

## UNIT-4

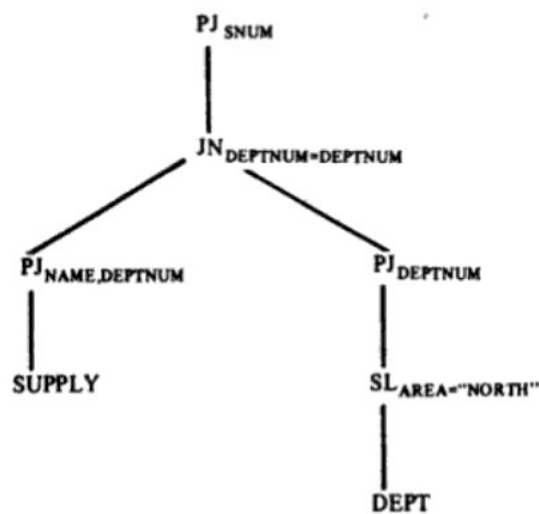
### 4.1 TRANSFORMING GLOBAL QUERIES INTO FRAGMENT QUERIES

#### A. Canonical Expression of a Fragment Query

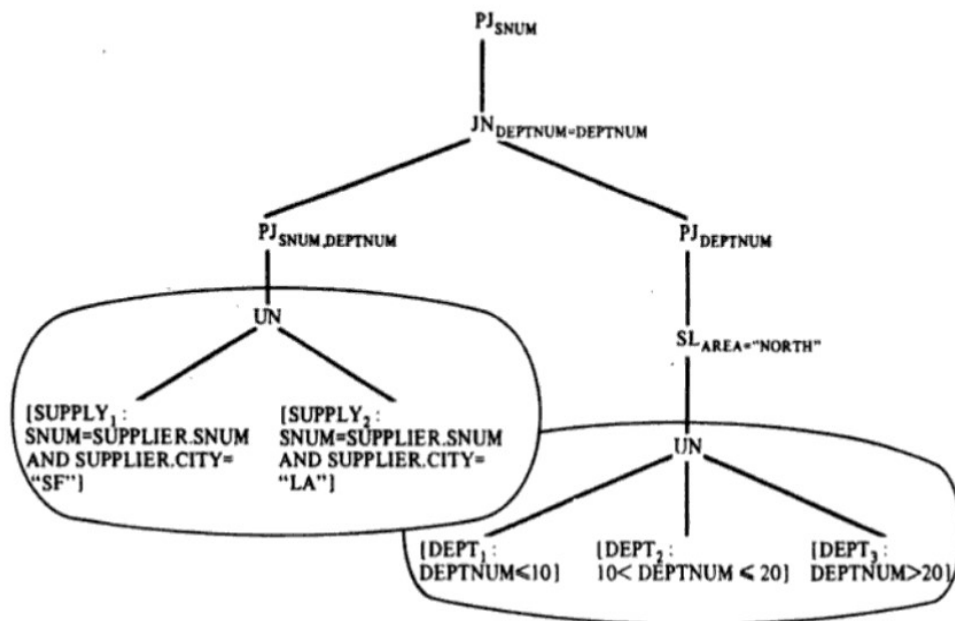
A canonical expression for the algebraic expression over the global schema is obtained as follows:

“Reconstruction of global relation is obtained from fragments. We map an operator tree on the global schema, in which collecting the fragments of all global relations into temporary relations and then applying to these temporaries the global query. “

**Example:** Operator tree of Q1 as follows



The following figure shows the canonical expression of above tree



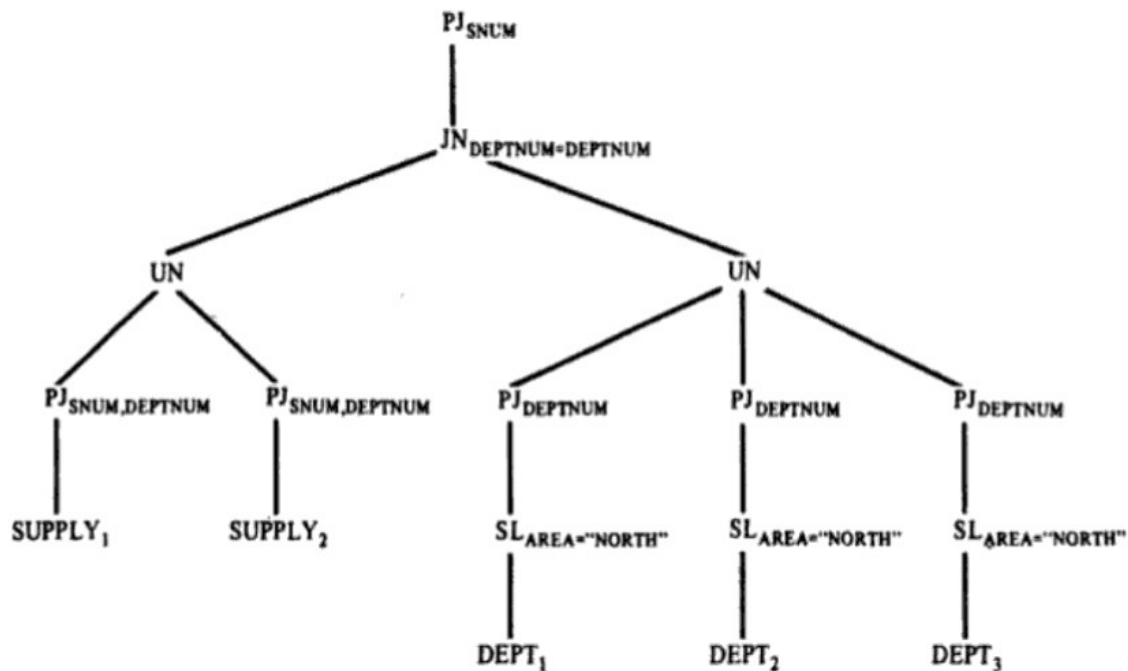
(a) Canonical form of query Q1.

We apply following criteria 1 and criteria 2 for further simplification of above canonical expression.

Criterion 1 : Use Idempotence of selection and projection to generate appropriate selections and projections for each operand relation

Criterion 2 : Push selections and projections down in the tree as far as possible.

After applying these two criterions, the resulting operator are described as follows:



(b) Pushing selections and projections down in the operator tree.

## B. Algebra of Qualified Relations

A qualified relation is a relation extended by a qualification and the qualification possessed by all the tuples of the relation.

Qualified relation can be denoted by a pair  $[R : q_R]$ , where  $R$  is a relation called the **body of the qualified relation** and  $q_R$  is a **predicate called the qualification of the qualified relation**. Horizontal fragments are typical examples of **qualified relations**.

The following rules define the result of applying the operations of relational algebra to qualified relations:

$$\text{Rule 1} \quad \mathbf{SL}_F[R : q_R] \Rightarrow [\mathbf{SL}_F R : F \text{ AND } q_R]$$

This rule states that the application of a selection  $\mathbf{SL}_F$  to a qualified relation  $[R : q_R]$  produces a qualified relation having the relation  $\mathbf{SL}_F R$  as its body and the predicate  $F \text{ AND } q_R$  as its qualification. The extension of the qualification to  $F \text{ AND } q_R$  after a selection reflects the fact that  $F$  holds on all the selected tuples as well as  $q_R$ .

$$\text{Rule 2} \quad \mathbf{PJ}_A[R : q_R] \Rightarrow [\mathbf{PJ}_A R : q_R]$$

The qualification of the result of a projection remains unchanged, even if the projection eliminates some of the attributes upon which the qualification was evaluated.

$$\text{Rule 3} \quad [R : q_R] \mathbf{CP} [S : q_S] \Rightarrow [R \mathbf{CP} S : q_R \text{ AND } q_S]$$

The extension of the qualification to  $q_R \text{ AND } q_S$  is rather intuitive; notice that the two qualifications apply to disjoint attributes of  $R \mathbf{CP} S$ .

$$\text{Rule 4} \quad [R : q_R] \mathbf{DF} [S : q_S] \Rightarrow [R \mathbf{DF} S : q_R]$$

However, the following problem arises with rule 4: consider the intersection operation of traditional relational algebra, defined as  $R \mathbf{IN} S = R \mathbf{DF} (R \mathbf{DF} S)$ . From the definition, it is easy to show that intersection is commutative; i.e.,

$$R \mathbf{IN} S = S \mathbf{IN} R.$$

If we apply the definition of extended algebra, we come to a surprising result:

$$\begin{aligned} [R : q_R] \mathbf{IN} [S : q_S] &\Rightarrow [R : q_R] \mathbf{DF} ([R : q_R] \mathbf{DF} [S : q_S]) \Rightarrow \\ [R : q_R] \mathbf{DF} [R \mathbf{DF} S : q_R] &\Rightarrow [R \mathbf{DF} (R \mathbf{DF} S) : q_R] \Rightarrow \\ [R \mathbf{IN} S : q_R] & \end{aligned} \quad (5.4a)$$

$$\begin{aligned} [S : q_S] \mathbf{IN} [R : q_R] &\Rightarrow [S : q_S] \mathbf{DF} ([S : q_S] \mathbf{DF} [R : q_R]) \Rightarrow \\ [S : q_S] \mathbf{DF} [S \mathbf{DF} R : q_S] &\Rightarrow [S \mathbf{DF} (S \mathbf{DF} R) : q_S] \Rightarrow \\ [S \mathbf{IN} R : q_S] & \end{aligned} \quad (5.4b)$$

In fact, the result that we would like to find is

$$[R \mathbf{IN} S : q_R \text{ AND } q_S]$$

$$\text{Rule 5} \quad [R : q_R] \mathbf{UN} [S : q_S] \Rightarrow [R \mathbf{UN} S : q_R \text{ OR } q_S]$$

The union is extended by taking the disjunction of qualifications.

Given the above rules, let us see how derived operations of algebra, such as join and semi-joins, are extended. We have two additional rules:

$$\text{Rule 6} \quad [R : q_R] \mathbf{JN}_F [S : q_S] \Rightarrow [R \mathbf{JN}_F S : q_R \text{ AND } q_S \text{ AND } F]$$

$$\text{Rule 7} \quad [R : q_R] \mathbf{SJ}_F [S : q_S] \Rightarrow [R \mathbf{SJ}_F S : q_R \text{ AND } q_S \text{ AND } F]$$

We give following two criterion 3 and 4 based on qualification of relations in further simplification.

**Criterion 3.** Push selections down to the leaves of the tree, and then apply them using the algebra of qualified relations; substitute the selection result with the empty relation if the qualification of the result is contradictory.

**Criterion 4.** Use the algebra of qualified relations to evaluate the qualification of operands of joins; substitute the subtree, including the join and its operands, with the empty relation if the qualification of the result of the join is contradictory.

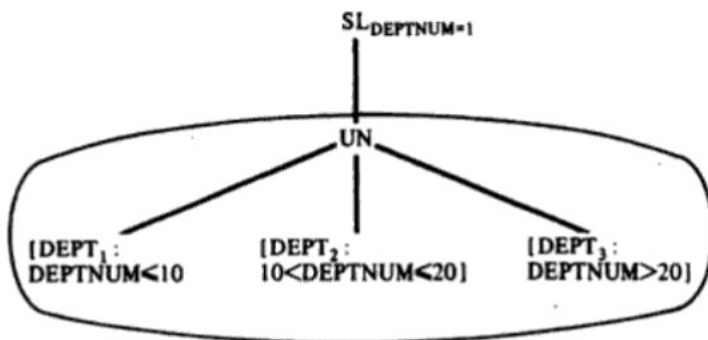
These two criteria are used in simplification of horizontally fragmented relations and joins between horizontally fragmented relations.

### C. Simplification of Horizontally Fragmented Relations

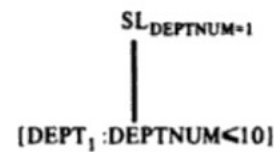
We show this kind of simplification with an example. Let us consider query Q3 on relation *DEPT*, which is horizontally fragmented:

$$Q3 : SL_{DEPTNUM=1} DEPT$$

The canonical form of the query is shown in Figure 5.5a. We push the selection toward the leaves by distributing it with respect to union. Then we apply selections according to criterion 3; the qualification of results of the selections on *DEPT*<sub>2</sub> and *DEPT*<sub>3</sub> is a contradiction;



(a) Canonical form of query Q3.



(b) Simplification of query Q3.

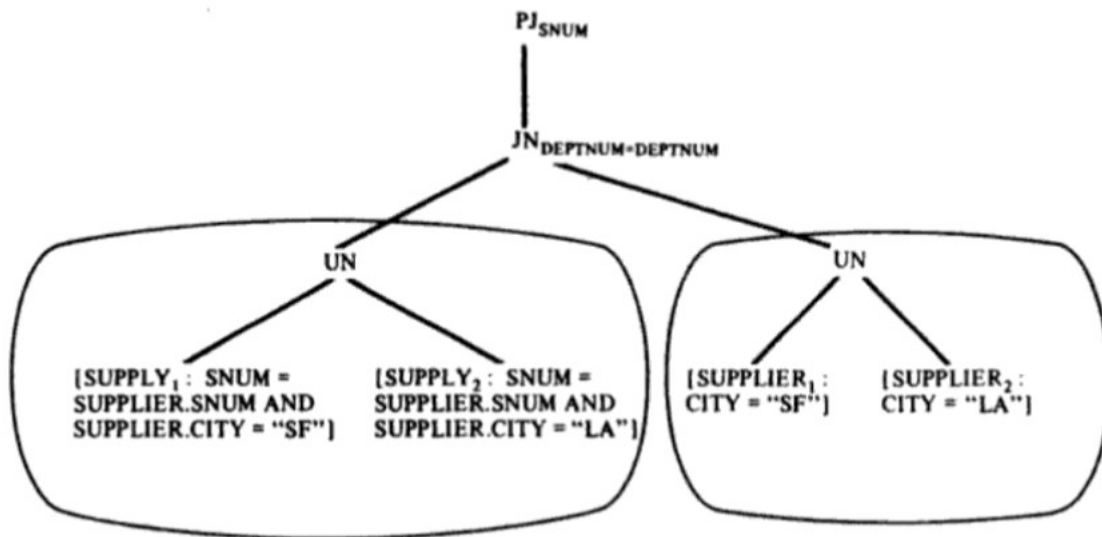
### D. Simplification of Joins between Horizontally Fragmented Relations

Let us consider, for simplicity, the join between two fragmented relations *R* and *S*. There are two distinct possibilities of joining them; the first one requires collecting all the fragments of *R* and *S* before performing the join. The second one consists of performing the join between fragments and then collecting all the results into the same result relation; we refer to this second case as “distributed join.”

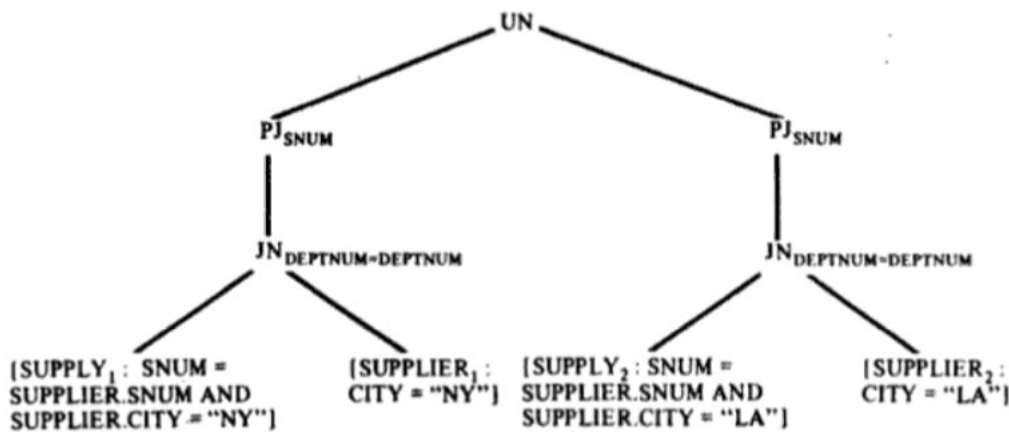
**Criterion 5.** In order to distribute joins which appear in the global query, unions (representing fragment collections) must be pushed up, beyond the joins that we want to distribute.

Building a join graph requires, then, applying criterion 5 (for distributing the join) followed by criterion 4

Figure 5.6a shows the canonical form of the query. We recall that the fragmentation of *SUPPLY* is derived from the fragmentation of *SUPPLIER*; i.e., each tuple of *SUPPLY* is stored either in fragment *SUPPLY*<sub>1</sub>, if it refers to a supplier of San Francisco, or in fragment *SUPPLY*<sub>2</sub>, if it refers to a supplier of Los Angeles. Applying criterion 5, we push the two unions up beyond the join; thus, we generate four joins between fragments. We then apply criterion 4, and we discover that two of them are intrinsically empty because their qualification is contradictory. The empty joins are those of *SUPPLIER*<sub>1</sub> (in "SF") with *SUPPLY*<sub>2</sub> (of "LA" suppliers), and likewise of *SUPPLIER*<sub>2</sub> (in "LA") with *SUPPLY*<sub>1</sub> (of "SF" suppliers). Thus, the operator tree reduces to that of Figure 5.6b. Assuming that fragments with the same index are placed at the same site (i.e., that data about *SUPPLY* are stored together with data about *SUPPLIERS*), this operator tree corresponds to an efficient way of evaluating the query, because each join is local to one site.



(a) Canonical form of query Q4.



(b) Distributed join for query Q4.

**Figure 5.8** Simplification of joins between horizontally fragmented relations.

## E. Using Inference for Further Simplifications

Assume that the following knowledge is available to the query optimizer:

1. The North area includes only departments 1 to 10.
2. Orders from departments 1 to 10 are all addressed to suppliers of San Francisco.

We use the above knowledge to “infer” contradictions that allow eliminating subexpressions.

- a. From 1 above, we can write the following implications:

$$AREA = \text{“North”} \Rightarrow \text{NOT } (10 < DEPTNUM \leq 20)$$

$$AREA = \text{“North”} \Rightarrow \text{NOT } (DEPTNUM > 20)$$

Using criterion 3, we apply the selection to fragments  $DEPT_1$ ,  $DEPT_2$ , and  $DEPT_3$  and evaluate the qualification of the results. By virtue of the above implications, two of them are contradictory. This allows us to eliminate the subexpressions for fragments  $DEPT_2$  and  $DEPT_3$ . Thus, the operator tree of Figure 5.4b reduces to that of Figure 5.7a.

- b. We then apply criterion 5 for distributing the join; in principle, we would need to join the subtree including  $DEPT_1$  with both subtrees including  $SUPPLY_1$  and  $SUPPLY_2$ . But from 1 above, we know that:

---

to join the subtree including  $DEPT_1$  with both subtrees including  $SUPPLY_1$  and  $SUPPLY_2$ . But from 1 above, we know that:

$$AREA = \text{“North”} \Rightarrow DEPTNUM \leq 10$$

and from 2 above we know that:

$$DEPTNUM \leq 10 \Rightarrow$$

$$\text{NOT } (SNUM = SUPPLIER.SNUM \text{ AND } SUPPLIER.CITY = \text{“LA”})$$

By applying criterion 4, it is possible to deduce that only the subtree including  $SUPPLY_1$  needs to be joined. The final operator tree for query Q1 is shown in Figure 5.7b.

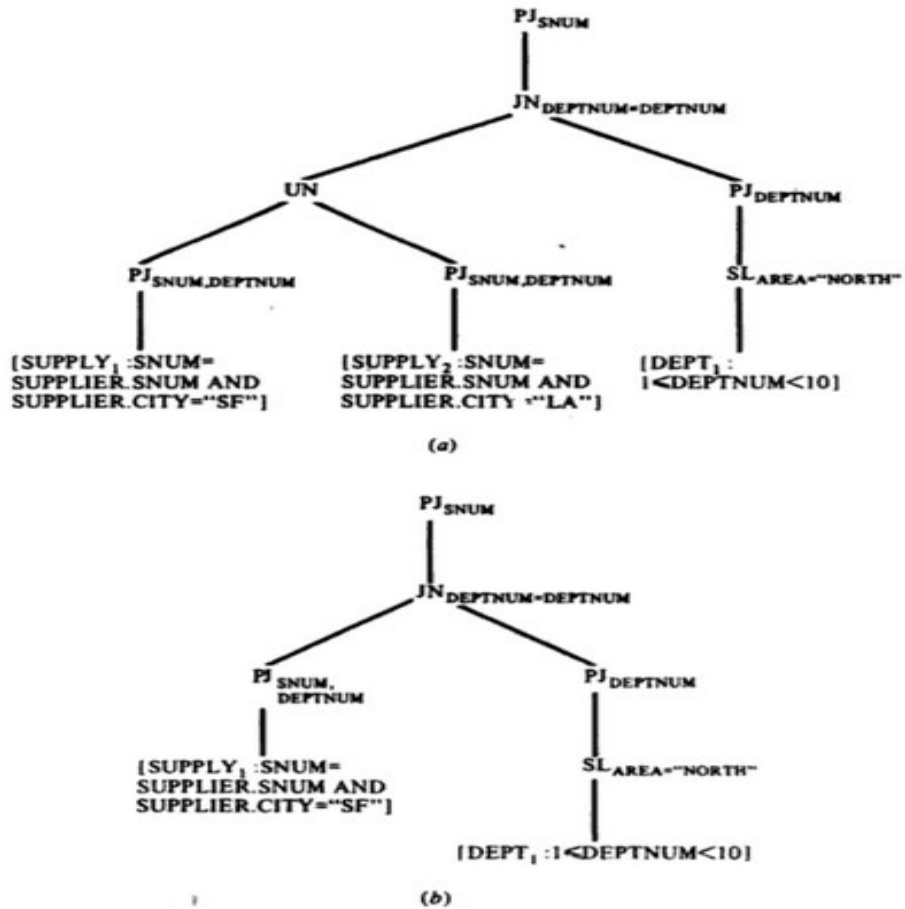


Figure 5.7 Simplification of an operator tree using inference.

### F. Simplification of Vertically Fragmented Relations

Simplification is to determine a proper subset of the fragments which is sufficient for answering the query, then eliminate all other fragments from the query expression, as well as the joins which are used in the inverse of the fragmentation schema for reconstructing global relations.

For the following query, simplification of vertical fragmented relations are as follows:

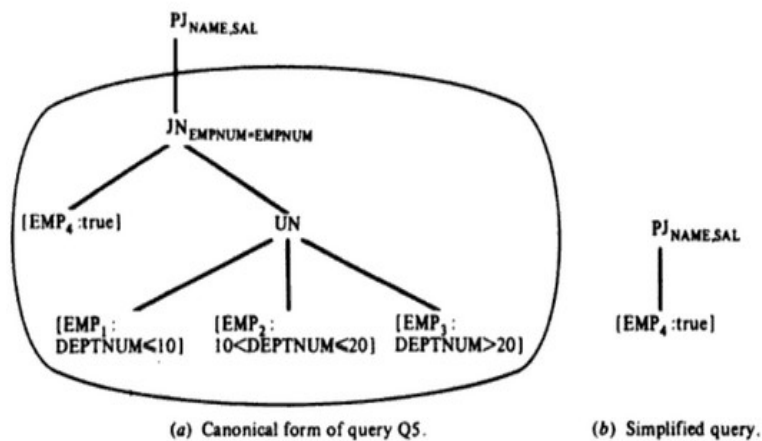


Figure 5.8 Simplification of vertically fragmented relations.

## 4.2 EQUIVALENCE TRANSFORMATIONS FOR QUERIES:

It is possible to transform most SQL queries into equivalent expressions of relational algebra.

### 4.2.1 Operator Tree of a query:

In order to have a more practical representation of queries, in which expression manipulation is easier to follow, we introduce **operator trees**. Let us consider query Q1, which requires the supplier number of suppliers that have issued a supply order in the North area of our company. Query Q1 corresponds to the following expression of the relational algebra:

$$Q1 : PJ_{SNUM} SL_{AREA="North"} (SUPPLY JN_{DEPTNUM=DEPTNUM} DEPT)$$

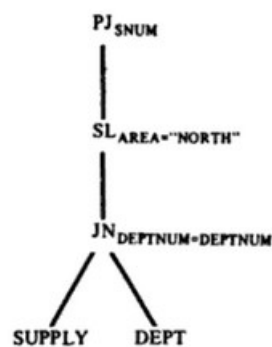


Figure 5.1 An operator tree for query Q1.

The operator tree of an expression of relational algebra can be regarded as the parse tree of the expression itself, assuming the following grammar:

$$\begin{aligned}
 R &\rightarrow \text{identifier} \\
 R &\rightarrow (R) \\
 R &\rightarrow un\_op R \\
 R &\rightarrow R bin\_op R \\
 un\_op &\rightarrow SL_F \mid PJ_A \\
 bin\_op &\rightarrow CP \mid UN \mid DF \mid JN_F \mid NJN_F \mid SJ_F \mid NSJ_F
 \end{aligned}$$

### 4.2.2 Equivalent transformations for the Relational Algebra

Equivalence transformations can be given systematically for small expressions, i.e., expressions of two or three operand relations. These transformations are classified into categories according to the type of the operators involved. Let  $U$  and  $B$  denote unary and binary algebraic operations, respectively. We have:

- **Commutativity of unary operations:**

$$U_1 U_2 R \leftrightarrow U_2 U_1 R$$



- **Commutativity** of operands of binary operations:

$$R B S \leftrightarrow S B R$$

- **Associativity** of binary operations:

$$R B (S B T) \leftrightarrow (R B S) B T$$

- **Idempotence** of unary operations:

$$U R \leftrightarrow U_1 U_2 R$$

- **Distributivity** of unary operations with respect to binary operations:

$$U(R B S) \rightarrow U(R) B U(S)$$

- **Factorization** of unary operations (this transformation is the inverse of distributivity):

$$U(R) B U(S) \rightarrow U(R B S)$$

The tables contain in each position a **validity indicator**. A validity indicator “Y” means that the property can always be applied; “N” means that it cannot be applied. For example, the validity indicator “Y” in the first row and first column of Table 5.1 means that the following transformation is correct

$$SL_{F1} SL_{F2} R \rightarrow SL_{F2} SL_{F1} R$$

where  $F1$  and  $F2$  are two generic selection specifications.

Validity indicators can also be “SNC,” specifying a condition which is necessary and sufficient for the application of the property. For example, the validity indicator  $SNC_1$  in the second row and first column of Table 5.1 means that the transformation

$$PJ_{A1} SL_{F2} R \rightarrow SL_{F2} PJ_{A1} R$$

**Table 5.1** Commutativity of unary operations

	$SL_{F2}$	$PJ_{A2}$
$SL_{F1}(*R)$ $\rightarrow *(SL_{F1}(R))$	Y	Y
$PJ_{A1}(*R)$ $\rightarrow *(PJ_{A1}(R))$	$SNC_1$	$SNC_2$

$$SNC_1 : Attr(F2) \subseteq A1$$

$$SNC_2 : A1 \equiv A2$$

**Table 5.2** Commutativity of operands and associativity of binary operations

	UN	DF	CP	$JN_F$	$SJ_F$
$R * S$ $\rightarrow S * R$	Y	N	Y	Y	N
$(R * S) * T$ $\rightarrow R * (S * T)$	Y	N	Y	$SNC_1$	N

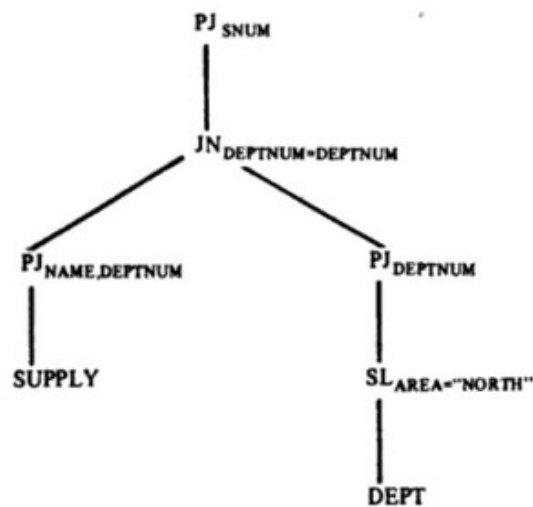
$$SNC_1 \text{ for } (R JN_{F1} S) JN_{F2} T \rightarrow R JN_{F1} (S JN_{F2} T) : Attr(F2) \subseteq Attr(S) \cup Attr(T)$$

**Table 5.3** Idempotence of unary operations

$PJ_A(R) \rightarrow PJ_{A1} PJ_{A2}(R)$	$SNC : A \equiv A1, A \subseteq A2$
$SL_F(R) \rightarrow SL_{F1} SL_{F2}(R)$	$SNC : F = F1 \wedge F2$

In nondistributed databases, general criteria have been given for applying equivalence transformations for the purpose of simplifying the execution of queries:

- Criterion 1.** Use idempotence of selection and projection to generate appropriate selections and projections for each operand relation.
- Criterion 2.** Push selections and projections down in the tree as far as possible.



**Figure 5.2** A modified operator tree for query Q1.

Figure 5.2 shows a modified operator tree for query Q1, in which the following transformations have been applied:

1. The selection is distributed with respect to the join; thus, the selection is applied directly to the *DEPT* relation.
2. Two new projection operations are generated and are distributed with respect to the join.

Notice that criterion 1 (use of idempotence) has been applied to the projection operation before distributing it; the selection operation, instead, has been directly moved to the *DEPT* relation.

### 4.2.3 Operator Graph and Determination of Common Sub-Expressions

An important issue in query expression is to discover its common sub-expressions; i.e., sub-expressions which appear more than once in the query; it shows that, saving the time and space if common sub-expressions are evaluated only once.

Ex:

Q2 :  $PJ_{EMP.NAME}((EMP JN_{DEPTNUM=DEPTNUM} SL_{MGRNUM=373} DEPT) DF$   
 $(SL_{SAL>35000} EMP JN_{DEPTNUM=DEPTNUM} SL_{MGRNUM=373} DEPT))$

In this, common subexpression is as follows:

$EMP \Join_{DEPTNUM=DEPTNUM} SL_{MGRNUM=373} DEPT$

Once common subexpressions are identified, we can use the following properties to further simplify an operator tree:

$R \text{ NJN } R \leftrightarrow R$

$R \text{ UN } R \leftrightarrow R$

$R \text{ DF } R \leftrightarrow \emptyset$

$R \text{ NJN } SL_F R \leftrightarrow SL_F R$

$R \text{ UN } SL_F R \leftrightarrow R$

$R \text{ DF } SL_F R \leftrightarrow SL_{\text{NOT } F} R$

$(SL_{F_1} R) \text{ NJN } (SL_{F_2} R) \leftrightarrow SL_{F_1 \text{ AND } F_2} R$

$(SL_{F_1} R) \text{ UN } (SL_{F_2} R) \leftrightarrow SL_{F_1 \text{ OR } F_2} R$

$(SL_{F_1} R) \text{ DF } (SL_{F_2} R) \leftrightarrow SL_{F_1 \text{ AND NOT } F_2} R$

### 4.3 DISTRIBUTED GROUPING AND AGGREGATE FUNCTION EVALUATION

Database applications having database access operations that can be expressed in query languages rather than specifying with relational algebra. Grouping tuples into disjoint subsets of relations and evaluating aggregate functions over them

We first give some examples, using SQL, of queries which use the above features.

Q6: Select AVG(QUAN)  
from SUPPLY  
where PNUM = "P1"

This query retrieves into a result relation having one attribute and one tuple the average quantity of supply orders for product "P1."

Q7: Select PNUM, SNUM, SUM(QUAN)  
from SUPPLY  
group by SNUM, PNUM

This query corresponds to partitioning the relation SUPPLY into groups having the same value of SNUM and PNUM (but different values of DEPTNUM and QUAN), evaluating for each such group the sum of the quantities, and retrieving SNUM, PNUM, and the sum of the quantities of each group into the result relation.

Q8: Select SNUM, PNUM, SUM(QUAN)  
from SUPPLY  
group by SNUM, PNUM  
having SUM(QUAN) > 300

#### Extensions with Relational Algebra

Relational algebra is extended with the following **group-by**  $GB_{G,AF} R$  such that:

- $G$  are the attributes which determine the grouping of  $R$ .
- $AF$  are aggregate functions to be evaluated on each group.
- $GB_{G,AF} R$  is a relation having:  
A relation schema made by the attributes of  $G$  and the aggregate functions of  $AF$ .

With the above operation, it is possible to write in algebra queries Q6, Q7, and Q8. We have

$$\begin{aligned} Q6 &: \mathbf{GB}_{\text{AVG}(\text{QUAN})} \mathbf{SL}_{\text{PNUM}=\text{"P1"}} \text{SUPPLY} \\ Q7 &: \mathbf{GB}_{\text{SNUM},\text{PNUM},\text{SUM}(\text{QUAN})} \text{SUPPLY} \\ Q8 &: \mathbf{SL}_{\text{SUM}(\text{QUANT})>300} \mathbf{GB}_{\text{SNUM},\text{PNUM},\text{SUM}(\text{QUAN})} \text{SUPPLY} \end{aligned}$$

Some comments are in order. The *G* part corresponds to the “group-by” clause, and the *AF* part corresponds to the aggregate functions whose computation is required. Typically in these queries attributes upon which grouping is made are also retrieved; thus, the attributes which appear in the “group-by” clause also appear in the “select” clause.

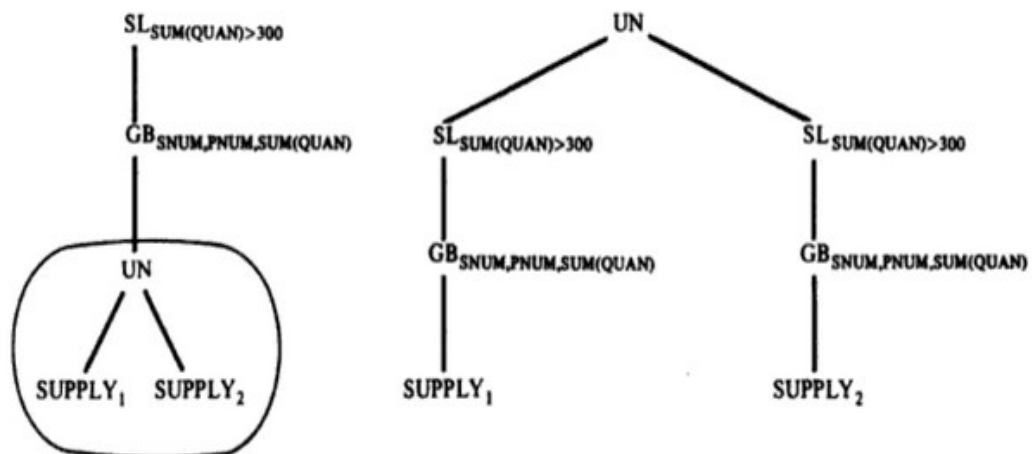
In query Q6, the *G* part is left empty, as the function is applied to all the tuples of *SUPPLY*. In query Q8, the usual selection operation is applied to the result of the *GB* operation; this selection corresponds to the “having” clause of SQL.

### Properties of Group By Operations:

In this section, we give an equivalence property for the new operation, and we discuss, in general, the possibility of evaluating aggregate functions in a distributed way.

The property in which we are interested is the distributivity of *GB* with respect to union:

$$\mathbf{GB}_{G,AF}(R_1 \cup R_2) \rightarrow (\mathbf{GB}_{G,AF}R_1) \cup (\mathbf{GB}_{G,AF}R_2), \quad \text{Val. ind.: } \text{SNC}$$



(a) Canonical form of query Q8.

(b) Distributed version of query Q8.

**Figure 5.10** A query with grouping and aggregate functions.

**Criterion 6.** In order to distribute grouping and aggregate function evaluations appearing in a global query, unions (representing fragment collections) must be pushed up, beyond the corresponding group-by operation.

An aggregate function for which it is possible to find the functions  $F_i$  and the expression  $E(F_i)$  is the function average

$$AVG(S) = \frac{SUM(SUM(S_1), SUM(S_2), \dots, SUM(S_n))}{SUM(COUNT(S_1), COUNT(S_2), \dots, COUNT(S_n))}$$

Similarly, we have

$$MIN(S) = MIN(MIN(S_1), MIN(S_2), \dots, MIN(S_n))$$

$$MAX(S) = MAX(MAX(S_1), MAX(S_2), \dots, MAX(S_n))$$

$$COUNT(S) = SUM(COUNT(S_1), COUNT(S_2), \dots, COUNT(S_n))$$

$$SUM(S) = SUM(SUM(S_1), SUM(S_2), \dots, SUM(S_n))$$

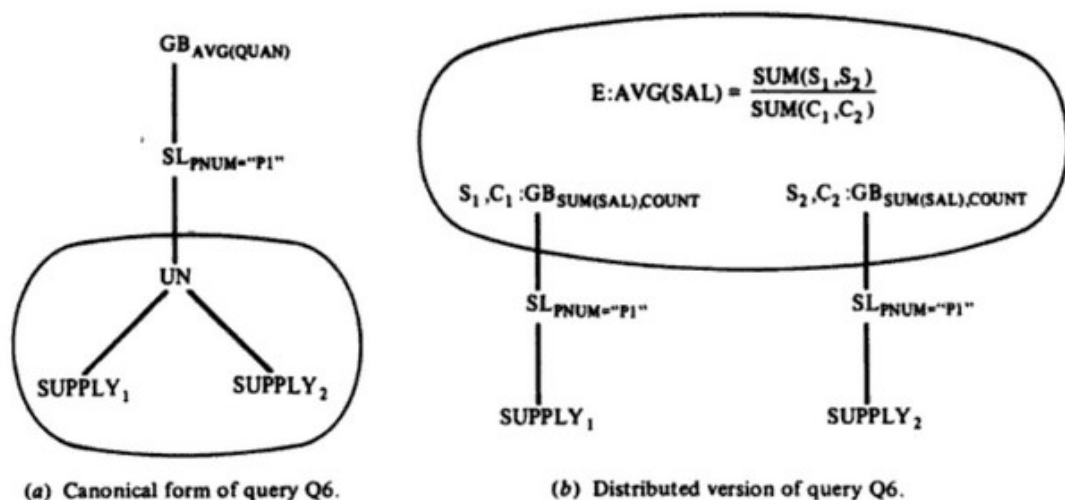


Figure 5.11 Distributed evaluation of aggregate functions.

#### 4.4 PARAMETRIC QUERIES

The selection criteria of queries include parameters whose values are not known when the query is compiled. When parametric queries are executed, the user provides values which are bound to parameters, thus parametric queries allow the repeated execution of queries for different values of the parameters and return different values at each execution.

##### 4.4.1 Simplification of Parametric queries and extension of algebra

As an example of a parametric query, let us consider query Q9, selecting tuples of the global relation *DEPT* having given department numbers. Let the selection on *DEPTNUM* be parametric:

$$Q9 : SL_{DEPTNUM=\$X \text{ OR } DEPTNUM=\$Y} SUPPLY$$

At run time, the actual values are assigned to the parameters  $\$X$  and  $\$Y$  by the program which issues the query.

The simplification of parametric queries has some distinguishing features. Consider the above query, whose canonical form is shown in Figure 5.12a.

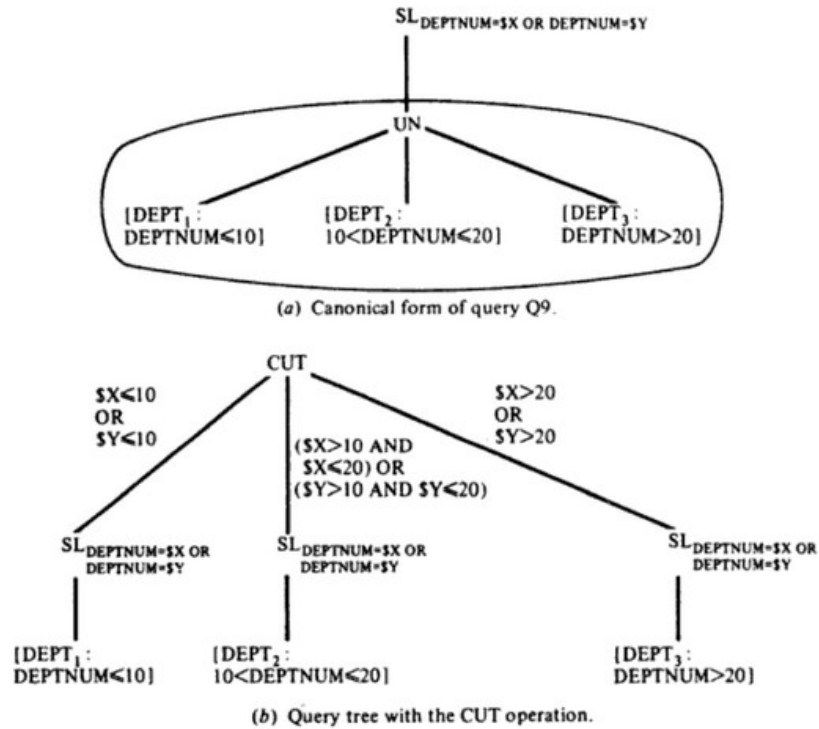


Figure 5.12 Use of CUT in a parametric query.

time, we don't know which of the fragments of the global relation *DEPT* will be addressed. However, we know that at most two of them will be involved in the query. At run time, when the query will be activated, it will also be possible to know, given the value of the parameters, which ones of the fragments will be involved in the query. This shows that part of the simplification of a parametric query can be done at compile time, but part of it needs to be done at run time.

Figure 5.12b shows the new operator tree for the parametric query. In this example, the selection is pushed below the union operation, applying criterion 2, and then the union is substituted by the **CUT** operation. The three formulas for the **CUT** operation (one for each operand of the union operation) are

$$\begin{aligned}
 F1 &: \$X \leq 10 \text{ OR } \$Y \leq 10 \\
 F2 &: (\$X > 10 \text{ AND } \$X \leq 20) \text{ OR } (\$Y > 10 \text{ AND } \$Y \leq 20) \\
 F3 &: \$X > 20 \text{ OR } \$Y > 20
 \end{aligned}$$

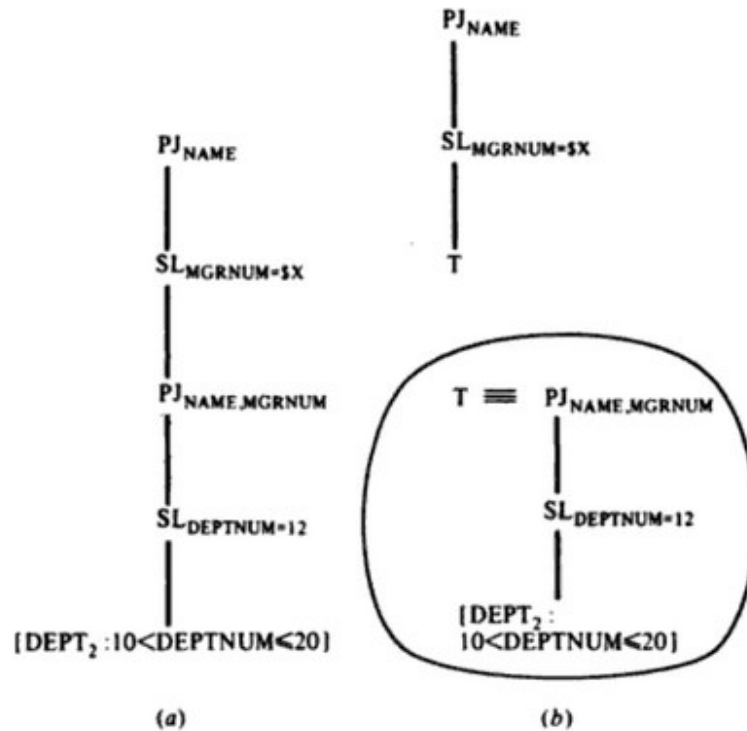
At run time, the **CUT** operation is the first one to be evaluated; the query, for  $\$X = 1$  and  $\$Y = 13$ , reduces to the first and second branch of the operator tree; for  $\$X = 1$  and  $\$Y = 5$ , the query reduces to the first branch of the tree.

#### 4.4.2 Using temporaries in Multiple Activations of Parametric Queries

For the case of repeated executions and lower the costs, it might be useful to build temporary relations at the site of origin of the query, which store a superset of the data required at each iteration.

$$Q10 : PJ_{NAME} SL_{MGRNUM=\$X \text{ AND } DEPTNUM=12} EMP$$

First, we modify the query so that the operator tree of Figure 5.13a is obtained. The general goal of this transformation is to insulate subtrees which are **not parametric**, as they will constitute the required temporaries.



**Figure 5.13** Use of temporary relations for parametric queries.

$PJ_{NAME, MGRNUM} SL_{DEPTNUM=12} EMP$

We then call  $T$  the temporary given by the above expression, and divide the operator tree into two pieces in correspondence to  $T$  (see Figure 5.13b).  $T$  represents the temporary relation that will be used for repeated executions of query Q10.