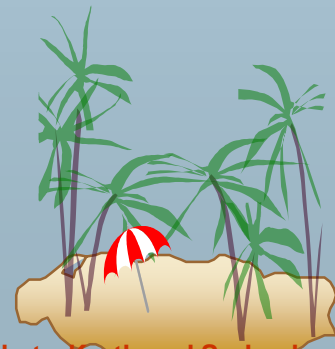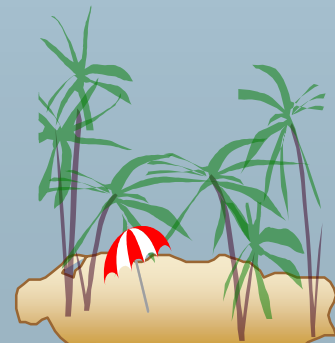# Concurrency Control

- Transaction Concept
- Transaction State
- Implementation of Atomicity and Durability
- Concurrent Executions
- Serializability
- Recoverability

# Transaction Concept

- A *transaction* is a *unit* of program execution that accesses and possibly updates various data items.

- A transaction must see a consistent database.

- During transaction execution the database may be inconsistent.

- When the transaction is committed, the database must be consistent.

- Two main issues to deal with:
  - ★ Failures of various kinds, such as hardware failures and system crashes
  - ★ Concurrent execution of multiple transactions

# ACID Properties

To preserve integrity of data, the database system must ensure:
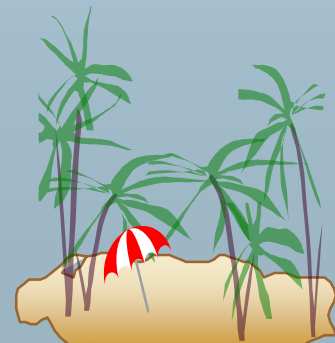
- **Atomicity.**  Either all operations of the transaction are properly reflected in the database or none are.

- **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.
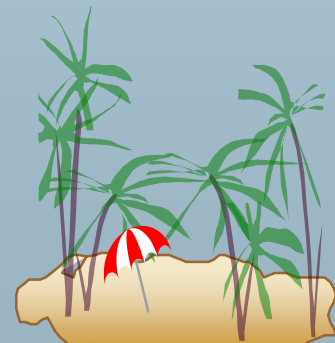
# Example of Fund Transfer

- 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.
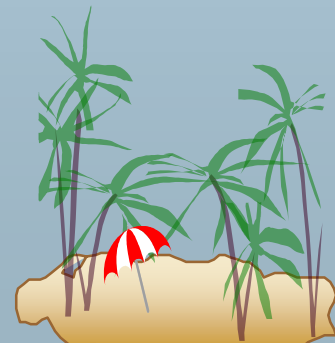
# Example of Fund Transfer (Cont.)

■ 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.
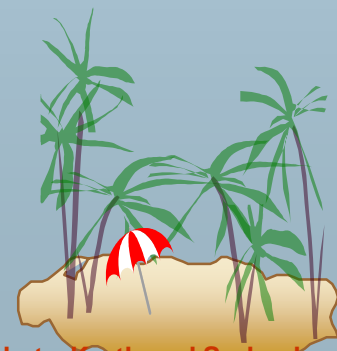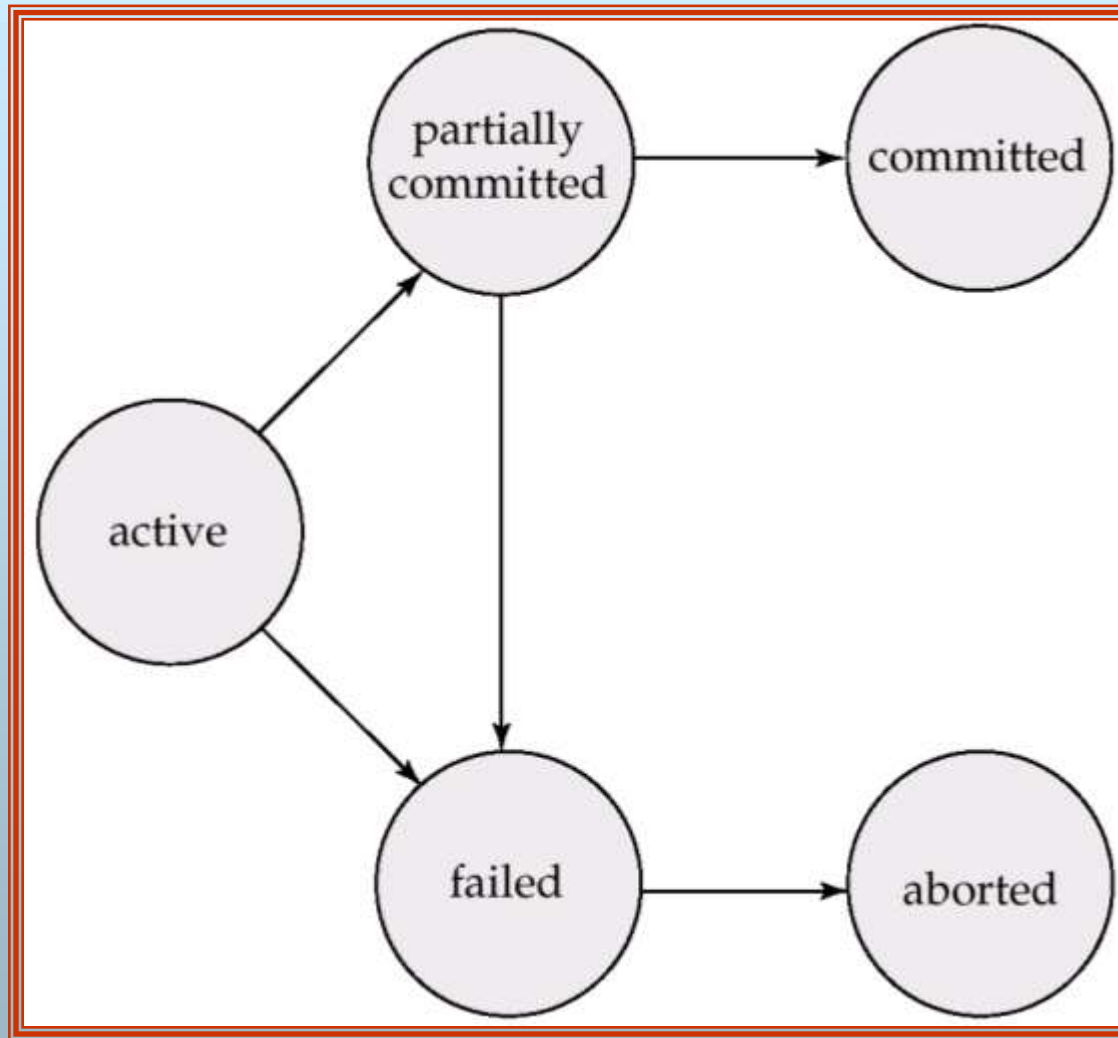
# Transaction State

- **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

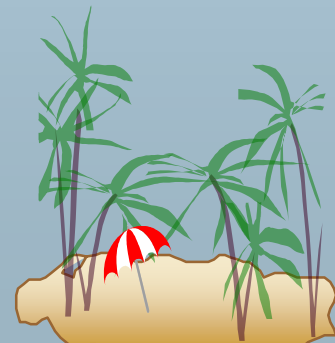- **Committed,** after *successful completion*.
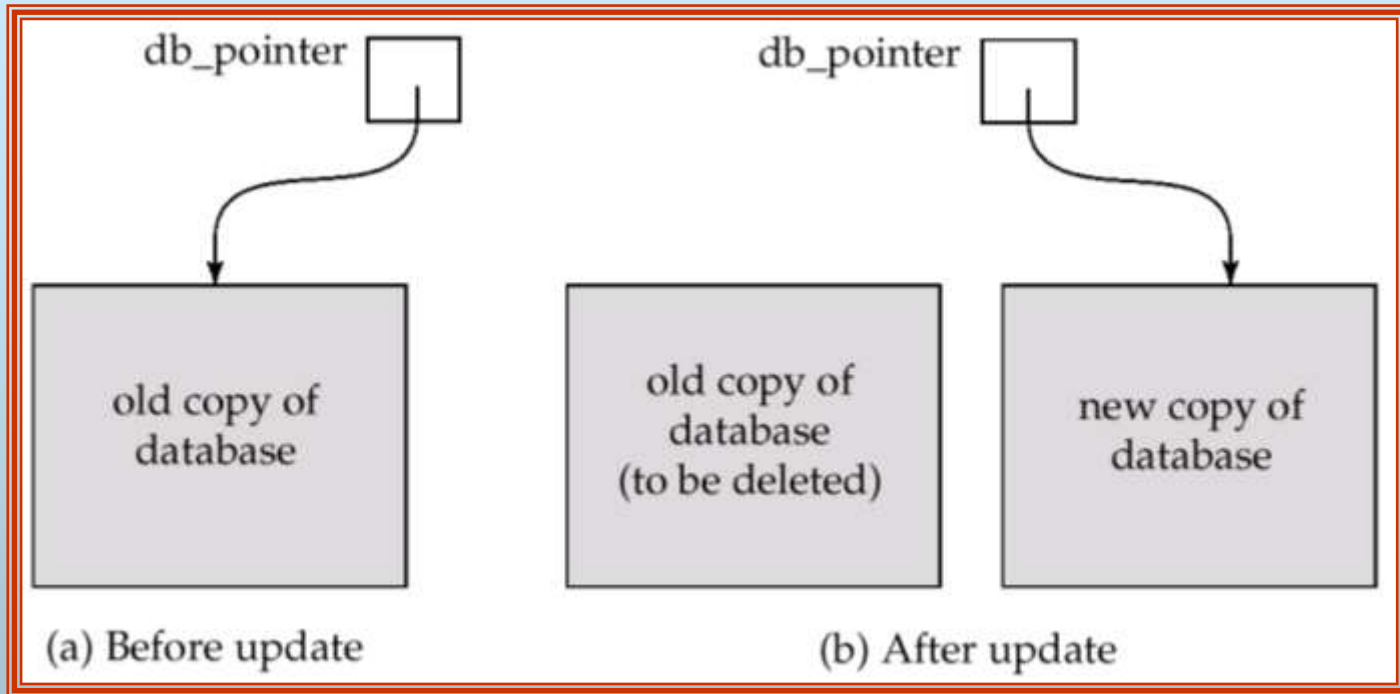
# Implementation of Atomicity and Durability

- The recovery-management component of a database system implements the support for atomicity and durability.

- The *shadow-database* scheme:
  - ★ assume that only one transaction is active at a time.
  - ★ a pointer called db_pointer always points to the current consistent copy of the database.
  - ★ all updates are made on a *shadow copy* of the database, and **db_pointer** is made to point to the updated shadow copy only after the transaction reaches partial commit and all updated pages have been flushed to disk.
  - ★ in case transaction fails, old consistent copy pointed to by **db_pointer** can be used, and the shadow copy can be deleted.

# Implementation of Atomicity and Durability (Cont.)

The shadow-database scheme:
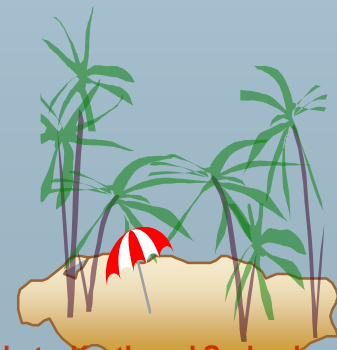


(a) Before update          (b) After update

- Assumes disks to not fail

- Useful for text editors, but extremely inefficient for large databases: executing a single transaction requires copying the *entire* database.  Will see better schemes in Chapter 17.
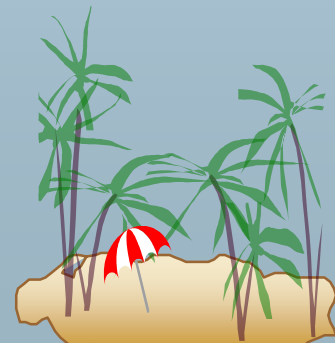
# Concurrent Executions

- 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
  - ★ Will study in Chapter 14, after studying notion of correctness of concurrent executions.

# 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

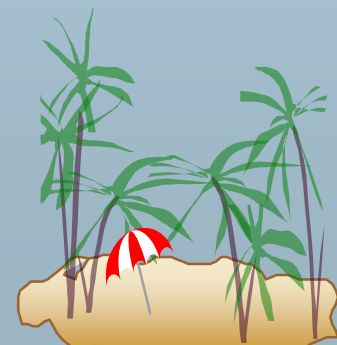★ must preserve the order in which the instructions appear in each individual transaction.

# 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 $T_2$.

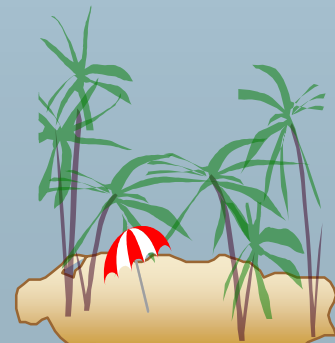| $T_1$ | $T_2$ |
|---|---|
| read($A$) | |
| $A := A - 50$ | |
| write ($A$) | |
| read($B$) | |
| $B := B + 50$ | |
| write($B$) | |
| | read($A$) |
| | $temp := A * 0.1$ |
| | $A := A - temp$ |
| | write($A$) |
| | read($B$) |
| | $B := B + temp$ |
| | write($B$) |

# 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.

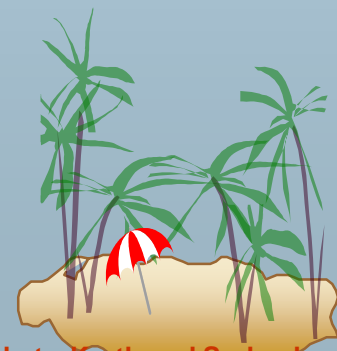| $T_1$ | $T_2$ |
|---|---|
| read($A$) | |
| $A := A - 50$ | |
| write($A$) | |
| | read($A$) |
| | $temp := A * 0.1$ |
| | $A := A - temp$ |
| | write($A$) |
| read($B$) | |
| $B := B + 50$ | |
| write($B$) | |
| | read($B$) |
| | $B := B + temp$ |
| | write($B$) |

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

# Example Schedules (Cont.)

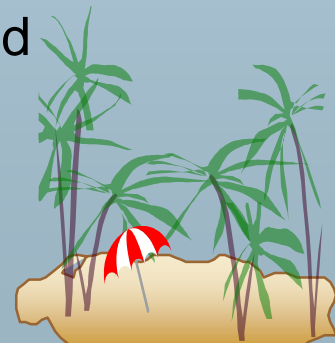■ The following concurrent schedule (Schedule 4 in the text) does not preserve the value of the the sum $A + B$.

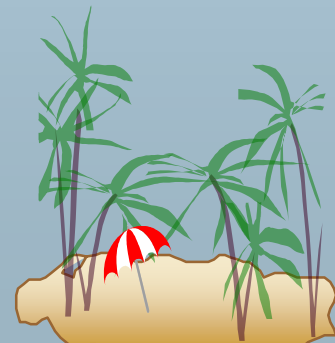| $T_1$ | $T_2$ |
|---|---|
| read($A$)<br>$A := A - 50$ | |
| | read($A$)<br>$temp := A * 0.1$<br>$A := A - temp$<br>write($A$)<br>read($B$) |
| write($A$)<br>read($B$)<br>$B := B + 50$<br>write($B$) | |
| | $B := B + temp$<br>write($B$) |

# 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. $I_i =$ **read**$(Q)$, $I_j =$ **read**$(Q)$.   $I_i$ and $I_j$ don't conflict.
  2. $I_i =$ **read**$(Q)$,  $I_j =$ **write**$(Q)$.  They conflict.
  3. $I_i =$ **write**$(Q)$, $I_j =$ **read**$(Q)$.   They conflict
  4. $I_i =$ **write**$(Q)$, $I_j =$ **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.
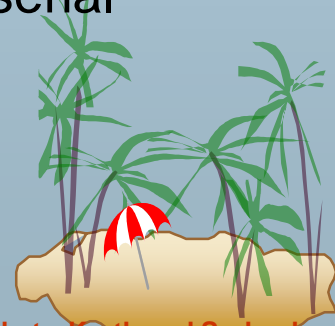
# Conflict Serializability (Cont.)

- 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:

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

We are unable to swap instructions in the above schedule to obtain either the serial schedule $< T_3, T_4 >$, or the serial schedule $< T_4, T_3 >$.
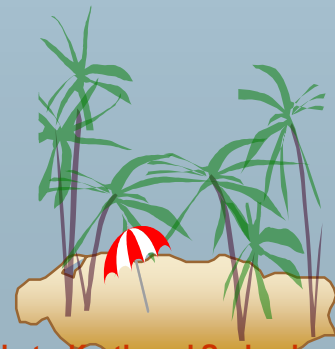
# Conflict Serializability (Cont.)

■ 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.

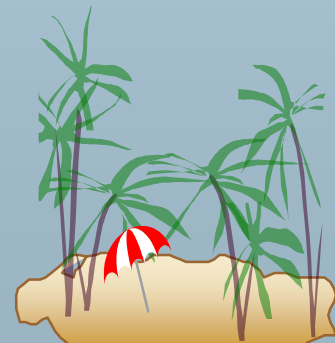| $T_1$ | $T_2$ |
|---|---|
| read($A$) | |
| write($A$) | |
| | read($A$) |
| | write($A$) |
| read($B$) | |
| write($B$) | |
| | read($B$) |
| | write($B$) |

# 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 $S'$.

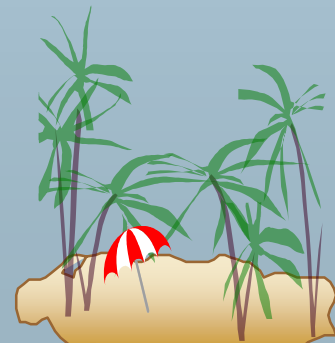As can be seen, view equivalence is also based purely on **reads** and **writes** alone.

# View Serializability (Cont.)

- A schedule $S$ is **view serializable** 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.**

# Recoverability

Need to address the effect of transaction failures on concurrently running transactions.

- **Recoverable schedule** — if a transaction $T_j$ reads a data items previously written by a transaction $T_i$, the commit operation of $T_i$ appears before the commit operation of $T_j$.

- The following schedule (Schedule 11) is not recoverable if $T_9$ commits immediately after the read

| $T_8$ | $T_9$ |
|---|---|
| read($A$) | |
| write($A$) | |
| | read($A$) |
| read($B$) | |

- If $T_8$ should abort, $T_9$ would have read (and possibly shown to the user) an inconsistent database state. Hence database must ensure that schedules are recoverable.
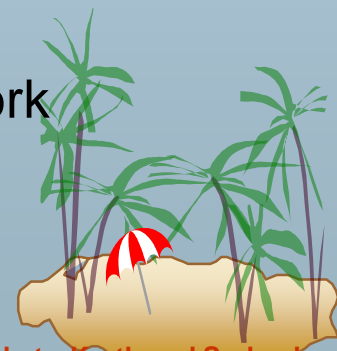
# Recoverability (Cont.)

■ **Cascading rollback** – a single transaction failure leads to a series of transaction rollbacks.  Consider the following schedule where none of the transactions has yet committed (so the schedule is recoverable)

| $T_{10}$ | $T_{11}$ | $T_{12}$ |
|---|---|---|
| read(A)<br>read(B)<br>write(A) | | |
| | read(A)<br>write(A) | |
| | | read(A) |

If $T_{10}$ fails, $T_{11}$ and $T_{12}$ must also be rolled back.

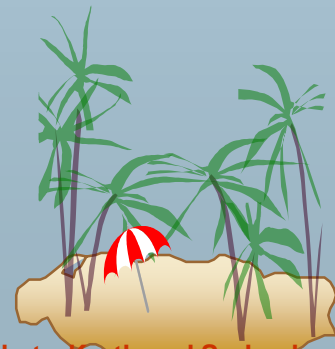■ Can lead to the undoing of a significant amount of work

# Recoverability (Cont.)

■ **Cascadeless schedules** — cascading rollbacks cannot occur; for each pair of transactions $T_i$ and $T_j$ such that $T_j$ reads a data item previously written by $T_i$, the commit operation of $T_i$ appears before the read operation of $T_j$.

■ Every cascadeless schedule is also recoverable

■ It is desirable to restrict the schedules to those that are cascadeless
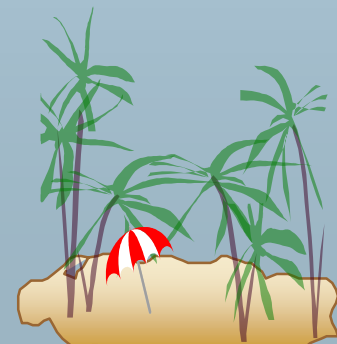
# Concurrency Control

- Lock-Based Protocols

- Deadlock Handling

# 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.
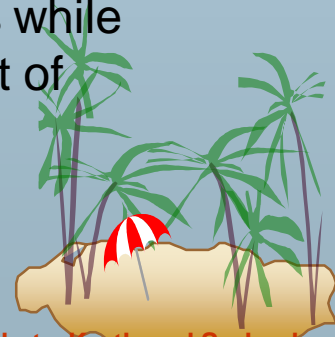
# Lock-Based Protocols (Cont.)

- 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.
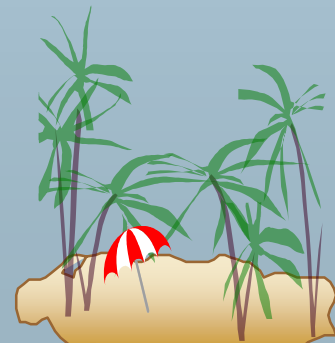
# Pitfalls of Lock-Based Protocols

■ Consider the partial schedule

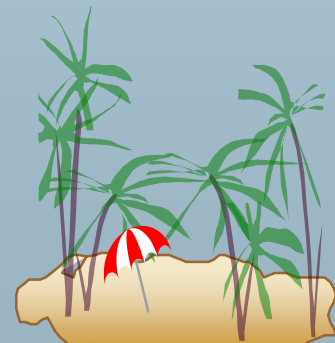| $T_3$ | $T_4$ |
|---|---|
| lock-x($B$) | |
| read($B$) | |
| $B := B - 50$ | |
| write($B$) | |
| | lock-s($A$) |
| | read($A$) |
| | lock-s($B$) |
| lock-x($A$) | |

■ 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.
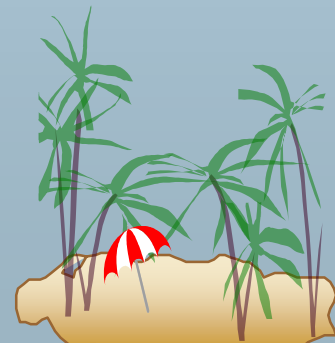
# Pitfalls of Lock-Based Protocols (Cont.)

- 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.
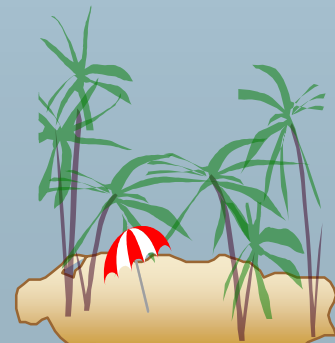
# The Two-Phase Locking Protocol

- 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.
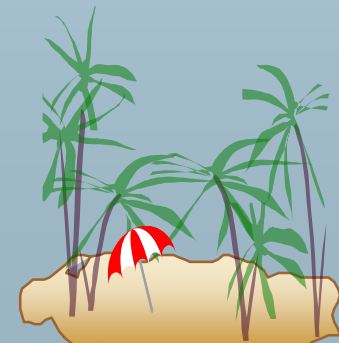
# The Two-Phase Locking Protocol (Cont.)

- 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.
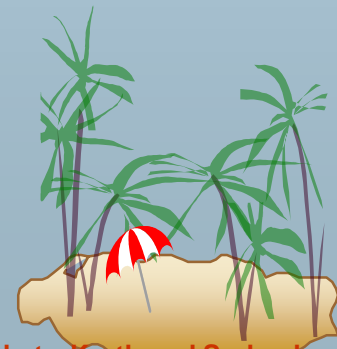
# Deadlock Handling

■ Consider the following two transactions:

$T_1$:    write ($X$)          $T_2$:    write($Y$)

     write($Y$)               write($X$)
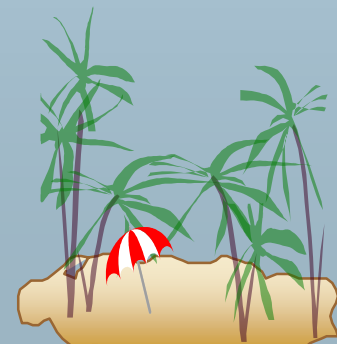
■ Schedule with deadlock

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

# 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).
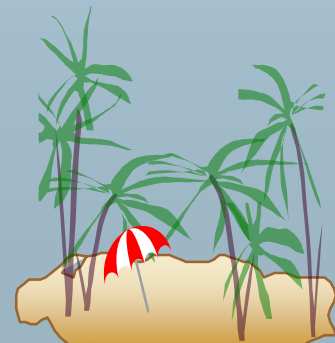
# More Deadlock Prevention Strategies

- 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.

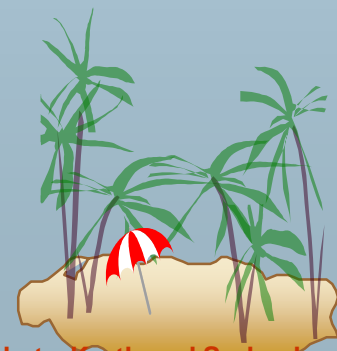  - ★ may be fewer rollbacks than *wait-die* scheme.

# 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.
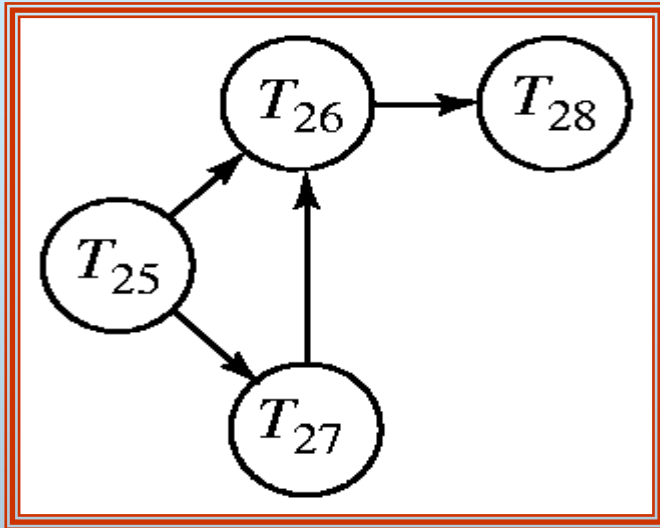
# Deadlock Detection

- Deadlocks can be described as a *wait-for graph*, which consists of a pair $G = (V,E)$,
  - ★ $V$ is a set of vertices (all the transactions in the system)
  - ★ $E$ is a set of edges; each element is an ordered pair $T_i \rightarrow T_j$.

- If $T_i \rightarrow T_j$ is in $E$, then there is a directed edge from $T_i$ to $T_j$, implying that $T_i$ is waiting for $T_j$ to release a data item.

- When $T_i$ requests a data item currently being held by $T_j$, then the edge $T_i \ T_j$ is inserted in the wait-for graph. This edge is removed only when $T_j$ is no longer holding a data item needed by $T_i$.

- The system is in a deadlock state if and only if the wait-for graph has a cycle. Must invoke a deadlock-detection algorithm periodically to look for cycles.
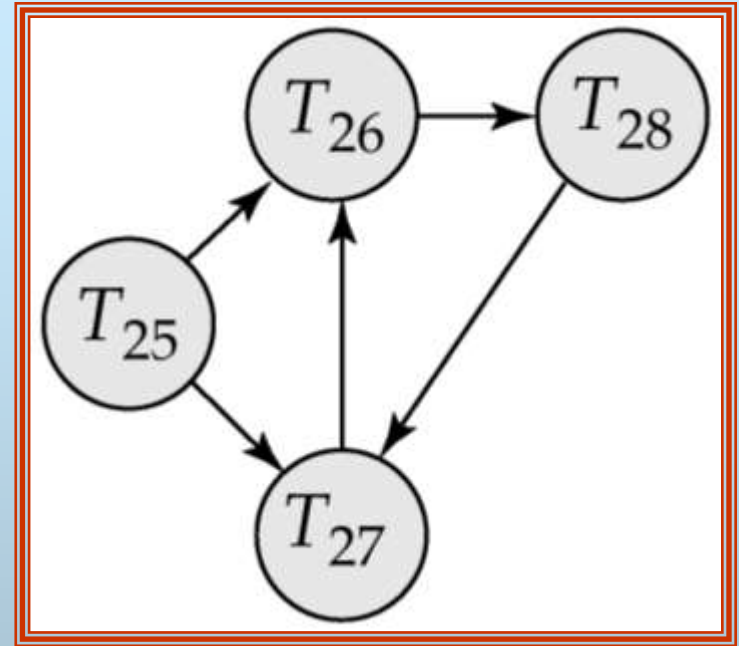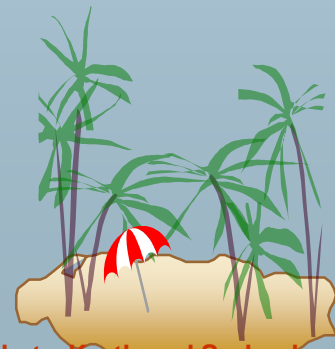
# Deadlock Detection (Cont.)



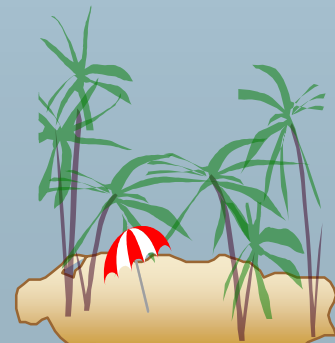Wait-for graph without a cycle

Wait-for graph with a cycle
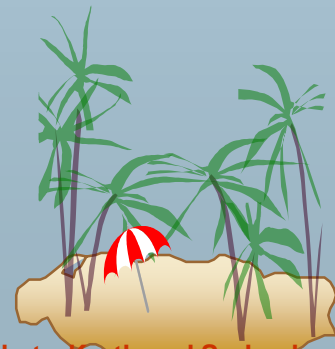
# Deadlock Recovery

- When deadlock is detected :
  - ★ Some transaction will have to rolled back (made a victim) to break deadlock. Select that transaction as victim that will incur minimum cost.
  - ★ Rollback -- determine how far to roll back transaction
    - ➢ Total rollback: Abort the transaction and then restart it.
    - ➢ More effective to roll back transaction only as far as necessary to break deadlock.
  - ★ Starvation happens if same transaction is always chosen as victim. Include the number of rollbacks in the cost factor to avoid starvation

# Recovery System

- Failure Classification

- Storage Structure

- Recovery and Atomicity

- Log-Based Recovery

- Shadow Paging

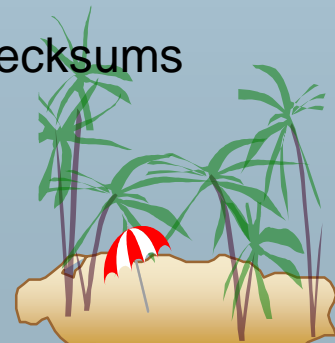# Failure Classification

- **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
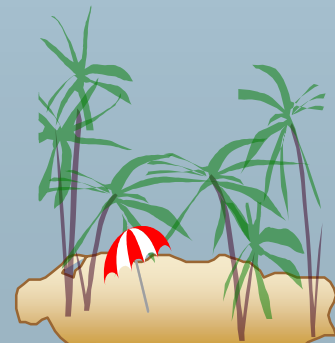  - ★ Destruction is assumed to be detectable: disk drives use checksums to detect failures
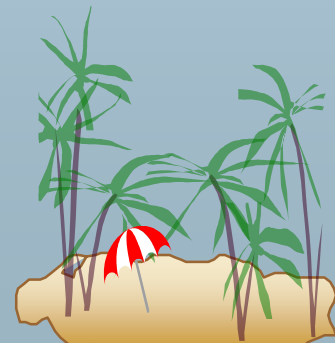
# Recovery Algorithms

- 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
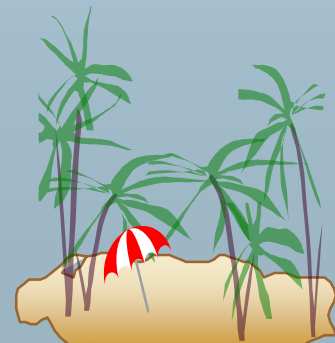
# Recovery and Atomicity

- 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.
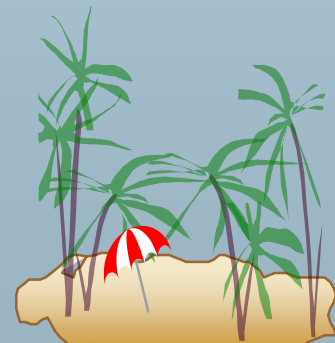
# Recovery and Atomicity (Cont.)

■ 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
  $<T_i$ **start**$>$log record
- *Before $T_i$* executes **write**$(X)$, a log record $<T_i, X, V_1, V_2>$ 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 it last statement, the log record $<T_i$ **commi**t$>$ 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 **write**s to after partial commit.

- Assume that transactions execute serially

- Transaction starts by writing $<T_i$ **start**$>$ record to log.

- A **write**$(X)$ operation results in a log record $<T_i, X, V>$ 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, $<T_i$ **commit**$>$ 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 $<T_i$ **start**$>$ and$<T_i$ **commit**$>$ 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** ($A$)                                 $T_1$ : **read** ($C$)

     $A$: - $A$ - 50                                      $C$:- $C$- 100

     **Write** ($A$)                                      **write** ($C$)

     **read** ($B$)

     $B$:-  $B$ + 50

     **write** ($B$)

# Deferred Database Modification (Cont.)

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

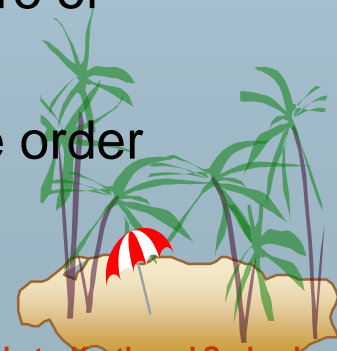| | | |
|---|---|---|
| $<T_0$ start> | $<T_0$ start> | $<T_0$ start> |
| $<T_0, A, 950>$ | $<T_0, A, 950>$ | $<T_0, A, 950>$ |
| $<T_0, B, 2050>$ | $<T_0, B, 2050>$ | $<T_0, B, 2050>$ |
| | $<T_0$ commit> | $<T_0$ commit> |
| | $<T_1$ start> | $<T_1$ start> |
| | $<T_1, C, 600>$ | $<T_1, C, 600>$ |
| | | $<T_1$ commit> |
| (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 $<T_0$ **commi**t> is present

(c) **redo**($T_0$) must be performed followed by redo($T_1$) since
        $<T_0$ **commit**> and $<T_i$ commit> 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 |
|-----|-------|--------|

$<T_0$ **start**$>$

$<T_0,$ A, 1000, 950$>$
$T_o,$ B, 2000, 2050

$\qquad\qquad\qquad A = 950$
$\qquad\qquad\qquad B = 2050$

$<T_0$ **commit**$>$
$<T_1$ **start**$>\quad x_1$
$<T_1,$ C, 700, 600$>$

$\qquad\qquad\qquad C = 600$

$\qquad\qquad\qquad\qquad\qquad\qquad B_B, B_C$

$<T_1$ **commit**$>$

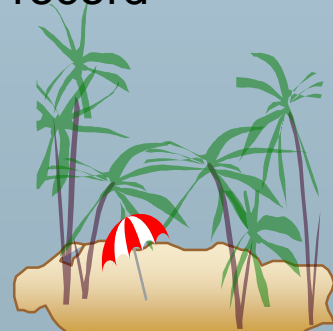$\qquad\qquad\qquad\qquad\qquad\qquad B_A$

- 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 <$T_i$ **start**>, but does not contain the record <$T_i$ **commit**>.
  - ★ Transaction $T_i$ needs to be redone if the log contains both the record <$T_i$ **start**> and the record <$T_i$ **commit**>.
- 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.

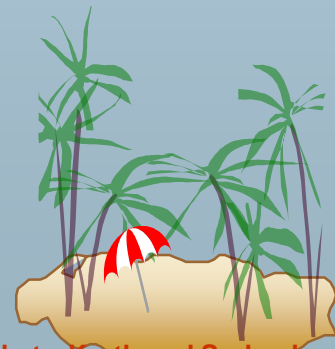| | | |
|---|---|---|
| $<T_0$ start$>$ | $<T_0$ start$>$ | $<T_0$ start$>$ |
| $<T_0$, A, 1000, 950$>$ | $<T_0$, A, 1000, 950$>$ | $<T_0$, A, 1000, 950$>$ |
| $<T_0$, B, 2000, 2050$>$ | $<T_0$, B, 2000, 2050$>$ | $<T_0$, B, 2000, 2050$>$ |
| | $<T_0$ commit$>$ | $<T_0$ commit$>$ |
| | $<T_1$ start$>$ | $<T_1$ start$>$ |
| | $<T_1$, C, 700, 600$>$ | $<T_1$, C, 700, 600$>$ |
| | | $<T_1$ commit$>$ |
| (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

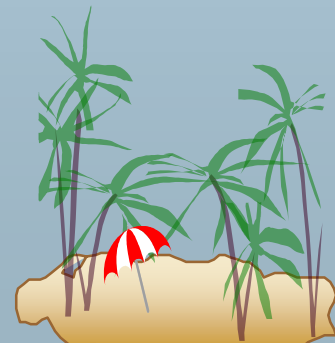> respectively. Then *C* is set to 600

# 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.
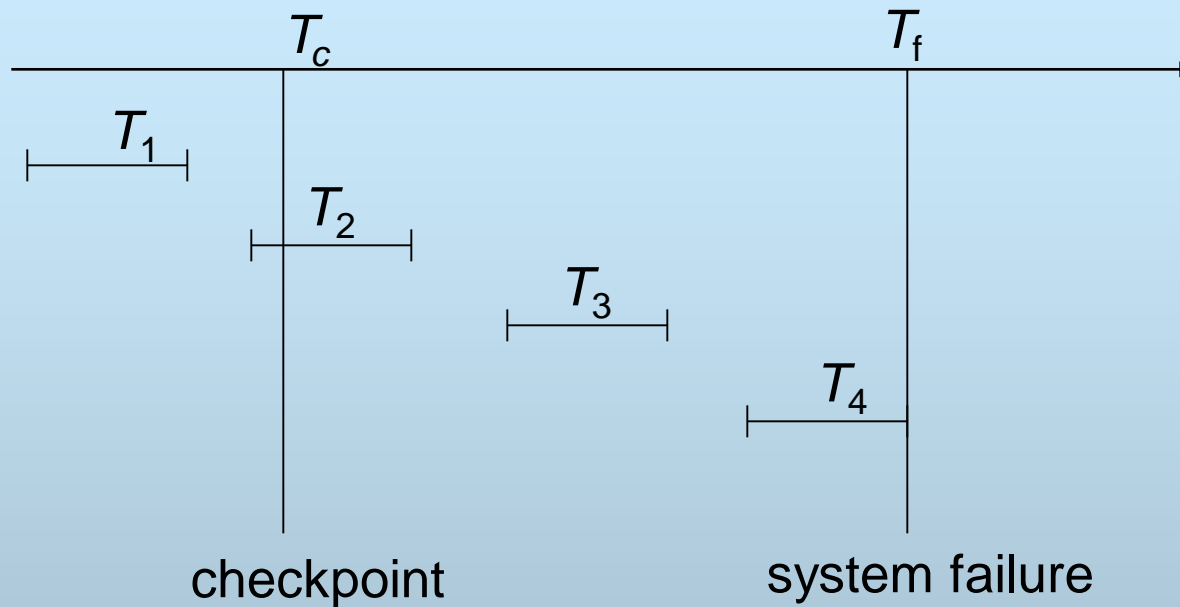    3. Write a log record < **checkpoint**> onto stable storage.

# 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 **star**t 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($T_i$)**.
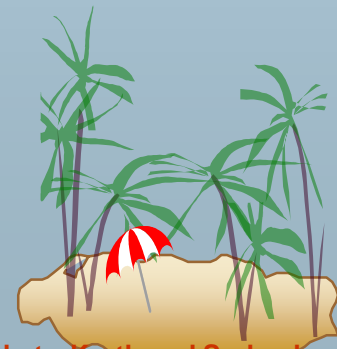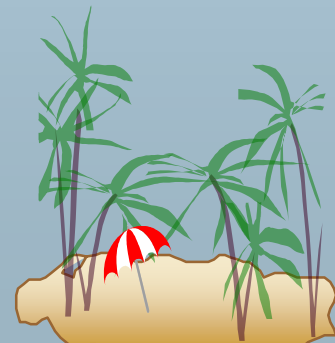
# Example of Checkpoints



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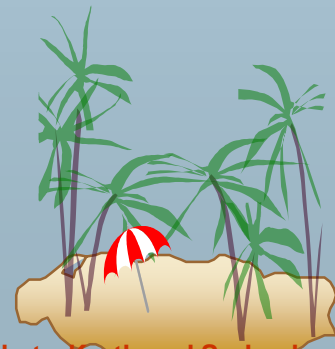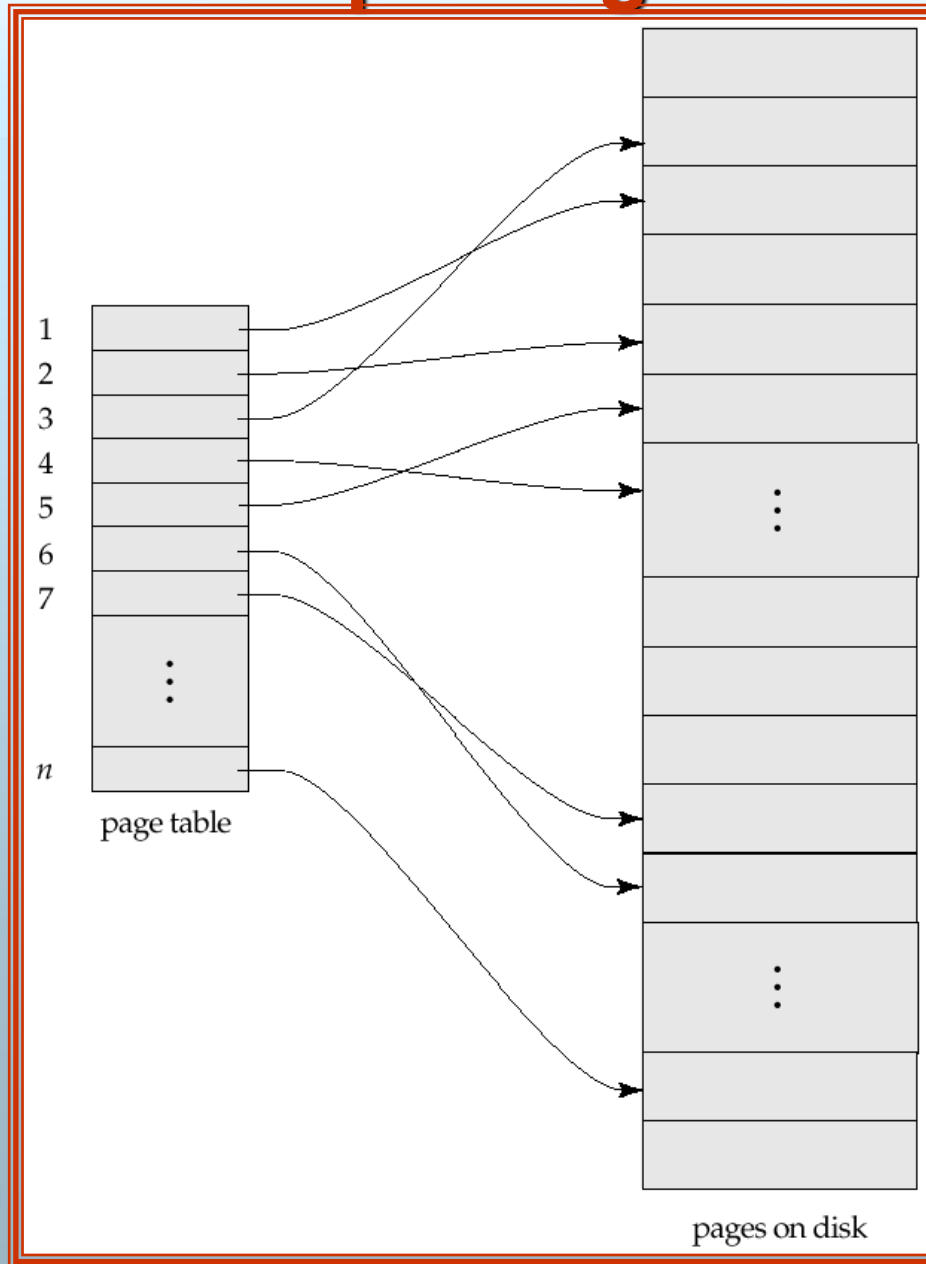- $T_2$ and $T_3$ redone.
- $T_4$ undone

# 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
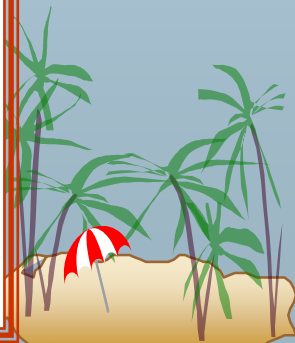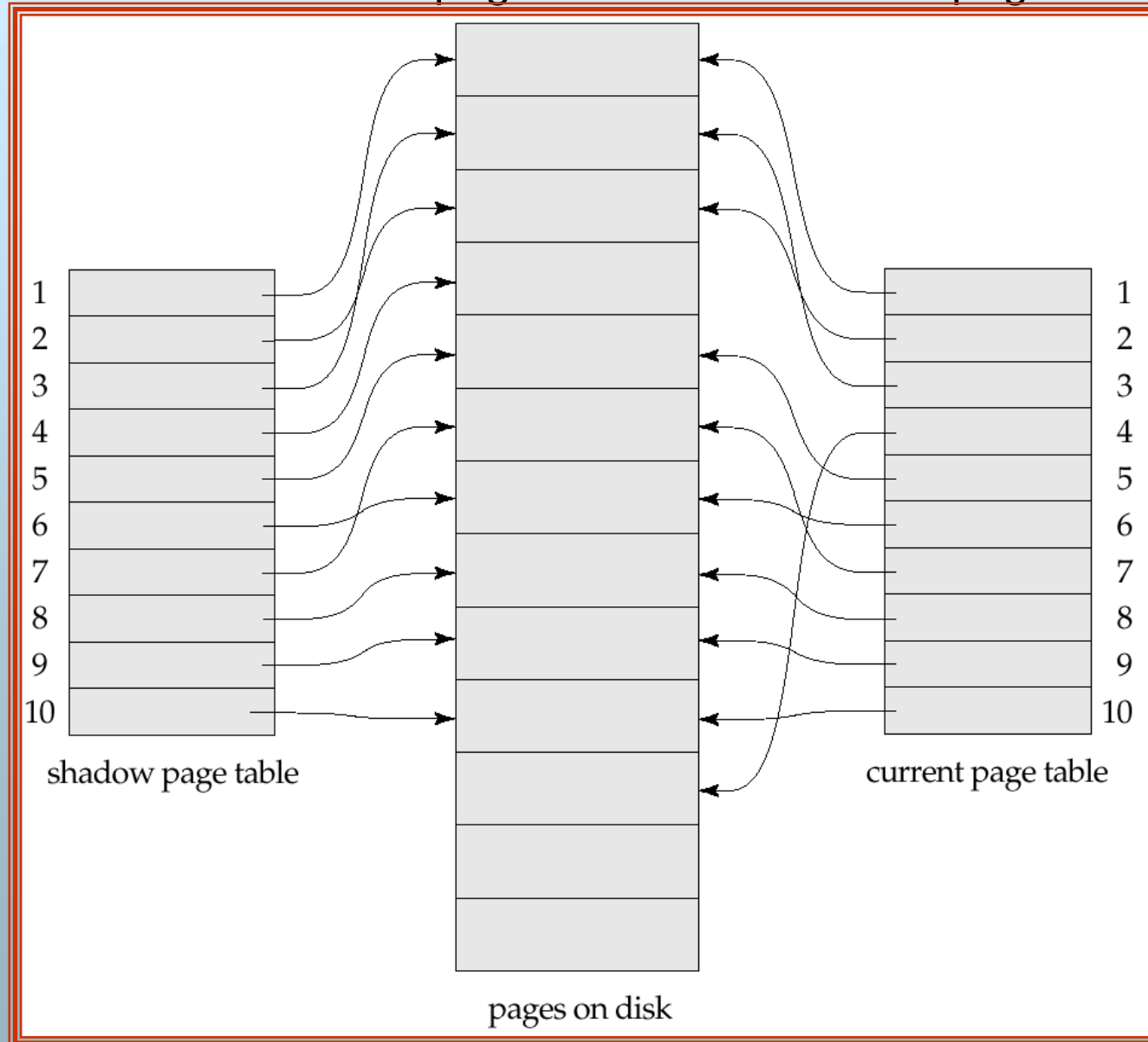  - ★ The update is performed on the copy

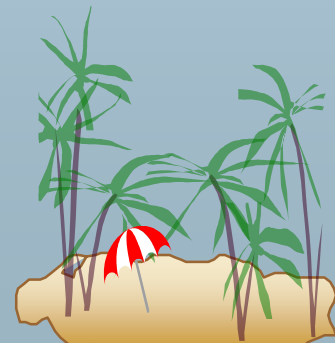# Sample Page Table

# Example of Shadow Paging

Shadow and current page tables after write to page 4



shadow page table

pages on disk

current page table

# 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).

# Show 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$^+$-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