Databases I.

RELATIONAL QUERY LANGUAGES - RELATIONAL ALGEBRA

#### Operations for the Relational Model

- Operations on relations can be expressed according to two basic formalisms:
  - Relational Algebra: queries are specified by applying specialized operators to relations
  - Relational Calculus: queries are expressed through logical formulas that must be verified by the tuples that are result of queries
- The two formalisms (under certain assumptions) are equivalent

Relational Algebra

- It consists of five basic operations:
  - Union
  - Difference
  - Cartesian Product
  - Projection
  - Selection
- These operations fully define the relational algebra

#### Relational Algebra

- <u>Each operation returns a relation as result</u>: it is thus possible to apply an operation to the result of another operation (*closure property*)
- There are some additional operations that can be expressed in terms of the five basic operations
  - Join
  - Intersection
  - division

5

#### Relational Algebra

- Such additional operations do not increase the expressive power of the set of the basic operations
- However, these additional operations are used as shorthand; the join is the most important of these operations
- With respect to the notation by name, it can be useful to introduce another operation, called renaming, that allows one to modify the attribute names

#### Relational Algebra

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

#### Relational Algebra - Union

- The union of relations R and S, denoted by R U S
- is the set of tuples that are in R, in S or in both
- The union of two relations can only be executed if the relations have the same degree; in addition the first attribute of R must be compatible with the first attribute of S, the second attribute of R must be compatible with the second attribute of S, and so forth
  - The duplicate tuples are eliminated
  - The degree of the resulting relation is the same as the degree of the input relations

#### Relational Algebra – Union

Example

Α	В	С
а	b	С
d	а	f
С	b	d
R		

D	E	F
b	g	а
d	а	f
S		

Α	В	С
а	b	С
d	а	f
С	b	d
b	g	а

 $R \cup S$ 

9

#### Relational Algebra - Difference

- The difference of relations R and S, denoted as R - S
   is the set of tuples that are in R and are not in S
  - is the set of tuples that are in R and are not in S
- The difference, like the union, can only be executed if the input relations have the same degrees and compatible attributes
- The degree of the resulting relation is the same as the degree of the input relations

### Relational Algebra – Difference Example

Α	В	С
a	b	С
d	а	f
С	b	d
П		

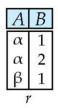
D	Е	F
b	g	а
d	а	f
S		

Α	В	С
а	b	С
С	b	d
R-S		

11

#### Intersection of two relations

Relation r, s:





 $r \cap s$ 

A B  $\alpha$  2

### Example Instances

- "Sailors" and "Reserves" relations for our examples. "bid"= boats. "sid": sailors
- 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.

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

<i>S</i> 1	<u>sid</u>	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
	31	lubber	8	55.5
	44	guppy	5	35.0
	58	rusty	10	35.0

#### Union, Intersection, Set-Difference

- All of these operations take two input relations, which must be <u>union-compatible</u>:
  - Same number of fields.
  - Corresponding' fields have the same type.
- What is the schema of 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

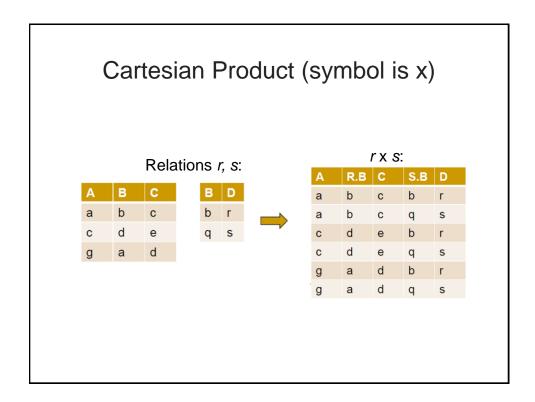
 $S1 \cap S2$ 

#### Cartesian Product (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.
  - 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)



Sample: R3 := R1 X R2

R1( A, B)
1 2
3 4

R2( B, C)
5 6
7 8
9 10

R3( R2.B R1.B, C 1 2 5 6 1 2 7 8 2 9 1 10 3 4 5 6 3 4 7 8 3 4 9 10

#### **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??, what are the consequences?)
  - Note: real systems typically don't do duplicate elimination unless the user explicitly asks for it. (Why not?)

sname	rating
yuppy	9
lubber	8
guppy	5
rusty	10

 $\pi_{sname,rating}(S2)$ 

age 35.0 55.5

 $\pi_{age}(S2)$ 

#### Selection

- Selects rows that satisfy selection condition.
- 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

 $\pi_{sname,rating}(\sigma_{rating>8}(S2))$ 

### Selection of rows (tuples)

The relation  $\mathbf{r}$ :

A	В	C	D
α	α	1	7
$\alpha$	β	5	7
β	β	12	3
β	β	23	10

Select tuples with A=B and D > 5

 $\sigma_{A=B \text{ and } D > 5}(r)$ 

A	В	C	D
α	α	1	7
β	β	23	10

# Projection=Selection of Columns (Attributes)

Relation r.

A	В	C
α	10	1
α	20	1
β	30	1
β	40	2

**Projection** (Select A and C)

$$\Pi_{A,C}$$
 (r)

$$\begin{array}{c|cccc}
A & C \\
\hline
\alpha & 1 \\
\alpha & 1 \\
\beta & 1 \\
\beta & 2
\end{array} = 
\begin{array}{c|cccc}
A & C \\
\hline
\alpha & 1 \\
\beta & 1 \\
\beta & 2
\end{array}$$

#### **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		11/12/96
31	lubber	8	55.5	58	103	11/12/96

$$S1 \bowtie_{S1.sid < R1.sid} R1$$

- Result schema same as that of cross-product.
- Fewer tuples than cross-product. Filters tuples not satisfying the join condition.
- Sometimes called a theta-join.

#### **Joins**

 <u>Equi-Join</u>: 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

$$\pi_{sid,..,age,bid,..}(Sl\bowtie_{sid}Rl)$$

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

#### Joining two relations - Natural Join

- Let r and s be relations on schemas R and S respectively.
  - Then, the "natural join" of relations R and S is a relation on schema  $R \cup S$  obtained as follows:
- Matches tuples with the same values for all common attributes, and retains only one copy of each common column

## Natural Join Example

Relations r, s:

A	В	C	D
α	1	α	a
β	2	γ	a
γ	4	β	b
α	1	γ	a
δ	2	β	b
		r	

B D E

1 a α
3 a β
1 a γ
2 b δ
3 b ε
s

Natural Join r ⋈ s

A	В	C	D	E
α	1	α	a	α
α	1	α	a	γ
α	1	γ	a	α
α	1	γ	a	γ
δ	2	β	b	δ

Relational algebra operations

Symbol (Name)	Example of Use			
σ (Selection)	σ <sub>salary&gt;=85000</sub> (instructor)			
(Selection)	Return rows of the input relation that satisfy the predicate.			
П (Projection)	Π <sub>ID,</sub> salary (instructor)			
(Projection)	Output specified attributes from all rows of the input relation. Remove duplicate tuples from the output.			
M	instructor ⋈ department			
(Natural Join)	Output pairs of rows from the two input relations that have the same value on all attributes that have the same name.			
×	$instructor \times department$			
(Cartesian Product)	Output all pairs of rows from the two input relations (regardless of whether or not they have the same values on common attributes)			
U (Union)	$\Pi_{name}(instructor) \cup \Pi_{name}(student)$			
(	Output the union of tuples from the two input relations.			

#### Find names of sailors who've reserved boat #103

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

Solution 2: 
$$\pi_{sname}(\sigma_{bid=103}(\text{Reserves} \bowtie Sailors))$$

## Find names of sailors who've reserved a red boat

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

$$\pi_{\mathit{sname}}((\sigma_{\mathit{color} = '\mathit{red}'}, \mathit{Boats}) \bowtie \mathsf{Re}\mathit{serves} \bowtie \mathit{Sailors})$$

A more efficient solution:

$$\pi_{\mathit{sname}}(\pi_{\mathit{sid}}((\pi_{\mathit{bid}}\sigma_{\mathit{color} = \mathit{red}}, \mathit{Boats}) \bowtie \mathsf{Res}) \bowtie \mathit{Sailors})$$

A query optimizer can find this, given the first solution!

## Find sailors 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} (Tempboats \bowtie Reserves \bowtie Sailors)
```

## Find sailors who've reserved a red <u>and</u> a green boat

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):

$$\rho \; (\textit{Tempred}, \pi_{\textit{sid}}((\sigma_{\textit{color} = '\textit{red}'}, \textit{Boats}) \bowtie \textit{Reserves}))$$
 
$$\rho \; (\textit{Tempgreen}, \pi_{\textit{sid}}((\sigma_{\textit{color} = '\textit{green}'}, \textit{Boats}) \bowtie \textit{Reserves}))$$
 
$$\pi_{\textit{sname}}((\textit{Tempred} \cap \textit{Tempgreen}) \bowtie \textit{Sailors})$$

#### Outer Join

- Often in joining two relations, a tuple in one relation does not have a matching tuple in the other relation; in other words, there is no matching value in the join attributes.
- We may want tuples from one of the relations to appear in the result even when there are no matching values in the other relation.
- The advantage of an Outer join is that information is preserved; that is, the Outer join preserves tuples that would have been lost by other types of join.
- R ⋈ S The LEFT OUTER JOIN is a join in which tuples from R that do not have matching values in the common attributes of S are also included in the result relation. Missing values in the second relation are set to null.
- R ⋈ S The RIGHT OUTER JOIN is a join in which tuples from S that do not have matching values in the common attributes of R are also included in the result relation. Missing values in the first relation are set to null.
- **R** ➤ **S** The **FULL OUTER JOIN** keeps all tuples in both relations, padding tuples with nulls when no matching tuples are found.

#### **Examples for JOIN operations**

We use *DreamHome* database relation schema:

Branch (branchNo, street, city, postcode)

Staff (staffNo, fName, IName, position, sex, DOB, salary, *branchNo*)

PropertyForRent (propertyNo, street, city, postcode, type, rooms, rent,

ownerNo, staffNo, branchNo)

Client (<u>clientNo</u>, fName, IName, telNo, prefType, maxRent, eMail)

PrivateOwner (ownerNo, fName, IName, address, telNo, eMail, password)

Viewing (<u>clientNo</u>, <u>propertyNo</u>, viewDate, comment)

Registration (clientNo, branchNo, staffNo, dateJoined)

**Ex. 1**. List the names and comments of all clients who have viewed a property for rent. (Equijoin operation)

 $(\Pi_{\text{clientNo, fName, IName}}(\text{Client})) \bowtie_{\text{Client.clientNo = Viewing.clientNo}} (\Pi_{\text{clientNo, propertyNo, comment}}(\text{Viewing}))$ 

**Ex. 2**. List the names and comments of all clients who have viewed a property for rent. (Natural join operation)

 $(\Pi_{\text{clientNo, fName, IName}}(\text{Client})) \ \bowtie \ (\Pi_{\text{clientNo, propertyNo, comment}}(\text{Viewing}))$ 

clientNo	fName	Name	propertyNo	comment
CR76	John	Kay	PG4	too remote
CR56	Aline	Stewart	PA14	too small
CR56	Aline	Stewart	PG4	
CR56	Aline	Stewart	PG36	
CR62	Mary	Tregear	PA14	no dining room

**Ex. 3.** Produce a relation consisting of the properties that have been viewed with comments and those that have not been viewed.

Note that properties PL94, PG21, and PG16 have no viewings, but these tuples are still contained in the result with nulls for the attributes from the Viewing relation. (Left outer join.)

propertyNo	street	city	clientNo	viewDate	comment
PA14 PA14 PL94 PG4 PG4 PG36 PG21 PG16	16 Holhead 16 Holhead 6 Argyll St 6 Lawrence St 6 Lawrence St 2 Manor Rd 18 Dale Rd 5 Novar Dr	Aberdeen Aberdeen London Glasgow Glasgow Glasgow Glasgow Glasgow	CR56 CR62 null CR76 CR56 CR56 null null	,	too small no dining room null too remote null null

#### **Aggregate Functions and Operations**

Aggregation function takes a collection of values and returns a single value as a result.

avg: average valuemin: minimum valuemax: maximum valuesum: sum of valuescount: number of values

■ Aggregate operation in relational algebra

$$G_{1,G2,\dots Gn}$$
  $\Gamma_{F1(A1),F2(A2),\dots Fn(An)}$  (E)

E is any relational-algebra expression

- $G_1$ ,  $G_2$  ...,  $G_n$  is a list of attributes on which to group (can be empty)
- Each F<sub>i</sub> is an aggregate function
- Each A<sub>i</sub> is an attribute name

#### Aggregate Operation - Example

Relation r.

Α	В	С
α	α	7
α	β	7
β	β	3
β	β	10

 $\Gamma_{\text{sum(c)}}(r)$ 

sum(c)

#### Aggregate, GROUPING Operation - Example

Find the average salary in each department

$$_{dept\_name}$$
  $\Gamma_{avg(salary)}$  (instructor)

#### *Instructo*r

ID	пате	dept_name	salary
76766	Crick	Biology	72000
45565	Katz	Comp. Sci.	75000
10101	Srinivasan	Comp. Sci.	65000
83821	Brandt	Comp. Sci.	92000
98345	Kim	Elec. Eng.	80000
12121	Wu	Finance	90000
76543	Singh	Finance	80000
32343	El Said	History	60000
58583	Califieri	History	62000
15151	Mozart	Music	40000
33456	Gold	Physics	87000
22222	Einstein	Physics	95000

#### Result:

dept_name	salary
Biology	72000
Comp. Sci.	77333
Elec. Eng.	80000
Finance	85000
History	61000
Music	40000
Physics	91000

#### **Aggregate Functions**

- Result of aggregation does not have a name
  - Can use rename operation to give it a name
  - For convenience, we permit renaming as part of aggregate operation

 $dept_name \Gamma$  **avg**(salary) **as** avg\_sal (instructor)

#### 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.