

Conjunctive Query Answering with *OWL 2 QL*

Stanislav Kikot and Roman Kontchakov and Michael Zakharyashev

Department of Computer Science and Information Systems
Birkbeck, University of London, U.K.
{kikot, roman, michael}@dcs.bbk.ac.uk

Abstract

We present a novel rewriting technique for conjunctive query answering over *OWL 2 QL* ontologies. In general, the obtained rewritings are not necessarily correct and can be of exponential size in the length of the query. We argue, however, that in most, if not all, practical cases the rewritings are correct and of polynomial size. Moreover, we prove some sufficient conditions, imposed on queries and ontologies, that guarantee correctness and succinctness. We also support our claim by experimental results.

Introduction

OWL 2 QL, one of three profiles of the Web Ontology Language *OWL 2*, was designed with the aim of supporting ontology-based data access (OBDA). The key idea is that data, ‘stored in a standard relational database management system (RDBMS), can be queried through an *OWL 2 QL* ontology via a simple rewriting mechanism, i.e., by rewriting the query into an SQL query that is then answered by the RDBMS, without any changes to the data’ (www.w3.org/TR/owl2-profiles). The rewritability property ensures, in particular, that the data complexity of answering queries over *OWL 2 QL* ontologies matches the complexity of database query answering, which is in AC^0 .

It has been observed, however, that the available ‘rewriting mechanisms’ for *OWL 2 QL* (Calvanese et al. 2007a; Pérez-Urbina, Motik, and Horrocks 2009; Rosati and Almatelli 2010; Chortaras, Trivela, and Stamou 2011; Gottlob, Orsi, and Pieris 2011) are actually not so ‘simple.’ In fact, the rewritten queries are often too long to be executed by modern RDBMSs, and the question whether ‘short’ rewritings exist has attracted considerable attention over the last two years. For example, Kikot, Kontchakov, and Zakharyashev (2011) showed that no polynomial algorithm can construct a ‘pure rewriting’ of a conjunctive query (CQ) q over an *OWL 2 QL* ontology \mathcal{T} . Here by a *pure rewriting* we mean any first-order (FO) rewriting with the same signature (predicates and constants) as q and \mathcal{T} , possibly with equality. On the other hand, Gottlob and Schwenck (2011) gave a polynomial-time (‘impure’) rewriting using addi-

tional constants and predicates. Optimised (but still exponential) pure rewritings to nonrecursive Datalog were suggested by Rosati and Almatelli (2010) and Gottlob, Orsi, and Pieris (2011). The combined approach of Kontchakov et al. (2010) was developed for *OWL 2 QL* without role inclusions; it uses a simple polynomial rewriting over the data expanded by applying the ontology axioms and introducing a small number of new individuals.

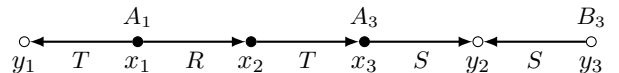
The diversity of approaches to query rewriting prompts another question: what is the type/shape/size of rewritings we should aim at to make OBDA with *OWL 2 QL* efficient? When trying to answer this question, we should bear in mind that (i) the OBDA paradigm relies on the proven efficiency of RDBMSs, but (ii) database query answering is not tractable in the size of queries (PSPACE-complete for FO queries and NP-complete for CQs). High efficiency of RDBMSs in practice appears to indicate only that answering real-world queries over real-world databases turns out to be tractable. As rewritings can turn a standard query to something ‘out of this world,’ a first rule of thumb could be as follows: *the rewritten query should look similar to the original one*. In this respect, as Gottlob and Schwenck (2011) remark, their polynomial rewriting is of rather ‘theoretical nature’ (it uses the extra constants, predicates and existentially quantified variables to encode, by making nondeterministic guesses, a relevant part of the chase, aka the canonical model; see the discussion in Conclusions for more details).

The aim of this paper is to investigate what causes exponentially long pure rewritings of CQs over *OWL 2 QL* ontologies and to check experimentally whether those ‘bad guys’ occur in real-world queries and ontologies. As a result of this analysis, we suggest some short (polynomial-size) rewritings that cover most, if not all, practical cases.

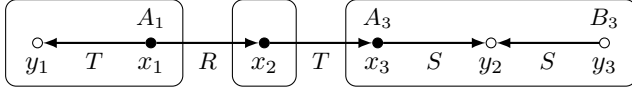
We think of a CQ as a labelled directed multigraph. For example, the query ‘find x_1, x_2, x_3 for which

$$\exists y_1, y_2, y_3 (A_1(x_1) \wedge A_3(x_3) \wedge B_3(y_3) \wedge T(x_1, y_1) \wedge R(x_1, x_2) \wedge T(x_2, x_3) \wedge S(x_3, y_2) \wedge S(y_3, y_2))$$

holds’ can be represented as the graph



To answer a CQ $q(\vec{x})$ over data \mathcal{A} and ontology \mathcal{T} , we seek homomorphisms from q to the structure, called the canonical model and obtained by expanding \mathcal{A} (extensional data) with knowledge in \mathcal{T} (intensional data). Thus, we have to find possible cuts of q into a number of pieces: some of them—say, of type (A)—are mapped to individuals in \mathcal{A} , while the others—of type (B)—may only have some of their terms (‘roots’) in \mathcal{A} , whereas the remaining part is implied by the knowledge in \mathcal{T} . For example, the query above can be cut, for a suitable \mathcal{T} , into 3 pieces of which only the middle one is of type (A), and the right piece has two roots x_3, y_3 :



If $q(\vec{x})$ does not have existentially quantified variables then the whole q forms the only possible cut of type (A), because the answer variables \vec{x} must be mapped to individuals in \mathcal{A} . If \mathcal{T} does not contain role inclusion axioms then, as shown by Kontchakov et al. (2010; 2011), every root determines a unique piece of type (B). However, in general, q may have exponentially many pieces of type (B). One can encode the intuition above as a pure positive existential rewriting q_e , which is, roughly speaking, a disjunction over all possible cuts of q , and so is exponential in $|q|$ in the worst case.

A slightly different approach to checking possible cuts of q is to consider, for every edge in q , whether it belongs to an (A) piece, or generates a (B) piece itself, or lies inside the (B) piece generated by some other edge. This ‘local’ view gives another rewriting, a conjunction of disjunctions q_c , which may still be exponential as the same edge may generate exponentially many distinct (B) pieces. (Two (B) pieces are different if their domains or roots are different.) Moreover, the new rewriting is not necessarily correct because some (B) pieces may not be realised together in the canonical model; we call such pieces conflicting.

We analyse—both theoretically and experimentally—conditions under which the number of (B) pieces is polynomial and they are not in conflict with each other. In particular, we develop techniques to ensure that rewritings are not affected by conflicting pieces. For example, one simple—but impure—approach involves a single fresh constant to represent all intensional objects in the canonical model (labelled nulls in the chase), but no extra variables.

We show that the number of (B) pieces generated by one edge in the query q is largely determined by a sophisticated interaction between role inclusions and inverse roles in the ontology \mathcal{T} , which can produce canonical models with very complex intensional parts. We give a sufficient condition (on q and \mathcal{T}) which guarantees that one edge in q can generate at most one piece of type (B) (though the number of ways this piece can be matched in the canonical model can still be exponential). This leads to our shortest pure rewriting, q_p , which can be constructed in polynomial time, but is correct only if q and \mathcal{T} satisfy the sufficient condition, which can also be checked in polynomial time. Trivial examples where this condition holds are ontologies without inverse roles, ontologies without role inclusions and qualified existential quantification, or those without positive occurrences

of existential quantifiers. Our experiments with a number of standard OWL 2 QL ontologies and queries demonstrate that the sufficient conditions always apply, the queries normally contain very few pieces of type (B), if any, and moreover, these pieces are never in conflict. Thus, in practice the rewritings q_c and q_p are both short and correct.

Omitted proofs can be found in the full version of the paper at www.dcs.bbk.ac.uk/~kikot.

OWL 2 QL

The language of OWL 2 QL contains *individual names* a_i , *concept names* A_i , and *role names* P_i ($i \geq 1$). *Roles* R , *basic concepts* B and *concepts* C are defined by the grammar:

$$\begin{aligned} R & ::= P_i \mid P_i^-, \\ B & ::= \perp \mid A_i \mid \exists R, \\ C & ::= B \mid \exists R.B \end{aligned}$$

A *TBox*, \mathcal{T} , is a finite set of *inclusions* of the form

$$\begin{aligned} B & \sqsubseteq C, & R_1 & \sqsubseteq R_2, \\ B_1 \sqcap B_2 & \sqsubseteq \perp, & R_1 \sqcap R_2 & \sqsubseteq \perp. \end{aligned}$$

(Note that concepts of the form $\exists R.B$ can only occur in the right-hand side of concept inclusions in OWL 2 QL. An inclusion $B' \sqsubseteq \exists R.B$ can be regarded as an abbreviation for three inclusions: $B' \sqsubseteq \exists R_B$, $\exists R_B \sqsubseteq B$ and $R_B \sqsubseteq R$, where R_B is a fresh role name.) An *ABox*, \mathcal{A} , is a finite set of *assertions* of the form $A_k(a_i)$ and $P_k(a_i, a_j)$. \mathcal{T} and \mathcal{A} together constitute the *knowledge base* (KB) $\mathcal{K} = (\mathcal{T}, \mathcal{A})$. The semantics for OWL 2 QL is defined in the usual way based on interpretations $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$; consult (Baader et al. 2003) for details. The set of individual names in \mathcal{A} will be denoted by $\text{ind}(\mathcal{A})$. For concepts or roles E_1 and E_2 , we write $E_1 \sqsubseteq_{\mathcal{T}} E_2$ if $\mathcal{T} \models E_1 \sqsubseteq E_2$; and we set $[E] = \{E' \mid E \sqsubseteq_{\mathcal{T}} E' \text{ and } E' \sqsubseteq_{\mathcal{T}} E\}$.

A *conjunctive query* (CQ) $q(\vec{x})$ is a first-order formula $\exists \vec{y} \varphi(\vec{x}, \vec{y})$, where φ is constructed, using \wedge , from atoms of the form $A_k(t_1)$ and $P_k(t_1, t_2)$, where each t_i is a *term* (an individual or a variable from \vec{x} or \vec{y}). Variables in \vec{x} are called *answer variables*, and those in \vec{y} *bound variables*. A tuple $\vec{a} \subseteq \text{ind}(\mathcal{A})$ is a *certain answer* to $q(\vec{x})$ over $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ if $\mathcal{I} \models q[\vec{a}]$ for all models \mathcal{I} of \mathcal{K} ; in this case we write $\mathcal{K} \models q[\vec{a}]$. To simplify notation, we will often identify q with the set of its atoms and use $P^-(t, t') \in q$ as a synonym of $P(t', t) \in q$; $\text{term}(q)$ is the set of terms in q . We call q *tree-shaped* if its primal graph ($\text{term}(q), \{\{t, t'\} \mid R(t, t') \in q\}$) is a tree.

Remark 1 Although the official OWL 2 QL contains the concept \top (for the whole domain), we do not consider it here as \top makes OBDA with OWL 2 QL *domain dependent* (Abiteboul, Hull, and Vianu 1995): take, for example, the query $A(x)$ over the ontology $\{\top \sqsubseteq A\}$. (But we shall use \top as an auxiliary symbol as it is convenient to regard $\exists R$ as an abbreviation for $\exists R.\top$.) To simplify presentation, we omit data properties and (ir)reflexivity constraints for roles.

Query answering over OWL 2 QL KBs is based on the fact that, for any consistent KB $\mathcal{K} = (\mathcal{T}, \mathcal{A})$, there is an interpretation $\mathcal{C}_{\mathcal{K}}$ such that, for all CQs $q(\vec{x})$ and $\vec{a} \subseteq \text{ind}(\mathcal{A})$,

we have $\mathcal{K} \models q[\vec{a}]$ iff $\mathcal{C}_{\mathcal{K}} \models q[\vec{a}]$. The interpretation $\mathcal{C}_{\mathcal{K}}$, called the *canonical model* of \mathcal{K} , can be constructed as follows. For each pair $[R], [B]$ with $\exists R.B$ in \mathcal{T} (recall that $\exists R.T$ is another way of writing $\exists R$), we introduce a fresh symbol $w_{[RB]}$ and call it the *witness for $\exists R.B$* . We write $\mathcal{K} \models C(w_{[RB]})$ if $\exists R^- \sqsubseteq_{\mathcal{T}} C$ or $B \sqsubseteq_{\mathcal{T}} C$. Define a *generating relation*, \rightsquigarrow , on the set of these witnesses together with $\text{ind}(\mathcal{A})$ by taking:

- $a \rightsquigarrow w_{[RB]}$ if $a \in \text{ind}(\mathcal{A})$, $[R]$ and $[B]$ are $\sqsubseteq_{\mathcal{T}}$ -minimal such that $\mathcal{K} \models \exists R.B(a)$ and there is no $b \in \text{ind}(\mathcal{A})$ with $\mathcal{K} \models R(a, b) \wedge B(b)$;
- $w_{[R'B']}\rightsquigarrow w_{[RB]}$ if $u \rightsquigarrow w_{[R'B']}$, for some u , $[R]$ and $[B]$ are $\sqsubseteq_{\mathcal{T}}$ -minimal such that $\mathcal{K} \models \exists R.B(w_{[R'B']})$ and it is not the case that $R' \sqsubseteq_{\mathcal{T}} R^-$ and $\mathcal{K} \models B(u)$.

If $a \rightsquigarrow w_{[R_1B_1]} \rightsquigarrow \dots \rightsquigarrow w_{[R_nB_n]}$, $n \geq 0$, then we say that a *generates* the path $aw_{[R_1B_1]} \dots w_{[R_nB_n]}$. Denote by $\text{path}_{\mathcal{K}}(a)$ the set of paths generated by a , and by $\text{tail}(\pi)$ the last element in $\pi \in \text{path}_{\mathcal{K}}(a)$. $\mathcal{C}_{\mathcal{K}}$ is defined by taking:

$$\Delta^{\mathcal{C}_{\mathcal{K}}} = \bigcup_{a \in \text{ind}(\mathcal{A})} \text{path}_{\mathcal{K}}(a), \quad a^{\mathcal{C}_{\mathcal{K}}} = a, \text{ for } a \in \text{ind}(\mathcal{A}),$$

$$A^{\mathcal{C}_{\mathcal{K}}} = \{\pi \in \Delta^{\mathcal{C}_{\mathcal{K}}} \mid \mathcal{K} \models A(\text{tail}(\pi))\},$$

$$P^{\mathcal{C}_{\mathcal{K}}} = \{(a, b) \in \text{ind}(\mathcal{A}) \times \text{ind}(\mathcal{A}) \mid \mathcal{K} \models P(a, b)\} \cup \\ \{(\pi, \pi \cdot w_{[RB]}) \mid \text{tail}(\pi) \rightsquigarrow w_{[RB]}, R \sqsubseteq_{\mathcal{T}} P\} \cup \\ \{(\pi \cdot w_{[RB]}, \pi) \mid \text{tail}(\pi) \rightsquigarrow w_{[RB]}, R \sqsubseteq_{\mathcal{T}} P^-\}.$$

$\Delta^{\mathcal{C}_{\mathcal{K}}} \setminus \text{ind}(\mathcal{A})$ is called the *tree part* of $\mathcal{C}_{\mathcal{K}}$. The following result is proved in a standard way (Kontchakov et al. 2010):

Theorem 2 For every OWL 2 QL KB $\mathcal{K} = (\mathcal{T}, \mathcal{A})$, every CQ $q(\vec{x})$ and every $\vec{a} \subseteq \text{ind}(\mathcal{A})$, $\mathcal{K} \models q[\vec{a}]$ iff $\mathcal{C}_{\mathcal{K}} \models q[\vec{a}]$.

Given a TBox \mathcal{T} , define a KB $\mathcal{K}_{\mathcal{T}} = (\mathcal{T}, \mathcal{A}_{\mathcal{T}})$ by taking

$$\mathcal{A}_{\mathcal{T}} = \{R(a_{RB}, b_{RB}), B(b_{RB}) \mid \\ \exists R.B \text{ is a } \mathcal{T}\text{-consistent concept in } \mathcal{T}\}.$$

In other words, $\mathcal{C}_{\mathcal{K}_{\mathcal{T}}}$ is the disjoint union of the canonical models for consistent $(\mathcal{T}, \{R(a_{RB}, b_{RB}), B(b_{RB})\})$. To simplify notation, we use $\mathcal{C}_{\mathcal{T}}$ in place of $\mathcal{C}_{\mathcal{K}_{\mathcal{T}}}$. We extend the generating relation \rightsquigarrow in $\mathcal{C}_{\mathcal{T}}$ by adding to it the pair $a_{RB} \rightsquigarrow b_{RB}$. The *RB-subtree* of $\mathcal{C}_{\mathcal{T}}$ has root a_{RB} and consists of the full subtree of $\mathcal{C}_{\mathcal{T}}$ with root b_{RB} extended with the edge (a_{RB}, b_{RB}) . (Note that there may be other branches starting from a_{RB} in $\mathcal{C}_{\mathcal{T}}$, for example, if $\exists R \sqsubseteq_{\mathcal{T}} \exists S$ and $R \not\sqsubseteq_{\mathcal{T}} S$.)

In the next section, we introduce the naïve exponential rewriting of CQs over OWL 2 QL ontologies sketched in the introduction, and analyse whether it can be made shorter.

Two Rewritings

Let $q(\vec{x})$ be a CQ and \mathcal{T} an OWL 2 QL ontology. Our task is to construct an FO query $q'(\vec{x})$, using the predicates and individuals of q , \mathcal{T} and $=$, such that, for any ABox \mathcal{A} and any $\vec{a} \subseteq \text{ind}(\mathcal{A})$, we have $(\mathcal{T}, \mathcal{A}) \models q[\vec{a}]$ iff $\mathcal{I}_{\mathcal{A}} \models q'[\vec{a}]$, where $\mathcal{I}_{\mathcal{A}}$ is the interpretation with domain $\text{ind}(\mathcal{A})$ given by the atoms in \mathcal{A} . Such a query q' is called a (*pure*) *rewriting* of

q and \mathcal{T} . As argued in the introduction, we are interested in rewritings q' that are (a) as short as possible (ideally, polynomial in $|q|$ and $|\mathcal{T}|$), (b) look similar to the original CQ q (in particular, built from the same predicates and terms as q and \mathcal{T}), and (c) can be constructed in reasonable time.

Without loss of generality, we will assume that the primal graph of q is connected. (If this is not the case, we consider each of the connected components separately.)

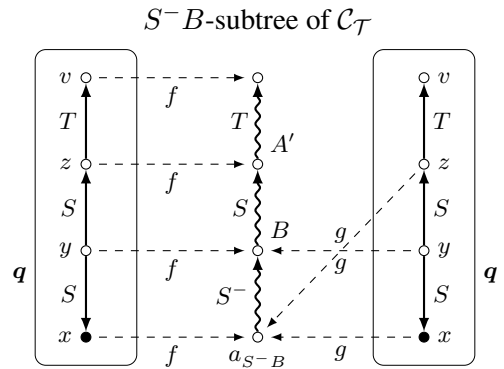
To understand the ingredients required for such a rewriting, suppose $q(\vec{x}) = \exists \vec{y} \varphi(\vec{x}, \vec{y})$, $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ and $\mathcal{C}_{\mathcal{K}} \models^a \varphi$, where \mathbf{a} is an assignment of elements of $\Delta^{\mathcal{C}_{\mathcal{K}}}$ to the variables in φ under which $\mathbf{a}(x) \in \text{ind}(\mathcal{A})$ for all $x \in \vec{x}$. Consider an atom $P(z, z') \in q$ with bound variables z, z' and assume that $\mathbf{a}(t) \in \text{ind}(\mathcal{A})$, for some $t \in \text{term}(q)$. The assignment \mathbf{a} can send z and z' to four different locations in $\mathcal{C}_{\mathcal{K}}$: (A) $\mathbf{a}(z), \mathbf{a}(z') \in \text{ind}(\mathcal{A})$; (B) $\mathbf{a}(z) \in \text{ind}(\mathcal{A})$, $\mathbf{a}(z') \notin \text{ind}(\mathcal{A})$; (B⁻) $\mathbf{a}(z) \notin \text{ind}(\mathcal{A})$, $\mathbf{a}(z') \in \text{ind}(\mathcal{A})$; (O) $\mathbf{a}(z), \mathbf{a}(z') \notin \text{ind}(\mathcal{A})$. Let us see how these alternatives can be reflected in our rewriting. In all of these cases, we need formulas of the form

$$\text{ext}_P(x, y) = \bigvee_{R \sqsubseteq_{\mathcal{T}} P} R(x, y),$$

$$\text{ext}_C(x) = \bigvee_{A \sqsubseteq_{\mathcal{T}} C} A(x) \vee \bigvee_{\exists R \sqsubseteq_{\mathcal{T}} C} \exists y R(x, y).$$

In case (A) we have $\mathcal{C}_{\mathcal{K}} \models^a P(z, z')$ iff $\mathcal{A} \models^a \text{ext}_P(z, z')$. Case (B) is possible only if, for some concept $\exists R.B$, we have $\mathbf{a}(z) \rightsquigarrow w_{[RB]}$, $R \sqsubseteq_{\mathcal{T}} P$ and the atoms of q ‘linked to’ z' can be mapped into the *RB-subtree* of $\mathcal{C}_{\mathcal{T}}$. To illustrate, consider an example.

Example 3 Let $\mathcal{T} = \{A \sqsubseteq \exists S^-.B, B \sqsubseteq \exists S.A', A' \sqsubseteq \exists T\}$ and $q(x) = \{S(y, x), S(y, z), T(z, v)\}$. The answer variable x must be mapped by \mathbf{a} to an ABox element. However, y can be mapped either to an ABox element or to the point $\mathbf{a}(x) \cdot w_{[S^-B]}$ provided that $\mathbf{a}(x)$ is an instance of A (and so of $\exists S^-.B$). In the latter case, we have two possibilities for mapping z : either to $\mathbf{a}(x) \cdot w_{[S^-B]} \cdot w_{[SA']}$, in which case we must further set $\mathbf{a}(v) = \mathbf{a}(x) \cdot w_{[S^-B]} \cdot w_{[SA']} \cdot w_{[T]}$, or to $\mathbf{a}(x)$, provided that $\mathbf{a}(x)$ is an instance of $\exists T$. In the picture below, these two (partial) maps are denoted by f and g .



The observations made in Example 3 are formalised in our central definition. Given a pair (t, t') of adjacent terms in (the graph of) q , we define a *tree witness for (t, t')* to be

a homomorphism f (with domain $\text{dom } f$) from the query

$$\mathbf{q}_f = \{S(s, s') \in \mathbf{q} \mid s, s' \in \text{dom } f\} \cup \{A(s) \in \mathbf{q} \mid s \in \text{dom } f \setminus f^{-1}(r_f)\}$$

to the RB -subtree of $\mathcal{C}_{\mathcal{T}}$, for some $\exists R.B$, with root $r_f = a_{RB}$ such that the following conditions hold:

- (t1) $\text{dom } f$ is the smallest set containing t, t' and such that if $s \in \text{dom } f \setminus f^{-1}(r_f)$ and $S(s, s') \in \mathbf{q}$ then $s' \in \text{dom } f$,
- (t2) $f(t) = r_f$ and if $s \in \text{dom } f \setminus f^{-1}(r_f)$ then s is bound.

Note that this notion of tree witness is different from the one used for the description logic $DL\text{-Lite}_{horn}^N$ by Kontchakov et al. (2010), where the structure of the canonical models ensured uniqueness of a tree witness if it existed. In Example 3, there are two tree witnesses, f and g , for (x, y) . As we shall see below, there may be exponentially many of them.

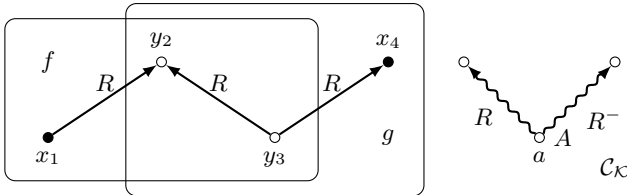
Returning back to case (B), we can say now that there must exist a tree witness f for (z, z') such that $\mathbf{a}(z)$ satisfies all $A(s) \in \mathbf{q}$ with $s \in f^{-1}(r_f)$. Case (B⁻) is symmetric, and in case (O) there must exist $R(t, t') \in \mathbf{q}$ for which (B) holds and $P(z, z')$ is ‘covered’ by a tree witness for (t, t') (as we assumed that $\mathbf{a}(t) \in \text{ind}(\mathcal{A})$, for some $t \in \text{term}(\mathbf{q})$).

This analysis suggests the following idea. Given a CQ \mathbf{q} and KB $\mathcal{K} = (\mathcal{T}, \mathcal{A})$, we guess pairs of adjacent terms (t, t') in \mathbf{q} that will be mapped to edges of the tree part of $\mathcal{C}_{\mathcal{K}}$ starting from $\text{ind}(\mathcal{A})$ (as in case (B) above) and, for each such pair (t, t') , we guess a tree witness for (t, t') . The part of the query that is not covered by the chosen tree witnesses will be mapped to $\text{ind}(\mathcal{A})$ (as in case (A)). The query representing these guesses is then evaluated over $\mathcal{I}_{\mathcal{A}}$. If \mathbf{q} has no answer variables, then we also have to take account of the case where the whole \mathbf{q} is mapped into the tree part of $\mathcal{C}_{\mathcal{K}}$. Unfortunately, this idea cannot be implemented in a straightforward way, as shown by the following example.

Example 4 Let $\mathcal{K} = (\mathcal{T}, \{A(a)\})$, where $\mathcal{T} = \{A \sqsubseteq \exists R, A \sqsubseteq \exists R^-\}$. Consider the query

$$\mathbf{q}(x_1, x_4) = \{R(x_1, y_2), R(y_3, y_2), R(y_3, x_4)\}$$

shown in the picture below alongside $\mathcal{C}_{\mathcal{K}}$.



We obviously have a tree witness f for (x_1, y_2) such that $\text{dom } f = \{x_1, y_2, y_3\}$ and $f^{-1}(r_f) = \{x_1, y_3\}$, and also a tree witness g for (x_4, y_3) with $\text{dom } g = \{x_4, y_3, y_2\}$ and $g^{-1}(r_g) = \{x_4, y_2\}$. Although these tree witnesses cover the whole query \mathbf{q} , they are only ‘realised’ in $\mathcal{C}_{\mathcal{K}}$ under *conflicting* maps: f sends x_1, y_3 to a and y_2 to $a \cdot w_{[R\top]}$, while g sends x_4, y_2 to a and y_3 to $a \cdot w_{[R^-\top]}$; in fact, $\mathcal{K} \not\models \mathbf{q}(a, a)$.

This example motivates the following definitions. Let Θ be the set of tree witnesses for \mathbf{q} and \mathcal{T} . We say that $f, g \in \Theta$ are *compatible* if $\text{dom } f \cap \text{dom } g \subseteq f^{-1}(r_f) \cap g^{-1}(r_g)$.

If f and g are incompatible and neither $\text{dom } f \subseteq \text{dom } g$ nor $\text{dom } g \subseteq \text{dom } f$, then we call f and g *conflicting*. In Example 4, $\text{dom } f \cap \text{dom } g = \{y_2, y_3\}$, and so f and g are conflicting.

We are now in a position to formulate a rewriting based on the idea discussed above. Given a tree witness $f \in \Theta$ for (t, t') with $r_f = a_{RB}$, we first define a *tree-witness formula* tw_f for f by taking

$$\text{tw}_f = \text{ext}_{\exists R.B}(t) \wedge \bigwedge_{s \in f^{-1}(r_f)} ((t = s) \wedge \bigwedge_{A(s) \in \mathbf{q}} \text{ext}_A(s)).$$

A (possibly empty) subset $\Xi \subseteq \Theta$ is called *consistent* if all pairs of tree witnesses in Ξ are compatible. Now, assuming that \vec{y} is a list of all bound variables in $\mathbf{q}(\vec{x})$, we set

$$\mathbf{q}_e(\vec{x}) = \text{detached}_{\mathbf{q}} \vee \bigvee_{\substack{\Xi \subseteq \Theta \\ \Xi \text{ consistent}}} \exists \vec{y} \left(\bigwedge_{f \in \Xi} \text{tw}_f \wedge \bigwedge_{\substack{A(s) \in \mathbf{q} \\ s \notin \text{dom } f \\ \text{for all } f \in \Xi}} \text{ext}_A(s) \wedge \bigwedge_{\substack{P(s, s') \in \mathbf{q} \\ \{s, s'\} \not\subseteq \text{dom } f \\ \text{for all } f \in \Xi}} \text{ext}_P(s, s') \right),$$

where

- $\text{detached}_{\mathbf{q}} = \perp$ if \mathbf{q} has at least one answer variable; otherwise $\text{detached}_{\mathbf{q}}$ is a disjunction of the sentences $\exists x \text{ext}_{\exists R.B}(x)$ such that there is a homomorphism from \mathbf{q} to the RB -subtree of $\mathcal{C}_{\mathcal{T}}$.

Clearly, \mathbf{q}_e is a positive existential formula built from the atoms occurring in \mathbf{q} and \mathcal{T} together with equality unifying certain terms; it contains the same variables and constants as the original CQ \mathbf{q} . Moreover, if we consider the predicates ext_E for concepts and roles as primitive, then \mathbf{q}_e is a union of conjunctive queries (UCQ), where each subquery can be thought of as the result of folding the respective pieces of \mathbf{q} into the tree witness formulas.

Theorem 5 For every ABox \mathcal{A} and every $\vec{a} \subseteq \text{ind}(\mathcal{A})$, we have $(\mathcal{T}, \mathcal{A}) \models \mathbf{q}[\vec{a}]$ iff $\mathcal{I}_{\mathcal{A}} \models \mathbf{q}_e[\vec{a}]$.

We illustrate the rewriting \mathbf{q}_e by a simple example.

Example 6 Suppose that $\mathbf{q}(x) = \{R_i(x, y_i) \mid i \leq n\}$ and $\mathcal{T} = \{A_i \sqsubseteq \exists R_i \mid i \leq n\}$. Each pair (x, y_i) gives rise to one tree witness f_i with $\text{tw}_{f_i} = A_i(x) \vee \exists y R_i(x, y)$, and $\mathbf{q}_e = \bigvee_{N \subseteq [0, n]} \exists \vec{y} (\bigwedge_{i \in N} \text{tw}_{f_i} \wedge \bigwedge_{j \notin N} R_j(x, y_j))$.

As all other known pure rewritings for OWL 2 QL, \mathbf{q}_e is of exponential size: $\mathcal{O}((n_{\mathcal{T}, \mathbf{q}} + 1)^{|\mathbf{q}|} \cdot |\mathcal{T}| \cdot |\mathbf{q}|^2)$, where $n_{\mathcal{T}, \mathbf{q}}$ is the maximum number of distinct tree witness formulas tw_f for tree witnesses containing a pair (t, t') of adjacent terms in \mathbf{q} . Two recent results may help to shed some light on whether this exponential blowup is unavoidable in OWL 2 QL.

One of them shows that no polynomial-time algorithm can construct *pure* rewritings for CQs over OWL 2 QL ontologies, unless $P = NP$ (Kikot, Kontchakov, and Zakharyashev 2011, Theorem 2). The idea of the proof is as follows. First, we encode any CNF $\chi = \bigwedge_{j=1}^m D_j$ over propositional variables p_1, \dots, p_n as an OWL 2 QL TBox,

\mathcal{T}_χ , containing the axioms, for $1 \leq i \leq n$, $1 \leq j \leq m$ and $k = 0, 1$,

$$A_{i-1} \sqsubseteq \exists P^-.X_i^k, \quad X_i^k \sqsubseteq A_i, \quad C_j \sqsubseteq \exists P.C_j,$$

$$X_i^0 \sqsubseteq \exists P.C_j \text{ if } \neg p_i \in D_j, \quad X_i^1 \sqsubseteq \exists P.C_j \text{ if } p_i \in D_j,$$

and consider the CQ

$$q(y_0) = A_0(y_0) \wedge \bigwedge_{i=1}^n P(y_i, y_{i-1}) \wedge A_n(y_n) \wedge$$

$$\bigwedge_{j=1}^m (P(y_n, z_0^j) \wedge \bigwedge_{i=1}^n P(z_{i-1}^j, z_i^j) \wedge C_j(z_n^j)).$$

Now, suppose $q'(y_0)$ is a rewriting of $q(y_0)$ and \mathcal{T}_χ that *does not use any constants*. Consider the ABox $\mathcal{A} = \{A_0(a)\}$. It is not hard to see that $(\mathcal{T}_\chi, \mathcal{A}) \models q[a]$ iff χ is satisfiable. On the other hand, checking whether $\mathcal{I}_\mathcal{A} \models q'[a]$ can be done in polynomial time in $|q'|$ because the domain of $\mathcal{I}_\mathcal{A}$ is a singleton. Thus, constructing the rewriting q' must be at least as hard as deciding satisfiability of χ . (See the Conclusions for a further discussion and results.)

The argument above does not go through if databases are assumed to have two special constants, say 0 and 1, which can be employed in rewritings q' . Indeed, as shown by Gottlob and Schwenk (2011), using 0 and 1, \neq and fresh predicates of arity $\mathcal{O}(\log(|q| \cdot |\mathcal{T}|))$, one can construct a nonrecursive Datalog rewriting q' for any given CQ q and OWL 2 QL ontology \mathcal{T} in polynomial time. Intuitively, q' uses the extra resources to encode a part of the canonical model, which is enough to provide all certain answers to q , to guess a map from q into this part and then check whether it is a homomorphism. As argued in the introduction, RDBMSs do not appear to be best suitable for tasks of that sort, although thorough experiments are required to confirm or refute this claim.

Very often, however, there are other ways to construct polynomial rewritings. Let us assume for a moment that the following condition holds:

(conf) *there are no conflicting tree witnesses in Θ .*

In this case, the query q_e can be transformed into the query

$$q_c(\vec{x}) = \text{detached}_q \vee$$

$$\exists \vec{y} \bigwedge \left[\left(\bigwedge_{\substack{\{t,t'\} \\ t \text{ and } t' \\ \text{adjacent}}} \left(\bigwedge_{A(s) \in q} A(s) \wedge \bigwedge_{R(t,t') \in q} R(t,t') \right) \vee \bigvee_{\substack{f \in \Theta \\ t,t' \in \text{dom } f}} \text{tw}_f \right) \right].$$

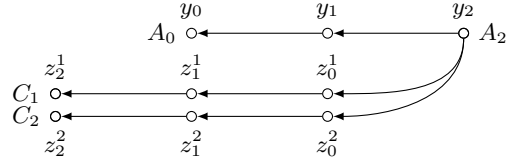
If q does not contain binary predicates then there cannot be any tree witnesses and we set $q_c = q_e$.

Theorem 7 *If \mathcal{T} and q satisfy (conf) then, for any ABox \mathcal{A} and $\vec{a} \subseteq \text{ind}(\mathcal{A})$, we have $(\mathcal{T}, \mathcal{A}) \models q[\vec{a}]$ iff $\mathcal{I}_\mathcal{A} \models q_c[\vec{a}]$.*

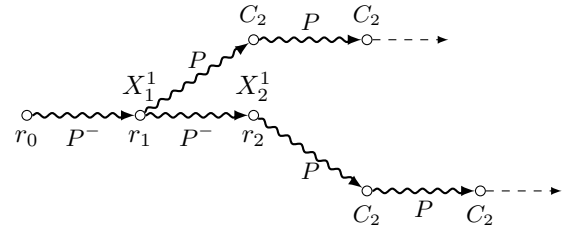
The size of q_c is $\mathcal{O}(n_{\mathcal{T},q} \cdot |\mathcal{T}| \cdot |q|^2)$. Thus, if every pair (t, t') of adjacent terms in q gives rise to polynomially many distinct tree-witness formulas tw_f and condition (conf) holds then q_c is a *polynomial* rewriting of q and \mathcal{T} . (If (conf) does not hold then q_c may return wrong answers.)

Example 8 The exponential query q_e in Example 6 reduces to polynomial $q_c = \exists \vec{y} \bigwedge_{i \leq n} (R_i(x, y_i) \vee \text{tw}_{f_i})$.

A good illuminative example of a CQ with exponentially many tree witness formulas is $q(y_0)$ over the TBox \mathcal{T}_χ constructed above for a CNF χ . One can show that the pairs (z_i^j, z_{i-1}^j) give rise to exponentially many different tree witness formulas tw_f for a suitable χ .



To illustrate, consider the CQ $q(y_0)$ for $n = m = 2$ and suppose that χ is such that the $P^-X_1^1$ -subtree of $\mathcal{C}_{\mathcal{T}_\chi}$ contains the fragment as shown in the picture below.



We can construct a tree witness f for (z_1^1, z_0^1) by taking $f(z_1^1) = r_0, f(z_0^1) = r_1, f(y_2) = r_2$, after which we have three different options for defining $f(z_2^1)$ and $f(z_0^2)$: go back to r_0 , go to r_1 and then take the C_2 -branch, or take the C_2 -branch starting from r_2 . The last two tree witnesses give the same tree witness formula, which is different from that given by the first tree witness. Imagine now some large n and m . It is this fact that makes it ‘hard’ to construct a rewriting for $q(y_0)$ and \mathcal{T}_χ .

The TBox \mathcal{T}_χ and CQ $q(y_0)$ involve a very complex interplay between role inclusions (or concepts of the form $\exists R.C$) and inverse roles, which appears to be rather artificial compared to how roles are used in real-world ontologies. The experiments to be reported later on in the paper demonstrate that real-world ontologies and queries generate very few tree witnesses, which are never in conflict, and so the rewriting q_c is both short and correct.

In the next section we analyse conflicting tree witnesses in more detail.

Conflicting Tree Witnesses

One simple way to tackle the problem of conflicting tree witnesses, identified in Example 4, would be to introduce a *fresh constant symbol*, say ν , representing all the non-ABox elements of the canonical models. We then use the following variant of the tree witness formula tw_f for (t, t') :

$$\text{tw}_f^\nu = \text{tw}_f \wedge \bigwedge_{s \in \text{dom } f \setminus f^{-1}(r_f)} (s = \nu).$$

We denote the resulting ‘impure’ rewriting by q_c^ν . Given an ABox \mathcal{A} , denote by $\mathcal{I}_\mathcal{A}^\nu$ the interpretation $\mathcal{I}_\mathcal{A}$ extended with a new domain element (interpreting) ν . Thus, ν is not involved in the interpretation of any predicate in $\mathcal{I}_\mathcal{A}$, and so, of any predicate ext_E .

Example 9 In the context of Example 4, tw_g^ν would contain the conjuncts $\text{ext}_{\exists R^-, \top}(x_4)$ and $(y_2 = x_4)$, which cannot be satisfied in \mathcal{I}_A^ν at the same time as the conjunct $(y_2 = \nu)$ of tw_f^ν . Thus, $\mathcal{I}_A^\nu \models \mathbf{q}_c^\nu(a, a)$.

The following result shows that \mathbf{q}_c^ν is a correct rewriting over this interpretation.

Theorem 10 For any ABox \mathcal{A} and any $\bar{a} \subseteq \text{ind}(\mathcal{A})$, we have $(\mathcal{T}, \mathcal{A}) \models \mathbf{q}[\bar{a}]$ iff $\mathcal{I}_A^\nu \models \mathbf{q}_c^\nu[\bar{a}]$.

The rewriting \mathbf{q}_c^ν can be viewed as a step in the direction of the combined approach (Lutz, Toman, and Wolter 2009; Kontchakov et al. 2010), though without expanding the ABox by applying the TBox axioms. Roughly, in both cases the RDBMS has to guess whether a bound variable is mapped to $\text{ind}(\mathcal{A})$ or to the tree part of the canonical model.

There is another way to ‘suppress conflicts’ in \mathbf{q}_c for the majority of practical cases, while keeping the resulting rewriting ‘pure.’ For a CQ \mathbf{q} , let \mathbf{q}_c^+ be the rewriting obtained from \mathbf{q}_c by replacing every tw_f in it with the following formula tw_f^+ :

$$\text{tw}_f^+ = \text{tw}_f \wedge \bigwedge_{f, g \text{ conflicting}} (\mathbf{q}_{g \setminus f})_c^+,$$

where $\mathbf{q}_{g \setminus f}$ denotes the restriction of the query \mathbf{q} to the set $\text{dom } g \setminus (\text{dom } f \setminus f^{-1}(r_f))$ (i.e., all terms in g that are not non-root terms in f) and $(\mathbf{q}_{g \setminus f})_c^+$, in turn, is the rewriting (as defined above) of the query $\mathbf{q}_{g \setminus f}$ in which all the variables in $f^{-1}(r_f) \cup g^{-1}(r_g)$ are regarded to be answer variables.

Example 11 In the context of Example 4, the rewriting \mathbf{q}_c^+ is constructed in two steps as follows: first, we obtain a representation of \mathbf{q}_c^+ as the following formula:

$$\begin{aligned} & \exists y_2, y_3 \left((R(x_1, y_2) \vee \text{tw}_f^+) \wedge \right. \\ & \left. (R(y_3, y_2) \vee \text{tw}_f^+ \vee \text{tw}_g^+) \wedge (R(y_3, x_4) \vee \text{tw}_g^+) \right), \end{aligned}$$

where

$$\begin{aligned} \text{tw}_f^+ &= \text{ext}_{\exists R}(x_1) \wedge (x_1 = y_3) \wedge (R(y_3, x_4))_c^+, \\ \text{tw}_g^+ &= \text{ext}_{\exists R^-(x_4)} \wedge (x_4 = y_2) \wedge (R(x_1, y_2))_c^+, \end{aligned}$$

with $R(y_3, x_4)$ and $R(x_1, y_2)$ being two new queries, which have only answer variables. Then, the rewritings of these two queries coincide with the queries themselves: $R(y_3, x_4)$ and $R(x_1, y_2)$, respectively. So, for example, the modified tree witness formula tw_f^+ for f contains the conjunct $R(y_3, x_4)$, which cannot be satisfied in the interpretation \mathcal{I}_A with $\mathcal{A} = \{A(a)\}$.

Theorem 12 Suppose that \mathbf{q} and \mathcal{T} satisfy the condition **(conf₁)** for every $f \in \Theta$, if there are $g, h \in \Theta$ such that the pairs f, g and f, h are both conflicting, then either $\text{dom } g \subseteq \text{dom } h$ or $\text{dom } h \subseteq \text{dom } g$.

Then, for every ABox \mathcal{A} and every $\bar{a} \subseteq \text{ind}(\mathcal{A})$, we have $(\mathcal{T}, \mathcal{A}) \models \mathbf{q}[\bar{a}]$ iff $\mathcal{I}_A \models \mathbf{q}_c^+[\bar{a}]$.

We do not know any cases where \mathbf{q}_c^+ is not correct. It is to be noted, however, that the rewriting \mathbf{q}_c^+ can be of exponential size even if every pair (t, t') in \mathbf{q} gives rise to at most one tree witness.

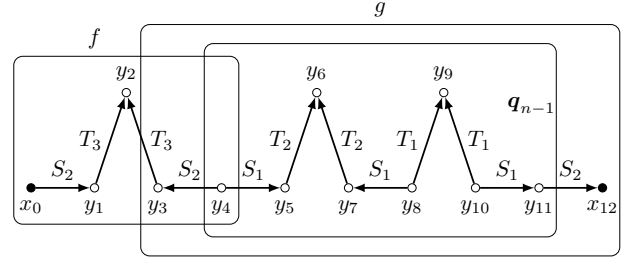
Example 13 Given a word $R_1 \dots R_m$ over roles, define the CQ

$$\mathbf{q}^{R_1 \dots R_m}(x_0, x_m) = \{R_i(x_{i-1}, x_i) \mid 1 \leq i \leq m\}.$$

Let σ_n be the following sequence of words of roles:

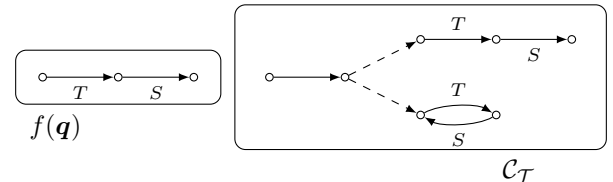
$$\sigma_1 = T_1 T_1^-, \quad \sigma_{n+1} = S_n T_{n+1} T_{n+1}^- S_n^- \sigma_n S_n, \text{ for } n \geq 1.$$

Now, consider the CQs $\mathbf{q}_n = \mathbf{q}^{\sigma_n}(x_0, x_n)$ and TBoxes $\mathcal{T}_n = \{A_n \sqsubseteq \exists R_n, A_{n+1} \sqsubseteq \exists S_n.A_n\}$. It is not hard to see that \mathbf{q}_n contains conflicting tree witnesses f and g as shown in the picture below, the query \mathbf{q}_{n-1} is a subquery of $\mathbf{q}_{g \setminus f}$, and so $(\mathbf{q}_{n-1})_c^+$ occurs in $(\mathbf{q}_n)_c^+$ four times. On the other hand, there is a constant c such that $|\mathbf{q}_n| = |\mathbf{q}_{n-1}| + c$.



Sufficient Conditions for Polynomial Rewriting

The rewriting \mathbf{q}_c is exponentially long if there are exponentially many tree witness formulas tw_f for some pair (t, t') of adjacent terms in \mathbf{q} . Each tw_f is determined by the RB -subtree of $\mathcal{C}_\mathcal{T}$ (containing the range of f), the domain $\text{dom } f$ of f and the set $f^{-1}(r_f)$; the terms in this set will be called the *roots* of f . Now we show that, for a large class of CQs \mathbf{q} and OWL 2 QL TBoxes \mathcal{T} , all tree witnesses for the same pair (t, t') in \mathbf{q} have the same domain and roots. As the number of distinct RB -subtrees of $\mathcal{C}_\mathcal{T}$ is polynomial in the size of \mathcal{T} , the rewriting \mathbf{q}_c in this case is also polynomial; moreover, we show that it can be constructed in polynomial time in $|\mathbf{q}|$ and $|\mathcal{T}|$. Roughly, Theorem 19 to be proved below demonstrates that this can be done if condition **(conf)** holds and there are no roles T, S and a tree witness f for which $f(\mathbf{q})$ and $\mathcal{C}_\mathcal{T}$ simultaneously contain fragments of the form



Here, by $f(\mathbf{q})$ we mean the quotient \mathbf{q}_f / \sim of \mathbf{q}_f modulo the equivalence relation \sim defined by taking $s \sim s'$ iff $f(s) = f(s')$. Roles T and S are called *adjacent* in $f(\mathbf{q})$ if $T(s_1, s_2), S(s_2, s_3) \in f(\mathbf{q})$, for some $s_i \in \text{term } f(\mathbf{q})$ such that $s_2 \notin f^{-1}(r_f)$.

We say that a role S is *forward* in \mathcal{T} if $u \rightsquigarrow v$ for all $(u, v) \in S^{\mathcal{C}_\mathcal{T}}$. If neither S nor its inverse S^- is forward then S is said to be a *twisty role* in \mathcal{T} .

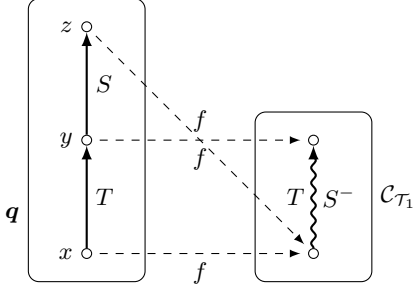
A tree witness f is called *perfect* if, for every pair T, S of adjacent roles in $f(\mathbf{q})$ such that S is twisty in \mathcal{T} , we have

(perf) $\mathcal{C}_{\mathcal{T}} \models \text{inv}(T, S) \wedge \text{suc}(T, S)$,

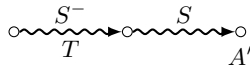
where

$$\begin{aligned} \text{inv}(T, S) &= \exists x, y (T(x, y) \wedge S(y, x)), \\ \text{suc}(T, S) &= \exists x, y, z (T(x, y) \wedge S(y, z) \wedge (x \neq z)). \end{aligned}$$

Example 14 Consider $\mathcal{T}_1 = \{A \sqsubseteq \exists T, T \sqsubseteq S^-\}$ and $\mathbf{q}(x) = \{T(x, y), S(y, z)\}$. We have $\mathcal{C}_{\mathcal{T}_1} \models \text{inv}(T, S)$ and $\mathcal{C}_{\mathcal{T}_1} \models \text{inv}(R, R^-)$, for every role R . Both T and S^- are forward roles in \mathcal{T}_1 . The homomorphism f shown below is a perfect tree witness for (x, y) :



Consider now $\mathcal{T}_2 = \{A \sqsubseteq \exists T, T \sqsubseteq S^-, \exists T^- \sqsubseteq \exists S.A'\}$. As $\mathcal{C}_{\mathcal{T}_2}$ contains the following fragment



S is a twisty role in \mathcal{T}_2 , and $\mathcal{C}_{\mathcal{T}_2} \models \text{inv}(T, S) \wedge \text{suc}(T, S)$. Thus, there is no perfect tree witness for (x, y) in \mathbf{q} and \mathcal{T}_2 , though there are two ‘imperfect’ tree witnesses.

It turns out that if all tree witnesses f for a pair of terms in \mathbf{q} are perfect then all of them have the same domain and roots, and so (i) the tree-witness formulas tw_f occur in the disjunctions for the same edges (t, t') and (ii) the tw_f may differ only in their $\text{ext}_{\exists R.B}(t)$ conjuncts. In the following example, exponentially many different tree witnesses give rise to the same tree witness formula.

Example 15 Let $\mathbf{q}(x) = \{S(x, y), R(y, z_i) \mid 1 \leq i \leq n\}$ and $\mathcal{T} = \{A \sqsubseteq \exists S, \exists S^- \sqsubseteq \exists R.B_1, \exists S^- \sqsubseteq \exists R.B_2\}$. There are 2^n (perfect) tree witnesses for (x, y) , as each z_i can be mapped either to a B_1 - or a B_2 -point in $\mathcal{C}_{\mathcal{T}}$. All these tree witnesses have the same domain (all terms in \mathbf{q}) and only one root (x). They define the same tree-witness formula $\text{ext}_{\exists S}(x)$.

The notion of universal tree witness we are about to introduce will serve as a compact representation for all such similar tree witnesses.

For a pair (t, t') of adjacent terms in \mathbf{q} , a tree-shaped CQ \mathbf{f} is called a *universal tree witness for (t, t') in \mathbf{q} and \mathcal{T}* if there is a partial homomorphism f from \mathbf{q} onto \mathbf{f} such that, for every tree witness g for (t, t') in \mathbf{q} and \mathcal{T} ,

- $\text{dom } g = \text{term}(\mathbf{f})$, and
- there exists a forward homomorphism f' from \mathbf{f} to $\mathcal{C}_{\mathcal{T}}$ such that $g(s) = f'(f(s))$, for every $s \in \text{term}(\mathbf{f})$,

where by a forward homomorphism we understand a homomorphism that preserves the distance from the root (recall that both \mathbf{f} and $\mathcal{C}_{\mathcal{T}}$ are tree-shaped).

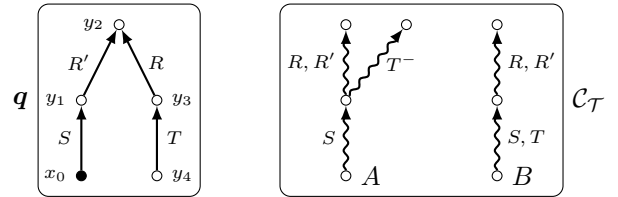
Lemma 16 (i) If all tree witnesses for a pair (t, t') of terms in \mathbf{q} and \mathcal{T} are perfect, then there is a unique (up to isomorphism) universal tree witness for (t, t') in \mathbf{q} and \mathcal{T} .

(ii) There is a polynomial-time algorithm which, given \mathbf{q} , \mathcal{T} and a pair (t, t') of adjacent terms in \mathbf{q} , checks whether all the tree witnesses for (t, t') in \mathbf{q} and \mathcal{T} are perfect, and if this is the case, returns a universal tree witness for (t, t') in \mathbf{q} and \mathcal{T} .

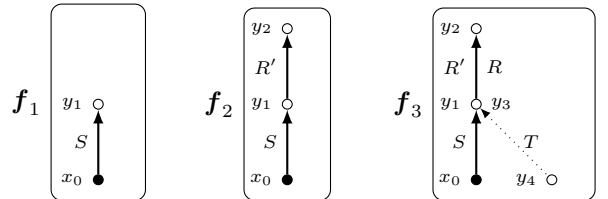
As a consequence of this lemma we obtain that, if all tree witnesses for (t, t') are perfect, then all of them share the same domain and roots.

Example 17 In Example 15, a universal tree witness for (x, y) coincides with the whole query $\mathbf{q}(x)$.

In the proof of Lemma 16, we construct a finite sequence \mathbf{f}_i of tree-shaped CQs \mathbf{f}_i such that the final \mathbf{f}_n is a universal tree witness for (t, t') . We begin by taking all the atoms with terms t, t' and then try to satisfy the tree witness condition (t1) by adding atoms with new terms to the current \mathbf{f}_i . Condition (perf), which is being checked during the construction, allows us to decide whether we have to add a new term to \mathbf{f}_i or reuse an existing one. To illustrate, we give one more example.



Example 18 Consider $\mathcal{T} = \{A \sqsubseteq \exists S.C, C \sqsubseteq \exists T^-, \exists S^- \sqsubseteq \exists R, R \sqsubseteq R', B \sqsubseteq \exists T', T' \sqsubseteq T, T' \sqsubseteq S\}$ and $\mathbf{q}(x_0) = \{S(x_0, y_1), R'(y_1, y_2), R(y_3, y_2), T(y_4, y_3)\}$ in the picture above. A universal tree witness for (x_0, y_1) must contain $S(x_0, y_1)$. As the role in $R'(y_1, y_2)$ is forward, the universal tree witness must contain this atom too. The role R in $R^-(y_2, y_3)$ is also forward, so we add this atom and identify y_1 with y_3 . This gives three approximations of the universal tree witness we are constructing:



In $T(y_4, y_3)$, the role T is twisty in \mathcal{T} . The canonical model $\mathcal{C}_{\mathcal{T}}$ suggests two alternative ways of treating $T(y_4, y_3)$: either to add it as it is, or to identify y_4 with root x_0 . Thus, no universal tree witness for (x_0, y_1) exists. This is also signalled by the fact that condition (perf) fails: by identifying y_1 and y_3 , we make S and T adjacent, while both $\text{suc}(S, T)$ and $\text{inv}(S, T)$ hold in $\mathcal{C}_{\mathcal{T}}$. Note that S and T were not adjacent in the original query \mathbf{q} ; that is why (perf) is checked for the image $f(\mathbf{q})$ of every tree witness f . As the number of tree witnesses can be exponential, this complication might suggest that (perf) cannot be checked efficiently. The proof

of Lemma 16 shows that this check can be done in polynomial time without constructing all tree witnesses.

Strictly speaking, a universal tree witness is not a tree witness in the sense of our original definition, but rather a convenient structure representing all tree witnesses for (t, t') , even if there are exponentially many of them. Universal tree witnesses can be used to further simplify the rewriting q_c . Namely, all formulas tw_f for (t, t') are merged into a single tree-witness formula tw_f for the universal tree witness f for (t, t') , which is defined by taking:

$$tw_f = \left[\bigvee_{\substack{\exists R.B \text{ such that there is} \\ \text{a homomorphism } h: \mathcal{F} \rightarrow \mathcal{C}_{\mathcal{T}} \\ h(t)=a_{RB} \\ h(t')=b_{RB}}} \text{ext}_{\exists R.B}(t) \right] \wedge \bigwedge_{\substack{s \text{ is root of } f \\ A(s) \in \mathcal{Q}}} ((t = s) \wedge \bigwedge_{A(s) \in \mathcal{Q}} A(s)).$$

As a universal tree witness is unique for each pair (t, t') of adjacent terms, let us denote it by $q_{(t,t')}$. We are now in a position to define our *polynomial* rewriting q_p for q and \mathcal{T} :

$$q_p(\vec{x}) = \text{detached}_q \vee \bigvee_{\substack{\{t,t'\} \\ t \text{ and } t' \\ \text{adjacent}}} \left[\left(\bigwedge_{\substack{A(s) \in \mathcal{Q} \\ t \text{ and } t' \\ \text{adjacent}}} \text{ext}_A(s) \wedge \bigwedge_{R(t,t') \in \mathcal{Q}} \text{ext}_R(t, t') \right) \vee \bigvee_{\substack{s \text{ and } s' \text{ adjacent} \\ t, t' \in \text{term}(q_{(s,s')})}} tw_{q_{(s,s')}} \right].$$

Theorem 19 *If all tree witnesses for q and \mathcal{T} are perfect and condition (conf) is satisfied then, for any ABox \mathcal{A} and any $\vec{a} \subseteq \text{ind}(\mathcal{A})$, we have $(\mathcal{T}, \mathcal{A}) \models q[\vec{a}]$ iff $\mathcal{I}_{\mathcal{A}} \models q_p[\vec{a}]$. Moreover, q_p is constructed in time polynomial in $|q|, |\mathcal{T}|$.*

Theorem 19 can be substantially sharpened by taking account of the ‘types’ of terms in q .

Example 20 Consider again the TBox \mathcal{T}_2 from Example 14 and the CQ $q_1(x) = \{T(x, y), S(y, z), A''(z)\}$. This time we have only one tree witness for (x, y) . It is not perfect; yet, it is type-perfect because the ‘typed’ *suc* formula $\exists x, y, z (T(x, y) \wedge S(y, z) \wedge A''(z) \wedge (x \neq z))$ does not hold in $\mathcal{C}_{\mathcal{T}_2}$.

If an ontology \mathcal{T} does not contain any twisty roles, then all tree witnesses in any CQ q over \mathcal{T} are clearly perfect. On the other hand, all examples of conflicting tree witnesses given above involve twisty roles. The following theorem shows that this is no accident:

Theorem 21 *Suppose that q is a CQ and \mathcal{T} an OWL 2 QL ontology without twisty roles. Then there are no conflicting tree witnesses for q and \mathcal{T} . Thus, the rewriting q_p for q and \mathcal{T} is correct and can be constructed in polynomial time.*

It may be worth noting that OWL 2 EL (Baader, Brandt, and Lutz 2005; 2008) ontologies satisfy this sufficient condition, and so a polynomial rewriting similar to q_p can also be used for conjunctive query answering over such ontologies provided that the ABoxes are complete (or saturated) with respect to the ontologies.

Experiments

To understand how the rewriting q_p looks like in real-world practice, we have run experiments with three known

| | | QuOnto size of UCQ | Requiem size of UCQ | Presto size of Datalog | Nyaya size of UCQ | tw/ utw | ext rules | non ext rules |
|---|----|--------------------------|---------------------------|------------------------------|-------------------------|------------|--------------|---------------------|
| A | Q1 | 783 | 402 | 69 | 249 | 8/1 | 53 | 3 |
| | Q2 | 1,812 | 103 | 52 | 94 | 1 | 30 | 3 |
| | Q3 | 4,763 | 104 | 55 | 104 | 0 | 30 | 1 |
| | Q4 | 7,251 | 492 | 93 | 456 | 4/1 | 42 | 3 |
| | Q5 | 78,885 | 624 | 71 | 624 | 0 | 36 | 1 |
| U | Q1 | 5 | 2 | 6 | 2 | 0 | 2 | 1 |
| | Q2 | 287 | 148 | 1 | 1 | 0 | 0 | 1 |
| | Q3 | 1,260 | 224 | 8 | 4 | 0 | 4 | 1 |
| | Q4 | 5,364 | 1,628 | 6 | 2 | 0 | 2 | 1 |
| | Q5 | 9,245 | 2,960 | 11 | 10 | 0 | 7 | 1 |
| | Q6 | 588 | 42 | 19 | 26 | 1 | 16 | 5 |
| | Q7 | 5,950 | 630 | 37 | 39 | 2 | 16 | 7 |
| S | Q1 | 6 | 6 | 7 | 6 | 0 | 6 | 1 |
| | Q2 | 204 | 160 | 3 | 2 | 0 | 2 | 1 |
| | Q3 | 1,194 | 480 | 5 | 4 | 0 | 4 | 1 |
| | Q4 | 1,632 | 960 | 5 | 4 | 0 | 4 | 1 |
| | Q5 | 11,487 | 2,880 | 7 | 8 | 0 | 6 | 1 |
| | Q6 | 588 | 564 | 12 | 36 | 1 | 24 | 3 |
| | Q7 | 118 | 106 | 14 | 29 | 1 | 19 | 5 |

OWL 2 QL ontologies—Adolena (A), University (U) and Stockexchange (S)—using the same conjunctive queries as in (Pérez-Urbina, Motik, and Horrocks 2009; Rosati and Almatelli 2010; Gottlob, Orsi, and Pieris 2011) and two new ones (Q6 and Q7); the ontologies and queries are available at www.dcs.bbk.ac.uk/~roman/query-rewriting. The results of the experiments are collected in the table above.

The most important conclusion we can draw from these experiments is that, in practice, queries and ontologies generate very few tree witnesses, if any; they are never in conflict with each other, the sufficient condition of Theorem 19 is satisfied, and so the rewritings q_c and q_p are both short and correct.

To describe how the rewritten queries look like and explain the figures in the table, assume first that a given CQ q contains no tree witnesses with respect to a given ontology \mathcal{T} . In this case, q_p is obtained from q by replacing each unary atom $A(s)$ with $\text{ext}_A(s)$ and each binary atom $R(s, s')$ with $\text{ext}_R(s, s')$. Roughly, the ext_E contain those concepts/roles that are located under E in the classification of \mathcal{T} by an ontology reasoner. It will be convenient for us to represent the ext_E as nonrecursive Datalog programs. For example, Adolena gives the rules:

$$\begin{aligned} \text{ext}_{\text{affects}}(x, y) &:- \text{affects}(x, y). \\ \text{ext}_{\text{affectedBy}}(x, y) &:- \text{isAffectedBy}(y, x). \\ \text{ext}_{\text{device}}(x) &:- \text{Device}(x). \\ &\dots \end{aligned}$$

The column ‘ext rules’ in the table refers to the number of such Datalog rules for the given CQ and ontology.

Consider now a query containing tree witnesses, for instance the following query Q4 over Adolena:

$$q(x) = \exists y (\text{Device}(x) \wedge \text{assistsWith}(x, y) \wedge \text{PhysicalAbility}(x, y)).$$

This CQ contains a tree witness f with root

$$f(x) = a_{\text{assistsWith, MovementAbility}}.$$

In this case, for each pair (s, s') of adjacent terms in q , we introduce a fresh predicate, say $\text{edge}_{s,s'}$, and define it by means of a nonrecursive Datalog program such as

$$\begin{aligned} \text{edge}_{x,y}(x, y) &:- \text{ext}_{\text{Device}}(x), \text{ext}_{\text{assistsWith}}(x, y), \\ &\quad \text{ext}_{\text{PhysicalAbility}}(x, y). \\ \text{edge}_{x,y}(x, y) &:- \text{ext}_{\text{Device}}(x), \\ &\quad \text{ext}_{\exists \text{assistsWith.MovementAbility}}(x). \end{aligned}$$

where the first rule corresponds to choosing both x and y among the ABox individuals and the second rule comes from a single universal tree witness for (x, y) . Then we replace by $\text{edge}_{s,s'}$ all the atoms in the CQ that are ‘covered’ by the terms s and s' . In our running example, we obtain

$$q(x) :- \text{edge}_{x,y}(x, y).$$

The column ‘non-ext rules’ in the table refers to the number of such rules (one rule for the whole q plus the rules defining the $\text{edge}_{s,s'}$); ‘tw’ gives the number of tree witnesses in the CQ, and ‘utw’ the number of universal tree witnesses.

In theory, a correct rewriting of a CQ q and an OWL 2 QL ontology \mathcal{T} can be exponential only if q and \mathcal{T} give rise to exponentially many tree witnesses, in which case the canonical model $\mathcal{C}_{\mathcal{T}}$ must be extremely complex. Our experiments indicate that, in practice, the contribution of tree witnesses does not look essential at all, especially in comparison with the contribution of the definitions of ext_E , which reflects the depth and width of the concept and role hierarchies in \mathcal{T} rather than the complexity of $\mathcal{C}_{\mathcal{T}}$. Note also that the same predicates ext_E are used in all queries, which makes these predicates an ideal target for optimisations. The ways to minimise the influence of these rules depend on how we store the data.

There are two main approaches to storing data in OBDA. Suppose first that an ABox \mathcal{A} is stored in a local database and the system has a certain degree of control over the data. In this case, one can saturate \mathcal{A} with the intensional data that is implied by the TBox axioms (more precisely, construct the ABox part of the canonical model). Having done so, we do not need the ext_E predicates any more and can replace them with the corresponding E . The ABox saturation (more precisely, a finite encoding of the whole canonical model) was suggested in the combined approach (Lutz, Toman, and Wolter 2009; Kontchakov et al. 2010). However, the downside of the ABox saturation is a significant increase of the storage space required for the data (and a slowdown of updates). One solution to this problem was found by Rodriguez Muro and Calvanese (2011a; 2011b). In a nutshell, the idea is to build a ‘semantic index’ by assigning numerical identifiers to the concept names, used in the ontology, in such a way that all subclasses of a given concept are associated with an interval (or a few intervals) of numbers. The semantic index allows one to encode the definitions of any of the ext_E predicates as a single query that selects all instances with concept identifiers falling into the respective interval(s). In this case, the classical database indexing techniques are employed to ensure efficiency of these interval queries.

In the other typical OBDA scenario, ABoxes do not come as sets of triples stored in a single database. Instead, the sets

of individuals that belong to concepts and roles are defined by means of queries (mappings) to a number of (relational) data sources (Lenzerini 2002; Calvanese et al. 2007b). Consider, for instance, the rules for the role affects above. One can clearly expect mappings to be defined in such a way that $\text{affects}(a, b) \in \mathcal{A}$ iff $\text{isAffectedBy}(b, a) \in \mathcal{A}$ in every ABox \mathcal{A} . This suggests that, in fact, there is no need for two separate rules, so that the predicate $\text{ext}_{\text{affects}}$ can be eliminated altogether (replaced by $\text{affects}(x, y)$, which halves the number of CQs produced). It is also quite feasible that information about all devices is stored in a single database relation and mappings for each of the 26 subclasses of the concept Device select appropriate devices from the same database relation. In such a case, every ABox \mathcal{A} defined by these mappings will be *complete* for all subclasses A of Device in the sense that $A(a) \in \mathcal{A}$ iff $(\mathcal{T}, \mathcal{A}) \models A(a)$. Therefore, with this information at hand, one can replace the 26 rules defining $\text{ext}_{\text{Device}}(x)$ with just a single rule $\text{ext}_{\text{Device}}(x) : \neg \text{Device}(x)$, or even eliminate the predicate $\text{ext}_{\text{Device}}(x)$ altogether.

Conclusions

In this paper, we considered pure (positive existential) rewritings of conjunctive queries over OWL 2 QL ontologies and analysed why such rewritings can be lengthy. We showed that the length of a rewriting is related to the number of tree witnesses in the query, which reflect how various parts of the query can be homomorphically mapped to the tree (‘intensional’) part of the canonical model. Thus, a rewriting can be lengthy if the original query is sufficiently long and the intensional part of the canonical model for the ontology is sufficiently complex. We proved that by restricting the interaction between inverse roles and role inclusion axioms in ontologies and queries, we can guarantee transparent polynomial rewritings. Moreover, we also demonstrated that real-world ontologies and queries contain very few tree witnesses, satisfy the above mentioned restrictions, and so enjoy polynomial rewritings.

Remark 22 When the final version of this paper was ready for submission, we obtained some new results that shed more light on the size of pure rewritings. Below is a brief summary of these results; for details consult the preliminary report (Kikot, Kontchakov, Podolskii and Zakharyashev 2012).

- (1) An exponential blow-up is unavoidable for pure positive existential rewritings (PE) and pure nonrecursive Datalog (NDL) rewritings; pure FO-rewritings can blow-up superpolynomially unless $\text{NP} \subseteq \text{P/poly}$.
- (2) Pure NDL-rewritings are in general exponentially more succinct than pure PE-rewritings.
- (3) Pure FO-rewritings can be superpolynomially more succinct than pure PE-rewritings.
- (4) Impure PE-rewritings can always be made polynomial, and so they are exponentially more succinct than pure PE-rewritings.

(1)–(3) are proved by first establishing connections between pure rewritings for CQs over OWL 2 QL ontologies and circuits for monotone Boolean functions, and then using known

lower bounds and separation results for the circuit complexity of such functions as $\text{CLIQUE}(n, k)$ ‘a graph with n nodes contains a k -clique’ and $\text{MATCHING}(2n)$ ‘a bipartite graph with n vertices in each part has a perfect matching’ (Razborov 1985; Borodin, von zur Gathen, and Hopcroft 1982; Raz and Wigderson 1992; Raz and McKenzie 1997)

The polynomial PE-rewriting in (4) is similar to the NDL-rewriting of Gottlob and Schwentick (2011): using two extra constants, $=$ and (polynomially-many) new existentially quantified variables, one can encode a relevant part of the canonical model of \mathcal{T} in the rewritten query. The difference between the resulting impure PE-rewritings and the exponential-size pure PE-rewritings is of the same kind as the difference between deterministic and nondeterministic Boolean circuits. As shown by Razborov (1985), no polynomial-size deterministic monotone circuit can compute $\text{CLIQUE}(n, k)$; however, it can be computed by a polynomial-size nondeterministic circuit (or a QBF), where the existentially quantified variables guess k vertices and the circuit checks whether they form a k -clique in the given graph. In the polynomial impure PE-rewriting (4), the extra constants, variables and $=$ are used to nondeterministically guess a part of the canonical model into which the query can be mapped. (We conjecture that there does not exist an impure polynomial-size rewriting if the number of new existentially quantified variables is bounded.)

Acknowledgments

The work on this paper was supported by the U.K. EPSRC grant EP/H05099X/1.

We are grateful to G. Gottlob, G. Orsi, R. Rosati and D. Tsarkov who helped us to run the experiments.

References

- Abiteboul, S.; Hull, R.; and Vianu, V. 1995. *Foundations of Databases*. Addison-Wesley.
- Baader, F.; Calvanese, D.; McGuinness, D.; Nardi, D.; and Patel-Schneider, P., eds. 2003. *The Description Logic Handbook: Theory, Implementation and Applications*. Cambridge University Press.
- Baader, F.; Brandt, S.; and Lutz, C. 2005. Pushing the EL envelope. In *Proc. of the 19th Int. Joint Conf. on Artificial Intelligence, IJCAI-05*, 364–369. Professional Book Center.
- Baader, F.; Brandt, S.; and Lutz, C. 2008. Pushing the EL envelope further. In Clark, K., and Patel-Schneider, P. F., eds., *Proc. of the OWLED 2008 DC Workshop on OWL: Experiences and Directions*.
- Borodin, A.; von zur Gathen, J.; and Hopcroft, J. E. 1982. Fast parallel matrix and GCD computations. In *Proc. of the 23rd Annual Symp. on Foundations of Computer Science, FOCS’82*, 65–71. IEEE Computer Society.
- Calvanese, D.; De Giacomo, G.; Lembo, D.; Lenzerini, M.; and Rosati, R. 2007a. Tractable reasoning and efficient query answering in description logics: The *DL-Lite* family. *J. of Automated Reasoning* 39(3):385–429.
- Calvanese, D.; De Giacomo, G.; Lembo, D.; Lenzerini, M.; Poggi, A.; and Rosati, R. 2007b. Ontology-based database access. In *Proc. of the 15th Ital. Conf. on Database Systems, SEBD 2007*, 324–331.
- Chortaras, A.; Trivela, D.; and Stamou, G. 2011. Goal-oriented query rewriting for OWL 2 QL. In *Proc. of the 24th Int. Workshop on Description Logics, DL 2011*, vol. 745 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Gottlob, G., and Schwentick, T. 2011. Rewriting ontological queries into small nonrecursive datalog programs. In *Proc. of the 24th Int. Workshop on Description Logics, DL 2011*, vol. 745 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Gottlob, G.; Orsi, G.; and Pieris, A. 2011. Ontological queries: Rewriting and optimization. In *Proc. of the 27th Int. Conf. on Data Engineering, ICDE 2011*, 2–13. IEEE Computer Society.
- Kikot, S.; Kontchakov, R.; and Zakharyashev, M. 2011. On (In)Tractability of OBDA with OWL 2 QL. In *Proc. of the 24th Int. Workshop on Description Logics, DL 2011*, vol. 745 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Kikot, S.; Kontchakov, R.; Podolskii, V.; and Zakharyashev, M. 2012. Exponential lower bounds and separation for query rewriting. *CoRR*, arXiv:1202.4193, 2012.
- Kontchakov, R.; Lutz, C.; Toman, D.; Wolter, F.; and Zakharyashev, M. 2010. The combined approach to query answering in DL-Lite. In *Proc. of the 12th Int. Conf. on Principles of Knowledge Representation and Reasoning, KR 2010*. AAAI Press.
- Kontchakov, R.; Lutz, C.; Toman, D.; Wolter, F.; and Zakharyashev, M. 2011. The combined approach to ontology-based data access. In *Proc. of the 20th Int. Joint Conf. on Artificial Intelligence, IJCAI-2011*, 2656–2661. AAAI Press.
- Lenzerini, M. 2002. Data integration: A theoretical perspective. In *Proc. of the 21st ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems, PODS 2002*, 233–246.
- Lutz, C.; Toman, D.; and Wolter, F. 2009. Conjunctive query answering in the description logic EL using a relational database system. In *Proc. of the 21st Int. Joint Conf. on Artificial Intelligence, IJCAI 2009*, 2070–2075. AAAI Press.
- Pérez-Urbina, H.; Motik, B.; and Horrocks, I. 2009. A comparison of query rewriting techniques for DL-Lite. In *Int. Workshop on Description Logics, DL 2009*, vol. 477 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Raz, R., and McKenzie, P. 1997. Separation of the monotone nc hierarchy. In *Proc. of the 38th Annual Symp. on Foundations of Computer Science, FOCS’97*, 234–243. IEEE Computer Society.
- Raz, R., and Wigderson, A. 1992. Monotone circuits for matching require linear depth. *J. ACM* 39(3):736–744.
- Razborov, A. 1985. Lower bounds for the monotone complexity of some Boolean functions. *Dokl. Akad. Nauk SSSR* 281(4):798–801.
- Rodríguez Muro, M., and Calvanese, D. 2011a. Dependencies to optimize ontology based data access. In *Proc.*

of the 24th Int. Workshop on Description Logics, DL 2011, vol. 745 of *CEUR Workshop Proceedings*. CEUR-WS.org.

Rodriguez Muro, M., and Calvanese, D. 2011b. Semantic index: Scalable query answering without forward chaining or exponential rewritings. In *Proc. of the 10th Int. Semantic Web Conf., ISWC 2011*.

Rosati, R., and Almatelli, A. 2010. Improving query answering over DL-Lite ontologies. In *Proc. of the 12th Int. Conf. on Principles of Knowledge Representation and Reasoning, KR 2010*. AAAI Press.