

## Query Optimization

In this chapter,<sup>1</sup> we will assume that the reader is already familiar with the strategies for query processing in relational DBMSs that we discussed in the previous chapter. The goal of query optimization is to select the best possible strategy for query evaluation. As we said before, the term *optimization* is a misnomer because the chosen execution plan may not always be the most optimal plan possible. The primary goal is to arrive at the most efficient and cost-effective plan using the available information about the schema and the content of relations involved, and to do so in a reasonable amount of time. Thus a proper way to describe **query optimization** would be that it is an activity conducted by a query optimizer in a DBMS to select the best available strategy for executing the query.

This chapter is organized as follows. In Section 19.1 we describe the notation for mapping of the queries from SQL into query trees and graphs. Most RDBMSs use an internal representation of the query as a tree. We present heuristics to transform the query into a more efficient equivalent form followed by a general procedure for applying those heuristics. In Section 19.2, we discuss the conversion of queries into execution plans. We discuss nested subquery optimization. We also present examples of query transformation in two cases: merging of views in Group By queries and transformation of Star Schema queries that arise in data warehouses. We also briefly discuss materialized views. Section 19.3 is devoted to a discussion of selectivity and result-size estimation and presents a cost-based approach to optimization. We revisit the information in the system catalog that we presented in Section 18.3.4 earlier and present histograms. Cost models for selection and join operation are presented in Sections 19.4 and 19.5. We discuss the join ordering problem, which is a critical one, in some detail in Section 19.5.3. Section 19.6 presents an example of cost-based optimization. Section 19.7 discusses some additional issues related to

<sup>1</sup>The substantial contribution of Rafi Ahmed to this chapter is appreciated.

query optimization. Section 19.8 is devoted to a discussion of query optimization in data warehouses. Section 19.9 gives an overview of query optimization in Oracle. Section 19.10 briefly discusses semantic query optimization. We end the chapter with a summary in Section 19.11.

## 19.1 Query Trees and Heuristics for Query Optimization

In this section, we discuss optimization techniques that apply heuristic rules to modify the internal representation of a query—which is usually in the form of a query tree or a query graph data structure—to improve its expected performance. The scanner and parser of an SQL query first generate a data structure that corresponds to an *initial query representation*, which is then optimized according to heuristic rules. This leads to an *optimized query representation*, which corresponds to the query execution strategy. Following that, a query execution plan is generated to execute groups of operations based on the access paths available on the files involved in the query.

One of the main **heuristic rules** is to apply SELECT and PROJECT operations *before* applying the JOIN or other binary operations, because the size of the file resulting from a binary operation—such as JOIN—is usually a multiplicative function of the sizes of the input files. The SELECT and PROJECT operations reduce the size of a file and hence should be applied *before* a join or other binary operation.

In Section 19.1.1, we reiterate the query tree and query graph notations that we introduced earlier in the context of relational algebra and calculus in Sections 8.3.5 and 8.6.5, respectively. These can be used as the basis for the data structures that are used for internal representation of queries. A *query tree* is used to represent a *relational algebra* or extended relational algebra expression, whereas a *query graph* is used to represent a *relational calculus expression*. Then, in Section 19.1.2, we show how heuristic optimization rules are applied to convert an initial query tree into an **equivalent query tree**, which represents a different relational algebra expression that is more efficient to execute but gives the same result as the original tree. We also discuss the equivalence of various relational algebra expressions. Finally, Section 19.1.3 discusses the generation of query execution plans.

### 19.1.1 Notation for Query Trees and Query Graphs

A **query tree** is a tree data structure that corresponds to an extended relational algebra expression. It represents the input relations of the query as *leaf nodes* of the tree, and it represents the relational algebra operations as internal nodes. An execution of the query tree consists of executing an internal node operation whenever its operands are available and then replacing that internal node by the relation that results from executing the operation. The order of execution of operations *starts at the leaf nodes*, which represents the input database relations for the query, and *ends at the root node*, which represents the final operation of the query. The execution terminates when the root node operation is executed and produces the result relation for the query.

(a)

 $\sigma_{P.Plocat}$ 

(b)

 $\pi_{P.Pnumb}$  $\sigma_{P.Dnum=}$ 

P

(c)

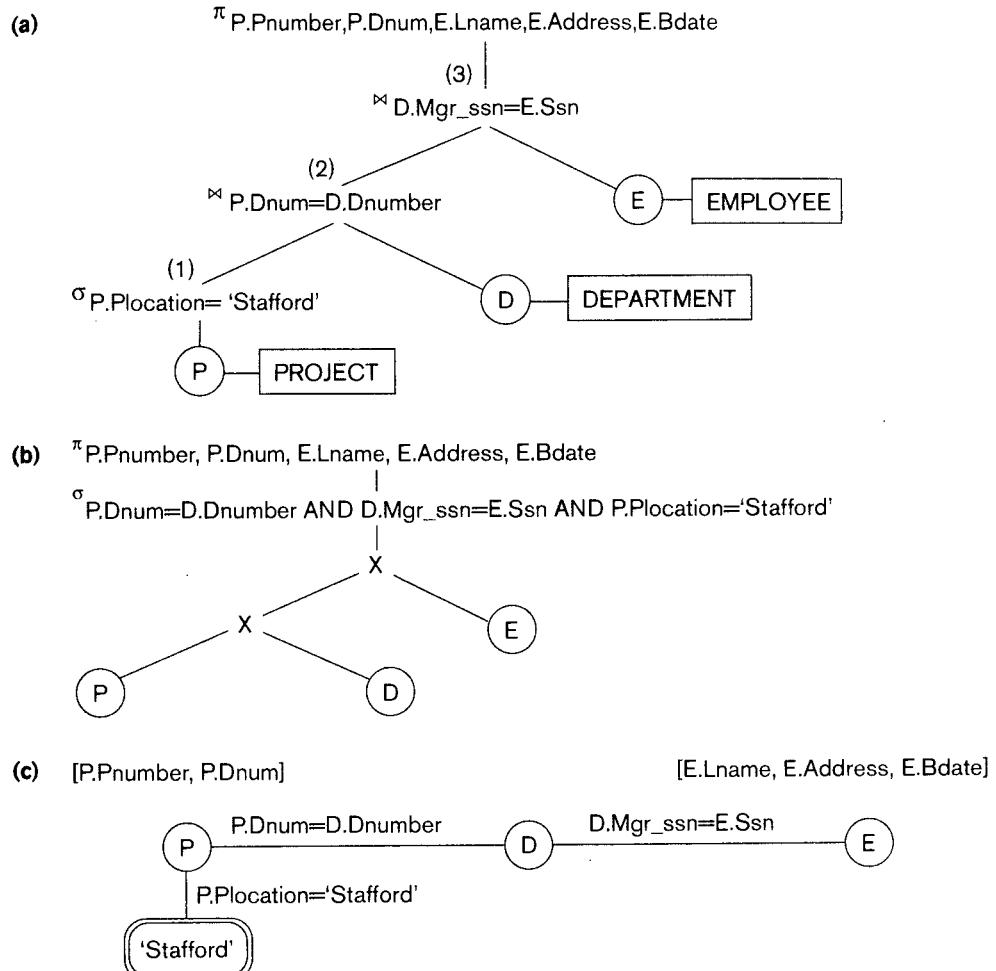
 $[P.Pnumbe$ 

'Sta

**Figure 19.1**  
Two query trees  
expression for Q

Figure 19.1(a)  
in Chapters 6  
number, the co  
name, address  
tional schema  
expression:

 $\pi_{Pnumber, D}$  $\bowtie_{Dnum}$



**Figure 19.1**  
Two query trees for the query Q2. (a) Query tree corresponding to the relational algebra expression for Q2. (b) Initial (canonical) query tree for SQL query Q2. (c) Query graph for Q2.

Figure 19.1(a) shows a query tree (the same as shown in Figure 6.9) for query Q2 in Chapters 6 to 8: For every project located in 'Stafford', retrieve the project number, the controlling department number, and the department manager's last name, address, and birthdate. This query is specified on the COMPANY relational schema in Figure 5.5 and corresponds to the following relational algebra expression:

$$\begin{aligned} \pi_{Pnumber, Dnum, Lname, Address, Bdate} & (((\sigma_{Plocation='Stafford'}(PROJECT)) \\ & \bowtie Dnum=Dnumber(DEPARTMENT)) \bowtie Mgr\_ssn=Ssn(EMPLOYEE)) \end{aligned}$$

This corresponds to the following SQL query:

```
Q2: SELECT P.Pnumber, P.Dnum, E.Lname, E.Address, E.Bdate
      FROM   PROJECT P, DEPARTMENT D, EMPLOYEE E
      WHERE P.Dnum=D.Dnumber AND D.Mgr_ssn=E.Ssn AND
             P.Plocation= 'Stafford';
```

In Figure 19.1(a), the leaf nodes P, D, and E represent the three relations PROJECT, DEPARTMENT, and EMPLOYEE, respectively, and the internal tree nodes represent the *relational algebra operations* of the expression. When this query tree is executed, the node marked (1) in Figure 19.1(a) must begin execution before node (2) because some resulting tuples of operation (1) must be available before we can begin executing operation (2). Similarly, node (2) must begin executing and producing results before node (3) can start execution, and so on.

As we can see, the query tree represents a specific order of operations for executing a query. A more neutral data structure for representation of a query is the **query graph** notation. Figure 19.1(c) (the same as shown in Figure 6.13) shows the query graph for query Q2. Relations in the query are represented by **relation nodes**, which are displayed as single circles. Constant values, typically from the query selection conditions, are represented by **constant nodes**, which are displayed as double circles or ovals. Selection and join conditions are represented by the graph **edges**, as shown in Figure 19.1(c). Finally, the attributes to be retrieved from each relation are displayed in square brackets above each relation.

The query graph representation does not indicate an order on which operations to perform first. There is only a single graph corresponding to each query.<sup>2</sup> Although some optimization techniques were based on query graphs such as those originally in the INGRES DBMS, it is now generally accepted that query trees are preferable because, in practice, the query optimizer needs to show the order of operations for query execution, which is not possible in query graphs.

### 19.1.2 Heuristic Optimization of Query Trees

In general, many different relational algebra expressions—and hence many different query trees—can be **semantically equivalent**; that is, they can represent the *same query and produce the same results*.<sup>3</sup>

The query parser will typically generate a standard **initial query tree** to correspond to an SQL query, without doing any optimization. For example, for a SELECT-PROJECT-JOIN query, such as Q2, the initial tree is shown in Figure 19.1(b). The CARTESIAN PRODUCT of the relations specified in the FROM clause is first applied; then the selection and join conditions of the WHERE clause are applied, followed by

the projection operation. This represents a relation because of the 'CARTESIAN PRODUCT' and contained 'PROJECT' operation. This canonical form is easily created by the optimizer and is efficient.

The optimizer uses expressions to represent the query tree. These heuristics, and can be used in an

**Example of a query tree for a project named 'Aqua'**

```
Q: SELL
   FROM  PROJECT P
   WHERE P.Pnumber = 1000
```

The initial query tree first creates a 'CARTESIAN PRODUCT' of the 'EMPLOYEE', 'PROJECT' and 'DEPARTMENT' relations. It never executes the 'PROJECT' operation. This is because for the 'Aqua' project, the birth date is after the year 2000, so it applies the 'SELECT' operation. The 'CARTESIAN PRODUCT' operation is not needed.

A further improvement is to ignore the 'PROJECT' information and only consider the 'SELECT' and 'DEPARTMENT' relations. We can skip the 'CARTESIAN PRODUCT' operation because it is only the attributes of the 'DEPARTMENT' relation that are selected, by including them in the 'SELECT' tree, as shown. This produces intermediate tuples (records) for the 'DEPARTMENT' relation.

<sup>2</sup>Hence, a query graph corresponds to a *relational calculus* expression as shown in Section 8.6.5.

<sup>3</sup>The same query may also be stated in various ways in a high-level query language such as SQL (see Chapters 7 and 8).

the projection on the SELECT clause attributes. Such a **canonical query tree** represents a relational algebra expression that is *very inefficient if executed directly*, because of the CARTESIAN PRODUCT ( $\times$ ) operations. For example, if the PROJECT, DEPARTMENT, and EMPLOYEE relations had record sizes of 100, 50, and 150 bytes and contained 100, 20, and 5,000 tuples, respectively, the result of the CARTESIAN PRODUCT would contain 10 million tuples of record size 300 bytes each. However, this canonical query tree in Figure 19.1(b) is in a simple standard form that can be easily created from the SQL query. It will never be executed. The heuristic query optimizer will transform this initial query tree into an equivalent **final query tree** that is efficient to execute.

The optimizer must include rules for *equivalence among extended relational algebra expressions* that can be applied to transform the initial tree into the final, optimized query tree. First we discuss informally how a query tree is transformed by using heuristics, and then we discuss general transformation rules and show how they can be used in an algebraic heuristic optimizer.

**Example of Transforming a Query.** Consider the following query Q on the database in Figure 5.5: *Find the last names of employees born after 1957 who work on a project named 'Aquarius'.* This query can be specified in SQL as follows:

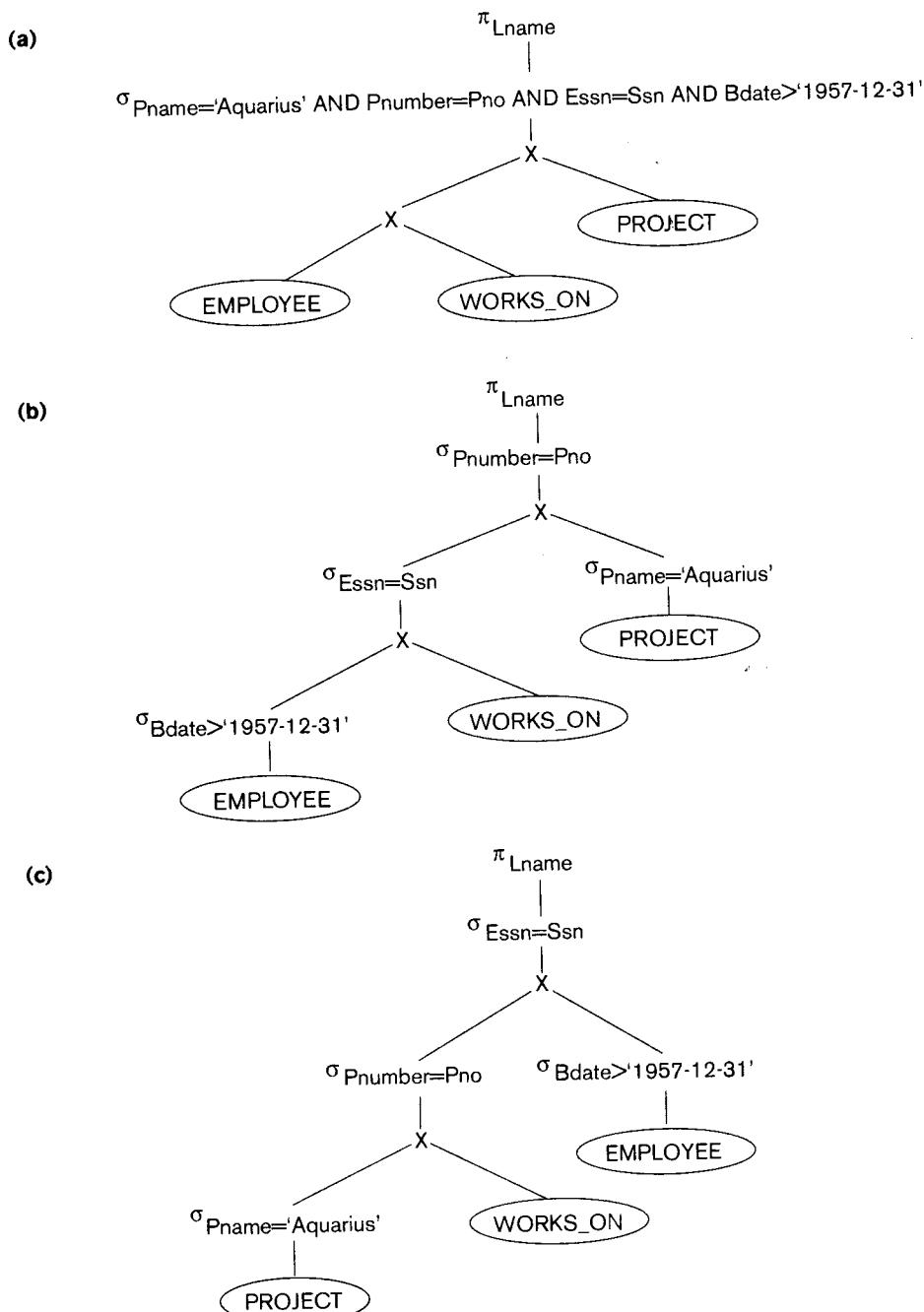
```
Q:   SELECT E.Lname
      FROM EMPLOYEE E, WORKS_ON W, PROJECT P
     WHERE P.Pname='Aquarius' AND P.Pnumber=W.Pno AND E.Essn=W.Ssn
          AND E.Bdate > '1957-12-31';
```

The initial query tree for Q is shown in Figure 19.2(a). Executing this tree directly first creates a very large file containing the CARTESIAN PRODUCT of the entire EMPLOYEE, WORKS\_ON, and PROJECT files. That is why the initial query tree is never executed, but is transformed into another equivalent tree that is efficient to execute. This particular query needs only one record from the PROJECT relation—for the 'Aquarius' project—and only the EMPLOYEE records for those whose date of birth is after '1957-12-31'. Figure 19.2(b) shows an improved query tree that first applies the SELECT operations to reduce the number of tuples that appear in the CARTESIAN PRODUCT.

A further improvement is achieved by switching the positions of the EMPLOYEE and PROJECT relations in the tree, as shown in Figure 19.2(c). This uses the information that Pnumber is a key attribute of the PROJECT relation, and hence the SELECT operation on the PROJECT relation will retrieve a single record only. We can further improve the query tree by replacing any CARTESIAN PRODUCT operation that is followed by a join condition as a selection with a JOIN operation, as shown in Figure 19.2(d). Another improvement is to keep only the attributes needed by subsequent operations in the intermediate relations, by including PROJECT ( $\pi$ ) operations as early as possible in the query tree, as shown in Figure 19.2(e). This reduces the attributes (columns) of the intermediate relations, whereas the SELECT operations reduce the number of tuples (records).

**Figure 19.2**

Steps in converting a query tree during heuristic optimization. (a) Initial (canonical) query tree for SQL query Q. Steps in converting a query tree during heuristic optimization. (a) Initial (canonical) query tree for SQL query Q.

**Figure 19.2 (c)**

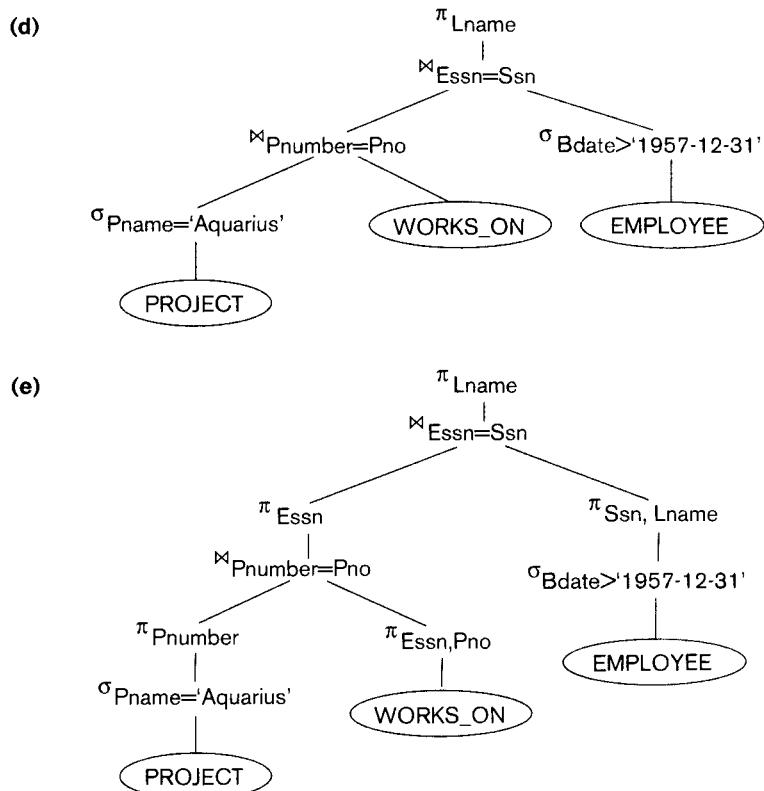
Steps in convert with JOIN opera

As the prece step into an must make tree. To do serve this eq

**General Tr**  
are many rule  
For query op  
tions and the  
butes in a di  
consider the  
nition of rela

**Figure 19.2 (continued)**

Steps in converting a query tree during heuristic optimization. (d) Replacing CARTESIAN PRODUCT and SELECT with JOIN operations. (e) Moving PROJECT operations down the query tree.



As the preceding example demonstrates, a query tree can be transformed step by step into an equivalent query tree that is more efficient to execute. However, we must make sure that the transformation steps always lead to an equivalent query tree. To do this, the query optimizer must know which transformation rules *preserve this equivalence*. We discuss some of these transformation rules next.

**General Transformation Rules for Relational Algebra Operations.** There are many rules for transforming relational algebra operations into equivalent ones. For query optimization purposes, we are interested in the meaning of the operations and the resulting relations. Hence, if two relations have the same set of attributes in a *different order* but the two relations represent the same information, we consider the relations to be equivalent. In Section 5.1.2 we gave an alternative definition of *relation* that makes the order of attributes unimportant; we will use this

definition here. We will state some transformation rules that are useful in query optimization, without proving them:

1. **Cascade of  $\sigma$ .** A conjunctive selection condition can be broken up into a cascade (that is, a sequence) of individual  $\sigma$  operations:

$$\sigma_{c_1} \text{ AND } c_2 \text{ AND } \dots \text{ AND } c_n(R) \equiv \sigma_{c_1}(\sigma_{c_2}(\dots(\sigma_{c_n}(R))\dots))$$

2. **Commutativity of  $\sigma$ .** The  $\sigma$  operation is commutative:

$$\sigma_{c_1}(\sigma_{c_2}(R)) \equiv \sigma_{c_2}(\sigma_{c_1}(R))$$

3. **Cascade of  $\pi$ .** In a cascade (sequence) of  $\pi$  operations, all but the last one can be ignored:

$$\pi_{\text{List}_1}(\pi_{\text{List}_2}(\dots(\pi_{\text{List}_n}(R))\dots)) \equiv \pi_{\text{List}_1}(R)$$

4. **Commuting  $\sigma$  with  $\pi$ .** If the selection condition  $c$  involves only those attributes  $A_1, \dots, A_n$  in the projection list, the two operations can be commuted:

$$\pi_{A_1, A_2, \dots, A_n}(\sigma_c(R)) \equiv \sigma_c(\pi_{A_1, A_2, \dots, A_n}(R))$$

5. **Commutativity of  $\bowtie$  (and  $\times$ ).** The join operation is commutative, as is the  $\times$  operation:

$$R \bowtie_c S \equiv S \bowtie_c R$$

$$R \times S \equiv S \times R$$

Notice that although the order of attributes may not be the same in the relations resulting from the two joins (or two Cartesian products), the *meaning* is the same because the order of attributes is not important in the alternative definition of relation.

6. **Commuting  $\sigma$  with  $\bowtie$  (or  $\times$ ).** If all the attributes in the selection condition  $c$  involve only the attributes of one of the relations being joined—say,  $R$ —the two operations can be commuted as follows:

$$\sigma_c(R \bowtie S) \equiv (\sigma_c(R)) \bowtie S$$

Alternatively, if the selection condition  $c$  can be written as  $(c_1 \text{ AND } c_2)$ , where condition  $c_1$  involves only the attributes of  $R$  and condition  $c_2$  involves only the attributes of  $S$ , the operations commute as follows:

$$\sigma_c(R \bowtie S) \equiv (\sigma_{c_1}(R)) \bowtie (\sigma_{c_2}(S))$$

The same rules apply if the  $\bowtie$  is replaced by a  $\times$  operation.

7. **Commuting  $\pi$  with  $\bowtie$  (or  $\times$ ).** Suppose that the projection list is  $L = \{A_1, \dots, A_n, B_1, \dots, B_m\}$ , where  $A_1, \dots, A_n$  are attributes of  $R$  and  $B_1, \dots, B_m$  are attributes of  $S$ . If the join condition  $c$  involves only attributes in  $L$ , the two operations can be commuted as follows:

$$\pi_L(R \bowtie S) \equiv (\pi_{A_1, \dots, A_n}(R)) \bowtie_c (\pi_{B_1, \dots, B_m}(S))$$

If the join condition  $c$  contains additional attributes not in  $L$ , these must be added to the projection list, and a final  $\pi$  operation is needed. For example, if attributes

$A_{n+1}, \dots$   
but are

$\pi_L(R \bowtie$

For  $\times, \bowtie$   
by repl.

3. **Commutativity,  $\bowtie$**

9. **Associativity;**  
any on

$(R \theta S)$

10. **Commutativity and  $-$ ;**  
expressions

$\sigma_c(R \Theta$

11. **The  $\pi$  expression**

$\pi_L(R \cup$

12. **Conversion correspondence**

$(\sigma_c(R \times$

13. **Pushing  $\sigma$  into  $\pi$**

$\sigma_c(R -$

However,

$\sigma_c(R \cap$

14. **Pushing  $\sigma$  into  $\bowtie$**

If in  $\tau$ ,

$\sigma_c(R \bowtie$

15. **Some remarks**

If  $S$  is

If the

There are other  
 $c$  can be converted  
rules from Boole

NOT ( $c_1 \wedge$

NOT ( $c_1 \vee$

useful in query

oken up into a

at the last one

ly those attri-  
be commuted:

ative, as is the  $\times$

me in the rela-  
s), the meaning  
the alternative

on condition  $c$   
—say,  $R$ —the

AND  $c_2$ ), where  
involves only

is  $L = \{A_1, \dots, A_n, B_1, \dots, B_m\}$  are  
in  $L$ , the two

must be added  
le, if attributes

$A_{n+1}, \dots, A_{n+k}$  of  $R$  and  $B_{m+1}, \dots, B_{m+p}$  of  $S$  are involved in the join condition  $c$  but are not in the projection list  $L$ , the operations commute as follows:

$$\pi_L(R \bowtie_c S) \equiv \pi_L((\pi_{A_1, \dots, A_n, A_{n+1}, \dots, A_{n+k}}(R)) \bowtie_c (\pi_{B_1, \dots, B_m, B_{m+1}, \dots, B_{m+p}}(S)))$$

For  $\times$ , there is no condition  $c$ , so the first transformation rule always applies by replacing  $\bowtie_c$  with  $\times$ .

8. **Commutativity of set operations.** The set operations  $\cup$  and  $\cap$  are commutative, but  $-$  is not.

9. **Associativity of  $\bowtie$ ,  $\times$ ,  $\cup$ , and  $\cap$ .** These four operations are individually associative; that is, if both occurrences of  $\theta$  stand for the same operation that is any one of these four operations (throughout the expression), we have:

$$(R \theta S) \theta T \equiv R \theta (S \theta T)$$

10. **Commuting  $\sigma$  with set operations.** The  $\sigma$  operation commutes with  $\cup$ ,  $\cap$ , and  $-$ . If  $\theta$  stands for any one of these three operations (throughout the expression), we have:

$$\sigma_c(R \theta S) \equiv (\sigma_c(R)) \theta (\sigma_c(S))$$

11. **The  $\pi$  operation commutes with  $\cup$ .**

$$\pi_L(R \cup S) \equiv (\pi_L(R)) \cup (\pi_L(S))$$

12. **Converting a  $(\sigma, \times)$  sequence into  $\bowtie$ .** If the condition  $c$  of a  $\sigma$  that follows a  $\times$  corresponds to a join condition, convert the  $(\sigma, \times)$  sequence into a  $\bowtie$  as follows:

$$(\sigma_c(R \times S)) \equiv (R \bowtie_c S)$$

13. **Pushing  $\sigma$  in conjunction with set difference.**

$$\sigma_c(R - S) = \sigma_c(R) - \sigma_c(S)$$

However,  $\sigma$  may be applied to only one relation:

$$\sigma_c(R - S) = \sigma_c(R) - S$$

14. **Pushing  $\sigma$  to only one argument in  $\cap$ .**

If in the condition  $\sigma_c$  all attributes are from relation  $R$ , then:

$$\sigma_c(R \cap S) = \sigma_c(R) \cap S$$

15. **Some trivial transformations.**

If  $S$  is empty, then  $R \cup S = R$

If the condition  $c$  in  $\sigma_c$  is true for the entire  $R$ , then  $\sigma_c(R) = R$ .

There are other possible transformations. For example, a selection or join condition  $c$  can be converted into an equivalent condition by using the following standard rules from Boolean algebra (De Morgan's laws):

$$\text{NOT } (c_1 \text{ AND } c_2) \equiv (\text{NOT } c_1) \text{ OR } (\text{NOT } c_2)$$

$$\text{NOT } (c_1 \text{ OR } c_2) \equiv (\text{NOT } c_1) \text{ AND } (\text{NOT } c_2)$$

Additional transformations discussed in Chapters 4, 5, and 6 are not repeated here. We discuss next how transformations can be used in heuristic optimization.

**Outline of a Heuristic Algebraic Optimization Algorithm.** We can now outline the steps of an algorithm that utilizes some of the above rules to transform an initial query tree into a final tree that is more efficient to execute (in most cases). The algorithm will lead to transformations similar to those discussed in our example in Figure 19.2. The steps of the algorithm are as follows:

1. Using Rule 1, break up any SELECT operations with conjunctive conditions into a cascade of SELECT operations. This permits a greater degree of freedom in moving SELECT operations down different branches of the tree.
2. Using Rules 2, 4, 6, and 10, 13, 14 concerning the commutativity of SELECT with other operations, move each SELECT operation as far down the query tree as is permitted by the attributes involved in the select condition. If the condition involves attributes from *only one table*, which means that it represents a *selection condition*, the operation is moved all the way to the leaf node that represents this table. If the condition involves attributes from *two tables*, which means that it represents a *join condition*, the condition is moved to a location down the tree after the two tables are combined.
3. Using Rules 5 and 9 concerning commutativity and associativity of binary operations, rearrange the leaf nodes of the tree using the following criteria. First, position the leaf node relations with the most restrictive SELECT operations so they are executed first in the query tree representation. The definition of *most restrictive* SELECT can mean either the ones that produce a relation with the fewest tuples or with the smallest absolute size.<sup>4</sup> Another possibility is to define the most restrictive SELECT as the one with the smallest selectivity; this is more practical because estimates of selectivities are often available in the DBMS catalog. Second, make sure that the ordering of leaf nodes does not cause CARTESIAN PRODUCT operations; for example, if the two relations with the most restrictive SELECT do not have a direct join condition between them, it may be desirable to change the order of leaf nodes to avoid Cartesian products.<sup>5</sup>
4. Using Rule 12, combine a CARTESIAN PRODUCT operation with a subsequent SELECT operation in the tree into a JOIN operation, if the condition represents a join condition.
5. Using Rules 3, 4, 7, and 11 concerning the cascading of PROJECT and the commuting of PROJECT with other operations, break down and move lists of projection attributes down the tree as far as possible by creating new PROJECT operations as needed. Only those attributes needed in the query result and in subsequent operations in the query tree should be kept after each PROJECT operation.

<sup>4</sup>Either definition can be used, since these rules are heuristic.

<sup>5</sup>Note that a CARTESIAN PRODUCT is acceptable in some cases—for example, if each relation has only a single tuple because each had a previous select condition on a key field.

6. Identify a single

In our example and 2 of the at after step 4; an operations in t. We may also g resulting from  $\pi_{E_{SSN}}$  because t

**Summary of** apply first the performing as and PROJECT and PROJECT and JOIN opera fewest tuples of similar operati nodes of the tre ing the rest of t

## 19.2 Cho

### 19.2.1 Alter

An execution p includes inform the algorithm's tree. As a simple relational alge

$\pi_{Fname, Lname}$

The query tree optimizer mig (assuming on the records in operation (ass of the JOIN re taken for exec although in ge

With material relation (that i tion can be con then read as in

6. Identify subtrees that represent groups of operations that can be executed by a single algorithm.

In our example, Figure 19.2(b) shows the tree in Figure 19.2(a) after applying steps 1 and 2 of the algorithm; Figure 19.2(c) shows the tree after step 3; Figure 19.2(d) after step 4; and Figure 19.2(e) after step 5. In step 6, we may group together the operations in the subtree whose root is the operation  $\pi_{E\text{ssn}}$  into a single algorithm. We may also group the remaining operations into another subtree, where the tuples resulting from the first algorithm replace the subtree whose root is the operation  $\pi_{E\text{ssn}}$ , because the first grouping means that this subtree is executed first.

**Summary of Heuristics for Algebraic Optimization.** The main heuristic is to apply first the operations that reduce the size of intermediate results. This includes performing as early as possible SELECT operations to reduce the number of tuples and PROJECT operations to reduce the number of attributes—by moving SELECT and PROJECT operations as far down the tree as possible. Additionally, the SELECT and JOIN operations that are most restrictive—that is, result in relations with the fewest tuples or with the smallest absolute size—should be executed before other similar operations. The latter rule is accomplished through reordering the leaf nodes of the tree among themselves while avoiding Cartesian products, and adjusting the rest of the tree appropriately.

## 19.2 Choice of Query Execution Plans

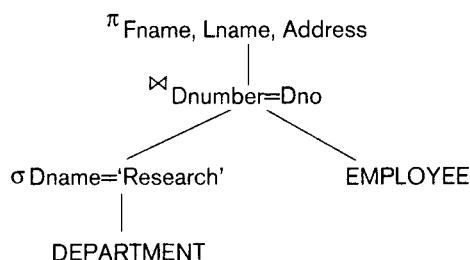
### 19.2.1 Alternatives for Query Evaluation

An execution plan for a relational algebra expression represented as a query tree includes information about the access methods available for each relation as well as the algorithms to be used in computing the relational operators represented in the tree. As a simple example, consider query Q1 from Chapter 7, whose corresponding relational algebra expression is

$$\pi_{Fname, Lname, Address}(\sigma_{Dname='Research'}(DEPARTMENT) \bowtie_{Dnumber=Dno} EMPLOYEE)$$

The query tree is shown in Figure 19.3. To convert this into an execution plan, the optimizer might choose an index search for the SELECT operation on DEPARTMENT (assuming one exists), an index-based nested-loop join algorithm that loops over the records in the result of the SELECT operation on DEPARTMENT for the join operation (assuming an index exists on the Dno attribute of EMPLOYEE), and a scan of the JOIN result for input to the PROJECT operator. Additionally, the approach taken for executing the query may specify a materialized or a pipelined evaluation, although in general a pipelined evaluation is preferred whenever feasible.

With **materialized evaluation**, the result of an operation is stored as a temporary relation (that is, the result is *physically materialized*). For instance, the JOIN operation can be computed and the entire result stored as a temporary relation, which is then read as input by the algorithm that computes the PROJECT operation, which

**Figure 19.3**

A query tree for query Q1.

would produce the query result table. On the other hand, with **pipelined evaluation**, as the resulting tuples of an operation are produced, they are forwarded directly to the next operation in the query sequence. We discussed pipelining as a strategy for query processing in Section 18.7. For example, as the selected tuples from DEPARTMENT are produced by the SELECT operation, they are placed in a buffer; the JOIN operation algorithm then consumes the tuples from the buffer, and those tuples that result from the JOIN operation are pipelined to the projection operation algorithm. The advantage of pipelining is the cost savings in not having to write the intermediate results to disk and not having to read them back for the next operation.

We discussed in Section 19.1 the possibility of converting query trees into equivalent trees so that the evaluation of the query is more efficient in terms of its execution time and overall resources consumed. There are more elaborate transformations of queries that are possible to optimize, or rather to “improve.” Transformations can be applied either in a heuristic-based or cost-based manner.

As we discussed in Sections 7.1.2 and 7.1.3, nested subqueries may occur in the WHERE clause as well as in the FROM clause of SQL queries. In the WHERE clause, if an inner block makes a reference to the relation used in the outer block, it is called a correlated nested query. When a query is used within the FROM clause to define a resulting or derived relation, which participates as a relation in the outer query, it is equivalent to a view. Both these types of nested subqueries are handled by the optimizer, which transforms them and rewrites the entire query. In the next two subsections, we consider these two variations of query transformation and rewriting with examples. We will call them nested subquery optimization and subquery (view) merging transformation. In Section 19.8, we revisit this topic in the context of data warehouses and illustrate star transformation optimizations.

### 19.2.2 Nested Subquery Optimization

We discussed nested queries in Section 7.1.2. Consider the query:

```

SELECT E1.Fname, E1.Lname
FROM EMPLOYEE E1
WHERE E1.Salary = ( SELECT MAX(Salary)
                    FROM EMPLOYEE E2)
  
```

In the above query, the outer block is a selection operation σ E1.Salary = MAX(E2.Salary). Evaluation of the inner query (SELECT MAX(Salary) FROM EMPLOYEE E2) produces a single value or scalar. The outer block is then evaluated to find all employees whose salary is maximum (if one exists). The result of the outer block is based on the scalar value produced by the inner block for both.

We discussed pipelined evaluation of the outer block, the inner block, and the outer block as a whole. The subquery is evaluated first, and the result is passed to the outer block.

Suppose in the above query, the inner query is a correlated query.

DEPARTMENT

Consider the query:

```

SELECT Fname, Lname
FROM EMPLOYEE
WHERE E1.Dnumber = Dnumber
  
```

In the above query, the outer block is a selection operation σ E1.Dnumber = Dnumber. The inner block is a projection operation π Fname, Lname. Depending on the cardinality of the outer relation, the optimizer may choose different strategies for evaluating the inner query. If the cardinality of the outer relation is small, the optimizer may convert the query into an equivalent query involving joins. If the cardinality of the outer relation is large, the optimizer may evaluate the inner query first and then use the result to evaluate the outer query.

```

SELECT Fname, Lname
FROM EMPLOYEE
WHERE E1.Dnumber = Dnumber
  
```

The process of transforming the query into an equivalent query involves rewriting the query. In the above query, the outer block is a selection operation σ E1.Dnumber = Dnumber. The inner block is a projection operation π Fname, Lname. The optimizer may choose to evaluate the inner query first and then use the result to evaluate the outer query. This is because the inner query is a simple projection operation, while the outer query is a selection operation. The optimizer may also choose to evaluate the outer query first and then use the result to evaluate the inner query. This is because the outer query is a selection operation, while the inner query is a simple projection operation. Other techniques for optimizing nested subqueries include star transformation and query merging.

~~view and avoid doing unnecessary computation. Sometimes an opposite situation happens. A view V is used in the query Q, and that view has been materialized as v; let us say the view includes  $R \bowtie S$ ; however, no access structures like indexes are available on v. Suppose that indexes are available on certain attributes, say, A of the component relation R and that the query Q involves a selection condition on A. In such cases, the query against the view can benefit by using the index on a component relation, and the view is replaced by its defining query; the relation representing the materialized view is not used at all.~~

### 19.3 Use of Selectivities in Cost-Based Optimization

A query optimizer does not depend solely on heuristic rules or query transformations; it also estimates and compares the costs of executing a query using different execution strategies and algorithms, and it then chooses the strategy with the *lowest cost estimate*. For this approach to work, accurate *cost estimates* are required so that different strategies can be compared fairly and realistically. In addition, the optimizer must limit the number of execution strategies to be considered; otherwise, too much time will be spent making cost estimates for the many possible execution strategies. Hence, this approach is more suitable for **compiled queries**, rather than ad-hoc queries where the optimization is done at compile time and the resulting execution strategy code is stored and executed directly at runtime. For **interpreted queries**, where the entire process shown in Figure 18.1 occurs at runtime, a full-scale optimization may slow down the response time. A more elaborate optimization is indicated for compiled queries, whereas a partial, less time-consuming optimization works best for interpreted queries.

This approach is generally referred to as **cost-based query optimization**.<sup>6</sup> It uses traditional optimization techniques that search the *solution space* to a problem for a solution that minimizes an objective (cost) function. The cost functions used in query optimization are estimates and not exact cost functions, so the optimization may select a query execution strategy that is not the optimal (absolute best) one. In Section 19.3.1, we discuss the components of query execution cost. In Section 19.3.2, we discuss the type of information needed in cost functions. This information is kept in the DBMS catalog. In Section 19.3.3, we describe histograms that are used to keep details on the value distributions of important attributes.

The decision-making process during query optimization is nontrivial and has multiple challenges. We can abstract the overall cost-based query optimization approach in the following way:

- For a given subexpression in the query, there may be multiple equivalence rules that apply. The process of applying equivalences is a cascaded one; it

<sup>6</sup>This approach was first used in the optimizer for the SYSTEM R in an experimental DBMS developed at IBM (Selinger et al., 1979).

- does not to cond  
 ■ It is nec  
 natives.  
 some c  
 for opt  
 ■ Approp  
 natives  
 ■ The sc  
 and inc  
 method  
 must c  
 ■ In a glo  
 blocks.<sup>7</sup>
- #### 19.3.1 Cost
- The cost of exe
1. Access c  
 and wri  
 ory buf  
 search u  
 tures o  
 indexes  
 contigu  
 access c
  2. Disk st  
 that are
  3. Compu  
 the reco  
 include  
 sort op  
 knowin
  4. Memo  
 memori
  5. Comm  
 from th  
 In disti  
 transfe  
 evaluat

<sup>7</sup>We do not discuss  
 Ahmed et al. (2000)

posite situation materialized as  $v$ ; key indexes are  $s$ , say,  $A$  of the relation on  $A$ . In this case, a composition representation does not have any limit and there is no definitive convergence. It is difficult to conduct this in a space-efficient manner.

- It is necessary to resort to some quantitative measure for evaluation of alternatives. By using the space and time requirements and reducing them to some common metric called cost, it is possible to devise some methodology for optimization.
- Appropriate search strategies can be designed by keeping the cheapest alternatives and pruning the costlier alternatives.
- The scope of query optimization is generally a query block. Various table and index access paths, join permutations (orders), join methods, group-by methods, and so on provide the alternatives from which the query optimizer must choose.
- In a global query optimization, the scope of optimization is multiple query blocks.<sup>7</sup>

### 19.3.1 Cost Components for Query Execution

The cost of executing a query includes the following components:

1. **Access cost to secondary storage.** This is the cost of transferring (reading and writing) data blocks between secondary disk storage and main memory buffers. This is also known as *disk I/O (input/output) cost*. The cost of searching for records in a disk file depends on the type of access structures on that file, such as ordering, hashing, and primary or secondary indexes. In addition, factors such as whether the file blocks are allocated contiguously on the same disk cylinder or scattered on the disk affect the access cost.
2. **Disk storage cost.** This is the cost of storing on disk any intermediate files that are generated by an execution strategy for the query.
3. **Computation cost.** This is the cost of performing in-memory operations on the records within the data buffers during query execution. Such operations include searching for and sorting records, merging records for a join or a sort operation, and performing computations on field values. This is also known as *CPU (central processing unit) cost*.
4. **Memory usage cost.** This is the cost pertaining to the number of main memory buffers needed during query execution.
5. **Communication cost.** This is the cost of shipping the query and its results from the database site to the site or terminal where the query originated. In distributed databases (see Chapter 23), it would also include the cost of transferring tables and results among various computers during query evaluation.

<sup>7</sup>We do not discuss global optimization in this sense in the present chapter. Details may be found in Ahmed et al. (2006).

For large databases, the main emphasis is often on minimizing the access cost to secondary storage. Simple cost functions ignore other factors and compare different query execution strategies in terms of the number of block transfers between disk and main memory buffers. For smaller databases, where most of the data in the files involved in the query can be completely stored in memory, the emphasis is on minimizing computation cost. In distributed databases, where many sites are involved (see Chapter 23), communication cost must be minimized. It is difficult to include all the cost components in a (weighted) cost function because of the difficulty of assigning suitable weights to the cost components. This is why some cost functions consider a single factor only—disk access. In the next section, we discuss some of the information that is needed for formulating cost functions.

### 19.3.2 Catalog Information Used in Cost Functions

To estimate the costs of various execution strategies, we must keep track of any information that is needed for the cost functions. This information may be stored in the DBMS catalog, where it is accessed by the query optimizer. First, we must know the size of each file. For a file whose records are all of the same type, the **number of records (tuples) ( $r$ )**, the (average) **record size ( $R$ )**, and the **number of file blocks ( $b$ )** (or close estimates of them) are needed. The **blocking factor ( $bfr$ )** for the file may also be needed. These were mentioned in Section 18.3.4, and we utilized them while illustrating the various implementation algorithms for relational operations. We must also keep track of the *primary file organization* for each file. The primary file organization records may be *unordered*, *ordered* by an attribute with or without a primary or clustering index, or *hashed* (static hashing or one of the dynamic hashing methods) on a key attribute. Information is also kept on all primary, secondary, or clustering indexes and their indexing attributes. The **number of levels ( $x$ )** of each multilevel index (primary, secondary, or clustering) is needed for cost functions that estimate the number of block accesses that occur during query execution. In some cost functions the **number of first-level index blocks ( $b_{fl}$ )** is needed.

Another important parameter is the **number of distinct values NDV** ( $A, R$ ) of an attribute in relation  $R$  and the attribute **selectivity** ( $sl$ ), which is the fraction of records satisfying an equality condition on the attribute. This allows estimation of the **selection cardinality** ( $s = sl \times r$ ) of an attribute, which is the *average* number of records that will satisfy an equality selection condition on that attribute.

Information such as the number of index levels is easy to maintain because it does not change very often. However, other information may change frequently; for example, the number of records  $r$  in a file changes every time a record is inserted or deleted. The query optimizer will need reasonably close but not necessarily completely up-to-the-minute values of these parameters for use in estimating the cost of various execution strategies. To help with estimating the size of the results of queries, it is important to have as good an estimate of the distribution of values as possible. To that end, most systems store a histogram.

### 19.3.3 Histo

Histograms are a common way to visualize data distribution. The assumption is that the data is continuous and can be divided into buckets. The width of each bucket is called the bin width, and the height of each bucket is called the frequency or count. The total number of buckets is called the bin count. The relationship between the bin width and the bin count is called the binning rule. Some common binning rules include equal-width binning, quantile binning, and density binning. Equal-width binning is the most common and easiest to understand. It divides the range of the data into equal-width intervals. Quantile binning divides the data into equal-sized groups based on their rank. Density binning divides the data into bins of equal density. The choice of binning rule depends on the nature of the data and the purpose of the histogram.

the access cost to  
and compare dif-  
block transfers  
where most of  
stored in memory,  
databases,  
on cost must be  
(weighted) cost  
to the cost com-  
actor only—disk  
that is needed for

keep track of any  
may be stored  
First, we must  
same type, the  
the **number of**  
**factor (bf)** for  
and we utilized  
elational oper-  
r each file. The  
ttribute with or  
; or one of the  
kept on all pri-  
s. The **number**  
(ring) is needed  
t occur during  
l index blocks

**V (A, R)** of an  
the fraction of  
estimation of  
the number of  
ite.

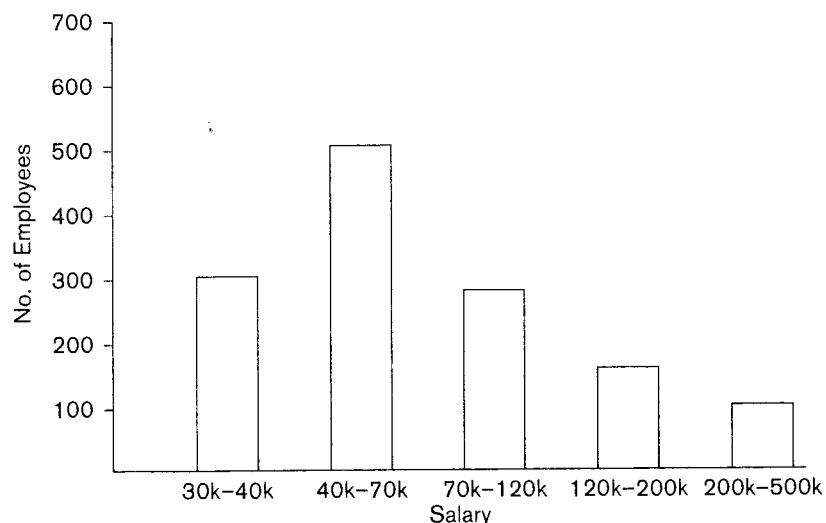
cause it does  
equently; for  
is inserted or  
essarily com-  
ng the cost of  
sults of que-  
values as pos-

### 19.3.3 Histograms

**Histograms** are tables or data structures maintained by the DBMS to record information about the distribution of data. It is customary for most RDBMSs to store histograms for most of the important attributes. Without a histogram, the best assumption is that values of an attribute are uniformly distributed over its range from high to low. Histograms divide the attribute over important ranges (called buckets) and store the total number of records that belong to that bucket in that relation. Sometimes they may also store the number of distinct values in each bucket as well. An implicit assumption is made sometimes that among the distinct values within a bucket there is a uniform distribution. All these assumptions are oversimplifications that rarely hold. So keeping a histogram with a finer granularity (i.e., larger number of buckets) is always useful. A couple of variations of histograms are common: in **equi-width** histograms, the range of values is divided into equal subranges. In **equi-height** histograms, the buckets are so formed that each one contains roughly the same number of records. Equi-height histograms are considered better since they keep fewer numbers of more frequently occurring values in one bucket and more numbers of less frequently occurring ones in a different bucket. So the uniform distribution assumption within a bucket seems to hold better. We show an example of a histogram for salary information in a company in Figure 19.4. This histogram divides the salary range into five buckets that may correspond to the important sub-ranges over which the queries may be likely because they belong to certain types of employees. It is neither an equi-width nor an equi-height histogram.

**Figure 19.4**

Histogram of salary in the relation EMPLOYEE.



## 19.4 Cost Functions for SELECT Operation

We now provide cost functions for the selection algorithms S1 to S8 discussed in Section 18.3.1 in terms of *number of block transfers* between memory and disk. Algorithm S9 involves an intersection of record pointers after they have been retrieved by some other means, such as algorithm S6, and so the cost function will be based on the cost for S6. These cost functions are estimates that ignore computation time, storage cost, and other factors. To reiterate, the following notation is used in the formulas hereafter:

- $C_{S_i}$ : Cost for method  $S_i$  in block accesses
- $r_X$ : Number of records (tuples) in a relation  $X$
- $b_X$ : Number of blocks occupied by relation  $X$  (also referred to as  $b$ )
- $bfr_X$ : Blocking factor (i.e., number of records per block) in relation  $X$
- $sl_A$ : Selectivity of an attribute  $A$  for a given condition
- $sA$ : Selection cardinality of the attribute being selected ( $= sl_A * r$ )
- $x_A$ : Number of levels of the index for attribute  $A$
- $b_{I1A}$ : Number of first-level blocks of the index on attribute  $A$
- $NDV(A, X)$ : Number of distinct values of attribute  $A$  in relation  $X$

*Note:* In using the above notation in formulas, we have omitted the relation name or attribute name when it is obvious.

- **S1—Linear search (brute force) approach.** We search all the file blocks to retrieve all records satisfying the selection condition; hence,  $C_{S1a} = b$ . For an *equality condition on a key attribute*, only half the file blocks are searched *on the average* before finding the record, so a rough estimate for  $C_{S1b} = (b/2)$  if the record is found; if no record is found that satisfies the condition,  $C_{S1b} = b$ .
- **S2—Binary search.** This search accesses approximately  $C_{S2} = \log_2 b + \lceil (s/bfr) \rceil - 1$  file blocks. This reduces to  $\log_2 b$  if the equality condition is on a unique (key) attribute, because  $s = 1$  in this case.
- **S3a—Using a primary index to retrieve a single record.** For a primary index, retrieve one disk block at each index level, plus one disk block from the data file. Hence, the cost is one more disk block than the number of index levels:  $C_{S3a} = x + 1$ .
- **S3b—Using a hash key to retrieve a single record.** For hashing, only one disk block needs to be accessed in most cases. The cost function is approximately  $C_{S3b} = 1$  for static hashing or linear hashing, and it is 2 disk block accesses for extendible hashing (see Section 16.8).
- **S4—Using an ordering index to retrieve multiple records.** If the comparison condition is  $>$ ,  $\geq$ ,  $<$ , or  $\leq$  on a key field with an ordering index, roughly half the file records will satisfy the condition. This gives a cost function of  $C_{S4} = x + (b/2)$ . This is a very rough estimate, and although it may be correct on the average, it may be inaccurate in individual cases. A more accurate estimate is possible if the distribution of records is stored in a histogram.
- **S5—Using a clustering index to retrieve multiple records.** One disk block is accessed at each index level, which gives the address of the first file disk

**ion**

S8 discussed in  
nary and disk.  
hey have been  
st function will  
nore computa-  
ing notation is

s b)

ion X

)

n X

e relation name

ne file blocks to  
 $C_{S1a} = b$ . For an  
are searched on  
 $C_{S1b} = (b/2)$  if  
dition,  $C_{S1b} = b$ .  
nately  $C_{S2} =$   
uality condition

For a primary  
disk block from  
the number of

ing, only one  
en is approxi-  
2 disk block

If the compari-  
index, roughly  
ost function of  
may be correct  
more accurate  
histogram.

One disk block  
e first file disk

block in the cluster. Given an equality condition on the indexing attribute,  $s$  records will satisfy the condition, where  $s$  is the selection cardinality of the indexing attribute. This means that  $\lceil (s/bfr) \rceil$  file blocks will be in the cluster of file blocks that hold all the selected records, giving  $C_{S5} = x + \lceil (s/bfr) \rceil$ .

- **S6—Using a secondary ( $B^+$ -tree) index.** For a secondary index on a key (unique) attribute, with an equality (i.e.,  $<\text{attribute} = \text{value}>$ ) selection condition, the cost is  $x + 1$  disk block accesses. For a secondary index on a nonkey (nonunique) attribute,  $s$  records will satisfy an equality condition, where  $s$  is the selection cardinality of the indexing attribute. However, because the index is nonclustering, each of the records may reside on a different disk block, so the (worst case) cost estimate is  $C_{S6a} = x + 1 + s$ . The additional 1 is to account for the disk block that contains the record pointers after the index is searched (see Figure 17.5). For range queries, if the comparison condition is  $>$ ,  $\geq$ ,  $<$ , or  $\leq$  and half the file records are assumed to satisfy the condition, then (very roughly) half the first-level index blocks are accessed, plus half the file records via the index. The cost estimate for this case, approximately, is  $C_{S6b} = x + (b_{I1}/2) + (r/2)$ . The  $r/2$  factor can be refined if better selectivity estimates are available through a histogram. The latter method  $C_{S6b}$  can be very costly. For a range condition such as  $v1 < A < v2$ , the selection cardinality  $s$  must be computed from the histogram or as a default, under the uniform distribution assumption; then the cost would be computed based on whether or not  $A$  is a key or nonkey with a  $B^+$ -tree index on  $A$ . (We leave this as an exercise for the reader to compute under the different conditions.)
- **S7—Conjunctive selection.** We can use either S1 or one of the methods S2 to S6 discussed above. In the latter case, we use one condition to retrieve the records and then check in the main memory buffers whether each retrieved record satisfies the remaining conditions in the conjunction. If multiple indexes exist, the search of each index can produce a set of record pointers (record ids) in the main memory buffers. The intersection of the sets of record pointers (referred to in S9) can be computed in main memory, and then the resulting records are retrieved based on their record ids.
- **S8—Conjunctive selection using a composite index.** Same as S3a, S5, or S6a, depending on the type of index.
- **S9—Selection using a bitmap index.** (See Section 17.5.2.) Depending on the nature of selection, if we can reduce the selection to a set of equality conditions, each equating the attribute with a value (e.g.,  $A = \{7, 13, 17, 55\}$ ), then a bit vector for each value is accessed which is  $r$  bits or  $r/8$  bytes long. A number of bit vectors may fit in one block. Then, if  $s$  records qualify,  $s$  blocks are accessed for the data records.
- **S10—Selection using a functional index.** (See Section 17.5.3.) This works similar to S6 except that the index is based on a function of multiple attributes; if that function is appearing in the SELECT clause, the corresponding index may be utilized.

**Cost-Based Optimization Approach.** In a query optimizer, it is common to enumerate the various possible strategies for executing a query and to estimate the costs for different strategies. An optimization technique, such as dynamic programming, may be used to find the optimal (least) cost estimate efficiently without having to consider all possible execution strategies. **Dynamic programming** is an optimization technique<sup>8</sup> in which subproblems are solved only once. This technique is applicable when a problem may be broken down into subproblems that themselves have subproblems. We will visit the dynamic programming approach when we discuss join ordering in Section 19.5.5. We do not discuss optimization algorithms here; rather, we use a simple example to illustrate how cost estimates may be used.

#### 19.4.1 Example of Optimization of Selection Based on Cost Formulas:

Suppose that the EMPLOYEE file in Figure 5.5 has  $r_E = 10,000$  records stored in  $b_E = 2,000$  disk blocks with blocking factor  $bfr_E = 5$  records/block and the following access paths:

1. A clustering index on Salary, with levels  $x_{\text{Salary}} = 3$  and average selection cardinality  $s_{\text{Salary}} = 20$ . (This corresponds to a selectivity of  $sl_{\text{Salary}} = 20/10000 = 0.002$ .)
2. A secondary index on the key attribute Ssn, with  $x_{\text{Ssn}} = 4$  ( $s_{\text{Ssn}} = 1$ ,  $sl_{\text{Ssn}} = 0.0001$ ).
3. A secondary index on the nonkey attribute Dno, with  $x_{\text{Dno}} = 2$  and first-level index blocks  $b_{I1\text{Dno}} = 4$ . There are  $NDV(\text{Dno}, \text{EMPLOYEE}) = 125$  distinct values for Dno, so the selectivity of Dno is  $sl_{\text{Dno}} = (1/NDV(\text{Dno}, \text{EMPLOYEE})) = 0.008$ , and the selection cardinality is  $s_{\text{Dno}} = (r_E * sl_{\text{Dno}}) = (r_E/NDV(\text{Dno}, \text{EMPLOYEE})) = 80$ .
4. A secondary index on Sex, with  $x_{\text{Sex}} = 1$ . There are  $NDV(\text{Sex}, \text{EMPLOYEE}) = 2$  values for the Sex attribute, so the average selection cardinality is  $s_{\text{Sex}} = (r_E/NDV(\text{Sex}, \text{EMPLOYEE})) = 5000$ . (Note that in this case, a histogram giving the percentage of male and female employees may be useful, unless the percentages are approximately equal.)

We illustrate the use of cost functions with the following examples:

- OP1:  $\sigma_{\text{Ssn}='123456789'}(\text{EMPLOYEE})$
- OP2:  $\sigma_{\text{Dno}>5}(\text{EMPLOYEE})$
- OP3:  $\sigma_{\text{Dno}=5}(\text{EMPLOYEE})$
- OP4:  $\sigma_{\text{Dno}=5 \text{ AND } \text{SALARY}>30000 \text{ AND } \text{Sex}='F'}(\text{EMPLOYEE})$

The cost of the brute force (linear search or file scan) option S1 will be estimated as  $C_{S1a} = b_E = 2000$  (for a selection on a nonkey attribute) or  $C_{S1b} = (b_E/2) = 1,000$ .

(average cost for S1 or method 1 chosen over method 2) either method 1 cost  $C_{S6b} = x_{\text{Dno}} * b_E = 2 * 2000 = 4000$ , the linear search cost  $C_{S6a} = 2 + 80 = 82$ , so

Finally, consider the cost of retrieving the record  $C_{S1a} = 2000$ . Using the condition  $3 + (2000/2) = 1003$ ,  $x_{\text{Sex}} + s_{\text{Sex}} = 1 + 5000 = 5001$ , secondary index on Dno = 5 is used. Condition (Salary is retrieved into memory) are included in the above; Dno has value 5 if we had an attribute 30332 and we have 5 records qualify. in a cost of 2,000.

## 19.5 Costs of Selections

To develop reasonable cost estimate for the selection. This is usually done by file to the size of the files, and it is computed relation R by  $|R|$ .

$$js = |(R \bowtie_C S)|$$

If there is no join operation PRODUCT. If n general,  $0 \leq js \leq r_A * s_B$ , we get

1. If A is a single record in the key of R. of S that

<sup>8</sup>For a detailed discussion of dynamic programming as a technique of optimization, the reader may consult an algorithm textbook such as Cormen et al. (2003).

r, it is common to  
and to estimate the  
dynamic program-  
ently without hav-  
**programming** is an  
once. This tech-  
subproblems that  
amming approach  
ssus optimization  
ow cost estimates

d

records stored in  
and the following

$\exists$  selection card-  
 $20/10000 = 0.002.$ )  
= 1,  $s_{Ssn} = 0.0001).$   
= 2 and first-level  
= 125 distinct val-  
 $EMPLOYEE) =$   
 $(r_E/NDV(Dno,$

$EMPLOYEE) =$   
linality is  $s_{Sex} =$   
se, a histogram  
e useful, unless

be estimated as  
 $(b_E/2) = 1,000$

eader may con-

(average cost for a selection on a key attribute). For OP1 we can use either method S1 or method S6a; the cost estimate for S6a is  $C_{S6a} = x_{Ssn} + 1 = 4 + 1 = 5$ , and it is chosen over method S1, whose average cost is  $C_{S1b} = 1,000$ . For OP2 we can use either method S1 (with estimated cost  $C_{S1a} = 2,000$ ) or method S6b (with estimated cost  $C_{S6b} = x_{Dno} + (b_{IDno}/2) + (r_E/2) = 2 + (4/2) + (10,000/2) = 5,004$ ), so we choose the linear search approach for OP2. For OP3 we can use either method S1 (with estimated cost  $C_{S1a} = 2,000$ ) or method S6a (with estimated cost  $C_{S6a} = x_{Dno} + s_{Dno} = 2 + 80 = 82$ ), so we choose method S6a.

Finally, consider OP4, which has a conjunctive selection condition. We need to estimate the cost of using any one of the three components of the selection condition to retrieve the records, plus the linear search approach. The latter gives cost estimate  $C_{S1a} = 2000$ . Using the condition ( $Dno = 5$ ) first gives the cost estimate  $C_{S6a} = 82$ . Using the condition ( $Salary > 30000$ ) first gives a cost estimate  $C_{S4} = x_{Salary} + (b_E/2) = 3 + (2000/2) = 1003$ . Using the condition ( $Sex = 'F'$ ) first gives a cost estimate  $C_{S6a} = x_{Sex} + s_{Sex} = 1 + 5000 = 5001$ . The optimizer would then choose method S6a on the secondary index on Dno because it has the lowest cost estimate. The condition ( $Dno = 5$ ) is used to retrieve the records, and the remaining part of the conjunctive condition ( $Salary > 30,000 \text{ AND } Sex = 'F'$ ) is checked for each selected record after it is retrieved into memory. Only the records that satisfy these additional conditions are included in the result of the operation. Consider the  $Dno = 5$  condition in OP3 above; Dno has 125 values and hence a  $B^+$ -tree index would be appropriate. Instead, if we had an attribute Zipcode in EMPLOYEE and if the condition were  $Zipcode = 30332$  and we had only five zip codes, bitmap indexing could be used to know what records qualify. Assuming uniform distribution,  $s_{Zipcode} = 2,000$ . This would result in a cost of 2,000 for bitmap indexing.

## 19.5 Cost Functions for the JOIN Operation

To develop reasonably accurate cost functions for JOIN operations, we must have an estimate for the size (number of tuples) of the file that results *after* the JOIN operation. This is usually kept as a ratio of the size (number of tuples) of the resulting join file to the size of the CARTESIAN PRODUCT file, if both are applied to the same input files, and it is called the **join selectivity (js)**. If we denote the number of tuples of a relation  $R$  by  $|R|$ , we have:

$$js = |(R \bowtie_c S)| / |(R \times S)| = |(R \bowtie_c S)| / (|R| * |S|)$$

If there is no join condition  $c$ , then  $js = 1$  and the join is the same as the CARTESIAN PRODUCT. If no tuples from the relations satisfy the join condition, then  $js = 0$ . In general,  $0 \leq js \leq 1$ . For a join where the condition  $c$  is an equality comparison  $R.A = S.B$ , we get the following two special cases:

1. If  $A$  is a key of  $R$ , then  $|(R \bowtie_c S)| \leq |S|$ , so  $js \leq (1/|R|)$ . This is because each record in file  $S$  will be joined with at most one record in file  $R$ , since  $A$  is a key of  $R$ . A special case of this condition is when attribute  $B$  is a *foreign key* of  $S$  that references the *primary key*  $A$  of  $R$ . In addition, if the foreign key  $B$

has the NOT NULL constraint, then  $js = (1/|R|)$ , and the result file of the join will contain  $|S|$  records.

2. If  $B$  is a key of  $S$ , then  $|(R \bowtie_c S)| \leq |R|$ , so  $js \leq (1/|S|)$ .

Hence a **simple formula** to use for join selectivity is:

$$js = 1 / \max(\text{NDV}(A, R), \text{NDV}(B, S))$$

Having an estimate of the join selectivity for commonly occurring join conditions enables the query optimizer to estimate the size of the resulting file after the join operation, which we call **join cardinality ( $jc$ )**.

$$jc = |(R_c S)| = js * |R| * |S|.$$

We can now give some sample *approximate* cost functions for estimating the cost of some of the join algorithms given in Section 18.4. The join operations are of the form:

$$R \bowtie_{A=B} S$$

where  $A$  and  $B$  are domain-compatible attributes of  $R$  and  $S$ , respectively. Assume that  $R$  has  $b_R$  blocks and that  $S$  has  $b_S$  blocks:

- **J1—Nested-loop join.** Suppose that we use  $R$  for the outer loop; then we get the following cost function to estimate the number of block accesses for this method, assuming *three memory buffers*. We assume that the blocking factor for the resulting file is  $bfr_{RS}$  and that the join selectivity is known:

$$C_{J1} = b_R + (b_R * b_S) + ((js * |R| * |S|) / bfr_{RS})$$

The last part of the formula is the cost of writing the resulting file to disk. This cost formula can be modified to take into account different numbers of memory buffers, as presented in Section 19.4. If  $n_B$  main memory buffer blocks are available to perform the join, the cost formula becomes:

$$C_{J1} = b_R + (\lceil b_R/(n_B - 2) \rceil * b_S) + ((js * |R| * |S|) / bfr_{RS})$$

- **J2—Index-based nested-loop join (using an access structure to retrieve the matching record(s)).** If an index exists for the join attribute  $B$  of  $S$  with index levels  $x_B$ , we can retrieve each record  $s$  in  $R$  and then use the index to retrieve all the matching records  $t$  from  $S$  that satisfy  $t[B] = s[A]$ . The cost depends on the type of index. For a secondary index where  $s_B$  is the selection cardinality for the join attribute  $B$  of  $S$ ,<sup>9</sup> we get:

$$C_{J2a} = b_R + (|R| * (x_B + 1 + s_B)) + ((js * |R| * |S|) / bfr_{RS})$$

For a clustering index where  $s_B$  is the selection cardinality of  $B$ , we get

$$C_{J2b} = b_R + (|R| * (x_B + (s_B/bfr_B))) + ((js * |R| * |S|) / bfr_{RS})$$

<sup>9</sup>Selection cardinality was defined as the average number of records that satisfy an equality condition on an attribute, which is the average number of records that have the same value for the attribute and hence will be joined to a single record in the other file.

For a pri

$$C_{J2c} = b_I$$

If a hash

$$C_{J2d} = b_h$$

where  $h$  given its hashing ically  $h$

- **J3—Sort cost fun**

$$C_{J3a} = b_s$$

If we mi formula

- **J4—Partitioned same ha and B of this join**

$$C_{J4} = 3 *$$

### 19.5.1 Join for S

We consider the join queries. In into these oper of doing exhaust Let us consider

#### Semi-Join

```
SELECT C
  FROM T1
 WHERE T1.B = T2.B
   AND T2.C = 'S'
```

Unnesting of t ion “S=” for s

```
SELECT C
  FROM T1
 WHERE T1.B = T2.B
```

result file of the join

ing join conditions  
g file after the join

imating the cost of  
ons are of the form:

pectively. Assume

r loop; then we get  
ck accesses for this  
he blocking factor  
known:

ng file to disk. This  
erent numbers of  
n memory buffer  
comes:

cture to retrieve  
tribute  $B$  of  $S$  with  
use the index to  
 $s[A]$ . The cost  
 $i_B$  is the selection

f  $B$ , we get

quality condition on  
attribute and

For a primary index, we get

$$C_{J2c} = b_R + (|R| * (x_B + 1)) + ((js * |R| * |S|) / bfr_{RS})$$

If a **hash key** exists for one of the two join attributes—say,  $B$  of  $S$ —we get

$$C_{J2d} = b_R + (|R| * h) + ((js * |R| * |S|) / bfr_{RS})$$

where  $h \geq 1$  is the average number of block accesses to retrieve a record, given its hash key value. Usually,  $h$  is estimated to be 1 for static and linear hashing and 2 for extendible hashing. This is an optimistic estimate, and typically  $h$  ranges from 1.2 to 1.5 in practical situations.

- **J3—Sort-merge join.** If the files are already sorted on the join attributes, the cost function for this method is

$$C_{J3a} = b_R + b_S + ((js * |R| * |S|) / bfr_{RS})$$

If we must sort the files, the cost of sorting must be added. We can use the formulas from Section 18.2 to estimate the sorting cost.

- **J4—Partition-hashed join (or just hash join).** The records of files  $R$  and  $S$  are partitioned into smaller files. The partitioning of each file is done using the same hashing function  $h$  on the join attribute  $A$  of  $R$  (for partitioning file  $R$ ) and  $B$  of  $S$  (for partitioning file  $S$ ). As we showed in Section 18.4, the cost of this join can be approximated to:

$$C_{J4} = 3 * (b_R + b_S) + ((js * |R| * |S|) / bfr_{RS})$$

### 19.5.1 Join Selectivity and Cardinality for Semi-Join and Anti-Join

We consider these two important operations, which are used when unnesting certain queries. In Section 18.1 we showed examples of subqueries that are transformed into these operations. The goal of these operations is to avoid the unnecessary effort of doing exhaustive pairwise matching of two tables based on the join condition. Let us consider the join selectivity and cardinality of these two types of joins.

#### Semi-Join

```
SELECT COUNT(*)
  FROM T1
 WHERE T1.X IN (SELECT T2.Y
                  FROM T2);
```

Unnesting of the query above leads to semi-join. (In the following query, the notation “ $S=$ ” for semi-join is nonstandard.)

```
SELECT COUNT(*)
  FROM T1, T2
 WHERE T1.X S= T2.Y;
```



per block. We can estimate the worst-case costs for the JOIN operation OP6 using the applicable methods J1 and J2 as follows:

1. Using method J1 with EMPLOYEE as outer loop:

$$\begin{aligned} C_{J1} &= b_E + (b_E * b_D) + ((js_{OP6} * r_E * r_D) / bfr_{ED}) \\ &= 2,000 + (2,000 * 13) + (((1/125) * 10,000 * 125)/4) = 30,500 \end{aligned}$$

2. Using method J1 with DEPARTMENT as outer loop:

$$\begin{aligned} C_{J1} &= b_D + (b_E * b_D) + ((js_{OP6} * r_E * r_D) / bfr_{ED}) \\ &= 13 + (13 * 2,000) + (((1/125) * 10,000 * 125)/4) = 28,513 \end{aligned}$$

3. Using method J2 with EMPLOYEE as outer loop:

$$\begin{aligned} C_{J2c} &= b_E + (r_E * (x_{Dnumber} + 1)) + ((js_{OP6} * r_E * r_D) / bfr_{ED}) \\ &= 2,000 + (10,000 * 2) + (((1/125) * 10,000 * 125)/4) = 24,500 \end{aligned}$$

4. Using method J2 with DEPARTMENT as outer loop:

$$\begin{aligned} C_{J2a} &= b_D + (r_D * (x_{Dno} + s_{Dno})) + ((js_{OP6} * r_E * r_D) / bfr_{ED}) \\ &= 13 + (125 * (2 + 80)) + (((1/125) * 10,000 * 125)/4) = 12,763 \end{aligned}$$

5. Using method J4 gives:

$$\begin{aligned} C_{J4} &= 3 * (b_D + b_E) + ((js_{OP6} * r_E * r_D) / bfr_{ED}) \\ &= 3 * (13 + 2,000) + 2,500 = 8,539 \end{aligned}$$

Case 5 has the lowest cost estimate and will be chosen. Notice that in case 2 above, if 15 memory buffer blocks (or more) were available for executing the join instead of just 3, 13 of them could be used to hold the entire DEPARTMENT relation (outer loop relation) in memory, one could be used as buffer for the result, and one would be used to hold one block at a time of the EMPLOYEE file (inner loop file), and the cost for case 2 could be drastically reduced to just  $b_E + b_D + ((js_{OP6} * r_E * r_D) / bfr_{ED})$  or 4,513, as discussed in Section 18.4. If some other number of main memory buffers was available, say  $n_B = 10$ , then the cost for case 2 would be calculated as follows, which would also give better performance than case 4:

$$\begin{aligned} C_{J1} &= b_D + (\lceil b_D/(n_B - 2) \rceil * b_E) + ((js * |R| * |S|) / bfr_{RS}) \\ &= 13 + (\lceil 13/8 \rceil * 2,000) + (((1/125) * 10,000 * 125)/4) = 28,513 \\ &= 13 + (2 * 2,000) + 2,500 = 6,513 \end{aligned}$$

As an exercise, the reader should perform a similar analysis for OP7.

### 19.5.3 Multirelation Queries and JOIN Ordering Choices

The algebraic transformation rules in Section 19.1.2 include a commutative rule and an associative rule for the join operation. With these rules, many equivalent join expressions can be produced. As a result, the number of alternative query trees grows very rapidly as the number of joins in a query increases. A query block that joins  $n$  relations will often have  $n - 1$  join operations, and hence can have a large number of different join orders. In general, for a query block that has  $n$  relations,

there are  $n!$  join orders; Cartesian products are included in this total number. Estimating the cost of every possible join tree for a query with a large number of joins will require a substantial amount of time by the query optimizer. Hence, some pruning of the possible query trees is needed. Query optimizers typically limit the structure of a (join) query tree to that of left-deep (or right-deep) trees. A **left-deep join tree** is a binary tree in which the right child of each non-leaf node is always a base relation. The optimizer would choose the particular left-deep join tree with the lowest estimated cost. Two examples of left-deep trees are shown in Figure 19.5(a). (Note that the trees in Figure 19.2 are also left-deep trees.) A **right-deep join tree** is a binary tree where the left child of every leaf node is a base relation (Figure 19.5(b)).

A **bushy join tree** is a binary tree where the left or right child of an internal node may be an internal node. Figure 19.5(b) shows a right-deep join tree whereas Figure 19.5(c) shows a bushy one using four base relations. Most query optimizers consider left-deep join trees as the preferred join tree and then choose one among the  $n!$  possible join orderings, where  $n$  is the number of relations. We discuss the join ordering issue in more detail in Sections 19.5.4 and 19.5.5. The left-deep tree has exactly one shape, and the join orders for  $N$  tables in a left-deep tree are given by  $N!$ . In contrast, the shapes of a bushy tree are given by the following recurrence relation (i.e., recursive function), with  $S(n)$  defined as follows:  $S(1) = 1$ .

$$S(n) = \sum_{i=1}^{n-1} S(i) * S(n - i)$$

The above recursive equation for  $S(n)$  can be explained as follows. It states that, for  $i$  between 1 and  $N - 1$  as the number of leaves in the left subtree, those leaves may be rearranged in  $S(i)$  ways. Similarly, the remaining  $N - i$  leaves in the right subtree can be rearranged in  $S(N - i)$  ways. The number of permutations of the bushy trees is given by:

$$P(n) = n! * S(n) = (2n - 2)!/(n - 1)!$$

Table 19.1 shows the number of possible left-deep (or right-deep) join trees and bushy join trees for joins of up to seven relations.

It is clear from Table 19.1 that the possible space of alternatives becomes rapidly unmanageable if all possible bushy tree alternatives were to be considered. In certain

**Figure 19.5**  
 (a) Two left-deep  
 (b) A right-deep  
 (c) A bushy query

cases like comp  
considering bus

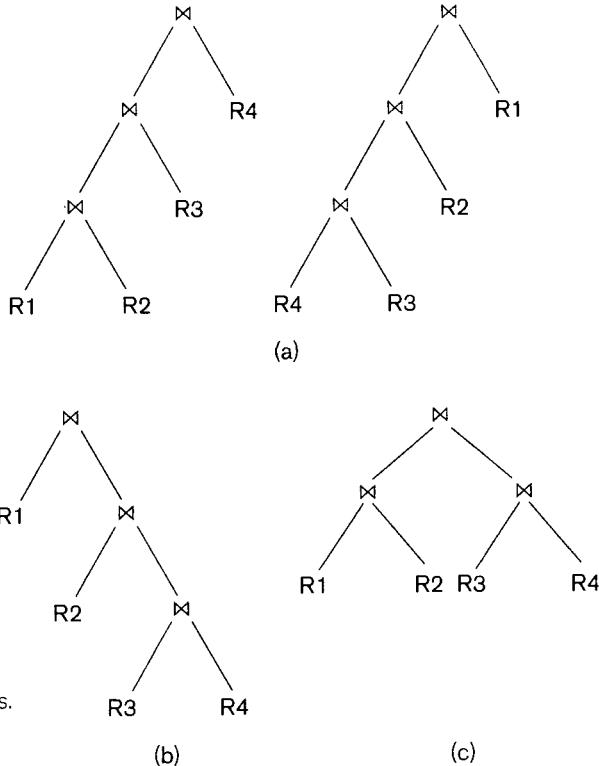
With left-deep t  
cutting a nested  
nested-loop jo  
amenable to pi  
left-deep tree it  
nested-loop mo  
probe the inner  
duced from the  
ing records for  
they could be  
is that having  
utilize any acc

If materializat  
results could b

**Table 19.1** Number of Permutations of Left-Deep and Bushy Join Trees of  $n$  Relations

No. of Relations $N$	No. of Left-Deep Trees $N!$	No. of Bushy Shapes $S(N)$	No. of Bushy Trees $(2N - 2)!/(N - 1)!$
2	2	1	2
3	6	2	12
4	24	5	120
5	120	14	1,680
6	720	42	30,240
7	5,040	132	665,280

<sup>12</sup>As a representa

**Figure 19.5**

- (a) Two left-deep join query trees.  
 (b) A right-deep join query tree.  
 (c) A bushy query tree.

cases like complex versions of snowflake schemas (see Section 29.3), approaches to considering bushy tree alternatives have been proposed.<sup>12</sup>

With left-deep trees, the right child is considered to be the inner relation when executing a nested-loop join, or the probing relation when executing an index-based nested-loop join. One advantage of left-deep (or right-deep) trees is that they are amenable to pipelining, as discussed in Section 18.7. For instance, consider the first left-deep tree in Figure 19.5(a) and assume that the join algorithm is the index-based nested-loop method; in this case, a disk page of tuples of the outer relation is used to probe the inner relation for matching tuples. As resulting tuples (records) are produced from the join of R1 and R2, they can be used to probe R3 to locate their matching records for joining. Likewise, as resulting tuples are produced from this join, they could be used to probe R4. Another advantage of left-deep (or right-deep) trees is that having a base relation as one of the inputs of each join allows the optimizer to utilize any access paths on that relation that may be useful in executing the join.

If materialization is used instead of pipelining (see Sections 18.7 and 19.2), the join results could be materialized and stored as temporary relations. The key idea from

<sup>12</sup>As a representative case for bushy trees, refer to Ahmed et al. (2014).

the optimizer's standpoint with respect to join ordering is to find an ordering that will reduce the size of the temporary results, since the temporary results (pipelined or materialized) are used by subsequent operators and hence affect the execution cost of those operators.

#### 19.5.4 Physical Optimization

For a given logical query plan based on the heuristics we have been discussing so far, each operation needs a further decision in terms of executing the operation by a specific algorithm at the physical level. This is referred to as **physical optimization**. If this optimization is based on the relative cost of each possible implementation, we call it cost-based physical optimization. The two sets of approaches to this decision making may be broadly classified as top-down and bottom-up approaches. In the **top-down** approach, we consider the options for implementing each operation working our way down the tree and choosing the best alternative at each stage. In the **bottom-up** approach, we consider the operations working up the tree, evaluating options for physical execution, and choosing the best at each stage. Theoretically, both approaches amount to evaluation of the entire space of possible implementation solutions to minimize the cost of evaluation; however, the bottom-up strategy lends itself naturally to pipelining and hence is used in commercial RDBMSs. Among the most important physical decisions is the ordering of join operations, which we will briefly discuss in Section 19.5.5. There are certain heuristics applied at the physical optimization stage that make elaborate cost computations unnecessary. These heuristics include:

- For selections, use index scans wherever possible.
  - If the selection condition is conjunctive, use the selection that results in the smallest cardinality first.
  - If the relations are already sorted on the attributes being matched in a join, then prefer sort-merge join to other join methods.
  - For union and intersection of more than two relations, use the associative rule; consider the relations in the ascending order of their estimated cardinalities.
  - If one of the arguments in a join has an index on the join attribute, use that as the inner relation.
  - If the left relation is small and the right relation is large and it has index on the joining column, then try index-based nested-loop join.
  - Consider only those join orders where there are no Cartesian products or where all joins appear before Cartesian products.

The following are only some of the types of physical level heuristics used by the optimizer. If the number of relations is small (typically less than 6) and, therefore, possible implementations options are limited, then most optimizers would elect to apply a cost-based optimization approach directly rather than to explore heuristics.

## 19.5.5 Dynan

We saw in Section 1.2 how to implement an *n*-way join. Every join can be implemented as a sequence of possible permutations of the inputs. This leads to bushy trees except for the preferred (over 1) permutation. First, they work with index-based nested loops, and finally with fully pipelined joins. Note that inner loops are not yet considered. The join algorithms presented here do not consider the join attribute. They are based on the assumption that all attributes are available at every node.

The common approach is to use heuristic approximation techniques to solve a problem that may be typical dynamic.

1. The stru-
  2. The valu
  3. The opti

Note that the so absolute optima the join order set  $r_2, r_3, r_4, r_5$ . The cost of each of the programming tasks is manageable. We know two solutions. Note that the cost would be 12 of the

$$r_1 \bowtie r_2 \bowtie \dots$$

The 6 (= 3!) possible options of the next join:

(temp1  $\bowtie$  r)

If we were to consider the 6

<sup>13</sup>For a detailed account on algorithm

<sup>14</sup>Based on Chai

### 19.5.5 Dynamic Programming Approach to Join Ordering

We saw in Section 19.5.3 that there are many possible ways to order  $n$  relations in an  $n$ -way join. Even for  $n = 5$ , which is not uncommon in practical applications, the possible permutations are 120 with left-deep trees and 1,680 with bushy trees. Since bushy trees expand the solution space tremendously, left-deep trees are generally preferred (over both bushy and right-deep trees). They have multiple advantages: First, they work well with the common algorithms for join, including nested-loop, index-based nested-loop, and other one-pass algorithms. Second, they can generate **fully pipelined plans** (i.e., plans where all joins can be evaluated using pipelining). Note that inner tables must always be materialized because in the join implementation algorithms, the entire inner table is needed to perform the matching on the join attribute. This is not possible with right-deep trees.

The common approach to evaluate possible permutations of joining relations is a greedy heuristic approach called dynamic programming. **Dynamic programming** is an optimization technique<sup>13</sup> where subproblems are solved only once, and it is applicable when a problem may be broken down into subproblems that themselves have subproblems. A typical dynamic programming algorithm has the following characteristics<sup>14</sup>:

1. The structure of an optimal solution is developed.
2. The value of the optimal solution is recursively defined.
3. The optimal solution is computed and its value developed in a bottom-up fashion.

Note that the solution developed by this procedure is an optimal solution and not the absolute optimal solution. To consider how dynamic programming may be applied to the join order selection, consider the problem of ordering a 5-way join of relations  $r_1, r_2, r_3, r_4, r_5$ . This problem has 120 ( $=5!$ ) possible left-deep tree solutions. Ideally, the cost of each of them can be estimated and compared and the best one selected. Dynamic programming takes an approach that breaks down this problem to make it more manageable. We know that for three relations, there are only six possible left-deep tree solutions. Note that if all possible bushy tree join solutions were to be evaluated, there would be 12 of them. We can therefore consider the join to be broken down as:

$$r_1 \bowtie r_2 \bowtie r_3 \bowtie r_4 \bowtie r_5 = (r_1 \bowtie r_2 \bowtie r_3) \bowtie r_4 \bowtie r_5$$

The 6 ( $=3!$ ) possible options of  $(r_1 \bowtie r_2 \bowtie r_3)$  may then be combined with the 6 possible options of taking the result of the first join, say,  $\text{temp1}$ , and then considering the next join:

$$(\text{temp1} \bowtie r_4 \bowtie r_5)$$

If we were to consider the 6 options for evaluating  $\text{temp1}$  and, for each of them, consider the 6 options of evaluating the second join  $(\text{temp1} \bowtie r_4 \bowtie r_5)$ , the possible

---

<sup>13</sup>For a detailed discussion of dynamic programming as a technique of optimization, the reader may consult an algorithm textbook such as Cormen et al. (2003).

<sup>14</sup>Based on Chapter 16 in Cormen et al. (2003).

solution space has  $6 * 6 = 36$  alternatives. This is where dynamic programming can be used to do a sort of greedy optimization. It takes the “optimal” plan for evaluating  $\text{temp1}$  and does *not* revisit that plan. So the solution space now reduces to only 6 options to be considered for the second join. Thus the total number of options considered becomes  $6 + 6$  instead of 120 (=5!) in the nonheuristic exhaustive approach.

The order in which the result of the join is generated is also important for finding the best overall order of joins since for using sort-merge join with the next relation, it plays an important role. The ordering beneficial for the next join is considered an **interesting join order**. This approach was first proposed in System R at IBM Research.<sup>15</sup> Besides the join attributes of the later join, System R also included grouping attributes of a later GROUP BY or a sort order at the root of the tree among interesting sort orders. For example, in the case we discussed above, the interesting join orders for the  $\text{temp1}$  relation will include those that match the join attribute(s) required to join with either  $r4$  or with  $r5$ . The dynamic programming algorithm can be extended to consider best join orders for each interesting sort order. The number of subsets of  $n$  relations is  $2^n$  (for  $n = 5$  it is 32;  $n = 10$  gives 1,024, which is still manageable), and the number of interesting join orders is small. The complexity of the extended dynamic programming algorithm to determine the optimal left-deep join tree permutation has been shown to be  $O(3^n)$ .

- 3. DEPAR
- 4. EMPLO

Assume that  $\text{tl}$   
tion. If we assi  
ated after each  
between PROJ1  
for the input  
according to F  
linear search).  
before the joi  
PROJ\_PLOC i  
two options. T  
shows the num

**Figure 19.6**  
Sample statistics  
(b) Table informa

(a)

## 19.6 Example to Illustrate Cost-Based Query Optimization

We will consider query Q2 and its query tree shown in Figure 19.1(a) to illustrate cost-based query optimization:

```

Q2: SELECT Pnumber, Dnum, Lname, Address, Bdate
      FROM PROJECT, DEPARTMENT, EMPLOYEE
      WHERE Dnum=Dnumber AND Mgr_ssn=Ssn AND
            Plocation='Stafford';
    
```

Suppose we have the information about the relations shown in Figure 19.6. The **LOW\_VALUE** and **HIGH\_VALUE** statistics have been normalized for clarity. The tree in Figure 19.1(a) is assumed to represent the result of the algebraic heuristic optimization process and the start of cost-based optimization (in this example, we assume that the heuristic optimizer does not push the projection operations down the tree).

(b)

The first cost-based optimization to consider is join ordering. As previously mentioned, we assume the optimizer considers only left-deep trees, so the potential join orders—without **CARTESIAN PRODUCT**—are:

1. PROJECT  $\bowtie$  DEPARTMENT  $\bowtie$  EMPLOYEE
2. DEPARTMENT  $\bowtie$  PROJECT  $\bowtie$  EMPLOYEE

(c)

<sup>15</sup>See the classic reference in this area by Selinger et al. (1979).

3. DEPARTMENT  $\bowtie$  EMPLOYEE  $\bowtie$  PROJECT
4. EMPLOYEE  $\bowtie$  DEPARTMENT  $\bowtie$  PROJECT

Assume that the selection operation has already been applied to the PROJECT relation. If we assume a materialized approach, then a new temporary relation is created after each join operation. To examine the cost of join order (1), the first join is between PROJECT and DEPARTMENT. Both the join method and the access methods for the input relations must be determined. Since DEPARTMENT has no index according to Figure 19.6, the only available access method is a table scan (that is, a linear search). The PROJECT relation will have the selection operation performed before the join, so two options exist—table scan (linear search) or use of the PROJ\_PLOC index—so the optimizer must compare the estimated costs of these two options. The statistical information on the PROJ\_PLOC index (see Figure 19.6) shows the number of index levels  $x = 2$  (root plus leaf levels). The index is nonunique

**Figure 19.6**

Sample statistical information for relations in Q2. (a) Column information.  
 (b) Table information. (c) Index information.

(a)	Table_name	Column_name	Num_distinct	Low_value	High_value
PROJECT	Plocation	200	1	200	
PROJECT	Pnumber	2000	1	2000	
PROJECT	Dnum	50	1	50	
DEPARTMENT	Dnumber	50	1	50	
DEPARTMENT	Mgr_ssn	50	1	50	
EMPLOYEE	Ssn	10000	1	10000	
EMPLOYEE	Dno	50	1	50	
EMPLOYEE	Salary	500	1	500	

(b)	Table_name	Num_rows	Blocks
PROJECT	2000	100	
DEPARTMENT	50	5	
EMPLOYEE	10000	2000	

(c)	Index_name	Uniqueness	Blevel*	Leaf_blocks	Distinct_keys
PROJ_PLOC	NONUNIQUE	1	4	200	
EMP_SSN	UNIQUE	1	50	10000	
EMP_SAL	NONUNIQUE	1	50	500	

\*Blevel is the number of levels without the leaf level.

(because Plocation is not a key of PROJECT), so the optimizer assumes a uniform data distribution and estimates the number of record pointers for each Plocation value to be 10. This is computed from the tables in Figure 19.6 by multiplying Selectivity \* Num\_rows, where Selectivity is estimated by 1/Num\_distinct. So the cost of using the index and accessing the records is estimated to be 12 block accesses (2 for the index and 10 for the data blocks). The cost of a table scan is estimated to be 100 block accesses, so the index access is more efficient as expected.

In the materialized approach, a temporary file TEMP1 of size 1 block is created to hold the result of the selection operation. The file size is calculated by determining the blocking factor using the formula Num\_rows/Blocks, which gives 2,000/100 or 20 rows per block. Hence, the 10 records selected from the PROJECT relation will fit into a single block. Now we can compute the estimated cost of the first join. We will consider only the nested-loop join method, where the outer relation is the temporary file, TEMP1, and the inner relation is DEPARTMENT. Since the entire TEMP1 file fits in the available buffer space, we need to read each of the DEPARTMENT table's five blocks only once, so the join cost is six block accesses plus the cost of writing the temporary result file, TEMP2. The optimizer would have to determine the size of TEMP2. Since the join attribute Dnumber is the key for DEPARTMENT, any Dnum value from TEMP1 will join with at most one record from DEPARTMENT, so the number of rows in TEMP2 will be equal to the number of rows in TEMP1, which is 10. The optimizer would determine the record size for TEMP2 and the number of blocks needed to store these 10 rows. For brevity, assume that the blocking factor for TEMP2 is five rows per block, so a total of two blocks are needed to store TEMP2.

Finally, the cost of the last join must be estimated. We can use a single-loop join on TEMP2 since in this case the index EMP\_SSN (see Figure 19.6) can be used to probe and locate matching records from EMPLOYEE. Hence, the join method would involve reading in each block of TEMP2 and looking up each of the five Mgr\_ssn values using the EMP\_SSN index. Each index lookup would require a root access, a leaf access, and a data block access ( $x + 1$ , where the number of levels  $x$  is 2). So, 10 lookups require 30 block accesses. Adding the two block accesses for TEMP2 gives a total of 32 block accesses for this join.

For the final projection, assume pipelining is used to produce the final result, which does not require additional block accesses, so the total cost for join order (1) is estimated as the sum of the previous costs. The optimizer would then estimate costs in a similar manner for the other three join orders and choose the one with the lowest estimate. We leave this as an exercise for the reader.

## 19.7 Additional Issues Related to Query Optimization

In this section, we will discuss a few issues of interest that we have not been able to discuss earlier.

### 19.7.1 Disk Usage

Most commercial DBMSs implement disk usage by the query cost model, which means they try to minimize the cost of disk access. There is some variation between different DBMSs.

#### ■ Oracle

EXPLAIN PLAN FOR <SQL>

The query may go into a table called PLAN\_TABLE or UTLXPLP.SQL respectively.

#### ■ IBM DB2

EXPLAIN

There is no plan command; instead, explain tables is used during the plan phase. The sample command is:

#### ■ SQL Server

SET STATISTICS SHOW

The above statistics output is present in the results of the above three operations.

#### ■ PostgreSQL

EXPLAIN ANALYZE

### 19.7.2 Size Estimation

In Sections 19.4 and 19.5, we discussed how size estimation is done. In this section, we consider the issue of estimating the size of the result of a query.

**Projection:** For a query  $R_1 \times R_2 \times \dots \times R_n \rightarrow R$ , the size of the result is  $|R|$ .

<sup>16</sup>We have just illustrated this for the Oracle DBMS.