# LinCQA: Faster Consistent Query Answering with Linear Time Guarantees

Xiating Ouyang [1]

*joint work with* Zhiwei Fan [1,2]    Paris Koutris [1]    Jef Wijsen [3]

University of Wisconsin–Madison [1]

Meta [2]

University of Mons [3]

SIGMOD, Seattle WA, June 18–23 2023

# Primary key constraint (violated)

- Metadata of `stackoverflow.com` as of 02/2021 from Stack Exchange Data Dump
- 551M rows, $\sim$400 GB

| Table | # of rows | inRatio | bSize | # of Attributes |
|---|---|---|---|---|
| Users | 14M | 0% | 1 | 14 |
| Posts | 53M | 0% | 1 | 20 |
| PostHistory | 141M | 0.001% | 4 | 9 |
| Badges | 40M | 0.58% | 941 | 4 |
| Votes | 213M | 30.9% | 1441 | 6 |

inconsistencyRatio = # facts violating PK constraint / # of rows

blockSize = max. # facts with the same PK

# Finding consistent answers

Course

| c_id | f_id |
|------|------|
| **CS 703** | 2 |
| **CS 703** | 5 |
| **CS 787** | 3 |
| **CS 787** | 5 |

CS_Faculty

| f_id | f_name |
|------|--------|
| **2** | Adam |
| **2** | Alice |
| **5** | Bob |

```
SELECT DISTINCT c_id
FROM Course, CS_Faculty
WHERE Course.f_id
      = CS_Faculty.f_id
```

# Finding consistent answers

Course

| c_id | f_id |
|------|------|
| **CS 703** | 2 |
| **CS 703** | 5 |
| **CS 787** | 3 |
| **CS 787** | 5 |

CS_Faculty

| f_id | f_name |
|------|--------|
| **2** | Adam |
| **2** | Alice |
| **5** | Bob |

```
SELECT DISTINCT c_id
FROM Course, CS_Faculty
WHERE Course.f_id
      = CS_Faculty.f_id
```

*Return all classes taught by a CS faculty*

{CS 703, CS 787} . . .

# Finding consistent answers

Course

| c_id | f_id |
|--------|------|
| **CS 703** | 2 |
| **CS 703** | 5 |
| **CS 787** | 3 |
| **CS 787** | 5 |

CS_Faculty

| f_id | f_name |
|------|--------|
| **2** | Adam |
| **2** | Alice |
| **5** | Bob |

```
SELECT DISTINCT c_id
FROM Course, CS_Faculty
WHERE Course.f_id
      = CS_Faculty.f_id
```

*Return all classes taught by a CS faculty*
{CS 703, CS 787} . . .

Data cleaning

# Finding consistent answers

Course

| c_id | f_id |
|------|------|
| **CS 703** | 2 |
| **CS 703** | 5 |
| **CS 787** | 3 |
| **CS 787** | 5 |

CS_Faculty

| f_id | f_name |
|------|--------|
| **2** | Adam |
| **2** | Alice |
| **5** | Bob |

```
SELECT DISTINCT c_id
FROM Course, CS_Faculty
WHERE Course.f_id
      = CS_Faculty.f_id
```

*Return all classes taught by a CS faculty*

{CS 703, CS 787} ...

Data cleaning $2 \times 2 \times 2 \times 1$ repairs

# Finding consistent answers

Course

| c_id | f_id |
|------|------|
| **CS 703** | 2 |
| **CS 703** | 5 |
| **CS 787** | 3 |
| **CS 787** | 5 |

CS_Faculty

| f_id | f_name |
|------|--------|
| **2** | Adam |
| **2** | Alice |
| **5** | Bob |

```
SELECT DISTINCT c_id
FROM Course, CS_Faculty
WHERE Course.f_id
      = CS_Faculty.f_id
```

*Return all classes taught by a CS faculty*

{CS 703, CS 787} . . .

Data cleaning                2 × 2 × 2 × 1 repairs

What are the answers  guaranteed to be returned on all repairs ?

# Finding consistent answers

Course

| c_id | f_id |
|------|------|
| **CS 703** | 2 |
| **CS 703** | 5 |
| **CS 787** | 3 |
| **CS 787** | 5 |

CS_Faculty

| f_id | f_name |
|------|--------|
| **2** | Adam |
| **2** | Alice |
| **5** | Bob |

```
SELECT DISTINCT c_id
FROM Course, CS_Faculty
WHERE Course.f_id
      = CS_Faculty.f_id
```

*Return all classes taught by a CS faculty*

$\{$CS 703, CS 787$\}$ ...

Data cleaning                    $2 \times 2 \times 2 \times 1$ repairs

What are the answers │ guaranteed to be returned on all repairs │?

**CS 703**

# Finding consistent answers

Course

| c_id | f_id |
|------|------|
| **CS 703** | 2 |
| **CS 703** | 5 |
| **CS 787** | 3 |
| **CS 787** | 5 |

CS_Faculty

| f_id | f_name |
|------|--------|
| **2** | Adam |
| **2** | Alice |
| **5** | Bob |

```
SELECT DISTINCT c_id
FROM Course, CS_Faculty
WHERE Course.f_id
      = CS_Faculty.f_id
```

*Return all classes taught by a CS faculty*

$\{CS\ 703,\ CS\ 787\}$ ...

Data cleaning        $2 \times 2 \times 2 \times 1$ repairs

What are the answers | guaranteed to be returned on <u>all</u> repairs |?

## CS 703

| Consistent Answer |

# Finding consistent answers without enumeration

Course

| c_id | f_id |
|------|------|
| **CS 703** | 2 |
| **CS 703** | 5 |
| **CS 787** | 3 |
| **CS 787** | 5 |

CS_Faculty

| f_id | f_name |
|------|--------|
| **2** | Adam |
| **2** | Alice |
| **5** | Bob |

```
SELECT DISTINCT c_id
FROM Course, CS_Faculty
WHERE Course.f_id
      = CS_Faculty.f_id  AND
  (all f_id's for the same c_id
  appear in CS_Faculty)
```

*The original query has a first-order rewriting*

# Finding consistent answers without enumeration

Course

| c_id | f_id |
|------|------|
| **CS 703** | 2 |
| **CS 703** | 5 |
| **CS 787** | 3 |
| **CS 787** | 5 |

CS_Faculty

| f_id | f_name |
|------|--------|
| **2** | Adam |
| **2** | Alice |
| **5** | Bob |

```
SELECT DISTINCT c_id
FROM Course, CS_Faculty
WHERE Course.f_id
      = CS_Faculty.f_id  AND
   (all f_id's for the same c_id
   appear in CS_Faculty)
```

*The original query has a first-order rewriting*

# Finding consistent answers without enumeration

| Course | |
| --- | --- |
| c_id | f_id |
| **CS 703** | 2 |
| **CS 703** | 5 |
| **CS 787** | 3 |
| **CS 787** | 5 |

| CS_Faculty | |
| --- | --- |
| f_id | f_name |
| **2** | Adam |
| **2** | Alice |
| **5** | Bob |

```
SELECT DISTINCT c_id
FROM Course, CS_Faculty
WHERE Course.f_id
      = CS_Faculty.f_id  AND
   (all f_id's for the same c_id
   appear in CS_Faculty)
```

*The original query* has a first-order rewriting

# Finding consistent answers without enumeration

Course

| c_id | f_id |
|------|------|
| **CS 703** | 2 |
| **CS 703** | 5 |
| **CS 787** | 3 |
| **CS 787** | 5 |

CS_Faculty

| f_id | f_name |
|------|--------|
| **2** | Adam |
| **2** | Alice |
| **5** | Bob |

```
SELECT DISTINCT c_id
FROM Course, CS_Faculty
WHERE Course.f_id
      = CS_Faculty.f_id  AND
  (all f_id's for the same c_id
  appear in CS_Faculty)
```

*The original query has a first-order rewriting*

For which $Q$ can the consistent answers be found efficiently?

Can we build a system to find the consistent answers?

For which $Q$ can the consistent answers be found efficiently?

Can we build a system to find the consistent answers?

$q() :\!\text{-}\ \text{Course}(\underline{x}, y), \text{CS\_Faculty}(\underline{y}, z)$

$q()$ :- $\text{Course}(\underline{x}, y), \text{CS\_Faculty}(\underline{y}, z)$

$q() :- \text{Course}(\underline{x}, y), \text{CS\_Faculty}(\underline{y}, z)$

Yannakakis [VLDB'81]

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

Yannakakis [VLDB'81]          Our result

consistent answer

The ~~answer~~ to every **Boolean** acyclic query can be computed in $O(|\mathbf{db}|)$.
$$\wedge$$
with a pair-pruning join tree (PPJT)

Yannakakis [VLDB'81]          Our result

consistent answer

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

$$\wedge$$

with a pair-pruning join tree (PPJT)

non-Boolean $\leq^P_T$ **Boolean**

# Pair-pruning join tree (PPJT)

A join tree rooted at some atom is a PPJT if

the root of every subtree is <u>unattacked</u> in the subtree

$$\text{Course}(\underline{x}, y)$$

$$\text{CS\_Faculty}(\underline{y}, z)$$

# Pair-pruning join tree (PPJT)

A join tree rooted at some atom is a PPJT if

> the root of every subtree is <u>unattacked</u> in the subtree

$$\text{Course}(\underline{x}, y)$$

$$\downarrow$$

$$\text{CS\_Faculty}(\underline{y}, z)$$

# Pair-pruning join tree (PPJT)

A join tree rooted at some atom is a PPJT if

the root of every subtree is <u>unattacked</u> in the subtree

# Pair-pruning join tree (PPJT)

A join tree rooted at some atom is a PPJT if

the root of every subtree is <u>unattacked</u> in the subtree

# Pair-pruning join tree (PPJT)

A join tree rooted at some atom is a PPJT if

the root of every subtree is <u>unattacked</u> in the subtree



Every acyclic query has a join tree, but not every acyclic query has a PPJT

# PPJT is a wide class

+ star/snowflake schema (e.g. TPC-H)
+ Every acyclic query in $\mathcal{C}_{\text{forest}}$ [FM, ICDT'05] has a PPJT

*Remove a primary key if some tuple with this primary key is "bad"*

$Course(\underline{x}, y)$

$Course_{join}() :- Course(x, y), \neg Course_{fkey}(x)$

$Course_{fkey}(x) :- Course(x, y), \neg Faculty_{join}(y)$

$\forall Child : Root_{fkey}(\bar{x}) :- Root(\underline{\bar{x}}, \bar{y}), \neg Child_{join}(\bar{a})$

$Faculty(\underline{y}, z, \text{"DB"})$

$Child_{join}(\bar{a}) :- Child(\underline{\bar{u}}, \bar{v}), \neg Child_{fkey}(\bar{v})$

$Faculty_{join}(x) :- Faculty(y, z, w), \neg Faculty_{fkey}(y)$

$Faculty_{fkey}(y) :- Faculty(y, z, w), w \neq \text{"DB"}$

also expressible in SQL!

runs in $O(N)$

*Remove a primary key if some tuple with this primary key is "bad"*

Course($\underline{x}$, $y$)

$\downarrow$

Faculty($\underline{y}$, $z$, "DB")



Course$_{join}$() :- Course($x$, $y$), ¬Course$_{fkey}$($x$)

Course$_{fkey}$($x$) :- Course($x$, $y$), ¬Faculty$_{join}$($y$)

∀Child : Root$_{fkey}$($\bar{x}$) :- Root($\underline{\bar{x}}$, $\bar{y}$), ¬Child$_{join}$($\bar{u}$)

Child$_{join}$($\bar{u}$) :- Child($\underline{\bar{u}}$, $\bar{v}$), ¬Child$_{fkey}$($\bar{v}$)

Faculty$_{join}$($y$) :- Faculty($y$, $z$, $w$), ¬Faculty$_{fkey}$($y$)

Faculty$_{fkey}$($y$) :- Faculty($y$, $z$, $w$), $w \neq$ "DB"

also expressible in SQL!

runs in $O(N)$

*Remove a primary key if some tuple with this primary key is "bad"*



Course($\underline{x}$, $y$)

$\downarrow$

Faculty($\underline{y}$, $z$, "DB")

Course$_{join}$() :- Course($x$, $y$), ¬Course$_{fkey}$($x$)

Course$_{fkey}$($x$) :- Course($x$, $y$), ¬Faculty$_{join}$($y$)

∀Child : Root$_{fkey}$($\bar{x}$) :- Root($\underline{\bar{x}}$, $\bar{y}$), ¬Child$_{join}$($\bar{u}$)

Child$_{join}$($\bar{u}$) :- Child($\underline{\bar{u}}$, $\bar{v}$), ¬Child$_{fkey}$($\bar{v}$)

Faculty$_{join}$($y$) :- Faculty($y$, $z$, $w$), ¬Faculty$_{fkey}$($y$)

Faculty$_{fkey}$($y$) :- Faculty($y$, $z$, $w$), $w \neq$ "DB"

also expressible in SQL!

runs in $O(N)$

*Remove a primary key if some tuple with this primary key is "bad"*



Course($\underline{x}$, $y$)

Faculty($\underline{y}$, $z$, "DB")

$\text{Course}_{join}() \text{ :- } \text{Course}(x, y), \neg \text{Course}_{fkey}(x)$

$\text{Course}_{fkey}(x) \text{ :- } \text{Course}(x, y), \neg \text{Faculty}_{join}(y)$

$\forall \text{Child} : \text{Root}_{fkey}(\bar{x}) \text{ :- } \text{Root}(\underline{\bar{x}}, \bar{y}), \neg \text{Child}_{join}(\bar{u})$

$\text{Child}_{join}(\bar{\alpha}) \text{ :- } \text{Child}(\underline{\bar{u}}, \bar{v}), \neg \text{Child}_{fkey}(\bar{u})$

$\text{Faculty}_{join}(y) \text{ :- } \text{Faculty}(y, z, w), \neg \text{Faculty}_{fkey}(y)$

$\text{Faculty}_{fkey}(y) \text{ :- } \text{Faculty}(y, z, w), w \neq \text{"DB"}$

also expressible in SQL!

runs in $O(N)$

*Remove a primary key if some tuple with this primary key is "bad"*



Course($\underline{x}$, $y$)

↓

Faculty($\underline{y}$, $z$, "DB")

Course$_{join}$() :- Course($x$, $y$), ¬Course$_{fkey}$($x$)

Course$_{fkey}$($x$) :- Course($x$, $y$), ¬Faculty$_{join}$($y$)

∀Child : Root$_{fkey}$($\bar{x}$) :- Root($\underline{x}$, $\bar{y}$), ¬Child$_{join}$($\bar{u}$)

Child$_{join}$($\bar{\alpha}$) :- Child($\underline{u}$, $\bar{v}$), ¬Child$_{fkey}$($\bar{u}$)

Faculty$_{join}$($y$) :- Faculty($y$, $z$, $w$), ¬Faculty$_{fkey}$($y$)

Faculty$_{fkey}$($y$) :- Faculty($y$, $z$, $w$), $w \neq$ "DB"

also expressible in SQL!

runs in $O(N)$

*Remove a primary key if some tuple with this primary key is "bad"*

Course($\underline{x}$, $y$)

$\downarrow$

Faculty($\underline{y}$, $z$, "DB")

$$\text{Course}_{join}() \text{ :- Course}(x, y), \neg\text{Course}_{fkey}(x)$$
$$\text{Course}_{fkey}(x) \text{ :- Course}(x, y), \neg\text{Faculty}_{join}(y)$$
$$\forall\text{Child : Root}_{fkey}(\bar{x}) \text{ :- Root}(\bar{x}, \bar{y}), \neg\text{Child}_{join}(\vec{\alpha})$$

$$\boxed{\text{Child}_{join}(\vec{\alpha}) \text{ :- Child}(\underline{\vec{u}}, \vec{v}), \neg\text{Child}_{fkey}(\vec{u})}$$
$$\text{Faculty}_{join}(y) \text{ :- Faculty}(y, z, w), \neg\text{Faculty}_{fkey}(y)$$
$$\text{Faculty}_{fkey}(y) \text{ :- Faculty}(y, z, w), w \neq \text{"DB"}$$

also expressible in SQL!
runs in $O(N)$

*Remove a primary key if some tuple with this primary key is "bad"*



$\text{Course}(\underline{x}, y)$

$\text{Faculty}(\underline{y}, z, \text{"DB"})$

$\text{Course}_{join}() :\!- \text{Course}(x, y), \neg\text{Course}_{fkey}(x)$

$\text{Course}_{fkey}(x) :\!- \text{Course}(x, y), \neg\text{Faculty}_{join}(y)$

$\forall \text{Child} : \text{Root}_{fkey}(\bar{x}) :\!- \text{Root}(\underline{\bar{x}}, \bar{y}), \neg\text{Child}_{join}(\bar{\alpha})$

$\boxed{\text{Child}_{join}(\bar{\alpha}) :\!- \text{Child}(\underline{\bar{u}}, \bar{v}), \neg\text{Child}_{fkey}(\bar{u})}$

$\text{Faculty}_{join}(y) :\!- \text{Faculty}(y, z, w), \neg\text{Faculty}_{fkey}(y)$

$\text{Faculty}_{fkey}(y) :\!- \text{Faculty}(y, z, w), w \neq \text{"DB"}$

*also expressible in SQL!*
*runs in $O(N)$*

*Remove a primary key if some tuple with this primary key is "bad"*

Course($\underline{x}$, $y$)

$\downarrow$

Faculty($\underline{y}$, $z$, "DB")



$\text{Course}_{join}() \text{ :- } \text{Course}(x, y), \neg\text{Course}_{fkey}(x)$

$\text{Course}_{fkey}(x) \text{ :- } \text{Course}(x, y), \neg\text{Faculty}_{join}(y)$

$\forall\text{Child} : \text{Root}_{fkey}(\vec{x}) \text{ :- } \text{Root}(\underline{\vec{x}}, \vec{y}), \neg\text{Child}_{join}(\vec{\alpha})$

$\boxed{\text{Child}_{join}(\vec{\alpha}) \text{ :- } \text{Child}(\underline{\vec{u}}, \vec{v}), \neg\text{Child}_{fkey}(\vec{u})}$

$\text{Faculty}_{join}(y) \text{ :- } \text{Faculty}(y, z, w), \neg\text{Faculty}_{fkey}(y)$

$\text{Faculty}_{fkey}(y) \text{ :- } \text{Faculty}(y, z, w), w \neq \text{"DB"}$

also expressible in SQL!

runs in $O(N)$

*Remove a primary key if some tuple with this primary key is "bad"*

Course($\underline{x}$, $y$)



$\mathsf{Course}_{join}() \text{ :- } \mathsf{Course}(x, y), \neg\mathsf{Course}_{fkey}(x)$

$\mathsf{Course}_{fkey}(x) \text{ :- } \mathsf{Course}(x, y), \neg\mathsf{Faculty}_{join}(y)$

$\forall \mathsf{Child} : \mathsf{Root}_{fkey}(\vec{x}) \text{ :- } \mathsf{Root}(\underline{\vec{x}}, \vec{y}), \neg\mathsf{Child}_{join}(\vec{\alpha})$

Faculty($\underline{y}$, $z$, "DB")



$\mathsf{Child}_{join}(\vec{\alpha}) \text{ :- } \mathsf{Child}(\underline{\vec{u}}, \vec{v}), \neg\mathsf{Child}_{fkey}(\vec{u})$

$\mathsf{Faculty}_{join}(y) \text{ :- } \mathsf{Faculty}(y, z, w), \neg\mathsf{Faculty}_{fkey}(y)$

$\mathsf{Faculty}_{fkey}(y) \text{ :- } \mathsf{Faculty}(y, z, w), w \neq \text{"DB"}$

also expressible in SQL!
runs in $O(N)$

*Remove a primary key if some tuple with this primary key is "bad"*

$\text{Course}(\underline{x}, y)$

$\downarrow$

$\text{Faculty}(\underline{y}, z, \text{"DB"})$

$\text{Course}_{join}() \;:\!\!-\; \text{Course}(x, y), \neg\text{Course}_{fkey}(x)$

$\text{Course}_{fkey}(x) \;:\!\!-\; \text{Course}(x, y), \neg\text{Faculty}_{join}(y)$

$\forall \text{Child} : \text{Root}_{fkey}(\vec{x}) \;:\!\!-\; \text{Root}(\underline{\vec{x}}, \vec{y}), \neg\text{Child}_{join}(\vec{\alpha})$

$\text{Child}_{join}(\vec{\alpha}) \;:\!\!-\; \text{Child}(\underline{\vec{u}}, \vec{v}), \neg\text{Child}_{fkey}(\vec{u})$

$\text{Faculty}_{join}(y) \;:\!\!-\; \text{Faculty}(y, z, w), \neg\text{Faculty}_{fkey}(y)$

$\text{Faculty}_{fkey}(y) \;:\!\!-\; \text{Faculty}(y, z, w), w \neq \text{"DB"}$

also expressible in SQL!

runs in $O(N)$

*Remove a primary key if some tuple with this primary key is "bad"*

Course($\underline{x}$, $y$)

$\downarrow$

Faculty($\underline{y}$, $z$, "DB")

$\text{Course}_{join}() :\!- \text{Course}(x, y), \neg\text{Course}_{fkey}(x)$

$\text{Course}_{fkey}(x) :\!- \text{Course}(x, y), \neg\text{Faculty}_{join}(y)$

$\forall\text{Child} : \text{Root}_{fkey}(\vec{x}) :\!- \text{Root}(\underline{\vec{x}}, \vec{y}), \neg\text{Child}_{join}(\vec{\alpha})$

$\text{Child}_{join}(\vec{\alpha}) :\!- \text{Child}(\underline{\vec{u}}, \vec{v}), \neg\text{Child}_{fkey}(\vec{u})$

$\text{Faculty}_{join}(y) :\!- \text{Faculty}(y, z, w), \neg\text{Faculty}_{fkey}(y)$

$\text{Faculty}_{fkey}(y) :\!- \text{Faculty}(y, z, w), w \neq \text{"DB"}$

also expressible in SQL!

runs in $O(N)$

*Remove a primary key if some tuple with this primary key is "bad"*

Course($\underline{x}$, $y$)

$\downarrow$

Faculty($\underline{y}$, $z$, "DB")

$\text{Course}_{join}()$ :- $\text{Course}(x, y), \neg\text{Course}_{fkey}(x)$

$\text{Course}_{fkey}(x)$ :- $\text{Course}(x, y), \neg\text{Faculty}_{join}(y)$

$\forall\text{Child} : \text{Root}_{fkey}(\vec{x})$ :- $\text{Root}(\underline{\vec{x}}, \vec{y}), \neg\text{Child}_{join}(\vec{\alpha})$

$\text{Child}_{join}(\vec{\alpha})$ :- $\text{Child}(\underline{\vec{u}}, \vec{v}), \neg\text{Child}_{fkey}(\vec{u})$

$\text{Faculty}_{join}(y)$ :- $\text{Faculty}(y, z, w), \neg\text{Faculty}_{fkey}(y)$

$\text{Faculty}_{fkey}(y)$ :- $\text{Faculty}(y, z, w), w \neq \text{"DB"}$

also expressible in SQL!

runs in $O(N)$

```
SELECT DISTINCT A1, A2 FROM T WHERE A3 = 42
```

Step 1 Evaluate directly

| A1 | A2 |
|----|----|
| a  | b  |
| x  | y  |
| ...| ...|

Step 2 Reduce to **Boolean** (using PPJT)

```
SELECT DISTINCT 1 FROM T WHERE A3 = 42 AND A1 = a AND A2 = b
```

if **yes**, then output $(a, b)$, otherwise continue

```
SELECT DISTINCT 1 FROM T WHERE A3 = 42 AND A1 = x AND A2 = y
```

...

$\xrightarrow{\text{LinCQA}}$ a single SQL/Datalog query

```
SELECT DISTINCT A1, A2 FROM T WHERE A3 = 42
```

Step 1 Evaluate directly

| A1 | A2 |
|----|----|
| a  | b  |
| x  | y  |
| ... | ... |

Step 2 Reduce to **Boolean** (using PPJT)

```
SELECT DISTINCT 1 FROM T WHERE A3 = 42 AND A1 = a AND A2 = b
```

if **yes**, then output $(a, b)$, otherwise continue

```
SELECT DISTINCT 1 FROM T WHERE A3 = 42 AND A1 = x AND A2 = y
```

...

$\xrightarrow{\text{LinCQA}}$ a single SQL/Datalog query

# From Boolean to non-Boolean

```
SELECT DISTINCT A1, A2 FROM T WHERE A3 = 42
```

**Step 1** Evaluate directly

| A1 | A2 |
|----|----|
| a | b |
| x | y |
| ... | ... |

**Step 2** Reduce to **Boolean** (using PPJT)

```
SELECT DISTINCT 1 FROM T WHERE A3 = 42 AND A1 = a AND A2 = b
```

if **yes**, then output $(a, b)$, otherwise continue

```
SELECT DISTINCT 1 FROM T WHERE A3 = 42 AND A1 = x AND A2 = y
```

...

$\xrightarrow{\text{LinCQA}}$ a single SQL/Datalog query

```
SELECT DISTINCT A1, A2 FROM T WHERE A3 = 42
```

Step 1  Evaluate directly

| A1 | A2 |
|----|----|
| a | b |
| x | y |
| ... | ... |

Step 2  Reduce to **Boolean** (using PPJT)

```
SELECT DISTINCT 1 FROM T WHERE A3 = 42 AND A1 = a AND A2 = b
```

if **yes**, then output $(a, b)$, otherwise continue

```
SELECT DISTINCT 1 FROM T WHERE A3 = 42 AND A1 = x AND A2 = y
```

...

$\xrightarrow{\text{LinCQA}}$ a single SQL/Datalog query

# From Boolean to non-Boolean

```
SELECT DISTINCT A1, A2 FROM T WHERE A3 = 42
```

Step 1 Evaluate directly

| A1 | A2 |
|----|----|
| a  | b  |
| x  | y  |
| ...| ...|

Step 2 Reduce to **Boolean** (using PPJT)

```
SELECT DISTINCT 1 FROM T WHERE A3 = 42 AND A1 = a AND A2 = b
```

if **yes**, then output $(a, b)$, otherwise continue

```
SELECT DISTINCT 1 FROM T WHERE A3 = 42 AND A1 = x AND A2 = y
```

...

$$\xrightarrow{\text{LinCQA}} \text{a single SQL/Datalog query}$$

| Acyclic $q$ | PPJT | Yannakakis [VLDB'81] |
|---|---|---|
| Boolean $q$ | $O(N)$ | $O(N)$ |
| non-Boolean $q$ | $O(N \cdot |\text{OUT}_{\text{inconsistent}}|)$ | $O(N \cdot |\text{OUT}|)$ |
| full $q$ (SELECT *) | $O(N + |\text{OUT}_{\text{consistent}}|)$ | $O(N + |\text{OUT}|)$ |

Consistent answers of common join queries can be computed with no asymptotic overhead

| Acyclic $q$ | PPJT | Yannakakis [VLDB'81] |
|---|---|---|
| Boolean $q$ | $O(N)$ | $O(N)$ |
| non-Boolean $q$ | $O(N \cdot |\text{OUT}_{\text{inconsistent}}|)$ | $O(N \cdot |\text{OUT}|)$ |
| full $q$ (SELECT *) | $O(N + |\text{OUT}_{\text{consistent}}|)$ | $O(N + |\text{OUT}|)$ |

*Consistent answers of common join queries can be computed with no asymptotic overhead*

| System | Target class | Interm. output | Backend |
|---|---|---|---|
| CAvSAT | * | SAT formula | SQL Server & MaxHS |
| Conquer | $\mathcal{C}_{\text{forest}}$ | SQL | SQL Server |
| Improved Conquesto | SJF **FO** | SQL | SQL Server |
| LinCQA | PPJT | SQL | SQL Server |

# Stackoverflow data

- Metadata of `stackoverflow.com` as of 02/2021 from Stack Exchange Data Dump
- 551M rows, 400 GB

| Table | # of rows | inRatio | bSize | # of Attributes |
|-------|-----------|---------|-------|-----------------|
| Users | 14M | 0% | 1 | 14 |
| Posts | 53M | 0% | 1 | 20 |
| PostHistory | 141M | 0.001% | 4 | 9 |
| Badges | 40M | 0.58% | 941 | 4 |
| Votes | 213M | 30.9% | 1441 | 6 |

# Experiments on Stackoverflow

$Q_1$ : Posts ⋈ Votes    $Q_2$ : Users ⋈ Badges    $Q_3$ : Users ⋈ Posts

$Q_4$ : Users ⋈ Posts ⋈ Comments

$Q_5$ : Posts ⋈ PostHistory ⋈ Votes ⋈ Comments

```
1  SELECT DISTINCT Posts.Title
2  FROM Posts, PostHistory, Votes, Comments
3  WHERE Posts.Tags LIKE "%SQL%"
4      AND Posts.id = PostHistory.PostId
5      AND Posts.id = Comments.PostId
6      AND Posts.id = Votes.PostId
7      AND Votes.BountyAmount > 100
8      AND PostHistory.PostHistoryTypeId = 2
9      AND Comments.score = 0
```

# Experiments on Stackoverflow

$Q_1$ : Posts ⋈ Votes     $Q_2$ : Users ⋈ Badges     $Q_3$ : Users ⋈ Posts

$Q_4$ : Users ⋈ Posts ⋈ Comments

$Q_5$ : Posts ⋈ PostHistory ⋈ Votes ⋈ Comments



| # poss. | 27578 | 145 | 38320 | 3925 | 1250 |
| # cons. | 27578 | 145 | 38320 | 3925 | 1245 |

# Concluding remarks

| Acyclic $q$ | PPJT | Yannakakis [VLDB'81] |
|---|---|---|
| Boolean $q$ | $O(N)$ | $O(N)$ |
| non-Boolean $q$ | $O(N \cdot |\text{OUT}_{\text{inconsistent}}|)$ | $O(N \cdot |\text{OUT}|)$ |
| full $q$ (SELECT *) | $O(N + |\text{OUT}_{\text{consistent}}|)$ | $O(N + |\text{OUT}|)$ |

# Concluding remarks

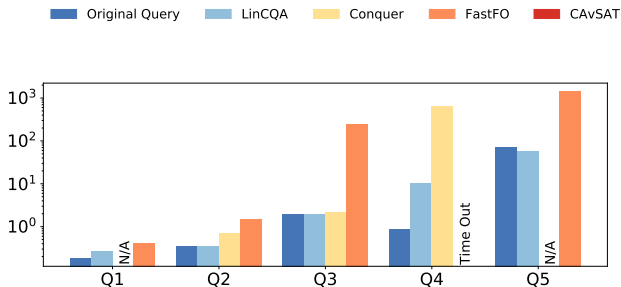| Acyclic $q$ | PPJT | Yannakakis [VLDB'81] |
|---|---|---|
| Boolean $q$ | $O(N)$ | $O(N)$ |
| non-Boolean $q$ | $O(N \cdot |\text{OUT}_{\text{inconsistent}}|)$ | $O(N \cdot |\text{OUT}|)$ |
| full $q$ (`SELECT *`) | $O(N + |\text{OUT}_{\text{consistent}}|)$ | $O(N + |\text{OUT}|)$ |

# Concluding remarks

| Acyclic $q$ | PPJT | Yannakakis [VLDB'81] |
|---|---|---|
| Boolean $q$ | $O(N)$ | $O(N)$ |
| non-Boolean $q$ | $O(N \cdot |\text{OUT}_{\text{inconsistent}}|)$ | $O(N \cdot |\text{OUT}|)$ |
| full $q$ (`SELECT *`) | $O(N + |\text{OUT}_{\text{consistent}}|)$ | $O(N + |\text{OUT}|)$ |



Thank you!