

CS 595 - Hot topics in database systems:

Data Provenance

I. Database Provenance

I.1 Provenance Models and Systems

Boris Glavic

September 12, 2012

Witness Lists - Example

Example

$$q = \bowtie_{itemId=id} (sales, items)$$

$$\langle s_1, i_1 \rangle \in \mathcal{W}(q, l)$$

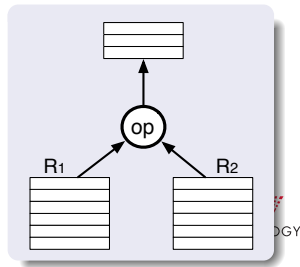
Example

$$[[\bowtie_{itemId=id} (sales, items)]]$$

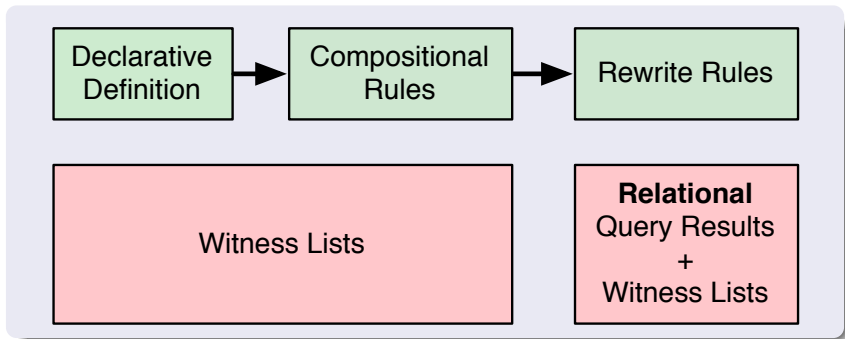
	shop	itemId	id	price
t_1	Migros	1	1	100
t_2	Migros	3	3	25
t_3	Coop	3	3	25

	sales	
	shop	itemId
s_1	Migros	1
s_2	Migros	3
s_3	Coop	3

	items	
	id	price
i_1	1	100
i_2	2	10
i_3	3	25



Approach Roadmap

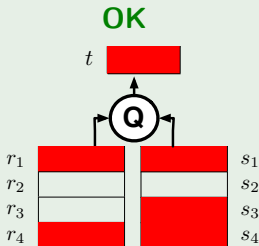


Definition Conditions

Condition (1)

- The provenance of t^x generates **exactly** t
- \Rightarrow Executing the query over the provenance returns t
- **(1):** $[[op(calPI(op, t, I))]] = \{t^x\}$

Example



Condition (1) - Example

Example

```

SELECT *
FROM sales, items
WHERE itemId = id

```

Example

sales		items		[[⋈ _{itemId=id} (sales, items)]]				
	shop	itemId	id	price	shop	itemId	id	price
<i>s</i> ₁	Migros	1	<i>i</i> ₁	100	<i>t</i> ₁	Migros	1	100
<i>s</i> ₂	Migros	3	<i>i</i> ₂	10	<i>t</i> ₂	Migros	3	25
<i>s</i> ₃	Coop	3	<i>i</i> ₃	25	<i>t</i> ₃	Coop	3	25

Condition (1) - Example

Example

- $\mathcal{PI}_1 = \{ \langle s_1, i_1 \rangle \}$: YES: $[[op(\mathcal{PI}_1)]] = \{ t_1 \}$

Example

	sales	
	shop	itemId
s_1	Migros	1
s_2	Migros	3
s_3	Coop	3

	items	
	id	price
i_1	1	100
i_2	2	10
i_3	3	25

	$[[\bowtie_{itemId=id}(sales, items)]]$			
	shop	itemId	id	price
t_1	Migros	1	1	100
t_2	Migros	3	3	25
t_3	Coop	3	3	25

Condition (1) - Example

Example

- $\mathcal{PI}_1 = \{ \langle s_1, i_1 \rangle \}$: **YES**: $[[op(\mathcal{PI}_1)]] = \{t_1\}$
- $\mathcal{PI}_2 = \{ \langle s_1, i_1 \rangle, \langle s_2, i_3 \rangle \}$: **NO**: $[[op(\mathcal{PI}_2)]] = \{t_1, t_2\}$
 - $s_2 + i_3$ irrelevant
- $\mathcal{PI}_3 = \{ \langle s_2, i_1 \rangle \}$: **NO**: $[[op(\mathcal{PI}_3)]] = \emptyset$

Example

		sales	
		shop	itemId
s_1		Migros	1
s_2		Migros	3
s_3		Coop	3

		items	
		id	price
i_1		1	100
i_2		2	10
i_3		3	25

		$[[\bowtie_{itemId=id} (sales, items)]]$			
		shop	itemId	id	price
t_1		Migros	1	1	100
t_2		Migros	3	3	25
t_3		Coop	3	3	25

Condition (1) - Example

Example

- $PI_1 = \{ \langle s_1, i_1 \rangle \}$: **YES**: $[[op(PI_1)]] = \{t_1\}$
- $PI_2 = \{ \langle s_1, i_1 \rangle, \langle s_2, i_3 \rangle \}$: **NO**: $[[op(PI_2)]] = \{t_1, t_2\}$
 - $s_2 + i_3$ irrelevant
- $PI_3 = \{ \langle s_2, i_1 \rangle \}$: **NO**: $[[op(PI_3)]] = \emptyset$
- $PI_4 = \{ \langle s_1, i_1 \rangle, \langle s_1, i_2 \rangle \}$: **YES**: $[[op(PI_4)]] = \{t_1\}$
 - i_2 irrelevant

Example

		sales	
		shop	itemId
s_1		Migros	1
s_2		Migros	3
s_3		Coop	3

		items	
		id	price
i_1		1	100
i_2		2	10
i_3		3	25

				$[[\bowtie_{itemId=id}(sales, items)]]$			
				shop	itemId	id	price
t_1				Migros	1	1	100
t_2				Migros	3	3	25
t_3				Coop	3	3	25

Definition Conditions Cont.

Condition (2)

- Each witness list *contributes* something to t
- Evaluating the query over the witness list \Rightarrow result not empty
- (2) $\forall w \in \mathcal{PI}(op, t, I) : [[op(w)]] \neq \emptyset$

Condition (2) - Example

Example

- $\mathcal{PI}_1 = \{ \langle s_1, i_1 \rangle \}$: YES: $[[op(\langle s_1, i_1 \rangle)]] = \{t_1\}$
- $\mathcal{PI}_2 = \{ \langle s_1, i_1 \rangle, \langle s_2, i_3 \rangle \}$: YES: $[[op(\langle s_2, i_3 \rangle)]] = \{t_2\}$
 - $s_2 + i_3$ irrelevant

Example

	sales	
	shop	itemId
s_1	Migros	1
s_2	Migros	3
s_3	Coop	3

	items	
	id	price
i_1	1	100
i_2	2	10
i_3	3	25

	[[$\bowtie_{itemId=id}$ (sales, items)]]			
	shop	itemId	id	price
t_1	Migros	1	1	100
t_2	Migros	3	3	25
t_3	Coop	3	3	25

Definition Conditions Cont.

Condition (3)

- Avoid subsumed witness lists
- (3) $w, w' \in \mathcal{W}(q, l) : w \prec w' \wedge w \in \mathcal{PI}(q, t, l) \Rightarrow w' \notin \mathcal{PI}(q, t, l)$
- $w \prec w' : (\forall i : w[i] = w'[i] \vee w[i] = \perp) \wedge (\exists i : w'[i] \neq \perp \wedge w[i] = \perp)$
- Similarity to Minimal elements for Why-provenance!

Condition (3) and Subsumption example

Example

$$w = \langle t_1, \perp \rangle$$
$$w'' = \langle t_1, t_3 \rangle$$

$$w' = \langle t_1, t_2 \rangle$$
$$w''' = \langle \perp, t_3 \rangle$$

Condition (3) and Subsumption example

Example

$$w = \langle t_1, \perp \rangle$$

$$w' = \langle t_1, t_2 \rangle$$

$$w'' = \langle t_1, t_3 \rangle$$

$$w''' = \langle \perp, t_3 \rangle$$

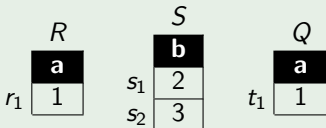
- $w \prec w'$: YES

Condition (3) and Subsumption example

Example

- $q = op(R, S) = (R - S)$
- $\mathcal{PI}(q, t, l) = \{ \langle r_1, \perp \rangle \}$

Example



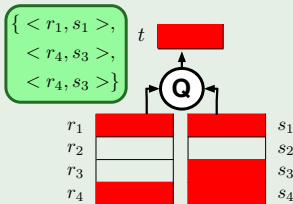
Definition Conditions Cont.

Condition (4)

- No contributing tuples omitted
- \Rightarrow Provenance is maximal set fulfilling conditions (1),(2),(3)
- (4): $\nexists \mathcal{P} \supset \mathcal{P}' \subseteq \mathcal{W}(q, I) : \mathcal{P}' \models (1), (2), (3)$

Example

OK



Definition Application

Example

Result	
a	
1	

t₁



```
SELECT DISTINCT R.a
FROM R, S
WHERE R.a = S.c;
```

R	
a	b
1	2
3	6

r₁

r₂

S	
c	d
1	2
1	3

s₁

s₂

Definition Application

Provenance

$$\{ \langle r_1, s_1 \rangle, \langle r_1, s_2 \rangle \}$$

Example

Result

	a
t_1	1



```
SELECT DISTINCT R.a
FROM R, S
WHERE R.a = S.c;
```



R

	a	b
r_1	1	2
r_2	3	6



S

	c	d
s_1	1	2
s_2	1	3

Definition Application

Provenance

$$\{ \langle r_1, s_1 \rangle, \langle r_1, s_2 \rangle \}$$

Conditions?

- 1 Returns t_1
- 2 Each witness list contributes
- 3 No subsumptions

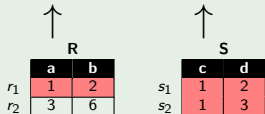
Example

Result	
a	
1	

t_1



```
SELECT DISTINCT R.a
FROM R, S
WHERE R.a = S.c;
```



Definition Application

Provenance

$$\{ \langle r_1, s_1 \rangle, \langle r_1, s_2 \rangle \}$$

Conditions?

- 1 Returns t_1
- 2 Each witness list contributes
- 3 No subsumptions
- 4 Cannot add anything without breaking 1, 2, or 3

Example

Result

a
1

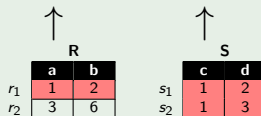
t_1



```

SELECT DISTINCT R.a
FROM R, S
WHERE R.a = S.c;

```



Definition Recap

Definition

Perm Influence Contribution Semantics (PI-CS) $\mathcal{PI}(op, t, I)$, the provenance of t^x from the result of op over instance I , is the unique subset of $\mathcal{W}(op, I)$ that fulfills the following conditions:

- (1): $[[op(\mathcal{PI}(op, t, I))]] = \{t^x\}$
- (2) $\forall w \in \mathcal{PI}(op, t, I) : [[op(w)]] \neq \emptyset$
- (3): $w, w' \in \mathcal{W}(q, I) : w \prec w' \wedge w \in \mathcal{PI}(q, t, I) \Rightarrow w' \notin \mathcal{PI}(q, t, I)$
- (4): $\exists \mathcal{P} \supset \mathcal{P}' \subseteq \mathcal{W}(q, I) : \mathcal{P}' \models (1), (2), (3)$

Provenance for Queries

Transitivity

- Recursive definition similar to Lineage
- Query = operator(subquery): $q = op(q_1)$
- Compute provenance for operator over result of subquery
- Substitute tuples in witness lists with their provenance according to subquery



Compositional Rules

4

Outline

1 Witness List based Provenance Models

- Introduction
- Provenance Model
- Compositional Rules
- Provenance Representation
- Provenance Generation through Query Rewrite
- Implementation
- Recap

Compositional Rules

- A more computable form than declarative definition
- Similar idea as for Lineage
- Prove equivalence with declarative definition

Compositional Rules Definition

Definition (Compositional Rules)

$$\mathcal{PI}(R, t) = \{ \langle t \rangle^n \mid t^n \in R \}$$

$$\mathcal{PI}(\sigma_C(q_1), t) = \mathcal{PI}(q_1, t)$$

$$\mathcal{PI}(\pi_A(q_1), t) = \{ w^n \mid w^n \in \mathcal{PI}(q_1, u) \wedge u.A = t \}$$

$$\mathcal{PI}(\alpha_{G,agg}(q_1), t) = \{ w^n \mid w^n \in \mathcal{PI}(q_1, u) \wedge u.G = t.G \} \\ \cup \{ \langle \perp \rangle \mid Q_1 = \emptyset \wedge \|G\| = 0 \}$$

$$\mathcal{PI}(q_1 \bowtie_C q_2, t) = \{ (w_1 \blacktriangleright w_2)^{n \times m} \mid w_1^n \in \mathcal{PI}(q_1, t.Q_1) \\ \wedge w_2^m \in \mathcal{PI}(q_2, t.Q_2) \}$$

$$\mathcal{PI}(q_1 \boxtimes_C q_2, t) =$$

$$\begin{cases} \{ (w \blacktriangleright \perp(q_2))^n \mid w^n \in \mathcal{PI}(q_1, t.Q_1) \} & \text{if } t \not\bowtie C \\ \mathcal{PI}(q_1 \bowtie_C q_2, t) & \text{else} \end{cases}$$

Compositional Rules Definition

Definition (Compositional Rules)

$$\mathcal{PI}(q_1 \cup q_2, t) = \{(w \blacktriangleright \perp(q_2))^n \mid w^n \in \mathcal{PI}(q_1, t)\} \\ \cup \{(\perp(q_1) \blacktriangleright w)^n \mid w^n \in \mathcal{PI}(q_2, t)\}$$

$$\mathcal{PI}(q_1 \cap q_2, t) = \{(w_1 \blacktriangleright w_2)^{n \times m} \mid w_1^n \in \mathcal{PI}(q_1, t) \\ \wedge w_2^m \in \mathcal{PI}(q_2, t)\}$$

$$\mathcal{PI}(q_1 - q_2, t) = \{(w \blacktriangleright \perp(q_2))^n \mid w^n \in \mathcal{PI}(q_1, t)\}$$

Compositional Semantics Example

Example

```

CREATE VIEW RevenueFirstQ
SELECT Shop, sum(Revenue) AS Revenue
FROM MonthlyRevenue
WHERE Month < 5
GROUP BY Shop

```

RevenueFirstQ

Shop	Revenue
New York	2265

 t_1 **MonthlyRevenue/Q₁**

Shop	Month	Revenue
New York	1	2247
New York	3	18
Wuppertal	5	9

 m_1/s_1 m_2/s_2 m_3

Compositional Semantics Example

Example

$$q = \alpha_{shop, sum(revenue)}(q_1)$$

$$q_1 = \sigma_{month < 5}(M)$$

RevenueFirstQ

Shop	Revenue
New York	2265

 t_1
MonthlyRevenue/Q₁

Shop	Month	Revenue
New York	1	2247
New York	3	18
Wuppertal	5	9

 m_1/s_1
 m_2/s_2
 m_3

Compositional Semantics Example

Example

$$q = \alpha_{shop, sum(revenue)}(q_1) \quad q_1 = \sigma_{month < 5}(M)$$

$$\mathcal{PI}(q, t_1) = \mathcal{PI}(\alpha_{shop, sum(revenue)}(q_1), t_1)$$

$$\mathcal{PI}(\alpha_{shop, sum(revenue)}(q_1), t) = \{w^n \mid w^n \in \mathcal{PI}(q_1, u) \\ \wedge u.shop = t.shop\}$$

$$\mathcal{PI}(q_1, t) = \mathcal{PI}(M, t)$$

RevenueFirstQ

Shop	Revenue
New York	2265

 t_1 MonthlyRevenue/ Q_1

Shop	Month	Revenue
New York	1	2247
New York	3	18
Wuppertal	5	9

 m_1/s_1 m_2/s_2 m_3

OGY

Compositional Semantics Example

Example

$$q = \alpha_{shop, sum(revenue)}(q_1) \quad q_1 = \sigma_{month < 5}(M)$$

$$\mathcal{PI}(q, t_1) = \mathcal{PI}(\alpha_{shop, sum(revenue)}(q_1), t_1)$$

$$\mathcal{PI}(\alpha_{shop, sum(revenue)}(q_1), t_1) = \{ \langle s_1 \rangle, \langle s_2 \rangle \}$$

$$\mathcal{PI}(q_1, s_1) = \langle m_1 \rangle$$

$$\mathcal{PI}(q_1, s_2) = \langle m_2 \rangle$$

RevenueFirstQ

Shop	Revenue
New York	2265

 t_1 MonthlyRevenue/ Q_1

Shop	Month	Revenue
New York	1	2247
New York	3	18
Wuppertal	5	9

 m_1/s_1 m_2/s_2 m_3

OGY

Compositional Semantics Example

Example

$$q = \alpha_{shop, sum(revenue)}(q_1) \quad q_1 = \sigma_{month < 5}(M)$$

$$\mathcal{PI}(q, t_1) = \mathcal{PI}(\alpha_{shop, sum(revenue)}(q_1), t_1)$$

$$\mathcal{PI}(\alpha_{shop, sum(revenue)}(q_1), t_1) = \{ \langle m_1 \rangle, \langle m_2 \rangle \}$$

$$\mathcal{PI}(q_1, s_1) = \langle m_1 \rangle$$

$$\mathcal{PI}(q_1, s_2) = \langle m_2 \rangle$$

RevenueFirstQ

Shop	Revenue
New York	2265

t_1

MonthlyRevenue/ Q_1

Shop	Month	Revenue
New York	1	2247
New York	3	18
Wuppertal	5	9

m_1/s_1

m_2/s_2

m_3

OGY

Compositional Semantics Example

Example

$$q = \alpha_{shop, sum(revenue)}(q_1)$$

$$q_1 = \sigma_{month < 5}(M)$$

$$PI(q, t_1) = \{ \langle m_1 \rangle, \langle m_2 \rangle \}$$

$$PI(\alpha_{shop, sum(revenue)}(q_1), t_1) = \{ \langle m_1 \rangle, \langle m_2 \rangle \}$$

$$PI(q_1, s_1) = \langle m_1 \rangle$$

$$PI(q_1, s_2) = \langle m_2 \rangle$$

RevenueFirstQ

Shop	Revenue
New York	2265

 t_1
MonthlyRevenue/ Q_1

Shop	Month	Revenue
New York	1	2247
New York	3	18
Wuppertal	5	9

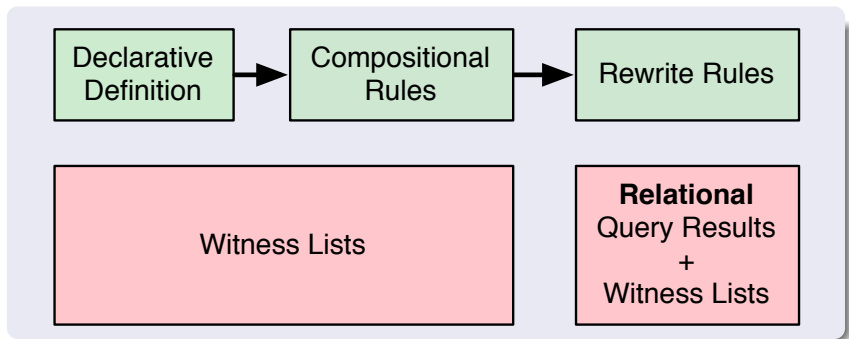
 m_1/s_1
 m_2/s_2
 m_3

OGY

Outline

- 1 Witness List based Provenance Models
 - Introduction
 - Provenance Model
 - Compositional Rules
 - Provenance Representation
 - Provenance Generation through Query Rewrite
 - Implementation
 - Recap

Approach Roadmap - Compositional Rules



Requirements for Provenance Representation

Problem Statement

- For each result tuple t
- Several *contributing* witness lists
- How to represent this information?

Requirements for Provenance Representation

Problem Statement

- For each result tuple t
- Several *contributing* witness lists
- How to represent this information?

Requirements

- 1 Comprehensible for human
 - Should help user to understand data
- 2 Query-able
 - Because provenance can be huge
- 3 Relationship data and provenance
 - Queries that combine both

Provenance Representation

Problem Statement

- For each result tuple t
- Several *contributing* witness lists
- How to represent this information?

Example (Provenance)

$$\{ \langle r_1, s_1 \rangle, \langle r_1, s_2 \rangle \}$$

Example

	Result		
t_1	<table><thead><tr><th>a</th></tr></thead><tbody><tr><td>1</td></tr></tbody></table>	a	1
a			
1			



```
SELECT DISTINCT R.a
FROM R, S
WHERE R.a = S.c;
```

	R		S									
r_1	<table><thead><tr><th>a</th><th>b</th></tr></thead><tbody><tr><td>1</td><td>2</td></tr></tbody></table>	a	b	1	2		<table><thead><tr><th>c</th><th>d</th></tr></thead><tbody><tr><td>1</td><td>2</td></tr></tbody></table>	c	d	1	2	
a	b											
1	2											
c	d											
1	2											
r_2	<table><tbody><tr><td>3</td><td>6</td></tr></tbody></table>	3	6		<table><tbody><tr><td>1</td><td>3</td></tr></tbody></table>	1	3					
3	6											
1	3											

Provenance Representation

Problem Statement

- For each result tuple t
- Several *contributing* witness lists
- How to represent this information?

Example (Provenance)

$$\{ \langle r_1, s_1 \rangle, \langle r_1, s_2 \rangle \}$$

Solution

- Provenance + normal tuples in *single* relation

Example

	Result		
t_1	<table><thead><tr><th>a</th></tr></thead><tbody><tr><td>1</td></tr></tbody></table>	a	1
a			
1			



```
SELECT DISTINCT R.a  
FROM R, S  
WHERE R.a = S.c;
```

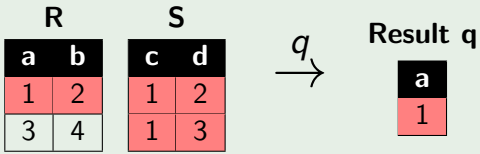
	R	S								
r_1	<table><thead><tr><th>a</th><th>b</th></tr></thead><tbody><tr><td>1</td><td>2</td></tr></tbody></table>	a	b	1	2	<table><thead><tr><th>c</th><th>d</th></tr></thead><tbody><tr><td>1</td><td>2</td></tr></tbody></table>	c	d	1	2
a	b									
1	2									
c	d									
1	2									
r_2	<table><tbody><tr><td>3</td><td>6</td></tr></tbody></table>	3	6	<table><tbody><tr><td>1</td><td>3</td></tr></tbody></table>	1	3				
3	6									
1	3									

Solution

Relation with Provenance + Normal Data

- **Data:** result tuple + all tuples from a witness list
 - Result tuple might have to be duplicated!
- **Schema:** + input attributes (renamed)
 - Generate self-explanatory names: P(a)

Example (Normal Query Result)



Solution

Relation with Provenance + Normal Data

- **Data:** result tuple + all tuples from a witness list
 - Result tuple might have to be duplicated!
- **Schema:** + input attributes (renamed)
 - Generate self-explanatory names: P(a)

Example (Provenance Representation)

R	
a	b
1	2
3	4

S	
c	d
1	2
1	3



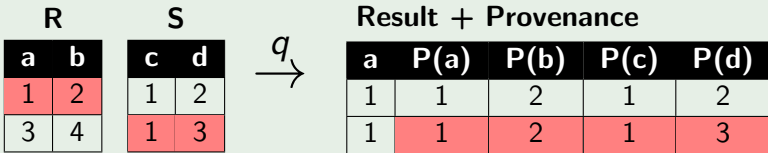
Result + Provenance				
a	P(a)	P(b)	P(c)	P(d)
1	1	2	1	2
1	1	2	1	3

Solution

Relation with Provenance + Normal Data

- **Data:** result tuple + all tuples from a witness list
 - Result tuple might have to be duplicated!
- **Schema:** + input attributes (renamed)
 - Generate self-explanatory names: P(a)

Example (Provenance Representation)



Representation Definition

Definition (Relational Provenance Representation)

The relational representation Q^{PI} for the PI-CS provenance of a query q is defined as:

$$Q^{PI} = \{(t \blacktriangleright w[1]' \blacktriangleright \dots \blacktriangleright w[n]')^m \mid t^p \in Q \wedge w^m \in \mathcal{PI}(q, t)\}$$

$$w[i]' = \begin{cases} w[i] & \text{if } w[i] \neq \perp \\ \text{null}(q_i) & \text{else} \end{cases}$$

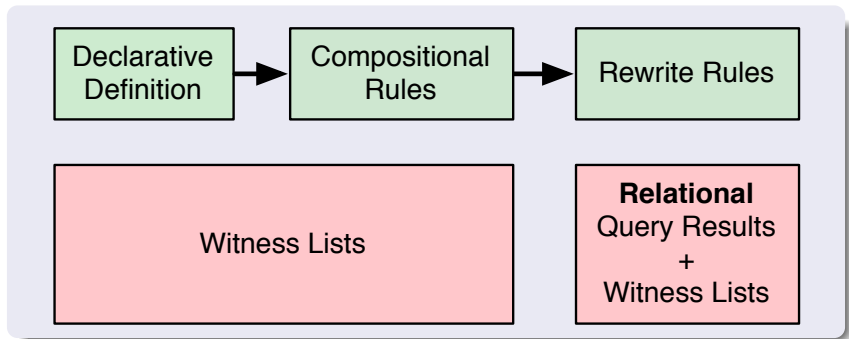
Discussion

- Comprehensible by human?
 - Yes, complete tuples
 - Better than alternatives, e.g., sets of tuple ids
- Query-able?
 - Standard relation \Rightarrow SQL
 - Same language for to query data + provenance!
- Relationship data + provenance
 - Explicitly modelled
- Disadvantages
 - Does not exploit overlap in provenance

Outline

- 1** Witness List based Provenance Models
 - Introduction
 - Provenance Model
 - Compositional Rules
 - Provenance Representation
 - Provenance Generation through Query Rewrite
 - Implementation
 - Recap

Approach Roadmap - Compositional Rules



Perm Provenance Generation

On-demand

- **Observation:** provenance not needed for every query

Perm Provenance Generation

On-demand

- **Observation:** provenance not needed for every query
- \Rightarrow Only generate when requested by the user

Query Rewrite

Approach

- Rewrite query $Q \rightarrow$ query Q^+
 - Q^+ computes provenance + original results of Q
 - by adding provenance to the inputs (duplicate attributes)
 - propagates provenance through the operations of the query
- Rewrite defined as rules
 - One rule for each “type of operation”
 - Query rewritten by recursive application to each operation

Example Query Rewrite

Example Query

```
SELECT sum(revenue) AS total, shop
FROM sales
GROUP BY shop;
```



Example (DB + Query Result)

Sales

shop	month	revenue
Migros	Jan	100
Migros	Feb	10
Migros	Mar	10
Coop	Jan	25
Coop	Feb	25

Result

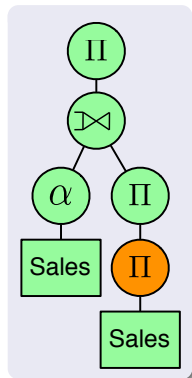
total	shop
120	Migros
50	Coop

Example Query Rewrite

Rewrite: 3 Provenance Attrs: 1

Example Query

```
SELECT total, shop, P(Q+)
FROM
  (SELECT sum(revenue) AS total, shop
   FROM sales GROUP BY shop) AS orig
LEFT OUTER JOIN
  (SELECT shop AS shop', P(orig+)
   FROM
     (SELECT shop, month, revenue,
             shop AS P(shop) ,
             month AS P(month) ,
             revenue AS P(revenue)
      FROM sales) AS sales+
  ) AS orig+
ON (shop = shop');
```



Example Query Rewrite

Rewrite: 3 Provenance Attrs: 2

Example Query

```

SELECT total, shop, P(Q+)
FROM
  (SELECT sum(revenue) AS total, shop
   FROM sales GROUP BY shop) AS orig
LEFT OUTER JOIN
  (SELECT shop AS shop', P(shop), P(month), P(revenue)
   FROM
     (SELECT shop, month, revenue,
            shop AS P(shop), month AS P(month),
            revenue AS P(revenue)
      FROM sales) AS sales+
  ) AS orig+
ON (shop = shop');

```


Example Rewritten Query Result

Example (Result + Provenance)

total	shop	P(shop)	P(month)	P(revenue)
120	Migros	Migros	Jan	100
120	Migros	Migros	Feb	10
120	Migros	Migros	Mar	10
50	Coop	Coop	Jan	25
50	Coop	Coop	Feb	25

Rewrite Rules Definitions

Definition (Rewrite Rules)

Structural Rewrite

$$q = R : \quad q^+ = \pi_{\mathbf{R}, \mathbf{R} \rightarrow \mathcal{N}(\mathbf{R})}(R) \quad (\mathbf{R1})$$

$$q = \sigma_C(q_1) : \quad q^+ = \sigma_C(q_1^+) \quad (\mathbf{R2})$$

$$q = \pi_A(q_1) : \quad q^+ = \pi_{A, \mathcal{P}(q^+)}(q_1^+) \quad (\mathbf{R3})$$

Provenance Attribute List Rewrite

$$\mathcal{P}(q^+) = \begin{cases} \mathcal{P}(q_1^+) & \text{if } q = \sigma_C(q_1) \mid \pi_A(q_1) \mid \alpha_{G,agg}(q_1) \mid \delta(q_1) \\ \mathcal{N}(R) & \text{if } q = R \\ \mathcal{P}(q_1^+) \blacktriangleright \mathcal{P}(q_2^+) & \text{else} \end{cases}$$

Rewrite Rules Definitions

Definition (Rewrite Rules)

Structural Rewrite

$$q = \alpha_{G,agg}(q_1) : \quad q^+ = \pi_{G,agg,\mathcal{P}(q^+)}(\alpha_{G,agg}(q_1) \bowtie_{G=nX} \pi_{G \rightarrow X,\mathcal{P}(q_1^+)}(q_1^+)) \quad (\text{R4})$$

$$q = \delta(q_1) : \quad q^+ = q_1^+ \quad (\text{R5})$$

Provenance Attribute List Rewrite

$$\mathcal{P}(q^+) = \begin{cases} \mathcal{P}(q_1^+) & \text{if } q = \sigma_C(q_1) \mid \pi_A(q_1) \mid \alpha_{G,agg}(q_1) \mid \delta(q_1) \\ \mathcal{N}(R) & \text{if } q = R \\ \mathcal{P}(q_1^+) \blacktriangleright \mathcal{P}(q_2^+) & \text{else} \end{cases}$$

Rewrite Rules Definitions

Definition (Rewrite Rules)

Structural Rewrite

$$q = q_1 \bowtie_C q_2 : \quad q^+ = \pi_{\mathbf{Q}_1, \mathbf{Q}_2, \mathcal{P}(q^+)}(q_1^+ \bowtie_C q_2^+) \quad (\mathbf{R6})$$

$$q = q_1 \Join_C q_2 : \quad q^+ = \pi_{\mathbf{Q}_1, \mathbf{Q}_2, \mathcal{P}(q^+)}(q_1^+ \Join_C q_2^+) \quad (\mathbf{R7})$$

Provenance Attribute List Rewrite

$$\mathcal{P}(q^+) = \begin{cases} \mathcal{P}(q_1^+) & \text{if } q = \sigma_C(q_1) \mid \pi_A(q_1) \mid \alpha_{G,agg}(q_1) \mid \delta(q_1) \\ \mathcal{N}(R) & \text{if } q = R \\ \mathcal{P}(q_1^+) \blacktriangleright \mathcal{P}(q_2^+) & \text{else} \end{cases}$$

Rewrite Rules Definitions

Definition (Rewrite Rules)

Structural Rewrite

$$q = q_1 \cup q_2 : \quad q^+ = (q_1^+ \times \text{null}(\mathcal{P}(q_2^+))) \\ \cup (\pi_{\mathbf{Q}_1, \mathcal{P}(q^+)}(q_2^+ \times \text{null}(\mathcal{P}(q_1^+)))) \quad (\mathbf{R8})$$

Provenance Attribute List Rewrite

$$\mathcal{P}(q^+) = \begin{cases} \mathcal{P}(q_1^+) & \text{if } q = \sigma_C(q_1) \mid \pi_A(q_1) \mid \alpha_{G,agg}(q_1) \mid \delta(q_1) \\ \mathcal{N}(R) & \text{if } q = R \\ \mathcal{P}(q_1^+) \blacktriangleright \mathcal{P}(q_2^+) & \text{else} \end{cases}$$

Rewrite Rules Definitions

Definition (Rewrite Rules)

Structural Rewrite

$$q = q_1 - q_2 : \quad q^+ = \pi_{\mathbf{Q}_1, \mathcal{P}(q^+)}(\delta(q_1 - q_2) \\ \bowtie_{\mathbf{Q}_1 =_n \mathbf{X}} \pi_{\mathbf{Q}_1 \rightarrow \mathbf{X}, \mathcal{P}(q_1^+)}(q_1^+) \quad (\mathbf{R10}) \\ \times \text{null}(\mathcal{P}(q_2^+)))$$

Provenance Attribute List Rewrite

$$\mathcal{P}(q^+) = \begin{cases} \mathcal{P}(q_1^+) & \text{if } q = \sigma_C(q_1) \mid \pi_A(q_1) \mid \alpha_{G,agg}(q_1) \mid \delta(q_1) \\ \mathcal{N}(R) & \text{if } q = R \\ \mathcal{P}(q_1^+) \blacktriangleright \mathcal{P}(q_2^+) & \text{else} \end{cases}$$

Outline

- 1 Witness List based Provenance Models
 - Introduction
 - Provenance Model
 - Compositional Rules
 - Provenance Representation
 - Provenance Generation through Query Rewrite
 - Implementation
 - Recap

Implementation

Perm

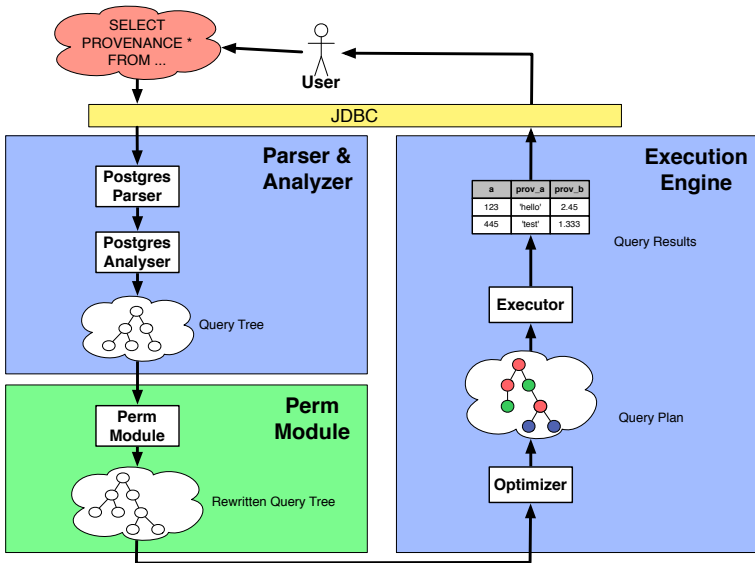
- modified *PostgreSQL* server
- *SQL-PLE*: language extension for provenance

Facts

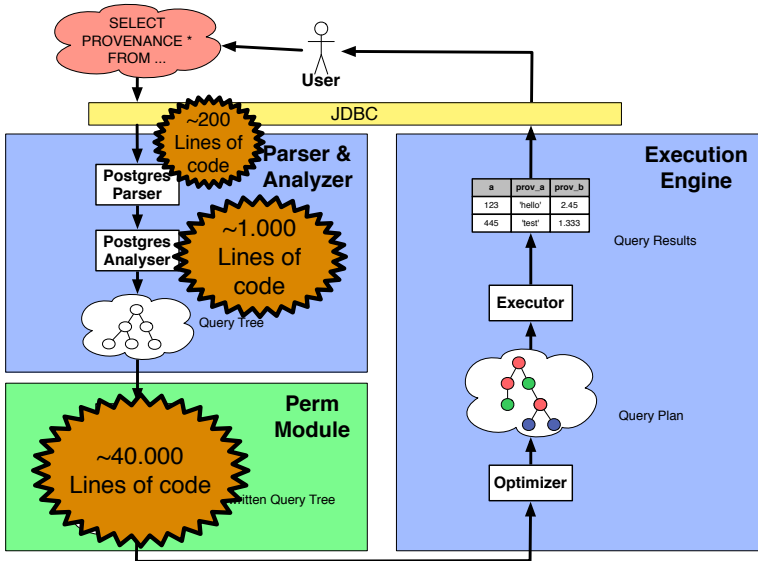
- Postgres Client Interfaces
 - JDBC
 - psql (command line)
 - ...
- Perm Module (Provenance)
 - implements query rewrites
- Open-source (<http://permdbms.sourceforge.net/>)

Implementation

Under the Hood



Under the Hood



Outline

- 1** Witness List based Provenance Models
 - Introduction
 - Provenance Model
 - Compositional Rules
 - Provenance Representation
 - Provenance Generation through Query Rewrite
 - Implementation
 - Recap

Recap

Perm Implementation

- PostgreSQL extension
- **PROVENANCE** keyword activates rewrite
- Query $q \Rightarrow$ Query q^+
 - q^+ is standard SQL/algebra query
 - Use unmodified optimizer of system

Provenance Model Comparison

Property	Why	Lin	PI-CS
Representation	Set of Set of Tuples	List of Set of Tuples	Set/Bag of List of Tuples
Language Support	USPJ	ASPJ-Set	ASPJ-Set + Nested sub-queries
Semantics	Set	Set + Bag*	Bag
Variants	Wit, Why, IWhy	Set/Bag	Influence + Copy
Definition	Decl. - Synt. - Decl./Synt.	Decl. + Synt.	Decl. + Synt.
Design Principles	Sufficiency - No false positives	Sufficiency + No false negatives + no false positives	Sufficiency + No false negatives + No false positives
Systems	-	WHIPS	Perm
Insensitivity	Yes - No - Yes	No	No

