Zhiwei Fan<sup>†</sup>, Paraschos Koutris<sup>†</sup>, Xiating Ouyang<sup>†</sup>, Jef Wijsen<sup>‡</sup> <sup>†</sup>University of Wisconsin – Madison, <sup>‡</sup>University of Mons

### **Consistent Query Answering (CQA)**

While data cleaning is widely adopted to repair the inconsistent/dirty data, finding the "right repair" remains a challenge. Alternatively, the idea of CQA is to compute the answers that are guaranteed to be returned in **all** repairs.



Fig. 1: An example of CQA. The middle layer shows the different sets of answers returned for each possible repair, and the bottom layer explains that the answer "A" is returned in **all** repairs.

$$\mathsf{CQA}_Q(\mathbf{db}^!) := \bigcap_{\mathbf{db}^{\checkmark} \text{ is a repair of } \mathbf{db}^!} Q(\mathbf{db}^{\checkmark})$$

We study inconsistent databases that could violate the primary key constraint: every primary key could correspond to multiple distinct tuples in the database.

**Problem**: CQA(Q), where Q is a SPJ query

**Input**: a database db that violates the primary key constraint

**Output**: the answers **guaranteed** to be returned by Q on all repairs of db

Course		Faculty		
course_id	faculty_id	faculty id	namo	aroa
<b>CS 703</b>	2			
CS 703	5	2	Adam	DB
03703	J	2	Alice	ML
MATH 770	3	5	Roh	DI
<b>MATH 770</b>	7	<b>J</b>		Г <b>L</b>
<u> </u>	0	9	Cathy	DB
63/8/	ð	Q	Carrol	DR
<b>CS 787</b>	9	•	Curror	

Executing the following blue SQL query on the database returns the inconsistent answers CS 703 and CS 787. The rewritten query Q' by adding the red segments would find the consistent answers (i.e., CS 703) that are returned in every repair of the database.

> SELECT DISTINCT Course.c\_id FROM Course, Faculty WHERE Course.f\_id = Faculty.f\_id AND (all f\_id's for the same c\_id appear in Faculty)

Fig. 2: An example of first-order (FO) rewriting Q'.

$$Q'(\mathbf{db}^{!}) = \bigcap_{\mathbf{db}^{\checkmark} \text{ is a repair of } \mathbf{db}^{!}} Q(\mathbf{db}^{\checkmark}) = \mathsf{CQA}_{Q}(\mathbf{db}^{!}).$$

And the classification on the Not all queries admit a first-order rewriting! rewritability remains an open problem.

# LINCQA: FASTER CONSISTENT QUERY ANSWERING WITH LINEAR TIME GUARANTEES

(alphabetical authorship)

## **Acyclic Queries in Linear Time**

Evaluating an acyclic query is a well-studied problem by using hash joins on the join tree. For example, the **Boolean** blue SQL query is acyclic:

q():-Course $(\underline{x}, y)$ , Faculty(y, z, w).

Faculty(y, z, w)Course(<u>x</u>,

Our result



Yannakakis [VLDB'81]

consistent answer

The answer to every **Boolean** acyclic query can be computed in  $O(|\mathbf{db}|)$ .

with a pair-pruning join tree (PPJT)

Queries with projection can be reduced to **Boolean** queries

### **Pair-pruning Join Tree (PPJT)**

We consider acyclic self-join-free SPJ queries (each table name occurs once).

**Definition**: A join tree rooted at some atom is a PPJT if

the root of every subtree is unattacked in the subtree. For example, q has a PPJT:



- Queries on star/snowflake schema (e.g. TPC-H, TPC-DS)
- Two tables
- acyclic queries in  $C_{\text{forest}}$  [Fuxman and Miller, SIGMOD'05]

**Proposition**: If a **Boolean** query q have a pair-pruning join tree (PPJT), then CQA(q) has a first-order rewriting that runs in linear time.







### LinCQA and the Rewriting

LinCQA is a query rewriter that takes as input a SQL query, output its firstorder rewriting.

python3 lincqa.py -i <input.sql/dlog> -o <output.sql/dlog> Using PPJT, the consistent answers can be computed in a bottom-up fashion. q():-Course $(\underline{x}, y)$ , Faculty(y, z, "DB").





 $Course_{join}() :- Course(x, y), \neg Course_{fkey}(x)$  $Course_{fkey}(x) := Course(x, y), \neg Faculty_{join}(y)$  $\forall \mathsf{Child} : \mathsf{Root}_{\mathit{fkey}}(\vec{x}) := \mathsf{Root}(\underline{\vec{x}}, \vec{y}), \neg \mathsf{Child}_{\mathit{join}}(\vec{\alpha})$ 

	×
	×
	×

Child<sub>join</sub>( $\vec{\alpha}$ ) :- Child( $\underline{\vec{u}}, \vec{v}$ ),  $\neg$ Child<sub>fkey</sub>( $\vec{u}$ ) Faculty<sub>join</sub>(y) :- Faculty(y, z, w),  $\neg$ Faculty<sub>fkey</sub>(y) Faculty<sub>*fkey*</sub>(y) :- Faculty(y, z, w),  $w \neq$  "DB"

For queries with projection:

SELECT DISTINCT A1, A2 FROM T WHERE A3 = 42we use PPJT to check for each potential answer  $(a, b) \dots$  in **one** program: SELECT DISTINCT 1 FROM T WHERE A3 = 42 AND A1 = a AND A2 = b if **yes**, then (a, b) is a consistent answer.

### Experiments

We used a 400GB StackOverflow dataset (among others) on SQL Server.



Consistent answers can be computed with no asymptotic overhead



# 14 20 Q5 1250 1245