## Schema Mappings and Data Exchange

Lecture #2

*EASSLC 2012* Southwest University August 2012

## The Relational Data Model (E.F. Codd – 1970)

- The Relational Data Model uses the mathematical concept of a relation as the formalism for describing and representing data.
- Question: What is a relation?
- Answer:
  - □ Formally, a relation is a subset of a cartesian product of sets.
  - □ Informally, a relation is a "table" with rows and columns.

CHECKING	Table
----------	-------

branch-name	account-no	customer-name	balance
Orsay	10991-06284	Abiteboul	\$3,567.53
Hawthorne	10992-35671	Hull	\$11,245.75
			•••

#### Relational Structures vs. Relational Databases

Relational Structure

 $A = (A, R_1, ..., R_m)$ 

- A is the **universe** of **A**
- R<sub>1</sub>,...,R<sub>m</sub> are the relations of A
- Relational Database

$$\mathbf{D} = (\mathsf{R}_1, \dots, \mathsf{R}_m)$$

- Thus, a relational database can be thought of as a relational structure without its universe.
  - □ And this causes some problems down the road ...

## Query Languages for the Relational Data Model

Codd introduced two different query languages for the relational data model:

- Relational Algebra, which is a procedural language.
  - It is an algebraic formalism in which queries are expressed by applying a sequence of operations to relations.
- Relational Calculus, which is a declarative language.
  - It is a logical formalism in which queries are expressed as formulas of first-order logic.

Codd's Theorem: Relational Algebra and Relational Calculus are "essentially equivalent" in terms of expressive power. (but what does this really mean?)

## The Five Basic Operations of Relational Algebra

- Group I: Three standard set-theoretic binary operations:
  - Union
  - Difference
  - Cartesian Product.
- Group II. Two special unary operations on relations:
  - Projection
  - Selection.
- Relational Algebra consists of all expressions obtained by combining these five basic operations in syntactically correct ways.

## **Relational Algebra**

- Definition: A relational algebra expression is a string obtained from relation schemas using union, difference, cartesian product, projection, and selection.
- Context-free grammar for relational algebra expressions:

 $E := R, S, ... | (E_1 \cup E_2) | (E_1 - E_2) | (E_1 \times E_2) | \pi_L(E) | \sigma_{\Theta}(E),$ where

- R, S, ... are relation schemas
- L is a list of attributes
- $\Theta$  is a condition.

Given TEACHES(fac-name,course,term) and ENROLLS(stud-name, course,term): To compute TAUGHT-BY(stud-name,course,term,fac-name)

- 1. ENROLLS  $\times$  TEACHES
- 2.  $\sigma_{\text{T.course} = \text{E.course} \land \text{T.term} = \text{E.term}}$  (ENROLLS × TEACHES)
- 3.  $\pi$  stud-name,E.course,E.term,fac-name ( $\sigma$  T.course = E.course  $\wedge$  T.term = E.term (ENROLLS  $\times$  TEACHES))

The result is ENROLLS  $\bowtie$  TEACHES.

# SQL vs. Relational Algebra

SQL	Relational Algebra
SELECT	Projection $\pi$
FROM	Cartesian Product $\times$
WHERE	Selection $\sigma$

Semantics of SQL via interpretation to Relational Algebra

SELECT 
$$R_{i1}$$
.A1, ...,  $R_{im}$ .A.m  
FROM  $R_1$ , ...,  $R_K$  =  $\pi_{Ri1.A1, ..., Rim.A.m} (\sigma_{\Psi} (R_1 \times ... \times R_K))$   
WHERE  $\Psi$ 

## **Relational Calculus**

- In addition to relational algebra, Codd introduced relational calculus.
- Relational calculus is a declarative database query language based on first-order logic.
- Relational calculus comes into two different flavors:
  - Tuple relational calculus
  - Domain relational calculus.

We will focus on domain relational calculus.

There is an easy translation between these two formalisms.

 Codd's main technical result is that relational algebra and relational calculus have "essentially" the same expressive power.

#### Relational Calculus (First-Order Logic for Databases)

- First-order variables: x, y, z, ...,  $x_1$ , ...,  $x_k$ ,...
  - □ They range over values that may occur in tables.
- Relation symbols: R, S, T, ... of specified arities (names of relations)
- Atomic (Basic) Formulas:
  - R(x<sub>1</sub>,...,x<sub>k</sub>), where R is a k-ary relation symbol (alternatively, (x<sub>1</sub>,...,x<sub>k</sub>) ∈ R; the variables need not be distinct)
  - □ (x op y), where op is one of =,  $\neq$ , <, >, ≤, ≥
  - □ (x op c), where c is a constant and op is one of =,  $\neq$ , <, >, ≤, ≥.
- Relational Calculus Formulas:
  - Every atomic formula is a relational calculus formula.
  - $\square$  If  $\phi$  and  $\psi$  are relational calculus formulas, then so are:
    - $(\phi \land \psi)$ ,  $(\phi \lor \psi)$ ,  $\neg \psi$ ,  $(\phi \rightarrow \psi)$  (propositional connectives)
    - $(\exists x \phi)$  (existential quantification)
    - ( $\forall \mathbf{x} \mathbf{\phi}$ ) (universal quantification).

#### Relational Calculus as a Database Query Language

Definition:

A relational calculus expression is an expression of the form

{  $(x_1,...,x_k): \phi(x_1,...,x_k)$  },

where  $\varphi(x_1,...,x_k)$  is a relational calculus formula with  $x_1,...,x_k$  as its free variables.

When applied to a relational database I, this relational calculus expression returns the k-ary relation that consists of all k-tuples (a<sub>1</sub>,...,a<sub>k</sub>) that make the formula "true" on I.

Example: The relational calculus expression  $\{(x,y): \exists z(E(x,z) \land E(z,y))\}$ returns the set P of all pairs of nodes (a,b) that are connected via a path of length 2. Natural Join in Relational Calculus

Given TEACHES(fac-name,course,term) and ENROLLS(stud-name, course,term):

Compute TAUGHT-BY(stud-name,course,term,fac-name)

In relational algebra:

<sup>π</sup> stud-name,E.course,E.term,fac-name ( $\sigma_{\text{T.course} = \text{E.course} \land \text{T.term} = \text{E.term}$  (ENROLLS × TEACHES))

In relational calculus:

{ (s,c,t,f): ENROLL(s,c,t)  $\land$  TEACHES(f,c,t) }

## Relational Algebra vs. Relational Calculus

Codd's Theorem (informal statement):

Relational Algebra and Relational Calculus have "essentially" the same expressive power, i.e., they can express the same queries.

Note: It is **not** true that for every relational calculus expression  $\varphi$ , there is an equivalent relational algebra expression E.

Examples:

- { ( $x_1,...,x_k$ ):  $\neg R(x_1,...,x_k)$  }
- { x: ∀y,z ENROLLS(x,y,z) }, where ENROLLS(s-name,course,term)

#### From Relational Calculus to Relational Algebra

Note: The previous relational calculus expression may produce different answers when we consider different domains over which the variables are interpreted.

Example: If the variables  $x_1, ..., x_k$  range over a domain D, then  $\{(x_1, ..., x_k): \neg R(x_1, ..., x_k)\} = D^k - R.$ 

Fact:

- The relational calculus expression { (x<sub>1</sub>,...,x<sub>k</sub>): ¬ R(x<sub>1</sub>,...,x<sub>k</sub>) } is not "domain independent".
- The relational calculus expression  $\{(x_1,...,x_k): S(x_1,...,x_k) \land \neg R(x_1,...,x_k)\}$  is "domain independent".

## Active Domain and Active Domain Interpretation

#### Definition:

- The active domain adom(I) of a relational database instance I is the set of all values that occur in the relations of I.
- Let  $\varphi(x_1,...,x_k)$  be a relational calculus formula and let I be a relational database instance. Then  $\varphi^{adom}(I)$

is the result of evaluating  $\phi(x_1,...,x_k)$  over adom(I) and I, that is,

- all variables and quantifiers are assumed to range over adom(I);
- the relation symbols in  $\phi$  are interpreted by the relations in I.

## Queries

Definition: Let **S** be a relational database schema.

A k-ary query on **S** is a function q defined on database instances over S such that if I is a database instance over S, then q(I) is a k-ary relation that is invariant under isomorphisms and has values among those occurring in the relations in I

(i.e., if h:  $I \rightarrow J$  is an isomorphism, then q(J) = h(q(I)).

#### Note:

- All "queries" that we have expressed in relational algebra and/or in relational calculus so far are queries in the above formal sense.
- In particular, a relational calculus expression of the form  $\{(x_1,...,x_k): \phi(x_1,...x_k)\}$  defines a k-ary query.

Theorem: The following are equivalent for a k-ary query q:

1. There is a relational algebra expression E such that q(I) = E(I), for every database instance I (in other words, g is expressible in relational algebra)

(in other words, q is expressible in relational algebra).

2. There is a relational calculus formula  $\psi$  such that q(I) =  $\psi^{adom}$  (I) (in other words, q is expressible in relational calculus under the active domain interpretation).

#### From Relational Algebra to Relational Calculus

(1)  $\Rightarrow$  (2) For every relational expression E, there is an equivalent relational calculus expression {( $x_1,...,x_k$ ):  $\phi(x_1,...,x_k)$ }.

Proof: By induction on the construction of rel. algebra expressions.

- If E is a relation R of arity k, then we take  $\{(x_1,...,x_k): E(x_1,...,x_k)\}$ .
- Assume E<sub>1</sub> and E<sub>2</sub> are expressible by {(x<sub>1</sub>,...,x<sub>k</sub>): φ<sub>1</sub>(x<sub>1</sub>,...,x<sub>k</sub>)} and by {(x<sub>1</sub>,...,x<sub>m</sub>): φ<sub>2</sub>(x<sub>1</sub>,...,x<sub>m</sub>)}. Then
   E<sub>1</sub> ∪ E<sub>2</sub> is expressible by {(x<sub>1</sub>,...,x<sub>k</sub>): φ<sub>1</sub>(x<sub>1</sub>,...,x<sub>k</sub>) ∨ φ<sub>2</sub>(x<sub>1</sub>,...,x<sub>k</sub>)}.
  - $E_1 E_2$  is expressible by { $(x_1,...,x_k)$ :  $\phi_1(x_1,...,x_k) \land \neg \phi_2(x_1,...,x_k)$ }.
  - $E_1 \times E_2$  is expressible by { $(x_1,...,x_k,y_1,...,y_m)$ :  $\phi_1 (x_1,...,x_k) \land \phi_2(y_1,...,y_m)$ }

#### From Relational Algebra to Relational Calculus

(1)  $\Rightarrow$  (2) For every relational expression E, there is an equivalent relational calculus expression {( $x_1,...,x_k$ ):  $\phi(x_1,...,x_k)$ }.

Proof: (continued)

- Assume that E is expressible by {(x<sub>1</sub>,...,x<sub>k</sub>): φ(x<sub>1</sub>,...,x<sub>k</sub>)}. Then
  - $\pi_{1,3}(E)$  is expressible by { $(x_1,x_3): (\exists x_2)(\exists x_4) ...(\exists x_k) \phi(x_1,...,x_k)$  }
  - σ<sub>Θ</sub>(E) is expressible by
     {(x<sub>1</sub>,...,x<sub>k</sub>): Θ\* ∧ φ(x<sub>1</sub>,...,x<sub>k</sub>)}, where Θ\* is the rewriting of Θ as
     a formula of relational calculus.

Proof (Sketch):

- $2. \Rightarrow 1.$ 
  - Show first that for every relational database schema S, there is a relational algebra expression E such that for every database instance I, we have that adom(I) = E(I).
  - Use the above fact and induction on the construction of relational calculus formulas to obtain a translation of relational calculus under the active domain interpretation to relational algebra.

- In this translation, the most interesting part is the simulation of the universal quantifier ∀ in relational algebra.
  - It uses the logical equivalence  $\forall y\psi \equiv \neg \exists y \neg \psi$
- As an illustration, consider  $\forall yR(x,y)$ .

$$\forall y R(x,y) \equiv \neg \exists y \neg R(x,y)$$

• adom(I) = 
$$\pi_1(\mathsf{R}) \cup \pi_2(\mathsf{R})$$

Rel.Calc. formula $\phi$	Relational Algebra Expression for $\phi^{adom}$
¬ R(x,y)	$(\pi_1(R) \cup \pi_2(R))  imes (\pi_1(R) \cup \pi_2(R)) - R$
∃y¬R(x,y)	$\pi_{_1}((\pi_{_1}(R)\cup\pi_{_2}(R)) imes(\pi_{_1}(R)\cup\pi_{_2}(R))$ - R)
¬∃y¬R(x,y)	$(\pi_{_{1}}(R) \cup \pi_{_{2}}(R)) - (\pi_{_{1}}((\pi_{_{1}}(R) \cup \pi_{_{2}}(R)) \times (\pi_{_{1}}(R) \cup \pi_{_{2}}(R)) - R))$

#### Remarks:

- The Equivalence Theorem is effective. Specifically, the proof of this theorem yields two algorithms:
  - an algorithm for translating from relational algebra to domain independent relational calculus, and
  - an algorithm from translating from domain independent relational calculus to relational algebra.
- Each of these two algorithms runs in linear time.

## Queries

Definition: Let **S** be a relational database schema.

- A k-ary query on S is a function q defined on database instances over S such that if I is a database instance over S, then q(I) is a k-ary relation on adom(I) that is invariant under isomorphisms (i.e., if h: I → J is an isomorphism, then q(J) = h(q(I)).
- A Boolean query on S is a function q defined on database instances over S such that if I is a database instance over S, then q(I) = 0 or q(I) = 1, and q(I) is invariant under isomorphisms.

**Example:** The following are Boolean queries on graphs:

- Given a graph E (binary relation), is the diameter of E at most 3?
- Given a graph E (binary relation), is E connected?

- The Query Evaluation Problem: Given a query q and a database instance I, find q(I).
- The Query Equivalence Problem: Given two queries q and q' of the same arity, is it the case that q = q' ?
   (i.e., is it the case that, for every database instance I, we have that q(I) = q'(I)?)
- The Query Containment Problem: Given two queries q and q' of the same arity, is it the case that q ⊆ q' ?
   (i.e., is it the case that, for every database instance I, we have that q(I) ⊆ q'(I)?)

- The Query Evaluation Problem is the main problem in query processing.
- The Query Equivalence Problem underlies query processing and optimization, as we often need to transform a given query to an equivalent one.
- The Query Containment Problem and Query Equivalence Problem are closely related to each other:
  - $\label{eq:q_add} \ \ \, \ \ \, q \equiv q' \ \, \mbox{if and only if } q \subseteq q' \ \, \mbox{and} \ \ q' \subseteq q.$
  - $q \subseteq q'$  if and only if  $q \equiv (q \land q')$ .

- Our goal is to investigate the algorithmic aspects of these problems for queries expressible in relational algebra/relational calculus.
- The questions we want to address are:
  - □ How can we measure the precise "difficulty" of these problems?
  - Are there "good" algorithms for solving these problems?
  - If not, are there special cases of these problems for which "good" algorithms exist?

Our study of these problems will use concepts and methods from two different, yet related, areas:

- Mathematical Logic:
  - Computability Theory and Undecidable Problems
- Computational Complexity Theory:
  - Complexity Classes and Complete Problems
  - □ In particular, the classes P and NP, and NP-complete problems.

## **Decision Problems and Languages**

- Definition (informal): A decision problem Q consists of a set of inputs and a question with a "yes" or "no" answer for each input.
   input x
   Q?
   0 ("no")
- Definition: A decision problem is
  - Decidable if there is an algorithm for solving it;
  - Undecidable if there is no algorithm for solving it.

#### **Undecidable Problems**

Theorem: The following problems are undecidable:

The Halting Problem (A. Turing – 1936): Given a Turing machine M and an input x, does M halt on x?

The Finite Validity Problem (B. Trakhtenbrot – 1949): Given a firstorder formula φ on graphs, is φ true on every finite graph?

## **Undecidable Problems**

- The Finite Validity Problem (B. Trakhtenbrot 1949): Given a firstorder sentence φ on graphs, is φ true on every finite graph?
- Examples of Finitely Valid Formulas:
  - $\Box \quad \forall \ \mathsf{x}(\mathsf{E}(\mathsf{x},\mathsf{x}) \to \exists \ \mathsf{y}\mathsf{E}(\mathsf{x},\mathsf{y}))$
  - $\Box \quad \forall \ \mathsf{x} \forall \ \mathsf{y}(\mathsf{E}(\mathsf{x},\mathsf{x}) \land \mathsf{x} = \mathsf{y} \to \mathsf{E}(\mathsf{y},\mathsf{y}))$
  - "if E is a total order, then E has a biggest element"
- Example of Non-Finitely Valid Formulas:

$$\Box \quad \forall \ x \ \forall \ y \ (\mathsf{E}(x,y) \to \mathsf{E}(y,x))$$

- □ ( $\forall x \exists y E(x,y)$ )  $\rightarrow$  ( $\exists y \forall x E(x,y)$ )
- The undecidability of the Finite Validity Problem means that there is no algorithm for telling formulas in the first group from formulas in the second group.

## The Reduction Method

- By now there is a vast library of undecidable problems.
- The Reduction Method is the main technique for establishing undecidability.
- Reduction Method: To show that a problem L\* is undecidable, it suffices to find an undecidable problem L and a computable function f such that for every input x, we have that

 $x\in L \quad \Leftrightarrow \quad f(x)\in L^*.$ 

- Such a function f is called a reduction of L to L\*
- $L \leq L^*$  means that there is a reduction of L to L\*.

## The Reduction Method

- The Halting Problem was the first fundamental decision problem shown to be undecidable.
- Many database problems have been shown to be undecidable via reductions from
  - The Halting Problem
  - or
  - The Finite Validity Problem

Undecidability of The Query Equivalence Problem

- The Query Equivalence Problem: Given two queries q and q' of the same arity, is it the case that q = q' ?
   (i.e., is q(I) = q'(I) on every database instance I?)
- Theorem: The Query Equivalence Problem for relational calculus queries is undecidable.
   Proof: Finite Validity Problem ≼ Query Equivalence Problem
  - □ To see, this let  $\psi^*$  be a fixed finitely valid relational calculus sentence (say,  $\forall x(E(x,x) \rightarrow \exists yE(x,y)))$ .
  - □ Then, for every relational calculus sentence  $\varphi$ , we have that  $\varphi$  is finitely valid  $\Leftrightarrow \varphi \equiv \psi^*$ .

#### Undecidability of the Query Containment Problem

The Query Containment Problem: Given two queries q and q' of the same arity, is it the case that q ⊆ q' ?
 (i.e., is q(I) ⊆ q'(I) on every database instance I?)

Corollary: The Query Containment Problem for relational calculus queries in undecidable.
 Proof: Query Equivalence ≼ Query Containment, since q ≡ q' ⇔ q ⊆ q' and q' ⊆ q.

Notice the chain of reductions:

Halting Problem  $\preccurlyeq$  Finite Validity  $\preccurlyeq$  Query Equiv.  $\preccurlyeq$  Query Cont.

## The Query Evaluation Problem

- The Query Evaluation Problem: Given a query q and a database instance I, find q(I).
- The Query Evaluation Problem for relational calculus queries is decidable, but, as we will see, it has high computational complexity.
- To understand the precise algorithmic difficulty of the Query Evaluation Problem, we need some basic notions and results from computational complexity.

## Decidable Problems and Computational Complexity



 Computational Complexity is the quantitative study of decidable problems.

"From these and other considerations grew our deep conviction that **there must be quantitative laws that govern the behavior of information and computing**. The results of this research effort were summarized in our first paper on this topic, which also named this new research area, "On the computational complexity of algorithms"."

J. Hartmanis, Turing Award Lecture, 1993

## **Computational Complexity Classes**

- Decidable problems are grouped together in computational complexity classes.
- Each computational complexity class consists of all problems that can be solved in a computational model under certain restrictions on the resources used to solve the problem.
- Examples of computational models:
  - Turing Machine TM (deterministic Turing machine)
  - Non-deterministic Turing machine NTM

• ...

- Examples of resources:
  - Amount of time needed to solve the problem
  - □ Amount of space (memory) needed to solve the problem.

•••••

## **Computational Complexity Classes**



There are many other complexity classes. For a comprehensive catalog, visit the Complexity Zoo at

qwiki.stanford.edu/wiki/Complexity\_Zoo

# Complexity of the Query Evaluation Problem

The Query Evaluation Problem for Relational Calculus: Given a relational calculus formula  $\varphi$  and a database instance I, find  $\varphi^{adom}(I)$ .

Theorem: The Query Evaluation Problem for Relational Calculus is PSPACE-complete.

**Proof:** We need to show that

- This problem is in PSPACE.
- This problem is PSPACE-hard.
  We start with the second task.

## Complexity of the Query Evaluation Problem

- Theorem: The Query Evaluation Problem for Relational Calculus is PSPACE-hard.
- Proof: Show that

Quantified Boolean Formulas  $\preccurlyeq_{D}$  Query Evaluation for Rel. Calc.

Given QBF  $\forall \mathbf{x}_1 \exists \mathbf{x}_2 \dots \forall \mathbf{x}_k \psi$ 

- Let V and P be two unary relation symbols
- Obtain  $\psi^*$  from  $\psi$  by replacing  $x_i$  by  $P(x_i)$ , and  $\neg x_i$  by  $\neg P(x_i)$
- Let I be the database instance with V = {0,1}, P={1}.
- Then the following statements are equivalent:
  - $\forall \mathbf{x_1} \exists \mathbf{x_2} \dots \forall \mathbf{x_k} \psi$  is true
  - $\forall x_1 (V(x_1) \rightarrow \exists x_2 (V(x_2) \land (... \forall x_k (V(x_k) \rightarrow \psi^*))...) \text{ is true on I.}$

## Complexity of the Query Evaluation Problem

- Theorem: The Query Evaluation Problem for Relational Calculus is in PSPACE.
   Proof (Hint): Let φ be a relational calculus formula ∀x<sub>1</sub>∃x<sub>2</sub> ... ∀x<sub>m</sub>ψ and let I be a database instance.
  - Exponential Time Algorithm: We can find  $\varphi^{adom}(I)$ , by exhaustively cycling over all possible interpretations of the  $x_i$ 's.

This runs in time  $O(n^m)$ , where n = |I| (size of I).

- A more careful analysis shows that this algorithm can be implemented in O(m·logn)-space.
  - Use m blocks of memory, each holding one of the n elements of adom(I) written in binary (so O(logn) space is used in each block).
  - Maintain also m counters in binary to keep track of the number of elements examined.

$\forall x_1$	$\exists x_2$	 $\forall \mathbf{x}_{m}$
a <sub>1</sub> in adom(I)	a <sub>2</sub> in adom(I)	 a <sub>m</sub> in adom(I)
written in binary	written in binary	written in binary

## Summary

 Relational Algebra and Relational Calculus have "essentially" the same expressive power.

- The Query Equivalence Problem for Relational Calculus in undecidable.
- The Query Containment Problem for Relational Calculus is undecidable.
- The Query Evaluation Problem for Relational Calculus is PSPACEcomplete.

# Sublanguages of Relational Calculus

 Question: Are there interesting sublanguages of relational calculus for which the Query Containment Problem and the Query Evaluation Problem are "easier" than the full relational calculus?

#### Answer:

- □ Yes, the language of conjunctive queries is such a sublanguage.
- Moreover, conjunctive queries are the most frequently asked queries against relational databases.

 Definition: A conjunctive query is a query expressible by a relational calculus formula in prenex normal form built from atomic formulas R(y<sub>1</sub>,...,y<sub>n</sub>), and ∧ and ∃ only.

{ ( $x_1,...,x_k$ ):  $\exists z_1 ... \exists z_m \chi(x_1,...,x_k, z_1,...,z_k)$  },

where  $\chi(x_1, ..., x_k, z_1, ..., z_k)$  is a conjunction of atomic formulas of the form  $R(y_1, ..., y_m)$ .

 Equivalently, a conjunctive query is a query expressible by a relational algebra expression of the form

 $\pi_X(\sigma_\Theta(\mathsf{R}_1 \times \ldots \times \mathsf{R}_n))$ , where

 $\Theta$  is a conjunction of equality atomic formulas (equijoin).

 Equivalently, a conjunctive query is a query expressible by an SQL expression of the form

SELECT <list of attributes>

FROM <list of relation names>

WHERE <conjunction of equalities>

- Definition: A conjunctive query is a query expressible by a relational calculus formula in prenex normal form built from atomic formulas R(y<sub>1</sub>,...,y<sub>n</sub>), and ∧ and ∃ only.
   { (x<sub>1</sub>,...,x<sub>k</sub>): ∃ z<sub>1</sub> ...∃ z<sub>m</sub> χ(x<sub>1</sub>, ...,x<sub>k</sub>, z<sub>1</sub>,...,z<sub>k</sub>) }
  - A conjunctive query can be written as a logic-programming rule:

 $Q(x_1,...,x_k) :-- R_1(\mathbf{u}_1), ..., R_n(\mathbf{u}_n)$ , where

- Each variable x<sub>i</sub> occurs in the right-hand side of the rule.
- Each u<sub>i</sub> is a tuple of variables (not necessarily distinct)
- The variables occurring in the right-hand side (the body), but not in the left-hand side (the head) of the rule are existentially quantified (but the quantifiers are not displayed).
- "," stands for conjunction.

Examples:

- □ Path of Length 2: (Binary query)  $\{(x,y): \exists z (E(x,z) \land E(z,y))\}$ 
  - As a relational algebra expression,  $\pi_{1,4}(\sigma_{\$2} = \$3} (E \times E))$

As a rule: q(x,y) :-- E(x,z), E(z,y)

- □ Cycle of Length 3: (Boolean query)  $\exists x \exists y \exists z(E(x,y) \land E(y,z) \land E(z,x))$ 
  - As a rule (the head has no variables)
     Q :-- E(x,z), E(z,y), E(z,x)

- Every relational join is a conjunctive query: P(A,B,C), R(B,C,D) two relation symbols
  - P  $\bowtie$  R = {(x,y,z,w): P(x,y,z)  $\land$  R(y,z,w)}
  - q(x,y,z,w) :-- P(x,y,z), R(y,z,w)
     (no variables are existentially quantified)
  - SELECT P.A, P.B, P.C, R.D
     FROM P, R
     WHERE P.B = R.B AND P.C = R.C
- Conjunctive queries are also known as SPJ-queries (SELECT-PROJECT-JOIN queries)

#### Conjunctive Query Evaluation and Containment

- Definition: Two fundamental problems about CQs
  - Conjunctive Query Evaluation (CQE):
     Given a conjunctive query q and an instance I, find q(I).
  - Conjunctive Query Containment (CQC):
    - Given two k-ary conjunctive queries q<sub>1</sub> and q<sub>2</sub>, is it true that q<sub>1</sub> ⊆ q<sub>2</sub>?
       (i.e., for every instance I, we have that q<sub>1</sub>(I) ⊆ q<sub>2</sub>(I))
    - Given two Boolean conjunctive queries q<sub>1</sub> and q<sub>2</sub>, is it true that q<sub>1</sub> ⊨ q<sub>2</sub>? (that is, for all I, if I ⊨ q<sub>1</sub>, then I ⊨ q<sub>2</sub>)?
       CQC is logical implication.

# CQE vs. CQC

- Recall that for relational calculus queries:
  - The Query Evaluation Problem is PSPACE-complete (combined complexity).
  - The Query Containment Problem is undecidable.
- Theorem: Chandra & Merlin, 1977
  - CQE and CQC are the "same" problem.
  - Moreover, each is an NP-complete problem.
- Question: What is the common link?
- Answer: The Homomorphism Problem