

# Certain Query Answering in Partially Consistent Databases

Sergio Greco    Fabian Pijcke    Jef Wijsen

Accepted at VLDB 2014



November 29, 2013

## Definition (Uncertain database)

A database in which primary keys need not be satisfied.

## Definition (Repair)

A maximal subset of tuples that satisfy primary keys.

### Example

R	First	Last	Birth	Sal	City	Country
	Ed	Smith	1960	50K	Acri	Italy
	Ed	Smith	1960	60K	Acri	Italy
	An	Allen	1970	40K	Mons	Belgium

S	City	Country	Stars
	Acri	Italy	**
	Acri	Italy	***
	Mons	Belgium	***
	Mons	France	***

## Definition (Uncertain database)

A database in which primary keys need not be satisfied.

## Definition (Repair)

A maximal subset of tuples that satisfy primary keys.

### Example

R	First	Last	Birth	Sal	City	Country
	Ed	Smith	1960	50K	Acri	Italy
	Ed	Smith	1960	60K	Acri	Italy
	An	Allen	1970	40K	Mons	Belgium

S	City	Country	Stars
	Acri	Italy	**
	Acri	Italy	***
	Mons	Belgium	***
	Mons	France	***

## Definition (Uncertain database)

A database in which primary keys need not be satisfied.

## Definition (Repair)

A maximal subset of tuples that satisfy primary keys.

### Example

R	First	Last	Birth	Sal	City	Country
	Ed	Smith	1960	50K	Acri	Italy
	Ed	Smith	1960	60K	Acri	Italy
	An	Allen	1970	40K	Mons	Belgium

S	City	Country	Stars
	Acri	Italy	**
	Acri	Italy	***
	Mons	Belgium	***
	Mons	France	***

## Definition (Uncertain database)

A database in which primary keys need not be satisfied.

## Definition (Repair)

A maximal subset of tuples that satisfy primary keys.

### Example

R	<u>First</u>	<u>Last</u>	Birth	Sal	City	Country
	Ed	Smith	1960	60K	Acri	Italy
	An	Allen	1970	40K	Mons	Belgium

  

S	<u>City</u>	Country	Stars
	Acri	Italy	***
	Mons	Belgium	***

## Definition (The problem $CERTAINTY(q)$ )

For a fixed boolean query  $q$  (parameter of the problem), the problem  $CERTAINTY(q)$  is the following:

**Input:** Uncertain database **db**.

**Question:** Is  $q$  true in every repair of **db**?

## Complexity classification problem

**Input:** Acyclic Boolean conjunctive query  $q$  without self-join.

**Question:** Determine the complexity class of  $CERTAINTY(q)$ , in particular:

- Is  $CERTAINTY(q)$  first-order definable, and hence solvable in SQL?
- Is  $CERTAINTY(q)$  in P?
- Is  $CERTAINTY(q)$  coNP-complete?

## Definition (The problem $CERTAINTY(q)$ )

For a fixed boolean query  $q$  (parameter of the problem), the problem  $CERTAINTY(q)$  is the following:

**Input:** Uncertain database **db**.

**Question:** Is  $q$  true in every repair of **db**?

## Complexity classification problem

**Input:** Acyclic Boolean conjunctive query  $q$  without self-join.

**Question:** Determine the complexity class of  $CERTAINTY(q)$ , in particular:

- Is  $CERTAINTY(q)$  first-order definable, and hence solvable in SQL?
- Is  $CERTAINTY(q)$  in P?
- Is  $CERTAINTY(q)$  coNP-complete?

## Example

R	<u>First</u>	<u>Last</u>	Birth	Sal	City	Country
	Ed	Smith	1960	50K	Acri	Italy
	Ed	Smith	1960	60K	Acri	Italy
	An	Allen	1970	40K	Mons	Belgium

S	<u>City</u>	Country	Stars
	Acri	Italy	**
	Acri	Italy	***
	Mons	Belgium	***
	Mons	France	***

$q_1 = \exists x \exists y S(\underline{x}, \text{Belgium}, y) \Rightarrow$  The answer to *CERTAINTY*( $q_1$ ) is "no"

$q_2 = \exists x \exists y \exists z R(\underline{\text{Ed}}, \underline{\text{Smith}}, 1960, x, \text{Acri}, z) \wedge S(\underline{\text{Acri}}, z, y)$

$\Rightarrow$  The answer to *CERTAINTY*( $q_2$ ) is "yes"

## Example

R	<u>First</u>	<u>Last</u>	Birth	Sal	City	Country
	Ed	Smith	1960	50K	Acri	Italy
	Ed	Smith	1960	60K	Acri	Italy
	An	Allen	1970	40K	Mons	Belgium

S	<u>City</u>	Country	Stars
	Acri	Italy	**
	Acri	Italy	***
	Mons	Belgium	***
	Mons	France	***

$q_1 = \exists x \exists y S(\underline{x}, \text{Belgium}, y) \Rightarrow$  The answer to  $CERTAINTY(q_1)$  is "no"

$q_2 = \exists x \exists y \exists z R(\underline{\text{Ed}}, \underline{\text{Smith}}, 1960, x, \text{Acri}, z) \wedge S(\underline{\text{Acri}}, z, y)$

$\Rightarrow$  The answer to  $CERTAINTY(q_2)$  is "yes"

## Example

R	<u>First</u>	<u>Last</u>	Birth	Sal	City	Country
	Ed	Smith	1960	50K	Acri	Italy
	An	Allen	1970	40K	Mons	Belgium

S	<u>City</u>	Country	Stars
	Acri	Italy	**
	Mons	France	***

$q_1 = \exists x \exists y S(\underline{x}, \text{Belgium}, y) \Rightarrow \text{The answer to } CERTAINTY(q_1) \text{ is "no"}$

$q_2 = \exists x \exists y \exists z R(\underline{\text{Ed}}, \underline{\text{Smith}}, 1960, x, \text{Acri}, z) \wedge S(\underline{\text{Acri}}, z, y)$   
 $\Rightarrow \text{The answer to } CERTAINTY(q_2) \text{ is "yes"}$

## Example

R	<u>First</u>	<u>Last</u>	Birth	Sal	City	Country
	Ed	Smith	1960	50K	Acri	Italy
	Ed	Smith	1960	60K	Acri	Italy
	An	Allen	1970	40K	Mons	Belgium

S	<u>City</u>	Country	Stars
	Acri	Italy	**
	Acri	Italy	***
	Mons	Belgium	***
	Mons	France	***

$q_1 = \exists x \exists y S(\underline{x}, \text{Belgium}, y) \Rightarrow$  The answer to *CERTAINTY*( $q_1$ ) is “no”

$q_2 = \exists x \exists y \exists z R(\underline{\text{Ed}}, \underline{\text{Smith}}, 1960, x, \text{Acri}, z) \wedge S(\underline{\text{Acri}}, z, y)$   
⇒ The answer to *CERTAINTY*( $q_2$ ) is “yes”

## Definition (First-order definable)

$CERTAINTY(q)$  is first-order definable if there exists a first-order sentence  $\phi$  such that for every uncertain database  $\mathbf{db}$ , the following are equivalent:

- $q$  is true in every repair of  $\mathbf{db}$ ; and
- $\phi$  is true in  $\mathbf{db}$ .

We call  $\phi$ , if it exists, a consistent FO rewriting.

## Example

$$q_3 = \exists x \exists y S(\underline{x}, \text{Italy}, y)$$

$$\phi_3 = \exists x \exists y ( S(\underline{x}, \text{Italy}, y) \wedge (\forall z_1 \forall z_2 S(\underline{x}, z_1, z_2) \rightarrow z_1 = \text{Italy}))$$

## Definition (First-order definable)

$CERTAINTY(q)$  is first-order definable if there exists a first-order sentence  $\phi$  such that for every uncertain database  $\mathbf{db}$ , the following are equivalent:

- $q$  is true in every repair of  $\mathbf{db}$ ; and
- $\phi$  is true in  $\mathbf{db}$ .

We call  $\phi$ , if it exists, a consistent FO rewriting.

## Example

$$q_3 = \exists x \exists y \ S(\underline{x}, \text{Italy}, y)$$

$$\phi_3 = \exists x \exists y ( S(\underline{x}, \text{Italy}, y) \wedge (\forall z_1 \forall z_2 \ S(\underline{x}, z_1, z_2) \rightarrow z_1 = \text{Italy}))$$

## Definition (The problem $CERTAINTY(q, \Sigma)$ )

- fixed boolean query  $q$
- fixed set  $\Sigma$  of constraints:
  - functional dependencies
  - at most one key-join dependency per relation
- Input: uncertain database **db** satisfying every constraint in  $\Sigma$
- Question: is  $q$  true in every repair of **db** ?

## Theorem

Given  $\Sigma$  (as above) and an acyclic boolean conjunctive query  $q$  without self-join, it is decidable whether  $CERTAINTY(q, \Sigma)$  is first-order definable. Furthermore, we can construct a consistent FO rewriting if it exists.

## Definition (The problem $CERTAINTY(q, \Sigma)$ )

- fixed boolean query  $q$
- fixed set  $\Sigma$  of constraints:
  - functional dependencies
  - at most one key-join dependency per relation
- Input: uncertain database  $\mathbf{db}$  satisfying every constraint in  $\Sigma$
- Question: is  $q$  true in every repair of  $\mathbf{db}$  ?

## Theorem

Given  $\Sigma$  (as above) and an acyclic boolean conjunctive query  $q$  without self-join, it is decidable whether  $CERTAINTY(q, \Sigma)$  is first-order definable. Furthermore, we can construct a consistent FO rewriting if it exists.

## Definition (The problem $CERTAINTY(q, \Sigma)$ )

- fixed boolean query  $q$
- fixed set  $\Sigma$  of constraints:
  - functional dependencies
  - at most one key-join dependency per relation
- Input: uncertain database **db** satisfying every constraint in  $\Sigma$
- Question: is  $q$  true in every repair of **db** ?

## Theorem

Given  $\Sigma$  (as above) and an acyclic boolean conjunctive query  $q$  without self-join, it is decidable whether  $CERTAINTY(q, \Sigma)$  is first-order definable. Furthermore, we can construct a consistent FO rewriting if it exists.

## Example

$$q_4 = \exists x \exists y \ S(\underline{x}, \text{Italy}, y), \Sigma_4 = \{S : \text{City}_1 \rightarrow \text{Country}_2\}$$

$$\phi_4 = \exists x \exists y ( S(\underline{x}, \text{Italy}, y) \wedge (\forall z_1 \forall z_2 \ S(\underline{x}, z_1, z_2) \rightarrow z_1 = \text{Italy}))$$

## Example

$$q_4 = \exists x \exists y \ S(\underline{x}, \text{Italy}, y), \Sigma_4 = \{S : \text{City}_1 \rightarrow \text{Country}_2\}$$

$$\phi_4 = \exists x \exists y ( \ S(\underline{x}, \text{Italy}, y) \wedge (\forall z_1 \forall z_2 \ S(\underline{x}, z_1, z_2) \rightarrow z_1 = \text{Italy}) )$$

## Example

$$q_5 = \exists u \exists v \exists w \exists x \exists y \exists z \ R(\underline{u}, \underline{v}, w, x, y, z) \wedge S(\underline{y}, z, ***)$$

There exists no consistent FO rewriting for  $q_5$ .

## Example

$$q_6 = \exists u \exists v \exists w \exists x \exists y \exists z \ R(\underline{u}, \underline{v}, w, x, y, z) \wedge S(\underline{y}, z, ***)$$
$$\Sigma_6 = \{R : \text{City}_5 \rightarrow \text{Country}_6\}$$

$$\phi = \exists u \exists v \exists w \exists x \exists y \exists z \left( R(\underline{u}, \underline{v}, v, x, y, z) \wedge \left( \begin{array}{l} \forall w \forall x \forall y \forall z (R(\underline{u}, \underline{v}, w, x, y, z) \rightarrow S(\underline{y}, z, ***)) \wedge \\ \forall z_1 \forall z_2 (S(\underline{y}, z_1, z_2) \rightarrow z_1 = z \wedge z_2 = ***) \end{array} \right) \right)$$

## Example

$$q_5 = \exists u \exists v \exists w \exists x \exists y \exists z R(\underline{u}, \underline{v}, w, x, y, z) \wedge S(\underline{y}, z, ***)$$

There exists no consistent FO rewriting for  $q_5$ .

## Example

$$q_6 = \exists u \exists v \exists w \exists x \exists y \exists z R(\underline{u}, \underline{v}, w, x, y, z) \wedge S(\underline{y}, z, ***)$$
$$\Sigma_6 = \{R : \text{City}_5 \rightarrow \text{Country}_6\}$$

$$\phi = \exists u \exists v \exists w \exists x \exists y \exists z \left( R(\underline{u}, \underline{v}, v, x, y, z) \wedge \left( \begin{array}{l} \forall w \forall x \forall y \forall z (R(\underline{u}, \underline{v}, w, x, y, z) \rightarrow S(\underline{y}, z, ***)) \wedge \\ \forall z_1 \forall z_2 (S(\underline{y}, z_1, z_2) \rightarrow z_1 = z \wedge z_2 = ***) \end{array} \right) \right)$$

## Example

$$q_5 = \exists u \exists v \exists w \exists x \exists y \exists z R(\underline{u}, \underline{v}, w, x, y, z) \wedge S(\underline{y}, z, ***)$$

There exists no consistent FO rewriting for  $q_5$ .

## Example

$$q_6 = \exists u \exists v \exists w \exists x \exists y \exists z R(\underline{u}, \underline{v}, w, x, y, z) \wedge S(\underline{y}, z, ***)$$
$$\Sigma_6 = \{R : \text{City}_5 \rightarrow \text{Country}_6\}$$

$$\phi = \exists u \exists v \exists w \exists x \exists y \exists z \left( R(\underline{u}, \underline{v}, v, x, y, z) \wedge \left( \begin{array}{l} \forall w \forall x \forall y \forall z (R(\underline{u}, \underline{v}, w, x, y, z) \rightarrow S(\underline{y}, z, ***)) \wedge \\ \forall z_1 \forall z_2 (S(\underline{y}, z_1, z_2) \rightarrow z_1 = z \wedge z_2 = ***) \end{array} \right) \right)$$

## Example

$$q_5 = \exists u \exists v \exists w \exists x \exists y \exists z R(\underline{u}, \underline{v}, w, x, y, z) \wedge S(\underline{y}, z, ***)$$

There exists no consistent FO rewriting for  $q_5$ .

## Example

$$q_6 = \exists u \exists v \exists w \exists x \exists y \exists z R(\underline{u}, \underline{v}, w, x, y, z) \wedge S(\underline{y}, z, ***)$$
$$\Sigma_6 = \{R : \text{City}_5 \rightarrow \text{Country}_6\}$$

$$\phi = \exists u \exists v \exists w \exists x \exists y \exists z \left( R(\underline{u}, \underline{v}, v, x, y, z) \wedge \left( \begin{array}{l} \forall w \forall x \forall y \forall z (R(\underline{u}, \underline{v}, w, x, y, z) \rightarrow S(\underline{y}, z, ***)) \wedge \\ \forall z_1 \forall z_2 (S(\underline{y}, z_1, z_2) \rightarrow z_1 = z \wedge z_2 = ***) \end{array} \right) \right)$$

## Original

```
SELECT R.Sal FROM R, S  
WHERE R.City = S.City AND R.Country = S.Country AND S.Stars = '***';
```

## Consistent FO rewriting, encoded in SQL

```
SELECT R1.Sal FROM R AS R1  
WHERE NOT EXISTS (  
    SELECT * FROM R AS R2  
    WHERE R2.First = R1.First AND R2.Last = R1.Last  
    AND (R2.Sal <> R1.Sal OR NOT EXISTS (  
        SELECT * FROM S AS S1  
        WHERE S1.City = R2.City AND S1.Country = R2.Country  
        AND S1.Stars = '***'  
        AND NOT EXISTS (  
            SELECT * FROM S AS S2  
            WHERE S2.City = S1.City  
            AND (S2.Country <> S1.Country OR S2.Stars <> '***')))));
```

