

Long Rewritings, Short Rewritings

S. Kikot¹, R. Kontchakov¹, V. Podolskii², and M. Zakharyashev¹

¹ Department of Computer Science and Information Systems
Birkbeck, University of London, U.K.

{kikot, roman, michael}@dcs.bbk.ac.uk

² Steklov Mathematical Institute, Moscow, Russia
podolskii@mi.ras.ru

1 Introduction

An ontology language \mathcal{L} is said to enjoy *FO-rewritability* if any conjunctive query (CQ) q over any ontology \mathcal{T} , given in \mathcal{L} , can be transformed into an FO-formula q' such that, for any data \mathcal{A} , all answers to q over the knowledge base $(\mathcal{T}, \mathcal{A})$ can be found by querying q' over \mathcal{A} using a standard relational database management system (RDBMS). Ontology languages with this property include the *OWL 2 QL* profile of *OWL 2*, which is based on description logics of the *DL-Lite* family [7, 16, 2], and fragments of Datalog[±] such as linear or sticky TGDs [5, 6]. Various rewriting techniques have been implemented in the systems QuOnto [1], REQUIEM [15], Presto [22], Nyaya [8], IQAROS³ and Quest.⁴

The idea of using languages with FO-rewritability for ontology-based data access (OBDA) relies on the empirical fact that RDBMSs are usually very efficient in practice. However, the first rewritings of CQs over *OWL 2 QL* ontologies [7, 15] turned out to be too lengthy even for modern RDBMSs. The attempts to employ various optimisation techniques still produced rewritings of exponential size in the worst case: $O((|\mathcal{T}| \cdot |q|)^{|q|})$ [22, 8, 20, 21]. The alternative two-step combined approach [14, 13]—first expand the data by applying the ontology axioms to the data and introducing (some of) the missing individuals, and only then rewrite the query over the expanded data—resulted in a simple polynomial rewriting only for the fragment of *OWL 2 QL* without role inclusions; for the full language, the rewriting remained exponential. Two seemingly contradictory results, presented at DL 2011, added more spice to the quest for short rewritings: [9] showed that one can construct, in polynomial time, a nonrecursive Datalog (NDL) rewriting for some fragments of Datalog[±] containing *OWL 2 QL*, while [11] argued that no FO-rewriting for *OWL 2 QL* can be constructed in polynomial time.

The aim of this paper is twofold. First, we investigate the worst-case size of FO- and NDL-rewritings for CQs over *OWL 2 QL* ontologies. We distinguish between ‘pure’ rewritings, which can use the signature of the original query and ontology as well as $=$, \neq (cf. [7]), and ‘impure’ rewritings, where other means such as new constants are allowed. Here is a summary of the obtained results:

³ <http://code.google.com/p/iqaros/>

⁴ <http://obda.inf.unibz.it/protege-plugin/quest/quest.html>

- (1) An exponential blow-up is unavoidable for pure positive existential (PE) rewritings and pure NDL-rewritings; pure FO-rewritings can blow-up super-polynomially 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- and NDL-rewritings can always be made polynomial, and so they are exponentially more succinct than pure PE- and NDL-rewritings, respectively.

(1)–(3) are proved by establishing connections between pure rewritings for CQs over *OWL 2 QL* ontologies and circuits for monotone Boolean functions. In a nutshell, we show that CQs and *OWL 2 QL* ontologies can encode such problems as the existence of a k -clique in a graph with n vertices whose edges are given by a single-element ABox. The polynomial PE-rewriting in (4) is similar to the NDL-rewriting of [9]: using two extra constants, $=$ and polynomially many new existentially quantified variables, one can guess a relevant part of the canonical model of \mathcal{T} in the rewritten query. The difference between the resulting impure PE-rewritings and exponential-size pure PE-rewritings is of the same kind as the difference between deterministic and nondeterministic Boolean circuits.

Our second aim is to analyse the causes behind long rewritings and whether they occur in real-world queries and ontologies. As a result, we suggest some short rewritings that cover most practical cases.

Omitted proofs can be found in [10] and the full version of [12].

2 Queries over *OWL 2 QL* Ontologies

The language of *OWL 2 QL* is defined by the following grammar:⁵

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

where the A_i are concept names and the P_i are role names. An *OWL 2 QL TBox*, \mathcal{T} , is a finite set of *inclusions* of the form $B \sqsubseteq C$, $R_1 \sqsubseteq R_2$, $B_1 \sqcap B_2 \sqsubseteq \perp$ and $R_1 \sqcap R_2 \sqsubseteq \perp$. Note that concepts of the form $\exists R.B$ can only occur in the right-hand side of concept inclusions. An *ABox*, \mathcal{A} , is a finite set of *assertions* of the form $A_k(a_i)$ and $P_k(a_i, a_j)$, where a_i, a_j are individual names. \mathcal{T} and \mathcal{A} together form the *knowledge base* (KB) $\mathcal{K} = (\mathcal{T}, \mathcal{A})$. The semantics is defined in the usual way, based on interpretations $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$ with domain $\Delta^{\mathcal{I}}$ and interpretation function $\cdot^{\mathcal{I}}$. The set of individuals in \mathcal{A} is denoted by $\text{ind}(\mathcal{A})$; $\mathcal{I}_{\mathcal{A}}$ is the interpretation with domain $\text{ind}(\mathcal{A})$ such that, for any concept or role E and any $\mathbf{a} \subseteq \text{ind}(\mathcal{A})$, we have $\mathcal{I}_{\mathcal{A}} \models E(\mathbf{a})$ iff $E(\mathbf{a}) \in \mathcal{A}$. We write $E_1 \sqsubseteq_{\mathcal{T}} E_2$ if $\mathcal{T} \models E_1 \sqsubseteq E_2$; and we set $[E] = \{E' \mid E \sqsubseteq_{\mathcal{T}} E' \text{ and } E' \sqsubseteq_{\mathcal{T}} E\}$.

⁵ We do not consider data properties, attributes and role (ir)reflexivity constraints.

A *conjunctive query* (CQ) $\mathbf{q}(\mathbf{x})$ is a formula $\exists \mathbf{y} \varphi(\mathbf{x}, \mathbf{y})$, where φ is a conjunction of atoms of the form $A_k(t_1)$ and $P_k(t_1, t_2)$, and each t_i is a *term* (an individual or a variable from \mathbf{x}, \mathbf{y}). A tuple $\mathbf{a} \subseteq \text{ind}(\mathcal{A})$ is a *certain answer* to $\mathbf{q}(\mathbf{x})$ over $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ if $\mathcal{I} \models \mathbf{q}(\mathbf{a})$ for all $\mathcal{I} \models \mathcal{K}$; then we write $\mathcal{K} \models \mathbf{q}(\mathbf{a})$.

Query answering over OWL2QL 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 $\mathbf{q}(\mathbf{x})$ and $\mathbf{a} \subseteq \text{ind}(\mathcal{A})$, we have $\mathcal{K} \models \mathbf{q}(\mathbf{a})$ iff $\mathcal{C}_{\mathcal{K}} \models \mathbf{q}(\mathbf{a})$. 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} (we assume $\exists R$ is just a shorthand for $\exists R.\top$), 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 with $\mathcal{K} \models \exists R.B(w_{[R'B']})$ and it is not the case that $R' \sqsubseteq_{\mathcal{T}} R^-$ and $\mathcal{K} \models B(u)$.

If $a \rightsquigarrow w_{[R_1 B_1]} \rightsquigarrow \dots \rightsquigarrow w_{[R_n B_n]}$, $n \geq 0$, then we say that a *generates the path* $aw_{[R_1 B_1]} \dots w_{[R_n B_n]}$. 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:

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

Theorem 1 ([7, 13]). *For every OWL2QL KB $\mathcal{K} = (\mathcal{T}, \mathcal{A})$, every CQ $\mathbf{q}(\mathbf{x})$ and every $\mathbf{a} \subseteq \text{ind}(\mathcal{A})$, $\mathcal{K} \models \mathbf{q}(\mathbf{a})$ iff $\mathcal{C}_{\mathcal{K}} \models \mathbf{q}(\mathbf{a})$.*

Given a CQ $\mathbf{q}(\mathbf{x})$ and a TBox \mathcal{T} , a first-order formula $\mathbf{q}'(\mathbf{x})$, possibly with = and \neq , is called an *FO-rewriting for $\mathbf{q}(\mathbf{x})$ and \mathcal{T}* if, for any ABox \mathcal{A} and any $\mathbf{a} \subseteq \text{ind}(\mathcal{A})$, we have $(\mathcal{T}, \mathcal{A}) \models \mathbf{q}(\mathbf{a})$ iff $\mathcal{I}_{\mathcal{A}} \models \mathbf{q}'(\mathbf{a})$. If \mathbf{q}' is an FO-rewriting of the form $\exists \mathbf{y} \varphi(\mathbf{x}, \mathbf{y})$, where φ is built from atoms using only \wedge and \vee , then we call $\mathbf{q}'(\mathbf{x})$ a *positive existential rewriting for $\mathbf{q}(\mathbf{x})$ and \mathcal{T}* (or a *PE-rewriting*, for short). We say that \mathbf{q}' is *pure* if it does not contain constants that do not occur in \mathbf{q} (such constants are interpreted by fresh individuals added to $\mathcal{I}_{\mathcal{A}}$). The *size* $|\mathbf{q}'|$ of \mathbf{q}' is the number of symbols in \mathbf{q}' .

We also consider rewritings in the form of nonrecursive Datalog queries. We remind the reader that a *Datalog program*, Π , is a finite set of Horn clauses $\forall \mathbf{x} (A_1 \wedge \dots \wedge A_m \rightarrow A_0)$, where each A_i is an atom of the form $P(t_1, \dots, t_l)$ and each t_j is either a variable from \mathbf{x} or a constant. A_0 is called the *head* of the clause, and A_1, \dots, A_m its *body*. All variables occurring in the head must also occur in the body. A predicate P *depends* on a predicate Q in Π if Π contains a clause whose head is P and whose body contains Q . Π is called *nonrecursive* if

this dependence relation for Π is acyclic. A *nonrecursive Datalog query* consists of a nonrecursive Datalog program Π and a *goal* G , which is just a predicate. A tuple $\mathbf{a} \subseteq \text{ind}(\mathcal{A})$ is a *certain answer* to (Π, G) over \mathcal{A} if $\Pi, \mathcal{A} \models G(\mathbf{a})$. The *size* $|\Pi|$ of Π is the number of symbols in it. We distinguish between *pure* and *impure* Datalog queries [3]. In a *pure query* (Π, G) , the clauses in Π do not contain constant symbols in their heads. One reason for considering only pure queries in OBDA is that impure ones can add new facts to the database that do not follow from the background ontology. Impure queries are known to be more succinct than pure ones.

Given a CQ $\mathbf{q}(\mathbf{x})$ and an *OWL 2 QL* TBox \mathcal{T} , a pure nonrecursive Datalog query (Π, G) is called a *nonrecursive Datalog rewriting for $\mathbf{q}(\mathbf{x})$ and \mathcal{T}* (or an *NDL-rewriting*, for short) if, for any ABox \mathcal{A} and any $\mathbf{a} \subseteq \text{ind}(\mathcal{A})$, we have $(\mathcal{T}, \mathcal{A}) \models \mathbf{q}(\mathbf{a})$ iff $\Pi, \mathcal{A} \models G(\mathbf{a})$. If Π does not contain constants that do not occur in \mathbf{q} then we say that the NDL-rewriting (Π, G) is *pure*.

3 Query Rewritings and Boolean Circuits

To establish results (1)–(3) on the size of rewritings mentioned in the introduction, we show how the problem of constructing circuits that compute monotone Boolean functions can be reduced to the problem of finding pure FO- and NDL-rewritings for CQs over *OWL 2 QL* ontologies.

Our reduction proceeds in three steps. First, we take any family f^1, f^2, \dots of *monotone* Boolean functions in NP, where $f^n: \{0, 1\}^n \rightarrow \{0, 1\}$, a polynomial p and a family $\mathbf{C}^1, \mathbf{C}^2, \dots$ of polynomial-size circuits such that $f^n(\boldsymbol{\alpha}) = 1$ iff $\mathbf{C}^n(\boldsymbol{\alpha}, \boldsymbol{\beta}) = 1$, for some $\boldsymbol{\beta} \in \{0, 1\}^{p(n)}$. Using the Tseitin transformation [23], we construct a polynomial-size CNF θ_{f^n} that computes f^n in the following sense:

Lemma 1. *If f^n is monotone then $\varphi_{f^n}^\alpha = (\bigwedge_{\alpha_j=0} \neg x_j) \wedge \theta_{f^n}$ is satisfiable iff $f^n(\boldsymbol{\alpha}) = 1$, for all $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n) \in \{0, 1\}^n$.*

Let φ_{f^n} be $\varphi_{f^n}^\alpha$ for $\boldsymbol{\alpha} = (0, \dots, 0)$. It should be clear that $\varphi_{f^n}^\alpha$ is obtained from φ_{f^n} by removing the clauses $\neg x_j$ for which the j th component of $\boldsymbol{\alpha}$ is 1.

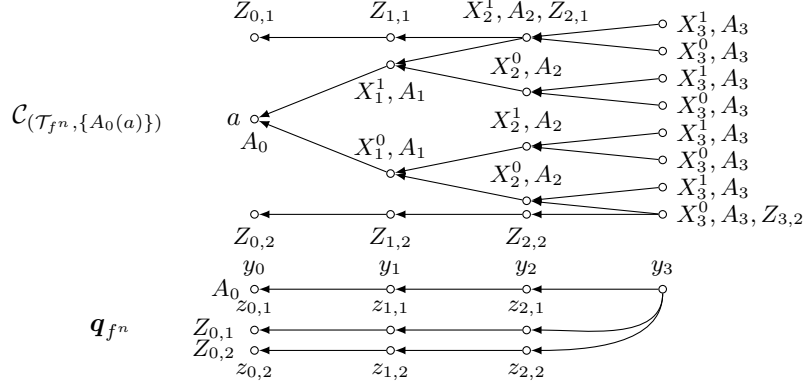
The second step is to encode the φ_{f^n} by means of TBoxes \mathcal{T}_{f^n} and CQs \mathbf{q}_{f^n} . Let p_1, \dots, p_N be the propositional variables and C_1, \dots, C_d the clauses in φ_{f^n} . Then \mathcal{T}_{f^n} contains the following concept inclusions, for $1 \leq i \leq N$, $1 \leq j \leq d$:

$$\begin{aligned} A_{i-1} \sqsubseteq \exists P^- . X_i^\ell, \quad X_i^\ell \sqsubseteq A_i, \quad \text{for } \ell = 0, 1, \quad X_i^0 \sqsubseteq Z_{i,j} \quad \text{if } \neg p_i \in C_j, \\ Z_{i,j} \sqsubseteq \exists P . Z_{i-1,j}, \quad X_i^1 \sqsubseteq Z_{i,j} \quad \text{if } p_i \in C_j, \\ A_0 \sqcap A_i \sqsubseteq \perp, \quad A_0 \sqcap \exists P \sqsubseteq \perp, \quad A_0 \sqcap Z_{i,j} \sqsubseteq \perp, \quad \text{if } (i, j) \notin \{(0, 1), \dots, (0, n)\}, \end{aligned}$$

and the CQ is defined as follows:

$$\mathbf{q}_{f^n} = \exists \mathbf{y} \exists \mathbf{z} \left[A_0(y_0) \wedge \bigwedge_{i=1}^N P(y_i, y_{i-1}) \wedge \bigwedge_{j=1}^d \left(P(y_N, z_{N-1,j}) \wedge \bigwedge_{i=1}^{N-1} P(z_{i,j}, z_{i-1,j}) \wedge Z_{0,j}(z_{0,j}) \right) \right],$$

where $\mathbf{y} = (y_0, \dots, y_N)$ and $\mathbf{z} = (z_{0,1}, \dots, z_{N-1,1}, \dots, z_{0,d}, \dots, z_{N-1,d})$. The size of \mathcal{T}_{f^n} and \mathbf{q}_{f^n} is $O(|\mathbf{C}^n|^2)$. Note that \mathcal{T}_{f^n} is acyclic and \mathbf{q}_{f^n} is tree-shaped and has no answer variables. The canonical model $\mathcal{C}_{(\mathcal{T}_{f^n}, \{A_0(a)\})}$ of $(\mathcal{T}_{f^n}, \{A_0(a)\})$ and the query \mathbf{q}_{f^n} are illustrated below.



For each $\alpha = (\alpha_1, \dots, \alpha_n) \in \{0, 1\}^n$, we define the following ABox:

$$\mathcal{A}_\alpha = \{A_0(a)\} \cup \{Z_{0,j}(a) \mid 1 \leq j \leq n \text{ and } \alpha_j = 1\}.$$

Lemma 2. $(\mathcal{T}_{f^n}, \mathcal{A}_\alpha) \models \mathbf{q}_{f^n}$ iff $\varphi_{f^n}^\alpha$ is satisfiable, for $\alpha \in \{0, 1\}^n$.

To complete our reduction, we show that rewritings for \mathbf{q}_{f^n} and \mathcal{T}_{f^n} can be turned into Boolean circuits computing f^n .

Lemma 3. (i) Suppose \mathbf{q}'_{f^n} is a pure FO- (PE-) rewriting for \mathcal{T}_{f^n} and \mathbf{q}_{f^n} . Then there is a (monotone) Boolean formula ψ_{f^n} computing f^n with $|\psi_{f^n}| \leq |\mathbf{q}'_{f^n}|$.

(ii) Suppose (Π_{f^n}, G) is a pure NDL-rewriting for \mathcal{T}_{f^n} and \mathbf{q}_{f^n} . Then there is a monotone Boolean circuit \mathbf{C}_{f^n} computing f^n with $|\mathbf{C}_{f^n}| \leq |\Pi_{f^n}|$.

The proof proceeds by eliminating quantifiers in the rewriting and replacing its predicates with propositional variables using the fact that, in ABoxes \mathcal{A}_α , these predicates can only be true on the individual a . Lemmas 1 and 2 ensure that the resulting Boolean formula or circuit computes f^n . The next lemma shows that circuits computing f^n can be turned into pure rewritings for \mathbf{q}_{f^n} and \mathcal{T}_{f^n} .

Lemma 4. (i) Suppose \mathbf{q}_n is an FO-sentence such that $(\mathcal{T}_{f^n}, \mathcal{A}_\alpha) \models \mathbf{q}_{f^n}$ iff $\mathcal{I}_{\mathcal{A}_\alpha} \models \mathbf{q}_n$, for all $\alpha \in \{0, 1\}^n$. Then there exists a pure FO-rewriting \mathbf{q}'_n for \mathbf{q}_{f^n} and \mathcal{T}_{f^n} with $|\mathbf{q}'_n| \leq |\mathbf{q}_n| + p(n)$, for a polynomial p .

(ii) Suppose (Π_n, G) is a pure NDL-query with a propositional goal G such that $(\mathcal{T}_{f^n}, \mathcal{A}_\alpha) \models \mathbf{q}_{f^n}$ iff $\Pi_n, \mathcal{A}_\alpha \models G$, for $\alpha \in \{0, 1\}^n$. Then there is a pure NDL-rewriting (Π'_n, G') for \mathbf{q}_{f^n} and \mathcal{T}_{f^n} with $|\Pi'_n| \leq |\Pi_n| + p(n)$, p a polynomial.

Now, results (1)–(3) formulated in the introduction can be obtained by applying Lemmas 1–4 to three concrete Boolean functions. For (1), we use the

function $\text{CLIQUE}_{n,k}$ of $n(n-1)/2$ variables e_{ij} , $1 \leq i < j \leq n$, which returns 1 iff the graph with vertices $\{1, \dots, n\}$ and edges $\{\{i, j\} \mid e_{ij} = 1\}$ contains a k -clique. A series of papers, started by Razborov's [19], gave an exponential lower bound for the size of monotone circuits computing $\text{CLIQUE}_{n,k}$: $2^{\Omega(\sqrt{k})}$ for $k \leq \frac{1}{4}(n/\log n)^{2/3}$. For monotone formulas, an even better lower bound is known: $2^{\Omega(k)}$ for $k = 2n/3$ [18]. The question whether $\text{CLIQUE}_{n,k}$ can be computed by a polynomial-size circuit is equivalent to whether $\text{NP} \subseteq \text{P/poly}$.

To show (2), we use the function GEN_{n^3} of n^3 variables x_{ijk} , $1 \leq i, j, k \leq n$, defined as follows. We say that 1 *generates* $k \leq n$ if either $k = 1$ or, for some i and j such that $x_{ijk} = 1$, 1 generates both i and j . $\text{GEN}_{n^3}(x_{111}, \dots, x_{nnn})$ returns 1 iff 1 generates n . It is clearly a monotone function computable by polynomial-size monotone circuits. On the other hand, any monotone formula computing GEN_{n^3} is of size at least 2^{n^ε} , for some $\varepsilon > 0$ [17].

For (3), we use the function MATCHING_{2n} of n^2 variables e_{ij} , $1 \leq i, j \leq n$, which returns 1 iff there is a *perfect matching* in a bipartite graph G with vertices $\{v_1^1, \dots, v_n^1, v_1^2, \dots, v_n^2\}$ and edges $\{\{v_i^1, v_j^2\} \mid e_{ij} = 1\}$, i.e., a subset E of edges in G such that every node in G occurs exactly once in E . An exponential lower bound $2^{\Omega(n)}$ for the size of monotone formulas computing this function was obtained in [18]. However, MATCHING_{2n} is computable by non-monotone formulas of size $n^{O(\log n)}$ [4].

The exponential bounds for the size of rewritings can be reduced to *polynomial* by using two extra constants, say 0 and 1, which are included in the domain of every $\mathcal{I}_{\mathcal{A}}$ (see [9] and [10] for polynomial-size NDL- and PE-rewritings, respectively). As these constants may not occur in the original query and intended ABoxes, we call such rewritings *impure* (in [9] and [10], 0, 1, = and polynomially-many fresh existentially quantified variables are used to guess some part of the canonical model). The difference between pure and impure rewritings is similar to the difference between deterministic circuits and nondeterministic circuits with additional existentially quantified input variables. For example, $\text{CLIQUE}_{n,k}$ is computed by a polynomial-size monotone nondeterministic circuit, but not by a monotone deterministic circuit of polynomial size [19].

4 Why are Pure Rewritings so Long?

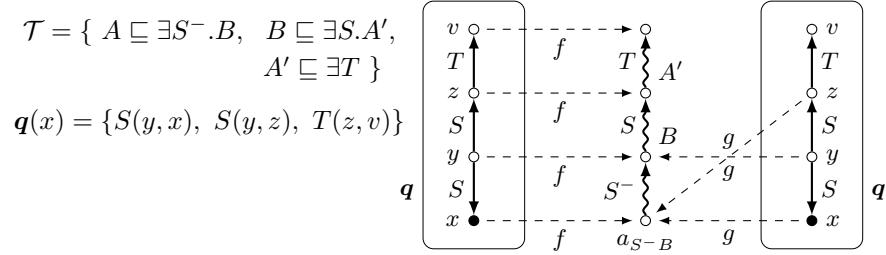
Let $\mathbf{q}(\mathbf{x}) = \exists \mathbf{y} \varphi(\mathbf{x}, \mathbf{y})$ be a connected CQ and \mathcal{T} a TBox. Let *term* \mathbf{q} be the set of terms in \mathbf{q} . Denote by $\mathcal{C}_{\mathcal{T}}$ the disjoint union of the canonical models for consistent $(\mathcal{T}, \{R(a_{RB}, b_{RB}), B(b_{RB})\})$, in which the generating relation is extended with $a_{RB} \rightsquigarrow b_{RB}$. 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}) . We also need the formulas

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

for concepts C and role names P . Suppose $\mathcal{C}_{\mathcal{K}} \models^{\mathbf{a}} \varphi$, $\mathcal{K} = (\mathcal{T}, \mathcal{A})$, 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 $x \in \mathbf{x}$. Consider an atom $P(t, t') \in \mathbf{q}$ with bound variables t, t' and assume that $\mathbf{a}(t_0) \in \text{ind}(\mathcal{A})$, for some $t_0 \in \text{term } \mathbf{q}$. The assignment \mathbf{a} can send t and t' to four different locations in $\mathcal{C}_{\mathcal{K}}$: (A) $\mathbf{a}(t), \mathbf{a}(t') \in \text{ind}(\mathcal{A})$; (B) $\mathbf{a}(t) \in \text{ind}(\mathcal{A})$, $\mathbf{a}(t') \notin \text{ind}(\mathcal{A})$; (B⁻) $\mathbf{a}(t) \notin \text{ind}(\mathcal{A})$, $\mathbf{a}(t') \in \text{ind}(\mathcal{A})$; (O) $\mathbf{a}(t), \mathbf{a}(t') \notin \text{ind}(\mathcal{A})$. Let us see how these alternatives can be reflected in a rewriting. In case (A) we have $\mathcal{C}_{\mathcal{K}} \models^{\mathbf{a}} P(t, t')$ iff $\mathcal{A} \models^{\mathbf{a}} \text{ext}_P(t, t')$. Case (B) is possible only if, for some concept $\exists R.B$, we have $\mathbf{a}(t) \rightsquigarrow w_{[RB]}$, $R \sqsubseteq_{\mathcal{T}} P$ and the atoms of \mathbf{q} ‘linked to’ t' can be mapped into the RB -subtree of $\mathcal{C}_{\mathcal{T}}$. To illustrate, consider an example.

Example 1. Let \mathcal{T} and \mathbf{q} be as in the picture below. 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 $\exists S^-B$. In the latter case, we have two ways of mapping z : either to $\mathbf{a}(x) \cdot w_{[S^-B]}w_{[SA']}$, in which case we must set $\mathbf{a}(v) = \mathbf{a}(x) \cdot w_{[S^-B]}w_{[SA']}w_{[T\top]}$, or to $\mathbf{a}(x)$, provided that $\mathbf{a}(x)$ is an instance of $\exists T$. Thus we have two (partial) maps f and g from \mathbf{q} into the S^-B -subtree of $\mathcal{C}_{\mathcal{T}}$ shown below.



These observations motivate our key definition. Given a pair (t, t') of adjacent terms in \mathbf{q} , a *tree witness*⁶ for (t, t') is a homomorphism f from the query $\mathbf{q}_f = \{E(s) \in \mathbf{q} \mid s \subseteq \text{dom } f, s \not\subseteq [t]_f\}$ to the RB -subtree of $\mathcal{C}_{\mathcal{T}}$, for some $\exists R.B$, such that $f(t) = a_{RB}$, $\text{dom } f$ is the smallest set containing t, t' for which $s' \in \text{dom } f$ whenever $S(s, s') \in \mathbf{q}$ with $s \in \text{dom } f \setminus [t]_f$, and all $s \in \text{dom } f \setminus [t]_f$ are bound variables in \mathbf{q} . Here \sim_f denotes the equivalence relation on $\text{dom } f$ defined by taking $s \sim_f s'$ iff $f(s) = f(s')$ and $[t]_f$ the equivalence class of t . In Example 1, f and g are two tree witnesses for (x, y) . Returning back to case (B), we can say now that there must exist a tree witness f for (t, t') such that $\mathbf{a}(t)$ satisfies the following *tree-witness formula* tw_f for f with $f(t) = a_{RB}$:

$$\text{tw}_f = \text{ext}_{\exists R.B}(t) \wedge \bigwedge_{s \in [t]_f} (t = s) \wedge \bigwedge_{E(s) \in \mathbf{q}, s \subseteq [t]_f} \text{ext}_{E(s)}.$$

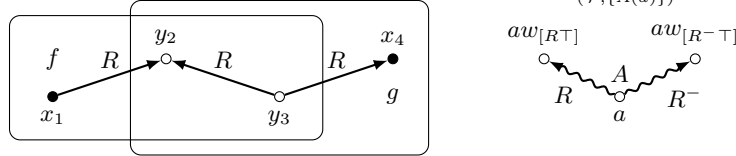
Case (B⁻) is symmetric, and in case (O) there must exist $S(s, s') \in \mathbf{q}$ for which (B) holds, with $P(t, t')$ being ‘covered’ by the tree witness for (s, s') .

This analysis suggests the following idea for a rewriting. We guess pairs of adjacent terms (t, t') in \mathbf{q} that will be mapped to edges of the tree part of $\mathcal{C}_{\mathcal{K}}$

⁶ A different notion of tree witness was used for $DL\text{-Lite}_{\text{horn}}^{\mathcal{N}}$ [13], where the structure of the canonical models ensured uniqueness of every tree witness.

and the tree witnesses that will cover them, as in cases (B), (B⁻) and (O). The part of \mathbf{q} that is not covered by these tree witnesses will be mapped to $\text{ind}(\mathcal{A})$, as in case (A). The query representing the guesses is then evaluated over $\mathcal{I}_{\mathcal{A}}$. The following example shows, however, that this idea needs a refinement.

Example 2. Let $\mathbf{q}(x_1, x_4) = \{R(x_1, y_2), R(y_3, y_2), R(y_3, x_4)\}$, shown below on the left, and $\mathcal{T} = \{A \sqsubseteq \exists R, A \sqsubseteq \exists R^-\}$.



Clearly, there is a tree witness f for (x_1, y_2) with $\text{dom } f = \{x_1, y_2, y_3\}$ and $[x_1]_f = \{x_1, y_3\}$, and a tree witness g for (x_4, y_3) with $\text{dom } g = \{x_4, y_3, y_2\}$ and $[x_4]_g = \{x_4, y_2\}$. Although these tree witnesses cover the whole query \mathbf{q} , they are only ‘realised,’ say, in the canonical model $\mathcal{C}_{(\mathcal{T}, \{A(a)\})}$ (shown above on the right) under *conflicting* maps: f sends x_1, y_3 to a and y_2 to $\text{aw}_{[R\top]}$, while g sends x_4, y_2 to a and y_3 to $\text{aw}_{[R^-\top]}$; in fact, $(\mathcal{T}, \{A(a)\}) \not\models \mathbf{q}(a, a)$.

Tree witnesses f and g , for (t, t') and (s, s') , respectively, are *compatible* if $\text{dom } f \cap \text{dom } g \subseteq [t]_f \cap [s]_g$. 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* (e.g., f and g in Example 2). A set Ξ of tree witnesses is called *consistent* if all pairs of tree witnesses in Ξ are compatible. Let

$$\mathbf{q}_e(\mathbf{x}) = \text{detached}_{\mathbf{q}} \vee \bigvee_{\Xi \text{ consistent}} \exists \mathbf{y} \left(\bigwedge_{f \in \Xi} \text{tw}_f \wedge \bigwedge_{\substack{E(\mathbf{s}) \in \mathbf{q} \\ \mathbf{s} \not\subseteq \text{dom } f, \text{ for all } f \in \Xi}} \text{ext}_E(\mathbf{s}) \right),$$

where $\text{detached}_{\mathbf{q}} = \perp$ if $\mathbf{x} \neq \emptyset$; otherwise it is a disjunction of the sentences $\exists \mathbf{x} \text{ext}_{\exists R.B}(\mathbf{x})$ such that there is a homomorphism from \mathbf{q} to the RB -subtree of $\mathcal{C}_{\mathcal{T}}$. The next theorem shows that \mathbf{q}_e is a pure PE-rewriting for \mathbf{q} and \mathcal{T} :

Theorem 2. *For all \mathcal{A} and $\mathbf{a} \subseteq \text{ind}(\mathcal{A})$, we have $(\mathcal{T}, \mathcal{A}) \models \mathbf{q}(\mathbf{a})$ iff $\mathcal{I}_{\mathcal{A}} \models \mathbf{q}_e(\mathbf{a})$.*

Example 3. Let $\mathbf{q}(x) = \{R_i(x, y_i) \mid i \leq n\}$ 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 $\mathbf{q}_e = \bigvee_{N \subseteq [0, n]} \exists \mathbf{y} (\bigwedge_{i \in N} \text{tw}_{f_i} \wedge \bigwedge_{j \notin N} R_j(x, y_j))$.

The size of \mathbf{q}_e is $O((n_{\mathcal{T}, \mathbf{q}} + 1)^{|\mathbf{q}|} \cdot |\mathcal{T}| \cdot |\mathbf{q}|^2)$, where $n_{\mathcal{T}, \mathbf{q}}$ is the number of distinct tree witnesses. Now, we observe that if

(conf) *there are no conflicting tree witnesses for \mathbf{q} and \mathcal{T}*

then \mathbf{q}_e can be transformed to the query

$$\mathbf{q}_c(\mathbf{x}) = \text{detached}_{\mathbf{q}} \vee \exists \mathbf{y} \bigwedge_{\substack{\{t, t'\} \\ t, t' \text{ adjacent}}} \left[\bigwedge_{E(\mathbf{s}) \in \mathbf{q}, \mathbf{s} \subseteq \{t, t'\}} \text{ext}_E(\mathbf{s}) \vee \bigvee_{\substack{f \text{ is a tree witness} \\ t, t' \in \text{dom } f}} \text{tw}_f \right]$$

(if \mathbf{q} has no binary predicates then $\mathbf{q}_c = \mathbf{q}_e$).

Theorem 3. *If \mathcal{T} and \mathbf{q} satisfy **(conf)** then, for any ABox \mathcal{A} and $\mathbf{a} \subseteq \text{ind}(\mathcal{A})$, we have $(\mathcal{T}, \mathcal{A}) \models \mathbf{q}(\mathbf{a})$ iff $\mathcal{I}_{\mathcal{A}} \models \mathbf{q}_c(\mathbf{a})$.*

The size of \mathbf{q}_c is $O(n_{\mathcal{T}, \mathbf{q}} \cdot |\mathcal{T}| \cdot |\mathbf{q}|^2)$. So, if **(conf)** holds and $n_{\mathcal{T}, \mathbf{q}}$ is polynomial then \mathbf{q}_c is a *polynomial* pure PE-rewriting of \mathbf{q} and \mathcal{T} . For instance, the exponential \mathbf{q}_e of Example 3 reduces to polynomial $\mathbf{q}_c = \exists \mathbf{y} \bigwedge_{i \leq n} (R_i(x, y_i) \vee \text{tw}_{f_i})$. Note, however, that the CQs and TBoxes used in Section 3 generate *exponentially many* distinct tree witnesses.

We show now that, for a large class of CQs \mathbf{q} and TBoxes \mathcal{T} , all tree witnesses f for each (t, t') in \mathbf{q} (even if there are exponentially many of them) can be represented by a polynomial formula. Observe that each tw_f is determined by a concept $\exists R.B$ (such that the RB -subtree of $\mathcal{C}_{\mathcal{T}}$ contains the range of f) and the equivalence relation \sim_f on $\text{dom } f$. As the number of concepts $\exists R.B$ is linear in $|\mathcal{T}|$, the rewriting \mathbf{q}_c may be regarded polynomial if we show that all tree witnesses for each (t, t') have the same equivalence relation. For example, 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 tree witnesses for (x, y) , as each z_i can be mapped either to a B_1 - or a B_2 -point in $\mathcal{C}_{\mathcal{T}}$, and yet they all define the same equivalence relation and the same tree-witness formula $\text{ext}_{\exists S}(x)$.

We formalise this intuition in the following definition. For a pair (t, t') of adjacent terms in \mathbf{q} , we call $\mathbf{f} = (\text{dom } \mathbf{f}, \sim_{\mathbf{f}})$ a *universal tree witness* for (t, t') if $\text{dom } \mathbf{f}$ is a subset of $\text{term } \mathbf{q}$ with $t, t' \in \text{dom } \mathbf{f}$ and $\sim_{\mathbf{f}}$ is an equivalence relation on $\text{dom } \mathbf{f}$ such that $\mathbf{q}_{\mathbf{f}} = \{E(\mathbf{s}) \in \mathbf{q} \mid \mathbf{s} \subseteq \text{dom } \mathbf{f}, \mathbf{s} \not\subseteq [t]_{\mathbf{f}}\} / \sim_{\mathbf{f}}$ is a tree-shaped query with root $[t]_{\mathbf{f}}$ and, for every tree witness g for (t, t') ,

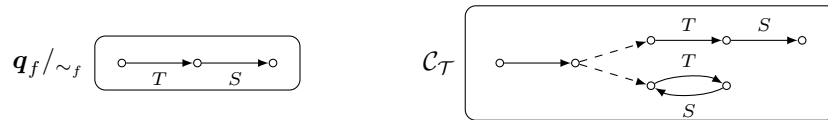
- $\text{dom } g = \text{dom } \mathbf{f}$ and there exists a homomorphism $h: \mathbf{q}_{\mathbf{f}} \rightarrow \mathcal{C}_{\mathcal{T}}$ that preserves the distance from the root and $g(\mathbf{s}) = h([s]_{\mathbf{f}})$, for every $\mathbf{s} \in \text{dom } \mathbf{f}$.

A universal tree witness \mathbf{f} for (t, t') is not a tree witness in the sense of our original definition, but rather a convenient structure representing all tree witnesses for (t, t') : we can merge the tree-witness formulas for (t, t') into one formula

$$\text{tw}_{\mathbf{f}} = \left[\bigvee_{\exists R.B \in \Phi_{\mathbf{f}}} \text{ext}_{\exists R.B}(t) \right] \wedge \bigwedge_{\mathbf{s} \in [t]_{\mathbf{f}}} (t = \mathbf{s}) \wedge \bigwedge_{E(\mathbf{s}) \in \mathbf{q}, \mathbf{s} \subseteq [t]_{\mathbf{f}}} \text{ext}_E(\mathbf{s}),$$

where $\Phi_{\mathbf{f}}$ is the set of concepts $\exists R.B$ such there is a homomorphism h from $\mathbf{q}_{\mathbf{f}}$ to the RB -subtree of $\mathcal{C}_{\mathcal{T}}$ with $h([t]_{\mathbf{f}}) = a_{RB}$.

We now identify a class of CQs and TBoxes for which a universal tree witness is unique (if exists) and can be constructed in time polynomial in $|\mathbf{q}|$ and $|\mathcal{T}|$. Intuitively, for each tree witness f , we disallow situations in which $\mathcal{C}_{\mathcal{T}}$ and the quotient $\mathbf{q}_{\mathbf{f}} / \sim_{\mathbf{f}}$ of $\mathbf{q}_{\mathbf{f}}$ simultaneously contain fragments of the form



We say that a role S is *forward* if $u \rightsquigarrow v$ for all $(u, v) \in S^{\mathcal{C}\tau}$. If neither S nor its inverse S^- is forward then S is said to be a *twisty role*. A tree witness f for (t, t') is called *perfect* if, for all $T(s_1, s_2), S(s_2, s_3) \in \mathbf{q}_f / \sim_f$ such that $s_2 \neq [t]_f$ and S is twisty, we have $\mathcal{C}\tau \not\models \text{inv}(T, S) \wedge \text{suc}(T, S)$, where

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

Lemma 5. *There is a polynomial-time algorithm which, given \mathbf{q} , \mathcal{T} and a pair (t, t') , checks whether all tree witnesses for (t, t') are perfect, and if this is the case, returns a unique universal tree witness $\mathbf{f}_{t, t'}$ for (t, t') .*

Note that even though there may be exponentially many tree witnesses for (t, t') , the algorithm checks whether they all are perfect in polynomial time. We are now in a position to define our *polynomial* PE-rewriting \mathbf{q}_p for \mathbf{q} and \mathcal{T} :

$$\mathbf{q}_p(\mathbf{x}) = \text{detached}_{\mathbf{q}} \vee \exists \mathbf{y} \bigwedge_{\substack{\{t, t'\} \\ t, t' \text{ adjacent}}} \left[\bigwedge_{E(s) \in \mathbf{q}, s \subseteq \{t, t'\}} \text{ext}_E(s) \vee \bigvee_{\substack{s, s' \text{ adjacent} \\ t, t' \in \text{dom } \mathbf{f}_{s, s'}}} \text{tw}_{\mathbf{f}_{s, s'}} \right].$$

Theorem 4. *If all tree witnesses for \mathbf{q} and \mathcal{T} are perfect and condition (conf) is satisfied then, for any ABox \mathcal{A} and any $\mathbf{a} \subseteq \text{ind}(\mathcal{A})$, we have $(\mathcal{T}, \mathcal{A}) \models \mathbf{q}(\mathbf{a})$ iff $\mathcal{I}_{\mathcal{A}} \models \mathbf{q}_p(\mathbf{a})$. Moreover, \mathbf{q}_p is constructed in time polynomial in $|\mathbf{q}|$ and $|\mathcal{T}|$.*

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

Theorem 5. *Let \mathcal{T} be an OWL2QL ontology without twisty roles. Then, for any CQ \mathbf{q} , there are no conflicting tree witnesses for \mathbf{q} and \mathcal{T} . Thus, \mathbf{q}_p is a pure PE-rewriting for \mathbf{q} and \mathcal{T} and it can be constructed in polynomial time.*

Note that OWL2EL ontologies satisfy this condition, and so a polynomial rewriting similar to \mathbf{q}_p can also be used for CQ answering over such ontologies provided that the ABoxes are complete with respect to the ontologies [12].

5 Conclusions

As we saw in Section 3, pure PE- and NDL-rewritings of CQs over OWL2QL ontologies are of exponential size in the worst case. The analysis in Section 4 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 intensional tree 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 inclusions in ontologies and queries, we can guarantee transparent polynomial rewritings. Moreover, as shown by a series of experiments [12], real-world ontologies and CQs contain very few tree witnesses, which are never in conflict, satisfy the above mentioned restrictions, and so enjoy pure, polynomial PE-rewritings.

References

1. Acciarri, A., Calvanese, D., De Giacomo, G., Lembo, D., Lenzerini, M., Palmieri, M., Rosati, R.: QUONTO: Querying ontologies. In: Proc. of the 20th Nat. Conf. on AI, AAAI. pp. 1670–1671 (2005)
2. Artale, A., Calvanese, D., Kontchakov, R., Zakharyashev, M.: The DL-Lite family and relations. *Journal of Artificial Intelligence Research (JAIR)* 36, 1–69 (2009)
3. Benedikt, M., Gottlob, G.: The impact of virtual views on containment. *PVLDB* 3(1), 297–308 (2010)
4. Borodin, A., von zur Gathen, J., Hopcroft, J.E.: Fast parallel matrix and gcd computations. In: Proc. of FOCS. pp. 65–71 (1982)
5. Cali, A., Gottlob, G., Lukasiewicz, T.: A general Datalog-based framework for tractable query answering over ontologies. In: Proc. of PODS. pp. 77–86 (2009)
6. Cali, A., Gottlob, G., Pieris, A.: Advanced processing for ontological queries. *PVLDB* 3(1), 554–565 (2010)
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. of Automated Reasoning* 39(3), 385–429 (2007)
8. Gottlob, G., Orsi, G., Pieris, A.: Ontological queries: Rewriting and optimization. In: Proc. of the IEEE Int. Conf. on Data Engineering, ICDE (2011)
9. Gottlob, G., Schwentick, T.: Rewriting ontological queries into small nonrecursive Datalog programs. In: Proc. of DL. vol. 745. CEUR-WS.org (2011)
10. Kikot, S., Kontchakov, R., Podolskii, V., Zakharyashev, M.: Exponential lower bounds and separation for query rewriting. CoRR, arXiv:1202.4193, 2012.
11. Kikot, S., Kontchakov, R., Zakharyashev, M.: On (in)tractability of OBDA with OWL 2 QL. In: Proc. of DL. vol. 745. CEUR-WS.org (2011)
12. Kikot, S., Kontchakov, R., Zakharyashev, M.: Conjunctive query answering with OWL 2 QL. In: Proc. of KR. AAAI Press (2012) (see www.dcs.bbk.ac.uk/~kikot)
13. Kontchakov, R., Lutz, C., Toman, D., Wolter, F., Zakharyashev, M.: The combined approach to query answering in DL-Lite. In: Proc. of KR. AAAI Press (2010)
14. Lutz, C., Toman, D., Wolter, F.: Conjunctive query answering in the description logic EL using a relational database system. In: Proceedings of the 21st Int. Joint Conf. on Artificial Intelligence, IJCAI 2009. pp. 2070–2075 (2009)
15. Pérez-Urbina, H., Motik, B., Horrocks, I.: A comparison of query rewriting techniques for DL-Lite. In: Proc. of DL. vol. 477. CEUR-WS.org (2009)
16. Poggi, A., Lembo, D., Calvanese, D., De Giacomo, G., Lenzerini, M., Rosati, R.: Linking data to ontologies. *J. on Data Semantics X*, 133–173 (2008)
17. Raz, R., McKenzie, P.: Separation of the monotone NC hierarchy. In: Proc. of FOCS. pp. 234–243 (1997)
18. Raz, R., Wigderson, A.: Monotone circuits for matching require linear depth. *J. ACM* 39(3), 736–744 (1992)
19. Razborov, A.: Lower bounds for the monotone complexity of some Boolean functions. *Dokl. Akad. Nauk SSSR* 281(4), 798–801 (1985)
20. Rodríguez-Muro, M., Calvanese, D.: Dependencies to optimize ontology based data access. In: Proc. of DL. vol. 745. CEUR-WS.org (2011)
21. Rodríguez-Muro, M., Calvanese, D.: Semantic index: Scalable query answering without forward chaining or exponential rewritings. In: Proc. of ISWC (2011)
22. Rosati, R., Almatelli, A.: Improving query answering over DL-Lite ontologies. In: Proc. of the 12th Int. Conf. KR. AAAI Press (2010)
23. Tseitin, G.: On the complexity of derivation in propositional calculus. In: *Automation of Reasoning 2: Classical Papers on Computational Logic 1967–1970* (1983)