

# CSE 544

# Theory of Query Languages

Tuesday, February 22<sup>nd</sup>, 2011

# Outline

- Conjunctive queries; containment
- Datalog
- Query complexity

# Conjunctive Queries

- A subset of Relational Calculus (=FO)
- Correspond to  
SELECT-DISTINCT-FROM-WHERE
- Most queries in practice are conjunctive
- Some optimizers handle only conjunctive queries  
- break larger queries into many CQs
- CQ' s have more positive theoretical properties  
than arbitrary queries

# Conjunctive Queries

- **Definition** A conjunctive query is defined by:

$\varphi ::= R(t_1, \dots, t_k) \mid t_i = t_j \mid \varphi \wedge \varphi' \mid \exists x. \varphi$

- missing are  $\forall, \vee, \neg$

# Conjunctive Queries, CQ

- Example of CQ

$$q(x,y) = \exists z.(R(x,z) \wedge \exists u.(R(z,u) \wedge R(u,y)))$$

$$q(x) = \exists z.\exists u.(R(x,z) \wedge R(z,u) \wedge R(u,y))$$

- Examples of non-CQ:

$$q(x,y) = \forall z.(R(x,z) \rightarrow R(y,z))$$

$$q(x) = T(x) \vee \exists z.S(x,z)$$

# Conjunctive Queries

- Any CQ query can be written as:

$$q(x_1, \dots, x_n) = \exists y_1. \exists y_2 \dots \exists y_p. (R_1(t_{11}, \dots, t_{1m}) \wedge \dots \wedge R_k(t_{k1}, \dots, t_{km}))$$

(i.e. all quantifiers are at the beginning)

- Same in **Datalog** notation:

Datalog rule

$$q(x_1, \dots, x_n) \text{ :- } R_1(t_{11}, \dots, t_{1m}), \dots, R_k(t_{k1}, \dots, t_{km})$$

head

body

# Examples

Employee(x), ManagedBy(x,y), Manager(y)

- Find all employees having the same manager as “Smith”:

$A(x) :- \text{ManagedBy}(\text{“Smith”}, y), \text{ManagedBy}(x, y)$

# Examples

Employee(x), ManagedBy(x,y), Manager(y)

- Find all employees having the same director as Smith:

$A(x) :- \text{ManagedBy}(\text{"Smith"}, y), \text{ManagedBy}(y, z), \text{ManagedBy}(x, u), \text{ManagedBy}(u, z)$

CQs are useful in practice



# CQ and SQL

CQ:

$A(x) :- \text{ManagedBy}(\text{"Smith"}, y), \text{ManagedBy}(x, y)$

SQL:

Notice  
"distinct"

```
select distinct m2.name  
from ManagedBy m1, ManagedBy m2  
where m1.name="Smith" AND  
      m1.manager=m2.manager
```

# CQ and SQL

- Are CQ queries precisely the `SELECT-DISTINCT-FROM-WHERE` queries ?

# CQ and SQL

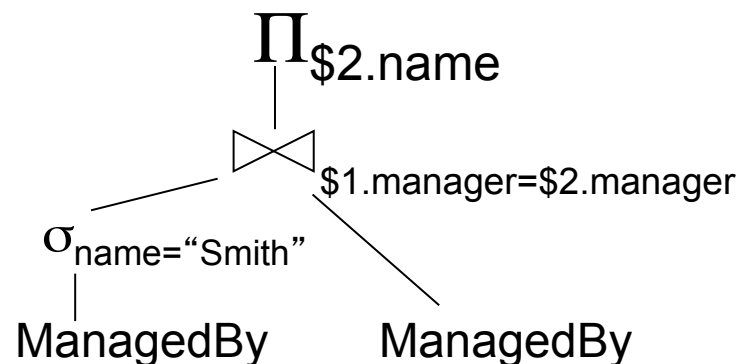
- Are CQ queries precisely the SELECT-DISTINCT-FROM-WHERE queries ?
- No: CQ queries do not allow  $<$ ,  $\leq$ ,  $\neq$
- But we can extend CQ with inequality predicates, and usually write the extended language as  $CQ^<$ , etc

# CQ and RA

## Relational Algebra:

- CQ correspond precisely to  $\sigma_C$ ,  $\Pi_A$ ,  $\times$  (missing:  $\cup$ ,  $-$ ) and where C has only =

$A(x) :- \text{ManagedBy}(\text{"Smith"}, y), \text{ManagedBy}(x, y)$



# Extensions of CQ

CQ<sup>≠</sup>

Find managers that manage at least 2 employees

$A(y) :- \text{ManagedBy}(x,y), \text{ManagedBy}(z,y), x \neq z$

# Extensions of CQ

CQ<sup><</sup>

Find employees earning more than their manager:

$A(y) :- \text{ManagedBy}(x,y), \text{Salary}(x,u), \text{Salary}(y,v), u > v$

# Extensions of CQ

CQ<sup>¬</sup>: negation applied only to one atom

Find people sharing the same office with Alice, but not the same manager:

```
A(y) :- Office("Alice",u), Office(y,u),  
        ManagedBy("Alice",x), ¬ManagedBy(x,y)
```

# Extensions of CQ

UCQ     Union of conjunctive queries

Datalog:

```
A(name) :- Employee(name, dept, age, salary), age > 50  
A(name) :- RetiredEmployee(name, address)
```

Datalog notation is very convenient  
for expressing unions (no need for  $\vee$  )



# Summary of Extensions of CQ

- CQ
  - CQ $\neq$
  - CQ $<$
  - UCQ
  - CQ $^{-1}$
- 
- Which of these classes contain only monotone queries ?

# Query Equivalence and Containment

- Justified by optimization needs
- Intensively studied since 1977

# Query Equivalence

- Queries  $q_1$  and  $q_2$  are **equivalent** if for every database  $\mathbf{D}$ ,  $q_1(\mathbf{D}) = q_2(\mathbf{D})$ .
- Notation:  $q_1 \equiv q_2$

# Query Containment

- Query  $q_1$  is **contained** in  $q_2$  if for every database  $\mathbf{D}$ ,  $q_1(\mathbf{D}) \subseteq q_2(\mathbf{D})$ .
- $q_1$  and  $q_2$  are **equivalent** if for every database  $\mathbf{D}$ ,  $q_1(\mathbf{D}) = q_2(\mathbf{D})$
- Notation:  $q_1 \subseteq q_2$ ,  $q_1 \equiv q_2$
- Obviously:  $q_1 \subseteq q_2$  and  $q_2 \subseteq q_1$  iff  $q_1 \equiv q_2$
- Conversely:  $q_1 \wedge q_2 \equiv q_1$  iff  $q_1 \subseteq q_2$

**We will study the containment problem only.**

# Examples of Query Containments

Is  $q_1 \subseteq q_2$  ?

$q_1(x) :- R(x,u), R(u,v), R(v,w)$   
 $q_2(x) :- R(x,u), R(u,v)$

# Examples of Query Containments

Is  $q_1 \subseteq q_2$  ?

$q_1(x) :- R(x,u), R(u,v), R(v,x)$   
 $q_2(x) :- R(x,u), R(u,x)$

# Examples of Query Containments

Is  $q_1 \subseteq q_2$  ?

$q_1(x) :- R(x,u), R(u,u)$

$q_2(x) :- R(x,u), R(u,v), R(v,w)$

# Examples of Query Containments

Is  $q_1 \subseteq q_2$  ?

$q_1(x) :- R(x,u), R(u, \text{"Smith"})$   
 $q_2(x) :- R(x,u), R(u,v)$



# Query Containment

- **Theorem** Query containment for FO is undecidable
- **Theorem** Query containment for CQ is decidable and NP-complete.

# Trakhtenbrot's Theorem

**Definition** A sentence  $\varphi$ , is called *finitely satisfiable* if there exists a finite database instance  $D$  s.t.  $D \models \varphi$

Satisfiable:

$\exists x. \exists y. \forall z. (R(x,z) \rightarrow R(y,z))$   
 $\exists x. \exists y. T(x) \vee \exists z. S(x,z)$

Unsatisfiable:

$\forall x. \forall y. \forall z. (R(x,y) \wedge R(x,z) \rightarrow y=z)$   
 $\wedge \exists y. \forall x. \text{not } R(x,y)$

**Theorem** The following problem is undecidable:  
Given FO sentence  $\varphi$ , check if  $\varphi$  is finitely satisfiable

# Query Containment

- **Theorem** Query containment for FO is undecidable
- **Proof:** By reduction from the finite satisfiability problem:
- Given a sentence  $\varphi$ , define two queries:  
 $q_1(x) = R(x)$ , and  $q_2(x) = R(x) \wedge \varphi$
- Then  $q_1 \subseteq q_2$  iff  $\varphi$  is not finitely satisfiable

# Query Containment Algorithm

How to check  $q_1 \subseteq q_2$

- **Canonical database** for  $q_1$  is:  
$$\mathbf{D}_{q_1} = (D, R_1^D, \dots, R_k^D)$$
  - $D$  = all variables and constants in  $q_1$
  - $R_1^D, \dots, R_k^D$  = the body of  $q_1$
- **Canonical tuple** for  $q_1$  is:  
 $t_{q_1}$  (the head of  $q_1$ )

# Examples of Canonical Databases

$q_1(x,y) :- R(x,u),R(v,u),R(v,y)$

- Canonical database:  $\mathbf{D}_{q_1} = (D, R^D)$

–  $D = \{x,y,u,v\}$

–  $R^D =$

x	u
v	u
v	y

- Canonical tuple:  $t_{q_1} = (x,y)$

# Examples of Canonical Databases

$q_1(x) :- R(x,u), R(u, \text{"Smith"}), R(u, \text{"Fred"}), R(u, u)$

- $D_{q_1} = (D, R)$ 
  - $D = \{x, u, \text{"Smith"}, \text{"Fred"}\}$
  - $R =$

x	u
u	"Smith"
u	"Fred"
u	u

- $t_{q_1} = (x)$

# Checking Containment

**Theorem:**  $q_1 \subseteq q_2$  iff  $t_{q_1} \in q_2(\mathbf{D}_{q_1})$ .

Example:

$q_1(x,y) :- R(x,u),R(v,u),R(v,y)$

$q_2(x,y) :- R(x,u),R(v,u),R(v,w),R(t,w),R(t,y)$

- $D=\{x,y,u,v\}$

- $R =$   $t_{q_1} = (x,y)$

- Yes,  $q_1 \subseteq q_2$

x	u
v	u
v	y

# Query Homomorphisms

- A **homomorphism**  $f : q_2 \rightarrow q_1$  is a function  $f: \text{var}(q_2) \rightarrow \text{var}(q_1) \cup \text{const}(q_1)$  such that:
  - $f(\text{body}(q_2)) \subseteq \text{body}(q_1)$
  - $f(t_{q_1}) = t_{q_2}$

**The Homomorphism Theorem**  $q_1 \subseteq q_2$  iff there exists a homomorphism  $f : q_2 \rightarrow q_1$



# Example of Query Homeomorphism

$$\text{var}(q_1) = \{x, u, v, y\}$$

$$\text{var}(q_2) = \{x, u, v, w, t, y\}$$

$$q_1(x, y) \text{ :- } R(x, u), R(v, u), R(v, y)$$

$$q_2(x, y) \text{ :- } R(x, u), R(v, u), R(v, w), R(t, w), R(t, y)$$

Therefore  $q_1 \subseteq q_2$

# Example of Query Homeomorphism

$$\text{var}(q_1) \cup \text{const}(q_1) = \{x, u, \text{"Smith"}\}$$

$$\text{var}(q_2) = \{x, u, v, w\}$$

$$q_1(x) \text{ :- } R(x, u), R(u, \text{"Smith"}), R(u, \text{"Fred"}), R(u, u)$$

$$q_2(x) \text{ :- } R(x, u), R(u, v), R(u, \text{"Smith"}), R(w, u)$$

Therefore  $q_1 \subseteq q_2$

# The Complexity

**Theorem** Checking containment of two CQ queries is NP-complete

# Containment for extensions of CQ

- CQ -- NP complete
- CQ<sup>≠</sup> -- ??
- CQ<sup><</sup> -- ??
- UCQ -- ??
- CQ<sup>∩</sup> -- ??
- Relational Calculus -- undecidable

# Query Containment for UCQ

$$q_1 \cup q_2 \cup q_3 \cup \dots \subseteq q_1' \cup q_2' \cup q_3' \cup \dots$$

Notice:  $q_1 \cup q_2 \cup q_3 \cup \dots \subseteq q$  iff  
 $q_1 \subseteq q$  and  $q_2 \subseteq q$  and  $q_3 \subseteq q$  and ....

**Theorem**  $q \subseteq q_1' \cup q_2' \cup q_3' \cup \dots$  iff there exists some  $k$  such that  $q \subseteq q_k'$

It follows that containment for UCQ is decidable, NP-complete.

# Query Containment for $CQ^<$

$q_1() :- R(x,y), R(y,x)$
$q_2() :- R(x,y), x \leq y$

$q_1 \subseteq q_2$  although there is no homomorphism !

To check containment do this:

- Consider all possible orderings of variables in  $q_1$
- For each of them check containment of  $q_1$  in  $q_2$
- If all hold, then  $q_1 \subseteq q_2$

Still decidable, but harder than NP: now in  $\Pi^P_2$

# Query Minimization

**Definition** A conjunctive query  $q$  is minimal if for every other conjunctive query  $q'$  s.t.  $q \equiv q'$ ,  $q'$  has at least as many predicates ( 'subgoals' ) as  $q$

Are these queries minimal ?

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

$q(x) :- R(x,y), R(y,z), R(x, 'Alice' )$

# Query Minimization

- Query minimization algorithm

Choose a subgoal  $g$  of  $q$

Remove  $g$ : let  $q'$  be the new query

We already know  $q \subseteq q'$  (why ?)

If  $q' \subseteq q$  then permanently remove  $g$

- Notice: the order in which we inspect subgoals doesn't matter



# Query Minimization In Practice

- No database system today performs minimization !!!
- Reason:
  - It's hard (NP-complete)
  - Users don't write non-minimal queries
- However, non-minimal queries arise when using views intensively

# Query Minimization for Views

```
CREATE VIEW HappyBoaters

SELECT DISTINCT E1.name, E1.manager
FROM Employee E1, Employee E2
WHERE E1.manager = E2.name
      and E1.boater= 'YES'
      and E2.boater= 'YES'
```

This query is minimal

# Query Minimization for Views

Now compute the Very-Happy-Boaters

```
SELECT DISTINCT H1.name  
FROM HappyBoaters H1, HappyBoaters H2  
WHERE H1.manager = H2.name
```

This query is also minimal

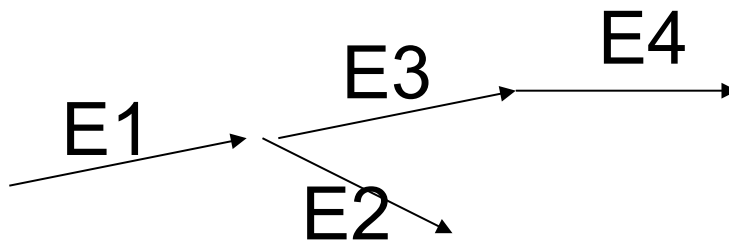
What happens in SQL when we run a query on a view ?

# Query Minimization for Views

## View Expansion

```
SELECT DISTINCT E1.name  
FROM Employee E1, Employee E2, Employee E3, Employee E4  
WHERE E1.manager = E2.name and E1.boater = 'YES' and E2.boater = 'YES'  
and E3.manager = E4.name and E3.boater = 'YES' and E4.boater = 'YES'  
and E1.manager = E3.name
```

This query is no longer minimal !



E2 is redundant

# Expressive Power of FO

- The following queries cannot be expressed in Relational Calculus (FO):
- Transitive closure:
  - $\forall x. \forall y.$  there exists  $x_1, \dots, x_n$  s.t.  
 $R(x, x_1) \wedge R(x_1, x_2) \wedge \dots \wedge R(x_{n-1}, x_n) \wedge R(x_n, y)$
- Parity: the number of edges in R is even

# Datalog

- Adds recursion, so we can compute transitive closure
- A datalog program (query) consists of several datalog rules:

$P_1(t_1) :- \text{body}_1$

$P_2(t_2) :- \text{body}_2$

$\vdots$

$P_n(t_n) :- \text{body}_n$

# Datalog

## Terminology:

- EDB = extensional database predicates
  - The database predicates
- IDB = intentional database predicates
  - The new predicates constructed by the program

# Datalog

Employee(x), ManagedBy(x,y), Manager(y)

All higher level managers that are employees:

EDBs

```
HMngr(x) :- Manager(x), ManagedBy(y,x), ManagedBy(z,y)
Answer(x) :- HMngr(x), Employee(x)
```

IDBs



# Datalog

Employee(x), ManagedBy(x,y), Manager(y)

All persons:

Person(x) :- Manager(x)  
Person(x) :- Employee(x)

Manager  $\cup$  Employee

# Unfolding non-recursive rules

Graph:  $R(x,y)$

$$\begin{aligned} P(x,y) &:- R(x,u), R(u,v), R(v,y) \\ A(x,y) &:- P(x,u), P(u,y) \end{aligned}$$

Can “unfold” it into:

$$A(x,y) :- R(x,u), R(u,v), R(v,w), R(w,m), R(m,n), R(n,y)$$

# Unfolding non-recursive rules

Graph:  $R(x,y)$

$$\begin{aligned} P(x,y) &:- R(x,y) \\ P(x,y) &:- R(x,u), R(u,y) \\ A(x,y) &:- P(x,y) \end{aligned}$$

Now the unfolding has a union:

$$A(x,y) :- R(x,y) \vee \exists u(R(x,u) \wedge R(u,y))$$

Non-recursive datalog = UCQ (why ?)

# Recursion in Datalog

## Example 1: transitive closure

Graph:  $R(x,y)$

Transitive closure:

$$\begin{aligned} P(x,y) &:- R(x,y) \\ P(x,y) &:- P(x,u), R(u,y) \end{aligned}$$

Transitive closure:

$$\begin{aligned} P(x,y) &:- R(x,y) \\ P(x,y) &:- P(x,u), P(u,y) \end{aligned}$$

# Recursion in Datalog

**Example 2:** same generation

Graph:  $R(x,y)$

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

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

$Sg(x,y) :- R(x,u), Sg(u,v), R(v,y)$

# Recursion in Datalog

**Example 3:** value of a Boolean circuit  
with And/Not gates

Graph:  $\text{And}(y,z,x)$ ,  $\text{Not}(y,x)$ ,  $\text{Value0}(x)$ ,  $\text{Value1}(x)$

```
isZero(x) :- Value0(x)
isOne(x)  :- Value1(x)
isZero(x) :- And(y,z,x),isZero(y)
isZero(x) :- And(y,z,x),isZero(z)
isOne(x)  :- And(y,z,x),isOne(y),isOne(z)
isZero(x) :- Not(y,x),isOne(y)
isOne(x)  :- Not(y,x),isZero(x)
```

# Semantics of a Datalog Program

- Let:
  - $R$  = the EDB predicates = some input instance  $D$
  - $S$  = the IDB predicates
- A datalog program  $P$  maps IDB predicate instances  $S$  to new IDB predicate instances  $S'$ :  $S' = P_D(S)$
- The function  $P_D$  is monotone (why ?)
- By Knaster-Tarski's fixpoint Theorem,  $P_D$  has a least fixpoint
- Definition: the meaning of  $P_D$  is the least fixpoint
- Alternatively: compute the meaning of  $P_D$  as follows:
  - $S_0 = \text{emptyset}$ ;  $S_{k+1} = P(S_k)$
  - The meaning is:  $P_D = S_0 \cup S_1 \cup S_2 \cup \dots S_n \cup \dots$

# Evaluation of Datalog Programs

- Naïve evaluation algorithm:
  - Start with  $S = \text{emptyset}$
  - Evaluate  $P$  on the EDB, and on the current IDB  $S$ ;  $\text{new-}S = P(S)$ ; note:  $S \subseteq \text{new-}S$  (why ?)
  - Replace  $S$  with  $\text{new-}S$ ; repeat;
  - Stop when  $S = \text{new-}S$
  - What is the complexity ?
- Semi-naïve evaluation algorithm:
  - Keep track of  $\Delta S = \text{new-}S - S$
  - Improve somewhat the computation of  $P(S)$

What is the complexity of the naïve/ semi-naïve algorithm ?



# Extensions of datalog with negation

- Problems with negation:

$S(x) \text{ :- } R(x), \text{ not } T(x)$

$T(x) \text{ :- } R(x), \text{ not } S(x)$

- There is no minimal fixpoint

# Extensions of datalog with negation

- **Solution 1: Stratified Datalog<sup>-</sup>**
  - Insist that the program be stratified: rules are partitioned into strata, and an IDB predicate that occurs only in strata  $\leq k$  may be negated in strata  $\geq k+1$
- **Solution 2: Inflationary-fixpoint Datalog<sup>-</sup>**
  - Run the naïve algorithm substituting:  $S = S \cup \text{new-S}$
  - Always terminates (why ?)
- **Solution 3: Partial-fixpoint Datalog<sup>-,\*</sup>**
  - Run the naïve algorithm, substituting  $S = \text{new-S}$
  - May not terminate (but we can detect that – how ?)

# Datalog<sup>-</sup>

**Example:** complement transitive closure

Graph:  $R(x,y)$

```
Tc(x,y) :- R(x,y)
Tc(x,y) :- Tc(x,u), R(u,y)
CTc(x,y) :- R(x,u), R(v,y), not Tc(x,y)
```

This is stratified datalog<sup>-</sup>

Challenge: express CTc in inflationary datalog<sup>-</sup>

# Expressive Power

**Theorem:**

stratified-datalog<sup>-</sup>  $\subsetneq$   
inflationary-datalog<sup>-</sup>  $\subsetneq$   
partial-datalog<sup>-</sup>

# Variants of Datalog

without recursion

with recursion

without  $\neg$

Non-recursive Datalog  
= UCQ

**Datalog**

with  $\neg$

Non-recursive Datalog $\neg$   
= FO

Datalog $\neg$   
(three variants)

# Non-recursive Datalog

- Union of Conjunctive Queries = UCQ
  - Containment is decidable, and NP-complete
- Non-recursive Datalog
  - Is equivalent to UCQ
  - Hence containment is decidable here too
  - Is it still NP-complete ?

# Non-recursive Datalog

- A non-recursive datalog:

$$\begin{aligned} T_1(x,y) & :- R(x,u), R(u,y) \\ T_2(x,y) & :- T_1(x,u), T_1(u,y) \\ & \dots \\ T_n(x,y) & :- T_{n-1}(x,u), T_{n-1}(u,y) \\ \text{Answer}(x,y) & :- T_n(x,y) \end{aligned}$$

- Its unfolding as a CQ:

$$\text{Anser}(x,y) \quad :- \quad R(x,u_1), R(u_1, u_2), R(u_2, u_3), \dots R(u_m, y)$$

- How big is this query ?

# Non-recursive Datalog

- A non-recursive datalog:

$$\begin{aligned} T_1(x,y) &:- R(x,u), R(u,y) \\ T_2(x,y) &:- T_1(x,u), T_1(u,y) \\ &\vdots \\ T_n(x,y) &:- T_{n-1}(x,u), T_{n-1}(u,y) \\ \text{Answer}(x,y) &:- T_n(x,y) \end{aligned}$$

- Its unfolding as a CQ:

$$\text{Anser}(x,y) \quad :- \quad R(x,u_1), R(u_1, u_2), R(u_2, u_3), \dots R(u_m, y)$$

- How big is this query ?

Although non-recursive-datalog = UCQ,  
the former is exponentially more concise



# Query Complexity

- Given a query  $\varphi$  in FO
- And given a model  $\mathbf{D} = (D, R_1^D, \dots, R_k^D)$
- What is the complexity of computing the answer  $\varphi(\mathbf{D})$

# Query Complexity

Vardi's taxonomy:

## **Data Complexity:**

- Fix  $\varphi$ . Compute  $\varphi(D)$  as a function of  $|D|$

## **Query Complexity:**

- Fix  $D$ . Compute  $\varphi(D)$  as a function of  $|\varphi|$

## **Combined Complexity:**

- Compute  $\varphi(D)$  as a function of  $|D|$  and  $|\varphi|$

Which is the most important in databases ?

# Example

$$\varphi(x) \equiv \exists u.(R(u,x) \wedge \forall y.(\exists v.S(y,v) \Rightarrow \neg R(x,y)))$$

R =

3	8
7	5
0	8
09	7
6	9
7	6
89	8
98	7
4	0

S =

43	4
5	58
8	6
9	79
6	67
4	7
6	8

How do we proceed ?

# General Evaluation Algorithm

**for** every subexpression  $\varphi_i$  of  $\varphi$ , ( $i = 1, \dots, m$ )  
    compute the answer to  $\varphi_i$  as a table  $T_i(x_1, \dots, x_n)$   
**return**  $T_m$

**Theorem.** If  $\varphi$  has  $k$  variables then one can compute  $\varphi(D)$  in time  $O(|\varphi|^*|D|^k)$

Data Complexity =  $O(|D|^k)$  = in PTIME

Query Complexity =  $O(|\varphi|^*c^k)$  = in EXPTIME

# General Evaluation Algorithm

Example:

$$\varphi(x) \equiv \exists u.(R(u,x) \wedge \forall y.(\exists v.S(y,v) \Rightarrow \neg R(x,y)))$$

$$\varphi_1(u,x) \equiv R(u,x)$$

$$\varphi_2(y,v) \equiv S(y,v)$$

$$\varphi_3(x,y) \equiv \neg R(x,y)$$

$$\varphi_4(y) \equiv \exists v.\varphi_2(y,v)$$

$$\varphi_5(x,y) \equiv \varphi_4(y) \Rightarrow \varphi_3(x,y)$$

$$\varphi_6(x) \equiv \forall y.\varphi_5(x,y)$$

$$\varphi_7(u,x) \equiv \varphi_1(u,x) \wedge \varphi_6(x)$$

$$\varphi_8(x) \equiv \exists u.\varphi_7(u,x) \equiv \varphi(x)$$

# Complexity

**Theorem.** If  $\varphi$  has  $k$  variables then one can compute  $\varphi(D)$  in time  $O(|\varphi|^*|D|^k)$

**Remark.** The number of variables matters !

# Paying Attention to Variables

- Compute all chains of length  $m$

$\text{Chain}_m(x,y) \text{ :- } R(x,u_1), R(u_1, u_2), R(u_2, u_3), \dots R(u_{m-1}, y)$

- We used  $m+1$  variables
- Can you rewrite it with fewer variables ?

# Counting Variables

- $FO^k = FO$  restricted to variables  $x_1, \dots, x_k$
- Write  $\text{Chain}_m$  in  $FO^3$ :

$$\text{Chain}_m(x,y) \text{ :- } \exists u.R(x,u) \wedge (\exists x.R(u, x) \wedge (\exists u.R(x,u) \dots \wedge (\exists u. R(u, y) \dots)))$$



# Query Complexity

- Note: it suffices to investigate boolean queries only
  - If non-boolean, do this:

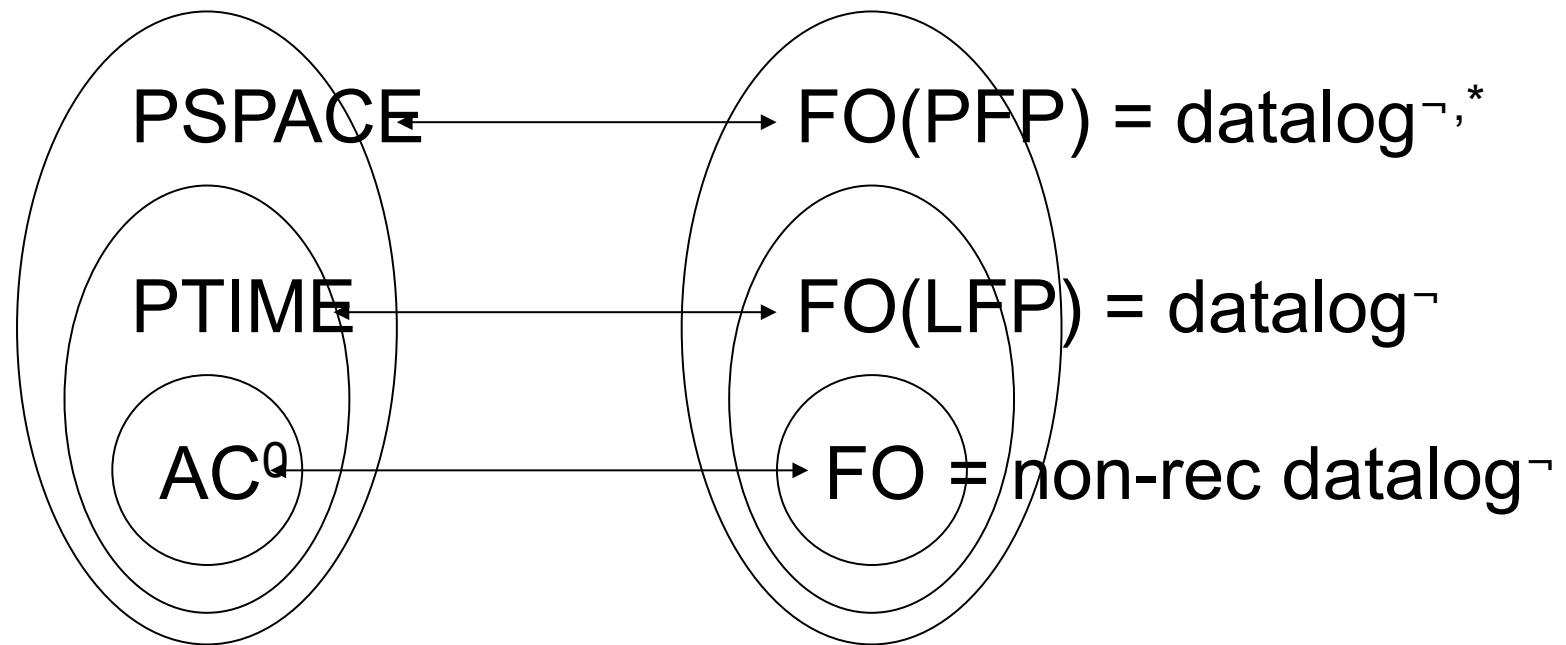
```
for  $a_1$  in  $D$ , ...,  $a_k$  in  $D$   
  if  $(a_1, \dots, a_k)$  in  $\varphi(D)$  /* this is a boolean query */  
    then output  $(a_1, \dots, a_k)$ 
```

# Computational Complexity Classes

Recall computational complexity classes:

- $AC^0$
- LOGSPACE
- NLOGSPACE
- PTIME
- NP
- PSPACE
- EXPTIME
- EXPSPACE
- (Kalmar) Elementary Functions
- Turing Computable functions

# Data Complexity of Query Languages



The more complex a QL, the harder it is to optimize