to be a problem in practice since the number of fork rewritings of realistic queries should be fairly small. And second, the number of splittings can be exponential and thus the same is true for the size of each  $\mathcal{T}_{q_r}^{\text{min}}$ . We expect that also this blowup will be moderate in practice. Moreover, in an optimized implementation one would not represent  $\mathcal{T}_{q_r}^{\text{min}}$  as a TBox, but rather check the existence of fork rewritings and splittings that give rise to concept inclusions in  $\mathcal{T}_{q_r}^{\text{min}}$  in a more direct way. This involves checking whether concepts of the form  $\exists r_i^{\nu(\mu(i))}.C_{q'_i|S_i}$  are derived, and the fact that there are only polynomially many different such concepts should thus be very relevant regarding performance.

## 5 Experiments

We have extended the *Grind* system [12] to support OMQs from  $(\mathcal{E}\mathcal{L}, \text{tqCQ})$  and  $(\mathcal{E}\mathcal{L}, \text{rCQ})$  instead of only from  $(\mathcal{E}\mathcal{L}, \text{AQ})$ , and conducted experiments with real-world ontologies and hand-crafted conjunctive queries. The system can be downloaded from <http://www.cs.uni-bremen.de/~hansen/grind>, together with the ontologies and queries, and is released under GPL. It outputs rewritings in the form of non-recursive Datalog queries. We have implemented the following optimization: given  $Q = (\mathcal{T}, \Sigma, q_0)$ , first compute all fork rewritings of  $q_0$ , rewrite away all variables outside of the core (in the same way in which tree parts of the query are removed in Section 3) to obtain a new OMQ  $(\mathcal{T}', \Sigma, q'_0)$ , and then test for each atom  $A(x) \in q'_0$  whether  $(\mathcal{T}', \Sigma, A(x))$  is FO-rewritable. It can be shown that, if this is the case, then  $Q$  is FO-rewritable, and it is also possible to transfer the actual rewritings. If this check fails, we go through the full construction described in the paper.

Experiments were carried out on a Linux (3.2.0) machine with a 3.5 GHz quad-core processor and 8 GB of RAM. For the experiments, we use (the  $\mathcal{E}\mathcal{L}$  part of) the ontologies ENVO, FBbi, SO, MOHSE, and not-galen. The first three ontologies are from the biology domain, and are available through Biportal<sup>3</sup>.

<sup>3</sup> <https://biportal.bioontology.org>

TBox	CI	CN	RN	Min CQ	Avg CQ	Max CQ	Avg AQ	Aborts
ENVO	1942	1558	7	0.2s	1.5s	7s	1s	0
FBbi	567	517	1	0.05s	0.5s	3s	0.3s	0
MOHSE	3665	2203	71	2s	10s	40s	6s	0
not-galen	4636	2748	159	6s	9s	28s	25s	2
SO	3160	2095	12	1s	19s	2m23s	4s	1

**Table 1.** TBox information and results of experiments

$q_1(x, y) = \text{Patient}(x) \wedge \text{shows}(x, y) \wedge \text{Endocarditis}(y)$
$q_2(w, x, y, z) = \text{Doctor}(w) \wedge \text{hasPersonPerforming}(x, w) \wedge \text{Surgery}(x) \wedge$ $\text{actsOn}(x, y) \wedge \text{Tissue}(y) \wedge \text{actsOn}(x, z) \wedge$ $\text{InternalOrgan}(z) \wedge \text{hasAlphaConnection}(y, z)$
$q_7(x) = \exists y, z \text{ Protein}(x) \wedge \text{contains}(x, y) \wedge \text{Tetracycline}(y) \wedge$ $\text{InternalOrgan}(z) \wedge \text{isActedOnSpecificallyBy}(z, y)$
$q_8(x) = \exists v, w, y, z \text{ Sulphonamide}(v) \wedge \text{serves}(v, w) \wedge \text{TumorMarkerRole}(w) \wedge$ $\text{NamedEnzyme}(x) \wedge \text{serves}(x, w) \wedge \text{actsOn}(x, z) \wedge \text{Liver}(z) \wedge$ $\text{TeichoicAcid}(y) \wedge \text{actsOn}(y, z)$
$q_{10}(x) = \exists y, z \text{ BodyStructure}(x) \wedge \text{isBetaConnectionOf}(x, y) \wedge \text{Brain}(y) \wedge$ $\text{IntrinsicallyNormalBodyStructure}(z) \wedge \text{isBetaConnectionOf}(z, y)$

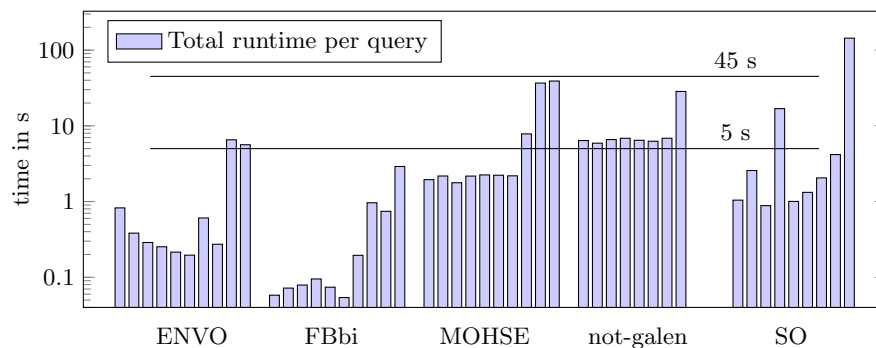
**Fig. 1.** Exemplary queries used for experiments with TBox not-galen.

MOHSE and not-galen are different versions of the GALEN ontology<sup>4</sup>, which describes medical terms. Some statistics is given in Table 1, namely the number of concept inclusions (CI), concept names (CN), and role names (RN) in each ontology. For each ontology, we hand-crafted 10 conjunctive queries (three tqCQs and seven rCQs), varying in size from 2 to 5 variables and showing several different topologies (see Fig. 1 for a sample).

The runtimes are reported in Table 1. Only three queries did not terminate in 30 minutes or exhausted the memory. For the successful ones, we list fastest (Min CQ), slowest (Max CQ), and average runtime (Avg CQ). For comparison, the Avg AQ column lists the time needed to compute FO-rewritings for all queries  $(\mathcal{T}, \Sigma, A(x))$  with  $A(x)$  an atom in  $q_0$ . This check is of course incomplete for FO-rewritability of  $Q$ , but can be viewed as a lower bound. A detailed picture of individual runtimes is given in Figure 2.

In summary, we believe that the outcome of our experiments is promising. While runtimes are higher than in the AQ case, they are still rather small given that we are dealing with an intricate static analysis task and that many parts of our system have not been seriously optimized. The queries with long runtimes or

<sup>4</sup> <http://www.opengalen.org/>



**Fig. 2.** Runtimes for individual OMQs, showing only non-aborting runs.

timeouts contain AQs that are not FO-rewritable which forces the decomposed algorithm implemented in Grind to enter a more expensive processing phase.

## 6 Conclusion

We remark that our approach can also be used to compute FO-rewritings of OMQs from  $(\mathcal{EL}, \text{CQ})$  even if the CQs are not rooted, as long as they are not Boolean (that is, as long as they contain at least one answer variable) and an algorithm for query containment in  $(\mathcal{EL}, \text{CQ})$  is also available. This follows from (a minor variation of) an observation from [5]: FO-rewritability of non-Boolean OMQs from  $(\mathcal{EL}, \text{CQ})$  can be polynomially reduced to a combination of containment in  $(\mathcal{EL}, \text{CQ})$  and FO-rewritability in  $(\mathcal{EL}, \text{rCQ})$ . As future work, it would be interesting to extend our approach to UCQs, to the extension of  $\mathcal{EL}$  with role hierarchies and domain and range restrictions, or even to the extension  $\mathcal{ELI}$  of  $\mathcal{EL}$  with inverse roles.

**Acknowledgements.** We acknowledge support by ERC grant 647289 ‘CODA’.

## References

1. Artale, A., Calvanese, D., Kontchakov, R., Zakharyashev, M.: The DL-Lite family and relations. *J. Artif. Intell. Res.* 36, pp. 1–69 (2009)
2. Baader, F., Horrocks, I., Lutz, C., Sattler, U.: An introduction to description logics. Cambridge University Press (2017)
3. Barceló, P., Berger, G., Pieris, A.: Containment for rule-based ontology-mediated queries (19 April 2017), available at <https://arxiv.org/abs/1703.07994> [cs.DB]
4. Bienvenu, M., ten Cate, B., Lutz, C., Wolter, F.: Ontology-based data access: A study through disjunctive datalog, CSP, and MMSNP. *J. ACM Trans. Database Syst.* 39(4), pp. 33:1–33:44 (2014)
5. Bienvenu, M., Hansen, P., Lutz, C., Wolter, F.: First order-rewritability and containment of conjunctive queries in Horn description logics. In: *Proc. of IJCAI*, pp. 965–971 (2016)

6. Bienvenu, M., Lutz, C., Wolter, F.: First order-rewritability of atomic queries in Horn description logics. In: Proc. of IJCAI, pp. 754–760 (2013)
7. Calvanese, D., De Giacomo, G., Lembo, D., Lenzerini, M., Rosati, R.: Tractable reasoning and efficient query answering in description logics: The DL-Lite family. *J. Autom. Reasoning* 39(3), pp. 385–429 (2007)
8. Deutsch, A., Popa, L., Tannen, V.: Physical data independence, constraints, and optimization with universal plans. In: Proc. of VLDB, pp. 459–470 (1999)
9. Eiter, T., Ortiz, M., Simkus, M., Tran, T., Xiao, G.: Query rewriting for Horn-SHIQ plus rules. In: Proc. of AAAI (2012)
10. Feier, C., Lutz, C., Kuusisto, A.: Rewritability in monadic disjunctive datalog, MMSNP, and expressive description logics. In: Proc. of ICDT (2017)
11. Gottlob, G., Kikot, S., Kontchakov, R., Podolskii, V.V., Schwentick, T., Zakharyashev, M.: The price of query rewriting in ontology-based data access. *J. Artif. Intell.* 213, pp. 42–59 (2014)
12. Hansen, P., Lutz, C., Seylan, I., Wolter, F.: Efficient query rewriting in the description logic EL and beyond. In: Proc. of IJCAI, pp. 3034–3040 (2015)
13. Kikot, S., Kontchakov, R., Zakharyashev, M.: Conjunctive query answering with OWL 2 QL. In: Proc. of KR (2012)
14. König, M., Leclère, M., Mugnier, M.: Query rewriting for existential rules with compiled preorder. In: Proc. of IJCAI, pp. 3106–3112 (2015)
15. König, M., Leclère, M., Mugnier, M., Thomazo, M.: Sound, complete and minimal UCQ-rewriting for existential rules. *Semantic Web* 6(5), pp. 451–475 (2015)
16. Lutz, C.: The complexity of conjunctive query answering in expressive description logics. In: Proc. of IJCAR, pp. 179–193 (2008)
17. Lutz, C.: Two upper bounds for conjunctive query answering in SHIQ. In: Proc. of DL (2008)
18. Lutz, C., Toman, D., Wolter, F.: Conjunctive query answering in the description logic EL using a relational database system. In: Proc. of IJCAI, pp. 2070–2075 (2009)
19. Lutz, C., Wolter, F.: Non-uniform data complexity of query answering in description logics. In: Proc. of KR (2012)
20. Motik, B., Cuenca Grau, B., Horrocks, I., Wu, Z., Fokoue, A., Lutz, C.: OWL 2 Web Ontology Language: Profiles. W3C recommendation (11 December 2012), available at <http://www.w3.org/TR/owl2-profiles/>
21. Pérez-Urbina, H., Motik, B., Horrocks, I.: Tractable query answering and rewriting under description logic constraints. *J. Applied Logic* 8(2), pp. 186–209 (2010)
22. Rodríguez-Muro, M., Calvanese, D.: High performance query answering over DL-Lite ontologies. In: Proc. of KR (2012)
23. Rosati, R., Almatelli, A.: Improving query answering over DL-Lite ontologies. In: Proc. of KR (2010)
24. Stefanoni, G., Motik, B.: Answering conjunctive queries over EL knowledge bases with transitive and reflexive roles. In: Proc. of AAAI, pp. 1611–1617 (2015)
25. Stefanoni, G., Motik, B., Horrocks, I.: Small datalog query rewritings for EL. In: Proc. of DL (2012)
26. Trivela, D., Stoilos, G., Chortaras, A., Stamou, G.B.: Optimising resolution-based rewriting algorithms for OWL ontologies. *J. Web Sem.* 33, pp. 30–49 (2015)