

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.l.cldress
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 cllstolner 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 Blain 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 IDs* 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

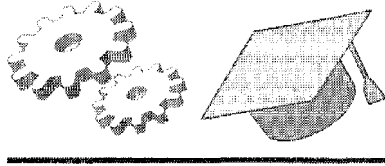
Exercise 21.1 Briefly answer the following questions:

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?

12. Explain how a company offering services on the Internet could use encryption-based techniques to make its order-entry process secure. Discuss the role of DES, AES, SSL, SET, and digital signatures. Search the Web to find out more about related techniques such as *electronic cash*.

Exercise 21.2 You are the DBA for the VeryFine Toy Company and create a relation called *Employees* with fields *ename*, *dept*, and *salary*. For authorization reasons, you also define views *EmployeeNames* (with *ename* as the only attribute) and *DeptInfo* with fields *dept* and *avgsalary*. The latter lists the average salary for each department.

1. Show the view definition statements for *EmployeeNames* and *DeptInfo*.
2. What privileges should be granted to a user who needs to know only average department salaries for the Toy and CS departments?
3. You want to authorize your secretary to fire people (you will probably tell him whom to fire, but you want to be able to delegate this task), to check on who is an employee, and to check on average department salaries. What privileges should you grant?
4. Continuing with the preceding scenario, you do not want your secretary to be able to look at the salaries of individuals. Does your answer to the previous question ensure this? Be specific: Can your secretary possibly find out salaries of *some* individuals (depending on the actual set of tuples), or can your secretary always find out the salary of any individual he wants to?
5. You want to give your secretary the authority to allow other people to read the *EmployeeNames* view. Show the appropriate command.
6. Your secretary defines two new views using the *EmployeeNames* view. The first is called *AtorNames* and simply selects names that begin with a letter in the range A to R. The second is called *HowManyNames* and counts the number of names. You are so pleased with this achievement that you decide to give your secretary the right to insert tuples into the *EmployeeNames* view. Show the appropriate command and describe what privileges your secretary has after this command is executed.
7. Your secretary allows Todd to read the *EmployeeNames* relation and later quits. You then revoke the secretary's privileges. What happens to Todd's privileges?
8. Give an example of a view update on the preceding schema that cannot be implemented through updates to *Employees*.
9. You decide to go on an extended vacation, and to make sure that emergencies can be handled, you want to authorize your boss Joe to read and modify the *Employees* relation and the *EmployeeNames* relation (and Joe must be able to delegate authority, of course, since he is too far up the management hierarchy to actually do any work). Show the appropriate SQL statements. Can Joe read the *DeptInfo* view?
10. After returning from your (wonderful) vacation, you see a note from Joe, indicating that he authorized his secretary Mike to read the *Employees* relation. You want to revoke Mike's SELECT privilege on *Employees*, but you do not want to revoke the rights you gave to Joe, even temporarily. Can you do this in SQL?
11. Later you realize that Joe has been quite busy. He has defined a view called *AllNames* using the view *EmployeeNames*, defined another relation called *StaffNames* that he has access to (but you cannot access), and given his secretary Mike the right to read from the *AllNames* view. Mike has passed this right on to his friend Susan. You decide that, even at the cost of annoying Joe by revoking some of his privileges, you simply have to take away Mike (and Susan's) rights to see your data. What REVOKE statement would you execute? What rights does Joe have on *Employees* after this statement is executed? What views are dropped as a consequence?



22

PARALLEL AND DISTRIBUTED DATABASES

- ☛ What is the motivation for parallel and distributed DBMSs?
- ☛ What are the alternative architectures for parallel database systems?
- ☛ How are pipelining and data partitioning used to gain parallelism?
- ☛ How are dataflow concepts used to parallelize existing sequential code?
- ☛ What are alternative architectures for distributed DBMSs?
- ☛ How is data distributed across sites?
- ☛ How can we evaluate and optimize queries over distributed data?
- ☛ What are the merits of synchronous vs. asynchronous replication?
- ☛ How are transactions managed in a distributed environment?
- Key concepts: parallel DBMS architectures; performance, speed-up and scale-up; pipelined versus data-partitioned parallelism, blocking; partitioning strategies; dataflow operators; distributed DBMS architectures; heterogeneous systems; gateway protocols; data distribution, distributed catalogs; serial joins, data shipping; synchronous versus asynchronous replication; distributed transactions, lock management, deadlock detection, two-phase commit, Presumed Abort

No man is an island, entire of itself; every man is a piece of the continent, a part of the main.

—John Donne

In this chapter we look at the issues of parallelism and data distribution in a DBMS. We begin by introducing parallel and distributed database systems in Section 22.1. In Section 22.2, we discuss alternative hardware configurations for a parallel DBMS. In Section 22.3, we introduce the concept of data partitioning and consider its influence on parallel query evaluation. In Section 22.4, we show how data partitioning can be used to parallelize several relational operations. In Section 22.5, we conclude our treatment of parallel query processing with a discussion of parallel query optimization.

The rest of the chapter is devoted to distributed databases. We present an overview of distributed databases in Section 22.6. We discuss some alternative architectures for a distributed DBMS in Section 22.7 and describe options for distributing data in Section 22.8. We describe distributed catalog management in Section 22.9, then in Section 22.10, we discuss query optimization and evaluation for distributed databases. In Section 22.11, we discuss updating distributed data, and finally, in Sections 22.12 to 22.14 we describe distributed transaction management.

22.1 INTRODUCTION

We have thus far considered centralized database management systems in which all the data is maintained at a single site and assumed that the processing of individual transactions is essentially sequential. One of the most important trends in databases is the increased use of parallel evaluation techniques and data distribution.

A parallel database system seeks to improve performance through parallelization of various operations, such as loading data, building indexes, and evaluating queries. Although data may be stored in a distributed fashion in such a system, the distribution is governed solely by performance considerations.

In a distributed database system, data is physically stored across several sites, and each site is typically managed by a DBMS capable of running independently of the other sites. The location of data items and the degree of autonomy of individual sites have a significant impact on all aspects of the system, including query optimization and processing, concurrency control, and recovery. In contrast to parallel databases, the distribution of data is governed by factors such as local ownership and increased availability, in addition to performance issues.

While parallelism is motivated by performance considerations, several distinct issues motivate data distribution:

- **Increased Availability:** If a site containing a relation goes down, the relation continues to be available if a copy is maintained at another site.
- **Distributed Access to Data:** An organization may have branches in several cities. Although analysts may need to access data corresponding to different sites, we usually find locality in the access patterns (e.g., a bank manager is likely to look up the accounts of customers at the local branch), and this locality can be exploited by distributing the data accordingly.
- **Analysis of Distributed Data:** Organizations want to examine all the data available to them, even when it is stored across multiple sites and on multiple database systems. Support for such integrated access involves many issues; even enabling access to widely distributed data can be a challenge.

22.2 ARCHITECTURES FOR PARALLEL DATABASES

The basic idea behind parallel databases is to carry out evaluation steps in parallel whenever possible, and there are many such opportunities in a relational DBMS; databases represent one of the most successful instances of parallel computing.

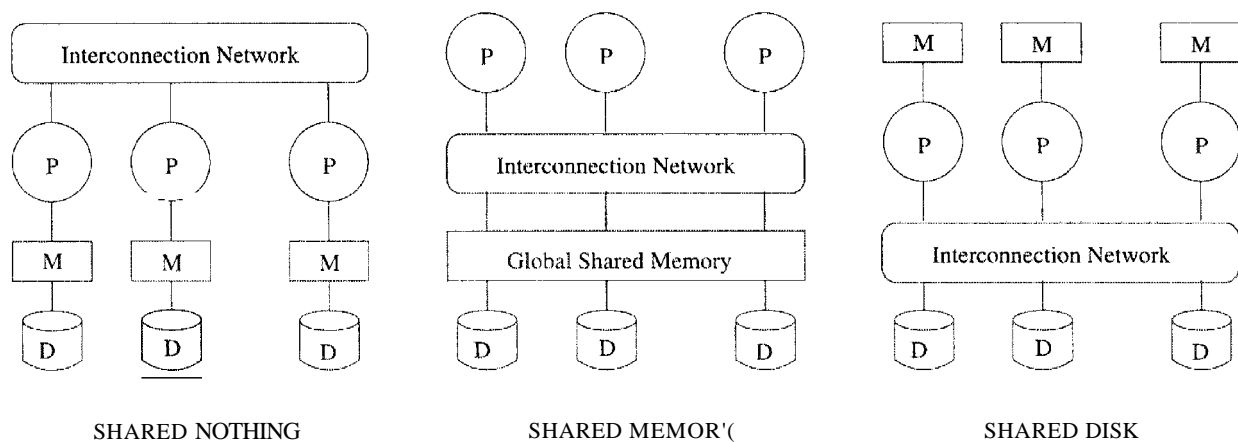


Figure 22.1 Physical Architectures for Parallel Database Systems

Three main architectures have been proposed for building parallel DBMSs. In a shared-memory system, multiple CPUs are attached to an interconnection network and can access a common region of main memory. In a shared-disk system, each CPU has a private memory and direct access to all disks through an interconnection network. In a shared-nothing system, each CPU has local main memory and disk space, but no two CPUs can access the same storage area; all communication between CPUs is through a network connection. The three architectures are illustrated in Figure 22.1.

The shared-memory architecture is closer to a conventional machine, and many commercial database systems have been ported to shared memory platforms with relative ease. Communication overhead is low, because main memory can be used for this purpose, and operating system services can be leveraged to utilize the additional CPUs. Although this approach is attractive for achieving moderate parallelism—a few tens of CPUs can be exploited in this fashion—memory contention becomes a bottleneck as the number of CPUs increases. The shared-disk architecture faces a similar problem because large amounts of data are shipped through the interconnection network.

The basic problem with the shared-memory and shared-disk architectures is interference: As more CPUs are added, existing CPUs are slowed down because of the increased contention for memory accesses and network bandwidth. It has been noted that even an average 1 percent slowdown per additional CPU means that the maximum speed-up is a factor of 37, and adding additional CPUs actually slows down the system; a system with 1000 CPUs is only 4 percent as effective as a *single-CPU* system! This observation has motivated the development of the shared-nothing architecture, which is now widely considered to be the best architecture for large parallel database systems.

The shared-nothing architecture requires more extensive reorganization of the DBMS code, but it has been shown to provide linear speed-up, in that the time taken for operations decreases in proportion to the increase in the number of CPUs and disks, and linear scale-up, in that performance is sustained if the number of CPUs and disks are increased in proportion to the amount of data. Consequently, ever-more-powerful parallel database systems can be built by taking advantage of rapidly improving performance for single-CPU systems and connecting as many CPUs as desired.

Speed-up and scale-up are illustrated in Figure 22.2. The speed-up curves show how, for a fixed database size, more transactions can be executed per second by adding CPUs. The scale-up curves show how adding more resources (in the form of CPUs) enables us to process larger problems. The first scale-up graph measures the number of transactions executed per second as the database size is increased and the number of CPUs is correspondingly increased. An alternative way to measure scale-up is to consider the time taken per transaction as more CPUs are added to process an increasing number of transactions per second; the goal here is to sustain the response time per transaction.

22.3 PARALLEL QUERY EVALUATION

In this section, we discuss parallel evaluation of a relational query in a DBMS with a shared-nothing architecture. While it is possible to consider parallel

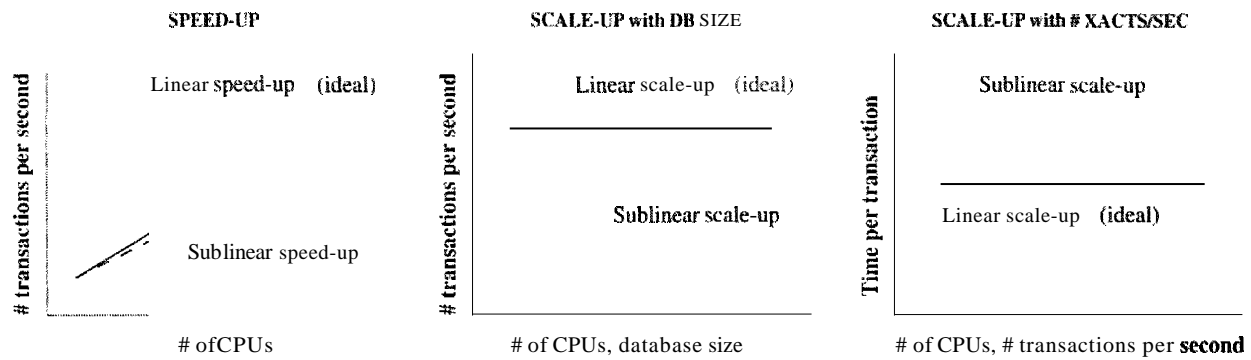


Figure 22.2 Speed-up and Scale-up

execution of multiple queries, it is hard to identify in advance which queries will run concurrently. So the emphasis has been on parallel execution of a single query.

A relational query execution plan is a graph of relational algebra operators, and the operators in a graph can be executed in parallel. If one operator consumes the output of a second operator, we have **pipelined parallelism** (the output of the second operator is worked on by the first operator as soon as it is generated); if not, the two operators can proceed essentially independently. An operator is said to **block** if it produces no output until it has consumed all its inputs. Pipelined parallelism is limited by the presence of operators (e.g., sorting or aggregation) that block.

In addition to evaluating different operators in parallel, we can evaluate each individual operator in a query plan in a parallel fashion. The key to evaluating an operator in parallel is to *partition* the input data; we can then work on each partition in parallel and combine the results. This approach is called **data-partitioned parallel evaluation**. By exercising some care, existing code for sequentially evaluating relational operators can be ported easily for data-partitioned parallel evaluation.

An important observation, which explains why shared-nothing parallel database systems have been very successful, is that database query evaluation is very amenable to data-partitioned parallel evaluation. The goal is to minimize data shipping by partitioning the data and structuring the algorithms to do most of the processing at individual processors. (We use *processor* to refer to a CPU together with its local disk.)

We now consider data partitioning and parallelization of existing operator evaluation code in more detail.

22.3.1 Data Partitioning

Partitioning a large dataset horizontally across several disks enables us to exploit the I/O bandwidth of the disks by reading and writing them in parallel. There are several ways to horizontally partition a relation. We can assign tuples to processors in a round-robin fashion, we can use hashing, or we can assign tuples to processors by ranges of field values. If there are n processors, the i th tuple is assigned to processor $i \bmod n$ in round-robin partitioning. Recall that round-robin partitioning is used in RAID storage systems (see Section 9.2). In hash partitioning, a hash function is applied to (selected fields of) a tuple to determine its processor. In range partitioning, tuples are sorted (conceptually), and n ranges are chosen for the sort key values so that each range contains roughly the same number of tuples; tuples in range i are assigned to processor i .

Round-robin partitioning is suitable for efficiently evaluating queries that access the entire relation. If only a subset of the tuples (e.g., those that satisfy the selection condition $age = 20$) is required, hash partitioning and range partitioning are better than round-robin partitioning because they enable us to access only those disks that contain matching tuples. (Of course, this statement assumes that the tuples are partitioned on the attributes in the selection condition; if $age = 20$ is specified, the tuples must be partitioned on age .) If range selections such as $15 < age < 25$ are specified, range partitioning is superior to hash partitioning because qualifying tuples are likely to be clustered together on a few processors. On the other hand, range partitioning can lead to data skew; that is, partitions with widely varying numbers of tuples across partitions or disks. Skew causes processors dealing with large partitions to become performance bottlenecks. Hash partitioning has the additional virtue that it keeps data evenly distributed even if the data grows and shrinks over time.

To reduce skew in range partitioning, the main question is how to choose the ranges by which tuples are distributed. One effective approach is to take samples from each processor, collect and sort all samples, and divide the sorted set of samples into equally sized subsets. If tuples are to be partitioned on age , the age ranges of the sampled subsets of tuples can be used as the basis for redistributing the entire relation.

22.3.2 Parallelizing Sequential Operator Evaluation Code

An elegant software architecture for parallel DBMSs enables us to readily parallelize existing code for sequentially evaluating a relational operator. The basic idea is to use parallel data streams. Streams (from different disks or

the output of other operators) are merged as needed to provide the inputs for a relational operator, and the output of an operator is split as needed to parallelize subsequent processing.

A parallel evaluation plan consists of a dataflow network of relational, merge, and split operators. The merge and split operators should be able to buffer some data and should be able to halt the operators producing their input data. They can then regulate the speed of the execution according to the execution speed of the operator that consumes their output.

As we will see, obtaining good parallel versions of algorithms for sequential operator evaluation requires careful consideration; there is no magic formula for taking sequential code and producing a parallel version. Good use of split and merge in a dataflow software architecture, however, can greatly reduce the effort of implementing parallel query evaluation algorithms, as we illustrate in Section 22.4.3.

22.4 PARALLELIZING INDIVIDUAL OPERATIONS

This section shows how various operations can be implemented in parallel in a shared-nothing architecture. We assume that each relation is horizontally partitioned across several disks, although this partitioning may or may not be appropriate for a given query. The evaluation of a query must take the initial partitioning criteria into account and repartition if necessary.

22.4.1 Bulk Loading and Scanning

We begin with two simple operations: *scanning* a relation and *loading* a relation. Pages can be read in parallel while scanning a relation, and the retrieved tuples can then be merged, if the relation is partitioned across several disks. More generally, the idea also applies when retrieving all tuples that meet a selection condition. If hashing or range partitioning is used, selection queries can be answered by going to just those processors that contain relevant tuples.

A similar observation holds for bulk loading. Further, if a relation has associated indexes, any sorting of data entries required for building the indexes during bulk loading can also be done in parallel (see later).

22.4.2 Sorting

A simple idea is to let each CPTJ sort the part of the relation that is on its local disk and then merge these sorted sets of tuples. The degree of parallelism is likely to be limited by the merging phase.

A better idea is to first redistribute all tuples in the relation using range partitioning. For example, if we want to sort a collection of employee tuples by salary, salary values range from 10 to 210, and we have 20 processors, we could send all tuples with salary values in the range 10 to 20 to the first processor, all in the range 21 to 30 to the second processor, and so on. (Prior to the redistribution, while tuples are distributed across the processors, we cannot assume that they are distributed according to salary ranges.)

Each processor then sorts the tuples assigned to it, using some sequential sorting algorithm. For example, a processor can collect tuples until its memory is full, then sort these tuples and write out a run, until all incoming tuples have been written to such sorted runs on the local disk. These runs can then be merged to create the sorted version of the set of tuples assigned to this processor. The entire sorted relation can be retrieved by visiting the processors in an order corresponding to the ranges assigned to them and simply scanning the tuples.

The basic challenge in parallel sorting is to do the range partitioning so that each processor receives roughly the same number of tuples; otherwise, a processor that receives a disproportionately large number of tuples to sort becomes a bottleneck and limits the scalability of the parallel sort. One good approach to range partitioning is to obtain a sample of the entire relation by taking samples at each processor that initially contains part of the relation. The (relatively small) sample is sorted and used to identify ranges with equal numbers of tuples. This set of range values, called a splitting vector, is then distributed to all processors and used to range partition the entire relation.

A particularly important application of parallel sorting is sorting the data entries in tree-structured indexes. Sorting data entries can significantly speed up the process of bulk-loading an index.

22.4.3 Joins

In this section, we consider how the join operation can be parallelized. We present the basic idea behind the parallelization and illustrate the use of the merge and split operators described in Section 22.3.2. We focus on parallel hash join, which is widely used, and briefly outline how sort-merge join can

be similarly parallelized. Other join algorithms can be parallelized as well, although not as effectively as these two algorithms.

Suppose that we want to join two relations, say, A and B , on the *age* attribute. We assume that they are initially distributed across several disks in some way that is not useful for the join operation; that is, the initial partitioning is not based on the join attribute. The basic idea for joining A and B in parallel is to decompose the join into a collection of k smaller joins. We can decompose the join by partitioning both A and B into a collection of k logical buckets or partitions. By using the same partitioning function for both A and B , we ensure that the union of the k smaller joins computes the join of A and B ; this idea is similar to intuition behind the partitioning phase of a sequential hash join, described in Section 14.4.3. Because A and B are initially distributed across several processors, the partitioning step itself can be done in parallel at these processors. At each processor, all local tuples are retrieved and hashed into one of k partitions, with the same hash function used at all sites, of course.

Alternatively, we can partition A and B by dividing the range of the join attribute *age* into k disjoint subranges and placing A and B tuples into partitions according to the subrange to which their *age* values belong. For example, suppose that we have 10 processors, the join attribute is *age*, with values from 0 to 100. Assuming uniform distribution, A and B tuples with $0 \leq \text{age} < 10$ go to processor 1, $10 \leq \text{age} < 20$ go to processor 2, and so on. This approach is likely to be more susceptible than hash partitioning to data skew (i.e., the number of tuples to be joined can vary widely across partitions), unless the subranges are carefully determined; we do not discuss how good subrange boundaries can be identified.

Having decided on a partitioning strategy, we can assign each partition to a processor and carry out a local join, using any join algorithm we want, at each processor. In this case, the number of partitions k is chosen to be equal to the number of processors n available for carrying out the join, and during partitioning, each processor sends tuples in the i th partition to processor i . After partitioning, each processor joins the A and B tuples assigned to it. Each join process executes sequential join code and receives input A and B tuples from several processors; a merge operator merges all incoming A tuples, and another merge operator merges all incoming B tuples. Depending on how we want to distribute the result of the join of A and B , the output of the join process may be split into several data streams. The network of operators for parallel join is shown in Figure 22.3. To simplify the figure, we assume that the processors doing the join are distinct from the processors that initially contain tuples of A and B and show only four processors.

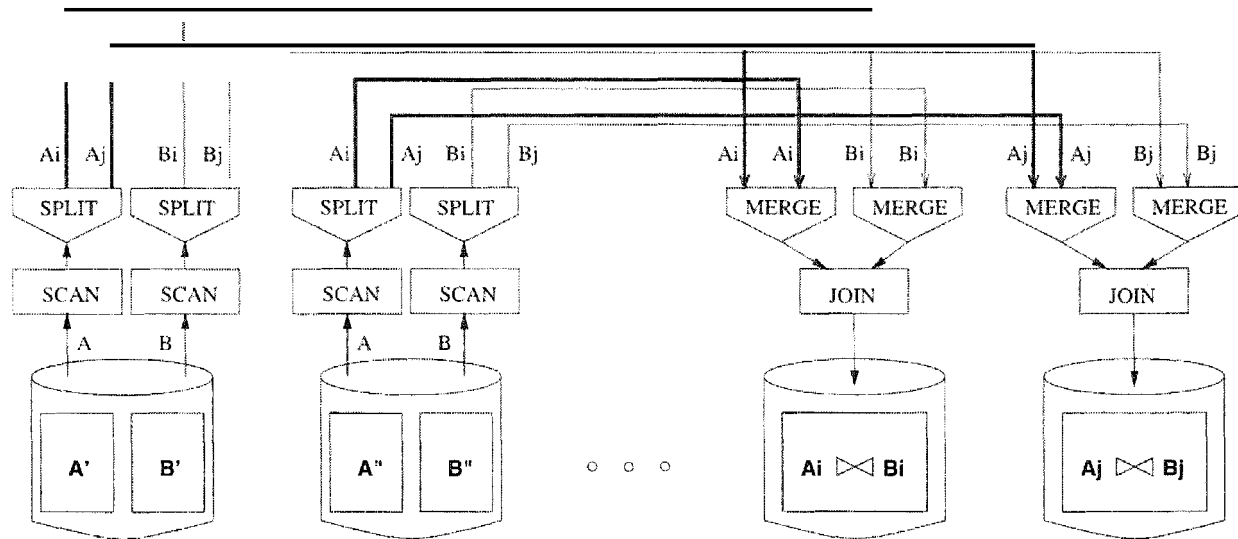


Figure 22.3 Dataflow Network of Operators for Parallel Join

If range partitioning is used, this algorithm leads to a parallel version of a sort-merge join, with the advantage that the output is available in sorted order. If hash partitioning is used, we obtain a parallel version of a hash join.

Improved Parallel Hash Join

A hash-based refinement of the approach offers improved performance. The main observation is that, if A and B are very large and the number of partitions k is chosen to be equal to the number of processors n , the size of each partition may still be large, leading to a high cost for each local join at the n processors.

An alternative is to execute the smaller joins $A_i \bowtie B_i$, for $i = 1 \dots k$, one after the other, but with each join executed in parallel using all processors. This approach allows us to utilize the total available main memory at all n processors in each join $A_i \bowtie B_i$ and is described in more detail as follows:

1. At each site, apply a hash function h_1 to partition the A and B tuples at this site into partitions $i = 1 \dots k$. Let A be the smaller relation. The number of partitions k is chosen such that each partition of A fits into the aggregate or combined memory of all n processors.
2. For $i = 1 \dots k$, process the join of the i th partitions of A and B . To compute $A_i \bowtie B_i$, do the following at every site:
 - (a) Apply a second hash function h_2 to all A_i tuples to determine where they should be joined and send tuple t to site $h_2(t)$.
 - (b) As A_i tuples arrive to be joined, add them to an in-memory hash table.

- (c) After all A_i tuples have been distributed, apply h_2 to B_i tuples to determine where they should be joined and send tuple t to site $h_2(t)$.
- (d) As B_i tuples arrive to be joined, probe the in-memory table of A_i tuples and output result tuples.

The use of the second hash function h_2 ensures that tuples are (more or less) uniformly distributed across all n processors participating in the join. This approach greatly reduces the cost for each of the smaller joins and therefore reduces the overall join cost. Observe that all available processors are fully utilized, even though the smaller joins are carried out one after the other.

The reader is invited to adapt the network of operators shown in Figure 22.3 to reflect the improved parallel join algorithm.

22.5 PARALLEL QUERY OPTIMIZATION

In addition to parallelizing individual operations, we can obviously execute different operations in a query in parallel and execute multiple queries in parallel. Optimizing a single query for parallel execution has received more attention; systems typically optimize queries without regard to other queries that might be executing at the same time.

Two kinds of interoperation parallelism can be exploited within a query:

- The result of one operator can be pipelined into another. For example, consider a left-deep plan in which all the joins use index nested loops. The result of the first (i.e., the bottommost) join is the outer relation tuples for the next join node. As tuples are produced by the first join, they can be used to probe the inner relation in the second join. The result of the second join can similarly be pipelined into the next join, and so on.
- Multiple independent operations can be executed concurrently. For example, consider a (bushy) plan in which relations A and B are joined, relations C and D are joined, and the results of these two joins are finally joined. Clearly, the join of A and B can be executed concurrently with the join of C and D .

An optimizer that seeks to parallelize query evaluation has to consider several issues, and we only outline the main points. The cost of executing individual operations in parallel (e.g., parallel sorting) obviously differs from executing them sequentially, and the optimizer should estimate operation costs accordingly.

Next, the plan that returns answers quickest may not be the plan with the least cost. For example, the cost of $A \bowtie B$ plus the cost of $C \bowtie D$ plus the cost of joining their results may be more than the cost of the cheapest left-deep plan. However, the time taken is the time for the more expensive of $A \bowtie B$ and $C \bowtie D$ plus the time to join their results. This time may be less than the time taken by the cheapest left-deep plan. This observation suggests that a parallelizing optimizer should not restrict itself to left-deep trees and should also consider *bushy* trees, which significantly enlarge the space of plans to be considered.

Finally, a number of parameters, such as available buffer space and the number of free processors, are known only at run-time. This comment holds in a multiuser environment even if only sequential plans are considered; a multiuser environment is a simple instance of interquery parallelism.

22.6 INTRODUCTION TO DISTRIBUTED DATABASES

As we observed earlier, data in a distributed database system is stored across several sites, and each site is typically managed by a DBMS that can run independent of the other sites. The classical view of a distributed database system is that the system should make the impact of data distribution transparent. In particular, the following properties are considered desirable:

- **Distributed Data Independence:** Users should be able to ask queries without specifying where the referenced relations, or copies or fragments of the relations, are located. This principle is a natural extension of physical and logical data independence; we discuss it in Section 22.8. Further, queries that span multiple sites should be optimized systematically in a cost-based manner, taking into account communication costs and differences in local computation costs. We discuss distributed query optimization in Section 22.10.
- **Distributed Transaction Atomicity:** Users should be able to write transactions that access and update data at several sites just as they would write transactions over purely local data. In particular, the effects of a transaction across sites should continue to be atomic; that is, all changes persist if the transaction commits and none persist if it aborts. We discuss this distributed transaction processing in Sections 22.11, 22.13, and 22.14.

Although most people would agree that these properties are in general desirable, in certain situations, such as when sites are connected by a slow long-distance network, these properties are not efficiently achievable. Indeed, it has been argued that when sites are globally distributed, these properties are not even desirable. The argument essentially is that the administrative overhead

of supporting a system with distributed data independence and transaction atomicity—in effect, coordinating all activities across all sites to support the view of the whole as a unified collection of data—is prohibitive, over and above DBMS performance considerations.

Keep these remarks about distributed databases in mind as we cover the topic in more detail in the rest of this chapter. There is no real consensus on what the design objectives of distributed databases should be, and the field is evolving in response to users' needs.

22.6.1 Types of Distributed Databases

If data is distributed but all servers run the same DBMS software, we have a homogeneous distributed database system. If different sites run under the control of different DBMSs, essentially autonomously, and are connected somehow to enable access to data from multiple sites, we have a heterogeneous distributed database system, also referred to as a multidatabase system.

The key to building heterogeneous systems is to have well-accepted standards for gateway protocols. A gateway protocol is an API that exposes DBMS functionality to external applications. Examples include ODBC and JDBC (see Section 6.2). By accessing database servers through gateway protocols, their differences (in capability, data format, etc.) are masked, and the differences between the different servers in a distributed system are bridged to a large degree.

Gateways are not a panacea, however. They add a layer of processing that can be expensive, and they do not completely mask the differences among servers. For example, a server may not be capable of providing the services required for distributed transaction management (see Sections 22.13 and 22.14), and even if it is capable, standardizing gateway protocols all the way down to this level of interaction poses challenges that have not yet been resolved satisfactorily.

Distributed data management, in the final analysis, comes at a significant cost in terms of performance, software complexity, and administration difficulty. This observation is especially true of heterogeneous systems.

22.7 DISTRIBUTED DBMS ARCHITECTURES

Three alternative approaches are used to separate functionality across different DBMS-related processes; these alternative distributed DBMS architectures are called *Client-Server*, *Collaborating Server*, and *Middleware*.

22.7.1 Client-Server Systems

A Client-Server system has one or more client processes and one or more server processes, and a client process can send a query to any server process. Clients are responsible for user-interface issues, and servers manage data and execute transactions. Thus, a client process could run on a personal computer and send queries to a server running on a mainframe.

This architecture has become very popular for several reasons. First, it is relatively simple to implement due to its clean separation of functionality and because the server is centralized. Second, expensive server machines are not underutilized by dealing with mundane user-interactions, which are now relegated to inexpensive client machines. Third, users can run a graphical user interface that they are familiar with, rather than the (possibly unfamiliar and unfriendly) user interface on the server.

While writing Client-Server applications, it is important to remember the boundary between the client and the server and keep the communication between them as set-oriented as possible. In particular, opening a cursor and fetching tuples one at a time generates many messages and should be avoided. (Even if we fetch several tuples and cache them at the client, messages must be exchanged when the cursor is advanced to ensure that the current row is locked.) Techniques to exploit client-side caching to reduce communication overhead have been studied extensively, although we do not discuss them further.

22.7.2 Collaborating Server Systems

The Client-Server architecture does not allow a single query to span multiple servers because the client process would have to be capable of breaking such a query into appropriate subqueries to be executed at different sites and then piecing together the answers to the subqueries. The client process would therefore be quite complex, and its capabilities would begin to overlap with the server; distinguishing between clients and servers becomes harder. Eliminating this distinction leads us to an alternative to the Client-Server architecture: a Collaborating Server system. We can have a collection of database servers, each capable of running transactions against local data, which cooperatively execute transactions spanning multiple servers.

When a server receives a query that requires access to data at other servers, it generates appropriate subqueries to be executed by other servers and puts the results together to compute answers to the original query. Ideally, the decom-

position of the query should be done using cost-based optimization, taking into account the cost of network communication as well as local processing costs.

22.7.3 Middleware Systems

The Middleware architecture is designed to allow a single query to span multiple servers, without requiring all database servers to be capable of managing such multi-site execution strategies. It is especially attractive when trying to integrate several legacy systems, whose basic capabilities cannot be extended.

The idea is that we need just one database server capable of managing queries and transactions spanning multiple servers; the remaining servers need to handle only local queries and transactions. We can think of this special server as a layer of software that coordinates the execution of queries and transactions across one or more independent database servers; such software is often called **middleware**. The middleware layer is capable of executing joins and other relational operations on data obtained from the other servers but, typically, does not itself maintain any data.

22.8 STORING DATA IN A DISTRIBUTED DBMS

In a distributed DBMS, relations are stored across several sites. Accessing a relation stored at a remote site incurs message-passing costs and, to reduce this overhead, a single relation may be *partitioned* or *fragmented* across several sites, with fragments stored at the sites where they are most often accessed or *replicated* at each site where the relation is in high demand.

22.8.1 Fragmentation

Fragmentation consists of breaking a relation into smaller relations or fragments and storing the fragments (instead of the relation itself), possibly at different sites. In horizontal fragmentation, each fragment consists of a subset of *rows* of the original relation. In vertical fragmentation, each fragment consists of a subset of *columns* of the original relation. Horizontal and vertical fragments are illustrated in Figure 22.4.

Typically, the tuples that belong to a given horizontal fragment are identified by a selection query; for example, employee tuples might be organized into fragments by city, with all employees in a given city assigned to the same fragment. The horizontal fragment shown in Figure 22.4 corresponds to Chicago. By storing fragments in the database site at the corresponding city, we achieve locality of reference—Chicago data is most likely to be updated and queried

TID	eid	name	city	age	sal
t1	53666	Jones	Madras	18	35
t2	53688	Snlith	Chicago	18	32
t3	53650	Sluith	Chicago	19	48
t4	53831	Madayan	Bombay	11	20
t5	53832	Guldu	BOlnbay	12	20

Vertical Fragment
Horizontal Fragment

Figure 22.4 Horizontal and Vertical Fragmentation

from Chicago, and storing this data in Chicago makes it local (and reduces communication costs) for most queries. Similarly, the tuples in a given vertical fragment are identified by a projection query. The vertical fragment in the figure results from projection on the first two columns of the employees relation.

When a relation is fragmented, we must be able to recover the original relation from the fragments:

- **Horizontal Fragmentation:** The union of the horizontal fragments must be equal to the original relation. Fragments are usually also required to be disjoint.
- **Vertical Fragmentation:** The collection of vertical fragments should be a lossless-join decomposition, as per the definition in Chapter 19.

To ensure that a vertical fragmentation is lossless-join, systems often assign a unique tuple id to each tuple in the original relation, as shown in Figure 22.4, and attach this id to the projection of the tuple in each fragment. If we think of the original relation as containing an additional tuple-id field that is a key, this field is added to each vertical fragment. Such a decomposition is guaranteed to be lossless-join.

In general, a relation can be (horizontally or vertically) fragmented, and each resulting fragment can be further fragmented. For simplicity of exposition, in the rest of this chapter, we assume that fragments are not recursively partitioned in this manner.

22.8.2 Replication

Replication means that we store several copies of a relation or relation fragment. An entire relation can be replicated at one or more sites. Similarly, one or more fragments of a relation can be replicated at other sites. For example, if a relation R is fragmented into R_1 , R_2 , and R_3 , there might be just one copy of R_1 , whereas R_2 is replicated at two other sites and R_3 is replicated at all sites.

The motivation for replication is twofold:

- **Increased Availability of Data:** If a site that contains a replica goes down, we can find the same data at other sites. Similarly, if local copies of remote relations are available, we are less vulnerable to failure of communication links.
- **Faster Query Evaluation:** Queries can execute faster by using a local copy of a relation instead of going to a remote site.

The two kinds of replication, called *synchronous* and *asynchronous* replication, differ primarily in how replicas are kept current when the relation is modified (see Section 22.11).

22.9 DISTRIBUTED CATALOG MANAGEMENT

Keeping track of data distributed across several sites can get complicated. We must keep track of how relations are fragmented and replicated—that is, how relation fragments are distributed across several sites and where copies of fragments are stored—in addition to the usual schema, authorization, and statistical information.

22.9.1 Naming Objects

If a relation is fragmented and replicated, we must be able to uniquely identify each replica of each fragment. Generating such unique names requires some care. If we use a global name-server to assign globally unique names, local autonomy is compromised; we want (users at) each site to be able to assign names to local objects without reference to names systemwide.

The usual solution to the naming problem is to use names consisting of several fields. For example, we could have:

- A *local name* field, which is the name assigned locally at the site where the relation is created. Two objects at different sites could have the same local name, but two objects at a given site cannot have the same local name.
- A *birth site* field, which identifies the site where the relation was created, and where information is maintained about all fragments and replicas of the relation.

These two fields identify a relation uniquely; we call the combination a global relation name. To identify a replica (of a relation or a relation fragment), we take the global relation name and add a *replica-id* field; we call the combination a global replica name.

22.9.2 Catalog Structure

A centralized system catalog can be used but is vulnerable to failure of the site containing the catalog. An alternative is to maintain a copy of a global system catalog, which describes all the data at every site. Although this approach is not vulnerable to a single-site failure, it compromises site autonomy, just like the first solution, because every change to a local catalog must now be broadcast to all sites.

A better approach, which preserves local autonomy and is not vulnerable to a single-site failure, was developed in the R^* distributed database project, which was a successor to the System R project at IBM. Each site maintains a local catalog that describes all copies of data stored at that site. In addition, the catalog at the birth site for a relation is responsible for keeping track of where replicas of the relation (in general, of fragments of the relation) are stored. In particular, a precise description of each replica's contents—a list of columns for a vertical fragment or a selection condition for a horizontal fragment—is stored in the birth site catalog. Whenever a new replica is created or a replica is moved across sites, the information in the birth site catalog for the relation must be updated.

To locate a relation, the catalog at its birth site must be looked up. This catalog information can be cached at other sites for quicker access, but the cached information may become out of date if, for example, a fragment is moved. We would discover that the locally cached information is out of date when we use it to access the relation, and at that point, we must update the cache by looking up the catalog at the birth site of the relation. (The birth site of a relation is recorded in each local cache that describes the relation, and the birth site never changes, even if the relation is moved.)

22.9.3 Distributed Data Independence

Distributed data independence means that users should be able to write queries without regard to how a relation is fragmented or replicated; it is the responsibility of the DBMS to compute the relation as needed (by locating suitable copies of fragments, joining the vertical fragments, and taking the union of horizontal fragments).

In particular, this property implies that users should not have to specify the full name for the data objects accessed while evaluating a query. Let us see how users can be enabled to access relations without considering how the relations are distributed. The *local name* of a relation in the system catalog (Section 22.9.1) is really a combination of a *user name* and a user-defined *relation name*. Users can give whatever names they wish to their relations, without regard to the relations created by other users. When a user writes a program or SQL statement that refers to a relation, he or she simply uses the relation name. The DBMS adds the user name to the relation name to get a local name, then adds the user's site-id as the (default) birth site to obtain a global relation name. By looking up the global relation name in the local catalog if it is cached there or in the catalog at the birth site, the DBMS can locate replicas of the relation.

A user may want to create objects at several sites or refer to relations created by other users. To do this, a user can create a synonym for a global relation name using an SQL-style command (although such a command is not currently part of the SQL:1999 standard) and subsequently refer to the relation using the synonym. For each user known at a site, the DBMS maintains a table of synonyms as part of the system catalog at that site and uses this table to find the global relation name. Note that a user's program runs unchanged even if replicas of the relation are moved, because the global relation name is never changed until the relation itself is destroyed.

Users may want to run queries against specific replicas, especially if asynchronous replication is used. To support this, the synonym mechanism can be adapted to also allow users to create synonyms for global replica names.

22.10 DISTRIBUTED QUERY PROCESSING

We first discuss the issues involved in evaluating relational algebra operations in a distributed database through examples and then outline distributed query optimization. Consider the following two relations:

Sailors(sid: integer, sname: string, rating: integer, age: real)
Reserves(sid: integer, bid: integer, day: date, rname: string)

As in Chapter 14, assume that each tuple of *Reserves* is 40 bytes long, that a page can hold 100 *Reserves* tuples, and that we have 1000 pages of such tuples. Similarly, assume that each tuple of *Sailors* is 50 bytes long, that a page can hold 80 *Sailors* tuples, and that we have 500 pages of such tuples.

To estimate the cost of an evaluation strategy, in addition to counting the number of page I/Os, we must count the number of pages sent from one site to another because communication costs are a significant component of overall cost in a distributed database. We must also change our cost model to count the cost of shipping the result tuples to the site where the query is posed from the site where the result is assembled! In this chapter, we denote the time taken to read one page from disk (or to write one page to disk) as t_d and the time taken to ship one page (from any site to another site) as t_s .

22.10.1 Nonjoin Queries in a Distributed DBMS

Even simple operations such as scanning a relation, selection, and projection are affected by fragmentation and replication. Consider the following query:

```
SELECT S.age
FROM   Sailors S
WHERE  S.rating > 3 AND S.rating < 7
```

Suppose that the *Sailors* relation is horizontally fragmented, with all tuples having a rating less than 5 at Shanghai and all tuples having a rating greater than 5 at Tokyo.

The DBMS must answer this query by evaluating it at both sites and taking the union of the answers. If the *SELECT* clause contained *AVG (S.age)*, combining the answers could not be done by simply taking the union—the DBMS must compute the sum and count of *age* values at the two sites and use this information to compute the average age of all sailors.

If the *WHERE* clause contained just the condition *S.rating > 6*, on the other hand, the DBMS should recognize that this query could be answered by just executing it at Tokyo.

As another example, suppose that the *Sailors* relation were vertically fragmented, with the *sid* and *rating* fields at Shanghai and the *sname* and *age* fields at Tokyo. No field is stored at both sites. This vertical fragmentation

would therefore be a lossy decomposition, except that a field containing the id of the corresponding Sailors tuple is included by the DBMS in both fragments! Now, the DBMS has to reconstruct the Sailors relation by joining the two fragments on the common tuple-id field and execute the query over this reconstructed relation.

Finally, suppose that the entire Sailors relation were stored at both Shanghai and Tokyo. We could answer any of the previous queries by executing it at either Shanghai or Tokyo. Where should the query be executed? This depends on the cost of shipping the answer to the query site (which may be Shanghai, Tokyo, or some other site) as well as the cost of executing the query at Shanghai and at Tokyo—the local processing costs may differ depending on what indexes are available on Sailors at the two sites, for example.

22.10.2 Joins in a Distributed DBMS

Joins of relations at different sites can be very expensive, and we now consider the evaluation options that must be considered in a distributed environment. Suppose that the Sailors relation were stored at London, and the Reserves relation were stored at Paris. We consider the cost of various strategies for computing $Sailor'S \bowtie Reserves$.

Fetch As Needed

We could do a page-oriented nested loops join in London with Sailors as the outer, and for each Sailors page, fetch all Reserves pages from Paris. If we cache the fetched Reserves pages in London until the join is complete, pages are fetched only once, but assume that Reserves pages are not cached, just to see how bad things can get. (The situation can get much worse if we use a tuple-oriented nested loops join!)

The cost is $500t_d$ to scan Sailors plus, for each Sailors page, the cost of scanning and shipping all of Reserves, which is $1000(td + ts)$. The total cost is therefore $500td + 500,000(td + ts)$.

In addition, if the query was not executed at the London site, we must add the cost of shipping the result to the query site; this cost depends on the size of the result. Because *sid* is a key for Sailors, the number of tuples in the result is 100,000 (the number of tuples in Reserves) and each tuple is $40 + 50 = 90$ bytes long; thus $4000/90 = 44$ result tuples fit on a page, and the result size is $100,000/44 = 2273$ pages. The cost of shipping the answer to another site, if necessary, is $2273 t_s$. In the rest of this section, we assume that the query is

posed at the site where the result is computed; if not, the cost of shipping the result to the query site must be added to the cost.

In this example, observe that, if the query site is not London or Paris, the cost of shipping the result is greater than the cost of shipping both *Sailors* and *Reserves* to the query site! Therefore, it would be cheaper to ship both relations to the query site and compute the join there.

Alternatively, we could do an index nested loops join in London, fetching all matching *Reserves* tuples for each *Sailors* tuple. Suppose we have an unclustered hash index on the *sid* column of *Reserves*. Because there are 100,000 *Reserves* tuples and 40,000 *Sailors* tuples, each sailor has on average 2.5 reservations. The cost of finding the 2.5 *Reservations* tuples that match a given *Sailors* tuple is $(1.2 + 2.5)t_d$, assuming 1.2 I/Os to locate the appropriate bucket in the index. The total cost is the cost of scanning *Sailors* plus the cost of finding and fetching matching *Reserves* tuples for each *Sailors* tuple, $500t_d + 40,000(3.7t_d + 2.5t_s)$.

Both algorithms fetch required *Reserves* tuples from a remote site as needed. Clearly, this is not a good idea; the cost of shipping tuples dominates the total cost even for a fast network.

Ship to One Site

We can ship *Sailors* from London to Paris and carry out the join there, ship *Reserves* to London and carry out the join there, or ship both to the site where the query was posed and compute the join there. Note again that the query could have been posed in London, Paris, or perhaps a third site, say, Timbuktu!

The cost of scanning and shipping *Sailors*, saving it at Paris, then doing the join at Paris is $500(2t_d + t_s) + 4500t_d$, assuming that the version of the sort-merge join described in Section 14.10 is used and we have an adequate number of buffer pages. In the rest of this section we assume that sort-merge join is the join method used when both relations are at the same site.

The cost of shipping *Reserves* and doing the join at London is $1000(2t(1 + t_s) + 4500t_d)$.

Selections and Bloom joins

Consider the strategy of shipping *Reserves* to London and computing the join at London. Some tuples in (the current instance of) *Reserves* do not join with any tuple in (the current instance of) *Sailors*. If we could somehow identify

Reserves tuples that are guaranteed not to join with any Sailors tuples, we could avoid shipping them.

Two techniques, *Semijoin* and *Bloomjoin*, have been proposed for reducing the number of Reserves tuples to be shipped. The first technique is called *Semijoin*. The idea is to proceed in three steps:

1. At London, compute the projection of Sailors onto the join columns (in this case just the *sid* field) and ship this projection to Paris.
2. At Paris, compute the natural join of the projection received from the first site with the Reserves relation. The result of this join is called the *reduction* of Reserves with respect to Sailors. Clearly, only those Reserves tuples in the reduction will join with tuples in the Sailors relation. Therefore, ship the reduction of Reserves to London, rather than the entire Reserves relation.
3. At London, compute the join of the reduction of Reserves with Sailors.

Let us compute the cost of using this technique for our example join query. Suppose we have a straightforward implementation of projection based on first scanning Sailors and creating a temporary relation with tuples that have only an *sid* field, then sorting the temporary and scanning the sorted temporary to eliminate duplicates. If we assume that the size of the *sid* field is 10 bytes, the cost of projection is $500t_d$ for scanning Sailors, plus $100t_d$ for creating the temporary, plus $400t_d$ for sorting it (in two passes), plus $100t_d$ for the final scan, plus $100t_d$ for writing the result into another temporary relation; a total of $1200t_d$. (Because *sid* is a key, no duplicates need be eliminated; if the optimizer is good enough to recognize this, the cost of projection is just $(500 + 100)t_d$.)

The cost of computing the projection and shipping it to Paris is therefore $1200t_d + 100t_s$. The cost of computing the reduction of Reserves is $3 \cdot (100 + 10(0)) = 3300t_d$, assuming that sort-merge join is used. (The cost does not reflect that the projection of Sailors is already sorted; the cost would decrease slightly if the refined sort-merge join exploited this.)

What is the size of the reduction? If every sailor holds at least one reservation, the reduction includes every tuple of Reserves! The effort invested in shipping the projection and reducing Reserves is a total waste. Indeed, because of this observation, we note that *Semijoin* is especially useful in conjunction with a selection on one of the relations. For example, if we want to compute the join of Sailors tuples with a *rating* greater than 8 with the Reserves relation, the size of the projection on *sid* for tuples that satisfy the selection would be just 20 percent of the original projection, that is, 20 pages.

Let us now continue the example join, with the assumption that we have the additional selection on *rating*. (The cost of computing the projection of Sailors goes down a bit, the cost of shipping it goes down to $20t_s$, and the cost of the reduction of Reserves also goes down a little, but we ignore these reductions for simplicity.) We assume that only 20 percent of the Reserves tuples are included in the reduction, thanks to the selection. Hence, the reduction contains 200 pages, and the cost of shipping it is $200t_s$.

Finally, at London, the reduction of Reserves is joined with Sailors, at a cost of $3 \cdot (200 + 500) = 2100t_d$. Observe that there are over 6500 page I/Os versus about 200 pages shipped, using this join technique. In contrast, to ship Reserves to London and do the join there costs $1000t_s$ plus $4500t_d$. With a high-speed network, the cost of Sernijoin may be more than the cost of shipping Reserves in its entirety, even though the shipping cost, itself is much less ($200t_s$ versus $1000t_s$).

The second technique, called Bloomjoin, is quite similar. The main difference is that a bit-vector is shipped in the first step, instead of the projection of Sailors. A bit-vector of (some chosen) size k is computed by hashing each tuple of Sailors into the range 0 to $k-1$ and setting bit i to 1 if some tuple hashes to i , and 0 otherwise. In the second step, the reduction of Reserves is computed by hashing each tuple of Reserves (using the *sid* field) into the range 0 to $k-1$, using the same hash function used to construct the bit-vector and discarding tuples whose hash value i corresponds to a 0 bit. Because no Sailors tuples hash to such an i , no Sailors tuple can join with any Reserves tuple that is not in the reduction.

The costs of shipping a bit-vector and reducing Reserves using the vector are less than the corresponding costs in Sernijoin. On the other hand, the size of the reduction of Reserves is likely to be larger than in Sernijoin; so, the costs of shipping the reduction and joining it with Sailors are likely to be higher.

Let us estimate the cost of this approach. The cost of computing the bit-vector is essentially the cost of scanning Sailors, which is $500t_d$. The cost of sending the bit-vector depends on the size we choose for the bit-vector, which is certainly smaller than the size of the projection; we take this cost to be $20t_s$, for concreteness. The cost of reducing Reserves is just the cost of scanning Reserves, $1000t_d$. The size of the reduction of Reserves is likely to be about the same as or a little larger than the size of the reduction in the Sernijoin approach; instead of 200, we will take this size to be 220 pages. (We assume that the selection on Sailors is included, to permit a direct comparison with the cost of Sernijoin.) The cost of shipping the reduction is therefore $220t_s$. The cost of the final join at London is $3 \cdot (500 + 220) = 2160t_d$.

Thus, in comparison to Semijoin, the shipping cost of this approach is about the same, although it could be higher if the bit-vector were not as selective as the projection of Sailors in terms of reducing Reserves. Typically, though, the reduction of Reserves is no more than 10 to 20 percent larger than the size of the reduction in SCJoin. In exchange for this slightly higher shipping cost, Bloornjoin achieves a significantly lower processing cost: less than $3700t_d$ versus more than $6500t_d$ for SCJoin. Indeed, Bloornjoin has a lower I/C cost and a lower shipping cost than the strategy of shipping all of Reserves to London! These numbers indicate why Bloornjoin is an attractive distributed join method; but the sensitivity of the method to the effectiveness of bit-vector hashing (in reducing Reserves) should be kept in mind.

22.10.3 Cost-Based Query Optimization

We have seen how data distribution can affect the implementation of individual operations, such as selection, projection, aggregation, and join. In general, of course, a query involves several operations, and optimizing queries in a distributed database poses the following additional challenges:

- Communication costs must be considered. If we have several copies of a relation, we must also decide which copy to use.
- If individual sites are run under the control of different DBMSs, the autonomy of each site must be respected while doing global query planning.

Query optimization proceeds essentially as in a centralized DBMS, as described in Chapter 12, with information about relations at remote sites obtained from the system catalogs. Of course, there are more alternative methods to consider for each operation (e.g., consider the new options for distributed joins), and the cost metric must account for communication costs as well, but the overall planning process is essentially unchanged if we take the cost metric to be the total cost of all operations. (If we consider response time, the fact that certain subqueries can be carried out in parallel at different sites would require us to change the optimizer as per the discussion in Section 22.5.)

In the overall plan, local manipulation of relations at the site where they are stored (to compute an intermediate relation to be shipped elsewhere) is encapsulated into a *suggested* local plan. The overall plan includes several such local plans, which we can think of as subqueries executing at different sites. While generating the global plan, the suggested local plans provide realistic cost estimates for the computation of the intermediate relations; the suggested local plans are constructed by the optimizer mainly to provide these local cost estimates. A site is free to ignore the local plan suggested to it if it is able to find a cheaper plan by using more current information in the local catalogs. Thus,

site autonomy is respected in the optimization and evaluation of distributed queries.

22.11 UPDATING DISTRIBUTED DATA

The classical view of a distributed DBMS is that it should behave just like a centralized DBMS from the point of view of a user; issues arising from distribution of data should be transparent to the user, although, of course, they must be addressed at the implementation level.

With respect to queries, this view of a distributed DBMS means that users should be able to ask queries without worrying about how and where relations are stored; we have already seen the implications of this requirement on query evaluation.

With respect to updates, this view means that transactions should continue to be atomic actions, regardless of data fragmentation and replication. In particular, all copies of a modified relation must be updated before the modifying transaction commits. We refer to replication with this semantics as synchronous replication; before an update transaction commits, it synchronizes all copies of modified data.

An alternative approach to replication, called asynchronous replication, has come to be widely used in commercial distributed DBMSs. Copies of a modified relation are updated only periodically in this approach, and a transaction that reads different copies of the same relation may see different values. Thus, asynchronous replication compromises distributed data independence, but it can be implemented more efficiently than synchronous replication.

22.11.1 Synchronous Replication

There are two basic techniques for ensuring that transactions see the same value regardless of which copy of an object they access. In the first technique, called voting, a transaction must write a majority of copies to modify an object and read at least enough copies to make sure that one of the copies is current. For example, if there are 10 copies and 7 copies are written by update transactions, then at least 4 copies must be read. Each copy has a version number, and the copy with the highest version number is current. This technique is not attractive in most situations because reading an object requires reading multiple copies; in most applications, objects are read much more frequently than they are updated, and efficient performance on reads is very important.

In the second technique, called read-any write-all, to read an object, a transaction can read anyone copy, but to write an object, it must write all copies. Reads are fast, especially if we have a local copy, but writes are slower, relative to the first technique. This technique is attractive when reads are much more frequent than writes, and it is usually adopted for implementing synchronous replication.

22.11.2 Asynchronous Replication

Synchronous replication comes at a significant cost. Before an update transaction can commit, it must obtain exclusive locks on all copies—assuming that the read-any write-all technique is used—of modified data. The transaction may have to send lock requests to remote sites and wait for the locks to be granted, and during this potentially long period, it continues to hold all its other locks. If sites or communication links fail, the transaction cannot commit until all the sites at which it has modified data recover and are reachable. Finally, even if locks are obtained readily and there are no failures, committing a transaction requires several additional messages to be sent as part of a *commit protocol* (Section 22.14.1).

For these reasons, synchronous replication is undesirable or even unachievable in many situations. Asynchronous replication is gaining in popularity, even though it allows different copies of the same object to have different values for short periods of time. This situation violates the principle of distributed data independence; users must be aware of which copy they are accessing, recognize that copies are brought up-to-date only periodically, and live with this reduced level of data consistency. Nonetheless, this seems to be a practical compromise that is acceptable in many situations.

Primary Site versus **Peer-to-Peer** Replication

Asynchronous replication comes in two flavors. In primary site asynchronous replication, one copy of a relation is designated the **primary** or **master** copy. Replicas of the entire relation or fragments of the relation can be created at other sites; these are **secondary** copies, and unlike the primary copy, they cannot be updated. A common technique for setting up primary and secondary copies is that users first register or publish the relation at the primary site and subsequently subscribe to a fragment of a registered relation from another (secondary) site.

In **peer-to-peer** asynchronous replication, more than one copy (although perhaps not all) can be designated as updatable, that is, a **master** copy. In addition to propagating changes, a conflict resolution strategy must be used to deal

with conflicting changes made at different sites. For example, Joe's age may be changed to 35 at one site and to 38 at another. Which value is 'correct'? Many more subtle kinds of conflicts can arise in peer-to-peer replication, and in general peer-to-peer replication leads to ad hoc conflict resolution. Some special situations in which peer-to-peer replication does not lead to conflicts arise quite often, and in such situations peer-to-peer replication is best utilized. For example:

- Each master is allowed to update only a fragment (typically a horizontal fragment) of the relation, and any two fragments updatable by different masters are disjoint. For example, it may be that salaries of German employees are updated only in Frankfurt, and salaries of Indian employees are updated only in Madras, even though the entire relation is stored at both Frankfurt and Madras.
- Updating rights are held by only one master at a time. For example, one site is designated a *backup* to another site. Changes at the master site are propagated to other sites and updates are not allowed at other sites (including the backup). But, if the master site fails, the backup site takes over and updates are now permitted at (only) the backup site.

We will not discuss peer-to-peer replication further.

Implementing Primary Site Asynchronous Replication

The main issue in implementing primary site replication is determining how changes to the primary copy are propagated to the secondary copies. Changes are usually propagated in two steps, called *Capture* and *Apply*. Changes made by committed transactions to the primary copy are somehow identified during the Capture step and subsequently propagated to secondary copies during the Apply step.

In contrast to synchronous replication, a transaction that modifies a replicated relation directly locks and changes only the primary copy. It is typically committed long before the Apply step is carried out. Systems vary considerably in their implementation of these steps. We present an overview of some of the alternatives.

Capture

The Capture step is implemented using one of two approaches. In log-based Capture, the log maintained for recovery purposes is used to generate a record of updates. Basically, when the log tail is written to stable storage, all log

records that affect replicated relations are also written to a separate change data table (**eDT**). Since the transaction that generated the update log record may still be active when the record is written to the **CDT**, it may subsequently abort. Update log records written by transactions that subsequently abort must be removed from the **eDT** to obtain a stream of updates due (only) to committed transactions. This stream can be obtained as part of the Capture step or subsequently in the Apply step if committed log records are added to the **eDT**; for concreteness, we assume that the committed update stream is obtained as part of the Capture step and that the **CDT** sent to the Apply step contains only update log records of committed transactions.

In procedural Capture, a procedure automatically invoked by the DBMS or an application program initiates the Capture process, which consists typically of taking a snapshot of the primary copy. A snapshot is just a copy of the relation as it existed at some instant in time. (A procedure that is automatically invoked by the DBMS, such as the one that initiates Capture, is called a *trigger*. We covered triggers in Chapter 5.)

Log-based Capture has a smaller overhead than procedural Capture and, because it is driven by changes to the data, results in a smaller delay between the time the primary copy is changed and the time that the change is propagated to the secondary copies. (Of course, this delay also depends on how the Apply step is implemented.) In particular, only changes are propagated, and related changes (e.g., updates to two tables with a referential integrity constraint between them) are propagated together. The disadvantage is that implementing log-based Capture requires a detailed understanding of the structure of the log, which is quite system specific. Therefore, a vendor cannot easily implement a log-based Capture mechanism that will capture changes made to data in another vendor's DBMS.

Apply

The Apply step takes the changes collected by the Capture step, which are in the **CDT** table or a snapshot, and propagates them to the secondary copies. This can be done by having the primary site continuously send the **CDT** or periodically requesting (the latest portion of) the **CDT** or a snapshot from the primary site. Typically, each secondary site runs a copy of the Apply process and 'pulls' the changes in the **eDT** from the primary site using periodic requests. The interval between such requests can be controlled by a timer or a user's application program. Once the changes are available at the secondary site, they can be applied directly to the replica.

In some systems, the replica need not be just a fragment of the original relation; it can be a view defined using SQL, and the replication mechanism is sufficiently sophisticated to maintain such a view at a remote site incrementally (by reevaluating only the part of the view affected by changes recorded in the CDT).

Log-based Capture in conjunction with continuous Apply minimizes the delay in propagating changes. It is the best combination in situations where the primary and secondary copies are both used as part of an operational DBMS and replicas must be as closely synchronized with the primary copy as possible. Log-based Capture with continuous Apply is essentially a less expensive substitute for synchronous replication. Procedural Capture and application-driven Apply offer the most flexibility in processing source data and changes before altering the replica; this flexibility is often useful in data warehousing applications where the ability to 'clean' and filter the retrieved data is more important than the currency of the replica.

Data Warehousing: An Example of Replication

Complex decision support queries that look at data from multiple sites are becoming very important. The paradigm of executing queries that span multiple sites is simply inadequate for performance reasons. One way to provide such complex query support over data from multiple sources is to create a copy of all the data at some one location and use the copy rather than going to the individual sources. Such a copied collection of data is called a data warehouse. Specialized systems for building, maintaining, and querying data warehouses have become important tools in the marketplace.

Data warehouses can be seen as one instance of asynchronous replication, in which copies are updated relatively infrequently. When we talk of replication, we typically call copies maintained under the control of a single DBMS, whereas with data warehousing, the original data may be on different software platforms (including database systems and as file systems) and even belong to different organizations. This distinction, however, is likely to become blurred as vendors adopt more 'open' strategies to replication. For example, some products already support the maintenance of replicas of relations stored in one vendor's DBMS in another vendor's DBMS.

We note that data warehousing involves more than just replication. We discuss other aspects of data warehousing in Chapter 2.5.

22.12 DISTRIBUTED TRANSACTIONS

In a distributed DBMS, a given transaction is submitted at some one site, but it can access data at other sites as well. In this chapter we refer to the activity of a transaction at a given site as a subtransaction. When a transaction is submitted at some site, the transaction manager at that site breaks it up into a collection of one or more subtransactions that execute at different sites, submits them to transaction managers at the other sites, and coordinates their activity.

We now consider aspects of concurrency control and recovery that require additional attention because of data distribution. As we saw in Chapter 16, there are many concurrency control protocols; in this chapter, for concreteness, we assume that Strict 2PL with deadlock detection is used. We discuss the following issues in subsequent sections:

- **Distributed Concurrency Control:** How can locks for objects stored across several sites be managed? How can deadlocks be detected in a distributed database?
- **Distributed Recovery:** Transaction atomicity must be ensured-----when a transaction commits, all its actions, across all the sites at which it executes, must persist. Similarly, when a transaction aborts, none of its actions must be allowed to persist.

22.13 DISTRIBUTED CONCURRENCY CONTROL

In Section 22.11.1, we described two techniques for implementing synchronous replication, and in Section 22.11.2, we discussed various techniques for implementing asynchronous replication. The choice of technique determines *which* objects are to be locked. *When* locks are obtained and released is determined by the concurrency control protocol. We now consider how lock and unlock requests are implemented in a distributed environment.

Lock management can be distributed across sites in many ways:

- **Centralized:** A single site is in charge of handling lock and unlock requests for all objects.
- **Primary Copy:** One copy of each object is designated the primary copy. All requests to lock or unlock a copy of this object are handled by the lock manager at the site where the primary copy is stored, regardless of where the copy itself is stored.

- **Fully Distributed:** Requests to lock or unlock a copy of an object stored at a site are handled by the lock manager at the site where the copy is stored.

The centralized scheme is vulnerable to failure of the single site that controls locking. The primary copy scheme avoids this problem, but in general, reading an object requires communication with two sites: the site where the primary copy resides and the site where the copy to be read resides. This problem is avoided in the fully distributed scheme, because locking is done at the site where the copy to be read resides. However, while writing, locks must be set at all sites where copies are modified in the fully distributed scheme, whereas locks need be set only at one site in the other two schemes.

Clearly, the fully distributed locking scheme is the most attractive scheme if reads are much more frequent than writes, as is usually the case.

22.13.1 Distributed Deadlock

One issue that requires special attention when using either primary copy or fully distributed locking is deadlock detection. (Of course, a deadlock prevention scheme can be used instead, but we focus on deadlock detection, which is widely used.) As in a centralized DBMS, deadlocks must be detected and resolved (by aborting some deadlocked transaction).

Each site maintains a local waits-for graph, and a cycle in a local graph indicates a deadlock. However, there can be a deadlock even if no local graph contains a cycle. For example, suppose that two sites, A and B, both contain copies of objects O1 and O2, and that the read-any write-all technique is used. T1, which wants to read O1 and write O2, obtains an S lock on O1 and an X lock on O2 at Site A, then requests an X lock on O2 at Site B. T2, which wants to read O2 and write O1, meanwhile, obtains an S lock on O2 and an X lock on O1 at Site B, then requests an X lock on O1 at Site A. As Figure 22.5 illustrates, T2 is waiting for T1 at Site A, and T1 is waiting for T2 at Site B; thus, we have a deadlock, which neither site can detect based solely on its local waits-for graph.

To detect such deadlocks, a distributed deadlock detection algorithm must be used. We describe three such algorithms.

The first algorithm, which is centralized, consists of periodically sending all local waits-for graphs to one site that is responsible for global deadlock detection. At this site, the global waits-for graph is generated by combining all the local graphs; the set of nodes is the union of nodes in the local graphs, and there is

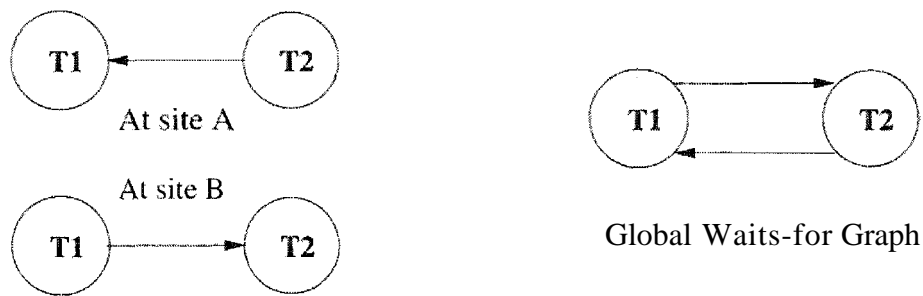


Figure 22.5 Distributed Deadlock

an edge from one node to another if there is such an edge in any of the local graphs.

The second algorithm, which is hierarchical, groups sites into a hierarchy. For instance, sites might be grouped by state, then by country, and finally into a single group that contains all sites. Every node in this hierarchy constructs a waits-for graph that reveals deadlocks involving only sites contained in (the subtree rooted at) this node. All sites periodically (e.g., every 10 seconds) send their local waits-for graph to the site responsible for constructing the waits-for graph for their state. The sites constructing waits-for graphs at the state level periodically (e.g., every minute) send the state waits-for graph to the site constructing the waits-for graph for their country. The sites constructing waits-for graphs at the country level periodically (e.g., every 10 minutes) send the country waits-for graph to the site constructing the global waits-for graph. This scheme is based on the observation that more deadlocks are likely across closely related sites than across unrelated sites, and it puts more effort into detecting deadlocks across related sites. All deadlocks are eventually detected, but a deadlock involving two different countries may take a while to detect.

The third algorithm is simple: If a transaction waits longer than some chosen time-out interval, it is aborted. Although this algorithm may cause many unnecessary restarts, the overhead of deadlock detection is (obviously!) low, and in a heterogeneous distributed database, if the participating sites cannot cooperate to the extent of sharing their waits-for graphs, it may be the only option.

A subtle point to note with respect to distributed deadlock detection is that delays in propagating local information might cause the deadlock detection algorithm to identify 'deadlocks' that do not really exist. Such situations, called **phantom deadlocks**, lead to unnecessary aborts. For concreteness, we discuss the centralized algorithm, although the hierarchical algorithm suffers from the same problem.

Consider a modification of the previous example. As before, the two transactions wait for each other, generating the local waits-for graphs shown in Figure 22.5, and the local waits-for graphs are sent to the global deadlock-detection site. However, T_2 is now aborted for reasons other than deadlock. (For example, T_2 may also be executing at a third site, where it reads an unexpected data value and decides to abort.) At this point, the local waits-for graphs have changed so that there is no cycle in the 'true' global waits-for graph. However, the constructed global waits-for graph will contain a cycle, and T_1 may well be picked as the victim!

22.14 DISTRIBUTED RECOVERY

Recovery in a distributed DBMS is more complicated than in a centralized DBMS for the following reasons:

- New kinds of failure can arise: failure of communication links and failure of a remote site at which a subtransaction is executing.
- Either all subtransactions of a given transaction must commit or none must commit, and this property must be guaranteed despite any combination of site and link failures. This guarantee is achieved using a commit protocol.

As in a centralized DBMS, certain actions are carried out as part of normal execution to provide the necessary information to recover from failures. A log is maintained at each site, and in addition to the kinds of information maintained in a centralized DBMS, actions taken as part of the commit protocol are also logged. The most widely used commit protocol is called *Two-Phase Commit (2PC)*. A variant called *2PC with Presumed Abort*, which we discuss next, has been adopted as an industry standard.

In this section, we first describe the steps taken during normal execution, concentrating on the commit protocol, and then discuss recovery from failures.

22.14.1 Normal Execution and Commit Protocols

During normal execution, each site maintains a log, and the actions of a subtransaction are logged at the site where it executes. The regular logging activity described in Chapter 18 is carried out and, in addition, a commit protocol is followed to ensure that all subtransactions of a given transaction either commit or abort uniformly. The transaction manager at the site where the transaction originated is called the coordinator for the transaction; transaction managers at sites where its subtransactions execute are called subordinates (with respect to the coordination of this transaction).

We now describe the **Two-Phase Commit (2PC)** protocol, in terms of the messages exchanged and the log records written. When the user decides to commit a transaction, the commit command is sent to the coordinator for the transaction. This initiates the 2PC protocol:

1. The coordinator sends a *prepare* message to each subordinate.
2. When a subordinate receives a *prepare* message, it decides whether to abort or commit its subtransaction. It force-writes an abort or prepare log record, and *then* sends a *no* or *yes* message to the coordinator. Note that a prepare log record is not used in a centralized DBMS; it is unique to the distributed commit protocol.
3. If the coordinator receives *yes* messages from all subordinates, it force-writes a commit log record and then sends a *commit* message to all subordinates. If it receives even one *no* message or receives no response from some subordinate for a specified time-out interval, it force-writes an abort log record, and then sends an *abort* message to all subordinates.¹
4. When a subordinate receives an *abort* message, it force-writes an abort log record, sends an *ack* message to the coordinator, and aborts the subtransaction. When a subordinate receives a *commit* message, it force-writes a commit log record, sends an *ack* message to the coordinator, and commits the subtransaction.
5. After the coordinator has received *ack* messages from all subordinates, it writes an end log record for the transaction.

The name *Two-Phase Commit* reflects the fact that two rounds of messages are exchanged: first a voting phase, then a termination phase, both initiated by the coordinator. The basic principle is that any of the transaction manager's involved (including the coordinator) can unilaterally abort a transaction, whereas there must be unanimity to commit a transaction. When a message is sent in 2PC, it signals a decision by the sender. To ensure that this decision survives a crash at the sender's site, the log record describing the decision is always forced to stable storage *before* the message is sent.

A transaction is officially committed at the time the coordinator's commit log record reaches stable storage. Subsequent failures cannot affect the outcome of the transaction; it is irrevocably committed. Log records written to record the commit protocol actions contain the type of the record, the transaction id, and the identity of the coordinator. A coordinator's commit or abort log record also contains the identities of the subordinates.

¹As an optimization, the coordinator need not send *abort* messages to subordinates who voted *no*,

22.14.2 Restart after a Failure

When a site comes back up after a crash, we invoke a recovery process that reads the log and processes all transactions executing the 2PC protocol at the time of the crash. The transaction manager at this site could have been the coordinator for some of these transactions and a subordinate for others. We do the following in the recovery process:

- If we have a commit or abort log record for transaction T , its status is clear; we redo or undo T , respectively. If this site is the coordinator, which can be determined from the commit or abort log record, we must periodically resend—because there may be other link or site failures in the system—a *commit* or *abort* message to each subordinate until we receive an *ack*. After we have received *acks* from all subordinates, we write an end log record for T .
- If we have a prepare log record for T but no commit or abort log record, this site is a subordinate, and the coordinator can be determined from the prepare record. We must repeatedly contact the coordinator site to determine the status of T . Once the coordinator responds with either *commit* or *abort*, we write a corresponding log record, redo or undo the transaction, and then write an end log record for T .
- If we have no prepare, commit, or abort log record for transaction T , T certainly could not have voted to commit before the crash; so we can unilaterally abort and undo T and write an end log record. In this case, we have no way to determine whether the current site is the coordinator or a subordinate for T . However, if this site is the coordinator, it might have sent a *prepare* message prior to the crash, and if so, other sites may have voted *yes*. If such a subordinate site contacts the recovery process at the current site, we now know that the current site is the coordinator for T , and given that there is no commit or abort log record, the response to the subordinate should be to abort T .

(Observe that, if the coordinator site for a transaction T fails, subordinates who voted *yes* cannot decide whether to commit or abort T until the coordinator site recovers; we say that T is blocked. In principle, the active subordinate sites could communicate among themselves, and if at least one of them contains an abort or commit log record for T , its status becomes globally known. If all communicate among themselves, all subordinates must be told the identity of the other subordinates at the time they are sent the *prepare* message. However, 2PC is still vulnerable to coordinator failure during recovery because even if all subordinates voted *yes*, the coordinator (who also has a vote!) may have decided to abort T , and this decision cannot be determined until the coordinator site recovers.

We covered how a site recovers from a crash, but what should a site that is involved in the COIUnit protocol do if a site that it is communicating with fails? If the current site is the coordinator, it should simply abort the transaction. If the current site is a subordinate, and it has not yet responded to the coordinator's *prepare* message, it can (and should) abort the transaction. If it is a subordinate and has voted *yes*, then it cannot unilaterally abort the transaction, and it cannot COIUnit either; it is blocked. It must periodically contact the coordinator until it receives a reply.

Failures of COIUnit communication links are seen by active sites as failure of other sites that they are communicating with, and therefore the solutions just outlined apply to this case as well.

22.14.3 Two-Phase Commit Revisited

Now that we examined how a site recovers from a failure, and saw the interaction between the 2PC protocol and the recovery process, it is instructive to consider how 2PC can be refined further. In doing so, we arrive at a more efficient version of 2PC, but equally important perhaps, we understand the role of the various steps of 2PC more clearly. Consider three basic observations:

1. The *ack* messages in 2PC are used to determine when a coordinator (or the recovery process at a coordinator site following a crash) can 'forget' about a transaction *T*. Until the coordinator knows that all subordinates are aware of the commit or abort decision for *T*, it must keep information about *T* in the transaction table.
2. If the coordinator site fails after sending out *prepare* messages but before writing a commit or abort log record, when it comes back up, it has no information about the transaction's commit status prior to the crash. However, it is still free to abort the transaction unilaterally (because it has not written a commit record, it can still cast a *no* vote itself). If another site inquires about the status of the transaction, the recovery process, as we have seen, responds with an *abort* message. Therefore, in the absence of information, a transaction is *presumed to have aborted*.
3. If a subtransaction does no updates, it has no changes to either redo or undo: in other words, its commit or abort status is irrelevant.

The first two observations suggest several refinements:

- When a coordinator aborts a transaction *T*, it can undo *T* and remove it from the transaction table immediately. After all removing *T* from the table results in a 'no information' state with respect to *T*, and the default

response (to an enquiry about T) in this state, which is *abort*, is the correct response for an aborted transaction.

- By the same token, if a subordinate receives an *abort* message, it need not send an *ack* message. The coordinator is not waiting to hear from subordinates after sending an *abort* message! If, for some reason, a subordinate that receives a *prepare* message (and voted *yes*) does not receive an *abort* or *commit* message for a specified time-out interval, it contacts the coordinator again. If the coordinator decided to abort, there may no longer be an entry in the transaction table for this transaction, but the subordinate receives the default *abort* message, which is the correct response.
- Because the coordinator is not waiting to hear from subordinates after deciding to abort a transaction, the names of subordinates need not be recorded in the abort log record for the coordinator.
- All abort log records (for the coordinator as well as subordinates) can simply be appended to the log tail, instead of doing a force-write. After all, if they are not written to stable storage before a crash, the default decision is to abort the transaction.

The third basic observation suggests some additional refinements:

- If a subtransaction does no updates (which can be easily detected by keeping a count of update log records), the subordinate can respond to a *prepare* message from the coordinator with a *reader* message, instead of *yes* or *no*. The subordinate writes no log records in this case.
- When a coordinator receives a *reader* message, it treats the message as a *yes* vote, but with the optimization that it does not send any more messages to the subordinate, because the subordinate's commit or abort status is irrelevant.
- If all subtransactions, including the subtransaction at the coordinator site, send a *reader* message, we do not need the second phase of the commit protocol. Indeed, we can simply remove the transaction from the transaction table, without writing any log records at any site for this transaction.

The Two-Phase Commit protocol with the refinements discussed in this section is called **Two-Phase Commit with Presumed Abort**.

22.14.4 Three-Phase Commit

A commit protocol called **Three-Phase Commit (3PC)** can avoid blocking even if the coordinator site fails during recovery. The basic idea is that, when

the coordinator sends out *prepare* messages and receives *yes* votes from all subordinates, it sends all sites a *precommit* message, rather than a *commit* message. When a sufficient number—more than the maximum number of failures that must be handled—of *acks* have been received, the coordinator force-writes a *commit* log record and sends a *commit* message to all subordinates. In 3PC, the coordinator effectively postpones the decision to commit until it is sure that enough sites know about the decision to commit; if the coordinator subsequently fails, these sites can communicate with each other and detect that the transaction must be committed—conversely, aborted, if none of them has received a *precommit* message—without waiting for the coordinator to recover.

The 3PC protocol imposes a significant additional cost during normal execution and requires that communication link failures do not lead to a network partition (wherein some sites cannot reach some other sites through any path) to ensure freedom from blocking. For these reasons, it is not used in practice.

22.15 REVIEW QUESTIONS

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

- Discuss the different motivations behind parallel and distributed databases. (Section 22.1)
- Describe the three main architectures for parallel DBMSs. Explain why the *shared-memory* and *shared-disk* approaches suffer from *interference*. What can you say about the *speed-up* and *scale-up* of the *shared-nothing* architecture? (Section 22.2)
- Describe and differentiate *pipelined parallelism* and *data-partitioned parallelism*. (Section 22.3)
- Discuss the following techniques for partitioning data: *round-robin*, *hash*, and *range*. (Section 22.3.1)
- Explain how existing code can be parallelized by introducing *split* and *merge* operators. (Section 22.3.2)
- Discuss how each of the following operators can be parallelized using data partitioning: *scanning*, *sorting*, *join*. Compare the use of sorting versus hashing for partitioning. (Section 22.4)
- What do we need to consider in optimizing queries for parallel execution? Discuss interoperation parallelism, left-deep trees versus bushy trees, and cost estimation. (Section 22.5)

- Define the terms *distributed data independence* and *distributed transaction atomicity*. Are these concepts supported in current commercial systems? Why not? What is the difference between *homogeneous* and *heterogeneous* distributed databases? (Section 22.6)
- Describe the three main architectures for distributed DBMSs. (Section 22.7)
 - 1\ relation can be distributed by *fragmenting* it or *replicating* it across several sites. Explain these concepts and how they differ. Also, distinguish between *horizontal* and *vertical* fragmentation. (Section 22.8)
- If a relation is fragmented and replicated, each partition needs a globally unique name called the *relation name*. Explain how such global names are created and the motivation behind the described approach to naming. (Section 22.9.1)
- Explain how metadata about such distributed data is maintained in a *distributed catalog*. (Section 22.9.2)
- Describe a naming scheme that supports distributed data independence. (Section 22.9.3)
- When processing queries in a distributed DBMS, the location of partitions of the relation needs to be taken into account. Discuss the alternatives when joining two relations that reside on different sites. In particular, explain and describe the motivation behind the *Semijoin* and *Blockjoin* techniques. (Section 22.10.2)
- What issues must be considered in optimizing queries over distributed data, in addition to where the data is located? (Section 22.10.3)
- What is the difference between *synchronous* and *asynchronous* replication? Why has asynchronous replication gained in popularity? (Section 22.11)
- Describe the *locking* and *Read-a'ny write-all* approaches to synchronous replication. (Section 22.11.1)
- Summarize the *peer-to-peer* and *primary site* approaches to asynchronous replication. (Section 22.11.2)
- In primary site replication, changes to the primary copy must be propagated to secondary copies. What is done in the *Capture* and *Apply* steps? Describe *log-based* and *procedural* approaches to Capture and compare them. What are the variations in scheduling the Apply step? Illustrate the use of asynchronous replication in a data warehouse. (Section 22.11.2)
- What is a *subtransaction*? (Section 22.12)

- What are the choices for managing locks in a distributed DBMS? (Section 22.13)
- Discuss deadlock detection in a distributed database. Contrast the *centralized*, *hierarchical*, and *time-out* approaches. (Section 22.13.1)
- Why is recovery in a distributed DBMS more complicated than in a centralized system? (Section 22.14)
- What is a *commit protocol* and why is it required in a distributed database? Describe and compare *Two-Phase* and *Three-Phase Commit*. What is *blocking*, and how does the Three-Phase protocol prevent it? Why is it nonetheless not used in practice? (Section 22.14)

EXERCISES

Exercise 22.1 Give brief answers to the following questions:

1. What are the similarities and differences between parallel and distributed database management systems?
2. Would you expect to see a parallel database built using a wide-area network? Would you expect to see a distributed database built using a wide-area network? Explain.
3. Define the terms *scale-up* and *speed-up*.
4. Why is a shared-nothing architecture attractive for parallel database systems?
5. The idea of building specialized hardware to run parallel database applications received considerable attention but has fallen out of favor. Comment on this trend.
6. What are the advantages of a distributed DBMS over a centralized DBMS?
7. Briefly describe and compare the Client-Server and Collaborating Servers architectures.
8. In the Collaborating Servers architecture, when a transaction is submitted to the DBMS, briefly describe how its activities at various sites are coordinated. In particular, describe the role of transaction managers at the different sites, the concept of *subtransactions*, and the concept of *distributed transaction atomicity*.

Exercise 22.2 Give brief answers to the following questions:

1. Define the terms *fragmentation* and *replication* in terms of where data is stored.
2. What is the difference between *synchronous* and *asynchronous* replication?
3. Define the term *distributed data independence*. What does this mean with respect to querying and updating data in the presence of data fragmentation and replication?
4. Consider the *voting* and *read-any write-all* techniques for implementing synchronous replication. What are their respective pros and cons?
5. Give an overview of how asynchronous replication can be implemented. In particular, explain the terms *Capture* and *Apply*.
6. What is the difference between log-based and procedural implementations of capture?
7. Why is giving database objects unique names more complicated in a distributed DBMS?

8. Describe a catalog organization that permits any replica (of an entire relation or a fragment) to be given a unique name and provides the naming infrastructure required for ensuring distributed data independence.
9. If information from remote catalogs is cached at other sites, what happens if the cached information becomes outdated? How can this condition be detected and resolved?

Exercise 22.3 Consider a parallel DBMS in which each relation is stored by horizontally partitioning its tuples across all disks:

Employees(eid: integer, did: integer, sal: real)
Departments(did: integer, mgrid: integer, budget: integer)

The *mgrid* field of Departments is the *eid* of the manager. Each relation contains 20-byte tuples, and the *sal* and *budget* fields both contain uniformly distributed values in the range 0 to 1 million. The Employees relation contains 100,000 pages, the Departments relation contains 5,000 pages, and each processor has 100 buffer pages of 4,000 bytes each. The cost of one page I/O is t_d , and the cost of shipping one page is t_s ; tuples are shipped in units of one page by waiting for a page to be filled before sending a message from processor i to processor j . There are no indexes, and all joins that are local to a processor are carried out using a sort-merge join. Assume that the relations are initially partitioned using a round-robin algorithm and that there are 10 processors.

For each of the following queries, describe the evaluation plan briefly and give its cost in terms of t_d and t_s . You should compute the total cost across all sites as well as the 'elapsed time' cost (i.e., if several operations are carried out concurrently, the time taken is the maximum over these operations).

1. Find the highest paid employee.
2. Find the highest paid employee in the department with *did* 55.
3. Find the highest paid employee over all departments with *budget* less than 100,000.
4. Find the highest paid employee over all departments with *budget* less than 300,000.
5. Find the average salary over all departments with *budget* less than 300,000.
6. Find the salaries of all managers.
7. Find the salaries of all managers who manage a department with a budget less than 300,000 and at least more than 100,000.
8. Print the *eids* of all employees, ordered by increasing salaries. Each processor is connected to a separate printer, and the answer can appear as several sorted lists, each printer by a different processor, as long as we can obtain a fully sorted list by concatenating the printed lists (in some order).

Exercise 22.4 Consider the same scenario as in Exercise 22.3, except that the relations are originally partitioned using range partitioning on the *sal* and *budget* fields.

Exercise 22.5 Repeat Exercises 22.3 and 22.4 with (i) 1 processor, and (ii) 10 processors.

Exercise 22.6 Consider the Employees and Departments relations described in Exercise 22.3. They are now stored in a distributed DBMS with all of Employees stored at Naples and all of Departments stored at Berlin. There are no indexes on these relations. The cost of various operations is as described in Exercise 22.3. Consider the query:

```

SELECT *
FROM   Employees E, Departments D
WHERE  E.eid = D.Ingrid

```

The query is posed at Delhi, and you are told that only 1 percent of employees are Managers. Find the cost of answering this query using each of the following plans:

1. Ship Departments to Naples, compute the query at Naples, then ship the result to Delhi.
2. Ship Employees to Berlin, compute the query at Berlin, then ship the result to Delhi.
3. COInput the query at Delhi by shipping both relations to Delhi.
4. COInput the query at Naples using BlockJoin; then ship the result to Delhi.
5. Compute the query at Berlin using BlockJoin; then ship the result to Delhi.
6. Compute the query at Naples using SortJoin; then ship the result to Delhi.
7. COInput the query at Berlin using SortJoin; then ship the result to Delhi.

Exercise 22.7 Consider your answers in Exercise 22.6. Which plan minimizes shipping costs? Is it necessarily the cheapest plan? Which do you expect to be the cheapest?

Exercise 22.8 Consider the Employees and Departments relations described in Exercise 22.3. They are now stored in a distributed DBMS with 10 sites. The Departments tuples are horizontally partitioned across the 10 sites by *did*, with the same number of tuples assigned to each site and no particular order to how tuples are assigned to sites. The Employees tuples are similarly partitioned, by *sal* ranges, with $sal \leq 100,000$ assigned to the first site, $100,000 < sal \leq 200,000$ assigned to the second site, and so on. In addition, the partition $sal \leq 100,000$ is frequently accessed and infrequently updated, and it is therefore replicated at every site. No other Employees partition is replicated.

1. Describe the best plan (unless a plan is specified) and give its cost:
 - (a) Compute the natural join of Employees and Departments by shipping all fragments of the smaller relation to every site containing tuples of the larger relation.
 - (b) Find the highest paid employee.
 - (c) Find the highest paid employee with salary less than 100,000.
 - (d) Find the highest paid employee with salary between 400,000 and 500,000.
 - (e) Find the highest paid employee with salary between 450,000 and 550,000.
 - (f) Find the highest paid manager for those departments stored at the query site.
 - (g) Find the highest paid manager.
2. Assuming the same data distribution, describe the sites visited and the locks obtained for the following update transactions, assuming that *synchronous* replication is used for the replication of Employees tuples with $sal \leq 100,000$:
 - (a) Give employees with salary less than 100,000 a 10 percent raise, with a maximum salary of 100,000 (i.e., the raise cannot increase the salary to more than 100,000).
 - (b) Give all employees a 10 percent raise. The conditions of the original partitioning of Employees must still be satisfied after the update.
3. Assuming the same data distribution, describe the sites visited and the locks obtained for the following update transactions, assuming that *asynchronous* replication is used for the replication of Employees tuples with $sal \leq 100,000$.

For all employees with salary less than 100,000 give them a 10 percent raise, with a maximum salary of 100,000.

Give all employees a 10 percent raise. After the update is completed, the conditions of the original partitioning of Employees must still be satisfied.

Exercise 22.9 Consider the Employees and Departments tables from Exercise 22.3. You are a DBA and you need to decide how to distribute these two tables across two sites, Manila and Nairobi. Your DBMS supports only unclustered B+ tree indexes. You have a choice between synchronous and asynchronous replication. For each of the following scenarios, describe how you would distribute them and what indexes you would build at each site. If you feel that you have insufficient information to make a decision, explain briefly.

1. Half the departments are located in Manila and the other half are in Nairobi. Department information, including that for employees in the department, is changed only at the site where the department is located, but such changes are quite frequent. (Although the location of a department is not included in the Departments schema, this information can be obtained from another table.)
2. Half the departments are located in Manila and the other half are in Nairobi. Department information, including that for employees in the department, is changed only at the site where the department is located, but such changes are infrequent. Finding the average salary for each department is a frequently asked query.
3. Half the departments are located in Manila and the other half are in Nairobi. Employees tuples are frequently changed (only) at the site where the corresponding department is located, but the Departments relation is almost never changed. Finding a given employee's manager is a frequently asked query.
4. Half the employees work in Manila and the other half work in Nairobi. Employees tuples are frequently changed (only) at the site where they work.

Exercise 22.10 Suppose that the Employees relation is stored in Madison and the tuples with $sal \leq 100,000$ are replicated at New York. Consider the following three options for lock management: all locks managed at a *single site*, say, Milwaukee; *primary copy* with Madison being the primary for Employees; and *fully distributed*. For each of the lock management options, explain what locks are set (and at which site) for the following queries. Also state from which site the page is read.

1. A query at Austin wants to read a page of Employees tuples with $sal \leq 50,000$.
2. A query at Madison wants to read a page of Employees tuples with $sal \leq 50,000$.
3. A query at New York wants to read a page of Employees tuples with $sal \leq 50,000$.

Exercise 22.11 Briefly answer the following questions:

1. Compare the relative merits of centralized and hierarchical deadlock detection in a distributed DBMS.
2. What is a *phantom deadlock*? Give an example.
3. Give an example of a distributed DBMS with three sites such that no two local waits-for graphs reveal a deadlock, yet there is a global deadlock.
4. Consider the following modification to a local waits-for graph: Add a new node T_{ext} , and for every transaction T_i that is waiting for a lock at another site, add the edge $T_i \rightarrow T_{ext}$. Also add an edge $T_{ext} \rightarrow T_i$ if a transaction executing at another site is waiting for T_i to release a lock at this site.

If there is a cycle in the modified local waits-for graph that does not involve T_{ext} , what can you conclude? If every cycle involves T_{ext} , what can you conclude?

Suppose that every site is assigned a unique integer. Whenever the local waits-for graph suggests that there might be a global deadlock, send the local waits-for graph to the site with the next higher site-id. At that site, combine the received graph with the local waits-for graph. If this combined graph does not indicate a deadlock, ship it on to the next site, and so on, until either a deadlock is detected or we are back at the site that originated this round of deadlock detection. Is this scheme guaranteed to find a global deadlock if one exists?

Exercise 22.12 Timestamp-based concurrency control schemes can be used in a distributed DBMS, but we must be able to generate globally unique, monotonically increasing timestamps without a bias in favor of any one site. One approach is to assign timestamps at a single site. Another is to use the local clock time and to append the site-id. A third scheme is to use a counter at each site. Compare these three approaches.

Exercise 22.13 Consider the multiple-granularity locking protocol described in Chapter 18. In a distributed DBMS, the site containing the root object in the hierarchy can become a bottleneck. You hire a database consultant who tells you to modify your protocol to allow only intention locks on the root and implicitly grant all possible intention locks to every transaction.

1. Explain why this modification works correctly, in that transactions continue to be able to set locks on desired parts of the hierarchy.
2. Explain how it reduces the demand on the root.
3. Why is this idea not included as part of the standard multiple-granularity locking protocol for a centralized DBMS?

Exercise 22.14 Briefly answer the following questions:

1. Explain the need for a commit protocol in a distributed DBMS.
2. Describe 2PC. Be sure to explain the need for force-writes.
3. Why are *ack* messages required in 2PC?
4. What are the differences between 2PC and 2PC with Presumed Abort?
5. Give an example execution sequence such that 2PC and 2PC with Presumed Abort generate an identical sequence of actions.
6. Give an example execution sequence such that 2PC and 2PC with Presumed Abort generate different sequences of actions.
7. What is the intuition behind 3PC? What are its pros and cons relative to 2PC?
8. Suppose that a site gets no response from another site for a long time. Can the first site tell whether the connecting link has failed or the other site has failed? How is such a failure handled?
9. Suppose that the coordinator includes a list of all subordinates in the *prepare* message. If the coordinator fails after sending out either an *abort* or *commit* message, can you suggest a way for active sites to terminate this transaction without waiting for the coordinator to recover? Assume that some but not all of the *abort* or *commit* messages from the coordinator are lost.

III Suppose that 2PC with Presumed Abort is used as the commit protocol. Explain how the system recovers from failure and deals with a particular transaction T in each of the following cases:

- (a) A subordinate site for T fails before receiving a *prepare* message.
- (b) A subordinate site for T fails after receiving a *prepare* message but before making a decision.
- (c) A subordinate site for T fails after receiving a *prepare* message and force-writing an abort log record but before responding to the *prepare* message.
- (d) A subordinate site for T fails after receiving a *prepare* message and force-writing a prepare log record but before responding to the *prepare* message.
- (e) A subordinate site for T fails after receiving a *prepare* message, force-writing an abort log record, and sending a *no* vote.
- (f) The coordinator site for T fails before sending a *prepare* message.
- (g) The coordinator site for T fails after sending a *prepare* message but before collecting all votes.
- (h) The coordinator site for T fails after writing an *abort* log record but before sending any further messages to its subordinates.
- (i) The coordinator site for T fails after writing a *commit* log record but before sending any further messages to its subordinates.
- (j) The coordinator site for T fails after writing an *end* log record. Is it possible for the recovery process to receive an inquiry about the status of T from a subordinate?

Exercise 22.15 Consider a heterogeneous distributed DBMS.

1. Define the terms *multidatabase system* and *gateway*.
2. Describe how queries that span multiple sites are executed in a multidatabase system. Explain the role of the gateway with respect to catalog interfaces, query optimization, and query execution.
3. Describe how transactions that update data at multiple sites are executed in a multidatabase system. Explain the role of the gateway with respect to lock management, distributed deadlock detection, Two-Phase Commit, and recovery.
4. Schemas at different sites in a multidatabase system are probably designed independently. This situation can lead to *semantic heterogeneity*; that is, units of measure may differ across sites (e.g., inches versus centimeters), relations containing essentially the same kind of information (e.g., employee salaries and ages) may have slightly different schemas, and so on. What impact does this heterogeneity have on the end user? In particular, comment on the concept of distributed data independence in such a system.

BIBLIOGRAPHIC NOTES

Work on parallel algorithms for sorting and various relational operations is discussed in the bibliographies for Chapters 13 and 14. Our discussion of parallel joins follows [220], and our discussion of parallel sorting follows [223]. DeWitt and Gray make the case that for future high performance database systems, parallelism will be the key [221]. Scheduling in parallel database systems is discussed in [522]. [496] contains a good collection of papers on query processing in parallel database systems.

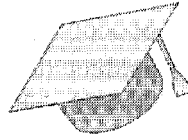
Textbook discussions of distributed databases include [78, 144, 580]. Good survey articles include [85], which focuses on concurrency control; [637], which is about distributed databases in general; and [785], which concentrates on distributed query processing. Two Inajal' projects in the area were 8DD-1 [636] and R* [777]. Fragmentation in distributed databases is considered in [157, 207]. Replication is considered in [11, 14, 137, 239, 238, 388, 385, 335, 552, 600]. For good overviews of current trends in asynchronous replication, see [234, 709, 772]. Papers on view maintenance mentioned in the bibliographic notes of Chapter 21 are also relevant in this context. Olston considers techniques for trading of performance versus precision in a replicated environment [571, 572, 573].

Query processing in the 8DD-1 distributed database is described in [88]. One of the notable aspects of 8DD-1 query processing was the extensive use of semijoins. Theoretical studies of Semijoins are presented in [83, 86, 414]. Query processing in R* is described in [667]. The R* query optimizer is validated in [500]; much of our discussion of distributed query processing is drawn from the results reported in this paper. Query processing in Distributed Ingres is described in [247]. Optimization of queries for parallel execution is discussed in [297, 323, 383]. Franklin, Jonsson, and Kossman discuss the trade-offs between *query shipping*, the more traditional approach in relational databases, and *data shipping*, which consists of shipping data to the client for processing and is widely used in object-oriented systems [284]. A good recent survey of distributed query processing techniques can be found in [450].

Concurrency control in the 8DD-1 distributed database is described in [91]. Transaction management in R* is described in [547]. Concurrency control in Distributed Ingres is described in [714]. [740] provides an introduction to distributed transaction management and various notions of distributed data independence. Optimizations for read-only transactions are discussed in [306]. Multiversion concurrency control algorithms based on timestamps were proposed in [620]. Timestamp-based concurrency control is discussed in [84, 356]. Concurrency control algorithms based on voting are discussed in [303, 318, 408, 452, 732]. The rotating primary copy scheme is described in [538]. Optimistic concurrency control in distributed databases is discussed in [660], and adaptive concurrency control is discussed in [488].

Two-Phase Commit was introduced in [466, 331]. 2PC with Presumed Abort is described in [546], along with an alternative called *2PC with Presumed Commit*. A variation of Presumed Commit is proposed in [465]. Three-Phase Commit is described in [692]. The deadlock detection algorithms in R* are described in [567]. Many papers discuss deadlocks, for example, [156, 243, 526, 632]. [441] is a survey of several algorithms in this area. Distributed clock synchronization is discussed by [464]. [333] argues that distributed data independence is not always a good idea, due to processing and administrative overheads. The ARIES algorithm is applicable for distributed recovery, but the details of how messages should be handled are not discussed in [544]. The approach taken to recovery in SDD-1 is described in [43]. [114] also addresses distributed recovery. [444] is a survey article that discusses concurrency control and recovery in distributed systems. [95] contains several articles on these topics.

Multidatabase systems are discussed in [10, 113, 230, 231, 242, 476, 485, 519, 520, 599, 641, 765, 797]; see [112, 486, 684] for surveys.



23

OBJECT-DATABASE SYSTEMS

- What are object-database systems and what new features do they support?
- What kinds of applications do they benefit?
- What kinds of data types can users define?
- What are abstract data types and their benefits?
- What is type inheritance and why is it useful?
- What is the impact of introducing object ids in a database?
- How can we utilize the new features in database design?
- What are the new implementation challenges?
- What differentiates object-relational and object-oriented DBMSs?
- Key concepts: user-defined data types, structured types, collection types; data abstraction, methods, encapsulation; inheritance, early and late binding of methods, collection hierarchies; object identity, reference types, shallow and deep equality

with Joseph M. Heisterstein
University of California-Berkeley

You know my methods, Watson. Apply them.

Arthur Conan Doyle, *The Memoirs of Sherlock Holmes*

Relational database systems support a small, fixed collection of data types (e.g., integers, dates, strings), which has proven adequate for traditional application domains such as administrative data processing. In many application domains, however, much more complex kinds of data must be handled. Typically this complex data has been stored in OS file systems or specialized data structures, rather than in a DBMS. Examples of domains with complex data include computer-aided design and modeling (CAD/CAM), multimedia repositories, and document management.

As the amount of data grows, the many features offered by a DBMS—for example, reduced application development time, concurrency control and recovery, indexing support, and query capabilities—become increasingly attractive and, ultimately, necessary. To support such applications, a DBMS must support complex data types. Object-oriented concepts strongly influenced efforts to enhance database support for complex data and led to the development of object-database systems, which we discuss in this chapter.

Object-database systems have developed along two distinct paths:

- **Object-Oriented Database Systems:** Object-oriented database systems are proposed as an alternative to relational systems and are aimed at application domains where complex objects play a central role. The approach is heavily influenced by object-oriented programming languages and can be understood as an attempt to add DBMS functionality to a programming language environment. The Object Database Management Group (ODMG) has developed a standard Object Data Model (ODM) and Object Query Language (OQL), which are the equivalent of the SQL standard for relational database systems.
- **Object-Relational Database Systems:** Object-relational database systems can be thought of as an attempt to extend relational database systems with the functionality necessary to support a broader class of applications and, in many ways, provide a bridge between the relational and object-oriented paradigms. The SQL:1999 standard extends SQL to incorporate support for the object-relational model of data.

We use acronyms for relational, object-oriented, and object-relational database management systems (RDBMS, OODBMS, ORJDBMS). In this chapter, we focus on ORDBMSs and emphasize how they can be viewed as a development of RDBMSs, rather than as an entirely different paradigm, as exemplified by the evolution of SQL:1999.

We concentrate on developing the fundamental concepts rather than presenting SQL:1999; some of the features we discuss are not included in SQL:1999.

Nonetheless, we have chosen to emphasize concepts relevant to SQL:1999 and its likely future extensions. We also try to be consistent with SQL:1999 for notation, although we occasionally diverge slightly for clarity. It is important to recognize that the main concepts discussed are common to both ORDBMSs and OODBMSs; we discuss how they are supported in the ODL/OQL standard proposed for OODBMSs in Section 23.9.

RDBMS vendors, including IBM, Informix, and Oracle, are adding ORDBMS functionality (to varying degrees) in their products, and it is important to recognize how the existing body of knowledge about the design and implementation of relational databases can be leveraged to deal with the ORDBMS extensions. It is also important to understand the challenges and opportunities these extensions present to database users, designers, and implementors.

In this chapter, Sections 23.1 through 23.6 introduce object-oriented concepts. The concepts discussed in these sections are common to both OODBMSs and ORDBMSs. We begin by presenting an example in Section 23.1 that illustrates why extensions to the relational model are needed to cope with some new application domains. This is used as a running example throughout the chapter. We discuss the use of type constructors to support user-defined structured data types in Section 23.2. We consider what operations are supported on these new types of data in Section 23.3. Next, we discuss data encapsulation and abstract data types in Section 23.4. We cover inheritance and related issues, such as method binding and collection hierarchies, in Section 23.5. We then consider objects and object identity in Section 23.6.

We consider how to take advantage of the new object-oriented concepts to do OODBMS database design in Section 23.7. In Section 23.8, we discuss some of the new implementation challenges posed by object-relational systems. We discuss ODL and OQL, the standards for OODBMSs, in Section 23.9, and then present a brief comparison of ORDBMSs and OODBMSs in Section 23.10.

23.1 MOTIVATING EXAMPLE

As a specific example of the need for object-relational systems, we focus on a new business data processing problem that is both harder and (in our view) more entertaining than the dollars and cents bookkeeping of previous decades. Today, companies in industries such as entertainment are in the business of selling *bits*; their basic corporate assets are not tangible products, but rather software artifacts such as video and audio.

We consider the fictional Dinky Entertainment Company, a large Hollywood conglomerate whose main assets are a collection of cartoon characters, espe-

cially the cuddly and internationally beloved Herbert the \Varill. Dinky has several Herbert the \Vornlfihns, many of which are shown in theaters around the world at any given time. Dinky also makes a good deal of money licensing Herbert's image, voice, and video footage for various purposes: action figures, video games, product endorsements, and so on. Dinky's database is used to manage the sales and leasing records for the various Herbert-related products, as well as the video and audio data that make up Herbert's many films.

23.1.1 New **Data** Types

The basic problem confronting Dinky's database designers is that they need support for considerably richer data types than is available in a relational DBMS:

- **User-defined data types:** Dinky's assets include Herbert's image, voice, and video footage, and these must be stored in the database. To handle these new types, we need to be able to represent richer structure. (See Section 23.2.) Further, we need special functions to manipulate these objects. For example, we may want to write functions that produce a compressed version of an image or a lower-resolution image. By hiding the details of the data structure through the functions that capture the behavior, we achieve *data abstraction*, leading to cleaner code design. (See Section 23.4.)
- **Inheritance:** As the number of data types grows, it is important to take advantage of the commonality between different types. For example, both compressed images and lower-resolution images are, at some level, just images. It is therefore desirable to *inherit* some features of image objects while defining (and later manipulating) compressed image objects and lower-resolution image objects. (See Section 23.5.)
- **Object Identity:** Given that some of the new data types contain very large instances (e.g., videos), it is important not to store copies of objects; instead, we must store *references*, or *pointers*, to such objects. In turn, this underscores the need for giving objects a unique *object identity*, which can be used to refer or 'point' to them from elsewhere in the data. (See Section 23.6.)

How might we address these issues in an RDBMS? We could store images, videos, and so on as BLOBs in current relational systems. A binary large object (BLOB) is just a long stream of bytes, and the DBMS's support consists of storing and retrieving BLOBs in such a manner that a user does not have to worry about the size of the BLOB; a BLOB can span several pages, unlike a traditional attribute. All further processing of the BLOB has to be done by the user's application program, in the host language in which the

The SQL/MM Standard: SQL/MM is an emerging standard that builds upon SQL:1999's new data types to define extensions of SQL:1999 that facilitate handling of complex multimedia data types. SQL/MM is a multipart standard. Part 1, SQL/MM Framework, identifies the SQL:1999 concepts that are the foundation for SQL/MM extensions. Each of the remaining parts addresses a specific type of complex data: Full Text, Spatial, Still Image, and Data Mining. SQL/MM anticipates that these new complex types can be used in columns of tables as field values.

Large Objects: SQL:1999 includes a new data type called LARGE OBJECT or LOB, with two variants called BLOB (binary large object) and CLOB (character large object). This standardizes the large object support found in many current relational DBMSs. LOBs cannot be included in primary keys, GROUP BY, or ORDER BY clauses. They can be compared using equality, inequality, and substring operations. A LOB has a locator that is essentially a unique id and allows LOBs to be manipulated without extensive copying.

LOBs are typically stored separately from the data records in whose fields they appear. IBM DB2, Informix, Microsoft SQL Server, Oracle 8, and Sybase ASE all support LOBs.

SQL code is embedded. This solution is not efficient because we are forced to retrieve all BLOBs in a collection even if most of them could be filtered out of the answer by applying user-defined functions (within the DBMS). It is not satisfactory from a data consistency standpoint either, because the semantics of the data now depends heavily on the host language application code and cannot be enforced by the DBMS.

As for structured types and inheritance, there is simply no support in the relational model. We are forced to map data with such complex structure into a collection of flat tables. (We saw examples of such mappings when we discussed the translation from ER diagrams with inheritance to relations in Chapter 2.)

This application clearly requires features not available in the relational model. As an illustration of these features, Figure 23.1 presents SQL:1999 DDL statements for a portion of Dinky's (HJ)31VIS schema used in subsequent examples. Although the DDL is very similar to that of a traditional relational system, some important distinctions highlight the new data modeling capabilities of an ORDBMS. A quick glance at the DDL statements is sufficient for now; we study them in detail in the next section, after presenting some of the basic

concepts that our sample application suggests are needed in a next-generation DBMS.

1. CREATE TABLE Frames
 (*framenno* integer, *image* jpeg_image, *category* integer);
2. CREATE TABLE Categories
 (*cid* integer, *name* text, *lease_price* float, *comments* text);
3. CREATE TYPE theater_t AS
 ROW(*tno* integer, *name* text, *address* text, *phone* text)
 REF IS SYSTEM GENERATED;
4. CREATE TABLE Theaters OF theater_t REF is tid SYSTEM GENERATED;
5. CREATE TABLE Nowshowing
 (*film* integer, *theater* REF(theater_t) SCOPE rrheaters, *start* date,
 end date);
6. CREATE TABLE FillIns
 (*filrno* integer, *title* text, *stars* VARCHAR(25) ARRAY [10]),
 director text, *budget* float);
7. CREATE TABLE Countries
 (*name* text, *boundary* polygon, *population* integer, *language* text);

Figure 23.1 SQL:1999 DDL Staternents for Dinky Schema

23.1.2 Manipulating the New Data

Thus far, we described the new kinds of data that must be stored in the Dinky database. We have not yet said anything about how to *use* these new types in queries, so let us study two queries that Dinky's database needs to support. The syntax of the queries is not critical; it is sufficient to understand what they express. We return to the specifics of the queries' syntax later.

Our first challenge comes from the Clog breakfast cereal company. Clog produces a cereal called Delirios and it wants to lease an image of Herbert the Worm in front of a sunrise to incorporate in the Delirios box design. A query to present a collection of possible images and their lease prices can be expressed in SQL-like syntax as in Figure 23.2. Dinky has a number of methods written in an imperative language like Java and registered with the database system. These methods can be used in queries in the same way as built-in methods, such as =, <, >, <=, >=, are used in a relational language like SQL. The *thumbnail* method in the *Select* clause produces a small version of its full-size input: image. The *is_sunrise* method is a boolean function that analyzes an image and returns *true* if the image contains a sunrise; the *is_herbert* method returns *true* if the image contains a picture of Herbert. The query produces the frame

code number, `irnage` thurnbnail, and price for all frames that contain Herbert and a sunrise.

```
SELECT F.fnuneno, thulnbail(F.irnage), C.lease_price
FROM   Fralnes F, Categories C
WHERE  F.category = C.cid AND is.Bllnrise(F.irnage) AND isJlerbert(F.inlage)
```

Figure 23.2 Extended SQL to Find Pictures of Herbert at Sunrise

The second challenge comes from Dinky's executives. They know that *Delirios* is exceedingly popular in the tiny country of Andorra, so they want to make sure that a number of Herbert films are playing at theaters near Andorra when the cereal hits the shelves. To check on the current state of affairs, the executives want to find the names of all theaters showing Herbert films within 100 kilometers of Andorra. Figure 23.3 shows this query in an SQL-like syntax.

```
SELECT N.theater-->name, N.theater->address, F.title
FROM   Nowshowing N, Films F, Countries C
WHERE  N.film = F.filmno AND
       overlaps(C.boundary, radius(N.theater->address, 100)) AND
       C.name = 'Andorra' AND 'Herbert the Worm' = F.stars[1]
```

Figure 23.3 Extended SQL to Find Herbert Films Playing near Andorra

The *theater* attribute of the *Nowshowing* table is a reference to an object in another table, which has attributes *name*, *address*, and *location*. This object referencing allows for the notation *N.theater->name* and *N.theater-->address*, each of which refers to attributes of the *theater_t* object referenced in the *Nowshowing* row *N*. The *stars* attribute of the *tUrns* table is a set of names of each film's stars. The *radius* method returns a circle centered at its first argument with radius equal to its second argument. The *overlaps* method tests for spatial overlap. *Nowshowing* and *Films* are joined by the equijoin clause, while *Nowshowing* and *Countries* are joined by the spatial overlap clause. The selections to 'Andorra' and films containing 'Herbert the Worm' complete the query.

These two object-relational queries are similar to SQL-92 queries but have some unusual features:

- **User-Defined Methods:** User-defined abstract types are manipulated via their methods, for example, *is_herbert* (Section 23.2).
- **Operators for Structured Types:** Along with the structured types available in the data model, ORDBMSs provide the natural methods for those types. For example, the *ARRAY* type supports the standard array

operation of accessing an array element by specifying the index; $F.stars[l]$ returns the first element of the array in the *stars* column of film F (Section 23.3).

- **Operators for Reference Types:** Reference types are *dereferenced* via an arrow (\rightarrow) notation (Section 23.6.2).

To summarize the points highlighted by our motivating example, traditional relational systems offer limited flexibility in the data types available. Data is stored in tables and the type of each field value is limited to a simple atomic type (e.g., integer or string), with a small, fixed set of such types to choose from. This limited type system can be extended in three main ways: user-defined abstract data types, structured types, and reference types. Collectively, we refer to these new types as **complex types**. In the rest of this chapter, we consider how a DBMS can be extended to provide support for defining new complex types and manipulating objects of these new types.

23.2 STRUCTURED DATA TYPES

SQL:1999 allows users to define new data types, in addition to the built-in types (e.g., integers). In Section 5.7.2, we discussed the definition of new *distinct* types. Distinct types stay within the standard relational model, since values of these types must be atomic.

SQL:1999 also introduced two **type constructors** that allow us to define new types with internal structure. Types defined using type constructors are called **structured types**. This takes us beyond the relational model, since field values need no longer be atomic:

- **ROW(t_1, \dots, t_n):** A type representing a row, or tuple, of n fields with fields $1, \dots, n$ of types t_1, \dots, t_n respectively.
- **base ARRAY [i]:** A type representing an array of (up to) i base-type items.

The `theater_t` type in Figure 23.1 illustrates the new ROW data type. In SQL:1999, the ROW type has a special role because every table is a collection of rows—every table is a set of rows or a multiset of rows. Values of other types can appear only as field values.

The *stars* field of table *Films* illustrates the new ARRAY type. It is an array of up to 10 elements, each of which is of type VARCHAR(25). Note that 10 is the maximum number of elements in the array; at any time, the array (unlike, say,

SQL:1999 Structured Data Types: Several commercial systems, including IBM DB2, Informix UDS, and Oracle 9i support the ROW and ARRAY constructors. The `listof`, `bagof`, and `setof` type constructors are not included in SQL:1999. Nonetheless, commercial systems support some of these constructors to varying degrees. Oracle supports nested relations and arrays, but does not support fully composing these constructors. Informix supports the `setof`, `bagof`, and `listof` constructors and allows them to be composed. Support in this area varies widely across vendors.

in C) can contain fewer elements. Since SQL:1999 does not support multidimensional arrays, *vector* might have been a more accurate name for the array constructor.

The power of type constructors comes from the fact that they can be composed. The following row type contains a field that is an array of at most 10 strings:

```
ROW(filmno: integer, stars: VARCHAR(25) ARRAY [10])
```

The row type in SQL:1999 is quite general; its fields can be of any SQL:1999 data type. Unfortunately, the array type is restricted; elements of an array cannot be arrays themselves. Therefore, the following definition is illegal:

```
(integer ARRAY [5]) ARRAY [10]
```

23.2.1 Collection Types

SQL:1999 supports only the ROW and ARRAY type constructors. Other SQL:1999 type constructors include

- `listof(base)`: A type representing a sequence of base-type items.
- `setof(base)`: A type representing a *set* of base-type items. Sets cannot contain duplicate elements.
- `bagof(base)`: A type representing a *bag* or *multiset* of base-type items.

Types using `listof`, `ARRAY`, `bagof`, or `setof` as the outermost type constructor are sometimes referred to as collection types or bulk data types.

The lack of support for these collection types is recognized as a weakness of SQL:1999's support for complex objects and it is quite possible that some of these collection types will be added in future revisions of the SQL standard.¹

23.3 OPERATIONS ON STRUCTURED DATA

The DBMS provides built-in methods for the types defined using type constructors. These methods are analogous to built-in operations such as addition and multiplication for atomic types such as integers. In this section we present the methods for various type constructors and illustrate how SQL queries can create and manipulate values with structured types.

23.3.1 Operations on Rows

Given an item i whose type is $\text{ROW}(n_1 t_1, \dots, n_n t_n)$, the field extraction method allows us to access an individual field n_k using the traditional dot notation $i.n_k$. If row constructors are nested in a type definition, dots may be nested to access the fields of the nested row; for example $i.n_k.m_l$. If we have a collection of rows, the dot notation gives us a collection as a result. For example, if i is a list of rows, $i.n_k$ gives us a list of items of type t_n ; if i is a set of rows, $i.n_k$ gives us a set of items of type t_n .

[This nested-dot notation is often called a **path expression**, because it describes a path through the nested structure.]

23.3.2 Operations on Arrays

Array types support an 'array index' method to allow users to access array items at a particular offset. A postfix 'square bracket' syntax is usually used. Since the number of elements can vary, there is an operator (CARDINALITY) that returns the number of elements in the array. The variable number of elements also motivates an operator to concatenate two arrays. The following example illustrates these operations on SQL:1999 arrays.

```
SELECT F.fullName, (F.stars || ['Brando', 'Pacino'])
FROM   FilmsF
WHERE  CARDINALITY(F.stars) < 3 AND F.stars[1]='Redford'
```

¹According to Jim Melton, the editor of the SQL:1999 standard, these collection types were considered for inclusion but omitted because some problems with their specifications were discovered too late for correction in the SQL:1999 time-frame.

For each *film* with Redford as the first star² and fewer than three stars, the result of the query contains the film's array of stars concatenated with the array containing the two elements 'Brando' and 'Pacino'. Observe how a value of type array (containing Brando and Pacino) is constructed through the use of square brackets in the *SELECT* clause.

23.3.3 Operations on Other Collection Types

Although only arrays are supported in SQL:1999, future versions of SQL are expected to support other collection types, and we consider what operations are appropriate over these types of data. We provide such operations. Our discussion is illustrative and not meant to be comprehensive. For example, one could additionally allow aggregate operators *count*, *sum*, *avg*, *max*, and *min* to be applied to any object of a collection type with an appropriate base type (e.g., *INTEGER*). One could also support operators for type conversions. For example, one could provide operators to convert a *multiset* object to a *set* object by eliminating duplicates.

Sets and Multisets

Set objects can be compared using the traditional set methods \subset , \subseteq , $=$, \supseteq , \supset . An item of type *setof*(*faa*) can be compared with an item of type *faa* using the \in method, as illustrated in Figure 23.3, which contains the comparison '*Herbert the Worm*' \in *F.stars*. Two set objects (having elements of the same type) can be combined to form a new object using the \cup , \cap , and $-$ operators.

Each of the methods for sets can be defined for multisets, taking the number of copies of elements into account. The \cup operation simply adds up the number of copies of an element, the \cap operation counts the lesser number of times a given element appears in the two input multisets, and $-$ subtracts the number of times a given element appears in the second multiset from the number of times it appears in the first multiset. For example, using multiset semantics $\cup (\{1,2,2,2\}, \{2,2,3\}) = \{1,2,2,2,2,3\}$, $\cap (\{1,2,2,2\}, \{2,2,3\}) = \{2,2\}$; and $\setminus (\{1,2,2,2\}, \{2,2,3\}) = \{1,2\}$.

Lists

Traditional list operations include *head*, which returns the first element; *tail*, which returns the list obtained by removing the first element; *prepend*, which

²Note that the first element in an SQL array has index value 1 (not 0, as in some languages).

takes an element and inserts it as the first element in a list; and *append*, which appends one list to another.

23.3.4 Queries Over Nested Collections

We now present some examples to illustrate how relations that contain nested collections can be queried, using SQL syntax. In particular, extensions of the relational model with nested sets and multisets have been widely studied and we focus on these collection types.

We consider a variant of the *Films* relation from Figure 23.1 in this section, with the *stars* field defined as a *set of* (`VARCHAR[25]`), rather than an array. Each tuple describes a film, uniquely identified by *filmno*, and contains a set (of stars in the film) as a field value.

Our first example illustrates how we can apply an aggregate operator to such a nested set. It identifies films with more than two stars by counting the number of stars; the `CARDINALITY` operator is applied once per *Films* tuple.³

```
SELECT F.filmno
FROM   Films F
WHERE  CARDINALITY(F.stars) > 2
```

Our second query illustrates an operation called *unnesting*. Consider the instance of *Films* shown in Figure 23.4; we have omitted the *director* and *budget* fields (included in the *Films* schema in Figure 23.1) for simplicity. A flat version of the same information is shown in Figure 23.5; for each film and star in the film, we have a tuple in *Films_flat*.

<i>filmno</i>	<i>title</i>	<i>stars</i>
98	Casablanca	{Bogart, Bergman}
54	Earth Vornes Are Juicy	{Herbert, Wanda}

Figure 23.4 A Nested Relation, Films

The following query generates the instance of *Films_flat* from *Films*:

```
SELECT F.filmno, F.title, S AS star
FROM   Films F, F.stars AS S
```

³SQL:1999 does not support set or multiset values, as we noted earlier. If it did, it would be natural to allow the `CARDINALITY` operator to be applied to a set-value to count the number of elements; we have used the operator in this spirit.

<i>filmno</i>	<i>title</i>	<i>star</i>
98	Casablanca	Bogart
98	Casablanca	Bergman
54	Earth Worms Are Juicy	Herbert
54	Earth Worms Are Juicy	Wanda

"Figure 23.5 A Flat Version, Films_flat

The variable F is successively bound to tuples in Films, and for each value of F , the variable S is successively bound to the set in the *stars* field of F . Conversely, we may want to generate the instance of Films from Films_flat. We can generate the Films instance using a generalized form of SQL's GROUP BY construct, as the following query illustrates:

```
SELECT  F.filmno, F.title, set_gen(F.star)
FROM    Films_flat F
GROUP BY F.filmno, F.title
```

This example introduces a new operator *set_gen*, to be used with GROUP BY, that requires some explanation. The GROUP BY clause partitions the Films_flat table by sorting on the *filmno* attribute; all tuples in a given partition have the same *filmno* (and therefore the same *title*). Consider the set of values in the *star* column of a given partition. In an SQL-92 query, this set must be summarized by applying an aggregate operator such as COUNT. Now that we allow relations to contain sets as field values, however, we can return the set of *star* values as a field value in a single answer tuple; the answer tuple also contains the *filmno* of the corresponding partition. The *set_gen* operator collects the set of *star* values in a partition and creates a set-valued object. This operation is called nesting. We can imagine similar generator functions for creating lists, and so on. However, such generators are not included in SQL:1999.

23.4 ENCAPSULATION AND ADTS

Consider the Frames table of Figure 23.1. It has a column *image* of type *jpeg_image*, which stores a compressed image representing a single frame of a film. The *jpeg_image* type is not one of the DBMS's built-in types and was defined by a user for the Dinky application to store image data compressed using the JPEG standard. As another example, the Countries table defined in Line 7 of Figure 23.1 has a column *boundary* of type *polygon*, which contains representations of the shapes of countries' outlines on a world map.

Allowing users to define arbitrary new data types is a key feature of ORDBMSs. The DBMS allows users to store and retrieve objects of type `jpeg_image`, just like an object of any other type, such as `integer`. New atomic data types usually need to have type-specific operations defined by the user who creates them. For example, one might define operations on an image data type such as `compress`, `rotate`, `shrink`, and `crop`. The combination of an atomic data type and its associated methods is called an abstract data type, or ADT. Traditional SQL comes with built-in ADTs, such as integers (with the associated arithmetic methods) or strings (with the equality, comparison, and LIKE methods). Object-relational systems include these ADTs and also allow users to define their own ADTs.

The label *abstract* is applied to these data types because the database system does not need to know how an ADT's data is stored nor how the ADT's methods work. It merely needs to know what methods are available and the input and output types for the methods. Hiding ADT internals is called encapsulation.⁴ Note that even in a relational system, atomic types such as integers have associated methods that encapsulate them. In the case of integers, the standard methods for the ADT are the usual arithmetic operators and comparators. To evaluate the addition operator on integers, the database system need not understand the laws of addition; it merely needs to know how to invoke the addition operator's code and what type of data to expect in return.

In an object-relational system, the simplification due to encapsulation is critical because it hides any substantive distinctions between data types and allows an ORDBMS to be implemented without anticipating the types and methods that users might want to add. For example, adding integers and overlaying images can be treated uniformly by the system, with the only significant distinctions being that different code is invoked for the two operations and differently typed objects are expected to be returned from that code.

23.4.1 Defining Methods

To register a new method for a user-defined data type, users must write the code for the method and then inform the database system about the method. The code to be written depends on the languages supported by the DBMS and, possibly, the operating system in question. For example, the ORDBMS may handle Java code in the Linux operating system. In this case, the method code must be written in Java and compiled into a Java bytecode file stored in a Linux file system. Then an SQL-style method registration command is given to the ODBMS so that it recognizes the new method:

⁴Some ORDBMSs actually refer to ADTs as opaque types because they are encapsulated and hence one cannot see their details.

Packaged ORDBMS Extensions: Developing a set of user-defined types and methods for a particular application—say, image management—can involve a significant amount of work and domain-specific expertise. As a result, most ORDBMS vendors partner with third parties to sell prepackaged sets of ADTs for particular domains. Informix calls these extensions *DataBlades*, Oracle calls them *Data Cartridges*, IBM calls them *DB2 Extenders*, and so on. These packages include the ADT method code, DDL scripts to automate loading the ADTs into the system, and in some cases specialized access methods for the data type. Packaged ADT extensions are analogous to the class libraries available for object-oriented programming languages: They provide a set of objects that together address a common task.

SQL:1999 has an extension called SQL/MM that consists of several independent parts, each of which specifies a type library for a particular kind of data. The SQL/MM parts for Full-Text, Spatial, Still Image, and Data Mining are available, or nearing publication.

```
CREATE FUNCTION is_sunrise(jpeg_image) RETURNS boolean
AS EXTERNAL NAME '/a/b/c/dinky.class' LANGUAGE 'java';
```

This statement defines the salient aspects of the method: the type of the associated ADT, the return type, and the location of the code. Once the method is registered, the DBMS uses a Java virtual machine to execute the code⁵. Figure 23.6 presents a number of method registration commands for our Dinky database.

1. CREATE FUNCTION thumbnail(jpeg_image) RETURNS jpeg_image
AS EXTERNAL NAME '/a/b/c/dinky.class' LANGUAGE 'java';
2. CREATE FUNCTION is_sunrise(jpeg_image) RETURNS boolean
AS EXTERNAL NAME '/a/b/e/dinky.class' LANGUAGE 'java';
3. CREATE FUNCTION isJernbert(jpeg_image) RETURNS boolean
AS EXTERNAL NAME '/a/b/c/dinky.class' LANGUAGE 'java';
4. CREATE FUNCTION radius(polygon, float) RETURNS polygon
AS EXTERNAL NAME '/a/b/c/dinky.class' LANGUAGE 'java';
5. CREATE FUNCTION overlaps(polygon, polygon) RETURNS boolean
AS EXTERNAL NAME '/a/b/c/dinky.class' LANGUAGE 'java';

Figure 23.6 Method Registration Commands for the Dinky Database

⁵In the case of non-portable compiled code written, for example, in a language like C++, the DBMS uses the operating system's dynamic linking facility to link the method code into the database system so that it can be invoked.

rrype definition statelments for the user-defined atornic data types in the Dinky scherna are given in Figure 23.7.

1. CREATE ABSTRACT DATA TYPE jpeg_image
 (*internallength* = VARIABLE, *input* = jpeg_in, *output* = jpeg_out);
2. CREATE ABSTRACT DATA TYPE polygon
 (*internallength* = VARIABLE, *input* = polyjn, *output* = poly_out);

Figure 23.7 Atomic Type Declaration Commands for Dinky Database

23.5 INHERITANCE

We considered the concept of inheritance in the context of the ER, model in Chapter 2 and discussed how ER diagrams with inheritance were translated into tables. In object-database systems, unlike relational systems, inheritance is supported directly and allows type definitions to be reused and refined very easily. It can be very helpful when modeling similar but slightly different classes of objects. In object-database systems, inheritance can be used in two ways: for reusing and refining types and for creating hierarchies of collections of similar but not identical objects.

23.5.1 Defining Types with Inheritance

In the Dinky database, we model movie theaters with the type `theater_t`. Dinky also wants their database to represent a new marketing technique in the theater business: the *theater-cafe*, which serves pizza and other meals while screening movies. rtheater-cafes require additional information to be represented in the database. In particular, a theater-cafe is just like a theater, but has an additional attribute representing the theater's *llenu*. Inheritance allows us to capture this 'specialization' explicitly in the database design with the followillg DDL staternent:

```
CREATE TYPE theatercafe_t UNDER theater_t (m,enu text);
```

This staternent creates a new type, `theatercafe_t`, which has the same attributes and rmethods as `theater_t`, plus one additional attribute *menu* of type `text`. Methods defined on `theater_t` apply to objects of type `theatercafe_t`, but not vice versa. We say that `theatercafe_t` inherits the attributes and rmethods of `theater_t`.

Note that the illherita,nce rnechanisrll is not rnerely a rnacro to shorten CREATE staternents. It creates an explicit relationship in the database between the subtype (`theatercafe_t`) and the supertype (`theater_t`): *An object of the*

subtype is also considered to be an object of the supertype. This treatment means that any operations that apply to the supertype (methods as well as query operators, such as projection or join) also apply to the subtype. This is generally expressed in the following principle:

The Substitution Principle: Given a supertype *A* and a subtype *B*, it is always possible to substitute an object of type *B* into a legal expression written for objects of type *A*, without producing type errors.

This principle enables easy code reuse because queries and methods written for the supertype can be applied to the subtype without modification.

Note that inheritance can also be used for atomic types, in addition to row types. Given a supertype `image_t` with methods `title()`, `number_of_colors()`, and `display()`, we can define a subtype `thumbnail_image_t` for all images that inherits the methods of `image_t`.

23.5.2 Binding Methods

In defining a subtype, it is sometimes useful to replace a method for the supertype with a new version that operates differently on the subtype. Consider the `image_t` type and the subtype `jpeg_image_t` from the Dinky database. Unfortunately, the `display()` method for standard images does not work for JPEG images, which are specially compressed. Therefore, in creating type `jpeg_image_t`, we write a special `display()` method for JPEG images and register it with the database system using the `CREATE FUNCTION` command:

```
CREATE FUNCTION display(jpeg_image) RETURNS jpeg_image
AS EXTERNAL NAME '/a/b/c/jpeg.class' LANGUAGE 'java.';
```

Registering a new method with the same name as an old method is called **overloading** the method name.

Because of overloading, the system must understand 'which method is intended in a particular expression. For example, when the system needs to invoke the `display()` method on an object of type `jpeg_image_t`, it uses the specialized `display` method. When it needs to invoke `display` on an object of type `image_t` that is not otherwise subtyped, it invokes the standard `display` method. The process of deciding which method to invoke is called **binding** the method to the object. In certain situations, this binding can be done when an expression is parsed (early binding), but in other cases the most specific type of an object cannot be known until runtime, so the method cannot be found until then (late binding). Late binding facilities add flexibility but can make it harder

for the user to reason about the methods that get invoked for a given query expression.

23.5.3 Collection Hierarchies

Type inheritance was invented for object-oriented programming languages, and our discussion of inheritance up to this point differs little from the discussion one might find in a book on an object-oriented language such as C++ or Java.

However, because database systems provide query languages over tabular data sets, the mechanisms from programming languages are enhanced in object databases to deal with tables and queries as well. In particular, in object-relational systems, we can define a table containing objects of a particular type, such as the *Theaters* table in the Dinky schema. Given a new subtype, such as *theatercafe_t*, we would like to create another table *Theater_cafes* to store the information about theater cafes. But, when writing a query over the *Theaters* table, it is sometimes desirable to ask the same query over the *theater_cafes* table; after all, if we project out the additional columns, an instance of the *Theater_cafes* table can be regarded as an instance of the *Theaters* table.

Rather than requiring the user to specify a separate query for each such table, we can inform the system that a new table of the subtype is to be treated as part of a table of the supertype, with respect to queries over the latter table. In our example, we can say

```
CREATE TABLE Theater_Cafes OF TYPE theatercafe_t UNDER Theaters;
```

This statement tells the system that queries over the *Theaters* table should actually be run over all tuples in both the *theaters* and *theater_Cafes* tables. In such cases, if the subtype definition involves method overloading, late-binding is used to ensure that the appropriate methods are called for each tuple.

In general, the *UNDER* clause can be used to generate an arbitrary tree of tables, called a collection hierarchy. Queries over a particular table *T* in the hierarchy are run over all tuples in *T* and its descendants. Sometimes, a user may want the query to run only on *T* and not on the descendants; additional syntax, for example, the keyword *ONLY*, can be used in the query's *FROM* clause to achieve this effect.

23.6 OBJECTS, OIDS, AND REFERENCE TYPES

In object-database systems, data objects can be given an object identifier (*oid*), which is some value that is unique in the database across time. The

OIDs: IBM DB2, Informix IJDS, and Oracle 9i support REF types.

DBMS is responsible for generating aids and ensuring that an oid identifies an object uniquely over its entire lifetime. In SOHle systems, all tuples stored in any table are objects and automatically assigned unique oids; in other systems, a user can specify the tables for which the tuples are to be assigned aids. Often, there are also facilities for generating oids for larger structures (e.g., tables) as well as smaller structures (e.g., instances of data values such as a copy of the integer 5 or a .JPEG image).

An object's aid can be used to refer to it from elsewhere in the data. An oid has a type similar to the type of a pointer in a programming language.

In SQL:1999 every tuple in a table can be given an aid by defining the table in terms of a structured type and declaring that a REF type is associated with it, as in the definition of the Theaters table in Line 4 of Figure 23.1. Contrast this with the definition of the Countries table in Line 7; Countries tuples do not have associated aids. (SQL:1999 also assigns 'oids' to large objects: This is the locator for the object.)

REF types have values that are unique identifiers or aids. SQL:1999 requires that a given REF type must be associated with a specific table. For example, Line 5 of Figure 23.1 defines a column *theater* of type REF(theater_t). The SCOPE clause specifies that items in this column are references to rows in the theaters table, which is defined in Line 4.

23.6.1 Notions of Equality

The distinction between reference types and reference-free structured types raises another issue: the definition of equality. Two objects having the same type are defined to be deep equal if and only if

1. The objects are of atomic type and have the same value.
2. The objects are of reference type and the *deep equals* operator is true for the two referenced objects.
3. The objects are of structured type and the *deep equals* operator is true for all the corresponding subparts of the two objects.

Two objects that have the same reference type are defined to be shallow equal if both refer to the same object (i.e., both references use the same aid). The

definition of shallow equality can be extended to objects of arbitrary type by taking the definition of deep equality and replacing *deep equals* by *shallow equals* in parts (2) and (3).

As an example, consider the complex objects ROW(538, *t89*, 6-3-97, 8-7-97) and ROW(538, *t33*, 6-3-97, 8-7-97), whose type is the type of rows in the table Nowshowing (Line 5 of Figure 23.1). These two objects are not shallow equal because they differ in the second attribute value. Nonetheless, they might be deep equal, if, for instance, the oids *t89* and *t33* refer to objects of type theater_t that have the same value; for example, tuple(54, 'Majestic', '115 King', '2556698').

While two deep equal objects may not be shallow equal, as the example illustrates, two shallow equal objects are always deep equal, of course. The default choice of deep versus shallow equality for reference types is different across systems, although typically we are given syntax to specify either semantics.

23.6.2 Dereferencing Reference Types

An item of reference type REF(basetype) is not the same as the basetype item to which it points. To access the referenced basetype item, a built-in deref() method is provided along with the REF type constructor. For example, given a tuple from the Nowshowing table, one can access the *name* field of the referenced theater_t object with the syntax Nowshowing.deref(theater).name. Since references to tuple types are common, SQL:1999 uses a Java-style arrow operator, which combines a postfix version of the dereference operator with a tuple-type dot operator. The name of the referenced theater can be accessed with the equivalent syntax Nowshowing.theater->name, as in Figure 23.3.

At this point we have covered all the basic type extensions used in the Dinky schema in Figure 23.1. The reader is invited to revisit the schema and examine the structure and content of each table and how the new features are used in the various sample queries.

23.6.3 URLs and DIDs in SQL:1999

It is instructive to note the differences between Internet URIs and the oids in object systems. First, oids uniquely identify a single object over all time (at least, until the object is deleted, when the oid is undefined), whereas the Web resource pointed at by an URL can change over time. Second, oids are simply identifiers and carry no physical information about the objects they identify—this makes it possible to change the storage location of an object without modifying pointers to the object. In contrast, URLs include network

addresses and often file-system names as well, meaning that if the resource identified by the URL has to move to another file or network address, then all links to that resource are either incorrect or require a ‘forwarding’ mechanism. Third, oids are automatically generated by the DBMS for each object, whereas URLs are user-generated. Since users generate URLs, they often embed semantic information into the URL via machine, directory, or file names; this can become confusing if the object's properties change over time.

For URLs, deletions can be troublesome: This leads to the notorious ‘404 Page Not Found’ error. For oids, SQL:1999 allows us to say REFERENCES ARE CHECKED as part of the SCOPE clause and choose one of several actions when a referenced object is deleted. This is a direct extension of referential integrity that covers oids.

23.7 DATABASE DESIGN FOR AN ORDBMS

The rich variety of data types in an ORDBMS offers a database designer many opportunities for a more natural or more efficient design. In this section we illustrate the differences between RDBMS and (O)RDBMS database design through several examples.

23.7.1 Collection Types and ADTs

Our first example involves several space probes, each of which continuously records a video. A single video stream is associated with each probe, and while this stream was collected over a certain time period, we assume that it is now a complete object associated with the probe. During the time period over which the video was collected, the probe's location was periodically recorded (such information can easily be piggy-backed onto the header portion of a video stream conforming to the MPEG standard). The information associated with a probe has three parts: (1) a *probe ID* that identifies a probe uniquely, (2) a *video stream*, and (3) a *location sequence* of $\langle \text{time}, \text{location} \rangle$ pairs. What kind of a database schema should we use to store this information?

An RDBMS Database Design

In an RDBMS, we must store each video stream as a BLOB and each location sequence as tuples in a table. A possible RDBMS database design follows:

```
Probes(pid: integer, time: timestamp, lat: real, long: real,
      camera: string, video: BLOB)
```

There is a single table named Probes and it has several rows for each probe. Each of these rows has the same *pid*, *camera* and *video* values, but different *time*, *tat*, and *long* values. (We have used latitude and longitude to denote location.) The key for this table can be represented as a functional dependency: $IJTLN \rightarrow CV$, where *N* stands for longitude. There is another dependency: $P \rightarrow CV$. This relation is therefore not in BCNF; indeed, it is not even in 3NF. We can decompose Probes to obtain a BCNF schema:

Probes_Loc(*pid*: integer, *time*: timestamp, *tat*: real, *long*: real)
Probes_Video(*pid*: integer, *camera*: string, *video*: BLOB)

This design is about the best we can achieve in an RDBMS. However, it suffers from several drawbacks.

First, representing videos as BLOBs means that we have to write application code in an external language to manipulate a video object in the database. Consider this query: "For probe 10, display the video recorded between 1:10 P.M. and 1:15 P.M. on May 10 1996." We must retrieve the entire video object associated with probe 10, recorded over several hours, to display a segment recorded over five minutes.

Next, the fact that each probe has an associated sequence of location readings is obscured, and the sequence information associated with a probe is dispersed across several tuples. A third drawback is that we are forced to separate the video information from the sequence information for a probe. These limitations are exposed by queries that require us to consider all the information associated with each probe; for example, "For each probe, print the earliest time at which it recorded, and the camera type." This query now involves a join of Probes_Loc and Probes_Video on the *pid* field.

An ORDBMS Database Design

An ORDBMS supports a much better solution. First, we can store the video as an ADT object and write methods that capture any special manipulation we wish to perform. Second, because we are allowed to store structured types such as lists, we can store the location sequence for a probe in a single tuple, along with the video information. This layout eliminates the need for joins in queries that involve both the sequence and video information. An ORDBMS design for our example consists of a single relation called Probes_AllInfo:

Probes_AllInfo(*pid*: integer, *locseq*: location_seq, *camera*: string;
video: mpeg_stream)

This definition involves two new types, `location_seq` and `mpeg_stream`. The `mpeg_stream` type is defined as an ADT, with a lmethod `display()` that takes a start time and an end time and displays the portion of the video recorded during that interval. This rmethod can be implemented efficiently by looking at the total recording duration and the total length of the video and interpolating to extract the segment recorded during the interval specified in the query.

Our first query in extended SQL using this `display` lmethod follows. We now retrieve only the required segment of the video rather than the entire video.

```
SELECT display(P.video, 1:10 Po M May 10 1996, 1:15 Po M May 10 1996)
FROM   Probes_AllInfo P
WHERE  Popid = 10
```

Now consider the `location_seq` type. We could define it as a `list` type, containing a list of ROW type objects:

```
CREATE TYPE location_seq listof
      (row (time: timestamp, lat: real, long: real))
```

Consider the `locseq` field in a row for a given probe. This field contains a list of rows, each of which has three fields. If the ORDBMS implements collection types in their full generality, we should be able to extract the `time` column from this list to obtain a list of timestamp values and apply the `MIN` aggregate operator to this list to find the earliest time at which the given probe recorded. Such support for collection types would enable us to express our second query thus:

```
SELECT   P.piel, MIN(P.locseq.time)
FROM     Probes_AllInfo P
```

Current ORDBMSs are not as general and clean as this example query suggests. For instance, the system may not recognize that projecting the `time` column from a list of rows gives us a list of timestamp values; or the system may allow us to apply an aggregate operator only to a table and not to a nested list value.

Continuing with our example, we may want to do specialized operations on our location sequences that go beyond the standard aggregate operators. For instance, we may want to define a lmethod that takes a time interval and computes the distance traveled by the probe during this interval. The code for this rmethod must understand details of a probe's trajectory and geospatial coordinate systems. For these reasons, we might choose to define `location_seq` as an ADT.

Clearly, an (ideal) ORDBMS gives us many useful design options that are not available in an RDBMS.

23.7.2 Object Identity

We now discuss some of the consequences of using reference types or aids. The use of aids is especially significant when the size of the object is large, either because it is a structured data type or because it is a big object such as an image.

Although reference types and structured types seem similar, they are actually quite different. For example, consider a structured type `my_theater tuple (integer, name text, address text, phone text)` and the reference type `theater ref (theater_t)` of Figure 23.1. There are important differences in the way that database updates affect these two types:

- **Deletion:** Objects with references can be affected by the deletion of objects that they reference, while reference-free structured objects are not affected by deletion of other objects. For example, if the `Theaters` table were dropped from the database, an object of type `theater` might change value to *null*, because the `theater_t` object it refers to has been deleted, while a similar object of type `my_theater` would not change value.
- **Update:** Objects of reference types change value if the referenced object is updated. Objects of reference-free structured types change value only if updated directly.
- **Sharing versus Copying:** An identified object can be referenced by multiple reference-type items, so that each update to the object is reflected in many places. To get a similar effect in reference-free types requires updating all 'copies' of an object.

There are also important storage distinctions between reference types and non-reference types, which might affect performance:

- **Storage Overhead:** Storing copies of a large value in multiple structured type objects may use much more space than storing the value once and referring to it elsewhere through reference type objects. This additional storage requirement can affect both disk usage and buffer management (if many copies are accessed at once).
- **Clustering:** The subparts of a structured object are typically stored together on disk. Objects with references may point to other objects that are far away on the disk, and the disk arm may require significant movement

OIDs and Referential Integrity: In SQL:1999, all the oids that appear in a column of a relation are required to reference the same target relation. This ‘scoping’ makes it possible to check oid references for ‘referential integrity’ just like foreign key references are checked. While current ORDBMS products supporting oids do not support such checks, it is likely that they will in future releases. This will make it much safer to use aids.

to assemble the object and its references together. Structured objects can thus be more efficient than reference types if they are typically accessed in their entirety.

Many of these issues also arise in traditional programming languages such as C or Pascal, which distinguish between the notions of referring to objects *by value* and *by reference*. In database design, the choice between using a structured type or a reference type typically includes consideration of the storage costs, clustering issues, and the effect of updates.

Object Identity versus Foreign Keys

Using an oid to refer to an object is similar to using a foreign key to refer to a tuple in another relation but not quite the same. An oid can point to an object of `theater_t` that is stored *anywhere* in the database, even in a field, whereas a foreign key reference is constrained to point to an object in a particular referenced relation. This restriction makes it possible for the DBMS to provide much greater support for referential integrity than for arbitrary aid pointers. In general, if an object is deleted while there are still oid-pointers to it, the best the DBMS can do is to recognize the situation by maintaining a reference count. (Even this limited support becomes impossible if oids can be copied freely.) Therefore, the responsibility for avoiding dangling references rests largely with the user if oids are used to refer to objects. This burdensome responsibility suggests that we should use oids with great caution and use foreign keys instead whenever possible.

23.7.3 Extending the ER Model

The ER model, as described in Chapter 2, is not adequate for ORDBMS design. We have to use an extended ER model that supports structured attributes (i.e., sets, lists, arrays as attribute values), distinguishes whether entities have object ids, and allows us to model entities whose attributes include methods. We illustrate these comments using an extended ER diagram to describe the

space probe data in Figure 23.8; our notational conventions are ad hoc and only for illustrative purposes.

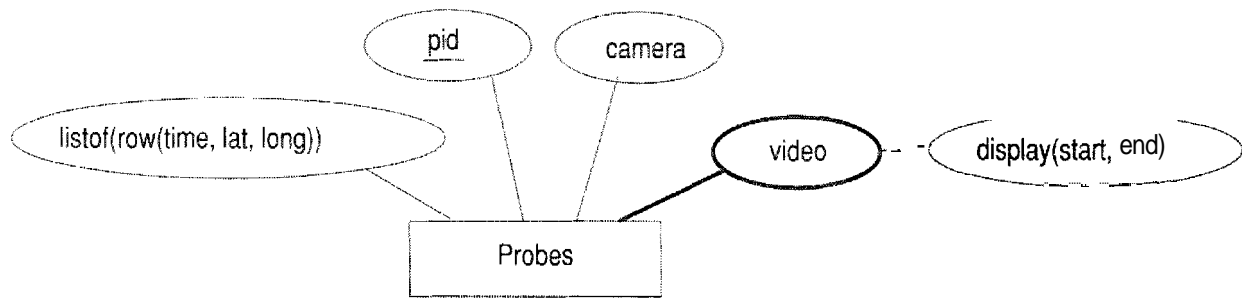


Figure 23.8 The Space Probe Entity Set

The definition of Probes in Figure 23.8 has two new aspects. First, it has a structured-type attribute `listof(row(time, lat, long))`; each value assigned to this attribute in a Probes entity is a list of tuples with three fields. Second, Probes has an attribute called *video* that is an abstract data type object, which is indicated by a dark oval for this attribute with a dark line connecting it to Probes. Further, this attribute has an 'attribute' of its own, which is a method of the *ADT*.

Alternatively, we could model each video as an entity by using an entity set called Videos. The association between Probes entities and Videos entities could then be captured by defining a relationship set that links them. Since each video is collected by precisely one probe and every video is collected by some probe, this relationship can be maintained by simply storing a reference to a probe object with each Videos entity; this technique is essentially the second translation approach from ER diagrams to tables discussed in Section 2.4.1.

If we also make Videos a weak entity set in this alternative design, we can add a referential integrity constraint that causes a Videos entity to be deleted when the corresponding Probes entity is deleted. More generally, this alternative design illustrates a strong similarity between storing references to objects and foreign keys; the foreign key mechanism achieves the same effect as storing oids, but in a controlled manner. If oids are used, the user must ensure that there are no dangling references when an object is deleted, with very little support from the DBMS.

Finally, we note that a significant extension to the ER model is required to support the design of nested collections. For example, if a location sequence is modeled as an entity, and we want to define an attribute of Probes that contains a set of such entities, there is no way to do this without extending the ER model. We do not discuss this point further at the level of ER diagrams, but consider an example next that illustrates when to use a nested collection.

23.7.4 Using Nested Collections

Nested collections offer great modeling power but also raise difficult design decisions. Consider the following way to model location sequences (other information about probes is omitted here to simplify the discussion):

Probes1(pid: integer, locseq: location_seq)

This is a good choice if the important queries in the workload require us to look at the location sequence for a particular probe, as in the query "For each probe, print the earliest time at which it recorded and the caluera type." On the other hand, consider a query that requires us to look at all location sequences: "Find the earliest time at which a recording exists for *lat*=5, *long*=90." This query can be answered more efficiently if the following schema is used:

Probes2(pid: integer, time: timestamp, lat: real, long: real)

The choice of schema must therefore be guided by the expected workload (as always). As another example, consider the following schema:

Can_Teach1(cid: integer, teacheTs: setof(ssn: string), sal: integer)

If tuples in this table are to be interpreted as "Course *cid* can be taught by any of the teachers in the *teacheTs* field, at a cost *sal*." then we have the option of using the following schema instead:

CanLTeach2(cid: integer, teachCT ssn: string, sal: integer)

A choice between these two alternatives can be made based on how we expect to query this table. On the other hand, suppose that tuples in CanLTeach1 are to be interpreted as "Course *cid* can be taught by the *tearnteacheTs*, at a combined cost of *sal*." CanLTeach2 is no longer a viable alternative. If we wanted to flatten Can_Teach1, we would have to use a separate table to encode teams:

Can_Teach2(cid: integer, team_id: oid, sal: integer)
Teams(tid: oid, ssn: string)

As these examples illustrate, nested collections are appropriate in certain situations, but this feature can easily be misused; nested collections should therefore be used with care.

23.8 ORDBMS IMPLEMENTATION CHALLENGES

The enhanced functionality of ORDBMSs raises several implementation challenges. Some of these are well understood and solutions have been implemented in products; others are subjects of current research. In this section we examine a few of the key challenges that arise in implementing an efficient, fully functional ORDBMS. Many more issues are involved than those discussed here; the interested reader is encouraged to revisit the previous chapters in this book and consider whether the implementation techniques described there apply naturally to ORDBMSs or not.

23.8.1 Storage and Access Methods

Since object-relational databases store new types of data, ORDBMS implementors need to revisit some of the storage and indexing issues discussed in earlier chapters. In particular, the system must efficiently store ADT objects and structured objects and provide efficient indexed access to both.

Storing Large ADT and Structured Type Objects

Large ADT objects and structured objects complicate the layout of data on disk. This problem is well understood and has been solved in essentially all ORDBMSs and OODBMSs. We present some of the main issues here.

User-defined ADTs can be quite large. In particular, they can be bigger than a single disk page. Large ADTs, like BLOBs, require special storage, typically in a different location on disk from the tuples that contain them. Disk-based pointers are maintained from the tuples to the objects they contain.

Structured objects can also be large, but unlike ADT objects, they often vary in size during the lifetime of a database. For example, consider the *stars* attribute of the *films* table in Figure 23.1. As the years pass, some of the 'bit actors' in an old movie may become famous.⁶ When a bit actor becomes famous, Dinky might want to advertise his or her presence in the earlier films. This involves an insertion into the *stars* attribute of an individual tuple in *films*. Because these bulk attributes can grow arbitrarily, flexible disk layout mechanisms are required.

⁶A famous example is Marilyn Monroe, who had a bit part in the Bette Davis classic *All About Eve*.

An additional complication arises with array types. Traditionally, array elements are stored sequentially on disk in a row-by-row fashion; for example

$$A_{11}, \dots, A_{1n}, A_{21}, \dots, A_{2n}, \dots, A_{m1}, \dots, A_{mn}$$

However, queries may often request subarrays that are not stored contiguously on disk (e.g., $A_{11}, A_{21}, \dots, A_{rn1}$). Such requests can result in a very high I/O cost for retrieving the subarray. To reduce the number of I/Os required, arrays are often broken into contiguous *chunks*, which are then stored in some order on disk. Although each chunk is some contiguous region of the array, chunks need not be row-by-row or column-by-column. For example, a chunk of size 4 might be $A_{11}, A_{12}, A_{21}, A_{22}$, which is a square region if we think of the array as being arranged row-by-row in two dimensions.

Indexing New Types

One important reason for users to place their data in a database is to allow for efficient access via indexes. Unfortunately, the standard RDBMS index structures support only equality conditions (B+ trees and hash indexes) and range conditions (B+ trees). An important issue for ORDBMSs is to provide efficient indexes for ADT methods and operators on structured objects.

Many specialized index structures have been proposed by researchers for particular applications such as cartography, genome research, multimedia repositories, Web search, and so on. An ORDBMS company cannot possibly implement every index that has been invented. Instead, the set of index structures in an ORDBMS should be user-extensible. Extensibility would allow an expert in cartography, for example, to not only register an ADT for points on a map (i.e., latitude-longitude pairs) but also implement an index structure that supports natural map queries (e.g., the R-tree, which matches conditions such as “Find me all theaters within 100 miles of Andorra”). (See Chapter 28 for more on R-trees and other spatial indexes.)

One way to make the set of index structures extensible is to publish an *access method interface* that lets users implement an index structure *outside* the DBMS. The index and data can be stored in a file system and the DBMS simply issues the *open*, *next*, and *close* iterator requests to the user’s external index code. Such functionality makes it possible for a user to connect a DBMS to a Web search engine, for example. A main drawback of this approach is that data in an external index is not protected by the DBMS’s support for concurrency and recovery. An alternative is for the ORDBMS to provide a generic ‘template’ index structure that is sufficiently general to encompass most index structures that users might invent. Because such a structure is implemented within the DBMS, it can support high concurrency and recovery. The *Gener-*

alized Search Tree (GiST) is such a structure. It is a template index structure based on B+ trees, which allows most of the tree index structures invented so far to be implemented with only a few lines of user-defined ADT code.

23.8.2 Query Processing

ADTs and structured types call for new functionality in processing queries in ORDBMSs. They also change a number of assumptions that affect the efficiency of queries. In this section we look at two functionality issues (user-defined aggregates and security) and two efficiency issues (method caching and pointer swizzling).

User-Defined Aggregation Functions

Since users are allowed to define new methods for their ADTs, it is not unreasonable to expect them to want to define new aggregation functions for their ADTs as well. For example, the usual SQL aggregates—COUNT, SUM, MIN, MAX, AVG—are not particularly appropriate for the image type in the Dinky schema.

Most ORDBMSs allow users to register new aggregation functions with the system. To register an aggregation function, a user must implement three methods, which we call *initialize*, *iterate*, and *terminate*. The *initialize* method initializes the internal state for the aggregation. The *iterate* method updates that state for every tuple seen, while the *terminate* method computes the aggregation result based on the final state and then cleans up. As an example, consider an aggregation function to compute the second-highest value in a field. The *initialize* call would allocate storage for the top two values, the *iterate* call would compare the current tuple's value with the top two and update the top two as necessary, and the *terminate* call would delete the storage for the top two values, returning a copy of the second-highest value.

Method Security

ADTs give users the power to add code to the DBMS; this power can be abused. A buggy or malicious ADT method can bring down the database server or even corrupt the database. The DBMS must have mechanisms to prevent buggy or malicious user code from causing problems. It may make sense to override these mechanisms for efficiency in production environments with vendor-supplied methods. However, it is important for the mechanisms to exist, if only to support debugging of ADT methods; otherwise method writers

would have to write bug-free code before registering their rmethods with the DBMS--not a very forgiving programming environment.

One mechanism to prevent problems is to have the user rmethods be *interpreted* rather than *compiled*. The DBMS can check that the rmethod is well behaved either by restricting the power of the interpreted language or by ensuring that each step taken by a rmethod is safe before executing it. Typical interpreted languages for this purpose include Java and the procedural portions of SQL:1999.

An alternative mechanism is to allow user methods to be compiled from a general-purpose programming language, such as C++, but to run those rmethods in a different address space than the DBMS. In this case, the DBMS sends explicit interprocess communications (IPCs) to the user method, which sends IPCs back in return. This approach prevents bugs in the user methods (e.g., stray pointers) from corrupting the state of the DBMS or database and prevents malicious methods from reading or modifying the DBMS state or database as well. Note that the user writing the method need not know that the DBMS is running the method in a separate process: The user code can be linked with a 'wrapper' that turns method invocations and return values into IPCs.

Method Caching

User-defined ADT methods can be very expensive to execute and can account for the bulk of the time spent in processing a query. During query processing, it may make sense to cache the results of methods, in case they are invoked multiple times with the same argument. Within the scope of a single query, one can avoid calling a method twice on duplicate values in a column by either sorting the table on that column or using a hash-based scheme such as that used for aggregation (see Section 14.6). An alternative is to maintain a *cache* of method inputs and matching outputs as a table in the database. Then, to find the value of a method on particular inputs, we essentially join the input tuples with the cache table. These two approaches can also be combined.

Pointer Swizzling

In some applications, objects are retrieved into memory and accessed frequently through their oids; dereferencing must be implemented very efficiently. Some systems maintain a table of oids of objects that are (currently) in memory. When an object () is brought into memory, they check each oid contained in θ and replace oids of in-memory objects by in-memory pointers to those objects. This technique, called *pointer swizzling*, makes references to in-memory objects very fast. The downside is that when an object is paged out,

23.9 OODBMS

In the introduction of this chapter, we defined an OODBMS as a programming language with support for persistent objects. While this definition reflects the origins of OODBMSs accurately, and to a certain extent the implementation focus of OODBMSs, the fact that OODBMSs support *collection types* (see Section 23.2.1) makes it possible to provide a query language over collections. Indeed, a standard has been developed by the Object Database Management Group and is called Object Query Language.

OQL is similar to SQL, with a SELECT--FROM--WHERE-style syntax (even GROUP BY, HAVING, and ORDER BY are supported) and many of the proposed SQL:1999 extensions. Notably, OQL supports structured types, including sets, bags, arrays, and lists. The OQL treatment of collections is more uniform than SQL:1999 in that it does not give special treatment to collections of rows; for example, OQL allows the aggregate operation COUNT to be applied to a list to compute the length of the list. OQL also supports reference types, path expressions, ADTs and inheritance, type extents, and SQL-style nested queries. There is also a standard Data Definition Language for OODBMSs (Object Data Language, or ODL) that is similar to the DDL subset of SQL but supports the additional features found in OODBMSs, such as ADT definitions.

23.9.1 The ODMG Data Model and ODL

The ODMG data model is the basis for an OODBMS, just like the relational data model is the basis for an RDBMS. A database contains a collection of objects, which are similar to entities in the ER model. Every object has a unique id, and a database contains collections of objects with similar properties; such a collection is called a class.

The properties of a class are specified using ODL and are of three kinds: attributes, relationships, and methods. Attributes have an atomic type or a structured type. ODL supports the **set**, **bag**, **list**, **array**, and **struct** type constructors; these are just **setof**, **bagof**, **listof**, **ARRAY**, and **ROW** in the terminology of Section 23.2.1.

Relationships have a type that is either a reference to an object or a collection of such references. A relationship captures how an object is related to one or more objects of the same class or of a different class. A relationship in the ODMG model is really just a binary relationship in the sense of the ER model. A relationship has a corresponding inverse relationship; intuitively, it is the relationship 'in the other direction.' For example, if a movie is being

Class = Interface + Implementation: Properly speaking, a class consists of an interface together with an implementation of the interface. An ODL interface definition is implemented in an OODBMS by translating it into declarations of the object-oriented language (e.g., C++, Smalltalk or Java) supported by the OODBMS. If we consider C++, for instance, there is a library of classes that implement the ODL constructs. There is also an Object Manipulation Language (OML) specific to the programming language (in our example, C++), which specifies how database objects are manipulated in the programming language. The goal is to seamlessly integrate the programming language and the database features.

shown at several theaters and each theater shows several movies, we have two relationships that are inverses of each other: *shownAt* is associated with the class of movies and is the set of theaters at which the given movie is being shown, and *nowShowing* is associated with the class of theaters and is the set of movies being shown at that theater.

Methods are functions that can be applied to objects of the class. There is no analog to methods in the ER or relational models.

The keyword **interface** is used to define a class. For each interface, we can declare an extent, which is the name for the current set of objects of that class. The extent is analogous to the instance of a relation and the interface is analogous to the schema. If the user does not anticipate the need to work with the set of objects of a given class—it is sufficient to manipulate individual objects—the extent declaration can be omitted.

The following ODL definitions of the Movie and Theater classes illustrate these concepts. (While these classes bear some resemblance to the Dinky database schema, the reader should not look for an exact parallel, since we have modified the example to highlight ODL features.)

```
interface Movie
  (extent Movies key movieName)
  { attribute date start;
    attribute date end;
    attribute string movieName;
    relationship Set<Theater> shownAt inverse Theater::nowShowing;
  }
```

The collection of database objects whose class is Movie is called Movies. No two objects in Movies have the same *movieName* value, as the key declaration

indicates. Each *Movie* is shown at a set of theaters and is shown during the specified period. (It would be more realistic to associate a different period with each theater, since a movie is typically played at different theaters over different periods. While we can define a class that captures this detail, we have chosen a simpler definition for our discussion.) A theater is an object of class *Theater*, defined as:

```
interface Theater
    (extent Theaters key theaterName)
    { attribute string theaterName;
      attribute string address;
      attribute integer ticketPrice;
      relationship Set(Movie) nowShowing inverse .IMovie::shownAt;
      float numshowing() raises(errorCountingMovies);
    }
```

Each theater shows several movies and charges the same ticket price for every movie. Observe that the *shownAt* relationship of *Movie* and the *nowShowing* relationship of *Theater* are declared to be inverses of each other. *Theater* also has a method *numshowing()* that can be applied to a theater object to find the number of movies being shown at that theater.

ODL also allows us to specify inheritance hierarchies, as the following class definition illustrates:

```
interface SpecialShow extends IMovie
    (extent SpecialShows)
    { attribute integer maximumAttendees;
      attribute string benefitCharity;
    }
```

An object of class *SpecialShow* is an object of class *Movie*, with some additional properties, as discussed in Section 23.5.

23.9.2 OQL

The ODMG query language OQL was deliberately designed to have syntax similar to SQL to make it easy for users familiar with SQL to learn OQL. Let us begin with a query that finds pairs of *Movies* and theaters such that the movie is shown at the theater and the theater is showing more than one movie:

```
SELECT Inname: M.movieName, tname: T.theaterName
```



```
FROM    Movies M, Theaters T
WHERE    T.numShowing() > 1
```

The SELECT clause indicates how we can give aliases to fields in the result: The two result fields are called *name* and *tname*. The part of this query that differs from SQL is the FROM clause. The variable *M* is bound in turn to each movie in the extent *Movies*. For a given movie *M*, we bind the variable *T* in turn to each theater in the collection *M.shownAt*. Thus, the use of the path expression *M.shownAt* allows us to easily express a nested query. The following query illustrates the grouping construct in OQL:

```
SELECT    T.ticketPrice,
          avgNum: AVG(SELECT P.T.numShowing() FROM partition P)
FROM      Theaters T
GROUP BY  T.ticketPrice
```

For each ticket price, we create a group of theaters with that ticket price. This group of theaters is the partition for that ticket price, referred to using the OQL keyword *partition*. In the SELECT clause, for each ticket price, we compute the average number of movies shown at theaters in the partition for that ticketPrice. OQL supports an interesting variation of the grouping operation that is missing in SQL:

```
SELECT    low, high,
          avgNum: AVG(SELECT P.T.numShowing() FROM partition P)
FROM      Theaters T
GROUP BY  low: T.ticketPrice < 5, high: T.ticketPrice >= 5
```

The GROUP BY clause now creates just two partitions called *low* and *high*. Each theater object *T* is placed in one of these partitions based on its ticket price. In the SELECT clause, *low* and *high* are boolean variables, exactly one of which is true in any given output tuple; *partition* is instantiated to the corresponding partition of theater objects. In our example, we get two result tuples. One of them has *low* equal to true and *avgNum* equal to the average number of movies shown at theaters with a low ticket price. The second tuple has *high* equal to true and *avgNum* equal to the average number of movies shown at theaters with a high ticket price.

The next query illustrates OQL support for queries that return collections other than set and multiset:

```
(SELECT    r.theaterName
FROM      Theaters T
ORDER BY  T.ticketPrice DESC) [0:4]
```

The `ORDER BY` clause makes the result a list of theater names ordered by ticket price. The elements of a list can be referred to by position, starting with position 0. Therefore, the expression `[0:4]` extracts a list containing the names of the five theaters with the highest ticket prices.

OQL also supports `DISTINCT`, `HAVING`, explicit nesting of subqueries, view definitions, and other SQL features.

23.10 COMPARING RDBMS, OODBMS, AND ORDBMS

Now that we have covered the main object-oriented DBMS extensions, it is time to consider the two main variants of object-databases, OODBMSs and ORDBMSs, and compare them with RDBMSs. Although we presented the concepts underlying object-databases, we still need to define the terms OODBMS and ORDBMS.

An **ORDBMS** is a relational DBMS with the extensions discussed in this chapter. (Not all ORDBMS systems support all the extensions in the general form that we have discussed them, but our concern in this section is the paradigm itself rather than specific systems.) An **OODBMS** is a programming language with a type system that supports the features discussed in this chapter and allows any data object to be persistent; that is, to survive across different program executions. Many current systems conform to neither definition entirely but are much closer to one or the other and can be classified accordingly.

23.10.1 RDBMS versus ORDBMS

Comparing an RDBMS with an OODBMS is straightforward. An RDBMS does not support the extensions discussed in this chapter. The resulting simplicity of the data model makes it easier to optimize queries for efficient execution, for example. A relational system is also easier to use because there are fewer features to master. On the other hand, it is less versatile than an OODBMS.

23.10.2 OODBMS versus ORDBMS: Similarities

OODBMSs and OJBMSs both support user-defined ADTs, structured types, object identity and reference types, and inheritance. Both support a query language for manipulating collection types. ORDBMSs support an extended form of SQL, and OODBMSs support ODL/OQL. The similarities are by no means accidental. ORDBMSs consciously try to add OODBMS features to an RDBMS, and OODBMSs in turn have developed query languages based on

relational query languages. Both OODBMSs and ORDBMSs provide DBMS functionality such as concurrency control and recovery.

23.10.3 OODBMS versus ORDBMS: Differences

The fundamental difference is really a philosophy that is carried all the way through: OODBMSs try to add DBMS functionality to a programming language, whereas ORDBMSs try to add richer data types to a relational DBMS. Although the two kinds of object-databases are converging in terms of functionality, this difference in their underlying philosophy (and for most systems, their implementation approach) has important consequences in terms of the issues emphasized in the design of these DBMSs and the efficiency with which various features are supported, as the following comparison indicates:

- OODBMSs aim to achieve seamless integration with a programming language such as C++, Java, or Smalltalk. Such integration is not an important goal for an ORDBMS. SQL:1999, like SQL-92, allows us to embed SQL commands in a host language, but the interface is very evident to the SQL programmer. (SQL:1999 also provides extended programming language constructs of its own, as we saw in Chapter 6.)
- An OODBMS is aimed at applications where an object-centric viewpoint is appropriate; that is, typical user sessions consist of retrieving a few objects and working on them for long periods, with related objects (e.g., objects referenced by the original objects) fetched occasionally. Objects may be extremely large and may have to be fetched in pieces; therefore, attention must be paid to buffering parts of objects. It is expected that most applications can cache the objects they require in memory, once the objects are retrieved from disk. Therefore, considerable attention is paid to making references to in-memory objects efficient. Transactions are likely to be of very long duration and holding locks until the end of a transaction may lead to poor performance; therefore, alternatives to Two-Phase Locking must be used.

An ORDBMS is optimized for applications in which large data collections are the focus, even though objects may have rich structure and be fairly large. It is expected that applications will retrieve data from disk extensively and optimizing disk access is still the main concern for efficient execution. Transactions are assumed to be relatively short and traditional RDBMS techniques are typically used for concurrency control and recovery.

- The query facilities of OQL are not supported efficiently in most OODBMSs, whereas the query facilities are the centerpiece of an ORDBMS. To some extent, this situation is the result of different concentrations of effort in the development of these systems. To a significant extent, it is also a

consequence of the systems' being optimized for very different kinds of applications.

23.11 REVIEW QUESTIONS

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

- Consider the extended Dinky example from Section 23.1. Explain how it motivates the need for each of the following object-database features: *User-defined structured types*, *abstract data types (ADTs)*, *inheritance*, and *object identity*. (Section 23.1)
- What are *structured data types*? What are *collection types*, in particular? Discuss the extent to which these concepts are supported in SQL:1999. What important type constructors are missing? What are the limitations on the ROW and ARRAY constructors? (Section 23.2)
- What kinds of operations should be provided for each of the structured data types? To what extent is such support included in SQL:1999? (Section 23.3)
- What is an *abstract data type*? How are methods of an abstract data type defined in an external programming language? (Section 23.4)
- Explain *inheritance* and how new types (called *subtypes*) extend existing types (called *supertypes*). What are *method overloading* and *late binding*? What is a *collection hierarchy*? Contrast this with inheritance in programming languages. (Section 23.5)
- How is an *object identifier (oid)* different from a record id in a relational DBMS? How is it different from a URL? What is a *reference type*? Define *deep* and *shallow* equality and illustrate them through an example. (Section 23.6)
- The multitude of data types in an ORDBMS allows us to design a more natural and efficient database schema but introduces some new design choices. Discuss ORDBMS database design issues and illustrate your discussion using an example application. (Section 23.7)
- Implementing an ORDBMS brings new challenges. The system must store large ADTs and structured types that might be very large. Efficient and extensible indexing mechanisms must be provided. Examples of new functionality include *user-defined aggregation functions* (we can define new aggregation functions for our ADTs) and *method security* (the system has to prevent user-defined methods from compromising the security of the DBMS). Examples of new techniques to increase performance include

method caching and *pointer swizzling*. The optimizer must know about the new functionality and use it appropriately. Illustrate each of these challenges through an example. (Section 23.8)

- Compare OODBMSs with ORDBMSs. In particular, compare OQL and SQL:1999 and discuss the underlying data model. (Sections 23.9 and 23.10)

EXERCISES

Exercise 23.1 Briefly answer the following questions:

1. What are the new kinds of data types supported in object-database systems? Give an example of each and discuss how the example situation would be handled if only an RDBMS were available.
2. What must a user do to define a new ADT?
3. Allowing users to define methods can lead to efficiency gains. Give an example.
4. What is late binding of methods? Give an example of inheritance that illustrates the need for dynamic binding.
5. What are collection hierarchies? Give an example that illustrates how collection hierarchies facilitate querying.
6. Discuss how a DBMS exploits encapsulation in implementing support for ADTs.
7. Give an example illustrating the nesting and unnesting operations.
8. Describe two objects that are deep equal but not shallow equal or explain why this is not possible.
9. Describe two objects that are shallow equal but not deep equal or explain why this is not possible.
10. Compare RDBMSs with ORDBMSs. Describe an application scenario for which you would choose an RDBMS and explain why. Similarly, describe an application scenario for which you would choose an ORDBMS and explain why.

Exercise 23.2 Consider the Dinky schema shown in Figure 23.1 and all related methods defined in the chapter. Write the following queries in SQL:1999:

1. How many films were shown at theater *tno* = 5 between January 1 and February 1 of 2002?
2. What is the lowest budget for a film with at least two stars?
3. Consider theaters at which a film directed by Steven Spielberg started showing on January 1, 2002. For each such theater, print the names of all countries within a 100-mile radius. (You can use the *overlap* and *radius* methods illustrated in Figure 23.2.)

Exercise 23.3 In a company database, you need to store information about employees, departments, and children of employees. For each employee, identified by *ssn*, you must record *years* (the number of years that the employee has worked for the company), *phone*, and *photo* information. There are two subclasses of employees: contract and regular. Salary is computed by invoking a method that takes *years* as a parameter; this method has a different

implementation for each subclass. Further, for each regular employee, you must record the name and age of every child. The most common queries involving children are similar to “Find the average age of Bob’s children” and “Print the names of all of Bob’s children.”

A photo is a large image object and can be stored in one of several image formats (e.g., gif, jpeg). You want to define a *display* method for image objects; display must be defined differently for each image format. For each department, identified by *dno*, you must record *dname*, *budget*, and *workers* information. *Workers* is the set of employees who work in a given department. Typical queries involving workers include, “Find the average salary of all workers (across all departments).”

1. Using extended SQL, design an ORDBMS schema for the company database. Show all type definitions, including method definitions.
2. If you have to store this information in an RDBMS, what is the best possible design?
3. Compare the ORDBMS and RDBMS designs.
4. If you are told that a common request is to display the images of all employees in a given department, how would you use this information for physical database design?
5. If you are told that an employee’s image must be displayed whenever any information about the employee is retrieved, would this affect your schema design?
6. If you are told that a common query is to find all employees who look similar to a given image and given code that lets you create an index over all images to support retrieval of similar images, what would you do to utilize this code in an ORDBMS?

Exercise 23.4 ORDBMSs need to support efficient access over collection hierarchies. Consider the collection hierarchy of Theaters and Theater_cafes presented in the Dinky example. In your role as a DBMS implementor (not a DBA), you must evaluate three storage alternatives for these tuples:

- All tuples for all kinds of theaters are stored together on disk in an arbitrary order.
 - All tuples for all kinds of theaters are stored together on disk, with the tuples that are from Theater_cafes stored directly after the last of the non-cafe tuples.
 - Tuples from Theater_cafes are stored separately from the rest of the (non-cafe) theater tuples.
1. For each storage option, describe a mechanism for distinguishing plain theater tuples from Theater_cafe tuples.
 2. For each storage option, describe how to handle the insertion of a new non-cafe tuple.
 3. Which storage option is most efficient for queries over all theaters? Over just Theater_cafes? In terms of the number of I/Os, how much more efficient is the best technique for each type of query compared to the other two techniques?

Exercise 23.5 Different ORDBMSs use different techniques for building indexes to evaluate queries over collection hierarchies. For our Dinky example, to index theaters by name there are two common options:

- Build one B+ tree index over Theaters.name and another B+ tree index over Theater_cafes.name.
- Build one B+ tree index over the union of Theaters.name and Theater_cafes.name.

1. Describe how to efficiently evaluate the following query using each indexing option (this query is over all kinds of theater tuples):

```
SELECT * FROM Theaters T WHERE T.name = 'Majestic'
```

Give an estimate of the number of I/Os required in the two different scenarios, assuming there are 1 million standard theaters and 1000 theater-cafes. Which option is more efficient?

2. Perform the same analysis over the following query:

```
SELECT * FROM Theater-cafes T WHERE T.name = 'Majestic'
```

3. For clustered indexes, does the choice of indexing technique interact with the choice of storage options? For unclustered indexes?

Exercise 23.6 Consider the following query:

```
SELECT thumbnail(image)
FROM Images I
```

Given that the *image* column may contain duplicate values, describe how to use hashing to avoid computing the *thumbnail* function more than once per distinct value in processing this query.

Exercise 23.7 You are given a two-dimensional, $n \times n$ array of objects. Assume that you can fit 100 objects on a disk page. Describe a way to layout (chunk) the array onto pages so that retrievals of square $m \times m$ subregions of the array are efficient. (Different queries request subregions of different sizes, i.e., different m values, and your arrangement of the array onto pages should provide good performance, on average, for all such queries.)

Exercise 23.8 An ORDBMS optimizer is given a single-table query with n expensive selection conditions, $\sigma_n(\dots(\sigma_1(T)))$. For each condition σ_i , the optimizer can estimate the cost c_i of evaluating the condition on a tuple and the reduction factor of the condition r_i . Assume that there are t tuples in T .

1. How many tuples appear in the output of this query?
2. Assuming that the query is evaluated as shown (without reordering selections), what is the total cost of the query? Be sure to include the cost of scanning the table and applying the selections.
3. In Section 23.8.2, it was asserted that the optimizer should reorder selections so that they are applied to the table in order of increasing rank, where $\text{rank}_i = (r_i - 1)/c_i$. Prove that this assertion is optimal. That is, show that no other ordering could result in a query of lower cost. (Hint: It may be easiest to consider the special case where $n = 2$ first and generalize from there.)

Exercise 23.9 ORDBMSs support references as a data type. It is often claimed that using references instead of key-foreign key relationships will give much higher performance for joins. This question asks you to explore this issue.

- Consider the following SQL:1999 DDL which only uses straight relational constructs:

```
CREATE TABLE R(rkey integer, rdata text);
CREATE TABLE S(skey integer, rfkey integer);
```

Assume that we have the following straightforward join query:

```
SELECT S.skey, H..relata
FROM   S, R
WHERE  S.rfkey = R.rkey
```

- Now consider the following SQL:1999 ORDBMS schema:

```
CREATE TYPE r_t AS ROW(rkey integer, rdata text);
CREATE TABLE R OF r_t REF is SYSTEM GENERATED;
CREATE TABLE S (skey integer, r REF(r_t) SCOPE R);
```

Assume we have the following query:

```
SELECT S.skey, S.r.rkey
FROM   S
```

What algorithm would you suggest to evaluate the pointer join in the ORDBMS schema? How do you think it will perform versus a relational join on the previous schema?

Exercise 23.1.0 Many object-relational systems support set-valued attributes using some variant of the setof constructor. For example, assuming we have a type `person_t`, we could have created the table `Films` in the Dinky Schema in Figure 23.1 as follows:

```
CREATE TABLE Films(filrno integer, title text, stars setof Person);
```

1. Describe two ways of implementing set-valued attributes. One way requires variable-length records, even if the set elements are all fixed-length.
2. Discuss the impact of the two strategies on optimizing queries with set-valued attributes.
3. Suppose you would like to create an index on the column `stars` in order to look up films by the name of the star that has starred in the film. For both implementation strategies, discuss alternative index structures that could help speed up this query.
4. What types of statistics should the query optimizer maintain for set-valued attributes? How do we obtain these statistics?

BIBLIOGRAPHIC NOTES

A number of the object-oriented features described here are based in part on fairly old ideas in the programming languages community. [42] provides a good overview of these ideas in a database context. Stonebraker's book [719] describes the vision of ORDBMSs embodied by his company's early product, Illustra (now a product of Informix). Current commercial DBMSs with object-relational support include Informix Universal Server, IBM DB/2 CS V2, and UniSQL. An new version of Oracle is scheduled to include ORDBMS features as well.

Many of the ideas in current object-relational systems came out of a few prototypes built in the 1980s, especially POSTGRES [723], Starburst [351], and O2 [218].

The idea of an object-oriented database was first articulated in [197], which described the GernStone prototype system. Other prototypes include DASDBS [657], EXODUS [130], nus [273], ObjectStore [463], ODE, [18] ORION [432], SHORE [129], and THOR [482]. O2 is actually an early example of a system that was beginning to merge the themes of ORDBMSs

and OODBMSs—it could fit in this list as well. [41] lists a collection of features that are generally considered to belong in an OODBMS. Current commercially available OODBMSs include GellStone, Itasca, O2, Objectivity, ObjectStore, Ontos, Poet, and Versant. [431] compares OODBMSs and RDBMSs.

Database support for ADTs was first explored in the INGRES and POSTGRES projects at U.C. Berkeley. The basic ideas are described in [716], including mechanisms for query processing and optimization with ADTs as well as extensible indexing. Support for ADTs was also investigated in the Dannstadt database system, [480]. Using the POSTGRES index extensibility correctly required intimate knowledge of DBMS-internal transaction mechanisms. Generalized search trees were proposed to solve this problem; they are described in [376], with concurrency and ARIES-based recovery details presented in [447]. [672] proposes that users must be allowed to define operators over ADT objects and properties of these operators that can be utilized for query optimization, rather than just a collection of methods.

Array chunking is described in [653]. Techniques for method caching and optimizing queries with expensive methods are presented in [37:3, 165]. Client-side data caching in a client-server OODBMS is studied in [283]. Clustering of objects on disk is studied in [741]. Work on nested relations was an early precursor of recent research on complex objects in OODBMSs and ORDBMSs. One of the first nested relation proposals is [504]. MVDs play an important role in reasoning about redundancy in nested relations; see, for example, [579]. Storage structures for nested relations were studied in [215].

Formal models and query languages for object-oriented databases have been widely studied; papers include [4, 56, 75, 125, 391, 392, 428, 578, 724]. [427] proposes SQL extensions for querying object-oriented databases. An early and elegant extension of SQL with path expressions and inheritance was developed in GEM [791]. There has been much interest in combining deductive and object-oriented features. Papers in this area include [44, 288, 495, 556, 706, 793]. See [3] for a thorough textbook discussion of formal aspects of object-orientation and query languages.

[433, 435, 721, 796] include papers on DBMSs that would now be termed object-relational and/or object-oriented. [794] contains a detailed overview of schema and database evolution in object-oriented database systems. A thorough presentation of SQL:1999 can be found in [525], and advanced features, including the object extensions, are covered in [523]. A short survey of new SQL:1999 features can be found in [2:37]. The incorporation of several SQL:1999 features into IBM DB2 is described in [128]. OQL is described in [141]. It is based to a large extent on the O2 query language, which is described, together with other aspects of O2, in the collection of papers [55].



24

DEDUCTIVE DATABASES

- ☛ What is the motivation for extending SQL with recursive queries?
- ☛ What important properties must recursive programs satisfy to be practical?
- ☛ What are least models and least fixpoints and how do they provide a theoretical foundation for recursive queries?
- ☛ What complications are introduced by negation and aggregate operations? How are they addressed?
- ☛ What are the challenges in efficient evaluation of recursive queries?
- **Key concepts:** Datalog, deductive databases, recursion, rules, inferences, safety, range-restriction; least model, declarative semantics; least fixpoint, operational semantics, fixpoint operator; negation, stratified programs; aggregate operators, multiset generation, grouping; efficient evaluation, avoiding repeated inferences, Seminaive fixpoint evaluation; pushing query selections, Magic Sets rewriting

For 'Is' and 'Is-Not' though with Rule and Line,
And 'Up-and-Down' by Logic I define,
Of all that one should care to fathom, I
Was never deep in anything but-----\Vine.

.. *Rubaiyat of Omar !*(hayyarn, Translated by Edward Fitzgerald

Relational database management systems have been enormously successful for C, administrative data processing. In recent years, however, as people have tried to

use database systems in increasingly complex applications, some important limitations of these systems have been exposed. For some applications, the query language and constraint definition capabilities have been found inadequate. As an example, some companies maintain a huge parts inventory database and frequently want to ask questions such as, “Are we running low on any parts needed to build a ZX600 sports car?” or “What is the total component and assembly cost to build a ZX600 at today’s part prices?” These queries cannot be expressed in SQL-92.

We begin this chapter by discussing queries that cannot be expressed in relational algebra or SQL and present a more powerful relational language called *Datalog*. Queries and views in SQL can be understood as **if-then** rules: “If some tuples exist in tables mentioned in the FROM clause that satisfy the conditions listed in the WHERE clause, then the tuple described in the SELECT clause is included in the answer.” Datalog definitions retain this if-then reading, with the significant new feature that definitions can be *recursive*, that is, a table can be defined in terms of itself. The SQL:1999 standard, the successor to the SQL-92 standard, requires support for recursive queries, and a large subset of the systems, notably IBM’s DB2 DBMS, already support them.

Evaluating Datalog queries poses some additional challenges, beyond those encountered in evaluating relational algebra queries, and we discuss some important implementation and optimization techniques developed to address these challenges. Interestingly, some of these techniques have been found to improve performance of even nonrecursive SQL queries and have therefore been implemented in several current relational DBMS products.

In Section 24.1, we introduce recursive queries and Datalog notation through an example. We present the theoretical foundations for recursive queries, least fixpoints and least models, in Section 24.2. We discuss queries that involve the use of negation or set-difference in Section 24.3. Finally, we consider techniques for evaluating recursive queries efficiently in Section 24.5.

24.1 INTRODUCTION TO RECURSIVE QUERIES

We begin with a simple example that illustrates the limits of SQL-92 queries and the power of recursive definitions. Let *Assembly* be a relation with three fields *part*, *subpart*, and *qty*. An example instance of *Assembly* is shown in Figure 24.1. Each tuple in *Assembly* indicates how many copies of a particular subpart are contained in a given part. The first tuple indicates, for example, that a trike contains three wheels. The *Assembly* relation can be visualized as a tree, as shown in Figure 24.2. A tuple is shown as an edge going from the part to the subpart, with the *qty* value as the edge label.

<i>part</i>	<i>subpart</i>	<i>qty</i>
trike	\vheel	3
trike	fralne	1
frarne	seat	1
franle	pedal	1
wheel	spoke	2
\vheel	tire	1
tire	run	1
tire	tube	1

Figure 24.1 An Instance of Assembly

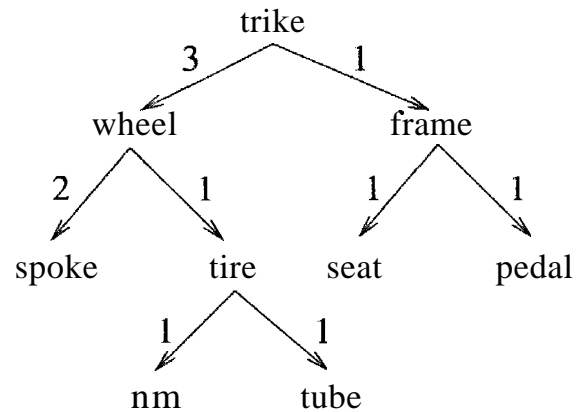


Figure 24.2 Assembly Instance Seen as a Tree

A natural question to ask is, "What are the components of a trike?" Rather surprisingly, this query is impossible to write in SQL-92. Of course, if we look at a given instance of the Assembly relation, we can write a 'query' that takes the union of the parts that are used in a trike. But such a query is not interesting---we want a query that identifies all components of a trike for *any* instance of Assembly, and such a query cannot be written in relational algebra or in SQL-92. Intuitively, the problem is that we are forced to join the Assembly relation with itself to recognize that *trike* contains *spoke* and *tire*, that is, to go one level down the Assembly tree. For each additional level, we need an additional join; two joins are needed to recognize that *trike* contains *rim*, which is a subpart of *tire*. Thus, the number of joins needed to identify all subparts of *trike* depends on the height of the Assembly tree, that is, on the given instance of the Assembly relation. No relational algebra query works for all instances; given any query, we can construct an instance whose height is greater than the number of joins in the query.

24.1.1 Datalog

We now define a relation called Components that identifies the components of every part. Consider the following program, or collection of rules:

```

Components(Part, Subpart) :- Assembly(Part, Subpart, Qty)
Components(Part, Subpart) :- Assembly(Part, Part2, Qty),
                                Components(Part2, Subpart)

```

These are rules in Datalog, a relational query language inspired by Prolog, the well-known logic programming language; indeed, the notation follows Prolog. The first rule should be read as follows:

For all values of Part, Subpart, and Qty,

if there is a tuple $\langle \text{Part}, \text{Subpart}, \text{Qty} \rangle$ in *Asscnbly*,
 then there MUST be a tuple $\langle \text{Part}, \text{Subpart} \rangle$ in *(;olnponents*.

The second rule should be read as follows:

For all values of *Part*, *Part2*, *Subpart*, and *Qty*,
 if there is a tuple $\langle \text{Part}, \text{Part2}, \text{Qty} \rangle$ in *Assernbly* and
 a tuple $\langle \text{Part2}, \text{Subpart} \rangle$ in *Components*,
 then there MUST be a tuple $\langle \text{Part}, \text{Subpart} \rangle$ in *COlnponents*.

The part to the right of the :- symbol is called the **body** of the rule, and the part to the left is called the **head** of the rule. The symbol :- denotes logical implication; if the tuples mentioned in the body exist in the database, it is implied that the tuple mentioned in the head of the rule must also be in the database. (Note that the body could be empty; in this case, the tuple mentioned in the head of the rule must be included in the database.) Therefore, if we are given a set of *Assenbly* and *Cornponents* tuples, each rule can be used to infer, or deduce, some new tuples that belong in *COlnponents*. This is why database systems that support Datalog rules are often called **deductive database systems**.

By assigning constants to the variables that appear in a rule, we can infer a specific *Coruponents* tuple. For example, by setting *Part*=*trike*, *Subpart*=*wheel*, and *Qty*=3, we can infer that $\langle \text{trike}, \text{wheel} \rangle$ is in *eoulponents*. Each rule is really a *ternplate* for making inferences: An inference is the use of a rule to generate a new tuple (for the relation in the head of the rule) by substituting constants for variables in such a way that every tuple in the rule body (after the substitution) is in the corresponding relation instance.

By considering each tuple in *Assenbly* in turn, the first rule allows us to infer that the set of tuples obtained by taking the projection of *Assenbly* onto its first two fields is in *CCHnponents*.

The second rule then allows us to combine previously discovered *Cornponents* tuples with *Assenbly* tuples to infer new *Cornponents* tuples. We can apply the second rule by considering the cross-product of *Assenbly* and (the current instance of) *Cornponents* and assigning values to the variables in the rule for each row of the cross-product, one row at a time. (Observe how the repeated use of the variable *Part2* prevents certain rows of the cross-product from contributing any new tuples; in effect, it specifies an equality join condition on *Assenbly* and *Cornponents*. The tuples obtained by one application of this rule are shown in Figure 24.3. (In addition, *COlnponents* contains the tuples obtained by applying the first rule; these are not shown.)

<i>part</i>	<i>subpart</i>
trike	spoke
trike	tire
trike	seat
trike	pedal
wheel	rhv
wheel	tube

Figure 24.3 Components Tuples Obtained by Applying the Second Rule Once

<i>part</i>	<i>subpart</i>
trike	<u>spoke</u>
trike	tire
<u>trike</u>	<u>seat</u>
trike	pedal
<u>wheel</u>	run
<u>wheel</u>	<u>tube</u>
<u>trike</u>	<u>rim</u>
trike	tube

Figure 24.4 Components Tuples Obtained by Applying the Second Rule Twice

The tuples obtained by a second application of this rule are shown in Figure 24.4. Note that each tuple shown in Figure 24.3 is reinferred. Only the last two tuples are new.

Applying the second rule a third time does not generate additional tuples. The set of Components tuples shown in Figure 24.4 includes all the tuples that can be inferred using the two Datalog rules defining Components and the given instance of Assembly. The components of a trike can now be obtained by selecting all Components tuples with the value *trike* in the first field.

Each application of a Datalog rule can be understood in terms of relational algebra. The first rule in our example program simply applies projection to the Assembly relation and adds the resulting tuples to the Components relation, which is initially empty. The second rule joins Assembly with Components and then does a projection. The result of each rule application is combined with the existing set of Components tuples using union.

The only Datalog operation that goes beyond relational algebra is the *repeated* application of the rules defining Components until no new tuples are generated. This repeated application of a set of rules is called the *fixpoint* operation, and we develop this idea further in the next section.

We conclude this section by rewriting the Datalog definition of Components using SQL:1999 syntax:

```
WITH RECURSIVE Components(Part, Subpart) AS
  (SELECT A1.Part, A1.Subpart FROM Assembly A1)
  UNION
  (SELECT A2.Part, C1.Subpart
   FROM Assembly A2, Components C1
```

```

WHERE  A2.Subpart = C1.Part)

SELECT * FROM Components C2

```

The WITH clause introduces a relation that is part of a query definition; this relation is similar to a view, but the scope of a relation introduced using WITH is local to the query definition. The RECURSIVE keyword signals that the table (in our example, Components) is recursively defined. The structure of the definition closely parallels the Datalog rules. Incidentally, if we wanted to find the components of a particular part, for example, *tTike*, we can simply replace the last line with the following:

```

SELECT * FROM Components C2
WHERE  C2.Part = 'trike'

```

24.2 THEORETICAL FOUNDATIONS

We classify the relations in a Datalog program as either output relations or input relations. Output relations are defined by rules (e.g., Components), and input relations have a set of tuples explicitly listed (e.g., Assembly). Given instances of the input relations, we must compute instances for the output relations. The meaning of a Datalog program is usually defined in two different ways, both of which essentially describe the relation instances for the output relations. Technically, a query is a selection over one of the output relations (e.g., all Components tuples *C* with *C.part* = *tTike*). However, the meaning of a query is clear once we understand how relation instances are associated with the output relations in a Datalog program.

The first approach to defining the semantics of a Datalog program, called the *least model semantics*, gives users a way to understand the program without thinking about how the program is to be executed. That is, the semantics is *declarative*, like the semantics of relational calculus, and not *operational* like relational algebra semantics. This is important because recursive rules make it difficult to understand a program in terms of an evaluation strategy.

The second approach, called the *least fixpoint semantics*, gives a conceptual evaluation strategy to compute the desired relation instances. This serves as the basis for recursive query evaluation in a DBMS. More efficient evaluation strategies are used in an actual implementation, but their correctness is shown by demonstrating their equivalence to the least fixpoint approach. The fixpoint semantics is thus operational and plays a role analogous to that of relational algebra semantics for nonrecursive queries.

24.2.1 Least Model Semantics

We want users to be able to understand a Datalog program by understanding each rule independent of other rules, with the meaning: *If the body is true, the head is also true.* This intuitive reading of a rule suggests that, given certain relation instances for the relation names that appear in the body of a rule, the relation instance for the relation mentioned in the head of the rule must contain a certain set of tuples. If a relation name R appears in the heads of several rules, the relation instance for R must satisfy the intuitive reading of all these rules. However, we do not want tuples to be included in the instance for R unless they are necessary to satisfy one of the rules defining R . That is, we want to compute only tuples for R that are supported by some rule for R .

To make these ideas precise, we need to introduce the concepts of models and least models. A model is a collection of relation instances, one instance for each relation in the program, that satisfies the following condition. For every rule in the program, whenever we replace each variable in the rule by a corresponding constant, the following holds:

If every tuple in the body (obtained by our replacement of variables with constants) is in the corresponding relation instance,

Then the tuple generated for the head (by the assignment of constants to variables that appear in the head) is also in the corresponding relation instance.

Observe that the instances for the input relations are given, and the definition of a model essentially restricts the instances for the output relations.

Consider the rule

$$\text{Components}(\text{Part}, \text{Subpart}) \text{ :- } \text{Assembly}(\text{Part}, \text{Part2}, \text{Qty}), \\ \text{Components}(\text{Part2}, \text{Subpart}).$$

Suppose we replace the variable *Part* by the constant *wheel*, *Part2* by *tire*, *Qty* by 1, and *Subpart* by *rim*:

$$\text{Components}(\text{wheel}, \text{rim}) \text{ :- } \text{Assembly}(\text{wheel}, \text{tire}, 1), \\ \text{Components}(\text{tire}, \text{rim}).$$

Let A be an instance of *Assembly* and C be an instance of *Components*. If A contains the tuple $\langle \text{wheel}, \text{tire}, 1 \rangle$ and C contains the tuple $\langle \text{tire}, \text{rim} \rangle$, then C must also contain the tuple $\langle \text{wheel}, \text{rim} \rangle$ for the pair of instances A and C .

to be a model. Of course, the instances A and C must satisfy the inclusion requirement just illustrated for *every* assignment of constants to the variables in the rule: If the tuples in the rule body are in A and C , the tuple in the head must be in C .

As an example, the instances of *Assembly* shown in Figure 24.1 and *Components* shown in Figure 24.4 together form a model for the *Components* program.

Given the instance of *Assembly* shown in Figure 24.1, there is no justification for including the tuple $\langle \text{spoke}, \text{pedal} \rangle$ to the *Components* instance. Indeed, if we add this tuple to the *components* instance in Figure 24.4, we no longer have a model for our program, as the following instance of the recursive rule demonstrates, since $\langle \text{wheel}, \text{pedal} \rangle$ is not in the *Components* instance:

Components(wheel, pedal) :- *Assembly*(wheel, spoke, 2),
 Components(spoke, pedal).

However, by also adding the tuple $\langle \text{wheel}, \text{pedal} \rangle$ to the *Components* instance, we obtain another model of the *Components* program. Intuitively, this is unsatisfactory since there is no justification for adding the tuple $\langle \text{spoke}, \text{pedal} \rangle$ in the first place, given the tuples in the *Assembly* instance and the rules in the program.

We address this problem by using the concept of a least model. A least model of a program is a model M such that for every other model M_2 of the same program, for each relation R in the program, the instance for R in M is contained in the instance of R in M_2 . The model formed by the instances of *Assembly* and *Components* shown in Figures 24.1 and 24.4 is the least model for the *Components* program with the given *Assembly* instance.

24.2.2 The Fixpoint Operator

A fixpoint of a function f is a value v such that the function applied to the value returns the same value, that is, $f(v) = v$. Consider a function applied to a set of values that also returns a set of values. For example, we can define *double* to be a function that multiplies every element of the input set by two and *double+* to be *double* \cup *identity*. Thus, *double*({1,2,5}) = {2,4,10}, and *double+*({1,2,5}) = {1,2,4,5,10}. The set of all even integers which happens to be an infinite set is a fixpoint of the function *double+*. Another fixpoint of the function *double+* is the set of all integers. The first fixpoint (the set of all even integers) is *smaller* than the second fixpoint (the set of all integers) because it is contained in the latter.

The least fixpoint of a function is the fixpoint that is smaller than every other fixpoint of that function. In general, it is not guaranteed that a function has a least fixpoint. For example, there may be two fixpoints, neither of which is smaller than the other. (Does *double* have a least fixpoint? What is it?)

Now let us turn to functions over sets of tuples, in particular, functions defined using relational algebra expressions. The Components relation can be defined by an equation of the form

$$\text{Components} = \pi_{1,5}(\text{Assembly} \bowtie_{2=1} \text{Components}) \cup \pi_{1,2}(\text{Assembly})$$

This equation has the form

$$\text{Components} = f(\text{Components}, \text{Assembly})$$

where the function f is defined using a relational algebra expression. For a given instance of the input relation Assembly, this can be simplified to

$$\text{Components} = f(\text{Components})$$

The least fixpoint of f is an instance of Components that satisfies this equation. Clearly the projection of the first two fields of the tuples in the given instance of the input relation Assembly must be included in the (instance that is the) least fixpoint of Components. In addition, any tuple obtained by joining Components with Assembly and projecting the appropriate fields must also be in Components.

A little thought shows that the instance of Components that is the least fixpoint of f can be computed using repeated applications of the Datalog rules shown in the previous section. Indeed, applying the two Datalog rules is identical to evaluating the relational expression used in defining Components. If an application generates Components tuples that are not in the current instance of the Components relation, the current instance cannot be the fixpoint. Therefore, we add the new tuples to Components and evaluate the relational expression (equivalently, the two Datalog rules) again. This process is repeated until every tuple generated is already in the current instance of Components. When applying the rules to the current set of tuples does not produce any new tuples, we have reached a fixpoint. If Components is initialized to the empty set of tuples, intuitively we infer only tuples that are necessary by the definition of a fixpoint, and the fixpoint computed is the least fixpoint.

24.2.3 Safe Datalog Programs

Consider the following program:

ComplexYarts(Part) :- Assembly(Part, Subpart, Qty), Qty > 2.

According to this rule, a complex part is defined to be any part that has more than two copies of any one subpart. For each part mentioned in the Assembly relation, we can easily check whether it is a complex part. In contrast, consider the following program:

```
PriceYarts(Part, Price) :-
    Assembly(Part, Subpart, Qty), Qty > 2.
```

This variation seeks to associate a price with each complex part. However, the variable *Price* does not appear in the body of the rule. This means that an infinite number of tuples must be included in any model of this program. To see this, suppose we replace the variable *Part* by the constant *trike*, *SubPart* by *wheel*, and *Qty* by 3. This gives us a version of the rule with the only remaining variable being *Price*:

```
PriceYarts(trike, Price) :- Assembly(trike, wheel, 3), 3 > 2.
```

Now, any assignment of a constant to *Price* gives us a tuple to be included in the output relation *Price_Parts*. For example, replacing *Price* by 100 gives us the tuple *Price_Parts(trike, 100)*. If the least model of a program is not finite, for even one instance of its input relations, then we say the program is unsafe.

Database systems disallow unsafe programs by requiring that every variable in the head of a rule also appear in the body. Such programs are said to be range-restricted, and every range-restricted Datalog program has a finite least model if the input relation instances are finite. In the rest of this chapter, we assume that programs are range-restricted.

24.2.4 Least Model = Least Fixpoint

Does a Datalog program always have a least model? (Or is it possible that there are two models, neither of which is contained in the other?) Similarly, does every Datalog program have a least fixpoint? What is the relationship between the least model and the least fixpoint of a Datalog program?

As we noted earlier, not every function has a least fixpoint. Fortunately, every function defined in terms of relational algebra expressions that do not contain set-difference is guaranteed to have a least fixpoint, and the least fixpoint can be computed by repeatedly evaluating the function. This tells us that every Datalog program has a least fixpoint and that it can be computed by repeatedly applying the rules of the program on the given instances of the input relations.

Further, every Datalog program is guaranteed to have a least model and the least model is equal to the least fixpoint of the program. These results (whose

proofs we do not discuss) provide the basis for Datalog query processing. 'Users can understand a program in terms of 'If the body is true, the head is also true,' thanks to the least model semantics. The DBMS can compute the answer by repeatedly applying the program rules, thanks to the least fixpoint semantics and the fact that the least model and the least fixpoint are identical.

24.3 RECURSIVE QUERIES WITH NEGATION

Unfortunately, once set-difference is allowed in the body of a rule, there may be no least model or least fixpoint for a program. Consider the following rules:

Big(Part):- **Assembly(Part, Subpart, Qty), Qty > 2,**
 NOT Small(Part) .
Small(Part) :- **Assembly(Part, Subpart, Qty), NOT Big(Part).**

These two rules can be thought of as an attempt to divide parts (those that are mentioned in the first column of the Assembly table) into two classes, Big and Small. The first rule defines Big to be the set of parts that use at least three copies of some subpart and are not classified as small parts. The second rule defines Small as the set of parts not classified as big parts.

If we apply these rules to the instance of Assembly shown in Figure 24.1, *trike* is the only part that uses at least three copies of some subpart. Should the tuple (*trike*) be in Big or Small? If we apply the first rule and then the second rule, this tuple is in Big. To apply the first rule, we consider the tuples in Assembly, choose those with Qty > 2 (which is just (*trike*)), discard those in the current instance of Small (both Big and Small are initially empty), and add the tuples that are left to Big. Therefore, an application of the first rule adds (*trike*) to Big. Proceeding similarly, we can see that if the second rule is applied before the first, (*trike*) is added to Small instead of Big.

This program has two fixpoints, neither of which is smaller than the other, as shown in Figure 24.5. (The first fixpoint has a Big tuple that does not appear in the second fixpoint; therefore, it is not smaller than the second fixpoint. The second fixpoint has a Small tuple that does not appear in the first fixpoint; therefore, it is not smaller than the first fixpoint. The order in which we apply the rules determines which fixpoint is computed; this situation is very unsatisfactory. We want users to be able to understand their queries without thinking about exactly how the evaluation proceeds.

The root of the problem is the use of NOT. When we apply the first rule, some inferences are disallowed because of the presence of tuples in Small. Parts

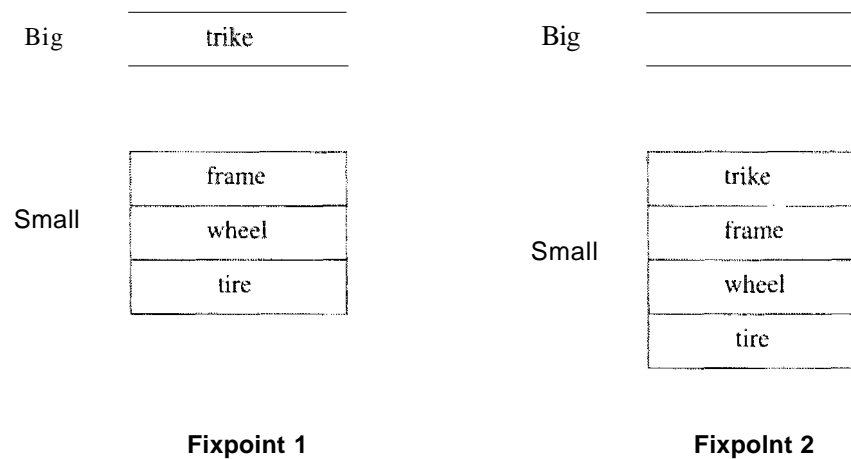


Figure 24.5 Two Fixpoints for the Big/Small Program

that satisfy the other conditions in the body of the rule are candidates for addition to Big; we remove the parts in Small from this set of candidates. Thus, some inferences that are possible if Small is empty (as it is before the second rule is applied) are disallowed if Small contains tuples (generated by applying the second rule before the first rule). Here is the difficulty: If NOT is used, the addition of tuples to a relation can *disallow* the inference of other tuples. Without NOT, this situation can never arise; the addition of tuples to a relation can *never* disallow the inference of other tuples.

Range-Restriction and Negation

If rules are allowed to contain NOT in the body, the definition of range-restriction must be extended to ensure that all range-restricted programs are safe. If a relation appears in the body of a rule preceded by NOT, we call this a *negated occurrence*. Relation occurrences in the body that are not negated are called *positive occurrences*. A program is range-restricted if every variable in the head of the rule appears in some positive relation occurrence in the body.

24.3.1 Stratification

A widely used solution to the problem caused by negation, or the use of NOT, is to impose certain syntactic restrictions on programs. These restrictions can be easily checked and programs that satisfy them have a natural meaning.

We say that a table T depends on a table S if some rule with T in the head contains S , or (recursively) contains a predicate that depends on S , in the body. A recursively defined predicate always depends on itself. For example, Big depends on Small (and on itself). Indeed, the tables Big and Small (re

mutually recursive, that is, the definition of *Big* depends on *Small* and vice versa. We say that a table *T* depends negatively on a table *S* if some rule with *T* in the head contains NOT *S*, or (recursively) contains a predicate that depends negatively on *S*, in the body.

Suppose we classify the tables in a program into strata or layers as follows. The tables that do not depend on any other tables are in stratum 0. In our *Big/Small* example, *Assembly* is the only table in stratum 0. Next, we identify tables in stratum 1; these are tables that depend only on tables in stratum 0 or stratum 1 and depend negatively only on tables in stratum 0. Higher strata are similarly defined: The tables in stratum *i* are those that do not belong to lower strata, depend only on tables in stratum *i* or lower strata, and depend negatively only on tables in lower strata. A stratified program is one whose tables can be classified into strata according to the above algorithm.

The *Big/Small* program is not stratified. Since *Big* and *Small* depend on each other, they must be in the same stratum. However, they depend negatively on each other, violating the requirement that a table can depend negatively only on tables in lower strata. Consider the following variant of the *Big/Small* program, in which the first rule has been modified:

```
Big2(Part) :- Assembly(Part, Subpart, Qty), Qty > 2.
Small2(Part) :- Assembly(Part, Subpart, Qty), NOT Big2(Part).
```

This program is stratified. *Small2* depends on *Big2* but *Big2* does not depend on *Small2*. *Assembly* is in stratum 0, *Big2* is in stratum 1, and *Small2* is in stratum 2.

A stratified program is evaluated stratum-by-stratum, starting with stratum 0. To evaluate a stratum, we compute the fixpoint of all rules defining tables in this stratum. When evaluating a stratum, any occurrence of NOT involves a table from a lower stratum, which has therefore been completely evaluated by now. The tuples in the negated table still disallow some inferences, but the effect is completely deterministic, given the stratum-by-stratum evaluation. In the example, *Big2* is computed before *Small2* because it is in a lower stratum than *Small2*: *<trike>* is added to *Big2*. Next, when we compute *Small2*, we recognize that *<trike>* is not in *Small2* because it is already in *Big2*.

Incidentally, note that the stratified *Big/Small* program is not even recursive. If we replace *Assembly* by *Components*, we obtain a recursive, stratified program: *Assembly* is in stratum 0, *Components* is in stratum 1, *Big2* is also in stratum 1, and *Small2* is in stratum 2.

Intuition behind Stratification

Consider the stratified version of the Big/Small program. The rule defining Big2 forces us to add $\langle trike \rangle$ to Big2 and it is natural to assume that $\langle trike \rangle$ is the only tuple in Big2, because we have no supporting evidence for any other tuple being in Big2. The minimal fixpoint computed by stratified fixpoint evaluation is consistent with this intuition. However, there is another minimal fixpoint: We can place every part in Big2 and make Small2 be empty. While this assignment of tuples to relations seems unintuitive, it is nonetheless a minimal fixpoint.

The requirement that programs be stratified gives us a natural order for evaluating rules. When the rules are evaluated in this order, the result is a unique fixpoint that is one of the minimal fixpoints of the program. The fixpoint computed by the stratified fixpoint evaluation usually corresponds well to our intuitive reading of a stratified program, even if the program has more than one minimal fixpoint.

For nonstratified Datalog programs, it is harder to identify a natural model from among the alternative minimal models, especially when we consider that the meaning of a program must be clear even to users who lack expertise in Datalog logic. Although considerable research has been done on identifying natural models for nonstratified programs, practical implementations of Datalog have concentrated on stratified programs.

Relational Algebra and Stratified Datalog

Every relational algebra query can be written as a range-restricted, stratified Datalog program. (Of course, not all Datalog programs can be expressed in relational algebra; for example, the Components program.) We sketch the translation from algebra to stratified Datalog by writing a Datalog program for each of the basic algebra operations, in terms of two example tables R and S, each with two fields:

Selection:	Result(Y) :- R(X,Y), X=c.
Projection:	Result(Y) :- R(X,Y).
Cross-product:	Result(X,Y,U,V) :- R(X,Y), S(U,V).
Set-difference:	Result(X,Y) :- R(X,Y), NOT S(U,V).
Join:	Result(X,Y) :- R(X,Y).
	Result(X,Y) :- S(X,Y).

We conclude our discussion of stratification by noting that SQL:1999 requires programs to be stratified. The stratified Big/Small program is shown below in SQL:1999 notation, with a final additional selection on Big2:

SQL:1999 and Datalog Queries: A Datalog rule is **linear** recursive if the body contains at most one occurrence of any table that depends on the table in the head of the rule. A linear recursive program contains only linear recursive rules. All linear recursive Datalog programs can be expressed using the recursive features of SQL:1999. However, these features are not in Core SQL.

```

WITH
  Big2(Part) AS
    (SELECT  A1.Part FROM  Assembly A1 WHERE  Qty > 2)
  Small2(Part) AS
    ((SELECT- A2.Part FROM  Assembly A2)
     EXCEPT
     (SELECT  B1.Part from1 Big2 B1))

SELECT  * FROM Big2 B2

```

24.4 FROM DATALOG TO SQL

To support recursive queries in SQL, we must take into account the features of SQL that are not found in Datalog. Two central SQL features missing in Datalog are (1) SQL treats tables as *multisets* of tuples, rather than sets, and (2) SQL permits grouping and aggregate operations.

The multiset semantics of SQL queries can be preserved if we do not check for duplicates after applying rules. Every relation instance, including instances of the recursively defined tables, is a multiset. The number of occurrences of a tuple in a relation is equal to the number of distinct inferences that generate this tuple.

The second point can be addressed by extending Datalog with grouping and aggregation operations. This must be done with multiset semantics in mind, as we now illustrate. Consider the following program:

```

NumPartsCPart, SUM((Qty))) :- AssemblyCPart, Subpart, Qty).

```

This program is equivalent to the SQL query

```

SELECT  A.Part, SUM (A.Qty)
FROM    Assembly A
GROUP BY A.Part

```


The angular brackets (...) notation was introduced in the LDL deductive system, one of the pioneering deductive database prototypes developed at IVCC in the late 1980s. We use it to denote *multiset generation*, or the creation of multiset-values. In principle, the rule defining NumParts is evaluated by first creating the temporary relation shown in Figure 24.6. We create the temporary relation by sorting on the *part* attribute (which appears on the left side of the rule, along with the (...) term) and collecting the multiset of *qty* values for each *part* value. We then apply the SUM aggregate to each multiset-value in the second column to obtain the answer, which is shown in Figure 24.7.

<i>part</i>	$\langle qty \rangle$
trike	{3,1}
frame	{1,1}
wheel	{2,1}
tire	{1,1}

Figure 24.6 Temporary Relation

<i>part</i>	SUM($\langle qty \rangle$)
trike	4
frame	2
wheel	3
tire	2

Figure 24.7 The Tuples in NumParts

The temporary relation shown in Figure 24.6 need not be materialized to compute NumParts; for example, SUM can be applied on-the-fly or Assembly can simply be sorted and aggregated as described in Section 14.6.

The use of grouping and aggregation, like negation, causes complications when applied to a partially computed relation. The difficulty is overcome by adopting the same solution used for negation, stratification. Consider the following program:¹

```
TotParts(Part, Subpart, SUM«(Qty)»):- BOM(Part, Subpart, Qty).
BOM(Part, Subpart, Qty) :- Assembly(Part, Subpart, Qty).
BOM(Part, Subpart, Qty) :- Assembly(Part, Part2, Qty2),
    BOM(Part2, Subpart, Qty3), Qty=Qty2*Qty3.
```

The idea is to count the number of copies of Subpart for each Part. By aggregating over BOM rather than Assembly, we count subparts at any level in the hierarchy instead of just immediate subparts. This program is a version of a well-known problem called *Bill-of-Materials* and variants of it are probably the most widely used recursive queries in practice.

The important point to note in this example is that we must wait until the relation BOM has been completely evaluated before we apply the TotParts rule. Otherwise, we obtain incomplete counts. This situation is analogous to the problem we faced with negation; we have to evaluate the negated relation

¹The reader should write this in SQL:1999 syntax, as a simple exercise.

SQL:1999 Cycle Detection: Safe Datalog queries that do not use arithmetic operations have finite answers and the fixpoint evaluation is guaranteed to halt. Unfortunately, recursive SQL queries may have infinite answer sets and query evaluation may not halt. There are two independent reasons for this: (1) the use of arithmetic operations to generate data values that are not stored in input tables of a query, and (2) multiset semantics for rule applications; intuitively, problems arise from *cycles* in the data. (To see this, consider the Components program on the Assembly instance shown in Figure 24.1 plus the tuple $\langle \text{tube}, \text{wheel}, 1 \rangle$.) SQL:1999 provides special constructs to check for such cycles.

completely before applying a rule that involves the use of NOT. If a program is stratified with respect to uses of $\langle \dots \rangle$ as well as NOT, stratified fixpoint evaluation gives us meaningful results.

There are two further aspects to this example. First, we must understand the cardinality of each tuple in BOM, based on the multiset semantics for rule application. Second, we must understand the cardinality of the multiset of Qty values for each $\langle Part, Subpart \rangle$ group in TotParts.

part	subpart	qty
trike	frame	1
trike	seat	1
frame	seat	1
frame	pedal	2
seat	cover	1

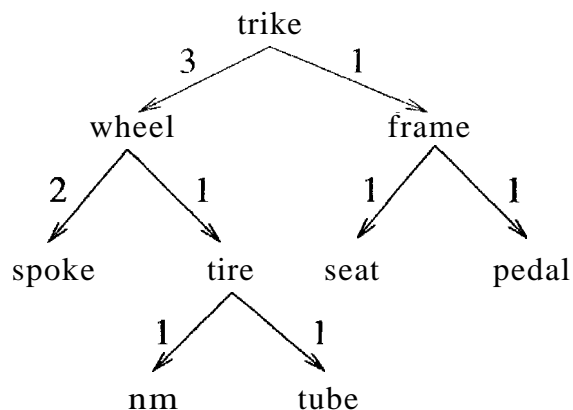


Figure 24.8 Another Instance of Assembly Figure 24.9 Assembly Instance Seen as a Graph

We illustrate these two points using the instance of Assembly shown in Figures 24.8 and 24.9. Applying the first BOM rule, we add (one copy of) every tuple in Assembly to BOM. Applying the second BOM rule, we add the following four tuples to BOM: $\langle \text{trike}, \text{seat}, 1 \rangle$, $\langle \text{trike}, \text{pedal}, 2 \rangle$, $\langle \text{trike}, \text{cover}, 1 \rangle$, and $\langle \text{frame}, \text{cover}, 1 \rangle$. (Observe that the tuple $\langle \text{trike}, \text{seat}, 1 \rangle$ was already in BOM because it was generated by applying the first rule; therefore, multiset semantics for rule application gives us two copies of this tuple. Applying the second BOM rule on the new tuples, we generate the tuple $\langle \text{trike}, \text{cover}, 1 \rangle$ (using the tuple $\langle \text{frame}, \text{cover}, 1 \rangle$ for BOM in the body of the rule): this is our second copy of the tuple. Applying the second rule again on this tuple does not generate any

tuples, and the cOInputation of the BOM relation is now cOInplete. The BaM instance at this stage is sho\vn in Figure 24.10.

<i>part</i>	<i>subpart</i>	<i>qty</i>
trike	frame	1
trike	seat	1
frame	seat	1
frame	pedal	2
seat	cover	1
trike	seat	1
trike	pedal	2
trike	cover	1
frame	cover	1
trike	cover	1

Figure 24.10 Instance of *BOM* Table

<i>part</i>	<i>subpart</i>	<i>qty</i>
trike	frame	{1}
trike	seat	{1,1}
trike	cover	{1,1}
trike	pedal	{2}
frame	seat	{1}
frame	pedal	{2}
seat	cover	{1}
frame	cover	{1}

Figure 24.11 Temporary Relation

Multiset grouping on this instance yields the temporary relation instance shown in Figure 24.11. (This step is only conceptual; the aggregation can be done on the fly without materializing this temporary relation.) Applying SUM to the multisets in the third column of this temporary relation gives us the instance for TotParts.

24.5 EVALUATING RECURSIVE QUERIES

The evaluation of recursive queries has been widely studied. While all the problems of evaluating nonrecursive queries continue to be present, the newly introduced fixpoint operation creates additional difficulties. A straightforward approach to evaluating recursive queries is to compute the fixpoint by repeatedly applying the rules as illustrated in Section 24.1.1. One application of all the program rules is called an iteration; we perform as many iterations as necessary to reach the least fixpoint. This approach has two main disadvantages:

- **Repeated Inferences:** As Figures 24.3 and 24.4 illustrate, inferences are repeated across iterations. That is, the same tuple is inferred repeatedly in *the same way*, using the same rule and the same tuples for tables in the body of the rule.
- **Unnecessary Inferences:** Suppose we want to find the components of only a *wheel*. Computing the entire Components table is wasteful and does not take advantage of information in the query.

In this section, we discuss how each of these difficulties can be overcome. We consider only Datalog programs without negation.

24.5.1 Fixpoint Evaluation without Repeated Inferences

Computing the fixpoint by repeatedly applying all rules is called Naive fixpoint evaluation. Naive evaluation is guaranteed to compute the least fixpoint, but every application of a rule repeats all inferences made by earlier applications of this rule. We illustrate this point using the following rule:

Components (Part, Subpart) :- Assembly (Part, Part2, Qty),
Components (Part2, Subpart).

When this rule is applied for the first time, after applying the first rule defining Components, the Components table contains the projection of Assembly on the first two fields. Using these Components tuples in the body of the rule, we generate the tuples shown in Figure 24.3. For example, the tuple (*wheel*, *rim*) is generated through the following inference:

```
Components(wheel, rim) :- Assembly(wheel, tire, 1),
                             Components(tire, rim).
```

When this rule is applied a second time, the Components table contains the tuples shown in Figure 24.3 in addition to the tuples that it contained before the first application. Using the Components tuples shown in Figure 24.3 leads to new inferences; for example,

```
Components(trike, rim) :- Assembly(trike, wheel, 3),
                             Components(wheel, rim).
```

However, every inference carried out in the first application of this rule is also repeated in the second application of the rule, since all the Assembly and Components tuples used in the first rule application are considered again. For example, the inference of $\langle wheel, rim \rangle$ shown above is repeated in the second application of this rule.

The solution to this repetition of inferences consists of remembering which inferences were carried out in earlier rule applications and not carrying them out again. We can ‘remember’ previously executed inferences efficiently by simply keeping track of which Components tuples were generated for the first time in the most recent application of the recursive rule. Suppose we keep track by introducing a new relation called *delta_Components* and storing just the newly generated Components tuples in it. Now, we can use only the tuples

in *delta_Components* in the next application of the recursive rule; any inference using other *Components* tuples should have been carried out in earlier rule applications.

This refinement of fixpoint evaluation is called *Seminaive fixpoint evaluation*. Let us trace *Seminaive fixpoint evaluation* on our example program. The first application of the recursive rule produces the *Components* tuples shown in Figure 24.3, just like *Naive fixpoint evaluation*, and these tuples are placed in *delta_Components*. In the second application, however, only *delta_Components* tuples are considered, which means that only the following inferences are carried out in the second application of the recursive rule:

```
Components(trike, rim) :- Assembly(trike, wheel, 3),
                           delta_Components(wheel, rim).
Components(trike, tube) :- Assembly(trike, wheel, 3),
                           delta_Components(wheel, tube).
```

Next, the bookkeeping relation *delta_Components* is updated to contain just these two *Components* tuples. In the third application of the recursive rule, only these two *delta_Components* tuples are considered and therefore no additional inferences can be made. The fixpoint of *Components* has been reached.

To implement *Seminaive fixpoint evaluation* for general Datalog programs, we apply all the recursive rules in a program together in an iteration. Iterative application of all recursive rules is repeated until no new tuples are generated in some iteration. To summarize how *Seminaive fixpoint evaluation* is carried out, there are two important differences with respect to *Naive fixpoint evaluation*:

- We maintain a *delta* version of every recursive predicate to keep track of the tuples generated for this predicate in the most recent iteration; for example, *delta_Components* for *Components*. The *delta* versions are updated at the end of each iteration.
- The original program rules are rewritten to ensure that every inference uses at least one *delta* tuple; that is, one tuple that was not known before the previous iteration. This property guarantees that the inference could not have been carried out in earlier iterations.

We do not discuss details of *Seminaive fixpoint evaluation* (such as the algorithm for rewriting program rules to ensure the use of a *delta* tuple in each inference).

24.5.2 Pushing Selections to Avoid Irrelevant Inferences

Consider a nonrecursive view definition. If we want only those tuples in the view that satisfy an additional selection condition, the selection can be added to the plan as a final operation, and the relational algebra transformations for commuting selections with other relational operators allow us to ‘push’ the selection ahead of more expensive operations such as cross-products (and joins). In effect, we restrict the computation by utilizing selections in the query specification. The problem is more complicated for recursively defined queries.

We use the following program as an example in this section:

```

SameLevel(81, 82)      Assembly(P1, 81, Q1),
                        Assembly(P1, 82, Q2),
SameLevel(81, 82)      Assembly(P1, 81, Qi),
                        SameLevel(P1, P2), Assembly(P2, 82, Q2).

```

Consider the tree representation of Assembly tuples illustrated in Figure 24.2. There is a tuple $\langle S1, S2 \rangle$ in SameLevel if there is a path from 81 to 82 that goes up a certain number of edges in the tree and then comes down the same number of edges.

Suppose we want to find all SameLevel tuples with the first field equal to *spoke*. Since SameLevel tuples can be used to compute other SameLevel tuples, we cannot just compute those tuples with *spoke* in the first field. For example, the tuple $\langle wheel, frame \rangle$ in SameLevel allows us to infer a SameLevel tuple with *spoke* in the first field:

```

SameLevel(spoke, seat) :- Assembly(wheel, spoke, 2),
                           SameLevel(wheel, frame),
                           Assembly(frame, seat, i),

```

Intuitively, we have to compute all SameLevel tuples whose first field contains a value on the path from *spoke* to the root in Figure 24.2. Each such tuple has the potential to contribute to answers for the given query. On the other hand, computing the entire SameLevel table is wasteful; for example, the SameLevel tuple $\langle tire, seat \rangle$ cannot be used to infer any answer to the given query (or, indeed, to infer any tuple that can in turn be used to infer an answer tuple). We define a new table, which we call *Magic_SameLevel*, such that each tuple in this table identifies a value *m* for which we have to compute all SameLevel tuples with *m* in the first column to answer the given query:

```

Magic_SameLevel(Pi) :- Magic_SameLevel(81), Assembly(P1, 81, Q1).
Magic_SameLevel(spoke) :-

```

Consider the tuples in `IVlagic_SameLevel`. Obviously we have $\langle spoke \rangle$. Using this `IVlagic_SameLevel` tuple and the Assembly tuple $\langle wheel, spoke, 2 \rangle$, we can infer that the tuple $\langle wheel \rangle$ is in `Magic_SameLevel`. Using this tuple and the Assembly tuple $\langle trike, wheel, 3 \rangle$, we can infer that the tuple $\langle trike \rangle$ is in `IVlagic_SameLevel`. Thus, `Magic_SameLevel` contains each node that is on the path from *spoke* to the root in Figure 24.2. The `Magic_SameLevel` table can be used as a filter to restrict the computation:

```
SameLevel(S1, S2) :- Magic_SameLevel(S1),
                    Assembly(P1, S1, Q1), Assembly(P2, S2, Q2).
SameLevel(S1, S2) :- Magic_SameLevel(S1), Assembly(P1, S1, Q1),
                    SameLevel(P1, P2), Assembly(P2, S2, Q2).
```

These rules together with the rules defining `rvlagic_SameLevel` give us a program for computing all `SameLevel` tuples with *spoke* in the first column. Notice that the new program depends on the query constant *spoke* only in the second rule defining `IVlagic_SameLevel`. Therefore, the program for computing all `SameLevel` tuples with *seat* in the first column, for instance, is identical except that the second `Magic_SameLevel` rule is

```
Magic_SameLevel(seat) :- .
```

The number of inferences made using the Magic program can be far fewer than the number of inferences made using the original program, depending on just how much the selection in the query restricts the computation.

24.5.3 The Magic Sets Algorithm

We illustrated the intuition behind the Magic Sets algorithm on the `SameLevel` program, which contains just one output relation and one recursive rule.

The intuition behind the rewriting is that the rows in the Magic tables correspond to the subqueries whose answers are relevant to the original query. By evaluating the rewritten program instead of the original program, we can restrict computation by intuitively pushing the selection condition in the query into the recursion.

The algorithm, however, can be applied to any Datalog program. The input to the algorithm consists of the program and a query pattern, which is a relation we want to query plus the fields for which a query will provide constants. The output of the algorithm is a rewritten program.

The Magic Sets program rewriting algorithm can be summarized as follows:

1. **Generate the Adorned Program:** In this step, the program is rewritten to make the pattern of queries and subqueries explicit.
2. **Add Magic Filters:** Modify each rule in the Adorned Program by adding a Magic condition to the body that acts as a filter on the set of tuples generated by this rule.
3. **Define the Magic Tables:** We create new rules to define the Magic tables. Intuitively, from each occurrence of a table R in the body of an Adorned Program rule, we obtain a rule defining the table Magic_R .

When a query is posed, we add the corresponding Magic tuple to the rewritten program and evaluate the least fixpoint of the program (using Sernaive evaluation).

We remark that the Magic Sets algorithm has turned out to be quite effective for computing correlated nested SQL queries, even if there is no recursion, and is used for this purpose in many commercial DBMSs, even systems that do not currently support recursive queries.

We now describe the three steps in the Magic Sets algorithm using the *SameLevel* program as a running example.

Adorned Program

We consider the query pattern SameLevel^{bf} . Thus, given a value c , we want to compute all rows in *SameLevel* in which c appears in the first column. We generate the Adorned Program P^{ad} from the given program P by repeatedly generating adorned versions of rules in P for every reachable query pattern, with the given query pattern as the only reachable pattern to begin with; additional reachable patterns are identified during the course of generating the Adorned Program as described next.

Consider a rule r whose head contains the same table as some reachable pattern. The adorned version of the rule depends on the order in which we consider the predicates in the body of the rule. To simplify our discussion, we assume that this is always left-to-right. First, we replace the head of the rule with the matching query pattern. After this step, the recursive *SameLevel* rule looks like this:

$$\text{SameLevel}^{bf}(S1, S2) \text{ :- } \text{Assembly}(P1, S1, Q1), \\ \text{SameLevel}(P1, P2), \text{Assembly}(P2, S2, Q2).$$

Next, we proceed left-to-right in the body of the rule until we encounter the first recursive predicate. All clauses that contain a constant or a variable that

appears to the left are marked *b* (for *bound*) and the rest are marked *f* (for *free*) in the query pattern for this occurrence of the predicate. We add this pattern to the set of reachable patterns and modify the rule accordingly:

$$\begin{aligned} \text{SameLevel}^{bf}(S1, S2) &:- \text{Assembly}(P1, S1, Q1), \\ &\quad \text{SameLevel}^{bf}(P1, P2), \text{Assembly}(P2, S2, Q2). \end{aligned}$$

If there are additional occurrences of recursive predicates in the body of the recursive rule, we continue (adding the query patterns to the reachable set and modifying the rule). (Of course, in linear recursive programs, there is at most one occurrence of a recursive predicate in a rule body.)

We repeat this until we have generated the adorned version of every rule in *P* for every reachable query pattern that contains the same table as the head of the rule. The result is the Adorned Program *pad*, which, in our example, is

$$\begin{aligned} \text{SameLevel}^{bf}(S1, S2) &:- \text{Assembly}(P1, S1, Q1), \\ &\quad \text{Assembly}(P1, S2, Q2). \\ \text{SameLevel}^{bf}(S1, S2) &:- \text{Assembly}(P1, S1, Q1), \\ &\quad \text{SameLevel}^{bf}(P1, P2), \text{Assembly}(P2, S2, Q2). \end{aligned}$$

In our example, there is only one reachable query pattern. In general, there can be several.²

Adding Magic Filters

Every rule in the Adorned Program is modified by adding a 'magic filter' predicate to obtain the rewritten program:

$$\begin{aligned} \text{SameLevel}^{bf}(S1, S2) &:- \text{Magic_SameLevel}^{bf}(S1), \\ &\quad \text{Assembly}(P1, S1, Q1), \text{Assembly}(P2, S2, Q2). \\ \text{SameLevel}^{bf}(S1, S2) &:- \text{Magic_SameLevel}^{bf}(S1), \\ &\quad \text{Assembly}(P1, S1, Q1), \text{SameLevel}^{bf}(P1, P2), \\ &\quad \text{Assembly}(P2, S2, Q2). \end{aligned}$$

The filter predicate is a copy of the head of the rule, 'with 'Magic' as a prefix for the table name and the variables in columns corresponding to *free* deleted, as illustrated in these two rules.

²As an example, consider a variant of the SameLevel program in which the variables *P1* and *P2* are interchanged in the body of the recursive rule (Exercise 24.5)

Defining Magic Filter Tables

Consider the Adorned Program after every rule has been modified as described. For each occurrence O of a recursive predicate in the body of a rule in this modified program, we generate a rule that defines a Magic predicate. The algorithm for generating this rule is as follows: (1) Delete everything to the right of occurrence $()$ in the body of the modified rule. (2) Add the prefix 'Magic' and delete the free columns of $()$. (3) Move O , with these changes, into the head of the rule.

From the recursive rule in our example, after steps (1) and (2) we get:

$$\begin{aligned} \text{SameLevel}^{bf}(S1, 82) &:- \text{Magic_SameLevel}^{bf}(S1), \\ &\quad \text{Assembly}(P1, S1, Q1), \text{Magic_SameLevel}^{bf}(P1). \end{aligned}$$

After step (3), we get:

$$\begin{aligned} \text{Magic_SameLevel}^{bf}(P1) &:- \text{Magic_SameLevel}^{bf}(S1), \\ &\quad \text{Assembly}(P1, S1, Q1). \end{aligned}$$

The query itself generates a row in the corresponding Magic table, for example, $\text{Magic_SameLevel}^{bf}(\text{seat})$.

24.6 REVIEW QUESTIONS

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

- Describe *Datalog* programs. Use an example Datalog program to explain why it is not possible to write recursive rules in SQL-92. (Section 24.1)
- Define the terms *rmodel* and *least model*. What can you say about least rmodels for Datalog programs? Why is this approach to defining the meaning of a Datalog program called *declarative*? (Section 24.2.1)
- Define the terms *fixpoint* and *least fixpoint*. What can you say about least fixpoints for Datalog programs? Why is this approach to defining the meaning of a Datalog program said to be *operational*? (Section 24.2.2)
- What is a *safe* program? Why is this property important? What is *range-restriction* and how does it ensure safety? (Section 24.2.3)
- What is the connection between least rmodels and least fixpoints for Datalog programs? (Section 24.2.4)

- Explain why programs with negation may not have a least model or least fixpoint. Extend the definition of *Tarski-Constraint* to programs with negation. (Section 24.3)
- What is a *stratified* program? How does stratification address the problem of identifying a desired fixpoint? Show how every relational algebra query can be written as a stratified Datalog program. (Section 24.3.1)
- Two important aspects of SQL, *multiset tables* and *aggregation with grouping*, are missing in Datalog. How can we extend Datalog to support these features? Discuss the interaction of these two new features and the need for stratification of aggregation. (Section 24.4)
- Define the terms *inference* and *iteration*. What are the two main challenges in efficient evaluation of recursive Datalog programs? (Section 24.5)
- Describe *Semi-naive fixpoint evaluation* and explain how it avoids repeated inferences. (Section 24.5.1)
- Describe the *Magic Sets* program transformation and explain how it avoids unnecessary inferences. (Sections 24.5.2 and 24.5.3)

EXERCISES

Exercise 24.1 Consider the Flights relation:

Flights(fno: integer, from: string, to: string, distance: integer,
departs: time, arrives: time)

Write the following queries in Datalog and SQL:1999 syntax:

1. Find the *fno* of all flights that depart from Madison.
2. Find the *fno* of all flights that leave Chicago after Flight 101 arrives in Chicago and no later than 1 hour after.
3. Find the *fno* of all flights that do not depart from Madison.
4. Find all cities reachable from Madison through a series of one or more connecting flights.
5. Find all cities reachable from Madison through a chain of one or more connecting flights, with no more than 1 hour spent on any connection. (That is, every connecting flight must depart within an hour of the arrival of the previous flight in the chain.)
6. Find the shortest time to fly from Madison to Madras, using a chain of one or more connecting flights.
7. Find the *fno* of all flights that do not depart from Madison or a city that is reachable from Madison through a chain of flights.

Exercise 24.2 Consider the definition of Components in Section 24.1.1. Suppose that the second rule is replaced by

Components(Part, Subpart) :- Components(Part, Part2),
 Components(Part2, Subpart).

1. If the modified program is evaluated on the Assembly relation in Figure 24.1, how many iterations does Naive fixpoint evaluation take and what Components facts are generated in each iteration?
2. Extend the given instance of Assembly so that Naive fixpoint iteration takes two more iterations.
3. Write this program in SQL:1999 syntax, using the WITH clause.
4. Write a program in Datalog syntax to find the part with the most distinct subparts; if several parts have the same maximum number of subparts, your query should return all these parts.
5. How would your answer to the previous part be changed if you also wanted to list the number of subparts for the part with the most distinct subparts?
6. Rewrite your answers to the previous two parts in SQL:1999 syntax.
7. Suppose that you want to find the part with the most subparts, taking into account the quantity of each subpart used in a part, how would you modify the Components program? (*Hint:* To write such a query you reason about the number of inferences of a fact. For this, you have to rely on SQL's maintaining as many copies of each fact as the number of inferences of that fact and take into account the properties of Seminaive evaluation.)

Exercise 24.3 Consider the definition of Components in Exercise 24.2. Suppose that the recursive rule is rewritten as follows for Seminaive fixpoint evaluation:

Components(Part, Subpart) :- delta_Components(Part, Part2, Qty),
 delta_Components(Part2, Subpart).

1. At the end of an iteration, what steps must be taken to update *delta_Components* to contain just the new tuples generated in this iteration? Can you suggest an index on Components that might help to make this faster?
2. Even if the *delta* relation is correctly updated, fixpoint evaluation using the preceding rule does not always produce all answers. Show an instance of Assembly that illustrates the problem.
3. Can you suggest a way to rewrite the recursive rule in terms of *delta_Components* so that Seminaive fixpoint evaluation always produces all answers and no inferences are repeated across iterations?
4. Show how your version of the rewritten program performs on the example instance of Assembly that you used to illustrate the problem with the given rewriting of the recursive rule.

Exercise 24.4 Consider the definition of SameLevel in Section 24.5.2 and the Assembly instance shown in Figure 24.1.

1. Rewrite the recursive rule for Seminaive fixpoint evaluation and show how Seminaive evaluation proceeds.
2. Consider the rules defining the relation Magic, with *spoke* as the query constant. For Seminaive evaluation of the 'Magic' version of the SameLevel program, all tuples in Magic are computed first. Show how Seminaive evaluation of the Magic relation proceeds.

3. After the Magic relation is computed, it can be treated as a fixed database relation, just like Assembly, in the Semaive fixpoint evaluation of the rules defining SameLevel in the 'Magic' version of the program. Rewrite the recursive rule for Semaive evaluation and show how Semaive evaluation of these rules proceeds.

Exercise 24.5 Consider the definition of SameLevel in Section 24.5.2 and a query in which the first argument is bound. Suppose that the recursive rule is rewritten as follows, leading to multiple binding patterns in the adorned program:

```
SameLevel(S1, S2) :- Assembly(P1, S1, Q1),
                    Assembly(P1, S2, Q2).
SameLevel(S1, S2) :- Assembly(P1, S1, Q1),
                    SameLevel(P2, P1), Assembly(P2, S2, Q2).
```

1. Show the adorned program.
2. Show the Magic program.
3. Show the Magic program after applying Semaive rewriting.
4. Construct an example instance of Assembly such that the evaluating the optimized program generates less than 1% of the facts generated by evaluating the original program (and finally selecting the query result).

Exercise 24.6 Again, consider the definition of SameLevel in Section 24.5.2 and a query in which the first argument is bound. Suppose that the recursive rule is rewritten as follows:

```
SameLevel(S1, S2) :- Assembly(P1, S1, Q1),
                    Assembly(P1, S2, Q2).
SameLevel(S1, S2) :- Assembly(P1, S1, Q1),
                    SameLevel(P1, R1), SameLevel(R1, P2), Assembly(P2, S2, Q2).
```

1. Show the adorned program.
2. Show the Magic program.
3. Show the Magic program after applying Semaive rewriting.
4. Construct an example instance of Assembly such that the evaluating the optimized program generates less than 1% of the facts generated by evaluating the original program (and finally selecting the query result).

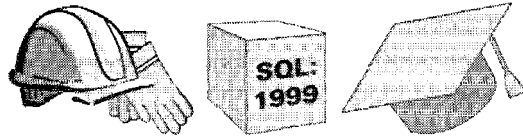
BIBLIOGRAPHIC NOTES

The use of logic as a query language is discussed in several papers [296, 537], "which arose out of influential workshops. Good textbook discussions of deductive databases can be found in [747, 3, 143, 794, 503]. [614] is a recent survey article that provides an overview and covers the major prototypes in the area, including LIDL [177], Glue-Nail! [214, 549] EKS-VI [758], Aditi [615], Coral [612], LOLA [804], and XSB [644].

The fixpoint semantics of logic programs (and deductive databases as a special case) is presented in [751], which also shows equivalence of the fixpoint semantics to a *least-model* semantics. The use of stratification to give a natural semantics to programs with negation was developed independently in [37, 154, 559, 752].

Efficient evaluation of deductive database queries has been widely studied, and [58] is a survey and comparison of several early techniques; [61] is a more recent survey. Serrinaive fixpoint evaluation was independently proposed several times; a good treatment appears in [54]. The Magic Sets technique is proposed in [57] and generalized to cover all deductive database queries without negation in [77]. The Alexander method [63] was independently developed and is equivalent to a variant of Magic Sets called *Supplementary Magic Sets* in [77]. [55] shows how Magic Sets offers significant performance benefits even for nonrecursive SQL queries. [67] describes a version of Magic Sets designed for SQL queries with correlation, and its implementation in the Starburst system (which led to its implementation in IBM's DB2 DBMS). [60] discusses how Magic Sets can be incorporated into a System R style cost-based optimization framework. The Magic Sets technique is extended to programs with stratified negation in [53, 76]. [12] compares Magic Sets with top-down evaluation strategies derived from Prolog.

[64] develops a program rewriting technique related to Magic Sets called *Magic Counting*. Other related methods that are not based on program rewriting but rather on run-time control strategies for evaluation include [226, 429, 756, 757]. The ideas in [226] have been developed further to design an *abstract machine* for logic program evaluation using tabling in [609, 727]; this is the basis for the XSB system [644].



25

DATA WAREHOUSING AND DECISION SUPPORT

- ☛ Why are traditional DBMSs inadequate for decision support?
- ☛ What is the multidimensional data model and what kinds of analysis does it facilitate?
- ☛ What SQL:1999 features support multidimensional queries?
- ☛ How does SQL:1999 support analysis of sequences and trends?
- ☛ How are DBMSs being optimized to deliver early answers for interactive analysis?
- ☛ What kinds of index and file organizations do OLAP systems require?
- ☛ What is data warehousing and why is it important for decision support?
- ☛ Why have materialized views become important?
- ☛ How can we efficiently maintain materialized views?
- ➡ **Key concepts:** OLAP, multidimensional model, dimensions, measures; roll-up, drill-down, pivoting, cross-tabulation, CUBE; WINDOW queries, frames, order; top N queries, online aggregation; bitmap indexes, join indexes; data warehouses, extract, refresh, purge; materialized views, incremental maintenance, maintaining warehouse views

Nothing is more difficult, and therefore more precious, than to be able to decide.

. Napoleon Bonaparte

Database management systems are widely used by organizations for maintaining data that documents their everyday operations. In applications that update such *operational data*, transactions typically make small changes (for example, adding a reservation or depositing a check) and a large number of transactions must be reliably and efficiently processed. Such online transaction processing (OLTP) applications have driven the growth of the DBMS industry in the past three decades and will doubtless continue to be important. DBMSs have traditionally been optimized extensively to perform well in such applications.

Recently, however, organizations have increasingly emphasized applications in which current and historical data is comprehensively analyzed and explored, identifying useful trends and creating summaries of the data, in order to support high-level decision making. Such applications are referred to as decision support. Mainstream relational DBMS vendors have recognized the importance of this market segment and are adding features to their products to support it. In particular, SQL has been extended with new constructs and novel indexing and query optimization techniques are being added to support complex queries.

The use of views has gained rapidly in popularity because of their utility in applications involving complex data analysis. While queries on views can be answered by evaluating the view definition when the query is submitted, precomputing the view definition can make queries run much faster. Carrying the motivation for precomputed views one step further, organizations can consolidate information from several databases into a *data warehouse* by copying tables from many sources into one location or materializing a view defined over tables from several sources. Data warehousing has become widespread, and many specialized products are now available to create and manage warehouses of data from multiple databases.

We begin this chapter with an overview of decision support in Section 25.1. We introduce the multidimensional model of data in Section 25.2 and consider database design issues in 25.2.1. We discuss the rich class of queries that it naturally supports in Section 25.3. We discuss how new SQL:1999 constructs allow us to express multidimensional queries in 25.3.1. In Section 25.4, we discuss SQL:1999 extensions that support queries over relations as ordered collections. We consider how to optimize for fast generation of initial answers in Section 25.5. The many query language extensions required in the OLAP environment prompted the development of MVS implementation techniques; we discuss these in Section 25.6. In Section 25.7, we examine the issues involved in creating and maintaining a data warehouse. From a technical standpoint, a key issue is how to maintain warehouse information (replicated tables or views) when the underlying source information changes. After covering the important role played by MVS in OLAP and warehousing in Section 25.8, we consider maintenance of materialized views in Sections 25.9 and 25.10.

25.1 INTRODUCTION TO DECISION SUPPORT

Organizational decision making requires a comprehensive view of all aspects of an enterprise, so many organizations created consolidated data warehouses that contain data drawn from several databases maintained by different business units together with historical and summary information.

The trend toward data warehousing is complemented by an increased emphasis on powerful analysis tools. Many characteristics of decision support queries make traditional SQL systems inadequate:

- The WHERE clause often contains many AND and OR conditions. As we saw in Section 14.2.3, OR conditions, in particular, are poorly handled in many relational DBMSs.
- Applications require extensive use of statistical functions, such as standard deviation, that are not supported in SQL-92. Therefore, SQL queries must frequently be embedded in a host language program.
- Many queries involve conditions over time or require aggregating over time periods. SQL-92 provides poor support for such time-series analysis.
- Users often need to pose several related queries. Since there is no convenient way to express these commonly occurring families of queries, users have to write them as a collection of independent queries, which can be tedious. Further, the DBMS has no way to recognize and exploit optimization opportunities arising from executing many related queries together.

Three broad classes of analysis tools are available. First, some systems support a class of stylized queries that typically involve group-by and aggregation operators and provide excellent support for complex boolean conditions, statistical functions, and features for time-series analysis. Applications dominated by such queries are called online analytic processing (OLAP). These systems support a querying style in which the data is best thought of as a multidimensional array and are influenced by end-user tools, such as spreadsheets, in addition to database query languages.

Second, some DBMSs support traditional SQL-style queries but are designed to also support OLAP queries efficiently. Such systems can be regarded as relational DBMSs optimized for decision support applications. Many vendors of relational DBMSs are currently enhancing their products in this direction and, over time, the distinction between specialized OLAP systems and relational DBMSs enhanced to support OLAP queries is likely to diminish.

The third class of analysis tools is motivated by the desire to find interesting or unexpected trends and patterns in large data sets rather than the complex

SQL:1999 and OLAP: In this chapter, we discuss a number of features introduced in SQL:1999 to support OLAP. In order not to delay publication of the SQL:1999 standard, these features were actually added to the standard through an *amendment* called SQL/OLAP.

query characteristics just listed. In exploratory data analysis, although an analyst can recognize an 'interesting pattern' when shown such a pattern, it is very difficult to formulate a query that captures the essence of an interesting pattern. For example, an analyst looking at credit-card usage histories may want to detect unusual activity indicating misuse of a lost or stolen card. A catalog merchant may want to look at customer records to identify promising customers for a new promotion; this identification would depend on income level, buying patterns, demonstrated interest areas, and so on. The amount of data in many applications is too large to permit manual analysis or even traditional statistical analysis, and the goal of data mining is to support exploratory analysis over very large data sets. We discuss data mining further in Chapter 26.

Clearly, evaluating OLAP or data mining queries over globally distributed data is likely to be excruciatingly slow. Further, for such complex analysis, often statistical in nature, it is not essential that the most current version of the data be used. The natural solution is to create a centralized repository of all the data; that is, a data warehouse. Thus, the availability of a warehouse facilitates the application of OLAP and data mining tools and, conversely, the desire to apply such analysis tools is a strong motivation for building a data warehouse.

25.2 OLAP: MULTIDIMENSIONAL DATA MODEL

OLAP applications are dominated by ad hoc, complex queries. In SQL terms, these are queries that involve group-by and aggregation operators. The natural way to think about typical OLAP queries, however, is in terms of a multidimensional data model. In this section, we present the multidimensional data model and compare it with a relational representation of data. In subsequent sections, we describe OLAP queries in terms of the multidimensional data model and consider some new implementation techniques designed to support such queries.

In the multidimensional data model, the focus is on a collection of numeric measures. Each measure depends on a set of dimensions. We use a running example based on sales data. The measure attribute in our example is *sales*. The dimensions are Product, Location, and Time. Given a product, a location;

and a time, we have at most one associated sales value. If we identify a product by a unique identifier *pid* and, similarly, identify location by *locid* and time by *timeid*, we can think of sales information as being arranged in a three-dimensional array Sales. This array is shown in Figure 25.1; for clarity, we show only the values for a single *locid* value, *locid* = 1, which can be thought of as a slice orthogonal to the *locid* axis.

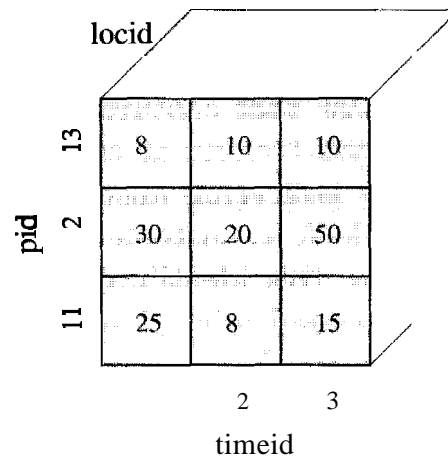


Figure 25.1 Sales: A Multidimensional Dataset

This view of data as a multidimensional array is readily generalized to more than three dimensions. In OLAP applications, the bulk of the data can be represented in such a multidimensional array. Indeed, some OLAP systems actually store data in a multidimensional array (of course, implemented without the usual programming language assumption that the entire array fits in memory). OLAP systems that use arrays to store multidimensional datasets are called multidimensional OLAP (MOLAP) systems.

The data in a multidimensional array can also be represented as a relation, as illustrated in Figure 25.2, which shows the same data as in Figure 25.1, with additional rows corresponding to the 'slice' *locid* = 2. This relation, which relates the dimensions to the measure of interest, is called the fact table.

Now let us turn to dimensions. Each dimension can have a set of associated attributes. For example, the Location dimension is identified by the *locid* attribute, which we used to identify a location in the Sales table. We assume that it also has attributes *country*, *state*, and *city*. We further assume that the Product dimension has attributes *pname*, *category*, and *price* in addition to the identifier *pid*. The *category* of a product indicates its general nature; for example, a product *pant* could have category value *apparel*. We assume that the Time dimension has attributes *date*, *week*, *month*, *quarter*, *year*, and *holiday_flag* in addition to the identifier *timeid*.

<i>locid</i>	<i>city</i>	<i>state</i>	<i>country</i>
1	Madison	WI	USA
2	Fresno	CA	USA
5	Chennai	TN	India

Locations

<i>pid</i>	<i>pname</i>	<i>category</i>	<i>price</i>
11	Lee Jeans	Apparel	25
12	Zord	Toys	18
13	Biro Pen	Stationery	2

Products

<i>pid</i>	<i>timeid</i>	<i>locid</i>	<i>sales</i>
11	1	1	25
11	2	1	8
11	3	1	15
12	1	1	30
12	2	1	20
12	3	1	50
13	1	1	8
13	2	1	10
13	3	1	10
11	1	2	35
11	2	2	22
11	3	2	10
12	1	2	26
12	2	2	45
12	3	2	20
13	1	2	20
13	2	2	40
13	3	2	5

Sales

Figure 25.2 Locations, Products, and Sales Represented as Helations

For each dimension, the set of associated values can be structured as a hierarchy. For example, cities belong to states, and states belong to countries. Dates belong to weeks and months, both weeks and months are contained in quarters, and quarters are contained in years. (Note that a week could span two months; therefore, weeks are not contained in months.) Some of the attributes of a dimension describe the position of a dimension value with respect to this underlying hierarchy of dimension values. The hierarchies for the Product, Location, and Time hierarchies in our example are shown at the attribute level in Figure 25.3.

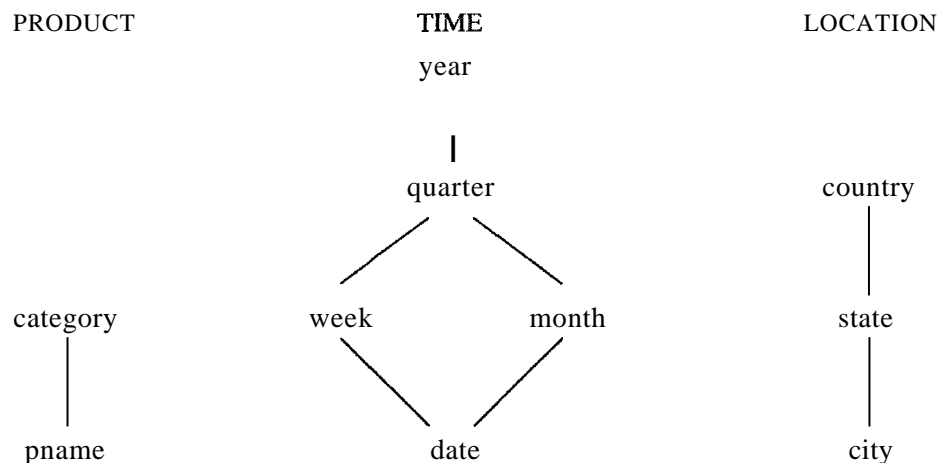


Figure 25.3 Dimension Hierarchies

Information about dimensions can also be represented as a collection of relations:

```

Locations(locid: integer, city: string, state: string, country: string)
Products(pid: integer, pname: string, category: string, price: real)
Times(timeid: integer, date: string, week: integer, month: integer,
      quarter: integer, year: integer, holiday_flag: boolean)

```

These relations are much smaller than the fact table in a typical OLAP application; they are called the **dimension** tables. OLAP systems that store all information, including fact tables, as relations are called relational OLAP (**ROLAP**) systems.

The Times table illustrates the attention paid to the Time dimension in typical OLAP applications. SQL's date and timestamp data types are not adequate; to support summarizations that reflect business operations, information such as fiscal quarters, holiday status, and so on is maintained for each time value.

25.2.1 Multidimensional Database Design

Figure 25.4 shows the tables in our running sales example. It suggests a star, centered at the fact table Sales; such a combination of a fact table and dimension tables is called a star schema. This schema pattern is very common in databases designed for OLAP. The bulk of the data is typically in the fact table, which has no redundancy; it is usually in BCNF. In fact, to minimize the size of the fact table, dimension identifiers (such as *pid* and *timeid*) are system-generated identifiers.

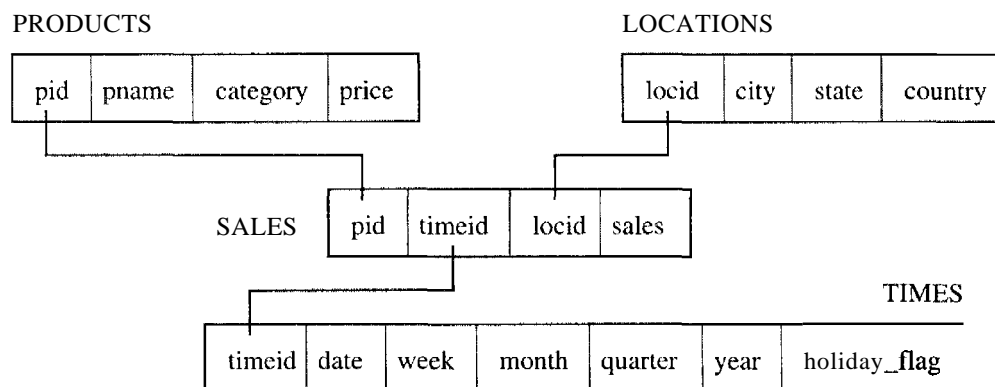


Figure 25.4 An Example of a Star Schema

Information about dimension values is maintained in the dimension tables. Dimension tables are usually not normalized. The rationale is that the dimension tables in a database used for OLAP are static and update, insertion, and deletion anomalies are not important. Further, because the size of the database is dominated by the fact table, the space saved by normalizing dimension tables is negligible. Therefore, minimizing the computation time for combining facts in the fact table with dimension information is the main design criterion, which suggests that we avoid breaking a dimension table into smaller tables (which might lead to additional joins).

Small response times for interactive querying are important in OLAP, and most systems support the materialization of summary tables (typically generated through queries using grouping). Ad hoc queries posed by users are answered using the original tables along with precomputed summaries. A very important design issue is which summary tables should be materialized to achieve the best use of available memory and answer commonly asked ad hoc queries with interactive response times. In current OLAP systems, deciding which summary tables to materialize may well be the most important design decision.

Finally, new storage structures and indexing techniques have been developed to support OLAP and they present the database designer with additional physical

design choices. We cover some of these implementation techniques in Section 25.6.

25.3 MULTIDIMENSIONAL AGGREGATION QUERIES

Now that we have seen the multidimensional model of data, let us consider how such data can be queried and manipulated. The operations supported by this model are strongly influenced by end user tools such as spreadsheets. The goal is to give end users who are not SQL experts an intuitive and powerful interface for common business-oriented analysis tasks. Users are expected to pose ad hoc queries directly, without relying on database application programmers.

In this section, we assume that the user is working with a multidimensional dataset and that each operation returns either a different presentation or a summary; the underlying dataset is always available for the user to manipulate, regardless of the level of detail at which it is currently viewed. In Section 25.3.1, we discuss how SQL:1999 provides constructs to express the kinds of queries presented in this section over tabular, relational data.

A very common operation is aggregating a measure over one or more dimensions. The following queries are typical:

- Find the total sales.
- Find total sales for each city.
- Find total sales for each state.

These queries can be expressed as SQL queries over the fact and dimension tables. When we aggregate a measure on one or more dimensions, the aggregated measure depends on fewer dimensions than the original measure. For example, when we compute the total sales by city, the aggregated measure is *total sales* and it depends only on the Location dimension, whereas the original *sales* measure depended on the Location, Time, and Product dimensions.

Another use of aggregation is to summarize at different levels of a dimension hierarchy. If we are given total sales per city, we can aggregate on the Location dimension to obtain sales per state. This operation is called roll-up in the OLAP literature. The inverse of roll-up is drill-down: Given total sales by state, we can ask for a more detailed presentation by drilling down on Location. We can ask for sales by city or just sales by city for a selected state (with sales presented on a per-state basis for the remaining states, as before). We can also drill down on a dimension other than Location. For example, we can ask

for total sales for each product for each state, drilling down on the Product dimension.

Another common operation is pivoting. Consider a tabular presentation of the Sales table. If we pivot it on the Location and Time dimensions, we obtain a table of total sales for each location for each time value. This information can be presented as a two-dimensional chart in which the axes are labeled with location and time values; the entries in the chart correspond to the total sales for that location and time. Therefore, values that appear in columns of the original presentation become labels of axes in the result presentation. The result of pivoting, called a cross-tabulation, is illustrated in Figure 25.5. Observe that in spreadsheet style, in addition to the total sales by year and state (taken together), we also have additional summaries of sales by year and sales by state.

	WI	CA	Total
1995	63	81	144
1996	38	107	145
1997	75	35	110
Total	176	223	399

Figure 25.5 Cross-Tabulation of Sales by Year and State

Pivoting can also be used to change the dimensions of the cross-tabulation; from a presentation of sales by year and state, we can obtain a presentation of sales by product and year.

Clearly, the OLAP framework makes it convenient to pose a broad class of queries. It also gives catchy names to some familiar operations: Slicing a dataset amounts to an equality selection on one or more dimensions, possibly also with some dimensions projected out. Dicing a dataset amounts to a range selection. These terms come from visualizing the effect of these operations on a cube or cross-tabulated representation of the data.

A Note on Statistical Databases

Many OLAP concepts are present in earlier work on statistical databases (SDBs), which are database systems designed to support statistical applications, although this connection has not been sufficiently recognized because of differences in application domains and terminology. The multidimensional data model, with the notions of a measure associated with dimensions and

classification hierarchies for dimension values, is also used in SDBs. OLAP operations such as roll-up and drill-down have counterparts in SDBs. Indeed, some implementation techniques developed for OLAP are also applied to SDBs.

Nonetheless, some differences arise from the different domains OLAP and SDBs were developed to support. For example, SDBs are used in socioeconomic applications, where classification hierarchies and privacy issues are very important. This is reflected in the greater complexity of classification hierarchies in SDBs, along with issues such as potential breaches of privacy. (The privacy issue concerns whether a user with access to summarized data can reconstruct the original, unsunsumarized data.) In contrast, OLAP has been aimed at business applications with large volumes of data and efficient handling of very large datasets has received more attention than in the SDB literature.

25.3.1 ROLLUP and CUBE in SQL:1999

In this section, we discuss how many of the query capabilities of the multidimensional model are supported in SQL:1999. Typically, a single OLAP operation leads to several closely related SQL queries with aggregation and grouping. For example, consider the cross-tabulation shown in Figure 25.5, which was obtained by pivoting the Sales table. To obtain the same information, we would issue the following queries:

```
SELECT  rr.year, l.state, SUM (S.sales)
FROM    Sales S, Times T, Locations L
WHERE   S.timeid=T.timeid AND S.locid=L.locid
GROUP BY T.year, l.state
```

This query generates the entries in the body of the chart (outlined by the dark lines). The summary column on the right is generated by the query:

```
SELECT  T.year, SUM (S.sales)
FROM    Sales S, Times T
WHERE   S.timeid = T.timeid
GROUP BY T.year
```

The summary row at the bottom is generated by the query:

```
SELECT  L.state, SUM (S.sales)
FROM    Sales S, Locations L
WHERE   S.locid=L.locid
GROUP BY L.state
```

The cumulative sum in the bottom-right corner of the chart is produced by the query:

```

SELECT    SUM (S.sales)
FROM      Sales S, Locations L
WHERE     S.locid=L.locid

```

The example cross-tabulation can be thought of as roll-up on the entire dataset (i.e., treating everything as one big group), on the Location dimension, on the Time dimension, and on the Location and Time dimensions together. Each roll-up corresponds to a single SQL query with grouping. In general, given a measure with k associated dimensions, we can roll up on any subset of these k dimensions; so we have a total of 2^k such SQL queries.

Through high-level operations such as pivoting, users can generate many of these 2^k SQL queries. Recognizing the commonalities between these queries enables more efficient, coordinated computation of the set of queries.

SQL:1999 extends the GROUP BY construct to provide better support for roll-up and cross-tabulation queries. The GROUP BY clause with the CUBE keyword is equivalent to a collection of GROUP BY statements, with one GROUP BY statement for each subset of the k dimensions.

Consider the following query:

```

SELECT    T.year, L.state, SUM (S.sales)
FROM      Sales S, Times T, Locations L
WHERE     S.timeid=T.timeid AND S.locid=L.locid
GROUP BY CUBE (T.year, L.state)

```

The result of this query, shown in Figure 25.6, is just a tabular representation of the cross-tabulation in Figure 25.5.

SQL:1999 also provides variants of GROUP BY that enable computation of subsets of the cross-tabulation computed using GROUP BY CUBE. For example, we can replace the grouping clause in the previous query with

```
GROUP BY ROLLUP (T.year, L.state)
```

In contrast to GROUP BY CUBE, we aggregate by all pairs of year and state values and by each year, and compute an overall sum for the entire dataset (the last row in Figure 25.6), but we do not aggregate for each state value. The result is identical to that shown in Figure 25.6, except that the rows with *null* in the *T.year* column and non-*null* values in the *L.state* column are not computed.

```
CUBE T.year, L.state, timeid BY SUM Sales
```

<i>T.year</i>	<i>L.state</i>	<i>SUM(S.sales)</i>
1995	WI	63
1995	CA	81
1995	<i>null</i>	144
1996	WI	38
1996	CA	107
1996	<i>null</i>	145
1997	WI	75
1997	CA	35
1997	<i>null</i>	110
<i>null</i>	WI	176
<i>null</i>	CA	223
<i>null</i>	<i>null</i>	399

Figure 25.6 The Result of GROUP BY CUBE on Sales

This query rolls up the table Sales on all eight subsets of the set {pid, locid, timeid} (including the empty subset). It is equivalent to eight queries of the form

```
SELECT    SUM (S.sales)
FROM      Sales S
GROUP BY grouping-list
```

The queries differ only in the *grouping-list*, which is some subset of the set {pid, locid, timeid}. We can think of these eight queries as being arranged in a lattice, as shown in Figure 25.7. The result tuples at a node can be aggregated further to compute the result for any child of the node. This relationship between the queries arising in a CUBE can be exploited for efficient evaluation.

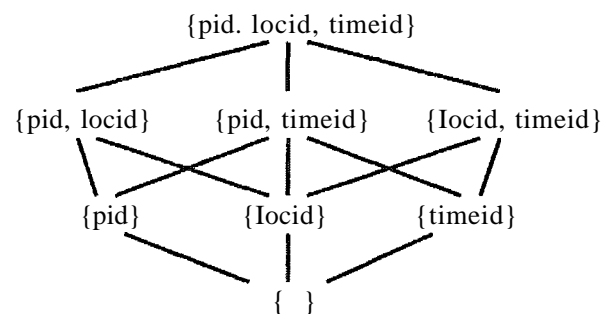


Figure 25.7 The Lattice of GROUP BY Queries in a CUBE Query

25.4 WINDOW QUERIES IN SQL:1999

The time dimension is very important in decision support and queries involving trend analysis have traditionally been difficult to express in SQL. To address this, SQL:1999 introduced a fundamental extension called a query window. Examples of queries that can be written using this extension, but are either difficult or impossible to write in SQL without it, include

1. Find total sales by month.
2. Find total sales by month for each city.
3. Find the percentage change in the total monthly sales for each product.
4. Find the top five products ranked by total sales.
5. Find the trailing n day moving average of sales. (For each day, we must compute the average daily sales over the preceding n days.)
6. Find the top five products ranked by cumulative sales, for every month over the past year.
7. Rank all products by total sales over the past year, and, for each product, print the difference in total sales relative to the product ranked behind it.

The first two queries can be expressed as SQL queries using GROUP BY over the fact and dimension tables. The next two queries can be expressed too, but are quite complicated in SQL-92. The fifth query cannot be expressed in SQL-92 if n is to be a parameter of the query. The last query cannot be expressed in SQL-92.

In this section, we discuss the features of SQL:1999 that allow us to express all these queries and, obviously, a rich class of similar queries.

The main extension is the WINDOW clause, which intuitively identifies an ordered 'window' of rows 'around' each tuple in a table. This allows us to apply a rich collection of aggregate functions to the window of a row and extend the row with the results. For example, we can associate the average sales over the past 3 days with every Sales tuple (each of which records 1 day's sales). This gives us a 3-day moving average of sales.

While there is some similarity to the GROUP BY and CUBE clauses, there are important differences as well. For example, like the WINDOW operator, GROUP BY allows us to create partitions of rows and apply aggregate functions such as SUM to the rows in a partition. However, unlike WINDOW, there is a single output row per partition, rather than one output row for each row, and each partition is an unordered collection of rows.

We now illustrate the window concept through an example:

```
SELECT L.state, T.month, AVG (S.sales) OVER W AS Inovavg
FROM   Sales S, Titles T, Locations L
WHERE  S.titleid=T.titleid AND S.locid=L.locid
WINDOW W AS (PARTITION BY L.state
              ORDER BY 'f.month'
              RANGE BETWEEN INTERVAL '1' MONTH PRECEDING
              AND INTERVAL '1' MONTH FOLLOWING)
```

The FROM and WHERE clauses are processed as usual to (conceptually) generate an intermediate table, which we refer to as Temp. Windows are created over the Temp relation.

There are three steps in defining a window. First, we define *partitions* of the table, using the PARTITION BY clause. In the example, partitions are based on the *L.state* column. Partitions are similar to groups created with GROUP BY, but there is a very important difference in how they are processed. To understand the difference, observe that the SELECT clause contains a column, *T.month*, which is not used to define the partitions; different rows in a given partition could have different values in this column. Such a column cannot appear in the SELECT clause in conjunction with grouping, but it is allowed for partitions. The reason is that there is one answer row for *each* row in a partition of Temp, rather than just one answer row per partition. The window around a given row is used to compute the aggregate functions in the corresponding answer row.

The second step in defining a window is to specify the *ordering* of rows within a partition. We do this using the ORDER BY clause; in the example, the rows within each partition are ordered by *T.month*.

The third step in window definition is to *frame* windows; that is, to establish the boundaries of the window associated with each row in terms of the ordering of rows within partitions. In the example, the window for a row includes the row itself plus all rows whose month value is within a month before or after; therefore, a row whose month value is June 2002 has a window containing all rows with month equal to May, June, or July 2002.

The answer row corresponding to a given row is constructed by first identifying its window. Then, for each answer column defined using a window aggregate function, we compute the aggregate using the rows in the window.

In our example, each row of Temp is essentially a row of Sales, tagged with extra details (about the location and time dimensions). There is one partition for each state and every row of Temp belongs to exactly one partition. Consider

a row for a store in Wisconsin. The row states the sales for a given product, in that store, at a certain time. The window for this row includes all rows that describe sales in Wisconsin within the previous or next 1 month and *movavg* is the average of sales (over all products) in Wisconsin within this period.

We note that the ordering of rows within a partition for the purposes of window definition does not extend to the table of answer rows. The ordering of answer rows is nondeterministic, unless, of course, we fetch them through a cursor and use `ORDER BY` to order the cursor's output.

25.4.1 Framing a Window

There are two distinct ways to frame a window in SQL:1999. The example query illustrated the `RANGE` construct, which defines a window based on the values in some column (*month* in our example). The ordering column has to be a numeric type, a datetime type, or an interval type since these are the only types for which addition and subtraction are defined.

The second approach is based on using the ordering directly and specifying how many rows before and after the given row are in its window. Thus, we could say

```
SELECT L.state, T.month, AVG (S.sales) OVER W AS Movavg
FROM   Sales S, Times T, Locations L
WHERE  S.timeid=T.timeid AND S.locid=L.locid
WINDOW W AS (PARTITION BY L.state
              ORDER BY T.month
              ROWS BETWEEN 1 PRECEDING AND 1 FOLLOWING)
```

If there is exactly one row in `T` for each `month`, this is equivalent to the previous query. However, if a given `month` has no rows or multiple rows, the two queries produce different results. In this case, the result of the second query is hard to understand because the windows for different rows do not align in a natural way.

The second approach is appropriate if, in terms of our example, there is exactly one row per `month`. Generalizing from this, it is also appropriate if there is exactly one row for every value in the sequence of ordering column values. Unlike the first approach, where the ordering has to be specified over a single (numeric, datetime, or interval type) column, the ordering can be based on a composite key.

We can also define windows that include all rows that are before a given row (UNBOUNDED PRECEDING) or all rows after a given row (UNBOUNDED FOLLOWING) within the row's partition.

25.4.2 New Aggregate Functions

While the standard aggregate functions that apply to multisets of values (e.g., SUM, AVG) can be used in conjunction with `ORDER BY`, there is a need for a new class of functions that operate on a *list* of values.

The RANK function returns the position of a row within its partition. If a partition has 15 rows, the first row (according to the ordering of rows in the window definition over this partition) has rank 1 and the last row has rank 15. The rank of intermediate rows depends on whether there are multiple (or no) rows for a given value of the ordering column.

Consider our running example. If the first row in the Wisconsin partition has the 11th January 2002, and the second and third rows both have the 1st February 2002, then their ranks are 1, 2, and 2, respectively. If the next row has 11th March 2002 its rank is 4.

In contrast, the DENSE_RANK function generates ranks without gaps. In our example, the four rows are given ranks 1, 2, 2, and 3. The only change is in the fourth row, whose rank is now 3 rather than 4.

The PERCENT_RANK function gives a measure of the relative position of a row within a partition. It is defined as (RANK-1) divided by the number of rows in the partition. CUME_DIST is similar but based on actual position within the ordered partition rather than rank.

25.5 FINDING ANSWERS QUICKLY

A recent trend, fueled in part by the popularity of the Internet, is an emphasis on queries for which a user wants only the first few, or the 'best' few, answers quickly. When users pose queries to a search engine such as AltaVista, they rarely look beyond the first or second page of results. If they do not find what they are looking for, they refine their query and resubmit it. The same phenomenon occurs in decision support applications and some DBMS products (e.g., DB2) already support extended SQL constructs to specify such queries. A related trend is that, for complex queries, users would like to see an approximate answer quickly and then have it be continually refined, rather than wait until the exact answer is available. We now discuss these two trends briefly.

25.5.1 Top N Queries

An analyst often wants to identify the top-selling handful of products, for example. We can sort by sales for each product and return answers in this order. If we have a billion products and the analyst is interested only in the top 10, this straightforward evaluation strategy is clearly wasteful. It is desirable for users to be able to explicitly indicate how many answers they want, making it possible for the DBMS to optimize execution. The following example query asks for the top 10 products ordered by sales in a given location and time:

```
SELECT    P.pid, P.pname, S.sales
FROM      Sales S, Products P
WHERE     S.pid=P.pid AND S.locid=1 AND S.timeid=3
ORDER BY  S.sales DESC
OPTIMIZE FOR 10 ROWS
```

The OPTIMIZE FOR N ROWS construct is not in SQL-92 (or even SQL:1999), but it is supported in IBM's DB2 product, and other products (e.g., Oracle 9i) have similar constructs. In the absence of a cue such as OPTIMIZE FOR 10 ROWS, the DBMS computes sales for all products and returns them in descending order by sales. The application can close the result cursor (i.e., terminate the query execution) after consulting 10 rows, but considerable effort has already been expended in computing sales for all products and sorting them.

Now let us consider how a DBMS can use the OPTIMIZE FOR cue to execute the query efficiently. The key is to somehow compute sales only for products that are likely to be in the top 10 by sales. Suppose that we know the distribution of sales values because we maintain a histogram on the *sales* column of the *Sales* relation. We can then choose a value of *sales*, say, *c*, such that only 10 products have a larger sales value. For those *Sales* tuples that meet this condition, we can apply the location and time conditions as well and sort the result. Evaluating the following query is equivalent to this approach:

```
SELECT    P.pid, P.pname, S.sales
FROM      Sales S, Products P
WHERE     S.pid=P.pid AND S.locid=1 AND S.timeid=3 AND S.sales > c
ORDER BY  S.sales DESC
```

This approach is, of course, much faster than the alternative of computing all product sales and sorting them, but there are some important problems to resolve:

1. *How do we choose the sales cutoff value *c*?* Histograms and other system statistics can be used for this purpose, but this can be a tricky issue. For

one thing, the statistics maintained by a DBMS are only approximate. For another, even if we choose the cutoff to reflect the top 10 sales values accurately, other conditions in the query may eliminate some of the selected tuples, leaving us with fewer than 10 tuples in the result.

2. *What if we have more than 10 tuples in the result?* Since the choice of the cutoff c is approximate, we could get more than the desired number of tuples in the result. This is easily handled by returning just the top 10 to the user. We still save considerably with respect to the approach of computing sales for all products, thanks to the conservative pruning of irrelevant sales information, using the cutoff c .
3. *What if we have fewer than 10 tuples in the result?* Even if we choose the sales cutoff c conservatively, we could still compute fewer than 10 result tuples. In this case, we can re-execute the query with a smaller cutoff value c_2 or simply re-execute the original query with no cutoff.

The effectiveness of the approach depends on how well we can estimate the cutoff and, in particular, on minimizing the number of tuples we obtain fewer than the desired number of result tuples.

25.5.2 Online Aggregation

Consider the following query, which asks for the average sales amount by state:

```
SELECT  L.state, AVG (S.sales)
FROM    Sales S, Locations L
WHERE   S.locid=L.locid
GROUP BY L.state
```

This can be an expensive query if Sales and Locations are large relations. We cannot achieve fast response times with the traditional approach of computing the answer in its entirety when the query is presented. One alternative, as we have seen, is to use precomputation. Another alternative is to compute the answer to the query when the query is presented but return an approximate answer to the user as soon as possible. As the computation progresses, the answer quality is continually refined. This approach is called online aggregation. It is very attractive for queries involving aggregation, because efficient techniques for computing and refining approximate answers are available.

Online aggregation is illustrated in Figure 25.8: For each state—the grouping criterion for our example query—the current value for average sales is displayed, together with a confidence interval. The entry for Alaska tells us that the








STATUS	PRIORITIZE	State	VG(sales)	Conf	Interval
		Alabama	5,232.5	97%	103.4
		Alaska	2,832.5	93%	132.2
		Arizona	6,432.5	98%	52.3
		Wyoming	4,243.5		

Figure 25.8 Online Aggregation

current estimate of average per-store sales in Alaska is \$2,832.50, and that this is within the range \$2,700.30 to \$2,964.70 with 93% probability. The status bar in the first column indicates how close we are to arriving at an exact value for the average sales and the second column indicates whether calculating the average sales for this state is a priority. Estimating average sales for Alaska is not a priority, but estimating it for Arizona is a priority. As the figure indicates, the DBMS devotes more system resources to estimating the average sales for high-priority states; the estimate for Arizona is much tighter than that for Alaska and holds with a higher probability. Users can set the priority for a state by clicking on the Prioritize button at any time during the execution. This degree of interactivity, together with the continuous feedback provided by the visual display, makes online aggregation an attractive technique.

To implement online aggregation, a DBMS must incorporate statistical techniques to provide confidence intervals for approximate answers and use non-blocking algorithms for the relational operators. An algorithm is said to block if it does not produce output tuples until it has consumed all its input tuples. For example, the sort-merge join algorithm blocks because sorting requires all input tuples before determining the first output tuple. Nested loops join and hash join are therefore preferable to sort-merge join for online aggregation. Similarly, hash-based aggregation is better than sort-based aggregation.

25.6 IMPLEMENTATION TECHNIQUES FOR OLAP

In this section we survey some implementation techniques motivated by the OLAP environment. The goal is to provide a feel for how OLAP systems differ from more traditional SQL systems; our discussion is far from comprehensive.

Beyond B+ Trees: Complex queries have motivated the addition of powerful indexing techniques to DBMSs. In addition to B+ tree indexes, Oracle 9i supports bitmap and join indexes and maintains these dynamically as the indexed relations are updated. Oracle 9i also supports indexes on expressions over attribute values, such as $10 * sal + bonus$. Microsoft SQL Server uses bitmap indexes. Sybase IQ supports several kinds of bitmap indexes, and may shortly add support for a linear hashing based index. Informix UDS supports R trees and Informix XPS supports bitmap indexes.

The mostly-read environment of OLAP systems makes the CPU overhead of maintaining indexes negligible and the requirement of interactive response times for queries over very large datasets makes the availability of suitable indexes very important. This combination of factors has led to the development of new indexing techniques. We discuss several of these techniques. We then consider file organizations and other OLAP implementation issues briefly.

We note that the emphasis on query processing and decision support applications in OLAP systems is being complemented by a greater emphasis on evaluating complex SQL queries in traditional SQL systems. Traditional SQL systems are evolving to support OLAP-style queries more efficiently, supporting constructs (e.g., CUBE and window functions) and incorporating implementation techniques previously found only in specialized OLAP systems.

25.6.1 Bitmap Indexes

Consider a table that describes customers:

`Customers(custid: integer, name: string, gender: boolean, rating: integer)`

The *rating* value is an integer in the range 1 to 5, and only two values are recorded for *gender*. Columns with few possible values are called sparse. We can exploit sparsity to construct a new kind of index that greatly speeds up queries on these columns.

The idea is to record values for sparse columns as a sequence of bits, one for each possible value. For example, a *gender* value is either 10 or 01; a 1 in the first position denotes male, and 1 in the second position denotes female. Similarly, 10000 denotes the *rating* value 1, and 00001 denotes the *rating* value 5.

If we consider the *gender* values for all rows in the Customers table, we can treat this as a collection of two bit vectors, one of which has the associated value M(ale) and the other the associated value F(emale). Each bit vector has one bit per row in the Customers table, indicating whether the value in that row is the value associated with the bit vector. The collection of bit vectors for a column is called a bitmap index for that column.

An example instance of the Customers table, together with the bitmap indexes for *gender* and *rating*, is shown in Figure 25.9.

<i>M</i>	<i>F</i>	<i>custid</i>	<i>name</i>	<i>gender</i>	<i>rating</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>
1	0	112	Joe	M	3	0	0	1	0	0
1	0	115	Rain	M	5	0	0	0	0	1
0	1	119	Sue	F	5	0	0	0	0	1
1	0	112	Woo	M	4	0	0	0	1	0

Figure 25.9 Bitmap Indexes on the Customers Relation

Bitmap indexes offer two important advantages over conventional hash and tree indexes. First, they allow the use of efficient bit operations to answer queries. For example, consider the query, "How many male customers have a rating of 5?" We can take the first bit vector for *gender* and do a bitwise AND with the fifth bit vector for *rating* to obtain a bit vector that has 1 for every male customer with rating 5. We can then count the number of 1s in this bit vector to answer the query. Second, bitmap indexes can be much more compact than a traditional B+ tree index and are very amenable to the use of compression techniques.

Bit vectors correspond closely to the rid-lists used to represent data entries in Alternative (3) for a traditional B+ tree index (see Section 8.2). In fact, we can think of a bit vector for a given *age* value, say, as an alternative representation of the rid-list for that value.

This suggests a way to combine bit vectors (and their advantages of bitwise processing) with B+ tree indexes: We can use Alternative (3) for data entries, using a bit vector representation of rid-lists. A caveat is that, if an rid-list is very small, the bit vector representation may be much larger than a list of rid values, even if the bit vector is compressed. Further, the use of compression leads to decompression costs, offsetting some of the computational advantages of the bit vector representation.

A more flexible approach is to use a standard list representation of the rid-list for some key values (intuitively, those that contain few elements) and a bit

vector representation for other key values (those that contain many elements, and therefore lend themselves to a compact bit vector representation).

This hybrid approach, which can easily be adapted to work with hash indexes as well as B+ tree indexes, has both advantages and disadvantages relative to a standard list of rids approach:

1. It can be applied even to columns that are not sparse; that is, in which all possible values can appear. The index levels (or the hashing scheme) allow us to quickly find the 'list' of rids, in a standard list or bit vector representation, for a given key value.
2. Overall, the index is more compact because we can use a bit vector representation for long rid lists. We also have the benefits of fast bit vector processing.
3. On the other hand, the bit vector representation of an rid list relies on a mapping from a position in the vector to an rid. (This is true of any bit vector representation, not just the hybrid approach.) If the set of rows is static, and we do not worry about inserts and deletes of rows, it is straightforward to ensure this by assigning contiguous rids for rows in a table. If inserts and deletes must be supported, additional steps are required. For example, we can continue to assign rids contiguously on a per-table basis and simply keep track of which rids correspond to deleted rows. Bit vectors can now be longer than the current number of rows, and periodic reorganization is required to compact the 'holes' in the assignment of rids.

25.6.2 Join Indexes

Computing joins with small response times is extremely hard for very large relations. One approach to this problem is to create an index designed to speed up specific join queries. Suppose that the Customers table is to be joined with a table called Purchases (recording purchases made by customers) on the *custid* field. We can create a collection of $\langle c, p \rangle$ pairs, where p is the rid of a Purchases record that joins with a Customers record with *custid* c .

This idea can be generalized to support joins over more than two relations. We discuss the special case of a star schema, in which the fact table is likely to be joined with several dimension tables. Consider a join query that joins fact table F with dimension tables D_1 and D_2 and includes selection conditions on column C_1 of table D_1 and column C_2 of table D_2 . We store a tuple $\langle r_1, c_1, r \rangle$ in the join index if r_1 is the rid of a tuple in table D_1 with value c_1 in column C_1 , r_2 is the rid of a tuple in table D_2 with value c_2 in column C_2 , and r is the rid of a tuple in the fact table F , and these three tuples join with each other.

Complex Queries: The IBM DB2 optimizer recognizes star join queries and performs rid-based semijoins (using Bloom filters) to filter the fact table. Then fact table rows are rejoined to the dimension tables. Complex (runt-time) dimension queries (called *snowflake queries*) are supported. DB2 also supports CUBE using smart algorithms that minimize sorts. Microsoft SQL Server optimizes star join queries extensively. It considers taking the cross-product of small dimension tables before joining with the fact table, the use of join indexes, and rid-based semijoins. Oracle 9i also allows users to create dimensions to declare hierarchies and functional dependencies. It supports the CUBE operator and optimizes star join queries by eliminating joins when no column of a dimension table is part of the query result. DBMS products have also been developed specifically for decision support applications, such as Sybase IQ.

The drawback of a join index is that the number of indexes can grow rapidly if several columns in each dimension table are involved in selections and joins with the fact table. An alternative kind of join index avoids this problem. Consider our example involving fact table F and dimension tables D_1 and D_2 . Let G_1 be a column of D_1 on which a selection is expressed in some query that joins D_1 with F . Conceptually, we now join F with D_1 to extend the fields of F with the fields of D_1 , and index F on the 'virtual field' G_1 : If a tuple of D_1 with value c_1 in column C_1 joins with a tuple of F with rid r , we add a tuple (C_1, r) to the join index. We create one such join index for each column of either D_1 or D_2 that involves a selection in some join with F ; C_1 is an example of such a column.

The price paid with respect to the previous version of join indexes is that join indexes created in this way have to be combined (rid intersection) to deal with the join queries of interest to us. This can be done efficiently if we make the new indexes *bitmap* indexes; the result is called a *bitmapped join index*. The idea works especially well if columns such as C_1 are sparse, and therefore well suited to bitmap indexing.

25.6.3 File Organizations

Since many OLAP queries involve just a few columns of a large relation, vertical partitioning becomes attractive. However, storing a relation column-wise can degrade performance for queries that involve several columns. An alternative in a mostly-read environment is to store the relation row-wise, but also store each column separately.