

**Lock Granularity:** Some database systems allow programmers to override the default mechanism for choosing a lock granularity. For example, Microsoft SQL Server allows users to select page locking instead of table locking, using the keyword **PAGLOCK**. IBM's DB2 UDB allows for explicit table-level locking.

a certain number of locks at that granularity, to start obtaining locks at the next higher granularity (e.g., at the page level). This procedure is called **lock escalation**.

## 17.6 CONCURRENCY CONTROL WITHOUT LOCKING

Locking is the most widely used approach to concurrency control in a DBMS, but it is not the only one. We now consider some alternative approaches.

### 17.6.1 Optimistic Concurrency Control

Locking protocols take a pessimistic approach to conflicts between transactions and use either transaction abort or blocking to resolve conflicts. In a system with relatively light contention for data objects, the overhead of obtaining locks and following a locking protocol must nonetheless be paid.

In optimistic concurrency control, the basic premise is that most transactions do not conflict with other transactions, and the idea is to be as permissive as possible in allowing transactions to execute. Transactions proceed in three phases:

1. **Read:** The transaction executes, reading values from the database and writing to a private workspace.
2. **Validation:** If the transaction decides that it wants to commit, the DBMS checks whether the transaction could possibly have conflicted with any other concurrently executing transaction. If there is a possible conflict, the transaction is aborted; its private workspace is cleared and it is restarted.
3. **Write:** If validation determines that there are no possible conflicts, the changes to data objects made by the transaction in its private workspace are copied into the database.

If, indeed, there are few conflicts, and validation can be done efficiently, this approach should lead to better performance than locking. If there are many

conflicts, the cost of repeatedly restarting transactions (thereby wasting the work they've done) hurts performance significantly.

Each transaction  $T_i$  is assigned a timestamp  $TS(T_i)$  at the beginning of its validation phase, and the validation criterion checks whether the timestamp ordering of transactions is an equivalent serial order. For every pair of transactions  $T_i$  and  $T_j$  such that  $TS(T_i) < TS(T_j)$ , one of the following validation conditions must hold:

1.  $T_i$  completes (all three phases) before  $T_j$  begins.
2.  $T_i$  completes before  $T_j$  starts its Write phase, and  $T_i$  does not write any database object read by  $T_j$ .
3.  $T_i$  completes its Read phase before  $T_j$  completes its Read phase, and  $T_i$  does not write any database object that is either read or written by  $T_j$ .

To validate  $T_j$ , we must check to see that one of these conditions holds with respect to each committed transaction  $T_i$  such that  $TS(T_i) < TS(T_j)$ . Each of these conditions ensures that  $T_j$ 's modifications are not visible to  $T_i$ .

Further, the first condition allows  $T_j$  to see some of  $T_i$ 's changes, but clearly, they execute completely in serial order with respect to each other. The second condition allows  $T_j$  to read objects while  $T_i$  is still modifying objects, but there is no conflict because  $T_j$  does not read any object modified by  $T_i$ . Although  $T_j$  might overwrite some objects written by  $T_i$ , all of  $T_i$ 's writes precede all of  $T_j$ 's writes. The third condition allows  $T_i$  and  $T_j$  to write objects at the same time and thus have even more overlap in time than the second condition, but the sets of objects written by the two transactions cannot overlap. Thus, no RW, WR, or WW conflicts are possible if any of these three conditions is met.

Checking these validation criteria requires us to maintain lists of objects read and written by each transaction. Further, while one transaction is being validated, no other transaction can be allowed to commit; otherwise, the validation of the first transaction might miss conflicts with respect to the newly committed transaction. The Write phase of a validated transaction must also be completed (so that its effects are visible outside its private workspace) before other transactions can be validated.

A synchronization mechanism such as a critical section can be used to ensure that at most one transaction is in its (combined) Validation/Write phases at any time. (When a process is executing a critical section in its code, the system suspends all other processes.) Obviously, it is important to keep these phases as short as possible in order to minimize the impact on concurrency. If copies of modified objects have to be copied from the private workspace, this

can make the Write phase long. An alternative approach (which carries the penalty of poor physical locality of objects, such as B+ tree leaf pages, that must be clustered) is to use a level of indirection. In this scheme, every object is accessed via a logical pointer, and in the Write phase, we simply switch the logical pointer to point to the version of the object in the private workspace, instead of copying the object.

Clearly, it is not the case that optimistic concurrency control has no overheads; rather, the locking overheads of lock-based approaches are replaced with the overheads of recording read-lists and write-lists for transactions, checking for conflicts, and copying changes from the private workspace. Similarly, the implicit cost of blocking in a lock-based approach is replaced by the implicit cost of the work wasted by restarted transactions.

## Improved Conflict Resolution<sup>1</sup>

Optimistic Concurrency Control using the three validation conditions described earlier is often overly conservative and unnecessarily aborts and restarts transactions. In particular, according to the validation conditions,  $T_i$  cannot write any object read by  $T_j$ . However, since the validation is aimed at ensuring that  $T_i$  logically executes before  $T_j$ , there is no harm if  $T_i$  writes all data items required by  $T_j$  before  $T_j$  reads them.

The problem arises because we have no way to tell when  $T_i$  wrote the object (relative to  $T_j$ 's reading it) at the time we validate  $T_j$ , since all we have is the list of objects written by  $T_i$  and the list read by  $T_j$ . Such false conflicts can be alleviated by a finer-grain resolution of data conflicts, using mechanisms very similar to locking.

The basic idea is that each transaction in the Read phase tells the DBMS about items it is reading, and when a transaction  $T_i$  is committed (and its writes are accepted), the DBMS checks whether any of the items written by  $T_i$  are being read by any (yet to be validated) transaction  $T_j$ . If so, we know that  $T_j$ 's validation must eventually fail. We can either allow  $T_i$  to discover this when it is validated (the die policy) or kill it and restart it immediately (the kill policy).

The details are as follows. Before reading a data item, a transaction enters an access entry in a hash table. The access entry contains the *transaction id*, a *data object id*, and a *modified* flag (initially set to false), and entries are hashed on the data object id. A temporary exclusive lock is obtained on the

---

<sup>1</sup>We thank Alexander Thomasian for writing this section.

hash bucket containing the entry, and the lock is held while the read data item is copied from the database buffer into the private 'workspace of the transaction.

During validation of  $T$  the hash buckets of all data objects accessed by  $T$  are again locked (in exclusive mode) to check if  $T$  has encountered any data conflicts.  $T$  has encountered a conflict if the *modified* flag is set to true in one of its access entries. (This assumes that the 'die' policy is being used; if the 'kill' policy is used,  $T$  is restarted when the flag is set to true.)

If  $T$  is successfully validated, we lock the hash bucket of each object modified by  $T$ , retrieve all access entries for this object, set the *modified* flag to true, and release the lock on the bucket. If the 'kill' policy is used, the transactions that entered these access entries are restarted. We then complete  $T$ 's Write phase.

It seems that the 'kill' policy is always better than the 'die' policy, because it reduces the overall response time and wasted processing. However, executing  $T$  to the end has the advantage that all of the data items required for its execution are prefetched into the database buffer, and restarted executions of  $T$  will not require disk I/O for reads. This assumes that the database buffer is large enough that prefetched pages are not replaced, and, more important, that access invariance prevails; that is, successive executions of  $T$  require the same data for execution. When  $T$  is restarted its execution time is much shorter than before because no disk I/O is required, and thus its chances of validation are higher. (Of course, if a transaction has already completed its Read phase once, subsequent conflicts should be handled using the 'kill' policy because all its data objects are already in the buffer pool.)

## 17.6.2 Timestamp-Based Concurrency Control

In lock-based concurrency control, conflicting actions of different transactions are ordered by the order in which locks are obtained, and the lock protocol extends this ordering on actions to transactions, thereby ensuring serializability. In optimistic concurrency control, a timestamp ordering is imposed on transactions and validation checks that all conflicting actions occurred in the same order.

Timestamps can also be used in another way: Each transaction can be assigned a timestamp at startup, and we can ensure, at execution time, that if action  $a_i$  of transaction  $T_i$  conflicts with action  $a_j$  of transaction  $T_j$ ,  $a_i$  occurs before  $a_j$  if  $TS(T_i) < TS(T_j)$ . If an action violates this ordering, the transaction is aborted and restarted.

To implement this concurrency control scheme, every database object  $O$  is given a read timestamp  $RTS(O)$  and a write timestamp  $WTS(O)$ . If transaction  $T$  wants to read object  $O$ , and  $TS(T) < WTS(O)$ , the order of this read with respect to the most recent write on  $O$  would violate the timestamp order between this transaction and the writer. Therefore,  $T$  is aborted and restarted *with a new, larger timestamp*. If  $TS(T) > WTS(O)$ ,  $T$  reads  $O$ , and  $RTS(O)$  is set to the larger of  $RTS(O)$  and  $TS(T)$ . (Note that a physical change—the change to  $RTS(O)$ —is written to disk and recorded in the log for recovery purposes, even on reads. This write operation is a significant overhead.)

Observe that if  $T$  is restarted with the same timestamp, it is guaranteed to be aborted again, due to the same conflict. Contrast this behavior with the use of timestamps in 2PL for deadlock prevention, where transactions are restarted with the *same* timestamp as before to avoid repeated restarts. This shows that the two uses of timestamps are quite different and should not be confused.

Next, consider what happens when transaction  $T$  wants to write object  $O$ :

1. If  $TS(T) < RTS(O)$ , the write action conflicts with the most recent read action of  $O$ , and  $T$  is therefore aborted and restarted.
2. If  $TS(T) < WTS(O)$ , a naive approach would be to abort  $T$  because its write action conflicts with the most recent write of  $O$  and is out of timestamp order. However, we can safely ignore such writes and continue. Ignoring outdated writes is called the Thomas Write Rule.
3. Otherwise,  $T$  writes  $O$  and  $WTS(O)$  is set to  $TS(T)$ .

## The Thomas Write Rule

We now consider the justification for the Thomas Write Rule. If  $TS(T) < WTS(O)$ , the current write action has, in effect, been made obsolete by the most recent write of  $O$ , which *follows* the current write according to the timestamp ordering. We can think of  $T$ 's write action as if it had occurred immediately *before* the most recent write of  $O$  and was never read by anyone.

If the Thomas Write Rule is not used, that is,  $T$  is aborted in case (2), the timestamp protocol, like 2PL, allows only conflict serializable schedules. If the Thomas Write Rule is used, some schedules are permitted that are not conflict serializable, as illustrated by the schedule in Figure 17.6.<sup>2</sup> Because  $T_2$ 's write follows  $T_1$ 's read and precedes  $T_1$ 's write of the same object, this schedule is not conflict serializable.

---

<sup>2</sup>In the other direction, 2PL permits some schedules that are not allowed by the timestamp algorithm with the Thomas Write Rule; see Exercise 17.7.

$T1$	$T2$
$R(A)$	
	$W(A)$
	Commit
$W(A)$	
Commit	

Figure 17.6 A Serializable Schedule 'rhat Is Not Conflict Serializable

The Thomas Write Rule relies on the observation that  $T2$ 's write is never seen by any transaction and the schedule in Figure 17.6 is therefore equivalent to the serializable schedule obtained by deleting this write action, which is shown in Figure 17.7.

$T1$	$T2$
$R(A)$	
	Commit
$W(A)$	
Commit	

Figure 17.7 A Conflict Serializable Schedule

## Recoverability

Unfortunately, the timestamp protocol just presented permits schedules that are not recoverable, as illustrated by the schedule in Figure 17.8. If  $TS(T1) = 1$  and  $TS(T2) = 2$ , this schedule is permitted by the timestamp protocol (with or without the Thomas Write Rule). The timestamp protocol can be modified to disallow such schedules by buffering all write actions until the transaction commits. In the example, when  $T1$  wants to write  $A$ ,  $WTS(A)$  is updated to reflect this action, but the change to  $A$  is not carried out immediately; instead, it is recorded in a private workspace, or buffer. When  $T2$  wants to read  $A$  subsequently, its timestamp is compared with  $WTS(A)$ , and the read is seen to be permissible. However,  $T2$  is blocked until  $T1$  completes. If  $T1$  commits, its change to  $A$  is copied from the buffer; otherwise, the changes in the buffer are discarded.  $T2$  is then allowed to read  $A$ .

This blocking of  $T2$  is similar to the effect of  $T1$  obtaining an exclusive lock on  $A$ . Nonetheless, even with this modification, the timestamp protocol permits some schedules not permitted by 2PL; the two protocols are not quite the same. (See Exercise 17.7.)

$T1$	$T2$
$W(A)$	$R(A)$
	$W(B)$
	Commit

Figure 17.8 An Unrecoverable Schedule

Because recoverability is essential, such a modification must be used for the timestamp protocol to be practical. Given the added overhead this entails, on top of the (considerable) cost of maintaining read and write timestamps, timestamp concurrency control is unlikely to beat lock-based protocols in centralized systems. Indeed, it has been used mainly in the context of distributed database systems (Chapter 22).

### 17.6.3 Multiversion Concurrency Control

This protocol represents yet another way of using timestamps, assigned at startup time, to achieve serializability. The goal is to ensure that a transaction never has to wait to read a database object, and the idea is to maintain several versions of each database object, each with a write timestamp, and let transaction  $T_i$  read the most recent version whose timestamp precedes  $TS(T_i)$ .

If transaction  $T_i$  wants to write an object, we must ensure that the object has not already been read by some other transaction  $T_j$  such that  $TS(T_i) < TS(T_j)$ . If we allow  $T_i$  to write such an object, its change should be seen by  $T_j$  for serializability, but obviously  $T_j$ , which read the object at some time in the past, will not see  $T_i$ 's change.

To check this condition, every object also has an associated read timestamp, and whenever a transaction reads the object, the read timestamp is set to the maximum of the current read timestamp and the reader's timestamp. If  $T_i$  wants to write an object  $O$  and  $TS(T_i) < RTS(O)$ ,  $T_i$  is aborted and restarted with a new, larger timestamp. Otherwise,  $T_i$  creates a new version of  $O$  and sets the read and write timestamps of the new version to  $TS(T_i)$ .

The drawbacks of this scheme are similar to those of timestamp concurrency control, and in addition, there is the cost of maintaining versions. On the other hand, reads are never blocked, which can be important for workloads dominated by transactions that only read values from the database.

**What Do Real Systems Do?** IBM DB2, Informix, Microsoft SQL Server, and Sybase ASE use Strict 2PL or variants (if a transaction requests a lower than SERIALIZABLE SQL isolation level; see Section 16.6). Microsoft SQL Server also supports modification timestamps so that a transaction can run without setting locks and validate itself (do-it-yourself Optimistic Concurrency Control!). Oracle 8 uses a multiversion concurrency control scheme in which readers never wait; in fact, readers never get locks and detect conflicts by checking if a block changed since they read it. All these systems support multiple-granularity locking, with support for table, page, and row level locks. All deal with deadlocks using waits-for graphs. Sybase ASE supports only table-level locks and aborts a transaction if a lock request fails---updates (and therefore conflicts) are rare in a data warehouse, and this simple scheme suffices.

## 17.7 REVIEW QUESTIONS

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

- When are two schedules *conflict equivalent*? What is a *conflict serializable* schedule? What is a *strict* schedule? (Section 17.1)
- What is a *precedence graph* or *serializability graph*? How is it related to conflict serializability? How is it related to two-phase locking? (Section 17.1)
- What does the *lock manager* do? Describe the *lock table* and *transaction table* data structures and their role in lock management. (Section 17.2)
- Discuss the relative merits of *lock upgrades* and *lock downgrades*. (Section 17.3)
- Describe and compare deadlock detection and deadlock prevention schemes. Why are detection schemes more commonly used? (Section 17.4)
- If the collection of database objects is not fixed, but can grow and shrink through insertion and deletion of objects, we must deal with a subtle complication known as the *phantom problem*. Describe this problem and the index locking approach to solving the problem. (Section 17.5.1)
- In tree index structures, locking higher levels of the tree can become a performance bottleneck. Explain why. Describe specialized locking techniques that address the problem, and explain why they work correctly despite not being two-phase. (Section 17.5.2)
- *Multiple-granularity locking* enables us to set locks on objects that contain other objects, thus implicitly locking all contained objects. Why is this approach important and how does it work? (Section 17.5.3)



- In *optimistic concurrency control*, no locks are set and transactions read and modify data objects in a private workspace. How are conflicts between transactions detected and resolved in this approach? (Section 17.6.1)
- In *timestamp-based concurrency control*, transactions are assigned a timestamp at startup; how is it used to ensure serializability? How does the *Thomas Write Rule* improve concurrency? (Section 17.6.2)
- Explain why timestamp-based concurrency control allows schedules that are not recoverable. Describe how it can be modified through *buffering* to disallow such schedules. (Section 17.6.2)
- Describe *multiversion concurrency control*. What are its benefits and disadvantages in comparison to locking? (Section 17.6.3)

## EXERCISES

Exercise 17.1 Answer the following questions:

1. Describe how a typical lock manager is implemented. Why must lock and unlock be atomic operations? What is the difference between a lock and a *latch*? What are *convoys* and how should a lock manager handle them?
2. Compare *lock downgrades* with upgrades. Explain why downgrades violate 2PL but are nonetheless acceptable. Discuss the use of *update* locks in conjunction with lock downgrades.
3. Contrast the timestamps assigned to restarted transactions when timestamps are used for deadlock prevention versus when timestamps are used for concurrency control.
4. State and justify the Thomas Write Rule.
5. Show that, if two schedules are conflict equivalent, then they are view equivalent.
6. Give an example of a serializable schedule that is not strict.
7. Give an example of a strict schedule that is not serializable.
8. Motivate and describe the use of locks for improved conflict resolution in Optimistic Concurrency Control.

Exercise 17.2 Consider the following classes of schedules: *serializable*, *conflict-serializable*, *view-serializable*, *recoverable*, *avoids-cascading-aborts*, and *strict*. For each of the following schedules, state which of the preceding classes it belongs to. If you cannot decide whether a schedule belongs in a certain class based on the listed actions, explain briefly.

The actions are listed in the order they are scheduled and prefixed with the transaction name. If a commit or abort is not shown, the schedule is incomplete; assume that abort or commit follow all the listed actions.

1. T1:R(X), T2:R(X), T1:W(X), T2:W(X)
2. T1:W(X), T2:R(Y), T1:R(Y), T2:R(X)

3. T1:R(X), T2:R(Y), T3:W(X), T2:R(X), T1:R(Y)
4. T1:R(X), T1:R(Y), T1:W(X), T2:R(Y), T3:W(Y), T1:W(X), T2:R(Y)
5. T1:R(X), T2:W(X), T1:W(X), T2:Abort, T1:Commit
6. T1:R(X), T2:W(X), T1:W(X), T2:Commit, T1:Commit
7. T1:W(X), T2:R(X), T1:W(X), T2:Abort, T1:Commit
8. T1:W(X), T2:R(X), T1:W(X), T2:Commit, T1:Commit
9. T1:W(X), T2:R(X), T1:W(X), T2:Commit, T1:Abort
10. T1:R(X), T1:W(X), T3:Commit, T1:W(Y), T1:Commit, T2:R(Y),  
T2:W(Z), T2:Commit
11. T1:R(X), T2:W(X), T2:Commit, T1:W(X), T1:Commit, T3:R(X), T3:Commit
12. T1:R(X), T2:W(X), T1:W(X), T3:R(X), T1:Commit, T2:Commit, T3:Commit

**Exercise 17.3** Consider the following concurrency control protocols: 2PL, Strict 2PL, Conservative 2PL, Optimistic, Timestamp without the Thomas Write Rule, Timestamp with the Thomas Write Rule, and Multiversion. For each of the schedules in Exercise 17.2, state which of these protocols allows it, that is, allows the actions to occur in exactly the order shown.

For the timestamp-based protocols, assume that the timestamp for transaction  $T_i$  is  $i$  and that a version of the protocol that ensures recoverability is used. Further, if the Thomas Write Rule is used, show the equivalent serial schedule.

**Exercise 17.4** Consider the following sequences of actions, listed in the order they are submitted to the DBMS:

- Sequence 81: T1:R(X), T2:W(X), T2:W(Y), T3:W(Y), T1:W(Y),  
T1:Commit, T2:Commit, T3:Commit
- Sequence 82: T1:R(X), T2:W(Y), T2:W(X), T3:W(Y), T1:W(Y),  
T1:Commit, T2:Commit, T3:Commit

For each sequence and for each of the following concurrency control mechanisms, describe how the concurrency control mechanism handles the sequence.

Assume that the timestamp of transaction  $T_i$  is  $i$ . For lock-based concurrency control mechanisms, add lock and unlock requests to the previous sequence of actions as per the locking protocol. The DBMS processes actions in the order shown. If a transaction is blocked, assume that all its actions are queued until it is rescheduled; the DBMS continues with the next action (according to the listed sequence) of an unblocked transaction.

1. Strict 2PL with timestamps used for deadlock prevention.
2. Strict 2PL with deadlock detection. (Show the waits-for graph in case of deadlock.)
3. Conservative (and Strict, i.e., with locks held until end-of-transaction) 2PL.
4. Optimistic concurrency control.
5. Timestamp concurrency control with buffering of reads and writes (to ensure recoverability) and the Thomas Write Rule.
6. Multiversion concurrency control.

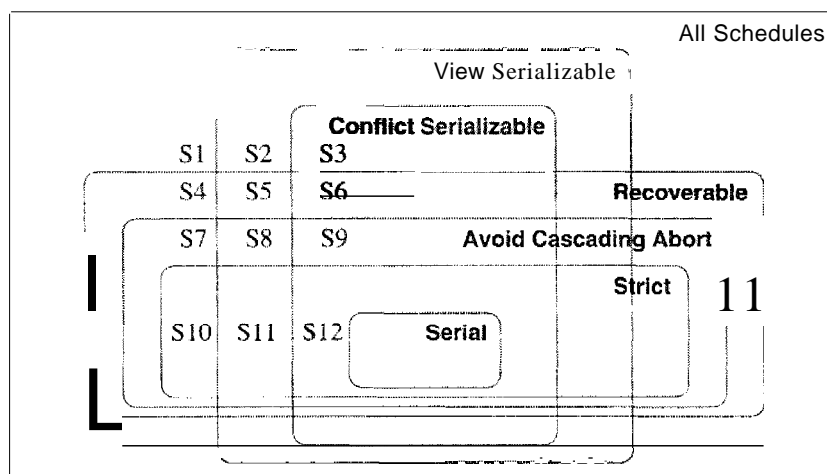


Figure 17.9 Venn Diagram for Classes of Schedules

**Exercise 17.5** For each of the following locking protocols, assuming that every transaction follows that locking protocol, state which of these desirable properties are ensured: serializability, conflict-serializability, recoverability, avoidance of cascading aborts.

1. Always obtain an exclusive lock before writing; hold exclusive locks until end-of-transaction. No shared locks are ever obtained.
2. In addition to (1), obtain a shared lock before reading; shared locks can be released at any time.
3. As in (2), and in addition, locking is two-phase.
4. As in (2), and in addition, all locks held until end-of-transaction.

**Exercise 17.6** The Venn diagram (from [76]) in Figure 17.9 shows the inclusions between several classes of schedules. Give one example schedule for each of the regions S1 through S12 in the diagram.

**Exercise 17.7** Briefly answer the following questions:

1. Draw a Venn diagram that shows the inclusions between the classes of schedules permitted by the following concurrency control protocols: *2PL*, *Strict 2PL*, *Conservative 2PL*, *Optimistic*, *Timestamp without the Thomas Write Rule*, *Timestamp with the Thomas Write Rule*, and *Multiversion*.
2. Give one example schedule for each region in the diagram.
3. Extend the Venn diagram to include serializable and conflict-serializable schedules.

**Exercise 17.8** Answer each of the following questions briefly. The questions are based on the following relational schema:

```
Emp(eid: integer, ename: string, age: integer, salary: real, did: integer)
Dept(did: integer, dname: string, floor: integer)
```

and on the following update command:

```
replace (salary = 1.1 * EMP.salary) where EMP.ename = 'Santa'
```

1. Give an example of a query that would conflict with this command (in a concurrency control sense) if both were run at the same time. Explain what could go wrong, and how locking tuples would solve the problem.
2. Give an example of a query or a command that would conflict with this command, such that the conflict could not be resolved by just locking individual tuples or pages but requires index locking.
3. Explain what index locking is and how it resolves the preceding conflict.

**Exercise 17.9** SQL supports four isolation-levels and two access-modes, for a total of eight combinations of isolation-level and access-mode. Each combination implicitly defines a class of transactions; the following questions refer to these eight classes:

1. For each of the eight classes, describe a locking protocol that allows only transactions in this class. Does the locking protocol for a given class make any assumptions about the locking protocols used for other classes? Explain briefly.
2. Consider a schedule generated by the execution of several SQL transactions. Is it guaranteed to be conflict-serializable? to be serializable? to be recoverable?
3. Consider a schedule generated by the execution of several SQL transactions, each of which has READ ONLY access-mode. Is it guaranteed to be conflict-serializable? to be serializable? to be recoverable?
4. Consider a schedule generated by the execution of several SQL transactions, each of which has SERIALIZABLE isolation-level. Is it guaranteed to be conflict-serializable? to be serializable? to be recoverable?
5. Can you think of a two-phase-based concurrency control scheme that can support the eight classes of SQL transactions?

**Exercise 17.10** Consider the tree shown in Figure 19.5. Describe the steps involved in executing each of the following operations according to the tree-index concurrency control algorithm discussed in Section 19.3.2, in terms of the order in which nodes are locked, unlocked, read, and written. Be specific about the kind of lock obtained and answer each part independently of the others, always starting with the tree shown in Figure 19.5.

1. Search for data entry 40\*.
2. Search for all data entries  $k^*$  with  $k \leq 40$ .
3. Insert data entry 62\*.
4. Insert data entry 40\*.
5. Insert data entries 62\* and 75\*.

**Exercise 17.11** Consider a database organized in terms of the following hierarchy of objects: The database itself is an object ( $D$ ), and it contains two files ( $F1$  and  $F2$ ), each of which contains 1000 pages ( $P1 \dots P1000$  and  $P1001 \dots P2000$ , respectively). Each page contains 100 records, and records are identified as  $p : i$ , where  $p$  is the page identifier and  $i$  is the slot of the record on that page.

Multiple-granularity locking is used, with  $S$ ,  $X$ ,  $IS$ ,  $IX$  and  $SIX$  locks, and database-level, file-level, page-level and record-level locking. For each of the following operations, indicate the sequence of lock requests that must be generated by a transaction that wants to carry out (just) these operations:

1. Read record P1200 : 5.
2. Read records P1200 : 98 through P1205 : 2.
3. Read all (records on all) pages in file F1.
4. Read pages P500 through P520.
5. Read pages P10 through P980.
6. Read all pages in *P1* and (based on the values read) modify 10 pages.
7. Delete record P1200 : 98. (This is a blind write.)
8. Delete the first record from each page. (Again, these are blind writes.)
9. Delete all records.

**Exercise 17.12** Suppose that we have only two types of transactions, *T1* and *T2*. Transactions preserve database consistency when run individually. We have defined several *integrity constraints* such that the DBMS never executes any SQL statement that brings the database into an inconsistent state. Assume that the DBMS does not perform *any* concurrency control. Give an example schedule of two transactions *T1* and *T2* that satisfies all these conditions, yet produces a database instance that is not the result of any serial execution of *T1* and *T2*.

## BIBLIOGRAPHIC NOTES

Concurrent access to B trees is considered in several papers, including [70,456,472,505,678]. Concurrency control techniques for Linear Hashing are presented in [240] and [543]. Multiple-granularity locking is introduced in [3:36] and studied further in [127, 449].

A concurrency control method that works with the ARIES recovery method is presented in [545]. Another paper that considers concurrency control issues in the context of recovery is [492]. Algorithms for building indexes without stopping the DBMS are presented in [548] and [9]. The performance of B tree concurrency control algorithms is studied in [704]. Performance of various concurrency control algorithms is discussed in [16, 729, 735]. A good survey of concurrency control methods and their performance is [734]. [455] is a comprehensive collection of papers on this topic.

Two-phase-based multiversion concurrency control is studied in [620]. Multiversion concurrency control algorithms are studied formally in [87]. Lock-based multiversion techniques are considered in [460]. Optimistic concurrency control is introduced in [457]. The use of access invariance to improve conflict resolution in high-contention environments is discussed in [281] and [280]. Transaction management issues for real-time database systems are discussed in [1, 15, 368, 382, 386, 448]. There is a large body of theoretical results on database concurrency control; [582, 89] offer thorough textbook presentations of this material.



# 18

---

## CRASH RECOVERY

- ☛ What steps are taken in the ARIES method to recover from a DBMS crash?
- ☛ How is the log maintained during normal operation?
- ☛ How is the log used to recover from a crash?
- ☛ What information in addition to the log is used during recovery?
- ☛ What is a checkpoint and why is it used?
- ☛ What happens if repeated crashes occur during recovery?
- ☛ How is media failure handled?
- ☛ How does the recovery algorithm interact with concurrency control?
- **Key concepts:** steps in recovery, analysis, redo, undo; ARIES, repeating history; log, LSN, forcing pages, WAL; types of log records, update, commit, abort, end, compensation; transaction table, lastLSN; dirty page table, recLSN; checkpoint, fuzzy checkpointing, master log record; media recovery; interaction with concurrency control; shadow paging

Hurnpty Durnpty sat on a wall.  
Hurnpty Durnpty had a great fall.  
All the King's horses and all the King's men  
Could not put Hurnpty together again.

—Old nursery rhyme

The recovery manager of a DBMS is responsible for ensuring two important properties of transactions: Atomicity and durability. It ensures *atomicity* by undoing the actions of transactions that do not commit and *durability* by making sure that all actions of committed transactions survive system crashes (e.g., a core dump caused by a bus error) and media failures (e.g., a disk is corrupted).

The recovery manager is one of the hardest components of a DBMS to design and implement. It must deal with a wide variety of database states because it is called on during system failures. In this chapter, we present the ARIES recovery algorithm, which is conceptually simple, works well with a wide range of concurrency control mechanisms, and is being used in an increasing number of database systems.

We begin with an introduction to ARIES in Section 18.1. We discuss the log, which is a central data structure in recovery, in Section 18.2, and other recovery-related data structures in Section 18.3. We complete our coverage of recovery-related activity during normal processing by presenting the Write-Ahead Logging protocol in Section 18.4, and checkpointing in Section 18.5.

We discuss recovery from a crash in Section 18.6. Aborting (or rolling back) a single transaction is a special case of Undo, discussed in Section 18.6.3. We discuss media failures in Section 18.7, and conclude in Section 18.8 with a discussion of the interaction of concurrency control and recovery and other approaches to recovery. In this chapter, we consider recovery only in a centralized DBMS; recovery in a distributed DBMS is discussed in Chapter 22.

## 18.1 INTRODUCTION TO ARIES

ARIES is a recovery algorithm designed to work with a steal, no-force approach. When the recovery manager is invoked after a crash, restart proceeds in three phases:

1. **Analysis:** Identifies dirty pages in the buffer pool (i.e., changes that have not been written to disk) and active transactions at the time of the crash.
2. **Redo:** Repeats all actions, starting from an appropriate point in the log, and restores the database state to what it was at the time of the crash.
3. **Undo:** Undoes the actions of transactions that did not commit, so that the database reflects only the actions of committed transactions.

Consider the simple execution history illustrated in Figure 18.1. When the system is restarted, the Analysis phase identifies *T1* and *T3* as transactions

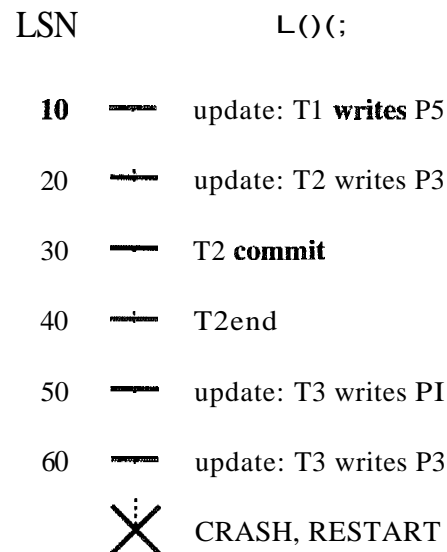


Figure 18.1 Execution History with a Crash

active at the time of the crash and therefore to be undone; *T2* as a committed transaction, and all its actions therefore to be written to disk; and *P1*, *P3*, and *P5* as potentially dirty pages. All the updates (including those of *T1* and *T3*) are reapplied in the order shown during the Redo phase. Finally, the actions of *T1* and *T3* are undone in reverse order during the Undo phase; that is, *T3*'s write of *P3* is undone, *T3*'s write of *P1* is undone, and then *T1*'s write of *P5* is undone.

Three main principles lie behind the ARIES recovery algorithm:

- **Write-Ahead Logging:** Any change to a database object is first recorded in the log; the record in the log must be written to stable storage before the change to the database object is written to disk.
- **Repeating History During Redo:** On restart following a crash, ARIES retraces all actions of the DBMS before the crash and brings the system back to the exact state that it was in at the time of the crash. Then, it undoes the actions of transactions still active at the time of the crash (effectively aborting them).
- **Logging Changes During Undo:** Changes made to the database while undoing a transaction are logged to ensure such an action is not repeated in the event of repeated (failures causing) restarts.

The second point distinguishes ARIES from other recovery algorithms and is the basis for much of its simplicity and flexibility. In particular, ARIES can support concurrency control protocols that involve locks of finer granularity than a page (e.g., record-level locks). The second and third points are also



**Crash Recovery:** IBM DB2, Informix, Microsoft SQL Server, Oracle 8, and Sybase I\SE all use a WAL scheme for recovery. IBM DB2 uses ARIES, and the others use schemes that are actually quite similar to ARIES (e.g., all changes are re-applied, not just the changes made by transactions that are 'winners') although there are several variations.

important in dealing with operations where redoing and undoing the operation are not exact inverses of each other. We discuss the interaction between concurrency control and crash recovery in Section 18.8, where we also discuss other approaches to recovery briefly.

## 18.2 THE LOG

The log, sometimes called the trail or journal, is a history of actions executed by the DBMS. Physically, the log is a file of records stored in stable storage, which is assumed to survive crashes; this durability can be achieved by maintaining two or more copies of the log on different disks (perhaps in different locations), so that the chance of all copies of the log being simultaneously lost is negligibly small.

The most recent portion of the log, called the log tail, is kept in main memory and is periodically forced to stable storage. This way, log records and data records are written to disk at the same granularity (pages or sets of pages).

Every log record is given a unique *id* called the log sequence number (LSN). As with any record id, we can fetch a log record with one disk access given the LSN. Further, LSNs should be assigned in monotonically increasing order; this property is required for the ARIES recovery algorithm. If the log is a sequential file, in principle growing indefinitely, the LSN can simply be the address of the first byte of the log record.<sup>1</sup>

For recovery purposes, every page in the database contains the LSN of the most recent log record that describes a change to this page. This LSN is called the pageLSN.

A log record is written for each of the following actions:

---

<sup>1</sup>In practice, various techniques are used to identify portions of the log that are 'too old' to be needed again to bound the amount of stable storage used for the log. Given such a bound, the log may be implemented as a 'circular' file, in which case the LSN may be the log record id plus a *wrap-count*.

- **Updating a Page:** After modifying the page, an *update* type record (described later in this section) is appended to the log tail. The pageLSN of the page is then set to the LSN of the update log record. (The page must be pinned in the buffer pool while these actions are carried out.)
- **Commit:** When a transaction decides to commit, it force-writes a *commit* type log record containing the transaction id. That is, the log record is appended to the log, and the log tail is written to stable storage, up to and including the commit record.<sup>2</sup> The transaction is considered to have committed at the instant that its commit log record is written to stable storage. (Some additional steps must be taken, e.g., relogging the transaction's entry in the transaction table; these follow the writing of the commit log record.)
- **Abort:** When a transaction is aborted, an *abort* type log record containing the transaction id is appended to the log, and Undo is initiated for this transaction (Section 18.6.3).
- **End:** As noted above, when a transaction is aborted or committed, some additional actions must be taken beyond writing the abort or commit log record. After all these additional steps are completed, an *end* type log record containing the transaction id is appended to the log.
- **Undoing an update:** When a transaction is rolled back (because the transaction is aborted, or during recovery from a crash), its updates are undone. When the action described by an update log record is undone, a *compensation log record*, or CLR, is written.

Every log record has certain fields: *prevLSN*, *transID*, and *type*. The set of all log records for a given transaction is maintained as a linked list going back in time, using the *prevLSN* field; this list must be updated whenever a log record is added. The *transID* field is the id of the transaction generating the log record, and the *type* field obviously indicates the type of the log record.

Additional fields depend on the type of the log record. We already mentioned the additional contents of the various log record types, with the exception of the update and compensation log record types, which we describe next.

## Update Log Records

The fields in an update log record are illustrated in Figure 18.2. The *pageID* field is the page id of the modified page; the *length* in bytes and the *offset* of the

<sup>2</sup>Note that this step requires the buffer manager to be able to selectively *force* pages to stable storage.

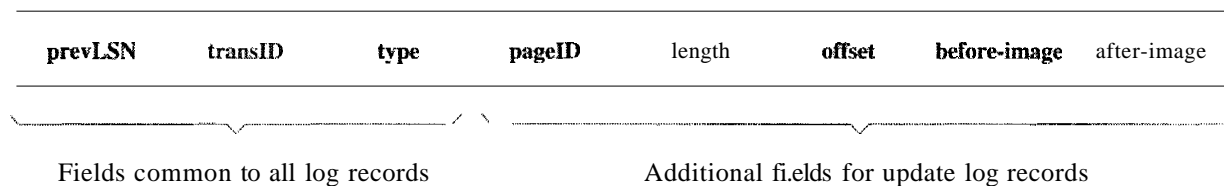


Figure 18.2 Contents of an Update Log Record

change are also included. The *before-image* is the value of the changed bytes before the change; the *after-image* is the value after the change. An update log record that contains both before- and after-images can be used to redo the change and undo it. In certain contexts, which we do not discuss further, we can recognize that the change will never be undone (or, perhaps, redone). A redo-only update log record contains just the after-image; similarly an undo-only update record contains just the before-image.

## Compensation Log Records

A compensation log record (CLR) is written just before the change recorded in an update log record  $U$  is undone. (Such an undo can happen during normal system execution when a transaction is aborted or during recovery from a crash.) A compensation log record  $C$  describes the action taken to undo the actions recorded in the corresponding update log record and is appended to the log tail just like any other log record. The compensation log record  $C$  also contains a field called *undoNextLSN*, which is the LSN of the next log record that is to be undone for the transaction that wrote update record  $U$ ; this field in  $C$  is set to the value of *prevLSN* in  $U$ .

As an example, consider the fourth update log record shown in Figure 18.3. If this update is undone, a CLR would be written, and the information in it would include the *transID*, *pageID*, *length*, *offset*, and *before-image* fields from the update record. Notice that the CLR records the (undo) action of changing the affected bytes back to the *before-image* value; thus, this value and the location of the affected bytes constitute the redo information for the action described by the CLR. The *undoNextLSN* field is set to the LSN of the first log record in Figure 18.3.

Unlike an update log record, a CLR describes an action that will never be *undone*, that is, we never undo an undo action. The reason is simple: An update log record describes a change made by a transaction during normal execution and the transaction may subsequently be aborted, whereas a CLR describes an action taken to rollback a transaction for which the decision to abort has already been made. Therefore, the transaction *must* be rolled back, and the

undo action described by the CLR is definitely required. This observation is very useful because it bounds the amount of space needed for the log during restart from a crash: The number of CLRs that can be written during Undo is no more than the number of update log records for active transactions at the time of the crash.

A CLR may be written to stable storage (following WAL, of course) but the undo action it describes may not yet have been written to disk when the system crashes again. In this case, the undo action described in the CLR is reapplied during the Redo phase, just like the action described in update log records.

For these reasons, a CLR contains the information needed to reapply, or redo, the change described but not to reverse it.

### 18.3 OTHER RECOVERY-RELATED STRUCTURES

In addition to the log, the following two tables contain important recovery-related information:

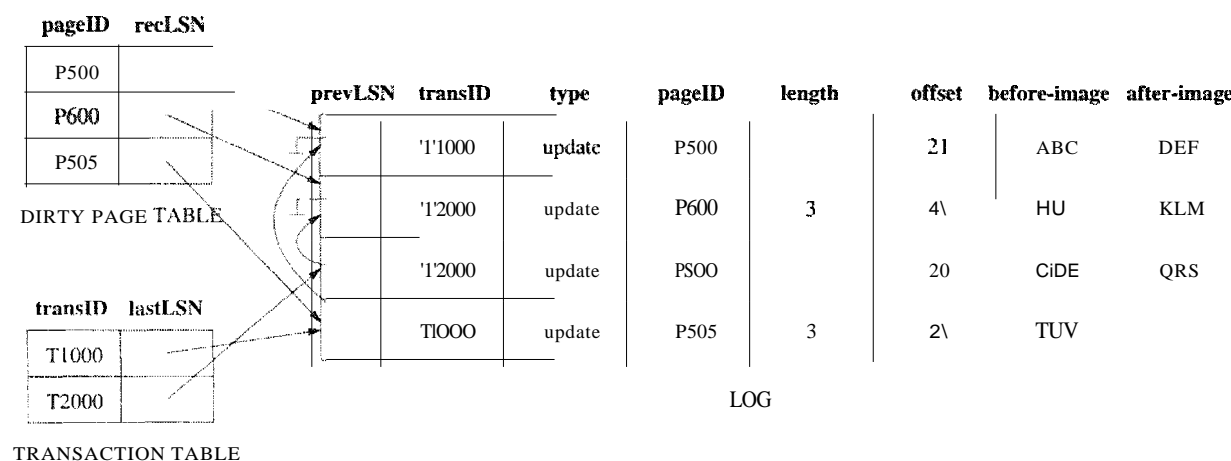
- **Transaction Table:** This table contains one entry for each active transaction. The entry contains (among other things) the transaction id, the status, and a field called **lastLSN**, which is the LSN of the most recent log record for this transaction. The status of a transaction can be that it is in progress, committed, or aborted. (In the latter two cases, the transaction will be removed from the table once certain 'clean up' steps are completed.)
- **Dirty page table:** This table contains one entry for each dirty page in the buffer pool, that is, each page with changes not yet reflected on disk. The entry contains a field **recLSN**, which is the LSN of the first log record that caused the page to become dirty. Note that this LSN identifies the earliest log record that might have to be redone for this page during restart from a crash.

During normal operation, these are maintained by the transaction manager and the buffer manager, respectively, and during restart after a crash, these tables are reconstructed in the Analysis phase of restart.

Consider the following simplistic example. Transaction T1000 changes the value of bytes 21 to 23 on page P500 from 'ABC' to 'DEF', transaction T2000 changes 'HIJ' to 'KLM' on page P600, transaction T2000 changes bytes 20 through 22 from 'GDE' to 'QRS' on page P500, then transaction T1000 changes 'TUV' to 'WXY' on page P505. The dirty page table, the transaction table,<sup>3</sup> and

---

<sup>3</sup>The status field is not shown in the figure for space reasons; all transactions are in progress.



**Figure 18.3** Instance of Log and Transaction Table

the log at this instant are shown in Figure 18.3. Observe that the log is shown growing from top to bottom; older records are at the top. Although the records for each transaction are linked using the prevLSN field, the log as a whole also has a sequential order that is important—for example, T2000's change to page P500 follows T1000's change to page P500, and in the event of a crash, these changes must be redone in the same order.

## 18.4 THE WRITE-AHEAD LOG PROTOCOL

Before writing a page to disk, every update log record that describes a change to this page must be forced to stable storage. This is accomplished by forcing all log records up to and including the one with LSN equal to the pageLSN to stable storage before writing the page to disk.

The importance of the WAL protocol cannot be overemphasized. WAL is the fundamental rule that ensures that a record of every change to the database is available while attempting to recover from a crash. If a transaction made a change and committed, the no-force approach means that some of these changes may not have been written to disk at the time of a subsequent crash. Without a record of these changes, there would be no way to ensure that the changes of a committed transaction survive crashes. Note that the definition of a *committed transaction* is effectively 'a transaction all of whose log records, including a commit record, have been written to stable storage'.

When a transaction is committed, the log tail is forced to stable storage, even if a no-force approach is being used. It is worth contrasting this operation with the actions taken under a force approach: If a force approach is used, all the pages modified by the transaction, rather than a portion of the log that includes all its records, must be forced to disk when the transaction commits. The set of

all changed pages is typically much larger than the log tail because the size of an update log record is close to (twice) the size of the changed bytes, which is likely to be much smaller than the page size. Further, the log is maintained as a sequential file, and all writes to the log are sequential writes. Consequently, the cost of forcing the log tail is much smaller than the cost of writing all changed pages to disk.

## 18.5 CHECKPOINTING;

A checkpoint is like a snapshot of the DBMS state, and by taking checkpoints periodically, as we will see, the DBMS can reduce the amount of work to be done during restart in the event of a subsequent crash.

Checkpointing in ARIES has three steps. First, a *begin\_checkpoint* record is written to indicate when the checkpoint starts. Second, an *end\_checkpoint* record is constructed, including in it the current contents of the transaction table and the dirty page table, and appended to the log. The third step is carried out after the *end\_checkpoint* record is written to stable storage: A special master record containing the LSN of the *begin\_checkpoint* log record is written to a known place on stable storage. While the *end\_checkpoint* record is being constructed, the DBMS continues executing transactions and writing other log records; the only guarantee we have is that the transaction table and dirty page table are accurate *as of the time of the begin\_checkpoint record*.

This kind of checkpoint, called a *fuzzy checkpoint*, is inexpensive because it does not require quiescing the System or writing out pages in the buffer pool (unlike some other forms of checkpointing). On the other hand, the effectiveness of this checkpointing technique is limited by the earliest recLSN of pages in the dirty pages table, because during restart we must redo changes starting from the log record whose LSN is equal to this recLSN. Having a background process that periodically writes dirty pages to disk helps to limit this problem.

When the System comes back up after a crash, the restart process begins by locating the most recent checkpoint record. For uniformity, the system always begins normal execution by taking a checkpoint, in which the transaction table and dirty page table are both empty.

## 18.6 RECOVERING FROM A SYSTEM CRASH

When the system is restarted after a crash, the recovery manager proceeds in three phases, as shown in Figure 18.4.

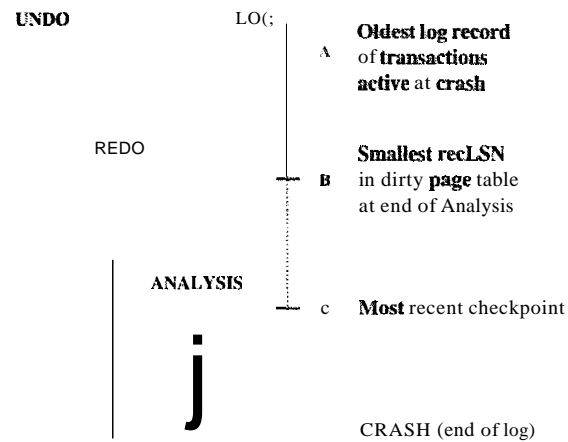


Figure 18.4 Three Phases of Restart in ARIES

The Analysis phase begins by examining the most recent begin\_checkpoint record, whose LSN is denoted *C* in Figure 18.4, and proceeds forward in the log until the last log record. The Redo phase follows Analysis and redoes all changes to any page that might have been dirty at the time of the crash; this set of pages and the starting point for Redo (the smallest recLSN of any dirty page) are determined during Analysis. The Undo phase follows Redo and undoes the changes of all transactions active at the time of the crash; again, this set of transactions is identified during the Analysis phase. Note that Redo reapplies changes in the order in which they were originally carried out; Undo reverses changes in the opposite order, reversing the most recent change first.

Observe that the relative order of the three points *A*, *B*, and *C* in the log may differ from that shown in Figure 18.4. The three phases of restart are described in more detail in the following sections.

### 18.6.1 Analysis Phase

The Analysis phase performs three tasks:

1. It determines the point in the log at which to start the Redo pass.
2. It determines (a conservative superset of the) pages in the buffer pool that were dirty at the time of the crash.
3. It identifies transactions that were active at the time of the crash and must be undone.

Analysis begins by examining the most recent begin\_checkpoint log record and initializing the dirty page table and transaction table to the copies of those structures in the next end\_checkpoint record. Thus, these tables are initialized to the set of dirty pages and active transactions at the time of the checkpoint.

(If additional log records are between the *begin\_checkpoint* and *end\_checkpoint* records, the tables must be adjusted to reflect the information in these records, but we omit the details of this step. See Exercise 18.9.) Analysis then scans the log in the forward direction until it reaches the end of the log:

- If an end log record for a transaction  $T$  is encountered,  $T$  is removed from the transaction table because it is no longer active.
- If a log record other than an end record for a transaction  $T$  is encountered, an entry for  $T$  is added to the transaction table if it is not already there. Further, the entry for  $T$  is modified:
  1. The *lastLSN* field is set to the LSN of this log record.
  2. If the log record is a commit record, the status is set to C, otherwise it is set to U (indicating that it is to be undone).
- If a redoable log record affecting page  $P$  is encountered, and  $P$  is not in the dirty page table, an entry is inserted into this table with page id  $P$  and *recLSN* equal to the LSN of this redoable log record. This LSN identifies the oldest change affecting page  $P$  that may not have been written to disk.

At the end of the Analysis phase, the transaction table contains an accurate list of all transactions that were active at the time of the crash—this is the set of transactions with status U. The dirty page table includes all pages that were dirty at the time of the crash but may also contain some pages that were written to disk. If an *end\_write* log record were written at the completion of each write operation, the dirty page table constructed during Analysis could be made more accurate, but in AHJES, the additional cost of writing *end\_write* log records is not considered to be worth the gain.

As an example, consider the execution illustrated in Figure 18.3. Let us extend this execution by assuming that  $T_{2000}$  commits, then  $T_{1000}$  modifies another page, say,  $P_{700}$ , and appends an update record to the log tail, and then the system crashes (before this update log record is written to stable storage).

The dirty page table and the transaction table, held in memory, are lost in the crash. The most recent checkpoint was taken at the beginning of the execution, with an empty transaction table and dirty page table; it is not shown in Figure 18.3. After examining this log record, which we assume is just before the first log record shown in the figure, Analysis initializes the two tables to be empty. Scanning forward in the log,  $T_{1000}$  is added to the transaction table; in addition,  $P_{500}$  is added to the dirty page table with *recLSN* equal to the LSN of the first shown log record. Similarly,  $T_{2000}$  is added to the transaction table and  $P_{600}$  is added to the dirty page table. There is no change based on the third log record, and the fourth record results in the addition of  $P_{505}$  to



the dirty page table. The commit record for T2000 (not in the figure) is not encountered, and T2000 is removed from the transaction table.

The Analysis phase is now complete, and it is recognized that the only active transaction at the time of the crash is T1000, with lastLSN equal to the LSN of the fourth record in Figure 18.3. The dirty page table reconstructed in the Analysis phase is identical to that shown in the figure. The update log record for the change to P700 is lost in the crash and not seen during the Analysis pass. Thanks to the WAL protocol, however, all is well—the corresponding change to page P700 cannot have been written to disk either!

Some of the updates may have been written to disk; for concreteness, let us assume that the change to P600 (and only this update) was written to disk before the crash. Therefore P600 is not dirty, yet it is included in the dirty page table. The pageLSN on page P600, however, reflects the write because it is now equal to the LSN of the second update log record shown in Figure 18.3.

## 18.6.2 Redo Phase

During the Redo phase, ARIES reapplies the updates of *all* transactions, committed or otherwise. Further, if a transaction was aborted before the crash and its updates were undone, as indicated by CLRs, the actions described in the CLRs are also reapplied. This repeating history paradigm distinguishes ARIES from other proposed WAL-based recovery algorithms and causes the database to be brought to the same state it was in at the time of the crash.

The Redo phase begins with the log record that has the smallest recLSN of all pages in the dirty page table constructed by the Analysis pass because this log record identifies the oldest update that may not have been written to disk prior to the crash. Starting from this log record, Redo scans forward until the end of the log. For each redoable log record (update or CLR) encountered, Redo checks whether the logged action must be redone. The action must be redone unless one of the following conditions holds:

- The affected page is not in the dirty page table.
- The affected page is in the dirty page table, but the recLSN for the entry is *greater than* the LSN of the log record being checked.
- The pageLSN (stored on the page, which must be retrieved to check this condition) is *greater than or equal* to the LSN of the log record being checked.

The first condition obviously implies that all changes to this page have been written to disk. Because the recLSN is the first update to this page that may

not have been written to disk, the second condition means that the update being checked was indeed propagated to disk. The third condition, which is checked last because it requires us to retrieve the page, also ensures that the update being checked was written to disk, because either this update or a later update to the page was written. (Recall our assumption that a write to a page is atomic; this assumption is important here!)

If the logged action must be redone:

1. The logged action is reapplied.
2. The pageLSN on the page is set to the LSN of the redone log record. No additional log record is written at this time.

Let us continue with the example discussed in Section 18.6.1. From the dirty page table, the smallest recLSN is seen to be the LSN of the first log record shown in Figure 18.3. Clearly, the changes recorded by earlier log records (there happen to be none in this example) have been written to disk. Now, Redo fetches the affected page, *P500*, and compares the LSN of this log record with the pageLSN on the page and, because we assumed that this page was not written to disk before the crash, finds that the pageLSN is less. The update is therefore reapplied; bytes 21 through 23 are changed to 'DEF', and the pageLSN is set to the LSN of this update log record.

Redo then examines the second log record. Again, the affected page, *P600*, is fetched and the pageLSN is compared to the LSN of the update log record. In this case, because we assumed that *P600* was written to disk before the crash, they are equal, and the update does not have to be redone.

The remaining log records are processed similarly, bringing the system back to the exact state it was in at the time of the crash. Note that the first two conditions indicating that a redo is unnecessary never hold in this example. Intuitively, they come into play when the dirty page table contains a very old recLSN, going back to before the most recent checkpoint. In this case, as Redo scans forward from the log record with this LSN, it encounters log records for pages that were written to disk prior to the checkpoint and therefore not in the dirty page table in the checkpoint. Some of these pages may be dirtied again after the checkpoint; nonetheless, the updates to these pages prior to the checkpoint need not be redone. Although the third condition alone is sufficient to recognize that these updates need not be redone, it requires us to fetch the affected page. The first two conditions allow us to recognize this situation without fetching the page. (The reader is encouraged to construct examples that illustrate the use of each of these conditions; see Exercise 18.8.)

At the end of the Redo phase, end type records are written for all transactions with status C, which are moved into the transaction table.

### 18.6.3 Undo Phase

The Undo phase, unlike the other two phases, scans backward from the end of the log. The goal of this phase is to undo the actions of all transactions active at the time of the crash, that is, to effectively abort them. This set of transactions is identified in the transaction table constructed by the Analysis phase.

#### The Undo Algorithm

Undo begins with the transaction table constructed by the Analysis phase, which identifies all transactions active at the time of the crash, and includes the LSN of the most recent log record (the lastLSN field) for each such transaction. Such transactions are called loser transactions. All actions of losers must be undone, and further, these actions must be undone in the reverse of the order in which they appear in the log.

Consider the set of lastLSN values for all loser transactions. Let us call this set ToUndo. Undo repeatedly chooses the largest (i.e., most recent) LSN value in this set and processes it, until ToUndo is empty. To process a log record:

1. If it is a CLR and the undoNextLSN value is not *null*, the undoNextLSN value is added to the set ToUndo; if the undoNextLSN is *null*, an end record is written for the transaction because it is completely undone, and the CLR is discarded.
2. If it is an update record, a CLR is written and the corresponding action is undone, as described in Section 18.2, and the prevLSN value in the update log record is added to the set ToUndo.

When the set ToUndo is empty, the Undo phase is complete. If restart is not complete, and the system can proceed with normal operations.

Let us continue with the scenario discussed in Sections 18.6.1 and 18.6.2. The only active transaction at the time of the crash was determined to be T1000. From the transaction table, we get the LSN of its most recent log record, which is the fourth update log record in Figure 18.3. The update is undone, and a CLR is written with undoNextLSN equal to the LSN of the first log record in the figure. The next record to be undone for transaction T1000 is the first log record in the figure. After this is undone, a CLR and an end log record for T1000 are written, and the Undo phase is complete.

In this example, undoing the action recorded in the first log record causes the action of the third log record, which is due to a committed transaction, to be overwritten and thereby lost! This situation arises because  $T_{2000}$  overwrote a data item written by  $T_{1000}$  while  $T_{1000}$  was still active; if Strict 2PL were followed,  $T_{2000}$  would not have been allowed to overwrite this data item.

## Aborting a Transaction

Aborting a transaction is just a special case of the Undo phase of Restart in which a single transaction, rather than a set of transactions, is undone. The example in Figure 18.5, discussed next, illustrates this point.

## Crashes during Restart

It is important to understand how the Undo algorithm presented in Section 18.6.3 handles repeated system crashes. Because the details of precisely how the action described in an update log record is undone are straightforward, we discuss Undo in the presence of system crashes using an execution history, shown in Figure 18.5, that abstracts away unnecessary detail. This example illustrates how aborting a transaction is a special case of Undo and how the use of CLRs ensures that the Undo action for an update log record is not applied twice.

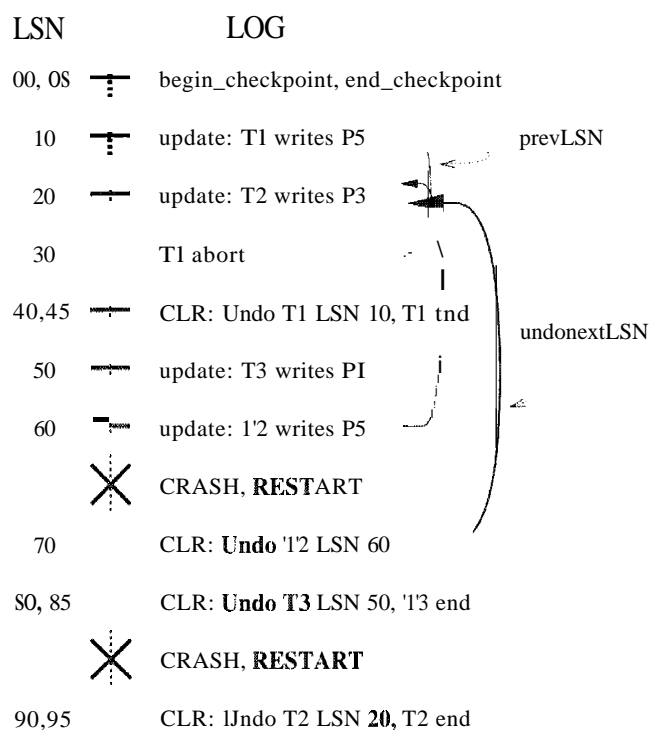


Figure 18.5 Example of Undo with Repeated Crashes

The log shows the order in which the DBMS executed various actions; note that the LSNs are in ascending order, and that each log record for a transaction has a `prevLSN` field that points to the previous log record for that transaction. We have not shown *null* `prevLSNs`, that is, some special value used in the `prevLSN` field of the first log record for a transaction to indicate that there is no previous log record. We also compacted the figure by occasionally displaying two log records (separated by a column) on a single line.

Log record (with LSN) 30 indicates that *T1* aborts. All actions of this transaction should be undone in reverse order, and the only action of *T1*, described by the update log record 10, is indeed undone as indicated by CLR, 40.

After the first crash, Analysis identifies *P1* (with `recLSN` 50), *P3* (with `recLSN` 20), and *P5* (with `recLSN` 10) as dirty pages. Log record 45 shows that *T1* is a completed transaction; hence, the transaction table identifies *T2* (with `lastLSN` 60) and *T3* (with `lastLSN` 50) as active at the time of the crash. The Redo phase begins with log record 10, which is the minimum `recLSN` in the dirty page table, and reapplies all actions (for the update and CLR records), as per the Redo algorithm presented in Section 18.6.2.

The redo set consists of LSNs 60, for *T2*, and 50, for *T3*. The Undo phase now begins by processing the log record with LSN 60 because 60 is the largest LSN in the ToUndo set. The update is undone, and a CLR, (with LSN 70) is written to the log. This CLR has `undoNextLSN` equal to 20, which is the `prevLSN` value in log record 60; 20 is the next action to be undone for *T2*. Now the largest remaining LSN in the ToUndo set is 50. The write corresponding to log record 50 is now undone, and a CLR, describing the change is written. This CLR has LSN 80, and its `undoNextLSN` field is *null* because 50 is the only log record for transaction *T3*. Therefore *T3* is completely undone, and an end record is written. Log records 70, 80, and 85 are written to stable storage before the system crashes a second time; however, the changes described by these records may not have been written, to disk.

When the system is restarted after the second crash, Analysis determines that the only active transaction at the time of the crash was *T2*; in addition, the dirty page table is identical to what it was during the previous restart. Log records 10 through 85 are processed again during Redo. (If some of the changes made during the previous Redo were written to disk, the pageLSN's on the affected pages are used to detect this situation and avoid writing these pages again.) The Undo phase considers the only LSN in the ToUndo set, 70, and processes it by adding the `undoNextLSN` value (20) to the ToUndo set. Next, log record 20 is processed by undoing *T2*'s write of page *P3*, and a CLR is written (LSN 90). Because 20 is the first of *T2*'s log records and therefore, the last of its records

to be undone—the undoNextLSN field in this CLR is *null*, an end record is written for  $T_2$ , and the ToUndo set is now empty.

Recovery is now complete, and normal execution can resume with the writing of a checkpoint record.

This example illustrated repeated crashes during the Undo phase. For completeness, let us consider what happens if the system crashes while Restart is in the Analysis or Redo phase. If a crash occurs during the Analysis phase, all the work done in this phase is lost, and on restart the Analysis phase starts afresh with the same information as before. If a crash occurs during the Redo phase, the only effect that survives the crash is that some of the changes made during Redo may have been written to disk prior to the crash. Restart starts again with the Analysis phase and then the Redo phase, and some update log records that were redone the first time around will not be redone a second time because the pageLSN is now equal to the update record's LSN (although the pages have to be fetched again to detect this).

We can take checkpoints during Restart to minimize repeated work in the event of a crash, but we do not discuss this point.

## 18.7 MEDIA RECOVERY

Media recovery is based on periodically making a copy of the database. Because copying a large database object such as a file can take a long time, and the DBMS must be allowed to continue with its operations in the meantime, creating a copy is handled in a manner similar to taking a fuzzy checkpoint.

When a database object such as a file or a page is corrupted, the copy of that object is brought up-to-date by using the log to identify and reapply the changes of committed transactions and undo the changes of uncommitted transactions (as of the time of the media recovery operation).

The begin\_checkpoint LSN of the most recent complete checkpoint is recorded along with the copy of the database object to minimize the work in reapplying changes of committed transactions. Let us compare the smallest recLSN of a dirty page in the corresponding end\_checkpoint record with the LSN of the begin\_checkpoint record and call the smaller of these two LSNs  $I$ . We observe that the actions recorded in all log records with LSNs less than  $I$  must be reflected in the copy. Thus, only log records with LSNs greater than  $I$  need be reapplied to the copy.

Finally, the updates of transactions that are incomplete at the time of Inedia recovery or that were aborted after the fuzzy copy was completed need to be undone to ensure that the page reflects only the actions of committed transactions. The set of such transactions can be identified as in the Analysis pass, and we omit the details.

## 18.8 OTHER APPROACHES AND INTERACTION WITH CONCURRENCY CONTROL

Like ARIES, the most popular alternative recovery algorithms also maintain a log of database actions according to the WAL protocol. A major distinction between ARIES and these variants is that the Redo phase in ARIES *repeats history*, that is, redoes the actions of *all* transactions, not just the non-losers. Other algorithms redo only the non-losers, and the Redo phase follows the Undo phase, in which the actions of losers are rolled back.

Thanks to the repeating history paradigm and the use of CLRs, ARIES supports fine-granularity locks (record-level locks) and logging of logical operations rather than just byte-level modifications. For example, consider a transaction  $T$  that inserts a data entry 15\* into a B+ tree index. Between the time this insert is done and the time that  $T$  is eventually aborted, other transactions may also insert and delete entries from the tree. If record-level locks are set rather than page-level locks, the entry 15\* may be on a different physical page when  $T$  aborts from the one that  $T$  inserted it into. In this case, the undo operation for the insert of 15\* must be recorded in logical terms because the physical (byte-level) actions involved in undoing this operation are not the inverse of the physical actions involved in inserting the entry.

Logging logical operations yields considerably higher concurrency, although the use of fine-granularity locks can lead to increased locking activity (because more locks must be set). Hence, there is a trade-off between different WAL-based recovery schemes. We chose to cover ARIES because it has several attractive properties, in particular, its simplicity and its ability to support fine-granularity locks and logging of logical operations.

One of the earliest recovery algorithms, used in the System R prototype at IBM<sup>1</sup>, takes a very different approach. There is no logging and, of course, no WAL protocol. Instead, the database is treated as a collection of pages and accessed through a page table, which maps page ids to disk addresses. When a transaction makes changes to a data page, it actually makes a copy of the page, called the shadow of the page, and changes the shadow page. The transaction copies the appropriate part of the page table and changes the entry for the changed page to point to the shadow, so that it can see the

changes; however, other transactions continue to see the original page table, and therefore the original page, until this transaction commits. Aborting a transaction is simple: Just discard its shadow versions of the page table and the data pages. Committing a transaction involves making its version of the page table public and discarding the original data pages that are superseded by shadow pages.

This scheme suffers from a number of problems. First, data becomes highly fragmented due to the replacement of pages by shadow versions, which may be located far from the original page. This phenomenon reduces data clustering and makes good garbage collection imperative. Second, the scheme does not yield a sufficiently high degree of concurrency. Third, there is a substantial storage overhead due to the use of shadow pages. Fourth, the process aborting a transaction can itself run into deadlocks, and this situation must be specially handled because the semantics of aborting an abort transaction gets murky.

For these reasons, even in System R, shadow paging was eventually superseded by WAL-based recovery techniques.

## 18.9 REVIEW QUESTIONS

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

- What are the advantages of the ARIES recovery algorithm? (Section 18.1)
- Describe the three steps in crash recovery in ARIES? What is the goal of the Analysis phase? The redo phase? The undo phase? (Section 18.1)
- What is the LSN of a log record? (Section 18.2)
- What are the different types of log records and when are they written? (Section 18.2)
- What information is maintained in the transaction table and the dirty page table? (Section 18.3)
- What is Write-Ahead Logging? What is forced to disk at the time a transaction commits? (Section 18.4)
- What is a fuzzy checkpoint? Why is it useful? What is a master log record? (Section 18.5)
- In which direction does the Analysis phase of recovery scan the log? At which point in the log does it begin and end the scan? (Section 18.6.1)
- Describe what information is gathered in the Analysis phase and how. (Section 18.6.1)



- In which direction does the Redo phase of recovery process the log? At which point in the log does it begin and end? (Section 18.6.2)
- What is a redoable log record? Under what conditions is the logged action redone? What steps are carried out when a logged action is redone? (Section 18.6.2)
- In which direction does the Undo phase of recovery process the log? At which point in the log does it begin and end? (Section 18.6.3)
- What are loser transactions? How are they processed in the Undo phase and in what order? (Section 18.6.3)
- Explain what happens if there are crashes during the Undo phase of recovery. What is the role of CLRs? What if there are crashes during the Analysis and Redo phases? (Section 18.6.3)
- How does a DBMS recover from media failure without reading the complete log? (Section 18.7)
- Record-level logging increases concurrency. What are the potential problems, and how does ARIES address them? (Section 18.8)
- What is shadow paging? (Section 18.8)

## EXERCISES

Exercise 18.1 Briefly answer the following questions:

1. How does the recovery manager ensure atomicity of transactions? How does it ensure durability?
2. What is the difference between stable storage and disk?
3. What is the difference between a system crash and a media failure?
4. Explain the WAL protocol.
5. Describe the steal and no-force policies.

Exercise 18.2 Briefly answer the following questions:

1. What are the properties required of LSNs?
2. What are the fields in an update log record? Explain the use of each field.
3. What are redoable log records?
4. What are the differences between update log records and CLRs?

Exercise 18.3 Briefly answer the following questions:

1. What are the roles of the Analysis, Redo, and Undo phases in ARIES?
2. Consider the execution shown in Figure 18.6.

LSN		LOG
00		begin_checkpoint
10	■	end_checkpoint
20	■	update: T1 writes P5
30	⊢	update: T2 writes P3
40	⊢	T2 commit
50		T2end
60	⊢	update: T3 writes P3
70		T1 abort
	⊗	CRASH, RESTART

Figure 18.6 Execution with a Crash

LSN		LOG
00	⊢	update: T1 writes P2
10	⊢	update: T1 writes P1
20	■	update: T2 writes P5
30		update: T3 writes P3
40		T3 commit
50	⊢	update: T2 writes P5
60	⊢	update: T2 writes P3
70	⊢	T2 abort

Figure 18.7 Aborting a Transaction

- What is done during Analysis? (Be precise about the points at which Analysis begins and ends and describe the contents of any tables constructed in this phase.)
- What is done during Redo? (Be precise about the points at which Redo begins and ends.)
- What is done during Undo? (Be precise about the points at which Undo begins and ends.)

Exercise 18.4 Consider the execution shown in Figure 18.7.

- Extend the figure to show prevLSN and lndonextLSN values.
- Describe the actions taken to rollback transaction T2.

LSN		LOG
00	—	begin_checkpoint
10		end_checkpoint
20	⊕	update: l'1 writes P1
30	⊕	update: l'2 writes P2
40	⊕	update: l'3 writes P3
50	—	l'2 commit
60	⊕	update: l'3 writes P2
70	⊕	l'2 end
80	—	update: l'1 writes P5
90	⊕	l'3 abort
	⊗	CRASH,RESTART

Figure 18.8 Execution with Multiple Crashes

3. Show the log after  $T2$  is rolled back, including all prevLSN and undonextLSN values in log records.

**Exercise 18.5** Consider the execution shown in Figure 18.8. In addition, the system crashes during recovery after writing two log records to stable storage and again after writing another two log records.

1. What is the value of the LSN stored in the master log record?
2. What is done during Analysis?
3. What is done during Redo?
4. What is done during Undo?
5. Show the log when recovery is complete, including all non-null prevLSN and undonextLSN values in log records.

**Exercise 18.6** Briefly answer the following questions:

1. How is checkpointing done in ARIES?
2. Checkpointing can also be done as follows: Quiesce the system so that only checkpointing activity can be in progress, write out copies of all dirty pages, and include the dirty page table and transaction table in the checkpoint record. What are the pros and cons of this approach versus the checkpointing approach of ARIES?
3. What happens if a second begin\_checkpoint record is encountered during the Analysis phase?
4. Can a second end\_checkpoint record be encountered during the Analysis phase?
5. Why is the use of CLRs important for the use of undo actions that are not the physical inverse of the original update?

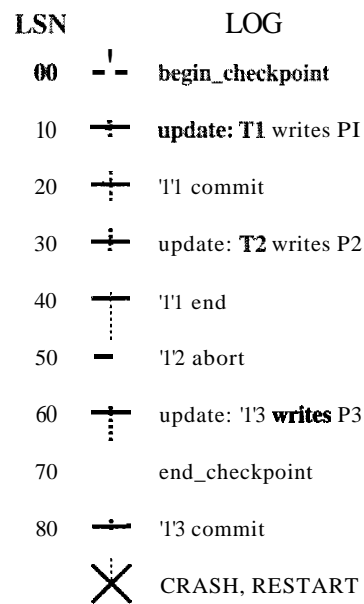


Figure 18.9 Log Records between Checkpoint Records

6. Give an example that illustrates how the paradigm of repeating history and the use of CLRs allow ARIES to support locks of finer granularity than a page.

Exercise 18.7 Briefly answer the following questions:

1. If the system fails repeatedly during recovery, what is the maximum number of log records that can be written (as a function of the number of update and other log records written before the crash) before restart completes successfully?
2. What is the oldest log record we need to retain?
3. If a bounded amount of stable storage is used for the log, how can we always ensure enough stable storage to hold all log records written during restart?

Exercise 18.8 Consider the three conditions under which a redo is unnecessary (Section 20.2.2).

1. Why is it cheaper to test the first two conditions?
2. Describe an execution that illustrates the use of the first condition.
3. Describe an execution that illustrates the use of the second condition.

Exercise 18.9 The description in Section 18.6.1 of the Analysis phase made the simplifying assumption that no log records appeared between the begin\_checkpoint and end\_checkpoint records for the most recent complete checkpoint. The following questions explore how such records should be handled.

1. Explain why log records could be written between the begin\_checkpoint and end\_checkpoint records.
2. Describe how the Analysis phase could be modified to handle such records.
3. Consider the execution shown in Figure 18.9. Show the contents of the end\_checkpoint record.
4. Illustrate your modified Analysis phase on the execution shown in Figure 18.9.

Exercise 18.10 Answer the following questions briefly:

1. Explain how Inedin recovery is handled in ARIES.
2. What are the pros and cons of using fuzzy durnps for media recovery?
3. What are the sirYlilarities and differences between checkpoints and fuzzy chunps?
4. Contrast ARIES with other WAL-based recovery schemes.
5. Contrast AHIES with shadow-page-based recovery.

## BIBLIOGRAPHIC NOTES

Our discussion of the ARIES recovery algorithm is based on [544]. [282] is a survey article that contains a very readable, short description of ARIES. [541, 545] also discuss ARIES. Fine-granularity locking increases concurrency but at the cost of lll0l'e locking activity; [542] suggests a technique based on LSNs for alleviating this problrYl. [458] presents a forlllal verification of ARIES.

[355] is an excellent survey that provides a broader treatrnt of recovery alorlthulls than our coverage, in which we chose to concentrate on one particular alorlthrn. [17] considers perfor-nance of concurrency control and recovery alorlthrlls, taking into account their interactions. The irnpact of recovery on concurrency control is also discussed in [769]. [625] contains a performance analysis of various recovery techniques. [236] cornpares recovery techniques for main rnerllory database systeulS, which are optirnized for the case that lll0st of the active data set fits in rnain Hlmemory.

[478] presents a description of a recovery algorithm based on write-ahead logging in which 'loser' transactions are first undone and then (only) transactions that corllnlltted before the crash are redone. Shadow paging is described in [493, 337]. A scheme that uses a cOlmbination of shadow paging and in-place updating is described in [624].

PART VI

---

# DATABASE DESIGN AND TUNING





# 19

---

## SCHEMA REFINEMENT AND NORMAL FORMS

- ☛ What problems are caused by redundantly storing information?
- ☛ What are functional dependencies?
- ☛ What are normal forms and what is their purpose?
- ☛ What are the benefits of BCNF and 3NF?
- ☛ What are the considerations in decomposing relations into appropriate normal forms?
- ☛ Where does normalization fit in the process of database design?
- ☛ Are there general dependencies useful in database design?
- **Key concepts:** redundancy, insert, delete, and update anomalies; functional dependency, Armstrong's Axioms; dependency closure, attribute closure; normal forms, BCNF, 3NF; decompositions, lossless-join, dependency-preservation; multivalued dependencies, join dependencies, inclusion dependencies, 4NF, 5NF

It is a melancholy truth that even great men have their poor relations.

Charles Dickens

Conceptual database design gives us a set of relation schemas and integrity constraints (ICs) that can be regarded as a good starting point for the final database design. This initial design must be refined by taking the ICs into account more fully than is possible with just the ER model constructs and also by considering performance criteria and typical workloads. In this chapter, we discuss how ICs can be used to refine the conceptual schema produced by



translating an ER model design into a collection of relations. Workload and performance considerations are discussed in Chapter 20.

We concentrate on an important class of constraints called *functional dependencies*. Other kinds of les, for example, *multivalued dependencies* and *join dependencies*, also provide useful information. They can sometimes reveal redundancies that cannot be detected using functional dependencies alone. We discuss these other constraints briefly.

This chapter is organized as follows. Section 19.1 is an overview of the schema refinement approach discussed in this chapter. We introduce functional dependencies in Section 19.2. In Section 19.3, we show how to reason with functional dependency information to infer additional dependencies from a given set of dependencies. We introduce normal forms for relations in Section 19.4; the normal form satisfied by a relation is a measure of the redundancy in the relation. A relation with redundancy can be refined by *decomposing it*, or replacing it with smaller relations that contain the same information but without redundancy. We discuss decompositions and desirable properties of decompositions in Section 19.5, and we show how relations can be decomposed into smaller relations in desirable normal forms in Section 19.6.

In Section 19.7, we present several examples that illustrate how relational schemas obtained by translating an ER model design can nonetheless suffer from redundancy, and we discuss how to refine such schemas to eliminate the problems. In Section 19.8, we describe other kinds of dependencies for database design. We conclude with a discussion of normalization for our case study, the Internet shop, in Section 19.9.

## 19.1 INTRODUCTION TO SCHEMA REFINEMENT

We now present an overview of the problems that schema refinement is intended to address and a refinement approach based on decompositions. Redundant storage of information is the root cause of these problems. Although decomposition can eliminate redundancy, it can lead to problems of its own and should be used with caution.

### 19.1.1 Problems Caused by Redundancy

Storing the same information redundantly, that is, in more than one place within a database, can lead to several problems:

- **Redundant Storage:** Some information is stored repeatedly.

- ii **Update Anomalies:** If one copy of such repeated data is updated, an inconsistency is created unless all copies are similarly updated.
- ii **Insertion Anomalies:** It may not be possible to store certain information unless some other, unrelated, information is stored as well.
- ii **Deletion Anomalies:** It may not be possible to delete certain information without losing some other, unrelated, information as well.

Consider a relation obtained by translating a variant of the Hourly\_Emps entity set from Chapter 2:

Hourly\_Emps(*ssn*, *name*, *lot*, *rating*, *hourly\_wages*, *hours\_worked*)

In this chapter, we omit attribute type information for brevity, since our focus is on the grouping of attributes into relations. We often abbreviate an attribute name to a single letter and refer to a relation schema by a string of letters, one per attribute. For example, we refer to the Hourly\_Emps schema as *SNLRWH* (*W* denotes the *hourly\_wages* attribute).

The key for Hourly\_Emps is *ssn*. In addition, suppose that the *hourly\_wages* attribute is determined by the *rating* attribute. That is, for a given *rating* value, there is only one permissible *hourly\_wages* value. This IC is an example of a *functional dependency*. It leads to possible redundancy in the relation Hourly\_Emps, as illustrated in Figure 19.1.

	<i>name</i>	<i>lot</i>	<i>rating</i>	<i>hourly_wages</i>	<i>hours_worked</i>
123-22-3666	Attishoo	48	8	10	40
231-31-5368	Sruiley	22	8	10	30
131-24-3650	Srilethurst	35	5	7	30
434-26-3751	Guldu	35	5	7	32
612-67-4134	Madayan	35	8	10	40

Figure 19.1 An Instance of the Hourly\_Emps Relation

If the same value appears in the *rating* column of two tuples, the IC tells us that the same value must appear in the *hourly\_wages* column as well. This redundancy has the same negative consequences as before:

- ii **Redundant Storage:** the rating value 8 corresponds to the hourly wage 10, and this association is repeated three times.
- **Update Anomalies:** The *hourly\_wages* in the first tuple could be updated without making a similar change in the second tuple.

- *Insertion Anomalies:* We cannot insert a tuple for an employee unless we know the hourly wage for the employee's rating value.
- *Deletion Anomalies:* If we delete all tuples with a given rating value (e.g., we delete the tuples for Smithurst and Guldu) we lose the association between that *rating* value and its *hourly\_wage* value.

Ideally, we want schemas that do not permit redundancy, but at the very least we want to be able to identify schemas that do allow redundancy. Even if we choose to accept a schema with some of these drawbacks, perhaps owing to performance considerations, we want to make an informed decision.

## Null Values

It is worth considering whether the use of *null* values can address some of these problems. As we will see in the context of our example, they cannot provide a complete solution, but they can provide some help. In this chapter, we do not discuss the use of *null* values beyond this one example.

Consider the example *Hourly\_Emps* relation. Clearly, *null* values cannot help eliminate redundant storage or update anomalies. It appears that they can address insertion and deletion anomalies. For instance, to deal with the insertion anomaly example, we can insert an *employee* tuple with *null* values in the hourly wage field. However, *null* values cannot address all insertion anomalies. For example, we cannot record the hourly wage for a rating unless there is an employee with that rating, because we cannot store a null value in the *ssn* field, which is a primary key field. Similarly, to deal with the deletion anomaly example, we might consider storing a tuple with *null* values in all fields except *rating* and *hourly\_wages* if the last tuple with a given *rating* would otherwise be deleted. However, this solution does not work because it requires the *ssn* value to be *null*, and primary key fields cannot be *null*. Thus, *null* values do not provide a general solution to the problems of redundancy, even though they can help in some cases.

### 19.1.2 Decompositions

Intuitively, redundancy arises when a relational schema forces an association between attributes that is not natural. Functional dependencies (and, for that matter, other keys) can be used to identify such situations and suggest refinements to the schema. The essential idea is that many problems arising from redundancy can be addressed by replacing a relation with a collection of 'smaller' relations.

A decomposition of a relation schema  $R$  consists of replacing the relation schema by two (or more) relation schemas that each contain a subset of the attributes of  $R$  and together include all attributes in  $R$ . Intuitively, we want to store the information in any given instance of  $R$  by storing projections of the instance. This section examines the use of decompositions through several examples.

We can decompose Hourly\_Emps into two relations:

Hourly\_Emps2(ssn, name, lot, rating, hours\_worked)  
 Wages(rating, hourly\_wages)

The instances of these relations corresponding to the instance of Hourly\_Emps relation in Figure 19.1 is shown in Figure 19.2.

<u>ssn</u>	name	lot	rating	hours_worked
123-22-3666	Attishoo	48	8	40
231-31-5368	Sluiley	22	8	30
131-24-3650	Smethurst	35	5	30
434-26-3751	Guldu	35	5	32
612-67-4134	Madayan	35	8	40

rating	hourly_wages
8	10
5	7

Figure 19.2 Instances of Hourly\_Emps2 and Wages

Note that we can easily record the hourly wage for any rating simply by adding a tuple to Wages, even if no employee with that rating appears in the current instance of Hourly\_Emps. Changing the wage associated with a rating involves updating a single Wages tuple. This is more efficient than updating several tuples (as in the original design), and it eliminates the potential for inconsistency.

### 19.1.3 Problems Related to Decomposition

Unless we are careful, decomposing a relation schema can create more problems than it solves. Two important questions must be asked repeatedly:

1. Do we need to decompose a relation?

## 2. What problems (if any) does a given decomposition cause?

To help with the first question, several *normal forms* have been proposed for relations. If a relation schema is in one of these normal forms, we know that certain kinds of problems cannot arise. Considering the normal form of a given relation schema can help us to decide whether or not to decompose it further. If we decide that a relation schema must be decomposed further, we must choose a particular decomposition (i.e., a particular collection of smaller relations to replace the given relation).

With respect to the second question, two properties of decompositions are of particular interest. The *lossless-join* property enables us to recover any instance of the decomposed relation from corresponding instances of the smaller relations. The *dependency-preservation* property enables us to enforce any constraint on the original relation by simply enforcing the same constraints on each of the smaller relations. That is, we need not perform joins of the smaller relations to check whether a constraint on the original relation is violated.

From a performance standpoint, queries over the original relation may require us to join the decomposed relations. If such queries are common, the performance penalty of decomposing the relation may not be acceptable. In this case, we may choose to live with some of the problems of redundancy and not decompose the relation. It is important to be aware of the potential problems caused by such residual redundancy in the design and to take steps to avoid them (e.g., by adding some checks to application code). In some situations, decomposition could actually *improve* performance. This happens, for example, if most queries and updates examine only one of the decomposed relations, which is smaller than the original relation. We do not discuss the impact of decompositions on query performance in this chapter; this issue is covered in Section 20.8.

Our goal in this chapter is to explain some powerful concepts and design guidelines based on the theory of functional dependencies. A good database designer should have a firm grasp of normal forms and what problems they (do or do not) alleviate, the technique of decomposition, and potential problems with decompositions. For example, a designer often asks questions such as these: Is a relation in a given normal form? Is a decomposition dependency-preserving? Our objective is to explain when to raise these questions and the significance of the answers.

## 19.2 FUNCTIONAL DEPENDENCIES

A functional dependency (FD) is a kind of **lc** that generalizes the concept of a *key*. Let  $R$  be a relation schema and let  $X$  and  $Y$  be nonempty sets of attributes in  $R$ . We say that an instance  $r$  of  $R$  satisfies the FD  $X \rightarrow Y$ <sup>1</sup> if the following holds for every pair of tuples  $t_1$  and  $t_2$  in  $r$ :

If  $t_1.X = t_2.X$ , then  $t_1.Y = t_2.Y$ .

We use the notation  $t_1.X$  to refer to the projection of tuple  $t_1$  onto the attributes in  $X$ , in a natural extension of our TRC notation (see Chapter 4)  $t.a$  for referring to attribute  $a$  of tuple  $t$ . An FD  $X \rightarrow Y$  essentially says that if two tuples agree on the values in attributes  $X$ , they must also agree on the values in attributes  $Y$ .

Figure 19.3 illustrates the meaning of the FD  $AB \rightarrow C$  by showing an instance that satisfies this dependency. The first two tuples show that an FD is not the same as a key constraint: Although the FD is not violated,  $AB$  is clearly not a key for the relation. The third and fourth tuples illustrate that if two tuples differ in either the  $A$  field or the  $B$  field, they can differ in the  $C$  field without violating the FD. On the other hand, if we add a tuple  $(a_1, b_1, c_2, d_1)$  to the instance shown in this figure, the resulting instance would violate the FD; to see this violation, compare the first tuple in the figure with the new tuple.

$A$	$B$	$C$	$D$
a1	b1	c1	d1
a1	b1	c1	d2
a1	b2	c2	d1
a2	b1	c3	d1

Figure 19.3 An Instance that Satisfies  $AB \rightarrow C$

Recall that a *legal* instance of a relation must satisfy all specified **lc**s, including all specified FDs. As noted in Section 3.2, **lc**s must be identified and specified based on the semantics of the real-world enterprise being modeled. By looking at an instance of a relation, we might be able to tell that a certain FD does *not* hold. However, we can never deduce that an FD *does* hold by looking at one or more instances of the relation, because an FD, like other **lc**s, is a statement about *all* possible legal instances of the relation.

---

<sup>1</sup> $X \rightarrow Y$  is read as  $X$  *functionally determines*  $Y$ , or simply as  $X$  *determines*  $Y$ .

A primary key constraint is a special case of an FD. The attributes in the key play the role of  $X$ , and the set of all attributes in the relation plays the role of  $Y$ . Note, however, that the definition of an FD does not require that the set  $X$  be minimal; the additional minimality condition must be met for  $X$  to be a key. If  $X \rightarrow Y$  holds, where  $Y$  is the set of all attributes, and there is some (strictly contained) subset  $V$  of  $X$  such that  $V \rightarrow Y$  holds, then  $X$  is a *superkey*.

In the rest of this chapter, we see several examples of FDs that are not key constraints.

### 19.3 REASONING ABOUT FDS

Given a set of FDs over a relation schema  $R$ , typically several additional FDs hold over  $R$  whenever all of the given FDs hold. As an example, consider:

Workers(*ssn*, *name*, *lot*, *did*, *since*)

We know that  $ssn \rightarrow did$  holds, since *ssn* is the key, and FD  $did \rightarrow lot$  is given to hold. Therefore, in any legal instance of *Workers*, if two tuples have the same *ssn* value, they must have the same *did* value (from the first FD), and because they have the same *did* value, they must also have the same *lot* value (from the second FD). Therefore, the FD  $ssn \rightarrow lot$  also holds on *Workers*.

We say that an FD  $f$  is implied by a given set  $F$  of FDs if  $f$  holds on every relation instance that satisfies all dependencies in  $F$ ; that is,  $f$  holds whenever all FDs in  $F$  hold. Note that it is not sufficient for  $f$  to hold on some instance that satisfies all dependencies in  $F$ ; rather,  $f$  must hold on *every* instance that satisfies all dependencies in  $F$ .

#### 19.3.1 Closure of a Set of FDs

The set of all FDs implied by a given set  $F$  of FDs is called the closure of  $F$ , denoted as  $F^+$ . An important question is how we can infer, or compute, the closure of a given set  $F$  of FDs. The answer is simple and elegant. The following three rules, called Armstrong's Axioms, can be applied repeatedly to infer all FDs implied by a set  $F$  of FDs. We use  $X$ ,  $Y$ , and  $Z$  to denote *sets* of attributes over a relation schema  $R$ :

- Reflexivity: If  $X \supseteq Y$ , then  $X \rightarrow Y$ .
- Augmentation: If  $X \rightarrow Y$ , then  $XZ \rightarrow YZ$  for any  $Z$ .
- Transitivity: If  $X \rightarrow Y$  and  $Y \rightarrow Z$ , then  $X \rightarrow Z$ .

**Theorem 1** *Armstrong's Axioms are **sound**, in that they generate only FDs in  $F^+$  when applied to a set  $F$  of FDs. They are also **complete**, in that repeated application of these rules will generate all FDs in the closure  $F^+$ .*

The soundness of Armstrong's Axioms is straightforward to prove. Completeness is harder to show; see Exercise 19.17.

It is convenient to use some additional rules while reasoning about  $F^+$ :

- **Union:** If  $X \rightarrow Y$  and  $X \rightarrow Z$ , then  $X \rightarrow YZ$ .
- **Decomposition:** If  $X \rightarrow YZ$ , then  $X \rightarrow Y$  and  $X \rightarrow Z$ .

These additional rules are not essential; their soundness can be proved using Armstrong's Axioms.

To illustrate the use of these inference rules for FDs, consider a relation schema  $ABC$  with FDs  $A \rightarrow B$  and  $B \rightarrow C$ . In a trivial FD, the right side contains only attributes that also appear on the left side; such dependencies always hold due to reflexivity. Using reflexivity, we can generate all trivial dependencies, which are of the form:

$$X \rightarrow Y, \text{ where } Y \subseteq X, X \subseteq ABC, \text{ and } Y \subseteq ABC.$$

From transitivity we get  $A \rightarrow C$ . From augmentation we get the nontrivial dependencies:

$$AC \rightarrow BC, AB \rightarrow AC, AB \rightarrow C.$$

As another example, we use a more elaborate version of Contracts:

$$\text{Contracts}(\text{contractid}, \text{supplierid}, \text{projectid}, \text{deptid}, \text{partid}, \text{qty}, \text{value})$$

We denote the schema for Contracts as  $CSJDPQV$ . The meaning of a tuple is that the contract with *contractid*  $C$  is an agreement that supplier  $S(\text{supplierid})$  will supply  $Q$  items of *part*? (*partid*) to project  $J$  (*projectid*) associated with department  $D$  (*deptid*); the value  $V$  of this contract is equal to *value*.

The following facts are known to hold:

1. The contract id  $C$  is a key:  $C \rightarrow CSJDPQV$ .
2. A project purchases a given part using a single contract:  $(J, P) \rightarrow C$ .



3. A department purchases at most one part from a supplier:  $SD \rightarrow P$ .

Several additional FDs hold in the closure of the set of given FDs:

From  $JP \rightarrow C$ ,  $C \rightarrow CSJDPQV$ , and transitivity, we infer  $JP \rightarrow CSJDPQV$ .

From  $SD \rightarrow P$  and augmentation, we infer  $SDJ \rightarrow JP$ .

From  $SDJ \rightarrow JP$ ,  $JP \rightarrow CSJDPQV$ , and transitivity, we infer  $SDJ \rightarrow CSJDPQV$ . (Incidentally, while it may appear tempting to do so, we *cannot* conclude  $SD \rightarrow CSJDPQV$ , canceling  $J$  on both sides. FD inference is not like arithmetic multiplication!)

We can infer several additional FDs that are in the closure by using augmentation or decomposition. For example, from  $C \rightarrow CSJDPQV$ , using decomposition, we can infer:

$C \rightarrow C$ ,  $C \rightarrow S$ ,  $C \rightarrow J$ ,  $C \rightarrow D$ , and so forth

Finally, we have a number of trivial FDs from the reflexivity rule.

### 19.3.2 Attribute Closure

If we just want to check whether a given dependency, say,  $X \rightarrow Y$ , is in the closure of a set  $F$  of FDs, we can do so efficiently without computing  $F^+$ . We first compute the attribute closure  $X^+$  with respect to  $F$ , which is the set of attributes  $A$  such that  $X \rightarrow A$  can be inferred using the Armstrong Axioms. The algorithm for computing the attribute closure of a set  $X$  of attributes is shown in Figure 19.4.

```

closure = X;
repeat until there is no change: {
    if there is an FD  $U \rightarrow V$  in  $F$  such that  $U \subseteq \text{closure}$ ,
        then set  $\text{closure} = \text{closure} \cup V$ 
}
```

Figure 19.4 Computing the Attribute Closure of Attribute Set  $X$

**Theorem 2** The algorithm shown in Figure 19.4 computes the attribute closure  $X^+$  of the attribute set  $X$  with respect to the set of FDs  $F$ .

The proof of this theorem is considered in Exercise 19.15. This algorithm can be modified to find keys by starting with set  $X$  containing a single attribute and stopping as soon as *closure* contains all attributes in the relation schema. By varying the starting attribute and the order in which the algorithm considers FDs, we can obtain all candidate keys.

## 19.4 NORMAL FORMS

Given a relation schema, we need to decide whether it is a good design or we need to decompose it into smaller relations. Such a decision must be guided by an understanding of what problems, if any, arise from the current schema. To provide such guidance, several normal forms have been proposed. If a relation schema is in one of these normal forms, we know that certain kinds of problems cannot arise.

The normal forms based on FDs are *first normal form (1NF)*, *second normal form (2NF)*, *third normal form (3NF)*, and *Boyce-Codd normal form (BCNF)*. These forms have increasingly restrictive requirements: Every relation in BCNF is also in 3NF, every relation in 3NF is also in 2NF, and every relation in 2NF is in 1NF. A relation is in first normal form if every field contains only atomic values, that is, no lists or sets. This requirement is implicit in our definition of the relational model. Although some of the newer database systems are relaxing this requirement, in this chapter we assume that it always holds. 2NF is mainly of historical interest. 3NF and BCNF are important from a database design standpoint.

While studying normal forms, it is important to appreciate the role played by FDs. Consider a relation schema  $R$  with attributes  $ABC$ . In the absence of any FDs, any set of ternary tuples is a legal instance and there is no potential for redundancy. On the other hand, suppose that we have the FD  $A \rightarrow B$ . Now if several tuples have the same  $A$  value, they must also have the same  $B$  value. This potential redundancy can be predicted using the FD information. If more detailed FDs are specified, we may be able to detect more subtle redundancies as well.

We primarily discuss redundancy revealed by FD information. In Section 19.8, we discuss more sophisticated FDs called *multivalued dependencies* and *join dependencies* and normal forms based on them.

### 19.4.1 Boyce-Codd Normal Form

Let  $R$  be a relation schema,  $F$  be the set of FDs given to hold over  $R$ ,  $X$  be a subset of the attributes of  $R$ , and  $A$  be an attribute of  $R$ .  $R$  is in Boyce-Codd

normal form if, for every  $F1)X \rightarrow A$  in  $F$ , one of the following statements is true:

- $A \in X$ ; that is, it is a trivial FD, or
- $X$  is a superkey.

Intuitively, in a BCNF relation, the only nontrivial dependencies are those in which a key determines some attribute(s). Therefore, each tuple can be thought of as an entity or relationship, identified by a key and described by the remaining attributes. Rent (in [425]) puts this colorfully, if a little loosely: "Each attribute must describe [an entity or relationship identified by] the key, the whole key, and nothing but the key." If we use ovals to denote attributes or sets of attributes and draw arcs to indicate FDs, a relation in BCNF has the structure illustrated in Figure 19.5, considering just one key for simplicity. (If there are several candidate keys, each candidate key can play the role of KEY in the figure, with the other attributes being the ones not in the chosen candidate key.)

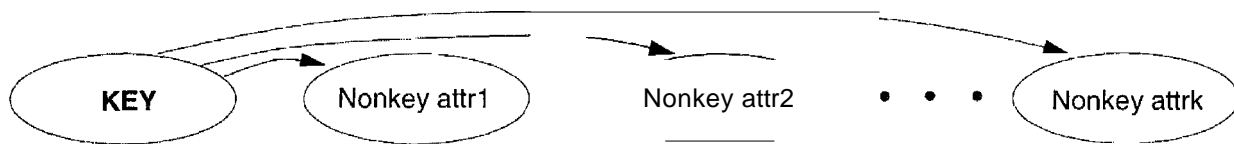


Figure 19.5 FDs in a BCNF Relation

BCNF ensures that no redundancy can be detected using FD information alone. It is thus the most desirable normal form (from the point of view of redundancy) if we take into account only FD information. This point is illustrated in Figure 19.6.

$X$	$Y$	$A$
$x$	$y_1$	$a$
$x$	$y_2$	?

Figure 19.6 Instance Illustrating BCNF

This figure shows (two tuples in) an instance of a relation with three attributes  $X$ ,  $Y$ , and  $A$ . There are two tuples with the same value in the  $X$  column. Now suppose that we know that this instance satisfies an FD  $X \rightarrow A$ . We can see that one of the tuples has the value  $a$  in the  $A$  column. What can we infer about the value in the  $A$  column in the second tuple? Using the FD, we can conclude that the second tuple also has the value  $a$  in this column. (Note that this is really the only kind of inference we can make about values in the fields of tuples by using FDs.)

But is this situation not an example of redundancy? We appear to have stored the value  $a$  twice. Can such a situation arise in a BCNF relation? The answer is No! If this relation is in BCNF, because  $A$  is distinct from  $X$ , it follows that  $X$  must be a key. (Otherwise, the FD  $X \rightarrow A$  would violate BCNF.) If  $X$  is a key, then  $Y_1 = Y_2$ , which means that the two tuples are identical. Since a relation is defined to be a *set* of tuples, we cannot have two copies of the same tuple and the situation shown in Figure 19.6 cannot arise.

Therefore, if a relation is in BCNF, every field of every tuple records a piece of information that cannot be inferred (using only FDs) from the values in all other fields in (all tuples of) the relation instance.

### 19.4.2 Third Normal Form

Let  $R$  be a relation schema,  $F$  be the set of FDs given to hold over  $R$ ,  $X$  be a subset of the attributes of  $R$ , and  $A$  be an attribute of  $R$ .  $R$  is in third normal form if, for every FD  $X \rightarrow A$  in  $F$ , one of the following statements is true:

- $A \in X$ ; that is, it is a trivial FD, or
- $X$  is a superkey, or
- $A$  is part of some key for  $R$ .

The definition of 3NF is similar to that of BCNF, with the only difference being the third condition. Every BCNF relation is also in 3NF. To understand the third condition, recall that a key for a relation is a *minimal* set of attributes that uniquely determines all other attributes.  $A$  must be part of a key (any key, if there are several). It is not enough for  $A$  to be part of a superkey, because the latter condition is satisfied by every attribute! Finding all keys of a relation schema is known to be an NP-complete problem, and so is the problem of determining whether a relation schema is in 3NF.

Suppose that a dependency  $X \rightarrow A$  causes a violation of 3NF. There are two cases:

- $X$  is a *proper subset of some key*  $K$ . Such a dependency is sometimes called a **partial dependency**. In this case, we store  $(X, A)$  pairs redundantly. As an example, consider the Reserves relation with attributes  $SBDC$  from Section 19.7.4. The only key is  $SEI$ , and we have the FD  $S \rightarrow C$ . We store the credit card number for a sailor as many times as there are reservations for that sailor.
- $X$  is *not a proper subset of any key*. Such a dependency is sometimes called a **transitive dependency**, because it means we have a chain of

dependencies  $X \rightarrow A$ . The problem is that we cannot associate an  $X$  value with a  $K$  value unless we also associate an  $A$  value with an  $X$  value. As an example, consider the *Hourly\_Emps* relation with attributes *SNLRWH* from Section 19.7.1. The only key is  $S$ , but there is an FD  $R \rightarrow W$ , which gives rise to the chain  $S \rightarrow R \rightarrow W$ . The consequence is that we cannot record the fact that employee  $S$  has rating  $R$  without knowing the hourly wage for that rating. This condition leads to insertion, deletion, and update anomalies.

Partial dependencies are illustrated in Figure 19.7, and transitive dependencies are illustrated in Figure 19.8. Note that in Figure 19.8, the set  $X$  of attributes may or may not have some attributes in common with  $KEY$ ; the diagram should be interpreted as indicating only that  $X$  is not a subset of  $KEY$ .

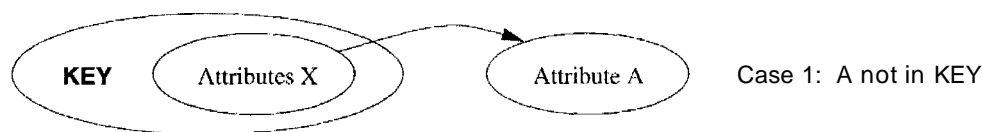


Figure 19.7 Partial Dependencies

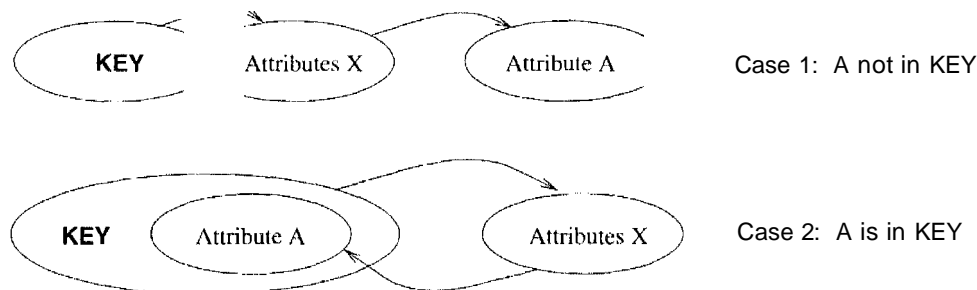


Figure 19.8 Transitive Dependencies

The motivation for 3NF is rather technical. By making an exception for certain dependencies involving key attributes, we can ensure that every relation schema can be decomposed into a collection of 3NF relations using only decompositions that have certain desirable properties (Section 19.5). Such a guarantee does not exist for BCNF relations; the 3NF definition weakens the BCNF requirements just enough to make this guarantee possible. We may therefore compromise by settling for a 3NF design. As we see in Chapter 20, we may sometimes accept this compromise (or even settle for a non-3NF schema) for other reasons as well.

Unlike BCNF, however, some redundancy is possible with 3NF. The problems associated with partial and transitive dependencies persist if there is a nontrivial dependency  $X \rightarrow A$  and  $X$  is not a superkey, even if the relation is in 3NF because  $A$  is part of a key. To understand this point, let us revisit the *Reserves*

relation with attributes *SEDe* and the FD  $S \rightarrow C$ , which states that a sailor uses a unique credit card to pay for reservations.  $S$  is not a key, and  $C$  is not part of a key. (In fact, the only key is *SED*.) Hence, this relation is not in 3NF;  $(S, C)$  pairs are stored redundantly. However, if we also know that credit cards uniquely identify the owner, we have the FD  $C \rightarrow S$ , which means that *CBD* is also a key for Reserves. Therefore, the dependency  $S \rightarrow C$  does not violate 3NF, and Reserves is in 3NF. Nonetheless, in all tuples containing the same  $S$  value, the same  $(S, C)$  pair is redundantly recorded.

For completeness, we remark that the definition of second normal form is essentially that partial dependencies are not allowed. Thus, if a relation is in 3NF (which precludes both partial and transitive dependencies), it is also in 2NF.

## 19.5 PROPERTIES OF DECOMPOSITIONS

Decomposition is a tool that allows us to eliminate redundancy. As noted in Section 19.1.3, however, it is important to check that a decomposition does not introduce new problems. In particular, we should check whether a decomposition allows us to recover the original relation, and whether it allows us to check integrity constraints efficiently. We discuss these properties next.

### 19.5.1 Lossless-Join Decomposition

Let  $R$  be a relation schema and let  $F$  be a set of FDs over  $R$ . A decomposition of  $R$  into two schemas with attribute sets  $X$  and  $Y$  is said to be a lossless-join decomposition with respect to  $F$  if, for every instance  $r$  of  $R$  that satisfies the dependencies in  $F$ ,  $\pi_X(r) \bowtie \pi_Y(r) = r$ . In other words, we can recover the original relation from the decomposed relations.

This definition can easily be extended to cover a decomposition of  $R$  into more than two relations. It is easy to see that  $r \subseteq \pi_X(r) \bowtie \pi_Y(r)$  always holds. In general, though, the other direction does not hold. If we take projections of a relation and recombine them using natural join, we typically obtain some tuples that were not in the original relation. This situation is illustrated in Figure 19.9.

By replacing the instance  $r$  shown in Figure 19.9 with the instances  $\pi_{SP}(r)$  and  $\pi_{PI}(r)$ , we lose some information. In particular, suppose that the tuples in  $r$  denote relationships. We can no longer tell that the relationships  $(s_1, p_1, d_3)$  and  $(s_3, p_1, d_1)$  do not hold. The decomposition of schema *SPD* into *SP* and *PI* is therefore lossy if the instance  $r$  shown in the figure is legal, that is, if this

<i>S</i>	<i>P</i>	<i>D</i>		<i>S</i>	<i>P</i>	<i>P</i>	<i>D</i>	<i>S</i>	<i>P</i>	<i>D</i>
s1	pl.	ell		s1	pI	p1	d1	81	pI	ell
s2	p2	d2		s2	p2	p2	d2	82	p2	d2
s3	pl.	d3		s3	pI	pI	d3	83	pI	d3
								s1	pI	d3
								s3	pI	ell

Instance *r*                       $\pi_{SP}(r)$                        $\pi_{PD}(r)$                        $\pi_{SP}(r) \bowtie \pi_{PD}(r)$

Figure 19.9 Instances Illustrating Lossy Decompositions

instance could arise in the enterprise being modeled. (Observe the similarities between this example and the Contracts relationship set in Section 2.5.3.)

*All decompositions used to eliminate redundancy must be lossless.* The following simple test is very useful:

**Theorem 3** *Let  $R$  be a relation and  $F$  be a set of FDs that hold over  $R$ . The decomposition of  $R$  into relations with attribute sets  $R_1$  and  $R_2$  is lossless if and only if  $F$  contains either the FD  $R_1 \twoheadrightarrow R_2$  or the FDR  $R_1 \twoheadrightarrow R_2$ .*

In other words, the attributes common to  $R_1$  and  $R_2$  must contain a key for either  $R_1$  or  $R_2$ .<sup>2</sup> If a relation is decomposed into more than two relations, an efficient (time polynomial in the size of the dependency set) algorithm is available to test whether or not the decomposition is lossless, but we will not discuss it.

Consider the Hourly\_Emps relation again. It has attributes *SNLRWH*, and the FD  $R \twoheadrightarrow W$  causes a violation of 3NF. We dealt with this violation by decomposing the relation into *SNLRH* and *RW*. Since  $R$  is common to both decomposed relations and  $R \twoheadrightarrow W$  holds, this decomposition is lossless-join.

This example illustrates a general observation that follows from Theorem 3:

If an FD  $X \twoheadrightarrow Y$  holds over a relation  $R$  and  $X \cap Y$  is empty, the decomposition of  $R$  into  $R - Y$  and  $XY$  is lossless.

$X$  appears in both  $R - Y$  (since  $X \cap Y$  is empty) and  $XY$ , and it is a key for  $XY$ .

<sup>2</sup>See Exercise 19.19 for a proof of Theorem 3. Exercise 19.11 illustrates that the 'only if' claim depends on the assumption that only functional dependencies can be specified as integrity constraints.

Another important observation, which we state without proof, has to do with repeated decomposition. Suppose that a relation  $R$  is decomposed into  $R_1$  and  $R_2$  through a lossless-join decomposition, and that  $R_1$  is decomposed into  $R_{11}$  and  $R_{12}$  through another lossless-join decomposition. Then, the decomposition of  $R$  into  $R_{11}$ ,  $R_{12}$ , and  $R_2$  is lossless-join; by joining  $R_{11}$  and  $R_{12}$ , we can recover  $R_1$ , and by then joining  $R_1$  and  $R_2$ , we can recover  $R$ .

### 19.5.2 Dependency-Preserving Decomposition

Consider the *Contracts* relation with attributes  $C, S, J, D, P, C, J, V$  from Section 19.3.1. The given FDs are  $C \rightarrow CSJDPQV$ ,  $JP \rightarrow C$ , and  $SD \rightarrow P$ . Because  $SD$  is not a key the dependency  $SD \rightarrow P$  causes a violation of BCNF.

We can decompose *Contracts* into two relations with schemas  $CSJDQV$  and  $SDP$  to address this violation; the decomposition is lossless-join. There is one subtle problem, however. We can enforce the integrity constraint  $JP \rightarrow C$  easily when a tuple is inserted into *Contracts* by ensuring that no existing tuple has the same  $JP$  values (as the inserted tuple) but different  $C$  values. Once we decompose *Contracts* into  $CSJDQV$  and  $SDP$ , enforcing this constraint requires an expensive join of the two relations whenever a tuple is inserted into  $CSJDQV$ . We say that this decomposition is not dependency-preserving.

Intuitively, a *dependency-preserving decomposition* allows us to enforce all FDs by examining a single relation instance on each insertion or modification of a tuple. (Note that deletions cannot cause violation of FDs.) To define dependency-preserving decompositions precisely, we have to introduce the concept of a projection of FDs.

Let  $R$  be a relation schema that is decomposed into two schemas with attribute sets  $X$  and  $Y$ , and let  $F$  be a set of FDs over  $R$ . The **projection of  $F$  on  $X$**  is the set of FDs in the closure  $F^+$  (not just  $F$ !) that involve only attributes in  $X$ . We denote the projection of  $F$  on attributes  $X$  as  $F_X$ . Note that a dependency  $U \rightarrow V$  in  $F^+$  is in  $F_X$  only if *all* the attributes in  $U$  and  $V$  are in  $X$ .

The decomposition of relation schema  $R$  with FDs  $F$  into schemas with attribute sets  $X$  and  $Y$  is dependency-preserving if  $(F_X \cup F_Y)^+ = F^+$ . That is, if we take the dependencies in  $F_X$  and  $F_Y$  and compute the closure of their union, we get back all dependencies in the closure of  $F$ . Therefore, we need to enforce only the dependencies in  $F_X$  and  $F_Y$ ; all FDs in  $F^+$  are then sure to be satisfied. To enforce  $F_X$ , we need to examine only relation  $X$  (on inserts to that relation). To enforce  $F_Y$ , we need to examine only relation  $Y$ .



To appreciate the need to consider the closure  $F^+$  while considering the projection of  $F$ , suppose that a relation  $R$  with attributes  $ABC$  is decomposed into relations with attributes  $AB$  and  $BC$ . The set  $F$  of FDs over  $R$  includes  $A \rightarrow B$ ,  $B \rightarrow C$ , and  $C \rightarrow A$ . Of these,  $A \rightarrow B$  is in  $F_{AB}$  and  $B \rightarrow C$  is in  $F_{BC}$ . But is this decomposition dependency-preserving? What about  $C \rightarrow A$ ? This dependency is not implied by the dependencies listed (thus far) for  $F_{AB}$  and  $F_{BC}$ .

The closure of  $F$  contains all dependencies in  $F$  plus  $A \rightarrow C$ ,  $B \rightarrow A$ , and  $C \rightarrow B$ . Consequently,  $F_{AB}$  also contains  $B \rightarrow A$ , and  $F_{BC}$  contains  $C \rightarrow B$ . Therefore,  $F_{AB} \cup F_{BC}$  contains  $A \rightarrow B$ ,  $B \rightarrow C$ ,  $B \rightarrow A$ , and  $C \rightarrow B$ . The closure of the dependencies in  $F_{AB}$  and  $F_{BC}$  now includes  $C \rightarrow A$  (which follows from  $C \rightarrow B$ ,  $B \rightarrow A$ , and transitivity). Thus, the decomposition preserves the dependency  $C \rightarrow A$ .

A direct application of the definition gives us a straightforward algorithm for testing whether a decomposition is dependency-preserving. (This algorithm is exponential in the size of the dependency set. A polynomial algorithm is available; see Exercise 19.9.)

We began this section with an example of a lossless-join decomposition that was not dependency-preserving. Other decompositions are dependency-preserving, but not lossless. A simple example consists of a relation  $ABC$  with FD  $A \rightarrow B$  that is decomposed into  $AB$  and  $BC$ .

## 19.6 NORMALIZATION

Having covered the concepts needed to understand the role of normal forms and decompositions in database design, we now consider algorithms for converting relations to BCNF or 3NF. If a relation schema is not in BCNF, it is possible to obtain a lossless-join decomposition into a collection of BCNF relation schemas. Unfortunately, there may be no dependency-preserving decomposition into a collection of BCNF relation schemas. However, there is always a dependency-preserving, lossless-join decomposition into a collection of 3NF relation schemas.

### 19.6.1 Decomposition into BCNF

We now present an algorithm for decomposing a relation schema  $R$  with a set of FDs into a collection of BCNF relation schemas:

1. Suppose that  $R$  is not in BCNF. Let  $X \subset R$ ,  $A$  be a single attribute in  $R$ , and  $X \rightarrow A$  be an FD that causes a violation of BCNF. Decompose  $R$  into  $R - A$  and  $XA$ .
2. If either  $R - A$  or  $XA$  is not in BCNF, decompose them further by a recursive application of this algorithm.

$R - A$  denotes the set of attributes other than  $A$  in  $R$ , and  $XA$  denotes the union of attributes in  $X$  and  $A$ . Since  $X \rightarrow A$  violates BCNF, it is not a trivial dependency; further,  $A$  is a single attribute. Therefore,  $A$  is not in  $X$ ; that is,  $X \cap A$  is empty. Therefore, each decomposition carried out in Step 1 is lossless-join.

The set of dependencies associated with  $R - A$  and  $XA$  is the projection of  $F$  onto their attributes. If one of the new relations is not in BCNF, we decompose it further in Step 2. Since a decomposition results in relations with strictly fewer attributes, this process terminates, leaving us with a collection of relations that are all in BCNF. Further, joining instances of the (two or more) relations obtained through this algorithm yields precisely the corresponding instance of the original relation (i.e., the decomposition into a collection of relations each of which is in BCNF is a lossless-join decomposition).

Consider the Contracts relation with attributes  $CSDJPQV$  and key  $C$ . We are given FDs  $JP \rightarrow C$  and  $3D \rightarrow P$ . By using the dependency  $3D \rightarrow P$  to guide the decomposition, we get the two schemas  $3DP$  and  $C5JDQV$ .  $3DP$  is in BCNF. Suppose that we also have the constraint that each project deals with a single supplier:  $.I \rightarrow S$ . This means that the schema  $CSJDQV$  is not in BCNF. So we decompose it further into  $JS$  and  $CJDQV$ .  $C \rightarrow JDQV$  holds over  $CJDQV$ ; the only other FDs that hold are those obtained from this PI by augmentation, and therefore all FDs contain a key in the left side. Thus, each of the schemas  $3DP$ ,  $JS$ , and  $CJDQV$  is in BCNF, and this collection of schemas also represents a lossless-join decomposition of  $CSJDQV$ .

The steps in this decomposition process can be visualized as a tree, as shown in Figure 19.10. The root is the original relation  $CSJDQV$ , and the leaves are the BCNF relations that result from the decomposition algorithm:  $3DP$ ,  $JS$ , and  $CJDQV$ . Intuitively, each internal node is replaced by its children through a single decomposition step guided by the FD shown just below the node.

## Redundancy in BCNF Revisited

The decomposition of  $CSJDQV$  into  $3DP$ ,  $JS$ , and  $CJDQV$  is not dependency-preserving. Intuitively, dependency  $JP \rightarrow C$  cannot be enforced without a join. One way to deal with this situation is to add a relation with attributes  $GP$ . In

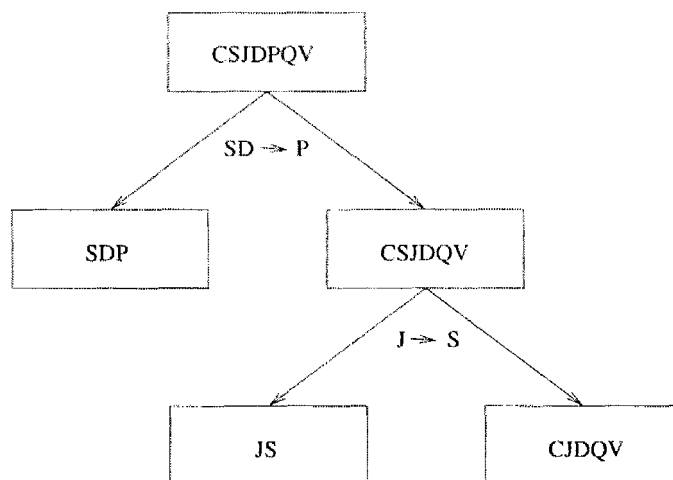


Figure 19.10 Decomposition of *CSJDPQV* into *SDP*, *JS*, and *CJDQV*

effect, this solution amounts to storing some information redundantly to make the dependency enforcement cheaper.

This is a subtle point: Each of the schemas *CJP*, *SDP*, *JS*, and *CJDQV* is in BCNF, yet some redundancy can be predicted by FD information. In particular, if we join the relation instances for *SDP* and *CJDQV* and project the result onto the attributes *CJP*, we must get exactly the instance stored in the relation with schema *CJP*. We saw in Section 19.4.1 that there is no such redundancy within a single BCNF relation. This example shows that redundancy can still occur across relations, even though there is no redundancy within a relation.

## Alternatives in Decomposing to BCNF

Suppose several dependencies violate BCNF. Depending on which of these dependencies we choose to guide the next decomposition step, we may arrive at quite different collections of BCNF relations. Consider *Contracts*. We just decomposed it into *SDP*, *JS*, and *CJDQV*. Suppose we choose to decompose the original relation *CSJDPQV* into *JS* and *CJDPQV*, based on the FD  $I \rightarrow S$ . The only dependencies that hold over *CJDPQV* are  $IP \rightarrow C$  and the key dependency  $C \rightarrow C.IDPQV$ . Since *IP* is a key, *CJDPQV* is in BCNF. Thus, the schemas *JS* and *CJDPQV* represent a lossless-join decomposition of *Contracts* into BCNF relations.

The lesson to be learned here is that the theory of dependencies can tell us when there is redundancy and give us clues about possible decompositions to address the problem, but it cannot discriminate among decomposition alternatives. A designer has to consider the alternatives and choose one based on the semantics of the application.

## BCNF and Dependency-Preservation

Sometimes, there simply is no decomposition into BCNF that is dependency-preserving. As an example, consider the relation schema  $SBD$ , in which a tuple denotes that sailor  $S$  has reserved boat  $B$  on date  $D$ . If we have the FDs  $SB \rightarrow D$  (a sailor can reserve a given boat for at most one day) and  $D \rightarrow B$  (on any given day at most one boat can be reserved),  $SBD$  is not in BCNF because  $D$  is not a key. If we try to decompose it, however, we cannot preserve the dependency  $SB \rightarrow D$ .

### 19.6.2 Decomposition into 3NF

Clearly, the approach we outlined for lossless-join decomposition into BCNF also gives us a lossless-join decomposition into 3NF. (Typically, we can stop a little earlier if we are satisfied with a collection of 3NF relations.) But this approach does not ensure dependency-preservation.

A simple modification, however, yields a decomposition into 3NF relations that is lossless-join and dependency-preserving. Before we describe this modification, we need to introduce the concept of a minimal cover for a set of FDs.

### Minimal Cover for a Set of FDs

A minimal cover for a set  $F$  of FDs is a set  $G$  of FDs such that:

1. Every dependency in  $G$  is of the form  $X \rightarrow A$ , where  $A$  is a single attribute.
2. The closure  $F^+$  is equal to the closure  $G^+$ .
3. If we obtain a set  $H$  of dependencies from  $G$  by deleting one or more dependencies or by deleting attributes from a dependency in  $G$ , then  $H^+ \neq F^+$ .

Intuitively, a minimal cover for a set  $F$  of FDs is an equivalent set of dependencies that is *minimal* in two respects: (1) Every dependency is as small as possible; that is, each attribute on the left side is necessary and the right side is a single attribute. (2) Every dependency in it is required for the closure to be equal to  $F^+$ .

As an example, let  $F$  be the set of dependencies:

$$A \rightarrow B, ABCD \rightarrow E, EF \rightarrow G, EF \rightarrow H, \text{ and } CDF \rightarrow EG.$$

First, let us rewrite  $CDF \rightarrow EG$  so that every right side is a single attribute:

$$ACDF \rightarrow E \text{ and } ACDF \rightarrow G,$$

Next consider  $ACDF \rightarrow G$ . This dependency is implied by the following FDs:

$$A \rightarrow B, ABCD \rightarrow E, \text{ and } EF \rightarrow G,$$

Therefore, we can delete it. Similarly, we can delete  $ACDF \rightarrow E$ . Next consider  $ABCD \rightarrow E$ . Since  $A \rightarrow B$  holds, we can replace it with  $ACD \rightarrow E$ . (At this point, the reader should verify that each remaining FD is minimal and required.) Thus, a minimal cover for  $F$  is the set:

$$A \rightarrow B, ACD \rightarrow E, EF \rightarrow G, \text{ and } EF \rightarrow H,$$

The preceding example illustrates a general algorithm for obtaining a minimal cover of a set  $F$  of FDs:

1. **Put the FDs in a Standard Form:** Obtain a collection  $G$  of equivalent FDs with a single attribute on the right side (using the decomposition axiom),
2. **Minimize the Left Side of Each FD:** For each FD in  $G$ , check each attribute in the left side to see if it can be deleted while preserving equivalence to  $F^+$ ,
3. **Delete Redundant FDs:** Check each remaining FD in  $G$  to see if it can be deleted while preserving equivalence to  $F^+$ ,

Note that the order in which we consider FDs while applying these steps could produce different minimal covers; there could be several minimal covers for a given set of FDs,

It's important, it is necessary to minimize the left sides of FDs *before* checking for redundant FDs. If these two steps are reversed, the final set of FDs could still contain some redundant FDs (i.e., not be a minimal cover), as the following example illustrates. Let  $F$  be the set of dependencies, each of which is already in the standard form:

$$ABCD \rightarrow E, E \rightarrow D, A \rightarrow B, \text{ and } AC \rightarrow D,$$

Observe that none of these FDs is redundant; if we checked for redundant FDs first, we would get the same set of FDs  $F$ . The left side of  $ABCD \rightarrow E$  can be replaced by  $AC$  while preserving equivalence to  $F^+$ , and we would stop here if we checked for redundant FDs in  $F$  before minimizing the left sides. However, the set of FDs we have is not a minimal cover:

$$AC \rightarrow E, E \twoheadrightarrow D, A \rightarrow B, \text{ and } AC \rightarrow D.$$

From transitivity, the first two FDs imply the last FD, which can therefore be deleted while preserving equivalence to  $F^+$ . The important point to note is that  $AC \rightarrow D$  becomes redundant only after we replace  $AB \twoheadrightarrow E$  with  $AC \rightarrow E$ . If we minimize left sides of FDs first and then check for redundant FDs, we are left with the first three FDs in the preceding list, which is indeed a minimal cover for  $F$ .

## Dependency-Preserving Decomposition into 3NF

Returning to the problem of obtaining a lossless-join, dependency-preserving decomposition into 3NF relations, let  $R$  be a relation with a set  $F$  of FDs that is a minimal cover, and let  $R_1, R_2, \dots, R_n$  be a lossless-join decomposition of  $R$ . For  $1 \leq i \leq n$ , suppose that each  $R_i$  is in 3NF and let  $F_i$  denote the projection of  $F$  onto the attributes of  $R_i$ . Do the following:

- Identify the set  $N$  of dependencies in  $F$  that is not preserved, that is, not included in the closure of the union of  $F_i$ s.
- For each FD  $X \rightarrow A$  in  $N$ , create a relation schema  $XA$  and add it to the decomposition of  $R$ .

Obviously, every dependency in  $F$  is preserved if we replace  $R$  by the  $R_i$ s plus the schemas of the form  $XA$  added in this step. The  $R_i$ s are given to be in 3NF. We can show that each of the schemas  $XA$  is in 3NF as follows: Since  $X \rightarrow A$  is in the minimal cover  $F$ ,  $Y \rightarrow A$  does not hold for any  $Y$  that is a strict subset of  $X$ . Therefore,  $X$  is a key for  $XA$ . Further, if any other dependencies hold over  $XA$ , the right side can involve only attributes in  $X$  because  $A$  is a single attribute (because  $X \rightarrow A$  is an FD in a minimal cover). Since  $X$  is a key for  $XA$ , none of these additional dependencies causes a violation of 3NF (although they might cause a violation of BCNF).

As an optimization, if the set  $N$  contains several FDs with the same left side, say,  $X \rightarrow A_1, X \rightarrow A_2, \dots, X \twoheadrightarrow A_n$ , we can replace them with a single equivalent FD  $X \rightarrow A_1 \dots A_n$ . Therefore, we produce one relation schema  $XA_1 \dots A_n$ , instead of several schemas  $XA_1, \dots, XA_n$ , which is generally preferable.

Consider the Contracts relation with attributes  $CSJDPQV$  and FDs  $JP \rightarrow C$ ,  $SD \rightarrow P$ , and  $J \twoheadrightarrow S$ . If we decompose  $CSJDPQV$  into  $SDIJ$  and  $CSJDQV$ , then  $SDP$  is in BCNF, but  $CSJDQV$  is not even in 3NF. So we decompose it further into  $JS$  and  $CJDQV$ . The relation schemas  $SDP$ ,  $JS$ , and  $CJDQV$  are in 3NF (in fact, in BCNF), and the decomposition is lossless-join. However,

the dependency  $JP \rightarrow C$  is not preserved. This problem can be addressed by adding a relation schema  $CJP$  to the decomposition.

## 3NF Synthesis

We assumed that the design process starts with an ER diagram, and that our use of FDs is primarily to guide decisions about decomposition. The algorithm for obtaining a lossless-join, dependency-preserving decomposition was presented in the previous section from this perspective-----a lossless-join decomposition into 3NF is straightforward, and the algorithm addresses dependency-preservation by adding extra relation schemas.

An alternative approach, called synthesis, is to take all the attributes over the original relation  $R$  and a minimal cover  $F$  for the FDs that hold over it and add a relation schema  $XA$  to the decomposition of  $R$  for each FD  $X \rightarrow A$  in  $F$ .

The resulting collection of relation schemas is in 3NF and preserves all FDs. If it is not a lossless-join decomposition of  $R$ , we can make it so by adding a relation schema that contains just those attributes that appear in some key. This algorithm gives us a lossless-join, dependency-preserving decomposition into 3NF and has polynomial complexity-----polynomial algorithms are available for computing minimal covers, and a key can be found in polynomial time (even though finding all keys is known to be NP-complete). The existence of a polynomial algorithm for obtaining a lossless-join, dependency-preserving decomposition into 3NF is surprising when we consider that testing whether a given schema is in 3NF is NP-complete.

As an example, consider a relation  $ABC$  with FDs  $F = \{A \rightarrow B, C \rightarrow B\}$ . The first step yields the relation schemas  $AB$  and  $BC$ . This is not a lossless-join decomposition of  $ABC$ ;  $AB \bowtie BC$  is  $B$ , and neither  $B \rightarrow A$  nor  $B \rightarrow C$  is in  $F^+$ . If we add a schema  $AC$ , we have the lossless-join property as well. Although the collection of relations  $AB, BC$ , and  $AC$  is a dependency-preserving, lossless-join decomposition of  $ABC$ , we obtained it through a process of *synthesis*, rather than through a process of repeated decomposition. We note that the decomposition produced by the synthesis approach heavily depends on the minimal cover used.

As another example of the synthesis approach, consider the Contracts relation with attributes  $CSJDPQV$  and the following FDs:

$$C \rightarrow CSJDPQV, SP \rightarrow C, SD \rightarrow P, \text{ and } J \rightarrow S.$$

This set of FDs is not a minimal cover, and so we must find one. We first replace  $G \rightarrow CSJDPQV$  with the FDs:

$$C \rightarrow S, C \rightarrow J, C \rightarrow IJ, C \twoheadrightarrow P, C \rightarrow Q, \text{ and } C \rightarrow V.$$

The FD  $C \rightarrow P$  is implied by  $C \rightarrow S$ ,  $C \twoheadrightarrow D$ , and  $SD \rightarrow P$ ; so we can delete it. The FD  $C \rightarrow S$  is implied by  $C \rightarrow J$  and  $J \twoheadrightarrow S$ ; so we can delete it. This leaves us with a minimal cover:

$$C \rightarrow J, C \rightarrow IJ, C \rightarrow Q, C \rightarrow V, JP \rightarrow C, SD \rightarrow P, \text{ and } J \twoheadrightarrow S.$$

Using the algorithm for ensuring dependency-preservation, we obtain the relational schema  $CJ, CD, CQ, CV, GJP, SDP$ , and  $JB$ . We can improve this schema by combining relations for which  $C$  is the key into  $CDJPQV$ . In addition, we have  $SDP$  and  $JS$  in our decomposition. Since one of these relations ( $CDJPQV$ ) is a superkey, we are done.

Comparing this decomposition with that obtained earlier in this section, we find they are quite close, with the only difference being that one of them has  $CDJPQV$  instead of  $CJP$  and  $CJDQV$ . In general, however, there could be significant differences.

## 19.7 SCHEMA REFINEMENT IN DATABASE DESIGN

We have seen how normalization can eliminate redundancy and discussed several approaches to normalizing a relation. We now consider how these ideas are applied in practice.

Database designers typically use a conceptual design methodology, such as ER design, to arrive at an initial database design. Given this, the approach of repeated decompositions to rectify instances of redundancy is likely to be the most natural use of PIs and normalization techniques.

In this section, we motivate the need for a schema refinement step following ER design. It is natural to ask whether we even need to decompose relations produced by translating an ER diagram. Should a good ER design not lead to a collection of relations free of redundancy problems? Unfortunately, ER design is a complex, subjective process, and certain constraints are not expressible in terms of ER diagrams. The examples in this section are intended to illustrate why decomposition of relations produced through ER design might be necessary.



### 19.7.1 Constraints on an Entity Set

Consider the Hourly\_Emps relation again. The constraint that attribute *ssn* is a key can be expressed as an FD:

$$\{ssn\} \rightarrow \{ssn, name, lot, rating, hourly\_wages, hours\_worked\}$$

For brevity, we write this FD as  $S \rightarrow SNLRWH$ , using a single letter to denote each attribute and omitting the set braces, but the reader should remember that both sides of an FD contain sets of attributes. In addition, the constraint that the *hourly\_wages* attribute is determined by the *rating* attribute is an FD:  $R \rightarrow W$ .

As we saw in Section 19.1.1, this FD led to redundant storage of rating wage associations. *It cannot be expressed in terms of the ER model. Only FDs that determine all attributes of a relation (i.e., key constraints) can be expressed in the ER model.* Therefore, we could not detect it when we considered Hourly\_Emps as an entity set during ER modeling.

We could argue that the problem with the original design was an artifact of a poor ER design, which could have been avoided by introducing an entity set called Wage\_Table (with attributes *rating* and *hourly\_wages*) and a relationship set Has\_Wages associating Hourly\_Emps and Wage\_Table. The point, however, is that we could easily arrive at the original design given the subjective nature of ER modeling. Having formal techniques to identify the problem with this design and guide us to a better design is very useful. The value of such techniques cannot be underestimated when designing large schemas—schemas with more than a hundred tables are not uncommon.

### 19.7.2 Constraints on a Relationship Set

The previous example illustrated how FDs can help to refine the subjective decisions made during ER design, but one could argue that the best possible ER diagram would have led to the same final set of relations. Our next example shows how functional information can lead to a set of relations unlikely to be arrived at solely through ER design.

We revisit an example from Chapter 2. Suppose that we have entity sets Parts, Suppliers, and Departments, as well as a relationship set Contracts that involves all of them. We refer to the schema for Contracts as *CQPSD*. A contract with contract id *C* specifies that a supplier *S* will supply some quantity *Q* of a part *P* to a department *J*. (We have added the contract id field *C* to the version of the Contracts relation discussed in Chapter 2.)

We might have a policy that a department purchases at most one part from any given supplier. Therefore, if there are several contracts between the same supplier and department, we know that the same part must be involved in all of them. This constraint is an FD,  $DS \rightarrow P$ .

Again we have redundancy and its associated problems. We can address this situation by decomposing Contracts into two relations with attributes *CQSD* and *3DP*. Intuitively, the relation *3DP* records the part supplied to a department by a supplier, and the relation *CQSD* records additional information about a contract. It is unlikely that we would arrive at such a design solely through ER modeling, since it is hard to formulate an entity or relationship that corresponds naturally to *CQSD*.

### 19.7.3 Identifying Attributes of Entities

This example illustrates how a careful examination of FDs can lead to a better understanding of the entities and relationships underlying the relational tables; in particular, it shows that attributes can easily be associated with the ‘wrong’ entity set during ER design. The ER diagram in Figure 19.11 shows a relationship set called *Works\_In* that is similar to the *Works\_In* relationship set of Chapter 2 but with an additional key constraint indicating that an employee can work in at most one department. (Observe the arrow connecting *Employees* to *Works\_In*.)

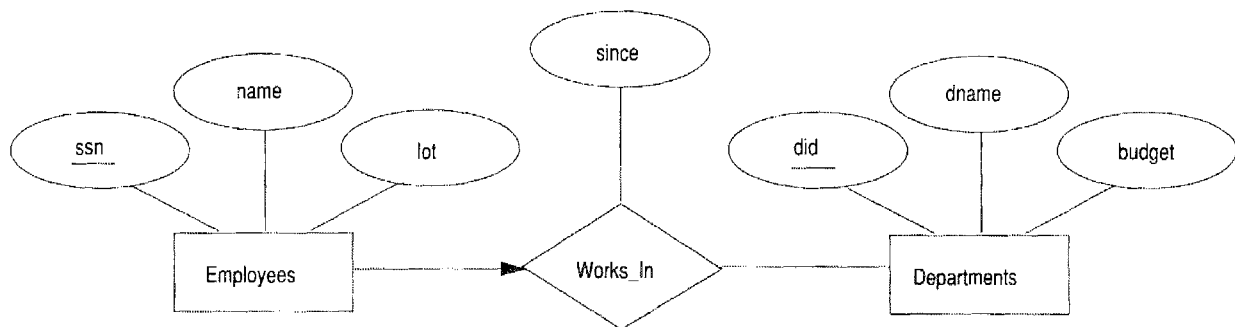


Figure 19.11. The *Works\_In* Relationship Set

Using the key constraint, we can translate this ER diagram into two relations:

*Workers*(*ssn*, *name*, *lot*, *d'id*, *since*)  
*Departments*(*did*, *dname*, *budget*)

The entity set *Employees* and the relationship set *Works\_In* are mapped to a single relation, *vWorkers*. This translation is based on the second approach discussed in Section 2.4.1.

Now suppose employees are assigned parking lots based on their department, and that all employees in a given department are assigned to the same lot. This constraint is not expressible with respect to the ER diagram of Figure 19.11. It is another example of an FD:  $did \rightarrow lot$ . The redundancy in this design can be eliminated by decomposing the Workers relation into two relations:

```
vWorkers2(ssn, name, did, since)
Dept_Lots(did, lot)
```

The new design has much to recommend it. We can change the lots associated with a department by updating a single tuple in the second relation (i.e., no update anomalies). We can associate a lot with a department even if it currently has no employees, without using *null* values (i.e., no deletion anomalies). We can add an employee to a department by inserting a tuple to the first relation even if there is no lot associated with the employee's department (i.e., no insertion anomalies).

Examining the two relations Departments and Dept\_Lots, which have the same key, we realize that a Departments tuple and a Dept\_Lots tuple with the same key value describe the same entity. This observation is reflected in the ER diagram shown in Figure 19.12.

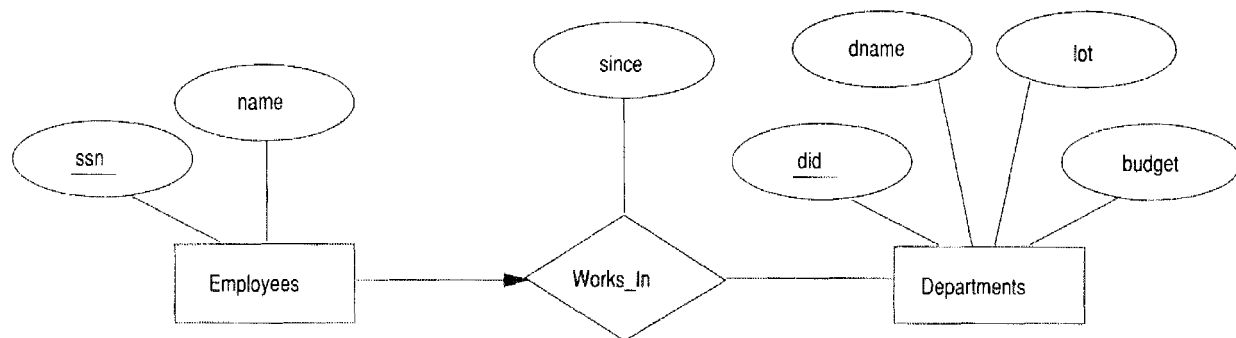


Figure 19.12 Refined Works\_In Relationship Set

Translating this diagram into the relational model would yield:

```
Workers2(ssn, name, did, since)
Departments(did, dname, budget, lot)
```

It seems intuitive to associate lots with employees; on the other hand, the lessons reveal that, in this example lots are really associated with departments. The subjective process of ER modeling could miss this point. The rigorous process of normalization would not.

### 19.7.4 Identifying Entity Sets

Consider a variant of the Reserves schema used in earlier chapters. Let Reserves contain attributes  $S$ ,  $B$ , and  $D$  as before, indicating that sailor  $S$  has a reservation for boat  $B$  on day  $D$ . In addition, let there be an attribute  $C$  denoting the credit card to which the reservation is charged. We use this example to illustrate how FD information can be used to refine an ER design. In particular, we discuss how FD information can help decide whether a concept should be modeled as an entity or as an attribute.

Suppose every sailor uses a unique credit card for reservations. This constraint is expressed by the FD  $S \rightarrow C$ . This constraint indicates that, in relation Reserves, we store the credit card number for a sailor as often as we have reservations for that sailor, and we have redundancy and potential update anomalies. A solution is to decompose Reserves into two relations with attributes  $SBD$  and  $SC$ . Intuitively, one holds information about reservations, and the other holds information about credit cards.

It is instructive to think about an ER design that would lead to these relations. One approach is to introduce an entity set called Credit\_Cards, with the sole attribute *cardno*, and a relationship set Has\_Card associating Sailors and Credit\_Cards. By noting that each credit card belongs to a single sailor, we can map Has\_Card and Credit\_Cards to a single relation with attributes  $SC$ . We would probably not model credit card numbers as entities if our main interest in card numbers is to indicate how a reservation is to be paid for; it suffices to use an attribute to model card numbers in this situation.

A second approach is to make *cardno* an attribute of Sailors. But this approach is not very natural—a sailor may have several cards, and we are not interested in all of them. Our interest is in the one card that is used to pay for reservations, which is best modeled as an attribute of the relationship Reserves.

A helpful way to think about the design problem in this example is that we first make *cardno* an attribute of Reserves and then refine the resulting tables by taking into account the FD information. (Whether we refine the design by adding *cardno* to the table obtained from Sailors or by creating a new table with attributes  $SC$  is a separate issue.)

## 19.8 OTHER KINDS OF DEPENDENCIES

FDs are probably the most common and important kind of constraint from the point of view of database design. However, there are several other kinds of dependencies. In particular, there is a well-developed theory for database

design using *multivalued dependencies* and *join dependencies*. By taking such dependencies into account, we can identify potential redundancy problems that cannot be detected using FDs alone.

This section illustrates the kinds of redundancy that can be detected using multivalued dependencies. Our main observation, however, is that simple guidelines (which can be checked using only FD reasoning) can tell us whether we even need to worry about complex constraints such as multivalued and join dependencies. We also comment on the role of *inclusion dependencies* in database design.

### 19.8.1 Multivalued Dependencies

Suppose that we have a relation with attributes *course*, *teacher*, and *book*, which we denote as *CTB*. The meaning of a tuple is that teacher *T* can teach course *C*, and book *B* is a recommended text for the course. There are no FDs; the key is *CTB*. However, the recommended texts for a course are independent of the instructor. The instance shown in Figure 19.13 illustrates this situation.

<i>course</i>	<i>teacher</i>	<i>book</i>
Physics101	Green	Mechanics
Physics101	Green	Optics
Physics101	Brown	Mechanics
Physics101	Brown	Optics
Math301	Green	Mechanics
Math301	Green	Vectors
Math301	Green	Geometry

Figure 19.13 BCNF Relation with Redundancy That Is Revealed by MVDs

Note three points here:

- The relation schema *CTB* is in BCNF; therefore we would not consider decomposing it further if we looked only at the FDs that hold over *(CTB)*.
- There is redundancy. The fact that Green can teach Physics101 is recorded once per recommended text for the course. Similarly, the fact that Optics is a text for Physics101 is recorded once per potential teacher.
- The redundancy can be eliminated by decomposing *CTB* into *CT* and *CE*.

The redundancy in this example is due to the constraint that the texts for a course are independent of the instructors, which cannot be expressed in terms

of FDs. This constraint is an example of a *multivalued dependency*, or MVD. Ideally, we should model this situation using two binary relationship sets, Instructors with attributes *CT* and Text with attributes *CB*. Because these are two essentially independent relationships, modeling them with a single ternary relationship set with attributes *CTE* is inappropriate. (See Section 2.5.3 for a further discussion of ternary versus binary relationships.) Given the subjectivity of ER design, however, we might create a ternary relationship. A careful analysis of the MVD information would then reveal the problem.

Let  $R$  be a relation schema and let  $X$  and  $Y$  be subsets of the attributes of  $R$ . Intuitively, the multivalued dependency  $X \twoheadrightarrow Y$  is said to hold over  $R$  if, in every legal instance  $r$  of  $R$ , each  $X$  value is associated with a set of  $Y$  values and this set is independent of the values in the other attributes.

Finally, if the MVD  $X \twoheadrightarrow Y$  holds over  $R$  and  $Z = R - XY$ , the following must be true for every legal instance  $r$  of  $R$ :

If  $t_1 \in r$ ,  $t_2 \in r$  and  $t_1.X = t_2.X$ , then there must be some  $t_3 \in r$  such that  $t_1.Y = t_3.Y$  and  $t_2.Z = t_3.Z$ ,

Figure 19.14 illustrates this definition. If we are given the first two tuples and told that the MVD  $X \twoheadrightarrow Y$  holds over this relation, we can infer that the relation instance must also contain the third tuple. Indeed, by interchanging the roles of the first two tuples—treating the first tuple as  $t_2$  and the second tuple as  $t_1$ —we can deduce that the tuple  $t_4$  must also be in the relation instance.

$X$	$Y$	$Z$	
$a$	$b_1$	$c_1$	tuple $t_1$
$a$	$b_2$	$c_2$	tuple $t_2$
$a$	$b_1$	$c_2$	tuple $t_3$
$a$	$b_2$	$c_1$	tuple $t_4$

Figure 19.14 Illustration of MVD Definition

This table suggests another way to think about MVDs: If  $X \twoheadrightarrow Y$  holds over  $R$ , then  $\pi_{YZ}(\sigma_{X=x}(R)) = \pi_Y(\sigma_{X=x}(R)) \times \pi_Z(\sigma_{X=x}(R))$  in every legal instance of  $R$ , for any value  $x$  that appears in the  $X$  column of  $R$ . In other words, consider groups of tuples in  $R$  with the same  $X$ -value. In each such group consider the projection onto the attributes  $YZ$ . This projection must be equal to the cross-product of the projections onto  $Y$  and  $Z$ . That is, for a given  $X$ -value, the  $Y$ -values and  $Z$ -values are independent. (From this definition it is easy to see that  $X \twoheadrightarrow Y$  will hold whenever  $X \twoheadrightarrow Z$  holds. If the FD  $X \rightarrow Y$

$Y$  holds, there is exactly one  $Y$ -value for a given  $X$ -value, and the conditions in the MVD definition hold trivially. The converse does not hold, as Figure 19.14 illustrates.)

Returning to our *CTB* example, the constraint that course texts are independent of instructors can be expressed as  $C \twoheadrightarrow T$ . In terms of the definition of MVDs, this constraint can be read as follows:

If (there is a tuple showing that)  $C$  is taught by teacher  $T$ ,  
 and (there is a tuple showing that)  $C$  has book  $B$  as text,  
 then (there is a tuple showing that)  $C$  is taught by  $T$  and has text  $B$ .

Given a set of FDs and MVDs, in general, we can infer that several additional FDs and MVDs hold. A sound and complete set of inference rules consists of the three Armstrong Axioms plus five additional rules. Three of the additional rules involve only MVDs:

- **MVD Complementation:** If  $X \twoheadrightarrow Y$ , then  $X \twoheadrightarrow R - XY$ .
- **MVD Augmentation:** If  $X \twoheadrightarrow Y$  and  $W \supseteq Z$ , then  $WX \twoheadrightarrow YZ$ .
- **MVD Transitivity:** If  $X \twoheadrightarrow Y$  and  $Y \twoheadrightarrow Z$ , then  $X \twoheadrightarrow (Z - Y)$ .

As an example of the use of these rules, since we have  $C \twoheadrightarrow T$  over *CTB*, MVD complementation allows us to infer that  $C \twoheadrightarrow CTB - CT$  as well, that is,  $C \twoheadrightarrow B$ . The remaining two rules relate FDs and MVDs:

- **Replication:** If  $X \rightarrow Y$ , then  $X \twoheadrightarrow Y$ .
- **Coalescence:** If  $X \twoheadrightarrow Y$  and there is a  $W$  such that  $W \cap Y$  is empty,  $W \rightarrow Z$ , and  $Y \supseteq Z$ , then  $X \rightarrow Z$ .

(Observe that replication states that every FD is also an MVD.)

## 19.8.2 Fourth Normal Form

Fourth Normal Form is a direct generalization of BCNF. Let  $R$  be a relation schema,  $X$  and  $Y$  be nonempty subsets of the attributes of  $R$ , and  $F$  be a set of dependencies that includes both FDs and MVDs.  $R$  is said to be in fourth normal form (4NF), if, for every nontrivial  $X \twoheadrightarrow Y$  that holds over  $R$ , one of the following statements is true:

- $Y \subseteq X$  or  $XY = R$ , or
- $X$  is a superkey.

In reading this definition, it is important to understand that the definition of a *key* has not changed.....the key must uniquely determine all attributes through FDs alone.  $X \twoheadrightarrow Y$  is a trivial MVD if  $Y \subseteq X \subseteq R$  or  $XY = R$ ; such MVDs always hold.

The relation  $CTB$  is not in 4NF because  $C \twoheadrightarrow T$  is a nontrivial MVD and  $C$  is not a key. We can eliminate the resulting redundancy by decomposing  $CTB$  into  $CT$  and  $CB$ ; each of these relations is then in 4NF.

To use MVD information fully, we must understand the theory of MVDs. However, the following result due to Date and Fagin identifies conditions—detected using only FD information!—under which we can safely ignore MVD information. That is, using MVD information in addition to the FD information will not reveal any redundancy. Therefore, if these conditions hold, we do not even need to identify all MVDs.

If a relation schema is in BCNF, and at least one of its keys consists of a single attribute, it is also in 4NF.

An important assumption is implicit in any application of the preceding result: *The set of FDs identified thus far is 'indeed the set of all FDs that hold over the relation.* This assumption is important because the result relies on the relation being in BCNF, which in turn depends on the set of FDs that hold over the relation.

We illustrate this point using an example. Consider a relation schema  $ABCD$  and suppose that the FD  $A \rightarrow BCD$  and the MVD  $B \twoheadrightarrow C$  are given. Considering only these dependencies, this relation schema appears to be a counterexample to the result. The relation has a simple key, appears to be in BCNF, and yet is not in 4NF because  $B \twoheadrightarrow C$  causes a violation of the 4NF conditions. Let us take a closer look.

$B$	$C$	$A$	$D$	
$b$	$c_1$	$a_1$	$d_1$	tuple $t_1$
$b$	$c_2$	$a_2$	$d_2$	tuple $t_2$
$b$	$c_1$	$a_2$	$d_2$	tuple $t_3$

Figure 19.15 Three Tuples from a Legal Instance of  $ABCD$

Figure 19.15 shows three tuples from an instance of  $ABCD$  that satisfies the given MVD  $B \twoheadrightarrow C$ . From the definition of an MVD, given tuples  $t_1$  and  $t_2$ , it follows that tuple  $t_3$  must also be included in the instance. Consider tuples  $t_2$  and  $t_3$ . From the given FD  $A \rightarrow BCD$  and the fact that these tuples have the



same  $A$ -value, we can deduce that  $c_1 = c_2$ . Therefore, we see that the FD  $B \rightarrow C$  must hold over  $ABCD$  whenever the FD  $A \rightarrow BCD$  and the MVD  $B \twoheadrightarrow C$  hold. If  $B \rightarrow C$  holds, the relation  $ABeD$  is not in **BCNF** (unless additional FDs make  $B$  a key)!

Thus, the apparent counterexample is really not a counterexample.....rather, it illustrates the importance of correctly identifying all FDs that hold over a relation. In this example,  $A \twoheadrightarrow BCD$  is not the only FD; the FD  $B \rightarrow C$  also holds but was not identified initially. Given a set of FDs and MVDs, the inference rules can be used to infer additional FDs (and MVDs); to apply the Date-Fagin result without first using the MVD inference rules, we must be certain that we have identified all the FDs.

In summary, the Date-Fagin result offers a convenient way to check that a relation is in 4NF (without reasoning about MVDs) if we are confident that we have identified all FDs. At this point, the reader is invited to go over the examples we have discussed in this chapter and see if there is a relation that is not in 4NF.

### 19.8.3 Join Dependencies

A join dependency is a further generalization of MVDs. A join dependency (JD)  $\bowtie \{R_1, \dots, R_n\}$  is said to hold over a relation  $R$  if  $R_1, \dots, R_n$  is a lossless-join decomposition of  $R$ .

An MVD  $X \twoheadrightarrow Y$  over a relation  $R$  can be expressed as the join dependency  $\bowtie \{XV, X(R, -Y)\}$ . As an example, in the *GTB* relation, the MVD  $C \twoheadrightarrow T$  can be expressed as the join dependency  $\bowtie \{CT, CB\}$ .

Unlike FDs and MVDs, there is no set of sound and complete inference rules for JDs.

### 19.8.4 Fifth Normal Form

A relation schema  $R$  is said to be in fifth normal form (5NF) if, for every JD  $\bowtie \{R_1, \dots, R_n\}$  that holds over  $R$ , one of the following statements is true:

- $R_i = R$ , for some  $i$ , or
- The JD is implied by the set of those FDs over  $R$  in which the left side is a key for  $R$ .

The second condition deserves some explanation, since we have not presented inference rules for FDs and JDs taken together. Intuitively, we must be able to show that the decomposition of  $R$  into  $\{R_1, \dots, R_n\}$  is lossless-join whenever the key dependencies (FDs in which the left side is a key for  $R$ ) hold. (II)  $\bowtie \{R_1, \dots, R_n\}$  is a trivial JD if  $R_i = R$  for some  $i$ ; such a JD always holds.

The following result, also due to Date and Fagin, identifies conditions—again, detected using only FD information—under which we can safely ignore JD information:

If a relation schema is in 3NF and each of its keys consists of a single attribute, it is also in 5NF.

The conditions identified in this result are sufficient for a relation to be in 5NF but not necessary. The result can be very useful in practice because it allows us to conclude that a relation is in 5NF *without ever identifying the MVDs and JDs that may hold over the relation*.

## 19.8.5 Inclusion Dependencies

MVDs and JDs can be used to guide database design, as we have seen, although they are less common than FDs and harder to recognize and reason about. In contrast, inclusion dependencies are very intuitive and quite common. However, they typically have little influence on database design (beyond the ER design stage).

Informally, an inclusion dependency is a statement of the form that some columns of a relation are contained in other columns (usually of a second relation). A foreign key constraint is an example of an inclusion dependency; the referring column(s) in one relation must be contained in the primary key column(s) of the referenced relation. As another example, if  $R$  and  $S$  are two relations obtained by translating two entity sets that every  $R$  entity is also an  $S$  entity, we would have an inclusion dependency; projecting  $R$  on its key attributes yields a relation contained in the relation obtained by projecting  $S$  on its key attributes.

The main point to bear in mind is that we should not split groups of attributes that participate in an inclusion dependency. For example, if we have an inclusion dependency  $AB \subseteq CD$ , while decomposing the relation schema containing  $AB$ , we should ensure that at least one of the schemas obtained in the decomposition contains both  $A$  and  $B$ . Otherwise, we cannot check the inclusion dependency  $AB \subseteq CD$  without reconstructing the relation containing  $AB$ .

Most inclusion dependencies in practice are *key-based*, that is, involve only keys. Foreign key constraints are a good example of key-based inclusion dependencies. An ER diagram that involves ISA hierarchies (see Section 2.4.4) also leads to key-based inclusion dependencies. If all inclusion dependencies are key-based, we rarely have to worry about splitting attribute groups that participate in inclusion dependencies, since decompositions usually do not split the primary key. Note, however, that going from 3NF to BCNF always involves splitting some key (ideally not the primary key!), since the dependency guiding the split is of the form  $X \rightarrow A$  where  $A$  is part of a key.

## 19.9 CASE STUDY: THE INTERNET SHOP

Recall from Section 3.8 that DBDudes settled on the following schema:

Books(isbn: CHAR(10), title: CHAR(8), author: CHAR(80),  
 qty\_in\_stock: INTEGER, price: REAL, year\_published: INTEGER)  
 Customers(cid: INTEGER, name: CHAR(80), address: CHAR(200))  
 Orders(ordernum: INTEGER, isbn: CHAR(10), cid: INTEGER,  
 cardnum: CHAR(16), qty: INTEGER, order\_date: DATE, ship\_date: DATE)

DBDudes analyzes the set of relations for possible redundancy. The Books relation has only one key, (*isbn*), and no other functional dependencies hold over the table. Thus, Books is in BCNF. The Customers relation also has only one key, (*cid*), and no other functional dependencies hold over the table. Thus, Customers is also in BCNF.

DBDudes has already identified the pair  $\langle \text{ordernum}, \text{isbn} \rangle$  as the key for the Orders table. In addition, since each order is placed by one customer on one specific date with one specific credit card number, the following three functional dependencies hold:

$$\text{ordernum} \rightarrow \text{cid}, \text{ordernum} \rightarrow \text{order\_date}, \text{ and } \text{ordernum} \rightarrow \text{cardnum}$$

The experts at DBDudes conclude that Orders is not even in 3NF. (Can you see why?) They decide to decompose Orders into the following two relations:

Orders(ordernum, cid, order\_date, cardnum, and  
 Orderlists(ordernum, isbn, qty, ship\_date)

The resulting two relations, Orders and Orderlists, are both in BCNF, and the decomposition is lossless-join since *ordernum* is a key for (the new) Orders. The reader is invited to check that this decomposition is also dependency-preserving. For completeness, we give the SQL DDL for the Orders and Orderlists relations below:

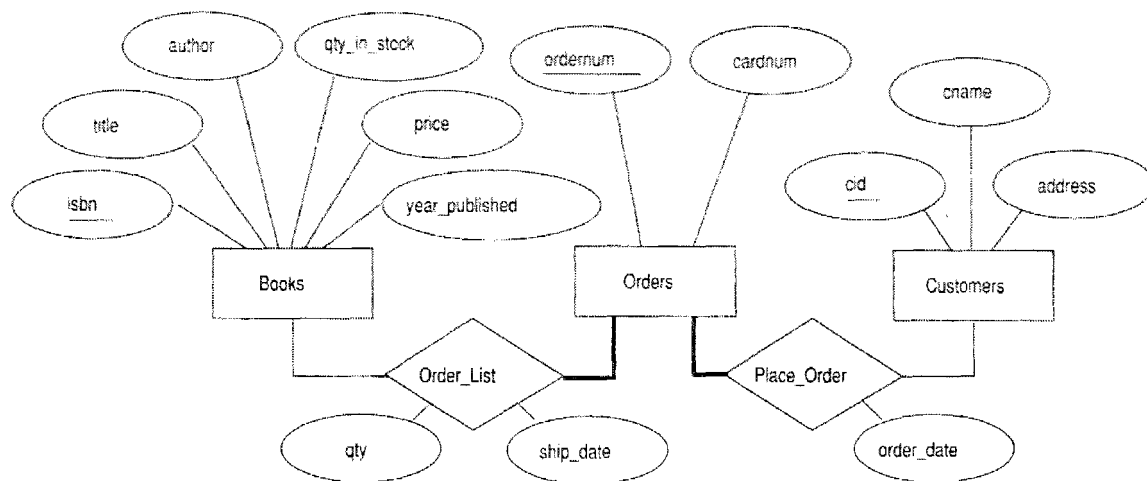


Figure 19.16 ER Diagram Reflecting the Final Design

```
CREATE TABLE Orders ( ordernum    INTEGER,
                      cid         INTEGER,
                      order_date  DATE,
                      cardnum     CHAR(16),
                      PRIMARY KEY (ordernum),
                      FOREIGN KEY (cid) REFERENCES Customers )
```

```
CREATE TABLE Orderlists ( ordernum    INTEGER,
                          isbn         CHAR(10),
                          qty          INTEGER,
                          ship_date   DATE,
                          PRIMARY KEY (ordernum, isbn),
                          FOREIGN KEY (isbn) REFERENCES Books)
```

Figure 19.16 shows an updated ER diagram that reflects the new design. Note that DBDudes could have arrived immediately at this diagram if they had made Orders an entity set instead of a relationship set right at the beginning. But at that time they did not understand the requirements completely, and it seemed natural to model Orders as a relationship set. This iterative refinement process is typical of real-life database design processes. As DBDudes has learned over time, it is rare to achieve an initial design that is not changed as a project progresses.

The DBDudes team celebrates the successful completion of logical database design and schema refinement by opening a bottle of champagne and charging it to B&N. After recovering from the celebration, they move on to the physical design phase.

## 19.10 REVIEW QUESTIONS

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

- Illustrate redundancy and the problems that it can cause. Give examples of *insert*, *delete*, and *update* anomalies. Can *null* values help address these problems? Are they a complete solution? (Section 19.1.1)
- What is a *decomposition* and how does it address redundancy? What problems may be caused by the use of decompositions? (Sections 19.1.2 and 19.1.3)
- Define *functional dependencies*. How are *primary keys* related to FDs? (Section 19.2)
- When is an FD *implied* by a set  $F$  of FDs? Define *Armstrong's Axioms*, and explain the statement that "they are a sound and complete set of rules for FD inference." (Section 19.3)
- What is the *dependency closure*  $F^+$  of a set  $F$  of FDs? What is the *attribute closure*  $X^+$  of a set of attributes  $X$  with respect to a set of FDs  $F$ ? (Section 19.3)
- Define 1NF, 2NF, 3NF, and BCNF. What is the motivation for putting a relation in BCNF? What is the motivation for 3NF? (Section 19.4)
- When is the decomposition of a relation schema  $R$  into two relation schemas  $X$  and  $Y$  said to be a *lossless-join* decomposition? Why is this property so important? Give a necessary and sufficient condition to test whether a decomposition is lossless-join. (Section 19.5.1)
- When is a decomposition said to be *dependency-preserving*? Why is this property useful? (Section 19.5.2)
- Describe how we can obtain a lossless-join decomposition of a relation into BCNF. Give an example to show that there may not be a dependency-preserving decomposition into BCNF. Illustrate how a given relation could be decomposed in different ways to arrive at several alternative decompositions, and discuss the implications for database design. (Section 19.6.1)
- Give an example that illustrates how a collection of relations in BCNF could have redundancy even though each relation, by itself, is free from redundancy. (Section 19.6.1)
- What is a *minimal cover* for a set of FDs? Describe an algorithm for computing the minimal cover of a set of FDs, and illustrate it with an example. (Section 19.6.2)

- Describe how the algorithm for lossless-join decomposition into BCNF can be adapted to obtain a lossless-join, dependency-preserving decomposition into 3NF. Describe the alternative *synthesis* approach to obtaining such a decomposition into 3NF. Illustrate both approaches using an example. (Section 19.6.2)
- Discuss how schema refinement through dependency analysis and normalization can improve schemas obtained through ER design. (Section 19.7)
- Define *multivalued dependencies*, *join dependencies*, and *inclusion dependencies*. Discuss the use of such dependencies for database design. Define 4NF and 5NF, and explain how they prevent certain kinds of redundancy that BCNF does not eliminate. Describe tests for 4NF and 5NF that use only FDs. What key assumption is involved in these tests? (Section 19.8)

## EXERCISES

Exercise 19.1 Briefly answer the following questions:

1. Define the term *functional dependency*.
2. Why are some functional dependencies called *trivial*?
3. Give a set of FDs for the relation schema  $R(A, B, C, D)$  with primary key  $AB$  under which  $R$  is in 1NF but not in 2NF.
4. Give a set of FDs for the relation schema  $R(A, B, C, D)$  with primary key  $AB$  under which  $R$  is in 2NF but not in 3NF.
5. Consider the relation schema  $R(A, B, C)$ , which has the FD  $B \rightarrow C$ . If  $A$  is a candidate key for  $R$ , is it possible for  $R$  to be in BCNF? If so, under what conditions? If not, explain why not.
6. Suppose we have a relation schema  $R(A, B, C)$  representing a relationship between two entity sets with keys  $A$  and  $B$ , respectively, and suppose that  $R$  has (among others) the FDs  $A \twoheadrightarrow B$  and  $B \rightarrow A$ . Explain what such a pair of dependencies means (i.e., what they imply about the relationship that the relation models).

Exercise 19.2 Consider a relation  $R$  with five attributes  $ABCDE$ . You are given the following dependencies:  $A \rightarrow B$ ,  $BE \rightarrow E$ , and  $ED \rightarrow A$ .

1. List all keys for  $R$ .
2. Is  $R$  in 3NF?
3. Is  $R$  in BCNF?

Exercise 19.3 Consider the relation shown in Figure 19.17.

1. List all the functional dependencies that this relation instance satisfies.
2. Assume that the value of attribute  $Z$  of the last record in the relation is changed from  $z_3$  to  $z_2$ . Now list all the functional dependencies that this relation instance satisfies.

Exercise 19.4 Assume that you are given a relation with attributes  $ABCD$ .

$X$	$Y$	$Z$
$x_1$	$y_1$	$z_1$
$x_1$	$y_1$	$z_2$
$x_2$	$y_1$	$z_1$
$x_2$	$y_1$	$z_3$

Figure 19.17 Relation for Exercise 19.3.

1. Assume that no record has NULL values. Write an SQL query that checks whether the functional dependency  $A \rightarrow B$  holds.
2. Assume again that no record has NULL values. Write an SQL assertion that enforces the functional dependency  $A \rightarrow B$ .
3. Let us now assume that records could have NULL values. Repeat the previous two questions under this assumption.

**Exercise 19.5** Consider the following collection of relations and dependencies. Assume that each relation is obtained through decomposition from a relation with attributes  $ABCDEFGHI$  and that all the known dependencies over relation  $ABCDEFGHI$  are listed for each question. (The questions are independent of each other, obviously, since the given dependencies over  $ABCDEFGHI$  are different.) For each (sub)relation: (a) State the strongest normal form that the relation is in. (b) If it is not in BCNF, decompose it into a collection of BCNF relations.

1.  $R_1(A, C, B, D, E)$ ,  $A \rightarrow B$ ,  $C \rightarrow D$
2.  $R_2(A, B, F)$ ,  $AC \rightarrow B$ ,  $B \rightarrow F$
3.  $R_3(A, D)$ ,  $D \rightarrow G$ ,  $G \rightarrow H$
4.  $R_4(D, C, H, G)$ ,  $A \rightarrow I$ ,  $I \rightarrow A$
5.  $R_5(A, I, C, B)$

**Exercise 19.6** Suppose that we have the following three tuples in a legal instance of a relation schema  $S$  with three attributes  $ABC$  (listed in order):  $(1,2,3)$ ,  $(4,2,3)$ , and  $(5,3,3)$ .

1. Which of the following dependencies can you infer does *not* hold over schema  $S$ ?  
(a)  $A \rightarrow B$ , (b)  $BC \rightarrow A$ , (c)  $B \rightarrow C$
2. Can you identify any dependencies that hold over  $S$ ?

**Exercise 19.7** Suppose you are given a relation  $R$  with four attributes  $ABCD$ . For each of the following sets of FDs, assuming those are the only dependencies that hold for  $R$ , do the following: (a) Identify the candidate key(s) for  $R$ . (b) Identify the best normal form that  $R$  satisfies (1NF, 2NF, 3NF, or BCNF). (c) If  $R$  is not in BCNF, decompose it into a set of BCNF relations that preserve the dependencies.

1.  $C \rightarrow D$ ,  $C \rightarrow A$ ,  $B \rightarrow C$
2.  $B \rightarrow C$ ,  $D \rightarrow A$
3.  $ABC \rightarrow D$ ,  $D \rightarrow A$
4.  $A \rightarrow B$ ,  $BC \rightarrow D$ ,  $A \rightarrow C$
5.  $AB \rightarrow C$ ,  $AB \rightarrow D$ ,  $C \rightarrow A$ ,  $D \rightarrow B$

**Exercise 19.8** Consider the attribute set  $R = ABCDEGH$  and the FD set  $F = \{AB \rightarrow C, AC \rightarrow B, AD \rightarrow E, B \rightarrow D, Be \rightarrow A, B \rightarrow G\}$ .

- For each of the following attribute sets, do the following: (i) Compute the set of dependencies that hold over the set and write down a minimal cover. (ii) Name the strongest normal form that is not violated by the relation containing these attributes. (iii) Decompose it into a collection of BCNF relations if it is not in BCNF.

(a)  $ABC$ , (b)  $ABCD$ , (c)  $ABCEG$ , (d)  $DC:BGH$ , (e)  $ACEH$

- Which of the following decompositions of  $R = ABCDEG$ , with the same set of dependencies  $F$ , is (a) dependency-preserving? (b) lossless-join?

(a)  $\{AB, BC, ABDE, EG\}$

(b)  $\{ABC, ACDE, ADG\}$

**Exercise 19.9** Let  $R$  be decomposed into  $R_1, R_2, \dots, R_n$ . Let  $F$  be a set of FDs on  $R$ .

- Define what it means for  $F$  to be preserved in the set of decomposed relations.
- Describe a polynomial-time algorithm to test dependency-preservation.
- Projecting the FDs stated over a set of attributes  $X$  onto a subset of attributes  $Y$  requires that we consider the closure of the FDs. Give an example where considering the closure is important in testing dependency-preservation, that is, considering just the given FDs gives incorrect results.

**Exercise 19.10** Suppose you are given a relation  $R(A,B,C,D)$ . For each of the following sets of FDs, assuming they are the only dependencies that hold for  $R$ , do the following: (a) Identify the candidate key(s) for  $R$ . (b) State whether or not the proposed decomposition of  $R$  into smaller relations is a good decomposition and briefly explain why or why not.

- $B \rightarrow C, D \rightarrow A$ ; decompose into  $BC$  and  $AD$ .
- $AB \rightarrow C, C \rightarrow A, C \rightarrow D$ ; decompose into  $ACD$  and  $Be$ .
- $A \rightarrow BC, C \rightarrow AD$ ; decompose into  $ABC$  and  $AD$ .
- $A \rightarrow B, B \rightarrow C, C \rightarrow D$ ; decompose into  $AB$  and  $ACD$ .
- $A \rightarrow B, B \rightarrow C, C \rightarrow D$ ; decompose into  $AB, AD$  and  $CD$ .

**Exercise 19.11** Consider a relation  $R$  that has three attributes  $ABC$ . It is decomposed into relations  $R_1$  with attributes  $AB$  and  $R_2$  with attributes  $Be$ .

- State the definition of a lossless-join decomposition with respect to this example. Answer this question concisely by writing a relational algebra equation involving  $R, R_1$ , and  $R_2$ .
- Suppose that  $B \twoheadrightarrow C$ . Is the decomposition of  $R$  into  $R_1$  and  $R_2$  lossless-join? Reconcile your answer with the observation that neither of the FDs  $H_1 \rightarrow R_1$  nor  $R_1 \rightarrow R_2$  hold, in light of the simple test offering a necessary and sufficient condition for lossless-join decomposition into two relations in Section 15.6.1.
- If you are given the following instances of  $R_1$  and  $R_2$ , what can you say about the instance of  $R$  from which these were obtained? Answer this question by listing tuples that are definitely in  $R$  and tuples that are possibly in  $R$ .

Instance of  $R_1 = \{(5,1), (6,1)\}$

Instance of  $R_2 = \{(1,8), (1,9)\}$

Can you say that attribute  $B$  definitely is or is not a key for  $R$ ?



**Exercise 19.12** Suppose that we have the following four tuples in a relation  $S$  with three attributes  $ABC$ :  $(1,2,3)$ ,  $(4,2,3)$ ,  $(5,3,3)$ ,  $(5,3,4)$ . Which of the following functional ( $\rightarrow$ ) and multivalued ( $\twoheadrightarrow$ ) dependencies can you infer does *not* hold over relation  $S$ ?

1.  $A \rightarrow B$
2.  $A \twoheadrightarrow B$
3.  $BC \rightarrow A$
4.  $BC \twoheadrightarrow A$
5.  $B \rightarrow C$
6.  $B \twoheadrightarrow C$

**Exercise 19.13** Consider a relation  $R$  with five attributes  $ABCDE$ .

1. For each of the following instances of  $R$ , state whether it violates (a) the FD  $BE \rightarrow D$  and (b) the MVD  $BE \twoheadrightarrow D$ :
  - (a)  $\{ \}$  (i.e., empty relation)
  - (b)  $\{(0,2,3,4,5), (2,0,3,5,5)\}$
  - (c)  $\{(0,2,3,4,5), (2,0,3,5,5), (0,2,3,4,6)\}$
  - (d)  $\{(0,2,3,4,5), (2,0,3,4,5), (0,2,3,6,5)\}$
  - (e)  $\{(0,2,3,4,5), (2,0,3,7,5), (0,2,3,4,6)\}$
  - (f)  $\{(0,2,3,4,5), (2,0,3,4,5), (0,2,3,6,5), (0,2,3,6,6)\}$
  - (g)  $\{(0,2,3,4,5), (0,2,3,6,5), (0,2,3,6,6), (0,2,3,4,6)\}$
2. If each instance for  $R$  listed above is legal, what can you say about the FD  $A \rightarrow B$ ?

**Exercise 19.14** JDs are motivated by the fact that sometimes a relation that cannot be decomposed into two smaller relations in a lossless-join manner can be so decomposed into three or more relations. An example is a relation with attributes *supplier*, *part*, and *project*, denoted  $SPJ$ , with no FDs or MVDs. The JD  $\{SP, PJ, JS\}$  holds.

From the JD, the set of relation schemes  $SP$ ,  $PJ$ , and  $JS$  is a lossless-join decomposition of  $SPJ$ . Construct an instance of  $SPJ$  to illustrate that no two of these schemes suffice.

**Exercise 19.15** Answer the following questions

1. Prove that the algorithm shown in Figure 19.4 correctly computes the attribute closure of the input attribute set  $X$ .
2. Describe a linear-time (in the size of the set of FDs, where the size of each FD is the number of attributes involved) algorithm for finding the attribute closure of a set of attributes with respect to a set of FDs. Prove that your algorithm correctly computes the attribute closure of the input attribute set.

**Exercise 19.16** Let us say that an FD  $X \rightarrow Y$  is *simple* if  $Y$  is a single attribute.

1. Replace the FD  $AB \rightarrow CD$  by the smallest equivalent collection of simple FDs.
2. Prove that every FD  $X \rightarrow Y$  in a set of FDs  $F$  can be replaced by a set of simple FDs such that  $F^+$  is equal to the closure of the new set of FDs.

**Exercise 19.17** Prove that Armstrong's Axioms are sound and complete for FD inference. That is, show **that** repeated application of these axioms on a set  $F$  of FDs produces exactly the dependencies in  $P^+$ .

**Exercise 19.18** Consider a relation  $R$  with attributes  $ABCDE$ . Let the following FDs be given:  $A \rightarrow BC$ ,  $BE \rightarrow E$ , and  $E \rightarrow DA$ . Similarly, let  $S$  be a relation with attributes  $ABCDE$  and let the following FDs be given:  $A \rightarrow BC$ ,  $B \rightarrow E$ , and  $E \rightarrow DA$ . (Only the second dependency differs from those that hold over  $R$ .) You do not know whether or which other (join) dependencies hold.

1. Is  $R$  in BCNF?
2. Is  $R$  in 4NF?
3. Is  $R$  in 5NF?
4. Is  $S$  in BCNF?
5. Is  $S$  in 4NF?
6. Is  $S$  in 5NF?

**Exercise 19.19** Let  $R$  be a relation schema with a set  $F$  of FDs. Prove that the decomposition of  $R$  into  $H_1$  and  $R_2$  is lossless-join if and only if  $P^+$  contains  $H_1 \cap R_2 \rightarrow R_1$  or  $R_1 \cap R_2 \rightarrow R_2$ .

**Exercise 19.20** Consider a schema  $R$  with FDs  $F$  that is decomposed into schemas with attributes  $X$  and  $Y$ . Show that this is dependency-preserving if  $F' \subseteq (F_X \cup F_Y)^+$ .

**Exercise 19.21** Prove that the optimization of the algorithm for lossless-join, dependency-preserving decomposition into 3NF relations (Section 19.6.2) is correct.

**Exercise 19.22** Prove that the 3NF synthesis algorithm produces a lossless-join decomposition of the relation containing all the original attributes.

**Exercise 19.23** Prove that an MVD  $X \twoheadrightarrow Y$  over a relation  $R$  can be expressed as the join dependency  $\bowtie \{XY, X(R - Y)\}$ .

**Exercise 19.24** Prove that, if  $R$  has only one key, it is in BCNF if and only if it is in 3NF.

**Exercise 19.25** Prove that, if  $R$  is in 3NF and every key is simple, then  $R$  is in 4NF.

**Exercise 19.26** Prove these statements:

1. If a relation schema is in BCNF and at least one of its keys consists of a single attribute, it is also in 4NF.
2. If a relation schema is in 3NF and each key has a single attribute, it is also in 5NF.

**Exercise 19.27** Give an algorithm for testing whether a relation schema is in BCNF. The algorithm should be polynomial in the size of the set of given FDs. (The 'size' is the sum over all FDs of the number of attributes that appear in the FD.) Is there a polynomial algorithm for testing whether a relation schema is in 3NF?

**Exercise 19.28** Give an algorithm for testing whether a relation schema is in BCNF. The algorithm should be polynomial in the size of the set of given FDs. (The 'size' is the sum over all FDs of the number of attributes that appear in the FD.) Is there a polynomial algorithm for testing whether a relation schema is in 3NF?

**Exercise 19.29** Prove that the algorithm for decomposing a relation schema with a set of FDs into a collection of BCNF relation schemas as described in Section 19.6.1 is correct (i.e., it produces a collection of BCNF relations, and is lossless-join) and terminates.

## BIBLIOGRAPHIC NOTES

Textbook presentations of dependency theory and its use in database design include [3, 45, 501, 509, 747]. Good survey articles on the topic include [755, 415].

FDs were introduced in [187], along with the concept of 3NF, and axioms for inferring FDs were presented in [38]. BCNF was introduced in [188]. The concept of a legal relation instance and dependency satisfaction are studied formally in [328]. FDs were generalized to semantic data models in [768].

Finding a key is shown to be NP-complete in [497]. Lossless-join decompositions were studied in [28, 502, 627]. Dependency-preserving decompositions were studied in [74]. [81] introduced minimal covers. Decomposition into 3NF is studied by [81, 98] and decomposition into BCNF is addressed in [742]. [412] shows that testing whether a relation is in 3NF is NP-complete. [253] introduced 4NF and discussed decomposition into 4NF. Fagin introduced other normal forms in [254] (project-join normal form) and [255] (domain-key normal form). In contrast to the extensive study of vertical decompositions, there has been relatively little formal investigation of horizontal decompositions. [209] investigates horizontal decomposition.

IVDs were discovered independently by Delobel [211], Fagin [253], and Zaniolo [789]. Axioms for FDs and MVDs were presented in [73]. [593] shows that there is no axiomatization for JDs, although [662] provides an axiomatization for a more general class of dependencies. The sufficient conditions for 4NF and 5NF in terms of FDs that were discussed in Section 19.8 are from [205]. An approach to database design that uses dependency information to construct sample relation instances is described in [508, 509].



# 20

## PHYSICAL DATABASE DESIGN AND TUNING

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

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

—Architect Frank Lloyd Wright

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

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

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

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

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

## 20.1 INTRODUCTION TO PHYSICAL DATABASE DESIGN

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

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

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

### 20.1.1 Database Workloads

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

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

For each query in the workload, we must identify

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

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

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

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

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

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

## 20.1.2 Physical Design and Tuning Decisions

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

### 1. Choice of indexes to create:

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

### 2. Tuning the conceptual schema:

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

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

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

### **20.1.3 Need for Database Tuning**

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

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

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

## **20.2 GUIDELINES FOR INDEX SELECTION**

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



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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

## 20.3 BASIC EXAMPLES OF INDEX SELECTION

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

Let us begin with a simple query:

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

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

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

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

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

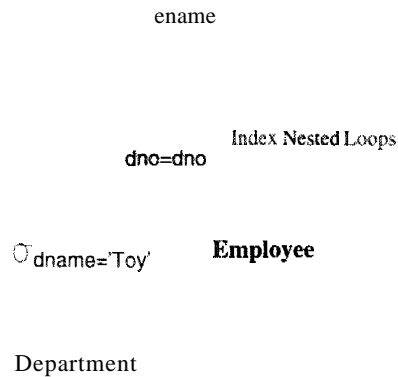


Figure 20.1 A Desirable Query Evaluation Plan

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

Our next query involves a range selection:

```

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

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

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

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

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

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

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

## 20.4 CLUSTERING AND INDEXING

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

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

We concluded that a good evaluation plan is to use an index on *dname* to retrieve Departments tuples satisfying the condition on *dname* and to find matching Employees tuples using an index on *dno*. Should these indexes be clustered? Given our assumption that the number of tuples satisfying 1). *dname*='Toy' is likely to be small, we should build an unclustered index on *dname*. On the other hand, Employees is the inner relation in an index nested loops join and *dno* is not a candidate key. This situation is a strong argument that the index on the *dno* field of Employees should be clustered. In fact, because the join consists of repeatedly posing equality selections on the *dno* field of the inner relation, this type of query is a stronger justification for making the index on *dno* clustered than a simple selection query such as the previous selection on

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

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

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

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

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

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

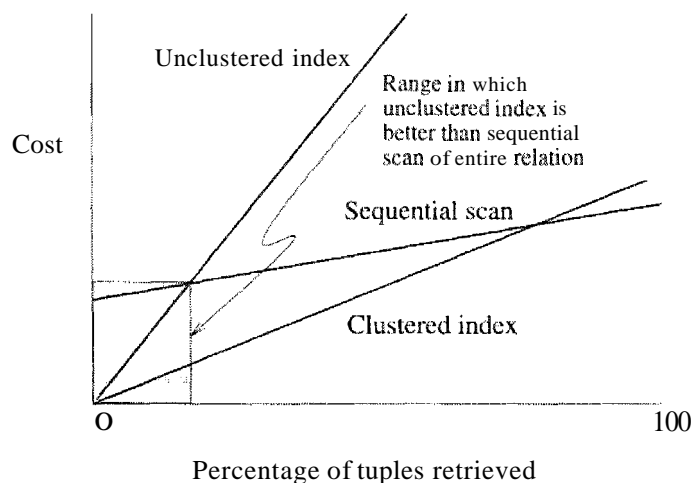


Figure 20.2 The Impact of Clustering

### 20.4.1 Co-clustering Two Relations

In our description of a typical database system architecture in Chapter 9, we explained how a relation is stored as a file of records. Although a file usually contains only the records of *some* one relation, *some* systems allow records from *more* than one relation to be stored in a single file. *Some* database user can request that the records from *two* relations be interleaved physically in this *channel*. This data layout is *sometimes* referred to as *co-clustering* the two relations. We now discuss when co-clustering can be beneficial.

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

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

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

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

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

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

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

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

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

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

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

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



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

To summarize co-clustering:

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

## 20.5 INDEXES THAT ENABLE INDEX-ONLY PLANS

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

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

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

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

The next query takes this idea a step further:

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

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

## 20.6 TOOLS TO ASSIST IN INDEX SELECTION

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

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

### 20.6.1 Automatic Index Selection

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

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

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

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

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

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

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

## 20.6.2 How Do Index Tuning Wizards Work?

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

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

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

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

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

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

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

## Candidate Index Selection

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

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

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

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

## Enumerating Index Configurations

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

## 20.7 OVERVIEW OF DATABASE TUNING

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

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

### 20.7.1 Tuning Indexes

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

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

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

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

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

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

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

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

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

## 20.7.2 Tuning the Conceptual Schema

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

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

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

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



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

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

### 20.7.3 Tuning Queries and Views

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

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

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

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

## 20.8 CHOICES IN TUNING THE CONCEPTUAL SCHEMA

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

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

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

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

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

### 20.8.1 Settling for a Weaker Normal Form

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

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

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

### 20.8.2 Denormalization

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

### 20.8.3 Choice of Decomposition

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

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

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

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

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

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

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

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

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

### 20.8.4 Vertical Partitioning of BCNF Relations

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

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

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

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

### 20.8.5 Horizontal Decomposition

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

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

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

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

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

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

## 20.9 CHOICES IN TUNING QUERIES AND VIEWS

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

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

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

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

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

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

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

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

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

This query is equivalent to

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

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

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

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

This query can be rewritten as

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

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

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

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



## 20.10 IMPACT OF CONCURRENCY

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

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

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

We now discuss techniques for achieving these goals.

### 20.10.1 Reducing Lock Durations

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

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

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

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

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

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

## 20.10.2 Reducing Hot Spots

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

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

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

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

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

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

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

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

## 20.11 CASE STUDY: THE INTERNET SHOP

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

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

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

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

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

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

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

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

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

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

### 20.11.1 Tuning the Database

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

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

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

## 20.12 DBMS BENCHMARKING

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

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

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

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

### 20.12.1 Well-Known DBMS Benchmarks

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

**Query Benchmarks:** The *Visconsin* benchmark is widely used for measuring the performance of simple relational queries. The *Set Query* benchmark measures the performance of a suite of more complex queries, and the *AS<sup>3</sup>A.P* benchmark measures the performance of a mixed workload of transactions, relational queries, and utility functions. The *TPC-I* benchmark is a suite of complex SQL queries intended to be representative of the (decision-support ap-

plication domain. The ()LAP Council also developed a benchmark for complex decision-support queries, including some queries that cannot be expressed easily in SQL; this is intended to measure systems for *online analytic processing (OLAP)*, which we discuss in Chapter 25, rather than traditional SQL systems. The Sequoia 2000 benchmark is designed to compare DBMS support for geographic information systems.

**Object-Database Benchmarks:** The OO1 and OO7 benchmarks measure the performance of object-oriented database systems. The Bucky benchmark measures the performance of object-relational database systems. (We discuss object-database systems in Chapter 23.)

## 20.12.2 Using a **Benchmark**

Benchmarks should be used with a good understanding of what they are designed to measure and the application environment in which a DBMS is to be used. When you use benchmarks to guide your choice of a DBMS, keep the following guidelines in mind:

- **How Meaningful is a Given Benchmark?** Benchmarks that try to distill performance into a single number can be overly simplistic. A DBMS is a complex piece of software used in a variety of applications. A good benchmark should have a suite of tasks that are carefully chosen to cover a particular application domain and test DBMS features important for that domain.
- **How Well Does a Benchmark Reflect Your Workload?** Consider your expected workload and compare it with the benchmark. Give little weight to the performance of those benchmark tasks (i.e., queries and updates) that are similar to important tasks in your workload. Also consider how benchmark numbers are measured. For example, elapsed time for individual queries might be misleading if considered in a multiuser setting: A system may have higher elapsed times because of slower I/O. On a multiuser workload, given sufficient disks for parallel I/O, such a system might outperform a system with a lower elapsed time.
- **Create Your Own Benchmark:** Vendors often tweak their systems in ad hoc ways to obtain good numbers on important benchmarks. To counter this, create your own benchmark by modifying standard benchmarks slightly or by replacing the tasks in a standard benchmark with similar tasks from your workload.

## 20.13 REVIEW QUESTIONS

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

- || What are the components of a workload description? (Section 20.1.1)
- What decisions need to be made during physical design? (Section 20.1.2)
- || Describe six high-level guidelines for index selection. (Section 20.2)
- When should we create clustered indexes? (Section 20.4)
- What is co-clustering, and when should we use it? (Section 20.4.1)
- || What is an index-only plan, and how do we create indexes for index-only plans? (Section 20.5)
- Why is automatic index tuning a hard problem? Give an example. (Section 20.6.1)
- Give an example of one algorithm for automatic index tuning. (Section 20.6.2)
- || Why is database tuning important? (Section 20.7)
- || How do we tune indexes, the conceptual schema, and queries and views? (Sections 20.7.1 to 20.7.3)
- What are our choices in tuning the conceptual schema? What are the following techniques and when should we apply them: settling for a weaker normal form, denormalization, and horizontal and vertical decompositions. (Section 20.8)
- What choices do we have in tuning queries and views? (Section 20.9)
- What is the impact of locking on database performance? How can we reduce lock contention and hot spots? (Section 20.10)
- Why do we have standardized database benchmarks, and what common metrics are used to evaluate database systems? Can you describe a few popular database benchmarks? (Section 20.12)

## EXERCISES

Exercise 20.1 Consider the following BCNF schema for a portion of a simple corporate database (type information is not relevant to this question and is omitted):

Emp (eid, ename, addr, sal, age, yrs, deptid)  
 Dept (did, dname, floort, budget)



Suppose you know that the following queries are the six most common queries in the workload for this corporation and that all six are roughly equivalent in frequency and importance:

- List the id, name, and address of employees in a user-specified age range.
  - List the id, name, and address of employees who work in the department with a user-specified department name.
  - List the id and address of employees with a user-specified employee name.
  - List the overall average salary for employees.
  - List the average salary for employees of each age; that is, for each age in the data, list the age and the corresponding average salary.
  - List all the department information, ordered by department floor numbers.
1. Given this information, and assuming that these queries are more important than any updates, design a physical schema for the corporate database that will give good performance for the expected workload. In particular, decide which attributes will be indexed and whether each index will be a clustered index or an unclustered index. Assume that B+ tree indexes are the only index type supported by the DBMS and that both single- and multiple-attribute keys are permitted. Specify your physical design by identifying the attributes you recommend indexing on via clustered or unclustered B+ trees.
  2. Redesign the physical schema assuming that the set of important queries is changed to be the following:
    - List the id and address of employees with a user-specified employee name.
    - List the overall maximum salary for employees.
    - List the average salary for employees by department; that is, for each *deptid* value, list the *deptid* value and the average salary of employees in that department.
    - List the sum of the budgets of all departments by floor; that is, for each floor, list the floor and the sum.
    - Assume that this workload is to be tuned with an automatic index tuning wizard. Outline the main steps in the execution of the index tuning algorithm and the set of candidate configurations that would be considered.

**Exercise 20.2** Consider the following BCNF<sup>1</sup> relational schema for a portion of a university database (type information is not relevant to this question and is omitted):

*Prof*(*ssno*, *pname*, *office*, *age*, *sex*, *specialty*, *dept\_id*)  
*Dept*(*id*, *dname*, *budget*, *num\_majors*, *chair\_ssno*)

Suppose you know that the following queries are the five most common queries in the workload for this university and that all five are roughly equivalent in frequency and importance:

- List the names, ages, and offices of professors of a user-specified sex (male or female) who have a user-specified research specialty (e.g., *recursive query processing*). Assume that the university has a diverse set of faculty members, making it very uncommon for more than a few professors to have the same research specialty.
- List all the department information for departments with professors in a user-specified age range.
- List the *dept\_id*, department name, and chairperson name for departments with a user-specified number of majors.

- List the lowest budget for a department in the university.
- List all the information about professors who are department chairpersons.

These queries occur much more frequently than updates, so you should build whatever indexes you need to speed up these queries. However, you should not build any unnecessary indexes, as updates will occur (and would be slowed down by unnecessary indexes). Given this information, design a physical schema for the university database that will give good performance for the expected workload. In particular, decide which attributes should be indexed and whether each index should be a clustered index or an unclustered index. Assume that both B+ trees and hashed indexes are supported by the DBMS and that both single- and multiple-attribute index search keys are permitted.

1. Specify your physical design by identifying the attributes you recommend indexing on, indicating whether each index should be clustered or unclustered and whether it should be a B+ tree or a hashed index.
2. Assume that this workload is to be tuned with an automatic index tuning wizard. Outline the main steps in the algorithm and the set of candidate configurations considered.
3. Redesign the physical schema, assuming that the set of important queries is changed to be the following:
  - List the number of different specialties covered by professors in each department, by department.
  - Find the department with the fewest majors.
  - Find the youngest professor who is a department chairperson.

**Exercise 20.3** Consider the following BCNF relational schema for a portion of a company database (type information is not relevant to this question and is omitted):

*Project*(*pno*, *proj\_name*, *proj\_base\_dept*, *proj\_mgr*, *topic*, *budget*)  
*Manager*(*mid*, *mgr\_name*, *mgr\_dept*, *salary*, *age*, *sex*)

Note that each project is based in some department, each manager is employed in some department, and the manager of a project need not be employed in the same department (in which the project is based). Suppose you know that the following queries are the five most common queries in the workload for this university and all five are roughly equivalent in frequency and importance:

- List the names, ages, and salaries of managers of a user-specified sex (male or female) working in a given department. You can assume that, while there are many departments, each department contains very few project managers.
- List the names of all projects with managers whose ages are in a user-specified range (e.g., younger than 30).
- List the names of all departments such that a manager in this department manages a project based in this department.
- List the name of the project with the lowest budget.
- List the names of all managers in the same department as a given project.

These queries occur much more frequently than updates, so you should build whatever indexes you need to speed up these queries. However, you should not build any unnecessary indexes, as updates will occur (and would be slowed down by unnecessary indexes). Given

this information, design a physical schema for the company database that will give good performance for the expected workload. In particular, decide which attributes should be indexed and whether each index should be a clustered index or an unclustered index. Assume that both B+ trees and hashed indexes are supported by the DBMS, and that both single- and multiple-attribute index keys are permitted.

1. Specify your physical design by identifying the attributes you recommend indexing on, indicating whether each index should be clustered or unclustered and whether it should be a B+ tree or a hashed index.
2. Assume that this workload is to be tuned with an automatic index tuning wizard. Outline the main steps in the algorithm and the set of candidate configurations considered.
3. Redesign the physical schema assuming the set of important queries is changed to be the following:
  - Find the total of the budgets for projects managed by each manager; that is, list *proj\_mgr* and the total of the budgets of projects managed by that manager, for all values of *proj\_mgr*.
  - Find the total of the budgets for projects managed by each manager but only for managers who are in a user-specified age range.
  - Find the number of male managers.
  - Find the average age of managers.

**Exercise 20.4** The Globetrotters Club is organized into chapters. The president of a chapter can never serve as the president of any other chapter, and each chapter gives its president some salary. Chapters keep moving to new locations, and a new president is elected when (and only when) a chapter moves. This data is stored in a relation  $G(C, S, L, P)$ , where the attributes are chapters ( $C$ ), salaries ( $S$ ), locations ( $L$ ), and presidents ( $P$ ). Queries of the following form are frequently asked, and you must be able to answer them without computing a join: “Who was the president of chapter  $X$  when it was in location  $Y$ ?”

1. List the FDs that are given to hold over  $G$ .
2. What are the candidate keys for relation  $G$ ?
3. What normal form is the schema  $G$  in?
4. Design a good database schema for the club. (Remember that your design *must* satisfy the stated query requirement!)
5. What normal form is your good schema in? Give an example of a query that is likely to run slower on this schema than on the relation  $G$ .
6. Is there a lossless-join, dependency-preserving decomposition of  $G$  into BCNF?
7. Is there ever a good reason to accept something less than 3NF when designing a schema for a relational database? Use this example, if necessary adding further constraints, to illustrate your answer.

**Exercise 20.5** Consider the following BCNF relation, which lists the ids, types (e.g., nuts or bolts), and costs of various parts, along with the number available or in stock:

Parts (*pid*, *pname*, *cost*, *num\_avail*)

You are told that the following two queries are extremely important:

## Physical Database Design and Tuning

- Find the total number available by part type, for all types. (That is, the sum of the *num\_avail* value of all nuts, the sum of the *num\_avail* value of all bolts, and so forth)
  - List the *pids* of parts with the highest cost.
1. Describe the physical design that you would choose for this relation. That is, what kind of a file structure would you choose for the set of Parts records, and what indexes would you create?
  2. Suppose your customers subsequently complain that performance is still not satisfactory (given the indexes and file organization you chose for the Parts relation in response to the previous question). Since you cannot afford to buy new hardware or software, you have to consider a schema redesign. Explain how you would try to obtain better performance by describing the schema for the relation(s) that you would use and your choice of file organizations and indexes on these relations.
  3. How would your answers to the two questions change, if at all, if your system did not support indexes with multiple-attribute search keys?

Exercise 20.6 Consider the following BCNF relations, which describe employees and the departments they work in:

Emp (eid, sal, did)  
Dept (d'id, location, budget)

You are told that the following queries are extremely important:

- Find the location where a user-specified employee works.
  - Check whether the budget of a department is greater than the salary of each employee in that department.
1. Describe the physical design you would choose for this relation. That is, what kind of a file structure would you choose for these relations, and what indexes would you create?
  2. Suppose that your customers subsequently complain that performance is still not satisfactory (given the indexes and file organization that you chose for the relations in response to the previous question). Since you cannot afford to buy new hardware or software, you have to consider a schema redesign. Explain how you would try to obtain better performance by describing the schema for the relation(s) that you would use and your choice of file organizations and indexes on these relations.
  3. Suppose that your database system has very inefficient implementations of index structures. What kind of a design would you try in this case?

Exercise 20.7 Consider the following BCNF relations, which describe departments in a company and employees:

Dept(did, dname, location, managerid)  
Emp(eid, sal)

You are told that the following queries are extremely important:

- List the names and ids of managers for each department in a user-specified location, in alphabetical order by department name.
- Find the average salary of employees who manage departments in a user-specified location. You can assume that no one manages more than one department.

1. Describe the file structures and indexes that you would choose.
2. You subsequently realize that updates to these relations are frequent. Because indexes incur a high overhead, can you think of a way to improve performance on these queries without using indexes?

**Exercise 20.8** For each of the following queries, identify one possible reason why an optimizer might not find a good plan. Rewrite the query so that a good plan is likely to be found. Any available indexes or known constraints are listed before each query; assume that the relation schemas are consistent with the attributes referred to in the query.

1. An index is available on the *age* attribute:

```
SELECT E.dno
FROM   Employee E
WHERE  E.age=20 OR E.age=10
```

2. A B+ tree index is available on the *age* attribute:

```
SELECT E.dno
FROM   Employee E
WHERE  E.age<20 AND E.age>10
```

3. An index is available on the *age* attribute:

```
SELECT E.eIno
FROM   Employee E
WHERE  2*E.age<20
```

4. No index is available:

```
SELECT DISTINCT *
FROM   Employee E
```

5. No index is available:

```
SELECT  AVG (B.sal)
FROM    Employee E
GROUP BY E.dno
HAVING  E.dno=22
```

6. The *sid* in Reserves is a foreign key that refers to Sailors:

```
SELECT  S.sid
FROM    Sailors S, Reserves H
WHERE   S.sid=R.sid
```

**Exercise 20.9** Consider two ways to compute the names of employees who earn more than \$100,000 and whose age is equal to their manager's age. First, a nested query:

```
SELECT  E1.ename
FROM    Emp E1
WHERE   E1.sal > 100 AND E1.age = ( SELECT E2.age
                                   FROM   Emp E2, Dept D2
                                   WHERE  E1.dname = D2.dname
                                   AND   D2.mgr = E2.ename )
```

Second, a query that uses a view definition:

```

SELECT  E1.ename
FROM    Emp E1, MgrAge A
WHERE   E1.dname = A.dname AND E1.sal > 100 AND E1.age = A.age

```

```

CREATE VIEW MgrAge (dname, age)
AS SELECT D.dname, E.age
   FROM   Emp E, Dept D
   WHERE  D.mgr = E.ename

```

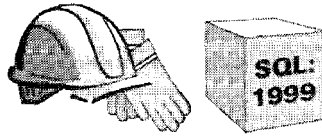
1. Describe a situation in which the first query is likely to outperform the second query.
2. Describe a situation in which the second query is likely to outperform the first query.
3. Can you construct an equivalent query that is likely to beat both these queries when every employee who earns more than \$100,000 is either 35 or 40 years old? Explain briefly.

## BIBLIOGRAPHIC NOTES

[658] is an early discussion of physical database design. [659] discusses the performance implications of normalization and observes that denormalization may improve performance for certain queries. The ideas underlying a physical design tool from IBM are described in [272]. The Microsoft AutoAdmin tool that performs automatic index selection according to a query workload is described in several papers [163, 164]. The DB2 Advisor is described in [750]. Other approaches to physical database design are described in [146, 639]. [679] considers *transaction tuning*, which we discussed only briefly. The issue is how an application should be structured into a collection of transactions to maximize performance.

The following books on database design cover physical design issues in detail; they are recommended for further reading. [274] is largely independent of specific products, although many examples are based on DB2 and Teradata systems. [779] deals primarily with DB2. Shasha and Bonnet give an in-depth, readable introduction to database tuning [104].

[334] contains several papers on benchmarking database systems and has accompanying software. It includes articles on the *AS<sup>3</sup>AP*, Set Query, TPC-A, TPC-B, Wisconsin, and OO1 benchmarks written by the original developers. The Bucky benchmark is described in [132], the OO1 benchmark is described in [131], and the TPC-D benchmark is described in [739]. The Sequoia 2000 benchmark is described in [720].



# 21

## SECURITY AND AUTHORIZATION

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

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

— William Congreve

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

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

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

## 21.1 INTRODUCTION TO DATABASE SECURITY

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

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

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



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

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

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

We use the following schemas in our examples:

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

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

## 21.2 ACCESS CONTROL

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

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

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

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

## 21.3 DISCRETIONARY ACCESS CONTROL

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

```
FROM   Sailors S
WHERE  S.age < 18
```

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

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

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

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

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

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

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

```
UPDATE Sailors S
SET     S.rating = 8
```

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

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

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

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

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

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

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

```
GRANT INSERT ON Sailors TO JMichael
```

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

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

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

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

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

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

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

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

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

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

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

Consider the following sequence, as another example:

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

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

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

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

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

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

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

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

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

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

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



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

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

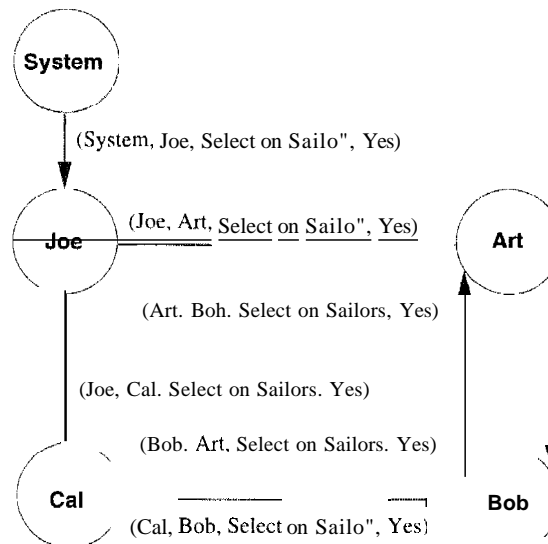


Figure 21.1 Example Authorization Graph

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

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

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

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

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

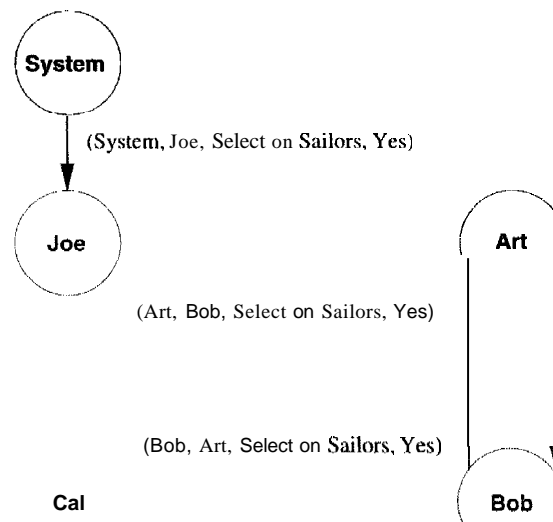


Figure 21.2 Example Authorization Graph during Revocation

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

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

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

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

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

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

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

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

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

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

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

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

## 21.4 MANDATORY ACCESS CONTROL

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

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

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

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

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

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

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

### 21.4.1 Multilevel Relations and Polyinstantiation

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

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

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

Figure 21.3 An Instance *B1* of Boats

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

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

This dilemma is resolved by effectively treating the security classification as part of the key. Thus, the insertion is allowed to continue, and the table instance is modified as shown in Figure 21.4.

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

Figure 21.4 Instance *B1* after Insertion

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

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

### 21.4.2 Covert Channels, DoD Security Levels

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

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

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

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

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

## 21.5 SECURITY FOR INTERNET APPLICATIONS

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

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

### 21.5.1 Encryption

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



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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

### 21.5.2 Certifying Servers: The SSL Protocol

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

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

( Verisign Amazon, http://www.amazon.com, *eA* )

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

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

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

<sup>2</sup>A browser uses the SSL protocol if the target URL begins with *https*.

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

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

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

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

### 21.5.3 Digital Signatures

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

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

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

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

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

## 21.6 ADDITIONAL ISSUES RELATED TO SECURITY

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

### 21.6.1 Role of the Database Administrator

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

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

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

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

## 21.6.2 Security in Statistical Databases

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

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