

# XPath for DL-Lite Ontologies

Egor V. Kostylev<sup>1</sup>, Juan L. Reutter<sup>2</sup>, and Domagoj Vrgoč<sup>3</sup>

<sup>1</sup> University of Oxford [egor.kostylev@cs.ox.ac.uk](mailto:egor.kostylev@cs.ox.ac.uk)

<sup>2</sup> PUC Chile [jreutter@ing.puc.cl](mailto:jreutter@ing.puc.cl)

<sup>3</sup> University of Edinburgh [domagoj.vrgoc@ed.ac.uk](mailto:domagoj.vrgoc@ed.ac.uk)

**Abstract.** Applications of description logics (DLs) such as OWL 2 and ontology-based data access (OBDA) require understanding of how to pose database queries over DL knowledge bases. While there have been many studies regarding traditional relational query formalisms such as conjunctive queries and their extensions, little attention has been paid to graph database queries, despite the fact that graph databases share the structure of interpretations with DLs; that is they describe essentially the same objects. In particular, not much is known about the interplay between DLs and XPath. The last is a powerful formalism for querying semistructured data: it is in the core of most practical query languages for XML trees, and it is also gaining popularity in theory and practice of graph databases. In this paper we make a step towards coupling knowledge bases and graph databases by studying how to answer powerful XPath-style queries over DL-Lite. We start with adapting the definition of XPath to the DL context, and then proceed to study the complexity of evaluating XPath queries over knowledge bases. Results show that, while query answering is undecidable for the full XPath, by carefully tuning the amount of negation allowed in the queries we can arrive to XPath fragments that have a potential to be used in practical applications.

## 1 Introduction

Satisfiability and model checking have long been two central problems in the knowledge representation community. But new applications of description logics (DLs for short), such as the Web Ontology Language (OWL) [21] and ontology-based data access (OBDA), are forcing us to develop algorithms solving more complex data manipulation and extraction tasks, and in particular, require understanding how to answer database-style queries over knowledge bases that are specified by DLs [14].

The literature on knowledge bases usually considers relational queries, mostly focusing on conjunctive queries (CQs) and their extensions with union [1, 7], forms of negation [15, 24], or aggregates [10, 18]. It is also natural to look at queries that are designed to operate over graph databases, since these share the same structure with DL knowledge bases; namely they use unary and binary predicates (that is concepts and roles in DL terminology) to represent graph nodes and edges. While the idea of using graph queries in knowledge bases is not new (see e.g., [8]), we are only starting to understand how such queries can be fine tuned for various ontology languages. So far the specific research on this topic has primarily been concerned with the class of *regular path queries* (RPQs; see e.g. [2]), one of the most basic graph query languages. We now know how to evaluate such queries over various description logics [9], how to deal with

some of their extensions [4, 5] and understand how to solve their containment problem in the presence of DL constraints [11].

There are of course many other languages for querying graph and semi-structured data, most notable amongst them being XPath. Originally designed to extract data from XML trees [27], XPath has recently been adapted to work over graph databases [19] and was shown to retain good evaluation properties while at the same time being more powerful than RPQs and many of their extensions. Moreover, XPath also subsumes navigational graph querying features of SPARQL [17], such as property paths, and other commercial graph query languages (see e.g. Neo4j [23]), and thus can also be seen as a querying primitive for more powerful languages. Given that XPath is capable of expressing virtually all relevant querying primitives for graphs, and was tried and tested by XML practitioners, it is natural to see if it will have the same impact as a language for DL knowledge bases.

In this paper we take a decisive step towards the implementation of graph query languages over ontologies and study the problem of evaluating XPath queries over DL. To that extent, we first adapt the XPath query language to work over ontologies, arriving to DLXPath—an expressive language designed specifically for DL knowledge bases. To get a flavour of the interplay between DLXPath and DLs, we then study the problem of evaluating DLXPath queries over  $DL-Lite_{\mathcal{R}}$  and  $DL-Lite_{core}$  ontologies. The choice for this particular family of DLs is twofold: on the one hand, it is simple enough to start such a research; on the other, this family is of the most importance in practice, since it underlies OWL 2 QL profile [21].

As a simple useful example of an DLXPath query which is not expressible by RPQs or any other query formalisms explored in the context of knowledge bases, we can give the following: ‘Can I fly from city A to city B making stops only in cities with UNESCO World Heritage Sites which are endangered?’ Here we assume that the considered ontology has binary predicates (roles) `HasDirectFlight` and `HasUNESCOsite`, as well as unary predicate (concept) `InDanger`.

Regardless of the choice of a particular DL, one cannot expect that the study of DLXPath query answering will come as a straightforward adaptation of known techniques. The language of DLXPath contains expressive primitives such as transitive closure and negation, that is, the operators which are known to bring in conceptual difficulties even when studied in isolation. For example, it has been shown that adding negation to conjunctive queries immediately leads to undecidability of query answering [15], while adding path operators usually implies jumping one or two exponents in the complexity of query evaluation algorithms [9]. Thus, one of the challenges when studying DLXPath is to fine-tune these operators to archive an optimal trade-off between good evaluation properties and expressive power of the query language. To this end, we will distinguish two flavours of the language: the *core* fragment, denoted by  $DLXPath_{core}$ , and the *regular* fragment, denoted by  $DLXPath_{reg}$ . The two are designed to match the duality present in the standard XPath query language [27], and while the core fragment allows transitive closure only over basic role names, the regular fragment lifts this restriction, allowing us to post more general path queries. Regarding negation, we will distinguish full fragments, which allow both unary and binary negation, *path-positive* fragments with only the unary one, and *positive* fragments without any negation.

As usual when gauging the usefulness of a new query language for ontologies, we start by considering *data complexity* of the query answering problem, where one assumes that the query and the *terminological knowledge* (TBox) are fixed, while the only input is the *assertional knowledge* (ABox). Although we show that for the most general case the problem is undecidable, by limiting the amount of negation we obtain expressive languages whose queries can be answered in CONP and even NLOGSPACE, both over  $DL\text{-Lite}_{\mathcal{R}}$  and  $DL\text{-Lite}_{core}$ . We then move to studying the *combined complexity*, where both the knowledge base and the query form the input. Here we obtain bounds ranging from NP-complete (thus matching the ones for ordinary CQs), to EXPTIME-complete, and undecidable for the full language. We would like to note that some of the results are obtained using the deep connection between DLs, DLXPath and *propositional dynamic logic*, thus providing some interesting new techniques for answering queries over ontologies.

**Proofs** Due to the space restrictions, we only give sketches of short proofs and ideas of longer ones. Full details of the proofs can be found in the appendix.

## 2 Preliminaries

The language of  $DL\text{-Lite}_{\mathcal{R}}$  (and  $DL\text{-Lite}_{core}$ ) [1, 7] contains *individuals*  $c_1, c_2, \dots$ , *concept names*  $A_1, A_2, \dots$ , and *role names*  $P_1, P_2, \dots$ . *Concepts*  $B$  and *roles*  $R$  are defined by the grammar

$$B ::= A_i \mid \exists R, \quad R ::= P_j \mid P_j^-.$$

A  $DL\text{-Lite}_{\mathcal{R}}$  TBox is a finite set of *concept* and *role inclusions* of the form

$$B_1 \sqsubseteq B_2, \quad B_1 \sqcap B_2 \sqsubseteq \perp, \quad R_1 \sqsubseteq R_2, \quad R_1 \sqcap R_2 \sqsubseteq \perp.$$

A  $DL\text{-Lite}_{core}$  TBox contains only concept inclusions. An ABox is a finite set of *assertions* of the form  $A_i(c_k)$  and  $P_j(c_k, c_\ell)$ . A *knowledge base* (KB) is a pair  $(\mathcal{T}, \mathcal{A})$ , where  $\mathcal{T}$  is a TBox and  $\mathcal{A}$  an ABox.

An *interpretation*  $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$  is a nonempty domain  $\Delta^{\mathcal{I}}$  of elements with an interpretation function  $\cdot^{\mathcal{I}}$  that assigns an element  $c_k^{\mathcal{I}} \in \Delta^{\mathcal{I}}$  to each individual  $c_k$ , a subset  $A_i^{\mathcal{I}}$  of  $\Delta^{\mathcal{I}}$  to each concept name  $A_i$ , and a binary relation  $P_j^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}}$  to each role name  $P_j$ . When dealing with  $DL\text{-Lite}$  it is usual to adopt the *unique name assumption* (UNA), and we do so here by requiring that  $c_k^{\mathcal{I}} \neq c_\ell^{\mathcal{I}}$ , for all individuals  $c_k \neq c_\ell$ . Our results, however, do not depend on UNA. The interpretation function  $\cdot^{\mathcal{I}}$  is extended to roles and concepts in the following standard way:

$$\begin{aligned} (\exists R)^{\mathcal{I}} &= \{d \in \Delta^{\mathcal{I}} \mid \text{there is } d' \in \Delta^{\mathcal{I}} \text{ with } (d, d') \in R^{\mathcal{I}}\}, \\ (P_j^-)^{\mathcal{I}} &= \{(d', d) \in \Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}} \mid (d, d') \in P_j^{\mathcal{I}}\}. \end{aligned}$$

The *satisfaction relation*  $\models$  for TBox inclusions and ABox assertions is also standard:

$$\begin{array}{lll}
\mathcal{I} \models B_1 \sqsubseteq B_2 & \text{iff} & B_1^{\mathcal{I}} \subseteq B_2^{\mathcal{I}}, \\
\mathcal{I} \models B_1 \sqcap B_2 \sqsubseteq \perp & \text{iff} & B_1^{\mathcal{I}} \cap B_2^{\mathcal{I}} = \emptyset, \\
\mathcal{I} \models R_1 \sqsubseteq R_2 & \text{iff} & R_1^{\mathcal{I}} \subseteq R_2^{\mathcal{I}}, \\
\mathcal{I} \models R_1 \sqcap R_2 \sqsubseteq \perp & \text{iff} & R_1^{\mathcal{I}} \cap R_2^{\mathcal{I}} = \emptyset, \\
\mathcal{I} \models A_i(c_k) & \text{iff} & c_k^{\mathcal{I}} \in A_i^{\mathcal{I}}, \\
\mathcal{I} \models P_j(c_k, c_\ell) & \text{iff} & (c_k^{\mathcal{I}}, c_\ell^{\mathcal{I}}) \in P_j^{\mathcal{I}}.
\end{array}$$

A KB  $\mathcal{K} = (\mathcal{T}, \mathcal{A})$  is *satisfiable* if there is an interpretation  $\mathcal{I}$  satisfying all inclusions of  $\mathcal{T}$  and assertions of  $\mathcal{A}$ . In this case we write  $\mathcal{I} \models \mathcal{K}$  and say that  $\mathcal{I}$  is a *model* of  $\mathcal{K}$ .

### 3 DLXPath: XPath for Knowledge Bases

As mentioned in the introduction, the family of *DL-Lite* was designed not only to keep satisfiability and model checking problems simple, but mainly to keep the complexity of conjunctive query answering the same as in the case of relational databases, while simultaneously maximising the expressive power of the ontological language. This allows to use *DL-Lite* as a foundation for practical data managing applications, such as OWL 2 QL and OBDA. However, this does not automatically mean that other useful query formalisms with good evaluation properties over databases also have good properties when posed over *DL-Lite* knowledge bases. Hence, each class of queries which can be useful in knowledge base applications requires a separate research on its computational properties.

In this paper we concentrate on an adaptation of XPath query language for XML trees to knowledge bases. Recently it was shown that (a version of) XPath can be successfully used for querying graph databases [19]. Every interpretation of a *DL-Lite* vocabulary can be seen as a graph, and hence every *DL-Lite* KB is an incomplete description of a graph. That is why we expect our adaptation DLXPath for querying knowledge bases to be useful in practical applications.

In what follows, we will consider several fragments of DLXPath. We start with DLXPath<sub>core</sub>, the fragment which corresponds to core XPath, the theoretical foundation of most practical languages for querying XML trees.<sup>4</sup>

**Definition 1.** Node formulas  $\varphi, \psi$  and path formulas  $\alpha, \beta$  of DLXPath<sub>core</sub> are expressions satisfying the grammar

$$\begin{array}{ll}
\varphi, \psi & ::= A \mid \neg\varphi \mid \varphi \wedge \psi \mid \varphi \vee \psi \mid \langle \alpha \rangle, \\
\alpha, \beta & ::= \varepsilon \mid R \mid [\varphi] \mid \alpha \cup \beta \mid \alpha \cdot \beta \mid \bar{\alpha} \mid R^+,
\end{array} \tag{1}$$

where  $A$  ranges over concept names and  $R$  ranges over roles (i.e., role names and their inverses).

Semantics  $\llbracket \cdot \rrbracket^{\mathcal{I}}$  of DLXPath<sub>core</sub> for an interpretation  $\mathcal{I}$  associates subsets of  $\Delta^{\mathcal{I}}$  to node formulas and binary relations on  $\Delta^{\mathcal{I}}$  to path formulas as given in Table 1.

<sup>4</sup> The subscript ‘core’ is used in this paper for two unrelated purposes. This matching is historical and accidental, but we decided to stay with conventional notation despite this undesired collision.

$$\begin{aligned}
\llbracket A \rrbracket^{\mathcal{I}} &= A^{\mathcal{I}} \\
\llbracket \neg\varphi \rrbracket^{\mathcal{I}} &= \Delta^{\mathcal{I}} \setminus \llbracket \varphi \rrbracket^{\mathcal{I}} \\
\llbracket \varphi \wedge \psi \rrbracket^{\mathcal{I}} &= \llbracket \varphi \rrbracket^{\mathcal{I}} \cap \llbracket \psi \rrbracket^{\mathcal{I}} \\
\llbracket \varphi \vee \psi \rrbracket^{\mathcal{I}} &= \llbracket \varphi \rrbracket^{\mathcal{I}} \cup \llbracket \psi \rrbracket^{\mathcal{I}} \\
\llbracket \langle \alpha \rangle \rrbracket^{\mathcal{I}} &= \{d \mid \text{there exists } d' \text{ such that } (d, d') \in \llbracket \alpha \rrbracket^{\mathcal{I}}\} \\
\llbracket \varepsilon \rrbracket^{\mathcal{I}} &= \{(d, d) \mid d \in \Delta^{\mathcal{I}}\} \\
\llbracket R \rrbracket^{\mathcal{I}} &= R^{\mathcal{I}} \\
\llbracket [\varphi] \rrbracket^{\mathcal{I}} &= \{(d, d) \mid d \in \llbracket \varphi \rrbracket^{\mathcal{I}}\} \\
\llbracket \alpha \cup \beta \rrbracket^{\mathcal{I}} &= \llbracket \alpha \rrbracket^{\mathcal{I}} \cup \llbracket \beta \rrbracket^{\mathcal{I}} \\
\llbracket \alpha \cdot \beta \rrbracket^{\mathcal{I}} &= \llbracket \alpha \rrbracket^{\mathcal{I}} \circ \llbracket \beta \rrbracket^{\mathcal{I}} \\
\llbracket \bar{\alpha} \rrbracket^{\mathcal{I}} &= (\Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}}) \setminus \llbracket \alpha \rrbracket^{\mathcal{I}} \\
\llbracket R^+ \rrbracket^{\mathcal{I}} &\text{ is the transitive closure of } R^{\mathcal{I}}
\end{aligned}$$

**Table 1.** Semantics of  $\text{DLXPath}_{\text{reg}}$ . The symbol ‘\’ stands for set-theoretic difference.

As usual when dealing with ontologies, our interest is not query answering for a particular interpretation, but computing those answers that are true in all possible models of the knowledge base. Formally, let  $\mathcal{K} = (\mathcal{T}, \mathcal{A})$  be a knowledge base and  $\alpha$  a  $\text{DLXPath}$  path formula. The *certain answers* of  $\alpha$  over  $\mathcal{K}$ , denoted  $\text{Certain}(\alpha, \mathcal{K})$ , is the set of all pairs  $(c_1, c_2)$  of individuals such that  $(c_1^{\mathcal{I}}, c_2^{\mathcal{I}}) \in \llbracket \alpha \rrbracket^{\mathcal{I}}$  for all models  $\mathcal{I}$  of  $\mathcal{K}$ . Similarly, one can define certain answers  $\text{Certain}(\varphi, \mathcal{K})$  for a  $\text{DLXPath}$  node formula  $\varphi$  as the set of all individuals  $c$  such that  $c^{\mathcal{I}} \in \llbracket \varphi \rrbracket^{\mathcal{I}}$ , for all  $\mathcal{I}$  models of  $\mathcal{K}$ . In the paper all of the results will be stated for path formulas (*queries* from here on), however, they remain unchanged for node formulas.

*Example 1.* Coming back to the example from the introduction, consider role names  $\text{HasDirectFlight}$ , which connects cities with a direct flight, and  $\text{HasUNESCOsite}$ , which connects cities with their UNESCO world heritage sites, as well as concept  $\text{InDanger}$  denoting that a particular site is endangered. Let KB  $\mathcal{K}$  represents the flight destination graph, as well as knowledge about heritage sites, partially explicitly in the ABox, and partially implicitly, by means of TBox inclusions. Then checking whether it is possible to fly from Edinburgh to a city that has an endangered UNESCO world heritage site is equivalent to checking if  $\text{Edinburgh}$  is in  $\text{Certain}(\varphi, \mathcal{K})$ , where  $\varphi$  is a node formula

$$\langle \text{HasDirectFlight}^+[\langle \text{HasUNESCOsite}[\text{InDanger}] \rangle] \rangle. \quad \square$$

As already noted, the formalism of core XPath is in the nutshell of the most practical query languages for XML trees. However, in [19] it was shown that the corresponding graph language cannot express certain properties that are deemed essential when querying graphs. Thus, besides  $\text{DLXPath}_{\text{core}}$  we also consider its generalisation called *regular DLXPath*, or  $\text{DLXPath}_{\text{reg}}$ , which extends the core fragment by allowing the use of transitive closure operator  $^+$  over arbitrary path formulas. Formally, path formulas of  $\text{DLXPath}_{\text{reg}}$  satisfy the grammar

$$\alpha, \beta := \varepsilon \mid R \mid [\varphi] \mid \alpha \cup \beta \mid \alpha \cdot \beta \mid \bar{\alpha} \mid \alpha^+,$$

while node formulas remain the same as for  $\text{DLXPath}_{\text{core}}$  in grammar (1). As expected, the semantics  $\llbracket \alpha^+ \rrbracket^{\mathcal{I}}$  over an interpretation  $\mathcal{I}$  is the transitive closure of  $\llbracket \alpha \rrbracket^{\mathcal{I}}$ .

*Example 2.* With full transitive closure we are now able to post more complex queries than in the core fragment. Consider again the knowledge base  $\mathcal{K}$  from Example 1. We can now ask if it is possible to fly from Liverpool to Jerusalem, making stopovers only in places that have an endangered UNESCO world heritage site by checking if the pair  $(\text{Liverpool}, \text{Jerusalem})$  is in  $\text{Certain}(\alpha, \mathcal{K})$ , where  $\alpha$  is a query

$$(\text{HasDirectFlight}[\langle \text{HasUNESCOsite}[\text{InDanger}] \rangle])^+.$$

Note that here we check if each of the cities along the path has an endangered site.  $\square$

Besides these two query languages, we will consider their fragments, which will be introduced as needed.

Before passing on to the complexity of  $\text{DLXPath}$  query evaluation, we briefly compare our languages with other formalisms. First,  $\text{DLXPath}$  is clearly incomparable with CQs. However, all tree-shaped CQs and unions of CQs with no more than two free variables can be written as  $\text{DLXPath}_{\text{core}}$  queries. On the other hand, all negation-free and  $^+$ -free  $\text{DLXPath}_{\text{reg}}$  queries can be written as unions of CQs, though with a cost of exponential blow-up. If we allow negation but only over concept and role names, then a query can be written as a union of CQs with safe negation. Second,  $\text{DLXPath}_{\text{reg}}$  queries without negation and node tests  $[\varphi]$  are *2-way regular path queries (2RPQs)*, the standard formalism for querying graph databases. If, additionally, role inverses are not allowed as path formulas, it becomes plain *RPQs*, that is, essentially, regular expressions. Finally,  $\text{DLXPath}_{\text{reg}}$  node formulas without path negation are nothing else but *propositional dynamic logic with converse (CPDL)* formulas; we will discuss this connection in more detail later on.

In the following sections we analyse the complexity of evaluating  $\text{DLXPath}$  queries over DL-Lite knowledge bases.

## 4 Data Complexity of $\text{DLXPath}$ Query Evaluation

As it is widely accepted in theory and proved in practice, the size of the query and TBox is usually much smaller than the size of the ABox (see e.g., [25] for discussion in the relational database context and [7] for DLs). This is why one usually considers *data complexity* of query answering, assuming that the TBox and the query are fixed, and only ABox is part of the input. In this section we study this problem for various fragments of  $\text{DLXPath}$ . Formally, let  $\mathcal{T}$  be a TBox and  $\alpha$  a  $\text{DLXPath}$  query. We are interested in the following family of problems.

CERTAIN ANSWERS  $(\alpha, \mathcal{T})$

**Input:** ABox  $\mathcal{A}$  and pair  $(c_1, c_2)$  of individuals

**Question:** Is  $(c_1, c_2) \in \text{Certain}(\alpha, (\mathcal{A}, \mathcal{T}))$ ?

As it was previously mentioned, data complexity of CQ query answering over *DL-Lite<sub>R</sub>* knowledge bases is the same as over relational databases, that is, in LOGSPACE [7]. At the first glance, one may expect a similar result in the case of DLXPath, where the complexity is NLOGSPACE-complete over graph databases [19]. However, the combination of the open world assumption and allowing negation in queries makes things quite different. In fact, the situation here is more like in the case of CQ with safe negation, where the complexity jump is dramatic: from polynomiality to undecidability [15].

**Theorem 1.** *There exists a DL-Lite<sub>core</sub> TBox  $\mathcal{T}$  and a DLXPath<sub>core</sub> query  $\alpha$  such that the problem CERTAIN ANSWERS  $(\alpha, \mathcal{T})$  is undecidable.*

The proof uses the same techniques as the proof of Theorem 1 in [15], that shows the undecidability of the problem of computing certain answers for conjunctive queries with safe negation over *DL-Lite<sub>core</sub>* knowledge bases.

Having this negative result, a natural direction is to search for DLXPath fragments that have decidable certain answers problem. Based on previous studies for XPath in XML and graph databases (see e.g., [3]), the most problematic primitive seems to be the negation  $\bar{\alpha}$  in path formulas. Indeed, we now show that removing this primitive from the syntax leads to decidability of the certain answers problem. We denote DLXPath<sub>core</sub><sup>path-pos</sup> the fragment of DLXPath<sub>core</sub> that does not allow the negation  $\bar{\alpha}$ , for  $\alpha$  a path formula, and define the class DLXPath<sub>reg</sub><sup>path-pos</sup> accordingly. We then have the following theorem.

**Theorem 2.** *There exists a DL-Lite<sub>core</sub> TBox  $\mathcal{T}$  and a DLXPath<sub>core</sub><sup>path-pos</sup> query  $\alpha$  such that the problem CERTAIN ANSWERS  $(\alpha, \mathcal{T})$  is CONP-hard. The problem is in CONP for any DL-Lite<sub>R</sub> TBox  $\mathcal{T}$  and DLXPath<sub>reg</sub><sup>path-pos</sup> query  $\alpha$ .*

*Proof (sketch).* We first sketch the CONP-hardness of the problem. Consider a language with role names  $P, T_1, T_2, T_3, F_1, F_2$ , and  $F_3$ , as well as concept name  $A$ . Fix a query

$$\alpha = \left[ P \left( \bigvee_{v_1 \neq u_1, v_2 \neq u_2, v_3 \neq u_3} \left( \bigwedge_{i=1,2,3} \langle S_i^{v_i} [A^{u_i}] \rangle \right) \right) \right],$$

where all  $v_i$  and  $u_i$  range over  $\{\top, \perp\}$ ;  $S_i^{v_i}$  is  $T_i$  if  $v_i = \top$  and  $F_i$  if  $v_i = \perp$ ;  $A^{u_i}$  is  $A$  if  $u_i = \top$  and  $\neg A$  if  $u_i = \perp$ .

Consider the complement of the NP-complete problem 3CNF-SAT whose input is a conjunction  $\Phi$  of clauses of the form  $\ell_1 \vee \ell_2 \vee \ell_3$ , where the  $\ell_i$  are literals, that is, propositional variables or their negations, and whose answer is ‘yes’ if  $\Phi$  is not satisfiable and ‘no’ otherwise.

For each variable  $x$  in  $\Phi$  the ABox  $\mathcal{A}$  uses an individual  $c_x$ . Also, for each clause  $\gamma$  in  $\Phi$  it uses an individual  $c_\gamma$ . Finally, it uses an individual  $c$ .

For each clause  $\gamma$  in  $\Phi$  the ABox  $\mathcal{A}$  contains the assertion  $P(c, c_\gamma)$ . Also, for each literal  $\ell_i$  in each clause  $\gamma = \ell_1 \vee \ell_2 \vee \ell_3$  of  $\Phi$  the ABox  $\mathcal{A}$  contains the assertion  $T_i(c_\gamma, c_x)$  if  $\ell_i = x$  and the assertion  $F_i(c_\gamma, c_x)$  if  $\ell_i = \bar{x}$ .

It is a technicality to show that  $(c, c)$  is a certain answer to  $\alpha$  over the KB  $(\mathcal{A}, \emptyset)$  if and only if  $\Phi$  is not satisfiable.

The CONP upper bound can be obtained by a careful analysis of the proof of Theorem 4. □

In fact, we can see that the reduction for hardness is done for the empty TBox, that is all the complexity is ‘contained’ in the formulation of the fixed query (which does not use the transitive closure operator  $^+$ ).

While this theorem looks positive in light of the general undecidability result, the complexity might still be too high for practical applications, as it leads to algorithms which run in exponential time in the size of the ABox. To lower the complexity and obtain tractable algorithms, one could also consider fragments of DLXPath that do not allow any form of negation, neither in node nor in path formulas. We denote such fragments of  $\text{DLXPath}_{\text{core}}$  and  $\text{DLXPath}_{\text{reg}}$  by  $\text{DLXPath}_{\text{core}}^{\text{pos}}$  and  $\text{DLXPath}_{\text{reg}}^{\text{pos}}$ , respectively. This last fragment is nothing else but *nested 2RPQs* [22], a language that has already been studied for *DL-Lite* knowledge bases [4], where an NLOGSPACE tight complexity bound was shown for the problem of computing certain answers for both *DL-Lite*<sub>core</sub> and *DL-Lite* <sub>$\mathcal{R}$</sub> . Furthermore, by adapting known results from graph databases, it is not difficult to extend this result to show that the problem is already NLOGSPACE hard even for  $\text{DLXPath}_{\text{core}}^{\text{pos}}$  queries (see e.g., [2] for a survey on such techniques).

Note that from this result we conclude that nesting and inverse in regular expressions, as well as fixed *DL-Lite* <sub>$\mathcal{R}$</sub>  TBoxes do not increase the tractable data complexity of query evaluation.

## 5 Combined complexity of DLXPath Query Evaluation

Even if data complexity is the most important measure in practice, *combined complexity*, that is complexity under the assumption that both TBox, ABox and the query are given as input, allows us to get a better understanding of the query answering problem, and often provides a blueprint of how to solve the problem in practice. That is why we continue our study in this direction. Formally, we consider the following family of problems, where  $X$  ranges over  $\{\text{core}, \mathcal{R}\}$ ,  $y$  over  $\{\text{core}, \text{reg}\}$ , and  $z$  is either nothing, or ‘path-pos’, or ‘pos’.

*DL-Lite* <sub>$X$</sub>  CERTAIN ANSWERS OVER  $\text{DLXPath}_y^z$   
**Input:** *DL-Lite* <sub>$X$</sub>  KB  $\mathcal{K}$ ,  $\text{DLXPath}_y^z$  query  $\alpha$ ,  
and pair  $(c_1, c_2)$  of individuals  
**Question:** Is  $(c_1, c_2) \in \text{Certain}(\alpha, \mathcal{K})$ ?

As it immediately follows from Theorem 1, the problem remains undecidable for full  $\text{DLXPath}_{\text{core}}$  and, hence, for full  $\text{DLXPath}_{\text{reg}}$ . However, the last result also follows from the connection of DLXPath with *propositional dynamic logic (PDL)* and some well known properties of PDL. Since this connection will be heavily used in the remainder of the paper we now discuss it in more detail.

As already mentioned, node formulas of  $\text{DLXPath}_{\text{reg}}^{\text{path-pos}}$  are, essentially, PDL with converse (CPDL) formulas (see [16] for a good introduction on the topic). Before continuing, we briefly explain these results using DL terminology. Note that in PDL community a different notation is used: for example, interpretations are called *Kripke structures*, concept names are called *propositional letters* or *variables*, role names are *atomic programs*, and inverse operator is *converse*.



More formally, the syntax of CPDL extends plain PDL with the inverse operator  $R^-$  on role names, and is thus equivalent to adding transitive closure and inverse to node formulas of  $\text{DLXPath}_{\text{reg}}^{\text{path-pos}}$  (PDL formulas follow the same syntax but are not allowed to use roles of form  $R^-$ ). The standard problem in the PDL community is the *satisfiability* of a node formula  $\varphi$ ; that is, checking whether there exists an interpretation  $\mathcal{I}$  and an element  $d$  in its domain such that  $\varphi$  holds in  $d$  (written,  $\mathcal{I}, d \models \varphi$ ).

Lutz et al. showed that the satisfiability of PDL formulas extended with arbitrary path negation ( $\text{PDL}^\neg$  for short) is undecidable [20], which already implies the undecidability of CERTAIN ANSWERS OVER  $\text{DLXPath}_{\text{reg}}$ : a  $\text{PDL}^\neg$  formula  $\varphi$  is satisfiable if and only if the  $\text{DLXPath}_{\text{reg}}$  query  $[\neg\varphi]$  has the certain answer  $(c, c)$  over the empty KB with individual  $c$  in the language. Note, however, that it does not immediately imply undecidability in data complexity shown in Theorem 1.

Turning our attention to the certain answers problem for  $\text{DLXPath}^{\text{path-pos}}$  queries, we again use the connection with the theory of PDL. It is well-known that the satisfiability problem for PDL is EXPTIME-complete [16]. It remains in EXPTIME even if these formulas are allowed to use the transitive closure operator  $^+$  over role names. The same holds for CPDL [16] and  $\text{PDL}^{(\neg)}$  – the extension of PDL which allows path negation in a limited form—only on role names [20]. Similarly to the previous undecidability result, the EXPTIME lower bounds for these classes of PDL formulas already imply EXPTIME-hardness of the certain answers problem for  $\text{DLXPath}_{\text{core}}^{\text{path-pos}}$ , even for empty KBs. Also, in [4] hardness was established even for  $\text{DLXPath}_{\text{reg}}^{\text{pos}}$ .

The upper bound is, however, much more challenging. To deal with  $\text{DL-Lite}_{\mathcal{R}}$  knowledge bases, in particular with role inclusions, we join the aforementioned extensions of PDL with inverse and negation on role names, and consider the language  $\text{CPDL}^{(\neg)}$  whose node formulas obey the same grammar (1) as PDL (and  $\text{DLXPath}$ ), and path formulas are defined as follows:

$$\alpha, \beta := \varepsilon \mid R \mid [\varphi] \mid \alpha \cup \beta \mid \alpha \cdot \beta \mid \overline{R} \mid \alpha^+.$$

We need the following result to establish complexity of the certain answers problem.

**Theorem 3.** *Checking satisfiability of a  $\text{CPDL}^{(\neg)}$  node formula can be done in EXPTIME.*

The proof makes use of ideas from [20] and [26]. Although it is not strictly related to description logics and knowledge representation, we state the result explicitly as we believe it might be of interest to the PDL community.

Of course, satisfiability results do not transfer directly to query answering over KBs. However, widening and recasting the ideas from [12] and [13], we obtain the desired upper bound.

**Lemma 1.** *The problem  $\text{DL-Lite}_{\mathcal{R}}$  CERTAIN ANSWERS OVER  $\text{DLXPath}_{\text{reg}}^{\text{path-pos}}$  is in EXPTIME.*

Summing up, we obtain the following theorem (the hardness results are from [16] and [4] and included just for completeness).

**Theorem 4.** *The problem  $DL\text{-Lite}_{\mathcal{R}}$  CERTAIN ANSWERS OVER  $DLXPath_{\text{reg}}^{\text{path-pos}}$  is EXPTIME-complete. It remains EXPTIME-hard for  $DLXPath_{\text{reg}}^{\text{pos}}$  queries with  $DL\text{-Lite}_{\text{core}}$  KBs, and for  $DLXPath_{\text{core}}^{\text{path-pos}}$  even with empty KBs.*

The only remaining fragment is that of  $DLXPath_{\text{core}}^{\text{pos}}$  queries, that is, the restriction of the  $DLXPath_{\text{core}}$  language that use neither binary nor unary negation. In the previous section we saw that data complexity of answering  $DLXPath^{\text{pos}}$  queries is the same as answering RPQs and 2RPQs, regardless of whether we made use of the regular or the core fragment. For combined complexity, the case is now different. We have already seen that query answering remains EXPTIME-hard for  $DLXPath^{\text{pos}}$ . We now show that the restriction to the core fragment limits the complexity by almost one exponential.

**Theorem 5.** *The problem  $DL\text{-Lite}_{\mathcal{R}}$  CERTAIN ANSWERS OVER  $DLXPath_{\text{core}}^{\text{pos}}$  is NP-complete. It remains NP-hard for  $DL\text{-Lite}_{\text{core}}$ .*

*Proof (sketch).* For designing an NP algorithm we rethink the ideas from the proof of Theorem 4.4 in [3], where a similar problem was solved in the context of XML. First, we show that a pair  $(c_1, c_2)$  is a certain answer to a query  $\alpha$  over a satisfiable KB  $\mathcal{K}$  if and only if  $(c_1, c_2) \in \llbracket \alpha \rrbracket^{Can(\mathcal{K})}$ , where  $Can(\mathcal{K})$  is the canonical model of  $\mathcal{K}$  (see e.g., [7] for definition). Second, if  $(c_1, c_2)$  is an answer, then this is witnessed by a sub-interpretation of  $Can(\mathcal{K})$  of polynomial size. Hence, the algorithm can just guess this sub-interpretation.

To show hardness we give a reduction from the NP-complete problem 3CNF-SAT. Suppose we are given a conjunction  $\Phi$  of clauses of the form  $\ell_1 \vee \ell_2 \vee \ell_3$ , where  $\ell_k$  are literals, that is propositional variables or their negations (we can assume that all literals in each clause are distinct). Let  $x_1, \dots, x_n$  be the variables in  $\Phi$ .

Consider the KB  $\mathcal{K}$  with the ABox  $\mathcal{A} = \{A(c)\}$  and the TBox  $\mathcal{T}$  consisting of the inclusions

$$\begin{aligned} A &\sqsubseteq \exists(X_1^v)^-, \text{ for } v \in \{\top, \perp\}, \\ \exists X_{i-1}^v &\sqsubseteq \exists(X_i^u)^-, \text{ for } v, u \in \{\top, \perp\}, 1 < i \leq n. \end{aligned}$$

For every variable  $x_i$  and truth value  $v \in \{\top, \perp\}$  let  $\psi_{x_i}^v$  be the node formula

$$\langle (X_n^\top \cup X_n^\perp) \cdots (X_{i+1}^\top \cup X_{i+1}^\perp) \cdot X_i^v \cdot (X_{i-1}^\top \cup X_{i-1}^\perp) \cdots (X_1^\top \cup X_1^\perp) \rangle.$$

Let  $\varphi$  be the node formula  $\langle \alpha[\psi] \rangle$ , where

$$\alpha = ((X_1^\top)^- \cup (X_1^\perp)^-) \cdot \dots \cdot ((X_n^\top)^- \cup (X_n^\perp)^-),$$

and  $\psi$  is a conjunction of node formulas of the following form, for each clause  $\gamma$  in  $\Phi$ , using variables  $x, y, z$ :

$$\bigvee_{\substack{(v_x, v_y, v_z) \text{ assignment of } (x, y, z) \\ \text{satisfying } \gamma}} \psi_x^{v_x} \wedge \psi_y^{v_y} \wedge \psi_z^{v_z}.$$

One can now show that  $(c, c)$  is a certain answer to  $[\varphi]$  over  $\mathcal{K}$  if and only if  $\Phi$  is satisfiable.  $\square$

Complexity	DLXPath <sup>pos</sup>	DLXPath <sup>path-pos</sup>	DLXPath
Data	NLOGSPACE-c	<b>CONP-c</b>	<b>undec.</b>
Combined	<b>NP-c*</b> / EXPTIME-c <sup>†</sup>	<b>EXPTIME-c</b>	<b>undec.</b>

**Table 2.** A summary of the complexity results. Here “-c” stands for “-complete” and “undec.” for “undecidable”. All the results hold for both  $DL-Lite_{core}$  and  $DL-Lite_{\mathcal{R}}$ , as well as for both  $DLXPath_{core}$  and  $DLXPath_{reg}$ ; the only exception is that \* holds just for  $DLXPath_{core}$  and <sup>†</sup> just for  $DLXPath_{reg}$ . The new results of this paper are set off in bold.

It is worth to mention, that the result holds even if the queries are not allowed to use the transitive closure operator <sup>+</sup> at all, being, essentially, very restricted form of unions of CQs but with the same complexity of query answering. Hence, the border between tractable and intractable combined complexity of the problem lies somewhere very close to here, leaving 2RPQs on one side and  $DLXPath_{core}^{pos}$  with CQs on the other.

## 6 Conclusions and Future Work

In this paper we conducted a detailed study about using the XPath language to query ontologies of the DL-Lite family. The results, summarised in Table 2, show that although the problem is generally undecidable, by limiting the amount of negation we can get decidable and even tractable fragments. The deep connection between XPath, DL and PDL allowed us to use ideas developed in other areas and gave insight on how query evaluation can be affected by certain aspects of the language. Although this connection was explored previously (see e.g., discussion in Chapters 13 and 14 in [6]), we believe that the time is ripe for a comprehensive survey describing how techniques from one area can be transferred to another and hope that the results of this paper can motivate such a survey.

From a practical point of view, we would like to find and test the classes of queries that are of particular practical interest and have either tractable general algorithms or reliable heuristics. Here we primarily want to tackle positive and path-positive fragments of XPath, as these queries were implemented and tested by the XML community, and many good heuristics have been developed over the years.

## References

1. A. Artale, D. Calvanese, R. Kontchakov, and M. Zakharyashev. The DL-Lite family and relations. *J. Artif. Intell. Res. (JAIR)*, 36:1–69, 2009.
2. P. Barceló. Querying graph databases. In *PODS*, pages 175–188, 2013.
3. M. Benedikt, W. Fan, and F. Geerts. XPath satisfiability in the presence of DTDs. *J. ACM*, 55(2), 2008.
4. M. Bienvenu, D. Calvanese, M. Ortiz, and M. Šimkus. Nested regular path queries in description logics. *arXiv preprint arXiv:1402.7122*, 2014.
5. M. Bienvenu, M. Ortiz, and M. Šimkus. Conjunctive regular path queries in lightweight description logics. In *IJCAI*, 2013.

6. P. Blackburn, J. F. A. K. v. Benthem, and F. Wolter. *Handbook of Modal Logic, Volume 3 (Studies in Logic and Practical Reasoning)*. Elsevier Science Inc., New York, NY, USA, 2006.
7. 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. *J. of Automated Reasoning*, 39(3):385–429, 2007.
8. D. Calvanese, G. De Giacomo, M. Lenzerini, and M. Vardi. Containment of conjunctive regular path queries with inverse. In *7th International Conference on Principles of Knowledge Representation and Reasoning (KR)*, pages 176–185, 2000.
9. D. Calvanese, T. Eiter, and M. Ortiz. Answering regular path queries in expressive description logics: An automata-theoretic approach. In *AAAI*, pages 391–396, 2007.
10. D. Calvanese, E. Kharlamov, W. Nutt, and C. Thorne. Aggregate queries over ontologies. In R. Elmasri, M. Doerr, M. Brochhausen, and H. Han, editors, *ONISW*, pages 97–104. ACM, 2008.
11. D. Calvanese, M. Ortiz, and M. Šimkus. Containment of regular path queries under description logic constraints. In *IJCAI*, pages 805–812, 2011.
12. G. De Giacomo and M. Lenzerini. Boosting the Correspondence between Description Logics and Propositional Dynamic Logics. In *AAAI*, pages 205–212, 1994.
13. G. De Giacomo and M. Lenzerini. TBox and ABox Reasoning in Expressive Description Logics. In *KR*, pages 316–327, 1996.
14. B. Glimm, C. Ogbuji, S. Hawke, I. Herman, B. Parsia, A. Polleres, and A. Seaborne. SPARQL 1.1 entailment regimes, 2013. W3C Recommendation 21 March 2013, <http://www.w3.org/TR/2013/REC-sparql11-entailment-20130321/>.
15. V. Gutiérrez-Basulto, Y. A. Ibáñez-García, R. Kontchakov, and E. V. Kostylev. Conjunctive queries with negation over dl-lite: A closer look. In *RR*, pages 109–122, 2013.
16. D. Harel, D. Kozen, and J. Tiuryn. *Dynamic Logic*. MIT Press, 2000.
17. S. Harris and A. Seaborne. Sparql 1.1 query language. *W3C working draft*, 14, 2010.
18. E. V. Kostylev and J. L. Reutter. Answering counting aggregate queries over ontologies of the DL-Lite family. In *AAAI*, 2013.
19. L. Libkin, W. Martens, and D. Vrgoč. Querying graph databases with XPath. In *ICDT*, pages 129–140, 2013.
20. C. Lutz and D. Walther. PDL with Negation of Atomic Programs. *Journal of Applied Non-Classical Logics*, 15(2):189–213, 2005.
21. B. Motik, B. Cuenca Grau, I. Horrocks, Z. Wu, A. Fokoue, and C. Lutz. OWL 2 Web Ontology Language Profiles (2nd Edition), 2012. W3C Recommendation.
22. J. Pérez, M. Arenas, and C. Gutierrez. nSPARQL: A navigational language for RDF. *Web Semantics: Science, Services and Agents on the World Wide Web*, 8(4):255 – 270, 2010.
23. I. Robinson, J. Webber, and E. Eifrem. *Graph databases*. ” O’Reilly Media, Inc.”, 2013.
24. R. Rosati. The limits of querying ontologies. In *Proc. of the 11th Int. Conf. on Database Theory (ICDT)*, volume 4353 of *LNCS*, pages 164–178. Springer, 2007.
25. M. Y. Vardi. The complexity of relational query languages (extended abstract). In *STOC*, pages 137–146, 1982.
26. M. Y. Vardi and P. Wolper. Automata-Theoretic Techniques for Modal Logics of Programs. *J. Comput. Syst. Sci.*, 32(2):183–221, 1986.
27. XML Path Language (XPath) 2.0 (Second Edition). [www.w3.org/TR/xpath20](http://www.w3.org/TR/xpath20), 2010.

## Appendix

In this appendix we give the proofs for all theorems of the paper. For presentational purposes of the paper, some of the proofs refer to constructions from other proofs described after, but of course there are no circular references.

**Theorem 1.** *There exists a DL-Lite<sub>core</sub> TBox  $\mathcal{T}$  and a DLXPath<sub>core</sub> query  $\alpha$  such that the problem CERTAIN ANSWERS  $(\alpha, \mathcal{T})$  is undecidable.*

*Proof.* In this proof we use intersection operator  $\cap$  for DLXPath path formulas. This does not increase the expressive power of the language, since the intersection  $\alpha_1 \cap \alpha_2$  can be written as  $\bar{\alpha}_1 \cup \bar{\alpha}_2$ .

The proof is by reduction of the halting problem for deterministic Turing machines. In particular, given a Turing machine  $M$ , we construct a TBox  $\mathcal{T}$  and a DLXPath<sub>core</sub> query  $\alpha$  such that  $M$  does not accept an input  $w$  encoded as an ABox  $\mathcal{A}_w$  and an element  $c_1$  iff  $(c_1, c_1) \in \text{Certain}(\alpha, (\mathcal{T}, \mathcal{A}_w))$ . (Note that  $\mathcal{T}$  and  $\alpha$  depend on  $M$  but not on  $w$ ). Applying this construction to a fixed deterministic *universal* Turing machine, that is, a machine that accepts its input  $w$  iff the Turing machine encoded by  $w$  accepts the empty input, we obtain the required undecidability result.

Let  $M = (T, Q, q_0, q_1, \delta)$  be a deterministic Turing machine, where  $T$  is an alphabet (containing the blank symbol  $\sqcup$ ),  $Q$  a set of states,  $q_0 \in Q$  and  $q_1 \in Q$  are an initial and an accepting state, respectively, and  $\delta: Q \times T \rightarrow Q \times T \times \{-1, 1\}$  is a transition function. Computations of  $M$  can be thought of as sequences of configurations, with each configuration determined by the contents of all (infinitely many) cells of the tape, the state and the head position. We are going to encode a computation by domain elements arranged, roughly speaking, into a two-dimensional grid: one dimension is the tape and the other is time.

More precisely, we use a role  $T$  to point to the representation of the next cell on the tape (within the same configuration) and a role  $S$  to point to the representation of the same cell in a successive configuration. Concepts  $C_a$ , for  $a \in T$ , encode the contents of cells in the sense that a domain element belongs to  $C_a$  if the respective cell contains the symbol  $a$ . We use concepts  $H_q$ , for  $q \in Q$ , to indicate both the position of the head and the current state: a domain element belongs to  $H_q$  if the respective cell is under the head and the machine is in state  $q$ . We also use a concept  $H_\emptyset$  to mark all other cells on the tape (that is, cells that are not under the head of the machine). Finally, roles  $P_{qa}$ , for  $q \in Q$  and  $a \in T$ , are used to encode transitions; concepts  $D_\sigma$  and roles  $T_{\emptyset\sigma}$ , for  $\sigma \in \{-1, +1\}$ , to propagate the no-head marker backwards and forwards along the tape; and role  $T_0$  to make sure the tape is initially blank beyond the input word.

Let  $\varphi$  be DLXPath<sub>core</sub> node formula

$$\langle S^* T^* (\alpha_1 \cup \alpha_2 \cup \alpha_3 \cup \alpha_4 \cup \alpha_5) \rangle,$$

where

$$\begin{aligned} \alpha_1 &= [S^- T S \cap \bar{T}], \\ \alpha_2 &= \bigcup_{\delta(q,a)=(q',a',\sigma)} [H_q \wedge C_a \wedge \langle S[T^\sigma \cap \bar{P}_{q'a'}] \rangle], \end{aligned}$$

$$\begin{aligned}
\alpha_3 &= \bigcup_{a \in \Gamma} [H_\emptyset \wedge C_a \wedge \langle S[\neg C_a] \rangle], \\
\alpha_4 &= \bigcup_{\sigma \in \{-1, +1\}} [D_\sigma \wedge \langle T^\sigma \cap \bar{T}_{\emptyset\sigma} \rangle], \\
\alpha_5 &= [T_0 \cap \bar{T}],
\end{aligned}$$

and  $T^\sigma(y, z)$  stands for  $T(y, z)$  if  $\sigma = +1$  and for  $T(z, y)$  if  $\sigma = -1$ . Consider also a TBox  $\mathcal{T}$  containing the following concept inclusions:

$$\begin{aligned}
\exists T &\sqsubseteq \exists S, & \exists T_0^- &\sqsubseteq \exists T_0 \cap C_-, \\
\exists P_{qa}^- &\sqsubseteq H_q, & \exists P_{qa} &\sqsubseteq C_a, & \text{for } q \in Q \text{ and } a \in \Gamma, \\
H_q &\sqsubseteq D_\sigma, & \exists T_{\emptyset\sigma}^- &\sqsubseteq D_\sigma \cap H_\emptyset, & \text{for } q \in Q \text{ and } \sigma \in \{-1, +1\}, \\
H_{q_1} &\sqsubseteq \perp.
\end{aligned}$$

For every input  $w = a_1 \dots a_n \in \Gamma^*$ , we take the following ABox  $\mathcal{A}_w$ :

$$H_{q_0}(c_1), \quad C_{a_i}(c_i) \text{ and } T(c_i, c_{i+1}), \text{ for } 1 \leq i \leq n, \quad T_0(c_n, c_{n+1}).$$

It is now a matter of technicality to show that  $(c_1, c_1) \in \text{Certain}([\varphi], (\mathcal{T}, \mathcal{A}_w))$  iff  $M$  does not accept  $w$ .  $\square$

**Theorem 2.** *There exists a DL-Lite<sub>core</sub> TBox  $\mathcal{T}$  and a DLXPath<sub>core</sub><sup>path-pos</sup> query  $\alpha$  such that the problem CERTAIN ANSWERS  $(\alpha, \mathcal{T})$  is CONP-hard. The problem is in CONP for any DL-Lite<sub>R</sub> TBox  $\mathcal{T}$  and DLXPath<sub>reg</sub><sup>path-pos</sup> query  $\alpha$ .*

*Proof.* Let us start with the CONP-hardness of the problem. The proof will use the empty TBox. Fix the following query (note that it does not use the transitive closure operator  $^+$ ):

$$\begin{aligned}
\alpha = [P & \\
& (\langle T_1[\neg A] \rangle \wedge \langle T_2[\neg A] \rangle \wedge \langle T_3[\neg A] \rangle) \vee \\
& (\langle T_1[\neg A] \rangle \wedge \langle T_2[\neg A] \rangle \wedge \langle F_3[A] \rangle) \vee \\
& (\langle T_1[\neg A] \rangle \wedge \langle F_2[A] \rangle \wedge \langle T_3[\neg A] \rangle) \vee \\
& (\langle T_1[\neg A] \rangle \wedge \langle F_2[A] \rangle \wedge \langle F_3[A] \rangle) \vee \\
& (\langle F_1[A] \rangle \wedge \langle T_2[\neg A] \rangle \wedge \langle T_3[\neg A] \rangle) \vee \\
& (\langle F_1[A] \rangle \wedge \langle T_2[\neg A] \rangle \wedge \langle F_3[A] \rangle) \vee \\
& (\langle F_1[A] \rangle \wedge \langle F_2[A] \rangle \wedge \langle T_3[\neg A] \rangle) \vee \\
& (\langle F_1[A] \rangle \wedge \langle F_2[A] \rangle \wedge \langle F_3[A] \rangle) \vee ] ,
\end{aligned}$$

where  $P, T_1, T_2, T_3, F_1, F_2, F_3$  are role names and  $A$  is a concept name.

Consider the complement of the NP-complete problem 3CNF-SAT which input is a conjunction  $\Phi$  of clauses of the form  $\ell_1 \vee \ell_2 \vee \ell_3$ , where the  $\ell_i$  are literals, that is, propositional variables or their negations, and which output is ‘yes’ if  $\Phi$  is not satisfiable and ‘no’ otherwise.

For each variable  $x$  in  $\Phi$  the ABox  $\mathcal{A}$  uses an individual  $c_x$ . Also, for each clause  $\gamma$  in  $\Phi$  it uses an individual  $c_\gamma$ . Finally, it uses an individual  $c$ .

For each clause  $\gamma$  in  $\Phi$  the ABox  $\mathcal{A}$  contains the assertion  $P(c, c_\gamma)$ . Also, for each literal  $\ell_i$  in each clause  $\gamma = \ell_1 \vee \ell_2 \vee \ell_3$  of  $\Phi$  the ABox  $\mathcal{A}$  contains the assertion  $T_i(c_\gamma, c_x)$  if  $\ell_i = x$  or the assertion  $F_i(c_\gamma, c_x)$  if  $\ell_i = \bar{x}$ .

It is a matter of technicality to show that the pair  $(c, c)$  is a certain answer to  $\alpha$  over the KB  $(\mathcal{A}, \emptyset)$  if and only if  $\Phi$  is not satisfiable.

Next we carefully analyse the algorithm in the proof of Theorem 4 (which is based on the proof of Theorem 3) to obtain an CONPdata complexity algorithm which decides the problem.

Looking at the proof of Theorem 3 note that there are two sources of exponentiality in the construction of the Büchi tree automaton  $\mathcal{B}_{\varphi_{\mathcal{K}, \alpha}}$ , for  $\varphi_{\mathcal{K}, \alpha}$  as in the proof of Theorem 4. First, the number of Hintikka sets  $\mathcal{H}_{\varphi_{\mathcal{K}, \alpha}}$  can be exponential, because it is roughly the number of all subsets of the Fisher-Ladner closure of the formula. Second, each state of the automaton has a component from  $2^{\mathcal{P}_{\square}(\varphi_{\mathcal{K}, \alpha})}$ , where  $\mathcal{P}_{\square}(\varphi_{\mathcal{K}, \alpha}) = \{[\alpha]\psi, [\beta]\eta \mid [\alpha]\psi, [\beta]\eta \in cl(\varphi_{\mathcal{K}, \alpha})\}$ .

To deal with the first source, we need the following claim. Let  $\ell$  be the fixed number of role inclusions in the TBox  $\mathcal{T}$ .

*Claim 5.1.* If the formula  $\varphi_{\mathcal{K}, \alpha}$  has a Hintikka tree then it has a Hintikka tree satisfying the following conditions:

1. if for an individual  $c_i$  and node formula  $\psi$  there exists a node  $x$  such that  $C_i \in T(x)^{(1)}$  and  $\psi \in T(x)^{(1)}$  then for any node  $y$  such that  $C_i \in T(y)^{(1)}$  it holds that  $\psi \in T(y)^{(1)}$ ;
2. for any node  $x$  such that  $C_i \notin T(x)^{(1)}$  for any individual  $c_i$ , there exist at most  $\ell$  roles  $R$  such that  $\langle R \rangle C_i \in T(x)^{(1)}$ .

This claim follows immediately from the construction in the proof of Proposition 11 in [20]. Having this claim at hand we conclude that when constructing the Büchi tree automaton  $\mathcal{B}_{\varphi_{\mathcal{K}, \alpha}}$  we may consider not all Hintikka sets, but only those that satisfy the conditions of the claim. Since the TBox and the query are fixed, we can guess which formulas hold and construct the automaton based on the polynomial number of Hintikka sets.

To deal with the second source, note that by the construction all the states on a run of the automaton have the same  $\mathcal{P}$ . Hence, again, it is enough to guess this  $\mathcal{P}$  and build only the corresponding part of the automaton.  $\square$

**Theorem 3.** *Checking satisfiability of a CPDL<sup>( $\neg$ )</sup> node formula can be done in EXPTIME.*

*Proof.* This prove goes the same lines as the proof of existence of EXPTIME algorithm for checking satisfiability of PDL<sup>( $\neg$ )</sup> node formula from [20].

In this proof we will use the shortcut  $\langle \alpha \rangle \varphi$  for  $\langle \alpha \cdot [\varphi] \rangle$  and  $[\alpha] \varphi$  for  $\neg \langle \alpha \rangle \neg \varphi$ .<sup>5</sup> In fact, we even replace  $\langle \alpha \rangle$  with  $\langle \alpha \rangle \varphi$  and  $[\alpha]$  with  $[\alpha] \varphi$  in the syntax of PDL node formulas, which clearly does not change the expressive power and conciseness of formulas.

<sup>5</sup> Note that the standard PDL syntax (e.g., in [16]) uses  $\varphi?$  instead of  $[\varphi]$ , and  $[\alpha]\varphi$  instead of  $[\alpha]\varphi$ . Since this work is influenced by the XPath query language we opted to stay faithful to its syntax and use it throughout the paper.

Hence, we can assume that all our node formulas are in the *negation normal form*, that is the unary negation appears only in front of concept names.

Let  $\Pi_0 = \{R, \neg R \mid R \text{ is a role name or inverse}\}$ . It will be convenient to look at any CPDL<sup>( $\neg$ )</sup> path formula  $\alpha$  not as regular expression over the alphabet  $\Pi_0 \cup \{[\varphi] \mid \varphi \text{ is a node formula}\}$  but as an equivalent nondeterministic finite state automata over this alphabet with the set of states  $Q_\alpha$ , initial state  $q_\alpha$ , transition function  $\delta_\alpha$  and set of final states  $F_\alpha$ . Given a state  $q$  in  $Q_\alpha$  denote  $\alpha_q$  the automata obtained from  $\alpha$  by declaring  $q$  the initial state.

*Fisher-Ladner closure*  $cl(\varphi)$  of a CPDL<sup>( $\neg$ )</sup> node formula  $\varphi$  is the smallest set closed under the following conditions:

- (C1)  $\varphi \in cl(\varphi)$ ,
- (C2) if  $\psi_1 \wedge \psi_2 \in cl(\varphi)$  or  $\psi_1 \vee \psi_2 \in cl(\varphi)$  then  $\psi_1, \psi_2 \in cl(\varphi)$ ,
- (C3) if  $\psi \in cl(\varphi)$  then  $\neg\psi \in cl(\varphi)$  (from here onwards we assume that  $\neg\neg\psi$  denotes  $\psi$ , as well as  $\neg\neg S$  denotes  $S$  for  $S \in \Pi_0$ ),
- (C4) if  $\langle\alpha\rangle\psi \in cl(\varphi)$  then  $\psi \in cl(\varphi)$ ,  $\psi' \in cl(\varphi)$  for all  $[\psi']$  in the active alphabet of  $\alpha$ , and  $\langle\alpha_q\rangle\psi \in cl(\varphi)$  for all  $q$  in  $Q_\alpha$ ,
- (C5) if  $\lceil\alpha\rceil\psi \in cl(\varphi)$  then  $\psi \in cl(\varphi)$ ,  $\psi' \in cl(\varphi)$  for all  $[\psi']$  in the active alphabet of  $\alpha$ , and  $\lceil\alpha_q\rceil\psi \in cl(\varphi)$  for all  $q$  in  $Q_\alpha$ .

It is straightforward to check that the size of the Fisher-Ladner closure of a node formula  $\varphi$  is linear in the length of  $\varphi$ .

Given a CPDL<sup>( $\neg$ )</sup> node formula  $\varphi$ , a subset  $\Psi$  of  $cl(\varphi)$  is a *Hintikka set* for  $\varphi$  if the following conditions hold:

- (H1) if  $\psi_1 \wedge \psi_2 \in \Psi$  then  $\psi_1 \in \Psi$  and  $\psi_2 \in \Psi$ ,
- (H2) if  $\psi_1 \vee \psi_2 \in \Psi$  then  $\psi_1 \in \Psi$  or  $\psi_2 \in \Psi$ ,
- (H3)  $\psi \in \Psi$  if and only if  $\neg\psi \notin \Psi$ ,
- (H4) if  $\lceil\alpha\rceil\psi \in \Psi$  and  $q_\alpha \in F_\alpha$  then  $\psi \in \Psi$ ,
- (H5) if  $\lceil\alpha\rceil\psi \in \Psi$  and  $q \in \delta_\alpha(q_\alpha, [\eta])$  then  $\neg\eta \in \Psi$  or  $\lceil\alpha_q\rceil\psi \in \Psi$ .

Denote  $\mathcal{H}_\varphi$  the set of all Hintikka sets for  $\varphi$ .

Given a CPDL<sup>( $\neg$ )</sup> node formula  $\varphi$ , let  $\epsilon_1, \dots, \epsilon_k$  be all the formulas of the form  $\langle\alpha\rangle\psi$  in  $cl(\varphi)$ , and let  $\Lambda_\varphi$  be the set of triples  $\mathcal{H}_\varphi \times (\Pi_0 \cup \{\perp\}) \times \{0, \dots, k\}$ . For every triple  $\lambda$  from  $\Lambda_\varphi$  denote  $\lambda^{(i)}$ ,  $1 \leq i \leq 3$ , the  $i$ -th component of  $\lambda$ . A  $(k+1)$ -tuple  $(\lambda, \lambda_1, \dots, \lambda_k)$  of triples from  $\Lambda_\varphi$  is a *matching* if and only if for all  $1 \leq i \leq k$  the following holds:

- (M1) if  $\epsilon_i = \langle\alpha\rangle\psi \in \lambda^{(1)}$  then there are  $\psi_1, \dots, \psi_n \in \lambda^{(1)}$ ,  $n \geq 0$ , such that
  - (a) either  $\delta_\alpha(q_\alpha, [\psi_1] \cdots [\psi_n]) \cap F_\alpha \neq \emptyset$ ,  $\psi \in \lambda^{(1)}$ ,  $\lambda_i^{(2)} = \perp$ , and  $\lambda_i^{(3)} = 0$ ,
  - (b) or there exists  $S \in \Pi_0$  and  $q \in Q_\alpha$  such that  $q \in \delta_\alpha(q_\alpha, [\psi_1] \cdots [\psi_n]S)$ ,  $\epsilon_j = \langle\alpha_q\rangle\psi \in \lambda_i^{(2)}$ , and  $\lambda_i^{(3)} = j$ ;
- (M2a) if  $\lceil\alpha\rceil\psi \in \lambda^{(1)}$ ,  $q \in Q_\alpha$  and  $S \in \Pi_0$  are such that  $q \in \delta_\alpha(q_\alpha, S)$  and  $S = \lambda_i^{(2)}$  then  $\lceil\alpha_q\rceil\psi \in \lambda_i^{(1)}$ ,
- (M2b) if  $\lceil\alpha\rceil\psi \in \lambda_i^{(1)}$ ,  $q \in Q_\alpha$  and  $S \in \Pi_0$  are such that  $q \in \delta_\alpha(q_\alpha, S)$  and  $S = (\lambda_i^{(2)})^-$  then  $\lceil\alpha_q\rceil\psi \in \lambda^{(1)}$ .



The condition (M2b) is the one which differs the construction from the construction in [20].

An *infinite  $k$ -ary  $M$ -tree* for a set  $M$  and number  $k$  is a mapping from  $\{1, \dots, k\}^*$  to  $M$ .

Given a  $\text{CPDL}^{(\neg)}$  node formula  $\varphi$  with  $\epsilon_1, \dots, \epsilon_k$  set all the formulas of the form  $\langle \alpha \rangle \psi$  in  $cl(\varphi)$ , a  $k$ -ary  $\Lambda_\varphi$ -tree  $T$  is a *Hintikka tree* for  $\varphi$  if and only if the following conditions hold for all nodes  $x, y \in \{1, \dots, k\}^*$ :

- (T1)  $\varphi \in T(\varepsilon)$ ,
- (T2) the tuple  $(T(x), T(x1), \dots, T(xk))$  is a matching,
- (T3) there is no  $\epsilon_i \in T(x)^{(1)}$  with  $\gamma_1 \gamma_2 \dots \in \{1, \dots, k\}^\omega$  such that  $\gamma_1 = i$  and  $\gamma_{\ell+1} = T(x\gamma_1 \dots \gamma_\ell)^{(3)}$  for all  $\ell \geq 1$ ,
- (T4) if  $\lceil \alpha \rceil \psi, \lceil \beta \rceil \eta \in T(x)^{(1)}$ ,  $S \in \Pi_0$ ,  $q_1 \in Q_\alpha$ , and  $q_2 \in Q_\beta$  are such that  $q_1 \in \delta_\alpha(q_\alpha, S)$ ,  $q_2 \in \delta_\beta(q_\beta, \neg S)$  and  $\lceil \alpha_{q_1} \rceil \psi \notin T(y)^{(1)}$  then  $\lceil \beta_{q_2} \rceil \eta \in T(y)^{(1)}$ .

The  $\text{PDL}^{(\neg)}$  version of the following claim is proved in [20] (Proposition 11). The proof carries verbatim to the case of  $\text{CPDL}^{(\neg)}$ .

*Claim 5.2.* A  $\text{CPDL}^{(\neg)}$  node formula is satisfiable if and only if it has a Hintikka tree.

Given a  $\text{CPDL}^{(\neg)}$  node formula  $\varphi$  with  $k$  formulas of the form  $\langle \alpha \rangle \psi$  in  $cl(\varphi)$ , denote  $\mathcal{P}_\square(\varphi) = \{\{\lceil \alpha \rceil \psi, \lceil \beta \rceil \eta\} \mid \lceil \alpha \rceil \psi, \lceil \beta \rceil \eta \in cl(\varphi)\}$ . Then the  $k$ -ary Büchi tree automata  $\mathcal{B}_\varphi$  over the alphabet  $\Lambda_\varphi$  is defined as follows:

1. the set  $Q$  of states is a subset of  $\Lambda_\varphi \times 2^{\mathcal{P}_\square(\varphi)} \times \{\emptyset, \uparrow\}$ , such that for each state  $((\Psi, S, j), p, s)$  the following conditions hold:
  - (a) if  $\lceil \alpha \rceil \psi, \lceil \beta \rceil \eta \in \Psi$  then  $\{\lceil \alpha \rceil \psi, \lceil \beta \rceil \eta\} \in p$ ,
  - (b) if  $\{\lceil \alpha \rceil \psi, \lceil \beta \rceil \eta\} \in p$ ,  $S \in \Pi_0$ ,  $q_1 \in \delta_\alpha(q_\alpha, S)$ ,  $q_2 \in \delta_\beta(q_\beta, \neg S)$ , and  $\lceil \alpha_{q_1} \rceil \psi \notin \Psi$  then  $\lceil \beta_{q_2} \rceil \eta \in \Psi$ ;
2. the set  $I$  of initial states is  $\{((\Psi, S, j), p, s) \mid \varphi \in \Psi \text{ and } s = \emptyset\}$ ;
3. the transition function  $\delta$  is defined as follows:  $((\lambda, p, s), (\Psi, S, j), (\lambda_1, p_1, s_1), \dots, (\lambda_k, p_k, s_k)) \in \delta$  if and only if for all  $1 \leq i \leq k$ 
  - (a)  $\lambda = (\Psi, S, j)$ ,
  - (b)  $p_i = p$ ,
  - (c)  $(\lambda, \lambda_1, \dots, \lambda_k)$  is a matching, and
  - (d) if either  $s = \emptyset$ ,  $\lambda_i^{(3)} \neq 0$  and  $\epsilon_i \in \Psi$ , or  $s = \uparrow$ ,  $\lambda^{(3)} = i$  and  $\lambda_i^{(3)} \neq 0$ , then  $s_i = \uparrow$ ; otherwise,  $s_i = \emptyset$ ;
4. the set  $F$  of accepting states is  $\{(\lambda, p, s) \mid s = \emptyset\}$ .

It is straightforward to check that the size of the automata  $\mathcal{B}_\varphi$  is exponential in the length of  $\varphi$ .

Similarly to Claim 6, the  $\text{PDL}^{(\neg)}$  version of the following claim is proved in [20] (Proposition 15). Again, the proof carries verbatim to the case of  $\text{CPDL}^{(\neg)}$ .

*Claim 5.3.* A tree  $T$  is a Hintikka tree for a  $\text{CPDL}^{(\neg)}$  node formula  $\varphi$  if and only if  $T$  is in the language of  $\mathcal{B}_\varphi$ .

Now the statement of the theorem follows from the fact that emptiness of a tree Büchi automaton can be checked in polynomial (quadratic) time [26] and the observation that the automaton obtained by this construction is exponential in the size of the formula.  $\square$

**Theorem 4.** *The problem DL-Lite<sub>R</sub> CERTAIN ANSWERS OVER DLXPath<sub>reg</sub><sup>path-pos</sup> is EXPTIME-complete. It remains EXPTIME-hard for DLXPath<sub>reg</sub><sup>pos</sup> queries with DL-Lite<sub>core</sub> KBs, and for DLXPath<sub>core</sub><sup>path-pos</sup> even with empty KBs.*

*Proof.* The EXPTIME-hardness immediately follows from the results of [16] and [4]. That is why next we concentrate on an algorithm.

The proof is by polynomial reduction to the satisfiability of CPDL<sup>( $\neg$ )</sup> formulas, which is in EXPTIME by Theorem 3. It adopts the ideas from [12] to our needs. In the proof we use  $\varphi \Rightarrow \psi$  as a shortcut for  $\neg\varphi \vee \psi$ , as well as the shortcuts  $\langle \alpha \rangle \varphi$  and  $\lceil \alpha \rceil \varphi$  as in the proof of Theorem 3.

Let  $\mathcal{K} = (\mathcal{A}, \mathcal{T})$  be a knowledge base,  $\alpha$  a DLXPath<sub>core</sub><sup>path-pos</sup> query and  $(c', c'')$  a pair of individuals. Next we construct a CPDL<sup>( $\neg$ )</sup> formula  $\varphi_{\mathcal{K}, \alpha}$  which is satisfiable if and only if  $(c', c'') \notin \text{Certain}(\alpha, (\mathcal{A}, \mathcal{T}))$ .

The formula  $\varphi_{\mathcal{K}, \alpha}$  uses concept and role names of  $\mathcal{K}$ . Besides this,  $\varphi_{\mathcal{K}, \alpha}$  uses a special role name *Create* and for every individual  $c_i$  in  $\mathcal{A}$  a concept name  $C_i$ . We start with the part of  $\varphi_{\mathcal{K}, \alpha}$  which corresponds to the ABox  $\mathcal{A}$ :

$$\varphi_{\mathcal{A}} = \left( \bigwedge_{\substack{i \neq j \\ c_i, c_j \text{ individuals in } \mathcal{A}}} C_i \Rightarrow \neg C_j \right) \wedge \left( \bigwedge_{A(c_i) \in \mathcal{A}} C_i \Rightarrow A \right) \wedge \left( \bigwedge_{R(c_i, c_j) \in \mathcal{A}} C_i \Rightarrow \langle R \rangle C_j \right).$$

Note, that if we drop UNA, then the first part of  $\varphi_{\mathcal{A}}$  should be dropped as well. We continue with the part of  $\varphi_{\mathcal{K}, \alpha}$  corresponding to the TBox  $\mathcal{T}$ :

$$\begin{aligned} \varphi_{\mathcal{T}} = & \left( \bigwedge_{B_i \sqsubseteq B_j \in \mathcal{T}} \delta(B_i) \Rightarrow \delta(B_j) \right) \wedge \left( \bigwedge_{B_i \sqcap B_j \sqsubseteq \perp \in \mathcal{T}} \delta(B_i) \Rightarrow \neg \delta(B_j) \right) \\ & \wedge \left( \bigwedge_{R_i \sqsubseteq R_j \in \mathcal{T}} \lceil \bar{R}_i \cup R_j \rceil \top \right) \wedge \left( \bigwedge_{R_i \sqcap R_j \sqsubseteq \perp \in \mathcal{T}} \lceil \bar{R}_i \cup \bar{R}_j \rceil \top \right), \end{aligned}$$

where  $\delta(B)$  is  $\langle R \rangle$  if  $B = \exists R$  and  $B$  otherwise.

Having  $\varphi_{\mathcal{A}}$  and  $\varphi_{\mathcal{T}}$  components at hand, we next define the formula corresponding to the whole KB  $\mathcal{K}$  with respect to the query  $\alpha$ :

$$\begin{aligned} \varphi_{\mathcal{K}} = & \left( \bigwedge_{c_i \text{ individual in } \mathcal{A}} \langle \text{Create} \rangle C_i \right) \wedge \lceil \beta \rceil (\varphi_{\mathcal{A}} \wedge \varphi_{\mathcal{T}}) \\ & \wedge \left( \bigwedge_{\substack{\psi \in \text{cl}(\varphi_{\mathcal{A}} \wedge \varphi_{\mathcal{T}} \wedge \langle \alpha \rangle \top) \\ c_i \text{ individual in } \mathcal{A}}} \langle \beta \rangle (C_i \wedge \psi) \Rightarrow \lceil \beta \rceil (C_i \Rightarrow \psi) \right), \end{aligned}$$

where

$$\beta = \left( \text{Create} \cup \text{Create}^- \cup \left( \bigcup_{R \text{ role name in } \mathcal{K}} (R \cup R^-) \right) \right)^*.$$

Finally, to take into account the testing individuals  $c'$  and  $c''$ ,

$$\varphi_{\mathcal{K},\alpha} = \varphi_{\mathcal{K}} \wedge \neg(\text{Create}[C']\alpha[C''])\top.$$

Similarly as in [12] one can check that the formula  $\varphi_{\mathcal{K},\alpha}$  has the required property.

□

**Theorem 5.** *The problem DL-Lite<sub>R</sub> CERTAIN ANSWERS OVER DLXPath<sub>core</sub><sup>pos</sup> is NP-complete. It remains NP-hard for DL-Lite<sub>core</sub>.*

*Proof.* We start with the upper bound. In what follows we will use the following standard notions.

The (*non-oblivious*) chase  $\text{Chase}(\mathcal{K})$  of a (satisfiable) KB  $\mathcal{K} = (\mathcal{T}, \mathcal{A})$  is a set of assertions which is a result of exhaustive (possibly infinite) application the following rules, starting with  $\mathcal{A}$  (by saying that a set of atoms contains  $P^-(c_1, c_2)$  we mean that it contains  $P(c_2, c_1)$ , and saying that it contains  $\exists R(c)$  we mean that there exists  $c'$  such that  $R(c, c')$  is contained in this set):

- add  $A_i(c)$  if  $\mathcal{T}$  contains  $B \sqsubseteq A_i$  and  $\text{Chase}(\mathcal{K})$  contains  $B(c)$ ;
- add  $R_2(c_1, c_2)$  if  $\mathcal{T}$  contains  $R_1 \sqsubseteq R_2$  and  $\text{Chase}(\mathcal{K})$  contains  $R_1(c_1, c_2)$ ;
- add  $R(c, c^{\text{new}})$ , for a fresh  $c^{\text{new}}$ , if  $\mathcal{T}$  contains  $B \sqsubseteq \exists R$  and  $\text{Chase}(\mathcal{K})$  contains  $B(c)$  but does not contain  $\exists R(c)$ .

The *canonical model*  $\text{Can}(\mathcal{K})$  of a satisfiable KB  $\mathcal{K}$  is the model which interprets all individuals in the  $\text{Chase}(\mathcal{K})$  by themselves,  $c \in A^{\text{Can}(\mathcal{K})}$  if and only if  $A(c) \in \text{Chase}(\mathcal{K})$ , and  $(c_1, c_2) \in P^{\text{Can}(\mathcal{K})}$  if and only if  $P(c_1, c_2) \in \text{Chase}(\mathcal{K})$ . It is well known that for any model  $\mathcal{I}$  of  $\mathcal{K}$  there exists a homomorphism from  $\text{Can}(\mathcal{K})$  to  $\mathcal{I}$ .

The following claim can be proved by straightforward induction on the structure of the DLXPath<sub>core</sub><sup>pos</sup> query.

*Claim 5.4.* A pair of individuals  $(c_1, c_2)$  is a certain answer to a DLXPath<sub>core</sub><sup>pos</sup> query  $\alpha$  over a satisfiable KB  $\mathcal{K}$  if and only if  $(c_1, c_2) \in \llbracket \alpha \rrbracket^{\text{Can}(\mathcal{K})}$ .

Essentially, this claim says that we can concentrate only on the canonical model. Since it can be quite big and even infinite, next we show that for positive query answering it is always enough to guess just a polynomial part of the canonical model. In what follows we need one more definition.

Given a DLXPath<sub>core</sub><sup>pos</sup> query which does not use unary  $\vee$  primitive, binary  $\cup$  primitive, concept names, as well as the transitive closure operator  $+$  but may use the reflexive transitive closure operator  $*$  over roles, a *skeleton*  $\mathcal{S}_\alpha$  of  $\alpha$  is a directed unordered tree represented the structure of  $\alpha$  in the sense that besides the root it has a distinguished node called *drain*, each edge is labelled by either  $R$  or  $R^*$  according to  $\alpha$  and branching in this tree represents nesting  $[\varphi]$  in  $\alpha$ .

Let now  $\alpha$  be a DLXPath<sub>core</sub><sup>pos</sup> query and  $\mathcal{K} = (\mathcal{T}, \mathcal{A})$  be a DL-Lite<sub>R</sub> KB. Without loss of generality we may assume that

- the KB  $\mathcal{K}$  is satisfiable (if not, it can be easily checked in polynomial time),
- the vocabulary does not contain any concept names (they can be simulated by role names in the straightforward manner),
- $\alpha$  does not use  $\forall$  and  $\cup$  primitives (our algorithm can always guess which alternative to use), and
- $\alpha$  does not use the transitive closure operator  $+$ , but may use the reflexive transitive closure operator  $*$  over roles (these primitives are interreducible).

For clarity of exposition we also assume that  $\mathcal{K}$  uses only one individual  $c$ , the generalisation to arbitrary case is a matter of technicality.

It is straightforward to show that  $(c, c) \in \llbracket \alpha \rrbracket^{Can(\mathcal{K})}$  if and only if there exists an embedding  $\gamma$  of the skeleton  $\mathcal{S}_\alpha$  to  $Can(\mathcal{K})$ , that is a mapping preserving labels in the way that

- if  $R$  is a label of  $(a, a')$  in  $\mathcal{S}_\alpha$  then  $R^{Can(\mathcal{K})}(\gamma(a), \gamma(a'))$  holds,
- if  $R^*$  is a label of  $(a, a')$  in  $\mathcal{S}_\alpha$  then there exist a sequence  $c_1, \dots, c_n$ ,  $n \geq 1$ , without repetitions such that  $c_1 = \gamma(a)$ ,  $c_n = \gamma(a')$  and  $R^{Can(\mathcal{K})}(\gamma(c_i), \gamma(c_{i+1}))$  holds for all  $i$ ,

which maps the root and the drain to  $c$ . Note that given  $\gamma$ , under our assumptions the sequence  $c_1, \dots, c_n$  (without repetitions) satisfying the second requirement of embedding is unique for every edge  $(a, a')$  labelled  $R^*$ . We denote it  $\Gamma(a, a')$ .

Having this fact at hand, it is enough to show that if such an embedding  $\gamma$  exists then there exists an embedding such that the maximal distance from  $c$  to the image of a node from the skeleton in  $Can(\mathcal{K})$  is of polynomial size. Indeed, if it holds then we can guess all the polynomial number of branches in  $Can(\mathcal{K})$  which leads to images of all the nodes from skeleton and check whether the query holds on this polynomially sized sub-interpretation of the canonical model. That is why the rest of the proof is devoted to reducing this maximal distance for a given embedding. For this, we reuse the *shortcutting* technique, developed in [3] for a similar problem in the XML context.

Let  $c'$  be the image  $\gamma(a')$  of a node in  $\mathcal{S}_\alpha$  with the maximal distance from  $c$  in  $Can(\mathcal{K})$ . The branch  $c = c_1, \dots, c_m = c'$  in the canonical model can be partitioned in the following way:  $c_i$  and  $c_{i+1}$  are in the same part if and only if

- both of them are not images of any nodes of  $\mathcal{S}_\alpha$  by  $\gamma$ , and
- $\{(a, a') \mid c_i \in \Gamma(a, a')\} = \{(a, a') \mid c_{i+1} \in \Gamma(a, a')\}$ .

That is, every part of this partitioning consists of a continuous sequence of nodes on the branch from  $c$  to  $c'$ , and every image of a node of the skeleton (which is on the branch) forms its own part. However, we are interested in those parts which contains more than  $|\mathcal{T}|$  elements. If such a part does not exist, we are done, because the length  $n$  of the branch is bounded by  $|\mathcal{T}| \times |\mathcal{S}_\alpha|$ . If, contrary, such a part exists, then it contains two elements  $c_i, c_j$ ,  $i > j$ , with the same type in  $Can(\mathcal{K})$  (by *type* of an element  $d$  in an interpretation  $\mathcal{I}$  we mean the set  $\{(\exists R) \mid \text{there exists } d' \text{ such that } R^{\mathcal{I}}(d, d')\}$ ). Since the types are the same, we can simply replace the tree-like section of the canonical model starting in  $c_i$  with the section starting in  $c_j$ . After such a replacement the canonical model stays the same (since it is infinite), but the distance from  $c$  to  $c'$  strictly decreases.

Applying such a procedure to the canonical model  $Can(\mathcal{K})$  and embedding  $\gamma$  while possible we arrive to an embedding with desired properties. It means that the pair  $(c, c)$  is a certain answer to  $\alpha$  over satisfiable  $\mathcal{K}$  if and only if there exists a polynomially sized sub-interpretation of the canonical model of  $\mathcal{K}$  which witness  $\alpha$  for  $(c, c)$ .

To show hardness we give a reduction from 3CNF-SAT. Suppose we are given a conjunction  $\Phi$  of clauses of the form  $\ell_1 \vee \ell_2 \vee \ell_3$ , where the  $\ell_k$  are literals, that is, propositional variables or their negations (we can assume that all literals in each clause are distinct). Let  $x_1, \dots, x_n$  be the variables in  $\Phi$ .

Consider the KB  $\mathcal{K}$  with the ABox  $\mathcal{A} = \{A(c)\}$  and the TBox  $\mathcal{T}$  consisting of the inclusions

$$\begin{aligned} A &\sqsubseteq \exists(X_1^v)^-, \text{ for } v \in \{\top, \perp\}, \\ \exists X_{i-1}^v &\sqsubseteq \exists(X_i^u)^-, \text{ for } v, u \in \{\top, \perp\} \text{ and } 1 < i \leq n. \end{aligned} \quad (2)$$

For every variable  $x_i$  and truth value  $v \in \{\top, \perp\}$  let  $\psi_{x_i}^v$  be the node formula:

$$\langle\langle (X_n^\top \cup X_n^\perp) \cdots (X_{i+1}^\top \cup X_{i+1}^\perp) \cdot X_i^v \cdot (X_{i-1}^\top \cup X_{i-1}^\perp) \cdots (X_1^\top \cup X_1^\perp) \rangle\rangle.$$

Let  $\varphi$  be a node formula  $\langle\alpha[\psi]\rangle$ , where

$$\alpha = ((X_1^\top)^- \cup (X_1^\perp)^-) \cdot \dots \cdot ((X_n^\top)^- \cup (X_n^\perp)^-),$$

and  $\psi$  is a conjunction of node formulas of the following form, for each clause  $\gamma$  in  $\Phi$ , using variables  $x, y, z$ :

$$\bigvee_{\substack{(v_x, v_y, v_z) \text{ assignment of } (x, y, z) \\ \text{satisfying } \gamma}} \psi_x^{v_x} \wedge \psi_y^{v_y} \wedge \psi_z^{v_z}.$$

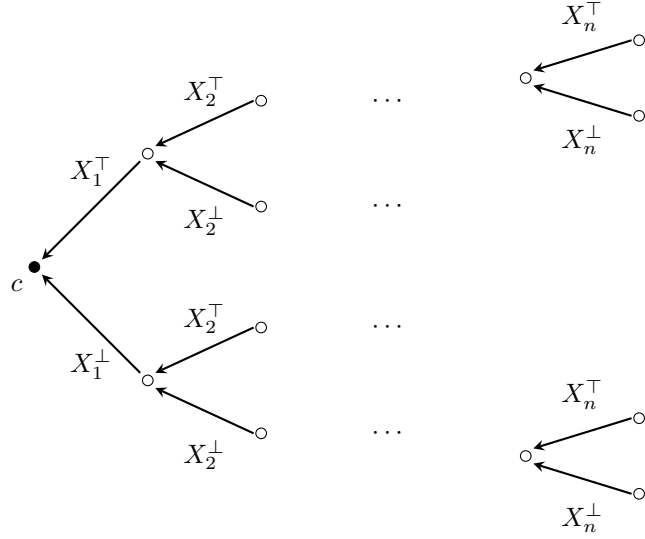
We now show that the pair  $(c, c)$  is a certain answer to  $[\varphi]$  over  $\mathcal{K}$  if and only if  $\Phi$  is satisfiable.

Suppose first that  $\Phi$  is satisfiable and let  $\sigma : \{x_1, \dots, x_n\} \rightarrow \{\top, \perp\}$  be a truth assignment such that  $\sigma(\Phi)$  is true. Take any model  $\mathcal{I} = (\Delta^\mathcal{I}, \cdot^\mathcal{I})$  of  $\mathcal{K}$ . Then using assertions from (2) we can find elements  $d_0, d_1, d_2, \dots, d_n \in \Delta^\mathcal{I}$  such that:

$$\begin{aligned} d_0 &= c^\mathcal{I}, \text{ and} \\ (X_i^{\sigma(x_i)})^\mathcal{I}(d_i, d_{i-1}), &\text{ for all } 1 \leq i \leq n. \end{aligned} \quad (3)$$

Then by definition  $(c^\mathcal{I}, d_n) \in \llbracket \alpha \rrbracket^\mathcal{I}$ . To show that  $c^\mathcal{I} \in \llbracket \varphi \rrbracket^\mathcal{I}$  it suffices to check that  $d_n \in \llbracket \psi \rrbracket^\mathcal{I}$ . To see this take any clause  $\gamma$  in  $\Phi$  using variables  $x, y, z$  and consider the formula  $\psi_x^{\sigma(x)} \wedge \psi_y^{\sigma(y)} \wedge \psi_z^{\sigma(z)}$ . The facts (3) above immediately imply that  $d_n \in \llbracket \psi_x^{\sigma(x)} \rrbracket^\mathcal{I}$ . Analogously we get that  $d_n \in \llbracket \psi_y^{\sigma(y)} \rrbracket^\mathcal{I}$  and  $d_n \in \llbracket \psi_z^{\sigma(z)} \rrbracket^\mathcal{I}$ . From this it follows that  $d_n \in \llbracket \psi \rrbracket^\mathcal{I}$  and therefore  $c^\mathcal{I} \in \llbracket \varphi \rrbracket^\mathcal{I}$ . Since  $\mathcal{I}$  was chosen arbitrary we conclude that  $(c, c)$  is a certain answer to  $[\varphi]$  over  $\mathcal{K}$ .

Conversely, assume that  $(c, c)$  is a certain answer to  $[\varphi]$  over  $\mathcal{K}$ . Consider the canonical model  $\mathcal{I}_0$  of  $\mathcal{K}$ . Therefore we can view  $\mathcal{I}_0$  as a complete (inverted) binary tree of height  $n$ , as illustrated in the following image.



As before we depict the fact that  $(X_i^v)^{\mathcal{I}_0}(d, d')$  holds by having an  $X_i^v$  labelled arrow between nodes representing  $d$  and  $d'$ . Since  $(c, c)$  is a certain answer to  $[\varphi]$  over  $\mathcal{K}$  we have that  $c^{\mathcal{I}_0} \in \llbracket \varphi \rrbracket^{\mathcal{I}_0}$ . Therefore there exist elements  $d_0, d_1, \dots, d_n$  in  $\mathcal{I}_0$  and truth values  $v_1, \dots, v_n$  such that:

$$\begin{aligned}
 d_0 &= c^{\mathcal{I}_0}, \\
 (X_i^{v_i})^{\mathcal{I}_0}(d_i, d_{i-1}), & \text{ for all } 1 \leq i \leq n, \text{ and} \\
 d_n &\in \llbracket \psi \rrbracket^{\mathcal{I}_0}.
 \end{aligned}$$

Define now  $\sigma(x_i) := v_i$  for all  $i = 1, \dots, n$ . We claim that  $\sigma$  is a satisfying assignment for  $\Phi$ . To see this take any clause  $\gamma$  in  $\Phi$  and assume that  $\gamma$  uses variables  $x, y, z$ . Since  $d_n \in \llbracket \psi \rrbracket^{\mathcal{I}_0}$  this implies that there is an assignment  $v_x, v_y, v_z$  of variables  $x, y, z$  satisfying  $\gamma$  and such that the formula  $\psi_x^{v_x} \wedge \psi_y^{v_y} \wedge \psi_z^{v_z}$  is true at  $d_n$ . In particular  $\psi_x^{v_x}$  is true at  $d_n$  which implies that we can reach  $c^{\mathcal{I}_0}$  using an  $X_i^{v_x}$  labelled edge, where  $x = x_i$ . Since every path in  $\mathcal{I}_0$  has a unique label we conclude that  $v_x = v_i$  and similarly for  $y$  and  $z$ . Since  $\gamma$  was arbitrary we conclude that  $\Phi$  is satisfiable.  $\square$