- After completing this chapter, you should be able to
 - enumerate and explain the operations of relational algebra (there is a core of 5 relational algebra operators),
 - write relational algebra queries of the type join-select-project,
 - ▷ discuss correctness and equivalence of given relational algebra queries.

Relational Algebra



Example Database (recap)

	STU	DENTS	
<u>SID</u>	FIRST	LAST	EMAIL
101	Ann	Smith	
102	Michael	Jones	(null)
103	Richard	Turner	
104	Maria	Brown	

EXERCISES									
<u>CAT</u>	<u>ENO</u>	TOPIC MAXPT							
Н	1	Rel.Alg.	10						
Н	2	SQL	10						
М	1	SQL	14						

	RESULTS									
SID	<u>SID CAT ENO</u> POINTS									
101		Н	1	10						
101		H	2	8						
101		М	1	12						
102		H	1	9						
102		H	2	9						
102		М	1	10						
103		Н	1	5						
103		М	1	7						

Relational Algebra (1)

- **Relational algebra (RA)** is a **query language** for the relational model with a solid theoretical foundation.
- Relational algebra is *not* visible at the user interface level (not in any commercial RDBMS, at least).
- However, almost any RDBMS uses RA to represent queries internally (for query optimization and execution).
- Knowledge of relational algebra will help in understanding SQL and relational database systems in general.

Relational Algebra (2)

• One particular operation of relational algebra is **selection.**

Many operations of the relational algebra are denoted as greek letters. Selection is σ (sigma).

• For example, the operation $\sigma_{\rm SID=101}$ selects all tuples in the input relation which have the value 101 in column SID.



Relational Algebra (3)

• Since the output of any RA operation is some relation *R* again, *R* may be the input for another RA operation.

The operations of RA nest to arbitrary depth such that complex queries can be evaluated. The final results will always be a relation.

• A query is a term (or expression) in this relational algebra.



Relational Algebra (4)

- There are some difference between the two **query languages** RA and SQL:
 - ▷ Null values are usually excluded in the definition of relational algebra, except when operations like **outer join** are defined.
 - ▷ Relational algebra treats relations as sets, *i.e.*, duplicate tuples will never occur in the input/output relations of an RA operator.

Remember: In SQL, relations are **multisets** (bags) and may contain duplicates. Duplicate elimination is explicit in SQL (SELECT DISTINCT).

Selection (1)

Selection

The **selection** σ_{φ} selects a subset of the tuples of a relation, namely those which satisfy **predicate** φ . Selections acts like a filter on a set.



Selection (2)

• A simple **selection predicate** φ has the form

 $\langle Term \rangle \langle ComparisonOperator \rangle \langle Term \rangle.$

- (*Term*) is an expression that can be **evaluated to a data value** for a given tuple:
 - > an attribute name,
 - b a constant value,
 - ▷ an expression built from attributes, constants, and data type operations like +, -, *, /.

Selection (3)

• $\langle ComparisonOperator \rangle$ is

$$\triangleright = (equals), \neq (not equals),$$

- \triangleright < (less than), > (greater than), \leqslant , \geqslant ,
- ▷ or other data type-dependent predicates (*e.g.*, LIKE).
- Examples for simple selection predicates:
 - ▷ LAST = 'Smith'
 - \triangleright POINTS \geqslant 8
 - \triangleright POINTS = MAXPT.

Selection (4)

```
• \sigma_{\varphi}(R) may be imlemented as:
  "Naive" selection
     create a new temporary relation T;
     foreach t \in R do
        p \leftarrow \varphi(t);
        if p then
          insert t into T;
        fi
     od
     return T;
```

 If index structures are present (e.g., a B-tree index), it is possible to evaluate σ_φ(R) without reading every tuple of R.

Selection (5)

A few corner cases



Selection (6)

• $\sigma_{\varphi}(R)$ corresponds to the following SQL query:

• A different relational algebra operation called **projection** corresponds to the SELECT clause. Source of confusion.



Selection (7)

• More **complex selection predicates** may be performed using the **Boolean connectives:**

 $\triangleright \ \varphi_1 \land \varphi_2 \ (\text{``and''}), \quad \varphi_1 \lor \varphi_2 \ (\text{``or''}), \quad \neg \varphi_1 \ (\text{``not''}).$

• Note:
$$\sigma_{\varphi_1 \wedge \varphi_2}(R) = \sigma_{\varphi_1}(\sigma_{\varphi_2}(R)).$$

• The selection predicate must permit evaluation for each input tuple in **isolation.** A predicate may not refer to other tuples.

Projection (1)

Projection

The **projection** π_L eliminates all attributes (columns) of the input relation but those mentioned in the **projection list** *L*.



Projection (2)

- The projection $\pi_{A_{i_1},\ldots,A_{i_k}}(R)$ produces for each input tuple $(A_1: d_1,\ldots,A_n: d_n)$ an output tuple $(A_{i_1}: d_{i_1},\ldots,A_{i_k}: d_{i_k})$
- π may be used to **reorder columns.**

" σ discards rows, π discards columns."

• DB slang: "All attributes not in L are projected away."

Projection (3)

• In general, the **cardinalities** of the input and output relations **are not equal**.



Projection (4)

```
• \pi_{A_{i_1},\ldots,A_{i_k}}(R) may be imlemented as:
   "Naive" projection
     create a new temporary relation T:
     foreach t = (A_1 : d_1, \ldots, A_n : d_n) \in R do
         u \leftarrow (A_{i_1} : d_{i_1}, \ldots, A_{i_k} : d_{i_k});
         insert u into T:
     od
     eliminate duplicate tuples in T;
      return T;
```

 The necessary duplicate elimination makes π_L one of the more costly operations in RDBMSs. Thus, query optimizers try hard to "prove" that the duplicate eliminaton step is not necessary.

Projection (5)

- If RA is used to simulate SQL, the format of the projection list is often generalized:
 - Attribute renaming:

$$\pi_{B_1\leftarrow A_{i_1},\ldots,B_k\leftarrow A_{i_k}}(R)$$
 .

Computations (e.g., string concatenation via || or arithmetics via +,-,...) to derive the value in **new columns**, e.g.:

 $\pi_{\text{SID,NAME} \leftarrow \text{FIRST} ||}$, || LAST (STUDENTS).

Projection (6)

• $\pi_{A_1,\ldots,A_k}(R)$ corresponds to the SQL query:

SELECT DISTINCT
$$A_1, \ldots, A_k$$

FROM R

• $\pi_{B_1 \leftarrow A_1,...,B_k \leftarrow A_k}(R)$ is equivalent to the SQL query:

SELECT DISTINCT A_1 [AS] B_1, \ldots, A_k [AS] B_k FROM R

Selection vs. Projection



Combining Operations (1)

- Since the result of any relational algebra operation is a relation again, this intermediate result may be the input of a subsequent RA operation.
- Example: retrieve the exercises solved by student with ID 102:

 $\pi_{ ext{CAT,ENO}}(\sigma_{ ext{SID}=102}(ext{RESULTS}))$.

• We can think of the intermediate result to be stored in a **named temporary relation** (or as a macro definition):

$$ext{S102} \leftarrow \sigma_{ ext{SID}=102}(ext{RESULTS}); \\ \pi_{ ext{CAT,ENO}}(ext{S102})$$

Combining Operations (2)

• Composite RA expressions are typically depicted as **operator trees**:



• In these trees, computation proceeds **bottom-up**. The evaluation order of sibling branches is not pre-determined.

Combining Operations (3)

• SQL-92 permits the **nesting of queries** (the result of a SQL query may be used in a place of a relation name):

Nested SQL Query		
SELECT DISTINCT	CAT, ENG	D
FROM	(SELECT	*
	FROM	RESULTS
	WHERE	SID = 102) AS S102

• Note that this is *not* the typical style of SQL querying.

Combining Operations (4)

• Instead, a single SQL query is equivalent to an RA operator tree containing σ , π , and (multiple) × (see below):

SELECT-FROM	-WHERE Block	
SELECT DI	STINCT CAT, ENO	
FROM	RESULTS	
WHERE	SID = 102	

• Really complex queries may be constructed step-by-step (using SQL's **view** mechanism), S102 may be used like a relation:

```
SQL View Definition

CREATE VIEW S102

AS SELECT *

FROM RESULTS

WHERE SID = 102
```

Relational Algebra



4. Outer Join

Cartesian Product (1)

- In general, queries need to combine information from **several tables.**
- In RA, such queries are formulated using $\times,$ the Cartesian product.

Cartesian Product

The **Cartesian product** $R \times S$ of two relations R, S is computed by concatenating each tuple $t \in R$ with each tuple $u \in S$.

Cartesian Product (2)



• Since attribute names must be unique within a tuple, the Cartesian product may only be applied if R, S **do not share any attribute names.** (This is no real restriction because we have π .)

Cartesian Product (3)

• If $t = (A_1 : a_1, ..., A_n : a_n)$ and $u = (B_1 : b_1, ..., B_m : b_m)$, then $t \circ u = (A_1 : a_1, ..., A_n : a_n, B_1 : b_1, ..., B_m : b_m)$.

Cartesian Product: Nested Loops

```
create a new temporary relation T;
foreach t \in R do
foreach u \in S do
insert t \circ u into T;
od
od
return T:
```

Cartesian Product and Renaming

• *R* × *S* may be computed by the equivalent SQL query (SQL does not impose the unique column name restriction, a column *A* of relation *R* may uniquely be identified by *R*.*A*):

Cartesian	Pro	oduct in SQL
SELECT	*	
FROM	R,	S

▷ In RA, this is often formalized by means of of a **renaming operator** $\rho_X(R)$. If $sch(R) = (A_1 : D_1, ..., A_n : D_n)$, then $\rho_X(R) \equiv \pi_{XA_1 \leftarrow A_1} \xrightarrow{XA_2 \leftarrow A_n} (R)$.

Join (1)

- The intermediate result generated by a Cartesian product may be quite large in general (|R| = n, |S| = m ⇒ |R × S| = n * m).
- Since the combination of **Cartesian product and selection** in queries is common, a special operator **join** has been introduced.

Join The **(theta-)join** $R \Join_{\theta} S$ between relations R, S is defined as $R \Join_{\theta} S \equiv \sigma_{\theta}(R \times S).$ The **join predicate** θ may refer to attribute names of Rand S.

Join (2)

$\rho_{S}(\text{STUDENTS}) \bowtie_{S.\text{SID}=R.\text{SID}} \rho_{R}(\text{RESULTS})$

S.SID	S.FIRST	S.LAST	S. EMAIL	R.SID	R.CAT	R.ENO	R.POINTS
101	Ann	Smith		101	Н	1	10
101	Ann	Smith		101	Н	2	8
101	Ann	Smith		101	М	1	12
102	Michael	Jones	(null)	102	Н	1	9
102	Michael	Jones	(null)	102	Н	2	9
102	Michael	Jones	(null)	102	М	1	10
103	Richard	Turner		103	Н	1	5
103	Richard	Turner		103	М	1	7

• Note: student Maria Brown does not appear in the join result.

Join (3)

• $R \Join_{\theta} S$ can be evaluated by "folding" the procedures for σ, \times :

```
Nested Loop Join
```

```
create a new temporary relation T:
foreach t \in R do
  foreach \mu \in S do
     if \theta(t \circ u) then
        insert t \circ u into T:
     fi
   od
od
return T:
```

Join (4)

• Join combines tuples from two relations **and** acts like a filter: tuples without join partner are removed.

Note: if the join is used to follow a foreign key relationship, then no tuples are filtered:

Join follows a foreign key relationship (dereference)

RESULTS $\bowtie_{\text{SID}=S.\text{SID}} \pi_{S.\text{SID}\leftarrow\text{SID},\text{FIRST},\text{LAST},\text{EMAIL}}(\text{STUDENTS})$

• There are join variants which act like **filters only:** left and right **semijoin** (\ltimes, \rtimes) :

$$R\ltimes_{ heta}S \quad \equiv \quad \pi_{sch(R)}(R\Join_{ heta}S)$$
 ,

or do not filter at all: outer-join (see below).

Natural Join

• The **natural join** provides another useful abbreviation ("RA macro").

In the natural join $R \bowtie S$, the join predicate θ is defined to be an **conjunctive equality comparison of attributes sharing the same name** in R, S.

Natural join handles the necessary attribute renaming and projection.

Natural Join

Assume R(A, B, C) and S(B, C, D). Then:

$$R \bowtie S = \pi_{A,B,C,D}(\sigma_{B=B' \land C=C'}(R \times \pi_{B' \leftarrow B,C' \leftarrow C,D}(S)))$$

(Note: shared columns occur once in the result.)

Joins in SQL (1)

• In SQL, $R \Join_{\theta} S$ is normally written as

Join in SQL	("class	ic" an	d SQL-92))
SELECT FROM WHERE	* R,S θ	or	SELECT FROM	* R JOIN S ON $ heta$

• Note: the left query is **exactly** the SQL equivalent of $\sigma_{\theta}(R \times S)$ we have seen before.

SQL is a **declarative language:** it is the task of the SQL optimizer to infer that this query may be evaluated using a join instead of a Cartesian product.

Algebraic Laws (1)

- A significant number **algebraic laws** hold for join which are heavily utilized by the query optimizer.
- Example: selection push-down.

If predicate φ refers to attributes in S only, then

$$\sigma_{\varphi}(R \bowtie S) \equiv R \bowtie \sigma_{\varphi}(S)$$

Selection push-down

Why is selection push-down considered one of the most significant algebraic optimizations?

• (Such effficiency considerations are the subject of "Datenbanken II.")

A Common Query Pattern (1)

• The following operator tree structure is **very** common:



 Join all tables needed to answer the query, (2) select the relevant tuples, (3) project away all irrelevant columns.

A Common Query Pattern (2)

• The select-project-join query

$$\pi_{A_1,\ldots,A_k}(\sigma_{\varphi}(R_1 \bowtie_{\theta_1} R_2 \bowtie_{\theta_2} \cdots \bowtie_{\theta_{n-1}} R_n))$$

has the obvious SQL equivalent

SELECT DISTINCT FROM	A_1,\ldots,A_k R_1,\ldots,R_n
WHERE	arphi
AND	$ heta_1$ and \cdots and $ heta_{n-1}$

• It is a common source of errors to forget a join condition: think of the scenario R(A, B), S(B, C), T(C, D) when attributes A, D are relevant for the query output.

	STUDENTS							RESULTS			
<u>SID</u>	FIF	lST	LAST	EM	AIL		<u>SID</u>	CAT	<u>ENO</u>	POINTS	
101	A	lnn	Smith	1			101	Н	1	10	
102	Micha	ael	Jones	s (nu	11)		101	Н	2	8	
103	Richa	ard	Turner	-			101	М	1	12	
104	Mar	ia	Brown				102	Н	1	9	
						,	102	Н	2	9	
	Е	XERC	ISES				102	М	1	10	
CAT	ENO		TOPTC	MAXPT			103	Н	1	5	
Н	1	Rel	Alg.	10			103	М	1	7	
н	2		SQL	10							
М	1		SQL	14							

Relational Algebra Quiz (Level: Novice)

Formulate equivalent queries in RA

- Print all homework results for Ann Smith (print exercise number and points).
- ② Who has got the maximum number of points for a homework? Print full name and homework number.
- ③ (Who has got the maximum number of points *for all* homework exercises?)

- Sometimes it is necessary to refer to more than one tuple of the same relation at the same time.
 - Example: "Who got more points than the student with ID 101 for any of the exercises?"
 - ▷ Two answer this query, we need to compare two tuples *t*, *u* of the relation RESULTS:
 - (1) tuple t corresponding to the student with ID 101,
 - ② tuple u, corresponding to the same exercise as the tuple t, in which u.POINTS > t.POINTS.

• This requires a generalization of the select-project-join query pattern, in which **two instances of the same relation are joined** (the attributes in at least one instances must be renamed first):

$$S := \varrho_X(\text{RESULTS}) \underset{X.\text{CAT}=Y.\text{CAT} \land X.\text{ENO}=Y.\text{ENO}}{\bowtie} \varrho_Y(\text{RESULTS})$$
$$\pi_{X.\text{SID}}(\sigma_{X.\text{POINTS}>Y.\text{POINTS} \land Y.\text{SID}=101})(S)$$

• Such joins are commonly referred to as **self joins.**

Relational Algebra

Overview

- 1. Introduction; Selection, Projection
- 2. Cartesian Product, Join
- 3. Set Operations
- 4. Outer Join

- Relations are sets (of tuples). The "usual" family of binary **set operations** can also be applied to relations.
- It is a requirement, that both input relations have the same schema.

Set Operations

The set operations of relational algebra are $R \cup S$, $R \cap S$, and R - S (union, intersection, difference).

Set Operations (2)



204

Set Operations (3)

• $R \cup S$ may be implemented as follows:

Union

```
create a new temporary relation T;
foreach t \in R do
  insert t into T:
od
foreach t \in S do
  insert t into T:
od
remove duplicates in T;
return T:
```

Set Operations (4)

• R - S may be implemented as follows:

Difference

```
create a new temporary relation T;
foreach t \in R do
  remove \leftarrow false:
  foreach u \in S do
     remove \leftarrow remove or (t = u);
  od
  if not(remove) then
     insert t into T:
  fi
od
return T:
```

Union (1)

• In RA queries, a typical application for \cup is case analysis.

Example: Grading

Union (2)

• In SQL, \cup is directly supported: keyword UNION.

UNION may be placed between two SELECT-FROM-WHERE blocks:

```
SOL'S UNION
  SELECT SID, 'A' AS GRADE
  FROM
        RESULTS
  WHERE CAT = 'M' AND ENO = '1' AND POINTS >= 12
UNION
  SELECT SID, 'B' AS GRADE
  FROM RESULTS
  WHERE CAT = 'M' AND ENO = '1'
  AND POINTS >= 10 AND POINTS < 12
UNION
```

Set Difference (1)

• **Note:** the RA operators *σ*, *π*, ×, ⋈, ∪ are **monotic** by definition, *e.g.*:

$$R\subseteq S \implies \sigma_arphi(R)\subseteq\sigma_arphi(S)$$
 .

- Then it follows that every query *Q* that exclusively uses the above operators behaves monotonically:
 - ▷ Let I_1 be a database state, and let $I_2 = I_1 \cup \{t\}$ (database state after insertion of tuple *t*).
 - ▷ Then every tuple u contained in the answer to Q in state I₁ is also contained in the anser to Q in state I₂.

Database insertion never invalidates a correct answer.

Set Difference (2)

- If we pose non-monotonic queries, e.g.,
 - ▷ "Which student has not solved any exercise?"
 - ▷ "Who got the most points for Homework 1?"
 - ▷ "Who has solved all exercises in the database?"

then it is obvious that $\sigma, \pi, \times, \bowtie, \cup$ are **not sufficient** to formulate the query. Such queries require **set difference** (-).

A non-monotonic query

"Which student has not solved any exercise? (Print full name (FIRST, LAST)." (Example database tables repeated on next slide.)

	STUDENTS							RESULTS				
<u>SID</u>	FIF	lST	LAST	[EMAI	IL.		<u>SID</u>	<u>CAT</u>	<u>ENO</u>	POINTS	
101	A	lnn	Smith	1		•		101	Н	1	10	
102	Micha	ael	Jones	3 (null)		101	Н	2	8	
103	Richa	ard	Turner	:		•		101	М	1	12	
104	Mar	ria	Brown	n				102	Н	1	9	
<u> </u>								102	Н	2	9	
	Е	XERC	ISES					102	М	1	10	
CAT	ENO		TOPTC	MAX	(PT			103	Н	1	5	
Н	1	Re			10			103	М	1	7	
н	2		SQL		10							
M	1		SQL		14							

Set Difference (3)

A correct solution?

 $\pi_{\text{FIRST,LAST}}(\text{STUDENTS} \bowtie_{\text{SID} \neq \text{SID2}} \pi_{\text{SID2} \leftarrow \text{SID}}(\text{RESULTS}))$

A correct solution?

 $\pi_{\text{SID,FIRST,LAST}}(\text{STUDENTS} - \pi_{\text{SID}}(\text{RESULTS}))$

Correct solution!

Set Operations and Complex Selections

• Note that the availability of ∪, - (and ∩) renders **complex** selection predicates superfluous:

Predicate Simplification Rules

$$\begin{array}{rcl}
\sigma_{\varphi_1 \land \varphi_2}(Q) & \stackrel{\rightarrow}{=} & \sigma_{\varphi_1}(Q) \cap \sigma_{\varphi_2}(Q) \\
\sigma_{\varphi_1 \lor \varphi_2}(Q) & = & \sigma_{\varphi_1}(Q) \cup \sigma_{\varphi_2}(Q) \\
\sigma_{\neg \varphi}(Q) & = & Q - \sigma_{\varphi}(Q)
\end{array}$$

RDBMS implement complex selection predicates anyway

Why?

Relational Algebra Quiz (Level: Intermediate)

• The "RA quiz" below refers to the Homework DB. Schema:

RESULTS (<u>SID</u> \rightarrow STUDENTS,(<u>CAT</u>, <u>ENO</u>) \rightarrow EXERCISES, POINTS) STUDENTS (<u>SID</u>,FIRST,LAST,EMAIL) EXERCISES (CAT,ENO,TOPIC,MAXPT)

Formulate equivalent queries in RA

① Who got the **most** points (of all students) for Homework 1?

2 Which students solved **all** exercises in the database?

Find RA expressions that translate between the two

Two alternative representations of the homework, midterm exam, and final totals of the students are:

RES	RESULTS_1					RESULTS_2			
STUDENT	H	М	F		STUDENT	CAT	PCT		
Jim Ford	95	60	75		Jim Ford	Н	95		
Ann Smith	80	90	95		Jim Ford	М	60		
					Jim Ford	F	75		
					Ann Smith	Н	80		
					Ann Smith	М	90		
					Ann Smith	F	95		

Summary



 Derived (and thus redundant) operations: Theta-Join ⋈_θ, Natural Join ⋈, Semi-Join ⋈, Renaming ϱ, and Intersection ∩.

Relational Algebra

Overview

- 1. Introduction; Selection, Projection
- 2. Cartesian Product, Join
- 3. Set Operations
- 4. Outer Join

Outer Join (1)

• Join (M) eliminates tuples without partner:



• The **left outer join** preserves all tuples in its **left** argument, even if a tuple does not team up with a partner in the join:

$$\begin{array}{c|cccc} A & B \\ \hline a_1 & b_1 \\ a_1 & b_2 \end{array} \xrightarrow{M} \begin{array}{c} B & C \\ \hline b_2 & c_2 \\ b_3 & c_3 \end{array} = \begin{array}{c} A & B & C \\ \hline a_1 & b_1 & (null) \\ a_2 & b_2 & c_2 \end{array}$$

Outer Join (2)

• The right outer join preserves all tuples in its right argument:



• The full outer join preserves all tuples in both arguments:

Outer Join (3)

• Example: Prepare a full homework results report, including those students who did not hand in any solution at all:

STUDENTS $\bowtie_{\text{SID}=\text{SID}'} \pi_{\text{SID}' \leftarrow \text{SID}, \text{ENO}, \text{POINTS}}(\sigma_{\text{CAT}='\text{H}'}(\text{RESULTS}))$						
SID	FIRST	LAST	EMAIL	SID'	ENO	POINTS
101	Ann	Smith		101	1	10
101	Ann	Smith		101	2	8
102	Michael	Jones	(null)	102	1	9
102	Michael	Jones	(null)	102	2	9
103	Richard	Turner		103	1	5
104	Maria	Brown		(null)	(null)	(null)