

Optimal Nonrecursive Datalog Rewritings of Linear TGDs and Bounded (Hyper)Tree-Width Queries

M. Bienvenu¹, S. Kikot², R. Kontchakov², V. Ryzhikov³, and M. Zakharyashev²

¹ CNRS & University of Montpellier, France (meghyn@lirmm.fr)

² Birkbeck, University of London, UK ({kikot, roman, michael}@dcs.bbk.ac.uk)

³ Free University of Bozen-Bolzano, Italy (ryzhikov@inf.unibz.it)

Abstract. Our concern is answering ontology-mediated queries (\mathcal{O}, q) , where \mathcal{O} is a set of linear tgds and q a conjunctive query (CQ) of bounded hypertree width. Assuming that the arity of predicates is bounded, we show that polynomial-size nonrecursive Datalog rewritings can be constructed and executed in (i) LOGCFL for OMQs with ontologies of bounded existential depth; (ii) NL for OMQs with ontologies of bounded depth and CQs whose hypertree decompositions have a bounded number of leaves; (iii) LOGCFL for OMQs with acyclic CQs whose join trees have a bounded number of leaves.

1 Introduction

As shown in [3, 4, 13], the optimal combined complexity (LOGCFL and NL) of answering ontology-mediated queries (OMQs) with *OWL 2 QL* ontologies of bounded depth and conjunctive queries (CQs) of bounded treewidth can be achieved by means of rewriting them into nonrecursive datalog (NDL) queries, though not via positive-existential rewritings. (Note that in these cases the complexity of OMQs matches the complexity of evaluating the underlying CQs.) Our recent experiments have demonstrated that such NDL rewritings, reformulated as Spark SQL queries with views, are efficiently executed by Apache Spark taking advantage of their parallelisable structure.

The aim of this paper is to extend the above mentioned results to ontologies and CQs with predicates of *arbitrary fixed arity*. We consider ontologies that consist of linear TGDs (linear existential rules or atomic-hypothesis rules) [2, 6, 11, 12], which are instances of finite unification sets [18, 15, 16]. Our interest in this problem is also motivated by the system ETAP [5] designed to answer natural language questions by translating them into SPARQL and executing—along with background knowledge—over RDF data extracted from texts. To illustrate, suppose that the data contains the atoms $Purchased(j, c)$ and $Car(c)$ representing the sentence ‘John purchased a car’. To answer the question ‘has a car been sold?’ ETAP utilises the ontology rules (with omitted universal quantifiers)

$$\begin{aligned} Purchased(x, y) &\rightarrow \exists v z (Purchase(v) \wedge hasAgent1(v, x) \wedge \\ &\quad hasObject(v, y) \wedge hasAgent2(v, z)), \\ Purchase(v) \wedge hasAgent1(v, x) \wedge hasObject(v, y) \wedge hasAgent2(v, z) &\rightarrow \\ &\quad \exists v' (Sale(v') \wedge hasAgent1(v', z) \wedge hasObject(v', x) \wedge hasAgent2(v', y)), \end{aligned}$$

where v and v' represent the acts of purchase and sale, respectively. The rules are clearly beyond the limitations of *OWL 2 QL*; however, the knowledge they represent can also be captured by means of linear TGDs with ternary predicates:

$$\text{Purchased}(x, y) \rightarrow \exists z \text{Purchase}(x, y, z), \quad \text{Purchase}(x, y, z) \rightarrow \text{Sale}(z, y, x),$$

which are enough to answer the query $\exists xyz (\text{Car}(y) \wedge \text{Sale}(x, y, z))$.

We classify OMQs $Q = (\mathcal{O}, q)$ with linear TGDs and predicates of any *fixed arity* $n < \omega$ along three axes: (1) the existential depth d of \mathcal{O} , that is, the maximal depth of Skolem terms in the chases of \mathcal{O} over arbitrary data (cf. [8]), (2) the hypertree width t of q , and (3) the number ℓ of leaves in the tree underlying a hypertree decomposition of q . Thus, $\text{OMQ}(p_1, p_2, p_3)$ denotes the class of OMQs in which parameter (i) is bounded by $p_i \in \mathbb{N} \cup \{\infty\}$. We show that, for any fixed $d, t, \ell < \omega$, answering OMQs in the classes $\text{OMQ}(d, t, \infty)$ and $\text{OMQ}(\infty, 1, \ell)$ can be done in LOGCFL (for combined complexity) by means of NDL-rewritings, and even in NL for $\text{OMQ}(d, t, \ell)$. On the other hand, one can show that answering OMQs in $\text{OMQ}(\infty, t, \ell)$, for $t, \ell \geq 2$, is NP-hard by observing that the sequence of tree-shaped CQs from the proof of [3, Theorem 20] is of path width 2. Thus, we obtain a full classification of the classes $\text{OMQ}(p_1, p_2, p_3)$, for $p_i \in \mathbb{N} \cup \{\infty\}$, with respect to combined complexity.

2 Preliminaries

Ontology-mediated queries. Let Σ be a *relational schema* with the maximum arity $ar(\Sigma)$ of its predicates bounded by n . By writing $P(\mathbf{x})$, for a predicate name P and an n -tuple \mathbf{x} of variables (with possible repetitions), we mean that P is n -ary. By writing $\gamma(\mathbf{x})$, we mean that the free variables of formula γ are \mathbf{x} , where \mathbf{x} contains no repetitions. If the meaning is clear from the context, we use set-theoretic notation for lists.

A *data instance*, \mathcal{D} , over Σ is any finite set of ground atoms $P(\mathbf{a})$ with predicate symbols P from Σ . We denote by $\text{ind}(\mathcal{D})$ the set of individual constants in \mathcal{D} . An *ontology* is any finite set, \mathcal{O} , of sentences of the form

$$\forall \mathbf{x} (\gamma_0(\mathbf{x}) \rightarrow \exists \mathbf{y} \gamma_1(\mathbf{x}', \mathbf{y})) \quad \text{and} \quad \forall \mathbf{x} (\gamma_0(\mathbf{x}) \rightarrow \gamma_2(\mathbf{x}')),$$

where γ_0, γ_1 and γ_2 are atoms with predicate symbols from Σ and $\mathbf{x}' \subseteq \mathbf{x}$, for disjoint sets \mathbf{x} and \mathbf{y} of variables. When writing rules, we omit the universal quantifiers.

An *ontology-mediated query* (OMQ) $Q(\mathbf{x})$ is a pair $(\mathcal{O}, q(\mathbf{x}))$, in which \mathcal{O} is an ontology and $q(\mathbf{x})$ a *conjunctive query* (CQ), that is, a formula of the form $\exists \mathbf{y} \varphi(\mathbf{x}, \mathbf{y})$, where φ is a conjunction of atoms $P(\mathbf{z})$ over Σ with $\mathbf{z} \subseteq \mathbf{x} \cup \mathbf{y}$. A tuple $\mathbf{a} \in \text{ind}(\mathcal{D})^{|\mathbf{x}|}$ is a *certain answer* to $Q(\mathbf{x})$ over \mathcal{D} if $\mathfrak{M} \models q(\mathbf{a})$, for every model \mathfrak{M} of $\mathcal{O} \cup \mathcal{D}$; in this case we write $\mathcal{O}, \mathcal{D} \models q(\mathbf{a})$. If the list \mathbf{x} of *answer variables* is empty, a *certain answer* to Q over \mathcal{D} is ‘yes’ if $\mathfrak{M} \models q$, for every model \mathfrak{M} of $\mathcal{O} \cup \mathcal{D}$, and ‘no’ otherwise. OMQs and CQs without answer variables are called *Boolean*. We often regard CQs as *sets* of their atoms. We abuse notation and use sets of variables in place of sequences assuming that they are ordered in some (fixed) way. Also, given $\mathbf{c} \in \text{ind}(\mathcal{D})^{|\mathbf{z}|}$ and $z \in \mathbf{z}$, we write $\mathbf{c}(z)$ to refer to the component of \mathbf{c} that corresponds to z .

Canonical models. An important property of tgds is the fact [1] that, for any \mathcal{O} and \mathcal{D} , there is a (possibly infinite) *canonical* (or *universal*) *model* $\mathcal{C}_{\mathcal{O}, \mathcal{D}}$ such that, for every

CQ $q(\mathbf{x})$ and $\mathbf{a} \in \text{ind}(\mathcal{D})^{|\mathbf{x}|}$, we have $\mathcal{O}, \mathcal{D} \models q(\mathbf{a})$ iff $\mathfrak{C}_{\mathcal{O}, \mathcal{D}} \models q(\mathbf{a})$. Such a canonical model can be constructed by the following (*oblivious*) *chase procedure* that, intuitively, ‘repairs’ \mathcal{D} with respect to \mathcal{O} (though not in the most economical way). With each rule ϱ of the form $\gamma_0(\mathbf{x}) \rightarrow \exists \mathbf{y} \gamma_1(\mathbf{x}', \mathbf{y})$, where $\mathbf{x} = (x_1, \dots, x_n)$, $\mathbf{y} = (y_1, \dots, y_k)$ and $k > 0$, we associate the k -tuple $\mathbf{s}_\varrho = (s_\varrho^1, \dots, s_\varrho^k)$ of distinct n -ary *Skolem function* symbols. An *application of ϱ to \mathcal{D}* under a map $h: \mathbf{x} \rightarrow \text{ind}(\mathcal{D})$ such that $h(\gamma_0) \in \mathcal{D}$ adds $h'(\gamma_1)$ to \mathcal{D} , where h' is defined by taking $h'(x_i) = h(x_i)$, for $1 \leq i \leq n$, and $h'(y_j) = s_\varrho^j(h(\mathbf{x}))$, for $1 \leq j \leq k$. An application of a rule $\gamma_0(\mathbf{x}) \rightarrow \gamma_2(\mathbf{x}')$ to \mathcal{D} under such an h adds $h(\gamma_2)$ to \mathcal{D} . The *chase algorithm* applies these two rules exhaustively to \mathcal{O} and \mathcal{D} in a breadth-first manner. More precisely, we set $\mathfrak{C}_{\mathcal{O}, \mathcal{D}}^0 = \mathcal{D}$ and say that the atoms in $\mathfrak{C}_{\mathcal{O}, \mathcal{D}}^0$ are of (*derivation*) *level* 0. Assuming that $\mathfrak{C}_{\mathcal{O}, \mathcal{D}}^{n-1}$ has already been constructed, we define $\mathfrak{C}_{\mathcal{O}, \mathcal{D}}^n$ as follows. Take some enumeration of all distinct pairs (ϱ_i, h_i) such that $\varrho_i \in \mathcal{O}$ with γ_i on the left-hand side is applicable to $\mathfrak{C}_{\mathcal{O}, \mathcal{D}}^{n-1}$ under h_i . If none of the atoms in the $h_i(\gamma_i)$ is of level $n-1$, then we set $\mathfrak{C}_{\mathcal{O}, \mathcal{D}}^n = \mathfrak{C}_{\mathcal{O}, \mathcal{D}}^{n-1}$. Otherwise, we apply the ϱ_i under h_i to $\mathfrak{C}_{\mathcal{O}, \mathcal{D}}^{n-1}$ one after the other and say that the newly added atoms are of (*derivation*) *level* n ; the resulting extension of $\mathfrak{C}_{\mathcal{O}, \mathcal{D}}^{n-1}$ is denoted by $\mathfrak{C}_{\mathcal{O}, \mathcal{D}}^n$. The *canonical model* $\mathfrak{C}_{\mathcal{O}, \mathcal{D}}$ is then the union of all $\mathfrak{C}_{\mathcal{O}, \mathcal{D}}^n$, for $n < \omega$.

The domain $\Delta^{\mathfrak{C}_{\mathcal{O}, \mathcal{D}}}$ of $\mathfrak{C}_{\mathcal{O}, \mathcal{D}}$ consists of terms built from the constants in \mathcal{D} using Skolem functions s_ϱ^j , for $\varrho \in \mathcal{O}$. The *depth* of such a term is the maximal number of nested occurrences of function symbols in it. We say that \mathcal{O} is of *depth* $k \leq \omega$ if k is the minimal ordinal such that $\Delta^{\mathfrak{C}_{\mathcal{O}, \mathcal{D}}}$ contains no terms of depth $> k$, for any data \mathcal{D} . For an ontology \mathcal{O} and a ground atom $P(\mathbf{a})$, we set $\text{term}_{\mathcal{O}}(P(\mathbf{a})) = \Delta^{\mathfrak{C}_{\mathcal{O}, \mathcal{D}} \setminus \{P(\mathbf{a})\}} \setminus \mathbf{a}$. We denote by $\text{term}_{\mathcal{O}}$ the union of $\text{term}_{\mathcal{O}}(P(\mathbf{a}))$, for all possible (up to renaming the constants) atoms $P(\mathbf{a})$ with predicates in \mathcal{O} (assuming that distinct $P(\mathbf{a})$ do not share constants). It should be clear that \mathcal{O} is of finite depth iff $\text{term}_{\mathcal{O}}$ is finite. By counting the number of possible linear derivations of Skolem terms, we see that, for \mathcal{O} of depth k , $|\text{term}_{\mathcal{O}}| \leq (ar(\Sigma)^{ar(\Sigma)}|\mathcal{O}|)^k$. We assume that any constant a occurring in $\text{term}_{\mathcal{O}}$ has a *twin variable* \tilde{a} and denote by $\tilde{\mathbf{a}}$ the result of replacing all constants in \mathbf{a} with their twin variables. Given a tuple $\mathbf{b} \subseteq \text{ind}(\mathcal{D})$, we denote by $\mathbf{a}/\mathbf{b}(\tilde{\mathbf{a}})$ the substitution that maps each a in \mathbf{a} to the corresponding $b(\tilde{a})$ in \mathbf{b} .

NDL-rewritings. A *datalog program*, Π , is a finite set of Horn clauses of the form $\forall \mathbf{z} (\gamma_0 \leftarrow \gamma_1 \wedge \dots \wedge \gamma_m)$, where each γ_i is an atom $Q(\mathbf{y})$ with $\mathbf{y} \subseteq \mathbf{z}$ or an equality $(z = z')$ with $z, z' \in \mathbf{z}$. (As usual, we omit $\forall \mathbf{z}$ from clauses.) The atom γ_0 is the *head* of the clause, and $\gamma_1, \dots, \gamma_m$ its *body*. All variables in the head must occur in the body, and $=$ can only occur in the body. The predicates in the heads of clauses in Π are *IDB predicates*, the rest (including $=$) *EDB predicates*. A predicate Q *depends* on P in Π if Π has a clause with Q in the head and P in the body. Π is a *nonrecursive datalog (NDL) program* if the (directed) *dependence graph* of the dependence relation is acyclic. The size $|\Pi|$ of Π is the number of symbols in it. An *NDL query* is a pair $(\Pi, G(\mathbf{x}))$, where Π is an NDL program and G a predicate. A tuple $\mathbf{a} \in \text{ind}(\mathcal{D})^{|\mathbf{x}|}$ is an *answer to $(\Pi, G(\mathbf{x}))$* over a data instance \mathcal{D} if $G(\mathbf{a})$ holds in the first-order structure with domain $\text{ind}(\mathcal{D})$ obtained by closing \mathcal{D} under the clauses in Π ; in this case we write $\Pi, \mathcal{D} \models G(\mathbf{a})$. The problem of checking whether \mathbf{a} is an answer to $(\Pi, G(\mathbf{x}))$ over \mathcal{D} is called the *query evaluation problem*. The *depth* of $(\Pi, G(\mathbf{x}))$ is the length, $d(\Pi, G)$, of the longest directed path in the dependence graph for Π starting from G .

An NDL query $(II, G(\mathbf{x}))$ is an *NDL-rewriting of an OMQ* $Q(\mathbf{x}) = (\mathcal{O}, \mathbf{q}(\mathbf{x}))$ in case $\mathcal{O}, \mathcal{D} \models \mathbf{q}(\mathbf{a})$ iff $II, \mathcal{D} \models G(\mathbf{a})$, for any \mathcal{D} and any $\mathbf{a} \in \text{ind}(\mathcal{D})^{|\mathbf{x}|}$. Every OMQ is known to have an NDL-rewriting [2, 6].

Tree decomposition. A *tree decomposition* of a CQ \mathbf{q} with variables $\text{var}(\mathbf{q})$ is a pair (T, λ) of an (undirected) tree $T = (V, E)$ and $\lambda: V \rightarrow 2^{\text{var}(\mathbf{q})}$ such that

- for any atom $P(\mathbf{z}) \in \mathbf{q}$, there exists $v \in V$ with $\mathbf{z} \subseteq \lambda(v)$;
- for any variable z in \mathbf{q} , the set of vertices $\{v \in V \mid z \in \lambda(v)\}$ is connected in T .

We call $\lambda(v)$ the *bag for v* . The *width* of (T, λ) is $\max_{v \in V} |\lambda(v)| - 1$. The *treewidth of \mathbf{q}* is the minimum width over all tree decompositions of \mathbf{q} . It is known [9] that, for CQs of bounded arity n , the notions of bounded treewidth and *bounded hypertree width* [10] are interchangeable. Indeed, in this case, every hypertree decomposition of width t induces a tree decomposition of width $t \cdot n$. A CQ \mathbf{q} is called *acyclic* if it has a *join tree* whose nodes are the atoms of \mathbf{q} and, whenever atoms γ_1 and γ_2 share a variable, this variable occurs in all atoms along the (unique) path in the tree linking γ_1 and γ_2 . It is known [10] that a CQ \mathbf{q} is acyclic iff \mathbf{q} is of hypertree width 1.

3 NL and LogCFL Fragments of NDL

In this section we present two classes of NDL queries that enjoy NL- and LOGCFL-complete evaluation. First, observe that if the number of variables in each clause of an NDL query is bounded, then the size of its grounding (obtained by replacing variables by all possible combinations of constants) is polynomial. So, evaluation of such NDL queries is tractable. However, if we bound the number of variables in clauses of NDL rewritings of OMQs, then we will also effectively impose a bound on the number of answer variables in their CQs. To avoid this limitation, we treat answer variables of the CQs (and the predicate positions they occur in) differently from all other variables in the NDL-rewritings. Intuitively, answer variables get their values fixed by a candidate certain answer and thus do not cause an exponential blowup of groundings.

Formally, an NDL query $(II, G(\mathbf{x}))$ is called *ordered* if each of its IDB predicates Q has a fixed list of variables $\mathbf{x}_Q \subseteq \mathbf{x}$, the *parameters of Q* , such that

- the parameters of G are \mathbf{x} and, in every clause, the parameters of the head include all the parameters of the predicates in the body;
- the parameters \mathbf{x}_Q of each Q occupy the last $|\mathbf{x}_Q|$ positions in every occurrence of Q in II ; they can, however, occur in other positions too.

The *width* $w(II, G)$ of an ordered $(II, G(\mathbf{x}))$ is the maximum number of non-parameter variables in a clause of II . Observe that Boolean NDL queries are trivially ordered (their IDB predicates have no parameters), and the width of such queries is simply the maximum number of variables in a clause of II . As all the NDL-rewritings we construct are ordered, with their parameters being the answer variables, in the sequel we will consider only ordered NDL queries. We say that a class of NDL queries is of *bounded width* if there is $w > 0$ such that $w(II, G) \leq w$, for all $(II, G(\mathbf{x}))$ in the class. As we observed above, evaluation of NDL queries of bounded width is P-complete.

Our first subclass of NDL queries is based on linear rules. An NDL program is *linear* [1] if the body of its every clause contains at most one IDB predicate.

Theorem 1. *Evaluation of linear NDL queries of bounded width is NL-complete for combined complexity.*

Our second subclass was inspired by semi-unbounded fan-in circuits. Recall that the class LOGCFL of problems reducible in logarithmic space to context-free languages can equivalently be defined in terms of L-uniform families of semi-unbounded fan-in circuits (where OR-gates have arbitrarily many inputs, and AND-gates two inputs) of polynomial size and logarithmic depth. Alternatively, LOGCFL can be defined using *nondeterministic auxiliary pushdown automata* (NAuxPDAs) [7], which are nondeterministic Turing machines with an additional work tape constrained to operate as a pushdown store. Sudborough [17] proved that LOGCFL coincides with the class of problems that are solved by NAuxPDAs in logarithmic space and polynomial time (the space on the pushdown tape is not subject to the logarithmic bound). Moreover, there is an algorithm that, given a semi-unbounded fan-in circuit C and an input, computes the output using an NAuxPDA in logarithmic space in the size of C and exponential time in the depth of C [19, pp. 392–397]. Using these results, it can be shown that any $(II, G(x))$ with *at most two atoms* in the body of any clause can be evaluated on a data instance \mathcal{D} by an NAuxPDA in space $\log |II| + w(II, G) \cdot \log |\mathcal{D}|$ and time $2^{O(d(II, G))}$ (thus, in LOGCFL provided the query width is bounded and its depth is logarithmic).

In the rewritings we propose in Sections 5 and 7, however, the number of atoms in the clauses is not bounded. We require the following to generalise the idea. A function ν from the predicate names in II to non-negative integers \mathbb{N} is called a *weight function* for an NDL query $(II, G(x))$ if, for any clause $Q(z) \leftarrow P_1(z_1) \wedge \dots \wedge P_k(z_k)$ in II , we have

$$\nu(Q) > 0 \quad \text{and} \quad \nu(Q) \geq \nu(P_1) + \dots + \nu(P_k),$$

Note that $\nu(P)$ can be 0 for an EDB predicate P . To illustrate, we consider NDL queries with the following dependency graphs:



The one on the left has a weight function bounded by the number of predicates (i.e., linear in the size of the query); intuitively, this function corresponds to the number of directed paths from a vertex to the leaves. In contrast, any NDL query with the dependency graph on the right can only have a weight function whose values (numbers of paths) are exponential. Linear NDL queries have weight functions bounded by 1.

Let e_{II} be the maximum number of EDB predicates in a clause of II . The *skinny depth* $\text{sd}(II, G)$ of $(II, G(x))$ is the minimum value of

$$2d(II, G) + \log \nu(G) + \log e_{II}$$

over possible weight functions ν . One can show, using Huffman coding, that any NDL query $(II, G(x))$ can be transformed into an equivalent skinny NDL query $(II', G(x))$ of depth not exceeding $\text{sd}(II, G)$ and such that $|II'| = O(|II|^2)$ and $w(II', G) \leq w(II, G)$. We say that a class of NDL queries has *logarithmic skinny depth* if there is $c > 0$ such that $\text{sd}(II, G) \leq c \log |II|$, for all $(II, G(x))$ in the class. We now obtain:

Theorem 2. *Evaluation of NDL queries of logarithmic skinny depth and bounded width is LOGCFL-complete for combined complexity.*

3.1 NDL Rewritings over (Complete) Data

We say that a data instance \mathcal{D} is *complete for an ontology* \mathcal{O} if $\mathcal{O}, \mathcal{D} \models P(\mathbf{a})$ implies $P(\mathbf{a}) \in \mathcal{D}$, for any ground atom $P(\mathbf{a})$, where P in Σ and $\mathbf{a} \subseteq \text{ind}(\mathcal{D})$. An NDL query $(\Pi, G(\mathbf{x}))$ is an *NDL-rewriting of an OMQ* $\mathbf{Q}(\mathbf{x}) = (\mathcal{O}, \mathbf{q}(\mathbf{x}))$ over complete data in case $\mathcal{O}, \mathcal{D} \models \mathbf{q}(\mathbf{a})$ iff $\Pi, \mathcal{D} \models G(\mathbf{a})$, for any \mathcal{D} complete for \mathcal{O} and any $\mathbf{a} \subseteq \text{ind}(\mathcal{D})$.

Given an NDL-rewriting $(\Pi, G(\mathbf{x}))$ of $\mathbf{Q}(\mathbf{x})$ over complete data, we denote by Π^* the result of replacing each EDB predicate P in Π with a fresh IDB predicate P^* of the same arity and adding the clauses $P^*(z) \leftarrow \gamma$ for every atom γ with a predicate symbol from \mathcal{O} such that $\mathcal{O} \models \gamma \rightarrow P(z)$, where z is a tuple of variables (with possible repetitions). Clearly, $(\Pi^*, G(\mathbf{x}))$ is an NDL-rewriting of $\mathbf{Q}(\mathbf{x})$ over arbitrary data instances and $|\Pi^*| \leq |\Pi| + \text{ar}(\Sigma)^{\text{ar}(\Sigma)} \cdot |\mathcal{O}|^2$.

We say that a class of OMQs is *skinny-reducible* if there are $c > 0$ and $w > 0$ and an L^{LOGCFL} -transducer that, given any OMQ $\mathbf{Q}(\mathbf{x})$ in the class, computes its NDL-rewriting $(\Pi, G(\mathbf{x}))$ over complete data with $\text{sd}(\Pi, G) \leq c \log |\Pi|$ and $w(\Pi, G) \leq w$. Theorem 2 and the transformation $*$ give the following:

Corollary 1. *Answering OMQs is in LOGCFL for combined complexity for any skinny-reducible class.*

The transformation $*$, however, does not preserve linearity because it replaces occurrences of EDB predicates P by IDB predicates P^* . A more involved ‘linear’ construction is given in the proof of the following, where a possible increase of the width is due to the ‘replacement’ of atoms $P(z)$ by atoms γ whenever $\mathcal{O} \models \gamma \rightarrow P(z)$:

Lemma 1. *Fix any $w > 0$. There is an L^{NL} -transducer that, for a linear NDL-rewriting $(\Pi, G(\mathbf{x}))$ of an OMQ $\mathbf{Q}(\mathbf{x})$ over complete data with $w(\Pi, G) \leq w$, computes its linear NDL-rewriting $(\Pi', G(\mathbf{x}))$ over arbitrary data with $w(\Pi', G) \leq w + \text{ar}(\Sigma)$.*

4 Conditional Rewritings

Let $\mathbf{Q}(\mathbf{x}) = (\mathcal{O}, \mathbf{q}(\mathbf{x}))$ be an OMQ with an ontology of finite depth. Intuitively, we recursively split $\mathbf{q}(\mathbf{x})$ into subqueries \mathbf{q}_D based on subtrees D of a tree decomposition of \mathbf{q} and combine the rewritings of \mathbf{q}_D into a rewriting of \mathbf{q} . To guarantee ‘compatibility’ of the rewritings of the subqueries, we take account of the *types* of points on the boundaries of the \mathbf{q}_D . So, for each D and each type \mathbf{w} , we take a fresh IDB predicate $G_D^{\mathbf{w}}$ to represent the *conditional rewriting* of \mathbf{q}_D provided that its boundary satisfies the type. We now give formal definitions.

A *type* is a partial map \mathbf{s} from the variables of \mathbf{q} to $\text{term}_{\mathcal{O}} \cup \{\varepsilon\}$; its domain is denoted by $\text{dom}(\mathbf{s})$. The unique partial type with $\text{dom}(\varepsilon) = \emptyset$ is denoted by ε . We use types to represent how variables are mapped into the canonical model: $\mathbf{s}(z) = \varepsilon$ means that z is mapped to an individual constant and $\mathbf{s}(z) = f(\mathbf{a})$, for a Skolem term $f(\mathbf{a})$, means that z is mapped to an element of the form $f(\mathbf{c})$, for some $\mathbf{c} \subseteq \text{ind}(\mathcal{D})$. Given a type \mathbf{s} and a tuple $\mathbf{z} = (z_1, \dots, z_n) \subseteq \text{dom}(\mathbf{s})$, we denote the tuple $(\mathbf{s}(z_1), \dots, \mathbf{s}(z_n))$ by $\mathbf{s}(\mathbf{z})$. A type \mathbf{s} is *compatible* with a bag $\lambda(t)$ if $\mathbf{s}(x) = \varepsilon$, for all $x \in \mathbf{x} \cap \text{dom}(\mathbf{s})$, and, for every $S(\mathbf{z}) \in \mathbf{q}$ with $\mathbf{z} \subseteq \lambda(t) \cap \text{dom}(\mathbf{s})$, one of the following applies:

- (d) $\mathbf{s}(\mathbf{z}) \subseteq \{\varepsilon\}$;

- (b) there is $P(\mathbf{a})$ such that $\mathbf{s}(z) \subseteq \text{term}_{\mathcal{O}}(P(\mathbf{a})) \cup \{\varepsilon\}$ but neither $\mathbf{s}(z) \subseteq \{\varepsilon\}$ nor $\mathbf{s}(z) \subseteq \text{term}_{\mathcal{O}}(P(\mathbf{a}))$;
- (i) there is $P(\mathbf{a})$ such that $\mathbf{s}(z) \subseteq \text{term}_{\mathcal{O}}(P(\mathbf{a}))$ and $S(\mathbf{s}(z)) \in \mathfrak{C}_{\mathcal{O},\{P(\mathbf{a})\}}$.

Given a type \mathbf{s} , we take a tuple of variables $\text{var}(\mathbf{s})$ that contains, for $z \in \text{dom}(\mathbf{s}) \setminus \mathbf{x}$,

variable z , if $\mathbf{s}(z) = \varepsilon$, and variables $\tilde{\mathbf{a}}$, if $\mathbf{s}(z) \in \text{term}_{\mathcal{O}}(P(\mathbf{a}))$.

Denote the answer variables that occur in $\text{dom}(\mathbf{s})$ by $\mathbf{x}_{\mathbf{s}}$. Our rewritings use conjunctions $\text{At}^{\mathbf{s}}(\text{var}(\mathbf{s}), \mathbf{x}_{\mathbf{s}})$ of the following formulas, for all $S(z) \in \mathbf{q}$ with $z \subseteq \text{dom}(\mathbf{s})$:

(d') $S(z)$ if $\mathbf{s}(z) \subseteq \{\varepsilon\}$;

(b') the disjunction
$$\bigvee_{g: \mathbf{z}' \rightarrow \mathbf{a}} \left[P(\tilde{\mathbf{a}}) \wedge \bigwedge_{z \in \mathbf{z}' \text{ and } g(z)=a} (z = \tilde{a}) \right]$$

over *grounding functions* $g: \mathbf{z}' \rightarrow \mathbf{a}$ such that $\mathbf{z}' = \{z \in \mathbf{z} \mid \mathbf{s}(z) = \varepsilon\} \neq \emptyset$, $\mathbf{z}'' = \{z \in \mathbf{z} \mid \mathbf{s}(z) \in \text{term}_{\mathcal{O}}(P(\mathbf{a}))\} \neq \emptyset$ and $\mathfrak{C}_{\mathcal{O},\{P(\mathbf{a})\}}$ contains the result of replacing \mathbf{z}' and \mathbf{z}'' in $S(z)$ by $g(\mathbf{z}')$ and $\mathbf{s}(\mathbf{z}'')$, respectively;

(i') $P(\tilde{\mathbf{a}})$ if $\mathbf{s}(z) \subseteq \text{term}_{\mathcal{O}}(P(\mathbf{a}))$.

Strictly speaking, the resulting rewritings will not be NDL programs because of disjunctions in (b'), but we can get rid of them using an extra predicate and (if needed) the construction from the proof of Lemma 1 keeping the size and the execution time polynomial.

5 LOGCFL Rewritings for $\text{OMQ}(d, t, \infty)$

We now construct skinny-reducible NDL rewritings for the CQs of bounded treewidth.

Theorem 3. *For any $d \geq 0$ and $t \geq 1$, the class $\text{OMQ}(d, t, \infty)$ is skinny-reducible.*

Fix a connected CQ $\mathbf{q}(\mathbf{x})$ and a tree decomposition (T, λ) of its Gaifman graph. Let D be a subtree of T . The *size* of D is the number of nodes in it. A node t of D is called *boundary* if T has an edge $\{t, t'\}$ with $t' \notin D$. We denote by ∂D the union of all $\lambda(t) \cap \lambda(t')$ for boundary nodes t of D and its neighbours t' in T outside D . The *degree* $\text{deg}(D)$ of D is the number of its boundary nodes (so, the only subtree of T of degree 0 is T itself). We say that a node t *splits* D into subtrees D_1, \dots, D_k if the D_i partition D without t : each node of D except t belongs to exactly one D_i .

Lemma 2. *Let D be a subtree of T of size $n > 1$.*

If $\text{deg}(D) = 2$, then there is a node t splitting D into subtrees of size $\leq n/2$ and degree ≤ 2 and, possibly, one subtree of size $< n - 1$ and degree 1.

If $\text{deg}(D) \leq 1$, then there is t splitting D into subtrees of size $\leq n/2$ and degree ≤ 2 .

We define recursively a set \mathfrak{R} of subtrees of T , a binary ‘predecessor’ relation \prec on \mathfrak{R} , and a function β on \mathfrak{R} indicating the bag of the splitting node. We begin by adding T to \mathfrak{R} . Take any $D \in \mathfrak{R}$ that has not been split yet. If D is of size 1, then $\beta(D) = \lambda(t)$ for the only node t of D . Otherwise, by Lemma 2, we find a node t in D

that splits it into D_1, \dots, D_k . We set $\beta(D) = \lambda(t)$ and, for $1 \leq i \leq k$, add D_i to \mathfrak{R} and set $D_i \prec D$; then, we apply the procedure to each of D_1, \dots, D_k . For each $D \in \mathfrak{R}$, we recursively define a set of atoms

$$\mathbf{q}_D = \{S(\mathbf{z}) \in \mathbf{q} \mid \mathbf{z} \subseteq \beta(D)\} \cup \bigcup_{D' \prec D} \mathbf{q}_{D'}.$$

Let \mathbf{x}_D be the set of variables from \mathbf{x} that occur in \mathbf{q}_D . By the definition of tree decomposition, $\mathbf{q}_T = \mathbf{q}$ and $\mathbf{x}_T = \mathbf{x}$.

We now define an NDL-rewriting of $\mathbf{Q}(\mathbf{x}) = (\mathcal{O}, \mathbf{q}(\mathbf{x}))$. Fix $D \in \mathfrak{R}$ and a type \mathbf{w} with $\text{dom}(\mathbf{w}) = \partial D$. Let $G_D^{\mathbf{w}}(\text{var}(\mathbf{w}), \mathbf{x}_D)$ be a fresh IDB predicate with parameters \mathbf{x}_D . As we described above, a node is selected in D to split it into smaller trees (provided that it contains more than one node). We extend the type \mathbf{w} to cover the variables $\beta(D)$ of the selected bag: more precisely, we consider types \mathbf{s} with $\text{dom}(\mathbf{s}) = \beta(D)$ such that they are compatible with bag $\beta(D)$ and agree with \mathbf{w} on their common domain. Observe that, if D' is a subtree resulting from splitting D , then the domain of the extended type, $\mathbf{s} \cup \mathbf{w}$, includes $\partial D'$, and thus $\partial D'$ coincides with the domain of the restriction of $\mathbf{s} \cup \mathbf{w}$ to $\partial D'$, denoted $(\mathbf{s} \cup \mathbf{w}) \upharpoonright_{\partial D'}$. Now, for each type \mathbf{s} with $\text{dom}(\mathbf{s}) = \beta(D)$ such that \mathbf{s} is compatible with bag $\beta(D)$ and agrees with \mathbf{w} on their common domain, the NDL program Π_Q^{LOG} contains

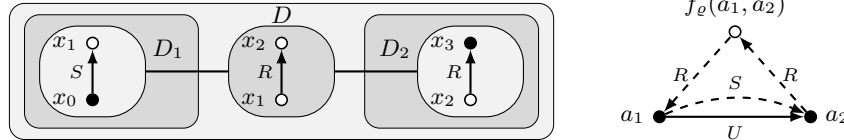
$$G_D^{\mathbf{w}}(\text{var}(\mathbf{w}), \mathbf{x}_D) \leftarrow \text{At}^{\mathbf{s}}(\text{var}(\mathbf{s}), \mathbf{x}_s) \wedge \bigwedge_{D' \prec D} G_{D'}^{(\mathbf{s} \cup \mathbf{w}) \upharpoonright_{\partial D'}}(\text{var}((\mathbf{s} \cup \mathbf{w}) \upharpoonright_{\partial D'}), \mathbf{x}_{D'}).$$

By induction on \prec , one can now show that $(\Pi_Q^{\text{LOG}}, G_T^{\varepsilon})$ is a rewriting of $\mathbf{Q}(\mathbf{x})$.

Example 1. Let $\mathbf{q}(x_0, x_3) = \exists x_1 x_2 (S(x_0, x_1) \wedge R(x_1, x_2) \wedge R(x_2, x_3))$ and \mathcal{O} consist of the following linear rules:

$$\begin{aligned} \varrho: \quad U(x, y) &\rightarrow \exists v T(x, v, y), & T(x, v, y) &\rightarrow R(v, x), \\ T(x, v, y) &\rightarrow R(y, v), & T(x, v, y) &\rightarrow S(x, y). \end{aligned}$$

The subtree structure of the tree decomposition of $\mathbf{q}(x_0, x_3)$ and the canonical model are as follows:



The goal predicate for the rewriting of $\mathbf{q}(x_0, x_3)$ is $G_D^{\varepsilon}(x_0, x_3)$ with parameters x_0 and x_3 . For the type \mathbf{s} for the middle bag sending x_1 to ε and x_2 to $f_e(a_1, a_2)$, we have

$$G_D^{\varepsilon}(x_0, x_3) \leftarrow G_{D_1}^{x_1 \mapsto \varepsilon}(x_1, x_0) \wedge U(\tilde{a}_1, \tilde{a}_2) \wedge (x_1 = \tilde{a}_2) \wedge G_{D_2}^{x_2 \mapsto f_e(a_1, a_2)}(\tilde{a}_1, \tilde{a}_2, x_3),$$

where \tilde{a}_1 and \tilde{a}_2 are the twin variables in $\text{var}(\mathbf{s})$. Note that the type for $G_{D_2}^{x_2 \mapsto f_e(a_1, a_2)}$ has no non-twin variables, and we have the following rule for this predicate

$$G_{D_2}^{x_2 \mapsto f_e(a_1, a_2)}(\tilde{a}_1, \tilde{a}_2, x_3) \leftarrow U(\tilde{a}_1, \tilde{a}_2) \wedge (x_3 = \tilde{a}_1).$$

Lemma 3. For any \mathcal{D} complete for \mathcal{O} , any predicate G_D^w and any $\mathbf{b} \in \text{ind}(\mathcal{D})^{|\text{var}(w)|+|\mathbf{x}_D|}$, we have $\Pi_Q^{\text{LOG}}, \mathcal{D} \models G_D^w(\mathbf{b})$ iff there is a homomorphism $h: \mathbf{q}_D \rightarrow \mathcal{C}_{\mathcal{O}, \mathcal{D}}$ such that

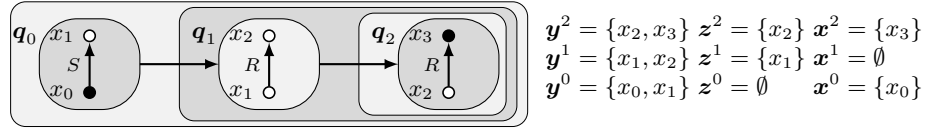
$$h(z) = \begin{cases} \mathbf{b}(z), & \text{for all } z \in \mathbf{x}_D \text{ and all } z \in \partial D \text{ with } w(z) = \varepsilon, \\ w(z)[\mathbf{a}/\mathbf{b}(\tilde{\mathbf{a}})], & \text{for all } z \in \partial D \text{ with } w(z) \in \text{term}_{\mathcal{O}}(P(\mathbf{a})). \end{cases}$$

6 NL Rewritings for OMQ(d, t, ℓ)

For OMQs based upon bounded leaf queries and bounded depth ontologies, we establish the following theorem:

Theorem 4. Let $d \geq 0$, $t \geq 1$ and $\ell \geq 2$ be fixed. There is an L^{NL} -transducer that, given any OMQ in $\text{OMQ}(d, t, \ell)$, constructs its polynomial-size linear NDL-rewriting of width $\leq \ell(t+1)$.

Let \mathcal{O} be an ontology of finite depth d and $q(\mathbf{x})$ a CQ with a tree decomposition (T, λ) of width $\leq t$ having $\leq \ell$ leaves. Fix one of the nodes of T as root, and let M be the maximum distance to a leaf from the root. For $0 \leq n \leq M$, by an n -slice we mean the set of all nodes of T located at distance n from the root. Denote by \mathbf{y}^n the union of all bags $\lambda(t)$ for a node t in the n -slice. For $1 \leq n \leq M$, let \mathbf{z}^n be the union of all $\lambda(t) \cap \lambda(t')$ for a node t in the n -slice and its predecessor t' in T (which is in $(n-1)$ -slice), and let $\mathbf{z}^0 = \emptyset$. By definition, $\mathbf{z}^{n+1} \subseteq \mathbf{y}^{n+1} \cap \mathbf{y}^n$ and, clearly, $|\mathbf{z}^n| \leq |\mathbf{y}^n| \leq \ell(t+1)$. Denote by $\mathbf{q}_n(\mathbf{z}_{\exists}^n, \mathbf{x}^{\geq n})$ the query consisting of all atoms $S(z)$ of q with $z \subseteq \bigcup_{k \geq n} \mathbf{y}^k$, where $\mathbf{z}_{\exists}^n = \mathbf{z}^n \setminus \mathbf{x}$ and $\mathbf{x}^{\geq n} = \mathbf{x} \cap \bigcup_{k \geq n} \mathbf{y}^k$. These queries and sets of variables for the CQ from Example 1 are shown below:



A type for \mathbf{z}^n is a total map w from \mathbf{z}^n to $\text{term}_{\mathcal{O}} \cup \{\varepsilon\}$. Likewise, a type for \mathbf{y}^n is a total map s from \mathbf{y}^n to $\text{term}_{\mathcal{O}} \cup \{\varepsilon\}$. We say s compatible with \mathbf{y}^n if it is compatible with every bag $\lambda(t)$ in the n -slice.

Consider the NDL program Π_Q^{LIN} defined as follows. For every $0 \leq n < M$ and every type w for \mathbf{z}^n , we introduce a new IDB predicate $G_n^w(\text{var}(w), \mathbf{x}^{\geq n})$ with parameters $\mathbf{x}^{\geq n}$. For each type s for \mathbf{y}^n such that s is compatible with \mathbf{y}^n and agrees with w on \mathbf{z}^n , the program Π_Q^{LIN} contains the clause

$$G_n^w(\text{var}(w), \mathbf{x}^{\geq n}) \leftarrow \text{At}^s(\text{var}(s), \mathbf{x}_s) \wedge G_{n+1}^{s \upharpoonright \mathbf{z}^{n+1}}(\text{var}(s \upharpoonright_{\mathbf{z}^{n+1}}), \mathbf{x}^{\geq n+1}).$$

For every type w for \mathbf{z}^M and every type s for \mathbf{y}^M such that s is compatible with \mathbf{y}^M and agrees with w on \mathbf{z}^M , we include the clause

$$G_M^w(\text{var}(w), \mathbf{x}^{\geq M}) \leftarrow \text{At}^s(\text{var}(s), \mathbf{x}_s).$$

Finally, we use G_0^ε with parameters \mathbf{x} as the goal predicate (note that $\mathbf{z}^0 = \emptyset$, and so the domain of any type for \mathbf{z}^0 is empty).

Lemma 4. For any \mathcal{D} complete for \mathcal{O} , any predicate G_n^w , any $\mathbf{b} \in \text{ind}(\mathcal{D})^{|\text{var}(w)|+|\mathbf{x}^{\geq n}|}$, we have $\Pi_{\mathcal{Q}}^{\text{LIN}}, \mathcal{D} \models G_n^w(\mathbf{b})$ iff there is a homomorphism $h: \mathbf{q}_n \rightarrow \mathcal{C}_{\mathcal{O}, \mathcal{D}}$ such that

$$h(z) = \begin{cases} \mathbf{b}(z), & \text{for all } z \in \mathbf{x}^{\geq n} \text{ and all } z \in \mathbf{z}_{\exists}^n \text{ with } \mathbf{w}(z) = \varepsilon, \\ \mathbf{w}(z)[\mathbf{a}/\mathbf{b}(\tilde{\mathbf{a}})], & \text{for all } z \in \mathbf{z}_{\exists}^n \text{ with } \mathbf{w}(z) \in \text{term}_{\mathcal{O}}(P(\mathbf{a})). \end{cases}$$

It should be clear that $\Pi_{\mathcal{Q}}^{\text{LIN}}$ is a linear NDL program of width $\leq \ell(t+1)$ and containing $\leq |\mathbf{q}| \cdot |\text{term}_{\mathcal{O}}|^{\ell(t+1)}$ predicates. Moreover, it takes only logarithmic space to store a type \mathbf{w} , which allows us to show that $\Pi_{\mathcal{Q}}^{\text{LIN}}$ can be computed by an L^{NL} -transducer. We apply Lemma 1 to obtain an NDL-rewriting for arbitrary data instances, and then use Theorem 1 to conclude that the resulting program can be evaluated in NL.

7 LOGCFL Rewritings for $\text{OMQ}(\infty, 1, \ell)$

Unlike the previous two classes, answering OMQs from the class $\text{OMQ}(\infty, 1, \ell)$ can be harder—LOGCFL-complete—than evaluating their CQs, which can be done in NL.

Theorem 5. For any fixed $\ell \geq 2$, the class $\text{OMQ}(\infty, 1, \ell)$ is skinny-reducible.

For OMQs with ontologies of unbounded depth and acyclic CQs whose join trees have a bounded number of leaves, our rewriting uses the notion of Skolem witness that generalises tree witnesses [14].

Let $\mathbf{Q}(\mathbf{x}) = (\mathcal{O}, \mathbf{q}(\mathbf{x}))$ be an OMQ, let $\mathfrak{s} = (\mathfrak{s}_r^1, \dots, \mathfrak{s}_r^n, \mathfrak{s}_i)$ be a tuple of disjoint sets of variables in $\mathbf{q}(\mathbf{x})$ such that $\mathfrak{s}_i \neq \emptyset$ and $\mathfrak{s}_i \cap \mathbf{x} = \emptyset$, and let $\mathfrak{s}_r = \mathfrak{s}_r^1 \cup \dots \cup \mathfrak{s}_r^n$,

$$\mathbf{q}_{\mathfrak{s}} = \{ S(z) \in \mathbf{q} \mid z \subseteq \mathfrak{s}_r \cup \mathfrak{s}_i \text{ and } z \not\subseteq \mathfrak{s}_r \}.$$

If $\mathbf{q}_{\mathfrak{s}}$ is a minimal subset of \mathbf{q} containing every atom of \mathbf{q} with at least one variable from \mathfrak{s}_i and such that there is a homomorphism $h: \mathbf{q}_{\mathfrak{s}} \rightarrow \mathcal{C}_{\mathcal{O}, \{P(\mathbf{a})\}}$ with $\mathbf{a} = (a_1, \dots, a_n)$ and $h^{-1}(a_j) = \mathfrak{s}_r^j$ for $1 \leq j \leq n$, then we call \mathfrak{s} a *Skolem witness for $\mathbf{Q}(\mathbf{x})$ generated by $P(\mathbf{a})$* . Intuitively, \mathfrak{s} identifies a minimal subset of \mathbf{q} that can be mapped to the *Skolem part* of the canonical model $\mathcal{C}_{\mathcal{O}, \{P(\mathbf{a})\}}$ consisting of Skolem terms: the variables in \mathfrak{s}_r are mapped to constants from \mathbf{a} and the variables in \mathfrak{s}_i to Skolem terms in $\text{term}_{\mathcal{O}}(P(\mathbf{a}))$.

The logarithmic-depth NDL-rewriting for $\text{OMQ}(\infty, 1, \ell)$ is based on the following:

Lemma 5. Every tree T of size n has a node splitting it into subtrees of size $\leq \lceil n/2 \rceil$.

Let $\mathbf{Q}(\mathbf{x}_0) = (\mathcal{O}, \mathbf{q}_0(\mathbf{x}_0))$ be an OMQ with an acyclic CQ having a join tree T_0 . We repeatedly apply Lemma 5 to decompose the CQ into smaller and smaller subqueries. Formally, for an acyclic CQ \mathbf{q} , we denote by $\gamma_{\mathbf{q}}$ a vertex in the join tree T for \mathbf{q} that satisfies the condition of Lemma 5. Let Ω be the smallest set containing $\mathbf{q}_0(\mathbf{x}_0)$ and the following CQs, for every $\mathbf{q}(\mathbf{x}) \in \Omega$ with at least one existentially quantified variable:

- (1) the CQs $\mathbf{q}_i(\mathbf{x}_i)$ corresponding to the connected components T_i with root $\gamma_{\mathbf{q}_i}$ adjacent to $\gamma_{\mathbf{q}}$ of the result of removing $\gamma_{\mathbf{q}}$ from T , where \mathbf{x}_i consists of the restriction of \mathbf{x} to the variables in \mathbf{q}_i together with the common variables of $\gamma_{\mathbf{q}_i}$ and $\gamma_{\mathbf{q}}$;

- (2) for each Skolem witness \mathfrak{s} for $(\mathcal{O}, \mathbf{q}(\mathbf{x}))$ with $\mathfrak{s}_r \neq \emptyset$ and $\gamma_{\mathbf{q}} \in \mathbf{q}_{\mathfrak{s}}$, the CQs $\mathbf{q}_1^{\mathfrak{s}}(\mathbf{x}_1^{\mathfrak{s}}), \dots, \mathbf{q}_k^{\mathfrak{s}}(\mathbf{x}_k^{\mathfrak{s}})$ that correspond to the connected components $T_i^{\mathfrak{s}}$ of the results of removing $\mathbf{q}_{\mathfrak{s}}$ from T (note that $\mathbf{q}_{\mathfrak{s}}$ is connected in T), where each $\mathbf{x}_i^{\mathfrak{s}}$ is the set of variables in $\mathbf{x} \cup \mathfrak{s}_r$ that occur in $\mathbf{q}_i^{\mathfrak{s}}$.

The NDL program Π_Q^{SW} uses IDB predicates $G_{\mathbf{q}}(\mathbf{x})$, for $\mathbf{q}(\mathbf{x}) \in \mathfrak{Q}$, whose parameters are the variables in \mathbf{x}_0 that occur in $\mathbf{q}(\mathbf{x})$. For each $\mathbf{q}(\mathbf{x}) \in \mathfrak{Q}$ that has no existentially quantified variables, we include the clause $G_{\mathbf{q}}(\mathbf{x}) \leftarrow \mathbf{q}(\mathbf{x})$. For any $\mathbf{q}(\mathbf{x}) \in \mathfrak{Q}$ with existential variables, we include

$$G_{\mathbf{q}}(\mathbf{x}) \leftarrow \gamma_{\mathbf{q}} \wedge \bigwedge_{1 \leq i \leq n} G_{\mathbf{q}_i}(\mathbf{x}_i),$$

where $\mathbf{q}_1(\mathbf{x}_1), \dots, \mathbf{q}_n(\mathbf{x}_n)$ are the subqueries obtained by splitting \mathbf{q} by $\gamma_{\mathbf{q}}$ in (1), and, for any Skolem witness \mathfrak{s} of $(\mathcal{O}, \mathbf{q}(\mathbf{x}))$ with $\mathfrak{s}_r \neq \emptyset$ and $\gamma_{\mathbf{q}} \in \mathbf{q}_{\mathfrak{s}}$ and any $P(\mathbf{a})$ generating \mathfrak{s} , the clause

$$G_{\mathbf{q}}(\mathbf{x}) \leftarrow P(\tilde{\mathbf{a}}) \wedge \bigwedge_{z \in \mathfrak{s}_r^i} (z = \tilde{a}_j) \wedge \bigwedge_{1 \leq i \leq k} G_{\mathbf{q}_i^{\mathfrak{s}}}(\mathbf{x}_i^{\mathfrak{s}}),$$

where $\mathbf{q}_1^{\mathfrak{s}}, \dots, \mathbf{q}_k^{\mathfrak{s}}$ are the connected components of \mathbf{q} without $\mathbf{q}_{\mathfrak{s}}$. Finally, if \mathbf{q}_0 is Boolean, then we include $G_{\mathbf{q}_0} \leftarrow P(\tilde{\mathbf{a}})$ for all atoms $P(\mathbf{a})$ such that $\mathcal{O}, \{P(\mathbf{a})\} \models \mathbf{q}_0$.

Lemma 6. *For any OMQ with an acyclic CQ, any data \mathcal{D} complete for \mathcal{O} , any query $\mathbf{q}(\mathbf{x}) \in \mathfrak{Q}$ and any $\mathbf{b} \in \text{ind}(\mathcal{D})^{|\mathbf{x}|}$, we have $\Pi_Q^{\text{SW}}, \mathcal{D} \models G_{\mathbf{q}}(\mathbf{b})$ iff there is a homomorphism $h: \mathbf{q} \rightarrow \mathcal{C}_{\mathcal{O}, \mathcal{D}}$ with $h(\mathbf{x}) = \mathbf{b}$.*

Now fix $\ell > 1$ and consider $\mathbf{Q}(\mathbf{x}) = (\mathcal{O}, \mathbf{q}_0(\mathbf{x}))$ from the class $\text{OMQ}(\infty, 1, \ell)$ (remember that we have fixed arity n). The size of the program Π_Q^{SW} is polynomially bounded in $|\mathbf{Q}|$ since \mathbf{q}_0 has polynomially-many subtrees of $T_{\mathbf{q}_0}$ and $O(|\mathbf{q}_0|^{\ell})$ Skolem witnesses (there are at most $O(|\mathbf{q}_0|^{\ell} \cdot |\Sigma| \cdot n^n)$ pairs of a Skolem witness \mathfrak{s} and its generating atom $P(\mathbf{a})$). It is readily seen that the function ν defined by $\nu(G_{\mathbf{q}}) = |\mathbf{q}|$, for each $\mathbf{q} \in \mathfrak{Q}$, is a weight function for $(\Pi_Q^{\text{SW}}, G_{\mathbf{q}_0}(\mathbf{x}))$ with $\nu(G_{\mathbf{q}_0}) \leq |\mathbf{Q}|$. Moreover, by Lemma 5, $d(\Pi_Q^{\text{SW}}, G_{\mathbf{q}_0}) \leq \log \nu(G_{\mathbf{q}_0}) + 1$; also, $w(\Pi_Q^{\text{SW}}, G_{\mathbf{q}_0}) \leq \ell + 1$. Finally, we note that, since the number of leaves is bounded, it is in NL to decide whether a vertex satisfies the conditions of Lemma 5, and in LOGCFL to decide whether $\mathcal{O}, \{P(\mathbf{a})\} \models \mathbf{q}(\mathbf{a})$, for bounded-leaf acyclic CQs $\mathbf{q}(\mathbf{x})$ (see the full version¹), or whether a (logspace) representation of a possible Skolem witness is indeed a Skolem witness. This allows us to show that $(\Pi_Q^{\text{SW}}, G_{\mathbf{q}_0}(\mathbf{x}))$ can be generated by an L^{LOGCFL} -transducer. By Corollary 1, the obtained NDL-rewritings can be evaluated in LOGCFL.

8 Conclusion

We presented NDL rewritings for three classes of OMQs with CQs of bounded (hyper)-tree width and ontologies given as linear TGDs. These NDL rewritings can be constructed and evaluated in LOGCFL, NL and LOGCFL, respectively (provided that the arity of predicates is bounded). Since the three upper bounds match the lower bounds inherited from the OWL 2 QL setting [4], the proposed rewritings are theoretically optimal.

¹ <http://www.dcs.bbk.ac.uk/~kikot/DL17-1-full.pdf>

References

1. Abiteboul, S., Hull, R., Vianu, V.: Foundations of Databases. Addison-Wesley (1995)
2. Baget, J.-F., Leclère, M., Mugnier, M.-L., Salvat, E.: On rules with existential variables: Walking the decidability line. *Artificial Intelligence* 175(9–10), 1620–1654 (2011)
3. Bienvenu, M., Kikot, S., Kontchakov, R., Podolskii, V.V., Ryzhikov, V., Zakharyashev, M.: The complexity of ontology-based data access with OWL 2 QL and bounded treewidth queries. In: Proc. of the 26th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS (2017)
4. Bienvenu, M., Kikot, S., Podolskii, V.V.: Tree-like queries in OWL 2 QL: succinctness and complexity results. In: Proc. of the 30th Annual ACM/IEEE Symposium on Logic in Computer Science, LICS 2015. pp. 317–328. IEEE Computer Society (2015)
5. Boguslavsky, I., Dikonov, V., Iomdin, L., Lazursky, A., Sizov, V., Timoshenko, S.: Semantic analysis and question answering: a system under development. In: Computational Linguistics and Intellectual Technologies. Papers from the Annual International Conference Dialogue. p. 21. No. 14 (2015)
6. Cali, A., Gottlob, G., Lukasiewicz, T.: A general datalog-based framework for tractable query answering over ontologies. *Journal of Web Semantics* 14, 57–83 (2012)
7. Cook, S.A.: Characterizations of pushdown machines in terms of time-bounded computers. *Journal of the ACM* 18(1), 4–18 (1971)
8. Cuenca Grau, B., Horrocks, I., Krötzsch, M., Kupke, C., Magka, D., Motik, B., Wang, Z.: Acyclicity notions for existential rules and their application to query answering in ontologies. *Journal of Artificial Intelligence Research* 47, 741–808 (2013)
9. Gottlob, G., Greco, G., Leone, N., Scarcello, F.: Hypertree decompositions: Questions and answers. In: Proc. of the 35th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS (2016)
10. Gottlob, G., Leone, N., Scarcello, F.: Hypertree decompositions and tractable queries. *Journal of Computer and System Sciences* 64(3), 579–627 (2002)
11. Gottlob, G., Manna, M., Pieris, A.: Polynomial Rewritings for Linear Existential Rules, pp. 2992–2998. In: Proc. of the 24th Int. Joint Conf. on Artificial Intelligence, IJCAI (2015)
12. Gottlob, G., Orsi, G., Pieris, A.: Query rewriting and optimization for ontological databases. *ACM Transactions on Database Systems* 39(3), 25:1–25:46 (2014)
13. Kikot, S., Kontchakov, R., Podolskii, V., Zakharyashev, M.: On the succinctness of query rewriting over shallow ontologies. In: Proc. of the Joint Meeting of the 23rd EACSL Annual Conf. on Computer Science Logic (CSL 2014) and the 29th Annual ACM/IEEE Symposium on Logic in Computer Science (LICS 2014). pp. 57:1–57:10. ACM (2014)
14. Kikot, S., Kontchakov, R., Zakharyashev, M.: Conjunctive query answering with OWL 2 QL. In: Proc. of the 13th Int. Conf. on Principles of Knowledge Representation and Reasoning (KR 2012). pp. 275–285. AAAI (2012)
15. König, M., Leclère, M., Mugnier, M.-L., Thomazo, M.: Sound, complete and minimal UCQ-rewriting for existential rules. *Semantic Web* 6(5), 451–475 (2015)
16. König, M., Leclère, M., Mugnier, M.-L.: Query rewriting for existential rules with compiled preorder. In: Proc. of the 24th Int. Joint Conf. on Artificial Intelligence, IJCAI (2015)
17. Sudborough, I.H.: On the tape complexity of deterministic context-free languages. *Journal of the ACM* 25(3), 405–414 (1978)
18. Thomazo, M.: Compact rewriting for existential rules. In: Proc. of the 23rd Int. Joint Conf. on Artificial Intelligence, IJCAI (2013)
19. Venkateswaran, H.: Properties that characterize LOGCFL. *Journal of Computer and System Sciences* 43(2), 380–404 (1991)