

Optimization Properties for Classes of Conjunctive Regular Path Queries

Alin Deutsch and Val Tannen

University of Pennsylvania,
200 South 33rd Str. 19104 Philadelphia, PA, USA
{adeutsch, val}@saul.cis.upenn.edu

Abstract. We are interested in the theoretical foundations of the optimization of conjunctive regular path queries (CRPQs). The basic problem here is deciding query containment both in the absence and presence of constraints. Containment without constraints for CRPQs is EXPSPACE-complete, as opposed to only NP-complete for relational conjunctive queries. Our past experience with implementing similar algorithms suggests that staying in PSPACE might still be useful. Therefore we investigate the complexity of containment for a hierarchy of fragments of the CRPQ language. The classifying principle of the fragments is the expressivity of the regular path expressions allowed in the query atoms. For most of these fragments, we give matching lower and upper bounds for containment in the absence of constraints. We also introduce for every fragment a naturally corresponding class of constraints in whose presence we show both decidability and undecidability results for containment in various fragments. Finally, we apply our results to give a complete algorithm for rewriting with views in the presence of constraints for a fragment that contains Kleene-star and disjunction.

1 Introduction

Semistructured data models and query languages [1] have become a very active area of interesting research in databases. In this paper we are interested in semistructured query languages, more precisely in theoretical foundations of query optimization for such languages. We concentrate on two computational problems:

- The problem of query equivalence (more generally, query containment), with or without integrity constraints.
- The problem of rewriting queries to make (some) use of views, again with or without integrity constraints.

For queries on relational, complex values, dictionary and OO data, these problems can be solved nicely and uniformly with a strengthening of the classical ideas on tableaux and chase. (See the chase extension in [23] and the chase & backchase technique for rewriting with views in [9].) Although the problems have theoretically intractable lower bounds, these bounds are in terms of query

and constraint size. It turns out that these techniques are in fact practical for practical-size queries and constraints [22]. Our experience with implementing them suggests that a necessary condition for practicality is the ability to decide containment in polynomial space. Can this be done for semistructured languages?

At the theoretical core of such languages lie the *conjunctive regular path queries* (CRPQs) of [13,6]. Here is an example:

$$Q(Y, Z) \leftarrow \mathbf{start} \ (a^*|b).c \ X, \ X \ a.b^* \ Y, \ X \ c^* \ Z$$

This is interpreted in a graph whose edge labels are taken from a set containing a, b, c while **start** is a constant node. The query returns the set of pairs (Y, Z) of nodes such that for some node X there are paths $\mathbf{start} \rightarrow X$, $X \rightarrow Y$, $X \rightarrow Z$ whose labels belong to the regular languages $(a^*|b).c$, $a.b^*$, c^* respectively.

However, containment of general CRPQs is EXPSPACE-complete [6,13]! Therefore, in this paper we pay attention to restricted fragments of CRPQs. This is an approach validated by practice: typical users exploit only a fraction of the expressive power of regular expressions. This is based on the experiences of users of the semistructured query language StruQL [12], but also of the XML query language XML-QL [8], and it is supported by the restrictions on path expressions imposed by the XPath standard [27]. Here is a very simple example of query optimization in such a fragment. Consider the query

$$A(N) \leftarrow \mathbf{start} * X, \ X \ \mathit{name.John} \ Y, \ X * \mathit{.tells.name} \ N$$

which returns the names of persons who find out a secret from somebody connected directly or indirectly to John. Assume that the following view is materialized

$$\mathit{Divulge}(V, W) \leftarrow \mathbf{start} * U, \ U \ \mathit{name} \ V, \ U \ \mathit{tells} \ W$$

and the following integrity constraint holds

$$(\mathit{tellAll}) \ \forall X, Y \ [\ X * \mathit{.tells} * Y \rightarrow X \ \mathit{tells} \ Y]$$

saying that our database models a society in which whenever two of its members share a secret, eventually everybody connected to them shares that secret. *Under this constraint*, the query A can be equivalently rewritten to use $\mathit{Divulge}$:

$$A'(N) \leftarrow \mathit{Divulge}(X, Y), \ X \ \mathit{John} \ Z, \ Y \ \mathit{name} \ N.$$

Depending on the storage schema, A' may be cheaper to evaluate. In the extended version [11] we show a detailed example of how the methods we develop in this paper succeed in finding this rewriting.

To study various fragments of CRPQs we develop a novel technique. [6,13] use automata-theoretic techniques but here we will try something different: reductions to problems formulated in the relational setting. The fragments for which we prove upper bounds and decidability results are such that we can translate queries and dependencies into relational versions, over a special relational schema. For example, the query Q shown above translates to the following union of relational conjunctive queries:

$$Q'(y, z) = C_1(y, z) \cup C_2(y, z)$$

$$C_1(y, z) \leftarrow a^*(\mathbf{start}, w_1), \ c(w_1, X), \ a(x, w_2), \ b^*(w_2, y), \ c^*(x, z)$$

$$C_2(y, z) \leftarrow b(\mathbf{start}, w_1), \ c(w_1, X), \ a(X, w_2), \ b^*(w_2, y), \ c^*(x, Z)$$

We think of C_1, C_2 as ordinary relational conjunctive queries over a schema containing a, a^*, b^*, c, c^* . A priori a and a^* etc., are independent binary relation symbols, but we interpret them only in relational instances in which certain relational constraints hold. The constraints are first-order and they say, for example, that a^* is transitive, reflexive, and includes a . Of course, transitive closure itself cannot be expressed in first-order logic. It is therefore remarkable that first-order reasoning suffices for some of the semistructured language fragments we consider in this paper. However, we also provide undecidability results that together with the aforementioned EXPSPACE lower bound [6] show some of the theoretical limits of what can be done about optimization in semistructured languages.

Organization of the Remainder of This Paper. In section 2 we define the classes (language fragments) of queries and dependencies under study here, as well as their translation into relational correspondents. In section 3 we summarize our results and discuss some related work. Section 4 contains our results on upper bounds for pure query containment while section 5 contains the corresponding lower bound results. Section 6 presents our results on deciding (or not!) containment of queries in the presence of dependencies. Section 7 extends the chase & backchase technique [9] to two of the fragments we study. We conclude in section 8.

Due to space limitations, we have relegated some proofs and a worked example to the full paper [11], which also contains the extension to unions and disjunction of the chase. Although this extension, in the form that we need, has not—apparently—been published previously, it will not surprise anyone with an understanding of the classical chase.

2 Queries and Constraints

Databases. Let \mathcal{L} be a set of *labels*. For technical reasons we assume that \mathcal{L} is infinite, but of course only a finite number of labels will occur in a given database, query, or constraint. A *semistructured database* is a finite directed graph whose edges are \mathcal{L} -labeled. Equivalently, we can be given a set N (the nodes of the graph) and a finite set of labels from \mathcal{L} , each interpreted as a non-empty binary relation on N .

A word about **constants** denoting nodes. The upper bound and decidability results do *not*, as stated, assume the presence of such constants. Equalities between distinct constants cause the usual problem [2] and our results can be extended straightforwardly to deal with this. For clarity of exposition we have omitted this extension. On the other hand, some of our examples and even some of the constructions used in lower bounds and undecidability results do use constants denoting nodes. Such use is in fact inessential and is made for the same reasons of clarity.

Queries: CRPQs. A *conjunctive regular path query (CRPQ)* [13,6] has the general form

$$Q(x_1, \dots, x_n) \leftarrow A_1, \dots, A_m \quad (1)$$

Here the *atoms* (*conjuncts*) A_i are either equalities $y = z$ or *regular path atoms* of the form $y R z$ where R is a regular expressions defined by ¹

$$R ::= l \mid _ \mid R^* \mid R_1.R_2 \mid (R_1 \mid R_2) \quad (2)$$

$$* \stackrel{\text{def}}{=} _ *$$

where l ranges over labels in \mathcal{L} and $_$ means *any* (single) label. Of course, each *distinguished* variables x_j must also occur in the right hand side. As indicated, we follow [13] in using the *shorthand* $*$ for $_*$.

If B is a semistructured database, an atom $x R y$ is *satisfied* by a valuation that maps x, y to nodes s, t in B if there is a path from s to t in B which spells out a word in the language denoted by the regular expression R . We extend this definition of atom satisfaction to give semantics to whole CRPQs in the way that is usual for conjunctive queries. Query containment is also defined as usual.

Unions of CRPQs. In spite of being called “conjunctive”, CRPQs contain implicit forms of disjunction, most glaringly because of the \mid operator in regular expressions. In fact, we are naturally led to consider *unions of CRPQs* as the class of queries of interest. It is easy to see that the EXPSPACE upper bound on containment [13,6] still holds for unions of CRPQs.

Containment and Dependencies. Much of the early relational database theory dealt with conjunctive (tableau) queries and *embedded dependencies* [2] which are logical assertions of the special form

$$\forall \mathbf{x} [C_1(\mathbf{x}) \rightarrow \exists \mathbf{y} C_2(\mathbf{x}, \mathbf{y})] \quad (3)$$

where C_1, C_2 are conjunctions of relational atoms or (in C_2) equalities ². Such dependencies are tightly related to containment assertions [28]. Given two (type-compatible) conjunctive queries Q_1, Q_2 it is easy to construct an embedded dependency that is equivalent (in each database instance) to the containment $Q_1 \subseteq Q_2$. It is equally easy to construct an equivalent containment assertion from any given embedded dependency.

In this paper we will consider several classes of queries, and for each of them we will identify a class of dependencies (constraints) that has this kind of tight correspondence with the containment of queries from the associated class.

Add Disjunction: DEDs. Generalizing from conjunctive queries to *unions of conjunctive queries*, we consider the associated class of *disjunctive embedded dependencies* (DEDs) which are logical assertions of the form

$$\forall \mathbf{x} [C_1(\mathbf{x}) \rightarrow \bigvee_{i=1}^m \exists \mathbf{y}_i C_{2,i}(\mathbf{x}, \mathbf{y}_i)] \quad (4)$$

where $C_1, C_{2,i}$ are as in (3). We don’t need disjunction in the premise of the implication because it is equivalent to conjunctions of DEDs. We have the following

¹ We ask the reader to distinguish between the \mid in regular expressions and the *meta* use of \mid as part of the BNF for the syntax.

² The notation \mathbf{x} abbreviates x_1, \dots, x_n .

tight correspondence: the containment of two unions of conjunctive queries is equivalent to a finite number of DEDs, and a single DED is equivalent to the containment of a conjunctive query into a union of conjunctive queries.

A DED is *full* if it does not have existentially quantified variables. The *chase* [3] can be extended to DEDs, giving a decision procedure for containment of unions of conjunctive queries under a set of full DEDs (see [15] for a partial treatment and the extended version of this paper [11] for a sketch of the results we use.)

Semistructured Constraints: DERPDS. As with DEDs, we define the class of dependencies that corresponds to unions of conjunctive regular path queries (CRPQs). We call such dependencies *disjunctive embedded regular path dependencies* (*DERPDs*) and they are defined as assertions that have the same logical form as DEDs, see (4), but in which $C_1, C_{2,i}$ are conjunctions of regular path atoms $x R y$ or equalities. The definition for satisfaction of a given DERPDS in a given semistructured database follows from the usual meaning of logical connectives and quantifiers and from the satisfaction for regular path atoms given earlier.

When the regular expressions are restricted to single labels in \mathcal{L} , CRPQs are equivalent to the usual conjunctive queries and DERPDS to just DEDs seen over a relational schema consisting of binary symbols from \mathcal{L} .

Examples. DERPDS can express a large variety of constraints on semistructured data. As we saw, they generalize most relational dependencies of interest. In addition we can express constraints similar to the ones DTDs [26] specify for XML. The first two below say that “any person has exactly one social security number”. The third says that “telephone numbers can only be of two (if any) kinds, voice or fax” while the fourth is a kind of generalized join-like dependency.

$$\begin{aligned} &\forall x [\text{start} * . \text{person } x \rightarrow \exists y x \text{ ssn } y] \\ &\forall x \forall y_1 \forall y_2 [\text{start} * . \text{person } x \wedge x \text{ ssn } y_1 \wedge x \text{ ssn } y_2 \rightarrow y_1 = y_2] \\ &\forall x \forall y \forall z [x \text{ telNo } y \wedge y - z \rightarrow y \text{ voice } z \vee y \text{ fax } z] \\ &\forall x \forall y \forall z [x \text{ child } y \wedge y \text{ child } z \rightarrow z \text{ grandparent } x] \end{aligned}$$

Fragments: F -Queries and F -Dependencies. Since containment of CRPQs is EXPSpace-complete [6] we study *fragments* of the language defined by restricting the regular expressions allowed in atoms (conjuncts). The simplest fragment, allowing just labels and concatenation, is equivalent to conjunctive queries over binary relations. Between these and general CRPQs we consider the fragments described by the table below. For any fragment F , we call the corresponding queries *F -queries*. Applying the same restriction to the atoms that appear in dependencies, we define corresponding classes of DERPDS, calling the respective constraints *F -dependencies*. The correspondence discussed above, between containment assertions and dependencies, continues to hold for each fragment F . The fragments called W and Z have technical importance but their definitions did not suggest anything better than choosing these arbitrary names.

Fragment name	Regular expressions syntax
conj. queries	$R \rightarrow l \mid R_1.R_2$
$(*)$	$R \rightarrow l \mid * \mid R_1.R_2$
$(*,)$	$R \rightarrow l \mid * \mid R_1.R_2 \mid (R_1 R_2)$
$(*, -)$	$R \rightarrow l \mid - \mid * \mid R_1.R_2$
(l^*)	$R \rightarrow l \mid l^* \mid R_1.R_2$
$(*, -, l^*,)$	$R \rightarrow l \mid * \mid - \mid l^* \mid R_1.R_2 \mid (R_1 R_2)$
W	$R \rightarrow l \mid - \mid S^* \mid R_1.R_2 \mid (R_1 R_2)$ $S \rightarrow l \mid - \mid S_1.S_2$
Z	$R \rightarrow S \mid S^*$ $S \rightarrow l \mid S_1.S_2 \mid (S_1 S_2)$
CRPQs	$R \rightarrow l \mid - \mid R^* \mid R_1.R_2 \mid (R_1 R_2)$

Containment of	Upper bound	Containment of	Lower bound
conjunctive queries	NP [7]	conjunctive queries	NP [7]
$(*)$ -queries	NP [13] or corollary 1	$(*)$ -queries	\downarrow
unions of conj. queries	Π_2^P [24]	unions of conj. queries	Π_2^P [24]
unions of $(*,)$ -queries	\uparrow	$(*,)$ -queries	Π_2^P remark 2
unions of $(*, -)$ -queries		$(*, -)$ -queries	Π_2^P theorem 3
unions of (l^*) -queries		(l^*) -queries	Π_2^P theorem 3
unions of $(*, -, l^*,)$ -queries		$(*, -, l^*,)$ -queries	\downarrow
unions of W queries	?	W queries	PSPACE theorem 4
unions of Z queries	\uparrow	Z queries	EXPSpace
unions of CRPQs	EXPSpace [13,6]	CRPQs	[6] and remark 3 \downarrow

Fig. 1. Upper and lower bounds for containment

First-Order Relational Translation. At the core of our technique is a translation of semistructured queries and dependencies into first-order logic, namely into (unions of) conjunctive queries and DEDs *over a special relational schema* that includes l and l^* as well as $-$ and $*$ as separate binary relation symbols. A priori these symbols are independent, but we will try to capture *some* of the Kleene star semantics through relational dependencies.

Our translation is designed for the $(*, _, l^*, |)$ -fragment only. It relies essentially on the fact that in this fragment concatenation and $|$ are *not* nested inside Kleene stars.

The first thing we do is **translate away** $|$. Using the equivalence $(a|b).c = (a.c)|(b.c)$ we move $|$ in the outermost position in the $(*, _, l^*, |)$ -regular expressions. Then, we note that $Q \leftarrow \dots, x R_1 | R_2 y, \dots$ is equivalent to $Q_1 \cup Q_2$ where $Q_i \leftarrow \dots, x R_i y, \dots$. For dependencies, we note that $x R_1 | R_2 y$ is equivalent to $x R_1 y \vee x R_2 y$ after which logical equivalences bring the disjunctions out. A disjunction in the premise of the implication in a dependency is equivalent to a conjunction (a set) of dependencies. To summarize:

Remark 1. By translating away the $|$, any $(*, _, l^*, |)$ -query becomes an equivalent union of $(*, _, l^*)$ -queries. Similarly, any $(*, _, l^*, |)$ -dependency becomes an equivalent set of $(*, _, l^*)$ -dependencies.

Next, we translate any $(*, _, l^*)$ -queries and dependencies into (relational) conjunctive queries and DEDs over the special schema

$$\mathcal{L}\text{-Rel} \stackrel{\text{def}}{=} \{l \mid l \in \mathcal{L}\} \cup \{l^* \mid l \in \mathcal{L}\} \cup \{-, *\} \cup \{N\}$$

in which all symbols are binary relations with the exception of N which is unary (the need for N is explained below).

The **translation** $\mathcal{T}(Q)$ of a $(*, _, l^*)$ -query Q is defined by translating its conjuncts according to the rules (for each binary r in $\mathcal{L}\text{-Rel}$)

$$\begin{aligned} \mathcal{T}(x r y) &= r(x, y) \\ \mathcal{T}(x R_1.R_2 y) &= \mathcal{T}(x R_1 v), \mathcal{T}(v R_2 y) \end{aligned}$$

The variable v is (implicitly) existentially quantified and so it must be fresh each time its rule is applied. For example, $Q(x, y) \leftarrow x a.*b y$, translates to $Q'(x, y) \leftarrow a(x, z), *(z, u), b(u, y)$.

The **translation** $\mathcal{T}(d)$ of a $(*, _, l^*)$ -dependency d is defined similarly. The presence of concatenation in the conclusion of the implication in d will add existentially quantified variables, while the presence of concatenation in the premise of the implication in d will add universally quantified variables.

Example of Translation. Let d be the dependency

$$\forall x \forall y [x (a|b).* y \rightarrow \exists z y *. (a|b) z]$$

It translates to the following set of two DEDs

$$\begin{aligned} \forall x \forall y \forall u [x a u \wedge u * y \rightarrow [\exists z_1 \exists v_1 y * v_1 \wedge v_1 a z_1] \vee [\exists z_2 \exists v_2 y * v_2 \wedge v_2 b z_2]] \\ \forall x \forall y \forall u [x b u \wedge u * y \rightarrow [\exists z_1 \exists v_1 y * v_1 \wedge v_1 a z_1] \vee [\exists z_2 \exists v_2 y * v_2 \wedge v_2 b z_2]] \end{aligned}$$

Now, $\mathcal{T}(Q)$ is a relational query and $\mathcal{T}(d)$ is a relational dependency, both over the schema $\mathcal{L}\text{-Rel}$. However, we will use them not over arbitrary instances of $\mathcal{L}\text{-Rel}$ but only over instances that satisfy specific sets of relational dependencies. To deal with the various fragments, we consider two such sets

The Σ_* Dependencies:

$$\begin{array}{ll}
(\text{node}_l) & \forall x \forall y [l(x, y) \rightarrow N(x) \wedge N(y)] & (\text{node}_*) & \forall x \forall y [* (x, y) \rightarrow N(x) \wedge N(y)] \\
(\text{base}) & \forall x \forall y [l(x, y) \rightarrow *(x, y)] & (\text{refl}_*) & \forall x [N(x) \rightarrow *(x, x)] \\
(\text{trans}_*) & \forall x \forall y \forall z [* (x, y) \wedge *(y, z) \rightarrow *(x, z)]
\end{array}$$

where l ranges over \mathcal{L} . (This is an infinite set of dependencies but of course only finitely many matter for a given database, query, or dependency.) Here we see how we use N : we want the chase with (refl_*) to apply only to variables x that are already present.

The Σ_{l^*} Dependencies: are obtained by replacing $*$ with l^* in Σ_* above.

The intention behind these dependencies is to narrow the gap between the semistructured meaning of the Kleene star and the arbitrary interpretation that could be given to the relational schema $\mathcal{L}\text{-Rel}$. We can associate directly to each semistructured database a relational $\mathcal{L}\text{-Rel}$ -instance that satisfies $\Sigma_* \cup \Sigma_{l^*}$ (call it a $\Sigma_* \cup \Sigma_{l^*}$ -instance). But this will not cover $\Sigma_* \cup \Sigma_{l^*}$ -instances containing pairs of distinct nodes which are not connected by any path with labels from \mathcal{L} . Of course, it is not possible to close the gap this way, since transitive closure is not first-order definable. It is therefore remarkable that first-order reasoning suffices for some of the semistructured language fragments we consider in this paper.

Full Dependencies. Relational dependencies (3) and DEDs (4) are called *full* when they do not have existentially quantified variables. In the case of DERPDS fullness must be more complicated because concatenation in regular expressions introduces an implicit existential. Here we take a very simple approach.

Let d be an $(*, -, l^*, |)$ -dependency and let $\mathcal{T}(d)$ be the set of DEDs into which d translates. We say that d is a *full dependency* if each DED in $\mathcal{T}(d)$ is full.

3 Summary of Results

Containment for F -Queries. We summarize in figure 1 our new results on the complexity of deciding containment for queries in the various fragments, putting them in the context of known results.

The upper bounds are for containment of unions of F -queries, with the remarkable exception of the $(*)$ -fragment for which containment of $(*)$ -queries is in NP, just like containment of conjunctive queries. This was already shown in [13]. Motivated by the study of containment *under dependencies*, the new technique introduced here reproves, along the way, this NP bound, see corollary 1³.

Our new upper bound result is that containment of unions of $(*, -, l^*, |)$ -queries is in Π_2^P (theorem 2). It can be seen from the lower bounds table that this is a tight bound.

³ Although we do not consider it explicitly here, the fragment obtained by adding just $-$ to labels and concatenation is easily seen to yield no surprises: containment is still in NP.

We have tried to state our lower bound results in their strongest form, for F -queries rather than unions of F -queries. It is not surprising to see a Π_2^P lower bound in the presence of $|$. This does not follow directly from [24] but we have

Remark 2. The Π_2^P -hardness proof in [24] for the lower bound on containment of unions of conjunctive queries can be adapted to containment of F -queries provided that F includes $|$.

It is surprising however what happens in the absence of $|$. While containment of $(*)$ -queries is in NP, we show in theorem 3 that containment of (l^*) -queries is Π_2^P -hard. Moreover a simple variation of the same proof applies to the $(*, -)$ -fragment. Therefore, we find that the increase in complexity does not stem from the mere presence of the Kleene star in the query, but from the *interaction* between l and l^* or between $-$ and $-^*$.

A more liberal nesting of regular expressions within the Kleene star increases complexity. If we allow concatenation inside the Kleene star, we get the W-fragment, for which we show in theorem 4 a PSPACE lower bound on containment. We don't know (mainly because of difficulties with a relational translation) if this bound is tight, which is why we put a question mark in the corresponding upper bound entry. If in addition we allow disjunction within the Kleene star, we obtain the Z-fragment which is as bad as general CRPQs:

Remark 3. The EXPSPACE-hardness proof in [6] applies to containment of Z-queries.

Containment of	Under what constraints	Decidable?
conjunctive queries	full relational dependencies	YES ([3])
unions of conjunctive queries	full DEDs	YES ([11])
unions of $(*,)$ -queries	full $(*,)$ -deps.	YES (theorem 5)
unions of $(*, -)$ - queries	full V-deps.	YES (theorem 6)
$(*, -)$ -query in union of $(*, -)$ -queries	full DEDs over special models	NO (theorem 7)
(l^*) -query in union of (l^*) -queries	full DEDs	NO (theorem 8)

Fig. 2. Results for containment under dependencies

Containment under Dependencies. The chase technique in classical relational theory gives us the decidability of containment of conjunctive queries under full dependencies [2]. Decidability extends straightforwardly to containment of unions of conjunctive queries under full DEDs (see the full paper [11]).

This nice situation for relational languages contrasts with the situation for semistructured languages, as summarized in figure 2. The general problem studied is containment of unions of F -queries under full F -dependencies. It turns out that even containment of (l^*) -queries under just full DEDs is undecidable (theorem 8).

There is some good news, as our technique carries through in theorem 5 to prove decidability for the $(*, |)$ -fragment.

We leave open the general problem corresponding to the $(*, _)$, but we have two partial results that suggest that the problem might be complicated. We show decidability in the case of a restricted class of $(*, _)$ -dependencies, that we call *V-dependencies* (definition in section 6). And we show *undecidability* with just DEDs in the case of a class of *special models* (also defined in section 6).

In our two undecidability proofs, just like in the proof of theorem 3, we make essential use of the *interaction* between l and l^* or between $_$ and $_*$.

Rewriting with Views under Dependencies. Given a set V of views, a set D of dependencies expressing integrity constraints, and a query Q , we are interested in finding “rewritings” Q' which mention some of the views (but may still contain labels from Q) and are *exactly* equivalent to Q .

We do not study this problem in its full generality, but rather we look at extending to some of the F -fragments the *chase&backchase* (C&B) algorithm that we introduced in [9]. This algorithm relies on the chase with dependencies. In view of the undecidability results we have obtained for other F -fragments, we have looked at rewriting with views only for the $(*)$ - and $(*, |)$ -fragments.

In theorem 9 we show that (essentially) the C&B algorithm is complete for the $(*)$ -fragment, in the sense that it finds all rewritings that are *minimal* in a precise sense.

For the $(*, |)$ -fragment we extend the original C&B algorithm to account for disjunction, and we prove that this extended version is also complete.

Related Work. Perhaps the closest in spirit is [4], which gives an EXPTIME-complete decision procedure for containment of queries and constraints expressed in a different fragment of CRPQs, which corresponds to description logics. This fragment allows unrestricted regular expressions in the conjuncts, but restricts the shape of the query graph (thus being incompatible with our classification principle for query fragments). The corresponding dependencies allow unrestricted regular path expressions and even cardinality constraints, but have restricted shape and in particular cannot express functional dependencies. As a matter of fact, [14] shows that, when adding functional dependencies to a generalization of description logics called the Guarded Fragment of first order logic, satisfiability (and hence containment) becomes undecidable. None of our query fragments is contained in description logics.

The class of $(*)$ -queries was introduced in [13] (under the name of “simple StruQL₀ queries”) as a class of semistructured queries using transitive closure and whose containment problem is in NP. The decision procedure was based on an automata-theoretic argument which was applicable to CRPQs with arbitrary regular path expressions.

[19,20] study the expressivity and satisfiability of queries over tree structures, in formalisms that are equivalent to MSO. Classes of tree structures are given as grammars, which can be viewed as constraints on their structure in a broader sense.

[5] gives a complete algorithm for finding rewritings of regular path queries (i.e. single-conjunct CRPQs) with views defined by regular path queries. The path expressions allowed in the conjunct are unrestricted, but no constraints are taken into account, and only complete rewritings are obtained (that is, rewritings mentioning only views). [17] addresses the problem of finding arbitrary rewritings of regular path queries, and [16] gives an algorithm for the related problem of *answering* regular path queries using incomplete views.

4 Upper Bounds

(*)-Queries. Recall that a $(*)$ -query is a CRPQ whose atoms allow only regular expressions built from labels, $*$, and their concatenation. For example, $Q(x, y) \leftarrow x a. * . b y, y c x$ is a $(*)$ -query, as opposed to $Q'(y) \leftarrow x a.b^*.c y$ (because of b^*) and $Q''(y) \leftarrow x a|b y$ (because of $|$).

We have shown in section 2 how to translate any $(*)$ -query into a conjunctive query $\mathcal{T}(Q)$ over the schema $\mathcal{L}\text{-Rel}$. While not obvious, it turns out that reasoning about $\mathcal{T}(Q)$ under the set of dependencies Σ_* introduced in section 2 suffices (see [11] for proof):

Proposition 1. *Let Q_1, Q_2 be two $(*)$ -queries. The containment $Q_1 \subseteq Q_2$ is valid if and only if $\Sigma_* \models \mathcal{T}(Q_1) \subseteq \mathcal{T}(Q_2)$.*

Next, we observe that the dependencies in Σ_* are full hence the chase with them terminates, giving a decision procedure for $\Sigma_* \models \mathcal{T}(Q_1) \subseteq \mathcal{T}(Q_2)$ [3,2]. We denote with $\text{chase}_{\Sigma_*}(Q)$ the result of chasing the query Q with the dependencies in Σ_* .

Theorem 1. *The $(*)$ -query Q_1 is contained in the $(*)$ -query Q_2 if and only if there exists a containment mapping (see [2]) from $\mathcal{T}(Q_2)$ into $\text{chase}_{\Sigma_*}(\mathcal{T}(Q_1))$.*

Corollary 1. (see also [13]) *$(*)$ -query containment is NP-complete.*

Proof: First notice that the size of $\mathcal{T}(Q)$ is linear in that of Q . The time to chase is polynomial in the size of the queries, but exponential in the maximum size of a dependency and the maximum arity of the relations in the schema [3]. However, the dependencies in Σ_* have fixed size and the maximum arity of a relation in the schema is 2. The upper bound follows noting that the containment mapping can be found in NP. For the lower bound, note that the proof of NP-hardness for containment of conjunctive queries in [7] holds even if all relations are binary. •

Example. Consider $Q_1(x_1, x_3) \leftarrow x_1 a x_2, x_2 b.c x_3$ and $Q_2(y_1, y_2) \leftarrow y_1 * .a. * y_2$. It is easy to see that Q_1 is contained in Q_2 . We show how we infer this using theorem 1. The translation to conjunctive queries yields $TQ_1 = \mathcal{T}(Q_1)$ and $TQ_2 = \mathcal{T}(Q_2)$, with $TQ_1(x_1, x_3) \leftarrow a(x_1, x_2), b(x_2, u_1), c(u_1, x_3)$ and $TQ_2(y_1, y_2) \leftarrow *(y_1, v_1), a(v_1, v_2), *(v_2, y_2)$. Note that there is no containment mapping from TQ_2 to TQ_1 as the latter contains no $*$ -atoms to serve as image for the former's $*$ -atoms. But by chasing TQ_1 with (node_a) and then with (refl_*) , we

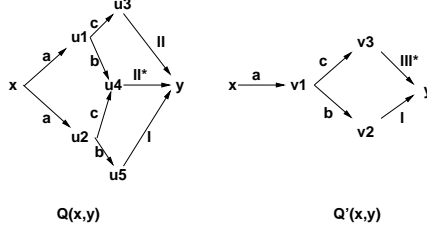


Fig. 3. Counterexample queries for proposition 2

obtain $Q'(x_1, x_3) \leftarrow N(x_1), N(x_2), *(x_1, x_1), a(x_1, x_2), b(x_2, u_1), c(u_1, x_3)$, thus creating an image for TQ_2 's conjunct $*(y_1, v_1)$. We continue chasing with $(base_*^b)$, then $(base_*^c)$ and $(trans_*)$, obtaining $Q''(x_1, x_3) \leftarrow N(x_1), N(x_2), *(x_1, x_1), a(x_1, x_2), b(x_2, u_1), c(u_1, x_3), *(x_2, u_1), *(u_1, x_3), *(x_2, x_3)$. Now $\{y_1 \mapsto x_1, y_2 \mapsto x_3, v_1 \mapsto x_1, v_2 \mapsto x_2\}$ is a containment mapping from $\mathcal{T}(Q_2)$ into Q'' . There are further applicable chase steps, omitted here as they only add new atoms and hence do not affect the existence of the containment mapping. •

Unions of $(*, \neg, l^*, |)$ -Queries. The idea we have just used to handle $(*)$ -queries is easily extended to $(*, |)$ -queries (giving a Π_2^P procedure), but how about other fragments? Can we deal with (l^*) -queries using their relational translation and the set Σ_{l^*} of dependencies defined in section 2? The answer is negative, which is surprising given the syntactic similarity of the $(*)$ - with the (l^*) -fragment.

Proposition 2. *There exist (l^*) -queries Q, Q' such that $Q \subseteq Q'$ but $\Sigma_{l^*} \not\models \mathcal{T}(Q) \subseteq \mathcal{T}(Q')$.*

Proof: Here are the queries (see figure 3 for possibly helpful graph representations of Q, Q'):

$Q(x, y) \leftarrow x a u_1, x a u_2, u_1 c u_3, u_1 b u_4, u_2 b u_5, u_2 c u_4, u_3 l l y, u_4 l l^* y, u_5 l y$
 $Q'(x, y) \leftarrow x a v_1, v_1 b v_2, v_1 c v_3, v_2 l y, v_3 l l l^* y$

To see that Q is contained in Q' , observe that $ll^* = l \cup lll^*$ and Q is equivalent to the union of queries $Q_1 \cup Q_2$ where Q_1, Q_2 are obtained by replacing the conjunct $u_4 ll^* y$ with $u_4 l y$, respectively $u_4 lll^* y$ in Q . But both Q_1, Q_2 are contained in Q' , as witnessed by the containment mappings $\{v_1 \mapsto u_1, v_2 \mapsto u_4, v_3 \mapsto u_3\}$ and $\{v_1 \mapsto u_2, v_2 \mapsto u_5, v_3 \mapsto u_4\}$. Intuitively, for any instance I , and any mapping from Q to I , depending on whether $u_4 ll^* y$ in Q is satisfied by a path of length 1 or at least 2, $v_1 c v_3$ in Q' is satisfied by the same path which satisfies either $u_1 c u_3$ or $u_2 c u_4$, respectively.

On the other hand, according to the chase theorem [2], $\mathcal{T}(Q)$ is not contained in $\mathcal{T}(Q')$ under Σ_{l^*} because there is no containment mapping from $\mathcal{T}(Q')$ into $chase_{\Sigma_{l^*}}(\mathcal{T}(Q))$. (Intuitively, what Σ_{l^*} does not capture is the *minimality* of l^* : it only states that l^* contains the reflexive transitive closure of l , but it doesn't rule out pairs of nodes that aren't reachable via a path of l -edges. Instances containing such a pair (s, t) are counterexamples for the containment: conjunct

$u_4 ll^* y$ in Q is satisfied by the endpoints of the path $r \xrightarrow{l} s \xrightarrow{l^*} t \xrightarrow{l} q$ even if s has no outgoing l -edge, while $v_3 lll^* y$ in Q' is not.) •

A simple variation of the counterexample above applies to $(*, _)$ -queries. In any case, if the same idea would have applied it would have given us NP algorithms, and we show in theorem 3 that containment for both the (l^*) - and the $(*, _)$ -fragment is Π_2^P -hard! Therefore, we will take another route towards a containment test.

We start from the observation that $\Sigma_* \cup \Sigma_{l^*}$ is sufficient in deciding containment of Q_1 in Q_2 in the restricted case in which Q_1 contains no Kleene star (no $*$ or l^*), and Q_2 is a $(*, _, l^*)$ -query. We call Q_1 *star-free*.

Proposition 3. *The star-free query Q_1 is contained in the $(*, _, l^*)$ -query Q_2 if and only if there is a containment mapping from $\mathcal{T}(Q_2)$ to $\text{chase}_{\Sigma_* \cup \Sigma_{l^*}}(\mathcal{T}(Q_1))$.*

A proof sketch is given in the full paper [11]. Next we show how to use proposition 3 to decide containment even if Q_1 is a proper $(*, _, l^*)$ -query.

In the rest of this section l will denote either a label in \mathcal{L} or the symbol $_$. Observe that for any $l \in \mathcal{L} \cup \{-\}$, $l^* = \bigcup_{0 \leq p} l^p$, where l^p is short for the concatenation of p successive l 's. More generally, let $Q(l_1^*, \dots, l_n^*)$ be a $(*, _, l^*)$ -query in which (l_1^*, \dots, l_n^*) are all the *occurrences* of starred symbols (the l_i 's are not necessarily distinct). Such a query is equivalent to an infinite union of star-free queries:

$$Q(l_1^*, \dots, l_n^*) = \bigcup_{0 \leq p_1, \dots, 0 \leq p_n} Q(l_1^{p_1}, \dots, l_n^{p_n})$$

The key to our containment test is that this infinite union can be replaced with a finite one. For any $(*, _, l^*)$ -query Q let $\text{sfs}(Q)$ be the *star-free size* of Q , defined as the count of all occurrences of non-Kleene-starred labels in Q . For example, for $Q(x, y) \leftarrow x a.b^* y, y * .c z$ we have $\text{sfs}(Q) = 2$.

Proposition 4. *Let Q_1, Q_2 be two $(*, _, l^*)$ -queries and let $k \stackrel{\text{def}}{=} \text{sfs}(Q_2) + 1$. Then, $Q_1 \subseteq Q_2$ if and only if*

$$\bigcup_{0 \leq p_1 \leq k, \dots, 0 \leq p_n \leq k} Q_1(l_1^{p_1}, \dots, l_n^{p_n}) \subseteq Q_2$$

The proof is given in the full paper [11]. We can now give our decision procedure for containment of unions of $(*, _, l^*, |)$ -queries, which has four steps:

Step 1: We first translate away the $|$, obtaining finite unions U_1, U_2 of $(*, _, l^*)$ -queries.

Step 2: Next we use proposition 4 to obtain from U_1 a finite union of star-free queries SF_1 , which must be checked for containment in U_2 ⁴.

Step 3: Containment of SF_1 in U_2 is decided using the following easy result:

⁴ An alternative way of obtaining SF_1 is by chasing the queries in U_1 with $\forall x \forall y [l^*(x, y) \rightarrow x = y \vee \exists z l(x, z) \wedge l^*(z, y)]$ (and similarly for $*$). This would result in a non-terminating chase, which we however could stop after sufficiently many steps.

Proposition 5. *The union of star-free queries $\bigcup_{i=1}^n Q_i$ is contained in the union of $(*, -, l^*)$ -queries $\bigcup_{j=1}^m Q'_j$ if and only if for every $1 \leq i \leq n$ there is a $1 \leq j \leq m$ such that $Q_i \subseteq Q'_j$.*

Step 4: Finally, checking each star-free Q_i for containment in Q'_j is done using proposition 3.

The upper bound for this algorithm is straightforward, proven in [11]:

Theorem 2. *Containment of unions of $(*, -, l^*, |)$ -queries is in Π_2^P .*

5 Lower Bounds

(l^*) -Queries, $(*, -)$ -Queries. The $|$ operator corresponds to the union and containment for unions of conjunctive queries is Π_2^P -complete [24]. But it turns out that even in the absence of $|$ we have Π_2^P -hardness results, with completely different proofs:

Theorem 3. *Containment of (l^*) -queries is Π_2^P -hard. Containment of $(*, -)$ -queries is Π_2^P -hard.*

As we pointed out in figure 1, the Π_2^P lower bound for containment of $(*, -, l^*, |)$ -queries follows (independently) from three sources: the two lower bounds in the previous theorem and the one in remark 2.

W-Queries. The following result shows that a more liberal nesting of regular path expressions within the Kleene star is problematic in terms of complexity of containment. If we allow concatenations of labels within the Kleene star, we obtain the W-fragment, whose lower bound for containment is PSPACE (a proof is in the full paper [11]):

Theorem 4. *Containment of W-queries is PSPACE-hard.*

As pointed out in remark 3, a bit more nesting yields EXPSPACE-hardness!

6 Containment under Dependencies

The $(*, |)$ -Fragment. This is where our technique of relational translation is most effective. First recall that by translating $|$ away, any union of $(*, |)$ -queries is equivalent to a union of $(*)$ -queries. Recall also that any set C of $(*, |)$ -dependencies is translated into a set $\mathcal{T}(C)$ of DEDs. By definition, “the dependencies in C are full” means that the DEDs in $\mathcal{T}(C)$ are full.

Since the DEDs in Σ_* are all full, the fact that containment of unions of $(*, |)$ -queries under full $(*, |)$ -dependencies is decidable follows from our extension of the chase to DEDs [11] and the following result:

Theorem 5. *Let C be a set of full $(*, |)$ -dependencies, and U_1, U_2 two unions of $(*, |)$ -queries. Let the equivalent unions of $(*)$ -queries be $\bigcup_{i=1}^n Q_i$, respectively $\bigcup_{j=1}^m Q'_j$. Then U_1 is contained in U_2 under C if and only if for every $1 \leq i \leq n$ there exists $1 \leq j \leq m$ such that $\mathcal{T}(Q_i)$ is contained in $\mathcal{T}(Q'_j)$ under $\Sigma_* \cup \mathcal{T}(C)$.*

The proof exploits the work we already did in section 4 and is omitted.

The $(*, _)$ -Fragment. As stated, this problem is open. However, we have two variations of it, one decidable, the other one, surprisingly, not.

Variation 1: V-Dependencies. Consider a subclass of full $(*, _)$ -dependencies, called V-dependencies, which disallow

- occurrences of the wildcard $_$ in the premise of the implication, and
- occurrences of $*$ in the conclusion of the implication (see formula (4)).

Theorem 6. *Containment of unions of $(*, _)$ -queries under full V-dependencies is decidable.*

The proof is omitted. The decision procedure is basically the same as the one for deciding containment of unions of $(*, _, l^*)$ -queries without dependencies: consider only a finite union of star-free queries, and check containment chasing with Σ_* and (as only difference from that case) with the translation of the V-dependencies.

Variation 2: Attributed Models. Suppose now that we restrict the full $(*, _)$ -dependencies even more, forcing their atoms to be star-free. We obtain precisely the full DEDs. But assume that we allow a special class of semistructured databases, in which the data graph can be “adorned” by attaching attributes to its nodes. More precisely, attributed models have schema $\mathcal{L}\text{-Rel} \cup \mathcal{A}$, where \mathcal{A} is a set of binary relations names, called *attributes*, who are disjoint from \mathcal{L} . The only difference between an attribute and a label is that the former is not included in the interpretation of $_$, while the latter is ⁵.

Theorem 7. *Containment of a $(*, _)$ -query in a union of $(*, _)$ -queries under full DEDs, but over attributed models, is undecidable.*

The proof is omitted, but very similar to that of theorem 8.

The (l^*) -Fragment. Surprisingly, this problem is undecidable, despite the syntactic similarity of the (l^*) and $(*)$ -fragments. We show a stronger undecidability result, which holds even if the dependencies are star-free, thus corresponding to purely relational full DEDs.

Theorem 8. *The containment of an l^* -query in a union of l^* -queries in the presence of full DEDs is undecidable.*

Proof: By reduction from the following problem: Given context-free grammar $G = (\Sigma, N, S, P)$ where Σ is the set of terminals (containing at least two symbols), N the nonterminals, $S \in N$ the start symbol, $P \subseteq N \times (\Sigma \cup N)^*$ the productions, and $L(G)$ the language generated by G , it is undecidable if $L(G) = \Sigma^*$ [18].

⁵ This model is similar in spirit to the XML data model and XPath specification [27], where attribute nodes are not reachable by navigation along the *child* axis.

The Reduction. Given context-free grammar $G = (\Sigma, N, S, P)$, we construct an instance of containment as follows:

$$Q \subseteq_D Q_S \cup Q_{cyc} \cup \bigcup_{\sigma_1 \neq \sigma_2 \in \Sigma} Q_{\sigma_1, \sigma_2}$$

$$\begin{aligned} Q() &\leftarrow \mathbf{b} H^* \mathbf{e} \quad (\mathbf{b}, \mathbf{e} \text{ constants}, H \notin \Sigma \cup N), \\ Q_{\sigma_1, \sigma_2}() &\leftarrow x \sigma_1 y, x \sigma_2 y \quad (\sigma_1, \sigma_2 \in \Sigma), \end{aligned}$$

$$\begin{aligned} Q_S() &\leftarrow \mathbf{b} S \mathbf{e}, \\ Q_{cyc}() &\leftarrow x H.H^* x \end{aligned}$$

D consists of the following full, star-free DERPDS (DEDs):

$$\begin{aligned} (\text{fn}) \quad &\forall x, y, z [x H y \wedge x H z \rightarrow y = z], & (\text{inj}) \quad &\forall x, y, z [y H x \wedge z H x \rightarrow y = z] \\ (\text{sym}) \quad &\forall x, y [x H y \rightarrow \bigvee_{\sigma \in \Sigma} x \sigma y], (d_p) \quad &\forall x_0, \dots, x_k [\bigwedge_{i=1}^k x_{i-1} M_i x_i \rightarrow x_0 N x_k] \\ & & & \text{(for every } p = N \rightarrow M_1 \dots M_k \in P) \end{aligned}$$

Now observe that (5) below holds for any binary queries. In addition, we claim that (6) holds as well, by construction, implying that it is a reduction.

$$Q \not\subseteq_D Q_S \cup Q_{cyc} \cup \bigcup_{\sigma_1 \neq \sigma_2 \in \Sigma} Q_{\sigma_1, \sigma_2} \quad \Leftrightarrow \quad (5)$$

$$\exists I [I \models D \wedge Q(I) \neq \emptyset \wedge Q_S(I) = \emptyset \wedge Q_{cyc}(I) = \emptyset \wedge \bigwedge_{\sigma_1 \neq \sigma_2 \in \Sigma} Q_{\sigma_1, \sigma_2}(I) = \emptyset] \quad \Leftrightarrow \quad (6)$$

$$\exists (w \in \Sigma^*) w \notin L(G)$$

The intuition behind the claim is that the instance I which is a counterexample for containment encodes a word w from $\Sigma^* \setminus L(G)$.

Proof of Claim (6): \Rightarrow : Assuming the standard interpretation of H^* , $Q(I) \neq \emptyset$ implies that there exists a path of H -edges from \mathbf{b} to \mathbf{e} in I . $Q_{cyc}(I) = \emptyset$ implies that all paths of H -edges are simple (no cycles). $I \models (\text{fn}) \wedge (\text{inj})$ implies that there is a unique (simple) path of H -edges from \mathbf{b} to \mathbf{e} which we call *the H -chain*. $I \models (\text{sym})$ says that every H -edge has in parallel with it a σ -edge for some symbol $\sigma \in \Sigma$, and it follows from $\bigwedge_{\sigma_1, \sigma_2} Q_{\sigma_1, \sigma_2}(I) = \emptyset$ that this edge is unique. The H -chain thus corresponds to a string w in Σ^* , of length equal to that of the H -chain. Each H -edge along the chain corresponds to a position in w . We make the following *subclaim*: let x, y be the source, respectively target nodes of a subchain of the H -chain, and let u be the corresponding substring of w . Let N be any nonterminal such that there exists a derivation of u in G starting from N . Then there is an N -edge from x to y in I . The subclaim is shown by induction on the length of the derivation, and it uses the fact that $I \models \bigwedge_{p \in P} (d_p)$. Together with $Q_S(I) = \emptyset$, the subclaim implies that there is no derivation of w in G starting from the start symbol S , in other words $w \notin L(G)$.

\Leftarrow : Starting from w , build the minimal model I consisting of (i) an H -chain of length $|w|$, (ii) the corresponding parallel edges spelling w , and (iii) for every subchain from node x to node y corresponding to the substring u of w , and every nonterminal N from which there is a derivation of u in G , add an N -edge from x to y . (i) implies $(I \models (\text{fn}) \wedge (\text{inj})) \wedge Q_{cyc}(I) = \emptyset \wedge Q(I) \neq \emptyset$, (ii) ensures

$(I \models (\text{symp})) \wedge \bigwedge Q_{\sigma_1, \sigma_2}(I) = \emptyset$ and (iii) guarantees $I \models \bigwedge_{p \in P}(d_p)$. $w \notin L(G)$ and the minimality of the model enforce $Q_S(I) = \emptyset$. •

7 Rewriting with Views under Dependencies

[9] introduces the *chase&backchase* (C&B) algorithm for rewriting queries with views under dependencies⁶. Due to space constraints we can only sketch here the idea and we omit proofs. The strategy of the C&B algorithm is to reduce the problem of rewriting with views to the problem of rewriting under dependencies. If V_i is a view name and QV_i the query that defines it, we capture V_i by writing a pair of inclusion dependencies that essentially say $V_i \subseteq QV_i$ and $QV_i \subseteq V_i$. Denote the set of all such pairs of dependencies with D_V and let us also assume that we rewrite under an additional set D of dependencies.

The C&B algorithm on a query Q has two phases. First the **chase** of Q with $D \cup D_V$. The dependencies in D_V that apply are full, so if those in D are full too, the chase will terminate, with a query we call the **universal plan** UP because it explicitly mentions all views that can be used to answer Q . The second phase is the **backchase** which considers all *subqueries* of the universal plan UP (subsets of its conjuncts, mentioning all distinguished variables). The output of the algorithm is the set of those subqueries equivalent to Q for whom the removal of any conjunct compromises this equivalence. We call such queries *minimal* rewritings of Q ⁷. The subqueries of the universal plan are tested for equivalence to Q again by chasing (see [11] for an illustration on our motivating example).

The (*)-Fragment. The C&B algorithm applies almost directly here. We should point out that the views may not be binary relations and therefore the rewritings we obtain will not correspond to pure (*)-queries, but rather may contain relational atoms with the view names. We have the following *completeness* result for the algorithm:

Theorem 9. *Let Q be a (*)-query, V be a set of (*)-views and D be a set of full (*)-dependencies. Let $E = \Sigma_* \cup \mathcal{T}(D) \cup \mathcal{T}(D_V)$ and let $UP = \text{chase}_E(\mathcal{T}(Q))$ (the chase terminates).*

Then, for any minimal rewriting Q' of Q with V , $\mathcal{T}(Q')$ is a subquery of UP .

The (*,|)-Fragment. In this case the query and views translate to unions of conjunctive queries and the (*,|)-dependencies translate to DEDs. If we plug into the C&B method the extended chase with DEDs (see [11]), we obtain a *union* of universal plans U_1, \dots, U_n after the chase phase. Each U_i plan is backchased yielding a set of minimal subqueries S_i . Every entry in the cartesian product

⁶ This is done in [9] for *path-conjunctive* queries and dependencies, which generalize the relational setting to a data model with dictionaries and complex values that also captures the OO setting.

⁷ Under a *monotonic cost assumption* minimal queries are cheaper.

$S_1 \times \dots \times S_n$ corresponds to a set of queries whose union is a rewriting of $\mathcal{T}(Q)$. We call this extension of the C&B algorithm the *disjunctive C&B* algorithm. We say that a union of queries is *reduced* if all members are minimal and none of them is contained in another member. The following result implies that the disjunctive C&B algorithm is *complete* for the $(*, |)$ -fragment.

Theorem 10. *Given a $(*, |)$ -query Q , and one of its reduced rewritings $Q' = Q'_1 \cup \dots \cup Q'_m$, for every $1 \leq j \leq m$, there is some $1 \leq i \leq n$ such that $\mathcal{T}(Q'_j)$ is a subquery of U_i .*

8 Conclusions

In this work, we propose a classification of conjunctive regular queries (CRPQ) and the associated constraint languages by the expressivity of the regular path expressions allowed in the conjuncts. We have studied the complexity of containment, with or without integrity constraints for the various fragments proposed. For certain fragments we have also studied the completeness of a specific kind of algorithm (chase & backchase) for rewriting with views under constraints.

A subtle observation that can be made based on the results we have obtained is that $_$ is “more” than the union (the $|$ actually) of the labels that occur in a given context. Indeed, one might attempt to contradict the decidability for the $(*, |)$ -fragment by reducing $(*, _)$ -queries and $_$ -dependencies to $(*, |)$ -queries and $_$ -dependencies, using a translation like $_ = l_1 | \dots | l_n | f$ where l_1, \dots, l_n are all the labels mentioned in the queries and dependencies and f is a fresh label. This attempt fails because it does not capture the equivalence $* = \bigcup_{n \geq 0} _^n$, which in turn is essential for the undecidability result. Of course, the correct translation an infinite disjunction of labels takes us out of the languages considered here.

We conclude that as a query language feature regular expressions are surprisingly “naughty”, in the sense that adding supposedly innocuous operators to some fragments causes surprising increases in complexity. (For example, adding either $*$ or $_$ to the fragment of conjunctive queries does not affect complexity of containment –still NP–, but adding *both* raises the complexity to Π_2^P .)

We are leaving some interesting problems open. One is the upper bound on containment in the W fragment. Another open problem is the decidability of containment under constraints in the $(*, _)$ -fragment. The reader can see that several open questions can be formulated about rewriting with views in certain fragments.

Since the submission of this work, we have applied our results to conjunctive queries over XML documents with XPath [27] expressions in their conjuncts [10].

References

1. S. Abiteboul, P. Buneman, and D. Suciu. *Data on the Web: From Relations to Semistructured Data and XML*. Morgan Kaufman, 1999.
2. S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.

3. C. Beeri and M. Vardi. A proof procedure for data dependencies. *JACM*, 31(4), 1984.
4. D. Calvanese, G. De Giacomo, and M. Lenzerini. On the Decidability of Query Containment under Constraints In *PODS*, 1998.
5. D. Calvanese, G. De Giacomo, M. Lenzerini, and M. Vardi. Rewriting of Regular Expressions and Regular Path Queries. In *PODS*, 1999.
6. D. Calvanese, G. De Giacomo, M. Lenzerini, and M. Vardi. Containment of conjunctive regular path queries with inverse. In *KR*, 2000.
7. Ashok Chandra and Philip Merlin. Optimal implementation of conjunctive queries in relational data bases. In *STOC*, 1977.
8. A. Deutsch, M. Fernandez, D. Florescu, A. Levy, and D. Suciu. A Query Language for XML. In *WWW8*, 1999.
9. Alin Deutsch, Lucian Popa, and Val Tannen. Physical Data Independence, Constraints and Optimization with Universal Plans. In *VLDB*, 1999.
10. A. Deutsch and V. Tannen. Containment and Integrity Constraints for XPath Fragments. In *KRDB 2001*.
11. A. Deutsch and V. Tannen. Optimization Properties for Classes of Conjunctive Regular Path Queries. Technical Report MS-CIS-01-20, University of Pennsylvania, 2001. Available from <http://db.cis.upenn.edu/cgi-bin/Person.perl?adeutsch>
12. M. Fernandez, D. Florescu, J. Kang, A. Levy, and D. Suciu. Strudel: A web-site management system. In *SIGMOD*, 1997.
13. Daniela Florescu, Alon Y. Levy, and Dan Suciu. Query containment for conjunctive queries with regular expressions. In *PODS*, 1998.
14. E. Grädel. On the restraining power of guards. *J. of Symbolic Logic*, 64, 1999.
15. Gösta Grahne and Alberto O. Mendelzon. Tableau techniques for querying information sources through global schemas. In *ICDT*, 1999.
16. G. Grahne and A. Thomo. An optimization technique for answering regular path queries. In *WebDB*, 2000.
17. G. Grahne and A. Thomo. Algebraic rewritings for regular path queries. *ICDT'01*.
18. J. Hopcroft and J. Ullman. *Introduction to automata theory, languages, and computation*. Addison-Wesley, 1979.
19. Frank Neven and Thomas Schwentick. Query automata. In *PODS*, 1999.
20. Frank Neven and Thomas Schwentick. Expressive and efficient pattern languages for tree-structured data. In *PODS*, 2000.
21. C. Papadimitriou. *Computational Complexity*. Addison-Wesley, 1994.
22. L. Popa, A. Deutsch, A. Sahuguet, and V. Tannen. A Chase Too Far? *SIGMOD* 2000.
23. Lucian Popa and Val Tannen. An equational chase for path-conjunctive queries, constraints, and views. In *ICDT*, 1999.
24. Yehoshua Sagiv and Mihalis Yannakakis. Equivalences among relational expressions with the union and difference operators. *Journal of the ACM*, 27, 1980.
25. P. van Emde Boas. The convenience of tilings. In A. Sorbi(Ed.) *Complexity, Logic, and Recursion Theory*, pp. 331–363, 2000.
26. W3C. Extensible Markup Language (XML) 1.0. W3C Recommendation 10-February-1998. Available from <http://www.w3.org/TR/1998/REC-xml-19980210>.
27. W3C. XML Path Language (XPath) 1.0. W3C Recommendation 16 November 1999. Available from <http://www.w3.org/TR/xpath>.
28. M. Yannakakis and C. Papadimitriou. Algebraic dependencies. *JCSS*, 25, 1982.