



#### Perfect Recall: Database Systems (→ 165A)

- A Database System (DBS) consists of a Database (DB) and a Database Management System (DBMS)
- A Database is a (typically very large) integrated collection of interrelated data which are stored in files.
- Data can come from commercial or scientific applications and (usually) represent some abstraction/piece of the modeled real world.
- E.g, a scientific database might contain information about known biological, chemical, astronomical entities, lab experiments, etc
- A Database Management System is a collection of software packages designed to store, access, and manage databases. It provides users and applications with an environment that is convenient and efficient to use.



Ludaescher ECS289E-W05 Tonics in Scientific Data Management

#### **Relational Database Model**

- Think of a relational DB as a number of *tables*, each have a particular *schema*:
  - Course(Instructor, Name, Quarter, Department)
- The *table/relation name* "Course", identifies which table we are talking about.
- The attribute/column name (e.g., "Instructor") corresponds to the "column header"
- Elements aka *instances* or *tuples* of a table/relation can be written, e.g., as follows: Course("Gertz", "ECS165A", "W-2005", "CS"). Course("Ludaescher", "ECS289F", "W-2005", "CS"). ...

Example						
	Course					
	Instructor	Name		Quarter	Department	
	Gertz	ECS165A		W-2005	CS	
	Ludaescher	ECS289F		W-2005	CS	
•	<ul> <li>The same in <i>Datalog notation</i> – as a set of <i>facts</i>: course('Ludaescher', 'ECS289F', 'W-2005', 'CS').</li> <li>course(,,,).</li> </ul>					

Ludaescher, ECS289F-W05, Topics in Scientific Data Managemen

escher, ECS289F-W05, Topics in Scientific Data Man

# Hmm.. looks like a Spreadsheet ...

- ... but there are differences.
- What are they?

er, ECS289F-W05, Topics in Sc







## **Query Languages**

- Databases can be queried!
- We state a question, usually in terms of the given database schema, about the stored data.
- **Query languages** such as Datalog and SQL (Structured Query Language) are *declarative* (just say what you're interested in) – you do not need to give the details *how* to retrieve the data, but can focus on the *what* (to retrieve).

## Question

• What's the difference between keyword-based search and querying a database?



But watch out

ther, ECS289F-W05, Topics in Scientific Data Ma

 ... some recent work in the database community on "keyword search in databases"...



er, ECS289F-W05, Topics in Scientific Data M

#### **DATALOG: Examples of Relational Operations** Relational operations have concise representations! Examples: sel(X,Y) :- p(X,Y), X=a, not X=Y. % SELECT some tuples from p(X,Y) proj(X) :- p(X,Y). % PROJECT on the first argument $\texttt{join}(\texttt{X},\texttt{Y},\texttt{Z}) \ := \ \texttt{p}(\texttt{X},\texttt{Y}), \ \texttt{q}(\texttt{Y},\texttt{Z}). \qquad \ \texttt{\%} \ \texttt{JOIN} \ \texttt{p}(\texttt{A},\texttt{B}), \ \texttt{q}(\texttt{C},\texttt{D}) \ \texttt{s.t.} \ \texttt{B=C}$ prod(X,Y) :- p(X), q(Y). % PRODUCT of p(X) and q(Y)intersect(X) :- p(X), q(X). % INTERSECTION of p(X), q(X) diff(X) := p(X), not q(X).% SET-DIFFERENCE: $p(X) \setminus q(X)$ union(X) :- p(X). union(X) :- q(X). % UNION of p(X), ... % ... and q(X) Rules have a "logical reading" (i.e., rules are formulas): $\begin{array}{ll} \forall X \;(\; \texttt{diff}(X) \; \leftarrow \; \texttt{p}(X) \land \neg \, \texttt{q}(X) \;). \\ \forall X \;(\; \texttt{union}(X) \; \leftarrow \; \texttt{p}(X) \lor \texttt{q}(X) \;). \end{array}$

## What is a Query?

- A <u>query expression</u> e.g. in SQL or in Datalog denotes a query (but we still don't know what a query is...)
- A <u>query</u> is a (generic\*) mapping *f* from instances of an input schema (EDB) to instances of an output schema (IDB):
   *f* : inst(EDB) → inst(IDB)

er, ECS289F-W05, Topics in Scientific Data M

• Note: Different query expressions can denote the same query (mapping). Example...?

#### What is a Query?

 A query is a <u>generic</u> mapping *f* from instances of an input schema (EDB) to instances of an output schema (IDB):
 *f* : *inst(EDB)* → *inst(IDB)*

\*generic: invariant under renamings r, i.e.,

- f(r(I)) = r(f(I)) for all database instances *I* of the schema EDB
- Examples: Consider EBD = {p(X), emp(N,S)}. Which of the following are generic?
  - f\_even: "T" if  $\mid \{x \mid p(x) \text{ is in DB } I\} \mid \text{ is even}$
  - f\_jeff: { (N,S) | emp(N,S) in DB I, N = "Jeff" }

#### **Problem**

• How can one *evaluate* DATALOG queries? That is, given a database *instance* (= a set of *facts*), how can one obtain the *answer* to a given query (=*rule* or *set of rules*) ?



Ludaescher, ECS289F-W05, Topics in Scientific Data Managemen

cher, ECS289F-W05, Topics in Scientific Data Manageme



#### **Example: Transitive Closure**

e(a,b). e(b,c). e(c,d). % FACTS (EDB-relation: e/2) 

% RULES (IDB-relation: tc/2)

#### 0 Ø

- 1 {e(a, b), e(b, c), e(c, d)}
- $\begin{array}{l} \{e(a, b), e(b, c), e(c, d)\} \cup \{tc(a, b), tc(b, c), tc(c, d)\} \\ 2 \{e(a, b), e(b, c), e(c, d)\} \cup \{tc(a, b), tc(b, c), tc(c, d), tc(a, c), tc(b, d)\} \\ 4 \{e(a, b), e(b, c), e(c, d)\} \cup \{tc(a, b), tc(b, c), tc(c, d), tc(a, c), tc(b, d), tc(a, d)\} \end{array}$

 $\Rightarrow$  if the longest path in e/2 has length n, then O(n) rounds are needed!

(Exercise: how about the following rule?  $tc(X,Y) \leftarrow tc(X,Z), tc(Z,Y).$ )

## **DATALOG: Minimal Model Semantics**

Rules can be seen as first-order formulas:

her, ECS289F-W05, Topics in Scientific Data Managemen

- $p(X) := q(X, Y) :\Leftrightarrow \forall X(p(X) \leftarrow \exists Yq(X, Y)).$
- A model of a program P (=  $Fact \cup Rules$ ) is an interpretation of the relations which satisfies all rules. (Here, we're dealing with Herbrand interpretations and models, i.e., which interpret

(nee, we be doming with reporting the prevalues and modes, i.e., which memory constants "syntactically". Then, an interpretation can be written as a set of *true* tuples; all other tuples are regarded as *false*.)

- From the possible different models find the "intuitive", "intended" models, i.e., which make true only what is "strictly necessary", given the rules of P ⇒ minimal models:
  A model M of P is minimal, if there is no other model M' ⊂ M.



