# Repair Localization for Query Answering from Inconsistent Databases

THOMAS EITER and MICHAEL FINK
Technische Universität Wien
and
GIANLUIGI GRECO
Università della Calabria
and
DOMENICO LEMBO
SAPIENZA Università di Roma

---

Query answering from inconsistent databases amounts to finding "meaningful" answers to queries posed over database instances that do not satisfy integrity constraints specified over their schema. A declarative approach to this problem relies on the notion of repair, i.e., a database that satisfies integrity constraints and is obtained from the original inconsistent database by "minimally" adding and/or deleting tuples. Consistent answers to a user query are those answers that are in the evaluation of the query over each repair. Motivated by the fact that computing consistent answers from inconsistent databases is in general intractable, the present paper investigates techniques that allow to localize the difficult part of the computation on a small fragment of the database at hand, called "affected" part. Based on a number of localization results, an approach to query answering from inconsistent data is presented, in which the query is evaluated over each of the repairs of the affected part only, augmented with the part that is not affected. Single query results are then suitably recombined. For some relevant settings, techniques are also discussed to factorize repairs into components that can be processed independently of one another, thereby guaranteeing exponential gain w.r.t. the basic approach, which is not based on localization. The effectiveness of the results is demonstrated for consistent query answering over expressive schemas, based on logic programming specifications as proposed in the literature.

---

## 1. INTRODUCTION

A database is inconsistent if it does not satisfy the integrity constraints specified over its schema. This may happen for different reasons [Arenas et al. 2003]; for instance, when pre-existing data are re-organized under a new schema that has integrity constraints describing semantic aspects of the new scenario. This is particularly challenging in the context of data integration, where a number of data sources, heterogeneous and widely distributed, must be presented to the user as if they were a single (virtual) centralized database, which is often equipped with a rich set of constraints expressing important semantic properties of the application at hand. Since, in general, the integrated sources are autonomous, the data resulting from the integration are likely to violate these constraints.

One of the main issues arising when dealing with inconsistent databases is establishing the answers which have to be returned to a query issued over the database schema.

**Example 1.1** Consider a database schema $\chi_0$ providing information about soccer teams of the 2006/07 edition of the U.E.F.A. Champions League. The schema consists of the relation predicates $player(Pcode, Pname, Pteam)$, $team(Tcode, Tname, Tleader)$, and $coach(Ccode, Cname, Cteam)$. The associated constraints $\Sigma_0$ specify that the keys of $player$, $team$, and $coach$, are the sets of attributes $\{Pcode, Pteam\}$, $\{Tcode\}$, and $\{Ccode, Cteam\}$, respectively, and that a coach can neither be a player nor a team leader.

Consider the following inconsistent database $D_0$ for $\chi_0$ (possibly built by integrating some autonomous data sources):

$player^{D_0}:$

| 10 | Totti | RM |
|---|---|---|
| 9 | Ronaldinho | BC |

$team^{D_0}:$

| RM | Roma | 10 |
|---|---|---|
| BC | Barcelona | 8 |
| RM | Real Madrid | 10 |

$coach^{D_0}:$

| 7 | Capello | RM |
|---|---|---|

$D_0$ violates the key constraint on $team$, witnessed by the facts $team(RM, Roma, 10)$ and $team(RM, Real\ Madrid, 10)$, which coincide on $Tcode$ but differ on $Tname$. In such a situation, it is not clear what answers should be returned to a query over $D_0$ asking, for instance, for the names of teams, or for the pairs formed by team code and team leader. □

The standard approach to remedy the existence of conflicts in the data is through data cleaning [Bouzeghoub and Lenzerini 2001]. This approach is procedural in nature, and is based on domain-specific transformation mechanisms applied to the data. One of its problems is incomplete information on how certain conflicts should be resolved [Staworko et al. 2006]. This typically happens in systems which are not tailored for business logic support at the enterprise level, like systems for information integration on-demand. Here, data cleaning may be insufficient even if only few inconsistencies are present in the data.

In the last years, an alternative declarative approach has been investigated which builds on the notion of a *repair* for an inconsistent database [Arenas et al. 1999]. Roughly speaking, a repair is a new database which satisfies the constraints in the schema and minimally differs from the original one. The suitability of a possible repair depends on the underlying semantics adopted for the inconsistent database, and on the kinds of integrity constraints which are allowed on the schema. Importantly, in general, not a single but multiple repairs might be possible; therefore, the standard way of answering a user query is to compute the answers which are true in every possible repair, called *consistent answers* in the literature.

**Example 1.2** Recall that in our scenario, the database $D_0$ for $\chi_0$ violates the key constraint on $team$, witnessed by $team(RM, Roma, 10)$ and $team(RM, Real\ Madrid, 10)$.

$$player^{R_1}: \begin{array}{|c|c|c|} \hline 10 & Totti & RM \\ \hline 9 & Ronaldinho & BC \\ \hline \end{array} \qquad team^{R_1}: \begin{array}{|c|c|c|} \hline RM & Roma & 10 \\ \hline BC & Barcelona & 8 \\ \hline \end{array} \qquad coach^{R_1}: \begin{array}{|c|c|c|} \hline 7 & Capello & RM \\ \hline \end{array}$$

$$player^{R_2}: \begin{array}{|c|c|c|} \hline 10 & Totti & RM \\ \hline 9 & Ronaldinho & BC \\ \hline \end{array} \qquad team^{R_2}: \begin{array}{|c|c|c|} \hline BC & Barcelona & 8 \\ \hline RM & Real Madrid & 10 \\ \hline \end{array} \qquad coach^{R_2}: \begin{array}{|c|c|c|} \hline 7 & Capello & RM \\ \hline \end{array}$$

Fig. 1. Repairs of $D_0$.

According to [Arenas et al. 1999], a repair results by removing exactly one of these facts from $D_0$. Hence, there are two repairs only, say $R_1$ and $R_2$, which are as shown in Figure 1. Accordingly, the consistent answer to the query asking for the names of the teams is $\{(Barcelona)\}$, while the consistent answers to the query asking for pairs of team code and team leader are $\{(RM, 10), (BC, 8)\}$. □

Query answering in the presence of inconsistent data (a.k.a. consistent query answering) has been the subject of a large body of research (for a survey on this topic, see [Bertossi and Chomicki 2003], and for a discussion on relevant issues in the area see [Chomicki 2007]) and some prototype implementations of systems which fit the semantic repair framework are available [Fuxman et al. 2005; Chomicki et al. 2004b; Leone et al. 2005]. Basically, these systems differ in the kinds of constraints and queries they are able to deal with. Indeed, depending on these two ingredients, the data complexity of consistent query answering ranges from polynomial-time over co-NP up to $\Pi_2^P$ (see e.g., [Calì et al. 2003a; Chomicki and Marcinkowski 2005; Chomicki 2007]).

## 1.1 Contributions

In this paper, we elaborate techniques for consistent query answering in highly-expressive settings. Given that in these cases query answering is unlikely to be feasible in polynomial time, our main research interest is to devise an approach that allows to localize the "difficult" part of the computation in a small fragment of the database at hand.

The basic intuition of this approach is that resolving constraint violations in inconsistent databases does not generally require to deal with the whole set of facts. For instance, in Example 1.1 inconsistency may be fixed by just looking at the (few) tuples conflicting on the key. However, there are many interesting cases for which devising some similar strategies is not as simple as above and, therefore, it appears relevant to assess under which circumstances a localization approach can be pursued and when localized repair computation can be exploited to optimize consistent query answering. In this respect, our overall contribution is twofold in nature.

First, we attack the problem from a theoretic point of view. We provide a unifying view of previous approaches to query answering from inconsistent data, we shed light on the interaction between integrity constraint violation and the structure of repairs, and we study localization and factorization of consistent query answering. Specifically,

1) We present a formal framework for consistent query answering which is, to large extent, independent of a commitment to a specific definition of repair, but is based on a common setting of repair semantics: the repairs of the database are characterized by the minimal (non-preferred) databases from a space of candidate repairs with a preference order. Our setting generalizes previous proposals in the literature, such as set-inclusion based orderings [Fagin et al. 1983; Arenas et al. 1999; 2003; Barceló and Bertossi 2003; Bravo and Bertossi 2003; Calì et al. 2003a; 2003b; Chomicki and Marcinkowski 2005; Greco et al.

3

2003], cardinality-based orderings [Arenas et al. 2003; Lin and Mendelzon 1998], and weighted-based orderings [Lin 1996].

2) We investigate some locality properties for repairing inconsistent databases, aiming to isolate in the data those facts that will possibly be touched by a repair, called the "affected part" of the database and the facts that for sure will be not, called the "safe part" of the database. Specifically, we establish localization results for different classes of constraints:

- The first class, $\mathbf{C}_0$, contains all constraints of the form $\forall \vec{x} \alpha(\vec{x}) \supset \phi(\vec{x})$, where $\alpha(\vec{x})$ is a nonempty conjunction of atoms over database relations and $\phi(\vec{x})$ is a disjunction of built-in literals. These constraints are semantically equivalent to denial constraints [Chomicki et al. 2004a].
- The second class, $\mathbf{C}_1$, allows more general constraints of the form $\forall \vec{x} \alpha(\vec{x}) \supset \beta(\vec{x}) \vee \phi(\vec{x})$, where $\alpha(\vec{x})$ and $\phi(\vec{x})$ are as above and $\beta(\vec{x})$ is a disjunction of atoms over database relations.
- The third class, $\mathbf{C}_2$, has similar constraints $\forall \vec{x} \alpha(\vec{x}) \supset \beta(\vec{x}) \vee \phi(\vec{x})$; here $\alpha(\vec{x})$ may be empty, representing unconditional logical truth, but $\beta(\vec{x})$ may have at most one atom.
- The fourth class is the class of all universal constraints in clausal form. Thus, semantically, this class captures all universal constraints.

3) We propose a *repair localization approach* to query answering from inconsistent databases, in which the query is first evaluated over each of the repairs of the affected part only, augmented with the safe part, and then results are suitably recombined. Also, we investigate techniques for factorizing repairs into components that can be processed independently of each other. For some relevant settings, these techniques guarantee an exponential gain compared to the basic approach. Indeed, by looking at the actual forms of inconsistencies occurring in the underlying database instance, we identify settings for which the number of localized repairs to be considered for consistent query answering is linear in the size of the affected part, while generic database instances have data complexity from co-NP to $\Pi_2^P$ (see e.g., [Calì et al. 2003a; Chomicki and Marcinkowski 2005; Chomicki 2007; Greco et al. 2003]).

Secondly, our contribution is practical. Indeed, based on the above localization results, we develop strategies to consistent query answering relying on existing technologies offered by stable model engines and relational DBMSs. Resembling several proposals in the literature, our techniques make use of logic programs to solve inconsistency. However, we limit their usage to the affected part of the data. This approach is useful to localize the difficult part of the computation and to overcome the lack of scalability of current (yet still improving) implementations of stable model engines such as DLV [Leone et al. 2006] or Smodels [Simons et al. 2002]. Specifically:

4) We propose a formal model of inconsistency resolution via logic programming specification, which abstracts from several proposals in the literature [Arenas et al. 2003; Barceló and Bertossi 2003; Bertossi et al. 2002; Bravo and Bertossi 2003; Calì et al. 2003b; Greco et al. 2003]. Results obtained on this model are applicable to all such approaches.

5) We discuss an architecture that recombines the repairs of the affected part with the safe part of an inconsistent database, interleaving a stable model and a relational database engine. This is driven by the fact that database engines are geared towards efficient processing of large data sets, and thus help to achieve scalability. In this architecture, the database engine has to "update" the consistent answers to a certain query each time a new

repair is computed by the stable model engine. To further improve this strategy, a technique for simultaneously processing a (large) group of repairs in the DBMS is proposed. Basically, it consists in a marking and query rewriting strategy for compiling the reasoning tasks needed for consistent query answering into a relational database engine.

6) Finally, we assess the effectiveness of our approach in a suite of experiments. They have been carried out on a prototype implementation in which the stable model engine DLV is coupled with the DBMS PostgreSQL. The experimental results show that the implementation scales reasonably well.

We observe that our results on localization extend and generalize previous localization results which have been utilized (sometimes tacitly) for particular repair orderings and classes of constraints, for instance, for denial constraints and repairs which are closest to the original database measured by set symmetric difference [Chomicki et al. 2004a]. Also, our results can be exploited for efficient implementation of consistent query answering techniques in general, independent of a logic-programming based approach.

The rest of this paper is organized as follows. Section 2 introduces the notation for the relational data model and for logic programs used throughout the paper. Section 3 defines the formal framework for consistent query answering from inconsistent databases. Localization properties in database repairs and their exploitation to optimize consistent query answering are discussed in Section 4 and Section 5, respectively. The logic specification for consistent query answering is presented in Section 6, together with an architecture that interleaves a DBMS and a stable model engine. Section 7 considers other approaches to consistent query answering, and Section 8 reports results of our experimental activity. The final Section 9 concludes the paper.

Some further details of our techniques have been moved to an on-line appendix, which also contains further examples and experiments.

## 2. PRELIMINARIES

### 2.1 Data Model

We assume a countable infinite database domain $\mathcal{U}$ whose elements are referenced by constants $c_1$, $c_2$, ... under the *unique name assumption*, that is, different constants denote different real-world objects.

A *relational schema* (or simply *schema*) $\chi$ is a pair $\langle \Psi, \Sigma \rangle$, where:

- $\Psi$ is a finite set of relation (predicate) symbols, each with an associated positive arity.
- $\Sigma$ is a finite set of *integrity constraints* (ICs) expressed on the relation symbols in $\Psi$. We consider here universally quantified constraints [Abiteboul et al. 1995], i.e., first-order sentences of the form

$$\forall \vec{x} \, A_1(\vec{x}_1) \wedge \cdots \wedge A_l(\vec{x}_l) \supset B_1(\vec{y}_1) \vee \cdots \vee B_m(\vec{y}_m) \vee \phi_1(\vec{z}_1) \vee \cdots \vee \phi_n(\vec{z}_n), \quad (1)$$

where $l + m > 0$, $n \geq 0$, the $A_i(\vec{x}_i)$ and the $B_j(\vec{y}_j)$ are atoms over $\Psi$, the $\phi_k(\vec{z}_k)$ are atoms or negated atoms over possible built-in relations like equality ($=$), inequality ($\neq$), etc., $\vec{x}$ is a list of all variables occurring in the formula, and the $\vec{x}_i$, $\vec{y}_j$, and $\vec{z}_k$ are lists of variables from $\vec{x}$ and constants from $\mathcal{U}$.[1] The conjunction left of "$\supset$" is the *body* of the constraint, and the disjunction right of "$\supset$" its *head*.

---

[1] The condition $l + m > 0$ excludes constraints involving only built-in relations, which are irrelevant from a schema modeling perspective.

In the rest of the paper, $\chi = \langle \Psi, \Sigma \rangle$ denotes a relational schema. Since all variables in (1) are universally quantified, we omit quantifiers in constraints.

Note that (1) is a clausal normal form for arbitrary universal constraints on a relational schema. However, since existential quantification is not allowed, referential constraints such as *foreign-key constraints* cannot be expressed. We pay special attention to the following subclasses of constraints:

- Constraints with only built-in relations in the head (i.e., $m = 0$ in (1)). The class of these constraints, which we denote by $\mathbf{C}_0$, is a clausal normal form of *denial constraints* [Chomicki et al. 2004a], also called *generic constraints* in [Bertossi and Chomicki 2003]. This class (semantically) includes:
  - *key constraints* $p(\vec{x}, \vec{y}) \wedge p(\vec{x}, \vec{z}) \supset y_i{=}z_i$, for $1 \le i \le n$,
  - *functional dependencies* $p(\vec{x}, \vec{y}, \vec{v}) \wedge p(\vec{x}, \vec{z}, \vec{w}) \supset y_i{=}z_i$, for $1 \le i \le n$, and
  - *exclusion dependencies* $p_1(\vec{v}, \vec{y}) \wedge p_2(\vec{w}, \vec{z}) \supset y_1{\neq}z_1 \vee \cdots \vee y_n{\neq}z_n$,
  where $\vec{y} = y_1, \ldots, y_n$ and $\vec{z} = z_1, \ldots, z_n$.
- Constraints with non-empty body (i.e., $l > 0$ in (1)). We denote the class of these constraints, which permit conditional generation of tuples in the database, by $\mathbf{C}_1$. Note that $\mathbf{C}_0 \subseteq \mathbf{C}_1$ (since $l + m > 0$). The class $\mathbf{C}_1$ includes, for instance, *inclusion dependencies* of the form $p_1(\vec{x}) \supset p_2(\vec{x})$.
- Constraints with at most one database atom in the head (i.e., $m \le 1$ in (1)). We denote the class of these constraints, which we call non-disjunctive, by $\mathbf{C}_2$. Beyond denials, such constraints also allow to enforce the (unconditional) presence of a tuple. Parts of the database may be protected from modifications in this way. Note that $\mathbf{C}_0 \subseteq \mathbf{C}_2$, while $\mathbf{C}_1$ and $\mathbf{C}_2$ are incomparable.

**Example 2.1** In our example, the schema $\chi_0$ is the tuple $\langle \Psi_0, \Sigma_0 \rangle$, where $\Psi_0$ consists of the ternary relation symbols *player*, *team*, and *coach*, and $\Sigma_0$ can be defined as follows:

$$\begin{aligned}
\sigma_1 &: player(x, y, z) \wedge player(x, y', z) \supset y{=}y', \\
\sigma_2 &: team(x, y, z) \wedge team(x, y', z') \supset y{=}y', \\
\sigma_3 &: team(x, y, z) \wedge team(x, y', z') \supset z{=}z', \\
\sigma_4 &: coach(x, y, z) \wedge coach(x, y', z) \supset y{=}y', \\
\sigma_5 &: coach(x, y, z) \wedge player(x', y', z) \supset x{\neq}x', \\
\sigma_6 &: coach(x, y, z) \wedge team(z, y', x') \supset x{\neq}x'.
\end{aligned}$$

Here $\sigma_1$–$\sigma_4$ are key constraints, while $\sigma_5$ and $\sigma_6$ encode that, for any given team, the coach is neither a player nor a team leader. Note that all these constraints are in $\mathbf{C}_0$. □

For a set of relation symbols $\Psi$ as above, $\mathcal{F}(\Psi)$ denotes the set of all facts $r(\vec{t})$, where $r \in \Psi$ has arity $n$ and $\vec{t} = (c_1, \ldots, c_n) \in \mathcal{U}^n$ is an $n$-tuple of constants from $\mathcal{U}$. A *database instance* (or simply *database*) for $\Psi$ is any finite set $D \subseteq \mathcal{F}(\Psi)$. The extension of relation $r$ in $D$ is the set of tuples $r^D = \{\vec{t} \mid r(\vec{t}) \in D\}$. We denote by $D(\Psi)$ the set of all databases for $\Psi$. For any relation schema $\chi = \langle \Psi, \Sigma \rangle$, in abuse of notation, $\mathcal{F}(\chi)$ and $D(\chi)$ denote $\mathcal{F}(\Psi)$ and $D(\Psi)$, respectively, and a database for $\chi$ is a database for $\Psi$.

A constraint $\sigma$ is *ground*, if it is variable-free. For any such $\sigma$, $facts(\sigma)$ denotes the set of all facts $p(\vec{t}) \in \mathcal{F}(\chi)$ occurring in $\sigma$, and for any set $\Sigma$ of ground constraints, $facts(\Sigma) = \bigcup_{\sigma \in \Sigma} facts(\sigma)$. For any constraint $\sigma = \alpha(\vec{x})$, we denote by $ground(\sigma)$ the set of its *ground instances* $\theta(\alpha(\vec{x}))$, where $\theta$ is any substitution of the variables $\vec{x}$ by constants from $\mathcal{U}$. For any set of constraints $\Sigma$, $ground(\Sigma) = \bigcup_{\sigma \in \Sigma} ground(\sigma)$.

Given $D \subseteq \mathcal{F}(\Psi)$, where $\Psi = \{r_1, \ldots, r_n\}$, $D$ *satisfies* a constraint $\sigma$, denoted $D \models \sigma$, if $\sigma$ is true on the relational structure $(\mathcal{U}, r_1^D, \ldots, r_n^D, c_1^D, c_2^D, \ldots)$ where $c_i^D = c_i$, for all $c_i \in \mathcal{U}$ (i.e., each $\sigma' \in ground(\sigma)$ evaluates to true), and *violates* $\sigma$ otherwise; $D$ *satisfies* (or is *consistent with*) a set of constraints $\Sigma$, denoted $D \models \Sigma$, if $D \models \sigma$ for every $\sigma \in \Sigma$, and *violates* $\Sigma$ otherwise. Finally, a relational schema $\chi = \langle \Psi, \Sigma \rangle$ is *consistent*, if there exists a database $D$ for $\chi$ that is consistent with $\Sigma$, otherwise $\chi$ is inconsistent.

**Example 2.2** Consider the constraint $\sigma_2$ in $\Sigma_0$, and its ground instance

$$team(RM, Roma, 10) \land team(RM, Real\ Madrid, 10) \supset Roma{=}Real\ Madrid.$$

Clearly, this instance does not evaluate true on the relational structure associated with $D_0$, which therefore violates $\Sigma_0$. □

## 2.2 Datalog$^{\lor,\lnot}$ Programs and Queries

*Syntax.* A *Datalog$^{\lor,\lnot}$* rule $\rho$ is an expression of the form

$$a_1 \lor \ldots \lor a_n \leftarrow b_1, \ldots, b_k,\ not\ b_{k+1}, \ldots,\ not\ b_{k+m} \tag{2}$$

where $a_i$, $b_j$ are atoms in a relational first-order language $\mathcal{L}$, "*not*" is *negation as failure*, and "," is conjunction. If $k = m = 0$, then $\rho$ is a *fact* and "$\leftarrow$" is omitted. The part left of "$\leftarrow$" is the *head* of $\rho$, denoted $head(\rho)$, and the part right of "$\leftarrow$" the *body* of $\rho$, denoted $body(\rho)$. We assume *safety*, i.e., each variable occurring in $\rho$ occurs in some $b_i$, $1 \le i \le k$, whose predicate is not a built-in relation. Built-in relations may occur only in the body.

A *Datalog$^{\lor,\lnot}$ program* $\mathcal{P}$ is a finite set of Datalog$^{\lor,\lnot}$ rules. Important restrictions are *normal programs*, Datalog$^{\lnot}$, where $n = 1$ for all rules, *stratified normal* programs, Datalog$^{\lnot s}$, and *non-recursive* programs, defined as follows. Each Datalog$^{\lnot}$ program $\mathcal{P}$ has a *dependency graph* $G(\mathcal{P}) = \langle V, E \rangle$, where $V$ are the predicates occurring in $\mathcal{P}$ and $E$ contains an arc $r \to s$ if $r$ occurs in $head(\rho)$ and $s$ in $body(\rho)$ for some rule $\rho \in \mathcal{P}$. Moreover, if $s$ occurs under negation, the arc is labeled with '$*$.' Then $\mathcal{P}$ is *stratified*, if $G(\mathcal{P})$ has no cycle with an arc labeled '$*$,' and *non-recursive*, if $G(\mathcal{P})$ is acyclic.

*Semantics.* The semantics of a Datalog$^{\lor,\lnot}$ program $\mathcal{P}$ is defined via its *grounding* $ground(\mathcal{P})$ w.r.t. $\mathcal{L}$ (usually, the language generated by $\mathcal{P}$), which consists of all ground instances of rules in $\mathcal{P}$ possible with constant symbols from $\mathcal{L}$. Let $B_{\mathcal{L}}$ be the set of all ground atoms with a predicate and constant symbols in $\mathcal{L}$. A *(Herbrand) interpretation for* $\mathcal{P}$ is any subset $I \subseteq B_{\mathcal{L}}$; an atom $p(\vec{c}) \in B_{\mathcal{L}}$ is true in $I$, if $p(\vec{c}) \in I$, and false in $I$ otherwise. A ground rule (2) is *satisfied* by $I$, if either some $a_i$ or $b_{k+j}$ is true in $I$, or some $b_i$, $1 \le i \le k$, is false in $I$. Finally, $I$ is a model of $\mathcal{P}$, if $I$ satisfies all rules in $ground(\mathcal{P})$.

The *stable model semantics* [Gelfond and Lifschitz 1991] assigns *stable models* to any Datalog$^{\lor,\lnot}$ program $\mathcal{P}$ as follows. If $\mathcal{P}$ is "*not*"-free, its stable models are its minimal models, where a model $M$ of $\mathcal{P}$ is minimal, if no $N \subset M$ is a model of $\mathcal{P}$. If $\mathcal{P}$ has negation, $M$ is a stable model of $\mathcal{P}$, if $M$ is a minimal model of the *reduct* $\mathcal{P}$ w.r.t. $M$, which results from $ground(\mathcal{P})$ by deleting *(i)* each rule $\rho$ with a literal *not* $p(\vec{c})$ in the body such that $p(\vec{c}) \in M$, and *(ii)* the negative literals from all remaining rules.

We denote by $\mathrm{SM}(\mathcal{P})$ the set of stable models of $\mathcal{P}$. Note that for "*not*"-free programs, minimal models and stable models coincide, and that positive disjunction-free (resp. stratified) programs have a unique stable model [Gelfond and Lifschitz 1991].

7

*Queries.* A *Datalog$^{\vee,\neg}$ query* $Q$ over a schema $\chi = \langle \Psi, \Sigma \rangle$ is a pair $\langle q, \mathcal{P} \rangle$, where $\mathcal{P}$ is a Datalog$^{\vee,\neg}$ program such that every $p \in \Psi$ occurs in $\mathcal{P}$ only in rule bodies, and $q$ occurs in some rule head of $\mathcal{P}$ but not in $\Psi$. The *arity of* $Q$ is the arity of $q$. Given any database $D$ for $\chi$, the *evaluation of* $Q$ *over* $D$, is $Q[D] = \{(c_1, \dots, c_n) \mid q(c_1, \dots, c_n) \in M$, for each $M \in SM(\mathcal{P} \cup D)\}$. Note that as for $Q$, any non-recursive $\mathcal{P}$ can be rewritten to a *union of conjunctive queries*, i.e., a set of rules (2) where $n = 1$ and $m = 0$, with the same head predicate $q$ which does not occur in rule bodies. For further background on Datalog$^{\vee,\neg}$ and queries, see [Abiteboul et al. 1995; Eiter et al. 1997].

**Example 2.3** In our ongoing example, we may consider a query $Q$ that asks for the codes of all players and team leaders, and that is formally written as $Q = \langle q, \mathcal{P} \rangle$ where $\mathcal{P} = \{q(x) \leftarrow player(x, y, z), q(x) \leftarrow team(v, w, x)\}$. $Q$ has arity 1. Note that $\mathcal{P}$ is a union of conjunctive queries. $\square$

## 3. CONSISTENT QUERY ANSWERING FRAMEWORK

### 3.1 A General Framework for Database Repairs

Let us assume that $\chi = \langle \Psi, \Sigma \rangle$ is given together with a (possibly inconsistent) database $D$ for $\chi$. Following a common approach in the literature on inconsistent databases [Arenas et al. 1999; Greco et al. 2003; Calì et al. 2003a; Chomicki 2007], we next define the semantics of querying $D$ in terms of its *repairs*. Specifically, we present a generalization of previous approaches where the way of repairing a database is chosen according to an arbitrary preorder on databases satisfying some conditions.

We suppose that $\leq_D$ is a (fixed) preorder (i.e., a reflexive and transitive binary relation) on $D(\chi)$, and denote by $<_D$ the induced preference order (i.e., an irreflexive and transitive binary relation) given by $R_1 <_D R_2$, if $R_1 \leq_D R_2 \wedge R_2 \not\leq_D R_1$. We call $R_1 <_D$-*preferred* to $R_2$ in this case. A repair for $D$ is now defined in terms of a minimal element under $<_D$.

**Definition 3.1 (Repair)** Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, and let $\leq_D$ be a (fixed) preorder on $D(\chi)$. A database $R \in D(\chi)$ is a $\leq_D$-*repair* (simply, *repair*) *for* $D$ *w.r.t.* $\chi$, if

(1) $R \models \Sigma$, and
(2) $R$ is minimal in $D(\chi)$ w.r.t. $<_D$, i.e., there is no $R' \in D(\chi)$ such that $R' \models \Sigma$ and $R'$ is $<_D$-preferred to $R$.

The set of all repairs for $D$ w.r.t. $\chi$ is denoted by $rep_{\chi}^{\leq_D}(D)$. When clear from the context, the subscript $\chi$ may be dropped. Similarly, the superscript $\leq_D$ is omitted. $\square$

The definition of repair relies on a general notion of preorder on databases. The method for consistent query answering presented in the following is based on abstract properties of the induced preference order, which we refer to as set inclusion proximity, disjoint preference expansion and disjunctive split. The property of *set inclusion proximity* is as follows:

(SIP) For any databases $R_1, R_2$, and $D$, $\triangle(R_1, D) \subset \triangle(R_2, D)$ implies $R_1 <_D R_2$,

where $\triangle(A, B) = (A \setminus B) \cup (B \setminus A)$ is symmetric set difference. Informally, this property effects that a database $R$ satisfying the constraints can be a repair only if there is no way to establish consistency with $\Sigma$ by touching merely a strict subset of facts compared to $R$.

The properties *disjoint preference expansion* and *disjunctive split* are as follows:

(DPE) If $R_1 <_{D_1} R_1'$ and $R_2, D_2$ are disjoint from $R_1, R_1'$, and $D_1$ (i.e., $(R_1 \cup R_1' \cup D_1) \cap (R_2 \cup D_2) = \emptyset$), then $R_1 \cup R_2 <_{D_1 \cup D_2} R_1' \cup R_2$.

(DIS) If $R_1 <_D R_2$, then for every database $R$ it holds that either $R_1 \cap R <_{D \cap R} R_2 \cap R$ or $R_1 \setminus R <_{D \setminus R} R_1' \setminus R$ (or both).

Loosely speaking, (DPE) says that preference must be invariant under adding new facts, while (DIS) says that preference must uniformly stem from disjoint "components."

The prototypical preorder $\leq_D$ is given by $R_1 \leq_D R_2$ iff $\triangle(R_1, D) \subseteq \triangle(R_2, D)$ [Arenas et al. 1999; 2003; Barceló and Bertossi 2003; Bravo and Bertossi 2003; Chomicki et al. 2004a; Greco et al. 2003; Fuxman and Miller 2007]. Intuitively, each repair of $D$ is then obtained by properly adding and deleting facts from $D$ in order to satisfy constraints in $\Sigma$, as long as we "minimize" such changes. The following proposition is easy to prove.

**Proposition 3.1** *The prototypical preorder satisfies properties (SIP), (DPE), and (DIS).*

In our examples we refer to the prototypical preorder. Notice that a variety of repair semantics are either defined in terms of a preorder satisfying the above properties or can be characterized by such a preorder, beside those based on the prototypical preorder discussed above, including set-inclusion based ordering [Fagin et al. 1983; Calì et al. 2003a], cardinality-based ordering [Arenas et al. 2003; Lin and Mendelzon 1998], weight-based orderings [Lin 1996], as well as refinements with priority levels. An interesting special case of weight-based ordering is the lexicographic preference, where $R_1$ is preferred to $R_2$ w.r.t. $D$ if the first fact in a total ordering of $\mathcal{F}(\chi)$ on which $R_1$ and $R_2$ repair $D$ differently belongs to $R_2$. However, we point out that our methods and results for query answering can also be extended to other preference orderings under certain conditions (see Section 9).

## 3.2 Constructible Repairs and Safe Constraints

An important aspect is that constraints might enforce that *any* set of facts $R$ for $\chi = \langle \Psi, \Sigma \rangle$ which satisfies $\Sigma$ must be infinite, and thus $\chi$ is inconsistent, i.e., no $D \in D(\chi)$ satisfies $\Sigma$. A simple example is where $\Sigma = \{ \forall x \, p(x) \}$. Semantically, this is commonly avoided by requesting domain-independence of constraints [Ullman 1989], which syntactically is ensured by *safety*, i.e., each variable occurring in the head of a constraint must also occur in its body. Notice that major classes of constraints including key constraints, functional dependencies, exclusion dependencies, inclusion dependencies of the form $p_1(\vec{x}) \supset p_2(\vec{x})$, or denial constraints fulfill safety. Together with (SIP), safety of constraints ensures that any database $D$ has a repair if this is possible at all. For any $R \subseteq \mathcal{F}(\chi)$, we denote by $adom(R, \chi)$ the *active domain* of $R$ and $\chi$, i.e., the set of constants occurring in $R$ and $\Sigma$.

**Proposition 3.2** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, where all constraints in $\Sigma$ are safe. Suppose that $<_D$ satisfies (SIP). Then, every repair $R \in rep(D)$ involves only constants from $adom(D, \chi)$, and some repair exists if $\chi$ is consistent.*

Notice that, for an arbitrary preference order, no repair may exist, even if $\chi$ is consistent. In the rest of this article, unless stated otherwise, we assume safe constraints.

Finite repairs can also be ensured for unsafe constraints in which variables violating safety are guarded by built-in relations, such as for $D = \emptyset$ w.r.t. $\chi = \langle \{p\}, \{p(x) \lor x > 100\} \rangle$, assuming that $\mathcal{U}$ are the natural numbers. As this example shows, repairs may

go beyond the active domain. However, this is prevented if built-ins involve only equality and inequality. We have a result similar to Proposition 3.2 (cf. Appendix A for proofs).

**Proposition 3.3** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$ where no built-in relations occur in $\Sigma$ except $=$ and $\neq$. Suppose that $<_D$ satisfies (SIP). Then, every repair $R \in rep(D)$ involves only constants from $adom(D, \chi)$, and some repair exists if $\chi$ is consistent.*

### 3.3 Queries and Consistent Answers

The notion of repair is crucial for the definition of the semantics of querying inconsistent databases. We conclude this section by formalizing this aspect.

**Definition 3.2** Let $Q$ be a non-recursive Datalog$^\neg$ query. For any database $D \in D(\chi)$, the set of *consistent answers to $Q$ w.r.t.* $D$ is the set of tuples $ans_c(Q, D) = \{\vec{t} \mid \vec{t} \in Q[R], \text{ for each } R \in rep_\chi(D)\}$. □

Informally, a tuple $\vec{t}$ is a consistent answer if it is a consequence under standard certainty semantics for each possible repair of the database $D$. Note that in real applications, a query language subsumed by non-recursive Datalog$^\neg$ is often adopted.

**Example 3.1** Recall that in our scenario, repairs for the database $D_0$ for $\chi_0$ are shown in Fig. 1. For the query $Q = \langle q, \mathcal{P} \rangle$, where $\mathcal{P} = \{q(x) \leftarrow player(x, y, z), q(x) \leftarrow team(v, w, x)\}$, we thus obtain $ans_c(Q, D_0) = \{(8), (9), (10)\}$. For the query $Q' = \langle q, \{q(y) \leftarrow team(x, y, z)\} \rangle$, we have $ans_c(Q', D_0) = \{(Barcelona)\}$, while for $Q'' = \langle q', \{q'(x, z) \leftarrow team(x, y, z)\} \rangle$, we have $ans_c(Q'', D_0) = \{(RM, 10), (BC, 8)\}$. □

## 4. LOCALITY PROPERTIES FOR REPAIRING INCONSISTENT DATABASES

In this section, we investigate how to localize inconsistency in a given database $D$, that is, how to narrow down the set of facts in $D$ to a part which is "affected" by inconsistency and repair, and how to obtain the repairs of $D$ from the repairs of this affected part. To this end, we introduce the notion of a *repair envelope*. Informally, a repair envelope is a set of facts $E$ such that the repairs of $D$ touch only facts in $E$ and are given by the repairs of $D \cap E$ plus the "unaffected" ("safe") part of $D$, i.e., the portion of $D$ which is outside the envelope. More formally:

**Definition 4.1** Let $D$ be a database for a relational schema $\chi = \langle \Psi, \Sigma \rangle$. A set $E$ of facts over $\Psi$ is a repair envelope for $D$ if it fulfills the following conditions:

$$\triangle(R, D) \subseteq E, \text{ for all } R \in rep_\chi(D), \tag{3}$$

$$rep_\chi(D) = \{R \cup (D \setminus E) \mid R \in rep_\chi(D \cap E)\}. \tag{4}$$

The repairs of $D$ can then be fully localized to the repairs of $D \cap E$, which in practice may be much smaller than $D$. In fact, as will be shown in this section, for constraints $\mathbf{C}_0$ the set of all facts involved in constraint violations, denoted $C$ (formally defined in the beginning of Section 4.1), is always a repair envelope, and for constraints $\mathbf{C}_1$ and $\mathbf{C}_2$, a closure $C^*$ of $C$ under syntactic conflict propagation, i.e., under co-occurrence of facts in violated constraints (cf. Definition 4.4), is a repair envelope. Such a closure, as we will explain in detail in the following, takes care of facts that "indirectly" participate in constraint violations. In general, however, a repair envelope needs not be a superset of $C^*$ or $C$. Figure 2 shows the different sets.

**Example 4.1** Recall that $team(RM, Roma, 10) \wedge team(RM, Real\ Madrid, 10) \supset Roma = Real\ Madrid$ witnesses in Example 2.2 a violation of the key of $team$; it is the only ground constraint violated by $D_0$. Here, $C = \{team(RM, Roma, 10), team(RM, Real\ Madrid, 10)\}$, and since the constraints are of type $\mathbf{C}_0$, it is a repair envelope for $D$. The database $D \cap C = C$ has the two repairs $R_1 = \{team(RM, Roma, 10)\}$ and $R_2 = \{team(RM,\ Real\ Madrid, 10)\}$; therefore, according to (4), $D$ has the two repairs $R_1 \cup (D_0 \setminus C)$ and $R_2 \cup (D_0 \setminus C)$, which are those shown in Figure 1. $\square$

The following example shows that taking only $C$ into account is not always sufficient.

**Example 4.2** Let $D = \{s(a)\}$ for $\chi = \langle \Psi, \Sigma \rangle$, where $\Sigma = \{s(a) \supset r(a), r(a) \supset p(a) \vee q(a), r(a) \wedge p(a) \supset a \neq a\}$. In this case $C^* = \{s(a), r(a), p(a), q(a)\}$ and the constraints are of type $\mathbf{C}_1$, hence, $C^*$ is a repair envelope for $D$. The database $D \cap C^* = D$ has the two repairs $R_1 = \emptyset$ and $R_2 = \{s(a), r(a), q(a)\}$. Note that $\triangle(R_2, D) = \{r(a), q(a)\} \nsubseteq \{s(a), r(a)\} = C$. Therefore, $C$ violates (3) and is not a repair envelope.

Note that a repair envelope always exists, since the set of all facts is a trivial repair envelope. As for localizing the computation of $rep(D)$, only Condition (4) is relevant (if $E$ satisfies it, then so does every $E'$ such that $E' \cap D = E \cap D$, in particular $E' = E \cap D$). Condition (3), however, allows to bound the answer to certain queries. In particular, for monotone queries $Q$, we have that $Q[D \setminus E] \subseteq ans(Q, D) \subseteq Q[D \cup E]$.

For general constraints, $C^*$ is not always a repair envelope. However, we show that it is a *weak repair envelope*:

**Definition 4.2** Let $D$ be a database for a relational schema $\chi = \langle \Psi, \Sigma \rangle$. A set $E$ of facts over $\Psi$ is a weak repair envelope for $D$ if it fulfills Condition (3) and instead of (4), the relaxed equation
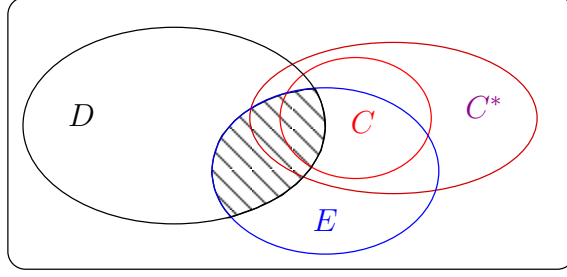
$$rep_\chi(D) = \{(R \cap E) \cup (D \setminus E) \mid R \in rep_\chi(D \cap E)\} \tag{5}$$

That is, the repairs of $D$ are obtained by constraining the repairs of $D \cap E$ to the repair envelope. This is necessary since facts outside the envelope might be added to such repairs (see Example 4.4). However, this can only occur in presence of certain disjunctions.

Despite the difference that $E$ is either a repair envelope or a weak repair envelope, we define affected database and safe database w.r.t. $E$ as follows.

**Definition 4.3** Let $D$ be a database for a relational schema $\chi = \langle \Psi, \Sigma \rangle$ and $E$ a (weak) repair envelope. Then $D \cap E$ is the affected part of $D$ (w.r.t. $\chi$), or simply the "affected database", and $D \setminus E$ is the safe part of $D$ (w.r.t. $\chi$), or simply the "safe database". $\square$

Note that we mostly refer to specific kinds of envelopes in the paper, namely the conflict set $C$ and the conflict closure $C^*$. We proceed as follows. After formally defining $C$ and $C^*$ and establishing some auxiliary results, we show that $C^*$ is a weak repair envelope in general. We then prove that it is a repair envelope under restrictions, in particular for $\mathbf{C}_1$ and $\mathbf{C}_2$ constraints. This envelope may be further decreased. Indeed, we prove that $C$ is a repair envelope for $\mathbf{C}_0$ constraints. In fact, the results for special constraints are stronger and establish 1-1 correspondences between repairs of $D \cap E$ and repairs of $D$. Some of the proofs are omitted here but can be found in Appendix B.

*Conflict set $C$* : facts occurring in $ground(\Sigma)$ violated in $D$
*Conflict closure $C^*$* : syntactic conflict propagation from $C$ by $\Sigma$
*Repair envelope $E$* : safe bound on tuple changes with local repairs (hatched)

Fig. 2.    Localization of database repair

## 4.1   General Constraints

Let $D$ be a database for a relational schema $\chi = \langle \Psi, \Sigma \rangle$. The *conflict set* for $D$ w.r.t. $\chi$ is the set of facts $C_\chi(D) = \{p(\vec{t}) \mid \exists \sigma \in ground(\Sigma), p(\vec{t}) \in facts(\sigma), D \not\models \sigma\}$, i.e., $C_\chi(D)$ is the set of facts occurring in the ground instances of $\Sigma$ which are violated by $D$. In the following, if clear from the context, $D$ and/or the subscript $\chi$ will be dropped.

Figure 2 shows that the conflict set may contain both facts in $D$ (as in Example 4.1) and facts in $\mathcal{F}(\chi)$ that do not belong to $D$. For example, let $D = \{p(a)\}$, and let $\chi$ contain the dependency $p(x) \supset q(x)$. Then $C = \{p(a), q(a)\}$.

For defining conflict propagation, we first introduce the following notion. Two facts $p(\vec{t}), p'(\vec{t'})$ in $\mathcal{F}(\chi)$ are *constraint-bounded in $\chi$*, if there exists some $\sigma \in ground(\Sigma)$ such that all constants occurring in $facts(\sigma)$ are from $adom(D, \chi)$, and $\{p(\vec{t}), p'(\vec{t'})\} \subseteq facts(\sigma)$. (Note that by assumed safety of constraints and the results of Section 3.2, we only need to consider $adom(D, \chi)$.) We now generalize the notion of conflict set.

**Definition 4.4 (Conflict closure)** Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$. Then, the *conflict closure* for $D$, denoted by $C_\chi^*(D)$, is the least set $A \supseteq C_\chi(D)$ which contains every fact $p(\vec{t})$ constraint-bounded in $\chi$ with some fact $p'(\vec{t'}) \in A$.   $\square$

We omit $D$ and/or the subscript $\chi$ if clear from the context. Intuitively, $C^*$ contains, besides facts from $C$, facts which possibly must be touched by repair in turn to avoid new inconsistency with $\Sigma$ caused by previous repairing actions. For example, assume that $\chi$ contains the constraints $p(x) \supset q(x)$ and $q(x) \supset s(x)$. Then, for $D = \{p(a)\}$, we have that $C = \{p(a), q(a)\}$ and $C^* = C \cup \{s(a)\}$. As shown in Figure 2, $C^*$ may add to $C$ both facts inside and outside $D$, but may also coincide with $C$, as in our example if $s(a)$ was in $D$.

Towards a proof that $C^*$ is a weak repair envelope, we need some preliminary technical results. For $D$ and $\chi = \langle \Psi, \Sigma \rangle$, consider the following two sets of ground constraints:

*(i)* $\Sigma_\chi^a(D) = \{\sigma \in ground(\Sigma) \mid facts(\sigma) \cap C^* \neq \emptyset\}$ consists of all ground constraints in which at least one fact from $C^*$ occurs;

*(ii)* $\Sigma_\chi^s(D) = \{\sigma \in ground(\Sigma) \mid facts(\sigma) \not\subseteq C^*\}$ consists of all ground constraints in which at least one fact occurs which is *not* in $C^*$.

As usual, $\chi$ and/or $D$ will be omitted. We first show that $\Sigma^a \cup \Sigma^s$ is a special partitioning of $ground(\Sigma)$.

**Proposition 4.1  (Separation)** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$. Then, (1) $facts(\Sigma^a) = C^*$, (2) $facts(\Sigma^s) \cap C^* = \emptyset$, (3) $\Sigma^a \cap \Sigma^s = \emptyset$, and (4) $\Sigma^a \cup \Sigma^s = ground(\Sigma)$.*

The separation property allows us to shed light on the structure of repairs.

**Proposition 4.2  (Safe database)** *Let $D$ be any database for $\chi = \langle \Psi, \Sigma \rangle$. Then, for each repair $R \in rep(D)$ it holds that $R \setminus C^* = D \setminus C^*$.*

Informally, the above proposition shows that $D \setminus C^*$ is a safe portion of $D$, in the sense that tuples of $D$ outside the conflict closure will not be touched by repair.

Prior to the main result of this subsection, we establish the following lemma:

**Lemma 4.3** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, and let $A = D \cap C^*$ and $\chi^a = \langle \Psi, \Sigma^a \rangle$. Then, for each $S \subseteq D \setminus C^*$, the following holds:*

(1)  *for each $R \in rep_\chi(A \cup S)$, it holds that $R \cap C^* \in rep_{\chi^a}(A)$;*

(2)  *for each $R \in rep_{\chi^a}(A)$ there exists a set of facts $S' \subseteq \mathcal{F}(\chi)$ such that $S' \cap C^* = \emptyset$, and $(R \cup S') \in rep_\chi(A \cup S)$.*

In other words, item (1) in the lemma above shows how to obtain a repair of the database $A = D \cap C^*$ w.r.t. $\chi^a$, from a repair, computed w.r.t. $\chi$, of $A$ augmented with any subset $S$ of the safe database $D \setminus C^*$. Conversely, item (2) shows how to obtain a repair of $A \cup S$ w.r.t. $\chi$, from a repair of $A$ w.r.t. $\chi^a$. Notice that repairing $A$ w.r.t. $\chi^a$, and not w.r.t. $\chi$, is necessary for the lemma above to hold, since for a repair $R \in rep_\chi(A \cup S)$, it does not hold in general that $R \cap C^* \in rep_\chi(A)$. Also, repairing $A$ w.r.t. $\chi^a$ avoids repairing constraints in $\Sigma^s$ not satisfied by $A$.

Armed with the above concepts and results, we state the main theorem of this subsection.

**Theorem 4.4** *Every database $D$ for $\chi = \langle \Psi, \Sigma \rangle$ has $C^*$ as a weak repair envelope.*

PROOF.  We first show that for each $R \in rep_\chi(D)$ then $\triangle(R, D) \subseteq C^*$, as specified by condition (3) in Section 4, where we pose $E = C^*$. Assume by contradiction that there exists a fact $f \in \triangle(R, D)$ such that $f \notin C^*$. By Proposition 4.1 it follows that there exists no $\sigma \in \Sigma^a$ such that $f \in facts(\sigma)$. Then, if $f \in R \setminus D$, it is easy to see that $R \setminus \{f\} \models \Sigma$, but by property (SIP) we have that $R \setminus \{f\} <_D R$, thus contradicting the assumption that $R \in rep_\chi(D)$. Analogously, if $f \in D \setminus R$, it is easy to see that $R \cup \{f\} \models \Sigma$, but by property (SIP) we have that $R \cup \{f\} <_D R$, thus again contradicting the assumption.

We now prove that $rep_\chi(D)$ coincides with the set defined by Equation (5) in Section 4, where we pose $E = C^*$. To this aim, we show that (i) for every $R \in rep_\chi(D)$, there exists some $R' \in rep_\chi(D \cap C^*)$ such that $R = (R' \cap C^*) \cup (D \setminus C^*)$, and (ii) for every $R \in rep_\chi(D \cap C^*)$ there exists some $R' \in rep_\chi(D)$ such that $R' = (R \cap C^*) \cup (D \setminus C^*)$.

(i) We first apply Item 1 of Lemma 4.3 for $S = D \setminus C^*$ and obtain $R \cap C^* \in rep_{\chi^a}(D \cap C^*)$ (notice that in Lemma 4.3 $A = D \cap C^*$, $\chi^a = \langle \Psi, \Sigma^a \rangle$, and since we pose $S = D \setminus C^*$, we have that $A \cup S = D$). We then apply Item 2 for $S = \emptyset$, and we obtain that there exists $S'$ such that $S' \cap C^* = \emptyset$ and $R' = (R \cap C^*) \cup S' \in rep_\chi(A)$. Since

13

$S' \cap C^* = \emptyset$, we also have that $R' \cap C^* = R \cap C^*$, and from Proposition 4.2, it follows that $R = (R \cap C^*) \cup (D \setminus C^*)$. Therefore, $R = (R' \cap C^*) \cup (D \setminus C^*)$.

(ii) Similarly, applying first Item 1 of Lemma 4.3 for $S = \emptyset$ and then Item 2 for $S = D \setminus C^*$, we obtain that there exists $S'$ such that $S' \cap C^* = \emptyset$, and $R' = (R \cap C^*) \cup S' \in rep_\chi(D)$. Then, from Proposition 4.2, it follows that $S' = D \setminus C^*$. We thus easily obtain that $R' = (R \cap C^*) \cup (D \setminus C^*)$.  $\square$

For computing repairs for an inconsistent database $D$, we can thus proceed as follows:
  1. compute the conflict closure $C^*$ (w.r.t. $\Sigma$);
  2. compute the repairs of $A = D \cap C^*$ (w.r.t. $\Sigma$);
  3. intersect each repair obtained with $C^*$; and
  4. for each such set, take the union with $D \setminus C^*$.

In fact, as shown below, it is sufficient to consider $\Sigma^a$ in Step 2 instead of $\Sigma$. We note that $C^*$ is polynomially computable w.r.t. data volume from $C$ by transitive closure computation (simply use a Datalog program over $C$), but omit discussing efficient algorithms.

A drawback of the approach in general, however, is that in Step 2, facts outside $C^*$ might be included in a repair of $A$, which are stripped off subsequently in Step 3.

**Example 4.3** Consider $D = \{p(a)\}$ for $\chi = \langle \Psi, \{p(a), q(a)\} \rangle$. In this case, $C = C^* = \{q(a)\}$, $A = D \cap C^* = \emptyset$, and $D \setminus C^* = D$. We have $rep(A) = \{\{p(a), q(a)\}\}$ and $\{p(a), q(a)\} \cap C^* = \{q(a)\}$; $p(a)$ is stripped off from the repair of $A$.  $\square$

In this example, the repair of $A$ added a fact outside $C^*$ but from the safe part of $D$, which doesn't hurt. The following example shows that facts outside $C^* \cup D$ may be added.

**Example 4.4** Consider $D = \{r(a), p(a)\}$ for $\chi = \langle \Psi, \Sigma \rangle$, where $\Sigma = \{r(a) \supset p(a) \vee q(a), r(a), s(a)\}$. Then, $C^* = \{s(a)\}$ and $D \cap C^* = \emptyset$ has two repairs, viz. $R_1 = \{r(a), s(a), p(a)\}$ and $R_2 = \{r(a), s(a), q(a)\}$. Note that if $C^*$ were a repair envelope, then according to Condition (4), $R_2 \cup (D \setminus C^*) = \{r(a), s(a), p(a), q(a)\}$ would have to be a repair of $D$, which is incorrect. Note also, that this time $q(a)$ has been added in $R_2$ which was neither in $D$ nor in $C^*$.  $\square$

We remark that in Example 4.4, $\Sigma$ contains constraints from both the classes $\mathbf{C}_1$ and $\mathbf{C}_2$, but not from a single class. As we show in the next subsection, the effects in Example 4.4 can not happen under restriction to a single class, and $C^*$ is always a repair envelope.

We finally provide the result below that follows from Theorem 4.4, and remarks that repairing basically depends on $\Sigma^a$.

**Corollary 4.5** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, and let $\chi' = \langle \Psi, \Sigma' \rangle$ be such that $\Sigma^a_{\chi'}(D) = \Sigma^a_\chi(D)$. Then $rep_\chi(D) = rep_{\chi'}(D)$.*

Then, we can modify or prune constraints "outside" $\Sigma^a$ in arbitrary manner, e.g., for optimization purposes. As we show in the next subsection, this makes $C^*$ a repair envelope, rather than a weak repair envelope, in several cases in which $\Sigma$ contains general constraints.

### 4.2 Special Constraints

In this section, we consider the constraint classes $\mathbf{C}_i$ which have been introduced in Section 2, and determine repair envelopes for them.

**4.2.1 Constraints $\mathbf{C}_1$ and $\mathbf{C}_2$.** Recall that $\mathbf{C}_1$ constraints have nonempty bodies, and thus cannot unconditionally enforce the inclusion of facts to a database instance.

**Proposition 4.6** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$ such that $\Sigma \subseteq \mathbf{C}_1$. Then, each repair $R$ of $A = D \cap C^*$ w.r.t. $\chi$ satisfies $R \subseteq C^*$.*

PROOF. By Item 1 of Lemma 4.3, for $S = \emptyset$, each $R \in rep_\chi(A)$ gives rise to a repair $R' = R \cap C^*$ of $A$ w.r.t. $\chi^a = \langle \Psi, \Sigma^a \rangle$. By Item 2 of Lemma 4.3, for $S = \emptyset$, $R'$ in turn gives rise to a repair $R''$ of $A$ w.r.t. $\chi$ of the form $R'' = R' \cup S'$ such that $S' \cap C^* = \emptyset$. Since clearly $S' \models \Sigma^s$, property (DPE) implies that $S'$ is a repair of $S = \emptyset$ w.r.t. $\langle \Psi, \Sigma^s \rangle$. Since each constraint in $\Sigma^s$ has a nonempty body, it follows by (SIP) that $S' = \emptyset$. Hence $R \cap C^*$ is a repair of $A$ w.r.t. $\chi$. Now if $R \not\subseteq C^*$ held, then $\triangle(R'', A) \subset \triangle(R, A)$ would hold, which by (SIP) implies $R'' <_D R$. This is a contradiction. $\square$

Recall that $\mathbf{C}_2$ are the non-disjunctive constraints, i.e., every constraint has at most one database atom in the head.

**Proposition 4.7** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, where $\Sigma \subseteq \mathbf{C}_2$. Then (i) every repair $R$ of $A = D \cap C^*$ satisfies $R \subseteq D \cup C^*$, and (ii) for every repairs $R, R'$ of $A$, $R \cap (D \setminus C^*) = R' \cap (D \setminus C^*)$.*

PROOF. By the argument in the proof of Proposition 4.6, every $R \in rep(A)$ gives rise to some $R'' \in rep(A)$ of the form $R'' = (R \cap C^*) \cup S'$ such that $S' \cap C^* = \emptyset$ and $S'$ is a repair of $S = \emptyset$ w.r.t. $\langle \Psi, \Sigma^s \rangle$. As each constraint in $\Sigma^s$ is non-disjunctive, there is the least (w.r.t. $\subseteq$) set of facts $\mathcal{F}$ such that $\mathcal{F} \models \Sigma^s$ (in essence, $\Sigma^s$ is a Horn theory), and $\mathcal{F} \subseteq S'$ must hold; by (SIP), $\mathcal{F} = S'$. Now if $R \not\subseteq C^* \cup D$ held, then $\triangle(R'', A) \subset \triangle(R, A)$ would hold ($\mathcal{F} \subseteq R$ must hold, and thus $R'' \subseteq R$), which by (SIP) means $R'' <_D R$. This is a contradiction, and proves (i). Item (ii) holds as $R \cap (D \setminus C^*) = \mathcal{F}$ for each $R \in rep(A)$. $\square$

The proposition above allows us to exploit Theorem 4.4 in a constructive way for many significant classes of constraints, for which it implies a bijection between the repairs of a database $D$, and the repairs of the affected part $A = D \cap C^*$.

**Corollary 4.8** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$ where $\Sigma \subseteq \mathbf{C}_i$, for $i \in \{1, 2\}$. Then, $C^*$ is a repair envelope for $D$. In fact, there exists a bijection $\mu : rep(D) \to rep(D \cap C^*)$, such that for every $R \in rep(D)$, $R = \mu(R) \cup (D \setminus C^*)$.*

By this result, the repairs of a database $D$ can be computed by avoiding Step 3 of the procedure given in Section 4.1. Note also that by the above corollary and Proposition 4.1 and Corollary 4.5, we can make $C^*$ a repair envelope for an arbitrary relational schema $\chi = \langle \Psi, \Sigma \rangle$, if we can modify $\Sigma$ to constraints $\Sigma'$ from $\mathbf{C}_1$ or $\mathbf{C}_2$ while preserving the affected constraints, i.e., $\Sigma_\chi^a(D) = \Sigma_{\langle \Psi, \Sigma' \rangle}^a(D)$. Technically, this can be exploited in different ways, e.g., by dropping constraints, adding ground instances of constraints, rewriting constraints by modifying the built-in part (in fact, only semantic equivalence of affected ground constraints is needed), etc.

We also remark that $C^*$ may be decreased to a smaller repair envelope, by taking tuple generating constraints into account. For example, if $p(a)$ belongs to each repair (e.g., enforced by a constraint), $p(a)$ can be removed from the repair envelope. If there is another constraint $p(x) \supset q(x)$, also $q(a)$ can be removed. Exploring this is left for further study.

4.2.2 *Constraints* $\mathbf{C}_0$. Recall that constraints in $\mathbf{C}_0$ have only built-in relations in the head. Notably, the repairs of a database with integrity constraints from this class are computable by focusing on the immediate conflicts in the database, without the need of computing the conflict closure set, which may be onerous in general. Furthermore, repairs always do only remove tuples from relations, but never include new tuples. We will next formally prove these properties, starting with the following proposition.

**Proposition 4.9** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, $\Sigma \subseteq \mathbf{C}_0$, and let $A = D \cap C^*$. Then,*

*(1)* $C \subseteq D$;

*(2) for each $R \in rep(A)$, (i) $R \subseteq A$, (ii) $\triangle(R, A) \subseteq C$, (iii) $A \backslash C \subseteq R$, and (iv) $R \cap C \in rep(C)$;*

*(3) for each $R \in rep(C)$, $R \cup (A \setminus C) \in rep(A)$.*

Note that Proposition 4.9 shows that each repair of the conflict set $C$ just removes tuples from $C$ (take $D = C$ in Item 2.(ii)). Furthermore, because $C \subseteq D$, we can compute $C$ efficiently by suitable SQL statements which express constraint violations. We are now ready to prove that under $\mathbf{C}_0$ constraints, we can use $C$ instead of $C^*$ as a repair envelope, and thus avoid the onerous construction of $C^*$. In fact, we prove a more general result.

**Theorem 4.10** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$ where $\Sigma \subseteq \mathbf{C}_0$. Then, every set of facts $E \supseteq C$ is a repair envelope for $D$. Moreover, there exists a bijection $\nu : rep(D) \rightarrow rep(D \cap E)$, such that for each $R \in rep(D)$, $R = \nu(R) \cup (D \setminus E)$.*

PROOF. By Corollary 4.8, there is a bijection $\mu : rep(D) \rightarrow rep(A)$, where $A = D \cap C^*$, such that the repairs of $D$ are given by $\mu(R) \cup (D \setminus C^*)$, for all $R \in rep(A)$. Items 1 and 2.(iv) of Proposition 4.9 and the fact that each repair $R$ of $C$ satisfies $R \subseteq C$, imply that all repairs of $A$ are given by $(A \setminus C) \cup R$, where $R \in rep(C)$. Hence, the mapping $\nu : rep(D) \rightarrow rep(C)$ given by $\nu(R) = \mu(R) \cap C$ is a bijection such that

$$
\begin{aligned}
R &= \mu(R) \cup (D \setminus C^*) \\
&= \nu(R) \cup (A \setminus C) \cup (D \setminus C^*) \\
&= \nu(R) \cup ((D \cap C^*) \setminus C) \cup (D \setminus C^*) \;=\; \nu(R) \cup (D \setminus C)
\end{aligned}
$$

This proves the result for $E = C$. For general $E \supseteq C$, we note that $D' = D \cap E$ and $D$ have the same conflict set; hence, there exists a bijection $\nu' : rep(D \cap E) \rightarrow rep(C)$ such that $R = \nu'(R) \cup (D' \setminus C)$, for each $R \in rep(D')$. This implies a bijection $\nu'' : rep(D) \rightarrow rep(D \cap E)$ of the given form. $\square$

Consequently, in this setting we can compute the repairs of a database $D$ as follows:

1. compute $C$,
2. compute the repairs $R$ of $C$ (where $R \subseteq C \subseteq D$), and
2. take for each such repair $R$ the union with $D \setminus C$.

An example of application of the above procedure has been given in Example 4.1. The fact that every $E \supseteq C$ is a repair envelope gives convenient flexibility to modify the SQL statements for computing $C$ (i.e., one may sensibly simplify conditions in the SQL statements such that they are easier to evaluate but might infer more tuples).

## 5. QUERY ANSWERING THROUGH LOCALIZED REPAIRS

The localization properties discussed in the previous section may be used to optimize consistent query answering from an inconsistent database $D$. Indeed, based on them, one may conceive an optimization procedure consisting of the following three steps:

*Focusing Step.* Localize inconsistency in $D$, and single out facts that are affected by repair, and facts that are not, i.e., compute the (weak) repair envelope $E$ and the affected database $D \cap E$ and the safe database $D \setminus E$.

*Decomposition Step.* Compute repairs of the affected database, and obtain from them repairs of $D$ (by suitably incorporating the safe database).

*Recombination Step.* Recombine the repairs of $D$ for computing the consistent query answers, i.e., evaluate the query over local repairs augmented with the safe database and compose consistent query answers from these query results.

In situations in which the size of the affected database is much smaller than the size of the database $D$, computing the repairs of the affected database is significantly faster than the naive computation, which just aims at changing tuples "randomly" in the database, and does not in general rely on a focusing strategy. Moreover, localizing the inconsistency can be carried out easily by evaluating the constraints issued over the schema (by means of suitable SQL statements). Focusing and decomposition have been amply discussed in Section 4. In this section, we address the issue of efficient recombination, by illustrating three (increasingly elaborate) approaches for its implementation:

- A basic methodology for evaluating the query over local repairs augmented with the safe database is discussed in Section 5.1. This methodology linearly scales w.r.t. the number of repairs, but possibly exponentially w.r.t. the size of the affected database. Yet, it can be applied to arbitrary queries and constraints.

- A more elaborate approach is then discussed in Section 5.2, based on the concept of *repair factorization*. Roughly, it aims at decomposing a repair envelope into disjoint components such that the repairs of $D$ can be efficiently obtained from their repairs. Note that here two aspects are crucial:
  *(i)* The ability to identify such components. We propose sufficient conditions for factorization, based on which components can be identified in polynomial time for practical relevant classes of constraints (e.g., $\mathbf{C}_0$ and $\mathbf{C}_1$).
  *(ii)* The ability to combine the component repairs. We define in Section 5.2.1 two kinds of relevant components: *singular* components, which though inconsistent can be repaired in one particular rather than in all the possible ways to compute consistent answers to a certain query $Q$; and *decomposable* components, whose repairs can be processed in an independent (parallel) manner. When all the components of a factorization fall in one of these two categories (which can be efficiently checked by analyzing both $Q$ and $D$), consistent query answering is polynomial.

- Eventually, in Section 5.2.2 we show how to combine repair factorization with other techniques. In particular, we present a *query grounding* strategy which reduces a general query to a set of ground queries with the effect that more components can become singular while decomposability of components is unaffected. We thus enlarge the class of queries for which consistent answers are computable in polynomial time, and in fact single out a purely syntactic condition for tractability.

Towards the combined application of these techniques, we may want to implement recombination via the repair factorization approach for constraints in $\mathbf{C}_0$ or $\mathbf{C}_1$, for instance. If some non-singular component emerges, then we may exploit the query grounding strategy provided that all other components are decomposable. If some component remains neither singular nor decomposable, then we eventually have to resort to the basic method.

## 5.1 Recombination Step

Let us now consider the problem of evaluating a query $Q$ issued over an inconsistent database $D$ for $\chi$, i.e., to compute $ans_c(Q, D)$. Recall that according to the definition in Section 2, a tuple $\vec{t}$ belongs to $ans_c(Q, D)$ if $\vec{t}$ is in the evaluation of $Q$ over every repair of $D$, i.e., $ans_c(Q, D) = \{\vec{t} \mid \vec{t} \in Q[R] \text{ for each } R \in rep(D)\} = \bigcap_{R \in rep(D)} Q[R]$. The following proposition, which is immediate from the definitions, states how we can exploit repair envelopes for localization in query answering.

**Proposition 5.1** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, and let $Q$ be a non-recursive Datalog$^\neg$ query. Let $E$ be a set of facts, and let $A = D \cap E$ and $S = D \setminus E$. Then*

$$ans_c(Q, D) = \bigcap_{R \in rep(A)} Q[\varepsilon(R) \cup S], \tag{6}$$

*where (i) $\varepsilon(R) = R$ if $E$ is a repair envelope for $D$, and (ii) $\varepsilon(R) = R \cap E$ if $E$ is a weak repair envelope for $D$.*

By the results from above, we can always apply (ii) of (6) with $E = C^*$, and for $\mathbf{C}_1$ or $\mathbf{C}_2$ constraints apply always (i) of (6) with $E = C^*$. Furthermore, for $\mathbf{C}_0$, we can apply always (i) of (6) with $E = C$. Since in this case $C \subseteq D$, we can rewrite (6) to $ans_c(Q, D) = \bigcap_{R \in rep(C)} Q[R \cup S]$.

In the light of the equations above, query answering can be carried out by "locally" repairing the affected database, and evaluating the query over each local repair augmented with the safe portion of the data. While this approach has the advantage of localizing the inefficient repair computation (repair checking is already co-NP-hard, cf. [Chomicki 2007] and references therein) on a fragment $A$ of the database $D$, its implementation leads to an algorithm for consistent query answering which linearly scales w.r.t. the number of repairs, but possibly exponentially w.r.t. the size of the affected database. Indeed, such an algorithm computes consistent answers to the query by one evaluation of the query per repair, but in general the number of repairs may be exponential in the number of constraint violations, respectively in the size of the affected database. Actually, this is the best one may asymptotically expect to achieve for general inconsistent databases and universal constraints, unless P = NP, given that consistent query answering is $\Pi_2^P$-hard in such a setting for non-recursive Datalog$^\neg$ queries (cf. [Chomicki 2007] and references therein).

Hence, it is particularly relevant to assess whether some smarter strategies can be conceived for special settings, in order to have an algorithm that both implements localized repair computation and linearly scales w.r.t. the size of the database. More precisely, the aim is at localization strategies that, for the recombination step, exploit situations where the number of repairs which need to be considered for consistent query answering is linear in the size of the affected database. In the next subsection we formally elaborate such a strategy based on repair factorization, whereas a logic-programming based implementation is described in Section 6.

## 5.2 Repair Factorization

In this section, we present a technique that factorizes repairs into independent components. The basic idea is to partition the affected part $A = D \cap E$ of the database $D$ w.r.t. a repair envelope $E$ into disjoint subparts $A_1, \ldots, A_m$, such that the repairs of $A$ are obtained by combining the repairs of $A_1, \ldots, A_m$ in all possible ways. Given a repair envelope $E$ for $D$ and $\chi$, a partitioning $E_1, \ldots, E_m$ of $E$ is a *factorization* of $E$ for $D$ and $\chi$, if

$$rep(D) = \{(D \setminus E) \cup R_1 \cup \cdots \cup R_m \mid R_i \in rep(D \cap E_i), 1 \leq i \leq m\}. \tag{7}$$

Towards sufficient conditions for factorization, we define a repair-compliant partitioning as follows.

**Definition 5.1** Let $E$ be a repair envelope for a database $D$ for a schema $\chi = \langle \Psi, \Sigma \rangle$. A partitioning $E_1, \ldots, E_m$ of $E$ is *repair-compliant*, if (1) it is *constraint-bounded*, i.e., constraint-bounded facts from $E$ belong to the same component $E_i$, and (2) for all $R \in rep_\chi(D \cap E)$ and $R_i \in rep_\chi(D \cap E_i)$, $1 \leq i \leq m$, $R \setminus E = R_i \setminus E_i$. $\square$

**Example 5.1** Consider a schema with relations $r(A, B)$ and $s(B, C)$ which have the keys $A$ and $B$, respectively. Let $D = \{r(a_1, b_1), r(a_1, b_2), r(a_2, b_1), r(a_2, b_3), r(a_3, b_1), s(b_1, c_1), s(b_1, c_2), s(b_3, c_3)\}$. Its conflict set is $C = D \setminus \{r(a_3, b_1), s(b_3, c_3)\}$, which is a repair envelope. Note that the safe part of $D$ is $S = \{r(a_3, b_1), s(b_3, c_3)\}$. The partitioning $C_{r_1} = \{r(a_1, b_1), r(a_1, b_2)\}, C_{r_2} = \{r(a_2, b_1), r(a_2, b_3)\}$ and $C_s = \{s(b_1, c_1), s(b_1, c_2)\}$ of $C$ is repair-compliant: It is easily verified that it is constraint-bounded (constraint-bounded facts are exactly the pairs of facts in each partition). Moreover for each repair $R \in rep(D \cap E)$ we have $R \setminus E = \emptyset$, since it consists in dropping one fact of each conflicting pair. Therefore, the same is true for each partition, i.e., $R \setminus C_{r_1} = \emptyset$ for $R \in rep(D \cap C_{r_1})$, $R \setminus C_{r_2} = \emptyset$ for $R \in rep(D \cap C_{r_2})$, and $R \setminus C_s = \emptyset$ for $R \in rep(D \cap C_s)$. $\square$

By means of a repair-compliant partitioning, we can factorize the repair of $A = D \cap E$ into the repair of the (mutually disjoint) parts $A_i = A \cap E_i = D \cap E_i$ of $A$, for $i = 1, \ldots, m$. The repairs for each $A_i$ are confined to $F \cup E_i$ for a fixed set of facts $F$, and by the abstract properties (SIP), (DPE), and (DIS) of the preference ordering, they can be easily combined with the repairs for all other parts $A_j$, as shown next.

**Theorem 5.2 (Factorization)** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, and let $E$ be a repair envelope for $D$. Then, every repair-compliant partitioning $E_1, \ldots, E_m$ of $E$ is a factorization of $E$ for $D$ and $\chi$.*

PROOF. We need to show that $rep(D) = \{(D \setminus E) \cup R_1 \cup \cdots \cup R_m \mid R_i \in rep(D \cap E_i), 1 \leq i \leq m\}$. Since $E$ is a repair envelope for $D$ and $\chi$, we know that $rep(D) = \{(D \setminus E) \cup R \mid R \in rep(D \cap E)\}$. Hence, it is sufficient to prove that:

($\subseteq$) $R \in rep(D \cap E)$ implies $R = R_1 \cup \ldots \cup R_m$ and $R_i \in rep(D \cap E_i)$ for $1 \leq i \leq m$;
($\supseteq$) every $R \in \{R_1 \cup \ldots \cup R_m \mid R_i \in rep(D \cap E_i), 1 \leq i \leq m\}$ is a repair of $D \cap E$.

($\subseteq$) Let $R \in rep(D \cap E)$. Then, by repair-compliance of $E_1, \ldots E_m$, $R = F \cup R_E$, where $F = R \setminus E$ and $R_E \subseteq E$. Consider $R_i = F \cup (R_E \cap E_i)$ for $1 \leq i \leq m$. It remains to show that $R_i \in rep(D \cap E_i)$ for $1 \leq i \leq m$. Towards a contradiction first assume that $R_i \not\models \Sigma$ for some $1 \leq i \leq m$. Then, there exists some $\sigma \in ground(\Sigma)$ such

19

that $R_i \not\models \sigma$. Thus, $R_i \models body(\sigma)$, which implies $R \models body(\sigma)$, and $R_i \not\models head(\sigma)$. However, $R \models head(\sigma)$ must hold since $R \models \Sigma$ by hypothesis. This means that there exists a head atom $B(\vec{y})$ of $\sigma$ which is true in $R$. Since $R_i \not\models head(\sigma)$, none of the built-in predicates of $\sigma$ is true and therefore $B(\vec{y})$ is a fact such that $B(\vec{y}) \in E_j$, $j \neq i$. Since the partitioning is constraint-bounded, it follows that $body(\sigma) \subseteq F$ and $head(\sigma) \cap E_i = \emptyset$. Thus no repair of the form $F \cup R'_{E_i}$ of $D \cap E_i$ such that $R'_{E_i} \subseteq E_i$ can exist, a contradiction to the repair compliance of $E_1, \ldots, E_m$. This proves $R_i \models \Sigma$ for every $i = 1, \ldots, m$.

Consequently, $R_i \in rep(D \cap E_i)$ iff there exists no $R'_i \in rep(D \cap E_i)$ such that $R'_i <_{D \cap E_i} R_i$ and $R'_i \models \Sigma$. Assume such an $R'_i$ exists. Then $R'_i = F \cup R'_{E_i}$ and thus by (DIS) $R'_{E_i} <_{D \cap E_i} R_{E_i}$. By (DPE) we would conclude for $R'_E = (R \cap E_1) \cup \ldots \cup R'_{E_i} \cup \ldots \cup (R \cap E_m)$, that $R'_E <_{D \cap E} R_E$, which implies $R' <_{D \cap E} R$ for $R' = F \cup R'_E$. Furthermore, $R' \models \Sigma$. (Otherwise there exists some $\sigma \in ground(\Sigma)$ such that $R' \models body(\sigma)$ and $R' \not\models head(\sigma)$, while $R \models \sigma$. We can conclude that $body(\sigma) \subseteq R'_i$, and since $R'_i \models \Sigma$ we obtain $R' \models head(\sigma)$, a contradiction). Together with $R' <_{D \cap E} R$, however, this contradicts $R \in rep(D \cap E)$. Hence, $R_i \in rep(D \cap E_i)$, for $1 \leq i \leq m$.

($\supseteq$) Let $R \in \{R_1 \cup \ldots \cup R_m \mid R_i \in rep(D \cap E_i), 1 \leq i \leq m\}$. We show that $R \in rep(D \cap E)$. Towards a contradiction suppose $R \not\models \Sigma$, i.e., $R \not\models \sigma$ for some $\sigma \in ground(\Sigma)$. By definition of a repair-compliant partitioning, we conclude that $R = F \cup (R_i \cap E_i) \ldots \cup (R_m \cap E_m)$, where $F = R_i \setminus E_i$ for any $1 \leq i \leq m$. Consequently, $R \models body(\sigma)$ implies $R_i \models body(\sigma)$ for some $1 \leq i \leq m$ by constraint-boundedness. However, $R_i \not\models head(\sigma)$ (otherwise $R \models head(\sigma)$), which contradicts $R_i \in rep(D \cap E_i)$. Hence, $R \models \Sigma$.

It remains to show that there is no $R' \in rep(D \cap E)$ such that $R' <_{D \cap E} R$. Assume the contrary and let $F = R' \setminus E$. Then by (DIS) (disjunctive split), either ($i$) $R' \cap E_i <_{D \cap E_i} R \cap E_i$ or ($ii$) $R' \setminus E_i <_{(D \cap E) \setminus E_i} R \setminus E_i$ holds for each $i = 1, \ldots, m$. Case ($i$) leads to a contradiction with $R_i \in rep(D \cap E_i)$, however, since it implies $F \cup (R' \cap E_i) <_{D \cap E_i} R_i$ and $F \cup (R' \cap E_i) \models \Sigma$ (otherwise $R' \not\models \Sigma$). So ($ii$) must hold for every $i = 1, \ldots, m$. As shown by the recursive argument below, it follows that $R' \setminus E <_{(D \cap E) \setminus E} R \setminus E$, which however, by repair-compliance of $E_1, \ldots, E_m$, is equivalent to $F <_\emptyset F$, a contradiction. To see this, note that we can apply (DIS) to $R' \setminus E_i <_{(D \cap E) \setminus E_i} R \setminus E_i$ w.r.t. $E_j$ for any $1 \leq j \neq i \leq m$, and arrive in a similar situation as above: either ($i'$) $(R' \setminus E_i) \cap E_j <_{((D \cap E) \setminus E_i) \cap E_j} (R \setminus E_i) \cap E_j = R' \cap E_j <_{D \cap E_j} R \cap E_j$, or ($ii'$) $(R' \setminus E_i) \setminus E_j <_{((D \cap E) \setminus E_i) \setminus E_j} (R \setminus E_i) \setminus E_j$. Now ($i'$) leads to a contradiction as in ($i$), and therefore ($ii'$) must hold. Iterating this argument $m - 1$ times yields $R' \setminus (E_1 \cup \ldots \cup E_m) <_{(D \cap E) \setminus (E_1 \cup \ldots \cup E_m)} R \setminus (E_1 \cup \ldots \cup E_m)$, which is equivalent $R' \setminus E <_{(D \cap E) \setminus E} R \setminus E$. This proves $R \in rep(D \cap E)$. $\square$

Note that Condition (2) of Definition 5.1 is trivially satisfied for $\mathbf{C}_0$ constraints. Furthermore, it is immaterial for $\mathbf{C}_1$ constraints under the standard envelope $E = C^*$.

**Proposition 5.3** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, and let $E$ be a repair envelope for $D$. If either (1) $\Sigma \subseteq \mathbf{C}_1$ and $E = C^*$ or (2) $\Sigma \subseteq \mathbf{C}_0$, then every constraint-bounded partitioning $E_1, \ldots, E_m$ of $E$ is repair-compliant.*

Thus, for the practically important classes of constraints $\mathbf{C}_1$ and $\mathbf{C}_0$, repair-compliant partitionings, and thus factorizations, can be obtained by a constraint-bounded partitioning of $C^*$, respectively by a constraint-bounded partitioning of any repair envelope. Conse-

quently, for $\mathbf{C}_0$ and the canonical envelope $E = C$, Equation (7) can be rewritten to:

$$rep(D) = \{(D \setminus C) \cup R_1 \cup \cdots \cup R_m \mid R_i \in rep(C_i), 1 \le i \le m\}, \qquad (8)$$

where $C_1, \ldots, C_m$ is a constraint-bounded partition of $C$.

**Example 5.2** Let $\chi$ consist of the relation $p(x, y, z)$ and the functional dependency $f$ : $p(x, y, z) \wedge p(x, y', z') \supset z = z'$ (which is of class $\mathbf{C}_0$), and consider the database $D = \{p(a_i, b_j, c_k) \mid 1 \le i \le m \wedge 1 \le j, k \le \ell\}$. The conflict set $C$ consists of all tuples in $D$, since each pair of facts of the form $p(a_i, b_j, c_k)$ and $p(a_i, b_{j'}, c_{k'})$ with $k \ne k'$ witnesses a violation of $f$. The partitioning $C_1, \ldots, C_m$ of $C$, where $C_i = \{p(a_i, b_j, c_k) \in C\}$, $1 \le i \le m$, is constraint-bounded and thus, by Proposition 5.3, repair-compliant, and by Theorem 5.2, a factorization. Every $C_i$ has $\ell$ repairs, while $D$ has $\ell^m$ repairs in total. In particular, the repairs of $D$ are of the form $R_1 \cup \cdots \cup R_m$, where each $R_i$ is a repair for $C_i$, according to Equation (8). □

We finally remark that under particular preference relations, Condition (2) for repair-compliance (see Definition 5.1) might be relaxed. For instance, in case of the prototypical preorder $\le_D$, i.e., set inclusion w.r.t. symmetric difference, it is sufficient that the repairs of $D \cap E_i$ coincide outside $E_i$ on a fixed part: for all $1 \le i, j \le m$, $R_i \in rep(D \cap E_i)$ and $R_j \in rep(D \cap E_j)$ implies $R_i \setminus E_i = R_j \setminus E_j$.

Furthermore, we note that we can compute efficiently repair-compliant partitionings of arbitrary repair envelopes for $\mathbf{C}_0$ constraints and of the standard envelope $E = C^*$ for $\mathbf{C}_1$ constraints, for instance, using techniques for computing the connected components of a graph. Note that each $E_i$ is a union of connected components of the graph with nodes in $E$ and edges between each pair of constraint-bounded facts. In this respect, point out that techniques exploiting different graph (and hypergraph) representations of conflicts in data have been introduced and used in [Arenas et al. 2001; Chomicki and Marcinkowski 2005].

5.2.1 *Recombination of Independent Factors.* We are now in the position to show how the notion of factorization can be used to optimize query answering from inconsistent databases. To this end, we proceed in two directions:

– First, for a user query $Q$, we investigate when some of the components of a factorization $E_1, \ldots, E_m$ do not entail repairs in all possible ways to answer $Q$. Intuitively, this happens when a single repair of the affected part $D \cap E_i$ is sufficient for answering $Q$.

– Second, we investigate how to improve the naive usage of Equation (7), by discussing scenarios where consistent answers to $Q$ can be obtained by independently processing the different components, rather than combining their repairs in all possible ways.

We focus here on non-recursive Datalog queries $Q = \langle q, \mathcal{P} \rangle$. Since they can be effectively unfolded to a union of conjunctive queries with a single head predicate $q$, we assume that queries are already in this form, i.e., $\mathcal{P} = \{\rho_1, \ldots, \rho_n\}$, where $head(\rho_j) = q(\vec{t}_j)$. We denote by $n(Q)$ the number of conjunctive queries, i.e., rules in $\mathcal{P}$ and by $v(Q)$ the maximal number of variables appearing in any $\rho_j$.

In order to understand whether for the component $E_i$ not all repairs are needed, we consider a *rewriting* of $Q$ w.r.t. $E_i$. It aims at determining whether a particular repair for $E_i$ is sufficient to answer $Q$. Roughly, for each $\rho_j$ in $\mathcal{P}$ a test rule $\rho_{i,j}$ is created, which for each repair $R$ for $E_i$ yields a cautious overestimate of the result of $\rho_j$ over the safe

database plus $R$; the overestimates for all $\rho_j$ are then collected into an overestimate of the result of $Q$. If over all repairs $R$ for $E_i$ a single minimal overestimate exists, then we can simply use the corresponding $R$ in computing the consistent answer of $Q$ w.r.t. $D$.

To define $\rho_{i,j}$, we introduce three subsets of $body(\rho_j)$: a set $\beta_{i,j}^{in}$ of atoms which may change value in different repairs, a set $\beta_{i,j}^{out}$ of non built-in atoms unaffected by all repairs, and a set $\beta_{i,j}^{bin}$ of built-in atoms that are connected to atoms in these two sets. In detail,

$$
\begin{aligned}
\beta_{i,j}^{in} &= \{p(\vec{x}) \in body(\rho_j) \mid \exists \theta : p(\vec{x}\theta) \in E_i\}, \\
\beta_{i,j}^{out} &= \{p(\vec{x}) \in body(\rho_j) \mid \forall \theta : p(\vec{x}\theta) \notin E\}, \\
\beta_{i,j}^{bin} &= \{\phi(\vec{x}) \in body(\rho_j) \mid \forall x \in \vec{x} \exists p(\vec{x}') \in \beta_{i,j}^{in} \cup \beta_{i,j}^{out} \ \wedge x \in \vec{x}'\},
\end{aligned}
$$

where $\theta$ ranges over ground substitutions, $p \in \Psi$ (i.e., non built-in), and $\phi$ denotes a built-in predicate; note that $\beta_{i,j}^{out} = \beta_{i',j}^{out}$ for all $E_i$ and $E_{i'}$. In words, for each $\rho_j$ in $\mathcal{P}$, $\beta_{i,j}^{in}$ are the body atoms in $\rho_j$ that match with some fact in $E_i$; $\beta_{i,j}^{out}$ are the body atoms of $\rho_j$ that cannot be matched with any fact in the envelope $E$ via ground variable substitution; and $\beta_{i,j}^{bin}$ are the built-in (body) atoms in $\rho_j$ which join some atom in $\beta_{i,j}^{out}$ or $\beta_{i,j}^{in}$.[2]

**Example 5.3** Consider a schema $\chi$ with a ternary relation $r_1$, two binary relations $r_2$ and $r_3$, and a unary relation $r_4$, where the first argument for each relation is the key. Assume also that relations $r_1$ and $r_2$ are disjoint, i.e., it holds that $r_2(x,y) \wedge r_3(x,y') \supset y \neq y'$. Let $D = \{r_1(a,b,c),\ r_1(a,b,d),\ r_2(c,b),\ r_3(c,b),\ r_4(b)\}$. Its conflict set is $C = D \setminus \{r_4(b)\}$, which is a repair envelope. The partitioning $E_1 = \{r_1(a,b,c), r_1(a,b,d)\}$, $E_2 = \{r_2(c,b),\ r_3(c,b)\}$ of $C$ is repair-compliant (it is indeed constraint-bounded), and therefore is a factorization of $E$ for $D$ and $\chi$. Consider the query $Q = \langle q, \mathcal{P}\rangle$, where $\mathcal{P} = \{\ q(x) \leftarrow r_1(x,y,z), r_2(w,y),\ q(x) \leftarrow r_1(x,y,z), r_3(w,y), r_4(y)\ \}$, i.e., such that $n(Q) = 2$ and $v(Q) = 4$. Then, we easily obtain $\beta_{1,1}^{in} = \beta_{1,2}^{in} = \{r_1(x,y,z)\}$, $\beta_{2,1}^{in} = \{r_2(w,y)\}$, $\beta_{2,2}^{in} = \{r_3(w,y)\}$, $\beta_{1,1}^{out} = \beta_{2,1}^{out} = \emptyset$, $\beta_{1,2}^{out} = \beta_{2,2}^{out} = \{r_4(y)\}$, $\beta_{1,1}^{bin} = \beta_{1,2}^{bin} = \beta_{2,1}^{bin} = \beta_{2,2}^{bin} = \emptyset$. $\square$

Armed with the above notions, we now define the rewriting $Q_i$ of $Q$ w.r.t. $E_i$. The body of the test rule $\rho_{i,j}$ consists simply of all atoms in $\beta_{i,j}^{in}$, $\beta_{i,j}^{out}$, and $\beta_{i,j}^{bin}$. Its head contains all head variables of $\rho_j$ that occur in $\beta_{i,j}^{in}$, plus further join variables from atoms in $\beta_{i,j}^{in}$ to atoms in the body of $\rho_j$ outside $\beta_{i,j}^{in} \cup \beta_{i,j}^{out}$; intuitively, the join variables provide additional context information. The head of $\rho_{i,j}$ also contains an identifier $d_j$ that marks the contribution of $\rho_{i,j}$ to the result of $Q_i$; finally, since technically each head $\rho_{i,j}$ must have the same predicate, but the resulting lists of head variables for $\rho_{i,j}$ may have different lengths, we pad all lists to the same length using $d_j$. In detail,

**Definition 5.2** Let $E_1, \ldots, E_m$ be a factorization of a repair envelope $E$ for a database $D$, let $Q$ be a non-recursive (unfolded) Datalog query as above, and let $d_1, \ldots, d_{n(Q)} \notin \mathcal{U}$ be fresh constants. We define $Q_i$, the rewriting of $Q$ w.r.t. $E_i$, as follows:

$$
Q_i = \langle q_i, \mathcal{P}_i \rangle, \quad \mathcal{P}_i = \{\ \rho_{i,j} \mid 1 \le j \le n(Q)\ \},
$$

where

---

[2]Note that we consider safe queries. We could keep all built-in predicates if we allowed for unsafe rules or if we added respective domain predicates, ranging over the active domain, to make the rule safe.

(1) $head(\rho_{i,j}) = q_i(d_j, \vec{u}_{i,j}, \vec{v}_{i,j}, \vec{d_j})$ and $body(\rho_{i,j}) = \beta^{in}_{i,j} \cup \beta^{out}_{i,j} \cup \beta^{bin}_{i,j}$ $(=: \beta_{i,j})$;

(2) the arity of $q_i$ is $v(Q) + 1$;

(3) $\vec{u}_{i,j}$ are the variables from $\vec{t}_j$ (in any order), i.e., the variables in $head(\rho_j) = q(\vec{t}_j)$, which occur in some $p(\vec{x}) \in \beta^{in}_{i,j}$;

(4) $\vec{v}_{i,j}$ are the variables (in any order) not occurring in $\vec{u}_{i,j}$ but in some $p(\vec{x}) \in \beta^{in}_{i,j}$ and some $p'(\vec{y}) \in body(\rho_j) \setminus \beta^{in}_{i,j}$, unless $p'(\vec{y}) \in \beta^{out}_{i,j}$ or $p'$ is built-in, and all variables from $\vec{y}$ occur in $\beta^{in}_{i,j}$;

(5) $\vec{d_j} = d_j, \ldots, d_j$ is a padding to the arity of $q_i$;   □

Note that in Condition (2), we use $v(Q) + 1$ for simplicity (smaller values are possible). In Condition (4), the join variables from $\beta^{in}_{i,j}$ to atoms outside $\beta^{in}_{i,j} \cup \beta^{out}_{i,j}$ are collected in $\vec{v}_j$; as an optimization, built-in atoms having all their variables in $\beta^{in}_{i,j}$ can be omitted.

**Example 5.4** In the setting of Example 5.3, we have $Q_1 = \langle q_1, \{\rho_{1,1}, \rho_{1,2}\}\rangle$, with $\rho_{1,1} = q_1(d_1, x, y, d_1, d_1) \leftarrow r_1(x, y, z)$, and $\rho_{1,2} = q_1(d_2, x, y, d_2, d_2) \leftarrow r_1(x, y, z), r_4(y)$, and $Q_2 = \langle q_2, \{\rho_{2,1}, \rho_{2,2}\}\rangle$, with $\rho_{2,1} = q_2(d_1, y, d_1, d_1, d_1) \leftarrow r_2(w, y)$, and $\rho_{2,2} = q_2(d_2, y, d_2, d_2, d_2) \leftarrow r_3(w, y), r_4(y)$.   □

**Example 5.5** Continuing Example 5.1, the partitioning $C_{r_1}, C_{r_2}, C_s$ of $C$ is repair-compliant, and thus by Theorem 5.2 a factorization. The query $Q = \langle q, \{q(x) \leftarrow r(x, y), s(y, z).\}\rangle$ has $n(Q) = 1$ and $v(Q) = 3$. For $C_{r_1}$ and $C_{r_2}$, we obtain $\beta^{in}_{r_1,1} = \beta^{in}_{r_2,1} = \{r(x, y)\}$, and for $C_s$, $\beta^{in}_{s,1} = \{s(y, z)\}$. Furthermore, $\beta^{out}_{r_1,1} = \beta^{out}_{r_2,1} = \beta^{out}_{s,1} = \emptyset$ and $\beta^{bin}_{r_1,1} = \beta^{bin}_{r_2,1} = \beta^{bin}_{s,1} = \emptyset$. We thus have $Q_{r_1} = \langle q_{r_1}, \{q_{r_1}(d_1, x, y, d_1) \leftarrow r(x, y)\}\rangle$ and $Q_{r_2} = \langle q_{r_2}, \{q_{r_2}(d_1, x, y, d_1) \leftarrow r(x, y)\}\rangle$, while $Q_s = \langle q_s, \{q_s(d_1, y, d_1, d_1) \leftarrow s(y, z)\}\rangle$.   □

Recall that the rewriting $Q_i$ should allow to identify components for which a single repair is sufficient to evaluate the original query $Q$. As we show now, this is the case if considering all possible repairs $R$ for a component, $Q_i$ has a single minimal (w.r.t. set inclusion) result over the safe database plus $R$. We call such components *singular*.

**Definition 5.3** Let $E_1 \ldots, E_m \subseteq E$ be a factorization of a repair envelope $E$ for a database $D$. Let $Q$ be a non-recursive (unfolded) Datalog query, and let $Q_i = \langle q_i, \mathcal{P}_i\rangle$ be a rewriting of $Q$ w.r.t. $E_i$, $1 \leq i \leq m$. Denote by $amin(Q_i)$ the set of the evaluations $Q_i[(D \setminus E) \cup R]$, where $R \in rep(D \cap E_i)$, which are minimal w.r.t. set inclusion. We then call $E_i$ *singular*, if $R \setminus E_i = R' \setminus E_i$ for all $R, R' \in rep(D \cap E_i)$ and $|amin(Q_i)| = 1$. □

**Example 5.6** For the setting of Example 5.4, we have $D \setminus E = \{r_4(b)\}$, and $rep(D \cap E_1) = \{R_1, R_2\}$, with $R_1 = \{r_1(a, b, c)\}$, and $R_2 = \{r_1(a, b, d)\}$, $rep(D \cap E_2) = \{R_3, R_4\}$, with $R_3 = \{r_2(c, b)\}$ and $R_4 = \{r_3(c, b)\}$. It is easy to see that $Q_1[D \setminus E \cup R_1] = Q_1[D \setminus E \cup R_2] = \{(d_1, a, b, d_1, d_1), (d_2, a, b, d_2, d_2)\}$, and therefore $|amin(Q_1)| = 1$, i.e., $E_1$ is singular. Furthermore, $Q_2[D \setminus E \cup R_3] = \{(d_1, b, d_1, d_1, d_1)\}$, $Q_2[D \setminus E \cup R_4] = \{(d_2, b, d_2, d_2, d_2)\}$. Hence $|amin(Q_1)| \neq 1$, i.e., $E_2$ is not singular.   □

**Example 5.7** In our Example 5.5, $D \cap C_{r_1}$ has the two repairs $\{r(a_1, b_1)\}$ and $\{r(a_1, b_2)\}$, and $amin(Q_{r_1}) = \{\{(d_1, a_1, b_1, d_1), (d_1, a_3, b_1, d_1)\}, \{(d_1, a_1, b_2, d_1), (d_1, a_3, b_1, d_1)\}\}$. Similarly, $D \cap C_{r_2}$ has two repairs, $\{r(a_2, b_1)\}$ and $\{r(a_2, b_3)\}$, and $amin(Q_{r_2}) =$

$\{\{(d_1, a_2, b_1, d_1), (d_1, a_3, b_1, d_1)\}, \{(d_1, a_2, b_3, d_1), (d_1, a_3, b_1, d_1)\}\}$. Finally, $D \cap C_s$ has the two repairs $\{s(b_1, c_1)\}$ and $\{s(b_1, c_2)\}$, but $amin(Q_s)$ is the singleton $\{\{(d_1, b_1, d_1, d_1), (d_1, b_3, d_1, d_1)\}\}$, i.e., $C_s$ is a singular component. $\square$

For a singular component $E_i$ it is sufficient to consider a repair $R$ that yields the minimal evaluation w.r.t. $Q_i$, i.e., such that $Q_i[(D \setminus E) \cup R] \in amin(Q_i)$. Hence, we can pick one such repair for each singular component and add it to the safe database, obtaining a set of facts that has not to be altered and can safely be combined with repairs of other components without altering the consistent answers of $Q$. We call such a set a *query-safe part*.

**Definition 5.4** Let $E_1, \ldots, E_m$ be a factorization of a repair envelope $E$ for a database $D$, where $E_1, \ldots, E_\ell$ are singular components, and let $Q$ be a non-recursive Datalog query. We call $S_Q = (D \setminus E) \cup R_1 \cup \cdots \cup R_\ell$ a *query-safe part* of $D$ w.r.t. $Q$, if each $R_i$, $1 \le i \le \ell$, is an arbitrary repair of $D \cap E_i$, such that $Q_i[(D \setminus E) \cup R_i] \in amin(Q_i)$. $\square$

In Example 5.6, $E_1$ is the only singular component, and $R_1 = \{r_1(a, b, c)\}$ is one such repair of $D \cap E_1$. Hence, $S_Q = \{r_1(a, b, c), r_4(b)\}$ is a query-safe part of $D$. Similarly in Example 5.7, $C_s$ is a singular component, and $S_Q = \{r(a_3, b_1), s(b_3, c_3), s(b_1, c_2)\}$ is a query-safe part. We then have the following result.

**Proposition 5.4** *Let $E_1, \ldots, E_m$ be a factorization of a repair envelope $E$ for a database $D$, let $Q$ be a non-recursive Datalog query, and let $S_Q = (D \setminus E) \cup R_1 \cup \cdots \cup R_\ell$ be a query-safe part of $D$ w.r.t. $Q$. Then,*

$$ans_c(Q, D) = \bigcap_{R_{\ell+1} \in rep(D \cap E_{\ell+1})} \cdots \bigcap_{R_m \in rep(D \cap E_m)} Q[S_Q \cup R_{\ell+1} \cdots \cup R_m]. \quad (9)$$

If all components $E_i$ are singular ($\ell = m$), query answering can be carried out by considering an arbitrary query-safe part of $D$. In this ideal case, the cost for query answering amounts to checking for all $1 \le i \le m$ that $|amin(Q_i)| = 1$ and that $R \setminus E_i = R' \setminus E_i$ for all $R, R' \in rep(D \cap E_i)$ (recall that this is trivial for $\mathbf{C}_0$ constraints and the standard envelope $C^*$ in case of $\mathbf{C}_1$ constraints) as well as eventually computing $Q[S_Q]$. Note that checking for singular components $E_i$ and determining a repair $R_i$, such that $Q_i[(D \setminus E) \cup R_i] \in amin(Q_i)$, can be carried out by processing the components independently of each other and is polynomial if the local repairs can be computed in polynomial time. For each component $E_i$, the effort generally depends on the number of its repairs.

Moreover, *irrelevant* components can be easily detected by syntactic checks: If $\beta_{i,j}^{in} = \emptyset$ for $1 \le j \le v(Q)$, then trivially $|amin(Q_i)| = 1$. In this case, if $R \subseteq (D \setminus E) \cup E_i$ for all $R \in rep(D \cap E_i)$, we can even tolerate inconsistency, i.e., we do not need to repair the component for consistent query answering and can skip it in the query-safe part $S_Q$.

**Example 5.8** Concluding Example 5.6, we apply Equation (9) as follows:

$$ans_c(Q, D) = \bigcap_{R \in rep(D \cap E_2)} Q[S_Q \cup R] = \bigcap_{R \in \{\, \{r_2(c,b)\}, \{r_3(c,b)\} \,\}} Q[\{r_1(a, b, c), r_4(b)\} \cup R],$$

and we get $ans_c(Q, D) = \{(b)\}$; this is the correct result. We can proceed analogously also for Example 5.7, where there are two non-singular components. However, we will show below that in this case a further optimization is possible. $\square$

Even if non-singular components are present, it may be possible to "parallelize" consistent query answering without a need for recombination. A simple case is if there is just a single atom in the body of the query such that for a set of components this is the only atom that unifies with facts from these components. Then, the query can be independently evaluated over these components since there is no "interference" between the components w.r.t. the query evaluation, like a join in the query that unifies with facts belonging to different components in this set. We call such a set of components *decomposable*.

**Definition 5.5** Let $E_1, \ldots, E_m$ be a factorization of a repair envelope $E$ for $D$. A set of components $E_1, \ldots, E_\ell$ is *decomposable* w.r.t. query $Q$, if they satisfy:

(1) $\beta_{i,k}^{in} = \beta_{j,k}^{in}$, for every $1 \le i, j \le \ell$ and $1 \le k \le n(Q)$;

(2) $|\beta_{i,k}^{in}| = 1$, for every $1 \le i \le \ell$ and $1 \le k \le n(Q)$;

(3) $R_i \setminus E_i = R_j \setminus E_j$, for every $R_i \in rep(D \cap E_i)$, $R_j \in rep(D \cap E_j)$, $1 \le i, j \le \ell$. $\quad\Box$

Query answering over decomposable components can be parallelized, and can be combined with singular components as follows.

**Theorem 5.5** *Let $E_1, \ldots, E_m$ be a factorization of a repair envelope $E$ for a database $D$, and let $Q$ be a non-recursive Datalog query. Suppose that $E_1, \ldots, E_\ell$ are singular components, with query safe part $S_Q$, and that $E_{\ell+1}, \ldots E_k$ is a set of decomposable components w.r.t. query $Q$. Then*

$$ans_c(Q, D) = \bigcap_{R_{k+1} \in rep(D \cap E_{k+1})} \cdots \bigcap_{R_m \in rep(D \cap E_m)}$$
$$\bigcup_{i=\ell+1}^{k} \Big( \bigcap_{R_i \in rep(D \cap E_i)} Q[S_Q \cup R_i \cup R_{k+1} \cup \cdots \cup R_m] \Big). \quad (10)$$

PROOF. We first show that

$$Q[X \cup R_{\ell+1} \cup \cdots \cup R_k] = \bigcup_{i=\ell+1}^{k} Q[X \cup R_i], \quad\quad\quad (11)$$

for every $R_j \in rep(D \cap E_j)$, $\ell < j \le k$, and for every set of facts $X$. We have $Q = \langle q, \mathcal{P} \rangle$, where $\mathcal{P} = \{\rho_1, \ldots, \rho_n\}$ is a set of "*not*"-free rules $\rho_j$. Consider any ground instance $\rho_j'$ of $\rho_j$, and an atom $p(\bar{t})$ occurring in the body of $\rho_j'$ that is satisfied by $X \cup R_{\ell+1} \cup \cdots \cup R_k$. Then either $p(\bar{t}) \in X$ or $p(\bar{t}) \in R_h$, for some $\ell < h \le k$. Furthermore, since $E_{\ell+1}, \ldots, E_k$ are decomposable w.r.t. $Q$, in the case where $p(\bar{t}) \in R_h \setminus (X \cup \bigcap_{i=\ell+1}^{k} R_i)$, there is no atom $p'(\bar{t}')$ in the body of $\rho_j'$ which belongs to $R_{h'} \setminus R_h$, for some $\ell < h' \ne h \le k$. Indeed, by Condition (3), $R_h' \setminus R_h = R_{h'} \cap E_{h'}$, and thus $p'(\bar{t}') \in E_{h'}$ would hold, while $p(\bar{t}) \in E_h$ holds. Conditions (1) and (2) would imply that $p(\bar{t})$ and $p'(\bar{t}')$ are instances of the same atom in the body of $\rho_j$, and thus $p(\bar{t}) = p'(\bar{t}')$. This contradicts $E_h \cap E_{h'} = \emptyset$.

As a consequence, the body of $\rho_j'$ is satisfied by $X \cup R_{\ell+1} \cdots \cup R_k$ iff it is satisfied by $X \cup R_h$, for some $\ell + 1 \le h \le k$. Since $\rho_j$ is non-disjunctive and "*not*"-free, it follows that $Q^{(j)}[X \cup R_{\ell+1} \cdots \cup R_k] = \bigcup_{i=\ell+1}^{k} Q^{(j)}[X \cup R_i]$, where $Q^{(j)} = \langle q, \{\rho_j\} \rangle$, and that

$$Q[X \cup R_{\ell+1} \cdots \cup R_k] = \bigcup_{j=1}^{n} Q^{(j)}[X \cup R_{\ell+1} \cdots \cup R_k] = \bigcup_{j=1}^{n} \bigcup_{i=\ell+1}^{k} Q^{(j)}[X \cup R_i]$$

25

$$= \bigcup_{i=\ell+1}^{k} \bigcup_{j=1}^{n} Q^{(j)}[X \cup R_i] = Q[X \cup R_{\ell+1}] \cup \cdots \cup Q[X \cup R_k].$$

This proves (11). To conclude the proof, by Proposition 5.4 and setting $X = S_Q \cup R_{k+1} \cup \cdots \cup R_m$ the consistent answers to $Q$ w.r.t. $D$ are:

$$ans_c(Q, D) = \bigcap_{R_{\ell+1} \in rep(D \cap E_{\ell+1})} \cdots \bigcap_{R_m \in rep(D \cap E_m)} Q[S_Q \cup R_{\ell+1} \cdots \cup R_m]$$

By Equation (11), we then get:

$$ans_c(Q, D) = \bigcap_{R_{\ell+1} \in rep(D \cap E_{\ell+1})} \cdots \bigcap_{R_m \in rep(D \cap E_m)} (Q[X \cup R_{\ell+1}] \cup \cdots \cup Q[X \cup R_k]),$$

from which the result follows by Boolean algebra (recall that for any sets $A, B_1, \ldots, B_k$, it holds that $\bigcap_{B \in \{B_1, \ldots, B_k\}} (A \cup B) = A \cup \bigcap_{B \in \{B_1, \ldots, B_k\}} B$). $\square$

By this result, we can take any (but not necessarily all) singular components, any decomposable components, and parallelize query answering. As mentioned above, singular components can always be identified efficiently given their local repairs. Furthermore, also a maximal set of decomposable components can be identified efficiently given the local repairs. In fact, all maximal such sets—which are pairwise disjoint—can be determined efficiently (note that the relation $E_{i_1} \equiv_Q E_{i_2}$ iff the set $E_{i_1}, E_{i_2}$ is decomposable w.r.t. $Q$, is an equivalence relation). In particular, if the factorization stems from a repair-compliant partitioning, then Condition (3) is always fulfilled, and the effort depends only on the check for (1) and (2). Roughly, in total the computational effort by exploiting singular and decomposable components reduces from the global combination of all local repairs of the components to taking the union of the results of local consistent query answering.

**Example 5.9** In Example 5.7, $C_{r_1}$ and $C_{r_2}$ are non-singular components. Since $n(Q) = 1$ and for both, $C_{r_1}$ and $C_{r_2}$, we have $\beta_{r_i,1}^{in} = \{r(x,y)\}$, $1 \leq i \leq 2$, the set $\{C_{r_1}, C_{r_2}\}$ is decomposable w.r.t. $Q$. By Theorem 5.5, the query $Q = \langle q, \{q(x) \leftarrow r(x,y), s(y,z).\} \rangle$ can be evaluated independently over $C_{r_1}$ and $C_{r_2}$, taking, e.g., the query-safe part $S_Q = \{r(a_3, b_1), s(b_3, c_3), s(b_1, c_2)\}$, into account. Specifically, for $C_{r_1}$, we must compute $Q[S_Q \cup \{r(a_1, b_1)\}] \cap Q[S_Q \cup \{r(a_1, b_2)\}]$, which yields $\{(a_3)\}$. For $C_{r_2}$, we must compute $Q[S_Q \cup \{r(a_2, b_1)\}] \cap Q[S_Q \cup \{r(a_2, b_3)\}]$, which yields $\{(a_3), (a_2)\}$. Therefore, $ans_c(Q, D) = \{(a_3), (a_2)\}$. As can be checked, this is the correct result. $\square$

Importantly, by virtue of Proposition 5.4 and Theorem 5.5, one can consistently answer queries in polynomial time which do not belong to any class shown to be tractable in the literature [Chomicki et al. 2004b; 2004a; Fuxman et al. 2005; Grieco et al. 2005]. We exemplify this for the scenario discussed in the following Example 5.10. The query there contains a non-key to non-key join, i.e., a join between non-key positions, and hence is not in any class shown to be tractable (see also Section 7). Nonetheless, we will show how to compute consistent answers for this query in polynomial time. We illustrate this in two steps: First, we exemplify this by exploiting the techniques from above for *restricted inputs*; it would not be clear how to do this from previous results either (which hinge on conditions on the query and the constraints, but are insensitive to the actual data). Second, we combine the factorization techniques with a simple grounding strategy, which then

yields a polynomial evaluation method for the query over *arbitrary inputs*. We finally extend this approach to a class of queries (of slightly different form) for which consistent query answering is tractable.

**Example 5.10** Inspired by a typical information system supporting university administration, we consider a schema with the relations $student(IDS, First, Last, Address)$, $prof(IDP, First, Last)$, and $dean(IDD, First, Last)$, where $IDS$, $IDP$, and $IDD$ are the respective keys. As a further constraint, we have the inclusion dependency $dean(x, y, z) \supset prof(x, y, z)$.

Suppose we want to know the identifiers of professors having their last name in common with a student. The query $Q = \langle q, \{q(x_2) \leftarrow student(x_1, y_1, z, w_1), prof(x_2, y_2, z).\} \rangle$ extracts them; note that it involves a non-key to non-key join. Now let $D$ be a database with the conflict set $C = \{student(0815, johann, meier, addr1), student(0815, hans, meier, addr1), prof(4711, markus, schmidt), prof(4711, mark, schmidt), dean(1111, egbertus, neumann), prof(1111, egbertus, neumann))\}$. $C$ is a repair envelope, which can be readily factorized into three components, two of them containing a pair of tuples from $E$ with equal key value over $prof$ and $student$, respectively, and one containing the remaining facts of $C$. As easily verified, all these components are singular. (This would not hold for the canonical repair envelope $C^*$.)

If $D'$ is a database yielding the conflict set $C' = \{student(0815, johann, meier, addr1), student(0815, johann, maier, addr1), student(4711, bodo, schmied, addr2), student(4711, bodo, schmid, addr2)\}$, and professors called $meier$, $maier$, $schmied$, and $schmid$, respectively, exist. Then, obviously $E' = C'$ is a repair envelope, and we can, e.g., again factorize into components containing tuples with equal key values (two this time). In this case the respective components are not singular, but decomposable.

Hence, we observe that as long as violations are restricted to only one of the relations $student$ or $prof$, we end up with singular and decomposable components such that Theorem 5.5 remains effective with a single set of decomposable components. The same is true for database instances where violations affect both relations, but where all violations w.r.t. one of the relations end up in singular components. □

We also remark that, as long as violations are restricted to only one of the relations or all components are singular, computing consistent answers to the query in Example 5.10 is feasible in polynomial time, even if we had a further relation $staff$ of the same structure as $prof$ in the schema and the query extended to a union of conjunctive queries by adding the rule $q(x_2) \leftarrow student(x_1, y_1, z, w_1), staff(x_2, y_2, z)$. As well, we could add the exclusion dependency $dean(x_1, y_1, z_1) \land staff(x_2, y_2, z_2) \supset x_1 \neq x_2 \lor y_1 \neq y_2 \lor z_1 \neq z_2$ to the scenario and still would obtain tractability for consistently answering the query. Observe that in all these scenarios, in order to obtain a repair envelope, we do not need to take the entire conflict closure into account. Rather we can restrict the envelope to constraint bounded facts over tuples of constants occurring in the conflict set (or simply to the conflict set, as done in Example 5.10). Furthermore, the resulting envelope can trivially be factorized into components with identical key values.

For arbitrary violations however, conflicts may interfere w.r.t. the query. In our student and professor example, too many non-singular and non-decomposable components might emerge, such that Theorem 5.5 can not be applied to establish tractability of the query, even if local repairs are computable in polynomial time. Nevertheless, tractability of this

query can be established by means of a refined factorization strategy, as shown next.

5.2.2 *Factorization and Query Grounding.* Our basic factorization approach might be combined with other techniques in order to get further tractability results. One such technique is to reduce general non-recursive queries $Q = \langle q, \mathcal{P} \rangle$, where $q$ has arity $k$, to ground (Boolean) queries by means of unification: for a tuple $\vec{c} = (c_1, \ldots, c_k)$ of constants, unify each rule $\rho \in \mathcal{P}$ with $q$ in the head with $q(\vec{c})$, i.e., apply to $\rho$ a substitution $\theta$ of constants to $\rho$'s head variables such that $head(\rho\theta) = q(\vec{c})$; if no such $\theta$ exists (e.g., for $q(b, X) \leftarrow p(X, c)$ and $q(a, b)$), simply delete $\rho$. Denote the resulting query by $Q_{\vec{c}}$;[3] then

$$ans_c(Q, D) = \{\vec{c} = (c_1, \ldots, c_k) \mid c_1, \ldots, c_k \text{ occur in } D \text{ or } Q, \ ans_c(Q_{\vec{c}}, D) \neq \emptyset\}. \quad (12)$$

That is, we can parallelize query answering for each tuple.

We can view $Q_{\vec{c}}$ as a refinement of $Q$, which is equivalently obtained by adding equality literals $x = c_i$ in the rule bodies in $Q$. As for consistent query answering, this does not affect decomposability of components $E_i$ (since this property does not depend on built-ins), and preserves their singularity, i.e., whenever $E_i$ is singular w.r.t. $Q$, then $E_i$ is also singular w.r.t. $Q_{\vec{c}}$. Thus, the number of singular components can only increase, and in benign cases only few components that are non-singular nor decomposable remain. In these cases, we may exploit Equation (12) and Theorem 5.5 to compute consistent query answers in polynomial time.

In the light of this important observation, we review the scenario in Example 5.10.

**Example 5.11** Reconsider the schema and $Q = \langle q, \{q(x_2) \leftarrow student(x_1, y_1, z, w_1), prof(x_2, y_2, z).\}\rangle$ in Example 5.10. Assume that we have arbitrary violations in both the relations *student* and *prof*. Consider for the tuple $\vec{c} = c_1$ the query $Q_{\vec{c}} = \langle q, \{q(c_1) \leftarrow student(x_1, y_1, z, w_1), prof(c_1, y_2, z).\}\rangle$. Then, all components in the discussed factorization for *prof*, except at most one, are singular w.r.t. $Q_{\vec{c}}$, and all components for *student* are singular or (jointly) decomposable. Indeed, the only component for *prof* which may not be singular is the one involving tuples with key $c_1$ (if such a component exists). Using Theorem 5.5, we can evaluate $Q_{\vec{c}}$ in polynomial time w.r.t. the database $D$. By evaluating $Q_{\vec{c}}$ for each constant $c_1$ occurring in $D$, we thus can compute the consistent query answer of $Q$ w.r.t. $D$ as in Equation (12) in polynomial time. Hence, the query $Q$ is tractable. □

This example illustrates that queries of the form of—and in a constraint setting as in—our student and professor example are polynomial in data complexity for arbitrary inputs. This will be formally established by Proposition 5.6.

As a further optimization, we note here that it is possible to do even better than evaluating query $Q_{\vec{c}}$, for each possible tuple $\vec{c} = (c_1, \ldots, c_k)$ in the output. In fact, rather than fixing the output of $Q$ at all positions $i$ to $c_i$, we may fix them one by one. After each step, we test whether sufficiently many components are singular w.r.t. to the modified query $Q'$ (i.e., only few components are neither decomposable nor singular, which means that $Q'$ can be polynomially evaluated); if not, then we fix the next position. This is repeated until the answer is yes or no position remains.

The resulting query $Q'$ covers all output tuples for $Q$ that have at the respective positions the values fixed in $Q'$. Similarly, we construct a modified query $Q''$ for a possible tuple $\vec{c}$ for $Q$ that is not covered by $Q'$; by repeating this process, we will collect a list of

---

[3]Notice that if $n = 0$ (i.e., $\vec{c}$ is void), $ans_c(Q, D) \neq \emptyset$ means that $Q$ is consistently true, false otherwise.

queries $Q', Q'', \ldots$ that we can evaluate using Theorem 5.5, such that the union of all their answers will give us $ans_c(Q, D)$. Noticeably, the singularity tests for the modified queries can be done without evaluating their rewritings as in Definition 5.2. A detailed study and refinement of this strategy remains for future work.

We conclude this section by noting that the tractability result in Example 5.10 can be generalized, and identify an entire class of queries for which we can establish that consistent query answering is polynomial in data complexity without looking at the inconsistencies in the actual database, i.e., we single out a purely "syntactic" condition for tractability.

Let $\mathcal{Q}_{1k\exists}$ denote the class of all conjunctive queries $Q$ (without built-in literals) over schemas that have at most one key constraint per relation, such that except in at most one atom, key positions are always head variables or constants in $Q$. Note that the query in Example 5.10 satisfies this condition.

**Proposition 5.6** *For $\mathcal{Q}_{1k\exists}$, consistent query answering is polynomial in data complexity.*

Note that $\mathcal{Q}_{1k\exists}$ does not fall into any tractable query class in the literature (cf. Section 7).[4] The class may be further generalized, e.g., by allowing between relations $r$ and $s$, where $r$ occurs in the query and $s$ does not, limited exclusion dependencies and inclusion dependencies of the form $r(\vec{x}) \supset s(\vec{x})$ or $s(\vec{x}) \supset r(\vec{x})$ under further syntactic restrictions.

## 6. LOGIC PROGRAMMING FOR CONSISTENT QUERY ANSWERING

According to several proposals in the literature, consistent answers from inconsistent databases can be computed by encoding the constraints in the schema by means of a Datalog program using unstratified negation or disjunction, in such a way that the stable models of this program map to the repairs of the database. A framework that abstracts from several logic programming formalizations in the literature (such as [Greco et al. 2003; Arenas et al. 2003; Barceló and Bertossi 2003]) is introduced next.[5]

**Definition 6.1** Let $Q = \langle q, \mathcal{P} \rangle$ be a non-recursive Datalog$^\neg$ query over $\chi = \langle \Psi, \Sigma \rangle$. A *logic specification for querying $\chi$ with $Q$* is a (safe) Datalog$^{\vee,\neg}$ program $\Pi_\chi(Q) = \Pi_\Sigma \cup \Pi_Q$ such that, for a given $D \in D(\chi)$,

(1)  $rep_\chi(D) \rightleftharpoons \text{SM}(\Pi_\Sigma \cup D)$, and

(2)  $ans_c(Q, D) = Q'[D]$, where $Q' = \langle q, \Pi_\chi(Q) \rangle$, i.e., $ans_c(Q, D) = \{\vec{t} \mid q(\vec{t}) \in M$ for each $M \in \text{SM}((\Pi_\Sigma \cup \Pi_Q) \cup D)\}$, where $\Pi_Q$ is a non-recursive safe Datalog$^\neg$ program,

and $\rightleftharpoons$ denotes a polynomial-time computable correspondence between two sets.  $\square$

In the above definition, $\Pi_\Sigma$ is that portion of $\Pi_\chi(Q)$ that encodes the integrity constraints in $\Sigma$, whereas $\Pi_Q$ represents an encoding of the logic program $\mathcal{P}$ in the user query $Q$ (examples of instantiations of the above logic framework are given in Appendix E).

---

[4]A dichotomy result in [Fuxman and Miller 2007] would suggest that some queries in this class are co-NP-complete (e.g., queries with only nonkey to nonkey joins). However, the proof there assumes that queries have at least one variable in the key positions of each atom, which is not the case for the considered queries.

[5]Other logic formalizations proposed in the data integration setting [Lembo et al. 2002; Bertossi et al. 2002; Calì et al. 2003b; Bravo and Bertossi 2003] also fit in our framework, provided that the *retrieved global database* [Lenzerini 2002] is computed. Notice also that other logic-based approaches to data integration, based on abductive logic programming [Arieli et al. 2004] and ID-logic [Nuffelen et al. 2004], do not fit this framework.
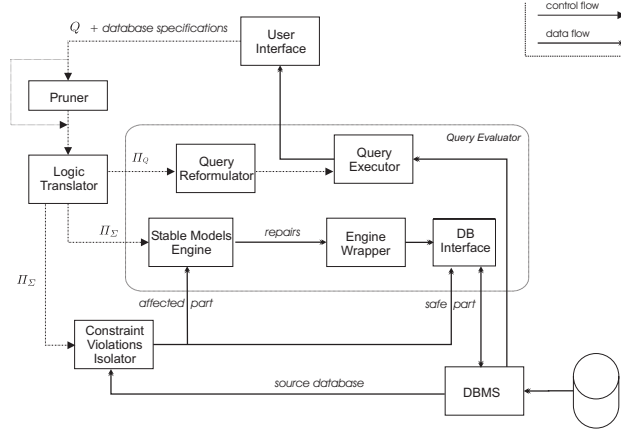
Fig. 3.    System Architecture.

Encoding repair computation by means of logic programs has some attractive features. An important one is that Datalog$^{\vee,\neg}$ programs serve as *executable logical specifications of repair*, and thus provide a language for expressing repair policies in a fully declarative manner rather than in a procedural way. In fact, extensions to the Datalog$^{\vee,\neg}$ language that allow, for instance, to handle priorities and weight constraints [Leone et al. 2006; Simons et al. 2002], provide a useful set of constructs for expressing also more involved criteria that repairs should satisfy, which possibly have to be customized to a particular application scenario (as in [Arenas et al. 2003]).

However, with current (yet still improving) implementations of stable model engines, such as DLV [Leone et al. 2006] or Smodels [Simons et al. 2002], query evaluation over large data sets quickly becomes infeasible because of lacking scalability. The source of complexity in evaluating the program $\Pi_\chi(Q)$ lies in the conflict resolution module $\Pi_\Sigma$. Indeed, while $\Pi_Q$, which is in general a non-recursive Datalog$^\neg$ program, can be evaluated in polynomial time with respect to underlying databases (data complexity) [Dantsin et al. 2001], $\Pi_\Sigma$ is in general a Datalog$^{\vee,\neg}$ program [Greco et al. 2003], whose evaluation data complexity is at the second level of the polynomial hierarchy [Dantsin et al. 2001].

## 6.1    General Architecture for Repair Compilation

The localization properties discussed in Section 4 and Section 5 may be used to optimize consistent query answering from inconsistent databases. Indeed, computing the repairs for $D$ may be done in practice by evaluating the program $\Pi_\Sigma$ only over the affected part of the database $D$, rather than on the whole $D$ as obtained by a straight evaluation of the program $\Pi_\chi(Q)$ over $D$ (Item 1 in Definition 6.1). We thus propose an approach to optimize query answering that implements the strategies in Equation (6) and Equation (10). In practice, we just need an architecture in which a stable model engine used to retrieve one repair at a time is interfaced with a DBMS that evaluates the query over the repair augmented with the safe part of $D$. Figure 3 shows a concrete architecture, whose components have the following functionalities:

- *Pruner:* It takes the user query $Q$ and the schema $\chi$, and produces an equivalent specification (w.r.t. $Q$), stripping off relations and constraints irrelevant for answering $Q$. This

30

is a preprocessing step, which is not discussed in detail here.

- *Logic Translator:* It takes the specification of $\chi$ relevant for $Q$ returned by the Pruner, and produces the logic program $\Pi_\chi(Q) = \Pi_\Sigma \cup \Pi_Q$, according to some encoding proposed in the literature. In our tests, we used the mapping in [Calì et al. 2003b; Grieco et al. 2005].

- *Constraint Violations Isolator:* It is responsible of processing the program $\Pi_\Sigma$ to produce a set of SQL views isolating the safe and the affected parts of the database at hand. When strategies in Section 5.2 are to be applied (cf. Equation (10)), the affected part is in fact provided in terms of a factorization.

- *Stable Models Engine:* It takes as input the affected database (or, in fact, each of the components involved in a factorization, when recombination is going to be implemented via Equation (10)) and computes the repairs using the program $\Pi_\Sigma$. In our implementation, we used the DLV system [Leone et al. 2006].

- *Engine Wrapper:* It wraps the output of the Stable Models Engine, by asking the engine for one repair at time. In our implementation, this is done with the JAVA Wrapper module available for DLV.[6] If the constraints are not in the class $\mathbf{C}_1$, it also has to filter from any repair the facts that are not in the envelope $E$ — see condition $(ii)$ of Proposition 5.1.

- *DB Interface:* It does the interfacing between the Stable Models Engine and the DBMS, in which it stores the repair computed by the Stable Model Engine—in fact, the safe part is not modified in this process. After a new repair is stored, it notifies the Query Executor.

- *Query Reformulator:* It takes the user query and transforms it in a suitable set of SQL statements that can be executed directly over the DBMS.

- *Query Executor:* It is responsible for implementing the recombination step. As discussed in Section 5, based on the query $Q$ and the database $D$, with safe part $S$, it may choose to apply either the repair factorization strategy in Equation (10), possibly with some further optimizations as the one discussed in Section 5.2.2, or the basic approach in Equation (6). As for the strategy in Equation (6), the module stores in the DBMS the result of the execution in a table. When the first repair of the affected part (intersected with $E$), say $R_1$, is processed, the table is initialized with the result of the evaluation of $Q$ over $R_1 \cup S$. Then, for each other repair $R_i$, the table is updated by filtering those tuples that do not occur in the answer to $Q$ over $R_i \cup$. After the last repair is computed, the table is returned to the user. A similar strategy is applied for Equation (10), with the major difference that now the process is repeated for each component involved in the factorization rather than once for the whole affected part. Eventually, to recombine the results of the evaluation over each component (as to implemented Equation (10)), some further temporary tables are used in the DBMS.

Note that in the case where $D$ is consistent, query processing resorts to standard query evaluation over the DBMS, with some overhead for checking constraint violations by the *Constraint Violations Isolator*. In fact, in this case, the *Query Executor* module evaluates the query directly over $S = D$, since no repair is produced by the *Stable Models Engine*.

## 6.2 Grouped Repair Computation

As discussed above, the *Query Executor* module implements the recombination step by executing some SQL statements, for *each* repair computed by the *Stable Models Engine*.

---

[6] *http://www.mat.unical.it/wrapper/index.html*

$player_m{}^{\mathcal{M}_{D_0}}:$

| 10 | Totti | RM | $'11'$ |
|---|---|---|---|
| 9 | Ronaldinho | BC | $'11'$ |

$team_m{}^{\mathcal{M}_{D_0}}:$

| RM | Roma | 10 | $'10'$ |
|---|---|---|---|
| BC | Barcelona | 8 | $'11'$ |
| RM | Real Madrid | 10 | $'01'$ |

$coach_m{}^{\mathcal{M}_{D_0}}:$

| 7 | Capello | RM | $'11'$ |
|---|---|---|---|

Fig. 4.   The database of our running example after marking.

As a further optimization we next consider the idea of grouping these repairs in such a way that a single SQL statement may be evaluated over more than one repair at time. This can be done using a marking strategy, and independently on whether recombination is implemented by means of the basic approach in Section 5.1 or by repair factorization. Indeed, in the former case, we may think of grouping all the repairs of the affected part, while in the latter case, for each component $E_i$ involved in the factorization, we may think of grouping all the repairs of $D \cap E_i$ (see Equation (10)).

Let $R_1, ..., R_n$ be the repairs which we want to simultaneously process on the DBMS, indexed using the order in which the *Stable Models Engine* computes them. In each relation $s$, we add an auxiliary attribute *mark*, leading to a new relation $s_m$. The values for *mark* are strings of bits $0, 1$. To each fact $s(\vec{t}\,) \in D$, we associate a mark $b =' b_1 \dots b'_n$ such that, for every $1 \leq i \leq n$, $b_i = 1$ if $s(\vec{t}\,)$ belongs to $R_i$, and $b_i = 0$ otherwise. The marked tuple $\vec{t}, b$ is stored in the corresponding relation $s_m$. The extensions of all $s_m$ constitute the *marked* database, denoted by $D_m$. Note that the facts in the safe database can be marked without preprocessing: their mark is $'11 \dots 1'$, since they belong to every repair $R_i$. In our running example, the marked database derived from the repairs in Figure 1 is shown in Figure 4. In a first approximation, the marked database may be considered as having its tables altered with an extra column which stores the mark.

A non-recursive Datalog$^\neg$ query $Q = \langle q, \mathcal{P} \rangle$ is reformulated into an SQL query over $D_m$ by first normalizing the rules in $\mathcal{P}$ and then converting each rule $r$ into a separate SQL query $SQL_r$. Let $r : h(\vec{x}_0) \leftarrow B(\vec{x})$ be a safe rule of form

$$p_0(\vec{x}_0) \leftarrow p_1(\vec{x}_1), ..., p_l(\vec{x}_l), \ not\ p_{l+1}(\vec{x}_{l+1}), ..., \ not\ p_{l+k}(\vec{x}_{l+k}),^7 \qquad (13)$$

and let $t_{i,j}$ denote the $j$-th term in $p_i(\vec{x}_i) = p_i(t_{i,1}, \dots, t_{i,k_i})$, where $0 \leq i \leq l + k$ and $1 \leq j \leq k_i$. Then, we associate with $r$ a normalized rule $r'$ obtained from it as follows:

(1) Replace each $t_{i,j}$ by a new variable $y_{i,j}$.
(2) if $t_{i,j}$ is a constant $c$, then add the equality atom $y_{i,j} = c$ to the body;
(3) if $t_{i,j}$ is a variable $x$, then add the equality atom $y_{i,j} = y_{i',j'}$ to the body, where $t_{i',j'}$ is the first occurrence of $x$ in the body of $r$ (from left to right), except for $i = i'$ and $j = j'$. (Note that safety of $r$ guarantees $0 \leq i' \leq l$.)

Informally, $SQL_r$ selects tuples for the head predicate of $r$, thereby respecting not only the join conditions given by the body of $r$, but also the marks of the joined tuples. Marks corresponding to negative literals are inverted and missing tuples (which belong to no repair) are viewed as marked by $'0 \dots 0'$ (see Appendix D for details).

Eventually, all rules, $r_1, \dots, r_\ell$, defining the same predicate $h$ of arity $n$ are collected into a view by the SQL statement $SQL_h$:

CREATE VIEW $h_m(a_1, \dots, a_n, mark)$ AS

---

[7]For the sake of simplicity and w.l.o.g. we assume that variables occurring in $\vec{x}_0$ are all distinct.

```
SELECT a_1, ..., a_n, SUMBIT(mark)
FROM (SQL_{r_1} UNION ... UNION SQL_{r_ℓ})
GROUP BY a_1, ... a_n,
```

where `SUMBIT` denotes an aggregate function that, given $m$ marks (i.e., bit strings), returns the mark given by bitwise OR. By such a view for the query predicate $q$, denoted $q_m$, the consistent answers to the query $Q$ are obtained through the statement $SQL_Q$:

```
SELECT a_1, ... a_n FROM q_m WHERE mark =' 1...1'.
```

It computes the consistent query answers by selecting the facts that are true in all repairs.

**Example 6.1** The query in our running example has two rules: $r_1 : q(x) \leftarrow player(x, y, z)$ and $r_2 : q(x) \leftarrow team(v, w, x)$. Their normalized versions are:

$$r_1' : q(y_{0,1}) \leftarrow player(y_{1,1}, y_{1,2}, y_{1,3}), y_{0,1} = y_{1,1};$$
$$r_2' : q(y_{0,1}) \leftarrow team(y_{1,1}, y_{1,2}, y_{1,3}), y_{0,1} = y_{1,3}.$$

Thus, they translate into corresponding SQL statements $SQL_{r_1}$ and $SQL_{r_2}$:

```
SELECT player_m.Pcode AS a_1,          SELECT team_m.Tleader AS a_1,
       player_m.mark AS mark,                 team_m.mark AS mark,
FROM player_m;                         FROM team_m;
```

Finally, a view for the query predicate $q$ and the final query $SQL_Q$ are expressed as:

```
CREATE VIEW q_m(a_1, mark) AS
       SELECT a_1, SUMBIT(mark)
       FROM (SQL_{r_1} UNION SQL_{r_2})
       GROUP BY a_1;

SELECT a_1 FROM q_m WHERE mark =' 11';
```

$SQL_Q$ yields on $D_m$ the tuples (8), (9), and (10); they are the consistent answers to $Q$. □

The query $SQL_Q$ has the following property (the proof is given in Appendix D).

**Proposition 6.1** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, let $Q$ be a non-recursive Datalog$^\neg$ query over it, and let $R_1, ..., R_n$ be databases such that $R_i = R_i' \cap E$, where $E$ is a weak repair envelope for $D$ and $R_i'$ is a repair for $A = D \cap E$. Then, $SQL_Q$ computes on $D_m$ the set of tuples $\bigcap_{i=1}^n \{\vec{t} \mid \vec{t} \in Q[R_i \cup S]\}$, for $S = D \setminus E$.*

Note that when $R_1', ..., R_n'$ are all repairs for $A$, then the tuples computed by $SQL_Q$ are the consistent answers to $Q$ w.r.t. $D$ — see, again, Equation (6).

A limitation to the scalability of the marking strategy is that all safe tuples must be marked with $'11...1'$, since they belong to each repair. However, we can avoid this, and evaluate a reformulated query on a database instance in which only affected tuples have been marked. In more detail, with each relation symbol $r$, we associate two predicate symbols $r_{safe}$ and $r_{aff}$, which are intended to store the tuples that occur in the safe and the affected part of $r^D$, respectively. Also, we construct the database instance $A'$ by replacing each relation symbol $r$ in $A_D$ with $r_{aff}$, and the database instance $S'$ by replacing each relation symbol $r$ in $S_D$ with $r_{safe}$, i.e., we have that $r_{aff}^{A'} = \{\vec{t} \mid r(\vec{t}) \in A\}$ and $r_{safe}^{S'} = \{\vec{t} \mid r(\vec{t}) \in S\}$. Then, given a query $Q = \langle q, \mathcal{P} \rangle$, where $\mathcal{P}$ is assumed to be normalized, over a schema $\chi = \langle \Psi, \Sigma \rangle$, we proceed as follows:

- for each rule $r : h(\vec{x}_0) \leftarrow B(\vec{x})$ of form (13) belonging to $\mathcal{P}$, we replace each atom $p_j(\vec{x}_j)$ of its positive body, i.e., $1 \leq j \leq l$, by $p_{aff_j}(\vec{x}_j) \vee p_{safe_j}(\vec{x}_j)$;
- we rewrite the resulting rule body into disjunctive normal form $B_1(\vec{x}) \vee \cdots \vee B_n(\vec{x})$;
- we replace in $B_i(\vec{x})$ each negative literal $not\ p_j(\vec{x}_j)$ with a relation $p_j \in \Psi$ by the literals $not\ p_{aff_j}(\vec{x_j}),\ not\ p_{safe_j}(\vec{x}_j)$; let $B_i'(\vec{x})$ be the result;
- we replace $r$ with the rules $r_i : h(\vec{x}_0) \leftarrow B_i'(\vec{x})$, for $1 \leq i \leq n$;
- in the SQL statement $SQL_{r_i}$ for $r_i$, we replace every $p_{safe_{j\,m}}$ by $p_{safe_j}$, and $p_{safe_j}.mark$ by $'1\ldots1'$.

One can show that the SQL reformulation of the query $Q$ as described above, denoted $SQL_Q'$, yields over the partially marked database $S' \cup A_m'$ the same result as $SQL_Q$ over the fully marked database $D_m$. That is, for the reformulation $SQL_Q'$ only the affected tuples have to be marked. Notice that $SQL_Q'$ is exponential in the size of $Q$ (more precisely, in the number of atoms). However, as commonly agreed in the database community, the overhead in query complexity usually pays off the advantage gained in data complexity. With this approach, the additional space depends only on the size of $A$ but not on the size of $S$. For example, for 10 constraint violations involving two tuples each, the required marking space is $2*10*2^{10}$ bits = 2.5 KB, independently of the size of $D$. Furthermore, by allotting 5 MB (=$2*20*2^{20}$ bits) marking space, the technique may scale up to 20 constraint violations, involving two tuples each.

Further optimizations concerning the marking strategy may be carried out, in particular DBMS dependent techniques can be deployed, but are beyond the scope of this paper.

Here we conclude by noticing that throughout this section, we have fixed the length of the markings to coincide with the total number $n$ of repairs of the affected part. Thus, according to Proposition 6.1, the DBMS can be queried just once for recombining the results of the repairs with the safe part and for getting consistent answers. However, the length $n$ of the marks can be fixed independently of the actual number of repairs. Indeed, if $R_1, \ldots, R_m$ are the repairs of the affected part and if $n < m$, then applying Proposition 6.1 $\lceil m/n \rceil$ times (once for each group of $n$ repairs as they are incrementally computed by the *Stable Models Engine*) and intersecting the partial results is sufficient. In the extreme case where $n = 1$, this approach would amount to standard evaluation without markings. Section 8 extensively discusses benefits of this grouping strategy and suggests appropriate values for the length of marking strings.

## 7. OTHER APPROACHES TO CONSISTENT QUERY ANSWERING

Efficient computation of consistent answers to queries on inconsistent databases has received increasing attention recently; see, e.g., [Arenas et al. 1999; Fuxman and Miller 2007; Grieco et al. 2005; Chomicki and Marcinkowski 2005; Chomicki et al. 2004a]. The quoted works single out settings in which this task is feasible with polynomial data complexity, by imposing suitable restrictions on both the form of the constraints in the database schema and on the query language. Briefly, they differ from our work as follows.

- The papers [Arenas et al. 1999; Fuxman and Miller 2007; Grieco et al. 2005] considered only repairs according to the prototypical repair ordering $\leq_D$ introduced in [Arenas et al. 1999]. In contrast, our results cover a generic range of repair orderings, and may be extended to repair semantics based on preference orderings violating the properties in Section 3.1; e.g., Chomicki and Marcinkowski [2005] consider repairs in which a

smallest (in terms of inclusion) set of tuples is deleted from the database but no tuples are added. For such repairs, Proposition 4.2, Lemma 4.3, and Theorem 4.4 can be established similarly.

– Next, all papers above present methods for consistent query answering that hinge on the query and the constraints, but not on the actual database. The papers [Arenas et al. 1999; Fuxman and Miller 2007; Grieco et al. 2005] employ *first-order query rewriting* as the main approach, where the user query $q$ is rewritten into a new query $q'$ expressed in first-order logic, such that the evaluation of $q'$ over every underlying databases $D$ returns the consistent answers to $q$ w.r.t. $D$. Chomicki et al. [2005; 2004a] use a constraint violation hypergraph to reason efficiently about consistent query answers, which also works on an arbitrary database $D$. In contrast, our localization approach looks at the actual database $D$ and exploits the structure of the data and possible repairs. In case of benign violations, consistent query answering can be done efficiently using our approach, also when the first-order rewriting and the hypergraph approach as presented are not applicable.

– Furthermore, our techniques can handle very general and powerful forms of integrity constraints and/or queries that do not fall within the settings studied previously in the literature, which usually make quite limiting assumptions.

In the following we briefly comment on the above mentioned papers, starting from those based on first-order rewriting. The first results in this direction were given by Arenas et al. [1999]. The method proposed in that paper was proved to be sound and complete for queries expressed in quantifier-free first-order logic without disjunction, in the presence of binary universal integrity constraints (a limited fragment of our constraint classes $\mathbf{C}_1$ and $\mathbf{C}_2$). Celle at al. [2000] proposed an extension of the above technique that applies to a slightly more general class of (still binary universal) integrity constraints, and described an implementation of their algorithm. However, the setting considered in that work is still very restrictive.

More recently, Fuxman and Miller [2007] singled out a class of first-order rewritable queries, called $\mathcal{C}_{forest}$, for database schemas containing only key constraints (at most one for each relational predicate). Roughly speaking, $\mathcal{C}_{forest}$ contains conjunctive queries for which joins involving non-key positions must satisfy a particular acyclicity condition. Furthermore, the queries must not contain self-joins, i.e., repeated relation symbols, nor non-full nonkey-to-key joins, i.e., joins between relation symbols $r$ and $s$ that involve a non-key position in $r$ and a strict subset of the key positions of $s$ (and vice versa). The given query rewriting algorithm has been adapted in the ConQuer system [Fuxman et al. 2005] to deal directly with Select-Project-Join queries expressed in SQL; moreover, also aggregate expressions were allowed.

Grieco et al. [2005] extended the class $\mathcal{C}_{forest}$ by allowing some forms of non-full nonkey-to-key joins and also considered additional exclusion dependencies in the database schema. They gave sufficient conditions and algorithms for consistent query answering via first-order rewriting for this setting. To ensure rewritability, only limited interaction between the query atoms and every exclusion dependency between relation symbols $r$ and $s$ is allowed: $r$ and $s$ cannot both occur in the query; if $r$ or $s$ occurs in the query, then the dependency must involve subsets of the keys of $r$ and $s$; other exclusion dependencies between $r$ and $q$ and between $s$ and $p$, respectively, where $p$ and $q$ are (not necessarily distinct) relation symbols occurring in the query, are forbidden.

As mentioned above, the approach of Chomicki et al. [2004a; 2005] to consistent query answering is not based on first-order query rewriting, but constructs a hypergraph that represents the conflicts in the database, and exploits this conflict hypergraph to reason about the consistent query answers by considering independent sets. The technique enables consistent query answering in polynomial time (in data complexity) for specific combinations of denial constraints and queries: projection-free queries in relational algebra (each variable is an output-variable) in presence of arbitrary denial constraints [Chomicki et al. 2004a], and closed simple conjunctive queries (where projections are allowed, but no joins) in presence of functional dependencies over different relations [Chomicki and Marcinkowski 2005]. For the latter queries, all components of the natural factorization of the standard repair envelope w.r.t. any database are clearly singular. Hence, Proposition 5.4 re-establishes that this class of queries can be handled in polynomial time. The approach of Chomicki et al. is implemented in the Hippo System [Chomicki et al. 2004a; 2004b].

The above results are interesting from a practical point of view, and are often supported by experimental validations on large amounts of inconsistent data [Fuxman et al. 2005; Chomicki et al. 2004a; Grieco et al. 2005]. However, they apply to quite narrow settings, with a particular repair semantics and very limited sets of constraints and/or queries.

Our optimization techniques, instead, cover a number of repair orderings from the literature, and address large classes of queries and constraints. Our repair factorization strategy can be applied, for instance, to a setting with generic unions of conjunctive queries on database schemas that contain generic key constraints and exclusion dependencies (but also inclusion dependencies falling in the class $C_1$). It does not impose a priori restrictions on the structure of the schema and the query, but instead conditions on the interaction between the actual inconsistency in the database and the query. For sufficiently benign violation behavior, consistent query answering in polynomial time will be achieved; this may be good enough if non-benign violations rarely occur. Guaranteed worst case polynomial time behavior can be enforced by limiting, like in other approaches, the structure of the query and the constraints.

Finally, we point out that our optimization techniques may be also combined with the other approaches. For example, we can exploit factorization to single out cases that are not first-order rewritable in general, but are so in the light of the actual conflicts that we localize in the database. This approach seems to be promising, as it enables the exploitation of consolidated relational database technology also in some settings that cannot be handled by current first-order rewriting techniques.

## 8. EXPERIMENTAL RESULTS

In this section, we present experimental results for evaluating the effectiveness of our approach and, specifically, the benefits of the localization techniques discussed in the paper.

### 8.1 Benchmark Databases and Compared Methods

As already mentioned, Hippo [Chomicki et al. 2004a; 2004b] and ConQuer [Fuxman et al. 2005; Fuxman and Miller 2007] are two noticeable prototype systems for consistent query answering from inconsistent databases. These systems are tailored for specific settings where this task is tractable and manage very specific classes of queries and constraints. For this reason, their performances have been tested on ad-hoc created benchmark databases; Fuxman et al. [2005] mainly generated synthetic data for the TCP-H specifications over a schema containing only single keys on relation symbols, and used queries of the

class $\mathcal{C}_{forest}$, encoded into SQL, with aggregate expressions. Chomicki et al. [2004a] considered project-free queries over ternary relations and functional dependencies of the form $p(x, y, z) \wedge p(x, y', z') \supset z = z'$ over each such relation.

To assess the effectiveness of our localization approach, we need some novel scenario as it is not directly comparable with Hippo and Conquer, whose incomparability similarly requested specific benchmarks and data. Indeed, our techniques are designed for and can be used also in settings more general than those addressed by those systems (see Section 7).

On the other hand, if we focus on the classes of queries for which Hippo and ConQuer have been designed, it will come as no surprise that our approach pays in efficiency for its generality and expressiveness. And, in fact, we envisage an integrated architecture that switches to these more specialized and efficient systems whenever the query and the constraints fall in one of the classes they are able to deal with. To test our framework and the factorization techniques discussed in Section 5, we thus proceed as follows:

- We first present experimental results for our running example (on football teams). The results show the advantages of focusing the computation by making use of the techniques discussed in Section 6.2, through the system described in Section 6.1, with respect to direct evaluating, through the DLV system, the logic specification for querying the inconsistent database with the query at hand.

- We then focus on a test suite over the database schema $\chi_f^2$ used in [Chomicki et al. 2004a], which contains two ternary relations $r_1$ and $r_2$, and the functional dependencies $r_i(x, y, z) \wedge r_i(x, y', z') \supset z = z'$, with $i = 1, 2$, but we consider queries that involve projections, so that the system Hippo is not applicable. In this experiment, we show scalability of our approach with respect to a growing number of tuples in conflicts, and the advantage of combining repair factorization with markings.

- We discuss the impact of the number of atoms in the query on the performance of the localization approach, by considering the schema $\chi_f^N$ which generalizes $\chi_f^2$ with an increasing number of predicates.

- Finally, we focus on a test suite that is based on the schema reported in Example 5.10, which refers to a typical information system supporting university administration. Experiments with this scenario aim to show the benefit of the factorization technique and the achievable scalability, in a context that is beyond the scope of applicability of other approaches in the literature focused on tractable settings of consistent query answering.

For the schemas, we generated random data following the idea of tuning the size of the safe part and the number of conflicts as in [Chomicki et al. 2004a; Fuxman et al. 2005].

All experiments have been carried out on a 1.6GHz Pentium IV with 512MB memory, by assessing the time needed for consistent query answering when the DLV system computes repairs of the affected part only, plus the time required for the recombination of the results in PostgreSQL. We always used a repair semantics based on the prototypical preorder $\leq_D$.

## 8.2 The Football Teams Example

We next discuss the performances of our approach and, specifically, its scaling w.r.t. the size of the safe database, on a simple scenario. For our running example, we built a synthetic data set $D_{FT}$, such that tuples in $coach$ and $team$ satisfy the key constraints issued on these relations, while tuples in $player$ violate the corresponding key constraint. Each violation consists of two facts that coincide on $Pcode$ but differ on either $Pname$ or $Pteam$;
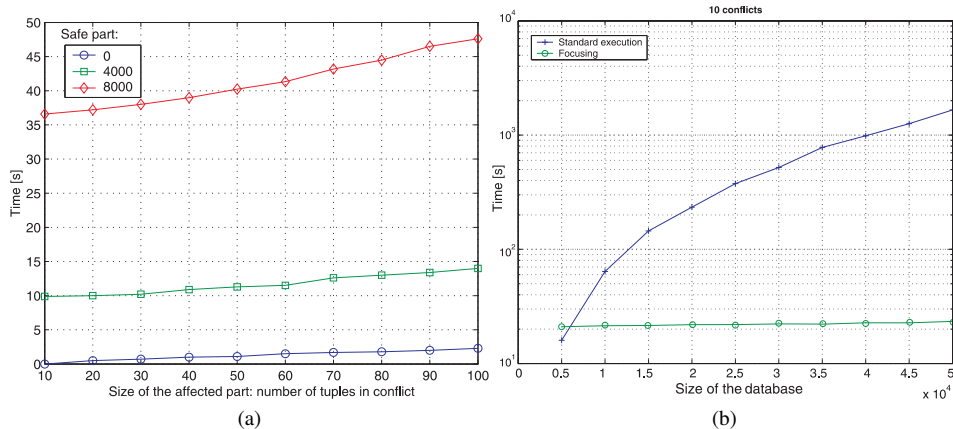
Fig. 5. Football Team. (a) Execution time in DLV system w.r.t. size of the affected part. (b) Comparison with the optimization method.

these facts constitute the affected part of $D_{FT}$. For our experiments, we consider the query $Q = \langle q, \mathcal{P} \rangle$ where $\mathcal{P} = \{q(x) \leftarrow player(x, y, z); \ q(x) \leftarrow team(v, w, x)\}$, and we encode our problem into a Datalog$^{\neg}$ program $\Pi_{\chi_0}(Q)$ in the line of [Calì et al. 2003b; Grieco et al. 2005] (the encoding used is the one given in Appendix E.1, in which we get rid of the encoding for the mapping). We first measure the execution time of the program $\Pi_{\chi_0}(Q)$ in DLV depending on the size of the affected part, while the size of the safe part is fixed to the values(i) 0, (ii) 4000, and (iii) 8000, respectively. We stress that values for the execution time of the DLV system refer to query answering with non-ground queries.

The results for this experiment, reported in Figure 5.(a), show that the DLV system scales well w.r.t. the size of the affected part. Still the big size of the safe part appears to be the most limiting factor for an efficient implementation. Indeed, only 8000 facts (in absence of conflicting tuples) would require more than 35 seconds for consistent query answering.

The performance degradation under varying database size is further stressed in Figure 5.(b), which shows a comparison (in log-scale) between consistent query answering using a single DLV program and the optimization approach proposed in this paper. As for the optimization approach, values on execution time include the cost of computing repairs of the affected database only, plus marking and evaluating the associated SQL query over marked relations. Specifically, we considered 10 violations and a marking string of $2^{10}$ bits, such that issuing one query over the database is sufficient to recombine the repairs of the affected part with the safe part. Interestingly, the growth of the running time of our optimization method under a varying database size is negligible.

### 8.3 Scalability Assessment

In experiments below, we assessed the relevance of the strategy for grouping repair computation by focusing on the database $\chi_f^2$. Indeed, so far, we have assumed that the marking string is sufficient for storing all repairs for the affected part and, therefore, the DBMS has been queried just once for recombining the results of the localized repairs with the safe part only. But, the reader may at this point wonder whether this approach is more efficient than processing each repair sequentially (one at a time).

Figure 6.(a) answers the above question positively. It reports the time needed for answer-
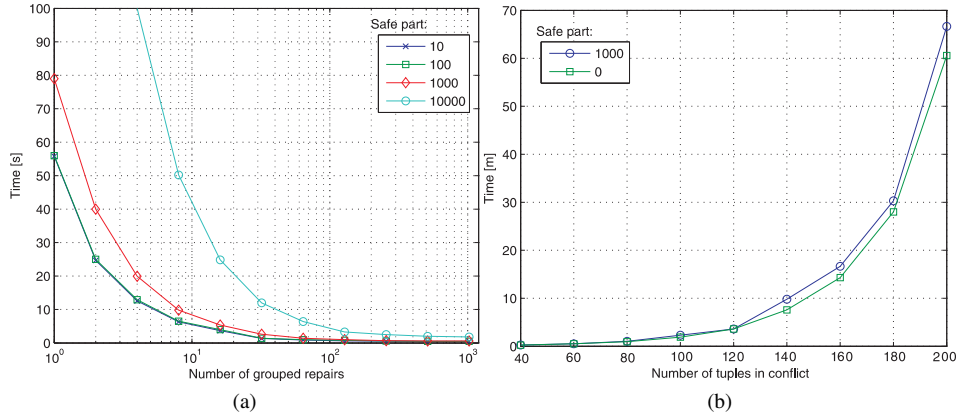
Fig. 6. Query answering over $\chi^2_f$. (a) Optimization method w.r.t. number $n$ of grouped repairs, for a fixed number of conflicts. (b) Optimization method w.r.t. the size of the affected part, for $n = 2^9$.

ing the query $Q_f = \langle q_f, \mathcal{P}_f \rangle$ where $\mathcal{P}_f = \{q_f(y_1) \leftarrow r_1(x, y_1, z_1), r_2(x, y_2, z_2)\}$ w.r.t. the number $n$ of repairs that are grouped and processed simultaneously on the DBMS. Specifically, we fixed 10 conflicts in the data (each involving two inconsistent tuples). Hence, for $n = 1$, we sequentially process each repair, while for $n = 2^{10}$, all the repairs are combined in the DBMS at the same time. The advantage of grouping repairs is evident, specifically by considering the scaling of the curves for different sizes of the safe part. Actually, we note that as the markings may grow exponentially with the size of the affected part, processing all the repairs at the same time is generally infeasible, since the length of the marking strings may exceed the maximum size allowed by the DBMS. For instance, with the bit string type of the PostgreSQL system, we may store marks up to $n = 2^{25}$ bits. In fact, to scale up to a few Gigabytes we may resort to the *large objects* facilities of the system, or we may use well-known commercial DBMSs that provide embedded support for dealing with large binary objects. We point out that in our experiments we did not exploit such features, which instead may profitably be used within an engineered version of the prototype; indeed, in the following we used an incremental evaluation approach, eventually by grouping up to 1000 repairs at time—which is a value we experienced to be appropriate for the use with PostgreSQL.

In a second set of experiments over the query $Q_f$, we measured again the scalability w.r.t. the number of conflicts. In particular, we augmented the number of conflicts up to 100, and we fixed the marking string to $2^9$ bits. The results shown in Figure 6.(b) evidence the exponential scaling in the number of conflicting tuples; this is indeed the best scaling we can expect for inconsistent databases in general, as the problem is co-NP-hard.

In fact, it is interesting to assess whether some nicer scaling can be obtained by applying the factorization strategy discussed in Section 5.2. In this respect, we notice that the setting above is such that our factorization strategy can be applied. Indeed, the setting we are considering is basically the one described in Example 5.2, where each component contains only those facts witnessing a violation of the functional dependency over each of the two relations $r_1$ and $r_2$. Specifically, in our experiments, we fixed the structure of each component to contain 20 tuples and 1000 repairs (any pair of these tuples witnesses a violation of the dependency), and we generated some synthetic data for increasingly large number
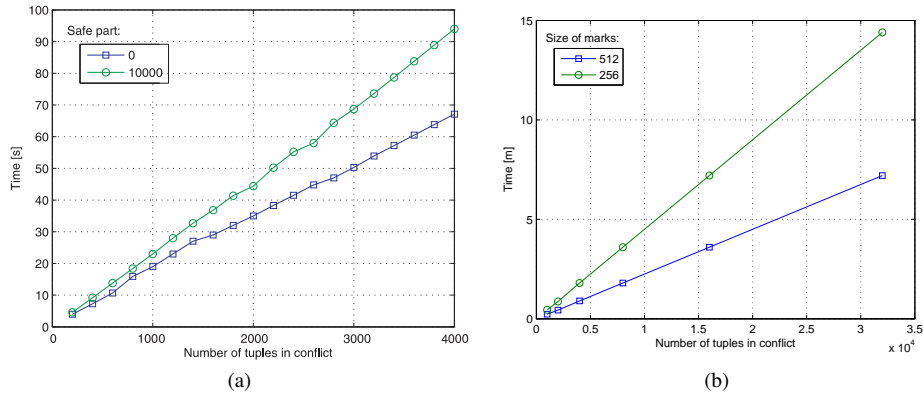
39

Fig. 7. Factorization strategy. Answering $Q_f$ over $\chi_f^2$ w.r.t. the size of the affected part: (a) Dependency on the safe part. (b) Dependency on the marking string for $Q_f'$.

of independent components. In addition to the factorization strategy, we still exploit the grouping repair approach, by fixing the number of repairs simultaneously processed to $2^{10}$.

The results obtained by applying the recombination strategy in Equation (10) are shown in Figure 7.(a). Given the ability of independently processing the components, the scaling is now linear in the number of components and, hence, in the size of the whole affected part. In fact, query answering is feasible for a much larger number of constraint violations.

A similar experiment has been repeated for the query $Q_f' = \langle q_f', \mathcal{P}_f' \rangle$ where $\mathcal{P}_f' = \{q_f'(y_1) \leftarrow r_1(x, y_1, z_1)\}$. After fixing the safe part to 10,000 tuples, we repeated the experiment up to a very large number of conflicts considering two different values for the parameter $n$ which bounds the number of simultaneously processed repairs. Note in Figure 7.(b) the linear scaling in the number of tuples in conflicts and the benefit from combining marking with the factorization strategy. In fact, this linear scaling is again due to the fact that components can independently be processed, so that the cost for query answering basically amounts to computing all repairs for each of these components and store them as marks into the database. As an extreme scenario, we note that the linear behavior of these latter tasks has been confirmed up to a million of tuples, for which the computation was about 9 hours.

Finally, the setting $\chi_f^2$ has been generalized. We also considered the database $\chi_f^N$ and the query $Q = \langle q, \mathcal{P} \rangle$ where $\mathcal{P} = \{q(y_1) \leftarrow r_1(x, y_1, z_1), r_2(x, y_2, z_2), ..., r_N(x, y_N, z_N)\}$, for 10 constraint violations per relation and $2^{10}$ repairs simultaneously processed. In this scenario, we performed some experiments to assess the dependence of query answering on the number of atoms $N$ in the query. The results are reported in Figure 8.(a), which shows (as discussed in Section 6.2) an exponential dependency.

We conclude this overview on the behavior of the localization strategies presented in this paper by considering the schema in Example 5.10. As a first case, we consider the query $Q' = \langle q, \{q(w_1) \leftarrow student(x_1, y_1, z, w_1), prof(x_2, y_2, z).\} \rangle$, which is a slight modification of the query $Q$ in such example and which asks for addresses of those students that might possibly be involved in family relationships with professors.

Note that a full-inclusion dependency is issued over the schema (in addition to the keys), and that in $Q'$ there is a join involving only non-key positions, and also projection. Com-
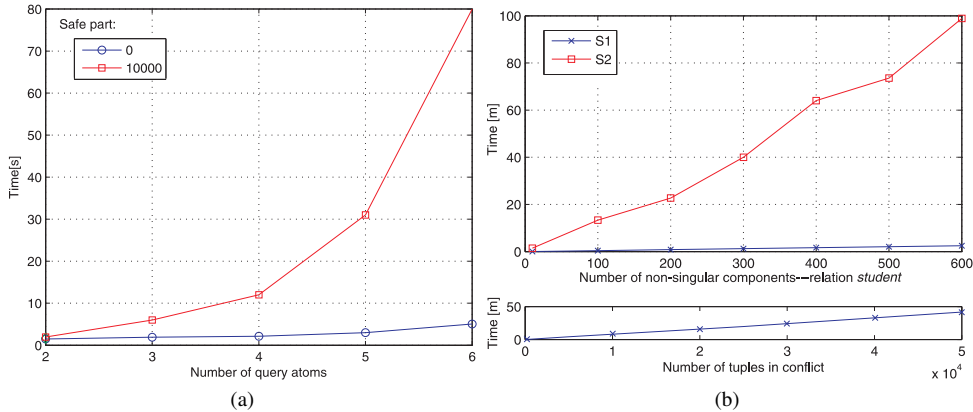
Fig. 8. (a) Query answering over $\chi_f^N$, $N > 2$. (b) Results for Example 5.10.

puting the consistent answers to such a query is co-NP-hard; indeed, the reader may check that $Q'$ can be used to encode the MONOTONE-3SAT problem with minor modifications over the construction presented in [Fuxman and Miller 2007] (which is originally from [Chomicki and Marcinkowski 2005]) to prove hardness of consistent query answering for queries over two distinct body literals.

In these experiments, the safe part is fixed to 10,000 tuples (9,850 students and 150 professors), while the affected part has been generated according to two different scenarios:

(S1) The components over the relation $prof$ are singular (in total we fixed 50 such components), while the components over $student$ are non-singular but decomposable (experiments have been conducted for different component sizes). Each component consists of 10 conflicting tuples. This scenario models an information system where data pertaining to professors are managed centrally, so that inconsistencies over last names may hardly emerge (last names of professors often act in practice as an identifier). However, data about students are assumed to be unreliable; they may, for instance, result from naive integration of autonomous data sources residing in different faculties or schools;

(S2) In this scenario, five non-singular components are associated with $prof$, each one containing two tuples in conflict. Hence, this time, we assume that $prof$ is also the result of some integration task (e.g., data wrapped from the Web); and, beside 50 professors in conflict over first names, five additional conflicts over last names emerge.

In all the scenarios above, repair grouping is exploited with $n = 512$.

The results for this set of experiments are reported in Figure 8.(b). We note that in the lower diagram, in the case of (S1) the scaling is basically linear in the number of conflicting tuples in the relation $student$ and, in fact, in the number of non-singular components.

On the other hand, processing the non-singular components associated with $prof$ in (S2) causes a performance degradation depending on the number of non-singular components of $student$. Given the bound on the number of these components, the scaling is again linear, but with faster growth rate which exponentially depends on the number of professors in conflicts over last names, as for it emerges from the upper diagram in Figure 8.(b).

For a further remark on this experiment, consider again the original query $Q$ introduced in Example 5.10, which asks for last name of professors possibly involved in family rela-

41

tionships with students. This query does not fall in any tractable class of queries singled out in the literature (cf. Section 7). However, according to our discussion in Section 5.2.2, the grounding technique coupled with factorization may indeed ensure tractability. Basically, for each possible last name of professor, say $c$, we have to compute the result of the ground query: $Q_c = \langle q, \{q(c) \leftarrow student(x_1, y_1, z, w_1), prof(c, y_2, z).\} \rangle$. Given that now we focus on precisely one component for $prof$, performances for answering $Q_c$ are even faster than those registered for (S1). And, the global performances for evaluating $Q$ will linearly scale in the number of distinct last names of professors in the data.

In closing this section, we would like to summarize the lesson learned from our experiments, which might give some guidelines for further investigations into consistent query answering from inconsistent database.

On the one hand, our activity has certainly provided some bad news. First, the scalability of "pure" logic-programming based approaches for consistent query answering is in many cases not suited for real world applications. This is related to the intrinsic complexity of the problem, which is the prize to be paid for the generality of these approaches. And, second, if the data can not be factorized (e.g., by means of techniques in Section 5.2), then there is little chance to answer queries over large data sets (actually, large affected databases), even when logic-programming based approaches are used in combination with our grouping and marking techniques. On the other hand, there are two good news:

– First, in those scenarios where the size of the affected part is not very large, marking strategies can be very effective to support consistent query answering even if the data can not be factorized. Indeed, there are substantial benefits w.r.t. basic approaches where safe and affected data are not distinguished in the computation. Clearly enough, if these scenarios fall in one of the syntactically tractable classes already known in the literature, then applying the proper rewriting is definitively the best choice. However, for general classes of queries and constraints, our techniques represent a viable option, given the infeasibility of directly applying logic-programming based approaches. We point out that we experienced such kind of scenarios in applications related to data integration, and arguably they can be found in several different contexts (e.g., re-engineering of legacy systems, semantic data retrieval over the web, etc). In this regard, setting up a benchmark suite of real datasets coming from practical scenarios should be a primary goal for the community.

– Second, localization and factorization strategies may support well consistent query answering up to a large number of conflicts if the analysis of the data reveals some nice structure leading to the possibility of isolating singular and/or decomposable components. A polynomial scaling can be obtained even in some interesting settings that cannot be dealt with by approaches only focused on limiting the form of the input queries and integrity constraints (roughly, localization and factorization techniques can isolate easy instances for hard problems, on the basis of data inspection). This perspective complements previous results in the literature. In particular, if all cases in which non-singular components are jointly decomposable, then a linear scaling can eventually be achieved for common classes of constraints and queries. However, when non-decomposable components emerge, this nice scaling tends to deteriorate (exponentially in the number of components). While this trend can not be avoided in the limit (since in the extreme case where a linear number of these components affects the data, query answering is actually

42

co-NP-hard), the factorization strategy might be refined, in particular when a small number of non-singular and non-decomposable components are in the data. In fact, specific elaborations in this direction constitute an interesting avenue for further research.

## 9. CONCLUSION

For optimizing logic-programming based query answering from inconsistent databases, we have presented a repair localization approach. In this approach, repairs are conceptually confined to a repair envelope, which intuitively comprises the part of the database affected by inconsistency, and then recombined with the unaffected (safe) part before determining the query result. We have investigated this approach in a generic framework accommodating different classes of integrity constraints (including denial constraints [Chomicki et al. 2004a]), and preference orderings for repairs from the literature (see Section 3.1). We then have discussed how this approach can be fruitfully utilized for query answering using logic programming specifications, where a logic programming engine and a DBMS are combined, such that tremendous performance gains are achieved.

While motivated by logic programming specifications, our localization results are not bound to such a setting and are, in fact, applicable to any realization of consistent query answering. Furthermore, the generic form of preferences, constraints, and repair envelopes allows to instantiate the results to many different concrete settings in practice.

The work presented here can be extended in different directions. As for localization and query answering, our results may be extended to repair semantics based on preference orderings violating the properties in Section 3.1; e.g., the one in [Chomicki and Marcinkowski 2005]. Furthermore, for any negation-free query $Q$ the intersection of the answers $Q[R]$ on all repairs is equivalent to answering it only on the repairs which are minimal under set-inclusion, i.e., do not contain any other repair properly. If an ordering $\leq_D$ fails to satisfy (SIP), (DPE), and (DIS), we may characterize the repairs w.r.t. $\leq_D$ that are minimal under set inclusion as the repairs under another ordering $\leq_D'$ which satisfies these properties. An example is the ordering $R_1 \sqsubseteq_D R_2$ iff $R_1 \cap D \supseteq R_2 \cap D$ [Calì et al. 2003a; 2003b], which violates (SIP). Intuitively, in such an ordering, violations are preferably repaired by adding tuples to $D$, and minimally deleting tuples, in case adding tuples is not sufficient to repair inconsistency. We can use here the ordering $R_1 \sqsubseteq_D' R_2$ iff $R_1 \sqsubseteq_D R_2 \wedge (R_1 \cap D = R_2 \cap D \Rightarrow R_1 \backslash D \subseteq R_2 \backslash D)$ instead for answering $Q$.

Other approaches considered consistent query answering under the perspective of modifying values in the database rather than entire tuples [Wijsen 2005]. Due to the different semantics considered in these works, such repairs are not immediately captured by our framework. A study of respective extensions is left for future work.

Another extension of the results here is from a single database to a data integration system $\mathcal{I} = \langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$, where $\mathcal{G}$ is the global schema, $\mathcal{S}$ is the schema of the various sources, and $\mathcal{M}$ is the mapping establishing the relationship between $\mathcal{G}$ and $\mathcal{S}$ [Lenzerini 2002]. As briefly discussed in Appendix E.1 and more in detail in [Eiter et al. 2005], the results developed here can be readily adapted for a Global-As-View (GAV) setting in which $\mathcal{M}$ is given by stratified Datalog queries, and for constraints on the global schema falling in the classes considered in this paper. They can be further extended to other GAV settings, e.g., as in [Lembo et al. 2002; Calì et al. 2003b], and certain Local-As-View (LAV) settings, e.g., as in [Bertossi et al. 2002; Bravo and Bertossi 2003].

In fact, most of the research reported here has been carried out within the EU project

INFOMIX on advanced data integration for expressive schemas using logic programming. However, the INFOMIX system is not the implementation of all results in this paper. For more information about the project, see [Leone et al. 2005].

ELECTRONIC APPENDIX

The electronic appendix for this article can be accessed in the ACM Digital Library by visiting the following URL: `http://www.acm.org/pubs/citations/journals/tods/2008-V-N/p1-URLend`.

Acknowledgments

REFERENCES

ABITEBOUL, S., HULL, R., AND VIANU, V. 1995. *Foundations of Databases*. Addison Wesley.

ARENAS, M., BERTOSSI, L., AND CHOMICKI, J. 2001. Scalar aggregation in fd-inconsistent databases. In *Proc. of the 8th Int. Conf. on Database Theory (ICDT'01)*. Springer, 39–53.

ARENAS, M., BERTOSSI, L., AND CHOMICKI, J. 1999. Consistent query answers in inconsistent databases. In *Proc. of the 18th ACM SIGACT SIGMOD Symp. on Principles of Database Systems (PODS'99)*. 68–79.

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.

ARIELI, O., DENECKER, M., NUFFELEN, B. V., AND BRUYNOOGHE, M. 2004. Coherent integration of databases by abductive logic programming. *J. of Artificial Intelligence Research 21*, 245–286.

BARCELÓ, P. AND BERTOSSI, L. 2003. Logic programs for querying inconsistent databases. In *Proc. of the 5th Int. Symp. on Practical Aspects of Declarative Languages (PADL'03)*. 208–222.

BERTOSSI, L. AND CHOMICKI, J. 2003. Query answering in inconsistent databases. In *Logics for Emerging Applications of Databases*, J. Chomicki, R. van der Meyden, and G. Saake, Eds. Springer, Chapter 2, 43–83.

BERTOSSI, L., CHOMICKI, J., CORTES, A., AND GUTIERREZ, C. 2002. Consistent answers from integrated data sources. In *Proc. of the 6th Int. Conf. on Flexible Query Answering Systems (FQAS'02)*. 71–85.

BERTOSSI, L., HUNTER, A., AND SCHAUB, T., Eds. 2005. *Inconsistency Tolerance [result from a Dagstuhl seminar]*. LNCS, vol. 3300, Springer.

BOUZEGHOUB, M. AND LENZERINI, M. 2001. Introduction to the special issue on data extraction, cleaning, and reconciliation. *Information Systems 26,* 8, 535–536.

BRAVO, L. AND BERTOSSI, L. 2003. Logic programming for consistently querying data integration systems. In *Proc. of the 18th Int. Joint Conf. on Artificial Intelligence (IJCAI'03)*. 10–15.

BRAVO, L. AND BERTOSSI, L. 2005. Disjunctive deductive databases for computing certain and consistent answers to queries from mediated data integration systems. *J. of Applied Logic 3,* 1, 329–367.

CALÌ, A., LEMBO, D., AND ROSATI, R. 2003a. On the decidability and complexity of query answering over inconsistent and incomplete databases. In *Proc. of the 22nd ACM SIGACT SIGMOD Symp. on Principles of Database Systems (PODS'03)*. 260–271.

CALÌ, A., LEMBO, D., AND ROSATI, R. 2003b. Query rewriting and answering under constraints in data integration systems. In *Proc. of the 18th Int. Joint Conf. on Artificial Intelligence (IJCAI'03)*. 16–21.

CELLE, A. AND BERTOSSI, L. 2000. Querying inconsistent databases: Algorithms and implementation. In *Proc. of the Int. Conf. on Computational Logic (CL'00)*. 942–956.

CHOMICKI, J. 2007. Consistent query answering: Five easy pieces. In *Proc. of the 11th Int. Conf. on Database Theory (ICDT'07)*, T. Schwentick and D. Suciu, Eds. LNCS, Springer, 1–17.

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 *Proc. of the 13th ACM Conf. on Information and Knowledge Management (CIKM'04)*. ACM Press, 417–426.

CHOMICKI, J., MARCINKOWSKI, J., AND STAWORKO, S. 2004b. Hippo: A system for computing consistent answers to a class of SQL queries. In *Proc. of the 9th Int. Conf. on Extending Database Technology (EDBT'04)*. LNCS, vol. 2992, Springer, 841–844.

DANTSIN, E., EITER, T., GOTTLOB, G., AND VORONKOV, A. 2001. Complexity and expressive power of logic programming. *ACM Computing Surveys 33,* 3, 374–425.

EITER, T., FINK, M., GRECO, G., AND LEMBO, D. 2003. Efficient evaluation of logic programs for querying data integration systems. In *Proc. of the 19th Int. Conf. on Logic Programming (ICLP'03)*, C. Palamidessi, Ed. LNCS vol. 2916, Springer, 163–177.

EITER, T., FINK, M., GRECO, G., AND LEMBO, D. 2005. Optimization methods for logic-based query answering from inconsistent data integration systems. Tech. Rep. INFSYS RR-1843-05-05, TU Wien. (July).

EITER, T., GOTTLOB, G., AND MANNILA, H. 1997. Disjunctive Datalog. *ACM Transactions on Database Systems 22,* 3, 364–418.

FAGIN, R., ULLMAN, J. D., AND VARDI, M. Y. 1983. On the semantics of updates in databases. In *Proc. of the 2nd ACM SIGACT SIGMOD Symp. on Principles of Database Systems (PODS'83)*. 352–365.

FUXMAN, A., FAZLI, E., AND MILLER, R. J. 2005. ConQuer: Efficient management of inconsistent databases. In *Proc. of the ACM SIGMOD Int. Conf. on Management of Data*. 155–166.

FUXMAN, A. AND MILLER, R. J. 2007. First-order query rewriting for inconsistent databases. *J. of Computer and System Sciences 73,* 4, 610–635.

GELFOND, M. AND LIFSCHITZ, V. 1991. Classical Negation in Logic Programs and Disjunctive Databases. *New Generation Computing 9*, 365–385.

GRECO, G., GRECO, S., AND ZUMPANO, E. 2003. A logical framework for querying and repairing inconsistent databases. *IEEE Trans. on Knowledge and Data Engineering 15,* 6, 1389–1408.

GRIECO, L., LEMBO, D., RUZZI, M., AND ROSATI, R. 2005. Consistent query answering under key and exclusion dependencies: Algorithms and experiments. In *Proc. of the 14th Int. Conf. on Information and Knowledge Management (CIKM'05)*. 792–799.

KIFER, M. AND LOZINSKII, E. L. 1992. A logic for reasoning with inconsistency. *J. of Automated Reasoning 9,* 2, 179–215.

KOWALSKI, R. A. AND DADRI, F. 1990. Logic programming with exceptions. In *Proc. of the 7th Int. Conf. on Logic Programming (ICLP'90)*. 490–504.

LEMBO, D., LENZERINI, M., AND ROSATI, R. 2002. Source inconsistency and incompleteness in data integration. In *Proc. of the 9th Int. Workshop on Knowledge Representation meets Databases (KR'02)*.

LENZERINI, M. 2002. Data integration: A theoretical perspective. In *Proc. of the 21st ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems (PODS'02)*. 233–246.

LEONE, N., EITER, T., FABER, W., FINK, M., GOTTLOB, G., GRECO, G., IANNI, G., KALKA, E., LEMBO, D., LENZERINI, M., LIO, V., NOWICKI, B., ROSATI, R., RUZZI, M., STANISZKIS, W., AND TERRACINA, G. 2005. The INFOMIX system for advanced integration of incomplete and inconsistent data. In *Proc. of the ACM SIGMOD Int. Conf. on Management of Data*. 915–917.

LEONE, N., PFEIFER, G., FABER, W., EITER, T., GOTTLOB, G., PERRI, S., AND SCARCELLO, F. 2006. The DLV System for Knowledge Representation and Reasoning. *ACM Trans. Computational Logic 7,* 3, 499–562.

LIN, J. 1996. Integration of weighted knowledge bases. *Artificial Intelligence 83,* 2, 363–378.

LIN, J. AND MENDELZON, A. O. 1998. Merging databases under constraints. *Int. J. of Coop. Information Systems 7,* 1, 55–76.

NUFFELEN, B. V., CORTÉS-CALABUIG, A., DENECKER, M., ARIELI, O., AND BRUYNOOGHE, M. 2004. Data integration using ID-logic. In *Proc. of the 16th Int. Conf. on Advanced Information Systems Engineering (CAiSE'04)*, A. Persson and J. Stirna, Eds. LNCS, vol. 3084, Springer, 67–81.

SIMONS, P., NIEMELÄ, I., AND SOININEN, T. 2002. Extending and Implementing the Stable Model Semantics. *Artificial Intelligence 138*, 181–234.

STAWORKO, S., CHOMICKI, J., AND MARCINKOWSKI, J. 2006. Preference-driven querying of inconsistent relational databases. In *EDBT Workshops*, T. Grust, et al., Eds. LNCS, vol. 4254, Springer, 318–335.

ULLMAN, J. D. 1989. *Principles of Database and Knowledge Base Systems*. Computer Science Press.

WIJSEN, J. 2005. Database repairing using updates. *ACM Transactions on Database Systems 30,* 3, 722–768.

# Repair Localization for Query Answering from Inconsistent Databases

THOMAS EITER and MICHAEL FINK
Technische Universität Wien
and
GIANLUIGI GRECO
Università della Calabria
and
DOMENICO LEMBO
SAPIENZA Università di Roma

## A. PROOFS FOR SECTION 3

**Proposition 3.2** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, where all constraints in $\Sigma$ are safe. Suppose that $<_D$ satisfies (SIP). Then, every repair $R \in rep(D)$ involves only constants from $adom(D, \chi)$, and some repair exists if $\chi$ is consistent.*

PROOF. Let $R$ be an arbitrary database of $\chi$ consistent with $\Sigma$. Let $R'$ result from $R$ by removing every fact containing some constant $c \notin adom(D, \chi)$. We show that $R'$ is a repair. We first show that $R' \models \Sigma$. Towards a contradiction, assume that $R' \not\models \Sigma$. Hence, there exists a ground instance $\sigma^g$ of some constraint $\sigma \in \Sigma$ of form $A_1(\vec{c_1}) \wedge \cdots \wedge A_l(\vec{c_l}) \supset B_1(\vec{d_1}) \vee \cdots \vee B_m(\vec{d_m}) \vee \phi_1(\vec{e_1}) \vee \cdots \vee \phi_n(\vec{e_n})$ which is violated by $R'$, i.e., (i) $A_1(\vec{c_1}), \ldots, A_l(\vec{c_l}) \in R'$, (ii) $B_1(\vec{d_1}), \ldots, B_m(\vec{d_m}) \notin R'$, and (iii) $\phi_1(\vec{e_1}) \vee \cdots \vee \phi_n(\vec{e_n})$ is false. Since $R \models \sigma^g$, by construction of $R'$ we have $B_j(\vec{d_j}) \in R \setminus R'$ for some $j \in \{1, \ldots, m\}$ and thus $\vec{d_j}$ contains some constant $c \notin adom(D, \chi)$. It follows that some variable occurring in the head of $\sigma$ does not occur in the body of $\sigma$; that is, $\sigma$ is not safe, which is a contradiction. Since $R$ and $R'$ differ only for facts outside $D$, we have that $D \setminus R = D \setminus R'$, and since $R' \subset R$, we have that $R' \setminus D \subset R \setminus D$. Therefore, $\triangle(R', D) \subset \triangle(R, D)$, and thus by (SIP) $R' <_D R$. The $<_D$-minimality of repairs implies that each database in $rep(D)$ involves only constants from $adom(D, \chi)$. Furthermore, by consistency of $\chi$ and the fact that each sequence $R_1 >_D R_2 >_D \cdots R_i >_D \cdots$ of databases $R_i$ on $adom(D, \chi)$ must be finite, one such repair $R$ must exist. $\square$

**Proposition 3.3** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$ where no built-in relations occur in $\Sigma$ except $=$ and $\neq$. Suppose that $<_D$ satisfies (SIP). Then, every repair $R \in rep(D)$*

*involves only constants from $adom(D, \chi)$, and some repair exists if $\chi$ is consistent.*

PROOF. Following the argumentation in the proof of Proposition 3.2, consider a ground instance $\sigma^g$ of some $\sigma \in \Sigma$, which is violated by $R'$. Then, some atom $B_j(\vec{c}_j) \in R \setminus R'$ in the head of $\sigma^g$ exists such that $\vec{c}_j = c_{j,1}, \ldots, c_{j,n_j}$ contains some constant $c_{j,h} \notin adom(D, \chi)$ and the respective variable $y_{j,h}$ in the atom $B_j(\vec{y}_j)$ in $\sigma$ does not occur in the body of $\sigma$ (notice that if $y_{j,h}$ would occur in the body, $c_{j,h}$ might not be outside $adom(D, \chi)$, because the head of $\sigma^g$ is satisfied in $R'$, which contains only constants from $adom(D, \chi)$). Since all built-in literals in $\sigma$ are equalities and inequalities, there are infinitely many constants $c$ such that for the ground instance $\sigma_c^g$ of $\sigma$ which differs from $\sigma^g$ only by substitution of $y_{j,h}$ with $c$, all built-in literals evaluate to false. Since $\sigma_c^g$ and $\sigma^g$ have the same body and $R \models \sigma_c^g$, $R$ must contain a fact in which $c$ occurs. This means that $R$ is infinite, which is a contradiction. □

## B. PROOFS FOR SECTION 4

**Proposition 4.1 (Separation)** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$. Then, (1) $facts(\Sigma^a) = C^*$, (2) $facts(\Sigma^s) \cap C^* = \emptyset$, (3) $\Sigma^a \cap \Sigma^s = \emptyset$, and (4) $\Sigma^a \cup \Sigma^s = ground(\Sigma)$.*

PROOF. By definition, $\sigma \in \Sigma^a$ contains at least one fact $p(\vec{t})$ from $C^*$; any other fact in $\sigma$ is constraint-bounded in $\chi$ with $p(\vec{t})$, and hence it also must be in $C^*$. This proves $facts(\Sigma^a) \subseteq C^*$. Consider now any fact $p(\vec{t}) \in C^*$. The minimality of $C^*$ implies that there exist facts $f_1, \ldots, f_n$ in $C^*$ such that $f_1 \in C$, $f_n = p(\vec{t})$, and $f_{i+1}$ is constraint-bounded to $f_i$, for each $i \in \{1, \ldots, n-1\}$; i.e., $f_i, f_{i+1} \in facts(\sigma_i)$ for some $\sigma_i \in ground(\Sigma)$. Each $\sigma_i$ then belongs to $\Sigma^a$, and thus $p(\vec{t}) \in facts(\Sigma^a)$. This proves $C^* \subseteq facts(\Sigma^a)$, and therefore (1) holds. As for (2), assume by contradiction that some $\sigma \in \Sigma^s$ with $facts(\sigma) \cap C^* \neq \emptyset$ exists. Then, from Definition 4.4 it follows that $facts(\sigma) \subseteq C^*$, which contradicts $\sigma \in \Sigma^s$. Item (3) is straightforward from (1) and (2). Finally, in order to prove (4), we suppose that there exists $\sigma \in ground(\Sigma)$ such that $\sigma \notin \Sigma^s$ and $\sigma \notin \Sigma^a$, but this means that $facts(\sigma) \cap C^* = \emptyset$ and $facts(\sigma) \subseteq C^*$, which is an obvious contradiction. □

**Proposition 4.2 (Safe database)** *Let $D$ be any database for $\chi = \langle \Psi, \Sigma \rangle$. Then, for each repair $R \in rep(D)$ it holds that $R \setminus C^* = D \setminus C^*$.*

PROOF. Towards a contradiction, suppose that there exists a repair $R \in rep(D)$ such that $R \setminus C^* \neq D \setminus C^*$. Let $R' = (R \cap C^*) \cup (D \setminus C^*)$ (notice that $R' \neq R$ only if $R \setminus C^* \neq D \setminus C^*$). Consider any $\sigma \in ground(\Sigma)$. By Proposition 4.1, either (i) $\sigma \in \Sigma^a$ or (ii) $\sigma \in \Sigma^s$. In case (i), $R' \models \sigma$: by Proposition 4.1 (a), $facts(\sigma) \subseteq C^*$, and therefore $R' \models \sigma$ iff $R' \cap C^* \models \sigma$, which is true, since $R' \cap C^* = R \cap C^*$ and $R \models \sigma$ (because $R \in rep(D)$). In case (ii), again $R' \models \sigma$: by Proposition 4.1 (b), $facts(\sigma) \cap C^* = \emptyset$ and therefore $R' \models \sigma$ iff $R' \setminus C^* \models \sigma$, which is true, since $R' \setminus C^* = D \setminus C^*$ and $D \models \sigma$. It follows that $R' \models \Sigma$. Furthermore, it is easy to show that $\triangle(R', D) \subset \triangle(R, D)$. Indeed, $R' \setminus D = ((R \cap C^*) \cup (D \setminus C^*)) \setminus D = (R \cap C^*) \setminus D \subseteq R \setminus D$, and also $D \setminus R' = D \setminus ((R \cap C^*) \cup (D \setminus C^*)) = (D \setminus (R \cap C^*)) \cap (D \setminus (D \setminus C^*)) = (D \setminus (R \cap C^*)) \cap (D \cap C^*) \subseteq D \setminus R$. From (SIP), it follows that $R' <_D R$. This contradicts $R \in rep(D)$. □

**Lemma 4.3** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, and let $A = D \cap C^*$ and $\chi^a = \langle \Psi, \Sigma^a \rangle$. Then, for each $S \subseteq D \setminus C^*$, the following holds:*

*(1) for each $R \in rep_\chi(A \cup S)$, it holds that $R \cap C^* \in rep_{\chi^a}(A)$;*

(2) *for each $R \in rep_{\chi^a}(A)$ there exists a set of facts $S' \subseteq \mathcal{F}(\chi)$ such that $S' \cap C^* = \emptyset$, and $(R \cup S') \in rep_\chi(A \cup S)$.*

PROOF. (1) Let $R \in rep_\chi(A \cup S)$, and let $R' = R \cap C^*$. Since $R \models \Sigma$, then $R \models ground(\Sigma)$, therefore, from Proposition 4.1 it follows that $R' \models \Sigma^a$, while $R \setminus R' = R \setminus C^* \models \Sigma^s$. Assume $R' \notin rep_{\chi^a}(A)$. Since $R' \models \Sigma^a$, there must exist some $R'' \in rep_{\chi^a}(A)$ such that $R'' <_A R'$. Since $R'' \models \Sigma^a$ and $R \setminus R' \models \Sigma^s$, we have that $R'' \cup (R \setminus R') \models \Sigma$. Since the conflict closure of $D$ w.r.t. $\chi$ is the same as w.r.t. $\chi^a$ and since $A$ is contained in its conflict closure w.r.t. $\chi$, by Proposition 4.2, we have $R'' \setminus C^* = A \setminus C^* = \emptyset$, and therefore $R'' \subseteq C^*$. As a consequence, $(R'' \cup R' \cup A) \cap ((R \setminus R') \cup S) = \emptyset$ (notice that $(R'' \cup R' \cup A) \subseteq C^*$ whereas $(R \setminus R') \cap S) = \emptyset$). Then, by (DPE), it follows that $R'' \cup (R \setminus R') <_{A \cup S} R' \cup (R \setminus R') = R$. This contradicts that $R \in rep_\chi(A \cup S)$.

(2) We choose as $S'$ an arbitrary repair for $S$ w.r.t. $\chi^s = \langle \Psi, \Sigma^s \rangle$, i.e., $S' \in rep_{\chi^s}(S)$ (notice that $S$ may violate $\Sigma^s$, and therefore in general $S' \neq S$). We first show that $S' \cap C^* = 0$. Let us write $S' = S'_a \cup S'_s$, where $S'_a = S' \cap C^*_{\chi^s}(S)$ and $S'_s = S' \setminus C^*_{\chi^s}(S)$. By Proposition 4.2, we have that $S'_s = S \setminus C^*_{\chi^s}(S)$, and therefore, since $S \subseteq D \setminus C^*$, $S'_s \cap C^* = \emptyset$. Also, it is easy to see that $C^* \cap C^*_{\chi^s}(S) = \emptyset$, and therefore $S'_a \cap C^* = \emptyset$. We thus conclude that $S' \cap C^* = \emptyset$. Now, we concentrate on proving that $R \cup S' \in rep_\chi(A \cup S)$.

Since $R \in rep_{\chi^a}(A)$, from Proposition 4.2, and from the fact that $C^*_{\chi^a}(A) = C^*$ (in computing $C^*_{\chi^a}(A)$ and $C^*$ we start from the same conflict set and we close such set w.r.t. the same set of constraints) and $A \subseteq C^*$, it follows that $R \setminus C^* = A \setminus C^* = \emptyset$, and therefore $R \subseteq C^*$. Then, it is easy to see that from $R \models \Sigma^a$, $S' \models \Sigma^s$, and $\Sigma = \Sigma^a \cup \Sigma^s$, it follows that $R \cup S' \models \Sigma$. Assume now by contradiction that $R \cup S' \notin rep_\chi(A \cup S)$, then there must exist some $R''$ consistent with $\Sigma$ such that $R'' <_{A \cup S} R \cup S'$. We can write $R'' = R''_a \cup R''_s$, where $R''_a = R'' \cap C^*$ and $R''_s = R'' \setminus C^*$. Let us now apply property (DIS) with $R = C^*$, and obtain that either $R''_a <_A R \cap C^*$ or $R''_s <_S S'$. From Proposition 4.1, it follows that $R''_a \models \Sigma^a$ and $R''_s \models \Sigma^s$, but this contradicts the assumptions that $R \in rep_{\chi^a}(A)$ and $S' \in rep_{\chi^s}(S)$. This proves that $R \cup S' \in rep_\chi(A \cup S)$. $\square$

**Corollary 4.5** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, and let $\chi' = \langle \Psi, \Sigma' \rangle$ be such that $\Sigma^a_{\chi'}(D) = \Sigma^a_\chi(D)$. Then $rep_\chi(D) = rep_{\chi'}(D)$.*

PROOF. We prove that $rep_\chi(D) \subseteq rep_{\chi'}(D)$. The converse can be proved analogously. Assume by contradiction that there exists $R \in rep_\chi(D)$ such that $R \notin rep_{\chi'}(D)$. This means that either (a) $R \not\models \Sigma'$ or (b) there exists $R'$ such that $R' \models \Sigma'$ and $R' <_D R$. Consider first case (a). Since $R \not\models \Sigma'$ iff $R \not\models \Sigma^a_{\chi'}(D) \cup \Sigma^s_{\chi'}(D)$ and since $R \models \Sigma^a_\chi(D)$ (which is equal to $\Sigma^a_{\chi'}(D)$), there exists some $\sigma \in \Sigma^s_{\chi'}(D)$ such that $R \not\models \sigma$. By Theorem 4.4, $R = (R'' \cap C^*_\chi(D)) \cup (D \setminus C^*_\chi(D))$, where $R'' \in rep_\chi(D \cap C^*_\chi(D))$. Since $facts(\Sigma^a_\chi(D)) = facts(\Sigma^a_{\chi'}(D))$, by Proposition 4.1 $C^*_\chi(D) = C^*_{\chi'}(D) = C^*$. Furthermore, since $\sigma \in \Sigma^s_{\chi'}(D)$, we have $facts(\sigma) \cap C^* = \emptyset$. Therefore, $R$ can violate $\sigma$ only if $D \setminus C^*$ violates $\sigma$, but this is a contradiction. In case (b), we can show similarly as in case (a) that $R' \models \Sigma'$ implies $R' \models \Sigma$. Then, $R' <_D R$ contradicts $R \in rep_\chi(D)$. $\square$

**Corollary 4.8** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$ where $\Sigma \subseteq \mathbf{C}_i$, for $i \in \{1, 2\}$. Then, $C^*$ is a repair envelope for $D$. In fact, there exists a bijection $\mu : rep(D) \to rep(D \cap C^*)$, such that for every $R \in rep(D)$, $R = \mu(R) \cup (D \setminus C^*)$.*

PROOF. The result for $\Sigma \subseteq \mathbf{C}_1$ (resp. $\Sigma \subseteq \mathbf{C}_2$) follows from Theorem 4.4 by applying Proposition 4.6 (resp. Proposition 4.7). Note that for each $R \in rep(D \cap C^*)$, when $\Sigma \subseteq \mathbf{C}_1$, $R \cap C^* = R$, whereas when $\Sigma \subseteq \mathbf{C}_2$, $(R \cap C^*) \cup (D \setminus C^*) = R \cup (D \setminus C^*)$. $\square$

**Proposition 4.9** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, $\Sigma \subseteq \mathbf{C}_0$, and let $A = D \cap C^*$. Then,*

*(1)* $C \subseteq D$;

*(2)* *for each $R \in rep(A)$, (i) $R \subseteq A$, (ii) $\triangle(R, A) \subseteq C$, (iii) $A \backslash C \subseteq R$, and (iv) $R \cap C \in rep(C)$;*

*(3)* *for each $R \in rep(C)$, $R \cup (A \setminus C) \in rep(A)$.*

PROOF. 1) By definition, $C$ is the set of facts occurring in any constraint $\sigma \in ground(\Sigma)$ violated in $D$. Since each $\sigma$ is of the form $\bigwedge_{i=1}^{l} A_i(\vec{c}_i) \supset \bigvee_{k=1}^{n} \phi_k(\vec{d}_k)$, it can be violated only if all the body facts are in $D$. That is, $C \subseteq D$.

2) Let $R \in rep(A)$. We first show (i). Assume towards a contradiction that $R \not\subseteq A$ and consider $R' = R \cap A$. From the fact that $R \models \Sigma$, $R' \subseteq R$, and that each $\sigma \in \Sigma$ is of the form $\bigwedge_{i=1}^{l} A_i(\vec{x}_i) \supset \bigvee_{k=1}^{n} \phi_k(\vec{z}_k)$, it follows that $R' \models \Sigma$, therefore (SIP) would raise a contradiction. We now show (ii). Assume towards a contradiction that $\triangle(R, A) \not\subseteq C$. Since $R \subseteq A$, this implies that there exists some $p(\vec{t}) \in A \setminus R$ such that $p(\vec{t}) \notin C$. By minimality of $R$, $p(\vec{t})$ occurs in the body of at least one constraint in $ground(\Sigma)$ of the form $\bigwedge_{i=1}^{l} a_i \supset \bigvee_{k=1}^{n} \phi_k$. No such constraint, however, is violated in $A$ Hence, $R \cup \{p(\vec{t})\} \models \Sigma$, which by (SIP) implies that $R \notin rep(A)$; this is a contradiction. Therefore, (ii) holds. From (i) and (ii), follows that $A \setminus C \subseteq R$; this proves (iii). To show (iv), suppose towards a contradiction that $R \cap C \notin rep(C)$. Then, it is easy to see that $R \cap C \models \Sigma$, therefore some $R' \in rep(C)$ must exist such that $R' <_C R$. Since all constraints have only built-ins in their heads, $R' \subseteq C$. But then $(R \setminus C) \cup R' <_A R$ contradicts $R \in rep(A)$.

3) Let $R' \in rep(C)$. Item 2.(ii) for $D=C$ (where $A=C$) implies $R' \subseteq C$. Since $R' \models \Sigma$, we must have $R = (A \setminus C) \cup R' \models \Sigma$; otherwise, suppose $R \not\models \sigma$ for some $\sigma \in ground(\Sigma)$. Then $\sigma$ must contain a fact $p(\vec{t})$ from $A \setminus C$. Since $\sigma$ is from $\mathbf{C}_0$ and Item 1 implies $R \subseteq D$, also $D \not\models \sigma$. But this means $p(\vec{t}) \in C$, a contradiction. From Item 2 and (DIS) (split by $C$), we can conclude that no $R'' \in rep(A)$ exists such that $R'' <_A R$. Hence, $R \in rep(A)$. $\square$

## C. PROOFS FOR SECTION 5

**Proposition 5.3** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, and let $E$ be a repair envelope for $D$. If either (1) $\Sigma \subseteq \mathbf{C}_1$ and $E = C^*$ or (2) $\Sigma \subseteq \mathbf{C}_0$, then every constraint-bounded partitioning $E_1, \ldots, E_m$ of $E$ is repair-compliant.*

PROOF. Since $E_1, \ldots, E_m$ is constraint-bounded, what remains to show is that for all $R \in rep(D \cap E)$ and $R_i \in rep(D \cap E_i)$, $1 \leq i \leq m$, $R \setminus E = R_i \setminus E_i$. We show that $R \setminus E = R_i \setminus E_i = \emptyset$.

Case (1): $E = C^*$. By Proposition 4.6, $R \in rep(D \cap E)$ implies $R \subseteq E$, hence $R \setminus E = \emptyset$. Towards a contradiction assume that there exists $R_i \in rep(D \cap E_i)$, such that $R_i \setminus E_i \neq \emptyset$ for some $1 \leq i \leq m$. Consider $R'_i = R_i \cap E_i$. Clearly (since $R'_i \subset R_i$), we have by (SIP) that $R'_i <_{D \cap E_i} R_i$. We show $R'_i \models \Sigma$. Assume $R'_i \not\models \Sigma$, i.e., there exists a ground constraint $\sigma \in ground(\Sigma)$ such that $R'_i \not\models \sigma$. Thus, $R'_i \models body(\sigma)$, which implies $R_i \models body(\sigma)$, and $R'_i \not\models head(\sigma)$. However, $R_i \models head(\sigma)$ must hold since $R_i \models \Sigma$ by hypothesis. This means that there exists a head atom $B(\vec{y})$ of $\sigma$ which is true in $R_i$. Since $R'_i \not\models head(\sigma)$, none of the built-in predicates of $\sigma$ is true and therefore $B(\vec{y})$ is a fact. Furthermore, $body(\sigma) \neq \emptyset$ since $\sigma \in \mathbf{C}_1$, and thus, $B(\vec{y}) \notin E$, since $E_1, \ldots, E_m$

is constraint-bounded. Consequently, $facts(\sigma) \cap E \neq \emptyset$ holds, as well as $facts(\sigma) \not\subseteq E$. Therefore, $\sigma \notin \Sigma^s$ and $\sigma \notin \Sigma^a$, which is a contradiction to Proposition 4.1. This proves $R'_i \models \Sigma$. Together with $R'_i <_{D \cap E_i} R_i$, we arrive at a contradiction to $R_i \in rep(D \cap E_i)$. Hence, $R_i \setminus E_i = \emptyset$ must hold.

Case (2): $\Sigma \subseteq \mathbf{C}_0$. Towards a contradiction assume that there exists $R \in rep(D \cap E)$ such that $R \setminus E \neq \emptyset$. Consider $R' = R \cap E$. Again (since $R' \subset R$), we have that $R' <_{D \cap E} R$. Furthermore, $R \models \Sigma$ implies $R' \models \Sigma$. Indeed, consider any $\sigma \in ground(\Sigma)$ such that $R' \models body(\sigma)$. Then $R \models body(\sigma)$, and, by hypothesis, $R \models head(\sigma)$. Since $\Sigma \subseteq \mathbf{C}_0$, one of the built-in atoms in the head of $\sigma$ is true. Thus, $R' \models head(\sigma)$. This shows $R' \models \Sigma$, which raises a contradiction to $R \in rep(D \cap E)$. This proves that $R \setminus E = \emptyset$. Along the same argumentation line, we can show that every $R_i \in rep(D \cap E_i)$ fulfills $R_i \setminus E_i = \emptyset$. Hence, $R \setminus E = R_i \setminus E_i = \emptyset$ holds for all $R \in rep(D \cap E)$ and $R_i \in rep(D \cap E_i)$, $1 \leq i \leq m$. $\quad\square$

**Proposition 5.4** *Let $E_1, \ldots, E_m$ be a factorization of a repair envelope $E$ for a database $D$, let $Q$ be a non-recursive Datalog query, and let $S_Q = (D \setminus E) \cup R_1 \cup \cdots \cup R_\ell$ be a query-safe part of $D$ w.r.t. $Q$. Then,*

$$ans_c(Q, D) = \bigcap_{R_{\ell+1} \in rep(D \cap E_{\ell+1})} \cdots \bigcap_{R_m \in rep(D \cap E_m)} Q[S_Q \cup R_{\ell+1} \cdots \cup R_m].$$

PROOF. Let $E_1, \ldots, E_m$ be a factorization of $E$ for $D$. From Equation (7), it holds that:

$$rep(D) = \{(D \setminus E) \cup R_1 \cup \cdots \cup R_m \mid R_i \in rep(D \cap E_i), 1 \leq i \leq m\}.$$

Moreover, from Proposition 5.1, $ans_c(Q, D) = \bigcap_{R \in rep(D \cap E)} Q[R \cup S]$, where $S = D \setminus E$. Thus, $ans_c(Q, D)$ can be equivalently written as:

$$\bigcap_{R_1 \in rep(D \cap E_1)} \cdots \bigcap_{R_m \in rep(D \cap E_m)} Q[(D \setminus E) \cup R_1 \cup \cdots \cup R_m].$$

Consider now repairs $R'_i \in rep(D \cap E_i)$ and $R''_h \in rep(D \cap E_h)$, with $1 \leq h \leq \ell$ and $1 \leq i \leq m$. We claim that $Q[X \cup R'_h] \subseteq Q[X \cup R''_h]$, where $X = (D \setminus E) \cup \bigcup\{R'_j \mid 1 \leq j \leq m, j \neq h\}$, holds if $R'_h = R_h$ (i.e., a repair of $D \cap E_h$ such that $Q_h[(D \setminus E) \cup R_h] \in amin(Q_h)$) and $R''_h$ is an arbitrary repair of $D \cap E_h$.

To prove the claim, consider an arbitrary ground instance $\rho'_j$ of some rule $\rho_j$, $1 \leq j \leq v(Q)$ of $Q$, such that $body(\rho'_j)$ is satisfied by $X \cup R'_h$ and therefore $head(\rho'_j)$ is included in $Q[X \cup R'_h]$. Let $q_h(\vec{t'}) \leftarrow \beta'_{h,j}$ be the corresponding ground instance of the rewritten rule $q_h(\vec{t}) \leftarrow \beta_{h,j}$ (obtained with the same variable bindings). Then, $(D \setminus E) \cup R_h$ satisfies $\beta'_{h,j}$ (note that $\beta'_{h,j} \subseteq body(\rho'_j)$). By singularity of $E_h$, we conclude that there is a (possibly different) ground instance $q_h(\vec{t''}) \leftarrow \beta''_{h,j}$ of $q_h(\vec{t}) \leftarrow \beta_{h,j}$, such that $\vec{t'} = \vec{t''}$ and $(D \setminus E) \cup R''_h$ satisfies $\beta''_{h,j}$.

We now consider the ground instance $\rho''_j$ of $\rho_j$ in which all variables occurring in $\beta^{in}_{h,j}$ get the same value as in $\beta''_{h,j}$, and all other variables get the same value as in $\rho'_j$. Obviously, every atom in $\beta^{in}_{h,j}$ is instantiated in $\rho''_j$ as in $\beta^{in}_{h,j}''$, and is thus satisfied by $X \cup R''_h$. Now consider any atom $p(\vec{x})$ from $body(\rho_j) \setminus \beta^{in}_{h,j}$. Each variable $x$ in $\vec{x}$ that occurs in $\beta^{in}_{h,j}$ also occurs, by definition, in $\vec{t}$; hence, $x$ has in $\rho''_j$ the same value as in $\rho'_j$. The latter holds, by definition, also for every variable $x$ in $\vec{x}$ that does not occur in $\beta^{in}_{h,j}$. Hence, $p(\vec{x})$ is

instantiated in $\rho_j''$ to the same atom as in $\rho_j'$, denoted by $p(\vec{x}')$. Since $X \cup R_h'$ satisfies $p(\vec{x}')$ and $p(\vec{x}) \notin \beta_{h,j}^{in}$ implies $p(\vec{x}') \notin E_h$, it follows that $X \cup R_h''$ satisfies $p(\vec{x}')$. In conclusion, $X \cup R_h''$ satisfies $body(\rho_j'')$; this means that $head(\rho_j')$ is in $Q[X \cup R_h'']$.

Now by construction, $head(\rho_j'') = head(\rho_j')$, as each variable in $head(\rho_j)$ occurring in $\beta_{h,j}^{in}$ also occurs in $\vec{t}$. This shows that $head(\rho_j')$ is also in $Q[X \cup R_h'']$, and proves the claim.

From the claim, we conclude that:

$$\bigcap_{R_h'' \in rep(D \cap E_h)} Q[(D \setminus E) \cup R_1' \cup \cdots R_h'' \cdots \cup R_m'] = Q[(D \setminus E) \cup R_1' \cup \cdots R_h \cdots \cup R_m'].$$

Hence, $ans_c(Q, D)$ can be also computed as:

$$\bigcap_{R_{\ell+1} \in rep(D \cap E_{\ell+1})} \cdots \bigcap_{R_m \in rep(D \cap E_m)} Q[(D \setminus E) \cup R_1 \cup \cdots \cup R_m].$$

The result follows by letting $S_Q = (D \setminus E) \cup R_1 \cup \cdots \cup R_\ell$. $\square$

**Proposition 5.6** *For $\mathcal{Q}_{1k\exists}$, consistent query answering is polynomial in data complexity.*

PROOF. For computing consistent answers to queries $Q$ in $\mathcal{Q}_{1k\exists}$, we take (partial) ground instances $Q_{\vec{c}}$ of $Q$ and proceed as in Example 5.10. More precisely, let $Q = \langle q, \{\rho\} \rangle$ be a query in $\mathcal{Q}_{1k\exists}$, and let $E_1, \ldots, E_m$ be the natural factorization of the repair envelope $E = C$ for $\mathbf{C}_0$ (cf. Theorem 4.10; equally well, we may use $E = C^*$ here), i.e., each component $E_i$ contains the facts over one relation symbol $p_i$ that clash on the key of $p_i$ (each $p_i$ has some—either explicitly or implicitly given—key).

Let, for any substitution $\theta$ of constants for the head variables in $Q$, be $Q_{\vec{c}} = \langle q, \{\rho\theta\} \rangle$ the (partial) ground instance of the query $Q$, where $head(\rho\theta) = q(\vec{c})$. Suppose that $body(\rho\theta) = a_1, \ldots, a_k$ where $a_i = p_i(t_{i,1}, \ldots, t_{i,n_i})$ (each $t_{i,j}$ being a constant or variable), and that $1, 2, ..., \ell_i (\leq n_i)$ are the positions of the key attributes of $p_i$. Furthermore, suppose that $a_k$ has a non-head variable at a key position, i.e., some $t_{k,j}, 1 \leq j \leq \ell_i$, is a variable that does not occur in $head(\rho\theta)$. Note that no other $a_i$ has this property since $Q$ is in $\mathcal{Q}_{1k\exists}$ and therefore satisfies the following condition: $(*)$ except in at most one atom, key positions are always head variables or constants. Then, the following holds:

(1) at most $k-1$ components $E_j, 1 \leq j \leq m$, have a match with some atom $a_i, 1 \leq i < k$, i.e., for some substitution $\sigma$ of constants for the variables in $a_i$ we have $a_i\sigma \in E_j$;

(2) the components $E_j$ having a match with $a_k$ but with no $a_i, 1 \leq i < k$, are jointly decomposable; and

(3) all other components (which have no match with $a_1, \ldots, a_k$) are singular.

Indeed, as for 1), by condition $(*)$ of $\mathcal{Q}_{1k\exists}$, every atom $a_i$, where $1 \leq j < k$, has constants at all key positions of $p_i$. Suppose that $E_j$ and $E_{j'}$ both match with $a_i$ via substitutions $\sigma$ and $\sigma'$, respectively. Then $a_i\sigma$ and $a_i\sigma'$ coincide on the key attributes of $p_i$, and hence $a_i\sigma, a_i\sigma'$ both belong to $E_j$ and $E_{j'}$; this implies that $E_j = E_{j'}$, from which 1) follows.

As for 2), $n(Q_{\vec{c}}) = 1$ and for each such component $E_j$, we have $\beta_{j,1}^{in} = \{a_k\}$ and furthermore $R_j \setminus E_j = \emptyset$ for each repair $R_j$ of $D \cap E_j$. Hence, the conditions 1)–3) of decomposability w.r.t. $Q_{\vec{c}}$ are satisfied.

Finally, for each $E_j$ in 3), $\beta_{j,1}^{in} = \emptyset$ holds; since $R_j \subseteq E_j$ for each repair $R_j$ of $D \cap E_j$, the component $E_j$ is trivially singular (irrelevant) w.r.t. $Q_{\vec{c}}$. This establishes 1)–3) above.

Clearly, each component $E_j$, $1 \le j \le m$, can be easily classified into the correct group 1), 2), or 3) in polynomial time, and $rep(D \cap E_j) = \{\{a\} \mid a \in D \cap E_j\}$ is easily computed. Therefore, we can compute $ans_c(Q_{\vec{c}}, D)$ via Equation (10) in Theorem 5.5 in polynomial time under data complexity (the components $E_j$ in 1) have polynomially many combinations of their repairs, and the innermost intersection in (10) has a linear number of terms). Hence, $Q_{\vec{c}}$ has polynomial data complexity.

The argument is easily adapted to the case where no atom has a non-head variable at a key position (however, we may assume without loss of generality that such an atom exists). Polynomial data complexity of the query $Q$ then follows by Equation (12). $\quad\square$

## D. GROUPED REPAIR COMPUTATION

This section gives the technical details on grouped repair computation by means of evaluating an SQL query over a marked database, as discussed less formally in Section 6.2.

### D.1 Query Reformulation

We first show how a non-recursive Datalog$^\neg$ query $Q$ can be reformulated into an SQL query whose evaluation over the marked database returns the consistent answers to $Q$, i.e, those that are true in any repair of $\{R_1, ..., R_n\}$.

In query reformulation, we use the following functions ANDBIT, INVBIT, and SUMBIT, which can be build as *user-defined functions* (ANDBIT and INVBIT) and *aggregate operators* (SUMBIT) in many relational DBMSs, such as PostgreSQL:

—ANDBIT is a binary function that takes as its input two bit strings $'a_1 \ldots a'_n$ and $'b_1 \ldots b'_n$ and returns $'c_1 \ldots c'_n$, where $c_i = a_i \wedge b_i$ is the Boolean "and," $i = 1, \ldots, n$;

—INVBIT is a unary function that takes as its input a bit string $'a_1 \ldots a'_n$ and returns $'c_1 \ldots c'_n$, where $c_i = \neg a_i$ is the Boolean complement, $i = 1, \ldots, n$;

—SUMBIT is an aggregate function such that given $m$ strings of form $'b_{i,1} \ldots b'_{i,n}$, $i = 1, \ldots, m$, it returns $'c_1 \ldots c'_n$, where $c_j = b_{1,j} \vee \ldots \vee b_{m,j}$ is the Boolean "or," $j = 1, \ldots, n$.

Let $Q = \langle q, \mathcal{P} \rangle$ be a non-recursive Datalog$^\neg$ query of arity $n$, where $\mathcal{P}$ consists of normalized rules $r : h(\vec{x}') \leftarrow B(\vec{y}'), e(\vec{z})$, where $e(\vec{z})$ are all the equality atoms introduced in normalization, as described in Section 6.2. Let $a_i$, $1 \le i \le n$, be pairwise distinct identifiers for the attributes of a predicate of arity $n$. Then, each $r$ is translated into the following SQL statement $SQL_r$ (notice that, in the statements below, each relation symbol $p_i$ occurring in $r$ is transformed into the corresponding marked symbol $p_{i_m}$):

```
SELECT p_{i'_m}.a_{j'} AS a_j          (for each atom y_{0,j} = y_{i',j'} in e(z))
       c AS a_j,                        (for each atom y_{0,j} = c in e(z))
       (p_{1_m}.mark ANDBIT ... ANDBIT p_{l_m}.mark ANDBIT INVBIT(n_p_{l+1_m}.mark) ANDBIT
       ... ANDBIT INVBIT(n_p_{l+k_m}.mark)) AS mark
FROM p_{1_m}, ..., p_{l_m}, SQL_{r,l+1}, ..., SQL_{r,l+k}
WHERE p_{i_m}.a_j = p_{i'_m}.a_{j'},    (for each atom y_{i,j} = y_{i',j'} in e(z), 0 < i ≤ l)
      p_{i_m}.a_j = c,                  (for each atom y_{i,j} = c in e(z), 0 < i ≤ l)
      n_p_{i_m}.a_j = p_{i'_m}.a_{j'},  (for each atom y_{i,j} = y_{i',j'} in e(z), l < i)
      n_p_{i_m}.a_j = c                 (for each atom y_{i,j} = c in e(z), l < i).
```

where each $SQL_{r,h}$, $l < h \le l + k$, is a subquery of form:

```
( SELECT * FROM p_{h_m}
  UNION
  SELECT p_{i'_m}.a_{j'} AS a_j,                    (for each atom y_{h,j} = y_{i',j'}  in  e(z⃗))
         c AS a_j,                                  (for each atom y_{h,j} = c  in  e(z⃗))
         '0...0' AS mark
  FROM p_{1_m}, ..., p_{l_m}
  WHERE p_{i_m}.a_j = p_{i'_m}.a_{j'},              (for each atom y_{i,j} = y_{i',j'}  in  e(z⃗), 0 < i ≤ l)
        p_{i_m}.a_j = c,                            (for each atom y_{i,j} = c  in  e(z⃗), 0 < i ≤ l)
        ROW(a_1, ..., a_{k_h}) NOT IN (SELECT a_1, ..., a_{k_h} FROM p_{h_m})
) AS n_p_{h_m}.
```

Roughly speaking, in the statement $SQL_r$, the ANDBIT operator allows us to obtain the mark $'b_1, \ldots, b'_n$ of each tuple $\vec{t}$ computed for the relation predicate $h$, according to rule $r$. More precisely, for $i \in \{1, \ldots, n\}$, $b_i = 1$ if $\vec{t}$ is in the repair $R_i \in rep(A)$, $b_i = 0$ otherwise. Moreover, for each negative literal $not\ p_{h_m}(\vec{y}_h)$, the marks must be inverted, where missing tuples (which do not belong to any repair, and thus would be marked $'0 \ldots 0'$) must be taken into account.

To this aim, $SQL_{r,h}$ singles out the tuples returned by the positive body of the rule $r$, projects them on the attributes that are in join with the attributes in $p_{h_m}$, and returns, with mark $'0 \ldots 0'$, those which do not occur in $p_{h_m}$ (taking then the union with the tuples in $p_{h_m}$ itself). The operator INVBIT guarantees that, for each such tuple, the mark returned by $SQL_r$ is the one computed in the positive part of the query (in these cases indeed the negative literal is satisfied in every repair). Note that safety of the rule $r$ ensures that the two queries in $SQL_{r,h}$ have the same arity.

All rules, $r_1, \ldots, r_\ell$, defining the same predicate $h$ of arity $n$, are collected into a view by the SQL statement $SQL_h$:

```
CREATE VIEW  h_m(a_1, ..., a_n, mark)  AS
       SELECT a_1, ..., a_n, SUMBIT(mark)
       FROM (SQL_{r_1} UNION ... UNION SQL_{r_ℓ})
       GROUP BY a_1, ... a_n.
```

Finally, the consistent answers to the query $Q = \langle q, \mathcal{P} \rangle$ are obtained through the statement $SQL_Q$:

```
SELECT a_1, ... a_n FROM q_m WHERE mark =' 1...1',
```

where $q_m$ is the view predicate defined by the statement $SQL_q$.

**Example D.1** Let us consider a query $Q = \langle q, \mathcal{P} \rangle$ asking for players that are not team leaders. Here $\mathcal{P}$ contains two rules, of which one defines an auxiliary (thus intensional) predicate $leader$:

$$r_1 : q(x) \leftarrow player(x, y, z),\ not\ leader(x);$$
$$r_2 : leader(x) \leftarrow team(v, w, x).$$

The use of negation is reflected in $SQL_{r'_1}$ (let $r'_1$, $r'_2$ be the normalized versions of $r_1$, $r_2$):

```
SELECT  player_m.Pcode AS a_1,
        (player_m.mark ANDBIT INVBIT(n_leader_m.mark)) AS mark,
FROM  player_m,
      (SELECT player_m.Pcode AS a_1, '00' AS mark
       FROM player_m WHERE ROW(player_m.Pcode) NOT IN (SELECT a_1 FROM leader_m)
       UNION SELECT * FROM leader_m) AS n_leader_m
WHERE  n_leader_m.a_1 = player_m.Pcode;
```

The use of an auxiliary predicate causes the creation of two views: one for each intensional predicate. The respective SQL statements $SQL_q$ and $SQL_{leader}$, resemble the statement $SQL_q$ in Example 6.1, however, each of them just depends on a single SQL query ($SQL_{r'_1}$ and $SQL_{r'_2}$, respectively). Moreover, one can show that $SQL_{r'_2}$ and $SQL_Q$ are equal to the corresponding queries of Example 6.1.

Hence, as easily retraced, the answer to the query $Q$ consists of the tuple $(9)$, as expected. $\square$

We next show the correctness of the above transformation $SQL_Q$. Recall that any predicate names in a Datalog$^{\vee,\neg}$ program $\mathcal{P}$ are called *extensional* (*EDB predicates*), if they occur only in the bodies of the rules in $\mathcal{P}$, and *intensional* (*IDB predicates*) otherwise.

**Proposition 6.1** *Let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, let $Q$ be a non-recursive Datalog$^{\neg}$ query over it, and let $R_1, ..., R_n$ be databases such that $R_i = R'_i \cap E$, where $E$ is a weak repair envelope for $D$ and $R'_i$ is a repair for $A = D \cap E$. Then, $SQL_Q$ computes on $D_m$ the set of tuples $\bigcap_{i=1}^{n}\{\vec{t} \mid \vec{t} \in Q[R_i \cup S]\}$, for $S = D \setminus E$.*

PROOF. We first show that normalization of a Datalog$^{\neg}$ query does not change query semantics. In particular, let $r$ be a safe Datalog$^{\neg}$ rule and let $r'$ denote the normalized rule as introduced in Section 6.2. We show that there is a one-to-one correspondence between relevant ground instances, that is, between ground instances of $r$ and ground instances of $r'$ that satisfy the equality conditions.

($\subseteq$) Let $r(\vec{t}) = p_0(\vec{t_0}) \leftarrow p_1(\vec{t_1}), \ldots, p_l(\vec{t_l}), \text{not } p_{l+1}(\vec{t_{l+1}}), \ldots, \text{not } p_{l+k}(\vec{t_{l+k}})$ be a ground instance of $r$. Consider $r'(\vec{t}) = p_0(\vec{t_0}) \leftarrow p_1(\vec{t_1}), \ldots, p_l(\vec{t_l}), \text{not } p_{l+1}(\vec{t_{l+1}}), \ldots,$ $\text{not } p_{l+k}(\vec{t_{l+k}}), e(\vec{t_z})$, where $\vec{t_z}$ is obtained by substituting $t_{i,j}$ for every variable $y_{i,j}$ in $\vec{z}$, such that $0 \le i \le l + k$ and $1 \le j \le k_i$. Then, $r'(\vec{t})$ is a ground instance of $r'$, since every occurrence of a variable $y_{i,j}$ in $e(\vec{z})$ is substituted uniformly. Moreover, $e(\vec{t_z})$ is true (otherwise we arrive at a contradiction to our hypothesis since then by construction of $e(\vec{z})$ either $x_{i,j} = c$ and $c \ne t_{i,j}$ or $x_{i,j} = x_{i',j'}$ and $t_{i,j} \ne t_{i',j'}$, for some $0 \le i, i' \le l + k$, $1 \le j \le k_i$, and $1 \le j' \le k_{i'}$).

($\supseteq$) Let $r'(\vec{t}) = p_0(\vec{t_0}) \leftarrow p_1(\vec{t_1}), \ldots, p_l(\vec{t_l}), \text{not } p_{l+1}(\vec{t_{l+1}}), \ldots, \text{not } p_{l+k}(\vec{t_{l+k}}), e(\vec{t_z})$ be a ground instance of $r'$, such that $e(\vec{t_z})$ is true. Then, $r(\vec{t}) = p_0(\vec{t_0}) \leftarrow p_1(\vec{t_1}), \ldots, p_l(\vec{t_l}), \text{not } p_{l+1}(\vec{t_{l+1}}), \ldots, \text{not } p_{l+k}(\vec{t_{l+k}})$ is a ground instance of $r$, since the truth of $e(\vec{t_z})$ implies, by construction of $e(\vec{z})$, that for all $0 \le i, i' \le l + k$, $1 \le j \le k_i$, and $1 \le j' \le k_{i'}$, if $x_{i,j} = c$ then $t_{i,j} = c$, and if $x_{i,j} = x_{i',j'}$ then $t_{i,j} = t_{i',j'}$.

An immediate consequence is that $\text{SM}(\mathcal{P}) = \text{SM}(\mathcal{P}')$ for a program $\mathcal{P}$ and the program $\mathcal{P}'$, obtained by replacing each rule in $\mathcal{P}$ by its normalization. Thus, w.l.o.g. we just consider normalized Datalog$^{\neg}$ queries.

Let $Q = \langle q, \mathcal{P} \rangle$ be a non-recursive Datalog$^{\neg}$ query with query predicate $q$ of arity $n$, and $\mathcal{P}$ its normalized program. Since $\mathcal{P}$ is non-recursive, there exists an enumeration of its intensional predicates, such that the following holds: Let $o(p)$, $1 \le o(p) \le |\text{IDB}|$, the enumeration index assigned to an intensional predicate $p$, where $|\text{IDB}|$ denotes the number if IDB predicates. Then $(i)$ $o(p) \ne o(p')$ if $p \ne p'$; and $(ii)$ for every rule $r \in \mathcal{P}$ such that $head(r) = p(\vec{x})$, if $p'(\vec{x'}) \in body(r)$, then $p'$ is an extensional predicate or $o(p') < o(p)$.

Furthermore, $\mathcal{P}$ is constraint-free. Hence, for any database, i.e., for any finite set of facts $F$, the program $F \cup \mathcal{P}$ has a unique stable model $M$, i.e., $\text{SM}(\mathcal{P}) = \{M\}$. In particular, let $D$ be a database for $\chi = \langle \Psi, \Sigma \rangle$, let $R_1, ..., R_n$ be $n$ databases such that $R_i = R'_i \cap E$, where $E$ is a weak repair envelope for $D$, $A = D \cap E$, and $R'_i$ is a repair for $A$, and let

$S = D \setminus E$. Then, $S \cup R_i \cup \mathcal{P}$, $1 \leq i \leq n$, has a unique stable model which we will denote by $M_i$.

Towards a proof of Proposition 6.1, let $D_m$ be the marked database built by processing $R_1, \ldots R_n$, and consider any enumeration, $o$, of the IDB predicates of $\mathcal{P}$ w.r.t. $D$ that satisfies $(i)$ and $(ii)$. We first show that for $1 \leq i \leq n$ and every intensional predicate $p$, there exists a tuple $\langle \vec{t}, m \rangle$ in the view relation $p_m$ such that $m(i) = 1$, i.e., its mark at position $i$ is 1, iff $p(\vec{t}) \in M_i$. The proof is by induction on $o(p)$.

Induction base: $o(p) = 1$.

($\subseteq$) Let $\langle \vec{t}, m \rangle \in p_m$ and $m(i) = 1$. Then by the definition of $p_m$, it holds that $m = \texttt{SUMBIT}(m_1, \ldots, m_j)$, $j \geq 1$, and $m_k(i) = 1$ for at least one mark $m_k$, $1 \leq k \leq j$. Let $m'$ be any of the marks $m_k$ such that $m_k(i) = 1$. Then, $\langle \vec{t}, m' \rangle \in SQL_{r_j}$, for some $1 \leq j \leq l$, where $r_1, \ldots, r_l$ are all the rules in $\mathcal{P}$ with head predicate $p$. Let

$$r_j = p(\vec{y}_0) \leftarrow p_1(\vec{y}_1), \ldots, p_{j_l}(\vec{y}_{j_l}), not\ p_{j_{l+1}}(\vec{y}_{j_{l+1}}), \ldots, not\ p_{j_{l+k}}(\vec{y}_{j_{l+k}}), e(\vec{z}).$$

Since $o(p) = 1$, by Condition $(ii)$ we conclude that $body(r_j)$ only contains extensional predicates. By the definition of $SQL_{r_j}$, there exist tuples $\langle \vec{t}_1, m_1' \rangle, \ldots, \langle t_{j_{l+k}}, m_{j_{l+k}}' \rangle$ that cause the selection of $\langle \vec{t}, m' \rangle$, such that

(1) $\langle \vec{t}_h, m_h' \rangle \in p_{h_m}$ and $m_h'(i) = 1$, for $1 \leq h \leq j_l$,

(2) $\langle \vec{t}_h, m_h' \rangle \in SQL_{r,h}$ and $m_h'(i) = 0$, for $j_{l+1} \leq h \leq j_{l+k}$,

(3) the equality conditions $e(\vec{t}_z)$ evaluate to true under the corresponding substitutions $\vec{y}_0|_{\vec{t}}$ and $\vec{y}_h|_{\vec{t}_h}$, for $1 \leq h \leq j_{l+k}$.

For $1 \leq h \leq j_l$, we conclude from Condition 1 that $p_h(\vec{t}_h) \in S \cup R_i$ since $p_h$ is extensional. For every $j_{l+1} \leq h \leq j_{l+k}$, Condition 2 implies, by construction of $SQL_{r,h}$, that either $\langle \vec{t}_h, m_h' \rangle \in p_{h_m}$ and $m_h'(i) = 0$, or that there does not exist a mark $m_h$, such that $\langle \vec{t}_h, m_h \rangle \in p_{h_m}$. In both cases, we conclude that $p_h(\vec{t}_h) \notin S \cup R_i$ and thus, since $p_h$ is extensional, that $p_h(\vec{t}_h) \notin M_i$. Condition 3 guarantees that $r_j(\vec{t'}) = p(\vec{t}) \leftarrow p_1(\vec{t}_1), \ldots, p_{j_l}(\vec{t}_{j_l}), not\ p_{j_{l+1}}(\vec{t}_{j_{l+1}}), \ldots, not\ p_{j_{l+k}}(\vec{t}_{j_{l+k}}), e(\vec{t}_z)$ is a ground instance of $r_j$ and $e(\vec{t}_z)$ is true. Thus, $r_j(\vec{t'})$ is a ground instance of rule $r_j \in \mathcal{P}$, such that $M_i \models body(r_j(\vec{t'}))$. Consequently, $M_i \models head(r_j(\vec{t'}))$, i.e., $p(\vec{t}) \in M_i$.

($\supseteq$) Let $p(\vec{t}) \in M_i$. Since $p$ is intensional, there exists a ground rule $r(\vec{t'}) = p(\vec{t}) \leftarrow p_1(\vec{t}_1), \ldots, p_l(\vec{t}_l), not\ p_{l+1}(\vec{t}_{l+1}), \ldots, not\ p_{l+k}(\vec{t}_{l+k}), e(\vec{t}_z)$ in $\mathcal{P}$, such that $p_h(\vec{t}_h) \in M_i$, for $1 \leq h \leq l$, and $p_h(\vec{t}_h) \notin M_i$, for $l + 1 \leq h \leq l + k$. Furthermore, since $o(p) = 1$, $p_h$ is extensional for $1 \leq h \leq l + k$. Thus, for $1 \leq h \leq l$ there exist tuples $\langle \vec{t}_h, m_h \rangle \in p_{h_m}$ such that $m_h(i) = 1$, while for $l + 1 \leq h \leq l + k$ either (a) $\langle \vec{t}_h, m_h \rangle \in p_{h_m}$ such that $m_h(i) = 0$, or (b) no mark $m_h$ exists such that $\langle \vec{t}_h, m_h \rangle \in p_{h_m}$. Since $e(\vec{t}_z)$ is true for the given ground instance, in Case (a) $\langle \vec{t}_h, m_h \rangle$ is also in $SQL_{r,h}$; in Case (b) $\langle \vec{t}_h, {}'0 \ldots 0' \rangle \in SQL_{r,h}$. Thus, for $l + 1 \leq h \leq l + k$, in any case there exists a tuple $\langle \vec{t}_h, m_h \rangle \in SQL_{r,h}$, such that $m_h(i) = 0$. Consequently, $\langle \vec{t}, m \rangle \in SQL_r$, where $m = m_1\ \texttt{ANDBIT}\ \ldots\ \texttt{ANDBIT}\ m_l\ \texttt{ANDBIT}\ \texttt{INVBIT}(m_{l+1})\ \texttt{ANDBIT}\ \ldots\ \texttt{ANDBIT}\ \texttt{INVBIT}(m_{l+k})$, i.e., $m(i) = 1$. By the definition of view $p_m$, we conclude that $\langle \vec{t}, m' \rangle \in p_m$ and $m'(i) = 1$.

Induction hypothesis: $\langle \vec{t}, m \rangle \in p_m$ and $m(i) = 1$ iff $p(\vec{t}) \in M_i$, for all intensional predicates $p$ such that $o(p) < n$.

Induction step: $o(p) = n$.

($\subseteq$) Let $\langle \vec{t}, m \rangle \in p_m$ and $m(i) = 1$. Then by the definition of $p_m$, it holds that $m = \texttt{SUMBIT}(m_1, \ldots, m_j)$, $j \geq 1$, and $m_k(i) = 1$ for at least one mark $m_k$, $1 \leq k \leq j$. Let $m'$

be any of the marks $m_k$ such that $m_k(i) = 1$. Then, $\langle \vec{t}, m' \rangle \in SQL_{r_j}$, for some $1 \le j \le l$, where $r_1, \ldots, r_l$ are all the rules in $\mathcal{P}$ with head predicate $p$. Let

$$r_j = p(\vec{y}_0) \leftarrow p_1(\vec{y}_1), \ldots, p_{j_l}(\vec{y}_{j_l}), not\ p_{j_{l+1}}(\vec{y}_{j_{l+1}}), \ldots, not\ p_{j_{l+k}}(\vec{y}_{j_{l+k}}), e(\vec{z}).$$

Since $o(p) = n$, by Condition $(ii)$ we conclude that $body(r_j)$ just contains extensional predicates or intensional predicates $p'$ such that $o(p') < n$. By the definition of $SQL_{r_j}$, there exist tuples $\langle \vec{t}_1, m'_1 \rangle, \ldots, \langle \vec{t}_{j_{l+k}}, m'_{j_{l+k}} \rangle$ that cause the selection of $\langle \vec{t}, m' \rangle$, such that

(1) $\langle \vec{t}_h, m'_h \rangle \in p_{h_m}$ and $m'_h(i) = 1$, for $1 \le h \le j_l$,
(2) $\langle \vec{t}_h, m'_h \rangle \in SQL_{r,h}$ and $m'_h(i) = 0$, for $j_{l+1} \le h \le j_{l+k}$,
(3) the equality conditions $e(\vec{t}_z)$ evaluate to true under the corresponding substitutions $\vec{y}_0|_{\vec{t}}$ and $\vec{y}_h|_{\vec{t}_h}$, for $1 \le h \le j_{l+k}$.

For $1 \le h \le j_l$, we conclude from Condition 1 that $p_h(\vec{t}_h) \in M_i$ since $p_h$ is either extensional, i.e., $p_h(\vec{t}_h) \in S \cup R_i$, or intensional with $o(p_h) < n$, i.e., the induction hypothesis applies. For every $j_{l+1} \le h \le j_{l+k}$, Condition 2 implies, by construction of $SQL_{r,h}$, that either (a) $\langle \vec{t}_h, m'_h \rangle \in p_{h_m}$ and $m'_h(i) = 0$, or that (b) there does not exist a mark $m_h$, such that $\langle \vec{t}_h, m_h \rangle \in p_{h_m}$. We show that $p_h(\vec{t}_h) \notin M_i$ follows for $j_{l+1} \le h \le j_{l+k}$. If $p_h$ is extensional, in both cases, we conclude that $p_h(\vec{t}_h) \notin M_i$, since $p_h(\vec{t}_h) \notin S \cup R_i$. For intensional $p_h$, towards a contradiction, assume that $p_h(\vec{t}_h) \in M_i$. We know that $o(p_h) < n$ and conclude by the induction hypothesis, that there exists a mark $m_h$, such that $m_h(i) = 1$ and $\langle \vec{t}_h, m_h \rangle \in p_{h_m}$. This contradicts (b), so (a) has to be the case, i.e., also $\langle \vec{t}_h, m'_h \rangle \in p_{h_m}$, and $m_h \ne m'_h$. However, this cannot be the case by construction of the view $p_m$: tuples which just differ in the mark attribute are grouped and their mark is computed by the aggregate function SUMBIT. As a consequence, $\langle \vec{t}_h, m_h \rangle \in p_{h_m}$ and $\langle \vec{t}_h, m'_h \rangle \in p_{h_m}$ implies $m_h = m'_h$, contradiction. This proves $p_h(\vec{t}_h) \notin M_i$ for $j_{l+1} \le h \le j_{l+k}$. Eventually, Condition 3 guarantees that $r_j(\vec{t'}) = p(\vec{t'}) \leftarrow p_1(\vec{t}_1), \ldots, p_{j_l}(\vec{t}_{j_l}), not\ p_{j_{l+1}}(\vec{t}_{j_{l+1}}), \ldots, not\ p_{j_{l+k}}(\vec{t}_{j_{l+k}}), e(\vec{t}_z)$ is a ground instance of $r_j$ and $e(\vec{t}_z)$ is true. Thus, $r_j(\vec{t'})$ is a ground instance of rule $r_j \in \mathcal{P}$, such that $M_i \models body(r_j(\vec{t'}))$. Consequently, $M_i \models head(r_j(\vec{t'}))$ has to hold, i.e., $p(\vec{t}) \in M_i$.

$(\supseteq)$ Let $p(\vec{t}) \in M_i$. Since $p$ is intensional, there exists a ground rule $r(\vec{t'}) = p(\vec{t'}) \leftarrow p_1(\vec{t}_1), \ldots, p_l(\vec{t}_l), not\ p_{l+1}(\vec{t}_{l+1}), \ldots, not\ p_{l+k}(\vec{t}_{l+k}), e(\vec{t}_z)$ in $\mathcal{P}$, such that $p_h(\vec{t}_h) \in M_i$, for $1 \le h \le l$, and $p_h(\vec{t}_h) \notin M_i$, for $l+1 \le h \le l+k$. Furthermore, since $o(p) = n$, by Condition $(ii)$ $p_h$ is either extensional or $o(p_h) < n$, for $1 \le h \le l+k$. Thus, for $1 \le h \le l$ there exist tuples $\langle \vec{t}_h, m_h \rangle \in p_{h_m}$ such that $m_h(i) = 1$, while for $l+1 \le h \le l+k$ either (a) $\langle \vec{t}_h, m_h \rangle \in p_{h_m}$ such that $m_h(i) = 0$, or (b) no mark $m_h$ exists such that $\langle \vec{t}_h, m_h \rangle \in p_{h_m}$. Since $e(\vec{t}_z)$ is true for the given ground instance, in Case (a) $\langle \vec{t}_h, m_h \rangle$ is also in $SQL_{r,h}$; in Case (b) $\langle \vec{t}_h, '0 \ldots 0' \rangle \in SQL_{r,h}$. Thus, for $l+1 \le h \le l+k$, in any case there exists a tuple $\langle \vec{t}_h, m_h \rangle \in SQL_{r,h}$, such that $m_h(i) = 0$. Consequently, $\langle \vec{t}, m \rangle \in SQL_r$, where $m = m_1$ ANDBIT $\ldots$ ANDBIT $m_l$ ANDBIT INVBIT$(m_{l+1})$ ANDBIT $\ldots$ ANDBIT INVBIT$(m_{l+k})$, i.e., $m(i) = 1$. By the definition of view $p_m$, we conclude that $\langle \vec{t}, m' \rangle \in p_m$ and $m'(i) = 1$.

This proves $\langle \vec{t}, m \rangle \in p_m$ and $m(i) = 1$ iff $p(\vec{t}) \in M_i$ for every (intensional) predicate $p$ in $\mathcal{P}$.

We now prove Proposition 6.1, that is, that the set of tuples computed by $SQL_Q$ on the marked database $D_m$ coincides with the set $\bigcap_{i=1}^{n} \{\vec{t} \mid \vec{t} \in Q[R_i \cup S]\}$.

App–11

($\subseteq$) Let $\vec{t} \in SQL_Q$. Then, $\langle \vec{t}, '1 \dots 1' \rangle \in q_m$, and, hence, $q(\vec{t}\,) \in M_i$ for $1 \leq i \leq n$. Since $\mathrm{SM}(R_i \cup S \cup \mathcal{P}) = \{M_i\}$ for $1 \leq i \leq n$, by definition we obtain $\vec{t} \in Q[R_i \cup S]$ for $1 \leq i \leq n$, which implies $\vec{t} \in \bigcap_{i=1}^{n} \{\vec{t} \mid \vec{t} \in Q[R_i \cup S]\}$.

($\supseteq$) Let $\vec{t} \in \bigcap_{i=1}^{n} \{\vec{t} \mid \vec{t} \in Q[R_i \cup S]\}$. Then, $q(\vec{t}\,) \in M_i$ since $\mathrm{SM}(R_i \cup S \cup \mathcal{P}) = \{M_i\}$, for $1 \leq i \leq n$. Thus, there exist marks $m_i$, such that $\langle \vec{t}, m_i \rangle \in q_m$ and $m_i(i) = 1$, for $1 \leq i \leq n$. By the definition of view $q_m$ (in particular by the definition of SUMBIT), it follows that $m_i =' 1 \dots 1'$, for $1 \leq i \leq n$, i.e., $\langle \vec{t}, '1 \dots 1' \rangle \in q_m$. Consequently, $\vec{t} \in SQL_Q$. $\square$

Clearly, the SQL statements $SQL_r$, $SQL_h$, and $SQL_Q$ can be optimized (which will be done by the DBMS anyway), and we do not consider optimization here. We remark that the final query, $SQL_Q$, could also be integrated into the view definition, $SQL_q$, for the query predicate $q$. By keeping the query definition $SQL_Q$ separate, however, other query semantics can easily be expressed; for instance, possibilistic query semantics, which selects those tuples which are computed by the query with respect to at least one repair, is obtained by replacing the condition in the WHERE clause by $\mathbf{mark} \neq \ '0 \dots 0'$. We finally remark that the reformulation technique is amenable to other semantics of negation in queries as well. Specifically, a more cautious semantics for negation can also be accomplished with slight modifications: For a negative ground literal to be true, it has to evaluate to true in *all* repairs, i.e., the respective tuple has to be marked $'0 \dots 0'$ (or absent) in the corresponding marked table.

## E.  EXAMPLES OF LOGIC PROGRAM SPECIFICATIONS

In this section, we discuss some approaches for consistent query answering in inconsistent database that rely on the use of logic programming. The notion of repair adopted in the papers described below relies on the prototypical, natural preorder $\leq_D$, originally introduced in [Arenas et al. 1999], for which $R_1 \leq_D R_2$ iff $\triangle(R_1, D) \subseteq \triangle(R_2, D)$. The only exception is [Calì et al. 2003b], which uses the ordering $R_1 \sqsubseteq_D R_2$ iff $R_1 \cap D \supseteq R_2 \cap D$ [Calì et al. 2003a; 2003b]. However, for the set of integrity constraints and queries considered in [Calì et al. 2003b], the adoption of a different repair ordering is of no concern; as discussed in Section 9, we can use $R_1 \sqsubseteq'_D R_2$ instead of $R_1 \sqsubseteq_D R_2$ for answering a negation-free queries, like a union of conjunctive queries as in [Calì et al. 2003b]; under this ordering, the repairs coincide with the repairs under the canonical ordering $\leq_D$. Logic programs that we devise in this section refer to the football team scenario introduced in Example 1.1.

Since some of the techniques analyzed below have been applied to a data integration setting, we first recall some formal notions on data integration systems.

*Data Integration Systems.* Data integration systems are systems in charge of uniformly providing users with data residing at different sources, according to some mapping assertions. More formally, a data integration system $\mathcal{I}$ may be viewed as a triple $\langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$, where $\mathcal{G}$ is the *global schema*, which specifies the global elements exported to users, $\mathcal{S}$ is the *source schema*, which describes the structure of the data sources in the system, and $\mathcal{M}$ is the *mapping*, which establishes the relationship between the sources and the global schema [Lenzerini 2002]. Classical approaches for specifying the mapping are the *Global-As-View (GAV)* approach, which requires that every element of the global schema is associated with a view over the sources, so that its the meaning is given in terms of

the source data, and the *Local-As-View* (LAV) approach, which conversely requires the sources to be defined as views over the global schema. A third and more general approach is *Global-Local-As-View* (GLAV), which captures both LAV and GAV. It allows for mapping assertions in which a view over the source schema is put in correspondence with a view over the sources[Lenzerini 2002]. From a technical point of view, GLAV basically raises the same issues as LAV.

We point out that all notions and techniques provided in the present paper can be easily generalized to GAV data integration systems. Indeed, the mapping specification in GAV systems provide a means for populating the global schema with one (possibly inconsistent) global database instance, which can be obtained by simply evaluating the views in the mapping over the source data. This database is also called *retrieved global database* [Lenzerini 2002]. The semantics of a GAV data integration system may be then given in terms of the only retrieved global database (in this case the mapping is called exact), or it may be given in terms of all global database instances that contain the retrieved global database (in this case the mapping is called sound). In both cases, repairing a data integration system amounts to repairing the global retrieved database, and therefore our results can be straightforwardly applied to such a setting.

We further note that our techniques can be applied to some LAV data integration proposals in the literature (as discussed in Appendix E.3).

### E.1  Logic programs with unstratified negation

The paper [Calì et al. 2003b] addresses the repair problem in *GAV data integration systems* in which key constraints are issued over the global schema, and presents a technique for consistent query answering based on the use of Datalog$^{\neg}$.

More precisely, according to [Calì et al. 2003b], given a data integration system $\mathcal{I} = \langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$, key constraints in $\mathcal{G}$ can be encoded into a suitable Datalog$^{\neg}$ program $\mathcal{P}_{KD}$, whereas views in the mapping, which are expressed as union of conjunctive queries, are cast into a Datalog program $\mathcal{P}_M$. Consistent answers to a union of conjunctive queries $Q$ over $\mathcal{I}$ w.r.t. a source database $D$ are returned by the evaluation of the Datalog$^{\neg}$ query $\langle q, \mathcal{P}_Q \cup \mathcal{P}_{KD} \cup \mathcal{P}_M \rangle$, where $\langle q, \mathcal{P}_Q \rangle$ is the Datalog encoding of the query $Q$.

In the following, we provide the logic program produced by the above technique for our running example (suitably adapted to a data integration scenario). To this aim, we exploit an extension of the algorithm of [Calì et al. 2003b], provided in [Grieco et al. 2005], that allows for dealing with the exclusion dependencies specified in Example 1.1. We assume to have a data integration system $\mathcal{I}_0$ such that the global schema is equal to the relational schema $\chi_0$ given in Example 1.1, the source schema consists of the relation $s_1$, of arity 4, and the relations $s_2$, $s_3$, and $s_4$, all of arity 4, whereas the mapping, denoted $\mathcal{M}_0$, is defined by the Datalog program $player(x,y,z) \leftarrow s_1(x,y,z,w); team(x,y,z) \leftarrow s_2(x,y,z); team(x,y,z) \leftarrow s_3(x,y,z); coach(x,y,z) \leftarrow s_4(x,y,z)$. Then, the logic program for consistent query answering over $\mathcal{I}_0$, denoted $\Pi_{\mathcal{I}_0}(Q)$, is as follows.

$$q(x) \leftarrow player(x,y,z)$$
$$q(x) \leftarrow team(v,w,x)$$
$$player_D(x,y,z) \leftarrow s_1(x,y,z,w)$$
$$team_D(x,y,z) \leftarrow s_2(x,y,z)$$

$$team_D(x,y,z) \leftarrow s_3(x,y,z)$$
$$coach_D(x,y,z) \leftarrow s_4(x,y,z)$$
$$player(x,y,z) \leftarrow player_D(x,y,z) \ , \ not \ \overline{player}(x,y,z)$$
$$\overline{player}(x,y,z) \leftarrow player(x,w,z) \ , \ player_D(x,y,z) \ , \ y \neq w$$
$$team(x,y,z) \leftarrow team_D(x,y,z) \ , \ not \ \overline{team}(x,y,z)$$
$$\overline{team}(x,y,z) \leftarrow team(x,v,w) \ , \ team_D(x,y,z) \ , \ y \neq v$$
$$\overline{team}(x,y,z) \leftarrow team(x,v,w) \ , \ team_D(x,y,z) \ , \ z \neq w$$
$$coach(x,y,z) \leftarrow coach_D(x,y,z) \ , \ not \ \overline{coach}(x,y,z)$$
$$\overline{coach}(x,y,z) \leftarrow coach(x,w,z) \ , \ coach_D(x,y,z) \ , \ y \neq w$$
$$\overline{player}(x,y,z) \leftarrow player_D(x,y,z) \ , \ coach(x,w,z)$$
$$\overline{coach}(x,y,z) \leftarrow coach_D(x,y,z) \ , \ team(z,w,x)$$
$$\overline{coach}(x,y,z) \leftarrow coach_D(x,y,z) \ , \ player(x,w,z)$$
$$\overline{team}(x,y,z) \leftarrow team_D(x,y,z) \ , \ coach(z,w,x)$$

In the above program, $\mathcal{P}_Q$ consists of the first two rules; $\mathcal{P}_M$ comprises the rules ranging from the 3rd to the 6th. The program $\mathcal{P}_{KD}$ contains the rules ranging from the 7th to the 13th, whereas the last four rules, which we denote by $\mathcal{P}_{ED}$, encode exclusion dependencies. Informally, for each global relation $r$, the above program contains (i) a relation $r_D$ that represents $r^{ret(\mathcal{I},D)}$; (ii) a relation $r$ that represents a subset of $r^{ret(\mathcal{I},D)}$ that is consistent with the key constraints and the exclusion dependencies for $r$; (iii) an auxiliary relation $\overline{r}$. We can easily see that $\Pi_{\mathcal{I}_0}(Q) = \mathcal{P}_M \cup \Pi_{\chi_0}(Q)$, where $\Pi_{\chi_0}(Q)) = \Pi_{\Sigma_0} \cup \Pi_Q$ is the logic specification for querying the relational database $\chi_0$, in which $\Pi_{\Sigma_0} = \mathcal{P}_{KD} \cup \mathcal{P}_{ED}$ and $\Pi_Q = \mathcal{P}_Q$.

We point out that in [Calì et al. 2003b], together with key constraints, also (existentially-quantified) inclusion dependencies in the global schema $\mathcal{G}$ are considered. In this respect, a query reformulation technique is given that, on the basis of inclusion dependencies on $\mathcal{G}$, rewrites the user query $Q$ into a new union of conjunctive queries $Q_{ID}$, again expressed over the global schema, in such a way that the consistent answers to $Q$ over $\mathcal{I}$ w.r.t. a source database $D$ coincide with consistent answers to $Q_{ID}$ over $\mathcal{I}'$ w.r.t. $D$, where $\mathcal{I}'$ is obtained from $\mathcal{I}$ by dropping the inclusion dependencies of $\mathcal{G}$. In other words, after computing $Q_{ID}$, it is possible to proceed as if inclusion dependencies were not specified on the global schema, i.e., by providing the logic specification for querying $\mathcal{I}'$ with $Q_{ID}$ described above. Hence, after the first reformulation, the problem of computing consistent answers in the above setting and our problem coincide.

## E.2 Logic programs with exceptions

A specification of database repairs for consistent query answering in inconsistent databases exploiting logic programs with exceptions (LPEs) is presented in [Arenas et al. 2003]. We recall that this sort of programs, firstly introduced by [Kowalski and Dadri 1990], contains both *default rules*, i.e., classic clauses with classic negation in the body literals, and *exception rules*, i.e., clauses with negative heads whose conclusion overrides conclusions of default ones. Actually, [Arenas et al. 2003] presents an extension of LPEs for accommodating both negative defaults and extended disjunctive exceptions whose semantics is given in terms of *e-answer sets*, and shows how these models are, in fact, in correspondence with

standard stable models of a suitable standard disjunctive logic program.

In more detail, the transformation in [Arenas et al. 2003] associates to each relation $p$ in the database schema a new relation $p'$ corresponding to its repaired version, and defines $\Pi_\Sigma$ to contain three set of rules: *(i)* triggering exceptions, *(ii)* stabilizing exceptions, and *(iii)* persistence defaults. Let us, for instance, consider our running example. Then, triggering exception rules are as follows.

$$\neg\, player'(x,y,z) \vee \neg\, player'(x,y_1,z) \leftarrow player(x,y,z),\ player(x,y_1,z),\ y \neq y_1.$$
$$\neg\, team'(x,y,z) \vee \neg\, team'(x,y_1,z_1) \leftarrow team(x,y,z),\ team(x,y_1,z_1),\ y \neq y_1$$
$$\neg\, team'(x,y,z) \vee \neg\, team'(x,y_1,z_1) \leftarrow team(x,y,z),\ team(x,y_1,z_1),\ z \neq z_1$$
$$\neg\, coach'(x,y,z) \vee \neg\, coach'(x,y_1,z) \leftarrow coach(x,y,z),\ coach(x,y_1,z),\ y \neq y_1$$
$$\neg\, coach'(x,y,z) \vee \neg\, player'(x,y_1,z) \leftarrow coach(x,y,z),\ player(x,y_1,z)$$
$$\neg\, coach'(x,y,z) \vee \neg\, team'(z,y_1,x) \leftarrow coach(x,y,z),\ team(z,y_1,x)$$

The above rules represent a suitable rewriting of the integrity constraints that encodes the basic way of repairing each inconsistency. For example, a conflict on a key is resolved by deleting one of the tuples that cause the conflict, i.e., by not including this tuple in the extension of the corresponding primed predicate. Notice that, in the case of (universally quantified) inclusion dependencies, it is possible to have repairs by adding tuples. For instance, the constraint $p(x,y) \supset q(x,y)$ would be repaired with the rule $\neg\, p'(x,y) \vee q'(x,y) \leftarrow p(x,y), not\ q(x,y)$.

Stabilizing exception rules and persistence defaults have been introduced for technical reasons. Indeed, rules of the former kind state that each integrity constraint must be eventually satisfied in the repair while rules of the latter kind impose that by default each relation $p'$ contains the facts in $p$.

Given the rewriting $\Pi_\Sigma$, the user query can be simply issued over the primed relations, i.e., the program $\Pi_Q$ is easily obtained by substituting in the user query (suitably expressed in Datalog) each predicate $p$ with its repaired version $p'$.

## E.3 Programs with Annotation Constants

The paper [Barceló and Bertossi 2003] proposes to specify database repairs by means of disjunctive normal programs under the stable model semantics. To this aim, suitable annotations are used in an extra argument introduced in each (non built-in) predicate of the logic program, for marking the operations of insertion and deletion of tuples required in the repair process. The idea of annotating predicates has been inspired by the Annotated Predicate Calculus [Kifer and Lozinskii 1992], a non-classical logic in which inconsistencies may be accommodated without trivializing reasoning. The values used in [Barceló and Bertossi 2003] for the annotations are:

—$\mathbf{t_d}$ and $\mathbf{f_d}$, which indicate whether, before the repair, a given tuple is in the database or not, respectively;

—$\mathbf{t_a}$ and $\mathbf{f_a}$, which represent advisory values that indicate how to resolve possible conflicts, i.e., a tuple annotated with $\mathbf{t_a}$ (resp. $\mathbf{f_a}$) has to be inserted (resp. deleted) in the database;

—$\mathbf{t^*}$ and $\mathbf{f^*}$, which indicate whether a given tuple is in the repaired database or not, respectively.

For instance, the annotated logic program used for solving the conflicts on the key of the relation $player$ in our running example is as follows:

$$player(x, y, z, \mathbf{t}^*) \leftarrow player(x, y, z, \mathbf{t_d})$$
$$player(x, y, z, \mathbf{t}^*) \leftarrow player(x, y, z, \mathbf{t_a})$$
$$player(x, y, z, \mathbf{f}^*) \leftarrow not\ player(x, y, z, \mathbf{t_d})$$
$$player(x, y, z, \mathbf{f}^*) \leftarrow player(x, y, z, \mathbf{f_a})$$
$$player(x, y, z, \mathbf{f_a}) \vee player(x, y_1, z, \mathbf{f_a}) \leftarrow player(x, y, z, \mathbf{t}^*),\ player(x, y_1, z, \mathbf{t}^*),\ y \neq y_1.$$

Furthermore, each fact in the original database is assumed to be annotated by $\mathbf{t_d}$.

Intuitively, the last rule says that when the key of the relation *player* is violated (body of the rule), the database instance has to be repaired according to one of the two alternatives shown in the head. Possible interaction between different constraints are then taken into account by the other rules, which force the repair process to continue and stabilize in a state in which all the integrity constraints hold. Indeed, annotations $\mathbf{t}^*$ and $\mathbf{f}^*$ can feed back rules of the last kind, until consistency is restored. This should be evident if we consider also a constraint of the form $coach(x, y, z) \supset player(x, y, z)$ (we disregard exclusion dependencies of our running example for a while). This constraint is repaired with the rule

$$coach(x, y, z, \mathbf{f_a}) \vee player(x, y, z, \mathbf{t_a}) \leftarrow coach(x, y, z, \mathbf{t}^*), player(x, y, z, \mathbf{f}^*),$$

besides the rules for the predicate *coach* that compute facts with annotations $\mathbf{t}^*$ (resp. $\mathbf{f}^*$) from facts annotated by $\mathbf{t_d}$ or $\mathbf{t_a}$ (resp. $\mathbf{f_d}$ or $\mathbf{f_a}$).

The program $\Pi_Q$ is then computed by reformulating the original query according to the annotations: in our running example, we have

$$q(x) \leftarrow player(x, y, z, \mathbf{t_a}) \vee (player(x, y, z, \mathbf{t_d}) \wedge \neg player(x, y, z, \mathbf{f_a}))$$
$$q(x) \leftarrow team(v, w, x, \mathbf{t_a}) \vee (team(v, w, x, \mathbf{t_d}) \wedge \neg team(v, w, x, \mathbf{f_a})).$$

The above rewriting is proposed in [Barceló and Bertossi 2003] for the setting of a single database. In the line of the discussion of data integration systems in the beginning of Appendix E, this technique can be straightforwardly extended to work in GAV data integration systems. An interesting, more complex generalization to the LAV setting appears instead in [Bravo and Bertossi 2003; 2005]. Since in LAV each source relation is associated with a query over the global schema, an exact specification of which data of the sources fit the global schema is actually missing. In general, given a source database, several different ways of populating the global schema according to the mapping may exist. Hence, not a single but multiple retrieved global databases must be taken into account for repairing. According to [Bravo and Bertossi 2003; 2005], the repairs are defined as those consistent global databases which have a minimal (under set inclusion) symmetric difference to one of the minimal (again, under set inclusion) retrieved global databases. In other words, each such retrieved global database is repaired by adopting the classic preorder of [Arenas et al. 1999]. These repairs can be obtained from the stable models of a suitable disjunctive logic program, which comprises rules for the encoding of integrity constraints constructed as in [Barceló and Bertossi 2003], and specific rules for computing the minimal retrieved global databases.

## F. FURTHER EXPERIMENTS

### F.1 Assessing the Need of Localization

We conducted a set of experiments to assess the importance of localization approaches, even in those situations that involve very simple logic programs for computing consistent answers. To this aim, we considered the database schemas $\chi_k$, which contains a relation of the form $p(x, y)$, where $x$ is the key; and, $\chi_e$, which contains relations of the forms $p(x, y)$ and $q(v, w)$, with an exclusion dependency between attributes $x$ and $v$.
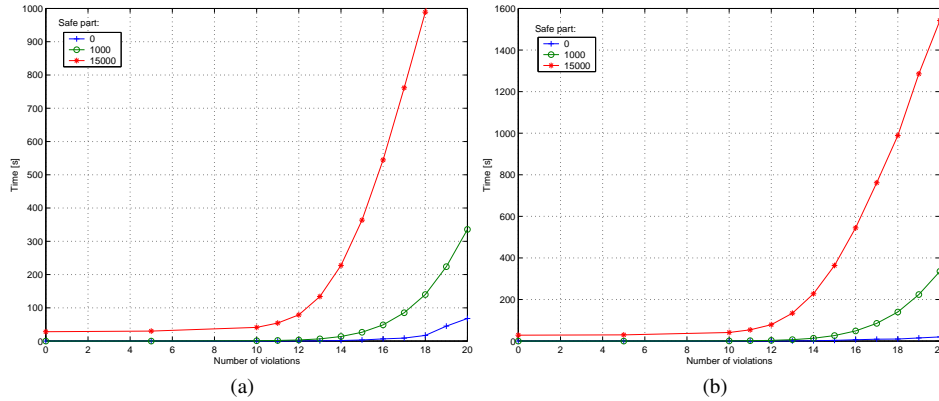


(a)                     (b)

Fig. 9. Stable model computation time in DLV system w.r.t. number of conflicts. (a) One Key. (b) One exclusion dependency.

In Figure 9, we report the time needed in the DLV system for computing the stable models of the logic encoding for repairing the two databases w.r.t. the number of conflicts, for different database sizes (where the sizes of the safe part are printed). This first set of experiments is particularly interesting, since the cost of computing all the stable models is a reasonable lower-bound for the cost of computing consistent query answers, given that most of the state-of-art answer set engines provide support for "Boolean" query answering, that is, for deciding whether a given ground fact is entailed in any/all models, but not for computing non-ground queries. From Figure 9, we observe that DLV scales exponentially in the number of conflicts, since repairs are in general exponential in this number. Moreover, since the size of each repair is about the size of the database, the number of processed tuples in the DLV system is, in turn, exponential in the size of the whole database.

The results of Figure 9 stimulated the development of techniques for computing consistent answers even to non-ground queries in stable models engines. Moreover, our preliminary investigations showed that most of the current answer set engines are not well-suited for data base applications since they do not offer primitives for interfacing with databases, for instance, for importing and exporting relations or views. And, in fact, in our first experiments it was necessary to write wrappers that interface the output of the answer set engines and provides I/O functionalities. The DLV system, however, provides some interfacing modules to automatically access a relational DBMS by means of standard ODBC calls, and more importantly, provides support for non-ground queries akin to our techniques. Still, the need for instantiating the logic program for consistent query answering over large data sets makes the use of these systems infeasible in practice.
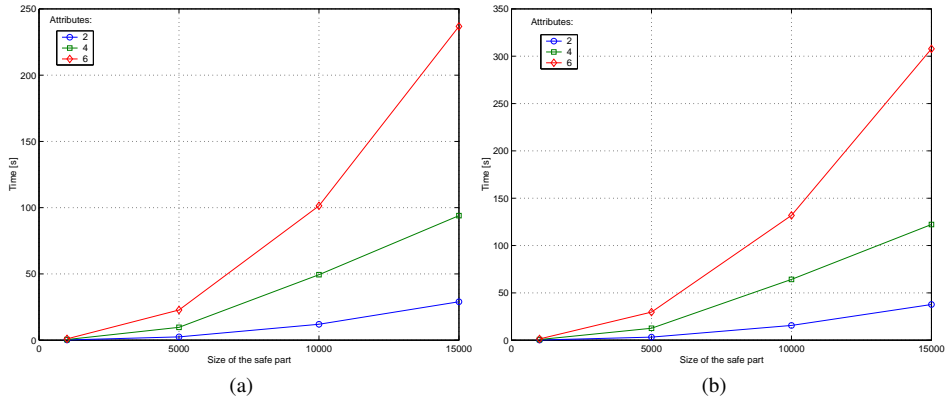
Fig. 10. Execution time in DLV system w.r.t. size of the safe part, for different numbers of attributes in relation. (a) One key. (b) One exclusion dependency.
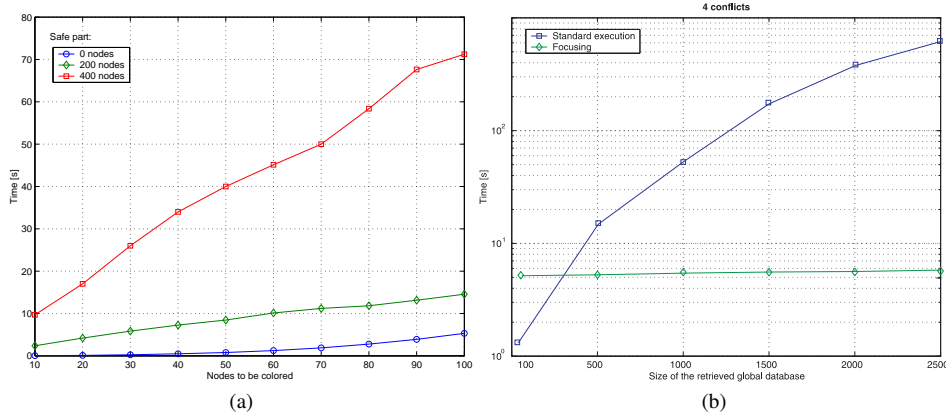


Fig. 11. 3Coloring. (a) Execution time in DLV w.r.t. number of nodes (i.e., conflicts). (b) Comparison with the optimization method.

Indeed, in a second set of experiments, we tested the scaling of DLV in answering non-ground queries. Figure 10 reports the results for evaluating in DLV some non-ground queries on the two databases $\chi_k$ and $\chi_e$. Interestingly, the support for non-ground queries appears to be quite powerful, since the system scales well in the size of the input database for a fixed number of conflicts. However, the performance is not suited for real database applications. In fact, for 15000 tuples it requires more than 200 seconds for computing answers. Moreover, the curves rapidly increase if the number of attributes (arities of the relations) grows. This behavior does not correspond to the intrinsic complexity of the problem instances, which can be formally proven to be solvable in polynomial time.

In fact, a careful analysis of the execution time showed that most of the time spent by DLV is for instantiating the logic program with the whole database. Hence, our localization approach to query answering may help speed-up performances, by reducing the size of the program to be instantiated in DLV and, hence, the time needed for the execution.

## F.2   3-Coloring

As a further example, we encoded the classical NP-complete graph 3-coloring problem into a consistent query answering problem over a database $\chi_{3c}$ containing the relations $edge(Node,Node)$ and $colored(Node, Color)$, where the attribute $Node$ is established to be the key for $colored$. Then, for a database $D_{3c}$ for $\chi_{3c}$, we fixed a number of nodes and generated facts in $edge$ producing a graph; moreover, for each node $i$ in the affected part, we generated three facts: $colored(i, red)$, $colored(i, blue)$, $colored(i, yellow)$. Clearly $D_{3c}$ is inconsistent with the key constraint on the relation $colored$, and each node creates three conflicts. Moreover, in each repair of $D_{3c}$ only one of the three facts involved in each constraint violation can be maintained.

Now, consider the query $q \leftarrow edge(x, y), colored(x, C), colored(y, C)$. As easily seen, it evaluates to true on $D_{3c}$ iff there is no legal 3-coloring for the graph in $D_{3c}$.

We encoded the problem of establishing consistent answers to the query $q$ above over the relational schema $\chi_{3c}$ into a Datalog$^\neg$ program, according to the encoding proposed in [Calì et al. 2003a].

Figure 11.(a) reports the execution time in DLV for different values of the size of the affected part, while Figure 11.(b) reports the comparison with our approach. Again, the advantage of the localization technique is evident when the size of the database increases.