

 The picture can't be displayed.

# **Chapter 9 & 10: Transactions Recovery System and Concurrency Control**



# Transaction Concept

- A **transaction** is a *unit* of program execution that accesses and possibly updates various data items.
- Most database system ensures that users accessing s database and manipulating data do not make the inconsistent i.e. A transaction must see a consistent database. They accomplish this by using the concept of transaction.
- During transaction execution the database may be inconsistent. When the transaction is committed, the database must be consistent.
- Beside ensuring consistency, transactions have many other properties which are enforced by database systems. For example, a transaction cannot be executed partially- it is either executed in its entirety or not at all.
- Transactions allow a group of statements to be executed as one logical “atomic” action.

# Transaction Concept

- Transactions allow multiple users to access and update a database simultaneously by guaranteeing that their actions will not interfere with each other.
- Finally, if a transaction executes successfully, the database system guarantees that the changes made by transaction will be reflected in the database even if the database system crashes immediately after its execution.
- Two main issues to deal with:
  - ★ Failures of various kinds, such as hardware failures and system crashes
  - ★ Concurrent execution of multiple transactions
- Transaction in SQL

START TRANSACTION;

SQL Statement;

COMMIT;

# ACID Properties

To preserve integrity of data, the database system must ensure:

- **Atomicity.** Either all operations of the transaction are properly reflected in the database or none are.
- **Consistency.** Execution of a transaction in isolation preserves the consistency of the database.
- **Isolation.** Although multiple transactions may execute concurrently, each transaction must be unaware of other concurrently executing transactions. Intermediate transaction results must be hidden from other concurrently executed transactions.
  - ★ That is, for every pair of transactions  $T_i$  and  $T_j$ , it appears to  $T_i$  that either  $T_j$  finished execution before  $T_i$  started, or  $T_j$  started execution after  $T_i$  finished.
- **Durability.** After a transaction completes successfully, the changes it has made to the database persist, even if there are system failures.

# Example of Fund Transfer

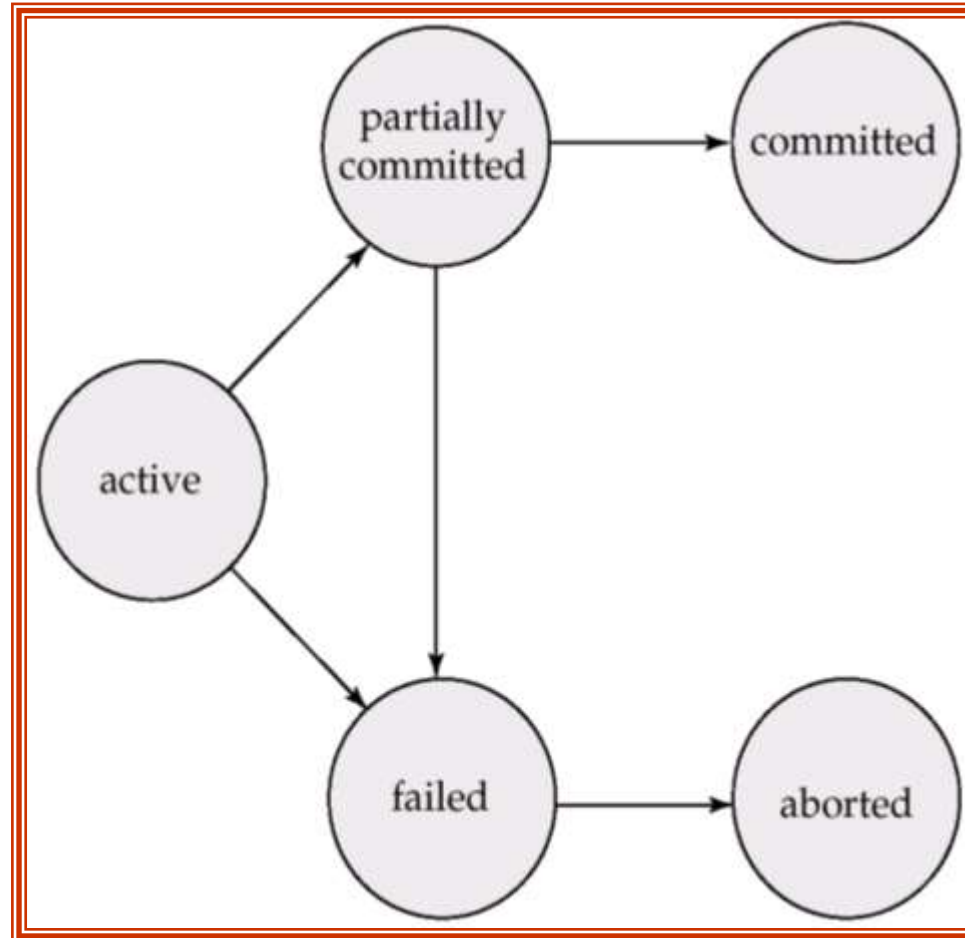
- Transaction to transfer \$50 from account  $A$  to account  $B$ :
  1. **read**( $A$ )
  2.  $A := A - 50$
  3. **write**( $A$ )
  4. **read**( $B$ )
  5.  $B := B + 50$
  6. **write**( $B$ )
- Consistency requirement – the sum of  $A$  and  $B$  is unchanged by the execution of the transaction.
- Atomicity requirement — if the transaction fails after step 3 and before step 6, the system should ensure that its updates are not reflected in the database, else an inconsistency will result.

# Example of Fund Transfer (Cont.)

- Durability requirement — once the user has been notified that the transaction has completed (i.e., the transfer of the \$50 has taken place), the updates to the database by the transaction must persist despite failures.
- Isolation requirement — if between steps 3 and 6, another transaction is allowed to access the partially updated database, it will see an inconsistent database (the sum  $A + B$  will be less than it should be).  
Can be ensured trivially by running transactions *serially*, that is one after the other. However, executing multiple transactions concurrently has significant benefits, as we will see.

# Transaction State

- **Active**, the initial state; the transaction stays in this state while it is executing
- **Partially committed**, after the final statement has been executed.
- **Failed**, after the discovery that normal execution can no longer proceed.
- **Aborted**, after the transaction has been rolled back and the database restored to its state prior to the start of the transaction. Two options after it has been aborted:
  - ★ restart the transaction – only if no internal logical error
  - ★ kill the transaction
- **Committed**, after *successful completion*.



# Concurrent Executions

- The process of managing simultaneous operations on the database without having them interfere with one another.
- DBMSs are required to allow simultaneous, (concurrent), multi-user access
  - multiple users within the one application accessing the same files.
  - users in different applications accessing the same files.
- Multiple transactions are allowed to run concurrently in the system.  
Advantages are:
  - ★ **increased processor and disk utilization**, leading to better transaction *throughput*: one transaction can be using the CPU while another is reading from or writing to the disk
  - ★ **reduced average response time** for transactions: short transactions need not wait behind long ones.



- *Concurrency control schemes* – mechanisms to achieve isolation, i.e., to control the interaction among the concurrent transactions in order to prevent them from destroying the consistency of the database
- Without centralized concurrency control, two types of problem can occur:
- The lost update problem.
- The so-called “dirty-read” problem, also described as allowing a transaction to view the partial, (uncommitted), results of another transaction.

# Concurrency Control in a Multi-user Environment The 'lost update' problem - sample scenario

- Consider the relation:

PRODUCT(prod-code,description,unit-price, quantity-onhand)

- Assume a particular record in the file is:

|     |              |       |    |
|-----|--------------|-------|----|
| P45 | Board-marker | 25.50 | 35 |
|-----|--------------|-------|----|

- Concurrent User 1, (in the Warehouse), has just received a shipment of twenty P45s.
- Concurrent User 2, (a salesman), has just received an order for ten P45s.
- CU1, (Concurrent User 1), and CU2 initiate transactions which each read the same P45 record. In each CWA (current work area), the field Q-O-H, Quantity-On-Hand, shows 35.
- CU2's transaction reduces Q-O-H by 10, and rewrites the record, (overwriting the old record), with Q-O-H of 25.
- CU1's transaction, a moment later, increments Q-O-H by 20 to 55, and rewrites the record, overwriting the existing record, (the one just written by CU2's transaction).
- Quantity-On-Hand thus shows 55, when it should show
  - $35 - 10 + 20$ , i.e. 45. CU2's update has been "lost".

# Concurrency Control in a Multi-user Environment The “dirty-read” problem - similar sample scenario

- Consider again the relation:  
PRODUCT(prod-code,description,unit-price, quantity-onhand)
- Assume again a particular record in the file is:  

|     |              |       |    |
|-----|--------------|-------|----|
| P45 | Board-marker | 25.50 | 35 |
|-----|--------------|-------|----|
- Concurrent User 1, (in the Warehouse), has just received a shipment of twenty P54s. (Note: P54s – not P45s.)
- Concurrent User 2, (a salesman), has just received an order for ten P45s.
- CU1 initiates the Warehouse update transaction, but by mistake keys in “P45”, instead of “ P54”. The transaction reads the P45 record with Q-O-H of 35, updates it to 55, and rewrites it (in a shared memory area for example).
- A moment later, CU2’s transaction reads CU1’s update showing the value of Q-O-H as 55.
- Just then, CU1, realizing the error, aborts the transaction, which executes rollback.
- A moment later, CU2’s transaction, having decremented QO- H, rewrites it as 45.
- Quantity-On-Hand thus shows 45, when it should show  $35 - 10 = 25$ . By viewing the uncommitted results of CU1’s

# Schedules

- *Schedules* – sequences that indicate the chronological order in which instructions of concurrent transactions are executed
  - ★ a schedule for a set of transactions must consist of all instructions of those transactions
  - ★ must preserve the order in which the instructions appear in each individual transaction.

## Example Schedules

- Let  $T_1$  transfer \$50 from  $A$  to  $B$ , and  $T_2$  transfer 10% of the balance from  $A$  to  $B$ . The following is a serial schedule (Schedule 1 in the text), in which  $T_1$  is followed by  $T_2$ .

| $T_1$   | $T_2$   |
|---|---|
| read( $A$ )<br>$A := A - 50$<br>write ( $A$ )<br>read( $B$ )<br>$B := B + 50$<br>write( $B$ ) | read( $A$ )<br>$temp := A * 0.1$<br>$A := A - temp$<br>write( $A$ )<br>read( $B$ )<br>$B := B + temp$<br>write( $B$ ) |

## Example Schedule (Cont.)

- Let  $T_1$  and  $T_2$  be the transactions defined previously. The following schedule (Schedule 3 in the text) is not a serial schedule, but it is *equivalent* to Schedule 1.

| $T_1$                                | $T_2$   |
|--------------------------------------|---|
| read(A)<br>$A := A - 50$<br>write(A) | read(A)<br>$temp := A * 0.1$<br>$A := A - temp$<br>write(A) |
| read(B)<br>$B := B + 50$<br>write(B) | read(B)<br>$B := B + temp$<br>write(B)                      |

In both Schedule 1 and 3, the sum  $A + B$  is preserved.

## Example Schedules (Cont.)

- The following concurrent schedule (Schedule 4 in the text) does not preserve the value of the the sum  $A + B$ .

| $T_1$   | $T_2$   |
|---|---|
| $\text{read}(A)$<br>$A := A - 50$   | $\text{read}(A)$<br>$\text{temp} := A * 0.1$<br>$A := A - \text{temp}$<br>$\text{write}(A)$<br>$\text{read}(B)$ |
| $\text{write}(A)$<br>$\text{read}(B)$<br>$B := B + 50$<br>$\text{write}(B)$ | $B := B + \text{temp}$<br>$\text{write}(B)$   |

# Serializability

- Basic Assumption – Each transaction preserves database consistency.
- Thus serial execution of a set of transactions preserves database consistency.
- A (possibly concurrent) schedule is serializable if it is equivalent to a serial schedule. Different forms of schedule equivalence give rise to the notions of:
  1. **conflict serializability**
  2. **view serializability**
- We ignore operations other than **read** and **write** instructions, and we assume that transactions may perform arbitrary computations on data in local buffers in between reads and writes. Our simplified schedules consist of only **read** and **write** instructions.



# Conflict Serializability

- Instructions  $I_i$  and  $I_j$  of transactions  $T_i$  and  $T_j$  respectively, **conflict** if and only if there exists some item  $Q$  accessed by both  $I_i$  and  $I_j$ , and at least one of these instructions wrote  $Q$ .
  1.  $I_i = \text{read}(Q)$ ,  $I_j = \text{read}(Q)$ .  $I_i$  and  $I_j$  don't conflict.
  2.  $I_i = \text{read}(Q)$ ,  $I_j = \text{write}(Q)$ . They conflict.
  3.  $I_i = \text{write}(Q)$ ,  $I_j = \text{read}(Q)$ . They conflict
  4.  $I_i = \text{write}(Q)$ ,  $I_j = \text{write}(Q)$ . They conflict
- Intuitively, a conflict between  $I_i$  and  $I_j$  forces a (logical) temporal order between them. If  $I_i$  and  $I_j$  are consecutive in a schedule and they do not conflict, their results would remain the same even if they had been interchanged in the schedule.

## Conflict Serializability (Cont.)

- If a schedule  $S$  can be transformed into a schedule  $S'$  by a series of swaps of non-conflicting instructions, we say that  $S$  and  $S'$  are **conflict equivalent**.
- We say that a schedule  $S$  is **conflict serializable** if it is conflict equivalent to a serial schedule
- Example of a schedule that is not conflict serializable:

| $T_3$                | $T_4$                |
|----------------------|----------------------|
| <b>read</b> ( $Q$ )  |                      |
| <b>write</b> ( $Q$ ) | <b>write</b> ( $Q$ ) |

We are unable to swap instructions in the above schedule to obtain either the serial schedule  $\langle T_3, T_4 \rangle$ , or the serial schedule  $\langle T_4, T_3 \rangle$ .

# Conflict Serializability (Cont.)

- Schedule 3 below can be transformed into Schedule 1, a serial schedule where  $T_2$  follows  $T_1$ , by series of swaps of non-conflicting instructions. Therefore Schedule 3 is conflict serializable.

| $T_1$                       | $T_2$                       |
|-----------------------------|-----------------------------|
| read( $A$ )<br>write( $A$ ) |                             |
|                             | read( $A$ )<br>write( $A$ ) |
| read( $B$ )<br>write( $B$ ) |                             |
|                             | read( $B$ )<br>write( $B$ ) |

# View Serializability

- Let  $S$  and  $S'$  be two schedules with the same set of transactions.  $S$  and  $S'$  are **view equivalent** if the following three conditions are met:
  1. For each data item  $Q$ , if transaction  $T_i$  reads the initial value of  $Q$  in schedule  $S$ , then transaction  $T_i$  must, in schedule  $S'$ , also read the initial value of  $Q$ .
  2. For each data item  $Q$  if transaction  $T_i$  executes **read**( $Q$ ) in schedule  $S$ , and that value was produced by transaction  $T_j$  (if any), then transaction  $T_i$  must in schedule  $S'$  also read the value of  $Q$  that was produced by transaction  $T_j$ .
  3. For each data item  $Q$ , the transaction (if any) that performs the final **write**( $Q$ ) operation in schedule  $S$  must perform the final **write**( $Q$ ) operation in schedule  $S'$ .

As can be seen, view equivalence is also based purely on **reads** and **writes** alone.

## View Serializability (Cont.)

- A schedule  $S$  is **view serializable** if it is view equivalent to a serial schedule.
- Every conflict serializable schedule is also view serializable.
- Schedule 9 (from text) — a schedule which is view-serializable but *not* conflict serializable.

| $T_3$        | $T_4$        | $T_6$        |
|--------------|--------------|--------------|
| read( $Q$ )  | write( $Q$ ) |              |
| write( $Q$ ) |              |              |
|              |              | write( $Q$ ) |

- Every view serializable schedule that is not conflict serializable has **blind writes**.

## Other Notions of Serializability

- Schedule 8 (from text) given below produces same outcome as the serial schedule  $\langle T_1, T_5 \rangle$ , yet is not conflict equivalent or view equivalent to it.

| $T_1$                                | $T_5$                                |
|--------------------------------------|--------------------------------------|
| read(A)<br>$A := A - 50$<br>write(A) |                                      |
|                                      | read(B)<br>$B := B - 10$<br>write(B) |
| read(B)<br>$B := B + 50$<br>write(B) |                                      |
|                                      | read(A)<br>$A := A + 10$<br>write(A) |

- Determining such equivalence requires analysis of operations other than read and write.

# Lock-Based Protocols

- A lock is a mechanism to control concurrent access to a data item
- Data items can be locked in two modes :
  1. *exclusive* (X) mode. Data item can be both read as well as written. X-lock is requested using **lock-X** instruction.
  2. *shared* (S) mode. Data item can only be read. S-lock is requested using **lock-S** instruction.
- Lock requests are made to concurrency-control manager. Transaction can proceed only after request is granted.

# Lock-Based Protocols (Cont.)

## ■ Lock-compatibility matrix

|   | S     | X     |
|---|-------|-------|
| S | true  | false |
| X | false | false |

- A transaction may be granted a lock on an item if the requested lock is compatible with locks already held on the item by other transactions
- Any number of transactions can hold shared locks on an item, but if any transaction holds an exclusive on the item no other transaction may hold any lock on the item.
- If a lock cannot be granted, the requesting transaction is made to wait till all incompatible locks held by other transactions have been released. The lock is then granted.



# Lock-Based Protocols (Cont.)

- Example of a transaction performing locking:

```
 $T_2$ : lock-S(A);  
      read (A);  
      unlock(A);  
      lock-S(B);  
      read (B);  
      unlock(B);  
      display(A+B)
```

- Locking as above is not sufficient to guarantee serializability — if  $A$  and  $B$  get updated in-between the read of  $A$  and  $B$ , the displayed sum would be wrong.
- A **locking protocol** is a set of rules followed by all transactions while requesting and releasing locks. Locking protocols restrict the set of possible schedules.

# Pitfalls of Lock-Based Protocols

- Consider the partial schedule

| $T_3$         | $T_4$         |
|---------------|---------------|
| lock-X( $B$ ) |               |
| read( $B$ )   |               |
| $B := B - 50$ |               |
| write( $B$ )  |               |
|               | lock-S( $A$ ) |
|               | read( $A$ )   |
|               | lock-S( $B$ ) |
| lock-X( $A$ ) |               |

- Neither  $T_3$  nor  $T_4$  can make progress — executing **lock-S( $B$ )** causes  $T_4$  to wait for  $T_3$  to release its lock on  $B$ , while executing **lock-X( $A$ )** causes  $T_3$  to wait for  $T_4$  to release its lock on  $A$ .
- Such a situation is called a **deadlock**.
  - ★ To handle a deadlock one of  $T_3$  or  $T_4$  must be rolled back and its locks released.

# Pitfalls of Lock-Based Protocols (Cont.)

- The potential for deadlock exists in most locking protocols. Deadlocks are a necessary evil.
- **Starvation** is also possible if concurrency control manager is badly designed. For example:
  - ★ A transaction may be waiting for an X-lock on an item, while a sequence of other transactions request and are granted an S-lock on the same item.
  - ★ The same transaction is repeatedly rolled back due to deadlocks.
- Concurrency control manager can be designed to prevent starvation.

# The Two-Phase Locking Protocol

- This is a protocol which ensures conflict-serializable schedules.
- Phase 1: Growing Phase
  - ★ transaction may obtain locks
  - ★ transaction may not release locks
- Phase 2: Shrinking Phase
  - ★ transaction may release locks
  - ★ transaction may not obtain locks
- The protocol assures serializability. It can be proved that the transactions can be serialized in the order of their **lock points** (i.e. the point where a transaction acquired its final lock).

# The Two-Phase Locking Protocol (Cont.)

- Two-phase locking *does not* ensure freedom from deadlocks
- Cascading roll-back is possible under two-phase locking. To avoid this, follow a modified protocol called **strict two-phase locking**. Here a transaction must hold all its exclusive locks till it commits/aborts.
- **Rigorous two-phase locking** is even stricter: here *all* locks are held till commit/abort. In this protocol transactions can be serialized in the order in which they commit.
- Let us consider the transaction:

R(A) W(A) R(B) W(B) R(C) W(C) R(D) W(D)

**strict two-phase locking:**

**X (A)** R(A) W(A) **X(B)** R(B) W(B) **X(C)** R(C) W(C) **X(D)** R(D) W(D)  
**U(A)U(B)U(C) U(D)**

**Rigorous two-phase locking**

**X (A) X(B) X(C) X(D)** R(A) W(A) R(B) W(B) R(C) W(C) R(D) W(D)  
**U(A)U(B)U(C) U(D)**

# The Two-Phase Locking Protocol (Cont.)

- There can be conflict serializable schedules that cannot be obtained if two-phase locking is used.
- However, in the absence of extra information (e.g., ordering of access to data), two-phase locking is needed for conflict serializability in the following sense:

Given a transaction  $T_i$  that does not follow two-phase locking, we can find a transaction  $T_j$  that uses two-phase locking, and a schedule for  $T_i$  and  $T_j$  that is not conflict serializable.

# Implementation of Locking

- A **Lock manager** can be implemented as a separate process to which transactions send lock and unlock requests
- The lock manager replies to a lock request by sending a lock grant messages (or a message asking the transaction to roll back, in case of a deadlock)
- The requesting transaction waits until its request is answered
- The lock manager maintains a datastructure called a **lock table** to record granted locks and pending requests
- The lock table is usually implemented as an in-memory hash table indexed on the name of the data item being locked

# Timestamp-Based Protocols

- Each transaction is issued a timestamp when it enters the system. If an old transaction  $T_i$  has time-stamp  $TS(T_i)$ , a new transaction  $T_j$  is assigned time-stamp  $TS(T_j)$  such that  $TS(T_i) < TS(T_j)$ .
- The protocol manages concurrent execution such that the time-stamps determine the serializability order.
- In order to assure such behavior, the protocol maintains for each data  $Q$  two timestamp values:
  - ★ **W-timestamp**( $Q$ ) is the largest time-stamp of any transaction that executed **write**( $Q$ ) successfully.
  - ★ **R-timestamp**( $Q$ ) is the largest time-stamp of any transaction that executed **read**( $Q$ ) successfully.



# Timestamp-Based Protocols (Cont.)

- The timestamp ordering protocol ensures that any conflicting **read** and **write** operations are executed in timestamp order.
- Suppose a transaction  $T_i$  issues a **read**( $Q$ )
  1. If  $TS(T_i) \leq \mathbf{W}\text{-timestamp}(Q)$ , then  $T_i$  needs to read a value of  $Q$  that was already overwritten. Hence, the **read** operation is rejected, and  $T_i$  is rolled back.
  2. If  $TS(T_i) \geq \mathbf{W}\text{-timestamp}(Q)$ , then the **read** operation is executed, and  $\mathbf{R}\text{-timestamp}(Q)$  is set to the maximum of  $\mathbf{R}\text{-timestamp}(Q)$  and  $TS(T_i)$ .

# Timestamp-Based Protocols (Cont.)

- Suppose that transaction  $T_i$  issues **write**( $Q$ ).
- If  $TS(T_i) < R\text{-timestamp}(Q)$ , then the value of  $Q$  that  $T_i$  is producing was needed previously, and the system assumed that that value would never be produced. Hence, the **write** operation is rejected, and  $T_i$  is rolled back.
- If  $TS(T_i) < W\text{-timestamp}(Q)$ , then  $T_i$  is attempting to write an obsolete value of  $Q$ . Hence, this **write** operation is rejected, and  $T_i$  is rolled back.
- Otherwise, the **write** operation is executed, and  $W\text{-timestamp}(Q)$  is set to  $TS(T_i)$ .

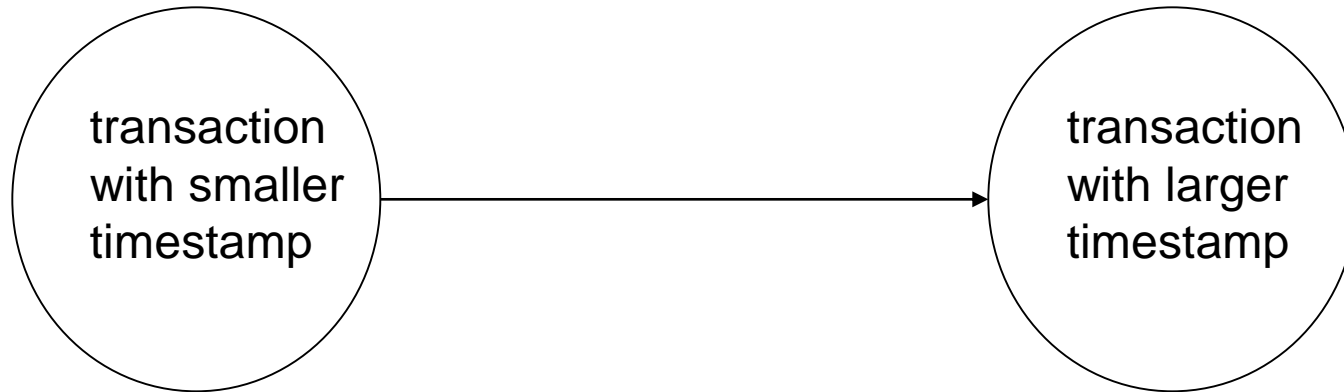
# Example Use of the Protocol

A partial schedule for several data items for transactions with timestamps 1, 2, 3, 4, 5

| $T_1$       | $T_2$                | $T_3$                        | $T_4$ | $T_5$                        |
|-------------|----------------------|------------------------------|-------|------------------------------|
| read( $Y$ ) | read( $Y$ )          | write( $Y$ )<br>write( $Z$ ) |       | read( $X$ )                  |
| read( $X$ ) | read( $X$ )<br>abort | write( $Z$ )<br>abort        |       | read( $Z$ )                  |
|             |                      |                              |       | write( $Y$ )<br>write( $Z$ ) |

# Correctness of Timestamp-Ordering Protocol

- The timestamp-ordering protocol guarantees serializability since all the arcs in the precedence graph are of the form:



Thus, there will be no cycles in the precedence graph

- Timestamp protocol ensures freedom from deadlock as no transaction ever waits.
- But the schedule may not be cascade-free, and may not even be recoverable.

# Recoverability and Cascade Freedom

## ■ Problem with timestamp-ordering protocol:

- ★ Suppose  $T_i$  aborts, but  $T_j$  has read a data item written by  $T_i$
- ★ Then  $T_j$  must abort; if  $T_j$  had been allowed to commit earlier, the schedule is not recoverable.
- ★ Further, any transaction that has read a data item written by  $T_j$  must abort
- ★ This can lead to cascading rollback --- that is, a chain of rollbacks

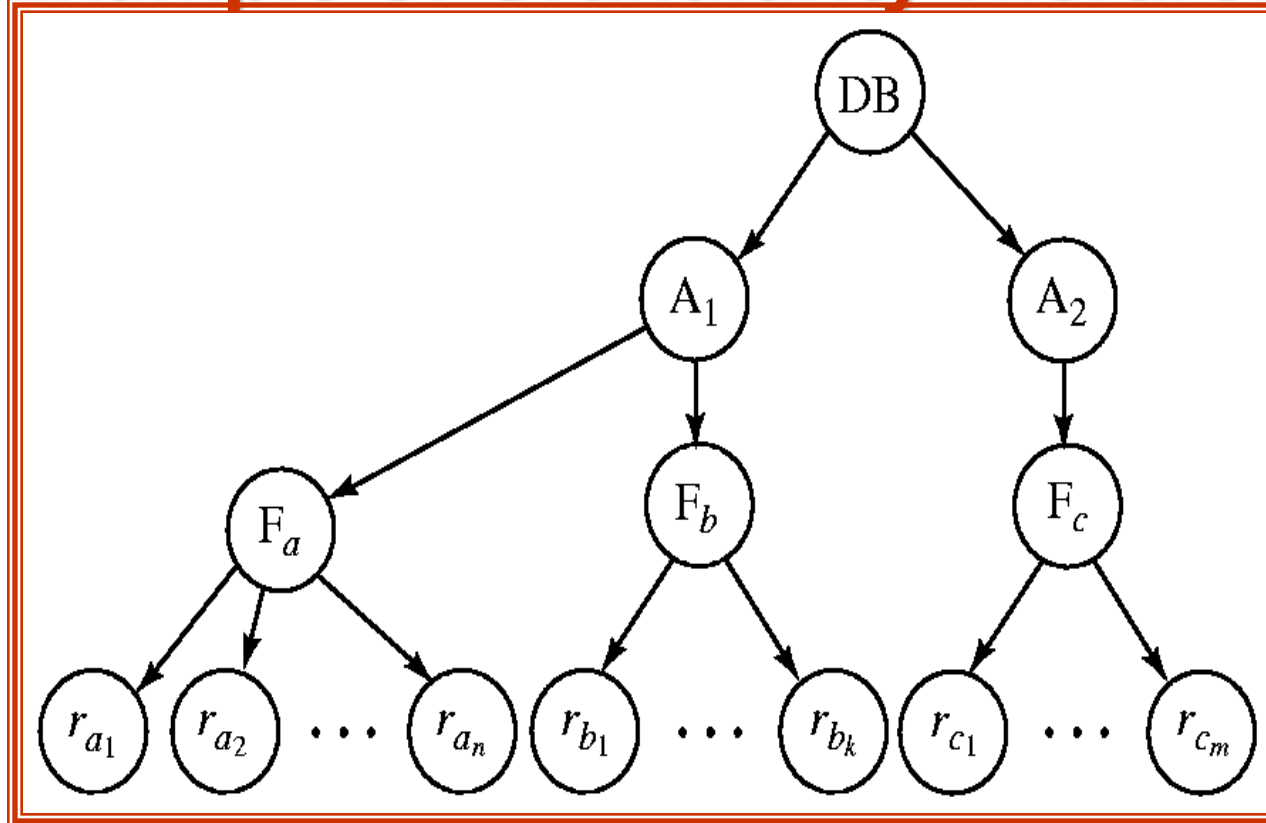
## ■ Solution:

- ★ A transaction is structured such that its writes are all performed at the end of its processing
- ★ All writes of a transaction form an atomic action; no transaction may execute while a transaction is being written
- ★ A transaction that aborts is restarted with a new timestamp

# Multiple Granularity

- Allow data items to be of various sizes and define a hierarchy of data granularities, where the small granularities are nested within larger ones
- Can be represented graphically as a tree (but don't confuse with tree-locking protocol)
- When a transaction locks a node in the tree *explicitly*, it *implicitly* locks all the node's descendents in the same mode.
- **Granularity of locking** (level in tree where locking is done):
  - ★ *fine granularity* (lower in tree): high concurrency, high locking overhead
  - ★ *coarse granularity* (higher in tree): low locking overhead, low concurrency

# Example of Granularity Hierarchy



The highest level in the example hierarchy is the entire database.  
The levels below are of type *area*, *file* and *record* in that order.

# Deadlock Handling

- Consider the following two transactions:

$T_1$ :    write ( $X$ )  
          write( $Y$ )

$T_2$ :    write( $Y$ )  
          write( $X$ )

- Schedule with deadlock

| $T_1$  | $T_2$  |
|--|--|
| <b>lock-X</b> on $X$<br>write ( $X$ )<br><br>wait for <b>lock-X</b> on $Y$ | <b>lock-X</b> on $Y$<br>write ( $X$ )<br>wait for <b>lock-X</b> on $X$ |



# Deadlock Handling

- System is deadlocked if there is a set of transactions such that every transaction in the set is waiting for another transaction in the set.
- **Deadlock prevention** protocols ensure that the system will *never* enter into a deadlock state. Some prevention strategies :
  - ★ Require that each transaction locks all its data items before it begins execution (predeclaration).
  - ★ Impose partial ordering of all data items and require that a transaction can lock data items only in the order specified by the partial order (graph-based protocol).

# More Deadlock Prevention Strategies

- Following schemes use transaction timestamps for the sake of deadlock prevention alone.
- **wait-die** scheme — non-preemptive
  - ★ older transaction may wait for younger one to release data item. Younger transactions never wait for older ones; they are rolled back instead.
  - ★ a transaction may die several times before acquiring needed data item
- **wound-wait** scheme — preemptive
  - ★ older transaction *wounds* (forces rollback) of younger transaction instead of waiting for it. Younger transactions may wait for older ones.
  - ★ may be fewer rollbacks than *wait-die* scheme.

# Deadlock prevention (Cont.)

- Both in *wait-die* and in *wound-wait* schemes, a rolled back transactions is restarted with its original timestamp. Older transactions thus have precedence over newer ones, and starvation is hence avoided.
- **Timeout-Based Schemes :**
  - ★ a transaction waits for a lock only for a specified amount of time. After that, the wait times out and the transaction is rolled back.
  - ★ thus deadlocks are not possible
  - ★ simple to implement; but starvation is possible. Also difficult to determine good value of the timeout interval.

# Recovery System

# Failure Classification

## ■ Transaction failure :

- ★ **Logical errors:** transaction cannot complete due to some internal error condition
- ★ **System errors:** the database system must terminate an active transaction due to an error condition (e.g., deadlock)

## ■ **System crash:** a power failure or other hardware or software failure causes the system to crash.

- ★ **Fail-stop assumption:** non-volatile storage contents are assumed to not be corrupted by system crash
  - Database systems have numerous integrity checks to prevent corruption of disk data

## ■ **Disk failure:** a head crash or similar disk failure destroys all or part of disk storage

- ★ Destruction is assumed to be detectable: disk drives use checksums to detect failures

# Recovery Algorithms

- Recovery algorithms are techniques to ensure database consistency and transaction atomicity and durability despite failures
  - ★ Focus of this chapter
- Recovery algorithms have two parts
  1. Actions taken during normal transaction processing to ensure enough information exists to recover from failures
  2. Actions taken after a failure to recover the database contents to a state that ensures atomicity, consistency and durability

# Storage Structure

## ■ Volatile storage:

- ★ does not survive system crashes
- ★ examples: main memory, cache memory

## ■ Nonvolatile storage:

- ★ survives system crashes
- ★ examples: disk, tape, flash memory,  
non-volatile (battery backed up) RAM

## ■ Stable storage:

- ★ a mythical form of storage that survives all failures
- ★ approximated by maintaining multiple copies on distinct nonvolatile media

# Recovery and Atomicity

- Modifying the database without ensuring that the transaction will commit may leave the database in an inconsistent state.
- Consider transaction  $T_i$  that transfers \$50 from account  $A$  to account  $B$ ; goal is either to perform all database modifications made by  $T_i$  or none at all.
- Several output operations may be required for  $T_i$  (to output  $A$  and  $B$ ). A failure may occur after one of these modifications have been made but before all of them are made.



# Recovery and Atomicity (Cont.)

- To ensure atomicity despite failures, we first output information describing the modifications to stable storage without modifying the database itself.
- We study two approaches:
  - ★ **log-based recovery**, and
  - ★ **shadow-paging**
- We assume (initially) that transactions run serially, that is, one after the other.

# Log-Based Recovery

- A **log** is kept on stable storage.
  - ★ The log is a sequence of **log records**, and maintains a record of update activities on the database.
- When transaction  $T_i$  starts, it registers itself by writing a  $\langle T_i \text{ start} \rangle$  log record
- Before  $T_i$  executes **write**( $X$ ), a log record  $\langle T_i, X, V_1, V_2 \rangle$  is written, where  $V_1$  is the value of  $X$  before the write, and  $V_2$  is the value to be written to  $X$ .
  - ★ Log record notes that  $T_i$  has performed a write on data item  $X_j$ .  $X_j$  had value  $V_1$  before the write, and will have value  $V_2$  after the write.
- When  $T_i$  finishes its last statement, the log record  $\langle T_i \text{ commit} \rangle$  is written.
- We assume for now that log records are written directly to stable storage (that is, they are not buffered)
- Two approaches using logs
  - ★ Deferred database modification
  - ★ Immediate database modification

# Deferred Database Modification

- The **deferred database modification** scheme records all modifications to the log, but defers all the **writes** to after partial commit.
- Assume that transactions execute serially
- Transaction starts by writing  $\langle T_i, \text{start} \rangle$  record to log.
- A **write**( $X$ ) operation results in a log record  $\langle T_i, X, V \rangle$  being written, where  $V$  is the new value for  $X$ 
  - ★ Note: old value is not needed for this scheme
- The write is not performed on  $X$  at this time, but is deferred.
- When  $T_i$  partially commits,  $\langle T_i, \text{commit} \rangle$  is written to the log
- Finally, the log records are read and used to actually execute the previously deferred writes.

# Deferred Database Modification (Cont.)

- During recovery after a crash, a transaction needs to be redone if and only if both  $\langle T_i \text{ start} \rangle$  and  $\langle T_i \text{ commit} \rangle$  are there in the log.
- Redoing a transaction  $T_i$  (**redo**  $T_i$ ) sets the value of all data items updated by the transaction to the new values.
- Crashes can occur while
  - ★ the transaction is executing the original updates, or
  - ★ while recovery action is being taken
- example transactions  $T_0$  and  $T_1$  ( $T_0$  executes before  $T_1$ ):

$T_0$ : **read** ( $A$ )

$A$ : -  $A - 50$

**Write** ( $A$ )

**read** ( $B$ )

$B$ :-  $B + 50$

**write** ( $B$ )

$T_1$  : **read** ( $C$ )

$C$ :-  $C - 100$

**write** ( $C$ )

# Deferred Database Modification (Cont.)

- Below we show the log as it appears at three instances of time.

|                                     |                                      |                                      |
|-------------------------------------|--------------------------------------|--------------------------------------|
| $\langle T_0 \text{ start} \rangle$ | $\langle T_0 \text{ start} \rangle$  | $\langle T_0 \text{ start} \rangle$  |
| $\langle T_0, A, 950 \rangle$       | $\langle T_0, A, 950 \rangle$        | $\langle T_0, A, 950 \rangle$        |
| $\langle T_0, B, 2050 \rangle$      | $\langle T_0, B, 2050 \rangle$       | $\langle T_0, B, 2050 \rangle$       |
|                                     | $\langle T_0 \text{ commit} \rangle$ | $\langle T_0 \text{ commit} \rangle$ |
|                                     | $\langle T_1 \text{ start} \rangle$  | $\langle T_1 \text{ start} \rangle$  |
|                                     | $\langle T_1, C, 600 \rangle$        | $\langle T_1, C, 600 \rangle$        |
|                                     |                                      | $\langle T_1 \text{ commit} \rangle$ |
| (a)                                 | (b)                                  | (c)                                  |

- If log on stable storage at time of crash is as in case:
  - (a) No redo actions need to be taken
  - (b) redo( $T_0$ ) must be performed since  $\langle T_0 \text{ commit} \rangle$  is present
  - (c) redo( $T_0$ ) must be performed followed by redo( $T_1$ ) since  $\langle T_0 \text{ commit} \rangle$  and  $\langle T_1 \text{ commit} \rangle$  are present

# Immediate Database Modification

- The **immediate database modification** scheme allows database updates of an uncommitted transaction to be made as the writes are issued
  - ★ since undoing may be needed, update logs must have both old value and new value
- Update log record must be written *before* database item is written
  - ★ We assume that the log record is output directly to stable storage
  - ★ Can be extended to postpone log record output, so long as prior to execution of an **output**( $B$ ) operation for a data block  $B$ , all log records corresponding to items  $B$  must be flushed to stable storage
- Output of updated blocks can take place at any time before or after transaction commit
- Order in which blocks are output can be different from the order in which they are written.

# Immediate Database Modification Example

| Log   | Write      | Output     |
|---|------------|------------|
| $\langle T_0 \text{ start} \rangle$           |            |            |
| $\langle T_0, A, 1000, 950 \rangle$           |            |            |
| $T_0, B, 2000, 2050$                          |            |            |
|   | $A = 950$  |            |
|   | $B = 2050$ |            |
| $\langle T_0 \text{ commit} \rangle$          |            |            |
| $\langle T_1 \text{ start} \rangle \quad X_1$ |            |            |
| $\langle T_1, C, 700, 600 \rangle$            |            |            |
|   | $C = 600$  |            |
|   |            | $B_B, B_C$ |
| $\langle T_1 \text{ commit} \rangle$          |            | $B_A$      |

■ Note:  $B_X$  denotes block containing X.

# Immediate Database Modification (Cont.)

- Recovery procedure has two operations instead of one:
  - ★ **undo**( $T_i$ ) restores the value of all data items updated by  $T_i$  to their old values, going backwards from the last log record for  $T_i$
  - ★ **redo**( $T_i$ ) sets the value of all data items updated by  $T_i$  to the new values, going forward from the first log record for  $T_i$
- Both operations must be **idempotent**
  - ★ That is, even if the operation is executed multiple times the effect is the same as if it is executed once
    - Needed since operations may get re-executed during recovery
- When recovering after failure:
  - ★ Transaction  $T_i$  needs to be undone if the log contains the record  $\langle T_i \text{ start} \rangle$ , but does not contain the record  $\langle T_i \text{ commit} \rangle$ .
  - ★ Transaction  $T_i$  needs to be redone if the log contains both the record  $\langle T_i \text{ start} \rangle$  and the record  $\langle T_i \text{ commit} \rangle$ .
- Undo operations are performed first, then redo operations.



# Immediate DB Modification Recovery Example

Below we show the log as it appears at three instances of time.

|                                      |                                      |                                      |
|--------------------------------------|--------------------------------------|--------------------------------------|
| $\langle T_0 \text{ start} \rangle$  | $\langle T_0 \text{ start} \rangle$  | $\langle T_0 \text{ start} \rangle$  |
| $\langle T_0, A, 1000, 950 \rangle$  | $\langle T_0, A, 1000, 950 \rangle$  | $\langle T_0, A, 1000, 950 \rangle$  |
| $\langle T_0, B, 2000, 2050 \rangle$ | $\langle T_0, B, 2000, 2050 \rangle$ | $\langle T_0, B, 2000, 2050 \rangle$ |
|                                      | $\langle T_0 \text{ commit} \rangle$ | $\langle T_0 \text{ commit} \rangle$ |
|                                      | $\langle T_1 \text{ start} \rangle$  | $\langle T_1 \text{ start} \rangle$  |
|                                      | $\langle T_1, C, 700, 600 \rangle$   | $\langle T_1, C, 700, 600 \rangle$   |
|                                      |                                      | $\langle T_1 \text{ commit} \rangle$ |
| (a)                                  | (b)                                  | (c)                                  |

Recovery actions in each case above are:

- (a) undo ( $T_0$ ): B is restored to 2000 and A to 1000.
- (b) undo ( $T_1$ ) and redo ( $T_0$ ): C is restored to 700, and then A and B are set to 950 and 2050 respectively.
- (c) redo ( $T_0$ ) and redo ( $T_1$ ): A and B are set to 950 and 2050 respectively. Then C is set to 600

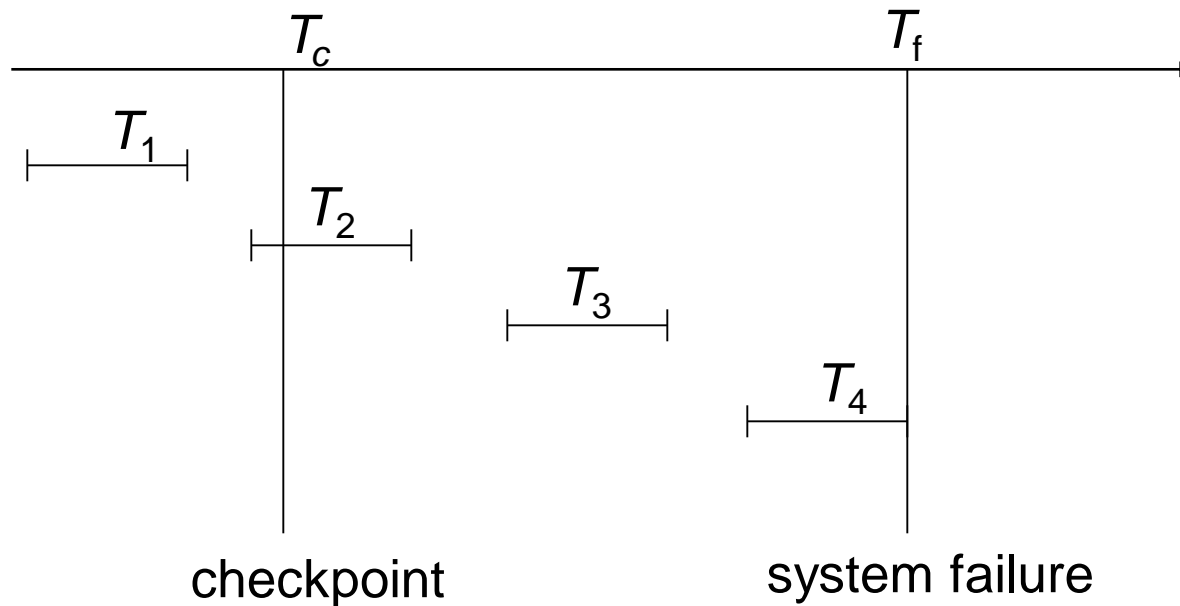
# Checkpoints

- Problems in recovery procedure as discussed earlier :
  1. searching the entire log is time-consuming
  2. we might unnecessarily redo transactions which have already
  3. output their updates to the database.
- Streamline recovery procedure by periodically performing **checkpointing**
  1. Output all log records currently residing in main memory onto stable storage.
  2. Output all modified buffer blocks to the disk.
  3. Write a log record < **checkpoint** > onto stable storage.

# Checkpoints (Cont.)

- During recovery we need to consider only the most recent transaction  $T_i$  that started before the checkpoint, and transactions that started after  $T_i$ .
  1. Scan backwards from end of log to find the most recent **<checkpoint>** record
  2. Continue scanning backwards till a record **< $T_i$  start>** is found.
  3. Need only consider the part of log following above **start** record. Earlier part of log can be ignored during recovery, and can be erased whenever desired.
  4. For all transactions (starting from  $T_i$  or later) with no **< $T_i$  commit>**, execute **undo( $T_i$ )**. (Done only in case of immediate modification.)
  5. Scanning forward in the log, for all transactions starting from  $T_i$  or later with a **< $T_i$  commit>**, execute **redo( $T_i$ )**.

# Example of Checkpoints

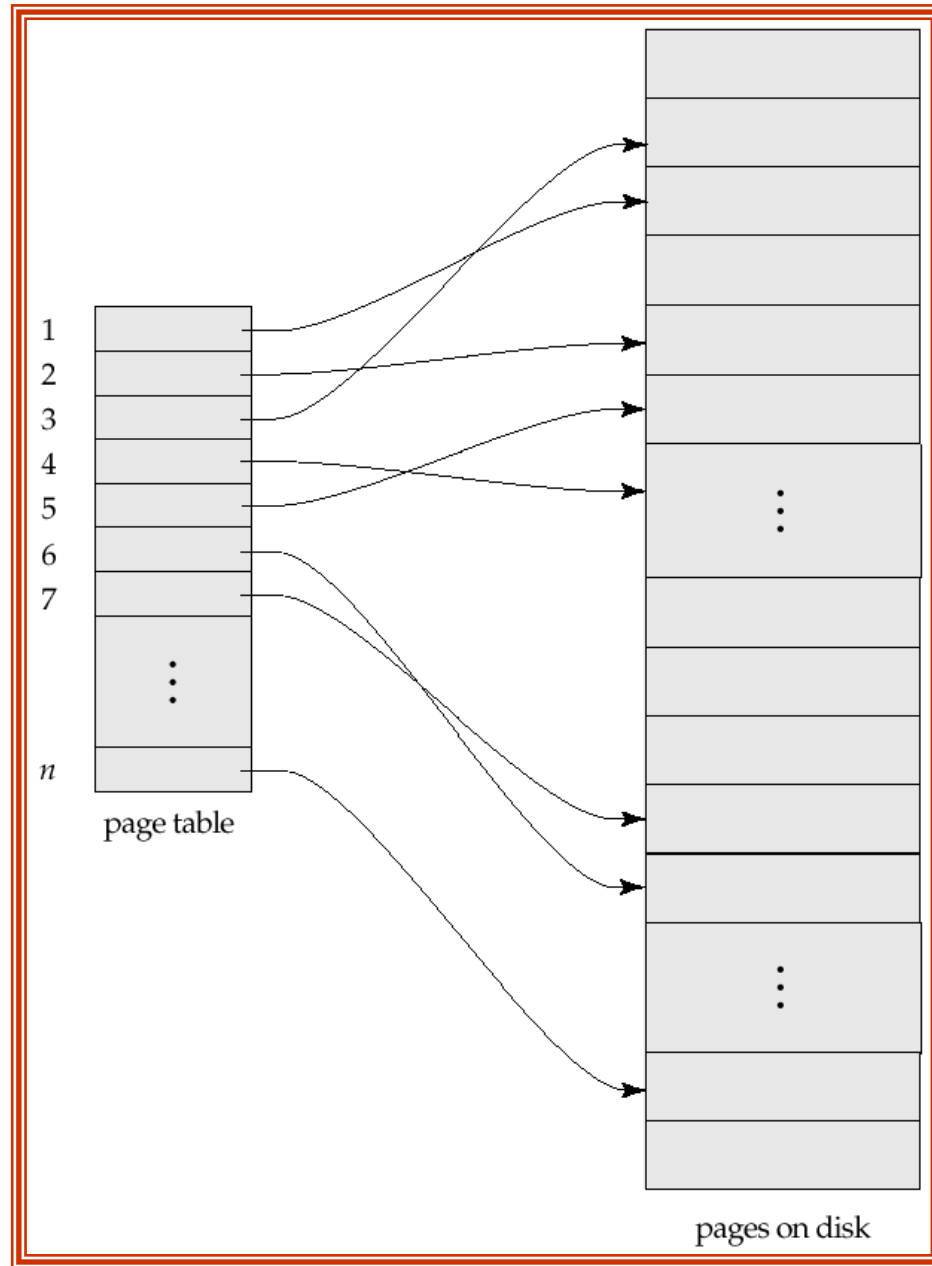


- $T_1$  can be ignored (updates already output to disk due to checkpoint)
- $T_2$  and  $T_3$  redone.
- $T_4$  undone

# Shadow Paging

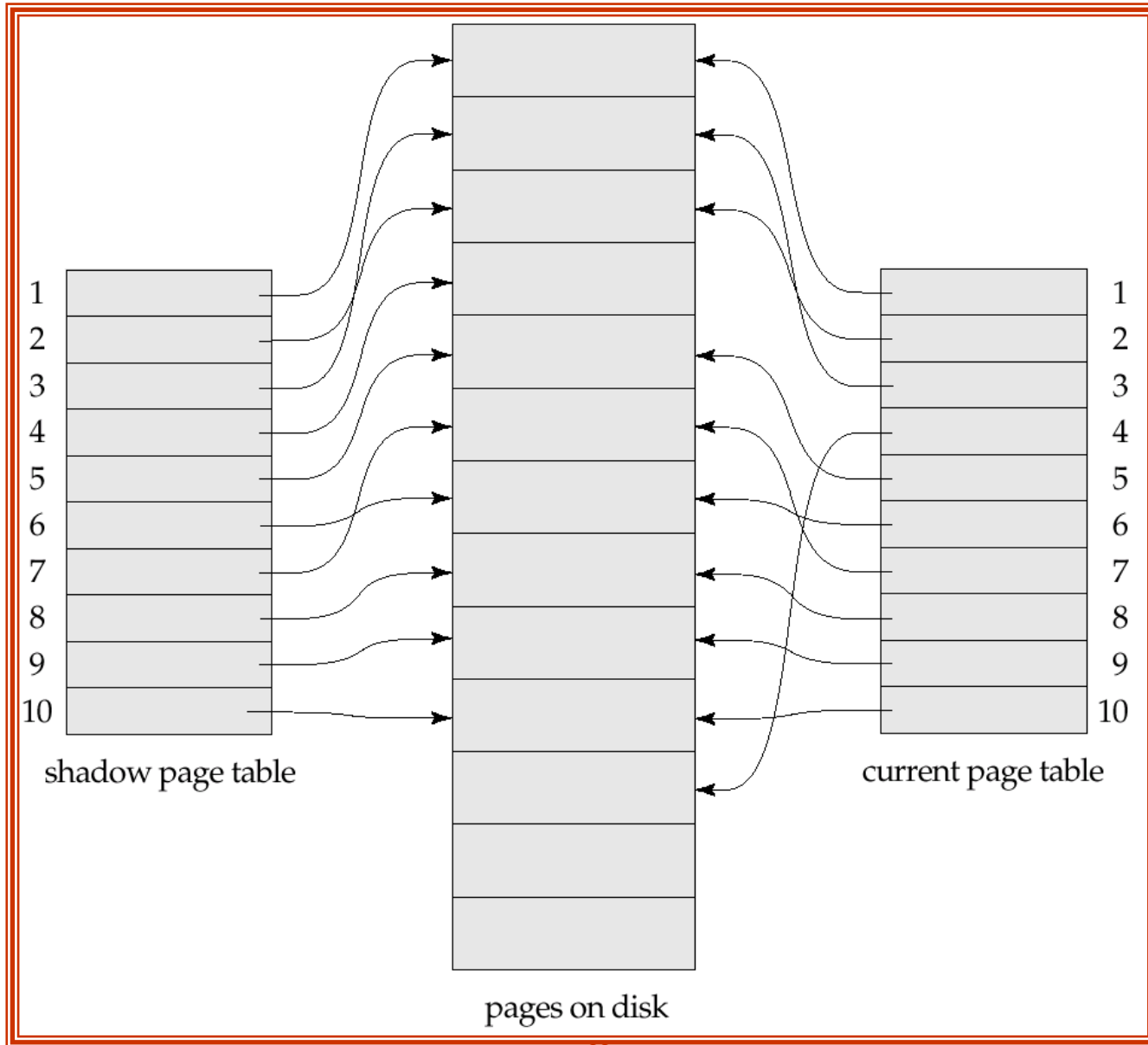
- **Shadow paging** is an alternative to log-based recovery; this scheme is useful if transactions execute serially
- Idea: maintain *two* page tables during the lifetime of a transaction – the **current page table**, and the **shadow page table**
- Store the shadow page table in nonvolatile storage, such that state of the database prior to transaction execution may be recovered.
  - ★ Shadow page table is never modified during execution
- To start with, both the page tables are identical. Only current page table is used for data item accesses during execution of the transaction.
- Whenever any page is about to be written for the first time
  - ★ A copy of this page is made onto an unused page.
  - ★ The current page table is then made to point to the copy
  - ★ The update is performed on the copy

# Sample Page Table



# Example of Shadow Paging

Shadow and current page tables after write to page 4



# Shadow Paging (Cont.)

- To commit a transaction :
  1. Flush all modified pages in main memory to disk
  2. Output current page table to disk
  3. Make the current page table the new shadow page table, as follows:
    - ★ keep a pointer to the shadow page table at a fixed (known) location on disk.
    - ★ to make the current page table the new shadow page table, simply update the pointer to point to current page table on disk
- Once pointer to shadow page table has been written, transaction is committed.
- No recovery is needed after a crash — new transactions can start right away, using the shadow page table.
- Pages not pointed to from current/shadow page table should be freed (garbage collected).



# Shadow Paging (Cont.)

- Advantages of shadow-paging over log-based schemes
  - ★ no overhead of writing log records
  - ★ recovery is trivial
- Disadvantages :
  - ★ Copying the entire page table is very expensive
    - Can be reduced by using a page table structured like a B<sup>+</sup>-tree
      - No need to copy entire tree, only need to copy paths in the tree that lead to updated leaf nodes
  - ★ Commit overhead is high even with above extension
    - Need to flush every updated page, and page table
  - ★ Data gets fragmented (related pages get separated on disk)
  - ★ After every transaction completion, the database pages containing old versions of modified data need to be garbage collected
  - ★ Hard to extend algorithm to allow transactions to run concurrently
    - Easier to extend log based schemes

# Recovery With Concurrent Transactions

- We modify the log-based recovery schemes to allow multiple transactions to execute concurrently.
  - ★ All transactions share a single disk buffer and a single log
  - ★ A buffer block can have data items updated by one or more transactions
- We assume concurrency control using strict two-phase locking;
  - ★ i.e. the updates of uncommitted transactions should not be visible to other transactions
    - Otherwise how to perform undo if T1 updates A, then T2 updates A and commits, and finally T1 has to abort?
- Logging is done as described earlier.
  - ★ Log records of different transactions may be interspersed in the log.
- The checkpointing technique and actions taken on recovery have to be changed
  - ★ since several transactions may be active when a checkpoint is performed.

# Recovery With Concurrent Transactions (Cont.)

- Checkpoints are performed as before, except that the checkpoint log record is now of the form  
    < **checkpoint**  $L$  >  
where  $L$  is the list of transactions active at the time of the checkpoint
  - ★ We assume no updates are in progress while the checkpoint is carried out (will relax this later)
- When the system recovers from a crash, it first does the following:
  1. Initialize *undo-list* and *redo-list* to empty
  2. Scan the log backwards from the end, stopping when the first <**checkpoint**  $L$ > record is found.  
For each record found during the backward scan:
    - 👉 if the record is < $T_i$  **commit**>, add  $T_i$  to *redo-list*
    - 👉 if the record is < $T_i$  **start**>, then if  $T_i$  is not in *redo-list*, add  $T_i$  to *undo-list*
  3. For every  $T_i$  in  $L$ , if  $T_i$  is not in *redo-list*, add  $T_i$  to *undo-list*

# Recovery With Concurrent Transactions (Cont.)

- At this point *undo-list* consists of incomplete transactions which must be undone, and *redo-list* consists of finished transactions that must be redone.
- Recovery now continues as follows:
  1. Scan log backwards from most recent record, stopping when  $\langle T_i \text{ start} \rangle$  records have been encountered for every  $T_i$  in *undo-list*.
    - During the scan, perform **undo** for each log record that belongs to a transaction in *undo-list*.
  2. Locate the most recent **<checkpoint L>** record.
  3. Scan log forwards from the **<checkpoint L>** record till the end of the log.
    - During the scan, perform **redo** for each log record that belongs to a transaction on *redo-list*

# Example of Recovery

- Go over the steps of the recovery algorithm on the following log:

$\langle T_0 \text{ start} \rangle$

$\langle T_0, A, 0, 10 \rangle$

$\langle T_0 \text{ commit} \rangle$

$\langle T_1 \text{ start} \rangle$

$\langle T_1, B, 0, 10 \rangle$

$\langle T_2 \text{ start} \rangle$

*/\* Scan in Step 4 stops here \*/*

$\langle T_2, C, 0, 10 \rangle$

$\langle T_2, C, 10, 20 \rangle$

$\langle \text{checkpoint } \{T_1, T_2\} \rangle$

$\langle T_3 \text{ start} \rangle$

$\langle T_3, A, 10, 20 \rangle$

$\langle T_3, D, 0, 10 \rangle$

$\langle T_3 \text{ commit} \rangle$

# Log Record Buffering

- **Log record buffering**: log records are buffered in main memory, instead of being output directly to stable storage.
  - ★ Log records are output to stable storage when a block of log records in the buffer is full, or a **log force** operation is executed.
- Log force is performed to commit a transaction by forcing all its log records (including the commit record) to stable storage.
- Several log records can thus be output using a single output operation, reducing the I/O cost.

# Log Record Buffering (Cont.)

- The rules below must be followed if log records are buffered:
  - ★ Log records are output to stable storage in the order in which they are created.
  - ★ Transaction  $T_i$  enters the commit state only when the log record  $\langle T_i \text{ commit} \rangle$  has been output to stable storage.
  - ★ Before a block of data in main memory is output to the database, all log records pertaining to data in that block must have been output to stable storage.
    - This rule is called the **write-ahead logging** or **WAL** rule
      - Strictly speaking WAL only requires undo information to be output

# Backups

## ■ Backups are needed to recover from media failure

- ★ The transaction log and entire contents of the database is written to secondary storage (often tape)
- ★ Time consuming, and often requires down time

## ■ Backups frequency

- ★ Frequent enough that little information is lost
- ★ Not so frequent as to cause problems
- ★ Every day (night) is common

## ■ Backup storage



# Backup Modes

## ■ Hot backup

- ★ allows backup of the database while the database is running and available to users.
- ★ performance degrades during the backup period
- ★ takes longer than a cold backup

## ■ Cold backup

- ★ requires database shutdown before backup begins
- ★ physical files are backed up while shutdown
- ★ database is unavailable to users during backup period
- ★ faster than a hot backup

# Backup Types

- Complete (Full)
  - ★ copy all database and related files
  - ★ delete the archive log files
- Cumulative (Differential)
  - ★ copy blocks that have changed since last full backup  
or
  - ★ copy all archive log files generated since last full backup
- Incremental
  - ★ copy blocks that have change since the last partial backup  
or
  - ★ copy all log files generated since last partial backup
- Complete (Copy)
  - ★ copy all target data
  - ★ Don't include the set in backup set logic