

CS 595 - Hot topics in database systems:

## **Data Provenance**

I. Database Provenance

I.1 Provenance Models and Systems

Boris Glavic

October 31, 2012







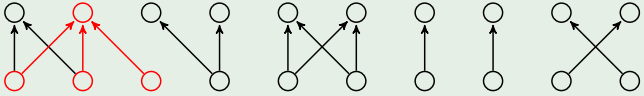




# Model Independent Storage Requirements

- Each of the  $m$  outputs
- depends on up to all of the  $n$  inputs

## Example



























# Transaction Granularity

## Example

**Transaction:** provenance is tuples in database state before transaction execution

```

BEGIN TRANSACTION
INSERT INTO persons (SELECT name, salary
                     FROM employee, pay
                     WHERE (SSN = employee));
UPDATE persons SET salary = salary + 1000
        WHERE job = 'consultant';
COMMIT

```



# Transaction Granularity

## Example

**SQL query block:** Track for each individual `SELECT-FROM-WHERE-...` block

```
BEGIN TRANSACTION
INSERT INTO persons (SELECT name, salary
                    FROM employee, pay
                    WHERE (SSN = employee));
UPDATE persons SET salary = salary + 1000
   WHERE job = 'consultant';
COMMIT
```

# Transaction Granularity

## Example

**Single operators:** E.g., join and projection in `INSERT`

```
BEGIN TRANSACTION
INSERT INTO persons (SELECT name, salary
                     FROM employee, pay
                     WHERE (SSN = employee));
UPDATE persons SET salary = salary + 1000
WHERE job = 'consultant';
COMMIT
```

# Transaction Granularity

## Example

**Individual expressions:** E.g., provenance  $a = b$  in **WHERE** condition  
 $a = b \wedge c > 5$

```
BEGIN TRANSACTION
INSERT INTO persons (SELECT name, salary
                     FROM employee, pay
                     WHERE (SSN = employee));
UPDATE persons SET salary = salary + 1000
      WHERE job = 'consultant';
COMMIT
```

# Example Transformation Granularity

## Example

```
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
```















# Datamodel Mismatch

## Rationale

- Provenance datamodel often  $\neq$  data model of transformation
  - Lineage is list of relations
  - Provenance Polynomials are formulas over tuple variables
  - Causality is set of tuples (but from different relations)
  - ...

# Datamodel Mismatch

## Rationale

- Provenance datamodel often  $\neq$  data model of transformation
  - Lineage is list of relations
  - Provenance Polynomials are formulas over tuple variables
  - Causality is set of tuples (but from different relations)
  - ...
- $\Rightarrow$  Map it to transformation data model
  - That what *Perm* does
  - That is what *DBNotes* does

# Datamodel Mismatch

## Rationale

- Provenance datamodel often  $\neq$  data model of transformation
  - Lineage is list of relations
  - Provenance Polynomials are formulas over tuple variables
  - Causality is set of tuples (but from different relations)
  - ...
- $\Rightarrow$  Map it to transformation data model
  - That what *Perm* does
  - That is what *DBNotes* does
- $\Rightarrow$  Live with different data models
  - Store provenance separately
    - $\Rightarrow$  Querying becomes a problem
  - Use datamodel expressive enough to store both data and provenance

BY











# Example Use more expressive data model

## Example

- Causality is set of tuples from different relation (causes)
- $\Rightarrow$  e.g., XML can express such sets
- map data and provenance into XML

```

<DB>
  <Person>
    <Tuple id="1"><Attr>Peter</Attr><Attr>1</Attr></Tuple>
    <Tuple id="2"><Attr>Heinz</Attr><Attr>1</Attr></Tuple>
  </Person>
  <Address>
    <Tuple id="1"><Attr>1</Attr><Attr>Toronto</Attr>
      <Attr>52 Bloor Street</Attr></Tuple>
  </Address>
  <QueryResult query="q">
    <Tuple id="1"><Attr>Peter</Attr></Tuple>
    <Tuple id="2"><Attr>Heinz</Attr></Tuple>
  </QueryResult>
  <Provenance tuple="1" query="q">
    <TupleRef relation="Person" tuple="1"/>
    <TupleRef relation="Address" tuple="1"/>
  </Provenance>
</DB>

```

# Outline

## 1 Provenance Storage

- Introduction
- Compression Methods for Provenance
- Index Structures for Provenance
- Recap

# Overview

- Given that provenance is large
- How to save storage space?

## Approaches

- 1 Choose coarser granularity
  - ⇒ loose information
  - ⇒ sometime positive: Information overload

# Overview

- Given that provenance is large
- How to save storage space?

## Approaches

- 1 Choose coarser granularity
  - ⇒ loose information
  - ⇒ sometime positive: Information overload
- 2 Choose compact provenance model
  - ⇒ may lose information
  - ⇒ restrict to important information





# Compression for Provenance

## Goals for Compressing Provenance

- 1 Reduce storage size
  - That is why we are doing that at all
- 2 Efficient compression/decompression
  - E.g., Compression useless if it takes days to compress
- 3 Lossless
  - or only loose unimportant information
- 4 Show Compressed representation to user
  - Compressed provenance as a summary
  - Reduce information overload
- 5 Execute queries over compressed representation (partially?)
  - Save cost of compression/decompression
  - Best possible outcome: query performance increase

# Choose Compression Techniques

- ① Generic Compression Techniques
  - E.g., Dictionary compression (LZ77)

# Choose Compression Techniques

## 1 Generic Compression Techniques

- E.g., Dictionary compression (LZ77)

## 2 Exploiting structure of provenance

- Repeating part in the provenance
  - $\Rightarrow$  Specialized dictionary compression
- Structure in provenance imposed by transformation
  - E.g., provenance of  $R \bowtie S$  always has tuples from  $R$  combined with tuples from  $S$

# Generic Compression

## Properties

- Good compression rates
- No need to adapt to provenance
- Output incomprehensible for human
- Compression/decompression quite expensive
- Queries over compressed provenance unlikely
  - Unclear how to extract parts of provenance

## Example

try:

```
echo "(r1 x s1) + t1" | gzip -cf
```

result:

```
?m|P?(2T?P(6?T?V(1????o?
```

BY

# Factoring out Common Parts from Provenance

## Overlap in provenance

- Recall alternatives
- Reuse of subquery results

## Example

$$q = \pi_a(R \bowtie_{a=d} S) \times U$$

	<b>R</b>				<b>S</b>		<b>U</b>		<b>Q</b>	
<i>r</i> <sub>1</sub>	<b>a</b>	<b>b</b>	<b>c</b>		<b>d</b>	<i>s</i> <sub>1</sub>	<b>e</b>	<i>t</i> <sub>1</sub>	<b>a</b>	<b>e</b>
	1	1	1		1		Hello		1	Hello
<i>r</i> <sub>2</sub>	2	3	2		2	<i>s</i> <sub>2</sub>	Test	<i>t</i> <sub>2</sub>	1	Test
<i>r</i> <sub>3</sub>	5	1	1					<i>t</i> <sub>3</sub>	2	Hello
								<i>t</i> <sub>4</sub>	2	Test

# Factoring out Common Parts from Provenance

## Overlap in provenance

- Recall alternatives
- Reuse of subquery results

## Example

$$q = \pi_a(R \bowtie_{a=d} S) \times U$$

### Provenance Polynomials

$$\mathbb{N}[I](q, t_1) = (r_1 \times s_1) \times u_1$$

$$\mathbb{N}[I](q, t_2) = (r_1 \times s_1) \times u_2$$

R		
a	b	c
$r_1$	1	1
$r_2$	2	3
$r_3$	5	1

S	
d	
$s_1$	1
$s_2$	2

U	
e	
$u_1$	Hello
$u_2$	Test

Q	
a	e
$t_1$	1
$t_2$	1
$t_3$	2
$t_4$	2

## Factoring out Common Parts from Provenance

## Overlap in provenance

- Recall alternatives
- Reuse of subquery results

## Example

$$q = \pi_a(R \bowtie_{a=d} S) \times U$$

## Why-provenance

$$\text{Why}(q, t_1) = \{\{r_1, s_1, u_1\}\}$$

$$\text{Why}(q, t_2) = \{\{r_1, s_1, u_2\}\}$$

	R		
	a	b	c
$r_1$	1	1	1
$r_2$	2	3	2
$r_3$	5	1	1

	S
	d
$s_1$	1
$s_2$	2

	U
	e
$u_1$	Hello
$u_2$	Test

	Q	
	a	e
$t_1$	1	Hello
$t_2$	1	Test
$t_3$	2	Hello
$t_4$	2	Test





# How to Detect Common Elements?

## Brute force

- Use binary or text representation of provenance
- Find common substrings (bit-patterns)
- Simple approach quadratic in size of provenance (of two results)

## Improvements

- Take equivalences on provenance into account
  - E.g., Why-provenance witnesses are sets  $\Rightarrow \{a, b\} = \{b, a\}$
  - E.g., Equivalences on polynomials
 
$$r_1 \times (s_1 + s_2) = (r_1 \times s_1) + (r_1 \times s_2)$$
- Use knowledge about query to improve matching performance
  - E.g.,  $(R \bowtie S) \times U$ : the provenance of each result tuple from  $R \bowtie S$  will be repeated for each tuple in  $U$

# Factorization of Provenance

## Properties

- Compression rate depends on query
- No need to adapt to provenance
  - compress provenance as text
  - without adaptation we may loose opportunities
- human readable: e.g., graph representation
- Queries over provenance: only small changes
- General approach of structural matching to detect to expensive
- Unclear how to integrate with provenance computation on demand

# Factoring out Structural Provenance Parts

## Rationale

- Analyze query to predetermine structure of provenance
- Factor out the structural parts from the representation

## Example

- $q = \pi_a(R \bowtie_{b=d} S)$
- Provenance of a result tuple (Set semantics is sum of multiplications of tuple from  $R$  and tuple from  $S$ )
- $\Rightarrow$  No need to store addition and multiplication operations if we know query

$$\mathbb{N}[I](q, t) = r_1 \times s_1 + r_1 \times s_2 + r_3 \times s_5 \quad \Rightarrow \quad r_1 s_1 r_1 s_2 r_3 s_5$$

BY

# Factoring out Structural Provenance Parts

## How to store?

- Compressed representation only interpretable with query!
- $\Rightarrow$  Store query with compressed representation
  - Easy:  $\pi_a(R \bowtie_{b=d} S), r_1 s_1 r_1 s_2 r_3 s_5$
- Alternatively store instructions for decompression
  - May be more compact
  - E.g., pattern  $(\sum e_i \times e_{i+1})$

# Factoring out Structural Provenance Parts

## Properties

- Compression rate depends on query, usually low constant factor
- Need query to interpret it
- Need representation for alternative patterns
- Not really human readable
- Queries over provenance: some adaptations
- After static analysis of query its simple
- More or less clear how to integrate with on-demand provenance tracking

3Y

# Symbolic and Declarative Representations of Provenance

- So far: explicitly listing tuples (input data) in provenance
- The amount of input data in provenance may be huge
- ⇒ Find more compact representations

## Example

- The provenance of each result tuple  $t$
- contains all tuples from  $S$  with  $b < a$
- Assume  $|S| = 1,000,000$

```
SELECT *  
FROM R  
WHERE EXISTS (SELECT * FROM S WHERE S.b < R.a)
```

BY

# Symbolic and Declarative Representations of Provenance

## Rationale

- Queries are concise representations of large number of tuples
- $\Rightarrow$  Replace part of the provenance with an expression to compute it

## Example

- $PI(q, t) = \{ \langle r_1, s_2 \rangle, \langle r_1, s_5 \rangle, \dots \}$
- contains all tuples from  $S$  with  $b < a$
- $\Rightarrow PI(q, t) = \{ \langle r_1, \sigma_{R.a < b}(S) \rangle \}$

```

SELECT *
FROM R
WHERE EXISTS (SELECT * FROM S WHERE S.b < R.a)
  
```

BY



# Symbolic and Declarative Representations of Provenance

## Use case: Queries over Provenance

- Assume user runs query over Perm provenance to retrieve provenance of specific tuple  $t$
- Provenance is stored as queries
- ⇒ We can delay the generation of actual provenance to when its needed
- ⇒ Improve performance of queries

## Challenges

- How to decide when to use queries vs. actual data?
- Dynamically interpreting and executing query require significant changes to system

# Symbolic and Declarative Representations of Provenance

## Properties

- Compression Rates can be significant
- Human readable: depends on user background
- Advantage for query processing over provenance
- Integrate with query engines is hard
- How to determine when to use symbolic representation?



# Provenance Model Specific Methods

## Factorization of Provenance Polynomials

- Factorization of Polynomials  $r_1 \times r_2 + r_1 \times r_3 = r_1(r_2 + r_3)$
- Equivalent polynomials can sometimes differ exponentially in size!
- Find different factorization to save space

# Provenance Model Specific Methods

## Factorization of Provenance Polynomials

- Factorization of Polynomials  $r_1 \times r_2 + r_1 \times r_3 = r_1(r_2 + r_3)$
- Equivalent polynomials can sometimes differ exponentially in size!
- Find different factorization to save space

## Example

- E.g.,  $\pi_1(R_1 \times R_2 \times \dots \times R_n)$
- provenance is  $(r_{11} \times r_{21} \times \dots \times r_{n1}) + (r_{12} \times r_{21} \times \dots \times r_{n1}) + \dots$
- Size of provenance polynomial:  $|R_1| \times |R_2| \times \dots \times |R_n|$

BY

# Provenance Model Specific Methods

## Factorization of Provenance Polynomials

- Factorization of Polynomials  $r_1 \times r_2 + r_1 \times r_3 = r_1(r_2 + r_3)$
- Equivalent polynomials can sometimes differ exponentially in size!
- Find different factorization to save space

## Example

- E.g.,  $\pi_1(R_1 \times R_2 \times \dots \times R_n)$
- provenance is  $(r_{11} \times r_{21} \times \dots \times r_{n1}) + (r_{12} \times r_{21} \times \dots \times r_{n1}) + \dots$
- Size of provenance polynomial:  $|R_1| \times |R_2| \times \dots \times |R_n|$
- Equivalent polynomials:  
 $(r_{11} + r_{12} + \dots + r_{1m_1}) \times (r_{21} + r_{22} + \dots) \times \dots$
- Size:  $|R_1| + \dots + |R_n|$

BY

# Outline

- 1 Provenance Storage
  - Introduction
  - Compression Methods for Provenance
  - Index Structures for Provenance
  - Recap

# Overview

## Why Index Structures for Provenance?

- Improve performance of specific access patterns
- Specialized index structures needed if access patterns differ from regular data
  - E.g., path queries on provenance

# Overview

## Why Index Structures for Provenance?

- Improve performance of specific access patterns
- Specialized index structures needed if access patterns differ from regular data
  - E.g., path queries on provenance

## Access Patterns

- Forward queries:
  - Give me all output data items that are derived from input data item  $x$
- Backward queries:
  - Give me all input data items that are in the provenance of output data item  $x$
- For sets?

BY



# Overview

## Why Index Structures for Provenance?

- Improve performance of specific access patterns
- Specialized index structures needed if access patterns differ from regular data
  - E.g., path queries on provenance

## Note

- Use generic data processing model:
  - Atomic unit of input and output data: **data item**
  - Transformations: DAG (directed acyclic graph)





# Interval Encoding of Provenance

## Provenance Model

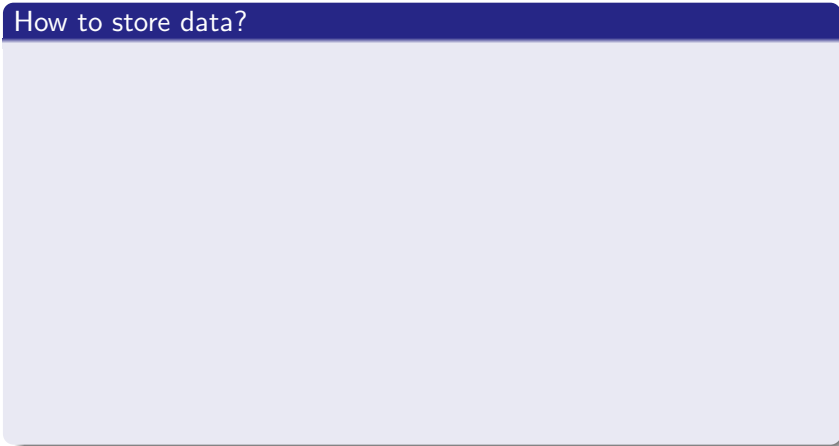
- Tree (directed acyclic graph (DAG))
  - Nodes are data items and transformations
- Stored in relational DB

## Queries

- **Forward** find all ancestors of node
- **Backward** find all descendants of node

# Interval Encoding of Provenance

How to store data?



BY

# Interval Encoding of Provenance

## How to store data?

- ① Edge relation ( $node_1, node_2$ )
  - Storage cost is small
  - Queries are recursive!

# Interval Encoding of Provenance

## How to store data?

- ① Edge relation ( $node_1, node_2$ )
  - Storage cost is small
  - Queries are recursive!
- ② Transitive closure of edge relation
  - Storage costs are high!
  - Queries are simple

# Interval Encoding of Provenance

## How to store data?

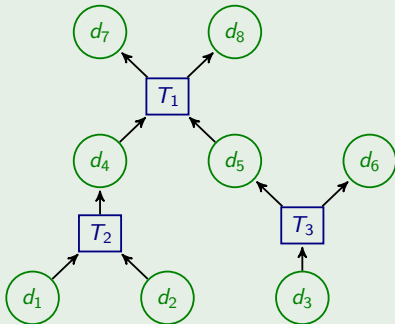
- 1 Edge relation ( $node_1, node_2$ )
  - Storage cost is small
  - Queries are recursive!
- 2 Transitive closure of edge relation
  - Storage costs are high!
  - Queries are simple
- 3 Path relation
  - Give each path an identifier
  - Store triples
    - (path\_id, position, node)
  - ⇒ Storage costs are high, queries cheaper



# Example Provenance DAG Storage

## Example (Edge Relation)

<b>from</b>	<b>to</b>
$d_1$	$T_2$
$d_2$	$T_2$
$d_3$	$T_3$
$T_2$	$d_4$
$T_3$	$d_5$
$T_3$	$d_6$
$d_4$	$T_1$
$d_5$	$T_1$
$T_1$	$d_7$
$T_1$	$d_8$



# Example Provenance DAG Storage

## Example (Edge Relation)

### Backward query for $d_7$

```

WITH RECURSIVE reachable(from,to)
AS (
  SELECT * FROM edges
  UNION ALL
  SELECT l.from, r.to
  FROM reachable l,
       reachable r
  WHERE l.to = r.from
)
SELECT from
FROM reachable
WHERE to = d7;

```

BY

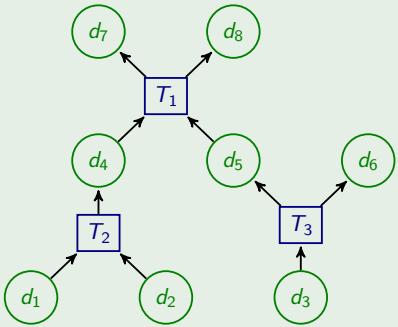




# Example Provenance DAG Storage

## Example (Path Relation)

path_id	order	node
1	1	$d_1$
1	2	$T_2$
1	3	$d_4$
1	4	$T_1$
1	5	$d_7$
2	1	$d_2$
2	2	$T_2$
...	...	...



# Example Provenance DAG Storage

## Example (Path Relation)

### Backward query for $d_7$

```
SELECT from  
FROM path l, path r  
WHERE l.path_id = r.path_id  
      AND r.node =  $d_7$   
      AND l.position < r.position
```



# Discussion of Storage Alternatives

- Both presented alternatives are unsatisfactory
- Number of nodes in graph:  $n$
- Edge Relation
  - Storage Size:  $n$
  - Recursive querying
- Transitive Closure
  - Storage size:  $n^2$
  - No recursive querying
  - Query table of size  $n^2$
- Paths
  - Storage size:  $O(f^h)$  with  $h$  = height and  $f$  is fan-out
  - No recursive querying
  - Query table of exponential size

# Interval Encoding of Trees

## Rationale

- Represent nodes as numeric intervals  $(l, r)$
- Interval inclusion determines parent child relationship
  - If  $c$  is child of  $p$  then
    - $p.l \leq c.l$  and  $c.r \leq p.r$
  - $\Rightarrow$  Reverse condition for reverse check
  - $\Rightarrow$  Same condition for descendent instead of child

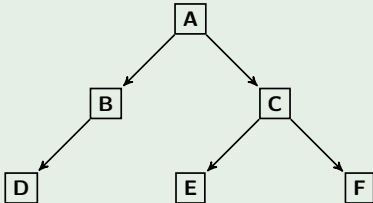


# Interval Encoding of Trees

## Rationale

- Represent nodes as numeric intervals  $(l, r)$
- Interval inclusion determines parent child relationship
  - If  $c$  is child of  $p$  then
    - $p.l \leq c.l$  and  $c.r \leq p.r$
  - $\Rightarrow$  Reverse condition for reverse check
  - $\Rightarrow$  Same condition for descendent instead of child

## Example



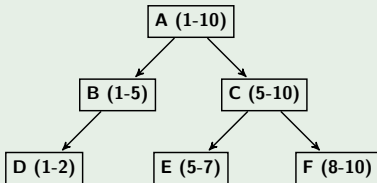
BY

# Interval Encoding of Trees

## Rationale

- Represent nodes as numeric intervals  $(l, r)$
- Interval inclusion determines parent child relationship
  - If  $c$  is child of  $p$  then
    - $p.l \leq c.l$  and  $c.r \leq p.r$
  - $\Rightarrow$  Reverse condition for reverse check
  - $\Rightarrow$  Same condition for descendent instead of child

## Example



# Interval Encoding of Trees

## Example

<b>node</b>	<b>l</b>	<b>r</b>
<i>A</i>	1	10
<i>B</i>	1	5
<i>C</i>	5	10
<i>D</i>	1	2
<i>E</i>	5	7
<i>F</i>	8	10

# Interval Encoding of Trees

## Example

### Backward query for $d_7$

```

SELECT from
FROM nodes a, nodes d
WHERE a.l < d.l AND d.r < a.r
      a.node =  $d_7$ 

```

# Interval Encoding of Trees

## Example

### Direct parent query for $d_7$

```

SELECT from
FROM nodes a, nodes d
WHERE a.l < d.l AND d.r < a.r
      a.node =  $d_7$ 
      AND NOT EXISTS (
        SELECT *
        FROM nodes m
        WHERE a.l < m.l AND m.r < a.r
              AND m.l < d.l AND d.r < m.r
      )

```

# Properties of Interval Encoding for Trees

- Storage size:  $n$
- Ancestor/Decedent Queries:  $O(n^2)$  (no index)
- Parent/Child Queries:  $O(n^3)$  (no index)

# Properties of Interval Encoding for Trees

## Comparison

- Storage
  - $O(n)$  (edge, interval),  $O(n^2)$  (edge transitive),  $O(f^h)$  (paths)





# Properties of Interval Encoding for Trees

## Comparison

- Storage
  - $O(n)$  (edge, interval),  $O(n^2)$  (edge transitive),  $O(f^h)$  (paths)
- Ancestor/Decedent Queries
  - $O(n * h)$  (edge),  $O(n^2)$  (interval, edge transitive),  $O(f^{2h})$  (paths)
  - **edge recursive!**
  - “interval” has inequality condition while “edge transitive” has equality
- Parent/Child Queries
  - $O(n)$  (edge),  $O(f^h)$  (paths),  $O(n^2)$  (edge transitive),  $O(n^3)$  (interval)

# Interval Encoding for Provenance DAGs

## Rationale

- Interval encoding is advantageous
- Only works for trees
  - One-dimensional intervals not enough!
- ... but provenance is DAG
- ⇒ need extension that deals with DAGs

# Interval Encoding for Provenance DAGs

## Rationale

- Interval encoding is advantageous
- **Only works for trees**
  - One-dimensional intervals not enough!
- **... but provenance is DAG**
- $\Rightarrow$  need extension that deals with DAGs

## N-dimensional Encoding

- Number of dimensions determined by graph structure
- **Sometimes even DAGs can be encoded using just one dimension**

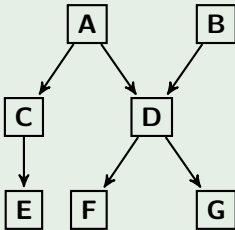
# Turning Graphs into Trees

## Idea

- Replicate nodes with more than one parent
- A node with  $n$  parents
  - Create  $n$  nodes with one parent
- ⇒ Parent/child relationships are the same
- ⇒ Can use interval encoding

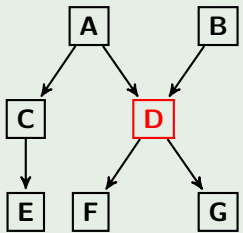
# Example Transform To Tree

## Example



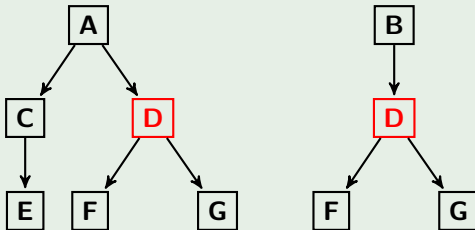
# Example Transform To Tree

## Example



# Example Transform To Tree

## Example



# Discussion Transformation to Tree

- Result can be interval encoded
- Size may be much larger
- Not all transformations necessary (overlapping)
- ⇒ Need method to test whether subgraph is interval encodable
- ⇒ Avoid “tree-ifying” subgraphs under replicated nodes



# When is a DAG Interval Encodable?

## Incomparability Graph $I_G$

- Model which nodes in a graph are incomparable
  - $\Rightarrow$  have no parent/child relationship
- Can be used to determine whether DAG is interval encodable
- Same node as original graph
- Edge between two node  $\Rightarrow$  nodes are incomparable

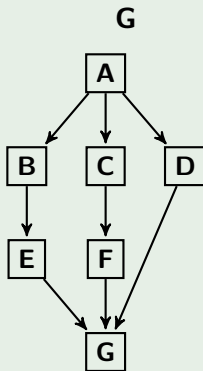
## Transitive orientable

- $I_G$  is transitively orientable if edges can be directed ...
  - Edges  $(u, v)$  and  $(v, w)$
  - $\Rightarrow$  there exists edge  $(u, w)$
- If  $I_G$  is transitively orientable  $\Rightarrow$  interval-encodable

BY

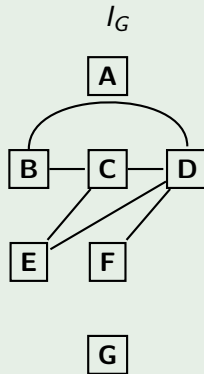
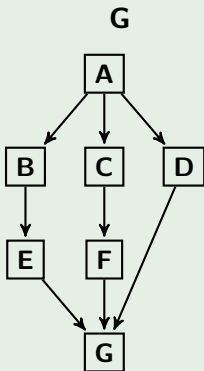
Example  $I_G$ 

## Example



# Example $I_G$

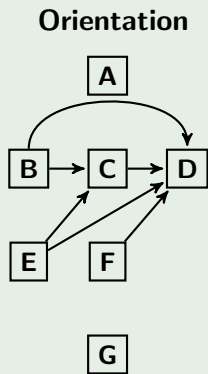
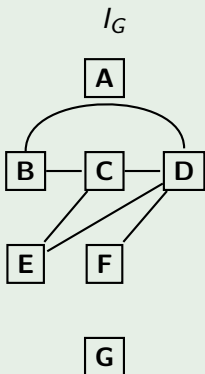
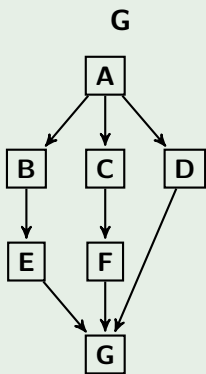
## Example



BY

# Example $I_G$

## Example



BY

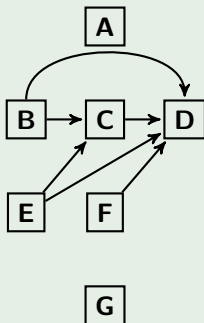
# Derive an Interval Encoding from $I_G$

## Approach

- Create two total orders  $L$  and  $L'$ 
  - Order  $L$  is created by traversing the nodes
    - From sources to sinks
    - Never visit a node before visiting its parents
    - If multiple nodes are ok  $\Rightarrow$  traverse in order of **orientation**
  - Order  $L'$  is created by using **reversed orientation**
- Assign each node  $n$  an interval based on its position in the orders
  - Let  $P_L(n)$  be the position in the order  $L$
  - Assign node  $n$  the interval  $[-P_L(n), P'_L(n)]$

# Example Derive an Interval Encoding

## Example



### Order $L$

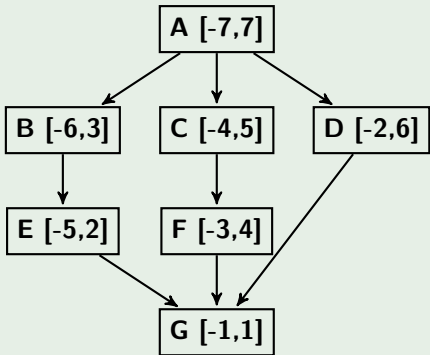
7 6 5 4 3 2 1  
 A > B > E > C > F > D > G

### Order $L'$

7 6 5 4 3 2 1  
 A > D > C > F > B > E > G

# Example Derive an Interval Encoding

## Example



BY

# Complete Approach for Interval Encoding

- 1 Detect maximal interval-encodable subgraphs (using  $I_G$ 's)
- 2 Replace these subgraphs with new nodes
- 3 Transform into tree
- 4 Substitute new nodes with original sub-graphs
- 5 Recompute  $I_G$  and create interval-encoding



# Conclusions Interval Encoding

## Advantages

- Index structure for provenance DAGs
- Supports efficient backward and forward queries
- Lower storage costs than alternatives

## Broader Perspective

- One provenance-specific access pattern
- ⇒ Improves performance and storage for this pattern
- E.g., provenance for relational queries
  - Non-typical access pattern for normal data access



# Open Problems for Indexing Provenance

- Integrate indices with automatic provenance generation
  - Construct index during provenance generation
  - Use index during provenance generation
- Incremental maintenance of indices
  - Most provenance index type are not incrementally maintainable
- Indices for both data and its provenance
- Combine indexing with compression





# Recap

## Provenance Compression

- Goals:
  - Reduce Size ;-)
  - Lossless?
  - Use as summary
  - Efficient compression/decompression (queries)
  - Integrate with generation and querying
- Approaches:
  - Generic Compression
    - Good compression
    - Expensive
    - Hard to evaluate queries over
    - Not human readable

BY



# Recap

## Provenance Compression

- Goals:

- Reduce Size ;-)
- Lossless?
- Use as summary
- Efficient compression/decompression (queries)
- Integrate with generation and querying

- Approaches:

- Use knowledge about transformation to reduce storage costs
  - Compression reasonable
  - Cost mostly paid upfront
  - Needs transformation to interpret it
  - Not human readable
  - Good for querying



# Recap

## Provenance Compression

- Goals:
  - Reduce Size ;-)
  - Lossless?
  - Use as summary
  - Efficient compression/decompression (queries)
  - Integrate with generation and querying
  
- Approaches:
  - Store expressions instead of actual provenance
    - Good compression if provenance follows pattern
    - Unclear how to determine expressions that represent provenance
    - Query processing needs run-time interpretation and execution of expressions to reconstruct provenance
    - Expressions can be good summaries if user understands expression language well

BY

# Recap

## Provenance Indexing

- What kind of access to support
  - Path-queries
  - Forward queries (which data items are derived from  $x$ )
  - Backward queries (which data items are in provenance of  $x$ )
  - For single data item  $x$  or sets?
- Approaches
  - Interval encoding of paths
    - For backward and forward queries
    - Less storage than storing all paths
    - Better performance than recursive queries
  - Adapted Information Retrieval Index Structure
    - Efficient forward and backward lookup for single items



# Literature II



A. Kementsietsidis and M. Wang.

Provenance Query Evaluation: What's so Special about it?.

In CIKM '09: Proceeding of the 18th Conference on Information and Knowledge Management, 681–690, 2009.



D. Olteanu and J. Zvodn.

On factorisation of provenance polynomials.

In TaPP '11: 3rd USENIX Workshop on the Theory and Practice of Provenance, 2011.



Y. Amsterdamer, D. Deutch, T. Milo, and V. Tannen.

On provenance minimization.

In Proceedings of the 30th Symposium on Principles of Database Systems (PODS), 141–152, 2011.