

# The Relational Algebra and Relational Calculus

5DV119 — Introduction to Database Management

Umeå University

Department of Computing Science

Stephen J. Hegner

hegner@cs.umu.se

<http://www.cs.umu.se/~hegner>

# The Roots of SQL

- It can scarcely be said that SQL has a clean and simple design.
- Rather, SQL is based upon the blending of many ideas, and has evolved over a long period of time.
- Nevertheless, SQL has its roots in two ideal query languages.

**Relational Algebra:** A *procedural language* grounded in basic operations on relations.

- Widely used in algorithms for query optimization.

**Relational Calculus:** A *declarative language* grounded in first-order predicate logic.

- To understand better the capabilities and limitations of SQL, as well as for other reasons, it is therefore useful to study these two languages.
- They are part of almost all basic courses on database management given by faculties of science and technology at the university level.

# Overview of the Relational Algebra

- The relational algebra is a *procedural* query language on relations.
- Its basic operations have one of the following forms:

Relation  $\longrightarrow$  Relation

Relation  $\times$  Relation  $\longrightarrow$  Relation

- It therefore provides a basic computational model of how queries in SQL may be evaluated by a DBMS.
- It is often used in the internal representation of queries for the query optimizer in real relational DBMSs.
- A basic knowledge of the relational algebra can thus be very helpful in understanding why certain query operations are very expensive (in terms of time and computational resources) relative to others.

# Overview of the Operations of the Relational Algebra

- The relational algebra is defined in terms of three kinds of operations on relations:

## Operations specific to relations:

Projection:  $\text{Relation} \rightarrow \text{Relation}$ : Trim some columns from a relation.

Selection:  $\text{Relation} \rightarrow \text{Relation}$ : Trim some rows from a relation.

Join:  $\text{Relation} \times \text{Relation} \rightarrow \text{Relation}$ :

Combine two relations by matching values.

## The three fundamental set-theoretic operations:

all  $\text{Relation} \times \text{Relation} \rightarrow \text{Relation}$

Union:  $X \cup Y =$  all elements in either  $X$  or  $Y$ .

Intersection:  $X \cap Y =$  all elements in both  $X$  and  $Y$ .

Difference:  $X \setminus Y$  or  $X - Y =$  all elements in  $X$  which are not in  $Y$ .

## A special operation of the form $\text{Relation} \rightarrow \text{Relation}$ :

Attribute renaming: Change the names of some attributes of a relation.

# Projection

- The projection operation takes a “vertical” slice of a relation by dropping some columns while retaining others.
- The *projection* operator is represented by the lowercase Greek letter  $\pi$ , with the subscript identifying the columns to be retained.

$$\pi_{\{A_1, A_2, \dots, A_k\}}(R)$$

- The semantics of this expression are exactly those of the following SQL query.

```
SELECT DISTINCT A1, A2, . . . , Ak
FROM R;
```

- This is a formal operation on sets; duplicates are not part of the model.
- Often, the set brackets are dropped in the subscript.

$$\pi_{A_1, A_2, \dots, A_k}(R)$$

- If the attribute names are single letters, even the commas are sometimes dropped.

$$\pi_{A_1 A_2 \dots A_k}(R)$$

# Selection

- The selection operation takes a “horizontal” slice of a relation by dropping some rows while retaining others.
- The *selection* operator is represented by the lowercase Greek letter  $\sigma$ , with the subscript containing an expression which identifies the rows to be retained.

$$\sigma_{\varphi}(R)$$

- The semantics of this expression are exactly those of the following SQL query.

```
SELECT DISTINCT *  
FROM R  
WHERE  $\varphi$ ;
```

- The expression  $\varphi$  is often written in a more formal, logical style than that used by SQL.

Example:

$$\sigma_{((DNo=5)\wedge(Salary\geq 30000))}(R)$$

# Combining Expressions in the Relational Algebra

- The operations in the relational algebra themselves produce relations as results.
- Therefore, they may be composed.

**Example:**  $\pi_{A_1, A_2, \dots, A_k}(\sigma_{\varphi}(R))$  has the same meaning as

```
SELECT DISTINCT A1, A2, . . . , Ak
FROM R
WHERE  $\varphi$ ;
```

- Typing rules must be observed, since it is the composition of two distinct operations.

**Example:** While

$$\pi_{\text{LName, SSN}}(\sigma_{\text{Salary} \geq 30000}(\text{Employee}))$$

makes perfect sense,

$$\sigma_{\text{Salary} \geq 30000}(\pi_{\text{LName, SSN}}(\text{Employee}))$$

does not.

# Assignment Programs in the Relational Algebra

- Instead of composing operations in functional notation, queries in the relational algebra may be expressed as a sequence of assignment statements.

**Example:** The functional composition

$$\pi_{\text{LName,SSN}}(\sigma_{\text{Salary} \geq 30000}(\text{Employee}))$$

may also be expressed as the program of assignments

$$X_1 \leftarrow \sigma_{\text{Salary} \geq 30000}(\text{Employee})$$

$$X_2 \leftarrow \pi_{\text{LName,SSN}}(X_1)$$

with  $X_2$  as the final result.

- It is often easier to read and follow such sequence of assignments than to read and follow a complex functional composition.



# Join

- The join is a binary operation represented by the “bowtie” symbol  $\bowtie$ .
- It is basically the inner join of SQL.
- There are, however, a number of variants depending upon the subscript (or lack thereof).
- The expression

$$R_1 \bowtie_{\varphi} R_2$$

has the semantics of the SQL expression

```
SELECT *  
FROM R_1 JOIN R_2 ON ( $\varphi$ );
```

provided  $\varphi$  is represented in the correct way.

**Example:**  $\text{Employee} \bowtie_{(\text{DNo}=\text{DNumber})} \text{Department}$

has the meaning of

```
SELECT *  
FROM Employee JOIN Department ON (DNo=DNumber);
```

## Further Join Conventions

- Multiple conditions may be shown in various ways:

Employee  $\bowtie_{(DNo=DNumber)\wedge(Super\_SSN=Mgr\_SSN)}$  Department

Employee  $\bowtie_{\{(DNo=DNumber),(Super\_SSN=Mgr\_SSN)\}}$  Department

Employee  $\bowtie_{(DNo=DNumber),(Super\_SSN=Mgr\_SSN)}$  Department

- These all have the meaning of

```
SELECT *
FROM   Employee JOIN Department
        ON ((DNo=DNumber) AND (Super_SSN=Mgr_SSN));
```

- Other logical connectives:


Employee  $\bowtie_{(DNo=DNumber)\vee(Super\_SSN=Mgr\_SSN)}$  Department

has the meaning of

```
SELECT *
FROM   Employee JOIN Department
        ON ((DNo=DNumber) OR (Super_SSN=Mgr_SSN));
```

but is not a construction which occurs often in practice.

# Natural and Cross Joins

- The natural join is indicated by the absence of any subscripts on  $\bowtie$ .
-  The textbook uses the symbol  $*$  for natural join, although this notation is rather dated (and was used to denote inner join in early literature).
- Thus, the following two expressions are equivalent.

Department  $\bowtie$  Dept\_Locations  
Department  $*$  Dept\_Locations

with the same meaning as

```
SELECT *  
FROM   Department NATURAL JOIN Dept_Locations;
```

- Note that  $\bowtie_{\emptyset}$  is the *cross join*, with no matches. ( $\emptyset = \{\} =$  empty set.)
- Thus, Department  $\bowtie_{\emptyset}$  Dept\_Locations has the meaning of

```
SELECT *  
FROM   Department JOIN Dept_Locations ON (TRUE);
```

- This cross join (or *Cartesian product*) is also denoted  
Department  $\times$  Dept\_Locations.

# Theta Join

- Theta joins may be specified in the relational algebra in the obvious way.

**Query:** Find those employees who have an older dependent.

Employee  $\bowtie_{(SSN=ESSN) \wedge (Employee.BDate > Dependent.BDate)}$  Dependent

is equivalent to:

```
SELECT DISTINCT LName, FName, MInit, SSN
FROM Employee JOIN Dependent
ON ((SSN=ESSN)
AND (Employee.BDate > Dependent.BDate));
```

- which is equivalent to:

```
SELECT DISTINCT LName, FName, MInit, SSN
FROM Employee JOIN Dependent ON (SSN=ESSN)
WHERE (Employee.BDate > Dependent.BDate);
```

# Renaming

- Recall that it is sometimes necessary to have multiple copies of the same relation.

Query: Find the name of the supervisor of each employee.

```
SELECT E.LName, E.FName, E.MInit, S.LName, S.FName, S.MInit
FROM   Employee AS E JOIN Employee AS S
      ON (E.Super_SSN=S.SSN);
```

- In the relational algebra, there is a *rename* operation for this.
- There are two main formats:
  - $\rho_{R'}(R)$  returns a copy of  $R$  named  $R'$ , with the same attribute names.
  - $\rho_{R'}(A'_1, A'_2, \dots, A'_k)(R)$  returns a copy of  $R$  named  $R'$ , with the the attributes renamed to  $A'_1, A'_2, \dots, A'_k$ .
- Name qualifiers are used as in SQL.
- However, the original relation does not require a qualifier.

# Renaming Examples

**Query:** Find the name of the supervisor of each employee.

- The above query as a sequence of steps in the relational algebra, with  $X_3$  the answer, using each of the renaming conventions:

$$X_1 \leftarrow \rho_E(\text{Employee})$$

$$X_2 \leftarrow \rho_S(\text{Employee})$$

$$X_3 \leftarrow X_1 \bowtie_{(E.\text{Super\_SSN}=S.\text{SSN})} X_2$$

$$X_4 \leftarrow \pi_{E.\text{LName},E.\text{FName},E.\text{MInit},S.\text{LName},S.\text{FName},S.\text{MInit}}(X_3)$$

$$X_1 \leftarrow \rho_{S(\text{FName}',\text{MInit}'.\text{LName}',\text{SSN}',\text{BDate}',\text{Address}',\text{Sex}',\text{Salary}',\text{Super\_SSN}',\text{DNo}')}(Employee)$$

$$X_2 \leftarrow Employee \bowtie_{(\text{Super\_SSN}=\text{SSN}')} X_1$$

$$X_3 \leftarrow \pi_{\text{LName},\text{FName},\text{MInit},\text{LName}',\text{FName}',\text{MInit}'}(X_2)$$

## Another Renaming Example

**Query:** Find the Name and SSN of those employees who work on exactly one project.

- The above query as a sequence of steps in the relational algebra, with  $X_7$  the answer:

$X_1 \leftarrow \rho_W(\text{Works\_On})$  -- Copy of Works\_On

$X_2 \leftarrow \text{Works\_On} \bowtie_{(PNo \neq W.PNo) \wedge (ESSN = W.ESSN)} X_1$

$X_3 \leftarrow \rho_{X_a(SSN)}(\pi_{ESSN}(X_2))$  -- Employees who work on  $> 1$  projects

$X_4 \leftarrow \pi_{SSN}(\text{Employee}) \setminus \rho_{X_b(SSN)}(\pi_{ESSN}(\text{Works\_On}))$  -- Employees who work on  $< 1$  projects

$X_5 \leftarrow \pi_{SSN}(\text{Employee}) \setminus (X_3 \cup X_4)$  -- Employees who work on  $= 1$  project

$X_6 \leftarrow X_5 \bowtie \text{Employee}$

$X_7 \leftarrow \pi_{LName, FName, MInit, SSN}(X_6)$  -- Add the names

# Set Operations

- The following set operations are considered part of the relational algebra:

**Union:**  $X \cup Y =$  all elements in either  $X$  or  $Y$ .

**Intersection:**  $X \cap Y =$  all elements in both  $X$  and  $Y$ .

**Difference:**  $X \setminus Y$  or  $X - Y =$  all elements in  $X$  which are not in  $Y$ .

- They may only be applied when the elements in each set are of the same type.
  - If they are tuples, they have the same number of columns.
  - The attributes for matching columns must be of the same type.



## Recall Division in SQL

- The division operation has already been seen in the following SQL example:

**Query:** Find all employees who work on every project which Alicia Zeyala (999887777) also works on. Exclude Alicia herself.

**Recall the strategy:** Find all employees E for which there is no project P which Alicia works on but E does not work on.

```
SELECT DISTINCT LName, FName, MInit, SSN
FROM Employee JOIN Works_On ON (SSN=ESSN)
WHERE NOT EXISTS (SELECT PNo
                  FROM Works_On
                  WHERE (ESSN='999887777')
                  EXCEPT (SELECT PNo
                            FROM Works_On
                            WHERE (SSN=ESSN)))
AND (SSN<>'999887777');
```

- This operation may be formalized within the relational algebra.

# Formalization of Division via Example

| Works_On    |            |
|-------------|------------|
| <u>ESSN</u> | <u>PNo</u> |

| PList      |
|------------|
| <u>PNo</u> |

- Consider the schema as shown to the right.

**Query:** Find the SSNs of those employees in Works\_On who work on every project in PList.

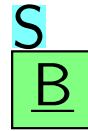
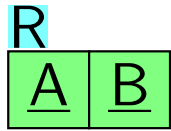
- Here is an assignment program in the relational algebra which provides a solution:

$X_1 \leftarrow \pi_{\text{ESSN}}(\text{Works\_On})$     -- Workers: employees who work on some project  
 $X_2 \leftarrow X_1 \times \text{PList}$     -- Every worker works on every project in PList  
 $X_3 \leftarrow X_2 \setminus \text{Works\_On}$     -- The “Does\_Not\_Work\_On” relation  
 $X_4 \leftarrow \pi_{\text{ESSN}}(X_3)$     -- Workers who do not work on some project in PList  
 $X_5 \leftarrow X_1 \setminus X_4$     -- Employees who work on every project in PList

- As a single expression:

$$\pi_{\text{ESSN}}(\text{Works\_On}) \setminus (\pi_{\text{ESSN}}(\pi_{\text{ESSN}}(\text{Works\_On}) \times \text{PList}) \setminus \text{Works\_On})$$

# Formalization of Division



**Query:** Find the  $A$ 's in  $R$  which are associated with every  $B$  in  $S$ .

- Here is an assignment program in the relational algebra which provides a solution:

$$\begin{aligned} X_1 &\leftarrow \pi_A(R) && \text{-- all } A\text{'s} \\ X_2 &\leftarrow X_1 \bowtie S && \text{-- } A \times B \\ X_3 &\leftarrow X_2 \setminus R && \text{-- } (A \times B) \setminus R \\ X_4 &\leftarrow \pi_A(X_3) && \text{-- } A\text{'s not associated with some } B \\ X_5 &\leftarrow X_1 \setminus X_4 && \text{-- } A\text{'s associated with every } B \end{aligned}$$

- As a single expression:

$$\pi_A(R) \setminus (\pi_A(\pi_A(R) \bowtie S) \setminus R)$$

- This division is written  $R \div S$ .
- This extends easily to  $R[\mathbf{A}]$ ,  $S[\mathbf{B}]$ , with sets  $\mathbf{A}$ ,  $\mathbf{B}$  of attributes satisfying  $\mathbf{B} \subseteq \mathbf{A}$ .

# Additional Operations of the Relational Algebra

- Many additional operations may be added to the relational algebra to make it as powerful as SQL, including:
  - Aggregation and grouping operators
  - Outer join
  - Recursive closure operations
- These are relatively straightforward to define, but will not be pursued further in this course.

# Declarative Query Languages

**Procedural:** A query language is *procedural* if it indicates explicitly how to compute its result.

- The relational algebra is a procedural query language.

**Declarative:** A query language is *declarative* if it indicates *what* to compute without requiring any indication of *how*.

- A great advantage of the relational model of data is that it admits a fully declarative query language.
- This means that the query language may be decoupled completely from the procedural model of computation.
  - This is particularly important for the support of non-technical users.
- For other data models, including object-oriented models, such a decoupling is difficult, if possible at all.
- The declarative query language for the relational model is called the *relational calculus*, and will be examined briefly.

# Propositional Logic

- The relational calculus is based upon first-order mathematical logic, which in turn is based upon propositional logic.
- Familiarity with propositional logic is assumed, including:

Connectives:  $\vee, \wedge, \neg, \Rightarrow$ .

- $(A \Rightarrow B)$  is *defined* to mean  $((\neg A) \vee B)$ .

Well-formed formulas (WFFs):  $(A \wedge ((\neg B) \vee C) \Rightarrow (D \vee E))$

DeMorgan's Laws:  $(\neg(A \wedge B)) = ((\neg A) \vee (\neg B))$   
 $(\neg(A \vee B)) = ((\neg A) \wedge (\neg B))$

# The Tuple Relational Calculus

- The specific relational calculus presented here is called the *tuple relational calculus*.

**Tuple variables:** The tuple relational calculus works with *tuple variables*.

- Each tuple variable has a *type* which is one of the relations in the schema.
  - $R(t)$  declares tuple  $t$  to be of type  $R$ .

**Example:** Employee( $e$ ).

- The value for a specific attribute is retrieved using standard notation.
  - $t.A$  retrieves the  $A$ -value of tuple variable  $t$ .

**Example:**  $e$ .Salary.

- Call an expression such as  $t.A$  a *tuple-field variable*.
- For those familiar with first-order predicate logic, each  $t.A$  corresponds (roughly) to a variable.

# The Tuple Relational Calculus — 2

**Quantifiers:** Quantifiers are used in expressions in the calculus.

$\forall$ : For all.

$\exists$ : There exists.

- Queries are of the form

$$\{t_1.A_1, t_2.A_2, \dots, t_k.A_k \mid \varphi\}$$

in which:

- Each  $t_i.A_i$  is a tuple-field variable.
- $\varphi$  is a logical formula in which exactly the elements of  $\{t_1, t_2, \dots, t_k\}$  are *free* (not within the scope of any quantifier).
- Rather than present a long formal syntax of well formedness, a number of examples will be used to illustrate the various constructions and techniques.



## Examples in the Tuple Relational Calculus

**Query:** Find the name and SSN of those employees who work on some project.

$$\{e.LName, e.FName, e.MInit, e.SSN \mid Employee(e) \\ \wedge (\exists w)(Works\_On(w) \wedge (e.SSN = w.ESSN))\}$$

**Query:** Find the name and SSN of those employees who work on the ProductX project.

$$\{e.LName, e.FName, e.MInit, e.SSN \mid Employee(e) \\ \wedge (\exists w)(\exists p)(Works\_On(w) \wedge Project(p) \wedge (p.PName = 'ProductX') \\ \wedge (p.PNumber = w.PNo) \wedge (e.SSN = w.ESSN))\}$$

**Query:** Find the name and SSN of those employees who work on every project.

$$\{e.LName, e.FName, e.MInit, e.SSN \mid Employee(e) \\ \wedge (\forall p)(Project(p) \Rightarrow \\ (\exists w)(Works\_On(w) \wedge (e.SSN = w.ESSN) \wedge (p.PNumber = w.PNo)))\}$$

- Note how easy and natural division is in the tuple relational calculus!

## Examples in the Tuple Relational Calculus — 2

**Query:** Find the name and SSN of those employees who work on exactly one project.

$$\begin{aligned} & \{e.Lname, e.FName, e.MInit, e.SSN \mid \text{Employee}(e) \\ & \quad \wedge (\exists w)(\text{Works\_On}(w) \wedge (e.SSN = w.ESSN)) \\ & \quad \wedge (\forall w_1)(\forall w_2)((\text{Works\_On}(w_1) \wedge \text{Works\_On}(w_2) \wedge (w_1.ESSN = w_2.ESSN) \\ & \quad \quad \wedge (w_1.ESSN = e.SSN)) \Rightarrow (w_1.PNo = w_2.PNo))\} \end{aligned}$$

**Query:** Find the name and SSN of those employees who do not work on any project.

$$\begin{aligned} & \{e.Lname, e.FName, e.MInit, e.SSN \mid \text{Employee}(e) \\ & \quad \wedge (\neg(\exists w)(\text{Works\_On}(w) \wedge (e.SSN = w.ESSN)))\} \end{aligned}$$

or

$$\begin{aligned} & \{e.Lname, e.FName, e.MInit, e.SSN \mid \text{Employee}(e) \\ & \quad \wedge ((\forall w)(\text{Works\_On}(w) \Rightarrow (e.SSN \neq w.ESSN)))\} \end{aligned}$$

## Examples in the Tuple Relational Calculus — 3

**Query:** Find the name and SSN of those employees who work on at least two distinct projects.

$$\{e.Lname, e.FName, e.MInit, e.SSN \mid \text{Employee}(e) \\ \wedge (\exists w_1)(\exists w_2)((\text{Works\_On}(w_1) \wedge \text{Works\_On}(w_2) \wedge (w_1.ESSN = w_2.ESSN) \\ \wedge (w_1.ESSN = e.SSN)) \wedge (w_1.PNo \neq w_2.PNo))\}$$

**Query:** Find the name and SSN of those employees who work on exactly two distinct projects.

$$\{e.Lname, e.FName, e.MInit, e.SSN \mid \text{Employee}(e) \\ \wedge (\exists w_1)(\exists w_2)((\text{Works\_On}(w_1) \wedge \text{Works\_On}(w_2) \wedge (w_1.ESSN = w_2.ESSN) \\ \wedge (w_1.ESSN = e.SSN)) \wedge (w_1.PNo \neq w_2.PNo)) \\ \wedge (\forall w_1)(\forall w_2)(\forall w_3)(\text{Works\_On}(w_1) \wedge \text{Works\_On}(w_2) \wedge \text{Works\_On}(w_3) \\ \wedge (w_1.ESSN = w_2.ESSN) \wedge (w_1.ESSN = w_3.ESSN) \wedge (w_1.ESSN = e.SSN)) \\ \Rightarrow ((w_1.PNo = w_2.PNo) \vee (w_1.PNo = w_3.PNo) \vee (w_2.PNo = w_3.PNo))\}$$

# Remarks about Queries in the Relational Calculus

⚠  $(\neg\forall)$  and  $(\neg\exists)$  are ambiguous and incorrect, and should never be used.

**Question:** What does  $(\forall e)(\neg\exists w)$  mean?

- Write  $(\forall e)\neg(\exists w)$  if that is what is meant.
- Recall that negation “flips” quantifiers.
  - $\neg((\forall x)(\varphi))$  is equivalent to  $(\exists x)((\neg\varphi))$ .
  - Think about it for simple examples.
  - Similarly  $\neg((\exists x)(\varphi))$  is equivalent to  $(\forall x)((\neg\varphi))$ .
- Keep in mind that  $\varphi_1 \Rightarrow \varphi_2$  is defined to mean  $(\neg\varphi_1)\vee\varphi_2$ .
- The value of a variable must always be defined in one of two ways.
  - By nature of lying within the scope of a quantifier
    - $(\forall e)(\exists w)(e \text{ and } w \text{ are bound here.})$ .
  - By nature of being in the argument list of a query.
    - $e$  and  $w$  arguments in  $\{e.A, w.B \mid \text{Employee}(e)\wedge\text{Works\_On}(w)\wedge\langle\text{some formula}\rangle\}$ .
- The type of each variable in the argument list must be defined in the formula of the query.

# The Expressive Power of the Algebra and Calculus

- A major, nontrivial result is the following:

**Theorem:** The relational algebra and the tuple relational calculus have the same expressive power.  $\square$

- This means that there is no loss of expressive power in using an entirely declarative language for querying relational databases.
- There has been substantial debate, with supporting research for both sides, for the relative merits of declarative query languages versus procedural query languages.
- In any case, SQL is blend of the two, with choices made for historical rather than scientific reasons.
- Still, it is very useful to be aware of the distinction between these two flavors of query expression.

# Safe Queries

**Theorem:** The relational algebra and the tuple relational calculus have the same expressive power.  $\square$

- There is one restriction which must be imposed for this result to hold: the queries must be *safe*.
- Roughly speaking, a query is safe if it can only return answers whose attribute values occur in the database being queried.

**Example of an unsafe query:** Give the set of all numbers which are not the salary of some employee.

- A query in the relational algebra is always safe.
- A query in the tuple relational calculus is guaranteed to be safe if every tuple variable in the argument list is bound to a type.
  - This is guaranteed in the formalism which has been developed here.
- Unsafe queries can arise in an alternative called the *domain relational calculus*, which is essentially standard first-order logic.
- The domain calculus will not be considered here.