

Complexity of Consistent Query Answering in Databases Under Cardinality-Based and Incremental Repair Semantics

Andrei Lopatenko^{1,*} and Leopoldo Bertossi²

¹ Free University of Bozen-Bolzano
Faculty of Computer Science Bozen-Bolzano, Italy
lopatenko@inf.unibz.it

² Carleton University, School of Computer Science
Ottawa, Canada
bertossi@scs.carleton.ca

Abstract. A database D may be inconsistent wrt a given set IC of integrity constraints. Consistent Query Answering (CQA) is the problem of computing from D the answers to a query that are consistent wrt IC . Consistent answers are invariant under all the *repairs* of D , i.e. the consistent instances that minimally depart from D . Three classes of repair have been considered in the literature: those that minimize set-theoretically the set of tuples in the symmetric difference; those that minimize the changes of attribute values, and those that minimize the cardinality of the set of tuples in the symmetric difference. The latter class has not been systematically investigated. In this paper we obtain algorithmic and complexity theoretic results for CQA under this cardinality-based repair semantics. We do this in the usual, static setting, but also in a dynamic framework where a consistent database is affected by a sequence of updates, which may make it inconsistent. We also establish comparative results with the other two kinds of repairs in the dynamic case.

1 Introduction

The purpose of *consistent query answering* (CQA) is to compute query answers that are consistent with certain integrity constraints (ICs) that the database as a whole may fail to satisfy. Consistent answers have been characterized as those that are invariant under minimal forms of restoration of the consistency of the database [1, 5]. A particular and first notion of minimal restoration of consistency was captured in [1] in terms of database *repairs*, i.e. consistent database instances that share the schema with the original database, but differ from the latter by a *minimal set of whole tuples under set inclusion*. In this paper we call this semantics “the S-repair semantics”, for being set oriented. In [5, 15, 1, 7, 3, 9], complexity bounds for CQA under the S-repair semantics have been reported.

Two other repair semantics naturally arise and have been considered in the literature. The *A-repair semantics* is based on changing in a minimal way *attribute*

* Current affiliation: Google, Zürich.

values in database tuples in order to restore consistency. CQA under the A-repair semantics has also been investigated [29, 14, 4, 12]. The *C-repair semantics* is based on repairs of the original database that minimize the *cardinality* of the set of tuples by which the instances differ [2]. This semantics has received much less attention so far.

Example 1. Consider a database schema $P(X, Y, Z)$ with the functional dependency $X \rightarrow Y$. The inconsistent instance $D = \{P(a, b, c), P(a, c, d), P(a, c, e)\}$, seen as a set of ground atoms, has two S-repairs, $D_1 = \{P(a, b, c)\}$ and $D_2 = \{P(a, c, d), P(a, c, e)\}$, because the symmetric set differences with D , $\Delta(D, D_1)$ and $\Delta(D, D_2)$, are minimal under set inclusion. However, only for D_2 the cardinality $|\Delta(D, D_2)|$ of the symmetric set difference is minimum; and D_2 is the only C-repair.

The query $P(x, y, z)$ has consistent answers (a, c, d) and (a, c, e) under the C-repair semantics (they are classic answers in the only C-repair), but none under the S-repair semantics (the two S-repairs share no classic answers). \square

The consistent query answers under C-repairs form a superset of the consistent answers under S-repairs, because every C-repair is also an S-repair. Actually, in situations where the S-repair semantics does not give any consistent answers, the C-repair semantics may return answers. These answers could be further filtered out according to other criteria at a post-processing stage. For example, in the extreme case where there is only one database tuple in semantic conflict with a possibly large set of other tuples, the existence of an S-repair containing the only conflicting tuple would easily lead to an empty set of consistent answers. The C-repair semantics would not allow such a repair (c.f. Example 3 below).

Furthermore, the C-repair semantics has the interesting property that CQA, a form of *cautious* or *certain* reasoning (declaring true what is true in *all* repairs), and its *brave* or *possible* version (i.e. true in *some* repair), are mutually reducible in polynomial time and share the same data complexity. This is established in Section 3 by proving first some useful graph-theoretic lemmas about maximum independent sets that are interesting in themselves, and have a wider applicability in the context of CQA.

In [2], C-repairs were specified using disjunctive logic programs with stable model semantics [17] and weak cardinality constraints [6]. In this paper, applying the graph-theoretic techniques and results mentioned above, we obtain the first non-trivial complexity results for CQA under the C-repair semantics. Our emphasis is on CQA, as opposed to computing or checking specific repairs.

All the complexity bounds on CQA given so far in the literature, no matter which repair semantics is chosen, consider *the static case*: Given a snapshot of a database, a set of integrity constraints, and a query, the problems are the computation and verification of consistent answers to the query. In this paper

we also take into account dynamic aspects of data, studying the complexity of CQA when the consistency of a database may be affected by update actions.

Example 2. (example 1 continued) The C-repair $D_2 = \{P(a, c, d), P(a, c, e)\}$ is obviously consistent, however after the execution of the update operation $insert(P(a, f, d))$ it becomes inconsistent. In this case, the only C-repair of $D_2 \cup \{P(a, f, d)\}$ is D_2 itself. So, CQA from $D_2 \cup \{P(a, f, d)\}$ amounts to classic query answering from D_2 . However, if we start from the consistent instance $D' = \{P(a, c, d)\}$, executing the same update operation leads to two C-repairs, D' and also $\{P(a, f, d)\}$, and now CQA from $D' \cup \{P(a, f, d)\}$ is different from classic query answering from D' , because two repairs have to be considered. \square

Understanding and handling CQA in a dynamic setting is crucial for its applicability. Incremental methods should be developed, since it would be inefficient to compute a materialized repair of the database or a consistent answer to a query from scratch after every update.

While we think that the right repair semantics may be application dependent, being able to compare the possible semantics in terms of complexity may also shed some light on what may be the repair semantics of choice. This comparison should consider both static and incremental CQA, because a specific semantics might be better than others in terms of complexity when the database is affected by certain updates. In this paper we compare the C-repair semantics with the S- and A-repair semantics mentioned before, and both in the static and incremental settings.

In Section 3 we prove that static CQA under C-repairs is $P^{NP(\log(n))}$ -hard for denial constraints and ground atomic queries; which contrasts with the *PTIME* result for S-repairs in [9]. On the other side, in Section 4, we prove that incremental CQA, i.e. CQA in the dynamic setting, under the C-repair semantics is in *PTIME* for denial constraints and conjunctive queries; and that the same problem under S-repairs is *coNP*-hard (in data).

The naive algorithms for incremental CQA under the C-repair semantics are polynomial in data, but exponential in the size of the update sequence. In consequence, we also study the *parameterized complexity* [10, 13] of incremental CQA under the C-repair semantics, being the parameter the size of the update sequence. We establish that the problem is *fixed parameter tractable* (FPT).

For establishing comparisons with the C-repair semantics, we obtain new results on the static and incremental complexity both under the classic, i.e. S-repair semantics, and the A-repair semantics. We prove, for the former, that incremental CQA is *coNP*-hard; whereas for the latter, static and incremental CQA become both P^{NP} -hard in data.

We concentrate on relational databases and denial integrity constraints, which include most of the constraints found in applications where inconsistencies naturally arise, e.g. census-like databases [4], experimental samples databases, biological databases, etc. Complexity results refer to data complexity. For complexity theory we refer to [26]; and to [13] for parameterized complexity. Proofs of the results in this paper can be found in [22].

2 Semantics for Consistent Query Answering

A relational database instance D is a finite set of ground atoms $R(\bar{t})$ (also called *database tuples*¹), where R is a relation in the schema \mathcal{D} , and \bar{t} is a finite sequence of constants from the domain \mathcal{U} . A database atom is of the form $R(\bar{t})$, where R is a predicate in \mathcal{D} , and \bar{t} may contain constants or variables. A database literal is a database atom or a negation of a database atom. With $\Delta(D', D)$ we denote the symmetric difference $(D' \setminus D) \cup (D \setminus D')$ between instances D, D' , conceived both as sets of ground atoms.

The relational schema \mathcal{D} determines a first-order language $L(\mathcal{D})$ based on the relation names, the elements of \mathcal{U} , and extra built-in predicates. In the language $L(\mathcal{D})$, integrity constraints are sentences, and queries are formulas, usually with free variables. We assume in this paper that sets IC of ICs are always consistent in the sense that they are simultaneously satisfiable as first-order sentences. A database is *consistent* wrt to a given set of integrity constraints IC if the sentences in IC are all true in D , denoted $D \models IC$. An answer to a query $Q(\bar{x})$, with free variables \bar{x} , is a tuple \bar{t} that makes Q true in D when the variables in \bar{x} are interpreted as the corresponding values in \bar{t} , denoted $D \models Q[\bar{t}]$.

Definition 1. For a database D , integrity constraints IC , and a partial order $\preceq_{D,S}$ over databases that depends on the original database D and a repair semantics \mathcal{S} , a *repair of D wrt IC under \mathcal{S}* is an instance D' such that: (a) D' has the same schema and domain as D ; (b) $D' \models IC$; and (c) there is no D'' satisfying (a) and (b), such that $D'' \prec_{D,S} D'$, i.e. $D'' \preceq_{D,S} D'$ and not $D' \preceq_{D,S} D''$. The set of all repairs is denoted with $Rep(D, IC, \mathcal{S})$. \square

The class $Rep(D, IC, \mathcal{S})$ depends upon the semantics \mathcal{S} , that determines the partial order \preceq and the way repairs can be obtained, e.g. by allowing both insertions and deletions of whole database tuples [1], or deletions of them only [9], or only changes of attribute values [29, 4, 12], etc. (c.f. Definition 2.) We summarize here the most common repair semantics.

Definition 2. (a) *S-repair semantics* [1]: $D' \preceq_{D,S} D''$ iff $\Delta(D', D) \subseteq \Delta(D'', D)$.
 (b) *C-repair semantics*: $D' \preceq_{D,C} D''$ iff $|\Delta(D', D)| \leq |\Delta(D'', D)|$.
 (c) *A-repair semantics*: $D' \preceq_{D,A} D''$ iff $f(D, D') \leq f(D, D'')$, where f is a fixed numerical aggregation function over differences of attribute values. \square

More details about the A-repair semantics can be found in Section 4.3. Particular cases of A-repairs can be found in [14, 12], where the aggregation function to be minimized is the number of all attribute changes; and in [4], where the function is the overall quadratic difference obtained from the changes in numerical attributes between the original database and the repair. S-repairs and C-repairs are “tuple-based”, in the sense that consistency is restored by inserting and/or deleting whole database tuples; whereas A-repairs are obtained by changing attributes values in existing tuples only.

¹ We also use the term *tuple* to refer to a finite sequence $\bar{t} = (c_1, \dots, c_n)$ of constants of the database domain \mathcal{U} , but a *database tuple* is a ground atomic sentence with predicate in \mathcal{D} (excluding built-ins predicates, like comparisons).

In Example 1, attribute-based repairs could be $\{P(a, c, c), P(a, c, d), P(a, c, e)\}$, suggesting that we made a mistake in the second argument of the first tuple, but also $\{P(a, b, c), P(a, b, d), P(a, b, e)\}$. If the aggregate function in Definition 2(c) is the number of changes in attribute values, the former would be a repair, but not the latter. A-repairs may not be S- or C-repairs if the changes of attribute values have to be simulated via deletions followed by insertions.

Definition 3. Let D be a database, IC a set of ICs, and $Q(\bar{x})$ a query. (a) A ground tuple \bar{t} is a *consistent answer* to Q wrt IC under semantics \mathcal{S} if for every $D' \in Rep(D, IC, \mathcal{S})$, $D' \models Q[\bar{t}]$. (b) $Cqa(Q, D, IC, \mathcal{S})$ is the set of consistent answers to Q in D wrt IC under semantics \mathcal{S} . If Q is a sentence (a boolean query), $Cqa(Q, D, IC, \mathcal{S}) := \{yes\}$ when $D' \models Q$ for every $D' \in Rep(D, IC, \mathcal{S})$, and $Cqa(Q, D, IC, \mathcal{S}) := \{no\}$, otherwise. (c) $CQA(Q, IC, \mathcal{S}) := \{(D, \bar{t}) \mid \bar{t} \in Cqa(Q, D, IC, \mathcal{S})\}$ is the *decision problem of consistent query answering*. \square

Denial constraints are integrity constraints expressed by $L(\mathcal{D})$ -sentences of the form $\forall \bar{x} \neg (A_1 \wedge \dots \wedge A_m \wedge \gamma)$, where each A_i is a database atom and γ is a conjunction of comparison atoms. In particular, functional dependencies (FDs), e.g. $\forall x \forall y \forall z \neg (R(x, y) \wedge R(x, z) \wedge y \neq z)$, are denial constraints. For denial ICs, tuple-based repairs are obtained by tuple deletions only [9].

3 Complexity of CQA Under the C-Repair Semantics

As a consequence of the specification of C-repairs as the stable models of disjunctive logic programs with non-prioritized weak constraints [2] and the results in [6], we obtain that an upper bound on the data complexity of CQA under the C-repair semantics is the class $\Delta_3^P(\log(n))$.

In [3], *conflict graphs* were first introduced to study the complexity of CQA for aggregate queries wrt FDs under the S-repair semantics. They have as vertices the database tuples; and edges connect two tuples that simultaneously violate a FD. There is a one-to-one correspondence between S-repairs of the database and the set-theoretically *maximal* independent sets in the conflict graph. Similarly, there is a one-to-one correspondence between C-repairs and *maximum* independent sets in the same graph (but now they are maximum in cardinality).

Conflict graphs for databases wrt general denial constraints become *conflict hypergraphs* [9] that have as vertices the database tuples, and as hyperedges the (set theoretically minimal) collections of tuples that simultaneously violate one of the denial constraints. The size of the hypergraph (including vertices and hyperedges) is polynomial in the size of the database, because we have a fixed set of denial constraints. The correspondence for conflict graphs between repairs and independent sets – maximum or maximal depending on the semantics – still holds for hypergraphs, where an independent set in an hypergraph is a set of vertices that does not contain any hyperedges [9].

Notice that, unless an IC forces a particular tuple not to belong to the database,² every tuple in the original database belongs to some S-repair, but

² We do not consider in this work such *non generic* ICs [5].

not necessarily to a C-repair (c.f. Example 1, where the tuple $P(a, b, c)$ does not belong to the only C-repair).

In consequence, testing membership of vertices to some maximum independent set becomes a relevant for C-repairs. The complexity of this problem will determine the complexity of CQA under the C-repair semantics. For this purpose we will use some graph-theoretic constructions and lemmas about maximum independent sets, whose proofs use a self-reducibility property of independent sets that can be expressed as follows: For any graph G and vertex v , every maximum independent set that contains v (meaning maximum among the independent sets that contain v) consists of vertex v together with a maximum independent set of the graph G' that is obtained from G by deleting all vertices adjacent to v .

To keep the presentation simpler, we concentrate mostly on conflicts graphs and FDs. However, the results obtained carry over to denial constraints and their hypergraphs. Notice, as a motivation for the next lemmas, that a ground atomic query is consistently true when it belongs, as a database tuple, i.e. as a vertex in the conflict graph, to all the maximum independent sets of the conflict graph.

Lemma 1. Consider a graph G and a vertex v in it. (a) For the graph G' obtained by adding a new vertex v' that is connected only to the neighbors of v , the following properties are equivalent: 1. There is a maximum independent set of G containing v . 2. v belongs to every maximum independent set of G' . 3. The sizes of maximum independent sets in G and G' differ by one. (b) There is a graph G' extending G that can be constructed in logarithmic space, such that v belongs to all maximum independent sets of G iff v belongs to some maximum independent set of G' . \square

From this lemma and the membership to $FP^{NP(\log(n))}$ of computing the size of a maximum clique in a graph [21], we obtain

Lemma 2. The problems of deciding for a vertex in a graph if it belongs to some maximum independent set and if it belongs to all maximum independent sets are both in $P^{NP(\log(n))}$. \square

Theorem 1. For functional dependencies and ground atomic queries, CQA under the C-repair semantics belongs to $P^{NP(\log(n))}$. \square

Considering the maximum independent sets, i.e. C-repairs, as a collection of possible worlds, the previous lemma shows a close connection between the *certain* C-repair semantics (true in *every* repair), that is the basis for CQA, and the *possible* C-repair semantics (true in *some* repair). CQA under these semantics and functional dependencies are polynomially reducible to each other; actually also for negations of ground atomic queries.

Lemma 3. The following problems are mutually *LOGSPACE*-reducible to each other: (1) *Certain positive*: Given a vertex v and a graph G , decide if v belongs to every maximum independent set of G . (2) *Certain negative*: Given a vertex v and a graph G , decide if all the maximum independent sets of G do not contain v . (3) *Possible negative*: Given a vertex v and a graph G , decide if

there is a maximum independent set of G that does not contain v . (4) *Possible positive*: Given a vertex v and a graph G , decide if v belongs to at least one maximum independent set of G . \square

Since the negation $\neg R(\bar{t})$ of a ground atomic query $R(\bar{t})$ is consistently true wrt the C-repair semantics iff the vertex corresponding to $R(\bar{t})$ in the conflict graph does not belong to any maximum independent set, using Lemma 3 we can extend Theorem 1 to conjunctions of literals.³ Actually, since Lemmas 1, 2 and 3 still hold for hypergraphs, we obtain

Theorem 2. For denial constraints and queries that are conjunctions of literals, CQA under the C-repair semantics belongs to $P^{NP(log(n))}$. \square

Now we will represent the maximum independent sets of a graph as C-repairs of an inconsistent database wrt a denial constraint. This is interesting, because conflict graphs for databases wrt denial constraints are, as indicate before, actually conflict hypergraphs.

Lemma 4. There is a fixed database schema \mathcal{D} and a denial constraint φ in $L(\mathcal{D})$, such that for every graph G , there is an instance D over \mathcal{D} , whose C-repairs wrt φ are in one-to-one correspondence with the maximum independent sets of G . Furthermore, D can be built in polynomial time in the size of G . \square

From Lemma 4 and the $P^{NP(log(n))}$ -completeness of determining the size of a maximum clique [21], we obtain

Theorem 3. Determining the size of a C-repair for denial constraints is complete for $FP^{NP(log(n))}$. \square

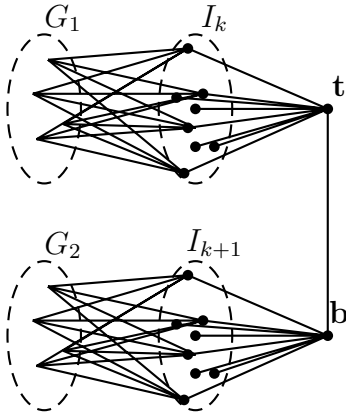


Fig. 1. The block $B_k(G, t)$

In order to obtain hardness for CQA under the C-repair semantics, we need

to construct the block graph $B_k(G, t)$ (c.f. Figure 1), consisting of two copies G_1, G_2 of G , and two internally disconnected subgraphs I_k, I_{k+1} , with k and $k + 1$ vertices, resp. Every vertex in G (G') is connected to every vertex in I_k (resp. I_{k+1}).

Lemma 5. Given a graph G and a number k , a graph $B_k(G, t)$ can be computed in polynomial time in the size of G , where t is a distinguished vertex in it that belongs to all its maximum independent sets iff the cardinality of a maximum independent set of G is equal to k . \square

³ This can also be obtained, less directly, from the closure of $P^{NP(log(n))}$ under complement.

Lemma 6. Deciding if a vertex belongs to all maximum independent sets of a graph is $P^{NP(\log(n))}$ -hard. \square

This result can be proved by reduction from the following $P^{NP(\log(n))}$ -complete decision problem [21]: Given a graph G and an integer k , is the size of a maximum clique in G equivalent to $0 \pmod k$? G is reduced to a graph G' that is built by combining a number of versions of the block construction in Figure 1. Now, the graph G' used in Lemma 6 can be represented according to Lemma 4 as a database consistency problem, and in this way we obtain

Theorem 4. For denial constraints, CQA under the C-repair semantics for queries that are conjunctions of ground literals is $P^{NP(\log(n))}$ -complete. \square

This theorem still holds for ground atomic queries, which is interesting, because for this kind of queries and denial constraints CQA under the S-repair semantics is in $PTIME$ [9].

4 Incremental Complexity of CQA

Assume that we have a consistent database instance D wrt to IC . D may become inconsistent after the execution of an update sequence U composed of operations of the forms $insert(R(\bar{t}))$, $delete(R(\bar{t}))$, meaning insert/delete tuple $R(\bar{t})$ into/from D , or $change(R(\bar{t}), A, a)$, for changing value of attribute A in $R(\bar{t})$ to a , with $a \in \mathcal{U}$. We are interested in whether we can find consistent query answers from the possibly inconsistently updated database $U(D)$ more efficiently by taking into account the previous consistent database state.

Definition 4. For a consistent database D wrt IC , and a sequence U of update operations U_1, \dots, U_m , *incremental consistent query answering* for query Q is CQA for Q wrt IC from instance $U(D)$, that results from applying U to D . \square

Update sequences U will be atomic, in the sense that they are completely executed or not. This allows us to concentrate on “minimized” versions of update sequences, e.g. containing only insertions and/or attribute changes when dealing with denial constraints, because deletions do not cause any violations. We are still interested in data complexity, i.e. wrt the size $|D|$ of the original database. In particular, m is fixed, and usually small wrt $|D|$.

A notion of incremental complexity has been introduced in [23], and also in [20] under the name of *dynamic complexity*. There, the instance that is updated can be arbitrary, and the question is about the complexity for the updated version when information about the previous instance can be used. In our case, we are assuming that the initial database is consistent. As opposed to [23, 20], where new incremental or dynamic complexity classes are introduced, we appeal to those classic complexity classes found at a low level in the polynomial hierarchy.

4.1 Incremental Complexity: C-Repair Semantics

In contrast to static CQA for the C-repair semantics, it holds

Theorem 5. For the C-repair semantics, first-order boolean queries, denial constraints, and update sequences U of fixed length m applied to D , incremental CQA is in $P\text{TIME}$ in $|D|$. \square

The proof of this theorem provides an upper bound of $O(m \cdot n^m)$, that is polynomial in the size n of the initial database, but exponential in m , which makes the problem tractable in data, but with the size of the update sequence in the exponent. We are interested in determining if queries can be consistently answered in time $O(f(m) \times n^c)$, for a constant c and a function $f(m)$ depending only on m . In this way we isolate the complexity introduced by U .

The area of parameterized complexity studies this kind of problems [19, 25]. A decision problem with inputs of the form (I, p) , where p is a distinguished parameter of the input, is *fixed parameter tractable*, and by definition belongs to the class FPT [10], if it can be solved in time $O(f(|p|) \cdot |I|^c)$, where c and the hidden constant do not depend on $|p|$ or $|I|$ and f does not depend on $|I|$.

Definition 5. Given a query Q , ICs IC , and a ground tuple \bar{t} , *parameterized incremental CQA* is the decision problem $CQA^p(Q, IC) := \{(D, U, \bar{t}) \mid D \text{ is an instance, } U \text{ an update sequence, } \bar{t} \text{ is consistent answer to } Q \text{ in } U(D)\}$, whose parameter is U , and consistency of answers refers to C-repairs of $U(D)$. \square

We keep Q and IC fixed in the problem definition because, except for the parameter U , we are interested in data complexity.

Theorem 6. For functional dependencies and queries that are conjunctions of literals, parameterized incremental CQA is in FPT . \square

The *vertex cover problem*, of deciding if graph G has a vertex cover (VC) of size no bigger than k , belongs to the class FPT , i.e. there is a polynomial time parameterized algorithm $VC(G, k)$ for it [10]; actually one that runs in time $O(1.2852^k + k \cdot n)$, being n the size of G [8].

The algorithm whose existence is claimed in Theorem 6 is as follows: Let G be the conflict graph associated to the database obtained after the insertion of m tuples. By binary search, calling each time $VC(G, -)$, it is possible to determine the size of a minimum VC for G . This gives us the minimum number of tuples that have to be removed in order to restore consistency; and can be done in time $O(\log(m) \cdot (1.2852^m + m \cdot n))$, where n is the size of the original database. In order to determine if a tuple $R(\bar{t})$ belongs to every maximum independent set, i.e. if it is consistently true, compute the size of a minimum VC for $G \setminus \{R(\bar{t})\}$. The two numbers are the same iff the answer is *yes*. The total time is still $O(\log(m) \cdot (1.2852^m + m \cdot n))$, which is linear in the size of the original database. The same algorithm applies if, in addition to tuple insertions, we also have changes of attribute values in the update part; of course, still under the C-repair semantics.

Theorem 6 uses the membership to FPT of the VC problem, which we apply to conflict graphs for functional dependencies. However, the result can be extended to denials constraints and their conflict hypergraphs. In our case, the maximum

size of an hyperedge is the maximum number of database atoms in a denial constraint, which is determined by the fixed database schema. If this number is d , then we are in the presence of the so-called *d-hitting set problem*, consisting in finding the size of a minimum hitting set for an hypergraph with hyperedges bounded in size by d . This problem is in *FPT* [24].

Theorem 7. For denial constraints and queries that are conjunctions of literals, parameterized incremental CQA is in *FPT*. \square

Using the reductions in Section 3, this result can be extended to incremental CQA under the *possible C-repair* semantics.

4.2 Incremental Complexity: S-Repair Semantics

Incremental CQA for non-quantified conjunctive queries under denial constraints belongs to *PTIME*, which can be established by applying the algorithm in [9] for the static case to $U(D)$.

However, for quantified conjunctive queries the situation may change. Actually, by reduction from static CQA for conjunctive queries and denial ICs under the S-repair semantics, which is *coNP*-hard [9], we obtain

Theorem 8. Under the S-repair semantics, incremental CQA for conjunctive queries and denial constraints is *coNP*-hard. \square

We can see that, for denial constraints, static CQA under the C-repair semantics seems to be harder than under the S-repair semantics ($P^{NP(\log(n))}$ - vs. *coNP*-hard). On the other side, incremental CQA under the S-repair semantics seems to be harder than under the C-repair semantics (*coNP*-hard vs. *PTIME*). The reason is that for the C-repair semantics the cost of a repair cannot exceed the size of the update, whereas for the S-repair semantics the cost of a repair may be unbounded wrt the size of an update.

Example 3. Consider a schema $R(\cdot), S(\cdot)$ with the denial constraint $\forall x \forall y \neg (R(x) \wedge S(y))$; and the consistent database $D = \{R(1), \dots, R(n)\}$, with an empty table for S . After the update $U = \text{insert}(S(0))$, the database becomes inconsistent, and the S-repairs are $\{R(1), \dots, R(n)\}$ and $\{S(0)\}$. However, only the former is a C-repair, and is at a distance 1 from the original instance, i.e. as the size of the update. However, the second S-repair is at a distance n . \square

4.3 Incremental Complexity: A-Repair Semantics

Before addressing the problem of incremental complexity, we give a complexity lower bound for the *weighted* version of static CQA for the A-repair semantics. In this case, we have a numerical weight function w defined on triples of the form $(R(\bar{t}), A, \text{newValue})$, where $R(\bar{t})$ is a database tuple stored in the database, A is an attribute of R , and newValue is a new value for A in $R(\bar{t})$. The *weighted A-repair semantics* (wA-repair semantics) is just a particular case of Definition 2(c), where the distance is given by an aggregation function g applied to the set of numbers $\{w(R(\bar{t}), A, \text{newValue}) \mid R(\bar{t}) \in D\}$.

Typically, g is the sum, and the weights are $w(R(\bar{t}), A, newValue) = 1$ if $R(\bar{t})[A]$ is different from $newValue$, and 0 otherwise, where $R(\bar{t})[A]$ is the projection of database tuple $R(\bar{t})$ on attribute A , i.e. just the number of changes is counted [14]. In [4], g is still the sum, but w is given by $w(R(\bar{t}), A, newValue) = \alpha_A \cdot (R(\bar{t})[A] - newValue)^2$, where α_A is a coefficient introduced to capture the relative importance of attribute A or scale factors. In these cases, w does not depend on D . However, if the weight function w depended on the size of D , w should become part of the input for the decision problem of CQA.

Theorem 9. Static CQA for ground atomic queries and denial constraints under the wA-repair semantics is P^{NP} -hard. \square

In order to obtain a hardness result in the incremental case and for denial constraints (for which we are assuming update sequences do not contain tuple deletions), we can use the kind of A-repairs introduced in [4].

Theorem 10. Incremental CQA for atomic queries and denial constraints under the wA-repair semantics is P^{NP} -hard. \square

These results still hold for tuple insertions as update actions, the fixed weight function that assigns value 1 to every change, and the sum as aggregation function. In case we have numerical values as in [4] or a bounded domain, we can obtain as in [4, theorem 4(b)] that the problems in Theorems 9 and 10 belong both to Π_2^P .

Under the A-repair semantics, if the update sequence consist of *change* actions, then we can obtain polynomial time incremental CQA under the additional condition that the set of attribute values than can be used to restore consistency is bounded in size, independently from the database (or its active domain). Such an assumption can be justified in several applications, like in census-like databases that are corrected according to inequality-free denial constraints that force the new values to be taken at the border of a database independent region [4]; and also in applications where denial constraints, this time containing inequalities, force the attribute values to be taken in a finite, pre-specified set. The proof is similar to the one of Theorem 5, and the polynomial bound now also depends on the size of the set of candidate values.

Theorem 11. For a database independent and bounded domain of attribute values, incremental CQA under the A-repair semantics, for first-order boolean queries, denial constraints, and update sequences containing only *change* actions is in $PTIME$ in the size of the original database. \square

Now, we present a lower bound for CQA under the A-repair semantics for *first-order ICs* and tuple deletions, which now may affect their satisfaction.

Lemma 7. For any planar graph G with vertices of degree at most 4, there exists a regular graph G' of degree 4 that is 4-colorable, such that G' is 3-colorable iff G is 3-colorable. G' can be built in polynomial time in $|G|$. \square

Notice that graph G , due to its planarity, is 4-colorable. The graph G' , is an extension of graph G that may not be planar, but preserves 4-Colorability. We use

the construction in Lemma 7 as follows: Given any planar graph G of degree 4, construct graph G' as in the lemma, which is regular of degree 4 and 4-colorable. Its 4-colorability is encoded as a database problem with a fixed set of first-order constraints. Since G' is 4-colorable, the database is consistent. Furthermore, G' uses all the 4 colors in the official table of colors, as specified by the ICs. In the update part, deleting one of the colors leaves us with the problem of coloring G' with only three colors (under an A-repair semantics only changes of colors are allowed to restore consistency), which is possible iff the original graph G is 3-colorable. Deciding about the latter problem is *NP*-complete [16]. We obtain

Theorem 12. For ground atomic queries, first-order ICs, and update sequences consisting of tuple deletions, incremental CQA under the A-repair semantics is *coNP*-hard. \square

To obtain this result it is good enough to use the sum as the aggregation function and the weight function that assigns 1 to each change. Clearly, this lower bound also applies to update sequences containing any combination of *insert*, *delete*, *change*.

5 Conclusions

The dynamic scenario for consistent query answering that considers possible updates on a database had not been considered before in the literature. Doing incremental CQA on the basis of the original database and the sequence of updates is an important and natural problem. Developing algorithms that take into account previously obtained consistent answers that are possible cached and the updates at hand is a crucial problem for making CQA scale up for real database applications. Much research is still needed in this direction.

In this paper we have concentrated mostly on complexity bounds for this problem under different semantics. When we started obtaining results for incremental CQA under repairs that differ from the original instance by a minimum number of tuples, i.e. C-repairs, we realized that this semantics had not been sufficiently explored in the literature in the static version of CQA, and that a full comparison was not possible. In the first part of this paper we studied the complexity of CQA for the C-repair semantics and denial constraints. In doing so, we developed graph-theoretic techniques for polynomially reducing each of the certain and possible (or cautious and brave) C-repair semantics for CQA to the other. A similar result does not hold for the S-repair semantics, conjunctive queries, and denial constraints: CQA (under the *certain* semantics) is *coNP*-complete [9], but is in *PTIME* for the *possible* semantics.

The complexity of CQA in a P2P setting was studied in [18], including a form a cardinality-based repairs. However, a different semantics is used, which makes it difficult to compare results. Actually, in that setting it is possible that repairs do not exist, whereas in our case, since S-repairs always exist [1], also C-repairs exist. The complexity result for CQA in [18], that seems to be shared by C- and S-repairs, is obtained on the basis of the complexity of checking the existence of repairs (a problem that in our case is trivial).

The C-repair semantics can be generalized considering weights on tuples. Under denial constraints, this means that it may be more costly to remove certain tuples than others to restore consistency. More precisely, database tuples $R(\bar{t})$ have associated numerical costs $w(R(\bar{t}))$, that become part of the input for the CQA decision problem. Now, the partial order between instances is given by $D_1 \preceq_{D,wC} D_2$ iff $|D\Delta D_1|_w \leq |D\Delta D_2|_w$, where, for a set of database tuples S , $|S|_w$ is the sum of the weights of the elements of S . It can be proved that CQA for ground atomic queries wrt denial constraints under this semantics belongs to P^{NP} [22, proposition 5].

Furthermore, it is possible to reduce CQA under the C-repair semantics to CQA under *least-squares* A-repairs semantics that minimizes the sum of the quadratic differences between numerical values [4], which is a particular case of the general semantics studied in Section 4.3.

Theorem 13. Given a database schema \mathcal{D} , a set IC of denial constraints in $L(\mathcal{D})$, and a ground atomic query $Q \in L(\mathcal{D})$, there are a schema \mathcal{D}' with some fixable numerical attributes, a set IC' of ICs in $L(\mathcal{D}')$, and a query $Q' \in L(\mathcal{D}')$, such that: For every database D over \mathcal{D} , there is a database D' over \mathcal{D}' that can be computed from D in $LOGSPACE$ (in data) for which it holds: Q is consistently true wrt IC in D under the C-repairs semantics iff Q' is consistently true wrt to IC' in D' under the least-squares A-repair semantics. \square

This result also applies to other numerical A-repair semantics as discussed in [4], and is about data complexity. For fixed \mathcal{D}, IC, Q, D , also fixed \mathcal{D}', IC', Q' can be obtained in $LOGSPACE$ from \mathcal{D}, IC, Q . Theorem 13, together with Theorem 4, allows us to obtain a simple proof of the $P^{NP(\log n)}$ -hardness of the least-squares repair semantics. In [4], P^{NP} -hardness is obtained for the latter as a better lower bound, but the proof is more complex. This theorem can be extended to the weighted C-repair semantics if integer numerical weights are used.

Our results show that the incremental complexity is lower than the static one in several useful cases, but sometimes the complexity cannot be lowered. It is a subject of ongoing work the development of concrete and explicit algorithms for incremental CQA.

We obtained the first results about fixed parameter tractability for incremental CQA, where the input, for a fixed database schema, can be seen as formed by the original database and the update sequence, whose length is the relevant parameter. This problem requires additional investigation. In particular, the parameterized complexity of incremental CQA under the S- and A-repair semantics has to be investigated, and a more complete picture still has to emerge.

It would be interesting to examine the area of CQA in general from the point of view of parameterized complexity, including the static case. Natural candidates to be a parameter in the classic, static setting could be: (a) the number of inconsistencies in the database, (b) the degree of inconsistency, i.e. the maximum number of violations per database tuple, (c) complexity of inconsistency, i.e. the length of the longest path in the conflict graph or hypergraph. These parameters may be practically significant, since in many applications, like census application [4], inconsistencies are “local”.

We considered a version of incremental CQA that assumes that the database is already consistent before updates are executed, a situation that could have been achieved because no previous updates violated the given semantic constraints or a repaired version was chosen before the new updates were executed.

We are currently investigating the dynamic case of CQA in the frameworks of *dynamic complexity* [20, 28] and *incremental complexity* as introduced in [23]. In this case we start with a database D that is not necessarily consistent on which a sequence of basic update operations U_1, U_2, \dots, U_m is executed. A clever algorithm for CQA may create or update intermediate data structures at each atomic update step, to help obtain answers at subsequent steps. We are interested in the complexity of CQA after a sequence of updates, when the data structures created by the query answering algorithm at previous states are themselves updatable and accessible.

Acknowledgments. Research supported by NSERC, and EU projects: Knowledge Web, Interop and Tones. L. Bertossi is Faculty Fellow of IBM Center for Advanced Studies (Toronto Lab.). L. Bertossi appreciates the hospitality and support of Enrico Franconi and the KRDB group in Bolzano. We are grateful to Jan Chomicki, Jörg Flum, and anonymous referees for many useful comments.

References

- [1] Arenas, M., Bertossi, L. and Chomicki, J. Consistent Query Answers in Inconsistent Databases. *Proc. ACM Symposium on Principles of Database Systems (PODS 99)*, ACM Press, 1999, pp. 68-79.
- [2] Arenas, M., Bertossi, L. and Chomicki, J. Answer Sets for Consistent Query Answering in Inconsistent Databases. *Theory and Practice of Logic Programming*, 2003, 3(4-5):393-424.
- [3] Arenas, M., Bertossi, L., Chomicki, J., He, X., Raghavan, V. and Spinrad, J. Scalar Aggregation in Inconsistent Databases. *Theoretical Computer Science*, 2003, 296:405-434.
- [4] Bertossi, L., Bravo, L., Franconi, E. and Lopatenko, A. Fixing Numerical Attributes under Integrity Constraints. *Proc. Tenth International Symposium on Database Programming Languages (DBPL 05)*, Springer LNCS 3774, 2005, pp. 262-278.
- [5] Bertossi, L. and Chomicki, J. Query Answering in Inconsistent Databases. In *Logics for Emerging Applications of Databases*. Springer, 2003, pp. 43-83.
- [6] Buccafurri, F., Leone, N. and Rullo, P. Enhancing Disjunctive Datalog by Constraints. *IEEE Transactions on Knowledge and Data Engineering*, 2000, 12(5):845-860.
- [7] Cali, A., Lembo, D. and Rosati, R. Complexity of Query Answering over Inconsistent and Incomplete Databases. *Proc. ACM Symposium on Principles of Database Systems (PODS 03)*, ACM Press, 2003, pp. 260-271.
- [8] Chen, J., Kanj, I. and Jia, W. Vertex Cover: Further Observations and Further Improvements. In *Proc. 25th International Workshop on Graph-Theoretic Concepts in Computer Science (WG 99)*, Springer LNCS 1665, 1999, pp. 313-324.
- [9] Chomicki, J. and Marcinkowski, J. Minimal-Change Integrity Maintenance using Tuple Deletions. *Information and Computation*, 2005, 197(1-2):90-121.
- [10] Downey, R.G. and Fellows, M.R. *Parameterized Complexity*. Springer, Monographs in Computer Science, 1999.

- [11] Eiter, T. and Gottlob, G. On the Complexity of Propositional Knowledge Base Revision, Updates, and Counterfactuals. *Artificial Intelligence*, 1992, 57(2-3):227-270.
- [12] Flesca, S., Furfaro, F. Parisi, F. Consistent Query Answers on Numerical Databases under Aggregate Constraints. *Proc. Tenth International Symposium on Database Programming Languages (DBPL 05)*, Springer LNCS 3774, 2005, pp. 279-294.
- [13] Flum, J. and Grohe, M. *Parameterized Complexity Theory*. Texts in Theoretical Computer Science, Springer Verlag, 2006.
- [14] Franconi, E., Laureti Palma, A., Leone, N., Perri, S. and Scarcello, F. Census Data Repair: a Challenging Application of Disjunctive Logic Programming. In *Proc. Logic for Programming, Artificial Intelligence, and Reasoning (LPAR 01)*, Springer LNCS 2250, 2001, pp. 561-578.
- [15] Fuxman, A. and Miller, R. First-Order Query Rewriting for Inconsistent Databases. *Proc. International Conference on Database Theory (ICDT 05)*, Springer LNCS 3363, 2004, pp. 337-351.
- [16] Garey, M., Johnson, D. and Stockmeyer, L. Some Simplified NP-Complete Graph Problems. *Theoretical Computer Science*, 1976, 1(3):237-267.
- [17] Gelfond, M. and Lifschitz, V. Classical Negation in Logic Programs and Disjunctive Databases. *New Generation Computing*, 1991, 9:365-385.
- [18] Greco, G. and Scarcello, F. On the Complexity of Computing Peer Agreements for Consistent Query Answering in Peer-to-Peer Data Integration Systems. *Proc. International Conference on Information and Knowledge Management (CIKM 05)*, ACM Press, 2005, pp. 36-43.
- [19] Grohe, M. Parameterized Complexity for the Data-base Theorist. *SIGMOD Record*, 2002, 31(4):86-96.
- [20] Immerman, N. *Descriptive Complexity*. Graduate Texts in Computer Science. Springer, 1999.
- [21] Krentel, M. The Complexity of Optimization Problems. *J. Computer and Systems Sciences*, 1988, 36:490-509.
- [22] Lopatenko, A. and Bertossi, L. Complexity of Consistent Query Answering in Databases under Cardinality-Based and Incremental Repair Semantics. Corr Archiv paper cs.DB/0604002.
- [23] Miltersen, P.B., Subramanian, S., Vitter, J.S. and Tamassia, R. Complexity Models for Incremental Computation. *Theoretical Computer Science*, 1994, 130(1):203-236.
- [24] Niedermeier, R. and Rossmannith, P. An Efficient Fixed-Parameter Algorithm for 3-Hitting Set. *Journal of Discrete Algorithms*, 2003, 1(1):89-102.
- [25] Papadimitriou, C.H. and Yannakakis, M. On the Complexity of Database Queries. *J. Comput. Syst. Sci.*, 1999, 58(3):407-427.
- [26] Papadimitriou, C. *Computational Complexity*. Addison-Wesley, 1994.
- [27] Robertson, N., Sanders, D.P., Seymour, P. and Thomas, R. Efficiently Four-Coloring Planar Graphs. In *Proc. 28th ACM Symposium on the Theory of Computing (STOC 96)*, ACM Press, 1996, pp. 571-575.
- [28] Weber, V. and Schwentick, T. Dynamic Complexity Theory Revisited. *Proc. Annual Symposium on Theoretical Aspects of Computer Science (STACS 05)*, Springer LNCS 3404, 2005, pp. 256-268.
- [29] Wijsen, J. Condensed Representation of Database Repairs for Consistent Query Answering. *Proc. International Conference on Database Theory (ICDT 03)*, Springer LNCS 2572, 2003, pp. 378-393.