

## Introduction to Transaction Processing Concepts and Theory

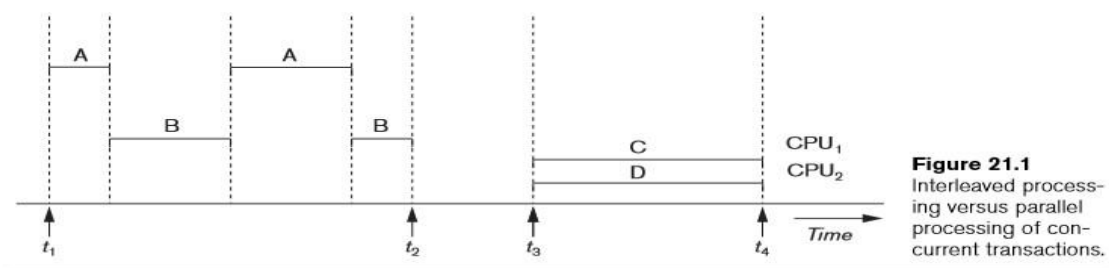
The concept of transaction provides a mechanism for describing logical units of database processing. Transaction processing systems are systems with large databases and hundreds of concurrent users executing database transactions. Examples of such systems include airline reservations, banking, credit card processing, online retail purchasing, stock markets, supermarket checkouts, and many other applications.

### ***Introduction to Transaction Processing***

In this section we discuss the concepts of concurrent execution of transactions and recovery from transaction failures.

### **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. Database systems used in banks, insurance agencies, stock exchanges, supermarkets, and many other applications are multiuser systems. In these systems, hundreds or thousands of users are typically operating on the database by submitting transactions concurrently to the system.



**Figure 21.1**  
Interleaved processing versus parallel processing of concurrent transactions.

Multiple users can access databases—and use computer systems—simultaneously because of the concept of multiprogramming, which allows the operating system of the computer to execute multiple programs—or processes—at the same time. A single central processing unit (CPU) can only 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 21.1, 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 21.1.

### Transactions, Database Items, 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**; otherwise it is known as a **read-write transaction**.

A **database** is basically represented as a collection of named data items. The size of a data item is called its **granularity**. A data item can be a database record, but it can also be a larger unit such as a whole disk block, or even a smaller unit such as an individual field (attribute) value of some record in the database.

The basic database access operations that a transaction can include are as follows:

- 1) `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`.
- 2) `write_item(X)`. Writes the value of program variable `X` into the database item named `X`.

Executing a `read_item(X)` command includes the following steps:

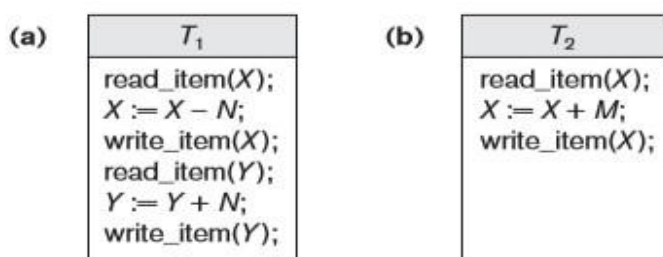
1. Find the address of the disk block that contains item `X`.
2. Copy that disk block into a buffer in main memory (if that disk block is not already in some main memory buffer).
3. Copy item `X` from the buffer to the program variable named `X`.

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).

Why Concurrency Control Is Needed

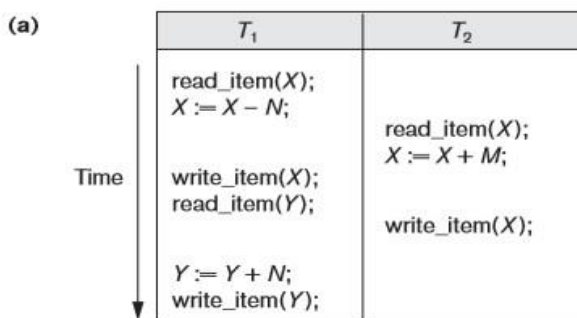
Several problems can occur when concurrent transactions execute in an uncontrolled manner. We illustrate some of these problems by referring to a much simplified airline reservations database in which a record is stored for each airline flight. Each record includes the number of reserved seats on that flight as a named (uniquely identifiable) data item, among other information. Figure 21.2(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 21.2(b) shows a simpler transaction  $T_2$  that just reserves  $M$  seats on the first flight ( $X$ ) referenced in transaction  $T_1$ .

**Figure 21.2**

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

Next we discuss the types of problems we may encounter with these two simple transactions if they run concurrently.

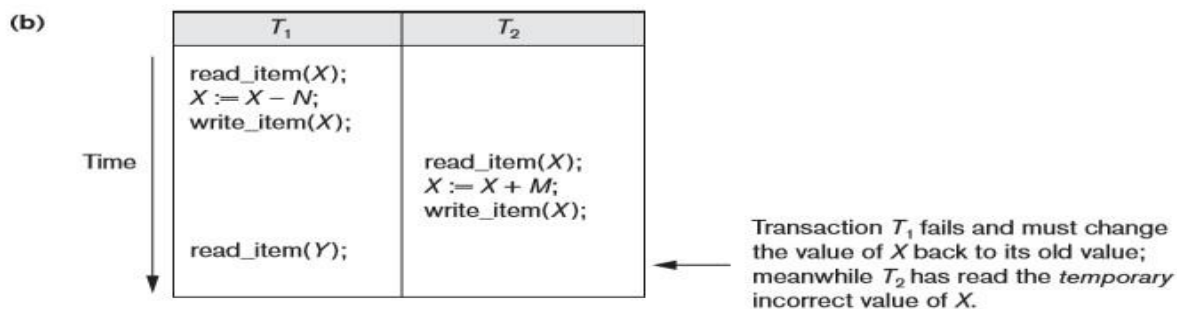
1) **The Lost Update Problem.** This problem occurs when two transactions that access the same database items have their operations interleaved in a way that makes the value of some database items incorrect. Suppose that transactions  $T_1$  and  $T_2$  are submitted at approximately the same time, and suppose that their operations are interleaved as shown in Figure 21.3(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  $M=4$  ( $T_2$  reserves 4 seats on  $X$ ), the final result should be  $X=79$ . However, in the interleaving of operations shown in Figure 21.3(a), it is  $X=84$  because the update in  $T_1$  that removed the five seats from  $X$  was lost.

**Figure 21.3**

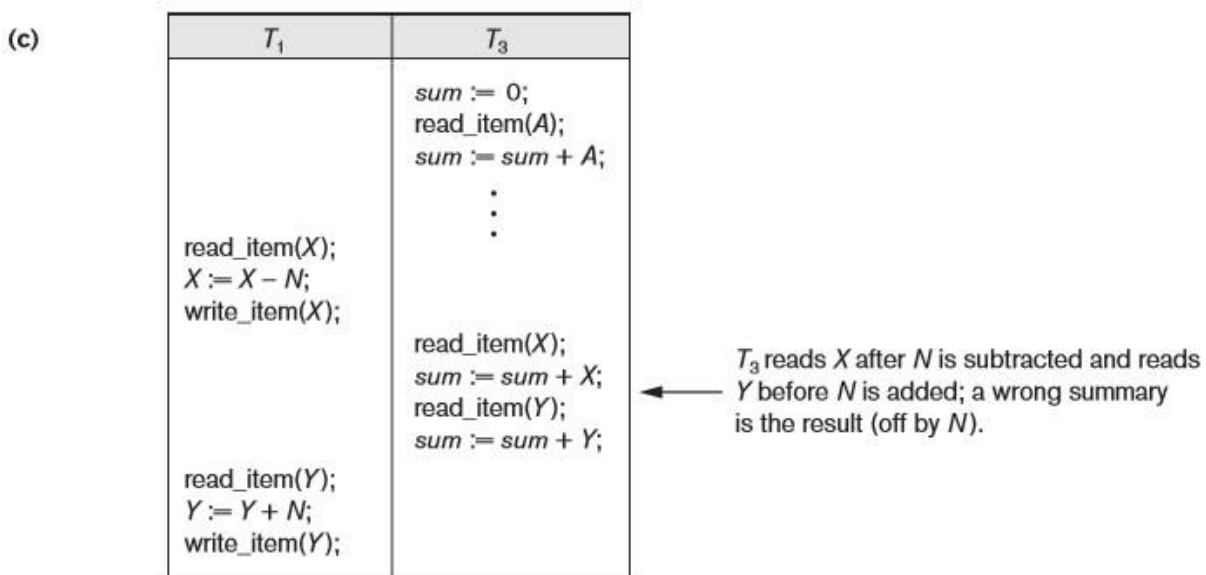
Some problems that occur when concurrent execution is uncontrolled. (a) The lost update problem. (b) The temporary update problem. (c) The incorrect summary problem.

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

**2) 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. Meanwhile, the updated item is accessed (read) by another transaction before it is changed back to its original value. Figure 21.3(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**.



**3) The Incorrect Summary Problem.** If one transaction is calculating an aggregate summary function on a number of database items while other transactions are updating some of these items, 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 21.3(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.



**4) The Unrepeatable Read Problem.** Another problem that may occur is called unrepeatable read, where a transaction T reads the same 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, and it may end up reading a different value for the item.

### **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 that the transaction does not have any effect on the database or any other transactions. In the first case, the transaction is said to be committed, whereas in the second case, the transaction is aborted. The DBMS must not permit some operations of a transaction T to be applied to the database while other operations of T are not, because the whole transaction is a logical unit of database processing. If a transaction fails after executing some of its operations but before executing all of them, the operations already executed must be undone and have no lasting effect.

**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. 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 could be programmed in the transaction itself, and in such a case would not be considered as a transaction failure.



4. **Concurrency control enforcement.** The concurrency control method (see Chapter 22) may decide to abort a transaction because it violates serializability (see Section 21.5), or it may abort one or more transactions to resolve a state of deadlock among several transactions (see Section 22.1.3). Transactions aborted because of serializability violations or deadlocks are typically restarted automatically at a later time.
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.

### ***Transaction and System Concepts***

In this section we discuss additional concepts relevant to transaction processing.

#### **Transaction States and Additional Operations**

A **transaction is an atomic unit of work** that should either be completed in its **entirety or not done at all**. For recovery purposes, the system needs to keep track of when each transaction starts, terminates, and commits or aborts. Therefore, the recovery manager of the DBMS needs to keep 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 21.5) 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 21.4**

State transition diagram illustrating the states for transaction execution.

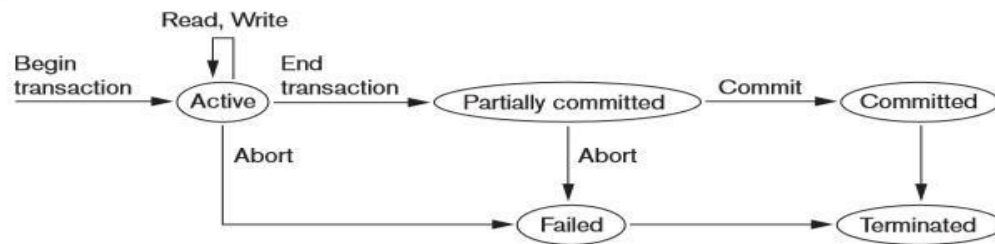


Figure 21.4 shows a state transition diagram that illustrates how a transaction moves through its execution states. A transaction goes into an **active state** immediately after it starts execution, where it can execute its READ and WRITE operations. When the transaction ends, it moves to the **partially 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). Once this check is successful, the transaction is said to have reached its **commit** point and enters the **committed state**. When a transaction is committed, it has concluded its execution successfully and all its changes must be recorded permanently in the database, even if a system failure occurs.

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

### The System Log

To be able to recover from failures that affect transactions, the system maintains a **log** to keep track of all transaction operations that affect the values of database items, as well as other transaction information that may be needed to permit recovery from failures.

The following are the types of entries—called **log-records**—that are written to the log file and the corresponding action for each log record. In these entries, *T* refers to a unique transaction-id that is generated automatically by the system for each transaction and that 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.

**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 must be permanently recorded in the database. The transaction then writes a commit record [commit, T] into the log. If a system failure occurs, we can 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.

***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 DBMS.

The following are the **ACID** properties:

- **Atomicity.** A transaction is an atomic unit of processing; it should either be performed in its entirety or not performed at all.
- **Consistency preservation.** A transaction should be consistency preserving, meaning that if it is completely executed from beginning to end without interference from other transactions, it should take the database from one consistent state to another.
- **Isolation.** A transaction should appear as though it is being executed in isolation from other transactions, even though many transactions are executing concurrently. 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.



**Characterizing Schedules Based on Recoverability**

When transactions are executing concurrently in an interleaved fashion, then the order of execution of operations from all the various transactions is known as a **schedule** (or history).

**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*.

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

Similarly, the schedule for Figure 21.3(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.

Characterizing Schedules Based on Recoverability

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  $(X)$ ).

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  $S_e$  and  $S_d$  that follow:

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

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

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

)

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

*Characterizing Schedules Based on Serializability*

In the previous section, we characterized schedules based on their recoverability properties. Now we characterize the types of schedules that are always considered to be correct when concurrent transactions are executing. Such schedules are known as **serializable schedules**.

Suppose that two users—for example, two airline reservations agents—submit to the DBMS transactions  $T_1$  and  $T_2$  in Figure 21.2 at approximately the same time. If no interleaving of operations is permitted, there are only two possible outcomes:

1. Execute all the operations of transaction  $T_1$  (in sequence) followed by all the operations of transaction  $T_2$  (in sequence).
2. Execute all the operations of transaction  $T_2$  (in sequence) followed by all the operations of transaction  $T_1$  (in sequence).

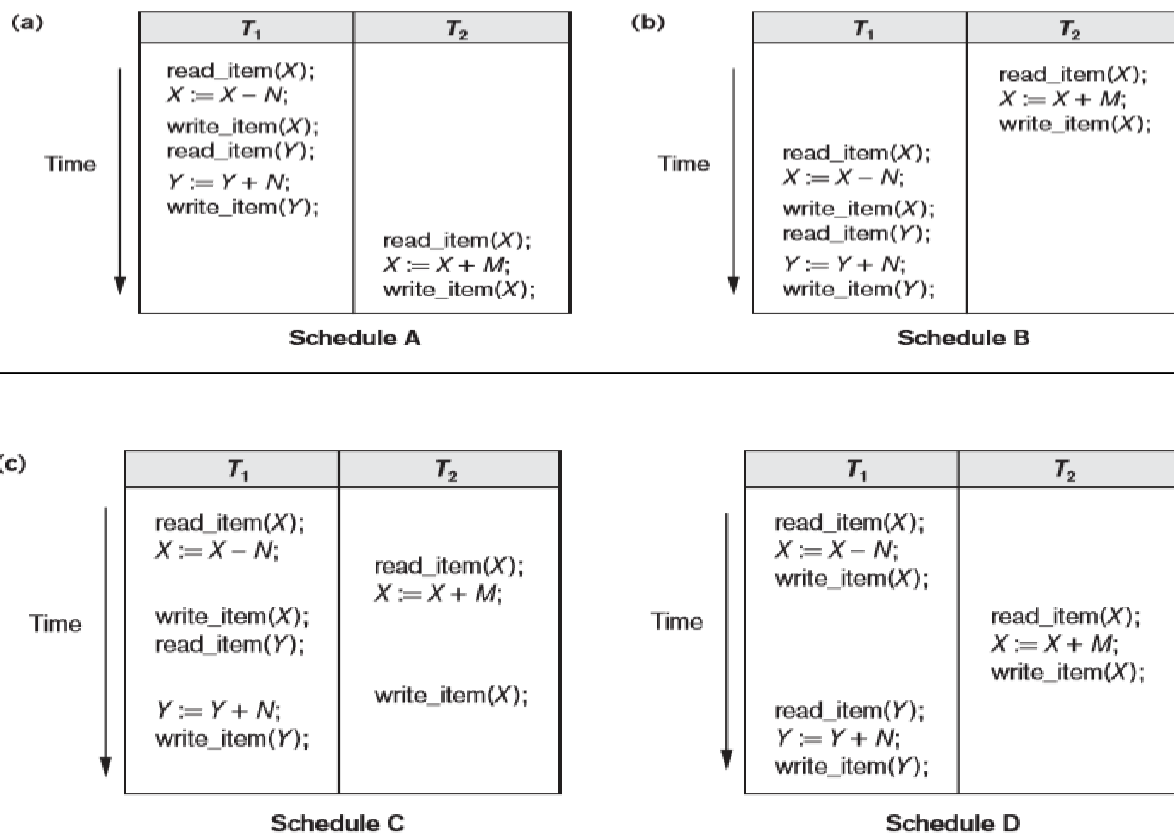
These two schedules—called **serial schedules**—are shown in Figure 21.5(a) and (b), respectively.

If interleaving of operations is allowed, there will be many possible orders in which the system can execute the individual operations of the transactions. Two possible schedules are shown in Figure 21.5(c). The concept of serializability of schedules is used to identify which schedules are correct when transaction executions have interleaving of their operations in the schedules. This section defines serializability and discusses how it may be used in practice.

**Figure 21.5**

Examples of serial and nonserial schedules involving transactions  $T_1$  and  $T_2$ . (a) Serial schedule A:  $T_1$  followed by  $T_2$ . (b) Serial schedule B:  $T_2$  followed by  $T_1$ .

(c) Two nonserial schedules C and D with interleaving of operations.



**Serial, Nonserial, and Conflict-Serializable Schedules**

Schedules A and B in Figure 21.5(a) and (b) are called **serial** because the operations of each transaction are executed consecutively, without any interleaved operations from the other transaction. In a serial schedule, entire transactions are performed in serial order: T1 and then T2 in Figure 21.5(a), and T2 and then T1 in Figure 21.5(b). Schedules C and D in Figure 21.5(c) are called **nonserial** because each sequence interleaves operations from the two transactions.

Formally, a schedule S is **serial** if, for every transaction T participating in the schedule, all the operations of T are executed consecutively in the schedule; otherwise, the schedule is called **nonserial**. Therefore, in a serial schedule, only one transaction at a time is active—the commit (or abort) of the active transaction initiates execution of the next transaction. No interleaving occurs in a serial schedule.

The problem with serial schedules is that they limit concurrency by prohibiting interleaving of operations. In a serial schedule, if a transaction waits for an I/O operation to complete, we cannot switch the CPU processor to another transaction, thus wasting valuable CPU processing time. Additionally, if some transaction T is quite long, the other transactions must wait for T to complete all its operations before starting. Hence, serial schedules are considered unacceptable in practice.

consider the schedules in Figure 21.5, and assume that the initial values of database items are  $X = 90$  and  $Y = 90$  and that  $N = 3$  and  $M = 2$ . After executing transactions T1 and T2, we would expect the database values to be  $X = 89$  and  $Y = 93$ , according to the meaning of the transactions. Sure enough, executing either of the serial schedules A or B gives the correct results. Now consider the nonserial schedules C and D. Schedule C (which is the same as Figure 21.3(a)) gives the results  $X = 92$  and  $Y = 93$ , in which the X value is erroneous, whereas schedule D gives the correct results.

Schedule C gives an erroneous result because of the lost update problem. Transaction T2 reads the value of X before it is changed by transaction T1, so only the effect of T2 on X is reflected in the database. The effect of T1 on X is lost, overwritten by T2, leading to the incorrect result for item X. However, some nonserial schedules give the correct expected result, such as schedule D. We would like to determine which of the nonserial schedules always give a correct result and which may give erroneous results. The concept used to characterize schedules in this manner is that of serializability of a schedule

The definition of **serializable schedule** is as follows: A schedule S of n transactions is **serializable** if it is equivalent to some serial schedule of the same n transactions.

There are several ways to define **schedule equivalence**. Two schedules are called **result equivalent** if they produce the same final state of the database. However, two different schedules may accidentally produce the same final state. For example, in Figure 21.6, schedules S1 and S2 will produce the same final database state if they execute on a database with an initial value of  $X = 100$ ; however, for other initial values of X, the schedules are not result equivalent.

**Figure 21.6**

Two schedules that are result equivalent for the initial value of  $X = 100$  but are not result equivalent in general.

$S_1$	$S_2$
<code>read_item(X);</code> <code>X := X + 10;</code> <code>write_item(X);</code>	<code>read_item(X);</code> <code>X := X * 1.1;</code> <code>write_item(X);</code>

The definition of **conflict equivalence** of schedules is as follows: Two schedules are said to be **conflict equivalent** if the order of any two conflicting operations is the same in both schedules. Two operations in a schedule are said to conflict if they belong to different transactions, access the same database item, and either both are `write_item` operations or one is a `write_item` and the other a `read_item`.

Using the notion of conflict equivalence, we define a schedule  $S$  to be **conflict serializable** if it is (conflict) equivalent to some serial schedule  $S$ . In such a case, we can reorder the nonconflicting operations in  $S$  until we form the equivalent serial schedule  $S$ . According to this definition, schedule  $D$  in Figure 21.5(c) is equivalent to the serial schedule  $A$  in Figure 21.5(a).

Schedule  $C$  in Figure 21.5(c) is not equivalent to either of the two possible serial schedules  $A$  and  $B$ , and hence is **not serializable**.

### Testing for Conflict Serializability of a Schedule

There is a simple algorithm for determining whether a particular schedule is conflict serializable or not.

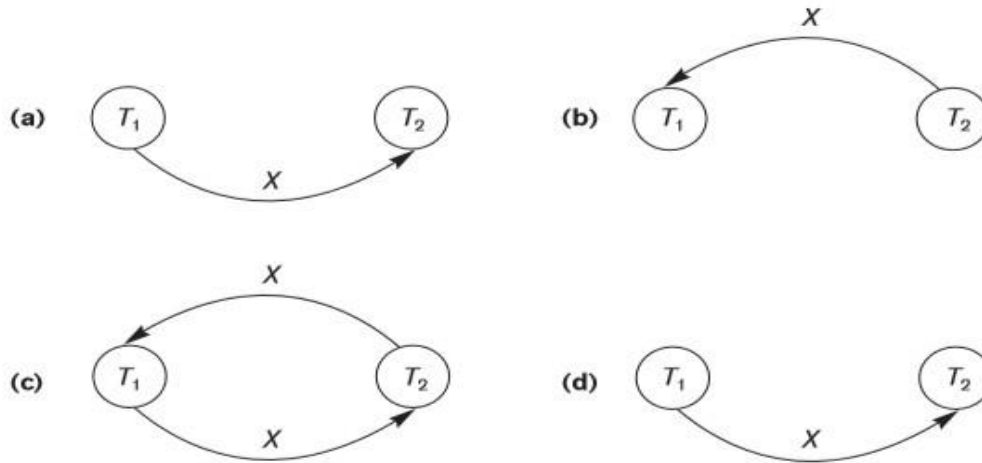
Algorithm 21.1 can be used to test a schedule for conflict serializability. The algorithm looks at only the `read_item` and `write_item` operations in a schedule to construct a precedence graph (or serialization graph), which is a directed graph  $G = (N, E)$  that consists of a set of nodes  $N = \{T_1, T_2, \dots, T_n\}$  and a set of directed edges  $E = \{e_1, e_2, \dots, e_m\}$ . There is one node in the graph for each transaction  $T_i$  in the schedule. Each edge  $e_i$  in the graph is of the form  $(T_j \rightarrow T_k)$ ,  $1 \leq j \leq n$ ,  $1 \leq k \leq n$ , where  $T_j$  is the starting node of  $e_i$  and  $T_k$  is the ending node of  $e_i$ . Such an edge from node  $T_j$  to node  $T_k$  is created by the algorithm if one of the operations in  $T_j$  appears in the schedule before some conflicting operation in  $T_k$ .

#### **Algorithm 21.1.** Testing Conflict Serializability of a Schedule $S$

1. For each transaction  $T_i$  participating in schedule  $S$ , create a node labeled  $T_i$  in the precedence graph.
2. For each case in  $S$  where  $T_j$  executes a `read_item(X)` after  $T_i$  executes a `write_item(X)`, create an edge  $(T_i \rightarrow T_j)$  in the precedence graph.
3. For each case in  $S$  where  $T_j$  executes a `write_item(X)` after  $T_i$  executes a `read_item(X)`, create an edge  $(T_i \rightarrow T_j)$  in the precedence graph.
4. For each case in  $S$  where  $T_j$  executes a `write_item(X)` after  $T_i$  executes a `write_item(X)`, create an edge  $(T_i \rightarrow T_j)$  in the precedence graph.
5. The schedule  $S$  is serializable if and only if the precedence graph has no cycles.



The precedence graph is constructed as described in Algorithm 21.1. If there is a **cycle** in the precedence graph, schedule S is **not (conflict) serializable**; if there is **no cycle**, S is **serializable**.



**Figure 21.7**

Constructing the precedence graphs for schedules A to D from Figure 21.5 to test for conflict serializability. (a) Precedence graph for serial schedule A. (b) Precedence graph for serial schedule B. (c) Precedence graph for schedule C (not serializable). (d) Precedence graph for schedule D (serializable, equivalent to schedule A).

### How Serializability Is Used for Concurrency Control

A serializable schedule gives the benefits of concurrent execution without giving up any correctness. In practice, it is quite difficult to test for the serializability of a schedule. The interleaving of operations from concurrent transactions—which are usually executed as processes by the operating system—is typically determined by the operating system scheduler, which allocates resources to all processes. Factors such as system load, time of transaction submission, and priorities of processes contribute to the ordering of operations in a schedule. Hence, it is difficult to determine how the operations of a schedule will be interleaved beforehand to ensure serializability.

### View Equivalence and View Serializability

**view equivalence**: less restrictive definition of equivalence of schedules is called **view equivalence**.

**view serializability**: Two schedules S and S' are said to be view equivalent if the following three conditions hold:

1. The same set of transactions participates in S and S', and S and S' include the same operations of those transactions.
2. For any operation  $r_i(X)$  of  $T_i$  in S, if the value of X read by the operation has been written by an operation  $w_j(X)$  of  $T_j$  (or if it is the original value of X before the schedule started), the same condition must hold for the value of X read by operation  $r_i(X)$  of  $T_i$  in S'.
3. If the operation  $w_k(Y)$  of  $T_k$  is the last operation to write item Y in S, then  $w_k(Y)$  of  $T_k$  must also be the last operation to write item Y in S'.

A schedule S is said to be view serializable if it is view equivalent to a serial schedule.

The definitions of conflict serializability and view serializability are similar if a condition known as the constrained write assumption (or **no blind writes**) holds on all transactions in the schedule.

A blind write is a write operation in a transaction T on an item X that is not dependent on the value of X, so it is not preceded by a read of X in the transaction T.

The definition of view serializability is less restrictive than that of conflict serializability under the unconstrained write assumption, where the value written by an operation  $w_i(X)$  in  $T_i$  can be independent of its old value from the database. This is possible when blind writes are allowed, and it is illustrated by the following schedule Sg of three transactions  $T_1: r_1(X); w_1(X); T_2: w_2(X);$  and  $T_3: w_3(X):$

Sg:  $r_1(X); w_2(X); w_1(X); w_3(X); c_1; c_2; c_3;$

In Sg the operations  $w_2(X)$  and  $w_3(X)$  are blind writes, since  $T_2$  and  $T_3$  do not read the value of X. The schedule Sg is view serializable, since it is view equivalent to the serial schedule  $T_1, T_2, T_3$ . However, Sg is not conflict serializable, since it is not conflict equivalent to any serial schedule.

### Other Types of Equivalence of Schedules

Serializability of schedules is sometimes considered to be too restrictive as a condition for ensuring the correctness of concurrent executions. Some applications can produce schedules that are correct by satisfying conditions less stringent than either conflict serializability or view serializability. An example is the type of transactions known as debit-credit transactions—for example, those that apply deposits and withdrawals to a data item whose value is the current balance of a bank account. The semantics of debit-credit operations is that they update the value of a data item X by either subtracting from or adding to the value of the data item. Because addition and subtraction operations are commutative—that is, they can be applied in any order—it is possible to produce correct schedules that are not serializable. For example, consider the following transactions, each of which may be used to transfer an amount of money between two bank accounts:

$T_1: r_1(X); X := X - 10; w_1(X); r_1(Y); Y := Y + 10; w_1(Y);$

$T_2: r_2(Y); Y := Y - 20; w_2(Y); r_2(X); X := X + 20; w_2(X);$

Consider the following nonserializable schedule Sh for the two transactions:

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

With the additional knowledge, or semantics, that the operations between each  $r_i(I)$  and  $w_i(I)$  are commutative, we know that the order of executing the sequences consisting of (read, update, write) is not important as long as each (read, update, write) sequence by a particular transaction  $T_i$  on a particular item I is not interrupted by conflicting operations. Hence, the schedule Sh is considered to be correct even though it is not serializable.

***Transaction Support in SQL***

In this section, we give a brief introduction to transaction support in SQL. The basic definition of an SQL transaction is similar to our already defined concept of a transaction. That is, it is a logical unit of work and is guaranteed to be atomic. A single SQL statement is always considered to be atomic—either it completes execution without an error or it fails and leaves the database unchanged.

With SQL, there is no explicit `Begin_Transaction` statement. Transaction initiation is done implicitly when particular SQL statements are encountered. However, every transaction must have an explicit end statement, which is either a `COMMIT` or a `ROLLBACK`. Every transaction has certain characteristics attributed to it. These characteristics are specified by a `SET TRANSACTION` statement in SQL. The characteristics are the access mode, the diagnostic area size, and the isolation level.

The access mode can be specified as `READ ONLY` or `READ WRITE`. The default is `READ WRITE`, unless the isolation level of `READ UNCOMMITTED` is specified (see below), in which case `READ ONLY` is assumed. A mode of `READ WRITE` allows select, update, insert, delete, and create commands to be executed. A mode of `READ ONLY`, as the name implies, is simply for data retrieval.

The diagnostic area size option, `DIAGNOSTIC SIZE n`, specifies an integer value `n`, which indicates the number of conditions that can be held simultaneously in the diagnostic area. These conditions supply feedback information (errors or exceptions) to the user or program on the `n` most recently executed SQL statement.

The isolation level option is specified using the statement `ISOLATION LEVEL <isolation>`, where the value for `<isolation>` can be `READ UNCOMMITTED`, `READ COMMITTED`, `REPEATABLE READ`, or `SERIALIZABLE`.<sup>15</sup> The default isolation level is `SERIALIZABLE`, although some systems use `READ COMMITTED` as their default. The use of the term `SERIALIZABLE` here is based on not allowing violations that cause dirty read, unrepeatable read, and phantoms,<sup>16</sup> and it is thus not identical to the way serializability was defined earlier in Section 21.5. If a transaction executes at a lower isolation level than `SERIALIZABLE`, then one or more of the following three violations may occur:

1. **Dirty read.** A transaction `T1` may read the update of a transaction `T2`, which has not yet committed. If `T2` fails and is aborted, then `T1` would have read a value that does not exist and is incorrect.
2. **Nonrepeatable read.** A transaction `T1` may read a given value from a table. If another transaction `T2` later updates that value and `T1` reads that value again, `T1` will see a different value.
3. **Phantoms.** A transaction `T1` may read a set of rows from a table, perhaps based on some condition specified in the SQL `WHERE`-clause. Now suppose that a transaction `T2` inserts a new row that also satisfies the `WHERE`-clause condition used in `T1`, into the table used by `T1`. If `T1` is repeated, then `T1` will see a phantom, a row that previously did not exist.

Table 21.1 summarizes the possible violations for the different isolation levels. An entry of Yes indicates that a violation is possible and an entry of No indicates that it is not possible. `READ UNCOMMITTED` is

the most forgiving, and SERIALIZABLE is the most restrictive in that it avoids all three of the problems mentioned above.

**Table 21.1** Possible Violations Based on Isolation Levels as Defined in SQL

Isolation Level	Type of Violation		
	Dirty Read	Nonrepeatable Read	Phantom
READ UNCOMMITTED	Yes	Yes	Yes
READ COMMITTED	No	Yes	Yes
REPEATABLE READ	No	No	Yes
SERIALIZABLE	No	No	No

A sample SQL transaction might look like the following:

```
EXEC SQL WHENEVER SQLERROR GOTO UNDO;
EXEC SQL SET TRANSACTION
    READ WRITE
    DIAGNOSTIC SIZE 5
    ISOLATION LEVEL SERIALIZABLE;
EXEC SQL INSERT INTO EMPLOYEE (Fname, Lname, Ssn, Dno, Salary)
    VALUES ('Robert', 'Smith', '991004321', 2, 35000);
EXEC SQL UPDATE EMPLOYEE
    SET Salary = Salary * 1.1 WHERE Dno = 2;
EXEC SQL COMMIT;
GOTO THE_END;
UNDO: EXEC SQL ROLLBACK;
THE_END: ... ;
```

The above transaction consists of first inserting a new row in the EMPLOYEE table and then updating the salary of all employees who work in department 2. If an error occurs on any of the SQL statements, the entire transaction is rolled back. This implies that any updated salary (by this transaction) would be restored to its previous value and that the newly inserted row would be removed.

## Concurrency Control in Databases

In this chapter we discuss a number of concurrency control techniques that are used to ensure the non interference or isolation property of concurrently executing transactions. Most of these techniques ensure serializability of schedules.

### *Two-Phase Locking Techniques for Concurrency Control*

Some of the main techniques used to control concurrent execution of transactions are based on the concept of locking data items. A **lock** is a variable associated with a data item that describes the status of the item with respect to possible operations that can be applied to it.

### Types of Locks and System Lock Tables

Several types of locks are used in concurrency control. To introduce locking concepts gradually, first we discuss **binary locks**, which are simple, but are also too restrictive for database concurrency control purposes, and so are not used in practice. Then we discuss **shared/exclusive** locks—also known as **read/write** locks—which provide more general locking capabilities and are used in practical database locking schemes. we describe an additional type of lock called a **certify lock**, and show how it can be used to improve performance of locking protocols.

**Binary Locks.** A binary lock can have two states or values: locked and unlocked (or 1 and 0, for simplicity). A distinct lock is associated with each database item X. If the value of the lock on X is 1, item X cannot be accessed by a database operation that requests the item. If the value of the lock on X is 0, the item can be accessed when requested, and the lock value is changed to 1. We refer to the current value (or state) of the lock associated with item X as  $\text{lock}(X)$ .

Two operations,  $\text{lock\_item}$  and  $\text{unlock\_item}$ , are used with binary locking. A transaction requests access to an item X by first issuing a  $\text{lock\_item}(X)$  operation. If  $\text{LOCK}(X) = 1$ , the transaction is forced to wait. If  $\text{LOCK}(X) = 0$ , it is set to 1 (the transaction locks the item) and the transaction is allowed to access item X. When the transaction is through using the item, it issues an  $\text{unlock\_item}(X)$  operation, which sets  $\text{LOCK}(X)$  back to 0 (unlocks the item) so that X may be accessed by other transactions. Hence, a binary lock enforces mutual exclusion on the data item. A description of the  $\text{lock\_item}(X)$  and  $\text{unlock\_item}(X)$  operations is shown in Figure 22.1.

```

lock_item(X):
  B:  if LOCK(X) = 0          (* item is unlocked *)
      then LOCK(X) ← 1      (* lock the item *)
      else
        begin
          wait (until LOCK(X) = 0
              and the lock manager wakes up the transaction);
          go to B
        end;
unlock_item(X):
  LOCK(X) ← 0;              (* unlock the item *)
  if any transactions are waiting
    then wakeup one of the waiting transactions;
  
```

**Figure 22.1**  
Lock and unlock operations for binary locks.



If the simple binary locking scheme described here is used, every transaction must obey the following rules:

1. A transaction  $T$  must issue the operation  $\text{lock\_item}(X)$  before any  $\text{read\_item}(X)$  or  $\text{write\_item}(X)$  operations are performed in  $T$ .
2. A transaction  $T$  must issue the operation  $\text{unlock\_item}(X)$  after all  $\text{read\_item}(X)$  and  $\text{write\_item}(X)$  operations are completed in  $T$ .
3. A transaction  $T$  will not issue a  $\text{lock\_item}(X)$  operation if it already holds the lock on item  $X$ .<sup>1</sup>
4. A transaction  $T$  will not issue an  $\text{unlock\_item}(X)$  operation unless it already holds the lock on item  $X$ .

**Shared/Exclusive (or Read/Write) Locks:** The preceding binary locking scheme is too restrictive for database items because at most, one transaction can hold a lock on a given item. We should allow several transactions to access the same item  $X$  if they all access  $X$  for reading purposes only. This is because read operations on the same item by different transactions are not conflicting. However, if a transaction is to write an item  $X$ , it must have exclusive access to  $X$ .

For this purpose, a different type of lock called a **multiple-mode** lock is used. In this scheme—called **shared/exclusive** or **read/write** locks—there are **three locking operations:  $\text{read\_lock}(X)$ ,  $\text{write\_lock}(X)$ , and  $\text{unlock}(X)$** . A lock associated with an item  $X$ ,  $\text{LOCK}(X)$ , now has three possible states: read-locked, write-locked, or unlocked. A read-locked item is also called share-locked because other transactions are allowed to read the item, whereas a write-locked item is called exclusive-locked because a single transaction exclusively holds the lock on the item.

When we use the shared/exclusive locking scheme, the system must enforce the following rules:

1. A transaction  $T$  must issue the operation  $\text{read\_lock}(X)$  or  $\text{write\_lock}(X)$  before any  $\text{read\_item}(X)$  operation is performed in  $T$ .
2. A transaction  $T$  must issue the operation  $\text{write\_lock}(X)$  before any  $\text{write\_item}(X)$  operation is performed in  $T$ .
3. A transaction  $T$  must issue the operation  $\text{unlock}(X)$  after all  $\text{read\_item}(X)$  and  $\text{write\_item}(X)$  operations are completed in  $T$ .
4. A transaction  $T$  will not issue a  $\text{read\_lock}(X)$  operation if it already holds a read (shared) lock or a write (exclusive) lock on item  $X$ . This rule may be relaxed, as we discuss shortly.
5. A transaction  $T$  will not issue a  $\text{write\_lock}(X)$  operation if it already holds a read (shared) lock or write (exclusive) lock on item  $X$ . This rule may also be relaxed, as we discuss shortly.
6. A transaction  $T$  will not issue an  $\text{unlock}(X)$  operation unless it already holds a read (shared) lock or a write (exclusive) lock on item  $X$ .

**Conversion of Locks.** Sometimes it is desirable to relax conditions 4 and 5 in the preceding list in order to allow **lock conversion**; that is, a transaction that already holds a lock on item  $X$  is allowed under certain conditions to **convert the lock from one locked state to another**. For example, it is possible for a transaction  $T$  to issue a  $\text{read\_lock}(X)$  and then later to **upgrade** the lock by issuing a  $\text{write\_lock}(X)$  operation. If  $T$  is the only transaction holding a read lock on  $X$  at the time it issues the  $\text{write\_lock}(X)$  operation, the lock can be upgraded; otherwise, the transaction must wait. It is also possible for a transaction  $T$  to issue a  $\text{write\_lock}(X)$  and then later to **downgrade** the lock by issuing a  $\text{read\_lock}(X)$  operation.

```

read_lock(X):
B: if LOCK(X) = "unlocked"
    then begin LOCK(X) ← "read-locked";
              no_of_reads(X) ← 1
            end
    else if LOCK(X) = "read-locked"
    then no_of_reads(X) ← no_of_reads(X) + 1
    else begin
        wait (until LOCK(X) = "unlocked"
              and the lock manager wakes up the transaction);
        go to B
      end;
write_lock(X):
B: if LOCK(X) = "unlocked"
    then LOCK(X) ← "write-locked"
    else begin
        wait (until LOCK(X) = "unlocked"
              and the lock manager wakes up the transaction);
        go to B
      end;
unlock (X):
  if LOCK(X) = "write-locked"
  then begin LOCK(X) ← "unlocked";
            wakeup one of the waiting transactions, if any
          end
  else if LOCK(X) = "read-locked"
  then begin
        no_of_reads(X) ← no_of_reads(X) - 1;
        if no_of_reads(X) = 0
        then begin LOCK(X) = "unlocked";
                  wakeup one of the waiting transactions, if any
                end
      end;
  end;

```

**Figure 22.2**  
Locking and unlocking  
operations for two-  
mode (read-write or  
shared-exclusive)  
locks.

### Guaranteeing Serializability by Two-Phase Locking

#### Two-Phase Locking

A transaction is said to follow the two-phase locking protocol if all locking operations (read\_lock, write\_lock) precede the first unlock operation in the transaction. Such a transaction can be divided into **two phases**.

**Growing (first) phase:** an expanding or growing (first) phase, during which new locks on items can be acquired but none can be released.

**shrinking (second) phase :** a shrinking (second) phase, during which existing locks can be released but no new locks can be acquired.

If lock conversion is allowed, then **upgrading** of locks (from read-locked to write-locked) must be done during the expanding phase, and **downgrading** of locks (from write-locked to read-locked) must be done in the shrinking phase. Hence, a read\_lock(X) operation that downgrades an already held write lock on X can appear only in the shrinking phase.

Transactions T1 and T2 in Figure 22.3(a) do not follow the two-phase locking protocol because the write\_lock(X) operation follows the unlock(Y) operation in T1, and similarly the write\_lock(Y) operation follows the unlock(X) operation in T2.

If we enforce two-phase locking, the transactions can be rewritten as T1 and T2, as shown in Figure 22.4. Now, the schedule shown in Figure 22.3(c) is not permitted for T1 and T2 (with their modified order of locking and unlocking operations).

(c)

	$T_1$	$T_2$
Time	read_lock(Y); read_item(Y); unlock(Y);   write_lock(X); read_item(X); $X := X + Y$ ; write_item(X); unlock(X);	read_lock(X); read_item(X); unlock(X); write_lock(Y); read_item(Y); $Y := X + Y$ ; write_item(Y); unlock(Y);

Result of schedule  $S$ :  
 $X=50, Y=50$   
 (nonserializable)

**Figure 22.3** Transactions that do not obey two-phase locking. (a) Two transactions  $T_1$  and  $T_2$ . (b) Results of possible serial schedules of  $T_1$  and  $T_2$ . (c) A nonserializable schedule  $S$  that uses locks.

$T_1'$	$T_2'$
read_lock(Y); read_item(Y); write_lock(X); unlock(Y) read_item(X); $X := X + Y$ ; write_item(X); unlock(X);	read_lock(X); read_item(X); write_lock(Y); unlock(X) read_item(Y); $Y := X + Y$ ; write_item(Y); unlock(Y);

**Figure 22.4**  
Transactions  $T_1'$  and  $T_2'$ , which are the same as  $T_1$  and  $T_2$  in Figure 22.3, but follow the two-phase locking protocol. Note that they can produce a deadlock.

Variations(Types of two phase locking) of two-phase locking (2PL).

a) Basic, b) Conservative, c) Strict, and d) Rigorous Two-Phase Locking:

There are a number of variations of two-phase locking (2PL). The technique just described is known as **basic 2PL**.

A variation known as **conservative 2PL** (or static 2PL) requires a transaction to lock all the items it accesses before the transaction begins execution, by predeclaring its read-set and write-set. 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 it writes. If any of the predeclared items needed cannot be locked, the transaction does not lock any item; instead, it waits until all the items are available for locking. **Conservative 2PL is a deadlock-free protocol.**

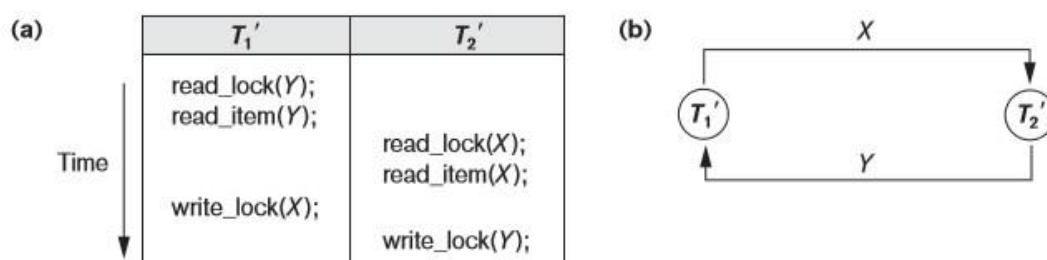
In practice, the most popular variation of 2PL is **strict 2PL**, which guarantees strict schedules. In this variation, a transaction T does not release any of its exclusive (write) locks until after it commits or aborts. Hence, no other transaction can read or write an item that is written by T unless T has committed, leading to a strict schedule for recoverability. **Strict 2PL is not deadlock-free.**

A more restrictive variation of strict 2PL is **rigorous 2PL**, which also guarantees strict schedules. In this variation, a transaction T does not release any of its locks (exclusive or shared) until after it commits or aborts, and so it is easier to implement than strict 2PL.

### Dealing with Deadlock and Starvation

Deadlock occurs when each transaction T in a set of two or more transactions is waiting for some item that is locked by some other transaction T in the set. Hence, each transaction in the set is in a waiting queue, waiting for one of the other transactions in the set to release the lock on an item. But because the other transaction is also waiting, it will never release the lock.

A simple example is shown in Figure 22.5(a), where the two transactions T<sub>1</sub> and T<sub>2</sub> are **deadlocked** in a partial schedule; T<sub>1</sub> is in the waiting queue for X, which is locked by T<sub>2</sub>, while T<sub>2</sub> is in the waiting queue for Y, which is locked by T<sub>1</sub>. Meanwhile, neither T<sub>1</sub> nor T<sub>2</sub> nor any other transaction can access items X and Y.



**Figure 22.5**

Illustrating the deadlock problem. (a) A partial schedule of  $T_1'$  and  $T_2'$  that is in a state of deadlock. (b) A wait-for graph for the partial schedule in (a).

**Deadlock Prevention Protocols.** One way to prevent deadlock is to use a deadlock prevention protocol. One deadlock prevention protocol, which is used in conservative two-phase locking, requires that every transaction lock all the items it needs in advance (which is generally not a practical assumption)—if any of the items cannot be obtained, none of the items are locked. Rather, the transaction waits and then tries again to lock all the items it needs. Obviously this solution further limits concurrency. A second protocol, which also limits concurrency, involves ordering all the items in the database and making sure that a transaction that needs several items will lock them according to that order. This requires that the programmer (or the system) is aware of the chosen order of the items, which is also not practical in the database context.

A number of other deadlock prevention schemes have been proposed that make a decision about what to do with a transaction involved in a possible deadlock situation.

Some of these techniques use the concept of transaction **timestamp**  $TS(T)$ , which is a unique identifier assigned to each transaction. The timestamps are typically based on the order in which transactions are started; hence, if transaction  $T_1$  starts before transaction  $T_2$ , then  $TS(T_1) < TS(T_2)$ .

**Two schemes** that prevent deadlock are called **wait-die** and **woundwait**. Suppose that transaction  $T_i$  tries to lock an item  $X$  but is not able to because  $X$  is locked by some other transaction  $T_j$  with a conflicting lock. The rules followed by these schemes are:

- **Wait-die.** If  $TS(T_i) < TS(T_j)$ , then ( $T_i$  older than  $T_j$ )  $T_i$  is allowed to wait; otherwise ( $T_i$  younger than  $T_j$ ) abort  $T_i$  ( $T_i$  dies) and restart it later with the same timestamp.

- **Wound-wait.** If  $TS(T_i) < TS(T_j)$ , then ( $T_i$  older than  $T_j$ ) abort  $T_j$  ( $T_i$  wounds  $T_j$ ) and restart it later with the same timestamp; otherwise ( $T_i$  younger than  $T_j$ )  $T_i$  is allowed to wait.

Another group of protocols that prevent deadlock do not require timestamps. These include the **no waiting (NW)** and **cautious waiting (CW) algorithms**.

In the **no waiting** algorithm, if a transaction is unable to obtain a lock, it is immediately aborted and then restarted after a certain time delay without checking whether a deadlock will actually occur or not. In this case, no transaction ever waits, so no deadlock will occur.

The **cautious waiting** algorithm was proposed to try to reduce the number of needless aborts/restarts. Suppose that transaction  $T_i$  tries to lock an item  $X$  but is not able to do so because  $X$  is locked by some other transaction  $T_j$  with a conflicting lock. The cautious waiting rules are as follows:

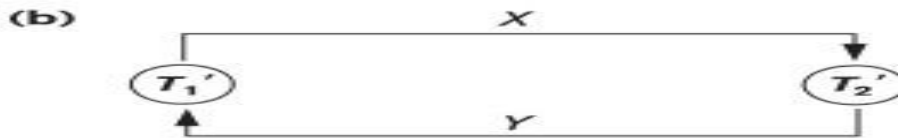
- **Cautious waiting.** If  $T_j$  is not blocked (not waiting for some other locked item), then  $T_i$  is blocked and allowed to wait; otherwise abort  $T_i$ .

**Deadlock Detection:** A second, more practical approach to dealing with deadlock is deadlock detection, where the system checks if a state of deadlock actually exists.



A simple way to **detect a state of deadlock** is for the system to construct and maintain a **wait-for graph**. One node is created in the wait-for graph for each transaction that is currently executing. Whenever a transaction  $T_i$  is waiting to lock an item  $X$  that is currently locked by a transaction  $T_j$ , a directed edge ( $T_i \rightarrow T_j$ ) is created in the wait-for graph. When  $T_j$  releases the lock(s) on the items that  $T_i$  was waiting for, the directed edge is dropped from the wait-for graph. **We have a state of deadlock if and only if the wait-for graph has a cycle.**

Figure 22.5(b)(below figure) shows the wait-for graph for the (partial) schedule shown in Figure 22.5(a).



If the system is in a state of deadlock, some of the transactions causing the deadlock must be aborted. Choosing which transactions to abort is known as **victim selection**.

**Timeouts.** Another simple scheme to deal with deadlock is the use of timeouts. This method is practical because of its low overhead and simplicity. In this method, if a transaction waits for a period longer than a system-defined timeout period, the system assumes that the transaction may be deadlocked and aborts it—regardless of whether a deadlock actually exists or not.

**Starvation.** Another problem that may occur when we use locking is **starvation**, which occurs when a transaction cannot proceed for an indefinite period of time while other transactions in the system continue normally. This may occur if the waiting scheme for locked items is unfair, giving priority to some transactions over others. One solution for starvation is to have a fair waiting scheme, such as using a **first-come-first-served queue**; transactions are enabled to lock an item in the order in which they originally requested the lock. Another scheme allows some transactions to have priority over others but increases the priority of a transaction the longer it waits, until it eventually gets the highest priority and proceeds. Starvation can also occur because of victim selection if the algorithm selects the same transaction as victim repeatedly, thus causing it to abort and never finish execution. The algorithm can use higher priorities for transactions that have been aborted multiple times to avoid this problem. **The wait-die and wound-wait schemes discussed previously avoid starvation**, because they restart a transaction that has been aborted with its same original timestamp, so the possibility that the same transaction is aborted repeatedly is slim.

### ***Concurrency Control Based on Timestamp Ordering***

A different approach that guarantees serializability involves using transaction timestamps to order transaction execution for an equivalent serial schedule. we discuss how serializability is enforced by ordering transactions based on their timestamps.

### Timestamps

a timestamp is a unique identifier created by the DBMS to identify a transaction. Typically, timestamp values are assigned in the order in which the transactions are submitted to the system, so a timestamp can be thought of as the transaction start time. We will refer to the timestamp of transaction T as  $TS(T)$ . Concurrency control techniques based on timestamp ordering do not use locks; hence, deadlocks cannot occur.

Timestamps can be generated in several ways. One possibility is to use a **counter** that is incremented each time its value is assigned to a transaction. The transaction timestamps are numbered 1, 2, 3, ... in this scheme. A computer counter has a finite maximum value, so the system must periodically reset the counter to zero when no transactions are executing for some short period of time.

Another way to implement timestamps is to use the **current date/time** value of the system clock and ensure that no two timestamp values are generated during the same tick of the clock.

### The Timestamp Ordering Algorithm

The idea for this scheme is to order the transactions based on their timestamps. A schedule in which the transactions participate is then serializable, and the only equivalent serial schedule permitted has the transactions in order of their timestamp values. This is called **timestamp ordering (TO)**. In timestamp ordering, however, the schedule is equivalent to the particular serial order corresponding to the order of the transaction timestamps. The algorithm must ensure that, for each item accessed by conflicting operations in the schedule, the order in which the item is accessed does not violate the timestamp order. To do this, the algorithm associates with each database item X two timestamp (TS) values:

1. **read\_TS(X)**. The read timestamp of item X is the largest timestamp among all the timestamps of transactions that have successfully read item X—that is,  $read\_TS(X) = TS(T)$ , where T is the youngest transaction that has read X successfully.
2. **write\_TS(X)**. The write timestamp of item X is the largest of all the timestamps of transactions that have successfully written item X—that is,  $write\_TS(X) = TS(T)$ , where T is the youngest transaction that has written X successfully.

**Basic Timestamp Ordering (TO):** Whenever some transaction T tries to issue a  $read\_item(X)$  or a  $write\_item(X)$  operation, the basic TO algorithm compares the timestamp of T with  $read\_TS(X)$  and  $write\_TS(X)$  to ensure that the timestamp order of transaction execution is not violated. If this order is violated, then transaction T is aborted and resubmitted to the system as a new transaction with a new timestamp. If T is aborted and rolled back, any transaction T1 that may have used a value written by T must also be rolled back. Similarly, any transaction T2 that may have used a value written by T1 must also be rolled back, and so on. This effect is known as **cascading rollback**.

We first describe the basic TO algorithm here. The concurrency control algorithm must check whether conflicting operations violate the timestamp ordering in the following two cases:

**1. Whenever a transaction T issues a write\_item(X) operation, the following is checked:**

a. If  $\text{read\_TS}(X) > \text{TS}(T)$  or if  $\text{write\_TS}(X) > \text{TS}(T)$ , then abort and roll back T and reject the operation. This should be done because some younger transaction with a timestamp greater than  $\text{TS}(T)$ —and hence after T in the timestamp ordering—has already read or written the value of item X before T had a chance to write X, thus violating the timestamp ordering.

b. If the condition in part (a) does not occur, then execute the write\_item(X) operation of T and set  $\text{write\_TS}(X)$  to  $\text{TS}(T)$ .

**2. Whenever a transaction T issues a read\_item(X) operation, the following is checked:**

a. If  $\text{write\_TS}(X) > \text{TS}(T)$ , then abort and roll back T and reject the operation. This should be done because some younger transaction with timestamp greater than  $\text{TS}(T)$ —and hence after T in the timestamp ordering—has already written the value of item X before T had a chance to read X.

b. If  $\text{write\_TS}(X) \leq \text{TS}(T)$ , then execute the read\_item(X) operation of T and set  $\text{read\_TS}(X)$  to the larger of  $\text{TS}(T)$  and the current  $\text{read\_TS}(X)$ .

**Strict Timestamp Ordering (TO):** A variation of basic TO called strict TO ensures that the schedules are both strict (for easy recoverability) and (conflict) serializable. In this variation, a transaction T that issues a read\_item(X) or write\_item(X) such that  $\text{TS}(T) > \text{write\_TS}(X)$  has its read or write operation delayed until the transaction T that wrote the value of X (hence  $\text{TS}(T) = \text{write\_TS}(X)$ ) has committed or aborted. To implement this algorithm, it is necessary to simulate the locking of an item X that has been written by transaction T until T is either committed or aborted. This algorithm does not cause deadlock, since T waits for T only if  $\text{TS}(T) > \text{TS}(T)$ .

**Thomas's Write Rule:** A modification of the basic TO algorithm, known as Thomas's write rule, does not enforce conflict serializability, but it rejects fewer write operations by modifying the checks for the write\_item(X) operation as follows:

1. If  $\text{read\_TS}(X) > \text{TS}(T)$ , then abort and roll back T and reject the operation.

2. If  $\text{write\_TS}(X) > \text{TS}(T)$ , then do not execute the write operation but continue processing. This is because some transaction with timestamp greater than  $\text{TS}(T)$ —and hence after T in the timestamp ordering—has already written the value of X. Thus, we must ignore the write\_item(X) operation of T because it is already outdated and obsolete. Notice that any conflict arising from this situation would be detected by case (1).

3. If neither the condition in part (1) nor the condition in part (2) occurs, then execute the write\_item(X) operation of T and set  $\text{write\_TS}(X)$  to  $\text{TS}(T)$ .

---

**Multiversion Concurrency Control Techniques**

Other protocols for concurrency control keep the old values of a data item when the item is updated. These are known as multiversion concurrency control, because several versions (values) of an item are maintained. When a transaction requires access to an item, an appropriate version is chosen to maintain the serializability of the currently executing schedule, if possible. The idea is that some read operations that would be rejected in other techniques can still be accepted by reading an older version of the item to maintain serializability. When a transaction writes an item, it writes a new version and the old version(s) of the item are retained. Some multiversion concurrency control algorithms use the concept of view serializability rather than conflict serializability. An obvious drawback of multiversion techniques is that more storage is needed to maintain multiple versions of the database items.

Several multiversion concurrency control schemes have been proposed. We discuss two schemes here, one based on **timestamp ordering** and the other based on **2PL**.

**Multiversion Technique Based on Timestamp Ordering**

In this method, several versions  $X_1, X_2, \dots, X_k$  of each data item  $X$  are maintained. For each version, the value of version  $X_i$  and the following two timestamps are kept:

1.  $read\_TS(X_i)$ . The read timestamp of  $X_i$  is the largest of all the timestamps of transactions that have successfully read version  $X_i$ .
2.  $write\_TS(X_i)$ . The write timestamp of  $X_i$  is the timestamp of the transaction that wrote the value of version  $X_i$ .

Whenever a transaction  $T$  is allowed to execute a  $write\_item(X)$  operation, a new version  $X_{k+1}$  of item  $X$  is created, with both the  $write\_TS(X_{k+1})$  and the  $read\_TS(X_{k+1})$  set to  $TS(T)$ . Correspondingly, when a transaction  $T$  is allowed to read the value of version  $X_i$ , the value of  $read\_TS(X_i)$  is set to the larger of the current  $read\_TS(X_i)$  and  $TS(T)$ .

To ensure serializability, the following rules are used:

1. If transaction  $T$  issues a  $write\_item(X)$  operation, and version  $i$  of  $X$  has the highest  $write\_TS(X_i)$  of all versions of  $X$  that is also less than or equal to  $TS(T)$ , and  $read\_TS(X_i) > TS(T)$ , then abort and roll back transaction  $T$ ; otherwise, create a new version  $X_j$  of  $X$  with  $read\_TS(X_j) = write\_TS(X_j) = TS(T)$ .
2. If transaction  $T$  issues a  $read\_item(X)$  operation, find the version  $i$  of  $X$  that has the highest  $write\_TS(X_i)$  of all versions of  $X$  that is also less than or equal to  $TS(T)$ ; then return the value of  $X_i$  to transaction  $T$ , and set the value of  $read\_TS(X_i)$  to the larger of  $TS(T)$  and the current  $read\_TS(X_i)$ .

**Multiversion Two-Phase Locking Using Certify Locks**

In this multiple-mode locking scheme, there are **three locking modes** for an item: **read**, **write**, and **certify**, instead of just the two modes (**read**, **write**) discussed previously. Hence, the state of  $LOCK(X)$  for an item  $X$  can be one of **read-locked**, **write locked**, **certify-locked**, or **unlocked**. In the standard locking

scheme, with only read and write locks, a write lock is an exclusive lock. We can describe the relationship between read and write locks in the standard scheme by means of the lock compatibility table shown in Figure 22.6(a). An entry of Yes means that if a transaction T holds the type of lock

(a)		Read	Write
Read	Yes	No	
Write	No	No	

(b)		Read	Write	Certify
Read	Yes	Yes	No	
Write	Yes	No	No	
Certify	No	No	No	

**Figure 22.6**

Lock compatibility tables.  
 (a) A compatibility table for read/write locking scheme.  
 (b) A compatibility table for read/write/certify locking scheme.

specified in the column header on item X and if transaction T requests the type of lock specified in the row header on the same item X, then T can obtain the lock because the locking modes are compatible. On the other hand, an entry of No in the table indicates that the locks are not compatible, so T must wait until T releases the lock.

In the standard locking scheme, once a transaction obtains a write lock on an item, no other transactions can access that item. However, once T is ready to commit, it must obtain a certify lock on all items that it currently holds write locks on before it can commit. The certify lock is not compatible with read locks, so the transaction may have to delay its commit until all its write-locked items are released by any reading transactions in order to obtain the certify locks. Once the certify locks—which are exclusive locks—are acquired, the committed version X of the data item is set to the value of version X, version X is discarded, and the certify locks are then released. The lock compatibility table for this scheme is shown in Figure 22.6(b).

In this multiversion 2PL scheme, reads can proceed concurrently with a single write operation—an arrangement not permitted under the standard 2PL schemes.

### ***Validation (Optimistic) Concurrency Control Techniques***

In optimistic concurrency control techniques, also known as **validation** or **certification** techniques, no checking is done while the transaction is executing. Several theoretical concurrency control methods are based on the validation technique

At the end of transaction execution, a validation phase checks whether any of the transaction's updates violate serializability. Certain information needed by the validation phase must be kept by the system. If serializability is not violated, the transaction is committed and the database is updated from the local copies; otherwise, the transaction is aborted and then restarted later.



There are three phases for this concurrency control protocol:

1. **Read phase.** A transaction can read values of committed data items from the database. However, updates are applied only to local copies (versions) of the data items kept in the transaction workspace.
2. **Validation phase.** Checking is performed to ensure that serializability will not be violated if the transaction updates are applied to the database.
3. **Write phase.** If the validation phase is successful, the transaction updates are applied to the database; otherwise, the updates are discarded and the transaction is restarted.

The validation phase for  $T_i$  checks that, for each such transaction  $T_j$  that is either committed or is in its validation phase, one of the following conditions holds:

1. Transaction  $T_j$  completes its write phase before  $T_i$  starts its read phase.
2.  $T_i$  starts its write phase after  $T_j$  completes its write phase, and the read\_set of  $T_i$  has no items in common with the write\_set of  $T_j$ .
3. Both the read\_set and write\_set of  $T_i$  have no items in common with the write\_set of  $T_j$ , and  $T_j$  completes its read phase before  $T_i$  completes its read phase.

### ***Granularity of Data Items and Multiple Granularity Locking***

All concurrency control techniques assume that the database is formed of a number of named data items. A database item could be chosen to be one of the following:

: ■ A database record ■ A field value of a database record ■ A disk block ■ A whole file ■ The whole database

The granularity can affect the performance of concurrency control and recovery.

### **Granularity Level Considerations for Locking**

The size of data items is often called the data item **granularity**. **Fine granularity** refers to small item sizes, whereas **coarse granularity** refers to large item sizes. Several tradeoffs must be considered in choosing the data item size.

First, notice that the larger the data item size is, the lower the degree of concurrency permitted.

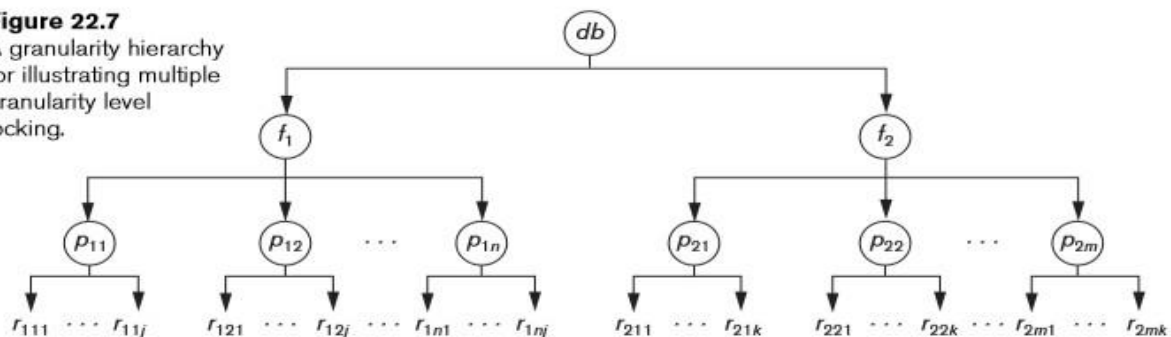
On the other hand, the smaller the data item size is, the more the number of items in the database. Because every item is associated with a lock, the system will have a larger number of active locks to be handled by the lock manager.

**Multiple Granularity Level Locking**

Since the best granularity size depends on the given transaction, it seems appropriate that a database system should support multiple levels of granularity, where the granularity level can be different for various mixes of transactions. Figure 22.7 shows a simple granularity hierarchy with a database containing two files, each file containing several disk pages, and each page containing several records. This can be used to illustrate a multiple granularity level 2PL protocol, where a lock can be requested at any level. However, additional types of locks will be needed to support such a protocol efficiently.

**Figure 22.7**

A granularity hierarchy for illustrating multiple granularity level locking.



Consider the following scenario, with only shared and exclusive lock types, that refers to the example in Figure 22.7. Suppose transaction T1 wants to update all the records in file f1, and T1 requests and is granted an exclusive lock for f1. Then all of f1's pages (p11 through p1n)—and the records contained on those pages—are locked in exclusive mode. This is beneficial for T1 because setting a single file-level lock is more efficient than setting n page-level locks or having to lock each individual record. Now suppose another transaction T2 only wants to read record r1nj from page p1n of file f1; then T2 would request a shared record-level lock on r1nj. However, the database system (that is, the transaction manager or more specifically the lock manager) must verify the compatibility of the requested lock with already held locks. One way to verify this is to traverse the tree from the leaf r1nj to p1n to f1 to db. If at any time a conflicting lock is held on any of those items, then the lock request for r1nj is denied and T2 is blocked and must wait. This traversal would be fairly efficient.

To make multiple granularity level locking practical, additional types of locks, called **intention locks**, are needed. The idea behind intention locks is for a transaction to indicate, along the path from the root to the desired node, what type of lock (shared or exclusive) it will require from one of the node's descendants. There are three types of intention locks:

1. Intention-shared (IS) indicates that one or more shared locks will be requested on some descendant node(s).
2. Intention-exclusive (IX) indicates that one or more exclusive locks will be requested on some descendant node(s).

3. Shared-intention-exclusive (SIX) indicates that the current node is locked in shared mode but that one or more exclusive locks will be requested on some descendant node(s).

The compatibility table of the three intention locks, and the shared and exclusive locks, is shown in Figure 22.8. Besides the introduction of the three types of intention locks, an appropriate locking protocol must be used. The **multiple granularity locking (MGL)** protocol consists of the following rules:

1. The lock compatibility (based on Figure 22.8) must be adhered to.
2. The root of the tree must be locked first, in any mode.
3. A node  $N$  can be locked by a transaction  $T$  in S or IS mode only if the parent node  $N$  is already locked by transaction  $T$  in either IS or IX mode.
4. A node  $N$  can be locked by a transaction  $T$  in X, IX, or SIX mode only if the parent of node  $N$  is already locked by transaction  $T$  in either IX or SIX mode.
5. A transaction  $T$  can lock a node only if it has not unlocked any node (to enforce the 2PL protocol).

6. A transaction  $T$  can unlock a node,  $N$ , only if none of the children of node  $N$  are currently locked by  $T$ .

Rule 1 simply states that conflicting locks cannot be granted. Rules 2, 3, and 4 state the conditions when a transaction may lock a given node in any of the lock modes. Rules 5 and 6 of the MGL protocol enforce 2PL rules to produce serializable schedules.

	IS	IX	S	SIX	X
IS	Yes	Yes	Yes	Yes	No
IX	Yes	Yes	No	No	No
S	Yes	No	Yes	No	No
SIX	Yes	No	No	No	No
X	No	No	No	No	No

**Figure 22.8**

Lock compatibility matrix for multiple granularity locking.

## Introduction to Database Recovery Protocols

In this chapter we discuss some of the techniques that can be used for database recovery from failures. This chapter presents additional concepts that are relevant to recovery protocols, and provides an overview of the various database recovery algorithms.

### **Recovery Concepts**

#### **Recovery Outline and Categorization of Recovery Algorithms**

Recovery from transaction failures usually means that the database is restored to the most recent consistent state just before the time of failure. To do this, the system must keep information about the changes that were applied to data items by the various transactions. This information is typically kept in the **system log**.

A typical strategy for recovery may be summarized informally as follows:

1. If there is extensive damage to a wide portion of the database due to catastrophic failure, such as a disk crash, the recovery method restores a past copy of the database that was backed up to archival storage (typically tape or other large capacity offline storage media) and reconstructs a more current state by reapplying or redoing the operations of committed transactions from the backed up log, up to the time of failure.
2. When the database on disk is not physically damaged, and a non catastrophic failure of types 1 through 4 in Section 21.1.4 has occurred, the recovery strategy is to identify any changes that may cause an inconsistency in the database. For example, a transaction that has updated some database items on disk but has not been committed needs to have its changes reversed by undoing its write operations. It may also be necessary to redo some operations in order to restore a consistent state of the database; for example, if a transaction has committed but some of its write operations have not yet been written to disk. For non catastrophic failure, the recovery protocol does not need a complete archival copy of the database. Rather, the entries kept in the online system log on disk are analyzed to determine the appropriate actions for recovery.

Conceptually, we can distinguish two main techniques for recovery from noncatastrophic transaction failures: **deferred update** and **immediate update**.

- 1) The **deferred update** techniques do not physically update the database on disk until after a transaction reaches its commit point; then the updates are recorded in the database.
- 2) In the **immediate update** techniques, the database may be updated by some operations of a transaction before the transaction reaches its commit point.

In the general case of immediate update, both undo and redo may be required during recovery. This technique, known as the **UNDO/REDO** algorithm.

#### **Caching (Buffering) of Disk Blocks**

The recovery process is often closely intertwined with operating system functions— in particular, the buffering of database disk pages in the DBMS main memory cache. Typically, multiple disk pages that include the data items to be updated are cached into main memory buffers and then updated in memory before being written back to disk. The caching of disk pages is traditionally an operating system function, but because of its importance to the efficiency of recovery procedures, it is handled by the DBMS by calling low-level operating systems routines.

In general, it is convenient to consider recovery in terms of the database disk pages (blocks). Typically a collection of in-memory buffers, called the **DBMS cache**, is kept under the control of the

DBMS for the purpose of holding these buffers. A **directory** for the cache is used to keep track of which database items are in the buffers.

Two main strategies can be employed when flushing a modified buffer back to disk. The first strategy, known as **in-place updating**, writes the buffer to the same original disk location, thus overwriting the old value of any changed data items on disk. Hence, a single copy of each database disk block is maintained. The second strategy, known as **shadowing**, writes an updated buffer at a different disk location, so multiple versions of data items can be maintained, but this approach is not typically used in practice.

In general, the old value of the data item before updating is called the **before image (BFIM)**, and the new value after updating is called the **after image (AFIM)**. If shadowing is used, both the BFIM and the AFIM can be kept on disk; hence, it is not strictly necessary to maintain a log for recovering.

### **Write-Ahead Logging, Steal/No-Steal, and Force/No-Force**

When in-place updating is used, it is necessary to use a log for recovery .

In this case, the recovery mechanism must ensure that the BFIM of the data item is recorded in the appropriate log entry and that the log entry is flushed to disk before the BFIM is overwritten with the AFIM in the database on disk. This process is generally known as **write-ahead logging**, and is necessary to be able to UNDO the operation if this is required during recovery .

A **REDO-type log entry** includes the new value (AFIM) of the item written by the operation since this is needed to redo the effect of the operation from the log (by setting the item value in the database on disk to its AFIM). The **UNDO-type log entries** include the old value (BFIM) of the item since this is needed to undo the effect of the operation from the log (by setting the item value in the database back to its BFIM).

Standard DBMS recovery terminology includes the terms **steal/no-steal** and **force/no-force**, which specify the rules that govern when a page from the database can be written to disk from the cache:

1. If a cache buffer page updated by a transaction cannot be written to disk before the transaction commits, the recovery method is called a **no-steal approach**. The pin-unpin bit will be used to indicate if a page cannot be written back to disk. On the other hand, if the recovery protocol allows writing an updated buffer before the transaction commits, it is called **steal**. Steal is used when the DBMS cache (buffer) manager needs a buffer frame for another transaction and the buffer manager replaces an existing page that had been updated but whose transaction has not committed. The no-steal rule means that UNDO will never be needed during recovery, since a committed transaction will not have any of its updates on disk before it commits.
2. If all pages updated by a transaction are immediately written to disk before the transaction commits, it is called a **force approach**. Otherwise, it is called **no-force**. The force rule means that REDO will never be needed during recovery, since any committed transaction will have all its updates on disk before it is committed.

To permit recovery when in-place updating is used, the appropriate entries required for recovery must be permanently recorded in the log on disk before changes are applied to the database. For example, consider the following **write-ahead logging (WAL) protocol** for a recovery algorithm that requires both UNDO and REDO:

1. The before image of an item cannot be overwritten by its after image in the database on disk until all UNDO-type log records for the updating transaction—up to this point—have been force-written to disk.
2. The commit operation of a transaction cannot be completed until all the REDO-type and UNDO-type log records for that transaction have been force written to disk.

**Checkpoints in the System Log and Fuzzy Checkpointing**

Another type of entry in the log is called a **checkpoint**. A [checkpoint, list of active transactions] record is written into the log periodically at that point when the system writes out to the database on disk all DBMS buffers that have been modified. As a consequence of this, all transactions that have their [commit, T] entries in the log before a [checkpoint] entry do not need to have their WRITE operations redone in case of a system crash, since all their updates will be recorded in the database on disk during checkpointing.

The recovery manager of a DBMS must decide at what intervals to take a checkpoint. The interval may be measured in time—say, every  $m$  minutes—or in the number  $t$  of committed transactions since the last checkpoint, where the values of  $m$  or  $t$  are system parameters. Taking a checkpoint consists of the following actions:

1. Suspend execution of transactions temporarily.
2. Force-write all main memory buffers that have been modified to disk.
3. Write a [checkpoint] record to the log, and force-write the log to disk.
4. Resume executing transactions.

**fuzzy checkpointing.** In this technique, the system can resume transaction processing after a [begin\_checkpoint] record is written to the log without having to wait for step 2 to finish.

**Transaction Rollback and Cascading Rollback**

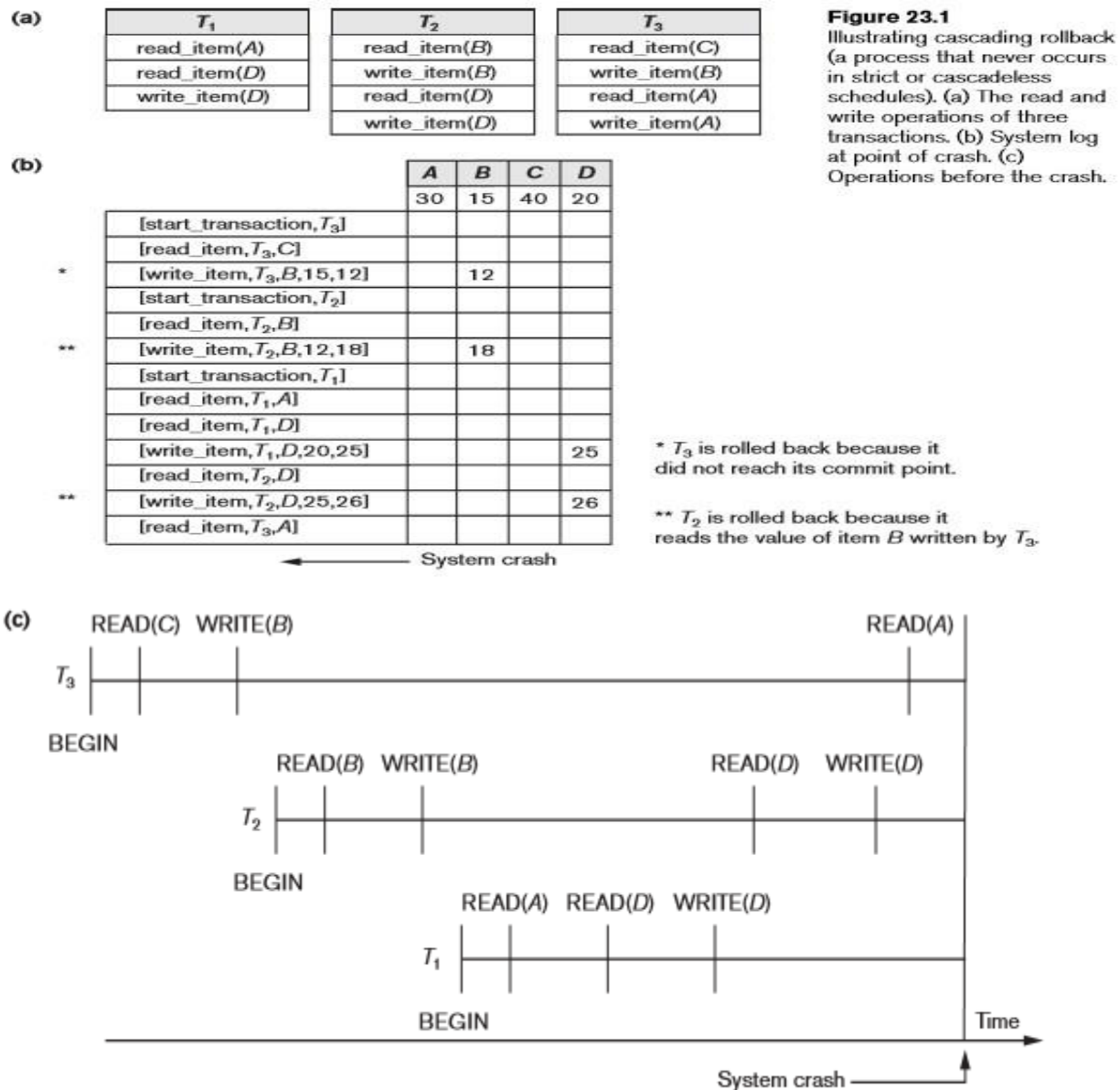
If a transaction fails for whatever reason after updating the database, but before the transaction commits, it may be necessary to **roll back** the transaction. If any data item values have been changed by the transaction and written to the database, they must be restored to their previous values (BFIMs).

If a transaction  $T$  is rolled back, any transaction  $S$  that has, in the interim, read the value of some data item  $X$  written by  $T$  must also be rolled back. Similarly, once  $S$  is rolled back, any transaction  $R$  that has read the value of some data item  $Y$  written by  $S$  must also be rolled back; and so on. This phenomenon is called **cascading rollback**, and can occur when the recovery protocol ensures recoverable schedules but does not ensure strict or cascadeless schedules.

Figure 23.1 shows an example where cascading rollback is required. The read and write operations of three individual transactions are shown in Figure 23.1(a). Figure 23.1(b) shows the system log at the point of a system crash for a particular execution schedule of these transactions. The values of data items  $A, B, C$ , and  $D$ , which are used by the transactions, are shown to the right of the system log entries. We assume that the original item values, shown in the first line, are  $A=30, B=15, C=40$ , and  $D=20$ . At the point of system failure, transaction  $T_3$  has not reached its conclusion and must be rolled back. The WRITE operations of  $T_3$ , marked by a single  $*$  in Figure 23.1(b), are the  $T_3$  operations that are undone during transaction rollback. Figure 23.1(c) graphically shows the operations of the different transactions along the time axis.

We must now check for cascading rollback. From Figure 23.1(c) we see that transaction  $T_2$  reads the value of item  $B$  that was written by transaction  $T_3$ ; this can also be determined by examining the log. Because  $T_3$  is rolled back,  $T_2$  must now be rolled back, too. The WRITE operations of  $T_2$ , marked by  $**$  in the log, are the ones that are undone. Note that only write\_item operations need to be undone during transaction rollback; read\_item operations are recorded in the log only to determine whether cascading rollback of additional transactions is necessary.



**Figure 23.1**

Illustrating cascading rollback (a process that never occurs in strict or cascadeless schedules). (a) The read and write operations of three transactions. (b) System log at point of crash. (c) Operations before the crash.

\*  $T_3$  is rolled back because it did not reach its commit point.

\*\*  $T_2$  is rolled back because it reads the value of item B written by  $T_3$ .

### Transaction Actions That Do Not Affect the Database

In general, a transaction will have actions that do not affect the database, such as generating and printing messages or reports from information retrieved from the database. If a transaction fails before completion, we may not want the user to get these reports, since the transaction has failed to complete.

**NO-UNDO/REDO Recovery Based on Deferred Update**

The idea behind deferred update is to defer or postpone any actual updates to the database on disk until the transaction completes its execution successfully and reaches its commit point.

During transaction execution, the updates are recorded only in the log and in the cache buffers. After the transaction reaches its commit point and the log is force-written to disk, the updates are recorded in the database. If a transaction fails before reaching its commit point, there is no need to **undo** any operations because the transaction has not affected the database on disk in any way. Therefore, only **REDO type log** entries are needed in the log, which include the new **value (AFIM)** of the item written by a write operation. The **UNDO-type log** entries are not needed since no undoing of operations will be required during recovery.

We can state a typical **deferred update protocol** as follows:

1. A transaction cannot change the database on disk until it reaches its commit point.
2. A transaction does not reach its commit point until all its REDO-type log entries are recorded in the log and the log buffer is force-written to disk.

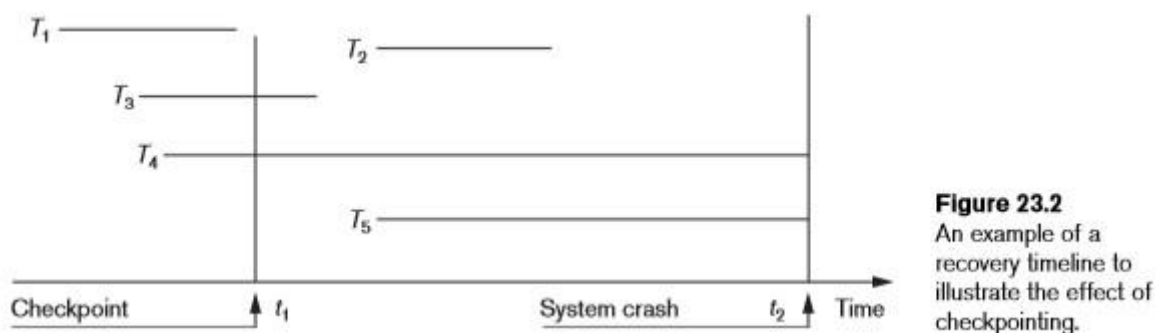
Assuming that [checkpoint] entries are included in the log, a possible recovery algorithm for this case, which we call **RDU\_M (Recovery using Deferred Update in a Multiuser environment)**, is given next.

**Procedure RDU\_M (NO-UNDO/REDO with checkpoints).** Use two lists of transactions maintained by the system: the committed transactions  $T$  since the last checkpoint (commit list), and the active transactions  $T$  (active list). REDO all the WRITE operations of the committed transactions from the log, in the order in which they were written into the log. The transactions that are active and did not commit are effectively canceled and must be resubmitted.

The **REDO procedure** is defined as follows:

**Procedure REDO (WRITE\_OP).** Redoing a write\_item operation WRITE\_OP consists of examining its log entry [write\_item,  $T$ ,  $X$ , new\_value] and setting the value of item  $X$  in the database to new\_value, which is the after image (AFIM).

Figure 23.2 illustrates a timeline for a possible schedule of executing transactions



When the checkpoint was taken at time  $t_1$ , transaction  $T_1$  had committed, whereas transactions  $T_3$  and  $T_4$  had not. Before the system crash at time  $t_2$ ,  $T_3$  and  $T_2$  were committed but not  $T_4$  and  $T_5$ . According to the RDU\_M method, there is no need to redo the write\_item operations of transaction  $T_1$ —or any transactions committed before the last checkpoint time  $t_1$ . The write\_item operations of  $T_2$  and  $T_3$  must be redone, however, because both transactions reached their commit points after the last checkpoint. Recall that the log is force-written before committing a transaction. Transactions  $T_4$  and  $T_5$  are ignored: They are effectively canceled or rolled back because none of their write\_item operations were recorded in the database on disk under the deferred update protocol.

***Recovery Techniques Based on Immediate Update***

In these techniques, when a transaction issues an update command, the database on disk can be updated immediately, without any need to wait for the transaction to reach its commit point. Notice that it is not a requirement that every update be applied immediately to disk; it is just possible that some updates are applied to disk before the transaction commits.

Provisions must be made for undoing the effect of update operations that have been applied to the database by a failed transaction. This is accomplished by rolling back the transaction and undoing the effect of the transaction's write\_item operations. Therefore, the UNDO-type log entries, which include the old value (BFIM) of the item, must be stored in the log. Because UNDO can be needed during recovery, these methods follow a steal strategy for deciding when updated main memory buffers can be written back to disk. Theoretically, we can distinguish two main categories of immediate update algorithms. If the recovery technique ensures that all updates of a transaction are recorded in the database on disk before the transaction commits, there is never a need to REDO any operations of committed transactions. This is called the UNDO/NO-REDO recovery algorithm. In this method, all updates by a transaction must be recorded on disk before the transaction commits, so that REDO is never needed. Hence, this method must utilize the force strategy for deciding when updated main memory buffers are written back to disk.

For a strict schedule, UNDO of an operation requires changing the item back to its old value (BFIM).

**Procedure RIU\_M (UNDO/REDO with checkpoints).**

1. Use two lists of transactions maintained by the system: the committed transactions since the last checkpoint and the active transactions.
2. Undo all the write\_item operations of the active (uncommitted) transactions, using the UNDO procedure. The operations should be undone in the reverse of the order in which they were written into the log.
3. Redo all the write\_item operations of the committed transactions from the log, in the order in which they were written into the log, using the REDO procedure defined earlier.

**The UNDO procedure is defined as follows:**

**Procedure UNDO(WRITE\_OP).**

Undoing a write\_item operation write\_op consists of examining its log entry [write\_item, T, X, old\_value, new\_value] and setting the value of item X in the database to old\_value, which is the before image (BFIM). Undoing a number of write\_item operations from one or more transactions from the log must proceed in the reverse order from the order in which the operations were written in the log.

***Shadow Paging***

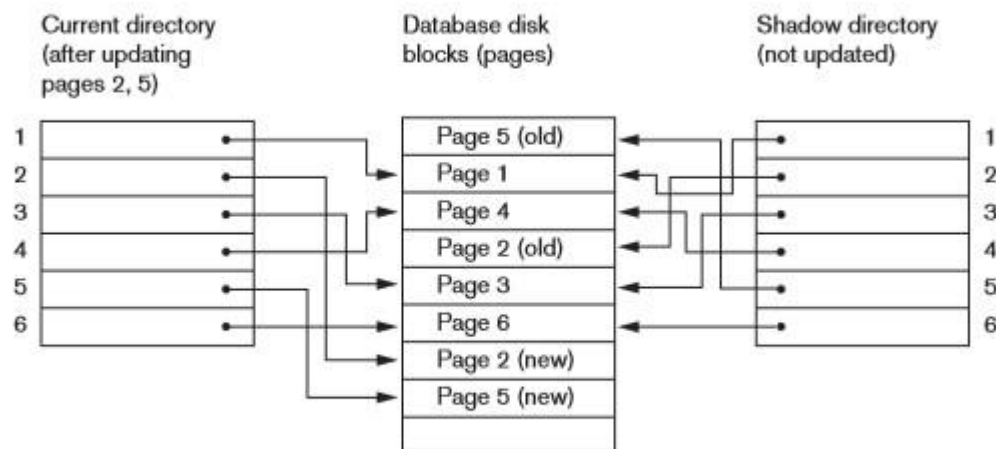
This recovery scheme does not require the use of a log in a single-user environment. In a multiuser environment, a log may be needed for the concurrency control method. **Shadow paging** considers the database to be made up of a number of fixed-size disk pages (or disk blocks)—say, n—for recovery purposes. A **directory** with n entries is constructed, where the i<sup>th</sup> entry points to the i<sup>th</sup> database page on disk. The directory is kept in main memory if it is not too large, and all references—reads or writes—to database pages on disk go through it. When a transaction begins executing, the **current directory**—whose entries point to the most recent or current database pages on disk—is copied

into a **shadow directory**. The shadow directory is then saved on disk while the current directory is used by the transaction.

During transaction execution, the shadow directory is never modified. When a write\_item operation is performed, a new copy of the modified database page is created, but the old copy of that page is not overwritten. Instead, the new page is written elsewhere—on some previously unused disk block. The current directory entry is modified to point to the new disk block, whereas the shadow directory is not modified and continues to point to the old unmodified disk block. Figure 23.4 illustrates the concepts of shadow and current directories. For pages updated by the transaction, two versions are kept. The **old version** is referenced by the **shadow directory** and the **new version** by the **current directory**.

**Figure 23.4**

An example of shadow paging.



<sup>5</sup>The directory is similar to the page table maintained by the operating system for each process.

### 23.7 Database Backup and Recovery from Catastrophic Failures

So far, all the techniques we have discussed apply to noncatastrophic failures. A key assumption has been that the system log is maintained on the disk and is not lost as a result of the failure. Similarly, the shadow directory must be stored on disk to allow recovery when shadow paging is used. The recovery techniques we have discussed use the entries in the system log or the shadow directory to recover from failure by bringing the database back to a consistent state.

The recovery manager of a DBMS must also be equipped to handle more catastrophic failures such as disk crashes. The main technique used to handle such crashes is a database backup, in which the whole database and the log are periodically copied onto a cheap storage medium such as magnetic tapes or other large capacity offline storage devices. In case of a catastrophic system failure, the latest **backup copy** can be reloaded from the tape to the disk, and the system can be restarted.

Data from critical applications such as banking, insurance, stock market, and other databases is periodically backed up in its entirety and moved to physically separate safe locations. Subterranean storage vaults have been used to protect such data from flood, storm, earthquake, or fire damage. Events like the 9/11 terrorist attack in New York (in 2001) and the Katrina hurricane disaster in New Orleans (in 2005) have created a greater awareness of disaster recovery of business-critical databases.

To avoid losing all the effects of transactions that have been executed since the last backup, it is customary to back up the system log at more frequent intervals than full database backup by periodically copying it to magnetic tape. The system log is usually substantially smaller than the database itself and hence can be backed up more frequently. Therefore, users do not lose all transactions they have performed since the last database backup. All committed transactions recorded in the portion of the system log that has been backed up to tape can have their effect on the database redone. A new log is started after each database backup. Hence, to recover from disk failure, the database is first recreated on disk from its latest backup copy on tape. Following that, the effects of all the committed transactions whose operations have been recorded in the backed-up copies of the system log are reconstructed.