

**I N F S Y S  
R E S E A R C H  
R E P O R T**



**INSTITUT FÜR INFORMATIONSSYSTEME  
ARBEITSBEREICH WISSENSBASIERTE SYSTEME**

**REPAIR LOCALIZATION  
FOR QUERY ANSWERING  
FROM INCONSISTENT DATABASES**

**Thomas Eiter      Michael Fink      Gianluigi Greco  
Domenico Lembo**

**INFSYS RESEARCH REPORT 1843-07-01**

**JANUARY 2007**

Institut für Informationssysteme  
AB Wissensbasierte Systeme  
Technische Universität Wien  
Favoritenstrasse 9-11  
A-1040 Wien, Austria  
Tel: +43-1-58801-18405  
Fax: +43-1-58801-18493  
sek@kr.tuwien.ac.at  
www.kr.tuwien.ac.at





## REPAIR LOCALIZATION FOR QUERY ANSWERING FROM INCONSISTENT DATABASES

Thomas Eiter<sup>1</sup> and Michael Fink<sup>1</sup> and Gianluigi Greco<sup>2</sup> and Domenico Lembo<sup>3</sup>

**Abstract.** 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 classes of queries and constraints, 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.

---

<sup>1</sup>Institute of Information Systems, Knowledge-Based Systems Group, TU Vienna, Favoritenstraße 9-11, A-1040 Vienna, Austria. Email: (eiter | michael)@kr.tuwien.ac.at

<sup>2</sup>Dept. of Mathematics, Univ. Calabria, Via Pietro Bucci 30B, I-87036 Rende, Italy. Email: ggreco@mat.unical.it

<sup>3</sup>DIS University of Roma “La Sapienza”, Via Salaria 113, I-00198 Roma, Italy. Email: lembo@dis.uniroma1.it

**Acknowledgements:** This work has been partially supported by the European Commission FET Programme Projects IST-2002-33570 INFOMIX and IST-2001-37004 WASP, and the Austrian Science Fund (FWF) project P18019-N04.

Copyright © 2007 by the authors

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Contributions . . . . .	2
<b>2</b>	<b>Preliminaries</b>	<b>4</b>
2.1	Data Model . . . . .	4
2.2	Datalog <sup>V,¬</sup> Programs and Queries . . . . .	6
<b>3</b>	<b>Consistent Query Answering Framework</b>	<b>7</b>
3.1	A General Framework for Database Repairs . . . . .	7
3.2	Constructible Repairs and Safe Constraints . . . . .	8
3.3	Queries and Consistent Answers . . . . .	8
<b>4</b>	<b>Locality Properties for Repairing Inconsistent Databases</b>	<b>9</b>
4.1	General Constraints . . . . .	10
4.2	Special Constraints . . . . .	14
4.2.1	Constraints $C_1$ and $C_2$ . . . . .	14
4.2.2	Constraints $C_0$ . . . . .	15
<b>5</b>	<b>Query Answering through Localized Repairs</b>	<b>16</b>
5.1	Recombination Step . . . . .	17
5.2	Repair Factorization . . . . .	17
5.2.1	Recombination of Independent Factors . . . . .	19
<b>6</b>	<b>Logic Programming for Consistent Query Answering</b>	<b>21</b>
6.1	General Architecture for Repair Compilation . . . . .	21
6.2	Grouped Repair Computation . . . . .	23
<b>7</b>	<b>Experimental Results</b>	<b>24</b>
7.1	Benchmark Databases and Compared Methods . . . . .	25
7.2	The Football Teams Example . . . . .	26
7.3	Scalability Assessment . . . . .	27
<b>8</b>	<b>Discussion and Conclusion</b>	<b>29</b>
<b>A</b>	<b>Proofs for Section 3</b>	<b>33</b>
<b>B</b>	<b>Proofs for Section 5</b>	<b>34</b>
<b>C</b>	<b>Grouped Repair Computation</b>	<b>37</b>
C.1	Query Reformulation . . . . .	37
C.2	Scaling the Technique . . . . .	43

<b>D</b>	<b>Examples of Logic Program Specifications</b>	<b>43</b>
D.1	Logic programs with unstratified negation . . . . .	44
D.2	Logic programs with exceptions . . . . .	45
D.3	Programs with Annotation Constants . . . . .	46
<b>E</b>	<b>Further Experiments</b>	<b>47</b>
E.1	Assessing the Need of Localization . . . . .	47
E.2	3-Coloring . . . . .	49

# 1 Introduction

A database is inconsistent if it does not satisfy the integrity constraints specified over its schema. This may happen for different reasons [4]; 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  $\mathcal{S}_0$  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  $\mathcal{S}_0$  (possibly built by integrating some autonomous data sources):

$player^{D_0}$ :	<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td style="padding: 2px;">10</td><td style="padding: 2px;">Totti</td><td style="padding: 2px;">RM</td></tr> <tr><td style="padding: 2px;">9</td><td style="padding: 2px;">Ronaldo</td><td style="padding: 2px;">BC</td></tr> </table>	10	Totti	RM	9	Ronaldo	BC
10	Totti	RM					
9	Ronaldo	BC					

$team^{D_0}$ :	<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td style="padding: 2px;">RM</td><td style="padding: 2px;">Roma</td><td style="padding: 2px;">10</td></tr> <tr><td style="padding: 2px;">BC</td><td style="padding: 2px;">Barcelona</td><td style="padding: 2px;">8</td></tr> <tr><td style="padding: 2px;">RM</td><td style="padding: 2px;">Real Madrid</td><td style="padding: 2px;">10</td></tr> </table>	RM	Roma	10	BC	Barcelona	8	RM	Real Madrid	10
RM	Roma	10								
BC	Barcelona	8								
RM	Real Madrid	10								

$coach^{D_0}$ :	<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td style="padding: 2px;">7</td><td style="padding: 2px;">Capello</td><td style="padding: 2px;">RM</td></tr> </table>	7	Capello	RM
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 [10]. 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 [41]. 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 [3]. 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  $\mathcal{S}_0$  violates the key constraint on  $team$ , witnessed by  $team(RM, Roma, 10)$  and  $team(RM, Real Madrid, 10)$ .

A repair results by removing exactly one of these facts. 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  $\{(RM, 10), (BC, 8)\}$ .  $\square$

$player^{R_1}$ :	<table style="border-collapse: collapse; text-align: center;"> <tr><td style="padding: 2px 10px;">10</td><td style="padding: 2px 10px;">Totti</td><td style="padding: 2px 10px;">RM</td></tr> <tr><td style="padding: 2px 10px;">9</td><td style="padding: 2px 10px;">Ronaldinho</td><td style="padding: 2px 10px;">BC</td></tr> </table>	10	Totti	RM	9	Ronaldinho	BC	$team^{R_1}$ :	<table style="border-collapse: collapse; text-align: center;"> <tr><td style="padding: 2px 10px;">RM</td><td style="padding: 2px 10px;">Roma</td><td style="padding: 2px 10px;">10</td></tr> <tr><td style="padding: 2px 10px;">BC</td><td style="padding: 2px 10px;">Barcelona</td><td style="padding: 2px 10px;">8</td></tr> </table>	RM	Roma	10	BC	Barcelona	8	$coach^{R_1}$ :	<table style="border-collapse: collapse; text-align: center;"> <tr><td style="padding: 2px 10px;">7</td><td style="padding: 2px 10px;">Capello</td><td style="padding: 2px 10px;">RM</td></tr> </table>	7	Capello	RM
10	Totti	RM																		
9	Ronaldinho	BC																		
RM	Roma	10																		
BC	Barcelona	8																		
7	Capello	RM																		
$player^{R_2}$ :	<table style="border-collapse: collapse; text-align: center;"> <tr><td style="padding: 2px 10px;">10</td><td style="padding: 2px 10px;">Totti</td><td style="padding: 2px 10px;">RM</td></tr> <tr><td style="padding: 2px 10px;">9</td><td style="padding: 2px 10px;">Ronaldinho</td><td style="padding: 2px 10px;">BC</td></tr> </table>	10	Totti	RM	9	Ronaldinho	BC	$team^{R_2}$ :	<table style="border-collapse: collapse; text-align: center;"> <tr><td style="padding: 2px 10px;">BC</td><td style="padding: 2px 10px;">Barcelona</td><td style="padding: 2px 10px;">8</td></tr> <tr><td style="padding: 2px 10px;">RM</td><td style="padding: 2px 10px;">Real Madrid</td><td style="padding: 2px 10px;">10</td></tr> </table>	BC	Barcelona	8	RM	Real Madrid	10	$coach^{R_2}$ :	<table style="border-collapse: collapse; text-align: center;"> <tr><td style="padding: 2px 10px;">7</td><td style="padding: 2px 10px;">Capello</td><td style="padding: 2px 10px;">RM</td></tr> </table>	7	Capello	RM
10	Totti	RM																		
9	Ronaldinho	BC																		
BC	Barcelona	8																		
RM	Real Madrid	10																		
7	Capello	RM																		

Figure 1: Repairs of  $D_0$ .

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 [7], and for a discussion on relevant issues in the area see [15]) and some prototype implementations of systems which fit the semantic repair framework are available [9, 26, 18, 35]. Basically, these systems differ in the kinds of constraints and queries they are able to deal with. Indeed, depending on these two ingredients, the complexity of consistent query answering ranges from polynomial-time over co-NP up to  $\Sigma_2^P$  (see, e.g., [13, 16]).

## 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 on a small fragment of the database to 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 [24, 3, 4, 6, 11, 13, 14, 16, 29], cardinality-based orderings [4, 38], and weighted-based orderings [37].

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,  $C_0$ , 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 [17].
- The second class,  $C_1$ , 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,  $C_2$ , has similar constraints  $\forall \vec{x} \alpha(\vec{x}) \supset \beta(\vec{x}) \vee \phi(\vec{x})$ ; here  $\alpha(\vec{x})$  may be empty 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 classes of queries and constraints, these techniques guarantee an exponential gain compared to the basic approach.

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 DBMS. 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 [36] or Smodels [40]. Specifically:

4) We propose a formal model of inconsistency resolution via logic programming specification, which abstracts from several proposals in the literature [4, 6, 8, 11, 14, 29]. 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 [17]. Also, our results can be exploited for efficient implementation of consistent query answering techniques in general, independent of a logic-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 DBMS and stable model engines. Finally, Section 7 reports results of our experimental activity, and Section 8 concludes the paper with a brief discussion.

Some proofs as well as 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  $c_1, c_2, \dots$  under the *unique name assumption*, that is, different constants denote different real-world objects.

A *relational schema* (or simply *schema*)  $\mathcal{S}$  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 [1], i.e., first-order sentences of the form

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

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  $\mathcal{U}$ .<sup>1</sup> The conjunction left of “ $\supset$ ” is the *body* of the constraint, and the disjunction right of “ $\supset$ ” its *head*.

In the rest of the paper,  $\mathcal{S} = \langle \Psi, \Sigma \rangle$  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. 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* [17], also called *generic constraints* in [7]. This class (semantically) includes:

- *key constraints*  $p(\vec{x}, \vec{y}) \wedge p(\vec{x}, \vec{z}) \supset y_i = z_i$ , for  $1 \leq i \leq n$ ,
- *functional dependencies*  $p(\vec{x}, \vec{y}, \vec{v}) \wedge p(\vec{x}, \vec{z}, \vec{w}) \supset y_i = z_i$ , for  $1 \leq i \leq n$ , and
- *exclusion dependencies*  $p_1(\vec{v}, \vec{y}) \wedge p_2(\vec{w}, \vec{z}) \supset y_1 \neq z_1 \vee \dots \vee y_n \neq z_n$ ,

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

where  $\vec{y} = y_1, \dots, y_n$  and  $\vec{z} = z_1, \dots, z_n$ .

- 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.,  $m \leq 1$  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  $\mathcal{S}_0$  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:

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

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  $\mathbf{C}_0$ .  $\square$

For a set of relation symbols  $\Psi$  as above,  $\mathcal{F}(\Psi)$  denotes the set of all facts  $r(t)$ , where  $r \in \Psi$  has arity  $n$  and  $t = (c_1, \dots, c_n) \in \mathcal{U}^n$  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  $r^D = \{t \mid r(t) \in D\}$ . We denote by  $D(\Psi)$  the set of all databases for  $\Psi$ . For any relation schema  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ , in abuse of notation,  $\mathcal{F}(\mathcal{S})$  and  $D(\mathcal{S})$  denote  $\mathcal{F}(\Psi)$  and  $D(\Psi)$ , respectively, and a database for  $\mathcal{S}$  is a database for  $\Psi$ .

A constraint  $\sigma$  is *ground*, if it is variable-free. For any such  $\sigma$ ,  $\text{facts}(\sigma)$  denotes the set of all facts  $p(t) \in \mathcal{F}(\mathcal{S})$  occurring in  $\sigma$ , and for any set  $\Sigma$  of ground constraints,  $\text{facts}(\Sigma) = \bigcup_{\sigma \in \Sigma} \text{facts}(\sigma)$ . For any constraint  $\sigma = \alpha(\vec{x})$ , we denote by  $\text{ground}(\sigma)$  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  $\Sigma$ ,  $\text{ground}(\Sigma) = \bigcup_{\sigma \in \Sigma} \text{ground}(\sigma)$ .

Given  $D \subseteq \mathcal{F}(\Psi)$ , where  $\Psi = \{r_1, \dots, r_n\}$ ,  $D$  *satisfies* a constraint  $\sigma$ , denoted  $D \models \sigma$ , if  $\sigma$  is true on the relational structure  $(\mathcal{U}, r_1^D, \dots, r_n^D, c_1^D, c_2^D, \dots)$  where  $c_i^D = c_i$ , for all  $c_i \in \mathcal{U}$  (i.e., each  $\sigma' \in \text{ground}(\sigma)$  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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$  is *consistent*, if there exists a database  $D$  for  $\mathcal{S}$  that is consistent with  $\Sigma$ , otherwise  $\mathcal{S}$  is inconsistent.

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

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

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

## 2.2 Datalog<sup>∨,¬</sup> Programs and Queries

**Syntax** A Datalog<sup>∨,¬</sup> rule  $\rho$  is an expression of the form

$$a_1 \vee \dots \vee a_n \leftarrow b_1, \dots, b_k, \text{ not } b_{k+1}, \dots, \text{ not } b_{k+m} \quad (2)$$

where  $a_i, b_j$  are atoms in a relational first-order language  $\mathcal{L}$ . Here, “not” is *negation as failure* and “ $\vee$ ” 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  $b_i$ ,  $1 \leq i \leq k$ , whose predicate is not a built-in relation. Built-in relations may occur only in the body.

A Datalog<sup>∨,¬</sup> program  $\mathcal{P}$  is a finite set of Datalog<sup>∨,¬</sup> rules. Important restrictions are *normal programs*, Datalog<sup>¬</sup>, where  $n = 1$  for all rules, *stratified normal programs*, Datalog<sup>¬s</sup>, and *non-recursive programs* as follows. Each Datalog<sup>¬</sup> 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  $r \rightarrow s$  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<sup>∨,¬</sup> 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  $b_i$ ,  $1 \leq i \leq k$ , 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* [28] assigns *stable models* to any Datalog<sup>∨,¬</sup> 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.  $I$ , which results from  $ground(\mathcal{P})$  by deleting (i) each rule  $\rho$  with a literal *not*  $p(\vec{c})$  in the body such that  $p(\vec{c}) \in I$ , and (ii) the negative literals from all remaining rules.

We denote by  $SM(\mathcal{P})$  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 [28].

**Queries** A Datalog<sup>∨,¬</sup> query  $Q$  over a schema  $\mathcal{S} = \langle \Psi, \Sigma \rangle$  is a pair  $\langle q, \mathcal{P} \rangle$ , where  $\mathcal{P}$  is a Datalog<sup>∨,¬</sup> 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  $\mathcal{S}$ , the *evaluation* of  $Q$  over  $D$ , is  $Q[D] = \{(c_1, \dots, c_n) \mid q(c_1, \dots, c_n) \in M, \text{ for each } M \in SM(\mathcal{P} \cup D)\}$ . 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<sup>∨,¬</sup> and queries, see [1, 22].

**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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$  is given together with a (possibly inconsistent) database  $D$  for  $\mathcal{S}$ . Following a common approach in the literature on inconsistent databases [3, 29, 13, 15], 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 preorder (i.e., a reflexive and transitive binary relation) on  $D(\mathcal{S})$ , 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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ , and let  $\leq_D$  be a preorder on  $D(\mathcal{S})$ . Then, a database  $R \in D(\mathcal{S})$  is a *repair for  $D$  w.r.t.  $\mathcal{S}$* , if

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

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

The definition of repair relies on a general notion of preorder on databases. The method for consistent query answering presented in the next sections 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  $D$ ,  $\Delta(R_1, D) \subset \Delta(R_2, D)$  implies  $R_1 <_D R_2$ ,

where  $\Delta(A, B) = (A \setminus B) \cup (B \setminus A)$  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  $\Delta(R_1, D) \subseteq \Delta(R_2, D)$  [3, 4, 6, 11, 17, 29, 27]. 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).*

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 [24, 13], cardinality-based ordering [4, 38], weight-based orderings [37], 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  $\mathcal{F}(S)$  on which  $R_1$  and  $R_2$  repair  $D$  differently belongs to  $R_2$ .

However, we point out that our method and results for query answering can also be extended to other preference orderings under certain conditions (see Section 8).

### 3.2 Constructible Repairs and Safe Constraints

An important aspect is that constraints might enforce that *any* set of facts  $R$  for  $\mathcal{S} = \langle \Psi, \Sigma \rangle$  which satisfies  $\Sigma$  must be infinite, and thus  $\mathcal{S}$  is inconsistent, i.e., no  $D \in D(\mathcal{S})$  satisfies  $\Sigma$ . A simple example is where  $\Sigma = \{\forall x p(x)\}$ . Semantically, this is commonly avoided by requesting domain-independence of constraints [42], 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 (proofs of the propositions below are given in Appendix A). For any  $R \subseteq \mathcal{F}(S)$ , we denote by  $\text{adom}(R, S)$  the *active domain* of  $R$  and  $S$ , i.e. the set of constants occurring in  $R$  and  $\Sigma$ .

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

Notice that, for a generic preference order, existence of a repair is not always guaranteed, even if  $\mathcal{S}$  is consistent.

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.  $\mathcal{S} = \langle \{p\}, \{p(x) \vee x > 100\} \rangle$ , assuming that  $\mathcal{U}$  are the natural numbers. As this example shows, repairs may in this case go beyond the active domain. However, this is prevented if built-ins involve only equality and inequality. We have here a result similar to Proposition 3.2.

**Proposition 3.3** *Let  $D$  be a database for  $\mathcal{S} = \langle \Psi, \Sigma \rangle$  where no built-in relations occur in  $\Sigma$  except  $=$  and  $\neq$ . Suppose that  $<_D$  satisfies (SIP). Then, every repair  $R \in \text{rep}(D)$  involves only constants from  $\text{adom}(D, S)$ , and some repair exists if  $\mathcal{S}$  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<sup>-</sup> query. For any database  $D \in D(\mathcal{S})$ , the set of consistent answers to  $Q$  w.r.t.  $D$  is the set of tuples  $\text{ans}(Q, D) = \{t \mid t \in Q[R], \text{ for each } R \in \text{rep}(D)\}$ .*

Informally, a tuple  $t$  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<sup>-</sup> is often adopted.

**Example 3.1** Recall that in our scenario, repairs for the database  $D_0$  for  $\mathcal{S}_0$  are shown in Fig. 1. For the query  $Q = \langle q, \mathcal{P} \rangle$ , where  $\mathcal{P} = \{q(x) \leftarrow \text{player}(x, y, z), q(x) \leftarrow \text{team}(v, w, x)\}$ , we thus obtain  $\text{ans}(Q, D_0) = \{(8), (9), (10)\}$ . For the query  $Q' = \langle q, \{q(y) \leftarrow \text{team}(x, y, z)\} \rangle$ , we have  $\text{ans}(Q', D_0) = \{\text{Barcelona}\}$ , while for  $Q'' = \langle q', \{q'(x, z) \leftarrow \text{team}(x, y, z)\} \rangle$ , we have  $\text{ans}(Q'', D_0) = \{(RM, 10), (BC, 8)\}$ .  $\square$

## 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, given  $\mathcal{S}$  and  $D$ ,  $E$  has to fulfill the conditions

$$\Delta(R, D) \subseteq E, \text{ for all } R \in \text{rep}_{\mathcal{S}}(D), \quad (3)$$

$$\text{rep}_{\mathcal{S}}(D) = \{R \cup (D \setminus E) \mid R \in \text{rep}_{\mathcal{S}}(D \cap E)\}. \quad (4)$$

The repair of  $D$  can then be fully localized to the repair of  $D \cap E$ , which in practice may be much smaller than  $D$ . In fact, as shown below, for constraints  $\mathbf{C}_0$  the set of all facts witnessing inconsistency, denoted  $C$ , is always a repair envelope, and for constraints  $\mathbf{C}_1$  and  $\mathbf{C}_2$ , a closure  $C^*$  of  $C$  under syntactic conflict propagation 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. Figure 2 shows the different sets.

**Example 4.1** Recall that  $\text{team}(RM, Roma, 10) \wedge \text{team}(RM, Real Madrid, 10) \supset Roma = Real Madrid$  witnesses in Example 2.2 a violation of the key of *team*; it is the only ground constraint violated by  $D_0$ . Since the constraints are of type  $\mathbf{C}_0$ , the set  $C = \{\text{team}(RM, Roma, 10), \text{team}(RM, Real Madrid, 10)\}$  is a repair envelope for  $D$ . The database  $D \cap C = C$  has the two repairs  $R_1 = \{\text{team}(RM, Roma, 10)\}$  and  $R_2 = \{\text{team}(RM, Real Madrid, 10)\}$ ; 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.  $\square$

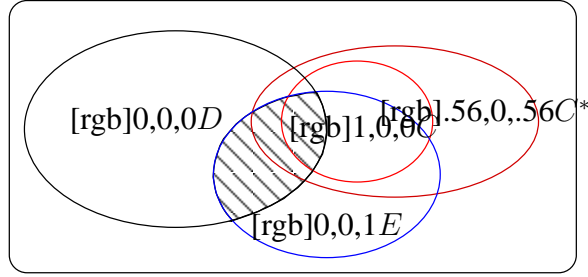
Note that a repair envelope always exists, since the set of all facts is a trivial repair envelope. As for localizing the computation of  $\text{rep}(D)$ , only condition (4) is relevant (if  $E$  satisfies it, then so does every  $E'$  such that  $E' \cap D = E \cap D$ , in particular  $E' = E \cap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

For general constraints,  $C^*$  is not always a repair envelope. However, we show that it is a *weak repair envelope*  $E$ , which has to fulfill, instead of (4), the relaxed equation

$$\text{rep}_{\mathcal{S}}(D) = \{(R \cap E) \cup (D \setminus E) \mid R \in \text{rep}_{\mathcal{S}}(D \cap E)\} \quad (5)$$

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 4.3). 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 call the set  $D \cap E$  the affected part of  $D$  (w.r.t.  $\mathcal{S}$ ), or simply affected database, and we call  $D \setminus E$  the safe part of  $D$  (w.r.t.  $\mathcal{S}$ ), or simply “safe” database.



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)

Figure 2: Localization of database repair

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  $C_1$  and  $C_2$  constraints. This envelope may be further decreased. Indeed, we prove that  $C$  is a repair envelope for  $C_0$  constraints. In fact, the results for special constraints are stronger and establish 1-1 correspondences between repairs of  $D \cap E$  and repairs of  $D$ .

#### 4.1 General Constraints

Let  $D$  be a database for a relational schema  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ . The *conflict set* for  $D$  w.r.t.  $\mathcal{S}$  is the set of facts  $C_{\mathcal{S}}(D) = \{p(t) \mid \exists \sigma \in ground(\Sigma), p(t) \in facts(\sigma), D \not\models \sigma\}$ , i.e.,  $C_{\mathcal{S}}(D)$  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  $\mathcal{S}$  will be dropped.

Figure 2 shows that the conflict set may contain both facts in  $D$  (as in Example 4.1) and facts in  $\mathcal{F}(\mathcal{S})$  that do not belong to  $D$ . For example, let  $D = \{p(a)\}$ , and let  $\mathcal{S}$  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  $p(t), p'(t')$  in  $\mathcal{F}(\mathcal{S})$  are *constraint-bounded in  $\mathcal{S}$* , if there exists some  $\sigma \in ground(\Sigma)$  such that all constants occurring in  $facts(\sigma)$  are from  $adom(D, \mathcal{S})$ , and  $\{p(t), p'(t')\} \subseteq facts(\sigma)$ . (Note that by assumed safety of constraints and the results of Section 3.2, we only need to consider  $adom(D, \mathcal{S})$ .) We now generalize the notion of conflict set.

**Definition 4.1 (Conflict closure)** Let  $D$  be a database for  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ . Then, the *conflict closure* for  $D$ , denoted by  $C_{\mathcal{S}}^*(D)$ , is the least set  $A \supseteq C_{\mathcal{S}}(D)$  which contains every fact  $p(t)$  constraint-bounded in  $\mathcal{S}$  with some fact  $p'(t') \in A$ .  $\square$

We omit  $D$  and/or the subscript  $\mathcal{S}$  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  $\mathcal{S}$  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  $D$ . In the example above, for instance,  $C$  and  $C^*$  would be the same 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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ , consider the following two sets of ground constraints:

- (i)  $\Sigma^a_{\mathcal{S}}(D) = \{\sigma \in \text{ground}(\Sigma) \mid \text{facts}(\sigma) \cap C^* \neq \emptyset\}$  consists of all ground constraints in which at least one fact from  $C^*$  occurs;
- (ii)  $\Sigma^s_{\mathcal{S}}(D) = \{\sigma \in \text{ground}(\Sigma) \mid \text{facts}(\sigma) \not\subseteq C^*\}$  consists of all ground constraints in which at least one fact occurs which is *not* in  $C^*$ .

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

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

*Proof.* By definition,  $\sigma \in \Sigma^a$  contains at least one fact  $p(t)$  from  $C^*$ ; any other fact in  $\sigma$  is constraint-bounded in  $\mathcal{S}$  with  $p(t)$ , and hence it also must be in  $C^*$ . This proves  $\text{facts}(\Sigma^a) \subseteq C^*$ . Consider now any fact  $p(t) \in C^*$ . The minimality of  $C^*$  implies that there exist facts  $f_1, \dots, f_n$  in  $C^*$  such that  $f_1 \in C$ ,  $f_n = p(t)$ , and  $f_{i+1}$  is constraint-bounded to  $f_i$ , for each  $i \in \{1, \dots, n-1\}$ ; i.e.,  $f_i, f_{i+1} \in \text{facts}(\sigma_i)$  for some  $\sigma_i \in \text{ground}(\Sigma)$ . Each  $\sigma_i$  then belongs to  $\Sigma^a$ , and thus  $p(t) \in \text{facts}(\Sigma^a)$ . This proves  $C^* \subseteq \text{facts}(\Sigma^a)$ , and therefore (1) holds. As for (2), assume by contradiction that some  $\sigma \in \Sigma^s$  with  $\text{facts}(\sigma) \cap C^* \neq \emptyset$  exists. Then, from Definition 4.1 it follows that  $\text{facts}(\sigma) \subseteq C^*$ , 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  $\sigma \in \text{ground}(\Sigma)$  such that  $\sigma \notin \Sigma^s$  and  $\sigma \notin \Sigma^a$ , but this means that  $\text{facts}(\sigma) \cap C^* = \emptyset$  and  $\text{facts}(\sigma) \not\subseteq C^*$ , which is an obvious contradiction.  $\square$

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

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

*Proof.* Towards a contradiction, suppose that there exists a repair  $R \in \text{rep}(D)$  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  $\sigma \in \text{ground}(\Sigma)$ . 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),  $\text{facts}(\sigma) \subseteq C^*$ , 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  $R \in \text{rep}(D)$ ). In case (ii), again  $R' \models \sigma$ : by Proposition 4.1 (b),  $\text{facts}(\sigma) \cap C^* = \emptyset$  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  $\Delta(R', D) \subset \Delta(R, D)$ . 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  $R \in \text{rep}(D)$ .  $\square$

Informally, the above lemma 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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ , and let  $A = D \cap C^*$  and  $\mathcal{S}^a = \langle \Psi, \Sigma^a \rangle$ . Then, for each  $S \subseteq D \setminus C^*$ , the following holds:*

1. for each  $R \in \text{rep}_{\mathcal{S}}(A \cup S)$ , it holds that  $R \cap C^* \in \text{rep}_{\mathcal{S}^a}(A)$ ;



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

*Proof.* (1) Let  $R \in \text{rep}_{\mathcal{S}}(A \cup S)$ , and let  $R' = R \cap C^*$ . Since  $R \models \Sigma$ , then  $R \models \text{ground}(\Sigma)$ , therefore, from Proposition 4.1 it follows that  $R' \models \Sigma^a$ , while  $R \setminus R' = R \setminus C^* \models \Sigma^s$ . Assume  $R' \notin \text{rep}_{\mathcal{S}^a}(A)$ . Since  $R' \models \Sigma^a$ , there must exist some  $R'' \in \text{rep}_{\mathcal{S}^a}(A)$  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.  $\mathcal{S}$  is the same as w.r.t.  $\mathcal{S}^a$  and since  $A$  is contained in its conflict closure w.r.t.  $\mathcal{S}$ , 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  $R \in \text{rep}_{\mathcal{S}}(A \cup S)$ .

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

Since  $R \in \text{rep}_{\mathcal{S}^a}(A)$ , from Proposition 4.2, and from the fact that  $C_{\mathcal{S}^a}^*(A) = C^*$  (in computing  $C_{\mathcal{S}^a}^*(A)$  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  $R \cup S' \notin \text{rep}_{\mathcal{S}}(A \cup S)$ , 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  $R \in \text{rep}_{\mathcal{S}^a}(A)$  and  $S' \in \text{rep}_{\mathcal{S}^s}(S)$ . This proves that  $R \cup S' \in \text{rep}_{\mathcal{S}}(A \cup S)$ .  $\square$

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.  $\mathcal{S}^a$ , from a repair, computed w.r.t.  $\mathcal{S}$ , 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.  $\mathcal{S}$ , from a repair of  $A$  w.r.t.  $\mathcal{S}^a$ . Notice that repairing  $A$  w.r.t.  $\mathcal{S}^a$ , and not w.r.t.  $\mathcal{S}$ , is necessary for the lemma above to hold, since for a repair  $R \in \text{rep}_{\mathcal{S}}(A \cup S)$ , it does not hold in general that  $R \cap C^* \in \text{rep}_{\mathcal{S}^a}(A)$ . Also, repairing  $A$  w.r.t.  $\mathcal{S}^a$  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** *Let  $D$  be a database for  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ . Then,  $C^*$  is a weak repair envelope for  $D$ .*

*Proof.* We first show that for each  $R \in \text{rep}_{\mathcal{S}}(D)$  then  $\Delta(R, D) \subseteq C^*$ , as specified by condition (3) in Section 4, where we pose  $E = C^*$ . Assume by contradiction that there exists a fact  $f \in \Delta(R, D)$  such that  $f \notin C^*$ . By Proposition 4.1 it follows that there does not exist any  $\sigma \in \Sigma^a$  such that  $f \in \text{facts}(\sigma)$ . 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  $R \in \text{rep}_{\mathcal{S}}(D)$ . 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  $\text{rep}_{\mathcal{S}}(D)$  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  $R \in \text{rep}_{\mathcal{S}}(D)$ , there exists some  $R' \in \text{rep}_{\mathcal{S}}(D \cap C^*)$  such that  $R = (R' \cap C^*) \cup (D \setminus C^*)$ , and (ii) for every  $R \in \text{rep}_{\mathcal{S}}(D \cap C^*)$  there exists some  $R' \in \text{rep}_{\mathcal{S}}(D)$  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  $R \cap C^* \in \text{rep}_{\mathcal{S}^a}(D \cap C^*)$  (notice that in Lemma 4.3  $A = D \cap C^*$ ,  $\mathcal{S}^a = \langle \Psi, \Sigma^a \rangle$ , 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  $R' = (R \cap C^*) \cup S' \in \text{rep}_{\mathcal{S}}(A)$ . Since  $S' \cap C^* = \emptyset$ , 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  $R' = (R \cap C^*) \cup S' \in \text{rep}_{\mathcal{S}}(D)$ . 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  $C^*$ ;
2. compute the repairs of  $A = D \cap C^*$ ;
3. intersect each repair obtained with  $C^*$ ; and
4. for each such set, take the union with  $D \setminus C^*$ .

A drawback of this approach 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 4.2** Consider  $D = \{p(a)\}$  for  $\mathcal{S} = \langle \Psi, \{p(a), q(a)\} \rangle$ . In this case,  $C = C^* = \{q(a)\}$ ,  $A = D \cap C^* = \emptyset$ , and  $D \setminus C^* = D$ . We have  $\text{rep}(A) = \{\{p(a), q(a)\}\}$  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 4.3** Consider  $D = \{r(a), p(a)\}$  where for  $\mathcal{S} = \langle \Psi, \Sigma \rangle$  where  $\Sigma = \{r(a) \supset p(a) \vee q(a), r(a)\}$ . Then,  $C^* = \emptyset$  and  $D \cap C^* = \emptyset$  has two repairs, viz.  $R_1 = \{r(a), p(a)\}$  and  $R_2 = \{r(a), q(a)\}$ . According to (4),  $R_2 \cup (D \setminus C^*) = \{r(a), p(a), q(a)\}$  would have to be a repair of  $D$ , which is incorrect. Note that the constraint  $r(a) \supset p(a) \vee q(a)$  can be satisfied by including  $q(a)$ , which was neither in  $D$  nor in  $C^*$ .  $\square$

Note that in Example 4.3,  $\Sigma$  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 4.3 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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ , and let  $\mathcal{S}' = \langle \Psi, \Sigma' \rangle$  be such that  $\Sigma_{\mathcal{S}'}^a(D) = \Sigma_{\mathcal{S}}^a(D)$ . Then  $\text{rep}_{\mathcal{S}}(D) = \text{rep}_{\mathcal{S}'}(D)$ .*

*Proof.* We prove that  $\text{rep}_{\mathcal{S}}(D) \subseteq \text{rep}_{\mathcal{S}'}(D)$ . The converse can be proved analogously. Assume by contradiction that there exists  $R \in \text{rep}_{\mathcal{S}}(D)$  such that  $R \notin \text{rep}_{\mathcal{S}'}(D)$ . 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  $R \not\models \Sigma_{\mathcal{S}'}^a(D) \cup \Sigma_{\mathcal{S}'}^s(D)$  and since  $R \models \Sigma_{\mathcal{S}}^a(D)$  (which is equal to  $\Sigma_{\mathcal{S}'}^a(D)$ ), there exists some  $\sigma \in \Sigma_{\mathcal{S}'}^s(D)$  such that  $R \not\models \sigma$ . By Theorem 4.4,  $R = (R'' \cap C_{\mathcal{S}}^*(D)) \cup (D \setminus C_{\mathcal{S}}^*(D))$ , where  $R'' \in \text{rep}_{\mathcal{S}}(D \cap C_{\mathcal{S}}^*(D))$ . Since  $\text{facts}(\Sigma_{\mathcal{S}}^a(D)) = \text{facts}(\Sigma_{\mathcal{S}'}^a(D))$ , by Proposition 4.1  $C_{\mathcal{S}}^*(D) = C_{\mathcal{S}'}^*(D) = C^*$ . Furthermore, since  $\sigma \in \Sigma_{\mathcal{S}'}^s(D)$ , we have  $\text{facts}(\sigma) \cap C^* = \emptyset$ . 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  $R \in \text{rep}_{\mathcal{S}}(D)$ .  $\square$

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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$  such that  $\Sigma \subseteq \mathbf{C}_1$ . Then, each repair  $R$  of  $A = D \cap C^*$  w.r.t.  $\mathcal{S}$  satisfies  $R \subseteq C^*$ .*

*Proof.* By Item 1 of Lemma 4.3, for  $S = \emptyset$ , each  $R \in \text{rep}_{\mathcal{S}}(A)$  gives rise to a repair  $R' = R \cap C^*$  of  $A$  w.r.t.  $\mathcal{S}^a = \langle \Psi, \Sigma^a \rangle$ . By Item 2 of Lemma 4.3, for  $S = \emptyset$ ,  $R'$  in turn gives rise to a repair  $R''$  of  $A$  w.r.t.  $\mathcal{S}$  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.  $\mathcal{S}$ . Now if  $R \not\subseteq C^*$  would hold, then  $\Delta(R'', A) \subset \Delta(R, A)$  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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ , 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  $R \in \text{rep}(A)$  gives rise to some  $R'' \in \text{rep}(A)$  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.  $\langle \Psi, \Sigma^s \rangle$ . Since 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  $R \not\subseteq D \cup C^*$  would hold, then  $\Delta(R'', A) \subset \Delta(R, A)$  would hold (note that  $\mathcal{F} \subseteq R$  must hold, and thus  $R'' \subseteq R$ ), which by (SIP) means  $R'' <_D R$ . This is a contradiction, and proves (i). Item (ii) holds since  $R \cap (D \setminus C^*) = \mathcal{F}$  for each  $R \in \text{rep}(A)$ .  $\square$

The propositions 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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$  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  $\mu : \text{rep}(D) \rightarrow \text{rep}(D \cap C^*)$ , such that for every  $R \in \text{rep}(D)$ ,  $R = \mu(R) \cup (D \setminus C^*)$ .*

*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  $R \in \text{rep}(D \cap C^*)$ , when  $\Sigma \subseteq \mathbf{C}_1$ ,  $R \cap C^* = R$ , whereas when  $\Sigma \subseteq \mathbf{C}_2$ ,  $(R \cap C^*) \cup (D \setminus C^*) = R \cup (D \setminus C^*)$ .  $\square$

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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ , if we can modify  $\Sigma$  to constraints  $\Sigma'$  from  $\mathbf{C}_1$  or  $\mathbf{C}_2$  while preserving the affected constraints, i.e.,  $\Sigma_{\mathcal{S}}^a(D) = \Sigma_{\langle \Psi, \Sigma' \rangle}^a(D)$ . 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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ ,  $\Sigma \subseteq \mathbf{C}_0$ , and let  $A = D \cap C^*$ . Then,*

1.  $C \subseteq D$ ;
2. for each  $R \in \text{rep}(A)$ , (i)  $R \subseteq A$ , (ii)  $\Delta(R, A) \subseteq C$ , (iii)  $A \setminus C \subseteq R$ , and (iv)  $R \cap C \in \text{rep}(C)$ ;
3. for each  $R \in \text{rep}(C)$ ,  $R \cup (A \setminus C) \in \text{rep}(A)$ .

*Proof.* 1) By definition,  $C$  is the set of facts occurring in any constraint  $\sigma \in \text{ground}(\Sigma)$  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  $R \in \text{rep}(A)$ . 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  $\Delta(R, A) \not\subseteq C$ . 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  $\text{ground}(\Sigma)$  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  $R \notin \text{rep}(A)$ ; 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  $R \cap C \notin \text{rep}(C)$ . Then, it is easy to see that  $R \cap C \models \Sigma$ , therefore some  $R' \in \text{rep}(C)$  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  $R \in \text{rep}(A)$ .

3) Let  $R' \in \text{rep}(C)$ . 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  $\sigma \in \text{ground}(\Sigma)$ . Then  $\sigma$  must contain a fact

$p(t)$  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  $p(t) \in C$ , a contradiction. From Item 2 and (DIS) (split by  $C$ ), we can conclude that no  $R'' \in \text{rep}(A)$  exists such that  $R'' <_A R$ . Hence,  $R \in \text{rep}(A)$ .  $\square$

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)). 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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$  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  $\nu : \text{rep}(D) \rightarrow \text{rep}(D \cap E)$ , such that for each  $R \in \text{rep}(D)$ ,  $R = \nu(R) \cup (D \setminus E)$ .*

*Proof.* By Corollary 4.8, there is a bijection  $\mu : \text{rep}(D) \rightarrow \text{rep}(A)$ , where  $A = D \cap C^*$ , such that the repairs of  $D$  are given by  $\mu(R) \cup (D \setminus C^*)$ , for all  $R \in \text{rep}(A)$ . 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  $R \in \text{rep}(C)$ . Hence, the mapping  $\nu : \text{rep}(D) \rightarrow \text{rep}(C)$  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  $\nu' : \text{rep}(D \cap E) \rightarrow \text{rep}(C)$  such that  $R = \nu'(R) \cup (D' \setminus C)$ , for each  $R \in \text{rep}(D')$ . This implies a bijection  $\nu'' : \text{rep}(D) \rightarrow \text{rep}(D \cap E)$  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. Notice that because  $C \subseteq D$ , we can compute  $C$  efficiently by suitable SQL statements which express constraint violations. The fact that every  $E \supseteq C$  is a repair envelope gives convenient flexibility to modify the statements in case (allow 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 answers.

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. We now address the issue of efficient recombination.

## 5.1 Recombination Step

Let us now consider the problem of evaluating a query  $Q$  issued over an inconsistent database  $D$  for  $\mathcal{S}$ , i.e., to compute  $ans(Q, D)$ . Recall that according to the definition in Section 2, a tuple  $\vec{t}$  belongs to  $ans(Q, D)$  if  $\vec{t}$  is in the answer to  $Q$  on every repair of  $D$ , i.e.,  $ans(Q, D) = \{\vec{t} \mid \vec{t} \in Q[R] \text{ for each } R \in rep(D)\} = \bigcap_{R \in rep(D)} Q[R]$ . 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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ . Let  $E$  be a set of facts, and let  $A = D \cap E$  and  $S = D \setminus E$ . Then*

$$ans(Q, D) = \bigcap_{R \in rep(A)} Q[\chi(R) \cup S], \quad (6)$$

where (i)  $\chi(R) = R$  if  $E$  is a repair envelope for  $D$ , and (ii)  $\chi(R) = R \cap E$  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  $ans(Q, D) = \bigcap_{R \in rep(C)} Q[R \cup S]$ .

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 (co-NP) computation 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. Actually, this is the best one may asymptotically expect to achieve for general inconsistent databases, unless  $P = NP$ , given that consistent query answering is co-NP-hard.

Hence, it is particularly relevant to assess whether some smarter strategies can be conceived for special classes of queries and constraints, in order to have an algorithm that both implements localized repair computation and linearly scales w.r.t. the size of the database.

## 5.2 Repair Factorization

In this section, we present a technique that factorizes repairs into independent components (proofs of theorems and propositions are given in Appendix B). 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  $A_1, \dots, A_m$ , such that the repairs of  $A$

are obtained by combining the repairs of  $A_1, \dots, A_m$  in all possible ways. Given a repair envelope  $E$  for  $D$  and  $\mathcal{S}$ , a partitioning  $E_1, \dots, E_m$  of  $E$  is a *factorization* of  $E$  for  $D$  and  $\mathcal{S}$ , if

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

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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ . A partitioning  $E_1, \dots, E_m$  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  $R \in \text{rep}(D \cap E)$  and  $R_i \in \text{rep}(D \cap E_i)$ ,  $1 \leq i \leq m$ ,  $R \setminus E = R_i \setminus E_i$ .*

By means of a repair-compliant partitioning, we can factorize the repair of  $A$  into the repair of the (mutually disjoint) parts  $A_i = A \cap E_i = D \cap E_i$  of  $A$ , for  $i = 1, \dots, m$ . 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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ , and let  $E$  be a repair envelope for  $D$ . Then, every repair-compliant partitioning  $E_1, \dots, E_m$  of  $E$  is a factorization of  $E$  for  $D$  and  $\mathcal{S}$ .*

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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ , 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  $E_1, \dots, E_m$  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. Consequently, for  $\mathbf{C}_0$  and the canonical envelope  $E = C$ , Equation (7) can be rewritten to:

$$\text{rep}(D) = \{(D \setminus C) \cup R_1 \cup \dots \cup R_m \mid R_i \in \text{rep}(C_i), 1 \leq i \leq m\}. \quad (8)$$

**Example 5.1** Let  $\mathcal{S}$  consist of the relation  $p(x, y, z)$  and the functional dependency  $f : p(x, y, z) \wedge p(x, y', z') \supset z = z'$ , and consider the database  $D = \{p(a_i, b_j, c_k) \mid 1 \leq i \leq m \wedge 1 \leq j, k \leq \ell\}$ . 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  $k \neq k'$  witnesses a violation of  $f$ . The partitioning  $C_1, \dots, C_m$  of  $C$ , where  $C_i = \{p(a_i, b, c) \in C\}$ ,  $1 \leq i \leq m$ , is constraint-bounded and thus, by Proposition 5.3, repair-compliant. Every  $C_i$  has  $\ell$  repairs, while  $D$  has  $\ell^m$  repairs in total. In particular, the repairs of  $D$  are of the form  $R_1 \cup \dots \cup R_m$ , where each  $R_i$  is a repair for  $C_i$ , according to Equation (8).  $\square$

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  $\leq_D$ , 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  $1 \leq i, j \leq m$ ,  $R_i \in \text{rep}(D \cap E_i)$  and  $R_j \in \text{rep}(D \cap E_j)$  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 facts which are constraint-bounded. In this respect, we point out that techniques exploiting graph (and hypergraph) representations of conflicts in data have also been introduced and used in [2, 16].

### 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  $E_1, \dots, E_m$ , although they are inconsistent, need not be repaired to answer  $Q$ . Intuitively, this happens when all the repairs of a given component  $E_i$  are indistinguishable as far as answering  $Q$  is concerned.
- Second, we investigate how to improve on the naive exploitation of equation (7), by discussing scenarios where the answer to  $Q$  can be obtained by independently processing the different components, rather than combining their repairs in all the 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.,  $\mathcal{P} = \{\rho_1, \dots, \rho_n\}$ , where  $head(\rho_j) = q(\vec{t}_j)$ , and each predicate in  $body(\rho_j)$ ,  $1 \leq j \leq n$ , is from the schema. We denote by  $n(Q) = n$  the number of rules in  $\mathcal{P}$  and by  $a(Q)$  the maximal number of variables appearing in any  $\rho_j$ .

**Example 5.2** 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, and thus by Theorem 5.2 is a factorization. For the query  $Q = \langle q, \{q(x) \leftarrow r(x, y), s(y, z)\} \rangle$ , we have  $n(Q) = 1$  and  $a(Q) = 3$ .  $\square$

We now formalize scenarios where inconsistencies in some components can be tolerated.

**Definition 5.2** Let  $E_1, \dots, E_m \subseteq E$  be a factorization of a repair envelope  $E$  for a database  $D$ , and let  $Q$  be a non-recursive (unfolded) Datalog query as above. Let  $c$  be a new constant symbol not belonging to the universal database domain  $\mathcal{U}$ . For each component  $E_i$ , we define  $Q_i = \langle q_i, \mathcal{P}_i \rangle$ ,  $\mathcal{P}_i = \{q_i(j, \vec{t}'_j) \leftarrow atoms(E_i, j) \mid 1 \leq j \leq n(Q)\}$ , where  $atoms(E_i, j)$  is the set of atoms  $p(\vec{x}) \in body(\rho_j)$  such that  $E_i$  contains a fact over  $p$ , and  $\vec{t}'_j$  is obtained as follows: substitute  $c$  in  $\vec{t}_j$  for each variable not appearing in  $atoms(E_i, j)$ ; append all variables from  $atoms(E_i, j)$  occurring in  $body(\rho_j) \setminus atoms(E_i, j)$ ; and pad the resulting list to length  $a(Q)$ , using  $c$ . Finally, let  $ans(Q_i, j) = \{\{(j, \vec{s}) \in Q_i[(D \setminus E) \cup R]\} \mid R \in rep(D \cap E_i)\}$ .  $\square$

**Example 5.3** In Example 5.2, we have  $atoms(C_{r_1}, 1) = atoms(C_{r_2}, 1) = \{r(x, y)\}$  and  $atoms(C_s, 1) = \{s(y, z)\}$ . We thus have  $Q_{r_1} = \langle q_{r_1}, \{q_{r_1}(1, x, y, c) \leftarrow r(x, y)\} \rangle$  and  $Q_{r_2} = \langle q_{r_2}, \{q_{r_2}(1, x, y, c) \leftarrow r(x, y)\} \rangle$ , while  $Q_s = \langle q_s, \{q_s(1, c, y, c) \leftarrow s(y, z)\} \rangle$ .  $C_{r_1}$  has the two repairs  $\{r(a_1, b_1)\}$  and  $\{r(a_1, b_2)\}$ , and  $ans(Q_{r_1}, 1) = \{\{(1, a_1, b_1, c)\}, \{(1, a_1, b_2, c)\}\}$ . Similarly,  $C_{r_2}$  has two repairs,  $\{r(a_2, b_1)\}$  and  $\{r(a_2, b_3)\}$ , and  $ans(Q_{r_2}, 1) = \{\{(1, a_2, b_1, c)\}, \{(1, a_2, b_3, c)\}\}$ . Finally,  $C_s$  has the two repairs  $\{s(b_1, c_1)\}$  and  $\{s(b_1, c_2)\}$ , but  $ans(Q_s, 1)$  is the singleton  $\{\{(1, c, b_1, c)\}\}$ .  $\square$



Notice that a component  $E_i$  such that  $|ans(Q_i, j)| = 1$  for  $1 \leq j \leq n(Q)$ , while possibly inconsistent, does not need to be actually repaired as far as answering  $Q$  is concerned. In the following, every such component  $E_i$  is called *singular*.

Guided by this observation, given a factorization  $E_1, \dots, E_m$  of a repair envelope  $E$  for  $D$  where  $E_1, \dots, E_\ell$  are the singular components, we call  $S_Q = (D \setminus E) \cup R_1 \cup \dots \cup R_\ell$  a *query-safe part* of  $D$  w.r.t.  $Q$ , if each  $R_i$ ,  $1 \leq i \leq \ell$ , is an arbitrary repair of  $D \cap E_i$ . In Example 5.2,  $C_s$  is the only singular component, and  $R = \{s(b_1, c_2)\}$  is a repair of  $C_s = D \cap C_s$ ; hence,  $S_Q = \{r(a_3, b_1), s(b_3, c_3), s(b_1, c_1), s(b_1, c_2)\}$  is a query-safe part of  $D$ .

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

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

If all components  $E_i$  are singular, 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 that  $|ans(Q_i, j)| = 1$  for all  $1 \leq i \leq m$  and  $1 \leq j \leq n(Q)$ , which can be efficiently carried out by processing the components independently of each other.

Interestingly, even if non-singular components are present, ‘‘parallelizing’’ query answering without a need for recombination may be possible.

**Definition 5.3** A factorization  $E_1, \dots, E_m$  of a repair envelope  $E$  for  $D$  is *decomposable* w.r.t. query  $Q$ , if its non-singular components  $E_\ell, \dots, E_m$  satisfy

- (1)  $atoms(E_i, k) = atoms(E_j, k)$ , for every  $\ell \leq i, j \leq m$  and  $1 \leq k \leq n(Q)$ ; and,
- (2)  $|atoms(E_i, k)| = 1$ , for every  $\ell \leq i \leq m$  and  $1 \leq k \leq n(Q)$ .
- (3)  $R_i \setminus E_i = R_j \setminus E_j$ , for every  $R_i \in rep(D \cap E_i)$ ,  $R_j \in rep(D \cap E_j)$ ,  $\ell \leq i, j \leq m$ .  $\square$

**Proposition 5.5** *Let  $E_1, \dots, E_m$  be a factorization of a repair envelope  $E$  for  $D$  which is decomposable w.r.t. query  $Q$  having the non-singular components  $E_\ell, \dots, E_m$ . Then*

$$ans(Q, D) = \bigcup_{i=\ell}^m \left( \bigcap_{R_i \in rep(D \cap E_i)} Q[S_Q \cup R_i] \right) \quad (10)$$

**Example 5.4** In our example, the non-singular components are  $C_{r_1}$  and  $C_{r_2}$ . Since  $atoms(C_{r_1}, 1) = atoms(C_{r_2}, 1) = \{r(x, y)\}$  ( $n(Q) = 1$ ), the factorization  $C_{r_1}, C_{r_2}, C_s$  is decomposable w.r.t.  $Q$ . By Proposition 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 the query-safe part  $S_Q = \{r(a_3, b_1), s(b_3, c_3), s(b_1, c_1), s(b_1, c_2)\}$  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,  $ans(Q, D) = \{(a_3), (a_2)\}$ . As can be checked, this is the correct result.  $\square$

In closing this section, we note that by virtue of Propositions 5.4 and 5.5 one can efficiently answer a number of queries which do not fall within any of the tractable classes proposed in the literature so far [18, 17, 26, 30]. Examples of such queries and a discussion on the performances of this strategy are given in Section 7.

## 6 Logic Programming for Consistent Query Answering

According to several proposals in the literature, consistent answers over 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 [29, 4, 6]) is introduced next.<sup>2</sup>

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

1.  $rep_{\mathcal{S}}(D) \rightleftharpoons \text{SM}(\Pi_{\Sigma} \cup D)$ , and
2.  $ans(Q, D) = Q'[D]$ , where  $Q' = \langle q, \Pi_{\mathcal{S}}(Q) \rangle$ , i.e.,  $ans(Q, D) = \{t \mid q(t) \in M \text{ for each } M \in \text{SM}((\Pi_{\Sigma} \cup \Pi_Q) \cup D)\}$ , where  $\Pi_Q$  is a non-recursive safe Datalog<sup>¬</sup> program,

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

In the above definition,  $\Pi_{\Sigma}$  is that portion of  $\Pi_{\mathcal{S}}(Q)$  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 D).

Encoding repair computation by means of logic programs has some attractive features. An important one is that Datalog<sup>∨,¬</sup> 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<sup>∨,¬</sup> language that allow, for instance, to handle priorities and weight constraints [36, 40], 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 [4]).

However, with current (yet still improving) implementations of stable model engines, such as DLV [36] or Smodels [40], query evaluation over large data sets quickly becomes infeasible because of lacking scalability. The source of complexity in evaluating the program  $\Pi_{\mathcal{S}}(Q)$  lies in the conflict resolution module  $\Pi_{\Sigma}$ . Indeed, while  $\Pi_Q$ , which is in general a non-recursive Datalog<sup>¬</sup> program, can be evaluated in polynomial time with respect to underlying databases (data complexity) [19],  $\Pi_{\Sigma}$  is in general a Datalog<sup>∨,¬</sup> program [29], whose evaluation data complexity is at the second level of the polynomial hierarchy [19].

### 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  $\Pi_{\mathcal{S}}(Q)$  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 time is interfaced and 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:

<sup>2</sup>Other logic formalizations proposed in the data integration setting [33, 8, 14, 11] also fit in our framework, provided that the *retrieved global database* is computed [34]. Notice also that other logic-based approaches to data integration, based on abductive logic programming [5] and ID-logic [39], do not fit this framework.

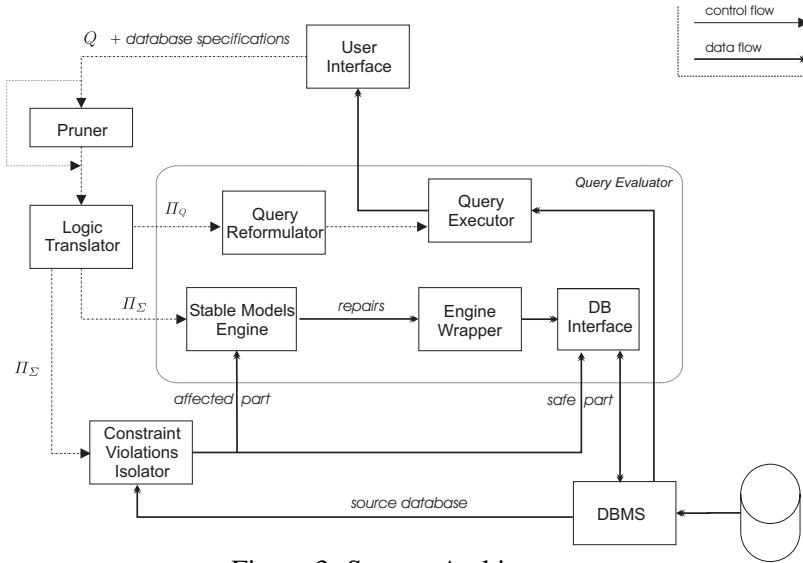


Figure 3: System Architecture.

- *Pruner*: It takes the user query  $Q$  and the schema  $\mathcal{S}$ , and produces an equivalent specification (w.r.t.  $Q$ ), stripping off relations and constraints irrelevant for answering  $Q$ . This is a preprocessing step, which is not discussed in detail here.

- *Logic Translator*: It takes the specification of  $\mathcal{S}$  relevant for  $Q$  returned by the Pruner, and produces the logic program  $\Pi_{\mathcal{S}}(Q) = \Pi_{\Sigma} \cup \Pi_Q$ , according to some encoding proposed in the literature. In our tests, we used the mapping in [14, 30].

- *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, it is also responsible for computing a factorization.

- *Stable Models Engine*: It takes as input the affected database and computes the repairs using the program  $\Pi_{\Sigma}$ . In our implementation, we used the DLV system [36].

- *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<sup>3</sup>. In the case the constraints are not in the class  $\mathbf{C}_1$ , it is also responsible of filtering 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 both the safe part of the database and the repair computed by the Stable Model Engine. After that a new repair is stored, it notifies the query executor module.

- *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 executing the reformulated query. Specifically, it implements the strategies described in Equation (6) (for its adaptations to the classes  $\mathbf{C}_1$  and  $\mathbf{C}_0$ ) and Equation (10), using

<sup>3</sup><http://www.mat.unical.it/wrapper/index.html>

$player_m^{M_{D_0}}$ :	<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>10</td><td><i>Totti</i></td><td><i>RM</i></td><td>'11'</td></tr> <tr><td>9</td><td><i>Ronaldinho</i></td><td><i>BC</i></td><td>'11'</td></tr> </table>	10	<i>Totti</i>	<i>RM</i>	'11'	9	<i>Ronaldinho</i>	<i>BC</i>	'11'
10	<i>Totti</i>	<i>RM</i>	'11'						
9	<i>Ronaldinho</i>	<i>BC</i>	'11'						
$coach_m^{M_{D_0}}$ :	<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>7</td><td><i>Capello</i></td><td><i>RM</i></td><td>'11'</td></tr> </table>	7	<i>Capello</i>	<i>RM</i>	'11'				
7	<i>Capello</i>	<i>RM</i>	'11'						

$team_m^{M_{D_0}}$ :	<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td><i>RM</i></td><td><i>Roma</i></td><td>10</td><td>'10'</td></tr> <tr><td><i>BC</i></td><td><i>Barcelona</i></td><td>8</td><td>'11'</td></tr> <tr><td><i>RM</i></td><td><i>Real Madrid</i></td><td>10</td><td>'01'</td></tr> </table>	<i>RM</i>	<i>Roma</i>	10	'10'	<i>BC</i>	<i>Barcelona</i>	8	'11'	<i>RM</i>	<i>Real Madrid</i>	10	'01'
<i>RM</i>	<i>Roma</i>	10	'10'										
<i>BC</i>	<i>Barcelona</i>	8	'11'										
<i>RM</i>	<i>Real Madrid</i>	10	'01'										

Figure 4: The database of our running example after marking.

standard envelopes. 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 query over  $R_1$  plus the safe part  $S$ . Then, for each other repair  $R_i$ , the table is updated by filtering those tuples that do not occur in the answers of the query over  $R_i$  plus the safe part. 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 computation of a decomposable partitioning is needed.

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

We next consider the idea of grouping the repairs computed by the *Stable Models Engine* in a way such that a single query may evaluate more than one repair at time. This can be done using a marking strategy.

Let  $R_1, \dots, R_n$  be the repairs of the affected part (intersected, if needed, with the envelope  $E$ ) 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  $s(t) \in D$ , we associate a mark  $b = 'b_1 \dots b_n'$  such that, for every  $1 \leq i \leq n$ ,  $b_i = 1$  if  $s(t)$  belongs to  $R_i$ , and  $b_i = 0$  otherwise. The marked tuple  $t, b$  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...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<sup>-</sup> 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  $SQ L_r$ . Informally,  $SQ L_r$  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 do not belong to any repair) are considered as marked by '0...0'. The details are given in Appendix C.

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  $SQ L_h$ :

```
CREATE VIEW  $h_m(a_1, \dots, a_n, mark)$  AS
SELECT  $a_1, \dots, a_n, \text{SUMBIT}(mark)$ 
FROM ( $SQ L_{r_1} \text{ UNION } \dots \text{ UNION } SQ L_{r_\ell}$ )
GROUP BY  $a_1, \dots, a_n,$ 
```

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

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

It computes the query answers by selecting the facts which evaluate to true in all repairs.

**Example 6.1** The query in our running query 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:

$$\begin{aligned} r'_1 : q(y_{0,1}) &\leftarrow player(y_{1,1}, y_{1,2}, y_{1,3}), y_{0,1} = y_{1,1}; \\ r'_2 : q(y_{0,1}) &\leftarrow team(y_{1,1}, y_{1,2}, y_{1,3}), y_{0,1} = y_{1,3}. \end{aligned}$$

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 the Appendix C).

**Proposition 6.1** *Let  $D$  be a database for  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ , let  $Q$  be a non-recursive Datalog<sup>-</sup> query over it, and let  $R_1, \dots, R_n$  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  $\bigcap_{i=1}^n \{t \mid t \in Q[R_i \cup S]\}$ , for  $S = D \setminus E$ .*

Note that when  $R'_1, \dots, R'_n$  are all repairs for  $A$ , then the tuples computed by  $SQL_Q$  are the consistent answer 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  $'11 \dots 1'$ , 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. For details on scaling the technique this way, we also refer to the Appendix. 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.

## 7 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.

## 7.1 Benchmark Databases and Compared Methods

Hippo [17, 18] and ConQuer [26, 27] are two most noticeable prototype systems available in the literature for consistent query answering from inconsistent databases. These systems focus on specific classes of constraints and queries over which consistent answers can be efficiently computed. Indeed, Hippo is able to deal with queries definable in relational algebra without projection under denial constraints, and encodes inconsistencies by a conflict hypergraph built from constraint violations [2, 16]. Instead, ConQuer considers database schemas under key dependencies, and deals with user queries from a subclass of conjunctive queries, possibly enriched with aggregates. On this class, ConQuer is shown to be very efficient for providing consistent answers over large databases with many inconsistent tuples (up to the 50% of the whole database).

Since these systems are tailored to efficiently manage specific classes of queries and constraints, their performances have been tested on some ad-hoc created benchmark databases. Specifically, [26] mainly generated syntectic data for the TCP-H specifications over a schema containing primary keys only and used queries with aggregate expressions, while [17] considered project-free queries over tables having attributes  $x, y, z$  correlated by the functional dependency  $p(x, y, z) \wedge p(x, y', z') \supset z = z'$ .

In the same way as Hippo and ConQuer are not directly comparable with each other, thereby having required the definition of specific benchmarks problems and data, some novel scenario has to be proposed to assess the effectiveness of our localization approach. Indeed, our techniques are designed for more general settings than those addressed by Hippo and ConQuer, and therefore they can be used in scenarios that go beyond the scope of such systems. For instance, neither Hippo nor ConQuer can answer queries involving projections, when an integrity constraint which is not a key (such as a general functional dependency or an exclusion dependency) is issued on the schema. But, these constraints often occur in database design; in fact, exclusion dependencies are typical for database schemes stemming from ER-models or other conceptual data modeling languages, and are widely used in applications in which the schema is given in terms of an ontology.

On the other hand, if we focus on the class of queries for which Hippo and ConQuer have been respectively 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.

Therefore, to test our more general framework as well as the factorization techniques discussed in Section 5, we proceed as follows:

- We first present a global picture of our approach by considering results for our running example (on football teams).
- We then focus on a test suite over the database schema  $\mathcal{S}_f^2$  used in [17], which contains two relations of the form  $r_i(x_i, y_i, z_i)$ , with  $i = 1, 2$ , and the functional dependency  $r_i(x, y, z) \wedge r_i(x, y', z') \supset z = z'$ , but we consider queries that involve projections, so that the system Hippo is not applicable. Also, we consider the database  $\mathcal{S}_e$ , which contains a predicate of the form  $p(x, y)$  and a predicate of the form  $q(v, w)$ , with an exclusion dependency between attributes  $x$  and  $v$ .
- We also discuss the impact of the number of atoms involved in the query on the performances of the localization approach, by considering the schema  $\mathcal{S}_f^N$ , obtained by generalizing  $\mathcal{S}_f^2$ , to an increasingly large number of predicates.

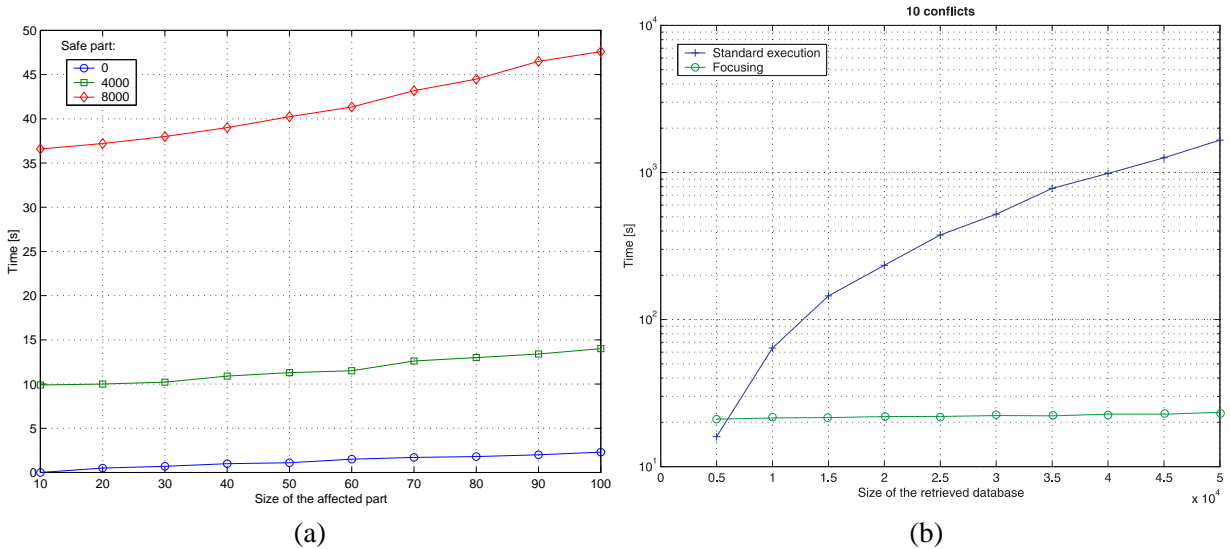


Figure 5: Football Team. (a) Execution time in DLV system w.r.t. size of the affected part. (b) Comparison with the optimization method.

For the schemas above, we generated some random data according to the ideas, described in [17, 26], of tuning the size of the safe part and the number of conflicts.

All experiments have been carried out on a 1.6GHz Pentium IV with 512MB memory, by assessing the time needed for query answering when the DLV system computes repairs of the affected part only, plus the time required for the recombination of the results in the DBMS PostgreSQL. In the experiments, the approach proposed in the present paper is compared with the approach in which DLV is used to evaluate the whole logic specification for querying the inconsistent databases. In both cases, the logic programming encoding we used is inspired by [14, 30].

## 7.2 The Football Teams Example

We next present an overview of the performances of our approach and, specifically, its scaling w.r.t. the size of the safe database, by considering 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*; note that 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<sup>-</sup> program  $\Pi_{S_0}(Q)$  in the line of [14, 30] (the encoding used is the one provide in Appendix D.1, in which we get rid off the encoding for the mapping). We first measure the execution time of the program  $\Pi_{S_0}(Q)$  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 conflict) would require more than 35 second for consistent query answering.

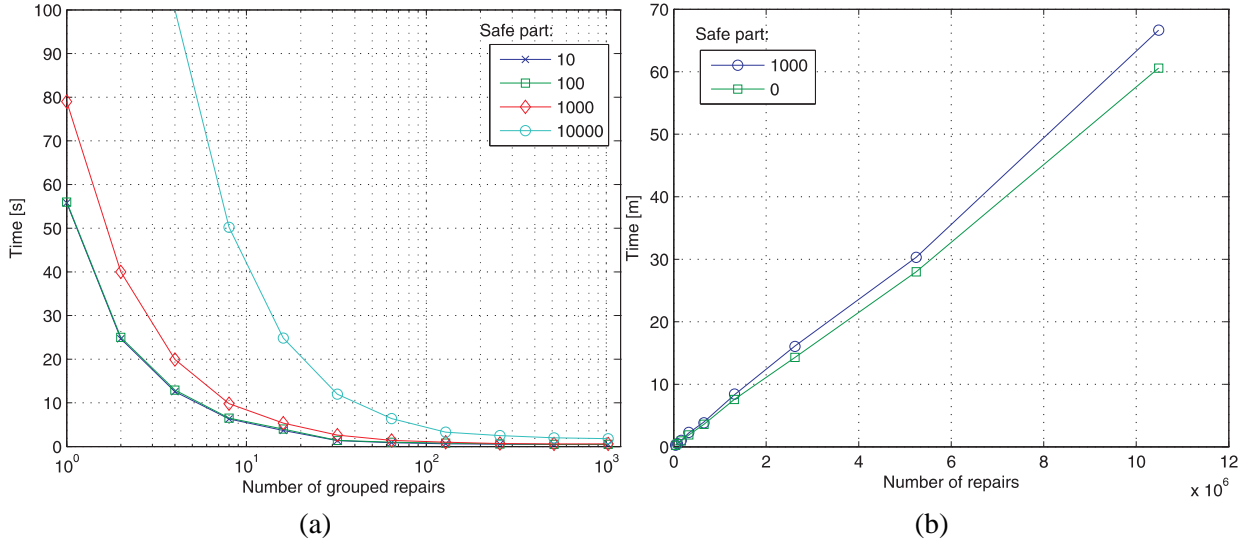


Figure 6: Query answering over  $\mathcal{S}_f^2$ . (a) Optimization method w.r.t. number  $n$  of grouped repairs, for a fixed number of conflicts. (b) Optimization method w.r.t. the size of the affected part, for  $n = 2^9$ .

The performance degradation under varying size database 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, so 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.

### 7.3 Scalability Assessment

In a first series of experiments, we assessed the relevance of the strategy for grouping repair computation by focusing on the database  $\mathcal{S}_f^2$ . Indeed, so far, we have assumed that the marking string suffices 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 answering 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.

A second set of experiments has been devoted to assess the scalability w.r.t. the number of conflicts. Thus, we reconsidered the above scenario, but we augmented the number of conflicts up to 100, and we fixed the marking string to  $2^9$  bits. The results are shown in Figure 6.(b). We notice the decent scaling w.r.t. the size of the safe part and that, as expected, the time needed for answering a query linearly depends on the total number of repairs. Note that, since this number is exponential in the number of constraint violation,



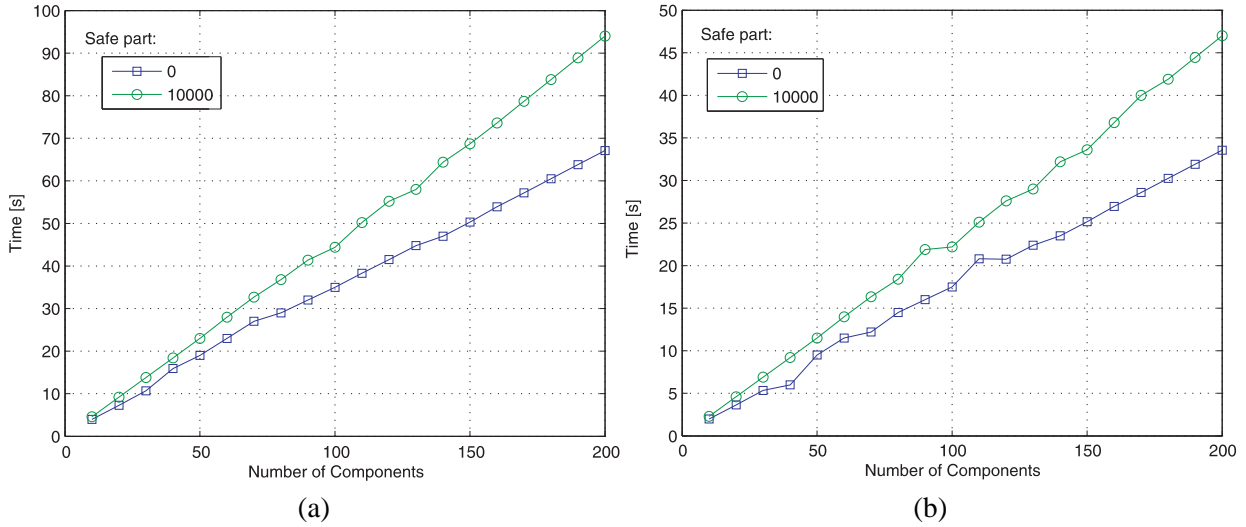


Figure 7: Factorization strategy. Answering (a)  $Q_f$  over  $\mathcal{S}_f^2$  w.r.t. the size of the affected part (b)  $Q'_f$  over  $\mathcal{S}_f^2$ .

the resulting curve may be in some scenarios exponential in the size of the affected part. Yet, this is the best one expect to achieve for general inconsistent databases for which the consistent query answering problem is known to be 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 6.2. In this respect, we notice that the query  $Q_f$  above and the constraints over  $\mathcal{S}_f^2$  are such that our factorization strategy can be applied. Indeed, the setting we are considering is basically the one described in Example 5.1, 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 1000 tuples and 1000 repairs (any pair of these tuples witnesses a violation of the dependency), and we generated some syntectic data for increasingly large number of independent components. The parameter  $n$ , bounding the number of repairs simultaneously processed, is fixed 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 (results are reported up to  $200^{1000}$  repairs).

A similar experiment has been repeated for the query  $Q'_f = \langle q'_f, \mathcal{P}'_f \rangle$  where  $\mathcal{P}'_f = \{q'_f(y_1) \leftarrow r_1(x, y_1, z_1)\}$ . Figure 7.(b) shows that  $Q'_f$  has half the answering time of  $Q_f$ , given that answering  $Q'_f$  does not need to resolve conflicts in the relation  $r_2$ .

As a further example to test the factorization technique, we considered the database  $\mathcal{S}_e$ , which contains a predicate of the form  $p(x, y)$  and a predicate of the form  $q(v, w)$ , with a constraint imposing that attributes  $x$  and  $v$  are disjoint (exclusion dependency). Over it, we evaluated the query  $Q_e = \langle q_e, \mathcal{P}_e \rangle$  where  $\mathcal{P}_e = \{q_e(y) \leftarrow p(x, y), q(x, w), q_e(w) \leftarrow p(x, y), q(x, w)\}$ . Notice that each component in a factorization contains tuples of the form  $p(a, y)$  and  $q(a, w)$  witnessing the violation of the exclusion dependency between  $p$  and  $q$ . Moreover, repairing each component always yields two repairs (one selecting tuples over  $p$  and the other selecting tuples over  $q$ ). Figure 8.(a) reports timing for consistent query answering, where the safe part consists of 10000 tuples.

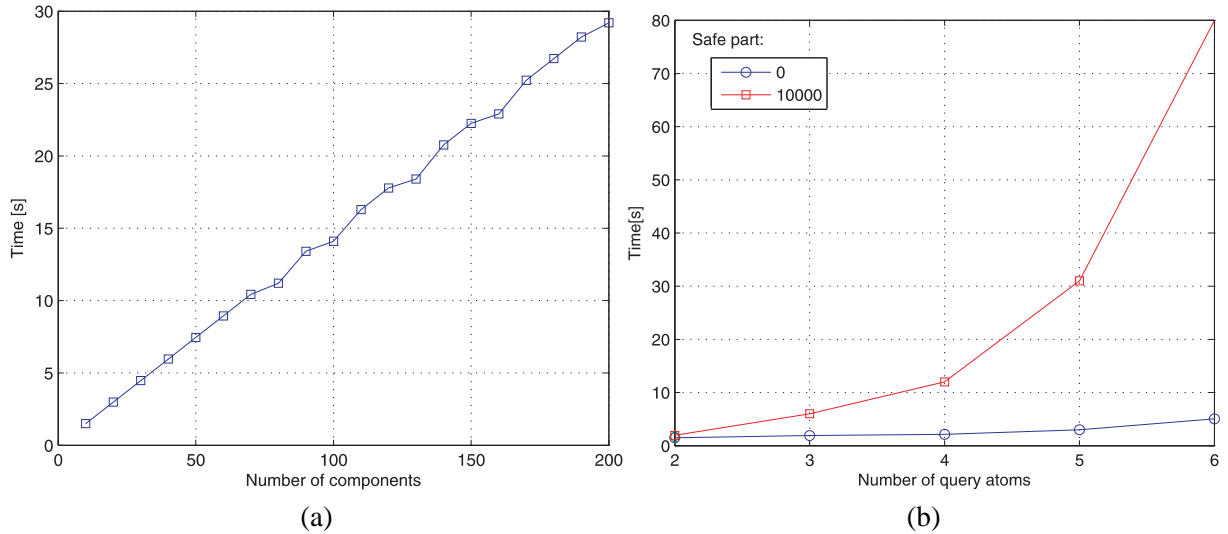


Figure 8: (a) Factorization strategy over  $S_e$ . (b) Query answering over  $S_f^N$ ,  $N > 2$ .

Finally, we also considered the database  $S_f^N$  and the query  $Q = \langle q, \mathcal{P} \rangle$  where  $\mathcal{P} = \{q(y_1) \leftarrow r_1(x, y_1, z_1), r_2(x, y_2, z_2), \dots, r_N(x, y_N, z_N)\}$ , 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, which shows an exponential dependency.

## 8 Discussion and Conclusion

For optimizing logic-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 [17]), 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. For example, Chomicki and Marcinkowski [16] 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.

Furthermore, answering a negation-free query  $Q$  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 set-inclusion minimal repairs w.r.t.  $\leq_D$  as repairs under an ordering  $\leq'_D$  satisfying these properties. An example is the ordering  $R_1 \sqsubseteq_D R_2$  iff  $R_1 \cap D \supseteq R_2 \cap D$  [13, 14], which violates (SIP). We can use here the ordering  $R_1 \sqsubseteq'_D R_2$  iff  $R_1 \sqsubseteq_D R_2 \wedge (R_1 \cap D = R_2 \cap D \Rightarrow R_1 \setminus D \subseteq R_2 \setminus D)$  instead.

Other approaches considered consistent query answering under the perspective of modifying values in the database rather than entire tuples [43]. 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  $\mathcal{S}$  [34]. As briefly discussed in Appendix D.1 and more in detail in [21], 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 [33, 14], and certain Local-As-View (LAV) settings, e.g. as in [8, 11].

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 [35].

## References

- [1] ABITEBOUL, S., HULL, R., AND VIANU, V. 1995. *Foundations of Databases*. Addison Wesley.
- [2] ARENAS, M., BERTOSSI, L., AND CHOMICKI, J. 2001. Scalar aggregation in fd-inconsistent databases. In *International Conference on Database Theory (ICDT-2001)*. Springer Verlag, 39–53.
- [3] ARENAS, M., BERTOSSI, L. E., AND CHOMICKI, J. 1999. Consistent query answers in inconsistent databases. In *Proc. of the 18th ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems (PODS'99)*. 68–79.
- [4] ARENAS, M., BERTOSSI, L. E., AND CHOMICKI, J. 2003. Answer sets for consistent query answering in inconsistent databases. *Theory and Practice of Logic Programming* 3, 4, 393–424. arXiv.org paper cs.DB/0207094.
- [5] ARIELI, O., DENECKER, M., NUFFELEN, B. V., AND BRUYNNOOGHE, M. 2004. Coherent integration of databases by abductive logic programming. *Journal of Artificial Intelligence Research* 21, 245–286.
- [6] BARCELÓ, P. AND BERTOSSI, L. E. 2003. Logic programs for querying inconsistent databases. In *Proc. of Practical Aspects of Declarative Languages (PADL'03)*. 208–222.
- [7] 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-Verlag, Chapter 2, 43–83.
- [8] BERTOSSI, L. E., 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 2002)*. 71–85.

- [9] BERTOSSI, L. E., HUNTER, A., AND SCHAUB, T., Eds. 2005. *Inconsistency Tolerance [result from a Dagstuhl seminar]*. Lecture Notes in Computer Science, vol. 3300. Springer.
- [10] BOUZEGHOUB, M. AND LENZERINI, M. 2001. Introduction to the special issue on data extraction, cleaning, and reconciliation. *Information Systems* 26, 8, 535–536.
- [11] 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 2003)*. 10–15.
- [12] BRAVO, L. AND BERTOSSI, L. 2005. Disjunctive deductive databases for computing certain and consistent answers to queries from mediated data integration systems. *Journal of Applied Logic* 3, 1, 329–367.
- [13] 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 SIGART Symp. on Principles of Database Systems (PODS 2003)*. 260–271.
- [14] 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 2003)*. 16–21.
- [15] CHOMICKI, J. 2007. Consistent query answering: Five easy pieces. In *Proceedings of the 11th International Conference on Database Theory (ICDT 2007)*, T. Schwentick and D. Suciu, Eds. LNCS. Springer, Berlin Heidelberg, Germany, 1–17.
- [16] CHOMICKI, J. AND MARCINKOWSKI, J. 2005. Minimal-change integrity maintenance using tuple deletions. *Information and Computation* 197, 1-2, 90–121. Earlier version: arXiv.org paper cs.DB/0212004.
- [17] CHOMICKI, J., MARCINKOWSKI, J., AND STAWORKO, S. 2004a. Computing consistent query answers using conflict hypergraphs. In *Proc. 13th ACM Conference on Information and Knowledge Management (CIKM-2004)*. ACM Press, 417–426.
- [18] CHOMICKI, J., MARCINKOWSKI, J., AND STAWORKO, S. 2004b. Hippo: A system for computing consistent answers to a class of sql queries. In *Proceedings 9th International Conference on Extending Database Technology (EDBT-2004)*. Number 2992 in LNCS. Springer, 841–844.
- [19] DANTSIN, E., EITER, T., GOTTLOB, G., AND VORONKOV, A. 2001. Complexity and Expressive Power of Logic Programming. *ACM Computing Surveys* 33, 3, 374–425.
- [20] EITER, T., FINK, M., GRECO, G., AND LEMBO, D. 2003. Efficient Evaluation of Logic Programs for Querying Data Integration Systems. In *Proceedings 19th International Conference on Logic Programming (ICLP 2003)*, C. Palamidessi, Ed. Number 2916 in LNCS. Springer, 163–177.
- [21] 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, Institut für Informationssysteme, Technische Universität Wien, A-1040 Vienna, Austria. July. Extends [20].
- [22] EITER, T., GOTTLOB, G., AND MANNILA, H. 1997. Disjunctive Datalog. *ACM Transactions on Database Systems* 22, 3 (September), 364–418.

- [23] FAGIN, R., KOLAITIS, P., MILLER, R. J., AND POPA, L. 2005. Data exchange: Semantics and query answering. *Theoretical Computer Science* 336, 89–124.
- [24] 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.
- [25] FRIEDMAN, M., LEVY, A. Y., AND MILLSTEIN, T. D. 1999. Navigational plans for data integration. In *Proceedings Sixteenth National Conference on Artificial Intelligence and Eleventh Conference on Innovative Applications of Artificial Intelligence (AAAI/IAAI)*. AAAI Press / The MIT Press, 67–73.
- [26] FUXMAN, A., FAZLI, E., AND MILLER, R. J. 2005. ConQuer: Efficient management of inconsistent databases. In *SIGMOD Conference*. 155–166.
- [27] FUXMAN, A. AND MILLER, R. J. 2005. First-order query rewriting for inconsistent databases. In *Proceedings of the 10th International Conference on Database Theory (ICDT 2005)*, T. Eiter and L. Libkin, Eds. Number 3363 in LNCS. Springer, 337–351.
- [28] GELFOND, M. AND LIFSCHITZ, V. 1991. Classical Negation in Logic Programs and Disjunctive Databases. *New Generation Computing* 9, 365–385.
- [29] GRECO, G., GRECO, S., AND ZUMPARNO, E. 2003. A logical framework for querying and repairing inconsistent databases. *IEEE Trans. on Knowledge and Data Engineering* 15, 6, 1389–1408.
- [30] GRIECO, L., LEMBO, D., ROSATI, R., AND RUZZI, M. 2005. Consistent query answering under key and exclusion dependencies: algorithms and experiments. submitted for publication to an international conference.
- [31] KIFER, M. AND LOZINSKII, E. L. 1992. A logic for reasoning with inconsistency. *Journal of Automated Reasoning* 9, 2, 179–215.
- [32] 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.
- [33] 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 (KRDB 2002)*. <http://ceur-ws.org/Vol-54/>.
- [34] LENZERINI, M. 2002. Data integration: A theoretical perspective. In *Proc. of the 21st ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems (PODS 2002)*. 233–246.
- [35] LEONE, N., EITER, T., FABER, W., FINK, M., GOTTLÖB, 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 *Proceedings ACM SIGMOD/PODS 2005 Conference, June 13-16, 2005, Baltimore, Maryland*. ACM, 915–917.
- [36] LEONE, N., PFEIFER, G., FABER, W., EITER, T., GOTTLÖB, G., PERRI, S., AND SCARCELLO, F. 2006. The DLV System for Knowledge Representation and Reasoning. *ACM Transactions on Computational Logic* 7, 3, 499–562.

- [37] LIN, J. 1996. Integration of weighted knowledge bases. *Artificial Intelligence* 83, 2, 363–378.
- [38] LIN, J. AND MENDELZON, A. O. 1998. Merging databases under constraints. *Int. J. of Cooperative Information Systems* 7, 1, 55–76.
- [39] NUFFELEN, B. V., CORTÉS-CALABUIG, A., DENECKER, M., ARIELI, O., AND BRUYNOOGHE, M. 2004. Data integration using id-logic. In *Proceedings 16th International Conference on Advanced Information Systems Engineering (CAiSE 2004)*, A. Persson and J. Stirna, Eds. Lecture Notes in Computer Science, vol. 3084. Springer, 67–81.
- [40] SIMONS, P., NIEMELÄ, I., AND SOININEN, T. 2002. Extending and Implementing the Stable Model Semantics. *Artificial Intelligence* 138, 181–234. Smodels home page: <http://www.tcs.hut.fi/Software/smodels/>.
- [41] STAWORKO, S., CHOMICKI, J., AND MARCINKOWSKI, J. 2006. Preference-driven querying of inconsistent relational databases. In *EDBT Workshops*, T. Grust, H. Höpfnér, A. Illarramendi, S. Jablonski, M. Mesiti, S. Müller, P.-L. Patranjan, K.-U. Sattler, M. Spiliopoulou, and J. Wijsen, Eds. Lecture Notes in Computer Science, vol. 4254. Springer, 318–335.
- [42] ULLMAN, J. D. 1989. *Principles of Database and Knowledge Base Systems*. Computer Science Press.
- [43] WIJSEN, J. 2005. Database repairing using updates. *ACM Transactions on Database Systems* 30, 3, 722–768.

## A Proofs for Section 3

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

*Proof.* Let  $R$  be an arbitrary database of  $\mathcal{S}$  consistent with  $\Sigma$ . Let  $R'$  result from  $R$  by removing every fact containing some constant  $c \notin \text{adom}(D, \mathcal{S})$ . 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  $A_1(\vec{c}_1) \wedge \dots \wedge A_l(\vec{c}_l) \supset B_1(\vec{d}_1) \vee \dots \vee B_m(\vec{d}_m) \vee \phi_1(\vec{e}_1) \vee \dots \vee \phi_n(\vec{e}_n)$  which is violated by  $R'$ , i.e., (i)  $A_1(\vec{c}_1), \dots, A_l(\vec{c}_l) \in R'$ , (ii)  $B_1(\vec{d}_1), \dots, B_m(\vec{d}_m) \notin R'$ , and (iii)  $\phi_1(\vec{e}_1) \vee \dots \vee \phi_n(\vec{e}_n)$  is false. Since  $R \models \sigma^g$ , by construction of  $R'$  we have  $B_j(\vec{d}_j) \in R \setminus R'$  for some  $j \in \{1, \dots, m\}$  and thus  $\vec{d}_j$  contains some constant  $c \notin \text{adom}(D, \mathcal{S})$ . 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,  $\Delta(R', D) \subset \Delta(R, D)$ , and thus by (SIP)  $R' <_D R$ . The  $<_D$ -minimality of repairs implies that each database in  $\text{rep}(D)$  involves only constants from  $\text{adom}(D, \mathcal{S})$ . Furthermore, by consistency of  $\mathcal{S}$  and the fact that each sequence  $R_1 >_D R_2 >_D \dots R_i >_D \dots$  of databases  $R_i$  on  $\text{adom}(D, \mathcal{S})$  must be finite, one such repair  $R$  must exist.  $\square$

**Proposition 3.3** *Let  $D$  be a database for  $\mathcal{S} = \langle \Psi, \Sigma \rangle$  where no built-in relations occur in  $\Sigma$  except  $=$  and  $\neq$ . Suppose that  $<_D$  satisfies (SIP). Then, every repair  $R \in \text{rep}(D)$  involves only constants from  $\text{adom}(D, \mathcal{S})$ , and some repair exists if  $\mathcal{S}$  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  $\vec{c}_j = c_{j,1}, \dots, c_{j,n_j}$  contains some constant  $c_{j,h} \notin \text{adom}(D, \mathcal{S})$  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  $\text{adom}(D, \mathcal{S})$ , because the head of  $\sigma^g$  is satisfied in  $R'$ , which contains only constants from  $\text{adom}(D, \mathcal{S})$ ). 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.  $\square$

## B Proofs for Section 5

**Theorem 5.2 (Factorization)** *Let  $D$  be a database for  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ . Let  $E_1, \dots, E_m$  be a repair-compliant partitioning of a repair envelope  $E$  for  $D$ . Then,  $E_1, \dots, E_m$  is a factorization of  $E$  for  $D$  and  $\mathcal{S}$ .*

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

( $\subseteq$ )  $R \in \text{rep}(D \cap E)$  implies  $R = R_1 \cup \dots \cup R_m$  and  $R_i \in \text{rep}(D \cap E_i)$  for  $1 \leq i \leq m$ ;

( $\supseteq$ ) every  $R \in \{R_1 \cup \dots \cup R_m \mid R_i \in \text{rep}(D \cap E_i), 1 \leq i \leq m\}$  is a repair of  $D \cap E$ .

( $\subseteq$ ) Let  $R \in \text{rep}(D \cap E)$ . Then, by repair-compliance of  $E_1, \dots, E_m$ ,  $R = F \cup R_E$ , 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  $R_i \in \text{rep}(D \cap E_i)$  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  $\sigma \in \text{ground}(\Sigma)$  such that  $R_i \not\models \sigma$ . Thus,  $R_i \models \text{body}(\sigma)$ , which implies  $R \models \text{body}(\sigma)$ , and  $R_i \not\models \text{head}(\sigma)$ . However,  $R \models \text{head}(\sigma)$  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  $R_i \not\models \text{head}(\sigma)$ , 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  $\text{body}(\sigma) \subseteq F$  and  $\text{head}(\sigma) \cap E_i = \emptyset$ . This excludes the existence of a repair of the form  $F \cup R'_{E_i}$  of  $D \cap E_i$ , such that  $R'_{E_i} \subseteq E_i$ , a contradiction to the repair compliance of  $E_1, \dots, E_m$ . This proves  $R_i \models \Sigma$  for every  $i = 1, \dots, m$ .

Consequently,  $R_i \in \text{rep}(D \cap E_i)$  iff there exists no  $R'_i \in \text{rep}(D \cap E_i)$  such that  $R'_i <_{D \cap E_i} R_i$  and  $R'_i \models \Sigma$ . Assume such an  $R'_i$  would exist. 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  $R'_E = (R \cap E_1) \cup \dots \cup R'_{E_i} \cup \dots \cup (R \cap E_m)$ , 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  $\sigma \in \text{ground}(\Sigma)$  such that  $R' \models \text{body}(\sigma)$  and  $R' \not\models \text{head}(\sigma)$ , while  $R \models \sigma$ . We can conclude that  $\text{body}(\sigma) \subseteq R'_i$ , and since  $R'_i \models \Sigma$  we obtain  $R' \models \text{head}(\sigma)$ , a contradiction.) Together with  $R' <_{D \cap E} R$ , however, this contradicts  $R \in \text{rep}(D \cap E)$ . Hence,  $R_i \in \text{rep}(D \cap E_i)$ , for  $1 \leq i \leq m$ .

( $\supseteq$ ) Let  $R \in \{R_1 \cup \dots \cup R_m \mid R_i \in \text{rep}(D \cap E_i), 1 \leq i \leq m\}$ . We show that  $R \in \text{rep}(D \cap E)$ . Towards a contradiction suppose  $R \not\models \Sigma$ , i.e.,  $R \not\models \sigma$  for some  $\sigma \in \text{ground}(\Sigma)$ . By definition of a repair-compliant partitioning, we conclude that  $R = F \cup (R_1 \cap E_1) \cup \dots \cup (R_m \cap E_m)$ , where  $F = R \setminus E$  for any  $1 \leq i \leq m$ . Consequently,  $R \models \text{body}(\sigma)$  implies  $R_i \models \text{body}(\sigma)$  for some  $1 \leq i \leq m$  by constraint-boundedness. However,  $R_i \not\models \text{head}(\sigma)$  (otherwise  $R \models \text{head}(\sigma)$ ), which contradicts  $R_i \in \text{rep}(D \cap E_i)$ . Hence,  $R \models \Sigma$ .

It remains to show that there is no  $R' \in \text{rep}(D \cap E)$  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  $i = 1, \dots, m$ . Case (i) leads to a contradiction with  $R_i \in \text{rep}(D \cap E_i)$ , 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  $i = 1, \dots, m$ . 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  $E_1, \dots, E_m$ , is equivalent to  $F <_{\emptyset} F$ , a contradiction. To see this, note that we can apply (DIS) to  $R' \setminus E_i <_{(D \cap E) \setminus E_i} R \setminus E_i$  wrt.  $E_j$  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  $R' \setminus (E_1 \cup \dots \cup E_m) <_{(D \cap E) \setminus (E_1 \cup \dots \cup E_m)} R \setminus (E_1 \cup \dots \cup E_m)$ , which is equivalent  $R' \setminus E <_{(D \cap E) \setminus E} R \setminus E$ . This proves  $R \in \text{rep}(D \cap E)$ .  $\square$

**Proposition 5.3** *Let  $D$  be a database for  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ , 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  $E_1, \dots, E_m$  of  $E$  is repair-compliant.*

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

Case (1):  $E = C^*$ . By Proposition 4.6,  $R \in \text{rep}(D \cap E)$  implies  $R \subseteq E$ , hence  $R \setminus E = \emptyset$ . Towards a contradiction assume that there exists  $R_i \in \text{rep}(D \cap E_i)$ , 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  $\sigma \in \text{ground}(\Sigma)$  such that  $R'_i \not\models \sigma$ . Thus,  $R'_i \models \text{body}(\sigma)$ , which implies  $R_i \models \text{body}(\sigma)$ , and  $R'_i \not\models \text{head}(\sigma)$ . However,  $R_i \models \text{head}(\sigma)$  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  $R'_i \not\models \text{head}(\sigma)$ , none of the built-in predicates of  $\sigma$  is true and therefore  $B(\vec{y})$  is a fact. Furthermore,  $\text{body}(\sigma) \neq \emptyset$  since  $\sigma \in \mathbf{C}_1$ , and thus,  $B(\vec{y}) \notin E$ , since  $E_1, \dots, E_m$  is constraint-bounded. Consequently,  $\text{facts}(\sigma) \cap E \neq \emptyset$  holds, as well as  $\text{facts}(\sigma) \not\subseteq E$ . 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  $R_i \in \text{rep}(D \cap E_i)$ . Hence,  $R_i \setminus E_i = \emptyset$  must hold.

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

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

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



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

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

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

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

Consider, now, two repairs  $R'_h$  and  $R''_h$  of  $E_h$ , with  $1 \leq h \leq \ell$ . Given that  $E_h$  is singular, it holds that:  $Q[(D \setminus E) \cup R_1 \cup \dots \cup R'_h \cup \dots \cup R_m] = Q[(D \setminus E) \cup R_1 \cup \dots \cup R''_h \cup \dots \cup R_m]$ . Therefore, for each singular component  $E_h$ , we can take an arbitrary repair  $R_h$  and conclude that:

$$\bigcap_{\tilde{R}_h \in \text{rep}(D \cap E_h)} Q[(D \setminus E) \cup R_1 \cup \dots \cup \tilde{R}_h \cup \dots \cup R_m] = Q[(D \setminus E) \cup R_1 \cup \dots \cup R_h \cup \dots \cup R_m].$$

Hence,  $\text{ans}(Q, D)$  can be also computed as:

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

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

**Proposition 5.5** *Let  $E_1, \dots, E_m$  be a factorization of a repair envelope  $E$  for  $D$  which is decomposable w.r.t. query  $Q$  having the non-singular components  $E_\ell, \dots, E_m$ . Then*

$$\text{ans}(Q, D) = \bigcup_{i=\ell}^m \left( \bigcap_{R_i \in \text{rep}(D \cap E_i)} Q[S_Q \cup R_i] \right)$$

*Proof.* Let  $E_1, \dots, E_m$  be a factorization of a repair envelope  $E$  for  $D$ , and let  $E_\ell, \dots, E_m$  be its non-singular components. We first show that

$$Q[S_Q \cup R_{\ell+1} \cup \dots \cup R_m] = Q[S_Q \cup R_{\ell+1}] \cup \dots \cup Q[S_Q \cup R_m], \quad (11)$$

for every  $R_i \in \text{rep}(D \cap E_i)$ ,  $1 \leq i \leq m$ . We have  $Q = \langle q, \mathcal{P} \rangle$ , where  $\mathcal{P} = \{\rho_1, \dots, \rho_n\}$  is a set of “not”-free rules  $\rho_j$ . Consider any ground instance  $\rho'_j$  of  $\rho_j$ , and an atom  $p_i(\bar{t}_i)$  occurring in the body of  $\rho'_j$  that is satisfied by  $S_Q \cup R_\ell \cup \dots \cup R_m$ . Then either  $p_i(\bar{t}_i) \in S_Q$  or  $p_i(\bar{t}_i) \in R_h$ , for some  $\ell \leq h \leq m$ . Furthermore, since  $E_1, \dots, E_m$  is decomposable, in the case where  $p_i(\bar{t}_i) \in R_h \setminus (S_Q \cup \bigcap_{k=\ell}^m R_k)$ , there is no atom  $p_k(\bar{t}_k)$  in the body of  $\rho'_j$  which belongs to  $R_{h'} \setminus R_h$ , for some  $\ell \leq h' \neq h \leq m$ . Indeed, by Condition (3),  $R'_h \setminus R_h = R_{h'} \cap E_{h'}$ , and thus  $p_k(\bar{t}_k) \in E_{h'}$  would hold, while  $p_i(\bar{t}_i) \in E_h$  holds. Conditions (1) and (2) would imply that  $p_i(\bar{t}_i)$  and  $p_k(\bar{t}_k)$  are instances of the same atom in the body of  $\rho_j$ , and thus  $p_i(\bar{t}_i) = p_k(\bar{t}_k)$ . This, however, contradicts that  $E_h \cap E_{h'} = \emptyset$ .

As a consequence, the body of  $\rho'_j$  is satisfied by  $S_Q \cup R_\ell \cdots \cup R_m$  if and only if it is satisfied by  $S_Q \cup R_h$ , for some  $\ell \leq h \leq m$ . Since  $\rho_j$  is non-disjunctive and “not”-free, we can conclude from this that  $Q^{(j)}[S_Q \cup R_\ell \cdots \cup R_m] = \bigcup_{i=\ell}^m Q^{(j)}[S_Q \cup R_i]$ , where  $Q^{(j)} = \langle q, \{\rho_j\} \rangle$ , and that

$$\begin{aligned} Q[S_Q \cup R_\ell \cdots \cup R_m] &= \bigcup_{j=1}^n Q^{(j)}[S_Q \cup R_\ell \cdots \cup R_m] = \bigcup_{j=1}^n \bigcup_{i=\ell}^m Q^{(j)}[S_Q \cup R_i] \\ &= \bigcup_{i=\ell}^m \bigcup_{j=1}^n Q^{(j)}[S_Q \cup R_i] = Q[S_Q \cup R_{\ell+1}] \cup \cdots \cup Q[S_Q \cup R_m]. \end{aligned}$$

This proves (11). To conclude the proof, we observe that from Proposition 5.4, the answer to a non-recursive Datalog query  $Q$  w.r.t.  $D$  is:

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

By Equation (11), we then get:

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

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

## C 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.

### C.1 Query Reformulation

We first show how a non-recursive Datalog<sup>-</sup> query  $Q$  can be reformulated into an SQL query whose evaluation over the marked database returns the answers to  $Q$  that are true in any repair of  $\{R_1, \dots, R_n\}$ . Let  $r : h(\vec{x}_0) \leftarrow B(\vec{x})$  be a safe rule of form

$$p_0(\vec{x}_0) \leftarrow p_1(\vec{x}_1), \dots, p_l(\vec{x}_l), \text{ not } p_{l+1}(\vec{x}_{l+1}), \dots, \text{ not } p_{l+k}(\vec{x}_{l+k}).^4 \quad (12)$$

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  $1 \leq j \leq k_i$ . 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;

<sup>4</sup>For the sake of simplicity and w.l.o.g. we assume that variables occurring in  $\vec{x}_0$  are all distinct.

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  $0 \leq i' \leq l$ .)

In query reformulation, we furthermore 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  $'a_1 \dots a'_n$  and  $'b_1 \dots b'_n$  and returns  $'c_1 \dots c'_n$ , where  $c_i = a_i \wedge b_i$  is the Boolean “and,”  $i = 1, \dots, n$ ;
- INVBIT is a unary function that takes as its input a bit string  $'a_1 \dots a'_n$  and returns  $'c_1 \dots c'_n$ , where  $c_i = \neg a_i$  is the Boolean complement,  $i = 1, \dots, n$ ;
- SUMBIT is an aggregate function such that given  $m$  strings of form  $'b_{i,1} \dots b'_{i,n}$ ,  $i = 1, \dots, m$ , it returns  $'c_1 \dots c'_n$ , where  $c_j = b_{1,j} \vee \dots \vee b_{m,j}$  is the Boolean “or,”  $j = 1, \dots, n$ .

Let  $Q = \langle q, \mathcal{P} \rangle$  be a non-recursive Datalog<sup>-</sup> 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. Let  $a_i$ ,  $1 \leq i \leq n$ , 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  $e(\vec{z})$ )
       $c$  AS  $a_j$ ,                     (for each atom  $y_{0,j} = c$  in  $e(\vec{z})$ )
      ( $p_{1,m}.mark$  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  $p_{1,m}, \dots, p_{l,m}, SQL_{r,l+1}, \dots, SQL_{r,l+k}$ 
WHERE  $p_{i,m}.a_j = p_{i',m}.a_{j'}$ ,      (for each atom  $y_{i,j} = y_{i',j'}$  in  $e(\vec{z})$ ,  $0 < i \leq l$ )
       $p_{i,m}.a_j = c$ ,                (for each atom  $y_{i,j} = c$  in  $e(\vec{z})$ ,  $0 < i \leq l$ )
       $n-p_{i,m}.a_j = p_{i',m}.a_{j'}$ ,   (for each atom  $y_{i,j} = y_{i',j'}$  in  $e(\vec{z})$ ,  $l < i$ )
       $n-p_{i,m}.a_j = c$                 (for each atom  $y_{i,j} = c$  in  $e(\vec{z})$ ,  $l < i$ ).

```

where each  $SQL_{r,h}$ ,  $l < h \leq l + k$ , 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  $e(\vec{z})$ )
         $c$  AS  $a_j$ ,                    (for each atom  $y_{h,j} = c$  in  $e(\vec{z})$ )
        '0 ... 0' AS  $mark$ 
FROM  $p_{1,m}, \dots, p_{l,m}$ 
WHERE  $p_{i,m}.a_j = p_{i',m}.a_{j'}$ ,    (for each atom  $y_{i,j} = y_{i',j'}$  in  $e(\vec{z})$ ,  $0 < i \leq l$ )
       $p_{i,m}.a_j = c$ ,                (for each atom  $y_{i,j} = c$  in  $e(\vec{z})$ ,  $0 < i \leq l$ )
      ROW( $a_1, \dots, a_{k_h}$ ) NOT IN (SELECT  $a_1, \dots, a_{k_h}$  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  $'b_1, \dots, b'_n$  of each tuple  $t$  computed for the relation predicate  $h$ , according to rule  $r$ . More precisely, for  $i \in \{1, \dots, n\}$ ,  $b_i = 1$  if  $t$  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...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...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,  $r_1, \dots, r_\ell$ , 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
  SELECT  $a_1, \dots, a_n, SUMBIT(mark)$ 
  FROM ( $SQL_{r_1}$  UNION ... UNION  $SQL_{r_\ell}$ )
  GROUP BY  $a_1, \dots, a_n$ .
```

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

```
SELECT  $a_1, \dots, a_n$  FROM  $q_m$  WHERE  $mark = '1 \dots 1'$ ,
```

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

**Example C.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*:

$$\begin{aligned} r_1 : q(x) &\leftarrow player(x, y, z), not leader(x); \\ r_2 : leader(x) &\leftarrow team(v, w, x). \end{aligned}$$

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<sup>V,  $\neg$</sup>  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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ , let  $Q$  be a non-recursive Datalog <sup>$\neg$</sup>  query over it, and let  $R_1, \dots, R_n$  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  $\bigcap_{i=1}^n \{t \mid t \in Q[R_i \cup S]\}$ , 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  $r(t) = p_0(t_0) \leftarrow p_1(t_1), \dots, p_l(t_l), \text{not } p_{l+1}(t_{l+1}), \dots, \text{not } p_{l+k}(t_{l+k})$  be a ground instance of  $r$ . Consider  $r'(t) = p_0(t_0) \leftarrow p_1(t_1), \dots, p_l(t_l), \text{not } p_{l+1}(t_{l+1}), \dots, \text{not } p_{l+k}(t_{l+k}), e(t_z)$ , where  $t_z$  is obtained by substituting  $t_{i,j}$  for every variable  $y_{i,j}$  in  $\vec{z}$ , such that  $0 \leq i \leq l+k$  and  $1 \leq j \leq k_i$ . Then,  $r'(t)$  is a ground instance of  $r'$ , since every occurrence of a variable  $y_{i,j}$  in  $e(\vec{z})$  is substituted uniformly. Moreover,  $e(t_z)$  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  $c \neq t_{i,j}$  or  $x_{i,j} = x_{i',j'}$  and  $t_{i,j} \neq t_{i',j'}$ , for some  $0 \leq i, i' \leq l+k$ ,  $1 \leq j \leq k_i$ , and  $1 \leq j' \leq k_{i'}$ ).

( $\supseteq$ ) Let  $r'(t) = p_0(t_0) \leftarrow p_1(t_1), \dots, p_l(t_l), \text{not } p_{l+1}(t_{l+1}), \dots, \text{not } p_{l+k}(t_{l+k}), e(t_z)$  be a ground instance of  $r'$ , such that  $e(t_z)$  is true. Then,  $r(t) = p_0(t_0) \leftarrow p_1(t_1), \dots, p_l(t_l), \text{not } p_{l+1}(t_{l+1}), \dots, \text{not } p_{l+k}(t_{l+k})$  is a ground instance of  $r$ , since the truth of  $e(t_z)$  implies, by construction of  $e(\vec{z})$ , that for all  $0 \leq i, i' \leq l+k$ ,  $1 \leq j \leq k_i$ , and  $1 \leq j' \leq k_{i'}$ , 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  $o(p)$ ,  $1 \leq o(p) \leq |\text{IDB}|$ , the enumeration index assigned to an intensional predicate  $p$ , where  $|\text{IDB}|$  denotes the number of IDB predicates. Then (i)  $o(p) \neq o(p')$  if  $p \neq p'$ ; and (ii) for every rule  $r \in \mathcal{P}$  such that  $\text{head}(r) = p(\vec{x})$ , if  $p'(\vec{x}') \in \text{body}(r)$ , 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  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ , let  $R_1, \dots, R_n$  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  $R_1, \dots, R_n$ , and consider any enumeration,  $o$ , of the IDB predicates of  $\mathcal{P}$  wrt.  $D$  that satisfies (i) and (ii). We first show that for  $1 \leq i \leq n$  and every intensional predicate  $p$ , there exists a tuple  $\langle t, m \rangle$  in the view relation  $p_m$  such that  $m(i) = 1$ , i.e., its mark at position  $i$  is 1, iff  $p(t) \in M_i$ . The proof is by induction on  $o(p)$ .

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

( $\subseteq$ ) Let  $\langle t, m \rangle \in p_m$  and  $m(i) = 1$ . Then by the definition of  $p_m$ , it holds that  $m = \text{SUMBIT}(m_1, \dots, m_j)$ ,  $j \geq 1$ , 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,  $\langle t, m' \rangle \in SQL_{r_j}$ , for some  $1 \leq j \leq l$ , where  $r_1, \dots, r_l$  are all the rules in  $\mathcal{P}$  with head predicate  $p$ . Let

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

Since  $o(p) = 1$ , by Condition (ii) we conclude that  $\text{body}(r_j)$  only contains extensional predicates. By the definition of  $SQL_{r_j}$ , there exist tuples  $\langle t_1, m'_1 \rangle, \dots, \langle t_{j_{l+k}}, m'_{j_{l+k}} \rangle$  that cause the selection of  $\langle t, m' \rangle$ , such that

1.  $\langle t_h, m'_h \rangle \in p_{h_m}$  and  $m'_h(i) = 1$ , for  $1 \leq h \leq j_l$ ,
2.  $\langle t_h, m'_h \rangle \in SQL_{r,h}$  and  $m'_h(i) = 0$ , for  $j_{l+1} \leq h \leq j_{l+k}$ ,
3. the equality conditions  $e(t_z)$  evaluate to true under the corresponding substitutions  $\vec{y}_0|_t$  and  $\vec{y}_h|_{t_h}$ , for  $1 \leq h \leq j_{l+k}$ .

For  $1 \leq h \leq j_l$ , we conclude from Condition 1 that  $p_h(t_h) \in S \cup R_i$  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  $\langle t_h, m'_h \rangle \in p_{h_m}$  and  $m'_h(i) = 0$ , or that there does not exist a mark  $m_h$ , such that  $\langle t_h, m_h \rangle \in p_{h_m}$ . In both cases, we conclude that  $p_h(t_h) \notin S \cup R_i$  and thus, since  $p_h$  is extensional, that  $p_h(t_h) \notin M_i$ . Condition 3 guarantees that  $r_j(t') = p(t) \leftarrow p_1(t_1), \dots, p_{j_l}(t_{j_l}), \text{not } p_{j_{l+1}}(t_{j_{l+1}}), \dots, \text{not } p_{j_{l+k}}(t_{j_{l+k}}), e(t_z)$  is a ground instance of  $r_j$  and  $e(t_z)$  is true. Thus,  $r_j(t')$  is a ground instance of rule  $r_j \in \mathcal{P}$ , such that  $M_i \models \text{body}(r_j(t'))$ . Consequently,  $M_i \models \text{head}(r_j(t'))$ , i.e.,  $p(t) \in M_i$ .

( $\supseteq$ ) Let  $p(t) \in M_i$ . Since  $p$  is intensional, there exists a ground rule  $r(t') = p(t) \leftarrow p_1(t_1), \dots, p_l(t_l), \text{not } p_{l+1}(t_{l+1}), \dots, \text{not } p_{l+k}(t_{l+k}), e(t_z)$  in  $\mathcal{P}$ , such that  $p_h(t_h) \in M_i$ , for  $1 \leq h \leq l$ , and  $p_h(t_h) \notin M_i$ , 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  $\langle t_h, m_h \rangle \in p_{h_m}$  such that  $m_h(i) = 1$ , while for  $l+1 \leq h \leq l+k$  either (a)  $\langle t_h, m_h \rangle \in p_{h_m}$  such that  $m_h(i) = 0$ , or (b) no mark  $m_h$  exists such that  $\langle t_h, m_h \rangle \in p_{h_m}$ . Since  $e(t_z)$  is true for the given ground instance, in Case (a)  $\langle t_h, m_h \rangle$  is also in  $SQL_{r,h}$ ; in Case (b)  $\langle t_h, '0 \dots 0' \rangle \in SQL_{r,h}$ . Thus, for  $l+1 \leq h \leq l+k$ , in any case there exists a tuple  $\langle t_h, m_h \rangle \in SQL_{r,h}$ , such that  $m_h(i) = 0$ . Consequently,  $\langle t, m \rangle \in SQL_r$ , where  $m = m_1 \text{ ANDBIT } \dots \text{ ANDBIT } m_l \text{ ANDBIT INVBIT}(m_{l+1}) \text{ ANDBIT } \dots \text{ ANDBIT INVBIT}(m_{l+k})$ , i.e.,  $m(i) = 1$ . By the definition of view  $p_m$ , we conclude that  $\langle t, m' \rangle \in p_m$  and  $m'(i) = 1$ .

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

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

( $\subseteq$ ) Let  $\langle t, m \rangle \in p_m$  and  $m(i) = 1$ . Then by the definition of  $p_m$ , it holds that  $m = \text{SUMBIT}(m_1, \dots, m_j)$ ,  $j \geq 1$ , 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,  $\langle t, m' \rangle \in SQL_{r_j}$ , for some  $1 \leq j \leq l$ , where  $r_1, \dots, r_l$  are all the rules in  $\mathcal{P}$  with head predicate  $p$ . Let

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

Since  $o(p) = n$ , by Condition (ii) we conclude that  $\text{body}(r_j)$  just contains extensional predicates or intensional predicates  $p'$  such that  $o(p') < n$ . By the definition of  $SQL_{r_j}$ , there exist tuples  $\langle t_1, m'_1 \rangle, \dots, \langle t_{j_{l+k}}, m'_{j_{l+k}} \rangle$  that cause the selection of  $\langle t, m' \rangle$ , such that

1.  $\langle t_h, m'_h \rangle \in p_{h_m}$  and  $m'_h(i) = 1$ , for  $1 \leq h \leq j_l$ ,
2.  $\langle t_h, m'_h \rangle \in SQL_{r,h}$  and  $m'_h(i) = 0$ , for  $j_{l+1} \leq h \leq j_{l+k}$ ,
3. the equality conditions  $e(t_z)$  evaluate to true under the corresponding substitutions  $\vec{y}_0|_t$  and  $\vec{y}_h|_{t_h}$ , for  $1 \leq h \leq j_{l+k}$ .

For  $1 \leq h \leq j_l$ , we conclude from Condition 1 that  $p_h(t_h) \in M_i$  since  $p_h$  is either extensional, i.e.,  $p_h(t_h) \in S \cup R_i$ , or intensional with  $o(p_h) < n$ , i.e., the induction hypothesis applies. For every  $j_{l+1} \leq h \leq j_{l+k}$ , Condition 2 implies, by construction of  $SQL_{r,h}$ , that either (a)  $\langle t_h, m'_h \rangle \in p_{h_m}$  and  $m'_h(i) = 0$ ,

or that (b) there does not exist a mark  $m_h$ , such that  $\langle t_h, m_h \rangle \in p_{h_m}$ . We show that  $p_h(t_h) \notin M_i$  follows for  $j_{l+1} \leq h \leq j_{l+k}$ . If  $p_h$  is extensional, in both cases, we conclude that  $p_h(t_h) \notin M_i$ , since  $p_h(t_h) \notin S \cup R_i$ . For intensional  $p_h$ , towards a contradiction, assume that  $p_h(t_h) \in M_i$ . 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  $\langle t_h, m_h \rangle \in p_{h_m}$ . This contradicts (b), so (a) has to be the case, i.e., also  $\langle t_h, m'_h \rangle \in p_{h_m}$ , and  $m_h \neq m'_h$ . 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,  $\langle t_h, m_h \rangle \in p_{h_m}$  and  $\langle t_h, m'_h \rangle \in p_{h_m}$  implies  $m_h = m'_h$ , contradiction. This proves  $p_h(t_h) \notin M_i$  for  $j_{l+1} \leq h \leq j_{l+k}$ . Eventually, Condition 3 guarantees that  $r_j(t') = p(t) \leftarrow p_1(t_1), \dots, p_{j_l}(t_{j_l}), \text{not} p_{j_{l+1}}(t_{j_{l+1}}), \dots, \text{not} p_{j_{l+k}}(t_{j_{l+k}}), e(t_z)$  is a ground instance of  $r_j$  and  $e(t_z)$  is true. Thus,  $r_j(t')$  is a ground instance of rule  $r_j \in \mathcal{P}$ , such that  $M_i \models \text{body}(r_j(t'))$ . Consequently,  $M_i \models \text{head}(r_j(t'))$  has to hold, i.e.,  $p(t) \in M_i$ .

( $\supseteq$ ) Let  $p(t) \in M_i$ . Since  $p$  is intensional, there exists a ground rule  $r(t') = p(t) \leftarrow p_1(t_1), \dots, p_l(t_l), \text{not} p_{l+1}(t_{l+1}), \dots, \text{not} p_{l+k}(t_{l+k}), e(t_z)$  in  $\mathcal{P}$ , such that  $p_h(t_h) \in M_i$ , for  $1 \leq h \leq l$ , and  $p_h(t_h) \notin M_i$ , for  $l+1 \leq h \leq l+k$ . Furthermore, since  $o(p) = n$ , by Condition (ii)  $p_h$  is either extensional or  $o(p_h) < n$ , for  $1 \leq h \leq l+k$ . Thus, for  $1 \leq h \leq l$  there exist tuples  $\langle t_h, m_h \rangle \in p_{h_m}$  such that  $m_h(i) = 1$ , while for  $l+1 \leq h \leq l+k$  either (a)  $\langle t_h, m_h \rangle \in p_{h_m}$  such that  $m_h(i) = 0$ , or (b) no mark  $m_h$  exists such that  $\langle t_h, m_h \rangle \in p_{h_m}$ . Since  $e(t_z)$  is true for the given ground instance, in Case (a)  $\langle t_h, m_h \rangle$  is also in  $SQL_{r,h}$ ; in Case (b)  $\langle t_h, '0 \dots 0' \rangle \in SQL_{r,h}$ . Thus, for  $l+1 \leq h \leq l+k$ , in any case there exists a tuple  $\langle t_h, m_h \rangle \in SQL_{r,h}$ , such that  $m_h(i) = 0$ . Consequently,  $\langle t, m \rangle \in SQL_r$ , where  $m = m_1 \text{ ANDBIT } \dots \text{ ANDBIT } m_l \text{ ANDBIT } \text{INVBIT}(m_{l+1}) \text{ ANDBIT } \dots \text{ ANDBIT } \text{INVBIT}(m_{l+k})$ , i.e.,  $m(i) = 1$ . By the definition of view  $p_m$ , we conclude that  $\langle t, m' \rangle \in p_m$  and  $m'(i) = 1$ .

This proves  $\langle t, m \rangle \in p_m$  and  $m(i) = 1$  iff  $p(t) \in M_i$  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  $\bigcap_{i=1}^n \{t \mid t \in Q[R_i \cup S]\}$ .

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

( $\supseteq$ ) Let  $t \in \bigcap_{i=1}^n \{t \mid t \in Q[R_i \cup S]\}$ . Then,  $q(t) \in M_i$  since  $\text{SM}(R_i \cup S \cup \mathcal{P}) = \{M_i\}$ , for  $1 \leq i \leq n$ . Thus, there exist marks  $m_i$ , such that  $\langle t, m_i \rangle \in q_m$  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.,  $\langle t, '1 \dots 1' \rangle \in q_m$ . Consequently,  $t \in SQL_Q$ .  $\square$

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  $\text{mark} \neq '0 \dots 0'$ . 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.

## C.2 Scaling the Technique

As mentioned in Section 6.2, a limitation for the scalability of the marking strategy is that safe tuples are required to be marked with the string '11...1', which appears unnecessary given that we know they belong to all the repairs. In fact, this overhead can be avoided as described below. 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  $r_{aff}^{A'} = \{t \mid r(t) \in A\}$  and  $r_{safe}^{S'} = \{t \mid r(t) \in S\}$ . Then, given a query  $Q = \langle q, \mathcal{P} \rangle$ , where  $\mathcal{P}$  is assumed to be normalized, over a schema  $\mathcal{S} = \langle \Psi, \Sigma \rangle$ , we proceed as follows:

- for each rule  $r : h(\vec{x}_0) \leftarrow B(\vec{x})$  of form (12) 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  $B_1(\vec{x}) \vee \dots \vee B_n(\vec{x})$ ;
- 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  $r_i$ , replace every  $p_{safe_{j_m}}$  by  $p_{safe_j}$ , and  $p_{safe_j.mark}$  by '1...1'.

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  $2 \cdot 10 \cdot 2^{10}$  bits = 2.5 KB, independently of the size of  $D$ . Furthermore, by allotting 5 MB ( $= 2 \cdot 20 \cdot 2^{20}$  bits) marking space, the technique may scale up to 20 constraint violations, involving two tuples each.

## D 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 [3], for which  $R_1 \leq_D R_2$  iff  $\Delta(R_1, D) \subseteq \Delta(R_2, D)$ . The only exception is [14], which uses the ordering  $R_1 \sqsubseteq_D R_2$  iff  $R_1 \cap D \supseteq R_2 \cap D$  [13, 14]. However, for the set of integrity constraints and queries considered in [14], the adoption of a different repair ordering is of no concern; as discussed in Section 8, 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 [14]; 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 [34]. 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 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 [34, 23]. A simpler version of GLAV mapping is the one in [25], which allows for combining pure LAV and GAV mapping assertions. 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* [34]. 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 D.3).

## D.1 Logic programs with unstratified negation

The paper [14] 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  $\text{Datalog}^-$ .

More precisely, according to [14], 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  $\text{Datalog}^-$  program  $\mathcal{P}_{KD}$ , whereas views in the mapping, which are expressed as union of conjunctive queries, are cast into a  $\text{Datalog}$  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  $\text{Datalog}^-$  query  $\langle q, \mathcal{P}_Q \cup \mathcal{P}_{KD} \cup \mathcal{P}_M \rangle$ , where  $\langle q, \mathcal{P}_Q \rangle$  is the  $\text{Datalog}$  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 [14], provided in [30], 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  $\mathcal{S}_0$  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  $\text{Datalog}$  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.

$$\begin{aligned}
q(x) &\leftarrow \text{player}(x, y, z) \\
q(x) &\leftarrow \text{team}(v, w, x) \\
\text{player}_D(x, y, z) &\leftarrow s_1(x, y, z, w) \\
\text{team}_D(x, y, z) &\leftarrow s_2(x, y, z) \\
\text{team}_D(x, y, z) &\leftarrow s_3(x, y, z) \\
\text{coach}_D(x, y, z) &\leftarrow s_4(x, y, z) \\
\overline{\text{player}}(x, y, z) &\leftarrow \text{player}_D(x, y, z), \text{ not } \overline{\text{player}}(x, y, z) \\
\overline{\text{player}}(x, y, z) &\leftarrow \text{player}(x, w, z), \text{ player}_D(x, y, z), y \neq w \\
\overline{\text{team}}(x, y, z) &\leftarrow \text{team}_D(x, y, z), \text{ not } \overline{\text{team}}(x, y, z) \\
\overline{\text{team}}(x, y, z) &\leftarrow \text{team}(x, v, w), \text{ team}_D(x, y, z), y \neq v \\
\overline{\text{team}}(x, y, z) &\leftarrow \text{team}(x, v, w), \text{ team}_D(x, y, z), z \neq w \\
\overline{\text{coach}}(x, y, z) &\leftarrow \text{coach}_D(x, y, z), \text{ not } \overline{\text{coach}}(x, y, z) \\
\overline{\text{coach}}(x, y, z) &\leftarrow \text{coach}(x, w, z), \text{ coach}_D(x, y, z), y \neq w \\
\overline{\text{player}}(x, y, z) &\leftarrow \text{player}_D(x, y, z), \text{ coach}(x, w, z) \\
\overline{\text{coach}}(x, y, z) &\leftarrow \text{coach}_D(x, y, z), \text{ team}(z, w, x) \\
\overline{\text{coach}}(x, y, z) &\leftarrow \text{coach}_D(x, y, z), \text{ player}(x, w, z) \\
\overline{\text{team}}(x, y, z) &\leftarrow \text{team}_D(x, y, z), \text{ coach}(z, w, x)
\end{aligned}$$

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  $r^{\text{ret}(\mathcal{I}, D)}$ ; (ii) a relation  $r$  that represents a subset of  $r^{\text{ret}(\mathcal{I}, D)}$  that is consistent with the key constraints and the exclusion dependencies for  $r$ ; (iii) an auxiliary relation  $\bar{r}$ . We can easily see that  $\Pi_{\mathcal{I}_0}(Q) = \mathcal{P}_M \cup \Pi_{\mathcal{S}_0}(Q)$ , where  $\Pi_{\mathcal{S}_0}(Q) = \Pi_{\Sigma_0} \cup \Pi_Q$  is the logic specification for querying the relational database  $\mathcal{S}_0$ , in which  $\Pi_{\Sigma_0} = \mathcal{P}_{KD} \cup \mathcal{P}_{ED}$  and  $\Pi_Q = \mathcal{P}_Q$ .

We point out that in [14], 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 a way such 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.

## D.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 [4]. We recall that this sort of programs, firstly introduced by [32], 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, [4] 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 [4] 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.

$$\begin{aligned}
\neg \text{player}'(x, y, z) \vee \neg \text{player}'(x, y_1, z) &\leftarrow \text{player}(x, y, z), \text{player}(x, y_1, z), y \neq y_1. \\
\neg \text{team}'(x, y, z) \vee \neg \text{team}'(x, y_1, z_1) &\leftarrow \text{team}(x, y, z), \text{team}(x, y_1, z_1), y \neq y_1 \\
\neg \text{team}'(x, y, z) \vee \neg \text{team}'(x, y_1, z_1) &\leftarrow \text{team}(x, y, z), \text{team}(x, y_1, z_1), z \neq z_1 \\
\neg \text{coach}'(x, y, z) \vee \neg \text{coach}'(x, y_1, z) &\leftarrow \text{coach}(x, y, z), \text{coach}(x, y_1, z), y \neq y_1 \\
\neg \text{coach}'(x, y, z) \vee \neg \text{player}'(x, y_1, z) &\leftarrow \text{coach}(x, y, z), \text{player}(x, y_1, z) \\
\neg \text{coach}'(x, y, z) \vee \neg \text{team}'(z, y_1, x) &\leftarrow \text{coach}(x, y, z), \text{team}(z, y_1, x)
\end{aligned}$$

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  $p(x, y) \supseteq q(x, y)$  would be repaired with the rule  $\neg p'(x, y) \vee q'(x, y) \leftarrow p(x, y), \text{not } q(x, y)$ .

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  $p'$ .

### D.3 Programs with Annotation Constants

The paper [6] 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 [31], a non-classical logic in which inconsistencies may be accommodated without trivializing reasoning. The values used in [6] 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:

$$\begin{aligned}
\text{player}(x, y, z, \mathbf{t}^*) &\leftarrow \text{player}(x, y, z, \mathbf{t}_d) \\
\text{player}(x, y, z, \mathbf{t}^*) &\leftarrow \text{player}(x, y, z, \mathbf{t}_a) \\
\text{player}(x, y, z, \mathbf{f}^*) &\leftarrow \text{not player}(x, y, z, \mathbf{t}_d) \\
\text{player}(x, y, z, \mathbf{f}^*) &\leftarrow \text{player}(x, y, z, \mathbf{f}_a) \\
\text{player}(x, y, z, \mathbf{f}_a) \vee \text{player}(x, y_1, z, \mathbf{f}_a) &\leftarrow \text{player}(x, y, z, \mathbf{t}^*), \text{player}(x, y_1, z, \mathbf{t}^*), y \neq y_1.
\end{aligned}$$

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  $\text{coach}(x, y, z) \supset \text{player}(x, y, z)$  (we disregard exclusion dependencies of our running example for a while). This constraint is repaired with the rule

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

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

$$\begin{aligned}
q(x) &\leftarrow \text{player}(x, y, z, \mathbf{t}_a) \vee (\text{player}(x, y, z, \mathbf{t}_d) \wedge \neg \text{player}(x, y, z, \mathbf{f}_a)) \\
q(x) &\leftarrow \text{team}(v, w, x, \mathbf{t}_a) \vee (\text{team}(v, w, x, \mathbf{t}_d) \wedge \neg \text{team}(v, w, x, \mathbf{f}_a)).
\end{aligned}$$

The above rewriting is proposed in [6] for the setting of a single database. In the line of the discussion of data integration systems in the beginning of Appendix D, 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 [11, 12]. 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 [11, 12], 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 [3]. 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 [6], and specific rules for computing the minimal retrieved global databases.

## E Further Experiments

### E.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  $\mathcal{S}_k$ , which contains a relation of the form  $p(x, y)$ , where  $x$  is the key; and,  $\mathcal{S}_e$ , which contains relations of the forms  $p(x, y)$  and  $q(v, w)$ , with an exclusion dependency between attributes  $x$  and  $v$ .

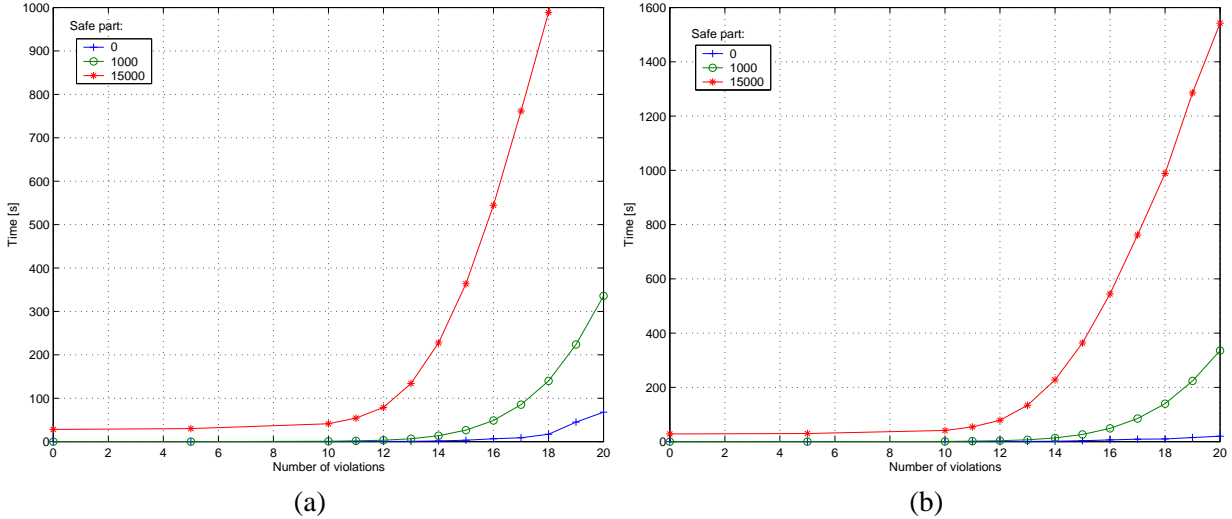


Figure 9: Stable model computation time in DLV system w.r.t. number of conflicts. (a) One Key. (b) One exclusion dependency.

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 unfeasible in practice.

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  $\mathcal{S}_k$  and  $\mathcal{S}_e$ . 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

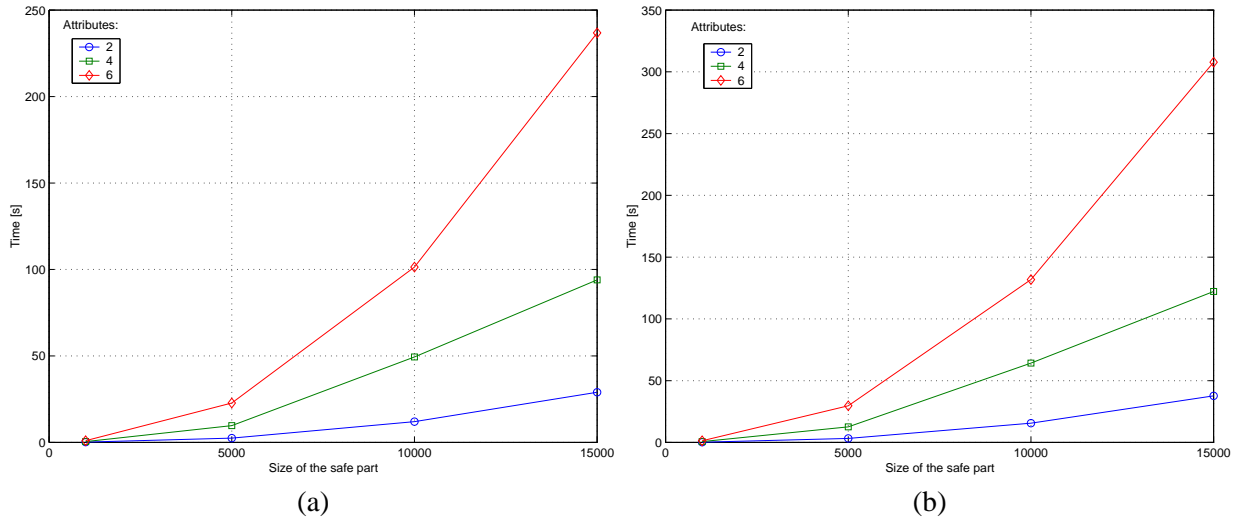


Figure 10: Execution time in DLV system w.r.t. size of the safe part, for different numbers of attributes in relation. (a) One key. (b) One exclusion dependency.

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.

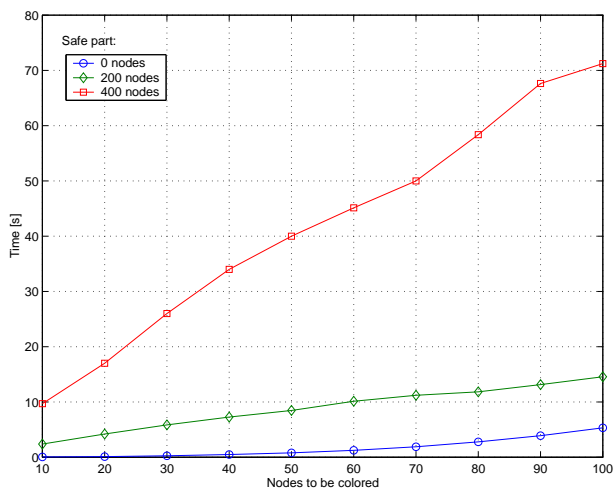
## E.2 3-Coloring

As a further example, we encoded the classical NP-complete graph 3-coloring problem into a consistent query answering problem over a database  $\mathcal{S}_{3c}$  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  $\mathcal{S}_{3c}$ , 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.

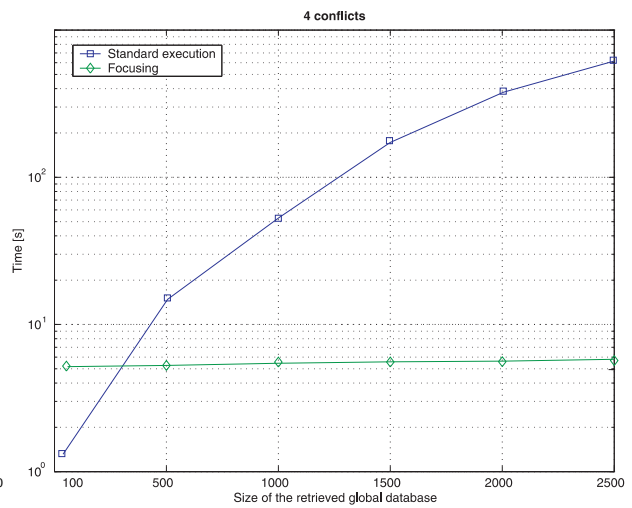
Now, consider the query  $q \leftarrow edge(x, y), colored(x, C), colored(y, C)$ . As easily seen, it evaluates to true on  $D_{3c}$  iff there is no legal 3-coloring for the graph in  $D_{3c}$ .

We encoded the problem of establishing consistent answers to the query  $q$  above over the relational schema  $\mathcal{S}_{3c}$  into a Datalog<sup>-</sup> program, according to the encoding proposed in [13].

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.



(a)



(b)

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