



# The Complexity of Ontology-Based Data Access with OWL 2 QL and Bounded Treewidth Queries

Meghyn Bienvenu, Stanislav Kikot, Roman Kontchakov, Vladimir V Podolskii, Vladislav Ryzhikov, Michael Zakharyashev

## ► To cite this version:

Meghyn Bienvenu, Stanislav Kikot, Roman Kontchakov, Vladimir V Podolskii, Vladislav Ryzhikov, et al.. The Complexity of Ontology-Based Data Access with OWL 2 QL and Bounded Treewidth Queries. PODS: Principles of Database Systems, Jun 2017, Chicago, United States. hal-01632638

**HAL Id: hal-01632638**

**<https://inria.hal.science/hal-01632638>**

Submitted on 10 Nov 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# The Complexity of Ontology-Based Data Access with OWL 2 QL and Bounded Treewidth Queries

Meghyn Bienvenu<sup>1</sup>    Stanislav Kikot<sup>2</sup>    Roman Kontchakov<sup>2</sup>  
 Vladimir V. Podolskii<sup>3</sup>    Vladislav Ryzhikov<sup>4</sup>  
 Michael Zakharyashev<sup>2</sup>

<sup>1</sup> CNRS & University of Montpellier, France

<sup>2</sup> Birkbeck, University of London, UK

<sup>3</sup> Steklov Mathematical Institute & National Research University  
 Higher School of Economics, Moscow, Russia

<sup>4</sup> Free University of Bozen-Bolzano, Italy

## Abstract

Our concern is the overhead of answering *OWL 2 QL* ontology-mediated queries (OMQs) in ontology-based data access compared to evaluating their underlying tree-shaped and bounded treewidth conjunctive queries (CQs). We show that OMQs with bounded-depth ontologies have nonrecursive datalog (NDL) rewritings that can be constructed and evaluated in **LOGCFL** for combined complexity, even in **NL** if their CQs are tree-shaped with a bounded number of leaves, and so incur no overhead in complexity-theoretic terms. For OMQs with arbitrary ontologies and bounded-leaf CQs, NDL-rewritings are constructed and evaluated in **LOGCFL**. We show experimentally feasibility and scalability of our rewritings compared to previously proposed NDL-rewritings. On the negative side, we prove that answering OMQs with tree-shaped CQs is not fixed-parameter tractable if the ontology depth or the number of leaves in the CQs is regarded as the parameter, and that answering OMQs with a fixed ontology (of infinite depth) is **NP**-complete for tree-shaped and **LOGCFL** for bounded-leaf CQs.

**Keywords:** Ontology-based data access; ontology-mediated query; query rewriting; combined & parameterised complexity.

## 1 Introduction

The main aim of ontology-based data access (OBDA) [49, 42] is to facilitate access to complex data for non-expert end-users. The ontology, given by a logical theory  $\mathcal{T}$ , provides a unified conceptual view of one or more data sources, so the users do not have to know the actual structure of the data and can formulate their queries in the vocabulary of the ontology, which is connected to the data schema by a mapping  $\mathcal{M}$ . The instance  $\mathcal{M}(\mathcal{D})$  obtained by applying  $\mathcal{M}$  to a given dataset  $\mathcal{D}$  is interpreted under the open-world assumption, and additional facts can be *inferred* using the domain knowledge provided by the ontology. A certain answer to a query  $q(\mathbf{x})$  over  $\mathcal{D}$  is any tuple of constants  $\mathbf{a}$  such that  $\mathcal{T}, \mathcal{M}(\mathcal{D}) \models q(\mathbf{a})$ . OBDA is closely related to querying incomplete databases under (ontological) constraints, data integration [19], and data exchange [2].

In the classical approach to OBDA [12, 49], the computation of certain answers is reduced to standard database query evaluation: given an ontology-mediated query (OMQ)  $Q = (\mathcal{T}, q(\mathbf{x}))$ , one constructs a first-order (FO) query  $q'(\mathbf{x})$ , called a rewriting of  $Q$ , such that, for all datasets  $\mathcal{D}$  and mappings  $\mathcal{M}$ ,

$$\mathcal{T}, \mathcal{M}(\mathcal{D}) \models q(\mathbf{a}) \quad \text{iff} \quad \mathcal{I}_{\mathcal{M}(\mathcal{D})} \models q'(\mathbf{a}), \quad (1)$$

where  $\mathcal{I}_{\mathcal{M}(\mathcal{D})}$  is the FO-structure comprised of the atoms in  $\mathcal{M}(\mathcal{D})$ . When the form of  $\mathcal{M}$  is appropriately restricted (e.g.,  $\mathcal{M}$  is a GAV mapping), one can further unfold  $q'(\mathbf{x})$  using  $\mathcal{M}$  to

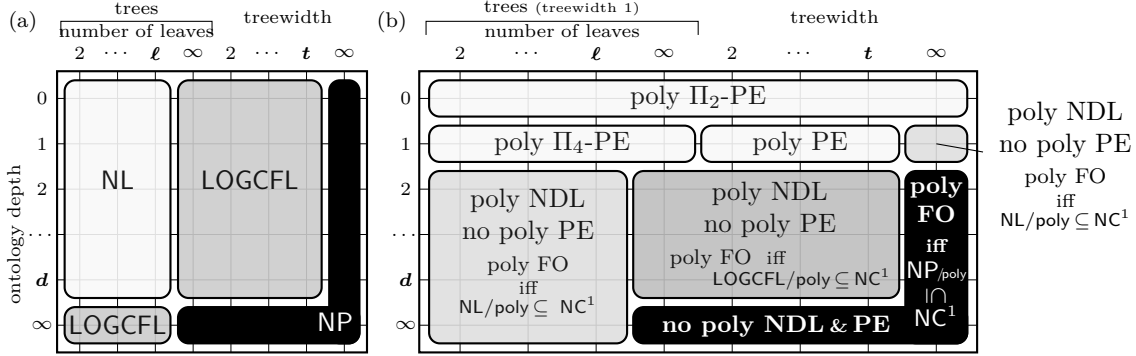


Figure 1: OMQ answering in *OWL 2 QL* (a) combined complexity and (b) the size of rewritings.

obtain an FO-query that can be evaluated directly over the original dataset  $\mathcal{D}$  (so there is no need to materialise  $\mathcal{M}(\mathcal{D})$ ).

For reduction (1) to hold for all OMQs, it is necessary to restrict the expressivity of  $\mathcal{T}$  and  $\mathbf{q}$ . The *DL-Lite* family of description logics [12] was specifically designed to ensure (1) for OMQs with conjunctive queries (CQs)  $\mathbf{q}$ . Other ontology languages with this property include linear and sticky tuple-generating dependencies (tgds) [9, 10], and the *OWL 2 QL* profile [44] of the W3C-standardised Web Ontology Language *OWL 2*, the focus of this work. Like many other ontology languages, *OWL 2 QL* admits only unary and binary predicates, but arbitrary relational instances can be queried due to the mapping. Various types of FO-rewritings  $\mathbf{q}'(\mathbf{x})$  have been developed and implemented for the preceding languages [49, 46, 40, 53, 14, 20, 52, 37, 27, 43, 39], and a few mature OBDA systems have emerged, including pioneering MASTRO [11], commercial Stardog [47] and Ultrawrap [54], and the Optique platform [23] with the query answering engine Ontop [50, 41].

Our concern here is the overhead of OMQ answering—i.e., checking whether the left-hand side of (1) holds—compared to evaluating the underlying CQs. At first sight, there is no apparent difference between the two problems when viewed through the lens of computational complexity: OMQ answering is in  $\text{AC}^0$  for data complexity by (1) and **NP**-complete for combined complexity [12], which in both cases corresponds to the complexity of evaluating CQs in the relational setting. Further analysis revealed, however, that answering OMQs is already **NP**-hard for combined complexity when the underlying CQs are tree-shaped (acyclic) [36], which sharply contrasts with the well-known LOGCFL-completeness of evaluating bounded treewidth CQs [61, 13, 26]. This surprising difference motivated a systematic investigation of the combined complexity of OMQ answering along two dimensions: (i) the query topology (treewidth  $t$  of CQs, and the number  $\ell$  of leaves in tree-shaped CQs), and (ii) the existential depth  $d$  of ontologies (i.e., the length of the longest chain of labelled nulls in the chase on any data). The resulting landscape, displayed in Fig. 1 (a) (under the assumption that datasets are given as RDF graphs and  $\mathcal{M}$  is the identity) [12, 36, 34, 5], indicates three tractable cases:

**OMQ( $d, t, \infty$ ):** ontologies of depth  $\leq d$  coupled with CQs of treewidth  $\leq t$  (for fixed  $d, t$ );

**OMQ( $d, 1, \ell$ ):** ontologies of depth  $\leq d$  with tree-shaped CQs with  $\leq \ell$  leaves (for fixed  $d, \ell$ );

**OMQ( $\infty, 1, \ell$ ):** ontologies of arbitrary depth and tree-shaped CQs with  $\leq \ell$  leaves (for fixed  $\ell$ ).

Observe in particular that when the ontology depth is bounded by a fixed constant, the complexity of OMQ answering is precisely the same as for evaluating the underlying CQs. If we place no restriction on the ontology, then tractability of tree-shaped queries can be recovered by bounding the number of leaves, but we have LOGCFL rather than the expected NL.

While the results in Fig. 1(a) appear to answer the question of the additional cost incurred by adding an *OWL 2 QL* ontology, they only tell part of the story. Indeed, in the context of classical

rewriting-based OBDA [49], it is not the abstract complexity of OMQ answering that matters, but the cost of computing and evaluating OMQ rewritings. Fig. 1(b) summarises what is known about the size of positive existential (PE), nonrecursive datalog (NDL) and FO-rewritings [35, 25, 34, 5]. Thus, we see, for example, that PE-rewritings for OMQs from  $\text{OMQ}(d, t, \infty)$  can be of super-polynomial size, and so are not computable and evaluable in polynomial time, even though Fig. 1(a) shows that such OMQs can be answered in LOGCFL. The same concerns  $\text{OMQ}(d, 1, \ell)$  and  $\text{OMQ}(\infty, 1, \ell)$ , which can be answered in NL and LOGCFL, respectively, but do not enjoy polynomial-size PE-rewritings. Moreover, our experiments show that standard rewriting engines exhibit exponential behaviour on OMQs drawn from  $\text{OMQ}(1, 1, 2)$  lying in the intersection of the three tractable classes.

Our first aim is to show that the positive complexity results in Fig. 1(a) can in fact be achieved using query rewriting. To this end, we develop NDL-rewritings for the three tractable cases that can be computed and evaluated by algorithms of optimal combined complexity. In theory, such algorithms are known to be space efficient and highly parallelisable. We demonstrate practical efficiency of our optimal NDL-rewritings by comparing them with the NDL-rewritings produced by Clipper [20], Presto [53] and Rapid [14], using a sequence of OMQs from the class  $\text{OMQ}(1, 1, 2)$ .

Our second aim is to understand the contribution of the ontology depth and the number of leaves in tree-shaped CQs to the complexity of OMQ answering. (As follows from Fig. 1 (a), if these parameters are unbounded, this problem is harder than evaluating the underlying CQs unless  $\text{LOGCFL} = \text{NP}$ .) Unfortunately, it turns out that answering OMQs with ontologies of finite depth and tree-shaped CQs is not fixed-parameter tractable if either the ontology depth or the number of leaves in CQs is regarded as a parameter. More precisely, we prove that the problem is  $W[2]$ -hard in the former case and  $W[1]$ -hard in the latter. These results suggest that the ontology depth and the number of leaves are inherently in the exponent of the size of the input in any OMQ answering algorithm.

Finally, we revisit the NP- and LOGCFL-hardness results for OMQs with tree-shaped CQs. The known NP and LOGCFL lower bounds have been established using sequences  $(\mathcal{T}_n, \mathbf{q}_n)$  of OMQs, where the depth of  $\mathcal{T}_n$  grows with  $n$  [36, 5]. One might thus hope to make answering OMQs with tree-shaped CQs easier by restricting the ontology signature, size, or even by fixing the whole ontology, which is very relevant for applications as a typical OBDA scenario has users posing different queries over the same ontology. Our third main result is that this is not the case: we present ontologies  $\mathcal{T}_{\dagger}$  and  $\mathcal{T}_{\ddagger}$  of infinite depth such that answering OMQs  $(\mathcal{T}_{\dagger}, \mathbf{q})$  with tree-shaped  $\mathbf{q}$  and  $(\mathcal{T}_{\ddagger}, \mathbf{q})$  with linear  $\mathbf{q}$  is NP- and LOGCFL-hard for query complexity, respectively. We also show that no algorithm can construct FO-rewritings of the OMQs  $(\mathcal{T}_{\dagger}, \mathbf{q})$  in polynomial time unless  $\text{P} = \text{NP}$ , even though polynomial-size FO-rewritings of these OMQs do exist.

The paper is organised as follows. We begin in Section 2 by introducing the *OWL 2 QL* ontology language and key notions like OMQ answering and query rewriting. In Section 3, we first identify fragments of NDL which can be evaluated in LOGCFL or NL, and then we use these results to develop NDL-rewritings of optimal combined complexity for the three tractable cases. Section 4 concerns the parameterised complexity of OMQ answering with tree-shaped CQs. For ontologies of finite depth, we show  $W[2]$ -hardness (resp.  $W[1]$ -hardness) when the ontology depth (resp. number of leaves) is taken as the parameter. For the infinite depth case, we show in Section 5 that NP-hardness applies even for a fixed ontology. The final section of the paper presents preliminary experiments comparing our new rewritings to those produced by existing rewriting engines and discusses possible directions for future work.

## 2 Preliminaries

An *OWL 2 QL ontology* (*TBox* in description logic),  $\mathcal{T}$ , is a finite set of sentences (*axioms*) of the forms

$$\begin{aligned} \forall x (\tau(x) \rightarrow \tau'(x)), & \quad \forall x (\tau(x) \wedge \tau'(x) \rightarrow \perp), \\ \forall xy (\varrho(x, y) \rightarrow \varrho'(x, y)), & \quad \forall xy (\varrho(x, y) \wedge \varrho'(x, y) \rightarrow \perp), \\ \forall x \varrho(x, x), & \quad \forall x (\varrho(x, x) \rightarrow \perp), \end{aligned}$$

where  $\tau(x)$  and  $\varrho(x, y)$  are defined, using unary predicates  $A$  and binary predicates  $P$ , by the grammars

$$\begin{aligned} \tau(x) &::= \top \mid A(x) \mid \exists y \varrho(x, y), \\ \varrho(x, y) &::= \top \mid P(x, y) \mid P(y, x). \end{aligned}$$

When writing ontology axioms, we omit the universal quantifiers and denote by  $\mathbf{R}_{\mathcal{T}}$  the set of binary predicates  $P$  occurring in  $\mathcal{T}$  and their inverses  $P^-$ , assuming that  $P^{--} = P$ . For every  $\varrho \in \mathbf{R}_{\mathcal{T}}$ , we take a fresh unary predicate  $A_{\varrho}$  and add  $A_{\varrho}(x) \leftrightarrow \exists y \varrho(x, y)$  to  $\mathcal{T}$  (where, as usual,  $\varphi \leftrightarrow \psi$  is an abbreviation for  $\varphi \rightarrow \psi$  and  $\psi \rightarrow \varphi$ ). The resulting ontology is said to be in *normal form*, and we assume, without loss of generality, that all our ontologies are in normal form.

A *data instance*,  $\mathcal{A}$ , is a finite set of unary or binary ground atoms (called an *ABox* in description logic). We denote by  $\text{ind}(\mathcal{A})$  the set of individual constants in  $\mathcal{A}$  and write  $\varrho(a, b) \in \mathcal{A}$  if  $P(a, b) \in \mathcal{A}$  and  $\varrho = P$ , or  $P(b, a) \in \mathcal{A}$  and  $\varrho = P^-$ . We say that  $\mathcal{A}$  is *complete* for an ontology  $\mathcal{T}$  if  $\mathcal{T}, \mathcal{A} \models S(\mathbf{a})$  implies  $S(\mathbf{a}) \in \mathcal{A}$ , for any ground atom  $S(\mathbf{a})$  with  $\mathbf{a} \subseteq \text{ind}(\mathcal{A})$ .<sup>1</sup>

A *conjunctive query* (CQ)  $\mathbf{q}(\mathbf{x})$  is a formula of the form  $\exists \mathbf{y} \varphi(\mathbf{x}, \mathbf{y})$ , where  $\varphi$  is a conjunction of atoms  $S(\mathbf{z})$  all of whose variables are among  $\text{var}(\mathbf{q}) = \mathbf{x} \cup \mathbf{y}$ . We assume, without loss of generality, that CQs contain no constants. We often regard a CQ as the set of its atoms. With every CQ  $\mathbf{q}$ , we associate its *Gaifman graph*  $\mathcal{G}$  whose vertices are the variables of  $\mathbf{q}$  and whose edges are the pairs  $\{u, v\}$  such that  $P(u, v) \in \mathbf{q}$ , for some  $P$ . We call  $\mathbf{q}$  *connected* if  $\mathcal{G}$  is connected, *tree-shaped* if  $\mathcal{G}$  is a tree, and *linear* if  $\mathcal{G}$  is a tree with two leaves.

An *ontology-mediated query* (OMQ) is a pair  $\mathbf{Q}(\mathbf{x}) = (\mathcal{T}, \mathbf{q}(\mathbf{x}))$ , where  $\mathcal{T}$  is an ontology and  $\mathbf{q}(\mathbf{x})$  a CQ. A tuple  $\mathbf{a} \subseteq \text{ind}(\mathcal{A})$  is a *certain answer* to  $\mathbf{Q}(\mathbf{x})$  over a data instance  $\mathcal{A}$  if  $\mathcal{I} \models \mathbf{q}(\mathbf{a})$  for all models  $\mathcal{I}$  of  $\mathcal{T}$  and  $\mathcal{A}$ ; in this case we write  $\mathcal{T}, \mathcal{A} \models \mathbf{q}(\mathbf{a})$ . If  $\mathbf{x} = \emptyset$ , then a certain answer to  $\mathbf{Q}$  over  $\mathcal{A}$  is ‘yes’ if  $\mathcal{T}, \mathcal{A} \models \mathbf{q}$  and ‘no’ otherwise. The *OMQ answering problem* (for a class of OMQs) is to decide whether  $\mathcal{T}, \mathcal{A} \models \mathbf{q}(\mathbf{a})$  holds, given an OMQ  $\mathbf{Q}(\mathbf{x})$  (in the class),  $\mathcal{A}$  and  $\mathbf{a} \subseteq \text{ind}(\mathcal{A})$ . If  $\mathcal{T}$ ,  $\mathbf{q}(\mathbf{x})$ , and  $\mathcal{A}$  are regarded as input, we speak about *combined complexity* of OMQ answering; if  $\mathcal{A}$  and  $\mathcal{T}$  are regarded as fixed, we speak about *query complexity*.

Every consistent *knowledge base* (KB)  $(\mathcal{T}, \mathcal{A})$  has a *canonical model* (or *chase* in database theory) [1]  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$  with the property that  $\mathcal{T}, \mathcal{A} \models \mathbf{q}(\mathbf{a})$  iff  $\mathcal{C}_{\mathcal{T}, \mathcal{A}} \models \mathbf{q}(\mathbf{a})$ , for all CQs  $\mathbf{q}(\mathbf{x})$  and  $\mathbf{a} \subseteq \text{ind}(\mathcal{A})$ . In our constructions, we use the following definition of  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ , where without loss of generality we assume that  $\mathcal{T}$  contains no binary predicates  $P$  with  $\mathcal{T} \models \forall xy P(x, y)$ . The domain,  $\Delta^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ , consists of  $\text{ind}(\mathcal{A})$  and the *witnesses* (or *labelled nulls*) of the form  $w = a\varrho_1 \dots \varrho_n$ , for  $n \geq 1$ , such that

- $a \in \text{ind}(\mathcal{A})$  and  $\mathcal{T}, \mathcal{A} \models \exists y \varrho_1(a, y)$ ;
- $\mathcal{T} \not\models \varrho_i(x, x)$ , for  $1 \leq i \leq n$ ;
- $\mathcal{T} \models \exists x \varrho_i(x, y) \rightarrow \exists z \varrho_{i+1}(y, z)$  but  $\mathcal{T} \not\models \varrho_i(x, y) \rightarrow \varrho_{i+1}(y, x)$ , for  $1 \leq i < n$ .

We denote by  $\mathbf{W}_{\mathcal{T}}$  the set of words  $\varrho_1 \dots \varrho_n \in \mathbf{R}_{\mathcal{T}}^*$  satisfying the last two conditions. Every  $a \in \text{ind}(\mathcal{A})$  is interpreted in  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$  by itself, and unary and binary predicates are interpreted as follows:

- $\mathcal{C}_{\mathcal{T}, \mathcal{A}} \models A(u)$  iff either  $u \in \text{ind}(\mathcal{A})$  and  $\mathcal{T}, \mathcal{A} \models A(u)$ , or  $u = w\varrho$  with  $\mathcal{T} \models \exists y \varrho(y, x) \rightarrow A(x)$ ;

<sup>1</sup>If the meaning is clear from the context, we use set-theoretic notation for lists.

- $\mathcal{C}_{\mathcal{T}, \mathcal{A}} \models P(u, v)$  iff one of the three conditions holds: (i)  $u, v \in \text{ind}(\mathcal{A})$  and  $\mathcal{T}, \mathcal{A} \models P(u, v)$ ;
- (ii)  $u = v$  and  $\mathcal{T} \models P(x, x)$ ; (iii)  $\mathcal{T} \models \varrho(x, y) \rightarrow P(x, y)$  and either  $v = u\varrho$  or  $u = v\varrho^-$ .

We say that  $\mathcal{T}$  is of *depth* 0 if it does not contain any axioms with  $\exists$  on the right-hand side, excepting the normalisation axioms<sup>2</sup>. Otherwise, we say that  $\mathcal{T}$  is of *depth*  $0 < d < \infty$  if  $d$  is the maximum length of the words in  $\mathbf{W}_{\mathcal{T}}$ , and it is of *depth*  $\infty$  if  $\mathbf{W}_{\mathcal{T}}$  is infinite. (Note that the depth of  $\mathcal{T}$  is computable in NL; cf. [24, 8] for related results on chase termination for tgds.)

An FO-formula  $q'(x)$ , possibly with equality, is an *FO-rewriting of an OMQ*  $Q(x) = (\mathcal{T}, q(x))$  if, for *any* data instance  $\mathcal{A}$  and any tuple  $a \subseteq \text{ind}(\mathcal{A})$ ,

$$\mathcal{T}, \mathcal{A} \models q(a) \quad \text{iff} \quad \mathcal{I}_{\mathcal{A}} \models q'(a), \quad (2)$$

where  $\mathcal{I}_{\mathcal{A}}$  is the FO-structure over the domain  $\text{ind}(\mathcal{A})$  such that  $\mathcal{I}_{\mathcal{A}} \models S(a)$  iff  $S(a) \in \mathcal{A}$ , for any ground atom  $S(a)$ . If  $q'(x)$  is a positive existential formula, we call it a *PE-rewriting of*  $Q(x)$ . A PE-rewriting whose matrix is a  $\Pi_k$ -formula (with respect to  $\wedge$  and  $\vee$ ) is called a  $\Pi_k$ -rewriting. The size  $|q'|$  of  $q'$  is the number of symbols in it.

We also consider rewritings in the form of nonrecursive datalog queries. A *datalog program*,  $\Pi$ , is a finite set of Horn clauses  $\forall z (\gamma_0 \leftarrow \gamma_1 \wedge \dots \wedge \gamma_m)$ , where each  $\gamma_i$  is an atom  $Q(y)$  with  $y \subseteq z$  or an equality ( $z = z'$ ) with  $z, z' \in z$ . (As usual, we omit  $\forall z$  from clauses.) The atom  $\gamma_0$  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  $\Pi$  are *IDB predicates*, the rest (including  $=$ ) *EDB predicates*. A predicate  $Q$  *depends* on  $P$  in  $\Pi$  if  $\Pi$  has a clause with  $Q$  in the head and  $P$  in the body.  $\Pi$  is a *nonrecursive datalog (NDL) program* if the (directed) *dependence graph* of the dependence relation is acyclic.

An *NDL query* is a pair  $(\Pi, G(x))$ , where  $\Pi$  is an NDL program and  $G(x)$  a predicate. A tuple  $a \subseteq \text{ind}(\mathcal{A})$  is an *answer to*  $(\Pi, G(x))$  over a data instance  $\mathcal{A}$  if  $G(a)$  holds in the first-order structure with domain  $\text{ind}(\mathcal{A})$  obtained by closing  $\mathcal{A}$  under the clauses in  $\Pi$ ; in this case we write  $\Pi, \mathcal{A} \models G(a)$ . The problem of checking whether  $a$  is an answer to  $(\Pi, G(x))$  over  $\mathcal{A}$  is called the *query evaluation problem*. The *depth* of  $(\Pi, G(x))$  is the length,  $d(\Pi, G)$ , of the longest directed path in the dependence graph for  $\Pi$  starting from  $G$ . NDL queries are *equivalent* if they have exactly the same answers over any data instance.

An NDL query  $(\Pi, G(x))$  is an *NDL-rewriting of an OMQ*  $Q(x) = (\mathcal{T}, q(x))$  over complete data instances in case  $\mathcal{T}, \mathcal{A} \models q(a)$  iff  $\Pi, \mathcal{A} \models G(a)$ , for any complete  $\mathcal{A}$  and any  $a \subseteq \text{ind}(\mathcal{A})$ . Rewritings over arbitrary data instances are defined by dropping the completeness condition. Given an NDL-rewriting  $(\Pi, G(x))$  of  $Q(x)$  over complete data instances, we denote by  $\Pi^*$  the result of replacing each predicate  $S$  in  $\Pi$  with a fresh IDB predicate  $S^*$  of the same arity and adding the clauses

$$\begin{array}{ll} A^*(x) \leftarrow \tau(x), & \text{if } \mathcal{T} \models \tau(x) \rightarrow A(x), \\ P^*(x, y) \leftarrow \varrho(x, y), & \text{if } \mathcal{T} \models \varrho(x, y) \rightarrow P(x, y), \\ P^*(x, x) \leftarrow \top(x), & \text{if } \mathcal{T} \models P(x, x), \end{array}$$

where  $\top(x)$  is an EDB predicate for the active domain [32]. Clearly,  $(\Pi^*, G(x))$  is an NDL-rewriting of  $Q(x)$  over arbitrary data instances and  $|\Pi^*| \leq |\Pi| + |\mathcal{T}|^2$ .

Finally, we remark that, without loss of generality, we can (and will) assume that our ontologies  $\mathcal{T}$  do not contain  $\perp$ . Indeed, we can always incorporate into rewritings subqueries that check whether the left-hand side of an axiom with  $\perp$  holds and output all tuples of constants if this is the case [9].

### 3 Optimal NDL-Rewritings

To construct theoretically optimal NDL-rewritings for OMQs in the three tractable classes, we first identify two types of NDL queries whose evaluation problems are in NL and LOGCFL for

<sup>2</sup>This somewhat awkward definition of depth 0 ontologies is due to the use of normalisation axioms, which may introduce unnecessary words on length 1 in  $\mathbf{W}_{\mathcal{T}}$ .

combined complexity.

### 3.1 NL and LOGCFL fragments of NDL

To simplify the analysis of non-Boolean NDL queries, it is convenient to regard certain variables as parameters to be instantiated with constants from the candidate answer. Formally, an NDL query  $(\Pi, G(x_1, \dots, x_n))$  is called *ordered* if each of its IDB predicates  $Q$  comes with fixed variables  $x_{i_1}, \dots, x_{i_k}$  ( $1 \leq i_1 < \dots < i_k \leq n$ ), called the *parameters of  $Q$* , such that (i) every occurrence of  $Q$  in  $\Pi$  is of the form  $Q(y_1, \dots, y_m, x_{i_1}, \dots, x_{i_k})$ , (ii) the parameters of  $G$  are  $x_1, \dots, x_n$ , and (iii) parameters of the head of every clause include all the parameters of the predicates in the body. Observe that Boolean NDL queries are trivially ordered. The *width*  $w(\Pi, G)$  of an ordered  $(\Pi, G)$  is the maximal number of non-parameter variables in a clause of  $\Pi$ .

**Example 1.** The NDL query  $(\Pi, G(x))$ , where

$$\Pi = \{ G(x) \leftarrow R(x, y) \wedge Q(x), \quad Q(x) \leftarrow R(y, x) \},$$

is ordered with parameter  $x$  and width 1 (the conditions do not restrict the EDB predicate  $R$ ). Replacing  $Q(x)$  by  $Q(y)$  in the first clause yields a query that is not ordered in view of (i). A further swap of  $Q(x)$  in the second clause with  $Q(y)$  would satisfy (i) but not (iii).

As all the NDL-rewritings we construct are ordered, with their parameters being the answer variables, from now on we only consider ordered NDL queries.

Given an NDL query  $(\Pi, G(\mathbf{x}))$ , a data instance  $\mathcal{A}$  and a tuple  $\mathbf{a}$  with  $|\mathbf{x}| = |\mathbf{a}|$ , the  *$\mathbf{a}$ -grounding  $\Pi_{\mathcal{A}}^{\mathbf{a}}$  of  $\Pi$  on  $\mathcal{A}$*  is the set of ground clauses obtained by first replacing each parameter in  $\Pi$  by the corresponding constant from  $\mathbf{a}$ , and then performing the standard grounding [17] of  $\Pi$  using the constants from  $\mathcal{A}$ . The size of  $\Pi_{\mathcal{A}}^{\mathbf{a}}$  is bounded by  $|\Pi| \cdot |\mathcal{A}|^{w(\Pi, G)}$ , and so we can check whether  $\Pi, \mathcal{A} \models G(\mathbf{a})$  holds in time  $\text{poly}(|\Pi| \cdot |\mathcal{A}|^{w(\Pi, G)})$ .

#### 3.1.1 Linear NDL in NL

An NDL program is *linear* [1] if the body of its every clause contains at most one IDB predicate.

**Theorem 2.** *For any  $w > 0$ , evaluation of linear NDL queries of width  $\leq w$  is NL-complete for combined complexity.*

*Proof.* Let  $(\Pi, G(\mathbf{x}))$  be a linear NDL query. Deciding whether  $\Pi, \mathcal{A} \models G(\mathbf{a})$  is reducible to finding a path to  $G(\mathbf{a})$  from a certain set  $X$  in the grounding graph  $\mathfrak{G}$  constructed as follows. The vertices of  $\mathfrak{G}$  are the IDB atoms of  $\Pi_{\mathcal{A}}^{\mathbf{a}}$ , and  $\mathfrak{G}$  has an edge from  $Q(\mathbf{c})$  to  $Q'(\mathbf{c}')$  iff  $\Pi_{\mathcal{A}}^{\mathbf{a}}$  contains  $Q'(\mathbf{c}') \leftarrow Q(\mathbf{c}) \wedge S_1(\mathbf{c}_1) \wedge \dots \wedge S_k(\mathbf{c}_k)$  with  $S_i(\mathbf{c}_i) \in \mathcal{A}$ , for  $1 \leq i \leq k$  (we assume  $\mathcal{A}$  contains all  $c = c$ , for  $c \in \text{ind}(\mathcal{A})$ ). The set  $X$  consists of all vertices  $Q(\mathbf{c})$  with IDB predicates  $Q$  being of in-degree 0 in the dependency graph of  $\Pi$  for which there is a clause  $Q(\mathbf{c}) \leftarrow S_1(\mathbf{c}_1) \wedge \dots \wedge S_k(\mathbf{c}_k)$  in  $\Pi_{\mathcal{A}}^{\mathbf{a}}$  with  $S_i(\mathbf{c}_i) \in \mathcal{A}$  ( $1 \leq i \leq k$ ). Bounding the width of  $(\Pi, G)$  ensures that  $\mathfrak{G}$  is of polynomial size and can be constructed by a deterministic Turing machine with read-only input, write-once output and logarithmic-size work tapes.  $\square$

The transformation  $*$  of NDL-rewritings over complete data instances into NDL-rewritings over arbitrary data instances does not preserve linearity. A more involved construction is given in the proof of the following:

**Lemma 3.** *Fix any  $w > 0$ . There is an  $\text{L}^{\text{NL}}$ -transducer that, for any linear NDL-rewriting  $(\Pi, G(\mathbf{x}))$  of an OMQ  $Q(\mathbf{x})$  over complete data instances with  $w(\Pi, G) \leq w$ , computes a linear NDL-rewriting  $(\Pi', G(\mathbf{x}))$  of  $Q(\mathbf{x})$  over arbitrary data instances such that  $w(\Pi', G) \leq w + 1$ .*

We note that a possible increase of the width by 1 is due to the ‘replacement’ of unary atoms  $A(z)$  by binary atoms  $\varrho(y, z)$  whenever  $\mathcal{T} \models \exists y \varrho(y, z) \rightarrow A(z)$ .

### 3.1.2 Skinny NDL in LOGCFL

The complexity class LOGCFL can be defined using *nondeterministic auxiliary pushdown automata* (NAuxPDAs) [15], which are nondeterministic Turing machines with an additional work tape constrained to operate as a pushdown store. Sudborough [57] 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). It is known that LOGCFL can equivalently be defined in terms of logspace-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. Moreover, there is an algorithm that, given such a circuit  $\mathbf{C}$ , computes the output using an NAuxPDA in logarithmic space in the size of  $\mathbf{C}$  and exponential time in the depth of  $\mathbf{C}$  [60, pp. 392–397].

Similarly to the restriction on the circuits for LOGCFL, we call an NDL query  $(\Pi, G)$  *skinny* if the body of any clause in  $\Pi$  has at most two atoms.

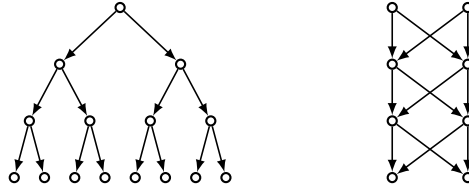
**Lemma 4.** *For any skinny  $(\Pi, G(\mathbf{x}))$  and any data instance  $\mathcal{A}$ , query evaluation can be done by an NAuxPDA in space  $\log |\Pi| + w(\Pi, G) \cdot \log |\mathcal{A}|$  and time  $2^{O(d(\Pi, G))}$ .*

*Proof.* Using the atoms of the grounding  $\Pi_{\mathcal{A}}^a$  as gates and inputs, we define a monotone Boolean circuit  $\mathbf{C}$  as follows: its output is  $G(\mathbf{a})$ ; for every atom  $\gamma$  in the head of a clause in  $\Pi_{\mathcal{A}}^a$ , we take an OR-gate whose output is  $\gamma$  and inputs are the bodies of the clauses with head  $\gamma$ ; for every such body, we take an AND-gate whose inputs are the atoms in the body. We set input  $\gamma$  to 1 iff  $\gamma \in \mathcal{A}$ . Clearly,  $\mathbf{C}$  is a semi-unbounded fan-in circuit of depth  $O(d(\Pi, G))$  with  $O(|\Pi| \cdot |\mathcal{A}|^{w(\Pi, G)})$  gates. Having observed that our  $\mathbf{C}$  can be computed by a deterministic logspace Turing machine, we conclude that the query evaluation problem can be solved by an NAuxPDA in the required space and time.  $\square$

Observe that Lemma 4 holds for NDL queries with any *bounded* number of atoms, not only two. In the rewritings we propose in Sections 3.2 and 3.4, however, the number of atoms in the clauses is not bounded by a constant. We require the following notion to generalise skinny programs. A function  $\nu$  from the predicate names in  $\Pi$  to  $\mathbb{N}$  is called a *weight function* for an NDL query  $(\Pi, G(\mathbf{x}))$  if

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

for any clause  $Q(\mathbf{z}) \leftarrow P_1(\mathbf{z}_1) \wedge \cdots \wedge P_k(\mathbf{z}_k)$  in  $\Pi$ . Note that  $\nu(P)$  can be 0 for an EDB predicate  $P$ . To illustrate, we consider NDL queries with the following dependency graphs:



The NDL on the left has a weight function bounded by the number of predicates, and so, such weight functions are 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. Also observe that linear NDL queries have weight functions bounded by 1.

We show, using Huffman coding, that any NDL query  $(\Pi, G(\mathbf{x}))$  can be transformed into an equivalent skinny NDL query whose depth increases linearly in addition to the logarithms of the weight function and the number  $e_{\Pi}$  of EDB predicates in a clause. We call the minimum (over possible weight functions  $\nu$ ) value of  $2d(\Pi, G) + \log \nu(G) + \log e_{\Pi}$  the *skinny depth* of  $(\Pi, G)$  and denote it by  $\text{sd}(\Pi, G)$ .

**Lemma 5.** *Any NDL query  $(\Pi, G(\mathbf{x}))$  is equivalent to a skinny NDL query  $(\Pi', G(\mathbf{x}))$  such that  $|\Pi'| = O(|\Pi|^2)$ ,  $d(\Pi', G) \leq \text{sd}(\Pi, G)$ , and  $w(\Pi', G) \leq w(\Pi, G)$ .*

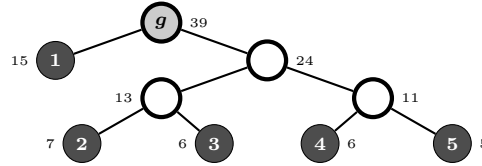


*Proof.* Let  $\nu$  be a weight function such that  $\text{sd}(\Pi, G) = 2\text{d}(\Pi, G) + \log \nu(G) + \log e_\Pi$ . Without loss of generality, we will assume that  $\nu(E) = 0$ , for EDB predicates  $E$ . First, we split clauses into their EDB and IDB components: each  $Q(\mathbf{z}) \leftarrow \varphi(\mathbf{z}')$  is replaced by  $Q(\mathbf{z}) \leftarrow Q_E(\mathbf{z}_E) \wedge Q_I(\mathbf{z}'_I)$  and  $Q_\alpha(\mathbf{z}_\alpha) \leftarrow \varphi_\alpha(\mathbf{z}'_\alpha)$ , for  $\alpha \in \{E, I\}$ , where  $Q_E$  and  $Q_I$  are fresh predicates, and  $\varphi_E(\mathbf{z}'_E)$  and  $\varphi_I(\mathbf{z}'_I)$  are conjunctions of the EDB and IDB predicates in  $\varphi$ , respectively. The depth of the resulting NDL query  $(\Pi_*, G(\mathbf{x}))$  is  $2\text{d}(\Pi, G)$ . Next, each clause  $Q_E(\mathbf{z}_E) \leftarrow \varphi_E(\mathbf{z}'_E)$  in  $\Pi_*$  is replaced by  $\leq e_\Pi - 1$  clauses with at most two atoms in the body, which results in an NDL query of depth not exceeding  $2\text{d}(\Pi, G) + \log e_\Pi$ . In the rest of the proof, we concentrate on the part  $\Pi_\dagger$  of  $\Pi_*$  comprising clauses that have predicates  $Q$  and  $Q_I$  in their heads (thus making the  $Q_E$  EDB predicates). The weight function for  $(\Pi_\dagger, G(\mathbf{x}))$  is obtained by extending  $\nu$  as follows: we set  $\nu(Q_I) = \nu(Q)$  and  $\nu(Q_E) = 0$ , for each  $Q$ .

Next, by induction on  $\text{d}(\Pi_\dagger, G)$ , we show that there is an equivalent skinny NDL query  $(\Pi'_\dagger, G(\mathbf{x}))$  of the required size and width and such that  $\text{d}(\Pi'_\dagger, G) \leq \text{d}(\Pi_\dagger, G) + \log \nu(G)$ . We take  $\Pi'_\dagger = \Pi_\dagger$  if  $\text{d}(\Pi_\dagger, G) = 0$ . Otherwise, let  $\psi$  be a clause of the form  $G(\mathbf{z}) \leftarrow P_1(\mathbf{z}_1) \wedge \dots \wedge P_k(\mathbf{z}_k)$  in  $\Pi_\dagger$ , for  $k > 2$ . Since, by the construction of  $\Pi_\dagger$ , if a clause has an EDB predicate, then  $k = 2$ . So, the  $P_i$  are IDB predicates and  $\nu(G) \geq \nu(P_i) > 0$ . Suppose that, for each  $i$  ( $1 \leq i \leq k$ ), we have an NDL query  $(\Pi'_i, P_i)$  equivalent to  $(\Pi_\dagger, P_i)$  with

$$\text{d}(\Pi'_i, P_i) \leq \text{d}(\Pi_\dagger, P_i) + \log \nu(P_i) \leq \text{d}(\Pi_\dagger, G) - 1 + \log \nu(P_i). \quad (3)$$

Construct the Huffman tree [30] for the alphabet  $\{1, \dots, k\}$ , where the frequency of  $i$  is  $\nu(P_i)/\nu(G)$ . For example, for  $\nu(G) = 39$ ,  $\nu(P_1) = 15$ ,  $\nu(P_2) = 7$ ,  $\nu(P_3) = 6$ ,  $\nu(P_4) = 6$  and  $\nu(P_5) = 5$ , we obtain the following tree:



In general, the Huffman tree is a binary tree with  $k$  leaves  $1, \dots, k$ , a root  $g$  and  $k - 2$  internal nodes and such that the length of the path from  $g$  to any leaf  $i$  is bounded by  $\lceil \log(\nu(G)/\nu(P_i)) \rceil$ . For each internal node  $v$  of the tree, we take a predicate  $P_v(\mathbf{z}_v)$ , where  $\mathbf{z}_v$  is the union of  $\mathbf{z}_u$  for all descendants  $u$  of  $v$ ; for the root  $g$ , we take  $P_g(\mathbf{z}_g) = G(\mathbf{z})$ . Let  $\Pi'_\psi$  be the extension of the union of the  $\Pi'_i$  ( $1 \leq i \leq k$ ) with clauses  $P_v(\mathbf{z}_v) \leftarrow P_{u_1}(\mathbf{z}_{u_1}) \wedge P_{u_2}(\mathbf{z}_{u_2})$ , for each  $v$  with immediate successors  $u_1$  and  $u_2$ . The number of the new clauses is  $k - 1$ . By (3), we have:

$$\begin{aligned} \text{d}(\Pi'_\psi, G) &\leq \max_i \{ \lceil \log(\nu(G)/\nu(P_i)) \rceil + \text{d}(\Pi'_i, P_i) \} \\ &\leq \max_i \{ \log(\nu(G)/\nu(P_i)) + \text{d}(\Pi_\dagger, G) + \log \nu(P_i) \} = \text{d}(\Pi_\dagger, G) + \log \nu(G). \end{aligned}$$

Let  $\Pi'_\dagger$  be the result of applying this transformation to each clause in  $\Pi_\dagger$  with head  $G(\mathbf{z})$  and more than two atoms in the body.

Finally, we add to  $\Pi'_\dagger$  the clauses with the  $Q_E$  predicates and denote the result by  $\Pi'$ . It is readily seen that  $(\Pi', G)$  is as required; in particular,  $|\Pi'| = O(|\Pi|^2)$ .  $\square$

We now use Lemmas 4 and 5 to obtain the following:

**Theorem 6.** *For every  $c > 0$  and  $w > 0$ , evaluation of NDL queries  $(\Pi, G(\mathbf{x}))$  of width at most  $w$  and such that  $\text{sd}(\Pi, G) \leq c \log |\Pi|$  is in LOGCFL for combined complexity.*

We say that a class of OMQs is *skinny-reducible* if, for some fixed  $c > 0$  and  $w > 0$ , there is an  $\mathcal{L}^{\text{LOGCFL}}$ -transducer that, given any OMQ  $Q(\mathbf{x})$  in the class, computes its NDL-rewriting  $(\Pi, G(\mathbf{x}))$  over complete data instances such that  $\text{sd}(\Pi, G) \leq c \log |\Pi|$  and  $w(\Pi, G) \leq w$ . Theorem 6 and the transformation \* give the following:

**Corollary 7.** *For any skinny-reducible class, the OMQ answering problem is in LOGCFL for combined complexity.*

In the following subsections, we will exploit the results obtained above to construct optimal NDL-rewritings for the three classes of tractable OMQs. Appendix A.6 gives concrete examples of our rewritings.

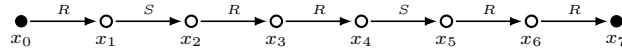
### 3.2 LOGCFL rewritings for $\text{OMQ}(d, t, \infty)$

Recall (see, e.g., [22]) that a *tree decomposition* of an undirected graph  $\mathcal{G} = (V, E)$  is a pair  $(T, \lambda)$ , where  $T$  is an (undirected) tree and  $\lambda$  a function from the nodes of  $T$  to  $2^V$  such that

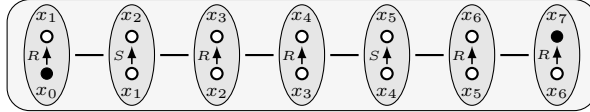
- for every  $v \in V$ , there exists a node  $t$  with  $v \in \lambda(t)$ ;
- for every  $e \in E$ , there exists a node  $t$  with  $e \subseteq \lambda(t)$ ;
- for every  $v \in V$ , the nodes  $\{t \mid v \in \lambda(t)\}$  induce a connected subgraph of  $T$  (called a *subtree* of  $T$ ).

We call the set  $\lambda(t) \subseteq V$  a *bag* for  $t$ . The *width* of  $(T, \lambda)$  is  $\max_{t \in T} |\lambda(t)| - 1$ . The *treewidth* of a graph  $\mathcal{G}$  is the minimum width over all tree decompositions of  $\mathcal{G}$ . The *treewidth* of a CQ is the treewidth of its Gaifman graph.

**Example 8.** Consider the CQ  $q(x_0, x_7)$  depicted below (black nodes represent answer variables):



Its natural tree decomposition of treewidth 1 is based on the chain  $T$  of 7 vertices shown as bags below:



In this section, we prove the following:

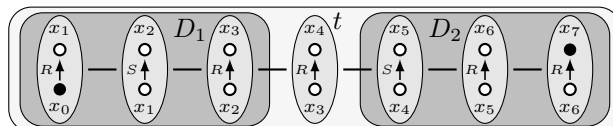
**Theorem 9.** *For any fixed  $d \geq 0$  and  $t \geq 1$ , the class  $\text{OMQ}(d, t, \infty)$  is skinny-reducible.*

In a nutshell, we split recursively a given CQ  $q$  into sub-CQs  $q_D$  based on subtrees  $D$  of the tree decomposition of  $q$ , and combine their rewritings into a rewriting of  $q$ . To guarantee compatibility of these rewritings, we use ‘boundary conditions’  $w$  that describe the types of points on the boundaries of the  $q_D$  and, for each possible boundary condition  $w$ , we define recursively a fresh IDB predicate  $G_D^w$ . We now formalise the construction and illustrate it using the CQ from Example 8.

Fix a connected CQ  $q(x)$  and a tree decomposition  $(T, \lambda)$  of its Gaifman graph  $\mathcal{G} = (V, E)$ . 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  $t' \notin D$ . The *degree*  $\deg(D)$  of  $D$  is the number of its boundary nodes ( $T$  itself is the only subtree of  $T$  of degree 0). 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 10** ([5]). *Let  $D$  be a subtree of  $T$  of size  $n > 1$ . If  $\deg(D) = 2$ , then there is a node  $t$  splitting  $D$  into subtrees of size  $\leq n/2$  and degree  $\leq 2$  and, possibly, one subtree of size  $< n - 1$  and degree 1. If  $\deg(D) \leq 1$ , then there is  $t$  splitting  $D$  into subtrees of size  $\leq n/2$  and degree  $\leq 2$ .*

In Example 8,  $t$  splits  $T$  into  $D_1$  and  $D_2$  as follows:



We define recursively a set  $\mathfrak{D}$  of subtrees of  $T$ , a binary ‘predecessor’ relation  $\prec$  on  $\mathfrak{D}$ , and a function  $\sigma$  on  $\mathfrak{D}$  indicating the splitting node. We begin by adding  $T$  to  $\mathfrak{D}$ . Take any  $D \in \mathfrak{D}$  that has not been split yet. If  $D$  is of size 1, then  $\sigma(D)$  is the only node of  $D$ . Otherwise, by Lemma 10, we find a node  $t$  in  $D$  that splits it into  $D_1, \dots, D_k$ . We set  $\sigma(D) = t$  and, for  $1 \leq i \leq k$ , add  $D_i$  to  $\mathfrak{D}$  and set  $D_i \prec D$ ; then, we apply the procedure recursively to each of  $D_1, \dots, D_k$ . In Example 8 with  $t$  splitting  $T$ , we have  $\sigma(T) = t$ ,  $D_1 \prec T$  and  $D_2 \prec T$ .

For each  $D \in \mathfrak{D}$ , we recursively define a set of atoms

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

By the definition of tree decomposition,  $\mathbf{q}_T = \mathbf{q}$ . Denote by  $\mathbf{x}_D$  the subset of  $\mathbf{x}$  that occurs in  $\mathbf{q}_D$ . In Example 8,  $\mathbf{x}_T = \{x_0, x_7\}$ ,  $\mathbf{x}_{D_1} = \{x_0\}$  and  $\mathbf{x}_{D_2} = \{x_7\}$ . Let  $\partial D$  be the union of all  $\lambda(t) \cap \lambda(t')$  for boundary nodes  $t$  of  $D$  and its neighbours  $t'$  in  $T$  *outside*  $D$ . In our example,  $\partial T = \emptyset$ ,  $\partial D_1 = \{x_3\}$  and  $\partial D_2 = \{x_4\}$ .

Let  $\mathcal{T}$  be an ontology of depth  $\leq \mathbf{d}$ . A *type* is a partial map  $\mathbf{w}$  from  $V$  to  $\mathbf{W}_{\mathcal{T}}$ ; its domain is denoted by  $\text{dom}(\mathbf{w})$ . The unique partial type with  $\text{dom}(\mathbf{w}) = \emptyset$  is denoted by  $\varepsilon$ . We use types to represent how variables are mapped into  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ , with  $\mathbf{w}(z) = w$  indicating that  $z$  is mapped to an element of the form  $aw$  (for some  $a \in \text{ind}(\mathcal{A})$ ), and with  $\mathbf{w}(z) = \varepsilon$  that  $z$  is mapped to an individual constant. We say that a type  $\mathbf{w}$  is *compatible* with a bag  $t$  if, for all  $y, z \in \lambda(t) \cap \text{dom}(\mathbf{w})$ , we have

- if  $z \in \mathbf{x}$ , then  $\mathbf{w}(z) = \varepsilon$ ;
- if  $A(z) \in \mathbf{q}$ , then either  $\mathbf{w}(z) = \varepsilon$  or  $\mathbf{w}(z) = w_\varrho$  with  $\mathcal{T} \models \exists y \varrho(y, x) \rightarrow A(x)$ ;
- if  $P(y, z) \in \mathbf{q}$ , then one of the three conditions holds: (i)  $\mathbf{w}(y) = \mathbf{w}(z) = \varepsilon$ ; (ii)  $\mathbf{w}(y) = \mathbf{w}(z)$  and  $\mathcal{T} \models P(x, x)$ ; (iii)  $\mathcal{T} \models \varrho(x, y) \rightarrow P(x, y)$  and either  $\mathbf{w}(z) = \mathbf{w}(y)\varrho$  or  $\mathbf{w}(y) = \mathbf{w}(z)\varrho^-$ .

In the sequel we abuse notation and use sets of variables in place of sequences assuming that they are ordered in some (fixed) way. For example, we use  $\mathbf{x}_D$  for a tuple of variables in the set  $\mathbf{x}_D$  (ordered in some way). Also, given a tuple  $\mathbf{a} \in \text{ind}(\mathcal{A})^{|\mathbf{x}_D|}$  and  $x \in \mathbf{x}_D$ , we write  $\mathbf{a}(x)$  to refer to the component of  $\mathbf{a}$  that corresponds to  $x$  (that is, the component with the same index).

We now define an NDL-rewriting of  $\mathbf{Q}(\mathbf{x}) = (\mathcal{T}, \mathbf{q}(\mathbf{x}))$ . For any  $D \in \mathfrak{D}$  and type  $\mathbf{w}$  with  $\text{dom}(\mathbf{w}) = \partial D$ , let  $G_D^{\mathbf{w}}(\partial D, \mathbf{x}_D)$  be a fresh IDB predicate with parameters  $\mathbf{x}_D$  (note that  $\partial D$  and  $\mathbf{x}_D$  may be not disjoint). For each type  $\mathbf{s}$  with  $\text{dom}(\mathbf{s}) = \lambda(\sigma(D))$  such that  $\mathbf{s}$  is compatible with  $\sigma(D)$  and agrees with  $\mathbf{w}$  on their common domain, the NDL program  $\Pi_{\mathbf{Q}}^{\text{LOG}}$  contains

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

where  $(\mathbf{s} \cup \mathbf{w}) \upharpoonright \partial D'$  is the restriction of the union  $\mathbf{s} \cup \mathbf{w}$  to  $\partial D'$  (since  $\text{dom}(\mathbf{s} \cup \mathbf{w})$  covers  $\partial D'$ , the domain of the restriction is  $\partial D'$ ), and  $\text{At}^{\mathbf{s}}$  is the conjunction of

- (a)  $A(z)$ , for  $A(z) \in \mathbf{q}$  with  $\mathbf{s}(z) = \varepsilon$ , and  $P(y, z)$ , for  $P(y, z) \in \mathbf{q}$  with  $\mathbf{s}(y) = \mathbf{s}(z) = \varepsilon$ ;
- (b)  $y = z$ , for  $P(y, z) \in \mathbf{q}$  with  $\mathbf{s}(y) \neq \varepsilon$  or  $\mathbf{s}(z) \neq \varepsilon$ ;
- (c)  $A_\varrho(z)$ , for  $z$  with  $\mathbf{s}(z) = \varrho w$ , for some  $w$ .

The conjuncts in (a) ensure that atoms all of whose variables are assigned  $\varepsilon$  hold in the data instance. The conjuncts in (b) ensure that if one variable in a binary atom is not mapped to  $\varepsilon$ , then the images of both its variables share the same initial individual. Finally, the conjuncts in (c) ensure that if a variable is to be mapped to  $aw$ , then  $aw$  is indeed in the domain of  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ .

**Example 11.** With the query in Example 8, consider now the following ontology  $\mathcal{T}$ :

$$\begin{aligned} P(x, y) &\rightarrow S(x, y), & A_P(x) &\leftrightarrow \exists y P(x, y), \\ P(x, y) &\rightarrow R(y, x), & A_{P-}(x) &\leftrightarrow \exists y P(y, x) \end{aligned}$$

(the remaining normalisation axioms are omitted). Since  $\lambda(t) = \{x_3, x_4\}$ , there are two types compatible with  $t$  that can contribute to the rewriting:  $\mathbf{s}_1 = \{x_3 \mapsto \varepsilon, x_4 \mapsto \varepsilon\}$  and  $\mathbf{s}_2 = \{x_3 \mapsto \varepsilon, x_4 \mapsto P^-\}$ . So we have  $\text{At}^{\mathbf{s}_1} = R(x_3, x_4)$  and  $\text{At}^{\mathbf{s}_2} = A_{P^-}(x_4) \wedge (x_3 = x_4)$ . Thus, the predicate  $G_T^\varepsilon$  is defined by two clauses with the head  $G_T^\varepsilon(x_0, x_7)$  and the following bodies:

$$\begin{aligned} & G_{D_1}^{x_3 \mapsto \varepsilon}(x_3, x_0) \wedge R(x_3, x_4) \wedge G_{D_2}^{x_4 \mapsto \varepsilon}(x_4, x_7), \\ & G_{D_1}^{x_3 \mapsto \varepsilon}(x_3, x_0) \wedge A_{P^-}(x_4) \wedge (x_3 = x_4) \wedge G_{D_2}^{x_4 \mapsto P^-}(x_4, x_7), \end{aligned}$$

for  $\mathbf{s}_1$  and  $\mathbf{s}_2$ , respectively. Although  $\{x_3 \mapsto P, x_4 \mapsto \varepsilon\}$  is also compatible with  $t$ , its predicate  $G_{D_1}^{x_3 \mapsto P}$  will have no definition in the rewriting, and hence can be omitted. The same is true of the other compatible types  $\{x_3 \mapsto \varepsilon, x_4 \mapsto R\}$  and  $\{x_3 \mapsto R^-, x_4 \mapsto \varepsilon\}$ .

By induction on  $\prec$ , one can now show that  $(\Pi_Q^{\text{LOG}}, G_T^\varepsilon)$  is a rewriting of  $\mathbf{Q}(\mathbf{x})$ ; see Appendix A.3 for details.

Now fix  $\mathbf{d}$  and  $\mathbf{t}$ , and consider  $\mathbf{Q}(\mathbf{x}) = (\mathcal{T}, \mathbf{q}(\mathbf{x}))$  from  $\text{OMQ}(\mathbf{d}, \mathbf{t}, \infty)$ . Let  $T$  be a tree decomposition of  $\mathbf{q}$  of treewidth  $\leq \mathbf{t}$ ; we may assume without loss of generality that  $T$  has at most  $|\mathbf{q}|$  nodes. We take the following weight function:  $\nu(G_D^\mathbf{w}) = |D|$ , where  $|D|$  is the size of  $D$ , that is, the number of nodes in it. Clearly,  $\nu(G_T^\varepsilon) \leq |\mathbf{Q}|$ . By Lemma 10, we have

$$\begin{aligned} \mathbf{w}(\Pi_Q^{\text{LOG}}, G_T^\varepsilon) &\leq \max_D |\partial D \cup \lambda(\sigma(D))| \leq 3(\mathbf{t} + 1), \\ \text{sd}(\Pi_Q^{\text{LOG}}, G_T^\varepsilon) &\leq 4 \log |T| + 2 \log |\mathbf{Q}| \leq 6 \log |\mathbf{Q}|. \end{aligned}$$

Since  $|\mathcal{D}| \leq |T|^2$  and there are at most  $|\mathcal{T}|^{2d(\mathbf{t}+1)}$  options for  $\mathbf{w}$ , there are polynomially many predicates  $G_D^\mathbf{w}$ , and so  $\Pi_Q^{\text{LOG}}$  is of polynomial size. Thus, by Corollary 7, the constructed NDL-rewriting over arbitrary data instances can be evaluated in LOGCFL. Finally, we note that a tree decomposition of treewidth  $\leq \mathbf{t}$  can be computed using an  $\text{L}^{\text{LOGCFL}}$ -transducer [26], and so the NDL-rewriting can also be constructed by an  $\text{L}^{\text{LOGCFL}}$ -transducer.

The obtained NDL-rewriting shows that answering OMQs  $(\mathcal{T}, \mathbf{q}(\mathbf{x}))$  with  $\mathcal{T}$  of finite depth  $\mathbf{d}$  and  $\mathbf{q}$  of treewidth  $\mathbf{t}$  over any data instance  $\mathcal{A}$  can be done in time

$$\text{poly}(|\mathcal{T}|^{d\mathbf{t}}, |\mathbf{q}|, |\mathcal{A}|^{\mathbf{t}}). \quad (4)$$

Indeed, we can evaluate  $(\Pi_Q^{\text{LOG}}, G_T^\varepsilon(\mathbf{x}))$  in time polynomial in  $|\Pi_Q^{\text{LOG}}|$  and  $|\mathcal{A}|^{\mathbf{w}(\Pi_Q^{\text{LOG}}, G_T^\varepsilon)}$ , which are bounded by a polynomial in  $|\mathcal{T}|^{2d(\mathbf{t}+1)}$ ,  $|\mathbf{q}|$  and  $|\mathcal{A}|^{2(\mathbf{t}+1)}$ .

### 3.3 NL rewritings for $\text{OMQ}(\mathbf{d}, 1, \ell)$

**Theorem 12.** *Let  $\mathbf{d} \geq 0$  and  $\ell \geq 2$  be fixed. There is an  $\text{L}^{\text{NL}}$ -transducer that, given an OMQ  $\mathbf{Q} = (\mathcal{T}, \mathbf{q}(\mathbf{x}))$  in  $\text{OMQ}(\mathbf{d}, 1, \ell)$ , constructs its polynomial-size linear NDL-rewriting of width  $\leq 2\ell$ .*

Let  $\mathcal{T}$  be an ontology of finite depth  $\mathbf{d}$ , and let  $\mathbf{q}(\mathbf{x})$  be a tree-shaped CQ with at most  $\ell$  leaves. Fix one of the variables of  $\mathbf{q}$  as root, and let  $M$  be the maximal distance to a leaf from the root. For  $0 \leq n \leq M$ , let  $\mathbf{z}^n$  denote the set of all variables of  $\mathbf{q}$  at distance  $n$  from the root; clearly,  $|\mathbf{z}^n| \leq \ell$ . We call the  $\mathbf{z}^n$  *slices* of  $\mathbf{q}$  and observe that they satisfy the following: for every  $P(z, z') \in \mathbf{q}$  with  $z \neq z'$ , there exists  $n < M$  such that

$$\text{either } z \in \mathbf{z}^n \text{ and } z' \in \mathbf{z}^{n+1} \quad \text{or} \quad z' \in \mathbf{z}^n \text{ and } z \in \mathbf{z}^{n+1}.$$

For  $0 \leq n \leq M$ , let  $\mathbf{q}_n(\mathbf{z}^n, \mathbf{x}^n)$  be the query consisting of all atoms  $S(\mathbf{z})$  of  $\mathbf{q}$  such that  $\mathbf{z} \subseteq \bigcup_{n \leq k \leq M} \mathbf{z}^k$ , where  $\mathbf{x}^n$  is the subset of  $\mathbf{x}$  that occurs in  $\mathbf{q}_n$  and  $\mathbf{z}^n = \mathbf{z}^n \setminus \mathbf{x}$ .

By a *type for slice  $\mathbf{z}^n$* , we mean a total map  $\mathbf{w}$  from  $\mathbf{z}^n$  to  $\mathbf{W}_\mathcal{T}$ . Analogously to Section 3.2, we define the notions of types compatible with slices. Specifically, we call  $\mathbf{w}$  *locally compatible* with  $\mathbf{z}^n$  if for every  $z \in \mathbf{z}^n$ :

- if  $z \in \mathbf{x}$ , then  $\mathbf{w}(z) = \varepsilon$ ;

- if  $A(z) \in \mathbf{q}$ , then either  $\mathbf{w}(z) = \varepsilon$  or  $\mathbf{w}(z) = w\rho$  with  $\mathcal{T} \models \exists y \rho(y, x) \rightarrow A(x)$ ;
- if  $P(z, z) \in \mathbf{q}$ , then either  $\mathbf{w}(z) = \varepsilon$  or  $\mathcal{T} \models P(x, x)$ .

If  $\mathbf{w}, \mathbf{s}$  are types for  $\mathbf{z}^n$  and  $\mathbf{z}^{n+1}$ , respectively, then we say  $(\mathbf{w}, \mathbf{s})$  is *compatible* with  $(\mathbf{z}^n, \mathbf{z}^{n+1})$  if  $\mathbf{w}$  is locally compatible with  $\mathbf{z}^n$ ,  $\mathbf{s}$  is locally compatible with  $\mathbf{z}^{n+1}$ ,

- for every  $P(z, z') \in \mathbf{q}$  with  $z \in \mathbf{z}^n$  and  $z' \in \mathbf{z}^{n+1}$ , one of the three condition holds:  $\mathbf{w}(z) = \mathbf{s}(z') = \varepsilon$ , or  $\mathbf{w}(z) = \mathbf{s}(z')$  with  $\mathcal{T} \models P(x, x)$ , or  $\mathcal{T} \models \rho(x, y) \rightarrow P(x, y)$  with either  $\mathbf{s}(z') = \mathbf{w}(z)\rho$  or  $\mathbf{w}(z) = \mathbf{s}(z')\rho^-$ .

Consider the NDL program  $\Pi_Q^{\text{LIN}}$  defined as follows. For every  $0 \leq n < M$  and every pair of types  $(\mathbf{w}, \mathbf{s})$  that is compatible with  $(\mathbf{z}^n, \mathbf{z}^{n+1})$ , we include the clause

$$G_n^{\mathbf{w}}(\mathbf{z}_{\exists}^n, \mathbf{x}^n) \leftarrow \text{At}^{\mathbf{w} \cup \mathbf{s}}(\mathbf{z}^n, \mathbf{z}^{n+1}) \wedge G_{n+1}^{\mathbf{s}}(\mathbf{z}_{\exists}^{n+1}, \mathbf{x}^{n+1}),$$

where  $\mathbf{x}^n$  are the parameters of  $G_n^{\mathbf{w}}$  and  $\text{At}^{\mathbf{w} \cup \mathbf{s}}(\mathbf{z}^n, \mathbf{z}^{n+1})$  is the conjunction of atoms (a)–(c) as defined in Section 3.2, for the union  $\mathbf{w} \cup \mathbf{s}$ . For every type  $\mathbf{w}$  locally compatible with  $\mathbf{z}^M$ , we include the clause

$$G_M^{\mathbf{w}}(\mathbf{z}_{\exists}^M, \mathbf{x}^M) \leftarrow \text{At}^{\mathbf{w}}(\mathbf{z}^M).$$

(Recall that  $\mathbf{z}^M$  is a disjoint union of  $\mathbf{z}_{\exists}^M$  and  $\mathbf{x}^M$ .) We use  $G$  with parameters  $\mathbf{x}$  as the goal predicate and include  $G(\mathbf{x}) \leftarrow G_0^{\mathbf{w}}(\mathbf{z}_{\exists}^0, \mathbf{x})$  for every predicate  $G_0^{\mathbf{w}}$  occurring in the head of one of the preceding clauses.

By induction on  $n$ , we show in Appendix A.4 that  $(\Pi_Q^{\text{LIN}}, G(\mathbf{x}))$  is a rewriting of  $(\mathcal{T}, \mathbf{q}(\mathbf{x}))$  over complete data instances. It should be clear that  $\Pi_Q^{\text{LIN}}$  is a linear NDL program of width  $\leq 2\ell$  and containing  $\leq |\mathbf{q}| \cdot |\mathcal{T}|^{2d\ell}$  predicates. Moreover, it takes only logarithmic space to store a type  $\mathbf{w}$ , which allows us to show that  $\Pi_Q^{\text{LIN}}$  can be computed by an  $\text{L}^{\text{NL}}$ -transducer. We apply Lemma 3 to obtain an NDL-rewriting for arbitrary data instances, and then use Theorem 2 to conclude that the resulting program can be evaluated in NL.

The obtained NDL-rewriting shows that answering OMQs  $(\mathcal{T}, \mathbf{q}(\mathbf{x}))$  with  $\mathcal{T}$  of finite depth  $d$  and tree-shaped  $\mathbf{q}$  with  $\ell$  leaves over any data  $\mathcal{A}$  can be done in time

$$\text{poly}(|\mathcal{T}|^{d\ell}, |\mathbf{q}|, |\mathcal{A}|^{\ell}). \quad (5)$$

Indeed,  $(\Pi_Q^{\text{LIN}}, G(\mathbf{x}))$  can be evaluated in time polynomial in  $|\Pi_Q^{\text{LIN}}|$  and  $|\mathcal{A}|^{w(\Pi_Q^{\text{LIN}}, G)}$ , which are bounded by a polynomial in  $|\mathcal{T}|^{2d\ell}$ ,  $|\mathbf{q}|$  and  $|\mathcal{A}|^{2\ell}$ .

### 3.4 LOGCFL rewritings for OMQ( $\infty, 1, \ell$ )

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

**Theorem 13.** *For any fixed  $\ell \geq 2$ , OMQ( $\infty, 1, \ell$ ) is skinny-reducible.*

For OMQs with bounded-leaf CQs and ontologies of unbounded depth, our rewriting uses the notion of tree witness [37]. Consider an OMQ  $\mathbf{Q}(\mathbf{x}) = (\mathcal{T}, \mathbf{q}(\mathbf{x}))$ . Let  $\mathbf{t} = (\mathbf{t}_r, \mathbf{t}_i)$  be a pair of disjoint sets of variables in  $\mathbf{q}$  such that  $\mathbf{t}_i \neq \emptyset$  but  $\mathbf{t}_i \cap \mathbf{x} = \emptyset$ . Set

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

If  $\mathbf{q}_{\mathbf{t}}$  is a minimal subset of  $\mathbf{q}$  containing every atom of  $\mathbf{q}$  with a variable from  $\mathbf{t}_i$  and such that there is a homomorphism  $h: \mathbf{q}_{\mathbf{t}} \rightarrow \mathcal{C}_{\mathcal{T}, \{A_{\rho}(a)\}}$  with  $h^{-1}(a) = \mathbf{t}_r$ , we call  $\mathbf{t}$  a *tree witness for  $\mathbf{Q}(\mathbf{x})$  generated by  $\rho$* . Intuitively,  $\mathbf{t}$  identifies a minimal subset of  $\mathbf{q}$  that can be mapped to the tree-shaped part of the canonical model consisting of labelled nulls: the variables in  $\mathbf{t}_r$  are mapped to an individual constant, say,  $a$ , at the root of a tree and the  $\mathbf{t}_i$  are mapped to the labelled nulls of the form  $aw$ , for some  $w \in \mathbf{W}_{\mathcal{T}}$  that begins with  $\rho$ . Note that the same tree witness can be generated by different  $\rho$ .

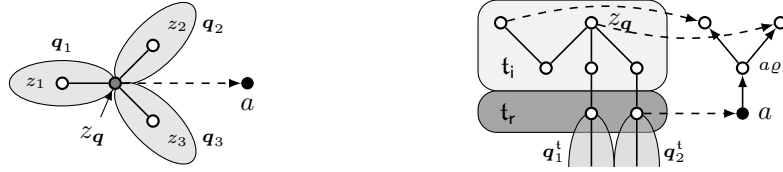
The logarithmic-depth NDL-rewriting for OMQs from OMQ( $\infty, 1, \ell$ ) is based on the following observation:

**Lemma 14** ([34]). *Every tree  $T$  of size  $n$  has a node splitting it into subtrees of size  $\leq \lceil n/2 \rceil$ .*

Let  $\mathbf{Q}(\mathbf{x}_0) = (\mathcal{T}, \mathbf{q}_0(\mathbf{x}_0))$  be an OMQ with a tree-shaped CQ. We will repeatedly apply Lemma 14 to decompose the CQ into smaller and smaller subqueries. Formally, for a tree-shaped CQ  $\mathbf{q}$ , we denote by  $z_{\mathbf{q}}$  a vertex in the Gaifman graph  $\mathcal{G}$  of  $\mathbf{q}$  that satisfies the condition of Lemma 14; if  $|\text{var}(\mathbf{q})| = 2$  and  $\mathbf{q}$  has at least one existentially quantified variable, then we assume that  $z_{\mathbf{q}}$  is such. Let  $\Omega$  be the smallest set that contains  $\mathbf{q}_0(\mathbf{x}_0)$  and the following CQs, for every  $\mathbf{q}(\mathbf{x}) \in \Omega$  with existentially quantified variables:

- for each  $z_i$  adjacent to  $z_{\mathbf{q}}$  in  $\mathcal{G}$ , the CQ  $\mathbf{q}_i(\mathbf{x}_i)$  comprising all binary atoms with both  $z_i$  and  $z_{\mathbf{q}}$ , and all atoms whose variables cannot reach  $z_{\mathbf{q}}$  in  $\mathcal{G}$  without passing by  $z_i$ , where  $\mathbf{x}_i$  is the set of variables in  $\mathbf{x} \cup \{z_{\mathbf{q}}\}$  that occur in  $\mathbf{q}_i$ ;
- for each tree witness  $\mathbf{t}$  for  $(\mathcal{T}, \mathbf{q}(\mathbf{x}))$  with  $\mathbf{t}_r \neq \emptyset$  and  $z_{\mathbf{q}} \in \mathbf{t}_i$ , the CQs  $\mathbf{q}_1^{\mathbf{t}}(\mathbf{x}_1^{\mathbf{t}}), \dots, \mathbf{q}_k^{\mathbf{t}}(\mathbf{x}_k^{\mathbf{t}})$  that correspond to the connected components of the set of atoms of  $\mathbf{q}$  that are not in  $\mathbf{q}_{\mathbf{t}}$ , where each  $\mathbf{x}_i^{\mathbf{t}}$  is the set of variables in  $\mathbf{x} \cup \mathbf{t}_r$  that occur in  $\mathbf{q}_i^{\mathbf{t}}$ .

The two cases are depicted below:



Note that  $\mathbf{t}_r \neq \emptyset$  ensures that part of the query without  $\mathbf{q}_{\mathbf{t}}$  is mapped onto individual constants.

The NDL program  $\Pi_{\mathbf{Q}}^{\text{Tw}}$  uses IDB predicates  $G_{\mathbf{q}}(\mathbf{x})$ , for  $\mathbf{q}(\mathbf{x}) \in \Omega$ , whose parameters are the variables in  $\mathbf{x}_0$  that occur in  $\mathbf{q}(\mathbf{x})$ . For each  $\mathbf{q}(\mathbf{x}) \in \Omega$ , if it has no existentially quantified variables, then we include the clause  $G_{\mathbf{q}}(\mathbf{x}) \leftarrow \mathbf{q}(\mathbf{x})$ . Otherwise, we include the clause

$$G_{\mathbf{q}}(\mathbf{x}) \leftarrow \bigwedge_{S(z) \in \mathbf{q}, z \subseteq \{z_{\mathbf{q}}\}} S(z) \wedge \bigwedge_{1 \leq i \leq n} G_{\mathbf{q}_i}(\mathbf{x}_i),$$

where  $\mathbf{q}_1(\mathbf{x}_1), \dots, \mathbf{q}_n(\mathbf{x}_n)$  are the subqueries induced by the neighbours of  $z_{\mathbf{q}}$  in  $\mathcal{G}$ , and, for each tree witness  $\mathbf{t}$  for  $(\mathcal{T}, \mathbf{q}(\mathbf{x}))$  with  $\mathbf{t}_r \neq \emptyset$  and  $z_{\mathbf{q}} \in \mathbf{t}_i$  and for every  $\varrho$  generating  $\mathbf{t}$ , the following clause

$$G_{\mathbf{q}}(\mathbf{x}) \leftarrow A_{\varrho}(z_0) \wedge \bigwedge_{z \in \mathbf{t}_r \setminus \{z_0\}} (z = z_0) \wedge \bigwedge_{1 \leq i \leq k} G_{\mathbf{q}_i^{\mathbf{t}}}(\mathbf{x}_i^{\mathbf{t}}),$$

where  $z_0$  is any variable in  $\mathbf{t}_r$  and  $\mathbf{q}_1^{\mathbf{t}}, \dots, \mathbf{q}_k^{\mathbf{t}}$  are the connected components of  $\mathbf{q}$  without  $\mathbf{q}_{\mathbf{t}}$ . Finally, if  $\mathbf{q}_0$  is Boolean, then we include clauses  $G_{\mathbf{q}_0} \leftarrow A(x)$  for all unary predicates  $A$  such that  $\mathcal{T}, \{A(a)\} \models \mathbf{q}_0$ .

The program  $\Pi_{\mathbf{Q}}^{\text{Tw}}$  is inspired by a similar construction from [34]. By adapting the proof, we can show that  $(\Pi_{\mathbf{Q}}^{\text{Tw}}, G_{\mathbf{q}_0}(\mathbf{x}_0))$  is indeed a rewriting; see Appendix A.5.

Now fix  $\ell > 1$  and consider  $\mathbf{Q}(\mathbf{x}) = (\mathcal{T}, \mathbf{q}_0(\mathbf{x}))$  from the class  $\text{OMQ}(\infty, 1, \ell)$ . The size of the program  $\Pi_{\mathbf{Q}}^{\text{Tw}}$  is polynomially bounded in  $|\mathbf{Q}|$  since  $\mathbf{q}_0$  has  $O(|\mathbf{q}_0|^{\ell})$  tree witnesses and tree-shaped subqueries. It is readily seen that the function  $\nu$  defined by setting  $\nu(G_{\mathbf{q}}) = |\mathbf{q}|$ , for each  $\mathbf{q} \in \Omega$ , is a weight function for  $(\Pi_{\mathbf{Q}}^{\text{Tw}}, G_{\mathbf{q}_0}(\mathbf{x}))$  with  $\nu(G_{\mathbf{q}_0}) \leq |\mathbf{Q}|$ . Moreover, by Lemma 14,  $d(\Pi_{\mathbf{Q}}^{\text{Tw}}, G_{\mathbf{q}_0}) \leq \log \nu(G_{\mathbf{q}_0}) + 1$ ; and clearly,  $w(\Pi_{\mathbf{Q}}^{\text{Tw}}, G_{\mathbf{q}_0}) \leq \ell + 1$ . By Corollary 7, the obtained NDL-rewritings can be evaluated in LOGCFL. 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 14, and in LOGCFL to decide whether  $\mathcal{T}, \{A(a)\} \models \mathbf{q}_0$  [5] or whether a (logspace) representation of a possible tree witness is indeed a tree witness. This allows us to show that  $(\Pi_{\mathbf{Q}}^{\text{Tw}}, G_{\mathbf{q}_0}(\mathbf{x}))$  can be generated by an  $L^{\text{LOGCFL}}$ -transducer.

It also follows that answering OMQs  $(\mathcal{T}, \mathbf{q}(\mathbf{x}))$  with a tree-shaped CQ with  $\ell$  leaves over any data instance  $\mathcal{A}$  can be done in time

$$\text{poly}(|\mathcal{T}|, |\mathbf{q}|^{\ell}, |\mathcal{A}|^{\ell}). \quad (6)$$

Indeed,  $(\Pi_Q^{\text{Tw}}, G(\mathbf{x}))$  can be evaluated in time polynomial in  $|\Pi_Q^{\text{Tw}}|$  and  $|\mathcal{A}|^{w(\Pi_Q^{\text{Tw}}, G)}$ , which are bounded by polynomials in  $|\mathcal{T}|$ ,  $|\mathbf{q}|^\ell$  and  $|\mathcal{A}|^\ell$ , respectively.

## 4 Parameterised complexity

The upper bounds (4) and (6) for the time required to evaluate NDL-rewritings of OMQs from  $\text{OMQ}(\mathbf{d}, 1, \infty)$  and  $\text{OMQ}(\infty, 1, \ell)$  contain  $\mathbf{d}$  and  $\ell$  in the exponent of  $|\mathcal{T}|$  and  $|\mathbf{q}|$ . Moreover, if we allow  $\mathbf{d}$  and  $\ell$  to grow while keeping CQs tree-shaped, the combined complexity of OMQ answering will jump to NP; see Fig. 1(a). In this section, we regard  $\mathbf{d}$  and  $\ell$  as parameters and show that answering tree-shaped OMQs is not fixed-parameter tractable.

### 4.1 Ontology Depth

Consider the following problem *pDepth-TREEOMQ*:

**Instance:** an OMQ  $Q = (\mathcal{T}, \mathbf{q})$  with  $\mathcal{T}$  of finite depth and tree-shaped Boolean CQ  $\mathbf{q}$ .

**Parameter:** the depth of  $\mathcal{T}$ .

**Problem:** decide whether  $\mathcal{T}, \{A(a)\} \models \mathbf{q}$ .

**Theorem 15.** *pDepth-TREEOMQ is  $W[2]$ -hard.*

*Proof.* The proof is by reduction of the problem *p-HITTINGSET*, which is known to be  $W[2]$ -complete [22]:

**Instance:** a hypergraph  $H = (V, E)$  and  $k \in \mathbb{N}$ .

**Parameter:**  $k$ .

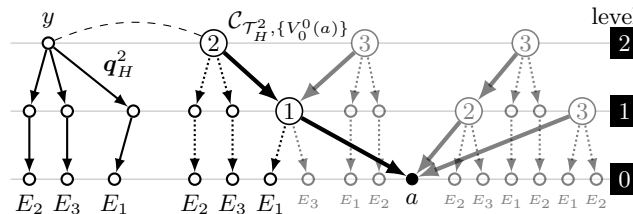
**Problem:** decide whether there is  $A \subseteq V$  such that  $|A| = k$  and  $e \cap A \neq \emptyset$ , for every  $e \in E$ . (Such a set  $A$  of vertices is called a *hitting set of size  $k$* .) Suppose that  $H = (V, E)$  is a hypergraph with vertices  $V = \{v_1, \dots, v_n\}$  and hyperedges  $E = \{e_1, \dots, e_m\}$ . Let  $\mathcal{T}_H^k$  be the (normal form of an) ontology with the following axioms, for  $1 \leq l \leq k$ :

$$\begin{aligned} V_i^{l-1}(x) &\rightarrow \exists z (P(z, x) \wedge V_{i'}^l(z)), & \text{for } 0 \leq i < i' \leq n, \\ V_i^l(x) &\rightarrow E_j^l(x), & \text{for } v_i \in e_j, e_j \in E, \\ E_j^l(x) &\rightarrow \exists z (P(x, z) \wedge E_j^{l-1}(z)), & \text{for } 1 \leq j \leq m. \end{aligned}$$

Let  $\mathbf{q}_H^k$  be a tree-shaped Boolean CQ with the following atoms, for  $1 \leq j \leq m$ :

$$P(y, z_j^{k-1}), \quad P(z_j^l, z_j^{l-1}) \text{ for } 1 \leq l < k, \quad \text{and } E_j^0(z_j^0).$$

The first axiom of  $\mathcal{T}_H^k$  generates a tree of depth  $k$ , with branching ranging from  $n$  to 1, such that the points  $w$  of level  $k$  are labelled with subsets  $X \subseteq V$  of size  $k$  that are read off the path from the root to  $w$ . The CQ  $\mathbf{q}_H^k$  is a star with rays corresponding to the hyperedges of  $H$ . The second and third axioms generate ‘pendants’ ensuring that, for any hyperedge  $e$ , the central point of the CQ can be mapped to a point with a label  $X$  iff  $X$  and  $e$  have a common vertex. The canonical model of  $(\mathcal{T}_H^k, \{V_0^0(a)\})$  and the CQ  $\mathbf{q}_H^k$ , for  $H = (V, \{e_1, e_2, e_3\})$  with  $V = \{1, 2, 3\}$ ,  $e_1 = \{1, 3\}$ ,  $e_2 = \{2, 3\}$  and  $e_3 = \{1, 2\}$ , is shown below:



Points  $\textcircled{i}$  at level  $l$  belong to  $V_i^l$ . In Appendix B.1 we prove that  $\mathcal{T}_H^k, \{V_0^0(a)\} \models \mathbf{q}_H^k$  iff  $H$  has a hitting set of size  $k$ . In the example above,  $\{1, 2\}$  is a hitting set of size 2, which corresponds to the homomorphism from  $\mathbf{q}_H^2$  into the part of  $\mathcal{T}_H^2, \{V_0^0(a)\}$  shown in black.  $\square$

By Theorem 9, OMQs  $(\mathcal{T}, \mathbf{q})$  from  $\text{OMQ}(\mathbf{d}, 1, \infty)$  can be answered (via NDL-rewriting) over a data instance  $\mathcal{A}$  in time  $\text{poly}(|\mathcal{T}|^{\mathbf{d}}, |\mathbf{q}|, |\mathcal{A}|)$ . Theorem 15 shows that no algorithm can do this in time  $f(\mathbf{d}) \cdot \text{poly}(|\mathcal{T}|, |\mathbf{q}|, |\mathcal{A}|)$ , for any computable function  $f$ , unless  $W[2] = \text{FPT}$ .

## 4.2 Number of Leaves

Next we consider the problem  $p\text{Leaves-TREEOMQ}$ :

**Instance:** an OMQ  $\mathbf{Q} = (\mathcal{T}, \mathbf{q})$  with  $\mathcal{T}$  of finite depth and tree-shaped Boolean CQ  $\mathbf{q}$ .

**Parameter:** the number of leaves in  $\mathbf{q}$ .

**Problem:** decide whether  $\mathcal{T}, \{A(a)\} \models \mathbf{q}$ .

**Theorem 16.**  $p\text{Leaves-TREEOMQ}$  is  $W[1]$ -hard.

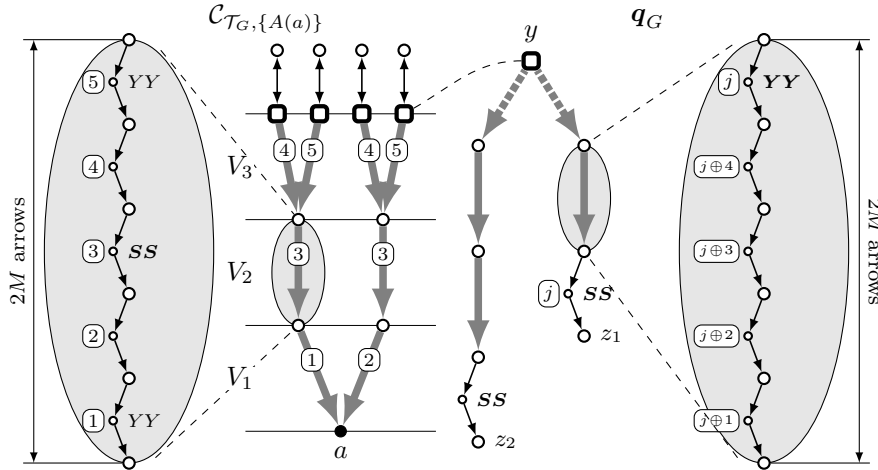
*Proof.* The proof is by reduction of the following  $W[1]$ -complete PARTITIONEDCLIQUE problem [21]:

**Instance:** a graph  $G = (V, E)$  whose vertices are partitioned into  $p$  sets  $V_1, \dots, V_p$ .

**Parameter:**  $p$ , the number of partitions.

**Problem:** decide whether  $G$  has a clique of size  $p$  containing one vertex from each  $V_i$ .

Consider a graph  $G = (V, E)$  with  $V = \{v_1, \dots, v_M\}$  partitioned into  $V_1, \dots, V_p$ . The ontology  $\mathcal{T}_G$  will create a tree rooted at  $A(a)$  whose every branch corresponds to selecting one vertex from each  $V_i$ . Each branch has length  $(p \cdot 2M) + 1$  and consists of  $p$  ‘blocks’ of length  $2M$ , plus an extra edge at the end (used for padding). Each block corresponds to an enumeration of  $V$ , with positions  $2j$  and  $2j + 1$  being associated with  $v_j$ . In the  $i$ th block of a branch, we will select a vertex  $v_{j_i}$  from  $V_i$  by marking the positions  $2j_i$  and  $2j_i + 1$  with the binary predicate  $S$ ; we also mark the positions of the neighbours of  $v_{j_i}$  in  $G$  with the predicate  $Y$ . We use the unary predicate  $B$  to mark the end of the  $p$ th block (square nodes in the picture below). The left side of the picture illustrates the construction for  $p = 3$ , where  $V_1 = \{v_1, v_2\}$ ,  $V_2 = \{v_3\}$ ,  $V_3 = \{v_4, v_5\}$ , and  $E = \{\{v_1, v_3\}, \{v_3, v_5\}\}$ .



Since vertices are enumerated in the same order in every block, to check whether the selected vertex  $v_{j_i}$  for  $V_i$  is a neighbour of the vertices selected from  $V_{i+1}, \dots, V_p$ , it suffices to check that positions  $2j_i$  and  $2j_i + 1$  in blocks  $i + 1, \dots, p$  are marked  $YY$ . Moreover, the distance between the positions of a vertex in consecutive blocks is always  $2M - 2$ . The idea is thus to construct a CQ  $\mathbf{q}_G$  (right side of the picture) which, starting from a variable labelled  $B$  (mapped to the end of a  $p$ th block), splits into  $p - 1$  branches, with the  $i$ th branch checking for a sequence of  $i$  evenly-spaced  $YY$  markers leading to an  $SS$  marker. The distance from the end of the  $p$ th block (marked  $B$ ) to the positions  $2j_i$  and  $2j_i + 1$  in the  $p$ th block (where the first  $YY$  should occur) depends on the choice of  $v_{j_i}$ . We thus add an outgoing edge at the end of the  $p$ th block, which can be navigated in both directions, to be able to ‘consume’ any *even number* of query atoms preceding the first  $YY$ .



The Boolean CQ  $\mathbf{q}_G$  looks as follows (for readability, we use atoms with star-free regular expressions):

$$B(y) \wedge \bigwedge_{1 \leq i < p} (U^{2M-2} \cdot (YY \cdot U^{2M-2})^i \cdot SS)(y, z_i),$$

and the ontology  $\mathcal{T}_G$  contains the following axioms:

$$\begin{aligned} A(x) &\rightarrow \exists y L_j^1(x, y), & \text{for } v_j \in V_1, \\ \exists z L_j^k(z, x) &\rightarrow \exists y L_j^{k+1}(x, y), & \text{for } 1 \leq k < 2M, v_j \in V, \\ \exists z L_j^{2M}(z, x) &\rightarrow \exists y L_{j'}^1(x, y), & \text{for } v_j \in V_i, v_{j'} \in V_{i+1}, \\ L_j^k(x, y) &\rightarrow S(y, x), & \text{for } k \in \{2j, 2j+1\}, \\ L_j^k(x, y) &\rightarrow Y(y, x), & \text{for } \{v_j, v_{j'}\} \in E \text{ and } k \in \{2j', 2j'+1\}, \\ L_j^k(x, y) &\rightarrow U(y, x), & \text{for } 1 \leq k \leq 2M, v_j \in V, \\ \exists z L_j^{2M}(z, x) &\rightarrow B(x), & \text{for } v_j \in V_p, \\ B(x) &\rightarrow \exists y (U(x, y) \wedge U(y, x)). \end{aligned}$$

We prove in the appendix that  $\mathcal{T}_G, \{A(a)\} \models \mathbf{q}_G$  iff  $G$  has a clique containing one vertex from each set  $V_i$ .  $\square$

By (6), OMQs  $(\mathcal{T}, \mathbf{q})$  from  $\text{OMQ}(\infty, 1, \ell)$  can be answered (via NDL-rewriting) over a data instance  $\mathcal{A}$  in time  $\text{poly}(|\mathcal{T}|, |\mathbf{q}|^\ell, |\mathcal{A}|^\ell)$ . Theorem 16 shows that no algorithm can do this in time  $f(\ell) \cdot \text{poly}(|\mathcal{T}|, |\mathbf{q}|, |\mathcal{A}|)$ , for any computable function  $f$ , unless  $W[1] = \text{FPT}$ .

One may consider various other types of parameters that can hopefully reduce the complexity of OMQ answering. Obvious candidates are the size of ontology, the size of ontology signature or the number of role inclusions in ontologies. (Indeed, it is shown in [6] that in the absence of role inclusions, tree-shaped OMQ answering is tractable.) Unfortunately, bounding any of these parameters does not make OMQ answering easier, as we establish in Section 5 that already one *fixed* ontology makes the problem NP-hard for tree-shaped CQs and LOGCFL-hard for linear ones.

## 5 OMQs with a Fixed Ontology

In a typical OBDA scenario [33], users are provided with an ontology in a familiar signature (developed by a domain expert) with which they formulate their queries. Thus, it is of interest to identify the complexity of answering tree-shaped OMQs  $(\mathcal{T}, \mathbf{q})$  with a fixed  $\mathcal{T}$  of infinite depth (see Fig. 1). Surprisingly, we show that the problem is NP-hard even when both  $\mathcal{T}$  and  $\mathcal{A}$  are fixed (in the database setting, answering tree-shaped CQs is in LOGCFL for combined complexity).

**Theorem 17.** *There is an ontology  $\mathcal{T}_\dagger$  such that answering OMQs of the form  $(\mathcal{T}_\dagger, \mathbf{q})$  with Boolean tree-shaped CQs  $\mathbf{q}$  is NP-hard for query complexity.*

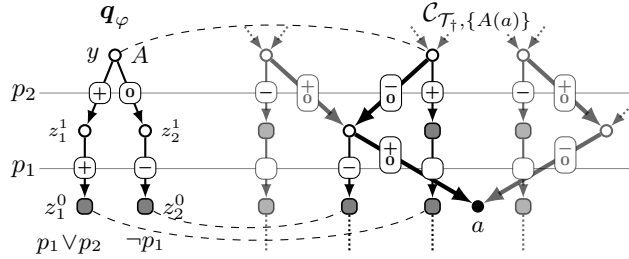
*Proof.* The proof is by reduction of SAT. Given a CNF  $\varphi$  with variables  $p_1, \dots, p_k$  and clauses  $\chi_1, \dots, \chi_m$ , take a Boolean CQ  $\mathbf{q}_\varphi$  with  $A(y)$  and, for  $1 \leq j \leq m$ , the following atoms with  $z_j^k = y$ :

$$\begin{aligned} P_+(z_j^l, z_j^{l-1}), & \quad \text{if } p_l \text{ occurs in } \chi_j \text{ positively,} \\ P_-(z_j^l, z_j^{l-1}), & \quad \text{if } p_l \text{ occurs in } \chi_j \text{ negatively,} \\ P_0(z_j^l, z_j^{l-1}), & \quad \text{if } p_l \text{ does not occur in } \chi_j, \\ B_0(z_j^0). & \end{aligned}$$

Thus,  $\mathbf{q}_\varphi$  is a star with centre  $A(y)$  and  $m$  rays encoding the  $\chi_j$  by the binary predicates  $P_+$ ,  $P_-$  and  $P_0$ . Let  $\mathcal{T}_\dagger$  be an ontology with the axioms

$$\begin{aligned} A(x) &\rightarrow \exists y (P_+(y, x) \wedge P_0(y, x) \wedge B_-(y) \wedge A(y)), \\ B_-(y) &\rightarrow \exists x' (P_-(y, x') \wedge B_0(x')), \\ A(x) &\rightarrow \exists y (P_-(y, x) \wedge P_0(y, x) \wedge B_+(y) \wedge A(y)), \\ B_+(y) &\rightarrow \exists x' (P_+(y, x') \wedge B_0(x')), \\ B_0(x) &\rightarrow \exists y (P_+(x, y) \wedge P_-(x, y) \wedge P_0(x, y) \wedge B_0(y)). \end{aligned}$$

Intuitively,  $(\mathcal{T}_\dagger, \{A(a)\})$  generates an infinite binary tree whose nodes of depth  $n$  represent all  $2^n$  truth assignments to  $n$  propositional variables. The CQ  $\mathbf{q}_\varphi$  can only be mapped along a branch of this tree towards its root  $a$ , with the image of  $y$ , the centre of the star, giving a satisfying assignment for  $\varphi$ . Each non-root node of the tree also starts an infinite ‘sink’ branch of  $B_0$ -nodes, where the remainder of the ray for  $\chi_j$  can be mapped as soon as one of its literals is satisfied. We show in Appendix C.1 that  $\mathcal{T}_\dagger, \{A(a)\} \models \mathbf{q}_\varphi$  iff  $\varphi$  is satisfiable. To illustrate, the CQ  $\mathbf{q}_\varphi$  for  $\varphi = (p_1 \vee p_2) \wedge \neg p_1$  and a fragment of the canonical model  $\mathcal{C}_{\mathcal{T}_\dagger, \{A(a)\}}$  are shown below:



Here,  $\bullet$  are the points in  $B_0$  and the labels on arrows indicate the subscripts of the binary predicates  $P$  (the empty label means all three:  $+$ ,  $-$  and  $0$ ); predicates  $A$ ,  $B_+$ ,  $B_-$  are not shown in  $\mathcal{C}_{\mathcal{T}_\dagger, \{A(a)\}}$ .  $\square$

The proof above uses OMQs  $\mathbf{Q}_\varphi = (\mathcal{T}_\dagger, \mathbf{q}_\varphi)$  over a data instance with a single individual constant. Thus:

**Corollary 18.** *No polynomial-time algorithm can construct FO- or NDL-rewritings for the OMQs  $\mathbf{Q}_\varphi$  unless  $P = NP$ .*

*Proof.* Indeed, if a polynomial-time algorithm could find a rewriting  $\mathbf{q}'_\varphi$  of  $\mathbf{Q}_\varphi$ , then we would be able to check whether  $\varphi$  is satisfiable in polynomial time by evaluating  $\mathbf{q}'_\varphi$  over the data instance  $\{A(a)\}$ .  $\square$

Curiously enough, Corollary 18 can be complemented with the following theorem:

**Theorem 19.** *The  $\mathbf{Q}_\varphi$  have polynomial FO-rewritings.*

*Proof.* Define  $\mathbf{q}'_\varphi$  as the FO-sentence

$$\forall xy ((x = y) \wedge A(x) \wedge \varphi^*) \vee \exists xy ((x \neq y) \wedge \mathbf{q}^*_\varphi(x, y)),$$

where  $\varphi^*$  is  $\top$  if  $\varphi$  is satisfiable and  $\perp$  otherwise, and  $\mathbf{q}^*_\varphi(x, y)$  is the polynomial-size FO-rewriting of  $\mathbf{Q}_\varphi$  over data with *at least 2* constants [25, Corollary 14]. Recall that the proof of Theorem 17 shows that, if  $\mathcal{A}$  has a single constant,  $a$ , and there is a homomorphism from  $\mathbf{q}_\varphi$  to  $\mathcal{C}_{\mathcal{T}_\dagger, \mathcal{A}}$ , then  $A(a) \in \mathcal{A}$  and  $\varphi$  is satisfiable. Thus, the first disjunct of  $\mathbf{q}'_\varphi$  is an FO-rewriting of  $\mathbf{Q}_\varphi$  over data instances with a single constant; the case of at least 2 constants follows from [25, Corollary 14].  $\square$

Whether the OMQs  $\mathbf{Q}_\varphi$  have a polynomial-size PE- or NDL-rewritings remains open. We have only managed to construct a modification  $\bar{\mathbf{q}}_\varphi(x)$  of  $\mathbf{q}_\varphi$  with the following interesting properties (details are given in Appendix C.2). Let  $\mathfrak{T}$  be the class of data instances representing finite binary trees with root  $a$  whose edges are labelled with  $P_+$  and  $P_-$ , and some of whose leaves are labelled with  $B_0$ . Let  $\mathcal{QL}$  be any query language such that, for every  $\mathcal{QL}$ -query  $\Phi(x)$  and every  $\mathcal{A} \in \mathfrak{T}$ , the answer to  $\Phi(a)$  over  $\mathcal{A}$  can be computed in time polynomial in  $|\Phi|$  and  $|\mathcal{A}|$ . Typical examples of  $\mathcal{QL}$  are modal-like languages such as certain fragments of XPath [38] or description logic instance queries [4].

**Theorem 20.** *The OMQs  $(\mathcal{T}_\dagger, \bar{\mathbf{q}}_\varphi(x))$  do not have polynomial-size rewritings in  $\mathcal{QL}$  unless  $\text{NP} \subseteq \text{P/poly}$ .*

To our surprise, Theorem 20 is not applicable to PE.<sup>3</sup>

**Theorem 21.** *Evaluating PE-queries over trees in  $\mathfrak{T}$  is NP-hard.*

Finally, we consider bounded-leaf CQs (whose evaluation is NL-complete in the database setting) with fixed ontology and data.

**Theorem 22.** *There is an ontology  $\mathcal{T}_\dagger$  such that answering OMQs of the form  $(\mathcal{T}_\dagger, \mathbf{q})$  with Boolean linear CQs  $\mathbf{q}$  is LOGCFL-hard for query complexity.*

The proof is by reduction of the recognition problem for the hardest LOGCFL language  $\mathcal{L}$  [29, 56]. We construct an ontology  $\mathcal{T}_\dagger$  and a logspace transducer that converts the words  $w$  in the alphabet of  $\mathcal{L}$  to linear CQs  $\mathbf{q}_w$  such that  $w \in \mathcal{L}$  iff  $\mathcal{T}_\dagger, \{A(a)\} \models \mathbf{q}_w$ .

## 6 Experiments & Conclusions

The main positive result of this paper is the development of theoretically optimal NDL-rewritings for three classes  $\text{OMQ}(\mathbf{d}, \mathbf{t}, \infty)$ ,  $\text{OMQ}(\mathbf{d}, 1, \ell)$ ,  $\text{OMQ}(\infty, 1, \ell)$  of OMQs. It was known that answering such OMQs is tractable, but the proofs employed elaborate algorithms tailored for each of the three cases. We have shown that the optimal complexity can be achieved *via NDL-rewriting*, thus reducing OMQ answering to standard query evaluation. This result is practically relevant as many user queries are tree-shaped (see, e.g., [48] for evidence in the RDF setting), and indeed, recent tools for query formulation over ontologies (like [55]) produce tree-shaped CQs. Moreover, the majority of important real-world OWL2 ontologies are of finite depth; see [16] for statistics. In the context of OBDA, OWL2 QL ontologies are often built starting from the database schemas (bootstrapping [31]), which typically do not contain cycles such as ‘every manager is managed by a manager.’ For example, the NPD FactPages ontology,<sup>4</sup> designed to facilitate querying the datasets of the Norwegian Petroleum Directorate, is of depth 5.

The starting point of our research was the observation that standard query rewriting systems tend to produce suboptimal rewritings of the OMQs in these three classes. This is obviously so for UCQ-rewriters [49, 46, 14, 27, 43, 39]. However, this is also true of more elaborate PE-rewriters (which use disjunctions inside conjunctions) [50, 58] whose rewritings in theory can be of superpolynomial size; see Fig. 1(b). Surprisingly, even NDL-rewriters such as Clipper [20], Presto [53] and Rapid [14] do not fare much better in practice. To illustrate, we generated three sequences of OMQs in the class  $\text{OMQ}(1, 1, 2)$  (lying in the intersection of  $\text{OMQ}(\mathbf{d}, \mathbf{t}, \infty)$ ,  $\text{OMQ}(\mathbf{d}, 1, \ell)$  and  $\text{OMQ}(\infty, 1, \ell)$ ) with the ontology from Example 11 and linear CQs of up to 15 atoms as in Example 8 (which are associated with words from  $\{R, S\}^*$ ). By Fig. 1(a), answering these OMQs can be done in NL. The barcharts in Fig. 2 show the number of clauses in their NDL-rewritings produced by Clipper, Presto and Rapid, as well as by our algorithms LIN, LOG and Tw from Sections 3.2–3.4, respectively. The first three NDL-rewritings display a clear exponential growth, with Clipper and Rapid failing to produce rewritings for longer CQs. In contrast, our rewritings grow linearly in accord with theory.

<sup>3</sup>This result might be known but we could not find it in the literature, and so provide a proof in Appendix C.3.

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

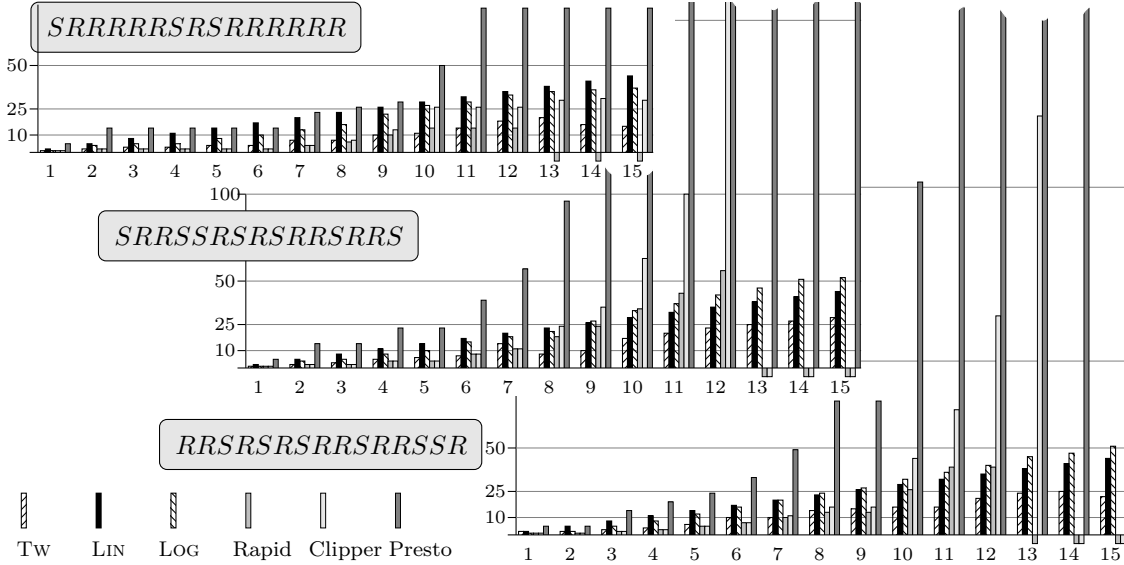


Figure 2: The size of NDL-rewritings produced by different algorithms.

We evaluated the rewritings over a few randomly generated data instances using off-the-shelf datalog engine RDBFox [45]. The experiments (details are in the appendix) show that our rewritings are usually executed faster than those produced by Clipper, Presto and Rapid.

The version of RDBFox we used did not seem to take advantage of the structure of the NL/LOGCFL rewritings by simply materialising all the predicates without using magic sets or optimising programs before execution. It would be interesting to see whether the nonrecursiveness and parallelisability of our rewritings can be utilised to produce efficient execution plans. One could also investigate whether our rewritings can be efficiently implemented using views in standard DBMSs.

Our rewriting algorithms are based on the same idea: pick a point splitting the given CQ into sub-CQs, rewrite the sub-CQs recursively, and then formulate rules that combine the resulting rewritings. The difference between the algorithms is in the choice of the splitting points, which determines the execution plans for OMQs and has a big impact on their performance. The experiments show that none of the three splitting strategies systematically outperforms the others. This suggests that execution times may be dramatically improved by employing an ‘adaptable’ splitting strategy that would work similarly to query execution planners in DBMSs and use statistical information about the relational tables to generate efficient NDL programs. For example, one could first define a ‘cost function’ on some set of alternative rewritings that roughly estimates their evaluation time and then construct a rewriting minimising this function. Such a performance-oriented approach was introduced and exploited in [7], where the target language for OMQ rewritings was joins of UCQs (unions of CQs). Other optimisation techniques for removing redundant rules or sub-queries from rewritings [53, 50, 28, 39] or exploiting the emptiness of certain predicates [59] are also relevant here. In the context of OBDA with relational databases and mappings, integrity constraints [52, 51] and the structure of mappings [18] are particularly important for optimisation.

Having observed that (i) the ontology depth and (ii) the number of leaves in tree-shaped CQs occur in the exponent of our upper bounds for the complexity of OMQ answering algorithms, we regarded (i) and (ii) as parameters and investigated the parameterised complexity of the OMQ answering problem. We proved that the problem is  $W[2]$ -hard in the former case and  $W[1]$ -hard in the latter (it remains open whether these lower bounds are tight). Furthermore, we established that answering OMQs with a fixed ontology (of infinite depth) is NP-complete for tree-shaped CQs and LOGCFL-complete for linear CQs, which dashed hopes of taming intractability by restricting

the ontology size, signature, etc. One remaining open problem is whether answering OMQs with a fixed ontology and tree-shaped CQs is fixed-parameter tractable if the number of leaves is regarded as the parameter.

A more general avenue for future research is to extend the study of succinctness and optimality of rewritings to suitable ontology languages with predicates of higher-arity, such as linear and sticky tgds.

## 7 Acknowledgements

This work was supported by the French ANR grant 12-JS02-007-01 ‘PAGODA: Practical Algorithms for Ontology-Based Data Access’, the UK EPSRC grant EP/M012670 ‘iTract: Islands of Tractability in Ontology-Based Data Access’, the Russian Foundation for Basic Research grant MK-7312.2016.1, and the Russian Academic Excellence Project 5-100. We thank the developers of Clipper and Rapid for making their systems freely available and Riccardo Rosati for the opportunity to conduct experiments with Presto.

## A Proofs for Section 3

### A.1 Lemma 3

LEMMA 3. *Fix any  $w > 0$ . There is an  $L^{NL}$ -transducer that, for any linear NDL-rewriting  $(\Pi, G(\mathbf{x}))$  of an OMQ  $\mathbf{Q}(\mathbf{x})$  over complete data instances with  $w(\Pi, G) \leq w$ , computes a linear NDL-rewriting  $(\Pi', G(\mathbf{x}))$  of  $\mathbf{Q}(\mathbf{x})$  over arbitrary data instances such that  $w(\Pi', G) \leq w + 1$ .*

*Proof.* Let  $(\Pi, G(\mathbf{x}))$  be a linear NDL-rewriting of the OMQ  $\mathbf{Q}(\mathbf{x}) = (\mathcal{T}, \mathbf{q}(\mathbf{x}))$  over complete data instances such that  $w(\Pi, G) \leq w$ . We will replace every clause  $\lambda$  in  $\Pi$  by a set of clauses  $\lambda^*$  defined as follows. Suppose  $\lambda$  is of the form

$$Q(\mathbf{z}) \leftarrow I \wedge EQ \wedge E_1 \wedge \dots \wedge E_n,$$

where  $I$  is the only IDB body atom in  $\lambda$ ,  $EQ$  contains all equality body atoms, and  $E_1, \dots, E_n$  are the EDB body atoms not involving equality. For every atom  $E_i$ , we define a set  $v(E_i)$  of atoms by taking

$$\begin{aligned} v(E_i) &= \{B(z) \mid \mathcal{T} \models B(x) \rightarrow A(x)\} \cup \\ &\quad \{\varrho(y_i, z) \mid \mathcal{T} \models \exists y \varrho(y, x) \rightarrow A(x)\}, & \text{if } E_i = A(z), \\ v(E_i) &= \{\varrho(z, z') \mid \mathcal{T} \models \varrho(x, y) \rightarrow P(x, y)\}, & \text{if } E_i = P(z, z'), \end{aligned}$$

where  $y_i$  is a fresh variable not occurring in  $\lambda$ ; we assume  $P^-(z, z')$  coincides with  $P(z', z)$ , for all binary predicates  $P$ . Intuitively,  $v(E_i)$  captures all atoms that imply  $E_i$  with respect to  $\mathcal{T}$ . Then  $\lambda^*$  consists of the following clauses:

$$\begin{aligned} Q_0(\mathbf{z}_0) &\leftarrow I, \\ Q_{i+1}(\mathbf{z}_i) &\leftarrow H_i(\mathbf{z}_i) \wedge E'_i, \text{ for } 1 \leq i \leq n \text{ and } E'_i \in v(E_i), \\ Q(\mathbf{z}) &\leftarrow H_{n+1}(\mathbf{z}_n) \wedge EQ, \end{aligned}$$

where  $\mathbf{z}_i$  is the restriction of  $\mathbf{z}$  to variables occurring in  $I$  if  $i = 0$  and in  $Q_i(\mathbf{z}_i)$  and  $E'_i$  except for  $y_i$  if  $i > 0$  (note that  $\mathbf{z}_n = \mathbf{z}$ ). Let  $\Pi'$  be the program obtained from  $\Pi$  by replacing each clause  $\lambda$  by the set of clauses  $\lambda^*$ . By construction,  $\Pi'$  is a linear NDL program and its width cannot exceed  $w(\Pi, G) + 1$  (the possible increase of 1 is due to the replacement of unary atoms  $A(z)$  by binary atoms  $\varrho(y_i, z)$ ).

We now argue that  $(\Pi', G(\mathbf{x}))$  is a rewriting of  $\mathbf{Q}(\mathbf{x})$  over arbitrary data instances. It can be easily verified that  $(\Pi', G(\mathbf{x}))$  is equivalent to  $(\Pi'', G(\mathbf{x}))$ , where NDL program  $\Pi''$  is obtained

from  $\Pi$  by replacing each clause  $Q(\mathbf{z}) \leftarrow I \wedge EQ \wedge E_1 \wedge \dots \wedge E_n$  by the (possibly exponentially larger) set of clauses of the form

$$Q(\mathbf{z}) \leftarrow I \wedge EQ \wedge E'_1 \wedge \dots \wedge E'_n,$$

for all  $E'_i \in v(E_i)$  and  $1 \leq i \leq n$ . It thus suffices to show that  $(\Pi'', G(\mathbf{x}))$  is a rewriting of  $Q(\mathbf{x})$  over arbitrary data instances.

First suppose that  $\mathcal{T}, \mathcal{A} \models \mathbf{q}(\mathbf{a})$ , where  $\mathcal{A}$  is an arbitrary data instance. Let  $\mathcal{A}'$  be the complete data instance obtained from  $\mathcal{A}$  by adding the ground atoms:

$$\begin{aligned} P(a, b) & \text{ if } \varrho(a, b) \in \mathcal{A} \text{ and } \mathcal{T} \models \varrho(x, y) \rightarrow P(x, y); \\ A(a) & \text{ if } B(a) \in \mathcal{A} \text{ and } \mathcal{T} \models B(x) \rightarrow A(x); \\ A(a) & \text{ if } \varrho(a, b) \in \mathcal{A} \text{ and } \mathcal{T} \models \exists y \varrho(y, x) \rightarrow A(x). \end{aligned}$$

(We write  $\varrho(a, b) \in \mathcal{A}$  for  $P(a, b) \in \mathcal{A}$  if  $\varrho = P$  and for  $P(b, a)$  if  $\varrho = P^-$ .) Clearly,  $\mathcal{T}, \mathcal{A}' \models \mathbf{q}(\mathbf{a})$ , so we must have  $\Pi, \mathcal{A}' \models G(\mathbf{a})$ . A simple inductive argument (on the order of derivation of ground atoms) shows that whenever a clause  $Q(\mathbf{z}) \leftarrow I \wedge EQ \wedge E_1 \wedge \dots \wedge E_n$  is applied using a substitution  $\mathbf{c}$  for the variables in the body to derive  $Q(\mathbf{c}(\mathbf{z}))$  using  $\Pi$ , we can find a corresponding clause  $Q(\mathbf{z}) \leftarrow I \wedge EQ \wedge E'_1 \wedge \dots \wedge E'_n$  and a substitution  $\mathbf{c}'$  extending  $\mathbf{c}$  (on the fresh variables  $y_i$ ) that allows us to derive  $Q(\mathbf{c}'(\mathbf{z}))$  using  $\Pi''$ . Indeed,

- if  $E_i = A(z)$ , then  $A(\mathbf{c}(z)) \in \mathcal{A}'$ , so there must exist either a unary ground atom  $B(\mathbf{c}(z)) \in \mathcal{A}$  such that  $\mathcal{T} \models B(x) \rightarrow A(x)$  or a binary ground atom  $\varrho(a, \mathbf{c}(z)) \in \mathcal{A}$ , for some  $a \in \text{ind}(\mathcal{A})$ , such that  $\mathcal{T} \models \exists y \varrho(y, x) \rightarrow A(x)$ ; in the latter case, we set  $\mathbf{c}'(y_i) = a$ ;
- similarly, if  $E_i = P(z, z')$ , then there must exist a binary ground atom  $\varrho(\mathbf{c}(z), \mathbf{c}(z')) \in \mathcal{A}$  such that  $\mathcal{T} \models \varrho(x, y) \rightarrow P(x, y)$ .

It then suffices to choose  $Q(\mathbf{z}) \leftarrow I \wedge EQ \wedge E'_1 \wedge \dots \wedge E'_n$  with atoms  $E'_i$  whose form match that of the ground atoms in  $\mathcal{A}$  corresponding to  $E_i$ .

For the converse direction, it suffices to observe that  $\Pi \subseteq \Pi''$ .

To complete the proof, we note that it is in NL to decide whether an atom belongs to  $v(E_i)$ , and thus we can construct the program  $\Pi'$  by means of an  $\mathbb{L}^{\text{NL}}$ -transducer.  $\square$

## A.2 Theorem 6

Next, we combine the transformation in Lemma 5 with the established complexity in Lemma 4 to obtain the combined complexity upper bound:

**THEOREM 6.** *For every  $c > 0$  and  $\mathbf{w} > 0$ , evaluation of NDL queries  $(\Pi, G(\mathbf{x}))$  of width at most  $\mathbf{w}$  and such that  $\text{sd}(\Pi, G) \leq c \log |\Pi|$  is in LOGCFL for combined complexity.*

*Proof.* By Lemma 5,  $(\Pi, G)$  is equivalent to a skinny NDL query  $(\Pi', G)$  such that  $|\Pi'| = O(|\Pi|^2)$ ,  $\mathbf{w}(\Pi', G) \leq \mathbf{w}$ , and  $\text{d}(\Pi', G) \leq \text{sd}(\Pi, G)$ . By Lemma 4, query evaluation for  $(\Pi', G)$  over  $\mathcal{A}$  is done by an NauxPDA in space  $\log |\Pi'| + \mathbf{w}(\Pi', G) \cdot \log |\mathcal{A}| = O(\log |\Pi| + \log |\mathcal{A}|)$  and time  $2^{O(\text{d}(\Pi', G))} \leq |\Pi|^{O(1)}$ .  $\square$

## A.3 Log-rewritings

**Lemma 23.** *For any complete data instance  $\mathcal{A}$ , any  $D \in \mathfrak{D}$ , any type  $\mathbf{w}$  with  $\text{dom}(\mathbf{w}) = \partial D$  and any tuples  $\mathbf{b} \in \text{ind}(\mathcal{A})^{|\partial D|}$  and  $\mathbf{a} \in \text{ind}(\mathcal{A})^{|\mathbf{x}_D|}$ , we have  $\Pi_Q^{\text{Log}}, \mathcal{A} \models G_D^{\mathbf{w}}(\mathbf{b}, \mathbf{a})$  iff there is a homomorphism  $h: \mathbf{q}_D \rightarrow \mathcal{C}_{\mathcal{T}, \mathcal{A}}$  such that*

$$h(x) = \mathbf{a}(x), \quad \text{for } x \in \mathbf{x}_D, \quad \text{and} \quad h(z) = \mathbf{b}(z)\mathbf{w}(z), \quad \text{for } z \in \partial D. \quad (7)$$

*Proof.* ( $\Rightarrow$ ) The proof is by induction on  $\prec$ . For the basis of induction, let  $D$  be of size 1. By the definition of  $\Pi_Q^{\text{LOG}}$ , there exists a type  $\mathbf{s}$  such that  $\text{dom}(\mathbf{s}) = \lambda(\sigma(D))$  and  $\mathbf{w}$  agrees with  $\mathbf{s}$  on  $\partial D$  and a respective tuple  $\mathbf{c} \in \text{ind}(\mathcal{A})^{|\lambda(\sigma(D))|}$  such that  $\mathbf{c}(z) = \mathbf{b}(z)$ , for all  $z \in \partial D$ , and  $\mathbf{c}(x) = \mathbf{a}(x)$ , for all  $x \in \mathbf{x}_D$ , and  $\Pi_Q^{\text{LOG}}, \mathcal{A} \models \text{At}^{\mathbf{s}}(\mathbf{c})$ . Then, for any atom  $S(\mathbf{z}) \in \mathbf{q}_D$ , we have  $\mathbf{z} \subseteq \lambda(\sigma(D))$ , whence  $\mathcal{C}_{\mathcal{T}, \mathcal{A}} \models S(h(\mathbf{z}))$  as  $\mathbf{w}$  agrees with  $\mathbf{s}$  on  $\partial D$ .

For the inductive step, suppose that we have  $\Pi_Q^{\text{LOG}}, \mathcal{A} \models G_D^{\mathbf{w}}(\mathbf{b}, \mathbf{a})$ . By the definition of  $\Pi_Q^{\text{LOG}}$ , there exists a type  $\mathbf{s}$  such that  $\text{dom}(\mathbf{s}) = \lambda(\sigma(D))$  and  $\mathbf{w}$  agrees with  $\mathbf{s}$  on their common domain and a respective tuple  $\mathbf{c} \in \text{ind}(\mathcal{A})^{|\lambda(\sigma(D))|}$  such that  $\mathbf{c}(z) = \mathbf{b}(z)$ , for all  $z \in \partial D$ , and  $\mathbf{c}(x) = \mathbf{a}(x)$ , for all  $x \in \mathbf{x}_D$ , and

$$\Pi_Q^{\text{LOG}}, \mathcal{A} \models \text{At}^{\mathbf{s}}(\mathbf{c}) \wedge \bigwedge_{D' \prec D} G_{D'}^{(\mathbf{s} \cup \mathbf{w}) \upharpoonright \partial D'}(\mathbf{b}_{D'}, \mathbf{a}_{D'}),$$

where  $\mathbf{b}_{D'}$  and  $\mathbf{a}_{D'}$  are the restrictions of  $\mathbf{b} \cup \mathbf{c}$  to  $\partial D'$  and of  $\mathbf{a}$  to  $\mathbf{x}_{D'}$ , respectively. By the induction hypothesis, for any  $D' \prec D$ , there is a homomorphism  $h_{D'}: \mathbf{q}_{D'} \rightarrow \mathcal{C}_{\mathcal{T}, \mathcal{A}}$  such that (7) is satisfied.

Let us show that the  $h_{D'}$  agree on common variables. Suppose that  $z$  is shared by  $\mathbf{q}_{D'}$  and  $\mathbf{q}_{D''}$  for  $D' \prec D$  and  $D'' \prec D$ . By the definition of tree decomposition, for every  $z \in V$ , the nodes  $\{t \mid z \in \lambda(t)\}$  induce a connected subtree of  $T$ , and so  $z \in \lambda(\sigma(D)) \cap \lambda(t') \cap \lambda(t'')$ , where  $t'$  and  $t''$  are the unique neighbours of  $\sigma(D)$  lying in  $D'$  and  $D''$ , respectively. Since  $\mathbf{w}' = (\mathbf{w} \cup \mathbf{s}) \upharpoonright \partial D'$  and  $\mathbf{w}'' = (\mathbf{w} \cup \mathbf{s}) \upharpoonright \partial D''$  are the restrictions of  $\mathbf{w} \cup \mathbf{s}$ , we have  $\mathbf{w}'(z) = \mathbf{w}''(z)$ . This implies that

$$h_{D'}(z) = \mathbf{c}(z)\mathbf{w}'(z) = \mathbf{c}(z)\mathbf{w}''(z) = h_{D''}(z).$$

Now we define  $h$  on every  $z$  in  $\mathbf{q}_D$  by taking

$$h(z) = \begin{cases} h_{D'}(z) & \text{if } z \in \lambda(t), \\ & \text{for } t \in D' \text{ and } D' \prec D, \\ \mathbf{c}(z) \cdot (\mathbf{w} \cup \mathbf{s})(z), & \text{if } z \in \lambda(\sigma(D)). \end{cases}$$

It follows that  $h$  is well defined,  $h$  satisfies (7) and that  $h$  is a homomorphism from  $\mathbf{q}_D$  to  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ . Indeed, take an atom  $S(\mathbf{z}) \in \mathbf{q}_D$ . Then either  $\mathbf{z} \subseteq \lambda(\sigma(D))$ , in which case  $\mathcal{C}_{\mathcal{T}, \mathcal{A}} \models S(h(\mathbf{z}))$  since  $\mathbf{w}$  is compatible with  $\sigma(D)$  and  $\Pi_Q^{\text{LOG}}, \mathcal{A} \models \text{At}^{\mathbf{s}}(\mathbf{c})$ , or  $S(\mathbf{z}) \in \mathbf{q}_{D'}$  for some  $D' \prec D$ , in which case we use the fact that  $h$  extends a homomorphism  $h_{D'}$ .

( $\Leftarrow$ ) The proof is by induction on  $\prec$ . Fix  $D$  and  $\mathbf{w}$  such that  $|\mathbf{w}| = |\partial D|$ . Take tuples  $\mathbf{b} \in \text{ind}(\mathcal{A})^{|\partial D|}$  and  $\mathbf{a} \in \text{ind}(\mathcal{A})^{|\mathbf{x}_D|}$ , and a homomorphism  $h: \mathbf{q}_D \rightarrow \mathcal{C}_{\mathcal{T}, \mathcal{A}}$  satisfying (7). Define a type  $\mathbf{s}$  and a tuple  $\mathbf{c} \in \text{ind}(\mathcal{A})^{|\lambda(\sigma(D))|}$  by taking, for all  $z \in \lambda(\sigma(D))$ ,

$$\mathbf{s}(z) = \mathbf{w} \text{ and } \mathbf{c}(z) = \mathbf{a}, \quad \text{if } h(z) = \mathbf{a}\mathbf{w}, \text{ for } \mathbf{a} \in \text{ind}(\mathcal{A}).$$

By definition,  $\text{dom}(\mathbf{s}) = \lambda(\sigma(D))$  and, by (7),  $\mathbf{s}$  and  $\mathbf{w}$  agree on the common domain. For the inductive step, for each  $D' \prec D$ , let  $h_{D'}$  be the restriction of  $h$  to  $\mathbf{q}_{D'}$  and let  $\mathbf{b}_{D'}$  and  $\mathbf{a}_{D'}$  be the restrictions of  $\mathbf{b} \cup \mathbf{c}$  to  $\partial D'$  and of  $\mathbf{a}$  to  $\mathbf{x}_{D'}$ , respectively. By the inductive hypothesis,  $\Pi_Q^{\text{LOG}}, \mathcal{A} \models G_{D'}^{\mathbf{w}'}(\mathbf{b}_{D'}, \mathbf{a}_{D'})$ . (This argument is not needed for the basis of induction.) Since  $h$  is a homomorphism, we have  $\Pi_Q^{\text{LOG}}, \mathcal{A} \models \text{At}^{\mathbf{s}}(\mathbf{c})$ , whence,  $\Pi_Q^{\text{LOG}}, \mathcal{A} \models G_D^{\mathbf{w}}(\mathbf{b}, \mathbf{a})$ .  $\square$

It follows that answering OMQs  $\mathbf{Q}(\mathbf{x}) = (\mathcal{T}, \mathbf{q}(\mathbf{x}))$  with  $\mathcal{T}$  of finite depth  $d$  and  $\mathbf{q}$  of treewidth  $t$  over any data instance  $\mathcal{A}$  can be done in time

$$\text{poly}(|\mathcal{T}|^{dt}, |\mathbf{q}|, |\mathcal{A}|^t). \quad (4)$$

Indeed, we can evaluate  $(\Pi_Q^{\text{LOG}}, G_T^{\mathbf{e}}(\mathbf{x}))$  in time polynomial in  $|\Pi_Q^{\text{LOG}}|$  and  $|\mathcal{A}|^{w(\Pi_Q^{\text{LOG}}, G_T^{\mathbf{e}})}$ , which are bounded by a polynomial in  $|\mathcal{T}|^{2d(t+1)}$ ,  $|\mathbf{q}|$  and  $|\mathcal{A}|^{2(t+1)}$ .

## A.4 Lin-rewritings

**Lemma 24.** *For any complete data instance  $\mathcal{A}$ , any predicate  $G_n^w$ , any  $\mathbf{a} \in \text{ind}(\mathcal{A})^{|\mathbf{x}^n|}$  and  $\mathbf{b} \in \text{ind}(\mathcal{A})^{|\mathbf{z}_{\exists}^n|}$ , we have  $\Pi_Q^{\text{LIN}}, \mathcal{A} \models G_n^w(\mathbf{b}, \mathbf{a})$  iff there is a homomorphism  $h: \mathbf{q}_n \rightarrow \mathcal{C}_{\mathcal{T}, \mathcal{A}}$  such that*

$$h(x) = \mathbf{a}(x), \quad \text{for } x \in \mathbf{x}^n, \quad \text{and} \quad h(z) = \mathbf{b}(z)\mathbf{w}(z), \quad \text{for } z \in \mathbf{z}_{\exists}^n. \quad (8)$$

*Proof.* The proof is by induction on  $n$ .

For the base case ( $n = M$ ), first suppose that we have  $\Pi_Q^{\text{LIN}}, \mathcal{A} \models G_M^w(\mathbf{b}, \mathbf{a})$ . The only rule in  $\Pi_Q^{\text{LIN}}$  with head predicate  $G_M^w$  is  $G_M^w(\mathbf{z}_{\exists}^M, \mathbf{x}^M) \leftarrow \text{At}^w(\mathbf{z}^M)$  with  $\mathbf{z}^M = \mathbf{z}_{\exists}^M \uplus \mathbf{x}^M$ , which is equivalent to

$$G_M^w(\mathbf{z}_{\exists}^M, \mathbf{x}^M) \leftarrow \bigwedge_{z \in \mathbf{z}^M} \left( \bigwedge_{\substack{A(z) \in \mathbf{q} \\ \mathbf{w}(z) = \varepsilon}} A(z) \wedge \bigwedge_{\substack{P(z, z) \in \mathbf{q} \\ \mathbf{w}(z) = \varepsilon}} P(z, z) \wedge \bigwedge_{\mathbf{w}(z) = \varrho w} A_{\varrho}(z) \right). \quad (9)$$

So the body of this rule must be satisfied when  $\mathbf{b}$  and  $\mathbf{a}$  are substituted for  $\mathbf{z}_{\exists}^M$  and  $\mathbf{x}^M$  respectively. Moreover, by local compatibility of  $\mathbf{w}$  with  $\mathbf{z}^M$ , we know that  $\mathbf{w}(x) = \varepsilon$  for every  $x \in \mathbf{x}^M$ . It follows that

- $A(\mathbf{a}(x)) \in \mathcal{A}$  for every  $A(x) \in \mathbf{q}$  such that  $x \in \mathbf{x}^M$ ;
- $A(\mathbf{b}(z)) \in \mathcal{A}$  for every  $A(z) \in \mathbf{q}$  such that  $z \in \mathbf{z}_{\exists}^M$  and  $\mathbf{w}(z) = \varepsilon$ ;
- $P(\mathbf{a}(x), \mathbf{a}(x)) \in \mathcal{A}$  for every  $P(x, x) \in \mathbf{q}$  such that  $x \in \mathbf{x}^M$ ;
- $P(\mathbf{b}(z), \mathbf{b}(z)) \in \mathcal{A}$  for every  $P(z, z) \in \mathbf{q}$  such that  $z \in \mathbf{z}_{\exists}^M$  and  $\mathbf{w}(z) = \varepsilon$ ;
- $A_{\varrho}(z) \in \mathcal{A}$  for every  $z \in \mathbf{z}^M$  with  $\mathbf{w}(z) = \varrho w$ .

Now let  $h^M$  be the unique mapping from  $\mathbf{z}^M$  to  $\Delta^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$  satisfying (8). First note that  $h^M$  is well-defined, since by the last item, if  $\mathbf{w}(z) = \varrho w$ , then we have  $A_{\varrho}(z) \in \mathcal{A}$  and  $\varrho w \in \mathbf{W}_{\mathcal{T}}$ , so  $\mathbf{b}(z)\varrho w$  belongs to  $\Delta^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ . To show that  $h^M$  is a homomorphism of  $\mathbf{q}_M$  into  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ , first recall that the atoms of  $\mathbf{q}_M$  are of two types:  $A(z)$  or  $P(z, z)$ , with  $z \in \mathbf{z}^M$ . Take some  $A(z) \in \mathbf{q}_M$ . If  $\mathbf{w}(z) = \varepsilon$ , then we immediately obtain either  $A(h^M(z)) = A(\mathbf{a}(z)) \in \mathcal{A}$  or  $A(h^M(z)) = A(\mathbf{b}(z)) \in \mathcal{A}$ , depending on whether  $z \in \mathbf{z}_{\exists}^M$  or in  $\mathbf{x}^M$ . Otherwise, if  $\mathbf{w}(z) \neq \varepsilon$ , then the local compatibility of  $\mathbf{w}$  with  $\mathbf{z}^M$  means that the final letter  $\varrho$  in  $\mathbf{w}(z)$  is such that  $\mathcal{T} \models \exists y \varrho(y, x) \rightarrow A(x)$ , hence  $h^M(z) = \mathbf{b}(z)\mathbf{w}(z) \in \Delta^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ . Finally, suppose that  $P(z, z) \in \mathbf{q}$ . The local compatibility of  $\mathbf{w}$  with  $\mathbf{z}^M$  ensures that either  $\mathbf{w}(z) = \varepsilon$  or  $\mathcal{T} \models P(x, x)$ . In the former case, we have either  $P(\mathbf{a}(z), \mathbf{a}(z)) \in \mathcal{A}$  or  $P(\mathbf{b}(z), \mathbf{b}(z)) \in \mathcal{A}$ , depending again on whether  $z \in \mathbf{z}_{\exists}^M$  or  $z \in \mathbf{x}^M$ . In the latter case,  $(h^M(z), h^M(z)) \in P^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ .

For the other direction, ( $\Leftarrow$ ), of the base case, suppose that the mapping  $h^M$  given by (8) defines a homomorphism from  $\mathbf{q}_M$  into  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ . We therefore have:

- $\mathbf{a}(x) \in \Delta^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$  for every  $A(x) \in \mathbf{q}$  with  $x \in \mathbf{x}^M$ ;
- $\mathbf{b}(z)\mathbf{w}(z) \in \Delta^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$  for every  $A(z) \in \mathbf{q}$  with  $z \in \mathbf{z}_{\exists}^M$ ;
- $(\mathbf{a}(x), \mathbf{a}(x)) \in P^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$  for every  $P(x, x) \in \mathbf{q}$  such that  $x \in \mathbf{x}^M$ ;
- $(\mathbf{b}(z), \mathbf{b}(z)) \in P^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$  for every  $P(z, z) \in \mathbf{q}$  such that  $z \in \mathbf{z}_{\exists}^M$ ;
- $\mathcal{T}, \mathcal{A} \models \exists y \varrho(\mathbf{b}(z), y)$  for every  $z \in \mathbf{z}_{\exists}^M$  with  $\mathbf{w}(z) = \varrho w$  (for otherwise  $\mathbf{b}(z)\mathbf{w}(z)$  would not belong to the domain of  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ ).

The first two items, together with completeness of the data instance  $\mathcal{A}$ , ensure that all atoms in

$$\{A(z) \mid A(z) \in \mathbf{q}, z \in \mathbf{z}^M, \mathbf{w}(z) = \varepsilon\}$$



are present in  $\mathcal{A}$  when  $\mathbf{b}$  and  $\mathbf{a}$  substituted for  $\mathbf{z}_{\exists}^M$  and  $\mathbf{x}^M$ , respectively. The third and fourth items, again together with completeness of  $\mathcal{A}$ , ensure the presence of the atoms in

$$\{P(z, z) \mid P(z, z) \in \mathbf{q}, z \in \mathbf{z}^M, \mathbf{w}(z) = \varepsilon\}.$$

Finally, the fifth item plus completeness of  $\mathcal{A}$  ensure that  $\mathcal{A}$  contains all atoms in

$$\{A_{\varrho}(z) \mid z \in \mathbf{z}^M, \mathbf{w}(z) = \varrho w\}.$$

It follows that the body of the unique rule for  $G_M^w$  is satisfied when  $\mathbf{b}$  and  $\mathbf{a}$  are substituted for  $\mathbf{z}_{\exists}^M$  and  $\mathbf{x}^M$  respectively, and thus  $\Pi_Q^{\text{LIN}}, \mathcal{A} \models G_M^w(\mathbf{b}, \mathbf{a})$ .

For the induction step, assume that the statement has been shown to hold for all  $n \leq k+1 \leq M$ , and let us show that it holds when  $n = k$ . For the first direction,  $(\Rightarrow)$ , suppose  $\Pi_Q^{\text{LIN}}, \mathcal{A} \models G_k^w(\mathbf{b}, \mathbf{a})$ . It follows that there exists a pair of types  $(\mathbf{w}, \mathbf{s})$  compatible with  $(\mathbf{z}^k, \mathbf{z}^{k+1})$  and an assignment  $\mathbf{c}$  of individuals from  $\mathcal{A}$  to the variables in  $\mathbf{z}^k \cup \mathbf{z}^{k+1}$  such that  $\mathbf{c}(x) = \mathbf{a}(x)$  for all  $x \in (\mathbf{z}^k \cup \mathbf{z}^{k+1}) \cap \mathbf{x}$ , and  $\mathbf{c}(z) = \mathbf{b}(z)$  for all  $z \in \mathbf{z}_{\exists}^k$ , and such that every atom in the body of the clause

$$G_k^w(\mathbf{z}_{\exists}^k, \mathbf{x}^k) \leftarrow \text{At}^{\mathbf{w} \cup \mathbf{s}}(\mathbf{z}^k, \mathbf{z}^{k+1}) \wedge G_{k+1}^s(\mathbf{z}_{\exists}^{k+1}, \mathbf{x}^{k+1})$$

is entailed from  $\Pi_Q^{\text{LIN}}, \mathcal{A}$  when the individuals in  $\mathbf{c}$  are substituted for  $\mathbf{z}^k \cup \mathbf{z}^{k+1}$ . Recall that  $\text{At}^{\mathbf{w} \cup \mathbf{s}}(\mathbf{z}^k, \mathbf{z}^{k+1})$  is the conjunction of the following atoms, for  $z, z' \in \mathbf{z}^k \cup \mathbf{z}^{k+1}$ :

- $A(z)$ , if  $A(z) \in \mathbf{q}$  and  $(\mathbf{w} \cup \mathbf{s})(z) = \varepsilon$ ,
- $P(z, z')$ , if  $P(z, z') \in \mathbf{q}$  and  $(\mathbf{w} \cup \mathbf{s})(z) = (\mathbf{w} \cup \mathbf{s})(z') = \varepsilon$ ,
- $z = z'$ , if  $P(z, z') \in \mathbf{q}$  and either  $(\mathbf{w} \cup \mathbf{s})(z) \neq \varepsilon$  or  $(\mathbf{w} \cup \mathbf{s})(z') \neq \varepsilon$ ,
- $A_{\varrho}(z)$ , if  $(\mathbf{w} \cup \mathbf{s})(z)$  is of the form  $\varrho w$ .

In particular, we have  $\Pi_Q^{\text{LIN}}, \mathcal{A} \models G_{k+1}^s(\mathbf{c}(\mathbf{z}_{\exists}^{k+1}), \mathbf{c}(\mathbf{x}^{k+1}))$ . By the induction hypothesis, there exists a homomorphism  $h^{k+1}: \mathbf{q}_{k+1} \rightarrow \mathcal{C}_{\mathcal{T}, \mathcal{A}}$  such that  $h^{k+1}(z) = \mathbf{c}(z)\mathbf{s}(z)$  for every  $z \in \mathbf{z}_{\exists}^{k+1} \cup \mathbf{x}^{k+1}$ . Define a mapping  $h^k$  from  $\text{var}(\mathbf{q}_k)$  to  $\Delta^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$  by setting  $h^k(z) = h^{k+1}(z)$  for every variable  $z \in \text{var}(\mathbf{q}_{k+1})$ , setting  $h^k(x) = \mathbf{a}(x)$  for every  $x \in \mathbf{z}^k \cap \mathbf{x}$ , and setting  $h^k(z) = \mathbf{b}(z)\mathbf{w}(z)$  for every  $z \in \mathbf{z}^k$ . Using the same argument as was used in the base case, we can show that  $h^k$  is well-defined. For atoms from  $\mathbf{q}_k$  involving only variables from  $\mathbf{q}_{k+1}$ , we can use the induction hypothesis to conclude that they are satisfied under  $h^k$ , and for atoms only involving variables from  $\mathbf{z}^k$ , we can argue as in the base case. It thus remains to handle role atoms that contain one variable from  $\mathbf{z}^k$  and one variable from  $\mathbf{z}^{k+1}$ . Consider such an atom  $P(z, z') \in \mathbf{q}_k$ , for  $z \in \mathbf{z}^k$  and  $z' \in \mathbf{z}^{k+1}$ . If  $\mathbf{w}(z) = \mathbf{s}(z') = \varepsilon$ , then the atom  $P(z, z')$  appears in the body of the clause we are considering. It follows that  $\Pi_Q^{\text{LIN}}, \mathcal{A} \models P(\mathbf{c}(z), \mathbf{c}(z'))$ , hence  $(\mathbf{c}(z), \mathbf{c}(z')) \in P^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ . It then suffices to note that  $\mathbf{c}$  agrees with  $\mathbf{a}$  and  $\mathbf{b}$  on the variables in  $\mathbf{z}^k$ . Next suppose that either  $\mathbf{w}(z) \neq \varepsilon$  or  $\mathbf{s}(z') \neq \varepsilon$ . It follows that the clause body contains  $z = z'$ , hence  $\mathbf{c}(z) = \mathbf{c}(z')$ . As  $(\mathbf{w}, \mathbf{s})$  is compatible with  $(\mathbf{z}^k, \mathbf{z}^{k+1})$ , one of the following must hold: either

- (a)  $\mathbf{s}(z') = \mathbf{w}(z)$  and  $\mathcal{T} \models P(x, x)$
- (b) or  $\mathcal{T} \models \varrho(x, y) \rightarrow P(x, y)$  and either  $\mathbf{s}(z') = \mathbf{w}(z)\varrho$  or  $\mathbf{w}(z) = \mathbf{s}(z')\varrho^-$ .

We give the argument in the case where  $z \in \mathbf{z}_{\exists}^k$  (the argument is entirely similar if  $z \in \mathbf{x}^k$ ). If (a) holds, then

$$(h^k(z), h^k(z')) = (\mathbf{b}(z)\mathbf{w}(z), \mathbf{c}(z')\mathbf{s}(z')) = (\mathbf{b}(z)\mathbf{w}(z), \mathbf{c}(z')\mathbf{w}(z)) \in P^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$$

since  $\mathcal{T} \models P(x, x)$  and  $\mathbf{c}(z') = \mathbf{c}(z) = \mathbf{b}(z)$ . If the first option of (b) holds, then

$$(h^k(z), h^k(z')) = (\mathbf{b}(z)\mathbf{w}(z), \mathbf{c}(z')\mathbf{s}(z')) = (\mathbf{b}(z)\mathbf{w}(z), \mathbf{c}(z')\mathbf{w}(z)\varrho) \in P^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$$

since  $\mathcal{T} \models \varrho(x, y) \rightarrow P(x, y)$  and  $\mathbf{c}(z') = \mathbf{c}(z) = \mathbf{b}(z)$ . If the second option of (b) holds, then

$$(h^k(z), h^k(z')) = (\mathbf{b}(z)\mathbf{w}(z), \mathbf{c}(z')\mathbf{s}(z')) = (\mathbf{b}(z)\mathbf{s}(z')\varrho^-, \mathbf{c}(z')\mathbf{s}(z')) \in P^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$$

since  $\mathcal{T} \models \varrho(x, y) \rightarrow P(x, y)$ .

For the converse direction, ( $\Leftarrow$ ), of the induction step, let  $\mathbf{w}$  be a type that is locally compatible with  $\mathbf{z}^k$ , let  $\mathbf{a} \in \text{ind}(\mathcal{A})^{|\mathbf{x}^k|}$ ,  $\mathbf{b} \in \text{ind}(\mathcal{A})^{|\mathbf{z}_{\exists}^k|}$ , and let  $h^k: \mathbf{q}_k \rightarrow \mathcal{C}_{\mathcal{T}, \mathcal{A}}$  be a homomorphism satisfying

$$h^k(x) = \mathbf{a}(x), \quad \text{for } x \in \mathbf{x}^k, \quad \text{and} \quad h^k(z) = \mathbf{b}(z)\mathbf{w}(z), \quad \text{for } z \in \mathbf{z}_{\exists}^k. \quad (10)$$

We let  $\mathbf{c}$  for  $\mathbf{z}^{k+1}$  be defined by setting  $\mathbf{c}(z)$  equal to the unique individual  $c$  such that  $h(z)$  is of the form  $cw$  (for some  $w \in \mathbf{W}_{\mathcal{T}}$ ), and let  $\mathbf{s}$  be the unique type for  $\mathbf{z}^{k+1}$  satisfying  $h(z) = \mathbf{c}(z)\mathbf{s}(z)$  for every  $z \in \mathbf{z}^{k+1}$ ; in other words, we obtain  $\mathbf{s}(z)$  from  $h(z)$  by omitting the initial individual name  $\mathbf{c}(z)$ . Note that since  $\mathbf{x}^{k+1} \subseteq \mathbf{x}^k$ , we have  $\mathbf{a}(x) = \mathbf{c}(x)$  for every  $x \in \mathbf{x}^{k+1}$ . It follows from the fact that  $h^k$  is a homomorphism that  $\mathbf{s}$  is locally compatible with  $\mathbf{z}^{k+1}$  and that, for every role atom  $P(z, z') \in \mathbf{q}_k$  with  $z \in \mathbf{z}^k$  and  $z' \in \mathbf{z}^{k+1}$ , one of the following holds: (i)  $\mathbf{w}(z) = \mathbf{s}(z') = \varepsilon$ , (ii)  $\mathbf{w}(z) = \mathbf{s}(z')$  and  $\mathcal{T} \models P(x, x)$ , (iii)  $\mathcal{T} \models \varrho(x, y) \rightarrow P(x, y)$  and either  $\mathbf{s}(z') = \mathbf{w}(z)\varrho$  or  $\mathbf{w}(z) = \mathbf{s}(z')\varrho^-$ . Thus, the pair of types  $(\mathbf{w}, \mathbf{s})$  is compatible with  $(\mathbf{z}^k, \mathbf{z}^{k+1})$ , and so the following rule appears in  $\Pi_{\mathbf{Q}}^{\text{LIN}}$ :

$$G_k^{\mathbf{w}}(\mathbf{z}_{\exists}^k, \mathbf{x}^k) \leftarrow \text{At}^{\mathbf{w} \cup \mathbf{s}}(\mathbf{z}^k, \mathbf{z}^{k+1}) \wedge G_{k+1}^{\mathbf{s}}(\mathbf{z}_{\exists}^{k+1}, \mathbf{x}^{k+1}),$$

where we recall that  $\text{At}^{\mathbf{w} \cup \mathbf{s}}(\mathbf{z}^k, \mathbf{z}^{k+1})$  is the conjunction of the following atoms, for  $z, z' \in \mathbf{z}^k \cup \mathbf{z}^{k+1}$ :

- $A(z)$ , if  $A(z) \in \mathbf{q}$  and  $(\mathbf{w} \cup \mathbf{s})(z) = \varepsilon$ ,
- $P(z, z')$ , if  $P(z, z') \in \mathbf{q}$  and  $(\mathbf{w} \cup \mathbf{s})(z) = (\mathbf{w} \cup \mathbf{s})(z') = \varepsilon$ ,
- $z = z'$ , if  $P(z, z') \in \mathbf{q}$  and either  $(\mathbf{w} \cup \mathbf{s})(z) \neq \varepsilon$  or  $(\mathbf{w} \cup \mathbf{s})(z') \neq \varepsilon$ ,
- $A_{\varrho}(z)$ , if  $(\mathbf{w} \cup \mathbf{s})(z)$  is of the form  $\varrho w$ .

It follows from Equation (10) and the fact that  $h^k$  is a homomorphism that each of the ground atoms obtained by taking an atom from  $\text{At}^{\mathbf{w} \cup \mathbf{s}}(\mathbf{z}^k, \mathbf{z}^{k+1})$  and substituting  $\mathbf{a}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$  for  $\mathbf{x}^k$ ,  $\mathbf{z}_{\exists}^k$  and  $\mathbf{z}^{k+1}$ , respectively, is present in  $\mathcal{A}$ . By applying the induction hypothesis to the predicate  $G_{k+1}^{\mathbf{s}}$  and the homomorphism  $h^{k+1}: \mathbf{q}_{k+1} \rightarrow \mathcal{C}_{\mathcal{T}, \mathcal{A}}$  obtained by restricting  $h^k$  to  $\text{var}(\mathbf{q}_{k+1})$ , we obtain that  $\Pi_{\mathbf{Q}}^{\text{LIN}}, \mathcal{A} \models G_{k+1}^{\mathbf{s}}(\mathbf{c}(\mathbf{z}_{\exists}^{k+1}), \mathbf{a}(\mathbf{x}^{k+1}))$ . Since for the considered substitution, all body atoms are entailed, we can conclude that  $\Pi_{\mathbf{Q}}^{\text{LIN}}, \mathcal{A} \models G_k^{\mathbf{w}}(\mathbf{b}, \mathbf{a})$ .  $\square$

It follows that answering OMQs  $\mathbf{Q}(\mathbf{x}) = (\mathcal{T}, \mathbf{q}(\mathbf{x}))$  with  $\mathcal{T}$  of finite depth  $d$  and tree-shaped  $\mathbf{q}$  with  $\ell$  leaves over any data instance  $\mathcal{A}$  can be done in time

$$\text{poly}(|\mathcal{T}|^{d\ell}, |\mathbf{q}|, |\mathcal{A}|^{\ell}). \quad (5)$$

Indeed,  $(\Pi_{\mathbf{Q}}^{\text{LIN}}, G(\mathbf{x}))$  can be evaluated in time polynomial in  $|\Pi_{\mathbf{Q}}^{\text{LIN}}|$  and  $|\mathcal{A}|^{\text{w}(\Pi_{\mathbf{Q}}^{\text{LIN}}, G)}$ , which are bounded by a polynomial in  $|\mathcal{T}|^{2d\ell}$ ,  $|\mathbf{q}|$  and  $|\mathcal{A}|^{2\ell}$ .

## A.5 Tw-rewritings

**Lemma 25.** *For any OMQ  $\mathbf{Q}(\mathbf{x}_0) = (\mathcal{T}, \mathbf{q}_0(\mathbf{x}_0))$  with a tree-shaped CQ, any complete data instance  $\mathcal{A}$ , any  $\mathbf{q}(\mathbf{x}) \in \mathfrak{Q}$  and  $\mathbf{a} \in \text{ind}(\mathcal{A})^{|\mathbf{x}|}$ , we have  $\Pi_{\mathbf{Q}}^{\text{TW}}, \mathcal{A} \models G_{\mathbf{q}}(\mathbf{a})$  iff there exists a homomorphism  $h: \mathbf{q} \rightarrow \mathcal{C}_{\mathcal{T}, \mathcal{A}}$  such that  $h(\mathbf{x}) = \mathbf{a}$ .*

*Proof.* An inspection of the definition of the set  $\mathfrak{Q}$  shows that every  $\mathbf{q}(\mathbf{x}) \in \mathfrak{Q}$  is a tree-shaped query having at least one answer variable, with the possible exception of the original query  $\mathbf{q}_0(\mathbf{x}_0)$ , which may be Boolean.

Just as we did for subtrees in Section 3.2, we associate a binary relation on the queries in  $\mathfrak{Q}$  by setting  $\mathbf{q}'(\mathbf{x}') \prec \mathbf{q}(\mathbf{x})$  whenever  $\mathbf{q}'(\mathbf{x}')$  was introduced when applying one of the two decomposition conditions on p. 13 to  $\mathbf{q}(\mathbf{x})$ . The proof is by induction on the subqueries in  $\mathfrak{Q}$ , according to  $\prec$ . We will start by establishing the statement for all queries in  $\mathfrak{Q}$  other than  $\mathbf{q}_0(\mathbf{x}_0)$ , and afterwards, we will complete the proof by giving an argument for  $\mathbf{q}_0(\mathbf{x}_0)$ .

For the basis of induction, take some  $\mathbf{q}(\mathbf{x}) \in \mathfrak{Q}$  that is minimal in the ordering induced by  $\prec$ , which means that  $\text{var}(\mathbf{q}) = \mathbf{x}$ . Indeed, if there is an existentially quantified variable, then the first decomposition rule will give rise to a ‘smaller’ query (in particular, if  $|\text{var}(\mathbf{q})| = 2$ , then although the ‘smaller’ query may have the same atoms, the selected existential variable will become an answer variable). For the first direction,  $(\Rightarrow)$ , suppose that  $\Pi_Q^{\text{Tw}}, \mathcal{A} \models G_{\mathbf{q}}(\mathbf{a})$ . By definition,  $G_{\mathbf{q}}(\mathbf{x}) \leftarrow \mathbf{q}(\mathbf{x})$  is the only clause with head predicate  $G_{\mathbf{q}}$ . Thus, all atoms in the ground CQ  $\mathbf{q}(\mathbf{a})$  are present in  $\mathcal{A}$ , and hence the desired homomorphism exists. For the converse direction,  $(\Leftarrow)$ , suppose there is a homomorphism  $h: \mathbf{q}(\mathbf{x}) \rightarrow \mathcal{C}_{\mathcal{T}, \mathcal{A}}$  such that  $h(\mathbf{x}) = \mathbf{a}$ . It follows that every atom in the ground CQ  $\mathbf{q}(\mathbf{a})$  is entailed from  $\mathcal{T}, \mathcal{A}$ . Completeness of  $\mathcal{A}$  ensures that all of the ground atoms in  $\mathbf{q}(\mathbf{a})$  are present in  $\mathcal{A}$ , and thus we can apply the clause  $G_{\mathbf{q}}(\mathbf{x}) \leftarrow \mathbf{q}(\mathbf{x})$  to derive  $G_{\mathbf{q}}(\mathbf{a})$ .

For the induction step, let  $\mathbf{q}(\mathbf{x}) \in \mathfrak{Q}$  with  $\text{var}(\mathbf{q}) \neq \mathbf{x}$  and suppose that the claim holds for all  $\mathbf{q}'(\mathbf{x}') \in \mathfrak{Q}$  with  $\mathbf{q}'(\mathbf{x}') \prec \mathbf{q}(\mathbf{x})$ . For the first direction,  $(\Rightarrow)$ , suppose  $\Pi_Q^{\text{Tw}}, \mathcal{A} \models G_{\mathbf{q}}(\mathbf{a})$ . There are two cases, depending on which type of clause was used to derive  $G_{\mathbf{q}}(\mathbf{a})$ .

- Case 1:  $G_{\mathbf{q}}(\mathbf{a})$  was derived by an application of the following clause:

$$G_{\mathbf{q}}(\mathbf{z}) \leftarrow \bigwedge_{A(z_q) \in \mathbf{q}} A(z_q) \ \wedge \ \bigwedge_{P(z_q, z_q) \in \mathbf{q}} P(z_q, z_q) \ \wedge \ \bigwedge_{1 \leq i \leq n} G_{\mathbf{q}_i}(\mathbf{x}_i),$$

where  $\mathbf{q}_1(\mathbf{x}_1), \dots, \mathbf{q}_n(\mathbf{x}_n)$  are the subqueries induced by the neighbours of  $z_q$  in the Gaifman graph  $\mathcal{G}$  of  $\mathbf{q}$ . Then there exists a substitution  $\mathbf{c}$  for the variables in the body of this rule that coincides with  $\mathbf{a}$  on  $\mathbf{z}$  and is such that the ground atoms obtained by applying  $\mathbf{c}$  to the variables in the body are all entailed from  $\Pi_Q^{\text{Tw}}, \mathcal{A}$ . In particular,  $\Pi_Q^{\text{Tw}}, \mathcal{A} \models G_{\mathbf{q}_i}(\mathbf{c}(\mathbf{x}_i))$  for every  $1 \leq i \leq n$ . We can apply the induction hypothesis to the  $\mathbf{q}_i(\mathbf{x}_i)$  to obtain homomorphisms  $h_i: \mathbf{q}_i \rightarrow \mathcal{C}_{\mathcal{T}, \mathcal{A}}$  such that  $h_i(\mathbf{x}_i) = \mathbf{c}(\mathbf{x}_i)$ . Let  $h$  be the mapping from  $\text{var}(\mathbf{q})$  to  $\Delta^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$  defined by taking  $h(z) = h_i(z)$ , for  $z \in \text{var}(\mathbf{q}_i)$ . Note that  $h$  is well-defined since  $\text{var}(\mathbf{q}) = \bigcup_{i=1}^n \text{var}(\mathbf{q}_i)$ , and the  $\mathbf{q}_i$  have no variable in common other than  $z_q$ , which is sent to  $\mathbf{c}(z_q)$  by every  $h_i$ . To see why  $h$  is a homomorphism from  $\mathbf{q}$  to  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ , observe that

$$\mathbf{q} = \bigcup_{i=1}^n \mathbf{q}_i \cup \{A(z_q) \in \mathbf{q}\} \cup \{P(z_q, z_q) \in \mathbf{q}\}.$$

By the definition of  $h$ , all atoms in  $\bigcup_{i=1}^n \mathbf{q}_i$  hold under  $h$ . If  $A(z_q) \in \mathbf{q}$ , then  $A(\mathbf{c}(z_q))$  is entailed from  $\Pi_Q^{\text{Tw}}, \mathcal{A}$ , and hence is present in  $\mathcal{A}$ . Similarly, we can show that for every  $P(z_q, z_q) \in \mathbf{q}$ , the ground atom  $P(\mathbf{c}(z_q), \mathbf{c}(z_q))$  belongs to  $\mathcal{A}$ . It follows that all of these atoms hold in  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$  under  $h$ . Finally, we recall that  $\mathbf{c}$  coincides with  $\mathbf{a}$  on  $\mathbf{x}$ , so we have  $h(\mathbf{x}) = \mathbf{a}$ , as required.

- Case 2:  $G_{\mathbf{q}}(\mathbf{a})$  was derived by an application of the following clause, for a tree witness  $\mathbf{t}$  for  $(\mathcal{T}, \mathbf{q}(\mathbf{x}))$  generated by  $\varrho$  with  $\mathbf{t}_r \neq \emptyset$  and  $z_q \in \mathbf{t}_r$ :

$$G_{\mathbf{q}}(\mathbf{x}) \leftarrow A_{\varrho}(z_0) \ \wedge \ \bigwedge_{z \in \mathbf{t}_r \setminus \{z_0\}} (z = z_0) \ \wedge \ \bigwedge_{1 \leq i \leq k} G_{\mathbf{q}_i^{\mathbf{t}}}(\mathbf{x}_i^{\mathbf{t}}),$$

where  $\mathbf{q}_1^{\mathbf{t}}, \dots, \mathbf{q}_k^{\mathbf{t}}$  are the connected components of  $\mathbf{q}$  without  $\mathbf{q}_{\mathbf{t}}$  and  $z_0$  is some variable in  $\mathbf{t}_r$ . There must exist a substitution  $\mathbf{c}$  for the variables in the body of this rule that coincides with  $\mathbf{a}$  on  $\mathbf{x}$  and is such that the ground atoms obtained by applying  $\mathbf{c}$  to the

variables in the body are all entailed from  $\Pi_Q^{\text{Tw}}, \mathcal{A}$ . In particular, for every  $1 \leq i \leq k$ , we have  $\Pi_Q^{\text{Tw}}, \mathcal{A} \models G_{q_i^t}(\mathbf{c}(\mathbf{x}_i^t))$ . We can apply the induction hypothesis to the  $q_i^t(\mathbf{z}_i^t)$  to find homomorphisms  $h_1, \dots, h_k$  of  $q_1^t, \dots, q_k^t$  into  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$  such that  $h_i(\mathbf{x}_i^t) = \mathbf{c}(\mathbf{x}_i^t)$ . Since  $\mathbf{t}$  is a tree witness for  $(\mathcal{T}, \mathbf{q}(\mathbf{x}))$  generated by  $\varrho$ , there exists a homomorphism  $h_t$  of  $q_t$  into  $\mathcal{C}_{\mathcal{T}, \{A_\varrho(a)\}}$  with  $\mathbf{t}_r = h_t^{-1}(a)$  and such that  $h_t(z)$  begins by  $a\varrho$  for every  $z \in \mathbf{t}_t$ . Now take  $z_0 \in \mathbf{t}_r$  such that  $A_\varrho(z_0)$  is the atom in the clause body (recall that  $\mathbf{t}_r \neq \emptyset$ ), and so  $\Pi_Q^{\text{Tw}}, \mathcal{A} \models A_\varrho(\mathbf{c}(z_0))$ , which means that  $A_\varrho(\mathbf{c}(z_0))$  must appear in  $\mathcal{A}$ . It follows that for every element in  $\mathcal{C}_{\mathcal{T}, \{A_\varrho(a)\}}$  of the form  $a\varrho w$ , there exists a corresponding element  $\mathbf{c}(z_0)\varrho w$  in  $\Delta^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ . We now define a mapping  $h$  from  $\text{var}(\mathbf{q})$  to  $\Delta^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$  as follows:

$$h(z) = \begin{cases} h_i(z), & \text{for every } z \in \text{var}(\mathbf{q}_i^t), \\ \mathbf{c}(z_0)\varrho w, & \text{if } z \in \mathbf{t}_t \text{ and } h_t(z) = a\varrho w, \\ \mathbf{c}(z_0) & \text{if } z \in \mathbf{t}_r. \end{cases}$$

Every variable in  $\text{var}(\mathbf{q})$  occurs in  $\mathbf{t}_r \cup \mathbf{t}_t$  or in exactly one of the  $q_i^t$ , and so is assigned a unique value by  $h$ . Note that although  $\mathbf{t}_r \cap \text{var}(\mathbf{q}_i^t)$  is not necessarily empty, due to the equality atoms, we have  $h(z) = h(z')$ , for all  $z, z' \in \mathbf{t}_r$ , and so the function is well-defined. We claim that  $h$  is a homomorphism from  $\mathbf{q}$  into  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ . Clearly, the atoms occurring in some  $q_i^t$  are preserved under  $h$ . Now consider some unary atom  $A(z)$  with  $z \in \mathbf{t}_t$ . Then  $h(z) = \mathbf{c}(z_0)\varrho w$ , where  $h_t(z) = a\varrho w$ . Since  $h_t$  is a homomorphism, we know that  $w$  ends with a role  $\sigma$  such that  $\mathcal{T} \models \exists y \sigma(y, x) \rightarrow A(x)$ . It follows that  $h(z)$  also ends with  $\sigma$ , and thus  $h(z) \in A^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ . Next, consider a binary atom  $P(z, z')$ , where at least one of  $z$  and  $z'$  belongs to  $\mathbf{t}_t$ . As  $h_t$  is a homomorphism, either

- $\mathcal{T} \models \sigma(x, y) \rightarrow P(x, y)$ , for some  $\sigma$ , such that  $h_t(z') = h_t(z)\sigma$  or  $h_t(z) = h_t(z')\sigma^-$ ,
- or  $\mathcal{T} \models P(x, x)$  and  $h_t(z') = h_t(z)$ .

We also know that  $\mathbf{c}(z) = \mathbf{c}(z_0)$  for all  $z \in \mathbf{t}_r$ , hence  $h(z) = h(z_0)$  for all  $z \in \mathbf{t}_r$ . It follows that in the former case we have  $h(z') = h(z)\sigma$  or  $h(z) = h(z')\sigma^-$  with  $\mathcal{T} \models \sigma(x, y) \rightarrow P(x, y)$ . In the latter case, we have  $h(z') = h(z)$  with  $\mathcal{T} \models P(x, x)$ . Thus,  $P(z, z')$  is preserved under  $h$ . Finally, since  $\mathbf{c}$  coincides with  $\mathbf{a}$  on  $\mathbf{x}$ , we have  $h(\mathbf{x}) = \mathbf{a}$ .

For the converse direction, ( $\Leftarrow$ ), of the induction step, suppose that  $h$  is a homomorphism of  $\mathbf{q}$  into  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$  such that  $h(\mathbf{x}) = \mathbf{a}$ . There are two cases to consider, depending on where  $h$  maps the ‘splitting’ variable  $z_q$ .

- Case 1:  $h(z_q) \in \text{ind}(\mathcal{A})$ . Let  $q_1(\mathbf{x}_1), \dots, q_n(\mathbf{x}_n)$  be the subqueries of  $\mathbf{q}(\mathbf{x})$  induced by the neighbours of  $z_q$  in  $\mathcal{G}$ . Recall that  $\mathbf{x}_i$  consists of  $z_q$  and the variables in  $\text{var}(\mathbf{q}_i) \cap \mathbf{x}$ . By restricting  $h$  to  $\text{var}(\mathbf{q}_i)$ , we obtain, for each  $1 \leq i \leq n$ , a homomorphism of  $q_i(\mathbf{x}_i)$  into  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$  that maps  $z_q$  to  $h(z_q)$  and  $\text{var}(\mathbf{q}_i) \cap \mathbf{x}$  to  $\mathbf{a}(\text{var}(\mathbf{q}_i) \cap \mathbf{x})$ . Consider  $\mathbf{a}^*$  defined by taking  $\mathbf{a}^*(x) = \mathbf{a}(x)$  for every  $x \in \text{var}(\mathbf{q}_i) \cap \mathbf{x}$  and  $\mathbf{a}^*(z_q) = h(z_q)$ . By the induction hypothesis, for every  $1 \leq i \leq n$ , we have  $\Pi_Q^{\text{Tw}}, \mathcal{A} \models G_{q_i}(\mathbf{a}^*(\mathbf{x}_i))$ . Next, since  $h$  is a homomorphism, we must have  $h(z_q) \in A^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$  whenever  $A(z_q) \in \mathbf{q}$  and  $(h(z_q), h(z_q)) \in P^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$  whenever  $P(z_q, z_q) \in \mathbf{q}$ . Since  $\mathcal{A}$  is a complete data instance,  $A(h(z_q)) \in \mathcal{A}$  for every  $A(z_q) \in \mathbf{q}$  and  $P(h(z_q), h(z_q))$  for every  $P(z_q, z_q) \in \mathbf{q}$ . We have thus shown that, under the substitution  $\mathbf{a}^*$ , every atom in the body of the clause

$$G_q(z) \leftarrow \bigwedge_{A(z_q) \in \mathbf{q}} A(z_q) \wedge \bigwedge_{P(z_q, z_q) \in \mathbf{q}} P(z_q, z_q) \wedge \bigwedge_{1 \leq i \leq n} G_{q_i}(\mathbf{x}_i),$$

is entailed from  $\Pi_Q^{\text{Tw}}, \mathcal{A}$ . It follows that we must also have  $\Pi_Q^{\text{Tw}}, \mathcal{A} \models G_q(\mathbf{a})$ .

- Case 2:  $h(z_q) \notin \text{ind}(\mathcal{A})$ . Then  $h(z_q)$  is of the form  $b\varrho w$ , for some  $\varrho$ . Let  $V$  be the smallest subset of  $\text{var}(\mathbf{q})$  that contains  $z_q$  and satisfies the following closure property:

- if  $z \in V$ ,  $h(z) \notin \text{ind}(\mathcal{A})$  and  $\mathbf{q}$  contains an atom with  $z$  and  $z'$ , then  $z' \in V$ .

Let  $V'$  consist of all variables  $z$  in  $V$  such that  $h(z) \notin \text{ind}(\mathcal{A})$ . We observe that  $h(z)$  begins by  $b\rho$  for every  $z \in V'$  and  $h(z) = b$  for every  $z \in V \setminus V'$ . Define  $\mathbf{q}_V$  as the CQ comprising all atoms in  $\mathbf{q}$  whose variables are in  $V$  and which contain at least one variable from  $V'$ ; the answer variables of  $\mathbf{q}_V$  are  $V \setminus V'$ . By replacing the initial  $b$  by  $a$  in the mapping  $h$ , we obtain a homomorphism  $h_V$  of  $\mathbf{q}_V$  into  $\mathcal{C}_{\mathcal{T}, \{A_\rho(a)\}}$  with  $V \setminus V' = h_V^{-1}(a)$ . It follows that  $\mathbf{t} = (\mathbf{t}_r, \mathbf{t}_i)$  with  $\mathbf{t}_r = V \setminus V'$  and  $\mathbf{t}_i = V'$  is a tree witness for  $(\mathcal{T}, \mathbf{q}(\mathbf{x}))$  generated by  $\rho$  (and  $\mathbf{q}_i = \mathbf{q}_V$ ). Moreover,  $\mathbf{t}_r \neq \emptyset$  because  $\mathbf{q}$  has at least one answer variable. This means that the program  $\Pi_Q^{\text{Tw}}$  contains the following clause

$$G_{\mathbf{q}}(\mathbf{x}) \leftarrow A_\rho(z_0) \wedge \bigwedge_{z \in \mathbf{t}_r \setminus \{z_0\}} (z = z_0) \wedge \bigwedge_{1 \leq i \leq k} G_{\mathbf{q}_i^t}(\mathbf{x}_i^t),$$

where  $\mathbf{q}_1^t, \dots, \mathbf{q}_k^t$  are the connected components of  $\mathbf{q}$  without  $\mathbf{q}_i$  and  $z_0 \in \mathbf{t}_r$ . Recall that the query  $\mathbf{q}_i^t$  has answer variables  $\mathbf{x}_i^t = \text{var}(\mathbf{q}_i^t) \cap (\mathbf{x} \cup \mathbf{t}_r)$ . Let  $\mathbf{a}^*$  be the substitution for  $\mathbf{x} \cup \mathbf{t}_r$  such that  $\mathbf{a}^*(x) = \mathbf{a}(x)$  for  $x \in \mathbf{x}$  and  $\mathbf{a}^*(z) = h(z)$  for  $z \in \mathbf{t}_r$ . Then, for every  $1 \leq i \leq k$ , there exists a homomorphism  $h_i$  from  $\mathbf{q}_i^t$  to  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$  such that  $h_i(x) = \mathbf{a}^*(x)$  for every  $x \in \mathbf{x}_i^t$ . By the induction hypothesis, we obtain  $\Pi_Q^{\text{Tw}}, \mathcal{A} \models G_{\mathbf{q}_i^t}(\mathbf{a}^*(\mathbf{x}_i^t))$ . Next, since  $h(z) = b$  for every  $z \in \mathbf{t}_r$ , we have  $\mathbf{a}^*(z) = \mathbf{a}^*(z')$  for every  $z, z' \in \mathbf{t}_r$ . Moreover, the presence of the element  $b\rho$  in  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$  means that  $\mathcal{T}, \mathcal{A} \models A_\rho(b)$ . Since  $\mathcal{A}$  is a complete data instance, we have  $A_\rho(b) \in \mathcal{A}$ . It follows that under the substitution  $\mathbf{a}^*$ , all atoms in the body of the clause under consideration are entailed by  $\Pi_Q^{\text{Tw}}, \mathcal{A}$ . Therefore, we must also have  $\Pi_Q^{\text{Tw}}, \mathcal{A} \models G_{\mathbf{q}}(\mathbf{a})$ .

We have thus shown the lemma for all queries  $\mathbf{Q}$  other than  $\mathbf{q}_0(\mathbf{x}_0)$ . Let us now turn to  $\mathbf{q}_0(\mathbf{x}_0)$ .

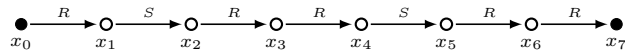
For the first direction, ( $\Rightarrow$ ), suppose  $\Pi_Q^{\text{Tw}}, \mathcal{A} \models G_{\mathbf{q}_0}(\mathbf{a})$ . There are four cases, depending on which type of clause was used to derive  $G_{\mathbf{q}_0}(\mathbf{a})$ . We skip the first three cases, which are identical to those considered in the base case and induction step, and focus instead on the case in which  $G_{\mathbf{q}_0}(\mathbf{a})$  was derived using a clause of the form  $G_{\mathbf{q}_0} \leftarrow A(x)$  with  $A$  a unary predicate such that  $\mathcal{T}, \{A(a)\} \models \mathbf{q}_0$ . In this case, there must exist some  $b \in \text{ind}(\mathcal{A})$  such that  $\mathcal{T}, \mathcal{A} \models A(b)$ . By completeness of  $\mathcal{A}$ , we obtain  $A(b) \in \mathcal{A}$ . Since  $\mathcal{T}, \{A(a)\} \models \mathbf{q}_0$ , we get  $\mathcal{T}, \mathcal{A} \models \mathbf{q}_0$ , which implies the existence of a homomorphism from  $\mathbf{q}_0$  into  $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ .

For the converse direction, ( $\Leftarrow$ ), suppose that there is a homomorphism  $h: \mathbf{q}_0 \rightarrow \mathcal{C}_{\mathcal{T}, \mathcal{A}}$  such that  $h(\mathbf{x}_0) = \mathbf{a}$ . We focus on the case in which  $\mathbf{q}_0$  is Boolean ( $\mathbf{x}_0 = \emptyset$ ) and none of the variables in  $\mathbf{q}_0$  is mapped to an individual constant (the other cases can be handled exactly as in the induction basis and induction step). In this case, there must exist an individual constant  $b$  and some  $\rho$  such that  $h(z)$  begins by  $b\rho$  for every  $z \in \text{var}(\mathbf{q}_0)$ . It follows that  $\mathcal{T}, \{A_\rho(a)\} \models \mathbf{q}_0$ , since the mapping  $h'$  defined by setting  $h'(z) = a\rho w$  whenever  $h(z) = b\rho w$  is a homomorphism from  $\mathbf{q}_0$  to  $\mathcal{C}_{\mathcal{T}, \{A_\rho(a)\}}$ . It follows that  $\Pi_Q^{\text{Tw}}$  contains the clause  $G_{\mathbf{q}_0} \leftarrow A_\rho(x)$ . Since  $b\rho$  occurs in  $\Delta^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ , we have  $\mathcal{T}, \mathcal{A} \models A_\rho(b)$ . By completeness of  $\mathcal{A}$ ,  $A_\rho(b) \in \mathcal{A}$ , and so by applying the clause  $G_{\mathbf{q}_0} \leftarrow A_\rho(x)$ , we obtain  $\Pi_Q^{\text{Tw}}, \mathcal{A} \models G_{\mathbf{q}_0}$ .  $\square$

## A.6 Rewritings Zoo

In this section, we put together the rewritings from Sections 3.2–3.4 for the OMQ given in Examples 8 and 11.

Consider the CQ  $\mathbf{q}(x_0, x_7)$  depicted below (black nodes represent answer variables)



and the following ontology  $\mathcal{T}$  in normal form:

$$\begin{array}{ll}
P(x, y) \rightarrow S(x, y), & P(x, y) \rightarrow R(y, x), \\
A_P(x) \leftrightarrow \exists y P(x, y), & A_{P^-}(x) \leftrightarrow \exists y P(y, x), \\
A_R(x) \leftrightarrow \exists y R(x, y), & A_{R^-}(x) \leftrightarrow \exists y R(y, x), \\
A_S(x) \leftrightarrow \exists y S(x, y) & A_{S^-}(x) \leftrightarrow \exists y S(y, x).
\end{array}$$

### A.6.1 UCQ rewriting

The 9 CQs below form a UCQ rewriting of the OMQ  $\mathbf{Q}(x_0, x_7) = (\mathcal{T}, \mathbf{q}(x_0, x_7))$  over complete data instances given as an NDL program with goal predicate  $G$ :

$$\begin{aligned}
G(x_0, x_7) \leftarrow & [R(x_0, x_1) \wedge S(x_1, x_2) \wedge R(x_2, x_3)] \wedge \\
& [R(x_3, x_4) \wedge S(x_4, x_5) \wedge R(x_5, x_6)] \wedge R(x_6, x_7), \\
G(x_0, x_7) \leftarrow & [A_{P^-}(x_0) \wedge R(x_0, x_3)] \wedge \\
& [R(x_3, x_4) \wedge S(x_4, x_5) \wedge R(x_5, x_6)] \wedge R(x_6, x_7), \\
G(x_0, x_7) \leftarrow & [R(x_0, x_3) \wedge A_P(x_3)] \wedge \\
& [R(x_3, x_4) \wedge S(x_4, x_5) \wedge R(x_5, x_6)] \wedge R(x_6, x_7), \\
G(x_0, x_7) \leftarrow & [R(x_0, x_1) \wedge S(x_1, x_2) \wedge R(x_2, x_3)] \wedge \\
& [A_{P^-}(x_3) \wedge R(x_3, x_6)] \wedge R(x_6, x_7), \\
G(x_0, x_7) \leftarrow & [R(x_0, x_1) \wedge S(x_1, x_2) \wedge R(x_2, x_3)] \wedge \\
& [R(x_3, x_6) \wedge A_P(x_6)] \wedge R(x_6, x_7), \\
G(x_0, x_7) \leftarrow & [A_{P^-}(x_0) \wedge R(x_0, x_3)] \wedge \\
& [A_{P^-}(x_3) \wedge R(x_3, x_6)] \wedge R(x_6, x_7), \\
G(x_0, x_7) \leftarrow & [A_{P^-}(x_0) \wedge R(x_0, x_3)] \wedge \\
& [R(x_3, x_6) \wedge A_P(x_6)] \wedge R(x_6, x_7), \\
G(x_0, x_7) \leftarrow & [R(x_0, x_3) \wedge A_P(x_3)] \wedge \\
& [A_{P^-}(x_3) \wedge R(x_3, x_6)] \wedge R(x_6, x_7), \\
G(x_0, x_7) \leftarrow & [R(x_0, x_3) \wedge A_P(x_3)] \wedge \\
& [R(x_3, x_6) \wedge A_P(x_6)] \wedge R(x_6, x_7).
\end{aligned}$$

We note that a UCQ rewriting over all data instances would in addition contain variants of the CQs above with each of the predicates  $R$  and  $S$  replaced by  $P$  (with arguments swapped appropriately).

The UCQ rewriting above can be obtained by transforming the following PE-formula into UCQ form:

$$\begin{aligned}
& [(R(x_0, x_1) \wedge S(x_1, x_2) \wedge R(x_2, x_3)) \\
& \quad \vee (A_{P^-}(x_0) \wedge R(x_0, x_3)) \vee (R(x_0, x_3) \wedge A_P(x_3))] \\
\wedge & [(R(x_3, x_4) \wedge S(x_4, x_5) \wedge R(x_5, x_6)) \\
& \quad \vee (A_{P^-}(x_3) \wedge R(x_5, x_6)) \vee (R(x_3, x_6) \wedge A_P(x_6))] \\
\wedge & R(x_6, x_7).
\end{aligned}$$

(Intuitively, each of the two sequences  $RSR$  in the query can be derived in three possible ways: from  $RSR$ , from  $A_{P^-}R$  and from  $RA_P$ ).

### A.6.2 Log-rewriting

As explained in Example 11, we split  $T$  into  $D_1$  and  $D_2$  and obtain two rules:

$$\begin{aligned} G_T^\varepsilon(x_0, x_7) &\leftarrow G_{D_1}^{x_3 \mapsto \varepsilon}(x_3, x_0) \wedge R(x_3, x_4) \wedge G_{D_2}^{x_4 \mapsto \varepsilon}(x_4, x_7), \\ G_T^\varepsilon(x_0, x_7) &\leftarrow G_{D_1}^{x_3 \mapsto \varepsilon}(x_3, x_0) \wedge A_{P^-}(x_4) \wedge (x_3 = x_4) \wedge G_{D_2}^{x_4 \mapsto P^-}(x_4, x_7). \end{aligned}$$

Next, we split each of  $D_1$  and  $D_2$  into single-atom subqueries, which yields the following rules:

$$\begin{aligned} G_{D_1}^{x_3 \mapsto \varepsilon}(x_3, x_0) &\leftarrow (x_0 = x_1) \wedge A_{P^-}(x_1) \wedge (x_1 = x_2) \wedge R(x_2, x_3), \\ G_{D_1}^{x_3 \mapsto \varepsilon}(x_3, x_0) &\leftarrow R(x_0, x_1) \wedge (x_1 = x_2) \wedge A_P(x_2) \wedge (x_2 = x_3), \\ G_{D_1}^{x_3 \mapsto \varepsilon}(x_3, x_0) &\leftarrow R(x_0, x_1) \wedge S(x_1, x_2) \wedge R(x_2, x_3), \\ G_{D_2}^{x_4 \mapsto \varepsilon}(x_4, x_7) &\leftarrow (x_4 = x_5) \wedge A_P(x_5) \wedge (x_5 = x_6) \wedge R(x_6, x_7), \\ G_{D_2}^{x_4 \mapsto \varepsilon}(x_4, x_7) &\leftarrow S(x_4, x_5) \wedge R(x_5, x_6) \wedge R(x_6, x_7), \\ G_{D_2}^{x_4 \mapsto P^-}(x_4, x_7) &\leftarrow A_{P^-}(x_4) \wedge (x_4 = x_5) \wedge R(x_5, x_6) \wedge R(x_6, x_7). \end{aligned}$$

Note that in each case we consider only those types that give rise to predicates that have definitions in the rewriting. The resulting NDL rewriting with goal  $G_T^\varepsilon$  consists of 8 rules. Note, however, that the rewriting illustrated above is a slight simplification of the definition given in Section 3.2: here, for the leaves of the tree decomposition, we directly use the atoms  $\text{At}^s$  instead of including a rule  $G_D^w(\partial D, \mathbf{x}_D) \leftarrow \text{At}^s$  in the rewriting. This simplification clearly does not affect the width of the NDL query or the choice of weight function.

### A.6.3 Lin-rewriting

We assume that  $x_0$  is the root, which makes  $x_7$  the only leaf of the query. (Note that we could have chosen another variable, say  $x_3$ , as the root, with  $x_0$  and  $x_7$  the two leaves.) So, the top-level rule is

$$G(x_0, x_7) \leftarrow G_0^{x_0 \mapsto \varepsilon}(x_0, x_7).$$

We then move along the query and consider the variables  $x_1$ ,  $x_2$  and  $x_3$ . The possible ways of mapping these variables to the canonical model give rise to the following 7 rules:

$$\begin{aligned} G_0^{x_0 \mapsto \varepsilon}(x_0, x_7) &\leftarrow R(x_0, x_1) \wedge P_1^{x_1 \mapsto \varepsilon}(x_1, x_7), \\ G_0^{x_0 \mapsto \varepsilon}(x_0, x_7) &\leftarrow (x_0 = x_1) \wedge A_{P^-}(x_1) \wedge G_1^{x_1 \mapsto P^-}(x_1, x_7), \\ G_1^{x_1 \mapsto \varepsilon}(x_1, x_7) &\leftarrow S(x_1, x_2) \wedge G_2^{x_2 \mapsto \varepsilon}(x_2, x_7), \\ G_1^{x_1 \mapsto \varepsilon}(x_1, x_7) &\leftarrow (x_1 = x_2) \wedge A_P(x_2) \wedge G_2^{x_2 \mapsto P}(x_2, x_7), \\ G_1^{x_1 \mapsto P^-}(x_1, x_7) &\leftarrow A_{P^-}(x_1) \wedge (x_1 = x_2) \wedge G_2^{x_2 \mapsto \varepsilon}(x_2, x_7), \\ G_2^{x_2 \mapsto \varepsilon}(x_2, x_7) &\leftarrow R(x_2, x_3) \wedge G_3^{x_3 \mapsto \varepsilon}(x_3, x_7), \\ G_2^{x_2 \mapsto P}(x_2, x_7) &\leftarrow A_P(x_2) \wedge (x_2 = x_3) \wedge G_3^{x_3 \mapsto \varepsilon}(x_3, x_7). \end{aligned}$$

Next, we move to the variables  $x_4$ ,  $x_5$  and  $x_6$ , which give similar 7 rules:

$$\begin{aligned}
G_3^{x_3 \mapsto \varepsilon}(x_3, x_7) &\leftarrow R(x_3, x_4) \wedge P_4^{x_4 \mapsto \varepsilon}(x_4, x_7), \\
G_3^{x_3 \mapsto \varepsilon}(x_3, x_7) &\leftarrow (x_3 = x_4) \wedge A_{P^-}(x_4) \wedge G_4^{x_4 \mapsto P^-}(x_4, x_7), \\
G_4^{x_4 \mapsto \varepsilon}(x_4, x_7) &\leftarrow S(x_4, x_5) \wedge G_5^{x_5 \mapsto \varepsilon}(x_5, x_7), \\
G_4^{x_4 \mapsto \varepsilon}(x_4, x_7) &\leftarrow (x_4 = x_5) \wedge A_P(x_5) \wedge G_5^{x_5 \mapsto P}(x_5, x_7), \\
G_4^{x_4 \mapsto P^-}(x_4, x_7) &\leftarrow A_{P^-}(x_4) \wedge (x_4 = x_5) \wedge G_5^{x_5 \mapsto \varepsilon}(x_5, x_7), \\
G_5^{x_5 \mapsto \varepsilon}(x_5, x_7) &\leftarrow R(x_5, x_6) \wedge G_6^{x_6 \mapsto \varepsilon}(x_6, x_7), \\
G_5^{x_5 \mapsto P}(x_5, x_7) &\leftarrow A_P(x_5) \wedge (x_5 = x_6) \wedge G_6^{x_6 \mapsto \varepsilon}(x_6, x_7).
\end{aligned}$$

Finally, the last variable can only be mapped to a constant in the data instance, which yields a single rule:

$$G_6^{x_6 \mapsto \varepsilon}(x_6, x_7) \leftarrow R(x_6, x_7).$$

Note that, like in the previous case, we consider only those types that give rise to predicates with definitions (and ignore the dead-ends in the construction).

#### A.6.4 Tw-rewriting

We begin by splitting the query roughly in the middle, that is, we choose  $x_3$  and consider two subqueries:

$$\begin{aligned}
q_{03}(x_0, x_3) &= \exists x_1 x_2 (R(x_0, x_1) \wedge S(x_1, x_2) \wedge R(x_2, x_3)) \\
&\text{and} \\
q_{37}(x_3, x_7) &= \exists x_4 x_5 x_6 (R(x_3, x_4) \wedge S(x_4, x_5) \wedge \\
&\quad R(x_5, x_6) \wedge R(x_6, x_7)).
\end{aligned}$$

Since there is no tree witness  $\mathbf{t}$  for  $(\mathcal{T}, q(x_0, x_7))$  that contains  $x_3$  in  $\mathbf{t}_i$ , we have only one top-level rule:

$$G_{07}(x, y) \leftarrow G_{03}(x_0, x_3) \wedge G_{37}(x_3, x_7).$$

Next, we focus on  $q_{03}$  and choose  $x_1$  as the splitting variable. In this case, there is a tree witness  $\mathbf{t}^1$  with  $\mathbf{t}_i^1 = \{x_1\}$  and  $\mathbf{t}_r^1 = \{x_0, x_2\}$ , and so we obtain two rules for  $G_{03}$ :

$$\begin{aligned}
G_{03}(x_0, x_3) &\leftarrow R(x_0, x_1) \wedge G_{13}(x_1, x_3), \\
G_{03}(x_0, x_3) &\leftarrow A_{P^-}(x_0) \wedge (x_0 = x_2) \wedge R(x_2, x_3).
\end{aligned}$$

The subquery  $q_{13}(x_1, x_3) = \exists x_2 (S(x_1, x_2) \wedge R(x_2, x_3))$  contains two atoms and is split at  $x_2$ . Since there is a tree witness  $\mathbf{t}^2$  for  $(\mathcal{T}, q_{13}(x_1, x_3))$  with  $\mathbf{t}_i^2 = \{x_2\}$  and  $\mathbf{t}_r^2 = \{x_1, x_3\}$ , we obtain two rules:

$$\begin{aligned}
G_{13}(x_1, x_3) &\leftarrow S(x_1, x_2) \wedge R(x_2, x_3), \\
G_{13}(x_1, x_3) &\leftarrow A_P(x_1) \wedge (x_1 = x_3).
\end{aligned}$$

By applying the same procedure to  $q_{37}(x_3, x_7)$ , we get the following five rules:

$$\begin{aligned}
G_{37}(x_3, x_7) &\leftarrow G_{35}(x_3, x_5) \wedge G_{57}(x_5, x_7), \\
G_{37}(x_5, x_7) &\leftarrow R(x_3, x_4) \wedge A_P(x_4) \wedge (x_4 = x_6) \wedge R(x_6, x_7), \\
G_{35}(x_3, x_5) &\leftarrow R(x_3, x_5) \wedge S(x_5, x_7), \\
G_{35}(x_3, x_5) &\leftarrow A_{P^-}(x_3) \wedge (x_3 = x_5), \\
G_{57}(x_3, x_5) &\leftarrow R(x_3, x_4) \wedge R(x_4, x_7).
\end{aligned}$$



Note that the rewriting illustrated above is slightly simpler than the definition in Section 3.4: here, we directly use the atoms of  $\mathbf{q}(\mathbf{x})$  instead of including a rule  $G_{\mathbf{q}}(\mathbf{x}) \leftarrow \mathbf{q}(\mathbf{x})$ , for each  $\mathbf{q}(\mathbf{x})$  without existentially quantified variables. This simplification clearly does not affect the width of the NDL query and the choice of weight function.

## B Proofs for Section 4

### B.1 Theorem 15

THEOREM 15. *pDepth-TREEOMQ is  $W[2]$ -hard.*

*Proof.* We show that  $\mathcal{T}_H^k, \{V_0^0(a)\} \models \mathbf{q}_H^k$  iff  $H$  has a hitting set of size  $k$ . Denote by  $\mathcal{C}$  the canonical model of  $(\mathcal{T}_H^k, \{V_0^0(a)\})$ . For convenience of reference to the points of the canonical model we assume that  $\mathcal{T}_H^k$  contains the following axioms:

$$\begin{aligned} V_i^{l-1}(x) &\rightarrow \exists z v_{i'}^l(x, z) \text{ and} \\ v_{i'}^l(x, z) &\rightarrow P(z, x) \wedge V_{i'}^l(z), & \text{for } 0 \leq i < i' \leq n, \\ V_i^l(x) &\rightarrow E_j^l(z), & \text{for } v_i \in e_j, e_j \in E, \\ E_j^l(x) &\rightarrow \exists z \eta_j^l(x, z) \text{ and} \\ \eta_j^l(x, z) &\rightarrow P(x, z) \wedge E_j^{l-1}(z), & \text{for } 1 \leq j \leq m. \end{aligned}$$

We show that  $\mathcal{C} \models \mathbf{q}_H^k$  iff  $H$  has a hitting set of size  $k$ .

( $\Rightarrow$ ) Suppose  $h: \mathbf{q}_H^k \rightarrow \mathcal{C}$  is a homomorphism. Note that  $\mathcal{C}$  satisfies the following properties: (i)  $w \in E_j^0$  iff  $w = av_{i_1}^1 v_{i_2}^2 \dots v_{i_s}^s \eta_j^s \eta_j^{s-1} \dots \eta_j^1$  where  $v_{j_s} \in e_j$  and (ii) all points in  $\Delta^{\mathcal{C}}$  have at most one  $P$ -predecessor. By starting with some  $E_j^0$  atom and applying first (i) and then iterating (ii), we conclude that  $h(y) = av_{i_1}^1 \dots v_{i_k}^k$  for some  $1 \leq i_1 < i_2 < \dots < i_k \leq n$ . We claim that  $\{v_{i_1}, v_{i_2}, \dots, v_{i_k}\}$  is a hitting set in  $H$ . Indeed, for every branch  $j$  of  $\mathbf{q}_H^k$ , there is  $1 \leq s \leq k$  such that this branch is mapped on  $\mathcal{C}$  in the following way:

$$\begin{aligned} h(z_j^l) &= av_{i_1}^1 v_{i_2}^2 \dots v_{i_l}^l, & s \leq l \leq k-1, \\ h(z_j^l) &= av_{i_1}^1 v_{i_2}^2 \dots v_{i_s}^s \eta_j^s \eta_j^{s-1} \dots \eta_j^{l+1}, & 0 \leq l < s, \end{aligned}$$

with  $v_{i_s} \in e_j$ . This can be shown by induction on  $l$  from 0 to  $k-1$  using (i) to prove the base of induction and (ii) to prove the induction step. Therefore, for every  $j$ , there exists  $s$  such that  $v_{i_s} \in e_j$ .

( $\Leftarrow$ ) Suppose  $\{v_{i_1}, v_{i_2}, \dots, v_{i_k}\}$  is a hitting set in  $H$ . We construct a homomorphism  $h$  from  $\mathbf{q}_H^k$  to  $\mathcal{C}$ . First, we set  $h(y) = av_{i_1}^1 \dots v_{i_k}^k$ . Then, for each  $1 \leq j \leq m$ , we find  $s$  such  $v_{i_s} \in e_j$  and define  $h$  as follows:

$$\begin{aligned} h(z_j^l) &= av_{i_1}^1 v_{i_2}^2 \dots v_{i_l}^l, & s \leq l \leq k-1, \\ h(z_j^l) &= av_{i_1}^1 v_{i_2}^2 \dots v_{i_s}^s \eta_j^s \eta_j^{s-1} \dots \eta_j^{l+1}, & 0 \leq l < s. \end{aligned}$$

It should be clear that  $h$  is indeed a homomorphism. □

### B.2 Theorem 16

THEOREM 16. *pLeaves-TREEOMQ is  $W[1]$ -hard.*

*Proof.* We prove that  $\mathcal{T}_G, \{A(a)\} \models \mathbf{q}_G$  iff  $G$  has a clique containing one vertex from each set  $V_i$ .

We start with some preliminaries. First note we assume that the final axiom in  $\mathcal{T}_G$  (which uses the syntactic sugar  $\wedge$ ) is actually given by the following three axioms (where  $P$  is a fresh binary predicate):

$$\begin{aligned} B(x) &\rightarrow \exists y P(x, y), \\ P(x, y) &\rightarrow U(x, y), \\ P(x, y) &\rightarrow U(y, x). \end{aligned}$$

To simplify notation, we will abbreviate  $\mathcal{C}_{\mathcal{T}_G, \{A(a)\}}$  by  $\mathcal{C}$ , and for every  $1 \leq j \leq M$ , we let  $\mathbf{w}(v_j) = L_j^1 L_j^2 \dots L_j^{2M}$ . Observe that for every  $v_{j_1} \in V_1, v_{j_2} \in V_2, \dots, v_{j_p} \in V_p$ , the element  $aw(v_{j_1})\mathbf{w}(v_{j_2}) \dots \mathbf{w}(v_{j_p})$  belongs to  $\Delta^{\mathcal{C}}$ . Further, observe that if  $aw \in \Delta^{\mathcal{C}}$  with  $|w| = 2M \cdot p$ , then there exist  $v_{j_1} \in V_1, v_{j_2} \in V_2, \dots, v_{j_p} \in V_p$  such that  $w = \mathbf{w}(v_{j_1})\mathbf{w}(v_{j_2}) \dots \mathbf{w}(v_{j_p})$ .

( $\Rightarrow$ ) Suppose that  $\mathcal{T}_G, \{A(a)\} \models \mathbf{q}_G$ , and let  $h$  be a homomorphism of  $\mathbf{q}_G$  into  $\mathcal{C}$ . Note that because of the atom  $B(y)$ , the variable  $y$  must be sent by  $h$  to an element occurring at the end of the  $p$ th block. As noted above, every such element takes the form

$$aw(v_{j_1})\mathbf{w}(v_{j_2}) \dots \mathbf{w}(v_{j_p})$$

where  $v_{j_1} \in V_1, v_{j_2} \in V_2, \dots, v_{j_p} \in V_p$ . We claim that  $\{v_{j_1}, \dots, v_{j_p}\}$  is a clique in  $G$ . To see why, consider the  $i$ th branch of  $\mathbf{q}_G$ , compactly represented as follows:

$$(U^{2M-2} \cdot (YY \cdot U^{2M-2})^i \cdot SS)(y, z_i)$$

By examining the axioms, we see that starting from the first occurrence of  $YY$ , every  $U$  and  $Y$  atom takes us one step closer to  $a$  (prior to the first  $YY$ , we may go back and forth on the extra  $P$ -edge leaving from  $h(y)$ ). It follows that  $SS$  must be mapped within the  $p$ -th block of the selected branch, and since  $S$  is present only at positions  $2j_{p-i}$  and  $2j_{p-i} + 1$  of the block, we must have  $h(z_i) = aw(v_{j_1}) \dots \mathbf{w}(v_{j_{p-i-1}}) L_{j_{p-i}}^1 \dots L_{j_{p-i}}^{2j_{p-i}-1}$ . As the distance between consecutive occurrences of  $YY$  (and between the final  $YY$  and the  $SS$ ) is  $2M - 2$ , it follows that all  $YY$  blocks occur at positions  $2j_{p-i}$  and  $2j_{p-i} + 1$  of blocks  $p - i + 1, \dots, p$ , which implies that  $v_{j_{p-i+1}}, \dots, v_{j_p}$  are neighbours of  $v_{j_i}$  in  $G$ . Since  $\mathbf{q}_G$  contains branches for every  $1 \leq i < p$ , the selected vertices  $v_{j_1}, \dots, v_{j_p}$  are all neighbours in  $G$ , and  $G$  contains a clique with the required properties.

( $\Leftarrow$ ) Suppose that  $v_{j_1} \in V_1, \dots, v_{j_p} \in V_p$  form a clique. We construct a homomorphism  $h$  of  $\mathbf{q}_G$  into  $\mathcal{C}$ . First, set  $h(y) = aw$  where  $w = \mathbf{w}(v_{j_1})\mathbf{w}(v_{j_2}) \dots \mathbf{w}(v_{j_p})$  and observe that the atom  $B(y)$  is satisfied by this assignment. We will use  $w[\ell, \ell']$  to denote the subword of  $w$  beginning with the  $\ell$ th symbol of  $w$  and ending with the  $\ell'$ th symbol (note that  $w = |2M \cdot p|$ , so  $w = w[1, 2M \cdot p]$ ). Next, consider the  $i$ th branch of the query, which connects  $y$  to  $z_i$ , and let  $y_0, y_1, \dots, y_{2M(i+1)}$  be the variables lying between  $y$  and  $z_i$  with  $y_0 = y$  and  $z_i = y_{2M(i+1)}$ . For  $0 \leq k \leq 2j_{p-i}$ , we set  $h(y_k) = h(y)$  if  $k$  is even, and set  $h(y_k) = h(y)P$  otherwise. Observe that because  $P$  is included in both  $U$  and  $U^-$ , we satisfy all binary atoms between variables from  $\{y_0, \dots, y_{2j_{p-i}}\}$ . For  $2j_{p-i} < k \leq 2M(i+1)$ , we set

$$h(y_k) = aw[1, 2M \cdot p - (k - 2j_{p-i})].$$

Note that, in particular, this yields

$$\begin{aligned} h(y_{2M(i+1)-2}) &= aw[1, 2M(p - i - 1) + 2j_{p-i} + 2], \\ h(y_{2M(i+1)-1}) &= aw[1, 2M(p - i - 1) + 2j_{p-i} + 1], \\ h(y_{2M(i+1)}) &= aw[1, 2M(p - i - 1) + 2j_{p-i}], \end{aligned}$$

so the final two  $S$ -atoms in the branch are satisfied by  $h$ . It is easy to see that all  $U$ -atoms between variables from  $y_{2j_{p-i}}, \dots, y_{2M(i+1)}$  are also satisfied. Finally, using the fact that vertices  $v_{j_{p-i+1}}, \dots, v_{j_p}$  are neighbours of  $v_{j_{p-i}}$ , we can show that all of the  $Y$ -atoms in the  $i$ th branch are satisfied by  $h$ . As we have constructed a homomorphism from  $\mathbf{q}_G$  into  $\mathcal{C}$ , we can conclude  $\mathcal{T}_G, \{A(a)\} \models \mathbf{q}_G$ .  $\square$

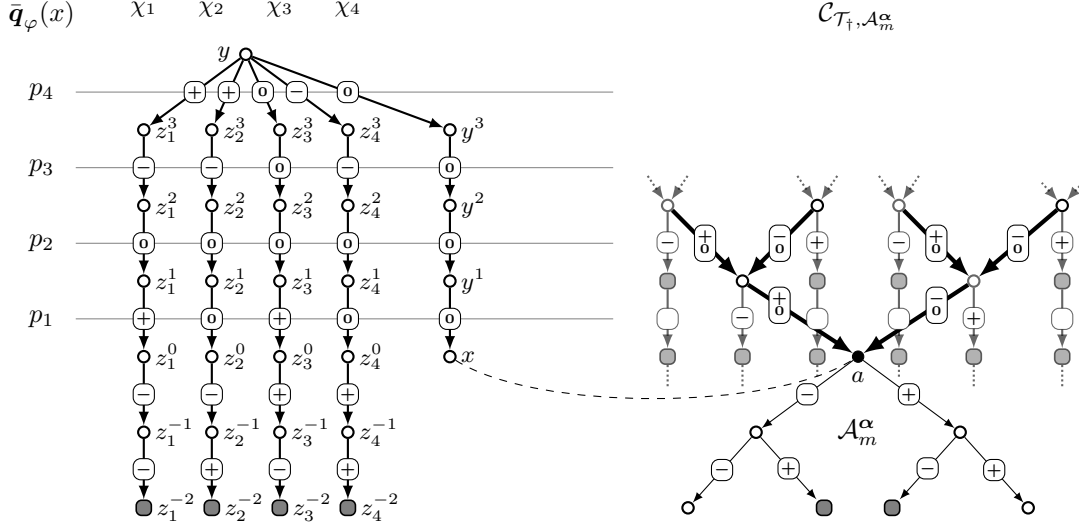


Figure 3: Example of  $\bar{q}_\varphi(x)$  and  $\mathcal{C}_{\mathcal{T}_+, \mathcal{A}_m^\alpha}$  for  $\varphi = \chi_1 \wedge \dots \wedge \chi_4$  with  $\chi_1 = (p_1 \vee \neg p_3 \vee p_4)$ ,  $\chi_2 = (\neg p_3 \wedge p_4)$ ,  $\chi_3 = p_1$ ,  $\chi_4 = (\neg p_3 \vee \neg p_4)$  and  $\alpha = (0, 1, 1, 0)$

## C Proofs for Section 5

### C.1 Theorem 17

**THEOREM 17.** *There is an ontology  $\mathcal{T}_\dagger$  such that answering OMQs of the form  $(\mathcal{T}_\dagger, \mathbf{q})$  with Boolean tree-shaped CQs  $\mathbf{q}$  is NP-hard for query complexity.*

*Proof.* We assume that  $\mathcal{T}_\dagger$  consists of the following axioms:

$$\begin{aligned}
& A(x) \rightarrow \exists y v_+(x, y) \\
& v_+(x, y) \rightarrow P_+(y, x) \wedge P_0(y, x) \wedge B_-(y) \wedge A(y), \\
& B_-(x) \rightarrow \exists y \eta_-(x, y) \\
& \eta_-(x, y) \rightarrow P_-(x, y) \wedge B_0(y), \\
& A(x) \rightarrow \exists y v_-(x, y) \\
& v_-(x, y) \rightarrow P_-(y, x) \wedge P_0(y, x) \wedge B_+(y) \wedge A(y), \\
& B_+(x) \rightarrow \exists y \eta_+(x, y) \\
& \eta_+(x, y) \rightarrow P_+(x, y) \wedge B_0(y), \\
& B_0(x) \rightarrow \exists y \eta_0(x, y) \\
& \eta_0(x, y) \rightarrow P_+(x, y) \wedge P_-(x, y) \wedge P_0(x, y) \wedge B_0(y).
\end{aligned}$$

Let  $\mathcal{C}$  be the canonical model of  $(\mathcal{T}_\dagger, \{A(a)\})$ . We prove that  $\mathcal{C} \models \mathbf{q}_\varphi$  iff  $\varphi$  is satisfiable.

( $\Rightarrow$ ) Suppose  $h$  is a homomorphism from  $\mathbf{q}_\varphi$  to  $\mathcal{C}$  and  $h(z_j^k) = h(y) = a\varrho_1 \dots \varrho_n$ , for some roles  $\varrho_l$ . Since  $A(y) \in \mathbf{q}_\varphi$ , it follows that  $\varrho_l \in \{v_+, v_-\}$ . Moreover, because of the structure of  $\mathcal{C}$ , without any loss of generality we may assume that  $n = k$ . Define a valuation  $\nu: \{p_1, \dots, p_k\} \rightarrow \{\mathbf{t}, \mathbf{f}\}$  by taking  $\nu(p_l) = \mathbf{t}$  if  $\varrho_l = v_-$ ,  $\nu(p_l) = \mathbf{f}$ , if  $\varrho_l = v_+$ . We claim that  $\nu$  makes  $\varphi$  true. To verify that the clause  $\chi_j$  is satisfied, consider a number  $1 \leq s \leq k$ , such that the  $j$ th branch of the query is mapped on  $\mathcal{C}$  in the following way:

$$\begin{aligned}
h(z_j^l) &= a\varrho_1 \dots \varrho_l, & s \leq l \leq k, \\
h(z_j^l) &= a\varrho_1 \dots \varrho_s \gamma_1 \dots \gamma_{s-l}, & 0 \leq l < s,
\end{aligned}$$

for some roles  $\gamma_1 \dots \gamma_{s-l}$  with  $\gamma_1 \in \{\eta_-, \eta_+\}$  and  $\gamma_i = \eta_0$  for  $2 \leq i \leq s-l$ . Such  $s$  and the roles  $\gamma_i$  exist, because the  $P$ -atoms in  $\mathcal{C}$  are directed towards the root if they cover  $v$ -atoms, and away from the root if they cover  $\eta$ -atoms ( $s \geq 1$  since  $B_0(z_j^0) \in \mathbf{q}_\varphi$ ). Clearly,  $\mathcal{T}_\dagger \models \gamma_1(x, y) \rightarrow P_+(x, y)$  iff  $\rho_s = v_-$  iff  $\nu(p_s) = \mathbf{t}$  and  $\mathcal{T}_\dagger \models \gamma_1(x, y) \rightarrow P_-(x, y)$  iff  $\rho_s = v_+$  iff  $\nu(p_s) = \mathbf{f}$ . It follows that either  $P_+(z_j^s, z_j^{s-1}) \in \mathbf{q}_\varphi$  and  $\nu(p_s) = \mathbf{t}$ , or  $P_-(z_j^s, z_j^{s-1}) \in \mathbf{q}_\varphi$  and  $\nu(p_s) = \mathbf{f}$ . In either case,  $\chi_j$  contains a literal with  $p_s$  satisfied by  $\nu$ .

( $\Leftarrow$ ) Suppose a valuation  $\nu: \{p_1, \dots, p_k\} \rightarrow \{\mathbf{t}, \mathbf{f}\}$  satisfies  $\varphi$ . Consider the sequence of roles  $\varrho_1 \dots \varrho_k$ , such that for  $1 \leq l \leq k$  we have  $\varrho_l = v_+$ , if  $\nu(p_l) = \mathbf{f}$ , and  $\varrho_l = v_-$ , if  $\nu(p_l) = \mathbf{t}$ . We claim that there exists a homomorphism  $h$  from  $\mathbf{q}_\varphi$  to  $\mathcal{C}$ . First, let  $h(y) = a\varrho_1 \dots \varrho_k$ . To map the  $j$ th branch of the query, consider the maximal  $1 \leq s \leq k$ , such that a  $p_s$ -literal (positive or negative) makes  $\chi_j$  true. Set

$$\begin{aligned} h(z_j^l) &= a\varrho_1 \dots \varrho_l, & s \leq l \leq k-1, \\ h(z_j^l) &= a\varrho_1 \dots \varrho_s \gamma_1 \dots \gamma_{s-l}, & 0 \leq l < s, \end{aligned}$$

where  $\gamma_1 = \eta_+$  if  $p_s$  occurs positively,  $\gamma_1 = \eta_-$  if  $p_s$  occurs negatively and  $\gamma_i = \eta_0$  for  $i \geq 2$ . That  $z_j^l$ , for  $s \leq l \leq k-1$ , are mapped correctly follows from the maximality of  $s$ . That  $z_j^l$  is mapped correctly for  $l = s-1$  follows from the fact that  $p_s$  occurs in  $\chi_j$  positively iff  $P_+(z_j^s, z_j^{s-1}) \in \mathbf{q}_\varphi$  iff  $\nu(p_s) = \mathbf{t}$  iff  $\varrho_s = v_-$  iff  $\gamma_1 = \eta_+$  (similarly for negative  $p_s$ ). Finally,  $z_j^l$  is mapped correctly for  $0 \leq l < s-1$  since the sequence of roles  $\gamma_2 \dots \gamma_{s-l}$  can embed any  $P_+$ ,  $P_-$ , or  $P_0$  roles, and  $B_0$  concept. Thus,  $h$  is a homomorphism from  $\mathbf{q}_\varphi$  to  $\mathcal{C}$ .  $\square$

## C.2 Theorem 20

We need several intermediate results and definitions before we present the proof in the end of the section. Suppose  $\varphi$  is a propositional formula in CNF having  $k$  variables  $p_1, \dots, p_k$  and  $m$  clauses  $\chi_1, \dots, \chi_m$ . We assume that  $m = 2^\ell$ . We associate with every such  $\varphi$  a CQ  $\bar{\mathbf{q}}_\varphi(x)$  with one answer variable  $x$  and the following atoms, where  $1 \leq j \leq m$ ,  $1 \leq l \leq k$ , and  $z_j^k = y^k$ :

$$\begin{aligned} &P_0(y^1, x), \dots, P_0(y^k, y^{k-1}), \\ &P_+(z_j^l, z_j^{l-1}) && \text{if } \chi_j \text{ contains } p_l, \\ &P_-(z_j^l, z_j^{l-1}), && \text{if } \chi_j \text{ contains } \neg p_l, \\ &P_0(z_j^l, z_j^{l-1}), && \text{if } \chi_j \text{ contains no occurrence of } p_l. \end{aligned}$$

Then, for  $0 \leq l \leq \ell-1$ ,

$$\begin{aligned} &P_-(z_j^{-l}, z_j^{-l-1}), && \text{if the } l\text{th bit of } (j-1)_2 \text{ is } 0, \\ &P_+(z_j^{-l}, z_j^{-l-1}), && \text{if the } l\text{th bit of } (j-1)_2 \text{ is } 1, \\ &B_0(z_j^{-\ell}). \end{aligned}$$

See an example in Fig. 3. For any  $\alpha \in \{0, 1\}^m$ , define a data instance  $\mathcal{A}_m^\alpha$  as the full binary tree of depth  $\ell$  (and so  $m = 2^\ell$  leaves) on the binary predicates  $P_-$  (for the left child) and  $P_+$  (for the right child);  $\mathcal{A}_m^\alpha$  contains  $A(a)$  for the root  $a$  of the tree and, for every  $i$ th leaf  $b_i$  of the tree,  $B_0(b_i) \in \mathcal{A}_m^\alpha$  iff  $\alpha_i = 1$ .

Denote by  $f_\varphi: \{0, 1\}^m \rightarrow \{0, 1\}$  the *monotone* function such that  $f_\varphi(\alpha) = 1$  iff the CNF  $\varphi^{-\alpha}$ , which is obtained from  $\varphi$  by removing all conjuncts  $\chi_i$  with  $\alpha_i = 1$ , is satisfiable. It is readily checked that we have

**Lemma 26.** *For any  $\alpha \in \{0, 1\}^m$ ,*

$$\mathcal{T}_\dagger, \mathcal{A}_m^\alpha \models \bar{\mathbf{q}}_\varphi(a) \quad \text{iff} \quad f_\varphi(\alpha) = 1.$$

Let  $\mathcal{QL}$  be any query language such that, for any  $\mathcal{QL}$ -query  $\Phi(x)$  and any  $\mathcal{A}_m^\alpha$ , the answer to  $\Phi(a)$  over  $\mathcal{A}_m^\alpha$  can be computed in time  $\text{poly}(|\Phi|, m)$ .

**Theorem 27.** *The OMQ  $(\mathcal{T}_\dagger, \bar{q}_\varphi(x))$  does not have a polynomial-size rewriting in  $\mathcal{QL}$  unless  $\text{NP} \subseteq \text{P/poly}$ .*

*Proof.* Take any sequence of CNFs  $\varphi_n$  of polynomial size in  $n$  such that  $f_{\varphi_n}$  is NP-hard [25, Sec. 3]. Suppose there is a  $\mathcal{QL}$ -rewriting  $\Phi_n$  of  $(\mathcal{T}_\dagger, \bar{q}_\varphi(x))$  of polynomial size. By adapting the proof of  $\text{P} \subseteq \text{P/poly}$  [3, Theorem 6.6] to the algorithm that checks  $\mathcal{A}_m^\alpha \models \Phi_n(a)$ , we obtain a sequence of polynomial-size circuits computing  $f_{\varphi_n}$ , from which  $\text{NP} \subseteq \text{P/poly}$ .  $\square$

### C.3 Theorem 21

THEOREM 21. *Evaluating PE-queries over trees in  $\mathfrak{T}$  is NP-hard.*

More precisely, we are going to prove:

**Theorem 28.** *The evaluation problem for PE-queries over data instances of the form  $\mathcal{A}_m^\alpha$  is NP-hard.*

*Proof.* Let  $\varphi_k$ ,  $k \geq 1$ , be the 3-CNF with all possible  $m = O(k^3)$  clauses of  $k$  variables. Without loss of generality, we will assume that the number of clauses in  $\varphi_k$  is actually  $m = 2^\ell$ , for some  $\ell$ . We construct a PE-query  $\mathbf{q}_m(x)$  such that, for any  $\alpha \in \{0, 1\}^m$ , we have  $\mathcal{A}_m^\alpha \models \mathbf{q}_m(a)$  iff the CNF  $\varphi_k^{-\alpha}$  is satisfiable, and the size of  $\mathbf{q}_m$  is polynomial in  $m$  (and  $k$ ).

The query  $\mathbf{q}_m(x)$  takes the form

$$\mathbf{q}_m(x) = \exists \mathbf{z} (\mathbf{r}(x, \mathbf{z}) \wedge \mathbf{s}(x, \mathbf{z}) \wedge \mathbf{t}(x, \mathbf{z})),$$

where the subqueries (without quantified variables)  $\mathbf{r}$ ,  $\mathbf{s}$  and  $\mathbf{t}$  and the variables  $\mathbf{z}$  are defined as follows. Among the variables  $\mathbf{z}$ , there are variables  $z_1, \dots, z_m$  corresponding to the leaves of  $\mathcal{A}_m^\alpha$ , variables  $x_1, \dots, x_k$  corresponding to the propositional variables of  $\varphi_k$ , and variables  $x'_1, \dots, x'_k$  corresponding to their negations (there are other auxiliary variables which will be introduced later on).

Now we will describe the subqueries  $\mathbf{r}$ ,  $\mathbf{s}$ ,  $\mathbf{t}$  of  $\mathbf{q}_m$ . The subquery  $\mathbf{r}$  expresses that the variables  $z_1, \dots, z_m$  indeed correspond to the clauses of  $\varphi_k$ ; it takes the form  $\mathbf{r} = \bigwedge_{i=1}^m \mathbf{r}_i$ . Each  $\mathbf{r}_i$  corresponds to a leaf of  $\mathcal{A}_m^\alpha$ . Consider a path from the root  $a$  to this  $i$ th leaf. Let  $P_1, \dots, P_\ell$  be the sequence of labels on the edges of this path, that is, each  $P_i$  is either  $P_-$  or  $P_+$ . Then

$$\mathbf{r}_i = P_1(x, y_i^1) \wedge P_2(y_i^1, y_i^2) \wedge \dots \wedge P_\ell(y_i^{\ell-1}, z_i),$$

where  $y_i^1, \dots, y_i^{\ell-1}$  are variables among  $\mathbf{z}$ .

The subquery  $\mathbf{s}$  encodes that the variables  $x_1, \dots, x_k$  and  $x'_1, \dots, x'_k$  correspond to an arbitrary Boolean assignment. It is of the form  $\mathbf{s} = \bigwedge_{i=1}^k \mathbf{s}_i$ , and each  $\mathbf{s}_i$  is the following:

$$\begin{aligned} P_\pm(x, u_i^1) \wedge P_\pm(u_i^1, u_i^2) \wedge \dots \wedge P_\pm(u_i^{\ell-2}, u_i^{\ell-1}) \wedge \\ \left[ (P_\pm(u_i^{\ell-1}, x_i) \wedge P_\pm(x'_i, u_i^{\ell-1}) \wedge B_0(x_i)) \vee \right. \\ \left. (P_\pm(u_i^{\ell-1}, x'_i) \wedge P_\pm(x_i, u_i^{\ell-1}) \wedge B_0(x'_i)) \right], \end{aligned}$$

where  $u_i^1, \dots, u_i^{\ell-1}$  are variables among  $\mathbf{z}$  and  $P_\pm(x, y) = P_-(x, y) \vee P_+(x, y)$ .

The last subquery  $\mathbf{t}$  encodes that the assignment given by  $x_1, \dots, x_k$  and  $x'_1, \dots, x'_k$  satisfies the CNF given by  $z_1, \dots, z_m$ . The formula  $\mathbf{t}$  has the following form:  $\mathbf{t} = \bigwedge_{i=1}^m \mathbf{t}_i$ . Suppose the clause  $z_i$  is a disjunction of literals  $l_{i,1}, l_{i,2}$  and  $l_{i,3}$ , where each  $l_{i,n}$  is among  $x_1, \dots, x_k$  and  $x'_1, \dots, x'_k$ . Then

$$\mathbf{t}_i = B_0(z_i) \vee B_0(l_{i,1}) \vee B_0(l_{i,2}) \vee B_0(l_{i,3}).$$

It is easy to see that  $\mathbf{q}_m$  is satisfiable over a given  $\mathcal{A}_m^\alpha$  iff  $\mathcal{A}_m^\alpha$  corresponds to a satisfiable 3-CNF  $\varphi_k^{-\alpha}$ . Thus we have reduced the 3-SAT problem to the problem of evaluating  $\mathbf{q}_m$  over  $\mathcal{A}_m^\alpha$ . Since 3-SAT is NP-complete, we thus have shown NP-hardness of our query evaluation problem.  $\square$

## C.4 Theorem 22

**THEOREM 22.** *There is an ontology  $\mathcal{T}_{\ddagger}$  such that answering OMQs of the form  $(\mathcal{T}_{\ddagger}, \mathbf{q})$  with Boolean linear CQs  $\mathbf{q}$  is LOGCFL-hard for query complexity.*

*Proof.* Our proof encodes the hardest LOGCFL language  $\mathcal{L}$  [29] as formulated in [56]. The language  $\mathcal{L}$  enjoys the following property: for every language  $\mathcal{L}'$  over the alphabet  $\Sigma'$  in LOGCFL, there exists a logspace transducer  $\tau$  converting words over  $\Sigma'$  to the words over the alphabet  $\Sigma$  of  $\mathcal{L}$  in the sense that  $w \in \mathcal{L}'$  iff  $\tau(w) \in \mathcal{L}$ . We construct an ontology  $\mathcal{T}_{\ddagger}$  and a logspace transducer that converts the words  $w \in \Sigma^*$  to linear Boolean CQs  $\mathbf{q}_w$  such that

$$w \in \mathcal{L} \quad \text{iff} \quad \mathcal{T}_{\ddagger}, \{A(a)\} \models \mathbf{q}_w.$$

To explain the construction, we begin with a simpler context-free language. Let  $\Sigma_0 = \{a_1, b_1, a_2, b_2\}$  be an alphabet and  $B_0$  be the context-free language generated by the following grammar:

$$S \rightarrow SS, \quad S \rightarrow \epsilon, \quad S \rightarrow a_1 S b_1, \quad S \rightarrow a_2 S b_2.$$

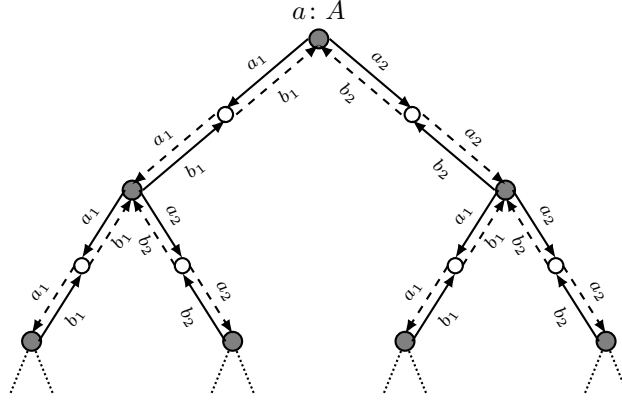
With each word  $w = c_0 \dots c_n$  over  $\Sigma_0$  we associate conjunction  $\gamma_w(u_0, v_0, \dots, u_n, v_n, u_{n+1})$  of the following atoms:

$$R_{c_0}(u_0, v_0), S_{c_0}(v_0, u_1), R_{c_1}(u_1, v_1), S_{c_1}(v_1, u_2), \dots, R_{c_n}(u_n, v_n), S_{c_n}(v_n, u_{n+1}),$$

where  $R_c$  and  $S_c$  are binary predicates, for  $c \in \Sigma_0$ . Let  $\mathcal{T}_0$  contain the following axioms, for  $i = 1, 2$ :

$$D(x) \rightarrow \exists y (R_{a_i}(x, y) \wedge S_{b_i}(y, x) \wedge \exists z (S_{a_i}(y, z) \wedge R_{b_i}(z, y) \wedge D(z))). \quad (11)$$

An initial part of the canonical model of  $(\mathcal{T}_0, \{A(a), D(a)\})$  encoded by these axioms is shown below:



(each large gray node belongs to  $D$ , each solid arrow with label  $c$  belongs to  $R_c$  and each dashed arrow with label  $c$  to  $S_c$ , for  $c \in \Sigma_0$ ). Let  $\mathbf{q}_w^A$  be the following linear Boolean CQ:

$$A(u_0) \wedge \gamma_w(u_0, v_0, \dots, u_n, v_n, u_{n+1}) \wedge A(u_{n+1}).$$

The following claim can readily be verified:

**Proposition 29.** *For every  $w \in \Sigma_0^*$ , we have  $w \in B_0$  iff  $\mathcal{T}_0, \{A(a), D(a)\} \models \mathbf{q}_w^A$ .*

The language  $B_0$  is, however, not LOGCFL-hard. We now reproduce the definition of the hardest LOGCFL language  $\mathcal{L}$  from [56], which uses  $B_0$  as a basis of the construction. Let  $\Sigma = \Sigma_0 \cup \{[, ], \#\}$ , for distinct symbols  $[, ]$ , and  $\#$  not in  $\Sigma_0$ . Then set

$$\begin{aligned} \mathcal{L} = \{ & [x_1 y_1 z_1][x_2 y_2 z_2] \dots [x_k y_k z_k] \mid k \geq 1, \\ & x_i \in (\Sigma_0 \cup \{\#\})^* \{\#\} \cup \{\epsilon\} \text{ and} \\ & z_i \in \{\epsilon\} \cup \{\#\} (\Sigma_0 \cup \{\#\})^*, \text{ for all } i \leq k, \text{ and } y_1 y_2 \dots y_k \in B_0 \}. \end{aligned}$$

To explain the intuition, following [56], let a string of symbols of the form  $[w_1\#w_2\#\dots\#w_n]$ , where  $w_i \in \Sigma^*$  for all  $i$ , be called a *block* and let each of the substrings  $w_i$  be called a *choice*. Then,  $\mathcal{L}$  is the set of all strings of blocks such that there exists a sequence of choices, one from each block, which is in the base language  $B_0$ . The reader should notice that a choice (possibly of the empty string) must be made from each block. For example,

$$[a_1a_2\#b_2b_1] \notin \mathcal{L}, \quad (12)$$

$$[a_1a_2\#b_2b_1][b_2b_1] \in \mathcal{L}, \quad (13)$$

$$[a_1a_2\#b_2b_1][a_1b_1] \notin \mathcal{L}, \quad (14)$$

$$[\#a_1a_2\#b_2b_1][a_1b_1] \in \mathcal{L}. \quad (15)$$

We say that a word  $w$  over  $\Sigma$  is *block-formed* if the following conditions are satisfied:

- the word begins with  $[$  and ends with  $]$ ,
- after each  $[$  there is no  $[$  before  $]$ ;
- each non-final  $]$  is followed immediately by  $[\$ ;
- between each pair of matching  $[$  and  $]$  there is at least one symbol.

With these definitions at hand, we first describe a logspace transducer that, given a word  $w$  over  $\Sigma$ , returns a linear Boolean CQ  $\mathbf{q}_w$  with binary predicates  $R_c$  and  $S_c$ , for  $c \in \Sigma$ , and unary predicates  $A$  and  $E$ . If the word  $w = c_0 \dots c_n$  is block-formed, then  $\mathbf{q}_w$  consists of the following atoms:

$$A(u_0) \wedge \gamma_w(u_0, v_0, \dots, u_n, v_n, u_{n+1}) \wedge A(u_{n+1}).$$

Otherwise, the transducer returns a query that consists of a prefix of  $A(u_0) \wedge \gamma_w(u_0, v_0, \dots, u_n, v_n, u_{n+1})$  and ends in  $E(u_i)$ , for some  $i$ , which will indicate an error (as all queries containing  $E$  will be false in  $\mathcal{T}_{\ddagger}, \{A(a)\}$ ). It is straightforward to verify that the required transducer can be implemented in  $\mathsf{L}$ .

Let  $\mathcal{T}_{\ddagger}$  contain the two axioms (11) and the following axioms:

$$A(x) \rightarrow D(x), \quad (16)$$

$$D(x) \rightarrow \exists y (R_{[}(x, y) \wedge S_{[}(y, x)), \quad (17)$$

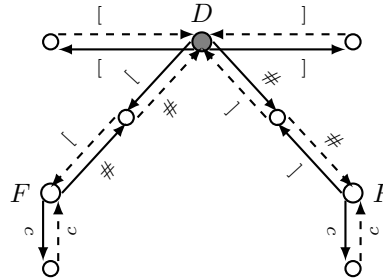
$$D(x) \rightarrow \exists y (R_{\#}(x, y) \wedge S_{\#}(y, x) \wedge \exists z (S_{[}(y, z) \wedge R_{\#}(z, y) \wedge F(z))), \quad (18)$$

$$D(x) \rightarrow \exists y (R_{]}(x, y) \wedge S_{]}(y, x)), \quad (19)$$

$$D(x) \rightarrow \exists y (R_{\#}(x, y) \wedge S_{]}(y, x) \wedge \exists z (S_{\#}(y, z) \wedge R_{]}(z, y) \wedge F(z))), \quad (20)$$

$$F(x) \rightarrow \exists y (R_c(x, y) \wedge S_c(y, x)), \quad \text{for } c \in \Sigma_0 \cup \{\#\}. \quad (21)$$

The four additional branches of the canonical model of  $(\mathcal{T}_{\ddagger}, \{A(a)\})$  at each point in  $D$  are shown below:



(the labels  $D$  and  $F$  are indicated next to the nodes, and, as before, each solid arrow with label  $c$  belongs to  $R_c$  and each dashed arrow with label  $c$  to  $S_c$ , for  $c \in \Sigma_0$ ; to avoid clutter, only one pair of  $c$ -arrows is shown at the bottom).

Let  $\mathbf{q}_w^D$  be defined identically to  $\mathbf{q}_w^A$  except that the two occurrences of  $A$  are replaced by  $D$ . The following property is established similarly to Proposition 29:

**Proposition 30.** *For any block-formed word  $w \in \Sigma^*$ ,*

$$w = [x, \text{ for } x \in (\Sigma_0 \cup \{\#\})^* \{\#\} \cup \{\epsilon\}, \quad \text{iff} \quad \{(17), (18), (21)\}, \{D(d)\} \models \mathbf{q}_w^D.$$

*For any block-formed word  $[w \in \Sigma^*$ ,*

$$w = z], \text{ for } z \in \{\epsilon\} \cup \{\#\}(\Sigma_0 \cup \{\#\})^*, \quad \text{iff} \quad \{(19), (20), (21)\}, \{D(d)\} \models \mathbf{q}_w^D.$$

With these properties established, it can readily be verified that  $\mathcal{T}_{\ddagger}, \{A(a)\} \models \mathbf{q}_w$  iff  $w \in \mathcal{L}$ . Consider a block-formed word  $w \in \Sigma^*$ . Let  $[w_1 \# w_2 \# \dots \# w_n]$  be its  $m$ -th block and  $w_j = y_m$  (that is,  $w_j$  is the segment of the  $B_0$ -word in this block). By Proposition 30, the subtree generated by (18) matches the (translation of)  $[w_1 \# \dots \# w_{j-1} \#$ , whereas the subtree generated by (20) matches  $\# w_{j+1} \# \dots \# w_n]$ . By Proposition 29, the  $w_j$  itself is mapped into the main tree generated by (11). Note that (17) and (19) are needed for the case when  $j = 1$  and  $j = n$ , respectively. Finally, observe that (the translation of)  $w$  has to be mapped starting from  $a$  (the root of the tree) and ending at  $a$ , and that the tree of the canonical model does not contain concept  $E$ , so only a block-formed  $w$  can be mapped to the canonical model. In particular,  $\mathcal{T}_{\ddagger}, \{A(a)\} \not\models \mathbf{q}_w$  for  $w$  of (12) and (14), and  $\mathcal{T}_{\ddagger}, \{A(a)\} \models \mathbf{q}_w$  for  $w$  of (13) and (15).  $\square$

## D Experiments

### D.1 Computing rewritings

We computed 6 types of rewritings for linear queries similar to those in Example 8 and a fixed ontology from Example 11. The first three rewritings were obtained by running executables of Rapid [14], Clipper [20] and Presto [53] with a 15 minute timeout on a desktop machine. The other three rewritings are rewritings LIN, LOG and TW described in Sections 3.3, 3.2 and 3.4 respectively.

We considered the following three sequences:

$$RRSRSRSSRRSSR, \quad (\text{Sequence 1})$$

$$SRRRRRSRRRRRR, \quad (\text{Sequence 2})$$

$$SRRSSRSRRSSRS, \quad (\text{Sequence 3})$$

For each of the three sequences, we consider the line-shaped queries with 1–15 atoms formed by their prefixes. Table 1 presents the sizes of the different types of rewritings.



Table 1: The size (number of clauses) of different types of rewritings for the three sequences of queries (– indicates timeout after 15 minutes)

no. of atoms	Sequence 1 <i>RRSRSSRRSSRRSSR</i>						Sequence 2 <i>SSRRSSRRSSRRSSRR</i>						Sequence 3 <i>SSRRSSRRSSRRSSRR</i>					
	Rapid	Clipper	Presto	LIN	LOG	Tw	Rapid	Clipper	Presto	LIN	LOG	Tw	Rapid	Clipper	Presto	LIN	LOG	Tw
1	1	1	5	2	1	1	1	1	5	2	1	1	1	1	5	2	1	1
2	1	1	5	5	2	0	2	2	14	5	4	2	2	2	14	5	4	2
3	2	2	14	8	5	3	2	2	14	8	5	3	2	2	14	8	5	3
4	3	3	19	11	8	4	2	2	14	11	6	3	4	4	23	11	8	5
5	5	5	24	14	12	6	2	2	14	14	8	4	4	4	23	14	10	6
6	7	7	33	17	16	10	2	2	14	17	10	4	8	8	39	17	15	7
7	10	11	49	20	20	10	4	4	23	20	13	7	11	11	57	20	18	14
8	13	16	77	23	24	14	6	7	29	23	16	7	18	24	96	23	21	8
9	13	16	77	26	27	15	10	13	50	26	22	10	24	35	183	26	27	10
10	26	44	203	29	32	16	14	26	83	29	27	11	34	63	356	29	33	17
11	39	72	329	32	36	16	14	26	83	32	29	14	43	100	356	32	37	20
12	39	126	329	35	40	21	14	26	83	35	33	18	56	302	1028	35	42	23
13	–	241	959	38	45	24	–	30	83	38	35	20	–	–	1712	38	46	25
14	–	–	959	41	47	25	–	31	83	41	36	16	–	–	1712	41	51	27
15	–	–	2723	44	51	22	–	30	83	44	37	15	–	–	5108	44	52	29

Table 2: Generated datasets

dataset	$V$	$p$	$q$	avg. degree of vertices	no. of atoms
1.ttl	1 000	0.050	0.050	50	61 498
2.ttl	5 000	0.002	0.004	10	64 157
3.ttl	10 000	0.002	0.004	20	256 804
4.ttl	20 000	0.002	0.010	40	1 027 028

## D.2 Datasets

We used Erdős-Rényi random graphs with independent parameters  $V$  (number of vertices),  $p$  (probability of an  $R$ -edge) and  $q$  (probability of concepts  $A$  and  $B$  at a given vertex). Note that we intentionally did not introduce any  $S$ -edges. The last parameter, the average degree of a vertex, is  $V \cdot p$ . Table 2 summarises the parameters of the datasets.

## D.3 Evaluating rewritings

We evaluated all obtained rewritings on the datasets in Section D.2 using RDFox triplestore [45] with 999-second timeout. The materialisation time and other relevant statistics are given in Tables 3, 4, and 5.

## D.4 Discussion

Note that the three types of rewritings suggested in this paper give rise to three different rewriting strategies for linear queries. Let us compare how the execution time depends on the exact rewriting strategy. We see in Table 3 that for most queries in Sequence 1 the LIN rewriting shows the best performance, while for Sequences 2 and 3 algorithms LOG and Tw\* are the winners (Tables 4 and 5). Note also that even within a single sequence the results may vary with the number of atoms.

All three rewriting algorithms are based upon a common idea: given a query, pick a point (or a set of points) that would split the query into subqueries, then rewrite these subqueries recursively, and then include rules that join the results into the rewriting of the initial query. However, there is a liberty in the choice of this point, and our rewritings are essentially different in this strategy. Thus, different rewritings generate NDL programs which are related to each other like different execution plans for CQs. Taking into account that we use highly unbalanced data (empty  $S$  versus dense  $R$ ) and that RDFox just materialises all of the predicates of the program without

Table 3: Evaluating rewritings on RDFox - 1

data-set	query size	evaluation time (sec)							no. of answers	no. of generated tuples						
		Rapid	Clipper	Presto	Lin	Log	Tw	Tw*		Rapid	Clipper	Presto	Lin	Log	Tw	Tw*
1.ttl	1	0.021	0.019	0.034	0.049	0.017	0.016	0.01	61390	61390	61390	122780	61449	61390	61390	61390
	2	0.675	0.694	0.706	0.898	0.505	0.652	0.698	976789	976789	976789	1038179	1041822	1038179	976789	976789
	3	0.058	0.053	0.125	0.013	0.112	0.01	0.012	2956	2956	2956	64394	3054	64394	3004	3004
	4	0.204	0.201	0.314	0.087	0.675	0.76	0.12	212213	212213	212213	273710	283409	1314797	1189061	212272
	5	0.12	0.114	0.314	0.014	0.576	0.696	0.064	2956	2956	2956	64453	3150	1105636	976837	3004
	6	0.266	0.248	0.685	0.093	0.266	0.768	0.124	212213	212213	212213	273710	292815	337479	1198455	218710
	7	0.271	0.242	1.11	0.008	0.243	0.687	0.05	2956	2956	2956	64453	3246	125361	982797	3148
	8	0.412	0.377	1.406	0.084	0.904	0.944	0.186	212213	212213	212213	273710	302221	1659409	1410727	431100
	9	3.117	3.337	12.713	3.376	2.941	2.405	1.633	998945	998945	998945	1060442	2927979	2684359	2435551	1455913
	10	1.079	1.102	18.432	0.012	0.607	0.76	0.166	8374	8374	10760	69871	12573	1178714	1203649	224057
	11	2.246	1.984	48.311	0.385	0.945	1.075	0.371	436000	436000	436000	497497	836876	1618743	1663534	664174
	12	13.693	30.032	>999	8.129	6.867	5.922	5.28	999998	999998	1000000	—	5311314	4439352	3217262	2241208
	13	—	6.810	560.206	0.027	0.616	0.946	0.274	20985	—	—	24839	82482	38200	553821	1234421
	14	—	—	913.387	0.013	0.358	0.819	0.27	0	—	—	—	61497	48	312723	1201459
	15	—	—	>999	0.032	0.394	0.994	0.33	2000	—	—	—	—	70277	376713	1417786
2.ttl	1	0.02	0.022	0.039	0.02	0.019	0.017	0.008	64103	64103	64103	128206	64125	64103	64103	64103
	2	0.273	0.305	0.321	0.29	0.297	0.275	0.466	809731	809731	809731	873834	874112	873834	809731	809731
	3	0.03	0.028	0.06	0.011	0.058	0.01	0.013	427	427	427	64561	489	64561	458	458
	4	0.057	0.054	0.103	0.032	0.448	0.315	0.035	8778	8778	8778	72934	74004	947164	818531	8800
	5	0.05	0.046	0.128	0.014	0.423	0.301	0.03	427	427	427	64583	551	938875	809762	458
	6	0.08	0.074	0.27	0.035	0.084	0.316	0.038	8778	8778	8778	72934	75103	77253	819648	9490
	7	0.089	0.080	0.378	0.008	0.078	0.295	0.024	427	427	427	64583	613	68546	810647	551
	8	0.136	0.125	0.467	0.029	0.434	0.322	0.037	8778	8778	8778	72934	76202	1085362	828448	18334
	9	0.202	0.254	1.179	0.369	0.554	0.391	0.102	105853	105853	105853	170009	1020363	1190249	933295	123190
	10	0.174	0.204	2.341	0.011	0.461	0.321	0.052	11	11	438	64167	506	943097	819428	9354
	11	0.192	0.259	4.726	0.036	0.473	0.336	0.053	651	651	9396	64807	74922	944210	820354	11271
	12	0.244	0.699	24.778	0.396	1.034	0.509	0.15	8058	8058	113179	72214	1004735	1940300	934269	124420
	13	—	0.629	20.555	0.015	0.244	0.458	0.084	0	—	438	64156	502	209915	820373	10321
	14	—	—	25.243	0.014	0.153	0.350	0.081	0	—	—	64156	31	200962	820106	10722
	15	—	—	66.916	0.032	0.172	0.335	0.072	0	—	—	64156	64543	265087	828884	19522
3.ttl	1	0.131	0.094	0.225	0.101	0.096	0.14	0.032	256699	256699	256699	513398	256756	256699	256699	256699
	2	2.933	2.946	3.017	2.955	3.053	2.929	3.039	6379932	6379932	6379932	6636631	6638131	6636631	6379932	6379932
	3	0.206	0.175	0.519	0.03	0.499	0.029	0.034	1217	1217	1217	257963	1310	257963	1264	1264
	4	0.399	0.424	0.927	0.171	4.003	3.419	0.231	67022	67022	67022	323825	327716	6961626	6447011	67079
	5	0.36	0.357	1.112	0.036	4.133	3.396	0.179	1217	1217	1217	258020	1405	6895915	6379979	1264
	6	0.632	0.57	1.806	0.169	0.836	3.425	0.228	67022	67022	67022	323825	331647	363640	6450931	69782
	7	0.631	0.581	2.981	0.035	0.756	3.255	0.156	1217	1217	1217	258020	1499	296711	6382460	1405
	8	0.925	0.876	3.739	0.159	4.377	3.405	0.278	67022	67022	67022	323825	335578	7546184	6518010	136975
	9	1.949	2.275	14.564	4.063	5.251	4.169	1.169	1678668	1678668	1678668	1935471	8613829	9225201	8196944	1815899
	10	1.24	1.377	35.109	0.049	4.731	3.571	0.342	60	60	1277	256863	1389	6936178	6449555	68557
	11	1.403	1.798	60.858	0.249	4.846	3.607	0.343	11498	11498	77811	268301	341459	6949160	6462905	85267
	12	1.697	5.413	572.53	4.355	10.128	6.693	1.645	305640	305640	1951654	562443	8780232	15626926	8438115	2058532
	13	—	4.382	484.969	0.082	1.762	4.926	0.599	0	—	1277	256803	1377	917117	6453717	72776
	14	—	—	575.487	0.063	1.115	3.972	0.584	0	—	—	256803	47	850309	6452195	73900
	15	—	—	>999	0.177	1.011	3.585	0.501	0	—	—	—	257974	1107065	6519217	140979
4.ttl	1	0.433	0.451	1.037	0.495	0.439	0.456	0.165	1026526	1026526	1026526	2053052	1026774	1026526	1026526	1026526
	2	27.549	28.088	28.329	27.011	29.532	32.331	31.34	49364886	49364886	49364886	50391412	50404311	50391412	49364886	49364886
	3	2.067	2.409	3.657	0.159	4.087	0.161	0.162	13103	13103	13103	1039882	13613	1039882	13356	13356
	4	4.866	5.438	9.511	1.37	38.919	31.188	2.746	1286991	1286991	1286991	2314018	2353661	52718280	50652125	1287239
	5	4.061	4.032	10.374	0.209	42.943	33.064	2.142	13103	13103	13103	1040130	14119	51444898	49365139	13356
	6	6.909	7.133	16.249	1.443	7.767	36.268	2.782	1286991	1286991	1286991	2314018	2393145	2952225	50691250	1313261
	7	6.614	6.277	23.7	0.243	8.586	29.098	2.02	13103	13103	13103	1040130	14625	1665376	49391598	14115
	8	11.441	10.923	29.1	1.880	54.813	29.426	3.669	1286991	1286991	1286991	2314018	2432629	56098445	51978489	2600996
	9	46.704	50.668	193	76.169	102.055	66.464	33.63	58753514	58753514	58753514	59780541	114973160	114837395	110717131	61339643
	10	14.348	15.503	462	0.375	43.347	30.008	4.694	19966	19966	33014	1046993	35359	52103362	50698955	1321716
	11	19.593	20.907	821	2.843	44.410	31.061	5.319	1872159	1872159	3051184	2899186	4397556	53986724	52602849	3224788
	12	71.354	182.499	>999	172.822	237.478	179.12	90.04	79939048	79939048	120229590	—	199083489	242500074	189429768	140064931
	13	—	54.497	>999	0.562	22.345	44.427	7.105	22474	—	—	53717	—	5686759	50759705	1382714
	14	—	—	>999	0.550	12.462	36.259	7.493	0	—	—	—	253	4356739	50704606	1353393
	15	—	—	>999	1.211	11.315	30.709	7.028	12165	—	—	—	1064542	5395902	52014512	2652797

Table 4: Evaluating rewritings on RDFS - 2

data-set	query size	evaluation time (sec)							no. of answers	no. of generated tuples							
		Rapid	Clipper	Presto	LIN	LOG	Tw	Tw*		Rapid	Clipper	Presto	LIN	LOG	Tw	Tw*	
1.ttl	1	0.009	0.005	0.005	0.005	0.005	0.005	0.007	0	0	0	48	0	0	0	0	
	2	0.009	0.008	0.021	0.05	0.012	0.008	0.007	59	59	59	61508	64406	118	59	59	
	3	0.083	0.058	0.077	0.9	0.093	0.732	0.058	3584	3584	3584	65033	1092161	65033	980373	3584	
	4	2.363	4.049	2.301	8.32	0.11	0.723	0.073	57571	57571	57571	119020	2204964	119079	1034419	57630	
	5	97	92	102	13.599	2	14.272	2.718	59000	59000	59000	120449	3265393	1097297	2035848	59059	
	6	>999	>999	>999	17.882	19	13.881	42.914	59000	-	-	-	4324393	1162212	2039373	62584	
	7	129	122	>999	0.384	0.25	0.749	0.344	2832	2832	2832	-	156824	132259	1030122	6464	
	8	>999	>999	>999	10.963	2	1.82	21.399	55991	-	-	-	3352724	304347	1302623	268322	
	9	162	158	>999	0.395	0.21	0.722	0.344	2832	2832	2832	-	156920	187155	1040255	5895	
	10	>999	>999	>999	11.118	2	12.928	39.21	55991	-	-	-	3362130	1220806	2104667	68937	
	11	>999	>999	>999	20.217	4	14.611	>999	59000	-	-	-	5920653	2251570	2342243	-	
	12	>999	>999	>999	31.648	21	19.079	>999	59000	-	-	-	8714382	3361965	4165789	-	
	13	-	>999	>999	34.395	46	193.512	>999	59000	-	-	-	9783393	3429574	4198870	-	
	14	-	>999	>999	39.818	223	190.334	>999	59000	-	-	-	10842393	1509563	4130571	-	
	15	-	>999	>999	49.391	232	226.827	>999	59000	-	-	-	11901393	1594164	4420495	-	
2.ttl	1	0.007	0.007	0.005	0.007	0.007	0.005	0.004	0	0	0	0	31	0	0	0	
	2	0.01	0.01	0.028	0.027	0.011	0.008	0.008	22	22	22	64147	64543	44	22	22	
	3	0.025	0.025	0.041	0.345	0.046	0.313	0.024	256	256	256	64381	879372	64381	809987	256	
	4	0.135	0.136	0.169	4.798	0.055	0.297	0.023	3300	3300	3300	67425	9329702	67447	813053	3322	
	5	1.314	1.278	1.824	39.195	0.513	4.714	0.122	34474	34474	34474	98599	33935400	908352	9240858	34496	
	6	13.597	13.652	19.52	119.212	0.698	4.606	0.178	106742	106742	106742	170867	59117304	1044957	9313360	106998	
	7	1.396	1.34	18.91	0.116	0.102	0.326	0.028	248	248	248	64404	214761	129190	815625	535	
	8	1.572	1.987	20.58	2.518	0.095	0.364	0.069	3478	3478	3478	67634	2968573	199843	825309	12300	
	9	1.397	1.554	35.15	0.118	0.076	0.333	0.033	248	248	248	64404	214823	132187	813759	728	
	10	1.636	2.634	233	2.591	0.639	4.45	0.069	3478	3478	3478	67634	2969672	976875	9245685	4871	
	11	1.677	12.024	895	30.575	0.98	4.434	0.66	35382	35382	35382	99538	26328037	1823608	9285127	44313	
	12	2.009	143	>999	128.532	1.756	5.666	7.999	106895	106895	106895	-	71017728	2184441	10358119	1010563	
	13	-	>999	>999	243.656	2.559	47.098	5.121	110000	-	-	-	115653199	2742932	34483363	145486	
	14	-	>999	>999	325.755	2.866	50.997	12.028	110000	-	-	-	151038934	1448087	35282112	111224	
	15	-	>999	>999	433.438	26.903	54.518	133.512	110000	-	-	-	176515562	9102348	35442252	118515	
3.ttl	1	0.009	0.01	0.009	0.011	0.009	0.011	0.009	0	0	0	0	47	0	0	0	
	2	0.023	0.02	0.115	0.145	0.022	0.019	0.019	57	57	57	256813	257974	114	57	57	
	3	0.123	0.127	0.249	3.364	0.315	3.212	0.136	1462	1462	1462	258218	6668553	258218	6381394	1462	
	4	1.992	1.93	3.072	85.844	0.345	3.21	0.122	36260	36260	36260	293016	86686553	293073	6416249	36317	
	5	47	56	76.8	967	7.09	70.117	1.898	452502	452502	452502	709258	187656175	7089247	86439255	452559	
	6	>999	>999	>999	>999	9.996	73.99	3.965	570000	-	-	-	-	7464849	86558158	571462	-
	7	47	51	>999	1.591	0.736	3.47	0.181	2125	2125	2125	-	883690	518306	6413768	3634	
	8	77	99	>999	60.365	0.667	3.601	1.842	53191	53191	53191	-	22657990	862422	6536462	120327	
	9	50	56	>999	1.885	0.473	3.496	0.223	2125	2125	2125	-	883784	553583	6419638	3446	
	10	79	142	>999	59.019	7.999	67.145	2.083	53191	53191	53191	-	22661921	7401781	86497805	58664	
	11	81	>999	>999	>999	10.862	68.956	50.812	516631	516631	-	-	-	14275796	87027128	587987	-
	12	116	>999	>999	>999	26.218	112.098	306.304	570000	570000	-	-	-	16280643	95308112	8298971	-
	13	-	>999	>999	>999	45.19	>999	785.247	570000	-	-	-	-	27255415	-	1026838	-
	14	-	>999	>999	>999	74.691	>999	>999	570000	-	-	-	-	9092721	-	-	-
	15	-	>999	>999	>999	>999	>999	>999	-	-	-	-	-	-	-	-	-
4.ttl	1	0.026	0.027	0.027	0.035	0.026	0.047	0.029	0	0	0	0	253	0	0	0	
	2	0.068	0.067	0.5	0.543	0.078	0.069	0.07	248	248	248	1027022	1040241	496	248	248	
	3	0.992	0.99	1.483	33.62	1.98	30.768	0.976	12651	12651	12651	1039425	51050537	1039425	49377537	12651	
	4	60.836	69.126	65.671	M	2.175	30.532	1.272	609193	609193	609193	1635967	-	1636215	49974327	609441	
	5	>999	>999	>999	>999	85	>999	60.335	4947136	-	-	-	-	55339044	-	4947384	-
	6	>999	>999	>999	>999	287	>999	261.562	4960000	-	-	-	-	56390837	-	4972651	-
	7	>999	>999	>999	63	5	31.839	3.118	62572	-	-	-	10949093	2141879	50070886	75476	
	8	>999	>999	>999	>999	13	37.121	273.336	2435666	-	-	-	-	6151203	53696984	3723153	-
	9	>999	>999	>999	61	5	31.899	5.725	62572	-	-	-	10949599	2739031	50050255	76176	
	10	>999	>999	>999	>999	131	>999	319.902	2435666	-	-	-	-	58829172	-	2487953	-
	11	>999	>999	>999	M	214	>999	-	4960000	-	-	-	-	111363802	-	-	-
	12	>999	>999	>999	M	>999	>999	-	-	-	-	-	-	-	-	-	-
	13	-	>999	>999	M	>999	>999	-	-	-	-	-	-	-	-	-	-
	14	-	>999	>999	M	>999	>999	-	-	-	-	-	-	-	-	-	-
	15	-	>999	>999	M	>999	>999	-	-	-	-	-	-	-	-	-	-

Table 5: Evaluating rewritings on RDFox - 3

data-set	query size	evaluation time (sec)							no. of answers	no. of generated tuples						
		Rapid	Clipper	Presto	LIN	LOG	Tw	Tw*		Rapid	Clipper	Presto	LIN	LOG	Tw	Tw*
1.ttl	1	0.004	0.003	0.003	0.004	0.003	0.021	0.003	0	0	0	0	48	0	0	0
	2	0.006	0.006	0.017	0.022	0.008	0.014	0.005	59	59	59	61508	64406	118	59	59
	3	0.053	0.06	0.065	0.849	0.087	0.69	0.053	3584	3584	3584	65033	1092161	65033	980373	3584
	4	0.012	0.01	0.074	0.01	0.009	0.008	0.008	2	2	2	61499	3176	168	109	109
	5	0.011	0.009	0.07	0.008	0.009	0.008	0.009	0	0	0	61497	48	166	59	59
	6	0.018	0.015	0.139	0.023	0.09	0.677	0.055	2	2	2	61499	64560	65203	980434	3704
	7	0.017	0.015	0.145	0.01	0.087	0.68	0.057	0	0	0	61497	144	65190	980480	3691
	8	0.025	0.034	0.339	0.026	0.044	0.009	0.008	2	2	135	61499	73966	129565	170	286
	9	0.025	0.034	0.433	0.01	0.034	0.009	0.008	0	0	2	61497	240	65530	109	214
	10	0.035	0.086	0.549	0.029	0.026	0.015	0.015	2	2	135	61499	83372	67690	12950	13114
	11	0.034	0.086	4.445	1.164	0.54	0.765	0.221	133	0	2	61630	1684864	1095576	1227962	251278
	12	0.048	0.223	4.877	0.013	0.137	0.699	0.062	2	2	135	61499	4082	192211	983694	4115
	13	-	-	13.007	0.141	0.153	0.79	0.175	133	-	-	61630	380205	226874	1228297	229396
	14	-	-	382.922	3.738	0.878	1.166	0.318	1967	-	-	63464	3842746	1282299	1809813	270081
	15	-	-	307.184	0.017	0.36	0.771	0.224	11	-	-	61508	16542	242156	1228610	252720
2.ttl	1	0.004	0.004	0.004	0.004	0.004	0.004	0.003	0	0	0	0	31	0	0	0
	2	0.006	0.006	0.02	0.023	0.009	0.006	0.006	22	22	22	64147	64543	44	22	22
	3	0.022	0.019	0.04	0.339	0.045	0.29	0.019	256	256	256	64381	879372	64381	809987	256
	4	0.013	0.011	0.047	0.01	0.01	0.009	0.009	0	0	0	64156	490	75	53	53
	5	0.012	0.011	0.044	0.008	0.009	0.01	0.009	0	0	0	64156	31	75	22	22
	6	0.02	0.016	0.081	0.027	0.042	0.304	0.021	0	0	0	64156	64543	64456	810009	300
	7	0.018	0.015	0.094	0.011	0.041	0.297	0.024	0	0	0	64156	93	64465	810040	309
	8	0.025	0.036	0.182	0.027	0.053	0.01	0.009	0	0	0	64156	65642	129037	75	119
	9	0.026	0.037	0.215	0.012	0.03	0.009	0.009	0	0	0	64156	155	64611	53	106
	10	0.038	0.091	0.327	0.028	0.029	0.014	0.013	0	0	0	64156	66741	65120	1393	1468
	11	0.036	0.09	1.467	0.345	0.358	0.314	0.055	0	0	0	64156	906286	879949	818896	9218
	12	0.052	0.268	1.868	0.014	0.106	0.294	0.03	0	0	0	64156	494	193096	810468	385
	13	-	-	4.579	0.032	0.119	0.359	0.051	0	-	-	64156	74216	193944	819532	10495
	14	-	-	26.213	0.38	0.454	0.37	0.123	0	-	-	64156	995998	1008466	819319	9523
	15	-	-	26.689	0.017	0.209	0.352	0.063	0	-	-	64156	502	198540	819067	9293
3.ttl	1	0.009	0.009	0.009	0.01	0.007	0.008	0.008	0	0	0	0	47	0	0	0
	2	0.019	0.017	0.104	0.111	0.02	0.017	0.016	57	57	57	256813	257974	114	57	57
	3	0.11	0.135	0.233	3.244	0.274	3.109	0.113	1462	1462	1462	258218	6668549	258218	6381394	1462
	4	0.038	0.034	0.277	0.034	0.026	0.027	0.028	0	0	0	256803	1314	161	104	104
	5	0.036	0.036	0.275	0.025	0.024	0.031	0.027	0	0	0	256803	47	161	57	57
	6	0.063	0.056	0.663	0.128	0.298	3.122	0.133	0	0	0	256803	257974	258379	6381451	1576
	7	0.061	0.062	0.709	0.032	0.287	3.101	0.132	0	0	0	256803	141	258369	6381498	1566
	8	0.094	0.153	1.425	0.138	0.297	0.03	0.027	0	0	0	256803	261905	516433	161	275
	9	0.098	0.15	1.819	0.037	0.156	0.03	0.027	0	0	0	256803	235	258660	104	208
	10	0.143	0.399	2.478	0.15	0.148	0.049	0.048	0	0	0	256803	265836	259504	5473	5634
	11	0.141	0.368	12.374	3.343	3.315	3.397	0.384	0	0	0	256803	6866425	6670079	6452693	72865
	12	0.21	1.136	15.915	0.051	0.576	3.133	0.171	0	0	0	256803	1326	773580	6382718	1730
	13	-	-	35.05	0.194	0.623	3.521	0.341	0	-	-	256803	329484	776449	6451652	74135
	14	-	-	399.257	3.948	3.982	3.344	0.558	0	-	-	256803	8461907	7190771	6463879	74581
	15	-	-	388.289	0.06	1.34	3.213	0.378	0	-	-	256803	1377	803755	6452448	73026
4.ttl	1	0.026	0.025	0.025	0.039	0.025	0.024	0.025	0	0	0	0	253	0	0	0
	2	0.064	0.069	0.471	0.522	0.064	0.07	0.064	248	248	248	1027022	1040241	496	248	248
	3	0.929	0.938	1.404	28.325	1.857	28.103	0.945	12651	12651	12651	1039425	51050537	1039425	49377537	12651
	4	0.198	0.173	1.617	0.157	0.095	0.129	0.135	4	4	4	1027031	13800	753	505	505
	5	0.182	0.174	1.617	0.144	0.094	0.143	0.138	0	0	0	1027027	253	749	248	248
	6	0.327	0.312	4.729	0.64	1.913	28.148	1	4	4	4	1027031	1040479	1040182	49377789	13151
	7	0.308	0.325	4.721	0.222	1.98	27.908	1.106	0	0	0	1027027	759	1040183	49378038	13152
	8	0.504	0.778	9.217	0.675	1.278	0.158	0.129	4	4	236	1027031	1079963	2080575	757	1249
	9	0.522	0.835	12.456	0.266	0.705	0.14	0.131	0	0	4	1027027	1265	1041493	505	1002
	10	0.782	2.174	15.698	0.738	0.66	0.288	0.253	4	4	236	1027031	1119447	1055223	52295	53040
	11	0.76	2.077	93.286	30.477	30.641	29.507	3.476	232	0	4	1027259	54927712	51065747	50689528	1325139
	12	1.083	6.03	114.063	0.354	3.362	28.046	1.329	4	4	236	1027031	15222	3107857	49391554	14314
	13	-	-	253.131	1.64	3.442	30.217	3.913	232	-	-	1027259	2499217	3173640	50730474	1353321
	14	-	-	>999	74.607	35.483	30.531	5.52	10972	-	-	-	117902759	53931133	52556376	1368984
	15	-	-	>999	0.454	10.929	29.497	3.763	1	-	-	-	35953	3754770	50690218	1326126

using magic sets or optimising the program before executions, the performance naturally depends on how we split the query into subqueries in the rewriting algorithm.

In the paper, we described three simple complexity-motivated splitting strategies. Our experiments show that none of them is always the best and the execution time may be dramatically improved by using an ‘adaptable’ splitting strategy which would work similarly to a query execution planner in database management systems and use statistical information about the data to generate a quickly executable NDL program.

The difference in performance between different types of optimal rewritings made us investigate its causes. For example, we noticed that the TW-rewriting of the query with 3 atoms of Sequence 3

$$\begin{aligned} G(x, y) &\leftarrow S(x, z) \wedge P_{13}(z, y), \\ P_{13}(x, y) &\leftarrow R(x, z) \wedge R(z, y), \\ G(x, y) &\leftarrow A_P(x) \wedge R(x, y) \end{aligned}$$

takes as long as 28 seconds to execute on the fourth dataset because it needs so much time to materialise  $P_{13}$ , which has around  $6 \cdot 10^6$  triples. On the other hand, if we remove this predicate by substituting its definition into the first rule, we obtain the rewriting

$$\begin{aligned} G(x, y) &\leftarrow S(x, z) \wedge R(x, v) \wedge R(v, y), \\ G(x, y) &\leftarrow A_P(x) \wedge R(x, y), \end{aligned}$$

which is executed in 0.945 seconds. This substitution could be done automatically by a clever NDL engine, but not performed by RDBFox. Thus, we made an attempt to ‘improve’ the TW-rewriting by getting rid in this fashion of all predicates that are defined by a single rule and occur not more than twice in the bodies of the rules. However, though the rewriting  $TW^*$  thus obtained shows a much better performance on Sequences 1 and 3 (see Tables 3 and 5), it is not always so on Sequence 2 (Table 4). This observation suggests that our rewriting could be executed faster on a more advanced NDL engine than RDBFox which would carry out such substitutions depending on the cardinality of EDBs.

## References

- [1] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [2] M. Arenas, P. Barceló, L. Libkin, and F. Murlak. *Foundations of Data Exchange*. Cambridge University Press, 2014.
- [3] S. Arora and B. Barak. *Computational Complexity: A Modern Approach*. Cambridge University Press, New York, NY, USA, 1st edition, 2009.
- [4] F. Baader, D. Calvanese, D. McGuinness, D. Nardi, and P. Patel-Schneider, editors. *The Description Logic Handbook: Theory, Implementation and Applications*. Cambridge University Press, 2003.
- [5] M. Bienvenu, S. Kikot, and V. V. Podolskii. 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*, pages 317–328. IEEE Computer Society, 2015.
- [6] M. Bienvenu, M. Ortiz, M. Simkus, and G. Xiao. Tractable queries for lightweight description logics. In *Proc. of the 23rd Int. Joint Conf. on Artificial Intelligence (IJCAI 2013)*, pages 768–774. IJCAI/AAAI, 2013.
- [7] D. Bursztyn, F. Goasdoué, and I. Manolescu. Teaching an RDBMS about ontological constraints. *PVLDB*, 9(12):1161–1172, 2016.

- [8] M. Calautti, G. Gottlob, and A. Pieris. Chase termination for guarded existential rules. In *Proc. of the 34th ACM Symposium on Principles of Database Systems, PODS 2015*, pages 91–103, 2015.
- [9] A. Calì, G. Gottlob, and T. Lukasiewicz. A general datalog-based framework for tractable query answering over ontologies. *Journal of Web Semantics*, 14:57–83, 2012.
- [10] A. Calì, G. Gottlob, and A. Pieris. Towards more expressive ontology languages: The query answering problem. *Artificial Intelligence*, 193:87–128, 2012.
- [11] D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, A. Poggi, M. Rodriguez-Muro, R. Rosati, M. Ruzzi, and D. F. Savo. The MASTRO system for ontology-based data access. *Semantic Web*, 2(1):43–53, 2011.
- [12] D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati. Tractable reasoning and efficient query answering in description logics: the *DL-Lite* family. *Journal of Automated Reasoning*, 39(3):385–429, 2007.
- [13] C. Chekuri and A. Rajaraman. Conjunctive query containment revisited. *Theoretical Computer Science*, 239(2):211–229, 2000.
- [14] A. Chortaras, D. Trivela, and G. Stamou. Optimized query rewriting for OWL 2 QL. In *Proc. of CADE-23*, volume 6803 of *LNCs*, pages 192–206. Springer, 2011.
- [15] S. A. Cook. Characterizations of pushdown machines in terms of time-bounded computers. *Journal of the ACM*, 18(1):4–18, 1971.
- [16] B. Cuenca Grau, I. Horrocks, M. Krötzsch, C. Kupke, D. Magka, B. Motik, and Z. Wang. Acyclicity notions for existential rules and their application to query answering in ontologies. *Journal of Artificial Intelligence Research (JAIR)*, 47:741–808, 2013.
- [17] E. Dantsin, T. Eiter, G. Gottlob, and A. Voronkov. Complexity and expressive power of logic programming. *ACM Computing Surveys*, 33(3):374–425, 2001.
- [18] F. Di Pinto, D. Lembo, M. Lenzerini, R. Mancini, A. Poggi, R. Rosati, M. Ruzzi, and D. F. Savo. Optimizing query rewriting in ontology-based data access. In *Proc. of the 16th Int. Conf. on Extending Database Technology (EDBT 2013)*, pages 561–572. ACM, 2013.
- [19] A. Doan, A. Y. Halevy, and Z. G. Ives. *Principles of Data Integration*. Morgan Kaufmann, 2012.
- [20] T. Eiter, M. Ortiz, M. Šimkus, T.-K. Tran, and G. Xiao. Query rewriting for Horn-SHIQ plus rules. In *Proc. of the 26th AAAI Conf. on Artificial Intelligence (AAAI 2012)*, pages 726–733. AAAI, 2012.
- [21] M. R. Fellows, D. Hermelin, F. A. Rosamond, and S. Vialette. On the parameterized complexity of multiple-interval graph problems. *Theoretical Computer Science*, 410(1):53–61, 2009.
- [22] J. Flum and M. Grohe. *Parameterized Complexity Theory*. Texts in Theoretical Computer Science. An EATCS Series. Springer, 2006.
- [23] M. Giese, A. Soylu, G. Vega-Gorgojo, A. Waaler, P. Haase, E. Jiménez-Ruiz, D. Lanti, M. Rezk, G. Xiao, Ö. Özçep, and R. Rosati. Optique: Zooming in on big data. *IEEE Computer*, 48(3):60–67, 2015.
- [24] T. Gogacz and J. Marcinkowski. All-instances termination of chase is undecidable. In *Proc. of the 41st Int. Colloquium Automata, Languages, and Programming (ICALP 2014), Part II*, volume 8573 of *Lecture Notes in Computer Science*, pages 293–304. Springer, 2014.

- [25] G. Gottlob, S. Kikot, R. Kontchakov, V. V. Podolskii, T. Schwentick, and M. Zakharyashev. The price of query rewriting in ontology-based data access. *Artificial Intelligence*, 213:42–59, 2014.
- [26] G. Gottlob, N. Leone, and F. Scarcello. Computing LOGCFL certificates. In *Proc. of the 26th Int. Colloquium on Automata, Languages and Programming (ICALP-99)*, volume 1644 of *Lecture Notes in Computer Science*, pages 361–371. Springer, 1999.
- [27] G. Gottlob, G. Orsi, and A. Pieris. Ontological queries: Rewriting and optimization. In *Proc. of ICDE 2011*, pages 2–13. IEEE Computer Society, 2011.
- [28] G. Gottlob, G. Orsi, and A. Pieris. Query rewriting and optimization for ontological databases. *ACM Transactions on Database Systems (TODS)*, 39(3):25, 2014.
- [29] S. A. Greibach. The hardest context-free language. *SIAM J. Comput.*, 2(4):304–310, 1973.
- [30] D. A. Huffman. A method for the construction of minimum-redundancy codes. *Proceedings of the Institute of Radio Engineers*, 40(9):1098–1101, 1952.
- [31] E. Jiménez-Ruiz, E. Kharlamov, D. Zheleznyakov, I. Horrocks, C. Pinkel, M. G. Skjæveland, E. Thorstensen, and J. Mora. BootOX: Bootstrapping OWL 2 ontologies and R2RML mappings from relational databases. In *Proc. of the ISWC 2015 Posters & Demonstrations Track at the 14th Int. Semantic Web Conf. (ISWC-2015)*, volume 1486 of *CEUR Workshop Proceedings*. CEUR-WS, 2015.
- [32] M. Kaminski, Y. Nenov, and B. Cuenca Grau. Datalog rewritability of Disjunctive Datalog programs and non-Horn ontologies. *Artificial Intelligence*, 236:90–118, 2016.
- [33] E. Kharlamov, D. Hovland, E. Jiménez-Ruiz, D. Lanti, H. Lie, C. Pinkel, M. Rezk, M. G. Skjæveland, E. Thorstensen, G. Xiao, D. Zheleznyakov, and I. Horrocks. Ontology based access to exploration data at Statoil. In *Proc. of the 14th Int. Semantic Web Conf. (ISWC 2015), Part II*, volume 9367 of *Lecture Notes in Computer Science*, pages 93–112. Springer, 2015.
- [34] S. Kikot, R. Kontchakov, V. Podolskii, and M. Zakharyashev. 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)*, pages 57:1–57:10. ACM, 2014.
- [35] S. Kikot, R. Kontchakov, V. V. Podolskii, and M. Zakharyashev. Exponential lower bounds and separation for query rewriting. In *Proc. of the 39th Int. Colloquium on Automata, Languages and Programming (ICALP 2012)*, volume 7392 of *Lecture Notes in Computer Science*, pages 263–274. Springer, 2012.
- [36] S. Kikot, R. Kontchakov, and M. Zakharyashev. On (in)tractability of OBDA with OWL 2 QL. In *Proc. of the 24th Int. Workshop on Description Logics (DL 2011)*, volume 745, pages 224–234. CEUR-WS, 2011.
- [37] S. Kikot, R. Kontchakov, and M. Zakharyashev. Conjunctive query answering with OWL 2 QL. In *Proc. of the 13th Int. Conf. on Principles of Knowledge Representation and Reasoning (KR 2012)*, pages 275–285. AAAI, 2012.
- [38] C. Koch. Processing queries on tree-structured data efficiently. In *Proc. of the 25th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS 2006)*, pages 213–224. ACM, 2006.
- [39] M. König, M. Leclère, M.-L. Mugnier, and M. Thomazo. Sound, complete and minimal UCQ-rewriting for existential rules. *Semantic Web*, 6(5):451–475, 2015.

- [40] R. Kontchakov, C. Lutz, D. Toman, F. Wolter, and M. Zakharyashev. The combined approach to query answering in DL-Lite. In *Proc. of the 12th Int. Conf. on Principles of Knowledge Representation and Reasoning (KR 2010)*, pages 247–257. AAAI Press, 2010.
- [41] R. Kontchakov, M. Rezk, M. Rodriguez-Muro, G. Xiao, and M. Zakharyashev. Answering SPARQL queries over databases under OWL 2 QL entailment regime. In *Proc. of the 13th Int. Semantic Web Conf. (ISWC 2014), Part I*, volume 8796 of *Lecture Notes in Computer Science*, pages 552–567. Springer, 2014.
- [42] M. Lenzerini. Ontology-based data management. *ACM SIGMOD Blog*, May 2013.
- [43] J. Mora, R. Rosati, and Ó. Corcho. Kyrie2: query rewriting under extensional constraints in ELHIO. In *Proc. of the 13th Int. Semantic Web Conf. (ISWC 2014)*, volume 8796 of *Lecture Notes in Computer Science*, pages 568–583. Springer, 2014.
- [44] B. Motik, B. Cuenca Grau, I. Horrocks, Z. Wu, A. Fokoue, and C. Lutz. *OWL 2 Web Ontology Language Profiles*. W3C Recommendation, 2012. Available at <http://www.w3.org/TR/owl2-profiles/>.
- [45] Y. Nenov, R. Piro, B. Motik, I. Horrocks, Z. Wu, and J. Banerjee. RDFox: A highly-scalable RDF store. In *Proc. of the 14th Int. Semantic Web Conf. (ISWC 2015), Part II*, volume 9367 of *Lecture Notes in Computer Science*, pages 3–20. Springer, 2015.
- [46] H. Pérez-Urbina, B. Motik, and I. Horrocks. A comparison of query rewriting techniques for DL-Lite. In *Proc. of the 22nd Int. Workshop on Description Logics (DL 2009)*, volume 477 of *CEUR Workshop Proceedings*. CEUR-WS, 2009.
- [47] H. Pérez-Urbina, E. Rodríguez-Díaz, M. Grove, G. Konstantinidis, and E. Sirin. Evaluation of query rewriting approaches for OWL 2. In *Proc. of SSWS+HPCSW 2012*, volume 943 of *CEUR Workshop Proceedings*. CEUR-WS, 2012.
- [48] F. Picalausa and S. Vansumneren. What are real SPARQL queries like? In *Proc. of the Int. Workshop on Semantic Web Information Management (SWIM)*. ACM, 2011.
- [49] A. Poggi, D. Lembo, D. Calvanese, G. De Giacomo, M. Lenzerini, and R. Rosati. Linking data to ontologies. *Journal on Data Semantics*, X:133–173, 2008.
- [50] M. Rodriguez-Muro, R. Kontchakov, and M. Zakharyashev. Ontology-based data access: Ontop of databases. In *Proc. of the 12th Int. Semantic Web Conf. (ISWC 2013), Part I*, volume 8218 of *Lecture Notes in Computer Science*, pages 558–573. Springer, 2013.
- [51] M. Rodriguez-Muro, R. Kontchakov, and M. Zakharyashev. Query rewriting and optimisation with database dependencies in Ontop. In *Informal Proc. of the 26th Int. Workshop on Description Logics (DL 2013)*, volume 1014 of *CEUR Workshop Proceedings*, pages 917–929. CEUR-WS, 2013.
- [52] R. Rosati. Prexto: Query rewriting under extensional constraints in DL-Lite. In *Proc. of the 9th Extended Semantic Web Conf. (EWSW 2012)*, volume 7295 of *Lecture Notes in Computer Science*, pages 360–374. Springer, 2012.
- [53] R. Rosati and A. Almatelli. Improving query answering over DL-Lite ontologies. In *Proc. of the 12th Int. Conf. on Principles of Knowledge Representation and Reasoning (KR 2010)*, pages 290–300. AAAI Press, 2010.
- [54] J. F. Sequeda, M. Arenas, and D. P. Miranker. OBDA: query rewriting or materialization? In practice, both! In *Proc. of the 13th Int. Semantic Web Conf. (ISWC 2014), Part I*, volume 8796 of *Lecture Notes in Computer Science*, pages 535–551. Springer, 2014.



- [55] A. Soylu, M. Giese, E. Jimenez-Ruiz, G. Vega-Gorgojo, and I. Horrocks. Experiencing optiquevqs: A multi-paradigm and ontology-based visual query system for end users. *Universal Access in the Information Society*, 15(1):129–152, 2016.
- [56] I. H. Sudborough. A note on tape-bounded complexity classes and linear context-free languages. *Journal of the ACM*, 22(4):499–500, Oct. 1975.
- [57] I. H. Sudborough. On the tape complexity of deterministic context-free languages. *Journal of the ACM*, 25(3):405–414, 1978.
- [58] M. Thomazo. Compact rewritings for existential rules. In *Proc. of the 23rd Int. Joint Conf. on Artificial Intelligence (IJCAI 2013)*. IJCAI/AAAI, 2013.
- [59] T. Venetis, G. Stoilos, and V. Vassalos. Rewriting minimisations for efficient ontology-based query answering. In *Proc. of the 28th Int. Conf. on Tools with Artificial Intelligence (ICTAI 2016)*, pages 1095–1102. IEEE, 2016.
- [60] H. Venkateswaran. Properties that characterize LOGCFL. *Journal of Computer and System Sciences*, 43(2):380–404, 1991.
- [61] M. Yannakakis. Algorithms for acyclic database schemes. In *Proc. of the 7th Int. Conf. on Very Large Data Bases (VLDB)*, pages 82–94. IEEE Computer Society, 1981.