

Conjunctive Query Containment and Answering Under Description Logic Constraints

DIEGO CALVANESE

Free University of Bozen-Bolzano

and

GIUSEPPE DE GIACOMO and MAURIZIO LENZERINI

Università di Roma “La Sapienza”

Query containment and query answering are two important computational tasks in databases. While query answering amounts to computing the result of a query over a database, query containment is the problem of checking whether, for every database, the result of one query is a subset of the result of another query.

In this article, we deal with unions of conjunctive queries, and we address query containment and query answering under description logic constraints. Every such constraint is essentially an inclusion dependency between concepts and relations, and their expressive power is due to the possibility of using complex expressions in the specification of the dependencies, for example, intersection and difference of relations, special forms of quantification, regular expressions over binary relations. These types of constraints capture a great variety of data models, including the relational, the entity-relationship, and the object-oriented model, all extended with various forms of constraints. They also capture the basic features of the ontology languages used in the context of the Semantic Web.

We present the following results on both query containment and query answering. We provide a method for query containment under description logic constraints, thus showing that the problem is decidable, and analyze its computational complexity. We prove that query containment is undecidable in the case where we allow inequalities in the right-hand-side query, even for very simple constraints and queries. We show that query answering under description logic constraints can be

This work has been partially supported by the EU funded FP6-7603 FET Project Thinking OntologiES (TONES), by project HYPER, funded by IBM through a Shared University Research (SUR) Award grant, and by MIUR FIRB 2005 project “Tecnologie Orientate alla Conoscenza per Aggregazioni di Imprese in Internet” (TOCAL.IT).

This article is an extended and revised version of an article published in the *Proceedings of the 17th ACM SIGACT SIGMOD SIGART Symposium on Principles of Database Systems (PODS'98)*. Authors' addresses: D. Calvanese, Faculty of Computer Science, Free University of Bozen-Bolzano, Piazza Domenicani 3, 39100 Bolzano, Italy; email: calvanese@inf.unibz.it; G. De Giacomo and M. Lenzerini, Dipartimento di Informatica e Sistemistica, Università di Roma “La Sapienza”, Via Ariosto 25, 00185, Roma, Italy; email: [degiacomo,lenzerini@dis.uniroma1.it](mailto:{degiacomo,lenzerini}@dis.uniroma1.it).

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or direct commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org. © 2008 ACM 1529-3785/2008/06-ART22 \$5.00 DOI 10.1145/1352582.1352590 <http://doi.acm.org/10.1145/1352582.1352590>

reduced to query containment, and illustrate how such a reduction provides upper-bound results with respect to both combined and data complexity.

Categories and Subject Descriptors: I.2.4 [**Artificial Intelligence**]: Knowledge Representation Formalisms and Methods—*Representation languages*; F.4.1 [**Formal Languages**]: Mathematical Logic—*Computational logic, modal logic*; H.2.3 [**Database Management**]: Languages—*Query languages*

General Terms: Algorithms, Theory

Additional Key Words and Phrases: Description logics, conjunctive queries, query containment, computational complexity

ACM Reference Format:

Calvanese, D., De Giacomo, G., and Lenzerini, M. 2008. Conjunctive query containment and answering under description logic constraints. *ACM Trans. Comput. Logic* 9, 3, Article 22 (June 2008), 31 pages. DOI = 10.1145/1352582.1352590 <http://doi.acm.org/10.1145/1352582.1352590>

1. INTRODUCTION

Query containment and query answering are two important computational tasks in databases. While query answering amounts to compute the result of a query over a database, query containment is the problem of checking whether for every database, the result of one query is a subset of the result of another query.¹ Many articles point out that checking containment is a relevant task in several contexts, including information integration [Ullman 1997], query optimization [Abiteboul et al. 1995; Aho et al. 1979a], (materialized) view maintenance [Gupta and Mumick 1995], data warehousing [Widom 1995], constraint checking [Gupta et al. 1994], and semantic caching [Amir et al. 2003].

In this article, we deal with query containment and query answering under integrity constraints, or simply constraints.

The former is the problem of checking whether containment between two queries holds for every database satisfying a given set of constraints. This problem arises in those situation where one wants to check query containment relative to a database schema specified with a rich data definition language. For example, in the case of information integration, queries are often to be compared relatively to (interschema) constraints, which are used to declaratively specify the “glue” between two source schemas, and between one source schema and the global schema [Calvanese et al. 1998b; Hull 1997; Ullman 1997; Catarci and Lenzerini 1993; Levy et al. 1995; Lenzerini 2002; Halevy 2001].

The complexity of query containment in the absence of constraints has been studied in various settings. In Chandra and Merlin [1977], NP-completeness has been established for conjunctive queries, and in Chekuri and Rajaraman [1997] a multiparameter analysis has been performed for the same case, showing that the intractability is due to certain types of cycles in the queries. In Klug [1988] and van der Meyden [1998], Π_2^P -completeness of containment of conjunctive queries with inequalities was proved, and in Sagiv and Yannakakis [1980] the case of queries with the union and difference operators was studied. For

¹We refer to the set semantics of query containment. Bag semantics is studied, for example, in Ioannidis and Ramakrishnan [1995]

various classes of Datalog queries with inequalities, decidability and undecidability results were presented in Chaudhuri and Vardi [1992], van der Meyden [1998], Bonatti [2004], and Calvanese et al. [2003], respectively.

Query containment under constraints has also been the subject of several investigations. For example, decidability of conjunctive query containment was investigated in Aho et al. [1979b] under functional and multi valued dependencies, in Johnson and Klug [1984] under functional and inclusion dependencies, in Chan [1992], Levy and Rousset [1996], and Levy and Suciu [1997] under constraints representing *is-a* hierarchies and complex objects, and in Dong and Su [1996] in the case of constraints represented as Datalog programs. Undecidability was proved in Calvanese and Rosati [2003] for recursive queries under inclusion dependencies. Several results on containment of XML queries under constraints expressed as DTDs were reported in Neven and Schwentick [2003] and Wood [2003].

Query answering under constraints is the problem of computing the answers to a query over an incomplete database relatively to a set of constraints van der Meyden [1998]. Since an incomplete database is partially specified, this task amounts to computing the tuples that satisfy the query in every database that conforms to the partial specification, and satisfies the constraints. It is well known in the database literature that there is a tight connection between the problems of conjunctive query containment and conjunctive query answering Chandra and Merlin [1977]. Since this relationship holds also in the presence of constraints, most of the results reported above apply to query answering as well. In this article, we concentrate mainly on query containment, and address query answering only in Section 5.

In this article,² we address query containment and answering in a setting where the following occur.

- The schema is constituted by concepts (unary relations) and relations as basic elements, and by a set of constraints expressed in a variant of description logics [Baader et al. 2003]. Every constraint is an inclusion of the form $\alpha_1 \subseteq \alpha_2$, where α_1 and α_2 are complex expressions built by using intersection and difference of relations, special forms of quantification, regular expressions over binary relation, and number restrictions (i.e., cardinality constraints imposing limitations on the number of tuples in a certain relation in which an object may appear). The constraints express essentially inclusion dependencies between concepts and relations, and their expressive power is due to the possibility of using complex expressions in the specification of the dependencies. It can be shown that our formalism is able to capture a great variety of data models, including the relational, the entity-relationship, and the object-oriented model, all extended with various forms of constraints. The relevance of the constraints dealt with in this paper is also testified by the large interest that the semantic Web community expresses toward description logics. Indeed, several articles

²This article is an improved and extended version of part of Calvanese et al. [1998a]. In that article, we overlooked a technical detail in the encoding of the containment problem, which was pointed out in Horrocks et al. [2000]. Here we present a corrected version of the original encoding.

point out that ontologies play a key role in developing semantic Web tools [Gruber 1993], and description logics are regarded as the main formalisms for the specification of ontologies in this context [Patel-Schneider et al. 2004]. The results presented in this article are also relevant as a formal analysis on querying ontologies. In particular, \mathcal{DLR}_{reg} can be considered as one of the most expressive description logics studied in the literature [Calvanese and De Giacomo 2003]

- Queries are formed as disjunctions of conjunctive queries whose atoms are concepts and relations, and therefore can express non-recursive Datalog programs. Regular expressions are confined to concepts and relations, and do not constitute atoms per se.
- An incomplete database is specified as a set of facts asserting that a specific object is an instance of a concept, or that a specific tuple of objects is an instance of a relation. As we said before, an incomplete database \mathcal{D} is intended to provide a partial specification of a database, in the sense that a database conforming to \mathcal{D} contains all facts explicitly asserted in \mathcal{D} , and may contain additional instances of concepts and relations.

We observe that, given the form of constraints and queries allowed in our approach, none of the previous results can be applied to get decidability/undecidability of query containment and query answering in our setting.

We present the following results on both query containment and query answering:

- (1) We provide a method for query containment under description logic constraints, thus showing that the problem is decidable, and analyze its computational complexity. This result is obtained by resorting to a translation of the schema and the containment to be checked into a particular propositional dynamic logic (PDL) formula, and then verifying the unsatisfiability of the formula. The technique is justified by the fact that reasoning about the schema itself (without considering the queries) is optimally done within the framework of PDL [De Giacomo and Lenzerini 1996].
- (2) We prove that query containment is undecidable in the case where we allow inequalities in the right-hand side query, even for very simple constraints and queries.
- (3) We show that query answering under description logic constraints can be reduced to query containment, and illustrate how such a reduction provides upper bound results with respect to both combined and data complexity.

The article is organized as follows. In Section 2, we present the formalism used to express both the constraints in the schema, and the queries. In Section 3, we deal with query containment. In particular, in Section 3.1 we describe the logic CPDL_g , which will be used for deciding query containment, in Section 3.2 we describe the reduction of query containment to unsatisfiability in CPDL_g , in Section 3.3 we prove its correctness, and in Section 3.4 we analyze the complexity bounds for checking containment of queries. In Section 4, we show

undecidability of query containment in the presence of inequalities. In Section 5, we deal with query answering, and in Section 6 we conclude the article.

2. SCHEMAS AND QUERIES IN \mathcal{DLR}_{reg}

To specify database schemas and queries, we use the logical language \mathcal{DLR}_{reg} , inspired by Catarci and Lenzerini [1993] and Calvanese et al. [1995], belonging to the family of (expressive) description logics [Calvanese et al. 2001b; Baader et al. 2003]. The language is based on the relational model, in the sense that a schema \mathcal{S} describes the properties of a set of relations, while a query for \mathcal{S} denotes a relation that is supposed to be computed from any database conforming to \mathcal{S} . A schema is specified in terms of a set of assertions on relations, which express the constraints that must be satisfied by every conforming database. Note that such a notion of schema corresponds to that of TBox in description logics [Baader et al. 2003].

2.1 Schemas

The basic elements of \mathcal{DLR}_{reg} are *concepts* (unary relations), *n-ary relations*, and *regular expressions* built over projections of relations on two of their components.³

We assume to deal with a finite set of atomic concepts and relations, denoted by A and \mathbf{P} , respectively. Each atomic relation has an associated *arity*, which is an integer greater than 1. We use C to denote arbitrary concepts, \mathbf{R} to denote arbitrary relations, and E to denote regular expressions, respectively, built according to the following syntax:

$$\begin{aligned} C &::= \top_1 \mid A \mid \neg C \mid C_1 \sqcap C_2 \mid \exists E.C \mid \exists[\$i]\mathbf{R} \mid (\leq k[\$i]\mathbf{R}), \\ \mathbf{R} &::= \top_n \mid \mathbf{P} \mid (\$i/n:C) \mid \neg\mathbf{R} \mid \mathbf{R}_1 \sqcap \mathbf{R}_2, \\ E &::= \mathbf{R}_{\$i,\$j} \mid \varepsilon \mid E_1 \circ E_2 \mid E_1 \sqcup E_2 \mid E^*, \end{aligned}$$

where i, j are positive integers that intuitively denote components of relations, n is a positive integer greater than 1 denoting the arity of a relation, and k is a nonnegative integer. Actually, we restrict the attention to concepts and relations that are *well-typed*, which means that

- only relations of the same arity n are combined to form expressions of type $\mathbf{R}_1 \sqcap \mathbf{R}_2$ (which inherit the arity n), and
- $i \leq n$ whenever i denotes a component of a relation of arity n .

A \mathcal{DLR}_{reg} schema is constituted by a finite set of *assertions*, of the form

$$\begin{aligned} C_1 &\sqsubseteq C_2, \\ \mathbf{R}_1 &\sqsubseteq \mathbf{R}_2, \end{aligned}$$

where \mathbf{R}_1 and \mathbf{R}_2 are of the same arity. Note that our notion of schema corresponds to that of TBox in description logics [Baader et al. 2003].

³We could include in the logic also domains, that is, sets of values such as integer, string, etc. However, for the sake of simplicity, we do not consider this aspect in this work.

$$\begin{aligned}
\top_1^{\mathcal{I}} &= \Delta^{\mathcal{I}} \\
A^{\mathcal{I}} &\subseteq \Delta^{\mathcal{I}} \\
(\neg C)^{\mathcal{I}} &= \Delta^{\mathcal{I}} \setminus C^{\mathcal{I}} \\
(C_1 \sqcap C_2)^{\mathcal{I}} &= C_1^{\mathcal{I}} \cap C_2^{\mathcal{I}} \\
(\exists E.C)^{\mathcal{I}} &= \{d \in \Delta^{\mathcal{I}} \mid \exists d' \in C^{\mathcal{I}}. (d, d') \in E^{\mathcal{I}}\} \\
(\exists [\$i]\mathbf{R})^{\mathcal{I}} &= \{d \in \Delta^{\mathcal{I}} \mid \exists (d_1, \dots, d_n) \in \mathbf{R}^{\mathcal{I}}. d_i = d\} \\
(\leq k [\$i]\mathbf{R})^{\mathcal{I}} &= \{d \in \Delta^{\mathcal{I}} \mid \#\{(d_1, \dots, d_n) \in \mathbf{R}_1^{\mathcal{I}} \mid d_i = d\} \leq k\} \\
\top_n^{\mathcal{I}} &\subseteq (\Delta^{\mathcal{I}})^n \\
\mathbf{P}^{\mathcal{I}} &\subseteq \top_n^{\mathcal{I}} \\
(\$i/n : C)^{\mathcal{I}} &= \{(d_1, \dots, d_n) \in \top_n^{\mathcal{I}} \mid d_i \in C^{\mathcal{I}}\} \\
(\neg \mathbf{R})^{\mathcal{I}} &= \top_n^{\mathcal{I}} \setminus \mathbf{R}^{\mathcal{I}} \\
(\mathbf{R}_1 \sqcap \mathbf{R}_2)^{\mathcal{I}} &= \mathbf{R}_1^{\mathcal{I}} \cap \mathbf{R}_2^{\mathcal{I}} \\
(\mathbf{R} |_{\$,i,\$,j})^{\mathcal{I}} &= \{(x_i, x_j) \mid (x_1, \dots, x_n) \in \mathbf{R}^{\mathcal{I}}\} \\
\varepsilon^{\mathcal{I}} &= \{(x, x) \mid x \in \Delta^{\mathcal{I}}\} \\
(E_1 \circ E_2)^{\mathcal{I}} &= E_1^{\mathcal{I}} \circ E_2^{\mathcal{I}} \\
(E_1 \sqcup E_2)^{\mathcal{I}} &= E_1^{\mathcal{I}} \cup E_2^{\mathcal{I}} \\
(E^*)^{\mathcal{I}} &= (E^{\mathcal{I}})^*
\end{aligned}$$

Fig. 1. Semantic rules for \mathcal{DLR}_{reg} (\mathbf{P} , \mathbf{R} , \mathbf{R}_1 , and \mathbf{R}_2 have arity n).

The semantics of \mathcal{DLR}_{reg} is specified through the notion of interpretation. An *interpretation* $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$ of a \mathcal{DLR}_{reg} schema \mathcal{S} and a set \mathcal{C} (of constants to be used in queries) is constituted by an *interpretation domain* $\Delta^{\mathcal{I}}$ and an *interpretation function* $\cdot^{\mathcal{I}}$ that assigns

- to each constant c in \mathcal{C} , an element $c^{\mathcal{I}}$ of $\Delta^{\mathcal{I}}$ under the *unique name assumption*, that is, for any two distinct constants c_1 and c_2 of $\Delta^{\mathcal{I}}$, we have that $c_1^{\mathcal{I}} \neq c_2^{\mathcal{I}}$;
- to each atomic concept A , a subset $A^{\mathcal{I}}$ of $\Delta^{\mathcal{I}}$;
- to each relation \mathbf{P} of arity n , a subset $\mathbf{P}^{\mathcal{I}}$ of $(\Delta^{\mathcal{I}})^n$.

The interpretation function is extended to arbitrary concepts, arbitrary relations, and regular expressions in such a way that the conditions in Figure 1 are satisfied.

We observe that \top_1 denotes the interpretation domain, while \top_n , for $n > 1$, does *not* denote the n -Cartesian product of the domain, but only a subset of it that covers all relations of arity n . It follows from this property that the “ \neg ” constructor on relations is used to express difference of relations, rather than complement. The constructors “ \neg ” and “ \sqcap ” on concepts have the usual meaning of negation and conjunction, respectively. The expression $\exists[\$i]\mathbf{R}$ denotes the projection of relation \mathbf{R} on its i th component. An expression of the form $(\leq k [\$i]\mathbf{R})$ is called a *number restriction*, and imposes a limit on the number of times an object can participate to relation \mathbf{R} as the i th component. The expression $(\$i/n : C)$ represents those tuples of arity n whose i th component is an instance of concept C . Finally, E are regular expressions build over projections of relations on two of their components.

In what follows, we abbreviate $\neg\exists E.C$ with $\forall E.C$, and $(\$i/n : C)$ with $(\$i : C)$ when n is clear from the context.

An interpretation \mathcal{I} *satisfies* an assertion $C_1 \sqsubseteq C_2$ (respectively, $\mathbf{R}_1 \sqsubseteq \mathbf{R}_2$) if $C_1^{\mathcal{I}} \subseteq C_2^{\mathcal{I}}$ (respectively, $\mathbf{R}_1^{\mathcal{I}} \subseteq \mathbf{R}_2^{\mathcal{I}}$). An interpretation that satisfies all assertions in a schema \mathcal{S} is called a *model* of \mathcal{S} . It is easy to see that a model of a schema \mathcal{S} actually corresponds to a database conforming to \mathcal{S} , that is, a database satisfying all the constraints represented by \mathcal{S} . A schema is *satisfiable* if it admits a model. A schema \mathcal{S} *logically implies* an inclusion assertion $C_1 \sqsubseteq C_2$ (respectively, $\mathbf{R}_1 \sqsubseteq \mathbf{R}_2$) if for every model \mathcal{I} of \mathcal{S} we have that $C_1^{\mathcal{I}} \subseteq C_2^{\mathcal{I}}$ (respectively, $\mathbf{R}_1^{\mathcal{I}} \subseteq \mathbf{R}_2^{\mathcal{I}}$).

It can be shown that \mathcal{DLR}_{reg} is able to capture a great variety of data models with many forms of constraints. For example, we obtain the entity-relationship model (including *is-a* relations on both entities and relations) in a straightforward way [Calvanese et al. 1995], and an object-oriented data model (extended with several types of constraints), by restricting the use of existential and universal quantifications in concept expressions, by restricting the attention to binary relations, and by eliminating negation, disjunction, and regular expressions. Compared with the relational model, the following observations point out the kinds of constraints that can be expressed using \mathcal{DLR}_{reg} .

- Assertions directly express a special case of typed inclusion dependencies, namely, the one where no projection of relations is used.
- Unary inclusion dependencies are easily expressible by means of the $\exists[\$2]\mathbf{P}$ construct. For example, $\exists[\$2]\mathbf{P}_1 \sqsubseteq \exists[\$3]\mathbf{P}_2$ is a unary inclusion dependency between attribute 2 of \mathbf{P}_1 and attribute 3 of \mathbf{P}_2 .
- Existence and exclusion dependencies are expressible by means of \exists and \neg , respectively.
- A limited form of functional dependencies can be expressed by means of $(\leq 1[\$i]\mathbf{R})$. For example, $\top_1 \sqsubseteq (\leq 1[\$i]\mathbf{P})$ specifies that attribute i functionally determines all other attributes of \mathbf{P} .
- The possibility of constructing complex expressions provides a special form of view definition. Indeed, the two assertions $\mathbf{P} \sqsubseteq \mathbf{R}$, $\mathbf{R} \sqsubseteq \mathbf{P}$ (where \mathbf{R} is a complex expression) is a view definition for \mathbf{P} . Notably, views can be freely used in assertions (even with cyclic references), and, therefore, all the above discussed constraints can be imposed not only on atomic relations, but also on views. These features make our logic particularly suited for expressing inter-schema relationships in the context of information integration [Calvanese et al. 1998b], where it is crucial to be able to state that a certain concept of a schema corresponds (by means of inclusion or equivalence) to a view in another schema.
- Finally, regular expressions can be profitably used to represent in the schema inductively defined structures such as sequences and lists, imposing complex conditions on them.

One of the distinguishing features of \mathcal{DLR}_{reg} is that it is equipped with a method for checking logical implication. Indeed, \mathcal{DLR}_{reg} shares EXPTIME-completeness of schema satisfiability and logical implication with many expressive description logics [Calvanese et al. 2001b; Baader et al. 2003] (see below).

We point out that \mathcal{DLR}_{reg} supports only special forms of functional and inclusion dependencies. Hence the undecidability result of implication for (general) functional and inclusion dependencies taken together, shown in Mitchell [1983] and Chandra and Vardi [1985], does not apply.

We finally observe that \mathcal{DLR}_{reg} does not enjoy the *finite model property*, similarly to all description logics that include the ability to predicate at least on binary relations, stating their functionality and navigating them in both directions [Calvanese and De Giacomo 2003].

2.2 Queries

A query q for a \mathcal{DLR}_{reg} schema is a union of conjunctive queries, written (using Datalog notation) in the form

$$q(\vec{x}) \leftarrow \text{conj}_1(\vec{x}, \vec{y}_1, \vec{c}_1) \vee \dots \vee \text{conj}_m(\vec{x}, \vec{y}_m, \vec{c}_m),$$

where each $\text{conj}_i(\vec{x}, \vec{y}_i, \vec{c}_i)$ is a conjunction of *atoms*, and \vec{x}, \vec{y}_i (respectively, \vec{c}_i) are exactly the variables (respectively, constants) appearing in the i th conjunction.⁴ Each atom has one of the forms $C(t)$ or $\mathbf{R}(\vec{t})$, where

- t and \vec{t} are constants or variables in $\vec{x}, \vec{y}_i, \vec{c}_i$, and
- C and \mathbf{R} are, respectively, concepts and relations expressions over S .

The number of variables of \vec{x} is called the *arity* of q , that is, the arity of the relation denoted by the query q .

We observe that the atoms in the queries are arbitrary \mathcal{DLR}_{reg} concepts and relations, freely used in the assertions of the schema. This distinguishes our approach with respect to Donini et al. [1998] and Levy and Rousset [1996], where no constraints can be expressed in the schema on the relations that appear in the queries. Note also that regular expressions may appear in concepts and relations. However, they cannot be used to form binary atoms specifying that two variables or constants are connected by a path described by the regular expression [Calvanese et al. 2000].

Given an interpretation \mathcal{I} of a schema S , a query q for S of arity n is interpreted as the set $q^{\mathcal{I}}$ of n -tuples (o_1, \dots, o_n) , with each $o_i \in \Delta^{\mathcal{I}}$, such that, under the assignment of o_i to x_i , for $i \in \{1, \dots, n\}$, the formula

$$\exists \vec{y}_1. \text{conj}_1(\vec{x}, \vec{y}_1, \vec{c}_1) \vee \dots \vee \exists \vec{y}_m. \text{conj}_m(\vec{x}, \vec{y}_m, \vec{c}_m)$$

evaluates to true in \mathcal{I} [Enderton 1972].

If q and q' are two queries (of the same arity) for S , we say that q is *contained in q' with respect to S* , denoted $S \models q \subseteq q'$, if $q^{\mathcal{I}} \subseteq q'^{\mathcal{I}}$ for every model \mathcal{I} of S . The query q is *satisfiable with respect to S* , denoted $S \models q \equiv \emptyset$, if there is a model \mathcal{I} of S such that $q^{\mathcal{I}} \neq \emptyset$.

Query containment (with respect to a schema) is the problem of checking whether $S \models q \subseteq q'$, where S, q , and q' are given as input. *Query satisfiability*

⁴Note that the assumption that the variables \vec{x} appear in all conjunctions is not a limitation, since we can always add to a conjunction new atoms $\top_1(x)$, one for each variable x that does not originally appear in the conjunction.

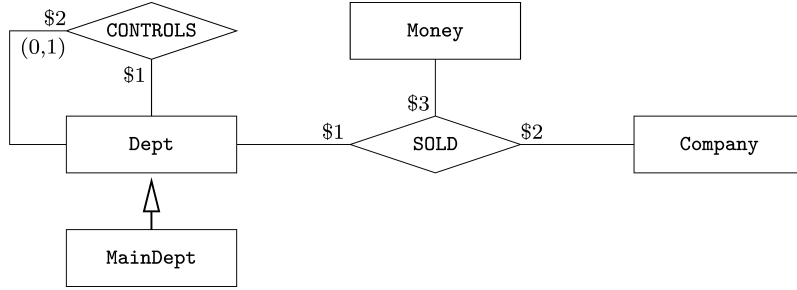


Fig. 2. The entity-relationship diagram for the example in Section 2.3.

(with respect to a schema) is the problem of checking whether $S \not\models q \equiv \emptyset$, where S and q are given as input.

We provide now an alternative characterization of query containment, on which we base the technique presented in Section 3. Given a $D\mathcal{L}\mathcal{R}_{reg}$ schema S and two queries for S

$$\begin{aligned}
 q(\vec{x}) &\leftarrow conj_1(\vec{x}, \vec{y}_1, \vec{c}_1) \vee \dots \vee conj_m(\vec{x}, \vec{y}_m, \vec{c}_m), \\
 q'(\vec{x}) &\leftarrow conj'_1(\vec{x}, \vec{y}'_1, \vec{c}'_1) \vee \dots \vee conj'_{m'}(\vec{x}, \vec{y}'_{m'}, \vec{c}'_{m'}),
 \end{aligned}$$

we have that $S \models q \subseteq q'$ iff there is no model \mathcal{I} of S such that, when assigning suitable objects in $\Delta^{\mathcal{I}}$ to $\vec{x}, \vec{y}_1, \dots, \vec{y}_m$, the formula

$$\begin{aligned}
 &(conj_1(\vec{x}, \vec{y}_1, \vec{c}_1) \vee \dots \vee conj_m(\vec{x}, \vec{y}_m, \vec{c}_m)) \wedge \\
 &\neg \exists \vec{z}_1. conj'_1(\vec{x}, \vec{z}_1, \vec{c}'_1) \wedge \dots \wedge \neg \exists \vec{z}_{m'}. conj'_{m'}(\vec{x}, \vec{z}_{m'}, \vec{c}'_{m'})
 \end{aligned}$$

evaluates to true in \mathcal{I} . In other words, $S \models q \subseteq q'$ if and only if there is no model of S that makes the formula

$$\begin{aligned}
 &(conj_1(\vec{a}, \vec{b}_1, \vec{c}_1) \vee \dots \vee conj_m(\vec{a}, \vec{b}_m, \vec{c}_m)) \wedge \\
 &\neg \exists \vec{z}_1. conj'_1(\vec{a}, \vec{z}_1, \vec{c}'_1) \wedge \dots \wedge \neg \exists \vec{z}_{m'}. conj'_{m'}(\vec{a}, \vec{z}_{m'}, \vec{c}'_{m'})
 \end{aligned}$$

true, where $\vec{a}, \vec{b}_1, \dots, \vec{b}_m$ are Skolem constants, that is, constants not appearing elsewhere for which the unique name assumption does not hold.

2.3 Example

Consider an application where the departments of a given company can be controlled by other departments, and sold to companies. Every department is controlled by at most one department, and by at least one main department, possibly indirectly. A main department is not controlled by any department. If a main department is sold, then all the departments controlled by it are also sold. Finally, if a department is sold, then there is a main department controlling it, either directly or indirectly, that is also sold.

The basic concepts and relations are shown in Figure 2 in the form of an entity-relationship diagram. The specification of the application in $D\mathcal{L}\mathcal{R}_{reg}$ makes use of the concepts Dept, MainDept, Money, Company, and the relations CONTROLS, SOLD. In particular, CONTROLS(x, y) means that department x has control over department y , and SOLD(x, y, z) means that department

x has been sold to company y at price z . The schema \mathcal{S} is constituted by the following assertions:

$$\begin{aligned} \text{SOLD} &\sqsubseteq (\$1:\text{Dept}) \sqcap (\$2:\text{Company}) \sqcap (\$3:\text{Money}), \\ \text{CONTROLS} &\sqsubseteq (\$1:\text{Dept}) \sqcap (\$2:\text{Dept}), \\ \text{Dept} &\sqsubseteq (\leq 1[\$2]\text{CONTROLS}) \sqcap \exists(\text{CONTROLS}|_{\$2,\$1})^*.\text{MainDept}, \\ \text{MainDept} &\sqsubseteq \text{Dept} \sqcap \neg\exists[\$2]\text{CONTROLS}, \\ \text{MainDept} \sqcap \exists[\$1]\text{SOLD} &\sqsubseteq \forall(\text{CONTROLS}|_{\$1,\$2})^*.\exists[\$1]\text{SOLD}, \\ \text{Dept} \sqcap \exists[\$1]\text{SOLD} &\sqsubseteq \exists(\text{CONTROLS}|_{\$2,\$1})^*.(\text{MainDept} \sqcap \exists[\$1]\text{SOLD}). \end{aligned}$$

The first two assertions are used to specify the types of the attributes of the relations. The third and the fourth assertions specify the basic properties of Dept and MainDept . It is easy to see that such assertions imply that, in all the models of \mathcal{S} , the set of CONTROLS links starting from an instance m of MainDept form a tree (which we call CONTROLS -tree) with root m . The role of the transitive closure $(\text{CONTROLS}|_{\$2,\$1})^*$ and the number restrictions is crucial for correctly representing the above property in the schema. Finally, the last two assertions, each one stating inclusions between views, specify the company policy for selling departments. Note again the use of the transitive closure for this purpose.

We now consider two queries for the schema \mathcal{S} . The first query, called q , is used to retrieve all departments that control two departments that have been sold to the same company and such that one of them controls the other. The second query, called q' , retrieves all the departments that have been sold. The queries q and q' are defined as follows:

$$\begin{aligned} q(x) &\leftarrow \text{CONTROLS}(x, y) \wedge \text{SOLD}(y, z, z_1) \wedge \text{CONTROLS}(x, w) \wedge \text{SOLD}(w, z, z_2) \wedge \\ &\quad \text{CONTROLS}(y, w), \\ q'(x) &\leftarrow \text{Dept}(x) \wedge \text{SOLD}(x, z, z'). \end{aligned}$$

The schema \mathcal{S} imposes that (i) the CONTROLS relation is typed, so that x in q is a department; (ii) each department is controlled by at most one department; (iii) when a department is sold, there is a main department (possibly indirectly) controlling it that is also sold; (iv) when a main department is sold, then all departments it (possibly indirectly) controls are also sold. From (i)–(iv) it follows that a department controlling one that is sold, is sold as well. Therefore, in every model of \mathcal{S} each x satisfying q also satisfies q' , and hence $\mathcal{S} \models q \subseteq q'$.

Following the same arguments, it is easy to see that the following query q'' is unsatisfiable with respect to \mathcal{S} :

$$q''(x) \leftarrow \text{CONTROLS}(x, y) \wedge \text{SOLD}(y, z_1, z_2) \wedge \neg \text{SOLD}(x, w_1, w_2).$$

3. CHECKING QUERY CONTAINMENT

We address the problem of deciding, given a schema \mathcal{S} and two queries q and q' of the same arity, whether $\mathcal{S} \models q \subseteq q'$. To do so, we make use of a reduction of query containment to a problem of unsatisfiability in a variant of Propositional Dynamic Logic, called CPDL_g .

The reduction we present is based on the combination of several techniques, some of which stem from work in the description logic community, and some

of which are novel to the present research. First, the general idea of reducing reasoning in description logics to reasoning in propositional dynamic logics has been widely used since Schild [1991] and De Giacomo and Lenzerini [1994] to show decidability and complexity results [Calvanese and De Giacomo 2003]. Then, we exploit *reification* (cf. construction of Φ_S in Section 3.2) to deal with relations of arbitrary arity in a formalism whose semantic structures are based on graphs, that is, unary and binary relations only. Such a technique was introduced in Catarci and Lenzerini [1993] and De Giacomo and Lenzerini [1995] and is at the basis of the tight relationship between description logics and conceptual models for databases [Calvanese et al. 1999; Borgida and Brachman 2003]. Another technique we use is the encoding of formulae denoting single objects (here called *name-formulae*, a.k.a. nominals) within standard propositional dynamic logics. These techniques were introduced in De Giacomo and Lenzerini [1994, 1996] to deal with knowledge on single individuals in expressive description logics. Here, we directly build on the results in De Giacomo and Lenzerini [1996] and use the technique to deal with what is essentially a “disjunction of ABoxes,” instead of a single ABox (cf. construction of Φ_{aux} and of $\bigvee_{j=1}^m \Phi_{conj_j}$ in Section 3.2). The core technique to deal with containment is a novelty of this research. It is based on realizing that we can restrict our attention to models in which the only cycles are those formed by explicitly named individuals, and hence cyclic queries can only be satisfied through these individuals (cf. construction of $\bigwedge_{j=1}^m \neg \Phi_{conj_j}$ in Section 3.2).

In the next subsection, we introduce CPDL_g . Then we present the reduction, prove its correctness, and analyze the computational complexity of the resulting containment algorithm.

3.1 The Propositional Dynamic Logic CPDL_g

Propositional dynamic logics are specific modal logics originally proposed as a formal system for reasoning about computer program schemas [Fischer and Ladner 1979]. Since then PDLs have been studied extensively and extended in several ways (see, e.g., Kozen and Tiuryn [1990] for a survey).

Here, we make use of CPDL_g (studied in De Giacomo and Lenzerini [1996] in the context of description logics), which is an extension of Converse PDL [Kozen and Tiuryn 1990] with *graded modalities* [Fattorosi-Barnaba and De Caro 1985]. The syntax of CPDL_g is as follows (A denotes an *atomic formula*, ϕ an arbitrary *formula*, p an *atomic program*, and r an arbitrary *program*):

$$\begin{aligned} \phi &::= A \mid \neg\phi \mid \phi_1 \wedge \phi_2 \mid \langle r \rangle \phi \mid [p]_{\leq k} \phi \mid [p^-]_{\leq k} \phi, \\ r &::= p \mid \varepsilon \mid r_1; r_2 \mid r_1 \cup r_2 \mid r^* \mid \phi? \mid r^- . \end{aligned}$$

We use the standard abbreviations, namely, T for true, F for false, \vee for disjunction, \Rightarrow for material implication, and $[r]\phi$ for $\neg\langle r \rangle\neg\phi$.

As usual for PDLs, the semantics of CPDL_g is based on *Kripke structures* $\mathcal{M} = (S, \cdot^{\mathcal{M}})$, where S is a *set of states* and $\cdot^{\mathcal{M}}$ is a mapping interpreting each atomic formula A as a subset $A^{\mathcal{M}}$ of S and each atomic program p as a binary relation $p^{\mathcal{M}}$ over S . The mapping $\cdot^{\mathcal{M}}$ is then extended to arbitrary formulae and programs in such a way that the conditions in Figure 3 are satisfied.

$$\begin{aligned}
A^{\mathcal{M}} &\subseteq S \\
(\neg\phi)^{\mathcal{M}} &= S \setminus \phi^{\mathcal{M}} \\
(\phi_1 \wedge \phi_2)^{\mathcal{M}} &= \phi_1^{\mathcal{M}} \cap \phi_2^{\mathcal{M}} \\
\langle r \rangle \phi^{\mathcal{M}} &= \{s \mid \exists s'. (s, s') \in r^{\mathcal{M}} \wedge s' \in \phi^{\mathcal{M}}\} \\
([p]_{\leq k} \phi)^{\mathcal{M}} &= \{s \mid \#\{s' \mid (s, s') \in p^{\mathcal{M}} \wedge s' \in \phi^{\mathcal{M}}\} \leq k\} \\
([p^-]_{\leq k} \phi)^{\mathcal{M}} &= \{s \mid \#\{s' \mid (s', s) \in p^{\mathcal{M}} \wedge s' \in \phi^{\mathcal{M}}\} \leq k\} \\
p^{\mathcal{M}} &\subseteq S \times S \\
\varepsilon^{\mathcal{M}} &= \{(x, x) \mid x \in S\} \\
(r_1; r_2)^{\mathcal{M}} &= r_1^{\mathcal{M}} \circ r_2^{\mathcal{M}} \\
(r_1 \cup r_2)^{\mathcal{M}} &= r_1^{\mathcal{M}} \cup r_2^{\mathcal{M}} \\
(r^*)^{\mathcal{M}} &= (r^{\mathcal{M}})^* = \bigcup_{i \geq 0} (r^{\mathcal{M}})^i \\
(\phi?)^{\mathcal{M}} &= \{(s, s) \mid s \in \phi^{\mathcal{M}}\} \\
(r^-)^{\mathcal{M}} &= \{(s, s') \mid (s', s) \in r^{\mathcal{M}}\}
\end{aligned}$$

Fig. 3. Semantic rules for CPDL_g .

It can be shown that CPDL_g has typical properties of PDLs, in particular the *connected-model property* (if a formula has a model, then it has one that is connected when viewing it as a graph), the *tree-model property* (if a formula has a model, then it has one that is a tree when viewing it as an undirected graph), and *EXPTIME-completeness* of checking satisfiability of a formula (with the assumption that numbers in graded modalities are represented in unary) [De Giacomo and Lenzerini 1996; Calvanese et al. 2001b; Baader et al. 2003].

In the following we make use of two notions related to CPDL_g , namely, the Fischer-Ladner closure of a formula and the prefix of a program. Given a CPDL_g formula ϕ , the *Fischer-Ladner closure* $CL(\phi)$ of ϕ is defined as the smallest set C of formulae containing ϕ and such that

$$\begin{array}{ll}
\neg\phi' \in C & \text{implies } \phi' \in C, \\
\phi' \in C & \text{implies } \neg\phi' \in C \quad (\text{if } \phi' \neq \neg\phi''), \\
\phi_1 \wedge \phi_2 \in C & \text{implies } \phi_1, \phi_2 \in C, \\
\langle r \rangle \phi' \in C & \text{implies } \phi' \in C, \\
[p]_{\leq k} \phi \in C & \text{implies } \phi \in C, \\
[p^-]_{\leq k} \phi \in C & \text{implies } \phi \in C, \\
\langle r_1; r_2 \rangle \phi' \in C & \text{implies } \langle r_1 \rangle \langle r_2 \rangle \phi' \in C, \\
\langle r_1 \cup r_2 \rangle \phi' \in C & \text{implies } \langle r_1 \rangle \phi', \langle r_2 \rangle \phi' \in C, \\
\langle r^* \rangle \phi' \in C & \text{implies } \langle r \rangle \langle r^* \rangle \phi' \in C, \\
\langle \phi? \rangle \phi'' \in C & \text{implies } \phi' \in C.
\end{array}$$

Given a CPDL_g program r , the set $Pre(r)$ of *prefixes* of r is defined as follows:

$$\begin{aligned}
Pre(p) &= \{\varepsilon, p\}, \\
Pre(p^-) &= \{\varepsilon, p^-\}, \\
Pre(r_1; r_2) &= Pre(r_1) \cup \{r_1; r'_2 \mid r'_2 \in Pre(r_2)\}, \\
Pre(r_1 \cup r_2) &= Pre(r_1) \cup Pre(r_2), \\
Pre(r_1^*) &= \{r_1^*; r'_1 \mid r'_1 \in Pre(r_1)\}, \\
Pre(\phi?) &= \{\varepsilon, \phi?\}.
\end{aligned}$$

	$\sigma(\top_n) = \top_n$
	$\sigma(\mathbf{P}) = \mathbf{P}$
$\sigma(\top_1) = \top_1$	$\sigma((i/n:C)) = \top_n \wedge [f_i]\sigma(C)$
$\sigma(A) = A$	$\sigma(\neg\mathbf{R}) = \top_n \wedge \neg\sigma(\mathbf{R})$
$\sigma(\neg C) = \top_1 \wedge \neg\sigma(C)$	$\sigma(\mathbf{R}_1 \sqcap \mathbf{R}_2) = \sigma(\mathbf{R}_1) \wedge \sigma(\mathbf{R}_2)$
$\sigma(C_1 \sqcap C_2) = \sigma(C_1) \wedge \sigma(C_2)$	$\sigma(\varepsilon) = \varepsilon$
$\sigma(\exists E.C) = \langle \sigma(E) \rangle \sigma(C)$	$\sigma(R _{\$i, \$j}) = f_i^- ; \sigma(R)? ; f_j$
$\sigma(\exists[\$i]\mathbf{R}) = \langle f_i^- \rangle \sigma(\mathbf{R})$	$\sigma(E_1 \circ E_2) = \sigma(E_1); \sigma(E_2)$
$\sigma(\leq k[\$i]\mathbf{R}) = [f_i^-]_{\leq k} \sigma(\mathbf{R})$	$\sigma(E_1 \sqcup E_2) = \sigma(E_1) \cup \sigma(E_2)$
	$\sigma(E^*) = \sigma(E)^*$

 Fig. 4. Mapping $\sigma(\cdot)$ from \mathcal{DLR}_{reg} to CPDL_g .

3.2 Reduction of Query Containment to Unsatisfiability in CPDL_g

Our aim is to reduce query containment to a problem of unsatisfiability in CPDL_g . To this end, we construct a CPDL_g formula starting from an instance of the query containment problem. More precisely, if we have to check whether there is no model of S that makes the formula

$$\begin{aligned} & (\text{conj}_1(\vec{\mathbf{a}}, \vec{\mathbf{b}}_1, \vec{\mathbf{c}}_1) \vee \dots \vee \text{conj}_m(\vec{\mathbf{a}}, \vec{\mathbf{b}}_m, \vec{\mathbf{c}}_m)) \wedge \\ & \neg \exists \vec{\mathbf{z}}_1. \text{conj}'_1(\vec{\mathbf{a}}, \vec{\mathbf{z}}_1, \vec{\mathbf{c}}'_1) \wedge \dots \wedge \neg \exists \vec{\mathbf{z}}_{m'}. \text{conj}'_{m'}(\vec{\mathbf{a}}, \vec{\mathbf{z}}_{m'}, \vec{\mathbf{c}}'_{m'}) \end{aligned}$$

true, where $\vec{\mathbf{a}}, \vec{\mathbf{b}}_1, \dots, \vec{\mathbf{b}}_m$ are Skolem constants, we check the unsatisfiability of the CPDL_g formula

$$\Phi_{S \neq q \subseteq q'} = \Phi_S \wedge \left(\bigvee_{j=1}^m \Phi_{\text{conj}_j} \right) \wedge \left(\bigwedge_{j=1}^{m'} \neg \Phi_{\text{conj}'_j} \right) \wedge \Phi_{aux},$$

constructed as described below.

3.2.1 Φ_S : Encoding of S . Φ_S is the translation of S into a CPDL_g formula that is based on *reification* of n -ary relations, that is, a tuple in a model of S is represented in a model of $\Phi_{S \neq q \subseteq q'}$ by a state having one functional link f_i for each tuple component $\$i$. Let n_{max} be the maximum arity of relations in S . The formula Φ_S makes use of the mapping $\sigma(\cdot)$ from \mathcal{DLR}_{reg} expressions to CPDL_g formulae defined in Figure 4. The atomic formula \top_1 denotes those states that represent objects, while each atomic formula \top_n , with $n \in \{2, \dots, n_{max}\}$, denotes those states that represent tuples of arity n . We denote with U the program $(create \cup f_1 \cup \dots \cup f_{n_{max}} \cup create^- \cup f_1^- \cup \dots \cup f_{n_{max}}^-)^*$, where $create, f_1, \dots, f_{n_{max}}$ are all atomic programs used in $\Phi_{S \neq q \subseteq q'}$. Hence, U is the so-called *master modality* [Blackburn et al. 2001], which in our case, due to the connected-model property of CPDL_g , represents the universal accessibility relation. Therefore, for a given interpretation, $[U]\phi$ expresses that ϕ holds in every state, and $\langle U \rangle \phi$ expresses that ϕ holds in some state.

Φ_S is the conjunction of the following formulae:

$$[U](\top_1 \vee \dots \vee \top_{n_{max}}), \quad (1)$$

$$[U][f_i]_{\leq 1} \top, \quad \text{for each } i \in \{1, \dots, n_{max}\}, \quad (2)$$

$$[U](\top_n \equiv \langle f_1 \rangle \top_1 \wedge \dots \wedge \langle f_n \rangle \top_1 \wedge [f_{n+1}] \mathbf{F}), \quad \text{for each } n \in \{2, \dots, n_{max}\}, \quad (3)$$

$$[U]([f_i] \mathbf{F} \Rightarrow [f_{i+1}] \mathbf{F}), \quad \text{for each } i \in \{1, \dots, n_{max}\}, \quad (4)$$

$$[U](A \Rightarrow \top_1), \quad \text{for each atomic concept } A, \quad (5)$$

$$[U](\mathbf{P} \Rightarrow \top_n), \quad \text{for each atomic relation } \mathbf{P} \text{ of arity } n, \quad (6)$$

$$[U](\sigma(C_1) \Rightarrow \sigma(C_2)), \quad \text{for each assertion } C_1 \sqsubseteq C_2 \text{ in } S, \quad (7)$$

$$[U](\sigma(\mathbf{R}_1) \Rightarrow \sigma(\mathbf{R}_2)), \quad \text{for each assertion } \mathbf{R}_1 \sqsubseteq \mathbf{R}_2 \text{ in } S. \quad (8)$$

Formula (1) above expresses that each state represents an object or a tuple of arity between 2 and n_{max} . Formula (2) expresses that all programs f_i are functional (i.e., deterministic). Formulae (3) and (4) express that the states representing tuples of arity n are exactly those connected through programs f_1, \dots, f_n to states representing objects, and not connected via programs f_i , with $i > n$, to any state. Formulae (5) and (6) express that states satisfying atomic propositions corresponding to atomic concepts (respectively, atomic relations of arity n) are states representing objects (respectively, tuples of arity n). Finally, formulae (7) and (8) encode the assertions in S .

3.2.2 Φ_{conj_j} : Encoding of Each $conj_j(\vec{\mathbf{a}}, \vec{\mathbf{b}}, \vec{\mathbf{c}}_j)$. For each $j \in \{1, \dots, m\}$, the encoding Φ_{conj_j} of $conj_j(\vec{\mathbf{a}}, \vec{\mathbf{b}}_j, \vec{\mathbf{c}}_j)$ makes use of special atomic propositions, called *name-formulae*, whose distinguishing properties are specified by Φ_{aux} (see later). Specifically, one name-formula N_t is introduced for each term t in $\vec{\mathbf{a}}, \vec{\mathbf{b}}_j, \vec{\mathbf{c}}_j$, and one name-formula $N_{\vec{\mathbf{t}}}$ for each tuple $\vec{\mathbf{t}}$ such that for some \mathbf{R} , $\mathbf{R}(\vec{\mathbf{t}})$ appears in $conj_j(\vec{\mathbf{a}}, \vec{\mathbf{b}}_j, \vec{\mathbf{c}}_j)$. A name-formula assigns a name to a term t (respectively, tuple $\vec{\mathbf{t}}$), which allows for identifying in a model certain states which correspond to t (respectively, reified counterpart of $\vec{\mathbf{t}}$). The distinguishing properties of name-formulae guarantee that these states share some crucial properties that allow us to isolate a single state as a representative of t (respectively, $\vec{\mathbf{t}}$).

Once we have name-formulae in place, we define Φ_{conj_j} as the conjunction of the following formulae:

- (1) for each atom $C(t)$ in $conj_j(\vec{\mathbf{a}}, \vec{\mathbf{b}}_j, \vec{\mathbf{c}}_j)$

$$[U](N_t \Rightarrow \sigma(C)),$$
- (2) for each atom $\mathbf{R}(\vec{\mathbf{t}})$ in $conj_j(\vec{\mathbf{a}}, \vec{\mathbf{b}}_j, \vec{\mathbf{c}}_j)$

$$[U](N_{\vec{\mathbf{t}}} \Rightarrow \sigma(\mathbf{R})).$$

Intuitively, each formula (1) and (2) expresses that the states satisfying the name-formula corresponding to a term (respectively, tuple) appearing in an atom, satisfy also the formula corresponding to the predicate of the atom.

3.2.3 $\Phi_{conj'_j}$: Encoding of Each $\exists_{\vec{\mathbf{z}}_j}. conj'_j(\vec{\mathbf{a}}, \vec{\mathbf{z}}_j, \vec{\mathbf{c}}'_j)$. Now consider a $j \in \{1, \dots, m'\}$. We construct the formula $\Phi_{conj'_j}$ as a disjunction of formulae, one

for each possible partition of the variables \vec{z}_j in $\exists \vec{z}_j.conj'_j(\vec{a}, \vec{z}_j, \vec{c}'_j)$. More precisely, to build one such formula, we consider a partition π of the variables \vec{z}_j . Then, for each equivalence class in the partition we choose a variable as a representative, and substitute in $\exists \vec{z}_j.conj'_j(\vec{a}, \vec{z}_j, \vec{c}'_j)$ all other variables in the same equivalence class by the representative, thus obtaining a formula $\exists \vec{w}_\pi.conj'_j(\vec{a}, \vec{w}_\pi, \vec{c}'_j)$. Now, from such a formula we build a corresponding $CPDL_g$ formula by making use of a special graph, called a *tuple-graph*, which intuitively reflects the dependencies between variables and tuples resulting from the appearance of the variables in the atoms of $\exists \vec{w}_\pi.conj'_j(\vec{a}, \vec{w}_\pi, \vec{c}'_j)$.⁵ A tuple-graph is a directed graph with nodes labeled by $CPDL_g$ formulae and edges labeled by $CPDL_g$ programs, formed as follows:

- There is one node t for each term t in $\vec{a}, \vec{w}, \vec{c}_j$, and one node \vec{t} for each tuple \vec{t} such that $\mathbf{R}(\vec{t})$ appears in $\exists \vec{w}_\pi.conj'_j(\vec{a}, \vec{w}_\pi, \vec{c}'_j)$. Each node t is labeled by all $\sigma(C)$ such that $C(t)$ appears in $\exists \vec{w}_\pi.conj'_j(\vec{a}, \vec{w}_\pi, \vec{c}'_j)$. Each node \vec{t} is labeled by all $\sigma(\mathbf{R})$ such that $\mathbf{R}(\vec{t})$ appears in $\exists \vec{w}_\pi.conj'_j(\vec{a}, \vec{w}_\pi, \vec{c}'_j)$.
- There is one edge labeled by f_i from the node $\vec{t} = (t_1, \dots, t_n)$ to the node t_i , $i \in \{1, \dots, n\}$, for each tuple \vec{t} such that $\mathbf{R}(\vec{t})$ appears in $\exists \vec{w}_\pi.conj'_j(\vec{a}, \vec{w}_\pi, \vec{c}'_j)$.

Notice that dividing the variables \vec{z}_j in $\exists \vec{z}_j.conj'_j(\vec{a}, \vec{z}_j, \vec{c}'_j)$ in all possible ways into equivalence classes and replacing equivalent variables by one representative, corresponds to introducing in all possible ways equalities between variables. Such equalities allow us to take into account that a cycle in the tuple graph can in fact be eliminated, and become simply a chain, when different variables are assigned the same object. As will become clear in the following, the distinction between variables appearing in cycles in the tuple-graph and those that do not, is indeed necessary for the correctness of the proposed technique for query containment under constraints.

In the following, we call *formula-template* a $CPDL_g$ formula in which formula-placeholders occur that later will be substituted by actual formulae. From the tuple-graph G of $\exists \vec{w}_\pi.conj'_j(\vec{a}, \vec{w}_\pi, \vec{c}'_j)$, we build a $CPDL_g$ formula-template δ , and to do so we have to consider that in general the tuple-graph is composed of several connected components. For the i th connected component we build a formula-template δ_i by choosing a starting node t_0 (corresponding to a term) and performing a depth-first visit of the corresponding component and building the formula in a postorder fashion. We describe the construction by defining a visiting function V , which, given a node of the tuple-graph G , returns the corresponding formula-template, and as a side effect marks the nodes of the graph that it visits.

- If $u = t$, then $V(t)$ marks t , and returns the conjunction of
 - (i) t itself, used as a placeholder, and every formula labeling the node t ;
 - (ii) for each edge (\vec{t}, t) labeled by f_i (i.e., $t = t_i$ in \vec{t}) such that \vec{t} is not marked yet, the formula $\langle f_i^- \rangle V(\vec{t})$.

⁵The tuple-graph is similar to the graph used in Chekuri and Rajaraman [1997] to detect cyclic dependencies between variables.

- If $u = \vec{t} = (t_1, \dots, t_n)$, then $V(\vec{t})$ marks \vec{t} , and returns the conjunction of
- (i) \vec{t} itself, used as a placeholder, and every formula labeling the node \vec{t} ;
 - (ii) for each edge (\vec{t}, t_i) labeled by f_i , such that t_i is not marked yet, the formula $\langle f_i \rangle V(t_i)$;
 - (iii) for each edge (\vec{t}, t_i) labeled by f_i , such that t_i is already marked, the formula $\langle f_i \rangle t_i$.

Then the formula-template δ_i for the i th connected component is defined as $V(t_0)$, where t_0 is the starting node chosen for the visit.

The formula-template δ for the whole tuple-graph G , composed of $\ell \geq 1$ connected components, is

$$\langle U \rangle \delta_1 \wedge \dots \wedge \langle U \rangle \delta_\ell,$$

where $\delta_1, \dots, \delta_\ell$ are the formula-templates corresponding to all the connected components in the tuple-graph G .

We next define the notion of G -substitution for the formula-template δ , which allows us to obtain from δ a CPDL_G formula. Let G be the tuple-graph for $\exists \vec{w}_\pi \text{conj}'_j(\vec{a}, \vec{w}_\pi, \vec{c}'_j)$, and δ as above, a G -substitution for δ is a substitution for the placeholders of δ that replaces:

- (i) each placeholder \vec{t} by \top_n , where n is the arity of the tuple \vec{t} ;
- (ii) each placeholder t in \vec{a}, \vec{c}'_j by the name-formula N_t ;
- (iii) each placeholder w corresponding to a variable not occurring in a cycle in the tuple-graph by \top_1 ;
- (iv) each placeholder w corresponding to a variable occurring in a cycle in the tuple-graph by one of the name-formulae N_t corresponding to a term in $\vec{a}, \vec{b}_1, \dots, \vec{b}_m, \vec{c}_1, \dots, \vec{c}_m$ occurring in q .

Notice that there are several G -substitutions for a formula-template δ , one for each possible way of choosing (possibly Skolem) constants according to (iv). Hence, the number of such G -substitutions is $O(\ell_1^{\ell_2})$, where ℓ_1 is the number of variables and constants in q , and ℓ_2 is the number of variables occurring in a cycle in the tuple-graph G .

We are now ready to define the CPDL_G formula φ_π corresponding to a partition π of the variables \vec{z}_j : such a formula φ_π is the disjunction of all formulae obtained by applying to the formula-template δ a possible G -substitution (note that both δ and G are specific for the partition π .)

Since φ_π corresponds to one possible partition of the variables \vec{z}_j , we obtain the formula $\Phi_{\text{conj}'_j}$ as the disjunction of all formulae φ_π , one for each possible partition π of the variables \vec{z}_j . The number of such disjuncts is $O(2^{\ell_2})$, where ℓ_2 is the number of variables \vec{z}_j . Therefore, the total number of disjuncts for $\Phi_{\text{conj}'_j}$ is $O(\ell_1^{O(\ell_2)})$.

3.2.4 Φ_{aux} : Encoding of Constants and Variables. Let $\Phi' = \Phi_S \wedge (\bigvee_{j=1}^m \Phi_{\text{conj}'_j}) \wedge (\bigwedge_{j=1}^{m'} \neg \Phi_{\text{conj}'_j})$. Φ_{aux} is formed by the conjunction of

- (1) the formula $\langle create \rangle N$, for each name-formula N appearing in Φ' , which expresses the existence of a state satisfying N ;
- (2) one formula of the form $[U](N_{c_i} \Rightarrow \neg N_{c_j})$ for each pair of distinct constants c_i, c_j appearing in the queries (not Skolem constants);
- (3) for each name-formula $N_{\vec{t}}$ corresponding to a tuple $\vec{t} = (t_1, \dots, t_n)$ appearing in $\bigvee_{j=1}^m \Phi_{conj_j}$ (observe that we do not have any name-formula for tuples in $\bigwedge_{j=1}^{m'} \neg \Phi_{conj'_j}$)

$$[U](N_{\vec{t}} \equiv \langle f_1 \rangle N_{t_1} \wedge \dots \wedge \langle f_n \rangle N_{t_n} \wedge [f_{n+1}]F),$$

$$[U](N_{t_i} \Rightarrow (\langle f_i^- \rangle N_{\vec{t}} \wedge [f_i^-]_{\leq 1} N_{\vec{t}})), \quad \text{for each } i \in \{1, \dots, n\},$$

these conjuncts express the relationships between the name-formulae for tuples and those for tuple components;

- (4) one formula of the form $[U](N \wedge \phi \Rightarrow [U](N \Rightarrow \phi))$ for each name-formula N and each formula ϕ such that
 - $\phi \in CL(\Phi')$,
 - $\phi = \langle \bar{r} \rangle \phi'$ with $\langle r \rangle \phi' \in CL(\Phi')$, and
 - $\phi = \langle \bar{r}' ; p \rangle N'$ with $r' \in Pre(r)$, $p = f \mid f^-$, and r, f, N' occurring in $CL(\Phi')$,
 where \bar{r} is defined inductively as follows:

$$\begin{aligned} \bar{p} &= p; (\wedge_i \neg N_i)?, \\ \bar{r_1}; \bar{r_2} &= \bar{r_1}; \bar{r_2}, \\ \overline{r_1 \cup r_2} &= \bar{r_1} \cup \bar{r_2}, \\ \bar{r_1}^* &= \bar{r_1}^*, \\ \bar{\phi}^? &= \phi?. \end{aligned}$$

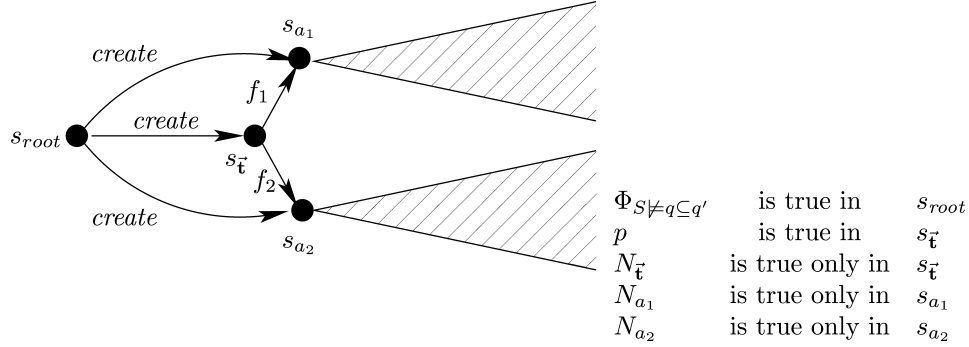
The role of the conjuncts (4) in Φ_{aux} is to enforce that, in every model of $\Phi_{S \not\subseteq q \subseteq q'}$, for each name-formula N , one representative state can be singled out among those satisfying N . This would be trivially obtained if we could force all these states to satisfy exactly the same formulae of the logic. Φ_{aux} forces a weaker condition, namely, that these states satisfy the same formulae in the finite set (whose size is polynomial wrt Φ') described above. Theorem 3.3 shows that this is sufficient for our purposes.

3.2.5 Example. We illustrate the encoding of the containment problem $S \models q \subseteq q'$ into unsatisfiability of the CPDL_g formula $\Phi_{S \not\subseteq q \subseteq q'}$ by means of a simple example.

Consider two queries

$$\begin{aligned} q(x_1, x_2) &\leftarrow p(x_1, x_2), \\ q'(x_1, x_2) &\leftarrow r(x_1, x_2, z), \end{aligned}$$

over a schema S such that $S \not\models q \subseteq q'$. Let Φ_S be the formula encoding S , Φ_{conj} the formula encoding $p(a_1, a_2)$, where a_1, a_2 are two Skolem constants, and Φ_{aux} as specified above. Figure 5 schematically shows a model of the formula $\Phi_{S \not\subseteq q \subseteq q'}$ that represents a counterexample to the containment. The model contains a state $s_{\vec{t}}$, where $\vec{t} = (a_1, a_2)$, in which p holds, that, being connected to

Fig. 5. A model of $\Phi_{S \neq q \subseteq q'}$.

s_{a_1} and s_{a_2} by means of f_1 and f_2 , respectively, represents the tuple $\vec{t} = (a_1, a_2)$, which satisfies p . The objects s_{a_1} and s_{a_2} satisfy the name-formulae N_{a_1} and N_{a_2} , respectively, and there are no other objects satisfying such name-formulae. From s_{a_1} and s_{a_2} the model has a tree-like structure. We will call such model pseudotree admissible (see later). Note that Φ_{conj} is satisfied due to the part of the model involving s_{a_1} , s_{a_2} , and $s_{\vec{t}}$, while Φ_{aux} is trivially satisfied since N_{a_1} , N_{a_2} , and $N_{\vec{t}}$ are true in just one state of the model.

Now, consider the formula-template for $\exists z.r(a_1, a_2, z)$:

$$\langle U \rangle (a_1 \wedge \langle f_1^- \rangle (\vec{t}' \wedge r \wedge \langle f_2 \rangle a_2 \wedge \langle f_3 \rangle z)),$$

where $\vec{t}' = (a_1, a_2, z)$. From such a formula-template we get the formula $\Phi_{conj'} = \langle U \rangle (N_{a_1} \wedge \langle f_1^- \rangle (\vec{t}' \wedge r \wedge \langle f_2 \rangle N_{a_2} \wedge \langle f_3 \rangle \top_1))$. Note that, $s_{\vec{t}}$ has no outgoing f_3 edge, and hence \top_3 is false in $s_{\vec{t}}$; also N_{a_1} and N_{a_2} are true, respectively, in s_{a_1} and s_{a_2} only. Hence, the formula $\neg \Phi_{conj'} = [U](N_{a_1} \Rightarrow [f_1^-](\top_3 \wedge r) \Rightarrow ([f_2]N_{a_2} \wedge [f_3] \neg \top_1))$ is true in s_{root} . Therefore, in the model there is no state that both represents a tuple (a_1, a_2, z) , for some z , and that satisfies r .

3.3 Correctness of the Reduction

We now prove the correctness of the reduction of query containment to satisfiability in CPDL_g based on the encoding presented above. We first prove soundness, that is, if $\Phi_{S \neq q \subseteq q'}$ is unsatisfiable, then $S \models q \subseteq q'$ (Theorem 3.2), and then completeness, that is, if $\Phi_{S \neq q \subseteq q'}$ is satisfiable, then $S \not\models q \subseteq q'$ (Theorems 3.3 and 3.4).

To prove soundness, we need a preliminary notion. We say that a tuple-graph G is *satisfied* in an interpretation \mathcal{M} for $\Phi_{S \neq q \subseteq q'}$ if there exists a G - \mathcal{M} -homomorphism η , that is, a mapping from the nodes of G to states of \mathcal{M} such that

- if a node u of G is a (possibly Skolem) constant, then $\eta(u) \in N_u^{\mathcal{M}}$;
- if a node u of G is labeled by a formula ϕ , then $\eta(u) \in \phi^{\mathcal{M}}$;
- if an edge (u, u') of G is labeled by a program f , then $(\eta(u), \eta(u')) \in f^{\mathcal{M}}$.

Given a formula-template δ corresponding to a tuple-graph G and a G -substitution θ of its placeholders, we denote by $\delta\theta$ the formula obtained from δ by substituting the placeholders according to θ .

LEMMA 3.1. *Let G be a tuple-graph, δ the corresponding formula-template, and \mathcal{M} a CPDL_g interpretation. If there exists a G -substitution θ such that $(\delta\theta)^\mathcal{M}$ is not empty, then G is satisfied in \mathcal{M} .*

PROOF. If $(\delta\theta)^\mathcal{M}$ is not empty, then it is possible to define a mapping η as follows: for each node t (respectively, \vec{t}) in G , let s_t (respectively, $s_{\vec{t}}$) be the state of \mathcal{M} that is used in satisfying $\delta\theta$ in the position corresponding to t (respectively, \vec{t}); then $\eta(t) = s_t$ (respectively, $\eta(\vec{t}) = s_{\vec{t}}$). We show that η is a G - \mathcal{M} -homomorphism:

- if t is a (possibly Skolem) constant, then in $\delta\theta$ we have substituted the placeholder t with N_t , and hence $s_t \in N_t^\mathcal{M}$;
- if t (respectively, \vec{t}) is labeled by ϕ in G , then in δ , in conjunction with the placeholder t there is the formula ϕ , and hence $s_t \in \phi^\mathcal{M}$ (respectively, $s_{\vec{t}} \in \phi^\mathcal{M}$).
- Finally, if an edge (\vec{t}, t) is labeled by a program f_i in G , then in δ either in conjunction with t we have $\langle f_i^- \rangle \vec{t}$, or in conjunction with \vec{t} we have $\langle f_i \rangle t$, and hence $(s_{\vec{t}}, t) \in f_i^\mathcal{M}$. \square

THEOREM 3.2. *Let S be a schema, q, q' two queries of the same arity, and $\Phi_{S \not\subseteq q \subseteq q'}$ the formula obtained as specified above. If $\Phi_{S \not\subseteq q \subseteq q'}$ is unsatisfiable, then $S \models q \subseteq q'$.*

PROOF. We show that, if $S \not\subseteq q \subseteq q'$ then $\Phi_{S \not\subseteq q \subseteq q'}$ is satisfiable. To this end, we consider a model $\mathcal{I} = (\Delta^\mathcal{I}, \cdot^\mathcal{I})$ of S that makes the following formula true:

$$(\text{conj}_1(\vec{\mathbf{a}}, \vec{\mathbf{b}}_1, \vec{\mathbf{c}}_1) \vee \dots \vee \text{conj}_m(\vec{\mathbf{a}}, \vec{\mathbf{b}}_m, \vec{\mathbf{c}}_m)) \wedge \\ \neg \exists \vec{\mathbf{z}}_1. \text{conj}'_1(\vec{\mathbf{a}}, \vec{\mathbf{z}}_1, \vec{\mathbf{c}}'_1) \wedge \dots \wedge \neg \exists \vec{\mathbf{z}}_{m'}. \text{conj}'_{m'}(\vec{\mathbf{a}}, \vec{\mathbf{z}}_{m'}, \vec{\mathbf{c}}'_{m'}).$$

From \mathcal{I} , we build a reified CPDL_g interpretation $\mathcal{M} = (S, \cdot^\mathcal{M})$ for $\Phi_{S \not\subseteq q \subseteq q'}$ as follows:

- (a) $S = \Delta^\mathcal{I} \cup \{s_{root}\} \cup \bigcup_{n \in \{2, \dots, n_{max}\}} \{s_{\vec{\mathbf{a}}} \mid \vec{\mathbf{d}} \in \mathbb{T}_n^\mathcal{I}\}$;
- (b) $\mathbb{T}_n^\mathcal{M} = \{s_{(d_1, \dots, d_n)} \mid (d_1, \dots, d_n) \in \mathbb{T}_n^\mathcal{I}\}$, for each $n \in \{2, \dots, n_{max}\}$;
- (c) $f_i^\mathcal{M} = \{(s_{(d_1, \dots, d_n)}, d_i) \mid (d_1, \dots, d_n) \in \mathbb{T}_n^\mathcal{I}\}$, for each $n \in \{2, \dots, n_{max}\}$ and each $i \in \{1, \dots, n\}$;
- (d) $\mathbf{P}^\mathcal{M} = \{s_{(d_1, \dots, d_n)} \mid (d_1, \dots, d_n) \in \mathbf{P}^\mathcal{I}\}$, for each atomic relation \mathbf{P} ;
- (e) $\mathbb{T}_1^\mathcal{M} = \Delta^\mathcal{I} \cup \{s_{root}\}$;
- (f) $A^\mathcal{M} = A^\mathcal{I}$, for each atomic concept A ;
- (g) $N_t^\mathcal{M} = \{d\}$, for the element domain $d = t^\mathcal{I}$, for each (possibly Skolem) constant t in $\vec{\mathbf{a}}, \vec{\mathbf{b}}_1, \dots, \vec{\mathbf{b}}_m, \vec{\mathbf{c}}_1, \dots, \vec{\mathbf{c}}_m, \vec{\mathbf{c}}'_1, \dots, \vec{\mathbf{c}}'_{m'}$;
- (h) $N_{\vec{t}}^\mathcal{M} = \{s_{\vec{\mathbf{a}}}\}$, for the tuple $\vec{\mathbf{d}} = \vec{t}^\mathcal{I}$ of element domains, for each tuple \vec{t} of (possibly Skolem) constants occurring in $\text{conj}_1(\vec{\mathbf{a}}, \vec{\mathbf{b}}_1, \vec{\mathbf{c}}_1) \vee \dots \vee \text{conj}_m(\vec{\mathbf{a}}, \vec{\mathbf{b}}_m, \vec{\mathbf{c}}_m)$;
- (i) $\text{create}^\mathcal{M} = \{(s_{root}, s) \mid s \in N^\mathcal{M}\}$, for some name formula N .

Next we show that $\Phi_{S \neq q \subseteq q'}$ is satisfiable by showing that $s_{root} \in \Phi_{S \neq q \subseteq q'}^{\mathcal{M}}$.

- By construction of \mathcal{M} , the interpretation of \top_n (items (b) and (e)), for $n \in \{1, \dots, n_{max}\}$, and the interpretation of f_i (item (c)), for $i \in \{2, \dots, n_{max}\}$, satisfy the conjuncts (1)–(4) of Φ_S , while the interpretation of the atoms \mathbf{P} and A (items (d) and (f)) corresponding to atomic relations and atomic concepts, respectively, satisfy the conjuncts (5)–(8), considering that \mathcal{I} is a model of S . Hence, we have that all states of \mathcal{M} , including s_{root} , satisfy Φ_S .
- By items (g) and (h), all name-formulae are interpreted as singletons; hence part (4) of the definition of Φ_{aux} is trivially satisfied; by items (g), (h), and (c), part (3) is satisfied; the interpretation of constants in \mathcal{I} guarantees that the corresponding name-formulae are disjoint in \mathcal{M} (item (g)); hence part (2) is also satisfied; finally, by item (i), all instances of the name-formulae are connected via *create* to s_{root} ; hence part (1) is satisfied as well.
- By items (f) and (d), conditions (1) and (2) of the definition of Φ_{conj_j} are satisfied if and only if $conj_j(\bar{\mathbf{a}}, \bar{\mathbf{b}}_j, \bar{\mathbf{c}}_j)$ is true in \mathcal{I} ; hence, considering that $conj_1(\bar{\mathbf{a}}, \bar{\mathbf{b}}_1, \bar{\mathbf{c}}_1) \vee \dots \vee conj_m(\bar{\mathbf{a}}, \bar{\mathbf{b}}_m, \bar{\mathbf{c}}_m)$ is true in \mathcal{I} , we have that $s_{root} \in (\bigvee_{j=1}^m \Phi_{conj_j})^{\mathcal{M}}$.

It remains to show that $s_{root} \notin \Phi_{conj'_j}^{\mathcal{M}}$, for each $j \in \{1, \dots, m'\}$. Suppose not, that is, suppose that $s_{root} \in \Phi_{conj'_j}^{\mathcal{M}}$, for some $j \in \{1, \dots, m'\}$. Then there exists a partition π of the variables $\bar{\mathbf{z}}_j$ in $\exists \bar{\mathbf{z}}_j. conj'_j(\bar{\mathbf{a}}, \bar{\mathbf{z}}_j, \bar{\mathbf{c}}'_j)$ such that $s_{root} \in \varphi_{\pi}^{\mathcal{M}}$, where φ_{π} is the formula obtained by considering the formula-template δ_{π} corresponding to the tuple-graph G associated to $\exists \bar{\mathbf{z}}_j. conj'_j(\bar{\mathbf{a}}, \bar{\mathbf{z}}_j, \bar{\mathbf{c}}'_j)$ and applying to it a G -substitution θ , that is, $\varphi_{\pi} = \delta_{\pi}\theta$. Then we have $s_{root} \in (\delta\theta)^{\mathcal{M}}$, and by Lemma 3.1 this implies that the tuple-graph G is satisfied in \mathcal{M} . Hence, by items (b)–(g) of the construction of \mathcal{M} , we would get that $\exists \bar{\mathbf{z}}_j. conj'_j(\bar{\mathbf{a}}, \bar{\mathbf{z}}_j, \bar{\mathbf{c}}'_j)$ is true in \mathcal{I} , contradicting the fact that $\neg \exists \bar{\mathbf{z}}_j. conj'_j(\bar{\mathbf{a}}, \bar{\mathbf{z}}_j, \bar{\mathbf{c}}'_j)$ is true in \mathcal{I} . \square

Next, we prove completeness of the reduction. In particular, we show that if $\Phi_{S \neq q \subseteq q'}$ is satisfiable, then it has a model of a specific form (called *pseudotree admissible*), from which a model of S that is a counterexample to the containment can be derived.

We need to introduce the following notions. We say that a model of $\Phi_{S \neq q \subseteq q'}$ is *tuple-admissible* if there is no pair of states that represent the same reified tuple. We say that a model of $\Phi_{S \neq q \subseteq q'}$ is *admissible* if it is tuple-admissible and each name-formula is true in exactly one state. We say that a model $\mathcal{M} = (S, \cdot^{\mathcal{M}})$ of $\Phi_{S \neq q \subseteq q'}$ is a *pseudotree-admissible model* if it is admissible and has the following form:

- it has a distinguished state s_{root} , and K not necessarily distinct states s_{N_1}, \dots, s_{N_K} , one for each name-formula N_i , such that $N_i^{\mathcal{M}} = \{s_{N_i}\}$;
- $create^{\mathcal{M}} = \{(s_{root}, s_{N_i}) \mid i \in \{1, \dots, K\}\}$;
- each maximal connected component of $\mathcal{M} \setminus (\{s_{root}\} \cup \{s_{N_i} \mid i \in \{1, \dots, K\}\})$ is a tree, when viewed as an undirected graph.

Notice that the subgraph induced by $\mathcal{M} \cap \{s_{N_i} \mid i \in \{1, \dots, K\}\}$ is an arbitrary graph, instead.

The following theorem shows that, with respect to satisfiability of $\Phi_{S \neq q \subseteq q'}$, one can restrict the attention to pseudotree-admissible models.

THEOREM 3.3. *Let S be a schema, q, q' two queries of the same arity, and $\Phi_{S \neq q \subseteq q'}$ the formula obtained as specified above. If $\Phi_{S \neq q \subseteq q'}$ is satisfiable then it has a pseudotree-admissible model.*

PROOF. If $\Phi_{S \neq q \subseteq q'}$ is satisfiable, then by the tree-model property of CPDL_g , $\Phi_{S \neq q \subseteq q'}$ admits a tree-model $\mathcal{M} = (S, \cdot^{\mathcal{M}})$, in which obviously there is no pair of states that represent the same reified tuple, and is therefore tuple admissible. Let $s_{root} \in \Phi^{\mathcal{M}}$ be the root of \mathcal{M} . We transform \mathcal{M} into a new model $\mathcal{M}' = (S', \cdot^{\mathcal{M}'})$ with $S' \subseteq S$, which is pseudotree admissible. \mathcal{M}' is obtained from \mathcal{M} by modifying the tree-structure of \mathcal{M} so that all name-formulae are interpreted as singletons.

In particular, to obtain \mathcal{M}' we proceed in two steps. First, we define a model $\mathcal{M}' = (S', \cdot^{\mathcal{M}'})$ by choosing, for each name-formula $N_i, i \in \{1, \dots, K\}$, a state s_{N_i} , among the states $s \in N_i^{\mathcal{M}}$ such that $(s_{root}, s) \in \text{create}^{\mathcal{M}}$, and then defining

- (a) $\text{create}^{\mathcal{M}'} = \{(s_{root}, s_{N_i}) \in \text{create}^{\mathcal{M}} \mid i \in \{1, \dots, K\}\};$
- (b) $f^{\mathcal{M}'} = (f^{\mathcal{M}} \setminus (\{(s_{N_i}, s) \in f^{\mathcal{M}} \mid s \in N_j^{\mathcal{M}}, i, j \in \{1, \dots, K\}\} \cup \{(s, s_{N_j}) \in f^{\mathcal{M}} \mid s \in N_i^{\mathcal{M}}, i, j \in \{1, \dots, K\}\})) \cup \{(s_{N_i}, s_{N_j}) \mid (s_{N_i}, s) \in f^{\mathcal{M}}, s \in N_j^{\mathcal{M}}, i, j \in \{1, \dots, K\}\},$
for each atomic program f corresponding to a tuple-component, that is, for all atomic programs except create ;
- (c) $A^{\mathcal{M}'} = A^{\mathcal{M}} \cap S'$, for each atomic formula A including name formulae;
- (d) $S' = \{s_{root}\} \cup \{s \in S \mid (s_{root}, s) \in \text{create}^{\mathcal{M}'} \circ (\bigcup_f (f^{\mathcal{M}'} \cup (f^-)^{\mathcal{M}'})^*)\}.$

The construction above chooses one representative s_{N_i} for each name-formula N_i among those connected through create to s_{root} (item (a)). Consider a pair N_1, N_2 of name-formulae related through f . This means that N_1 is a name-formula denoting a tuple $\bar{\mathbf{t}}$, and N_2 is a name-formula denoting a component t of $\bar{\mathbf{t}}$. Consider the representatives s_{N_1} and s_{N_2} . By conjuncts (2) of Φ_S , we have that s_{N_1} has a single f -successor s_2 , which is an instance of N_2 . Symmetrically, by conjuncts (3) of Φ_{aux} , we have that s_{N_2} has a single f -predecessor s_1 that is an instance of N_1 . Item (b) in the construction above disconnects s_{N_1} from s_2 and s_{N_2} from s_1 , thus cutting away the subtrees rooted at s_1 and s_2 (we are only interested in the maximally connected component of the Kripke structure that includes s_{root} , cf. item (d)), and connects s_{N_1} via f to s_{N_2} . Observe that, in doing so, we preserve the number of outgoing f -edges of s_{N_1} and of incoming f -edges of s_{N_2} .

It is possible to show, by using the construction in Lemma 5 of [De Giacomo and Lenzerini 1996],⁶ that for each $\phi \in \text{CL}(\Phi_{S \neq q \subseteq q'})$ and for each state $s \in S'$

⁶The construction in De Giacomo and Lenzerini [1996] was phrased in the Description Logic CTQ , and used to reduce ABox reasoning to satisfiability. CTQ and CPDL_g can be seen as a syntactic variant one of the other, and our handling of constants, through name-formulae, in CPDL_g is closely

we have that $s \in \phi^{\mathcal{M}'}$ if and only if $s \in \phi^{\mathcal{M}}$. Hence, since $\Phi_{S \not\equiv q \subseteq q'} \in CL(\Phi_{S \not\equiv q \subseteq q'})$, we get that $s_{root} \in (\Phi_{S \not\equiv q \subseteq q'})^{\mathcal{M}'}$.

From \mathcal{M}' we get \mathcal{M}'' by restricting the interpretation of each name-formula N_i to the representative s_{N_i} only, that is,

$$N_i^{\mathcal{M}''} = \{s_{N_i}\}.$$

Notice that $s_{root} \in \Phi_S^{\mathcal{M}''}$, since name-formulae do not appear in Φ_S . Also, $s_{root} \in (\bigvee_{j=1}^m \Phi_{conj_j})^{\mathcal{M}''}$, since $\bigvee_{j=1}^m \Phi_{conj_j}$ can be satisfied using only the representatives for name-formulae. Also, $s_{root} \Phi_{aux}^{\mathcal{M}''}$, in particular part (4) is trivially satisfied since in \mathcal{M}'' name-formulae are interpreted as singletons. Also, $s_{root} \in (\bigwedge_{j=1}^{m'} \neg \Phi_{conj'_j})^{\mathcal{M}''}$ since this holds already for \mathcal{M}' , and each $\neg \Phi_{conj'_j}$ is a disjunction in which name-formulae occur only negatively. Hence, restricting the extension of name-formulae does not restrict the extension of such a formula. \square

For pseudotree-admissible models, one can prove the “converse” of Theorem 3.2.

THEOREM 3.4. *Let S be a schema, q, q' two queries of the same arity, and $\Phi_{S \not\equiv q \subseteq q'}$ the formula obtained as specified above. If $\Phi_{S \not\equiv q \subseteq q'}$ has a pseudotree-admissible model then $S \not\equiv q \subseteq q'$.*

PROOF. We show how to construct from a pseudotree-admissible model \mathcal{M} of $\Phi_{S \not\equiv q \subseteq q'}$, a model \mathcal{I} of S in which there is a tuple $\vec{\mathbf{a}}$ of objects such that $\vec{\mathbf{a}} \in q^{\mathcal{I}}$ and $\vec{\mathbf{a}} \notin q'^{\mathcal{I}}$. \mathcal{I} is built as follows:

- $\Delta^{\mathcal{I}} = \top_1^{\mathcal{M}}$;
- $\mathbf{P}^{\mathcal{I}} = \{(s_1, \dots, s_n) \mid \exists s' \in \mathbf{P}^{\mathcal{M}}. ((s', s_i) \in f_i^{\mathcal{M}}, \text{ for } i \in \{1, \dots, n\})\}$, for each atomic relation \mathbf{P} of arity n , including \top_n ;
- $A^{\mathcal{I}} = A^{\mathcal{M}}$, for each atomic concept A , including \top_1 ;
- $t^{\mathcal{I}} = s \in N_t^{\mathcal{M}}$, for each constant and Skolem constant t in q and q' .

To show that \mathcal{I} does the job, we have to show that

- \mathcal{I} is a model of S ;
- $conj_1(\vec{\mathbf{a}}, \vec{\mathbf{b}}_1, \vec{\mathbf{c}}_1) \vee \dots \vee conj_m(\vec{\mathbf{a}}, \vec{\mathbf{b}}_m, \vec{\mathbf{c}}_m)$ is true in \mathcal{I} , that is, there is one $j \in \{1, \dots, m\}$ such that $conj_j(\vec{\mathbf{a}}, \vec{\mathbf{b}}_j, \vec{\mathbf{c}}_j)$ is true in \mathcal{I} ;
- $\neg \exists \vec{\mathbf{z}}_1. conj'_1(\vec{\mathbf{a}}, \vec{\mathbf{z}}_1, \vec{\mathbf{c}}'_1) \wedge \dots \wedge \neg \exists \vec{\mathbf{z}}_{m'}. conj'_{m'}(\vec{\mathbf{a}}, \vec{\mathbf{z}}_{m'}, \vec{\mathbf{c}}'_{m'})$ is true in \mathcal{I} , that is, for each $j \in \{1, \dots, m'\}$, we have that $\neg \exists \vec{\mathbf{z}}_j. conj'_j(\vec{\mathbf{a}}, \vec{\mathbf{z}}_j, \vec{\mathbf{c}}'_j)$ is true in \mathcal{I} .

To show that \mathcal{I} is a model of S we can exploit the fact that $\mathcal{M} = (S, \cdot^{\mathcal{I}})$ is a model of Φ_S and that, since it is admissible, there is no pair of states in S that represent the same reified tuple. By construction of \mathcal{I} it is easy to see that all assertions in S are true in \mathcal{I} .

related to handling ABoxes in CIQ , the only difference being that for constants in the ABoxes the unique name assumption is made, while here we do not make such an assumption. However, the unique name assumption played no role in the construction of De Giacomo and Lenzerini [1996]; hence that construction works in our case as well.

To show that there is one $j \in \{1, \dots, m\}$ such that $\text{conj}_j(\vec{\mathbf{a}}, \vec{\mathbf{b}}_j, \vec{\mathbf{c}}_j)$ is true in \mathcal{I} , we exploit that \mathcal{M} is an admissible model of $\Phi_{S \not\subseteq q \subseteq q'}$. Hence there is a $j \in \{1, \dots, m\}$ such that \mathcal{M} is an admissible model of Φ_{conj_j} , and since each name-formula is true in exactly one state, the claim easily follows.

It remains to show that for each $j \in \{1, \dots, m'\}$, we have that $\neg \exists \vec{\mathbf{z}}_j. \text{conj}'_j(\vec{\mathbf{a}}, \vec{\mathbf{z}}_j, \vec{\mathbf{c}}'_j)$ is true in \mathcal{I} . To do so, we show that, if for some assignment α of the variables $\vec{\mathbf{z}}_j = (z_1, \dots, z_n)$ to elements $\vec{\mathbf{d}} = (d_1, \dots, d_n)$ in $\Delta^{\mathcal{I}}$, we have that $\text{conj}'_j(\vec{\mathbf{a}}, \vec{\mathbf{z}}, \vec{\mathbf{c}}'_j)$ is true in \mathcal{I} , then we get a contradiction to \mathcal{M} being a model of $\neg \Phi_{\text{conj}'_j}$.

By considering which variables have been assigned by α to the same domain elements in $\vec{\mathbf{d}}$, we get a partition π of the variables in $\vec{\mathbf{z}}_j$. Corresponding to such a partition, we have considered in the construction of $\Phi_{\text{conj}'_j}$ the formula $\exists \vec{\mathbf{w}}_\pi. \text{conj}'_j(\vec{\mathbf{a}}, \vec{\mathbf{w}}_\pi, \vec{\mathbf{c}}'_j)$, obtained by replacing all variables in the same equivalence class by a representative. Observe that, as a result, distinct variables in $\vec{\mathbf{w}}_\pi$ are assigned distinct domain elements in $\vec{\mathbf{d}}$.

Observe that, since \mathcal{M} is a pseudotree-admissible model, the only cycles that can be present in \mathcal{M} are formed by the domain elements that interpret the name-formulae corresponding to the (Skolem) constants in $\vec{\mathbf{a}}, \vec{\mathbf{b}}_1, \dots, \vec{\mathbf{b}}_m, \vec{\mathbf{c}}_1, \dots, \vec{\mathbf{c}}_m$. Let G be the tuple-graph obtained from $\exists \vec{\mathbf{w}}_\pi. \text{conj}'_j(\vec{\mathbf{a}}, \vec{\mathbf{w}}_\pi, \vec{\mathbf{c}}'_j)$. Since \mathcal{I} keeps the pseudotree structure of \mathcal{M} , if G contains cycles, the assignment α must assign each variable w in a cycle to an element d that interprets one of the (Skolem) constants t in $\vec{\mathbf{a}}, \vec{\mathbf{b}}_1, \dots, \vec{\mathbf{b}}_m, \vec{\mathbf{c}}_1, \dots, \vec{\mathbf{c}}_m$. It follows, by definition of \mathcal{I} , that $d \in N_t^{\mathcal{M}}$.

Now, let δ be the formula-template corresponding to G . From α , we define a G -substitution θ for δ as the G -substitution that replaces (cf. item (iv) in the definition of G -substitution) each placeholder w corresponding to a variable occurring in a cycle in the tuple-graph by the name-formulae N_t , where w is assigned by α to d and $d \in N_t^{\mathcal{M}}$. By considering the construction of the formula-template δ from G and the fact that $\exists \vec{\mathbf{w}}_\pi. \text{conj}'_j(\vec{\mathbf{a}}, \vec{\mathbf{w}}_\pi, \vec{\mathbf{c}}'_j)$ is true in \mathcal{I} with the assignment α , it can be shown that $s_{\text{root}} \in (\delta\theta)^{\mathcal{M}}$. But then, since $\delta\theta$ is one of the disjuncts forming $\Phi_{\text{conj}'_j}$, we have that $s_{\text{root}} \in (\Phi_{\text{conj}'_j})^{\mathcal{M}}$. Hence we get a contradiction. \square

The following theorem, which is a consequence of Theorems 3.2, 3.3 and 3.4, summarizes the correctness of our reduction.

THEOREM 3.5. *Let S be a schema, q, q' two queries of the same arity, and $\Phi_{S \not\subseteq q \subseteq q'}$ the formula obtained as specified above. Then $S \models q \subseteq q'$ if and only if $\Phi_{S \not\subseteq q \subseteq q'}$ is unsatisfiable.*

3.4 Complexity of Query Containment

Theorem 3.5 directly provides an algorithm for checking query containment in our setting, and therefore shows that the problem is decidable. We analyze now the computational complexity of such an algorithm for query containment.⁷

⁷In the results presented here, we assume unary coding of numbers in number restrictions.

THEOREM 3.6. *Let S be a schema and q and q' two queries. Then, deciding whether $S \models q \subseteq q'$ can be done in time $2^{p(|S|+|q|+|q'| \cdot \ell_1^{\ell_2})}$, where $|S|$, $|q|$, and $|q'|$ are, respectively, the sizes of S , q , and q' , ℓ_1 is the sum of the number of variables and constants in q , and ℓ_2 is the number of existentially quantified variables in q' .*

PROOF. The correctness of the encoding of query containment $S \models q \subseteq q'$ into unsatisfiability of $\Phi_{S \not\models q \subseteq q'}$ is sanctioned in Theorem 3.5. With regard to complexity, since satisfiability in CPDL_g is EXPTIME -complete, it follows that query containment can be done in time $2^{p(|\Phi_{S \not\models q \subseteq q'}|)}$. It is easy to verify that $|\Phi_{S \not\models q \subseteq q'}| = O(|S| + |q| + |q'| \cdot \ell_1^{O(\ell_2)})$. \square

The previous theorem provides, for the problem of checking whether $S \models q \subseteq q'$, a single exponential upper bound in the size of S and of q , and a double exponential upper bound in the size of q' (note that $|q'|$ is an upper bound for ℓ_2). The single exponential upper bound in the size of S and of q is tight. Indeed, it follows from EXPTIME -hardness of satisfiability in CPDL_g (in fact plain PDL [Fischer and Ladner 1979]) and from the fact that any CPDL_g formula can be expressed as a \mathcal{DLR}_{reg} concept. EXPTIME -hardness in S holds even in the case where S does not contain regular expressions. Indeed, the formulae used in the EXPTIME -hardness proof of satisfiability in PDL [Fischer and Ladner 1979], can be expressed as assertions in \mathcal{DLR}_{reg} not involving regular expressions. It is still open whether the double-exponential upper bound in the size of q' is tight.

The double exponential upper bound in the size of q' is due to the exponential blowup in the size of $\Phi_{S \not\models q \subseteq q'}$. By analyzing the reduction presented in Section 3.2, one can observe that such an exponential blowup is only due to those existentially quantified variables in q' that appear inside a cycle in the tuple-graph for q' . Hence, when the tuple-graph for q' does not contain cycles, we have that $|\Phi_{S \not\models q \subseteq q'}| = O(|S| + |q| + |q'|)$, and query containment can be checked in time $2^{p(|S|+|q|+|q'|)}$. A relevant case when this occurs is when (the tuple-graph for) the query on the right-hand side has the structure of a tree.

COROLLARY 3.7. *Let S be a schema, q and q' two queries of the same arity, and let q' have the structure of a tree. Then deciding whether $S \models q \subseteq q'$ can be done in time $2^{p(|S|+|q|+|q'|)}$.*

Observe that this gives us an EXPTIME -completeness result for containment of an arbitrary query in a tree-structured one wrt a schema.

Query satisfiability can be considered as a special case of query containment. Indeed, given a schema S , a query q is satisfiable wrt S if and only if it is not contained in the empty query wrt S . The empty query can be expressed, for example, as $u(\vec{x}) \leftarrow \mathbf{P}(\vec{x}) \wedge \neg \mathbf{P}(\vec{x})$, where \vec{x} is a tuple of variables and \mathbf{P} is a new atomic relation, both of the same arity as q .

COROLLARY 3.8. *Let S be a schema, and q a query. Then deciding whether q is satisfiable wrt S can be done in time $2^{p(|S|+|q|)}$.*

Again, this result shows EXPTIME-completeness of query satisfiability wrt a schema.

4. UNDECIDABILITY OF CONTAINMENT OF QUERIES WITH INEQUALITIES

In this section we show that, if we allow for inequalities inside the queries, then query containment wrt a schema becomes undecidable. The proof of undecidability exploits a reduction from the unbounded *tiling problem* [Berger 1966]. An instance $\mathcal{T} = (\mathcal{D}, H, V)$ of the tiling problem is defined by a finite set \mathcal{D} of tile types, a horizontal adjacency relation $H \in \mathcal{D} \times \mathcal{D}$, and a vertical adjacency relation $V \in \mathcal{D} \times \mathcal{D}$, and consists in determining whether there exists a tiling of the first quadrant of the integer plane with tiles of type in \mathcal{D} such that the adjacency conditions are satisfied. As shown, for example, in Harel [1985] and van Emde Boas [1997], the tiling problem is well suited to show undecidability of variants of modal and dynamic logics, and the difficult part of the proof usually consists in enforcing that the tiles lie on an integer grid. To this end we exploit a query containing one inequality.

Formally, given an instance $\mathcal{T} = (\mathcal{D}, H, V)$ of the tiling problem, a \mathcal{T} -tiling is a total function $t : \mathbb{N} \times \mathbb{N} \rightarrow \mathcal{D}$, and such a tiling is *correct* if $(t(i, j), t(i + 1, j)) \in H$ and $(t(i, j), t(i, j + 1)) \in V$, for each $i, j \in \mathbb{N}$. We reduce the problem of checking whether there exists a correct \mathcal{T} -tiling to the problem of checking whether $\mathcal{S}_{\mathcal{T}} \models q_0 \subseteq q'_0$, for suitable schema $\mathcal{S}_{\mathcal{T}}$ and queries q_0 and q'_0 containing inequalities.

Consider an instance $\mathcal{T} = (\mathcal{D}, H, V)$ of the tiling problem with tile types $\mathcal{D} = \{D_1, \dots, D_k\}$. We construct a schema $\mathcal{S}_{\mathcal{T}}$ using the atomic concepts *Tile*, D_1, \dots, D_k and two binary atomic relations *Right* and *Up* as follows:

$$\text{Tile} \sqsubseteq D_1 \sqcup \dots \sqcup D_k, \quad (9)$$

$$D_i \sqsubseteq \text{Tile}, \quad \text{for each } i \in \{1, \dots, k\}, \quad (10)$$

$$D_i \sqsubseteq \neg D_j, \quad \text{for each } i, j \in \{1, \dots, k\}, i < j, \quad (11)$$

$$\text{Tile} \sqsubseteq (\leq 1 [\text{\$1}]\text{Right}) \sqcap (\leq 1 [\text{\$1}]\text{Up}), \quad (12)$$

$$\text{Tile} \sqsubseteq \exists[\text{\$1}](\text{Right} \sqcap (\text{\$2}:\text{Tile})) \sqcap \exists[\text{\$1}](\text{Up} \sqcap (\text{\$2}:\text{Tile})), \quad (13)$$

$$D_i \sqsubseteq (\bigsqcup_{(D_i, D_j) \in H} \neg \exists[\text{\$1}](\text{Right} \sqcap (\text{\$2}:\neg D_j))) \sqcap (\bigsqcup_{(D_i, D_j) \in V} \neg \exists[\text{\$1}](\text{Up} \sqcap (\text{\$2}:\neg D_j))), \quad \text{for each } i \in \{1, \dots, k\}. \quad (14)$$

Then define the Boolean queries q_0 and q'_0 as follows:

$$q_0() \leftarrow \text{Tile}(x),$$

$$q'_0() \leftarrow \text{Right}(x, y) \wedge \text{Up}(y, z) \wedge \text{Up}(x, y') \wedge \text{Right}(y', z') \wedge z \neq z'.$$

THEOREM 4.1. *Let \mathcal{T} be an instance of the tiling problem, $\mathcal{S}_{\mathcal{T}}$ a schema, and q_0 and q'_0 two queries defined as specified above. Then there is a correct \mathcal{T} -tiling if and only if $\mathcal{S}_{\mathcal{T}} \models q_0 \subseteq q'_0$.*

PROOF. “ \Rightarrow ” Let t be a correct \mathcal{T} -tiling. We construct an interpretation \mathcal{I}_t of $\mathcal{S}_{\mathcal{T}}$ as follows:

$$\begin{aligned}\Delta^{\mathcal{I}_t} &= \mathbb{N} \times \mathbb{N}, \\ \text{Tile}^{\mathcal{I}_t} &= \Delta^{\mathcal{I}_t}, \\ \mathcal{D}_h^{\mathcal{I}_t} &= \{(i, j) \in \Delta^{\mathcal{I}_t} \mid t(i, j) = D_h\}, \quad \text{for each } h \in \{1, \dots, k\}, \\ \text{Right}^{\mathcal{I}_t} &= \{((i, j), (i + 1, j)) \mid i, j \in \mathbb{N}\}, \\ \text{Up}^{\mathcal{I}_t} &= \{((i, j), (i, j + 1)) \mid i, j \in \mathbb{N}\}.\end{aligned}$$

It is immediate to verify that \mathcal{I}_t is a model of $\mathcal{S}_{\mathcal{T}}$ and that $q_0^{\mathcal{I}_t}$ is true while $q'_0{}^{\mathcal{I}_t}$ is false.

“ \Leftarrow ” Consider a model \mathcal{I} of $\mathcal{S}_{\mathcal{T}}$ in which q_0 is true and q'_0 is false. Then \mathcal{I} contains an instance o_0 of *Tile* and assertions (13) in $\mathcal{S}_{\mathcal{T}}$ force the existence of arbitrary long chains of instances of *Tile*, beginning with o_0 and connected one to the next by alternations of *Right* ^{\mathcal{I}} and *Up* ^{\mathcal{I}} . By assertions (12), *Right* and *Up* are functional for all instances of *Tile*, and since q'_0 is false in \mathcal{I} , these chains of objects form indeed a grid. By assertions (9) and (11), each such object is an instance of precisely one D_h . Hence, we can construct a tiling $t_{\mathcal{I}}$ by assigning to each object o of the grid, representing an element of the first quadrant, a unique tile type D_h . Considering also assertions (14), it is easy to show by induction on the length of the chain from o_0 to an instance o of *Tile* that the horizontal and vertical adjacency conditions for o are satisfied. Hence $t_{\mathcal{I}}$ is a correct \mathcal{T} -tiling. \square

The theorem above immediately implies undecidability of containment with respect to a schema of queries containing inequalities.

THEOREM 4.2. *Let \mathcal{S} be a schema, and q, q' two queries of the same arity that may contain atoms of the form $t \neq t'$. Then the query containment problem $\mathcal{S} \models q \subseteq q'$ is undecidable.*

The reduction used in the proof of Theorem 4.1 shows that query containment remains undecidable even in the restricted case where

- all relations in \mathcal{S} , q , and q' are binary,
- \mathcal{S} does not contain assertions on relations, and all assertions on concepts are of the form $A \sqsubseteq C$,
- \mathcal{S} does not contain regular expressions,
- q and q' do not contain union, or constants expressions, and
- there is a single inequality in q' , and no inequality in q .

5. QUERY ANSWERING

As we said in the introduction, it is well known in the database literature that there is a tight connection between the problems of conjunctive query containment and conjunctive query answering [Chandra and Merlin 1977]. Such a relationship has had a particular importance in settings of databases with incomplete information, such as those arising in information integration [Abiteboul and Duschka 1998; Lenzerini 2002], semistructured data [Calvanese et al.

2002], and description logics [Baader et al. 2003]. In this section we discuss query answering under description logic constraints, taking advantage of the results on query containment presented above. By query answering under description logics constraints, we mean to compute the answers to a query over an *incomplete database*, that is, a database that is partially specified and must satisfy all description logic constraints expressed in a schema.⁸

Given a \mathcal{DLR}_{reg} schema \mathcal{S} , we specify an incomplete database \mathcal{D} over \mathcal{S} by means of a set of facts, called *membership assertions*, of the form

$$C(a), \quad \mathbf{R}(\vec{\mathbf{a}}),$$

where C and \mathbf{R} are, respectively, a concept expression and a relation expression over \mathcal{S} , a is a constant, and $\vec{\mathbf{a}}$ is a tuple of constants of the same arity as \mathbf{R} . Note that such a notion of incomplete database corresponds to that of ABox in description logics [Baader et al. 2003].

An interpretation \mathcal{I} satisfies an assertion $C(a)$ if $a^{\mathcal{I}} \in C^{\mathcal{I}}$, and it satisfies an assertion $\mathbf{R}(\vec{\mathbf{a}})$ if $\vec{\mathbf{a}}^{\mathcal{I}} \in \mathbf{R}^{\mathcal{I}}$. We say that \mathcal{I} is a *model* of \mathcal{D} , if it satisfies all assertions in \mathcal{D} . An incomplete database \mathcal{D} is *satisfiable with respect to a schema* \mathcal{S} if there is an interpretation \mathcal{I} that is a model of both \mathcal{S} and \mathcal{D} . Intuitively, every such interpretation \mathcal{I} represents a complete database that is coherent with both \mathcal{D} , and the description logic constraints in \mathcal{S} .

Given a schema \mathcal{S} , an incomplete database \mathcal{D} over \mathcal{S} , and a query q for \mathcal{S} , the set of *certain answers* $cert(q, \mathcal{S}, \mathcal{D})$ of q with respect to \mathcal{S} and \mathcal{D} is the set of tuples $\vec{\mathbf{c}}$ of constants in \mathcal{D} that are answers to q for all complete databases coherent with \mathcal{D} and \mathcal{S} , that is, such that $\vec{\mathbf{c}} \in q^{\mathcal{I}}$, for all models \mathcal{I} of \mathcal{S} and \mathcal{D} .

Given a query

$$q(\vec{\mathbf{x}}) \leftarrow conj_1(\vec{\mathbf{x}}, \vec{\mathbf{y}}_1, \vec{\mathbf{c}}_1) \vee \dots \vee conj_m(\vec{\mathbf{x}}, \vec{\mathbf{y}}_m, \vec{\mathbf{c}}_m)$$

in order to check whether a tuple $\vec{\mathbf{c}}$ of constants is in $cert(q, \mathcal{S}, \mathcal{D})$, we can resort to query containment [Abiteboul and Duschka 1998]. In particular, let us define the Boolean (i.e., of arity 0) queries $Q_{\mathcal{D}}$ and $Q_{q, \vec{\mathbf{c}}}$ as follows:

$$\begin{aligned} Q_{\mathcal{D}}() &\leftarrow \bigwedge_{C(a) \in \mathcal{D}} C(a) \wedge \bigwedge_{\mathbf{R}(\vec{\mathbf{a}}) \in \mathcal{D}} \mathbf{R}(\vec{\mathbf{a}}), \\ Q_{q, \vec{\mathbf{c}}}() &\leftarrow conj_1(\vec{\mathbf{c}}, \vec{\mathbf{y}}_1, \vec{\mathbf{c}}_1) \vee \dots \vee conj_m(\vec{\mathbf{c}}, \vec{\mathbf{y}}_m, \vec{\mathbf{c}}_m). \end{aligned}$$

The first query $Q_{\mathcal{D}}$ is the conjunction of all facts in \mathcal{D} , while the second query $Q_{q, \vec{\mathbf{c}}}$ is obtained from q by replacing each variable in $\vec{\mathbf{x}}$ with the corresponding constant in $\vec{\mathbf{c}}$.

THEOREM 5.1. *Let \mathcal{S} be a schema, \mathcal{D} an incomplete database over \mathcal{S} , q a query for \mathcal{S} , and $\vec{\mathbf{c}}$ a tuple of constants in \mathcal{D} of the same arity as q . Then $\vec{\mathbf{c}} \in cert(q, \mathcal{S}, \mathcal{D})$ if and only if $\mathcal{S} \models Q_{\mathcal{D}} \subseteq Q_{q, \vec{\mathbf{c}}}$.*

PROOF. The result can be proved exactly as in Abiteboul and Duschka [1998]. \square

From Theorem 3.6 we immediately obtain the following complexity result.

⁸Note that, in the case in which we have complete information on the database, the constraints do not play any role in query answering, assuming that the database is consistent with them.

THEOREM 5.2. *Let S be a schema, \mathcal{D} an incomplete database over S , q a query for S , and \vec{c} a tuple of constants in \mathcal{D} of the same arity as q . Then deciding whether $\vec{c} \in \text{cert}(q, S, \mathcal{D})$ can be done in time $2^{p(|S|+|\mathcal{D}|+|q|\cdot d^\ell)}$, where $|S|$, $|\mathcal{D}|$, and $|q|$ are, respectively, the sizes of S , \mathcal{D} , and q , d is the number of constants in \mathcal{D} and q , and ℓ is the number of existentially quantified variables in q .*

Note that this means that, while query answering is double exponential in combined complexity, it is actually only single exponential in the number of constants in the database. It follows that our technique is exponential in data complexity, that is, the complexity measured only with respect to the size of \mathcal{D} . Note that this is the first result on data complexity on answering unions of conjunctive queries under such expressive description logic constraints.

Finally, it follows directly from the semantics that satisfiability of a given incomplete database \mathcal{D} with respect to a schema S can be rephrased as satisfiability of the query $Q_{\mathcal{D}}$ with respect to S . Thus, we obtain the following result.

COROLLARY 5.3. *Let S be a schema and \mathcal{D} an incomplete database over S . Then deciding whether \mathcal{D} is satisfiable with respect to S can be done in time $2^{p(|S|+|\mathcal{D}|)}$.*

In description logic jargon, this shows EXPTIME-completeness of TBox+ABox satisfiability in our setting. Observe that, since we allow for union of conjunctive queries on the left-hand side query in the containment, this result can be immediately extended to satisfiability of a TBox together with a *disjunction* of ABoxes [Calvanese et al. 2001a].

6. CONCLUSIONS

In this article we have introduced \mathcal{DLR}_{reg} , an expressive language for specifying database schemas and nonrecursive Datalog queries, and we have presented decidability (with complexity) and undecidability results of both the problem of checking query containment, and the problem of answering queries under the constraints expressed in the schema.

The query language considered in this article allows no form of recursion, not even the transitive closure of binary relations. It is our aim in the future to extend our analysis to the case where queries may contain regular expressions, in the spirit of Calvanese et al. [2000].

ACKNOWLEDGMENTS

We thank the anonymous reviewers for many valuable comments, which helped improve the article.

REFERENCES

- ABITEBOUL, S. AND DUSCHKA, O. 1998. Complexity of answering queries using materialized views. In *Proceedings of the 17th ACM SIGACT SIGMOD SIGART Symposium on Principles of Database Systems (PODS'98)*. 254–265.

- ABITEBOUL, S., HULL, R., AND VIANU, V. 1995. *Foundations of Databases*. Addison-Wesley, Reading, MA.
- AHO, A. V., SAGIV, Y., AND ULLMAN, J. D. 1979a. Efficient optimization of a class of relational expressions. *ACM Trans. Database Syst.* 4, 297–314.
- AHO, A. V., SAGIV, Y., AND ULLMAN, J. D. 1979b. Equivalence among relational expressions. *SIAM J. Comput.* 8, 218–246.
- AMIR, K., PARK, S., TEWARI, R., AND PADMANABHAN, S. 2003. Scalable template-based query containment checking for Web semantic caches. In *Proceedings of the 19th IEEE International Conference on Data Engineering (ICDE 2003)*. 493–504.
- BAADER, F., CALVANESE, D., MCGUINNESS, D., NARDI, D., AND PATEL-SCHNEIDER, P. F., Eds. 2003. *The Description Logic Handbook: Theory, Implementation and Applications*. Cambridge University Press, Cambridge, U.K.
- BERGER, R. 1966. The undecidability of the domino problem. *Mem. Amer. Math. Soc.* 66, 1–72.
- BLACKBURN, P., DE RIJKE, M., AND VENEMA, Y. 2001. *Modal Logic*. Cambridge Tracts in Theoretical Computer Science, vol. 53. Cambridge University Press, Cambridge, U.K.
- BONATTI, P. A. 2004. On the decidability of containment of recursive datalog queries—preliminary report. In *Proceedings of the 23rd ACM SIGACT SIGMOD SIGART Symposium on Principles of Database Systems (PODS 2004)*. 297–306.
- BORGIDA, A. AND BRACHMAN, R. J. 2003. Conceptual modeling with description logics. See [Baader et al. 2003], Chapter 10, 349–372.
- CALVANESE, D. AND DE GIACOMO, G. 2003. Expressive description logics. See [Baader et al. 2003], Chapter 5, 178–218.
- CALVANESE, D., DE GIACOMO, G., AND LENZERINI, M. 1995. Structured objects: Modeling and reasoning. In *Proceedings of the 4th International Conference on Deductive and Object-Oriented Databases (DOOD'95)*. Lecture Notes in Computer Science, vol. 1013. Springer, Berlin, Germany, 229–246.
- CALVANESE, D., DE GIACOMO, G., AND LENZERINI, M. 1998a. On the decidability of query containment under constraints. In *Proceedings of the 17th ACM SIGACT SIGMOD SIGART Symposium on Principles of Database Systems (PODS'98)*. 149–158.
- CALVANESE, D., DE GIACOMO, G., AND LENZERINI, M. 2001a. Identification constraints and functional dependencies in description logics. In *Proceedings of the 17th International Joint Conference on Artificial Intelligence (IJCAI 2001)*. 155–160.
- CALVANESE, D., DE GIACOMO, G., LENZERINI, M., AND NARDI, D. 2001b. Reasoning in expressive description logics. In *Handbook of Automated Reasoning*, A. Robinson and A. Voronkov, Eds. Vol. II. Elsevier Science Publishers, Amsterdam, The Netherlands, Chapter 23, 1581–1634.
- CALVANESE, D., DE GIACOMO, G., LENZERINI, M., NARDI, D., AND ROSATI, R. 1998b. Description logic framework for information integration. In *Proceedings of the 6th International Conference on the Principles of Knowledge Representation and Reasoning (KR'98)*. 2–13.
- CALVANESE, D., DE GIACOMO, G., LENZERINI, M., AND VARDI, M. Y. 2000. Containment of conjunctive regular path queries with inverse. In *Proceedings of the 7th International Conference on the Principles of Knowledge Representation and Reasoning (KR 2000)*. 176–185.
- CALVANESE, D., DE GIACOMO, G., LENZERINI, M., AND VARDI, M. Y. 2002. View-based query answering and query containment over semistructured data. In *Revised Papers of the 8th International Workshop on Database Programming Languages (DBPL 2001)*, G. Ghelli and G. Grahne, Eds. Lecture Notes in Computer Science, vol. 2397. Springer, Berlin, Germany, 40–61.
- CALVANESE, D., DE GIACOMO, G., AND VARDI, M. Y. 2003. Decidable containment of recursive queries. In *Proceedings of the 9th International Conference on Database Theory (ICDT 2003)*. Lecture Notes in Computer Science, vol. 2572. Springer, Berlin, Germany, 330–345.
- CALVANESE, D., LENZERINI, M., AND NARDI, D. 1999. Unifying class-based representation formalisms. *J. Artif. Intell. Res.* 11, 199–240.
- CALVANESE, D. AND ROSATI, R. 2003. Answering recursive queries under keys and foreign keys is undecidable. In *Proceedings of the 10th International Workshop on Knowledge Representation meets Databases (KRDB 2003)*. CEUR Electronic Workshop Proceedings. Available online at <http://ceur-ws.org/Vol-79/>.
- CATARCI, T. AND LENZERINI, M. 1993. Representing and using interschema knowledge in cooperative information systems. *J. Intell. Coop. Informat. Syst.* 2, 4, 375–398.

- CHAN, E. P. F. 1992. Containment and minimization of positive conjunctive queries in OODB's. In *Proceedings of the 11th ACM SIGACT SIGMOD SIGART Symposium on Principles of Database Systems* (PODS'92). 202–211.
- CHANDRA, A. K. AND MERLIN, P. M. 1977. Optimal implementation of conjunctive queries in relational data bases. In *Proceedings of the 9th ACM Symposium on Theory of Computing* (STOC'77). 77–90.
- CHANDRA, A. K. AND VARDI, M. Y. 1985. The implication problem for functional and inclusion dependencies is undecidable. *SIAM J. Comput.* 14, 3, 671–677.
- CHAUDHURI, S. AND VARDI, M. Y. 1992. On the equivalence of recursive and nonrecursive Datalog programs. In *Proceedings of the 11th ACM SIGACT SIGMOD SIGART Symposium on Principles of Database Systems* (PODS'92). 55–66.
- CHEKURI, C. AND RAJARAMAN, A. 1997. Conjunctive query containment revisited. In *Proceedings of the 6th International Conference on Database Theory* (ICDT'97). 56–70.
- DE GIACOMO, G. AND LENZERINI, M. 1994. Boosting the correspondence between description logics and propositional dynamic logics. In *Proceedings of the 12th National Conference on Artificial Intelligence* (AAAI'94). 205–212.
- DE GIACOMO, G. AND LENZERINI, M. 1995. What's in an aggregate: Foundations for description logics with tuples and sets. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence* (IJCAI'95). 801–807.
- DE GIACOMO, G. AND LENZERINI, M. 1996. TBox and ABox reasoning in expressive description logics. In *Proceedings of the 5th International Conference on the Principles of Knowledge Representation and Reasoning* (KR'96). 316–327.
- DONG, G. AND SU, J. 1996. Conjunctive query containment with respect to views and constraints. *Informat. Process. Lett.* 57, 2, 95–102.
- DONINI, F. M., LENZERINI, M., NARDI, D., AND SCHAEFER, A. 1998. *AL-log*: Integrating Datalog and description logics. *J. Intell. Informat. Syst.* 10, 3, 227–252.
- ENDERTON, H. B. 1972. *A Mathematical Introduction to Logic*. Academic Press, New York, NY.
- FATTOROSI-BARNABA, M. AND DE CARO, F. 1985. Graded modalities I. *Studia Logica* 44, 197–221.
- FISCHER, M. J. AND LADNER, R. E. 1979. Propositional dynamic logic of regular programs. *J. Comput. Syst. Sci.* 18, 194–211.
- GRUBER, T. R. 1993. Towards principles for the design of ontologies used for knowledge sharing. In *Formal Ontology in Conceptual Analysis and Knowledge Representation*, N. Guarino and R. Poli, Eds. Kluwer Academic Publishers, Dordrecht, The Netherlands.
- GUPTA, A. AND MUMICK, I. S. 1995. Maintenance of materialized views: Problems, techniques, and applications. *Bull. IEEE Comput. Soc. Tech. Comm. Data Eng.* 18, 2, 3–18.
- GUPTA, A., SAGIV, Y., ULLMAN, J. D., AND WIDOM, J. 1994. Constraint checking with partial information. In *Proceedings of the 13th ACM SIGACT SIGMOD SIGART Symposium on Principles of Database Systems* (PODS'94).
- HALEVY, A. Y. 2001. Answering queries using views: A survey. *Very Large Database J.* 10, 4, 270–294.
- HAREL, D. 1985. Recurring dominoes: Making the highly undecidable highly understandable. *Ann. Discr. Math.* 24, 51–72.
- HORROCKS, I., SATTLER, U., TESSARIS, S., AND TOBIES, S. 2000. How to decide query containment under constraints using a description logic. In *Proceedings of the 7th International Conference on Logic for Programming and Automated Reasoning* (LPAR 2000). Lecture Notes in Computer Science, vol. 1955. Springer, Berlin, Germany, 326–343.
- HULL, R. 1997. Managing semantic heterogeneity in databases: A theoretical perspective. In *Proceedings of the 16th ACM SIGACT SIGMOD SIGART Symposium on Principles of Database Systems* (PODS'97). 51–61.
- IOANNIDIS, Y. E. AND RAMAKRISHNAN, R. 1995. Containment of conjunctive queries: Beyond relations as sets. *ACM Trans. Database Syst.* 20, 3, 288–324.
- JOHNSON, D. S. AND KLUG, A. C. 1984. Testing containment of conjunctive queries under functional and inclusion dependencies. *J. Comput. Syst. Sci.* 28, 1, 167–189.
- KLUG, A. C. 1988. On conjunctive queries containing inequalities. *J. ACM* 35, 1, 146–160.

- KOZEN, D. AND TIURYN, J. 1990. Logics of programs. In *Handbook of Theoretical Computer Science—Formal Models and Semantics*, J. van Leeuwen, Ed. Elsevier Science Publishers, Amsterdam, The Netherlands, 789–840.
- LENZERINI, M. 2002. Data integration: A theoretical perspective. In *Proceedings of the 21st ACM SIGACT SIGMOD SIGART Symposium on Principles of Database Systems (PODS 2002)*. 233–246.
- LEVY, A. Y. AND ROUSSET, M.-C. 1996. CARIN: A representation language combining Horn rules and description logics. In *Proceedings of the 12th European Conference on Artificial Intelligence (ECAI'96)*. 323–327.
- LEVY, A. Y., SRIVASTAVA, D., AND KIRK, T. 1995. Data model and query evaluation in global information systems. *J. Intell. Informat. Syst.* 5, 121–143.
- LEVY, A. Y. AND SUCIU, D. 1997. Deciding containment for queries with complex objects. In *Proceedings of the 16th ACM SIGACT SIGMOD SIGART Symposium on Principles of Database Systems (PODS'97)*. 20–31.
- MITCHELL, J. C. 1983. The implication problem for functional and inclusion dependencies. *Informat. Contr.* 56, 154–173.
- NEVEN, F. AND SCHWENTICK, T. 2003. XPath containment in the presence of disjunction, DTDs, and variables. In *Proceedings of the 9th International Conference on Database Theory (ICDT 2003)*. 315–329.
- PATEL-SCHNEIDER, P., HAYES, P., AND HORROCKS, I. 2004. OWL Web Ontology Language semantics and abstract syntax. W3C Recommendation. Available online at <http://www.w3.org/TR/owl-semantics/>.
- SAGIV, Y. AND YANNAKAKIS, M. 1980. Equivalences among relational expressions with the union and difference operators. *J. ACM* 27, 4, 633–655.
- SCHILD, K. 1991. A correspondence theory for terminological logics: Preliminary report. In *Proceedings of the 12th International Joint Conference on Artificial Intelligence (IJCAI'91)*. 466–471.
- ULLMAN, J. D. 1997. Information integration using logical views. In *Proceedings of the 6th International Conference on Database Theory (ICDT'97)*. Lecture Notes in Computer Science, vol. 1186. Springer, Berlin, Germany, 19–40.
- VAN DER MEYDEN, R. 1998. Logical approaches to incomplete information. In *Logics for Databases and Information Systems*, J. Chomicki and G. Saake, Eds. Kluwer Academic Publishers, Dordrecht, The Netherlands, 307–356.
- VAN EMDE BOAS, P. 1997. The convenience of tilings. In *Complexity, Logic, and Recursion Theory*, A. Sorbi, Ed. Lecture Notes in Pure and Applied Mathematics, vol. 187. Marcel Dekker Inc., New York, NY, 331–363.
- WIDOM, J. ED. 1995. Special issue on materialized views and data warehousing. *Bull. IEEE Comput. Soc. Tech. Comm. Data Eng.* 18, 2.
- WOOD, P. T. 2003. Containment for XPath fragments under DTD constraints. In *Proceedings of the 9th International Conference on Database Theory (ICDT 2003)*. 300–314.

Received July 2005; Revised April 2006; accepted December 2006