Institute of Information and Communication and Technology Bangladesh University of Engineering and Technology



ICT 5103: Database Design and Management

Lecture 6

Instructor: Samin Rahman Khan

Relational Algebra

Chapter 4, Part A

© Database Management Systems 3ed, R. Ramakrishnan and J. Gehrke

Relational Query Languages

- **Query languages:** Allow manipulation and retrieval of data from a database.
- Relational model supports simple, powerful QLs:
 - \circ $\;$ Strong formal foundation based on logic.
 - Allows for much optimization.
- Query Languages != programming languages!
 - QLs not expected to be "Turing complete".
 - $\circ~$ QLs not intended to be used for complex calculations.
 - QLs support easy, efficient access to large data set

Formal Relational Query Languages

- Two mathematical Query Languages form the basis for "real" languages (e.g. SQL), and for implementation:
 - **Relational Algebra:** More operational, very useful for representing execution plans.
 - **Relational Calculus:** Lets users describe what they want, rather than how to compute it. (Non-operational, **declarative**.)

Preliminaries

- A query is applied to relation instances, and the result of a query is also a relation instance.
 - Schemas of input relations for a query are fixed (but query will run regardless of instance!)
 - The schema for the result of a given query is also fixed! Determined by definition of query language constructs.
- Positional vs. named-field notation:
 - Positional notation easier for formal definitions, named-field notation more readable.
 - \circ Both used in SQL

Example Instances

- "Sailors" and "Reserves" relations for our examples.
- We'll use positional or named field notation, assume that names of fields in query results are `inherited' from names of fields in query input relations.

Sailors(*sid:* integer, *sname:* string, *rating:* integer, *age:* real Boats(*bid:* integer, *bname:* string, *color:* string) Reserves(*sid:* integer, *bid:* integer, *day:* date)

R1	sid	bid	day
	22	101	10/10/96
	58	103	11/12/96

S1	sid	sname	rating	age
	22	dustin	7	45.0
	31	lubber	8	55.5
	58	rusty	10	35.0

S2	sid	sname	rating	age
-)	28	yuppy	9	35.0
L)	31	lubber	8	55.5
	44	guppy	5	35.0
	58	rusty	10	35.0

Relational Algebra

- Basic operations:
 - $\circ~$ Selection (5) Selects a subset of rows from relation.
 - **Projection** (π) Deletes unwanted columns from relation.
 - **Cross-product (X)** Allows us to combine two relations.
 - Set-difference (-) Tuples in reln. 1, but not in reln. 2.
 - **Union (U)** Tuples in reln. 1 and in reln. 2.
- Additional operations:
 - Intersection, **join**, division, renaming: Not essential, but (very!) useful.
- Since each operation returns a relation, operations can be composed! (Algebra is "closed".)

Projection

- Deletes attributes that are not in projection list.
- Schema of result contains exactly the fields in the projection list, with the same names that they had in the (only) input relation.
- Projection operator has to eliminate duplicates! (Why??)
 - Note: real systems typically don't do duplicate elimination unless the user explicitly asks for it. (Why not?)



Selection

- Selects rows that satisfy selection condition.
- No duplicates in result! (Why?)
- Schema of result identical to schema of (only) input relation.
- Result relation can be the input for another relational algebra operation! (Operator composition.)

sid	sname	rating	age
28	yuppy	9	35.0
58	rusty	10	35.0

$$\sigma_{rating>8}^{(S2)}$$

sname	rating
yuppy	9
rusty	10



Union, Intersection, Set-Difference

- All of these operations take two input relations, which must be **union-compatible**:
 - Same number of fields.
 - `Corresponding' fields have the same type.
- What is the schema of the result?

sid	sname	rating	age
22	dustin	7	45.0

$$S1 - S2$$

sid	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0
44	guppy	5	35.0
28	yuppy	9	35.0

 $S1 \cup S2$

sid	sname	rating	age
31	lubber	8	55.5
58	rusty	10	35.0

 $S1 \cap S2$

Cross-Product

- Each row of S1 is paired with each row of R1.
- Result schema has one field per field of S1 and R1, with field names `inherited' if possible.
 - \circ $\,$ Conflict: Both S1 and R1 have a field called sid.

(sid)	sname	rating	age	(sid)	bid	day
22	dustin	7	45.0	22	101	10/10/96
22	dustin	7	45.0	58	103	11/12/96
31	lubber	8	55.5	22	101	10/10/96
31	lubber	8	55.5	58	103	11/12/96
58	rusty	10	35.0	22	101	10/10/96
58	rusty	10	35.0	58	103	11/12/96

• **Renaming operator:** $\rho(C(1 \rightarrow sid1, 5 \rightarrow sid2), S1 \times R1)$

Joins

• **Condition Join:** $R \bowtie_C S = \sigma_C(R \times S)$

(sid)	sname	rating	age	(sid)	bid	day
22	dustin	7	45.0	58	103	11/12/96
31	lubber	8	55.5	58	103	11/12/96
$S1 \Join_{S1.sid < R1.sid} R1$						

- Result schema same as that of cross-product.
- Fewer tuples than cross-product, might be able to compute more efficiently
- Sometimes called a theta-join.

Joins

• Equi-Join: A special case of condition join where the condition c contains only equalities.

sid	sname	rating	age	bid	day		
22	dustin	7	45.0	101	10/10/96		
58	rusty	10	35.0	103	11/12/96		
$S1 \Join_{sid} R1$							

- Result schema similar to cross-product, but only one copy of fields for which equality is specified.
- Natural Join: Equijoin on all common fields.

Division

• Not supported as a primitive operator, but useful for expressing queries like:

Find sailors who have reserved all boats.

- Let A have 2 fields, x and y; B have only field y:
 - $\circ \quad A/B = \{ \langle x \rangle \mid \exists \langle x, y \rangle \in A \ \forall \ \langle y \rangle \in B \}$
 - i.e., A/B contains all x tuples (sailors) such that for every y tuple (boat) in B, there is an xy tuple in A.
 - Or: If the set of y values (boats) associated with an x value (sailor) in A contains all y values in B, the x value is in A/B.
- In general, x and y can be any lists of fields; y is the list of fields in B, and x y is the list of fields of A.

Examples of Division A/B



Expressing A/B Using Basic Operator

- Division is not essential op; just a useful shorthand.
 - (Also true of joins, but joins are so common that systems implement joins specially.)
- Idea: For A/B, compute all x values that are not `disqualified' by some y value in B.
 - x value is disqualified if by attaching y value from B, we obtain an xy tuple that is not in A.

Disqualified x values: $\pi_{X}((\pi_{X}(A) \times B) - A)$

A/B: $\pi_{X}(A)$ – all disqualified tuple

Find names of sailors who've reserved #103

• Solution 1: $\pi_{\text{sname}}((\sigma_{\text{bid}=103} \text{ Reserves}) \bowtie \text{ Sailors})$

Solution 2: ρ (Temp1, σ_{bid=103} Reserves)
ρ (Temp2, Temp1 ⋈ Sailors)

 π_{sname} (Temp2)

• Solution 3: $\pi_{\text{sname}}(\sigma_{\text{bid}=103} \text{ (Reserves } \bowtie \text{ Sailors))}$

Find names of sailors who've reserved a red boat

• Information about boat color only available in Boats; so need an extra join:

 $\pi_{sname}((\sigma_{color='red'} \text{ Boats}) \bowtie Reserves \bowtie Sailors)$

• A more efficient solution:

 $\pi_{\text{sname}}(\pi_{\text{sid}}((\pi_{\text{bid}}\sigma_{\text{color='red'}}\text{Boats}) \bowtie \text{Reserves}) \bowtie \text{Sailors})$ A query optimizer can find this, given the first solution!

Find sailor who've reserved a red or a green boat

• Can identify all red or green boats, then find sailors who've reserved one of these boats:

 $\rho(Tempboats, (\sigma_{color='red' \ \lor \ color='green'}Boats))$

 $\pi_{_{sname}}(\text{Tempboats} \bowtie \text{Reserves} \bowtie \text{Sailors})$

- Can also define Tempboats using union! (How?)
- What happens if \lor is replaced by \land in this query?

Find sailors who've reserved a red and a green boat

• Previous approach won't work! Must identify sailors who've reserved red boats, sailors who've reserved green boats, then find the intersection (note that sid is a key for Sailors):

$$\begin{split} &\rho(\text{Tempred}, \, \pi_{\text{sid}}((\sigma_{\text{color='red'}}\text{Boats}) \bowtie \, \text{Reserves})) \\ &\rho(\text{Tempgreen}, \, \pi_{\text{sid}}((\sigma_{\text{color='green''}}\text{Boats}) \bowtie \, \text{Reserves})) \\ &\pi_{\text{sname}}((\text{Tempred} \, \cap \, \text{Tempgreen}) \bowtie \, \text{Sailors}) \end{split}$$

Find the names of sailors who've reserved all boats

• Uses division; schemas of the input relations to / must be carefully chose:

$$\rho$$
 (Tempsids, ($\pi_{sid,bid}^{\text{Reserves}}$) / (π_{bid}^{Boats}))
 π_{sname} (Tempsids \Join Sailors)

• To find sailors who've reserved all 'Interlake' boats:

.....
$$\pi_{bid}(\sigma_{bname='Interlake'}^{mass}Boats)$$

Summary

- The relational model has rigorously defined query languages that are simple and powerful.
- Relational algebra is more operational; useful as internal representation for query evaluation plans.
- Several ways of expressing a given query; a query optimizer should choose the most efficient version.