

CS 595 - Hot topics in database systems:

Data Provenance

I. Database Provenance

I.1 Provenance Models and Systems

Boris Glavic

September 10, 2012

Tracing Rules

- Declarative definition nice, but ...
- how to compute provenance?

Example Aggregation

Query

```
SELECT shop, sum(price) AS rev
FROM sales
GROUP BY shop
```

$$q = \alpha_{shop, \text{sum}(\text{price})}(\text{sales})$$

Example

sales			
	shop	item	price
s_1	Migros	lawnmower	100
s_2	Migros	Shovel	25
s_3	Coop	Shovel	25

result		
	shop	rev
t_1	Migros	125
t_2	Coop	25

Example Aggregation

Query

```
SELECT shop, sum(price) AS rev
FROM sales
GROUP BY shop
```

$$q = \alpha_{shop, sum(price)}(sales)$$

$$q^{-1}(t) = \sigma_{shop=t.shop}(sales)$$

Example

sales			
	shop	item	price
s ₁	Migros	lawnmower	100
s ₂	Migros	Shovel	25
s ₃	Coop	Shovel	25

result	
	shop rev
t ₁	Migros 125
t ₂	Coop 25

Queries

- 1 Which operator tracing queries can be combined into a single one?
- 2 How to reorder operators to combine operators into segments that can be traced in one step?

SPJ Query

Query

- $q = \pi_A(\sigma_C(R_1 \bowtie_{C_1} \dots \bowtie_{C_{m-1}} R_m))$
- every SPJ query can be rewritten into this form!

Tracing query(ies)

- $Split_{R_1, \dots, R_m}(\sigma_{A=t \wedge C}(R_1 \bowtie_{C_1} \dots \bowtie_{C_{m-1}} R_m))$

Split operator

- $Split_{A_1, \dots, A_m}(R) = (\pi_{A_1}(R), \dots, \pi_{A_m}(R))$

ASPJ Discussion

- Each ASPJ-segment has to be traced on its own
- Tracing query needs access to inputs of segment
- ⇒ Need to store **intermediate** results for each segment
- ⇒ Or recompute large parts of the query several times

Outline

- 1** Lineage
 - Provenance Model
 - Compositional Tracing Rules
 - WHIPS Datawarehouse Implementation
 - Applications
 - Recap

Approach

- Tracing rules nice, but . . .
- how to implement computation?
- For views in data warehouse

Outline

1 Lineage

- Provenance Model
- Compositional Tracing Rules
- WHIPS Datawarehouse Implementation
- Applications
- Recap

View update - Propagate Deletions to Base Relations

Problem

- Consider SPJ view $V(D)$ over instance D (set semantics)
- How to delete t from view?
- ΔD : instance update that causes t to disappear from view
- ΔD is **exact** if only $V(D - \Delta D) = \Delta V = V(D) - \{t\}$
- $E = \Delta V - \{t\}$: **Side-effect**

Idea

- Provenance assumed to be available
 - In contrast to view maintenance could be ok compute on the fly!
- Use provenance to help us determine which inputs to delete from input

