

# *Consistent query answering via ASP from different perspectives: Theory and practice*

MARCO MANNA, FRANCESCO RICCA and GIORGIO TERRACINA

*Department of Mathematics, University of Calabria, Cosenza, Italy*  
(e-mail: {manna,ricca,terracina}@mat.unical.it)

*submitted 23 November 2010; revised 15 December 2010; accepted 24 January 2011*

---

## **Abstract**

A data integration system provides transparent access to different data sources by suitably combining their data, and providing the user with a unified view of them, called *global schema*. However, source data are generally not under the control of the data integration process; thus, integrated data may violate global integrity constraints even in the presence of locally consistent data sources. In this scenario, it may be anyway interesting to retrieve as much consistent information as possible. The process of answering user queries under global constraint violations is called *consistent query answering* (CQA). Several notions of CQA have been proposed, e.g., depending on whether integrated information is assumed to be *sound*, *complete*, *exact*, or a variant of them. This paper provides a contribution in this setting: it unifies solutions coming from different perspectives under a common Answer-Set Programming (ASP)-based core, and provides query-driven optimizations designed for isolating and eliminating inefficiencies of the general approach for computing consistent answers. Moreover, the paper introduces some new theoretical results enriching existing knowledge on the decidability and complexity of the considered problems. The effectiveness of the approach is evidenced by experimental results.

**KEYWORDS:** Answer-Set Programming, data integration, consistent query answering

---

## **1 Introduction**

The enormous amount of information dispersed over many data sources, often stored in different heterogeneous databases, has recently boosted the interest for data integration systems (Lenzerini 2002). Roughly speaking, a data integration system provides transparent access to different data sources by suitably combining their data, and providing the user with a unified view of them, called *global schema*. In many cases, the application domain imposes some consistency requirements on integrated data. For instance, it may be at least desirable to impose some integrity constraints (ICs), such as primary/foreign keys, on the global relations. It may be the case that data stored at the sources may violate global ICs when integrated, since in general data sources are not under the control of the data integration process.

The standard approach to this problem basically consists of explicitly modifying the data in order to eliminate IC violations (data cleaning). However, the explicit repair of data is not always convenient or possible. Therefore, when answering a user query, the system should be able to “virtually repair” relevant data (in the line of Arenas *et al.* 2003; Bertossi *et al.* 2005; Chomicki and Marcinkowski 2005), in order to provide consistent answers; this task is also called consistent query answering (CQA).

The database community has spent considerable efforts in this area; relevant research results have been obtained to clarify semantics, decidability, and complexity of data integration under constraints and, specifically, for CQA. In particular, several notions of CQA have been proposed (see Bertossi *et al.* 2005 for a survey), e.g., depending on whether the information in the database is assumed to be *sound*, *complete*, or *exact*. However, while efficient systems are already available for simple data integration scenarios, solutions being both scalable and comprehensive have not been implemented yet for CQA, mainly due to the fact that handling inconsistencies arising from constraints violation is inherently hard. Moreover, mixing different kinds of constraints (e.g., denial constraints and inclusion dependencies) on the same global database makes, often, the query answering process undecidable (Abiteboul *et al.* 1995; Cali *et al.* 2003a).

This paper provides some contributions in this setting. Specifically, it first starts from different state-of-the-art semantic perspectives (Arenas *et al.* 2003; Cali *et al.* 2003a; Chomicki and Marcinkowski 2005) and revisits them in order to provide a uniform, common core based on Answer-Set Programming (ASP) (Gelfond and Lifschitz 1988, 1991). Thus, it provides query-driven optimizations, in the light of the experience we gained in the INFOMIX (Leone *et al.* 2005) project in order to overcome the limitations observed in real-world scenarios. The main contributions of this paper can be summarized as follows:

- A theoretical analysis of considered semantics which extends previous results.
- The definition of a unified framework for CQA based on a purely declarative, logic-based approach which supports the most relevant semantics assumptions on source data. Specifically, the problem of CQA is reduced to cautious reasoning on (disjunctive) ASP programs with aggregates (Faber *et al.* 2010) automatically built from both the query and involved constraints.
- The definition of an optimization approach designed to (1) “localize” and limit the inefficient part of the computation of consistent answers to small fragments of the input, and (2) cast down the computational complexity of the repair process if possible.
- The implementation of the entire framework in a full-fledged prototype system.
- The capability of handling large amounts of data, typical of real-world data integration scenarios, using as internal query evaluator the DLV<sup>DB</sup> (Terracina, Leone, *et al.* 2008) system; indeed, DLV<sup>DB</sup> allows for mass-memory database evaluations and distributed data management features.

In order to assess the effectiveness of the proposed approach, we carried out experimental activities both on a real-world scenario and on synthetic data, comparing its behavior on different semantics and constraints.

The plan of the paper is as follows. Section 2 formally introduces the notion of CQA under different semantics and some new theoretical results on decidability and complexity for this problem. Section 3 first introduces a unified (general) solution to handle CQA via ASP, and then presents some optimizations. Section 4 describes the benchmark framework that we adopted in the tests and discusses on the obtained results. Finally, Section 5 compares related work and draws some conclusive considerations.

## 2 Data integration framework

In this paper, we exploit the data integration setting to point out motivations and challenges underlying CQA. However, as it will be clarified in the following, techniques and results provided in the paper hold also for a single database setting. We next formally describe the adopted data integration framework.

The following notation will be used throughout the paper. We always denote by  $\Gamma$  a countably infinite domain of totally ordered values; by  $t$  a tuple of values from  $\Gamma$ ; by  $X$  a variable; by  $\bar{x}$  a sequence  $X_1, \dots, X_n$  of (not necessarily distinct) variables; and by  $|\bar{x}| = n$  its length. Let  $\bar{x}, \bar{x}'$  be two sequences of variables, then we denote by  $\bar{x} - \bar{x}'$  the sequence obtained from  $\bar{x}$  by discarding a variable if it appears in  $\bar{x}'$ . Whenever all the variables of sequence  $\bar{x}$  appear in another sequence  $\bar{x}'$ , we simply write  $\bar{x} \leq \bar{x}'$ . Given a sequence  $\bar{x}$  and a set  $\pi \subseteq \{1, \dots, |\bar{x}|\}$ , we denote by  $\bar{x}^\pi$  the sequence obtained from  $\bar{x}$  by discarding a variable if its position is not in  $\pi$ . (Similarly, given a tuple  $t$  and a set  $\pi \subseteq \{1, \dots, |t|\}$ , we denote by  $t^\pi$  the tuple obtained from  $t$  by discarding a value if its position is not in  $\pi$ .) Moreover, we denote by  $\sigma(\bar{x})$  a conjunction of comparison atoms of the form  $X \odot X'$ , where  $\odot \in \{\leq, \geq, <, >, \neq\}$ , and by  $\ominus$  the symmetric difference operator between two sets.

A relational database schema is a pair  $\mathcal{R} = \langle names(\mathcal{R}), constr(\mathcal{R}) \rangle$ , where  $names(\mathcal{R})$  and  $constr(\mathcal{R})$  are the relation names and the ICs of  $\mathcal{R}$ , respectively. The arity of a given relation  $r \in names(\mathcal{R})$  is denoted by  $arity(r)$ . A database (instance) for  $\mathcal{R}$  is any set of facts (Abiteboul *et al.* 1995) of the form:

$$\mathcal{F} = \{r(t) : r \in names(\mathcal{R}) \wedge t \text{ is a tuple from } \Gamma \wedge |t| = arity(r)\}$$

In the following, we adopt the *unique name assumption*, and  $dom(\mathcal{F})$  denotes the subset of  $\Gamma$  containing all the values appearing in the facts of  $\mathcal{F}$ .

Let  $r_1, \dots, r_m \in names(\mathcal{R})$ , the set  $constr(\mathcal{R})$  contains ICs of the form:

- (1)  $\forall \bar{x}_1, \dots, \bar{x}_m \neg [ r_1(\bar{x}_1) \wedge \dots \wedge r_m(\bar{x}_m) \wedge \sigma(\bar{x}_1, \dots, \bar{x}_m) ]$  (*denial constraints* – DCs)
- (2)  $\forall \bar{x}_\forall [ r_1(\bar{x}_1) \rightarrow \exists \bar{x}_{\exists\exists} r_2(\bar{x}_2) ]$  (*inclusion dependencies* – INDs);

where  $arity(r_i) = |\bar{x}_i|$ , for each  $i$  in  $[1, \dots, m]$ . In particular, for INDs, we require that all the variables within an  $\bar{x}_i$  ( $1 \leq i \leq 2$ ) are distinct,  $\bar{x}_\forall \leq \bar{x}_1$ ,  $\bar{x}_\forall \leq \bar{x}_2$ , and  $\bar{x}_{\exists\exists} = \bar{x}_2 - \bar{x}_\forall$ . Note that if  $|\bar{x}_{\exists\exists}| = 0$ , then  $\bar{x}_\forall = \bar{x}_2 \leq \bar{x}_1$ . In the case, we are only interested in emphasizing the relation names involved in an IND, we simply write

$r_1(\bar{x}_1) \rightarrow r_2(\bar{x}_2)$  or  $r_1 \rightarrow r_2$ . A database  $\mathcal{F}$  is said to be *consistent* w.r.t.  $\mathcal{R}$  if all ICs are satisfied. A *conjunctive query*  $cq(\bar{x})$  over  $\mathcal{R}$  is a formula of the form:

$$\exists \bar{x}_{1\exists}, \dots, \bar{x}_{m\exists} \ r_1(\bar{x}_1) \wedge \dots \wedge r_m(\bar{x}_m) \wedge \sigma(\bar{x}_1 \dots, \bar{x}_m)$$

where  $\bar{x}_{i\exists} \leq \bar{x}_i$  for each  $i$  in  $[1, \dots, m]$ ,  $\bar{w} = \bar{x}_1 - \bar{x}_{1\exists}, \dots, \bar{x}_m - \bar{x}_{m\exists}$  are the *free variables* of  $q$ , and  $\bar{x}$  contains only and all the variables of  $\bar{w}$  (with no duplicates, and possibly in a different order). A *union of conjunctive queries*  $q(\bar{x})$  is a formula of the form  $cq_1(\bar{x}) \vee \dots \vee cq_n(\bar{x})$ . In the following, for simplicity, the term query refers to a union of conjunctive queries, if not differently specified. Given a database  $\mathcal{F}$  for  $\mathcal{R}$ , and a query  $q(\bar{x})$ , the *answer* to  $q$  is the set of  $n$ -tuples of values  $ans(q, \mathcal{F}) = \{t : \mathcal{F} \models q(t)\}$ .

### 2.1 Data integration model

A *data integration system* is formalized (Lenzerini 2002) as a triple  $\mathcal{I} = \langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$  where

- .  $\mathcal{G}$  is the *global schema*. A *global database* for  $\mathcal{I}$  is any database for  $\mathcal{G}$ ;
- .  $\mathcal{S}$  is the *source schema*. A *source database* for  $\mathcal{I}$  is any database consistent w.r.t.  $\mathcal{S}$ ;
- .  $\mathcal{M}$  is the *global-as-view (GAV) mapping* that associates each element  $g$  in  $names(\mathcal{G})$  with a union of conjunctive queries over  $\mathcal{S}$ .

Let  $\mathcal{F}$  be a source database for  $\mathcal{I}$ . The *retrieved global database* is

$$ret(\mathcal{I}, \mathcal{F}) = \{g(t) : g \in names(\mathcal{G}) \wedge t \in ans(q, \mathcal{F}) \wedge q \in \mathcal{M}(g)\}$$

for  $\mathcal{G}$  satisfying the mapping. Note that when source data are combined in a unified schema with its own ICs, the retrieved global database might be inconsistent.

In the following, when it is clear from the context, we use simply the symbol  $\mathcal{D}$  to denote the retrieved global database  $ret(\mathcal{I}, \mathcal{F})$ . In fact, all results provided in the paper hold for any database  $\mathcal{D}$  complying with some schema  $\mathcal{G}$  but possibly inconsistent w.r.t. the constraints of  $\mathcal{G}$ .

#### Example 1

Consider a bank association that desires to unify the databases of two branches. The first (source) database models managers by using a relation  $man(code, name)$  and employees by a relation  $emp(code, name)$ , where  $code$  is a primary key for both tables. The second database stores the same data in a relation  $employee(code, name, role)$ . Suppose that the data have to be integrated under a global schema with two relations  $m(code)$  and  $e(code, name)$ , where the global ICs are:

- $\forall X_1, X_2, X_3 \ \neg[e(X_1, X_2) \wedge e(X_1, X_3) \wedge X_2 \neq X_3]$ , namely,  $code$  is the key of  $e$ ;
- $\forall X_1 [m(X_1) \rightarrow \exists X_2 \ e(X_1, X_2)]$ , i.e., an IND imposing that each manager code must be an employee code as well.

The mapping is defined by the following Datalog rules (as usual, see Abiteboul et al. 1995):

$$e(X_c, X_n) :- emp(X_c, X_n). \qquad m(X_c) :- man(X_c, -).$$

$$e(X_c, X_n) := \text{employee}(X_c, X_n, -). \quad m(X_e) := \text{employee}(X_c, -, \text{'manager'}).$$

Assume that *emp* stores tuples (*'e1','john'*), (*'e2','mary'*), (*'e3','willy'*), *man* stores (*'e1','john'*), and *employee* stores (*'e1','ann','manager'*), (*'e2','mary','manager'*), (*'e3','rose','emp'*). It is easy to verify that although the source databases are consistent w.r.t. local constraints, the global database, obtained by evaluating the mapping, violates the key constraint on *e* as both *john* and *ann* have the same code *e1*, and both *willy* and *rose* have the same code *e3* in table *e*.

## 2.2 CQA under different semantics

In case a database  $\mathcal{D}$  violates ICs, one can still be interested in querying the “consistent” information originating from  $\mathcal{F}$ . One possibility is to “repair”  $\mathcal{D}$  (by inserting or deleting tuples) in such a way that all the ICs are satisfied. But there are several ways to “repair”  $\mathcal{D}$ . As an example, in order to satisfy an IND of the form  $r_1 \rightarrow r_2$  one might either remove violating tuples from  $r_1$  or insert new tuples in  $r_2$ . Moreover, the repairing strategy depends on the particular semantic assumption made on the data integration system. Semantic assumptions may range from (strict) soundness to (strict) completeness. Roughly speaking, completeness complies with the *closed world assumption* where missing facts are assumed to be false; on the contrary, soundness complies with the *open world assumption* where  $\mathcal{D}$  may be incomplete. We next define CQA under some relevant semantics, namely, *loosely-exact*, *loosely-sound*, and *CM-complete* (Arenas et al. 2003; Cali et al. 2003a; Chomicki and Marcinkowski 2005). More formally, let  $\Sigma$  denote a *semantics*, and  $\mathcal{D}$  a possibly inconsistent database for  $\mathcal{G}$ , a database  $\mathcal{B}$  is said to be a  $\Sigma$ -repair for  $\mathcal{D}$  if it is consistent w.r.t.  $\mathcal{G}$  and one of the following conditions holds:

- (1)  $\Sigma = \text{CM-complete}$ ,  $\mathcal{B} \subseteq \mathcal{D}$ , and  $\nexists \mathcal{B}' \subseteq \mathcal{D}$  such that  $\mathcal{B}'$  is consistent and  $\mathcal{B}' \supset \mathcal{B}$ ;
- (2)  $\Sigma = \text{loosely-sound}$ , and  $\nexists \mathcal{B}'$  such that  $\mathcal{B}'$  is consistent and  $\mathcal{B}' \cap \mathcal{D} \supset \mathcal{B} \cap \mathcal{D}$ ;
- (3)  $\Sigma = \text{loosely-exact}$ , and  $\nexists \mathcal{B}'$  such that  $\mathcal{B}'$  is consistent and  $\mathcal{B}' \ominus \mathcal{D} \subset \mathcal{B} \ominus \mathcal{D}$ .

The *CM-complete* semantics allows a minimal number of deletions in each repair to avoid empty repairs, if possible, but does not allow insertions. The *loosely-sound* semantics allows insertions and a minimal amount of deletions. Finally, the *loosely-exact* semantics allows both insertions and deletions by minimization of the symmetric difference between  $\mathcal{D}$  and the repairs.

### Definition 1

Let  $\mathcal{D}$  be a database for a schema  $\mathcal{G}$ , and  $\Sigma$  be a semantics. The *consistent answer* to a query  $q$  w.r.t.  $\mathcal{D}$  is the set  $\text{ans}_\Sigma(q, \mathcal{G}, \mathcal{D}) = \{t : t \in \text{ans}(q, \mathcal{B}) \text{ for each } \Sigma\text{-repair } \mathcal{B} \text{ for } \mathcal{D}\}$ . CQA is the problem of computing  $\text{ans}_\Sigma(q, \mathcal{G}, \mathcal{D})$ .

Observe that other semantics have been considered in the literature, such as *sound*, *complete*, *exact*, *loosely-complete*, etc. (Cali et al. 2003a); however, some of them are trivial for CQA. As an example, in the *exact* semantics, CQA makes sense only if the retrieved database is already consistent with the global constraints, whereas in the *complete* and *loosely-complete* semantics, CQA will always return a void answer.

Note that the semantics considered in this paper address a wide significant range of ways to repair the retrieved database which are also relevant for CQA.

*Example 2*

By following Example 1, the retrieved global database admits exactly the following repairs under the *CM-complete* semantics:  $\mathcal{B}_1 = \{e('e2', 'mary'), e('e1', 'john'), e('e3', 'willy'), m('e1'), m('e2')\}$ ;  $\mathcal{B}_2 = \{e('e2', 'mary'), e('e1', 'john'), e('e3', 'rose'), m('e1'), m('e2')\}$ ;  $\mathcal{B}_3 = \{e('e2', 'mary'), e('e1', 'ann'), e('e3', 'willy'), m('e1'), m('e2')\}$ ;  $\mathcal{B}_4 = \{e('e2', 'mary'), e('e1', 'ann'), e('e3', 'rose'), m('e1'), m('e2')\}$ . The query  $m(X)$ , asking for the list of manager codes, has then both  $e1$  and  $e2$  as consistent answers, whereas the query  $e(X, Y)$ , asking for the list of employees, has only  $e('e2', 'mary')$  as a consistent answer ( $e$  is the only tuple in each *CM-complete* repair).

**2.3 Restricted classes of ICs**

The problem of computing CQA, under general combinations of ICs, is undecidable (Abiteboul et al. 1995). However, restrictions on ICs to retain decidability and identify tractable cases can be imposed.

*Definition 2*

Let  $r$  be a relation name of arity  $n$ , and  $\pi$  be a set of  $m \leq n$  indices from  $I = \{1, \dots, n\}$ . A *key dependency* (KD) for  $r$  consists of a set of  $n - m$  DCs, exactly one for each index  $i \in I - \pi$ , of the form  $\forall \bar{x}_1, \bar{x}_2 \neg(r(\bar{x}_1) \wedge r(\bar{x}_2) \wedge \bar{x}_1^i \neq \bar{x}_2^i)$ , where no variable occurs twice in each  $\bar{x}_i$  ( $1 \leq i \leq 2$ ),  $|\bar{x}_1| = |\bar{x}_2| = n$ , the sequence  $\bar{x}_1^\pi$  exactly coincides with  $\bar{x}_2^\pi$ , and  $\bar{x}_1^i$  is distinct from  $\bar{x}_2^j$  for each  $j \in I - \pi$ . The set  $\pi$  is called the *primary key* of  $r$  and is denoted by  $key(r)$ . We assume that at most one KD is specified for each relation (Cali et al. 2003a). Finally, for each relation name  $r'$  such that no DC is explicitly specified for, we say, without loss of generality, that  $key(r') = \{1, \dots, arity(r')\}$ .

*Definition 3*

Given an IND  $d$  of the form  $\forall \bar{x}_\forall [ r_1(\bar{x}_1) \rightarrow \exists \bar{x}_\exists r_2(\bar{x}_2) ]$ , we denote by  $\pi_L^d \subseteq \{1, \dots, arity(r_1)\}$  and  $\pi_R^d \subseteq \{1, \dots, arity(r_2)\}$  the two sets of indices induced by the positions of the variables  $\bar{x}_\forall$  in  $\bar{x}_1$  and  $\bar{x}_\exists$ , respectively. More formally,  $\pi_L^d = \{i : \bar{x}_1^i \text{ is universally quantified in } d\}$  and  $\pi_R^d = \{i : \bar{x}_2^i \text{ is universally quantified in } d\}$ .

For example, let  $d$  denote the IND  $\forall X_1, X_2 [ r_1(X_1, X_3, X_2) \rightarrow \exists X_4 r_2(X_4, X_2, X_1) ]$ . We have that  $\pi_L^d = \{1, 3\}$  and  $\pi_R^d = \{2, 3\}$ .

*Definition 4*

An IND  $d$  is said to be (i) a *foreign key* (FK) if  $\pi_R^d = key(r_2)$  (Abiteboul et al. 1995); (ii) a *foreign superkey* (FSK) if  $\pi_R^d \supseteq key(r_2)$  (Levene and Vincent 2000); and (iii) *non-key-conflicting* (NKC) if  $\pi_R^d \not\supseteq key(r_2)$  (Cali et al. 2003a). (Note that each FK is an FSK.)

*Definition 5*

An FSK  $d$  of the form  $r_1 \rightarrow r_2$  is said to be *safe* (SFSK) if  $\pi_L^d \subseteq key(r_1)$ . In particular, if  $d$  is a *safe* FK we call it an SFK.

Table 1. Data complexity of CQA (distinguishing between cyclic and acyclic INDs)

DCs	INDs	Loosely-sound	Loosely-exact	CM-complete
No	Any	In PTIME *	In PTIME *	In PTIME †
KD	No	coNP-c *	coNP-c *	coNP-c †
KD	NKC	coNP-c *	$\Pi_2^p$ -c *	in $\Pi_2^p$ † / in coNP †
KD	SFSK	In $\Pi_2^p$ ‡	In $\Pi_2^p$ ‡	In $\Pi_2^p$ † / in coNP †
KD	Any	Undecidable *	Undecidable *	In $\Pi_2^p$ † / In coNP †
Any	Any	Undecidable ¶	Undecidable ¶	$\Pi_2^p$ -c † / coNP-c †

\* Cali et al. (2003a); † Chomicki and Marcinkowski (2005); ‡Section 2.4; ¶ Abiteboul et al. (1995).

For example, let  $d$  denote the FSK  $\forall X_1, X_2 [ r_1(X_1, X_3, X_2) \rightarrow \exists X_4 r_2(X_4, X_2, X_1) ]$  where  $key(r_2) = \{3\}$ . Thus, if  $key(r_1) = \{1, 3\}$ ,  $d$  is SFSK, whereas if  $key(r_1) = \{1, 2\}$ ,  $d$  is not SFSK.

Table 1 summarizes known and new results about the computability and complexity of CQA under relevant classes of ICs and the three semantic assumptions considered in this paper. In particular, given a query  $q$  (without comparison atoms if  $\Sigma \in \{loosely-sound, loosely-exact\}$ ), we refer to the decision problem of establishing whether a tuple from  $dom(\mathcal{D})$  belongs to  $ans_\Sigma(q, \mathcal{G}, \mathcal{D})$  or not. Note that Chomicki and Marcinkowski (2005) have proved the computability and complexity of CQA for the *CM-complete* semantics in the case of conjunctive queries with comparison predicates. However, since in such a setting, there is a finite number of repairs each of finite size, then their results straightforwardly hold for a union of conjunctive queries as well. New decidability and complexity results for CQA under KDs and SFSKs only, with  $\Sigma \in \{loosely-sound, loosely-exact\}$ , are proved in Section 2.4.

### 2.4 Loosely-exact and loosely-sound semantics under KD and SFSK

In this section, we provide new decidability and complexity results for CQA under both the loosely-exact and the loosely-sound semantics with KDs and SFSKs. In the rest of the section, we always denote by  $\mathcal{G}$  a schema containing KDs and SFSKs only; by  $\mathcal{D}$  a possibly inconsistent database for  $\mathcal{G}$ ; by  $q$  a union of conjunctive queries without comparison atoms; and by  $\Sigma \in \{loosely-exact, loosely-sound\}$ .

In the following, some of the proofs are not included due to space constraints; they can be found in the extended version of this paper (Manna et al. 2011).

We first show that in the aforementioned hypothesis, the size of each repair is finite.

#### Definition 6

Let  $\mathcal{B}$  be a  $\Sigma$ -repair for  $\mathcal{D}$  and  $i \geq 0$  be a natural number. We inductively define the sets  $\mathcal{B}^i$  as follows. (1) If  $i = 0$ , then  $\mathcal{B}^0 = \mathcal{B} \cap \mathcal{D}$ . (2) If  $i > 0$ , then  $\mathcal{B}^i \subseteq \mathcal{B} - (\mathcal{B}^0 \cup \dots \cup \mathcal{B}^{i-1})$  is arbitrarily chosen in such a way that its facts are necessary and sufficient for satisfying all the INDs in  $constr(\mathcal{G})$  that are violated in  $\mathcal{B}^0 \cup \dots \cup \mathcal{B}^{i-1}$ . Observe that  $\mathcal{B} = \bigcup_{i \geq 0} \mathcal{B}^i$  and that  $\mathcal{B}^i \cap \mathcal{B}^j = \emptyset$  for each  $j \neq i$ .



*Lemma 1*

Let  $\mathcal{B}$  be a  $\Sigma$ -repair for  $\mathcal{D}$ , then (1) the key of each fact in  $\mathcal{B}$  only contains values from  $\text{dom}(\mathcal{D})$ ; and (2)  $|\mathcal{B}|$  is finite.

*Proof*

(1) Let  $i > 0$  be a natural number. Let  $r_i(t_i)$  be a fact in  $\mathcal{B}^i$  such that there is an index  $j \in \text{key}(r_i)$  for which  $t_i^j \notin \text{dom}(\mathcal{B}^0)$ . Let  $r_{i-1}(t_{i-1})$  be one of the facts in  $\mathcal{B}^{i-1}$  that forces the presence of  $r_i(t_i)$  in  $\mathcal{B}^i$  for satisfying some IND, say  $d$ . (Note that by Definition 6, there must be at least one of such a fact because  $\mathcal{B}^i$  would otherwise violate condition 2, since  $r_i(t_i)$  would be unnecessary.) Moreover, since  $d$  is an SFSK, then there must exist an index  $k \in \text{key}(r_{i-1})$  such that  $t_i^j = t_{i-1}^k$ . Thus,  $r_{i-1}(t_{i-1})$  contains a value being not in  $\text{dom}(\mathcal{B}^0)$  inside its key as well as  $r_i(t_i)$ . Since  $i$  has been chosen arbitrarily, then value  $t_i^j$  has to be part of a fact of  $\mathcal{B}^0$ , which is clearly a contradiction.

(2) Since the key of each fact in  $\mathcal{B}$  can only contain values from  $\text{dom}(\mathcal{B}^0)$ , and  $|\text{dom}(\mathcal{B}^0)| \leq |\mathcal{B}^0| \cdot \alpha$  where  $\alpha = \max\{\text{arity}(g) : g \in \text{names}(\mathcal{G})\}$ , then  $|\mathcal{B}| \leq |\text{names}(\mathcal{G})| \cdot |\text{dom}(\mathcal{B}^0)|^\alpha \leq |\text{names}(\mathcal{G})| \cdot (\alpha \cdot |\mathcal{B}^0|)^\alpha \leq |\text{names}(\mathcal{G})| \cdot (\alpha \cdot |\mathcal{D}|)^\alpha$ .  $\square$

We next characterize representative databases for  $\Sigma$ -repairs.

*Definition 7*

Let  $\mathcal{B}$  be a  $\Sigma$ -repair for  $\mathcal{D}$ . We denote by  $\text{homo}(\mathcal{B})$  the (possibly infinite) set of databases defined in such a way that  $\mathcal{B}' \in \text{homo}(\mathcal{B})$  if and only if:

- $\mathcal{B}'$  can be obtained from  $\mathcal{B}$  by replacing each value (if any) that is not in  $\text{dom}(\mathcal{D})$  with a value from  $\Gamma - \text{dom}(\mathcal{D})$ ; and
- none of the values in  $\Gamma - \text{dom}(\mathcal{D})$  occurs twice in  $\mathcal{B}'$ .

Finally, we denote by  $h_{\mathcal{B}, \mathcal{B}'} : \text{dom}(\mathcal{B}') \rightarrow \text{dom}(\mathcal{B})$  the function (homomorphism) associating values in  $\text{dom}(\mathcal{B}')$  with values in  $\text{dom}(\mathcal{B})$ , where  $h_{\mathcal{B}, \mathcal{B}'}(\alpha) = \alpha$ , for each  $\alpha \in \text{dom}(\mathcal{D}) \cap \text{dom}(\mathcal{B}')$ .

Note that since (by Lemma 1) the key of each fact in  $\mathcal{B}$  only contains values from  $\text{dom}(\mathcal{D})$ , then  $|\mathcal{B}'| = |\mathcal{B}|$  holds.

For example, if  $\mathcal{B} = \{p(1, \varepsilon_1, \varepsilon_2), q(2, \varepsilon_2, \varepsilon_1)\}$  with  $\text{dom}(\mathcal{D}) = \{1, 2\}$  and  $\text{key}(p) = \text{key}(q) = \{1\}$ , then all of the following databases are in  $\text{homo}(\mathcal{B})$ :  $\{p(1, \varepsilon_1, \varepsilon_3), q(2, \varepsilon_2, \varepsilon_4)\}$ ,  $\{p(1, \varepsilon_4, \varepsilon_2), q(2, \varepsilon_3, \varepsilon_1)\}$ , and  $\{p(1, \varepsilon_5, \varepsilon_6), q(2, \varepsilon_7, \varepsilon_8)\}$ .

*Lemma 2*

If  $\mathcal{B}$  is a  $\Sigma$ -repair for  $\mathcal{D}$ , then each  $\mathcal{B}' \in \text{homo}(\mathcal{B})$  also is.

We next define the finite database  $\mathcal{D}^*$  having among its subsets a number of  $\Sigma$ -repairs sufficient for solving CQA.

*Definition 8*

Let  $c$  be a value in  $\Gamma - \text{dom}(\mathcal{D})$ . Consider the largest (possibly inconsistent) database, say  $C$ , constructible on the domain  $\text{dom}(\mathcal{D}) \cup \{c\}$  such that  $f \in C$  iff the value  $c$  does not appear in the key of  $f$ . Let  $\mathcal{N}$  be a fixed set of values arbitrarily chosen from  $\Gamma - \text{dom}(\mathcal{D})$  whose cardinality is equal to the number of occurrences of  $c$  in



$C$ . We denote by  $\mathcal{D}^*$  one possible database for  $\mathcal{G}$  obtained from  $C$  by replacing each occurrence of  $c$  with a value from  $\mathcal{N}$  in such a way that each value in  $\mathcal{N}$  occurs exactly once in  $\mathcal{D}^*$ . ( $|C| = |\mathcal{D}^*|$ .)

For example, if  $\text{dom}(\mathcal{D}) = \{1, 2\}$  and  $\mathcal{G} = \{p\}$  with  $\text{arity}(p) = 2$  and  $\text{key}(p) = \{1\}$ , then  $C = \{p(1, 1), p(1, 2), p(1, c), p(2, 1), p(2, 2), p(2, c)\}$ . Let us fix  $\mathcal{N} = \{\varepsilon_1, \varepsilon_2\}$ . Thus,  $\mathcal{D}^*$  has the following form:  $\{p(1, 1), p(1, 2), p(1, \varepsilon_1), p(2, 1), p(2, 2), p(2, \varepsilon_2)\}$ .

*Proposition 1*

The following hold:

- $|\mathcal{N}| = \sum_{g \in \mathcal{G}} (\text{arity}(g) - |\text{key}(g)|) \cdot |\text{dom}(\mathcal{D})|^{|\text{key}(g)|} \cdot (|\text{dom}(\mathcal{D})| + 1)^{\text{arity}(g) - |\text{key}(g)| - 1}$
- $|\mathcal{D}^*| \leq \sum_{g \in \mathcal{G}} (|\text{dom}(\mathcal{D})| + 1)^{\text{arity}(g)} \leq \sum_{g \in \mathcal{G}} (\text{arity}(g) \cdot |\mathcal{D}| + 1)^{\text{arity}(g)}$

*Lemma 3*

If  $\mathcal{B}$  is a  $\Sigma$ -repair for  $\mathcal{D}$ , then there exists  $\mathcal{B}' \in \text{homo}(\mathcal{B})$  such that  $\mathcal{B}' \subseteq \mathcal{D}^*$ .

*Lemma 4*

Let  $\mathcal{B}$  be a  $\Sigma$ -repair for  $\mathcal{D}$ ,  $\mathcal{B}' \in \text{homo}(\mathcal{B})$ ,  $q$  be a query, and  $t$  be a tuple of values from  $\text{dom}(\mathcal{D})$ . If  $t \in \text{ans}(q, \mathcal{B}')$ , then  $t \in \text{ans}(q, \mathcal{B})$ .

*Proof*

Let  $q_i$  be one of the conjunctions in  $q$ , if  $t \in \text{ans}(q_i, \mathcal{B}')$ , then there is a substitution  $\mu'$  from the variables of  $q_i$  to values in  $\Gamma$  such that  $\mathcal{B}' \models q_i(t)$ . But since, by Definition 7, each fact in  $\mathcal{B}'$  is univocally associated with a unique fact in  $\mathcal{B}$  by preserving the values in  $\text{dom}(\mathcal{D})$ , and since all the extra values in  $\mathcal{B}'$  are distinct, then there must also be a substitution  $\mu$  such that  $\mathcal{B} \models q_i(t)$ . In particular, let  $x$  be a variable in  $q_i$ , then we can define  $\mu$  in such a way that  $\mu(x) = h_{\mathcal{B}, \mathcal{B}'}(\mu'(x))$ , where  $h$  is the homomorphism from  $\mathcal{B}'$  to  $\mathcal{B}$  (see Definition 7). Clearly, if  $t \in \text{ans}(q_i, \mathcal{B}')$  for at least one  $q_i$  in  $q$ , then  $t \in \text{ans}(q, \mathcal{B}')$  too and, consequently,  $t \in \text{ans}(q, \mathcal{B})$ .  $\square$

The next theorem states the decidability of CQA under both the loosely-exact and the loosely-sound semantics with KDs and SFSKs only.

*Theorem 1*

Let  $\mathcal{B}$  be a  $\Sigma$ -repair for  $\mathcal{D}$ ,  $q$  a query, and  $t$  a tuple from  $\text{dom}(\mathcal{D})$ . Let  $\mathbb{B} \subseteq 2^{\mathcal{D}^*}$  denote the set of all  $\Sigma$ -repairs contained in  $\mathcal{D}^*$ . Then,  $t \in \text{ans}_{\Sigma}(q, \mathcal{G}, \mathcal{D})$  **iff**  $t \in \text{ans}(q, \mathcal{B}) \ \forall \mathcal{B} \in \mathbb{B}$ .

*Proof*

( $\Rightarrow$ ) We have to prove that if  $t \in \text{ans}_{\Sigma}(q, \mathcal{G}, \mathcal{D})$ , then  $t \in \text{ans}(q, \mathcal{B})$  for each  $\mathcal{B} \in \mathbb{B}$ , or equivalently, if  $t \notin \text{ans}(q, \mathcal{B})$  for some  $\mathcal{B} \in \mathbb{B}$ , then  $t \notin \text{ans}_{\Sigma}(q, \mathcal{G}, \mathcal{D})$ . This follows by the definition of  $\text{ans}_{\Sigma}(q, \mathcal{G}, \mathcal{D})$  and from the fact that  $\mathbb{B}$  only contains  $\Sigma$ -repairs.

( $\Leftarrow$ ) We have to prove that if  $t \in \text{ans}(q, \mathcal{B})$  for each  $\mathcal{B} \in \mathbb{B}$ , then  $t \in \text{ans}_{\Sigma}(q, \mathcal{G}, \mathcal{D})$ . Assume that  $t \in \text{ans}(q, \mathcal{B})$  for each  $\mathcal{B} \in \mathbb{B}$  but  $t \notin \text{ans}_{\Sigma}(q, \mathcal{G}, \mathcal{D})$ . This would entail that there is a repair  $\mathcal{B}_0$  such that  $t \notin \text{ans}(q, \mathcal{B}_0)$ . But since  $t \in \text{ans}(q, \mathcal{B}')$  for each  $\mathcal{B}' \in \text{homo}(\mathcal{B}_0)$  (by Lemma 4), and since  $\mathbb{B} \cap \text{homo}(\mathcal{B}_0)$  always contains a repair, say  $\mathcal{B}''$  (by Lemma 3), then we have a contradiction since  $t \notin \text{ans}(q, \mathcal{B}'')$  has to hold whereas we have assumed that  $t \in \text{ans}(q, \mathcal{B})$  for each  $\mathcal{B} \in \mathbb{B}$ .  $\square$

Decidability and complexity results, under KDs and SFSKs only, follow from Theorem 1.

*Corollary 1*

Let  $\mathcal{G}$  be a global schema containing KDs and SFSKs only,  $\mathcal{D}$  be a possibly inconsistent database for  $\mathcal{G}$ ,  $q$  be a query,  $\Sigma \in \{\textit{loosely-exact}, \textit{loosely-sound}\}$ , and  $t$  be a tuple of values from  $\text{dom}(\mathcal{D})$ . The problem of establishing whether  $t \in \text{ans}_\Sigma(q, \mathcal{G}, \mathcal{D})$  is in  $\Pi_2^p$  in data complexity.

*Proof*

It suffices to prove that the problem of establishing whether  $t \notin \text{ans}_\Sigma(q, \mathcal{G}, \mathcal{D})$  is in  $\Sigma_2^p$ . This can be done by (i) building  $\mathcal{D}^*$  and (ii) guessing  $\mathcal{B} \in 2^{\mathcal{D}^*}$  such that  $\mathcal{B}$  is a  $\Sigma$ -repair and  $t \notin \text{ans}(q, \mathcal{B})$ . Since, by Proposition 1,  $|\mathcal{D}^*| \in \mathcal{O}(|\mathcal{D}|^\alpha)$ , where  $\alpha = \max\{\text{arity}(g) : g \in \text{names}(\mathcal{G})\}$ , then Step (i) (enumerate the facts of  $\mathcal{D}^*$ ) can be done in polynomial time. Since checking that  $t \notin \text{ans}(q, \mathcal{B})$  can be done in PTIME. It remains to show that checking whether  $\mathcal{B}$  is a  $\Sigma$ -repair can be done in coNP.

**[loosely-exact]** If  $\Sigma = \textit{loosely-exact}$ , this task corresponds to checking that there is no consistent  $\mathcal{B}' \subseteq \mathcal{D} \cup \mathcal{B}$  such that  $\mathcal{B}' \ominus \mathcal{D} \subset \mathcal{B} \ominus \mathcal{D}$ , where this last task is doable in PTIME.

**[loosely-sound]** If  $\Sigma = \textit{loosely-sound}$ , this task corresponds to checking that there is no consistent  $\mathcal{B}' \subseteq \mathcal{D}^*$  such that  $\mathcal{B}' \cap \mathcal{D} \supset \mathcal{B} \cap \mathcal{D}$ , where this last task is doable in PTIME.

Then the thesis follows.  $\square$

### 2.5 Equivalence of CQA under loosely-exact and CM-complete semantics

In this section, we define some relevant cases in which CQA under loosely-exact and CM-complete semantics coincide.

In the following, some of the proofs are not included due to space constraints; they can be found in the extended version of this paper (Manna et al. 2011).

*Lemma 5*

Given a database  $\mathcal{D}$  for a schema  $\mathcal{G}$ , if  $\mathcal{B}$  is a CM-complete repair for  $\mathcal{D}$ , then it is a loosely-exact repair for  $\mathcal{D}$ .

*Corollary 2*

$$\text{ans}_{\textit{loosely-exact}}(q, \mathcal{G}, \mathcal{D}) \subseteq \text{ans}_{\textit{CM-complete}}(q, \mathcal{G}, \mathcal{D})$$

*Theorem 2*

There are cases where  $\text{ans}_{\textit{loosely-exact}}(q, \mathcal{G}, \mathcal{D}) \subset \text{ans}_{\textit{CM-complete}}(q, \mathcal{G}, \mathcal{D})$

*Proof*

By Chomicki and Marcinkowski (2005), stating that the two semantics are different, and by Corollary 2.  $\square$

*Proposition 2*

Let  $\mathcal{B}$  be a database consistent w.r.t. a set of ICs  $C$ .

(1) If  $C$  are DCs only, then each  $\mathcal{B}' \subset \mathcal{B}$  is consistent w.r.t.  $C$  as well.

(2) If  $C$  are INDs only, then  $\mathcal{B} \cup \mathcal{B}'$  is consistent w.r.t.  $C$  for each  $\mathcal{B}'$  consistent w.r.t.  $C$ .

### Theorem 3

Given a database  $\mathcal{D}$  for a schema  $\mathcal{G}$ , let  $\mathcal{B}$  be a loosely-exact repair for  $\mathcal{D}$ , and  $\overline{\mathcal{B}} = \mathcal{B} \cap \mathcal{D}$ . There is a CM-complete repair  $\mathcal{B}' \subseteq \overline{\mathcal{B}}$  for  $\mathcal{D}$  if at least one of the following holds: (1)  $\mathcal{G}$  contains DCs only (no INDs); (2)  $\mathcal{G}$  contains INDs only (no DCs); (3)  $\mathcal{G}$  contains KDs and FKs only, and  $\mathcal{D}$  is consistent w.r.t. KDs; and (4)  $\mathcal{G}$  contains KDs and SFKs only.

### Corollary 3

$ans_{loosely-exact}(q, \mathcal{G}, \mathcal{D}) = ans_{CM-complete}(q, \mathcal{G}, \mathcal{D})$  in the cases where Theorem 3 holds.

### Proposition 3

In general, Theorem 3 does not hold in case  $\mathcal{G}$  contains SFSKs and KDs only.

## 3 Computation of CQA via ASP

In this section, we show how to exploit ASP (Gelfond and Lifschitz 1988, 1991) for efficiently computing consistent answers to user queries under different semantic assumptions. ASP is a powerful logic programming paradigm allowing (in its general form) for disjunction in rule heads (Minker 1982) and non-monotonic negation in rule bodies. In the following, we assume that the reader is familiar with ASP with aggregates, and in particular, we adopt the DLV syntax (Leone *et al.* 2006; Faber *et al.* 2010).

The suitability of ASP for implementing CQA has been already recognized in the literature (Lenzerini 2002; Arenas *et al.* 2003; Bertossi *et al.* 2005; Chomicki and Marcinkowski 2005). The general approaches are based on the following idea: produce an ASP program  $P$  whose answer sets represent possible repairs, so that the problem of computing CQA corresponds to cautious reasoning on  $P$ . One of the hardest challenges in this context is the automatic identification of a program  $P$  considering a minimal number of repairs actually relevant to answering user queries.

In order to face these challenges, we first introduce a general encoding that unifies in a common core the solutions for CQA under the semantics considered in this paper. Then, based on this unified framework, we define optimization strategies precisely aiming at reducing the computational cost of CQA. This is done in several ways: (i) by casting down the original program to complexity-wise easier programs; (ii) by identifying portions of the database not requiring repairs at all, according to the query requirements; and (iii) by exploiting equivalence classes between some semantics in such a way to adopt optimized solutions.

We next present the general encoding first and, then, the optimizations.

### 3.1 General encoding

The general approach generates a program  $\Pi_{cqa}$  and a new query  $q_{cqa}$  obtained by rewriting both the constraints and the query  $q$  in such a way that CQA reduces to

cautious reasoning on  $\Pi_{cqa}$  and  $q_{cqa}$ . Recall that a union of conjunctive queries in ASP is expressed as a set of rules having the same head predicate with the same arity.

In what follows, we first present how to generate  $\Pi_{cqa}$  and  $q_{cqa}$  and then formally prove under which hypothesis cautious reasoning on such  $\Pi_{cqa}$  and  $q_{cqa}$  corresponds to CQA.

Given a database  $\mathcal{D}$  for a schema  $\mathcal{G}$  and a query  $q$  on  $\mathcal{G}$ , the ASP program  $\Pi_{cqa}$  is created by rewriting each IC belonging to  $constr(\mathcal{G})$  and  $q$  as follows:

*Denial constraints.* Let  $\Sigma \in \{CM\text{-complete, loosely-sound, loosely-exact}\}$ . For each DC of the form  $\forall \bar{x}_1, \dots, \bar{x}_m \neg [g_1(\bar{x}_1) \wedge \dots \wedge g_m(\bar{x}_m) \wedge \sigma(\bar{x}_1, \dots, \bar{x}_m)]$  in  $constr(\mathcal{G})$ , insert the following rule into  $\Pi_{cqa}$ :

$$\bullet g_1^c(\bar{x}_1) \vee \dots \vee g_m^c(\bar{x}_m) :- g_1(\bar{x}_1), \dots, g_m(\bar{x}_m), \sigma(\bar{x}_1, \dots, \bar{x}_m).$$

This rule states that in the presence of a violated DC, it must be guessed the tuple(s) to be removed in order to repair the database.

*Inclusion dependencies.* Let  $\Sigma = \{CM\text{-complete, loosely-exact}\}$ . For each IND  $d$  in  $constr(\mathcal{G})$  of the form  $\forall \bar{x}_1 [g_1(\bar{x}_1) \rightarrow \exists \bar{x}_{2\exists} g_2(\bar{x}_2)]$ , add the following rules into  $\Pi_{cqa}$ :

$$\begin{aligned} \bullet g_1^c(\bar{x}_1) :- g_1(\bar{x}_1), \#count\{\bar{x}_{2\exists} : g_2^c(\bar{x}_2)\} = \#count\{\bar{x}_{2\exists} : g_2(\bar{x}_2)\}. & \quad \text{if } |\bar{x}_{2\exists}| > 0 \\ \bullet g_1^c(\bar{x}_1) :- g_1(\bar{x}_1), g_2^c(\bar{x}_2). & \\ \bullet g_1^c(\bar{x}_1) :- g_1(\bar{x}_1), not\ g_2(\bar{x}_2). & \quad \text{if } |\bar{x}_{2\exists}| = 0 \end{aligned}$$

The first rule states that a tuple of  $g_1$  must be deleted iff either all the tuples in  $g_2$  previously referred to by  $g_1$  via  $d$  have been deleted due to the repairing process, or there is no tuple in  $g_2$  referred to by  $g_1$  via  $d$ . (This is done by comparing the total count of tuples in  $g_2$  and  $g_2^c$ .) Observe that if there is a cyclic set of INDs, the set of rules generated by this rewriting would contain recursive aggregates. Their semantics is described in Faber et al. (2010). The latter two rules replace the first one in the special case of  $|\bar{x}_{2\exists}| = 0$ .

*Repaired relations.* Let  $\Sigma \in \{CM\text{-complete, loosely-sound, loosely-exact}\}$ . For each relation name  $g \in names(\mathcal{G})$ , insert the following rule into  $\Pi_{cqa}$ :

$$\bullet g^r(\bar{x}) :- g(\bar{x}), not\ g^c(\bar{x}).$$

*Query rewriting.* Build  $q_{cqa}(\bar{x})$  from  $q(\bar{x})$  as follows:

- (1) If  $\Sigma = loosely\text{-sound}$ , then apply onto  $q$  the perfect rewriting algorithm that deals with INDs described in Cali et al. (2003b).<sup>1</sup>
- (2) For each atom  $g(\bar{y})$  in  $q$ , replace  $g(\bar{y})$  by  $g^r(\bar{y})$ .

The perfect rewriting introduced in Cali et al. (2003b) is intuitively described next. Given a query  $q(\bar{x})$  and a set of INDs, the algorithm iteratively computes a new query  $Q$  as follows. The query  $Q$  is first initialized with  $q$ ; then at each iteration it carries out the following two steps. (1) For each conjunction  $cq'$  in  $Q$ , and for each pair of atoms  $g_1, g_2$  in  $cq'$  that unify (i.e., for which there exists a substitution

<sup>1</sup> Observe that when  $\Sigma = loosely\text{-sound}$ , INDs are not encoded into logic rules.

transforming  $g_1$  into  $g_2$ ,  $g_1$  and  $g_2$  are substituted by one single unifying atom. (2) For each conjunction  $cq'$  in  $Q$ , and for each *applicable* IND  $d$  of the form  $g_1 \rightarrow g$  such that  $g$  is in  $cq'$ , it adds to  $Q$  a new conjunction  $cq''$  obtained from  $cq'$  by interpreting  $d$  as a rewriting rule on  $g$ , applied from right to left. The algorithm stops when no further modifications are possible on  $Q$  with the two steps above.

The following theorems show how and when cautious reasoning on  $\Pi_{cqa}$  and  $q_{cqa}$  correspond to CQA. First, we consider the CM-complete semantics.

In the following, some of the proofs are not included due to space constraints; they can be found in the extended version of this paper (Manna *et al.* 2011).

#### Theorem 4

Let  $\Sigma = CM\text{-complete}$ , let  $\mathcal{D}$  be a database for a schema  $\mathcal{G}$  with arbitrary DCs and (possibly cyclic) INDs, and let  $q$  be a union of conjunctive queries. Then  $t \in ans_{\Sigma}(q, \mathcal{G}, \mathcal{D})$  iff  $q_{cqa}(t)$  is a cautious consequence of the ASP program  $\mathcal{D} \cup \Pi_{cqa}$ .

#### Example 3

Consider again Example 2, the program (and the query built from  $q(X) :- m(X)$ ) under the *CM-complete* semantics obtained for it is:

- $e^c(X_c, X_n) \vee e^c(X_c, X'_n) :- e(X_c, X_n), e(X_c, X'_n), X_n \neq X'_n$ .
- $m^c(X_c) :- m(X_c), \#count\{X'_n : e^c(X_c, X'_n)\} = \#count\{X_n : e(X_c, X_n)\}$ .
- $e^r(X_c, X_n) :- e(X_c, X_n), not\ e^c(X_c, X_n)$ .
- $m^r(X_c) :- m(X_c), not\ m^c(X_c)$ .
- $q_{cqa}(X_c) :- m^r(X_c)$ .

When this program is evaluated on the database, we obtain four answer sets. It can be verified that all the answer sets contain  $m^r('e1')$  and  $m^r('e2')$  (i.e., they are cautious consequences of  $\Pi_{cqa}$ ), and thus, 'e1' and 'e2' are the consistent answers to the query.

#### Theorem 5

Let  $\Sigma = loosely\text{-sound}$ , let  $\mathcal{D}$  be a database for a schema  $\mathcal{G}$  with KDs (and exactly one key for each relation) and (possibly cyclic) NKC INDs, and let  $q$  be a union of conjunctive queries without comparison atoms.<sup>2</sup> Then  $t \in ans_{\Sigma}(q, \mathcal{G}, \mathcal{D})$  iff  $q_{cqa}(t)$  is a cautious consequence of the ASP program  $\mathcal{D} \cup \Pi_{cqa}$ .

The general encoding for the *loosely-exact* semantics is inherently more complex than the ones for *loosely-sound* and *CM-complete*, since both tuple deletions and tuple insertions are subject to minimization. As a consequence, we tackled the *loosely-exact* encoding by considering that there are common cases in which CQA under the *loosely-exact* semantics and the *CM-complete* semantics actually coincide (see Corollary 3). These cases can be easily checked, and thus, it is possible to handle the *loosely-exact* semantics with the encoding defined for the *CM-complete* case.

<sup>2</sup> Recall that equalities are expressed in terms of variables having the same name.

*Theorem 6*

Let  $\Sigma = \textit{loosely-exact}$ , and  $\mathcal{D}$  be a database for a schema  $\mathcal{G}$  such that one of the following holds: (i)  $\mathcal{G}$  contains DCs only (no INDs); (ii)  $\mathcal{G}$  contains INDs only (no DCs); (iii)  $\mathcal{G}$  contains KDs and FKs only, and  $\mathcal{D}$  is consistent w.r.t. KDs; and (iv)  $\mathcal{G}$  contains KDs and SFKs only. Given a union of conjunctive queries  $q$ , then  $t \in \textit{ans}_{\Sigma}(q, \mathcal{G}, \mathcal{D})$  iff  $q_{cqa}(t)$  is a cautious consequence of the ASP program  $\mathcal{D} \cup \Pi_{cqa}$ .

*Proof*

Follows from Corollary 3 and Theorem 5.  $\square$

**3.2 Optimized solution**

The strategy reported in the previous section is a general solution for solving the CQA problem but, in several cases, more efficient ASP programs can be produced. First of all, note that the general algorithm blindly considers all the ICs on the global schema, including those that have no effect on the specific query. Consequently, useless logic rules might be produced which may slow down program evaluation. Then, a very simple optimization may consist of considering relevant ICs only. However, there are several cases in which the complexity of CQA stays in **PTIME**, but disjunctive programs, for which cautious reasoning becomes a hard task (Eiter *et al.* 1997), are generated even in the presence of DCs only. This means that the evaluation of the produced logic programs might be much more expensive than required in those “easy” cases. In the following, we provide semantic-specific optimizations aiming to overcome such problems for the settings pointed out in Theorems 4–6.

Given a query  $q$  and an atom  $g$  in  $q$ , we define the set of *relevant indices* of  $g$  in  $q$ , say  $\textit{relevant}(q, g)$ , in such a way that an index  $i$  in  $[1, \dots, \textit{arity}(g)]$  belongs to  $\textit{relevant}(q, g)$  if at least one of the following holds for an occurrence  $g(X_1, \dots, X_n)$  of  $g$  in  $q$ :

- $X_i$  is not existentially quantified (it is a free variable, it is an output variable of  $q$ );
- $X_i$  is involved in some comparison atom (even if it is existentially quantified);
- $X_i$  appears more than once in the same conjunction;
- $X_i$  is a constant value;

If  $g$  does not appear in  $q$ , we say that  $\textit{relevant}(q, g) = \emptyset$ .

In the following, we denote by  $\pi$  a set of indices. Moreover, given a sequence of variables  $\bar{x}$  and a set  $\pi \subseteq \{1, \dots, |\bar{x}|\}$ , we denote by  $\bar{x}^{\pi}$  the sequence obtained from  $\bar{x}$  by discarding a variable if its position is not in  $\pi$ . Finally, given a relation name  $g$ , a set of indices  $\pi$ , and a label  $\ell$ , we denote by  $g^{\ell-\pi}(\bar{x}^{\pi})$  an auxiliary atom derived from  $g$ , marked by  $\ell$ , and using only variables in  $\bar{x}^{\pi}$ .

$\Sigma = \textit{loosely-sound}$ . The objective of this optimization is to single out, for each relation involved by the query, the set of attributes actually relevant to answer it and apply the necessary repairs only on them. As we show next, this may allow both to reduce (even to zero) the number of disjunctive rules needed to repair key violations and to reduce the cardinality of relations involved in such disjunctions.

Given a schema  $\mathcal{G}$  and a query  $q$ , perform the following steps for building the program  $\Pi_{cqa}$  and the query  $Q_{cqa}$ .

- (1) Apply the perfect rewriting algorithm that deals with INDs described in Cali *et al.* (2003b).
- (2) Let  $Q$  be the union of conjunctive queries obtained from  $q$  after Step 1. For each  $g \in \text{names}(\mathcal{G})$ , build the sets

$$\pi_R^g = \text{relevant}(Q, g) \quad \pi_S^g = \pi_R^g \cup \text{key}(g)$$

These two sets capture the fact that a key attribute is relevant for the repairing process, but it may not be strictly relevant for answering the query.

Observe that the perfect rewriting dealing with INDs must be applied *before* singling out relevant attributes. In fact,  $q$  may also depend, through INDs, on attributes of relations not explicitly mentioned in it. However, in the last step of this algorithm, the rewriting of the query is completed by substituting each relation in the query with its repaired (and possibly reduced) version.

- (3) For each  $g \in \text{names}(\mathcal{G})$  such that  $\pi_R^g \neq \emptyset$  and  $\text{key}(g) \not\subseteq \pi_R^g$ , add the following rules into  $\Pi_{cqa}$ :

- $g^{sr-\pi_S^g}(\bar{x}^{\pi_S^g}) :- g(\bar{x})$ .
  - $g^{c-\pi_S^g}(\bar{x}_1^{\pi_S^g}) \vee g^{c-\pi_S^g}(\bar{x}_2^{\pi_S^g}) :- g^{sr-\pi_S^g}(\bar{x}_1^{\pi_S^g}), g^{sr-\pi_S^g}(\bar{x}_2^{\pi_S^g}), \bar{x}_1^i \neq \bar{x}_2^i$ .
  - $g^{r-\pi_R^g}(\bar{x}^{\pi_R^g}) :- g^{sr-\pi_S^g}(\bar{x}^{\pi_S^g}), \text{not } g^{c-\pi_S^g}(\bar{x}^{\pi_S^g})$ .
- $\forall i \in \pi_S^g - \text{key}(g)$

Observe that if there exists at least one relevant non-key attribute for  $g$ , the repairing process can not be avoided; however, violations caused by irrelevant attributes only (i.e., not in  $\pi_S^g$ ) can be ignored, since the projection of  $g$  on  $\pi_S^g$  is still safe and sufficient for query answering purposes.

- (4) For each  $g \in \text{names}(\mathcal{G})$  such that  $\pi_R^g \neq \emptyset$  and  $\text{key}(g) \supseteq \pi_R^g$ , add the following rule into  $\Pi_{cqa}$ :

- $g^{r-\pi_R^g}(\bar{x}^{\pi_R^g}) :- g(\bar{x})$ .

Observe that if the relevant attributes of  $g$  are a subset of its key, the repair process of  $g$  for key violations through disjunction can be avoided at all. In fact, the projection of  $g$  on  $\pi_R^g$  is still safe and sufficient for query answering purposes. Moreover, for the same reason, it is not needed to take all the key of  $g$  into account.

- (5) For each atom of the form  $g(\bar{x})$  in  $Q$ , replace  $g(\bar{x})$  by  $g^{r-\pi_R^g}(\bar{x}^{\pi_R^g})$ .

$\Sigma = \text{CM-complete}$ . For the optimization of the CM-complete semantics, we exploit a graph which is used to navigate the query and the database in order to single out those relations and projections actually relevant for answering the query. Moreover, it allows to identify possible cycles generated by ICs which must be suitably handled; in fact, acyclic ICs induce a partial order among them and this information can be effectively exploited for the optimization. On the contrary, cyclic ICs must be handled in a more standard way.



Given a schema  $\mathcal{G}$  and a query  $q$ , build the directed labeled graph  $G_q = \langle N, A \rangle$  as follows:  $N = \{q\} \cup \text{names}(\mathcal{G})$ ;  $(g_1, g_2, c) \in A$  iff  $c$  is a DC in  $\text{constr}(\mathcal{G})$  involving both  $g_1$  and  $g_2$ ;  $(g_1, g_2, d) \in A$  iff  $d$  is an IND in  $\text{constr}(\mathcal{G})$  of the form  $g_1 \rightarrow g_2$ ; and  $(q, g, \varepsilon) \in A$  iff  $g$  appears in a conjunction of  $q$ . Perform the following steps for building the program  $\Pi_{cqa}$ :

- (1) Visit  $G_q$  starting from node  $q$ .
- (2) Discard unreachable nodes and update the sets  $N$  and  $A$ .
- (3) Partition the set  $N$  in  $(N_{cf}, N_{ncf})$  in such a way that a node  $n$  belongs to  $N_{cf}$  if it is not involved in any cycle ( $q$  always belongs to  $N_{cf}$ ). Contrariwise, a node  $n$  belongs to  $N_{ncf}$  if it is involved in some cycle.
- (4) For each node  $g \in N - \{q\}$ , compute the sets

$$\begin{aligned} \pi_R^g &= (\bigcup_{(g_L, g, d) \in A} \pi_R^d) \cup \text{relevant}(q, g); \\ \pi_S^g &= \pi_R^g \cup \text{key}(g), \text{ only if } g \text{ has exactly one primary key as DCs; } \pi_S^g = \emptyset \\ &\text{otherwise.} \end{aligned}$$

Here,  $\pi_R^g$  is the set of relevant variable indices of  $g$ , and  $\pi_S^g$  adds to  $\pi_R^g$ , the key of  $g$ .

Observe that Steps 1–4 implement a pre-processing phase in which relevant relations and their relevant indices are singled out, and each relevant relation is classified as cycle free or non-cycle free.

- (5) For each node  $g \in N_{cf}$ , if  $g$  has only one key as DCs, then add the following rules into  $\Pi_{cqa}$ :

$$\begin{aligned} &\bullet g^{\zeta, \pi_\chi^g}(\bar{x}^{\pi_\chi^g}) :- g(\bar{x}), g_1^{r, \pi_R^{d_1}}(\bar{x}_1^{\pi_R^{d_1}}), \dots, g_k^{r, \pi_R^{d_k}}(\bar{x}_k^{\pi_R^{d_k}}). \\ &\bullet g_i^{r, \pi_R^{d_i}}(\bar{x}_i^{\pi_R^{d_i}}) :- g_i^{r, \pi_R^{d_i}}(\bar{x}_i^{\pi_R^{d_i}}). \end{aligned} \quad \forall i \in [1, \dots, k] \text{ s.t. } \pi_R^{g_i} \supset \pi_R^{d_i}$$

where:

- (i)  $k \geq 0$  is the number of arcs in  $G_q$ , labeled by INDs, and outgoing from  $g$ ;
- (ii) the pair  $(\zeta, \chi)$  is either  $(r, R)$  or  $(sr, S)$ , according to whether  $\text{key}(g) \supseteq \pi_R^g$  or not, respectively. Intuitively, if  $\text{key}(g) \supseteq \pi_R^g$  holds, then the repair  $g^{r, \pi_R^g}$  of  $g$  can be directly computed; otherwise, the computation must first go through a semi-reparation step for computing  $g^{sr, \pi_S^g}$ . Intuitively, this semi-reparation step collects those tuples that violate no IND of the form  $g \rightarrow g_i$ , but that must be anyway processed in order to fix some key violation (see Steps 6–10).
- (iii) atom  $g_i^{r, \pi_R^{d_i}}$  is in the body of the first rule ( $1 \leq i \leq k$ ) only if both  $(g, g_i, d_i) \in A$  and  $d_i$  is an IND of the form  $g(\bar{x}) \rightarrow g_i(\bar{x}_i)$ . This atom is just a projection of  $g_i^{r, \pi_R^{d_i}}(\bar{x}_i^{\pi_R^{d_i}})$ .

- (6) For each node  $g \in N_{cf}$ , if  $g$  has only one primary key as DCs, and  $\text{key}(g) \subset \pi_R^g$ , and  $g$  has incoming arcs only from  $q$ , and all the relevant variables of  $g$  w.r.t.  $q$  are in the head of  $q$ , and each occurrence of  $g$  in  $q$  contains all of its relevant variables, then add the following rules into  $\Pi_{cqa}$  by considering that the key of  $g$  is defined by rules of the form  $\forall \bar{x}_1, \bar{x}_2 \neg [g(\bar{x}_1) \wedge g(\bar{x}_2) \wedge \bar{x}_1^i \neq \bar{x}_2^i]$ :

$$\bullet g^{c, \pi_S^g}(\bar{x}_1^{\pi_S^g}) :- g^{sr, \pi_S^g}(\bar{x}_1^{\pi_S^g}), g^{sr, \pi_S^g}(\bar{x}_2^{\pi_S^g}), \bar{x}_1^i \neq \bar{x}_2^i. \quad \forall i \in \pi_S^g - \text{key}(g)$$

- $g^{r-\pi_R^g}(\bar{x}_1^{\pi_R^g}) :- g^{sr-\pi_S^g}(\bar{x}_1^{\pi_S^g}), \text{ not } g^{c-\pi_S^g}(\bar{x}_1^{\pi_S^g}).$
- (7) For each node  $g \in N_{cf}$ , if  $g$  has only one primary key as DCs, and  $key(g) \not\supseteq \pi_R^g$ , and Case 6 does not apply, then add the following rules into  $\Pi_{cqa}$  by considering that the key is defined by rules of the form  $\forall \bar{x}_1, \bar{x}_2 \neg [g(\bar{x}_1) \wedge g(\bar{x}_2) \wedge \bar{x}_1^i \neq \bar{x}_2^i]$ :
- $g^{c-\pi_S^g}(\bar{x}_1^{\pi_S^g}) \vee g^{c-\pi_S^g}(\bar{x}_2^{\pi_S^g}) :- g^{sr-\pi_S^g}(\bar{x}_1^{\pi_S^g}), g^{sr-\pi_S^g}(\bar{x}_2^{\pi_S^g}), \bar{x}_1^i \neq \bar{x}_2^i.$
  - $g^{r-\pi_R^g}(\bar{x}_1^{\pi_R^g}) :- g^{sr-\pi_S^g}(\bar{x}_1^{\pi_S^g}), \text{ not } g^{c-\pi_S^g}(\bar{x}_1^{\pi_S^g}).$
- $\forall i \in \pi_S^g - key(g)$

Observe that, in this case, disjunctive rules are defined only on the set of relevant indices that are not in the key and that each  $g^{c-\pi_S^g}$  contains only the projection of deleted tuples on the set  $\pi_S^g$ .

Here, Steps 5–7 handle relations for which a key is defined and are classified as cycle free. In particular, if  $key(g) \supseteq \pi_R^g$  holds, key reparation can be avoided at all (and thus disjunctive rules too); otherwise, a semi-reparation step is required, but Step 6 identifies further cases in which even if key reparation is needed, disjunction can be still avoided. Finally, Step 7 handles all the other cases. Importantly, through Steps 5–7, we take into account only the minimal projections of involved relations in order to reduce as much as possible computational costs (and even disjunctive rules) not considering irrelevant attributes.

- (8) For each node  $g \in N_{ncf}$ , add the following rules into  $\Pi_{cqa}$ :
- $g^c(\bar{x}) :- g(\bar{x}), \text{ not } g_1^{r-\pi_R^d}(\bar{x}_1^{\pi_R^d}).$   
 $g_1^{r-\pi_R^d}(\bar{x}_1^{\pi_R^d}) :- g_1^{r-\pi_R^g}(\bar{x}_1^{\pi_R^g}).$   
 for each IND  $d$  of the form  $g(\bar{x}) \rightarrow g_1(\bar{x}_1)$  such that there is no cycle in  $G_g$  involving both  $g_1$  and  $g$ ;
  - $g^c(\bar{x}) :- g(\bar{x}), \#count\{\bar{x}_{1\exists} : g_1^c(\bar{x}_1)\} = \#count\{\bar{x}_{1\exists} : g_1(\bar{x}_1)\}.$   
 for each IND  $d$  of the form  $\forall \bar{x}_\forall [g(\bar{x}) \rightarrow \exists \bar{x}_{2\exists} g_1(\bar{x}_1)]$  such that  $g_1 \in N_{ncf}$ ;
  - $g^c(\bar{x}_1) \vee g^c(\bar{x}_2) :- g(\bar{x}_1), g(\bar{x}_2), \bar{x}_1^i \neq \bar{x}_2^i.$   $\forall i \in \pi$   
 where  $\pi = \{1, \dots, arity(g)\} - key(g)$  and the key of  $g$  is defined by DCs of the form  $\forall \bar{x}_1, \bar{x}_2 \neg [g(\bar{x}_1) \wedge g(\bar{x}_2) \wedge \bar{x}_1^i \neq \bar{x}_2^i]$ ;
  - $g^{r-\pi_R^g}(\bar{x}_1^{\pi_R^g}) :- g(\bar{x}), \text{ not } g^c(\bar{x}).$   
 if there is at least one node in  $N_{cf}$  with an arc to  $g$ , or  $g$  appears in  $q$ ;

- (9) For each DC of the form  $\forall \bar{x}_1, \dots, \bar{x}_m \neg [g_1(\bar{x}_1) \wedge \dots \wedge g_m(\bar{x}_m) \wedge \sigma(\bar{x}_1, \dots, \bar{x}_m)]$  involving at least two different relation names (entailing that each  $g_i \in N_{ncf}$ ), add the following rules into  $\Pi_{cqa}$ :

- $g_1^c(\bar{x}_1) \vee \dots \vee g_m^c(\bar{x}_m) :- g_1(\bar{x}_1), \dots, g_m(\bar{x}_m), \sigma(\bar{x}_1, \dots, \bar{x}_m).$

Steps 8 and 9 handle non-cycle free relations; the repairing process in this case mimics the standard rewriting, but projects relations on the relevant attributes whenever possible.

- (10) For each node  $g \in N_{cf}$ , if  $g$  is involved in DCs that do not form a primary key, then add the following rules into  $\Pi_{cqa}$ :
- $g^{sr}(\bar{x}) :- g(\bar{x}), g_1^{r-\pi_R^{d_1}}(\bar{x}_1^{\pi_R^{d_1}}), \dots, g_k^{r-\pi_R^{d_k}}(\bar{x}_k^{\pi_R^{d_k}}).$

- $g_i^{r-\pi_R^{d_i}}(\bar{x}_i^{\pi_R^{d_i}}) :- g_i^{r-\pi_R^{g_i}}(\bar{x}_i^{\pi_R^{g_i}}). \quad \forall i \in [1, \dots, k] \text{ s.t. } \pi_R^{g_i} \supset \pi_R^{d_i}$
- $g^c(\bar{x}_1) \vee \dots \vee g^c(\bar{x}_m) :- g^{sr}(\bar{x}_1), \dots, g^{sr}(\bar{x}_m), \sigma_d(\bar{x}_1, \dots, \bar{x}_m). \quad \forall d$
- $g^{r-\pi_R^g}(\bar{x}^{\pi_R^g}) :- g^{sr}(\bar{x}), \text{ not } g^c(\bar{x}).$

where:

- (i)  $k \geq 0$  is the number of arcs, labeled by INDs, and outgoing from  $g$ ;
- (ii) atom  $g_i^{r-\pi_R^{d_i}}$  is in the body of the first rule ( $1 \leq i \leq k$ ) iff both  $(g, g_i, d_i) \in A$  and  $d_i$  is an IND of the form  $g(\bar{x}) \rightarrow g_i(\bar{x}_i)$ ;
- (iii)  $d$  is a DC of the form  $\forall \bar{x}_1, \dots, \bar{x}_m \neg [g(\bar{x}_1) \wedge \dots \wedge g(\bar{x}_m) \wedge \sigma_d(\bar{x}_1, \dots, \bar{x}_m)]$ .

Step 10 handles the special case in which there is no key for a relation but DCs are defined (only) on it.

(11) For each atom of the form  $g(\bar{x})$  in  $q$ , replace  $g(\bar{x})$  by  $g^{r-\pi_R^g}(\bar{x}^{\pi_R^g})$ .

$\Sigma = \textit{loosely-exact}$ . In Section 3.1, we proved that there are common cases in which CQA under the *loosely-exact* semantics and the *CM-complete* semantics actually coincide. As a consequence, in these cases, all the optimizations defined for the *CM-complete* semantics apply also to the *loosely-exact* semantics.

### 4 Experiments

In this section, we present some of the experiments that we carried out to assess the effectiveness of our approach to CQA.

Testing has been performed by exploiting our complete system for data integration, which is intended to simplify both the integration system design and the querying activities by exploiting a user-friendly GUI. Indeed, this system both supports the user in designing the global schema and the mappings between global relations and source schemas, and allows to specify user queries over the global schema via a QBE-like interface. The query evaluation engine adopted for the tests is DLV<sup>DB</sup> (Terracina, De Francesco, et al. 2008) coupled, via ODBC, with a PostgreSQL DBMS where input data were stored. DLV<sup>DB</sup> is a DLP evaluator born as a database oriented extension of the well-known DLV system (Leone et al. 2006). It has been recently extended for dealing with unstratified negation, disjunction, and external function calls.

We first address tests on a real-world scenario and then report on tests for scalability issues on synthetic data.

#### 4.1 Tests on a real-world scenario

*Data set.* We have exploited the real-world data integration framework developed in the INFOMIX project (IST-2001-33570) (Leone et al. 2005), which integrates data from a real university context. In particular, considered data sources were available at the University of Rome “La Sapienza”. These comprise information on students, professors, curricula, and exams in various faculties of the university.

There are about 35 data sources in the application scenario, which are mapped into 12 global schema relations with 20 GAV mappings and 21 ICs. We call this data set **Infomix** in the following. (A detailed description of the **Infomix** framework is given in the extended version (Manna et al. 2011) of this paper.)

Besides the original source database instance (which takes about 16 MB on DBMS), we obtained bigger instances artificially. Specifically, we generated a number of copies of the original database; each copy is disjoint from the other ones but maintains the same data correlations between instances as the original database. This has been carried out by mapping each original attribute value to a new value having a copy-specific prefix.

Then, we considered two further data sets, namely, **Infomix-x-10** and **Infomix-x-50**, storing 10 copies (for a total amount of 160 MB of data) and 50 copies (800 MB) of the original database, respectively. It holds that  $\text{Infomix} \subset \text{Infomix-x-10} \subset \text{Infomix-x-50}$ .

*Compared methods and tested queries.* In order to assess the characteristics of the proposed optimizations, we measured the execution time of different queries with (i) the standard encoding (identified as **STD** in the following), (ii) a naïve optimization obtained by only removing relations not strictly needed for answering the queries (**OPT1** in the following), and (iii) the fully optimized encoding presented in Section 3 (**OPT2** in the following). Each of these cases has been evaluated for the three semantics considered in this paper. In order to isolate the impact of our optimizations, we disabled other optimizations (such as magic sets) embedded in the datalog evaluation engine. Clearly, such optimizations are complementary to our own and might further improve the overall performances.

We have considered six queries, fully specified in the extended version (Manna *et al.* 2011) of this paper, named Q1,..., Q6. In particular, Q2 involves key constraints only, while Q1 and Q3 involve both keys and acyclic INDs; specifically, Q3 involves an SFK while Q1 involves NKC INDs. Finally, Q4, Q5, and Q6 involve keys and cyclic NKC INDs.

*Results and discussion.* All tests have been carried out on an Intel Xeon X3430, 2.4 GHz, with 4 GB RAM, running Linux operating system. We set a time limit of 120 minutes after which query execution has been killed. Figures 1 and 2 show the obtained results for the *loosely-sound* and the *CM-complete* semantics. It is worth recalling that as we pointed out in Section 3.2, optimizations for the *loosely-exact* semantics are inherent to the equivalence classes to the *CM-complete* semantics discovered in this paper. As a consequence, we tested this semantics only on queries Q2 and Q3 for which such equivalence holds. Then, since the execution times of the optimized encoding coincide with the *CM-complete* graphs for queries Q2 and Q3, we do not report specific figures for them.

Analyzing the figures, we observe that the proposed optimizations do not introduce computational overhead and, in most cases, transform practically untractable queries in tractable ones; in fact, for all the tested queries, the execution time of the standard rewriting exceeded the time limit. **OPT1** helps mostly on the smallest data set; in fact, for **Infomix-x-10** it shows some gain in 33% of cases and only in two cases for **Infomix-x-50**.

As for the comparison among the optimized encodings, we can observe that if INDs are not involved by the query (Q2), the *loosely-sound* and the *CM-complete* optimizations have the same performances; this confirms theoretical expectations.

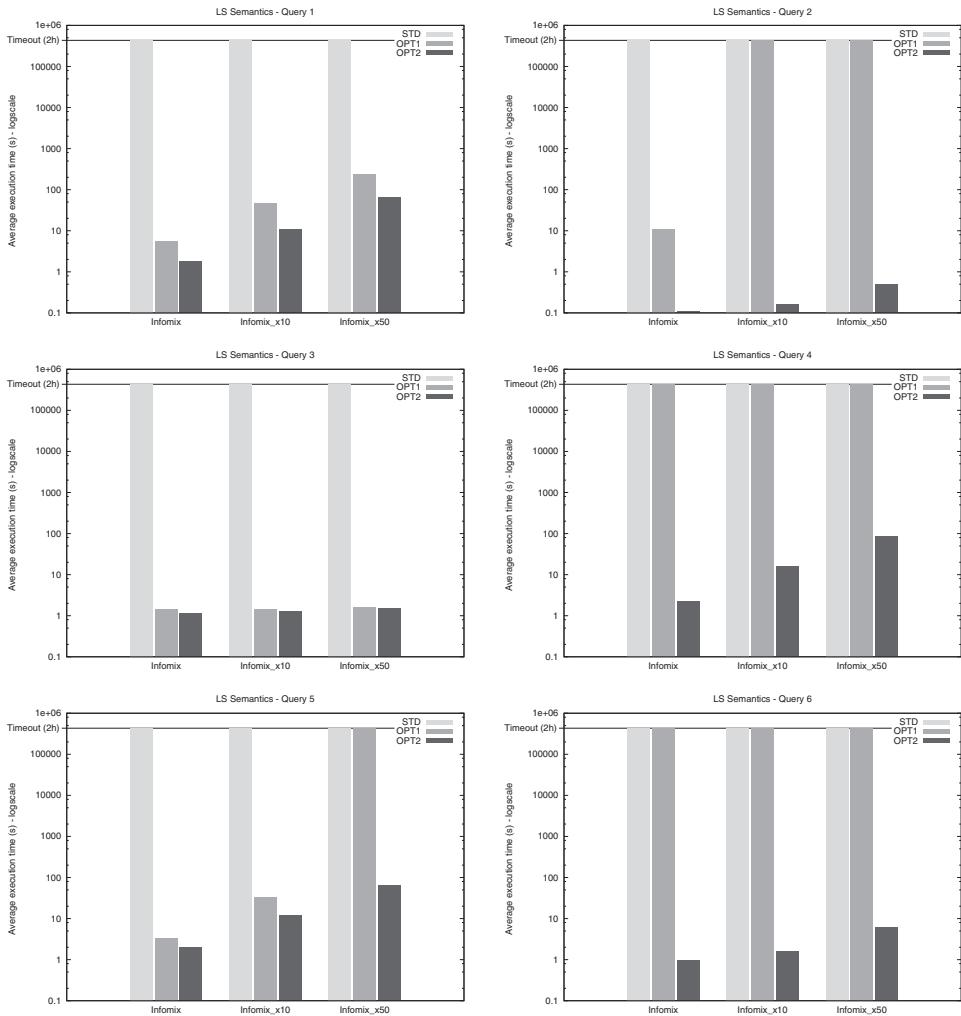


Fig. 1. Query evaluation execution times for the *loosely-sound* semantics.

When acyclic INDs are involved (Q1, Q3), the *loosely-sound* optimization performs slightly better because the *CM-complete* must choose the tuples to be deleted due to IND violations, whereas the *loosely-sound* semantics just works on the original data. Finally, when the involved INDs are cyclic (Q4, Q5, Q6), the performance of the *CM-complete* optimization further degrades w.r.t. the *loosely-sound* one because recursive aggregates must be exploited to choose deletions, and, thus, increases the complexity of query evaluation.

#### 4.2 Scalability analysis w.r.t. the number and kind of constraint violations

Since in the real-world scenario it emerged that the *CM-complete* semantics is more affected than the *loosely sound* one from the kind of involved constraints, we carried out a scalability analysis on this semantics, whose results are reported next.

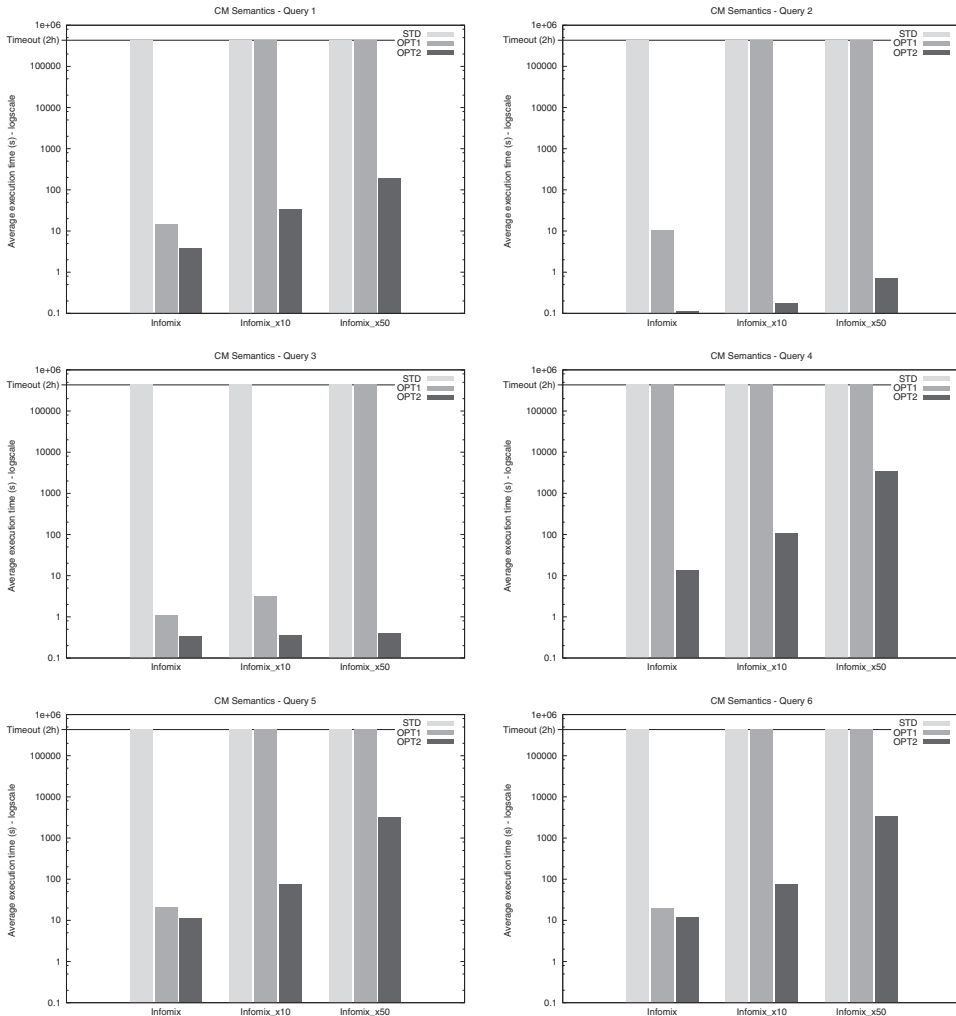


Fig. 2. Query evaluation execution times for the *CM-complete* semantics.

We considered a synthetic data set composed of three relations, named  $r_1, r_2,$  and  $r_3,$  over which we imposed different sets of ICs in order to analyze the scalability of our methods depending on the presence of keys and/or in the presence/absence of acyclic and cyclic INDs. In particular, we imposed the following key constraints:  $key(r_2) = \{1, 2\}, key(r_3) = \{1\},$  and we experimented with three different sets of INDs:  $NOINCL = \emptyset, ACYCLIC = \{r_1(X_1, X_2, X_3, X_4) \rightarrow r_2(X_2, X_5, X_3, X_6), r_1(X_1, X_2, X_3, X_4) \rightarrow r_3(X_1, X_5, X_6, X_7)\},$  and  $CYCLIC = ACYCLIC \cup \{r_2(X_1, X_2, X_3, X_4) \rightarrow r_1(X_5, X_6, X_7, X_2)\}.$  The employed query is:  $query(X1, X3) :- r_1(X1, X2, X3, X4), r_2(X2, X3, X5, X6)?$  We have randomly generated synthetic databases having a growing number of key violations on table  $r_2.$  The generation process progressively adds key violations to  $r_2$  by generating pairs of conflicting tuples; after an instance of  $r_2$  is obtained, tables  $r_1$  and  $r_3$  are generated by taking values from  $r_2$  in such a

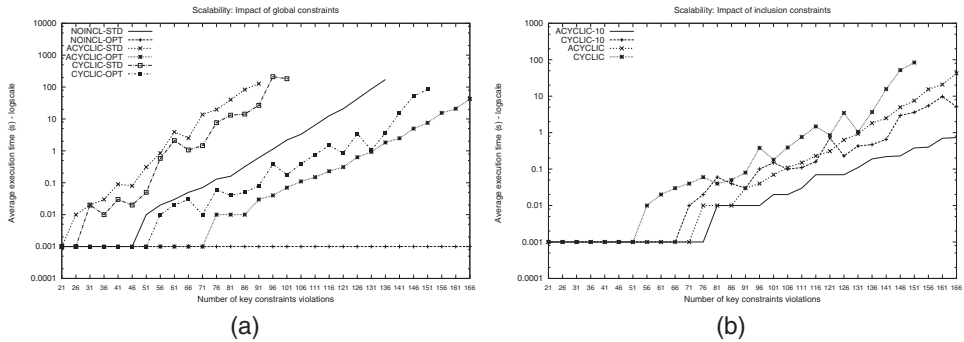


Fig. 3. Scalability analysis.

way that INDs are satisfied. In addition, for each tuple of  $r_3$ , a key-conflicting tuple is generated. In order to assess the impact of the number of INDs violations, for each database instance  $DB_x$ , containing  $x$  key violations on table  $r_2$ , we generated a  $DB_x-10$  instance where the 10% of tuples is (randomly) removed from tables  $r_1$  and  $r_3$  (causing INDs violations). We have generated six database instances per size (number of key violations on table  $r_2$ ), and plotted the time (averaged over the instances of the same size) in Figure 3.

In detail, Figure 3(a) shows the results for incrementally higher KD violations with no IND violations. Both standard and optimized encodings have been tested. Figure 3(b) compares the optimized encoding only, when the percentage of IND violations is 0% or 10%. Observe that, in general, even when there is no initial IND violation, the KD repairing process may induce some of them.

The analysis of these figures shows that even if cyclic INDs are generally harder, their scaling is almost the same as the acyclic ones. On the contrary, in the absence of INDs, the optimization may boost the performances (see the flat line in Fig. 3(a)). Figure 3(b) points out that when the number of IND violations increases, the performance may improve. This behavior is justified by the fact that tuple deletions due to IND repairs may, in their turn, remove KD violations. This reduces the number of disjunctions to be evaluated.

## 5 Related work and concluding remarks

From the 1990s – when the founding notions of *CQA* (Bry 1997), *GAV mapping* (Garcia-Molina et al. 1997; Tomasic et al. 1998; Goh et al. 1999), and *database repair* (Arenas et al. 1999) were introduced – *data integration* (Lenzerini 2002) and *inconsistent databases* (Bertossi et al. 2005) have been studied quite in depth.

Detailed characterizations of the main problems arising in a data integration system have been provided, taking into account different semantics, constraints, and query types (Arenas et al. 2003; Cali et al. 2003a, 2003b; Chomicki and Marcinkowski 2005; Grieco et al. 2005; Fuxman and Miller 2007; Eiter et al. 2008).



This paper provides a contribution in this scenario by extending the decidability boundaries for the *loosely-exact* semantics (as called in Cali *et al.* 2003a but firstly introduced by Arenas *et al.* 1999) and the *loosely-sound* semantics, in the case of both KDs and SFSK INDs.

A first proposal of an unifying framework for CQA in a data integration setting is presented in Cali *et al.* (2005) using first-order logic; it considers different semantics defined by interpreting the mapping assertions between the global and the local schemas of the data integration system. A common framework for computing repairs in a single database setting is proposed in Eiter *et al.* (2008); it covers a wide range of semantics relying on the general notion of pre-order for candidate repairs, but only universally quantified constraints are allowed. Moreover, the authors introduce an abstract logic programming framework to compute consistent answers. Finally, the authors propose an optimization strategy called factorization that, as will be clarified below, is orthogonal to our own.

This paper provides a contribution in this setting since it unifies different semantics, as in Cali *et al.* (2005) and Eiter *et al.* (2008), but also provides an algorithm that, given a retrieved database, a user query  $q$ , and a semantics, automatically composes an ASP program capable of computing the consistent answers to  $q$ . In particular, our ASP rewriting offers a natural, compact, and direct way for encoding even hard cases where the CQA problem belongs to the  $\Pi_2^p$  complexity class.

Theoretical studies gave rise to concrete implementations, most of which were conceived to operate on some specific semantics and/or constraint types (Arenas *et al.* 1999, 2003; Greco and Zumpano 2000; Greco *et al.* 2001; Cali *et al.* 2002, 2003b, 2004; Chomicki *et al.* 2004a, 2004b; Lembo 2004; Fuxman *et al.* 2005; Grieco *et al.* 2005; Leone *et al.* 2005; Fuxman and Miller 2007). As an example, in Leone *et al.* (2005), only the *loosely-sound* semantics was supported. In this paper, we provide both a unified framework based on ASP and a complete system supporting (i) all the three aforementioned significant semantics in the case of conjunctive queries and the most commonly used database constraints (KDs and INDs), (ii) specialized optimizations, and (iii) a user-friendly GUI.

Another general contribution of our work comes from a novel optimization technique that, after analyzing the query and localizing a minimal number of relevant ICs, tries to “simplify” their structure to reduce the number of database repairs – as they could be exponentially many (Arenas *et al.* 2001). Such technique could be classified as “vertical” due to the fact that it reduces (whenever possible) the arity of each active relation (with the effect, e.g., of decreasing the number of key conflicts) without looking at the data. It is orthogonal to other “horizontal” approaches, such as magic sets (Faber *et al.* 2007) and factorization (Eiter *et al.* 2008), which are based on data filtering strategies. In particular, a system exploiting ASP incorporating magic set techniques for CQA is described in Marileo and Bertossi (2010). Other approaches complementary to our own are based on first-order rewritings of the query (Arenas *et al.* 1999; Chomicki and Marcinkowski 2002; Cali *et al.* 2003b; Grieco *et al.* 2005; Fuxman and Miller 2007).

The combination of our optimizations with such approaches and further extensions of decidability boundaries for CQA are some of our future line of research.

### Acknowledgement

This work has been partially supported by the Calabrian Region under PIA (Pacchetti Integrati di Agevolazione industria, artigianato e servizi) project DLVSYS-TEM approved in BURC n. 20 parte III del 15/05/2009 - DR n. 7373 del 06/05/2009.

### References

- ABITEBOUL, S., HULL, R. AND VIANU, V. 1995. *Foundations of Databases: The Logical Level*. Addison-Wesley Longman Publishing Co., Boston, MA.
- ARENAS, M., BERTOSSI, L. AND CHOMICKI, J. 1999. Consistent query answers in inconsistent databases. In *Proceedings of PODS'99*. ACM, New York, 68–79.
- ARENAS, M., BERTOSSI, L. AND CHOMICKI, J. 2001. Scalar aggregation in FD-inconsistent databases. In *Proceedings of ICDT'01*, J. V. den Bussche and V. Vianu, Eds. Lecture Notes in Computer Science, vol. 1973. Springer, Berlin/ Heidelberg, 39–53.
- ARENAS, M., BERTOSSI, L. AND CHOMICKI, J. 2003. Answer sets for consistent query answering in inconsistent databases. *Theory and Practice of Logic Programming* 3, 4, 393–424.
- BERTOSSI, L. E., HUNTER, A. AND SCHAUB, T., Eds. 2005. *Inconsistency Tolerance*. Lecture Notes in Computer Science, vol. 3300. Springer, Berlin/Heidelberg.
- BRY, F. 1997. Query answering in information systems with integrity constraints. In *Proceedings of IICIS'97*, S. Jajodia, W. List, G. W. McGregor and L. Strous, Eds. Chapman & Hall, London, 113–130.
- CALÌ, A., CALVANESE, D., DE GIACOMO, G. AND LENZERINI, M. 2002. On the role of integrity constraints in data integration. *IEEE Data Engineering Bulletin* 25, 3, 39–45.
- CALÌ, A., CALVANESE, D., DE GIACOMO, G. AND LENZERINI, M. 2004. Data integration under integrity constraints. *Information Systems* 29, 2, 147–163.
- CALÌ, A., LEMBO, D. AND ROSATI, R. 2003a. On the decidability and complexity of query answering over inconsistent and incomplete databases. In *Proceedings of PODS'03*. ACM, New York, 260–271.
- CALÌ, A., LEMBO, D. AND ROSATI, R. 2003b. Query rewriting and answering under constraints in data integration systems. In *Proceedings of IJCAI'03*, G. Gottlob and T. Walsh, Eds. Morgan Kaufmann Publishers, San Francisco, CA, 16–21.
- CALÌ, A., LEMBO, D. AND ROSATI, R. 2005. A comprehensive semantic framework for data integration systems. *Journal of Algorithms* 3, 2, 308–328.
- CHOMICKI, J. AND MARCINKOWSKI, J. 2002. On the computational complexity of consistent query answers. *CoRR cs.DB/0204010*, 1–9.
- CHOMICKI, J. AND MARCINKOWSKI, J. 2005. Minimal-change integrity maintenance using tuple deletions. *Information and Computation* 197, 1–2, 90–121.
- CHOMICKI, J., MARCINKOWSKI, J. AND STAWORKO, S. 2004a. Computing consistent query answers using conflict hypergraphs. In *Proceedings of CIKM'04*, D. A. Grossman, L. Gravano, C. X. Zhai, O. Herzog and D. A. Evans, Eds. ACM, New York, 417–426.
- CHOMICKI, J., MARCINKOWSKI, J. AND STAWORKO, S. 2004b. Hippo: A system for computing consistent answers to a class of SQL queries. In *Advances in Database Technology – EDBT 2004*, E. Bertino, S. Christodoulakis, D. Plexousakis, V. Christophides, M. Koubarakis, K. Böhm and E. Ferrari, Eds. Lecture Notes in Computer Science, vol. 2992. Springer, Berlin/Heidelberg, 661–662.
- EITER, T., FINK, M., GRECO, G. AND LEMBO, D. 2008. Repair localization for query answering from inconsistent databases. *ACM Transactions on Database Systems* 33, 2, 10:1–10:51.

- EITER, T., GOTTLÖB, G. AND MANNILA, H. 1997. Disjunctive datalog. *ACM Transactions on Database Systems* 22, 3, 364–418.
- FABER, W., GRECO, G. AND LEONE, N. 2007. Magic sets and their application to data integration. *Journal of Computer and System Sciences* 73, 4, 584–609.
- FABER, W., PFEIFER, G. AND LEONE, N. 2010. Semantics and complexity of recursive aggregates in answer set programming. *Artificial Intelligence* In Press, Corrected Proof, 1–21.
- FUXMAN, A., FAZLI, E. AND MILLER, R. J. 2005. ConQuer: Efficient management of inconsistent databases. In *Proceedings of SIGMOD'05*, F. Özcan, Ed. ACM, New York, 155–166.
- FUXMAN, A. AND MILLER, R. J. 2007. First-order query rewriting for inconsistent databases. *Journal of Computer and System Sciences* 73, 4, 610–635. Special Issue: Database Theory 2005.
- GARCIA-MOLINA, H., PAPANIKOLAOU, Y., QUASS, D., RAJARAMAN, A., SAGIV, Y., ULLMAN, J., VASSALOS, V. AND WIDOM, J. 1997. The TSIMMIS approach to mediation: Data models and languages. *Journal of Intelligent Information Systems* 8, 2, 117–132.
- GELFOND, M. AND LIFSCHITZ, V. 1988. The stable model semantics for logic programming. In *Proceedings of ICLP/SLP'88*, R. A. Kowalski and K. A. Bowen, Eds. MIT Press, Cambridge, MA, 1070–1080.
- GELFOND, M. AND LIFSCHITZ, V. 1991. Classical negation in logic programs and disjunctive databases. *New Generation Computing* 9, 3–4, 365–385.
- GOH, C. H., BRESSAN, S., MADNICK, S. AND SIEGEL, M. 1999. Context interchange: New features and formalisms for the intelligent integration of information. *ACM Transactions on Information Systems* 17, 3, 270–293.
- GRECO, G., GRECO, S. AND ZUMPANO, E. 2001. A logic programming approach to the integration, repairing and querying of inconsistent databases. In *Proceedings of ICLP'01*, P. Codognot, Ed. Lecture Notes in Computer Science, vol. 17. Springer, Berlin/Heidelberg, 348–364.
- GRECO, S. AND ZUMPANO, E. 2000. Querying inconsistent databases. In *Proceedings of LPAR'00*, M. Parigot and A. Voronkov, Eds. Springer-Verlag, Berlin/Heidelberg, 308–325.
- GRIECO, L., LEMBO, D., ROSATI, R. AND RUZZI, M. 2005. Consistent query answering under key and exclusion dependencies: Algorithms and experiments. In *Proceedings of CIKM'05*, O. Herzog, H.-J. Schek, N. Fuhr, A. Chowdhury and W. Teiken, Eds. ACM, New York, 792–799.
- LEMBO, D. 2004. *Dealing with Inconsistency and Incompleteness in Data Integration*. PhD Thesis, Dipartimento di Informatica e Sistemistica, Università di Roma “La Sapienza”.
- LENZERINI, M. 2002. Data integration: A theoretical perspective. In *Proceedings of PODS'02*, L. Popa, Ed. ACM, New York, 233–246.
- LEONE, N., GRECO, G., IANNI, G., LIO, V., TERRACINA, G., EITER, T., FABER, W., FINK, M., GOTTLÖB, G., ROSATI, R., LEMBO, D., LENZERINI, M., RUZZI, M., KALKA, E., NOWICKI, B. AND STANISZKIS, W. 2005. The INFOMIX system for advanced integration of incomplete and inconsistent data. In *Proceedings of SIGMOD'05*, F. Özcan, Ed. ACM, New York, 915–917.
- LEONE, N., PFEIFER, G., FABER, W., EITER, T., GOTTLÖB, G., PERRI, S. AND SCARCELLO, F. 2006. The DLV system for knowledge representation and reasoning. *ACM Transactions on Computational Logic* 7, 3, 499–562.
- LEVENE, M. AND VINCENT, M. W. 2000. Justification for inclusion dependency normal form. *IEEE Transactions on Knowledge and Data Engineering* 12, 2, 281–291.
- MANNA, M., RICCA, F. AND TERRACINA, G. 2011. Consistent query answering via ASP from different perspectives: Theory and practice. *CoRR cs.DB/1107.4570*, 1–30.

- MARILEO, M. C. AND BERTOSSI, L. E. 2010. The consistency extractor system: Answer set programs for consistent query answering in databases. *Data and Knowledge Engineering* 69, 6, 545–572.
- MINKER, J. 1982. On indefinite data bases and the closed world assumption. In *Proceedings of CADE'82*, D. W. Loveland, Ed. Lecture Notes in Computer Science, vol. 138. Springer, Berlin/Heidelberg, 292–308.
- TERRACINA, G., DE FRANCESCO, E., PANETTA, C. AND LEONE, N. 2008. Enhancing a DLP system for advanced database applications. In *Proceedings of RR'08*, D. Calvanese and G. Lausen, Eds. Lecture Notes in Computer Science, vol. 5341. Springer, Berlin/Heidelberg, 119–134.
- TERRACINA, G., LEONE, N., LIO, V. AND PANETTA, C. 2008. Experimenting with recursive queries in database and logic programming systems. *Theory and Practice of Logic Programming* 8, 2, 129–165.
- TOMASIC, A., RASCHID, L. AND VALDURIEZ, P. 1998. Scaling access to heterogeneous data sources with DISCO. *IEEE Transactions on Knowledge and Data Engineering* 10, 5, 808–823.