

EBERHARD KARLS

UNIVERSITÄT
TÜBINGEN



Explaining Missing Answers to SPJUA Queries

Melanie Herschel, Universität Tübingen, Germany

Mauricio A. Hernández, IBM Research - Almaden, USA

Conference on Very Large Databases, Singapore
September 14, 2010



Explaining Missing Data

U1_Email	Friend	U2_Email
john@univ.edu	Jane	jane@busy.com
jane@busy.com	John	john@univ.edu
jane@busy.com	Peter	peter@home.de
peter@home.de	Jane	jane@busy.com
john@univ.edu	\$name	\$email

U1_Email	Picture	PContributor
peter@home.de	pier39.jpg	John
john@univ.edu	\$pic	\$name

Query2

Pictures users are interested in

No user with that name?

No interests for user?

No pictures visible to user? ...

Query1

Pairs of connected users

Friend

User

UserInterest

Picture

PictureTag

UID1	UID2
U1	U2
U2	U3

UID	Email	Name
U1	john@univ.edu	John
U2	jane@busy.com	Jane
U3	peter@home.de	Peter

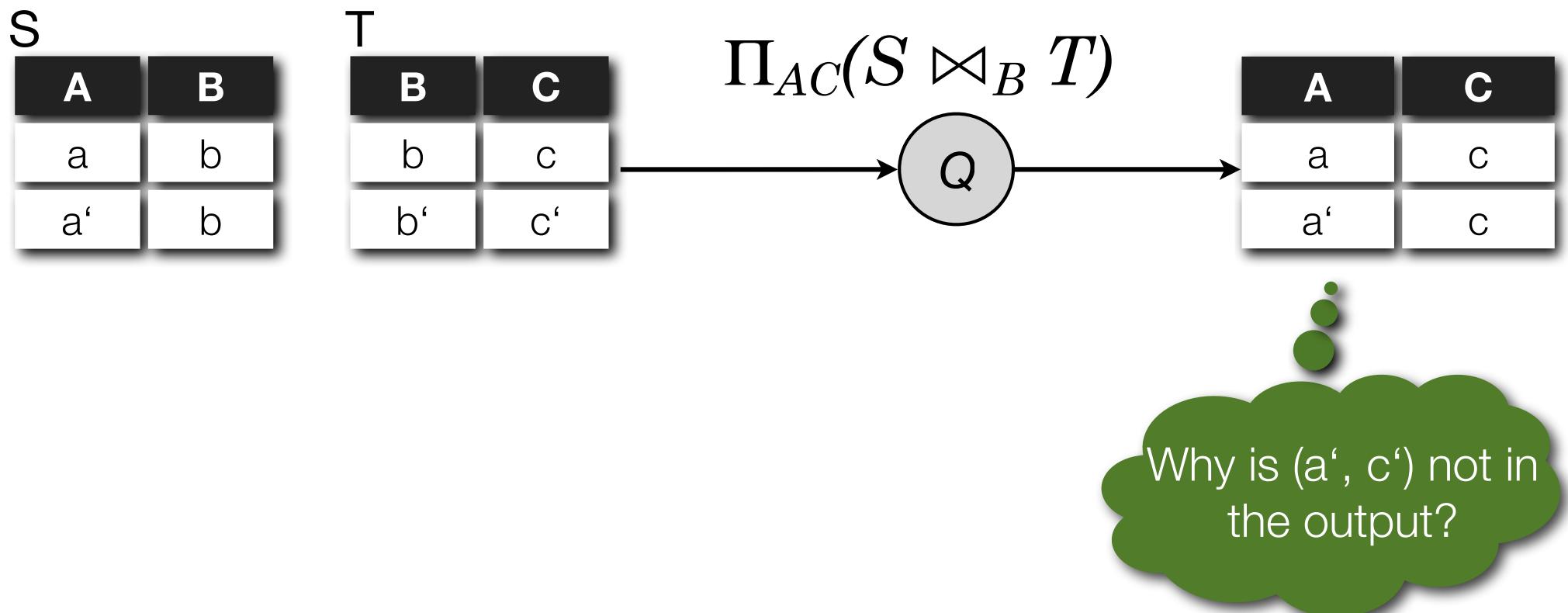
IID	UID
I1	U1
I4	U1
I2	U2
I3	U3
I4	U3

PID	UID	Picture	Visibility
P1	U1	goldengate.jpg	Friend
P2	U1	pier39.jpg	Public
P3	U2	market.jpg	Friend
P4	U3	winetasting.jpg	Public

PTID	PID	Category
PT1	P1	I3
PT2	P1	I1
PT3	P2	I4
PT4	P4	

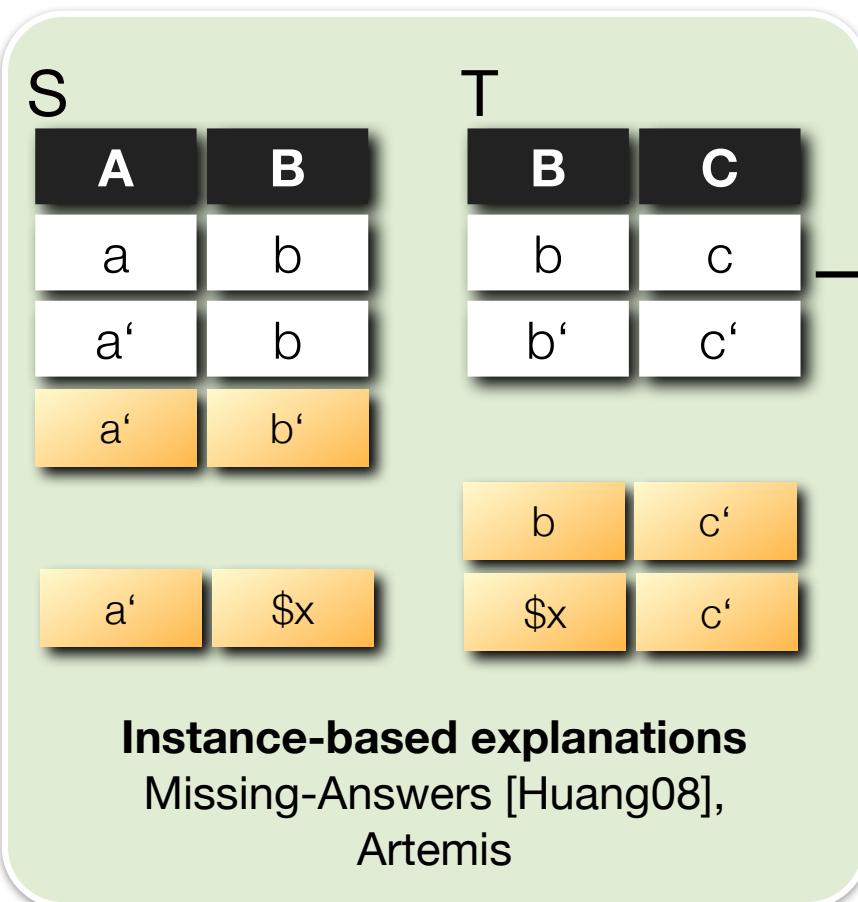
Explaining Missing Answers

Why is some data not in the result of a query Q?

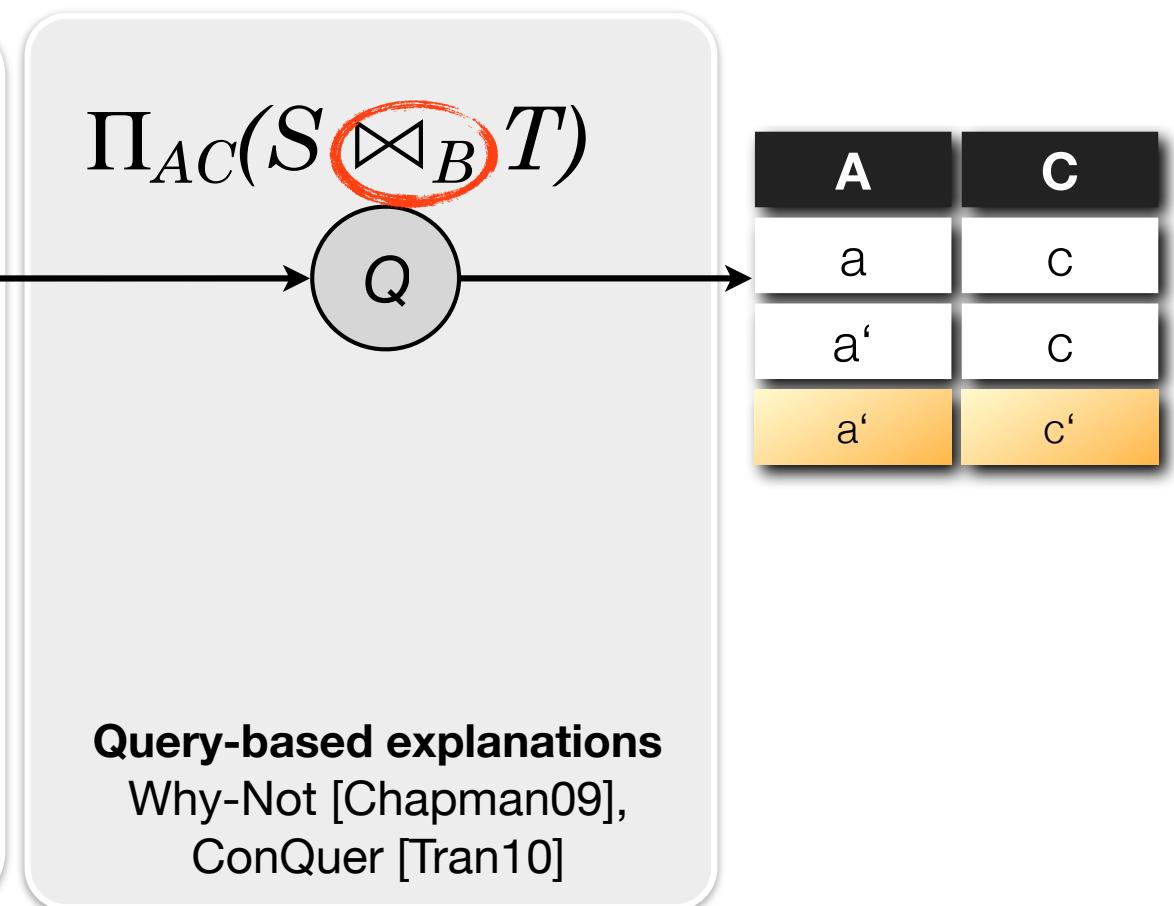


Explaining Missing Answers

Why is some data not in the result of a query Q ?



Instance-based explanations
Missing-Answers [Huang08],
Artemis

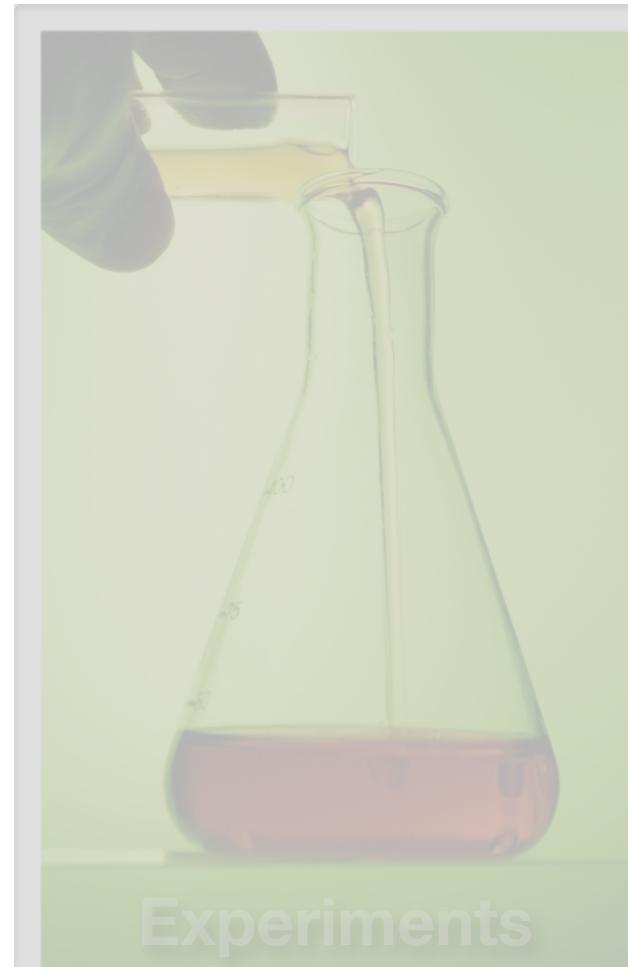


Query-based explanations
Why-Not [Chapman09],
ConQuer [Tran10]

Contributions

- **Artemis** algorithm
 - Explains a **set** of missing tuples over a **set** of queries that involve selection, projection, join, union, aggregation, and grouping (SPJUA).
 - Considers **side-effects**.
 - Guarantees on completeness and **correctness using a constraint solver**.
- **Framework** for instance-based explanation generation
- Comparative experimental **evaluation**

Agenda

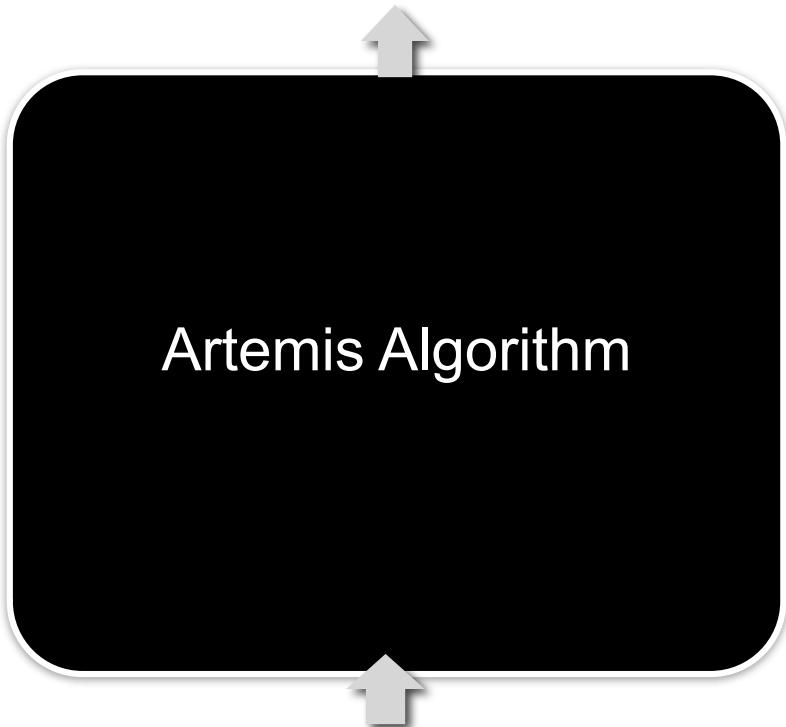


The Artemis Algorithm

For SPJU Queries



Set of explanations X



- 1) Source database D
- 2) A set of SPJU queries Q
- 3) A set of missing tuples E
- 4) *Further constraints*

The Artemis Algorithm

For SPJU Queries

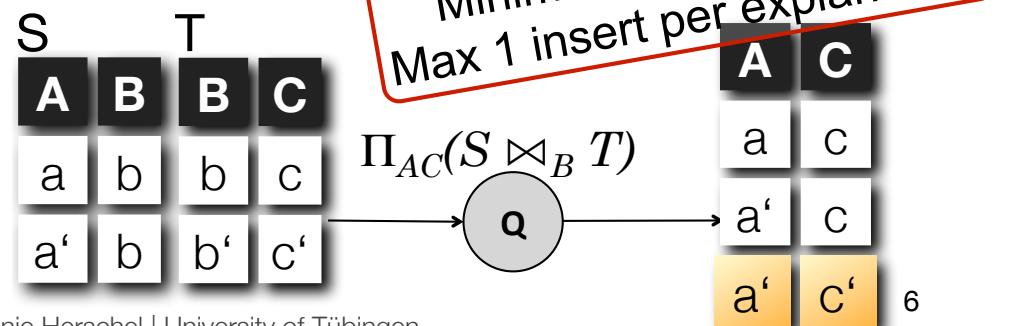
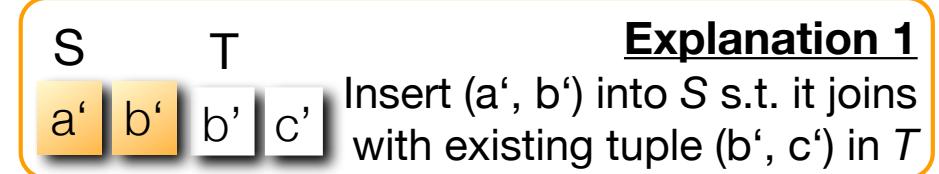


Set of explanations X



Artemis Algorithm

- 1) Source database D
- 2) A set of SPJU queries Q
- 3) A set of missing tuples E
- 4) Further constraints



The Artemis Algorithm

For SPJU Queries



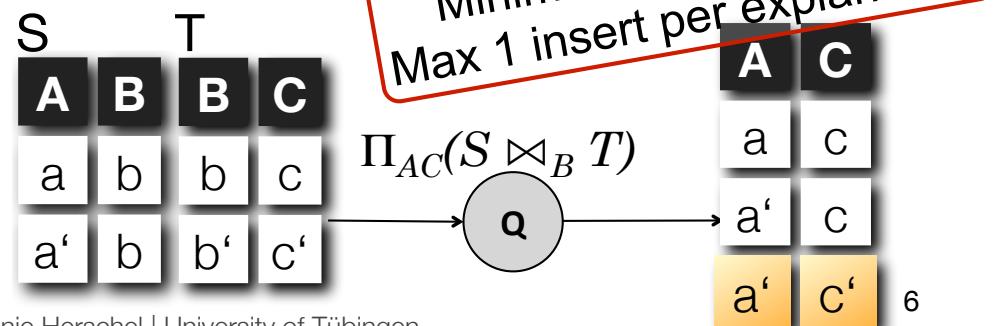
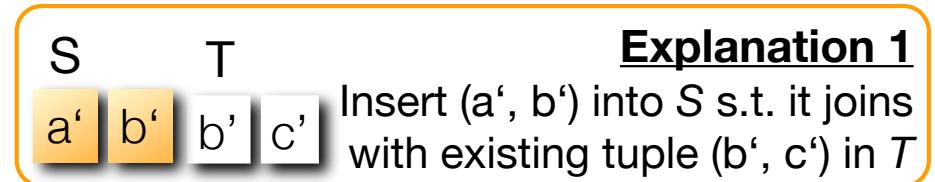
Set of explanations X



- (1) Compute generic witness
- (2) Create conditional tables (c-tables) for D
- (3) Compute c-tables of Q
- (4) Generate explanations
- (5) Filter and sort explanations



- 1) Source database D
- 2) A set of SPJU queries Q
- 3) A set of missing tuples E
- 4) Further constraints



The Artemis Algorithm

For SPJU Queries



S	T		
a'	b'	b'	c'

Explanation 1

Insert (a', b') into S s.t. it joins with existing tuple (b', c') in T



(5) Filter and sort explanations

(4) Generate explanations

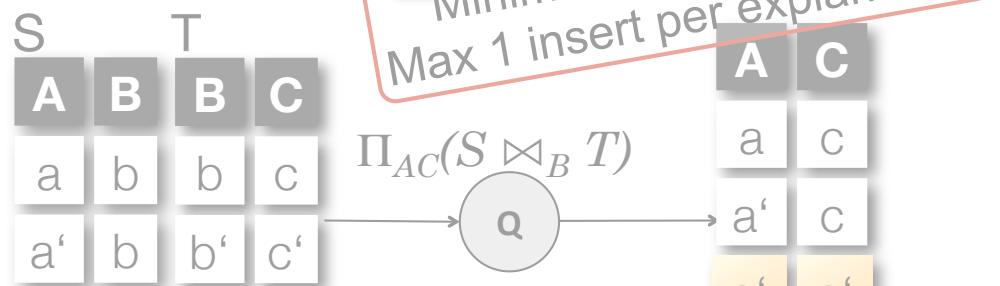
(3) Compute c-tables of Q

(2) Create conditional tables
(c-tables) for D

(1) Compute generic witness

$$E = \{(a', c')\}$$

$$Q: V(a', c') :- R(a', v_b), S(v_b, c')$$



The Artemis Algorithm

For SPJU Queries



S	T
a'	b'
b'	c'

Explanation 1

Insert (a', b') into S s.t. it joins with existing tuple (b', c') in T

(5) Filter and sort explanations

(4) Generate explanations

(3) Compute c-tables of \mathbf{Q}

(2) Create conditional tables
(c-tables) for \mathcal{D}

(1) Compute generic witness

Generic Witness:
 $R(a', \$x), S(\$x, c')$

S	T	
A	B	
a	b	
a'	b	
	b'	
	c	
	c'	

$\Pi_{AC}(S \bowtie_B T)$

q

Minimum #side-effects,
Max 1 insert per explanation

The Artemis Algorithm

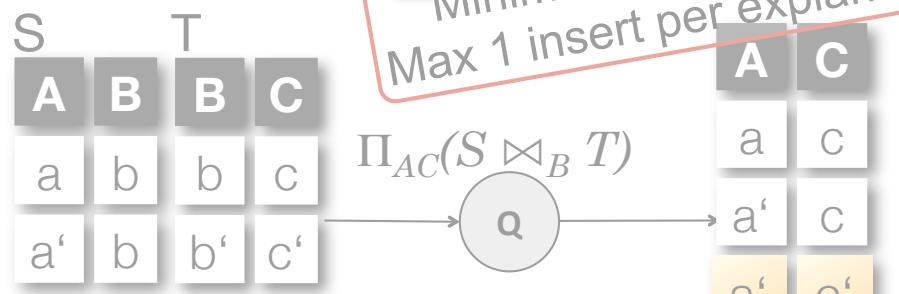
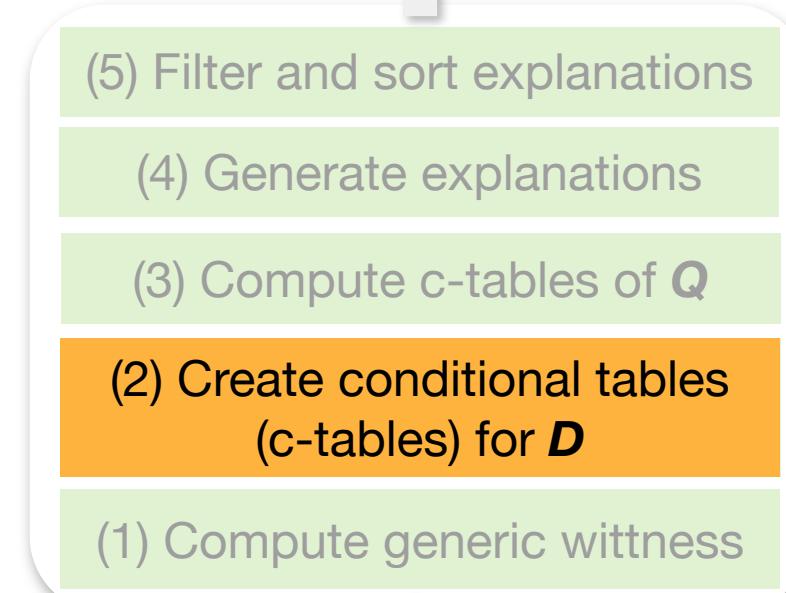
For SPJU Queries



S	T		
a'	b'	b'	c'

Explanation 1

Insert (a', b') into S s.t. it joins with existing tuple (b', c') in T



Minimum #side-effects,
Max 1 insert per explanation

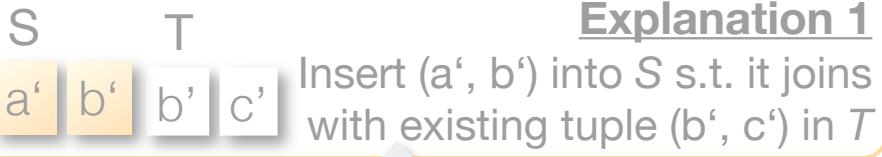
S^C		con
A	B	
a	b	TRUE
a'	b	TRUE
a'	\$x1	$\$x1 \neq b$

T^C		con
B	C	
b	c	TRUE
b'	c'	TRUE
\$x2	c'	$\$x2 \neq b'$

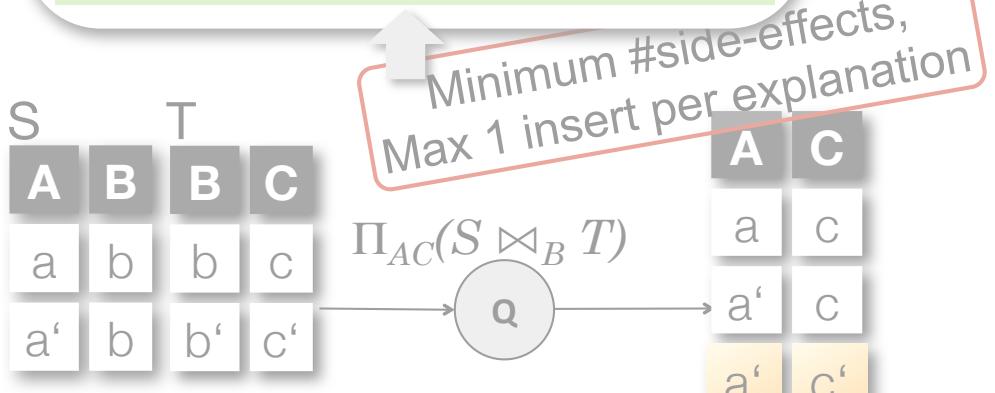
Generic Witness:
 $R(a', \$x), S(\$x, c')$

The Artemis Algorithm

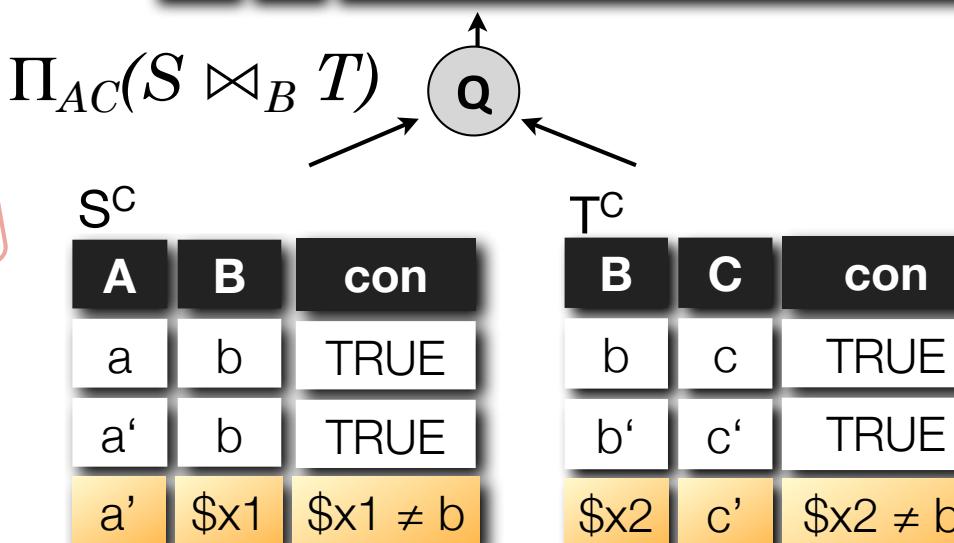
For SPJU Queries



- (5) Filter and sort explanations
- (4) Generate explanations
- (3) Compute c-tables of **Q**
- (2) Create conditional tables (c-tables) for **D**
- (1) Compute generic witness



A	C	con
a	c	TRUE
a'	c	TRUE
a	c'	$\$x_2 = b \wedge \$x_2 \neq b'$
a'	c'	$\$x_2 = b \wedge \$x_2 \neq b'$
a'	c'	$\$x_1 = b' \wedge \$x_1 \neq b$
a'	c'	$\$x_1 = \$x_2 \wedge \$x_2 \neq b' \wedge \$x_1 \neq b$



The Artemis Algorithm

For SPJU Queries



S	T		
a'	b'	b'	c'

Explanation 1

Insert (a', b') into S s.t. it joins with existing tuple (b', c') in T

(5) Filter and sort explanations

(4) Generate explanations

(3) Compute c-tables of Q

(2) Create conditional tables
(c-tables) for D

(1) Compute generic witness

S	T		
A	B	B	C
a	b	b	c
a'	b	b'	c'

Minimum #side-effects,
Max 1 insert per explanation

$$\Pi_{AC}(S \bowtie_B T)$$

A	C
a	c
a'	c
a'	c'

Constraint Satisfaction Problem

tuple (a', c') exists

AND

minimum number
of side-effects

side-effect
matches of
 (a', c')

A	C	con
a	c	TRUE
a'	c	TRUE
a	c'	$\$x_2 = b \wedge \$x_2 \neq b'$
a'	c'	$\$x_2 = b \wedge \$x_2 \neq b'$
a'	c'	$\$x_1 = b' \wedge \$x_1 \neq b$
a'	c'	$\$x_1 = \$x_2 \wedge \$x_2 \neq b' \wedge \$x_1 \neq b$

The Artemis Algorithm

For SPJU Queries



S	T		
a'	b'	b'	c'

Explanation 1

Insert (a', b') into S s.t. it joins with existing tuple (b', c') in T

(5) Filter and sort explanations

(4) Generate explanations

(3) Compute c-tables of Q

(2) Create conditional tables (c-tables) for D

(1) Compute generic witness

S	T		
A	B	B	C
a	b	b	c
a'	b	b'	c'

Minimum #side-effects,
Max 1 insert per explanation

$$\Pi_{AC}(S \bowtie_B T)$$

A	C
a	c
a'	c
a'	c'

side-effect
matches of
 (a', c')

Constraint Satisfaction Problem

$$\$x2 = b \wedge \$x2 \neq b'$$

AND

$$\$x2 = b \wedge \$x2 \neq b'$$

$$\$x1 = b' \wedge \$x1 \neq b$$

$$\$x1 = \$x2 \wedge \$x2 \neq b' \wedge \$x1 \neq b$$

A	C	con
a	c	TRUE
a'	c	TRUE
a	c'	$\$x2 = b \wedge \$x2 \neq b'$
a'	c'	$\$x2 = b \wedge \$x2 \neq b'$
a'	c'	$\$x1 = b' \wedge \$x1 \neq b$
a'	c'	$\$x1 = \$x2 \wedge \$x2 \neq b' \wedge \$x1 \neq b$

The Artemis Algorithm

For SPJU Queries



S	T		
a'	b'	b'	c'

Explanation 1

Insert (a', b') into S s.t. it joins with existing tuple (b', c') in T



- (5) Filter and sort explanations
- (4) Generate explanations
- (3) Compute c-tables of Q
- (2) Create conditional tables (c-tables) for D
- (1) Compute generic witness

S	T		
A	B	B	C
a	b	b	c
a'	b	b'	c'

Minimum #side-effects,
Max 1 insert per explanation

$$\Pi_{AC}(S \bowtie_B T)$$

A	C
a	c
a'	c
a'	c'

side-effect
matches of
 (a', c')

Output of constraint solver for...

1st match: $\$x_2 = b$,
1 side-effect

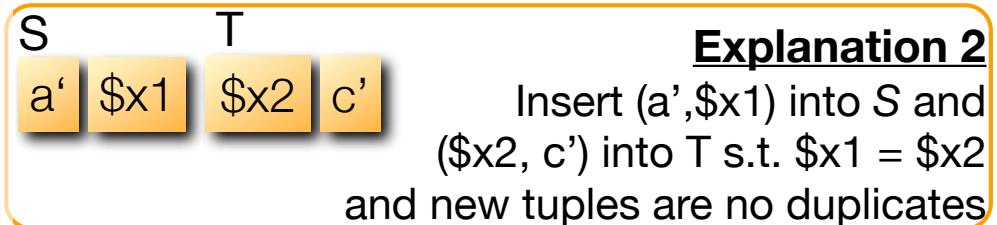
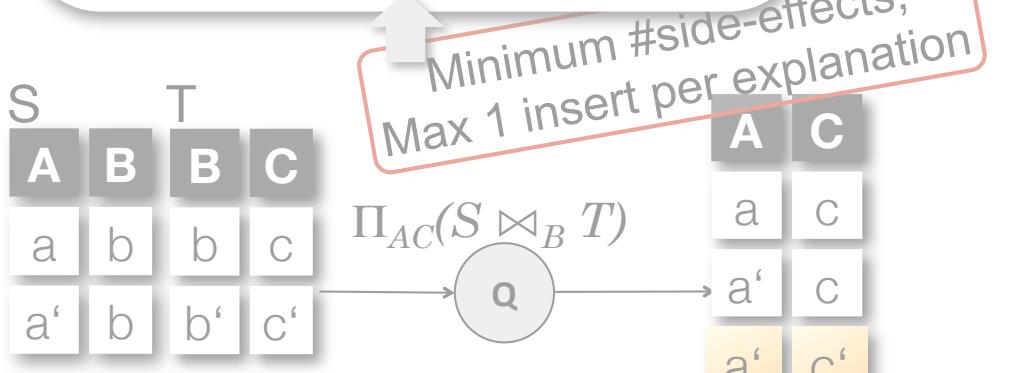
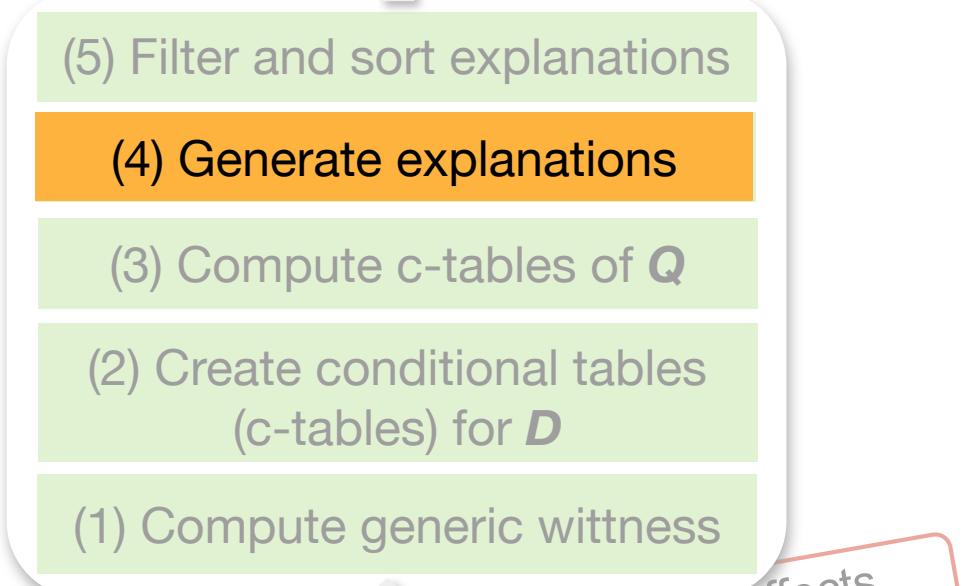
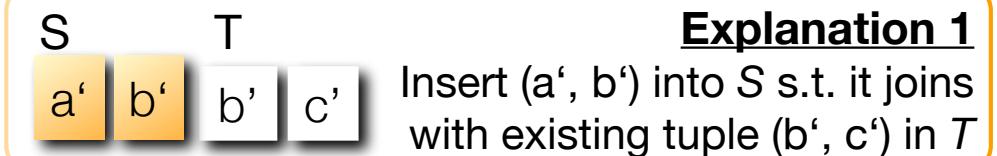
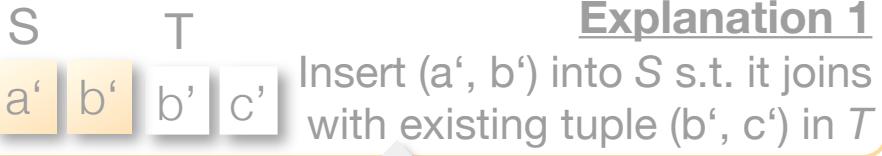
2nd match: $\$x_1 = b'$,
0 side-effects

3rd match: $\$x_1 = \x_2 , $\$x_1 \neq b$, $\$x_2 \neq b'$,
0 side-effects

A	C	con
a	c	TRUE
a'	c	TRUE
a	c'	$\$x_2 = b \wedge \$x_2 \neq b'$
a'	c'	$\$x_2 = b \wedge \$x_2 \neq b'$
a'	c'	$\$x_1 = b' \wedge \$x_1 \neq b$
a'	c'	$\$x_1 = \$x_2 \wedge \$x_2 \neq b' \wedge \$x_1 \neq b$

The Artemis Algorithm

For SPJU Queries



Output of constraint solver for...

1st match: ~~$\$x2 = b$,
1 side-effect~~ Too many side-effects

2nd match: $\$x1 = b'$,
0 side-effects

3rd match: $\$x1 = \$x2$, $\$x1 \neq b$, $\$x2 \neq b'$,
0 side-effects

The Artemis Algorithm

For SPJU Queries



S	T		
a'	b'	b'	c'

Explanation 1

Insert (a', b') into S s.t. it joins with existing tuple (b', c') in T

S	T		
a'	b'	b'	c'

Explanation 1

Insert (a', b') into S s.t. it joins with existing tuple (b', c') in T

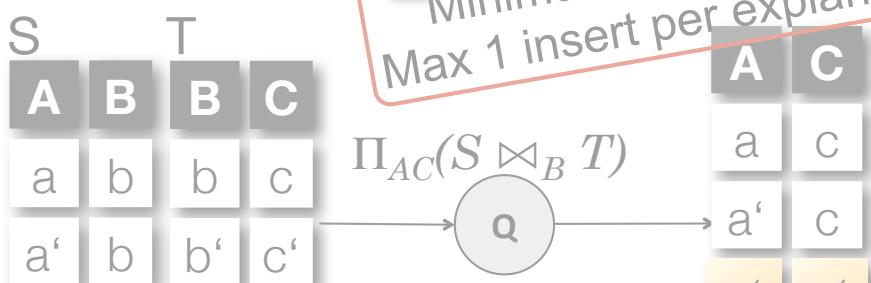
(5) Filter and sort explanations

(4) Generate explanations

(3) Compute c-tables of Q

(2) Create conditional tables (c-tables) for D

(1) Compute generic witness

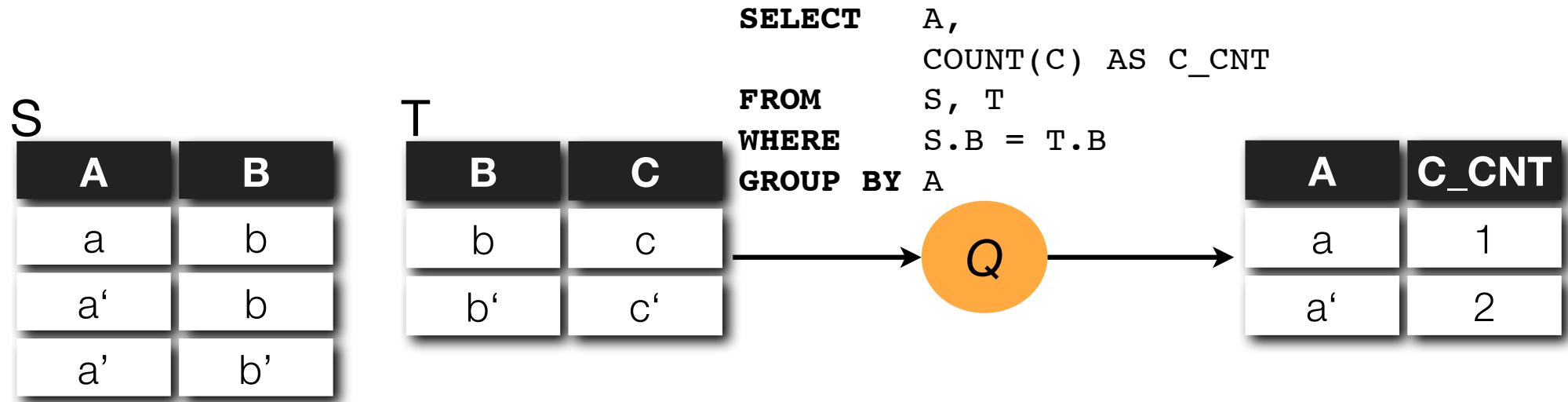


Minimum #side-effects,
Max 1 insert per explanation

A	C
a	c
a'	c

The Artemis Algorithm

Extension to Aggregation and Grouping



- Meet-in-the-middle approach
 - Given set of missing tuples E on an SPJA view, determine E' for SPJ view, s.t. applying aggregation and grouping on the SPJ view results in the original SPJA view.
 - Artemis algorithm for SPJ views for E' .
 - Extend returned explanations to include aggregation and grouping.

The Artemis Algorithm

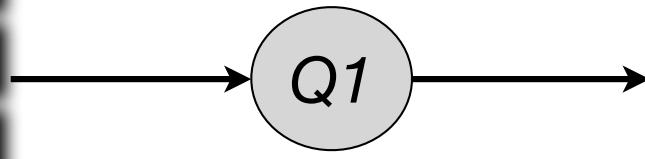
Extension to Aggregation and Grouping



S	T		
A	B	B	C
a	b	b	c
a'	b	b'	c'
a'	b'		

```

SELECT A, C
FROM S, T
WHERE S.B = T.B
    
```



10000 X

A	C
a	c
a'	c
a'	c'
b'	\$v

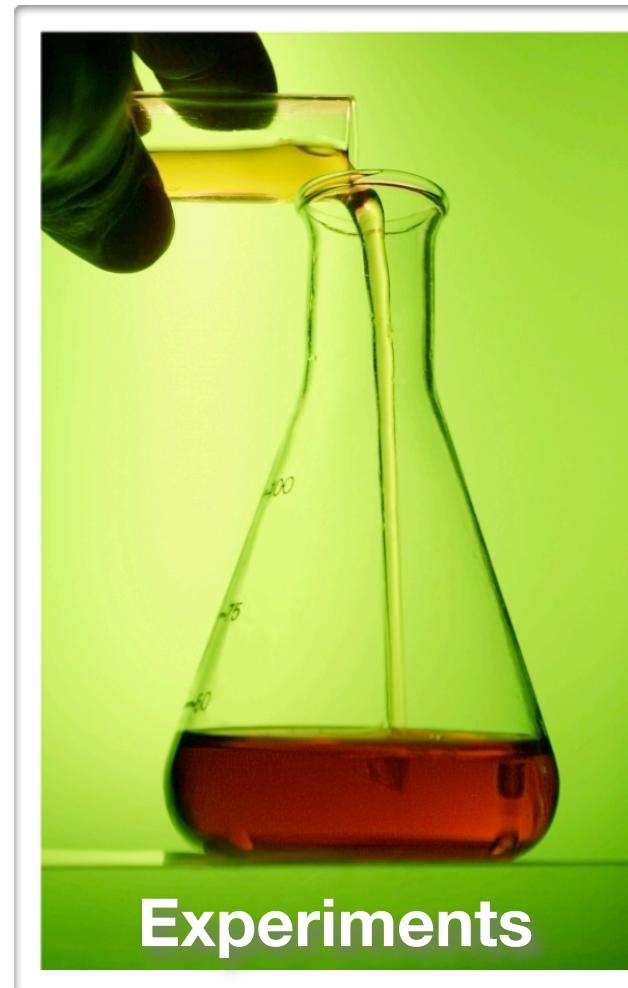
- Meet-in-the-middle approach
 - Given set of missing tuples E on an SPJA view, determine E' for SPJ view, s.t. applying aggregation and grouping on the SPJ view results in the original SPJA view.
 - Artemis algorithm for SPJ views for E' .
 - Extend returned explanations to include aggregation and grouping.

```

SELECT A,
       COUNT(C) AS C_CNT
FROM Q1
GROUP BY A
    
```

A	C_CNT
a	1
a'	2
b'	10000

Agenda



Experimental Setup

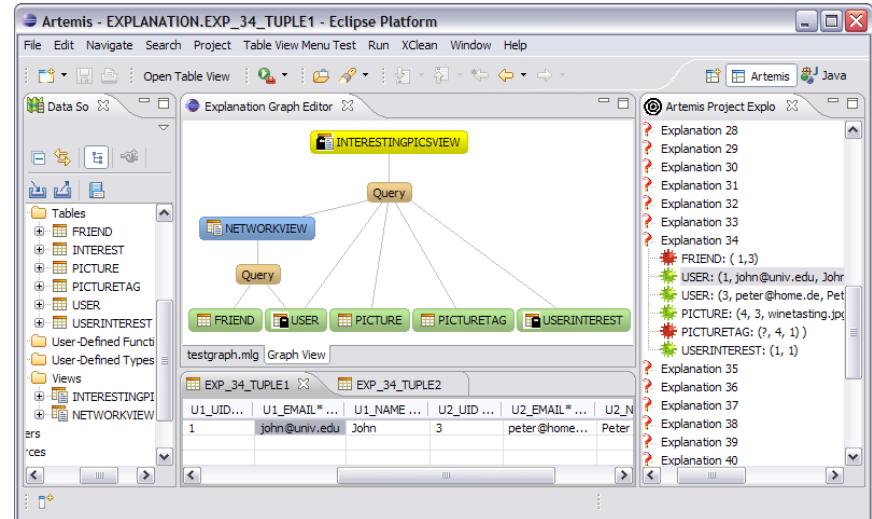


Implementation

- Eclipse Plugin [VLDB09].
- Artemis and Missing-Answers [VLDB08]
- Minion used as constraint solver for Artemis.
- IBM DB2 9.5 used as RDBMS.

Datasets

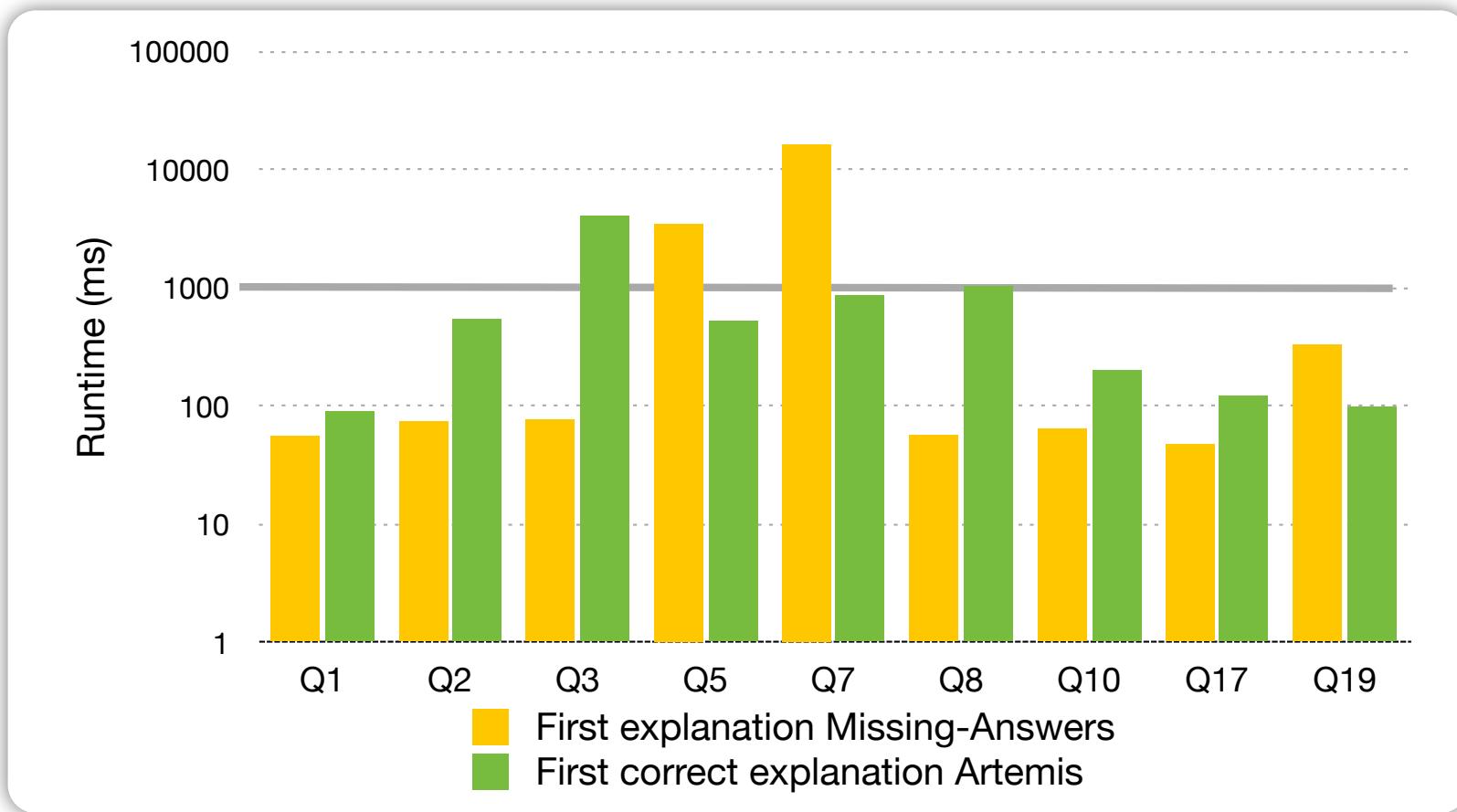
- TPCH
 - 10 MB of data
 - 9 queries (adaptations of TPCH queries limited to supported types of queries)
 - No insertions on Nation and Region.



TPC Transaction Processing
Performance Council



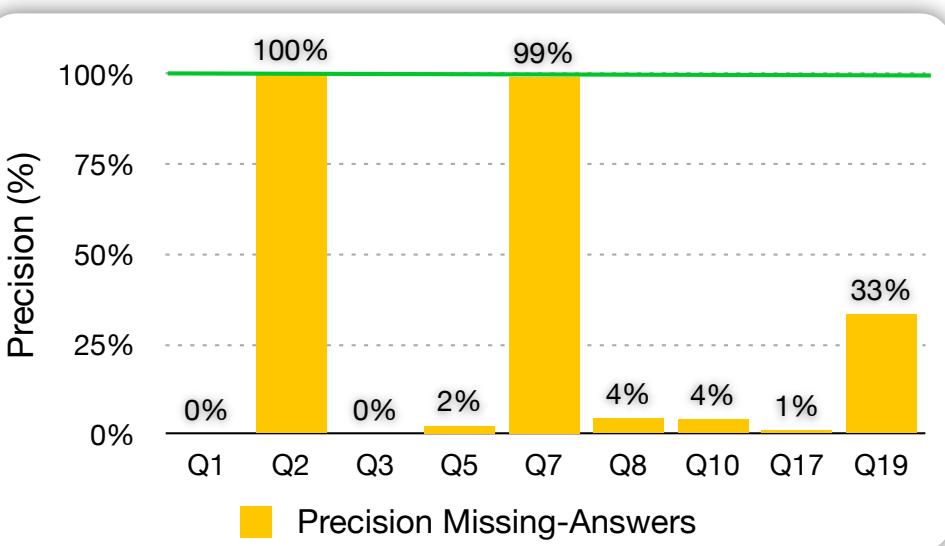
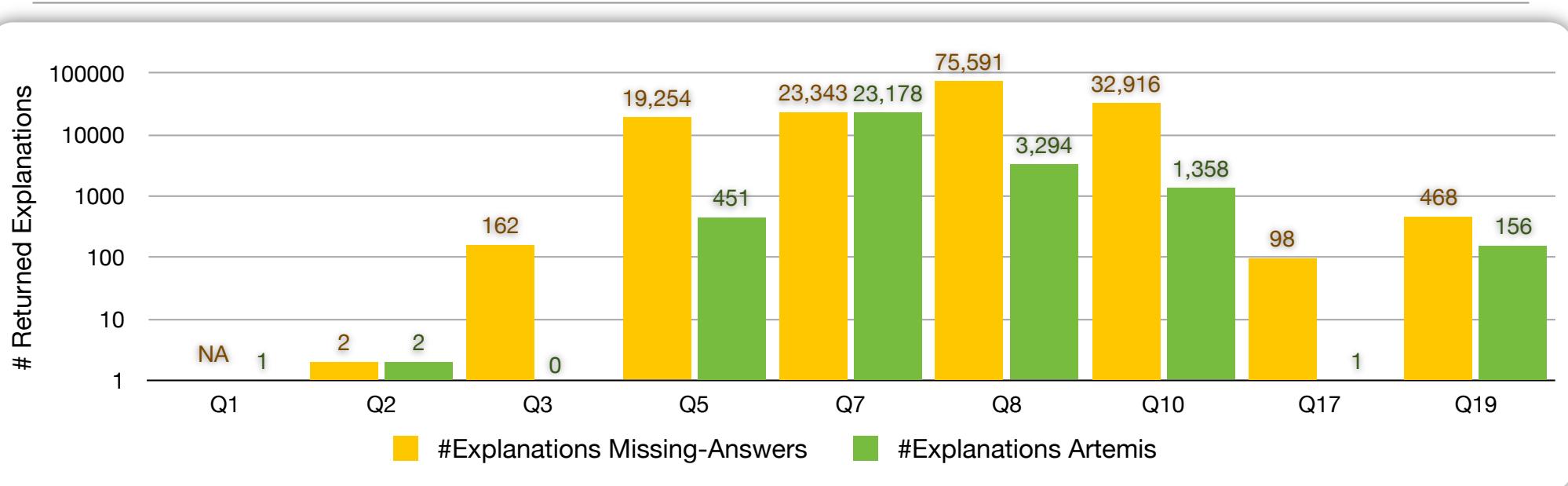
Runtime to First Correct Explanation



Artemis takes less than a second to find first correct explanation in most cases.

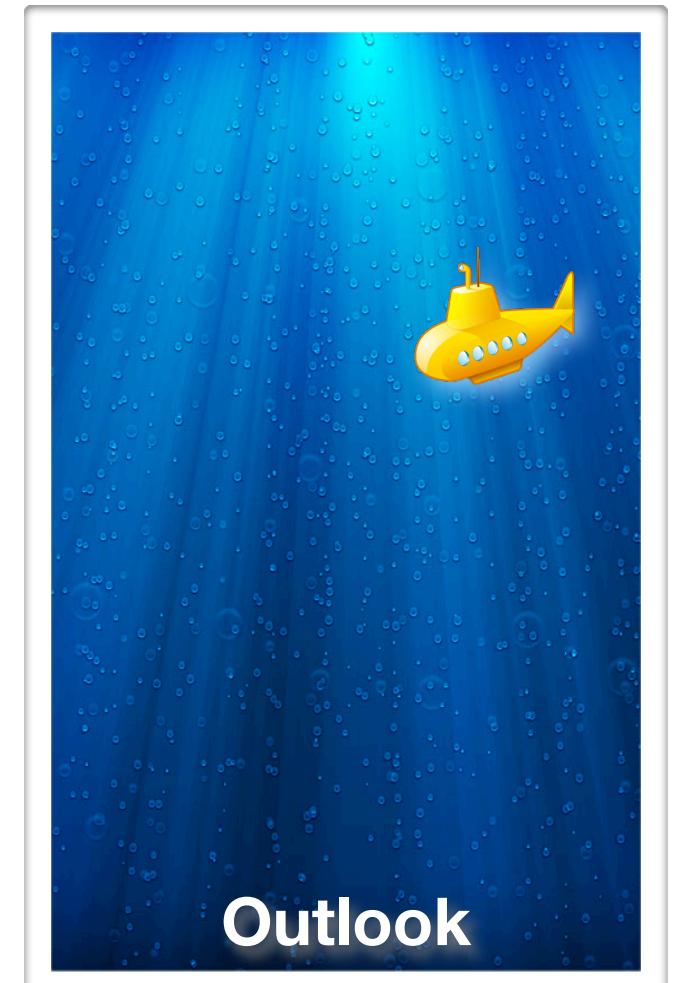
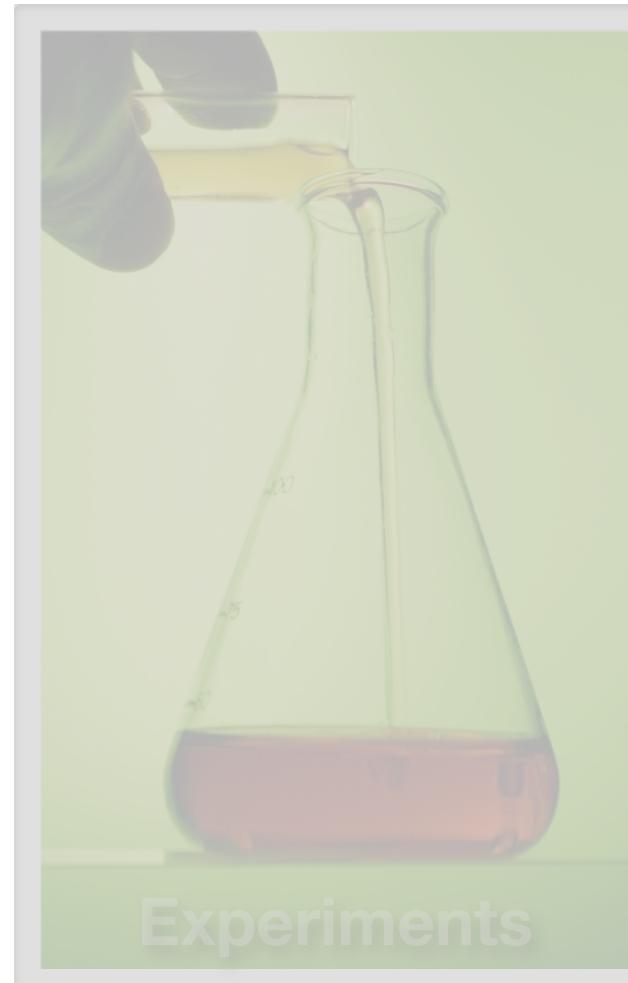
Missing-Answers usually faster, but returned explanation can be wrong.

Effectiveness



Number of unsatisfiable explanations can be substantial when using Missing-Answers.
Constraint solver makes Artemis run slower, but effectiveness significantly improved.

Agenda





Outlook

So far...

- **Framework** for instance-based explanation generation.
- **Artemis** algorithm (SPJUA queries, side-effects, correctness)
- Comparative experimental **evaluation**

In the future..

- Efficiency improvement
- Visualization to improve usability

The “Big Picture”

- Build a system to analyze, fix, and test data transformations



Outlook

So far...

- **Framework** for instance-based explanation generation.
- **Artemis** algorithm (SPJUA queries, side-effects, correctness)
- Comparative experimental **evaluation**

In the future..

- Efficiency improvement
- Visualization to improve usability

The “Big Picture”

- Build a system to analyze, fix, and test data transformations

