# Preference-Driven Querying of Inconsistent Relational Databases

Slawomir Staworko<sup>1</sup> Jan Chomicki<sup>1</sup> Jerzy Marcinkowski<sup>2</sup>

<sup>1</sup>Department of Computer Science and Engineering University at Buffalo, SUNY

> <sup>2</sup>Institute of Informatics Wroclaw University

IIDB, March 26, 2006

# Motivation

### Schema

Mgr(Name, Dept, Salary, Reports) Key<sub>1</sub>: Name Key<sub>2</sub>: Dept

### $Q_1$ : John earns more than Mary?

?-  $Mgr(John, ..., s_1, ...), Mgr(Mary, ..., s_2, ...), s_1 > s_2.$  $r \models Q_1$ , but is  $Q_1$  really *true*?

### Consistent Query Answers

### Repairs:

 $\begin{aligned} r_1 &= \{(Mary, R\&D, 40K, 3), (John, PR, 30K, 4)\} \\ r_2 &= \{(Mary, IT, 20K, 1), (John, R\&D, 10K, 2)\} \\ r_3 &= \{(Mary, IT, 20K, 1), (John, PR, 30K, 4)\} \\ Q_1 \text{ is not consistently true in } r! \end{aligned}$ 

### What if ???

The user knows: " $s_1$ ,  $s_2$  better than  $s_3$ "

# Motivation (cont.)

### Schema

Mgr(Name, Dept, Salary, Reports) Key<sub>1</sub> : Name Key<sub>2</sub> : Dept

### Data cleaning

- $s_1, s_2$  more reliable than  $s_3$ .
- the clean database:

 $r' = \left\{ \begin{array}{c} (\textit{Mary}, \textit{R\&D}, 40\textit{K}, 3), \\ (\textit{John}, \textit{R\&D}, 10\textit{K}, 2) \end{array} \right\}$ 

r' is inconsistent.

### Preferred Repairs and CQA

 $\begin{array}{l} \textit{Preferred repairs (maximizing reliablility):} \\ r_1 = \{(Mary, R\&D, 40K, 3), (John, PR, 30K, 4)\} \\ r_2 = \{(Mary, IT, 20K, 1), (John, R\&D, 10K, 2)\} \\ \hline r_3 = \{(Mary, IT, 20K, 1), (John, PR, 30K, 4)\} \end{array}$ 

 $\begin{aligned} & Q_2: \text{ Mary earns more for less?} \\ & ?- \text{ Mgr(Mary, _, s_1, r_1), Mgr(John, _, s_2, r_2), s_1 > s_2, r_1 < r_2.} \end{aligned}$ 

r

# Repairs and Consistent Query Answers

### Conflict graph:

- vertices = all tuples
- edges connect conflicting tuples

# $\begin{array}{c|c} R:A \rightarrow B \\ \hline A & B & C \\ \hline 1 & 1 & 1 \\ 1 & 2 & 1 \\ 3 & 3 & 3 \end{array}$



### Repair:

- a maximal consistent subset of the database
- Rep all repairs of the database
- Rep = MIS

### Consistent Query Answers:

answers present in every repair.

 $r_1 = \{(1,2,1), (3,3,3)\}$  $r_2 = \{(1,1,1), (3,3,3)\}$ 

# Priorities, Preferences, and Cleaning

### Priority ≻

(1, 2, 1)

(1, 1, 1)

- an acyclic orientation of the conflict graph
- ► ≻ is total when all edges are oriented

 $\omega_{\succ}(r) = \{t \in r | \neg \exists t' \in r.t' \succ t\}$ 

(3, 3, 3)

 $(1,2,1) \succ (1,1,1)$ 

### Preferred CQA

- ► A-Rep(≻), B-Rep(≻),... different families of preffered repairs w.r.t. ≻
- ➤ X-preferred consistent answers w.r.t. ≻ are the answers present in every X-preferred repair w.r.t ≻

### Database cleaning with a total $\succ$

- while  $\omega_{\succ}(r) \neq \emptyset$  do
  - 1. choose any  $x \in \omega_{\succ}(r)$
  - 2. add x to r'
  - 3. remove x from r with neighbors

# Basic Characterization of Preferred Repairs

### $(\mathcal{P}1)$ Non-emptiness

 $\mathcal{X}$ -Rep $(\succ) \neq \emptyset$ 

### $(\mathcal{P}2)$ Monotonicity

 $\begin{array}{c} \succ_1 \subseteq \succ_2 \\ \Downarrow \\ \mathcal{X}\text{-} \textit{Rep}(\succ_2) \subseteq \mathcal{X}\text{-} \textit{Rep}(\succ_1) \end{array}$ 

### $(\mathcal{P}3)$ Non-discrimination

 $\mathcal{X}$ - $Rep(\emptyset) = Rep$ 

### $(\mathcal{P}4)$ Categoricity

$$\succ$$
 is total  $\Rightarrow |\mathcal{X}\text{-}\mathsf{Rep}(\succ)| = 1$ 

### Trvial family $\mathcal{T}_1$ -Rep $(\succ)$ :

- $1^{\circ}$  if  $\succ$  is total then return the clean database
- 2° otherwise return Rep

$$T_1$$
-Rep satisfies  $\mathcal{P}1 - \mathcal{P}4$ .

# Optimal Use of Priorities



# *L*-*Rep*: Locally Optimal Repairs

### r' is locally optimal iff

no tuple  $x \in r'$  can be replaced with a tuple y such that:

 $y \succ x$ . (and the result is consistent)



 $\mathcal{L}$ -Rep satisfies  $\mathcal{P}1 - \mathcal{P}3$ 







# S-Rep: Semi-globally Optimal Repairs



 $\forall x \in X. y \succ x.$ 



 $\mathcal{S}$ -Rep satisfies  $\mathcal{P}1 - \mathcal{P}3$ 

 $\mathcal{S}\text{-}\mathit{Rep}$  is not categorical (not  $\mathcal{P}4)$ 



CI				
- SI:	awomi	ır S	tawc	nrko
- U.				

# *G*-*Rep*: Globally Optimal Repairs

### r' is globally optimal iff

no set  $X \subseteq r'$  can be replaced with a set Y such that:  $\forall x \in X. \exists y \in Y. y \succ x.$ 



### $\mathcal{G}\text{-}\textit{Rep}$ satisfies $\mathcal{P}1-\mathcal{P}4$



### Alternative characterization

 $\mathcal{G}\text{-}Rep = \ll\text{-maximal repairs}$  $r_1 \ll r_2 \Leftrightarrow \forall x \in r_1 \setminus r_2. \exists y \in r_2 \setminus r_1. y \succ x.$ 



	Panair Chack	Consistent Answers to		
	Repair Check	$\{\forall, \exists\}$ -free queries	conjunctive queries	
Rep	PTIME	PTIME	co-NP-complete	
<i>L</i> -Rep	PTIME	co-NP-complete	co-NP-complete	
S-Rep	PTIME	co-NP-complete	co-NP-complete	
G-Rep	co-NP-complete	$\Pi_p^2$ -complete	$\Pi_p^2$ -complete	

### $\mathcal{L}$ -Rep, $\mathcal{S}$ -Rep, and $\mathcal{G}$ -Rep

For one FD computing consistent answers to  $\{\exists, \forall\}$ -queries is PTIME.

Computing preferred CQA with any family of (semi-globally) optimal repairs satisfying  $\mathcal{P}1$  and  $\mathcal{P}2$  is co-NP-hard. (one atom and 2 FDs)

# C-Rep: Common optimal repairs

### Desired properties:

- optimality to enforce priority use
- monotonicity (P2) to prevent groundless elimination of repairs
- ▶ non-emptiness (𝒫1)

C-Rep - repairs common for all families of (globally) optimal repairs satisfying  $\mathcal{P}1$  and  $\mathcal{P}2$ 

### Database cleaning

- ▶ *r*′ := ∅
- while  $\omega_{\succ}(r) \neq \emptyset$  do
  - 1. choose any  $x \in \omega_{\succ}(r)$
  - 2. add x to r'
  - 3. remove x from r with neighbors
- ▶ return r'

- C-Rep satisfies  $\mathcal{P}1 \mathcal{P}4$
- $\blacktriangleright \ C\text{-}Rep \subseteq \mathcal{G}\text{-}Rep$
- C-Rep = G-Rep for priorities that cannot be extended to a cyclic orientation.
- Repair check: PTIME; CQA: co-NP-c

### Alternative characterization

 $r' \in C\text{-}Rep(\succ)$  iff r' can be a result of cleaning the database with  $\succ$ .

	Repair	Consistent Answers to		Possible
	Check	$\{\forall, \exists\}$ -free queries	conj. queries	Applications
Rep	PTIME	PTIME	co-NP-c	no priorities given
L-Rep	PTIME	co-NP-c		key (no duplicates)
S-Rep	PTIME	co-NP-c		one FD (duplicates)
G-Rep	co-NP-c	$\Pi_p^2$ -c		many FDs with
C-Rep	PTIME	co-NP-c		mutual conflicts

## **Related Work**

### S. Flesca, S. Greco, and E. Zumpano. Active Integrity Constraints.

$$S_\succ(r') = \{(x,y) \in r imes r \mid x \in r'\}$$
  
P- $Rep(\succ) = \{r' \in Rep \mid S_\succ(r') ext{ is maximal}$ 

- ► CQA: Π<sup>p</sup><sub>3</sub>-complete
- ▶ satisfies *P*1 and *P*3
- handles cyclic  $\succ$ , but then

T

▶ violates P2 and P4

### G. Greco and D. Lembo Data Integration with Preferences among Sources.

- repairing a relation by removing tuples has to be *justified* by removing similar tuples from other relations.
- ▶ satisfies P2, but not P1 (non-emptiness)
- ▶ *weakened* framework satisfies *P*1 but *P*2 is lost.