

Chapter 17: Recovery System

- 17.1 Failure Classification
- 17.2 Storage Structure
- 17.3 Recovery and Atomicity
- 17.4 Log-Based Recovery
- 17.5 Shadow Paging
- 17.6 Recovery With Concurrent Transactions
- 17.7 Buffer Management **skip**
- 17.8 Failure with Loss of Nonvolatile Storage
- 17.9 Advanced Recovery Techniques **skip**
- 17.10 Remote Backup Systems

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.
 - ★ Recovery algorithms are the 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.

Stable-Storage Implementation

- Maintain multiple copies of each block on separate disks.
 - ★ Copies can be at remote sites to protect against disasters such as fire or flooding.
- Failure during data transfer can still result in inconsistent copies. Block transfer can result in
 - ★ Successful completion.
 - ★ Partial failure: destination block has incorrect information.
 - ★ Total failure: destination block was never updated.
- Protecting storage media from failure during data transfer (one solution):
 - ★ Execute output operation as follows (assuming two copies of each block):
 1. Write the information onto the first physical block.
 2. When the first write successfully completes, write the same information onto the second physical block.
 3. The output is completed only after the second write successfully completes.

Stable-Storage Implementation (Cont.)

- Protecting storage media from failure during data transfer (cont.):
 - ★ Copies of a block may differ due to failure during output operation. To recover from failure:
 1. First find inconsistent blocks:
 2. *Expensive solution*: Compare the two copies of every disk block.
 3. *Better solution*:
 - Record in-progress disk writes on non-volatile storage (Non-volatile RAM or special area of disk).
 - Use this information during recovery to find blocks that may be inconsistent, and only compare copies of these.
 - Used in hardware RAID systems
 4. If either copy of an inconsistent block is detected to have an error (bad checksum), overwrite it by the other copy. If both have no error, but are different, overwrite the second block by the first block.

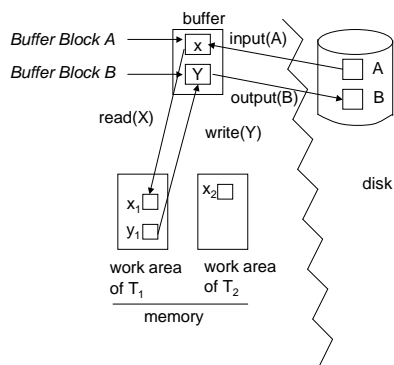
Data Access

- **Physical blocks** are those blocks residing on the disk.
- **Buffer blocks** are the blocks residing temporarily in main memory.
- Block movements between disk and main memory are initiated through the following two operations:
 - ★ **input(B)** transfers the physical block B to main memory.
 - ★ **output(B)** transfers the buffer block B to the disk, and replaces the appropriate physical block there.
- Each transaction T_i has its private work-area in which local copies of all data items accessed and updated by it are kept.
 - ★ T_i 's local copy of a data item X is called x_i .
- We assume, for simplicity, that each data item fits in, and is stored inside, a single block.

Data Access (Cont.)

- Transaction transfers data items between system buffer blocks and its private work-area using the following operations:
 - ★ **read(X)** assigns the value of data item X to the local variable x_i .
 - ★ **write(X)** assigns the value of local variable x_i to data item $\{X\}$ in the buffer block.
 - ★ Both these commands may necessitate the issue of an **input(B_x)** instruction before the assignment, if the block B_x in which X resides is not already in memory.
- Transactions
 - ★ Perform **read(X)** while accessing X for the first time.
 - ★ All subsequent accesses are to the local copy.
 - ★ After last access, transaction executes **write(X)**.
- **output(B_x)** need not immediately follow **write(X)**. System can perform the **output** operation when it deems fit.

Example of Data Access



Recovery and Atomicity

- Modifying the database without ensuring that the transaction will commit may leave the database in an inconsistent state.
- Consider transaction T that transfers \$50 from account A to account B ; goal is either to perform all database modifications made by T or none at all.
- Several output operations may be required for T (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 starts, it registers itself by writing a **< T start>** log record.
- Before T executes **write(X)**, a log record **< T , 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 has performed a write on data item X . X had value V_1 before the write, and will have value V_2 after the write.
- When T finishes its last statement, the log record **< T commit>** 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 **<T start>** record to log.
- A **write(X)** operation results in a log record **<T, 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 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 start>** and **<T commit>** are there in the log.
- Redoing a transaction *T* (**redo T**) 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 *T₀* and *T₁* (*T₀* executes before *T₁*):

<i>T₀</i> : read (A)	<i>T₁</i> : 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.

<T₀ start>	<T₀ start>	<T₀ start>
<T₀, A, 950>	<T₀, A, 950>	<T₀, A, 950>
<T₀, B, 2050>	<T₀, B, 2050>	<T₀, B, 2050>
	<T₀ commit>	<T₀ commit>
	<T₁ start>	<T₁ start>
	<T₁, C, 600>	<T₁, C, 600>
		<T₁ commit>
(a)	(b)	(c)

- If log on stable storage at time of crash is as in case:
 - No redo actions need to be taken.
 - redo(*T₀*) must be performed since **<T₀ commit>** is present.
 - redo(*T₀*) must be performed followed by redo(*T₁*) since **<T₀ commit>** and **<T₁ 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
<hr/>		
<T₀ start>		
<T₀, A, 1000, 950>		
T₀, B, 2000, 2050		
	A = 950	
	B = 2050	
<T₀ commit>		
<T₁ start> <i>x₁</i>		
<T₁, C, 700, 600>		
	C = 600	
		<i>B_B, B_C</i>
<T₁ commit>		<i>B_A</i>

- Note: *B_x* denotes block containing *X*.

Immediate Database Modification (Cont.)

- Recovery procedure has two operations instead of one:
 - ★ **undo(T)** restores the value of all data items updated by *T* to their old values, going backwards from the last log record for *T*.
 - ★ **redo(T)** sets the value of all data items updated by *T* to the new values, going forward from the first log record for *T*.
- 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* needs to be undone if the log contains the record **<T start>**, but does not contain the record **<T commit>**.
 - ★ Transaction *T* needs to be redone if the log contains both the record **<T start>** and the record **<T 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.

$\langle T_0 \text{ start} \rangle$	$\langle T_0 \text{ start} \rangle$	$\langle T_0 \text{ start} \rangle$
$\langle T_0, A, 1000, 950 \rangle$	$\langle T_0, A, 1000, 950 \rangle$	$\langle T_0, A, 1000, 950 \rangle$
$\langle T_0, B, 2000, 2050 \rangle$	$\langle T_0, B, 2000, 2050 \rangle$	$\langle T_0, B, 2000, 2050 \rangle$
	$\langle T_0 \text{ commit} \rangle$	$\langle T_0 \text{ commit} \rangle$
	$\langle T_1 \text{ start} \rangle$	$\langle T_1 \text{ start} \rangle$
	$\langle T_1, C, 700, 600 \rangle$	$\langle T_1, C, 700, 600 \rangle$
		$\langle T_1 \text{ commit} \rangle$

- (a) 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

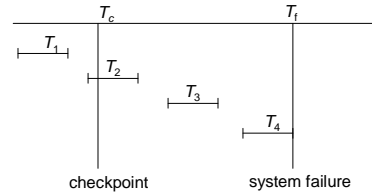
Checkpoints

- Problems in recovery procedure as discussed earlier:
 - Searching the entire log is time-consuming
 - We might unnecessarily redo transactions which have already
 - Output their updates to the database.
- Streamline recovery procedure by periodically performing **checkpointing**:
 - Output all log records currently residing in main memory onto stable storage.
 - Output all modified buffer blocks to the disk.
 - Write a log record $\langle \text{checkpoint} \rangle$ onto stable storage.

Checkpoints (Cont.)

- During recovery we need to consider only the latest transaction T that started before the checkpoint, and transactions that started after T .
 - Scan backwards from end of log to find the most recent $\langle \text{checkpoint} \rangle$ record.
 - Continue scanning backwards till a record $\langle T \text{ start} \rangle$ is found.
 - Need only consider the part of log following this start record. Earlier part of log can be ignored during recovery, and can be erased whenever desired.
 - For all transactions T' (starting from T or later) with no $\langle T' \text{ commit} \rangle$, execute $\text{undo}(T')$. (Done only in case of immediate modification.)
 - Scanning forward in the log, for all transactions T' (starting from T or later) with a $\langle T' \text{ commit} \rangle$, execute $\text{redo}(T')$.

Example of Checkpoints

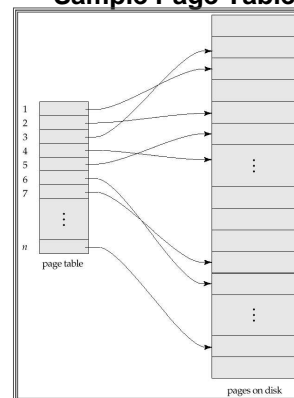


- T_1 can be ignored (updates already output to disk due to checkpoint)
- T_2 and T_3 redone.
- T_4 undone.

Shadow Paging

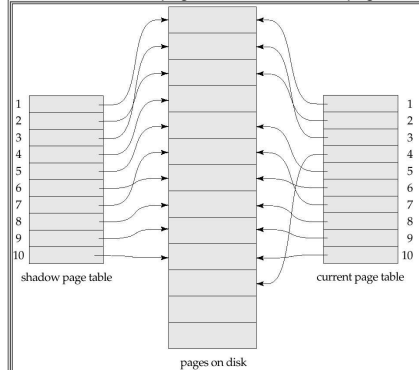
- Shadow paging** is an alternative to log-based recovery; this scheme is useful if transactions execute serially.
- Idea: maintain *two* page tables during the lifetime of a transaction – the **current page table**, and the **shadow page table**.
- Store the shadow page table in nonvolatile storage, such that state of the database prior to transaction execution may be recovered.
 - Shadow page table is never modified during execution
- Initially, 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.
 - ★ 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 commit>**, add *T* to *redo-list*
 - ⊕ if the record is **<T start>**, then if *T* is not in *redo-list*, add *T* to *undo-list*
 3. For every *T* in *L*, if *T* is not in *redo-list*, add *T* 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 **<T_i start>** 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:

```
<T0 start>
<T0, A, 0, 10>
<T0 commit>
<T1 start>
<T1, B, 0, 10>
<T2 start> /* Scan in Step 4 stops here */
<T2, C, 0, 10>
<T2, C, 10, 20>
<checkpoint {T1, T2}>
<T3 start>
<T3, A, 10, 20>
<T3, D, 0, 10>
<T3 commit>
```

Failure with Loss of Nonvolatile Storage

- So far we assumed no loss of non-volatile storage.
- Technique similar to checkpointing used to deal with loss of non-volatile storage
 - ★ Periodically **dump** the entire content of the database to stable storage
 - ★ No transaction may be active during the dump procedure; a procedure similar to checkpointing must take place
 - ✓ Output all log records currently residing in main memory onto stable storage.
 - ✓ Output all buffer blocks onto the disk.
 - ✓ Copy the contents of the database to stable storage.
 - ✓ Output a record **<dump>** to log on stable storage.
 - ★ To recover from disk failure
 - ✓ restore database from most recent dump.
 - ✓ Consult the log and redo all transactions that committed after the dump