



# Capturing Factorized Provenance (FDB 2022)

**Boris Glavic**

*Illinois Institute of Technology*

*Database Group*

Joint work with:

**Seokki Lee** (*University of Cincinnati*), **Bertram Ludäscher** (*UIUC*)

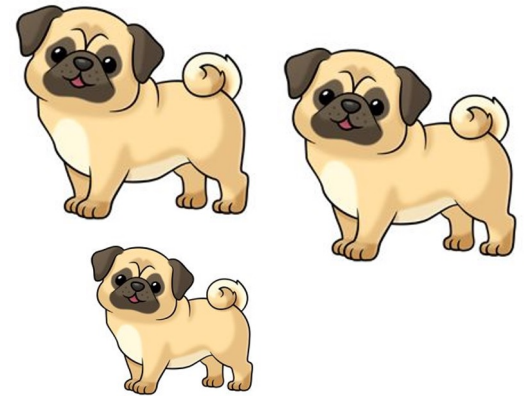
## Outline

- **Overview & Motivation**
- Provenance Graph Model
- Capturing Provenance
- Factorization
- Experiments
- Conclusions & Future work

# Debugging Datalog

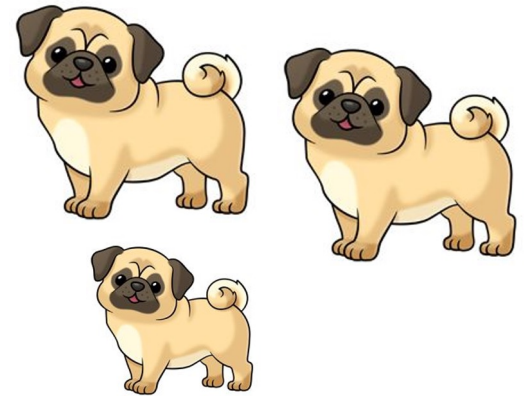
---

- **Datalog has seen a resurgence in academia & industry**
  - Programming distributed systems
  - Complex data-intensive computations
  - Analytics & ML over query results
- **Urgent need for debugging Datalog programs**
  - Why did my program produce this unexpected result?
  - Why did my program not produce this expected result?
  - Which rules are responsible for deriving this result?
  - Why did this rule derivation fail?
  - Would deleting & inserting a tuple change the result?



## PUG + PUGS

*Efficient capture and summarization for **why and why-not** provenance through a **provenance model developed** for queries with negation*



## PUG + PUGS

*Efficient capture and summarization for **why and why-not** provenance through a **provenance model developed** for queries with negation*

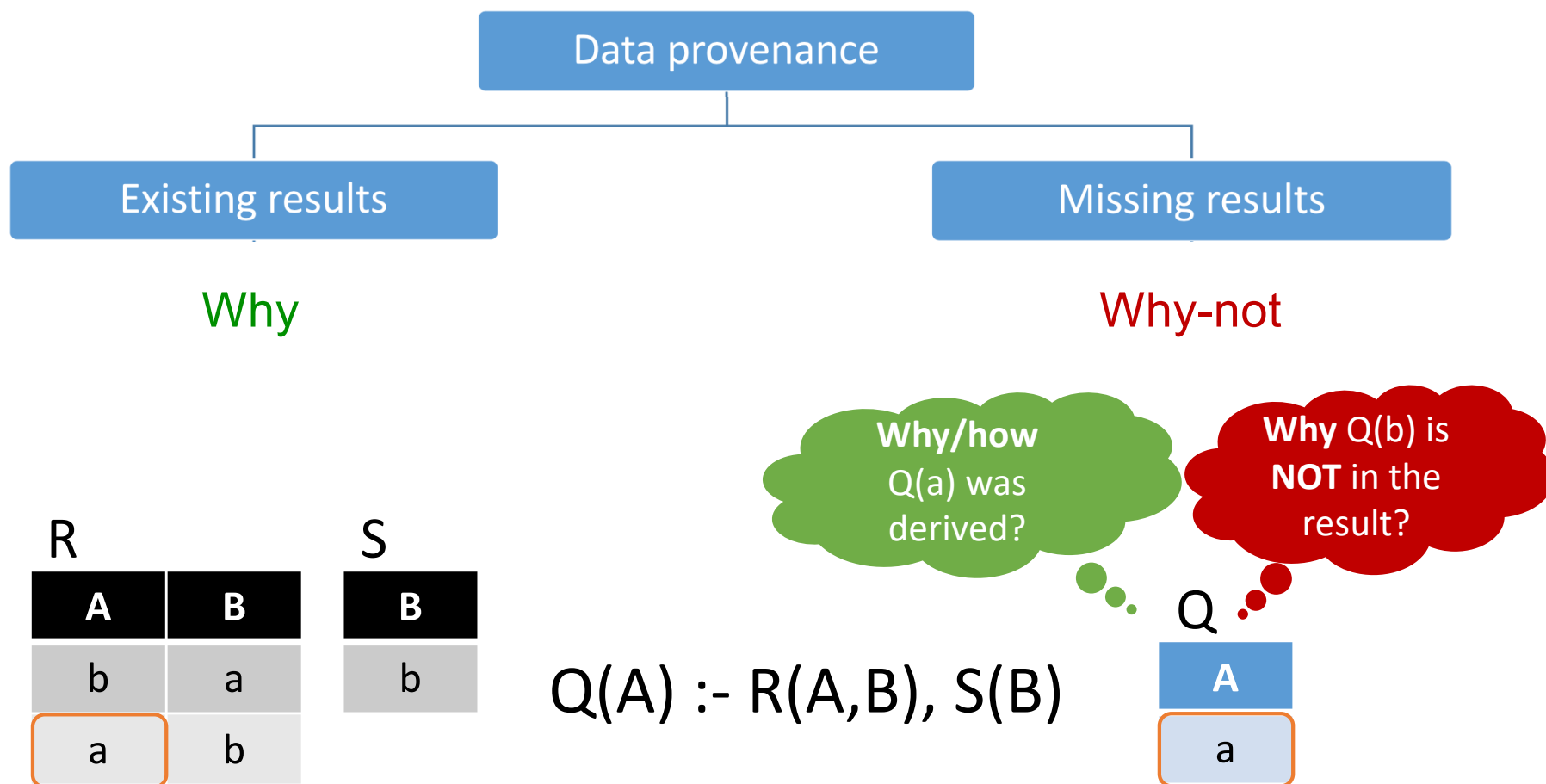
***This is FDB!***



*PUG utilizes a flat-relational encoding of provenance graphs which essentially corresponds to a factorized representation of provenance*

# Schism between Why / Why-not

- How the rules of a program did derive / failed to derive an existing / missing output from the input data



## Unifying Why / Why-not

---

- **Why** (provenance) and **why-not** (missing answers) have been mostly treated in isolation
- Why and why-not questions can be reduced to each other for a query language **L** if ...
  - for any query **Q**, its complement **Q<sup>C</sup>** is in **L**

$$t \notin Q(D) \equiv t \in Q^C(D)$$

$$Q(A) :- R(A, B)$$

$$Q^C(A) :- R(A, B), \text{ adom}(A), \text{ adom}(B)$$

## Outline

---

- Overview & Motivation
- **Provenance Graph Model**
- Capturing Provenance
- Factorization
- Experiments
- Conclusions & Future work



# Requirements

---

- **Syntax-driven**
  - Provenance structure aligns with program structure
- **Compatible with well-established provenance models**
  - Provenance polynomials for positive queries [1]
  - Dual polynomials [2]
- **Build-in support for sharing common subexpressions**

[1] T. Green, G. Karvounarakis, and V. Tannen. Provenance semirings. In PODS, pages 31–40, 2007.

[2] E.Grädel and V.Tannen.Semiring provenance for first-order model checking. arXiv preprint arXiv:1712.01980, 2017.



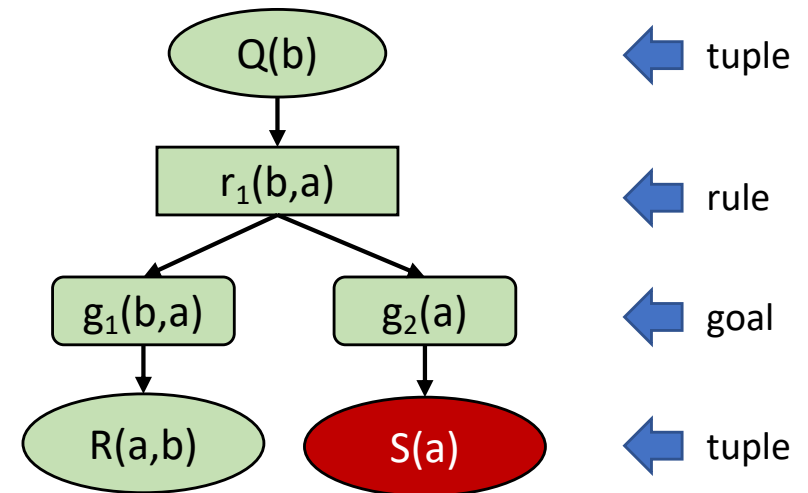
# Provenance Graph Model

R		S
A	B	B
b	a	b
a	b	

$r_1: Q(A) :- R(A,B), \neg S(B)$

Q

A
b



Green: success / existence  
Red: failure / absence



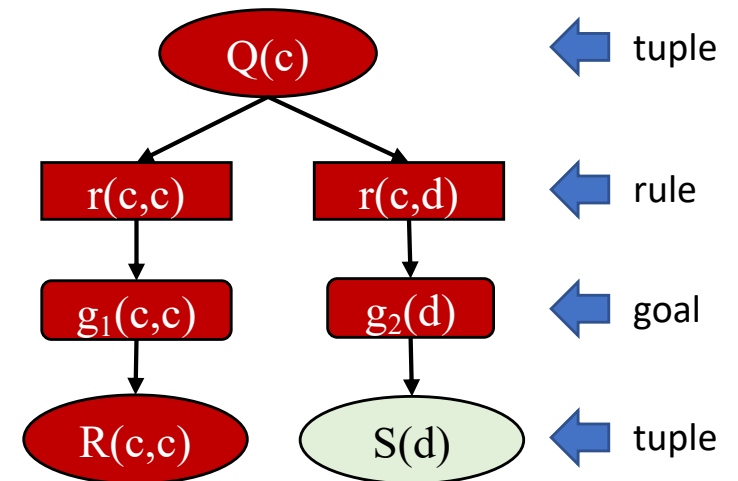
# Provenance Graph Model

R		S
A	B	B
d	c	d
c	d	

$r: Q(A) :- R(A,B), \neg S(B)$

Q
A
d

Why Q(c) is NOT in the output?



Green: success / existence  
Red: failure / absence

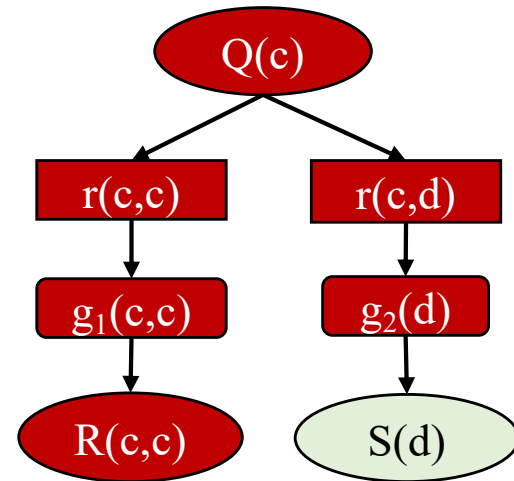
# Provenance Graph Model

R		S
A	B	B
d	c	d
c	d	

$r: Q(A) :- R(A,B), \neg S(B)$

Q
A
d

Why Q(c) is NOT in the output?





# Provenance Graph Model

R		S
A	B	B
d	c	d
c	d	

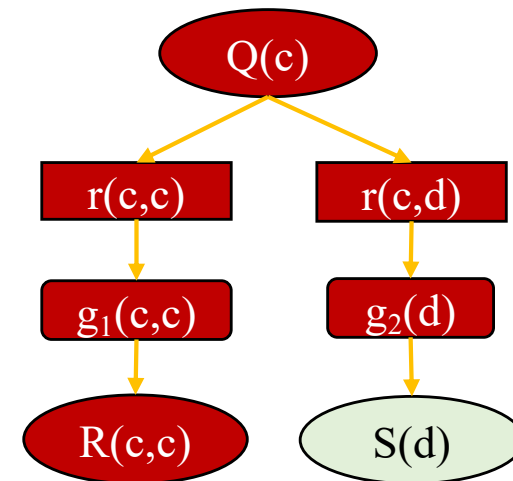


$r: Q(A) :- R(A,B), \neg S(B)$



Q
A
d

Why Q(c) is NOT in the output?





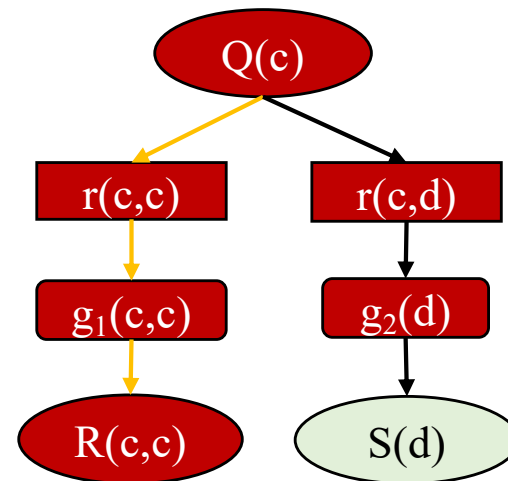
# Provenance Graph Model

R		S
A	B	B
d	c	d
c	d	

$r: Q(c) :- R(c,c), \neg S(c)$

Q
A
d

Why Q(c) is NOT in the output?



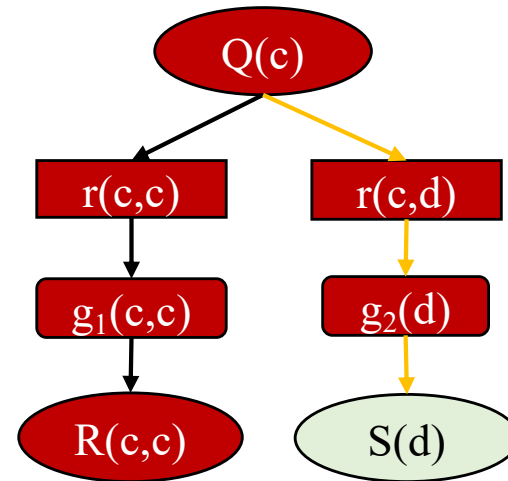
# Provenance Graph Model

R		S
A	B	B
d	c	d
c	d	

$r: Q(c) :- R(c,d), \neg S(d)$

Q
A
d

Why Q(c) is NOT in the output?





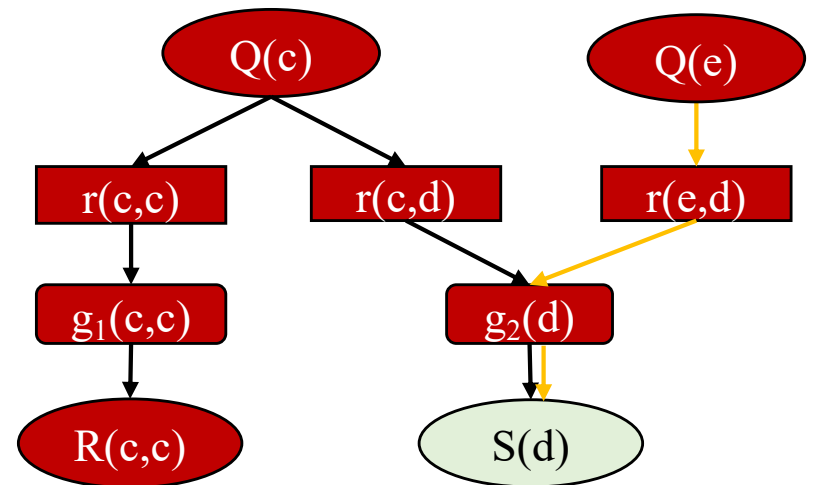
# Provenance Graph Model

R		S
A	B	B
d	c	d
c	d	
e	d	

$r: Q(e) :- R(e,d), \neg S(d)$

Q
A
d

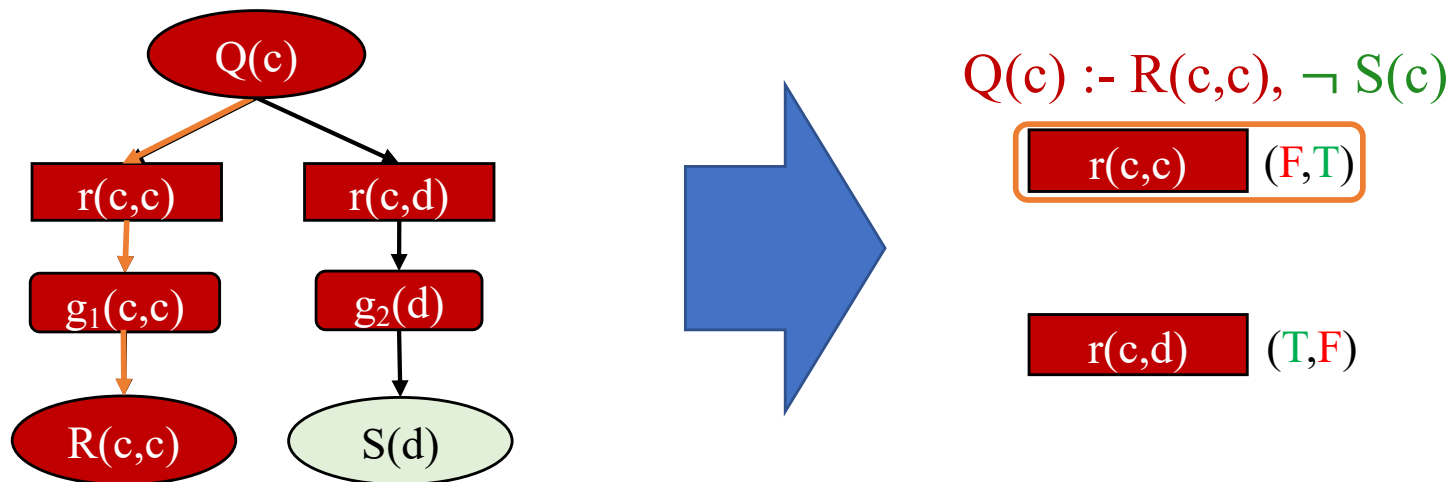
Why Q(c) is NOT in the output?





# Provenance Graph Model

- Equivalent model (annotated rule derivations)

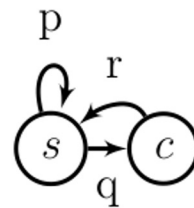


# Extracting Provenance Polynomials

$$r_2 : Q_{3hop}(X, Y) :- T(X, A), T(A, B), T(B, Y)$$

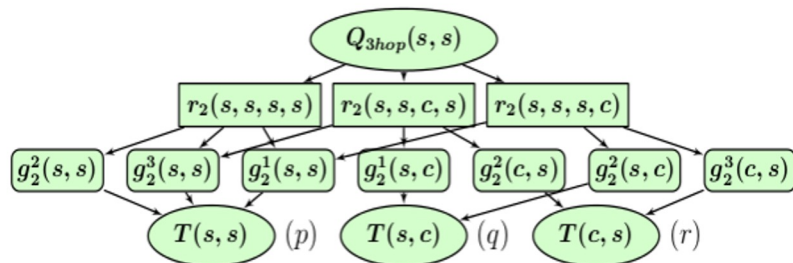
Relation Train

fromCity	toCity	$\mathbb{N}[X]$
seattle	seattle	$p$
seattle	chicago	$q$
chicago	seattle	$r$



Result of query  $Q_{3hop}$

X	Y	$\mathbb{N}[X]$
seattle	seattle	$p^3 + 2pqr$



$$r_2(g_1(p) \cdot g_2(p) \cdot g_3(p)) + r_2(g_1(q) \cdot g_2(r) \cdot g_3(p))$$

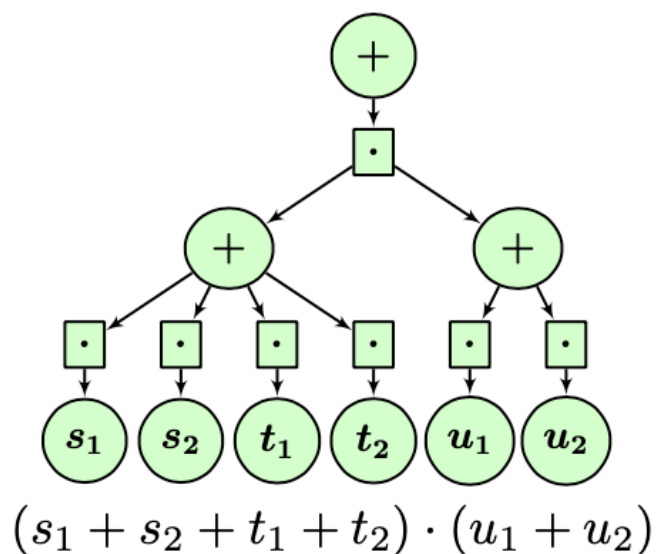
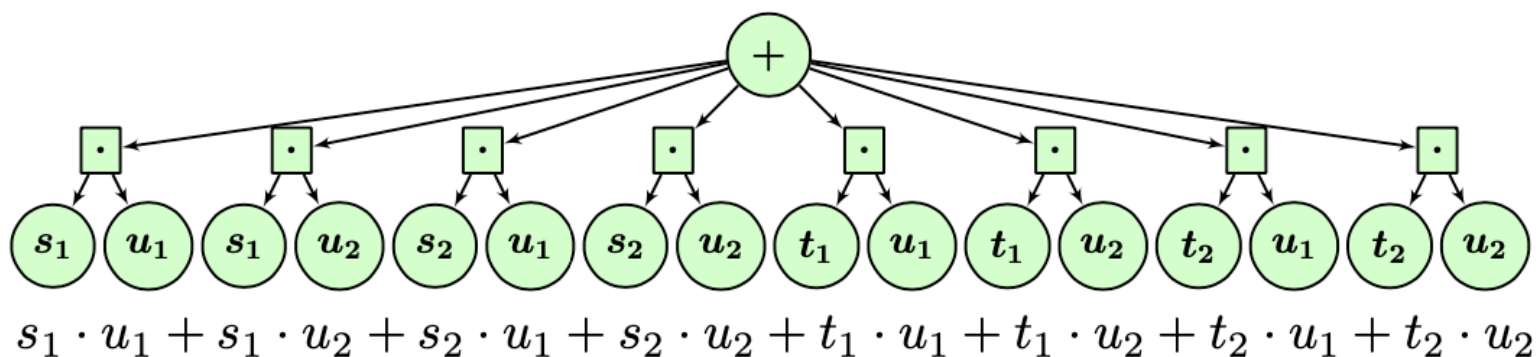
$$+ r_2(g_1(p) \cdot g_2(q) \cdot g_3(r))$$



$$p^3 + 2 \cdot (p \cdot q \cdot r)$$

$$[[p, p, p], [p, q, r], [p, q, r]]$$

# Factorized Provenance Graphs



## Outline

- Overview & Motivation
- Provenance Graph Model
- **Capturing Provenance**
- Factorization
- Experiments
- Conclusions & Future work

# Computing Provenance

---

## Problem definition

**input:** (missing) answers of interest

**output:** relevant subgraph for all answers of interest

## **Instrumentation** of the input Datalog program

- Use **firing rules** to capture rule derivations
- **Goal-oriented** approach
  - Limit the provenance based on the user question
  - Efficient bottom-up evaluation using relational engines
  - Compiling rewritten program into SQL

# Overview

---

## Rewriting steps

- **Unifying** the program with specification of outputs of interest
- **Statically annotate** program to indicate interest in success / failure
- Creating **firing rules**
- Checking **connectivity with joins**
- **Computing** the **edge** relation of the provenance graph

## Example query and provenance question

$r_1 : Q(X, Y) : \neg \text{Train}(X, Z), \text{Train}(Z, Y), \neg \text{Train}(X, Y)$

WHY  $Q(n, s)$ ?

# Computing Provenance

---

## Rewriting steps

- **Unifying** the program with specification of outputs of interest
- **Statically annotate** program to indicate interest in success / failure
- Creating **firing rules**
- Checking **connectivity with joins**
- **Computing** the **edge** relation of the provenance graph

# Computing Explanations

---

## Step 1: Unifying the program with the question

- **Limit computation** by binding variables to constants
- Propagating constants in the question throughout the program

## Example

- Binding variables:  $X = n$  (New York) and  $Y = s$  (Seattle)

$$r_1 : Q(X, Y) : \neg \text{Train}(X, Z), \text{Train}(Z, Y), \neg \text{Train}(X, Y)$$

+

WHY  $Q(n, s)$



$$r_1^{\langle X=n, Y=s \rangle} : Q(\underline{n}, \underline{s}) : \neg T(\underline{n}, Z), T(Z, \underline{s}), \neg T(\underline{n}, \underline{s})$$



# Computing Explanations

## Step 2: Annotating rule with success / failure

- **To determine the state of nodes**
- The question provides information about success / failure
- **T**True = Success / Exist                      **F**alse = Failed / Not Exist
- Propagate annotations throughout the unified program

## Example

- ▷ **WHY**  $Q(n,s)$  implies only successful goals need to be captured

$$r_1^{(X=n,Y=s)} : Q(n, s) :- T(n, Z), T(Z, s), \neg T(n, s)$$

+

**WHY**  $Q(n,s)$



$$r_1^{(X=n,Y=s),T} : Q(n, s)^T :- T(n, Z)^T, T(Z, s)^T, \neg T(n, s)^T$$

# Computing Explanations

## Step 3: Create firing rules

- **Capture successful / failed derivations** for the variable binding

### Example

- Create new head by adding the variable **Z** and use firing versions
- Invert the annotation **T** for the goal  $\neg T(n,s)$  in the firing version

$$\begin{array}{l}
 r_1^{(X=n, Y=s), T} : \boxed{Q(n, s)^T} : \underline{-T(n, Z)^T}, \underline{T(Z, s)^T}, \textcircled{\neg T(n, s)^T} \longrightarrow \begin{array}{l}
 F_{Q, T}(n, s) : - F_{r_1, T}(n, s, Z) \\
 \boxed{F_{r_1, T}(n, s, Z)} : \underline{-F_{T, T}(n, Z)}, \underline{F_{T, T}(Z, s)}, \textcircled{F_{T, F}(n, s)} \\
 F_{T, T}(n, Z) : - T(n, Z) \\
 F_{T, T}(Z, s) : - T(Z, s) \\
 \underline{\quad \quad \quad} \quad \quad \quad \underline{\quad \quad \quad}
 \end{array}
 \end{array}$$

# Computing Explanations

## Step 4: Connectivity joins

- Firing rules are **not sufficient to determine** which **subgraphs** of the provenance explain the outputs of interest
- Filter derivations by checking whether connectivity
- Check connectivity from the question node one hop at a time

## Example

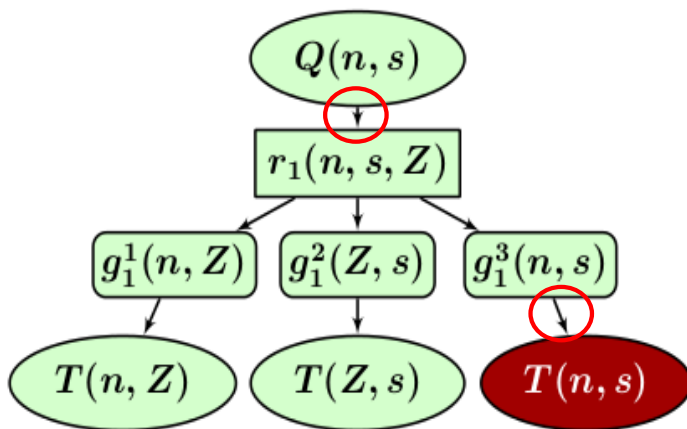
- No guarantees for the nodes in the red box
- Tuple node  $T(n, c)$  is only connected iff  $T(c, s)$  exists

$$\begin{array}{l}
 F_{Q,T}(n, s) :- F_{r_1,T}(n, s, Z) \\
 F_{r_1,T}(n, s, Z) :- F_{T,T}(n, Z), F_{T,T}(Z, s), F_{T,F}(n, s) \\
 F_{T,T}(n, Z) :- T(n, Z) \\
 F_{T,T}(Z, s) :- T(Z, s)
 \end{array}
 \quad \longrightarrow \quad
 \begin{array}{l}
 F_{Q,T}(n, s) :- F_{r_1,T}(n, s, Z) \\
 F_{r_1,T}(n, s, Z) :- F_{T,T}(n, Z), F_{T,T}(Z, s), F_{T,F}(n, s) \\
 FC_{r_2,r_1^1,T}(n, Z) :- T(n, Z), F_{r_1,T}(n, s, Z) \\
 FC_{r_2,r_1^2,T}(Z, s) :- T(Z, s), F_{r_1,T}(n, s, Z) \\
 F_{T,F}(n, s) :- \neg T(n, s)
 \end{array}$$

# Computing Explanations

## Step 5: Computing provenance subgraph edge relation

- Create edges for the provenance graph (explanation)
- Generate rules for the edge relation based on the rule binding information
- Use node identifier  $f_{r_1}^T(n, s, Z)$
- Type of the node, assignments to constants, success/failure state
- Each rule corresponds to a pattern in the graph



- Provenance graph structure -

$$\text{edge}(f_Q^T(n, s), f_{r_1}^T(n, s, Z)) : - F_{r_1, T}(n, s, Z)$$

$$\text{edge}(f_{r_1}^T(n, s, Z), f_{g_1^1}^T(n, Z)) : - F_{r_1, T}(n, s, Z)$$

$$\text{edge}(f_{g_1^1}^T(n, Z), f_T^T(n, Z)) : - F_{r_1, T}(n, s, Z)$$

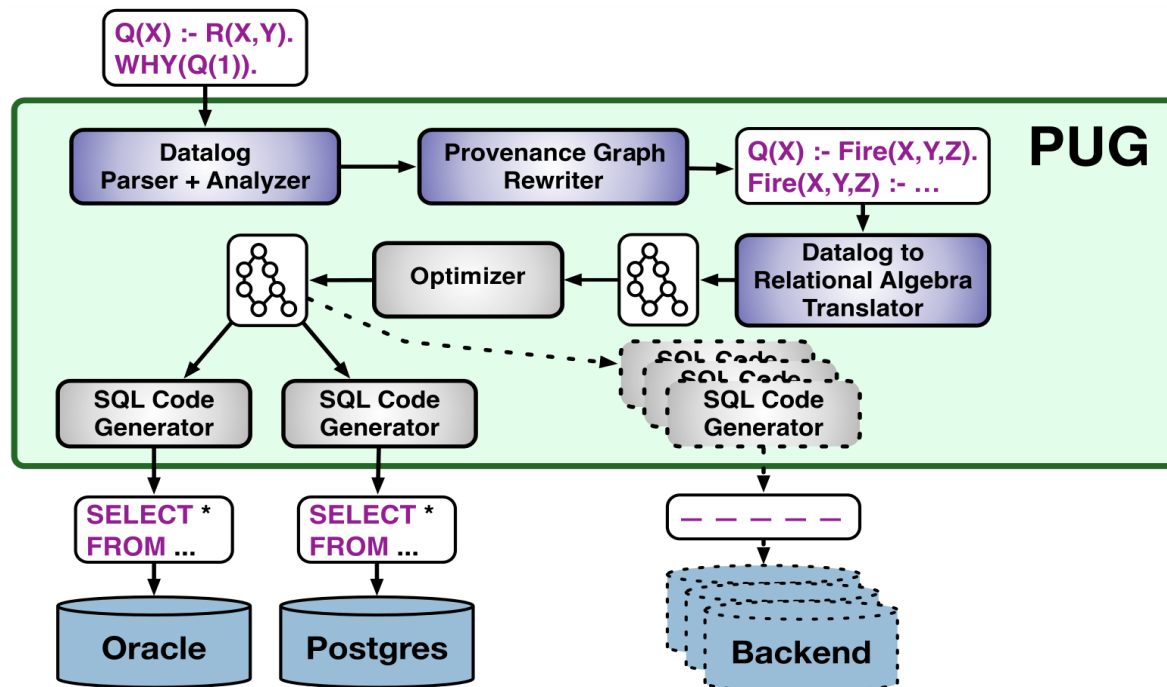
$$\text{edge}(f_{g_1^3}^T(n, s), f_T^F(n, s)) : - F_{r_1, T}(n, s, Z)$$

Example (partial) rules deriving the edge relation

# Implementation

## PUG (Provenance Unification through Graphs) architecture

- Extension of GProM supporting Datalog provenance
- GProM is a SQL+X – to – SQL optimizing compiler
- Relational algebra as IR



## Outline

- Overview & Motivation
- Provenance Graph Model
- Capturing Provenance
- **Factorization**
- Experiments
- Conclusions & Future work

# Recap

## Our provenance graphs can encode factorized polynomials

- What factorization we get is determined by the structure of the program

$$Q() :- S(X), U(X, Y)$$

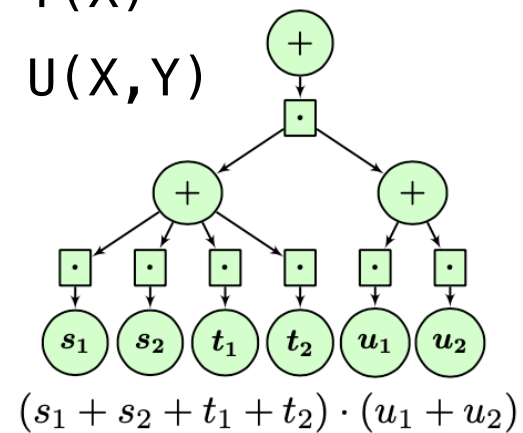
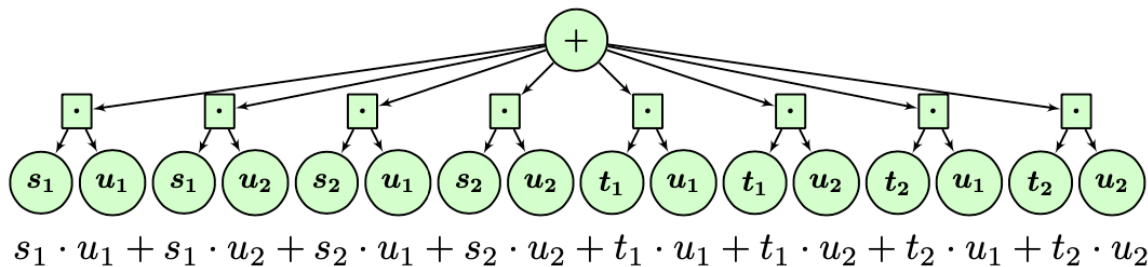
$$Q() :- T(X), U(X, Y)$$

$$Q() :- Q2(X), Q1(X)$$

$$Q1(X) :- S(X)$$

$$Q1(X) :- T(X)$$

$$Q2(X) :- U(X, Y)$$



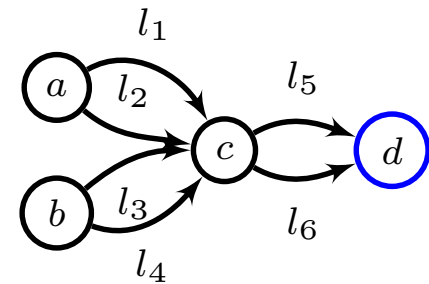
# Utilizing work on factorization

## Approach

- Determine worst-case optimal d-tree factorization
- Rewrite input query to produce this factorization
- Apply capture rewriting

D. Olteanu and J. Závodny`. Size bounds for factorised representations of query results. ACM Transactions on Database Systems (TODS), 40(1):2, 2015.

$r_3 : Q_{2\text{hop}}(X) :- H(Y, L_1, Z), H(Z, L_2, X)$







# Utilizing work on factorization

## Approach

- Determine worst-case optimal d-tree factorization
- Rewrite input query to produce this factorization
- Apply capture rewriting

$$r_5 : Q_{2hop}() :- Q_{L_1}(Z), Q_{L_2}(Z)$$

$$r_{5'} : Q_{L_1}(Z) :- H(Y, L_1, Z)$$

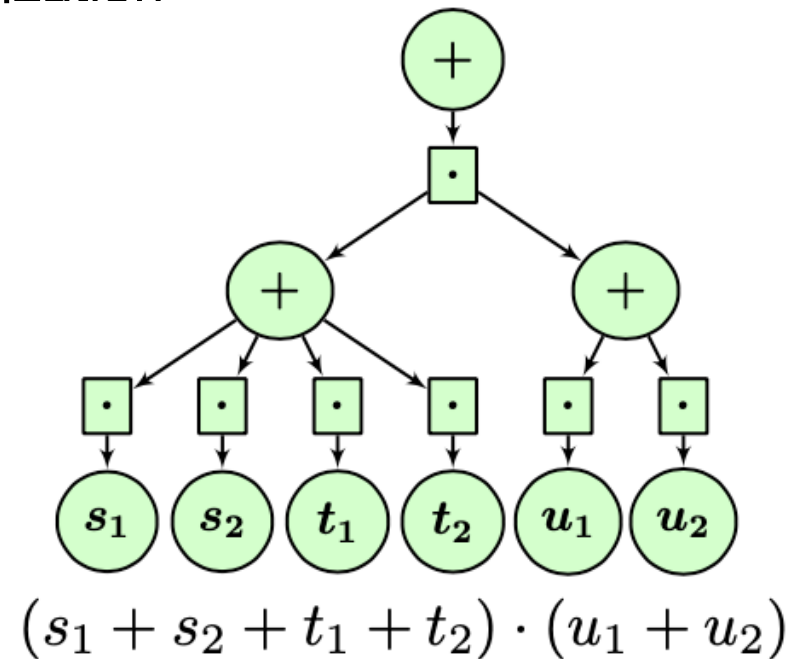
$$r_{5''} : Q_{L_2}(Z) :- H(Z, L_2, d)$$

$$\frac{\frac{\{Z\} L_1 \quad L_2 \quad \{Z\}}{\{Z, L_1\} Y}}$$

Relation H

S	L	E
a	$l_1$	c
a	$l_2$	c
b	$l_3$	c
b	$l_4$	c
c	$l_5$	d
c	$l_6$	d

$s_1$   
 $s_2$   
 $t_1$   
 $t_2$   
 $u_1$   
 $u_2$



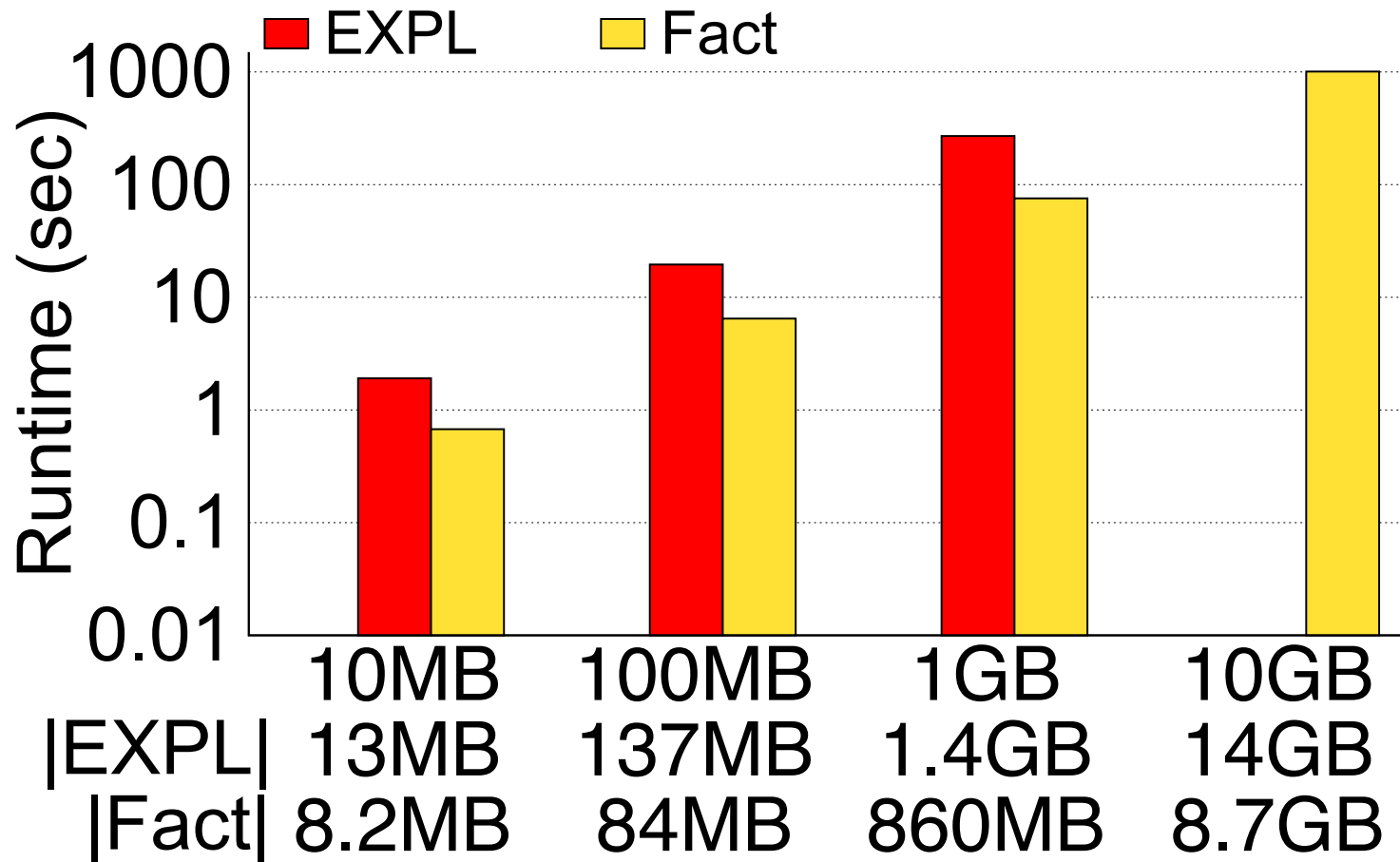
## Outline

- Overview & Motivation
- Provenance Graph Model
- Capturing Provenance
- Factorization
- **Experiments**
- Conclusions & Future work

# Experiments

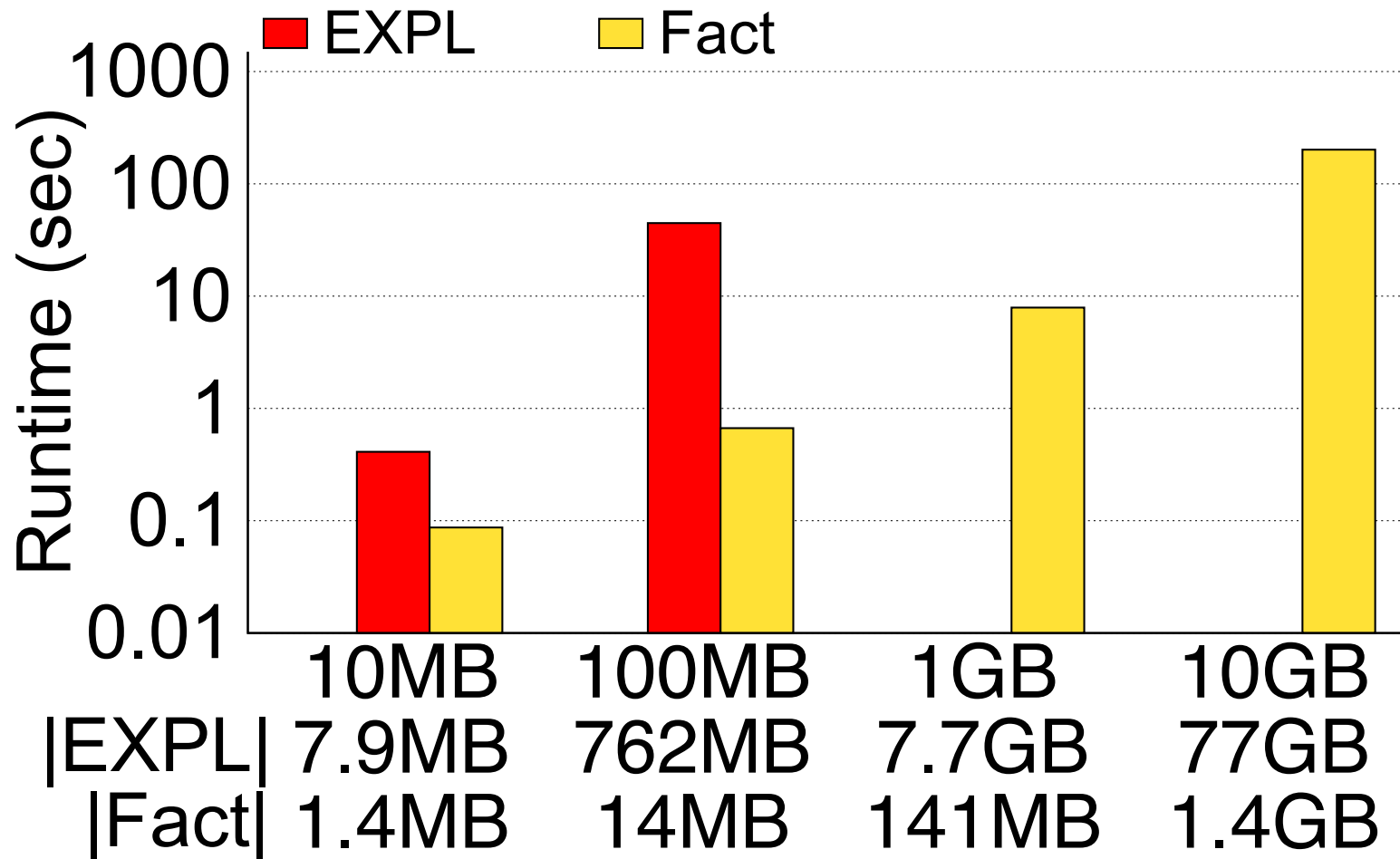
- TPC-H

```
ordDisc(X,Y) :- CUSTOMER(A,X,B,C,D,E,F,G),
                ORDERS(H,A,I,J,K,L,M,O,P),
                LINEITEM(H,Q,R,S,T,U,V,Y,W,Z,A',B',C',D',E',F')
```



# Experiments

- TPC-H `suppCust(N) :- SUPPLIER(A,B,C,N,D,E,F), CUSTOMER(G,H,I,N,J,K,L,M)`



## Outline

---

- Overview & Motivation
- Provenance Graph Model
- Capturing Provenance
- Factorization
- Experiments
- **Conclusions & Future work**

# Conclusions

---

- **Provenance graph model**

- Provenance structure aligns with program structure
- Compatible with well-established provenance models
- Provenance polynomials for positive queries + dual polynomials
- Build-in support for sharing common subexpressions
- Flat relational encoding as edge relation

- **Capturing Provenance**

- Incorporate user's provenance interest into the capture query
- Filter successful / failed assignments upfront (static analysis)
- Output is a query returning the edge relation of the graph

# Conclusions

---

- **Factorization**

- Program structure determines factorization

- **Approximate Summarization of Why-not Provenance**

- Use patterns to summarize provenance
- Use sampling to generate such summaries for very large why-not provenance graphs

## Future work

---

- **Expressiveness**

- Support aggregation and recursion

- **Efficiency**

- Leverage factorized DB techniques for aggregates?
- Factorizing missing answers (complement representations)?

- **Going beyond SQL / Datalog as the target language**

- What specialized algorithms & data structures would be beneficial?



*Questions?*



## Seokki Lee

**Assistant Professor (University of Cincinnati)**

CEAS - Elect Eng & Computer Science

### [Contact Information]

Office: Old Chemistry Building 613B

Email: [\[lee5sk\] at \[ucmail\] o \[uc\] o \[edu\]](mailto:lee5sk@ucmail.uc.edu)

Phone: 513-556-5795

Homepage: <https://researchdirectory.uc.edu/p/lee5sk>

Group page: <https://sites.google.com/view/ucdbg>

PUG: <https://github.com/IITDBGGroup/PUG>



# *Summarizing Provenance*

# Motivation (Example)

- Airbnb (bed and breakfast)

Listing (input)

Id	Name	Ptype	Rtype	NGroup	Neighbor
8403	central place	apt	shared	queen anne	east
9211	plum	apt	entire	ballard	adams
2445	cozy homebase	house	private	queen anne	west
8575	near SpaceNeedle	apt	shared	queen anne	lower
4947	seattle couch	condo	shared	downtown	first hill
2332	modern view	house	entire	queen anne	west

Availability (input)

Id	Date	Price
9211	2016-11-09	130
2445	2016-11-09	45
2332	2016-11-09	350
4947	2016-11-10	40



$r_1: AL(N,R) :- L(I, N, T, R, queen\ anne, E), A(I, 2016-11-09, P)$

“ What are available listings and the room types in Queen Anne on Nov 9<sup>th</sup>, 2016? ”



AvailableListings (output)

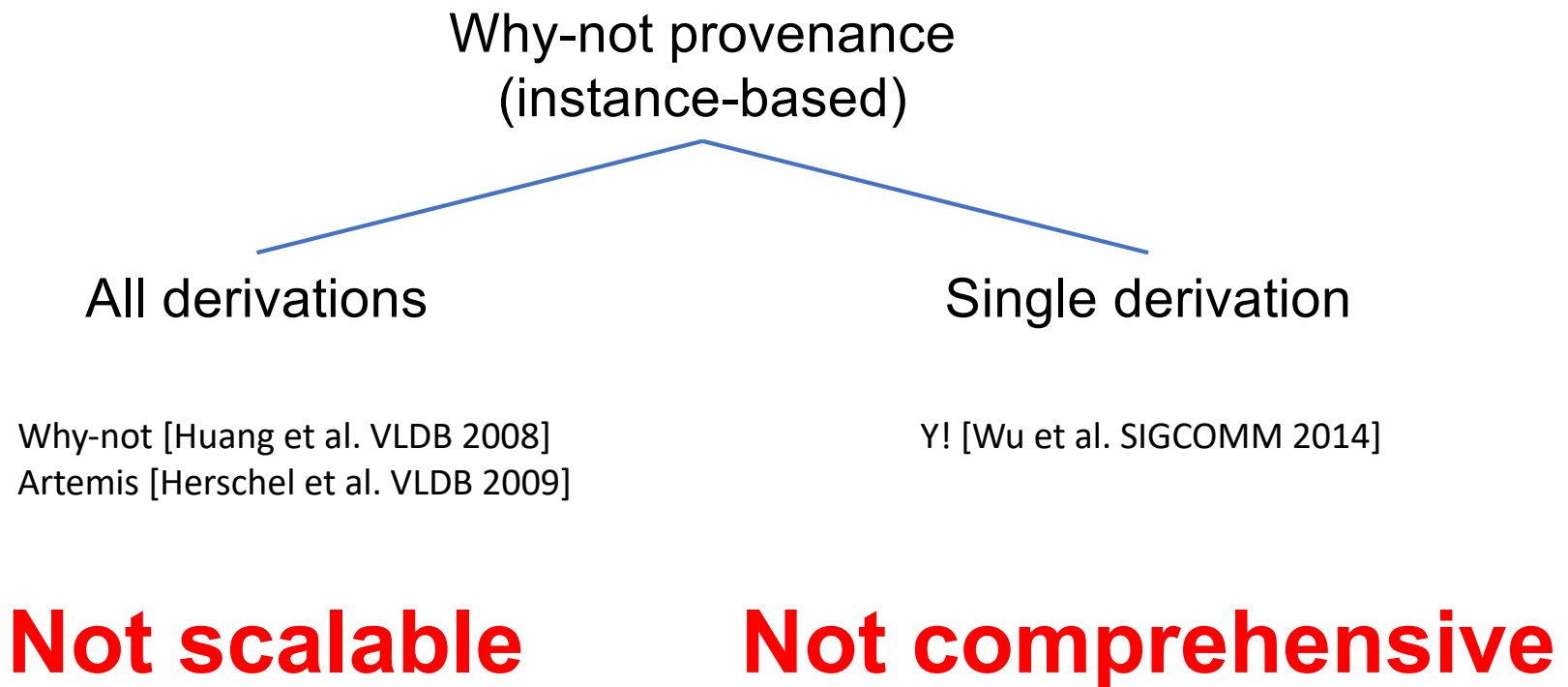
Name	Rtype
cozy homebase	private
modern view	entire

Why no shared room exists?



## Motivation (Example)

---



# Motivation (Example)

Listing (input)

Id	Name	Ptype	Rtype	NGroup	Neighbor
8403	central place	apt	shared	queen anne	east
9211	plum	apt	entire	ballard	adams
2445	cozy homebase	house	private	queen anne	west
8575	near SpaceNeedle	apt	shared	queen anne	lower
4947	seattle couch	condo	shared	downtown	first hill
2332	modern view	house	entire	queen anne	west

Availability (input)

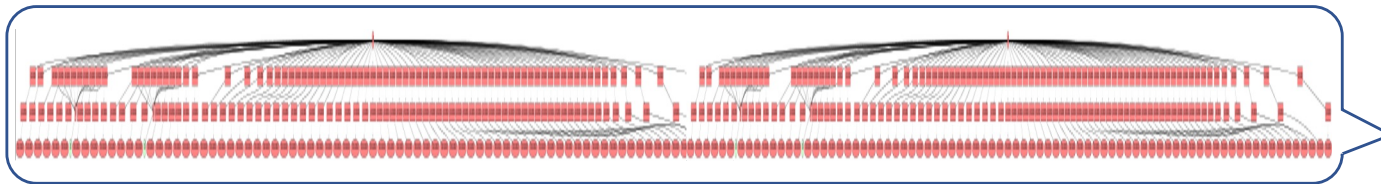
Id	Date	Price
9211	2016-11-09	130
2445	2016-11-09	45
2332	2016-11-09	350
4947	2016-11-10	40

Why no shared room exists?

$r_1: AL(N,R) :- L(I, N, T, R \text{ queen anne, E}), A(I, 2016-11-09, P)$

Attribute	Id	Name	Ptype	Rtype	NGroup	Neighbor	Date	Price
#Distinct Values	6	6	3	3	3	5	2	4

= 2160



~ $15 \cdot 10^{20}$  derivations over full dataset (~ 1.4M)

# Motivation (Example)

Listing (input)

Id	Name	Ptype	Rtype	NGroup	Neighbor
8403	central place	apt	shared	queen anne	east

Availability (input)

Id	Date	Price
----	------	-------

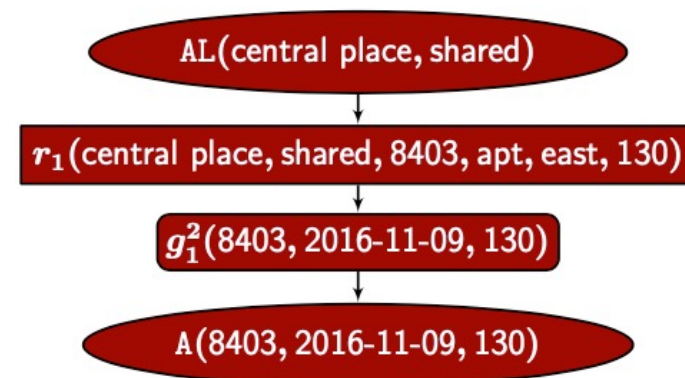
$r_1: AL(N,R) :- L(I, N, T, R, queen\ anne, E), A(I, 2016-11-09, P)$

AvailableListings (output)

Name	Rtype
cozy homebase	private
modern view	entire

Why no shared room exists?

The listing 'central place' has a shared room which is not available at \$130



# Motivation (Example)

Listing (input)

Id	Name	Ptype	Rtype	NGroup	Neighbor
8403	central place	apt	shared	queen anne	east
9211	plum	apt	entire	ballard	adams
2445	cozy homebase	house	private	queen anne	west
8575	near SpaceNeedle	apt	shared	queen anne	lower
4947	seattle couch	condo	shared	downtown	first hill
2332	modern view	house	entire	queen anne	west

Availability (input)

Id	Date	Price
9211	2016-11-09	130
2445	2016-11-09	45
2332	2016-11-09	350
4947	2016-11-10	40

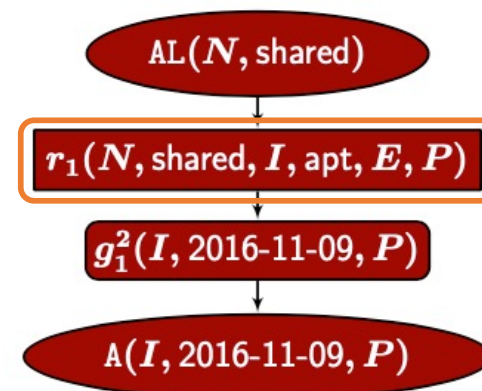
$r_1: AL(N,R) :- L(I, N, T, R, queen\ anne, E), A(I, 2016-11-09, P)$

AvailableListings (output)

Name	Rtype
cozy homebase	private
modern view	entire

Why no shared room exists?

“ All shared rooms of apartments in Queen Anne are not available at any price on Nov 9th, 2016 ”





# Summarizing Provenance

---

- Goals

- **Concise** (small size of explanations)
- **Complete** (covering all provenance)
- **Informative** (providing new insights)

- Challenge

- **Fullfilling all 3 elements** at the same time
- **Computing** summaries using **full why-not provenance**

# Summarizing Provenance

---

- Computing **top-k summaries using patterns**
  - Concise explanations
  - Meaningful (semantically)
- Integrating **sampling** into provenance capture process
  - Unbiased
  - Computing representative patterns
  - Calculating close enough approximate completeness of patterns



# Summarizing Provenance

---

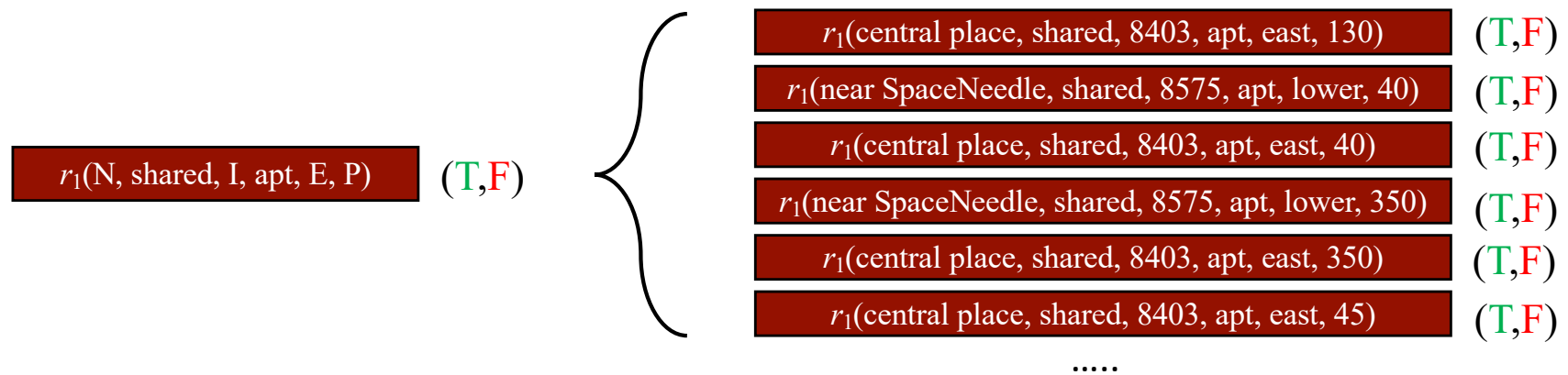
- What are patterns?

$r_1(N, \text{shared}, I, \text{apt}, E, P)$  (T,F)



# Summarizing Provenance

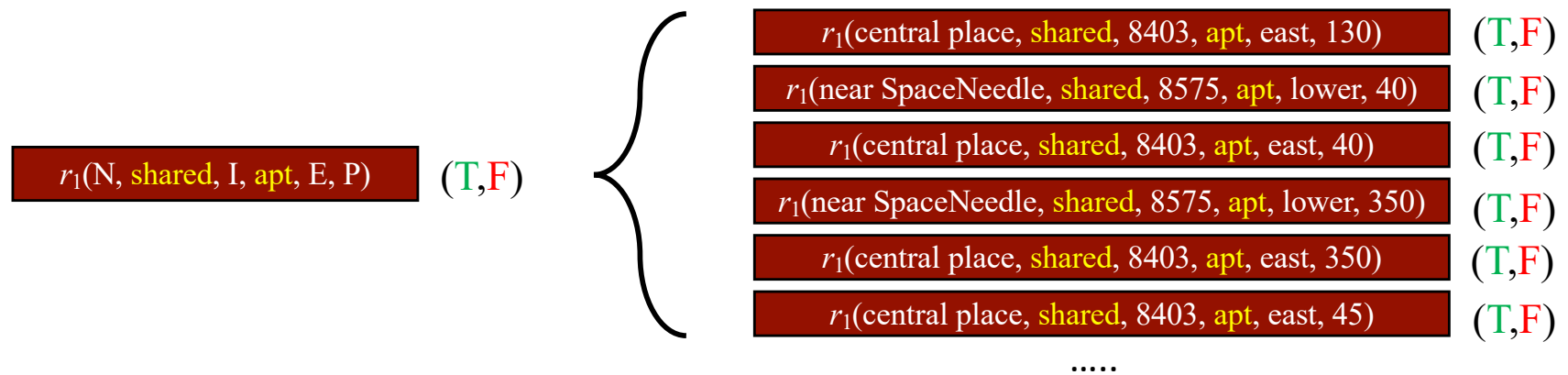
- What are patterns?





# Summarizing Provenance

- What are patterns?





# Summarizing Provenance

- What are provenance summaries?

		Summary	<i>k</i>	
$r_1(\text{N, shared, I, apt, E, P})$	( <b>T</b> , <b>F</b> )	$r_1(\text{N, shared, I, apt, E, 130})$	( <b>F</b> , <b>T</b> )	.....
$r_1(\text{central place, shared, I, apt, E, P})$	( <b>T</b> , <b>F</b> )	$r_1(\text{N, shared, I, house, east, P})$	( <b>F</b> , <b>F</b> )	.....
$r_1(\text{N, shared, I, condo, E, P})$	( <b>T</b> , <b>F</b> )	$r_1(\text{N, shared, I, house, E, 350})$	( <b>T</b> , <b>F</b> )	.....
<i>k</i>	.....	.....		



# Summarizing Provenance

- Quality metrics
  - Completeness (cp): **fraction** of provenance covered by a pattern

$r_1(\text{central place, shared, 8403, apt, east, 130})$

(T,F)  $\approx$

$r_1(\text{N, shared, I, apt, E, P})$

(T,F)

# Summarizing Provenance

- Quality metrics
  - Completeness (cp): **fraction** of provenance covered by a pattern
  - Informativeness (info): **degree of new information** from a pattern

Why no shared  
room exists?

$$r_1(N, \text{shared}, I, \text{apt}, E, P) \quad (\mathbf{T}, \mathbf{F}) \quad > \quad r_1(N, \text{shared}, I, R, E, P) \quad (\mathbf{T}, \mathbf{F})$$



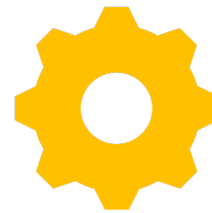
# Summarizing Provenance

- How to compute summaries
  - Heuristic using a sample of why-not provenance

Size of the summary ( $k$ )  
 $Q(A) :- R(A,B), \neg S(B)$



Input



Summarization  
process

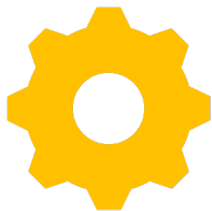


Provenance  
summary

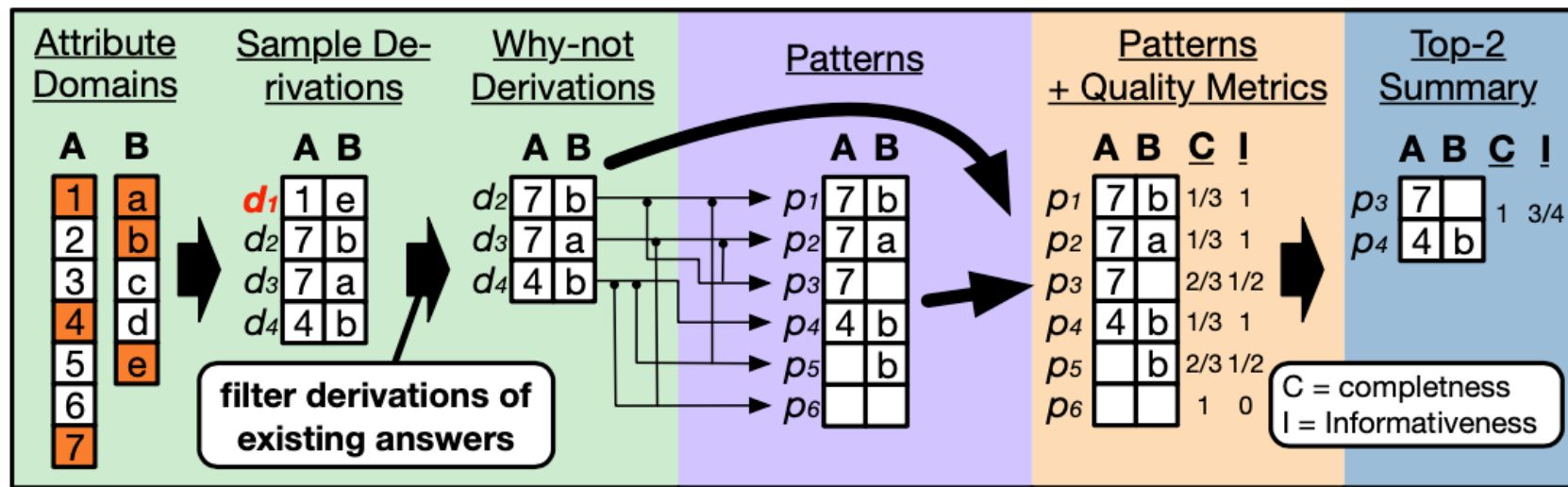


# Summarizing Provenance

- How to compute summaries



Summarization process



Sampling Why-not Provenance

Generating Patterns

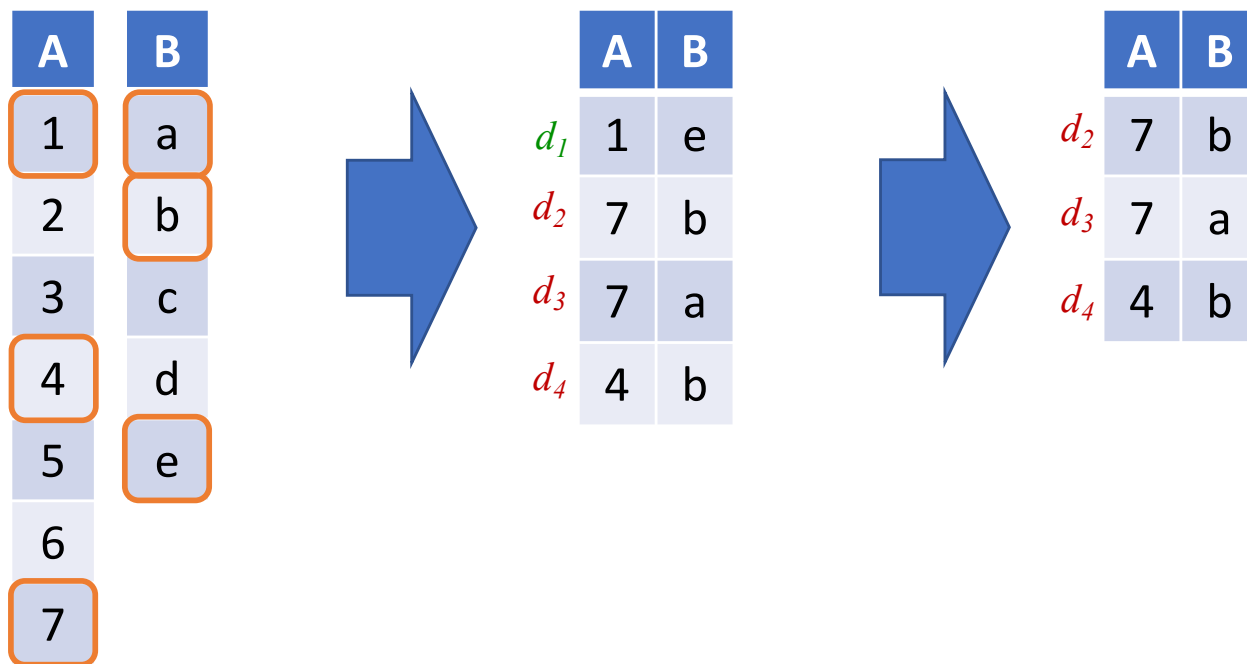
Measuring Quality

Selecting top-k



# Summarizing Provenance

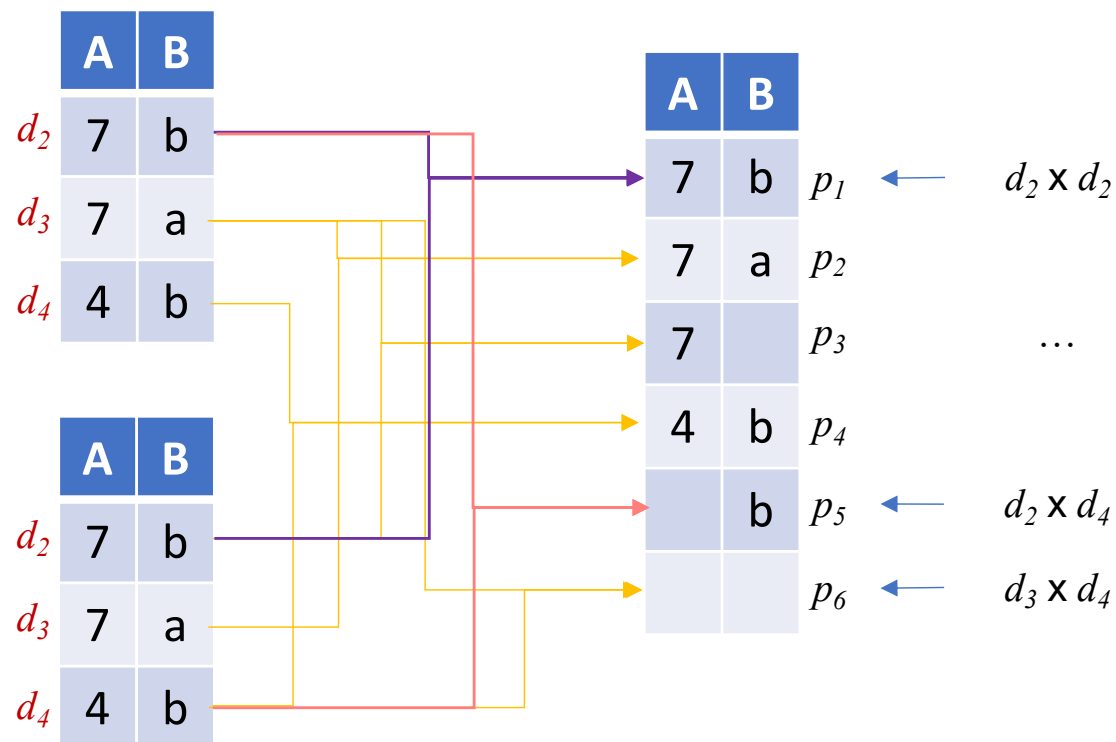
- Sampling why-not provenance
  - Generating a sample that is equivalent to uniform random sample
    - Without computing full why-not provenance
  - Batch sampling: generating a query that returns an unbiased sample





# Summarizing Provenance

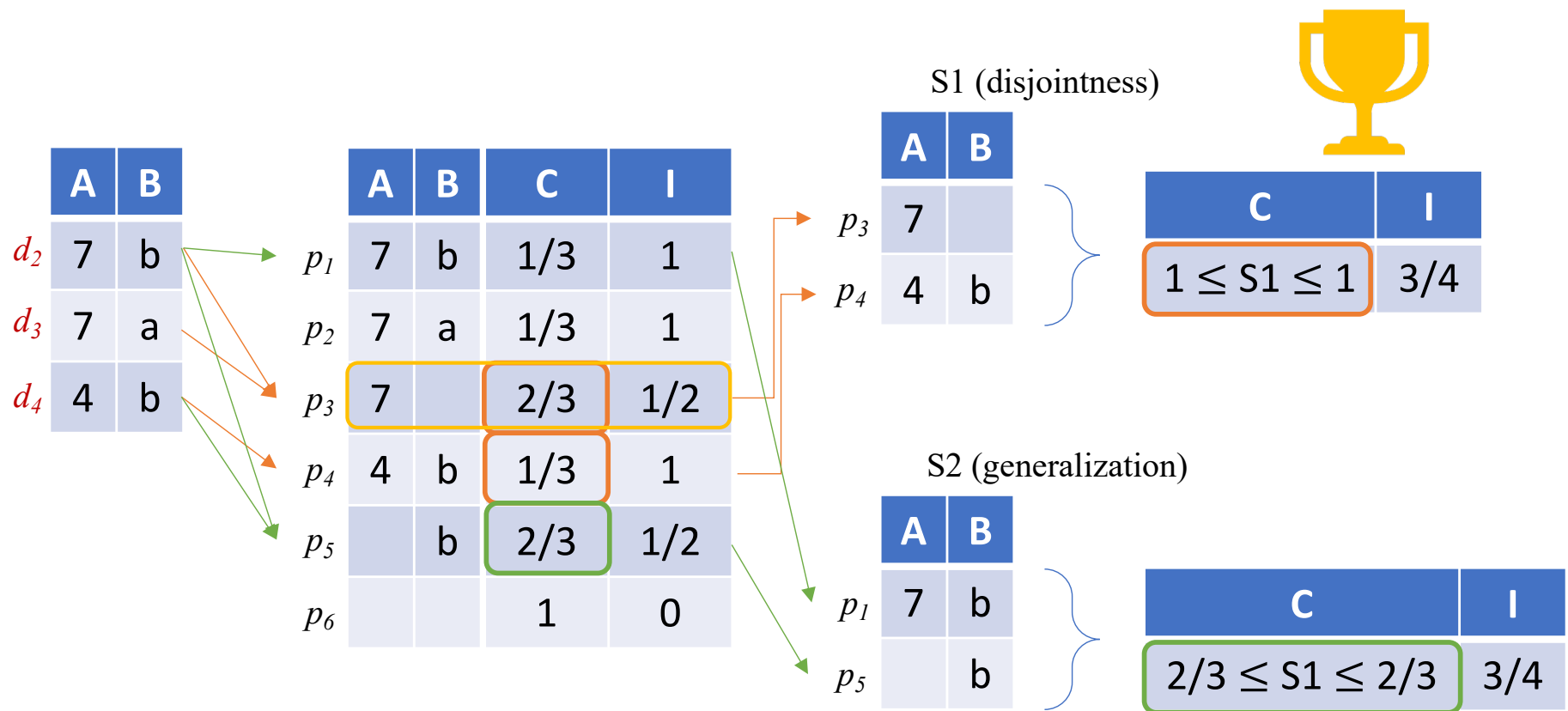
- Generating patterns
  - LCA (Lowest Common Ancestor)





# Summarizing Provenance

- Measuring quality and selecting top- $k$  patterns



C: completeness  
I: Informativeness

# *Experiments*

# Experiments

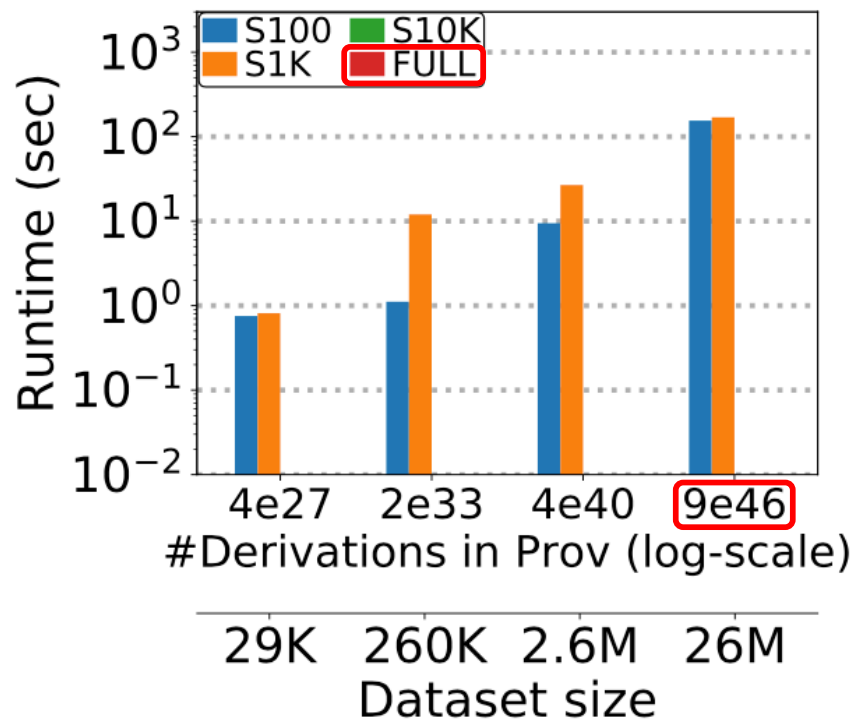
---

- **Performance** of computing summaries
- **Quality** of summaries
- **Comparison** with other approaches
  
- Datasets
  - 4 real-world datasets
  - TPCH
  
- Queries
  - Single rule through multiple rules
  - Negation and comparisons

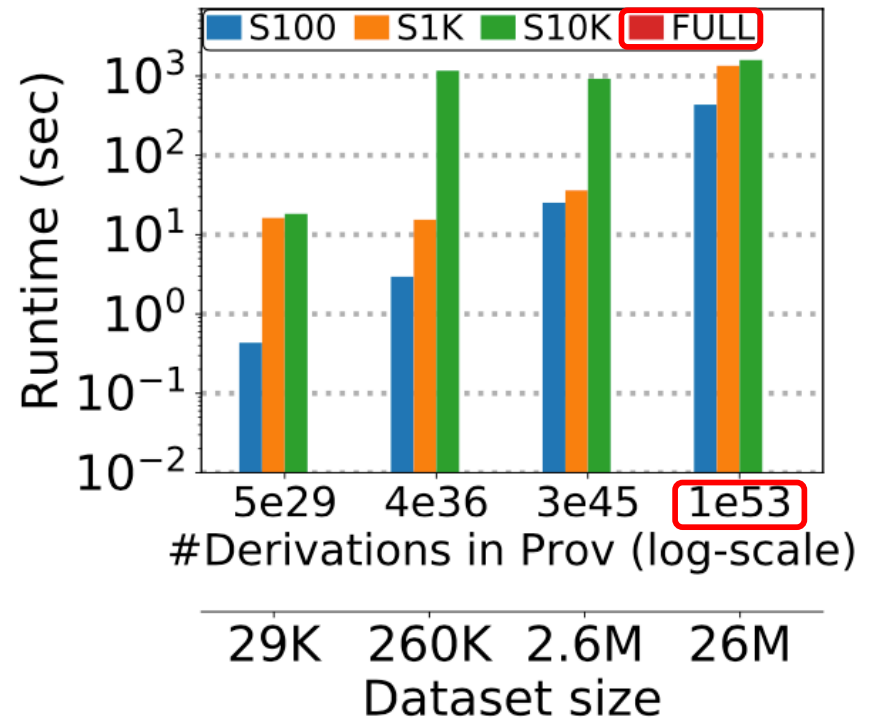
# Experiments

- Computing summaries for  $\sim 10^{50}$  derivations

12 variables + 6 goals + negation



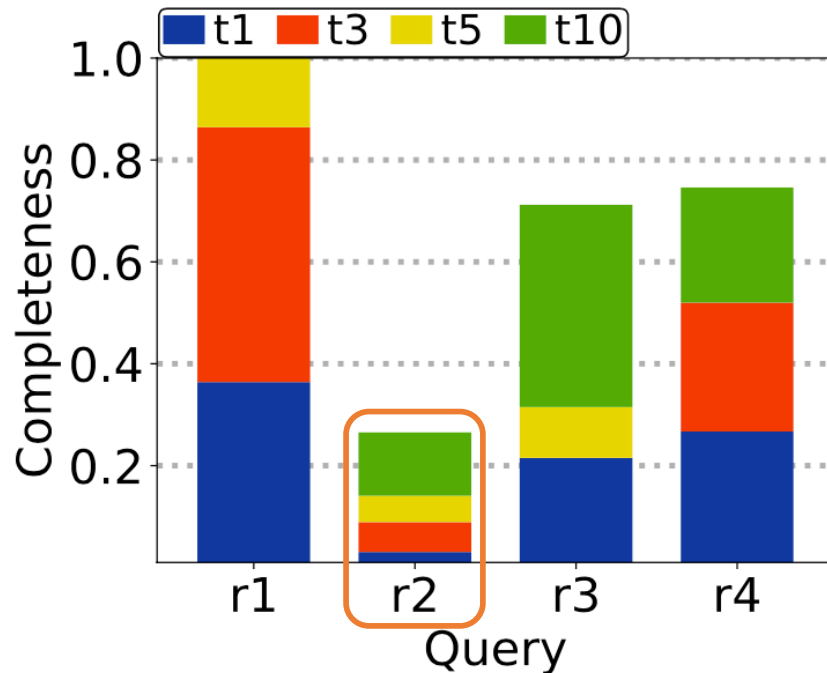
13 variables + 4 goals + multiple rules



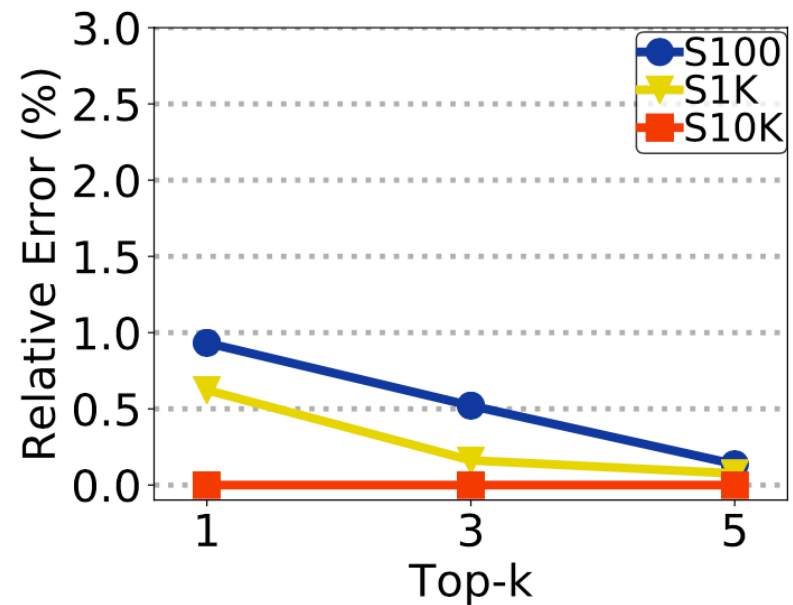


# Experiments

- Generating high-quality summaries



[Completeness comparison]



[Quality metric error caused by sampling]