

# Consistent Query Answering

Jan Chomicki  
University at Buffalo

University at Buffalo

- 1 Motivation
- 2 Basics
- 3 Computing CQA
- 4 Computational Complexity
- 5 Variants of CQA
- 6 Conclusions

## Database instance $D$ :

- a finite first-order **structure**
- the **information** about the world

# Integrity constraints (dependencies)

## Database instance $D$ :

- a finite first-order **structure**
- the **information** about the world

## Integrity constraints $\Sigma$ :

- first-order logic **formulas**
- the **properties** of the world

# Integrity constraints (dependencies)

## Database instance $D$ :

- a finite first-order **structure**
- the **information** about the world

## Integrity constraints $\Sigma$ :

- first-order logic **formulas**
- the **properties** of the world

Satisfaction of constraints:  $D \models \Sigma$

Formula **satisfaction** in a first-order structure.

# Integrity constraints (dependencies)

Database instance  $D$ :

- a finite first-order **structure**
- the **information** about the world

Integrity constraints  $\Sigma$ :

- first-order logic **formulas**
- the **properties** of the world

Satisfaction of constraints:  $D \models \Sigma$

Formula **satisfaction** in a first-order structure.

Consistent database:  $D \models \Sigma$

Name	City	Salary
Gates	Redmond	30M
Jobs	Cupertino	10M

Name  $\rightarrow$  City Salary

# Integrity constraints (dependencies)

## Database instance $D$ :

- a finite first-order **structure**
- the **information** about the world

## Integrity constraints $\Sigma$ :

- first-order logic **formulas**
- the **properties** of the world

Satisfaction of constraints:  $D \models \Sigma$

Formula **satisfaction** in a first-order structure.

Consistent database:  $D \models \Sigma$

Name	City	Salary
Gates	Redmond	30M
Jobs	Cupertino	10M

Name  $\rightarrow$  City Salary

**Inconsistent** database:  $D \not\models \Sigma$

Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

Name  $\rightarrow$  City Salary

# Whence Inconsistency?

## Sources of inconsistency:

- **integration** of independent data sources with overlapping data
- time lag of updates (**eventual** consistency)
- unenforced integrity constraints

# Whence Inconsistency?

## Sources of inconsistency:

- **integration** of independent data sources with overlapping data
- time lag of updates (**eventual** consistency)
- unenforced integrity constraints

## Eliminating inconsistency?

- not enough information, time, or money
- difficult, impossible or undesirable
- unnecessary: queries may be **insensitive** to inconsistency

# Whence Inconsistency?

## Sources of inconsistency:

- **integration** of independent data sources with overlapping data
- time lag of updates (**eventual** consistency)
- unenforced integrity constraints

## Eliminating inconsistency?

- not enough information, time, or money
- difficult, impossible or undesirable
- unnecessary: queries may be **insensitive** to inconsistency

## Living with inconsistency?

- **ignoring** inconsistency
- **modifying** the schema
- **exceptions** to constraints.

# Whence Inconsistency?

## Sources of inconsistency:

- **integration** of independent data sources with overlapping data
- time lag of updates (**eventual** consistency)
- unenforced integrity constraints

## Eliminating inconsistency?

- not enough information, time, or money
- difficult, impossible or undesirable
- unnecessary: queries may be **insensitive** to inconsistency

## Living with inconsistency?

- **ignoring** inconsistency
- **modifying** the schema
- **exceptions** to constraints.



Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

Name → City Salary

Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

Name → City Salary

```
SELECT Name  
FROM Employee  
WHERE Salary ≤ 25M
```



# Ignoring Inconsistency

Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

Name → City Salary

```
SELECT Name  
FROM Employee  
WHERE Salary ≤ 25M
```



Name
Gates
Jobs

# Ignoring Inconsistency

Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

Name → City Salary

```
SELECT Name  
FROM Employee  
WHERE Salary ≤ 25M
```



Name
Gates
Jobs

Query results not reliable.

## Decomposition into two relations:

- violators
- the rest

(De Bra, Paredaens [DBP83])



## Decomposition into two relations:

- violators
- the rest

(De Bra, Paredaens [DBP83])



Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

Name → City Salary

## Decomposition into two relations:

- violators
- the rest

(De Bra, Paredaens [DBP83])



Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

$Name \rightarrow City \text{ Salary}$

Jobs	Cupertino	10M
------	-----------	-----

$Name \rightarrow City \text{ Salary}$

Gates	Redmond	20M
Gates	Redmond	30M

$Name \rightarrow City \text{ Salary}$

## Weakening the constraints:

- functional dependencies  $\rightsquigarrow$  denial constraints

(Borgida [Bor85])



## Weakening the constraints:

- functional dependencies  $\rightsquigarrow$  denial constraints

(Borgida [Bor85])



Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

Name  $\rightarrow$  City Salary

## Weakening the constraints:

- functional dependencies  $\rightsquigarrow$  denial constraints

(Borgida [Bor85])



Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

**Name  $\rightarrow$  City Salary**



Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

**Name  $\rightarrow$  City Salary**  
**except Name='Gates'**

## Traditional view

- query results defined irrespective of integrity constraints
- query evaluation may be optimized in the presence of integrity constraints (semantic query optimization)

## Traditional view

- query results defined irrespective of integrity constraints
- query evaluation may be optimized in the presence of integrity constraints (semantic query optimization)

## Our view

- inconsistency leads to **uncertainty**
- query results may depend on integrity constraint satisfaction
- inconsistency may be eliminated (**repairing**) or tolerated (**consistent query answering**)

## Restoring consistency:

- insertion, deletion
- minimal change

## Restoring consistency:

- insertion, deletion
- minimal change

Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

Name → City Salary

## Restoring consistency:

- insertion, deletion
- minimal change

Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

Name → City Salary

Name	City	Salary
Gates	Redmond	30M
Jobs	Cupertino	10M

Name → City Salary

Name	City	Salary
Gates	Redmond	20M
Jobs	Cupertino	10M

Name → City Salary

Consistent query answer:

Query answer obtained in **every**  
**repair**.

(Arenas, Bertossi, Ch. [ABC99])



Consistent query answer:

Query answer obtained in **every**  
**repair**.

(Arenas, Bertossi, Ch. [ABC99])



Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

Name → City Salary

Consistent query answer:

Query answer obtained in **every** repair.

(Arenas, Bertossi, Ch. [ABC99])



Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

Name → City Salary

```
SELECT Name  
FROM Employee  
WHERE Salary ≤ 25M
```

Name
Jobs

Consistent query answer:

Query answer obtained in **every** repair.

(Arenas, Bertossi, Ch. [ABC99])



Name	City	Salary
Gates	Redmond	20M
Gates	Redmond	30M
Jobs	Cupertino	10M

Name  $\rightarrow$  City Salary

```
SELECT Name  
FROM Employee  
WHERE Salary  $\geq$  10M
```

Name
Gates
Jobs

## Formal definition

What constitutes reliable (**consistent**) information in an inconsistent database.

## Formal definition

What constitutes reliable (**consistent**) information in an inconsistent database.

## Algorithms

How to **compute** consistent information.

## Formal definition

What constitutes reliable (**consistent**) information in an inconsistent database.

## Algorithms

How to **compute** consistent information.

## Computational complexity analysis

- **tractable** vs. intractable classes of queries and integrity constraints
- tradeoffs: complexity vs. expressiveness.

## Formal definition

What constitutes reliable (**consistent**) information in an inconsistent database.

## Algorithms

How to **compute** consistent information.

## Computational complexity analysis

- **tractable** vs. intractable classes of queries and integrity constraints
- tradeoffs: complexity vs. expressiveness.

## Implementation

- preferably using **DBMS technology**.

## Formal definition

What constitutes reliable (**consistent**) information in an inconsistent database.

## Algorithms

How to **compute** consistent information.

## Computational complexity analysis

- **tractable** vs. intractable classes of queries and integrity constraints
- tradeoffs: complexity vs. expressiveness.

## Implementation

- preferably using **DBMS technology**.

## Applications

???

Repair  $D'$  of a database  $D$  w.r.t. the integrity constraints  $IC$ :

- $D'$ : over the same schema as  $D$
- $D' \models IC$
- symmetric difference between  $D$  and  $D'$  is **minimal**.

Repair  $D'$  of a database  $D$  w.r.t. the integrity constraints  $IC$ :

- $D'$ : over the same schema as  $D$
- $D' \models IC$
- symmetric difference between  $D$  and  $D'$  is **minimal**.

Consistent query answer to a query  $Q$  in  $D$  w.r.t.  $IC$ :

- an element of the result of  $Q$  in **every repair** of  $D$  w.r.t.  $IC$ .

Repair  $D'$  of a database  $D$  w.r.t. the integrity constraints  $IC$ :

- $D'$ : over the same schema as  $D$
- $D' \models IC$
- symmetric difference between  $D$  and  $D'$  is **minimal**.

Consistent query answer to a query  $Q$  in  $D$  w.r.t.  $IC$ :

- an element of the result of  $Q$  in **every repair** of  $D$  w.r.t.  $IC$ .

Another incarnation of the idea of **sure** query answers  
[Lipski: TODS'79].



## Example relation $R(A, B)$

- violates the dependency  $A \rightarrow B$
- has  $2^n$  repairs.

A	B
$a_1$	$b_1$
$a_1$	$c_1$
$a_2$	$b_2$
$a_2$	$c_2$
...	
$a_n$	$b_n$
$a_n$	$c_n$

$A \rightarrow B$

## Example relation $R(A, B)$

- violates the dependency  $A \rightarrow B$
- has  $2^n$  repairs.

A	B
$a_1$	$b_1$
$a_1$	$c_1$
$a_2$	$b_2$
$a_2$	$c_2$
...	
$a_n$	$b_n$
$a_n$	$c_n$

$A \rightarrow B$

It is impractical to apply the definition of CQA directly.

## Query Rewriting

Given a query  $Q$  and a set of integrity constraints  $IC$ , build a query  $Q^{IC}$  such that for every database instance  $D$

*the set of answers to  $Q^{IC}$  in  $D =$  the set of consistent answers to  $Q$  in  $D$  w.r.t.  $IC$ .*

## Query Rewriting

Given a query  $Q$  and a set of integrity constraints  $IC$ , build a query  $Q^{IC}$  such that for every database instance  $D$

*the set of answers to  $Q^{IC}$  in  $D =$  the set of consistent answers to  $Q$  in  $D$  w.r.t.  $IC$ .*

## Representing all repairs

Given  $IC$  and  $D$ :

- 1 build a space-efficient representation of all repairs of  $D$  w.r.t.  $IC$
- 2 use this representation to answer (many) queries.

## Query Rewriting

Given a query  $Q$  and a set of integrity constraints  $IC$ , build a query  $Q^{IC}$  such that for every database instance  $D$

*the set of answers to  $Q^{IC}$  in  $D =$  the set of consistent answers to  $Q$  in  $D$  w.r.t.  $IC$ .*

## Representing all repairs

Given  $IC$  and  $D$ :

- 1 build a space-efficient representation of all repairs of  $D$  w.r.t.  $IC$
- 2 use this representation to answer (many) queries.

## Logic programs

Given  $IC$ ,  $D$  and  $Q$ :

- 1 build a logic program  $P_{IC,D}$  whose models are the repairs of  $D$  w.r.t.  $IC$
- 2 build a logic program  $P_Q$  expressing  $Q$
- 3 use a logic programming system that computes the query atoms present in **all** models of  $P_{IC,D} \cup P_Q$ .

## Universal constraints

$$\forall. \neg A_1 \vee \dots \vee \neg A_n \vee B_1 \vee \dots \vee B_m$$

## Universal constraints

$$\forall. \neg A_1 \vee \dots \vee \neg A_n \vee B_1 \vee \dots \vee B_m$$

## Example

$$\forall. \neg \text{Par}(x) \vee \text{Ma}(x) \vee \text{Fa}(x)$$

## Universal constraints

$$\forall. \neg A_1 \vee \dots \vee \neg A_n \vee B_1 \vee \dots \vee B_m$$

## Denial constraints

$$\forall. \neg A_1 \vee \dots \vee \neg A_n$$

## Example

$$\forall. \neg Par(x) \vee Ma(x) \vee Fa(x)$$

## Universal constraints

$$\forall. \neg A_1 \vee \dots \vee \neg A_n \vee B_1 \vee \dots \vee B_m$$

## Denial constraints

$$\forall. \neg A_1 \vee \dots \vee \neg A_n$$

## Example

$$\forall. \neg \text{Par}(x) \vee \text{Ma}(x) \vee \text{Fa}(x)$$

## Example

$$\forall. \neg M(n, s, m) \vee \neg M(m, t, w) \vee s \leq t$$

## Universal constraints

$$\forall. \neg A_1 \vee \dots \vee \neg A_n \vee B_1 \vee \dots \vee B_m$$

## Denial constraints

$$\forall. \neg A_1 \vee \dots \vee \neg A_n$$

## Functional dependencies

$X \rightarrow Y$ :

- a **key** dependency in  $F$  if  $Y = U$
- a **primary-key** dependency: only one key exists

## Example

$$\forall. \neg Par(x) \vee Ma(x) \vee Fa(x)$$

## Example

$$\forall. \neg M(n, s, m) \vee \neg M(m, t, w) \vee s \leq t$$

## Universal constraints

$$\forall. \neg A_1 \vee \dots \vee \neg A_n \vee B_1 \vee \dots \vee B_m$$

## Denial constraints

$$\forall. \neg A_1 \vee \dots \vee \neg A_n$$

## Functional dependencies

$X \rightarrow Y$ :

- a **key** dependency in  $F$  if  $Y = U$
- a **primary-key** dependency: only one key exists

## Example

$$\forall. \neg Par(x) \vee Ma(x) \vee Fa(x)$$

## Example

$$\forall. \neg M(n, s, m) \vee \neg M(m, t, w) \vee s \leq t$$

## Example primary-key dependency

Name  $\rightarrow$  Address Salary

## Universal constraints

$$\forall. \neg A_1 \vee \dots \vee \neg A_n \vee B_1 \vee \dots \vee B_m$$

## Denial constraints

$$\forall. \neg A_1 \vee \dots \vee \neg A_n$$

## Functional dependencies

$X \rightarrow Y$ :

- a **key** dependency in  $F$  if  $Y = U$
- a **primary-key** dependency: only one key exists

## Inclusion dependencies

$R[X] \subseteq S[Y]$ :

- a **foreign key** constraint if  $Y$  is a key of  $S$

## Example

$$\forall. \neg Par(x) \vee Ma(x) \vee Fa(x)$$

## Example

$$\forall. \neg M(n, s, m) \vee \neg M(m, t, w) \vee s \leq t$$

## Example primary-key dependency

Name  $\rightarrow$  Address Salary

## Universal constraints

$$\forall. \neg A_1 \vee \dots \vee \neg A_n \vee B_1 \vee \dots \vee B_m$$

## Denial constraints

$$\forall. \neg A_1 \vee \dots \vee \neg A_n$$

## Functional dependencies

$X \rightarrow Y$ :

- a **key** dependency in  $F$  if  $Y = U$
- a **primary-key** dependency: only one key exists

## Inclusion dependencies

$R[X] \subseteq S[Y]$ :

- a **foreign key** constraint if  $Y$  is a key of  $S$

## Example

$$\forall. \neg Par(x) \vee Ma(x) \vee Fa(x)$$

## Example

$$\forall. \neg M(n, s, m) \vee \neg M(m, t, w) \vee s \leq t$$

## Example primary-key dependency

Name  $\rightarrow$  Address Salary

## Example foreign key constraint

$$M[Manager] \subseteq M[Name]$$

Constraints in **clausal form** (disjunctions of literals).

## Residues

- associated with single literals  $p(\bar{x})$  or  $\neg p(\bar{x})$  (one of each for every database relation  $p$ )
- for each literal and each constraint that contains a complementary literal (after renaming), the **local residue** is obtained by removing the complementary literal and the quantifiers for its associated variables
- for each literal, **global residue** = conjunction of local residues.

## Functional dependencies

$$(\forall x, y, z, y', z')(\neg E(x, y, z) \vee \neg E(x, y', z') \vee y = y')$$
$$(\forall x, y, z, y', z')(\neg E(x, y, z) \vee \neg E(x, y', z') \vee z = z')$$

## Query

$$E(x, y, z)$$

## Local residues

$$(\forall y', z')(\neg E(x, y', z') \vee y = y')$$
$$(\forall y', z')(\neg E(x, y', z') \vee z = z')$$

## Literal expansion

For every literal in the original query, construct the expanded version as the conjunction of this literal and its global residue.

## Iteration

The expansion step is iterated by replacing the literals in the residue by their expanded versions, until no changes occur.

## Query expansion

Replace the literals in the query by their final expanded versions.

## Functional dependencies

$$\begin{aligned} &(\forall x, y, z, y', z')(\neg E(x, y, z) \\ &\quad \vee \neg E(x, y', z') \vee y = y') \\ &(\forall x, y, z, y', z')(\neg E(x, y, z) \\ &\quad \vee \neg E(x, y', z') \vee z = z') \end{aligned}$$

## Query

$$E(x, y, z)$$

## Rewritten query

$$\begin{aligned} &E(x, y, z) \wedge (\forall y', z')(\neg E(x, y', z') \vee y = y') \\ &\quad \wedge (\forall y', z')(\neg E(x, y', z') \vee z = z') \end{aligned}$$

## Integrity constraints

$$(\forall x)(\neg P(x) \vee R(x))$$

$$\forall x)(\neg R(x) \vee S(x))$$

Literal	Residue	First expansion	Second (final) expansion
$R(x)$	$S(x)$	$R(x) \wedge S(x)$	$R(x) \wedge S(x)$
$P(x)$	$R(x)$	$P(x) \wedge R(x)$	$P(x) \wedge R(x) \wedge S(x)$
$\neg R(x)$	$\neg P(x)$	$\neg R(x) \wedge \neg P(x)$	$\neg R(x) \wedge \neg P(x)$
$\neg S(x)$	$\neg R(x)$	$\neg S(x) \wedge \neg R(x)$	$\neg S(x) \wedge \neg R(x) \wedge \neg P(x)$

(Arenas, Bertossi, Ch. [ABC99])

- Integrity constraints: **binary universal**
- Queries: **conjunctions** of literals (relational algebra:  $\sigma, \times, -$ )

## (Arenas, Bertossi, Ch. [ABC99])

- Integrity constraints: **binary universal**
- Queries: **conjunctions** of literals (relational algebra:  $\sigma, \times, -$ )

## (Fuxman, Miller [FM07])

- Integrity constraints: **primary key** functional dependencies
- Queries:  $C_{forest}$ 
  - a class of conjunctive queries ( $\pi, \sigma, \times$ )
  - no cycles
  - no non-key or non-full joins
  - no repeated relation symbols
  - no built-ins
- Generalization: conjunctive queries expressed as **rooted rules** (Wijzen [Wij07])

## SQL query

```
SELECT Name FROM Emp  
WHERE Salary  $\geq$  10K
```

## SQL query

```
SELECT Name FROM Emp
WHERE Salary ≥ 10K
```

## SQL rewritten query

```
SELECT e1.Name FROM Emp e1
WHERE e1.Salary ≥ 10K AND NOT EXISTS
  (SELECT * FROM EMPLOYEE e2
   WHERE e2.Name = e1.Name AND e2.Salary < 10K)
```

## SQL query

```
SELECT Name FROM Emp
WHERE Salary  $\geq$  10K
```

## SQL rewritten query

```
SELECT e1.Name FROM Emp e1
WHERE e1.Salary  $\geq$  10K AND NOT EXISTS
  (SELECT * FROM EMPLOYEE e2
   WHERE e2.Name = e1.Name AND e2.Salary < 10K)
```

## (Fuxman, Fazli, Miller [FM05a])

- **ConQuer**: a system for computing CQAs
- conjunctive ( $C_{forest}$ ) and aggregation SQL queries
- databases can be annotated with consistency indicators
- tested on TPC-H queries and medium-size databases

## Vertices

Tuples in the database.

(Gates, Redmond, 20M)

(Grove, Santa Clara, 10M)

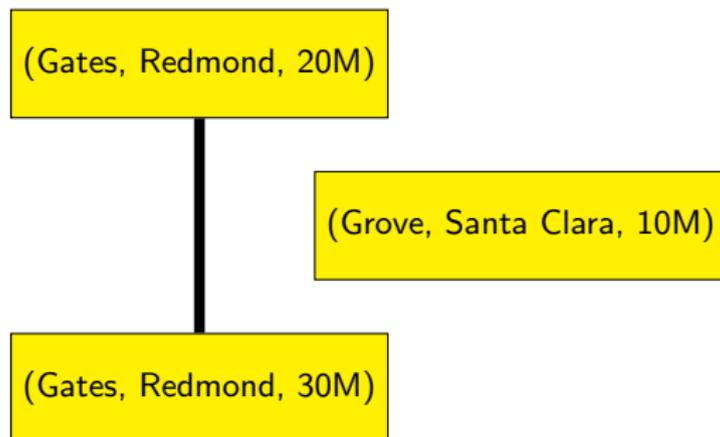
(Gates, Redmond, 30M)

## Vertices

Tuples in the database.

## Edges

Minimal sets of tuples violating a constraint.



## Vertices

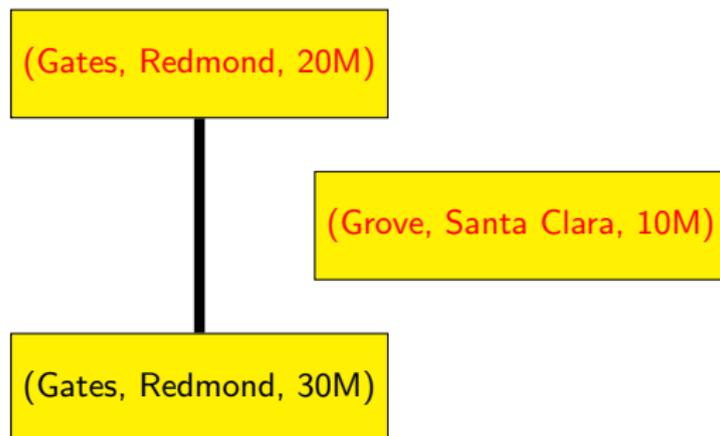
Tuples in the database.

## Edges

Minimal sets of tuples violating a constraint.

## Repairs

Maximal independent sets in the conflict graph.



## Vertices

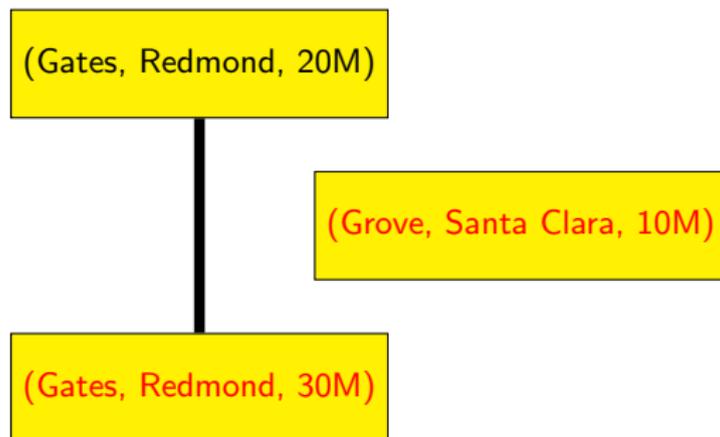
Tuples in the database.

## Edges

Minimal sets of tuples violating a constraint.

## Repairs

Maximal independent sets in the conflict graph.



## Algorithm HProver

INPUT: query  $\Phi$  a disjunction of ground literals, conflict hypergraph  $G$

OUTPUT: is  $\Phi$  false in some repair of  $D$  w.r.t.  $IC$ ?

ALGORITHM:

- 1  $\neg\Phi = P_1(t_1) \wedge \dots \wedge P_m(t_m) \wedge \neg P_{m+1}(t_{m+1}) \wedge \dots \wedge \neg P_n(t_n)$
- 2 find a consistent set of facts  $S$  such that
  - $S \supseteq \{P_1(t_1), \dots, P_m(t_m)\}$
  - for every fact  $A \in \{P_{m+1}(t_{m+1}), \dots, P_n(t_n)\}$ :  $A \notin D$  or there is an edge  $E = \{A, B_1, \dots, B_m\}$  in  $G$  and  $S \supseteq \{B_1, \dots, B_m\}$ .

## Algorithm HProver

INPUT: query  $\Phi$  a disjunction of ground literals, conflict hypergraph  $G$

OUTPUT: is  $\Phi$  false in some repair of  $D$  w.r.t.  $IC$ ?

ALGORITHM:

- ①  $\neg\Phi = P_1(t_1) \wedge \dots \wedge P_m(t_m) \wedge \neg P_{m+1}(t_{m+1}) \wedge \dots \wedge \neg P_n(t_n)$
- ② find a consistent set of facts  $S$  such that
  - $S \supseteq \{P_1(t_1), \dots, P_m(t_m)\}$
  - for every fact  $A \in \{P_{m+1}(t_{m+1}), \dots, P_n(t_n)\}$ :  $A \notin D$  or there is an edge  $E = \{A, B_1, \dots, B_m\}$  in  $G$  and  $S \supseteq \{B_1, \dots, B_m\}$ .

## (Ch., Marcinkowski, Staworko [CMS04])

- **Hippo**: a system for computing CQAs in PTIME
- quantifier-free queries and denial constraints
- only edges of the conflict hypergraph are kept in main memory
- optimization can eliminate many (sometimes all) database accesses in HProver
- tested for medium-size synthetic databases

## Specifying repairs as answer sets of logic programs

- (Arenas, Bertossi, Ch. [ABC03])
- (Greco, Greco, Zumpano [GGZ03])
- (Calì, Lembo, Rosati [CLR03b])

## Specifying repairs as answer sets of logic programs

- (Arenas, Bertossi, Ch. [ABC03])
- (Greco, Greco, Zumpano [GGZ03])
- (Calì, Lembo, Rosati [CLR03b])

## Example

$emp(x, y, z) \leftarrow emp_D(x, y, z), not\ dubious\_emp(x, y, z).$   
 $dubious\_emp(x, y, z) \leftarrow emp_D(x, y, z), emp(x, y', z'), y \neq y'.$   
 $dubious\_emp(x, y, z) \leftarrow emp_D(x, y, z), emp(x, y', z'), z \neq z'.$

## Specifying repairs as answer sets of logic programs

- (Arenas, Bertossi, Ch. [ABC03])
- (Greco, Greco, Zumpano [GGZ03])
- (Calì, Lembo, Rosati [CLR03b])

## Example

$emp(x, y, z) \leftarrow emp_D(x, y, z), not\ dubious\_emp(x, y, z).$   
 $dubious\_emp(x, y, z) \leftarrow emp_D(x, y, z), emp(x, y', z'), y \neq y'.$   
 $dubious\_emp(x, y, z) \leftarrow emp_D(x, y, z), emp(x, y', z'), z \neq z'.$

## Answer sets

- $\{emp(Gates, Redmond, 20M), emp(Grove, SantaClara, 10M), \dots\}$
- $\{emp(Gates, Redmond, 30M), emp(Grove, SantaClara, 10M), \dots\}$

## Logic Programs

- disjunction and classical negation
- checking whether an atom is in all answer sets is  $\Pi_2^P$ -complete
- dlv, smodels, ...

## Logic Programs

- disjunction and classical negation
- checking whether an atom is in all answer sets is  $\Pi_2^P$ -complete
- dlv, smodels, ...

## Scope

- arbitrary first-order queries
- universal constraints
- approach unlikely to yield tractable cases

## Logic Programs

- disjunction and classical negation
- checking whether an atom is in all answer sets is  $\Pi_2^P$ -complete
- `dlv`, `smodels`, ...

## Scope

- arbitrary first-order queries
- universal constraints
- approach unlikely to yield tractable cases

## INFOMIX (Eiter et al. [EFG03])

- combines CQA with data integration (GAV)
- uses `dlv` for repair computations
- optimization techniques: localization, factorization
- tested on small-to-medium-size legacy databases

## Theorem (Ch., Marcinkowski [CM05a])

*For primary-key functional dependencies and conjunctive queries, consistent query answering is **data-complete for co-NP**.*

## Theorem (Ch., Marcinkowski [CM05a])

For primary-key functional dependencies and conjunctive queries, consistent query answering is **data-complete for co-NP**.

## Proof.

**Membership:**  $S$  is a repair iff  $S \models IC$  and  $W \not\models IC$  if  $W = S \cup M$ .

**Co-NP-hardness:** reduction from MONOTONE 3-SAT.

- 1 Positive clauses  $\beta_1 = \phi_1 \wedge \dots \wedge \phi_m$ , negative clauses  $\beta_2 = \psi_{m+1} \wedge \dots \wedge \psi_l$ .
- 2 Database  $D$  contains two binary relations  $R(A, B)$  and  $S(A, B)$ :
  - $R(i, p)$  if variable  $p$  occurs in  $\phi_i$ ,  $i = 1, \dots, m$ .
  - $S(i, p)$  if variable  $p$  occurs in  $\psi_i$ ,  $i = m + 1, \dots, l$ .
- 3  $A$  is the primary key of both  $R$  and  $S$ .
- 4 Query  $Q \equiv \exists x, y, z. (R(x, y) \wedge S(z, y))$ .
- 5 There is an assignment which satisfies  $\beta_1 \wedge \beta_2$  iff there exists a repair in which  $Q$  is false.



## Theorem (Ch., Marcinkowski [CM05a])

For primary-key functional dependencies and conjunctive queries, consistent query answering is **data-complete for co-NP**.

## Proof.

**Membership:**  $S$  is a repair iff  $S \models IC$  and  $W \not\models IC$  if  $W = S \cup M$ .

**Co-NP-hardness:** reduction from MONOTONE 3-SAT.

- 1 Positive clauses  $\beta_1 = \phi_1 \wedge \dots \wedge \phi_m$ , negative clauses  $\beta_2 = \psi_{m+1} \wedge \dots \wedge \psi_l$ .
- 2 Database  $D$  contains two binary relations  $R(A, B)$  and  $S(A, B)$ :
  - $R(i, p)$  if variable  $p$  occurs in  $\phi_i$ ,  $i = 1, \dots, m$ .
  - $S(i, p)$  if variable  $p$  occurs in  $\psi_i$ ,  $i = m + 1, \dots, l$ .
- 3  $A$  is the primary key of both  $R$  and  $S$ .
- 4 Query  $Q \equiv \exists x, y, z. (R(x, y) \wedge S(z, y))$ .
- 5 There is an assignment which satisfies  $\beta_1 \wedge \beta_2$  iff there exists a repair in which  $Q$  is false.



**Q does not belong to  $C_{forest}$ .**

# Data complexity of CQA

	<i>Primary keys</i>	<i>Arbitrary keys</i>	<i>Denial</i>	<i>Universal</i>
$\sigma, \times, -$				
$\sigma, \times, -, \cup$				
$\sigma, \pi$				
$\sigma, \pi, \times$				
$\sigma, \pi, \times, -, \cup$				

# Data complexity of CQA

	<i>Primary keys</i>	<i>Arbitrary keys</i>	<i>Denial</i>	<i>Universal</i>
$\sigma, \times, -$	<b>PTIME</b>	<b>PTIME</b>		<b>PTIME: binary</b>
$\sigma, \times, -, \cup$				
$\sigma, \pi$				
$\sigma, \pi, \times$				
$\sigma, \pi, \times, -, \cup$				

- (Arenas, Bertossi, Ch. [ABC99])

# Data complexity of CQA

	<i>Primary keys</i>	<i>Arbitrary keys</i>	<i>Denial</i>	<i>Universal</i>
$\sigma, \times, -$	PTIME	PTIME	PTIME	PTIME: binary
$\sigma, \times, -, \cup$	PTIME	PTIME	PTIME	
$\sigma, \pi$	PTIME	co-NPC	co-NPC	
$\sigma, \pi, \times$	co-NPC	co-NPC	co-NPC	
$\sigma, \pi, \times, -, \cup$	co-NPC	co-NPC	co-NPC	

- (Arenas, Bertossi, Ch. [ABC99])
- (Ch., Marcinkowski [CM05a])

# Data complexity of CQA

	<i>Primary keys</i>	<i>Arbitrary keys</i>	<i>Denial</i>	<i>Universal</i>
$\sigma, \times, -$	PTIME	PTIME	PTIME	PTIME: binary
$\sigma, \times, -, \cup$	PTIME	PTIME	PTIME	
$\sigma, \pi$	PTIME	co-NPC	co-NPC	
$\sigma, \pi, \times$	co-NPC PTIME: $C_{forest}$	co-NPC	co-NPC	
$\sigma, \pi, \times, -, \cup$	co-NPC	co-NPC	co-NPC	

- (Arenas, Bertossi, Ch. [ABC99])
- (Ch., Marcinkowski [CM05a])
- (Fuxman, Miller [FM07])

# Data complexity of CQA

	<i>Primary keys</i>	<i>Arbitrary keys</i>	<i>Denial</i>	<i>Universal</i>
$\sigma, \times, -$	PTIME	PTIME	PTIME	PTIME: binary $\Pi_2^p$ -complete
$\sigma, \times, -, \cup$	PTIME	PTIME	PTIME	$\Pi_2^p$ -complete
$\sigma, \pi$	PTIME	co-NPC	co-NPC	$\Pi_2^p$ -complete
$\sigma, \pi, \times$	co-NPC PTIME: $C_{forest}$	co-NPC	co-NPC	$\Pi_2^p$ -complete
$\sigma, \pi, \times, -, \cup$	co-NPC	co-NPC	co-NPC	$\Pi_2^p$ -complete

- (Arenas, Bertossi, Ch. [ABC99])
- (Ch., Marcinkowski [CM05a])
- (Fuxman, Miller [FM07])
- (Staworko, Ph.D., 2007), (Staworko, Ch., 2008):
  - quantifier-free queries
  - co-NPC for full TGDs and denial constraints
  - PTIME for **acyclic** full TGDs, join dependencies and denial constraints

## Tuple-based repairs

- asymmetric treatment of insertion and deletion:
  - repairs by minimal deletions only (Ch., Marcinkowski [CM05a]): data possibly **incorrect** but **complete**
  - repairs by minimal deletions and arbitrary insertions (Cali, Lembo, Rosati [CLR03a]): data possibly **incorrect** and **incomplete**
- minimal **cardinality** changes (Lopatenko, Bertossi [LB07]), **more...**

## Tuple-based repairs

- asymmetric treatment of insertion and deletion:
  - repairs by minimal deletions only (Ch., Marcinkowski [CM05a]): data possibly **incorrect** but **complete**
  - repairs by minimal deletions and arbitrary insertions (Cali, Lembo, Rosati [CLR03a]): data possibly **incorrect** and **incomplete**
- minimal **cardinality** changes (Lopatenko, Bertossi [LB07]), **more...**

## Attribute-based repairs

- (A) **ground** and **non-ground** repairs (Wijsen [Wij05])
- (B) **project-join** repairs (Wijsen [Wij06])
- (C) repairs minimizing **Euclidean distance** (Bertossi et al. [BBFL08])
- (D) repairs of minimum **cost** (Bohannon et al. [BFFR05])

## Tuple-based repairs

- asymmetric treatment of insertion and deletion:
  - repairs by minimal deletions only (Ch., Marcinkowski [CM05a]): data possibly **incorrect** but **complete**
  - repairs by minimal deletions and arbitrary insertions (Cali, Lembo, Rosati [CLR03a]): data possibly **incorrect** and **incomplete**
- minimal **cardinality** changes (Lopatenko, Bertossi [LB07]), **more...**

## Attribute-based repairs

- (A) **ground** and **non-ground** repairs (Wijzen [Wij05])
- (B) **project-join** repairs (Wijzen [Wij06])
- (C) repairs minimizing **Euclidean distance** (Bertossi et al. [BBFL08])
- (D) repairs of minimum **cost** (Bohannon et al. [BFFR05])

## Computational complexity

- (A) and (B): similar to tuple based repairs
- (C) and (D): checking existence of a repair of cost  $< K$  NP-complete.

## Tuple-based repairs

- asymmetric treatment of insertion and deletion:
  - repairs by minimal deletions only (Ch., Marcinkowski [CM05a]): data possibly **incorrect** but **complete**
  - repairs by minimal deletions and arbitrary insertions (Cali, Lembo, Rosati [CLR03a]): data possibly **incorrect** and **incomplete**
- minimal **cardinality** changes (Lopatenko, Bertossi [LB07]), **more...**

## Attribute-based repairs

- (A) **ground** and **non-ground** repairs (Wijzen [Wij05])
- (B) **project-join** repairs (Wijzen [Wij06])
- (C) repairs minimizing **Euclidean distance** (Bertossi et al. [BBFL08])
- (D) repairs of minimum **cost** (Bohannon et al. [BFFR05])

## Computational complexity

- (A) and (B): similar to tuple based repairs
- (C) and (D): checking existence of a repair of cost  $< K$  NP-complete.

# The Need for Attribute-based Repairing

Tuple-based repairing leads to **information loss**.

# The Need for Attribute-based Repairing

Tuple-based repairing leads to **information loss**.

*EmpDept*

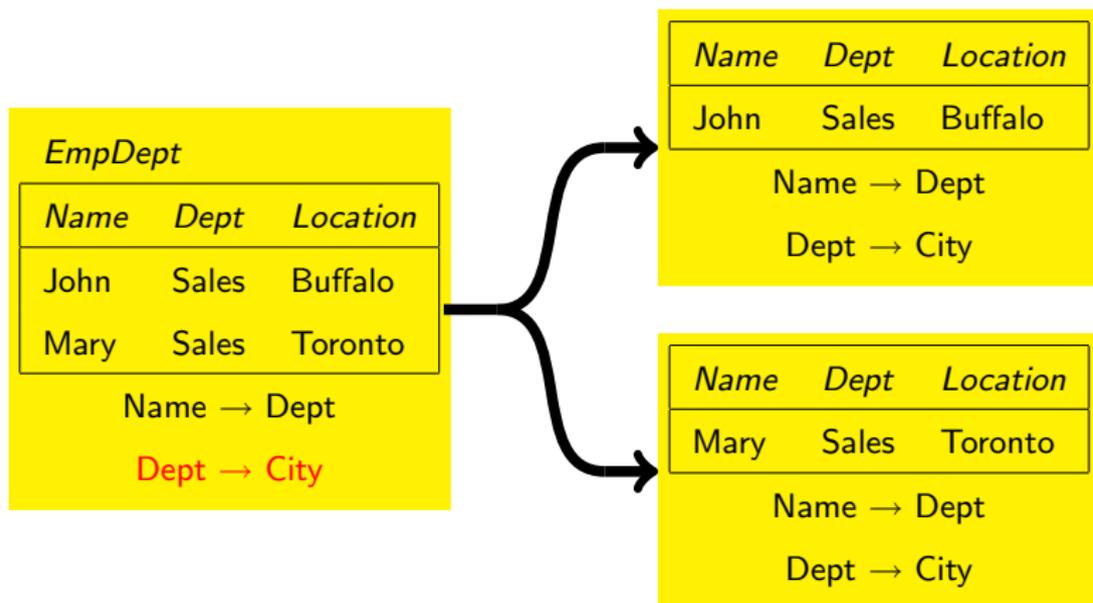
<i>Name</i>	<i>Dept</i>	<i>Location</i>
John	Sales	Buffalo
Mary	Sales	Toronto

Name  $\rightarrow$  Dept

Dept  $\rightarrow$  City

# The Need for Attribute-based Repairing

Tuple-based repairing leads to **information loss**.



Repair the **lossless join decomposition**:

$$\pi_{Name, Dept}(EmpDept) \bowtie \pi_{Dept, Location}(EmpDept)$$

# Attribute-based Repairs through Tuple-based Repairs (Wijsen [Wij06])

Repair the **lossless join decomposition**:

$$\pi_{Name, Dept}(EmpDept) \bowtie \pi_{Dept, Location}(EmpDept)$$

<i>Name</i>	<i>Dept</i>	<i>Location</i>
John	Sales	Buffalo
John	Sales	Toronto
Mary	Sales	Buffalo
Mary	Sales	Toronto

Name  $\rightarrow$  Dept  
Dept  $\rightarrow$  City

# Attribute-based Repairs through Tuple-based Repairs (Wijzen [Wij06])

Repair the **lossless join decomposition**:

$$\pi_{Name, Dept}(EmpDept) \bowtie \pi_{Dept, Location}(EmpDept)$$

<i>Name</i>	<i>Dept</i>	<i>Location</i>
John	Sales	Buffalo
John	Sales	Toronto
Mary	Sales	Buffalo
Mary	Sales	Toronto

Name  $\rightarrow$  Dept  
Dept  $\rightarrow$  City

<i>Name</i>	<i>Dept</i>	<i>Location</i>
John	Sales	Buffalo
Mary	Sales	Buffalo

Name  $\rightarrow$  Dept  
Dept  $\rightarrow$  City

<i>Name</i>	<i>Dept</i>	<i>Location</i>
John	Sales	Toronto
Mary	Sales	Toronto

Name  $\rightarrow$  Dept  
Dept  $\rightarrow$  City

(Andritsos, Fuxman, Miller [AFM06])

- potential **duplicates** identified and grouped into **clusters**
- **worlds**  $\approx$  **repairs**: one tuple from each cluster
- **world probability**: product of tuple probabilities
- **clean answers**: in the query result in some (supporting) world
- **clean answer probability**: sum of the probabilities of supporting worlds
  - **consistent** answer: clean answer **with probability 1**

(Andritsos, Fuxman, Miller [AFM06])

- potential **duplicates** identified and grouped into **clusters**
- **worlds**  $\approx$  **repairs**: one tuple from each cluster
- **world probability**: product of tuple probabilities
- **clean answers**: in the query result in some (supporting) world
- **clean answer probability**: sum of the probabilities of supporting worlds
  - **consistent** answer: clean answer **with probability 1**

## Salaries with probabilities

<i>Name</i>	<i>Salary</i>	<i>Prob</i>
Gates	20M	0.7
Gates	30M	0.3
Jobs	10M	0.5
Jobs	20M	0.5

## SQL query

```
SELECT Name  
FROM EmpProb e  
WHERE e.Salary > 15M
```

## SQL query

```
SELECT Name  
FROM EmpProb e  
WHERE e.Salary > 15M
```

## SQL rewritten query

```
SELECT e.Name,SUM(e.Prob)  
FROM EmpProb e  
WHERE e.Salary > 15M  
GROUP BY e.Name
```

## SQL query

```
SELECT Name
FROM EmpProb e
WHERE e.Salary > 15M
```

## SQL rewritten query

```
SELECT e.Name, SUM(e.Prob)
FROM EmpProb e
WHERE e.Salary > 15M
GROUP BY e.Name
```

### *EmpProb*

<i>Name</i>	<i>Salary</i>	<i>Prob</i>
Gates	20M	0.7
Gates	30M	0.3
Jobs	10M	0.5
Jobs	20M	0.5

Name → Salary

## SQL query

```
SELECT Name
FROM EmpProb e
WHERE e.Salary > 15M
```

## SQL rewritten query

```
SELECT e.Name,SUM(e.Prob)
FROM EmpProb e
WHERE e.Salary > 15M
GROUP BY e.Name
```

### *EmpProb*

<i>Name</i>	<i>Salary</i>	<i>Prob</i>
Gates	20M	0.7
Gates	30M	0.3
Jobs	10M	0.5
Jobs	20M	0.5

Name → Salary

```
SELECT e.Name,SUM(e.Prob)
FROM EmpProb e
WHERE e.Salary > 15M
GROUP BY e.Name
```



## SQL query

```
SELECT Name
FROM EmpProb e
WHERE e.Salary > 15M
```

## SQL rewritten query

```
SELECT e.Name,SUM(e.Prob)
FROM EmpProb e
WHERE e.Salary > 15M
GROUP BY e.Name
```

### *EmpProb*

<i>Name</i>	<i>Salary</i>	<i>Prob</i>
Gates	20M	0.7
Gates	30M	0.3
Jobs	10M	0.5
Jobs	20M	0.5

Name → Salary

```
SELECT e.Name,SUM(e.Prob)
FROM EmpProb e
WHERE e.Salary > 15M
GROUP BY e.Name
```

<i>Name</i>	<i>Prob</i>
Gates	1
Jobs	0.5

## Technology

- **practical methods** for CQA for subsets of SQL:
  - restricted conjunctive/aggregation queries, primary/foreign-key constraints
  - quantifier-free queries, denial constraints/acyclic TGDs/JDs
  - LP-based approaches for expressive query/constraint languages
- (slow) emergence of **generic** techniques
- implemented in **prototype systems**
- tested on **medium-size databases**

## Technology

- **practical methods** for CQA for subsets of SQL:
  - restricted conjunctive/aggregation queries, primary/foreign-key constraints
  - quantifier-free queries, denial constraints/acyclic TGDs/JDs
  - LP-based approaches for expressive query/constraint languages
- (slow) emergence of **generic** techniques
- implemented in **prototype systems**
- tested on **medium-size databases**

## The CQA Community

- over 30 active researchers
- over 100 publications (since 1999)
- at least 8 doctoral dissertations in Europe and North America
- 2007 SIGMOD Doctoral Dissertation Award (Ariel Fuxman)
- overview papers [BC03, Ber06, Cho07, CM05b]

# Taking Stock: Initial Progress

## “Blending in” CQA

- **data integration**: tension between repairing and satisfying source-to-target dependencies

## “Blending in” CQA

- **data integration**: tension between repairing and satisfying source-to-target dependencies

## Extensions

- **uncertainty**:
  - inconsistency leads to but cannot be reduced to uncertainty
  - repairs vs. possible worlds
  - probabilistic data
  - nulls: SQL conformance
- **priorities**:
  - preferred repairs
  - application: conflict resolution
- **XML**
  - notions of integrity constraint and repair
  - repair minimality based on tree edit distance?
- **aggregate** constraints

# Taking Stock: Largely Open Issues

## Applications

- no **deployed** applications
- repairing vs. CQA: data and query **characteristics**
- **heuristics** for CQA and repairing

## Applications

- no **deployed** applications
- repairing vs. CQA: data and query **characteristics**
- **heuristics** for CQA and repairing

## CQA in context

- taming the **semantic explosion**
- CQA and **data cleaning**
- CQA and **schema matching/mapping**

## Applications

- no **deployed** applications
- repairing vs. CQA: data and query **characteristics**
- **heuristics** for CQA and repairing

## CQA in context

- taming the **semantic explosion**
- CQA and **data cleaning**
- CQA and **schema matching/mapping**

## Foundations

- **repair checking**
- defining **measures** of consistency
- proving non-existence of **rewritings**

## Applications

- no **deployed** applications
- repairing vs. CQA: data and query **characteristics**
- **heuristics** for CQA and repairing

## CQA in context

- taming the **semantic explosion**
- CQA and **data cleaning**
- CQA and **schema matching/mapping**

## Foundations

- **repair checking**
- defining **measures** of consistency
- proving non-existence of **rewritings**



# Acknowledgments

Marcelo Arenas  
Alessandro Artale  
Leo Bertossi  
Loreto Bravo  
Andrea Cali  
Thomas Eiter  
Wenfei Fan  
Enrico Franconi  
Ariel Fuxman  
Gianluigi Greco  
Sergio Greco  
Phokion Kolaitis  
Domenico Lembo  
Maurizio Lenzerini  
Jerzy Marcinkowski  
Renée Miller  
Cristian Molinaro  
Riccardo Rosati  
Sławek Staworko  
David Toman  
Jef Wijsen





M. Arenas, L. Bertossi, and J. Chomicki.

Consistent Query Answers in Inconsistent Databases.

In *ACM Symposium on Principles of Database Systems (PODS)*, pages 68–79, 1999.



M. Arenas, L. Bertossi, and J. Chomicki.

Answer Sets for Consistent Query Answering in Inconsistent Databases.

*Theory and Practice of Logic Programming*, 3(4–5):393–424, 2003.



P. Andritsos, A. Fuxman, and R. Miller.

Clean Answers over Dirty Databases.

In *IEEE International Conference on Data Engineering (ICDE)*, 2006.



L. Bertossi, L. Bravo, E. Franconi, and A. Lopatenko.

The Complexity and Approximation of Fixing Numerical Attributes in Databases under Integrity Constraints.

*Information Systems*, 33(4-5):407–434, June-July 2008.



L. Bertossi and J. Chomicki.

Query Answering in Inconsistent Databases.

In J. Chomicki, R. van der Meyden, and G. Saake, editors, *Logics for Emerging Applications of Databases*, pages 43–83. Springer-Verlag, 2003.



L. Bertossi.

Consistent Query Answering in Databases.

*SIGMOD Record*, 35(2), June 2006.



P. Bohannon, M. Flaster, W. Fan, and R. Rastogi.

A Cost-Based Model and Effective Heuristic for Repairing Constraints by Value Modification.

In *ACM SIGMOD International Conference on Management of Data*, pages 143–154, 2005.



A. Borgida.

Language Features for Flexible Handling of Exceptions in Information Systems.

*ACM Transactions on Database Systems*, 10(4):565–603, 1985.



J. Chomicki.

Consistent Query Answering: Five Easy Pieces.

In *International Conference on Database Theory (ICDT)*, pages 1–17. Springer, LNCS 4353, 2007.

Keynote talk.



A. Cali, D. Lembo, and R. Rosati.

On the Decidability and Complexity of Query Answering over Inconsistent and Incomplete Databases.

In *ACM Symposium on Principles of Database Systems (PODS)*, pages 260–271, 2003.



A. Cali, D. Lembo, and R. Rosati.

Query Rewriting and Answering under Constraints in Data Integration Systems.

In *International Joint Conference on Artificial Intelligence (IJCAI)*, pages 16–21, 2003.



J. Chomicki and J. Marcinkowski.

Minimal-Change Integrity Maintenance Using Tuple Deletions.

*Information and Computation*, 197(1-2):90–121, 2005.



J. Chomicki and J. Marcinkowski.

On the Computational Complexity of Minimal-Change Integrity Maintenance in Relational Databases.

In L. Bertossi, A. Hunter, and T. Schaub, editors, *Inconsistency Tolerance*, pages 119–150. Springer-Verlag, 2005.



J. Chomicki, J. Marcinkowski, and S. Staworko.

Computing Consistent Query Answers Using Conflict Hypergraphs.

In *International Conference on Information and Knowledge Management (CIKM)*, pages 417–426. ACM Press, 2004.



P. De Bra and J. Paredaens.

Conditional Dependencies for Horizontal Decompositions.

In *International Colloquium on Automata, Languages and Programming (ICALP)*, pages 123–141, 1983.



T. Eiter, M. Fink, G. Greco, and D. Lembo.

Efficient Evaluation of Logic Programs for Querying Data Integration Systems.

In *International Conference on Logic Programming (ICLP)*, pages 163–177, 2003.



A. Fuxman and R. J. Miller.

ConQuer: Efficient Management of Inconsistent Databases.

In *ACM SIGMOD International Conference on Management of Data*, pages 155–166, 2005.



A. Fuxman and R. J. Miller.

First-Order Query Rewriting for Inconsistent Databases.

In *International Conference on Database Theory (ICDT)*, pages 337–351. Springer, LNCS 3363, 2005.

Full version to appear in JCSS.



A. Fuxman and R. J. Miller.

First-Order Query Rewriting for Inconsistent Databases.

*Journal of Computer and System Sciences*, 73(4):610–635, 2007.



G. Greco, S. Greco, and E. Zumpano.

A Logical Framework for Querying and Repairing Inconsistent Databases.

*IEEE Transactions on Knowledge and Data Engineering*, 15(6):1389–1408, 2003.



A. Lopatenko and L. Bertossi.

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

In *International Conference on Database Theory (ICDT)*, pages 179–193. Springer, LNCS 4353, 2007.



J. Wijsen.

Database Repairing Using Updates.

*ACM Transactions on Database Systems*, 30(3):722–768, 2005.



J. Wijsen.

Project-Join Repair: An Approach to Consistent Query Answering Under Functional Dependencies.

*In International Conference on Flexible Query Answering Systems (FQAS)*, 2006.



J. Wijsen.

On the Consistent Rewriting of Conjunctive Queries Under Primary Key Constraints.

*In International Workshop on Database Programming Languages*, pages 112–126.

Springer, LNCS 4797, 2007.