

CS 595 - Hot topics in database systems:  
**Data Provenance**

1. Introduction to Data Provenance

Boris Glavic

August 29, 2012

# Outline

- 1 Origin of Term
- 2 Relational Algebra Primer
- 3 What is Provenance?
- 4 Types of Provenance Information
- 5 Use Cases and Application Domains
- 6 Provenance Generation, Storage, and Querying
- 7 Recap

# Data Provenance

## Data Provenance

Information about the **creation process** and **origin** of data

# Why do we call it Provenance?

## Origin of the Term

- From art dealing

## Alternative Terms

- Lineage
- Data Pedigree

# Why do we call it Provenance?

## Origin of the Term

- From art dealing

## Alternative Terms

- Lineage for kings
- Data Pedigree

# Why do we call it Provenance?

## Origin of the Term

- From art dealing

## Alternative Terms

- Lineage for kings
- Data Pedigree for dogs

# Why do we call it Provenance?

## Origin of the Term

- From art dealing for pieces of art

## Alternative Terms

- Lineage for kings
- Data Pedigree for dogs

# Provenance in Art

## Given a piece of art

- How do we know ...
  - if it is authentic?
  - who created it?
  - if it has been altered?

## Example



Jan Van Eyck - Arnolfini  
Portrait





# Outline

- 1 Origin of Term
- 2 Relational Algebra Primer
- 3 What is Provenance?
- 4 Types of Provenance Information
- 5 Use Cases and Application Domains
- 6 Provenance Generation, Storage, and Querying
- 7 Recap

# Relational Algebra

- Formalizes queries over relational data
- Is an **algebra** over relations
  - Types of operators
  - An operator produces a **single output relation** from **one or more input relations**.
- **Relations**
  - A relation is a set of tuples with the same schema
  - Tuple is a list of values
- **Composable**
  - Output of an operator can be used as input to another operator!
  - $\Rightarrow$  Can build complex queries by combining simple operators

## Types of Operators

- Selection  $\sigma_C$
- Projection  $\pi_A$
- Joins
  - Theta-join  $\bowtie_C$
  - Cross-product  $\times$
  - Outer joins  $\Join, \Join, \Join$
- Aggregation with group-by  $\alpha_{agg,G}$
- Set Operations
  - Union  $\cup$
  - Intersection  $\cap$
  - Set difference  $-$
- Relation access  $R$

# Selection

## Signature

- $\sigma_C(R)$
- $R$  is the input relation
- $C$  is a logical condition (**selection condition**)
  - Logical operators: AND ( $\wedge$ ), OR ( $\vee$ ), and NOT ( $\neg$ )
  - Comparison operators:
    - E.g., equality ( $=$ ) or smaller equals ( $\leq$ )
    - Refer to constants and attributes
  - Calls to functions?
  - E.g.,  $name = 'Peter' \wedge salary \leq 1000$

# Selection

## Definition

$$[[\sigma_C(R)]] = \{t \mid t \in R \wedge t \models C\}$$

# Selection

## Example

$$Employee = \{(Peter, 100), (Heinz, 4000)\}$$

$$[[\sigma_{name='Peter'}(Employee)]] = \{(Peter, 100)\}$$

# Projection

## Signature

- $\pi_A(R)$
- $R$  is the input relation
- $A$  is a list of **projection expressions**
  - Attributes form  $R$
  - Functions calls and operators
    - E.g.,  $a + b$
  - Renaming,  $a \rightarrow b$



# Projection

## Definition

$$[[\pi_A(R)]] = \{t \mid \exists u \in R \wedge u.A = t\}$$

# Projection

## Example

$$Employee = \{(Peter, 100), (Heinz, 4000)\}$$

$$[[\pi_{salary}(Employee)]] = \{(100), (4000)\}$$

# Join

## Signature

- $R \bowtie_C S$
- $R, S$  are the input relations
- $C$  is a logical condition (**join condition**)
  - Same as selection condition
  - Only equality conditions and  $\wedge \Rightarrow$  **Equi-join**

# Join

## Definition

$$[[R \bowtie_C S]] = \{t \blacktriangleright t' \mid t \in R \wedge t' \in S \wedge t \blacktriangleright t' \models C\}$$

# Join

## Example

$$Employee = \{(Peter, 1), (Heinz, 2)\}$$

$$Department = \{(1, CS), (2, HR)\}$$

$$[[Employee \bowtie_{deptId=id} Department]] = \{(Peter, 1, 1, CS), \\ (Heinz, 2, 2, HR)\}$$

# Outer Joins: Left-outer Join

## Signature

- $R \bowtie_C S$
- $R, S$  are the input relations
- $C$  is a logical condition (**join condition**)

# Outer Joins: Left-outer Join

## Definition

$$\begin{aligned}
 [[R \bowtie_C S]] = & \{(t \blacktriangleright t') \mid t \in R \wedge t' \in S\} \\
 & \cup \{(t_1 \blacktriangleright \text{null}(S)) \mid t \in R \wedge (\nexists t' \in S : (t \blacktriangleright t') \models C)\}
 \end{aligned}$$

# Outer Joins: Left-outer Join

## Example

$$Employee = \{(Peter, 1), (Heinz, null)\}$$

$$Department = \{(1, CS), (2, HR)\}$$

$$[[Employee \bowtie_{deptId=id} Department]] = \{(Peter, 1, 1, CS), \\ (Heinz, null, null, null)\}$$



# Aggregation

## Signature

- $\alpha_{agg,G}(R)$
- $R$  is the input relation
- $agg$  list of aggregation functions
  - E.g.,  $sum(a)$  if  $a$  attribute of  $R$
- $G$  is a list of group-by expressions
  - Attributes
  - Operators and function expressions

# Aggregation

## Definition

$$[[\alpha_{G,agg}(R)]] = \{(t.G, res_1, \dots, res_m) \mid t \in R \\ \wedge \forall i \in \{1, m\} : res_i = agg_i(\pi_{b_i}(\sigma_{G=t.G}(R)))\}$$

- $b_i$  expression used as aggregation function input
  - E.g.,  $a$  for  $sum(a)$
- $res_i$  is result of computing aggregation function for a tuple

# Aggregation

## Example

$$Employee = \{(Peter, 1, 3000), (Heinz, 2, 4000), (Jule, 1, 2000)\}$$

$$[[\alpha_{sum(salary), depld}]] = \{(5000, 1), (4000, 2)\}$$

# Union

## Signature

- $R \cup S$
- $R$  and  $S$  are the input relations
- $R$  and  $S$  have to have same schema

# Union

## Definition

$$[[R \cup S]] = \{t \mid t \in R \vee t \in S\}$$

# Union

## Example

$$Employee = \{(Peter), (Heinz), (Jule)\}$$
$$Manager = \{(Peter), (Gertrud)\}$$
$$[[Employee \cup Manager]] = \{(Peter), (Heinz), (Jule), (Gertrud)\}$$

# Intersection

## Signature

- $R \cap S$
- $R$  and  $S$  are the input relations
- $R$  and  $S$  have to have same schema

# Intersection

## Definition

$$[[R \cap S]] = \{t \mid t \in R \wedge t \in S\}$$



# Intersection

## Example

$$Employee = \{(Peter), (Heinz), (Jule)\}$$
$$Manager = \{(Peter), (Gertrud)\}$$
$$[[Employee \cap Manager]] = \{(Peter)\}$$

# Set Difference

## Signature

- $R - S$
- $R$  and  $S$  are the input relations
- $R$  and  $S$  have to have same schema

# Set Difference

## Definition

$$[[R - S]] = \{t \mid t \in R \wedge t \notin S\}$$

# Set Difference

## Example

$$Employee = \{(Peter), (Heinz), (Jule)\}$$
$$Manager = \{(Peter), (Gertrud)\}$$
$$[[Employee - Manager]] = \{(Heinz), (Jule)\}$$

# Bag vs. Set semantics

- So far: Relations are sets (**Set semantics**)
  - $\Rightarrow$  A tuple appears at most one time

# Bag vs. Set semantics

- So far: Relations are sets (**Set semantics**)
  - $\Rightarrow$  A tuple appears at most one time
- This is different from SQL and database implementations
  - Tuple can appear more than once

# Bag vs. Set semantics

- So far: Relations are sets (**Set semantics**)
  - $\Rightarrow$  A tuple appears at most one time
- This is different from SQL and database implementations
  - Tuple can appear more than once
  - **In relation in DB only if no Primary key**

# Bag vs. Set semantics

- So far: Relations are sets (**Set semantics**)
  - $\Rightarrow$  A tuple appears at most one time
- This is different from SQL and database implementations
  - Tuple can appear more than once
  - **In relation in DB only if no Primary key**
  - This is called **Bag semantics**



# Bag vs. Set semantics

- So far: Relations are sets (**Set semantics**)
  - $\Rightarrow$  A tuple appears at most one time
- This is different from SQL and database implementations
  - Tuple can appear more than once
  - **In relation in DB only if no Primary key**
  - This is called **Bag semantics**
- Bag semantics
  - Formally: assign a multiplicity  $\geq 1$  to each tuple in a relation

# Bag vs. Set semantics cont.

- Why set semantics?
  - Cleaner formalism

# Bag vs. Set semantics cont.

- Why set semantics?
  - Cleaner formalism
- Why bag semantics?

# Bag vs. Set semantics cont.

- Why set semantics?
  - Cleaner formalism
- Why bag semantics?
  - Correctness

# Bag vs. Set semantics cont.

- Why set semantics?
  - Cleaner formalism
- Why bag semantics?
  - Correctness
    - (e.g., projecting on non-unique attribute, then aggregate)

# Bag vs. Set semantics cont.

- Why set semantics?
  - Cleaner formalism
- Why bag semantics?
  - Correctness
    - (e.g., projecting on non-unique attribute, then aggregate)
  - Performance

# Bag vs. Set semantics cont.

- Why set semantics?
  - Cleaner formalism
- Why bag semantics?
  - Correctness
    - (e.g., projecting on non-unique attribute, then aggregate)
  - Performance
    - Some operators require costly duplicate removal under set semantics

# Bag semantics: Notation

## How to write multiplicities

- Use power notation to express the multiplicity of a tuple
  - $t^n \in R$  denotes tuple  $t$  exists with multiplicity  $n$  in relation  $R$



# Bag semantics: Operators

## Duplicate Removal

- $\delta(R)$
- Returns a copy of  $R$  with all multiplicities set to one

# Bag semantics: Other operators

## Definitions

$$[[\pi_A(R)]] = \{t^n \mid n = \sum_{u^m \in R \wedge u.A=t} m\}$$

$$[[\sigma_C(R)]] = \{t^n \mid t^n \in R \wedge t \models C\}$$

$$[[\alpha_{G,agg}(R)]] = \{(t.G, res_1, \dots, res_m)^1 \mid t^n \in R \\ \wedge \forall i \in \{1, m\} : res_i = agg_i(\pi_{B_i}(\sigma_{G=t.G}(R)))\}$$

$$[[R \bowtie_C S]] = \{(t_1 \blacktriangleright t_2)^{n \times m} \mid t_1^n \in R \wedge t_2^m \in S \\ \wedge (t_1 \blacktriangleright t_2) \models C\}$$

$$[[R \Join_C S]] = \{(t_1 \blacktriangleright t_2)^{n \times m} \mid t_1^n \in R \wedge t_2^m \in S\} \\ \cup \{(t_1 \blacktriangleright null(S))^n \mid t_1^n \in R \\ \wedge (\nexists t_2 \in S : (t_1 \blacktriangleright t_2) \models C)\}$$

# Bag semantics: Other operators

## Definitions

$$[[R \cup S]] = \{t^{n+m} \mid t^n \in R \wedge t^m \in S\}$$

$$[[R \cap S]] = \{t^{\min(n,m)} \mid t^n \in R \wedge t^m \in S\}$$

$$[[R - S]] = \{t^{n-m} \mid t^n \in R \wedge t^m \in S\}$$

# Outline

- 1 Origin of Term
- 2 Relational Algebra Primer
- 3 What is Provenance?
  - Provenance in Data Processing
  - An Abstract View on Provenance
  - Running Example
- 4 Types of Provenance Information
- 5 Use Cases and Application Domains
- 6 Provenance Generation, Storage, and Querying

# Data Provenance

## Data Provenance

Information about the **creation process** and **origin** of data

# Provenance in Data Processing

## Given a piece of data

- How do we know ...
  - which data it is derived from?
  - which transformations (SQL) where used to create it?
  - who created it?
  - ...

## Example

result	
	shop      rev
t <sub>1</sub>	Migros    125
t <sub>2</sub>	Coop      25

# Provenance in Data Processing

## Given a piece of data

- How do we know ...
  - which data it is derived from?
  - which transformations (SQL) where used to create it?
  - who created it?
  - ...

## Example

Compute the **revenue** for each **shop** as **sum** of **prices** of **items** sold

## Example

result	
	shop rev
t <sub>1</sub>	Migros 125
t <sub>2</sub>	Coop 25



```
SELECT shop,
       sum(price) AS rev
FROM sales, items
WHERE itemId = id
GROUP BY shop
```



sales

	shop	itemId
s <sub>1</sub>	Migros	1
s <sub>2</sub>	Migros	3
s <sub>3</sub>	Coop	3



items

	id	price
i <sub>1</sub>	1	100
i <sub>2</sub>	2	10
i <sub>3</sub>	3	25

# Provenance in Data Processing

## Given a piece of data

- How do we know ...
  - which data it is derived from?
  - which transformations (SQL) where used to create it?
  - who created it?
  - ...

## Definition (Data Provenance)

Information about the **origin** and **creation process** of data.

## Example

result

	shop	rev
t <sub>1</sub>	Migros	125
t <sub>2</sub>	Coop	25



```
SELECT shop,
       sum(price) AS rev
FROM sales, items
WHERE itemId = id
GROUP BY shop
```



sales

	shop	itemId
s <sub>1</sub>	Migros	1
s <sub>2</sub>	Migros	3
s <sub>3</sub>	Coop	3



items

	id	price
i <sub>1</sub>	1	100
i <sub>2</sub>	2	10
i <sub>3</sub>	3	25



# Abstract View

## Data

- Structured? Schemata?
- Atomic units? (Data items)

# Abstract View

## Data

- Structured? Schemata?
- Atomic units? (Data items)

## Transformations

- Consume input data
- Produce output data
- Hierarchical composition?
- Fixed set of atomic operations?

# Abstract View

## Data

- Structured? Schemata?
- Atomic units? (Data items)

## Transformations

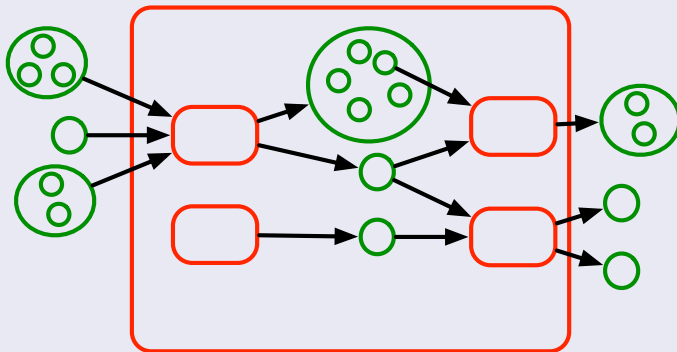
- Consume input data
- Produce output data
- Hierarchical composition?
- Fixed set of atomic operations?

## Provenance

Information about the **creation process** and **origin** of data

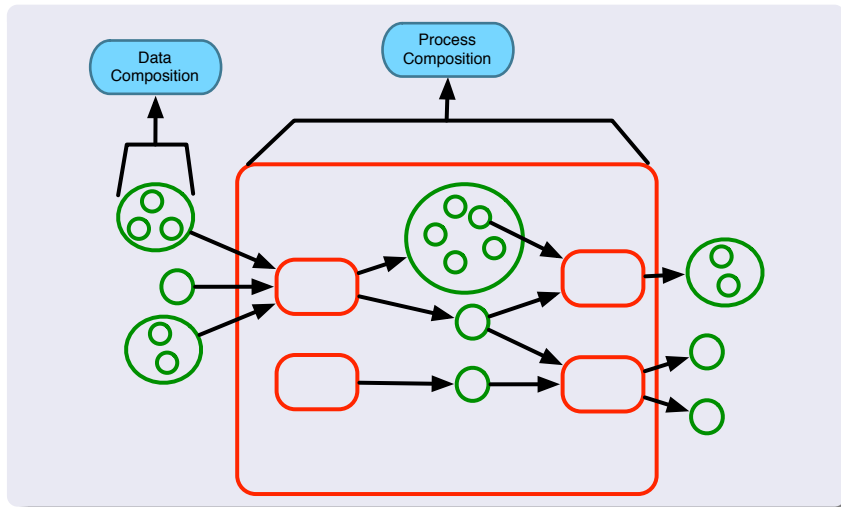
## An Abstract View on Provenance

## Abstract View



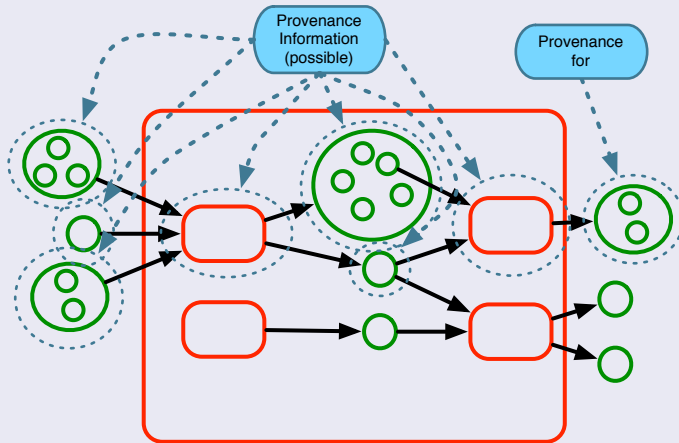
## An Abstract View on Provenance

## Abstract View



## An Abstract View on Provenance

## Abstract View



# Scenario

- You are an analyst for a garden supply shop
- You have to compute the first quarter revenue for each shop location
- Datawarehouse with sales data
- Use SQL to compute the required information from the warehouse

# Running Example

## Example (Input Data)

### Employee

SSN	Name	WorksFor
123	Peter Peterson	New York
342	Jane Janeson	New York
555	Heinz Heinzmann	Wuppertal

### Shop

Location	Budget
New York	1.000.000
Wuppertal	4.000

### Item

Id	Description	Price
1	Lawnmower	199
2	Fertilizer	32
3	Rake	9

### Sales

Employee	Item	Amount	Month
123	1	1	1
342	2	64	1
342	3	2	3
555	3	1	5



# Running Example

## Example (SalesTotal Query)

```
CREATE VIEW SalesTotal AS
SELECT Location AS Shop, Month, SSN AS Employee,
       Price * Amount AS Totalprice
FROM Employee E, Shop H, Item I, Sales S
WHERE E.WorksFor = H.Location
      AND E.SSN = S.Employee
      AND I.Id = S.Item
```

## Example (Results)

**SalesTotal**

Shop	Month	Employee	Totalprice
New York	1	123	199
New York	1	342	2048
New York	3	342	18
Wuppertal	5	555	9

# Running Example

## Example (MonthlyRevenue Query)

```
CREATE VIEW MonthlyRevenue
SELECT Shop, Month, sum(Totalprice) AS Revenue
FROM SalesTotal
GROUP BY Shop, Month
```

## Example (Results)

MonthlyRevenue

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

# Running Example

## Example (RevenueFirstQ Query)

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

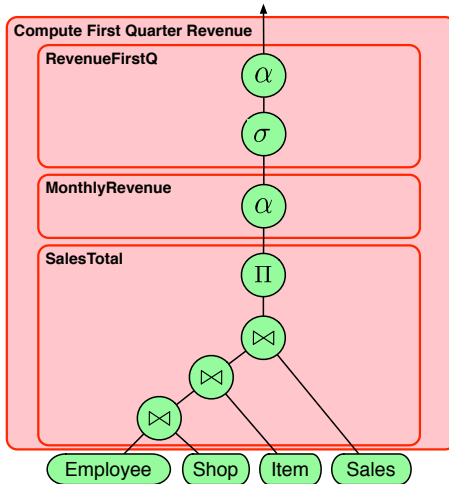
## Example (Results)

RevenueFirstQ

Shop	Revenue
New York	2265

## Running Example

# Running Example



# Tracing an Error

## Problem

- One result tuple of your query looks suspicious
- You expect the input data to be the culprit
- How to know which input data affected which output data

# Tracing an Error

## Problem

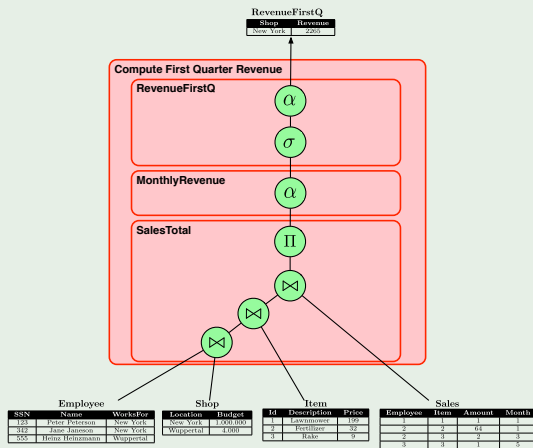
- One result tuple of your query looks suspicious
- You expect the input data to be the culprit
- How to know which input data affected which output data

## This is Data Provenance

## Running Example

## Example Data

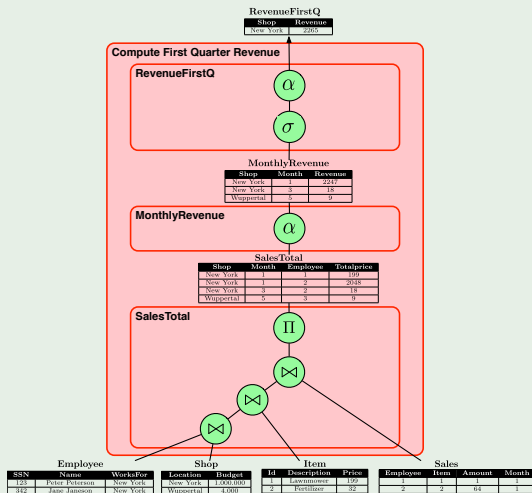
## Example



## Running Example

## Example Data

## Example





## Running Example

## Example Data

## Example

SalesTotal			
Shop	Month	Employee	Totalprice
New York	1	1	100
New York	1	2	2018
New York	3	2	18
Wuppertal	3	3	9

SalesTotal

EmployeeShop				
SSN	Name	WorksFor	Location	Budget
123	Peter Peterson	New York	New York	1,000,000
342	Jane Jansson	New York	New York	1,000,000
555	Heinz Heinemann	Wuppertal	Wuppertal	4,000

Employee		
SSN	Name	WorksFor
123	Peter Peterson	New York
342	Jane Jansson	New York
555	Heinz Heinemann	Wuppertal

Shop		
Location	Budget	
New York	1,000,000	
Wuppertal	4,000	

Item		
Id	Description	Price
1	Lawnmower	470
2	Fertilizer	32
3	Rake	9

Sales			
Employee	Item	Amount	Month
2	1	64	1
2	3	2	3
3	3	1	3

## Lessons learned

## Lessons learned

- Which inputs belong to provenance of outputs?

## Lessons learned

- Which inputs belong to provenance of outputs?
  - hard

## Lessons learned

- Which inputs belong to provenance of outputs?
  - hard
- Even if we know: How to get it?

## Lessons learned

- Which inputs belong to provenance of outputs?
  - hard
- Even if we know: How to get it?
- Manually?

## Lessons learned

- Which inputs belong to provenance of outputs?
  - **hard**
- Even if we know: How to get it?
- Manually?
  - Not reasonable for big data or complex query!

## Lessons learned

- Which inputs belong to provenance of outputs?
  - **hard**
- Even if we know: How to get it?
- Manually?
  - Not reasonable for big data or complex query!
- Need system that tracks it automatically!



# Outline

- 1 Origin of Term
- 2 Relational Algebra Primer
- 3 What is Provenance?
- 4 Types of Provenance Information
  - Data Provenance
  - Transformation Provenance
  - Other
- 5 Use Cases and Application Domains
- 6 Provenance Generation, Storage, and Querying

# Types of Provenance Information

## Provenance Types

- Data Provenance
- Transformation Provenance
- Additional Information

# Types of Provenance Information

## Provenance Types

- Data Provenance
  - From which **input data** is which **output data** derived from
- Transformation Provenance
- Additional Information

# Types of Provenance Information

## Provenance Types

- Data Provenance
  - From which **input data** is which **output data** derived from
- Transformation Provenance
  - Which **transformations** contributed in which way to which **output data**
- Additional Information

# Types of Provenance Information

## Provenance Types

- Data Provenance
  - From which **input data** is which **output data** derived from
- Transformation Provenance
  - Which **transformations** contributed in which way to which **output data**
- Additional Information
  - Execution environment (state of the world)
  - Involved Users

# Data Provenance

Which/How input data influences output data

# Data Provenance

## Which/How input data influences output data

- Data Granularity
  - Attribute value
  - Tuple
  - Relation

# Data Provenance

## Which/How input data influences output data

- Data Granularity
  - Attribute value
  - Tuple
  - Relation
- Transformation Granularity
  - Query with view unfolding
  - Query block
  - Algebra operator



# Data Provenance

## Which/How input data influences output data

- Data Granularity
  - Attribute value
  - Tuple
  - Relation
- Transformation Granularity
  - Query with view unfolding
  - Query block
  - Algebra operator
- “True” Data Dependencies?
  - Black-box: An output depends on all inputs
  - Fine-grained: Dependencies depending on how data is processed by transformation

# Data Granularity

## Example (Relation)

### RevenueFirstQ

Shop	Revenue
New York	2265

#### RevenueFirstQ

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

### MonthlyRevenue

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

# Data Granularity

## Example (Tuple)

### RevenueFirstQ

Shop	Revenue
New York	2265

#### RevenueFirstQ

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

### MonthlyRevenue

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

# Data Granularity

## Example (Attribute Value)

### RevenueFirstQ

Shop	Revenue
New York	2265

#### RevenueFirstQ

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

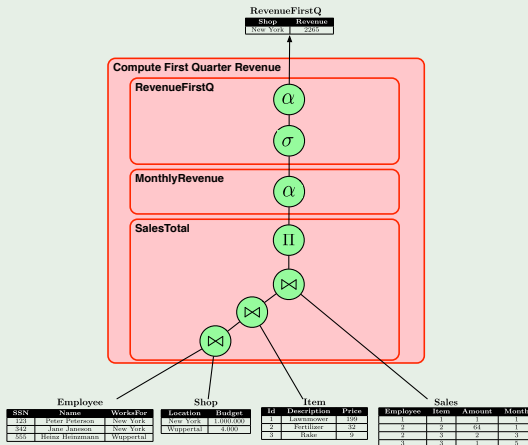
### MonthlyRevenue

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

## Data Provenance

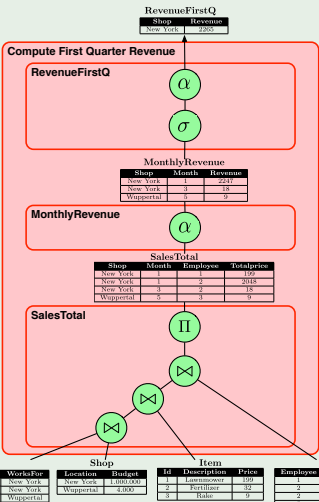
## Transformation Granularity

## Example



# Transformation Granularity

## Example

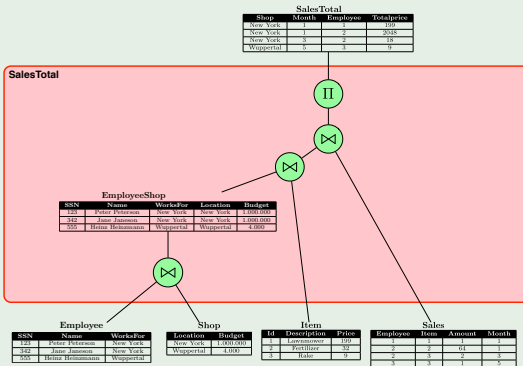


OGY

## Data Provenance

## Transformation Granularity

## Example



# Data Dependencies

## Example (Black-box)

### RevenueFirstQ

Shop	Revenue
New York	2265

#### RevenueFirstQ

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

### MonthlyRevenue

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



# Data Dependencies

## Example (Fine-grained)

### RevenueFirstQ

Shop	Revenue
New York	2265

#### RevenueFirstQ

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

### MonthlyRevenue

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

# Transformation Provenance

Which/How transformations contributed to output data

- Transformations that generated output (transitive?)
- Only the ones that had actual effect
- Workflow template/program vs. workflow run/execution

# Additional Information

## A small subset

- OS version

# Additional Information

## A small subset

- OS version
- Version of library linked against

# Additional Information

## A small subset

- OS version
- Version of library linked against
- Environment variables

# Additional Information

## A small subset

- OS version
- Version of library linked against
- Environment variables
- User that executed the process

# Additional Information

## A small subset

- OS version
- Version of library linked against
- Environment variables
- User that executed the process
- Current main memory content

# Additional Information

## A small subset

- OS version
- Version of library linked against
- Environment variables
- User that executed the process
- Current main memory content
- Room temperature?



# Additional Information

## A small subset

- OS version
- Version of library linked against
- Environment variables
- User that executed the process
- Current main memory content
- Room temperature?
- Geographical location

# Additional Information

## A small subset

- OS version
- Version of library linked against
- Environment variables
- User that executed the process
- Current main memory content
- Room temperature?
- Geographical location
- . . .

# Additional Information

## A small subset

- OS version
- Version of library linked against
- Environment variables
- User that executed the process
- Current main memory content
- Room temperature?
- Geographical location
- ...
- Butterfly that flapped in china

# Outline

- 1 Origin of Term
- 2 Relational Algebra Primer
- 3 What is Provenance?
- 4 Types of Provenance Information
- 5 Use Cases and Application Domains**
  - Use Cases
  - Debugging
  - Annotation Propagation
  - Deletion Propagation

# Use Cases

- Debugging (tracking the sources of errors)
- Propagating annotations
- Gain deeper understanding of data and transformations
  - Estimate quality, trust
- Improvement of other data processing technologies
  - Probabilistic databases
  - Deletion propagation
  - Testing

# Application Domains

- Complex database queries, e.g., datawarehousing
- E-science and curated databases
- Data integration/exchange
- Workflow systems

# Application Domains

- Complex database queries, e.g., datawarehousing
- E-science and curated databases
- Data integration/exchange
- Workflow systems
- $\Rightarrow$  Application domain with complex, multi-stage data processing
  - Map-Reduce style processing and its “frontends” like Pig
  - Simulations
  - ...

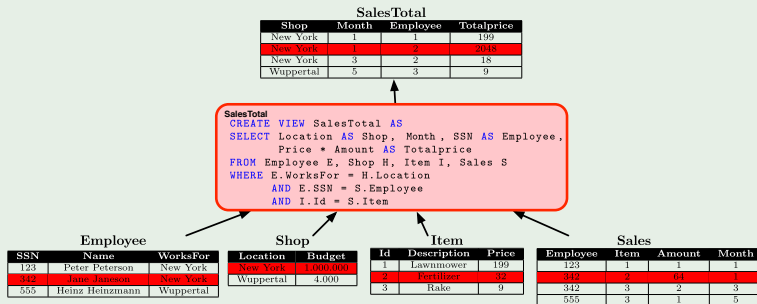
## Debugging

# Debugging

## Origin of Result Tuples

- Tuple in result suspicious/wrong/interesting

## Example



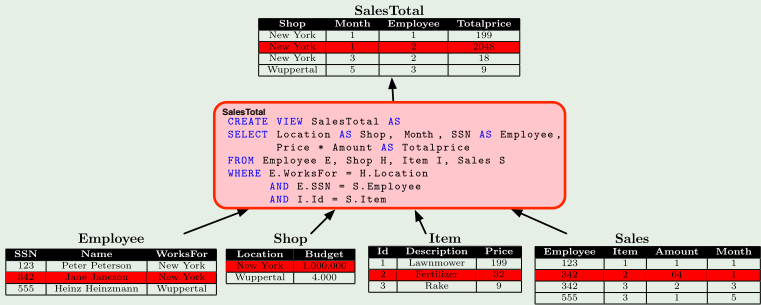


# Debugging

## Origin of Result Tuples

- Tuple in result suspicious/wrong/interesting
- Learn more by looking at relevant inputs (provenance)

## Example



# Debugging

## Approach

- 1 Identify tuples of interest ( $I$ )

# Debugging

## Approach

- 1 Identify tuples of interest ( $I$ )
- 2 Retrieve provenance
  - Need system that returns provenance for set  $I$

# Debugging

## Approach

- 1 Identify tuples of interest ( $I$ )
- 2 Retrieve provenance
  - Need system that returns provenance for set  $I$
  - How to represent this info?

# Debugging

## Approach

- ① Identify tuples of interest ( $I$ )
- ② Retrieve provenance
  - Need system that returns provenance for set  $I$
  - How to represent this info?
- ③ What if provenance large?
  - $\Rightarrow$  Query support? Visualization?

# Annotation Propagation

## Example

### EnzymeProduce

Enzyme	Gene
EC 1.1.1.1	ALB
EC 1.97.1.6	ALB

??  
??

```
CREATE VIEW EnzymeProduce AS
SELECT Enzyme, Name AS Gene
FROM Gene G, Enzyme E
WHERE G.Id = E.ProducedBy
```

### Gene

Id	Name
4q11-q13	ALB
18q21.3	BCL2

{a<sub>4</sub>}  
{}

### Enzyme

Enzyme	Weight	ProducedBy
EC 1.1.1.1	45	4q11-q13
EC 1.97.1.6	12	4q11-q13

{a<sub>1</sub>, a<sub>2</sub>}  
{a<sub>2</sub>, a<sub>3</sub>}

a<sub>1</sub> Necessary for red blood cells  
a<sub>3</sub> Unhealthy

a<sub>2</sub> Produced in liver  
a<sub>4</sub> Discovered by Edmond Hillary

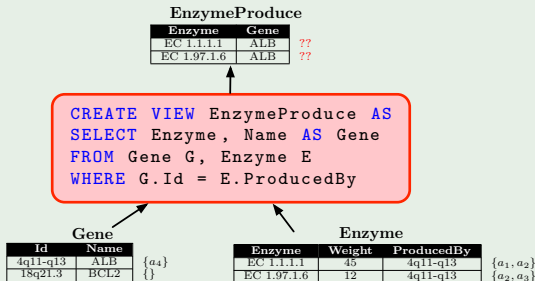
## Annotation Propagation

# Annotation Propagation

## Which annotations in query result?

- Find provenance for tuple
- Attach union of annotations in provenance

## Example



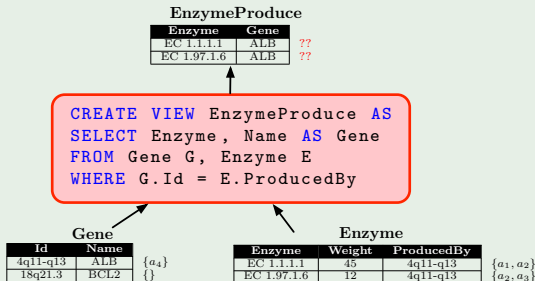
## Annotation Propagation

# Annotation Propagation

## For Example

- First result tuple
- **Provenance:** first tuples from Gene and Enzyme
- **Annotation:**  $a_1, a_2, a_4$

## Example





# Annotation Propagation - Caveats

## Potential Problems?

- What about negative influence?
- User should have control on propagation?
- What about annotations on
  - Attribute values
  - Spanning several tuples/relations/attributes

## Deletion Propagation

# Deletion Propagation

### Problem

- Given a *materialized view*

## Deletion Propagation

# Deletion Propagation

### Problem

- Given a *materialized view*
  - Query result stored as a table

## Deletion Propagation

# Deletion Propagation

### Problem

- Given a *materialized view*
  - Query result stored as a table
- How to update the view when input data changes

## Deletion Propagation

# Deletion Propagation

### Problem

- Given a *materialized view*
  - Query result stored as a table
- How to update the view when input data changes
  - Without recomputing the whole query

## Deletion Propagation

# Deletion Propagation

### Problem

- Given a *materialized view*
  - Query result stored as a table
- How to update the view when input data changes
  - Without recomputing the whole query
- **Deletion Propagation:** Update the view when input tuples are deleted?

# Deletion Propagation Example

## Example

```
CREATE VIEW ActiveCS AS
SELECT DISTINCT E.Name AS Emp
FROM Employee E, Project P, Assigned A
WHERE E.Id = A.Emp AND P.Name = A.Project
      AND Dep = CS
```

Employee	
	Id    Name
e <sub>1</sub>	1    Peter
e <sub>2</sub>	2    Gertrud
e <sub>2</sub>	3    Michael

Project	
	Name    Dep
p <sub>1</sub>	Server    CS
p <sub>2</sub>	Webpage    CS
p <sub>3</sub>	Fire CS    HR

Assigned	
	Project    Emp
a <sub>1</sub>	Server    1
a <sub>2</sub>	Server    2
a <sub>3</sub>	Webpage    2
a <sub>4</sub>	Fire CS    3

## Deletion Propagation

# Deletion Propagation Example

### Example

#### ActiveCS

	Emp
$t_1$	Peter
$t_2$	Gertrud

#### Employee

	Id	Name
$e_1$	1	Peter
$e_2$	2	Gertrud
$e_2$	3	Michael

#### Project

	Name	Dep
$p_1$	Server	CS
$p_2$	Webpage	CS
$p_3$	Fire CS	HR

#### Assigned

	Project	Emp
$a_1$	Server	1
$a_2$	Server	2
$a_3$	Webpage	2
$a_4$	Fire CS	3



## Deletion Propagation

## Deletion Propagation Example

## Example

- Delete tuple from Projects

## ActiveCS

Emp	
$t_1$	Peter
$t_2$	Gertrud

## Employee

	Id	Name
$e_1$	1	Peter
$e_2$	2	Gertrud
$e_2$	3	Michael

## Project

	Name	Dep
$p_1$	Server	CS
$p_2$	Webpage	CS
$p_3$	Fire CS	HR

## Assigned

	Project	Emp
$a_1$	Server	1
$a_2$	Server	2
$a_3$	Webpage	2
$a_4$	Fire CS	3

## Deletion Propagation

# Deletion Propagation Example

### Example

- What would be the effect on the view?

#### ActiveCS

	Emp
$t_1$	Peter
$t_2$	Gertrud

#### Employee

	Id	Name
$e_1$	1	Peter
$e_2$	2	Gertrud
$e_2$	3	Michael

#### Project

	Name	Dep
$p_1$	Server	CS
$p_2$	Webpage	CS
$p_3$	Fire CS	HR

#### Assigned

	Project	Emp
$a_1$	Server	1
$a_2$	Server	2
$a_3$	Webpage	2
$a_4$	Fire CS	3

## Deletion Propagation

## Deletion Propagation - Approach

## Assumption

- Assume we have provenance for each tuple
  - For now a set of input tuples
  - $P(t_1) = \{e_1, p_1, a_1\}$
  - $P(t_2) = \{e_1, p_1, p_2, a_1, a_2\}$
- Set of deleted tuples ( $D = \{p_1\}$ )

## Deletion Propagation

# Deletion Propagation - Approach

## Assumption

- Assume we have provenance for each tuple
  - For now a set of input tuples
  - $P(t_1) = \{e_1, p_1, a_1\}$
  - $P(t_2) = \{e_1, p_1, p_2, a_1, a_2\}$
- Set of deleted tuples ( $D = \{p_1\}$ )

## Approach

- Remove  $D$  from provenance
- Remove tuples without justification from view
  - Set provenance model to simple
  - Will learn later how this actually works

OGY

## Deletion Propagation

# Deletion Propagation Example

### Example

- $P(t_1) = \{e_1, p_1, a_1\} \rightarrow \{e_1, a_1\}$
- $P(t_2) = \{e_1, p_1, p_2, a_1, a_2\} \rightarrow \{e_1, p_2, a_1, a_2\}$

### ActiveCS

	Emp
$t_1$	Peter
$t_2$	Gertrud

### Employee

	Id	Name
$e_1$	1	Peter
$e_2$	2	Gertrud
$e_2$	3	Michael

### Project

	Name	Dep
$p_1$	Server	CS
$p_2$	Webpage	CS
$p_3$	Fire CS	HR

### Assigned

	Project	Emp
$a_1$	Server	1
$a_2$	Server	2
$a_3$	Webpage	2
$a_4$	Fire CS	3

# Outline

- 1 Origin of Term
- 2 Relational Algebra Primer
- 3 What is Provenance?
- 4 Types of Provenance Information
- 5 Use Cases and Application Domains
- 6 Provenance Generation, Storage, and Querying**
  - Provenance Generation
  - Provenance Storage
  - Provenance Querying

# Generation

## Manual vs. Automatic

- **Manual:** User has to provide provenance information
- **Automatic:** System generates provenance information automatically
- Design space: How much information has to be provided by the user or transformation developer?

## Lazy vs. Eager

- **Eager:** Generate provenance while the transformation is running
- **Lazy:** Generate provenance later once it is requested
- Tradeoff: Retrieval time vs. execution overhead

# Generation

## Approaches

- Run transformation in supervised environment that tracks provenance
- Instrument the transformations to produce provenance
- Record some information during execution and reconstruct provenance from this information



# Supervised Environment

## Idea

- Modify execution environment of transformations to capture provenance

## Considerations

- What provenance to capture?
- Which parts of system ...
  - Are accessible?
  - Are modifiable?
- Supervision for all or only some transformations

# Supervised Environment - Example

## Example

- Hadoop - Map/Reduce
- Modify the Hadoop system to
  - store relationships between input/output keys
  - for mappers and reducers
  - in HDFS?

# Supervised Environment - Discussion

## Advantages

- Can capture whatever provenance we want
- No modification to transformations

## Disadvantages

- **Intrusive** May have to re-implement whole system
- Overhead for transformation execution
- Parts of the system may not be accessible (e.g., web-service composition)

# Instrument Transformations

## Idea

- Modify the transformation to track its own provenance

## Considerations

- Transformation language expressive enough to compute its own provenance?
- How to represent provenance in the data model?

# Instrument Transformations - Example

## Example

- SQL queries
- Rewrite queries to produce their output + provenance information
- Possible?

# Instrument Transformations - Example

## Example

- SQL queries
- Rewrite queries to produce their output + provenance information
- Possible? **yes, later in course**

# Instrument Transformations - Example

## Example

- SQL queries
- Rewrite queries to produce their output + provenance information
- Possible? **yes, later in course**
- Build a middleware that does that over standard DBMS

# Instrument Transformations - Discussion

## Advantages

- **Non-intrusive:** Possible without changes to system
  - If we can gather enough information about transformation from outside
  - E.g., DBMS client
- No overhead if no provenance computed
- Same data model  $\Rightarrow$  Querying
- No manual changes to transformations

## Disadvantages

- Performance optimizations may be limited (overhead provenance computation)
- Data model may limit the provenance representation



# Reconstruction

## Idea

- Recover provenance from input + output data and knowledge about transformation

## Considerations

- Possible to know what's going on in the black box?
- Need to store extra information

# Reconstruction - Example

## Example

- Simple SQL query
- Write program to
  - Analyse query
  - Retrieve input and output data
  - Compute provenance

# Reconstruction - Discussion

## Advantages

- **Non-intrusive:** No changes to system
- No overhead for transformation
- No storage costs **or almost no storage**

## Disadvantages

- Not possible for complex operations
- Provenance generation may be more expensive

# Eager Generation

## Approach

- Generate provenance during transformation execution

## Considerations

- Overhead for transformation?
- How to trigger?

# Lazy Generation

## Approach

- Generate provenance on request

## Considerations

- Input/Output data still available?
- Transformation info available?

# Storage

- Provenance data can be orders of magnitude larger than input/output data
- ⇒ Be clever when to store what at which level of abstraction
- ⇒ Specialized compression for provenance
- ⇒ Index structured for provenance specific retrieval patterns

# Why is provenance large?

Simplified explanation:

# Why is provenance large?

## Simplified explanation:

- Input data: size  $N$



# Why is provenance large?

## Simplified explanation:

- Input data: size  $N$
- Output data: size  $M$

# Why is provenance large?

## Simplified explanation:

- Input data: size  $N$
- Output data: size  $M$
- Provenance is relationship between inputs and outputs

# Why is provenance large?

## Simplified explanation:

- Input data: size  $N$
- Output data: size  $M$
- Provenance is relationship between inputs and outputs
- $\Rightarrow$  Worst case:  $N \times M$

# Why is provenance large?

## Simplified explanation:

- Input data: size  $N$
- Output data: size  $M$
- Provenance is relationship between inputs and outputs
- $\Rightarrow$  Worst case:  $N \times M$
- Intermediate results?

# Why is provenance large?

## Simplified explanation:

- Input data: size  $N$
- Output data: size  $M$
- Provenance is relationship between inputs and outputs
- $\Rightarrow$  Worst case:  $N \times M$
- Intermediate results?
  - Transformation is tree with  $X$  nodes

# Why is provenance large?

## Simplified explanation:

- Input data: size  $N$
- Output data: size  $M$
- Provenance is relationship between inputs and outputs
- $\Rightarrow$  Worst case:  $N \times M$
- Intermediate results?
  - Transformation is tree with  $X$  nodes
  - $\Rightarrow \sim N \times M \times X$

# What to store and when?

## What?

- Only necessary level of detail
  - E.g., need attribute level provenance?
  - E.g., need provenance for intermediate results?

## When?

- Provenance for all transformations?
- Only for specific type?
- Only when requested by user?
- Only when triggering event happend?

# Compression

## Rationale

- Provenance large, but has overlap
  - Exploit overlap to compress
  - Information loss?
  - Access/querying without decompression
  - Tradeoff: speed vs. size

## Approaches

- Generic compression algorithms
  - Small size, slow?, probably no query
- Methods exploiting overlap being aware of provenance structure
  - Size less predictable, fast?, query may be possible



# Index structures

## Rationale

- Provenance querying needs efficient access to provenance data
- Traditional index structure useful?
- Can identify new access patterns?
  - Tree-path traversal?
- Static index or updates possible?

## Approaches

- E.g., adapt IR retrieval index structures

## Provenance Storage

# Example Storage - Provenance tables

- Provenance Table
- input TID's → output TID's

## Example

result	
	shop    rev
t <sub>1</sub>	Migros    125
t <sub>2</sub>	Coop    25



```

SELECT shop,
       sum(price) AS rev
FROM sales, items
WHERE itemId = id
GROUP BY shop
    
```



sales

	shop	itemId
s <sub>1</sub>	Migros	1
s <sub>2</sub>	Migros	3
s <sub>3</sub>	Coop	3



items

	id	price
i <sub>1</sub>	1	100
i <sub>2</sub>	2	10
i <sub>3</sub>	3	25

## Provenance Storage

# Example Storage - Provenance tables

- Provenance Table
- input TID's → output TID's

### Provenance

result	in
$t_1$	$s_1$
$t_1$	$s_2$
$t_1$	$i_1$
$t_1$	$i_3$
$t_2$	$s_3$
$t_2$	$i_3$

### Example

	shop	rev
$t_1$	Migros	125
$t_2$	Coop	25



```

SELECT shop,
       sum(price) AS rev
FROM sales, items
WHERE itemId = id
GROUP BY shop
    
```



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

# Querying

- Large amount of provenance information

# Querying

- Large amount of provenance information
- Query support to extract information
  - Focus on parts of interest
    - Backward: Which data contributed to output?
    - Forward: Which data is derived from input?
    - Transitive closure
  - Correlated with input/output data
  - Summarize, abstract

# Querying

- Large amount of provenance information
- Query support to extract information
  - Focus on parts of interest
    - Backward: Which data contributed to output?
    - Forward: Which data is derived from input?
    - Transitive closure
  - Correlated with input/output data
  - Summarize, abstract

## Example

For a subset of erroneous sales totals, which ones have been derived from input sales data from a shop in New York with a amount sold bigger than 100.

# Querying

## Approaches

- Extend query language for “normal” data
- New query language

# Querying

## Approaches

- Extend query language for “normal” data
  - Querying provenance in combination with “normal” data
  - Limitation on provenance representation
- New query language



# Querying

## Approaches

- Extend query language for “normal” data
  - Querying provenance in combination with “normal” data
  - Limitation on provenance representation
- New query language
  - Operations tailored for typical operations on provenance
  - “Re-inventing the wheel”

## Provenance Querying

## Querying Example

## Example

Provenance

result	in
$t_1$	$s_1$
$t_1$	$s_2$
$t_1$	$i_1$
$t_1$	$i_3$
$t_2$	$s_3$
$t_2$	$i_3$

## Example

result		
	shop	rev
$t_1$	Migros	125
$t_2$	Coop	25



```

SELECT shop,
       sum(price) AS rev
FROM sales, items
WHERE itemId = id
GROUP BY shop

```



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

# Querying Example

```
SELECT DISTINCT shop
FROM result r,
     Provenance p,
     items i
WHERE r.tid = p.result
      AND p.in = i.tid
      AND i.price > 90
```

## Example

Provenance

result	in
t <sub>1</sub>	s <sub>1</sub>
t <sub>1</sub>	s <sub>2</sub>
t <sub>1</sub>	i <sub>1</sub>
t <sub>1</sub>	i <sub>3</sub>
t <sub>2</sub>	s <sub>3</sub>
t <sub>2</sub>	i <sub>3</sub>

## Example

result

	shop	rev
t <sub>1</sub>	Migros	125
t <sub>2</sub>	Coop	25

↑

```
SELECT shop,
       sum(price) AS rev
FROM sales, items
WHERE itemId = id
GROUP BY shop
```

↑ sales                      ↑ items

	shop	itemId
s <sub>1</sub>	Migros	1
s <sub>2</sub>	Migros	3
s <sub>3</sub>	Coop	3

	id	price
i <sub>1</sub>	1	100
i <sub>2</sub>	2	10
i <sub>3</sub>	3	25

## Provenance Querying

## Querying Example

Prov query

$q_1$	shop
	Migros

## Example

Provenance

result	in
$t_1$	$s_1$
$t_1$	$s_2$
$t_1$	$i_1$
$t_1$	$i_3$
$t_2$	$s_3$
$t_2$	$i_3$

## Example

	shop	rev
$t_1$	Migros	125
$t_2$	Coop	25



```
SELECT shop,
       sum(price) AS rev
FROM sales, items
WHERE itemId = id
GROUP BY shop
```



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

# Outline

- 1 Origin of Term
- 2 Relational Algebra Primer
- 3 What is Provenance?
- 4 Types of Provenance Information
- 5 Use Cases and Application Domains
- 6 Provenance Generation, Storage, and Querying
- 7 Recap**

## Types of Provenance

- Data
- Transformation
- Other
- **Granularities**

## Generation, Storage, and Querying

- **Generation**
  - Manual vs. automatic
  - Eager vs. lazy
  - Supervised environment, instrumentation, reconstruction
- **Storage**
  - Compression, Indices, What to keep?
- **Querying**
  - Extending transformation language
  - Develop new query language

# Surveys I



[Boris Glavic and Renée J. Miller.](#)

Reexamining Some Holy Grails of Data Provenance.

In Tapp '11: 3rd unix workshop on the theory and practice of provenance, 2011.



[James Cheney, Laura Chiticariu, and Wang-Chiew Tan.](#)

Provenance in Databases: Why, How, and Where.

Foundations and Trends in Databases, 1(4):379–474, 2009.



[Robert Ikeda and Jennifer Widom.](#)

Data Lineage: A Survey.

Technical report, Stanford University, 2009.



[Juliana Freire, David Koop, Emanuele Santos, and Claudio T. Silva.](#)

Provenance for Computational Tasks: A Survey.

Computing in Science and Engineering, 10(3):11–21, 2008.



[Boris Glavic and Klaus R. Dittrich.](#)

Data Provenance: A Categorization of Existing Approaches.

In Btw '07: proceedings of datenbanksysteme in business, technologie und web, 227–241, 2007.



[Susan B. Davidson, Sarah Cohen-Boulakia, Anat Eyal, Bertram Ludscher, Timothy McPhillips, Shawn Bowers, and Juliana Freire.](#)

Provenance in Scientific Workflow Systems.

IEEE Data Engineering Bulletin, 32(4):44–50, 2007.

# Surveys II



Wang-Chiew Tan.

Provenance in Databases: Past, Current, and Future.  
[IEEE Data Engineering Bulletin](#), 30(4):3–12, 2007.



Yogesh L. Simmhan, Beth Plale, and Dennis Gannon.

A Survey of Data Provenance in e-science.  
[SIGMOD Record](#), 34(3):31–36, 2005.



Yogesh L. Simmhan, Beth Plale, and Dennis Gannon.

A Survey of Data Provenance Techniques.  
Technical report, Indiana University, Bloomington IN 47405, 2005.



Wang-Chiew Tan.

Research Problems in Data Provenance.  
[IEEE Data Engineering Bulletin](#), 27(4):42–52, 2004.