

support transaction processing. Chapter 18 describes the various techniques and Chapter 19 presents an overview of recovery techniques.

## 17.1 Introduction to Transaction Processing

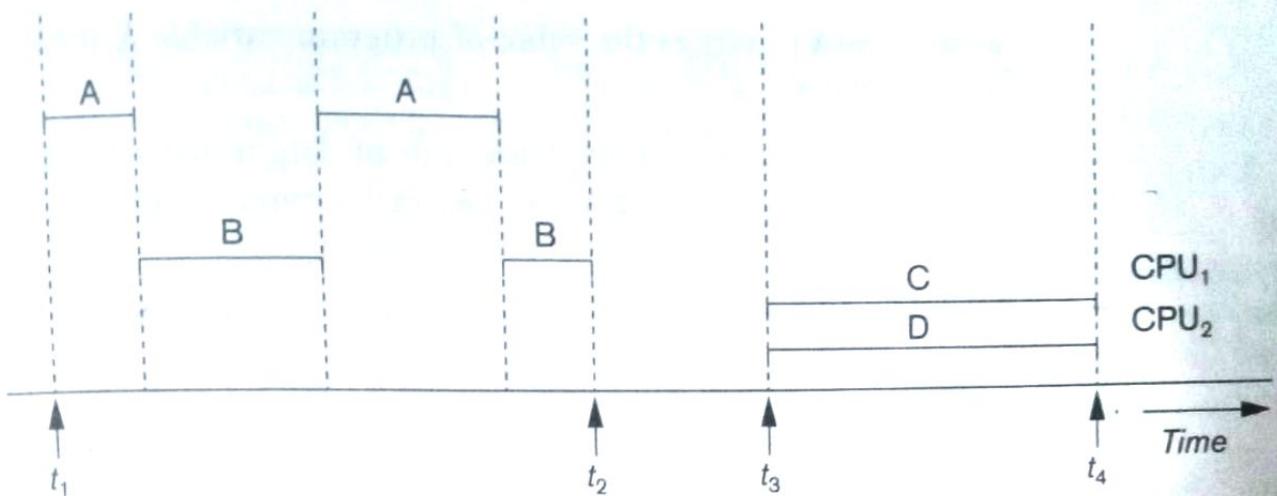
In this section we informally introduce the concepts of concurrent execution of transactions and recovery from transaction failures. Section 17.1.1 compares single-user and multiuser database systems and demonstrates how concurrent execution of transactions can take place in multiuser systems. Section 17.1.2 defines the concept of transaction and presents a simple model of transaction execution (based on read and write database operations) that is used to formalize concurrency control and recovery concepts. Section 17.1.3 shows by informal examples, why concurrency control techniques are needed in multiuser systems. Finally, Section 17.1.4 discusses why techniques are needed to permit recovery from failure by discussing the different ways in which transactions can fail while executing.

### 17.1.1 Single-User versus Multiuser Systems

One criterion for classifying a database system is according to the number of users who can use the system concurrently. A DBMS is **single-user** if at most one user at a time can use the system, and it is **multiuser** if many users can use the system—and hence access the database—concurrently. Single-user DBMSs are mostly restricted to personal computer systems; most other DBMSs are multiuser. For example, an airline reservations system is used by hundreds of travel agents and reservation clerks concurrently. Systems in banks, insurance agencies, stock exchanges, supermarkets, and the like are also operated on by many users who submit transactions concurrently to the system.

Multiple users can access databases—and use computer systems—simultaneously because of the concept of **multiprogramming**, which allows the computer to execute multiple programs—or **processes**—at the same time. If only a single central processing unit (CPU) exists, it can actually execute at most one process at a time. However, **multiprogramming operating systems** execute some commands from one process, then suspend that process and execute some commands from the next process, and so on. A process is resumed at the point where it was suspended whenever it gets its turn to use the CPU again. Hence, concurrent execution of processes is actually **interleaved**, as illustrated in Figure 17.1, which shows two processes A and B executing concurrently in an interleaved fashion. Interleaving keeps the CPU busy when a process requires an input or output (I/O) operation, such as reading a block from disk. The CPU is switched to execute another process rather than remaining idle during I/O time. Interleaving also prevents a long process from delaying other processes.

If the computer system has multiple hardware processors (CPUs), **parallel processing** of multiple processes is possible, as illustrated by processes C and D in Figure 17.1. Most of the theory concerning concurrency control in databases is



Interleaved processing versus parallel processing of concurrent transactions.

developed in terms of **interleaved concurrency**, so for the remainder of this chapter we assume this model. In a multiuser DBMS, the stored data items are the primary resources that may be accessed concurrently by interactive users or application programs, which are constantly retrieving information from and modifying the database.

### 17.1.2 Transactions, Read and Write Operations, and DBMS Buffers

A transaction is an executing program that forms a logical unit of database processing. A transaction includes one or more database access operations—these can include insertion, deletion, modification, or retrieval operations. The database operations that form a transaction can either be embedded within an application program or they can be specified interactively via a high-level query language such as SQL. One way of specifying the transaction boundaries is by specifying explicit **begin transaction** and **end transaction** statements in an application program; in this case, all database access operations between the two are considered as forming one transaction. A single application program may contain more than one transaction if it contains several transaction boundaries. If the database operations in a transaction do not update the database but only retrieve data, the transaction is called a **read-only transaction**.

The model of a database that is used to explain transaction processing concepts is much simplified. A **database** is basically represented as a collection of **named data items**. The size of a data item is called its **granularity**. It can be a field of some record in the database, or it may be a larger unit such as a record or even a whole disk block, but the concepts we discuss are independent of the data item granularity. Using this simplified database model, the basic database access operations that a transaction can include are as follows:

- **read\_item(X)**. Reads a database item named X into a program variable. To simplify our notation, we assume that the program variable is also named X.

- **write\_item( $X$ )**. Writes the value of program variable  $X$  into the database item named  $X$ .

As we discussed in Chapter 13, the basic unit of data transfer from disk to main memory is one block. Executing a **read\_item( $X$ )** command includes the following steps:

1. Find the address of the disk block that contains item  $X$ .
2. Copy that disk block into a buffer in main memory (if that disk block is not already in some main memory buffer).
3. Copy item  $X$  from the buffer to the program variable named  $X$ .

Executing a **write\_item( $X$ )** command includes the following steps:

1. Find the address of the disk block that contains item  $X$ .
2. Copy that disk block into a buffer in main memory (if that disk block is not already in some main memory buffer).
3. Copy item  $X$  from the program variable named  $X$  into its correct location in the buffer.
4. Store the updated block from the buffer back to disk (either immediately or at some later point in time).

Step 4 actually updates the database on disk. In some cases the buffer is not immediately stored to disk, in case additional changes are to be made to the buffer. Usually, the decision about when to store a modified disk block that is in a main memory buffer is handled by the recovery manager of the DBMS in cooperation with the underlying operating system. The DBMS will generally maintain a number of buffers in main memory that hold database disk blocks containing the database items being processed. When these buffers are all occupied, and additional database blocks must be copied into memory, some buffer replacement policy is used to choose which of the current buffers is to be replaced. If the chosen buffer has been modified, it must be written back to disk before it is reused.<sup>1</sup>

A transaction includes **read\_item** and **write\_item** operations to access and update the database. Figure 17.2 shows examples of two very simple transactions. The **read-set** of a transaction is the set of all items that the transaction reads, and the **write-set** is the set of all items that the transaction writes. For example, the read-set of  $T_1$  in Figure 17.2 is  $\{X, Y\}$  and its write-set is also  $\{X, Y\}$ .

Concurrency control and recovery mechanisms are mainly concerned with the database access commands in a transaction. Transactions submitted by the various users may execute concurrently and may access and update the same database items. If this concurrent execution is uncontrolled, it may lead to problems, such as an inconsistent database. In the next section we informally introduce some of the problems that may occur.

<sup>1</sup>We will not discuss buffer replacement policies here because they are typically discussed in operating systems textbooks.

(a)	$T_1$	(b)	$T_2$
	<pre>read_item(X); X := X - N; write_item(X); read_item(Y); Y := Y + N; write_item(Y);</pre>		<pre>read_item(X); X := X + M; write_item(X);</pre>

Two sample transactions  
(a) Transaction  $T_1$   
(b) Transaction  $T_2$

### 17.1.3 Why Concurrency Control Is Needed

Several problems can occur when concurrent transactions execute in an uncontrolled manner. We illustrate some of these problems by referring to a much simplified airline reservations database in which a record is stored for each airline flight. Each record includes the number of reserved seats on that flight as a *named data item*, among other information. Figure 17.2(a) shows a transaction  $T_1$  that *transfers*  $N$  reservations from one flight whose number of reserved seats is stored in the database item named  $X$  to another flight whose number of reserved seats is stored in the database item named  $Y$ . Figure 17.2(b) shows a simpler transaction  $T_2$  that just *reserves*  $M$  seats on the first flight ( $X$ ) referenced in transaction  $T_1$ .<sup>2</sup> To simplify our example, we do not show additional portions of the transactions, such as checking whether a flight has enough seats available before reserving additional seats.

When a database access program is written, it has the flight numbers, their dates, and the number of seats to be booked as parameters; hence, the same program can be used to execute many transactions, each with different flights and numbers of seats to be booked. For concurrency control purposes, a transaction is a *particular execution* of a program on a specific date, flight, and number of seats. In Figure 17.2(a) and (b), the transactions  $T_1$  and  $T_2$  are *specific executions* of the programs that refer to the specific flights whose numbers of seats are stored in data items  $X$  and  $Y$  in the database. Next we discuss the types of problems we may encounter with these transactions if they run concurrently.

**The Lost Update Problem.** This problem occurs when two transactions that access the same database items have their operations interleaved in a way that makes the value of some database items incorrect. Suppose that transactions  $T_1$  and  $T_2$  are submitted at approximately the same time, and suppose that their operations are interleaved as shown in Figure 17.3(a); then the final value of item  $X$  is incorrect because  $T_2$  reads the value of  $X$  *before*  $T_1$  changes it in the database, and hence the updated value resulting from  $T_1$  is lost. For example, if  $X = 80$  at the start (originally there were 80 reservations on the flight),  $N = 5$  ( $T_1$  transfers 5 seat reservations from the flight corresponding to  $X$  to the flight corresponding to  $Y$ ), and  $M = 4$ ,

<sup>2</sup> A similar, more commonly used example assumes a bank database, with one transaction doing a transfer of funds from account  $X$  to account  $Y$  and the other transaction doing a deposit to account  $X$ .

**Figure 17.3**

Some problems that occur when concurrent execution is uncontrolled. (a) The lost update problem. (b) The temporary update problem. (c) The incorrect summary problem.

(a)

$T_1$	$T_2$
read_item( $X$ ); $X := X - N;$	
write_item( $X$ ); read_item( $Y$ ); $Y := Y + N;$	read_item( $X$ ); $X := X + M;$ write_item( $X$ );

$$X = 80$$

$$N = 5$$

$$M = 4$$

Item  $X$  has an incorrect value because its update by  $T_1$  is *lost* (overwritten).

(b)

$T_1$	$T_2$
read_item( $X$ ); $X := X - N;$ write_item( $X$ );	
read_item( $Y$ );	read_item( $X$ ); $X := X + M;$ write_item( $X$ );

Transaction  $T_1$  fails and must change the value of  $X$  back to its old value; meanwhile  $T_2$  has read the *temporary* incorrect value of  $X$ .

(c)

$T_1$	$T_3$
	$sum := 0;$ read_item( $A$ ); $sum := sum + A;$ $\vdots$
read_item( $X$ ); $X := X - N;$ write_item( $X$ );  read_item( $Y$ ); $Y := Y + N;$ write_item( $Y$ );	read_item( $X$ ); $sum := sum + X;$ read_item( $Y$ ); $sum := sum + Y;$

$T_3$  reads  $X$  after  $N$  is subtracted and reads  $Y$  before  $N$  is added; a wrong summary is the result (off by  $N$ ).

( $T_2$  reserves 4 seats on  $X$ ), the final result should be  $X = 79$ ; but in the interleaving of operations shown in Figure 17.3(a), it is  $X = \underline{84}$  because the update in  $T_1$  that removed the five seats from  $X$  was *lost*.

**The Temporary Update (or Dirty Read) Problem.** This problem occurs when one transaction updates a database item and then the transaction fails for some reason (see Section 17.1.4). The updated item is accessed by another transaction before it is changed back to its original value. Figure 17.3(b) shows an example where  $T_1$  updates item  $X$  and then fails before completion, so the system must change  $X$  back to its original value. Before it can do so, however, transaction  $T_2$  reads the temporary value of  $X$ , which will not be recorded permanently in the database because of the failure of  $T_1$ . The value of item  $X$  that is read by  $T_2$  is called *dirty data* because it has been created by a transaction that has not completed and committed yet; hence, this problem is also known as the *dirty read problem*.

**The Incorrect Summary Problem.** If one transaction is calculating an aggregate summary function on a number of records while other transactions are updating some of these records, the aggregate function may calculate some values before they are updated and others after they are updated. For example, suppose that a transaction  $T_3$  is calculating the total number of reservations on all the flights; meanwhile, transaction  $T_1$  is executing. If the interleaving of operations shown in Figure 17.3(c) occurs, the result of  $T_3$  will be off by an amount  $N$  because  $T_3$  reads the value of  $X$  after  $N$  seats have been subtracted from it but reads the value of  $Y$  before those  $N$  seats have been added to it.

Another problem that may occur is called **unrepeatable read**, where a transaction  $T$  reads an item twice and the item is changed by another transaction  $T'$  between the two reads. Hence,  $T$  receives *different values* for its two reads of the same item. This may occur, for example, if during an airline reservation transaction, a customer inquires about seat availability on several flights. When the customer decides on a particular flight, the transaction then reads the number of seats on that flight a second time before completing the reservation.

#### 17.1.4 Why Recovery Is Needed

Whenever a transaction is submitted to a DBMS for execution, the system is responsible for making sure that either all the operations in the transaction are completed successfully and their effect is recorded permanently in the database, or the transaction does not affect the database or any other transactions. The DBMS must not permit some operations of a transaction  $T$  to be applied to the database while other operations of  $T$  are not. This may happen if a transaction **fails** after executing some of its operations but before executing all of them.

**Types of Failures.** Failures are generally classified as transaction, system, and media failures. There are several possible reasons for a transaction to fail in the middle of execution:

1. **A computer failure (system crash).** A hardware, software, or network error occurs in the computer system during transaction execution. Hardware crashes are usually media failures—for example, main memory failure.
2. **A transaction or system error.** Some operation in the transaction may cause it to fail, such as integer overflow or division by zero. Transaction failure may also occur because of erroneous parameter values or because of a logical programming error.<sup>3</sup> Additionally, the user may interrupt the transaction during its execution.
3. **Local errors or exception conditions detected by the transaction.** During transaction execution, certain conditions may occur that necessitate cancellation of the transaction. For example, data for the transaction may not be found. Notice that an exception condition,<sup>4</sup> such as insufficient account balance in a banking database, may cause a transaction, such as a fund withdrawal, to be canceled. This exception should be programmed in the transaction itself, and hence would not be considered a failure.
4. **Concurrency control enforcement.** The concurrency control method (see Chapter 18) may decide to abort the transaction, to be restarted later, because it violates serializability (see Section 17.5) or because several transactions are in a state of deadlock.
5. **Disk failure.** Some disk blocks may lose their data because of a read or write malfunction or because of a disk read/write head crash. This may happen during a read or a write operation of the transaction.
6. **Physical problems and catastrophes.** This refers to an endless list of problems that includes power or air-conditioning failure, fire, theft, sabotage, overwriting disks or tapes by mistake, and mounting of a wrong tape by the operator.

Failures of types 1, 2, 3, and 4 are more common than those of types 5 or 6. Whenever a failure of type 1 through 4 occurs, the system must keep sufficient information to recover from the failure. Disk failure or other catastrophic failures of type 5 or 6 do not happen frequently; if they do occur, recovery is a major task. We discuss recovery from failure in Chapter 19.

The concept of transaction is fundamental to many techniques for concurrency control and recovery from failures.

## 17.2 Transaction and System Concepts

In this section we discuss additional concepts relevant to transaction processing. Section 17.2.1 describes the various states a transaction can be in, and discusses additional relevant operations needed in transaction processing. Section 17.2.2 discusses

3. In general, a transaction should be thoroughly tested to ensure that it has no bugs (logical programming errors).

4. Exception conditions, if programmed correctly, do not constitute transaction failures.

the system log, which keeps information needed for recovery. Section 17.2.3 describes the concept of commit points of transactions, and why they are important in transaction processing.

### 17.2.1 Transaction States and Additional Operations

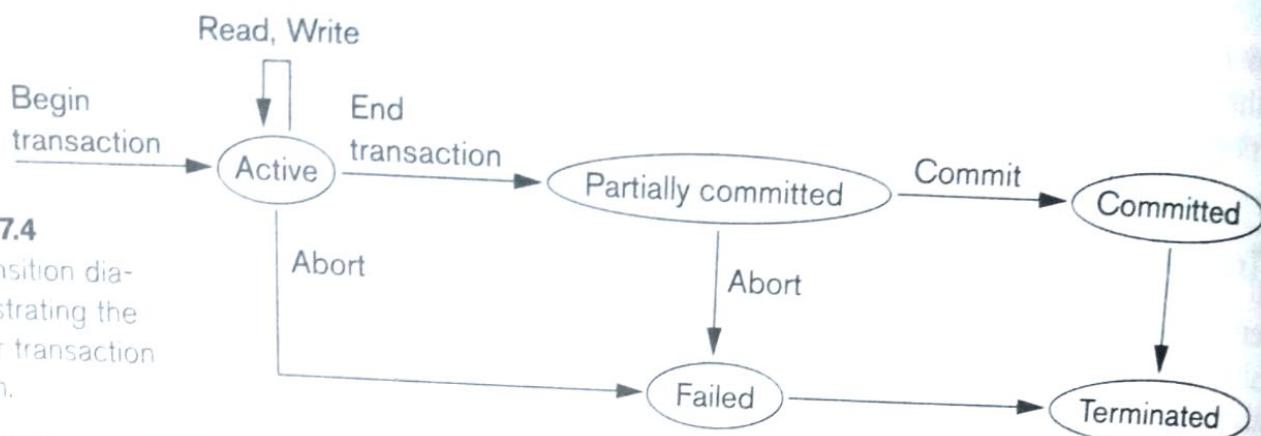
A transaction is an atomic unit of work that is either completed in its entirety or not done at all. For recovery purposes, the system needs to keep track of when the transaction starts, terminates, and commits or aborts (see Section 17.2.3). Therefore, the recovery manager keeps track of the following operations:

- BEGIN\_TRANSACTION. This marks the beginning of transaction execution.
- READ OR WRITE. These specify read or write operations on the database items that are executed as part of a transaction.
- END\_TRANSACTION. This specifies that READ and WRITE transaction operations have ended and marks the end of transaction execution. However, at this point it may be necessary to check whether the changes introduced by the transaction can be permanently applied to the database (committed) or whether the transaction has to be aborted because it violates serializability (see Section 17.5) or for some other reason.
- COMMIT\_TRANSACTION. This signals a *successful end* of the transaction so that any changes (updates) executed by the transaction can be safely committed to the database and will not be undone.
- ROLLBACK (or ABORT). This signals that the transaction has *ended unsuccessfully*, so that any changes or effects that the transaction may have applied to the database must be *undone*.

Figure 17.4 shows a state transition diagram that describes how a transaction moves through its execution states. A transaction goes into an **active state** immediately after it starts execution, where it can issue READ and WRITE operations. When the transaction ends, it moves to the **partially committed state**. At this point, some recovery protocols need to ensure that a system failure will not result in an inability to record the changes of the transaction permanently (usually by recording changes in the system log, discussed in the next section).<sup>5</sup> Once this check is successful, the transaction is said to have reached its **commit point** and enters the **committed state**. Commit points are discussed in more detail in Section 17.2.3. Once a transaction is committed, it has concluded its execution successfully and all its changes must be recorded permanently in the database.

However, a transaction can go to the **failed state** if one of the checks fails or if the transaction is aborted during its active state. The transaction may then have to be rolled back to undo the effect of its WRITE operations on the database. The **terminated state** corresponds to the transaction leaving the system. The transaction information that is maintained in system tables while the transaction has been

<sup>5</sup> Optimistic concurrency control (see Section 18.4) also requires that certain checks are made at this point to ensure that the transaction did not interfere with other executing transactions.

**Figure 17.4**

State transition diagram illustrating the states for transaction execution.

running is removed when the transaction terminates. Failed or aborted transaction may be *restarted* later—either automatically or after being resubmitted by the user—as brand new transactions.

### 17.2.2 The System Log /DBMS Journal

To be able to recover from failures that affect transactions, the system maintains a log<sup>6</sup> to keep track of all transaction operations that affect the values of database items. This information may be needed to permit recovery from failures. The log is kept on disk, so it is not affected by any type of failure except for disk or catastrophic failure. In addition, the log is periodically backed up to archival storage (tape) to guard against such catastrophic failures. The following lists the types of entries—called **log records**—that are written to the log and the action each performs. In these entries,  $T$  refers to a unique **transaction-id** that is generated automatically by the system and is used to identify each transaction:

- 1. **[start\_transaction,  $T$ ]**. Indicates that transaction  $T$  has started execution.
- 2. **[write\_item,  $T$ ,  $X$ ,  $old\_value$ ,  $new\_value$ ]**. Indicates that transaction  $T$  has changed the value of database item  $X$  from  $old\_value$  to  $new\_value$ .
- 3. **[read\_item,  $T$ ,  $X$ ]**. Indicates that transaction  $T$  has read the value of database item  $X$ .
- 4. **[commit,  $T$ ]**. Indicates that transaction  $T$  has completed successfully, and affirms that its effect can be committed (recorded permanently) to the database.
- 5. **[abort,  $T$ ]**. Indicates that transaction  $T$  has been aborted.

Protocols for recovery that avoid cascading rollbacks (see Section 17.4.2)—which include nearly all practical protocols—do not require that **READ** operations are written to the system log. However, if the log is also used for other purposes—such as auditing (keeping track of all database operations)—then such entries can be included. Additionally, some recovery protocols require simpler **WRITE** entries that do not include  $new\_value$  (see Section 17.4.2).

6. The log has sometimes been called the DBMS journal.

Notice that we are assuming that all permanent changes to the database occur within transactions, so the notion of recovery from a transaction failure amounts to either undoing or redoing transaction operations individually from the log. If the system crashes, we can recover to a consistent database state by examining the log and using one of the techniques described in Chapter 19. Because the log contains a record of every WRITE operation that changes the value of some database item, it is possible to **undo** the effect of these WRITE operations of a transaction  $T$  by tracing backward through the log and resetting all items changed by a WRITE operation of  $T$  to their `old_values`. **Redoing** the operations of a transaction may also be necessary if all its updates are recorded in the log but a failure occurs before we can be sure that all these `new_values` have been written permanently in the actual database on disk.<sup>7</sup> Redoing the operations of transaction  $T$  is applied by tracing forward through the log and setting all items changed by a WRITE operation of  $T$  to their `new_values`.

### 17.2.3 Commit Point of a Transaction

A transaction  $T$  reaches its **commit point** when all its operations that access the database have been executed successfully and the effect of all the transaction operations on the database have been recorded in the log. Beyond the commit point, the transaction is said to be **committed**, and its effect is assumed to be *permanently recorded* in the database. The transaction then writes a commit record [`commit, T`] into the log. If a system failure occurs, we search back in the log for all transactions  $T$  that have written a [`start_transaction, T`] record into the log but have not written their [`commit, T`] record yet; these transactions may have to be *rolled back* to undo their effect on the database during the recovery process. Transactions that have written their commit record in the log must also have recorded all their WRITE operations in the log, so their effect on the database can be *redone* from the log records.

Notice that the log file must be kept on disk. As discussed in Chapter 13, updating a disk file involves copying the appropriate block of the file from disk to a buffer in main memory, updating the buffer in main memory, and copying the buffer to disk. It is common to keep one or more blocks of the log file in main memory buffers until they are filled with log entries and then to write them back to disk only once, rather than writing to disk every time a log entry is added. This saves the overhead of multiple disk writes of the same log file block. At the time of a system crash, only the log entries that have been *written back to disk* are considered in the recovery process because the contents of main memory may be lost. Hence, *before* a transaction reaches its commit point, any portion of the log that has not been written to the disk yet must now be written to the disk. This process is called **force-writing** the log file before committing a transaction.

## 17.3 Desirable Properties of Transactions

Transactions should possess several properties, often called the **ACID** properties; they should be enforced by the concurrency control and recovery methods of the

<sup>7</sup> Undo and redo are discussed more fully in Chapter 19.

DBMS. The following are the ACID properties:

- **Atomicity.** A transaction is an atomic unit of processing; it is either performed in its entirety or not performed at all.
- **Consistency preservation.** A transaction is consistency preserving if its complete execution takes the database from one consistent state to another.
- **Isolation.** A transaction should appear as though it is being executed in isolation from other transactions. That is, the execution of a transaction should not be interfered with by any other transactions executing concurrently.
- **Durability or permanency.** The changes applied to the database by a committed transaction must persist in the database. These changes must not be lost because of any failure.

The atomicity property requires that we execute a transaction to completion. It is the responsibility of the transaction recovery subsystem of a DBMS to ensure atomicity. If a transaction fails to complete for some reason, such as a system crash in the midst of transaction execution, the recovery technique must undo any effects of the transaction on the database.

*niddle part*

The preservation of consistency is generally considered to be the responsibility of the programmers who write the database programs or of the DBMS module that enforces integrity constraints. Recall that a **database state** is a collection of all the stored data items (values) in the database at a given point in time. A **consistent state** of the database satisfies the constraints specified in the schema as well as any other constraints on the database that should hold. A database program should be written in a way that guarantees that, if the database is in a consistent state before executing the transaction, it will be in a consistent state after the *complete* execution of the transaction, assuming that *no interference with other transactions* occurs.

Isolation is enforced by the concurrency control subsystem of the DBMS.<sup>8</sup> If every transaction does not make its updates visible to other transactions until it is committed, one form of isolation is enforced that solves the temporary update problem and eliminates cascading rollbacks (see Chapter 19). There have been attempts to define the **level of isolation** of a transaction. A transaction is said to have level 0 (zero) isolation if it does not overwrite the dirty reads of higher-level transactions. Level 1 (one) isolation has no lost updates; and level 2 isolation has no lost updates and no dirty reads. Finally, level 3 isolation (also called *true isolation*) has, in addition to degree 2 properties, repeatable reads.

Finally, the durability property is the responsibility of the recovery subsystem of the DBMS. We will discuss how recovery protocols enforce durability and atomicity in Chapter 19.

<sup>8</sup> We will discuss concurrency control protocols in Chapter 18.

## 17.4 Characterizing Schedules Based on Recoverability

When transactions are executing concurrently in an interleaved fashion, then the order of execution of operations from the various transactions is known as a schedule (or history). In this section, first we define the concept of schedule, and then we characterize the types of schedules that facilitate recovery when failures occur. In Section 17.5, we characterize schedules in terms of the interference of participating transactions, leading to the concepts of serializability and serializable schedules.

### 17.4.1 Schedules (Histories) of Transactions

A schedule (or history)  $S$  of  $n$  transactions  $T_1, T_2, \dots, T_n$  is an ordering of the operations of the transactions subject to the constraint that, for each transaction  $T_i$  that participates in  $S$ , the operations of  $T_i$  in  $S$  must appear in the same order in which they occur in  $T_i$ . Note, however, that operations from other transactions  $T_j$  can be interleaved with the operations of  $T_i$  in  $S$ . For now, consider the order of operations in  $S$  to be a *total ordering*, although it is possible theoretically to deal with schedules whose operations form *partial orders* (as we discuss later).

For the purpose of recovery and concurrency control, we are mainly interested in the `read_item` and `write_item` operations of the transactions, as well as the `commit` and `abort` operations. A shorthand notation for describing a schedule uses the symbols  $r$ ,  $w$ ,  $c$ , and  $a$  for the operations `read_item`, `write_item`, `commit`, and `abort`, respectively, and appends as subscript the transaction id (transaction number) to each operation in the schedule. In this notation, the database item  $X$  that is read or written follows the  $r$  and  $w$  operations in parentheses. For example, the schedule of Figure 17.3(a), which we shall call  $S_a$ , can be written as follows in this notation:

$S_a: r_1(X); r_2(X); w_1(X); r_1(Y); w_2(X); w_1(Y);$

Similarly, the schedule for Figure 17.3(b), which we call  $S_b$ , can be written as follows, if we assume that transaction  $T_1$  aborted after its `read_item(Y)` operation:

$S_b: r_1(X); w_1(X); r_2(X); w_2(X); r_1(Y); a_1;$

Two operations in a schedule are said to **conflict** if they satisfy all three of the following conditions: (1) they belong to different transactions; (2) they access the same item  $X$ ; and (3) at least one of the operations is a `write_item(X)`. For example, in schedule  $S_a$ , the operations  $r_1(X)$  and  $w_2(X)$  conflict, as do the operations  $r_2(X)$  and  $w_1(X)$ , and the operations  $w_1(X)$  and  $w_2(X)$ . However, the operations  $r_1(X)$  and  $r_2(X)$  do not conflict, since they are both read operations; the operations  $w_2(X)$  and  $w_1(Y)$  do not conflict because they operate on distinct data items  $X$  and  $Y$ ; and the operations  $r_1(X)$  and  $w_1(X)$  do not conflict because they belong to the same transaction.

A schedule  $S$  of  $n$  transactions  $T_1, T_2, \dots, T_n$  is said to be a **complete schedule** if the following conditions hold:

1. The operations in  $S$  are exactly those operations in  $T_1, T_2, \dots, T_n$ , including a `commit` or `abort` operation as the last operation for each transaction in the schedule.

- { 2. For any pair of operations from the same transaction  $T_i$ , their order of appearance in  $S$  is the same as their order of appearance in  $T_i$ .
- 3. For any two conflicting operations, one of the two must occur before the other in the schedule.<sup>9</sup>

The preceding condition (3) allows for two *nonconflicting operations* to occur in the schedule without defining which occurs first, thus leading to the definition of a schedule as a **partial order** of the operations in the  $n$  transactions.<sup>10</sup> However, a total order must be specified in the schedule for any pair of conflicting operations (condition 3) and for any pair of operations from the same transaction (condition 2). Condition 1 simply states that all operations in the transactions must appear in the complete schedule. Since every transaction has either committed or aborted, a complete schedule will not contain any active transactions at the end of the schedule.

In general, it is difficult to encounter complete schedules in a transaction processing system because new transactions are continually being submitted to the system. Hence, it is useful to define the concept of the **committed projection**  $C(S)$  of a schedule  $S$ , which includes only the operations in  $S$  that belong to committed transactions—that is, transactions  $T_i$  whose commit operation  $c_i$  is in  $S$ .

#### 17.4.2 Characterizing Schedules Based on Recoverability

For some schedules it is easy to recover from transaction failures, whereas for other schedules the recovery process can be quite involved. Hence, it is important to characterize the types of schedules for which recovery is possible, as well as those for which recovery is relatively simple. These characterizations do not actually provide the recovery algorithm; they only attempt to theoretically characterize the different types of schedules.

First, we would like to ensure that, once a transaction  $T$  is committed, it should never be necessary to roll back  $T$ . The schedules that theoretically meet this criterion are called *recoverable schedules* and those that do not are called **nonrecoverable**, and hence should not be permitted. A schedule  $S$  is recoverable if no transaction  $T$  in  $S$  commits until all transactions  $T'$  that have written an item that  $T$  reads have committed. A transaction  $T$  **reads** from transaction  $T'$  in a schedule  $S$  if some item  $X$  is first written by  $T'$  and later read by  $T$ . In addition,  $T'$  should not have been aborted before  $T$  reads item  $X$ , and there should be no transactions that write  $X$  after  $T'$  writes it and before  $T$  reads it (unless those transactions, if any, have aborted before  $T$  reads  $X$ ).

Recoverable schedules require a complex recovery process as we shall see, but if sufficient information is kept (in the log), a recovery algorithm can be devised.

9. Theoretically, it is not necessary to determine an order between pairs of *nonconflicting operations*.

10. In practice, most schedules have a total order of operations. If parallel processing is employed, it is theoretically possible to have schedules with partially ordered nonconflicting operations.

The (partial) schedules  $S_a$  and  $S_b$  from the preceding section are both recoverable, since they satisfy the above definition. Consider the schedule  $S'_a$  given below, which is the same as schedule  $S_a$  except that two commit operations have been added to  $S_a$ :

$$S'_a: r_1(X); r_2(X); w_1(X); r_1(Y); w_2(X); c_2; w_1(Y); c_1;$$

$S'_a$  is recoverable, even though it suffers from the lost update problem. However, consider the two (partial) schedules  $S_c$  and  $S_d$  that follow:

$$S_c: r_1(X); w_1(X); r_2(X); r_1(Y); w_2(X); c_2; a_1;$$

$$S_d: r_1(X); w_1(X); r_2(X); r_1(Y); w_2(X); w_1(Y); c_1; c_2;$$

$$S_e: r_1(X); w_1(X); r_2(X); r_1(Y); w_2(X); w_1(Y); a_1; a_2;$$

$S_c$  is not recoverable because  $T_2$  reads item  $X$  from  $T_1$ , and then  $T_2$  commits before  $T_1$  commits. If  $T_1$  aborts after the  $c_2$  operation in  $S_c$ , then the value of  $X$  that  $T_2$  read is no longer valid and  $T_2$  must be aborted *after* it is committed, leading to a schedule that is not recoverable. For the schedule to be recoverable, the  $c_2$  operation in  $S_c$  must be postponed until after  $T_1$  commits, as shown in  $S_d$ ; if  $T_1$  aborts instead of committing, then  $T_2$  should also abort as shown in  $S_e$ , because the value of  $X$  it read is no longer valid.

In a recoverable schedule, no committed transaction ever needs to be rolled back. However, it is possible for a phenomenon known as **cascading rollback** (or cascading abort) to occur, where an *uncommitted* transaction has to be rolled back because it read an item from a transaction that failed. This is illustrated in schedule  $S_g$ , where transaction  $T_2$  has to be rolled back because it read item  $X$  from  $T_1$ , and  $T_1$  then aborted.

Because cascading rollback can be quite time-consuming—since numerous transactions can be rolled back (see Chapter 19)—it is important to characterize the schedules where this phenomenon is guaranteed not to occur. A schedule is said to be **cascadeless**, or to **avoid cascading rollback**, if every transaction in the schedule reads only items that were written by committed transactions. In this case, all items read will not be discarded, so no cascading rollback will occur. To satisfy this criterion, the  $r_2(X)$  command in schedules  $S_d$  and  $S_e$  must be postponed until after  $T_1$  has committed (or aborted), thus delaying  $T_2$  but ensuring no cascading rollback if  $T_1$  aborts.

Finally, there is a third, more restrictive type of schedule, called a **strict schedule**, in which transactions can *neither read nor write* an item  $X$  until the last transaction that wrote  $X$  has committed (or aborted). Strict schedules simplify the recovery process. In a strict schedule, the process of undoing a `write_item(X)` operation of an aborted transaction is simply to restore the **before image** (`old_value` or BFIM) of data item  $X$ . This simple procedure always works correctly for strict schedules, but it may not work for recoverable or cascadeless schedules. For example, consider schedule  $S_f$ :

$$S_f: w_1(X, 5); w_2(X, 8); a_1;$$