

# Towards a Data Complexity Classification of Ontology-Mediated Queries with Covering

O. Gerasimova<sup>1</sup>, S. Kikot<sup>2</sup>, and M. Zakharyashev<sup>3</sup>

<sup>1</sup> National Research University Higher School of Economics, Moscow, Russia

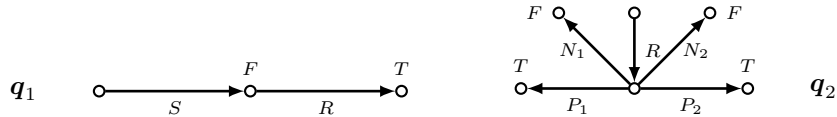
<sup>2</sup> University of Oxford, U.K.

<sup>3</sup> Birkbeck, University of London, U.K.

**Abstract.** We prove a number of new syntactic and semantic sufficient and necessary conditions for ontology-mediated queries (OMQs) with one covering axiom to be in the classes  $AC^0$  and NL for data complexity, and to be L-, NL- or P-hard. We also give two new examples of very simple CONP-complete OMQs.

## 1 Introduction

The problem we are concerned with in this paper originated from the observation that the NPD FactPages ontology<sup>4</sup>, used for testing ontology-based data access (OBDA) in industry [13, 12], contained covering axioms of the form  $A \sqsubseteq B_1 \sqcup \dots \sqcup B_n$ . It has been known since Schaerf's paper [18] that answering ontology-mediated queries (OMQs) with a covering axiom can be CONP-hard for data complexity, and so such OMQs are not suitable for OBDA systems in general. In fact, two interesting OMQs with covering can be extracted from [18]:  $Q_1 = (Cov_{\top}, q_1)$  and  $Q_2 = (Cov_{\top}, q_2)$ , where  $Cov_{\top} = \{\top \sqsubseteq F \sqcup T\}$  and the Boolean conjunctive queries (CQs)  $q_1$  and  $q_2$  are shown below.



Schaerf pointed out that answering these queries involves case analysis and showed that  $Q_2$  is CONP-hard for data complexity. It is readily seen that  $Q_1$  is NL-complete. The problem we started attacking in [9] was to find a complete syntactic or semantic classification of OMQs with covering according to their data complexity. The direct attack failed, and connections with other hard classification problems such as boundedness or linearisability of datalog programs (to be discussed in Section 3) indicate that a longer siege is needed. In this paper, we report on our recent findings.

We consider four ontologies with covering:  $Cov_{\top}$  as above,  $Cov_A = \{A \sqsubseteq F \sqcup T\}$ ,  $Cov_{\perp} = \{\top \sqsubseteq F \sqcup T, F \sqcap T \sqsubseteq \perp\}$ ,  $Cov_A^{\perp} = \{A \sqsubseteq F \sqcup T, F \sqcap T \sqsubseteq \perp\}$ , and we classify connected Boolean CQs  $q$  by the number of solitary occurrences of  $F$  and  $T$  in them, where a *solitary*  $F$  is any  $F(x) \in q$  with  $T(x) \notin q$ , and symmetrically for  $T$ . We prove a number of new syntactic and semantic sufficient and necessary conditions for

<sup>4</sup> <http://sws.ifi.uio.no/project/npd-v2/>

such OMQs with at most one solitary  $F$  to be in the complexity classes  $AC^0$  and NL, and to be L-, NL- or P-hard. We also give two new examples of very simple CONP-complete OMQs.

## 2 Preliminaries

In this paper, a *Boolean conjunctive query* (CQ) is any first-order (FO) sentence of the form  $q = \exists x \varphi(x)$ , where  $\varphi$  is a conjunction of unary or binary atoms  $P(y)$  with  $y \subseteq x$ . We often regard CQs as *sets* of their atoms, depict them as labelled digraphs, and assume that all of our CQs are *connected*. By *answering an ontology-mediated query* (OMQ)  $Q = (\mathcal{T}, q)$  with a TBox  $\mathcal{T}$  of the form defined above, we understand the problem of checking, given an ABox  $\mathcal{A}$  (a finite set of unary or binary ground atoms), whether  $q$  holds in every model of  $\mathcal{T} \cup \mathcal{A}$ , or  $\mathcal{T}, \mathcal{A} \models q$  in symbols. For every  $Q$ , this problem is clearly in CONP. It is in the complexity class  $AC^0$  if there is an FO-sentence  $q'$ , called an *FO-rewriting* of  $Q$ , such that  $\mathcal{T}, \mathcal{A} \models q$  iff  $\mathcal{A} \models q'$ , for any ABox  $\mathcal{A}$ .

A *datalog program*,  $\Pi$ , is a finite set of *rules*  $\forall x (\gamma_0 \leftarrow \gamma_1 \wedge \dots \wedge \gamma_m)$ , where each  $\gamma_i$  is an atom  $Q(y)$  with  $y \subseteq x$ . (As usual, we omit  $\forall x$ .) The atom  $\gamma_0$  is the *head* of the rule, and  $\gamma_1, \dots, \gamma_m$  its *body*. All the variables in the head must occur in the body. The predicates in the head of rules are *IDB predicates*, the rest *EDB predicates* [1].

A *datalog query* is a pair  $(\Pi, G)$ , where  $\Pi$  is a datalog program and  $G$  an 0-ary atom, the *goal*. The *answer* to  $(\Pi, G)$  over an ABox  $\mathcal{A}$  is ‘yes’ if  $G$  holds in the FO-structure obtained by closing  $\mathcal{A}$  under  $\Pi$ , in which case we write  $\Pi, \mathcal{A} \models G$ . A datalog query  $(\Pi, G)$  is a *datalog rewriting* of an OMQ  $Q = (\mathcal{T}, q)$  in case  $\mathcal{T}, \mathcal{A} \models q$  iff  $\Pi, \mathcal{A} \models G$ , for any ABox  $\mathcal{A}$ . The *answering problem* for  $(\Pi, G)$ —i.e., checking, given an ABox  $\mathcal{A}$ , whether  $\Pi, \mathcal{A} \models G$ —is clearly in P. Answering a datalog query with a *linear* program, whose rules have at most one IDB predicate in the body, can be done in NL. The NL upper bound also holds for datalog queries with a linear-stratified program defined as follows. A *stratified* program [1] is a sequence  $\Pi = (\Pi_0, \dots, \Pi_n)$  of datalog programs, called the *strata* of  $\Pi$ , such that each predicate in  $\Pi$  can occur in the head of a rule only in one stratum  $\Pi_i$  and can occur in the body of a rule only in strata  $\Pi_j$  with  $j \geq i$ . If, additionally, the body of each rule in  $\Pi$  contains at most one occurrence of a head predicate from the same stratum,  $\Pi$  is called *linear-stratified*. Every linear-stratified program can be converted to an equivalent linear datalog program [2].

## 3 Connections in High Places

We begin by putting our little problem into the context of more general investigations of (i) boundedness and linearisability of datalog programs and (ii) the data complexity of answering OMQs with expressive ontologies.

The decision problem whether a given datalog program is bounded has been a hot research topic in database theory since the late 1980s. Thus, it was shown that boundedness is undecidable already for linear datalog programs with binary IDB predicates [20] and single rule programs (aka *sirups*) [15]. On the other hand, deciding boundedness is 2EXPTIME-complete for *monadic* datalog programs [7, 4] and PSPACE-complete for linear monadic programs [7]; for linear sirups, it is even NP-complete [20].

The last two results are relevant for deciding FO-rewritability of OMQs  $(Cov_A, \mathbf{q})$ , where  $\mathbf{q}$  has a single solitary  $F$  and is called a 1- $CQ$ . Indeed, suppose that  $F(x)$  and  $T(y_1), \dots, T(y_n)$  are all the solitary occurrences of  $F$  and  $T$  in  $\mathbf{q}$ . Let  $\Pi_{\mathbf{q}}$  be a monadic datalog program with three rules

$$G \leftarrow F(x), \mathbf{q}', P(y_1), \dots, P(y_n), \quad (1)$$

$$P(x) \leftarrow T(x), \quad (2)$$

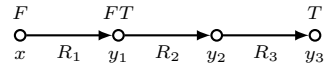
$$P(x) \leftarrow A(x), \mathbf{q}', P(y_1), \dots, P(y_n), \quad (3)$$

where  $\mathbf{q}' = \mathbf{q} \setminus \{F(x), T(y_1), \dots, T(y_n)\}$  and  $P$  is a fresh predicate symbol that never occurs in our ABoxes. Then, for any ABox  $\mathcal{A}$ , we have  $Cov_A, \mathcal{A} \models \mathbf{q}$  iff  $\Pi_{\mathbf{q}}, \mathcal{A} \models G$ . Thus, FO-rewritability of  $(Cov_A, \mathbf{q})$  is clearly related to boundedness of the sirup (3).

The problem of linearising datalog programs, that is, transforming them into equivalent linear datalog programs, which are known to be in NL for data complexity, has also attracted much attention [16, 17, 21, 2] after the Ullman and van Gelder pioneering paper [19].

The success of the OBDA paradigm spurred considerable interest in the data complexity of answering individual OMQs with expressive ontologies. Thus, by establishing a remarkable connection to CSPs, it was shown in [5] that deciding FO-rewritability and datalog rewritability of OMQs with a  $\mathcal{SHIU}$  ontology and a (Boolean) atomic query is NEXPTIME-complete. This result is obviously applicable to  $(Cov_A, \mathbf{q})$  with a *tree-shaped*  $\mathbf{q}$ . In [8], it was shown that checking rewritability of Boolean monadic *disjunctive* datalog programs into FO and monadic datalog can be done in, respectively, 2NEXPTIME and 3EXPTIME, which is applicable to all of our OMQs  $(Cov_A, \mathbf{q})$ .

An  $AC^0/NL/P$  trichotomy for the data complexity of answering OMQs with an  $\mathcal{EL}$  ontology and atomic query, which can be checked in EXPTIME, was established in [14]. This result is applicable to OMQs  $(Cov_A, \mathbf{q})$ , in which  $\mathbf{q}$  is an *F-tree* having a single solitary  $F(x)$  such that the binary atoms in  $\mathbf{q}$  form a ditree with root  $x$ . Indeed, denote by  $\mathcal{T}_{\mathbf{Q}}$  the  $\mathcal{EL}$  TBox with concept inclusions  $F \sqcap C_{\mathbf{q}} \sqsubseteq G'$ ,  $T \sqsubseteq P$  and  $A \sqcap C'_{\mathbf{q}} \sqsubseteq P$ , where  $C_{\mathbf{q}}$  is an  $\mathcal{EL}$ -concept representing  $\mathbf{q} \setminus \{F(x)\}$  with  $P$  for  $T$  (so for  $\mathbf{q}$  of the form



$C_{\mathbf{q}} = \exists R_1.(F \sqcap P \sqcap \exists R_2.\exists R_3.P)$ ). Then, for any ABox  $\mathcal{A}$  that does not contain  $G'$ , we have  $\Pi_{\mathbf{Q}}, \mathcal{A} \models G$  iff  $\mathcal{T}_{\mathbf{Q}}, \mathcal{A} \models \exists x G'(x)$ .

Yet, despite all of these efforts and results (implying, in view of the recent positive solution to the Feder-Vardi conjecture [6, 22], that there is a P/CONP dichotomy for OMQs with a  $\mathcal{SHIU}$  ontology and a (Boolean) atomic query, which is decidable in NEXPTIME), we are still lacking simple and transparent, in particular syntactic, conditions guaranteeing this or that data complexity or type of rewritability. Some results in this direction were obtained in [10, 11]. That a transparent classification of monadic sirups according to their data complexity has not been found so far and the close connection to CSPs indicate that this problem is extremely hard in general. Our aim below is to demonstrate that, for the OMQs with covering considered in this paper, syntactic conditions do exist but are difficult to find.

## 4 0-CQs

As observed in [9], if an OMQ  $Q$  (of the form defined above) is in  $AC^0$ , then  $q$  is an FO-rewriting of  $Q$ . By a 0-CQ we mean any CQ that does not contain a solitary  $F$ . A twin in a CQ  $q$  is any pair  $F(x), T(x) \in q$ . Here is an encouraging syntactic criterion of FO-rewritability for OMQs of the form  $(Cov_A^\perp, q)$ :

**Theorem 1.** (i) If  $q$  is a 0-CQ, then answering both  $(Cov_A^\perp, q)$  and  $(Cov_A, q)$  is in  $AC^0$ , with  $q$  being an FO-rewriting of these OMQs.

(ii) If  $q$  is not a 0-CQ and does not contain twins, then answering both  $(Cov_\top^\perp, q)$  and  $(Cov_\top, q)$  is L-hard.

*Proof.* (i) follows from [9, Theorem 1]. (ii) The proof is by reduction to the reachability problem for undirected graphs, which is known to be L-complete; see, e.g., [3]. Denote by  $q'$  the CQ obtained from  $q$  by gluing together all the variables  $x$  with  $F(x) \in q$  and also all the variables  $y$  with  $T(y) \in q$ . Thus,  $q'$  contains a single  $F$ -atom,  $F(x)$ , and a single  $T$ -atom,  $T(y)$ . Clearly, there is a homomorphism  $h: q \rightarrow q'$ . Let  $q'' = q' \setminus \{F(x), T(y)\}$ .

Suppose we are given an undirected graph  $G = (V, E)$  and two vertices  $s, t \in V$ . We regard  $G$  as a directed graph such that  $(u, v) \in E$  iff  $(v, u) \in E$ , for any  $u, v \in V$ . Now, we encode  $G$  by means of an ABox  $\mathcal{A}_G$  that is obtained from  $G$  as follows. For every edge  $e = (u, v) \in E$ , let  $q_e''$  be the set of atoms in  $q''$  with  $x$  renamed to  $u$ ,  $y$  to  $v$  and all other variables  $z$  to  $z_e$ . Then  $\mathcal{A}_G$  comprises all such  $q_e''$ , for  $e \in E$ , as well as  $F(s)$  and  $T(t)$ . We show that  $s \rightarrow_G t$  iff  $Cov_\top, \mathcal{A}_G \models q$ .

Suppose that  $s \rightarrow_G t$ , i.e., there exists a path  $s = v_0, \dots, v_n = t$  in  $G$  with  $e_i = (v_i, v_{i+1}) \in E$ , for  $i < n$ . Consider an arbitrary model  $\mathcal{I}$  of  $Cov_\top$  and  $\mathcal{A}_G$ . Since  $\mathcal{I} \models Cov_\top$ , and  $F(s)$  and  $T(t)$  are in  $\mathcal{A}_G$ , we can find some  $i < n$  such that  $\mathcal{I} \models F(v_i)$  and  $\mathcal{I} \models T(v_{i+1})$ . As  $q_{e_i}''$  is an isomorphic copy of  $q''$ , we obtain  $\mathcal{I} \models q'$ , and so  $\mathcal{I} \models q$ . Conversely, suppose  $s \not\rightarrow_G t$ . Define an interpretation  $\mathcal{I}$ , extending the ABox  $\mathcal{A}_G$ , by setting  $F^\mathcal{I}$  to be the set of objects in  $\mathcal{A}_G$  that are reachable from  $s$  and  $T^\mathcal{I}$  its complement. Clearly,  $\mathcal{I}$  is a model of  $Cov_\top$ . By the construction, the elements of the connected component of  $\mathcal{I}$  containing  $s$  cannot be instances of  $T$ , while the remaining elements of  $\mathcal{I}$  are not be instances of  $F$ . Since  $q$  is connected, it follows that  $\mathcal{I} \not\models q$ .

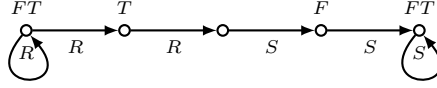
**Corollary 1.** An OMQ  $(Cov_A^\perp, q)$  is in  $AC^0$  iff  $q$  is a 0-CQ, which can be decided in linear time.

If twins can occur in CQs (that is,  $F$  and  $T$  are not necessarily disjoint), the picture becomes more complex. On one hand, we have the following criterion for OMQs  $(Cov_A, q)$  with a path CQ  $q$  whose variables  $x_0, \dots, x_n$  in  $q$  are ordered so that the binary atoms in  $q$  form a chain  $R_1(x_0, x_1), \dots, R_n(x_{n-1}, x_n)$ .

**Theorem 2 ([9]).** An OMQ  $(Cov_A, q)$  with a path CQ  $q$  is in  $AC^0$  iff  $q$  is a 0-CQ. If  $q$  contains both solitary  $F$  and  $T$ , then  $(Cov_A, q)$  is NL-hard.

On the other hand, this  $AC^0$ /NL criterion collapses for path CQs with loops:

**Proposition 1.** The OMQ  $(Cov_A, q)$ , where  $q$  is shown below, is in  $AC^0$ .



*Proof.* Follows from the sufficient condition of Theorem 4 (i).

Note that the CQ  $q$  above is *minimal* (not equivalent to any of its proper sub-CQs). Note also that, if a minimal 1-CQ  $q$  contains both a solitary  $F$  and a solitary  $T$ , then FO-rewritability of  $(Cov_A, q)$  implies that  $q$  contains at least one twin ( $FT$ ) and at least one  $y$  with  $T(y) \notin q$  and  $F(y) \notin q$  (which can be shown using Theorem 4 (i)).

## 5 1-CQs

1-CQs have exactly one solitary  $F$ . As observed in Section 3, we have the following:

**Theorem 3.** (i) Answering any OMQ  $(Cov_A, q)$  with a 1-CQ  $q$  can be done in P.

(ii) Answering any OMQ  $(Cov_A, q)$  with an  $F$ -tree  $q$  is either in  $AC^0$  or NL-complete or P-complete. The trichotomy can be decided in EXPTIME.

Theorem 3 (ii) was proved by a reduction to the  $AC^0$ /NL/P-trichotomy of [14]. It is to be noted, however, that applying the algorithm from [14] in our case is tricky because the input ontology  $\mathcal{T}_Q$  must first be converted to a normal form. As a result, we do not obtain transparent syntactic criteria on the shape of  $q$  that would guarantee that the OMQ  $(Cov_A, q)$  belongs to the desired complexity class (see examples below).

We now give a semantic sufficient condition for an OMQ with a 1-CQ to lie in NL. This condition uses ideas and constructions from [7, 14]. Let  $Q = (Cov_A, q)$  be an OMQ with a 1-CQ  $q$  having a solitary  $F(x)$ . Define by induction a class  $\mathfrak{K}_Q$  of ABoxes that will be called *cactuses for Q*. We start by setting  $\mathfrak{K}_Q := \{q\}$ , regarding  $q$  as an ABox, and then recursively apply to  $\mathfrak{K}_Q$  the following two rules:

- (bud) if  $T(y) \in \mathcal{A} \in \mathfrak{K}_Q$  with solitary  $T(y)$ , then we add to  $\mathfrak{K}_Q$  the ABox obtained by replacing  $T(y)$  in  $\mathcal{A}$  with  $(q \setminus F(x)) \cup \{A(x)\}$ , in which  $x$  is renamed to  $y$  and all of the other variables are given *fresh* names;
- (prune) if  $Cov_A, \mathcal{A}' \models Q$ , where  $\mathcal{A}' = \mathcal{A} \setminus \{T(y)\}$  and  $T(y)$  is solitary, we add to  $\mathfrak{K}_Q$  the ABox obtained by removing  $T(y)$  from  $\mathcal{A} \in \mathfrak{K}_Q$ .

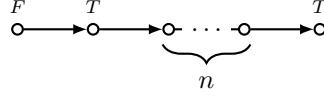
It is readily seen that, for any ABox  $\mathcal{A}'$ , we have  $Cov_A, \mathcal{A}' \models Q$  iff there exist  $\mathcal{A} \in \mathfrak{K}_Q$  and a homomorphism  $h: \mathcal{A} \rightarrow \mathcal{A}'$ . Denote by  $\mathfrak{K}_Q^\dagger$  the set of minimal cactuses in  $\mathfrak{K}_Q$  (that have no proper sub-cactuses in  $\mathfrak{K}_Q$ ).

For a cactus  $\mathcal{C} \in \mathfrak{K}_Q$ , we refer to the copies of (maximal subsets of)  $q$  that comprise  $\mathcal{C}$  as *segments*. The *skeleton*  $\mathcal{C}^s$  of  $\mathcal{C}$  is the ditree whose nodes are the segments  $\mathfrak{s}$  of  $\mathcal{C}$  and edges  $(\mathfrak{s}, \mathfrak{s}')$  mean that  $\mathfrak{s}'$  was attached to  $\mathfrak{s}$  by budding. The atoms  $T(y) \in \mathfrak{s}$  are called the *buds* of  $\mathfrak{s}$ . The *rank*  $r(\mathfrak{s})$  of  $\mathfrak{s}$  is defined by induction: if  $\mathfrak{s}$  is a leaf, then  $r(\mathfrak{s}) = 0$ ; for non-leaf  $\mathfrak{s}$ , we compute the maximal rank  $m$  of its children and then set

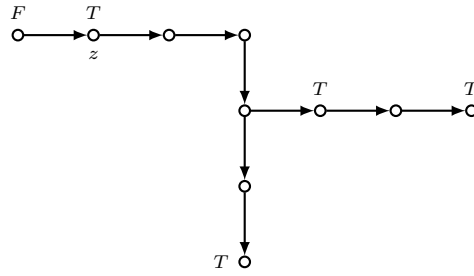
$$r(\mathfrak{s}) = \begin{cases} m + 1, & \text{if } \mathfrak{s} \text{ has } \geq 2 \text{ children of rank } m; \\ m, & \text{otherwise.} \end{cases}$$

The *width* of  $\mathcal{C}$  and  $\mathcal{C}^s$  is the rank of the root in  $\mathcal{C}^s$ . We say that  $\mathfrak{R}_Q^\dagger$  is of *width*  $k$  if it contains a cactus of width  $k$  but no cactus of greater width. The *depth* of  $\mathcal{C}$  and  $\mathcal{C}^s$  is the number of edges in the longest branch in  $\mathcal{C}^s$ .

We illustrate the definition by an example. Denote by  $q_{T^n T}$ , for  $n \geq 0$ , the 1-CQ shown below, where all the binary predicates are  $R$  and the  $n$  variables without labels do not occur in  $F$ - or  $T$ -atoms:



*Example 1.* Let  $Q = (\text{Cov}_\top, q_{T^1 T})$ . In the picture below, we show a cactus  $\mathcal{C}$  obtained by applying **(bud)** twice to  $q_{T^1 T}$  (with  $A = \top$  omitted):



One can check that  $\text{Cov}_\top, \mathcal{C} \setminus \{T(z)\} \models q_{T^1 T}$ , and so an application of **(prune)** will remove  $T(z)$  from  $\mathcal{C}$ . Using this observation, one can show that  $\mathfrak{R}_Q^\dagger$  is of width 1. On the other hand, if  $Q = (\text{Cov}_A, q_{T^1 T})$  then  $\mathfrak{R}_Q^\dagger$  is of unbounded width as follows from Theorem 6 below.

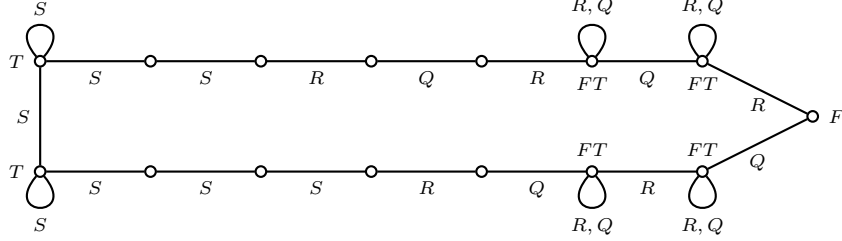
**Theorem 4.** Let  $Q = (\text{Cov}_A, q)$  be an OMQ with a 1-CQ  $q$ . Then

- (i)  $Q$  is in  $\text{AC}^0$  iff for every  $\mathcal{C} \in \mathfrak{R}_Q^\dagger$ , there is a homomorphism  $h: q \rightarrow \mathcal{C}$ ;
- (ii)  $Q$  is rewritable in linear datalog, and so is in NL, if  $\mathfrak{R}_Q^\dagger$  is of bounded width.

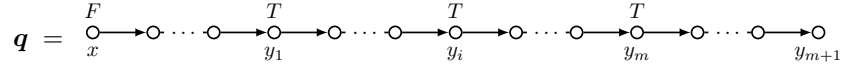
*Proof.* The proof of (i) is straightforward. To prove (ii), we represent the ABoxes in  $\mathfrak{R}_Q$  as words in a tree alphabet and construct a non-deterministic finite tree automaton  $\mathfrak{A}_Q$  such that  $\mathfrak{R}_Q^\dagger \subseteq L(\mathfrak{A}_Q) \subseteq \mathfrak{R}_Q$ . Then we show that, for  $\mathfrak{R}_Q^\dagger$  of finite width, the automaton  $\mathfrak{A}_Q$  can be transformed into a monadic linear-stratified datalog rewriting of  $Q$ . It can further be converted into a linear datalog rewriting at the expense of increasing the arity of IDBs in the program [2].

It is worth noting that, for  $Q = (\text{Cov}_\top, q)$  with  $q$  from Proposition 1,  $\mathfrak{R}_Q^\dagger$  consists of  $q$  and the cactus of depth 1, in which the only solitary  $T$  is removed by **(prune)**. Clearly, there is a homomorphism from  $q$  into this cactus, and so  $Q$  is FO-rewritable. However, for the 1-CQ  $q$  in the picture below (where all edges are bidirectional),  $(\text{Cov}_\top, q)$  is not FO-rewritable, but there is a homomorphism from  $q$  to both cactuses of depth 1. We do not know whether, in general, there is an upper bound  $N_q$  such that the existence of homomorphisms  $h: q \rightarrow \mathcal{C}$ , for all  $\mathcal{C} \in \mathfrak{R}_Q^\dagger$  of depth  $N_q$ , would ensure FO-rewritability of  $(\text{Cov}_A, q)$ . For 1-CQs  $q$  with a single solitary  $T$ , one can take  $N_q = |q| + 1$ . Neither do we know the exact complexity of deciding FO-rewritability of OMQs with 1-CQs.

As mentioned in Section 3, this problem is reducible to the boundedness problem for monadic datalog programs, which is known to be in 2EXPTIME.



Theorem 4 (ii) allows us to obtain a sufficient condition for linear-datalog rewritability of OMQs  $(Cov_A, \mathbf{q})$  with an  $F$ -path CQ  $\mathbf{q}$ , that is, a path CQ with a single solitary  $F$  at its root. We represent such a  $\mathbf{q}$  as shown in the picture below, which indicates *all* the solitary occurrences of  $F$  and  $T$ :



We require the following sub-CQs of  $\mathbf{q}$ :

- $\mathbf{q}_i$  is the suffix of  $\mathbf{q}$  that starts at  $y_i$ , but without  $T(y_i)$ , for  $1 \leq i \leq m$ ;
- $\mathbf{q}_i^*$  is the prefix of  $\mathbf{q}$  that ends at  $y_i$ , but without  $F(x)$  and  $T(y_i)$ , for  $1 \leq i \leq m$ ;
- $\mathbf{q}_{m+1}^*$  is  $\mathbf{q}$  without  $F(x)$ ,

and write  $f_i: \mathbf{q}_i \rightarrow \mathbf{q}$  if  $f_i$  is a homomorphism from  $\mathbf{q}_i$  into  $\mathbf{q}$  with  $f_i(y_i) = x$ .

**Theorem 5.** *If for each  $1 \leq i \leq m$  there exist  $f_i: \mathbf{q}_i \rightarrow \mathbf{q}$ , then  $(Cov_A, \mathbf{q})$  is rewritable into a linear datalog program, and so is NL-complete.*

For  $F$ -path CQs  $\mathbf{q}$  without twins, we extend Theorem 5 to a NL/P dichotomy (provided that  $NL \neq P$ ). Given such a CQ  $\mathbf{q}$ , we denote by  $N_{\mathbf{q}}$  the set of the numbers indicating the length of the path from  $x$  to each of the  $y_i$ ,  $i = 1, \dots, m + 1$ .

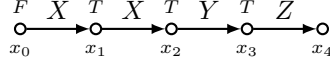
**Theorem 6.** *Let  $\mathbf{Q} = (Cov_A, \mathbf{q})$  be an OMQ where  $\mathbf{q}$  is an  $F$ -path CQ without twins having a single binary relation. The following are equivalent unless  $NL = P$ :*

- (i)  $\mathbf{Q}$  is NL-complete;
- (ii)  $\{0\} \cup N_{\mathbf{q}}$  is an arithmetic progression;
- (iii) there exist  $f_i: \mathbf{q}_i \rightarrow \mathbf{q}$  for every  $i = 1, \dots, m$ .

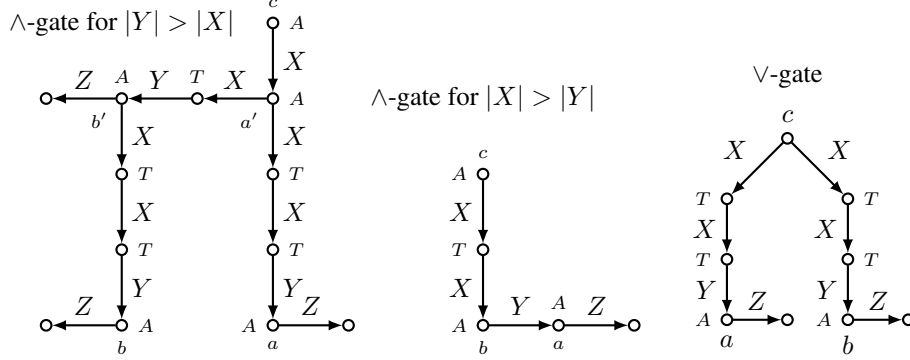
*If these conditions do not hold, then  $\mathbf{Q}$  is P-complete.*

*Proof.* It is readily seen, using Theorem 5, that (ii)  $\Rightarrow$  (iii)  $\Rightarrow$  (i). We prove the implication (i)  $\Rightarrow$  (ii). Suppose  $\{0\} \cup N_{\mathbf{q}}$  is not an arithmetic progression. We show that  $\mathbf{Q}$  is P-hard by reduction of the monotone circuit evaluation problem. Let  $X$  be the sub-CQ of  $\mathbf{q}$  between  $F(x)$  and the first  $T$  but without the  $F$  and  $T$  atoms, and let  $|X|$  be the number of atoms in  $X$ . Let  $Y$  be the first sub-CQ between neighbouring  $T$ -atoms without these  $T$ -atoms such that  $|Y| \neq |X|$ ,  $n$  the number of  $X$ -sub-CQs before  $Y$ , and let  $Z$  be the rest of the query without the first  $T$  atom; see the picture below where

$n = 2$ .



We distinguish between two cases:  $|Y| > |X|$  and  $|Y| < |X|$ . Depending on the case, we use the following gadgets for  $\wedge$ -gates, which are shown below for  $n = 2$  (it should be clear how to modify the construction for  $n > 2$ ); the  $\vee$ -gate gadget is the same for both cases.



Given any monotone Boolean circuit  $C$  and its input  $\alpha$ , we construct an ABox  $\mathcal{A}_{C,\alpha}$  by replacing  $\wedge$ - and  $\vee$ -gates with the inputs  $a$  and  $b$  and output  $c$  by the gadgets above, placing  $T$  on the input gates evaluated to 1 under  $\alpha$  and  $F$  on the output gate. We claim that  $Cov_A, \mathcal{A}_{C,\alpha} \models \mathbf{q}$  iff  $C(\alpha) = 1$ .

( $\Leftarrow$ ) It is not hard to show by induction that whenever a non-input gate  $g$  in  $C$  outputs 1 under  $\alpha$ , then the sub-ABox of  $\mathcal{A}_{C,\alpha}$  generated by  $g$  and with  $F$  placed on  $g$  validates  $\mathbf{Q}$ . For example, suppose  $T(a)$ ,  $T(b)$  and  $F(c)$  hold in the left-hand side gadget. Then either  $F(b')$  or  $T(b')$ . In the former case,  $\mathbf{q}$  is satisfied; so assume  $T(b')$ . We have either  $F(a')$  or  $T(a')$ . Now, in either case,  $\mathbf{q}$  is satisfied, as required.

( $\Rightarrow$ ) Suppose  $C(\alpha) = 0$ . We construct a model  $\mathcal{I}$  of  $Cov_A$  based on  $\mathcal{A}_{C,\alpha}$  by placing  $T$  or  $F$  to  $g$  and  $g'$  (if available) depending on the value of  $g$  in  $C$  under  $\alpha$ . It can be checked that  $\mathcal{I} \not\models \mathbf{q}$ . In particular, if  $|Y| > |X|$ ,  $c$  is the conjunction of  $b$  and  $a$ ,  $c = b = 0$  and  $a = 1$ , then if  $x_0$  goes to  $c$ ,  $x_1$  goes to  $a'$ , then  $x_2$  must go to the point below  $a'$ , and so  $x_3$  must go somewhere between  $a$  and the point above  $a$ , which is impossible. If  $a = 0$ ,  $b = 1$  and  $x_0$  goes to  $a'$  and  $x_1$  goes to the central  $T$ , then  $x_2$  must go somewhere between the central  $T$  and  $b'$ , which is also impossible. If  $|Y| < |X|$ ,  $c = b \wedge a$ ,  $b = d \vee e$ ,  $a = c = 0$ ,  $b = 1$ , then if  $x_0$  goes to  $c$ ,  $x_2$  goes to  $b$ , then  $x_3$  cannot go to  $a$  as  $a = 0$ , and so it must go into some  $X$ -segment of the gate  $b = d \vee e$  going out of  $b$ , which is again impossible.

Note that the proof of P-hardness in Theorem 6 does not go through for  $A = \top$ . Thus, for  $(Cov_{\top}, \mathbf{q}_{T1T})$ , we are in the framework of Example 1 and, by Theorem 4 (ii), this OMQ is in NL. In fact, we have the following NL/P dichotomy for the OMQs of the form  $(Cov_{\top}, \mathbf{q}_{TnT})$ .

**Theorem 7.** (i) *The OMQ  $(Cov_{\top}, \mathbf{q}_{T1T})$  is NL-complete (whereas  $(Cov_A, \mathbf{q}_{T1T})$  is P-complete).*



(ii) The OMQs  $(Cov_{\top}, \mathbf{q}_{TnT})$  (and  $(Cov_A, \mathbf{q}_{TnT})$ ), for  $n \geq 2$ , are P-complete.

*Proof.* The proof is similar to that of P-hardness in Theorem 6.

We now apply Theorem 4 (ii) to the class of  $TF$ -path CQs of the form

$$\mathbf{q}_{TF} = \begin{array}{c} T \qquad \qquad \qquad F \qquad \qquad \qquad T \qquad \qquad \qquad T \\ \circ \longrightarrow \circ \cdots \circ \longrightarrow \circ \longrightarrow \circ \cdots \circ \longrightarrow \circ \longrightarrow \circ \cdots \circ \longrightarrow \circ \longrightarrow \circ \cdots \circ \longrightarrow \circ \\ y_0 \qquad \qquad \qquad x \qquad \qquad \qquad y_1 \qquad \qquad \qquad y_m \qquad \qquad \qquad y_{m+1} \end{array}$$

where the  $T(y_i)$  and  $F(x)$  are all the solitary occurrences of  $T$  and  $F$  in  $\mathbf{q}_{TF}$ . We represent this CQ as

$$\mathbf{q}_{TF} = \{T(y_0)\} \cup \mathbf{q}_0 \cup \mathbf{q},$$

where  $\mathbf{q}_0$  is the sub-CQ of  $\mathbf{q}_{TF}$  between  $y_0$  and  $x$  with  $T(y_0)$  removed and  $\mathbf{q}$  is the same as in Theorem 5 (and  $\mathbf{q}_{m+1}^*$  is  $\mathbf{q}$  without  $F(x)$ ).

**Theorem 8.** *If  $\mathbf{q}$  satisfies the condition of Theorem 5 and there is a homomorphism  $h: \mathbf{q}_{m+1}^* \rightarrow \mathbf{q}_0$  such that  $h(x) = y_0$ , then answering  $(Cov_A, \mathbf{q}_{TF})$  is NL-complete.*

For example, the OMQ  $(Cov_A, \mathbf{q})$  with  $\mathbf{q}$  shown below is NL-complete:

$$\begin{array}{c} T \qquad \quad FT \qquad \quad F \qquad \quad T \\ \circ \longrightarrow \circ \longrightarrow \circ \longrightarrow \circ \longrightarrow \circ \end{array}$$

On the other hand, as shown in [9],  $(Cov_A, \mathbf{q})$  with  $\mathbf{q}$  of the form

$$\begin{array}{c} T \qquad \quad F \qquad \quad T \\ \circ \longrightarrow \circ \longrightarrow \circ \longrightarrow \circ \end{array}$$

is P-complete; in fact, it follows from the proof that  $(Cov_{\top}, \mathbf{q})$  is P-complete, too.

## 6 2-CQs

A 2-CQ has at least two solitary  $F$  and at least two solitary  $T$ . For example, as shown in [9], answering  $(Cov_A, \mathbf{q})$  with the following CQ  $\mathbf{q}$  is CONP-complete:

$$\begin{array}{c} T \qquad \quad T \qquad \quad F \qquad \quad F \\ \circ \longrightarrow \circ \longrightarrow \circ \longrightarrow \circ \longrightarrow \circ \end{array}$$

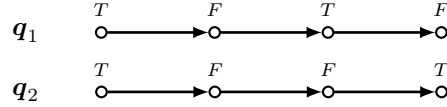
The proof can be generalised to the class of 2-2-CQs, which are path 2-CQs where all the  $F$  are located after all the  $T$ , and every occurrence of  $T$  or  $F$  is solitary. We represent any given 2-2-CQ  $\mathbf{q}$  as shown below

$$\begin{array}{c} T \qquad \quad T \qquad \quad F \qquad \quad F \\ \circ \longrightarrow \circ \longrightarrow \circ \longrightarrow \circ \longrightarrow \circ \\ p \qquad x \qquad r \qquad y \qquad s \qquad z \qquad u \qquad w \qquad v \end{array}$$

where  $p, r, u$  and  $v$  do not contain  $F$  and  $T$ , while  $s$  may contain solitary occurrences of both  $T$  and  $F$  (in other words, the  $T$  shown in the picture are the first two occurrences of  $T$  in  $\mathbf{q}$  and the  $F$  are the last two occurrences of  $F$  in  $\mathbf{q}$ ). Denote by  $\mathbf{q}_r$  the suffix of  $\mathbf{q}$  that starts from  $x$  but without  $T(x)$ ; similarly,  $\mathbf{q}_u$  is the suffix of  $\mathbf{q}$  starting from  $z$  but without  $F(z)$ . Denote by  $\mathbf{q}_r^-$  the prefix of  $\mathbf{q}$  that ends at  $y$  but without  $T(y)$ ; similarly,  $\mathbf{q}_u^-$  is the prefix of  $\mathbf{q}$  ending at  $w$  but without  $F(w)$ . Using the construction from [9], one can show the following:

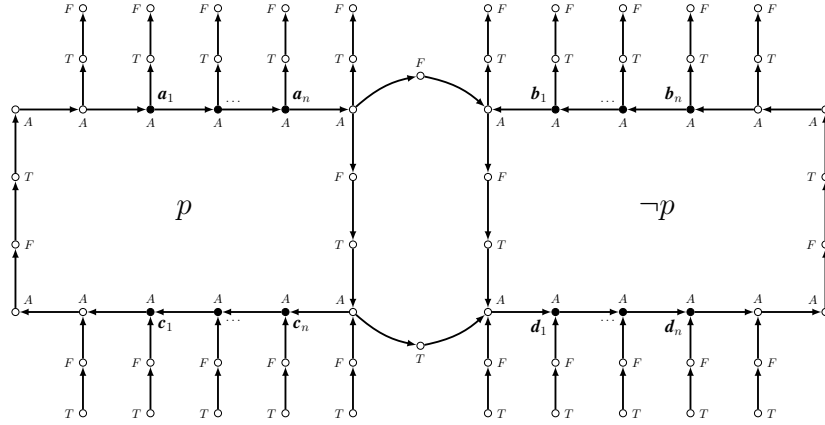
**Theorem 9.** Any OMQ  $(Cov_A, \mathbf{q})$  with a 2-2-CQ  $\mathbf{q}$  is CONP-complete provided the following conditions are satisfied: (i) there is no homomorphism  $h_1: \mathbf{q}_u \rightarrow \mathbf{q}_r$  with  $h_1(z) = x$ , and (ii) there is no homomorphism  $h_2: \mathbf{q}_r \rightarrow \mathbf{q}_u$  with  $h_2(y) = w$ .

Unfortunately, the proof does not generalise to the other two path 2-CQs with four variables, whose complexity has so far been open:



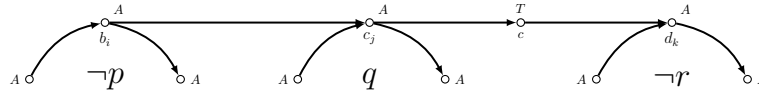
**Theorem 10.** The OMQs  $\mathbf{Q}_1 = (Cov_A, \mathbf{q}_1)$  and  $\mathbf{Q}_2 = (Cov_A, \mathbf{q}_2)$  are CONP-complete.

*Proof.* We prove CONP-hardness of  $\mathbf{Q}_1$  by reduction of the NP-complete 3SAT. Given a 3CNF  $\psi$ , we construct an ABox  $\mathcal{A}_\psi$  as follows. First, for every variable  $p$  occurring in  $\psi$ , we take the following  $p$ -gadget, where  $n$  is the number of clauses in  $\psi$ :



The key property of the  $p$ -gadget is that, for any model  $\mathcal{I}$  of  $Cov_A$  based on this gadget, if  $\mathcal{I} \not\models \mathbf{q}_1$  then either the  $A$ -points on left-hand side of the gadget are all in  $T^{\mathcal{I}}$  and the  $A$ -points on the right-hand side are all in  $F^{\mathcal{I}}$ , or the other way round. We refer to the  $a_i$  and  $b_i$  as  $p^\uparrow$ - and  $\neg p^\uparrow$ -contacts, and to the  $c_i$  and  $d_i$  as  $p^\downarrow$ - and  $\neg p^\downarrow$ -contacts, respectively.

Now, for every clause  $c = (l_1 \vee l_2 \vee l_3)$  in  $\psi$ , we add to the constructed gadgets for the variables in  $\psi$  the atoms  $R(u_{-l_1}^c, v_{l_2}^c)$ ,  $R(v_{l_2}^c, c)$ ,  $T(c)$ ,  $R(c, w_{l_3}^c)$ , where  $c$  is a new individual,  $u_{-l_1}^c$  a fresh  $\neg l_1^\uparrow$ -contact,  $v_{l_2}^c$  a fresh  $l_2^\downarrow$ -contact, and  $w_{l_3}^c$  a fresh  $l_3^\downarrow$ -contact. For example, for the clause  $c = (p \vee q \vee \neg r)$ , we obtain the fragment below:



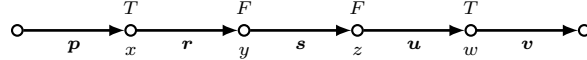
The resulting ABox  $\mathcal{A}_\psi$  is such that  $\psi$  is satisfiable iff  $Cov_A, \mathcal{A}_\psi \not\models \mathbf{q}_1$ .

For the OMQ  $\mathbf{Q}_2$ , we use (simplified)  $p$ -gadgets presented in Theorem 11 and connect them similarly to the picture above with  $T$  replaced by  $F$  at  $c$ .

We now sketch two generalisation of Theorem 10.

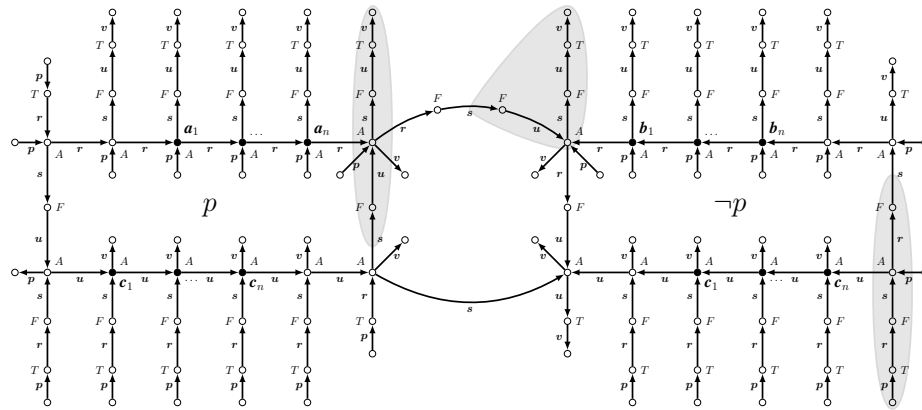
In Theorem 11 and 12, we assume that  $p$  and  $v$  do not contain  $F$  and  $T$ , while  $r$  and  $u$  may only contain solitary occurrences of  $T$  ( $F \notin r, u$ ), and  $s$  only solitary occurrences of  $F$  ( $T \notin s$ ).

**Theorem 11.** Any OMQ  $(Cov_A, q)$  with  $q$  of the form

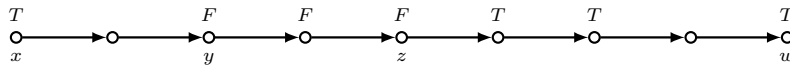


is CONP-complete provided the following conditions are satisfied: (i) there is no homomorphism  $h_1: r_t \rightarrow u$  with  $h_1(y) = w$ , and (ii) there is no homomorphism  $h_2: u_t \rightarrow r$  with  $h_2(z) = x$ , where  $r$  is the sub-CQ of  $q$  between  $x$  and  $y$  without  $T(x)$ ,  $F(y)$ , and similarly for  $u$ ,  $r_t$  is  $r$  with  $T(x)$  and  $u_t$  is  $u$  with  $T(w)$ .

The proof uses the following  $p$ -gadget:

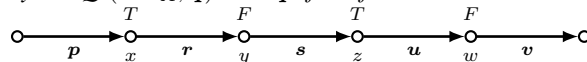


If one of the two conditions of this theorem is not satisfied then the given gadget will not work; see the shaded parts. For example, the complexity of the OMQ with the CQ below is still unknown:



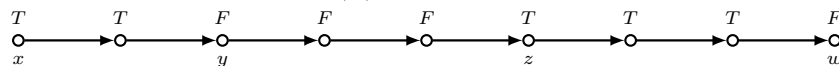
In Theorem 12, we use  $r^{ext} = r(x, y) \wedge T(y) \wedge s_1(y, y_1) \wedge F(y_1)$ , where  $s_1$  is the part of  $s$  such that  $s(y, z) = s_1(y, y_1) \wedge F(y_1) \wedge s_2(y_1, z)$  and  $s_1(y, y_1)$  does not contain any occurrences of  $F$ . In other words, the variable  $y_1$  corresponds to the first appearance of  $F$  in  $s$ , where  $s$  is the sub-CQ of  $q$  between  $y$  and  $z$  without  $F(y)$ ,  $T(z)$ .

**Theorem 12.** Any OMQ  $(Cov_A, q)$  with  $q$  of the form



is CONP-complete provided the following conditions hold: (i) there is no homomorphism  $g_1: r_t \rightarrow u$  with  $g_1(y) = w$ , and (ii) there is no homomorphism  $g_2: u \rightarrow r^{ext}$  with  $g_2(z) = x$  and  $g_2(w) = y_1$ .

The CQ below does not satisfy (ii), and its complexity is unknown:



It is also unknown whether there are OMQs with 2CQs in P.

**Acknowledgements.** The work of O. Gerasimova and M. Zakharyashev was carried out at the National Research University Higher School of Economics and supported by the Russian Science Foundation under grant 17-11-01294.

## References

1. Abiteboul, S., Hull, R., Vianu, V.: Foundations of Databases. Addison-Wesley (1995)
2. Afrati, F.N., Gergatsoulis, M., Toni, F.: Linearisability on datalog programs. *Theor. Comput. Sci.* 308(1-3), 199–226 (2003), [https://doi.org/10.1016/S0304-3975\(02\)00730-2](https://doi.org/10.1016/S0304-3975(02)00730-2)
3. Arora, S., Barak, B.: Computational Complexity: A Modern Approach. Cambridge University Press, New York, NY, USA, 1st edn. (2009)
4. Benedikt, M., ten Cate, B., Colcombet, T., Vanden Boom, M.: The complexity of boundedness for guarded logics. In: 30th Annual ACM/IEEE Symposium on Logic in Computer Science, LICS 2015, Kyoto, Japan, July 6-10, 2015. pp. 293–304. IEEE Computer Society (2015), <https://doi.org/10.1109/LICS.2015.36>
5. Bienvenu, M., ten Cate, B., Lutz, C., Wolter, F.: Ontology-based data access: A study through disjunctive datalog, CSP, and MMSNP. *ACM Transactions on Database Systems* 39(4), 33:1–44 (2014)
6. Bulatov, A.A.: A dichotomy theorem for nonuniform csp. In: Umans, C. (ed.) 58th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2017, Berkeley, CA, USA, October 15-17, 2017. pp. 319–330. IEEE Computer Society (2017), <https://doi.org/10.1109/FOCS.2017.37>
7. Cosmadakis, S.S., Gaifman, H., Kanellakis, P.C., Vardi, M.Y.: Decidable optimization problems for database logic programs (preliminary report). In: STOC. pp. 477–490 (1988)
8. Feier, C., Kuusisto, A., Lutz, C.: Rewritability in monadic disjunctive datalog, mmsnp, and expressive description logics (invited talk). In: Benedikt, M., Orsi, G. (eds.) 20th International Conference on Database Theory, ICDT 2017, March 21-24, 2017, Venice, Italy. LIPICs, vol. 68, pp. 1:1–1:17. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik (2017), <https://doi.org/10.4230/LIPICs.ICDT.2017.1>
9. Gerasimova, O., Kikot, S., Podolskii, V., Zakharyashev, M.: On the data complexity of ontology-mediated queries with a covering axiom. In: Proceedings of the 30th International Workshop on Description Logics (2017)
10. Hernich, A., Lutz, C., Ozaki, A., Wolter, F.: Schema.org as a description logic. In: Calvanese, D., Konev, B. (eds.) Proceedings of the 28th International Workshop on Description Logics, Athens, Greece, June 7-10, 2015. CEUR Workshop Proceedings, vol. 1350. CEUR-WS.org (2015), <http://ceur-ws.org/Vol-1350/paper-24.pdf>
11. Kaminski, M., Nenov, Y., Grau, B.C.: Datalog rewritability of disjunctive datalog programs and non-Horn ontologies. *Artif. Intell.* 236, 90–118 (2016), <http://dx.doi.org/10.1016/j.artint.2016.03.006>
12. Kharlamov, E., Hovland, D., Skjæveland, M.G., Bilidas, D., Jiménez-Ruiz, E., Xiao, G., Soyly, A., Lanti, D., Rezk, M., Zheleznyakov, D., Giese, M., Lie, H., Ioannidis, Y.E., Kotidis, Y., Koubarakis, M., Waaler, A.: Ontology based data access in statoil. *J. Web Sem.* 44, 3–36 (2017), <https://doi.org/10.1016/j.websem.2017.05.005>
13. Lanti, D., Rezk, M., Xiao, G., Calvanese, D.: The NPD benchmark: Reality check for OBDA systems. In: Alonso, G., Geerts, F., Popa, L., Barceló, P., Teubner, J., Ugarte, M., den Bussche, J.V., Paredaens, J. (eds.) Proceedings of the 18th International Conference on Extending Database Technology, EDBT 2015, Brussels, Belgium, March 23-27, 2015. pp. 617–628. OpenProceedings.org (2015), <https://doi.org/10.5441/002/edbt.2015.62>

14. Lutz, C., Sabellek, L.: Ontology-mediated querying with the description logic EL: trichotomy and linear datalog rewritability. In: Sierra, C. (ed.) Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017. pp. 1181–1187. ijcai.org (2017), <https://doi.org/10.24963/ijcai.2017/164>
15. Marcinkowski, J.: DATALOG sirups uniform boundedness is undecidable. In: Proceedings, 11th Annual IEEE Symposium on Logic in Computer Science, New Brunswick, New Jersey, USA, July 27-30, 1996. pp. 13–24. IEEE Computer Society (1996), <https://doi.org/10.1109/LICS.1996.561299>
16. Ramakrishnan, R., Sagiv, Y., Ullman, J.D., Vardi, M.Y.: Proof-tree transformation theorems and their applications. In: Proceedings of the eighth ACM SIGACT-SIGMOD-SIGART symposium on Principles of database systems. pp. 172–181. ACM (1989)
17. Saraiya, Y.P.: Linearizing nonlinear recursions in polynomial time. In: Silberschatz, A. (ed.) Proceedings of the Eighth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, March 29-31, 1989, Philadelphia, Pennsylvania, USA. pp. 182–189. ACM Press (1989), <http://doi.acm.org/10.1145/73721.73740>
18. Schaerf, A.: On the complexity of the instance checking problem in concept languages with existential quantification. *J. of Intelligent Information Systems* 2, 265–278 (1993)
19. Ullman, J.D., Gelder, A.V.: Parallel complexity of logical query programs. *Algorithmica* 3, 5–42 (1988), <https://doi.org/10.1007/BF01762108>
20. Vardi, M.Y.: Decidability and undecidability results for boundedness of linear recursive queries. In: Edmondson-Yurkkanan, C., Yannakakis, M. (eds.) Proceedings of the Seventh ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, March 21-23, 1988, Austin, Texas, USA. pp. 341–351. ACM (1988), <http://doi.acm.org/10.1145/308386.308470>
21. Zhang, W., Yu, C.T., Troy, D.: Necessary and sufficient conditions to linearize double recursive programs in logic databases. *ACM Trans. Database Syst.* 15(3), 459–482 (1990), <http://doi.acm.org/10.1145/88636.89237>
22. Zhuk, D.: A proof of CSP dichotomy conjecture. In: Umans, C. (ed.) 58th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2017, Berkeley, CA, USA, October 15-17, 2017. pp. 331–342. IEEE Computer Society (2017), <https://doi.org/10.1109/FOCS.2017.38>