

# CSE 544

# Relational Calculus

Lecture #2

January 11<sup>th</sup>, 2011

# Announcements

- HW1 is posted
- First paper review due next week
- Friday: project groups are due

# Announcements

- Guest lecture on Tuesday, 1/18, by Bill Howe  
*SQLShare: Smart Services for Ad Hoc Databases*
- Need to makeup **one** lecture this/next week.  
When ?
  - Friday 1/14, 12:30-2pm ?
  - Wednesday, 1/19, morning (9...12) ?
  - Friday, 1/21, late morning (11...1pm) ?

# Today's Agenda

Relational calculus:

- Domain Relational Calculus
- Non-recursive datalog (a reasonable abstraction of SQL)
- Relational algebra

They are equivalent and why we care

# Domain Relational Calculus

Given:

- A vocabulary:  $R_1, \dots, R_k$
- An arity,  $\text{ar}(R_i)$ , for each  $i=1, \dots, k$
- An infinite supply of variables  $x_1, x_2, x_3, \dots$
- Constants:  $c_1, c_2, c_3, \dots$

# Domain Relational Calculus

Terms  $t$  and Formulas  $\varphi$  are:

$$\begin{aligned} t ::= & x \mid a \quad /* \text{variable or constant} */ \\ \varphi ::= & R(t_1, \dots, t_k) \mid t_i = t_j \\ & \mid \varphi \wedge \varphi' \mid \varphi \vee \varphi' \mid \neg \varphi' \\ & \mid \forall x. \varphi \mid \exists x. \varphi \end{aligned}$$

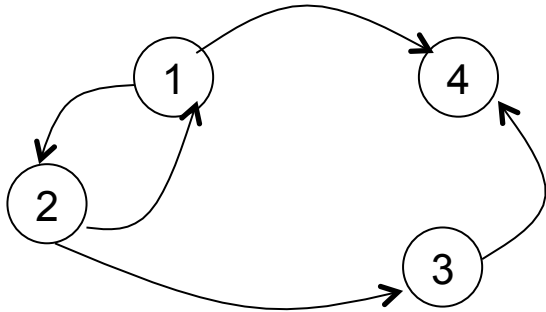
A query  $Q$  is any formula  $\varphi$ , usually written as:

$$Q = \{ \mathbf{x} \mid \varphi \}$$

where  $\mathbf{x} = (x_1, x_2, \dots)$  are the free variables in  $\varphi$ .

# Example: Querying a Graph

R encodes a graph



R=

1	2
2	1
2	3
1	4
3	4

What do these queries return ?

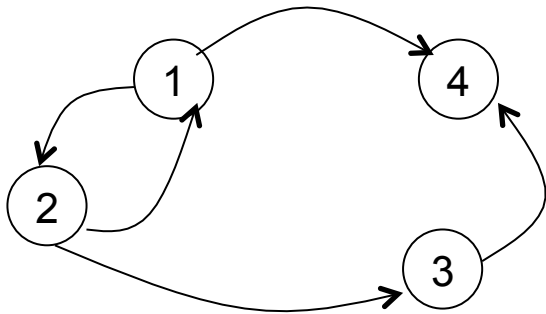
$$\{x \mid \exists y.R(x,y)\}$$

$$\{x \mid \exists y.\exists z.\exists u.(R(x,y) \wedge R(y,z) \wedge R(z,u))\}$$

$$\{(x,y) \mid \forall z.(R(x,z) \rightarrow R(y,z))\}$$

# Example: Querying a Graph

R encodes a graph



R=

1	2
2	1
2	3
1	4
3	4

What do these queries return ?

$$\{x \mid \exists y.R(x,y)\}$$

Nodes that have at least one child: {1,2,3}

$$\{x \mid \exists y.\exists z.\exists u.(R(x,y) \wedge R(y,z) \wedge R(z,u))\}$$

Nodes that have a great-grand-child: {1,2}

$$\{(x,y) \mid \forall z.(R(x,z) \rightarrow R(y,z))\}$$

Every child of x is a child of y:

{(1,1),(2,2),(3,1),(3,3),(4,1), (4,2), (4,3)}



# Semantics of a Formula

- A *database instance* (or a *model*) is  $\mathbf{D} = (D, R_1^D, \dots, R_k^D)$ 
  - $D$  = a set, called domain, or universe
  - $R_i^D \subseteq D \times D \times \dots \times D$ , ( $\text{ar}(R_i)$  times)  
 $i = 1, \dots, k$
- A *valuation* is a function  $s : \{x_1, x_2, \dots\} \rightarrow D$
- The semantics of a *formula*  $\varphi$  says when  $\mathbf{D} \models \varphi[s]$  holds

The domain  $D$  is not part of the database, yet we need it to define the semantics. Where ?

# Semantics of $\mathbf{D} \models \varphi[s]$

$\mathbf{D} \models (R(t_1, \dots, t_n)) [s]$

If  $(s(t_1), \dots, s(t_n)) \in R^D$

$\mathbf{D} \models (t = t') [s]$

If  $s(t) = s(t')$

$\mathbf{D} \models (\varphi \wedge \varphi') [s]$

If  $\mathbf{D} \models \varphi[s]$  and  $\mathbf{D} \models (\varphi') [s]$

$\mathbf{D} \models (\varphi \vee \varphi') [s]$

If  $\mathbf{D} \models \varphi[s]$  or  $\mathbf{D} \models \varphi'[s]$

$\mathbf{D} \models (\neg \varphi') [s]$

If it is not the case  $\mathbf{D} \models \varphi[s]$

# Semantics of $\mathbf{D} \models \varphi[s]$

Notation:  $s_{a/x}$  = the valuation  $s$  modified to map  $x$  to  $a$

$$\mathbf{D} \models (\forall x.\varphi) [s]$$

If for every constant  $a$  in  $D$ ,  
 $\mathbf{D} \models \varphi[s_{a/x}]$

$$\mathbf{D} \models (\exists x.\varphi) [s]$$

If for some constant  $a$  in  $D$ ,  
 $\mathbf{D} \models \varphi[s_{a/x}]$

# Semantics of $\mathbf{D} \models \varphi[s]$

Notation:  $s_{a/x}$  = the valuation  $s$  modified to map  $x$  to  $a$

$$\mathbf{D} \models (\forall x.\varphi) [s]$$

If for every constant  $a$  in  $D$ ,  
 $\mathbf{D} \models \varphi[s_{a/x}]$

$$\mathbf{D} \models (\exists x.\varphi) [s]$$

If for some constant  $a$  in  $D$ ,  
 $\mathbf{D} \models \varphi[s_{a/x}]$

Note: here we need the domain  $D$

# Semantics of a Query

The semantics of a *query*  $Q = \{ \mathbf{x} \mid \varphi \}$  is

$$Q(\mathbf{D}) = \{ s(\mathbf{x}) \mid \text{for all } s \text{ such that } \mathbf{D} \models \varphi[s] \}$$

# The drinkers-bars-beers example

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

What does this query compute ?

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

# The drinkers-bars-beers example

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

What does this query compute ?

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

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

# The drinkers-bars-beers example

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

What does this query compute ?

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



# The drinkers-bars-beers example

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

What does this query compute ?

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

Find drinkers that frequent only bars that serves some beer they like.

# The drinkers-bars-beers example

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

What does this query compute ?

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

# The drinkers-bars-beers example

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

What does this query compute ?

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

Find drinkers that frequent some bar that serves only beers they like.

# The drinkers-bars-beers example

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

What does this query compute ?

$x: \forall y. \text{Frequents}(x, y) \Rightarrow \forall z. (\text{Serves}(y, z) \Rightarrow \text{Likes}(x, z))$

# The drinkers-bars-beers example

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

What does this query compute ?

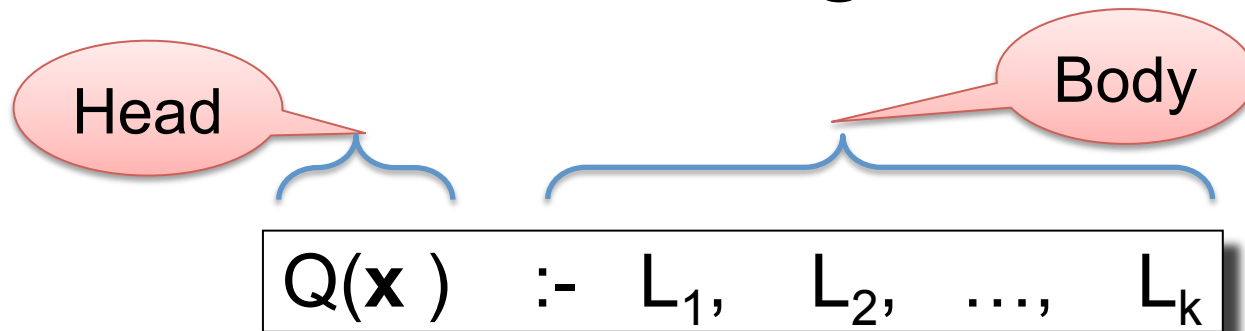
$x: \forall y. \text{Frequents}(x, y) \Rightarrow \forall z. (\text{Serves}(y, z) \Rightarrow \text{Likes}(x, z))$

Find drinkers that frequent only bars that serves only beer they like.

# Summary of Relational Calculus

- Same as First Order Logic
- See book for the two variants:
  - Domain relational calculus (what we discussed)
  - Tuple relational calculus
- This is a powerful, concise language

# Datalog Rule



$L_i$  = a *literal*, or *atom*, or *subgoal*; can be one of:

- $R(x_1, x_2, \dots)$  = relational atom
- $u = v$ , or  $u \neq v$ , where  $u, v$  = variables or constants

$S$  = a new symbol

$\mathbf{x} = (x_1, x_2, \dots)$  are *head variables*

Example:  $Q(x) \text{ :- Likes}(x,y), \text{Likes}(x,z), y \neq z$

# Semantics of Datalog Rule

$$Q(\mathbf{x}) \text{ :- } L_1, L_2, \dots, L_k \quad = \quad Q = \{ \mathbf{x} \mid \exists y_1. \exists y_2. \dots L_1 \wedge L_2 \wedge \dots \wedge L_k \}$$

$y_1, y_2, \dots$  are all non-head variables

The semantics is:

$$Q(\mathbf{D}) \text{ :- } \{ s(\mathbf{x}) \mid s = \text{valuation s.t. } \mathbf{D} \models L_1[s], \dots, \mathbf{D} \models L_k[s] \}$$

Note: a datalog rule is also called a *conjunctive query*



# Non-Recursive Datalog

```
S1(x1) :- body1  
S2(x2) :- body2  
S3(x3) :- body3  
...
```

Each symbol  $S_k$  may appear in  $\text{body}_{k+1}, \text{body}_{k+2}, \dots$ , but may not appear in  $\text{body}_1, \dots, \text{body}_k$ .

Example:

```
A(x,y) :- R(x,z),R(z,y)  
B(x,y) :- A(x,z),A(z,y)  
C(x,y) :- B(x,z),B(z,y)  
Q(x,y) :- C(x,z),C(z,y)
```

What does Q compute ?

# Non-Recursive Datalog

```
S1(x1) :- body1  
S2(x2) :- body2  
S3(x3) :- body3  
...
```

Each symbol  $S_k$  may appear in  $body_{k+1}, body_{k+2}, \dots$ , but may not appear in  $body_1, \dots, body_k$ .

Example:

```
A(x,y) :- R(x,z),R(z,y)  
B(x,y) :- A(x,z),A(z,y)  
C(x,y) :- B(x,z),B(z,y)  
Q(x,y) :- C(x,z),C(z,y)
```

What does Q compute ?

A: all pairs (x,y) connected by a path of length 16.

# Comparison

Is non-recursive datalog equivalent  
to the Relational Calculus ?

# Comparison

Is non-recursive datalog equivalent to the Relational Calculus ?

What about writing this query in non-recursive datalog ?

$$Q = \{x \mid \exists y. \text{Frequents}(x,y) \wedge \text{Serves}(y, \text{'Bud'}) \vee \exists z. \text{Frequents}(x,z) \wedge \text{Serves}(z, \text{'Miller'}) \}$$

# Comparison

Is non-recursive datalog equivalent to the Relational Calculus ?

What about writing this query in non-recursive datalog ?

$$Q = \{x \mid \exists y. \text{Frequents}(x,y) \wedge \text{Serves}(y, \text{'Bud'}) \vee \exists z. \text{Frequents}(x,z) \wedge \text{Serves}(z, \text{'Miller'}) \}$$
$$Q(x) \text{ :- Frequents}(x,y), \text{Serves}(y, \text{'Bud'})$$
$$Q(x) \text{ :- Frequents}(x,z), \text{Serves}(z, \text{'Miller'})$$

# Comparison

Is non-recursive datalog equivalent to the Relational Calculus ?

Now what about writing this query in non-recursive datalog ?

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

# Comparison

Is non-recursive datalog equivalent to the Relational Calculus ?

Now what about writing this query in non-recursive datalog ?

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

Impossible ! Proof on next slides.

# Monotone Queries

- Given two database instances  $\mathbf{D} = (D, R_1, \dots, R_k)$  and  $\mathbf{D}' = (D', R'_1, \dots, R'_k)$  we write  $\mathbf{D} \subseteq \mathbf{D}'$  if  $D \subseteq D', R_1 \subseteq R'_1, \dots, R_k \subseteq R'_k$ .
- In other words,  $\mathbf{D}'$  is obtained by adding tuples to  $\mathbf{D}$ .
- We say that  $Q$  is *monotone* if whenever  $\mathbf{D} \subseteq \mathbf{D}'$ ,  $Q(\mathbf{D}) \subseteq Q(\mathbf{D}')$ .



# Monotone Queries

**Fact.** Every datalog rule is monotone.

Why ?

**Fact.** Every datalog program is monotone.

Why ?

**Fact.** This query is not monotone:

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

Why ?

# Monotone Queries

**Fact.** Every datalog rule is monotone.

Why ?

Recall the semantics:  $Q(\mathbf{D}) = \{ s(x) \mid \mathbf{D} \models L_1[s], \dots, \mathbf{D} \models L_k[s] \}$   
If  $\mathbf{D} \models L_i[s]$  then  $\mathbf{D}' \models L_i[s]$  hence  $Q(\mathbf{D}) \subseteq Q(\mathbf{D}')$ .

**Fact.** Every datalog program is monotone.

Why ?

**Fact.** This query is not monotone:

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

Why ?

# Monotone Queries

**Fact.** Every datalog rule is monotone.

Why ?

Recall the semantics:  $Q(\mathbf{D}) = \{ s(x) \mid \mathbf{D} \models L_1[s], \dots, \mathbf{D} \models L_k[s] \}$   
If  $\mathbf{D} \models L_i[s]$  then  $\mathbf{D}' \models L_i[s]$  hence  $Q(\mathbf{D}) \subseteq Q(\mathbf{D}')$ .

**Fact.** Every datalog program is monotone.

Why ?

By induction on the number of rules

**Fact.** This query is not monotone:

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

Why ?

# Monotone Queries

**Fact.** Every datalog rule is monotone.

Why ?

Recall the semantics:  $Q(\mathbf{D}) = \{ s(x) \mid \mathbf{D} \models L_1[s], \dots, \mathbf{D} \models L_k[s] \}$   
If  $\mathbf{D} \models L_i[s]$  then  $\mathbf{D}' \models L_i[s]$  hence  $Q(\mathbf{D}) \subseteq Q(\mathbf{D}')$ .

**Fact.** Every datalog program is monotone.

Why ?

By induction on the number of rules

**Fact.** This query is not monotone:

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

Why ?

By adding a tuple to *Frequents* we may remove some answer

# Monotone Queries

Recall the relational calculus:

$$\begin{aligned} t ::= & x \mid a \quad /* \text{variable or constant} */ \\ \varphi ::= & R(t_1, \dots, t_k) \mid t_i = t_j \\ & \mid \varphi \wedge \varphi' \mid \varphi \vee \varphi' \mid \neg \varphi' \\ & \mid \forall x. \varphi \mid \exists x. \varphi \end{aligned}$$

What fragment of the relational calculus is monotone ?

# Monotone Queries

Recall the relational calculus:

$$\begin{aligned} t ::= & x \mid a \quad /* \text{variable or constant} */ \\ \varphi ::= & R(t_1, \dots, t_k) \mid t_i = t_j \\ & \mid \varphi \wedge \varphi' \mid \varphi \vee \varphi' \mid \neg \varphi' \\ & \mid \forall x. \varphi \mid \exists x. \varphi \end{aligned}$$

What fragment of the relational calculus is monotone ?

$$\begin{aligned} t ::= & x \mid a \quad /* \text{variable or constant} */ \\ \varphi ::= & R(t_1, \dots, t_k) \mid t_i = t_j \\ & \mid \varphi \wedge \varphi' \mid \varphi \vee \varphi' \mid \exists x. \varphi \end{aligned}$$

# Non-Recursive Datalog<sup>-</sup>

$$Q(\mathbf{x}) \text{ :- } L_1, L_2, \dots, L_k$$

$L_i$  = may also be  $\neg R(x_1, x_2, \dots)$

- $R(x_1, x_2, \dots)$  = relational atom
- $u = v$ , or  $u \neq v$ , where  $u, v$  = variables or constants

# Non-Recursive Datalog<sup>+</sup>

Example:

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



# Non-Recursive Datalog<sup>+</sup>

Example:

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

$$\begin{aligned} \text{FS}(x,y) &:- \text{Frequents}(x,z), \text{Serves}(z,y) \\ Q(x) &:- \text{Likes}(x,y), \neg \text{FS}(x,y) \end{aligned}$$

Rel. Calculus  $\rightarrow$  non-rec. Datalog<sup>-</sup>

**Theorem.** Every query in the Relational Calculus can be translated to non-recursive Datalog<sup>-</sup>

# Rel. Calculus $\rightarrow$ non-rec. Datalog<sup>-</sup>

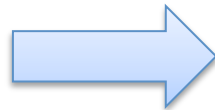
## Proof

First remove  $\forall$  by using De Morgan:  $\forall x.\varphi = \neg\exists x.\neg\varphi$

Then, inductively, for each subformula  $\varphi$   
create a datalog rule  $F(x):-\text{body}$   
that computes  $Q = \{ \mathbf{x} \mid \varphi \}$

# Rel. Calculus $\rightarrow$ non-rec. Datalog<sup>-</sup>

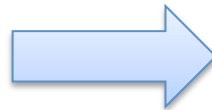
$R(t_1, \dots, t_k)$



$F(\mathbf{x}) \text{ :- } R(t_1, \dots, t_k)$

# Rel. Calculus $\rightarrow$ non-rec. Datalog<sup>-</sup>

$\varphi_1 \wedge \varphi_2$

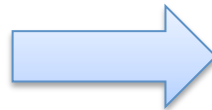


$F_1(\mathbf{x}_1) :- \dots$  /\* for  $\varphi_1$  \*/  
 $F_2(\mathbf{x}_2) :- \dots$  /\* for  $\varphi_2$  \*/

???

# Rel. Calculus $\rightarrow$ non-rec. Datalog<sup>-</sup>

$$\varphi_1 \wedge \varphi_2$$

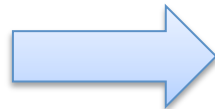


$$F_1(\mathbf{x}_1) \text{ :- } \dots \text{ /* for } \varphi_1 \text{ */}$$
$$F_2(\mathbf{x}_2) \text{ :- } \dots \text{ /* for } \varphi_2 \text{ */}$$

$$F(\mathbf{x}) \text{ :- } F_1(\mathbf{x}_1), F_2(\mathbf{x}_2)$$

# Rel. Calculus $\rightarrow$ non-rec. Datalog<sup>-</sup>

$\varphi_1 \vee \varphi_2$

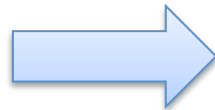


$F_1(\mathbf{x}_1) :- \dots /* \text{ for } \varphi_1 */$   
 $F_2(\mathbf{x}_2) :- \dots /* \text{ for } \varphi_2 */$

???

# Rel. Calculus $\rightarrow$ non-rec. Datalog<sup>-</sup>

$\varphi_1 \vee \varphi_2$



$F_1(\mathbf{x}_1) :- \dots$  /\* for  $\varphi_1$  \*/  
 $F_2(\mathbf{x}_2) :- \dots$  /\* for  $\varphi_2$  \*/

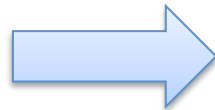
$F(\mathbf{x}) :- F_1(\mathbf{x}_1)$   
 $F(\mathbf{x}) :- F_2(\mathbf{x}_2)$



# Rel. Calculus $\rightarrow$ non-rec. Datalog $^-$

$F_1(\mathbf{x}) :- \dots /* \text{for } \varphi_1 */$

$\exists x. \varphi_1$

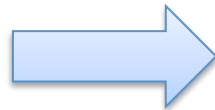


???

# Rel. Calculus $\rightarrow$ non-rec. Datalog<sup>-</sup>

$F_1(\mathbf{x})$  :- ... /\* for  $\varphi_1$  \*/

$\exists \mathbf{x}.\varphi_1$



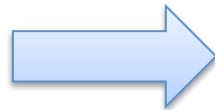
$F(\mathbf{y})$  :-  $F_1(\mathbf{x})$

where  $\mathbf{y} = \mathbf{x} - \{x\}$

# Rel. Calculus $\rightarrow$ non-rec. Datalog $^-$

$F_1(\mathbf{x})$  :- ... /\* for  $\varphi_1$  \*/

$\neg\varphi_1$



$F(\mathbf{x})$  :-  $\neg F_1(\mathbf{x})$

**End of Proof**

# Unsafe Queries

$Q(x) :- x=y$

What's wrong ?

$Q(x) :- \neg R(x,y)$

# Unsafe Queries

**Definition.** A datalog<sup>-</sup> rule  $S :- L_1, \dots, L_k$  is called safe, if every variable  $x$  appears in at least one positive relational atom.

Examples

$Q(x) :- T(x), S(y), x=y$

$Q(x) :- T(x), S(y), \neg R(x,y)$

We require that all rules of a non-recursive datalog<sup>-</sup> program to be safe.

# Unsafe Queries

But what about Relational Calculus ?

What does it mean for a query Q to be safe ?

Are these queries safe ?

$$Q = \{x \mid \forall y. \text{Frequents}(x,y)\}$$

$$Q = \{(x,y) \mid \text{Faculty}(x) \vee \text{Student}(y)\}$$

# Unsafe Queries

**Definition.** A query  $Q$  is domain independent, if for any two database instances  $\mathbf{D}$ ,  $\mathbf{D}'$ , such that they have the same relations  $R_1, \dots, R_k$  but possibly different domains  $D, D'$ ,  $Q(\mathbf{D}) = Q(\mathbf{D}')$ .

A safe datalog<sup>-</sup> rule is domain independent (obvious).  
Does the converse hold ?

Note: domain independent queries are also called safe queries, somewhat confusingly.

# Unsafe Queries

Explain why these queries  
are not domain independent:

$$Q = \{x \mid \forall y. \text{Frequents}(x,y)\}$$

$$Q = \{(x,y) \mid \text{Faculty}(x) \vee \text{Student}(y)\}$$



# Unsafe Queries

**Theorem.** The problem: “*given a relational query  $Q$  decide whether  $Q$  is safe*” is undecidable.

Bummer !...

Why doesn't the following work ?

Translate  $Q$  to a datalog<sup>-</sup> program,  
then check that each rule is safe.

Is this a decision procedure  
for domain independence ?

# Active Domain Semantics

**Definition.** The *active domain*  $ADom$  of a database instance  $\mathbf{D} = (D, R_1^D, \dots, R_k^D)$ , is the set of all constants in  $R_1^D, \dots, R_k^D$

The active domain  
can be computed  
by a boring  
relational query:

```
ADom(x) :- Frequents(x,y)
ADom(y) :- Frequents(x,y)
ADom(x) :- Likes(x,y)
ADom(y) :- Likes(x,y)
ADom(x) :- Serves(x,y)
ADom(y) :- Serves(x,y)
```

# Active Domain Semantics

Make the following two changes to standard semantics:

$$\mathbf{D} \models (\forall x.\varphi) [s]$$

If for every constant  $a$  in  $A\text{Dom}$ ,  
 $\mathbf{D} \models \varphi[s_{a/x}]$

$$\mathbf{D} \models (\exists x.\varphi) [s]$$

If for some constant  $a$  in  $A\text{Dom}$ ,  
 $\mathbf{D} \models \varphi[s_{a/x}]$

# Active Domain Semantics

**Theorem.** If  $Q$  is domain-independent, then its standard semantics coincides with the active domain semantics

Why ?

# Active Domain Semantics

**Theorem.** If  $Q$  is domain-independent, then its standard semantics coincides with the active domain semantics

Why ?

Let  $\mathbf{D} = (D, R_1^D, \dots, R_k^D)$ , be a database instance.

Denote  $\mathbf{D}' = (ADom, R_1^D, \dots, R_k^D)$

Since  $Q$  is d.i. we have  $Q(\mathbf{D}) = Q(\mathbf{D}')$

$Q(\mathbf{D}')$  is the same as the active-domain semantics on  $\mathbf{D}$

# Rel. Calculus $\rightarrow$ non-rec. Datalog $^-$

If  $Q$  is *domain independent*,  
then we want to translate it into  
a *safe* non-rec. datalog $^-$  program

$$\boxed{\neg \varphi_1} \quad \rightarrow \quad \boxed{F(\mathbf{x}) \text{ :- } \neg F_1(\mathbf{x})}$$

This rule is unsafe !  
How do we make it safe ?

# Rel. Calculus $\rightarrow$ non-rec. Datalog $^-$

If  $Q$  is *domain independent*,  
then we want to translate it into  
a *safe* non-rec. datalog $^-$  program



$F(\mathbf{x}) :- \text{Adom}(x_1), \text{Adom}(x_2), \dots, \text{Adom}(x_k), \neg F_1(\mathbf{x})$

Where  $\mathbf{x} = (x_1, x_2, \dots, x_k)$

# Discussion

- Non-recursive datalog<sup>-</sup> is a simple, friendly, very popular notation for queries
- Naturally compositional: sequence of rules
- Exposes a hierarchy of queries
  - Single rule datalog = conjunctive queries
  - Non-rec. datalog program = monotone-RC
  - Non-rec. datalog<sup>-</sup> program = RC
- No quantifiers ! Because all are  $\exists$



# Why We Care

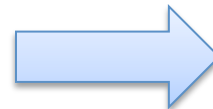
SQL is non-recursive datalog<sup>-</sup> !

$Q(\mathbf{x}) :- R(..), S(..), T(..)$



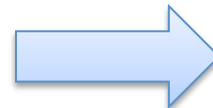
```
SELECT x
FROM R, S, T
WHERE ...
```

$Q1(\mathbf{x}) :- \dots$   
 $Q2(\mathbf{y}) :- \dots$   
 $Q(\mathbf{z}) :- Q1(\mathbf{x}), Q2(\mathbf{y})$



Subqueries in the  
FROM clause

$Q(\mathbf{x}) :- R(..), S(..), \neg T(..)$



Subquery in the  
WHERE clause

# Why We Care

- When we use a database system we need to write some complex queries
  - *The Masochists*: Find drinkers that frequent only bars that server no beer that they like
- Write it first in Relational Calculus
  - Easy exercise
- Translate it to Non-recursive datalog<sup>-</sup>
  - This is subtle, but we have seen it in detail
- Translate Non-recursive datalog<sup>-</sup> to SQL
  - Easy exercise

# Relational Algebra

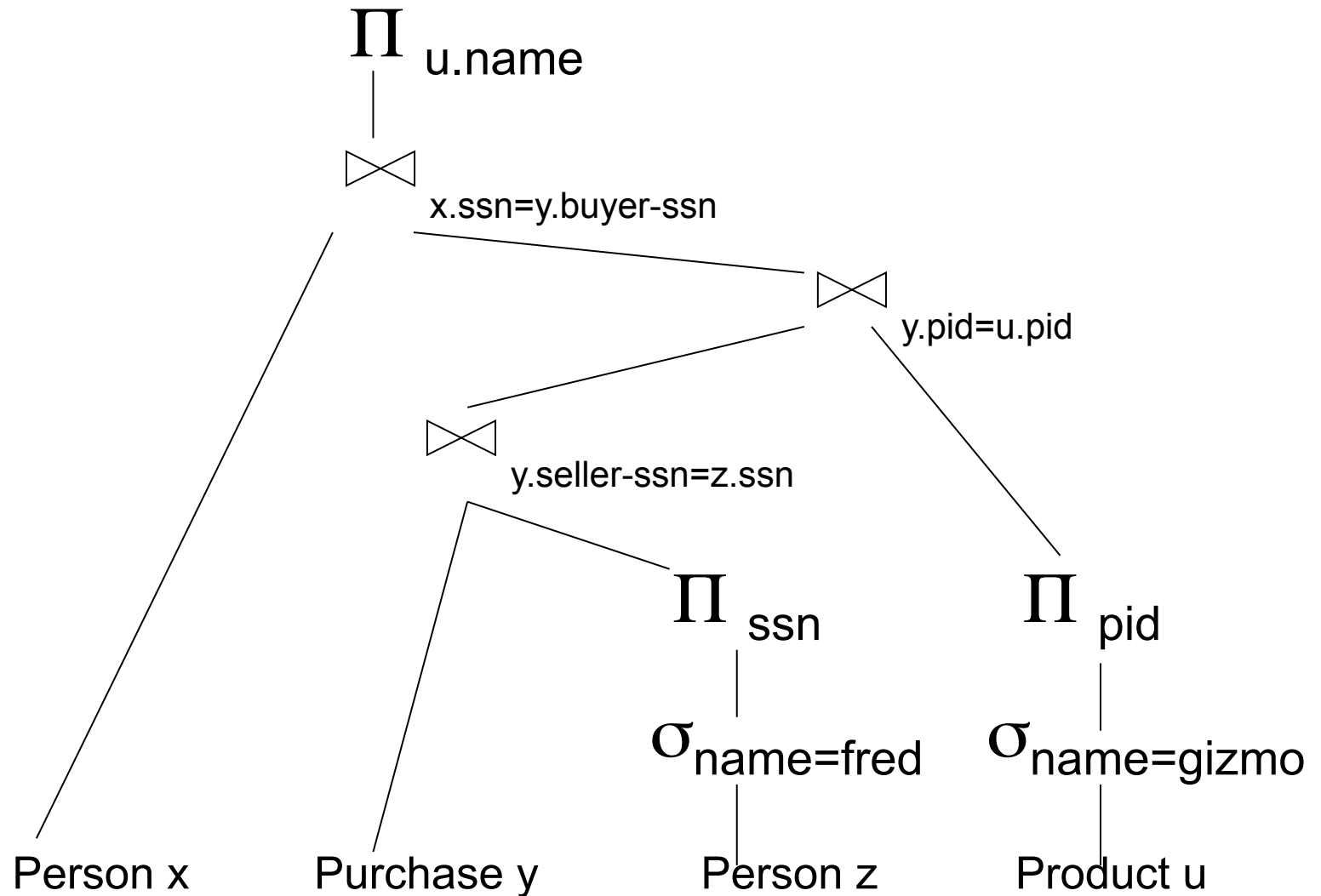
The Basic Five operators:

- Union:  $\cup$
- Difference:  $-$
- Selection:  $\sigma$
- Projection:  $\Pi$
- Join:  $\bowtie$

Plus: cartesian product  $\times$ , renaming  $\rho$

We will study RA in detail when we discuss Query Execution

# Complex RA Expressions



# Equivalence Theorem

**Theorem.** Every domain independent query in RC can be translated to safe non-recursive Datalog<sup>-</sup>

Done that.

**Theorem.** Every safe non-recursive Datalog<sup>-</sup> can be translated to Relational Algebra.

Easy exercise – to do at home

**Theorem.** Relational Algebra expressions can be translated to a domain independent RC query

Easy exercise – to do at home

# Discussion

- What is the monotone fragment of RA ?
  -
- What are the safe queries in RA ?
  -
- Where do we use RA (applications) ?
  - 
  -

# Discussion

- What is the monotone fragment of RA ?
  - No difference:  $\cup$ ,  $\sigma$ ,  $\Pi$ ,  $\bowtie$
- What are the safe queries in RA ?
  - All RA queries are safe
- Where do we use RA (applications) ?
  - Translating SQL (WHAT  $\rightarrow$  HOW)
  - Some new query languages (Pig-Latin)