

# CompSci 516

# Data Intensive Computing Systems

## Lecture 21

## Datalog

## Instructor: Sudeepa Roy

# Announcement

- HW3 due next Wednesday: 11/16

# Today

- Datalog
  - for **recursion** in database queries
- A quick look at Incremental View Maintenance (IVM)

# Reading Material: Datalog

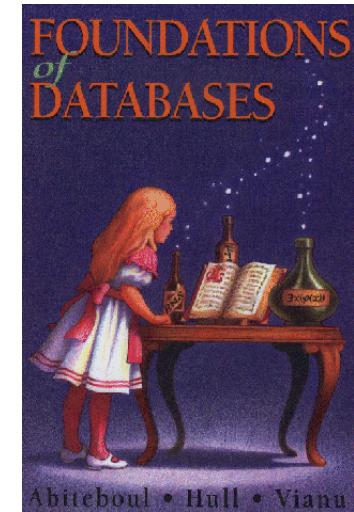
Optional:

1. The datalog chapters in the “Alice Book”

Foundations of Databases

Abiteboul-Hull-Vianu

Available online: <http://webdam.inria.fr/Alice/>



2. Datalog tutorial

SIGMOD 2011

“Datalog and Emerging Applications: An Interactive Tutorial”

# Brief History of Datalog

- Motivated by Prolog – started back in 1970-80's – then quiet for a long time
- A long argument in the Database community whether recursion should be supported in query languages
  - *"No practical applications of recursive query theory ... have been found to date"*—Michael Stonebraker, 1998  
*Readings in Database Systems, 3rd Edition* Stonebraker and Hellerstein, eds.
  - Recent work by Hellerstein et al. on Datalog-extensions to build networking protocols and distributed systems. [\[Link\]](#)

# Datalog is resurging!

- Number of papers and tutorials in DB conferences
- Applications in
  - data integration, declarative networking, program analysis, information extraction, network monitoring, security, and cloud computing
- Systems supporting datalog in both academia and industry:
  - Lixto (information extraction)
  - LogicBlox (enterprise decision automation)
  - Semmle (program analysis)
  - BOOM/Dedalus (Berkeley)
  - Coral
  - LDL++

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

## Recall our drinker example in RC (Lecture 4)

Find drinkers that frequent some bar that serves some beer they like.

$$Q(x) = \exists y. \exists z. \text{Frequents}(x, y) \wedge \text{Serves}(y, z) \wedge \text{Likes}(x, z)$$

Drinker example is from slides by Profs. Balazinska and Suciu  
and the [GUW] book

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

## Write it as a Datalog Rule

Find drinkers that frequent some bar that serves some beer they like.

RC:

$Q(x) = \exists y. \exists z. \text{Frequents}(x, y) \wedge \text{Serves}(y, z) \wedge \text{Likes}(x, z)$

Datalog:

$Q(x) :- \text{Frequents}(x, y), \text{Serves}(y, z), \text{Likes}(x, z)$

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

## Write it as a Datalog Rule

Find drinkers that frequent some bar that serves some beer they like.

RC:

$Q(x) = \exists y. \exists z. \text{Frequents}(x, y) \wedge \text{Serves}(y, z) \wedge \text{Likes}(x, z)$

Datalog:

$Q(x) :- \text{Frequents}(x, y), \text{Serves}(y, z), \text{Likes}(x, z)$

- Quick differences:
  - Uses “:-” not =
  - no need for  $\exists$  (assumed by default)
  - Use “,” on the right hand side (RHS)
  - Anything on RHS the of :- is assumed to be combined with  $\wedge$  by default
  - $\forall, \Rightarrow,$  not allowed – they need to use negation  $\neg$
  - Standard “Datalog” does not allow negation
  - Negation allowed in datalog with negation
- How to specify disjunction (OR /  $\vee$ )?

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

## Example: OR in Datalog

Find drinkers that (a) either frequent some bar that serves some beer they like, (b) or like beer “BestBeer”

RC:

$Q(x) = [\exists y. \exists z. \text{Frequents}(x, y) \wedge \text{Serves}(y, z) \wedge \text{Likes}(x, z)] \vee [\text{Likes}(x, "BestBeer")]$

Datalog:

$Q(x) :- \text{Frequents}(x, y), \text{Serves}(y, z), \text{Likes}(x, z)$   
 $Q(x) :- \text{Likes}(x, "BestBeer")$

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

## Example: OR in Datalog

Find drinkers that (a) either frequent some bar that serves some beer they like, (b) or like beer “BestBeer”, (c) or, frequent bars that “Joe” frequents

RC:

$$Q(x) = [\exists y. \exists z. \text{Frequents}(x, y) \wedge \text{Serves}(y, z) \wedge \text{Likes}(x, z)] \vee [\text{Likes}(x, \text{“BestBeer”})] \\ \vee [\exists w. \text{Frequents}(x, w) \wedge \text{Frequents}(\text{“Joe”}, w)]$$

Datalog:

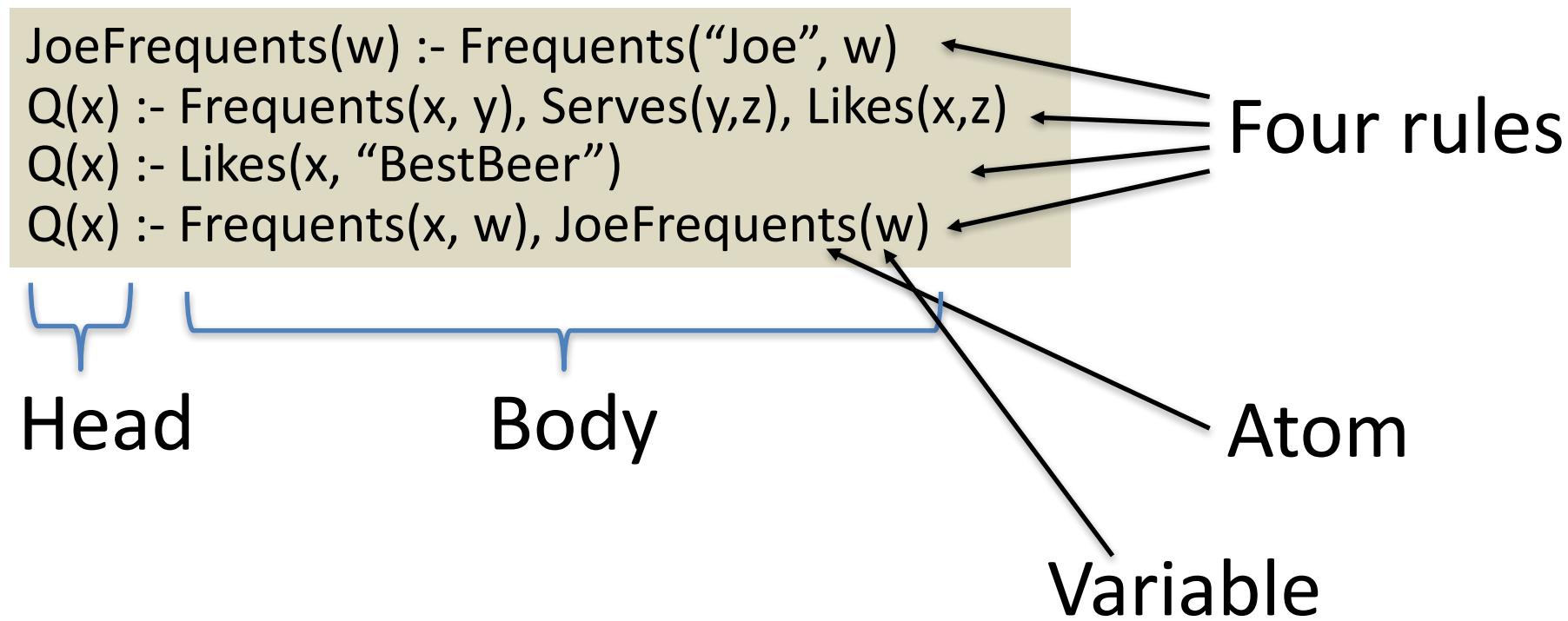
```
JoeFrequents(w) :- Frequents(“Joe”, w)
Q(x) :- Frequents(x, y), Serves(y, z), Likes(x, z)
Q(x) :- Likes(x, “BestBeer”)
Q(x) :- Frequents(x, w), JoeFrequents(w)
```

- To specify “OR”, write multiple rules with the same “Head”
- Next: terminology for Datalog

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

# Datalog Rules

- Each rule is of the form **Head :- Body**
- Each variable in the head of each rule must appear in the body of the rule



Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

# EDBs and IDBs

Tuple in an EDB or  
an IDB: a **FACT**

- **Extensional DataBases (EDBs)**

- Input relation names
- e.g. Likes, Frequents, Serves
- can only be on the RHS of a rule

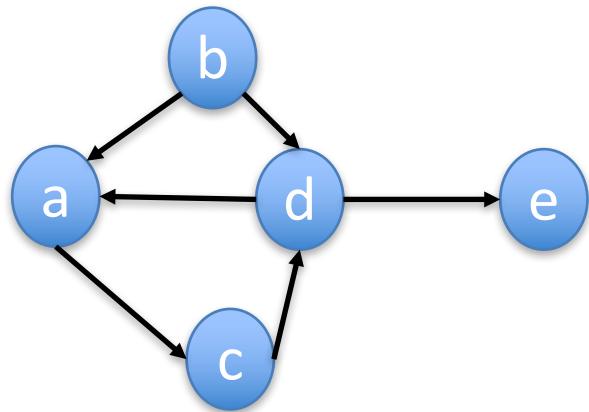
```
JoeFrequents(w) :- Frequents("Joe", w)
Q(x) :- Frequents(x, y), Serves(y,z), Likes(x,z)
Q(x) :- Likes(x, "BestBeer")
Q(x) :- Frequents(x, w), JoeFrequents(w)
```

either belongs to a  
given EDB relation,  
or is derived in an  
IDB relation

- **Intensional DataBases (IDBs)**

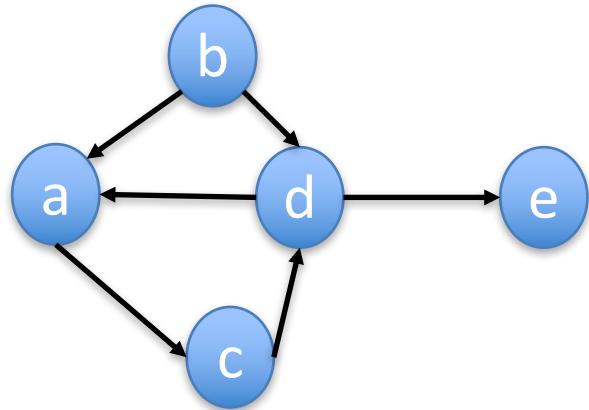
- Relations that are derived
- Can be intermediate or final output tables
- e.g. JoeFrequents, Q
- Can be on the LHS or RHS (e.g. JoeFrequents)

# Graph Example



v1	v2
a	c
b	a
b	d
c	d
d	a
d	e

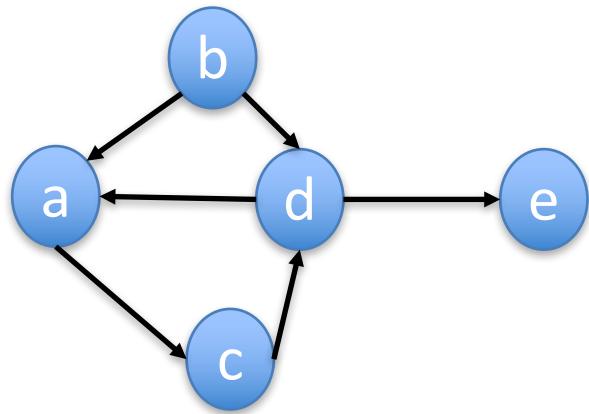
# Example 1



Write a Datalog program to find paths of length two (output start and finish vertices)

v1	v2
a	c
b	a
b	d
c	d
d	a
d	e

# Example 1

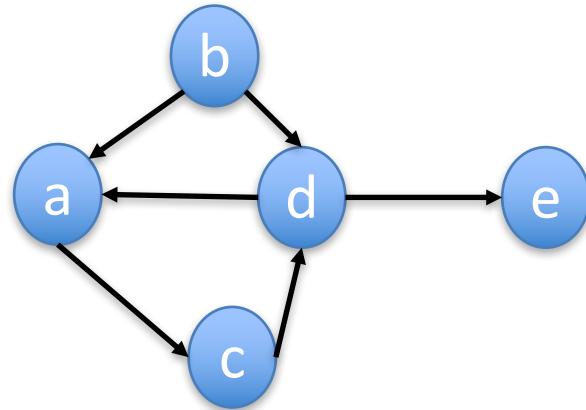


Write a Datalog program to find paths of length two (output start and finish vertices)

v1	v2
a	c
b	a
b	d
c	d
d	a
d	e

$P2(x, y) :- E(x, z), E(z, y)$

# Example 1: Execution



Write a Datalog program to find paths of length two (output start and finish vertices)

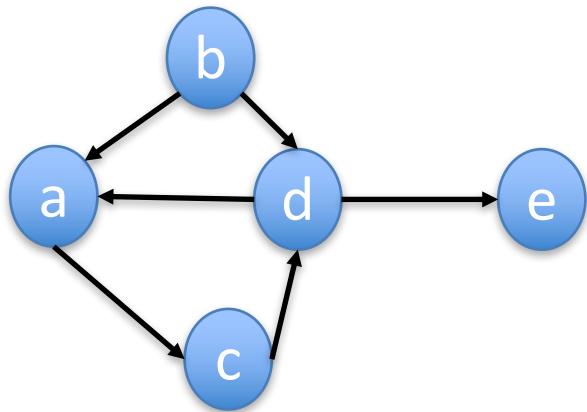
v1	v2
a	c
b	a
b	d
c	d
d	a
d	e

P2(x, y) :- E(x, z), E(z, y)

same as  $E \bowtie_{E.V2=E.V1} E$

P2	
v1	v2
a	d
b	c
b	e
c	a
c	e
d	c

# Example 2

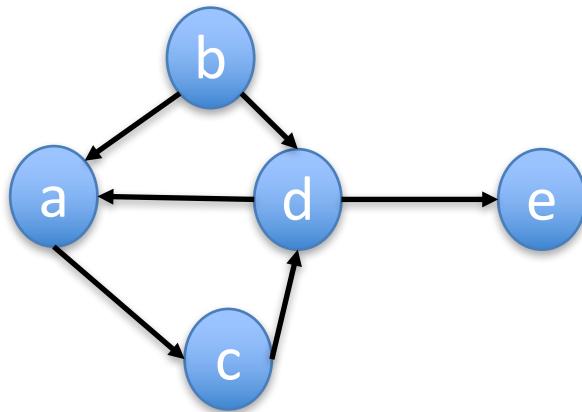


Write a Datalog program to find all pairs of vertices  $(u, v)$  such that  $v$  is reachable from  $u$

v1	v2
a	c
b	a
b	d
c	d
d	a
d	e

- Can you write a SQL/RA/RC query for reachability?

# Example 2

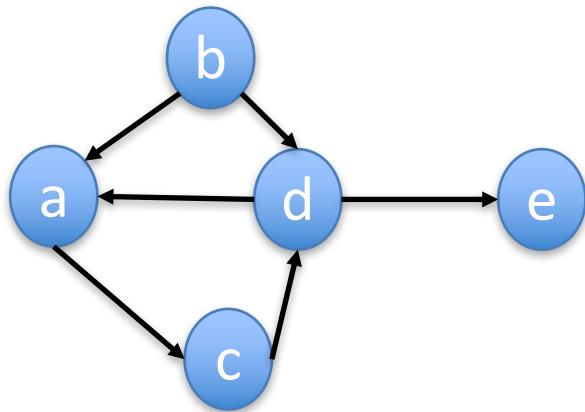


Write a Datalog program to find all pairs of vertices  $(u, v)$  such that  $v$  is reachable from  $u$

v1	v2
a	c
b	a
b	d
c	d
d	a
d	e

- Can you write a SQL/RA/RC query for reachability?
- NO - SQL/RA/RC cannot express reachability

# Example 2



Write a Datalog program to find all pairs of vertices  $(u, v)$  such that  $v$  is reachable from  $u$

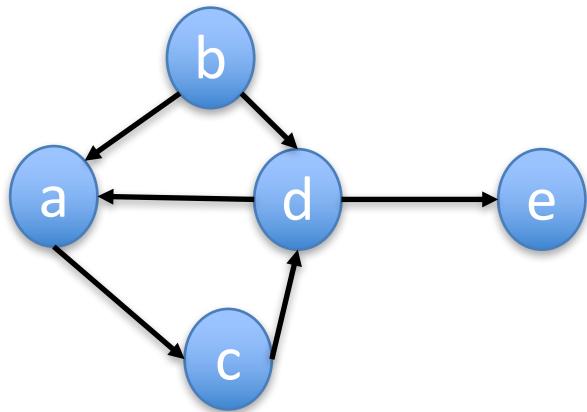
v1	v2
a	c
b	a
b	d
c	d
d	a
d	e

```

R(x, y) :- E(x, y)
R(x, y) :- E(x, z), R(z, y)
  
```

Option 1

# Example 2



Write a Datalog program to find all pairs of vertices  $(u, v)$  such that  $v$  is reachable from  $u$

```

R(x, y) :- E(x, y)
R(x, y) :- E(x, z), R(z, y)
  
```

non-linear

```

R(x, y) :- E(x, y)
R(x, y) :- R(x, z), R(z, y)
  
```

Option 1

linear

```

R(x, y) :- E(x, y)
R(x, y) :- R(x, z), E(z, y)
  
```

Option 2

Option 3

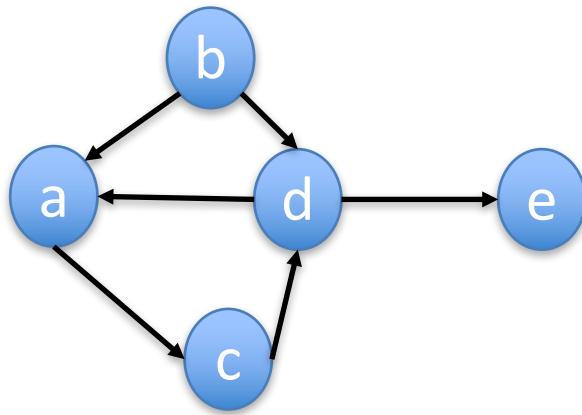
# Linear Datalog

- Linear rule
  - at most one atom in the body that is recursive with the head of the rule
  - e.g.  $R(x, y) :- E(x, z), R(z, y)$
- Linear datalog program
  - if all rules are linear
  - like linear recursion
- Top-down and bottom-up evaluation are possible
  - we will focus on bottom-up

Iteration 1

$R = E$

## Example 2: Execution



$R(x, y) :- E(x, y)$   
 $R(x, y) :- E(x, z), R(z, y)$

E	
V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

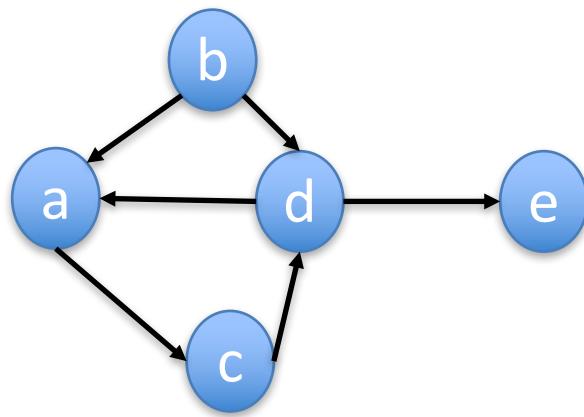
V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

Option 1

(vertices reachable in 1-hop by a direct edge)

Iteration 2

## Example 2: Execution



$R(x, y) :- E(x, y)$

$R(x, y) :- E(x, z), R(z, y)$

Option 1

E

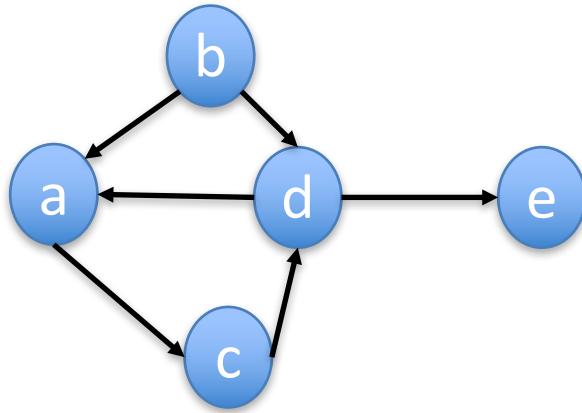
V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c

(vertices reachable in 2-hops)

Iteration 3

## Example 2: Execution



$R(x, y) :- E(x, y)$

$R(x, y) :- E(x, z), R(z, y)$

Option 1

(vertices reachable in 3-hops)

E

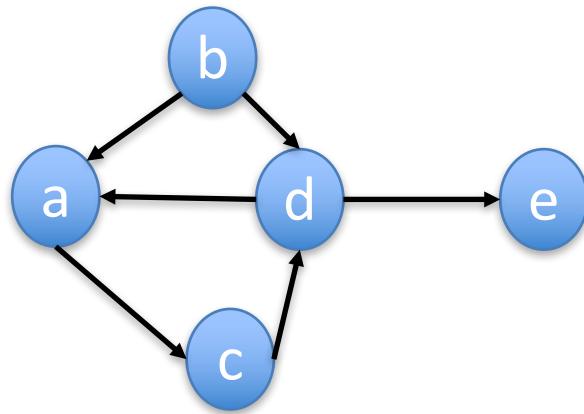
V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c
a	e
a	a
c	c
d	d

R

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c
a	e
a	a
c	c
d	d

Iteration 4

## Example 2: Execution



```
R(x, y) :- E(x, y)
R(x, y) :- E(x, z), R(z, y)
```

Option 1

R unchanged - stop

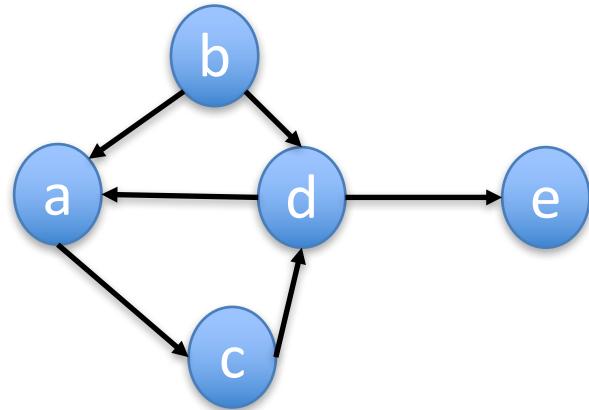
E

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c
a	e
a	a
c	c
d	d

R

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c
a	e
a	a
c	c
d	d

# Examples 3 and 4



v1	v2
a	c
b	a
b	d
c	d
d	a
d	e

Write a Datalog program to find all vertices reachable from **b**

```

R(x, y) :- E(x, y)
R(x, y) :- E(x, z), R(z, y)
QB(y) :- R(b, y)
  
```

Write a Datalog program to find all vertices u reachable from themselves  $R(u, u)$

```

R(x, y) :- E(x, y)
R(x, y) :- E(x, z), R(z, y)
Q(x) :- R(x, x)
  
```

# Termination of a Datalog Program

Q. A Datalog program always terminates – why?

# Termination of a Datalog Program

Q. A Datalog program always terminates – why?

- Because the values of the variables are coming from the “active domain” in the input relations (EDBs)
- Active domain = (finite) values from the (possibly infinite) domain appearing in the instance of a database
  - e.g. age can be any integer (infinite), but active domain is only finitely many in  $R(id, name, age)$
- Therefore the number of possible values in each of the IDBs is finite
- e.g. in the reachability example  $R(x, y)$ , the values of  $x$  and  $y$  come from  $\{a, b, c, d, e\}$ 
  - at most  $5 \times 5 = 25$  tuples possible in the IDB  $R(x, y)$
  - in any iteration, at least one new tuple is added in at least one IDB
  - Must stop after finite steps
  - e.g. the maximum number of iteration in the reachability example for any graph with five vertices is 25 (it was only 4 in our example)

# Bottom-up Evaluation of a Datalog Program

- Naïve evaluation
- Semi-naïve evaluation

# Naïve evaluation - 1

E

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

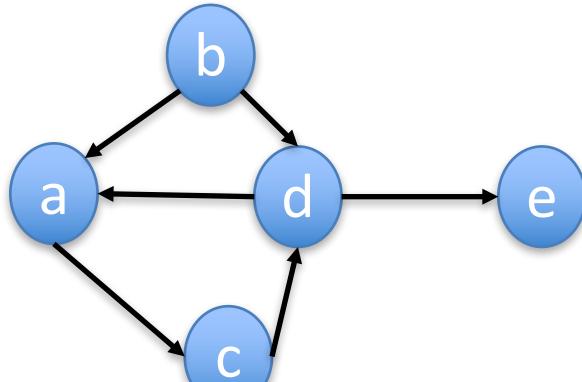
V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

Iteration 1:

$R = E = R_1$  (say)

In all subsequent iteration, check if any of the rules can be applied

Do union of all the rules with the same head IDB



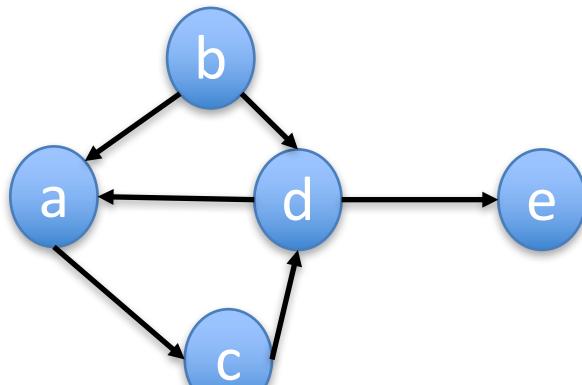
# Naïve evaluation - 2

E

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

Iteration 1:  
 $R = E = R_1$  (say)



Iteration 2:  
 $R = E \cup$   
 $E \bowtie R_1$   
=  $R_2$  (say)  
 $R_1 \neq R_2$   
so continue

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c

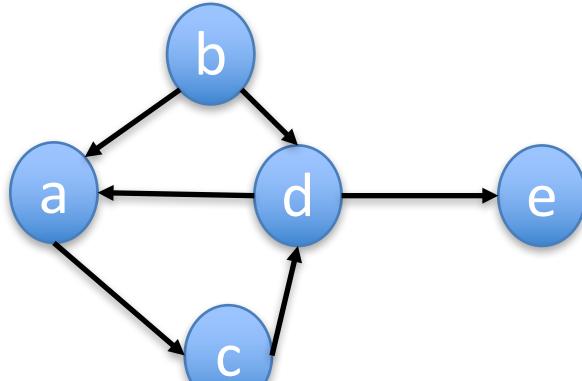
# Naïve evaluation - 3

E

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

Iteration 1:  
 $R = E = R_1$  (say)



Duke CS, Fall 2016

Iteration 2:  
 $R = E \cup$   
 $E \bowtie R_1$   
 $= R_2$  (say)  
 $R_1 \neq R_2$   
so continue

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
d	e
a	d
b	c
b	e
c	a
c	e
d	c

Iteration 3:  
 $R = E \cup$   
 $E \bowtie R_2$   
 $= R_3$  (say)  
 $R_2 \neq R_3$   
so continue

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c
a	e
a	a
c	c
d	d

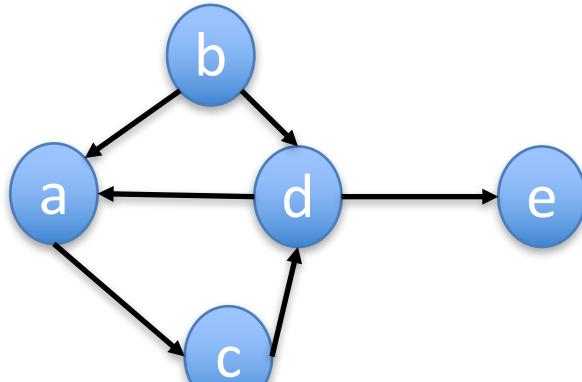
# Naïve evaluation - 4

E

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

Iteration 2:  
 $R = E \cup$   
 $E \bowtie R_1$   
 $= R_2$  (say)  
 $R_1 \neq R_2$   
so continue



Iteration 1:  
 $R = E = R_1$  (say)

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c

Iteration 3:  
 $R = E \cup$   
 $E \bowtie R_2$   
 $= R_3$  (say)

$R_2 \neq R_3$   
so continue

Iteration 4:  
 $R = E \cup$   
 $E \bowtie R_3$   
 $= R_4$  (say)

$R_3 = R_4$   
so STOP

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c
a	e
a	a
c	c
d	d

# Problem with Naïve Evaluation

- The same IDB facts are discovered again and again
  - e.g. in each iteration all edges in E are included in R
  - In the 2<sup>nd</sup>-4<sup>th</sup> iterations, the first six tuples in R are computed repeatedly
- Solution: Semi-Naïve Evaluation
- Work only with the new tuples generated in the previous iteration

# Semi-Naïve evaluation - 1

E

V1	V2	V1	V2
a	c	a	c
b	a	b	a
b	d	b	d
c	d	c	d
d	a	d	a
d	e	d	e

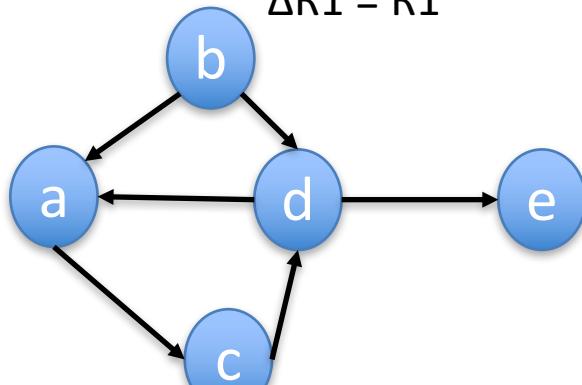
Initially:

$$R = \emptyset$$

Iteration 1:

$$R = E = R_1 \text{ (say)}$$

$$\Delta R_1 = R_1$$



# Semi-Naïve evaluation - 2

E

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

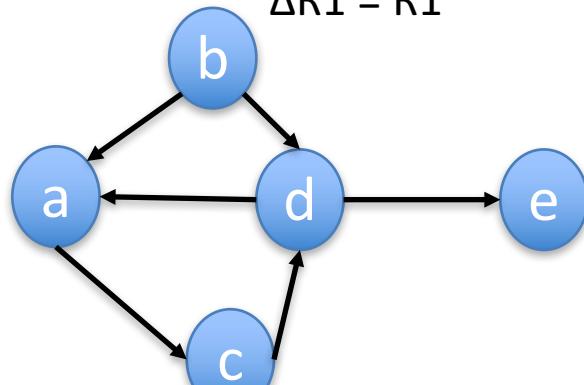
V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

Iteration 2:  
 $R = R1 \cup E \bowtie \Delta R1$   
= R2 (say)  
 $\Delta R2 = R2 - R1$   
 $\Delta R2 \neq \emptyset$   
so continue

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c

Initially:  
 $R = \emptyset$

Iteration 1:  
 $R = E = R1$  (say)  
 $\Delta R1 = R1$



# Semi-Naïve evaluation - 3

E

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

Iteration 2:  
 $R = R1 \cup E \bowtie \Delta R1$   
 $= R2$  (say)

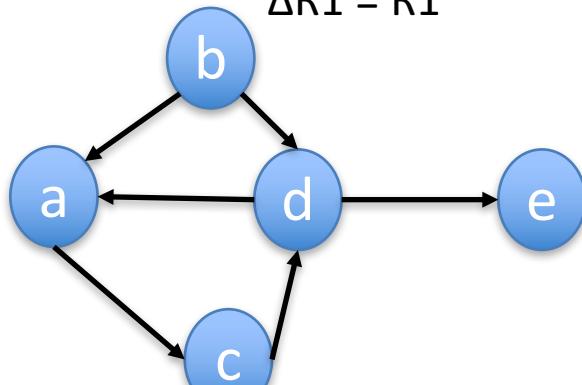
$\Delta R2 = R2 - R1$

$\Delta R2 \neq \Phi$   
so continue

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c

Initially:  
 $R = \Phi$

Iteration 1:  
 $R = E = R1$  (say)  
 $\Delta R1 = R1$



Duke CS, Fall 2016

Iteration 3:

$R = R2 \cup E \bowtie \Delta R2$   
 $= R3$  (say)

$\Delta R3 = R3 - R2$

$\Delta R3 \neq \Phi$   
so continue

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c
a	e
a	a
c	c
d	d

# Semi-Naïve evaluation - 4

E

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e

Iteration 2:  
 $R = R1 \cup E \bowtie \Delta R1$   
 $= R2$  (say)

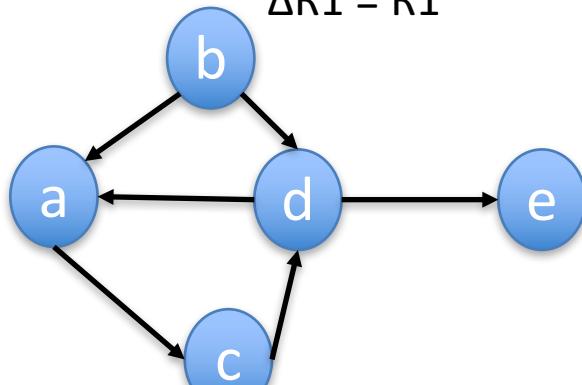
$\Delta R2 = R2 - R1$

$\Delta R2 \neq \emptyset$   
so continue

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c

Initially:  
 $R = \emptyset$

Iteration 1:  
 $R = E = R1$  (say)  
 $\Delta R1 = R1$



Iteration 3:  
 $R = R2 \cup E \bowtie \Delta R2$   
 $= R3$  (say)

$\Delta R3 = R3 - R2$

$\Delta R3 \neq \emptyset$   
so continue

V1	V2
a	c
b	a
b	d
c	d
d	a
d	e
a	d
b	c
b	e
c	a
c	e
d	c
a	e
a	a
c	c
d	d

Iteration 4:  
 $R = R3 \cup E \bowtie \Delta R3$   
 $= R4$  (say)

$\Delta R4 = R4 - R3$

$\Delta R = \emptyset$

**(CHECK ☺)**

**so STOP**

# Incremental View Maintenance (IVM)

- Why did the semi-naïve algorithm work?
- Because of the generic technique of Incremental View Maintenance (IVM)
- Suppose you have
  - a database  $D = (R_1, R_2, R_3)$
  - a query  $Q$  that gives answer  $Q(D)$
  - $D = (R_1, R_2, R_3)$  gets updated to  $D' = (R'_1, R'_2, R'_3)$
  - e.g.  $R'_1 = R_1 \cup \Delta R_1$  (insertion),  $R'_2 = R_2 - \Delta R_1$  (deletion) etc.

# Incremental View Maintenance (IVM)

- Why did the semi-naïve algorithm work?
- Because of the generic technique of Incremental View Maintenance (IVM)
- Suppose you have
  - a database  $D = (R_1, R_2, R_3)$
  - a query  $Q$  that gives answer  $Q(D)$
  - $D = (R_1, R_2, R_3)$  gets updated to  $D' = (R'_1, R'_2, R'_3)$
  - e.g.  $R'_1 = R_1 \cup \Delta R_1$  (insertion),  $R'_2 = R_2 - \Delta R_1$  (deletion) etc.
- **IVM:** Can you compute  $Q(D')$  using  $Q(D)$  and  $\Delta R_1, \Delta R_2, \Delta R_3$  without computing it from scratch (i.e. do not rerun the query  $Q$ )?

# IVM Example: Selection

V1	V2
a	c
b	a
d	a
c	d

R

$\sigma_{V1=b} R$

V1	V2
b	a

$\Delta R$

V1	V2
a	c
b	a
d	a
c	d
b	d
d	e

$R' = R \cup \Delta R$

$\sigma_{V1=b} R'$

V1	V2
b	a
b	d

V1	V2
b	a

$\sigma_{V1=b} R$

V1	V2
b	d

$\sigma_{V1=b} \Delta R$

- $\sigma_{V1=b}(R \cup \Delta R) = \sigma_{V1=b} R \cup \sigma_{V1=b} \Delta R$
- It suffices to apply the selection condition **only** on  $\Delta R$ 
  - and include with the original solution

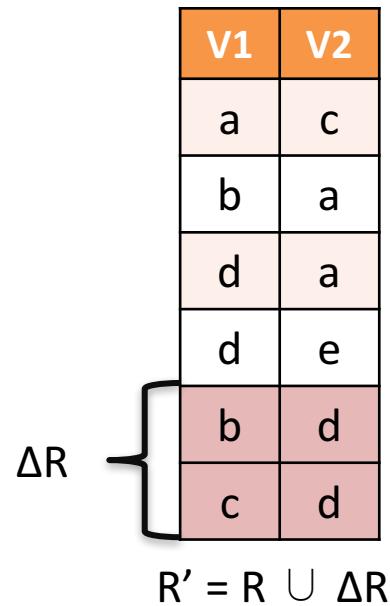
# IVM Example: Projection

V1	V2
a	c
b	a
d	a
d	e

R

$\pi_{V1} R$

V1
a
b
d



V1
a
b
d
c

$\pi_{V1} R'$

V1
a
b
d

$\pi_{V1} \Delta R$

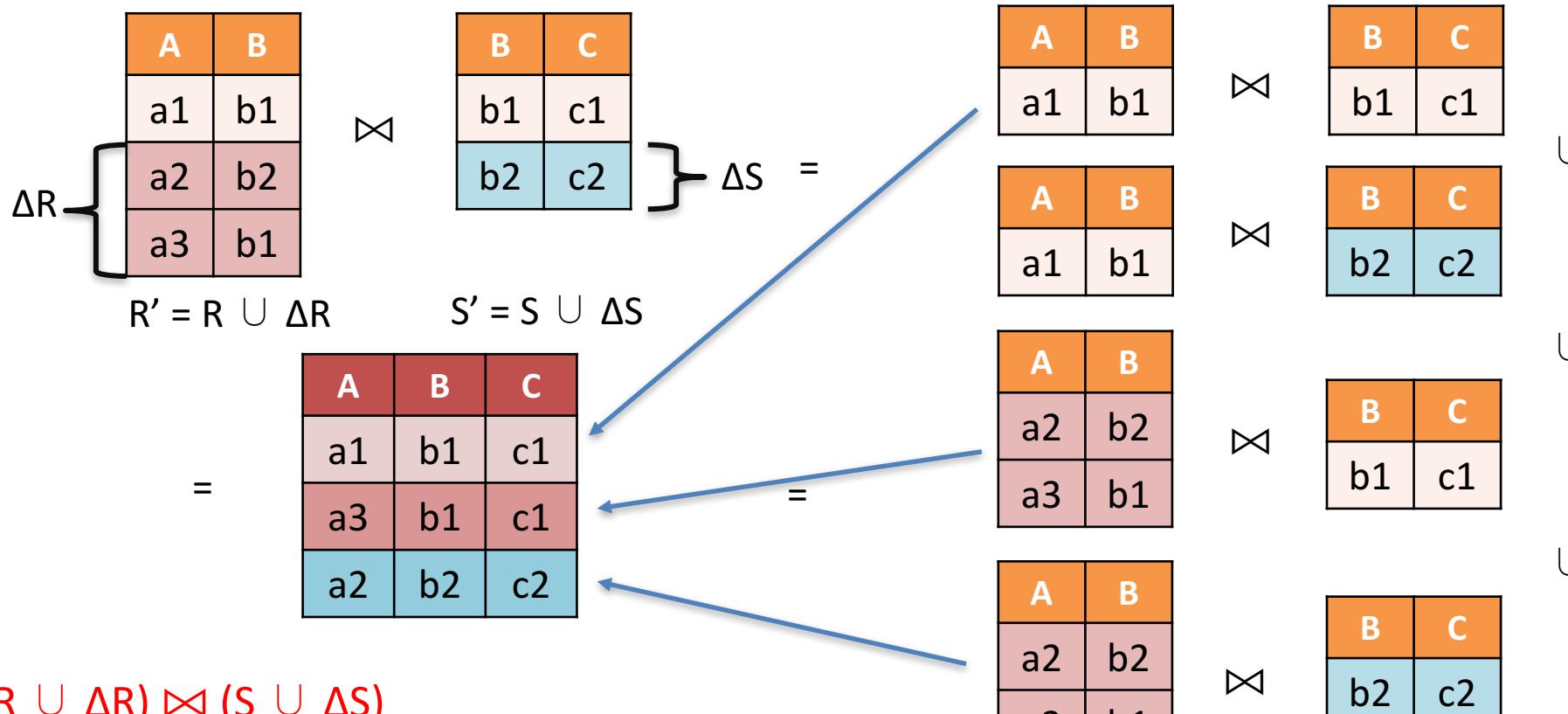
V1
b
c

$\pi_{V1} R$

- $\pi_{V1}(R \cup \Delta R) = \pi_{V1} R \cup \pi_{V1} \Delta R$
- It suffices to apply the projection condition **only** on  $\Delta R$ 
  - and include with the original solution

# IVM Example: Join

<table border="1" style="border-collapse: collapse; width: 100%; height: 100%;"> <thead> <tr> <th style="background-color: orange;">A</th> <th style="background-color: orange;">B</th> </tr> </thead> <tbody> <tr> <td>a1</td> <td>b1</td> </tr> </tbody> </table>	A	B	a1	b1	$\bowtie$	<table border="1" style="border-collapse: collapse; width: 100%; height: 100%;"> <thead> <tr> <th style="background-color: orange;">B</th> <th style="background-color: orange;">C</th> </tr> </thead> <tbody> <tr> <td>b1</td> <td>c1</td> </tr> </tbody> </table>	B	C	b1	c1	$=$	<table border="1" style="border-collapse: collapse; width: 100%; height: 100%;"> <thead> <tr> <th style="background-color: orange;">A</th> <th style="background-color: orange;">B</th> <th style="background-color: orange;">C</th> </tr> </thead> <tbody> <tr> <td>a1</td> <td>b1</td> <td>c1</td> </tr> </tbody> </table>	A	B	C	a1	b1	c1
A	B																	
a1	b1																	
B	C																	
b1	c1																	
A	B	C																
a1	b1	c1																



# IVM for Linear Datalog Rule

<table border="1" style="border-collapse: collapse; width: 100%; height: 100%;"> <thead> <tr> <th style="background-color: #f0a;">A</th><th style="background-color: #f0a;">B</th></tr> </thead> <tbody> <tr> <td>a1</td><td>b1</td></tr> </tbody> </table>	A	B	a1	b1	$\bowtie$	<table border="1" style="border-collapse: collapse; width: 100%; height: 100%;"> <thead> <tr> <th style="background-color: #f0a;">B</th><th style="background-color: #f0a;">C</th></tr> </thead> <tbody> <tr> <td>b1</td><td>c1</td></tr> </tbody> </table>	B	C	b1	c1	$=$	<table border="1" style="border-collapse: collapse; width: 100%; height: 100%;"> <thead> <tr> <th style="background-color: #f0a;">A</th><th style="background-color: #f0a;">B</th><th style="background-color: #f0a;">C</th></tr> </thead> <tbody> <tr> <td>a1</td><td>b1</td><td>c1</td></tr> </tbody> </table>	A	B	C	a1	b1	c1
A	B																	
a1	b1																	
B	C																	
b1	c1																	
A	B	C																
a1	b1	c1																

$\Delta R$

<table border="1" style="border-collapse: collapse; width: 100%; height: 100%;"> <thead> <tr> <th style="background-color: #f0a;">A</th><th style="background-color: #f0a;">B</th></tr> </thead> <tbody> <tr> <td>a1</td><td>b1</td></tr> </tbody> </table>	A	B	a1	b1	$\bowtie$	<table border="1" style="border-collapse: collapse; width: 100%; height: 100%;"> <thead> <tr> <th style="background-color: #f0a;">B</th><th style="background-color: #f0a;">C</th></tr> </thead> <tbody> <tr> <td>b1</td><td>c1</td></tr> </tbody> </table>	B	C	b1	c1
A	B									
a1	b1									
B	C									
b1	c1									
$\Delta R$	$S' = S$									

$R' = R \cup \Delta R$

$=$	<table border="1" style="border-collapse: collapse; width: 100%; height: 100%;"> <thead> <tr> <th style="background-color: #f0a;">A</th><th style="background-color: #f0a;">B</th><th style="background-color: #f0a;">C</th></tr> </thead> <tbody> <tr> <td>a1</td><td>b1</td><td>c1</td></tr> </tbody> </table>	A	B	C	a1	b1	c1
A	B	C					
a1	b1	c1					

$$(R \cup \Delta R) \bowtie (S \cup \Delta S)$$

$$= (R \bowtie S) \cup (R \bowtie \Delta S) \cup (\Delta R \bowtie S) \cup (\Delta R \bowtie \Delta S)$$

- $R(x, y) :- E(x, z), R(z, y)$ 
    - i.e.  $R_{\text{new}} = E \bowtie R$
  - But  $E$  is EDB
    - $\Delta E = \emptyset$
  - Therefore,
- $E \bowtie (R \cup \Delta R) = (E \bowtie R) \cup (E \bowtie \Delta R)$
- It suffices to join with the difference  $\Delta R$  and include in the result in the previous round  $E \bowtie R$
  - Advantage of having “linear rule”

# (Non-recursive) Datalog with Negation

- Recursion and negation together make Datalog execution complicated
  - there is a notion called “stratified semantic” for this purpose – compute IDB relations in strata/layers before taking a negation – **not covered in this class**
- We will only do **negation for non-recursive Datalog**

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

# Unsafe/Safe Datalog Rules

Find drinkers who like beer “BestBeer”

$Q(x) :- \text{Likes}(x, \text{"BestBeer"})$

Find drinkers who **DO NOT** like  
beer “BestBeer”

$Q(x) :- \neg \text{Likes}(x, \text{"BestBeer"})$

- What is the problem with this rule?
- What should this rule return?
  - names of all drinkers in the world?
  - names of all drinkers in the USA?
  - names of all drinkers in Durham?

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

# Problem with Negation in Datalog Rules

Find drinkers who like beer “BestBeer”

$Q(x) :- \text{Likes}(x, \text{"BestBeer"})$

Find drinkers who **DO NOT** like beer “BestBeer”

$Q(x) :- \neg \text{Likes}(x, \text{"BestBeer"})$

- What is the problem with this rule?
- Dependent on “domain” of drinkers
  - domain-dependent
  - infinite answers possible too..
    - keep generating “names”

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

# Problem with Negation in Datalog Rules

Find drinkers who like beer “BestBeer”

$Q(x) :- \text{Likes}(x, \text{"BestBeer"})$

Find drinkers who **DO NOT** like beer “BestBeer”

$Q(x) :- \neg \text{Likes}(x, \text{"BestBeer"})$

- Solution:
- Restrict to “active domain” of drinkers from the input *Likes* (or *Frequents*) relation
  - “domain-independence” – same finite answer always
- Becomes a “safe rule”



$Q(x) :- \text{Likes}(x, y), \neg \text{Likes}(x, \text{"BestBeer"})$

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

## RC → Datalog with negation → SQL (1/8)

Revisit example from Lecture 4

Query: Find drinkers that like some beer (so much) that they frequent all bars that serve it

$$Q(x) = \exists y. \text{Likes}(x, y) \wedge \forall z. (\text{Serves}(z, y) \Rightarrow \text{Frequents}(x, z))$$

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

## RC → Datalog with negation → SQL (2/8)

Revisit example from Lecture 4

Query: Find drinkers that like some beer so much that they frequent all bars that serve it

$$Q(x) = \exists y. \text{Likes}(x, y) \wedge \forall z. (\text{Serves}(z, y) \Rightarrow \text{Frequents}(x, z))$$

$P \Rightarrow Q$  same as  
 $\neg P \vee Q$

$$\equiv Q(x) = \exists y. \text{Likes}(x, y) \wedge \forall z. (\neg \text{Serves}(z, y) \vee \text{Frequents}(x, z))$$

Step 1: Replace  $\forall$  with  $\exists$  using de Morgan's Laws

$$Q(x) = \exists y. \text{Likes}(x, y) \wedge \neg \exists z. (\text{Serves}(z, y) \wedge \neg \text{Frequents}(x, z))$$

$\forall x P(x)$  same as  
 $\neg \exists x \neg P(x)$

$\neg(\neg P \vee Q)$  same as  
 $P \wedge \neg Q$

Ack: slides by Profs. Balazinska and Suciu

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

## RC → Datalog with negation → SQL (3/8)

Revisit example from Lecture 4

Query: Find drinkers that like some beer so much that they frequent all bars that serve it

$$Q(x) = \exists y. \text{Likes}(x, y) \wedge \forall z. (\text{Serves}(z, y) \Rightarrow \text{Frequents}(x, z))$$

Step 1: Replace  $\forall$  with  $\exists$  using de Morgan's Laws

$$Q(x) = \exists y. \text{Likes}(x, y) \wedge \neg \exists z. (\text{Serves}(z, y) \wedge \neg \text{Frequents}(x, z))$$

(new) Step 2: Make all subqueries domain independent

$$Q(x) = \exists y. \text{Likes}(x, y) \wedge \neg \exists z. (\text{Likes}(x, y) \wedge \text{Serves}(z, y) \wedge \neg \text{Frequents}(x, z))$$

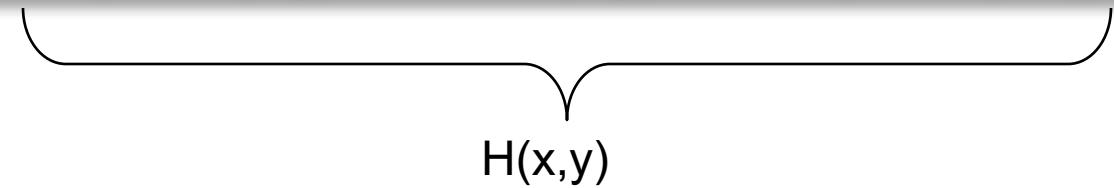
Ack: slides by Profs. Balazinska and Suciu

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

## RC → Datalog with negation → SQL (4/8)

Revisit example from Lecture 4

$$Q(x) = \exists y. \text{Likes}(x, y) \wedge \neg \exists z. (\text{Likes}(x, y) \wedge \text{Serves}(z, y) \wedge \neg \text{Frequents}(x, z))$$



(new) Step 3: Create a datalog rule for some subexpressions of the form  
 $\exists x \exists y \dots R(\dots) \wedge S(\dots) \wedge T(\dots) \wedge \dots$

$H(x, y) :- \text{Likes}(x, y), \text{Serves}(z, y), \neg \text{Frequents}(x, z)$   
 $Q(x) :- \text{Likes}(x, y), \neg H(x, y)$

Ack: slides by Profs. Balazinska and Suciu

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

## RC → Datalog with negation → SQL (5/8)

Revisit example from Lecture 4

```
H(x,y) :- Likes(x,y), Serves(z,y), not Frequents(x,z)
Q(x) :- Likes(x,y), not H(x,y)
```

Step 4: Write it in SQL

```
SELECT DISTINCT L.drinker FROM Likes L
WHERE .....
```

Ack: slides by Profs. Balazinska and Suciu

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

## RC → Datalog with negation → SQL (6/8)

Revisit example from Lecture 4

```
H(x,y) :- Likes(x,y), Serves(z,y), not Frequents(x,z)
Q(x) :- Likes(x,y), not H(x,y)
```

Step 4: Write it in SQL

```
SELECT DISTINCT L.drinker FROM Likes L
WHERE not exists
(SELECT * FROM Likes L2, Serves S
WHERE ... ....)
```

Ack: slides by Profs. Balazinska and Suciu

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

RC → Datalog with negation → SQL  
(7/8)

Revisit example from Lecture 4

```
H(x,y) :- Likes(x,y), Serves(z,y), not Frequents(x,z)
Q(x) :- Likes(x,y), not H(x,y)
```

Step 4: Write it in SQL

```
SELECT DISTINCT L.drinker FROM Likes L
WHERE not exists
  (SELECT * FROM Likes L2, Serves S
   WHERE L2.drinker=L.drinker and L2.beer=L.beer
     and L2.beer=S.beer
   and not exists (SELECT * FROM Frequents F
                  WHERE F.drinker=L2.drinker
                    and F.bar=S.bar))
```

Ack: slides by Profs. Balazinska and Suciu

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

## RC → Datalog with negation → SQL (8/8)

Revisit example from Lecture 4

```
H(x,y) :- Likes(x,y), Serves(z,y), not Frequents(x,z)
Q(x)   :- Likes(x,y), not H(x,y)
```

Unsafe rule

Sometimes can simplify the SQL query by using an unsafe datalog rule  
Correctness ensured by safe outermost rule

```
SELECT DISTINCT L.drinker FROM Likes L
WHERE not exists
  (SELECT * FROM Serves S
  WHERE L.beer=S.beer
  and not exists (SELECT * FROM Frequents F
  WHERE F.drinker=L.drinker
  and F.bar=S.bar))
```

Ack: slides by Profs. Balazinska and Suciu

# Optional/additional slides

## An Overview of Data Provenance with Annotations

Selected/adapted slides from the keynote by  
Prof. Val Tannen, EDBT 2010

(optional material: full slide deck is available on Val's webpage)

# Lineage

- [Cui-Widom-Wiener'00]
- Lineage:
  - Given a data item in the view
  - Determine the source data that produced it
  - The process by which it was produced

# Applications

- OLAP/OLAM (mining)
  - origin of anomalous data to verify reliability
- Scientific Databases
  - how answer was produced from raw data
- Online Network monitoring and Diagnosis system
  - identify faulty sensor from network monitors

# Applications

- Cleansed data feedback
  - Clean raw data and send report to sources
- Materialized view schema evolution
  - if view schema is changed (new column added), recomputation may not be necessary
- View update
  - translate view updates to base data updates

# Data Provenance

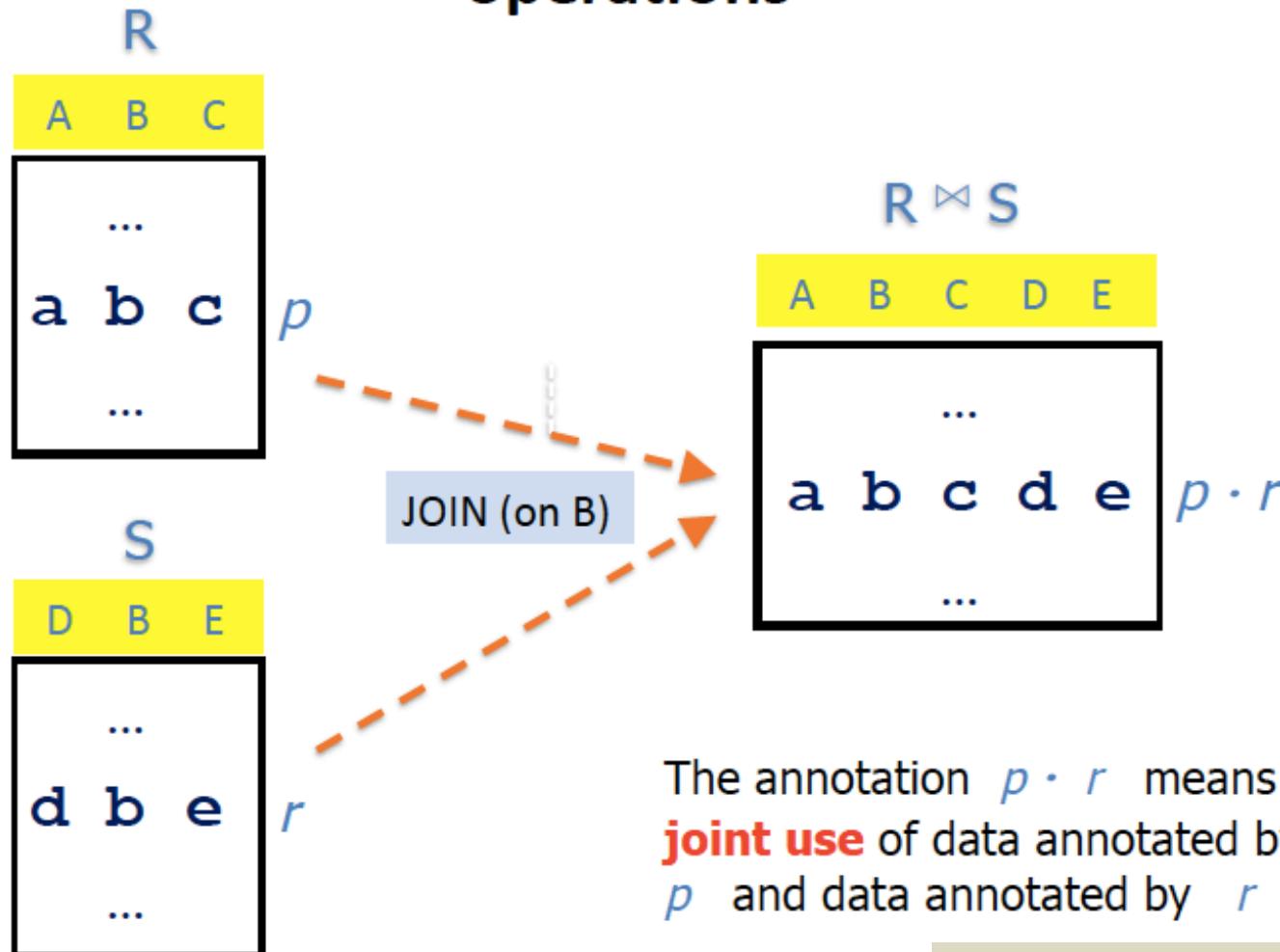
***provenance, n.***

*The fact of coming from some particular source or quarter; origin, derivation [Oxford English Dictionary]*

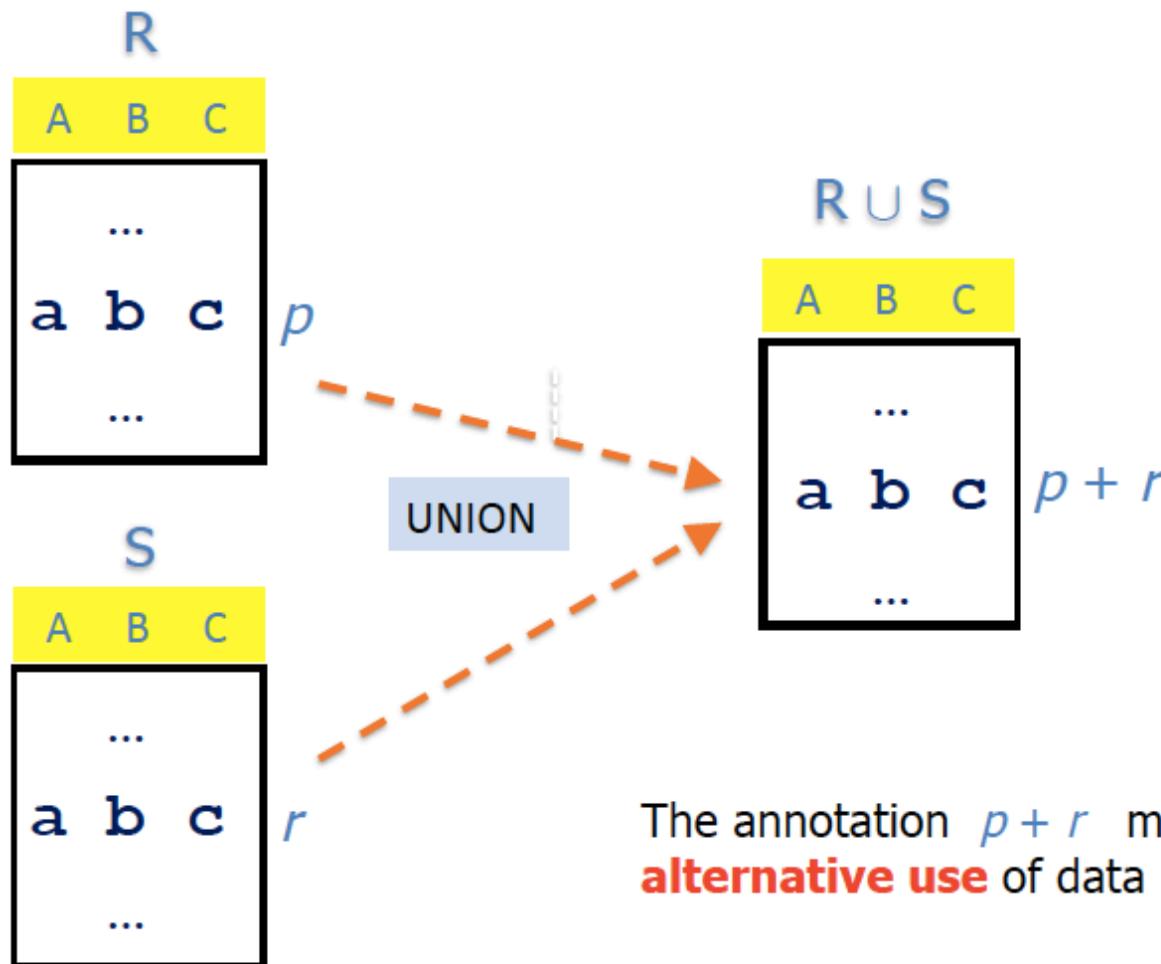
- **Data provenance** [BunemanKhannaTan 01]: aims to explain how a particular result (in an experiment, simulation, query, workflow, etc.) was derived.
- Most science today is **data-intensive**. Scientists, eg., biologists, astronomers, worry about data provenance all the time.

Slide by Val Tannen, EDBT 2010

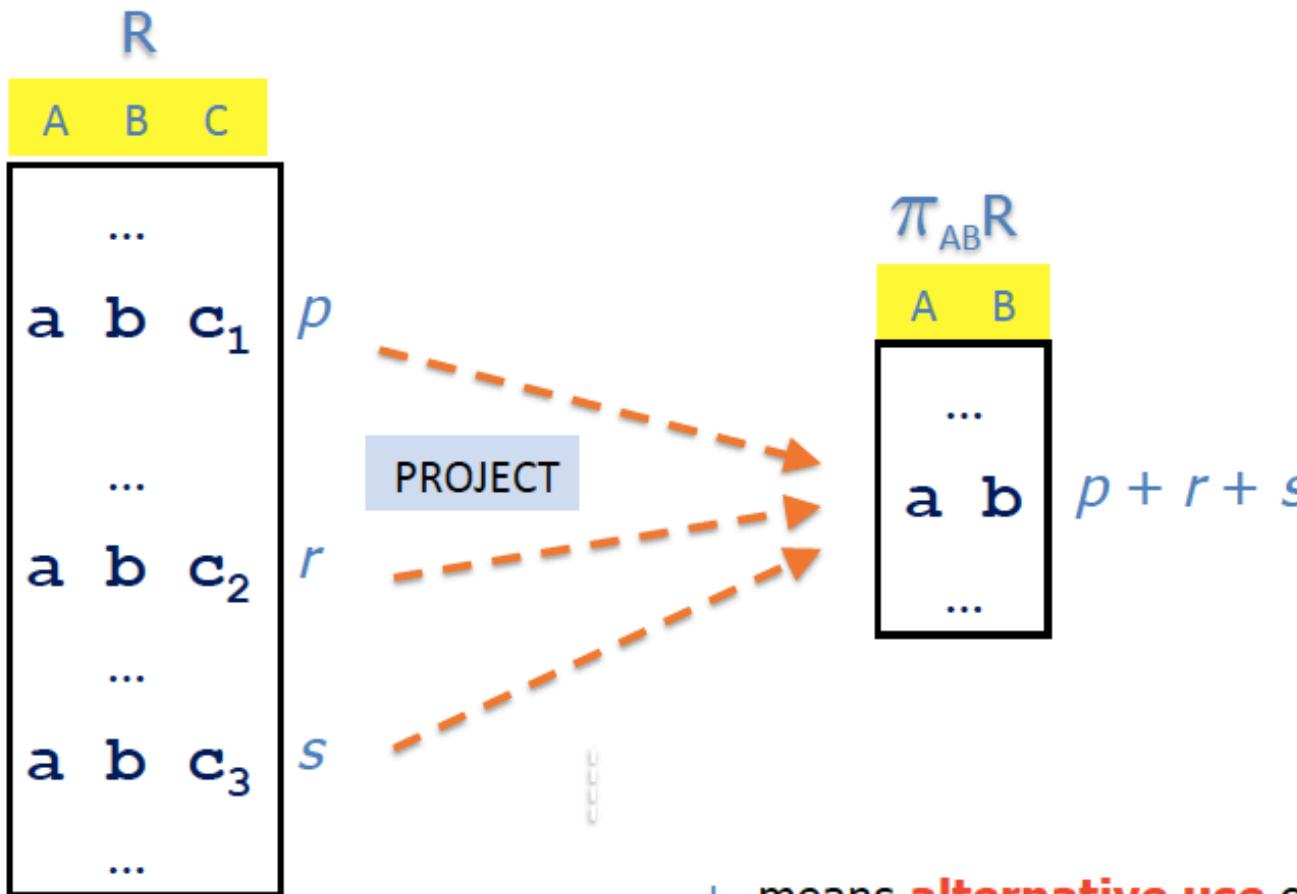
# Propagating annotations through database operations



## Another way to propagate annotations



# Another use of +



+ means **alternative use** of data

# An example in positive relational algebra (SPJU)

$R$	$Q = \sigma_{C=e} \pi_{AC} ( \pi_{AC} R \bowtie \pi_{BC} R \cup \pi_{AB} R \bowtie \pi_{BC} R )$																								
<table border="1"> <thead> <tr> <th>A</th><th>B</th><th>C</th></tr> </thead> <tbody> <tr> <td>a</td><td>b</td><td>c</td></tr> <tr> <td>d</td><td>b</td><td>e</td></tr> <tr> <td>f</td><td>g</td><td>e</td></tr> </tbody> </table> <p><math>p</math></p>	A	B	C	a	b	c	d	b	e	f	g	e	<table border="1"> <thead> <tr> <th>A</th><th>C</th></tr> </thead> <tbody> <tr> <td>a</td><td>c</td></tr> <tr> <td>a</td><td>e</td></tr> <tr> <td>d</td><td>c</td></tr> <tr> <td>d</td><td>e</td></tr> <tr> <td>f</td><td>e</td></tr> </tbody> </table> <p><math>(p \cdot p + p \cdot p) \cdot 0</math></p>	A	C	a	c	a	e	d	c	d	e	f	e
A	B	C																							
a	b	c																							
d	b	e																							
f	g	e																							
A	C																								
a	c																								
a	e																								
d	c																								
d	e																								
f	e																								
<table border="1"> <thead> <tr> <th>A</th><th>B</th><th>C</th></tr> </thead> <tbody> <tr> <td>a</td><td>b</td><td>c</td></tr> <tr> <td>d</td><td>b</td><td>e</td></tr> <tr> <td>f</td><td>g</td><td>e</td></tr> </tbody> </table> <p><math>r</math></p>	A	B	C	a	b	c	d	b	e	f	g	e	<table border="1"> <thead> <tr> <th>A</th><th>C</th></tr> </thead> <tbody> <tr> <td>a</td><td>c</td></tr> <tr> <td>a</td><td>e</td></tr> <tr> <td>d</td><td>c</td></tr> <tr> <td>d</td><td>e</td></tr> <tr> <td>f</td><td>e</td></tr> </tbody> </table> <p><math>p \cdot r \cdot 1</math></p>	A	C	a	c	a	e	d	c	d	e	f	e
A	B	C																							
a	b	c																							
d	b	e																							
f	g	e																							
A	C																								
a	c																								
a	e																								
d	c																								
d	e																								
f	e																								
<table border="1"> <thead> <tr> <th>A</th><th>B</th><th>C</th></tr> </thead> <tbody> <tr> <td>a</td><td>b</td><td>c</td></tr> <tr> <td>d</td><td>b</td><td>e</td></tr> <tr> <td>f</td><td>g</td><td>e</td></tr> </tbody> </table> <p><math>s</math></p>	A	B	C	a	b	c	d	b	e	f	g	e	<table border="1"> <thead> <tr> <th>A</th><th>C</th></tr> </thead> <tbody> <tr> <td>a</td><td>c</td></tr> <tr> <td>a</td><td>e</td></tr> <tr> <td>d</td><td>c</td></tr> <tr> <td>d</td><td>e</td></tr> <tr> <td>f</td><td>e</td></tr> </tbody> </table> <p><math>r \cdot p \cdot 0</math></p> <p><math>(r \cdot r + r \cdot s + r \cdot r) \cdot 1</math></p> <p><math>(s \cdot s + s \cdot r + s \cdot s) \cdot 1</math></p>	A	C	a	c	a	e	d	c	d	e	f	e
A	B	C																							
a	b	c																							
d	b	e																							
f	g	e																							
A	C																								
a	c																								
a	e																								
d	c																								
d	e																								
f	e																								

For selection we multiply  
with two special annotations, 0 and 1

# Provenance Example

Boolean query  $Q() :- \text{HasAsthma}(x), \text{Friend}(x, y), \text{Smoker}(y)$

HasAsthma	
$x_1$	Ann
$x_2$	Bob

Friend		
$y_1$	Ann	Joe
$y_2$	Ann	Tom
$y_3$	Bob	Tom

Smoker	
$z_1$	Joe
$z_2$	Tom

$$\text{Provenance } F_{Q,D} = x_1 y_1 z_1 + x_1 y_2 z_2 + x_2 y_3 z_2$$

- $x, y, z \in \{0, 1\}$
- 1 = present
- 0 = absent
- There are three alternative ways to derive the “true” answer

# Application to “Deletion Propagation”

optional slide

Boolean query  $Q: \exists x \exists y \text{ HasAsthma}(x) \wedge \text{Friend}(x, y) \wedge \text{Smoker}(y)$

HasAsthma	
$x_1$	Ann
$x_2$	Bob

Friend		
$y_1$	Ann	Joe
$y_2$	Ann	Tom
$y_3$	Bob	Tom

Smoker	
$z_1$	Joe
$z_2$	Tom

$$\text{Provenance } F_{Q,D} = x_1 y_1 z_1 + x_1 y_2 z_2 + x_2 y_3 z_2$$

[Green et al. '07]

- $x, y, z \in \{0, 1\}$
- What happens if Ann is deleted?
  - Does the answer change to false from true?

# Application to “Deletion Propagation”

optional slide

Boolean query  $Q: \exists x \exists y \text{ HasAsthma}(x) \wedge \text{Friend}(x, y) \wedge \text{Smoker}(y)$

HasAsthma	
x <sub>1</sub>	Ann
x <sub>2</sub>	Bob

Friend		
y <sub>1</sub>	Ann	Joe
y <sub>2</sub>	Ann	Tom
y <sub>3</sub>	Bob	Tom

Smoker	
z <sub>1</sub>	Joe
z <sub>2</sub>	Tom

$$\text{Provenance } F_{Q,D} = x_1 y_1 z_1 + x_1 y_2 z_2 + x_2 y_3 z_2$$

[Green et al. '07]

- $x, y, z \in \{0, 1\}$
- What happens if Ann is deleted?
  - Does the answer change to false from true?
- No need to re-evaluate the query
  - just plug in  $x_1 = 0$  and evaluate

0