

Recovery System

These slides are a modified version of the slides of the book “Database System Concepts” (Chapter 17), 5th Ed., [McGraw-Hill](#), by Silberschatz, Korth and Sudarshan. Original slides are available at www.db-book.com

Failures

- A computer system is subject to failures
- Causes are: disk failure, power outage, hardware or software errors,
- In any failure, information may be lost
- DBMS must take actions in advance to ensure that atomicity and durability properties of transactions are preserved
- **Recovery System:**
it can restore the database to the consistent state that existed before the failure

ASSUMPTIONS on failures:

- **System crash:** a power failure or other hardware or software failure causes the system to crash.
 - **Fail-stop assumption:** non-volatile storage contents are not corrupted by system crash
- **Disk failure:** a head crash or similar disk failure destroys all or part of disk storage
 - **Destruction is assumed to be detectable:** disk drivers use checksums to detect failures

Storage Structure

Resilience to failure classification:

■ **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
- **Information residing in stable storage is never lost!!!**
(theoretically cannot be guaranteed - it can be closely approximated by techniques that make data loss extremely unlikely)

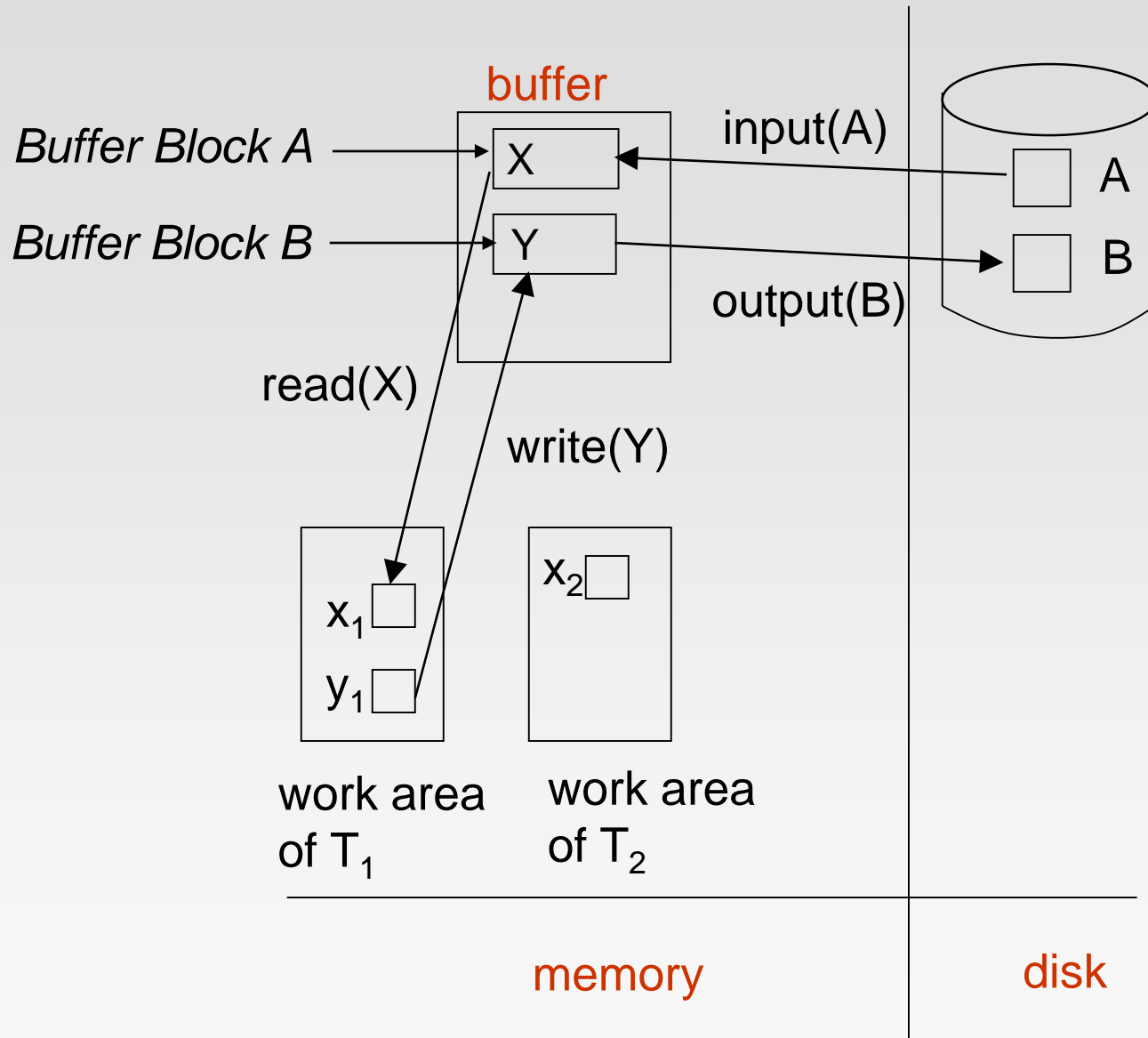
Data Access

- **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**
(the transaction aborts and some modifications have been made)

- **Performing the update of the database some time later than the commit of the transaction do not guarantee durability in case of faults**
(the transaction commits and the system crashes before all of the modifications are made)

Recovery Algorithms

- Recovery algorithms are techniques to ensure database transaction atomicity and durability despite failures
- 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 and durability

Recovery Algorithms (Cont.)

- To ensure atomicity despite failures, we first output information describing the modifications to stable storage without modifying the database itself.
- We study
 - **log-based recovery**

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 log record
 $\langle T_i \text{ start} \rangle$
- 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 (**partial commit of the transaction**), 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)

insert(X): log record $\langle T_i, X, \text{empty}, V_2 \rangle$

delete(X): log record $\langle T_i, X, V_1, \text{empty} \rangle$

Checkpoints

- Periodically performing **checkpointing**
 1. **Output all modified buffer blocks of partially committed transactions to the disk.**
 2. Restore all the original content of modified blocks of aborted transactions
 3. Write a log record < **checkpoint**> onto stable storage

Transactions are not allowed to execute any actions while a checkpoint is in progress.

Log-Based Recovery

Other possible approaches for the execution of the **output**(B_x) operation:

- 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 disk after partial commit.
 - the **output**(B_X) operation executed after the partial commit
- old value of X is not needed in the log file for this scheme
- 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
- The write is not performed on X , but it is deferred after the partial commit.
- When T_i partially commits, $\langle T_i, \text{commit} \rangle$ is written to the log

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 (**output(...)**) 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
$\langle T_0 \text{ start} \rangle$		
$\langle T_0, A, 1000, 950 \rangle$		
$\langle T_0, B, 2000, 2050 \rangle$		
	$A = 950$ $B = 2050$	
$\langle T_0 \text{ commit} \rangle$		
$\langle T_1 \text{ start} \rangle$		
$\langle T_1, C, 700, 600 \rangle$		
	$C = 600$	
$\langle T_1 \text{ commit} \rangle$		$B_B, B_C \leftarrow$ 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 output their updates to the database.

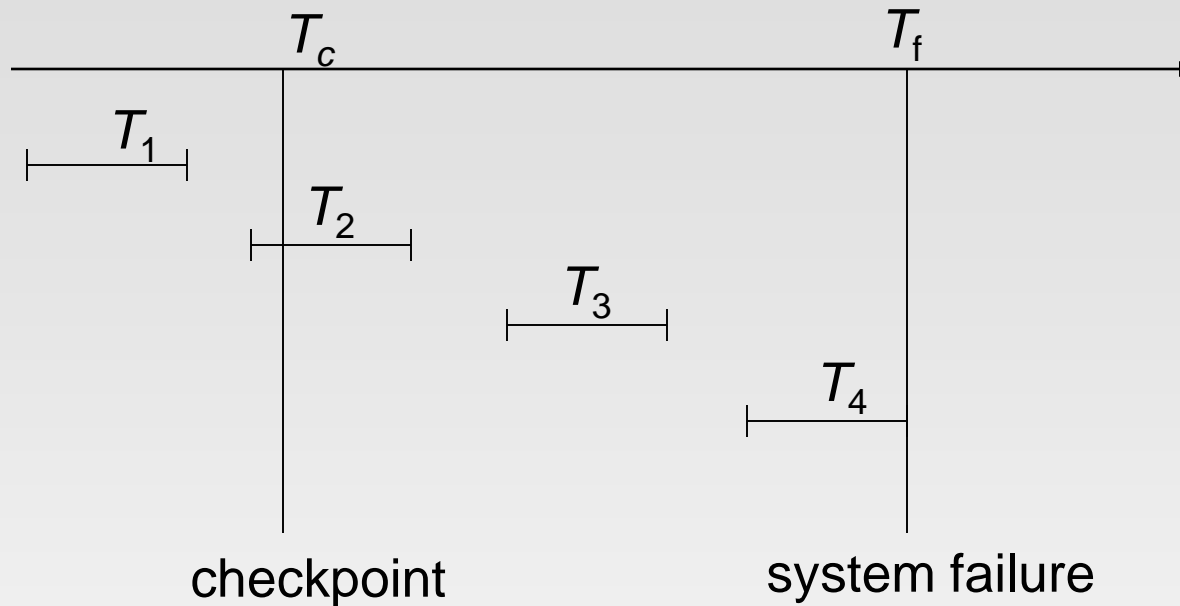
- Streamline recovery procedure by periodically performing **checkpointing**
 1. Output all **log records** currently residing in main memory onto stable storage.
 2. **Output all modified buffer blocks to the disk.**
 3. Write a log record < **checkpoint**> onto stable storage

Transactions are not allowed to execute any actions while a checkpoint is in progress.

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 for transaction in the checkpoint.
 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

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
- 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

- 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 L .
 - During the scan, perform **undo** for each log record that belongs to a transaction in *undo-list*.
 2. Scan log forwards from the $\langle T_i \text{ start} \rangle$ oldest record found at step 1 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:

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

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

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

$\langle T_1 \text{ start} \rangle$ /* Scan at step 1 comes up to here */

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

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

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

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

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

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

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

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

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

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.

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. **Log force** is performed to commit a transaction by forcing all its log records (including the commit record) 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

Database Buffering

- Database maintains an in-memory buffer of data blocks
 - When a new block is needed, if buffer is full an existing block needs to be removed from buffer
 - If the block chosen for removal has been updated, it must be output to disk
- If a block with uncommitted updates is output to disk, log records with undo information for the updates are output to the log on stable storage first
 - (Write ahead logging)
- No updates should be in progress on a block when it is output to disk.

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.

Recovering from Failure of Non-Volatile Storage

- To recover from disk failure
 - restore database from most recent dump.
 - Consult the log and redo all transactions that committed after the dump
 - Apply the Log Recovery

