

6

DATABASE APPLICATION DEVELOPMENT

- How do application programs connect to a DBMS?
- How can applications manipulate data retrieved from a DBMS?
- How can applications modify data in a DBMS?
- What are cursors?
- What is JDBC and how is it used?
- What is SQLJ and how is it used?
- What are stored procedures?
- **Key concepts:** Embedded SQL, Dynamic SQL, cursors; JDBC, connections, drivers, ResultSets, java.sql, SQLJ; stored procedures, SQL/PSM

He profits most who serves best.

-----Ivlotto for Rotary International

In Chapter 5, we looked at a wide range of SQL query constructs, treating SQL as an independent language in its own right. A relational DBMS supports an *interactive SQL* interface, and users can directly enter SQL commands. This simple approach is fine as long as the task at hand can be accomplished entirely with SQL commands. In practice, we often encounter situations in which we need the greater flexibility of a general-purpose programming language in addition to the data manipulation facilities provided by SQL. For example, we may want to integrate a database application with a nice graphical user interface, or we may want to integrate with other existing applications.

Applications that rely on the DBMS to manage data run as separate processes that connect to the DBMS to interact with it. Once a connection is established, SQL commands can be used to insert, delete, and modify data. SQL queries can be used to retrieve desired data, but we need to bridge an important difference in how a database system sees data and how an application program in a language like Java or C sees data: The result of a database query is a set (or multiset) of records, but Java has no set or multiset data type. This mismatch is resolved through additional SQL constructs that allow applications to obtain a handle on a collection and iterate over the records one at a time.

We introduce Embedded SQL, Dynamic SQL, and cursors in Section 6.1. Embedded SQL allows us to access data using static SQL queries in application code (Section 6.1.1); with Dynamic SQL, we can create the queries at run-time (Section 6.1.3). Cursors bridge the gap between set-valued query answers and programming languages that do not support set-values (Section 6.1.2).

The emergence of Java as a popular application development language, especially for Internet applications, has made accessing a DBMS from Java code a particularly important topic. Section 6.2 covers JDBC, a programming interface that allows us to execute SQL queries from a Java program and use the results in the Java program. JDBC provides greater portability than Embedded SQL or Dynamic SQL, and offers the ability to connect to several DBMSs without recompiling the code. Section 6.4 covers SQLJ, which does the same for static SQL queries, but is easier to program in than Java, with JDBC.

Often, it is useful to execute application code at the database server, rather than just retrieve data and execute application logic in a separate process. Section 6.5 covers stored procedures, which enable application logic to be stored and executed at the database server. We conclude the chapter by discussing our B&N case study in Section 6.6.

While writing database applications, we must also keep in mind that typically many application programs run concurrently. The transaction concept, introduced in Chapter 1, is used to encapsulate the effects of an application on the database. An application can select certain transaction properties through SQL commands to control the degree to which it is exposed to the changes of other concurrently running applications. We touch on the transaction concept at many points in this chapter, and, in particular, cover transaction-related aspects of JDBC. A full discussion of transaction properties and SQL's support for transactions is deferred until Chapter 16.

Examples that appear in this chapter are available online at

<http://www.cs.wisc.edu/-dbbook>

6.1 ACCESSING DATABASES **FROM** APPLICATIONS

In this section, we cover how SQL commands can be executed from within a program in a host language such as C or Java. The use of SQL commands within a host language program is called Embedded SQL. Details of Embedded SQL also depend on the host language. Although similar capabilities are supported for a variety of host languages, the syntax sometimes varies.

We first cover the basics of Embedded SQL with static SQL queries in Section 6.1.1. We then introduce cursors in Section 6.1.2. We discuss Dynamic SQL, which allows us to construct SQL queries at runtime (and execute them) in Section 6.1.3.

6.1.1 Embedded SQL

Conceptually, embedding SQL commands in a host language program is straightforward. SQL statements (i.e., not declarations) can be used wherever a statement in the host language is allowed (with a few restrictions). SQL statements must be clearly marked so that a preprocessor can deal with them before invoking the compiler for the host language. Also, any host language variables used to pass arguments into an SQL command must be declared in SQL. In particular, some special host language variables *must* be declared in SQL (so that, for example, any error conditions arising during SQL execution can be communicated back to the main application program in the host language).

There are, however, two complications to bear in mind. First, the data types recognized by SQL may not be recognized by the host language and vice versa. This mismatch is typically addressed by casting data values appropriately before passing them to or from SQL commands. (SQL, like other programming languages, provides an operator to cast values of one type into values of another type.) The second complication has to do with SQL being set-oriented, and is addressed using cursors (see Section 6.1.2. Commands operate on and produce tables, which are sets

In our discussion of Embedded SQL, we assume that the host language is C for concreteness, because minor differences exist in how SQL statements are embedded in different host languages.

Declaring Variables and Exceptions

SQL statements can refer to variables defined in the host program. Such host-language variables must be prefixed by a colon (:) in SQL statements and be declared between the commands EXEC SQL BEGIN DECLARE SECTION and EXEC

SQL END DECLARE SECTION. The declarations are similar to how they would look in a C program and, as usual in C, are separated by semicolons. For example, we can declare variables *c_sname*, *c_sid*, *c_rating*, and *c_age* (with the initial *c* used as a naming convention to emphasize that these are host language variables) as follows:

```
EXEC SQL BEGIN DECLARE SECTION
char c_sname[20];
long c_sid;
short c_rating;
float c_age;
EXEC SQL END DECLARE SECTION
```

The first question that arises is which SQL types correspond to the various C types, since we have just declared a collection of C variables whose values are intended to be read (and possibly set) in an SQL run-time environment when an SQL statement that refers to them is executed. The SQL-92 standard defines such a correspondence between the host language types and SQL types for a number of host languages. In our example, *c_sname* has the type CHARACTER(20) when referred to in an SQL statement, *c_sid* has the type INTEGER, *c_rating* has the type SMALLINT, and *c_age* has the type REAL.

We also need some way for SQL to report what went wrong if an error condition arises when executing an SQL statement. The SQL-92 standard recognizes two special variables for reporting errors, SQLCODE and SQLSTATE. SQLCODE is the older of the two and is defined to return some negative value when an error condition arises, without specifying further just what error a particular negative integer denotes. SQLSTATE, introduced in the SQL-92 standard for the first time, associates predefined values with several common error conditions, thereby introducing some uniformity to how errors are reported. One of these two variables *must* be declared. The appropriate C type for SQLCODE is long and the appropriate C type for SQLSTATE is char[6], that is, a character string five characters long. (Recall the null-terminator in C strings.) In this chapter, we assume that SQLSTATE is declared.

Embedding SQL Statements

All SQL statements embedded within a host program must be clearly marked, with the details dependent on the host language; in C, SQL statements must be prefixed by EXEC SQL. An SQL statement can essentially appear in any place in the host language program where a host language statement can appear.

As a simple example, the following Embedded SQL statement inserts a row, whose column values are based on the values of the host language variables contained in it, into the Sailors relation:

```
EXEC SQL
INSERT INTO Sailors VALUES (:c_sname, :c_sid, :c_rating, :c_age);
```

Observe that a semicolon terminates the command, as per the convention for terminating statements in C.

The SQLSTATE variable should be checked for errors and exceptions after each Embedded SQL statement. SQL provides the WHENEVER command to simplify this tedious task:

```
EXEC SQL WHENEVER [SQLERROR | NOT FOUND] [ CONTINUE | GOTO st'mt ]
```

The intent is that the value of SQLSTATE should be checked after each Embedded SQL statement is executed. If SQLERROR is specified and the value of SQLSTATE indicates an exception, control is transferred to *stmt*, which is presumably responsible for error and exception handling. Control is also transferred to *stmt* if NOT FOUND is specified and the value of SQLSTATE is 02000, which denotes NO DATA.

6.1.2 Cursors

A major problem in embedding SQL statements in a host language like C is that an *impedance mismatch* occurs because SQL operates on *sets* of records, whereas languages like C do not cleanly support a set-of-records abstraction. The solution is to essentially provide a mechanism that allows us to retrieve rows one at a time from a relation.

This mechanism is called a cursor. We can declare a cursor on any relation or on any SQL query (because every query returns a set of rows). Once a cursor is declared, we can open it (which positions the cursor just before the first row); fetch the next row; move the cursor (to the next row, to the row after the next *n*, to the first row, or to the previous row, etc., by specifying additional parameters for the FETCH command); or close the cursor. Thus, a cursor essentially allows us to retrieve the rows in a table by positioning the cursor at a particular row and reading its contents.

Basic Cursor Definition and Usage

Cursors enable us to examine, in the host language program, a collection of rows computed by an Embedded SQL statement:

- We usually need to open a cursor if the embedded statement is a SELECT (i.e.) a query). However, we can avoid opening a cursor if the answer contains a single row, as we see shortly.
- INSERT, DELETE, and UPDATE statements typically require no cursor, although some variants of DELETE and UPDATE use a cursor.

As an example, we can find the name and age of a sailor, specified by assigning a value to the host variable *c_sid*, declared earlier, as follows:

```
EXEC SQL SELECT S.sname, S.age
        INTO   :c_sname, :c_age
        FROM   Sailors S
        WHERE  S.sid = :c_sid;
```

The INTO clause allows us to assign the columns of the single answer row to the host variables *c_sname* and *c_age*. Therefore, we do not need a cursor to embed this query in a host language program. But what about the following query, which computes the names and ages of all sailors with a rating greater than the current value of the host variable *c_minrating*?

```
SELECT S.sname, S.age
FROM   Sailors S
WHERE  S.rating > :c_minrating
```

This query returns a collection of rows, not just one row. When executed interactively, the answers are printed on the screen. If we embed this query in a C program by prefixing the command with EXEC SQL, how can the answers be bound to host language variables? The INTO clause is inadequate because we must deal with several rows. The solution is to use a cursor:

```
DECLARE sinfo CURSOR FOR
SELECT S.sname, S.age
FROM   Sailors S
WHERE  S.rating > :c_minrating;
```

This code can be included in a C program, and once it is executed, the cursor *sinfo* is defined. Subsequently, we can open the cursor:

```
OPEN sinfo;
```

The value of *c_minrating* in the SQL query associated with the cursor is the value of this variable when we open the cursor. (The cursor declaration is processed at compile-time, and the OPEN command is executed at run-time.)

A cursor can be thought of as 'pointing' to a row in the collection of answers to the query associated with it. When a cursor is opened, it is positioned just before the first row. We can use the `FETCH` command to read the first row of cursor *sinfo* into host language variables:

```
FETCH sinfo INTO :c_sname, :c_age;
```

When the `FETCH` statement is executed, the cursor is positioned to point at the next row (which is the first row in the table when `FETCH` is executed for the first time after opening the cursor) and the column values in the row are copied into the corresponding host variables. By repeatedly executing this `FETCH` statement (say, in a `while`-loop in the C program), we can read all the rows computed by the query, one row at a time. Additional parameters to the `FETCH` command allow us to position a cursor in very flexible ways, but we do not discuss them.

How do we know when we have looked at all the rows associated with the cursor? By looking at the special variables `SQLCODE` or `SQLSTATE`, of course. `SQLSTATE`, for example, is set to the value `02000`, which denotes `NO DATA`, to indicate that there are no more rows if the `FETCH` statement positions the cursor after the last row.

When we are done with a cursor, we can close it:

```
CLOSE sinfo;
```

It can be opened again if needed, and the value of `: c_minrating` in the SQL query associated with the cursor would be the value of the host variable *c_minrating* at that time.

Properties of Cursors

The general form of a cursor declaration is:

```
DECLARE cursorname [INSENSITIVE] [SCROLL] CURSOR
        [WITH HOLD]
        FOR some query
        [ ORDER BY order-item-list ]
        [ FOR READ ONLY | FOR UPDATE ]
```

A cursor can be declared to be a read-only cursor (`FOR READ ONLY`) or, if it is a cursor on a base relation or an updatable view, to be an updatable cursor (`FOR UPDATE`). If it is updatable, simple variants of the `UPDATE` and

DELETE commands allow us to update or delete the row on which the cursor is positioned. For example, if *sinfa* is an updatable cursor and open, we can execute the following statement:

```
UPDATE Sailors S
SET      S.rating = S.rating - 1
WHERE    CURRENT of sinfo;
```

This Embedded SQL statement modifies the *rating* value of the row currently pointed to by cursor *sinfa*; similarly, we can delete this row by executing the next statement:

```
DELETE Sailors S
WHERE    CURRENT of sinfo;
```

A cursor is updatable by default unless it is a scrollable or insensitive cursor (see below), in which case it is read-only by default.

If the keyword SCROLL is specified, the cursor is scrollable, which means that variants of the FETCH command can be used to position the cursor in very flexible ways; otherwise, only the basic FETCH command, which retrieves the next row, is allowed.

If the keyword INSENSITIVE is specified, the cursor behaves as if it is ranging over a private copy of the collection of answer rows. Otherwise, and by default, other actions of some transaction could modify these rows, creating unpredictable behavior. For example, while we are fetching rows using the *sinfa* cursor, we might modify *rating* values in Sailor rows by concurrently executing the command:

```
UPDATE Sailors S
SET      S.rating = S.rating -
```

Consider a Sailor row such that (1) it has not yet been fetched, and (2) its original *rating* value would have met the condition in the WHERE clause of the query associated with *sinfa*, but the new *rating* value does not. Do we fetch such a Sailor row? If INSENSITIVE is specified, the behavior is as if all answers were computed and stored when *sinfo* was opened; thus, the update command has no effect on the rows fetched by *sinfa* if it is executed after *sinfo* is opened. If INSENSITIVE is not specified, the behavior is implementation dependent in this situation.

A holdable cursor is specified using the WITH HOLD clause, and is not closed when the transaction is conunitted. The motivation for this comes from long

transactions in which we access (and possibly change) a large number of rows of a table. If the transaction is aborted for any reason, the system potentially has to redo a lot of work when the transaction is restarted. Even if the transaction is not aborted, its locks are held for a long time and reduce the concurrency of the system. The alternative is to break the transaction into several smaller transactions, but remembering our position in the table between transactions (and other similar details) is complicated and error-prone. Allowing the application program to commit the transaction it initiated, while retaining its handle on the active table (i.e., the cursor) solves this problem: The application can commit its transaction and start a new transaction and thereby save the changes it has made thus far.

Finally, in what order do `FETCH` commands retrieve rows? In general this order is unspecified, but the optional `ORDER BY` clause can be used to specify a sort order. Note that columns mentioned in the `ORDER BY` clause cannot be updated through the cursor!

The order-item-list is a list of order-items; an order-item is a column name, optionally followed by one of the keywords `ASC` or `DESC`. Every column mentioned in the `ORDER BY` clause must also appear in the select-list of the query associated with the cursor; otherwise it is not clear what columns we should sort on. The keywords `ASC` or `DESC` that follow a column control whether the result should be sorted-with respect to that column-in ascending or descending order; the default is `ASC`. This clause is applied as the last step in evaluating the query.

Consider the query discussed in Section 5.5.1, and the answer shown in Figure 5.13. Suppose that a cursor is opened on this query, with the clause:

`ORDER BY minage ASC, rating DESC`

The answer is sorted first in ascending order by *minage*, and if several rows have the same *minage* value, these rows are sorted further in descending order by *rating*. The cursor would fetch the rows in the order shown in Figure 6.1.

<i>rating</i>	<i>minage</i>
8	25.5
3	25.5
7	35.0

Figure 6.1 Order in which Tuples Are Fetched

6.1.3 Dynamic SQL

Consider an application such as a spreadsheet or a graphical front-end that needs to access data from a DBMS. Such an application must accept commands from a user and, based on what the user needs, generate appropriate SQL statements to retrieve the necessary data. In such situations, we may not be able to predict in advance just what SQL statements need to be executed, even though there is (presumably) some algorithm by which the application can construct the necessary SQL statements once a user's command is issued.

SQL provides some facilities to deal with such situations; these are referred to as Dynamic SQL. We illustrate the two main commands, PREPARE and EXECUTE, through a simple example:

```
char c_sqlstring[] = {"DELETE FROM Sailors WHERE rating>5"};
EXEC SQL PREPARE readytogo FROM :c_sqlstring;
EXEC SQL EXECUTE readytogo;
```

The first statement declares the C variable *c_sqlstring* and initializes its value to the string representation of an SQL command. The second statement results in this string being parsed and compiled as an SQL command, with the resulting executable bound to the SQL variable *readytogo*. (Since *readytogo* is an SQL variable, just like a cursor name, it is not prefixed by a colon.) The third statement executes the command.

Many situations require the use of Dynamic SQL. However, note that the preparation of a Dynamic SQL command occurs at run-time and is run-time overhead. Interactive and Embedded SQL commands can be prepared once at compile-time and then re-executed as often as desired. Consequently you should limit the use of Dynamic SQL to situations in which it is essential.

There are many more things to know about Dynamic SQL—how we can pass parameters from the host language program to the SQL statement being prepared, for example—but we do not discuss it further.

6.2 AN INTRODUCTION TO JDBC

Embedded SQL enables the integration of SQL with a general-purpose programming language. As described in Section 6.1.1, a DBMS-specific preprocessor transforms the Embedded SQL statements into function calls in the host language. The details of this translation vary across DBMSs, and therefore even though the source code can be compiled to work with different DBMSs, the final executable works only with one specific DBMS.

```
while (true) {  
    #sql { FETCH :books INTO :title, :price, };  
    if (books.endFetch()) {  
        break;  
    }  
  
    // process the book  
}
```

6.5 STORED PROCEDURES

It is often important to execute some parts of the application logic directly in the process space of the database system. Running application logic directly at the database has the advantage that the amount of data that is transferred between the database server and the client issuing the SQL statement can be minimized, while at the same time utilizing the full power of the database server.

When SQL statements are issued from a remote application, the records in the result of the query need to be transferred from the database system back to the application. If we use a cursor to remotely access the results of an SQL statement, the DBMS has resources such as locks and memory tied up while the application is processing the records retrieved through the cursor. In contrast, a stored procedure is a program that is executed through a single SQL statement that can be locally executed and completed within the process space of the database server. The results can be packaged into one big result and returned to the application, or the application logic can be performed directly at the server, without having to transmit the results to the client at all.

Stored procedures are also beneficial for software engineering reasons. Once a stored procedure is registered with the database server, different users can re-use the stored procedure, eliminating duplication of efforts in writing SQL queries or application logic, and making code maintenance easy. In addition, application programmers do not need to know the database schema if we encapsulate all database access into stored procedures.

Although they are called stored *procedures*, they do not have to be procedures in a programming language sense; they can be functions.

6.5.1 Creating a Simple Stored Procedure

Let us look at the example stored procedure written in SQL shown in Figure 6.8. We see that stored procedures must have a name; this stored procedure

has the name 'ShowNumberOfOrders.' Otherwise, it just contains an SQL statement that is precompiled and stored at the server.

```
CREATE PROCEDURE ShowNumberOfOrders
SELECT C.cid, C.cname, COUNT(*)
FROM   Customers C, Orders o
WHERE  C.cid = O.cid
GROUP BY C.cid, C.cname
```

Figure 6.8 A Stored Procedure in SQL

Stored procedures can also have parameters. These parameters have to be valid SQL types, and have one of three different modes: IN, OUT, or INOUT. IN parameters are arguments to the stored procedure. OUT parameters are returned from the stored procedure; it assigns values to all OUT parameters that the user can process. INOUT parameters combine the properties of IN and OUT parameters: They contain values to be passed to the stored procedures, and the stored procedure can set their values as return values. Stored procedures enforce strict type conformance: If a parameter is of type INTEGER, it cannot be called with an argument of type VARCHAR.

Let us look at an example of a stored procedure with arguments. The stored procedure shown in Figure 6.9 has two arguments: book_isbn and addedQty. It updates the available number of copies of a book with the quantity from a new shipment.

```
CREATE PROCEDURE AddInventory (
    IN book_isbn CHAR(10),
    IN addedQty INTEGER)
UPDATE Books
SET    qty_in_stock = qtyjn_stock + addedQty
WHERE  bookisbn = isbn
```

Figure 6.9 A Stored Procedure with Arguments

Stored procedures do not have to be written in SQL; they can be written in any host language. As an example, the stored procedure shown in Figure 6.10 is a Java function that is dynamically executed by the database server whenever it is called by the client:

6.5.2 Calling Stored Procedures

Stored procedures can be called in interactive SQL with the CALL statement:

```
CREATE PROCEDURE RallkCustomers(IN number INTEGER)
LANGUAGE Java
EXTERNAL NAME 'file:///c:/storedProcedures/rank.jar'
```

Figure 6.10 A Stored Procedure in Java

```
CALL storedProcedureName(argument1, argument2, ... , argumentN);
```

In Embedded SQL, the arguments to a stored procedure are usually variables in the host language. For example, the stored procedure `AddInventory` would be called as follows:

```
EXEC SQL BEGIN DECLARE SECTION
char isbn[10];
long qty;
EXEC SQL END DECLARE SECTION

// set isbn and qty to some values
EXEC SQL CALL AddInventory(:isbn,:qty);
```

Calling Stored Procedures from JDBC

We can call stored procedures from JDBC using the `CallableStatement` class. `CallableStatement` is a subclass of `PreparedStatement` and provides the same functionality. A stored procedure could contain multiple SQL statements or a series of SQL statements—thus, the result could be many different `ResultSet` objects. We illustrate the case when the stored procedure result is a single `ResultSet`.

```
CallableStatement cstmt=
    con.prepareCall(" {call ShowNumberOfOrders}");
ResultSet rs = cstmt.executeQuery()
while (rs.next())
```

Calling Stored Procedures from SQLJ

The stored procedure 'ShowNumberOfOrders' is called as follows using SQLJ:

```
// create the cursor class
#sql !iterator CustomerInfo(int cid, String cname, int count);

// create the cursor
```

```

CustomerInfo customerinfo;

// call the stored procedure
#sql customerinfo = {CALL ShowNumberOfOrders};
while (customerinfo.nextO) {
    System.out.println(customerinfo.cid() + "," +
                        customerinfo.count());
}

```

6.5.3 SQL/PSM

All major database systems provide ways for users to write stored procedures in a simple, general purpose language closely aligned with SQL. In this section, we briefly discuss the SQL/PSM standard, which is representative of most vendor-specific languages. In PSM, we define modules, which are collections of stored procedures, temporary relations, and other declarations.

In SQL/PSM, we declare a stored procedure as follows:

```

CREATE PROCEDURE name (parameter1,..., parameterN)
    local variable declarations
    procedure code;

```

We can declare a function similarly as follows:

```

CREATE FUNCTION name (parameter1,..., parameterN)
    RETURNS sqlDataType
    local variable declarations
    function code;

```

Each parameter is a triple consisting of the mode (IN, OUT, or INOUT as discussed in the previous section), the parameter name, and the SQL datatype of the parameter. We can see very simple SQL/PSM procedures in Section 6.5.1. In this case, the local variable declarations were empty, and the procedure code consisted of an SQL query.

We start out with an example of a SQL/PSM function that illustrates the main SQL/PSM constructs. The function takes as input a customer identified by her *cid* and a year. The function returns the rating of the customer, which is defined as follows: Customers who have bought more than ten books during the year are rated 'two'; customer who have purchased between 5 and 10 books are rated 'one', otherwise the customer is rated 'zero'. The following SQL/PSM code computes the rating for a given customer and year.

```

CREATE PROCEDURE RateCustomer

```

Database Application Development

```
(IN custId INTEGER, IN year INTEGER)
RETURNS INTEGER
DECLARE rating INTEGER;
DECLARE numOrders INTEGER;
SET numOrders =
    (SELECT COUNT(*) FROM Orders O WHERE O.tid = custId);
IF (numOrders>10) THEN rating=2;
ELSEIF (numOrders>5) THEN rating=1;
ELSE rating=0;
END IF;
RETURN rating;
```

Let us use this example to give a short overview of some SQL/PSM constructs:

- We can declare local variables using the DECLARE statement. In our example, we declare two local variables: 'rating', and 'numOrders'.
- PSM/SQL functions return values via the RETURN statement. In our example, we return the value of the local variable 'rating'.
- We can assign values to variables with the SET statement. In our example, we assigned the return value of a query to the variable 'numOrders'.
- SQL/PSM has branches and loops. Branches have the following form:

```
IF (condition) THEN statements;
ELSEIF statements;

ELSEIF statements;
ELSE statements; END IF
```

Loops are of the form

```
LOOP
    statements;
END LOOP
```

- Queries can be used as part of expressions in branches; queries that return a single value can be assigned to variables as in our example above.
- We can use the same cursor statements as in Embedded SQL (OPEN, FETCH, CLOSE), but we do not need the EXEC SQL constructs, and variables do not have to be prefixed by a colon ':'.

We only gave a very short overview of SQL/PSM; the references at the end of the chapter provide more information.



20

PHYSICAL DATABASE DESIGN AND TUNING

- ☛ What is physical database design?
- ☛ What is a query workload?
- ☛ How do we choose indexes? What tools are available?
- ☛ What is co-clustering and how is it used?
- ☛ What are the choices in tuning a database?
- ☛ How do we tune queries and view?
- ☛ What is the impact of concurrency on performance?
- ☛ How can we reduce lock contention and hotspots?
- ☛ What are popular database benchmarks and how are they used?
- ➡ **Key concepts:** Physical database design, database tuning, workload, co-clustering, index tuning, tuning wizard, index configuration, hot spot, lock contention, database benchmark, transactions per second

Advice to a client who complained about rain leaking through the roof onto the dining table: "Move the table."

.....Architect Frank Lloyd Wright

The performance of a DBMS on commonly asked queries and typical update operations is the ultimate measure of a database design. A DBA can improve performance by identifying performance bottlenecks and adjusting some DBMS parameters (e.g., the size of the buffer pool or the frequency of checkpointing) or adding hardware to eliminate such bottlenecks. The first step in achieving

good performance, however, is to make good database design choices, which is the focus of this chapter.

After we design the *conceptual* and *external* schemas, that is, create a collection of relations and views along with a set of integrity constraints, we must address performance goals through physical database design, in which we design the *physical* schema. As user requirements evolve, it is usually necessary to tune, or adjust, all aspects of a database design for good performance.

This chapter is organized as follows. We give an overview of physical database design and tuning in Section 20.1. The most important physical design decisions concern the choice of indexes. We present guidelines for deciding which indexes to create in Section 20.2. These guidelines are illustrated through several examples and developed further in Sections 20.3. In Section 20.4, we look closely at the important issue of clustering; we discuss how to choose clustered indexes and whether to store tuples from different relations near each other (an option supported by some DBMSs). In Section 20.5, we emphasize how well-chosen indexes can enable some queries to be answered without ever looking at the actual data records. Section 20.6 discusses tools that can help the DBA to automatically select indexes.

In Section 20.7, we survey the main issues of database tuning. In addition to tuning indexes, we may have to tune the conceptual schema as well as frequently used query and view definitions. We discuss how to refine the conceptual schema in Section 20.8 and how to refine queries and view definitions in Section 20.9. We briefly discuss the performance impact of concurrent access in Section 20.10. We illustrate tuning on our Internet shop example in Section 20.11. We conclude the chapter with a short discussion of DBMS benchmarks in Section 20.12; benchmarks help evaluate the performance of alternative DBMS products.

20.1 INTRODUCTION TO PHYSICAL DATABASE DESIGN

Like all other aspects of database design, physical design must be guided by the nature of the data and its intended use. In particular, it is important to understand the typical workload that the database must support; the workload consists of a mix of queries and updates. Users also have certain requirements about how fast certain queries or updates must run or how many transactions must be processed per second. The workload description and users' performance requirements are the basis on which a number of decisions have to be made during physical database design.

Identifying Performance Bottlenecks: All commercial systems provide a suite of tools for monitoring a wide range of system parameters. These tools, used properly, can help identify performance bottlenecks and suggest aspects of the database design and application code that need to be tuned for performance. For example, we can ask the DBMS to monitor the execution of the database for a certain period of time and report on the number of clustered scans, open cursors, lock requests, checkpoints, buffer scans, average wait time for locks, and many such statistics that give detailed insight into a *snapshot* of the live system. In Oracle, a report containing this information can be generated by running a script called UTLBSTAT.SQL to initiate monitoring and a script UTLBSTAT.SQL to terminate monitoring. The system catalog contains details about the sizes of tables, the distribution of values in index keys, and the like. The plan generated by the DBMS for a given query can be viewed in a graphical display that shows the estimated cost for each plan operator. While the details are specific to each vendor, all major DBMS products on the market today provide a suite of such tools.

To create a good physical database design and tune the system for performance in response to evolving user requirements, a designer must understand the workings of a DBMS, especially the indexing and query processing techniques supported by the DBMS. If the database is expected to be accessed concurrently by many users, or is a *distributed database*, the task becomes more complicated and other features of a DBMS come into play. We discuss the impact of concurrency on database design in Section 20.10 and distributed databases in Chapter 22.

20.1.1 Database Workloads

The key to good physical design is arriving at an accurate description of the expected workload. A workload description includes the following:

1. A list of queries (with their frequency, as a ratio of all queries / updates).
2. A list of updates and their frequencies.
3. Performance goals for each type of query and update.

For each query in the workload, we must identify

- Which relations are accessed.
- Which attributes are retained (in the SELECT clause).

- Which attributes have selection or join conditions expressed on them (in the WHERE clause) and how selective these conditions are likely to be.

Similarly, for each update in the workload, we must identify

- Which attributes have selection or join conditions expressed on them (in the WHERE clause) and how selective these conditions are likely to be.
- The type of update (INSERT, DELETE, or UPDATE) and the updated relation.
- For UPDATE commands, the fields that are modified by the update.

Remember that queries and updates typically have parameters, for example, a debit or credit operation involves a particular account number. The values of these parameters determine selectivity of selection and join conditions.

Updates have a query component that is used to find the target tuples. This component can benefit from a good physical design and the presence of indexes. On the other hand, updates typically require additional work to maintain indexes on the attributes that they modify. Thus, while queries can only benefit from the presence of an index, an index may either speed up or slow down a given update. Designers should keep this trade-off in mind when creating indexes.

20.1.2 Physical Design and Tuning Decisions

Important decisions made during physical database design and database tuning include the following:

1. Choice of indexes to create:

- Which relations to index and which field or combination of fields to choose as index search keys.
- For each index, should it be clustered or unclustered?

2. *Tuning the conceptual schema:*

- *Alternative normalized schemas:* We usually have more than one way to decompose a schema into a desired normal form (BCNF or 3NF). A choice can be made on the basis of performance criteria.
- *Denormalization:* We might want to reconsider schema decompositions carried out for normalization, during the conceptual schema design process to improve the performance of queries that involve attributes from several previously decomposed relations.

- *Vertical partitioning:* Under certain circumstances we might want to further decompose relations to improve the performance of queries that involve only a few attributes.
 - *Views:* We might want to add some views to mask the changes in the conceptual schema from users.
3. *Query and transaction tuning:* Frequently executed queries and transactions might be rewritten to run faster.

In parallel or distributed databases, which we discuss in Chapter 22, there are additional choices to consider, such as whether to partition a relation across different sites or whether to store copies of a relation at multiple sites.

20.1.3 Need for Database Tuning

Accurate, detailed workload information may be hard to come by while doing the initial design of the system. Consequently, tuning a database after it has been designed and deployed is important—we must refine the initial design in the light of actual usage patterns to obtain the best possible performance.

The distinction between database design and database tuning is somewhat arbitrary. We could consider the design process to be over once an initial conceptual schema is designed and a set of indexing and clustering decisions is made. Any subsequent changes to the conceptual schema or the indexes, say, would then be regarded as tuning. Alternatively, we could consider some refinement of the conceptual schema (and physical design decisions affected by this refinement) to be part of the physical design process.

Where we draw the line between design and tuning is not very important, and we simply discuss the issues of index selection and database tuning without regard to when the tuning is carried out.

20.2 GUIDELINES FOR INDEX SELECTION

In considering which indexes to create, we begin with the list of queries (including queries that appear as part of update operations). Obviously, only relations accessed by some query need to be considered as candidates for indexing, and the choice of attributes to index is guided by the conditions that appear in the WHERE clauses of the queries in the workload. The presence of suitable indexes can significantly improve the evaluation plan for a query, as we saw in Chapters 8 and 12.

One approach to index selection is to consider the most important queries in turn, and, for each, determine which plan the optimizer would choose given the indexes currently on our list of (to be created) indexes. Then we consider whether we can arrive at a substantially better plan by adding more indexes; if so, these additional indexes are candidates for inclusion in our list of indexes. In general, range retrievals benefit from a B+ tree index, and exact-match retrievals benefit from a hash index. Clustering benefits range queries, and it benefits exact-match queries if several data entries contain the same key value.

Before adding an index to the list, however, we must consider the impact of having this index on the updates in our workload. As we noted earlier, although an index can speed up the query component of an update, all indexes on an updated attribute—in *any* attribute, in the case of inserts and deletes—must be updated whenever the value of the attribute is changed. Therefore, we must sometimes consider the trade-off of slowing some update operations in the workload in order to speed up some queries.

Clearly, choosing a good set of indexes for a given workload requires an understanding of the available indexing techniques, and of the workings of the query optimizer. The following guidelines for index selection summarize our discussion:

Whether to Index (Guideline 1): The obvious points are often the most important. Do not build an index unless some query—including the query components of updates—benefits from it. Whenever possible, choose indexes that speed up more than one query.

Choice of Search Key (Guideline 2): Attributes mentioned in a WHERE clause are candidates for indexing.

- An exact-match selection condition suggests that we consider an index on the selected attributes, ideally, a hash index.
- A range selection condition suggests that we consider a B+ tree (or ISAM) index on the selected attributes. A B+ tree index is usually preferable to an ISAM index. An ISAM index may be worth considering if the relation is infrequently updated, but we assume that a B+ tree index is always chosen over an ISAM index, for simplicity.

Multi-Attribute Search Keys (Guideline 3): Indexes with multiple-attribute search keys should be considered in the following two situations:

- A WHERE clause includes conditions on more than one attribute of a relation.

- They enable index-only evaluation strategies (i.e., accessing the relation can be avoided) for important queries. (This situation could lead to attributes being in the search key even if they do not appear in WHERE clauses.)

When creating indexes on search keys with multiple attributes, if range queries are expected, be careful to order the attributes in the search key to match the queries.

Whether to Cluster (Guideline 4): At most one index on a given relation can be clustered, and clustering affects performance greatly; so the choice of clustered index is important.

- As a rule of thumb, range queries are likely to benefit the most from clustering. If several range queries are posed on a relation, involving different sets of attributes, consider the selectivity of the queries and their relative frequency in the workload when deciding which index should be clustered.
- If an index enables an index-only evaluation strategy for the query it is intended to speed up, the index need not be clustered. (Clustering matters only when the index is used to retrieve tuples from the underlying relation.)

Hash versus Tree Index (Guideline 5): A B+ tree index is usually preferable because it supports range queries as well as equality queries. A hash index is better in the following situations:

- The index is intended to support index nested loops join; the indexed relation is the inner relation, and the search key includes the join columns. In this case, the slight improvement of a hash index over a B+ tree for equality selections is magnified, because an equality selection is generated for each tuple in the outer relation.
- There is a very important equality query, and no range queries, involving the search key attributes.

Balancing the Cost of Index Maintenance (Guideline 6): After drawing up a 'wishlist' of indexes to create, consider the impact of each index on the updates in the workload.

- If maintaining an index slows down frequent update operations, consider dropping the index.
- Keep in mind, however, that adding an index may well speed up a given update operation. For example, an index on employee IDs could speed up the operation of increasing the salary of a given employee (specified by ID).

20.3 BASIC EXAMPLES OF INDEX SELECTION

The following examples illustrate how to choose indexes during database design, continuing the discussion from Chapter 8, where we focused on index selection for single-table queries. The schemas used in the examples are not described in detail; in general, they contain the attributes named in the queries. Additional information is presented when necessary.

Let us begin with a simple query:

```
SELECT E.ename, D.dname
FROM   Employees E, Departments D
WHERE  D.dname='Toy' AND E.dno=D.dno
```

The relations mentioned in the query are Employees and Departments, and both conditions in the WHERE clause involve equalities. Our guidelines suggest that we should build hash indexes on the attributes involved. It seems clear that we should build a hash index on the *dname* attribute of Departments. But consider the equality *E.dno=D.dno*. Should we build an index (hash, of course) on the *dno* attribute of Departments or Employees (or both)? Intuitively, we want to retrieve Departments tuples using the index on *dname* because few tuples are likely to satisfy the equality selection *D.dname='Toy'*.¹ For each qualifying Departments tuple, we then find matching Employees tuples by using an index on the *dno* attribute of Employees. So, we should build an index on the *dno* field of Employees. (Note that nothing is gained by building an additional index on the *dno* field of Departments because Departments tuples are retrieved using the *dname* index.)

Our choice of indexes was guided by the query evaluation plan we wanted to utilize. This consideration of a potential evaluation plan is common while making physical design decisions. Understanding query optimization is very useful for physical design. We show the desired plan for this query in Figure 20.1.

As a variant of this query, suppose that the WHERE clause is modified to be *WHERE D.dname='Toy' AND E.dno=D.dno AND E.age=25*. Let us consider alternative evaluation plans. One good plan is to retrieve Departments tuples that satisfy the selection on *dname* and retrieve matching Employees tuples by using an index on the *dno* field; the selection on *age* is then applied on-the-fly. However, unlike the previous variant of this query, we do not really need to have an index on the *dno* field of Employees if we have an index on *age*. In this

¹This is only a heuristic. If *dname* is not the key, and we have no statistics to verify this claim, it is possible that several tuples satisfy this condition.

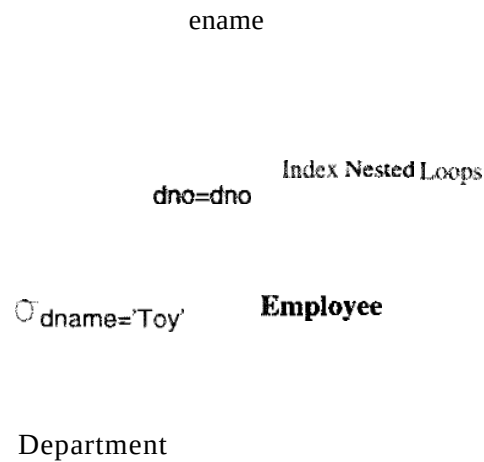


Figure 20.1 A Desirable Query Evaluation Plan

case we can retrieve Departments tuples that satisfy the selection on *dname* (by using the index on *dname*, as before), retrieve Employees tuples that satisfy the selection on *age* by using the index on *age*, and join these sets of tuples. Since the sets of tuples we join are small, they fit in memory and the join method is unimportant. This plan is likely to be somewhat poorer than using an index on *dno*, but it is a reasonable alternative. Therefore, if we have an index on *age* already (prompted by some other query in the workload), this variant of the sample query does not justify creating an index on the *dno* field of Employees.

Our next query involves a range selection:

```

SELECT E.ename, I, dname
FROM   Employees E, Departments D
WHERE  E.sal BETWEEN 10000 AND 20000
       AND E.hobby='Starups' AND E.dno=D.dno
  
```

This query illustrates the use of the BETWEEN operator for expressing range selections. It is equivalent to the condition:

$$10000 \leq E.sal \text{ AND } E.sal \leq 20000$$

The use of BETWEEN to express range conditions is recommended; it makes it easier for both the user and the optimizer to recognize both parts of the range selection.

Returning to the example query, both (nonjoin) selections are on the Employees relation. Therefore, it is clear that a plan in which Employees is the outer relation and Departments is the inner relation is the best, as in the previous query, and we should build a hash index on the *dno* attribute of Departments. But which index should we build on Employees? A B+ tree index on the *sal* attribute would help with the range selection, especially if it is clustered. A

hash index on the *hobby* attribute would help with the equality selection. If one of these indexes is available, we could retrieve Employees tuples using this index, retrieve matching Departments tuples using the index on *dno*, and apply all remaining selections and projections on-the-fly. If both indexes are available, the optimizer would choose the more selective index for the given query; that is, it would consider which selection (the range condition on *salary* or the equality on *hobby*) has fewer qualifying tuples. In general, which index is more selective depends on the data. If there are very few people with salaries in the given range and many people collect stamps, the B-tree index is best. Otherwise, the hash index on *hobby* is best.

If the query constants are known (as in our example), the selectivities can be estimated if statistics on the data are available. Otherwise, as a rule of thumb, an equality selection is likely to be more selective, and a reasonable decision would be to create a hash index on *hobby*. Sometimes, the query constants are not known—we might obtain a query by expanding a query on a view at runtime, or we might have a query in Dynamic SQL, which allows constants to be specified as *wild-card variables* (e.g., %X) and instantiated at runtime (see Sections 6.1.3 and 6.2). In this case, if the query is very important, we might choose to create a B+ tree index on *sal* and a hash index on *hobby* and leave the choice to be made by the optimizer at runtime.

20.4 CLUSTERING AND INDEXING

Clustered indexes can be especially important while accessing the inner relation in an index nested loops join. To understand the relationship between clustered indexes and joins, let us revisit our first example:

```
SELECT E.ename, D.dname
FROM   Employees E, Departments D
WHERE  D.dname='Toy' AND E.dno=D.dno
```

We concluded that a good evaluation plan is to use an index on *dname* to retrieve Departments tuples satisfying the condition on *dname* and to find matching Employees tuples using an index on *dno*. Should these indexes be clustered? Given our assumption that the number of tuples satisfying $D.dname = 'Toy'$ is likely to be small, we should build an unclustered index on *dname*. On the other hand, Employees is the inner relation in an index nested loops join and *dno* is not a candidate key. This situation is a strong argument that the index on the *dno* field of Employees should be clustered. In fact, because the join consists of repeatedly posing equality selections on the *dno* field of the inner relation, this type of query is a stronger justification for making the index on *dno* clustered than a simple selection query such as the previous selection on

hobby. (Of course, factors such as selectivities and frequency of queries have to be taken into account as well.)

The following example, very similar to the previous one, illustrates how clustered indexes can be used for sort-merge joins:

```
SELECT E.ename,D.dname
FROM   Employees E, Departments D
WHERE  E.hobby='Stamps' AND E.dno=D.dno
```

This query differs from the previous query in that the condition *E.hobby*='Stamps' replaces *D.dname*='Toy'. Based on the assumption that there are few employees in the Toy department, we chose indexes that would facilitate an indexed nested loops join with Departments as the outer relation. Now, let us suppose that many employees collect stamps. In this case, a block nested loops or sort-merge join might be more efficient. A sort-merge join can take advantage of a clustered B+ tree index on the *dno* attribute in Departments to retrieve tuples and thereby avoid sorting Departments. Note that an unclustered index is not useful---since all tuples are retrieved, performing one I/O per tuple is likely to be prohibitively expensive. If there is no index on the *dno* field of Employees, we could retrieve Employees tuples (possibly using an index on *hobby*, especially if the index is clustered), apply the selection *E.hobby*='Stamps' on-the-fly, and sort the qualifying tuples on *dno*.

As our discussion has indicated, when we retrieve tuples using an index, the impact of clustering depends on the number of retrieved tuples, that is, the number of tuples that satisfy the selection conditions that match the index. An unclustered index is just as good as a clustered index for a selection that retrieves a single tuple (e.g., an equality selection on a candidate key). As the number of retrieved tuples increases, the unclustered index quickly becomes more expensive than even a sequential scan of the entire relation. Although the sequential scan retrieves all tuples, each page is retrieved exactly once, whereas a page may be retrieved as often as the number of tuples it contains if an unclustered index is used. If blocked I/O is performed (as is common), the relative advantage of sequential scan versus an unclustered index increases further. (Blocked I/O also speeds up access using a clustered index, of course.)

We illustrate the relationship between the number of retrieved tuples, viewed as a percentage of the total number of tuples in the relation, and the cost of various access methods in Figure 20.2. We assume that the query is a selection on a single relation, for simplicity. (Note that this figure reflects the cost of writing out the result; otherwise, the line for sequential scan would be flat.)

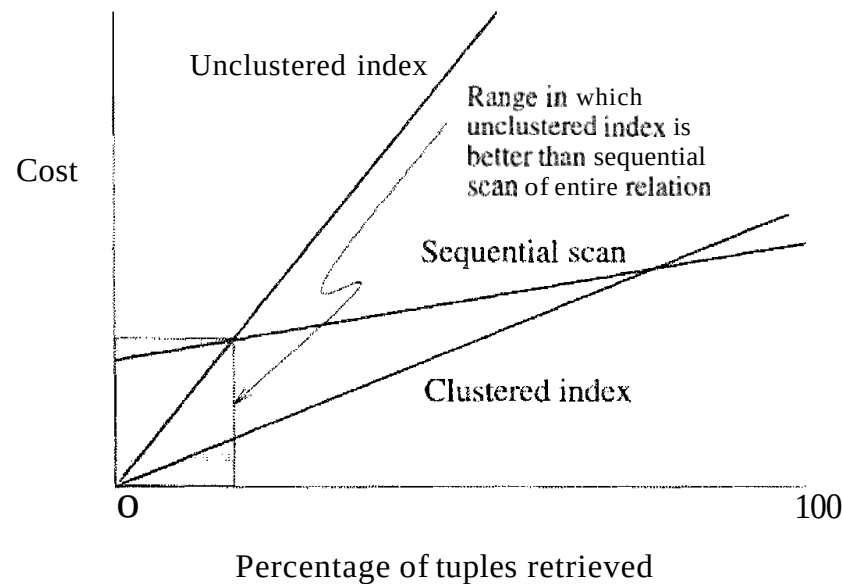


Figure 20.2 The Impact of Clustering

20.4.1 Co-clustering Two Relations

In our description of a typical database system architecture in Chapter 8, we explained how a relation is stored as a file of records. Although a file usually contains only the records of *SOIn8* one relation, *SCHue* systems allow records from more than one relation to be stored in a single file. The database user can request that the records from two relations be interleaved physically in this manner. This data layout is sometimes referred to as co-clustering the two relations. We now discuss when co-clustering can be beneficial.

As an example, consider two relations with the following schemas:

```
Parts(pid: integer, pname: string, cost: integer, supplierid: integer)
Assembly(partid: integer, componentid: integer, quantity: integer)
```

In this schema the *componentid* field of Assembly is intended to be the *pid* of some part that is used as a component in assembling the part with *pid* equal to *partid*. Therefore, the Assembly table represents a 1:N relationship between parts and their subparts; a part can have many subparts, but each part is the subpart of at most one part. In the Parts table, *pid* is the key. For composite parts (those assembled from other parts, as indicated by the contents of Assembly), the *cost* field is taken to be the cost of assembling the part from its subparts.

Suppose that a frequent query is to find the (immediate) subparts of all parts supplied by a given supplier:

```
SELECT P.piel, A.componentid
FROM Parts P, Assembly A
```

```
WHERE P.pid = A.partid AND P.supplierid = 'Acme'
```

A good evaluation plan is to apply the selection condition on Parts and then retrieve matching Assembly tuples through an index on the *partid* field. Ideally, the index on *partid* should be clustered. This plan is reasonably good. However, if such selections are common and we want to optimize them further, we can *co-cluster* the two tables. In this approach, we store records of the two tables together, with each Parts record *P* followed by all the Assembly records *A* such that *P.pid* = *A.partid*. This approach improves on storing the two relations separately and having a clustered index on *partid* because it does not need an index lookup to find the Assembly records that match a given Parts record. Thus, for each selection query, we save a few (typically two or three) index page I/Os.

If we are interested in finding the immediate subparts of *all* parts (i.e., the preceding query with no selection on *supplierid*), creating a clustered index on *partid* and doing an index nested loops join with Assembly as the inner relation offers good performance. An even better strategy is to create a clustered index on the *partid* field of Assembly and the *pid* field of Parts, then do a sort-merge join, using the indexes to retrieve tuples in sorted order. This strategy is comparable to doing the join using a co-clustered organization, which involves just one scan of the set of tuples (of Parts and Assembly, which are stored together in interleaved fashion).

The real benefit of co-clustering is illustrated by the following query:

```
SELECT P.pid,A.componentid
FROM   Parts P, Assembly A
WHERE  P.pid = A.partid AND P.cost=10
```

Suppose that many parts have *cost* = 10. This query essentially amounts to a collection of queries in which we are given a Parts record and want to find matching Assembly records. If we have an index on the *cost* field of Parts, we can retrieve qualifying Parts tuples. For each such tuple, we have to use the index on Assembly to locate records with the given *pid*. The index access for Assembly is avoided if we have a co-clustered organization. (Of course, we still require an index on the *cost* attribute of Parts tuples.)

Such an optimization is especially important if we want to traverse several levels of the part-subpart hierarchy. For example, a common query is to find the total cost of a part, which requires us to repeatedly carry out joins of Parts and Assembly. Incidentally, if we do not know the number of levels in the hierarchy in advance, the number of joins varies and the query cannot be expressed in SQL. The query can be answered by embedding an SQL statement

for the join inside an iterative host language program. How to express the query is orthogonal to our main point here, which is that co-clustering is especially beneficial when the join in question is carried out very frequently (either because it arises repeatedly in an important query such as finding total cost, or because the join query itself is asked frequently).

To summarize co-clustering:

- It can speed up joins, in particular key-foreign key joins corresponding to 1:N relationships.
- A sequential scan of either relation becomes slower. (In our example, since several Assembly tuples are stored in between consecutive Parts tuples, a scan of all Parts tuples becomes slower than if Parts tuples were stored separately. Similarly, a sequential scan of all Assembly tuples is also slower.)
- All inserts, deletes, and updates that alter record lengths become slower, thanks to the overheads involved in maintaining the clustering. (We do not discuss the implementation issues involved in co-clustering.)

20.5 INDEXES THAT ENABLE INDEX-ONLY PLANS

This section considers a number of queries for which we can find efficient plans that avoid retrieving tuples from one of the referenced relations; instead, these plans scan an associated index (which is likely to be much smaller). An index that is used (only) for index-only scans does *not* have to be clustered because tuples from the indexed relation are not retrieved.

This query retrieves the managers of departments with at least one employee:

```
SELECT D.mgr
FROM   Departments D, Employees E
WHERE  D.dno=E.dno
```

Observe that no attributes of Employees are retained. If we have an index on the *dno* field of Employees, the optimization of doing an index nested loops join using an index-only search for the inner relation is applicable. Given this variant of the query, the correct decision is to build an unclustered index on the *dno* field of Employees, rather than a clustered index.

The next query takes this idea a step further:

```
SELECT D.mgr, E.ename
FROM   Departments D, Employees E
WHERE  D.dno=E.dno
```

If we have an index on the *dno* field of *Employees*, we can use it to retrieve *Employees* tuples during the join (with *Departments* as the outer relation), but unless the index is clustered, this approach is not be efficient. On the other hand, suppose that we have a B+- tree index on $(dno, e'id)$. Now all the information we need about an *Employees* tuple is contained in the data entry for this tuple in the index. We can use the index to find the first data entry with a given *dno*; all data entries with the same *dno* are stored together in the index. (Note that a hash index on the composite key (dno, eid) cannot be used to locate an entry with just a given *dno*!) We can therefore evaluate this query using an index nested loops join with *Departments* as the outer relation and an index-only scan of the inner relation.

20.6 TOOLS TO ASSIST IN INDEX SELECTION

The number of possible indexes to consider building is potentially very large: For each relation, we can potentially consider all possible subsets of attributes as an index key; we have to decide on the ordering of the attributes in the index; and we also have to decide which indexes should be clustered and which unclustered. Many large applications—for example enterprise resource planning systems—create tens of thousands of different relations, and manual tuning of such a large schema is a daunting endeavor.

The difficulty and importance of the index selection task motivated the development of tools that help database administrators select appropriate indexes for a given workload. The first generation of such index tuning wizards, or index advisors, were separate tools outside the database engine; they suggested indexes to build, given a workload of SQL queries. The main drawback of these systems was that they had to replicate the database query optimizer's cost model in the tuning tool to make sure that the optimizer would choose the same query evaluation plans as the design tool. Since query optimizers change from release to release of a commercial database system, considerable effort was needed to keep the tuning tool and the database optimizer synchronized. The most recent generation of tuning tools are integrated with the database engine and use the database query optimizer to estimate the cost of a workload given a set of indexes, avoiding duplication of the query optimizer's cost model into an external tool.

20.6.1 Automatic Index Selection

We call a set of indexes for a given database schema an index configuration. We assume that a query workload is a set of queries over a database schema where each query has a frequency of occurrence assigned to it. Given a database schema and a workload, the cost of an index configuration is the expected

cost of running the queries in the workload given the index configuration taking the different frequencies of queries in the workload into account. Given a database schema and a query workload, we can now define the problem of automatic index selection as finding an index configuration with minimal cost. As in query optimization, in practice our goal is to find a *good* index configuration rather than the true optimal configuration.

Why is automatic index selection a hard problem? Let us calculate the number of different indexes with c attributes, assuming that the table has n attributes. For the first attribute in the index, there are n choices, for the second attribute $n - 1$, and thus for a c attribute index, there are overall $n \cdot (n - 1) \dots (n - c + 1) = \frac{n!}{(n - c)!}$ different indexes possible. The total number of different indexes with up to c attributes is

$$\sum_{i=1}^c \frac{n!}{(n - i)!}$$

For a table with 10 attributes, there are 10 different one-attribute indexes, 90 different two-attribute indexes, and 30240 different five-attribute indexes. For a complex workload involving hundreds of tables, the number of possible index configurations is clearly very large.

The efficiency of automatic index selection tools can be separated into two components: (1) the number of candidate index configurations considered, and (2) the number of optimizer calls necessary to evaluate the cost for a configuration. Note that reducing the search space of candidate indexes is analogous to restricting the search space of the query optimizer to left-deep plans. In many cases, the optimal plan is not left-deep, but among all left-deep plans there is usually a plan whose cost is close to the optimal plan.

We can easily reduce the time taken for automatic index selection by reducing the number of candidate index configurations, but the smaller the space of index configurations considered, the farther away the final index configuration is from the optimal index configuration. Therefore, different index tuning wizards prune the search space differently, for example, by considering only one- or two-attribute indexes.

20.6.2 How Do Index Tuning Wizards Work?

All index tuning wizards search a set of candidate indexes for an index configuration with lowest cost. Tools differ in the space of candidate index configurations they consider and how they search this space. We describe one representative algorithm; existing tools implement variants of this algorithm, but their implementations have the same basic structure.

The DB2 Index Advisor. The DB2 Index Advisor is a tool for automatic index recommendation given a workload. The workload is stored in the database system in a table called `ADVISE_WORKLOAD`. It is populated either (1) by SQL statements from the DB2 dynamic SQL statement cache, a cache for recently executed SQL statements, (2) with SQL statements from packages or groups of statically compiled SQL statements, or (3) with SQL statements from an online monitor called the Query Patroller. The DB2 Advisor allows the user to specify the maximum amount of disk space for new indexes and a maximum time for the computation of the recommended index configuration.

The DB2 Index Advisor consists of a program that intelligently searches a subset of index configurations. Given a candidate configuration, it calls the query optimizer for each query in the `ADVISE_WORKLOAD` table first in the `RECOMMEND_INDEXES` mode, where the optimizer recommends a set of indexes and stores them in the `ADVISE_INDEXES` table. In the `EVALUATE_INDEXES` mode, the optimizer evaluates the benefit of the index configuration for each query in the `ADVISE_WORKLOAD` table. The output of the index tuning step is a set of SQL DDL statements whose execution creates the recommended indexes.

The Microsoft SQL Server 2000 Index Tuning Wizard. Microsoft pioneered the implementation of a tuning wizard integrated with the database query optimizer. The Microsoft Index Tuning Wizard has three tuning modes that permit the user to trade off running time of the analysis and number of candidate index configurations examined: *fast*, *medium*, and *thorough*, with *fast* having the lowest running time and *thorough* examining the largest number of configurations. To further reduce the running time, the tool has a sampling mode in which the tuning wizard randomly samples queries from the input workload to speed up analysis. Other parameters include the maximum space allowed for the recommended indexes, the maximum number of attributes per index considered, and the tables on which indexes can be generated. The Microsoft Index Tuning Wizard also permits *table scaling*, where the user can specify an anticipated number of records for the tables involved in the workload. This allows users to plan for future growth of the tables.

Before we describe the index tuning algorithm, let us consider the problem of estimating the cost of a configuration. Note that it is not feasible to actually create the set of indexes in a candidate configuration and then optimize the query workload given the physical index configuration. Creation of even a single candidate configuration with several indexes might take hours for large databases and put considerable load on the database system itself. Since we want to examine a large number of possible candidate configurations, this approach is not feasible.

Therefore index tuning algorithms usually *simulate* the effect of indexes in a candidate configuration (unless such indexes already exist). Such what-if indexes look to the query optimizer like any other index and are taken into account when calculating the cost of the workload for a given configuration, but the creation of what-if indexes does not incur the overhead of actual index creation. Commercial database systems that support index tuning wizards using the database query optimizer have been extended with a module that permits the creation and deletion of what-if indexes with the necessary statistics about the indexes (that are used when estimating the cost of a query plan).

We now describe a representative index tuning algorithm. The algorithm proceeds in two steps, *candidate index selection* and *configuration enumeration*. In the first step, we select a set of candidate indexes to consider during the second step as building blocks for index configurations. Let us discuss these two steps in more detail.

Candidate Index Selection

We saw in the previous section that it is impossible to consider every possible index, due to the huge number of candidate indexes available for larger database schemas. One heuristic to prune the large space of possible indexes is to first tune each query in the workload independently and then select the union of the indexes selected in this first step as input to the second step.

For a query, let us introduce the notion of an indexable attribute, which is an attribute whose appearance in an index could change the cost of the query. An indexable attribute is an attribute on which the WHERE-part of the query has a condition (e.g., an equality predicate) or the attribute appears in a GROUP BY or ORDER BY clause of the SQL query. An admissible index for a query is an index that contains only indexable attributes in the query.

How do we select candidate indexes for an individual query? One approach is a basic enumeration of all indexes with up to k attributes. We start with all indexable attributes as single attribute candidate indexes, then add all com-

binations of two indexable attributes as candidate indexes, and repeat this procedure until a user-defined size threshold k . This procedure is obviously very expensive as we add overall $n + n \cdot (n - 1) + \dots + n \cdot (n - 1) \dots (n - k + 1)$ candidate indexes, but it guarantees that the best index with up to k attributes is among the candidate indexes. The references at the end of this chapter contain pointers to faster (but less exhaustive) heuristic search algorithms.

Enumerating Index Configurations

In the second phase, we use the candidate indexes to enumerate index configurations. As in the first phase, we can exhaustively enumerate all index configurations up to size k , this time combining candidate indexes. As in the previous phase, more sophisticated search strategies are possible that cut down the number of configurations considered while still generating a final configuration of high quality (i.e., low execution cost for the final workload).

20.7 OVERVIEW OF DATABASE TUNING

After the initial phase of database design, actual use of the database provides a valuable source of detailed information that can be used to refine the initial design. Many of the original assumptions about the expected workload can be replaced by observed usage patterns; in general, some of the initial workload specification is validated, and some of it turns out to be wrong. Initial guesses about the size of data can be replaced with actual statistics from the system catalogs (although this information keeps changing as the system evolves). Careful monitoring of queries can reveal unexpected problems; for example, the optimizer may not be using some indexes as intended to produce good plans.

Continued database tuning is important to get the best possible performance. In this section, we introduce three kinds of tuning: *tuning indexes*, *tuning the conceptual schema*, and *tuning queries*. Our discussion of index selection also applies to index tuning decisions. Conceptual schema and query tuning are discussed further in Sections 20.8 and 20.9.

20.7.1 Tuning Indexes

The initial choice of indexes may be refined for one of several reasons. The simplest reason is that the observed workload reveals that some queries and updates considered important in the initial workload specification are not very frequent. The observed workload may also identify some new queries and updates that *are* important. The initial choice of indexes has to be reviewed in light of this new information. Some of the original indexes may be dropped and

new ones added. The reasoning involved is similar to that used in the initial design.

It may also be discovered that the optimizer in a given system is not finding some of the plans that it was expected to. For example, consider the following query, which we discussed earlier:

```
SELECT D.Dname
FROM   Employees E, Departments D
WHERE  D.Dname='Toy' AND E.Dno=D.Dno
```

A good plan here would be to use an index on *Dname* to retrieve Departments tuples with *Dname* = 'Toy' and to use an index on the *Dno* field of Employees as the inner relation, using an index-only scan. Anticipating that the optimizer would find such a plan, we might have created an unclustered index on the *Dno* field of Employees.

Now suppose queries of this form take an unexpectedly long time to execute. We can ask to see the plan produced by the optimizer. (Most commercial systems provide a simple command to do this.) If the plan indicates that an index-only scan is not being used, but that Employees tuples are being retrieved, we have to rethink our initial choice of index, given this revelation about our system's (unfortunate) limitations. An alternative to consider here would be to drop the unclustered index on the *Dno* field of Employees and replace it with a clustered index.

Some other common limitations of optimizers are that they do not handle selections involving string expressions, arithmetic, or *null* values effectively. We discuss these points further when we consider query tuning in Section 20.9.

In addition to re-examining our choice of indexes, it pays to periodically reorganize some indexes. For example, a static index, such as an ISAM index, may have developed long overflow chains. Dropping the index and rebuilding it—if feasible, given the interrupted access to the indexed relation—can substantially improve access times through this index. Even for a dynamic structure such as a B+ tree, if the implementation does not merge pages on deletes, space occupancy can decrease considerably in some situations. This in turn makes the size of the index (in pages) larger than necessary, and could increase the height and therefore the access time. Rebuilding the index should be considered. Extensive updates to a clustered index might also lead to overflow pages being allocated, thereby decreasing the degree of clustering. Again, rebuilding the index may be worthwhile.

Finally, note that the query optimizer relies on statistics maintained in the SystCll catalogs. These statistics are updated only when a special utility program is run; be sure to run the utility frequently enough to keep the statistics reasonably current.

20.7.2 Tuning the Conceptual Schema

In the course of database design, we may realize that our current choice of relation Schema does not enable us meet our performance objectives for the given workload with any (feasible) set of physical design choices. If so, we may have to redesign our conceptual schema (and re-examine physical design decisions affected by the changes we make).

We may realize that a redesign is necessary during the initial design process or later, after the system has been in use for a while. Once a database has been designed and populated with tuples, changing the conceptual schema requires a significant effort in terms of mapping the contents of the relations affected. Nonetheless, it may be necessary to revise the conceptual schema in light of experience with the system. (Such changes to the schema of an operational system are sometimes referred to as schema evolution.) We now consider the issues involved in conceptual schema (re)design from the point of view of performance.

The main point to understand is that *our choice of conceptual schema should be guided by a consideration of the queries and updates in our workload* in addition to the issues of redundancy that motivate normalization (which we discussed in Chapter 19). Several options must be considered while tuning the conceptual schema:

- We may decide to settle for a 3NF design instead of a BCNF design.
- If there are two ways to decompose a given schema into 3NF or BCNF, our choice should be guided by the workload.
- Sometimes we might decide to further decompose a relation that is *already* in BCNF.
- In other situations, we might *denormalize*. That is, we might choose to replace a collection of relations obtained by a decomposition from a larger relation with the original (larger) relation, even though it suffers from some redundancy problems. Alternatively, we might choose to add some fields to certain relations to speed up some important queries, even if this leads to a redundant storage of some information (and, consequently, a schema that is in neither 3NF nor BCNF).

- This discussion of nonnalization has concentrated on the technique of *decomposition*, which amounts to vertical partitioning of a relation. Another technique to consider is *horizontal partitioning* of a relation, which would lead to having two relations with identical schemas. Note that we are not talking about physically partitioning the tuples of a single relation; rather, we want to create two distinct relations (possibly with different constraints and indexes on each).

Incidentally, when we redesign the conceptual schema, especially if we are tuning an existing database schema, it is worth considering whether we should create views to mask these changes from users for whom the original schema is more natural. We discuss the choices involved in tuning the conceptual schema in Section 20.8.

20.7.3 Tuning Queries and Views

If we notice that a query is running much slower than we expected, we have to examine the query carefully to find the problem. Some rewriting of the query, perhaps in conjunction with some index tuning, can often fix the problem. Similar tuning may be called for if queries on some view run slower than expected. We do not discuss view tuning separately; just think of queries on views as queries in their own right (after all, queries on views are expanded to account for the view definition before being optimized) and consider how to tune them.

When tuning a query, the first thing to verify is that the system uses the plan you expect it to use. Perhaps the system is not finding the best plan for a variety of reasons. Some common situations not handled efficiently by many optimizers follow:

- A selection condition involving *null* values.
- Selection conditions involving arithmetic or string expressions or conditions using the OR connective. For example, if we have a condition $E.age = 2 * J.age$ in the WHERE clause, the optimizer may correctly utilize an available index on $E.age$ but fail to utilize an available index on $J.age$. Replacing the condition by $E.age / 2 = J.age$ would reverse the situation.
- Inability to recognize a sophisticated plan such as an index-only scan for an aggregation query involving a GROUP BY clause. Of course, virtually no optimizer looks for plans outside the plan space described in Chapters 12 and 15, such as nonleft-deep join trees. So a good understanding of what an optimizer typically does is important. In addition, the more aware you are of a given system's strengths and limitations, the better off you are.

If the optimizer is not smart enough to find the best plan (using access methods and evaluation strategies supported by the DBMS), some systems allow users to guide the choice of a plan by providing hints to the optimizer; for example, users might be able to force the use of a particular index or choose the join order and join method. A user who wishes to guide optimization in this manner should have a thorough understanding of both optimization and the capabilities of the given DBMS. We discuss query tuning further in Section 20.9.

20.8 CHOICES IN TUNING THE CONCEPTUAL SCHEMA

We now illustrate the choices involved in tuning the conceptual schema through several examples using the following schemas:

```
Contracts(cid: integer, supplierid: integer, projectid: integer,
          deptid: integer, partid: integer, qty: integer, value: real)
Departments(did: integer, budget: real, annualreport: varchar)
Parts(pid: integer, cost: integer)
Projects(jid: integer, mng: char(20))
Suppliers(sid: integer, address: char(50))
```

For brevity, we often use the common convention of denoting attributes by a single character and denoting relation schemas by a sequence of characters. Consider the schema for the relation Contracts, which we denote as CSJDPQV, with each letter denoting an attribute. The meaning of a tuple in this relation is that the contract with *cid* C is an agreement that supplier S (with *sid* equal to *supplierid*) will supply Q items of part P (with *pid* equal to *partid*) to project J (with *jid* equal to *projectid*) associated with department D (with *deptid* equal to *did*), and that the value V of this contract is equal to *value*.²

There are two known integrity constraints with respect to Contracts. A project purchases a given part using a single contract; thus, there cannot be two distinct contracts in which the same project buys the same part. This constraint is represented using the FD $JP \rightarrow C$. Also, a department purchases at most one part from any given supplier. This constraint is represented using the FD $SD \rightarrow P$. In addition, of course, the contract ID C is a key. The meaning of the other relations should be obvious, and we do not describe them further because we focus on the Contracts relation.

²If this schema seems complicated, note that real-life situations often call for considerably more complex schemas!

20.8.1 Settling for a Weaker Normal Form

Consider the Contracts relation. Should we decompose it into smaller relations? Let us see what normal form it is in. The candidate keys for this relation are C and JP. (C is given to be a key, and JP functionally determines C.) The only nonkey dependency is $SD \rightarrow P$, and P is a *prime* attribute because it is part of candidate key JP. Thus, the relation is not in BCNF—because there is a nonkey dependency—but it is in 3NF.

By using the dependency $SD \rightarrow P$ to guide the decomposition, we get the two schemas SDP and CSJDQV. This decomposition is lossless, but it is not dependency-preserving. However, by adding the relation scheme CJP, we obtain a lossless-join, dependency-preserving decomposition into BCNF. Using the guideline that such a decomposition into BCNF is good, we might decide to replace Contracts by three relations with schemas CJP, SDP, and CSJDQV.

However, suppose that the following query is very frequently asked: Find the number of copies Q of part P ordered in contract C. This query requires a join of the decomposed relations CJP and CSJDQV (or SDP and CSJDQV), whereas it can be answered directly using the relation Contracts. The added cost for this query could persuade us to settle for a 3NF design and not decompose Contracts further.

20.8.2 Denormalization

The reasons motivating us to settle for a weaker normal form may lead us to take an even more extreme step: deliberately introduce some redundancy. As an example, consider the Contracts relation, which is in 3NF. Now, suppose that a frequent query is to check that the value of a contract is less than the budget of the contracting department. We might decide to add a budget field B to Contracts. Since *did* is a key for Departments, we now have the dependency $D \rightarrow B$ in Contracts, which means Contracts is not in 3NF any more. Nonetheless, we might choose to stay with this design if the motivating query is sufficiently important. Such a decision is clearly subjective and comes at the cost of significant redundancy.

20.8.3 Choice of Decomposition

Consider the Contracts relation again. Several choices are possible for dealing with the redundancy in this relation:

- We can leave Contracts as it is and accept the redundancy associated with its being in 3NF rather than BCNF.

- We might decide that we want to avoid the anomalies resulting from this redundancy by decomposing Contracts into BCNF using one of the following methods:
 - We have a lossless-join decomposition into PartInfo with attributes SDP and ContractInfo with attributes CSJDQV. As noted previously, this decomposition is not dependency-preserving, and to make it so would require us to add a third relation CJP, whose sole purpose is to allow us to check the dependency $JP \rightarrow C$.
 - We could choose to replace Contracts by just PartInfo and ContractInfo even though this decomposition is not dependency-preserving.

Replacing Contracts by just PartInfo and ContractInfo does not prevent us from enforcing the constraint $JP \rightarrow C$; it only makes this more expensive. We could create an assertion in SQL-92 to check this constraint:

```
CREATE ASSERTION checkDep
CHECK      ( NOT EXISTS
            (SELECT *
            FROM      PartInfo PI, ContractInfo CI
            WHERE      PI.supplierid=CI.supplierid
                      AND PI.deptid=CI.deptid
            GROUP BY  CI.projectid, PI.partid
            HAVING     COUNT (cid) > 1 ) )
```

This assertion is expensive to evaluate because it involves a join followed by a sort (to do the grouping). In comparison, the system can check that JP is a primary key for table CJP by maintaining an index on JP . This difference in integrity-checking cost is the motivation for dependency-preservation. On the other hand, if updates are infrequent, this increased cost may be acceptable; therefore, we might choose not to maintain the table CJP (and quite likely, an index on it).

As another example illustrating decomposition choices, consider the Contracts relation again, and suppose that we also have the integrity constraint that a department uses a given supplier for at most one of its projects: $SPQ \rightarrow V$. Proceeding as before, we have a lossless-join decomposition of Contracts into SDP and CSJDQV. Alternatively, we could begin by using the dependency $SPQ \rightarrow V$ to guide our decomposition, and replace Contracts with SPQV and CSJDPQ. We can then decompose CSJDPQ, guided by $SD \rightarrow P$, to obtain SDP and CSJDQ.

We now have two alternative lossless-join decompositions of Contracts into BCNF, neither of which is dependency-preserving. The first alternative is to

replace **Contracts** with the relations **SDP** and **CSJDQV**. The second alternative is to replace it with **SPQV**, **SDP**, and **CSJDQ**. The addition of **CJP** makes the second decomposition (but not the first) dependency-preserving. Again, the cost of maintaining the three relations **CJP**, **SPQV**, and **CSJDQ** (versus just **CSJDQV**) may lead us to choose the first alternative. In this case, enforcing the given FDs becomes more expensive. We might consider not enforcing them, but we then risk a violation of the integrity of our data.

20.8.4 Vertical Partitioning of BCNF Relations

Suppose that we have decided to decompose **Contracts** into **SDP** and **CSJDQV**. These schemas are in BCNF, and there is no reason to decompose them further from a normalization standpoint. However, suppose that the following queries are very frequent:

- Find the contracts held by supplier **S**.
- Find the contracts placed by department **D**.

These queries might lead us to decompose **CSJDQV** into **CS**, **CD**, and **CJQV**. The decomposition is lossless, of course, and the two important queries can be answered by examining much smaller relations. Another reason to consider such a decomposition is concurrency control *hot spots*. If these queries are common, and the most common updates involve changing the quantity of products (and the value) involved in contracts, the decomposition improves performance by reducing lock contention. Exclusive locks are now set mostly on the **CJQV** table, and reads on **CS** and **CD** do not conflict with these locks.

Whenever we decompose a relation, we have to consider which queries the decomposition might adversely affect, especially if the only motivation for the decomposition is improved performance. For example, if another important query is to find the total value of contracts held by a supplier, it would involve a join of the decomposed relations **CS** and **CJQV**. In this situation, we might decide against the decomposition.

20.8.5 Horizontal Decomposition

Thus far, we have essentially considered how to replace a relation with a collection of vertical decompositions. Sometimes, it is worth considering whether to replace a relation with two relations that have the same attributes as the original relation, each containing a subset of the tuples in the original. Intuitively, this technique is useful when different subsets of tuples are queried in very distinct ways.

For example, different rules may govern large contracts, which are defined as contracts with values greater than 10,000. (Perhaps, such contracts have to be awarded through a bidding process.) This constraint could lead to a number of queries in which Contracts tuples are selected using a condition of the form $value > 10,000$. One way to approach this situation is to build a clustered B+ tree index on the *value* field of Contracts. Alternatively, we could replace Contracts with two relations called LargeContracts and SmallContracts, with the obvious meaning. If this query is the only motivation for the index, horizontal decomposition offers all the benefits of the index without the overhead of index maintenance. This alternative is especially attractive if other important queries on Contracts also require clustered indexes (on fields other than *value*).

If we replace Contracts by two relations LargeContracts and SmallContracts, we could mask this change by defining a view called Contracts:

```
CREATE VIEW Contracts(cid, supplierid, projectid, deptid, partid, qty, value)
AS ((SELECT *
     FROM   LargeContracts)
   UNION
   (SELECT *
    FROM   SmallContracts))
```

However, any query that deals solely with LargeContracts should be expressed directly on LargeContracts and not on the view. Expressing the query on the view Contracts with the selection condition $value > 10,000$ is equivalent to expressing the query on LargeContracts but less efficient. This point is quite general: Although we can mask changes to the conceptual schema by adding view definitions, users concerned about performance have to be aware of the change.

As another example, if Contracts had an additional field *year* and queries typically dealt with the contracts in some one year, we might choose to partition Contracts by year. Of course, queries that involved contracts from more than one year might require us to pose queries against each of the decomposed relations.

20.9 CHOICES IN TUNING QUERIES AND VIEWS

The first step in tuning a query is to understand the plan used by the DBMS to evaluate the query. Systems usually provide some facility for identifying the plan used to evaluate a query. Once we understand the plan selected by the system, we can consider how to improve performance. We can consider a different choice of indexes or perhaps co-clustering two relations for join queries,

guided by our understanding of the old plan and a better plan that we want the DBLVIS to use. The details are similar to the initial design process.

One point worth making is that before creating new indexes we should consider whether rewriting the query achieves acceptable results with existing indexes. For example, consider the following query with an OR connective:

```
SELECT E.dno
FROM   Employees E
WHERE  E.hobby='Stamps' OR E.age==10
```

If we have indexes on both *hobby* and *age*, we can use these indexes to retrieve the necessary tuples, but an optimizer might fail to recognize this opportunity. The optimizer might view the conditions in the WHERE clause as a whole as not matching either index, do a sequential scan of Employees, and apply the selections on-the-fly. Suppose we rewrite the query as the union of two queries, one with the clause *WHERE E.hobby = 'Stamps'* and the other with the clause *WHERE E.age==10*. Now each query is answered efficiently with the aid of the indexes on *hobby* and *age*.

We should also consider rewriting the query to avoid some expensive operations. For example, including DISTINCT in the SELECT clause leads to duplicate elimination, which can be costly. Thus, we should omit DISTINCT whenever possible. For example, for a query on a single relation, we can omit DISTINCT whenever either of the following conditions holds:

- We do not care about the presence of duplicates.
- The attributes mentioned in the SELECT clause include a candidate key for the relation.

Sometimes a query with GROUP BY and HAVING can be replaced by a query without these clauses, thereby eliminating a sort operation. For example, consider:

```
SELECT   MIN (E.age)
FROM     Employees E
GROUP BY E.dno
HAVING   E.dno=102
```

This query is equivalent to

```
SELECT   MIN (E.age)
FROM     Employees E
WHERE    E.dno=102
```

Complex queries are often written in steps, using a temporary relation. We can usually rewrite such queries without the temporary relation to make them run faster. Consider the following query for computing the average salary of departments managed by Robinson:

```
SELECT  *
INTO    Ternp
FROM    EmployeesE, Departments D
WHERE   E.dno=D.dno AND D.mgrname='Robinson'
```

```
SELECT  T.dno, AVG (T.sal)
FROM    Ternp T
GROUP BY T.dno
```

This query can be rewritten as

```
SELECT  E.dno, AVG (E.sal)
FROM    Employees E, Departments D
WHERE   E.dno=D.dno AND D.mgrname='Robinson'
GROUP BY E.dno
```

The rewritten query does not materialize the intermediate relation, Ternp and is therefore likely to be faster. In fact, the optimizer may even find a very efficient index-only plan that never retrieves Employees tuples if there is a composite B+ tree index on (dno, sal). This example illustrates a general observation: *By rewriting queries to avoid unnecessary temporaries, we not only avoid creating the temporary relations, we also open up more optimization possibilities for the optimizer to explore.*

In some situations, however, if the optimizer is unable to find a good plan for a complex query (typically a nested query with correlation), it may be worthwhile to rewrite the query using temporary relations to guide the optimizer toward a good plan.

In fact, nested queries are a common source of inefficiency because many optimizers deal poorly with them, as discussed in Section 15.5. Whenever possible, it is better to rewrite a nested query without nesting and a correlated query without correlation. As already noted, a good reformulation of the query may require us to introduce new, temporary relations, and techniques to do so systematically (ideally, to be done by the optimizer) have been widely studied. (Often though, it is possible to rewrite nested queries without nesting or the use of temporary relations, as illustrated in Section 15.5.

20.10 IMPACT OF CONCURRENCY

In a system with many concurrent users, several additional points must be considered. Transactions obtain *locks* on the pages they access, and other transactions may be blocked waiting for locks on objects they wish to access.

We observed in Section 16.5 that blocking delays must be minimized for good performance and identified two specific ways to reduce blocking:

- Reducing the time that transactions hold locks.
- Reducing hot spots.

We now discuss techniques for achieving these goals.

20.10.1 Reducing Lock Durations

Delay Lock Requests: Tune transactions by writing to local program variables and deferring changes to the database until the end of the transaction. This delays the acquisition of the corresponding locks and reduces the time the locks are held.

Make Transactions Faster: The sooner a transaction completes, the sooner its locks are released. We have already discussed several ways to speed up queries and updates (e.g., using indexes, rewriting queries). In addition, a careful partitioning of the tuples in a relation and its associated indexes across a collection of disks can significantly improve concurrent access. For example, if we have the relation on one disk and an index on another, accesses to the index can proceed without interfering with accesses to the relation, at least at the level of disk reads.

Replace Long Transactions by Short Ones: Sometimes, just too much work is done within a transaction, and it takes a long time and holds locks a long time. Consider rewriting the transaction as two or more smaller transactions; holdable cursors (see Section 6.1.2) can be helpful in doing this. The advantage is that each new transaction completes quicker and releases locks sooner. The disadvantage is that the original list of operations is no longer executed atomically, and the application code must deal with situations in which one or more of the new transactions fail.

Build a Warehouse: Complex queries can hold shared locks for a long time. Often, however, these queries involve statistical analysis of business trends and it is acceptable to run them on a copy of the data that is a little out of date. This led to the popularity of *data warehouses*, which are databases that complement

the operational database by maintaining a copy of data used in complex queries (Chapter 25). Running these queries against the warehouse relieves the burden of long-running queries from the operational database.

Consider a Lower Isolation Level: In many situations, such as queries generating aggregate information or statistical summaries, we can use a lower SQL isolation level such as REPEATABLE READ or READ COMMITTED (Section 16.6). Lower isolation levels incur lower locking overheads, and the application programmer must make good design trade-offs.

20.10.2 Reducing Hot Spots

Delay Operations on Hot Spots: We already discussed the value of delaying lock requests. Obviously, this is especially important for requests involving frequently used objects.

Optimize Access Patterns: The pattern of updates to a relation can also be significant. For example, if tuples are inserted into the Employees relation in *eid* order and we have a B+ tree index on *eid*, each insert goes to the last leaf page of the B+ tree. This leads to hot spots along the path from the root to the rightmost leaf page. Such considerations may lead us to choose a hash index over a B+ tree index or to index on a different field. Note that this pattern of access leads to poor performance for ISAM indexes as well, since the last leaf page becomes a hot spot. This is not a problem for hash indexes because the hashing process randomizes the bucket into which a record is inserted.

Partition Operations on Hot Spots: Consider a data entry transaction that appends new records to a file (e.g., inserts into a table stored as a heap file). Instead of appending records one-per-transaction and obtaining a lock on the last page for each record, we can replace the transaction by several other transactions, each of which writes records to a local file and periodically appends a batch of records to the main file. While we do more work overall, this reduces the lock contention on the last page of the original file.

As a further illustration of partitioning, suppose we track the number of records inserted in a counter. Instead of updating this counter once per record, the preceding approach results in updating several counters and periodically updating the main counter. This idea can be adapted to many uses of counters, with varying degrees of effort. For example, consider a counter that tracks the number of reservations, with the rule that a new reservation is allowed only if the counter is below a maximum value. We can replace this by three counters, each with one-third the original maximum threshold, and three transactions that use these counters rather than the original. We obtain greater concurrency, but

have to deal with the case where one of the counters is at the maximum value but some other counter can still be incremented. Thus, the price of greater concurrency is increased complexity in the logic of the application code.

Choice of Index: If a relation is updated frequently, B+ tree indexes can become a concurrency control bottleneck, because all accesses through the index must go through the root. Thus, the root and index pages just below it can become hot spots. If the DBMS uses specialized locking protocols for tree indexes, and in particular, sets fine-granularity locks, this problem is greatly alleviated. Many current systems use such techniques.

Nonetheless, this consideration may lead us to choose an ISAM index in some situations. Because the index levels of an ISAM index are static, we need not obtain locks on these pages; only the leaf pages need to be locked. An ISAM index may be preferable to a B+ tree index, for example, if frequent updates occur but we expect the relative distribution of records and the number (and size) of records with a given range of search key values to stay approximately the same. In this case the ISAM index offers a lower locking overhead (and reduced contention for locks), and the distribution of records is such that few overflow pages are created.

Hashed indexes do not create such a concurrency bottleneck, unless the data distribution is very skewed and many data items are concentrated in a few buckets. In this case, the directory entries for these buckets can become a hot spot.

20.11 CASE STUDY: THE INTERNET SHOP

Revisiting our running case study, DBDudes considers the expected workload for the B&N Bookstore. The owner of the bookstore expects most of his customers to search for books by ISBN number before placing an order. Placing an order involves inserting one record into the Orders table and inserting one or more records into the Orderlists relation. If a sufficient number of books is available, a shipment is prepared and a value for the *ship_date* in the Orderlists relation is set. In addition, the available quantities of books in stock changes all the time, since orders are placed that decrease the quantity available and new books arrive from suppliers and increase the quantity available.

The DBDudes team begins by considering searches for books by ISBN. Since *isbn* is a key, an equality query on *isbn* returns at most one record. Therefore, to speed up queries from customers who look for books with a given ISBN, DBDudes decides to build an unclustered hash index on *isbn*.

Next, it considers updates to book quantities. To update the *qty_in_stock* value for a book, we must first search for the book by ISBN; the index on *isbn* speeds this up. Since the *qty_in_stock* value for a book is updated quite frequently, DBDudes also considers partitioning the Books relation vertically into the following two relations:

```
BooksQty(isbn, qty)
BookRest(isbn, title, author, price, year_published)
```

Unfortunately, this vertical partitioning slows down another very popular query: Equality search on ISBN to retrieve all information about a book now requires a join between BooksQty and BookRest. So DBDudes decides not to vertically partition Books.

DBDudes thinks it is likely that customers will also want to search for books by title and by author, and decides to add unclustered hash indexes on *title* and *author*—these indexes are inexpensive to maintain because the set of books is rarely changed even though the quantity in stock for a book changes often.

Next, DBDudes considers the Customers relation. A customer is first identified by the unique customer identification number. So the most common queries on Customers are equality queries involving the customer identification number, and DBDudes decides to build a clustered hash index on *cid* to achieve maximum speed for this query.

Moving on to the Orders relation, DBDudes sees that it is involved in two queries: insertion of new orders and retrieval of existing orders. Both queries involve the *ordernum* attribute as search key and so DBDudes decides to build an index on it. What type of index should this be—a B+ tree or a hash index? Since order numbers are assigned sequentially and correspond to the order date, sorting by *ordernum* effectively sorts by order date as well. So DBDudes decides to build a clustered B+ tree index on *ordernum*. Although the operational requirements mentioned until now favor neither a B+ tree nor a hash index, B&N will probably want to monitor daily activities and the clustered B+ tree is a better choice for such range queries. Of course, this means that retrieving all orders for a given customer could be expensive for customers with many orders, since clustering by *ordernum* precludes clustering by other attributes, such as *cid*.

The Orderlists relation involves mostly insertions, with an occasional update of a shipment date or a query to list all components of a given order. If Orderlists is kept sorted on *ordernum*, all insertions are appends at the end of the relation and thus very efficient. A clustered B+ tree index on *ordernum* maintains this sort order and also speeds up retrieval of all items for a given order. To update

a shipment date, we need to search for a tuple by *ordernum* and *isbn*. The index on *ordernum* helps here as well. Although an index on $\langle \textit{ordernum}, \textit{isbn} \rangle$ would be better for this purpose, insertions would not be as efficient as with an index on just *ordernum*; DBDudes therefore decides to index *Orderlists* on just *ordernum*.

20.11.1 Tuning the Database

Several months after the launch of the B&N site, DBDudes is called in and told that customer enquiries about pending orders are being processed very slowly. B&N has become very successful, and the *Orders* and *Orderlists* tables have grown huge.

Thinking further about the design, DBDudes realizes that there are two types of orders: *completed orders*, for which all books have already shipped, and *partially completed orders*, for which some books are yet to be shipped. Most customer requests to look up an order involve partially completed orders, which are a small fraction of all orders. DBDudes therefore decides to horizontally partition both the *Orders* table and the *Orderlists* table by *ordernum*. This results in four new relations: *NewOrders*, *OldOrders*, *NewOrderlists*, and *OldOrderlists*.

An order and its components are always in exactly one pair of relations—and we can determine which pair, old or new, by a simple check on *ordernum*—and queries involving that order can always be evaluated using only the relevant relations. Some queries are now slower, such as those asking for all of a customer's orders, since they require us to search two sets of relations. However, these queries are infrequent and their performance is acceptable.

20.12 DBMS BENCHMARKING

Thus far, we considered how to improve the design of a database to obtain better performance. As the database grows, however, the underlying DBMS may no longer be able to provide adequate performance, even with the best possible design, and we have to consider upgrading our system, typically by buying faster hardware and additional memory. We may also consider migrating our database to a new DBMS.

When evaluating DBMS products, performance is an important consideration. ADBIVIS is a complex piece of software, and different vendors may target their systems toward different market segments by putting more effort into optimizing certain parts of the system or choosing different system designs. For example, some systems are designed to run complex queries efficiently, while others are designed to run many simple transactions per second. Within

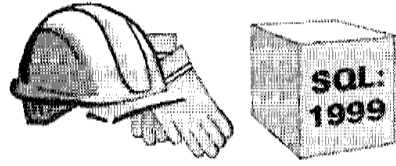
each category of systems, there are many competing products. To assist users in choosing a DBMS that is well suited to their needs, several performance benchmarks have been developed. These include benchmarks for measuring the performance of a certain class of applications (e.g., the TPC benchmarks) and benchmarks for measuring how well a DBMS performs various operations (e.g., the *Visconsin* benchmark).

Benchmarks should be portable, easy to understand, and scale naturally to larger problem instances. They should measure *peak performance* (e.g., *transactions per second*, or *tps*) as well as *price/performance ratios* (e.g., *\$/tps*) for typical workloads in a given application domain. The Transaction Processing Council (TPC) was created to define benchmarks for transaction processing and database systems. Other well-known benchmarks have been proposed by academic researchers and industry organizations. Benchmarks that are proprietary to a given vendor are not very useful for comparing different systems (although they may be useful in determining how well a given system would handle a particular workload).

20.12.1 Well-Known DBMS Benchmarks

Online Transaction Processing Benchmarks: The TPC-A and TPC-B benchmarks constitute the standard definitions of the *tps* and *\$/tps* measures. TPC-A measures the performance and price of a computer network in addition to the DBMS, whereas the TPC-B benchmark considers the DBMS by itself. These benchmarks involve a simple transaction that updates three data records, from three different tables, and appends a record to a fourth table. A number of details (e.g., transaction arrival distribution, interconnect method, system properties) are rigorously specified, ensuring that results for different systems can be meaningfully compared. The TPC-C benchmark is a more complex suite of transactional tasks than TPC-A and TPC-B. It models a warehouse that tracks items supplied to customers and involves five types of transactions. Each TPC-C transaction is much more expensive than a TPC-A or TPC-B transaction, and TPC-C exercises a much wider range of system capabilities, such as use of secondary indexes and transaction aborts. It has more or less completely replaced TPC-A and TPC-B as the standard transaction processing benchmark.

Query Benchmarks: The *Visconsin* benchmark is widely used for measuring the performance of simple relational queries. The *Set Query* benchmark measures the performance of a suite of more complex queries, and the *AS³AP* benchmark measures the performance of a mixed workload of transactions, relational queries, and utility functions. The *TPC-D* benchmark is a suite of complex SQL queries intended to be representative of the (decision-support ap-



21

SECURITY AND AUTHORIZATION

- ☛ What are the main security considerations in designing a database application?
- ☛ What mechanisms does a DBMS provide to control a user's access to data?
- ☛ What is discretionary access control and how is it supported in SQL?
- ☛ What are the weaknesses of discretionary access control? How are these addressed in mandatory access control?
- ☛ What are covert channels and how do they compromise mandatory access control?
- ☛ What must the DBA do to ensure security?
- ☛ What is the added security threat when a database is accessed remotely?
- ☛ What is the role of encryption in ensuring secure access? How is it used for certifying servers and creating digital signatures?
- **Key concepts:** security, integrity, availability; discretionary access control, privileges, GRANT, REVOKE; mandatory access control, objects, subjects, security classes, multilevel tables, polyinstantiation; covert channels, DoD security levels; statistical databases, inferring secure information; authentication for remote access, securing servers, digital signatures; encryption, public-key encryption. -

I know that's a secret, for it's whispered everywhere.

— William Congreve

The data stored in a DBMS is often vital to the business interests of the organization and is regarded as a corporate asset. In addition to protecting the intrinsic value of the data, corporations must consider ways to ensure privacy and control access to data that must not be revealed to certain groups of users for various reasons.

In this chapter, we discuss the concepts underlying access control and security in a DBMS. After introducing database security issues in Section 21.1, we consider two distinct approaches, called *discretionary* and *mandatory*, to specifying and managing access controls. An access control mechanism is a way to control the data accessible by a given user. After introducing access controls in Section 21.2, we cover discretionary access control, which is supported in SQL, in Section 21.3. We briefly cover mandatory access control, which is not supported in SQL, in Section 21.4.

In Section 21.6, we discuss some additional aspects of database security, such as security in a statistical database and the role of the database administrator. We then consider some of the unique challenges in supporting secure access to a DBMS over the Internet, which is a central problem in e-Commerce and other Internet database applications, in Section 21.5. We conclude this chapter with a discussion of security aspects of the Barnes and Noble case study in Section 21.7.

21.1 INTRODUCTION TO DATABASE SECURITY

There are three main objectives when designing a secure database application:

1. **Secrecy:** Information should not be disclosed to unauthorized users. For example, a student should not be allowed to examine other students' grades.
2. **Integrity:** Only authorized users should be allowed to modify data. For example, students may be allowed to see their grades, yet not allowed (obviously) to modify them.
3. **Availability:** Authorized users should not be denied access. For example, an instructor who wishes to change a grade should be allowed to do so.

To achieve these objectives, a clear and consistent security policy should be developed to describe what security measures must be enforced. In particular, we must determine what part of the data is to be protected and which users get access to which portions of the data. Next, the security mechanisms of the underlying DBMS and operating system, as well as external mechanisms,

such as securing access to buildings, must be utilized to enforce the policy. We emphasize that security measures must be taken at several levels.

Security leaks in the OS or network connections can circumvent database security mechanisms. For example, such leaks could allow an intruder to log on as the database administrator, 'with all the attendant DBMS access rights. Human factors are another source of security leaks. For example, a user may choose a password that is easy to guess, or a user who is authorized to see sensitive data may misuse it. Such errors account for a large percentage of security breaches. We do not discuss these aspects of security despite their importance because they are not specific to database management systems; our main focus is on database access control mechanisms to support a security policy.

We observe that views are a valuable tool in enforcing security policies. The view mechanism can be used to create a 'window' on a collection of data that is appropriate for some group of users. Views allow us to limit access to sensitive data by providing access to a restricted version (defined through a view) of that data, rather than to the data itself.

We use the following schemas in our examples:

```
Sailors(sid: integer, sname: string, rating: integer, age: real)
Boats(bid: integer, bname: string, color: string)
Reserves(sid: integer, bid: integer, day: date)
```

Increasingly, as database systems become the backbone of e-Commerce applications requests originate over the Internet. This makes it important to be able to authenticate a user to the database system. After all, enforcing a security policy that allows user Sam to read a table and Ebner to write the table is not of much use if Sam can masquerade as Ebner. Conversely, we must be able to assure users that they are communicating with a legitimate system (e.g., the real Amazon.com server, and not a spurious application intended to steal sensitive information such as a credit card number). While the details of authentication are outside the scope of our coverage, we discuss the role of authentication and the basic ideas involved in Section 21.5, after covering database access control mechanisms.

21.2 ACCESS CONTROL

A database for an enterprise contains a great deal of information and usually has several groups of users. Most users need to access only a small part of the database to carry out their tasks. Allowing users unrestricted access to all the

data can be undesirable, and a DBMS should provide mechanisms to control access to data.

A DBMS offers two main approaches to access control. Discretionary access control is based on the concept of access rights, or privileges, and mechanisms for giving users such privileges. A privilege allows a user to access some data object in a certain manner (e.g., to read or modify). A user who creates a database object such as a table or a view automatically gets all applicable privileges on that object. The DBMS subsequently keeps track of how these privileges are granted to other users, and possibly revoked, and ensures that at all times only users with the necessary privileges can access all object. SQL supports discretionary access control through the GRANT and REVOKE commands. The GRANT command gives privileges to users, and the REVOKE command takes away privileges. We discuss discretionary access control in Section 21.3.

Discretionary access control mechanisms, while generally effective, have certain weaknesses. In particular, a devious unauthorized user can trick an authorized user into disclosing sensitive data. Mandatory access control is based on systemwide policies that cannot be changed by individual users. In this approach each database object is assigned a *security class*, each user is assigned *clearance* for a security class, and rules are imposed on reading and writing of database objects by users. The DBMS determines whether a given user can read or write a given object based on certain rules that involve the security level of the object and the clearance of the user. These rules seek to ensure that sensitive data can never be 'passed on' to a user without the necessary clearance. The SQL standard does not include any support for mandatory access control. We discuss mandatory access control in Section 21.4.

21.3 DISCRETIONARY ACCESS CONTROL

SQL supports discretionary access control through the GRANT and REVOKE commands. The GRANT command gives users privileges to base tables and views. The syntax of this command is as follows:

```
GRANT privileges ON object TO users [WITH GRANT OPTION]
```

For our purposes object is either a base table or a view. SQL recognizes certain other kinds of objects, but we do not discuss them. Several privileges can be specified, including these:

- **SELECT:** The right to access (read) all columns of the table specified as the object, *including columns added later* through ALTER TABLE commands.

- `INSERT(column-name)`: The right to insert rows with (non-*null* or non-default) values in the named column of the table named as object. If this right is to be granted with respect to all columns, including columns that might be added later, we can simply use `INSERT`. The privileges `UPDATE(column-name)` and `UPDATE` are similar.
- `DELETE`: The right to delete rows from the table named as object.
- `REFERENCES(column-name)`: The right to define foreign keys (in other tables) that refer to the specified column of the table object. `REFERENCES` without a column name specified denotes this right with respect to all columns, including any that are added later.

If a user has a privilege with the `grant` option, he or she can **pass** it to another user (with or without the `grant` option) by using the `GRANT` command. A user who creates a base table automatically has all applicable privileges on it, along with the right to grant these privileges to other users. A user who creates a view has precisely those privileges on the view that he or she has on *everyone* of the views or base tables used to define the view. The user creating the view must have the `SELECT` privilege on each underlying table, of course, and so is always granted the `SELECT` privilege on the view. The creator of the view has the `SELECT` privilege with the `grant` option only if he or she has the `SELECT` privilege with the `grant` option on every underlying table. In addition, if the view is updatable and the user holds `INSERT`, `DELETE`, or `UPDATE` privileges (with or without the `grant` option) on the (single) underlying table, the user automatically gets the same privileges on the view.

Only the owner of a schema can execute the data definition statements `CREATE`, `ALTER`, and `DROP` on that schema. The right to execute these statements cannot be granted or revoked.

In conjunction with the `GRANT` and `REVOKE` commands, views are an important component of the security mechanisms provided by a relational DBMS. By defining views on the base tables, we can present needed information to a user while *hiding* other information that the user should not be given access to. For example, consider the following view definition:

```
CREATE VIEW ActiveSailors (name, age, day)
AS SELECT S.name, S.age, R.day
FROM   Sailors S, Reserves R
WHERE  S.sid = R.sid AND S.rating > 6
```

A user who can access `ActiveSailors` but not `Sailors` or `Reserves` knows the names of sailors who have reservations but cannot find out the *bids* of boats reserved by a given sailor.

Role-Based Authorization in SQL: Privileges are assigned to users (authorization IDs, to be precise) in SQL-92. In the real world, privileges are often associated with a user's job or *role* within the organization. Many DBMSs have long supported the concept of a role and allowed privileges to be assigned to roles. Roles can then be granted to users and other roles. (Of course, privileges can also be granted directly to users.) The SQL:1999 standard includes support for roles. Roles can be created and destroyed using the CREATE ROLE and DROP ROLE commands. Users can be granted roles (optionally, with the ability to pass the role on to others). The standard GRANT and REVOKE commands can assign privileges to (and revoke from) roles or authorization IDs.

What is the benefit of including a feature that many systems already support? This ensures that, over time, all vendors who comply with the standard support this feature. Thus, users can use the feature without worrying about portability of their application across DBMSs.

Privileges are assigned in SQL to authorization IDs, which can denote a single user or a group of users; a user must specify an authorization ID and, in many systems, a corresponding *password* before the DBMS accepts any commands from him or her. So, technically, *Joe*, *Michael*, and so on are authorization IDs rather than user names in the following examples.

Suppose that user Joe has created the tables Boats, Reserves, and Sailors. Some examples of the GRANT command that Joe can now execute follow:

```
GRANT INSERT, DELETE ON Reserves TO Yuppy WITH GRANT OPTION
GRANT SELECT ON Reserves TO Michael
GRANT SELECT ON Sailors TO Michael WITH GRANT OPTION
GRANT UPDATE (rating) ON Sailors TO Leah
GRANT REFERENCES (bid) ON Boats TO Bill
```

Yuppy can insert or delete Reserves rows and authorize someone else to do the same. Michael can execute SELECT queries on Sailors and Reserves, and he can pass this privilege to others for Sailors but not for Reserves. With the SELECT privilege, Michael can create a view that accesses the Sailors and Reserves tables (for example, the ActiveSailors view), but he cannot grant SELECT on ActiveSailors to others.

On the other hand, suppose that Michael creates the following view:

```
CREATE VIEW YoungSailors (sid, age, rating)
AS SELECT S.sid, S.age, S.rating
```



```

FROM    Sailors S
WHERE   S.age < 18

```

The only underlying table is `Sailors`, for which Michael has `SELECT` with the grant option. He therefore has `SELECT` with the grant option on `YoungSailors` and can pass on the `SELECT` privilege on `YoungSailors` to Eric and Guppy:

```
GRANT SELECT ON YoungSailors TO Eric, Guppy
```

Eric and Guppy can now execute `SELECT` queries on the view `YoungSailors`—note, however, that Eric and Guppy do *not* have the right to execute `SELECT` queries directly on the underlying `Sailors` table.

Michael can also define constraints based on the information in the `Sailors` and `Reserves` tables. For example, Michael can define the following table, which has an associated table constraint:

```

CREATE TABLE Sneaky (lnaxrating  INTEGER,
                      CHECK (maxrating >=
                             ( SELECT MAX (S.rating)
                               FROM   Sailors S )))

```

By repeatedly inserting rows with gradually increasing *maxrating* values into the `Sneaky` table until an insertion finally succeeds, Michael can find out the highest *rating* value in the `Sailors` table. This example illustrates why SQL requires the creator of a table constraint that refers to `Sailors` to possess the `SELECT` privilege on `Sailors`.

Returning to the privileges granted by Joe, Leah can update only the *rating* column of `Sailors` rows. She can execute the following command, which sets all ratings to 8:

```

UPDATE Sailors S
SET    S.rating = 8

```

However, she cannot execute the same command if the `SET` clause is changed to be `SET S.age = 25`, because she is not allowed to update the *age* field. A more subtle point is illustrated by the following command, which decrements the rating of all 'sailors:

```

UPDATE Sailors S
SET    S.rating = S.rating-1

```

Leah cannot execute this command because it requires the `SELECT` privilege on the *S.rating* column and Leah does not have this privilege.

Bill can refer to the *bid* column of **Boats** as a foreign key in another table. For example, Bill can create the **Reserves** table through the following command:

```
CREATE TABLE Reserves (sid    INTEGER,
                        bid    INTEGER,
                        day    DATE,
                        PRIMARY KEY (bid, day),
                        FOREIGN KEY (sid) REFERENCES Sailors ),
                        FOREIGN KEY (bid) REFERENCES Boats)
```

If Bill did not have the REFERENCES privilege on the *bid* column of **Boats**, he would not be able to execute this CREATE statement because the FOREIGN KEY clause requires this privilege. (A similar point holds with respect to the foreign key reference to **Sailors**.)

Specifying just the INSERT privilege (similarly, REFERENCES and other privileges) in a GRANT command is not the same as specifying SELECT(*column-name*) for each column currently in the table. Consider the following command over the **Sailors** table, which has columns *sid*, *sname*, *rating*, and *age*:

```
GRANT INSERT ON Sailors TO Michael
```

Suppose that this command is executed and then a column is added to the **Sailors** table (by executing an ALTER TABLE command). Note that Michael has the INSERT privilege with respect to the newly added column. If we had executed the following GRANT command, instead of the previous one, Michael would not have the INSERT privilege on the new column:

```
GRANT INSERT ON Sailors(sid), Sailors(sname), Sailors(rating),
Sailors(age), TO Michael
```

There is a complementary command to GRANT that allows the withdrawal of privileges. The syntax of the REVOKE command is as follows:

```
REVOKE [GRANT OPTION FOR ] privileges
      ON object FROM users {RESTRICT | CASCADE }
```

The command can be used to revoke either a privilege or just the grant option on a privilege (by using the optional GRANT OPTION FOR clause). One of the two alternatives, RESTRICT or CASCADE, must be specified; we see what this choice means shortly.

The intuition behind the GRANT command is clear: the creator of a base table or a view is given all the appropriate privileges with respect to it and is allowed

to pass these privileges—including the right to pass along a privilege—to other users. The REVOKE command is, as expected, intended to achieve the reverse: A user who has granted a privilege to another user may change his or her mind and want to withdraw the granted privilege. The intuition behind exactly what effect a REVOKE command has is complicated by the fact that a user may be granted the same privilege multiple times, possibly by different users.

When a user executes a REVOKE command with the CASCADE keyword, the effect is to withdraw the named privileges or grant option from all users who currently hold these privileges *solely* through a GRANT command that was previously executed by the same user who is now executing the REVOKE command. If these users received the privileges with the grant option and passed it along, those recipients in turn lose their privileges as a consequence of the REVOKE command, unless they received these privileges through an additional GRANT command.

We illustrate the REVOKE command through several examples. First, consider what happens after the following sequence of commands, where Joe is the creator of Sailors.

```
GRANT SELECT ON Sailors TO Art WITH GRANT OPTION    (executed by Joe)
GRANT SELECT ON Sailors TO Bob WITH GRANT OPTION    (executed by Art)
REVOKE SELECT ON Sailors FROM Art CASCADE           (executed by Joe)
```

Art loses the SELECT privilege on Sailors, of course. Then Bob, who received this privilege from Art, and only Art, also loses this privilege. Bob's privilege is said to be abandoned when the privilege from which it was derived (Art's SELECT privilege with grant option, in this example) is revoked. When the CASCADE keyword is specified, all abandoned privileges are also revoked (possibly causing privileges held by other users to become abandoned and thereby revoked recursively). If the RESTRICT keyword is specified in the REVOKE command, the command is rejected if revoking the privileges *just* from the users specified in the command would result in other privileges becoming abandoned.

Consider the following sequence, as another example:

```
GRANT SELECT ON Sailors TO Art WITH GRANT OPTION    (executed by Joe)
GRANT SELECT ON Sailors TO Bob WITH GRANT OPTION    (executed by Joe)
GRANT SELECT ON Sailors TO Bob WITH GRANT OPTION    (executed by Art)
REVOKE SELECT ON Sailors FROM Art CASCADE           (executed by Joe)
```

As before, Art loses the SELECT privilege on Sailors. But what about Bob? Bob received this privilege from Art, but he also received it independently

(coincidentally, directly from Joe). So Bob retains this privilege. Consider a third example:

```
GRANT SELECT ON Sailors TO Art WITH GRANT OPTION  (executed by Joe)
GRANT SELECT ON Sailors TO Art WITH GRANT OPTION  (executed by Joe)
REVOKE SELECT ON Sailors FROM Art CASCADE          (executed by Joe)
```

Since Joe granted the privilege to Art twice and only revoked it once, does Art get to keep the privilege? As per the SQL standard, no. Even if Joe absentmindedly granted the same privilege to Art several times, he can revoke it with a single REVOKE command.

It is possible to revoke just the grant option on a privilege:

```
GRANT SELECT ON Sailors TO Art WITH GRANT OPTION  (executed by Joe)
REVOKE GRANT OPTION FOR SELECT ON Sailors
      FROM Art CASCADE                             (executed by Joe)
```

This command would leave Art with the SELECT privilege on Sailors, but Art no longer has the grant option on this privilege and therefore cannot pass it on to other users.

These examples bring out the intuition behind the REVOKE command, and they highlight the complex interaction between GRANT and REVOKE commands. When a GRANT is executed, a privilege descriptor is added to a table of such descriptors maintained by the DBMS. The privilege descriptor specifies the following: the *grantor* of the privilege, the *grantee* who receives the privilege, the *granted privilege* (including the name of the object involved), and whether the grant option is included. When a user creates a table or view and 'automatically' gets certain privileges, a privilege descriptor with *system*, as the grantor is entered into this table.

The effect of a series of GRANT commands can be described in terms of an authorization graph in which the nodes are users—technically, they are authorization IDs—and the arcs indicate how privileges are passed. There is an arc from (the node for) user 1 to user 2 if user 1 executed a GRANT command giving a privilege to user 2; the arc is labeled with the descriptor for the GRANT command. A GRANT command has no effect if the same privileges have already been granted to the same grantee by the same grantor. The following sequence of commands illustrates the semantics of GRANT and REVOKE commands when there is a *cycle* in the authorization graph:

```
GRANT SELECT ON Sailors TO Art WITH GRANT OPTION  (executed by Joe)
GRANT SELECT ON Sailors TO Bob WITH GRANT OPTION  (executed by Art)
```

GRANT SELECT ON Sailors TO Art WITH GRANT OPTION *(executed by Bob)*
 GRANT SELECT ON Sailors TO Cal WITH GRANT OPTION *(executed by Joe)*
 GRANT SELECT ON Sailors TO Bob WITH GRANT OPTION *(executed by Cal)*
 REVOKE SELECT ON Sailors FROM Art CASCADE *(executed by Joe)*

The authorization graph for this example is shown in Figure 21.1. Note that we indicate how Joe, the creator of Sailors, acquired the SELECT privilege from the DBMS by introducing a *System* node and drawing an arc from this node to Joe's node.

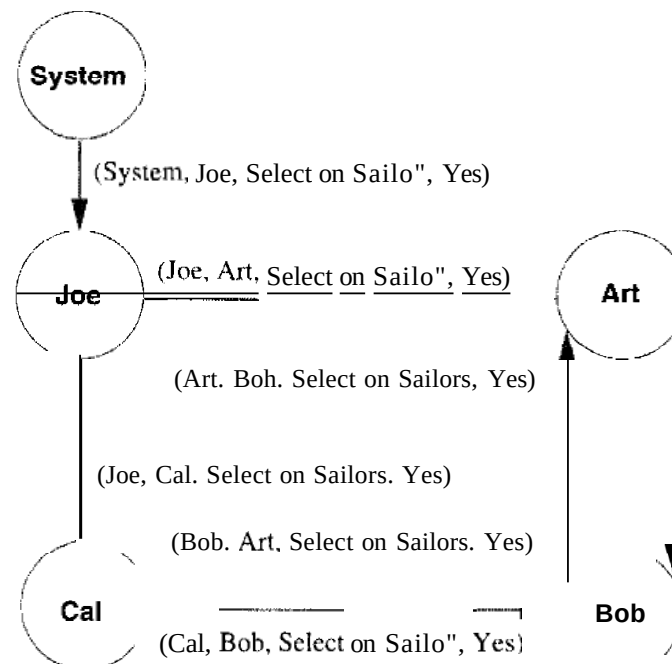


Figure 21.1 Example Authorization Graph

As the graph clearly indicates, Bob's grant to Art and Art's grant to Bob (of the same privilege) creates a cycle. Bob is subsequently given the same privilege by Cal, who received it independently from Joe. At this point Joe decides to revoke the privilege he granted Art.

Let us trace the effect of this revocation. The arc from Joe to Art is removed because it corresponds to the granting action that is revoked. All remaining nodes have the following property: *If node N has an outgoing arc labeled with a privilege, there is a path from the System node to node N in which each arc label contains the same privilege plus the grant option.* That is, any remaining granting action is justified by a privilege received (directly or indirectly) from the System. The execution of Joe's REVOKE command therefore stops at this point, with everyone continuing to hold the SELECT privilege on Sailors.

This result may seem unintuitive because Art continues to have the privilege only because he received it from Bob, and at the time that Bob granted the privilege to Art, he had received it only from Art. Although Bob acquired the privilege through Cal subsequently, should we not undo the effect of his grant

to Art when executing Joe's REVOKE command? The effect of the grant from Bob to Art is *not* undone in SQL. In effect, if a user acquires a privilege multiple times from different grantors, SQL treats each of these grants to the user as having occurred *before* that user passed on the privilege to other users. This implementation of REVOKE is convenient in many real-world situations. For example, if a manager is fired after passing on some privileges to subordinates (who may in turn have passed the privileges to others), we can ensure that only the manager's privileges are removed by first redoing all of the manager's granting actions and then revoking his or her privileges. That is, we need not recursively redo the subordinates' granting actions.

To return to the saga of Joe and his friends, let us suppose that Joe decides to revoke Cal's SELECT privilege as well. Clearly, the arc from Joe to Cal corresponding to the grant of this privilege is removed. The arc from Cal to Bob is removed as well, since there is no longer a path from System to Cal that gives Cal the right to pass the SELECT privilege on Sailors to Bob. The authorization graph at this intermediate point is shown in Figure 21.2.

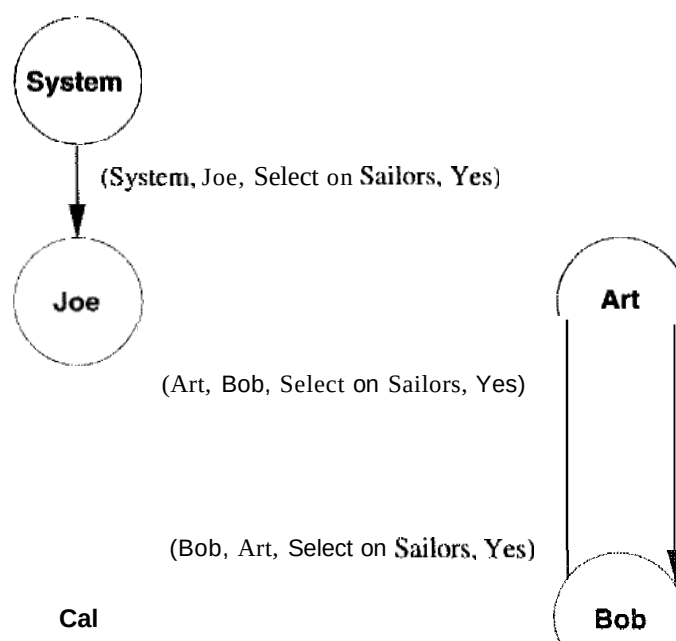


Figure 21.2 Example Authorization Graph during Revocation

The graph now contains two nodes (Art and Bob) for which there are outgoing arcs with labels containing the SELECT privilege on Sailors; therefore, these users have granted this privilege. However, although each node contains an incoming arc carrying the same privilege, *there is no such path from System to either of these nodes*; so these users' right to grant the privilege has been abandoned. We therefore remove the outgoing arcs as well. In general, these nodes might have other arcs incident on them, but in this example, they now have no incident arcs. Joe is left as the only user with the SELECT privilege on Sailors; Art and Bob have lost their privileges.

21.3.1 **Grant** and Revoke on Views and Integrity Constraints

The privileges held by the creator of a view (with respect to the view) change over time as he or she gains or loses privileges on the underlying tables. If the creator loses a privilege held with the grant option, users who were given that privilege on the view lose it as well. There are some subtle aspects to the GRANT and REVOKE commands when they involve views or integrity constraints. We consider some examples that highlight the following important points:

1. A view may be dropped because a SELECT privilege is revoked from the user who created the view.
2. If the creator of a view gains additional privileges on the underlying tables, he or she automatically gains additional privileges on the view.
3. The distinction between the REFERENCES and SELECT privileges is important.

Suppose that Joe created `Sailors` and gave Michael the SELECT privilege on it with the grant option, and Michael then created the view `YoungSailors` and gave Eric the SELECT privilege on `YoungSailors`. Eric now defines a view called `FineYoungSailors`:

```
CREATE VIEW FineYoungSailors (name, age, rating)
AS SELECT S.name, S.age, S.rating
FROM   YoungSailors S
WHERE  S.rating > 6
```

What happens if Joe revokes the SELECT privilege on `Sailors` from Michael? Michael no longer has the authority to execute the query used to define `YoungSailors` because the definition refers to `Sailors`. Therefore, the view `YoungSailors` is dropped (i.e., destroyed). In turn, `FineYoungSailors` is dropped as well. Both view definitions are removed from the system catalogs; even if, nevertheless, Joe decides to give back the SELECT privilege on `Sailors` to Michael, the views are gone and must be created afresh if they are required.

On a more happy note, suppose that everything proceeds as just described until Eric defines `FineYoungSailors`; then, instead of revoking the SELECT privilege on `Sailors` from Michael, Joe decides to also give Michael the INSERT privilege on `Sailors`. Michael's privileges on the view `YoungSailors` are upgraded to what he would have if he were to create the view *now*. He therefore acquires the INSERT privilege on `YoungSailors` as well. (Note that this view is updated.) What about Eric? His privileges are unchanged.

Whether or not Michael has the INSERT privilege on `YoungSailors` with the grant option depends on whether or not Joe gives him the INSERT privilege on

Sailors with the grant option. To understand this situation, consider Eric again. If Michael has the INSERT privilege on YoungSailors with the grant option, he can pass this privilege to Eric. Eric could then insert rows into the Sailors table because inserts on YoungSailors are effected by modifying the underlying base table, Sailors. Clearly, we do not want Michael to be able to authorize Eric to make such changes unless Michael has the INSERT privilege on Sailors with the grant option.

The REFERENCES privilege is very different from the SELECT privilege, as the following example illustrates. Suppose that Joe is the creator of Boats. He can authorize another user, say, Fred, to create Reserves with a foreign key that refers to the *bid* column of Boats by giving Fred the REFERENCES privilege with respect to this column. On the other hand, if Fred has the SELECT privilege on the *bid* column of Boats but not the REFERENCES privilege, Fred *cannot* create Reserves with a foreign key that refers to Boats. If Fred creates Reserves with a foreign key column that refers to *bid* in Boats and later loses the REFERENCES privilege on the *bid* column of boats, the foreign key constraint in Reserves is dropped; however, the Reserves table is *not* dropped.

To understand why the SQL standard chose to introduce the REFERENCES privilege rather than to simply allow the SELECT privilege to be used in this situation, consider what happens if the definition of Reserves specified the NO ACTION option with the foreign key. Joe, the owner of Boats, may be prevented from deleting a row from Boats because a row in Reserves refers to this Boats row. Giving Fred, the creator of Reserves, the right to constrain updates on Boats in this manner goes beyond, simply allowing him to read the values in Boats, which is all that the SELECT privilege authorizes.

21.4 MANDATORY ACCESS CONTROL

Discretionary access control mechanisms, while generally effective, have certain weaknesses. In particular they are susceptible to *Trojan horse* schemes whereby a devious unauthorized user can trick an authorized user into disclosing sensitive data. For example, suppose that student Ricky Dick wants to break into the grade tables of instructor Justin. Dick does the following:

- He creates a new table called MineAllMine and gives INSERT privileges on this table to Justin (who is blissfully unaware of all this attention, of course).
- He modifies the code of the IBIVIS application that Justin uses often to do a couple of additional things: first, read the Grades table, and next, write the result into MineAllMine.

Then he sits back and waits for the grades to be copied into MineAllMine and later undoes the modifications to the application to ensure that Justin does not somehow find out later that he has been cheated. Thus, despite the DBMS enforcing all discretionary access controls—only Justin's authorized code was allowed to access Grades—sensitive data is disclosed to an intruder. The fact that Dick could surreptitiously modify Justin's code is outside the scope of the DBMS's access control mechanism.

Mandatory access control mechanisms are aimed at addressing such loopholes in discretionary access control. The popular model for mandatory access control, called the Bell-LaPadula model, is described in terms of objects (e.g., tables, views, rows, columns), subjects (e.g., users, programs), security classes, and clearances. Each database object is assigned a *security class*, and each subject is assigned *clearance* for a security class; we denote the class of an object or subject A as $class(A)$. The security classes in a system are organized according to a partial order, with a most secure class and a least secure class. For simplicity, we assume that there are four classes: *top secret* (TS), *secret* (S), *confidential* (C), and *unclassified* (U). In this system, $TS > S > C > U$, where $A > B$ means that class A data is more sensitive than class B data.

The Bell-LaPadula model imposes two restrictions on all reads and writes of database objects:

1. **Simple Security Property:** Subject S is allowed to read object O only if $class(O) \geq class(S)$. For example, a user with TS clearance can read a table with C clearance, but a user with C clearance is not allowed to read a table with TS classification.
2. ***-Property:** Subject S is allowed to write object O only if $class(S) \leq class(O)$. For example, a user with S clearance can write only objects with S or TS classification.

If discretionary access controls are also specified, these rules represent additional restrictions. Therefore, to read or write a database object, a user must have the necessary privileges (obtained via GRANT commands) *and* the security classes of the user and the object must satisfy the preceding restrictions. Let us consider how such a mandatory control mechanism might have foiled Tricky Dick. If the Grades table could be classified as S , Justin could be given clearance for S , and Tricky Dick could be given a lower clearance (C). Dick can create objects of only C or lower classification; so the table MineAllMine can have at most the classification C . When the application program running on behalf of Justin (and therefore with clearance S) tries to copy Grades into MineAllMine, it is not allowed to do so because $class(MineAllMine) < class(application)$, and the *-Property is violated.

21.4.1 Multilevel Relations and Polyinstantiation

ro apply Inandatory access control policies in a relational DBMS, a security class must be assigned to each database object. The objects can be at the granularity of tables, rows, or even individual columnn values. Let us assU11le that each row is assigned a security class. This situation leads to the concept of a multilevel table, which is a table with the surprising property that users with different security clearances see a different collection of rows when they access the sarne table.

Consider the instance of the Boats table shown in Figure 21.3. Users with *S* and *TS* clearance get both rows in the answer when they ask to see all rows in Boats. A user with *C* clearance gets only the second row, and a user with *[J* clearance gets no rows.

<i>bid</i>	<i>bname</i>	<i>color</i>	Security Class
101	Salsa	Red	<i>S</i>
102	Pinto	Brown	<i>C</i>

Figure 21.3 An Instance *B1* of Boats

The Boats table is defined to have *bid* as the prirnary key. Suppose that a user with clearance *C* wishes to enter the row (101, *Picante*, *Scarlet*, *C*). We have a dilemrna:

- If the insertion is perlnitted, two distinct rows in the table have key 101.
- If the insertion is not pennitted because the priInary key constraint is vio-lated, the user trying to insert the new row, who has clearance *C*, can infer that there is a boat with *bid*=101 whose security class is higher than *C*. This situation cOlnpromises the principle that users should not be able to infer any infonnation about objects that have a higher security classification.

This dilerrlllla is resolved by effectively treating the security classification as part of the key. rrrhus, the insertion is allo\ved to continue, and the table instance is rmodified as shown in Figure 21.4.

<i>bid</i>	<i>bna'me</i>	<i>color</i>	Security Class
101	Salsa	Red	<i>S</i>
101	Picante	Scarlet	<i>C</i>
102	Pinto	Brown	<i>C</i>

Figure 21.4 Insta.nce 131 after Insertion

Users with clearance *C* or *[1]* see just the rows for Picante and Pinto, but users with clearance *S* or *TS* see all three rows. The two rows with *bid=101* can be interpreted in one of two ways: only the row with the higher classification (Salsa, with classification 8) actually exists, or both exist and their presence is revealed to users according to their clearance level. The choice of interpretation is up to application developers and users.

The presence of data objects that appear to have different values to users with different clearances (for example, the boat with *bid 101*) is called *polyinstantiation*. If we consider security classifications associated with individual columns, the intuition underlying polyinstantiation can be generalized in a straightforward manner, but some additional details must be addressed. We remark that the main drawback of mandatory access control schemes is their rigidity; policies are set by system administrators, and the classification mechanisms are not flexible enough. A satisfactory combination of discretionary and mandatory access controls is yet to be achieved.

21.4.2 Covert Channels, DoD Security Levels

Even if a DBMS enforces the mandatory access control scheme just discussed, information can flow from a higher classification level to a lower classification level through indirect means, called *covert channels*. For example, if a transaction accesses data at more than one site in a distributed DBMS, the actions at the two sites must be coordinated. The process at one site may have a lower clearance (say, *C*) than the process at another site (say, *S*), and both processes have to agree to commit before the transaction can be committed. This requirement can be exploited to pass information with an *S* classification to the process with a *C* clearance: The transaction is repeatedly invoked, and the process with the *C* clearance always agrees to commit, whereas the process with the *S* clearance agrees to commit if it wants to transmit a 1 bit and does not agree if it wants to transmit a 0 bit.

In this (admittedly tortuous) manner, information with an *S* clearance can be sent to a process with a *C* clearance as a stream of bits. This covert channel is an indirect violation of the intent behind the *-Property. Additional examples of covert channels can be found readily in statistical databases, which we discuss in Section 21.5.2.

DBMS vendors recently started implementing mandatory access control mechanisms (although they are not part of the SQL standard) because the United States Department of Defense (1991) requires such support for its systems. The DoD requirements can be described in terms of security levels *A*, *B*, *C*, and *D*, of which *A* is the most secure and *D* is the least secure.

Current Systems: Commercial RDBMSs are available that support discretionary controls at the *C2* level and mandatory controls at the *B1* level. IBM DB2, Informix, Microsoft SQL Server, Oracle 8, and Sybase ASE all support SQL's features for discretionary access control. In general, they do not support mandatory access control; Oracle offers a version of their product with support for mandatory access control.

Level *C* requires support for discretionary access control. It is divided into sublevels *C1* and *C2*; *C2* also requires some degree of accountability through procedures such as login verification and audit trails. Level *B* requires support for mandatory access control. It is subdivided into levels *B1*, *B2*, and *B3*. Level *B3* additionally requires the identification and elimination of covert channels. Level *B3* additionally requires maintenance of audit trails and the designation of a security administrator (usually, but not necessarily, the DBA). Level *A*, the most secure level, requires a mathematical proof that the security mechanism enforces the security policy!

21.5 SECURITY FOR INTERNET APPLICATIONS

When a DBMS is accessed from a secure location, we can rely upon a simple password mechanism for authenticating users. However, suppose our friend Sam wants to place an order for a book over the Internet. This presents some unique challenges: Sam is not even a known user (unless he is a repeat customer). From Amazon's point of view, we have an individual asking for a book and offering to pay with a credit card registered to Sam, but is this individual really Sam? From Sam's point of view, he sees a form asking for credit card information, but is this indeed a legitimate part of Amazon's site, and not a rogue application designed to trick him into revealing his credit card number?

This example illustrates the need for a more sophisticated approach to authentication than a simple password mechanism. Encryption techniques provide the foundation for modern authentication.

21.5.1 Encryption

The basic idea behind encryption is to apply an encryption algorithm to the data, using a user-specified or IJBA-Specified encryption key. The output of the algorithm is the encrypted version of the data. There is also a decryption algorithm which takes the encrypted data and a decryption key as input and then returns the original data. Without the correct decryption key, the decryption algorithm produces gibberish. The encryption and decryption

DES and AES: The DES standard, adopted in 1977, has a 56-bit encryption key. Over time, computers have become so fast that, in 1999, a special-purpose chip and a network of PCs were used to crack DES in under a day. The system was testing 245 billion keys per second when the correct key was found! It is estimated that a special-purpose hardware device can be built for under a billion dollars that can crack DES in under four hours. Despite growing concerns about its vulnerability, DES is still widely used. In 2000, a successor to DES, called the Advanced Encryption Standard (AES), was adopted as the new (symmetric) encryption standard. AES has three possible key sizes: 128, 192, and 256 bits. With a 128 bit key size, there are over $3 \cdot 10^{38}$ possible AES keys, which is on the order of 10^{24} more than the number of 56-bit DES keys. Assume that we could build a computer fast enough to crack DES in 1 second. This computer would compute for about 149 trillion years to crack a 128-bit AES key. (Experts think the universe is less than 20 billion years old.)

algorithms themselves are assumed to be publicly known, but one or both keys are secret (depending upon the encryption scheme).

In symmetric encryption, the encryption key is also used as the decryption key. The ANSI Data Encryption Standard (DES), which has been in use since 1977, is a well-known example of symmetric encryption. It uses an encryption algorithm that consists of character substitutions and permutations. The main weakness of symmetric encryption is that all authorized users must be told the key, increasing the likelihood of its becoming known to an intruder (e.g., by simple human error).

Another approach to encryption, called public-key encryption, has become increasingly popular in recent years. The encryption scheme proposed by Rivest, Shamir, and Adleman, called RSA, is a well-known example of public-key encryption. Each authorized user has a public encryption key, known to everyone, and a private decryption key, known only to him or her. Since the private decryption keys are known only to their owners, the weakness of DES is avoided.

A central issue for public-key encryption is how encryption and decryption keys are chosen. Technically, public-key encryption algorithms rely on the existence of one-way functions, whose inverses are computationally very hard to determine. The RSA algorithm, for example, is based on the observation that, although checking whether a given number is prime is easy, determining the prime factors of a nonprime number is extremely hard. (Determining the

Why RSA Works: The essential point of the scheme is that it is easy to compute d given e , p , and q , but *very* hard to compute d given just e and L . In turn, this difficulty depends on the fact that it is hard to determine the prime factors of L , which happen to be p and q . *A caveat:* Factoring is widely believed to be hard, but there is no proof that this is so. Nor is there a proof that factoring is the only way to crack RSA; that is, to comp d from e and L .

prime factors of a number with over 100 digits can take years of CPU time on the fastest available computers today.)

We now sketch the idea behind the RSA algorithm, assuming that the data to be encrypted is an integer I . To choose an encryption key and a decryption key for a given user, we first choose a very large integer L , larger than the largest integer we will ever need to encode.¹ We then select a number e as the encryption key and compute the decryption key d based on e and L ; how this is done is central to the approach, as we see shortly. Both L and e are made public and used by the encryption algorithm. However, d is kept secret and is necessary for decryption.

- The encryption function is $S = Ie \bmod L$.
- The decryption function is $I = Sd \bmod L$.

We choose L to be the product of two large (e.g., 1024-bit), distinct prime numbers, $p * q$. The encryption key e is a randomly chosen number between 1 and L that is relatively prime to $(p - 1) * (q - 1)$. The decryption key d is computed such that $d * e = 1 \bmod ((p - 1) * (q - 1))$. Given these choices, results in number theory can be used to prove that the decryption function recovers the original message from its encrypted version.

A very important property of the encryption and decryption algorithms is that the roles of the encryption and decryption keys can be reversed:

$$\text{decrypt}(d, (\text{encrypt}(e, I))) = I = \text{decrypt}(e, (\text{encrypt}(d, I)))$$

Since many protocols rely on this property, we henceforth simply refer to public and private keys (since both keys can be used for encryption as well as decryption).

¹A message that is to be encrypted is decomposed into blocks such that each block can be treated as an integer less than L .

While we introduced encryption in the context of authentication, we note that it is a fundamental tool for enforcing security. A DBMS can use *encryption* to protect information in situations where the normal security mechanisms of the DBMS are not adequate. For example, an intruder may steal tapes containing sensitive data or tap a communication line. By storing and transmitting data in an encrypted form, the DBMS ensures that such stolen data is not intelligible to the intruder.

21.5.2 Certifying Servers: The SSL Protocol

Suppose we associate a public key and a decryption key with Amazon. Anyone, say, user Sam, can send Amazon an order by encrypting the order using Amazon's public key. Only Amazon can decrypt this secret order because the decryption algorithm requires Amazon's private key, known only to Amazon.

This hinges on Sam's ability to reliably find out Amazon's public key. A number of companies serve as certification authorities, e.g., Verisign. Amazon generates a public encryption key e_A (and a private decryption key) and sends the public key to Verisign. Verisign then issues a certificate to Amazon that contains the following information:

(Verisign Amazon, <http://www.amazon.com>, e_A)

The certificate is encrypted using Verisign's own *private* key, which is known to (i.e., stored in) Internet Explorer, Netscape Navigator, and other browsers.

When Sam comes to the Amazon site and wants to place an order, his browser, running the SSL protocol,² asks the server for the Verisign certificate. The browser then validates the certificate by decrypting it (using Verisign's public key) and checking that the result is a certificate with the name Verisign, and that the URL it contains is that of the server it is talking to. (Note that an attempt to forge a certificate will fail because certificates are encrypted using Verisign's private key, which is known only to Verisign.) Next, the browser generates a random session key, encrypts it using Amazon's public key (which it obtained from the validated certificate and therefore trusts), and sends it to the Amazon server.

From this point on, the Amazon server and the browser can use the session key (which both know and are confident that only they know) and a *symmetric* encryption protocol like AES or DES to exchange securely encrypted messages: Messages are encrypted by the sender and decrypted by the receiver using the same session key. The encrypted messages travel over the Internet and may be

²A browser uses the SSL protocol if the target URL begins with *https*.

intercepted, but they cannot be decrypted without the session key. It is useful to consider why we need a session key; after all, the browser could simply have encrypted Sam's original request using Amazon's public key and sent it securely to the Amazon server. The reason is that, without the session key, the Amazon server has no way to securely send information back to the browser. A further advantage of session keys is that symmetric encryption is computationally much faster than public key encryption. The session key is discarded at the end of the session.

Thus, Sam can be assured that only Amazon can see the information he types into the form shown to him by the Amazon server and the information sent back to him in responses from the server. However, at this point, Amazon has no assurance that the user running the browser is actually Sam, and not Someone who has stolen Sam's credit card. Typically, merchants accept this situation, which also arises when a customer places an order over the phone.

If we want to be sure of the user's identity, this can be accomplished by additionally requiring the user to login. In our example, Sam must first establish an account with Amazon and select a password. (Sam's identity is originally established by calling him back on the phone to verify the account information or by sending email to an email address; in the latter case, all we establish is that the owner of the account is the individual with the given email address.) Whenever he visits the site and Amazon needs to verify his identity, Amazon redirects him to a login form *after* using SSL to establish a session key. The password typed in is transmitted securely by encrypting it with the session key.

One remaining drawback in this approach is that Amazon now knows Sam's credit card number, and he must trust Amazon not to misuse it. The Secure Electronic Transaction protocol addresses this limitation. Every customer must now obtain a certificate, with his or her own private and public keys, and every transaction involves the Amazon server, the customer's browser, and the server of a trusted third party, such as Visa for credit card transactions. The basic idea is that the browser encodes non-credit card information using Amazon's public key and the credit card information using Visa's public key and sends these to the Amazon server, which forwards the credit card information (which it cannot decrypt) to the Visa server. If the Visa server approves the information, the transaction goes through.

21.5.3 Digital Signatures

Suppose that Elmer, who works for Amazon, and Betsy, who works for McGraw-Hill, need to communicate with each other about inventory. Public key encryption can be used to create digital signatures for messages. That is, messages

can be encoded in such a way that, if Elmer gets a message supposedly from Betsy, he can verify that it is from Betsy (in addition to being able to decrypt the message) and, further, *prove* that it is from Betsy at McGraw-Hill, even if the message is sent from a Hotmail account when Betsy is traveling. Similarly, Betsy can authenticate the originator of messages from Elmer.

If Elmer encrypts messages for Betsy using her public key, and vice-versa, they can exchange information securely but cannot authenticate the sender. Someone who wishes to impersonate Betsy could use her public key to send a message to Elmer, pretending to be Betsy.

A clever use of the encryption scheme, however, allows Elmer to verify whether the message was indeed sent by Betsy. Betsy encrypts the message using her *private* key and then encrypts the result using Elmer's public key. When Elmer receives such a message, he first decrypts it using his private key and then decrypts the result using Betsy's public key. This step yields the original unencrypted message. Furthermore, Elmer can be certain that the message was composed and encrypted by Betsy because a forger could not have known her private key, and without it the final result would have been nonsensical, rather than a legible message. Further, because even Elmer does not know Betsy's private key, Betsy cannot claim that Elmer forged the message.

If authenticating the sender is the objective and hiding the message is not important, we can reduce the cost of encryption by using a message signature. A signature is obtained by applying a one-way function (e.g., a hashing scheme) to the message and is considerably smaller. We encode the signature as in the basic digital signature approach, and send the encoded signature together with the full, unencoded message. The recipient can verify the sender of the signature as just described, and validate the message itself by applying the one-way function and comparing the result with the signature.

21.6 ADDITIONAL ISSUES RELATED TO SECURITY

Security is a broad topic, and our coverage is necessarily limited. This section briefly touches on some additional important issues.

21.6.1 Role of the Database Administrator

The database administrator (DBA) plays an important role in enforcing the security-related aspects of a database design. In conjunction with the owners of the data, the DBA also contributes to developing a security policy. The DBA has a special account, which we call the system **account**, and is responsible

for the overall security of the system. In particular, the DBA deals with the following:

1. **Creating New Accounts:** Each new user or group of users must be assigned an authorization ID and a password. Note that application programs that access the database have the same authorization ID as the user executing the program.
2. **Mandatory Control Issues:** If the DBMS supports mandatory control—some customized systems for applications with very high security requirements (for example, military data) provide such support—the DBA must assign security classes to each database object and assign security clearances to each authorization ID in accordance with the chosen security policy.

The DBA is also responsible for maintaining the audit trail, which is essentially the log of updates with the authorization ID (of the user executing the transaction) added to each log entry. This log is just a minor extension of the log mechanism used to recover from crashes. Additionally, the DBA may choose to maintain a log of *all* actions, including reads, performed by a user. Analyzing such histories of how the DBMS was accessed can help prevent security violations by identifying suspicious patterns before an intruder finally succeeds in breaking in, or it can help track down an intruder after a violation has been detected.

21.6.2 Security in Statistical Databases

A statistical database contains specific information on individuals or events but is intended to permit only statistical queries. For example, if we maintained a statistical database of information about sailors, we would allow statistical queries about average ratings, maximum age, and so on, but not queries about individual sailors. Security in such databases poses new problems because it is possible to infer protected information (such as a sailor's rating) from answers to permitted statistical queries. Such inference opportunities represent covert channels that can compromise the security policy of the database.

Suppose that sailor Sneaky Pete wants to know the rating of Admiral Horntooter, the esteemed chairman of the sailing club, and happens to know that Horntooter is the oldest sailor in the club. Pete repeatedly asks queries of the form “How many sailors are there whose age is greater than X ?” for various values of X , until the answer is 1. Obviously, this sailor is Horntooter, the oldest sailor. Note that each of these queries is a valid statistical query and is permitted. Let the value of X at this point be, say, 65. Pete now asks the query, “What is the maximum rating of all sailors whose age is greater than

65?" Again, this query is permitted because it is a statistical query. However, the answer to this query reveals Horntooter's rating to Pete, and the security policy of the database is violated.

One approach to preventing such violations is to require that each query must involve at least some minimum number, say, N , of rows. With a reasonable choice of N , Pete would not be able to isolate the information about Horntooter, because the query about the maximum rating would fail. This restriction, however, is easy to overcome. By repeatedly asking queries of the form, "How many sailors are there whose age is greater than X ?" until the system rejects one such query, Pete identifies a set of N sailors, including Horntooter. Let the value of X at this point be 55. Now, Pete can ask two queries:

- "What is the sum of the ratings of all sailors whose age is greater than 55?" Since N sailors have age greater than 55, this query is permitted.
- "What is the sum of the ratings of all sailors, other than Horntooter, whose age is greater than 55, and sailor Pete?" Since the set of sailors whose ratings are added up now includes Pete instead of Horntooter, but is otherwise the same, the number of sailors involved is still N , and this query is also permitted.

From the answers to these two queries, say, A_1 and A_2 , Pete, who knows his rating, can easily calculate Horntooter's rating as $A_1 - A_2 + \text{Pete's rating}$.

Pete succeeded because he was able to ask two queries that involved many of the same sailors. The number of rows examined in common by two queries is called their intersection. If a limit were to be placed on the amount of intersection permitted between any two queries issued by the same user, Pete could be foiled. Actually, a truly fiendish (and patient) user can generally find out information about specific individuals even if the system places a minimum number of rows bound (N) and a maximum intersection bound (M) on queries, but the number of queries required to do this grows in proportion to N/M . We can try to additionally limit the total number of queries that a user is allowed to ask, but two users could still conspire to breach security. By maintaining a log of all activity (including read-only accesses), such query patterns can be detected, ideally before a security violation occurs. This discussion should make it clear, however, that security in statistical databases is difficult to enforce.

21.7 DESIGN CASE STUDY: THE INTERNET STORE

We return to our case study and our friends at DBI to consider security issues. There are three groups of users: customers, employees, and the owner of the book shop. (Of course, there is also the database administrator, who

has universal access to all data and is responsible for regular operation of the database system.)

The owner of the store has full privileges on all tables. Customers can query the Books table and place orders online, but they should not have access to other customers' records nor to other customers' orders. DBDudes restricts access in two ways. First, it designs a simple Web page with several forms similar to the page shown in Figure 7.1 in Chapter 7. This allows customers to submit a small collection of valid requests without giving them the ability to directly access the underlying DBMS through an SQL interface. Second, DBDudes uses the security features of the DBMS to limit access to sensitive data.

The Web page allows customers to query the Books relation by ISBN number, name of the author, and title of a book. The webpage also has two buttons. The first button retrieves a list of all of the customer's orders that are not completely fulfilled yet. The second button displays a list of all completed orders for that customer. Note that customers cannot specify actual SQL queries through the Web but only fill in some parameters in a form to instantiate an automatically generated SQL query. All queries generated through form input have a WHERE clause that includes the *cid* attribute value of the current customer, and evaluation of the queries generated by the two buttons requires knowledge of the customer identification number. Since all users have to log on to the website before browsing the catalog, the business logic (discussed in Section 7.7) must maintain state information about a customer (i.e., the customer identification number) during the customer's visit to the website.

The second step is to configure the database to limit access according to each user group's need to know. DBDudes creates a special customer account that has the following privileges:

```
SELECT ON Books, NewOrders, OldOrders, NewOrderlists, OldOrderlists
INSERT ON NewOrders, OldOrders, NewOrderlists, OldOrderlists
```

Employees should be able to add new books to the catalog, update the quantity of a book in stock, revise customer orders if necessary, and update all customer information *except the credit card information*. In fact, employees should not even be able to see a customer's credit card number. Therefore, DBDudes creates the following view:

```
CREATE VIEW CustomerInfo (cid,cname,address)
AS SELECT C.cid, C.cname, C.address
FROM Customers C
```

DBDudes gives the employee account the following privileges:

```

SELECT ON CustomerInfo, Books,
        NewOrders, ()IdOrders, NewOrderlists, Old()rderlists
INSERT ON CllstorerInfo, Books,
        Nc\vOrders, OldC)rders, NewOrderlists, ()ldOrderlists
UPDATE ON CustolnerInfo, Books,
        New()rders, OldOrders, NewOrderlists, Old()rderlists
DELETE ON Books, NewOrders, OldOrders, NewOrderlists, ()ldOrderlists

```

Observe that employees can modify CustomerInfo and even insert tuples into it. This is possible because they have the necessary privileges, and further, the view is updatable and insertable-into. While it seems reasonable that employees can update a customer's address, it does seem odd that they can insert a tuple into CllstorerInfo even though they cannot see related information about the customer (i.e., credit card number) in the Cllstomers table. The reason for this is that the store wants to be able to take orders from first-time customers over the phone without asking for credit card information over the phone. Employees can insert into CustomerInfo, effectively creating a new Customers record without credit card information, and customers can subsequently provide the credit card number through a Web interface. (Obviously, the order is not shipped until they do this.)

In addition, there are security issues when the user first logs on to the website using the customer identification number. Sending the number unencrypted over the Internet is a security hazard, and a secure protocol such as SSL should be used.

Companies such as CyberCash and DigiCash offer electronic commerce payment solutions, even including *electronic cash*. Discussion of how to incorporate such techniques into the website are outside the scope of this book.

21.8 REVIEW QUESTIONS

Answers to the review questions can be found in the listed sections.

- What are the main objectives in designing a secure database application? Explain the terms *secrecy*, *integrity*, *availability*, and *authentication*. (Section 21.1)
- Explain the terms *security policy* and *security mechanism* and how they are related. (Section 21.1)
- What is the main idea behind *discretionary access control*? What is the idea behind *mandatory access control*? What are the relative merits of these two approaches? (Section 21.2)

- Describe the privileges recognized in SQL? In particular, describe SELECT, INSERT, UPDATE, DELETE, and REFERENCES. For each privilege, indicate who acquires it automatically on a given table. (Section 21.3)
- How are the owners of privileges identified? In particular, discuss *authorization ID*s and *roles*. (Section 21.3)
- What is an *authorization graph*? Explain SQL's GRANT and REVOKE commands in terms of their effect on this graph. In particular, discuss what happens when users pass on privileges that they receive from someone else. (Section 21.3)
- Discuss the difference between having a privilege on a table and on a view defined over the table. In particular, how can a user have a privilege (say, SELECT) over a view without also having it on all underlying tables? Who must have appropriate privileges on all underlying tables of the view? (Section 21.3.1)
- What are *objects*, *subjects*, *security classes*, and *clearances* in mandatory access control? Discuss the Bell-LaPadula restrictions in terms of these concepts. Specifically, define the *simple security property* and the **-property*. (Section 21.4)
- What is a *Trojan horse* attack and how can it compromise discretionary access control? Explain how mandatory access control protects against Trojan horse attacks. (Section 21.4)
- What do the terms *multilevel table* and *polyinstantiation* mean? Explain their relationship, and how they arise in the context of mandatory access control. (Section 21.4.1)
- What are *covert channels* and how can they arise when both discretionary and mandatory access controls are in place? (Section 21.4.2)
- Discuss the DoD security levels for database systems. (Section 21.4.2)
- Explain why a simple password mechanism is insufficient for authentication of users who access a database remotely, say, over the Internet. (Section 21.5)
- What is the difference between *symmetric* and *public-key encryption*? Give examples of well-known encryption algorithms of both kinds. What is the main weakness of symmetric encryption and how is this addressed in public-key encryption? (Section 21.5.1)
- Discuss the choice of encryption and decryption keys in public-key encryption and how they are used to encrypt and decrypt data. Explain the role of *one-way functions*. What assurance do we have that the RSA scheme cannot be compromised? (Section 21.5.1)

- What are *certification authorities* and why are they needed? Explain how *certificates* are issued to sites and validated by a browser using the *SSL protocol*; discuss the role of the *session key*. (Section 21.5.2)
- If a user connects to a site using the SSL protocol, explain why there is still a need to login the user. Explain the use of SSL to protect passwords and other sensitive information being exchanged. What is the *secure electronic transaction protocol*? What is the added value over SSL? (Section 21.5.2)
- A *digital signature* facilitates secure exchange of messages. Explain what it is and how it goes beyond simply encrypting messages. Discuss the use of *message signatures* to reduce the cost of encryption. (Section 21.5.3)
- What is the role of the database administrator with respect to security? (Section 21.6.1)
- Discuss the additional security loopholes introduced in *statistical databases*. (Section 21.6.2)

EXERCISES

Exercise 21.1 Briefly answer the following questions:

1. Explain the intuition behind the two rules in the Bell-LaPadula model for mandatory access control.
2. Give an example of how covert channels can be used to defeat the Bell-LaPadula model.
3. Give an example of polyinstantiation.
4. Describe a scenario in which mandatory access controls prevent a breach of security that cannot be prevented through discretionary controls.
5. Describe a scenario in which discretionary access controls are required to enforce a security policy that cannot be enforced using only mandatory controls.
6. If a DBMS already supports discretionary and mandatory access controls, is there a need for encryption?
7. Explain the need for each of the following limits in a statistical database system:
 - (a) A maximum on the number of queries a user can pose.
 - (b) A minimum on the number of tuples involved in answering a query.
 - (c) A maximum on the intersection of two queries (i.e., on the number of tuples that both queries examine).
8. Explain the use of an audit trail, with special reference to a statistical database system.
9. What is the role of the DBA with respect to security?
10. Describe AES and its relationship to DES.
11. What is public-key encryption? How does it differ from the encryption approach taken in the Data Encryption Standard (DES), and in what ways is it better than DES?