

## **Module 5**

### **Chapter 1: Transaction Processing**

- 5.0 Introduction
- 5.1 Objectives
- 5.2 Introduction to Transaction Processing
  - 5.2.1 Single-User versus Multiuser Systems
  - 5.2.2 Transactions, Database Items, Read and Write Operations, and DBMS Buffers
  - 5.2.3 Why Concurrency Control Is Needed
  - 5.2.4 Why Recovery Is Needed
- 5.3 Transaction and System Concepts
  - 5.3.1 Transaction States and Additional Operations
  - 5.3.2 The System Log
  - 5.3.3 Commit Point of a Transaction:
  - 5.3.4 DBMS specific buffer Replacement policies
- 5.4 Desirable Properties of Transactions
- 5.5 Characterizing Schedules Based on Recoverability
- 5.6 Characterizing Schedules Based on Serializability
  - 5.6.1 Testing conflict serializability of a Schedule S
- 5.7 Transaction Support in SQL
- 5.8 Introduction to Concurrency Control
- 5.9 Two-Phase Locking Techniques for Concurrency Control
  - 5.9.1 Types of Locks and System Lock Tables
  - 5.9.2 Guaranteeing Serializability by Two-Phase Locking
- 5.10 Variations of Two-Phase Locking
- 5.11 Dealing with Deadlock and Starvation
- 5.11 Deadlock Detection.
- 5.13 Concurrency Control Based on Timestamp Ordering
  - 5.13.1 Timestamps
  - 5.13.2 The Timestamp Ordering Algorithm
- 5.14 Multiversion Concurrency Control Techniques
  - 5.14.1 Multiversion Technique Based on Timestamp Ordering
  - 5.14.2 Multiversion Two-Phase Locking Using Certify Locks
- 5.15 Validation (Optimistic) Concurrency Control Techniques
- 5.16 Granularity of Data Items and Multiple Granularity Locking
  - 5.16.1 Granularity Level Considerations for Locking
  - 5.16.2 Multiple Granularity Level Locking
- 5.17 Recovery Concepts

- 5.17.1 Recovery Outline and Categorization of Recovery Algorithms
- 5.17.2 Caching (Buffering) of Disk Blocks
- 5.17.3 Write-Ahead Logging, Steal/No-Steal, and Force/No-Force
- 5.17.4 Checkpoints in the System Log and Fuzzy Checkpointing
- 5.17.5 Transaction Rollback and Cascading Rollback
- 5.17.6 Transaction Actions That Do Not Affect the Database
- 5.18 NO-UNDO/REDO Recovery Based on Deferred Update
- 5.19 Recovery Techniques Based on Immediate Update
- 5.20 Shadow Paging
- 5.21 The ARIES Recovery Algorithm
- 5.22 Database Backup and Recovery from Catastrophic Failures
- 5.23 Assignment Questions
- 5.24 Expected Outcome
- 5.25 Further Reading

## 5.0 Introduction

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:

- airline reservations
- banking
- credit card processing,
- online retail purchasing,
- Stock markets, supermarket checkouts, and many other applications

These systems require high availability and fast response time for hundreds of concurrent users. A transaction is typically implemented by a computer program, which includes database commands such as retrievals, insertions, deletions, and updates.

## 5.1 Objectives

- ❖ To study transaction properties
- ❖ To study creation of schedule and maintaining schedule equivalence.
- ❖ To check whether the given schedule is serializable or not.
- ❖ To study protocols used for locking objects
- ❖ Differentiating between 2PL and Strict 2PL

## 5.2 Introduction to Transaction Processing

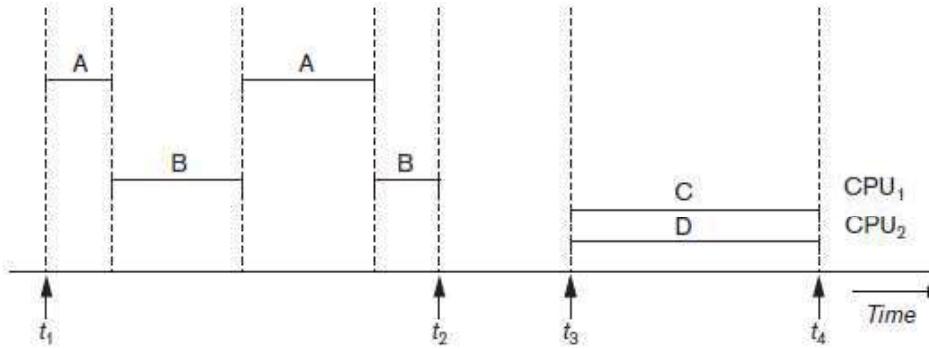
### 5.2.1 Single-User versus Multiuser Systems

- One criterion for classifying a database system is according to the number of users who can use the system **concurrently**

#### Single-User versus Multiuser Systems

- A DBMS is
- **single-user**
  - at most one user at a time can use the system
  - Eg: Personal Computer System
- **multiuser**
  - many users can use the system and hence access the database concurrently
  - Eg: Airline reservation database

- Concurrent access is possible because of **Multiprogramming**. Multiprogramming can be achieved by:
  - interleaved execution
  - Parallel Processing
- **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



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

- Figure 21.1, 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
- Most of the theory concerning concurrency control in databases is developed in terms of **interleaved concurrency**
- In a multiuser DBMS, the stored data items are the primary resources that may be accessed concurrently by interactive users or application programs, which are constantly retrieving information from and modifying the database.

## 5.2.2 Transactions, Database Items, Read and Write Operations, and DBMS

### Buffers

- A Transaction is an executing program that forms a logical unit of database processing
- It includes one or more DB access operations such as insertion, deletion, modification or retrieval operation.
- It can be either embedded within an application program using **begin transaction** and **end transaction** statements Or specified interactively via a high level query language such as SQL
- Transaction which do not update database are known as **read only transactions**.
- Transaction which do update database are known as **read write transactions**.
- 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
- Each data item has a unique name
- **Basic DB access operations that a transaction can include are:**
  - **read\_item(X)**: Reads a DB item named X into a program variable.
  - **write\_item(X)**: Writes the value of a program variable into the DB item named X
- **Executing read\_item(X) include the following steps:**
  1. Find the address of the disk block that contains item X
  2. Copy the block into a buffer in main memory
  3. Copy the item X from the buffer to program variable named X.
- **Executing write\_item(X) include the following steps:**
  1. Find the address of the disk block that contains item X
  2. Copy the disk block into a buffer in main memory
  3. Copy item X from program variable named X into its correct location in buffer.
  4. Store the updated disk block from buffer back to disk (either immediately or later).
- Decision of when to store a modified disk block is handled by **recovery manager** of the DBMS in cooperation with operating system.
- A DB cache includes a number of data buffers.
- When the buffers are all occupied a buffer replacement policy is used to choose one of the buffers to be replaced. EG: LRU

- A transaction includes `read_item` and `write_item` operations to access and update DB.

(a)	$T_1$	(b)	$T_2$	<b>Figure 21.2</b> Two sample transactions. (a) Transaction $T_1$ . (b) Transaction $T_2$ .
	<pre>read_item(X); X := X - N; write_item(X); read_item(Y); Y := Y + N; write_item(Y);</pre>		<pre>read_item(X); X := X + M; write_item(X);</pre>	

- The **read-set** of a transaction is the set of all items that the transaction reads
- The **write-set** is the set of all items that the transaction writes
- For example, the read-set of  $T_1$  in Figure 21.2 is  $\{X, Y\}$  and its write-set is also  $\{X, Y\}$ .

### 5.2.3 Why Concurrency Control Is Needed

- Several problems can occur when concurrent transactions execute in an uncontrolled manner
- Example:
  - We consider an Airline reservation DB
  - Each record is stored for an airline flight which includes Number of reserved seats among other information.
  - Types of problems we may encounter:
    - The Lost Update Problem
    - The Temporary Update (or Dirty Read) Problem
    - The Incorrect Summary Problem
    - The Unrepeatable Read Problem

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

- Transaction T1
  - 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.
- Transaction T2
  - reserves M seats on the first flight (X)

## 1. The Lost Update Problem

- occurs when two transactions that access the same DB items have their operations interleaved in a way that makes the value of some DB item incorrect
- Suppose that transactions T1 and T2 are submitted at approximately the same time, and suppose that their operations are interleaved as shown in Figure below

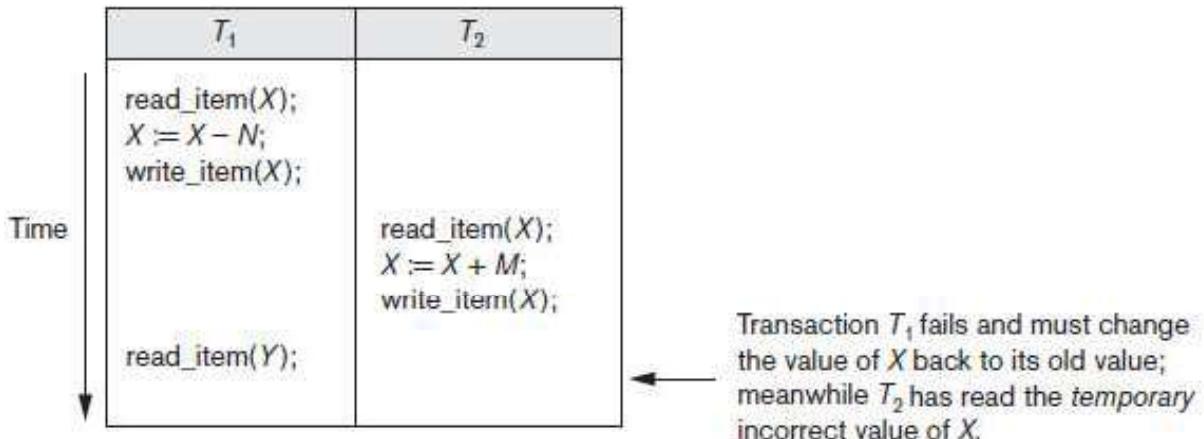
$T_1$	$T_2$
<p>Time ↓</p> <pre> read_item(X); X := X - N; write_item(X); read_item(Y); Y := Y + N; write_item(Y); </pre>	<pre> read_item(X); X := X + M; write_item(X); </pre>

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

- 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:
  - X = 80 at the start (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)
  - M = 4 ( $T_2$  reserves 4 seats on X)
  - The final result should be X = 79.
- The interleaving of operations shown in Figure is X = 84 because the update in  $T_1$  that removed the five seats from X was lost.

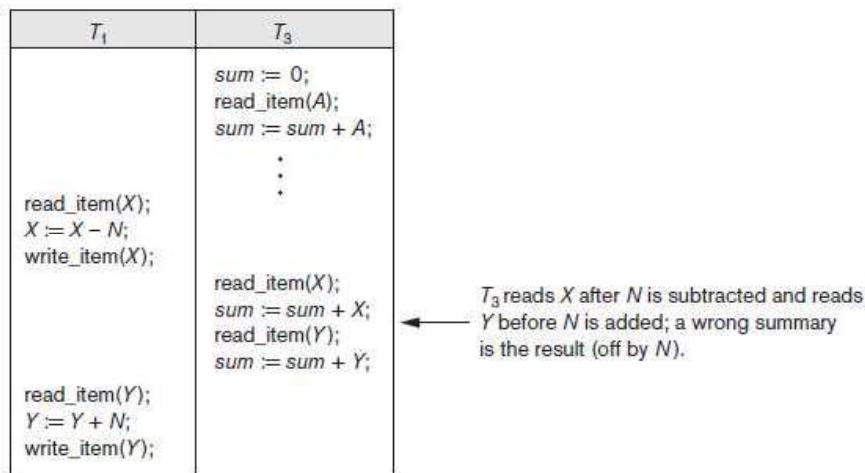
## 2. The Temporary Update (or Dirty Read) Problem

- occurs when one transaction updates a database item and then the transaction fails for some reason
- Meanwhile the updated item is accessed by another transaction before it is changed back to its original value



## 3. The Incorrect Summary Problem

- If one transaction is calculating an aggregate summary function on a number of db 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.



## 4. The Unrepeatable Read Problem

- Transaction T reads the same item twice and gets different values on each read, since the item was modified by another transaction  $T'$  between the two reads.
- 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.

### 5.2.4 Why Recovery Is Needed

- Whenever a transaction is submitted to a DBMS for execution, the system is responsible for making sure that either
  1. All the operations in the transaction are completed successfully and their effect is recorded permanently in the database or
  2. 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
- 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

#### 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
- Also occur because of erroneous parameter values

#### 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

#### 4. Concurrency control enforcement:

- The concurrency control may decide to abort a transaction because it violates serializability or several transactions are in a state of deadlock

#### 5. Disk failure:

- Some disk blocks may lose their data because of a read or write malfunction or because of a disk read/write head crash.

#### 6. Physical problems and catastrophes:

- refers to an endless list of problems that includes power or air-conditioning failure, fire, theft, overwriting disks or tapes by mistake
- Failures of types 1, 2, 3, and 4 are more common than those of types 5 or 6.
- Whenever a failure of type 1 through 4 occurs, the system must keep sufficient information to quickly recover from the failure.
- Disk failure or other catastrophic failures of type 5 or 6 do not happen frequently; if they do occur, recovery is a major task.

### 5.3 Transaction and System Concepts

#### 5.3.1 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 keeps track of start of a transaction, termination, commit or aborts.
  - **BEGIN\_TRANSACTION**: marks the beginning of transaction execution
  - **READ or WRITE**: specify read or write operations on the database items that are executed as part of a transaction
  - **END\_TRANSACTION**: specifies that READ and WRITE transaction operations have ended and marks the end of transaction execution
  - **COMMIT\_TRANSACTION**: 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**: 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**

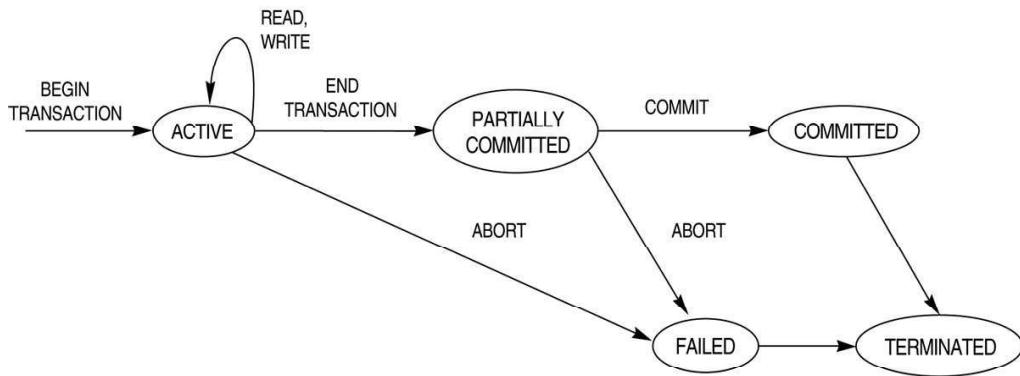


Figure: State transition diagram illustrating the states for transaction execution

- A transaction goes into **active state** immediately after it starts execution and can execute read and write operations.
- When the transaction ends it moves to **partially committed state**.
- At this end additional checks are done to see if the transaction can be committed or not. If these checks are successful the transaction is said to have reached commit point and enters **committed state**. All the changes are recorded permanently in the db.
- 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 operation.
- Terminated state corresponds to the transaction leaving the system. All the information about the transaction is removed from system tables.

### 5.3.2 The System Log

- **Log or Journal** keeps track of all transaction operations that affect the values of database items
- This information may be needed to permit recovery from transaction failures.
- The log is kept on disk, so it is not affected by any type of failure except for disk or catastrophic failure
- one (or more) main memory buffers hold the last part of the log file, so that log entries are first added to the main memory buffer
- When the **log buffer** is filled, or when certain other conditions occur, the log buffer is appended to the end of the log file on disk.

- In addition, the log is periodically backed up to archival storage (tape) to guard against such catastrophic 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.

### 5.3.3 Commit Point of a Transaction:

- **Definition a Commit Point:**
  - 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 has been recorded in the log.
  - Beyond the commit point, the transaction is said to be committed, and its effect is assumed to be permanently recorded in the database.
  - The transaction then writes an entry [commit,T] into the log.
- **Roll Back of transactions:**
  - Needed for transactions that have a [start\_transaction,T] entry into the log but no commit entry [commit,T] into the log.

### 5.3.4 DBMS specific buffer Replacement policies

#### Domain Separation(DS) method

- DBMS cache is divided into separate domains, each handles one type of disk pages and replacements within each domain are handled via basic LRU page replacement.
- LRU is a **static** algorithm and does not adopt to dynamically changing loads because the number of available buffers for each domain is predetermined.
- **Group LRU** adds dynamically load balancing feature since it gives each domain a priority and selects pages from lower priority level domain first for replacement.

**Hot Set Method:**

- This is useful in queries that have to scan a set of pages repeatedly.
- The hot set method determines for each db processing algorithm the set of disk pages that will be accessed repeatedly and it does not replace them until their processing is completed.

**The DBMIN method:**

- uses a model known as QLSM (Query Locality set model), which predetermines the pattern of page references for each algorithm for a particular db operation
- Depending on the type of access method, the file characteristics, and the algorithm used the QLSM will estimate the number of main memory buffers needed for each file involved in the operation.

## 5.4 Desirable Properties of Transactions

- Transactions should possess several properties, often called the **ACID** properties
  - A **Atomicity**: a transaction is an atomic unit of processing and it is either performed entirely or not at all.
  - C **Consistency Preservation**: a transaction should be consistency preserving that is it must take the database from one consistent state to another.
  - I **Isolation/Independence**: A transaction should appear as though it is being executed in isolation from other transactions, even though many transactions are executed concurrently.
  - D **Durability (or Permanency)**: if a transaction changes the database and is committed, the changes must never 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.
- The preservation of **consistency** is generally considered to be the responsibility of the programmers who write the database programs or of the DBMS module that enforces integrity constraints.
- The **isolation** property is enforced by the concurrency control subsystem of the DBMS. If every transaction does not make its updates (write operations) visible to other transactions until it is committed, one form of isolation is enforced that solves the temporary update problem and eliminates cascading rollbacks
- **Durability** is the responsibility of recovery subsystem.

## 5.5 Characterizing Schedules Based on Recoverability

- **schedule** (or **history**): the order of execution of operations from all the various transactions
- **Schedules (Histories) of Transactions:** A schedule S of n transactions  $T_1, T_2, \dots, T_n$  is a sequential ordering of the operations of the n transactions.
  - The transactions are interleaved
- 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)*
- **Conflicting operations:**
  - $r_1(X)$  conflicts with  $w_2(X)$
  - $r_2(X)$  conflicts with  $w_1(X)$
  - $w_1(X)$  conflicts with  $w_2(X)$
  - $r_1(X)$  do not conflicts with  $r_2(X)$

} Read write conflict      Write conflict

### Schedules classified on recoverability:

- **Recoverable schedule:**
  - One where no transaction needs to be rolled back.
  - 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.
  - Example:
    - $S_c: r_1(X); w_1(X); r_2(X); r_1(Y); w_2(X); c_2; a_1;$
    - $S_d: r_1(X); w_1(X); r_2(X); r_1(Y); w_2(X); w_1(Y); c_1; c_2;$
- **Cascadeless schedule:**
  - One where every transaction reads only the items that are written by committed transactions.
- **Schedules requiring cascaded rollback:**
  - A schedule in which uncommitted transactions that read an item from a failed transaction must be rolled back.
- **Strict Schedules:**
  - A schedule in which a transaction can neither read or write an item X until the last transaction that wrote X has committed.

## 5.6 Characterizing Schedules Based on Serializability

- schedules that are always considered to be correct when concurrent transactions are executing are known as **serializable** schedules
- Suppose that two users—for example, two airline reservations agents—submit to the DBMS transactions  $T_1$  and  $T_2$  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).

**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.

(a)	$T_1$	$T_2$	(b)	$T_1$	$T_2$
	$\begin{array}{l} \text{read\_item}(X); \\ X := X - N; \\ \text{write\_item}(X); \\ \text{read\_item}(Y); \\ Y := Y + N; \\ \text{write\_item}(Y); \end{array}$	$\begin{array}{l} \text{read\_item}(X); \\ X := X + M; \\ \text{write\_item}(X); \end{array}$		$\begin{array}{l} \text{read\_item}(X); \\ X := X - N; \\ \text{write\_item}(X); \\ \text{read\_item}(Y); \\ Y := Y + N; \\ \text{write\_item}(Y); \end{array}$	$\begin{array}{l} \text{read\_item}(X); \\ X := X + M; \\ \text{write\_item}(X); \end{array}$
<b>Schedule A</b>			<b>Schedule B</b>		
(c)	$T_1$	$T_2$	(c)	$T_1$	$T_2$
	$\begin{array}{l} \text{read\_item}(X); \\ X := X - N; \\ \text{write\_item}(X); \\ \text{read\_item}(Y); \\ Y := Y + N; \\ \text{write\_item}(Y); \end{array}$	$\begin{array}{l} \text{read\_item}(X); \\ X := X + M; \\ \text{write\_item}(X); \end{array}$		$\begin{array}{l} \text{read\_item}(X); \\ X := X - N; \\ \text{write\_item}(X); \end{array}$	$\begin{array}{l} \text{read\_item}(X); \\ X := X + M; \\ \text{write\_item}(X); \end{array}$
<b>Schedule C</b>			<b>Schedule D</b>		

- **Serial schedule:**
  - 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 schedule.
- **Serializable schedule:**
  - A schedule S is serializable if it is equivalent to some serial schedule of the same n transactions.
- **Result equivalent:**
  - Two schedules are called result equivalent if they produce the same final state of the database.
- **Conflict equivalent:**
  - Two schedules are said to be conflict equivalent if the order of any two conflicting operations is the same in both schedules.
- **Conflict serializable:**
  - A schedule S is said to be conflict serializable if it is conflict equivalent to some serial schedule S'.
- Being serializable is not the same as being serial
- Being serializable implies that the schedule is a correct schedule.
  - It will leave the database in a consistent state.
  - The interleaving is appropriate and will result in a state as if the transactions were serially executed, yet will achieve efficiency due to concurrent execution.

### 5.6.1 Testing conflict serializability of a Schedule S

For each transaction  $T_i$  participating in schedule S, create a node labeled  $T_i$  in the precedence graph.

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.

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.

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.

The schedule S is serializable if and only if the precedence graph has no cycles.

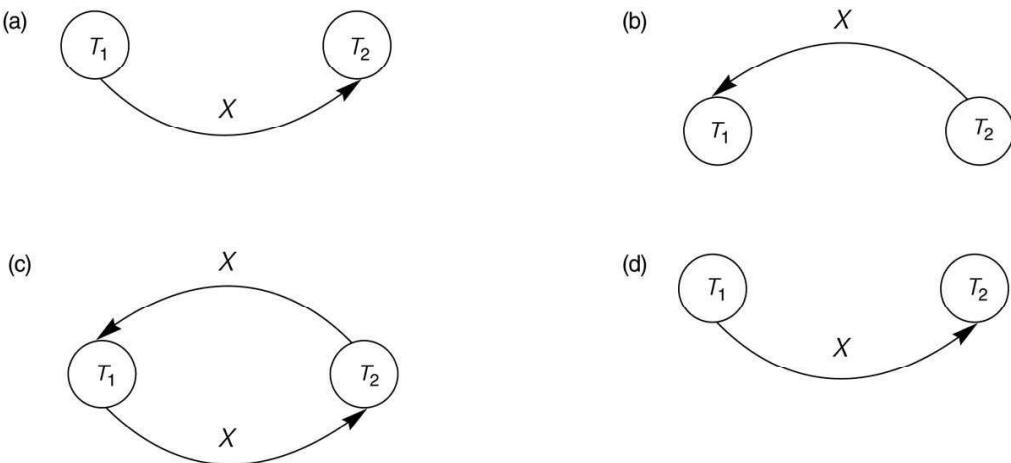


Fig: Constructing the precedence graphs for schedules A and D from fig 21.5 to test for conflict serializability.

- (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).

- Another example of serializability testing. (a) The READ and WRITE operations of three transactions  $T_1$ ,  $T_2$ , and  $T_3$ .

(a)	transaction $T_1$	transaction $T_2$	transaction $T_3$
	<pre>read_item (X); write_item (X); read_item (Y); write_item (Y);</pre>	<pre>read_item (Z); read_item (Y); write_item (Y); read_item (X); write_item (X);</pre>	<pre>read_item (Y); read_item (Z); write_item (Y); write_item (Z);</pre>

(b)

	transaction $T_1$	transaction $T_2$	transaction $T_3$
<i>Time</i>		read_item ( $Z$ ); read_item ( $Y$ ); write_item ( $Y$ );	
	read_item ( $X$ ); write_item ( $X$ );		read_item ( $Y$ ); read_item ( $Z$ );
	read_item ( $Y$ ); write_item ( $Y$ );	read_item ( $X$ );  write_item ( $X$ );	write_item ( $Y$ ); write_item ( $Z$ );

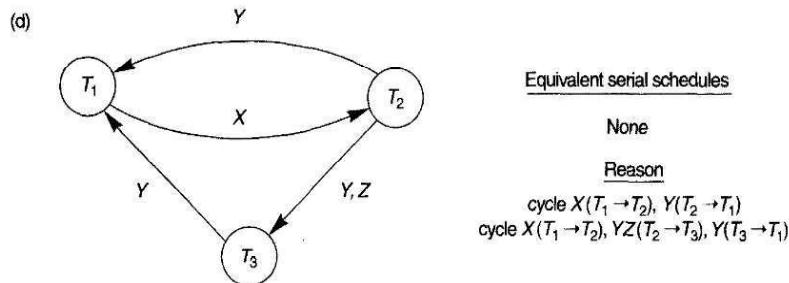
Schedule E

(c)

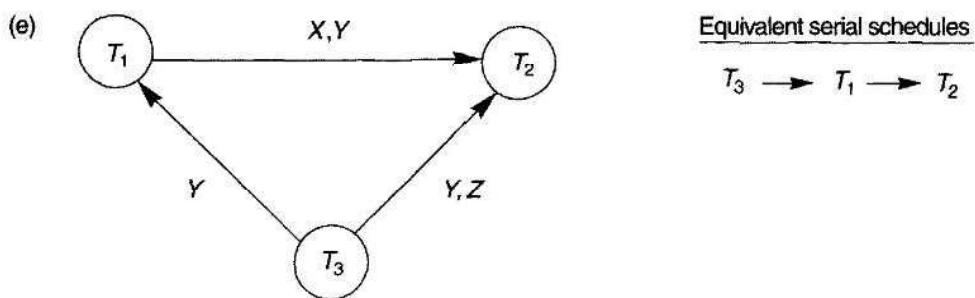
	transaction $T_1$	transaction $T_2$	transaction $T_3$
<i>Time</i>			read_item ( $Y$ ); read_item ( $Z$ );
	read_item ( $X$ ); write_item ( $X$ );	read_item ( $Z$ );  read_item ( $Y$ ); write_item ( $Y$ ); read_item ( $X$ ); write_item ( $X$ );	write_item ( $Y$ ); write_item ( $Z$ );
	read_item ( $Y$ ); write_item ( $Y$ );		

Schedule F

- Precedence graph for schedule E



- Precedence graph for schedule F



## 5.7 Transaction Support in SQL

- The basic definition of an SQL transaction 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
- Every transaction must have an explicit end statement, which is either a `COMMIT` or a `ROLLBACK`
- Every transaction has certain characteristics attributed to it and are specified by a `SET TRANSACTION` statement in SQL

- The characteristics are :
  - **The access mode**
    - can be specified as READ ONLY or READ WRITE
    - The default is READ WRITE
    - 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**
    - 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**
    - specified using the statement ISOLATION LEVEL <isolation>, where the value for <isolation> can be READ UNCOMMITTED, READ COMMITTED, REPEATABLE READ, or SERIALIZABLE
      - The default isolation level is SERIALIZABLE
      - The use of the term SERIALIZABLE here is based on not allowing violations that cause dirty read, unrepeatable read, and phantoms
      - 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  $T_1$  may read the update of a transaction  $T_2$ , which has not yet committed. If  $T_2$  fails and is aborted, then  $T_1$  would have read a value that does not exist and is incorrect.
        2. **Nonrepeatable read.** A transaction  $T_1$  may read a given value from a table. If another transaction  $T_2$  later updates that value and  $T_1$  reads that value again,  $T_1$  will see a different value.
        3. **Phantoms.** A transaction  $T_1$  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  $T_2$  inserts a new row that also satisfies the WHERE-clause condition used in  $T_1$ , into the table used by  $T_1$ . If  $T_1$  is repeated, then  $T_1$  will see a phantom, a row that previously did not exist.

**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

```

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 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.

## Chapter 2: Concurrency Control in Databases

### 5.8 Introduction to Concurrency Control

- Purpose of Concurrency Control
  - To enforce Isolation (through mutual exclusion) among conflicting transactions.
  - To preserve database consistency through consistency preserving execution of transactions.
  - To resolve read-write and write-write conflicts.
- Example:
  - In concurrent execution environment if T1 conflicts with T2 over a data item A, then the existing concurrency control decides if T1 or T2 should get the A and if the other transaction is rolled-back or waits.

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

- The concept of locking data items is one of the main techniques used for controlling the concurrent execution of transactions.
- A lock is a variable associated with a data item in the database. Generally there is a lock for each data item in the database.
- A lock describes the status of the data item with respect to possible operations that can be applied to that item.
- It is used for synchronizing the access by concurrent transactions to the database items.
- A transaction locks an object before using it
- When an object is locked by another transaction, the requesting transaction must wait

#### 5.9.1 Types of Locks and System Lock Tables

##### 1. Binary Locks

- A **binary lock** can have two **states or values**: locked and unlocked (or 1 and 0).
- 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 **lock( $X$ )**.
- Two operations, **lock\_item** and **unlock\_item**, are used with binary locking.
- A transaction requests access to an item  $X$  by first issuing a **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 **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.

**lock\_item( $X$ ):**

```

B: if  $\text{LOCK}(X) = 0$  (* item is unlocked *)
    then  $\text{LOCK}(X) \leftarrow 1$  (* lock the item *)
else
begin
    wait (until  $\text{LOCK}(X) = 0$ 
        and the lock manager wakes up the transaction);
    go to B
end;

```

**unlock\_item( $X$ ):**

```

 $\text{LOCK}(X) \leftarrow 0$ ; (* unlock the item *)
if any transactions are waiting
then wakeup one of the waiting transactions;

```

Fig: 2.1.1 Lock and unlock operations for binary locks.

- The lock\_item and unlock\_item operations must be implemented as indivisible units that is, no interleaving should be allowed once a lock or unlock operation is started until the operation terminates or the transaction waits
- The wait command within the lock\_item( $X$ ) operation is usually implemented by putting the transaction in a waiting queue for item  $X$  until  $X$  is unlocked and the transaction can be granted access to it
- Other transactions that also want to access  $X$  are placed in the same queue. Hence, the wait command is considered to be outside the lock\_item operation.
- It is quite simple to implement a binary lock; all that is needed is a binary-valued variable, LOCK, associated with each data item  $X$  in the database
- In its simplest form, each lock can be a record with three fields: <Data\_item\_name, LOCK, Locking\_transaction> plus a queue for transactions that are waiting to access the item
- 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 lock\_item( $X$ ) before any read\_item( $X$ ) or write\_item( $X$ ) operations are performed in  $T$ .
  2. A transaction  $T$  must issue the operation unlock\_item( $X$ ) after all read\_item( $X$ ) and write\_item( $X$ ) operations are completed in  $T$ .
  3. A transaction  $T$  will not issue a lock\_item( $X$ ) operation if it already holds the lock on item  $X$ .
  4. A transaction  $T$  will not issue an unlock\_item( $X$ ) operation unless it already holds the lock on item  $X$ .

## 2. Shared/Exclusive (or Read/Write) Locks

- binary locking scheme is too restrictive for database items because at most, one transaction can hold a lock on a given item
- should allow several transactions to access the same item  $X$  if they all access  $X$  for reading purposes only
- 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: **read\_lock( $X$ )**, **write\_lock( $X$ )**, and **unlock( $X$ )**.

- 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
- Method to implement read/write lock is to
  - keep track of the number of transactions that hold a shared (read) lock on an item in the lock table
  - Each record in the lock table will have four fields:  
 $\langle \text{Data\_item\_name}, \text{LOCK}, \text{No\_of\_reads}, \text{Locking\_transaction(s)} \rangle$ .
- If  $\text{LOCK}(X)$ =write-locked, the value of locking\_transaction(s) is a single transaction that holds the exclusive (write) lock on  $X$
- If  $\text{LOCK}(X)$ =read-locked, the value of locking transaction(s) is a list of one or more transactions that hold the shared (read) lock on  $X$ .

**read\_lock( $X$ ):**

```
B: if  $\text{LOCK}(X)$  = "unlocked"
    then begin  $\text{LOCK}(X) \leftarrow$  "read-locked";
             $\text{no\_of\_reads}(X) \leftarrow 1$ 
            end
    else if  $\text{LOCK}(X)$  = "read-locked"
        then  $\text{no\_of\_reads}(X) \leftarrow \text{no\_of\_reads}(X) + 1$ 
    else begin
        wait (until  $\text{LOCK}(X)$  = "unlocked"
              and the lock manager wakes up the transaction);
        go to B
        end;
```

**write\_lock( $X$ ):**

```
B: if  $\text{LOCK}(X)$  = "unlocked"
    then  $\text{LOCK}(X) \leftarrow$  "write-locked"
    else begin
        wait (until  $\text{LOCK}(X)$  = "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;

```

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

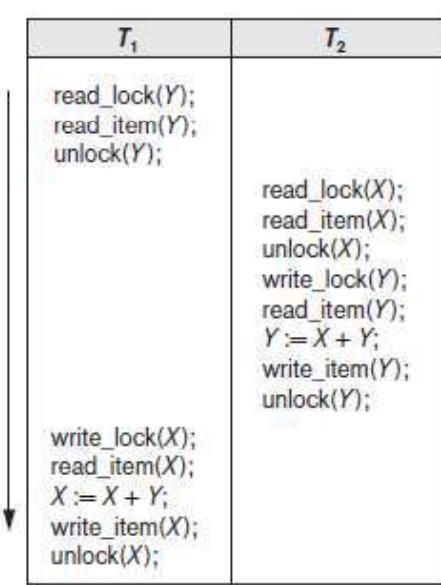
1. A transaction T must issue the operation `read_lock(X)` or `write_lock(X)` before any `read_item(X)` operation is performed in T.
2. A transaction T must issue the operation `write_lock(X)` before any `write_item(X)` operation is performed in T.
3. A transaction T must issue the operation `unlock(X)` after all `read_item(X)` and `write_item(X)` operations are completed in T.
4. A transaction T will not issue a `read_lock(X)` operation if it already holds a read (shared) lock or a write (exclusive) lock on item X.

### Conversion of Locks

- 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 `read_lock(X)` and then later to **upgrade** the lock by issuing a `write_lock(X)` operation
  - If T is the only transaction holding a read lock on X at the time it issues the `write_lock(X)` operation, the lock can be upgraded; otherwise, the transaction must wait

### 5.9.2 Guaranteeing Serializability by 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:
  - **Expanding or growing (first) phase**, during which new locks on items can be acquired but none can be released
  - **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.
- Transactions  $T_1$  and  $T_2$  in Figure 22.3(a) do not follow the two-phase locking protocol because the  $\text{write\_lock}(X)$  operation follows the  $\text{unlock}(Y)$  operation in  $T_1$ , and similarly the  $\text{write\_lock}(Y)$  operation follows the  $\text{unlock}(X)$  operation in  $T_2$ .

(a)	<table border="1" style="width: 100%; border-collapse: collapse;"> <tr> <th style="width: 50%; text-align: center; padding: 2px;"><math>T_1</math></th> <th style="width: 50%; text-align: center; padding: 2px;"><math>T_2</math></th> </tr> <tr> <td style="padding: 2px;"> <math>\text{read\_lock}(Y);</math>  <math>\text{read\_item}(Y);</math>  <math>\text{unlock}(Y);</math>  <math>\text{write\_lock}(X);</math>  <math>\text{read\_item}(X);</math>  <math>X := X + Y;</math>  <math>\text{write\_item}(X);</math>  <math>\text{unlock}(X);</math> </td> <td style="padding: 2px;"> <math>\text{read\_lock}(X);</math>  <math>\text{read\_item}(X);</math>  <math>\text{unlock}(X);</math>  <math>\text{write\_lock}(Y);</math>  <math>\text{read\_item}(Y);</math>  <math>Y := X + Y;</math>  <math>\text{write\_item}(Y);</math>  <math>\text{unlock}(Y);</math> </td> </tr> </table>	$T_1$	$T_2$	$\text{read\_lock}(Y);$ $\text{read\_item}(Y);$ $\text{unlock}(Y);$ $\text{write\_lock}(X);$ $\text{read\_item}(X);$ $X := X + Y;$ $\text{write\_item}(X);$ $\text{unlock}(X);$	$\text{read\_lock}(X);$ $\text{read\_item}(X);$ $\text{unlock}(X);$ $\text{write\_lock}(Y);$ $\text{read\_item}(Y);$ $Y := X + Y;$ $\text{write\_item}(Y);$ $\text{unlock}(Y);$	<b>(b)</b> Initial values: $X=20, Y=30$  Result serial schedule $T_1$ , followed by $T_2$ : $X=50, Y=80$  Result of serial schedule $T_2$ , followed by $T_1$ : $X=70, Y=50$
$T_1$	$T_2$					
$\text{read\_lock}(Y);$ $\text{read\_item}(Y);$ $\text{unlock}(Y);$ $\text{write\_lock}(X);$ $\text{read\_item}(X);$ $X := X + Y;$ $\text{write\_item}(X);$ $\text{unlock}(X);$	$\text{read\_lock}(X);$ $\text{read\_item}(X);$ $\text{unlock}(X);$ $\text{write\_lock}(Y);$ $\text{read\_item}(Y);$ $Y := X + Y;$ $\text{write\_item}(Y);$ $\text{unlock}(Y);$					
(c)		Figure 21.3 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				

- If we enforce two-phase locking, the transactions can be rewritten as  $T_1'$  and  $T_2'$  as shown in Figure 22.4.
- Now, the schedule shown in Figure 22.3(c) is not permitted for  $T_1'$  and  $T_2'$  (with their modified order of locking and unlocking operations) under the rules of locking because  $T_1'$  will issue its `write_lock(X)` before it unlocks item Y; consequently, when  $T_2'$  issues its `read_lock(X)`, it is forced to wait until  $T_1'$  releases the lock by issuing an `unlock(X)` in the schedule.

**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.

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

- If every transaction in a schedule follows the two-phase locking protocol, schedule guaranteed to be serializable
- Two-phase locking may limit the amount of concurrency that can occur in a schedule
- Some serializable schedules will be prohibited by two-phase locking protocol

## 5.10 Variations of Two-Phase Locking

- **Basic 2PL**
  - Technique described previously
- **Conservative (static) 2PL**
  - Requires a transaction to lock all the items it accesses before the transaction begins execution by predeclaring read-set and write-set
  - Its Deadlock-free protocol

- **Strict 2PL**

- guarantees strict schedules
- Transaction does not release exclusive locks until after it commits or aborts
- 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

- **Rigorous 2PL**

- guarantees strict schedules
- Transaction does not release any locks until after it commits or aborts
- easier to implement than strict 2PL

## 5.11 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  $T_1'$  and  $T_2'$  are deadlocked in a partial schedule;  $T_1'$  is in the waiting queue for  $X$ , which is locked by  $T_2'$ , while  $T_2'$  is in the waiting queue for  $Y$ , which is locked by  $T_1'$ . Meanwhile, neither  $T_1'$  nor  $T_2'$  nor any other transaction can access items  $X$  and  $Y$ .

(a)

$T_1'$	$T_2'$
<code>read_lock(<math>Y</math>);</code> <code>read_item(<math>Y</math>);</code>  <code>write_lock(<math>X</math>);</code>	<code>read_lock(<math>X</math>);</code> <code>read_item(<math>X</math>);</code>  <code>write_lock(<math>Y</math>);</code>

Time ↓

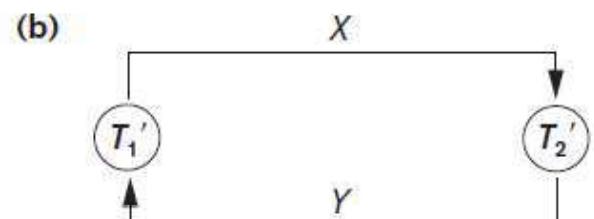


Figure 22.5 Illustrating the deadlock problem (a) A partial schedule of  $T_1'$  and  $T_2'$  that is in a state of deadlock (b) A wait-for graph for the partial schedule in (a)

## 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*. 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.
- 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
- Both approaches impractical
- 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  $TS(T_1) < TS(T_2)$ .
- The *older* transaction (which starts first) has the *smaller* timestamp value.
- Protocols based on a timestamp
  - Wait-die
  - Wound-wait
- 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.
- In wait-die, an older transaction is allowed to *wait for a younger transaction*, whereas a younger transaction requesting an item held by an older transaction is aborted and restarted.
- The wound-wait approach does the opposite: A younger transaction is allowed to *wait for an older one*, whereas an older transaction requesting an item held by a younger transaction *preempts* the younger transaction by aborting it.

- Both schemes end up aborting the *younger* of the two transactions (the transaction that started later) that *may be involved* in a deadlock, assuming that this will waste less processing.
- It can be shown that these two techniques are *deadlock-free*, since in wait-die, transactions only wait for younger transactions so no cycle is created.
- Similarly, in wound-wait, transactions only wait for older transactions so no cycle is created.
- Another group of protocols that prevent deadlock do not require timestamps. These include the
  - no waiting (NW) and
  - cautious waiting (CW) algorithms
- **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.
  - no transaction ever waits, so no deadlock will occur
  - this scheme can cause transactions to abort and restart needlessly
- **cautious waiting**
  - 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:
    - If  $T_j$  is not blocked (not waiting for some other locked item), then  $T_i$  is blocked and allowed to wait; otherwise abort  $T_i$ .
  - It can be shown that cautious waiting is deadlock-free, because no transaction will ever wait for another blocked transaction.

## 5.12 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.
- This solution is attractive if we know there will be little interference among the transactions—that is, if different transactions will rarely access the same items at the same time.

- This can happen if the transactions are short and each transaction locks only a few items, or if the transaction load is light.
- On the other hand, if transactions are long and each transaction uses many items, or if the transaction load is quite heavy, it may be advantageous to use a deadlock prevention scheme.
- 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.
- One problem with this approach is the matter of determining *when* the system should check for a deadlock.
- One possibility is to check for a cycle every time an edge is added to the wait-for graph, but this may cause excessive overhead.
- Criteria such as the number of currently executing transactions or the period of time several transactions have been waiting to lock items may be used instead to check for a cycle. Figure 22.5(b) 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**.
  - The algorithm for victim selection should generally avoid selecting transactions that have been running for a long time and that have performed many updates, and it should try instead to select transactions that have not made many changes (younger transactions).
- **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.

## 5.13 Concurrency Control Based on Timestamp Ordering

guarantees serializability using transaction timestamps to order transaction execution for an equivalent serial schedule

### 5.13.1 Timestamps

- **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.

### 5.13.2 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)**.
- 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,  $\text{read\_TS}(X) = \text{TS}(T)$ , 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,  $\text{write\_TS}(X) = \text{TS}(T)$ , where  $T$  is the *youngest transaction* that has written  $X$  successfully.

#### Basic Timestamp Ordering (TO).

- Whenever some transaction  $T$  tries to issue a  $\text{read\_item}(X)$  or a  $\text{write\_item}(X)$  operation, the **basic TO** algorithm compares the timestamp of  $T$  with  $\text{read\_TS}(X)$  and  $\text{write\_TS}(X)$  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  $T_1$  that may have used a value written by  $T$  must also be rolled back.

- Similarly, any transaction  $T_2$  that may have used a value written by  $T_1$  must also be rolled back, and so on. This effect is known as **cascading rollback** and is one of the problems associated with basic TO, since the schedules produced are not guaranteed to be recoverable.
- An *additional protocol* must be enforced to ensure that the schedules are recoverable, cascadeless, or strict.
- **The basic TO algorithm :**
  - 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  $\text{read\_TS}(X)$ .
  - Whenever the basic TO algorithm detects two *conflicting operations* that occur in the incorrect order, it rejects the later of the two operations by aborting the transaction that issued it. The schedules produced by basic TO are hence guaranteed to be *conflict serializable*

### 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  $\text{read\_item}(X)$  or  $\text{write\_item}(X)$  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  $\text{write\_item}(X)$  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  $\text{write\_item}(X)$  operation of  $T$  because it is already outdated and obsolete. Notice that any conflict arising from this situation would be detected by case (1).

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

## 5.14 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
- An obvious drawback of multiversion techniques is that more storage is needed to maintain multiple versions of the database items

### 5.14.1 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  $\text{write\_TS}(X_{k+1})$  and the  $\text{read\_TS}(X_{k+1})$  set to  $\text{TS}(T)$
- Correspondingly, when a transaction  $T$  is allowed to read the value of version  $X_i$ , the value of  $\text{read\_TS}(X_i)$  is set to the larger of the current  $\text{read\_TS}(X_i)$  and  $\text{TS}(T)$ .
- 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  $\text{write\_TS}(X_i)$  of all versions of  $X$  that is also *less than or equal to*  $\text{TS}(T)$ , and  $\text{read\_TS}(X_i) > \text{TS}(T)$ , then abort and roll back transaction  $T$ ; otherwise, create a new version  $X_j$  of  $X$  with  $\text{read\_TS}(X_j) = \text{write\_TS}(X_j) = \text{TS}(T)$ .
  2. If transaction  $T$  issues a `read_item(X)` operation, find the version  $i$  of  $X$  that has the highest  $\text{write\_TS}(X_i)$  of all versions of  $X$  that is also *less than or equal to*  $\text{TS}(T)$ ; then return the value of  $X_i$  to transaction  $T$ , and set the value of  $\text{read\_TS}(X_i)$  to the larger of  $\text{TS}(T)$  and the current  $\text{read\_TS}(X_i)$ .

### 5.14.2 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`
- Hence, the state of  $\text{LOCK}(X)$  for an item  $X$  can be one of `read-locked`, `writelocked`, `certify-locked`, or `unlocked`
- 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 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

(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.

- 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
- The idea behind multiversion 2PL is to allow other transactions  $T$  to read an item  $X$  while a single transaction  $T$  holds a write lock on  $X$
- This is accomplished by allowing two versions for each item  $X$ ; one version must always have been written by some committed transaction
- The second version  $X'$  is created when a transaction  $T$  acquires a write lock on the item

## 5.15 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
- In this scheme, updates in the transaction are not applied directly to the database items until the transaction reaches its end

- During transaction execution, all updates are applied to *local copies* of the data items that are kept for the transaction
- At the end of transaction execution, a **validation phase** checks whether any of the transaction's updates violate serializability.
- 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 idea behind optimistic concurrency control is to do all the checks at once; hence, transaction execution proceeds with a minimum of overhead until the validation phase is reached
- The techniques are called *optimistic* because they assume that little interference will occur and hence that there is no need to do checking during transaction execution.
- 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.

## 5.16 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

#### 5.16.1 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
- The larger the data item size is, the lower the degree of concurrency permitted.
- For example, if the data item size is a disk block, a transaction  $T$  that needs to lock a record  $B$  must lock the whole disk block  $X$  that contains  $B$  because a lock is associated with the whole data item (block). Now, if another transaction  $S$  wants to lock a different record  $C$  that happens to reside in the same block  $X$  in a conflicting lock mode, it is forced to wait. If the data item size was a single record, transaction  $S$  would be able to proceed, because it would be locking a different data item (record).
- 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. More lock and unlock operations will be performed, causing a higher overhead
- The best item size *depends on the types of transactions involved*.
- If a typical transaction accesses a small number of records, it is advantageous to have the data item granularity be one record
- On the other hand, if a transaction typically accesses many records in the same file, it may be better to have block or file granularity so that the transaction will consider all those records as one (or a few) data items

#### 5.16.2 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

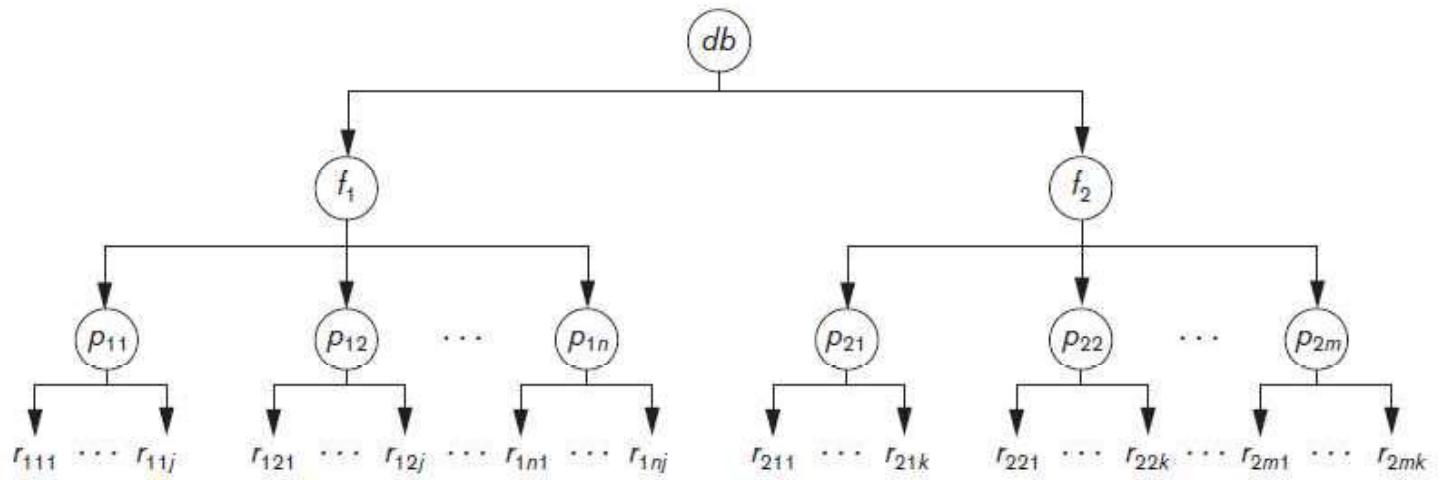


Figure 22.7 A granularity hierarchy for illustrating multiple granularity level locking

- 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.

	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.

- 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  $T$ .
- The multiple granularity level protocol is especially suited when processing a mix of transactions that include
  - (1) short transactions that access only a few items (records or fields) and
  - (2) long transactions that access entire files.

## Chapter 3: Introduction to Database Recovery Protocols

### 5.17 Recovery Concepts

#### 5.17.1 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**.
- Conceptually, we can distinguish two main techniques for recovery from noncatastrophic transaction failures: **deferred update** and **immediate update**.
  - 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
    - Before reaching commit, all transaction updates are recorded in the local transaction workspace or in the main memory buffers that the DBMS maintains
    - Before commit, the updates are recorded persistently in the log, and then after commit, the updates are written to the database on disk
    - If a transaction fails before reaching its commit point, it will not have changed the database in any way, so UNDO is not needed
    - It may be necessary to REDO the effect of the operations of a committed transaction from the log, because their effect may not yet have been recorded in the database on disk
    - Hence, deferred update is also known as the **NO-UNDO/REDO algorithm**
  - The **immediate update** techniques
    - the database *may be updated* by some operations of a transaction *before* the transaction reaches its commit point.
    - However, these operations must also be recorded in the log *on disk* by force-writing *before* they are applied to the database on disk, making recovery still possible

- If a transaction fails after recording some changes in the database on disk but before reaching its commit point, the effect of its operations on the database must be undone; that is, the transaction must be rolled back
- In the general case of immediate update, both *undo* and *redo* may be required during recovery.
- This technique, known as the **UNDO/REDO algorithm**, requires both operations during recovery, and is used most often in practice.

### 5.17.2 Caching (Buffering) of Disk Blocks

- 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
- This can be a table of <Disk\_page\_address, Buffer\_location, ... > entries.
- When the DBMS requests action on some item, first it checks the cache directory to determine whether the disk page containing the item is in the DBMS cache.
- If it is not, the item must be located on disk, and the appropriate disk pages are copied into the cache. It may be necessary to **replace** (or **flush**) some of the cache buffers to make space available for the new item.
- The entries in the DBMS cache directory hold additional information relevant to buffer management.
- Associated with each buffer in the cache is a **dirty bit**, which can be included in the directory entry, to indicate whether or not the buffer has been modified.
- When a page is first read from the database disk into a cache buffer, a new entry is inserted in the cache directory with the new disk page address, and the dirty bit is set to 0 (zero).
- As soon as the buffer is modified, the dirty bit for the corresponding directory entry is set to 1 (one)
- Additional information, such as the transaction id(s) of the transaction(s) that modified the buffer can also be kept in the directory
- When the buffer contents are replaced (flushed) from the cache, the contents must first be written back to the corresponding disk page *only if its dirty bit is 1*

- Another bit, called the **pin-unpin** bit, is also needed—a page in the cache is **pinned** (bit value 1 (one)) if it cannot be written back to disk as yet.
- 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.

### 5.17.3 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)
- In an UNDO/REDO algorithm, both types of log entries are combined. Additionally, when cascading rollback is possible, `read_item` entries in the log are considered to be UNDO-type entries
- 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.
- The deferred update (NO-UNDO) recovery scheme follows a *no-steal* approach.
  - However, typical database systems employ a *steal/no-force* strategy.
  - The *advantage of steal* is that it avoids the need for a very large buffer space to store all updated pages in memory.
  - The *advantage of no-force* is that an updated page of a committed transaction may still be in the buffer when another transaction needs to update it, thus eliminating the I/O cost to write that page multiple times to disk, and possibly to have to read it again from disk.
  - 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.

#### 5.17.4 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
- As part of checkpointing, the list of transaction ids for active transactions at the time of the checkpoint is included in the checkpoint record, so that these transactions can be easily identified during recovery.
- 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.
- The time needed to force-write all modified memory buffers may delay transaction processing because of step 1
- To reduce this delay, it is common to use a technique called **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.
- When step 2 is completed, an [end\_checkpoint, ...] record is written in the log with the relevant information collected during checkpointing.

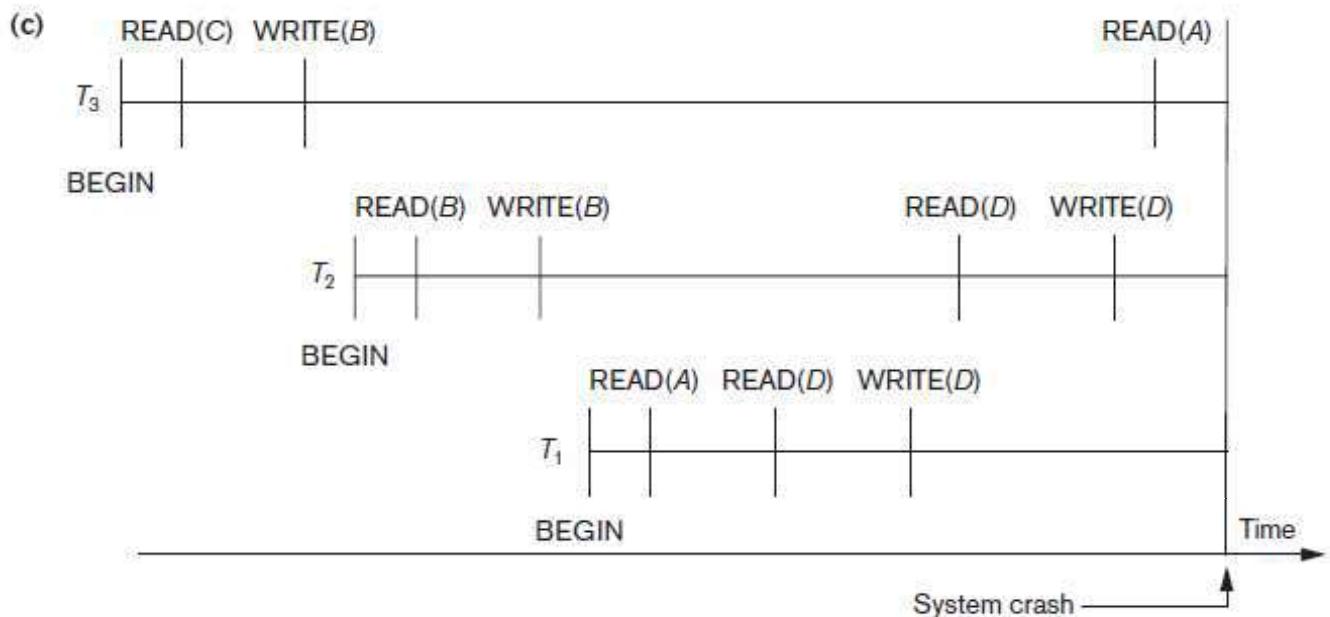
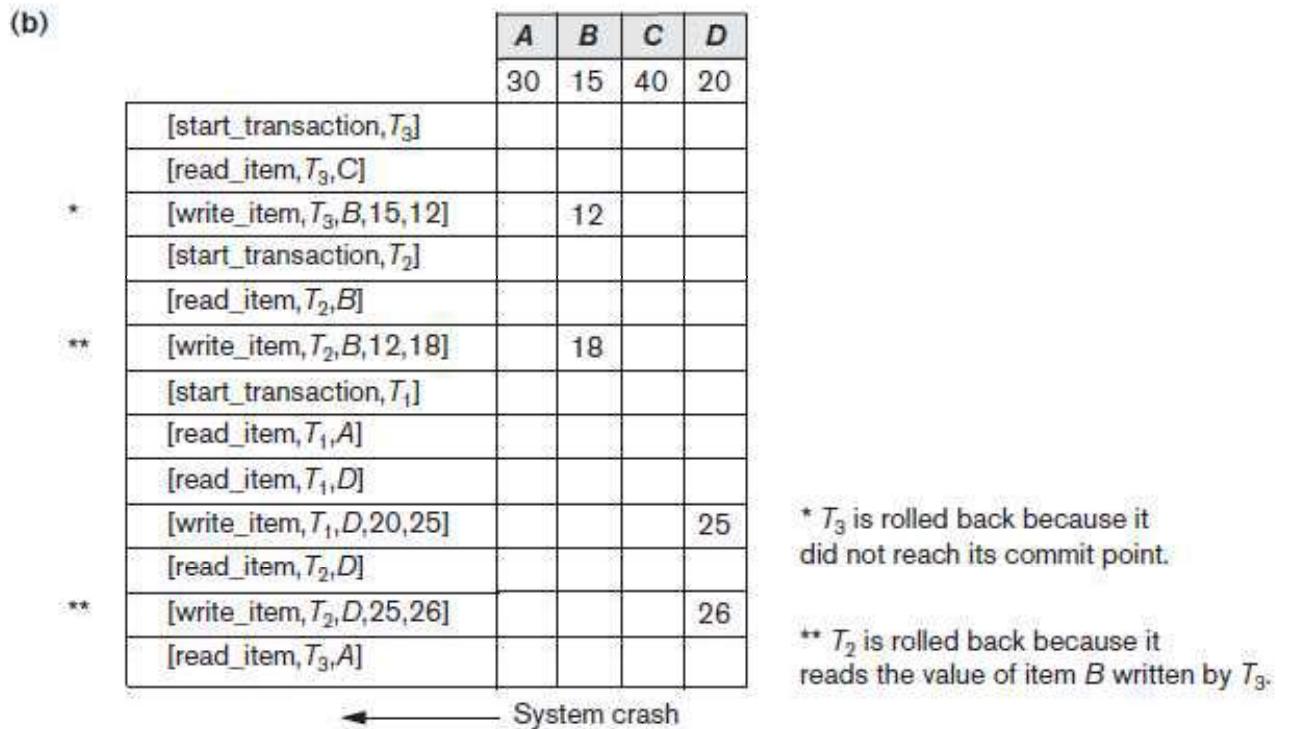
### 5.17.5 Transaction Rollback and Cascading Rollback

- 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)
- The undo-type log entries are used to restore the old values of data items that must be rolled back

- 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

(a)

$T_1$	$T_2$	$T_3$
read_item( $A$ )	read_item( $B$ )	read_item( $C$ )
read_item( $D$ )	write_item( $B$ )	write_item( $B$ )
write_item( $D$ )	read_item( $D$ )	read_item( $A$ )
		*
		write_item( $A$ )



**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.

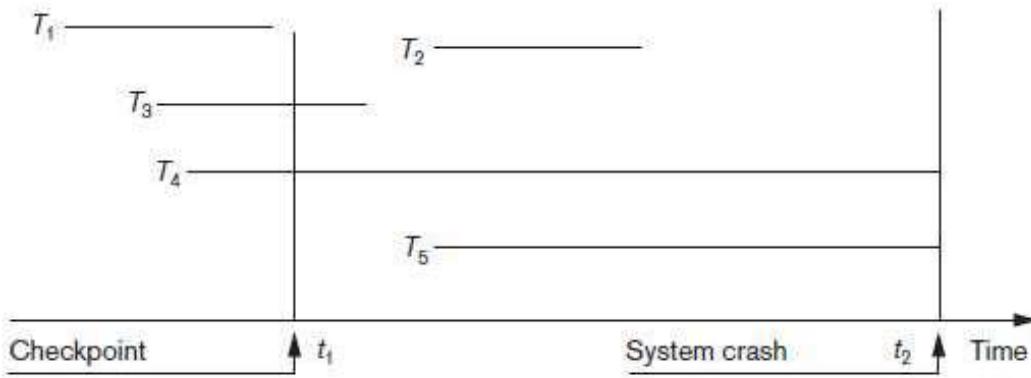
### 5.17.6 Transaction Actions That Do Not Affect the Database

- 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.
- If such erroneous reports are produced, part of the recovery process would have to inform the user that these reports are wrong, since the user may take an action based on these reports that affects the database.
- Hence, such reports should be generated only *after the transaction reaches its commit point*.
- A common method of dealing with such actions is to issue the commands that generate the reports but keep them as batch jobs, which are executed only after the transaction reaches its commit point. If the transaction fails, the batch jobs are canceled.

## 5.18 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 forcewritten 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.

- For multiuser systems with concurrency control, the concurrency control and recovery processes are interrelated.
- 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 as follows:
  - **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.



**Figure 23.2**  
An example of a recovery timeline to illustrate the effect of checkpointing.

- Figure 23.3 shows an example of recovery for a multiuser system that utilizes the recovery and concurrency control method

(a)

$T_1$	$T_2$	$T_3$	$T_4$
read_item(A)	read_item(B)	read_item(A)	read_item(B)
read_item(D)	write_item(B)	write_item(A)	write_item(B)
write_item(D)	read_item(D)	read_item(C)	read_item(A)
		write_item(D)	write_item(C)

(b)

[start_transaction, $T_1$ ]
[write_item, $T_1$ , D, 20]
[commit, $T_1$ ]
[checkpoint]
[start_transaction, $T_4$ ]
[write_item, $T_4$ , B, 15]
[write_item, $T_4$ , A, 20]
[commit, $T_4$ ]
[start_transaction, $T_2$ ]
[write_item, $T_2$ , B, 12]
[start_transaction, $T_3$ ]
[write_item, $T_3$ , A, 30]
[write_item, $T_2$ , D, 25]

System crash

$T_2$  and  $T_3$  are ignored because they did not reach their commit points.

$T_4$  is redone because its commit point is after the last system checkpoint.

**Figure 23.3**

An example of recovery using deferred update with concurrent transactions. (a) The READ and WRITE operations of four transactions. (b) System log at the point of crash.

- The method's main benefit is that transaction operations *never need to be undone*, for two reasons:
  1. A transaction does not record any changes in the database on disk until after it reaches its commit point—that is, until it completes its execution successfully. Hence, a transaction is never rolled back because of failure during transaction execution.
  2. A transaction will never read the value of an item that is written by an uncommitted transaction, because items remain locked until a transaction reaches its commit point. Hence, no cascading rollback will occur.

## 5.19 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.
- 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.
- If the transaction is allowed to commit before all its changes are written to the database, we have the most general case, known as the **UNDO/REDO recovery algorithm**. In this case, the **steal/no-force strategy** is applied.
- When concurrent execution is permitted, the recovery process again depends on the protocols used for concurrency control.
- The procedure RIU\_M (Recovery using Immediate Updates for a Multiuser environment) outlines a recovery algorithm for concurrent transactions with immediate update (UNDO/REDO recovery).
- **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:

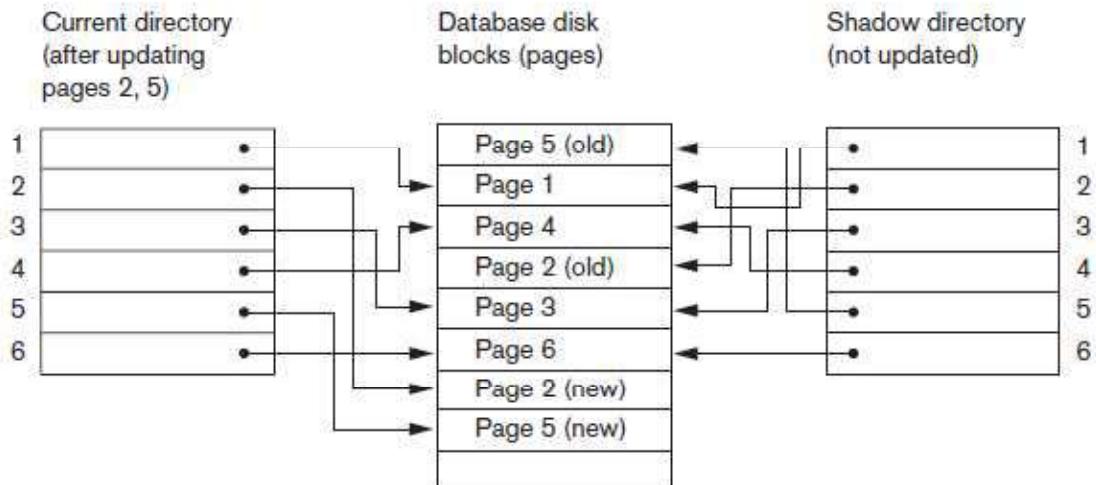
**Procedure UNDO (WRITE\_OP).** Undoing a write\_item operation write\_op consists of examining its log entry [write\_item,  $T$ ,  $X$ , old\_value, new\_value] and setting the value of item  $X$  in the database to old\_value, which is the before image (BFIM). Undoing a number of write\_item operations 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.

## 5.20 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 fixedsize disk pages (or disk blocks)—say,  $n$ —for recovery purposes
- A **directory** with  $n$  entries<sup>5</sup> is constructed, where the  $i$ th entry points to the  $i$ th 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.



- To recover from a failure during transaction execution, it is sufficient to free the modified database pages and to discard the current directory.
- The state of the database before transaction execution is available through the shadow directory, and that state is recovered by reinstating the shadow directory.
- Since recovery involves neither undoing nor redoing data items, this technique can be categorized as a NOUNDO/ NO-REDO technique for recovery.
- Disadvantage of shadow paging :
  - the updated database pages change location on disk
  - if the directory is large, the overhead of writing shadow directories to disk as transactions commit is significant
  - A further complication is how to handle **garbage collection** when a transaction commits
  - Another issue is that the operation to migrate between current and shadow directories must be implemented as an atomic operation.

## 5.21 The ARIES Recovery Algorithm

- It is used in many relational database-related products of IBM.
- ARIES uses a steal/no-force approach for writing, and it is based on three concepts:
  - write-ahead logging
  - repeating history during redo, and
  - logging changes during undo.

- **repeating history**, means that ARIES will retrace all actions of the database system prior to the crash to reconstruct the database state *when the crash occurred*. Transactions that were uncommitted at the time of the crash (active transactions) are undone.
- **logging during undo**, will prevent ARIES from repeating the completed undo operations if a failure occurs during recovery, which causes a restart of the recovery process.
- The ARIES recovery procedure consists of three main steps:
  1. Analysis
  2. REDO
  3. UNDO.
- **The analysis step**
  - identifies the dirty (updated) pages in the buffer and the set of transactions active at the time of the crash
  - The appropriate point in the log where the REDO operation should start is also determined
- **The REDO phase**
  - reapplies updates from the log to the database.
  - Certain information in the ARIES log will provide the start point for REDO, from which REDO operations are applied until the end of the log is reached
- **The UNDO phase**
  - the log is scanned backward and the operations of transactions that were active at the time of the crash are undone in reverse order.
- The information needed for ARIES to accomplish its recovery procedure includes the log, the Transaction Table, and the Dirty Page Table. Additionally, checkpointing is used.
- These tables are maintained by the transaction manager and written to the log during checkpointing.
- In ARIES, every log record has an associated **log sequence number (LSN)** that is monotonically increasing and indicates the address of the log record on disk.
- Each LSN corresponds to a *specific change* (action) of some transaction.
- Besides the log, two tables are needed for efficient recovery: the **Transaction Table** and the **Dirty Page Table**, which are maintained by the transaction manager.
- When a crash occurs, these tables are rebuilt in the analysis phase of recovery.

- The Transaction Table contains an entry for *each active transaction*, with information such as the transaction ID, transaction status, and the LSN of the most recent log record for the transaction.
- The Dirty Page Table contains an entry for each dirty page in the buffer, which includes the page ID and the LSN corresponding to the earliest update to that page.
- **Checkpointing** in ARIES consists of the following: writing a `begin_checkpoint` record to the log, writing an `end_checkpoint` record to the log, and writing *the LSN of the begin\_checkpoint record* to a special file.
- This special file is accessed during recovery to locate the last checkpoint information
- After a crash, the ARIES recovery manager takes over. Information from the last checkpoint is first accessed through the special file.
- The **analysis phase** starts at the `begin_checkpoint` record and proceeds to the end of the log
- The **REDO phase** follows next. To reduce the amount of unnecessary work, ARIES starts redoing at a point in the log where it knows (for sure) that previous changes to dirty pages *have already been applied to the database on disk*.

## 5.22 Database Backup and Recovery from Catastrophic 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 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.

- 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.

### 5.23 Assignment Questions

1. Explain properties of a transaction with state transition diagram.
2. Discuss the problems that can occur when concurrent transactions are executed.
3. Discuss the different types of failures. What is meant by catastrophic failure?
4. Discuss the actions taken by the read\_item and write\_item operations on a database.
5. What is two-phase locking protocol? How does it guarantee serializability?
6. What is a schedule? Explain with example serial, non serial and conflict serializable schedules.
7. Write short notes on
  1. Write ahead log protocol
  2. Time stamp Ordering
  3. Two phase locking protocol
8. Discuss the problems of deadlock and starvation, and the different approaches to dealing with these problems.
9. Describe the wait-die and wound-wait protocols for deadlock prevention.
10. Discuss the deferred update technique of recovery. What are the advantages and disadvantages of this technique?
11. Describe the shadow paging recovery technique.
12. Describe the three phases of the ARIES recovery method.

### 5.24 Expected Outcome

- ❖ To execute transactions by creating schedules
- ❖ To obtain equivalent and serializable schedules to avoid anomalies.
- ❖ To check whether the given schedule is serializable or not.
- ❖ To study locking protocols
- ❖ To improve resource utilization by applying various forms of locking protocol.

## 5.25 Further Reading

- ❖ <https://www.smartdraw.com/entity-relationship-diagram/>
  - ❖ [https://en.wikipedia.org/wiki/Database\\_normalization](https://en.wikipedia.org/wiki/Database_normalization)
  - ❖ [www.databasteknik.se/webbkursen/relalg-lecture](http://www.databasteknik.se/webbkursen/relalg-lecture)
  - ❖ [https://technet.microsoft.com/en-us/library/bb264565\(v=sql.90\).aspx](https://technet.microsoft.com/en-us/library/bb264565(v=sql.90).aspx)
  - ❖ [pages.cs.wisc.edu/~dbbook/openAccess/thirdEdition/.../Ch16\\_Overview\\_Xacts.pdf](http://pages.cs.wisc.edu/~dbbook/openAccess/thirdEdition/.../Ch16_Overview_Xacts.pdf)
-